

# STA302 notes

May 18, 2022

## Contents

<b>May 9: Lecture 1</b>	<b>1</b>
Syllabus . . . . .	1
Review . . . . .	2
Introduction to Regression . . . . .	3
Class Afterthoughts/Questions . . . . .	4
<b>May 11: Lecture 2</b>	<b>5</b>
Regression continued . . . . .	5
Inferences about the regression model . . . . .	6
<b>May 16: Lecture 3</b>	<b>9</b>
Analysis of variance (ANOVA) . . . . .	9
Multiple Linear Regression . . . . .	11
<b>May 18: Lecture 4</b>	<b>12</b>
More properties . . . . .	12
Multiple Linear Regression Continued . . . . .	12

## May 9: Lecture 1

### Syllabus

This is a course on linear regression. The focus is using R to do data analysis, and build the mathematical foundation for regression. We will understand how prediction works later, which is the foundation for data science.

### Marking

- 2 HW - 15% each, due June 1, June 15
- Test - 25% on May 25
- Exam - 45% during June 22-27

**Books** J. Sheather, A Modern Approach to Regression w/ R and D. Montgomery, Linear Regression Analysis.

## Review

**Definition 1.** A **sample space**  $S$  is the set of possible events. A **random variable** is a function  $X: S \rightarrow \mathbb{R}$  assigning a number to elements of the sample space.

Constants can also be pseudo random variables. These are called **degenerate random variables** that have a **degenerate distribution** since they have infinite cdf.

**Definition 2.** For an event  $A \subset S$ , we define the **indicator function**  $I_A$  as

$$I_A(s) = \begin{cases} 1, & s \in A \\ 0, & s \notin A \end{cases}$$

These are important since we later use them to create dummy variables in linear regression. When we write an inequality involving random variables, we mean that it holds for all elements of the sample space. I.e.  $X \geq Y \implies X(s) = Y(s), \forall s \in S$ .

**Example 1.** Consider  $S = \{1, 2, 3, 4, 5, 6\}$ . For  $s \in S$ ,  $X(s) = s$ , let  $Y(s) = X(s) + I_6(s)$ . Then  $Y = X$  for all  $s \in S$  except 6, where  $Y = 7$ ,  $X = 6$ .

**Definition 3.** **Discrete r.v.** are functions from a countable sample space, and **continuous r.v.** are functions from an uncountable sample space. There are also **mixture** random variables, which are continuous/discrete for different parts of the sample space. Random variables can be univariate and multivariate as well.

**Example 2.** The multinomial distribution is an example of a discrete multivariate random variable.

**Definition 4.** If  $X$  is a random variable, the p.d.f. is the derivative of the c.m.f. As well,  $\mathcal{P}(a \leq X \leq b) = \int_a^b f(x) dx$  where  $f(x)$  is pdf. Similar thing holds for discrete r.v.

**Proposition 1.** The expectation of two random variables is linear. For  $Z = aX + bY$ ,  $X, Y$  r.v., then  $E(Z) = aE(X) + bE(Y)$ .

**Definition 5.** The **variance** of  $X$  is  $V(X) = E(X - \mu_x)^2$ . The **sample variance**  $s^2 = \frac{\sum (x_i - \bar{x})^2}{n-1}$ . Note we divide by  $n-1$  so that it is an unbiased estimator (STA261).

Some properties:

- $V(X) \geq 0$
- $V(aX + b) = a^2 V(x)$
- $V(X) = E(X^2) - E(X)^2$
- $V(X) \leq E(X^2)$
- $\sigma_X = \sqrt{V(X)}$

**Note:** In linear regression, the variance of the predicted variable depends on the slope of regression line but not on the intercept (second property).

Let  $X_1, X_2, Y$  be r.v. and  $A$  be an event. Let  $Z = aX_1 + bX_2$ . Then

- $E(Z | A) = aE(X_1 | A) + bE(X_2 | A)$
- $E(Z | Y = y) = aE(X_1 | Y = y) + bE(X_2 | Y = y)$
- $E(Z | Y) = aE(X_1 | Y) + bE(X_2 | Y)$

**Proposition 2.** (Laws of Total Expectation and Variance)  $E(E(Y | X)) = E(Y)$  and  $V(X) = V(E(X | Y)) + E(V(X | Y))$ .

We will see that linear regression is a conditional r.v., and the above will be very useful. For  $X_1, \dots, X_n$  i.i.d. random variables,  $x_1 \dots x_n$  realizations, then  $\bar{x} = \frac{\sum x_i}{n}$ . The **sample average**  $\bar{X} = \frac{\sum X_i}{n}$  is a random variable. In general, any function of  $n$  i.i.d. random variables is a random variable, and called a **sampling statistic** that follows a **sampling distribution**.

**Theorem 1.** (Central Limit Theorem) For  $X_1, \dots, X_n$  i.i.d.  $f(x, \theta)$ ,  $E(X)$ ,  $V(x) < \infty$ , then  $\frac{\bar{X} - \mu}{\sigma/\sqrt{n}} \rightarrow N(0, 1)$  converges in distribution for sufficiently large  $n$ .

*Proof.* Proof with moment generating functions. □

**Example 3.** In the Cauchy distribution, this does not hold since it has infinite mean and variance.

**Definition 6.** The **covariance**  $\text{Cov}(X, Y) = E[(X - \mu_x)(Y - \mu_y)] = E(XY) - E(X)E(Y)$ . Covariance quantifies the relationship between two variables, i.e. how much one varies with the other. The **correlation**  $\text{Corr}(X, Y) = \frac{\text{Cov}(X, Y)}{\sqrt{V(X)V(Y)}}$ .

- Covariance is an inner product, variance is norm.
- $V(X + Y) = V(X) + V(Y) + 2\text{Cov}(X, Y)$ .
- If  $X \perp Y$ ,  $V(X + Y) = V(X) + V(Y)$ .
- In general,  $V(\sum_i X_i) = \sum_i V(X_i) + 2\sum_{i < j} \text{Cov}(X_i, X_j)$ .

These will be useful in regression, where we try to identify relationships between r.v.s.

## Definitions in statistics

In probability, we are given a mathematical model to work with. In statistics, we infer properties of a mathematical model. The steps of data analysis are: state the problem, identify what data is needed, decide on a model and collect data, clean data, estimate parameters of the model, and carry out appropriate tests, draw conclusions.

## Introduction to Regression

**Definition 7.** The **corelation coefficient**

$$\rho_{X,Y} = \frac{\sum_i (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_i (x_i - \bar{x})^2} \sqrt{\sum_i (y_i - \bar{y})^2}} = \frac{\text{Cov}(X, Y)}{s_x s_y}$$

The above value is somewhat like the  $\cos(\theta)$  between the vectors  $X, Y$ ; recall dot product. When we discuss correlation, we talk about linear relations only; the linear association between  $X, Y$ . We can see this by considering  $X$  and  $Y = X^2$ . Correlation is symmetric, it does not indicate the direction of the symmetry (which causes which/causation). Correlation only says the influence on the change of one variable when the other changes; think about moving along non-orthogonal vectors and projecting.

Galton investigated the effect of fathers heights on their sons height. Galton termed **regression** as a 'regression' of heights towards the mean; on average, heights of sons move towards the mean, so the average height across generations is the same.

In a linear regression, we assume there is a linear relation  $Y = \beta_0 + \beta_1 X + \epsilon$  between the random variables  $X, Y$  where  $\epsilon$  is an error random variable. The deviation not captured by linearity is incorporated to  $\epsilon$ . Given two values of  $X$ , it is not guaranteed that the value of  $Y$  is the same. But for a unique  $X$  we get **unique average**  $Y$ . We want  $E(Y | X = x) = \beta_0 + \beta_1 X$ ; the relationship between the mean of  $Y$  and a specific value of  $X$  is linear. Note  $E(\epsilon) = 0$ . We call  $X$  the **explanatory, predictor, independent** variable and  $Y$  as the **response, outcome, dependent** variable. Suppose we are given paired data  $(x_1, y_1), \dots, (x_n, y_n)$ . We try to fit a linear regression to model the relationship between  $X$  and  $Y$ :

$$Y = \beta_0 + \beta_1 X + \epsilon \text{ and want } E(Y | X = x) = \beta_0 + \beta_1 X$$

The values of  $\beta_0, \beta_1$  are not yet known and need to be estimated. In the sample, the error  $e_i$  replaces  $\epsilon_i$ . The line best predicting  $Y$  as  $X$  changes should minimize the squares of the errors  $e_i = y_i - \hat{y}_i$  where  $\hat{y}_i = b_0 + b_1 x_i$  where  $b_0, b_1$  are the intercept and slope of the regression line. We minimize the squares  $\sum_i e_i^2$ . The  $e_i$  are referred to as **residuals**; minimize residual sums squared. Note

$$RSS(b_0, b_1) = \sum_i e_i^2 = \sum_{i=1}^n (y_i - \hat{y}_i)^2 = \sum_{i=1}^n (y_i - b_0 - b_1 x_i)^2$$

**Aside:** What value of  $a$  minimizes (1)  $\sum |x_i - a|$ , and which minimizes (2)  $\sum (x_i - a)^2$ ? Answer: (1)  $a = \text{Med}(X)$ , (2)  $a = \bar{x}$ . We do not minimize the sum of the residuals, since this must always be 0. We minimize the RSS with respect to  $b_0, b_1$ .

$$\frac{\partial RSS}{\partial b_0} = -2 \sum_i (y_i - b_0 - b_1 x_i), \quad \frac{\partial RSS}{\partial b_1} = -2 \sum_i x_i (y_i - b_0 - b_1 x_i)$$

so setting these to 0, we get the **normal equations**

$$\sum_i y_i = b_0 n + b_1 \sum_i x_i, \quad \sum_i x_i y_i = b_0 \sum_i x_i + b_1 \sum_i x_i^2$$

Solving these, we get

$$\hat{\beta}_0 = b_0 = \bar{y} - \hat{\beta}_1 \bar{x}, \quad \hat{\beta}_1 = b_1 = \frac{\sum_i x_i y_i - n \bar{x} \bar{y}}{\sum_i x_i^2 - n \bar{x}^2} = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2} = \frac{S_{X,Y}}{S_X}$$

The intercept is the average value of the response when  $X = 0$ .

## Class Afterthoughts/Questions

When the errors have  $E(\epsilon = 0)$ , then  $V(\epsilon) = E(\epsilon^2) - E(\epsilon)^2 = E(\epsilon^2)$ . By minimizing this in the sample, we minimize the variance of the errors (?)

## May 11: Lecture 2

**Clarifying last class:**  $\hat{y}_i$  is the conditional mean of  $y_i$ . When this is true, then  $\sum_i e_i = 0$ . That is, we estimate  $\hat{y}_i$  so that  $\sum_i e_i = 0$ .

### Regression continued

We continue discussing linear regression; fitting a linear relation assuming it exists. The aim is to infer the true values of  $\beta_0, \beta_1$  by inspecting their sampling distributions. We also make some assumptions regarding the error terms; the properties of their distributions ( $\epsilon$  is r.v.).

#### Assumption: Linearity

The conditional mean of  $Y | X = x$  is linear with respect to  $X$ . However, the relationship  $E(Y | X)$  and  $X$  does not have to be linear, but the linearity assumption is linearity in the parameters. Our relationship must be realistic given the context; introducing linearity may produce unrealistic relationships.

**R simulation:** When generating random dataset, we set a seed so our results are reproducible. Always start with a seed in assignments. Note the  $Y$  variable is the transformation  $\beta_0 + \beta_1 \log X + \epsilon$ . Introducing linear relationship between  $X$  and  $Y$  is inaccurate. It is linear in the parameters  $\beta_0, \beta_1$  however.

**Qs:** Chaos in random number generation? Look up random number generation algorithms. How do we quantify linearity in a data set? Mostly with plots but is there better way?

#### Assumption: Independence

The errors  $\epsilon_i$  are independent. That is, the deviations from the mean are not related; they are i.i.d. r.v. This reduces predictive capabilities in some areas, but we can relax this assumption later (generalized least squares).

#### Assumption: Homoscedasticity (equal variance)

The error variance does not change depending on  $X$ . That is  $V(\epsilon | X = x) = \sigma^2$  and is independent of  $x$ . In the R codes, we see that variance of errors increases with  $X$ , which decreases predictive power as  $X$  increases. Moreover, this implies some of the variation in the errors is explained by  $X$ , which violates our assumption. Variance **cannot** depend of  $X$ .  $\epsilon \perp X$ . This is relaxed in GLS.

In multiple linear regression, we talk about the Gauss-Markov assumption, but we need to make some assumptions about how  $\epsilon_i$  is distributed in order to make inferences.

#### Assumption: Normality

$\epsilon \sim N(0, \sigma^2)$ . The previous assumptions are required to obtain the least squares estimates, but normality is not required. Under this assumption, we can make confidence intervals and tests, and have nice properties following from normal distribution.

There are more assumptions in general, but these are most important.

### More about variance of $\epsilon$

We have estimated  $\beta_0, \beta_1$  using least squares. However, we have another parameter to estimate;  $V(\epsilon) = \sigma^2$ . From afterthoughts,  $V(\epsilon) = E(\epsilon^2) = \sigma^2$ . We take the average of  $e_i^2$  using this, since we want summary measure. The mean residual squared (MRS) can be calculated as  $s^2 = \frac{\sum_i e_i^2}{n-2}$ . We show this estimator of  $E(\epsilon^2)$  is unbiased as homework; prove this!.

### Inferences about the regression model

#### Conditional expectation and variance of $\hat{\beta}_1$

Recall  $\beta_1 = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2}$

**Proposition 3.**  $\sum (x_i - \bar{x})(y_i - \bar{y}) = \sum_i (x_i - \bar{x})y_i$

*Proof.*

$$\begin{aligned} \sum (x_i - \bar{x})(y_i - \bar{y}) &= \sum (x_i y_i - \bar{x} y_i - \bar{y} x_i + \bar{x} \bar{y}) \\ &= \sum (x_i y_i - \bar{x} y_i) - n \bar{y} \bar{x} + n \bar{x} \bar{y} \\ &= \sum (x_i - \bar{x}) y_i \end{aligned}$$

□

A symmetric sum can be established for  $\sum_i (y_i - \bar{y})x_i$ . However, the above is needed to simplify conditional expectation calculations. We may also show  $\sum (x_i - \bar{x})x_i = \sum (x_i - \bar{x})^2$ . The idea of both of these proof is making the substitution  $n\bar{x} = \sum x_i$ .

**Proposition 4.**  $\sum (x_i - \bar{x})x_i = \sum (x_i - \bar{x})^2$

*Proof.*

$$\begin{aligned} \sum (x_i - \bar{x})x_i &= \sum (x_i^2 - \bar{x}x_i) \\ &= \sum (x_i^2 - 2\bar{x}x_i) + n\bar{x}^2 \\ &= \sum (x_i - \bar{x})^2 \end{aligned}$$

□

Other way of writing:  $\sum (x_i - \bar{x})^2 = \sum x_i^2 - n\bar{x}^2$ . Now, we calculate **conditional expectation of  $\hat{\beta}_1$**

$$E(\hat{\beta}_1 | X = x_i) = E\left(\frac{\sum (x_i - \bar{x})y_i}{\sum (x_i - \bar{x})^2} \mid X = x_i\right) = \frac{\sum (x_i - \bar{x})E(Y_i | X = x_i)}{\sum (x_i - \bar{x})^2}$$

Substituting  $E(Y_i | X_i = x) = \beta_0 + \beta_1 x$ , then

$$E(\hat{\beta}_1 | X = x_i) = \frac{\sum_i (x_i - \bar{x})\beta_0}{\sum (x_i - \bar{x})^2} + \frac{\sum_i (x_i - \bar{x})\beta_1 x_i}{\sum (x_i - \bar{x})^2} = \frac{\beta_1 \sum_i (x_i - \bar{x})^2}{\sum (x_i - \bar{x})^2} = \beta_1$$

Since  $\sum (x_i - \bar{x}) = \sum x_i - n\bar{x} = 0$  and by above prop.,  $\sum_i (x_i - \bar{x})x_i = \sum (x_i - \bar{x})^2$ . Therefore  $\hat{\beta}_1$  does not depend on  $X$ , and has expected value of  $\beta_1$ ; it is an unbiased estimator of  $\beta_1$ . That is,  $E(\hat{\beta}_1 | X = x_i) = E(\hat{\beta}_1) = \beta_1$ . Next, we may calculate  $V(\hat{\beta}_1)$ . First,  $V(Y_i | X = x_i) = \sigma^2$ , that is, the variance of the error.

$$V(\hat{\beta}_1 | X = x_i) = \left( \frac{\sum (x_i - \bar{x})y_i}{\sum (x_i - \bar{x})^2} \mid X = x_i \right) = \frac{\sum_i (x_i - \bar{x})^2 V(Y_i | X = x_i)}{(\sum_i (x_i - \bar{x})^2)^2} = \frac{\sigma^2}{\sum (x_i - \bar{x})^2} = \frac{\sigma^2}{S_{X,X}}$$

### Inferences for variance of $\hat{\beta}_1$

Since  $\epsilon_i \sim N(0, \sigma^2)$ , then  $Y_i | X \sim N(\beta_0 + \beta_1 X, \sigma^2)$ . Letting  $c_i = \frac{\sum (x_i - \bar{x})}{\sum (x_i - \bar{x})^2}$  then  $\hat{\beta}_1 = \sum c_i y_i$ . Observe that this is a **linear combination** of normally distributed random variables, so  $\hat{\beta}_1$  is normally distributed! Thus

$$\hat{\beta}_1 | X = x_i \sim N\left(\beta_1, \frac{\sigma^2}{S_{X,X}}\right)$$

We can construct a  $1 - \alpha$  confidence interval for  $\beta_1$  which has extremes  $\hat{\beta}_1 \pm Z_{1-\alpha/2} \frac{\sigma}{\sqrt{S_{X,X}}}$ . When  $\sigma^2$  is unknown, we construct a  $t$ -confidence using  $S^2 = \frac{\sum e_i^2}{n-2}$ . We therefore make a confidence interval with critical values

$$\hat{\beta}_1 \pm t_{1-\alpha/2, n-2} \frac{s^2}{\sqrt{S_{X,X}}}$$

Note our assumption of normality of errors.

**Clarification**  $S_{X,X} = \sum (x_i - \bar{x})^2$  and  $S_{X,Y} = \sum (x_i - \bar{x})(y_i - \bar{y})$ .

Recall, the **p-value** can be calculated as  $p = \mathcal{P}(Z \geq |z|)$  or  $p = \mathcal{P}(T \geq |t|)$  where  $z, t$  are the calculated test statistics. The p-value is the probability of obtaining a sample that provides strong evidence against the hypothesized value of  $H_0 : \beta_1$ , set by threshold  $\alpha$ .  $\alpha$  is the probability of making a type one error with repeated sampling.

**Example 4.**  $\sum x_i = 4035$ ,  $\sum y_i = 4041$ ,  $\sum e_i^2 = 4753.125$ ,  $\sum x_i^2 = 1005535$ ,  $\sum x_i y_i = 864910$ ,  $t_{0.975, 18} = 2.10$ .

We need to calculate  $\hat{\beta}_1, s, S_{X,X}$  from this information; recall  $\hat{\beta}_1 \pm t_{1-\alpha/2, n-2} \frac{s^2}{\sqrt{S_{X,X}}}$ . The interval becomes (0.18121, 0.33728). **Verify as homework.**

Do exercises from Montgomery (unassigned, do by chapter) and Sheather. Problems are similar to this, and this will appear on the midterm.

### Properties of $\beta_0$

The conditional expectation of  $\beta_0 | X$ . Since  $\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}$ . Using this,

$$E(\hat{\beta}_0 | X = x_i) = \frac{\sum E(y_i | X = x_i)}{n} - \beta_1 \bar{x} = \left( \frac{n\beta_0 + n\beta_1 \bar{x}}{n} \right) - \beta_1 \bar{x} = \beta_0$$

Therefore  $\hat{\beta}_0$  is an unbiased estimator of  $\beta_0$ . Now for the variance, (minor abuse of notation)

$$V(\hat{\beta}_0 | X = x_i) = V(\bar{y} - \hat{\beta}_1 \bar{x} | X = x_i) = V(\bar{y} | x_i) + \bar{x}^2 V(\hat{\beta}_1 | x_i) - 2\bar{x} \text{Cov}(\bar{y}, \hat{\beta}_1 | x_i)$$

Calculating each term separately,

$$V(\bar{y} | X = x_i) = V\left(\frac{\sum y_i}{n} | X = x_i\right) = \frac{\sum \sigma^2}{n^2} = \frac{\sigma^2}{n}$$

To calculate covariance term, we use substitutions involving  $\hat{\beta}_1 = \sum c_i y_i$  with  $c_i$  defined before the break

$$\text{Cov}(\bar{y}, \hat{\beta}_1 | X = x_i) = \text{Cov}\left(\frac{\sum_i y_i}{n}, \sum c_i y_i | X = x_i\right) = \frac{1}{n} \sum_i \text{Cov}(y_i, c_i y_i | X = x_i)$$

Recall  $\text{Cov}(X, aY) = a\text{Cov}(X, Y)$ . Also, given a particular  $x_i$ ,  $c_i$  is a constant.

$$= \frac{1}{n} \sum_i c_i \text{Cov}(y_i, y_i | X = x_i) = \frac{1}{n} \sum_i c_i V(y_i | X = x_i) = \frac{1}{n} \sum_i c_i \sigma^2 = 0$$

From last section,  $V(\hat{\beta}_1 | x_i) = \bar{x}^2 \frac{\sigma^2}{S_{X,X}}$ . Therefore

$$V(\hat{\beta}_0 | X = x_i) = \sigma^2 \left( \frac{1}{n} + \frac{\bar{x}^2}{S_{X,X}} \right), \text{ and } \hat{\beta}_0 | X = x_i \sim N\left(\beta_0, \sigma^2 \left( \frac{1}{n} + \frac{\bar{x}^2}{S_{X,X}} \right)\right)$$

Therefore the  $(1 - \alpha)$  confidence for  $\beta_0$  is

$$\hat{\beta}_0 \pm Z_{1-\alpha/2} \sigma \sqrt{\frac{1}{n} + \frac{\bar{x}^2}{S_{X,X}}}$$

(fill in when  $\sigma^2$  is unknown )

### Confidence interval for the regression line

Denote  $x^*, y^*$  as an observation in the future. I.e., an observation not currently in the sample, and we use the model built with the current observation. (Prediction) It can easily be shown that  $E(\hat{y}^* | X = x^*) = \beta_0 + \beta_1 x^*$ .  $X = x^*$  new observation,  $y^*$  unknown. As well,  $\hat{y}^*$  is the predicted value of  $y^*$  paired with  $x^*$ . Often, we are interested in calculating the variance of  $\hat{y}^* | X = x^*$  and confidence interval  $E(Y | X = x^*)$ . That is, calculate the variance and c.i. of the regression line. Note  $E(\hat{y}^* | X = x^*) = \beta_0 + \beta_1 x^* = E(Y | X = x^*)$  implies the sample regression is an unbiased estimator of the TRUE Linear relationship between  $X, Y$ . The variance can be calculated as

$$V(\hat{y}^* | X = x^*) = V(\hat{\beta}_0 + \hat{\beta}_1 x^*)$$

(etc fill in the slide) (edit this lol)

### Prediction error and interval

Assuming we fit a regression line between  $X, Y$  with some sample. If a new data point  $X = x^*$  is given, our predicted  $\hat{y}^*$  lies exactly on the line in the model we have fitted, but  $y^*$  associated with  $x^*$  may deviate from the line. How much does this  $y$  vary?  $y^* - \hat{y}^*$  is called the **prediction error** for  $X = x^*$ . We calculate its expectation and variance.

For expectation, the  $*$  is redundant, so we write  $E(y - \hat{y} | X = x^*)$ . We can easily show this is 0 since  $y - \hat{y} = 0$  (show this later) Therefore

$$V(y^* - \hat{y}^* | X = x^*) = V(y - \hat{y} | X = x^*) = \sigma^2 \left( 1 + \frac{1}{n} + \frac{(x^* - \bar{x})^2}{S_{X,X}} \right)$$



We just add the variance of  $y$  and variance of  $\hat{y}$  by expansion of variance and since  $\text{Cov}(\hat{y}, y) = 0$ . The observation  $y$  is independent of the previous sample by assumption. The prediction interval is built in the same way as before using  $t$  distribution.

**R simulation:**

The confidence interval is for the regression line. The prediction interval is for a new predicted value given  $x^*$ ; how far  $y^*$  can deviate from the predicted  $\hat{y}^*$ .

**Example 5.** Calculate summary measures for the production data (in slides hw)

## May 16: Lecture 3

**Clarification** In the derivations from last class, we used

$$\text{Cov}\left(\frac{\sum Y_i}{n}, \sum c_i Y_i \mid X = x_i\right) = \frac{1}{n} \sum \text{Cov}(y_i, c_i y_i \mid X = x_i)$$

since  $\text{Cov}(Y_i, Y_j) = 0$  by independence of  $Y_i, Y_j$ .

Understand theory and problem solving procedure for midterms. Data analysis will mostly be with R.

**Assignment Task 1**

The purpose of the assignment is using R for inference of parameters given simulated data. Use your student id as a seed. After data is generated, run the LM model. Repeating this procedure, get sampling distribution for  $\hat{\beta}_i, \sigma^2$ , and compare these to true variances.

**Analysis of variance (ANOVA)**

So far we have discussed inference about specific parameters, and hypothesis testing for their true values. For example, if we fail to reject  $H_0 : \beta_1 = 0$ , then there is no linear relationship between  $X, Y$ . In this case,  $Y = \beta_0 + \epsilon$ ,  $V(Y) = V(\epsilon) = \sigma^2$ , so  $\epsilon$  explains all the variance of  $Y$ . Usually,  $V(Y) = \beta_1^2 V(X) + \sigma^2$ , since  $X \perp \epsilon$ . Therefore when the above holds, part of the variance is given by  $V(X)$ . If most of the variation in  $Y$  is explained by  $X$ , then predictions are very accurate. We discuss this in ANOVA.

In the slides, points that are less scattered about the regression line have more of their variance explained by  $X$ . Explaining the variance of  $Y$  is very important; consider Anscombe's 4 datasets.

As the residual variance  $\sigma^2$  increases, the variation of  $Y$  is less explained by  $X$ . This increases prediction error. We want to answer how well the regression line might explain the variation we observe in the responses. ANOVA is another way of testing the significance of the regression line. The total variation of  $Y$  is explained by the **total sum of squares**, the numerator of  $s_Y$

$$SST = \sum (y_i - \bar{y})^2$$

This can be decomposed by

$$\sum (y_i - \bar{y})^2 = \sum (y_i - \hat{y}_i + \hat{y}_i - \bar{y})^2 = \sum (y_i - \hat{y}_i)^2 + \sum (\hat{y}_i - \bar{y})^2 + 2 \sum (y_i - \hat{y}_i)(\hat{y}_i - \bar{y})$$

Where the third term becomes

$$\sum (y_i - \hat{y}_i)(\hat{y}_i - \bar{y}) = \sum (\hat{y}_i(y_i - \hat{y}_i) - \bar{y}(y_i - \hat{y}_i)) = \sum \hat{y}_i e_i - \bar{y} \sum e_i = 0$$

Since  $\sum e_i = 0$  and  $\sum x_i e_i = 0$  by the second normal equation, which gives  $\sum \hat{y}_i e_i = 0$ . Hint:  $\sum (\beta_0 + \beta_1 x_i) e_i = \beta_0 \sum e_i + \beta_1 \sum x_i e_i$ . Therefore the total variation of  $Y$  can be divided into

$$\sum (y_i - \bar{y})^2 = \sum (y_i - \hat{y}_i)^2 + \sum (\hat{y}_i - \bar{y})^2$$

The term on the left is the **residual sum square**,  $(n-2)s^2$ . The second term explains the variance in  $\hat{y}_i$ , or the variation in fitted values from the regression. We may easily show  $\sum \frac{\hat{y}_i}{n} = \bar{y}$ . The second term on the right is the **regression sum squared**. The total variation in  $Y$  has been decomposed to come from the regression line, and from random errors.

**Degrees of Freedom.** This is the number of summed square normals. The proof for  $\frac{(n-1)s^2}{\sigma^2} \sim \chi_{n-1}^2$  shows where one of the ‘standard normal squares’ are lost. ( $s^2$  is sample variance). For each parameter we fix, we lose a degree of freedom. When  $\bar{y}$  is fixed, we are free to have  $n-1$  values, and are forced to choose one to get the fixed  $\bar{y}$ . That is,  $y_n$ , the  $n$ -th observation is fixed for a fixed  $\bar{y}$ . This is why sample variance,  $\sum (y_i - \bar{y})^2 / (n-1)$ , uses  $n-1$  degrees of freedom.

In the above SST, the **RSS**  $\sum (y_i - \hat{y}_i)^2$  has  $n-2$  degrees of freedom since  $\hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_i$  uses two estimated parameters. Since  $\sum (y_i - \bar{y})^2$  has  $n-1$  degrees of freedom, then the  $SS_{reg} \sum (\hat{y}_i - \bar{y})^2$  must have 1 degree of freedom. This follows since the sum depends only on  $\beta_1$  given fixed  $x_i$ :

$$\sum (\hat{y}_i - \bar{y})^2 = \sum (\hat{\beta}_0 + \hat{\beta}_1 x_i - \bar{y})^2 = \sum (\bar{y} - \hat{\beta}_1 \bar{x} + \hat{\beta}_1 x_i - \bar{y})^2 = \sum \hat{\beta}_1^2 (x_i - \bar{x})^2$$

We need degrees of freedom in order to test hypothesis. We will later show

$$\frac{SS_{reg}}{\sigma^2} \sim \chi_1^2, \quad \frac{RSS}{\sigma^2} \sim \chi_{n-2}^2$$

Under  $H_0 : \beta_0 = 0$  then  $F_0 \sim F_{1,n-2}$ . We want  $SS_{reg}$  as close to the SST as possible. The F-test here detects how close  $SS_{reg}$  is to TSS. The closer it is the bigger the value of  $F_0$ . We can show  $t_{n-2}^2 = F_{1,n-2}$ . We can also show

$$E(SS_{reg}) = \sigma^2 + S_{X,X} \beta_1^2$$

So when  $\beta_1 = 0$ , the regression sum squared have variance equal to  $\sigma^2$ . Below is an ANOVA table:

Sources of Variation	Df	Sum Sq	Mean Sq	F value	Pr(>F)
Regression	1	$SS_{reg}$	$MS_{reg} = \frac{SS_{reg}}{1}$	$F_0 = \frac{MS_{reg}}{MRSS}$	etc
Residuals	n-2	$RSS$	$MRSS_{reg} = \frac{RSS}{n-2}$		
Total	n-1	SST			

In general, the F-test measures whether the means of two groups measure significantly. The F statistic is the ratio of explained variance (regression model attributes to  $V(X)$ ) to unexplained variance (variance of  $e_i$ ). Under the null, our data reflects the intercept only model  $Y = \beta_0 + \epsilon$ , and we test the departure from this.

## The Coefficient of Determination

Another measure to assess whether the regression line explains enough of the variability in the response is the **coefficient of determination**,  $R^2$ . This gives the proportion of the total sample variability in the response that has been explained by the regression model.

$$R^2 = \frac{SS_{reg}}{SST} \text{ or } 1 - R^2 = \frac{RSS}{SST}$$

Note  $0 \leq R^2 \leq 1$ . If  $R^2$  is close to 1, it is an important predictor of  $Y$ . If it is close to 0, then it offers little predictive power for  $Y$ . In simple linear regression,  $\rho^2 = R^2$  where  $\rho$  is Pearson correlation coefficient.

## Categorical predictors

So far we have required  $X$  to be continuous. However,  $X$  could be categorical. ( $X$  smoking status vs.  $Y$  blood pressure). Here the predictor is binary and the output is continuous. How would we test if the mean blood pressure varies between these groups?

We did this in STA261 with a two-sample t-test, and by homoscedasticity we do one with equal variance. We may also use regression, by using **dummy variables** which are indicator variables. Setting 0 for non-smokers, 1 for smokers,

$$E(Y | X = 0) = \beta_0, E(Y | X = 1) = \beta_0 + \beta_1$$

Using ANOVA this is essentially a t-test.  $F_{1, n-2} \sim t_{n-1}^2$  so by squaring the  $t$  statistic we get  $F$  statistic; a significant  $F$  statistic indicate the change in means given by  $\beta_1$  is significant. Therefore using hypothesis test with ANOVA for  $\beta_1 = 0$ , we get a test for differing means.

The 'slope' becomes the change in average. We can say  $\beta_1$  reflects the average difference between two groups. The slope provides the magnitude of the difference, while the hypothesis test tells us whether the difference is statistically significant.

With categorical variables,  $R^2$  may be low but the test will give significance.

## Multiple Linear Regression

So far we have only had one predictor  $X$ , but we generalize to  $X_1, \dots, X_n$ . That is

$$Y = \beta_0 + \beta_1 X_1 + \dots + \beta_p X_p + \epsilon$$

This implies  $Y$  is related to  $X_1, \dots, X_p$  linearly. However, the predictor produces a  $p$ -dimensional subspace instead of a line. See image in 'Elements of Statistical Learning 2e'; with  $Y$  regressed on  $X_1, X_2$  we get a regression plane.

The conditional mean of  $Y$  is given by  $E(Y | X_1, \dots, X_p) = \beta_0 + \beta_1 X_1 + \dots + \beta_p X_p$ . For the sample dataset,

$$y_i = \beta_0 + \beta_1 x_{i,1} + \dots + \beta_p x_{i,p} + e_i$$

So we solve  $RSS(\beta_0, \dots, \beta_p) = \sum (y_i - \sum \beta_j x_{i,j})^2$  (dude fill this in)

## Matrix Notation

In order to simplify notation we use matrices. For this we write

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}$$

$\mathbf{Y}$  is an  $n \times 1$  vector,  $\mathbf{X}$  is an  $n \times (p+1)$  matrix, with the first column being a vector of 1s.  $\boldsymbol{\beta}$  is  $(p+1) \times 1$  vector,  $\boldsymbol{\epsilon}$  is  $n \times 1$  vector.

We denote the transpose of matrix  $\mathbf{A}$  as  $\mathbf{A}'$ . If  $\mathbf{A}$  is a square matrix with  $\mathbf{A} = \mathbf{A}'$  then it is symmetric (corresponds to self adjoint operator). If  $\mathbf{A}$  is invertible, we denote its inverse with  $\mathbf{A}^{-1}$ . A matrix is **orthogonal** if  $\mathbf{A}^{-1} = \mathbf{A}'$ ; column vectors are orthogonal. An **idempotent** matrix satisfies  $\mathbf{A}^2 = \mathbf{A}$ . Some important properties are that

$$(\mathbf{A} + \mathbf{B})' = \mathbf{A}' + \mathbf{B}' \text{ and } (\mathbf{AB})' = \mathbf{B}'\mathbf{A}'$$

**Example 6.** The projection matrix  $P : \mathbb{R}^n \rightarrow \mathbb{R}^n$  of rank  $p \leq n$  onto a subspace is a square matrix that is symmetric and idempotent.

## May 18: Lecture 4

### More properties

**Definition 8.** If  $Y = (Y_1, \dots, Y_n)$  is a random vector, then  $E(Y) = (E(Y_1), \dots, E(Y_n))$ . The **covariance matrix** of  $Y$  is denoted

$$V(Y) = \begin{pmatrix} V(Y_1) & \text{Cov}(Y_1, Y_2) & \dots & \text{Cov}(Y_1, Y_n) \\ \text{Cov}(Y_2, Y_1) & V(Y_2) & \dots & \text{Cov}(Y_2, Y_n) \\ \vdots & & \ddots & \vdots \\ \text{Cov}(Y_n, Y_1) & \text{Cov}(Y_n, Y_2) & \dots & V(Y_n) \end{pmatrix}$$

That is each entry  $a_{ij} = \text{Cov}(Y_i, Y_j)$ . It is created by  $\text{Cov}\{(Y - E(Y))(Y - E(Y))'\}$ , the outer product. If  $b$  is a vector, then  $V(b'Y) = b'V(Y)b$ .

### Multiple Linear Regression Continued

Above, we wrote  $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}$ , that is  $y_i = \beta_0 + \beta_1 x_{i,1} + \dots + \beta_p x_{i,p} + \epsilon_i$  in matrix form. Explicitly,

$$\begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix} = \begin{pmatrix} 1 & x_{1,1} & x_{1,2} & \dots & x_{1,p} \\ 1 & x_{2,1} & x_{2,2} & \dots & x_{2,p} \\ \vdots & & \ddots & & \vdots \\ 1 & x_{n,1} & x_{n,2} & \dots & x_{n,p} \end{pmatrix} \begin{pmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_p \end{pmatrix} + \begin{pmatrix} \epsilon_1 \\ \epsilon_2 \\ \vdots \\ \epsilon_n \end{pmatrix}$$

$\mathbf{Y}, \boldsymbol{\epsilon} \in \mathbb{R}^n$ ,  $\boldsymbol{\beta} \in \mathbb{R}^{p+1}$ , and  $\mathbf{X}$  is  $n \times (p+1)$  dimensional.

As before, we would like to minimize  $\sum_i^n e_i^2$  given values in  $X$ . This evaluates to the scalar

$$RSS(\boldsymbol{\beta}) = \sum_i^n e_i^2 = e'e = (Y - X\boldsymbol{\beta})'(Y - X\boldsymbol{\beta}) = Y'Y - 2Y'X\boldsymbol{\beta} + \boldsymbol{\beta}'X'X\boldsymbol{\beta}$$

Where  $Y'X\beta = \beta'X'Y$  since the transpose of a scalar is the same scalar. Note  $RSS : \mathbb{R}^{p+1} \rightarrow \mathbb{R}$  Differentiating with respect to  $\beta$ ,

$$\frac{\partial RSS}{\partial \beta} = \frac{\partial}{\partial \beta}(Y'Y - 2\beta'X'Y + \beta'X'X\beta) = -2X'Y + 2X'X\beta$$

Setting this to 0, we see  $\hat{\beta} = (X'X)^{-1}X'Y$ . In the case of simple LR,

$$X = \begin{pmatrix} 1 & x_1 \\ \vdots & \vdots \\ 1 & x_n \end{pmatrix}, Y = \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix} \Rightarrow X'X = \begin{pmatrix} n & \sum x_i \\ \sum x_i & \sum x_i^2 \end{pmatrix} = n \begin{pmatrix} 1 & \bar{x} \\ \bar{x} & \frac{1}{n} \sum x_i^2 \end{pmatrix}$$

We can compute  $\det X'X = n^2 \cdot \left( \frac{1}{n} \sum x_i^2 - \bar{x}^2 \right) = n \cdot \sum (x_i - \bar{x})^2 = n \cdot S_{X,X}$ . Therefore

$$(X'X)^{-1} = \begin{pmatrix} \frac{\sum x_i^2}{n \cdot S_{X,X}} & -\frac{\bar{x}}{S_{X,X}} \\ -\frac{\bar{x}}{S_{X,X}} & \frac{1}{S_{X,X}} \end{pmatrix}$$

Multiplying by  $\sigma^2$ , we see this is the **covariance matrix for  $\hat{\beta}_0, \hat{\beta}_1$** ;  $\text{Cov}(\hat{\beta}_0, \hat{\beta}_1) = \frac{-\sigma^2 \bar{x}}{S_{X,X}}$ . **Important for midterm!**

**Definition 9.** The **projection** of  $Y$  on  $X$  is given by  $\hat{Y} = X\hat{\beta} = X(X'X)^{-1}X'Y = HY$ . We call  $H$  the **hat** or **projection** matrix. Note it is  $n \times n$ , idempotent, and symmetric!

We let  $e = Y - \hat{Y} = Y - X(X'X)^{-1}X'Y = (I - H)Y$

**Proposition 5.**  $H$  and  $I - H$  are both idempotent.

Note that  $HX = X$ ; this is easily checked by tracing definition and cancelling inverses. We can partition the first  $k$  and last  $p + 1 - k$  columns of  $X$  into matrix  $[X_1, X_2]$ . Then  $HX = [HX_1, HX_2] = X = [X_1, X_2]$ . As well,  $\text{tr}(H) = p + 1$  and  $\dim \text{range } H = p + 1$ .

### Assumptions in Multiple LR

$E(Y | X) = X \cdot \beta$ . Linearity, independence, homoscedasticity, normality hold as assumptions for our model (same as before). We assume  $\epsilon \sim N(0, \sigma^2 I)$ . Then  $Y | X \sim N(X\beta, \sigma^2 I)$ . Now we discuss the distribution of  $\hat{\beta}$ .

$$E(\hat{\beta} | X) = E((X'X)^{-1}X'Y | X) = (X'X)^{-1}X'X\beta = \beta$$

so the estimator is consistent. For the variance, we carry out adjoints as in previous property

$$V(\hat{\beta} | X) = V((X'X)^{-1}X'Y | X) = (X'X)^{-1}X'\sigma^2 I X(X'X)^{-1} = (X'X)^{-1}\sigma^2$$

This is just the covariance matrix of  $\hat{\beta}$ ! Look back to our example above. That is

$$C = (X'X)^{-1} \Rightarrow c_{ij} = \sigma^2 \text{Cov}(\beta_i, \beta_j)$$

Least squares estimates are the **best linear unbiased estimators** according to the Gauss-Markov Theorem (which is stated later). The following assumptions are required for the theorem: (1) the errors

$\epsilon_i$  are independent, (2)  $E(\epsilon) = 0$ , (3)  $V(\epsilon) = \sigma^2$ . Note normality is **not** assumed.

As in simple LR, the  $\hat{\beta}_j$  are normally distributed;  $\hat{\beta}_j \sim N(\beta_j, \sigma^2 c_{j,j})$ . We can test hypotheses for  $\beta_j$  in the usual way. Given  $H_0 : \beta_j^0$ , then we can calculate  $Z = \frac{\hat{\beta}_j - \beta_j^0}{\sqrt{c_{j,j}}\sigma}$  and use a z-test. ANOVA for MLR will be next time.