

# Machine Learning 1 - Homework 2

Pascal Mattia Esser

September 20th, 2017

## 1 MAP solution for Linear Regression

### 1.1

given:

$$\mathcal{N}(x|\mu\sigma^2) = \frac{1}{(2\pi\sigma^2)^{1/2}} \exp\left(-\frac{1}{2\sigma^2}(x - \mu)^2\right) \quad (1)$$

a) we can write as product form:

$$p(\mathbf{x}|\mathbf{w}) = \prod_{i=1}^N \mathcal{N}(t_i|\mathbf{w}^T \Phi_i, \sigma^2) \quad (2)$$

$$= \prod_{i=1}^N \frac{1}{(2\pi\sigma^2)^{1/2}} \exp\left(-\frac{1}{2\sigma^2}(t_i - \mathbf{w}^T \Phi_i)^2\right) \quad (3)$$

b) write as matrix form

$$p(\mathbf{x}|\mathbf{w}) = \prod_{i=1}^N \frac{1}{(2\pi\sigma^2)^{1/2}} \exp\left(-\frac{1}{2\sigma^2}(t_i - \mathbf{w}^T \Phi_i)^2\right) \quad (4)$$

$$= \left(\frac{\beta}{2\pi}\right)^{N/2} \prod_{i=1}^N \exp\left(-\frac{\beta}{2}(t_i - \mathbf{w}^T \Phi_i)^2\right) \quad (5)$$

$$= \left(\frac{\beta}{2\pi}\right)^{N/2} \exp\left(-\frac{\beta}{2} \sum_{i=1}^N (t_i - \mathbf{w}^T \Phi_i)^2\right) \quad (6)$$

$$= \left(\frac{\beta}{2\pi}\right)^{N/2} \exp\left(-\frac{\beta}{2} \|\mathbf{t} - \mathbf{w}^T \Phi\|^2\right) \quad (7)$$

### 1.2

starting from the multivariate Gaussian distribution:

$$p(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{(2\pi|\boldsymbol{\Sigma}|)^{M/2}} \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu})\right) \quad (8)$$

with

$$p(\mathbf{w}) = \mathcal{N}(\mathbf{w}|\mathbf{0}, \alpha^{-1}\mathbf{I}) \quad (9)$$

you get

$$= \frac{1}{(2\pi|\alpha^{-1}\mathbf{I}|^2)^{M/2}} \exp\left(-\frac{1}{2}(\mathbf{w} - \mathbf{0})^T(\alpha^{-1}\mathbf{I})^{-1}(\mathbf{w} - \mathbf{0})\right) \quad (10)$$

$$= \left(\frac{\alpha}{2\pi}\right)^{M/2} \exp\left(-\frac{\alpha\mathbf{w}^T\mathbf{w}}{2}\right) \quad (11)$$

calculate the log

$$\ln(p(\mathbf{w})) = \frac{M}{2} \log\left(\frac{\alpha}{2\pi}\right) - \frac{\alpha\mathbf{w}^T\mathbf{w}}{2} \quad (12)$$

### 1.3

starting from a Bayes formula

$$p(\mathbf{w}|\mathbf{x}) = \frac{p(\mathbf{x}|\mathbf{w})p(\mathbf{w})}{p(\mathbf{x})} \quad (13)$$

and setting:

$$p(\mathbf{x}) = \int p(\mathbf{t}|\Phi, \mathbf{w}, \beta)p(\mathbf{w})d\mathbf{w} \quad (14)$$

gives us:

$$p(\mathbf{w}|\mathbf{t}, \mathbf{x}, \beta, \alpha) = \frac{\prod^N \mathcal{N}(t_i|\mathbf{w}^T\phi_i, \beta^{-1})\mathcal{N}(\mathbf{w}|\mathbf{0}, \alpha^{-1}\mathbf{I})}{\int p(\mathbf{t}|\Phi, \mathbf{w}, \beta)p(\mathbf{w})d\mathbf{w}} \quad (15)$$

rewrite the integral following [Bis06, p. 93]:

$$p(\mathbf{w}|\mathbf{t}, \mathbf{x}, \beta, \alpha) = \frac{\prod^N \mathcal{N}(t_i|\mathbf{w}^T\phi_i, \beta^{-1})\mathcal{N}(\mathbf{w}|\mathbf{0}, \alpha^{-1}\mathbf{I})}{\mathcal{N}(\mathbf{t}|\Phi\mathbf{0} + \mathbf{0}, \frac{1}{\beta}\mathbf{I} + \Phi(\alpha^{-1})^{-1}\Phi^T)} \quad (16)$$

$$= \frac{\prod^N \mathcal{N}(t_i|\mathbf{w}^T\phi_i, \beta^{-1})\mathcal{N}(\mathbf{w}|\mathbf{0}, \alpha^{-1}\mathbf{I})}{\mathcal{N}(\mathbf{0}, \frac{1}{\beta}\mathbf{I} + \frac{1}{\alpha}\Phi\Phi^T)} \quad (17)$$

### 1.4

for the logarithm, we know:

$$\ln p(\mathbf{w}|\mathbf{x}) = \ln p(\mathbf{x}|\mathbf{w}) + \ln p(\mathbf{w}) - \ln p(\mathbf{x}) \quad (18)$$

a) writing it in product form (keep  $\ln p(\mathbf{x})$  as term and don't expand it, as it does not depend on  $\mathbf{w}$ ):

$$\ln p(\mathbf{w}|\mathbf{x}) \quad (19)$$

$$= \ln\left(\left(\frac{\beta}{2\pi}\right)^{M/2} \prod_{i=1}^N \exp\left(-\frac{\beta}{2}(t_i - \mathbf{w}^T\phi_i)^2\right) + \left(\frac{\alpha}{2\pi}\right)^{M/2} \exp\left(-\frac{\alpha\mathbf{w}^T\mathbf{w}}{2}\right)\right) - \ln p(\mathbf{x}) \quad (20)$$

$$= \frac{M}{2} \ln \frac{\beta}{2\pi} + \left(-\frac{\beta}{2} \sum_{i=1}^N (t_i - \mathbf{w}^T\phi_i)^2\right) + \frac{M}{2} \log\left(\frac{\alpha}{2\pi}\right) - \frac{\alpha\mathbf{w}^T\mathbf{w}}{2} - \ln p(\mathbf{x}) \quad (21)$$

therms, depending on  $\mathbf{w}$  and substitute everything else in  $C$ :

$$\ln p(\mathbf{w}|\mathbf{x}) = -\frac{\beta}{2} \sum_{i=1}^N (t_i - \mathbf{w}^T \phi_i)^2 - \frac{\alpha \mathbf{w}^T \mathbf{w}}{2} + C \quad (22)$$

b) writing it in matrix form:

$$\ln p(\mathbf{w}|\mathbf{x}) \quad (23)$$

$$= \ln \left( \left( \frac{\beta}{2\pi} \right)^{M/2} \exp \left( -\frac{\beta}{2} \|\mathbf{t} - \mathbf{w}^T \Phi\|^2 \right) + \left( \frac{\alpha}{2\pi} \right)^{M/2} \exp \left( -\frac{\alpha \mathbf{w}^T \mathbf{w}}{2} \right) \right) - \ln p(\mathbf{x}) \quad (24)$$

$$= \frac{M}{2} \ln \frac{\beta}{2\pi} + \left( -\frac{\beta}{2} \|\mathbf{t} - \mathbf{w}^T \Phi\|^2 \right) + \frac{M}{2} \log \left( \frac{\alpha}{2\pi} \right) - \frac{\alpha \mathbf{w}^T \mathbf{w}}{2} - \ln p(\mathbf{x}) \quad (25)$$

therms, depending on  $\mathbf{w}$ :

$$\ln p(\mathbf{w}|\mathbf{x}) = \left( -\frac{\beta}{2} \|\mathbf{t} - \mathbf{w}^T \Phi\|^2 \right) - \frac{\alpha \mathbf{w}^T \mathbf{w}}{2} + C \quad (26)$$

were  $C$ , the terms not depending on  $\mathbf{w}$  is

$$C = \frac{M}{2} \ln \frac{\beta}{2\pi} + \frac{M}{2} \log \left( \frac{\alpha}{2\pi} \right) - \ln p(\mathbf{x}) \quad (27)$$

It is easier to calculate the MAP because you don't have to solve the integral in the evidence term  $p(\mathbf{x})$ , which can be hard analytically (or intractable).

## 1.5

take the derivative of the product form:

$$\frac{\partial}{\partial \mathbf{w}} \ln p(\mathbf{w}|\mathbf{x}) = \frac{\partial}{\partial \mathbf{w}} \left( -\frac{\beta}{2} \sum_{i=1}^N (t_i - \mathbf{w}^T \phi_i)^2 - \frac{\alpha \mathbf{w}^T \mathbf{w}}{2} \right) \quad (28)$$

$$= -\frac{\beta}{2} \sum_{i=1}^N \frac{\partial}{\partial \mathbf{w}} (t_i - \mathbf{w}^T \phi_i)^2 - \alpha \mathbf{I} \mathbf{w}^T \quad (29)$$

$$= -\beta \sum_{i=1}^N (t_i - \mathbf{w}^T \phi_i) \phi_i^T - \alpha \mathbf{I} \mathbf{w}^T \quad (30)$$

$$= \sum_{i=1}^N (\phi_i^T t_i) + \mathbf{w}^T \sum_{i=1}^N (\phi_i \phi_i^T + \lambda \mathbf{I}) \quad (31)$$

$$\mathbf{w} = \left( \sum_{i=1}^N (\phi_i^T \phi_i + \lambda \mathbf{I}) \right)^{-1} \sum_{i=1}^N (\phi_i^T t_i) \quad (32)$$

Take the derivative of the matrix form:

$$\frac{\partial}{\partial \mathbf{w}} \ln p(\mathbf{w}|\mathbf{x}) = \frac{\partial}{\partial \mathbf{w}} \left( -\frac{\beta}{2} \|\mathbf{t} - \mathbf{w}^T \Phi\|^2 \right) + -\frac{\alpha \mathbf{w}^T \mathbf{w}}{2} + C \quad (33)$$

$$= \left( -\frac{\beta}{2} \frac{\partial}{\partial \mathbf{w}} \|\mathbf{t} - \mathbf{w}^T \Phi\|^2 \right) - \frac{\partial}{\partial \mathbf{w}} \frac{\alpha \mathbf{w}^T \mathbf{w}}{2} \quad (34)$$

$$= \left( -\frac{\beta}{2} \frac{\partial}{\partial \mathbf{w}} (\mathbf{t} - \mathbf{w}^T \Phi)^T (\mathbf{t} - \mathbf{w}^T \Phi) \right) - \frac{\partial}{\partial \mathbf{w}} \frac{\alpha \mathbf{w}^T \mathbf{w}}{2} \quad (35)$$

$$= \left( -\frac{\beta}{2} \Phi^T \mathbf{I} (\mathbf{t} - \mathbf{w}^T \Phi) + \Phi^T \mathbf{I}^T (\mathbf{t} - \mathbf{w}^T \Phi) \right) - \alpha \mathbf{I} \mathbf{w}^T \quad (36)$$

$$= (\beta \Phi^T \mathbf{I} \mathbf{t}) + (\beta \Phi^T \mathbf{I} \Phi + \alpha \mathbf{I}) \mathbf{w}^T \quad (37)$$

set to 0 and solve for  $\mathbf{w}$ :

$$\mathbf{w}^T = (\beta \Phi^T \mathbf{I} \Phi + \alpha \mathbf{I})^{-1} (\beta \Phi^T \mathbf{I} \mathbf{t}) \quad (38)$$

$$\mathbf{w} = (\Phi^T \Phi + \lambda \mathbf{I})^{-1} (\Phi^T \mathbf{t}) \quad (39)$$

with  $\lambda = \alpha/\beta$

## 1.6

$\phi(\mathbf{x}) = 1$  is the offset or bias of the equation:

You should not use the same prior, because otherwise the same prior believe is applied on the offset of the function and the variables, that govern the shape of the function. taking this into account, we can write the new prior as:

$$p(\mathbf{w}) = \left( \frac{\alpha}{2\pi} \right)^{(N-1)/2} \exp \left( -\frac{\alpha \mathbf{w}^T \mathbf{w} - \mathbf{w}_0^2}{2} \right) \left( \frac{\gamma}{2\pi} \right)^{1/2} \exp \left( -\frac{\gamma \mathbf{w}_0^2}{2} \right) \quad (40)$$

here we have different distributions for  $\mathbf{w}$  and  $\mathbf{w}_0$ .

## 2 Probability distributions, likelihoods, and estimators

### 2.1 Describe distributions

#### 2.1.1 Bernoulli

- experiments, that ask for yes or now questions.
- it takes value 1 with probability  $p$  and 0 with  $p - 1$
- coin toss

#### 2.1.2 Beta

- contains values in the interval  $[0,1]$
- uses  $\alpha$  and  $\beta$  as shaping parameters for the from of the distribution
- betting rates in basketball for throwing vs. hitting

### 2.1.3 Poisson

- probability function over discrete probability values
- probability of a given number of events in a fixed time/space interval.
- known average rate and independently of the time of the last event
- probability of n mails per day

### 2.1.4 Gamma

- Poisson, where the waiting time between event is relevant
- mails are poisson distributed. now the waiting time till the 4th mail is gamma distributed

### 2.1.5 Gaussian

- distribution over variables for an sufficient large i.i.d. data set
- mean and standard deviation define the function
- hight distribution of humans

### 2.1.6 Log-Normal

- random variables, who logs are normal distributed
- distribution of the firing rate across a population of neurons in the hippocampus

## 2.2 Rain in Den Helder

given variables:

$$n_1 = 207 \tag{41}$$

$$N = 365 \tag{42}$$

$$n_0 = 158 \tag{43}$$

Assuming a Bernoulli distribution, we can write:

$$p(\mathbf{x}|\rho) = \prod_{t=1}^N p(r_t|\rho) \tag{44}$$

using Bernoulli:

$$p(r_t|\rho) = \rho^{[r_t=1]}(1 - \rho)^{[r_t=0]} \tag{45}$$

$$p(\mathbf{x}|\rho) = \prod_{t=1}^N \rho^{[r_t=1]}(1 - \rho)^{[r_t=0]} \tag{46}$$

**2.2.1**

we can write:

$$p(\mathbf{x}|\rho) = \prod_{i=1}^N \rho^{[r_i=1]} (1-\rho)^{[r_i=0]} \quad (47)$$

$$= \rho^{n_1} (1-\rho)^{n_0} \quad (48)$$

**2.2.2**

computing the log

$$\ln(p(\mathbf{x}|\rho)) = \ln \rho^{n_1} + \ln(1-\rho)^{n_0} \quad (49)$$

**2.2.3**

Take the derivatives

$$\frac{\partial}{\partial \rho} = \frac{n_1}{\rho} + \frac{n_0}{1-\rho} = 0 \quad (50)$$

solve for  $\rho$ :

$$\rho = \frac{n_1}{n_0} (1-\rho) \quad (51)$$

$$= \frac{n_1}{n_1 + n_0} \quad (52)$$

$$= \frac{n_1}{N} \quad (53)$$

use the numbers:

$$\rho = \frac{207}{365} \approx 0.5671 \quad (54)$$

**2.2.4**

rewrite with a known Beta prior:

$$p(\mathbf{x}|\rho) = \frac{\rho^{n_1} (1-\rho)^{n_0} \left( \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \rho^{a-1} (1-\rho)^{b-1} \right)}{p(\mathbf{x})} \quad (55)$$

compute the ln:

$$\ln p(\mathbf{x}|\rho) = n_1 \ln \rho + \ln n_0 (1-\rho) + \ln \left( \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \right) + (a-1) \ln \rho + (b-1) \ln(1-\rho) + C \quad (56)$$

take the derivative

$$\frac{\partial}{\partial \rho} = \frac{n_1}{\rho} - \frac{n_0}{1-\rho} + \frac{a-1}{\rho} - \frac{b-1}{1-\rho} \quad (57)$$

setting to 0 and solving for  $\rho$ :

$$\rho = \frac{n_1 + a - 1}{n_1 + a + b + n_0 - 2} \quad (58)$$

### 2.2.5

$$p(\rho|\mathbf{x}) = \frac{\rho^{n_1}(1-\rho)^{n_0} \left( \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \rho^{a-1}(1-\rho)^{b-1} \right)}{\int \rho'^{n_1}(1-\rho')^{n_0} \left( \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \rho'^{a-1}(1-\rho')^{b-1} \right) d\rho'} \quad (59)$$

### 2.2.6

starting from:

$$p(\rho|\mathbf{x}) = \frac{\rho^{n_1}(1-\rho)^{n_0} \left( \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \rho^{a-1}(1-\rho)^{b-1} \right)}{\int \rho'^{n_1}(1-\rho')^{n_0} \left( \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \rho'^{a-1}(1-\rho')^{b-1} \right) d\rho'} \quad (60)$$

canceling out

$$\frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \quad (61)$$

and solving for the denominator

$$\int \rho'^{n_1}(1-\rho')^{n_0} (\rho'^{a-1}(1-\rho')^{b-1}) d\rho' = 1 \quad (62)$$

rewrite:

$$\int \rho^{a+n_1-1}(1-\rho)^{b+n_0-1} = 1 \quad (63)$$

and compare with

$$\rho'^{a-1}(1-\rho')^{b-1} \quad (64)$$

gives

$$p(\rho) = B(a + n_1, b + n_0) \quad (65)$$

as the Beta distribution

## 2.3 staffing department of a maternity hospital

$$T = 14 \quad (66)$$

$$n = 4 + 7 + 3 + 0 + 2 + 2 + 1 + 5 + 4 + 4 + 3 + 3 + 2 + 3 = 43 \quad (67)$$

### 2.3.1

likelihood for one observation: set  $n = d$  for one observation

$$p(\mathbf{x}|\lambda) = \frac{\lambda^{d_i}}{d_i!} e^{-\lambda} \quad (68)$$

for all observations

$$p(\mathbf{x}|\lambda) = \prod \frac{\lambda^{d_i}}{d_i!} e^{-\lambda} \quad (69)$$

$$= \frac{1}{\prod d_i!} \lambda^n e^{-T\lambda} \quad (70)$$

### 2.3.2

likelihood for all observations:

$$\ln p(\mathbf{x}|\lambda) = \ln \left( \prod \frac{\lambda^{d_i}}{d_i!} e^{-\lambda} \right) \quad (71)$$

$$= \sum \ln \frac{\lambda^{d_i}}{d_i!} e^{-\lambda} \quad (72)$$

$$= -T\lambda + n \ln \lambda - \sum \ln(d_i!) \quad (73)$$

### 2.3.3

take the derivative of the log:

$$\frac{\partial}{\partial \lambda} = \frac{n}{\lambda} - T \quad (74)$$

set the derivative to zero and solve for  $\lambda$

$$\lambda = \frac{n}{T} \quad (75)$$

calculate with numbers:

$$\lambda = \frac{43}{14} \approx 3.0714 \quad (76)$$

### 2.3.4

Assume a Gamma prior for  $\lambda$

$$p(\lambda|\mathbf{x}) = \frac{\prod \left( \frac{\lambda^{d_i}}{d_i!} e^{-\lambda} \right) \frac{b^a}{\Gamma(a)} \lambda^{a-1} e^{-b\lambda}}{p(\mathbf{x})} \quad (77)$$

calculate the ln

$$\ln p(\lambda|\mathbf{x}) = \frac{n}{\lambda} - T + \ln \left( \frac{b^a}{\Gamma(a)} \right) + \lambda \ln(a-1) - b\lambda + C \quad (78)$$

calculate the derivative:

$$\frac{\partial}{\partial \lambda} = \frac{n}{\lambda} - T + \lambda \frac{1}{a-1} - b = 0 \quad (79)$$

solving for  $\lambda$  gives us:

$$\lambda = \frac{n+a-1}{T+b} \quad (80)$$



**2.3.5**

Write down the form of the posterior distribution for  $\lambda$

$$p(\lambda|\mathbf{x}) = \frac{\prod \left( \frac{\lambda^{d_i}}{d_i!} e^{-\lambda} \right) \frac{b^a}{\Gamma(a)} \lambda^{a-1} e^{-b\lambda}}{\int \prod \left( \frac{\lambda'^{d_i}}{d_i!} e^{-\lambda'} \right) \frac{b^a}{\Gamma(a)} \lambda'^{a-1} e^{-b\lambda'} d\lambda'} \quad (81)$$

**2.3.6**

same approach as in [subsubsection 2.2.6](#): start with

$$p(\lambda|\mathbf{x}) = \frac{\prod \left( \frac{\lambda^{d_i}}{d_i!} e^{-\lambda} \right) \frac{b^a}{\Gamma(a)} \lambda^{a-1} e^{-b\lambda}}{\int \prod \left( \frac{\lambda'^{d_i}}{d_i!} e^{-\lambda'} \right) \frac{b^a}{\Gamma(a)} \lambda'^{a-1} e^{-b\lambda'} d\lambda'} \quad (82)$$

rewrite as:

$$p(\lambda|\mathbf{x}) = \frac{\prod \frac{1}{d_i!} \lambda^T e^{-n\lambda} \frac{b^a}{\Gamma(a)} \lambda^{a-1} e^{-b\lambda}}{\int \prod \frac{1}{d_i!} \lambda'^T e^{-n\lambda'} \frac{b^a}{\Gamma(a)} \lambda'^{a-1} e^{-b\lambda'} d\lambda'} \quad (83)$$

after canceling out the constants we get for the denominator:

$$\int \lambda^n e^{-T\lambda} \lambda^{a-1} e^{-b\lambda} d\lambda = 1 \quad (84)$$

now bring it in the form of the initial prior,

$$\lambda^{a-1} e^{-b\lambda} \quad (85)$$

gives a gamma function of:

$$p(\lambda) = \Gamma(a + n, b + T) \quad (86)$$

**References**

- [Bis06] Christopher M. Bishop. *Pattern Recognition and Machine Learning (Information Science and Statistics)*. Secaucus, NJ, USA: Springer-Verlag New York, Inc., 2006. ISBN: 0387310738.