

# Machine Learning: Homework #3

Due on November 13, 2017 at 09:55am

*Professor Dr. Stephan Guennemann*

Marc Meissner

## Problem 1

Usually one considers the log likelihood  $\log p(x_1, \dots, x_n | \theta)$ . The next problems justifies this. In the lecture, we encountered the likelihood maximization problem

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

, where  $t$  and  $h$  denoted the number of tails and heads in a sequence of coin tosses, respectively. Compute the first and second derivative of this likelihood w.r.t.  $\theta$ . Then compute first and second derivative of the log likelihood  $\log \theta^t (1 - \theta)^h$ .

### Solution

First equation:

$$f(\theta) = \theta^t (1 - \theta)^h$$

$$\begin{aligned} \frac{\partial f(\theta)}{\partial \theta} &= t\theta^{t-1}(1 - \theta)^h - \theta^t h(1 - \theta)^{h-1} \\ &= \theta^{t-1}(1 - \theta)^{h-1}(t(1 - \theta) - h\theta) \end{aligned}$$

$$\begin{aligned} \frac{\partial^2 f(\theta)}{\partial \theta^2} &= \dots \\ &= [\theta^{t-2}(1 - \theta)^{h-2} \cdot ((t-1)(1 - \theta) - (h-1)\theta) \cdot (t(1 - \theta) - h\theta)] - \theta^{t-1}(1 - \theta)^{h-1}(t + h) \end{aligned}$$

Second equation:

$$g(\theta) = \ln[\theta^t (1 - \theta)^h] = t \ln \theta + h \ln(1 - \theta)$$

$$\frac{\partial g(\theta)}{\partial \theta} = \frac{t}{\theta} - \frac{h}{1 - \theta}$$

$$\frac{\partial^2 g(\theta)}{\partial \theta^2} = -\frac{t}{\theta^2} - \frac{h}{(1 - \theta)^2}$$

## Problem 2

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

### Solution

Given  $f(\theta)$  with local maximum  $\theta_{\max}$ . Then, for a certain  $\epsilon$ , it locally holds that:

$$f(\theta_{\max}) \geq f(\theta) \text{ for any } \theta \in [\theta_{\max} - \epsilon; \theta_{\max} + \epsilon]$$

Since  $g(\theta) = \ln f(\theta)$  is a monotonic function, the following property holds:

$$x_2 - x_1 \implies \ln(x_2) > \ln(x_1)$$

$$g(\theta_{\max}) = \ln f(\theta_{\max}) \geq \ln f(\theta) = g(\theta)$$

Thus, one can apply log-likelihood and preserve the position of the maximum while severely reducing the complexity of the solution.

### Problem 3

Show that  $\theta_{\text{MLE}}$  can be interpreted as a special case of  $\theta_{\text{MAP}}$  in the sense that there always exists a prior  $p(\theta)$  such that  $\theta_{\text{MLE}} = \theta_{\text{MAP}}$ .

#### Solution

Any constant prior (uniform distribution) should preserve the position of the maximum, since it just scales the distribution. Given  $p(\theta) = c$  and the definitions from the slides:

$$\theta_{\text{MAP}} = \arg \max_{\theta} p(D|\theta)p(\theta) = \arg \max_{\theta} p(D|\theta)c = \arg \max_{\theta} p(D|\theta) = \theta_{\text{MLE}}$$

### Problem 4

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  $E[\theta|D]$  (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  $\theta \leq \lambda \leq 1$ . 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

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

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

**Solution**

Solution.

**Problem 5**

Let  $X$  be Poisson distributed. Again, for  $n$  i.i.d. samples from  $X$ , determine the maximum likelihood estimate for  $\lambda$ . In class we also talked about avoiding overfitting of parameters via prior information. Compute the posterior distribution over  $\lambda$ , assuming a  $\text{Gamma}(\alpha, \beta)$  prior for it. Compute the MAP for  $\lambda$  under this prior. Show your work.

**Solution**

Solution.