

PROJECT I

part A)

The assumption we have made is that there exists a continuous function $f(x)$ and a normal distributed error $\epsilon \sim N(0, \sigma^2)$ which describes our data

$$y = f(x) + \epsilon$$

Function $f(x)$ is approximated with our model \tilde{y} where we minimized $(y - \tilde{y})^2$, with

$$\tilde{y} = X\beta$$

Show that the expectation value of y for a given element i is

$$\mathbb{E}(y_i) = \sum_j x_{ij} \beta_j = X_{i,*} \beta$$

And that its variance is

$$\text{Var}(y_i) = \sigma^2$$

Let us start the proof for the first element by introducing the expectation value for y .

$$\mathbb{E}[y] = \frac{1}{n} \sum_{i=0}^{p-1} y_i \Rightarrow \frac{1}{n} \sum_{i=0}^{p-1} (f(x_i) + \epsilon_i)$$

$$f(x) \approx X\beta = \tilde{y}$$

Our function f has been approximated through minimization of cost function $(y - \hat{y})^2$, also known as MSE.

$$C(\beta) = \frac{1}{n} \sum_{i=0}^{n-1} (y_i - \hat{y}_i)^2$$

As we see, the cost function is an extension of the mathematical expression for expectation of y . Now in order to find the optimal parameters β giving the best approximation of f , we have to compute the first derivative of $C(\beta)$ with respect to β .

$$\frac{\partial C(\beta)}{\partial \beta} = (X^T(y - X\beta))$$

Now in order to find the optimal parameters, we set the first derivative equal to zero.

$$\Rightarrow X^T(y - X\beta) = 0$$

$$X^T y - X^T X \beta = 0$$

Now if we assume that $X^T X$ is an invertible matrix, we have the solution : $\beta = (X^T X)^{-1} X^T y$

$$\epsilon = y - \hat{y} = y - X\beta$$

$$\Rightarrow X^T(y - X\beta) = 0 \text{ we have } X^T \epsilon = X^T(y - X\beta) = 0$$

The usual assumption is to write the expected value of residuals error as equal to zero as its density distribution has a true mean of 0 according to earlier assumption.

I wish to elaborate further on the value of residuals in the proof. If we take a look at one of our earlier:

$$X^T e = X^T(y - X\beta) = 0$$

$$\Rightarrow X^T e = X^T y - X^T X \beta$$

Now inserting the OLS expected value for β , we get

$$\Rightarrow X^T e = X^T y - X^T X (X^T X)^{-1} X^T y = 0$$

$$= X^T y - X^T y = 0$$

Thus as we see, it's not only a simple expectation, but the value of the RSS, $E(e) = 0$. We can now continue on with our proof.

(In the final version of the report, remember to include the respective set that each of the variables included belongs to, $\Rightarrow y \in \mathbb{R}^n, X \in \mathbb{R}^{n \times p}, \beta \in \mathbb{R}^p$, and the value of $\hat{\beta} = (X^T X)^{-1} X^T y$)

From earlier we have written that:

$$E[y] = \frac{1}{n} \sum_{i=0}^{n-1} y_i$$

Now given X , each of its rows is symbolized by writing $X_{i,:}$. So the expectation value of y for a given element can be written:

$$E(y_i) = E[X_{i,:} \beta] + E[e_i]$$

$$= X_{i,:} \beta + 0 = X_{i,:} \beta \quad \text{and} \quad \sum_{j=0}^{p-1} X_{i,j} \beta_j$$

Continuing on, we now wish to prove that the variance of y_i for any given element i is $\text{Var}(y_i) = \sigma^2$. From the lecture notes and "Pattern Recognition and Machine Learning" by Christopher M. Bishop, variance can be written in the following way.

$$\begin{aligned}\text{Var}(y_i) &= \mathbb{E}[(y_i - \mathbb{E}(y_i))^2] = \mathbb{E}(y_i^2) - (\mathbb{E}(y_i))^2 \\ &= \mathbb{E}[(X_{i,\beta} + \epsilon_i)^2] - (X_{i,\beta})^2 \\ &= \mathbb{E}[(X_{i,\beta})^2 + 2X_{i,\beta}\epsilon_i + \epsilon_i^2] - (X_{i,\beta})^2 \\ &= (X_{i,\beta})^2 + 2\mathbb{E}(\epsilon_i)X_{i,\beta} + \mathbb{E}(\epsilon_i^2) - (X_{i,\beta})^2\end{aligned}$$

If we now look back at the information we have available about the distribution of residual error, we see that $\mathbb{E}(\epsilon_i) = 0$ and $\mathbb{E}(\epsilon_i^2) = \sigma^2$

$$= \mathbb{E}(\epsilon_i^2) = \text{Var}(\epsilon_i) = \underline{\sigma^2}$$

Moving on, we now have to show that $\mathbb{E}(\hat{\beta}) = \beta$ using ordinary least squares expression for the optimal parameters $\hat{\beta}$. From before, we have that

$$\frac{\partial C(\beta)}{\partial \beta} = X^T(y - X\beta)$$

$$\Rightarrow X^T y = X^T X \beta$$

Assuming that the matrix $X^T X$ is invertible, we can write the expression for optimal OLS parameters:

$$\hat{\beta}_{\text{OLS}} = (X^T X)^{-1} X^T y$$

If the matrix $X^T X$ causes any problems, by being near-singular, we can rewrite it using Singular Value Decomposition. Back to the proof, we can now write:

$$\begin{aligned}\mathbb{E}(\hat{\beta}) &= \mathbb{E}((X^T X)^{-1} X^T y) = (X^T X)^{-1} X^T \mathbb{E}(y) \\ &= (X^T X)^{-1} X^T X \underline{\beta} = \underline{\beta}\end{aligned}$$

Thus our model using ordinary least squares regression is unbiased. In other words, this simply means that the expected value is equal to the true value of the parameter. Having obtained a proof for $\mathbb{E}(\hat{\beta})$, we can now finish this part by showing that the variance of parameters β is

$$\text{Var}(\hat{\beta}) = \sigma^2 (X^T X)^{-1}$$

$$\begin{aligned}\Rightarrow \text{Var}(\hat{\beta}) &= \mathbb{E}[(\hat{\beta} - \mathbb{E}(\hat{\beta}))^2] = \mathbb{E}[(\hat{\beta} - \mathbb{E}(\hat{\beta}))(\hat{\beta} - \mathbb{E}(\hat{\beta}))^T] \\ &= \mathbb{E}[(X^T X)^{-1} X^T y - \beta][((X^T X)^{-1} X^T y - \beta)^T] \\ &= (X^T X)^{-1} X^T y y^T X (X^T X)^{-1 T} - (X^T X)^{-1} X^T y \beta^T - \beta y^T X (X^T X)^{-1 T} \\ &\quad + \beta \beta^T\end{aligned}$$

Before continuing, we can pay attention to the fact that we can use $\mathbb{E}(yy^T) = X \beta \beta^T X^T + \sigma^2 I$.

$$\begin{aligned}& (X^T X)^{-1} X^T \mathbb{E}(yy^T) X (X^T X)^{-1 T} - (X^T X)^{-1} X^T y \beta^T - \beta y^T X (X^T X)^{-1 T} + \beta \beta^T \\ &= (X^T X)^{-1} X^T \mathbb{E}(yy^T) X (X^T X)^{-1} - (X^T X)^{-1} X^T \mathbb{E}(y) \beta^T - \beta \mathbb{E}(y^T) X (X^T X)^{-1} + \beta \beta^T \\ &= (X^T X)^{-1} X^T \mathbb{E}(yy^T) X (X^T X)^{-1} - (X^T X)^{-1} X^T X \beta \beta^T - \beta \beta^T X^T X (X^T X)^{-1} + \beta \beta^T\end{aligned}$$

$$= (x^T x)^{-1} x^T E(y y^T) x (x^T x)^{-1} - \underline{\beta \beta^T}$$

The lecture notes significantly shorten further calculations. We will however write full calculations.

$$\rightarrow (x^T x)^{-1} x^T (\beta \beta^T x^T + \sigma^2 I) x (x^T x)^{-1} - \beta \beta^T$$

$$= ((x^T x)^{-1} x^T x \beta \beta^T x^T + (x^T x)^{-1} x^T \sigma^2 I) x (x^T x)^{-1} - \beta \beta^T$$

$$= (x^T x)^{-1} x^T x \beta \beta^T x^T x (x^T x)^{-1} + (x^T x)^{-1} x^T \sigma^2 I x (x^T x)^{-1} - \beta \beta^T$$

$$= \beta \beta^T + \sigma^2 (x^T x)^{-1} - \beta \beta^T = \underline{\sigma^2 (x^T x)^{-1}}$$