

## Machine Learning Exercise Sheet 3

### Probabilistic Inference

---

Exercise sheets consist of two parts: In-class exercises and homework. The in-class exercises will be solved and discussed during the tutorial. The homework is for you to solve at home and further engage with the lecture content. There is no grade bonus and you do not have to upload any solutions. Note that the order of some exercises might have changed compared to last year's recordings.

---

#### In-class Exercises

Consider the probabilistic model

$$p(\mu \mid \alpha) = \mathcal{N}(\mu \mid 0, \alpha^{-1}) = \sqrt{\frac{\alpha}{2\pi}} \exp\left(-\frac{\alpha}{2}\mu^2\right)$$
$$p(x \mid \mu) = \mathcal{N}(x \mid \mu, 1) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{1}{2}(x - \mu)^2\right)$$

and a set of observations  $\mathcal{D} = \{x_1, \dots, x_N\}$  consisting of  $N$  samples  $x_i \in \mathbb{R}$ .

*Note:* We parametrize  $\mu \mid \alpha$  with the *precision* parameter  $\alpha = 1/\sigma^2$  instead of the usual variance  $\sigma^2$  because it leads to a nicer solution.

**Problem 1:** Derive the maximum likelihood estimate  $\mu_{\text{MLE}}$ . Show your work.

**Problem 2:** Derive the maximum a posteriori estimate  $\mu_{\text{MAP}}$ . Show your work.

**Problem 3:** Does there exist a prior distribution over  $\mu$  such that  $\mu_{\text{MLE}} = \mu_{\text{MAP}}$ ? Justify your answer.

**Problem 4:** Derive the posterior distribution  $p(\mu \mid \mathcal{D}, \alpha)$ . Show your work.

**Problem 5:** Derive the posterior predictive distribution  $p(x_{\text{new}} \mid \mathcal{D}, \alpha)$ . Show your work.

---

#### Homework

##### Optimizing Likelihoods: Monotonic Transforms

Usually we maximize the *log-likelihood*,  $\log p(x_1, \dots, x_n \mid \theta)$  instead of the likelihood. The next two problems provide a justification for this.

In the lecture, we encountered the likelihood maximization problem

$$\arg \max_{\theta \in [0,1]} \theta^t (1 - \theta)^h,$$

---

P1

$$\begin{aligned}
 \textcircled{1} \mu_{MLE} &= \arg \max_{\mu} \log p(D|\mu) \\
 &= \arg \max_{\mu} \log p(x_1, \dots, x_n | \mu) \\
 &= \arg \max_{\mu} \log \left( \prod_{i=1}^n p(x_i | \mu) \right) \\
 &= \arg \max_{\mu} \sum_{i=1}^n (\log p(x_i | \mu)) \\
 &= \arg \max_{\mu} \sum_{i=1}^n \left( \log \left( \frac{1}{\sqrt{2\pi}} \right) + \log \left( e^{-\frac{1}{2}(x_i - \mu)^2} \right) \right) \\
 &= \arg \max_{\mu} \sum_{i=1}^n \left( -\frac{1}{2} (x_i - \mu)^2 \right) \\
 &= \arg \max_{\mu} \sum_{i=1}^n \left( -\frac{1}{2} x_i^2 + \mu \sum_{i=1}^n x_i - \frac{n}{2} \mu^2 \right) \\
 &= \arg \max_{\mu} \sum_{i=1}^n \left( \mu \sum_{i=1}^n x_i - \frac{n}{2} \mu^2 \right) \\
 &= \arg \max_{\mu} \left( \mu \sum_{i=1}^n x_i - \frac{n}{2} \mu^2 \right) \\
 \frac{\partial \arg \max_{\mu} f(D|\mu)}{\partial \mu} &= \sum_{i=1}^n x_i - N\mu = 0 \\
 \mu_{MLE} &= \frac{1}{N} \sum_{i=1}^n x_i
 \end{aligned}$$

$$\begin{aligned}
 \textcircled{2} \mu_{MAP} &= \arg \max_{\mu} \log p(\mu | D, \alpha) \\
 &= \arg \max_{\mu} \log p(D|\mu) \cdot \underbrace{p(\mu|\alpha)}_{\text{prior}} \\
 &= \arg \max_{\mu} (\log p(D|\mu) + \log p(\mu|\alpha)) \\
 &= \arg \max_{\mu} \sum_{i=1}^n \left( -\frac{1}{2} (x_i - \mu)^2 \right) + \log p(\mu|\alpha) \\
 &= \arg \max_{\mu} \sum_{i=1}^n \left( -\frac{1}{2} (x_i - \mu)^2 \right) - \frac{\alpha}{2} \mu^2 \\
 \frac{\partial \arg \max_{\mu} p(\mu | D, \alpha)}{\partial \mu} &= \sum_{i=1}^n x_i - N\mu - \alpha\mu = 0 \\
 \mu &= \frac{1}{N + \alpha} \sum_{i=1}^n x_i
 \end{aligned}$$

$$\begin{aligned}
 \log p(\mu|\alpha) &= \log \left( \sqrt{\frac{\alpha}{2\pi}} \right) + \log \left( e^{-\frac{\alpha}{2}\mu^2} \right) \\
 &= \log \left( \sqrt{\frac{\alpha}{2\pi}} \right) - \frac{\alpha}{2} \mu^2
 \end{aligned}$$

$$\textcircled{3} \mu? \Rightarrow \mu_{MLE} = \mu_{MAP}$$

$$\begin{aligned}
 \mu_{MLE} &= \frac{1}{N} \sum_{i=1}^n x_i \\
 \mu_{MAP} &= \frac{1}{N + \alpha} \sum_{i=1}^n x_i
 \end{aligned}$$

When  $\alpha = 0$ ,  $\mu_{MLE} = \mu_{MAP}$   
 but  $\alpha = \frac{1}{\sigma^2}$   $\sigma^2$  never equal to  $\infty$   
 so impossible

① because  $\int_{-\infty}^{\infty} p(\mu|\alpha) d\mu = 1$

② ignore this limitation

We can use a uniform prior

$\alpha$  is very large

$\alpha \gg 0$   $\alpha \gg 0$   $\alpha \gg 0$   
 implies there is no assumption

$$\textcircled{4} p(\mu | D, \alpha) = \frac{p(D|\mu) p(\mu|\alpha)}{p(D)}$$

$$\begin{aligned}
 &\propto p(D|\mu) p(\mu|\alpha) \\
 &= \prod_{i=1}^n \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{1}{2}(x_i - \mu)^2\right) \cdot \sqrt{\frac{\alpha}{2\pi}} \exp\left(-\frac{\alpha}{2}\mu^2\right) \\
 &= \exp\left(-\frac{1}{2} \sum_{i=1}^n (x_i - \mu)^2 - \frac{\alpha}{2} \mu^2\right) \\
 &= \exp\left(-\frac{N + \alpha}{2} \mu^2 + \mu \sum_{i=1}^n x_i\right)
 \end{aligned}$$

$$\mathcal{N}(u|m, \beta^{-1}) = \sqrt{\frac{\beta}{2\pi}} \exp\left(-\frac{\beta}{2}(u - m)^2\right)$$

$$\propto \exp\left(-\frac{\beta}{2} u^2 + \beta m u\right)$$

$$\star p(\mu | x_1, \dots, x_n) = \mathcal{N}\left(m, \frac{1}{\beta}\right) = \mathcal{N}\left(\frac{1}{n + \alpha} \sum_{i=1}^n x_i, \frac{1}{n + \alpha}\right)$$

$$\beta = N + \alpha$$

$$\beta m = \sum_{i=1}^n x_i$$

$$m = \frac{1}{\beta} \sum_{i=1}^n x_i \Rightarrow m = \frac{1}{N + \alpha} \sum_{i=1}^n x_i$$

$$\textcircled{5} p(x | D, \alpha) = \int_{-\infty}^{\infty} p(x|\mu) \overset{D, \alpha}{p(\mu | D, \alpha)} d\mu$$

$$p(x|\mu) p(\mu | D, \alpha) d\mu$$

$$\textcircled{m} p(x) = \int p(x|\mu) p(\mu) d\mu$$

$$\begin{aligned}
 \downarrow \\
 \textcircled{a} \mu &\sim \mathcal{N}(m, \beta^{-1}) \\
 x &\sim \mathcal{N}(\mu, 1) \Rightarrow x = \mu + y \sim \mathcal{N}(m + y, 1) \\
 \mu &\sim \mathcal{N}(m, \beta^{-1}) \\
 y &\sim \mathcal{N}(0, 1)
 \end{aligned}$$

where  $t$  and  $h$  denoted the number of tails and heads in a sequence of coin tosses, respectively.

**Problem 6:** Compute the first and second derivative of this likelihood w.r.t.  $\theta$ . Then compute first and second derivative of the log-likelihood  $\log \theta^t (1 - \theta)^h$ .

**Problem 7:** Show that for *any* differentiable, positive function  $f(\theta)$  every local maximum of  $\log f(\theta)$  is also a local maximum of  $f(\theta)$ . Considering this and the previous exercise, what is your conclusion?

## Properties of MLE and MAP

**Problem 8:** Consider a Bernoulli random variable  $X$  and suppose we have observed  $m$  occurrences of  $X = 1$  and  $l$  occurrences of  $X = 0$  in a sequence of  $N = m + l$  Bernoulli experiments. We are only interested in the number of occurrences of  $X = 1$ —we will model this with a Binomial distribution with parameter  $\theta$ . A prior distribution for  $\theta$  is given by the Beta distribution with parameters  $a, b$ . Show that the posterior *mean* value  $\mathbb{E}[\theta \mid \mathcal{D}]$  (not the MAP estimate) of  $\theta$  lies between the prior mean of  $\theta$  and the maximum likelihood estimate for  $\theta$ .

To do this, show that the posterior mean can be written as  $\lambda$  times the prior mean plus  $(1 - \lambda)$  times the maximum likelihood estimate, with  $0 \leq \lambda \leq 1$ . This illustrates the concept of the posterior mean being a compromise between the prior distribution and the maximum likelihood solution.

The probability mass function of the Binomial distribution for some  $m \in \{0, 1, \dots, N\}$  is

$$p(x = m \mid N, \theta) = \binom{N}{m} \theta^m (1 - \theta)^{N-m}.$$

*Hint:* Identify the posterior distribution. You may then look up the mean rather than computing it.

**Problem 9:** Consider the following probabilistic model

$$\begin{aligned} p(\lambda \mid a, b) &= \text{Gamma}(\lambda \mid a, b) = \frac{b^a}{\Gamma(a)} \lambda^{a-1} \exp(-b\lambda) \\ p(x \mid \lambda) &= \text{Poisson}(x \mid \lambda) = \frac{\lambda^x \exp(-\lambda)}{x!} \end{aligned}$$

where  $a \in (1, \infty)$  and  $b \in (0, \infty)$ . We have observed a single data point  $x \in \mathbb{N}$ . Derive the maximum a posteriori (MAP) estimate of the parameter  $\lambda$  for the above probabilistic model. Show your work.

## Programming Task

**Problem 10:** Download the notebook `exercise_03_notebook.ipynb` from Moodle. Fill in the missing code and follow the instructions in the notebook.

*Note:* We suggest that you use Anaconda for installing Python and Jupyter, as well as for managing packages. We recommend that you use Python 3.

For more information on Jupyter notebooks, consult the Jupyter documentation.

P6

$$\begin{aligned} ① \quad \frac{d(\theta^t(1-\theta)^n)}{d\theta} &= t\theta^{t-1}(1-\theta)^n + \theta^t n(1-\theta)^{n-1} \\ ② \quad \frac{d^2}{d\theta^2} &= t(t-1)\theta^{t-2}(1-\theta)^n + t\theta^{t-1}n(1-\theta)^{n-1} \\ &\quad + t\theta^{t-1}n(1-\theta)^{n-1} + \theta^t n(n-1)(1-\theta)^{n-2} \\ ③ \quad f(\theta) &= \log \theta^t (1-\theta)^n = t \log \theta + n \log(1-\theta) \\ \frac{d f(\theta)}{d \theta} &= \frac{t}{\theta} - \frac{n}{1-\theta} \\ \frac{d^2 f(\theta)}{d \theta^2} &= -\frac{t}{\theta^2} - \frac{n}{(1-\theta)^2} \end{aligned}$$

P7  $g(x) = f(x) \quad h(x) = \log f(x)$

$$\begin{aligned} \frac{d g(x)}{d x} &= f'(x) = 0 \text{ at } f(x) \text{ max.} \\ \frac{d h(x)}{d x} &= \frac{1}{f(x)} = 0 \text{ at } f(x) \rightarrow \pm \infty \end{aligned}$$

$\hat{\theta}^*$  is maximum of  $g(x)$   
 $\Rightarrow g(\hat{\theta}^*) \geq g(\theta)$   
 $\exp g(\hat{\theta}^*) = f(\hat{\theta}^*)$   
 $\uparrow$   
 also maximum  
 For log and no log, maximum (and minimum is equal), so it's easy to use log for calculation

P7  $X \sim \text{Ber}$   $\begin{cases} X=1, & m \\ X=0, & L \end{cases} \quad m+L=N$

prior Distribution  $p(\theta) \sim \text{Beta}(a, b)$

posterior mean  $E[\theta|x]$  lies between mean of  $\theta$  and  $\theta_{MLE}$

$$p(\theta) = \text{Beta}(\theta|a, b) = \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \theta^{a-1} (1-\theta)^{b-1}$$

$$p(X=m|\theta) = \text{Binomial} = \binom{N}{m} \theta^m (1-\theta)^{N-m}$$

$$p(\theta|x) = \frac{p(x|\theta) \cdot p(\theta)}{p(x)} \propto \theta^m (1-\theta)^{N-m} \cdot \theta^{a-1} (1-\theta)^{b-1}$$

$$= \theta^{m+a-1} (1-\theta)^{N-m+b-1}$$

$$\Rightarrow p(\theta|x) \sim \text{Beta}(\theta|m+a, N-m+b)$$

$$\Rightarrow E(\text{Beta}) = \frac{m+a}{N+m+b} = \frac{m}{N+m+b} + \frac{a}{N+m+b}$$

$$\text{mean of } p(\theta) = \frac{a}{a+b} \quad \theta_{MLE} = \frac{m}{N}$$

$$\Rightarrow \lambda \cdot \frac{a}{a+b} + (1-\lambda) \frac{m}{N} = \frac{m+a}{N+m+b}$$

$$\lambda = \frac{a+b}{N+m+b} \quad 1-\lambda = \frac{N}{N+m+b}$$

P8  $p(\lambda|x, a, b) = \underset{\lambda}{\arg \max} \frac{p(x|\lambda) p(\lambda|a, b)}{p(x)} \propto \underset{\lambda}{\arg \max} p(x|\lambda) p(\lambda|a, b)$

$$= \underset{\lambda}{\arg \max} (\log p(x|\lambda) + \log p(\lambda|a, b))$$

$$= \underset{\lambda}{\arg \max} \log \lambda^x + (-\lambda) + \log \frac{b^a}{\Gamma(b)} + \log \lambda^{a-1} - b\lambda$$

$$= \underset{\lambda}{\arg \max} x \log \lambda - \lambda + (a-1) \log \lambda - b\lambda + C$$

$$\frac{\partial}{\partial \lambda} = \frac{x+a-1}{\lambda} - (b+1) = 0$$

$$\lambda_{MAP} = \frac{x+a-1}{b+1}$$