# Midterm 1: Math 6266

*Peter Williams*

## Section 1.1

*Exercise 1. Consider the linear regression model with mean zero, uncorrelated, heteroscedastic noise:*

$$Y_i = X_i^\intercal\theta + \varepsilon_i, \; for \; i = 1,..,n, \; E\varepsilon_i = 0, \; cov(\varepsilon_i, \varepsilon_j) = \begin{cases} \sigma_i^2, & \text{if } i = j \\ 0, & i \neq j \end{cases} \tag{1}$$

*Find expressions for the LSE and response estimator in this model*

To set up the problem, take $W^{-1} = diag\{\sigma_1^2,...,\sigma_n^2\}$, $W = diag\{\frac{1}{\sigma_1^2},...,\frac{1}{\sigma_n^2}\}$, $W^{1/2} = diag\{\sqrt{\frac{1}{\sigma_1^2}},...,\sqrt{\frac{1}{\sigma_n^2}}\}$, with $W^\intercal = W$, and $W^{1/2}W^{1/2} = W$, since they are diagonal matrices. Also we will use $w_i = \frac{1}{\sigma_i^2} = W_{ii}$.

Under heteroscedastic noise assumptions, we define the least squares estimator, denoted $\hat{\theta}$, as:

$$\hat{\theta} = \underset{\theta}{argmin} \sum_{i=1}^{n} w_i(Y_i - X_i^\intercal\theta)^2 = \underset{\theta}{argmin} \sum_{i=1}^{n} (\sqrt{w_i}Y_i - \sqrt{w_i}X_i^\intercal\theta)^2 = \underset{\theta}{argmin}||W^{1/2}Y - W^{1/2}X^\intercal\theta||^2$$

$$G(\theta) = ||W^{1/2}Y - W^{1/2}X^\intercal\theta||^2 = (W^{1/2}Y - W^{1/2}X^\intercal\theta)^\intercal(W^{1/2}Y - W^{1/2}X^\intercal\theta) = Y^\intercal WY - 2\theta^\intercal XWY + \theta^\intercal XWX^\intercal\theta$$

with gradient,

$$\nabla G(\theta) = -2XWY + 2XWX^\intercal\theta$$

Setting this expression equal to zero leads to estimator $\hat{\theta} = (XWX^\intercal)^{-1}XWY$, which leads to response estimator $\hat{Y} = X^\intercal\hat{\theta} = X^\intercal(XWX^\intercal)^{-1}XWY$.

*Exercise 2. Assume that $\varepsilon_i \sim N(0, \sigma_i^2)$ in the previous problem. What is known about the distribution of $\hat{\theta}$ and $\hat{Y}$?*

For $\hat{\theta}$, we have,

$$E[\hat{\theta}] = E[(XWX^\intercal)^{-1}XWY] = E[(XWX^\intercal)^{-1}XW(X^\intercal\theta^* + \varepsilon)] = E[\theta^*] + E[(XWX^\intercal)^{-1}XW\varepsilon] = \theta^*$$

indicating that $\hat{\theta}$ is unbiased. Further $\hat{\theta}$ is normally distributed, since is a linear transformation of $\varepsilon \sim N(0, W^{-1})$. Further we have,

$$Var(\hat{\theta}) = Var((XWX^\intercal)^{-1}XWY) = Var((XWX^\intercal)^{-1}XW(X^\intercal\theta^* + \varepsilon)) = Var((XWX^\intercal)^{-1}XW\varepsilon)) = ...$$

$$= (XWX^\intercal)^{-1}XWVar(\varepsilon)W^\intercal X^\intercal(XWX^\intercal)^{-1} = (XWX^\intercal)^{-1}XWX^\intercal(XWX^\intercal)^{-1} = (XWX^\intercal)^{-1} = Var(\hat{\theta})$$

For $\hat{Y}$ we have,

$$E[\hat{Y}] = E[X^\intercal(XWX^\intercal)^{-1}XWY] = E[X^\intercal(XWX^\intercal)^{-1}XW(X^\intercal\theta^* + \varepsilon)] = E[X^\intercal\theta^* + X^\intercal(XWX^\intercal)^{-1}XW\varepsilon] = E[X^\intercal\theta^*] = Y$$

and,

$$Var[\hat{Y}] = Var[X^\intercal(XWX^\intercal)^{-1}XWY] = Var[X^\intercal(XWX^\intercal)^{-1}XW(X^\intercal\theta^* + \varepsilon)] = Var[X^\intercal\theta^* + X^\intercal(XWX^\intercal)^{-1}XW\varepsilon] = ...$$

$$... = Var[X^\intercal(XWX^\intercal)^{-1}XW\varepsilon] = X^\intercal(XWX^\intercal)^{-1}XW \, Var(\varepsilon) \, W^\intercal X^\intercal(XWX^\intercal)^{-1}X = ...$$

$$= X^\intercal(XWX^\intercal)^{-1}XWX^\intercal(XWX^\intercal)^{-1}X = X^\intercal(XWX^\intercal)^{-1}X$$

*Now suppose additionally that $\sigma_i^2 \equiv \sigma^2 > 0$. What can be said about distribution of the estimator $\hat{\sigma^2}$?*

With $\sigma_i^2 \equiv \sigma^2 > 0$, we have $\hat{\sigma}^2 = \frac{||Y - X^\mathsf{T}\hat{\theta}||^2}{n-p} = \frac{||\hat{\varepsilon}||^2}{n-p}$. Further denote, $||\hat{\varepsilon}|| = ||Y - \hat{Y}|| = ||Y - \Pi Y|| = ||(I_n - \Pi)Y||$, also noting that $(I_n - \Pi)X^\mathsf{T} = X^\mathsf{T} - \Pi X^\mathsf{T} = X^\mathsf{T} - X^\mathsf{T}(XX^\mathsf{T})^{-1}XX^\mathsf{T} = X^\mathsf{T} - X^\mathsf{T} = 0$.

Then we have,

$$(n-p)E[\hat{\sigma}^2] = E||Y - X^\mathsf{T}\hat{\theta}||^2 = E||\hat{\varepsilon}||^2 = E[tr(\hat{\varepsilon}\hat{\varepsilon}^\mathsf{T})] = E[tr((I_n - \Pi)YY^\mathsf{T}(I_n - \Pi))] = \dots$$

,

$$\dots = E[tr((I_n - \Pi)(X^\mathsf{T}\theta^* + \varepsilon)(X^\mathsf{T}\theta^* + \varepsilon)^\mathsf{T}(I_n - \Pi))] = E[tr((I_n - \Pi)\varepsilon\varepsilon^\mathsf{T}(I_n - \Pi))] = tr((I_n - \Pi)E[\varepsilon\varepsilon^\mathsf{T}]) = \dots$$

Using the cylic property of the trace operator, the property that $(I_n - \Pi)(I_n - \Pi) = (I_n - \Pi)$, and the expectation $E[\varepsilon\varepsilon^\mathsf{T}] = \sigma^2 I_n$, leading to

$$\dots = \sigma^2 tr(I_n - \Pi) = \sigma^2(n-p) = (n-p)E[\hat{\sigma}^2]$$

Looking further at the distribution of $||Y - X^\mathsf{T}\hat{\theta}||^2 = \hat{\varepsilon}^\mathsf{T}\hat{\varepsilon}$, we have $\hat{\varepsilon}^\mathsf{T}\hat{\varepsilon} = ((I_n - \Pi)Y)^\mathsf{T}((I_n - \Pi)Y) = Y^\mathsf{T}(I_n - \Pi)Y = (X^\mathsf{T}\theta^* + \varepsilon)^\mathsf{T}(I_n - \Pi)(X^\mathsf{T}\theta^* + \varepsilon) = \varepsilon^\mathsf{T}(I_n - \Pi)\varepsilon$.

Since we know that $\varepsilon \sim N(0, \sigma^2 I_n)$, and further $\frac{\varepsilon^\mathsf{T}\varepsilon}{\sigma^2} \sim \chi^2(n)$, $(\frac{\varepsilon}{\sigma})^\mathsf{T}(I_n - \Pi)(\frac{\varepsilon}{\sigma}) \sim \chi^2(n-p)$, since we know from earlier that $(I_n - \Pi)$, is idempotent, with rank equal to $tr(I_n - \Pi) = tr(I_n) - tr(\Pi) = n - p$.


**Section 1.3**

*Exercise 4. Let $A \in R^{n \times n}$ be a matrix (corresponding to a linear map in $R^n$). Show that $A$ preserves length for all $x \in R^n$ iff it preserves the inner product. I.e. one needs to show the following:*

$||Ax|| = ||x|| \; \forall \; x \in R^n \iff (Ax)^\mathsf{T}(Ay) \; \forall \; x, y \in R^n$.

Take,

$$||x|| = \sqrt{x \cdot x} = \sqrt{x^\mathsf{T}x} \implies ||Ax|| = \sqrt{Ax \cdot Ax} = \sqrt{x^\mathsf{T}A^\mathsf{T}Ax} \implies$$

,

$$A^\mathsf{T}A = I_n = A^{-1}, \; A^\mathsf{T} = A^{-1}, ||Ax|| = ||x||$$

this implies $A$ is an orthogonal matrix, and further,

$$(Ax)^\mathsf{T}(Ay) = ||AxAy||^2 = x^\mathsf{T}A^\mathsf{T}Ay = x^\mathsf{T}y = ||xy||^2$$


*Exercise 5. (a) Let $x_0 \in R^n$ be some fixed vector, find a projection map on the subspace $span(x_0)$. Compare your result with matrix $\Pi$ (from section 1.3) for the case of $p = 1$.*

Let $x = span(x_0) = span(x_1, x_2, .., x_n)$, denote the subspace of interest, and $x_1, x_2, \dots$ are basis vectors and $y = (y_1, y_2, \dots, y_n)^\mathsf{T}$. The projection map is,

$$Proj_x(y) = \frac{<y \cdot x>}{<y \cdot y>}x = \sum_{i=1}^{n} \frac{<y_i \cdot x_i>}{<y_i \cdot y_i>}x_i$$

For the case $p = 1$, and $\Pi = X^\mathsf{T}(XX^\mathsf{T})^{-1}X, X^\mathsf{T} \in R^n$, we have,

$$\Pi y = \hat{y} = X^\mathsf{T}(XX^\mathsf{T})^{-1}Xy = X^\mathsf{T}\frac{Xy}{XX^\mathsf{T}} = \frac{\sum_i^n x_i y_i}{\sum_i^n x_i^2}(x_1, x_2, \dots, x_n)^\mathsf{T} = \frac{<X \cdot y>}{<y \cdot y>}X^\mathsf{T} = Proj_X(y)$$


*(b) Prove part 3) of Lemma 1.1 for an arbitrary orthogonal projection in $R^n$. Show $\forall h \in R^n$, $||h||^2 = ||\Pi h||^2 + ||h - \Pi h||^2$.*

Using the fact that $(I_n - \Pi)^\mathsf{T}(I_n - \Pi) = I_n - 2\Pi + \Pi = I_n - \Pi$, we have,

$||h||^2 = ||\Pi h||^2 + ||h - \Pi h||^2 = h^\mathsf{T}\Pi^\mathsf{T}\Pi h + h^\mathsf{T}(I_n - \Pi)^\mathsf{T}(I_n - \Pi)h = h^\mathsf{T}\Pi h + h^\mathsf{T}(I_n - \Pi)h = h^\mathsf{T}I_n h + h^\mathsf{T}\Pi h - h^\mathsf{T}\Pi h = ||h||^2$