

Homework 5  
Mouse Linos Spyridon  
February 2020

## Exercise 1

(a) Erlang Distribution:  $p(x) = \theta^2 x \exp(-\theta x) u(x)$   
where  $u(x) = \begin{cases} 1 & x \geq 0 \\ 0 & x < 0 \end{cases}$

- Given a set of  $N$ ,  $X = \{x_1, x_2, \dots, x_N\}$   
independent, measurements/data points, we can calculate  
the likelihood as the product of per-point  
pdf. The joint pdf of  $X$  is:

$$p(X; \theta) = \prod_{i=1}^N p(x_i)$$

The log-likelihood is:

$$L(X; \theta) = \ln(p(X; \theta)) = \sum_{i=1}^N \ln(p(x_i; \theta))$$

substituting:

$$L(\theta) = \sum_{i=1}^N \ln(\theta^2 x_i \exp(-\theta x_i)) \quad \text{for } x_i \geq 0$$

$$L(\theta) = \sum_{i=1}^N [\ln(\theta^2 x_i) + \ln(\exp(-\theta x_i))] = \sum_{i=1}^N [\ln(\theta^2 x_i) + (-\theta x_i)]$$

$$L(\theta) = \sum_{i=1}^N [\ln(\theta^2) + \ln(x_i) - \theta x_i] = 2N \ln(\theta) + \sum_{i=1}^N \ln(x_i) - \theta \sum_{i=1}^N x_i$$

in order to find the MLE we need to find  $\theta'$  such as

$$\frac{dL(\theta)}{d\theta} \Big|_{\theta=\theta'} = 0. \quad \text{So: } \frac{dL(\theta)}{d\theta} = 0 \Rightarrow$$

$$\Rightarrow (2N \ln(\theta'))' + 0 + (-\theta' \sum_{i=1}^N x_i)' = 0 \Rightarrow \frac{2N}{\theta'} - \sum_{i=1}^N x_i = 0$$

$$\Rightarrow \theta_{ML} = \theta' = \frac{2N}{\sum_{i=1}^N x_i}$$

(b) Now  $N=5$  /  $x_1=2$  /  $x_2=2.2$  /  $x_3=2.7$  /  $x_4=2.4$  /  $x_5=2.6$

First of all let's plug our values in our equation:

$$\hat{\theta} = \theta_{ML} = \frac{2 \cdot N}{\sum_{i=1}^N x_i} = \frac{2 \cdot 5}{2 + 2.2 + 2.7 + 2.4 + 2.6} = \frac{10}{11.9} = 0.8403$$

The Erlang distribution that best "explains the data" is

$$p(x) = \theta_{ML}^2 x \exp(-\theta_{ML} x) \nu(x)$$

$$p(x) = (0.84)^2 x \exp(-0.84 x) \nu(x)$$

Now in order to calculate the mean of the  $X$  RV. we have

$$E[X] = \int_{-\infty}^{+\infty} x p(x) dx = \int_0^{+\infty} x \cdot \theta^2 x \exp(-\theta x) dx$$

$$= \theta^2 \int_0^{+\infty} x^2 \exp(-\theta x) dx = \theta^2 \frac{2}{\theta^3} = \frac{2}{\theta}$$

If we want to take  $\theta = \theta_{ML}$  it would be  $\therefore E[X] = \frac{2}{\theta_{ML}} = 2.38$

### Exercise 1

(c) The pdf value of  $x_1' = 2.1$  would be

$$p(x_1') = 0,14^2 x_1' \exp(-0,14 x_1') \cdot 1 = 0,253912$$

Same:  $x_2' = 2.3$

$$p(x_2') = 0,14^2 x_2' \exp(-0,14 x_2') \cdot 1 = 0,235017$$

Same:  $x_3' = 2.9$

$$p(x_3') = (0,14)^2 x_3' \exp(-0,14 x_3') \cdot 1 = 0,179067$$

The joint probability of  $\mathcal{X} = \{x_1', x_2', x_3'\}$  would be:

$$\text{jpdf}(\mathcal{X}) = \text{jpdf}(\mathcal{X}; \theta_{ML}) = \prod_{i=1}^3 p(x_i'; \theta_{ML}) = p(x_1') p(x_2') p(x_3')$$

$$= 0,01068.$$

The ~~log~~ log likelihood:

$$L(\theta_{ML}) \text{ on } \mathcal{X} = \ln(0,01068) = -4,538562.$$

## Exercise 2

(a) We ~~now~~ now have the same scenario <sup>Erlang</sup>  $p(x) = \theta^2 x \exp(-\theta x) u(x)$  but we know a prior-probability for  $\theta \sim N_{\text{prior}}(\theta_0, \sigma_0^2)$  with known  $\theta_0/\sigma_0$ .

The MAP estimate is simply  $\theta_{\text{MAP}} = \arg \max p(x|\theta) p(\theta)$

Again, assuming independent measurements / data points  $(x_i)$

we have:  $\bullet p(x|\theta) = \prod_{i=1}^N p(x_i|\theta) \quad (1)$

$$\bullet p(\theta) = N(\theta_0, \sigma_0^2) \quad (2)$$

jointly:  $\theta_{\text{MAP}} = \arg \max \left( \prod_{i=1}^N p(x_i|\theta) N(\theta_0, \sigma_0^2) \right)$ , however

~~$p(x|\theta) p(\theta)$~~  we can monotonically restate the problem as its log equivalent.

$$\theta_{\text{MAP}} = \arg \max_{\theta} \sum_{i=1}^N \ln(p(x_i|\theta)) + \ln(p(\theta)) \quad (3)$$

~~Let~~ From (3):

Let  $f(\theta)$  be  $f(\theta) = \sum_{i=1}^N \ln(p(x_i|\theta)) + \ln(p(\theta))$ , then

$$f(\theta) = \ln \left[ \frac{1}{\sqrt{2\pi} \sigma_0} \exp\left(-\frac{(\theta - \theta_0)^2}{2\sigma_0^2}\right) \right] + \sum \ln p(x_i|\theta) = \ln \left[ \frac{1}{\sqrt{2\pi} \sigma_0} \exp\left(-\frac{(\theta - \theta_0)^2}{2\sigma_0^2}\right) \right]$$

$$+ \sum_{i=1}^N \ln(x_i) - \theta \sum_{i=1}^N x_i + 2N \ln(\theta)$$

## Exercise 2

(a) Taking the derivative wrt  $\theta$ , we get:

$$\frac{\partial f(\theta)}{\partial \theta} = -\frac{(\theta - \theta_0)}{\sigma_0^2} + \frac{W}{\theta} - \sum_{i=1}^N x_i$$

We want  $\theta_{\text{MAP}} = \arg\max f(\theta) \Rightarrow \theta'$  such as  $\left. \frac{\partial f(\theta)}{\partial \theta} \right|_{\theta=\theta'} = 0$

$$\frac{\partial f(\theta)}{\partial \theta} = 0 \Rightarrow \theta^2 + \left( \sigma_0^2 \sum_{i=1}^N x_i - \theta_0 \right) \theta - 2N\sigma_0^2 = 0$$

That means that:

$$\theta_{\text{MAP}} = \theta' = \frac{\theta_0 - \sigma_0^2 \sum_{i=1}^N x_i \pm \sqrt{\left( \sigma_0^2 \sum_{i=1}^N x_i - \theta_0 \right)^2 + 8N\sigma_0^2}}{2}$$

However not both of these solutions are correct. ~~to limit down to 1 solution we~~ notice that the "-" solution:

$$\theta_{\text{MAP}}^{--} = \frac{1}{2} \left[ \theta_0 - \sigma_0^2 \sum x_i - \sqrt{\left( \sigma_0^2 \sum x_i - \theta_0 \right)^2 + 8N\sigma_0^2} \right]$$

Let  $z = \theta_0 - \sigma_0^2 \sum x_i$  then

$$\theta_{\text{MAP}}^{--} = \frac{1}{2} \left[ z - \sqrt{z^2 + 8N\sigma_0^2} \right] \leq \frac{1}{2} \left[ z - \sqrt{z^2} - \sqrt{8N\sigma_0^2} \right]$$

$$\theta_{\text{MAP}}^{--} \leq \frac{1}{2} \left[ z - z - \sqrt{8N\sigma_0^2} \right] \leq \frac{1}{2} \left[ -\sqrt{8N\sigma_0^2} \right] \leq 0$$

Meaning the second solution leads to a negative value thus we dismiss it and keep the "+" one.



(b) we have that  $\theta_{MAP} = \frac{1}{2} \left( \theta_0 - \sigma^2 \sum_{i=1}^N x_i \right) + \sqrt{\left( \theta_0 - \sigma^2 \sum_{i=1}^N x_i \right)^2 + 8N\sigma^2}$

(i) In case where  $N \rightarrow \infty$ :

Let's recall the equation that lead us to  $\theta_{MAP}$ :

$$\frac{\partial \ell(\theta)}{\partial \theta} \Big|_{\theta=\theta_{MAP}} = 0 \Rightarrow -\frac{(\theta_{MAP} - \theta_0)}{\sigma^2} + \frac{2N}{\theta_{MAP}} - \sum_{i=1}^N x_i = 0 \quad (1)$$

if  $N \rightarrow \infty$  then by dividing with  $N$ :

$$(1) \Rightarrow -\frac{(\theta_{MAP} - \theta_0)}{N\sigma^2} + \frac{2}{\theta_{MAP}} - \frac{1}{N} \sum_{i=1}^N x_i = 0$$

$$\Rightarrow (N \rightarrow \infty) \quad 0 + \frac{2}{\theta_{MAP}} - \frac{1}{N} \sum_{i=1}^N x_i = 0$$

$$\Rightarrow \theta_{MAP} = \frac{2N}{\sum x_i} = \theta_{MLE}$$

This derives from the fact that MAP is nothing more than a "tradeoff" between the current data distribution and the prior we are presented with. For large values of  $N$  the "data" distribution has greater effect and leads  $\theta_{MAP}$  to converge to  $\theta_{MLE}$ .

(ii) In case  $\sigma_0^2 \gg$ , meaning the confidence of our prior knowledge on the distribution of the data is low we have:

$$\frac{\partial \mathcal{L}}{\partial \theta_{\text{MAP}}} = \sum_{i=1}^N x_i - \frac{(\theta_{\text{MAP}} - \theta_0)}{\sigma_0^2} = 0$$

$\sigma_0^2 \gg$  gives us:

$$\frac{\partial \mathcal{L}}{\partial \theta_{\text{MAP}}} = \sum_{i=1}^N x_i \Rightarrow \theta_{\text{MAP}} = \frac{\partial \mathcal{L}}{\sum_{i=1}^N x_i} = \theta_{\text{ML}}$$

That means that both in  $N \rightarrow \infty$  and  $\sigma_0^2 \gg$ , our MAP estimation will converge to the ML estimation

(iii) In case  $\sigma_0^2 \ll$ , meaning that we are very confident about the prior distribution of our data we have:

$$\frac{\partial \mathcal{L}}{\partial \theta_{\text{MAP}}} = \sum_{i=1}^N x_i - \frac{(\theta_{\text{MAP}} - \theta_0)}{\sigma_0^2} = 0$$

$$\frac{\partial \mathcal{L}}{\partial \theta_{\text{MAP}}} \sigma_0^2 = \sigma_0^2 \sum_{i=1}^N x_i - (\theta_{\text{MAP}} - \theta_0) = 0$$

$\sigma_0^2 \ll$  gives us:

$$0 = 0 - (\theta_{\text{MAP}} - \theta_0) = 0 \\ \Rightarrow \theta_{\text{MAP}} = \theta_0$$

that means our MAP estimation is strongly confident to our prior knowledge and will be the same as the known  $\theta_0$ , the prior expected value.

(C) Let  $\theta_0 = 1.1 / \sigma_0^2 = 1$

$N=5$  with  $x_1=2 / x_2=2.2 / x_3=2.7 / x_4=2.4 / x_5=2.6$

As solved in question (a).

~~$$\theta_{MAP} = \hat{\theta} = \frac{1}{2} \left[ z + \sqrt{z^2 + 8N\sigma_0^2} \right]$$~~

$$\theta_{MAP} = \hat{\theta} = \frac{1}{2} \left[ z + \sqrt{z^2 + 8N\sigma_0^2} \right]$$

where  $z = \theta_0 - \sigma_0^2 \sum x_i = 1.1 - 1 \sum x_i = 1.1 - 11.9 = -10.8$

$$\theta_{MAP} = \hat{\theta} = \frac{1}{2} \left[ -10.8 + \sqrt{116.64 + 8.5} \right] = 0.8577$$

The erlang formula that best explains the data's

$$\begin{aligned} p(x) &= \theta_{MAP}^2 \times \exp(-\theta_{MAP}x) u(x) \\ &= (0.8577)^2 \times \exp(-0.8577x) \quad \forall x \geq 0 \end{aligned}$$

The mean of RV  $X$  is:

$$\begin{aligned} E[X] &= \int_{-\infty}^{+\infty} x p(x) dx = \int_0^{+\infty} x \theta_{MAP}^2 \times \exp(-\theta_{MAP}x) dx \\ &= \theta_{MAP}^2 \frac{2}{\theta_{MAP}^3} = \frac{2}{\theta_{MAP}} = 2.33 \end{aligned}$$



### Exercise 3

We have the model  $x = \mu + n$  where  $(x, \mu, n) \in \mathbb{R}$  and a set of measurements  $Y = \{x_1, x_2, x_3, \dots, x_N\}$

Assuming prior knowledge on  $\mu$ , we claim that is close to  $\mu_0$ . As a ridge regression problem we know that:

$$\min J(\mu) = \sum_{n=1}^N (x_n - \mu)^2 \quad \text{subject to } (\mu - \mu_0)^2 \leq \rho$$

Assuming Gaussian Prior as well as Gaussian Pdf we can reformulate the Ridge Regression Problem to a MAP/Bayesian Inference one, where the prior acts as a regularizer. Keeping that in mind, we have:

$Y = \{x_1, x_2, \dots, x_N\}$  a set of independent observations following an ~~known~~ Gaussian pdf, with ~~known~~ known mean  $\mu$ .

The prior knowledge is also of a Gaussian prior with mean  $\mu_0$  that lies close to  $\mu$ . If we didn't have any prior knowledge the equivalent of this problem would be the unconstrained.

where  $J(\mu) = \sum_{n=1}^N (x_n - \mu)^2$ . However now we know LS that  $(\mu - \mu_0)^2 \leq \rho$ , where  $\rho \ll 1$ . That gives us the constraint LS  $\Rightarrow$  Ridge Regression problem:

$$\min L(\mu) = J(\mu) + \lambda (\mu - \mu_0)^2 - \rho$$

Let  $\lambda$  be the Lagrange multiplier corresponding to the constraint. Then we look for  $\mu_{RR}^1$  such as

$$\arg \min_{\mu_{RR}} L(\mu)$$

$$\frac{\partial L(\mu)}{\partial \mu} = \frac{\partial}{\partial \mu} \sum_{n=1}^N (x_n - \mu)^2 + \frac{\partial}{\partial \mu} \lambda ((\mu - \mu_0)^2 - \rho)$$

$$\frac{\partial L(\mu)}{\partial \mu} = 2N\mu - 2\sum_{n=1}^N x_n + 2\lambda(\mu - \mu_0)$$

we want

$$\left. \frac{\partial L(\mu)}{\partial \mu} \right|_{\mu = \mu_{RR}} = 0 \Rightarrow 2N\mu_{RR} - 2\sum x + 2\lambda\mu_{RR} - 2\lambda\mu_0 = 0$$

$$\Rightarrow \mu_{RR}^1 = \frac{\sum x_i + \lambda\mu_0}{N + \lambda}$$

Now if we choose  $\lambda = \sigma_0^2 / \sigma^2$ , where  $\sigma_0^2$  the variance of the prior and  $\sigma^2$  the variance of the assumed pdf of the data, we would get

$$\mu_{RR}^1 = \mu_{MAP}^1 = \mu_{BI}^1$$

(d) The pdf value for  $x_1' = 2.1$  is

$$p(x_1') = \theta_{\text{MAP}}^2 x_1 \exp(-\theta_{\text{MAP}} x_1) = 0,255$$

$$\text{// } x_2' = 2.3:$$

$$p(x_2') = 0,235$$

$$\text{// } x_3' = 2.9:$$

$$p(x_3') = 0,177$$

The joint probability of  $X = \{x_1', x_2', x_3'\}$  is

$$p_{\text{pdf}}(X) = \prod_{i=1}^3 p(x_i'; \theta_{\text{MAP}}) = 0,01060$$

The log likelihood:

$$L(\theta_{\text{MAP}}) = -4,546267$$

(e) By comparing these results with those of ex1 we can observe a slight per point difference in the probabilities but not a significant change. That means the  $\theta_{\text{MAP}}$  is close to  $\theta_{\text{ML}}$  thus the "prior" ~~the~~ knowledge is not confident enough to shift  $\theta_{\text{MAP}}$  close to  $\theta_0 = 1.1$