

1) (Chapter 28) Consider an unknown  $M$ -dimensional vector  $\mathbf{x} \sim f_{\mathbf{x}}(x)$  that we wish to infer from a collection of  $N$  independent and identically distributed observations  $\{\mathbf{y}_n, n = 1, 2, \dots, N\}$ . The conditional distribution of each observation  $\mathbf{y}_n$  given  $\mathbf{x}$  is uniform across all observations and denoted by  $\mathbf{y}_n | \mathbf{x} \sim f_{\mathbf{y}|\mathbf{x}}(y_n | x)$ . Show that the MAP estimator for  $\mathbf{x}$  given the  $N$  observations  $\{\mathbf{y}_1, \dots, \mathbf{y}_N\}$  can be found by solving

$$\hat{x}_{\text{MAP}} = \underset{x \in \mathbb{R}^M}{\operatorname{argmin}} \left\{ -\frac{1}{N} \sum_{n=1}^N \ln f_{\mathbf{y}|\mathbf{x}}(y_n | x) - \frac{1}{N} \ln f_{\mathbf{x}}(x) \right\}$$

**Solution:** Using the iid assumption we have that

$$f_{\{\mathbf{y}_n\}_{n=1}^N | \mathbf{x}}(\{\mathbf{y}_n\}_{n=1}^N | x) = \prod_{n=1}^N f_{\mathbf{y}_n | \mathbf{x}}(y_n | x) = \prod_{n=1}^N f_{\mathbf{y}|\mathbf{x}}(y_n | x)$$

Using Bayes' rule we have

$$\begin{aligned} \hat{x}_{\text{MAP}} &= \underset{x \in \mathbb{R}^M}{\operatorname{argmax}} f_{\mathbf{x} | \{\mathbf{y}_n\}_{n=1}^N}(x | \{\mathbf{y}_n\}_{n=1}^N) \\ &= \underset{x \in \mathbb{R}^M}{\operatorname{argmax}} \left\{ \frac{f_{\{\mathbf{y}_n\}_{n=1}^N | \mathbf{x}}(\{\mathbf{y}_n\}_{n=1}^N | x) \times f_{\mathbf{x}}(x)}{f_{\{\mathbf{y}_n\}_{n=1}^N}(\{\mathbf{y}_n\}_{n=1}^N)} \right\} \\ &= \underset{x \in \mathbb{R}^M}{\operatorname{argmax}} \left\{ f_{\{\mathbf{y}_n\}_{n=1}^N | \mathbf{x}}(\{\mathbf{y}_n\}_{n=1}^N | x) \times f_{\mathbf{x}}(x) \right\} \\ &= \underset{x \in \mathbb{R}^M}{\operatorname{argmax}} \left\{ \left( \prod_{n=1}^N f_{\mathbf{y}|\mathbf{x}}(y_n | x) \right) \times f_{\mathbf{x}}(x) \right\} \end{aligned}$$

Since the logarithm function is monotonically increasing, we can also write

$$\begin{aligned} \hat{x}_{\text{MAP}} &= \underset{x \in \mathbb{R}^M}{\operatorname{argmax}} \ln \left[ \left( \prod_{n=1}^N f_{\mathbf{y}|\mathbf{x}}(y_n | x) \right) \times f_{\mathbf{x}}(x) \right] \\ &= \underset{x \in \mathbb{R}^M}{\operatorname{argmax}} \left\{ \sum_{n=1}^N \ln f_{\mathbf{y}|\mathbf{x}}(y_n | x) + \ln f_{\mathbf{x}}(x) \right\} \\ &= \underset{x \in \mathbb{R}^M}{\operatorname{argmin}} \left\{ - \sum_{n=1}^N \ln f_{\mathbf{y}|\mathbf{x}}(y_n | x) - \ln f_{\mathbf{x}}(x) \right\} \end{aligned}$$

Normalizing by the sample size we conclude that

$$\hat{x}_{\text{MAP}} = \underset{x \in \mathbb{R}^M}{\operatorname{argmin}} \left\{ -\frac{1}{N} \sum_{n=1}^N \ln f_{\mathbf{y}|\mathbf{x}}(y_n | x) - \frac{1}{N} \ln f_{\mathbf{x}}(x) \right\}$$

□

2) (Chapter 29) We examined the correlation coefficient  $\rho_{xy}$  between two scalar random variables  $\{x, y\}$  in Prob. 3.13 in the text. The correlation between these random variables given a third scalar variable  $z$  is defined as follows. Let  $\tilde{\mathbf{x}} = \mathbf{x} - \hat{\mathbf{x}}$  denote the residual that remains from estimating  $\mathbf{x}$  from  $z$  using a linear regression model, say,  $\hat{\mathbf{x}} = \alpha z + \beta$ . Likewise, let  $\tilde{\mathbf{y}} = \mathbf{y} - \hat{\mathbf{y}}$  denote the residual

that remains from estimating  $\mathbf{y}$  from  $\mathbf{z}$  using a linear regression model as well, say,  $\hat{\mathbf{y}} = a\mathbf{z} + b$ . Then, the (conditional) correlation is defined by

$$\rho_{xy|z} \triangleq \rho_{\tilde{x}, \tilde{y}}$$

That is, it is equal to the standard correlation between the residual variables after the effect of  $\mathbf{z}$  has been removed using linear regression. Show that if  $\mathbf{x}$  and  $\mathbf{y}$  are conditionally independent of each other given  $\mathbf{z}$ , then  $\rho_{xy|z} = 0$ . Is the converse true?

**Solution:** We know that

$$\begin{aligned}\hat{\mathbf{x}} - \bar{x} &= \frac{\sigma_{xz}}{\sigma_z^2}(z - \bar{z}) \\ \hat{\mathbf{y}} - \bar{y} &= \frac{\sigma_{yz}}{\sigma_z^2}(z - \bar{z})\end{aligned}$$

so that the scalars  $(\alpha, \beta)$  and  $(a, b)$  are given by

$$\begin{aligned}\alpha &= \frac{\sigma_{xz}}{\sigma_z^2}, & \beta &= \bar{x} - \frac{\sigma_{xz}}{\sigma_z^2}\bar{z} \\ a &= \frac{\sigma_{yz}}{\sigma_z^2}, & b &= \bar{y} - \frac{\sigma_{yz}}{\sigma_z^2}\bar{z}\end{aligned}$$

Now given  $\mathbf{z}$ , the regression error  $\tilde{\mathbf{x}}$  is solely dependent on  $\mathbf{x}$  while  $\tilde{\mathbf{y}}$  is solely dependent on  $\mathbf{y}$  since

$$\begin{aligned}\tilde{\mathbf{x}} &= \mathbf{x} - \alpha\mathbf{z} - \beta \\ \tilde{\mathbf{y}} &= \mathbf{y} - a\mathbf{z} - b\end{aligned}$$

It follows that  $\tilde{\mathbf{x}}$  and  $\tilde{\mathbf{y}}$  are also conditionally independent of each other given  $\mathbf{z}$ , from which we conclude that  $\rho_{\tilde{x}, \tilde{y}} = 0$  and, therefore,  $\rho_{xy|z} = 0$ .

The converse is not true. If  $\rho_{xy|z} = 0$  then  $\rho_{\tilde{x}, \tilde{y}} = 0$ . However, we know that this latter condition does not imply that  $\tilde{\mathbf{x}}$  and  $\tilde{\mathbf{y}}$  are independent of each other. Now given by  $\mathbf{x}$  and  $\mathbf{y}$  can be expressed in terms of these errors as

$$\begin{aligned}\mathbf{x} &= \tilde{\mathbf{x}} + \alpha\mathbf{z} + \beta \\ \mathbf{y} &= \tilde{\mathbf{y}} + a\mathbf{z} + b\end{aligned}$$

we conclude that  $\mathbf{x}$  and  $\mathbf{y}$  are not necessarily independent of each other given  $\mathbf{z}$ . □

3) **(Chapter 31)** Consider  $N$  i.i.d. scalar random variables  $\{\mathbf{y}_n\}$  distributed according to  $\mathbf{y} \sim f_{\mathbf{y}}(y) = \frac{1}{2\sigma}e^{-|y|/\sigma}$ . Show that the ML estimate for  $\sigma$  is given by

$$\hat{\sigma} = \frac{1}{N} \sum_{n=1}^N |y_n|$$

Verify that this estimator is unbiased and asymptotically efficient.

**Solution:** The log-likelihood function is given by

$$\begin{aligned}
 \ell(\{y_n\}; \sigma) &= \ln \left[ \prod_{n=1}^N \frac{1}{2\sigma} e^{-|y_n|/\sigma} \right] \\
 &= \ln \left[ \left( \frac{1}{2\sigma} \right)^N \exp \left\{ - \sum_{n=1}^N \frac{|y_n|}{\sigma} \right\} \right] \\
 &= -N \ln(2\sigma) - \sum_{n=1}^N \frac{|y_n|}{\sigma}
 \end{aligned}$$

Differentiating relative to  $\sigma$  leads to

$$-N \frac{1}{\hat{\sigma}} + \sum_{n=1}^N \frac{|y_n|}{\hat{\sigma}^2} = 0$$

leading to

$$\hat{\sigma} = \frac{1}{N} \sum_{n=1}^N |y_n|$$

Next note that

$$\begin{aligned}
 \mathbb{E}|\mathbf{y}| &= \int_{-\infty}^{\infty} |y| f_{\mathbf{y}}(y) dy \\
 &= \frac{1}{2\sigma} \int_{-\infty}^{\infty} |y| e^{-|y|/\sigma} dy \\
 &= \frac{1}{2\sigma} \left[ \int_{-\infty}^0 -ye^{y/\sigma} dy + \int_0^{\infty} ye^{-y/\sigma} dy \right] \\
 &= \frac{1}{\sigma} \int_0^{\infty} ye^{-y/\sigma} dy \\
 &= \int_0^{\infty} \frac{y}{\sigma} e^{-y/\sigma} dy \\
 &= \sigma \int_0^{\infty} y' e^{-y'} dy', \quad y' = y/\sigma \\
 &= \sigma \times \Gamma(2) \\
 &= \sigma
 \end{aligned}$$

where we used the value of the gamma function at 2, which is equal to 1. It follows from the observation that  $\Gamma(n+1) = n!$  for integer values of  $n$ . Therefore, the proposed estimator is unbiased since

$$\mathbb{E} \hat{\sigma} = \frac{1}{N} \sum_{n=1}^N \mathbb{E} |y_n| = \sigma$$

Similarly, note that

$$\begin{aligned}
\mathbb{E}|\mathbf{y}|^2 &= \int_{-\infty}^{\infty} |y|^2 f_{\mathbf{y}}(y) dy \\
&= \frac{1}{2\sigma} \int_{-\infty}^{\infty} |y|^2 e^{-|y|/\sigma} dy \\
&= \frac{1}{2\sigma} \left[ \int_{-\infty}^0 y^2 e^{y/\sigma} dy + \int_0^{\infty} y^2 e^{-y/\sigma} dy \right] \\
&= \frac{1}{\sigma} \int_0^{\infty} y^2 e^{-y/\sigma} dy \\
&= \sigma \int_0^{\infty} \frac{y}{\sigma^2} e^{-y/\sigma} dy \\
&= \sigma^2 \int_0^{\infty} z^2 e^{-z} dz, \quad z = y/\sigma \\
&= \sigma^2 \times \Gamma(3) \\
&= \sigma^2 \times 2! \\
&= 2\sigma^2
\end{aligned}$$

Consequently, the mean-square-error of  $\hat{\sigma}$  is given by

$$\begin{aligned}
\text{MSE}(\hat{\sigma}) &= \mathbb{E}(\hat{\sigma} - \sigma)^2 \\
&= \text{var}(\hat{\sigma}) \\
&= \frac{1}{N^2} \sum_{n=1}^N \text{var}(|\mathbf{y}_n|) \\
&= \frac{1}{N^2} \sum_{n=1}^N \left( \mathbb{E}|\mathbf{y}_n|^2 - [\mathbb{E}|\mathbf{y}_n|]^2 \right) \\
&= \frac{1}{N^2} \sum_{n=1}^N \left( 2\sigma^2 - \sigma^2 \right) \\
&= \frac{\sigma^2}{N}
\end{aligned}$$

Observe that the variance tends to 0 as  $N \rightarrow \infty$  and, therefore, the ML estimator is asymptotically efficient.  $\square$

4) (Chapter 50) Consider two  $N \times M$  matrices  $D$  and  $H$ . Introduce the SVD representation

$$D^T H = U \Sigma V^T$$

where  $U$  and  $V$  are  $M \times M$  orthogonal, and  $\Sigma$  is diagonal. Show that the  $M \times M$  orthogonal matrix  $\Theta^*$  that solves

$$\Theta^* = \underset{\Theta \Theta^T = I_M}{\operatorname{argmin}} \|D - H\Theta\|_F$$

is given by  $\Theta^* = VU^T$ .

**Solution:** Using

$$\|A\|_F^2 = \text{Tr}(A^T A)$$

we write

$$\begin{aligned}
\|D - H\Theta\|_F^2 &= \text{Tr}\left((D - H\Theta)^T(D - H\Theta)\right) \\
&= \text{Tr}(D^T D) + \text{Tr}(\Theta^T H^T H\Theta) - 2\text{Tr}(D^T H\Theta) \\
&= \text{Tr}(D^T D) + \text{Tr}(\Theta\Theta^T H^T H) - 2\text{Tr}(D^T H\Theta) \\
&= \text{Tr}(D^T D) + \text{Tr}(H^T H) - 2\text{Tr}(D^T H\Theta) \quad (\text{since } \Theta \text{ is orthogonal}) \\
&= \|D\|_F^2 + \|H\|_F^2 - 2\text{Tr}(D^T H\Theta)
\end{aligned}$$

where we used the fact that  $\text{Tr}(AB) = \text{Tr}(BA)$ . Ignoring the terms that do not depend on  $\Theta$ , the optimization problem reduces to solving

$$\Theta^* = \underset{\Theta \Theta^T = I_M}{\operatorname{argmax}} \text{Tr}(D^T H\Theta)$$

The matrix  $D^T H$  has dimensions  $M \times M$ . We introduce its SVD:

$$D^T H = U \Sigma V^T$$

where  $U$  and  $V$  are  $M \times M$  orthogonal, and  $\Sigma$  is diagonal with entries  $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_M$ . Thus,

$$\begin{aligned}
\text{Tr}(D^T H\Theta) &= \text{Tr}(U \Sigma V^T \Theta) \\
&= \text{Tr}(\Sigma V^T \Theta U) \\
&\stackrel{\Delta}{=} \text{Tr}(\Sigma X) \\
&= \sum_{m=1}^M \sigma_m X_{mm}
\end{aligned}$$

where we introduced the matrix  $X = V^T \Theta U$  and denoted its diagonal entries by  $\{X_{mm}\}$ . It is obvious that  $X$  is an orthogonal matrix as well since

$$X X^T V^T \Theta U U^T \Theta^T V = I_M$$

and, moreover,  $X$  and  $\Theta$  define each other uniquely since

$$\Theta = V X U^T$$

Now recall from Prob. ?? that  $X_{mm} \leq 1$ . Therefore, it holds that

$$\sum_{m=1}^M \sigma_m X_{mm} \leq \sum_{m=1}^M \sigma_m$$

Equality holds if, and only if,  $X_{mm} = 1$  for all  $m$ . Since  $X$  is unitary and its rows have unit norm, this is only possible when  $X$  is the identity matrix. Therefore, the optimal choices for  $X$  is  $X = I_M$  from which we conclude that

$$\Theta^* = V U^T$$

□

5) **(Chapter 51)** Consider a collection of scalar measurements  $\{x(n)\}_{n=0}^{N-1}$  sampled independently from the Gaussian distribution  $\mathcal{N}(\theta^o, 1)$  with unknown mean  $\theta^o$ . We wish to estimate  $\theta^o$  by seeking the solution to a regularized least-squares problem with a regularization term that penalizes proximity to a prior estimate denoted by  $\theta_1$ , namely,

$$\theta^* \stackrel{\Delta}{=} \arg \min_{\theta \in \mathbb{R}} \left\{ \frac{1}{N} \sum_{n=0}^{N-1} (x(n) - \theta)^2 + \rho(\theta - \theta_1)^2 \right\}, \quad \rho > 0$$

(a) Let  $\bar{\theta} = \frac{1}{N} \sum_{n=0}^{N-1} x(n)$ . Verify that

$$\theta^* = \frac{1}{1+\rho} \bar{\theta} + \frac{\rho}{1+\rho} \theta_1$$

(b) Show that

$$\mathbb{E}(\theta^* - \theta^o)^2 = \left(\frac{1}{1+\rho}\right)^2 \frac{1}{N} + \left(\frac{\rho}{1+\rho}\right)^2 \mathbb{E}(\theta_1 - \theta^o)^2$$

(c) Verify that the optimal choice for the regularization parameter is

$$\rho = \frac{1/N}{\mathbb{E}(\theta_1 - \theta^o)^2}$$

for which the mean-square error expression in part (b) simplifies to

$$\mathbb{E}(\theta^* - \theta^o)^2 \approx \min \left\{ \frac{1}{N}, \mathbb{E}(\theta_1 - \theta^o)^2 \right\}$$

(d) Conclude that a large  $\rho$  is needed when the prior estimate has a small mean-square-error, while a small  $\rho$  is needed otherwise. In other words, conclude that the size of the regularization parameter depends on the quality of the initial estimate,  $\theta_1$ .

**Solution:**

(a) Differentiating relative to  $\theta$  and setting the gradient vector to 0 gives

$$-\frac{2}{N}(x(n) - \theta^*) + 2\rho(\theta^* - \theta_1) = 0$$

Solving for  $\theta^*$  leads to the desired expression.

(b) The MSE is given by

$$\mathbb{E}(\theta^* - \theta^o)^2 = \left(\frac{1}{1+\rho}\right)^2 \mathbb{E}(\theta^* - \theta^o)^2 + \left(\frac{\rho}{1+\rho}\right)^2 \mathbb{E}(\theta_1 - \theta^o)^2$$

where

$$\begin{aligned} \mathbb{E}(\theta^* - \theta^o)^2 &= \mathbb{E} \left| \frac{1}{N} \sum_{n=0}^{N-1} x(n) - \theta^o \right|^2 \\ &= \frac{1}{N^2} \mathbb{E} \left| \sum_{n=0}^{N-1} (x(n) - \theta^o) \right|^2 \\ &= \frac{1}{N^2} \times N \\ &= 1/N \end{aligned}$$

and the desired result follows.

(c) Differentiating the MSE expression over  $\rho$  gives

$$\frac{-2}{(1+\rho)^3} \frac{1}{N} + \frac{2\rho(1+\rho) - 2\rho^2}{(1+\rho)^3} \mathbb{E}(\theta_1 - \theta^o)^2 = 0$$

from which the desired result follows:

$$\rho = \frac{1/N}{\mathbb{E}(\theta_1 - \theta^o)^2}$$

Substituting into the MSE expression gives

$$\mathbb{E}(\boldsymbol{\theta}^* - \boldsymbol{\theta}^o)^2 = \frac{\frac{1}{N} \times \mathbb{E}(\boldsymbol{\theta}_1 - \boldsymbol{\theta}^o)^2}{\frac{1}{N} + \mathbb{E}(\boldsymbol{\theta}_1 - \boldsymbol{\theta}^o)^2}$$

When  $1/N$  is small we can approximate the above expression by

$$\mathbb{E}(\boldsymbol{\theta}^* - \boldsymbol{\theta}^o)^2 = \frac{\mathbb{E}(\boldsymbol{\theta}_1 - \boldsymbol{\theta}^o)^2}{1 + \mathbb{E}(\boldsymbol{\theta}_1 - \boldsymbol{\theta}^o)^2/(1/N)} \approx 1/N$$

On the other hand, when  $1/N$  is large we can write

$$\mathbb{E}(\boldsymbol{\theta}^* - \boldsymbol{\theta}^o)^2 = \frac{1/N}{(1/N)/\mathbb{E}(\boldsymbol{\theta}_1 - \boldsymbol{\theta}^o)^2 + 1} \approx \mathbb{E}(\boldsymbol{\theta}_1 - \boldsymbol{\theta}^o)^2$$

In other words, we can approximate the MSE expression by

$$\mathbb{E}(\boldsymbol{\theta}^* - \boldsymbol{\theta}^o)^2 \approx \min \left\{ \frac{1}{N}, \mathbb{E}(\boldsymbol{\theta}_1 - \boldsymbol{\theta}^o)^2 \right\}$$

(d) The value of the optimal  $\rho$  in part (c) has  $\mathbb{E}(\boldsymbol{\theta}_1 - \boldsymbol{\theta}^o)^2$  in the denominator, and the result follows. This means that the choice of the regularization parameter depends on the quality of the prior estimate,  $\boldsymbol{\theta}_1$ .

□

6) (Chapter 52) Consider feature vectors  $h \in \mathbb{R}^M$  and a collection of  $K$  classifiers (or experts) denoted by  $\{E_1(h), E_2(h), \dots, E_K(h)\}$ . Each feature vector  $h \in \mathbb{R}^M$  can belong to one of  $R$  classes denoted by  $r = 1, 2, \dots, R$ . Introduce an  $R \times K$  matrix  $\mathcal{E}$ , which summarizes the opinion of the experts about the class of  $h$ . Each row of index  $r$  corresponds to one of the labels, and each column of index  $c$  corresponds to one of the classifiers or experts. The entry  $\mathcal{E}_{rc}$  indicates the level of confidence that expert  $E_c$  has about feature  $h$  belonging to class  $r$ . For illustration purposes, we exhibit a matrix  $\mathcal{E}$  corresponding to  $R = 4$  labels and  $K = 5$  experts:

$$\mathcal{E} = \begin{bmatrix} \text{labels} & E_1 & E_2 & E_3 & E_4 & E_5 \\ \hline r = 1 & \mathcal{E}_{11} & \mathcal{E}_{12} & \mathcal{E}_{13} & \mathcal{E}_{14} & \mathcal{E}_{15} \\ r = 2 & \mathcal{E}_{21} & \mathcal{E}_{22} & \mathcal{E}_{23} & \mathcal{E}_{24} & \mathcal{E}_{25} \\ r = 3 & \mathcal{E}_{31} & \mathcal{E}_{32} & \mathcal{E}_{33} & \mathcal{E}_{34} & \mathcal{E}_{35} \\ r = 4 & \mathcal{E}_{41} & \mathcal{E}_{42} & \mathcal{E}_{43} & \mathcal{E}_{44} & \mathcal{E}_{45} \end{bmatrix}$$

(a) Explain how the  $K$ -NN classifier with  $R$  labels is a special case of this construction. What would the classifiers in this case be and how would the entries of  $\mathcal{E}$  be chosen?  
(b) Explain how the weighted  $K$ -NN classifier with  $R$  labels from Example 52.2 is a special case of this construction. What would the classifiers in this case be and how would the entries of  $\mathcal{E}$  be chosen?  
(c) Verify that in both cases the ultimate label for  $h$  is chosen in terms of the  $\infty$ -norm of matrix  $\mathcal{E}$  as follows:

$$r^* = \underset{1 \leq r \leq R}{\operatorname{argmax}} \|\mathcal{E}\|_\infty$$

**Solution:**

(a) For the  $K$ -NN classifier, each of the  $K$  neighbors defines one of the experts. Moreover, the value of  $\mathcal{E}_{rc}$  will be 1 if the  $c$ th neighbor belongs to class  $r$  and 0 otherwise.  
(b) For the weighted  $K$ -NN classifier with  $R$  labels, each of the  $K$  neighbors defines one of the experts. Moreover, the value of  $\mathcal{E}_{rc}$  will be  $w_c$  defined by (52.22c), while all other values on the  $c$ th column of  $\mathcal{E}$  will be 0.

(c) For the  $K$ –NN classifier, if we add the entries on each row  $r$  we find the number of votes that label  $r$  receives from the  $K$  neighbors. The row with the highest number of votes determines  $r^*$ . Likewise, for the weighted  $K$ –NN classifier, if we add the entries on each row  $r$  we find the overall weight that label  $r$  receives from the  $K$  neighbors. The row with the highest aggregate weight determines  $r^*$ . Both scenarios amount to determining  $r^*$  by determining the row of  $\mathcal{E}$  with the highest sum of its entries. This corresponds to the  $\infty$ –norm of  $\mathcal{E}$ .

□