

Statistical Machine Learning

Exercise sheet 2

Exercise 2.1 (Continuation of Ex 1.1) Let $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$, where $\mathbb{E}(\boldsymbol{\varepsilon}) = \mathbf{0}$, $\text{Cov}(\boldsymbol{\varepsilon}) = \sigma^2 \mathbf{I}$ and \mathbf{X} is a non-random full rank matrix of size $n \times p$. This setup contains the Gauss-Markov assumptions of a linear model.

- (a) Prove the Gauss-Markov theorem, i.e., $\hat{\boldsymbol{\beta}}$ is the best **linear unbiased** estimator (BLUE) of $\boldsymbol{\beta}$. “Best” in the sense that for all other linear unbiased estimators $\tilde{\boldsymbol{\beta}}$ of $\boldsymbol{\beta}$, $\text{Cov}(\tilde{\boldsymbol{\beta}}) - \text{Cov}(\hat{\boldsymbol{\beta}})$ is a positive semidefinite matrix.

Hints: Recall that an estimator $\tilde{\boldsymbol{\beta}}$ is linear if $\tilde{\boldsymbol{\beta}} = \mathbf{A}\mathbf{y}$, for some $\mathbf{A} \in \mathbb{R}^{p \times n}$. Notice that the matrix \mathbf{A} can be decomposed as $\mathbf{A} = \mathbf{B} + (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top$.

Solution: There are many variants of the proof of the Gauss Markov theorem given in different textbooks, many of which rely on elementary results from linear algebra and some “tricks”. We give here one such proof. First, note that the fact that we only consider estimators in the class of linear and unbiased estimators imposes the condition that $\mathbb{E}(\tilde{\boldsymbol{\beta}}) = \mathbb{E}(\mathbf{A}\mathbf{y}) = \mathbf{A}\mathbf{X}\boldsymbol{\beta} + \mathbb{E}(\mathbf{A}\boldsymbol{\varepsilon}) = \boldsymbol{\beta}$. This implies that $\mathbf{A}\mathbf{X} = \mathbf{I}$, the $p \times p$ identity matrix. Let $\mathbf{B} = \mathbf{A} - (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top$. We see that $\mathbf{B}\mathbf{X} = \mathbf{A}\mathbf{X} - (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{X} = \mathbf{I} - \mathbf{I} = \mathbf{0}$, so

$$\text{Cov}(\tilde{\boldsymbol{\beta}}) = \text{Cov}(\mathbf{A}\mathbf{y}) = \text{Cov}(\mathbf{A}\boldsymbol{\varepsilon}) = \sigma^2 \mathbf{A}\mathbf{A}^\top = \sigma^2 \mathbf{B}\mathbf{B}^\top + \sigma^2 (\mathbf{X}^\top \mathbf{X})^{-1},$$

where the last equality follows from $\mathbf{B}\mathbf{X} = \mathbf{0}$. Using the result from (b), we see that $\text{Cov}(\tilde{\boldsymbol{\beta}}) - \text{Cov}(\hat{\boldsymbol{\beta}}) = \sigma^2 \mathbf{B}\mathbf{B}^\top$, where $\mathbf{B}\mathbf{B}^\top$ is now a positive semidefinite matrix. The matrix difference is the zero matrix if and only if $\mathbf{B} = \mathbf{0}$, that is, when $\mathbf{A} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top$. So indeed, we have that $\hat{\boldsymbol{\beta}}$ gives the minimal variance in the class of all linear unbiased estimators.

- (b) Assume now that the errors $\boldsymbol{\varepsilon}$ are normally distributed. Prove that $\hat{\boldsymbol{\beta}}$ is the best estimator among **all unbiased** estimators. $\hat{\boldsymbol{\beta}}$ is then a uniformly minimum variance unbiased (UMVU) estimator.

Hint: Remember the Cramér–Rao bound.

Solution: We have

$$\frac{\partial^2 \text{RSS}}{\partial \boldsymbol{\beta} \partial \boldsymbol{\beta}^\top} = -2\mathbf{X}^\top \mathbf{X}.$$

The log-likelihood is $\ell(\boldsymbol{\beta}) = -\frac{1}{2\sigma^2} \text{RSS} + \text{constant}$, hence the Fisher information matrix is

$$I(\hat{\boldsymbol{\beta}}) = -\mathbb{E} \left(\frac{\partial^2 \ell}{\partial \boldsymbol{\beta} \partial \boldsymbol{\beta}^\top} \right) = \frac{1}{\sigma^2} \mathbf{X}^\top \mathbf{X}.$$

Exercise 2.2 (The regression function) Recall that we are interested in the predictive model $f^* : \mathbb{R}^p \rightarrow \mathbb{R}$ that minimizes the expected error for the ℓ^2 loss. i.e., we want to find the function f^* such that

$$\mathbb{E}[\ell\{Y, f^*(\mathbf{X})\}] = \mathbb{E}[\{Y - f^*(\mathbf{X})\}^2] = \min_{f: \mathbb{R}^p \rightarrow \mathbb{R}} \mathbb{E}[\{Y - f(\mathbf{X})\}^2].$$

- (a) Show that $f^*(\mathbf{x}) = \mathbb{E}(Y | \mathbf{X} = \mathbf{x})$.

Solution: There are many different ways to solve this exercise. We present here two solutions.

1st Solution: Let $m(\mathbf{x}) = \mathbb{E}(Y | \mathbf{X} = \mathbf{x})$. We see that

$$\begin{aligned} \mathbb{E}[L\{Y, f(\mathbf{X})\}] &= \mathbb{E}[\{Y - f(\mathbf{X})\}^2] \\ &= \mathbb{E}[\{Y - m(\mathbf{X}) + m(\mathbf{X}) - f(\mathbf{X})\}^2] \\ &= \mathbb{E}[\{Y - m(\mathbf{X})\}^2] + \mathbb{E}[\{m(\mathbf{X}) - f(\mathbf{X})\}^2] + 2\mathbb{E}[\{Y - m(\mathbf{X})\}\{m(\mathbf{X}) - f(\mathbf{X})\}] \\ &= \mathbb{E}[\{Y - m(\mathbf{X})\}^2] + \int_{\mathbb{R}^p} |m(\mathbf{x}) - f(\mathbf{x})|^2 P_{\mathbf{X}}(d\mathbf{x}), \end{aligned}$$

where $P_{\mathbf{X}}$ is the distribution of \mathbf{X} and the last equality holds because

$$\begin{aligned} \mathbb{E}[\{Y - m(\mathbf{X})\}\{m(\mathbf{X}) - f(\mathbf{X})\}] &= \mathbb{E}(\mathbb{E}[\{Y - m(\mathbf{X})\}\{m(\mathbf{X}) - f(\mathbf{X})\} | \mathbf{X}]) \\ &= \mathbb{E}(\{m(\mathbf{X}) - f(\mathbf{X})\}\{\mathbb{E}(Y | \mathbf{X}) - m(\mathbf{X})\}) \\ &= 0. \end{aligned}$$

Clearly, the L^2 loss is minimized when $f(\mathbf{x}) = m(\mathbf{x})$, proving that indeed, $f^*(\mathbf{x}) = \mathbb{E}(Y | \mathbf{X} = \mathbf{x})$.

2nd Solution: Write

$$\mathbb{E}[\{Y - f(\mathbf{X})\}^2] = \mathbb{E}(\mathbb{E}[\{Y - f(\mathbf{X})\}^2 | \mathbf{X}]) = \int_{\mathbb{R}^p} \mathbb{E}[\{Y - f(\mathbf{X})\}^2 | \mathbf{X} = \mathbf{x}] P_{\mathbf{X}}(d\mathbf{x}),$$

with $\mathbb{E}[\{Y - f(\mathbf{X})\}^2 | \mathbf{X} = \mathbf{x}] \geq 0$ for any $\mathbf{x} \in \mathbb{R}^p$, thus it suffices to minimize $\mathbb{E}[\{Y - f(\mathbf{x})\}^2 | \mathbf{X} = \mathbf{x}]$ over $f(\mathbf{x})$ (almost everywhere, with respect to the measure $\mathbb{P}(\mathbf{X})$). The problem is equivalent to finding the value of $c = c(\mathbf{x}) \in \mathbb{R}$ that minimizes $\mathbb{E}\{(Y - c)^2 | \mathbf{X} = \mathbf{x}\}$ almost everywhere. Write

$$\begin{aligned} \mathbb{E}\{(Y - c)^2 | \mathbf{X} = \mathbf{x}\} &= \mathbb{E}(Y^2 - 2cY + c^2 | \mathbf{X} = \mathbf{x}) \\ &= \mathbb{E}(Y^2 | \mathbf{X} = \mathbf{x}) - 2c\mathbb{E}(Y | \mathbf{X} = \mathbf{x}) + c^2. \end{aligned}$$

To minimize the previous expression, take the derivative with respect to c to find that the minimum is attained at $c = \mathbb{E}(Y | \mathbf{X} = \mathbf{x}) = m(\mathbf{x})$.

- (b) If we consider the ℓ^1 loss instead, i.e., $\ell(y, \hat{y}) = |y - \hat{y}|$, what is f^* ? (For simplicity suppose that $\mathbb{P}(Y | \mathbf{X})$ has a density.)

Solution: We want to find a function f^* that minimizes the expected loss $\mathbb{E}[|Y - f(\mathbf{X})|]$ over all functions $f : \mathbb{R}^p \rightarrow \mathbb{R}$. First, notice that

$$\mathbb{E}[|Y - f(\mathbf{X})|] = \mathbb{E}[\mathbb{E}[|Y - f(\mathbf{X})| | \mathbf{X}]]$$

Thus, using the same argument as in (a, 2nd solution), it suffices to find the constant c that minimizes $\mathbb{E}(|Y - c| \mid \mathbf{X} = \mathbf{x})$. Write

$$\begin{aligned}\mathbb{E}(|Y - c| \mid \mathbf{X}) &= \mathbb{E}\{(Y - c)I(Y - c \geq 0) + (c - Y)I(Y - c < 0) \mid \mathbf{X}\} \\ &= \mathbb{E}\{YI(Y - c \geq 0) \mid \mathbf{X}\} - c\mathbb{P}(Y \geq c \mid \mathbf{X}) \\ &\quad + c\mathbb{P}(Y < c \mid \mathbf{X}) - \mathbb{E}\{YI(Y - c < 0) \mid \mathbf{X}\} \\ &= \int_c^\infty yg(y \mid \mathbf{X})dy - c\mathbb{P}(Y \geq c \mid \mathbf{X}) + c\mathbb{P}(Y < c \mid \mathbf{X}) - \int_{-\infty}^c yg(y \mid \mathbf{X})dy,\end{aligned}$$

where g denotes the density of $Y \mid \mathbf{X}$. Thus

$$\begin{aligned}\frac{\partial}{\partial c}\mathbb{E}(|Y - c| \mid \mathbf{X}) &= -cg(c \mid \mathbf{X}) - \mathbb{P}(Y \geq c \mid \mathbf{X}) + cg(c \mid \mathbf{X}) \\ &\quad + \mathbb{P}(Y < c \mid \mathbf{X}) + cg(c \mid \mathbf{X}) - cg(c \mid \mathbf{X}) \\ &= -\mathbb{P}(Y \geq c \mid \mathbf{X}) + \mathbb{P}(Y < c \mid \mathbf{X}),\end{aligned}$$

and the zero of that derivative occurs when $\mathbb{P}(Y < c \mid \mathbf{X}) = \mathbb{P}(Y \geq c \mid \mathbf{X}) = 1/2$, i.e., for c equal to the conditional median $c = f^*(\mathbf{x}) = \text{median}(Y \mid \mathbf{X} = \mathbf{x})$. Note: the same result holds without assuming a density for $Y \mid \mathbf{X}$, it is just more difficult to prove.

Exercise 2.3 (Bias-variance tradeoff) In this exercise, we consider the expected ℓ^2 error of a random predictive model \hat{f}_n (depends on a training set \mathcal{D}_n), defined as

$$\mathbb{E} \left[\int_{\mathbb{R}^p} \{\hat{f}_n(\mathbf{x}) - f^*(\mathbf{x})\}^2 P_{\mathbf{X}}(d\mathbf{x}) \right]. \quad (1)$$

- (a) For any random predictive model \hat{f}_n and any fixed point $\mathbf{x}_0 \in \mathbb{R}^p$, prove that

$$\mathbb{E} \left[\{\hat{f}_n(\mathbf{x}_0) - f^*(\mathbf{x}_0)\}^2 \right] = \left[\text{bias}\{\hat{f}_n(\mathbf{x}_0)\} \right]^2 + \text{var}\{\hat{f}_n(\mathbf{x}_0)\}.$$

Solution: We have that

$$\begin{aligned}\mathbb{E} \left[\{\hat{f}_n(\mathbf{x}_0) - f^*(\mathbf{x}_0)\}^2 \right] &= \mathbb{E} \left[\{\hat{f}_n(\mathbf{x}_0) - \mathbb{E}\{\hat{f}_n(\mathbf{x}_0)\} + \mathbb{E}\{\hat{f}_n(\mathbf{x}_0)\} - f^*(\mathbf{x}_0)\}^2 \right] \\ &= \left[\text{bias}\{\hat{f}_n(\mathbf{x}_0)\} \right]^2 + \text{var}\{\hat{f}_n(\mathbf{x}_0)\}\end{aligned}$$

because

$$\mathbb{E} \left[\{\hat{f}_n(\mathbf{x}_0) - \mathbb{E}\{\hat{f}_n(\mathbf{x}_0)\}\} \{\mathbb{E}\{\hat{f}_n(\mathbf{x}_0)\} - f^*(\mathbf{x}_0)\} \right] = 0.$$

Note that in this expectation, $\mathbb{E}\{\hat{f}_n(\mathbf{x}_0)\} - f^*(\mathbf{x}_0)$ is non-random. Important: $f^*(\mathbf{x}_0)$ is just a constant here, it could be the conditional mean, as in Exercise 2.1(a), or the conditional median, as in Exercise 2.1(b), or something else.

- (b) Find a similar bias-variance decomposition for the expected ℓ^2 error (1).

Solution: By Fubini or Tonelli's theorem (where our function within the integrand is positive and measurable), we have that the expected ℓ^2 error is given by

$$\begin{aligned}\mathbb{E} \left[\int_{\mathbb{R}^p} \{\hat{f}_n(\mathbf{x}) - f^*(\mathbf{x})\}^2 P_{\mathbf{X}}(d\mathbf{x}) \right] &= \int_{\mathbb{R}^p} \mathbb{E} \left[\{\hat{f}_n(\mathbf{x}) - f^*(\mathbf{x})\}^2 \right] P_{\mathbf{X}}(d\mathbf{x}) \\ &= \int_{\mathbb{R}^p} \left[\text{bias}\{\hat{f}_n(\mathbf{x})\} \right]^2 P_{\mathbf{X}}(d\mathbf{x}) + \int_{\mathbb{R}^p} \text{var}\{\hat{f}_n(\mathbf{x})\} P_{\mathbf{X}}(d\mathbf{x}).\end{aligned}$$

We see that for the expected ℓ^2 error the squared bias and the variance integrate over the distribution of \mathbf{X} .

Exercise 2.4 (Ridge regression)

- (a) Consider the linear regression model

$$y_i = \beta_0 + \sum_{j=1}^p \beta_j x_{ij} + \varepsilon_i, \quad i = 1, \dots, n.$$

Define $\boldsymbol{\beta} = (\beta_1, \dots, \beta_p)^\top$ and the residuals as

$$r_i(\beta_0, \boldsymbol{\beta}) = y_i - \beta_0 - \sum_{j=1}^p \beta_j x_{ij}, \quad i = 1, \dots, n.$$

Show that the OLS estimator $\hat{\beta}_0 = \bar{y} - \sum_{j=1}^p \beta_j x_{.j}$ for any $\boldsymbol{\beta}$, where $x_{.j} = \frac{1}{n} \sum_{i=1}^n x_{ij}$. Hence deduce that

$$r_i(\hat{\beta}_0, \boldsymbol{\beta}) = y_i - \bar{y} - \sum_{j=1}^p \beta_j (x_{ij} - x_{.j}), \quad i = 1, \dots, n.$$

Discuss the implications of this result.

Solution: The OLS estimator $\hat{\beta}_0$ is found by minimizing the residual sum of squares with respect to β_0 , keeping $\boldsymbol{\beta}$ fixed:

$$\text{RSS} = (\mathbf{y} - \beta_0 \mathbf{1} - \mathbf{X}\boldsymbol{\beta})^\top (\mathbf{y} - \beta_0 \mathbf{1} - \mathbf{X}\boldsymbol{\beta}),$$

so $\frac{\partial}{\partial \beta_0} \text{RSS} = 2 \times \mathbf{1}^\top (\mathbf{y} - \beta_0 \mathbf{1} - \mathbf{X}\boldsymbol{\beta})$, and finding the root of the previous derivative gives $\hat{\beta}_0 = \bar{y} - \sum_{j=1}^p \beta_j x_{.j}$. The form of $r_i(\hat{\beta}_0, \boldsymbol{\beta})$ is then trivial. Note that this result also follows from the Frisch–Waugh–Lovell (FWL) theorem. We now discuss the significance of this result. Note that in both residuals formulations, the vector $\boldsymbol{\beta}$ is the same. Hence, by centering the response and predictor variables it is always possible to get rid of the intercept β_0 in the first equation. The least squares estimation of $\boldsymbol{\beta}$ is the same in both the model with an intercept and in the model without it but with centered response and covariates. Thus $\boldsymbol{\beta}$ can first be estimated from the centered response and covariates (using OLS or another method) and β_0 will then be estimated by $\hat{\beta}_0 = \bar{y} - \sum_{j=1}^p \hat{\beta}_j x_{.j}$. This motivates the definition of the ridge estimator in the next equation (equation (2)), compared to the definition with an unconstrained intercept: as β_0 is unconstrained we know $\hat{\beta}_0 = \bar{y} - \sum_{j=1}^p \hat{\beta}_j x_{.j}$ and $\hat{\boldsymbol{\beta}}$ is found solving equation (2), thus the two formulations of ridge will give the same estimates of $\boldsymbol{\beta}$ and β_0 .

- (b) Define the ridge regression estimator as a minimizer of the penalized residual sum of squares,

$$\hat{\boldsymbol{\beta}}(\lambda) = \underset{\boldsymbol{\beta}}{\operatorname{argmin}} \frac{1}{n} \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|^2 + \lambda \boldsymbol{\beta}^\top \boldsymbol{\beta}, \quad (2)$$

where $\lambda \geq 0$ is a parameter that controls the amount of shrinkage. Show that the ridge regression solution always exists, even if \mathbf{X} does not have full rank, and is given by

$$\hat{\beta}(\lambda) = (\mathbf{X}^\top \mathbf{X} + n\lambda \mathbf{I})^{-1} \mathbf{X}^\top \mathbf{y}.$$

Note that the ridge estimator is still linearly depending on the response \mathbf{y} , as for ordinary least squares.

Solution: Taking gradient of $\frac{1}{n} \|\mathbf{y} - \mathbf{X}\beta\|^2 + \lambda \beta^\top \beta$ with respect to β , setting it to 0, and rearranging the terms, gives the normal equations

$$(\mathbf{X}^\top \mathbf{X} + n\lambda \mathbf{I}) \hat{\beta}(\lambda) = \mathbf{X}^\top \mathbf{y}.$$

Now note that the matrix $\mathbf{X}^\top \mathbf{X} + n\lambda \mathbf{I}$ is always invertible: $\mathbf{X}^\top \mathbf{X}$ is positive semidefinite (since $\mathbf{x}^\top \mathbf{X}^\top \mathbf{X} \mathbf{x} = \|\mathbf{X}\mathbf{x}\|^2 \geq 0$) hence its eigenvalues $\alpha_i \geq 0$, and the eigenvalues of $\mathbf{X}^\top \mathbf{X} + n\lambda \mathbf{I}$ are $\alpha_i + n\lambda > 0$ for any $\lambda > 0$. Hence the ridge estimator always exists and is uniquely defined as

$$\hat{\beta}(\lambda) = (\mathbf{X}^\top \mathbf{X} + n\lambda \mathbf{I})^{-1} \mathbf{X}^\top \mathbf{y}.$$

(c) Show that the ridge regression estimator defined in (2) equals

$$\hat{\beta}(t) = \underset{\|\beta\|^2 \leq t}{\operatorname{argmin}} \|\mathbf{y} - \mathbf{X}\beta\|^2 \quad (3)$$

for a given $t = t(\lambda)$. *Hint: Use the Karush–Kuhn–Tucker (KKT) method.*

Solution: The constrained optimization problem in (3) can be solved by means of the Karush–Kuhn–Tucker (KKT) multiplier method, which minimizes a function subject to inequality constraints. The KKT multiplier method states that, under some regularity conditions (here satisfied), there exists a unique λ , called the multiplier, such that the solution $\hat{\beta}$ of the constrained minimization problem (3) satisfies the so-called KKT conditions. Define the Lagrangian of the problem as

$$\ell(\beta, \lambda) = \|\mathbf{y} - \mathbf{X}\beta\|^2 + \lambda (\|\beta\|^2 - t).$$

The KKT conditions are:

1. $\frac{\partial \ell}{\partial \beta}(\hat{\beta}, \lambda) = \mathbf{0}.$
2. $\lambda \geq 0.$
3. $\lambda(\|\hat{\beta}\|^2 - t) = 0.$
4. $\|\hat{\beta}\|^2 - t \leq 0.$

Now, suppose that the solution that minimizes the minimization problem in (2) is given by the ridge estimate $\hat{\beta}(\lambda/n)$. For $t = \|\hat{\beta}(\lambda/n)\|^2$ it is clear that $\hat{\beta}(\lambda/n)$ satisfies all KKT conditions. Hence, both the constrained optimization problem in (3) and the minimization problem in (2) have the same solution when $t = \|\hat{\beta}(\lambda/n)\|^2$.

Exercise 2.5 The Gauss-Markov Theorem makes the assumption that the training data is generated as $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$, \mathbf{X} is a non-random full rank matrix of size $n \times p$, where $\mathbb{E}[\boldsymbol{\varepsilon}] = \mathbf{0}$, $\text{Cov}(\boldsymbol{\varepsilon}) = \sigma^2 \mathbf{I}$.

- (a) Explain why the Gauss-Markov Theorem still holds for any random design matrix \mathbf{X} (in particular without assuming that the rows of \mathbf{X} are i.i.d.) provided we change the assumptions and assume that $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$ with $\mathbb{E}[\boldsymbol{\varepsilon}|\mathbf{X}] = \mathbf{0}$, $\text{Cov}(\boldsymbol{\varepsilon}|\mathbf{X}) = \sigma^2 \mathbf{I}$.

Solution: Recall that when the design matrix \mathbf{X} is non-random, the Gauss-Markov theorem indicates that for any unbiased linear estimator $\tilde{\boldsymbol{\beta}}$, there exists $\mathbf{K}_{\mathbf{X}}$ positive semi-definite such that

$$\text{Cov}(\tilde{\boldsymbol{\beta}}|\mathbf{X}) = \text{Cov}(\hat{\boldsymbol{\beta}}|\mathbf{X}) + \mathbf{K}_{\mathbf{X}}, \quad (\star)$$

where the covariances are conditional covariances and $\mathbf{K}_{\mathbf{X}}$ does not only depend on the choice of the estimator $\tilde{\boldsymbol{\beta}}$ but also on \mathbf{X} . To obtain a statement which is again expressed in terms of the marginal covariance of the estimators, we can use the variance decomposition formula

$$\text{Cov}(\tilde{\boldsymbol{\beta}}) = \mathbb{E}[\text{Cov}(\tilde{\boldsymbol{\beta}}|\mathbf{X})] + \text{Cov}(\mathbb{E}[\tilde{\boldsymbol{\beta}}|\mathbf{X}]),$$

so taking expectation in (\star) , we get

$$\begin{aligned} \text{Cov}(\tilde{\boldsymbol{\beta}}) &= \mathbb{E}[\text{Cov}(\tilde{\boldsymbol{\beta}}|\mathbf{X})] + \mathbb{E}[\mathbf{K}_{\mathbf{X}}] + \text{Cov}(\mathbb{E}[\tilde{\boldsymbol{\beta}}|\mathbf{X}]) \\ &= \text{Cov}(\hat{\boldsymbol{\beta}}) + \mathbb{E}[\mathbf{K}_{\mathbf{X}}] + \text{Cov}(\mathbb{E}[\tilde{\boldsymbol{\beta}}|\mathbf{X}]), \end{aligned}$$

where we use the fact that $\mathbb{E}[\hat{\boldsymbol{\beta}}|\mathbf{X}] = \boldsymbol{\beta}$ thus $\text{Cov}(\mathbb{E}[\hat{\boldsymbol{\beta}}|\mathbf{X}]) = \mathbf{0}$. Finally, both $\mathbb{E}[\mathbf{K}_{\mathbf{X}}]$ and $\text{Cov}(\mathbb{E}[\tilde{\boldsymbol{\beta}}|\mathbf{X}])$ are also positive semi-definite, which yields the Gauss-Markov theorem.

- (b) Let $\tilde{\boldsymbol{\beta}}$ be any linear unbiased estimator and let $\hat{\boldsymbol{\beta}}$ be the linear regression estimator (aka ordinary least squares estimator). Show that as a consequence of the Gauss-Markov theorem:

$$\forall \mathbf{x} \in \mathbb{R}^p, \quad \text{Var}(\mathbf{x}^\top \hat{\boldsymbol{\beta}}) \leq \text{Var}(\mathbf{x}^\top \tilde{\boldsymbol{\beta}}).$$

Solution: By the previous question, we have that

$$\text{Cov}(\tilde{\boldsymbol{\beta}}) = \text{Cov}(\hat{\boldsymbol{\beta}}) + \mathbf{K},$$

with $\mathbf{K} = \mathbb{E}[\mathbf{K}_{\mathbf{X}}]$. But this entails that, for any fixed $\mathbf{x} \in \mathbb{R}^p$, we have

$$\mathbf{x}^\top \text{Cov}(\tilde{\boldsymbol{\beta}}) \mathbf{x} = \mathbf{x}^\top \text{Cov}(\hat{\boldsymbol{\beta}}) \mathbf{x} + \mathbf{x}^\top \mathbf{K} \mathbf{x};$$

now since \mathbf{K} is positive semi-definite, we must have $\mathbf{x}^\top \mathbf{K} \mathbf{x} \geq 0$ and it is easy to check that $\text{Var}(\mathbf{x}^\top \tilde{\boldsymbol{\beta}}) = \mathbf{x}^\top \text{Cov}(\tilde{\boldsymbol{\beta}}) \mathbf{x}$ using the definitions of variance and covariance. Hence the result.

- (c) Consider now i.i.d. data (X_i, Y_i) with $Y_i = X_i^\top \boldsymbol{\beta} + \varepsilon_i$, $\mathbb{E}[\varepsilon_i|X_i] = 0$ and $\text{Var}(\varepsilon_i|X_i) = \sigma^2$. For data following this distribution, express the target function for the quadratic risk as a function of $\boldsymbol{\beta}$.

Solution: We know that for the quadratic risk $f^*(x) = \mathbb{E}[Y|X = x]$ and

$$f^*(X_i) = \mathbb{E}[Y_i|X_i] = \mathbb{E}[X_i^\top \beta + \varepsilon_i|X_i] = X_i^\top \beta + \mathbb{E}[\varepsilon_i|X_i] = X_i^\top \beta.$$

So $f^*(\mathbf{x}) = \beta^\top \mathbf{x}$.

- (d) Let $\hat{f} : \mathbf{x} \mapsto \mathbf{x}^\top \hat{\beta}$ and $\tilde{f} : \mathbf{x} \mapsto \mathbf{x}^\top \tilde{\beta}$ for $\tilde{\beta}$ some unbiased linear estimator based on \mathbf{X} and \mathbf{y} . Show that for any such \tilde{f} , if \mathcal{R} denotes the quadratic risk (i.e. the risk associated with the square loss), then we necessarily have $\mathbb{E}[\mathcal{R}(\hat{f})] \leq \mathbb{E}[\mathcal{R}(\tilde{f})]$. Show that the same inequality actually holds conditionally on the value of any new $X = \mathbf{x}$.

Solution: Note first that since $\tilde{\beta}$ is unbiased then, for any fixed \mathbf{x} .

$$\mathbb{E}[\tilde{f}(\mathbf{x})] = \mathbb{E}[\tilde{\beta}^\top \mathbf{x}] = \mathbb{E}[\tilde{\beta}]^\top \mathbf{x} = \beta^\top \mathbf{x} = f^*(\mathbf{x}),$$

and, since $\hat{\beta}$ is itself unbiased, $\mathbb{E}[\hat{f}(\mathbf{x})] = f^*(\mathbf{x})$. Now, this entails that

$$\text{Var}(\tilde{f}(\mathbf{x})) = \mathbb{E}[(\tilde{f}(\mathbf{x}) - f^*(\mathbf{x}))^2].$$

To connect this with the risk, we denote by $D_n = \{(X_i, Y_i)\}_{i=1..n}$ the training set and we use the result established on the quadratic risk in class:

$$\mathbb{E}[(\tilde{f}(\mathbf{x}) - Y)^2 | X = \mathbf{x}, D_n] = (\tilde{f}(\mathbf{x}) - f^*(\mathbf{x}))^2 + \mathbb{E}[(f^*(\mathbf{x}) - Y)^2 | X = \mathbf{x}, D_n]$$

that we can also write

$$\mathcal{R}(\tilde{f}(\mathbf{x})|\mathbf{x}) = (\tilde{f}(\mathbf{x}) - f^*(\mathbf{x}))^2 + \mathcal{R}(f^*(\mathbf{x})|\mathbf{x}),$$

by identifying the conditional risks of \tilde{f} and f^* . Now taking expectations w.r.t. to D_n on both sides we get

$$\mathbb{E}[\mathcal{R}(\tilde{f}(\mathbf{x})|\mathbf{x})] = \text{Var}(\tilde{f}(\mathbf{x})) + \mathcal{R}(f^*(\mathbf{x})|\mathbf{x}),$$

where $\mathbb{E}[\mathcal{R}(\tilde{f}(\mathbf{x})|\mathbf{x})]$ is the expected conditional risk at \mathbf{x} , with an expectation taken over the choice of the training data D_n only. Given that we have the same formula hold for \hat{f} and given that $\text{Var}(\mathbf{x}^\top \hat{\beta}) \leq \text{Var}(\mathbf{x}^\top \tilde{\beta})$ we have established that

$$\mathbb{E}[\mathcal{R}(\hat{f}(\mathbf{x})|\mathbf{x})] \leq \mathbb{E}[\mathcal{R}(\tilde{f}(\mathbf{x})|\mathbf{x})].$$

By taking expectations over the choice of $X = \mathbf{x}$, we get

$$\mathbb{E}[\mathcal{R}(\hat{f})] \leq \mathbb{E}[\mathcal{R}(\tilde{f})].$$