

MS1b: SDM - Problem Sheet 5

1. (Missing Features) Assume you have used some training data $D = \{x_i, y_i\}_{i=1}^n$ where $x_i \in \mathbb{R}^p$ and $y_i \in \{1, 2, \dots, K\}$ to learn a probabilistic classifier (say using maximum likelihood). We are interested in classifying a new input vector. However we have only been able to collect $p - 1$ features, say you have only $(x_i^1, \dots, x_i^{l-1}, x_i^{l+1}, \dots, x_i^p)$ and x_i^l is missing. Explain whether or not it is possible to use your classifier to classify this incomplete input vector in the two following scenarios:

- (a) When we consider a naïve Bayes model where

$$f(x | \phi_k) = \prod_{l=1}^p g(x^l | \phi_k^l);$$

i.e. conditional upon $Y = k$, you assume that the features are independent and feature x^l follows a distribution with density $g(x^l | \phi_k^l)$.

- (b) When we consider a QDA model; i.e.

$$f(x | \phi_k) = \mathcal{N}(x; \mu_k, \Sigma_k).$$

- (c) Generally speaking, which conditions on $f(x | \phi_k)$ are necessary to allow us to implement easily, that is without using any numerical integration scheme, a probabilistic classifier in presence of missing features?

2. (Bayesian classification) Consider some training data $D = \{(x_i, z_i), y_i\}_{i=1}^n$ where $(x_i, z_i) \in \mathbb{R}^p \times \mathbb{R}$ is the vector of inputs and $y_i \in \{0, 1\}$ the response. We adopt the following regression model for class k

$$Z = \beta_k^T X + \varepsilon$$

where $\varepsilon \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, \sigma_k^2)$ if $Y = k$. Hence we have for the class conditional density $f_k(z | x) = \mathcal{N}(z; \beta_k^T x, \sigma_k^2)$ so that the unconditional density of Z follows a so-called mixture of regressions model. Note that this model differs conceptually from the examples discussed in lectures as we do not model X . We adopt the notation $P(Y = k) = \pi_k$ and denote $\theta = (\pi_1, \beta_0, \beta_1, \sigma_0^2, \sigma_1^2)$ the set of unknown parameters.

- (a) Give an expression of the estimate $\hat{\theta}$ of θ maximizing the conditional log-likelihood

$$l(\theta) = \sum_{i=1}^n \log p(y_i, z_i | x_i, \theta).$$

What happens when $n < p$?

- (b) Consider a Bayesian approach with $\pi_1 \sim \text{Beta}(a, b)$ and

$$p(\sigma_0^2, \beta_0, \sigma_1^2, \beta_1) = p(\sigma_0^2, \beta_0) p(\sigma_1^2, \beta_1)$$

where $p(\sigma_k^2, \beta_k)$ satisfies a normal inverse-Gamma distribution; e.g.

$$\begin{aligned} p(\sigma_k^2, \beta_k) &= p(\sigma_k^2) p(\beta_k | \sigma_k^2) \\ &= \mathcal{IG}\left(\sigma_k^2; \frac{\nu}{2}, \frac{\kappa}{2}\right) \mathcal{N}(\beta_k; 0, \sigma_k^2 \Sigma) \end{aligned}$$

with some hyperparameters $(\nu, \kappa, \delta\Sigma)$ such that $\nu, \kappa > 0$ and Σ is a positive definite matrix. $\mathcal{IG}\left(\sigma^2; \frac{\nu}{2}, \frac{\kappa}{2}\right)$ denotes the inverse-Gamma density given by

$$\mathcal{IG}\left(\sigma^2; \frac{\nu}{2}, \frac{\kappa}{2}\right) = \frac{\left(\frac{\kappa}{2}\right)^{\frac{\nu}{2}}}{\Gamma\left(\frac{\nu}{2}\right)} (\sigma^2)^{-\frac{\nu}{2}-1} \exp\left(-\frac{\kappa}{2\sigma^2}\right).$$

The posterior distribution $p(\theta | D)$ satisfies

$$p(\theta | D) = p(\pi_1 | D) p(\sigma_0^2, \beta_0 | D) p(\sigma_1^2, \beta_1 | D).$$

Show that $p(\pi_1 | D)$ is a Beta distribution and $p(\sigma_k^2, \beta_k | D)$ a normal inverse-Gamma distribution.

(c) Given a new test data (x, z) , establish the expression of $p(z | D, x, y = k)$ and explain how you would use this expression to obtain a Bayesian classifier. What are the potential benefits of this approach over using $p(z | \hat{\theta}, x, y = k)$?

3. (Logistic Regression) Consider two-class data, (X, Y) with $Y \in \{-1, 1\}$ and $X = (B, Z)$ with $B \in \{0, 1\}$ and $Z \in \mathbb{R}^{p-1}$, so data vectors are made up of a binary variable and $p - 1$ continuous variables. Training data (X_i, Y_i) , $i = 1, \dots, n$ are available. Consider logistic regression of the log posterior odds, $\log P(Y = 1|x) - \log P(Y = -1|x) = \alpha + \beta^T x$ with $\beta = (\beta_b, \beta_z) \in \mathbb{R} \times \mathbb{R}^{p-1}$, so that $\alpha + \beta^T x \equiv \alpha + \beta_b b + \beta_z^T z$.

(a) Suppose there are no $B = 1$ outcomes in the class $Y = -1$ data (so $B_i = 0$ for all $i \in \{j : Y_j = -1\}$) but there are both $B = 0$ and $B = 1$ outcomes in the $Y = 1$ data. Show that the maximum likelihood estimate for β_b is $\hat{\beta}_b = \infty$.

(b) Problems of this kind arise when there exists a separating-hyperplane. Show that the assumptions of 3a do not imply the existence of a separating hyperplane.

4. Load the Vanveer gene expression data used in a previous practical. Make use of the 20 'best' genes (according to a marginal t-test) by using the following commands.

```
load(url("http://www.stats.ox.ac.uk/%7Eteh/MS1b/PracticalObjects.RData"))
vanv<- vanveer.4000[,2:21]
prog<- vanveer.4000[,1]
```

Your X matrix is thus `vanv` and the response Y is `prog`. Split the data into a test and training set (of equal size). Using logistic regression, plot a ROC curve.