

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

Assume we have an estimator  $\hat{\theta}$  for a parameter  $\theta$ . 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  $\mu$  and variance  $\sigma^2$ . 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  $\alpha$  is a parameter between 0 and 1.

$$\mathbb{E}[X_i] = \mu \quad \text{Var}[X_i] = \sigma^2$$

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

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

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

**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.  $Z$ .)

$$R_i = \exp(z_i) \Rightarrow \sum_{i=1}^C R_i = 1 = \sum_{i=1}^C \exp(z_i)$$

$$\mathbb{E}[D_{\text{KL}}(R||\hat{P})] = \mathbb{E}\left[\sum_i R_i \cdot \log \frac{R_i}{\hat{P}_i}\right] = \sum_i \mathbb{E}[R_i \cdot \log R_i] - \mathbb{E}[R_i \cdot \log \hat{P}_i]$$

$$= \sum_i R_i \cdot \log R_i - R_i \cdot \mathbb{E}[\log \hat{P}_i]$$

$$= \sum_i \exp(z_i) \cdot z_i - \exp(z_i) \cdot \mathbb{E}[\log \hat{P}_i], \quad \text{where} \quad \sum_{i=1}^C \exp(z_i) = 1$$

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

$$= \sum_i \exp(z_i) \cdot [z_i - E[\log(\hat{p}_i)] + \lambda] - \lambda$$

$$\nabla_z \mathcal{L}(z) = \sum_i \exp(z_i) \cdot [z_i - E[\log(\hat{p}_i)] + \lambda + 1] = 0$$

$$\Leftrightarrow z_i - E[\log(\hat{p}_i)] + \lambda + 1 = 0$$

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

$$\Rightarrow R_i = \exp(E[\log(\hat{p}_i)] - \lambda - 1)$$

$$\Rightarrow R_i = \frac{\exp(E[\log(\hat{p}_i)])}{\exp(1+\lambda)}$$

So we can let  $\exp(1+\lambda) = \sum_j \exp(E[\log \hat{p}_j])$

(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

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

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

$$\begin{aligned} E[\log R_i - \log \hat{P}_i] &= E\left[\log \frac{\exp(E[\log(\hat{p}_i)])}{\exp(1+\lambda)} - \log \hat{p}_i\right] \\ &= E[E[\log \hat{p}_i] - (1+\lambda) - \log \hat{p}_i] \\ &= -(1+\lambda) \end{aligned}$$

$$\begin{aligned} \text{Bias}(\hat{P}) + \text{Var}(\hat{P}) &= D_{\text{KL}}(P||R) + E[D_{\text{KL}}(R||\hat{P})] \\ &= \sum_i p_i \cdot \int \ln \frac{p_i}{R_i} + E\left[\sum_i R_i \cdot \int \ln \frac{R_i}{\hat{p}_i}\right] \\ &= \sum_i p_i \cdot \int \ln \frac{p_i}{R_i} + E\left(\sum_i R_i (\log R_i - \log \hat{p}_i)\right) \\ &= \sum_i p_i \cdot \int \ln \frac{p_i}{R_i} + E\left(\sum_i p_i (\log R_i - \log \hat{p}_i)\right) \end{aligned}$$

$$= E\left(\sum_i p_i \log p_i - p_i \log R_i + p_i \log R_i - p_i \log \hat{p}_i\right)$$

$$= E\left(\sum_i p_i \log p_i - p_i \log \hat{p}_i\right)$$

$$= E(DL(P \parallel \hat{P}))$$

$$= \text{error}(\hat{P})$$