

LECTURE 2, PART I: REVIEW OF LINEAR REGRESSION

Text references: 401 course material

Model Basics and Assumptions

- Recall our model building block from last time:

$$Y = r(X_1) + \varepsilon,$$

where $\mathbb{E}(\varepsilon|X_1) = 0$. The regression function here is $r(X_1) = \mathbb{E}(Y|X_1)$, or $r(x) = \mathbb{E}(Y|X_1 = x)$. (We write X_1 here to reflect the fact that we just have one predictor, i.e., $X_1 \in \mathbb{R}$)

- In linear regression, we predict Y from a linear function of X_1 , of the form $\beta_0 + \beta_1 X_1$. If we determine β_0, β_1 by minimizing mean squared error,

$$\text{MSE}(\beta_0, \beta_1) = \mathbb{E}[(Y - \beta_0 - \beta_1 X_1)^2],$$

then recall from last time that

$$\beta_1 = \frac{\text{Cov}(X_1, Y)}{\text{Var}(X_1)}, \quad \beta_0 = \mathbb{E}(Y) - \beta_1 \mathbb{E}(X_1)$$

Multiple Linear Regression

What happens now with p predictors, X_1, \dots, X_p ?

We want to model Y as a linear function

Using MSE again as our criterion,

The optimal coefficients are given by

We will refer to these as the **population regression coefficients**.

Our multiple (linear regression) model is

$$Y = r(\mathbf{X}) + \varepsilon.$$

1 What are really our *assumptions* when using linear regression? Recall from
2 your regression class,

3 _____

4 _____

5 _____

6 _____

7 _____

8 _____

9 _____

10 _____

11 [See [catexample.R](#) for a linear regression example]

Briefly, note that we can summarize these assumptions as

$$Y|\mathbf{X} \sim N(\beta_0 + \boldsymbol{\beta}^T \mathbf{X}, \sigma^2)$$

12 Remarks on these assumptions:

13 _____

14 _____

15 _____

1
2
3
4
5
6
7
8
9
10
11
12

Linear Regression Estimates from Data

In practice, we don't have access to the distributions of \mathbf{X}, Y so we cannot actually compute the population regression coefficients. Instead, we have, say, an iid sample of size n ,

$$(\mathbf{X}_1, Y_1), (\mathbf{X}_2, Y_2), \dots, (\mathbf{X}_n, Y_n)$$

from some distribution of \mathbf{X} and Y . Note that each $\mathbf{X}_i \in \mathbb{R}^p$ and $Y_i \in \mathbb{R}$, but for simplicity, *we will from now on follow the notational convention of most statistical texts (including AOS, Witten et al., and Shalizi starting in Chapter*

3) and drop the bold-faced notation on the covariates (so a subscript i will denote **observation** i).

Hence, we write our linear model as

Now define the data matrix or the *design matrix* \mathbb{X} as

Each subject corresponds to one row. **The number of columns of \mathbb{X} corresponds to the number of features plus 1 for the intercept, $q = p + 1$.** Now define the vectors $\vec{Y}, \vec{\epsilon}, \beta$ as

$$\vec{Y} = \begin{pmatrix} Y_1 \\ Y_2 \\ \vdots \\ Y_n \end{pmatrix} \quad \vec{\epsilon} = \begin{pmatrix} \epsilon_1 \\ \epsilon_2 \\ \vdots \\ \epsilon_n \end{pmatrix} \quad \beta = \begin{pmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_p \end{pmatrix}. \quad (1)$$

Then, we can write our linear model more concisely as

Under squared error loss, we minimize the residual sums of squares (RSS)

which gives as the **sample regression coefficients** (or just regression coefficients) $\hat{\beta}$:

Think: What's the dimension of the matrix A ? What's the distribution of $\hat{\beta}_i$? Why, under what assumptions?

Properties of Least Squares Estimates

Theorem 1 *The estimators satisfy the following properties (conditional on \mathbb{X})*

1. $\mathbb{E}(\hat{\beta}) = \beta.$

2. $\mathbb{V}(\hat{\beta}) = \sigma^2(\mathbb{X}^T\mathbb{X})^{-1}$

3. $\hat{\beta} \approx MVN(\beta, \sigma^2(\mathbb{X}^T\mathbb{X})^{-1}).$

4. *An approximate $1 - \alpha$ confidence interval for β_j is*

$$\hat{\beta}_j \pm z_{\alpha/2} \widehat{\text{se}}(\hat{\beta}_j) \quad (2)$$

where $\widehat{\text{se}}(\hat{\beta}_j)$ is the square root of the appropriate diagonal element of the matrix $\hat{\sigma}^2(\mathbb{X}^T\mathbb{X})^{-1}$.

Exercise: Prove the first two assertions¹

¹Recall: Let Y be a random vector. Denote the mean vector by μ and the covariance matrix by Σ . If a is a vector then

$$\mathbb{E}(a^T Y) = a^T \mu, \quad \mathbb{V}(a^T Y) = a^T \Sigma a.$$

If A is a matrix then

$$\mathbb{E}(AY) = A\mu, \quad \mathbb{V}(AY) = A\Sigma A^T.$$

Q: Why are the last two assertions useful? (Give examples)

2 _____

3 _____

4 _____

5 _____

6 _____

7 _____

8 _____

9 _____

10 _____

11 _____