

Machine Learning Exercise Sheet 3

Probabilistic Inference

We changed the submission settings: From now on only one member of a group must upload a submission on Moodle. Please make sure that the name and matriculation number of each group member is on the front page of the submission.

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

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

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

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

To solve this, we need to apply chain and product rule.

$$\frac{d}{d\theta} \theta^t (1 - \theta)^h = \theta^{t-1} (1 - \theta)^{h-1} ((1 - \theta)t - \theta h)$$

$$\frac{d^2}{d\theta^2} \theta^t (1 - \theta)^h = \theta^{t-2} (1 - \theta)^{h-2} \cdot ((1 - \theta)(t - 1) - \theta(h - 1)) \cdot ((1 - \theta)t - \theta h) - \theta^{t-1} (1 - \theta)^{h-1} (t + h)$$

Observe that the first factor of the first derivative is structurally the same as the original function. This can be used to ease some of the pain of these computations.

The product rule breaks our necks, which quickly renders the expressions long and confusing.

The logarithm decomposes the product into a sum. We only need to apply the chain rule on each of

Upload a single PDF file with your homework solution to Moodle by 03.11.2019, 23:59 CET. We recommend to typeset your solution (using L^AT_EX or Word), but handwritten solutions are also accepted. If your handwritten solution is illegible, it won't be graded and you waive your right to dispute that.

the summands. Do not forget the change of sign from taking the derivative of $1 - \theta$.

$$\begin{aligned} g(\theta) &:= \log \theta^t (1 - \theta)^h = t \log \theta + h \log(1 - \theta) \\ \frac{d}{d\theta} g(\theta) &= \frac{t}{\theta} - \frac{h}{1 - \theta} \\ \frac{d^2}{d\theta^2} g(\theta) &= - \left(\frac{t}{\theta^2} + \frac{h}{(1 - \theta)^2} \right) \end{aligned}$$

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?

Let θ^* be an arbitrary local maximum of $g(\theta) = \log f(\theta)$, i.e., for any θ in a small neighborhood of θ^* , we have that $g(\theta^*) \geq g(\theta)$. Since \exp is a monotonic transform, we also have

$$f(\theta^*) = \exp(g(\theta^*)) \geq \exp(g(\theta)) = f(\theta).$$

Hence, θ^* is also a maximum of f .

With the help of the previous exercise, we can now safely apply the logarithm and any maximum or minimum remains preserved (its position only, of course). Moreover, we have seen that the logarithmic domain can greatly simplify the computational effort to arrive at critical points. This also leads to improved numerical stability. Thus, it is often worth switching to the log domain when analyzing likelihoods.

Notice that the exercise left out a part of the argument: We only showed that a maximum of the log likelihood is also a maximum of the likelihood. We would still need to prove that taking the logarithm does not eliminate maxima of the likelihood. This is done by showing that monotonic transforms preserve critical points and observing that the logarithm is monotonic—we will not do this here.

Properties of MLE and MAP

Problem 3: You model a coin flip f as a Bernoulli distribution with a parameter θ

$$p(f \mid \theta) = \text{Bern}(f \mid \theta) = \theta^{\mathbb{I}[f=T]} (1 - \theta)^{\mathbb{I}[f=H]}.$$

That is, the probability of landing tails (T) is θ , and probability of heads (H) is $(1 - \theta)$ respectively.

Your prior on θ is a $\text{Beta}(6, 4)$ distribution

$$p(\theta) = \text{Beta}(\theta \mid 6, 4).$$

You observe $(M + N)$ coin flips, out of which M are tails and N are heads. After you do maximum a posteriori estimation of θ , you obtain the result $\theta_{\text{MAP}} = 0.75$.

Name any possible values of M and N that can lead to such result. Show your work.

Because Beta is a conjugate prior to binomial likelihood, the posterior is also a Beta

$$p(\theta \mid \mathcal{D}) = \text{Beta}(\theta \mid 6 + M, 4 + N)$$

The mode of a $\text{Beta}(a, b)$ distribution is

$$\theta_{MAP} = \frac{a - 1}{a + b - 2}$$

Which in our case means

$$\theta_{MAP} = \frac{6 + M - 1}{6 + M + 4 + N - 2} = \frac{M + 5}{M + N + 8}$$

so it must hold that

$$\begin{aligned} \frac{M + 5}{M + N + 8} &\stackrel{!}{=} \frac{3}{4} \\ 4(5 + M) &= 3(M + N + 8) \\ M &= 3N + 4 \end{aligned}$$

any pair (M, N) that satisfies this equation will produce the required result.

For example $N = 1$ and $M = 7$ works.

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 θ . A prior distribution for θ 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 θ lies between the prior mean of θ and the maximum likelihood estimate for θ .

To do this, show that the posterior mean can be written as λ 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.

For the given observations (and a, b hyper parameters for the prior beta distribution), we see that, just like in the lecture, the posterior is $\text{Beta}(m + a, l + b)$ -distributed. One gets the *expected* posterior mean for θ :

$$\mathbb{E}[\theta \mid \mathcal{D}] = \frac{m + a}{m + l + a + b} = \frac{m}{m + l + a + b} + \frac{a}{m + l + a + b}$$

But:

$$\frac{m}{m + l + a + b} = \underbrace{\frac{m + l}{m + l + a + b}}_{1-\lambda} \cdot \frac{m}{m + l}$$

Upload a single PDF file with your homework solution to Moodle by 03.11.2019, 23:59 CET. We recommend to typeset your solution (using L^AT_EX or Word), but handwritten solutions are also accepted. If your handwritten solution is illegible, it won't be graded and you waive your right to dispute that.

and

$$\frac{a}{m+l+a+b} = \frac{a+b}{\underbrace{m+l+a+b}_{\lambda}} \cdot \frac{a}{a+b}$$

producing what was asked, because $\frac{m}{m+l}$ is the maximum likelihood estimate and $\frac{a}{a+b}$ is the prior mean value of θ .

Programming Task

Problem 5: Load the notebook `exercise_03_notebook.ipynb` from Piazza. Fill in the missing code and run the notebook. Export (download) the evaluated notebook as PDF and add it to your submission.

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. Instructions for converting the Jupyter notebooks to PDF are provided within the notebook.

In-class Exercises

Consider the following probabilistic model

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

Note that here we don't parametrize the variance σ^2 , but rather specify the *precision* parameter $\alpha = 1/\sigma^2$.

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

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

Problem 6: Derive the maximum likelihood estimate μ_{MLE} . Show your work.

Our goal is to find

$$\begin{aligned}\mu_{\text{MLE}} &= \arg \max_{\mu \in \mathbb{R}} p(\mathcal{D} \mid \mu) \\ &= \arg \max_{\mu \in \mathbb{R}} \log p(\mathcal{D} \mid \mu)\end{aligned}$$

We solve this problem in two steps:

1. Write down & simplify the expression for $\log p(\mathcal{D} \mid \mu)$.
2. Solve $\frac{\partial}{\partial \mu} \log p(\mathcal{D} \mid \mu) \stackrel{!}{=} 0$ for μ .

$$\begin{aligned}
\log p(\mathcal{D} \mid \mu) &= \log p(x_1, \dots, x_N \mid \mu) \\
&= \log \left(\prod_{i=1}^N p(x_i \mid \mu) \right) && \text{iid assumption} \\
&= \sum_{i=1}^N \log p(x_i \mid \mu) \\
&= \sum_{i=1}^N \left[\log \left(\frac{1}{\sqrt{2\pi}} \right) + \log \left(\exp \left(-\frac{1}{2}(x_i - \mu)^2 \right) \right) \right] \\
&= \sum_{i=1}^N \left[-\frac{1}{2}(x_i - \mu)^2 \right] + \text{const.} \\
&= -\frac{1}{2} \sum_{i=1}^N (x_i^2 - 2x_i\mu + \mu^2) + \text{const.} \\
&= \left[-\frac{1}{2} \sum_{i=1}^N x_i^2 \right] + \left[\sum_{i=1}^N x_i\mu \right] - \left[\frac{1}{2} \sum_{i=1}^N \mu^2 \right] + \text{const.} \\
&= \mu \sum_{i=1}^N x_i - \frac{N}{2} \mu^2 + \text{const.}
\end{aligned}$$

Now compute the derivative and set it to zero.

$$\begin{aligned}
\frac{\partial}{\partial \mu} \log p(\mathcal{D} \mid \mu) &= \frac{\partial}{\partial \mu} \left(\mu \sum_{i=1}^N x_i - \frac{N}{2} \mu^2 + \text{const.} \right) \\
&= \sum_{i=1}^N x_i - N\mu \stackrel{!}{=} 0
\end{aligned}$$

Solving for μ we obtain

$$\mu_{\text{MLE}} = \frac{1}{N} \sum_{i=1}^N x_i$$

That is, μ_{MLE} is just the average of the datapoints.

Problem 7: Derive the maximum a posteriori estimate μ_{MAP} . Show your work.

Our goal is to find

$$\begin{aligned}\mu_{\text{MAP}} &= \arg \max_{\mu \in \mathbb{R}} p(\mu \mid \mathcal{D}, \alpha) \\ &= \arg \max_{\mu \in \mathbb{R}} \log p(\mu \mid \mathcal{D}, \alpha) \\ &= \arg \max_{\mu \in \mathbb{R}} [\log p(\mathcal{D} \mid \mu) + \log p(\mu \mid \alpha)]\end{aligned}$$

We solve this problem in two steps:

1. Write down & simplify the expression for $\log p(\mathcal{D} \mid \mu) + \log p(\mu \mid \alpha)$.
2. Solve $\frac{\partial}{\partial \mu} (\log p(\mathcal{D} \mid \mu) + \log p(\mu \mid \alpha)) \stackrel{!}{=} 0$ for μ .

$$\begin{aligned}\log p(\mu \mid \alpha) &= \log \left(\sqrt{\frac{\alpha}{2\pi}} \right) + \log \left(\exp \left(-\frac{\alpha}{2} \mu^2 \right) \right) \\ &= -\frac{\alpha}{2} \mu^2 + \text{const.}\end{aligned}$$

From the previous task, we know that

$$\log p(\mathcal{D} \mid \mu) = \mu \sum_{i=1}^N x_i - \frac{N}{2} \mu^2 + \text{const.}$$

Therefore, we get

$$\log p(\mathcal{D} \mid \mu) + \log p(\mu \mid \alpha) = \mu \sum_{i=1}^N x_i - \frac{N}{2} \mu^2 - \frac{\alpha}{2} \mu^2 + \text{const.}$$

Now compute the derivative and set it to zero.

$$\begin{aligned}\frac{\partial}{\partial \mu} (\log p(\mathcal{D} \mid \mu) + \log p(\mu \mid \alpha)) &= \frac{\partial}{\partial \mu} \left(\mu \sum_{i=1}^N x_i - \frac{N}{2} \mu^2 - \frac{\alpha}{2} \mu^2 + \text{const.} \right) \\ &= \sum_{i=1}^N x_i - N\mu - \alpha\mu \stackrel{!}{=} 0\end{aligned}$$

Solving for μ we obtain

$$\mu_{\text{MAP}} = \frac{1}{N + \alpha} \sum_{i=1}^N x_i$$

By comparing this to μ_{MLE} , we can understand the effect of a 0-mean Gaussian prior on our estimate of μ . Since $\alpha > 0$, we see that μ_{MAP} is always closer to zero than μ_{MLE} .

Problem 8: Does there exist a prior distribution over μ such that $\mu_{\text{MLE}} = \mu_{\text{MAP}}$? Justify your answer.

Let's compare the expressions for μ_{MLE} and μ_{MAP}

$$\mu_{\text{MLE}} = \frac{1}{N} \sum_{i=1}^N x_i \qquad \mu_{\text{MAP}} = \frac{1}{N + \alpha} \sum_{i=1}^N x_i$$

As α approaches zero ($\alpha \rightarrow 0$), μ_{MAP} gets closer to μ_{MLE} . As the *precision* of the prior distribution *decreases*, its *variance increases*. The prior distribution is getting more and more flat, thus being less informative and having a smaller effect on the posterior.

If we could set $\alpha = 0$, we would have a uniform prior on μ , and thus $\mu_{\text{MLE}} = \mu_{\text{MAP}}$. However, technically, we are not allowed to do that — since the distribution $p(\mu \mid \alpha)$ is defined over all of \mathbb{R} , it has to integrate to one ($\int_{-\infty}^{\infty} p(\mu \mid \alpha) d\mu = 1$).

We can ignore this restriction and assume that we have a uniform prior over μ . Such prior would be called *improper*. While in many cases it's fine to use an improper prior, it might lead to subtle problems in certain situations.

Problem 9: Derive the posterior distribution $p(\mu \mid \mathcal{D}, \alpha)$. Show your work.

We obtain the posterior distribution using Bayes formula

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

We know that the posterior distribution has to integrate to 1, but we don't know the normalizing constant. However, we know that it's proportional to $\exp(a\mu^2 + b\mu)$. This looks very similar to a normal distribution — we have a quadratic form inside the exponential.

How can we use this fact? Consider a normal distribution over μ with mean m and precision β

$$\begin{aligned}\mathcal{N}(\mu \mid m, \beta^{-1}) &= \sqrt{\frac{\beta}{2\pi}} \exp\left(-\frac{\beta}{2}(\mu - m)^2\right) \\ &\propto \exp\left(-\frac{\beta}{2}\mu^2 + \beta m\mu\right)\end{aligned}\quad (2)$$

If we find β and m such that Equations 1 and 2 are equal, we will know that our posterior $p(\mu \mid \mathcal{D}, \alpha)$ is a normal distribution with mean m and precision β .

First we observe that

$$\beta = N + \alpha$$

Now we need to find m such that

$$\begin{aligned}\beta m &= \sum_{i=1}^N x_i \\ m &= \frac{1}{\beta} \sum_{i=1}^N x_i \\ m &= \frac{1}{N + \alpha} \sum_{i=1}^N x_i\end{aligned}$$

Putting everything together we see that

$$p(\mu \mid \mathcal{D}, \alpha) = \mathcal{N}\left(\mu \mid \frac{1}{N + \alpha} \sum_{i=1}^N x_i, (N + \alpha)^{-1}\right)$$

Since the posterior is a normal distribution, its mean coincides with its mode — this means that $\mathbb{E}_{p(\mu \mid \mathcal{D}, \alpha)}[\mu] = \mu_{\text{MAP}}$. We can see that this is indeed the case, which is a good sanity check.

Problem 10: Derive the posterior predictive distribution $p(x_{\text{new}} \mid \mathcal{D}, \alpha)$. Show your work.

The posterior over μ is $p(\mu \mid \mathcal{D}, \alpha) = \mathcal{N}(\mu \mid m, \beta^{-1})$. Our goal is to find the *posterior predictive* distribution over the next sample $p(x_{\text{new}} \mid \mathcal{D}, \alpha)$. For brevity, we will drop the *new* subscript.

From the lecture we remember that thanks to the conditional independence assumption the posterior predictive is

$$p(x \mid \mathcal{D}, \alpha) = \int_{-\infty}^{\infty} p(x \mid \mu) p(\mu \mid \mathcal{D}, \alpha) d\mu$$

There are two (equivalent) ways to approach this problem.

Approach 1. Basically, we are modeling the following process

- We draw μ from the posterior distribution $\mu \sim \mathcal{N}(m, \beta^{-1})$.
- We draw x from the conditional distribution $x \sim \mathcal{N}(\mu, 1)$.

This process is identical to the following procedure

- We draw μ from the posterior distribution $\mu \sim \mathcal{N}(m, \beta^{-1})$.
- We draw y from the standard normal distribution $y \sim \mathcal{N}(0, 1)$.
- We calculate x as $\mu + y$.

Clearly, x is a sum of two *independent* normally distributed random variables. Hence, x also follows a normal distribution with mean $m + 0$ and precision $(\beta^{-1} + 1)^{-1}$.

$$p(x \mid \mathcal{D}, \alpha) = \mathcal{N}(x \mid m, \beta^{-1} + 1)$$

where m and β were computed in the previous problem.

Approach 2. We can directly look at the integral

$$\begin{aligned} p(x \mid \mathcal{D}, \alpha) &= \int_{-\infty}^{\infty} p(x \mid \mu) p(\mu \mid \mathcal{D}, \alpha) d\mu \\ &= \int_{-\infty}^{\infty} \mathcal{N}(x \mid \mu, 1) \mathcal{N}(\mu \mid m, \beta^{-1}) d\mu \\ &= \int_{-\infty}^{\infty} \mathcal{N}(x - \mu \mid 0, 1) \mathcal{N}(\mu \mid m, \beta^{-1}) d\mu \end{aligned}$$

This is a convolution of two Gaussian densities — the result is a Gaussian density as well

$$= \mathcal{N}(x \mid m, \beta^{-1} + 1)$$

You can find the proof on Wikipedia https://en.wikipedia.org/wiki/Sum_of_normally_distributed_random_variables#Proof_using_convolutions.

The two approaches are effectively identical, and both rely on two facts:

1. μ is the location parameter of the normal distribution. That means that if $p(x) = \mathcal{N}(x \mid \mu, \sigma^2)$ and $y = x + a$ (for a fixed $a \in \mathbb{R}$), then $p(y) = \mathcal{N}(y \mid \mu + a, \sigma^2)$.
2. the sum of two normally distributed RVs is a normally distributed RV