

# Machine Learning 1 - Homework 2

Selene Baez Santamaria

## 1 MAP solution for Linear Regression

In class we solved for the maximum likelihood estimator for linear regression with polynomial basis functions. In this exercise you will solve for the *maximum a posterior* (MAP) solution:  $\mathbf{w}_{\text{MAP}} = (\Phi^T \Phi + \lambda \mathbf{I})^{-1} \Phi^T \mathbf{t}$ . For this problem we assume  $N$  training vectors  $\{\mathbf{x}_n\}_{n=1}^N$ , each of which is mapped using basis functions to  $\phi_n$ . In the training set, the data come in input-output pairs, i.e.  $\{\mathbf{x}_n, t_n\}$ . We assume that one of the basis functions is the constant 1, and there are  $M - 1$  other basis function in  $\phi_n$ . We also have the following information:

- The regression prediction:  $y(\mathbf{x}_n, \mathbf{w}) = \mathbf{w}^T \phi_n$ .
- The likelihood function:  $p(t_n | \phi_n, \mathbf{w}, \beta) = \mathcal{N}(t_n | \mathbf{w}^T \phi_n, 1/\beta)$
- The prior over  $\mathbf{w}$ :  $p(\mathbf{w}) = \mathcal{N}(\mathbf{w} | \mathbf{0}, \mathbf{I}/\alpha)$ .  $\mathbf{I}$  is the identity matrix,  $\mathbf{0}$  is a vector of 0's.
- The data are iid (independently and identically distributed).

Answer the following:

1. Write down the likelihood  $p(\mathcal{D} | \theta)$  using a) a product over  $N$  and b) in vector/matrix form. Tip: You can answer both a) and b) in one set of equations by starting with a), then simplifying to get b). For b) make sure to define any matrices and vectors.

*Solution:*

Given that the likelihood function for a single training vector is a normal distribution, we can write its product as:

$$\begin{aligned}
p(\mathcal{D}|\boldsymbol{\theta}) &= p(\mathbf{t}|\boldsymbol{\Phi}, \mathbf{w}, \beta) \\
&= \prod_{n=0}^N p(t_n|\boldsymbol{\phi}_n, \mathbf{w}, \beta) \\
&= \prod_{n=1}^N \mathcal{N}(t_n|\mathbf{w}^T \boldsymbol{\phi}_n, 1/\beta) \\
&= \prod_{n=1}^N \frac{\beta^{1/2}}{(2\pi)^{1/2}} \exp(-\frac{\beta}{2}(t_n - \mathbf{w}^T \boldsymbol{\phi}_n)^2)
\end{aligned}$$

In order to reach a vector form, we solve the product:

$$\begin{aligned}
p(\mathcal{D}|\boldsymbol{\theta}) &= \frac{\beta^{N/2}}{(2\pi)^{N/2}} \prod_{n=1}^N \exp(-\frac{\beta}{2}(t_n - \mathbf{w}^T \boldsymbol{\phi}_n)^2) \\
&= \underbrace{\frac{\beta^{N/2}}{(2\pi)^{N/2}}}_{\text{constant term}} \exp(-\frac{\beta}{2} \sum (t_n - \mathbf{w}^T \boldsymbol{\phi}_n)^2)
\end{aligned}$$

Next, Using the fact that  $\sum_i x_i^2 = \mathbf{x}^T \mathbf{x}$ , we transform the  $t_n$  into a vector, and  $\boldsymbol{\phi}_n$  into a matrix, where

$$\boldsymbol{\Phi} = \begin{pmatrix} -\phi_1 - \\ -\phi_2 - \\ \vdots \\ -\phi_N - \end{pmatrix}$$

Thus:

$$\begin{aligned}
p(\mathcal{D}|\boldsymbol{\theta}) &= \frac{\beta^{N/2}}{(2\pi)^{N/2}} \exp(-\frac{\beta}{2}(\mathbf{t} - \boldsymbol{\Phi} \mathbf{w})^T (\mathbf{t} - \boldsymbol{\Phi} \mathbf{w})) \\
&= \mathcal{N}(\mathbf{t}|\boldsymbol{\Phi} \mathbf{w}, \frac{1}{\beta} \mathbf{I})
\end{aligned}$$

We can see that the product is also a Normal distribution.

2. Write down the prior  $p(\mathbf{w})$  (by expanding the expression for multivariate Gaussian distribution). Compute its log.

*Solution:*

The prior  $p(\mathbf{w})$  is given by a normal distribution with  $\mu = 0$  and  $\sigma = \frac{1}{\alpha}\mathbf{I}$ . Expanding over the general formula of the Normal Distribution we get:

$$\begin{aligned} p(\mathbf{w}) &= \mathcal{N}(\mathbf{w}|\mathbf{0}, \frac{1}{\alpha}\mathbf{I}) \\ &= \frac{\alpha^{D/2}}{(2\pi)^{D/2}} \exp(-\frac{\alpha}{2}\mathbf{w}^T\mathbf{w}) \end{aligned}$$

Where  $D$  is the dimensionality of  $\mathbf{w}$ .

Next, because  $\log(\frac{a}{b}) = \log(a) - \log(b)$ ,  $\log(ab) = \log(a) + \log(b)$  and  $\log(\exp(x)) = x$

$$\begin{aligned} \ln p(\mathbf{w}) &= \underbrace{\frac{D}{2} \ln \alpha - \frac{D}{2} \ln(2\pi)}_{\text{constant term}} + (-\frac{\alpha}{2}\mathbf{w}^T\mathbf{w}) \\ &= -\frac{\alpha}{2}\mathbf{w}^T\mathbf{w} + \mathcal{C} \end{aligned}$$

Where  $\mathcal{C}$  is a constant term.

3. Write down an expression for the posterior over  $\mathbf{w}$ . Remember this will involve applying Bayes rule to the prior, likelihood, and evidence. The evidence will require an integral. You do not need the analytic form for the evidence, but you need the correct variables and conditioning variables, e.g. something like  $p(a|b, c)$  where you define  $a$ ,  $b$ , and  $c$ .

*Solution:*

According to Bayes rule we have:

$$p(\mathbf{w}|\mathcal{D}) = \frac{p(\mathbf{w})p(\mathcal{D}|\mathbf{w})}{p(\mathcal{D})}$$

From previous questions we got the likelihood and prior for the Bayes rule:

$$\begin{aligned} p(\mathcal{D}|\boldsymbol{\theta}) &= p(\mathcal{D}|\mathbf{w}) \\ &= \mathcal{N}(\mathbf{t}|\boldsymbol{\Phi}\mathbf{w}, \frac{1}{\beta}\mathbf{I}) \end{aligned}$$

$$p(\mathbf{w}) = \mathcal{N}(\mathbf{w}|\mathbf{0}, \frac{1}{\alpha}\mathbf{I})$$

We express the evidence by marginalizing. Since this is a continuous distribution, it requires an integral over  $\mathbf{w}$

$$p(\mathcal{D}) = \int (p(\mathcal{D}|\mathbf{w})p(\mathbf{w}))d\mathbf{w}$$

Substituting in Bayes Rule, we get:

$$\begin{aligned} p(\mathbf{w}|\mathcal{D}) &= \frac{\mathcal{N}(\mathbf{w}|\mathbf{0}, \frac{1}{\alpha}\mathbf{I})\mathcal{N}(\mathbf{t}|\boldsymbol{\Phi}\mathbf{w}, \frac{1}{\beta}\mathbf{I})}{\int \mathcal{N}(\mathbf{w}|\mathbf{0}, \frac{1}{\alpha}\mathbf{I})\mathcal{N}(\mathbf{t}|\boldsymbol{\Phi}\mathbf{w}, \frac{1}{\beta}\mathbf{I})d\mathbf{w}} \\ &= \frac{\mathcal{N}(\mathbf{w}|\mathbf{0}, \frac{1}{\alpha}\mathbf{I})\mathcal{N}(\mathbf{t}|\boldsymbol{\Phi}\mathbf{w}, \frac{1}{\beta}\mathbf{I})}{\int p(\mathbf{t}|\boldsymbol{\Phi}, \mathbf{w}, \beta)p(\mathbf{w}|\alpha)d\mathbf{w}} \\ &= \frac{\mathcal{N}(\mathbf{w}|\mathbf{0}, \frac{1}{\alpha}\mathbf{I})\mathcal{N}(\mathbf{t}|\boldsymbol{\Phi}\mathbf{w}, \frac{1}{\beta}\mathbf{I})}{p(\mathbf{t}|\boldsymbol{\Phi}, \alpha, \beta)} \end{aligned}$$

4. Compute the log-posterior, both for the a) and b) likelihood forms from above. Collect everything that does not depend on  $\mathbf{w}$  into a constant  $C$ . What parts of the previous expression do not depend on  $\mathbf{w}$ ? Why

is finding the MAP much simpler than finding the full posterior distribution?

*Solution:*

$$\begin{aligned}
\ln p(\mathbf{w}|\mathcal{D}) &= \ln\left(\frac{p(\mathbf{w})p(\mathcal{D}|\mathbf{w})}{p(\mathcal{D})}\right) \\
&= \ln p(\mathbf{w}) + \ln p(\mathcal{D}|\mathbf{w}) - \underbrace{\ln p(\mathcal{D})}_{\text{independent of } \mathbf{w}} \\
&= -\frac{\alpha}{2}\mathbf{w}^T\mathbf{w} - \frac{\beta}{2}\sum_{n=1}^N(t_n - \mathbf{w}^T\phi_n)^2 + \mathcal{C} \\
&= -\frac{\alpha}{2}\mathbf{w}^T\mathbf{w} - \frac{\beta}{2}(\mathbf{t} - \Phi\mathbf{w})^T(\mathbf{t} - \Phi\mathbf{w}) + \mathcal{C}
\end{aligned}$$

We can see that finding MAP is simpler because the evidence term (found in the full posterior distribution) becomes a constant term independent of  $\mathbf{w}$ , when finding its log. Furthermore, when finding the derivative of the log, the constant term goes to zero and it is discarded

5. Solve for  $\mathbf{w}_{\text{MAP}}$  by a) taking the derivative of the log-posterior with respect to  $\mathbf{w}$ , b) setting it to 0, and c) solving for  $\mathbf{w}$ . Do this for both forms of likelihood.

*Solution:*

First we solve in the non-matrix form. Given that the log likelihood is:

$$\ln p(\mathbf{w}|\mathcal{D}) = -\frac{\alpha}{2}\mathbf{w}^T\mathbf{w} - \frac{\beta}{2}\sum_{n=1}^N(t_n - \mathbf{w}^T\phi_n)^2 + \mathcal{C}$$

We find its derivative and set it to 0:

$$\frac{\partial \ln p(\mathbf{w}|\mathcal{D})}{\partial \mathbf{w}} = -\alpha\mathbf{w} - \beta\sum_{n=1}^N(t_n - \mathbf{w}^T\phi_n)(-\phi_n) = 0$$

Solving for  $\mathbf{w}$

$$\begin{aligned}
\alpha \mathbf{w} &= \beta \sum_{n=1}^N (t_n - \mathbf{w}^T \phi_n) \phi_n \\
\alpha \mathbf{w} &= \beta \sum_{n=1}^N t_n \phi_n - \underbrace{\mathbf{w}^T \phi_n}_{\text{scalar}} \phi_n \\
\alpha \mathbf{w} &= \beta \sum_{n=1}^N t_n \phi_n - \phi_n \underbrace{\mathbf{w}^T \phi_n}_{\text{same as } \phi_n^T \mathbf{w}} \\
(\alpha \mathbf{I} + \beta \sum_{n=1}^N \phi_n \phi_n^T) \mathbf{w} &= \beta \sum_{n=1}^N t_n \phi_n
\end{aligned}$$

$$\mathbf{w}_{\text{MAP}} = (\alpha \mathbf{I} + \beta \sum_{n=1}^N \phi_n \phi_n^T)^{-1} \beta \sum_{n=1}^N t_n \phi_n$$

Secondly, we solve in a matrix form. Again, given the log likelihood is:

$$\begin{aligned}
\ln p(\mathbf{w}|\mathcal{D}) &= -\frac{\alpha}{2} \mathbf{w}^T \mathbf{w} - \frac{\beta}{2} (\mathbf{t} - \Phi \mathbf{w})^T (\mathbf{t} - \Phi \mathbf{w}) + \mathcal{C} \\
&= -\frac{\alpha}{2} \mathbf{w}^T \mathbf{w} - \frac{\beta}{2} \mathbf{w}^T \Phi^T \Phi \mathbf{w} + \beta \mathbf{w}^T \Phi^T \mathbf{t} + \mathcal{D}
\end{aligned}$$

The derivative is given by:

$$\frac{\partial \ln p(\mathbf{w}|\mathcal{D})}{\partial \mathbf{w}} = -\alpha \mathbf{w} - \beta \Phi^T \Phi \mathbf{w} + \beta \Phi^T \mathbf{t} = 0$$

Solving for  $\mathbf{w}_{\text{MAP}}$ :

$$(\alpha \mathbf{I} + \beta \Phi^T \Phi) \mathbf{w} = \beta \Phi^T \mathbf{t}$$

$$\begin{aligned}
\mathbf{w}_{\text{MAP}} &= (\alpha \mathbf{I} + \beta \Phi^T \Phi)^{-1} \beta \Phi^T \mathbf{t} \\
&= (\beta (\frac{\alpha}{\beta} \mathbf{I} + \Phi^T \Phi))^{-1} \beta \Phi^T \mathbf{t} \\
&= \frac{1}{\beta} (\lambda \mathbf{I} + \Phi^T \Phi)^{-1} \Phi^T \mathbf{t}
\end{aligned}$$

6. **BONUS** Our prior for  $\mathbf{w}$  assumes the same marginal distribution for each entry in  $\mathbf{w}$ , including that of the first basis function  $\phi_0 = 1$ . What is the role this basis function? Why should we avoid placing the same penalty/prior for this basis? Rewrite  $p(\mathbf{w})$  so that the first basis function has its own prior/penalty.

*Solution:*

## 2 Probability distributions, likelihoods, and estimators

For these questions you will be working with different probability density functions listed in the table below. The purpose of these questions is to practice working with a variety of PDFs and to make computing likelihoods, MLEs, etc. more natural. Note below the *indicator* notation  $[x = 0]$  (and  $[x = 1]$ ). The square brackets evaluate to 1 if the argument is true, and 0 otherwise. E.g. if  $x$  is 1, the  $[x = 0] = 0$  and  $[x = 1] = 1$  (here  $[x = 0]$  is lazy notation; in Python you would write  $x == 0$ , for example). We will use the notation a lot, both below and when we learn about classification.

Distribution	$p(x \theta)$	Range of $x$	Range of $\theta$
Bernoulli	$\theta^{[x=1]}(1 - \theta)^{[x=0]}$	$x \in \{0, 1\}$	$0 \leq \theta \leq 1$
Beta	$\frac{\Gamma(\theta_1 + \theta_0)}{\Gamma(\theta_1)\Gamma(\theta_0)} x^{\theta_1 - 1} (1 - x)^{\theta_0 - 1}$	$0 \leq x \leq 1$	$\theta_1 > 0, \theta_0 > 0$
Poisson	$\frac{\theta^x}{x!} e^{-\theta}$	$x \in \{0, 1, 2, \dots\}$	$\theta > 0$
Gamma	$\frac{\theta_0^{\theta_1}}{\Gamma(\theta_1)} x^{\theta_1 - 1} e^{-\theta_1 x}$	$x \geq 0$	$\theta_1 \geq 0, \theta_0 \geq 0$
Gaussian	$\frac{1}{\sqrt{2\pi}\theta_1} e^{-\frac{1}{2}\left(\frac{x - \theta_0}{\theta_1}\right)^2}$	$-\infty < x < \infty$	$-\infty < \theta_0 < \infty, \theta_1 > 0$

### Question 2.1

For each of the probability distributions above, write down their normalizing constants. Remember that  $\int p(x|\theta)dx = 1$  for continuous  $x$  and  $\sum_x p(x|\theta) = 1$  for discrete  $x$ .

*Solution:*

The normalizing constants are those terms which do not depend on  $x$ , but that help to maintain the sum of the probabilities for all  $x$  equal to 1. Hereby we state the constants for the distributions given:

1. Bernoulli

There is no normalizing constant in this case, since  $X \in \{1, 0\}$  and its probabilities are complementary.

2. Beta

$$\frac{\Gamma(\theta_1 + \theta_0)}{\Gamma(\theta_1)\Gamma(\theta_0)}$$

3. Poisson

$$\exp^\theta$$

4. Gamma

$$\frac{\theta_1^{\theta_0}}{\Gamma(\theta_0)}$$

5. Gaussian

$$\frac{1}{\sqrt{2\pi\theta_1}}$$

## Question 2.2

You live in Amsterdam and find that it rains quite a lot. You want to estimate the probability that it will rain any given day of the year. Every month for a year you count the number of days with rain, and you get the following (from January to December): 22,19,16,16,14,14,17,18,19,20,21,21 (for a grand total of 217 days with rain).<sup>1</sup> Let  $r_t$  be an observation for day  $t$  in the year;  $r_t = 1$  means there was some rain on day  $t$ ,  $r_t = 0$  means there was no rain. We want to estimate the parameter  $\rho$ , the probability of rain on any day of the year. We assume a Bernoulli distribution for the observations  $\{r_t\}_{t=1}^{365}$ , that is  $p(r_t|\rho) = \text{Bernoulli}(r_t|\rho)$ . To answer these questions, the number of days of rain per month is not important, only the total for the year is relevant. With this information, answer the following questions:

---

<sup>1</sup>Source: <http://www.amsterdam.climateps.com/>.



1. What is the likelihood for a single observation? For the entire set of observations?

*Solution:*

The likelihood of a single observation is:

$$p(r_t|\rho) = \rho^{r_t}(1 - \rho)^{1-r_t}$$

Then, we find the probability of the whole set, representing it as the product of single observations. Using matrix notation we get:

$$\begin{aligned} p(\mathbf{r}|\rho) &= \prod_{t=1}^T \rho^{r_t}(1 - \rho)^{1-r_t} \\ &= \rho^{\sum_{t=1}^T r_t} (1 - \rho)^{\sum_{t=1}^T 1-r_t} \\ &= \rho^{n_1} (1 - \rho)^{n_0} \end{aligned}$$

Where:

$$\begin{aligned} \underbrace{n_1}_{\text{days with rain}} &= \sum_{t=1}^T r_t \\ \underbrace{n_0}_{\text{days without rain}} &= \sum_{t=1}^T 1 - r_t \end{aligned}$$

2. Write the log-likelihood for the entire set of observations.

*Solution:*

Using the logarithm arithmetic rules, we get:

$$\begin{aligned} \ln p(\mathbf{r}|\rho) &= \ln \rho^{n_1} (1 - \rho)^{n_0} \\ &= \ln \rho^{n_1} + \ln (1 - \rho)^{n_0} \\ &= n_1 \ln \rho + n_0 \ln (1 - \rho) \end{aligned}$$

3. Solve for the MLE of  $\rho$ . Do it in general (with symbols for counts  $n_0$ ,  $n_1$  for days without and with rain) and for this specific case (plug-in the numbers).

*Solution:*

In order to compute the MLE, first we find the derivative of the log-likelihood, and set it to 0.

$$\frac{\partial \ln p(\mathbf{r}|\rho)}{\partial \rho} = \frac{n_1}{\rho} + \frac{n_0}{1-\rho}(-1) = 0$$

Then, we find the value of  $\rho$

$$\frac{n_1}{\rho} = \frac{n_0}{1-\rho}$$

$$n_1 - n_1\rho = n_0\rho$$

$$\rho = \frac{n_1}{N}$$

Plug in the number, we have:

$$\rho = 217/365$$

4. Assume a Beta prior for  $\rho$  with parameters  $a$  and  $b$ . What is the MAP for  $\rho$ ?

*Solution:*

We follow the same procedure, but this time for the posterior, given by:

$$\begin{aligned} f = \ln p(\rho|\mathbf{r}) &\propto \ln p(\rho|\mathbf{r}) + \ln p(\rho) \\ &= n_1 \ln \rho + n_0 \ln(1-\rho) + (a-1) \ln \rho + (b-1) \ln(1-\rho) \end{aligned}$$

We first find the derivative and set it to 0.

$$\frac{\partial f}{\partial \rho} = \frac{n_1}{\rho} - \frac{n_0}{1-\rho} + \frac{a-1}{\rho} - \frac{b-1}{1-\rho} = 0$$

Then we solve for  $\rho$

$$\rho = \frac{n_1 + a - 1}{N + a + b - 2}$$

5. Write the form of the posterior distribution for  $\rho$ ? You do not need to solve it analytically.

*Solution:*

Using Bayes Theorem, we have:

$$p(\rho|\mathbf{r}) = \frac{p(\mathbf{r}|\rho)p(\rho)}{p(\mathbf{r})}$$

Substituting with the specifics of this problem, we get:

$$\begin{aligned} p(\rho|\mathbf{r}) &= \frac{\rho^{n_1+a-1}(1-\rho)^{n_0+b-1}}{\int \rho^{n_1+a-1}(1-\rho)^{n_0+b-1}d\rho} \\ &= \frac{\Gamma(\mathcal{N} + a + b)}{\Gamma(n_1 + a)\Gamma(n_0 + b)} \rho^{n_1+a-1}(1-\rho)^{n_0+b-1} \\ &= \mathcal{B}(\rho|a + n_1, b + n_0) \end{aligned}$$

6. (Optional) Solve for the posterior distribution analytically. Hint: it is a Beta distribution.

*Solution:*

### Question 2.3

You work in the staffing department of a maternity hospital and part of your job is to determine the staffing requirements during the night shift at your hospital. This might mean the number of doctors and nurses at the hospital and the number of doctors on call (if there are more than the average number of deliveries). Your goal is to determine the distribution over the number of deliveries during the night shift  $d_t \in \{0, 1, 2, \dots\}$  ( $d$  for delivery count,  $t$  for time, the index of the night). With this you can compute the mean, the probability of more than 5 deliveries, etc. You collect data for two weeks, i.e.  $d_1, \dots, d_{14} = 4, 7, 3, 0, 2, 2, 1, 5, 4, 4, 3, 3, 2, 3$ . You assume the observations are explained by a Poisson distribution with parameter  $\lambda$  over the discrete delivery counts. With this information, answer the following questions:

1. What is the likelihood for a single observation? For the entire set of observations?

*Solution:*

The likelihood of a single observation is:

$$p(d_t|\lambda) = \frac{\lambda^{d_t}}{d_t!} \exp(-\lambda)$$

Then, we find the probability of the whole set, representing it as the product of single observations. Using matrix notation we get:

$$\begin{aligned} p(\mathbf{d}|\lambda) &= \prod_{t=1}^T \frac{\lambda^{d_t}}{d_t!} \exp(-\lambda) \\ &= \frac{\lambda^{\sum_{t=1}^T d_t}}{\prod_{t=1}^T d_t!} \exp(-T\lambda) \\ &= \frac{\lambda^n}{\prod_{t=1}^T d_t!} \exp(-T\lambda) \end{aligned}$$

2. Write the log-likelihood for the entire set of observations.

*Solution:*

Using the logarithm arithmetic rules, we get:

$$\ln p(\mathbf{d}|\lambda) = n \ln \lambda - T\lambda - \sum_{t=1}^T \ln(d_t!)$$

3. Solve for the MLE of  $\lambda$ . Do it in general and for this specific case (plug-in the numbers).

*Solution:*

$$\begin{aligned} f &= \ln p(\mathbf{d}|\lambda) \\ &= n \ln \lambda - T\lambda - \sum_{t=1}^T \ln(d_t!) \\ \frac{\partial f}{\partial \lambda} &= \frac{n}{\lambda} - T = 0 \\ \lambda &= \frac{n}{T} = \frac{43}{14} \end{aligned}$$

4. Assume a Gamma prior for  $\lambda$  with parameters  $a$  and  $b$ . What is the MAP estimate of  $\lambda$ ?

*Solution:*

$$\begin{aligned} f &= \ln p(\mathbf{d}|\lambda) + \ln p(\lambda) \\ &= n \ln \lambda - T\lambda - \sum_{t=1}^T \ln d_t! + (a-1) \ln \lambda - b\lambda + C \\ \frac{\partial f}{\partial \lambda} &= \frac{n}{\lambda} - T + \frac{(a-1)}{\lambda} - b = 0 \\ \lambda &= \frac{n+a-1}{T+b} \end{aligned}$$

5. Write the form of the posterior distribution for  $\lambda$ ? (You do not need to solve it analytically)

*Solution:*

$$\begin{aligned}
p(\lambda|\mathbf{d}) &= \frac{p(\mathbf{d}|\lambda)p(\lambda)}{\int p(\mathbf{d}|\lambda)p(\lambda)d\lambda} \\
&= \frac{\frac{\lambda^n}{\prod_{t=1}^T d_t!} \exp(-T\lambda) \frac{b^a}{\Gamma(a)} \lambda^{a-1} \exp(-b\lambda)}{\int \frac{\lambda^n}{\prod_{t=1}^T d_t!} \exp(-T\lambda) \frac{b^a}{\Gamma(a)} \lambda^{a-1} \exp(-b\lambda)d\lambda} \\
&= \frac{\lambda^{n+a-1} \exp(-(T+b)\lambda)}{\int \lambda^{n+a-1} \exp(-(T+b)\lambda)d\lambda} \\
&= \frac{(T+b)^{n+a}}{\Gamma(n+a)} \lambda^{n+a-1} \exp(-(T+b)\lambda) \\
&= \mathcal{G}(\lambda|a+n, b+T)
\end{aligned}$$

6. (Optional) Solve for the posterior distribution analytically. Hint: it is a Gamma distribution.

*Solution:*

#### Question 2.4

You have developed a blood test aimed at detecting a disease  $d \in \{0, 1\}$  (disease is absent ( $d = 0$ ) or present ( $d = 1$ )). The test measures the level of a specific indicator of the disease, that is it returns a real valued number relative to some baseline (so the levels can be both negative and positive – anywhere along the real line). Two models of the population are built: one for the patients with the disease, and another for the general population. Measurements tend to have a Gaussian shape, and we therefore model the entire population as a mixture of two Gaussians. That is,  $p(l) = p(d = 0)p(l|d = 0) + p(d = 1)p(l|d = 1)$ , where  $p(d)$  is the prior distribution of patients with and without the disease in the general population and  $p(l|d)$  are conditional Gaussian distributions, one for the patients with disease, and one for those without. Note: with this question and the previous two, we are simply applying rules of probability (with some algebra) to get the form of the posterior distribution; however, in this problem we are also classifying (since our target is the discrete label  $d$ ).

Assume we know  $p(d = 0) = \pi_0 = 0.999$  and  $p(d = 1) = \pi_1 = 0.001$  from previous experience. We do not know the parameters  $\mu_0, \sigma_0^2$  (the mean and variance of the disease-free population) nor  $\mu_1, \sigma_1^2$  (for the disease population). We measure levels  $\{l_n\}_{n=1}^N$  for  $N$  people, and we know that  $n \in \{D_0\}$  are the indices for the disease free patients and  $n \in \{D_1\}$  are the indices for the patients with the disease (i.e.  $D_0$  and  $D_1$  are non-intersecting sets of indices from 1 to  $N$ ). With this information, answer the following questions:

1. Write down the likelihood of the observations as a product over  $N$  level recordings. Hint: use indicator notation (like in the Bernoulli distribution) to distinguish between  $d_n = 0$  and  $d_n = 1$  in the likelihood.

*Solution:*

$$p(\mathbf{d}, \mathbf{l} | \boldsymbol{\theta}) = \prod_{n=1}^N (\pi_1 \mathcal{N}(l_n | \mu_1, \sigma_1^2))^{d_n} (\pi_0 \mathcal{N}(l_n | \mu_0, \sigma_0^2))^{1-d_n}$$

2. Write down the likelihood as a product over the likelihoods for  $\{D_0\}$  and  $\{D_1\}$ .

*Solution:*

$$p(\mathbf{d}, \mathbf{l} | \boldsymbol{\theta}) = \left( \prod_{n \in D_1} \pi_1 \mathcal{N}(l_n | \mu_1, \sigma_1^2) \right) \left( \prod_{n \in D_0} \pi_0 \mathcal{N}(l_n | \mu_0, \sigma_0^2) \right)$$

3. Compute the log-likelihood.

*Solution:*

$$\begin{aligned} \ln p(\mathbf{d}, \mathbf{l} | \boldsymbol{\theta}) &= \sum_{n \in D_1} (\ln \pi_1 + \ln \mathcal{N}(l_n | \mu_1, \sigma_1^2)) + \sum_{n \in D_0} (\ln \pi_0 + \ln \mathcal{N}(l_n | \mu_0, \sigma_0^2)) \\ &= \sum_{n \in D_1} \left( \ln \pi_1 + \frac{1}{2} \ln(2\pi) - \frac{1}{2} \ln(\sigma_1^2) - \frac{1}{2} \frac{(l_n - \mu_1)^2}{\sigma_1^2} \right) \\ &\quad + \sum_{n \in D_0} \left( \ln \pi_0 + \frac{1}{2} \ln(2\pi) - \frac{1}{2} \ln(\sigma_0^2) - \frac{1}{2} \frac{(l_n - \mu_0)^2}{\sigma_0^2} \right) \end{aligned}$$

4. Find the MLE for  $\mu_0$  and  $\sigma_0^2$ . Assume we can do the same for  $\mu_1$  and  $\sigma_1^2$ .

*Solution:*

$$\begin{aligned} \frac{\partial \ln p(\mathbf{d}, \mathbf{l} | \boldsymbol{\theta})}{\partial \mu_0} &= \sum_{n \in D_0} \left( -\frac{(l_n - \mu_0)}{\sigma_0^2} (-1) \right) = 0 \\ 0 &= \sum_{n \in D_0} (l_n - \mu_0) \\ 0 &= \sum_{n \in D_0} l_n - N_0 \mu_0 \\ \mu_0 &= \frac{1}{N_0} \sum_{n \in D_0} l_n \end{aligned}$$

$$\begin{aligned}\frac{\partial \ln p(\mathbf{d}, \mathbf{l} | \boldsymbol{\theta})}{\partial \sigma_0^2} &= \sum_{n \in D_0} \left( -\frac{1}{2\sigma_0^2} - \frac{1}{2} \frac{(l_n - \mu_0)^2}{(\sigma_0^2)^2} (-1) \right) = 0 \\ \frac{N_0}{2\sigma_0^2} &= \frac{1}{2} \frac{\sum_{n \in D_0} (l_n - \mu_0)^2}{(\sigma_0^2)^2} \\ \sigma_0^2 &= \frac{1}{N_0} \sum_{n \in D_0} (l_n - \mu_0)^2\end{aligned}$$

5. We now have our models. To make a prediction, solve for  $p(d = 1 | l_*)$ , where  $l_*$  is a level recorded for a new patient. Hint: use Bayes theorem.  
*Solution:*

$$\begin{aligned}p(d = 1 | l_*) &= \frac{\pi_1 \mathcal{N}(l_* | \mu_1, \sigma_1^2)}{\pi_1 \mathcal{N}(l_* | \mu_1, \sigma_1^2) + \pi_0 \mathcal{N}(l_* | \mu_0, \sigma_0^2)} \\ &= \frac{1}{1 + \frac{\pi_0 \mathcal{N}(l_* | \mu_0, \sigma_0^2)}{\pi_1 \mathcal{N}(l_* | \mu_1, \sigma_1^2)}} \\ &= \frac{1}{1 + \exp^{-a(l_*)}} \\ a(l_*) &= \ln\left(\frac{\pi_1 \mathcal{N}(l_* | \mu_1, \sigma_1^2)}{\pi_0 \mathcal{N}(l_* | \mu_0, \sigma_0^2)}\right)\end{aligned}$$

6. Reduce your solution to have the form of a sigmoid, i.e.

$$p(d = 1 | l_*) = \frac{1}{1 + e^{-a(l_*)}}.$$

*Solution:*

$$p(d = 1 | l_*) = \frac{1}{1 + \exp^{-a(l_*)}} = \frac{1}{1 + \frac{\pi_0 \mathcal{N}(l_* | \mu_0, \sigma_0^2)}{\pi_1 \mathcal{N}(l_* | \mu_1, \sigma_1^2)}}$$

Since we know the Normal distribution has a sigmoid form, this last solution has the desired form.