

Exercise Sheet 7

Exercise 1: Bias and Variance of Mean Estimators (20 P)

Assume we have an estimator $\hat{\theta}$ for a parameter θ . The bias of the estimator $\hat{\theta}$ is the difference between the true value for the estimator, and its expected value

$$\text{Bias}(\hat{\theta}) = \mathbb{E}[\hat{\theta} - \theta].$$

If $\text{Bias}(\hat{\theta}) = 0$, then $\hat{\theta}$ is called unbiased. The variance of the estimator $\hat{\theta}$ is the expected square deviation from its expected value

$$\text{Var}(\hat{\theta}) = \mathbb{E}[(\hat{\theta} - \mathbb{E}[\hat{\theta}])^2].$$

The mean squared error of the estimator $\hat{\theta}$ is

$$\text{Error}(\hat{\theta}) = \mathbb{E}[(\hat{\theta} - \theta)^2] = \text{Bias}(\hat{\theta})^2 + \text{Var}(\hat{\theta}).$$

Let X_1, \dots, X_N be a sample of i.i.d random variables. Assume that X_i has mean μ and variance σ^2 . Calculate the bias, variance and mean squared error of the mean estimator:

$$\hat{\mu} = \alpha \cdot \frac{1}{N} \sum_{i=1}^N X_i$$

where α is a parameter between 0 and 1.

Exercise 2: Bias-Variance Decomposition for Classification (30 P)

The bias-variance decomposition usually applies to regression data. In this exercise, we would like to obtain similar decomposition for classification, in particular, when the prediction is given as a probability distribution over C classes. Let $P = [P_1, \dots, P_C]$ be the ground truth class distribution associated to a particular input pattern. Assume a random estimator of class probabilities $\hat{P} = [\hat{P}_1, \dots, \hat{P}_C]$ for the same input pattern. The error function is given by the expected KL-divergence between the ground truth and the estimated probability distribution:

$$\text{Error} = \mathbb{E}[D_{\text{KL}}(P||\hat{P})] = \mathbb{E}\left[\sum_{i=1}^C P_i \log(P_i/\hat{P}_i)\right].$$

First, we would like to determine the mean of the class distribution estimator \hat{P} . We define the mean as the distribution that minimizes its expected KL divergence from the the class distribution estimator, that is, the distribution R that optimizes

$$\min_R \mathbb{E}[D_{\text{KL}}(R||\hat{P})].$$

(a) Show that the solution to the optimization problem above is given by

$$R = [R_1, \dots, R_C] \quad \text{where} \quad R_i = \frac{\exp \mathbb{E}[\log \hat{P}_i]}{\sum_j \exp \mathbb{E}[\log \hat{P}_j]} \quad \forall 1 \leq i \leq C.$$

(Hint: To implement the positivity constraint on R , you can reparameterize its components as $R_i = \exp(Z_i)$, and minimize the objective w.r.t. Z .)

(b) Prove the bias-variance decomposition

$$\text{Error}(\hat{P}) = \text{Bias}(\hat{P}) + \text{Var}(\hat{P})$$

where the error, bias and variance are given by

$$\text{Error}(\hat{P}) = \mathbb{E}[D_{\text{KL}}(P||\hat{P})], \quad \text{Bias}(\hat{P}) = D_{\text{KL}}(P||R), \quad \text{Var}(\hat{P}) = \mathbb{E}[D_{\text{KL}}(R||\hat{P})].$$

(Hint: as a first step, it can be useful to show that $\mathbb{E}[\log R_i - \log \hat{P}_i]$ does not depend on the index i .)

Exercise 3: Programming (50 P)

Download the programming files on ISIS and follow the instructions.

Exercise 1: Bias and Variance of Mean Estimators (20 P)

Assume we have an estimator $\hat{\theta}$ for a parameter θ . The bias of the estimator $\hat{\theta}$ is the difference between the true value for the estimator, and its expected value

$$\text{Bias}(\hat{\theta}) = E[\hat{\theta} - \theta].$$

If $\text{Bias}(\hat{\theta}) = 0$, then $\hat{\theta}$ is called unbiased. The variance of the estimator $\hat{\theta}$ is the expected square deviation from its expected value

$$\text{Var}(\hat{\theta}) = E[(\hat{\theta} - E[\hat{\theta}])^2].$$

The mean squared error of the estimator $\hat{\theta}$ is

$$\text{Error}(\hat{\theta}) = E[(\hat{\theta} - \theta)^2] = \text{Bias}(\hat{\theta})^2 + \text{Var}(\hat{\theta}).$$

Let X_1, \dots, X_N be a sample of i.i.d random variables. Assume that X_i has mean μ and variance σ^2 . Calculate the bias, variance and mean squared error of the mean estimator:

$$\hat{\mu} = \alpha \cdot \frac{1}{N} \sum_{i=1}^N X_i$$

where α is a parameter between 0 and 1.

$$\begin{aligned} \text{Bias}(\hat{\mu}) &= E[\hat{\mu} - \mu] = E\left[\alpha \cdot \frac{1}{N} \sum_{i=1}^N X_i - \mu\right] \\ &= \frac{\alpha}{N} \sum_{i=1}^N E[X_i] - \mu \\ &= \frac{\alpha}{N} \sum_{i=1}^N \mu - \mu \\ &= \frac{\alpha}{N} \cdot N \mu - \mu \\ &= (\alpha - 1) \mu \end{aligned}$$

$$\begin{aligned} \text{Var}[\hat{\mu}] &= \text{Var}\left[\frac{\alpha}{N} \sum_{i=1}^N X_i\right] = \left(\frac{\alpha}{N}\right)^2 \cdot \sum_{i=1}^N \text{Var}[X_i] \\ &= \left(\frac{\alpha}{N}\right)^2 \cdot \sum_{i=1}^N \sigma^2 \\ &= \frac{\alpha^2}{N^2} \cdot N \cdot \sigma^2 \\ &= \frac{\alpha^2}{N} \cdot \sigma^2 \end{aligned}$$

$$\text{Error}(\hat{\mu}) = \text{Bias}(\hat{\mu})^2 + \text{Var}[\hat{\mu}] = (\alpha - 1)^2 \mu^2 + \frac{\alpha^2}{N} \cdot \sigma^2$$

Exercise 2: Bias-Variance Decomposition for Classification (30 P)

The bias-variance decomposition usually applies to regression data. In this exercise, we would like to obtain similar decomposition for classification, in particular, when the prediction is given as a probability distribution over C classes. Let $P = [P_1, \dots, P_C]$ be the ground truth class distribution associated to a particular input pattern. Assume a random estimator of class probabilities $\hat{P} = [\hat{P}_1, \dots, \hat{P}_C]$ for the same input pattern. The error function is given by the expected KL-divergence between the ground truth and the estimated probability distribution:

$$\text{Error} = \mathbb{E}[D_{\text{KL}}(P||\hat{P})] = \mathbb{E}\left[\sum_{i=1}^C P_i \log(P_i/\hat{P}_i)\right].$$

First, we would like to determine the mean of the class distribution estimator \hat{P} . We define the mean as the distribution that minimizes its expected KL divergence from the the class distribution estimator, that is, the distribution R that optimizes

$$\min_R \mathbb{E}[D_{\text{KL}}(R||\hat{P})].$$

(a) Show that the solution to the optimization problem above is given by

$$R = [R_1, \dots, R_C] \quad \text{where} \quad R_i = \frac{\exp \mathbb{E}[\log \hat{P}_i]}{\sum_j \exp \mathbb{E}[\log \hat{P}_j]} \quad \forall 1 \leq i \leq C.$$

(Hint: To implement the positivity constraint on R , you can reparameterize its components as $R_i = \exp(Z_i)$, and minimize the objective w.r.t. Z .)

$$\begin{aligned} \min_R \mathbb{E}[D_{\text{KL}}(R||\hat{P})] &= \min_R \mathbb{E}\left[\sum_{i=1}^C R_i \cdot \log\left(\frac{R_i}{\hat{P}_i}\right)\right] \\ &= \min_R \mathbb{E}\left[\sum_{i=1}^C R_i \log(R_i) - \sum_{i=1}^C R_i \log(\hat{P}_i)\right] \\ &= \min_R \sum_{i=1}^C R_i \log(R_i) - \sum_{i=1}^C R_i \mathbb{E}[\log(\hat{P}_i)] \end{aligned}$$

这里E是对estimator求mean

(To make sure $R_i > 0 \quad \forall 1 \leq i \leq C$, we can reparameterize it as $R_i = e^{z_i}$)

$$\begin{aligned} &= \min_R \sum_{i=1}^C e^{z_i} \log(e^{z_i}) - \sum_{i=1}^C e^{z_i} \mathbb{E}[\log(\hat{P}_i)] \\ &= \min_R \sum_{i=1}^C z_i \cdot e^{z_i} - \sum_{i=1}^C e^{z_i} \mathbb{E}[\log(\hat{P}_i)] \end{aligned}$$

subject to $\sum_{i=1}^C R_i = 1 = \sum_{i=1}^C e^{z_i}$

$$\mathcal{L}(z, \lambda) = \sum_{i=1}^C z_i \cdot e^{z_i} - \sum_{i=1}^C e^{z_i} \mathbb{E}[\log(\hat{P}_i)] + \lambda \left(\sum_{i=1}^C e^{z_i} - 1\right)$$

$$\begin{aligned} \frac{\partial}{\partial z_i} \mathcal{L}(z, \lambda) &= e^{z_i} + z_i e^{z_i} - e^{z_i} \mathbb{E}[\log(\hat{P}_i)] + \lambda \cdot e^{z_i} \\ &= e^{z_i} \cdot (1 + z_i - \mathbb{E}[\log(\hat{P}_i)] + \lambda) = 0 \end{aligned}$$

$$\Rightarrow z_i = \mathbb{E}[\log(\hat{P}_i)] - 1 - \lambda \quad \textcircled{a}$$

$$\frac{\partial}{\partial \lambda} \mathcal{L}(z, \lambda) = \sum_{i=1}^C e^{z_i} - 1 = \sum_{i=1}^C e^{\mathbb{E}[\log(\hat{P}_i)] - 1 - \lambda} - 1 = 0$$

$$\sum_{i=1}^C e^{\mathbb{E}[\log(\hat{P}_i)]} = e^{1+\lambda}$$

$$1 + \lambda = \log\left[\sum_{i=1}^C e^{\mathbb{E}[\log(\hat{P}_i)]}\right] \quad \textcircled{b}$$

from 0.0: $z_i = E[\log(\hat{p}_i)] - 1 - \lambda$

$$= E[\log(\hat{p}_i)] - \log\left[\sum_{i=1}^c e^{E[\log(\hat{p}_i)]}\right]$$

$$R_i = e^{z_i} = e^{E[\log(\hat{p}_i)] - \log\left[\sum_{i=1}^c e^{E[\log(\hat{p}_i)]}\right]}$$

$$= \frac{e^{E[\log(\hat{p}_i)]}}{\sum_{i=1}^c e^{E[\log(\hat{p}_i)]}}$$

(b) Prove the bias-variance decomposition

$$\text{Error}(\hat{P}) = \text{Bias}(\hat{P}) + \text{Var}(\hat{P})$$

where the error, bias and variance are given by

$$\text{Error}(\hat{P}) = \mathbb{E}[D_{\text{KL}}(P||\hat{P})], \quad \text{Bias}(\hat{P}) = D_{\text{KL}}(P||R), \quad \text{Var}(\hat{P}) = \mathbb{E}[D_{\text{KL}}(R||\hat{P})].$$

(Hint: as a first step, it can be useful to show that $\mathbb{E}[\log R_i - \log \hat{P}_i]$ does not depend on the index i .)

$$\begin{aligned} \mathbb{E}[\log R_i - \log \hat{P}_i] &= \mathbb{E}\left[\log\left(\frac{e^{\mathbb{E}[\log(\hat{P}_i)]}}{\sum_{i=1}^C e^{\mathbb{E}[\log(\hat{P}_i)]}}\right) - \log \hat{P}_i\right] \\ &= \mathbb{E}\left[\mathbb{E}[\log(\hat{P}_i)] - \log\left(\sum_{i=1}^C e^{\mathbb{E}[\log(\hat{P}_i)]}\right) - \log \hat{P}_i\right] \\ &= \underbrace{\mathbb{E}\left[\mathbb{E}[\log(\hat{P}_i)] - \log \hat{P}_i\right]}_{=0} - \mathbb{E}\left[\log\left(\sum_{i=1}^C e^{\mathbb{E}[\log(\hat{P}_i)]}\right)\right] \\ &= -\mathbb{E}\left[\log\left(\sum_{i=1}^C e^{\mathbb{E}[\log(\hat{P}_i)]}\right)\right] \end{aligned}$$

from 2.a) @: $= -\mathbb{E}[1+\lambda] = -1-\lambda$ which is independent of index i

$$\begin{aligned} \text{Error}(\hat{P}) &= \mathbb{E}[D_{\text{KL}}(P||\hat{P})] = \mathbb{E}\left[\sum_{i=1}^C P_i \log\left(\frac{P_i}{\hat{P}_i}\right)\right] \\ &= \mathbb{E}\left[\sum_{i=1}^C P_i \log P_i - \sum_{i=1}^C P_i \log \hat{P}_i\right] \\ &= \mathbb{E}\left[\sum_{i=1}^C P_i \log P_i - \sum_{i=1}^C P_i \log R_i + \sum_{i=1}^C P_i \log R_i - \sum_{i=1}^C P_i \log \hat{P}_i\right] \\ &= \text{Bias}(\hat{P}) + \mathbb{E}\left[\sum_{i=1}^C P_i \log R_i - \sum_{i=1}^C P_i \log \hat{P}_i\right] \\ &= \text{Bias}(\hat{P}) + \sum_{i=1}^C \left[P_i \cdot \mathbb{E}[\log R_i - \log \hat{P}_i]\right] \\ &= \text{Bias}(\hat{P}) + \sum_{i=1}^C \left[P_i \cdot (-1-\lambda)\right] \\ &= \text{Bias}(\hat{P}) + (-1-\lambda) \sum_{i=1}^C P_i \\ &= \text{Bias}(\hat{P}) + (-1-\lambda) \sum_{i=1}^C R_i \quad \text{Because of } \sum_{i=1}^C P_i = \sum_{i=1}^C R_i = 1 \\ &= \text{Bias}(\hat{P}) + \sum_{i=1}^C \left[R_i \cdot \mathbb{E}[\log R_i - \log \hat{P}_i]\right] \\ &= \text{Bias}(\hat{P}) + \mathbb{E}\left[\sum_{i=1}^C R_i \cdot \log \frac{R_i}{\hat{P}_i}\right] \\ &= \text{Bias}(\hat{P}) + \mathbb{V}[\hat{P}] \end{aligned}$$