

Series 2, Oct 24th, 2016 (MLE)

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

Please hand in solutions until Friday Nov 4th.

("*" -exercises are a little bit more difficult, but still useful and doable)

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 θ . Our goal is to find this θ , 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(λ):** $p(x|\lambda) = \lambda e^{-\lambda x}$, $x > 0$, $\lambda > 0$, $\theta = \lambda$

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

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

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

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

Note that this is a discrete distribution.

- **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 it is, you can skip this problem now and return to it after the lecture on SVMs.

Problem 2 (Conjugate Priors):

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

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

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

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

Problem 3 (Bayesian linear regression):

Consider the linear regression model $\mathbf{Y} = \mathbf{X}\beta + \epsilon$ with n observations and p predictor variables. We model $\epsilon \sim N(0, \sigma^2 \mathbf{I})$ and $\beta \sim N(0, \Lambda^{-1})$.

What is the dimensionality of ϵ ? Of \mathbf{X} ? Of β ?

Show that the posterior distribution $P(\beta|\mathbf{Y}, \mathbf{X}, \sigma^2, \Lambda)$ is normal with mean $\mu_\beta = (\mathbf{X}^T \mathbf{X} + \sigma^2 \Lambda)^{-1} \mathbf{X}^T \mathbf{Y}$ and covariance matrix $\Sigma_\beta = \sigma^2 (\mathbf{X}^T \mathbf{X} + \sigma^2 \Lambda)^{-1}$.

What is the dimensionality of $(\mathbf{X}^T \mathbf{X} + \sigma^2 \Lambda)^{-1}$? Of μ_β ? Of Σ_β ?

Show that Ridge is just a special case of Bayesian linear regression where $\Lambda = \frac{\lambda}{\sigma^2} \mathbf{I}$.

Problem 4 (Kernels):

Which of the following are kernels? If it's not a kernel, why not?

$$k(\mathbf{x}, \mathbf{y}) = C, C > 0$$

$$k(\mathbf{x}, \mathbf{y}) = \tanh(a\mathbf{x}^T \mathbf{y}) - b$$

$$k(\mathbf{x}, \mathbf{y}) = (1 + \mathbf{x}^T \mathbf{y})^p, p \in \mathbb{N}$$

$$k(\mathbf{x}, \mathbf{y}) = \log(\|\mathbf{x} - \mathbf{y}\| + 1)$$

$$k(\mathbf{x}, \mathbf{y}) = \exp\left(-\frac{\|\mathbf{x} - \mathbf{y}\|_2^2}{h^2}\right)$$

Problem 5 (Gaussian Processes):

Write down the equation describing a general Gaussian process. What kind of distribution does \mathbf{Y} have? How do you know? Write down the distribution of \mathbf{Y} and explicitly determine its parameters.

What kind of kernel function do we have if $\Lambda = \lambda \mathbf{I}_p$?

What happens when we add a new observation y_{n+1} ? Write down the joint distribution of all the data including the new observation. Explicitly determine the mean and covariance.

Suppose Θ is a set of hyperparameters that determine the parameters \mathbf{C}_n , which in turn directly affect the distribution of \mathbf{Y} . Describe how we would in general optimize the hyperparameters.