

1 Multivariate Calculus

1.1 First Subtask

a)

Suppose we have

$$\mathbf{y} = \mathbf{x} - \boldsymbol{\mu}$$

Then we can simplify the original form to:

$$\nabla_{\boldsymbol{\mu}} \mathbf{y}^T \boldsymbol{\Sigma}^{-1} \mathbf{y}$$

According to chain rule:

$$\nabla_{\boldsymbol{\mu}} \mathbf{y}^T \boldsymbol{\Sigma}^{-1} \mathbf{y} = (\nabla_{\mathbf{y}} \mathbf{y}^T \boldsymbol{\Sigma}^{-1} \mathbf{y}) \nabla_{\boldsymbol{\mu}} \mathbf{y}$$

And we can derive the first part of the above form by:

$$\begin{aligned} & \nabla_{\mathbf{y}} \mathbf{y}^T \boldsymbol{\Sigma}^{-1} \mathbf{y} \\ &= \mathbf{y}^T \nabla_{\mathbf{y}} (\boldsymbol{\Sigma}^{-1} \mathbf{y}) + \nabla_{\mathbf{y}} \mathbf{y}^T (\boldsymbol{\Sigma}^{-1} \mathbf{y}) \\ &= 2 \mathbf{y}^T \boldsymbol{\Sigma}^{-1} \end{aligned}$$

by plugging in \mathbf{y} and $\nabla_{\boldsymbol{\mu}} \mathbf{y}$ we have the final solution:

$$\nabla_{\boldsymbol{\mu}} \mathbf{y}^T \boldsymbol{\Sigma}^{-1} \mathbf{y} = -2(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}$$

b)

It is easy to see that the function is real-valued function. So, $\frac{\delta f}{\delta \mathbf{q}} \in \mathbb{R}^{1 \times n}$

$$\left[\frac{\delta f}{\delta q_1}, \frac{\delta f}{\delta q_2}, \dots, \frac{\delta f}{\delta q_j}, \frac{\delta f}{\delta q_n} \right] = \left[\frac{p_1}{q_1}, \frac{p_2}{q_2}, \dots, \frac{p_j}{q_j}, \frac{p_n}{q_n} \right]$$

c)

We have

$$\mathbf{f} = \mathbf{W} \mathbf{x}, \quad \mathbf{f} \in \mathbb{R}^2, \quad \mathbf{W} \in \mathbb{R}^{2 \times 3}, \quad \mathbf{x} \in \mathbb{R}^3$$

So the derivative is of shape:

$$\frac{d\mathbf{f}}{d\mathbf{w}} \in \mathbb{R}^{2 \times (2 \times 3)}$$

Then, we have:

$$\frac{d\mathbf{f}}{d\mathbf{w}} = \begin{bmatrix} \frac{df_1}{dw} \\ \frac{df_2}{dw} \end{bmatrix}, \quad \frac{df_i}{dw} \in \mathbb{R}^{1 \times (2 \times 3)}$$

If, we expand the matrices and apply elementwise calculation, we should have the partial derivative:

$$\frac{\delta f_i}{\delta W_{iq}} = x_q$$

This allows us to compute the partial derivatives of f_i with respect to a row of W , which is given as

$$\frac{\delta f_i}{\delta W_{i,(1,2,3)}} = \mathbf{x}^T \in \mathbb{R}^{1 \times (1 \times 3)}$$

And for other rows:

$$\frac{\delta f_i}{\delta W_{k \neq i,(1,2,3)}} = \mathbf{0}^T \in \mathbb{R}^{1 \times (1 \times 3)}$$

So finally, we have:

$$\begin{aligned} \frac{\delta f_1}{\delta W} &= \begin{bmatrix} x_1 & x_2 & x_3 \\ 0 & 0 & 0 \end{bmatrix} \\ \frac{\delta f_2}{\delta W} &= \begin{bmatrix} 0 & 0 & 0 \\ x_1 & x_2 & x_3 \end{bmatrix} \end{aligned}$$

d)

We have:

$$W \in \mathbb{R}^{m \times k}, \quad x \in \mathbb{R}^{k \times 1}, \quad f \in \mathbb{R}, \quad \nabla_W f \in \mathbb{R}^{m \times k}$$

We can use the solutions derived above to calculate:

$$\nabla_W f = 2(\mu - Wx)^T \Sigma^{-1} \nabla_W(Wx)$$

Where we have $\nabla_W(Wx)$:

$$\nabla_W(Wx) = \begin{bmatrix} \frac{df_1}{dw} \\ \frac{df_2}{dw} \\ \dots \\ \frac{df_m}{dw} \end{bmatrix}, \quad \frac{df_i}{dw} \in \mathbb{R}^{1 \times (m \times k)}$$

Where

$$\frac{df_i}{dW} = \begin{bmatrix} \mathbf{0}^T \\ \dots \\ \mathbf{0}^T \\ x^T \\ \mathbf{0}^T \\ \dots \\ \mathbf{0}^T \end{bmatrix}$$

2 Probability Theory

2.1 First Subtask

a)

Even though criminal rate in population is relatively low (the prior), based on the observation, the chance of seeing a criminal under such situation (posterior) is very high.

b)

We have:

$$P(\text{criminal}) = \frac{1}{10^5}$$

$$P(\text{observedSituation}|\text{criminal}) = 0.8$$

$$P(\text{observedSituation}|\text{notCriminal}) = \frac{1}{10^6}$$

With this, according to the sum rule, we can calculate the probability of observing such situation:

$$\begin{aligned} P(\text{observedSituation}) &= P(\text{observedSituation}|\text{criminal})P(\text{criminal}) + \\ &\quad P(\text{observedSituation}|\text{notCriminal})P(\text{notCriminal}) \\ &= \frac{0.8}{10^5} + \frac{1 - \frac{1}{10^5}}{10^6} \end{aligned}$$

According to Bayes:

$$P(\text{criminal}|\text{observedSituation}) = \frac{P(\text{observedSituation}|\text{criminal})P(\text{criminal})}{P(\text{observedSituation})} = 0.89$$

c)

See solution is b).

d)

With the new evidence, the probability of making this observation becomes much higher, thus according to Bayes rule as shown in b), the belief in the man in question being a criminal lowers down. the probability of making this observation when the man is not a criminal is 1

2.2 Second Subtask

a)

Assume ρ satisfies i.i.d., and each observation also satisfies i.i.d., we have:

$$P(D|\rho) = \prod_{n=1}^N \prod_{i=1}^4 (\rho)^{x_{ni}}$$

b)

Since we have 4 out of 8 ♥ and also 4 out of 8 ♠, and none of other suits in the observations. The max likelihood solution is as follows:

$$\rho = [1/2, 0, 0, 1/2]$$

c)

Suppose we have M red observations in N:

$$P_{ML}(D|p) = p^M (1-p)^{N-M}$$

d)

Since we have 6 out of 8 in red, so the solution that maximize likelihood is 3/4.

e)

$$p = \frac{\rho_1 + \rho_3}{\rho_1 + \rho_2 + \rho_3 + \rho_4}$$

f)

$$P_{\text{posterior}}(p_{\text{bernoulli}}|D) = \frac{P_{\text{likelihood}}(D|p_{\text{bernoulli}})p_{\text{bernoulli}}}{P_{\text{evidence}}(D)}$$

Where $p_{\text{bernoulli}}$ is the prior, posterior, likelihood and evidence are pointed out in the form.

g)

$$p_{map} = \operatorname{argmax}_p P(D|p)P(p) = \operatorname{argmax}_p p^M (1-p)^{N-M} p^{\alpha-1} (1-p)^{\beta-1}$$

We find p such that the derivative of posterior:

$$\frac{d(p^M (1-p)^{N-M} p^{\alpha-1} (1-p)^{\beta-1})}{dp} = -p(N-M+\beta-1) + (M+\alpha-1)(1-p) = 0$$

Finally we have:

$$p = \frac{1 - \alpha - M}{2 - N - \beta - \alpha}$$