

Stat 461-561: Quiz 2

Friday 29th January 2008

• **Exercise 1.** Let $X_i \stackrel{\text{i.i.d.}}{\sim} g(x)$ and assume that we want to model these data using the parametrized family of probability density functions (pdf) $\{f(x|\theta); \theta \in \Theta\}$. Let θ_n be the Maximum Likelihood Estimate (MLE) for n observations; that is

$$\theta_n = \arg \max_{\theta \in \Theta} \sum_{i=1}^n \log f(X_i | \theta)$$

Under ‘suitable’ regularity assumptions, we have

$$\sqrt{n}(\theta_n - \theta^*) \xrightarrow{D} \mathcal{N}(0, \sigma^2)$$

Question 1.1: [1 point] Give the expression of the Kullback-Leibler divergence minimized in $\theta = \theta^*$.

Question 1.2: [1 point] Establish the expression of σ^2 .

See lecture notes.

• **Exercise 2.** Consider n (positive) observations X_1, \dots, X_n where $X_i \stackrel{\text{i.i.d.}}{\sim} g(x)$ with

$$g(x) = \begin{cases} \pi \lambda_1 \exp(-\lambda_1 x) + (1 - \pi) \lambda_2 \exp(-\lambda_2 x) & \text{for } x \geq 0 \\ 0 & \text{for } x < 0 \end{cases}$$

with $0 < \pi < 1$, $\lambda_1 > 0$ and $\lambda_2 > 0$. That is we have modeled the pdf $g(x)$ of the data by a mixture of two exponential distributions.

Given X_1, \dots, X_n , we are interested in estimating $\theta = (\pi, \lambda_1, \lambda_2)$ using the Expectation-Maximization (EM) algorithm.

Question 2.1: [1 point] Which set of latent variables can you introduce to implement the EM algorithm? Describe the associated statistical model.

We can associate to each observation a latent variable $Z_i \in \{1, 0\}$ such that

$$p(z_i = 1 | \theta) = 1 - p(z_i = 0 | \theta) = \pi$$

and

$$p(x_i, z_i | \theta) = \begin{cases} \pi \lambda_1 \exp(-\lambda_1 x) & \text{if } z_i = 1 \\ (1 - \pi) \lambda_2 \exp(-\lambda_2 x) & \text{if } z_i = 0. \end{cases}$$

Then the complete log-likelihood is given by

$$p(x_{1:n}, z_{1:n} | \theta) = \prod_{i=1}^n [\pi \lambda_1 \exp(-\lambda_1 x_i)]^{z_i} [(1 - \pi) \lambda_2 \exp(-\lambda_2 x_i)]^{1-z_i}$$

Question 2.2: [3 points] Establish the EM recursion giving the expression of the parameter estimate $\theta^{(k)}$ at iteration k given $\theta^{(k-1)}$ at iteration $k - 1$.

We have

$$\begin{aligned}
 Q(\theta, \theta^{(k)}) &= \sum \log p(x_{1:n}, z_{1:n} | \theta) \cdot p(z_{1:n} | \theta^{(k)}, x_{1:n}) \\
 &= \sum_{z_{1:n} \in \{0,1\}^n} \sum_{i=1}^n [z_i (\log \pi + \log \lambda_1 - \lambda_1 x_i) \\
 &\quad + (1 - z_i) (\log \pi + \log \lambda_2 - \lambda_2 x_i)] p(x_{1:n} | \theta^{(k)}, z_{1:n}) \\
 &= (\log \pi + \log \lambda_1) \sum_{i=1}^n p(z_i = 1 | \theta^{(k)}, x_i) - \lambda_1 \sum_{i=1}^n x_i p(z_i = 1 | \theta^{(k)}, x_i) \\
 &\quad + (\log(1 - \pi) + \log \lambda_2) \sum_{i=1}^n p(z_i = 0 | \theta^{(k)}, x_i) - \lambda_2 \sum_{i=1}^n x_i p(z_i = 0 | \theta^{(k)}, x_i).
 \end{aligned}$$

We obtain

$$\pi = \frac{\sum_{i=1}^n p(z_i = 1 | \theta^{(k)}, x_i)}{n}$$

and

$$\begin{aligned}
 \lambda_1 &= \frac{\sum_{i=1}^n p(z_i = 1 | \theta^{(k)}, x_i)}{\sum_{i=1}^n x_i p(z_i = 1 | \theta^{(k)}, x_i)}, \\
 \lambda_2 &= \frac{\sum_{i=1}^n p(z_i = 0 | \theta^{(k)}, x_i)}{\sum_{i=1}^n x_i p(z_i = 0 | \theta^{(k)}, x_i)}
 \end{aligned}$$

where

$$\begin{aligned}
 p(z_i = 1 | \theta^{(k)}, x_i) &= 1 - p(z_i = 0 | \theta^{(k)}, x_i) \\
 &= \frac{\pi^{(k)} \lambda_1^{(k)} \exp(-\lambda_1^{(k)} x_i)}{\pi^{(k)} \lambda_1^{(k)} \exp(-\lambda_1^{(k)} x_i) + (1 - \pi^{(k)}) \lambda_2^{(k)} \exp(-\lambda_2^{(k)} x_i)}.
 \end{aligned}$$

• **Exercise 3.** Consider the following simple polynomial regression model

$$y = \beta_0 + \beta_1 x + \beta_2 x^2 + \dots + \beta_k x^k + \varepsilon, \quad \varepsilon \sim \mathcal{N}(0, \sigma^2).$$

Question 3.1: [**2 points**] Assuming k is known, establish the Maximum likelihood estimate of $(\hat{\beta}_0, \hat{\beta}_1, \dots, \hat{\beta}_k, \hat{\sigma}^2)$ given n observations $\{x_i, y_i\}_{i=1, \dots, n}$.

We can rewrite the observations as

$$\mathbf{y} = \mathbf{X}\beta + \epsilon$$

where

$$\mathbf{y} = \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix}, \mathbf{X} = \begin{pmatrix} 1 & x_1 & \cdots & x_1^k \\ 1 & x_1 & \cdots & x_n^k \end{pmatrix},$$

$$\boldsymbol{\beta} = (\beta_0, \beta_1, \dots, \beta_k)^T$$

and $\boldsymbol{\epsilon} \sim \mathcal{N}(0, \sigma^2 I_n)$. So the log-likelihood is

$$l(\boldsymbol{\beta}, \sigma^2) = -\frac{n}{2} \log(2\pi) - \frac{n}{2} \log(\sigma^2) - \frac{(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})}{2\sigma^2}$$

and

$$\hat{\boldsymbol{\beta}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y},$$

$$\hat{\sigma}^2 = \frac{1}{n} (\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}})^T (\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}})$$

Question 3.2: [2 points] We now want to determine k through the Akaike Information Criterion. Establish that we have

$$AIC(k) = n(\log 2\pi + 1) + n \log \hat{\sigma}^2 + 2(k + 2)$$

We have

$$\begin{aligned} AIC(k) &= -2 \log \text{likelihood at MLE} + 2 (\text{number of parameters}) \\ &= n \log(2\pi) + n \log(\hat{\sigma}^2) + \frac{(\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}})^T (\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}})}{\hat{\sigma}^2} \\ &\quad + 2(k + 1) \\ &= n(\log(2\pi) + 1) + n \log(\hat{\sigma}^2) + 2(k + 2). \end{aligned}$$