

## Series 2, Oct 24th, 2016 (MLE)

Email questions 1,2 to: **Andrey Ignatov**  
ihnatova@student.ethz.ch  
Email questions 3,4,5 to: **David Zhao**  
zhaoda@student.ethz.ch

### Problem 1 (MLE):

Suppose we have a family of probability distributions  $p(x|\theta)$  that depends on parameters  $\theta = \theta_1, \theta_2, \dots, \theta_n$ , and the set  $X = \{x_i\}_{i=1}^N$  generated according to some concrete  $\theta$ . Our goal is to find this  $\theta$ , and to do this we are looking for  $\hat{\theta}$  that maximizes the likelihood of  $X$ :

$$\hat{\theta} = \arg \max_{\theta} L(X|\theta) = \arg \max_{\theta} \prod_{i=1}^N p(x_i|\theta). \quad (1)$$

If  $L$  is continuous and differentiable, its derivative in a stationary point is zero, and  $\hat{\theta}$  can be found from equation:  $\frac{\partial}{\partial \theta} L(X|\theta)|_{\theta=\hat{\theta}} = 0$ . Otherwise we have to apply some optimization techniques or try to find the solution analytically. In this exercise, you have to derive  $\hat{\theta}$  for the following distributions:

- **Exponential( $\lambda$ ):**  $p(x|\lambda) = \lambda e^{-\lambda x}$ ,  $x > 0$ ,  $\lambda > 0$ ,  $\theta = \lambda$

*Hint:* it is easier to maximize the logarithm of  $L$ .

*Solution.* The Likelihood of the dataset  $X$  is given by:

$$\begin{aligned} L(X|\lambda) &= \prod_{i=1}^N p(x_i|\lambda) = \lambda^N e^{-\lambda \sum_{i=1}^N x_i} \\ \log L(X|\lambda) &= N \log \lambda - \lambda \sum_{i=1}^N x_i \\ \frac{\partial}{\partial \lambda} \log L(X|\lambda) &= \frac{N}{\lambda} - \sum_{i=1}^N x_i = 0 \quad \rightarrow \quad \boxed{\hat{\lambda} = \frac{N}{\sum_{i=1}^N x_i}} \end{aligned}$$

Since all  $x_i \geq 0$ , the obtained  $\hat{\lambda}$  is also positive. Finally, we have to verify that  $\hat{\lambda}$  is a maximum, and for this purpose we will show that the second derivative  $\frac{\partial^2}{\partial \lambda^2} \log L(X|\lambda)$  is negative:

$$\frac{\partial^2}{\partial \lambda^2} \log L(X|\lambda) = -\frac{N}{\lambda^2} < 0.$$

- **Uniform[0,  $\theta$ ]:**  $p(x|\theta) = \begin{cases} \theta^{-1}, & x \in [0, \theta] \\ 0, & \text{else} \end{cases}$ ,  $\theta = \theta$

*Hint:*  $p(x|\theta)$  can be written as  $p(x|\theta) = \theta^{-1} \mathcal{I}_{[0, \theta]}(x)$ , where  $\mathcal{I}_{[0, \theta]}(x) = \begin{cases} 1, & x \in [0, \theta] \\ 0, & \text{else} \end{cases}$ .

*Solution.* The Likelihood of the dataset  $X$  is given by:

$$L(X|\theta) = \prod_{i=1}^N p(x_i|\theta) = \frac{1}{\theta^N} \mathcal{I}_{[0,\theta]}(x_1) \mathcal{I}_{[0,\theta]}(x_2) \dots \mathcal{I}_{[0,\theta]}(x_n)$$

Let  $\theta^* = x^* = \max\{x_i\}$ .

1) If  $\theta < \theta^*$ , then  $\mathcal{I}_{[0,\theta]}(x^*) = 0 \rightarrow L(X|\theta) = 0$ .

2) If  $\theta > \theta^*$ , then  $0 < L(X|\theta) = \frac{1}{\theta^N} < \frac{1}{\theta^{*N}} = L(X|\theta^*)$ .

Therefore,  $\boxed{\hat{\theta} = \theta^* = \max\{x_i\}}$

- **Bernoulli( $\theta$ ):**  $x \in \{0, 1\}$ ,  $\begin{cases} p(x=0|\theta) = 1-\theta \\ p(x=1|\theta) = \theta \end{cases}$  or  $p(x|\theta) = \theta^x (1-\theta)^{1-x}$ ,  $\theta = \theta$

Note that this is a discrete distribution.

*Solution.* The Likelihood of the dataset  $X$  is given by:

$$L(X|\theta) = \prod_{i=1}^N p(x_i|\theta) = \theta^{\sum_{i=1}^N x_i} (1-\theta)^{\sum_{i=1}^N 1-x_i} = \theta^S (1-\theta)^{N-S}, \quad S = \sum_{i=1}^N x_i$$

$$\log L(X|\theta) = S \log \theta + (N-S) \log (1-\theta)$$

$$\frac{\partial}{\partial \theta} \log L(X|\theta) = \frac{S}{\theta} - \frac{N-S}{1-\theta} = 0 \quad \rightarrow \quad \boxed{\hat{\theta} = \frac{S}{N} = \frac{\sum_{i=1}^N x_i}{N}}$$

$$\frac{\partial^2}{\partial \theta^2} \log L(X|\theta) = -\frac{S}{\theta^2} - \frac{N-S}{(1-\theta)^2} < 0.$$

- **Multinomial( $\mathbf{p}, m$ ):**  $\mathbf{x} = (x_1, x_2, \dots, x_k)$ ,  $\mathbf{p} = (p_1, p_2, \dots, p_k)$ ,  $\sum_{j=1}^k x_j = m$ ,  $\sum_{j=1}^k p_j = 1$ ,  $x_j \in \overline{0 \dots n}$

$$p(\mathbf{x}|\mathbf{p}) = \frac{m!}{x_1! x_2! \dots x_k!} p_1^{x_1} p_2^{x_2} \dots p_k^{x_k}, \quad \theta = \mathbf{p}.$$

*Hint:* you need to use the method of Lagrange multipliers here. If you don't know what is it, you can skip this problem now and return to it after the lecture on SVMs.

*Solution.* The Likelihood of the dataset  $X$  is given by:

$$L(X|\mathbf{p}) = \prod_{i=1}^N p(\mathbf{x}_i|\mathbf{p}) = \prod_{i=1}^N \frac{m!}{x_{i1}! x_{i2}! \dots x_{ik}!} p_1^{x_{i1}} p_2^{x_{i2}} \dots p_k^{x_{ik}}$$

$$\log L(X|\mathbf{p}) = \sum_{i=1}^N \sum_{j=1}^k x_{ij} \log p_j + C(X, n)$$

Since we have an additional constraint  $\sum_{j=1}^k p_j = 1$ , we have to use the method of Lagrange multipliers.

The Lagrangian for this problem will take the form:

$$\begin{aligned}\mathcal{L}(X, \mathbf{p}, \lambda) &= \sum_{i=1}^N \sum_{j=1}^k x_{ij} \log p_j + \lambda \left( \sum_{j=1}^k p_j - 1 \right) \\ \frac{\partial}{\partial \lambda} \mathcal{L}(X, \mathbf{p}, \lambda) &= \sum_{j=1}^k p_j - 1 = 0 \quad \rightarrow \quad \sum_{j=1}^k p_j = 1 \\ \frac{\partial}{\partial p_j} \mathcal{L}(X, \mathbf{p}, \lambda) &= \frac{\sum_{i=1}^N x_{ij}}{p_j} + \lambda = 0 \quad \rightarrow \quad \sum_{i=1}^N x_{ij} = -\lambda p_j\end{aligned}$$

Lets sum both sides of the last equation over  $j$ . Since  $\sum_{j=1}^k x_{ij} = m$ , and  $\sum_{j=1}^k p_j = 1$ :

$$Nm = -\lambda \quad \rightarrow \quad \boxed{p_j = \frac{\sum_{i=1}^N x_{ij}}{-\lambda} = \frac{\sum_{i=1}^N x_{ij}}{Nm}}$$

### Problem 2 (Conjugate Priors):

In general, the Bayesian Learning and MAP are more computationally expensive than MLE, since we have to integrate over  $\theta$  in these methods:

$$p(x|X) = \int_{\theta} p(x|\theta) p(\theta|X) d\theta \quad (2)$$

$$p(\theta|X) = \frac{p(X|\theta)p(\theta)}{p(X)} = \frac{p(\theta) \prod p(x_i|\theta)}{\int_{\theta} p(\theta) \prod p(x_i|\theta) d\theta} \quad (3)$$

The integration in 3 can be avoided if we can identify the type of the distribution  $p(\theta|X) \sim p(X|\theta)p(\theta)$ , and then the normalization constant  $\alpha = \int_{\theta} p(X|\theta)p(\theta) d\theta$  can be calculated using its known parameters. In particular, if  $p(\theta|X)$  is from the same distribution family as  $p(\theta)$ , then the prior distribution  $p(\theta)$  is called conjugate to  $p(X|\theta)$ . For example, the normal distribution is conjugate to the normal one.

In this exercise, you have to show that the gamma-distribution  $p_{\alpha,\beta}(\lambda) = \frac{\beta^\alpha}{\Gamma(\alpha)} \lambda^{\alpha-1} e^{-\beta\lambda}$  is conjugate to the exponential distribution  $p(x|\lambda) = \lambda e^{-\lambda x}$ .

*Solution.* Suppose  $\lambda \sim \text{Gamma}(\alpha, \beta)$ ,  $X \sim \text{Exponential}(\lambda)$ , then:

$$p(\lambda|X) \sim p(X|\lambda) p_{\alpha,\beta}(\lambda) = \prod_{i=1}^N p(x_i|\lambda) p_{\alpha,\beta}(\lambda) = \lambda^N e^{-\lambda \sum_{i=1}^N x_i} \frac{\beta^\alpha}{\Gamma(\alpha)} \lambda^{\alpha-1} e^{-\beta\lambda} \sim \lambda^{(N+\alpha)-1} e^{-\lambda(\sum_{i=1}^N x_i + \beta)},$$

Thus,  $p(\lambda|X) \sim \text{Gamma}(\tilde{\alpha}, \tilde{\beta})$ , where  $\tilde{\alpha} = N + \alpha$ ,  $\tilde{\beta} = \sum_{i=1}^N x_i + \beta$ .

### Problem 3 (Bayesian linear regression):

$\epsilon$  has dimension  $n \times 1$ ,  $\mathbf{X}$  has dimension  $n \times p$ ,  $\beta$  has dimension  $p \times 1$ .

See Bayesian linear regression section of Gaussian Processes tutorial slides

$(\mathbf{X}^T \mathbf{X} + \sigma^2 \mathbf{\Lambda})^{-1}$  has dimension  $p \times p$ ,  $\mu_\beta$  has dimension  $p \times 1$ ,  $\Sigma_\beta$  has dimension  $p \times p$

See Bayesian linear regression section of Gaussian Processes tutorial slides

**Problem 4 (Kernels):**

Yes, Yes, Yes, No, Yes

$k(\mathbf{x}, \mathbf{y}) = \log(\|\mathbf{x} - \mathbf{y}\| + 1)$  is not a kernel because it is not positive semi-definite. Consider two data points  $x, y$  where we let  $\log(\|\mathbf{x} - \mathbf{y}\| + 1) = 1$ . Then we can check that the kernel matrix  $K = \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix}$  has eigenvalues  $1, -1$ . Because not all eigenvalues are positive, the matrix is not PSD. Alternatively, let  $z = [1 \ -1]^T$  and we see that  $z^T K z = -2 < 0$ .

**Problem 5 (Gaussian Processes):**

See Gaussian Processes tutorial slides.