

Homework set 1

Machine Learning 1

Yke Rusticus

11306386

September 12, 2019

2 Multivariate Calculus

Question 2.1

2.1 a)

$$\nabla_{\boldsymbol{\mu}}(\mathbf{x} - \boldsymbol{\mu})^\top \Sigma^{-1}(\mathbf{x} - \boldsymbol{\mu}) = \nabla_{\boldsymbol{\mu}}(\mathbf{x}^\top - \boldsymbol{\mu}^\top) \Sigma^{-1}(\mathbf{x} - \boldsymbol{\mu}) \quad (1)$$

$$= \nabla_{\boldsymbol{\mu}}[\mathbf{x}^\top \Sigma^{-1} \mathbf{x} - \mathbf{x}^\top \Sigma^{-1} \boldsymbol{\mu} - \boldsymbol{\mu}^\top \Sigma^{-1} \mathbf{x} + \boldsymbol{\mu}^\top \Sigma^{-1} \boldsymbol{\mu}] \quad (2)$$

$$= \nabla_{\boldsymbol{\mu}}[-\mathbf{x}^\top \Sigma^{-1} \boldsymbol{\mu} - \boldsymbol{\mu}^\top \Sigma^{-1} \mathbf{x} + \boldsymbol{\mu}^\top \Sigma^{-1} \boldsymbol{\mu}] \quad (3)$$

$$= \nabla_{\boldsymbol{\mu}}[-\mathbf{x}^\top \Sigma^{-1} \boldsymbol{\mu} - (\boldsymbol{\mu}^\top \Sigma^{-1} \mathbf{x})^\top + \boldsymbol{\mu}^\top \Sigma^{-1} \boldsymbol{\mu}] \quad (4)$$

$$= \nabla_{\boldsymbol{\mu}}[-2\mathbf{x}^\top \Sigma^{-1} \boldsymbol{\mu} + \boldsymbol{\mu}^\top \Sigma^{-1} \boldsymbol{\mu}] \quad (5)$$

$$= \nabla_{\boldsymbol{\mu}}(-2\mathbf{x}^\top \Sigma^{-1} \boldsymbol{\mu}) + \nabla_{\boldsymbol{\mu}} \boldsymbol{\mu}^\top \Sigma^{-1} \boldsymbol{\mu} \quad (6)$$

$$= -2\mathbf{x}^\top \Sigma^{-1} + \boldsymbol{\mu}^\top (\Sigma^{-1} + (\Sigma^{-1})^\top) \quad (7)$$

$$= -2\mathbf{x}^\top \Sigma^{-1} + 2\boldsymbol{\mu}^\top \Sigma^{-1} \quad (8)$$

$$= 2(\boldsymbol{\mu} - \mathbf{x})^\top \Sigma^{-1} \quad (9)$$

Since $\mathbf{x}^\top \Sigma^{-1} \mathbf{x}$ does not depend on $\boldsymbol{\mu}$, it was left out of Eq. (3). In Eq. (5) and (8) we used $(\Sigma^{-1})^\top = \Sigma^{-1}$. The identities $\frac{\partial}{\partial \mathbf{x}} \mathbf{a}^\top \mathbf{x} = \mathbf{a}^\top$ and $\frac{\partial}{\partial \mathbf{x}} \mathbf{x}^\top \mathbf{B} \mathbf{x} = \mathbf{x}^\top (\mathbf{B} + \mathbf{B}^\top)$, as given in section 5.5 of “Mathematics for Machine Learning”, are used in Eq. (7).

2.1 b) Assuming $\log(\cdot) = \ln(\cdot)$ (not clearly specified in the exercise):

$$\nabla_{\mathbf{q}} - \mathbf{p}^\top \log \mathbf{q} = -\mathbf{p}^\top \nabla_{\mathbf{q}} \log \mathbf{q} \quad (10)$$

$$\left[\frac{\partial \log \mathbf{q}}{\partial \mathbf{q}} \right]_{i,j} = \frac{\partial \log q_i}{\partial q_j} = \frac{\delta_{i,j}}{q_i}, \quad (11)$$

where δ is the Kronecker delta. We then get the answer:

$$\nabla_{\mathbf{q}} - \mathbf{p}^\top \log \mathbf{q} = -\mathbf{p}^\top \begin{bmatrix} 1/q_1 & 0 & \cdots & 0 \\ 0 & 1/q_2 & \ddots & \vdots \\ \vdots & \ddots & \ddots & 0 \\ 0 & \cdots & 0 & 1/q_N \end{bmatrix}, \quad (12)$$

for $N \in \mathbb{N}$.

2.1 c) $\mathbf{W} \in \mathbb{R}^{2 \times 3}$ and $\mathbf{x} \in \mathbb{R}^3$, so $\mathbf{f} \in \mathbb{R}^2$.

This means we should get a final answer $\frac{\partial}{\partial \mathbf{W}} \mathbf{f} \in \mathbb{R}^{2 \times (3 \times 2)}$.

$$\frac{\partial}{\partial \mathbf{W}} \mathbf{f} = \begin{bmatrix} \frac{\partial}{\partial \mathbf{W}} f_1 \\ \frac{\partial}{\partial \mathbf{W}} f_2 \end{bmatrix} \quad (13)$$

$$\left[\frac{\partial f_i}{\partial \mathbf{W}} \right]_{1,j,k} = \frac{\partial f_i}{\partial W_{k,j}} \quad (14)$$

$$= \frac{\partial}{\partial W_{k,j}} W_{i,1}x_1 + W_{i,2}x_2 + W_{i,3}x_3. \quad (15)$$

Working this out gives a matrix of which the i 'th column is equal to \mathbf{x} , and the rest of the elements are zero. For our final answer we thus get:

$$\frac{\partial \mathbf{f}}{\partial \mathbf{W}} = \begin{bmatrix} \begin{bmatrix} x_1 & 0 \\ x_2 & 0 \\ x_3 & 0 \end{bmatrix} \\ \begin{bmatrix} 0 & x_1 \\ 0 & x_2 \\ 0 & x_3 \end{bmatrix} \end{bmatrix} \quad (16)$$

2.1 d)

$$f = (\boldsymbol{\mu} - \mathbf{W}\mathbf{x})^\top \Sigma^{-1} (\boldsymbol{\mu} - \mathbf{W}\mathbf{x}) \quad (17)$$

$$= (\boldsymbol{\mu} - \mathbf{y})^\top \Sigma^{-1} (\boldsymbol{\mu} - \mathbf{y}), \quad (18)$$

where $\mathbf{y} := \mathbf{W}\mathbf{x} \in \mathbb{R}^M$. We can then apply the chain rule:

$$\frac{\partial f}{\partial \mathbf{W}} = \frac{\partial f}{\partial \mathbf{y}} \frac{\partial \mathbf{y}}{\partial \mathbf{W}} \quad (19)$$

From similarity with exercise 2.1a we can see that

$$\frac{\partial f}{\partial \mathbf{y}} = 2(\mathbf{y} - \boldsymbol{\mu})^\top \Sigma^{-1}. \quad (20)$$

The last part is given by solving $\frac{\partial}{\partial \mathbf{W}} \mathbf{y}$, which is similar to exercise 2.1e. However now, instead of fixed matrix dimensions we have

$$\frac{\partial \mathbf{y}}{\partial \mathbf{W}} \in \mathbb{R}^{M \times (K \times M)}, \quad (21)$$

so we get

$$\frac{\partial \mathbf{y}}{\partial \mathbf{W}} = \begin{bmatrix} \begin{bmatrix} x_1 & 0 & \cdots & \cdots & 0 \\ \vdots & \vdots & & & \vdots \\ x_K & 0 & \cdots & \cdots & 0 \end{bmatrix} \\ \begin{bmatrix} 0 & x_1 & 0 & \cdots & 0 \\ \vdots & \vdots & \vdots & & \vdots \\ 0 & x_K & 0 & \cdots & 0 \end{bmatrix} \\ \vdots \\ \begin{bmatrix} 0 & \cdots & \cdots & 0 & x_1 \\ \vdots & & & \vdots & \vdots \\ 0 & \cdots & \cdots & 0 & x_K \end{bmatrix} \end{bmatrix}, \quad (22)$$

with M matrices on the vertical, each consisting of K rows and M columns. Combining Eqs. (20) and (22) gives us

$$\frac{\partial f}{\partial \mathbf{W}} = 2(\mathbf{W}\mathbf{x} - \boldsymbol{\mu})^\top \Sigma^{-1} \begin{bmatrix} [\mathbf{x} & \mathbf{0} & \cdots & \cdots & \mathbf{0}] \\ [\mathbf{0} & \mathbf{x} & \mathbf{0} & \cdots & \mathbf{0}] \\ & & \vdots & & \\ [\mathbf{0} & \cdots & \cdots & \mathbf{0} & \mathbf{x}] \end{bmatrix} \in \mathbb{R}^{1 \times (K \times M)}, \quad (23)$$

in which $\mathbf{x} = [x_1, x_2, \dots, x_K]^\top$ and $\mathbf{0} = [0, \dots, 0]^\top \in \mathbb{R}^M$. (Sorry, I was struggling to get the matrices to align properly.)

3 Probability Theory

Question 3.1

3.1 a) Most of the times when people experience a man climbing through a broken window a jewelry store with a bag over his shoulder, either on television or in stories, the man in question is a criminal. We learn from experience, so a man with this exact description would be classified by many people as a criminal.

3.1 b) Say, c stands for “the man being a criminal” and o stands for the described observation. The chance of the man being a criminal given our observation is given by Bayes’ theorem:

$$p(c|o) = \frac{p(o|c)p(c)}{p(o)} \quad (24)$$

3.1 c)

$$p(c) = 10^{-5} \quad (25)$$

$$p(o|c) = 0.8 \quad (26)$$

$$p(o|\neg c) = 10^{-6} \quad (27)$$

$$p(o) = p(o|c)p(c) + p(o|\neg c)p(\neg c) \quad (28)$$

$$p(\neg c) = 1 - p(c) \quad (29)$$

So, by combining Eqs. (24), (28) and (29) and using the given values gives:

$$p(c|o) = \frac{p(o|c)p(c)}{p(o|c)p(c) + p(o|\neg c)(1 - p(c))} \quad (30)$$

$$= \left(1 + \frac{p(o|\neg c)(1 - p(c))}{p(o|c)p(c)} \right)^{-1} \quad (31)$$

$$= 0.889 \quad (32)$$

$$\approx \frac{8}{9} \quad (33)$$

3.1 d) If some kids smashed multiple storefronts in the neighbourhood, then the man in question might be the jewelry store owner bringing his belongings to safety. In other words,

$p(o|\neg c)$ increases, although it is difficult to say how much. Say, it increases by a factor 10, i.e. $p(o|\neg c) = 10^{-5}$. Then Eq. (30) evaluates to $p(c|o) = 0.444 \approx \frac{4}{9}$. So in this case we would believe the man is (most probably) not a criminal, given the observation.

Question 3.2

3.2 a) The n 'th observation is $\mathbf{x}_n \in \mathbb{R}^4$. \mathcal{D} is the data set consisting of N independent observations $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N$, i.e. $\mathcal{D} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N\}$. The likelihood function for the observed data set is (Bishop Eq. 2.29):

$$p(\mathcal{D}|\boldsymbol{\rho}) = \prod_{k=1}^4 \rho_k^{m_k}, \quad (34)$$

where m_k is given by

$$m_k = \sum_{n=1}^N x_{n,k} \quad (35)$$

3.2 b) Maximum likelihood solution for $\boldsymbol{\rho}$:

$$\boldsymbol{\rho} = [1/2, 0, 0, 1/2] \quad (36)$$

This is based on the fact that half of the draws are hearts (i.e. $\rho_1 = 1$) and the other half are spades (i.e. $\rho_4 = 1$).

3.2 c) Define $r_n := x_{n,1} + x_{n,3}$, then

$$\text{if } r_n = 1 \rightarrow \text{the } n\text{'th observation is a red card} \quad (37)$$

$$\text{if } r_n = 0 \rightarrow \text{the } n\text{'th observation is a black card} \quad (38)$$

To avoid confusion, a capital P will now be used in $P(X)$ to denote the probability of an event X , instead of $p(X)$. The likelihood function for the card colors is (Bishop Eq. 2.5):

$$P(\mathcal{D}|p) = \prod_{n=1}^N P(r_n|p) = \prod_{n=1}^N p^{r_n} (1-p)^{1-r_n} \quad (39)$$

3.2 d) Using Bishop Eq. 2.7:

$$p_{ML} = \frac{1}{N} \sum_{n=1}^N r_n = \frac{6}{8} = 0.75 \quad (40)$$

3.2 e)

$$p = \rho_1 + \rho_3 \quad (41)$$

3.2 f)

$$\underbrace{P(p|\mathcal{D})}_{\text{posterior}} = \frac{\overbrace{P(\mathcal{D}|p)}^{\text{likelihood}} \overbrace{P(p)}^{\text{prior}}}{\underbrace{P(\mathcal{D})}_{\text{evidence}}} \quad (42)$$

3.2 g) We have

$$P(\mathcal{D}|p) = \prod_{n=1}^N p^{r_n} (1-p)^{1-r_n} \quad (43)$$

$$P(p) = \text{Beta}(p|\alpha, \beta) = \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} p^{\alpha-1} (1-p)^{\beta-1} \quad (44)$$

The log likelihood is

$$\mathcal{L} = \log(P(p|\mathcal{D})) \quad (45)$$

The MAP estimate p_{MAP} is given by

$$p_{\text{MAP}} = \arg \max_p \mathcal{L} \quad (46)$$

$$= \arg \max_p \log(P(\mathcal{D}|p)(p)) \quad (47)$$

$$= \arg \max_p \underbrace{\sum_{n=1}^N \log(p^{r_n} (1-p)^{1-r_n})}_{:= A} + \underbrace{\log\left(\frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} p^{\alpha-1} (1-p)^{\beta-1}\right)}_{:= B} \quad (48)$$

To find p_{MAP} , we need to set $\frac{\partial}{\partial p}(A + B) = \frac{\partial A}{\partial p} + \frac{\partial B}{\partial p} = 0$. For term A we have:

$$\frac{\partial A}{\partial p} = \sum_{n=1}^N \frac{\partial}{\partial p} \log(p^{r_n} (1-p)^{1-r_n}) \quad (49)$$

$$y := p^{r_n} (1-p)^{1-r_n} \quad (50)$$

$$\frac{\partial \log y}{\partial p} = \frac{\partial \log y}{\partial y} \frac{\partial y}{\partial p} \quad (51)$$

$$= \frac{1}{y} [r_n p^{r_n-1} (1-p)^{1-r_n} + p^{r_n} (r_n - 1) (1-p)^{-r_n}] \quad (52)$$

$$\frac{\partial A}{\partial p} = \sum_{n=1}^N \frac{1}{p^{r_n} (1-p)^{1-r_n}} [r_n p^{r_n} p^{-1} (1-p)^{1-r_n} + p^{r_n} (r_n - 1) \frac{(1-p)^{1-r_n}}{1-p}] \quad (53)$$

$$= \sum_{n=1}^N \left[\frac{r_n}{p} + \frac{r_n - 1}{1-p} \right] \quad (54)$$

$$= \frac{1}{p} \sum_{n=1}^N r_n + \frac{1}{1-p} \sum_{n=1}^N (r_n - 1) \quad (55)$$

For term B we have:

$$\frac{\partial B}{\partial p} = \frac{\partial}{\partial p} \log\left(\frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)}\right) + \frac{\partial}{\partial p} \log(p^{\alpha-1} (1-p)^{\beta-1}) \quad (56)$$

$$= \frac{\partial}{\partial p} \log(p^{\alpha-1} (1-p)^{\beta-1}) \quad (57)$$

$$= \frac{1}{p^{\alpha-1} (1-p)^{\beta-1}} [(\alpha-1)p^{\alpha-2} (1-p)^{\beta-1} + p^{\alpha-1} (\beta-1) (1-p)^{\beta-2}] \quad (58)$$

$$= \frac{(\alpha-1)p^{\alpha-1} p^{-1} (1-p)^{\beta-1}}{p^{\alpha-1} (1-p)^{\beta-1}} + \frac{p^{\alpha-1} (1-\beta) (1-p)^{\beta-1} (1-p)^{-1}}{p^{\alpha-1} (1-p)^{\beta-1}} \quad (59)$$

$$= \frac{\alpha-1}{p} + \frac{1-\beta}{1-p} \quad (60)$$

Setting $\frac{\partial A}{\partial p} + \frac{\partial B}{\partial p} = 0$ gives

$$0 = \frac{1}{p_{\text{MAP}}} \sum_{n=1}^N r_n + \frac{1}{1-p_{\text{MAP}}} \sum_{n=1}^N (r_n - 1) + \frac{\alpha - 1}{p_{\text{MAP}}} + \frac{1 - \beta}{1 - p_{\text{MAP}}} \quad (61)$$

$$= \sum_{n=1}^N r_n + \frac{p_{\text{MAP}}}{1 - p_{\text{MAP}}} \sum_{n=1}^N (r_n - 1) + \alpha - 1 + \frac{p_{\text{MAP}}(1 - \beta)}{1 - p_{\text{MAP}}} \quad (62)$$

$$p_{\text{MAP}} \left[\sum_{n=1}^N (r_n - 1) + 1 - \beta \right] = (1 - p_{\text{MAP}}) \left[1 - \alpha - \sum_{n=1}^N r_n \right] \quad (63)$$

$$= \left[1 - \alpha - \sum_{n=1}^N r_n \right] - p_{\text{MAP}} \left[1 - \alpha - \sum_{n=1}^N r_n \right] \quad (64)$$

$$p_{\text{MAP}} \left[\sum_{n=1}^N (r_n - 1) - \sum_{n=1}^N r_n + 2 - \beta - \alpha \right] = 1 - \alpha - \sum_{n=1}^N r_n \quad (65)$$

$$p_{\text{MAP}} = \frac{1 - \alpha - \sum_{n=1}^N r_n}{2 - N - \beta - \alpha} = \frac{\sum_{n=1}^N r_n + \alpha - 1}{N + \alpha + \beta - 2} \quad (66)$$

To encode the belief that there is an equal probability of drawing a red or black card, we would choose $\alpha = \beta$. The greater their values, the stronger the belief that $p \approx 0.5$.