

# MS1b

## Statistical Machine Learning and Data Mining

**Yee Whye Teh**  
Department of Statistics  
Oxford

<http://www.stats.ox.ac.uk/~teh/smldm.html>

### Course Structure

#### Lectures

- ▶ 1400-1500 Mondays in Math Institute L4.
- ▶ 1000-1100 Wednesdays in Math Institute L3.

#### Part C:

- ▶ 6 problem sheets.
- ▶ Classes: 1600-1700 Tuesdays (Weeks 3-8) in 1 SPR Seminar Room.
- ▶ Due Fridays week before classes at noon in 1 SPR.

#### MSc:

- ▶ 4 problem sheets.
- ▶ Classes: Tuesdays (Weeks 3, 5, 7, 9) in 2 SPR Seminar Room.
- ▶ Group A: 1400-1500, Group B: 1500-1600.
- ▶ Due Fridays week before classes at noon in 1 SPR.
- ▶ Practical: Week 5 and 7 (assessed) in 1 SPR Computing Lab.
- ▶ Group A: 1400-1600, Group B: 1600-1800.

1

3

### Course Information

- ▶ Course webpage:  
<http://www.stats.ox.ac.uk/~teh/smldm.html>
- ▶ Lecturer: Yee Whye Teh
- ▶ TA for Part C: Thibaut Lienart
- ▶ TA for MSc: Balaji Lakshminarayanan and Maria Lomeli
- ▶ Please subscribe to Google Group:  
<https://groups.google.com/forum/?hl=en-GB#!forum/smldm>
- ▶ Sign up for course using sign up sheets.

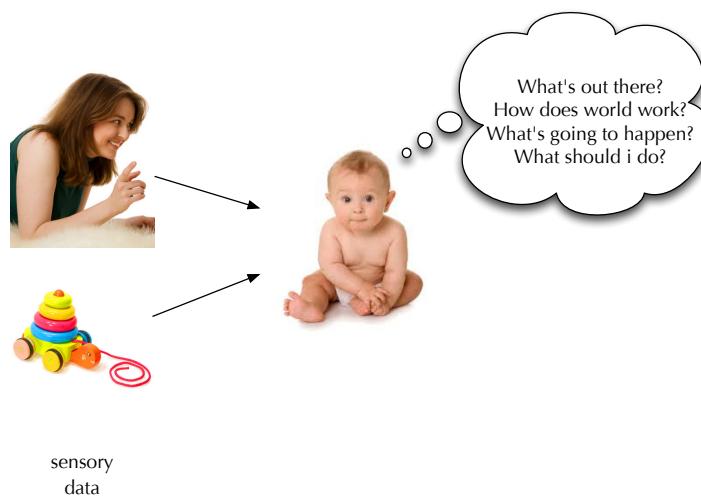
### Course Aims

1. Have ability to use the relevant R packages to analyse data, interpret results, and evaluate methods.
2. Have ability to identify and use appropriate methods and models for given data and task.
3. Understand the statistical theory framing machine learning and data mining.
4. Able to construct appropriate models and derive learning algorithms for given data and task.

2

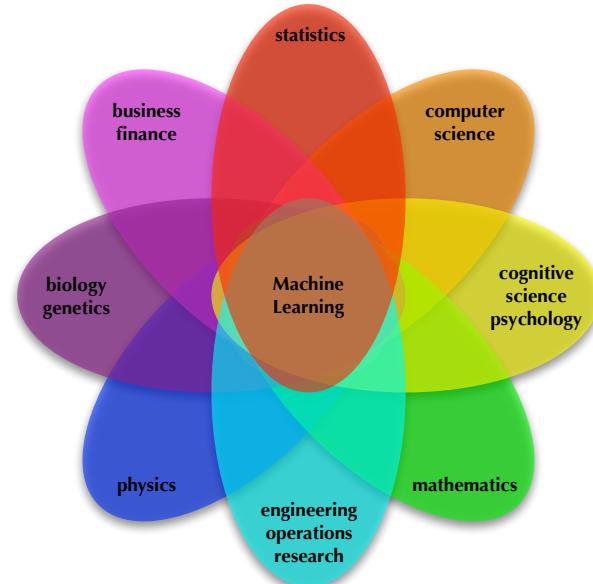
4

## What is Machine Learning?



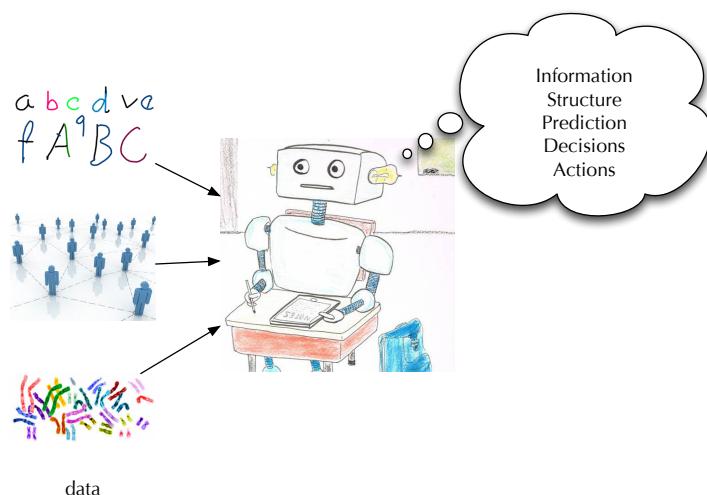
5

## What is Machine Learning?



7

## What is Machine Learning?



## What is the Difference?

### Traditional Problems in Applied Statistics

Well formulated question that we would like to answer.  
Expensive to gathering data and/or expensive to do computation.  
Create specially designed experiments to collect high quality data.

### Current Situation

#### Information Revolution

- ▶ Improvements in computers and data storage devices.
- ▶ Powerful data capturing devices.
- ▶ Lots of data with potentially valuable information available.

## What is the Difference?

### Data characteristics

- ▶ Size
- ▶ Dimensionality
- ▶ Complexity
- ▶ Messy
- ▶ Secondary sources

### Focus on generalization performance

- ▶ Prediction on new data
- ▶ Action in new circumstances
- ▶ Complex models needed for good generalization.

### Computational considerations

- ▶ Large scale and complex systems

## Applications of Machine Learning

### Business applications

- ▶ Help companies intelligently find information
- ▶ Credit scoring
- ▶ Predict which products people are going to buy
- ▶ Recommender systems
- ▶ Autonomous trading

### Scientific applications

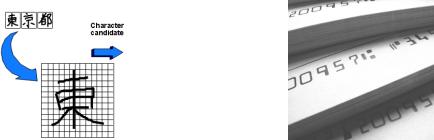
- ▶ Predict cancer occurrence/type and health of patients/personalized health
- ▶ Make sense of complex physical, biological, ecological, sociological models

9

11

## Applications of Machine Learning

### Pattern Recognition



- ▶ Sorting Cheques
- ▶ Reading License Plates
- ▶ Sorting Envelopes
- ▶ Eye/ Face/ Fingerprint Recognition

## Further Readings, News and Applications

Links are clickable in pdf. More recent news posted on course webpage.

- ▶ Leo Breiman: Statistical Modeling: The Two Cultures
- ▶ NY Times: R
- ▶ NY Times: Career in Statistics
- ▶ NY Times: Data Mining in Walmart
- ▶ NY Times: Big Data's Impact In the World
- ▶ Economist: Data, Data Everywhere
- ▶ McKinsey: Big data: The Next Frontier for Competition
- ▶ NY Times: Scientists See Promise in Deep-Learning Programs
- ▶ New Yorker: Is "Deep Learning" a Revolution in Artificial Intelligence?

10

12

## Types of Machine Learning

### Unsupervised Learning

Uncover structure hidden in ‘unlabelled’ data.

- ▶ Given network of social interactions, find communities.
- ▶ Given shopping habits for people using loyalty cards: find groups of ‘similar’ shoppers.
- ▶ Given expression measurements of 1000s of genes for 1000s of patients, find groups of functionally similar genes.

Goal: Hypothesis generation, visualization.

## Types of Machine Learning

### Semi-supervised Learning

A database of examples, only a small subset of which are labelled.

### Multi-task Learning

A database of examples, each of which has multiple labels corresponding to different prediction tasks.

### Reinforcement Learning

An agent acting in an environment, given rewards for performing appropriate actions, learns to maximize its reward.

13

15

## Types of Machine Learning

### Supervised Learning

A database of examples along with “labels” (task-specific).

- ▶ Given network of social interactions **along with their browsing habits**, predict what news might users find interesting.
- ▶ Given expression measurements of 1000s of genes for 1000s of patients **along with an indicator of absence or presence of a specific cancer**, predict if the cancer is present for a new patient.
- ▶ Given expression measurements of 1000s of genes for 1000s of patients **along with survival length**, predict survival time.

Goal: Prediction **on new examples**.

## OxWaSP

### Oxford-Warwick Centre for Doctoral Training in Statistics

- ▶ Programme aims to produce Europe's future research leaders in statistical methodology and computational statistics for modern applications.
- ▶ 10 fully-funded (UK, EU) students a year (1 international).
- ▶ Website for prospective students.
- ▶ **Deadline: January 24, 2014**

14

16

## Exploratory Data Analysis

### Notation

- ▶ Data consists of  $p$  measurements (variables/attributes) on  $n$  examples (observations/cases)
- ▶  $\mathbf{X}$  is a  $n \times p$ -matrix with  $\mathbf{X}_{ij} :=$  the  $j$ -th measurement for the  $i$ -th example

$$\mathbf{X} = \begin{bmatrix} x_{11} & x_{12} & \dots & x_{1j} & \dots & x_{1p} \\ x_{21} & x_{22} & \dots & x_{2j} & \dots & x_{2p} \\ \vdots & \vdots & \ddots & \vdots & \ddots & \vdots \\ x_{i1} & x_{i2} & \dots & x_{ij} & \dots & x_{ip} \\ \vdots & \vdots & \ddots & \vdots & \ddots & \vdots \\ x_{n1} & x_{n2} & \dots & x_{nj} & \dots & x_{np} \end{bmatrix}$$

- ▶ Denote the  $i$ th data item by  $x_i \in \mathbb{R}^p$ . (This is transpose of  $i$ th row of  $\mathbf{X}$ )
- ▶ Assume  $x_1, \dots, x_n$  are **independently and identically distributed** samples of a **random vector**  $\mathbf{X}$  over  $\mathbb{R}^p$ .

17

### Crabs Data I

```
## load package MASS containing the data
library(MASS)
## look at data
crabs

## assign predictor and class variables
Crabs <- crabs[,4:8]
Crabs.class <- factor(paste(crabs[,1],crabs[,2],sep=""))

## various plots
boxplot(Crabs)
hist(Crabs$FL,col='red',breaks=20,xname='Frontal Lobe Size (mm)')
hist(Crabs$RW,col='red',breaks=20,xname='Rear Width (mm)')
hist(Crabs$CL,col='red',breaks=20,xname='Carapace Length (mm)')
hist(Crabs$CW,col='red',breaks=20,xname='Carapace Width (mm)')
hist(Crabs$BD,col='red',breaks=20,xname='Body Depth (mm)')
plot(Crabs,col=unclass(Crabs.class))
parcoord(Crabs)
```

19

### Crabs Data ( $n = 200, p = 5$ )

Campbell (1974) studied rock crabs of the genus **leptograpsus**. One species, **L. variegatus**, had been split into two new species, previously grouped by colour, orange and blue. Preserved specimens lose their colour, so it was hoped that morphological differences would enable museum material to be classified.

Data are available on 50 specimens of each sex of each species, collected on sight at Fremantle, Western Australia. Each specimen has measurements on:

- ▶ the width of the frontal lobe  $FL$ ,
- ▶ the rear width  $RW$ ,
- ▶ the length along the carapace midline  $CL$ ,
- ▶ the maximum width  $CW$  of the carapace, and
- ▶ the body depth  $BD$  in mm.

in addition to colour (species) and sex.

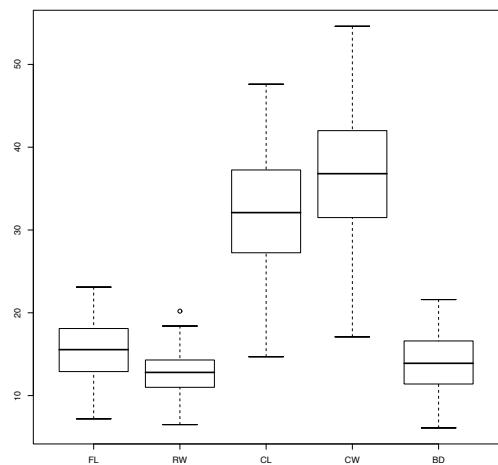
### Crabs data

	sp	sex	index	FL	RW	CL	CW	BD
1	B	M	1	8.1	6.7	16.1	19.0	7.0
2	B	M	2	8.8	7.7	18.1	20.8	7.4
3	B	M	3	9.2	7.8	19.0	22.4	7.7
4	B	M	4	9.6	7.9	20.1	23.1	8.2
5	B	M	5	9.8	8.0	20.3	23.0	8.2
6	B	M	6	10.8	9.0	23.0	26.5	9.8
7	B	M	7	11.1	9.9	23.8	27.1	9.8
8	B	M	8	11.6	9.1	24.5	28.4	10.4
9	B	M	9	11.8	9.6	24.2	27.8	9.7
10	B	M	10	11.8	10.5	25.2	29.3	10.3
11	B	M	11	12.2	10.8	27.3	31.6	10.9
12	B	M	12	12.3	11.0	26.8	31.5	11.4
13	B	M	13	12.6	10.0	27.7	31.7	11.4
14	B	M	14	12.8	10.2	27.2	31.8	10.9
15	B	M	15	12.8	10.9	27.4	31.5	11.0
16	B	M	16	12.9	11.0	26.8	30.9	11.4
17	B	M	17	13.1	10.6	28.2	32.3	11.0
18	B	M	18	13.1	10.9	28.3	32.4	11.2
19	B	M	19	13.3	11.1	27.8	32.3	11.3
20	B	M	20	13.9	11.1	29.2	33.3	12.1

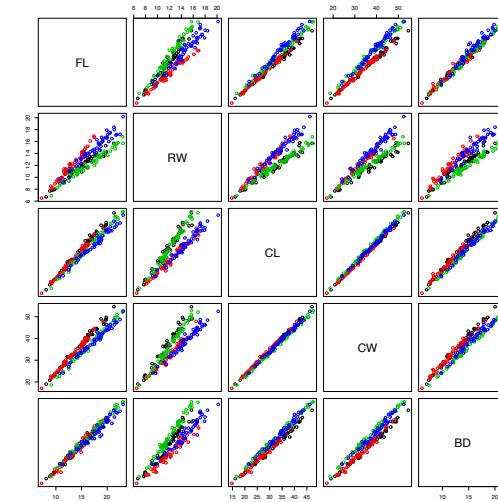
18

20

## Univariate Boxplots



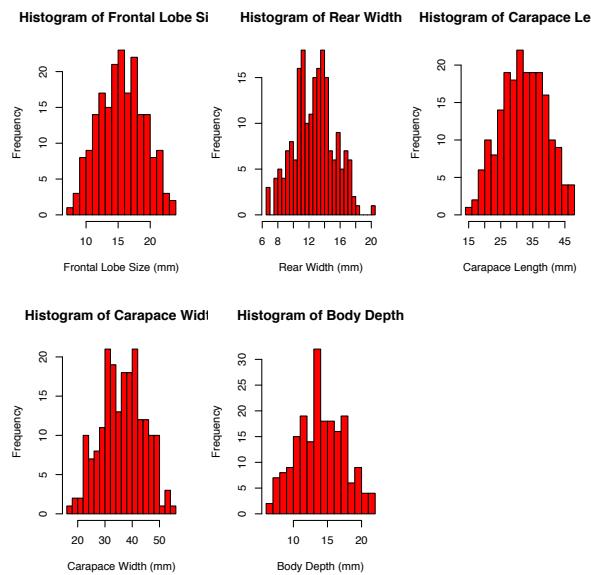
## Simple Pairwise Scatterplots



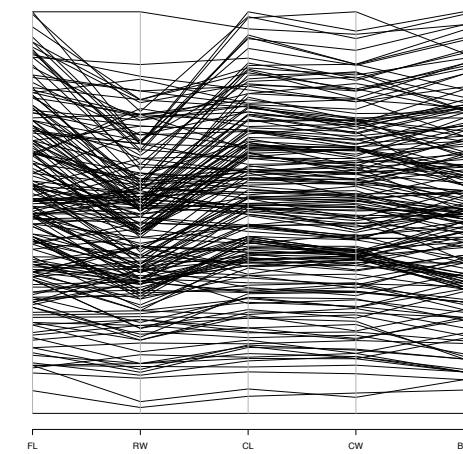
21

23

## Univariate Histograms



## Parallel Coordinate Plots



22

24

## Visualization and Dimensionality Reduction

These summary plots are helpful, but do not really help very much if the dimensionality of the data is high (a few dozen or thousands).

Visualizing higher-dimensional problems:

- ▶ We are constrained to view data in 2 or 3 dimensions
- ▶ Look for ‘interesting’ projections of  $\mathbf{X}$  into lower dimensions
- ▶ Hope that for large  $p$ , considering only  $k \ll p$  dimensions is just as informative.

### Dimensionality reduction

- ▶ For each data item  $x_i \in \mathbb{R}^p$ , find a lower dimensional representation  $z_i \in \mathbb{R}^k$  with  $k \ll p$ .
- ▶ Preserve as much as possible the interesting statistical properties/relationships of data items.

25

## Principal Components Analysis (PCA)

- ▶ Our data set is an iid sample of a random vector  $\mathbf{X} = [X_1 \dots X_p]^\top$ .
- ▶ For the  $1^{st}$  PC, we seek a derived variable of the form

$$Z_1 = v_{11}X_1 + v_{12}X_2 + \dots + v_{1p}X_p = v_1^\top \mathbf{X}$$

where  $v_1 = [v_{11}, \dots, v_{1p}]^\top \in \mathbb{R}^p$  are chosen to maximise  $\text{Var}(Z_1)$ .

To get a well defined problem, we fix

$$v_1^\top v_1 = 1.$$

- ▶ The  $2^{nd}$  PC is chosen to be orthogonal with the  $1^{st}$  and is computed in a similar way. It will have the largest variance in the remaining  $p - 1$  dimensions, etc.

27

## Principal Components Analysis (PCA)

- ▶ PCA considers interesting directions to be those with greatest **variance**.
- ▶ A **linear** dimensionality reduction technique:
- ▶ Finds an orthogonal basis  $v_1, v_2, \dots, v_p$  for the data space such that
  - ▶ The first principal component (PC)  $v_1$  is the direction of greatest variance of data.
  - ▶ The second PC  $v_2$  is the direction orthogonal to  $v_1$  of greatest variance, etc.
  - ▶ The subspace spanned by the first  $k$  PCs represents the ‘best’  $k$ -dimensional representation of the data.
  - ▶ The  $k$ -dimensional representation of  $x_i$  is:

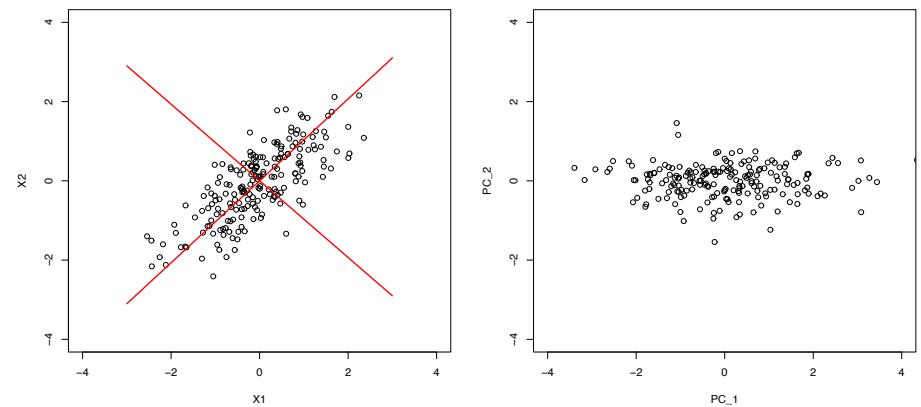
$$z_i = V^\top x_i = \sum_{\ell=1}^k v_\ell^\top x_i$$

where  $V \in \mathbb{R}^{p \times k}$ .

- ▶ For simplicity, we will assume from now on that our dataset is centred, i.e. we subtract the average  $\bar{x}$  from each  $x_i$ .

26

## Principal Components Analysis (PCA)



28

## Deriving the First Principal Component

- Maximise, subject to  $v_1^\top v_1 = 1$ :

$$\text{Var}(Z_1) = \text{Var}(v_1^\top X) = v_1^\top \text{Cov}(X)v_1 \approx v_1^\top S v_1$$

where  $S \in \mathbb{R}^{p \times p}$  is the sample covariance matrix, i.e.

$$S = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})(x_i - \bar{x})^\top = \frac{1}{n-1} \sum_{i=1}^n x_i x_i^\top = \frac{1}{n-1} \mathbf{X}^\top \mathbf{X}.$$

- Rewriting this as a constrained maximisation problem,

$$\mathcal{L}(v_1, \lambda_1) = v_1^\top S v_1 - \lambda_1 (v_1^\top v_1 - 1).$$

- The corresponding vector of partial derivatives yields

$$\frac{\partial \mathcal{L}(v_1, \lambda_1)}{\partial v_1} = 2Sv_1 - 2\lambda_1 v_1.$$

- Setting this to zero reveals the eigenvector equation, i.e.  $v_1$  must be an eigenvector of  $S$  and  $\lambda_1$  the corresponding eigenvalue.
- Since  $v_1^\top S v_1 = \lambda_1 v_1^\top v_1 = \lambda_1$ , the 1<sup>st</sup> PC must be the eigenvector associated with the largest eigenvalue of  $S$ .

29

## Properties of the Principal Components

- PCs are **uncorrelated**

$$\text{Cov}(X^\top v_i, X^\top v_j) \approx v_i^\top S v_j = 0 \text{ for } i \neq j.$$

- The **total sample variance** is given by

$$\sum_{i=1}^p S_{ii} = \lambda_1 + \dots + \lambda_p,$$

so the **proportion of total variance** explained by the  $k^{\text{th}}$  PC is

$$\frac{\lambda_k}{\lambda_1 + \lambda_2 + \dots + \lambda_p} \quad k = 1, 2, \dots, p$$

- $S$  is a real symmetric matrix, so eigenvectors (principal components) are orthogonal.
- Derived variables  $Z_1, \dots, Z_p$  have variances  $\lambda_1, \dots, \lambda_p$ .

31

## Deriving Subsequent Principal Components

- Proceed as before but include the additional constraint that the 2<sup>nd</sup> PC must be orthogonal to the 1<sup>st</sup> PC:

$$\mathcal{L}(v_2, \lambda_2, \mu) = v_2^\top S v_2 - \lambda_2 (v_2^\top v_2 - 1) - \mu (v_1^\top v_2).$$

- Solving this shows that  $v_2$  must be the eigenvector of  $S$  associated with the 2<sup>nd</sup> largest eigenvalue, and so on
- The eigenvalue decomposition of  $S$  is given by

$$S = V \Lambda V^\top$$

where  $\Lambda$  is a diagonal matrix with eigenvalues

$$\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_p \geq 0$$

and  $V$  is a  $p \times p$  orthogonal matrix whose columns are the  $p$  eigenvectors of  $S$ , i.e. the principal components  $v_1, \dots, v_p$ .

30

## R code

This is what we have had before:

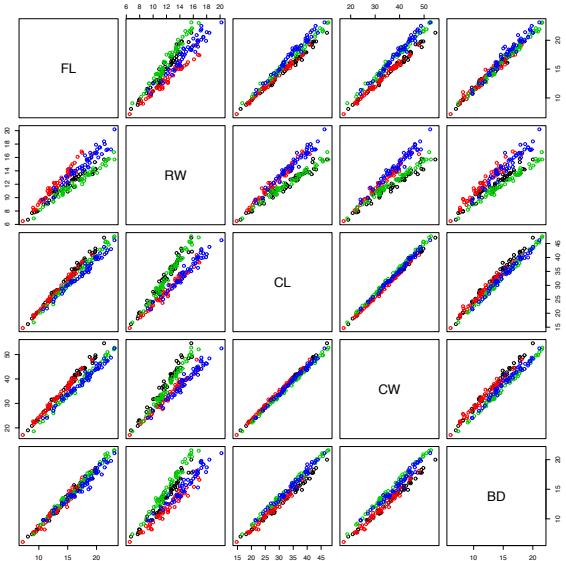
```
library(MASS)
Crabs <- crabs[, 4:8]
Crabs.class <- factor(paste(crabs[, 1], crabs[, 2], sep=""))
plot(Crabs, col=unclass(Crabs.class))
```

Now perform PCA with function `princomp`. (Alternatively, solve for the PCs yourself using `eigen` or `svd`).

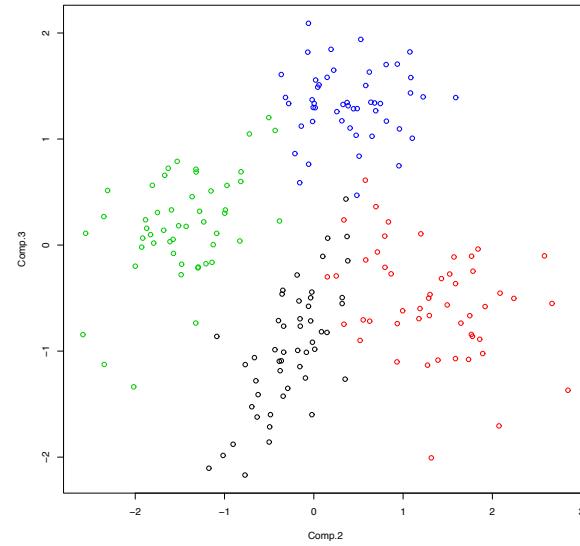
```
Crabs.pca <- princomp(Crabs, cor=FALSE)
plot(Crabs.pca)
pairs(predict(Crabs.pca), col=unclass(Crabs.class))
```

32

## Original Crabs Data



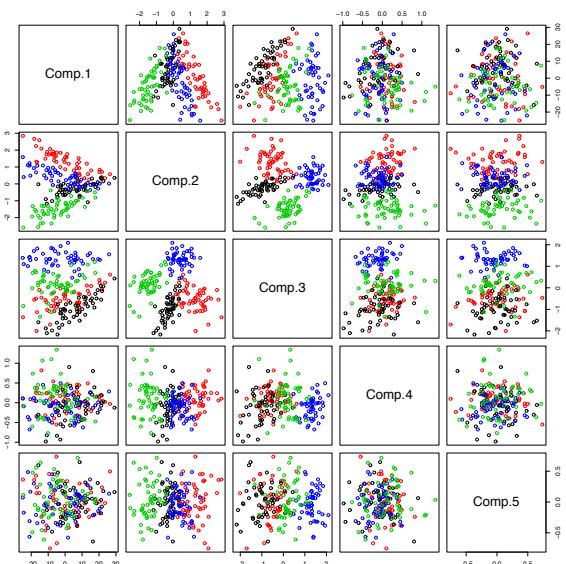
## PC 2 vs PC 3



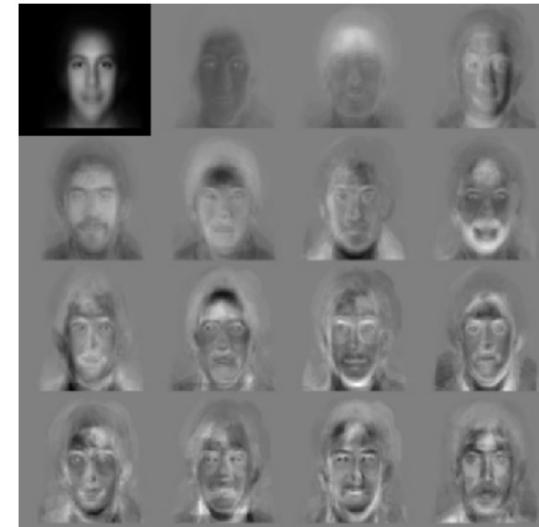
33

35

## PCA of Crabs Data



## PCA on Face Images

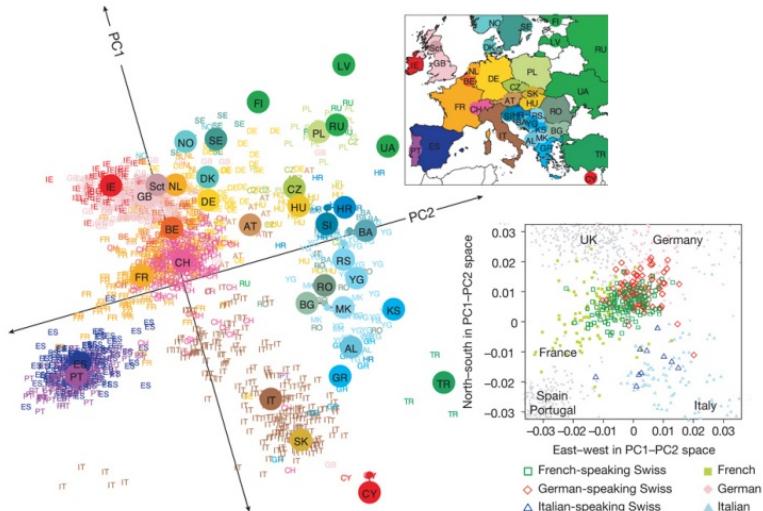


34

<http://vismod.media.mit.edu/vismod/demos/facerec/basic.html>

36

## PCA on European Genetic Variation



<http://www.nature.com/nature/journal/v456/n7218/full/nature07331.html>

37

## Eigenvalue Decomposition (EVD)

Eigenvalue decomposition plays a significant role in PCA. PCs are eigenvectors of  $S = \frac{1}{n-1} \mathbf{X}^T \mathbf{X}$  and PCA properties are derived from those of eigenvectors and eigenvalues.

- ▶ For any  $p \times p$  **symmetric** matrix  $S$ , there exists  $p$  eigenvectors  $v_1, \dots, v_p$  that are pairwise orthogonal and  $p$  associated eigenvalues  $\lambda_1, \dots, \lambda_p$  which satisfy the eigenvalue equation  $Sv_i = \lambda_i v_i \forall i$ .
- ▶  $S$  can be written as  $S = V\Lambda V^T$  where
  - ▶  $V = [v_1, \dots, v_p]$  is a  $p \times p$  orthogonal matrix
  - ▶  $\Lambda = \text{diag}\{\lambda_1, \dots, \lambda_p\}$
  - ▶ If  $S$  is a real-valued matrix, then the eigenvalues are real-valued as well,  $\lambda_i \in \mathbb{R} \forall i$
- ▶ To compute the PCA of a dataset  $\mathbf{X}$ , we can:
  - ▶ First estimate the covariance matrix using the sample covariance  $S$ .
  - ▶ Compute the EVD of  $S$  using the R command `eigen`.

39

## Comments on the use of PCA

- ▶ PCA commonly used to project data  $\mathbf{X}$  onto the first  $k$  PCs giving the  $k$ -dimensional view of the data that best preserves the first two moments.
- ▶ Although PCs are uncorrelated, scatterplots sometimes reveal structures in the data other than linear correlation.
- ▶ PCA commonly used for lossy compression of high dimensional data.
- ▶ Emphasis on variance is where the weaknesses of PCA stem from:
  - ▶ The PCs depend heavily on the units measurement. Where the data matrix contains measurements of vastly differing orders of magnitude, the PC will be greatly biased in the direction of larger measurement. It is therefore recommended to calculate PCs from  $\text{Corr}(\mathbf{X})$  instead of  $\text{Cov}(\mathbf{X})$ .
  - ▶ Robustness to outliers is also an issue. Variance is affected by outliers therefore so are PCs.

## Singular Value Decomposition (SVD)

Though the EVD does not always exist, the singular value decomposition is another matrix factorization technique that **always** exist, even for non-square matrices.

- ▶  $\mathbf{X}$  can be written as  $\mathbf{X} = \mathbf{U}\mathbf{D}\mathbf{V}^T$  where
  - ▶  $\mathbf{U}$  is an  $n \times n$  matrix with orthogonal columns.
  - ▶  $\mathbf{D}$  is a  $n \times p$  matrix with decreasing non-negative elements on the diagonal (the singular values) and zero off-diagonal elements.
  - ▶  $\mathbf{V}$  is a  $p \times p$  matrix with orthogonal columns.
- ▶ SVD can be computed using very fast and numerically stable algorithms. The relevant R command is `svd`.

38

40

## Some Properties of the SVD

- Let  $\mathbf{X} = \mathbf{UDV}^\top$  be the SVD of the  $n \times p$  data matrix  $\mathbf{X}$ .
- Note that

$$(n-1)\mathbf{S} = \mathbf{X}^\top \mathbf{X} = (\mathbf{UDV}^\top)^\top (\mathbf{UDV}^\top) = \mathbf{V}\mathbf{D}^\top \mathbf{U}^\top \mathbf{U}\mathbf{D}\mathbf{V}^\top = \mathbf{V}\mathbf{D}^\top \mathbf{D}\mathbf{V}^\top,$$

using orthogonality ( $\mathbf{U}^\top \mathbf{U} = \mathbf{I}_n$ ) of  $\mathbf{U}$ .

- The eigenvalues of  $\mathbf{S}$  are thus the diagonal entries of  $\frac{1}{n-1}\mathbf{D}^2$  and the columns of the orthogonal matrix  $\mathbf{V}$  are the eigenvectors of  $\mathbf{S}$ .
- We also have

$$\mathbf{XX}^\top = (\mathbf{UDV}^\top)(\mathbf{UDV}^\top)^\top = \mathbf{UDV}^\top \mathbf{V}\mathbf{D}^\top \mathbf{U}^\top = \mathbf{UDD}^\top \mathbf{U}^\top,$$

using orthogonality ( $\mathbf{V}^\top \mathbf{V} = \mathbf{I}_p$ ) of  $\mathbf{V}$ .

- SVD also gives the optimal low-rank approximations of  $\mathbf{X}$ :

$$\min_{\tilde{\mathbf{X}}} \|\tilde{\mathbf{X}} - \mathbf{X}\|^2 \quad \text{s.t. } \tilde{\mathbf{X}} \text{ has maximum rank } r < n, p.$$

This problem can be solved by keeping only the  $r$  largest singular values of  $\mathbf{X}$ , zeroing out the smaller singular values in the SVD.

41

43

## Biplots

Recall that  $\mathbf{X} = [X_1, \dots, X_p]^\top$  and  $\mathbf{X} = \mathbf{UDV}^\top$  is the SVD of the data matrix.

- The PC projection of  $x_i$  is:

$$z_i = \mathbf{V}^\top \mathbf{x}_i = \mathbf{D}\mathbf{U}_i^\top = [D_{11}U_{i1}, \dots, D_{kk}U_{ik}]^\top.$$

- The  $j$ th unit vector  $\mathbf{e}_j \in \mathbb{R}^p$  points in the direction of  $X_j$ . Its PC projection is  $\mathbf{V}_j^\top = \mathbf{V}^\top \mathbf{e}_j$ , the  $j$ th row of  $\mathbf{V}$ .
- The projection of the variable indicates the weighting each PC gives to the original variables.
- Dot products between the projections gives entries of the data matrix:

$$x_{ij} = \sum_{k=1}^p U_{ik} D_{kk} V_{jk} = \langle \mathbf{D}\mathbf{U}_i^\top, \mathbf{V}_j^\top \rangle.$$

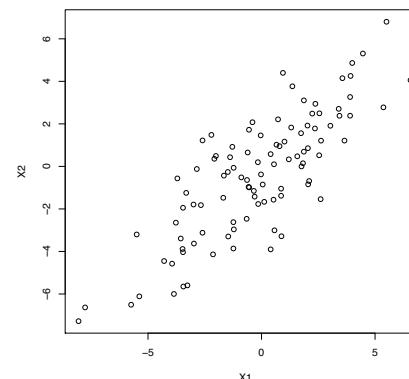
- Distance of projected points from projected variables gives original location.
- These relationships can be plotted in 2D by focussing on first two PCs.

## Biplots

- PCA plots show the data items (as rows of  $\mathbf{X}$ ) in the PC space.
- Biplots** allow us to visualize the **original variables** (as columns  $\mathbf{X}$ ) in the same plot.
- As for PCA, we would like the geometry of the plot to preserve as much of the covariance structure as possible.

42

## Biplots



44

## Biplots

- There are other projections we can consider for biplots:

$$x_{ij} = \sum_{k=1}^p U_{ik} D_{kk} V_{jk} = \langle D U_i^\top, V_j^\top \rangle = \langle D^{1-\alpha} U_i^\top, D^\alpha V_j^\top \rangle.$$

where  $0 \leq \alpha \leq 1$ . The  $\alpha = 1$  case has some nice properties.

- Covariance of the projected points is:

$$\frac{1}{n-1} \sum_{i=1}^n U_i^\top U_i = \frac{1}{n-1} I.$$

Projected points are uncorrelated and dimensions are equi-variance.

- The covariance between  $X_j$  and  $X_\ell$  is:

$$\text{Var}(X_j X_\ell) = \frac{1}{n-1} \langle D V_j^\top, D V_\ell^\top \rangle$$

So the angle between the projected variables gives the correlation.

- When using  $k < p$  PCs, quality depends on the proportion of variance explained by the PCs.

## Iris Data

50 sample from 3 species of iris: *iris setosa*, *versicolor*, and *virginica*

Each measuring the length and widths of both sepal and petals

Collected by E. Anderson (1935) and analysed by R.A. Fisher (1936)



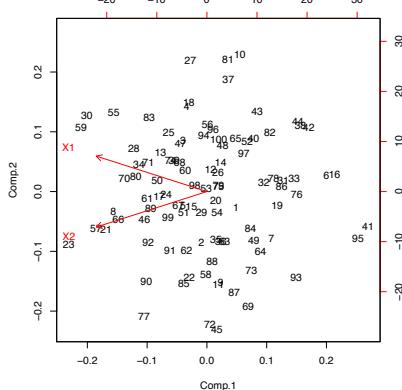
Using again function `princomp` and `biplot`.

```
iris1 <- iris
iris1 <- iris1[,-5]
biplot(princomp(iris1, cor=T))
```

45

47

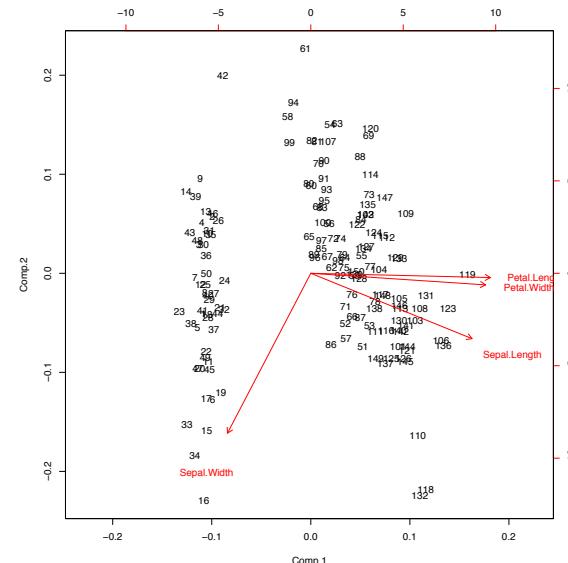
## Biplots



```
pc <- princomp(x)
biplot(pc, scale=0)
biplot(pc, scale=1)
```

46

## Iris Data



48

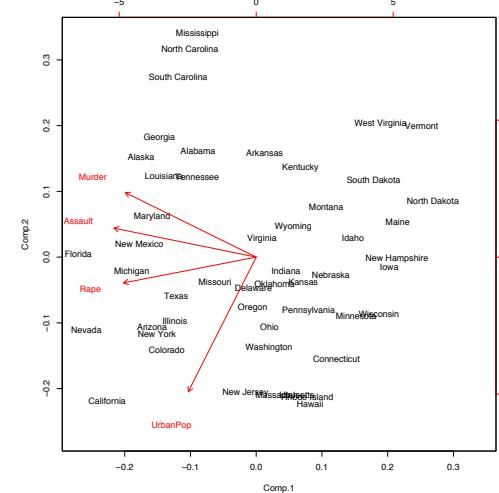
## US Arrests Data

This data set contains statistics, in arrests per 100,000 residents for assault, murder, and rape in each of the 50 US states in 1973. Also given is the percent of the population living in urban areas.

```
pairs(USArrests)
usarrests.pca <- princomp(USArrests, cor=T)
plot(usarrests.pca)

pairs(predict(usarrests.pca))
biplot(usarrests.pca)
```

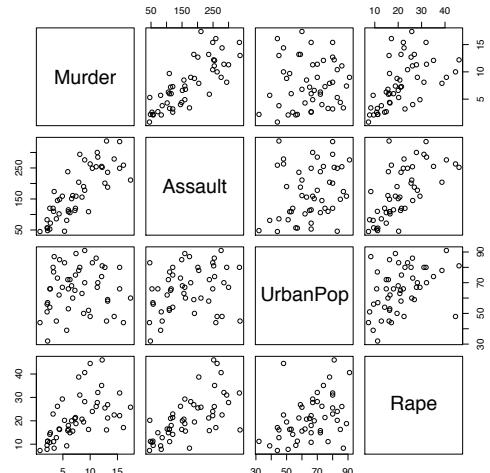
## US Arrests Data Biplot



49

51

## US Arrests Data Pairs Plot

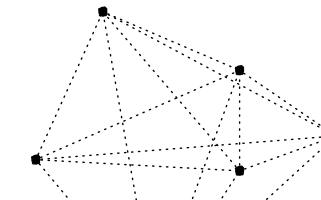


49

## Multidimensional Scaling

Suppose there are  $n$  points  $\mathbf{X}$  in  $\mathbb{R}^p$ , but we are only given the  $n \times n$  matrix  $\mathbf{D}$  of inter-point distances.

Can we reconstruct  $\mathbf{X}$ ?



50

52

## Multidimensional Scaling

Rigid transformations (translations, rotations and reflections) do not change inter-point distances so cannot recover  $\mathbf{X}$  exactly. However  $\mathbf{X}$  can be recovered up to these transformations!

- Let  $d_{ij} = \|\mathbf{x}_i - \mathbf{x}_j\|_2$  be the distance between points  $\mathbf{x}_i$  and  $\mathbf{x}_j$ .

$$\begin{aligned} d_{ij}^2 &= \|\mathbf{x}_i - \mathbf{x}_j\|_2^2 \\ &= (\mathbf{x}_i - \mathbf{x}_j)^\top (\mathbf{x}_i - \mathbf{x}_j) \\ &= \mathbf{x}_i^\top \mathbf{x}_i + \mathbf{x}_j^\top \mathbf{x}_j - 2\mathbf{x}_i^\top \mathbf{x}_j \end{aligned}$$

- Let  $\mathbf{B} = \mathbf{XX}^\top$  be the  $n \times n$  matrix of dot-products,  $b_{ij} = \mathbf{x}_i^\top \mathbf{x}_j$ . The above shows that  $\mathbf{D}$  can be computed from  $\mathbf{B}$ .
- Some algebraic exercise shows that  $\mathbf{B}$  can be recovered from  $\mathbf{D}$  if we assume  $\sum_{i=1}^n \mathbf{x}_i = \mathbf{0}$ .

53

## US City Flight Distances

We present a table of flying mileages between 10 American cities, distances calculated from our 2-dimensional world. Using  $\mathbf{D}$  as the starting point, metric MDS finds a configuration with the same distance matrix.

	ATLA	CHIG	DENV	HOUS	LA	MIAM	NY	SF	SEAT	DC
0	587	1212	701	1936	604	748	2139	2182	543	
587	0	920	940	1745	1188	713	1858	1737	597	
1212	920	0	879	831	1726	1631	949	1021	1494	
701	940	879	0	1374	968	1420	1645	1891	1220	
1936	1745	831	1374	0	2339	2451	347	959	2300	
604	1188	1726	968	2339	0	1092	2594	2734	923	
748	713	1631	1420	2451	1092	0	2571	2408	205	
2139	1858	949	1645	347	2594	2571	0	678	2442	
2182	1737	1021	1891	959	2734	2408	678	0	2329	
543	597	1494	1220	2300	923	205	2442	2329	0	

55

## US City Flight Distances

```
library(MASS)

us <- read.csv("http://www.stats.ox.ac.uk/
~teh/teaching/smldm/data/uscities.csv")

## use classical MDS to find lower dimensional views of the data
## recover X in 2 dimensions

us.classical <- cmdscale(d=us, k=2)

plot(us.classical)
text(us.classical, labels=names(us))
```

## Multidimensional Scaling

- If we knew  $\mathbf{X}$ , then SVD gives  $\mathbf{X} = \mathbf{UDV}^\top$ . As  $\mathbf{X}$  has rank  $k = \min(n, p)$ , we have at most  $k$  singular values in  $\mathbf{D}$  and we can assume  $\mathbf{U} \in \mathbb{R}^{n \times k}$ ,  $\mathbf{D} \in \mathbb{R}^{k \times p}$  and  $\mathbf{V} \in \mathbb{R}^{p \times p}$ .
- The eigendecomposition of  $\mathbf{B}$  is then:

$$\mathbf{B} = \mathbf{XX}^\top = \mathbf{UDV}^\top \mathbf{UDV}^\top = \mathbf{U}\Lambda\mathbf{U}^\top.$$

- This eigendecomposition can be obtained from  $\mathbf{B}$  without knowledge of  $\mathbf{X}$ !
- Let  $\tilde{\mathbf{x}}_i^\top = \mathbf{U}_i \Lambda^{\frac{1}{2}}$  be the  $i$ th row of  $\mathbf{U}\Lambda^{\frac{1}{2}}$ . Pad  $\tilde{\mathbf{x}}_i$  with 0s so that it has length  $p$ .

$$\tilde{\mathbf{x}}_i^\top \tilde{\mathbf{x}}_j = \mathbf{U}_i \Lambda^{\frac{1}{2}} \mathbf{U}_j^\top = b_{ij} = \mathbf{x}_i^\top \mathbf{x}_j$$

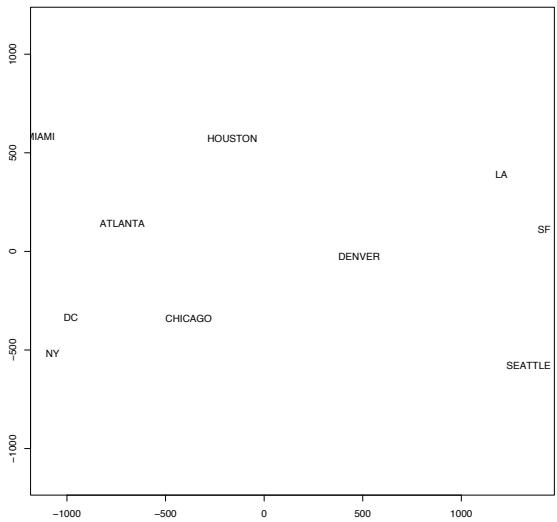
and we have found a set of vectors with dot-products given by  $\mathbf{B}$ .

- The vectors  $\tilde{\mathbf{x}}_i$  differs from  $\mathbf{x}_i$  only via the orthogonal matrix  $\mathbf{V}$  so are equivalent up to rotation and reflections.

54

56

## US City Flight Distances

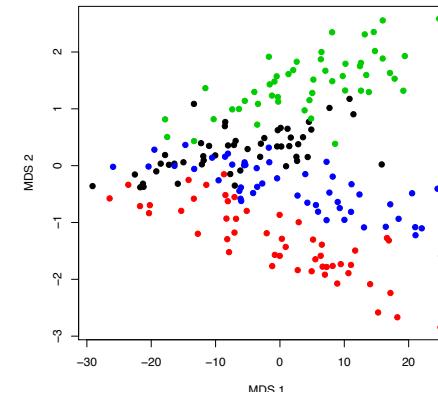


57

## Crabs Data

```
library(MASS)
Crabs <- crabs[,4:8]
Crabs.class <- factor(paste(crabs[,1],crabs[,2],sep=""))

crabsmds <- cmdscale(d= dist(Crabs),k=2)
plot(crabsmds, pch=20, cex=2,col=unclass(Crabs.class))
```



59

## Lower-dimensional Reconstructions

In classical MDS derivation, we used all eigenvalues in the eigendecomposition of  $\mathbf{B}$  to reconstruct

$$\tilde{x}_i = U_i \Lambda^{\frac{1}{2}}.$$

We can use only the largest  $k < \min(n, p)$  eigenvalues and eigenvectors in the reconstruction, giving the ‘best’  $k$ -dimensional view of the data.

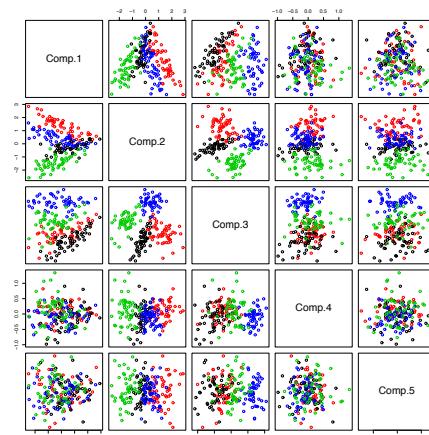
This is analogous to PCA, where only the largest eigenvalues of  $\mathbf{X}^\top \mathbf{X}$  are used, and the smallest ones effectively suppressed.

Indeed, PCA and classical MDS are duals and yield effectively the same result.

58

## Crabs Data

Compare with previous PCA analysis.  
Classical MDS solution corresponds to the first 2 PCs.



60

## Example: Language data

Presence or absence of 2867 homologous traits in 87 Indo-European languages.

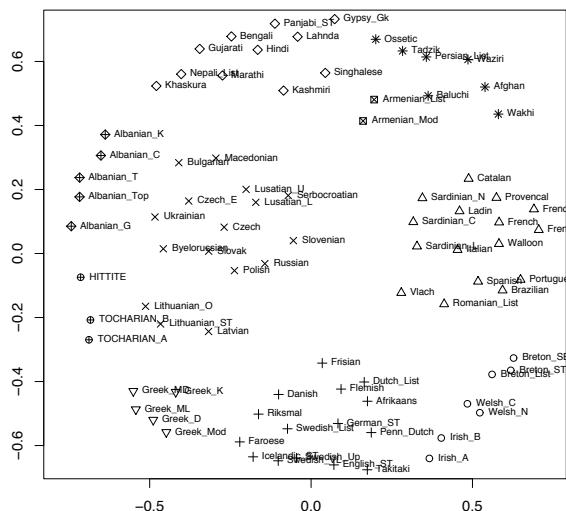
> X[1:15, 1:16]

	V1	V2	V3	V4	V5	V6	V7	V8	V9	V10	V11	V12	V13
Irish_A	0	0	0	0	1	0	0	0	0	0	0	0	0
Irish_B	0	0	0	0	1	0	0	0	0	0	0	0	0
Welsh_N	0	0	0	1	0	0	0	0	0	0	0	0	0
Welsh_C	0	0	0	1	0	0	0	0	0	0	0	0	0
Breton_List	0	0	0	0	1	0	0	0	0	0	0	0	0
Breton_SE	0	0	0	0	1	0	0	0	0	0	0	0	0
Breton_ST	0	0	0	0	1	0	0	0	0	0	0	0	0
Romanian_List	0	1	0	0	0	0	0	0	0	0	0	0	0
Vlach	0	1	0	0	0	0	0	0	0	0	0	0	0
Italian	0	1	0	0	0	0	0	0	0	0	0	0	0
Ladin	0	1	0	0	0	0	0	0	0	0	0	0	0
Provencal	0	1	0	0	0	0	0	0	0	0	0	0	0
French	0	1	0	0	0	0	0	0	0	0	0	0	0
Walloon	0	1	0	0	0	0	0	0	0	0	0	0	0
French_Creole_C	0	1	0	0	0	0	0	0	0	0	0	0	0

61

## Example: Language data

Using MDS with non-metric scaling.



62

## Varieties of MDS

Generally, MDS is a class of dimensionality reduction techniques which represents data points  $x_1, \dots, x_n \in \mathbb{R}^p$  in a lower-dimensional space  $z_1, \dots, z_n \in \mathbb{R}^k$  which tries to preserve inter-point (dis)similarities.

- It requires only the matrix  $\mathbf{D}$  of pairwise dissimilarities

$$d_{ij} = d(x_i, d_j).$$

For example we can use Euclidean distance  $d_{ij} = \|x_i - x_j\|_2$ . Other dissimilarities are possible. Conversely, it can use a matrix of similarities.

- MDS finds representations  $z_1, \dots, z_n \in \mathbb{R}^k$  such that

$$d(x_i, x_j) \approx \tilde{d}_{ij} = \tilde{d}(z_i, z_j),$$

where  $\tilde{d}$  represents dissimilarity in the reduced  $k$ -dimensional space, and differences in dissimilarities are measured by a **stress function**  $S(d_{ij}, \tilde{d}_{ij})$ .

63

## Varieties of MDS

Choices of (dis)similarities and stress functions lead to different objective functions and different algorithms.

- Classical - preserves similarities instead

$$S(\mathbf{Z}) = \sum_{i \neq j} (s_{ij} - \langle z_i - \bar{z}, z_j - \bar{z} \rangle)^2$$

- Metric Shepard-Kruskal

$$S(\mathbf{Z}) = \sum_{i \neq j} (d_{ij} - \|z_i - z_j\|_2)^2$$

- Sammon - preserves shorter distances more

$$S(\mathbf{Z}) = \sum_{i \neq j} \frac{(d_{ij} - \|z_i - z_j\|_2)^2}{d_{ij}}$$

- Non-Metric Shepard-Kruskal - ignores actual distance values, only ranks

$$S(\mathbf{Z}) = \min_{g \text{ increasing}} \sum_{i \neq j} (g(d_{ij}) - \|z_i - z_j\|_2)^2$$

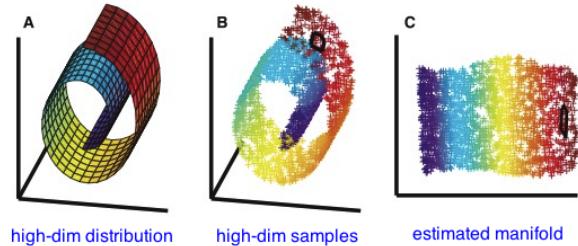
64

## Nonlinear Dimensionality Reduction

Two aims of different varieties of MDS:

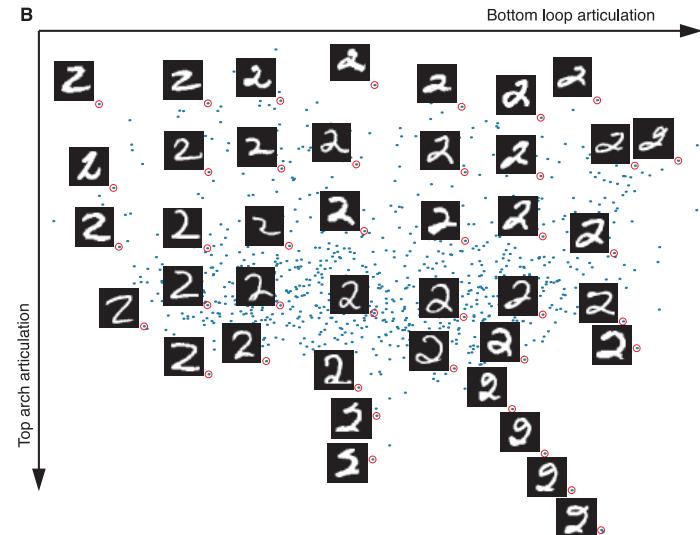
- ▶ To visualize the (dis)similarities among items in a dataset, where these (dis)similarities may not have Euclidean geometric interpretations.
- ▶ To perform **nonlinear** dimensionality reduction.

Many high-dimensional datasets exhibit low-dimensional structure (“live on a low-dimensional manifold”).



65

## Handwritten Characters

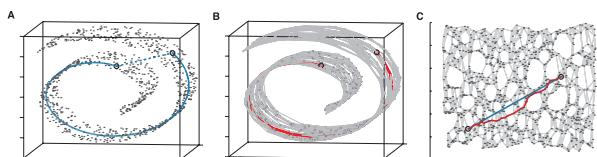


67

## Isomap

Isomap is a non-linear dimensional reduction technique based on classical MDS. Differs from other MDSs in its estimate of distances  $d_{ij}$ .

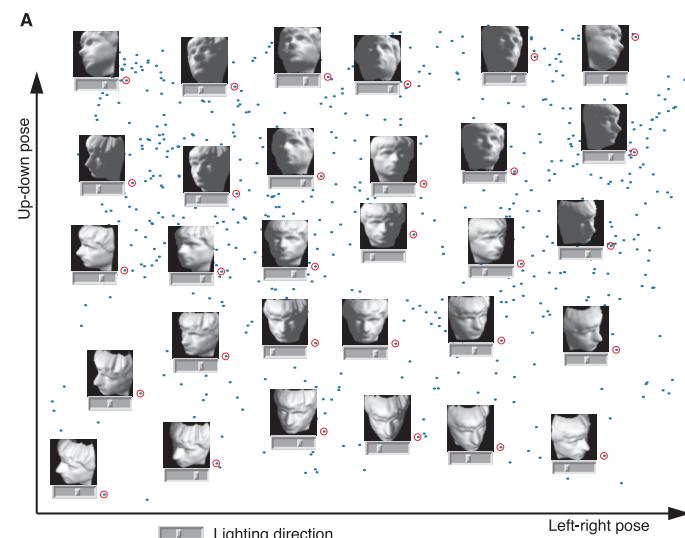
1. Calculate distances  $d_{ij}$  for  $i, j = 1, \dots, n$  between all data points, using the Euclidean distance.
2. Form a graph  $G$  with the  $n$  samples as nodes, and edges between the respective  $K$  nearest neighbours.
3. Replace distances  $d_{ij}$  by shortest-path distance on graph  $d_{ij}^G$  and perform classical MDS, using these distances.



Examples from Tenenbaum et al. (2000).

66

## Faces



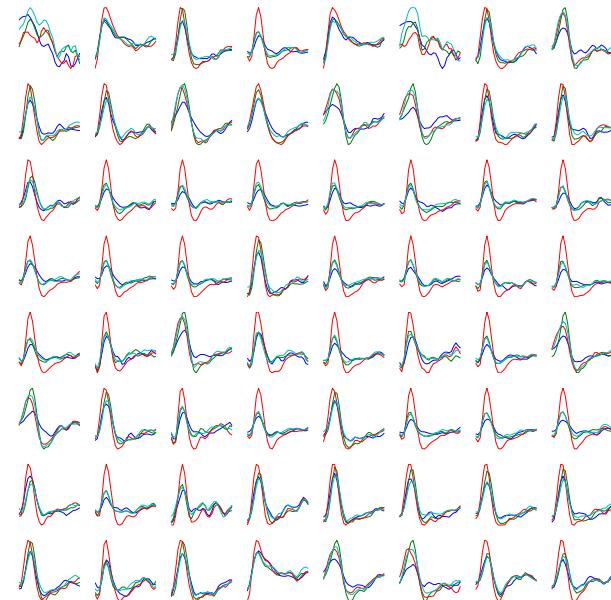
68

## Other Nonlinear Dimensionality Reduction Techniques

- ▶ Locally Linear Embedding.
- ▶ Laplacian Eigenmaps.
- ▶ Maximum Variance Unfolding.

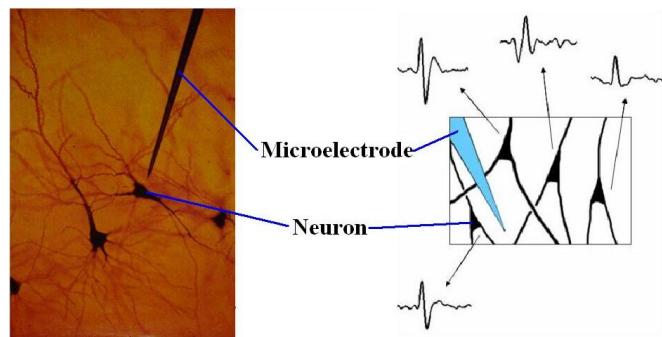
69

## Neural Spike Waveforms



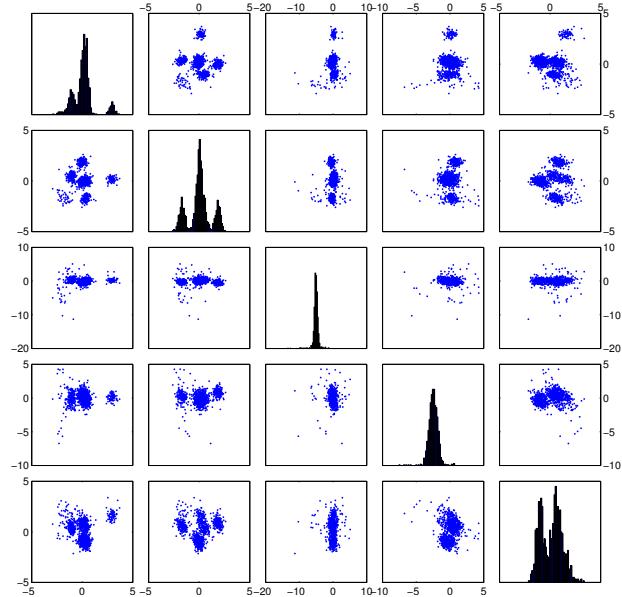
71

## Neural Electroencephalography (EEG)



70

## Pairs Plot of Principal Components



72

## Clustering

- ▶ Many datasets consist of multiple heterogeneous subsets. Cluster analysis is a range of methods that reveal this heterogeneity by discovering clusters of similar points.
- ▶ Model-based clustering:
  - ▶ Each cluster is described using a probability model.
- ▶ Model-free clustering:
  - ▶ Defined by similarity among points within clusters (dissimilarity among points between clusters).
- ▶ Partition-based clustering methods:
  - ▶ Allocate points into  $K$  clusters.
  - ▶ The number of cluster is usually fixed beforehand or investigated for various values of  $K$  as part of the analysis.
- ▶ Hierarchy-based clustering methods:
  - ▶ Allocate points into clusters and clusters into super-clusters forming a hierarchy.
  - ▶ Typically the hierarchy forms a binary tree (a dendrogram) where each cluster has two “children”.

73

## Measuring Dissimilarity

To find hierarchical clusters, we need some way to measure the dissimilarity between clusters

- ▶ Given two points  $x_i$  and  $x_j$ , it is straightforward to measure their dissimilarity, say  $d(x_i, x_j) = \|x_i - x_j\|_2$ .
- ▶ It is unclear however how to extend this to measure dissimilarity between clusters,  $D(C_i, C_j)$  for clusters  $C_i$  and  $C_j$ .

Many such proposals though no consensus as to which is best.

### (a) Single Linkage

$$D(C_i, C_j) = \min_{x,y} (d(x,y) | x \in C_i, y \in C_j)$$

### (b) Complete Linkage

$$D(C_i, C_j) = \max_{x,y} (d(x,y) | x \in C_i, y \in C_j)$$

### (c) Average Linkage

$$D(C_i, C_j) = \text{avg}_{x,y} (d(x,y) | x \in C_i, y \in C_j)$$

75

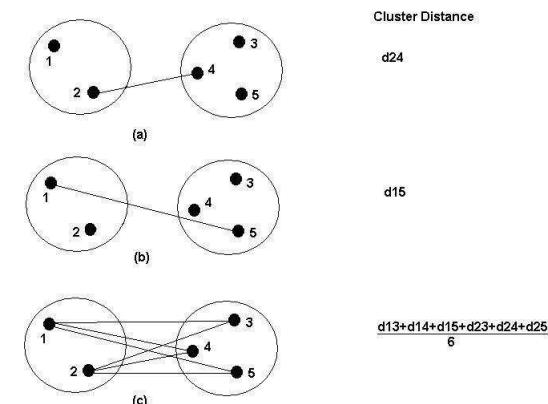
## Hierarchical Clustering

- ▶ Hierarchically structured data can be found everywhere (measurements of different species and different individuals within species), hierarchical methods attempt to understand data by looking for clusters.
- ▶ There are two general strategies for generating hierarchical clusters. Both proceed by seeking to minimize some measure of dissimilarity.
  - ▶ Agglomerative / Bottom-Up / Merging
  - ▶ Divisive / Top-Down / Splitting

**Hierarchical clusters** are generated where at each level, clusters are created by merging clusters at lower levels. This process can easily be viewed by a dendrogram/tree.

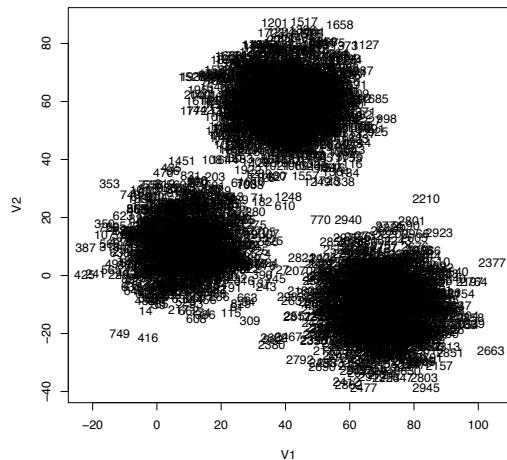
74

## Measuring Dissimilarity



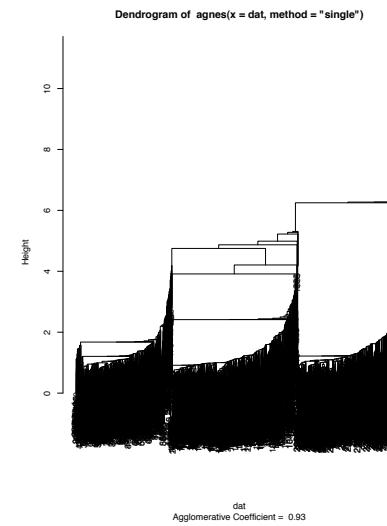
76

## Hierarchical Clustering on Artificial Dataset



77

## Hierarchical Clustering on Artificial Dataset



79

## Hierarchical Clustering on Artificial Dataset

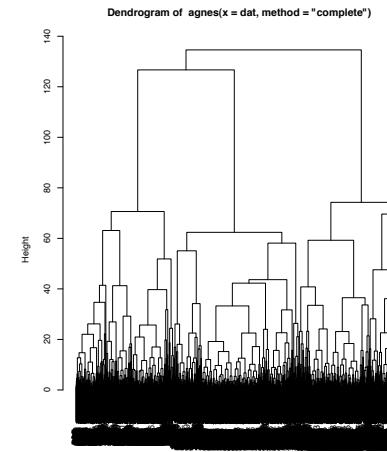
```
#start afresh
dat=xclara #3000 x 2
library(cluster)

#plot the data
plot(dat,type="n")
text(dat,labels=row.names(dat))

plot(agnes(dat,method="single"))
plot(agnes(dat,method="complete"))
plot(agnes(dat,method="average"))
```

78

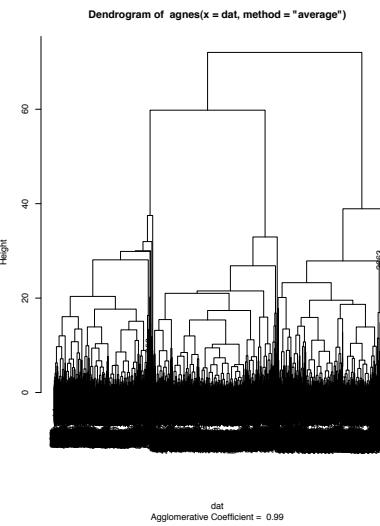
## Hierarchical Clustering on Artificial Dataset



80

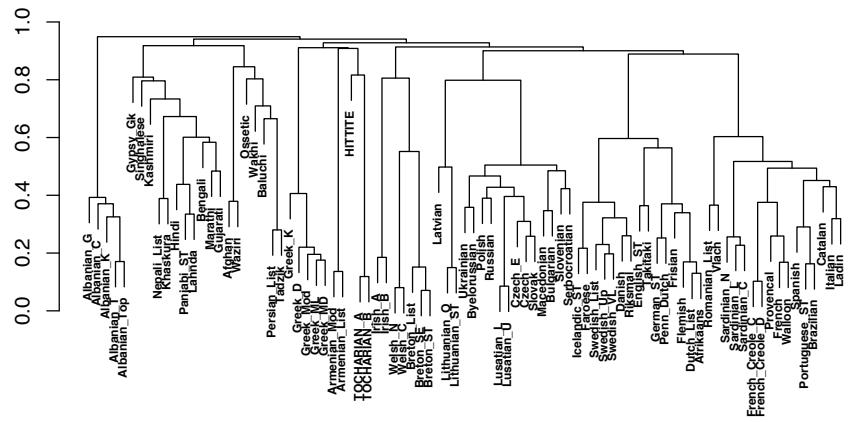
dat  
Agglomerative Coefficient = 0.99

## Hierarchical Clustering on Artificial Dataset



81

## Hierarchical Clustering on Indo-European Languages



83

## Using Dendograms

- ▶ Different ways of measuring dissimilarity result in different trees.
- ▶ Dendograms are useful for getting a feel for the structure of high-dimensional data though they don't represent distances between observations well.
- ▶ Dendograms show hierarchical clusters with respect to increasing values of dissimilarity between clusters, cutting a dendrogram horizontally at a particular height partitions the data into disjoint clusters which are represented by the vertical lines it intersects. Cutting horizontally effectively reveals the state of the clustering algorithm when the dissimilarity value between clusters is no more than the value cut at.
- ▶ Despite the simplicity of this idea and the above drawbacks, hierarchical clustering methods provide users with interpretable dendograms that allow clusters in high-dimensional data to be better understood.

82

## K-means

Partition-based methods seek to divide data points into a pre-assigned number of clusters  $C_1, \dots, C_K$  where for all  $k, k' \in \{1, \dots, K\}$ ,

$$C_k \subset \{1, \dots, n\}, \quad C_k \cap C_{k'} = \emptyset \quad \forall k \neq k', \quad \bigcup_{k=1}^K C_k = \{1, \dots, n\}.$$

For each cluster, represent it using a **prototype** or **cluster centre**  $\mu_k$ . We can measure the quality of a cluster with its **within-cluster deviance**

$$W(C_k, \mu_k) = \sum_{i \in C_k} \|x_i - \mu_k\|_2^2.$$

The overall quality of the clustering is given by the total within-cluster deviance:

$$W = \sum_{k=1}^K W(C_k, \mu_k).$$

The overall objective is to choose both the cluster centres and allocation of points to minimize the **objective function**.

84

## K-means

$$W = \sum_{k=1}^K \sum_{i \in C_k} \|x_i - \mu_k\|_2^2 = \sum_{i=1}^n \|x_i - \mu_{c_i}\|_2^2$$

where  $c_i = k$  if and only if  $i \in C_k$ .

- Given partition  $\{C_k\}$ , we can find the optimal prototypes easily by differentiating  $W$  with respect to  $\mu_k$ :

$$\frac{\partial W}{\partial \mu_k} = 2 \sum_{i \in C_k} (x_i - \mu_k) = 0 \quad \Rightarrow \mu_k = \frac{1}{|C_k|} \sum_{i \in C_k} x_i$$

- Given prototypes, we can easily find the optimal partition by assigning each data point to the closest cluster prototype:

$$c_i = \operatorname{argmin}_k \|x_i - \mu_k\|_2^2$$

But joint minimization over both is computationally difficult.

85

87

## K-means

Some notes about the K-means algorithm.

- The algorithm stops in a finite number of iterations.** Between steps 2 and 3,  $W$  either stays constant or it decreases, this implies that we never revisit the same partition. As there are only finitely many partitions, the number of iterations cannot exceed this.
- The K-means algorithm need not converge to global optimum.** K-means is a heuristic search algorithm so it can get stuck at suboptimal configurations. The result depends on the starting configuration. Typically perform a number of runs from different configurations, and pick best clustering.

## K-means

The K-means algorithm is a well-known method that **locally optimizes** the objective function  $W$ .

Iterative and alternating minimization.

- Randomly fix  $K$  cluster centres  $\mu_1, \dots, \mu_K$ .
- For each  $i = 1, \dots, n$ , assign each  $x_i$  to the cluster with the nearest centre,

$$c_i := \operatorname{argmin}_k \|x_i - \mu_k\|_2^2$$

- Set  $C_k := \{i : c_i = k\}$  for each  $k$ .
- Move cluster centres  $\mu_1, \dots, \mu_K$  to the average of the new clusters:

$$\mu_k := \frac{1}{|C_k|} \sum_{i \in C_k} x_i$$

- Repeat steps 2 to 4 until there is no more changes.
- Return the partition  $\{C_1, \dots, C_K\}$  and means  $\mu_1, \dots, \mu_K$  at the end.

86

## K-means on Crabs

Looking at the Crabs data again.

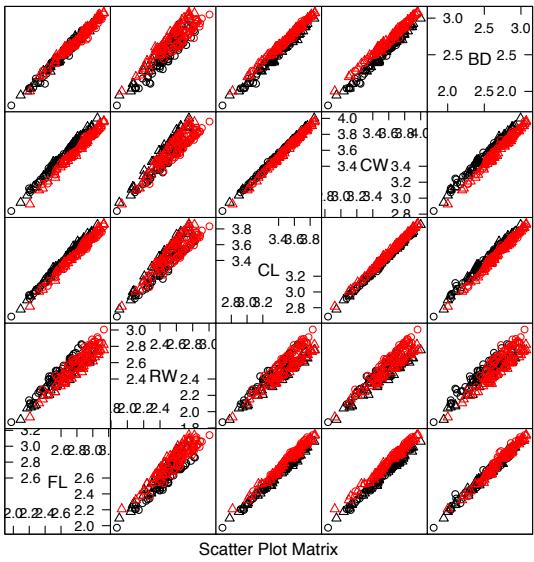
```
library(MASS)
library(lattice)
data(crabs)

splom(~log(crabs[, 4:8]),
      col=as.numeric(crabs[, 1]),
      pch=as.numeric(crabs[, 2]),
      main="circle/triangle is gender, black/red is species")
```

88

## K-means on Crabs

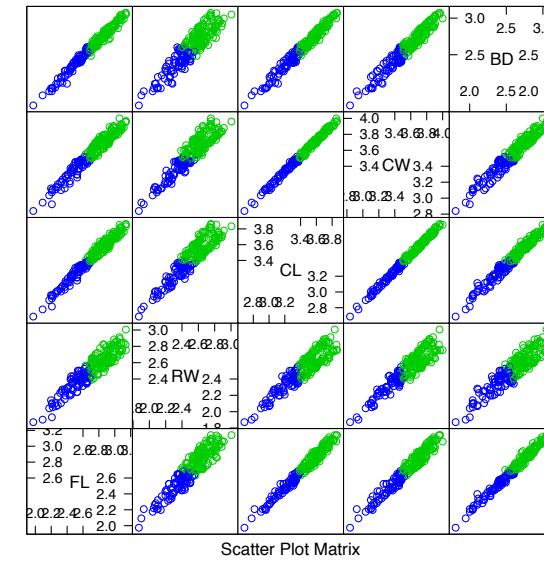
circle/triangle is gender, black/red is species



89

## K-means on Crabs

blue/green is cluster finds big/small



91

## K-means on Crabs

Apply K-means with 2 clusters and plot results.

```
cl <- kmeans( log(crabs[,4:8]), 2, nstart=1, iter.max=10)
splom(~log(crabs[,4:8]),
      col=cl$cluster+2,
      main="blue/green is cluster finds big/small")
```

90

## K-means on Crabs

'Whiten' or 'sphere'<sup>1</sup> the data using PCA.

```
pcp <- princomp( log(crabs[,4:8]) )
spc <- pcp$scores %*% diag(1/pcp$sdev)
splom( ~spc[,1:3],
      col=as.numeric(crabs[,1]),
      pch=as.numeric(crabs[,2]),
      main="circle/triangle is gender, black/red is species")
```

And apply K-means again.

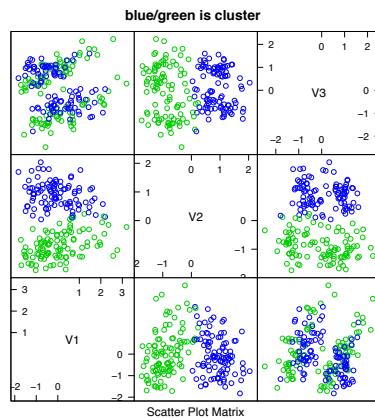
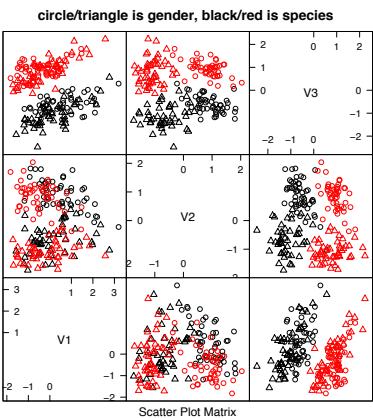
```
cl <- kmeans(spc, 2, nstart=1, iter.max=20)
splom( ~spc[,1:3],
      col=cl$cluster+2, main="blue/green is cluster")
```

---

<sup>1</sup>Apply a linear transformation so that covariance matrix is identity.

92

## K-means on Crabs



Discovers gender difference...

Results depends crucially on spherling the data first.

## K-means on Spike Waveforms

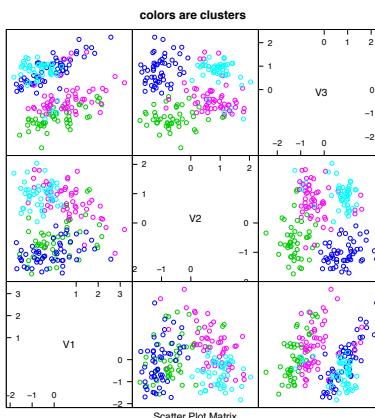
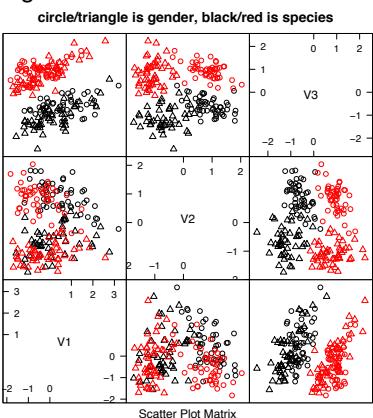
```
library(MASS)
library(lattice)
spikesPCA <- read.table("spikes.txt")
cl <- kmeans(data, 6, nstart=20)
splom(data, col=cl$cluster)
```

93

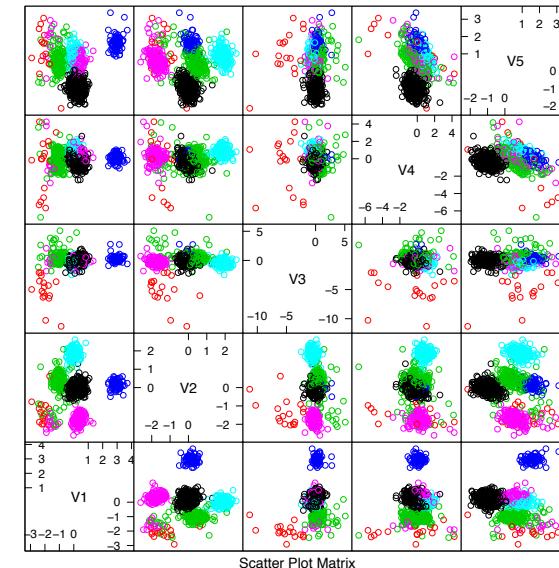
95

## K-means on Crabs

Using 4 cluster centers.



## K-means on Spike Waveforms



94

96

## Stochastic Optimization

- ▶ Each iteration of K-means requires a pass through whole dataset. In extremely large datasets, this can be computationally prohibitive.
- ▶ Stochastic optimization: update cluster means after assigning each data point to the closest cluster.
- ▶ Repeat for  $t = 1, 2, \dots$  until satisfactory convergence:
  1. Pick data item  $x_i$  either randomly or in order.
  2. Assign  $x_i$  to the cluster with the nearest centre,

$$c_i := \operatorname{argmin}_k \|x_i - \mu_k\|_2^2$$

3. Update cluster centre:

$$\mu_k := \mu_k + \alpha_t(x_i - \mu_k)$$

where  $\alpha_t > 0$  are **step sizes**.

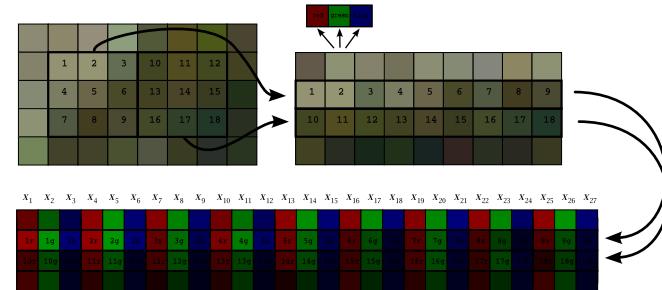
- ▶ Algorithm stochastically minimizes the objective function. Convergence requires slowly decreasing step sizes:

$$\sum_{t=1}^{\infty} \alpha_t = \infty$$

$$\sum_{t=1}^{\infty} \alpha_t^2 < \infty$$

## VQ Image Compression

$3 \times 3$  block VQ: View each block of  $3 \times 3$  pixels as single observation



## Vector Quantization

- ▶ A related algorithm developed in the signal processing literature for **lossy data compression**.
- ▶ If  $K \ll n$ , we can store the **codebook** of **codewords**  $\mu_1, \dots, \mu_K$ , and each vector  $x_i$  is encoded using  $c_i$ , which only requires  $\lceil \log K \rceil$  bits.
- ▶ As with K-means,  $K$  must be specified. Increasing  $K$  improves the quality of the compressed image but worsens the data compression rate, so there is a clear tradeoff.
- ▶ Some audio and video codecs use this method.
- ▶ Stochastic optimization algorithm for K-means was originally developed for VQ.

## VQ Image Compression

Original image (24 bits/pixel, uncompressed size 1,402 kB)



## VQ Image Compression

Codebook length 1024 (1.11 bits/pixel, total size 88kB)



101

## VQ Image Compression

Codebook length 16 (0.44 bits/pixel, total size 27kB)



103

## VQ Image Compression

Codebook length 128 (0.78 bits/pixel, total size 50kB)



102

## K-means Additional Comments

- ▶ **Sensitivity to distance measure.** Euclidean distance can be greatly affected by measurement unit and by strong correlations. Can use Mahalanobis distance,

$$\|x - y\|_M = \sqrt{(x - y)^\top M^{-1} (x - y)}$$

where  $M$  is positive semi-definite matrix, e.g. sample covariance.

- ▶ **Other partition based methods.** There are many other partition based methods that employ related ideas. For example K-medoids differs from K-means in requiring cluster centres  $\mu_i$  to be an observation  $x_i$ <sup>2</sup>, K-medians (use median in each dimension) and K-modes (use mode).
- ▶ **Determination of  $K$ .** The K-means objective will always improve with larger number of clusters  $K$ . Determination of  $K$  requires an additional **regularization** criterion. E.g., in DP-means<sup>3</sup>, use

$$W = \sum_{k=1}^K \sum_{i \in C_k} \|x_i - \mu_k\|_2^2 + \lambda K$$

<sup>2</sup>See also Affinity propagation.

<sup>3</sup>DP-means paper.

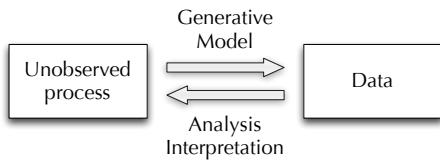
104

## Probabilistic Methods

- ▶ Algorithmic approach:



- ▶ Probabilistic modelling approach:



## Mixture Models

- ▶ Unknowns to learn given data are
  - ▶ **Parameters:**  $\pi_1, \dots, \pi_K, \phi_1, \dots, \phi_K$ , as well as
  - ▶ **Latent variables:**  $z_1, \dots, z_n$ .
- ▶ The joint probability over all cluster indicator variables  $\{Z_i\}$  are:

$$p_Z((z_i)_{i=1}^n) = \prod_{i=1}^n \pi_{z_i} = \prod_{i=1}^n \prod_{k=1}^K \pi_k^{1(z_i=k)}$$

- ▶ The joint density at observations  $X_i = x_i$  given  $Z_i = z_i$  are:
 
$$p_X((x_i)_{i=1}^n | (Z_i = z_i)_{i=1}^n) = \prod_{i=1}^n \prod_{k=1}^K f(x_i | \phi_k)^{1(z_i=k)}$$
- ▶ So the joint probability/density<sup>4</sup> is:
 
$$p_{X,Z}((x_i, z_i)_{i=1}^n) = \prod_{i=1}^n \prod_{k=1}^K (\pi_k f(x_i | \phi_k))^{1(z_i=k)}$$

<sup>4</sup>In this course we will treat probabilities and densities equivalently for notational simplicity. In general, the quantity is a density with respect to the product base measure, where the base measure is the counting measure for discrete variables and Lebesgue for continuous variables.

105

107

## Mixture Models

- ▶ Mixture models suppose that our dataset was created by sampling iid from  $K$  distinct populations (called **mixture components**).
- ▶ Typical samples in population  $k$  can be modelled using a distribution  $F(\phi_k)$  with density  $f(x|\phi_k)$ . For a concrete example, consider a Gaussian with unknown mean  $\phi_k$  and known symmetric covariance  $\sigma^2 I$ ,

$$f(x|\phi_k) = [2\pi\sigma^2]^{-\frac{p}{2}} \exp\left(-\frac{1}{2\sigma^2}\|x - \phi_k\|_2^2\right).$$

- ▶ Generative process: for  $i = 1, 2, \dots, n$ :

- ▶ First determine which population item  $i$  came from (independently):

$$Z_i \sim \text{Discrete}(\pi_1, \dots, \pi_K) \quad \text{i.e. } \mathbb{P}(Z_i = k) = \pi_k$$

where **mixing proportions** are  $\pi_k \geq 0$  for each  $k$  and  $\sum_{k=1}^K \pi_k = 1$ .

- ▶ If  $Z_i = k$ , then  $X_i = (X_{i1}, \dots, X_{ip})^\top$  is sampled (independently) from corresponding population distribution:

$$X_i | Z_i = k \sim F(\phi_k)$$

- ▶ We observe that  $X_i = x_i$  for each  $i$ , and would like to learn about the unknown parameters of the process.

## Mixture Models - Posterior Distribution

- ▶ Suppose we know the parameters  $(\pi_k, \phi_k)_{k=1}^K$ .
- ▶  $Z_i$  is a random variable, so the posterior distribution given data set  $\mathbf{X}$  tells us what we know about it:

$$Q_{ik} := p(Z_i = k | x_i) = \frac{p(Z_i = k, x_i)}{p(x_i)} = \frac{\pi_k f(x_i | \phi_k)}{\sum_{j=1}^K \pi_j f(x_i | \phi_j)}$$

where the marginal probability is:

$$p(x_i) = \sum_{j=1}^K \pi_j f(x_i | \phi_j)$$

- ▶ The posterior probability  $Q_{ik}$  of  $Z_i = k$  is called the **responsibility** of mixture component  $k$  for data point  $x_i$ .
- ▶ The posterior distribution **softly partitions** the dataset among the  $K$  components.

106

108

## Mixture Models - Maximum Likelihood

- ▶ How can we learn about the parameters  $\theta = (\pi_k, \phi_k)_{k=1}^K$  from data?
- ▶ Standard statistical methodology asks for the **maximum likelihood** estimator (MLE).
- ▶ The log likelihood is the log marginal probability of the data:

$$\ell((\pi_k, \phi_k)_{k=1}^K) := \log p((x_i)_{i=1}^n | (\pi_k, \phi_k)_{k=1}^K) = \sum_{i=1}^n \log \sum_{j=1}^K \pi_j f(x_i | \phi_j)$$

$$\begin{aligned}\nabla_{\phi_k} \ell((\pi_k, \phi_k)_{k=1}^K) &= \sum_{i=1}^n \frac{\pi_k f(x_i | \phi_k)}{\sum_{j=1}^K \pi_j f(x_i | \phi_j)} \nabla_{\phi_k} \log f(x_i | \phi_k) \\ &= \sum_{i=1}^n Q_{ik} \nabla_{\phi_k} \log f(x_i | \phi_k)\end{aligned}$$

- ▶ A difficult equation to solve, as  $Q_{ik}$  depends implicitly on  $\phi_k$ ...

109

## Mixture Models - Maximum Likelihood

$$\sum_{i=1}^n Q_{ik} \nabla_{\phi_k} \log f(x_i | \phi_k) = 0$$

- ▶ What if we ignore the dependence of  $Q_{ik}$  on the parameters?
- ▶ Taking the mixture of Gaussian with covariance  $\sigma^2 I$  as example,

$$\begin{aligned}& \sum_{i=1}^n Q_{ik} \nabla_{\phi_k} \left( -\frac{p}{2} \log(2\pi\sigma^2) - \frac{1}{2\sigma^2} \|x_i - \phi_k\|_2^2 \right) \\&= \frac{1}{\sigma^2} \sum_{i=1}^n Q_{ik} (x_i - \phi_k) = \frac{1}{\sigma^2} \left( (\sum_{i=1}^n Q_{ik} x_i) - \phi_k (\sum_{i=1}^n Q_{ik}) \right) = 0 \\&\phi_k^{MLE?} = \frac{\sum_{i=1}^n Q_{ik} x_i}{\sum_{i=1}^n Q_{ik}}\end{aligned}$$

110

## Mixture Models - Maximum Likelihood

- ▶ The estimate is a weighted average of data points, where the estimated mean of cluster  $k$  uses its responsibilities to data points as weights.

$$\phi_k^{MLE?} = \frac{\sum_{i=1}^n Q_{ik} x_i}{\sum_{i=1}^n Q_{ik}}$$

- ▶ Makes sense: Suppose we knew that data point  $x_i$  came from population  $z_i$ . Then  $Q_{iz_i} = 1$  and  $Q_{ik} = 0$  for  $k \neq z_i$  and:

$$\pi_k^{MLE} = \frac{\sum_{i:z_i=k} x_i}{\sum_{i:z_i=k} 1}$$

- ▶ Our best guess of the originating population is given by  $Q_{ik}$ .

111

## Mixture Models - Maximum Likelihood

- ▶ For the mixing proportions, we can similarly derive an estimator.
- ▶ Include a Lagrange multiplier  $\lambda$  to enforce constraint  $\sum_k \pi_k = 1$ .

$$\begin{aligned}& \nabla_{\log \pi_k} \left( \ell((\pi_k, \phi_k)_{k=1}^K) - \lambda(\sum_{k=1}^K \pi_k - 1) \right) \\&= \sum_{i=1}^n \frac{\pi_k f(x_i | \phi_k)}{\sum_{j=1}^K \pi_j f(x_i | \phi_j)} - \lambda \pi_k \\&= \sum_{i=1}^n Q_{ik} - \lambda \pi_k = 0 \\&\pi_k^{MLE?} = \frac{\sum_{i=1}^n Q_{ik}}{n}\end{aligned}$$

- ▶ Again makes sense: the estimate is simply (our best guess of) the proportion of data points coming from population  $k$ .

112

## Mixture Models - The EM Algorithm

- ▶ Putting all the derivations together, we get an iterative algorithm for learning about the unknowns in the mixture model.
- ▶ Start with some initial parameters  $(\pi_k^{(0)}, \phi_l^{(0)})_{k=1}^K$ .
- ▶ Iterate for  $t = 1, 2, \dots$ :
  - ▶ **Expectation Step:**

$$Q_{ik}^{(t)} := \frac{\pi_k^{(t-1)} f(x_i | \phi_k^{(t-1)})}{\sum_{j=1}^K \pi_j^{(t-1)} f(x_i | \phi_j^{(t-1)})}$$

- ▶ **Maximization Step:**

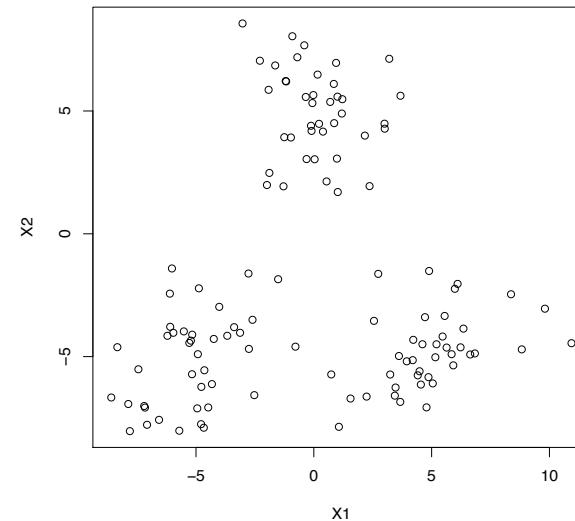
$$\pi_k^{(t)} = \frac{\sum_{i=1}^n Q_{ik}^{(t)}}{n} \quad \phi_k^{(t)} = \frac{\sum_{i=1}^n Q_{ik}^{(t)} x_i}{\sum_{i=1}^n Q_{ik}^{(t)}}$$

- ▶ Will the algorithm converge?
- ▶ What does it converge to?

113

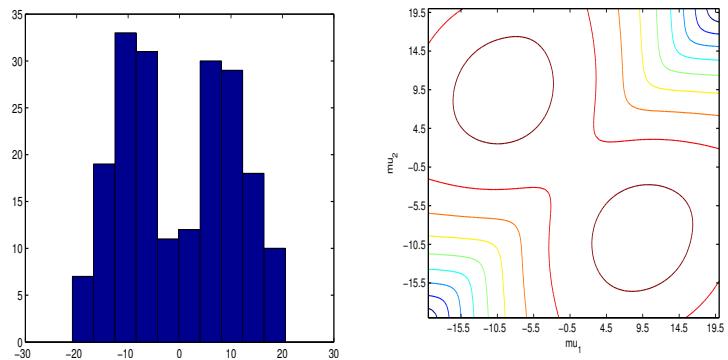
## Example: Mixture of 3 Gaussians

An example with 3 clusters.



115

## Likelihood Surface for a Simple Example

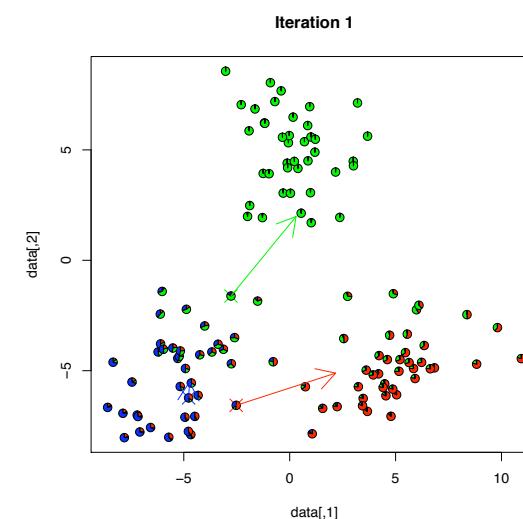


- (left)  $n = 200$  data points from a mixture of two 1D Gaussians with  $\pi_1 = \pi_2 = 0.5$ ,  $\sigma = 5$  and  $\mu_1 = 10, \mu_2 = -10$ .  
 (right) Log likelihood surface  $\ell(\mu_1, \mu_2)$ , all the other parameters being assumed known.

114

## Example: Mixture of 3 Gaussians

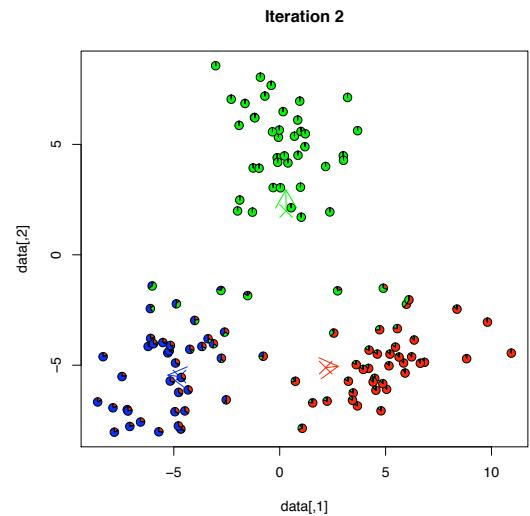
After 1st E and M step.



116

## Example: Mixture of 3 Gaussians

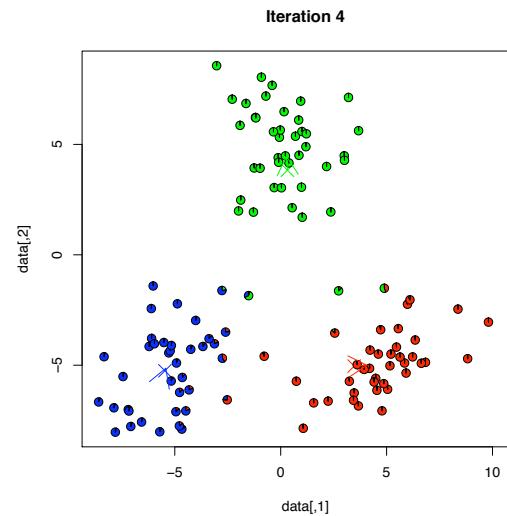
After 2nd E and M step.



117

## Example: Mixture of 3 Gaussians

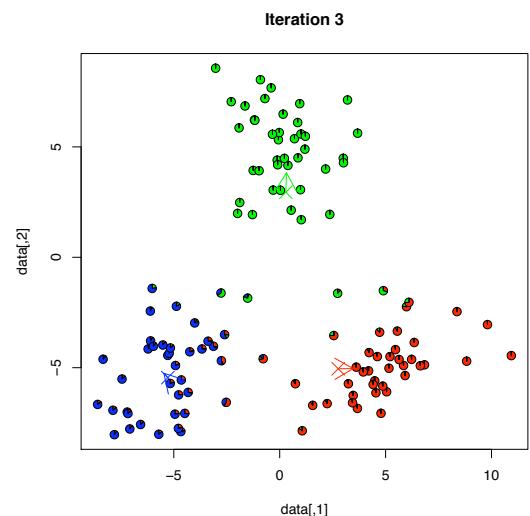
After 4th E and M step.



119

## Example: Mixture of 3 Gaussians

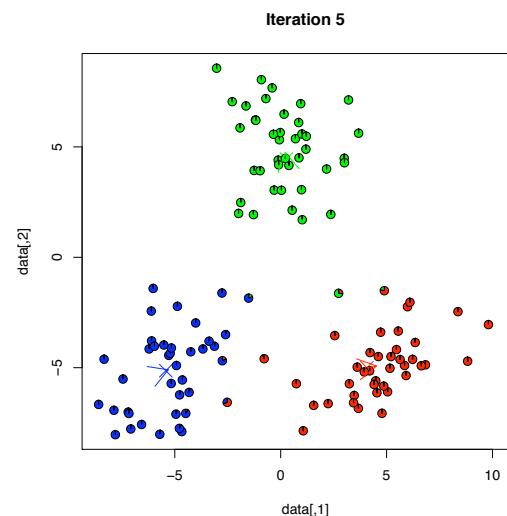
After 3rd E and M step.



118

## Example: Mixture of 3 Gaussians

After 5th E and M step.



120

## The EM Algorithm

- In a maximum likelihood framework, the objective function is the log likelihood,

$$\ell(\theta) = \sum_{i=1}^n \log \sum_{j=1}^K \pi_j f(x_i | \phi_j)$$

Direct maximization is not feasible.

- Consider another objective function  $\mathcal{F}(\theta, q)$  such that:

$$\begin{aligned}\mathcal{F}(\theta, q) &\leq \ell(\theta) \text{ for all } \theta, q, \\ \max_q \mathcal{F}(\theta, q) &= \ell(\theta)\end{aligned}$$

$\mathcal{F}(\theta, q)$  is a lower bound on the log likelihood.

- We can construct an alternating maximization algorithm as follows:

For  $t = 1, 2, \dots$  until convergence:

$$\begin{aligned}q^{(t)} &:= \operatorname{argmax}_q \mathcal{F}(\theta^{(t-1)}, q) \\ \theta^{(t)} &:= \operatorname{argmax}_\theta \mathcal{F}(\theta, q^{(t)})\end{aligned}$$

121

## EM Algorithm

- The lower bound we use is called the **variational free energy**.
- $q$  is a probability mass function for some distribution over  $(Z_i)$  and

$$\begin{aligned}\mathcal{F}(\theta, q) &= \mathbb{E}_q [\log p((x_i, z_i)_{i=1}^n) - \log q((z_i)_{i=1}^n)] \\ &= \mathbb{E}_q \left[ \left( \sum_{i=1}^n \sum_{k=1}^K \mathbb{1}(z_i = k) (\log \pi_k + \log f(x_i | \phi_k)) \right) - \log q(\mathbf{z}) \right] \\ &= \sum_{\mathbf{z}} q(\mathbf{z}) \left[ \left( \sum_{i=1}^n \sum_{k=1}^K \mathbb{1}(z_i = k) (\log \pi_k + \log f(x_i | \phi_k)) \right) - \log q(\mathbf{z}) \right]\end{aligned}$$

Using  $\mathbf{z} := (z_i)_{i=1}^n$  to shorten notation.

## EM Algorithm - Solving for $q$

- Introducing Lagrange multiplier to enforce  $\sum_{\mathbf{z}} q(\mathbf{z}) = 1$ , and setting derivatives to 0,

$$\begin{aligned}\nabla_{q(\mathbf{z})} \mathcal{F}(\theta, q) &= \sum_{i=1}^n \sum_{k=1}^K \mathbb{1}(z_i = k) (\log \pi_k + \log f(x_i | \phi_k)) - \log q(\mathbf{z}) - 1 - \lambda \\ &= \sum_{i=1}^n (\log \pi_{z_i} + \log f(x_i | \phi_{z_i})) - \log q(\mathbf{z}) - 1 - \lambda = 0 \\ q^*(\mathbf{z}) &= \frac{\prod_{i=1}^n \pi_{z_i} f(x_i | \phi_{z_i})}{\sum_{\mathbf{z}'} \prod_{i=1}^n \pi_{z'_i} f(x_i | \phi_{z'_i})} = \prod_{i=1}^n \frac{\pi_{z_i} f(x_i | \phi_{z_i})}{\sum_k \pi_k f(x_i | \phi_k)} = \prod_{i=1}^n p(z_i | x_i, \theta)\end{aligned}$$

- Optimal  $q^*$  is simply the posterior distribution.
- Plugging in optimal  $q^*$  into the variational free energy,

$$\mathcal{F}(\theta, q^*) = \sum_{i=1}^n \log \sum_{k=1}^K \pi_k f(x_i | \phi_k) = \ell(\theta)$$

123

## EM Algorithm - Solving for $\theta$

- Setting derivative with respect to  $\phi_k$  to 0,

$$\begin{aligned}\nabla_{\phi_k} \mathcal{F}(\theta, q) &= \sum_{\mathbf{z}} q(\mathbf{z}) \sum_{i=1}^n \mathbb{1}(z_i = k) \nabla_{\phi_k} \log f(x_i | \phi_k) \\ &= \sum_{i=1}^n q(z_i = k) \nabla_{\phi_k} \log f(x_i | \phi_k) = 0\end{aligned}$$

- This equation can be solved quite easily. E.g., for mixture of Gaussians,

$$\phi_k^* = \frac{\sum_{i=1}^n q(z_i = k) x_i}{\sum_{i=1}^n q(z_i = k)}$$

- If it cannot be solved exactly, we can use **gradient ascent** algorithm:

$$\phi_k^* = \phi_k + \alpha \sum_{i=1}^n q(z_i = k) \nabla_{\phi_k} \log f(x_i | \phi_k)$$

- This leads to **generalized EM algorithm**. Further extension using **stochastic optimization** method leads to **stochastic EM algorithm**.
- Similar derivation for optimal  $\pi_k$  as before.

122

124

## EM Algorithm

- ▶ Start with some initial parameters  $(\pi_k^{(0)}, \phi_l^{(0)})_{k=1}^K$ .
- ▶ Iterate for  $t = 1, 2, \dots$ :
- ▶ **Expectation Step:**

$$q^{(t)}(z_i = k) := \frac{\pi_k^{(t-1)} f(x_i | \phi_k^{(t-1)})}{\sum_{j=1}^K \pi_j^{(t-1)} f(x_i | \phi_j^{(t-1)})} = \mathbb{E}_{p(z_i | x_i, \theta^{(t-1)})} [\mathbb{1}(z_i = k)]$$

- ▶ **Maximization Step:**

$$\pi_k^{(t)} = \frac{\sum_{i=1}^n q^{(t)}(z_i = k)}{n} \quad \phi_k^{(t)} = \frac{\sum_{i=1}^n q^{(t)}(z_i = k) x_i}{\sum_{i=1}^n q^{(t)}(z_i = k)}$$

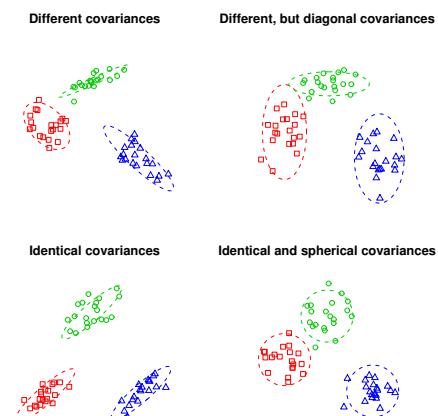
- ▶ Each step increases the log likelihood:

$$\ell(\theta^{(t-1)}) = \mathcal{F}(\theta^{(t-1)}, q^{(t)}) \leq \mathcal{F}(\theta^{(t)}, q^{(t)}) \leq \mathcal{F}(\theta^{(t)}, q^{(t+1)}) = \ell(\theta^{(t)}).$$

- ▶ Additional assumption, that  $\nabla_\theta^2 \mathcal{F}(\theta^{(t)}, q^{(t)})$  are negative definite with eigenvalues  $< -\epsilon < 0$ , implies that  $\theta^{(t)} \rightarrow \theta^*$  where  $\theta^*$  is a local MLE.

## Flexible Gaussian Mixture Models

- ▶ We can allow each cluster to have its own mean and covariance structure allows greater flexibility in the model.



125

127

## Notes on Probabilistic Approach and EM Algorithm

Some good things:

- ▶ Guaranteed convergence to locally optimal parameters.
- ▶ Formal reasoning of uncertainties, using both Bayes Theorem and maximum likelihood theory.
- ▶ Rich language of probability theory to express a wide range of generative models, and straightforward derivation of algorithms for ML estimation.

Some bad things:

- ▶ Can get stuck in local minima so multiple starts are recommended.
- ▶ Slower and more expensive than K-means.
- ▶ Choice of  $K$  still problematic, but rich array of methods for model selection comes to rescue.

## Probabilistic PCA

- ▶ A probabilistic model related to PCA has the following generative model: for  $i = 1, 2, \dots, n$ :
- ▶ Let  $k < n, p$  be given.
- ▶ Let  $Y_i$  be a  $k$ -dimensional normally distributed random variable with 0 mean and identity covariance:

$$Y_i \sim \mathcal{N}(0, I_k)$$

- ▶ We model the distribution of the  $i$ th data point given  $Y_i$  as a  $p$ -dimensional normal:

$$X_i \sim \mathcal{N}(\mu + LY_i, \sigma^2 I)$$

- ▶ where the parameters are a vector  $\mu \in \mathbb{R}^p$ , a matrix  $L \in \mathbb{R}^{p \times k}$  and  $\sigma^2 > 0$ .
- ▶ EM algorithm can be used for ML estimation, but PCA can more directly give a MLE (note this is not unique).
- ▶ Let  $\lambda_1 \geq \dots \geq \lambda_p$  be the eigenvalues of the sample covariance and let  $V \in \mathbb{R}^{p \times k}$  have columns given by the eigenvectors of the top  $k$  eigenvalues. Let  $R \in \mathbb{R}^{k \times k}$  be orthogonal. Then a MLE is:

$$\begin{aligned} \mu^{\text{MLE}} &= \bar{x} & (\sigma^2)^{\text{MLE}} &= \frac{1}{p-k} \sum_{j=k+1}^p \lambda_j \\ L^{\text{MLE}} &= V \text{diag}((\lambda_1 - (\sigma^2)^{\text{MLE}})^{\frac{1}{2}}, \dots, (\lambda_k - (\sigma^2)^{\text{MLE}})^{\frac{1}{2}}) R \end{aligned}$$

126

128

## Mixture of Probabilistic PCAs

- ▶ We have learnt two types of unsupervised learning techniques:
  - ▶ Dimensionality reduction, e.g. PCA, MDS, Isomap.
  - ▶ Clustering, e.g. K-means, linkage and mixture models.
- ▶ Probabilistic models allow us to construct more complex models from simpler pieces.
- ▶ Mixture of probabilistic PCAs allows both clustering and dimensionality reduction at the same time.

$$Z_i \sim \text{Discrete}(\pi_1, \dots, \pi_K)$$

$$Y_i \sim \mathcal{N}(0, I_d)$$

$$X_i | Z_i = k, Y_i = y_i \sim \mathcal{N}(\mu_k + Ly_i, \sigma^2 I_p)$$

- ▶ Allows flexible modelling of covariance structure without using too many parameters.

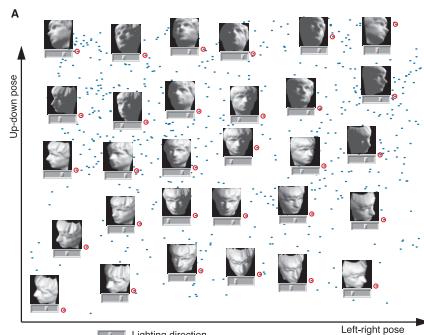
Ghahramani and Hinton 1996

129

131

## Mixture of Probabilistic PCAs

- ▶ PCA can reconstruct  $x$  given low dimensional embedding  $z$ , but is linear.
- ▶ Isomap is non-linear, but cannot reconstruct  $x$  given any  $z$ .



- ▶ We can learn a probabilistic mapping between the  $k$ -dimensional Isomap embedding space and the  $p$ -dimensional data space.
- ▶ Demo: [Using LLE instead of Isomap, and Mixture of factor analysers instead of Mixture of PPCAs.]

Teh and Roweis 2002

130

## Further Readings—Unsupervised Learning

- ▶ Hastie et al, Chapter 14.
- ▶ James et al, Chapter 10.
- ▶ Venables and Ripley, Chapter 11.
- ▶ Tukey, John W. (1980). We need both exploratory and confirmatory. *The American Statistician* 34 (1): 23-25.

## Supervised Learning

### Unsupervised learning:

- ▶ To “extract structure” and postulate hypotheses about data generating process from observations  $x_1, \dots, x_n$ .
- ▶ Visualize, summarize and compress data.

We have seen how response or grouping variables are used to validate the usefulness of the extracted structure.

### Supervised learning:

- ▶ In addition to the  $n$  observations of  $X$ , we also have a response variable  $Y \in \mathcal{Y}$ .
- ▶ Techniques for predicting  $Y$  given  $X$ .
  - ▶ Classification: discrete responses, e.g.  $\mathcal{Y} = \{+1, -1\}$  or  $\{1, \dots, K\}$ .
  - ▶ Regression: a numerical value is observed and  $\mathcal{Y} = \mathbb{R}$ .

Given training data  $(x_i, y_i), i = 1, \dots, n$ , the goal is to accurately predict the class or response  $Y$  on new observations of  $X$ .

132

## Regression Example: Boston Housing

The original data are 506 observations on 13 variables  $\mathbf{X}$ ; medv being the response variable  $\mathbf{Y}$ .

```

crim    per capita crime rate by town
zn      proportion of residential land zoned for lots
        over 25,000 sq.ft
indus   proportion of non-retail business acres per town
chas    Charles River dummy variable (= 1 if tract bounds river;
        0 otherwise)
nox     nitric oxides concentration (parts per 10 million)
rm      average number of rooms per dwelling
age     proportion of owner-occupied units built prior to 1940
dis     weighted distances to five Boston employment centers
rad     index of accessibility to radial highways
tax     full-value property-tax rate per USD 10,000
ptratio pupil-teacher ratio by town
b       1000(B - 0.63)^2 where B is the proportion of blacks by town
lstat   percentage of lower status of the population
medv   median value of owner-occupied homes in USD 1000's

```

133

## Regression Example: Boston Housing

```

> str(X)
'data.frame': 506 obs. of 13 variables:
 $ crim   : num  0.00632 0.02731 0.02729 0.03237 0.06905 ...
 $ zn     : num  18 0 0 0 0 12.5 12.5 12.5 12.5 ...
 $ indus  : num  2.31 7.07 7.07 2.18 2.18 2.18 7.87 7.87 7.87 ...
 $ chas   : int  0 0 0 0 0 0 0 0 0 ...
 $ nox   : num  0.538 0.469 0.469 0.458 0.458 0.458 0.524 0.524 ...
 $ rm    : num  6.58 6.42 7.18 7.00 7.15 ...
 $ age   : num  65.2 78.9 61.1 45.8 54.2 58.7 66.6 96.1 100 85.9
 $ dis   : num  4.09 4.97 4.97 6.06 6.06 ...
 $ rad   : int  1 2 2 3 3 3 5 5 5 5 ...
 $ tax   : num  296 242 242 222 222 311 311 311 311 ...
 $ ptratio: num  15.3 17.8 17.8 18.7 18.7 18.7 15.2 15.2 15.2 15.2 ...
 $ black  : num  397 397 393 395 397 ...
 $ lstat  : num  4.98 9.14 4.03 2.94 5.33 ...

```

```

> str(Y)
num[1:506] 24 21.6 34.7 33.4 36.2 28.7 22.9 27.1 16.5 18.9 ...

```

Goal: predict median house price  $\hat{Y}(\mathbf{X})$ , given 13 predictor variables  $\mathbf{X}$  of a new district.

## Classification Example: Lymphoma

We have gene expression measurements  $\mathbf{X}$  of  $n = 62$  patients for  $p = 4026$  genes. For each patient,  $\mathbf{Y}$  denotes one of two subtypes of cancer. Goal: predict cancer subtype  $\hat{Y}(\mathbf{X}) \in \{0, 1\}$ , given gene expressions of a new patient.

```

> str(X)
'data.frame': 62 obs. of 4026 variables:
 $ Gene 1   : num -0.344 -1.188 0.520 -0.748 -0.868 ...
 $ Gene 2   : num -0.953 -1.286 0.657 -1.328 -1.330 ...
 $ Gene 3   : num -0.776 -0.588 0.409 -0.991 -1.517 ...
 $ Gene 4   : num -0.474 -1.588 0.219 0.978 -1.604 ...
 $ Gene 5   : num -1.896 -1.960 -1.695 -0.348 -0.595 ...
 $ Gene 6   : num -2.075 -2.117 0.121 -0.800 0.651 ...
 $ Gene 7   : num -1.8755 -1.8187 0.3175 0.3873 0.0414 ...
 $ Gene 8   : num -1.539 -2.433 -0.337 -0.522 -0.668 ...
 $ Gene 9   : num -0.604 -0.710 -1.269 -0.832 0.458 ...
 $ Gene 10  : num -0.218 -0.487 -1.203 -0.919 -0.848 ...
 $ Gene 11  : num -0.340 1.164 1.023 1.133 -0.541 ...
 $ Gene 12  : num -0.531 0.488 -0.335 0.496 -0.358 ...

```

```

> str(Y)
num [1:62] 0 0 0 1 0 0 1 0 0 0 ...

```

135

## Decision Theory

- ▶ Suppose we made a prediction  $\hat{Y} \in \mathcal{Y}$  based on observation of  $\mathbf{X}$ .
- ▶ How good is the prediction? We can use a **loss function**  $L : \mathcal{Y} \times \mathcal{Y} \mapsto \mathbb{R}^+$  to formalize the quality of the prediction.
- ▶ Typical loss functions:

- ▶ **Misclassification loss** (or **0-1 loss**) for classification

$$L(Y, \hat{Y}) = \begin{cases} 0 & Y = \hat{Y} \\ 1 & Y \neq \hat{Y} \end{cases}.$$

- ▶ **Squared loss** for regression

$$L(Y, \hat{Y}) = (Y - \hat{Y})^2.$$

- ▶ Alternative loss functions are often useful (later). For example, **weighted misclassification error** often appropriate. Or **log-likelihood loss** (sometimes shortened as **log loss**)  $L(Y, \hat{p}) = -\log \hat{p}(Y)$ , where  $\hat{p}(k)$  is the estimated probability of class  $k \in \mathcal{Y}$ .

134

136

## Decision Theory

- For a given loss function  $L$ , the **risk**  $R$  of a learner is given by the expected loss

$$R(\hat{Y}) = \mathbb{E}(L(Y, \hat{Y}(X))),$$

where the expectation is with respect to the true (unknown) joint distribution  $(X, Y)$ .

- The risk is unknown, but we can estimate it by the **empirical risk**:

$$R(\hat{Y}) \approx R_n(\hat{Y}) = \frac{1}{n} \sum_{i=1}^n L(y_i, \hat{Y}(x_i)).$$

137

## The Bayes Classifier

- Consider the situation of the 0-1 loss.
- The risk simplifies to:

$$\begin{aligned}\mathbb{E}[L(Y, \hat{Y}(x))|X=x] &= \sum_{k=1}^K L(k, \hat{Y}(x))\mathbb{P}(Y=k|X=x) \\ &= 1 - \mathbb{P}(Y=\hat{Y}(x)|X=x)\end{aligned}$$

- The risk is minimized by choosing the class with the greatest posterior probability:

$$\begin{aligned}\hat{Y}(x) &= \arg \max_{k=1,\dots,K} \mathbb{P}(Y=k|X=x) = \arg \max_{k=1,\dots,K} \frac{\pi_k f_k(x)}{\sum_{k=1}^K \pi_k f_k(x)} \\ &= \arg \max_{k=1,\dots,K} \pi_k f_k(x).\end{aligned}$$

- The functions  $x \mapsto \pi_k f_k(x)$  are called **discriminant functions**. The function with maximum value determines the predicted class of  $x$ .

139

## The Bayes Classifier

- What is the optimal classifier if the joint distribution  $(X, Y)$  were known?
- The joint distribution  $f$  of  $X$  can be written as a mixture

$$f(X) = \sum_{k=1}^K f_k(X)\mathbb{P}(Y=k),$$

where, for  $k = 1, \dots, K$ ,

- the prior probabilities over classes are  $\mathbb{P}(Y=k) = \pi_k$
- and distributions of  $X$ , conditional on  $Y=k$ , is  $f_k(X)$ .

- The **Bayes classifier**  $\hat{Y}(X) \mapsto \{1, \dots, K\}$  is the one with minimum risk:

$$\begin{aligned}R(\hat{Y}) &= \mathbb{E}[L(Y, \hat{Y}(X))] = \mathbb{E}[\mathbb{E}[L(Y, \hat{Y}(x))|X=x]] \\ &= \int_X \mathbb{E}[L(Y, \hat{Y}(x))|X=x]f(x)dx\end{aligned}$$

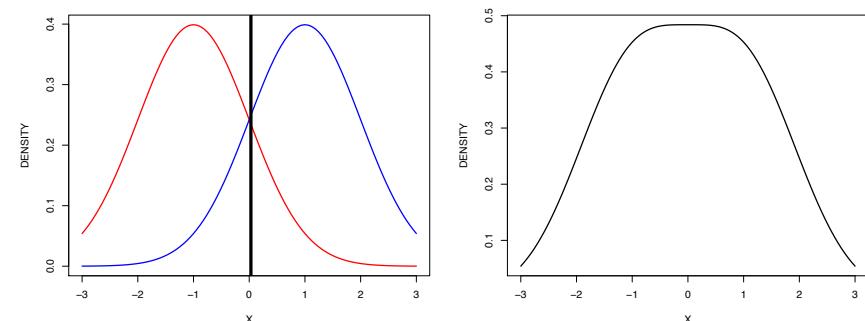
- The minimum risk attained by the Bayes classifier is called **Bayes risk**.
- Minimizing  $\mathbb{E}[L(Y, \hat{Y}(x))|X=x]$  separately for each  $x$  suffices.

138

## The Bayes Classifier

A simple two Gaussians example: Suppose  $X \sim \mathcal{N}(\mu_Y, 1)$ , where  $\mu_1 = -1$  and  $\mu_2 = 1$  and assume equal priors  $\pi_1 = \pi_2 = 1/2$ .

$$f_1(x) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{(x-(-1))^2}{2}\right) \quad \text{and} \quad f_2(x) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{(x-1)^2}{2}\right).$$

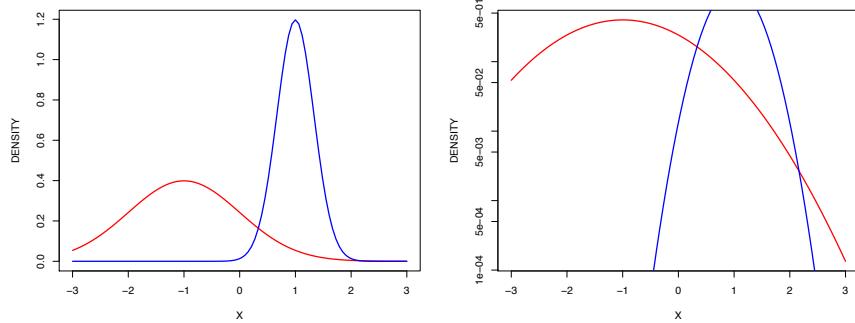


Optimal classification is  $\hat{Y}(x) = \arg \max_{k=1,\dots,K} \pi_k f_k(x) = \begin{cases} 1 & \text{if } x < 0, \\ 2 & \text{if } x \geq 0. \end{cases}$

140

## The Bayes Classifier

How do you classify a new observation  $x$  if now the standard deviation is still 1 for class 1 but 1/3 for class 2?



Looking at density in a log-scale, optimal classification is class 2 if and only if  $x \in [-0.39, 2.15]$ .

## Plug-in Classification

- The Bayes Classifier chooses the class with the greatest posterior probability

$$\hat{Y}(x) = \arg \max_{k=1,\dots,K} \pi_k f_k(x).$$

- Unfortunately, we usually know neither the conditional class probabilities nor the prior probabilities.
- We can estimate the joint distribution with:

- estimates  $\hat{\pi}_k$  for  $\pi_k$  and  $k = 1, \dots, K$  and
- estimates  $\hat{f}_k(x)$  of conditional class densities,

- The **plug-in classifiers** chooses the class

$$\hat{Y}(x) = \arg \max_{k=1,\dots,K} \hat{\pi}_k \hat{f}_k(x).$$

- **Linear Discriminant Analysis** will be an example of plug-in classification.

141

## Linear Discriminant Analysis

- LDA is the most well-known and simplest example of plug-in classification.
- Assume a multivariate Normal form for  $f_k(x)$  for each class  $k$ :

$$X|Y=k \sim \mathcal{N}(\mu_k, \Sigma),$$

- each class can have a **different mean**  $\mu_k$
- but all classes share the **same covariance**  $\Sigma$ .

- For an observation  $x$ ,

$$\begin{aligned} \log \mathbb{P}(Y=k|X=x) &= \kappa + \log \pi_k f_k(x) \\ &= \kappa + \log \pi_k - \frac{1}{2} (x - \mu_k)^\top \Sigma^{-1} (x - \mu_k) \end{aligned}$$

The quantity  $(x - \mu_k)^\top \Sigma^{-1} (x - \mu_k)$  is the square of the **Mahalanobis distance**. It gives the distance between  $x$  and  $\mu_k$  in the metric given by  $\Sigma$ .  
 ► If  $\Sigma = I_p$  and  $\pi_k = \frac{1}{K}$ ,  $\hat{Y}(x)$  simply chooses the class  $k$  with the nearest (in the Euclidean sense) mean.

143

## Linear Discriminant Analysis

- Expanding the **discriminant**  $(x - \mu_k)^\top \Sigma^{-1} (x - \mu_k)$ ,

$$\begin{aligned} \log \mathbb{P}(Y=k|x) &= \kappa + \log(\pi_k) - \frac{1}{2} (\mu_k^\top \Sigma^{-1} \mu_k - 2\mu_k^\top \Sigma^{-1} x + x^\top \Sigma^{-1} x) \\ &= \kappa + \log(\pi_k) - \frac{1}{2} \mu_k^\top \Sigma^{-1} \mu_k + \mu_k^\top \Sigma^{-1} x \end{aligned}$$

- Setting  $a_k = \log(\pi_k) - \frac{1}{2} \mu_k^\top \Sigma^{-1} \mu_k$  and  $b_k = \Sigma^{-1} \mu_k$ , we obtain

$$\log \mathbb{P}(Y=k|X=x) = \kappa + a_k + b_k^\top x$$

i.e. a **linear** discriminant function.

- Consider choosing class  $k$  over  $k'$ :

$$a_k + b_k^\top x > a_{k'} + b_{k'}^\top x \quad \Leftrightarrow \quad a_* + b_*^\top x > 0$$

where  $a_* = a_k - a_{k'}$  and  $b_* = b_k - b_{k'}$ .

- The Bayes classifier partitions  $\mathcal{X}$  into regions with the same class predictions via **separating hyperplanes**.
- The Bayes classifier under these assumptions is more commonly known as the **LDA classifier**.

142

144

## Parameter Estimation

- The final piece of the puzzle is to estimate the parameters of the LDA model.
- We can achieve this by maximum likelihood.
- EM algorithm is not needed here since the class variables  $y_j$  are observed.
- Let  $n_k = \#\{j : y_j = k\}$  be the number of observations in class  $k$ .

$$\ell(\pi, (\mu_k), \Sigma) = \kappa + \sum_{k=1}^K \sum_{j:y_j=k} \log \pi_k - \frac{1}{2} \left( \log |\Sigma| + (x_j - \mu_k)^\top \Sigma^{-1} (x_j - \mu_k) \right)$$

Then:

$$\hat{\pi}_k = \frac{n_k}{n} \quad \hat{\mu}_k = \frac{1}{n_k} \sum_{j:y_j=k} x_j$$

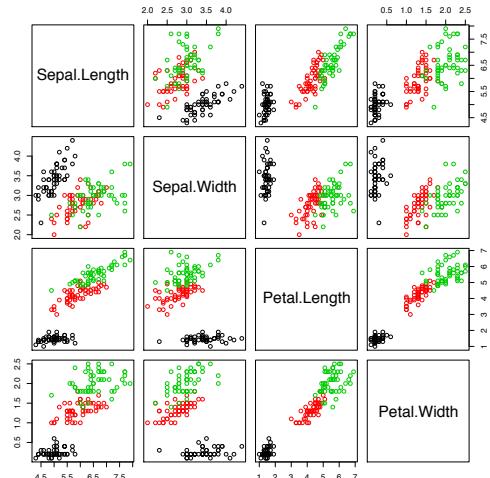
$$\hat{\Sigma} = \frac{1}{n} \sum_{k=1}^K \sum_{j:y_j=k} (x_j - \hat{\mu}_k)(x_j - \hat{\mu}_k)^\top$$

- Note: the ML estimate of  $\Sigma$  is not unbiased. For an unbiased estimate we need to divide by  $n - K$ .

145

## Iris Dataset

```
library(MASS)
data(iris)
##save class labels
ct <- rep(1:3, each=50)
##pairwise plot
pairs(iris[,1:4], col=ct)
```

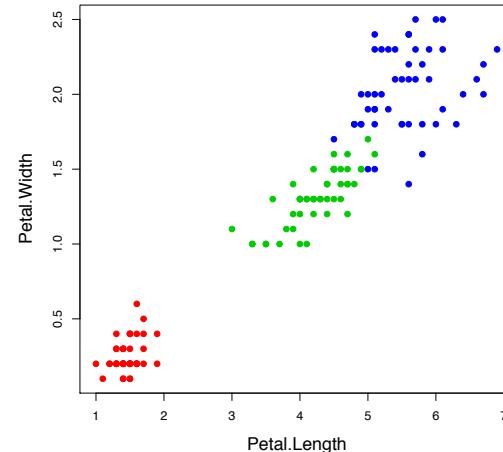


146

## Iris Dataset

Just focus on two predictor variables.

```
iris.data <- iris[, 3:4]
plot(iris.data, col=ct+1, pch=20, cex=1.5, cex.lab=1.4)
```



147

## Iris Dataset

Computing and plotting the LDA boundaries.

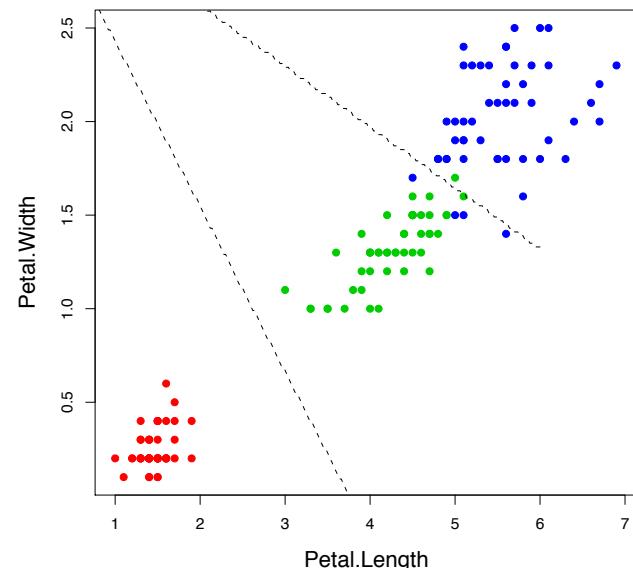
```
##fit LDA
iris.lda <- lda(x=iris.data, grouping=ct)

##create a grid for our plotting surface
x <- seq(-6, 6, 0.02)
y <- seq(-4, 4, 0.02)
z <- as.matrix(expand.grid(x, y), 0)
m <- length(x)
n <- length(y)

##classes are 1,2 and 3, so set contours at 1.5 and 2.5
iris.ldp <- predict(iris.lda, z)$class
contour(x, y, matrix(iris.ldp, m, n),
        levels=c(1.5, 2.5), add=TRUE, d=FALSE, lty=2)
```

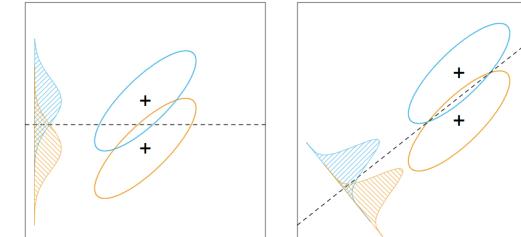
148

## Iris Dataset



149

## Fisher's Linear Discriminant Analysis



- ▶ Find a direction  $v \in \mathbb{R}^p$  to maximize the variance ratio

$$\frac{v^\top B v}{v^\top \Sigma v}$$

where

$$\Sigma = \frac{1}{n} \sum_{i=1}^n (x_i - \mu_{y_i})(x_i - \mu_{y_i})^\top \quad (\text{within class covariance})$$

$$B = \frac{1}{n} \sum_{k=1}^K n_k (\mu_{y_k} - \bar{x})(\mu_{y_k} - \bar{x})^\top \quad (\text{between class covariance})$$

$B$  has rank at most  $K - 1$ .

Figure from Hastie et al.

151

## Fisher's Linear Discriminant Analysis

- ▶ In LDA, data vectors are classified based on Mahalanobis distance from cluster means, which lie on a  $K - 1$  affine subspace.
- ▶ In measuring these distances, directions orthogonal<sup>5</sup> to the subspace can be ignored.
- ▶ Projecting data vectors onto the subspace can be viewed as a dimensionality reduction technique that preserves discriminative information about  $(y_i)_{i=1}^n$ .
- ▶ As with PCA, we can visualize the structure in the data by choosing an appropriate basis for the subspace and projecting data onto it.
- ▶ Choose a basis by finding directions that are separate classes best.

149

## Discriminant Coordinates

- ▶ To solve for the optimal  $v$ , we first reparameterize it as  $u = \Sigma^{\frac{1}{2}}v$ .

$$\frac{v^\top B v}{v^\top \Sigma v} = \frac{u^\top (\Sigma^{-\frac{1}{2}})^\top B \Sigma^{-\frac{1}{2}} u}{u^\top u} = \frac{u^\top B^* u}{u^\top u}$$

where  $B^* = (\Sigma^{-\frac{1}{2}})^\top B \Sigma^{-\frac{1}{2}}$ .

- ▶ The maximization over  $u$  is achieved by the first eigenvector  $u_1$  of  $B^*$ .
- ▶ We also look at the remaining eigenvectors  $u_l$  associated to the non-zero eigenvalues and defined the **discriminant coordinates** as  $v_l = \Sigma^{-\frac{1}{2}}u_l$ .
- ▶ The  $v_l$ 's span exactly the affine subspace spanned by  $(\Sigma^{-1}\mu_k)_{k=1}^K$  (these vectors are given as the “linear discriminants” in the R-function `lda`).

<sup>5</sup>Orthogonality defined in terms of the inner product corresponding to Mahalanobis distance:  
 $\langle x, y \rangle = x^\top \Sigma^{-1} y$ .

150

152

## Crabs Dataset

```
library(MASS)
data(crabs)

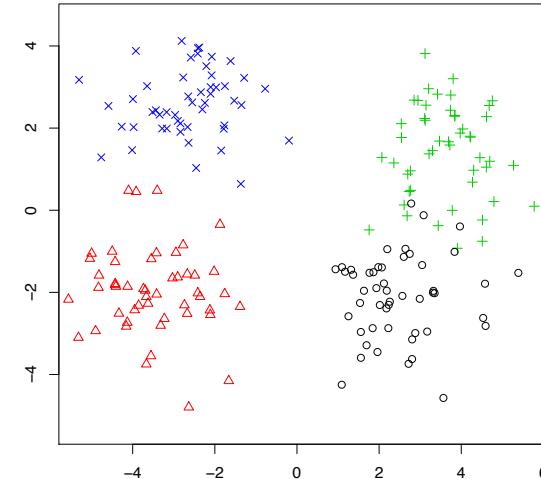
## numeric and text class labels
ct <- as.numeric(crabs[,1])-1+2*(as.numeric(crabs[,2])-1)

## Projection on Fisher's linear discriminant directions
print(cb.lda <- lda(log(crabs[,4:8]),ct))
```

153

## Crabs Dataset

```
cb.ldp <- predict(cb.lda)$x[,1:2]
eqscplot(cb.ldp,pch=ct+1,col=ct+1)
```



155

## Crabs Dataset

```
> > > > > > > Call:
lda(log(crabs[, 4:8]), ct)

Prior probabilities of groups:
 0    1    2    3 
0.25 0.25 0.25 0.25 

Group means:
      FL      RW      CL      CW      BD
0 2.564985 2.475174 3.312685 3.462327 2.441351
1 2.852455 2.683831 3.529370 3.649555 2.733273
2 2.672724 2.443774 3.437968 3.578077 2.560806
3 2.787885 2.489921 3.490431 3.589426 2.701580
```

```
Coefficients of linear discriminants:
      LD1        LD2        LD3
FL -31.217207 -2.851488 25.719750
RW -9.485303 -24.652581 -6.067361
CL -9.822169 38.578804 -31.679288
CW 65.950295 -21.375951 30.600428
BD -17.998493  6.002432 -14.541487
```

```
Proportion of trace:
   LD1    LD2    LD3
0.6891 0.3018 0.0091
```

154

## Crabs Dataset

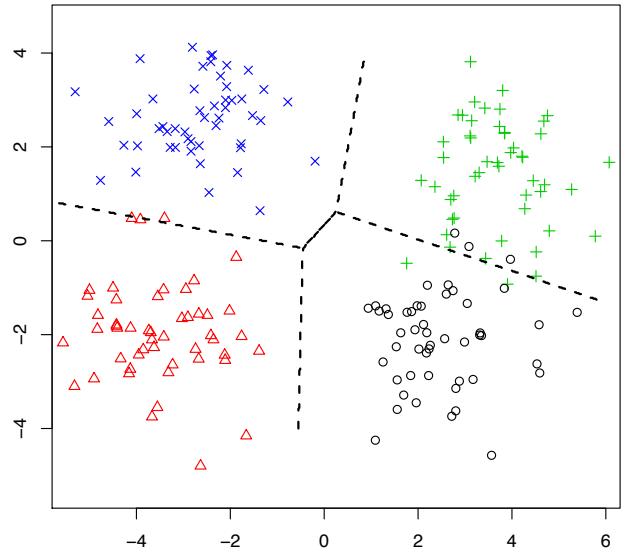
```
## display the decision boundaries
## take a lattice of points in LD-space
x <- seq(-6,6,0.02)
y <- seq(-4,4,0.02)
z <- as.matrix(expand.grid(x,y))
m <- length(x)
n <- length(y)

## perform LDA on first two discriminant directions
cb.lda <- lda(cb.lda,ct)
## predict onto the grid
cb.ldpp <- predict(cb.lda,z)$class

## classes are 0,1,2 and 3 so set contours
## at 0.5,1.5 and 2.5
contour(x,y,matrix(cb.ldpp,m,n),
         levels=c(0.5,2.5),
         add=TRUE,d=FALSE,lty=2,lwd=2)
```

156

## Crabs Dataset



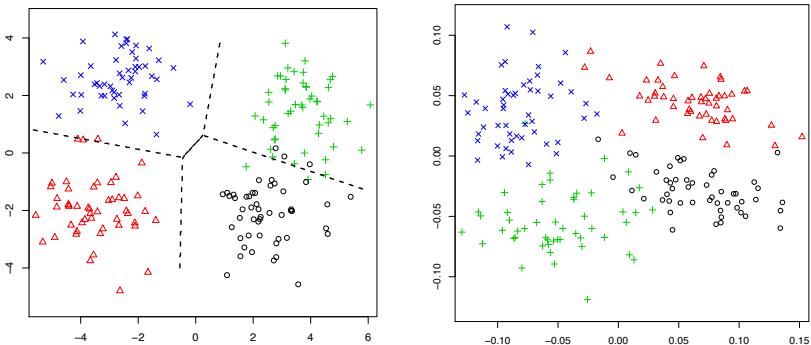
157

## Naïve Bayes

- ▶ Assume we are interested in classifying documents; e.g. scientific articles or emails.
  - ▶ A basic but standard model for text classification consists of considering a pre-specified dictionary of  $p$  words (including say physics, calculus.... or dollars, sex etc.) and summarizing each document  $i$  by a binary vector  $x_i$  where
- $$x_{ij} = \begin{cases} 1 & \text{if word } j \text{ is present in document} \\ 0 & \text{otherwise.} \end{cases}$$
- ▶ To implement a probabilistic classifier, we need to model  $f_k(x|\phi_k)$  for each class  $k = 1, \dots, K$ .

159

## Crabs Dataset



LDA separates the groups better.

158

## Naïve Bayes

- ▶ A Naïve Bayes approach ignores feature correlations and assumes  $f_k(x) = f(x|\phi_k)$  where
- $$f_k(x_i) = f(x_i|\phi_k) = \prod_{j=1}^p (\phi_{kj})^{x_{ij}} (1 - \phi_{kj})^{1-x_{ij}}$$
- ▶ Given dataset, the MLE is easily obtained
- $$\hat{\pi}_k = \frac{n_k}{n} \quad \hat{\phi}_{kj} = \frac{\sum_{i:y_i=k} x_{ij}}{n_k}$$
- ▶ One problem: if word  $j$  did not appear in documents labelled as class  $k$  then  $\hat{\phi}_{kj} = 0$  and
- $$\mathbb{P}(Y = k | X = x \text{ with } j\text{th entry equal to } 1) = 0$$
- i.e. we will never attribute a new document containing word  $j$  to class  $k$ .
- ▶ This problem is called **overfitting**, and is a major concern in modelling high-dimensional datasets common in machine learning.

160

## Generative and Discriminative Learning

- ▶ **Generative learning:** find parameters that **explains all the data.**

$$\theta^* = \operatorname{argmax}_{\theta} \sum_{i=1}^n \log p(x_i, y_i | \theta)$$

Examples: LDA, Naïve Bayes.

- ▶ Makes use of all the data.
- ▶ Flexible framework, can incorporate other tasks.
- ▶ Stronger modelling assumptions.

- ▶ **Discriminative learning:** find parameters that help to **predict relevant data.**

$$\theta^* = \operatorname{argmax}_{\theta} \sum_{i=1}^n \log p(y_i | x_i, \theta) \quad \text{or} \quad f^* = \operatorname{argmin}_f \sum_{i=1}^n L(y_i, f(x_i))$$

Examples: linear and logistic regression, rest of the course.

- ▶ Learns to perform better on the given task.
- ▶ Weaker modelling assumptions.
- ▶ Can overfitting more easily.

161

## Training and Test Performance

- ▶ **Training error** is the empirical risk

$$\frac{1}{n} \sum_{i=1}^n L(y_i, f(x_i))$$

For 0-1 loss in classification, this is the misclassification error on the training data, **which were used in learning**  $f(\underline{x})$ .

- ▶ **Test error** is the empirical risk on **new, previously unseen**, observations

$$\frac{1}{m} \sum_{i=1}^m L(y_i, f(x_i))$$

**which were NOT used in learning.**

- ▶ Test error is a much better gauge of how well learned function **generalizes** to new data.
- ▶ The test error is in general larger than the training error.

163

## Statistical Learning Theory

- ▶ We work with a joint distribution  $p^*(X, Y)$  over data vectors and labels.
- ▶ A learning algorithm constructs a function  $f(X)$  which predicts the label of  $X$ .
- ▶ Given a loss function  $L$ , the risk  $R$  of  $f(X)$  is

$$R(f) = \mathbb{E}_{X,Y}[L(Y, f(X))]$$

For classification, the best function  $f^*(X)$  is the Bayes classifier, achieving the minimum risk (Bayes risk).

- ▶ Hypothesis space  $\mathcal{H}$  is the space of functions under consideration.
- ▶ Find best function minimizing the risk:

$$\operatorname{argmin}_{f \in \mathcal{H}} \mathbb{E}_{X,Y}[L(Y, f(X))]$$

- ▶ **Empirical Risk Minimization:** minimize the empirical risk instead, since we typically do not know  $p^*(X, Y)$ .
- ▶ **Regularization:** Large hypothesis spaces can lead to overfitting,

$$\operatorname{argmin}_{f \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^n L(y_i, f(x_i)) + \lambda \|f\|_{\mathcal{H}}$$

162

## Logistic Regression

- ▶ Assume we have two classes  $\{+1, -1\}$ .
- ▶ Recall that the discriminant functions in LDA are linear. Assuming that data vectors in class  $k$  is modelled as  $\mathcal{N}(\mu_k, \Sigma)$ , choosing class  $+1$  over  $-1$  involves:

$$a_{+1} + b_{+1}^\top x > a_{-1} + b_{-1}^\top x \Leftrightarrow (a_{+1} - a_{-1}) + (b_{+1} - b_{-1})^\top x > 0$$

- ▶ If we care about minimizing classification errors, we can try to find  $a, b$  to minimize directly the average misclassification error (empirical risk associated with 0-1 loss):

$$\begin{aligned} & \operatorname{argmin}_{a,b} \frac{1}{n} \sum_{i=1}^n \begin{cases} 0 & \text{if } y_i = \operatorname{sign}(a + b^\top x) \\ 1 & \text{otherwise} \end{cases} \\ &= \operatorname{argmin}_{a,b} \frac{1}{n} \sum_{i=1}^n \frac{1}{2} - \frac{1}{2} \operatorname{sign}(y_i(a + b^\top x)) \end{aligned}$$

- ▶ An example of **Empirical Risk Minimization**. Unfortunately not typically possible to solve...

164

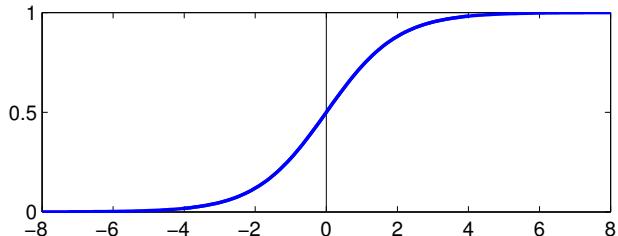
## Logistic Regression

- Logistic regression replaces the 0-1 loss with the log loss.
- A model parameterizing the conditional distribution of labels given data vectors:

$$p(Y = 1|X = x) = \frac{1}{1 + \exp(-(a + b^\top x))} =: s(a + b^\top x)$$

$$p(Y = -1|X = x) = \frac{1}{1 + \exp(+ (a + b^\top x))} = s(-a - b^\top x)$$

where  $s(\cdot)$  is the **logistic function**



165

## Logistic Regression

- Not possible to find optimal  $a, b$  analytically.
- For simplicity, absorb  $a$  as an entry in  $b$  by appending '1' into  $x$  vector.
- Objective function:

$$R_{\log}^{\text{emp}} = \frac{1}{n} \sum_{i=1}^n -\log s(y_i x_i^\top b)$$

### Logistic Function

$$s(-z) = 1 - s(z)$$

$$\nabla_z s(z) = s(z)s(-z)$$

$$\nabla_z \log s(z) = s(-z)$$

$$\nabla_z^2 \log s(z) = -s(z)s(-z)$$

- Differentiate wrt  $b$ :

$$\nabla_b R_{\log}^{\text{emp}} = \frac{1}{n} \sum_{i=1}^n -s(-y_i x_i^\top b) y_i x_i = \frac{1}{n} \sum_{i=1}^n -((.5 + .5y_i) - s(x_i^\top b)) x_i$$

$$\nabla_b^2 R_{\log}^{\text{emp}} = \frac{1}{n} \sum_{i=1}^n s(y_i x_i^\top b) s(-y_i x_i^\top b) x_i x_i^\top$$

167

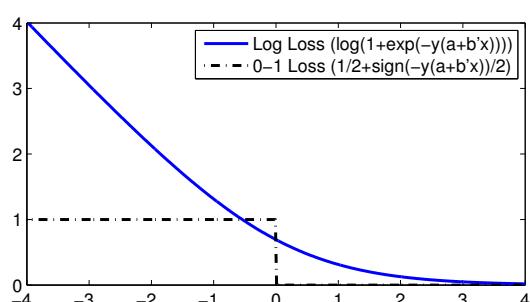
## Logistic Regression

- Consider maximizing the **conditional log likelihood**:

$$\ell(a, b) = \sum_{i=1}^n \log p(Y = y_i | X = x_i) = \sum_{i=1}^n -\log(1 + \exp(-y_i(a + b^\top x_i)))$$

- Equivalent to minimizing the empirical risk associated with the **log loss**:

$$R_{\log}^{\text{emp}} = \frac{1}{n} \sum_{i=1}^n \log(1 + \exp(-y_i(a + b^\top x_i)))$$



166

## Logistic Regression

- Second derivative is positive-definite: objective function is **convex** and there is a **single unique global minimum**.
- Many different algorithms can find optimal  $b$ , e.g.:

- Gradient descent:

$$b^{\text{new}} = b + \epsilon \frac{1}{n} \sum_{i=1}^n s(-y_i x_i^\top b) y_i x_i$$

- Stochastic gradient descent:

$$b^{\text{new}} = b + \epsilon_t \frac{1}{|I(t)|} \sum_{i \in I(t)} s(-y_i x_i^\top b) y_i x_i$$

where  $I(t)$  is a subset of the data at iteration  $t$ , and  $\epsilon_t \rightarrow 0$  slowly ( $\sum_t \epsilon_t = \infty, \sum_t \epsilon_t^2 < \infty$ ).

- Newton-Raphson:

$$b^{\text{new}} = b - (\nabla_b^2 R_{\log}^{\text{emp}})^{-1} \nabla_b R_{\log}^{\text{emp}}$$

This is also called **iterative reweighted least squares**.

- Conjugate gradient, LBGFS and other methods from numerical analysis.

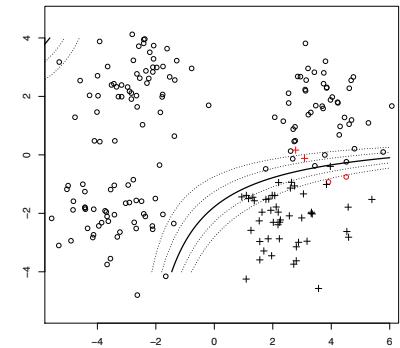
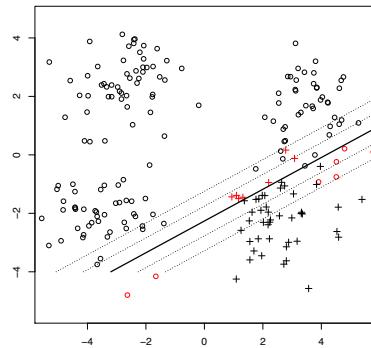
168

## Logistic Regression

Properties of logistic regression:

- ▶ Makes less modelling assumptions than LDA and naïve Bayes.
- ▶ Models only the conditional distribution of labels, not the marginal distribution of  $X$ .
- ▶ A linear method: decision boundary is a separating hyperplane.
- ▶ Logistic regression can be made **non-linear** by applying a non-linear transformation  $X \mapsto \phi(X)$ .
- ▶ Logistic regression is a simple example of a generalised linear model (GLM). Much statistical theory:
  - ▶ assessment of fit via deviance and plots,
  - ▶ interpretation of entries of  $b$  as **odds-ratios**,
  - ▶ fitting categorical data (sometimes called **multinomial logistic regression**),
  - ▶ well founded approaches to removing insignificant features (drop-in deviance test, Wald test),

## Crab Dataset

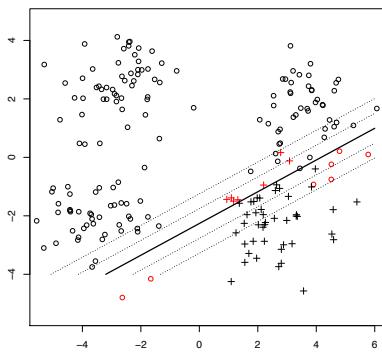
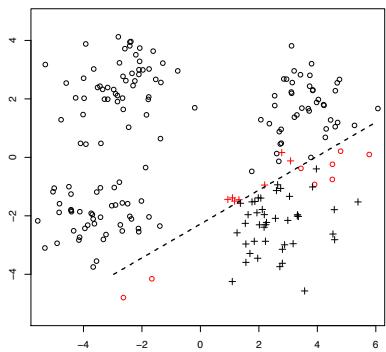


Comparing logistic regression with and without quadratic interactions.

169

171

## Crab Dataset



Comparing LDA and logistic regression.

170

## Crab Dataset

```
library(MASS)
## load crabs data
data(crabs)
ct <- as.numeric(crabs[,1])-1+2*(as.numeric(crabs[,2])-1)
## project into first two LD
cb.lda <- lda(log(crabs[,4:8]),ct)
cb.ldp <- predict(cb.lda)
x <- cb.ldp$x[,1:2]
y <- as.numeric(ct==0)
eqscplot(x,pch=2*y+1,col=y+1)

## visualize decision boundary
gx1 <- seq(-6,6,.02)
gx2 <- seq(-4,4,.02)
gx <- as.matrix(expand.grid(gx1,gx2))
gm <- length(gx1)
gn <- length(gx2)
gdf <- data.frame(LD1=gx[,1],LD2=gx[,2])

lda <- lda(x,y)
y.lda <- predict(lda,x)$class
eqscplot(x,pch=2*y+1,col=2-as.numeric(y==y.lda))
y.lda.grid <- predict(lda,gdf)$class
contour(gx1,gx2,matrix(y.lda.grid,gm,gn),
        levels=c(0.5), add=TRUE, d=FALSE, lty=2, lwd=2)
```

172

## Crab Dataset

```
## logistic regression
xdf <- data.frame(x)
logreg <- glm(y ~ LD1 + LD2, data=xdf, family=binomial)
y.lr <- predict(logreg,type="response")
eqscplot(x,pch=2*y+1,col=2-as.numeric(y==(y.lr>.5)))
y.lr.grid <- predict(logreg,newdata=gdf,type="response")
contour(gx1,gx2,matrix(y.lr.grid,gn,gn),
        levels=c(.1,.25,.75,.9), add=TRUE,d=FALSE,lty=3,lwd=1)
contour(gx1,gx2,matrix(y.lr.grid,gn,gn),
        levels=c(.5), add=TRUE,d=FALSE,lty=1,lwd=2)

## logistic regression with quadratic interactions
logreg <- glm(y ~ (LD1 + LD2)^2, data=xdf, family=binomial)
y.lr <- predict(logreg,type="response")
eqscplot(x,pch=2*y+1,col=2-as.numeric(y==(y.lr>.5)))
y.lr.grid <- predict(logreg,newdata=gdf,type="response")
contour(gx1,gx2,matrix(y.lr.grid,gn,gn),
        levels=c(.1,.25,.75,.9), add=TRUE,d=FALSE,lty=3,lwd=1)
contour(gx1,gx2,matrix(y.lr.grid,gn,gn),
        levels=c(.5), add=TRUE,d=FALSE,lty=1,lwd=2)
```

173

## Spam Dataset

Use logistic regression to predict spam/not spam.

```
library(kernlab)
data(spam)

## let Y=0 be non-spam and Y=1 be spam.
Y <- as.numeric(spam[, ncol(spam)])-1
X <- spam[, -ncol(spam)]

g1 <- glm(Y ~ ., data=X,family=binomial)
```

Which predictor variables seem to be important? Can for example check which ones are significant in the GLM.

```
> summary(g1)
Call:
glm(formula = Y ~ ., family = binomial, data = X)

Deviance Residuals:
    Min      1Q   Median      3Q     Max 
-4.127e+00 -2.030e-01 -1.967e-06  1.140e-01  5.364e+00
```

175

## Spam Dataset

```
> library(kernlab)
> data(spam)
> dim(spam)
[1] 4601 58

> spam[1:2,]
   make address all num3d our over remove internet order mail receive wil
1 0.00 0.64 0.64 0 0.32 0.00 0.00 0.00 0 0.00 0.00 0.6
2 0.21 0.28 0.50 0 0.14 0.28 0.21 0.07 0 0.94 0.21 0.7
   people report addresses free business email you credit your font num00c
1 0.00 0.00 0.00 0.32 0.00 1.29 1.93 0 0.96 0 0.00
2 0.65 0.21 0.14 0.14 0.07 0.28 3.47 0 1.59 0 0.43
   money hp hpl george num650 lab labs telnet num857 data num415 num85
1 0.00 0 0 0 0 0 0 0 0 0 0 0
2 0.43 0 0 0 0 0 0 0 0 0 0 0
   technology num1999 parts pm direct cs meeting original project re edu ta
1 0 0.00 0 0 0 0 0 0 0 0 0 0
2 0 0.07 0 0 0 0 0 0 0 0 0 0
   conference charSemicolon charRoundbracket charSquarebracket charExclamat
1 0 0 0.000 0 0 0.778
2 0 0 0.132 0 0 0.372
   charDollar charHash capitalAve capitalLong capitalTotal type
1 0.00 0.000 3.756 61 278 spam
2 0.18 0.048 5.114 101 1028 spam
```

174

## Spam Dataset

Coefficients:

	Estimate	Std. Error	z value	Pr(> z )
(Intercept)	-1.569e+00	1.420e-01	-11.044	< 2e-16 ***
make	-3.895e-01	2.315e-01	-1.683	0.092388 .
address	-1.458e-01	6.928e-02	-2.104	0.035362 *
all	1.141e-01	1.103e-01	1.035	0.300759
num3d	2.252e+00	1.507e+00	1.494	0.135168
our	5.624e-01	1.018e-01	5.524	3.31e-08 ***
over	8.830e-01	2.498e-01	3.534	0.000409 ***
remove	2.279e+00	3.328e-01	6.846	7.57e-12 ***
internet	5.696e-01	1.682e-01	3.387	0.000707 ***
order	7.343e-01	2.849e-01	2.577	0.009958 **
mail	1.275e-01	7.262e-02	1.755	0.079230 .
receive	-2.557e-01	2.979e-01	-0.858	0.390655
will	-1.383e-01	7.405e-02	-1.868	0.061773 .
people	-7.961e-02	2.303e-01	-0.346	0.729557
report	1.447e-01	1.364e-01	1.061	0.288855
addresses	1.236e+00	7.254e-01	1.704	0.088370 .
business	9.599e-01	2.251e-01	4.264	2.01e-05 ***
email	1.203e-01	1.172e-01	1.027	0.304533
you	8.131e-02	3.505e-02	2.320	0.020334 *
credit	1.047e+00	5.383e-01	1.946	0.051675 .

176

## Spam Dataset

```
your          2.419e-01 5.243e-02 4.615 3.94e-06 ***
font          2.013e-01 1.627e-01 1.238 0.215838
num000        2.245e+00 4.714e-01 4.762 1.91e-06 ***
money         4.264e-01 1.621e-01 2.630 0.008535 **
hp            -1.920e+00 3.128e-01 -6.139 8.31e-10 ***
hpl           -1.040e+00 4.396e-01 -2.366 0.017966 *
george        -1.177e+01 2.113e+00 -5.569 2.57e-08 ***
num650        4.454e-01 1.991e-01 2.237 0.025255 *
lab            -2.486e+00 1.502e+00 -1.656 0.097744 .
labs           -3.299e-01 3.137e-01 -1.052 0.292972
telnet         -1.702e-01 4.815e-01 -0.353 0.723742
num857        2.549e+00 3.283e+00 0.776 0.437566
data           -7.383e-01 3.117e-01 -2.369 0.017842 *
num415        6.679e-01 1.601e+00 0.417 0.676490
num85          -2.055e+00 7.883e-01 -2.607 0.009124 **
technology    9.237e-01 3.091e-01 2.989 0.002803 **
num1999       4.651e-02 1.754e-01 0.265 0.790819
parts          -5.968e-01 4.232e-01 -1.410 0.158473
pm             -8.650e-01 3.828e-01 -2.260 0.023844 *
direct         -3.046e-01 3.636e-01 -0.838 0.402215
cs              -4.505e+01 2.660e+01 -1.694 0.090333 .
meeting        -2.689e+00 8.384e-01 -3.207 0.001342 **
original       -1.247e+00 8.064e-01 -1.547 0.121978
project        -1.573e+00 5.292e-01 -2.973 0.002953 **
re              -7.923e-01 1.556e-01 -5.091 3.56e-07 ***
```

177

## Spam Dataset

How good is the classification?

```
> proba <- predict(gl,type="response")
> predicted_spam <- as.numeric( proba>0.5)
> table(predicted_spam,Y)
   Y
predicted_spam   0   1
      0 2666 194
      1 122 1619

> predicted_spam <- as.numeric( proba>0.99)
> table(predicted_spam,Y)
   Y
predicted_spam   0   1
      0 2776 1095
      1   12  718
```

So out of 730 emails marked as spam, 12 were actually not spam.  
Advantage of a probabilistic approach: probabilities give interpretable confidence to predictions.

179

## Spam Dataset

```
edu          -1.459e+00 2.686e-01 -5.434 5.52e-08 ***
table         -2.326e+00 1.659e+00 -1.402 0.160958
conference    -4.016e+00 1.611e+00 -2.493 0.012672 *
charSemicolon -1.291e+00 4.422e-01 -2.920 0.003503 **
charRoundbracket -1.881e-01 2.494e-01 -0.754 0.450663
charSquarebracket -6.574e-01 8.383e-01 -0.784 0.432914
charExclamation 3.472e-01 8.926e-02 3.890 0.000100 ***
charDollar     5.336e+00 7.064e-01 7.553 4.24e-14 ***
charHash       2.403e+00 1.113e+00 2.159 0.030883 *
capitalAve    1.199e-02 1.884e-02 0.636 0.524509
capitalLong   9.118e-03 2.521e-03 3.618 0.000297 ***
capitalTotal  8.437e-04 2.251e-04 3.747 0.000179 ***
```

---

Signif. codes: 0 '\*\*\*' 0.001 '\*\*' 0.01 '\*' 0.05 '.' 0.1 ' ' 1

(Dispersion parameter for binomial family taken to be 1)

Null deviance: 6170.2 on 4600 degrees of freedom  
Residual deviance: 1815.8 on 4543 degrees of freedom  
AIC: 1931.8

Number of Fisher Scoring iterations: 13

178

## Spam Dataset

Success rate is calculated on the same data that the GLM is trained on!  
Separate in training and test set.

```
n <- length(Y)
i <- sample( rep(c(TRUE,FALSE),each=n/2) ,round(n) ,replace=FALSE )
train <- (1:n)[i]
test  <- (1:n)[!i]
```

Fit only on training set and predict on both training and test set.

```
gl <- glm(Y[train] ~ ., data=X[train,],family=binomial)

proba_train <- predict(gl,newdata=X[train,],type="response")
proba_test  <- predict(gl,newdata=X[test,],type="response")

predicted_spam_train <- as.numeric(proba_train > 0.95)
predicted_spam_test  <- as.numeric(proba_test > 0.95)
```

180

## Spam Dataset

Results for training and test set:

```
> table(predicted_spam_train, Y[train])
predicted_spam_train 0 1
0 1403 354
1 11 567

> table(predicted_spam_test, Y[test])
predicted_spam_test 0 1
0 1346 351
1 28 541
```

It is no coincidence that test performance is worse than training performance.

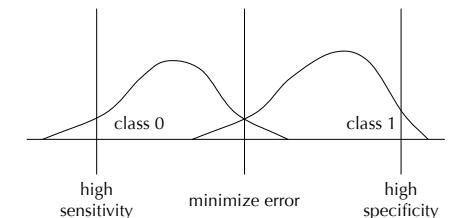
181

## Performance Measures

### ► Confusion matrix:

Prediction	True state	0	1
		# true negative	# false negative
0	# false positive	# true positive	
1			

- Accuracy:  $(TP + TN) / (TP + TN + FP + FN)$ .
- Error rate:  $(FP + FN) / (TP + TN + FP + FN)$ .
- Sensitivity (true positive rate):  $TP / (TP + FN)$ .
- Specificity (true negative rate):  $TN / (TN + FP)$ .
- Precision:  $TP / (TP + FP)$ .
- Recall:  $TP / (TP + FN)$ .
- F1: harmonic mean of precision and recall.



183

- As we vary the prediction threshold  $c$  from 0 to 1:
  - Specificity varies from 0 to 1.
  - Sensitivity goes from 1 to 0.

## Spam Dataset

Compare with LDA.

```
library(MASS)
lda_res <- lda(x=X[train], grouping=Y[train])

proba_lda <- predict(lda_res, newdata=X[test,])$posterior[,2]
predicted_spam_lda <- as.numeric(proba_lda > 0.95)

> table(predicted_spam_test, Y[test])
predicted_spam_test 0 1
0 1346 351
1 28 541

> table(predicted_spam_lda, Y[test])
predicted_spam_lda 0 1
0 1364 533
1 10 359
```

It seems as if LDA beats logistic regression here, but would need to adjust decision threshold to get proper comparison. Use **ROC curves**.

182

## ROC Curves

ROC curve plots sensitivity versus specificity as threshold varies.

```
cvec <- seq(0.001, 0.999, length=1000)
specif <- numeric(length(cvec))
sensit <- numeric(length(cvec))

for (cc in 1:length(cvec)){
  sensit[cc] <- sum( proba_lda > cvec[cc] & Y[test]==1) / sum(Y[test]==1)
  specif[cc] <- sum( proba_lda <=cvec[cc] & Y[test]==0) / sum(Y[test]==0)
}

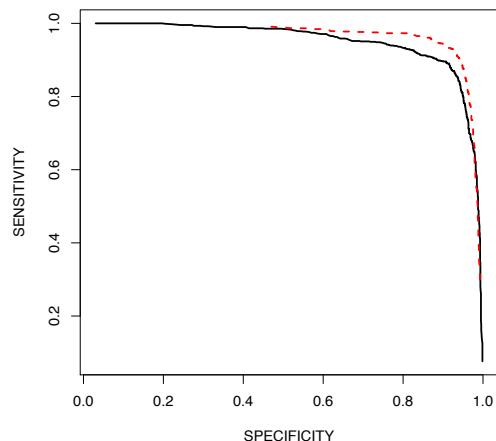
plot(specif, sensit, xlab="SPECIFICITY", ylab="SENSITIVITY", type="l", lwd=2)
```

184

## ROC Curves

ROC curve for LDA and logistic regression classification of spam dataset.

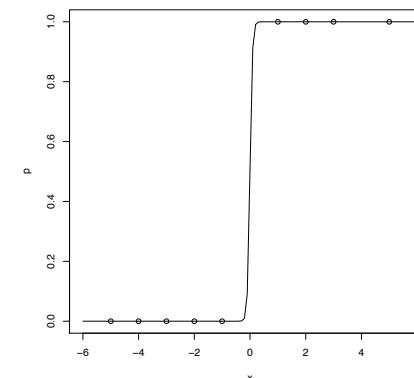
LDA = unbroken black line; LR = broken red line.



Obvious now that LR is better for this dataset than LDA, contrary to the first impression.

185

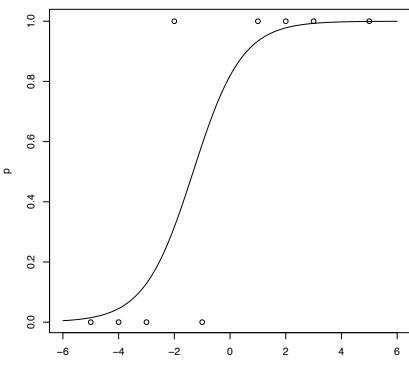
## Overfitting in Logistic Regression



```
dx <- c(-5, -4, -3, -2, -1, 1, 2, 3, 5)
d <- data.frame(dx)
x <- seq(-6, 6, .1)
y <- c(0, 0, 0, 0, 0, 1, 1, 1, 1)
lr <- glm(y ~ ., data=d, family=binomial)
p <- predict(lr, newdata=data.frame(dx=grid), type="response")
plot(x, p, type="l")
points(dx, y)
```

187

## Overfitting in Logistic Regression

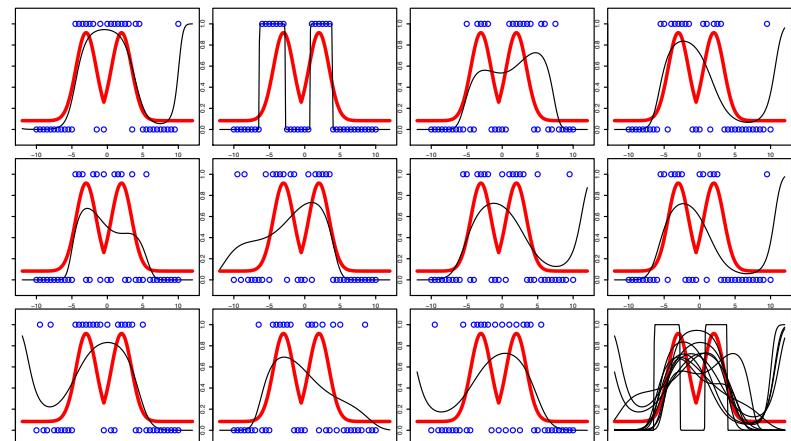


```
dx <- c(-5, -4, -3, -2, -1, 1, 2, 3, 5)
d <- data.frame(dx)
x <- seq(-6, 6, .1)
y <- c(0, 0, 0, 1, 0, 1, 1, 1, 1)
lr <- glm(y ~ ., data=d, family=binomial)
p <- predict(lr, newdata=data.frame(dx=x), type="response")
plot(x, p, type="l")
points(dx, y)
```

186

## Demo on Overfitting in Logistic Regression

True conditional probabilities in Red. Blue circles are training data, Black curve is predicted conditional probability. 11 datasets are sampled from true distribution and used to learn a logistic regression model with non-linear features  $\phi(x) = (1, x, x^2, \dots, x^{p-1})$ .



188

## Demo on Overfitting in Logistic Regression

```

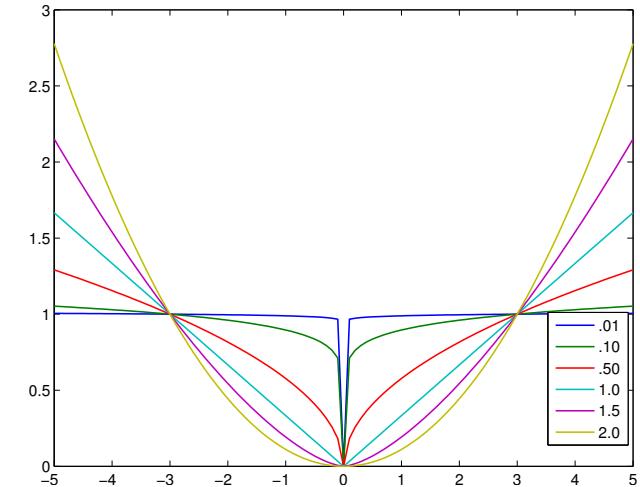
## true conditional probabilities
truep <- function(x) {
  return((pmax(exp(-(x-2)^2/4),exp(-(x+3)^2/4))+.1)/1.2)
}
## features are (x^i)
phi <- function(x,deg) {
  d <- matrix(0,length(x),deg+1)
  for (i in 0:deg) {
    d[,i+1] <- x ^ i
  }
  return (data.frame(d))
}
## demo learning logistic regression, with different datasets generated,
## and using different degree polynomials as features
demolearn <- function(trainx,testx,truep,deg) {
  trainp <- truep(trainx)
  testp <- truep(testx)
  par(mfrow=c(3,4),ann=FALSE,cex=.3,mar=c(1,1,1,1))
  predp <- matrix(0,length(testx),11)
  for (i in 1:11) {
    trainy <- as.numeric(runif(length(trainx)) < trainp)
    lr <- glm(trainy ~ .,data=phi(trainx,deg),family=binomial)
    predp[,i] <- predict(lr,newdata=phi(testx,deg),type="response")
    plot(testx,testp,type="l",lwd=3,col=2,ylim=c(-.1,1.1))
    lines(testx,predp[,i],type="l")
    points(trainx,trainy,pch=1,col=4,cex=2)
  }
  plot(testx,testp,type="l",lwd=3,col=2,ylim=c(-.1,1.1))
  for (i in 1:11) {
    lines(testx,predp[,i],type="l")
  }
  return(predp)
}

trainx <- seq(-10,10,.5)
testx <- seq(-12,12,.1)
pp <- demolearn(trainx,testx,truep,4)

```

189

## Regularization



$L_\rho$  regularization profile for different values of  $\rho$ .

191

## Regularization

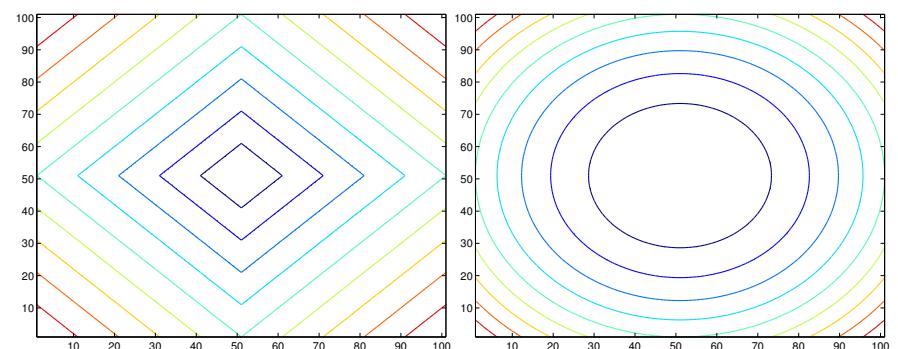
- ▶ Flexible models for high-dimensional problems require many parameters.
- ▶ With many parameters, learners can easily overfit to the noise in the training data.
- ▶ **Regularization:** Limit flexibility of model to prevent overfitting.
- ▶ Typically: add term **penalizing** large values of parameters  $\theta$ .

$$R^{\text{emp}}(\theta) + \lambda \|\theta\|_\rho^\rho = \frac{1}{n} \sum_{i=1}^n \log(1 + \exp(-y_i(a + b^\top x_i))) + \lambda \|b\|_\rho^\rho$$

where  $\rho \in [1, 2]$ , and  $\|z\|_\rho = (\sum_{j=1}^p |z_j|^\rho)^{1/\rho}$  is the  $L_\rho$  norm of  $b$  (also of interest when  $\rho \in [0, 1)$ , but is no longer a norm).

- ▶ Also known as **shrinkage** methods—parameters are shrunk towards 0.
- ▶ Typical cases are  $\rho = 2$  (Euclidean norm, **ridge regression**) and  $\rho = 1$  (**LASSO**). When  $\rho \leq 1$  it is called a **sparsity** inducing regularization.
- ▶  $\lambda$  is a **tuning parameter** (or **hyperparameter**) and controls the amount of regularization, and resulting complexity of the model.

## Regularization



$L_1$  and  $L_2$  norm contours.

192

## Sparsity Inducing Regularization

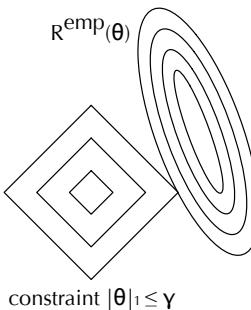
- Consider constrained optimization problem

$$\min_{\theta} R^{\text{emp}}(\theta) \text{ s.t. } \|\theta\|_1 < \gamma$$

- Lagrange multiplier  $\lambda > 0$  to enforce constraint,

$$\min_{\theta} R^{\text{emp}}(\theta) + \lambda(\|\theta\|_1 - \gamma)$$

- At the optimal value of  $\lambda$ , the parameter  $\theta$  is the one minimizing the regularized empirical risk objective.
- Conversely, given  $\lambda$ , there is a value of  $\gamma$  such that the corresponding optimal Lagrange multiplier is  $\lambda$ .
- Using  $L_1$  regularization, optimal  $\theta$  has  $\theta_2 = 0$ .
- Generally:  $L_1$  regularization leads to optimal solutions with many zeros, i.e. the regression function depends only on the (small) number of features with non-zero parameters.



193

## Demo on $L_1$ Regularized Logistic Regression

```

## true conditional probabilities
truep <- function(x) {
  return((pmax(exp(-(x-2)^2/4), exp(-(x+3)^2/4))+.1)/1.2)
}
## features are (x^i)
phi <- function(x, deg) {
  d <- matrix(0,length(x),deg+1)
  for (i in 0:deg) {
    d[,i+1] <- x ^ i
  }
  return (data.frame(d))
}
## demo L1 regularized learning of logistic regression,
## with different datasets generated, and using different
## degree polynomials as features

trainx <- seq(-10,10,.5)
testx <- seq(-12,12,.1)

demolearnL1 <- function(trainx,testx,truep,deg) {
  trainp <- truep(trainx)
  testp <- truep(testx)
  trainy <- as.numeric(runif(length(trainx)) < trainp)
  slr <- glmnet(as.matrix(phi(trainx,deg)),as.factor(trainy),
                family="binomial")
  s <- c(0,.0001,.001,.01,.05)
  predp <- predict(slr,newx=as.matrix(phi(testx,deg)),
                    s=s,types="response")
  par(mfrow=(1,2),mar=c(4,4,1,2))
  plot(testx,testp,type="l",col=2,lwd=3,ylim=c(-.1,1.1))
  points(trainx,trainy,pch=1,col=4,cex=2)
  for (i in 1:dim(predp)[2]) {
    lines(testx,predp[,i],type="l")
  }
  plot(slr,xvar="lambda")
  print(coef(slr,s))
  return(predp)
}

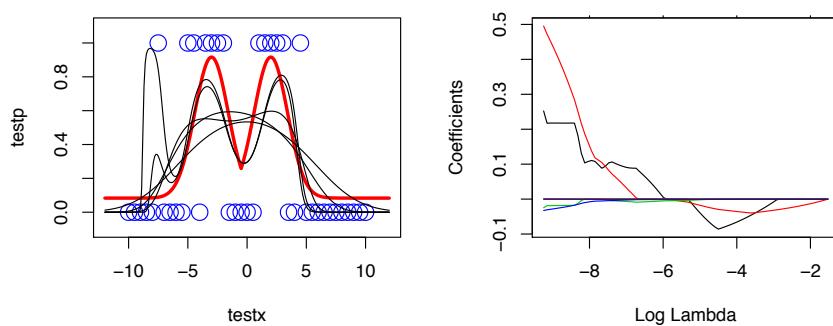
demolearnL1(trainx,testx,truep,10)

```

195

## Optimization

Use `glmnet` for regression with  $L_1$ ,  $L_2$  and combination regularization.



194

- Many more complex models in statistics and machine learning do not have analytic solutions to ML estimators.
- In most models parameters are learned by some numerical optimization technique.

$$\min_{\theta} F(\theta)$$

- How many minima are there?
- How do we find optimal  $\theta$ ?
- Are we guaranteed to find the global optimum  $\theta^*$ , rather just a local one?
- How efficiently can we solve for  $\theta$ ?
- What if there are constraints?

196

## Constrained Optimization

- Optimization problems with constraints, e.g.

$$\begin{aligned} \min_{\theta \in \mathbb{R}^d} \quad & F(\theta) \\ \text{subject to} \quad & g_i(\theta) \leq 0 \quad \text{for } i = 1, \dots, I \\ & h_j(\theta) = 0 \quad \text{for } j = 1, \dots, J \end{aligned}$$

where  $g_i$  enforce inequality constraints and  $h_j$  equality constraints.

- Can write this succinctly:

$$\begin{aligned} \min_{\theta \in \mathbb{R}^d} \quad & F(\theta) \\ \text{subject to} \quad & g(\theta) \preceq 0 \\ & h(\theta) = 0 \end{aligned}$$

where  $g : \mathbb{R}^d \rightarrow \mathbb{R}^I$  is a vector-valued function with  $g(\theta)_i = g_i(\theta)$ . Similarly  $h(\theta) : \mathbb{R}^d \rightarrow \mathbb{R}^J$ .  $x \preceq y$  iff  $x_i \leq y_i \forall i$ .

- These problems are called **programmes**.

197

A function  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  is **convex** if

$$f(\alpha x + (1 - \alpha)y) \leq \alpha f(x) + (1 - \alpha)f(y)$$

for all  $x, y \in \mathbb{R}^d$ ,  $\alpha \in [0, 1]$ .

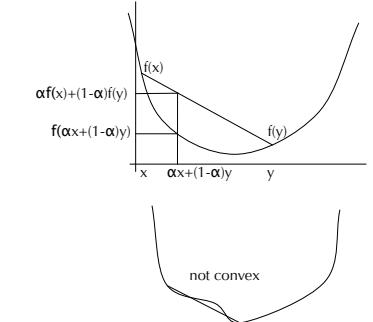
- For smooth functions: Equivalent to 2nd derivative (**Hessian**) being positive semidefinite.

- A programme is a **convex programme** if:

- $F(\theta)$  is convex,
- $g_i(\theta)$  is convex for each  $i$ ,
- $h(\theta) = A\theta + b$  is affine.

- Examples: linear, quadratic, semidefinite programming.

- Convex programmes have a unique minimum (typically), which can be efficiently found.



Boyd and Vandenberghe, Convex Optimization. 2004. MOOC right now.

199

## Constrained Optimization

$$\begin{aligned} \min_{\theta \in \mathbb{R}^d} \quad & F(\theta) \\ \text{subject to} \quad & g(\theta) \preceq 0 \\ & h(\theta) = 0 \end{aligned}$$

- We can enforce constraints by using **Lagrange multipliers** or **dual variables**  $\lambda \in \mathbb{R}^I$  and  $\kappa \in \mathbb{R}^J$ .
- The optimization problem can be written as a mini-max optimization of the Lagrangian:

$$\min_{\theta} \max_{\lambda \succeq 0, \kappa} \mathcal{L}(\theta, \lambda, \kappa) = \min_{\theta} \max_{\lambda \succeq 0, \kappa} F(\theta) + \lambda^\top g(\theta) + \kappa^\top h(\theta)$$

- Intuition: For any  $\theta$ , we have:

$$\max_{\lambda \succeq 0, \kappa} \mathcal{L}(\theta, \lambda, \kappa) = \begin{cases} +\infty & \text{if there is some unsatisfied constraint,} \\ F(\theta) & \text{if all constraints are satisfied.} \end{cases}$$

So the outer minimization over  $\theta$  results in the same optimization problem.

198

## Convex Duality

- Say the minimum is  $p^*$ , and occurred at  $\theta^*$ .

- The **dual programme** inverts the order of max and min:

$$p^* = \min_{\theta} \max_{\lambda \succeq 0, \kappa} \mathcal{L}(\theta, \lambda, \kappa) \geq \max_{\lambda \succeq 0, \kappa} \min_{\theta} \mathcal{L}(\theta, \lambda, \kappa) = d^*$$

where the dual optimum is  $d^*$ .

- **Karush-Kuhn-Tucker Theorem:** Subject to regularity conditions, a solution  $\theta^*$  is the optimal solution of a convex programme, if and only if there are  $\lambda^*$  and  $\kappa^*$  (the dual optimal solution) such that:

- **Primal feasible:**  $g(\theta^*) \preceq 0$ ,  $h(\theta^*) = 0$ .
- **Dual feasible:**  $\lambda^* \succeq 0$ .
- $(\theta^*, \lambda^*, \kappa^*)$  is a **saddle point** of  $\mathcal{L}$ : For every  $\theta, \lambda \succeq 0, \kappa$ , we have

$$\mathcal{L}(\theta^*, \lambda, \kappa) \leq \mathcal{L}(\theta^*, \lambda^*, \kappa^*) \leq \mathcal{L}(\theta, \lambda^*, \kappa^*)$$

$$\nabla_{\theta} \mathcal{L}(\theta^*, \lambda^*, \kappa^*) = \nabla_{\theta} F(\theta^*) + (\lambda^*)^\top \nabla_{\theta} g(\theta^*) + (\kappa^*)^\top \nabla_{\theta} h(\theta^*) = 0$$

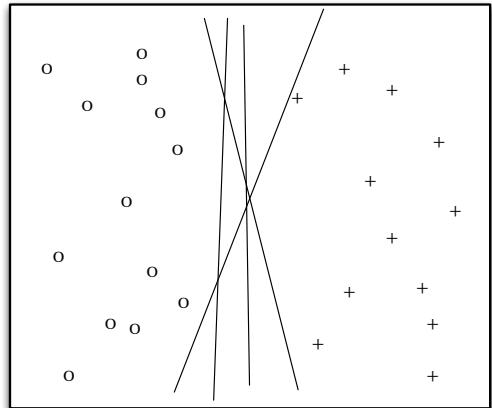
- **Complementary slackness:** For every  $i$ ,

$$\lambda_i^* g_i(\theta^*) = 0$$

200

## Linear Classification

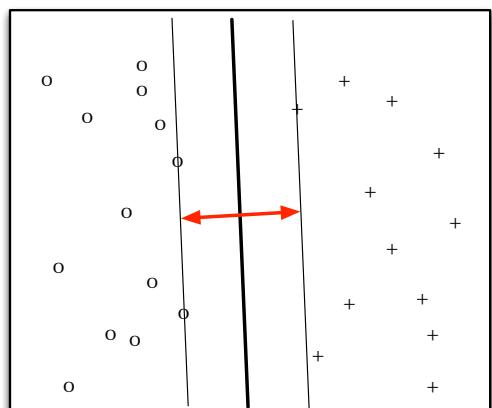
- A dataset with  $\{+1, -1\}$  labels is **linearly separable** if there is a hyperplane separating two classes.
- Typically there will be an infinite number of such **separating hyperplanes**.



201

## Maximum Margin Classification

- Good choice of separating hyperplane: one with **large margin**.
- Such a hyperplane will be defined by a number of data vectors close to the boundary—the **support vectors**, leading to a method called **support vector machines**.



202

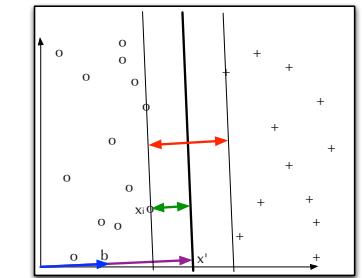
## Support Vector Machines

- A hyperplane can be parametrized as:

$$g(x) = a + b^\top x = 0$$

with the classification given by  $\text{sign}(g(x))$ .

- Distance and classification of a point  $x_i$  from hyperplane is  $y_i(a + b^\top x_i)/\|b\|$ .
- Multiplying  $a$  and  $b$  by  $c > 0$  does not affect result. Rescale such that margin (closest distance of data vectors to hyperplane) is  $1/\|b\|$ .



$$y_i(a + b^\top x_i)/\|b\| \geq 1/\|b\|$$

$$y_i(a + b^\top x_i) \geq 1$$

- Constrained optimization problem to solve for  $a, b$ :

$$\max_{a,b} \frac{1}{\|b\|} \quad \Leftrightarrow \quad \min_{a,b} \frac{1}{2} \|b\|^2$$

$$\text{subject to } y_i(a + b^\top x_i) \geq 1 \quad \text{subject to } y_i(a + b^\top x_i) \geq 1 \quad \text{for all } i$$

203

## Support Vector Machines

- Introduce Lagrange multipliers  $\lambda_i \geq 0$  to enforce constraints:

$$\min_{a,b} \max_{\lambda \geq 0} \mathcal{L}(a,b,\lambda) = \frac{1}{2} \|b\|^2 + \sum_{i=1}^n \lambda_i (1 - y_i(a + b^\top x_i))$$

- KKT optimality conditions:

Zero derivatives:

$$\nabla_a \mathcal{L}(a^*, b^*, \lambda^*) = - \sum_{i=1}^n \lambda_i^* y_i x_i = 0$$

$$\nabla_b \mathcal{L}(a^*, b^*, \lambda^*) = b^* - \sum_{i=1}^n \lambda_i^* y_i x_i = 0$$

Primal feasibility:

$$y_i(a^* + (b^*)^\top x_i) \geq 1$$

Dual feasibility:

$$\lambda_i^* \geq 0$$

Complementary slackness:

$$\lambda_i^* (1 - y_i(a^* + (b^*)^\top x_i)) = 0$$

204

## Support Vector Machines

- Substituting optimal  $a^*$  and  $b^*$  into Lagrangian leads to the **dual optimization problem**:

$$\begin{aligned} \max_{\lambda} \quad & \sum_{i=1}^n \lambda_i - \frac{1}{2} \sum_{i,j=1}^n \lambda_i \lambda_j y_i y_j (x_i)^\top (x_j) \\ \text{subject to} \quad & \sum_{i=1}^n \lambda_i y_i = 0 \\ & \lambda \succeq 0 \end{aligned}$$

**A quadratic programme.** Standard solvers can be used to find optimal  $\lambda^*$  in  $O(n^3)$  cost.

- Those vectors with  $\lambda_i > 0$  are called **support vectors**.
- Complementary slackness implies that if  $x_i$  does not lie on boundary, then  $\lambda_i = 0$ , i.e. not a support vector.
- Discriminant function is

$$g(x) = a^* + \sum_{i=1}^n \lambda_i^* y_i x_i^\top x$$

where  $a^*$  can be solved by noting that  $y_j g(x_j) = 1$  for a support vector  $x_j$ .

205

## Soft-Margin Support Vector Machines

- Introduce Lagrange multipliers  $\lambda_i \geq 0, \gamma_i \geq 0$  to enforce constraints:

$$\mathcal{L}(a, b, \xi, \lambda, \gamma) = \frac{1}{2} \|b\|^2 + C \sum_{i=1}^n \xi_i + \sum_{i=1}^n \lambda_i (1 - \xi_i - y_i(a + b^\top x_i)) - \sum_{i=1}^n \gamma_i \xi_i$$

- KKT optimality conditions:

Zero derivatives:  $\nabla_a \mathcal{L}(a^*, b^*, \xi^*, \lambda^*, \gamma^*) = - \sum_{i=1}^n \lambda_i^* y_i = 0$

$$\nabla_b \mathcal{L}(a^*, b^*, \xi^*, \lambda^*, \gamma^*) = b^* - \sum_{i=1}^n \lambda_i^* y_i x_i = 0$$

$$\nabla_{\xi_i} \mathcal{L}(a^*, b^*, \xi^*, \lambda^*, \gamma^*) = C - \lambda_i^* - \gamma_i^* = 0$$

Primal feasibility:  $y_i(a^* + (b^*)^\top x_i) \geq 1 - \xi_i^*$   
 $\xi_i^* \geq 0$

Dual feasibility:  $\lambda_i^* \geq 0$   
 $\gamma_i^* \geq 0$

Complementary slackness:  $\lambda_i^* (1 - \xi_i^* - y_i(a^* + (b^*)^\top x_i)) = 0$   
 $\gamma_i^* \xi_i^* = 0$

207

## Soft-Margin Support Vector Machines

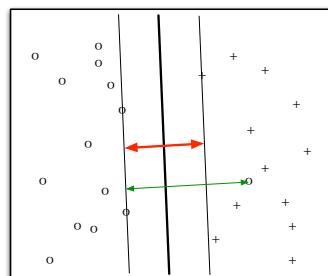
- For **non-linearly separable** datasets, we can allow for **margin violations**

$$\begin{aligned} \xi_i &= \begin{cases} 1 - y_i(a + b^\top x_i) & \text{if margin violated,} \\ 0 & \text{if not violated.} \end{cases} \\ &= \max(0, 1 - y_i(a + b^\top x_i)) \end{aligned}$$

- Penalizing violations by their magnitude,

$$\begin{aligned} \min_{a, b, \xi} \quad & \frac{1}{2} \|b\|^2 + C \sum_{i=1}^n \xi_i \\ \text{subject to} \quad & y_i(a + b^\top x_i) \geq 1 - \xi_i \\ & \xi_i \geq 0 \end{aligned}$$

where  $C$  is a tuning parameter.



206

## Soft-Margin Support Vector Machines

- Setting derivatives of primal variables to zero leads to the dual programme:

$$\begin{aligned} \max_{\lambda} \quad & \sum_{i=1}^n \lambda_i - \frac{1}{2} \sum_{i,j=1}^n \lambda_i \lambda_j y_i y_j (x_i)^\top (x_j) \\ \text{subject to} \quad & \sum_{i=1}^n \lambda_i y_i = 0 \\ & 0 \preceq \lambda \preceq C \end{aligned}$$

Only difference is the **box constraint** on  $\lambda_i \in [0, C]$ .

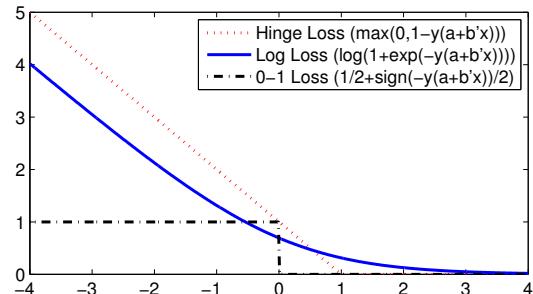
208

## Soft-Margin Support Vector Machines

- From primal programme, we can first minimize over  $\xi_i$ 's, leading to an unconstrained convex programme:

$$\min_{a,b} \frac{1}{2} \|b\|^2 + C \sum_{i=1}^n \max(0, 1 - y_i(a + b^\top x_i))$$

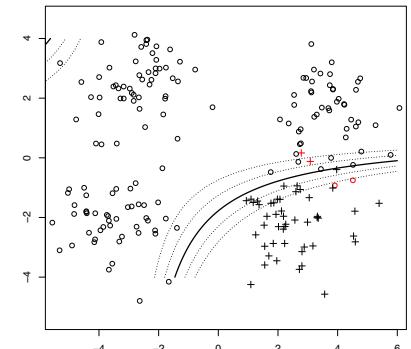
- Interpretation: regularized empirical risk minimization with the **hinge loss**.



209

## Nonlinear Methods

- Decision boundaries and regression functions often need to be nonlinear.
- One general approach: transform data  $x \mapsto \phi(x)$ .
- A **global** approach. Decisions and optimal parameters depend on **whole** training dataset.
- Alternative approach:  $p(Y = 1|X = x)$  or  $f(x)$  depends only on data cases in **local neighbourhood** of  $x$ .



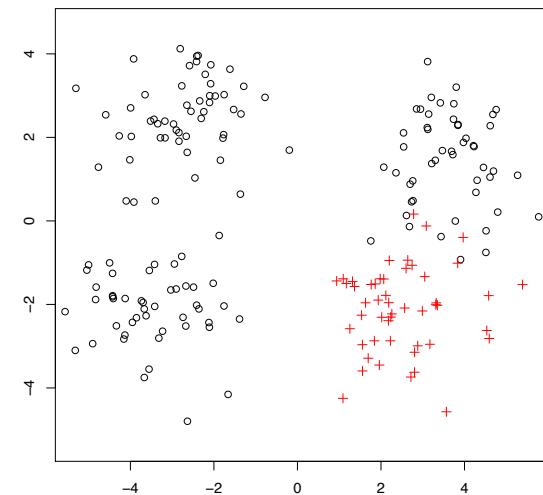
211

## Support Vector Machines – Discussion

- Multiclass classification:** If there are more than two classes, there are multiclass generalizations of the SVM.
- A simple practical idea: treat a multiclass problem as multiple binary classification problems.
  - One-vs-one: train  $K(K - 1)$  binary SVMs, for each pair of classes. At test time, predict class that got the most votes.
  - One-vs-rest: train  $K$  binary SVMs, one for each class vs all other classes. At test time, predict class with largest discriminant value  $a_k + b_k^\top x$ .
- Optimization for large scale problems:
  - Standard quadratic programme solvers not scalable.
  - Sequential minimal optimization (SMO): iterative solve pairs of  $\lambda_i$ 's.
  - Pegasos : stochastic gradient descent on regularized hinge loss objective.
- $L_2$  regularization controls overfitting.
- Not probabilistic and cannot produce uncertainty estimates.
- Statistical learning theory foundations.
- Further readings:
  - Bishop, Chapter 6.
  - Christopher Burgess, A Tutorial on Support Vector Machines for Pattern Recognition. 1998.

210

## Local Methods



212

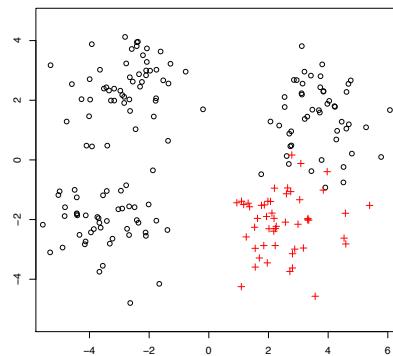
## k-Nearest Neighbours

- ▶ A simple, local, nonlinear, non-model-based, method.
- ▶ Prediction at a data vector  $x$  is simply determined by the  $k$  nearest neighbours  $ne_k(x)$  of  $x$  among the training set.
- ▶ Classification: predict the **majority vote** of the neighbours:

$$f^{\text{kNN}}(x) = \underset{l}{\operatorname{argmax}} \quad |\{j \in ne_k(x) : y_j = l\}|.$$

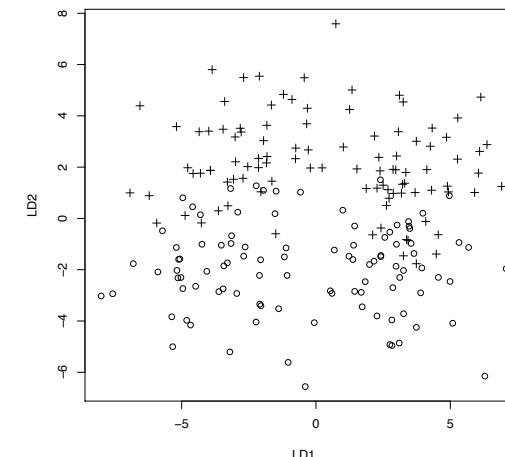
- ▶ Regression: predict the average among the neighbours:

$$f^{\text{kNN}}(x) = \frac{\sum_{j \in ne_k(x)} y_j}{\sum_{j \in ne_k(x)} 1}.$$



213

## k-Nearest Neighbour Demo



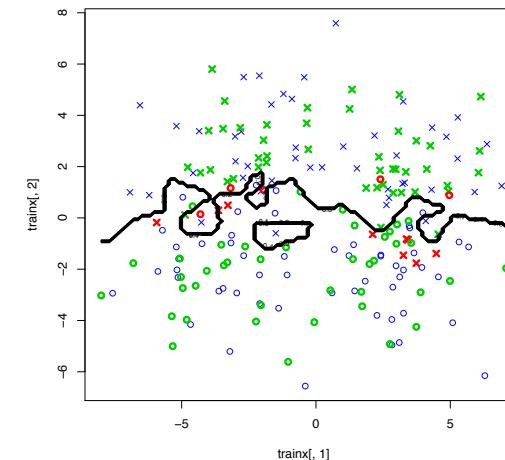
215

## k-Nearest Neighbours

- ▶ Nearest neighbours are simple and essentially model-free methods for classification.
- ▶ Weaker modelling assumptions than e.g. LDA, Naïve Bayes and logistic regression.
- ▶ These methods are not very useful for understanding relationships between attributes and class predictions.
- ▶ As **black box** classification methods however, they are often perform reasonably on real life problems and provide a good benchmark.
- ▶ Can break down in high-dimensional data:
  - ▶ Effectively, partitions input space into regions each containing  $k$  data points, and prediction in each region estimated separately.
  - ▶ In a space of dimension  $p \gg 0$ , number of regions needed is  $R = m^p$ , so size of dataset needed is  $km^p$ .

214

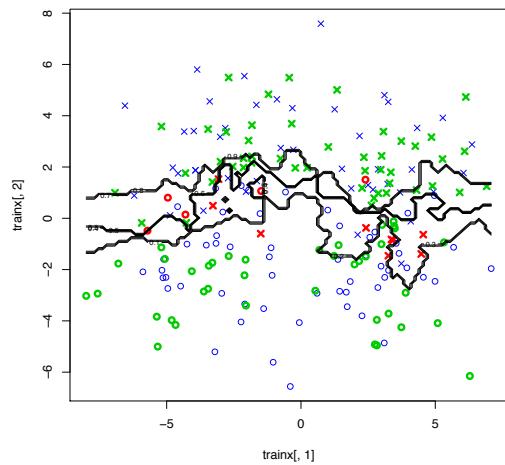
## k-Nearest Neighbour Demo



216

## Result of 1NN

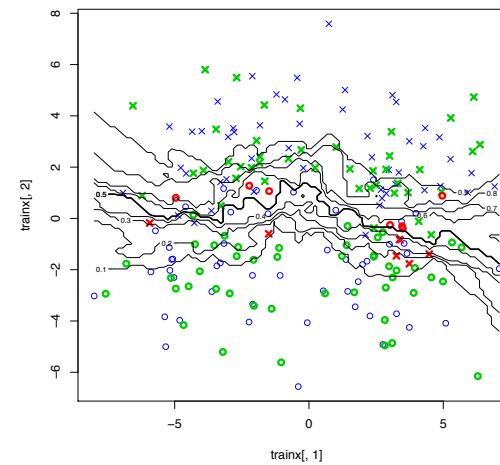
## k-Nearest Neighbour Demo



Result of 3NN

217

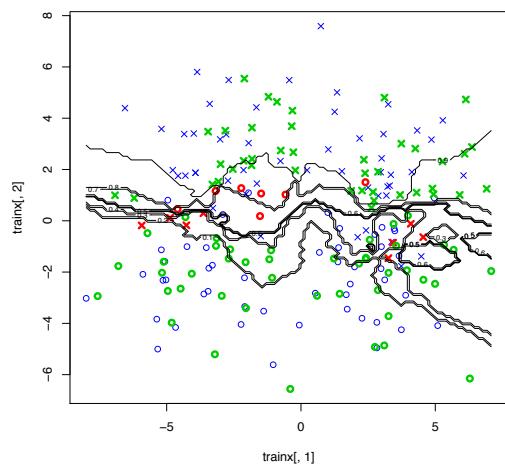
## k-Nearest Neighbour Demo



Result of 11NN

219

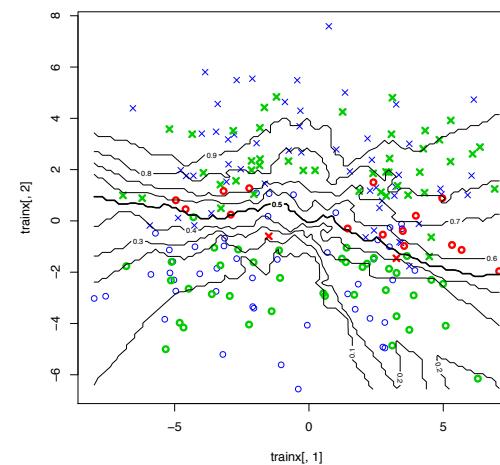
## k-Nearest Neighbour Demo



Result of 5NN

218

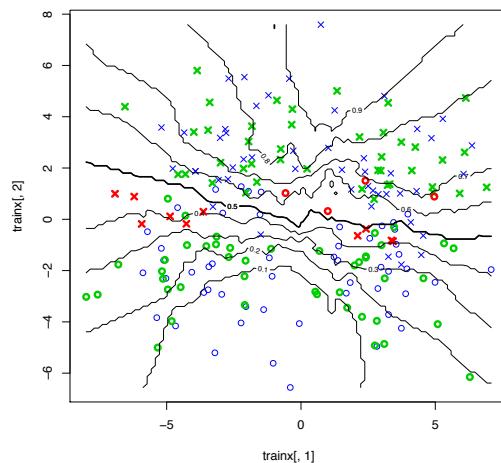
## k-Nearest Neighbour Demo



Result of 21NN

220

## k-Nearest Neighbour Demo



Result of 31NN

221

## k-Nearest Neighbour Demo – R Code I

```
library(MASS)
## load crabs data data(crabs)
ct <- as.numeric(crabs[,1])-1+2*(as.numeric(crabs[,2])-1)
## project into first two LD
cb.lda <- lda(log(crabs[,4:8]),ct)
cb.ldp <- predict(cb.lda)
x <- as.matrix(cb.ldp$x[,1:2])
y <- as.numeric(crabs[,2])-1
x <- x + rnorm(dim(x)[1]*dim(x)[2])*1.5
eqscplot(x,pch=2+y+1,col=1)

k <- 3

knn <- function(k,x,y,gridsize=100) {
  n      <- length(y)
  p      <- dim(x)[2]
  i      <- sample(rep(c(TRUE,FALSE),each=n/2),n,replace=FALSE)
  train <- (1:n)[i]
  test  <- (1:n)[i]
  trainx <- x[train,]
  trainy <- y[train]
  testx <- x[test,]
  testy <- y[test]

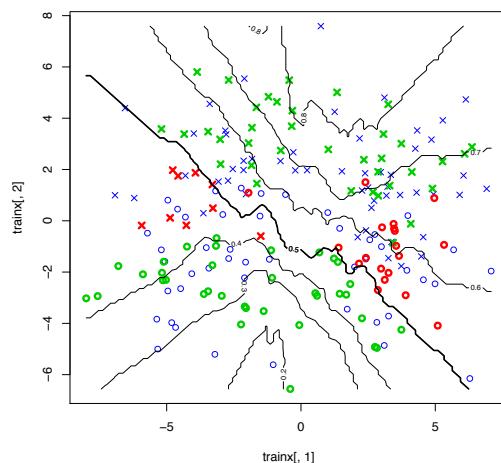
  trainn <- dim(trainx)[1]
  testn <- dim(testx)[1]

  gridx1 <- seq(min(x[,1]),max(x[,2]),length=gridsize)
  gridx2 <- seq(min(x[,2]),max(x[,2]),length=gridsize)
  gridx <- as.matrix(expand.grid(gridx1,gridx2))
  gridn <- dim(gridx)[1]

  # calculate distances, smart and intelligently.
  trainxx <- t((trainx*trainx) %*% matrix(1,p,1))
```

223

## k-Nearest Neighbour Demo



Result of 51NN

222

## k-Nearest Neighbour Demo – R Code II

```
testxx <- (testx*testx) %*% matrix(1,p,1)
gridxx <- (gridx*gridx) %*% matrix(1,p,1)
testtraindist <- matrix(1,testn,1) %*% trainxx +
  testxx %*% matrix(1,1,trainn) -
  2*(testx %*% t(trainx))
gridtraindist <- matrix(1,gridn,1) %*% trainxx +
  gridxx %*% matrix(1,1,trainn) -
  2*(gridx %*% t(trainx))

# predict
testp <- numeric(testn)
gridp <- numeric(gridn)
for (j in 1:testn) {
  nearestneighbors <- order(testtraindist[j,])[1:k]
  testp[j] <- mean(trainy[nearestneighbors])
}
for (j in 1:gridn) {
  nearestneighbors <- order(gridtraindist[j,])[1:k]
  gridp[j] <- mean(trainy[nearestneighbors])
}
predy <- as.numeric(testp>.5)

plot(trainx[,1],trainx[,2],pch=trainy+3+1,col=4,lwd=.5)
points(testx[,1],testx[,2],pch=testy+3+1,col=2+(predy==testy),lwd=3)
contour(gridx1,gridx2,matrix(gridp,gridsize,gridsize),
  levels=seq(.1,.9,1),lwd=5,add=TRUE)
contour(gridx1,gridx2,matrix(gridp,gridsize,gridsize),
  levels=c(.5),lwd=2,add=TRUE)
}
```

224

## Asymptotic Performance of 1NN

- Let  $(x_i, y_i)_{i=1}^n$  be training data where  $x_i \in \mathbb{R}^p$  and  $y_i \in \{1, 2, \dots, K\}$ .
- We define

$$\hat{y}_{\text{Bayes}}(x) = \arg \max_{l \in \{1, \dots, K\}} \pi_l f_l(x)$$

and

$$\hat{y}_{\text{1NN}}(x) = y \text{ (nearest neighbour of } x).$$

- The (optimal) Bayes risk and 1NN risk are:

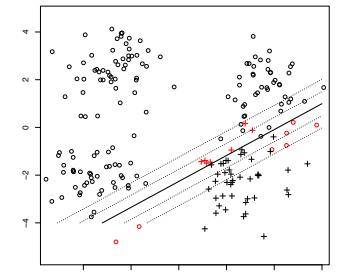
$$\begin{aligned} R_{\text{Bayes}} &= \mathbb{E}[\mathbb{I}(Y \neq \hat{y}_{\text{Bayes}}(X))] \\ R_{\text{1NN}} &= \mathbb{E}[\mathbb{I}(Y \neq \hat{Y}_{\text{1NN}}(X))] \end{aligned}$$

- As  $n \rightarrow \infty$ , we have the following powerful result

$$R_{\text{Bayes}} \leq R_{\text{1NN}} \leq 2R_{\text{Bayes}} - \frac{K}{K-1} R_{\text{Bayes}}^2.$$

## Non-linear Problems

- Linear methods (PCA, LDA, linear and logistic regression) are simple and effective techniques to learn from data “to first order”.
- To capture more intricate information from data, flexible, non-linear methods are often needed.
  - Explicit non-linear transformations  $x \mapsto \phi(x)$ .
  - Local methods like kNN.
- Kernel methods:** introduce non-linearities through **implicit** non-linear transforms, often local in nature.



225

227

## K-Nearest Neighbours – Discussion

- kNN is sensitive to distances: normalize data and find suitable metric.
- Choice of  $k$  important: controls flexibility of model.
- Computational cost of kNN is very high.
  - Need to store **all** training data.
  - Need to compare each test data vector to **all** training data.
  - Need **a lot of data** in high dimensions.
- Mitigation techniques:
  - Compute approximate nearest neighbours, using kd-trees, cover trees, random forests.
  - Apply K-means to data in each class, to reduce size of data (need to use large  $K$ ).

226

## The Kernel Method

- Back to the soft-margin SVM. The dual objective is:

$$\max_{\lambda} \sum_{i=1}^n \lambda_i - \frac{1}{2} \sum_{i,j=1}^n \lambda_i \lambda_j y_i y_j \phi(x_i)^\top \phi(x_j) \quad \text{subject to} \quad \begin{cases} \sum_{i=1}^n \lambda_i y_i = 0 \\ 0 \leq \lambda \leq C \end{cases}$$

- Suppose  $p = 2$ , and we would like to introduce quadratic non-linearities,

$$\phi(x_i) = (1, \sqrt{2}x_{i1}, \sqrt{2}x_{i2}, x_{i1}^2, x_{i2}^2, x_{i1}x_{i2})^\top$$

Then

$$\begin{aligned} \phi(x_i)^\top \phi(x_j) &= 1 + 2x_{i1}x_{j1} + 2x_{i2}x_{j2} + x_{i1}^2 x_{j1}^2 + x_{i2}^2 x_{j2}^2 + x_{i1}x_{i2}x_{j1}x_{j2} \\ &= (1 + x_i^\top x_j)^2 \end{aligned}$$

- Since only dot-products are needed in the objective function, non-linear transform need not be computed explicitly!
- Generally,  $m$ -order interactions can be implemented simply by  $\phi(x_i)^\top \phi(x_j) = (1 + x_i^\top x_j)^m$ . This is called a **polynomial kernel**.

228

## The Kernel Method

- The **Gram matrix** is the matrix of dot-products,  $B_{ij} = \phi(x_i)^\top \phi(x_j)$ .

$$B = \begin{pmatrix} -\phi(x_1)^\top & - \\ \vdots & \\ -\phi(x_i)^\top & - \\ \vdots & \\ -\phi(x_n)^\top & - \end{pmatrix} \times \begin{pmatrix} | & & | & & | \\ \phi(x_1) & \cdots & \phi(x_j) & \cdots & \phi(x_n) \\ | & & | & & | \end{pmatrix}$$

- Since  $B = \Phi\Phi^\top$ , it is symmetric and positive semidefinite.
- The Gram matrix is sufficient for training the soft-margin SVM.

$$\max_{\lambda} \quad \sum_{i=1}^n \lambda_i - \frac{1}{2} \sum_{i,j=1}^n \lambda_i \lambda_j y_i y_j B_{ij} \quad \text{subject to} \quad \begin{cases} \sum_{i=1}^n \lambda_i y_i = 0 \\ 0 \leq \lambda \leq C \end{cases}$$

229

## The Kernel Method

$$\kappa(x, x') = \langle \phi(x), \phi(x') \rangle$$

- We do not need to compute the features ever—the Gram matrix is sufficient for learning and prediction. The discriminant function (absorbing  $a$  into  $b$ ) is

$$g(x) = \sum_{i=1}^n \lambda_i^* y_i \phi(x_i)^\top \phi(x) = \sum_{i=1}^n \lambda_i^* y_i \kappa(x_i, x)$$

- The function  $\phi$  can be interpreted as non-linear **features** of our data vectors  $x \in \mathcal{X}$ .
- Generally, the Hilbert space can be **infinite-dimensional**, so we are effectively computing an infinite number of features of our data, and learning a SVM based on all features.
- There are an infinite number of parameters in the SVM—a **nonparametric** method.
- The  $L_2$  regularization of SVMs is very important to prevent overfitting.

231

## The Kernel Method

- A **kernel** is a function  $\kappa : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$  such that:
  - Symmetric:  $\kappa(x, x') = \kappa(x', x)$ .
  - Positive semidefinite: given any finite set  $\{x_i\}_{i=1}^n \subset \mathcal{X}$ , the matrix  $B \in \mathbb{R}^{n \times n}$  with entries  $B_{ij} = \kappa(x_i, x_j)$  is positive definite. Equivalently, for any  $c \in \mathbb{R}^n$ ,

$$\sum_{i=1}^n \sum_{j=1}^n c_i c_j \kappa(x_i, x_j) \geq 0$$

- **Mercer's Theorem:** if  $\kappa$  is continuous, symmetric and positive semidefinite, then there is a function  $\phi : \mathcal{X} \rightarrow \mathcal{H}$  into a Hilbert space  $\mathcal{H}$  with inner product  $\langle \cdot, \cdot \rangle$  such that

$$\kappa(x, x') = \langle \phi(x), \phi(x') \rangle$$

230

## Examples of Kernels

- Polynomial kernel:

$$\kappa(x, x') = (1 + x^\top x')^m$$

- **Gaussian, radial-basis function (RBF), or squared-exponential kernel:**

$$\kappa(x, x') = \exp\left(-\frac{1}{2} \|x - x'\|_M^2\right)$$

This leads to a discriminant function of form

$$g(x) = \sum_{i=1}^n \lambda_i^* y_i \exp\left(-\frac{1}{2} \|x_i - x\|_M^2\right)$$

A local method very similar to kNN.

- If  $\kappa_1$  and  $\kappa_2$  are both kernels, then so are kernels defined by

$$\kappa_3(x, x') = \kappa_1(x, x') + \kappa_2(x, x')$$

$$\kappa_4(x, x') = \kappa_1(x, x') \times \kappa_2(x, x')$$

232

## Kernel SVM Demo

```

library(MASS)
library(e1071)
## load crabs data, project onto LD space, add noise.
data(crabs)
ct <- as.numeric(crabs[,1])-1+2*(as.numeric(crabs[,2])-1)
cb.lda <- lda(log(crabs[,4:8]),ct)
cb.ldp <- predict(cb.lda)
x <- as.matrix(cb.ldp$x[,1:2])
y <- as.numeric(crabs[,2])-1
x <- x + rnorm(dim(x)[1]*dim(x)[2])*1.5
gridsize <- 100
xlim <- c(min(x[,1]),max(x[,1]))
ylim <- c(min(x[,2]),max(x[,2]))
gridx1 <- seq(xlim[1],xlim[2],length=gridsize)
gridx2 <- seq(ylim[1],ylim[2],length=gridsize)
gridx <- as.matrix(expand.grid(gridx1,gridx2))
gridn <- dim(gridx)[1]
plot(x,pch=2*y+1,col=1,xlim=xlim,ylim=ylim)

n <- length(y)
p <- dim(x)[2]
i <- sample(rep(c(TRUE,FALSE),each=n/2),n,replace=FALSE)
train <- (1:n)[i]
test <- (1:n)[!i]
trainx <- x[train,]
trainy <- y[train]
testx <- x[test,]
testy <- y[test]

svmdemo <- function(kernel,gamma=1,coef0=0,cost=1,degree=3) {
  model <- svm(trainx,trainy,kernel=kernel,gamma=gamma,coef0=coef0,degree=degree,cost=cost)
  gridp <- predict(model,gridx)
  predy <- as.numeric(predict(model,testx)>.5)

  plot(trainx[,1],trainx[,2],pch=trainy+3+1,col=4,lwd=.5,xlim=xlim,ylim=ylim)
  points(testx[,1],testx[,2],pch=testy+3+1,col=2+(predy==testy),lwd=3)
  contour(gridx1,gridx2, matrix(gridp,gridsize,gridsize),levels=seq(.1,.9,.1),lwd=.5,add=TRUE)
  contour(gridx1,gridx2, matrix(gridp,gridsize,gridsize),levels=c(.5),lwd=2,add=TRUE)
}

svmdemo()

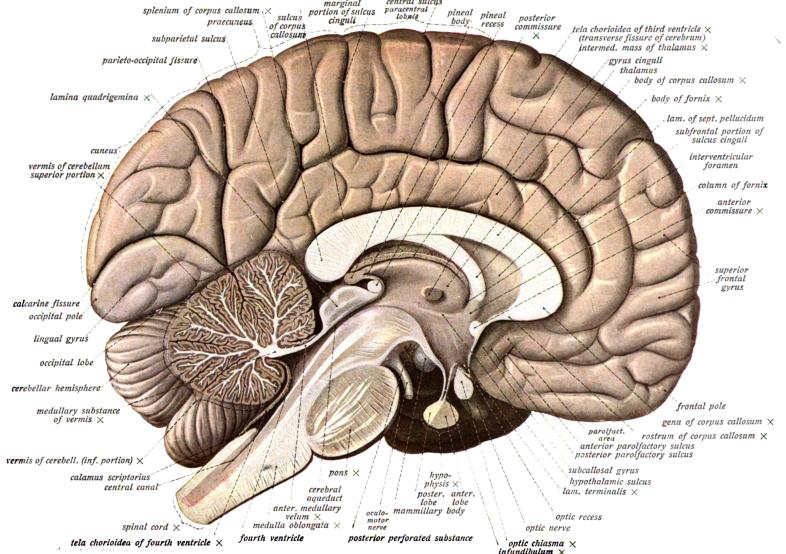
```

233

## Kernel Methods – Discussion

- ▶ The kernel method allows for very flexible and powerful machine learning models.
- ▶ Kernels can be defined over much more complex structures than vectors, e.g. graphs, strings.
- ▶ Can be hard to interpret.
- ▶  $O(n^3)$  computation and  $O(n^2)$  memory cost can be prohibitive.
- ▶ Further readings:
  - ▶ Bishop, Chapter 6.
  - ▶ Christopher Burgess, A Tutorial on Support Vector Machines for Pattern Recognition. 1998.
  - ▶ Rasmussen and Williams, Gaussian Processes for Machine Learning. 2006.

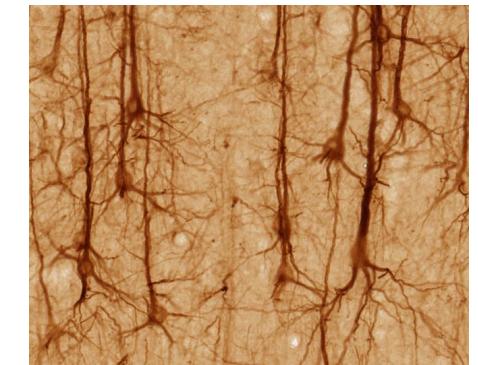
## The Brain



235

## The Brain

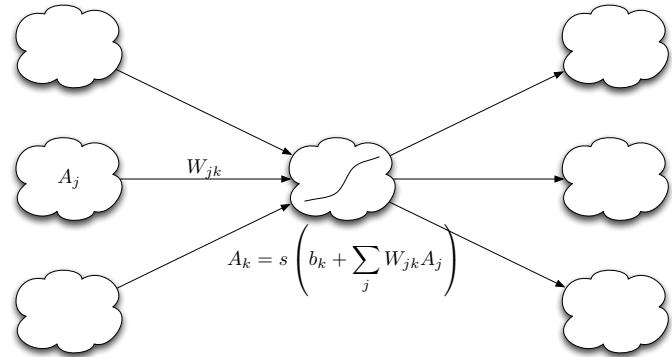
- ▶ Basic computational elements: neurons.
- ▶ Receives signals from other neurons via dendrites.
- ▶ Sends processed signals via axons.
- ▶ Axon-dendrite interactions at synapses.
- ▶  $10^{10} - 10^{11}$  neurons.
- ▶  $10^{14} - 10^{15}$  synapses.
- ▶ Connectionist architecture: the network and its structure govern the computations performed.



234

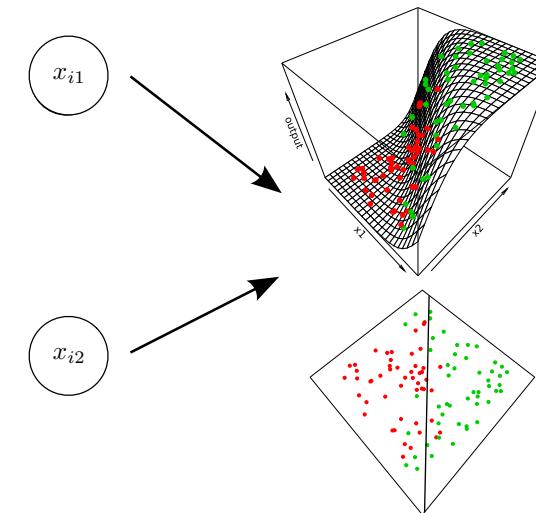
236

## A Simple Model of Neural Computations



237

## A Simple Model of Neural Computations



239

## Modelling Conditional Probabilities

- ▶ Data vectors  $x_i \in \mathbb{R}^p$ , binary labels  $y_i \in \{0, 1\}$ .
- ▶ **Inputs**  $x_{i1}, \dots, x_{ip}$
- ▶ **output**  $\hat{y}_i = p(Y = 1 | X = x_i)$
- ▶ **hidden unit activations**  $h_{i1}, \dots, h_{im}$ 
  - ▶ Compute **hidden unit activations**:

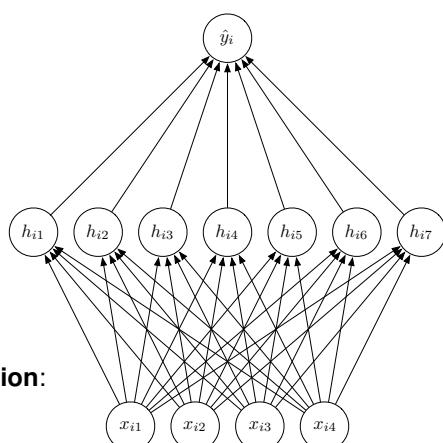
$$h_{ik} = s\left(b_k^h + \sum_{j=1}^p W_{jk}^h x_{ij}\right)$$

- ▶ Compute **output probability**:

$$\hat{y}_i = s\left(b^o + \sum_{k=1}^m W_k^o h_{ik}\right)$$

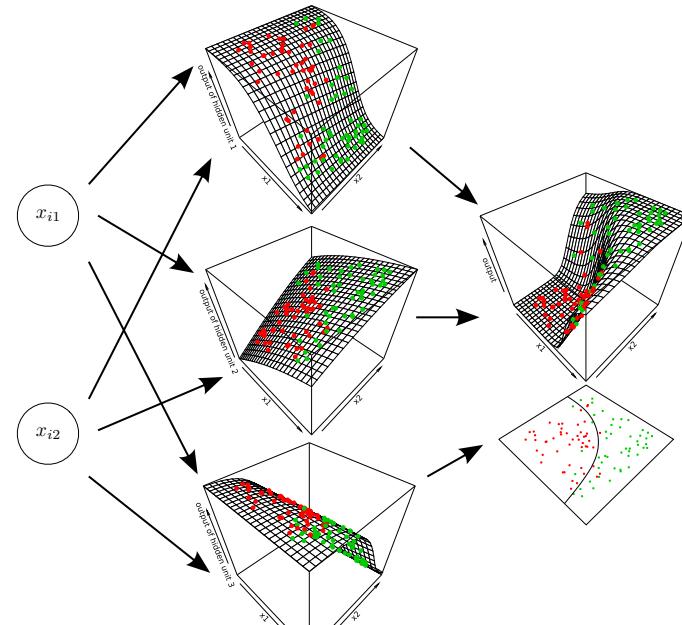
- ▶ Common nonlinear **activation function**:  
the logistic function

$$s(z) = \frac{1}{1 + \exp(-z)}$$



238

## A Simple Model of Neural Computations



240

## Training a Neural Network

- ▶ Objective function:  $L_2$ -regularized log loss

$$J = - \sum_{i=1}^n y_i \log \hat{y}_i + (1 - y_i) \log(1 - \hat{y}_i) + \frac{1}{2} \sum_{jk} C |W_{jk}^h|^2 + \frac{1}{2} \sum_k C |W_k^o|^2$$

where

$$\hat{y}_i = s \left( b^o + \sum_{k=1}^m W_k^o h_{ik} \right) \quad h_{ik} = s \left( b_k^h + \sum_{j=1}^p W_{jk}^h x_{ij} \right)$$

- ▶ Optimize parameters  $\{b_k^h, b^o, W_{jk}^h, W_k^o\}$  by gradient descent.

$$\frac{dJ}{dW_k^o} = CW_k^o + \sum_{i=1}^n \frac{dJ}{d\hat{y}_i} \frac{d\hat{y}_i}{dW_k^o} = CW_k^o + \sum_{i=1}^n (\hat{y}_i - y_i) h_{ik}$$

$$\frac{dJ}{dW_{jk}^h} = CW_{jk}^h + \sum_{i=1}^n \frac{dJ}{d\hat{y}_i} \frac{d\hat{y}_i}{dh_{ik}} \frac{dh_{ik}}{dW_{jk}^h} = CW_{jk}^h + \sum_{i=1}^n (\hat{y}_i - y_i) W_k^o h_{ik} (1 - h_{ik}) x_{ij}$$

- ▶ **Backpropagation:** gradients computed via chain rule, and propagated through the network backwards.
- ▶  $L_2$  regularization often called **weight decay**.

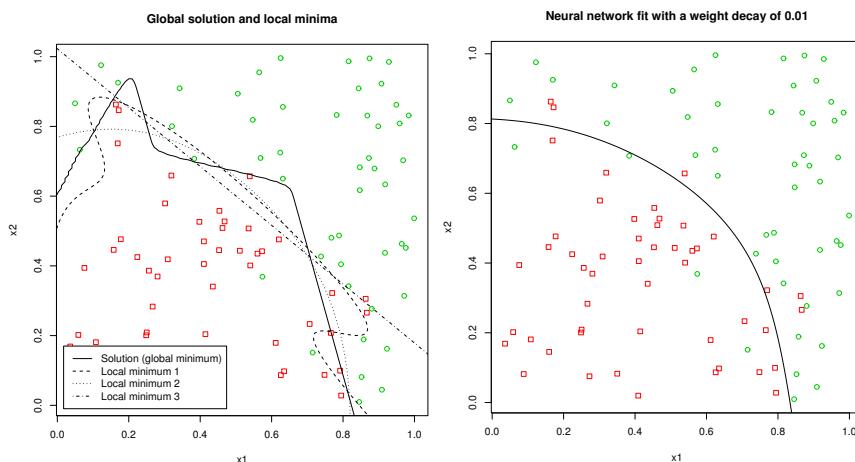
241

## Neural Networks – Discussion

- ▶ Nonlinear hidden units introduce modelling flexibility.
- ▶ As opposed to user introduced nonlinearities, kernel methods, kNNs, features are global, and learnt to maximize predictive performance.
- ▶ Neural networks with a single hidden layer can model arbitrarily complex functions (with enough hidden units).
- ▶ Highly flexible framework, with many variations to solve different learning problems and introduce domain knowledge.
- ▶ Optimization problem is not convex, and objective function can have many local optima, plateaus and ridges.
- ▶ On large scale problems, often use stochastic gradient descent, along with a whole host of techniques for optimization, regularization, and initialization.
- ▶ Strengths of neural networks:
  - ▶ Flexibility and generalization ability.
  - ▶ Computational efficiency, parallelizability.
- ▶ Recent developments, especially by Geoffrey Hinton, Yann LeCun, Yoshua Bengio, Andrew Ng and others. See also <http://deeplearning.net/>.

243

## Neural Networks



R package implementing neural networks with a single hidden layer: `nnet`.

242

## Neural Networks – Variations

- ▶ Other loss functions can be used, e.g. for regression:

$$\sum_{i=1}^n |y_i - \hat{y}_i|^2$$

For multiclass classification, use **softmax** outputs:

$$\hat{y}_{ik} = \frac{\exp(b_k^o + \sum_\ell W_{ik}^o h_{\ell k})}{\sum_{k'} \exp(b_{k'}^o + \sum_\ell W_{ik'}^o h_{\ell k'})} \quad L(y_i, \hat{y}_i) = \sum_{k=1}^K \mathbb{1}(y_i = k) \log \hat{y}_{ik}$$

- ▶ Other activation functions can be used, e.g. a recent popular one is called **rectified linear activation**:

$$s(z) = \log(1 + \exp(z))$$

- ▶ Multiple layers of hidden units can be used, called **multilayer perceptrons** (MLP) or **deep networks**.

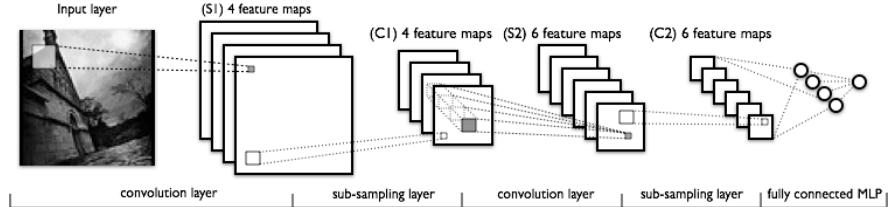
244

## Visual Object Recognition

The screenshot shows the IMAGENET search interface. At the top, there's a search bar with 'SEARCH' button, and links for 'Home', 'About', and 'Explore'. Below the search bar, it says '14,197,122 images, 21841 synsets indexed'. A message 'Not logged in. Login | Signup' is at the top right. The main content area is titled 'Cow' with a sub-description 'Mature female of mammals of which the male is called `bull''. It shows '1588 pictures' and '82.99% Popularity Percentile'. There are tabs for 'Treemap Visualization' (selected), 'Images of the Synset' (disabled), and 'Downloads'. Below these are two rows of small images of cows. To the left is a sidebar with a tree map visualization of the 'ImageNet 2011 Fall Release' structure, under the 'Cow' synset, showing various categories like 'plant, flora, plant life' (3232), 'geological formation, formation' (4486), etc.

245

## Deep Convolutional Neural Networks

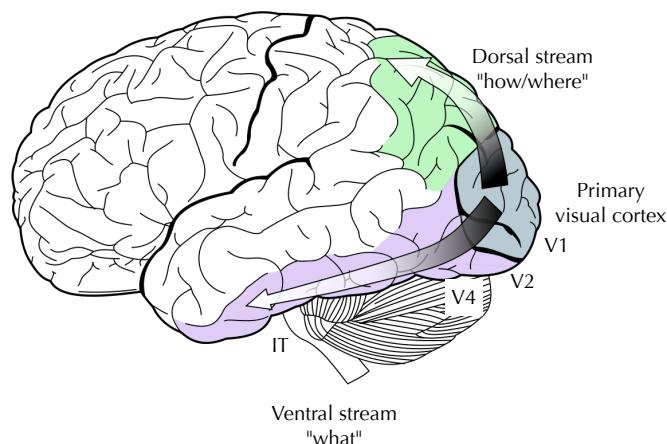


- ▶ Input is a 2D image,  $X \in \mathbb{R}^{p \times q}$ .
  - ▶ Convolution: detects simple object parts or features
- $$A^m = s(X * W^m)$$
- $$A_{jk}^m = s \left( b^m + \sum_{fg} X_{j-f, k-g} W_{fg}^m \right)$$
- ▶ Sub-sampling: incorporates local translation invariance by max-pooling
- $$B_{jk}^m = \max\{A_{fg}^m : |f - j| \leq w, |g - k| \leq h\}$$
- ▶ Learn features/parts of increasing complexity over multiple layers.

LeCun et al, Krizhevsky et al

247

## Visual Processing in the Brain



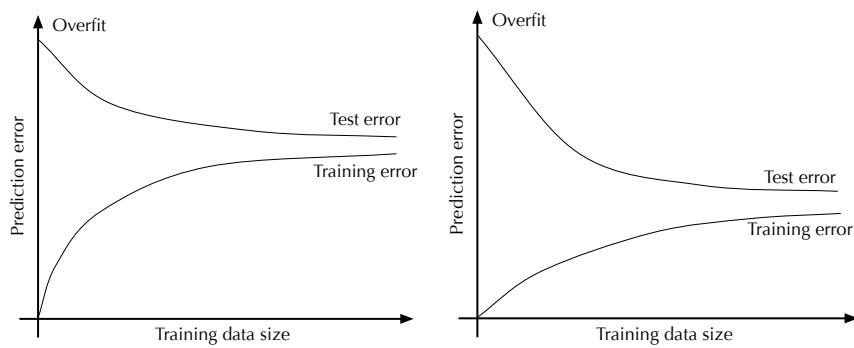
246

## Revisiting Learning Generalization

- ▶ Generalization ability is a central concept in machine learning.
- ▶ Splitting data into training and test sets allows us to estimate how well our methods are generalizing.
- ▶ Two important factors determining generalization ability:
  - ▶ Model complexity
  - ▶ Training data size
- ▶ To control overfitting, we need to regularize learning.
- ▶ Can we learn the tuning parameters as well?

248

## Learning Curves



Fixed model complexity, varying dataset size. Two models, of different complexities. Which is which?

249

## Bias-Variance Tradeoff

- ▶ A different perspective on generalization ability.
- ▶ Suppose we are in a regression setting, with

$$Y = f^*(X) + \mathcal{N}(0, \sigma^2)$$

- ▶ Given a dataset  $D = (x_i, y_i)_{i=1}^n$ , train a model  $f(x; \theta)$ .
- ▶ Estimated  $\hat{\theta}$  is a function of data set  $D$ .
- ▶ How will we do, averaging over data sets of size  $n$ ?

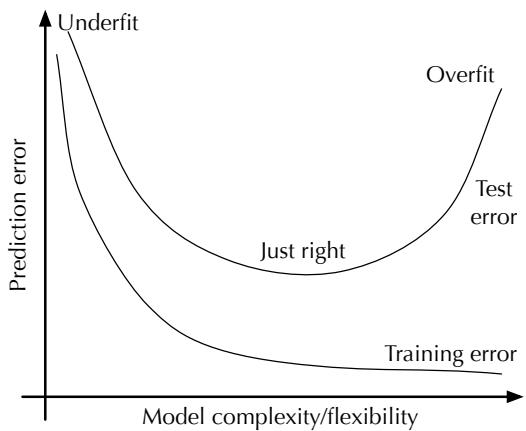
$$\begin{aligned} & \mathbb{E}_D[(Y - f(X; \hat{\theta}(D)))^2] \\ &= \underbrace{(\bar{f}(X) - f^*(X))^2}_{\text{bias}^2} + \underbrace{\mathbb{E}_D[(\bar{f}(X) - f(X; \hat{\theta}(D)))^2]}_{\text{variance}} + \underbrace{(Y - f^*(X))^2}_{\text{noise}} \end{aligned}$$

where  $\bar{f}(X) = \mathbb{E}_D[f(X; \hat{\theta}(D))]$  is average prediction (over data sets).

- ▶ **Noise:** intrinsic difficulty of regression problem.
- ▶ **Variance:** How variable is our method if given different datasets?
- ▶ **Bias:** How far is our average prediction away from the truth?

251

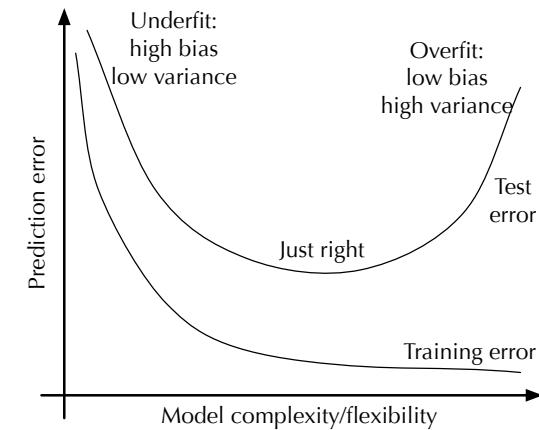
## Learning Curves



Fixed dataset size, varying model complexity.

250

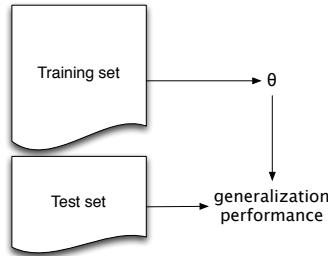
## Learning Curve



252

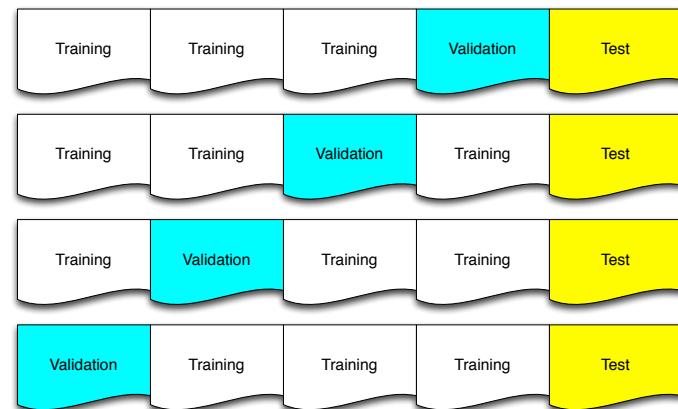
## Optimizing Hyperparameters and Model Complexity

- ▶ How can we optimize generalization ability, via optimizing choice of tuning parameters, model size, and learning parameters?
- ▶ Suppose we have split data into training/test set.
- ▶ Test set can be used to determine generalization ability, and used to choose best setting of tuning parameters/model size/learning parameters with best generalization.
- ▶ Once these meta-parameters are chosen, still important to determine generalization ability, but cannot use performance on test set to gauge this anymore!
- ▶ Idea: split data into 3 sets: training set, test set, and **validation set**.



253

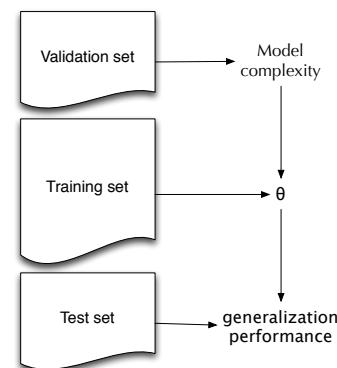
## Cross-Validation



255

## Validation Set

- ▶ For each combination of meta-parameters  $\gamma$ :
  - ▶ Train our model, obtaining model parameter  $\theta(\gamma)$ .
  - ▶ Evaluate  $\theta(\gamma)$  on validation set.
- ▶ Pick  $\gamma^*$  with best performance on validation set.
- ▶ Using  $\gamma^*$ , train on both training and validation set (**fold** the validation set into the training set), to obtain optimal  $\theta^*$ .
- ▶ Evaluate model with  $\gamma^*, \theta^*$  on test set, reporting generalization performance.
- ▶ Problem: if we have insufficient data, very wasteful to split into 3 subsets, and estimated generalization performance on validation set may be too noisy to effectively choose meta-parameters.
- ▶ Solution: **cross-validation**.



254

## Cross-Validation

- ▶ Basic approach:
  - ▶ Split training set into  $V$  folds.
  - ▶ For each  $\gamma$  and each  $v = 1, \dots, V$ :
    - ▶ Use fold  $v$  as validation set and the rest to train the model parameters  $\hat{\theta}_v$ .
- ▶ Choose  $\gamma^*$  which minimizes
 
$$R_v^{\text{emp}}(\gamma) = \frac{1}{|\text{Fold}(v)|} \sum_{i \in \text{Fold}(v)} L(y_i, \hat{Y}(x_i; \hat{\theta}_v))$$

$$\frac{1}{V} \sum_{v=1}^V R_v^{\text{emp}}(\gamma)$$
- ▶ Train model with meta-parameter  $\gamma^*$  on all training set.
- ▶ Report generalization performance on test set.
- ▶ Extreme case: **Leave-one-out (LOO)** cross validation: one data item per fold.
- ▶ Cross-validation can be computationally very expensive.

256

## Back to Maximum Likelihood

- Given a generative model

$$f(x, y = k) = \pi_k f_k(x)$$

- Using a generative modelling approach, we assume a parametric form for  $f_k(x) = f(x; \phi_k)$  and compute the MLE  $\hat{\theta}$  of  $\theta = (\pi_k, \phi_k)_{k=1}^K$  based on the training data  $\{x_i, y_i\}_{i=1}^n$ .
- We then use a plug-in approach to perform classification

$$p(Y = k | X = x, \hat{\theta}) = \frac{\hat{\pi}_k f(x; \hat{\phi}_k)}{\sum_{j=1}^K \hat{\pi}_j f(x; \hat{\phi}_j)}$$

- Even for simple models, this can prove difficult; e.g. for LDA,  $f(x; \phi_k) = \mathcal{N}(x; \mu_k, \Sigma)$ , and the MLE estimate of  $\Sigma$  is not full rank for  $p > n$ .
- One answer: simplify even further, e.g. using axis-aligned covariances, but this is usually too crude.
- Another answer: regularization.

257

## The Bayesian Learning Framework

- Bayes Theorem:** Given two random variables  $X$  and  $\Theta$ ,

$$p(\Theta|X) = \frac{p(X|\Theta)p(\Theta)}{p(X)}$$

- Likelihood:**  $p(X|\Theta)$
- Prior:**  $p(\Theta)$
- Posterior:**  $p(\Theta|X)$
- Marginal likelihood:**  $p(X) = \int p(X|\Theta)p(\Theta)d\Theta$
- Treat parameters as random variables, and process of learning is just computation of posterior  $p(\Theta|X)$ .
- Summarizing the posterior:
  - Posterior mode:**  $\hat{\theta}^{\text{MAP}} = \operatorname{argmax}_{\theta} p(\theta|X)$ . Maximum a posteriori.
  - Posterior mean:**  $\hat{\theta}^{\text{mean}} = \mathbb{E}[\Theta|X]$ .
  - Posterior variance:**  $\text{Var}[\Theta|X]$ .
- How to make decisions and predictions? Decision theory.
- How to compute posterior?

259

## Naïve Bayes

- Return to the spam classification example with two-class naïve Bayes

$$f(x_i; \phi_k) = \prod_{j=1}^p \phi_{kj}^{x_{ij}} (1 - \phi_{kj})^{1-x_{ij}}.$$

The MLE estimates are given by

$$\hat{\phi}_{kj} = \frac{\sum_{i=1}^n \mathbb{1}(x_{ij} = 1, y_i = k)}{n_k}, \quad \hat{\pi}_k = \frac{n_k}{n}$$

where  $n_k = \sum_{i=1}^n \mathbb{1}(y_i = k)$ .

- If a word  $j$  does not appear in class  $k$  by chance, but it does appear in a document  $x_*$ , then  $p(x_*|y_* = k) = 0$  and so posterior  $p(y_* = k|x_*) = 0$ .
- Worse things can happen: e.g., probability of document under all classes can be 0, so posterior is ill-defined.

258

## Simple Example: Coin Tosses

- A very simple example: We have a coin with probability  $\phi$  of coming up heads. Model coin tosses as iid Bernoullis, 1 =head, 0 =tail.
- Learn about  $\phi$  given dataset  $D = (x_i)_{i=1}^n$  of tosses.

$$f(D|\phi) = \phi^{n_1} (1 - \phi)^{n_0}$$

with  $n_j = \sum_{i=1}^n \mathbb{1}(x_i = j)$ .

- Maximum likelihood

$$\hat{\phi}^{\text{ML}} = \frac{n_1}{n}$$

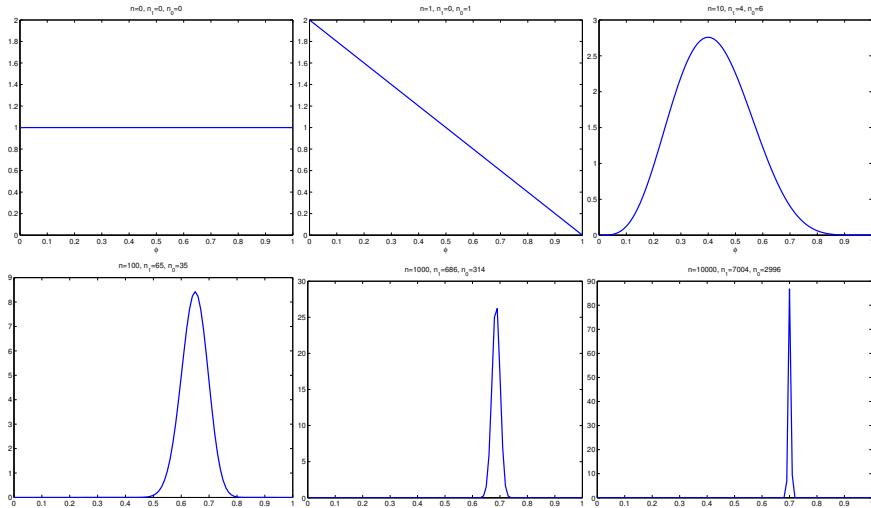
- Bayesian approach: treat unknown parameter as a random variable  $\Phi$ . Simple prior:  $\Phi \sim U[0, 1]$ . Posterior distribution:

$$p(\phi|D) = \frac{1}{Z} \phi^{n_1} (1 - \phi)^{n_0}, \quad Z = \int_0^1 \phi^{n_1} (1 - \phi)^{n_0} d\phi = \frac{(n+1)!}{n_1! n_0!}$$

Posterior is a  $\text{Beta}(n_1 + 1, n_0 + 1)$  distribution.

260

## Simple Example: Coin Tosses



Posterior becomes peaked at true value  $\phi^* = .7$  as dataset grows.

261

## Simple Example: Coin Tosses

- ▶ What about test data?
- ▶ The **posterior predictive distribution** is the conditional distribution of  $x_{n+1}$  given  $(x_i)_{i=1}^n$ :

$$\begin{aligned} p(x_{n+1}|(x_i)_{i=1}^n) &= \int_0^1 p(x_{n+1}|\phi, (x_i)_{i=1}^n)p(\phi|(x_i)_{i=1}^n)d\phi \\ &= \int_0^1 p(x_{n+1}|\phi)p(\phi|(x_i)_{i=1}^n)d\phi \\ &= (\hat{\phi}^{\text{mean}})^{x_{n+1}}(1 - \hat{\phi}^{\text{mean}})^{1-x_{n+1}} \end{aligned}$$

- ▶ We predict on new data by **averaging** the predictive distribution over the posterior. Accounts for uncertainty about  $\phi$ .

263

## Simple Example: Coin Tosses

- ▶ Posterior distribution captures all learnt information.

- ▶ Posterior mode:

$$\hat{\phi}^{\text{MAP}} = \frac{n_1}{n}$$

- ▶ Posterior mean:

$$\hat{\phi}^{\text{mean}} = \frac{n_1 + 1}{n + 2}$$

- ▶ Posterior variance:

$$\frac{1}{n+3} \hat{\phi}^{\text{mean}} (1 - \hat{\phi}^{\text{mean}})$$

- ▶ Asymptotically, for large  $n$ , variance decreases as  $1/n$  and is given by the inverse of Fisher's information.

- ▶ Posterior distribution converges to true parameter  $\phi^*$  as  $n \rightarrow \infty$ .

262

## Simple Example: Coin Tosses

- ▶ Posterior distribution is a known analytic form. In fact posterior distribution is in the same beta family as the prior.
- ▶ An example of a **conjugate prior**.
- ▶ A beta distribution  $\text{Beta}(a, b)$  with parameters  $a, b > 0$  is an exponential family distribution with density

$$p(\phi|a, b) = \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \phi^{a-1} (1-\phi)^{b-1}$$

where  $\Gamma(t) = \int_0^\infty u^{t-1} e^{-u} du$  is the gamma function.

- ▶ If the prior is  $\phi \sim \text{Beta}(a, b)$ , then the posterior distribution is

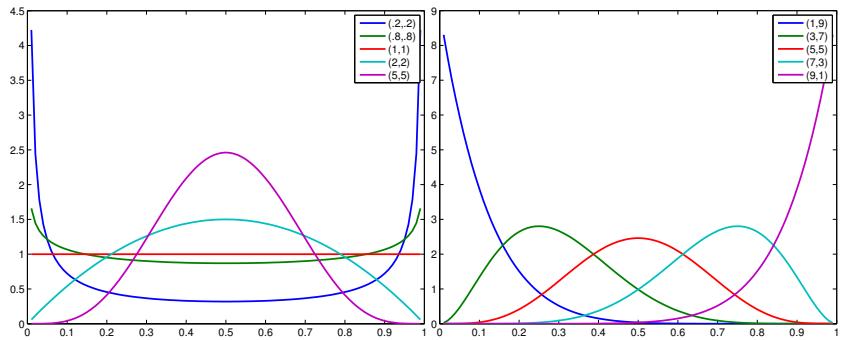
$$p(\phi|D, a, b) = \propto \phi^{a+n_1-1} (1-\phi)^{b+n_0-1}$$

so is  $\text{Beta}(a+n_1, b+n_0)$ .

- ▶ Hyperparameters  $a$  and  $b$  are **pseudo-counts**, an imaginary initial sample that reflects our prior beliefs about  $\phi$ .

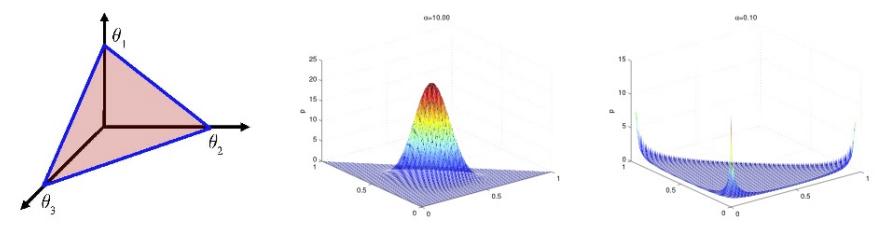
264

## Beta Distributions



265

## Dirichlet Distributions



267

## Bayesian Inference for Multinomials

- Suppose  $x_i \in \{1, \dots, K\}$  instead, and we model  $(x_i)_{i=1}^n$  as iid multinomials:

$$p(D|\pi) = \prod_{i=1}^n \pi_{x_i} = \prod_{k=1}^K \pi_k^{n_k}$$

with  $n_k = \sum_{i=1}^n \mathbb{1}(x_i = k)$  and  $\pi_k > 0$ ,  $\sum_{k=1}^K \pi_k = 1$ .

- The conjugate prior is the Dirichlet distribution.  $\text{Dir}(\alpha_1, \dots, \alpha_K)$  has parameters  $\alpha_k > 0$ , and density

$$p(\pi) = \frac{\Gamma(\sum_{k=1}^K \alpha_k)}{\prod_{k=1}^K \Gamma(\alpha_k)} \prod_{k=1}^K \pi_k^{\alpha_k - 1}$$

on the probability simplex  $\{\pi : \pi_k > 0, \sum_{k=1}^K \pi_k = 1\}$ .

- The posterior is also Dirichlet, with parameters  $(\alpha_k + n_k)_{k=1}^K$ .
- Posterior mean is

$$\hat{\pi}_k^{\text{mean}} = \frac{\alpha_k + n_k}{\sum_{j=1}^K \alpha_j + n_j}$$

266

## Text Classification with (Less) Naïve Bayes

- Under the Naïve Bayes model, the joint distribution of labels  $y_i \in \{1, \dots, K\}$  and data vectors  $x_i \in \{0, 1\}^p$  is

$$\begin{aligned} \prod_{i=1}^n p(x_i, y_i) &= \prod_{i=1}^n \prod_{k=1}^K \left( \pi_k \prod_{j=1}^p \phi_{kj}^{x_{ij}} (1 - \phi_{kj})^{1-x_{ij}} \right)^{\mathbb{1}(y_i=k)} \\ &= \prod_{k=1}^K \pi_k^{n_k} \prod_{j=1}^p \phi_{kj}^{n_{kj}} (1 - \phi_{kj})^{n_k - n_{kj}} \end{aligned}$$

where  $n_k = \sum_{i=1}^n \mathbb{1}(y_i = k)$ ,  $n_{kj} = \sum_{i=1}^n \mathbb{1}(y_i = k, x_{ij} = 1)$ .

- For conjugate prior, we can use  $\text{Dir}((\alpha_k)_{k=1}^K)$  for  $\pi$ , and  $\text{Beta}(a, b)$  for  $\phi_{kj}$  independently.
- Because the likelihood factorizes, the posterior distribution over  $\pi$  and  $(\phi_{kj})$  also factorizes, and posterior for  $\pi$  is  $\text{Dir}((\alpha_k + n_k)_{k=1}^K)$ , and for  $\phi_{kj}$  is  $\text{Beta}(a + n_{kj}, b + n_k - n_{kj})$ .

268

## Text Classification with (Less) Naïve Bayes

- ▶ For prediction give  $D = (x_i, y_i)_{i=1}^n$  we can calculate

$$p(x_0, y_0 = k|D) = p(y_0 = k|D)p(x_0|y_0 = k, D)$$

with

$$p(y_0 = k|D) = \frac{\alpha_k + n_k}{n + \sum_{l=1}^K \alpha_l}$$

$$p(x_{0j} = 1|y_0 = k, D) = \frac{a + n_{kj}}{a + b + n_k}$$

- ▶ Predicted class is

$$p(y_0 = k|x_0|D) = \frac{p(y_0 = k|D)p(x_0|y_0 = k, D)}{p(x_0|D)}$$

- ▶ Compared to ML plug-in estimator, pseudocounts help to regularize probabilities away from extreme values.

269

## Bayesian Learning – Discussion

- ▶ Clear separation between models, which frame learning problems and encapsulates prior information, and algorithms, which computes posteriors and predictions.
- ▶ Use probability theory to reason about uncertainties (both about latent variables and parameters).
- ▶ Bayesian computations — Most posteriors are intractable, and algorithms needed to efficiently approximate posterior:
  - ▶ Monte Carlo methods (Markov chain and sequential varieties).
  - ▶ Variational methods (variational Bayes, belief propagation etc).
- ▶ No optimization — no overfitting (!) but there can still be model misfit.
- ▶ Tuning parameters  $\Psi$  can be optimized (without need for cross-validation).

$$p(X|\Psi) = \int p(X|\theta)p(\theta|\Psi)d\theta$$

$$p(\Psi|X) = \frac{p(X|\Psi)p(\Psi)}{p(X)}$$

- ▶ Be Bayesian about  $\Psi$  — compute posterior.
- ▶ Type II maximum likelihood — find  $\Psi$  maximizing  $p(X|\Psi)$ .

271

## Bayesian Learning and Regularization

- ▶ Consider a Bayesian approach to logistic regression: introduce a multivariate normal prior for  $b$ , and uniform (improper) prior for  $a$ . The prior density is:

$$p(a, b) = (2\pi\sigma^2)^{-\frac{p}{2}} e^{-\frac{1}{2\sigma^2}\|b\|_2^2}$$

- ▶ The posterior is

$$p(a, b|D) \propto \exp \left( -\frac{1}{2\sigma^2} \|b\|_2^2 - \sum_{i=1}^n \log(1 + \exp(-y_i(a + b^\top x_i))) \right)$$

- ▶ The posterior mode is the parameters maximizing the above, equivalent to minimizing the  $L_2$ -regularized empirical risk.
- ▶ Regularized empirical risk minimization is (often) equivalent to having a prior and finding the maximum a posteriori (MAP) parameters.
  - ▶  $L_2$  regularization - multivariate normal prior.
  - ▶  $L_1$  regularization - multivariate Laplace prior.
- ▶ From a Bayesian perspective, the MAP parameters are just one way to summarize the posterior distribution.

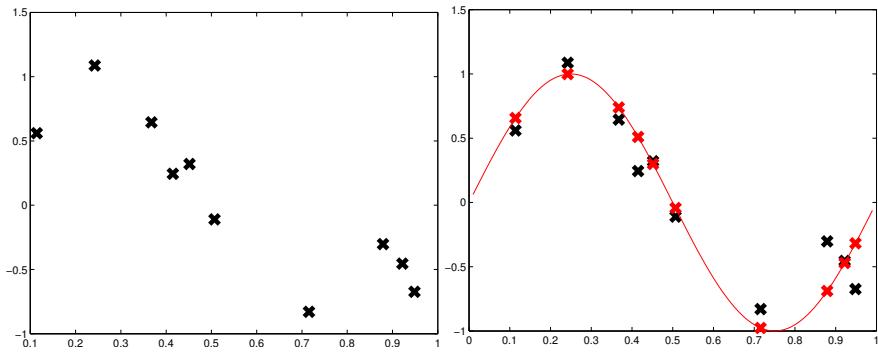
270

## Bayesian Learning – Further Readings

- ▶ Zoubin Ghahramani. Bayesian Learning. Graphical models. Videolectures.
- ▶ Gelman et al. Bayesian Data Analysis.
- ▶ Kevin Murphy. Machine Learning: a Probabilistic Perspective.
- ▶ E. T. Jaynes. Probability Theory: The Logic of Science.

272

## Gaussian Processes



- ▶ Suppose we are given a dataset consisting of  $n$  inputs  $\mathbf{x} = (x_i)_{i=1}^n$  and  $n$  outputs  $\mathbf{y} = (y_i)_{i=1}^n$ .
- ▶ Regression: learn the underlying function  $f(x)$ .

273

## Gaussian Processes

- ▶ The prior  $p(\mathbf{f})$  encodes our prior knowledge about the function.
- ▶ What properties of the function can we incorporate?

- ▶ Multivariate normal assumption:

$$\mathbf{f} \sim \mathcal{N}(0, G)$$

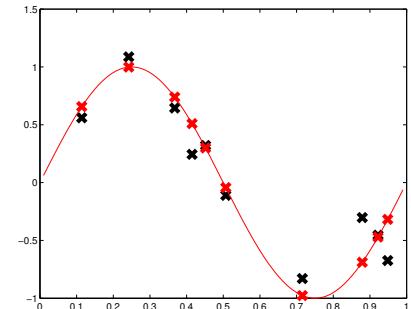
- ▶ Use a kernel function  $\kappa$  to define  $G$ :

$$G_{ij} = \kappa(x_i, x_j)$$

- ▶ Expect regression functions to be smooth: If  $x$  and  $x'$  are close by, then  $f(x)$  and  $f(x')$  have similar values, i.e. strongly correlated.

$$\begin{pmatrix} f(x) \\ f(x') \end{pmatrix} \sim \mathcal{N} \left( \begin{pmatrix} 0 \\ 0 \end{pmatrix}, \begin{pmatrix} \kappa(x, x) & \kappa(x, x') \\ \kappa(x', x) & \kappa(x', x') \end{pmatrix} \right)$$

In particular, want  
 $\kappa(x, x') \approx \kappa(x, x) = \kappa(x', x')$ .



- ▶ Model:

$$\mathbf{f} \sim \mathcal{N}(0, G)$$

$$y_i | f_i \sim \mathcal{N}(f_i, \sigma^2)$$

275

## Gaussian Processes

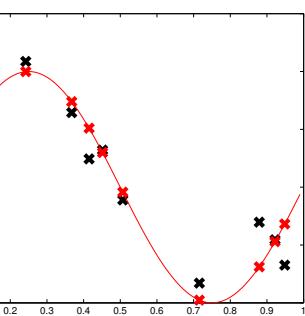
- ▶ We can model response as noisy version of an underlying function  $f(x)$ :

$$y_i | f(x_i) \sim \mathcal{N}(f(x_i), \sigma^2)$$

- ▶ Typical approach: parametrize  $f(x; \beta)$ , and learn  $\beta$ , e.g.,

$$f(x) = \sum_{j=1}^d \beta_j \phi_j(x)$$

- ▶ More direct approach: since  $f(x)$  is unknown, we take a Bayesian approach, introduce a prior over functions, and compute a posterior over functions.



- ▶ Instead of trying to work with the whole function, just work with the function values at the inputs

$$\mathbf{f} = (f(x_1), \dots, f(x_n))^\top$$

274

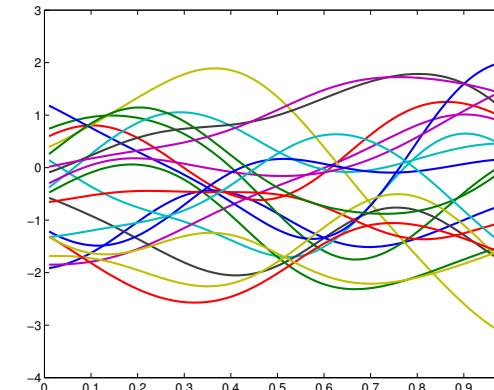
## Gaussian Processes

- ▶ What does a multivariate normal prior mean?
- ▶ Imagine  $\mathbf{x}$  forms a very dense grid of data space. Simulate prior draws

$$\mathbf{f} \sim \mathcal{N}(0, G)$$

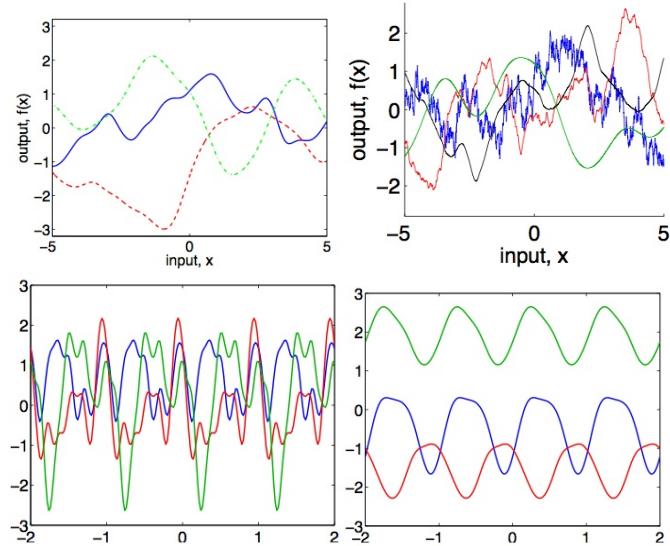
- ▶ Plot  $f_i$  vs  $x_i$  for  $i = 1, \dots, n$ .

- ▶ The prior over functions is called a **Gaussian process** (GP).



## Gaussian Processes

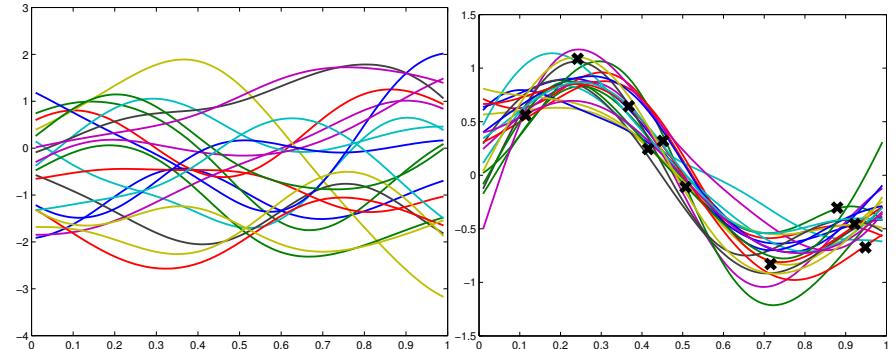
- Different kernels lead to different function characteristics.



Carl Rasmussen. Tutorial on Gaussian Processes at NIPS 2006.

277

## Gaussian Processes



279

## Gaussian Processes

$$\begin{aligned}\mathbf{f}|\mathbf{x} &\sim \mathcal{N}(0, G) \\ \mathbf{y}|\mathbf{f} &\sim \mathcal{N}(\mathbf{f}, \sigma^2 I)\end{aligned}$$

- Posterior distribution:

$$\mathbf{f}|\mathbf{y} \sim \mathcal{N}(G(G + \sigma^2 I)^{-1}\mathbf{y}, G - G(G + \sigma^2 I)G)$$

- Posterior predictive distribution: Suppose  $\mathbf{x}'$  is a test set. We can extend our model to include the function values  $\mathbf{f}'$  at the test set:

$$\begin{pmatrix} \mathbf{f} \\ \mathbf{f}' \end{pmatrix} | \mathbf{x}, \mathbf{x}' \sim \mathcal{N} \left( \begin{pmatrix} 0 \\ 0 \end{pmatrix}, \begin{pmatrix} K_{\mathbf{xx}} & K_{\mathbf{xx}'} \\ K_{\mathbf{x}'\mathbf{x}} & K_{\mathbf{x}'\mathbf{x}'} \end{pmatrix} \right)$$

$$\mathbf{y}|\mathbf{f} \sim \mathcal{N}(\mathbf{f}, \sigma^2 I)$$

where  $K_{\mathbf{zz}'}$  is matrix with  $ij$ th entry  $\kappa(z_i, z'_j)$ .  $K_{\mathbf{xx}} = G$ .

- Some manipulation of multivariate normals gives:

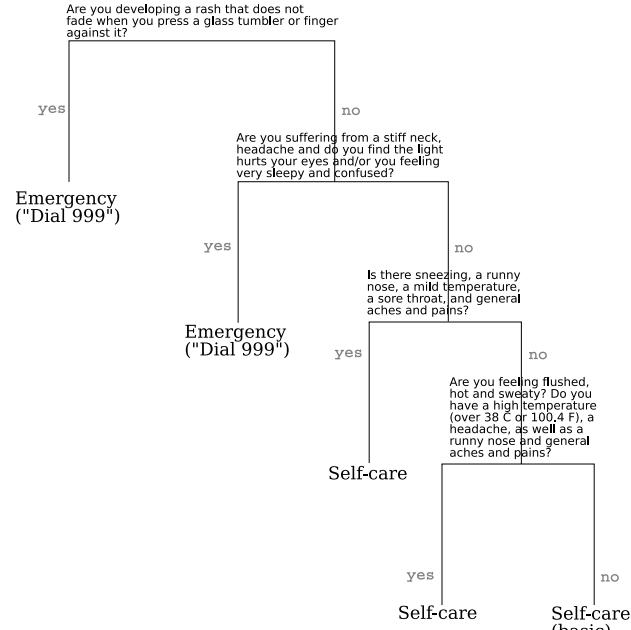
$$\mathbf{f}|\mathbf{y} \sim \mathcal{N} (K_{\mathbf{x}'\mathbf{x}}(K_{\mathbf{xx}} + \sigma^2 I)^{-1}\mathbf{y}, K_{\mathbf{x}'\mathbf{x}'} - K_{\mathbf{x}'\mathbf{x}}(K_{\mathbf{xx}} + \sigma^2 I)^{-1}K_{\mathbf{xx}'})$$

278

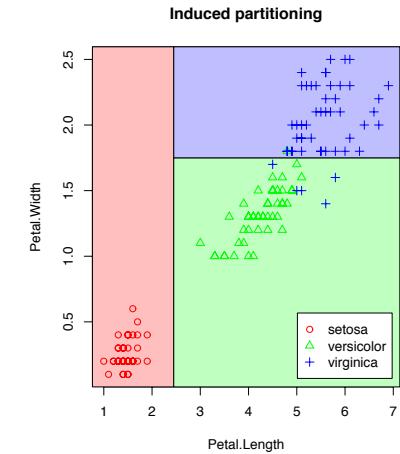
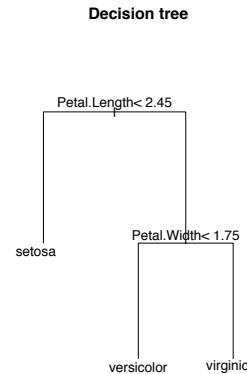
## Example: NHS Direct Self-help Guide

280

## Example: NHS Direct Self-help Guide



## Example: Iris Data



## Example: Iris Data

Sepal.Length	Sepal.Width	Petal.Length	Petal.Width	Species
4.4	3.2	1.3	0.2	setosa
5.9	3.0	5.1	1.8	virginica
6.3	3.3	6.0	2.5	virginica
5.3	3.7	1.5	0.2	setosa
5.5	2.5	4.0	1.3	versicolor
6.1	2.9	4.7	1.4	versicolor
6.1	3.0	4.9	1.8	virginica
5.7	2.8	4.5	1.3	versicolor
5.4	3.0	4.5	1.5	versicolor
4.8	3.4	1.6	0.2	setosa
4.6	3.1	1.5	0.2	setosa
4.9	3.1	1.5	0.2	setosa
6.4	2.9	4.3	1.3	versicolor
.....				

Previously seen Iris data set gives the measurements in centimeters of the variables sepal length and width and petal length and width, respectively, for 50 flowers from each of 3 species of iris. The species are Iris setosa, versicolor, and virginica.

## Decision Trees

- ▶ A decision tree is a hierarchically organized structure, with each node splitting the data space into two halves based on value of a feature.
- ▶ Equivalent to a partition of  $\mathbb{R}^p$  into  $R$  disjoint sets  $(\mathcal{R}_1, \dots, \mathcal{R}_R)$ , where each  $\mathcal{R}_j \subset \mathbb{R}^p$ .
- ▶ On each region  $\mathcal{R}_j$  the same decision/prediction is made  $\hat{f}(x) = \beta_j$  for all  $x \in \mathcal{R}_j$ .
- ▶ Some terminology:
  - ▶ **Parent** of a node  $c$  is the node which have an arrow pointing into  $c$ .
  - ▶ **Children** of a node  $c$  are those nodes which have node  $c$  as a parent.
  - ▶ **Root node** is the top node of the tree; the only node without parents.
  - ▶ **Leaf nodes** are nodes which do not have children.
  - ▶ **Stumps** are trees with just the root node and two leaf nodes.
  - ▶ A  **$K$ -ary tree** is a tree where each node (except for leaf nodes) has  $K$  children. Usually working with binary trees ( $K = 2$ ).
  - ▶ The **depth** of a tree is the maximal length of a path from the root node to a leaf node.

## Classification and Regression Trees

- For regression problems, the parameterized function is

$$\hat{f}(x) = \sum_{j=1}^R \beta_j \mathbb{1}_{[x \in \mathcal{R}_j]},$$

Using squared loss, optimal parameters are:

$$\hat{\beta}_j = \frac{\sum_i y_i \mathbb{1}_{[x_i \in \mathcal{R}_j]}}{\sum_i \mathbb{1}_{[x_i \in \mathcal{R}_j]}}$$

- For classification problems, the estimated probability of each class  $k$  in region  $\mathcal{R}_j$  is simply:

$$\hat{\beta}_{jk} = \frac{\sum_i \mathbb{1}(y_i = k) \mathbb{1}_{[x_i \in \mathcal{R}_j]}}{\sum_i \mathbb{1}_{[x_i \in \mathcal{R}_j]}}$$

- These estimates can be regularized as well.

285

## Growth Heuristic for Regression Trees

- Start with  $\mathcal{R}_1 = \mathbb{R}^p$ .
- For each feature  $j = 1, \dots, p$ , for each value  $v \in \mathbb{R}$  that we can split on:
  - Split data set:

$$I_< = \{i : x_{ij} < v\} \quad I_> = \{i : x_{ij} \geq v\}$$

- Estimate parameters:

$$\beta_< = \frac{\sum_{i \in I_<} y_i}{|I_<|} \quad \beta_> = \frac{\sum_{i \in I_>} y_i}{|I_>|}$$

- Quality of split is the squared loss:

$$\sum_{i \in I_<} (y_i - \beta_<)^2 + \sum_{i \in I_>} (y_i - \beta_>)^2$$

- Choose split with minimal loss.
- Recurse on both children, with datasets  $(x_i, y_i)_{i \in I_<}$  and  $(x_i, y_i)_{i \in I_>}$ .

287

## Partition Estimation

- Ideally, would like to find partition that achieves minimal risk: lowest mean-squared error for prediction or misclassification rate for classification.
- Number of potential partitions is too large to search exhaustively.
- 'Greedy' search heuristics for a good partition:
  - Start at root.
  - Determine best feature and value to split.
  - Recurse on children of node.
  - Stop at some point.

286

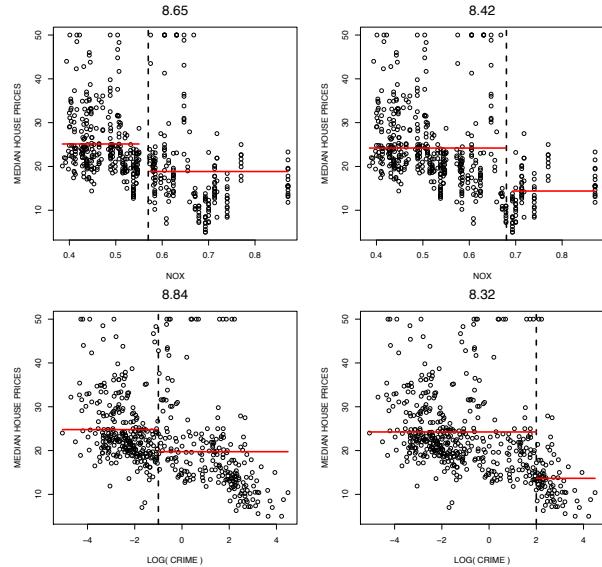
## Boston Housing Data

crim	per capita crime rate by town
zn	proportion of residential land zoned for lots over 25,000
indus	proportion of non-retail business acres per town
chas	Charles River dummy variable (= 1 if tract bounds river; 0 otherwise)
nox	nitric oxides concentration (parts per 10 million)
rm	average number of rooms per dwelling
age	proportion of owner-occupied units built prior to 1940
dis	weighted distances to five Boston employment centres
rad	index of accessibility to radial highways
tax	full-value property-tax rate per USD 10,000
ptratio	pupil-teacher ratio by town
b	$1000(B - 0.63)^2$ where $B$ is the proportion of blacks by town
lstat	percentage of lower status of the population
medv	median value of owner-occupied homes in USD 1000's

- Predict median house value.

288

## Boston Housing Data



289

## Growth Heuristics for Classification Trees

- For binary classification,

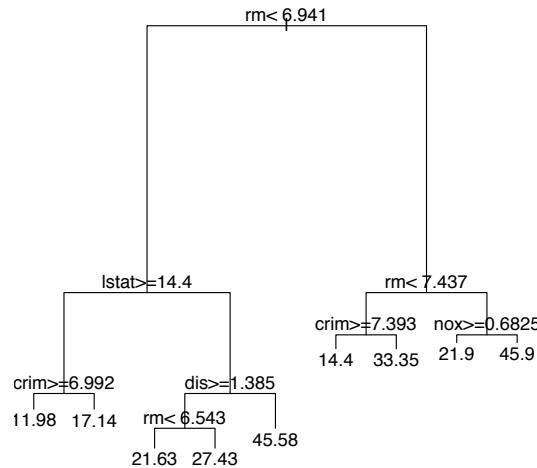
$$\hat{\beta}_{jl} = \frac{\sum_i \mathbb{1}(y_i = 1) \mathbb{1}_{[x_i \in \mathcal{R}_j]}}{\sum_i \mathbb{1}_{[x_i \in \mathcal{R}_j]}}$$

- A split is good if both sides are more “pure”, i.e.  $\hat{\beta}_{jl}$  is closer to 0 or 1.
- Different measures of node impurity:
  - Misclassification error:  $1 - \max\{\hat{\beta}_{jl}, 1 - \hat{\beta}_{jl}\}$ .
  - Gini Index:  $2\hat{\beta}_{jl}(1 - \hat{\beta}_{jl})$ .
  - Entropy:  $-\hat{\beta}_{jl} \log \hat{\beta}_{jl} - (1 - \hat{\beta}_{jl}) \log(1 - \hat{\beta}_{jl})$ .
- Typically prefer Gini and entropy: differentiable and produces purer nodes.
- Extension to multi-class:
  - Misclassification error:  $1 - \max_k \hat{\beta}_{jk}$ .
  - Gini Index:  $\sum_{k=1}^K \hat{\beta}_{jk}(1 - \hat{\beta}_{jk})$ .
  - Entropy:  $-\sum_{k=1}^K \hat{\beta}_{jk} \log \hat{\beta}_{jk}$ .
- Stops once a node has insufficient number of items, or is pure.

291

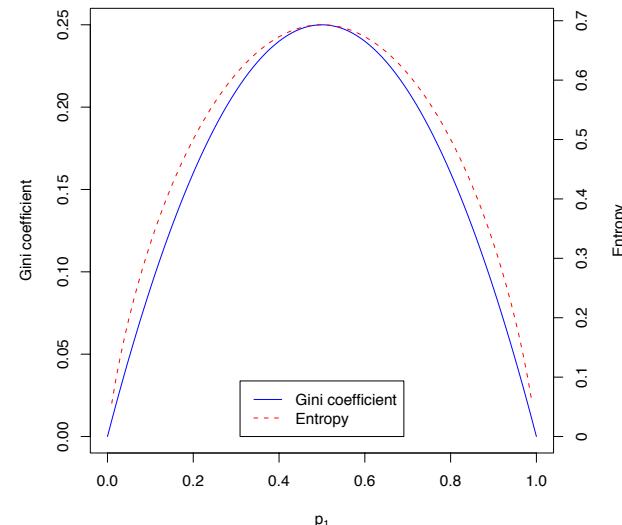
## Boston Housing Data

- Overall, the best first split is on variable `rm`, average number of rooms per dwelling.
- Final tree contains predictions in leaf nodes.



290

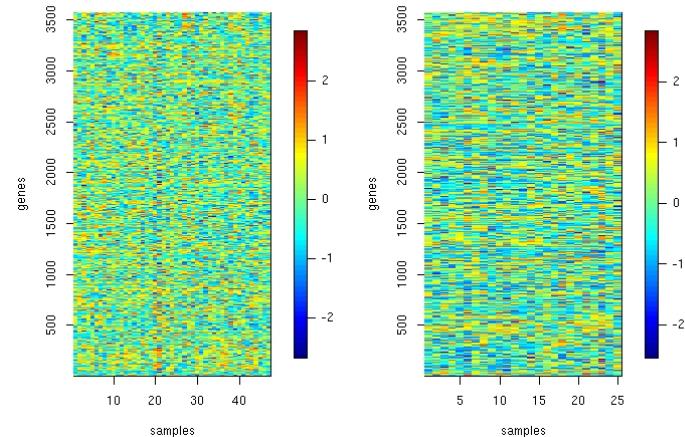
## Growth Heuristics for Classification Trees



Misclassification error?

292

## Example: Leukemia Prediction

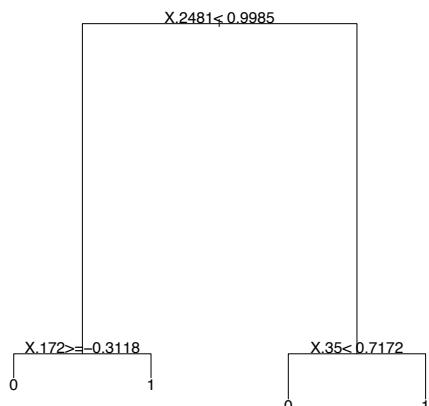


Leukemia Dataset: Expression values of 3541 genes for 47 patients with Leukemia ALL subtype (left) and 25 patients with AML (right).

293

## Example: Leukemia Prediction

- ▶ Tree found is of depth 2.
- ▶ Very interpretable as it selects 3 out of 4088 genes and bases prediction only on these.



294

## Example: Pima Indians Diabetes Dataset

- ▶ The subjects: women who were at least 21 years old, of Pima Indian heritage living near Phoenix, Arizona.
- ▶ Tested for diabetes according to World Health Organisation criteria.
- ▶ Features:
  - ▶ number of pregnancies (npreg)
  - ▶ plasma glucose concentration (glu)
  - ▶ diastolic blood pressure (bp)
  - ▶ tricep skin fold thickness (skin)
  - ▶ body mass index(bbi)
  - ▶ diabetes pedigree function (ped)
  - ▶ age (age)

295

## Example: Pima Indians Diabetes Dataset

```
> library(rpart)
> library(MASS)
> data(Pima.tr)
> rp <- rpart(Pima.tr[,8] ~ ., data=Pima.tr[,-8])
> rp
n= 200

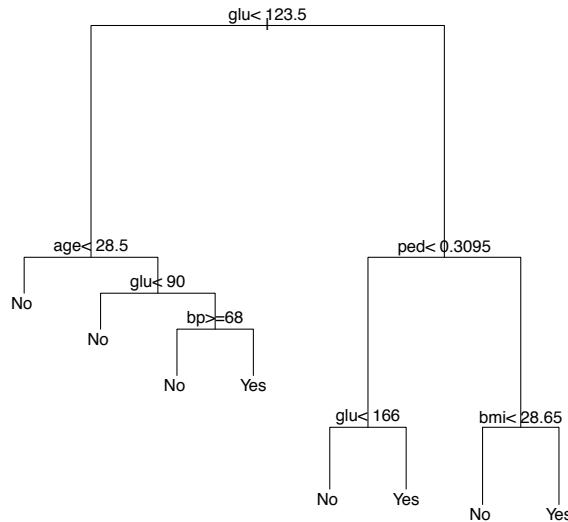
node), split, n, loss, yval, (yprob)
 * denotes terminal node

1) root 200 68 No (0.66000000 0.34000000)
  2) glu< 123.5 109 15 No (0.86238532 0.13761468)
    4) age< 28.5 74 4 No (0.94594595 0.05405405) *
    5) age>=28.5 35 11 No (0.68571429 0.31428571)
      10) glu< 90 9 0 No (1.00000000 0.00000000) *
      11) glu>=90 26 11 No (0.57692308 0.42307692)
        22) bp>=68 19 6 No (0.68421053 0.31578947) *
        23) bp< 68 7 2 Yes (0.28571429 0.71428571) *
  3) glu>=123.5 91 38 Yes (0.41758242 0.58241758)
    6) ped< 0.3095 35 12 No (0.65714286 0.34285714)
      12) glu< 166 27 6 No (0.77777778 0.22222222) *
      13) glu>=166 8 2 Yes (0.25000000 0.75000000) *
    7) ped>=0.3095 56 15 Yes (0.26785714 0.73214286)
    14) bmi< 28.65 11 3 No (0.72727273 0.27272727) *
    15) bmi>=28.65 45 7 Yes (0.15555556 0.84444444) *
```

296

## Example: Pima Indians Diabetes Dataset

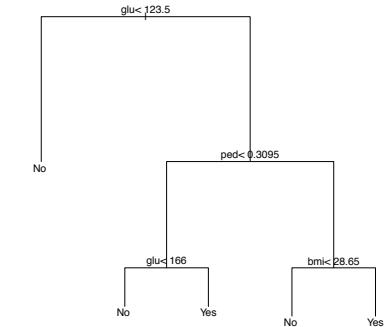
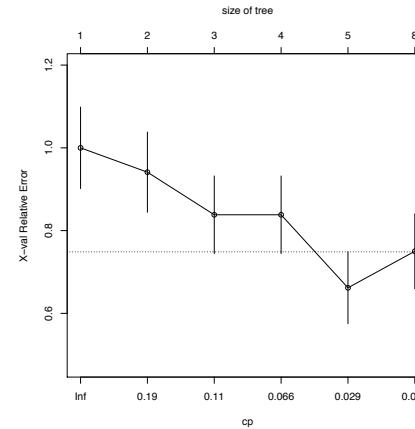
```
> plot(rp); text(rp)
```



297

## Model Complexity

```
> rp <- rpart(Pima.tr[,8] ~ ., data=Pima.tr[,-8],
               control=rpart.control(xval=10)) ## 10-fold CV
> plotcp(rp)
> rp2 <- prune.rpart(rp, .029)
> plot(rp2); text(rp2)
```



299

## Model Complexity

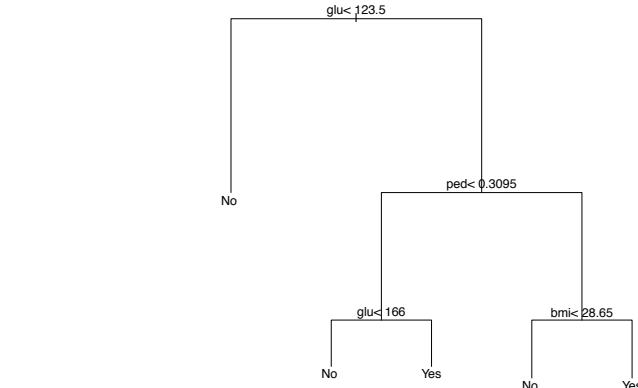
- When should tree growing be stopped?
- Will need to control complexity to prevent overfitting, and in general find optimal tree size with best predictive performance.
- Consider a regularized objective

$$R(T) + C \times \text{size}(T)$$

- Grow the tree from scratch and stop once the criterion objective starts to increase.
- First grow the full tree and prune nodes (starting at leaves), until the objective starts to increase.
- Second option is preferred as the choice of tree is less sensitive to “wrong” choices of split points and variables to split on in the first stages of tree fitting.
- Use cross validation to determine optimal  $C$ .

298

## Model Variability



- Is the tree ‘stable’ if training data were slightly different?

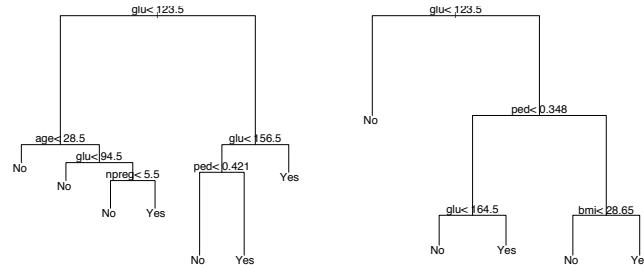
300

## Bootstrap

- The **bootstrap** is a way to assess the variance of estimators.
- Fit multiple trees, each on a **bootstrap sample**. This is a data set obtained by **sampling with replacement**  $n$  times from training set.

```
> n <- nrow(Pima.tr)
> bss <- sample(1:n, n , replace=TRUE)
> sort(bss)
[1] 2 4 4 5 6 7 9 10 11 12 12 12 12 13 13 15 15 20 ...

> tree_boot <- rpart(Pima.tr[bss,8] ~ ., data=Pima.tr[bss,-8],
control=rpart.control(xval=10)) ## 10-fold CV
```



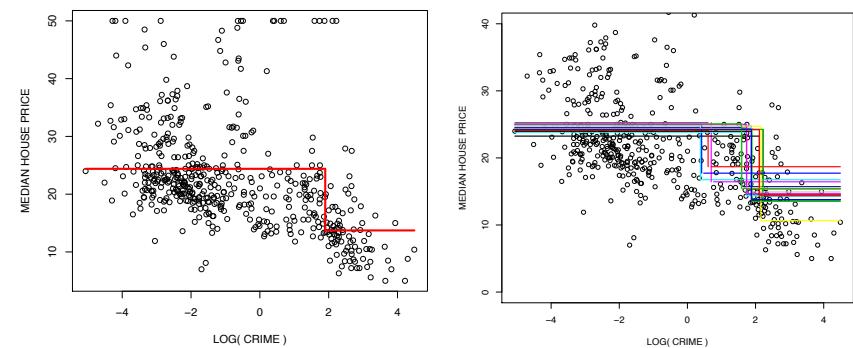
301

## Bootstrap for Regression Trees

- We fit a predictor  $\hat{f}(x)$  on the data  $(x_i, y_i)_{i=1}^n$ .
- Assess the variance of  $\hat{f}(x)$  by taking  $B = 20$  bootstrap samples of the original data, and obtaining bootstrap estimators

$$\hat{f}^b(x), \quad b = 1, \dots, B$$

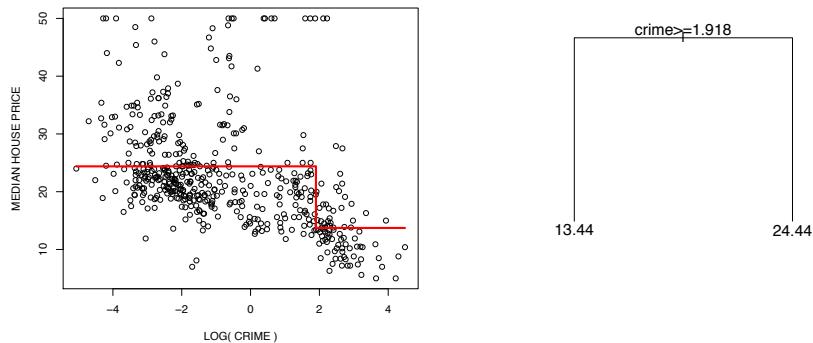
- Each tree  $\hat{f}^b$  is fitted on the resampled data  $(x_{j_i}, y_{j_i})_{i=1}^n$  where each  $j_i$  is chosen randomly from  $\{1, \dots, n\}$  with replacement.



303

## Bootstrap for Regression Trees

- Regression for Boston housing data.
- Predict median house prices based only on crime rate.
- Use decision **stump**—the simplest tree with a single split at root.



302

## Bagging

- Bagging (Bootstrap Aggregation)**: average across all  $B$  trees fitted on different bootstrap samples.

- For  $b = 1, \dots, B$ :

- Draw indices  $(j_1, \dots, j_n)$  from the set  $\{1, \dots, n\}$  with replacement.
- Fit the model, and form predictor  $\hat{f}^b(x)$  based on bootstrap sample

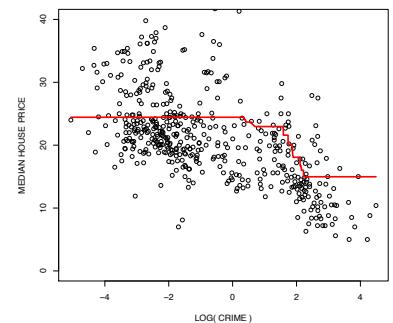
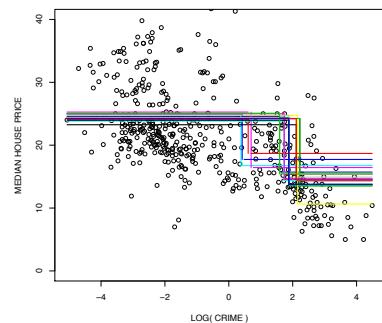
$$(x_{j_1}, y_{j_1}), \dots, (x_{j_n}, y_{j_n})$$

- Form bagged estimator

$$\hat{f}_{Bag}(x) = \frac{1}{B} \sum_{b=1}^B \hat{f}^b(x)$$

304

## Bagging



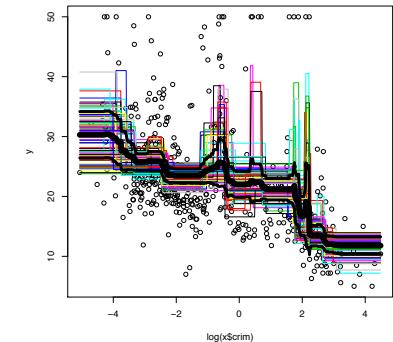
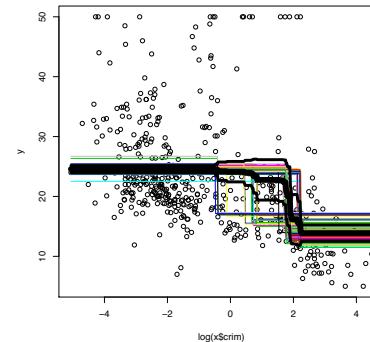
- ▶ Bagging smooths out the drop in the estimate of median house prices.
- ▶ Empirically, bagging seems to reduce the variance of  $\hat{f}$ , i.e.

$$\mathbb{E}[(\hat{f}(x) - \mathbb{E}[\hat{f}(x)])^2 | X = x] \geq \mathbb{E}[(\hat{f}_{Bag}(x) - \mathbb{E}[\hat{f}_{Bag}(x)])^2 | X = x]$$

305

## Variance Reduction in Bagging

- ▶ Deeper trees have higher complexity and variance.
- ▶ Compare bagging trees of depth 1 and 3.



307

## Variance Reduction in Bagging

- ▶ Suppose, in an ideal world, our estimators  $\hat{f}^b$  are based on independent samples of size  $n$  from the true joint distribution of  $X, Y$ .
- ▶ The aggregated estimator would then be

$$\hat{f}_{ag}(x) = \frac{1}{B} \sum_{b=1}^B \hat{f}^b(x) \rightarrow \bar{f}(x) = \mathbb{E}[\hat{f}(x)] \quad \text{as } B \rightarrow \infty$$

where expectation is with respect to datasets of size  $n$ .

- ▶ The squared-loss is:

$$\begin{aligned} \mathbb{E}[(Y - \hat{f}_{ag}(x))^2 | X = x] &= \mathbb{E}[(Y - \bar{f}(x))^2 | X = x] + \mathbb{E}[(\bar{f}(x) - \hat{f}_{ag}(x))^2 | X = x] \\ &\rightarrow \mathbb{E}[(Y - \bar{f}(x))^2 | X = x] \end{aligned}$$

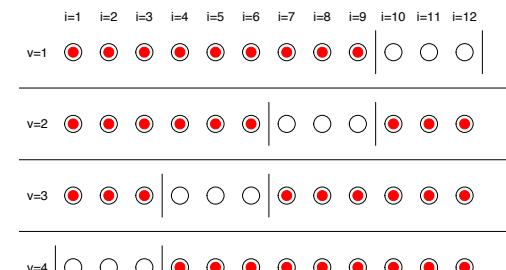
Aggregation reduces the squared loss by eliminating variance of  $\hat{f}(x)$ .

- ▶ In bagging, variance reduction still applies at the cost of a small increase in bias.
- ▶ Bagging is most useful for flexible estimators with high variance (and low bias).

306

## Out-of-bag Test Error Estimation

- ▶ How well does bagging do? Can we estimate generalization performance, and tune hyperparameters?
- ▶ Answer 1: cross-validation.

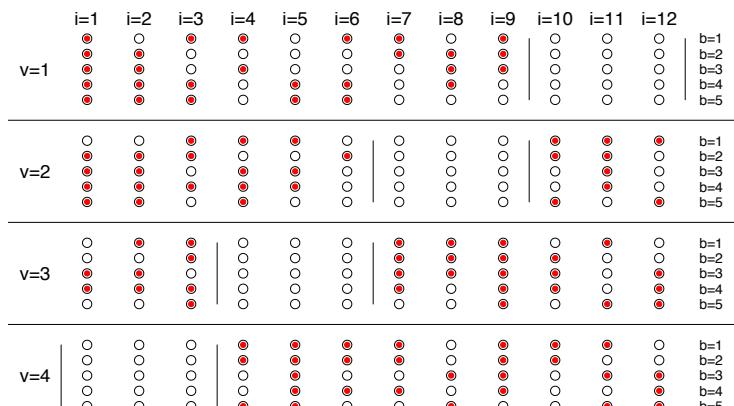


- ▶ For each  $v = 1, \dots, V$ ,
  - ▶ fit  $\hat{f}_{Bag}$  on the training samples.
  - ▶ predict on validation set.
- ▶ Compute the CV error by averaging the loss across all test observations.

308

## Out-of-bag Test Error Estimation

- But to fit  $\hat{f}_{Bag}$  on the training set for each  $v = 1, \dots, V$ , we need another set of  $B$  bootstrap samples!

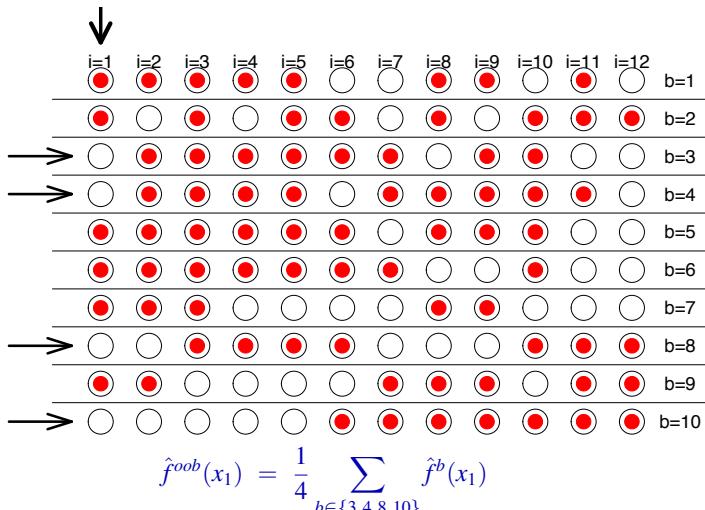


- Answer 2: Out-of-bag test error estimation.

309

## Out-of-bag Test Error Estimation

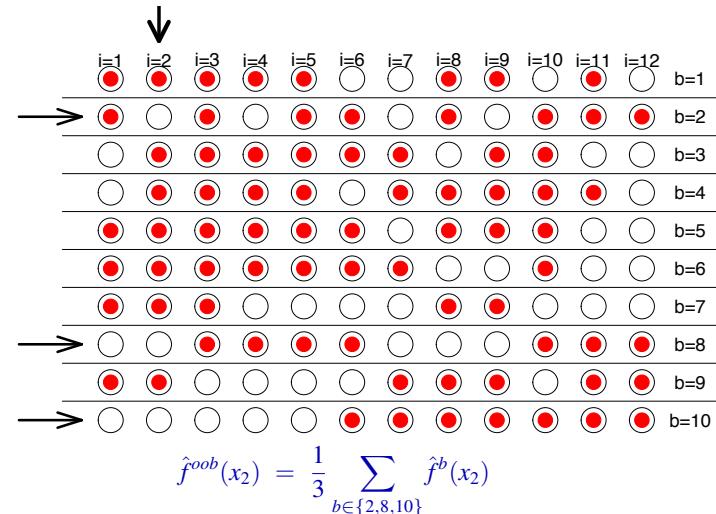
- Idea: test on the “unused” data points in each bootstrap iteration to estimate the test error.



310

## Out-of-bag Test Error Estimation

- Idea: test on the “unused” data points in each bootstrap iteration to estimate the test error.



311

## Out-of-bag Test Error Estimation

- For each  $i = 1, \dots, n$ , the out-of-bag sample is:

$$\tilde{B}_i = \{b : x_i \text{ is not in training set}\} \subseteq \{1, \dots, B\}.$$

- Construct the out-of-bag estimate:

$$\hat{f}^{oob}(x_i) = \frac{1}{|\tilde{B}_i|} \sum_{b \in \tilde{B}_i} \hat{f}^b(i_i)$$

- Estimate the test error as

$$\hat{R}_{test} = \frac{1}{n} \sum_{i=1}^n L(y_i, \hat{f}^{oob}(x_i))$$

312

## Out-of-bag Test Error Estimation

- We need  $|\tilde{B}_i|$  to be reasonably large for all  $i = 1, \dots, n$ .
- The probability  $\pi^{oob}$  of an observation NOT being included in a bootstrap sample  $(j_1, \dots, j_n)$  (and hence being ‘out-of-bag’) is:

$$\pi^{oob} = \prod_{i=1}^n \left(1 - \frac{1}{n}\right) \xrightarrow{n \rightarrow \infty} \exp(-1) \approx 0.367.$$

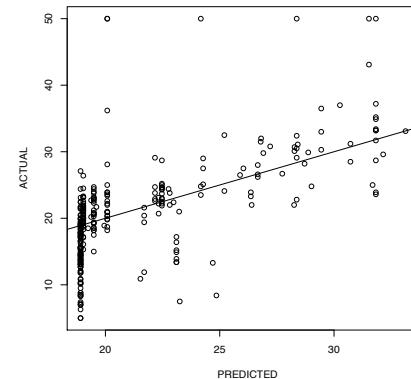
- Hence  $\mathbb{E}[|\tilde{B}_i|] \approx 0.367B$  for all  $i = 1, \dots, n$ .
- In practice, number of bootstrap samples  $B$  is typically between 200 and 1000, meaning that the number  $|\tilde{B}_i|$  of out-of-bag samples will be approximately in the range 70 – 350.
- The obtained test error estimate is asymptotically unbiased for large number  $B$  of bootstrap samples and large sample size  $n$ .

313

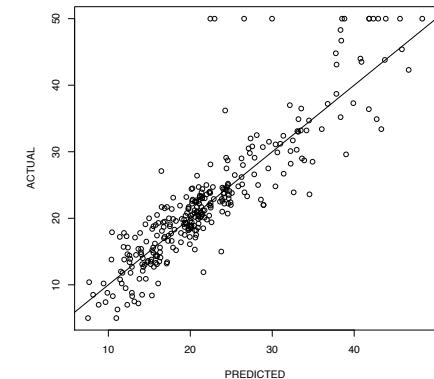
## Example: Boston Housing Dataset

```
plot(prediction_oob, Y, xlab="PREDICTED", ylab="ACTUAL")
```

For depth  $d = 1$ .



For depth  $d = 10$ .



315

## Example: Boston Housing Dataset

- Apply out of bag test error estimation to select optimal tree depth and assess performance of bagged trees for Boston Housing data.
- Use the entire dataset with  $p = 13$  predictor variables.

```
n <- nrow(BostonHousing) ## n samples
X <- BostonHousing[,-14]
Y <- BostonHousing[,14]
B <- 100
maxdepth <- 3
prediction_oob <- rep(0,length(Y)) ## vector with oob predictions
numbertrees_oob <- rep(0,length(Y)) ## number pf oob trees
for (b in 1:B) { ## loop over bootstrap samples
  subsample <- sample(1:n,n,replace=TRUE) ## "in-bag" samples
  outofbag <- (1:n)[-subsample] ## "out-of-bag" samples
  treeboot <- rpart(Y ~ ., data=X, subset=subsample,
    control=rpart.control(maxdepth=maxdepth,minsplit=2))
    ## predict on oob-samples
  prediction_oob[outofbag] <- prediction_oob[outofbag] +
    predict(treeboot, newdata=X[outofbag,])
  numbertrees_oob[outofbag] <- numbertrees_oob[outofbag] + 1
}
## final oob-prediction is average across all "out-of-bag" trees
prediction_oob <- prediction_oob / numbertrees_oob
```

314

## Example: Boston Housing Dataset

- Out-of-bag estimate of test error as a function of tree depth  $d$ :

tree depth $d$	1	2	3	4	5	10	30
single tree $\hat{f}$	60.7	44.8	32.8	31.2	27.7	26.5	27.3
bagged trees $\hat{f}_{Bag}$	43.4	27.0	22.8	21.5	20.7	20.1	20.1

- Without bagging, the optimal tree depth seems to be  $d = 10$ .
- With bagging, we could also take the depth up to  $d = 30$ .
- Bagging strongly improves performance.
- On the other hand, bagged trees cannot be displayed as nicely as single trees and some of the interpretability of trees is lost.
- Bagging does not always improve accuracy, but often does in practice.

316

## Random Forests and Extremely Randomized Trees

- ▶ **Random forests** are similar to bagged decision trees with a few key differences:
  - ▶ For each split point, the search is not over all  $p$  variables but just over  $mtry$  randomly chosen ones (where e.g.  $mtry = \lfloor p/3 \rfloor$ )
  - ▶ No pruning necessary. Trees can be grown until each node contains just very few observations (1 or 5).
  - ▶ Random forests tend to produce better predictions than bagging.
  - ▶ Results often not sensitive to the only tuning parameter  $mtry$ .
  - ▶ Implemented in `randomForest` library.
- ▶ Even more random methods, e.g. **extremely randomized trees**:
  - ▶ For each split point, sample  $mtry$  variables each with a value to split on, and pick the best one.
  - ▶ Often works even when  $mtry$  equals 1!
- ▶ Often produce state-of-the-art results, and top performing methods in machine learning competitions.

Breiman (2001)., Geurts et al (2006).

317

## Random Forests

TABLE 2  
*Test set misclassification error (%)*

Data set	Forest	Single tree
Breast cancer	2.9	5.9
Ionosphere	5.5	11.2
Diabetes	24.2	25.3
Glass	22.0	30.4
Soybean	5.7	8.6
Letters	3.4	12.4
Satellite	8.6	14.8
Shuttle $\times 10^3$	7.0	62.0
DNA	3.9	6.2
Digit	6.2	17.1

From Breiman: "Statistical Modelling: the two cultures".

## Random Forests and Nearest Neighbours

- ▶ Let  $P(x, x_i) \in [0, 1]$  be the proportion of trees for which a vector  $x$  falls into the same final leaf node as the training vector  $x_i$ .  $P(x, x_i)$  is a proximity value, and tends to be large when  $x$  and  $x_i$  are close by.
- ▶ If every leaf node contains the same number of training samples, the prediction of random forests (in regression mode) at  $x$  is:

$$\hat{f}^{RF}(x) = \frac{\sum_{i=1}^n P(x, x_i) y_i}{\sum_{i=1}^n P(x, x_i)},$$

which is a weighted (approximate) nearest neighbour scheme.

- ▶ If the nodes contain different number of original observations,  $P(x, x_i)$  is a weighted proportion of trees, where the weight of a tree is inversely proportional to the number of samples in the leaf node containing  $x$ .
- ▶ For classification, the prediction will be the weighted majority vote, where again weights are proportional to the proximities  $P(x, x_i)$ .
- ▶ kNN does not scale well to very large datasets, and computational geometry techniques for approximately finding nearest neighbours often rely on tree data structures, e.g. kd-trees, cover trees, ball trees. Random forests and other randomized trees can also be thought of similarly.

319

## Ensemble Methods

- ▶ Bagging and random forests are examples of **ensemble methods**, where predictions are based on an ensemble of many individual predictors.
- ▶ Bayesian posterior averaging can also be thought of as an ensemble method:

$$p(y|x, \text{Data}) = \int p(y|x, \theta)p(\theta|\text{Data})d\theta \\ \approx \frac{1}{B} \sum_{b=1}^B p(y|x, \theta^b)$$

where  $\theta^b$  are samples drawn from posterior  $p(\theta|\text{Data})$ .

- ▶ Many other ensemble learning methods: boosting, stacking, mixture of experts, Bayesian model combination etc.
- ▶ Often gives significant boost to predictive performance.

318

320

## Dropout Training of Neural Networks

- ▶ Neural network with single layer of hidden units:

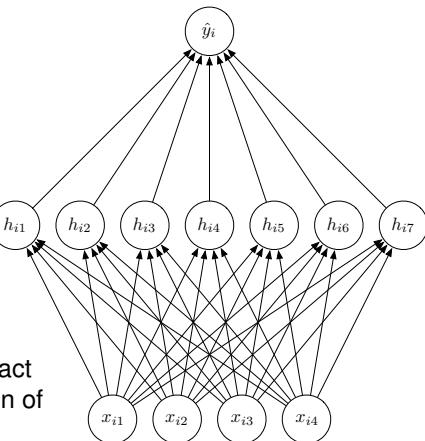
- ▶ **Hidden unit activations:**

$$h_{ik} = s \left( b_k^h + \sum_{j=1}^p W_{jk}^h x_{ij} \right)$$

- ▶ **Output probability:**

$$\hat{y}_i = s \left( b^o + \sum_{k=1}^m W_k^o h_{ik} \right)$$

- ▶ Large, overfitted, networks often have co-adapted hidden units.
- ▶ What each hidden unit learns may in fact be useless, e.g. predicting the negation of predictions from other units.
- ▶ Can prevent co-adaptation by randomly **dropping out** units from network.



Hinton et al (2012). 321

## Dropout Training of Neural Networks

Classification of phonemes in speech.

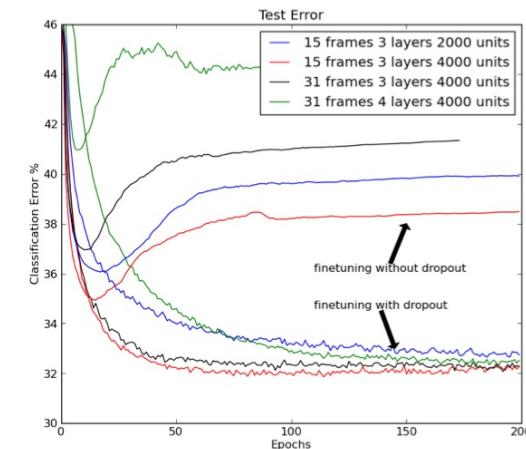


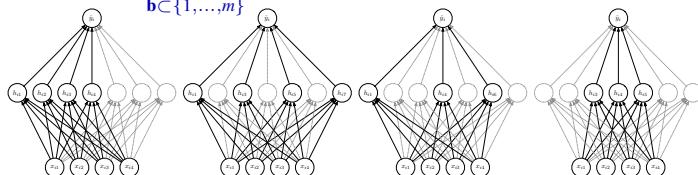
Figure from Hinton et al.

323

## Dropout Training of Neural Networks

- ▶ Model as an ensemble of networks:

$$p(y_i = 1|x_i, \theta) = \sum_{\mathbf{b} \subset \{1, \dots, m\}} q^{|\mathbf{b}|} (1 - q)^{n - |\mathbf{b}|} p(y_i = 1|x_i, \theta, \text{drop out units } \mathbf{b})$$



- ▶ **Weight-sharing** among all networks: each network uses a subset of the parameters of the full network (corresponding to the retained units).
- ▶ Training by stochastic gradient descent: at each iteration a network is sampled from ensemble, and its subset of parameters are updated.

322

## Boosting

- ▶ **Boosting** is an iterative ensemble learning technique. At iteration  $t$ , the predictor is (with  $0 < \nu < 1$ , typically small, say  $\nu = 0.1$ ):

$$\hat{f}_t(x) = \sum_{m=1}^t \nu \hat{g}_m(x)$$

- ▶ For regression,  $L_2$ -boosting works as follows:

1. Fit a first function to the data  $(x_i, y_i)_{i=1}^n$  with **base learner**, yielding  $\hat{g}_1(x)$ .
2. For  $t = 2, 3, \dots, T$  do:
  - 2.1 Compute residuals  $u_i = y_i - \hat{f}_{t-1}(x_i)$

$$u_i = y_i - \hat{f}_{t-1}(x_i)$$

- 2.2 Fit the residuals  $(x_i, u_i)_{i=1}^n$ , obtaining  $\hat{g}_t(x)$ .

- ▶ Boosting is a bias-reduction technique, as opposed to bagging and dropout.
- ▶ Boosting works well with simple base learners, e.g. decision stumps, with low variance and high bias.
- ▶ Implemented in the `mboost` library.

324

## Boosting

- ▶ Boosting can be viewed as functional gradient descent.
- ▶ Say we wish to minimize empirical risk with differentiable loss function,

$$R(f) = \frac{1}{n} \sum_{i=1}^n L(y_i - f(x_i))$$

- ▶ We can calculate  $\nabla_f R(f)$  and take a gradient descent step:

$$\hat{f}_t(x) = \hat{f}_{t-1}(x) - \nu \nabla_f R(\hat{f}_{t-1})$$

- ▶ We use a base learner to approximate  $\nabla_f R(\hat{f}_{t-1})$ .

1. With  $\hat{f}_0 \equiv 0$ , fit a first function to  $\{(x_i, -\nabla_f L(y_i, \hat{f}_0(x_i)))\}_{i=1}^n$ , yielding  $\hat{g}_1(\cdot)$ .
2. Take gradient descent step:

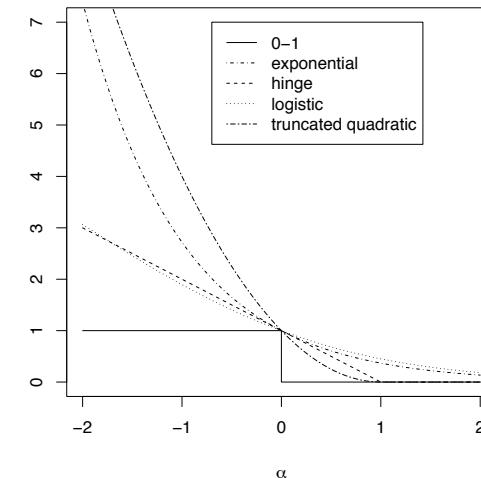
$$\hat{f}_1(x) = \nu \hat{g}_1(x)$$

3. For  $t = 2, 3, \dots, T$  do:

- Compute empirical gradient  $u_i = -\nabla_f L(y_i, \hat{f}_{t-1}(x_i))$ .
- Fit the gradient  $(x_i, u_i)_{i=1}^n$ , yielding  $\hat{g}_t(x)$ .
- Set

$$\hat{f}_t(x) = \hat{f}_{t-1}(x) + \nu \hat{g}_t(x)$$

## Boosting



325

327

## Boosting

- ▶ Obtain  $L_2$ -Boosting when using the quadratic loss function

$$R(f) = \frac{1}{n} \sum_{i=1}^n (y_i - f(x_i))^2.$$

- ▶ Obtain LogitBoost when using the logistic loss function (binary classification,  $y_i \in \{-1, 1\}$ ),

$$R(f) = \frac{1}{n} \sum_{i=1}^n \log(1 + \exp(-y_i f(x_i))).$$

- ▶ Obtain AdaBoost (the original boosting algorithm) when using the exponential loss (again  $y_i \in \{-1, 1\}$ )

$$R(f) = \frac{1}{n} \sum_{i=1}^n \exp(-y_i f(x_i)).$$

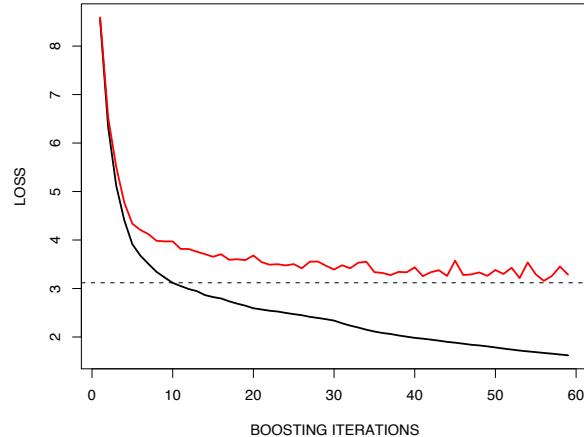
## Boosting

```
library(mboost)
n <- length(y)                                ## number of observations
Mvec <- 1:500                                   ## Mvec is vector with various stopping times
nM <- length(Mvec)                             ## number of possible stopping times
loss <- numeric(nM)                            ## loss contains the training error
losscv <- numeric(nM)                           ## losscv contains the validation error
for (mc in 1:nM){                               ## loop over stopping times (not efficient)
  yhat <- numeric(n)                            ## yhat are the fitted values
  yhatcv <- numeric(n)                           ## yhatcv the cross-validated fitted values
  M <- Mvec[mc]                                 ## use M iterations
  V <- 10                                       ## 10-fold cross validation
                                                 ## indCV contains the 'block' in 1,...,10
                                                 ## each observation falls into
  indCV <- sample(rep(1:V,each=ceiling(n/V)), n)
  for (cv in 1:V){                               ## loop over all blocks
    bb <- blackboost(y[indCV!=cv] ~ .,data=x[indCV!=cv,],
                      control=boost_control(msstop=M))
    ## predict the unused observations
    yhatcv[indCV==cv] <- predict(bb,x[indCV==cv,])
  }
  losscv[mc] <- sqrt(mean((y-yhatcv)^2))      ## CV test error
  bb <- blackboost(y ~ .,data=x,control=boost_control(msstop=M))
  yhat <- predict(bb,x)
  loss[mc] <- sqrt(mean((y-yhat)^2))           ## training error
}
```

## Boosting

Plot of validation error in red and training error in black as functions of iteration.

```
matplot( cbind(loss,losscv), type="p", lwd=2,col=c(1,2),lty=1)
abline(h= sqrt(mean(( predict(rf)-y)^2)),lwd=1,lty=2)
```



329

## Machine Learning

- ▶ The last 8 weeks has been a whirlwind tour of basic machine learning techniques:
  - ▶ Unsupervised learning: PCA, MDS, Isomap, K-means, mixture model, EM algorithm.
  - ▶ Supervised learning: LDA, naïve Bayes, kNN, decision trees, logistic regression, SVMs, kernel methods, Gaussian processes, neural networks, deep learning.
  - ▶ Conceptual framework: prediction, generalization, overfitting and regularization.
  - ▶ Theory: decision theory, statistical learning theory, Bayesian framework.
  - ▶ Cross validation, model selection, model averaging and model combination (bagging, random forests, boosting).
- ▶ Further explorations:
  - ▶ Machine learning summer schools, videolectures.net.
  - ▶ Conferences: NIPS, ICML, UAI, AISTATS.
  - ▶ Mailing list: ml-news.

330