

# Lecture 8: Linear Regression

May 24, 2012  
GENOME 560, Spring 2012

Su-In Lee, CSE & GS  
suinlee@uw.edu

1

## Goals

- Develop basic concepts of linear regression from a probabilistic framework
- Estimating parameters and hypothesis testing with linear models
- Linear regression in R

2

## Regression

- Technique used for the modeling and analysis of numerical data
- Exploits the relationship between two or more variables so that we can gain information about one of them through knowing values of the other
- Regression can be used for prediction, estimation, hypothesis testing, and modeling causal relationships

3

## Why Linear Regression?

- Suppose we want to model the outcome variable  $Y$  in terms of three predictors,  $X_1, X_2, X_3$

$$Y = f(X_1, X_2, X_3)$$

- Typically will not have enough data to try and directly estimate  $f$
- Therefore, we usually have to assume that it has some restricted form, such as **linear**

$$Y = X_1 + X_2 + X_3$$

4

## Regression Terminology

$$Y = X_1 + X_2 + X_3$$

Dependent Variable

Independent Variable

Outcome Variable

Predictor Variable

Response Variable

Explanatory Variable

Lung cancer risk

Genetic factor, smoking,  
diet, etc

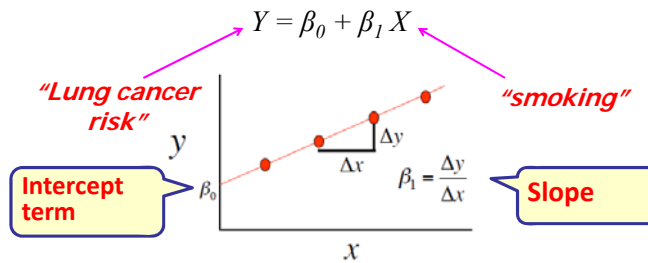
Expression level  
of gene X

Expression levels of X's TFs  
A, B and C

5

## Linear Regression is a Probabilistic Model

- Much of mathematics is devoted to studying variables that are deterministically related to one another



- But we're interested in understanding the relationship between variables related **in a nondeterministic fashion**

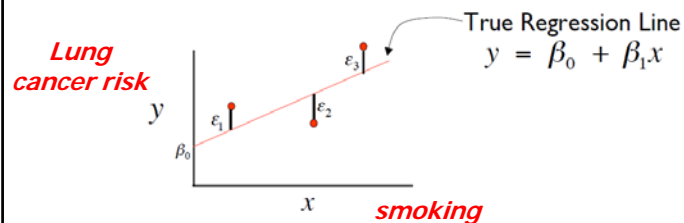
6

## A Linear Probabilistic Model

- **Definition:** There exists parameters  $\beta_0, \beta_1$  and  $\sigma^2$ , such that for any fixed value of the predictor variable  $X$ , the outcome variable  $Y$  is related to  $X$  through the model equation

$$Y = \beta_0 + \beta_1 X + \varepsilon,$$

where  $\varepsilon$  is a RV assumed to be  $N(0, \sigma^2)$



7

## Implications

- The **expected value of Y** is a linear function of X, but for fixed value  $x$ , the variable Y differs from its expected value by a **random amount**

**Variables & Symbols: How is  $x$  different from X?**

**Capital letter X:** a random variable

**Lower case letter x:** corresponding values

(i.e. the real numbers the RV  $X$  map into)

For example,

$X$ : Genotype of a certain locus

$x$ : 0, 1 or 2 (meaning AA, AG and GG, respectively)

8

## Implications

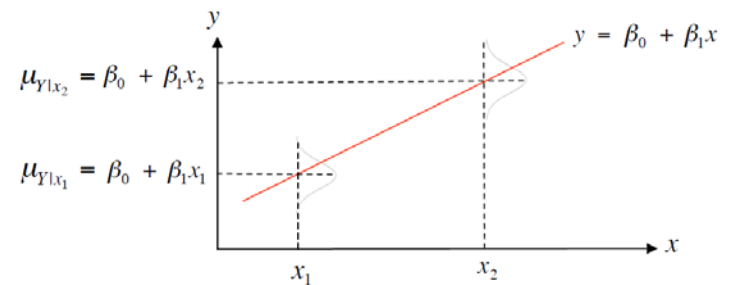
- The **expected value of  $Y$**  is a linear function of  $X$ , but for fixed value  $x$ , the variable  $Y$  differs from its expected value by a **random amount**
- Formally, let  $x^*$  denote a particular value of the predictor variable  $X$ , then our linear probabilistic model says:

$$E(Y | x^*) = \mu_{Y|x^*} = \text{mean value of } Y \text{ when } X \text{ is } x^*$$

$$V(Y | x^*) = \sigma_{Y|x^*}^2 = \text{variance of } Y \text{ when } X \text{ is } x^*$$

9

## Graphical Interpretation

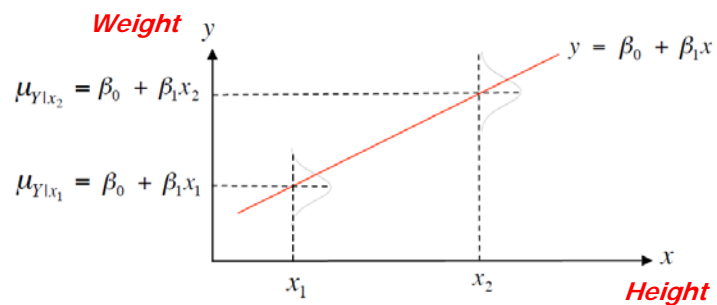


$$E(Y | x^*) = \mu_{Y|x^*} = \text{mean value of } Y \text{ when } X \text{ is } x^*$$

$$V(Y | x^*) = \sigma_{Y|x^*}^2 = \text{variance of } Y \text{ when } X \text{ is } x^*$$

10

## Graphical Interpretation



- Say that  $X = \text{height}$  and  $Y = \text{weight}$
- Then  $\mu_{Y|x=60}$  is the average weight for all individuals 60 inches tall in the population

11

## One More Example

- Suppose the relationship between the predictor variable height ( $X$ ) and outcome variable weight ( $Y$ ) is described by a simple linear regression model with true regression line

$$Y = 7.5 + 0.5X, \quad \varepsilon \sim N(0, \sigma^2) \text{ and } \sigma = 3$$

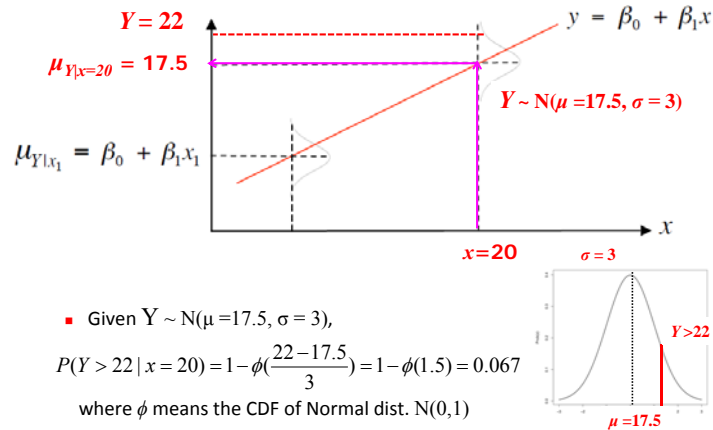
- Q1: What is the interpretation of  $\beta_1 = 0.5$ ?  
The expected change in weight ( $Y$ ) associated with a 1-unit increase in height ( $X$ )
- Q2: If  $x = 20$ , what is the expected value of  $Y$ ?

$$\mu_{Y|x=20} = 7.5 + 0.5(20) = 17.5$$

12

## One More Example

- Q3: If  $x = 20$ , what is  $P(Y > 22)$ ?



## Estimating Model Parameters

- Where are the parameters  $\beta_0$  and  $\beta_1$  from?

- Predicted**, or fitted, values are values of  $y$  predicted by plugging  $x_1, x_2, \dots, x_n$  into the estimated regression line:  $y = \beta_0 + \beta_1 x$

$$\hat{y}_1 = \beta_0 + \beta_1 x_1$$

$$\hat{y}_2 = \beta_0 + \beta_1 x_2$$

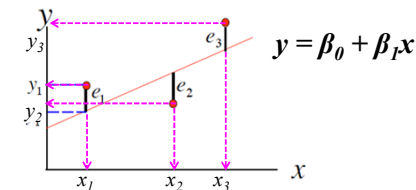
$$\hat{y}_3 = \beta_0 + \beta_1 x_3$$

- Residuals** are the deviations of observed (red dots) and predicted values (red line)

$$e_1 = y_1 - \hat{y}_1$$

$$e_2 = y_2 - \hat{y}_2$$

$$e_3 = y_3 - \hat{y}_3$$



14

## Residuals Are Useful!

- The error sum of squares (SSE) can tell us how well the line fits to the data

$$SSE = \sum_{i=1}^n (e_i)^2 = \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

$y = \beta_0 + \beta_1 x$

$y_1$

$y_2$

$y_3$

$\hat{y}_1$

$\hat{y}_2$

$\hat{y}_3$

$x_1$

$x_2$

$x_3$

$x$

$e_1$

$e_2$

$e_3$

- Least squares**

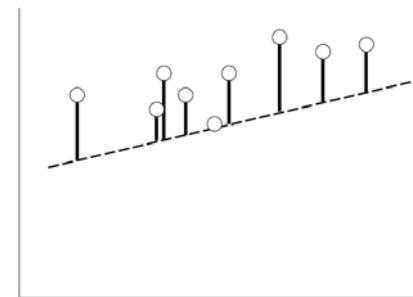
- Find  $\beta_0$  and  $\beta_1$  that minimizes SSE

$$f(\beta_0, \beta_1) = \sum_{i=1}^n [y_i - (\beta_0 + \beta_1 x_i)]^2$$

- Denote the solutions by  $\hat{\beta}_0$  and  $\hat{\beta}_1$

15

## Least Squares



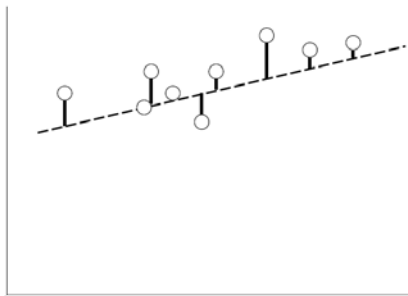
- Least squares**

- Find  $\beta_0$  and  $\beta_1$  that minimizes SSE

$$f(\beta_0, \beta_1) = \sum_{i=1}^n [y_i - (\beta_0 + \beta_1 x_i)]^2$$

16

## Least Squares



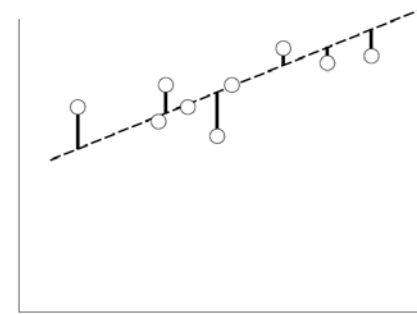
### Least squares

- Find  $\beta_0$  and  $\beta_1$  that minimizes SSE

$$f(\beta_0, \beta_1) = \sum_{i=1}^n [y_i - (\beta_0 + \beta_1 x_i)]^2$$

17

## Least Squares



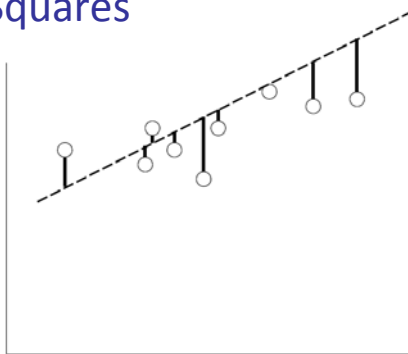
### Least squares

- Find  $\beta_0$  and  $\beta_1$  that minimizes SSE

$$f(\beta_0, \beta_1) = \sum_{i=1}^n [y_i - (\beta_0 + \beta_1 x_i)]^2$$

18

## Least Squares



### Least squares

- Find  $\beta_0$  and  $\beta_1$  that minimizes SSE

$$f(\beta_0, \beta_1) = \sum_{i=1}^n [y_i - (\beta_0 + \beta_1 x_i)]^2$$

19

## Coefficient of Determination

- Important statistic referred to as the coefficient of determination ( $R^2$ ):

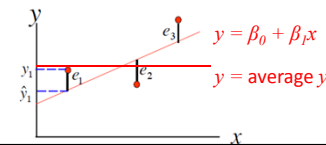
$$R^2 = 1 - \frac{SSE}{SST}$$

$$SSE = \sum_{i=1}^n (e_i)^2 = \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

Error Sum Squares

$$SST = \sum_{i=1}^n (y_i - \bar{y})^2$$

Error Sum Squares, when  $\beta_0 = \text{avg}(y)$  and  $\beta_1 = 0$



20

## Multiple Linear Regression

- Extension of the simple linear regression model to two or more independent variables

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_n x_n + \varepsilon$$

Expression = Baseline + Age + Tissue + Sex + Error

- Partial Regression Coefficients:**

$\beta_i$   $\equiv$  effect on the outcome variable when increasing the  $i^{\text{th}}$  predictor variable by 1 unit, **holding all other predictors constant**

21

## Categorical Independent Variables

- Qualitative variables are easily incorporated in regression framework through **dummy variables**
- Simple example: sex can be coded as 0/1
- What if my categorical variable contains three levels:

$$X_i = \begin{cases} 0 & \text{if AA} \\ 1 & \text{if AG} \\ 2 & \text{if GG} \end{cases}$$

- NO!

22

## Categorical Independent Variables

- Previous coding would result in **collinearity**
- Solution is to set up a series of dummy variable. In general for  $k$  levels you need  $(k-1)$  dummy variables

$$X_1 = \begin{cases} 1 & \text{if AA} \\ 0 & \text{otherwise} \end{cases} \quad X_2 = \begin{cases} 1 & \text{if AG} \\ 0 & \text{otherwise} \end{cases}$$

$X_i = \begin{cases} 0 & \text{if AA} \\ 1 & \text{if AG} \\ 2 & \text{if GG} \end{cases}$	if AA	$X_1$	$X_2$
	AA	1	0
	AG	0	1
	GG	0	0

23

## Hypothesis Testing: Model Utility Test

- The first thing we want to know after fitting a model is **whether any of the predictor variables ( $X$ 's) are significantly related to the outcome variable ( $Y$ ):**

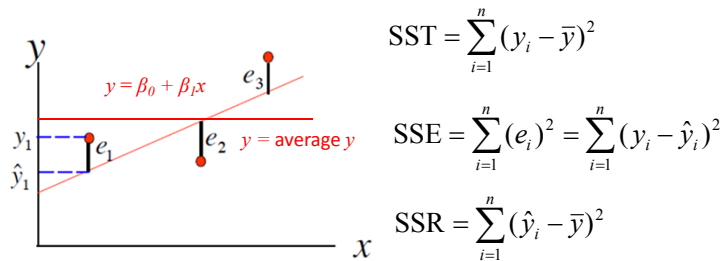
$$H_0 : \beta_1 = \beta_2 = \dots = \beta_k = 0$$

$$H_A : \text{At least one } \beta_i \neq 0$$

- Let's frame this in our ANOVA framework
- In ANOVA, we partitioned total variance (SST) into two components:
  - SSE (unexplained variation)
  - SSR (variation explained by linear model)

## Model Utility Test

- Partition total variance (SST) into two components:
  - SSE (unexplained variation)
  - SSR (variation explained by linear model)
- Let's consider  $n (=3)$  data points and  $k (=1)$  predictor model



25

## ANOVA Formulation of Model Utility Test

- Partition total variance (SST) into two components:
  - SSE (unexplained variation)
  - SSR (variation explained by linear model)

Source of Variation	df	Sum of Squares	MS	F
Regression	k	$SSR = \sum (\hat{y}_i - \bar{y})^2$	$\frac{SSR}{k}$	$\frac{MS_R}{MS_E}$
Error	$n-(k+1)$	$SSE = \sum (y_i - \hat{y}_i)^2$	$\frac{SSE}{n-(k+1)}$	
Total	$n-1$	$SST = \sum (y_i - \bar{y})^2$		

# data points - (# parameters in the model) Rejection Region :  $F_{\alpha, k, n-(k+1)}$

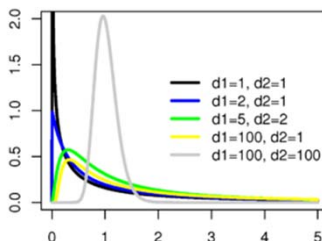
26

## ANOVA Formulation of Model Utility Test

- F-test statistic

$$F = \frac{MS_R}{MS_E} = \frac{SSR / k}{SSE / [n - (k + 1)]} = \frac{R^2}{1 - R^2} \cdot \frac{n - (k + 1)}{k}$$

Rejection Region :  $F_{\alpha, k, n-(k+1)}$



- Pick the distribution function, based on  $k$  and  $n-(k+1)$ .
- Choose the critical value based on  $\alpha$  ( $F_{\alpha, k, n-(k+1)}$ )
  - Say that  $\alpha = 0.05$
  - $\text{Prob}(F > F_{\alpha, k, n-(k+1)}) = 0.05$

27

## F Test For Subsets of Independent Variables

- A powerful tool in multiple regression analysis is the ability to compare two models

- For instance say we want to compare

$$\text{Full Model: } y = \beta_0 + \beta_1x_1 + \beta_2x_2 + \beta_3x_3 + \beta_4x_4 + \varepsilon$$

$$\text{Reduced Model: } y = \beta_0 + \beta_1x_1 + \beta_2x_2 + \varepsilon$$

- Again, another example of ANOVA

$SSE_R$  = error sum of squares for reduced model with  $l$  predictors

$SSE_F$  = error sum of squares for full model with  $k$  predictors

$$F = \frac{(SSE_R - SSE_F) / (k - l)}{SSE_F / [n - (k + 1)]}$$

28

## Example of Model Comparison

- We have a quantitative trait and want to test the effects at two markers, M1 and M2.

Full Model: Trait = Mean + M1 + M2 + (M1 × X2) + ε

Reduced Model: Trait = Mean + M1 + M2 + ε

$$F = \frac{(SSE_R - SSE_F)/(k-1)}{SSE_F/[n-(k+1)]} = \frac{(SSE_R - SSE_F)/(3-2)}{SSE_F/[100-(3+1)]}$$

$$= \frac{(SSE_R - SSE_F)}{SSE_F/96}$$

Rejection Region :  $F_{\alpha,1,96}$

29

## How To Do In R

- You can fit a least-squares regression using the function
  - `mm <- lsfit(x,y)`
- The coefficients of the fit are then given by
  - `mm$coef`
- The residuals are
  - `mm$residuals`
- And to print out the tests for zero slope just do
  - `ls.print (mm)`

30

## Input Data

- <http://www.cs.washington.edu/homes/suinlee/genome560/data/cats.txt>
- Data on fluctuating proportions of marked cells in marrow from heterozygous Safari cats
- Proportions of cells of one cell type in samples from cats (taken in our department many years ago). Column 1 is the ID number of the particular cat. You will want to plot the data from one cat.
  - For example cat 40004 is rows 1:17, 40005a is 18:31, 40005b is 32:47, 40006 is 48:65, 40665 is 66:83 and so on.

31

## Input Data

- 2<sup>nd</sup> column:** Time, in weeks from the start of monitoring, that the measurement from marrow is recorded.
- 3<sup>rd</sup> column:** Percent of domestic-type progenitor cells observed in a sample of cells at that time.
- 4<sup>th</sup> column:** Sample size at that time, i.e. the number of progenitor cells analyzed.

	40004	11	33	72
	40004	13	49	67
	40004	19	46	56
	40004	25	42	19
	40004	28	68	59
	40004	31	55	64
	40004	33	38	61
	40004	36	23	73
	40004	41	32	170
	40004	45	41	120
	40004	48	50	70
	40004	50	54	39
	40004	52	30	143
	40004	54	30	56
	40004	56	32	78
	40004	58	18	74
	40004	62	36	81
	40005a	14	34	65
	40005a	17	26	74
	40005a	23	21	73
	40005a	26	11	72
	40005a	29	19	77
	40005a	31	20	70
	40005a	34	13	56
	40005a	37	17	65