# Linear Regression-1

Prof. Asim Tewari
IIT Bombay

# Supervised learning

- Other factors:

1. Data heterogeneity: some algorithms require input to be numerical and scaled to similar range (eg. SVM, NN methods. On the other hand decision tree etc. Can handle heterogeneous data.

2. Redundancy in data (highly correlated input variables)

3. Presence of interactions and non-linearities

Eg. Of supervised learning :linear regression, logistic regression, SVM, naive Bayes, linear discriminant analysis, decision tree, k-nearest neighbor algo, neural network etc.

# Unsupervised learning

- Test of inferring a function to describe hidden structure.
- Eg. - Clustering (K-mean, mixture models, hierarchies clustering)

  - Anomaly detection etc.

# Regression

$$y = f(x)$$

$$(\bar{x}_i, y_i)$$

$$\bar{x}_i = \begin{Bmatrix} x_i^1 \\ x_i^2 \\ \vdots \\ x_i^p \end{Bmatrix}$$

- Extract a relationship from data

Learning arbitrary function is intractable

$H$ is a class of functions

$\Rightarrow$ parametric function  $\quad w \in R^d$

$\uparrow$

Parameters or the weights

$$h_w(x) = g(x, w)$$

Learning is to find $w$ from data $(x, y)$

# Regression

$$h_\omega(x) = g(x, \omega)$$

Case 1 : $g_1$ : $\omega_0 + \omega_1 x$ | $\omega_0, \omega_1$

Case 2 : $g_2$ : $\omega_0 + \omega_1 x + \omega_2 x^2$ | $\omega_0, \omega_1, \omega_2$

Case 3 : $g_3$ : $\omega_0 + \omega_1 x + \omega_2 x^2 + \omega_3 x^3$ | $\omega_0, \omega_1, \omega_2, \omega_3$
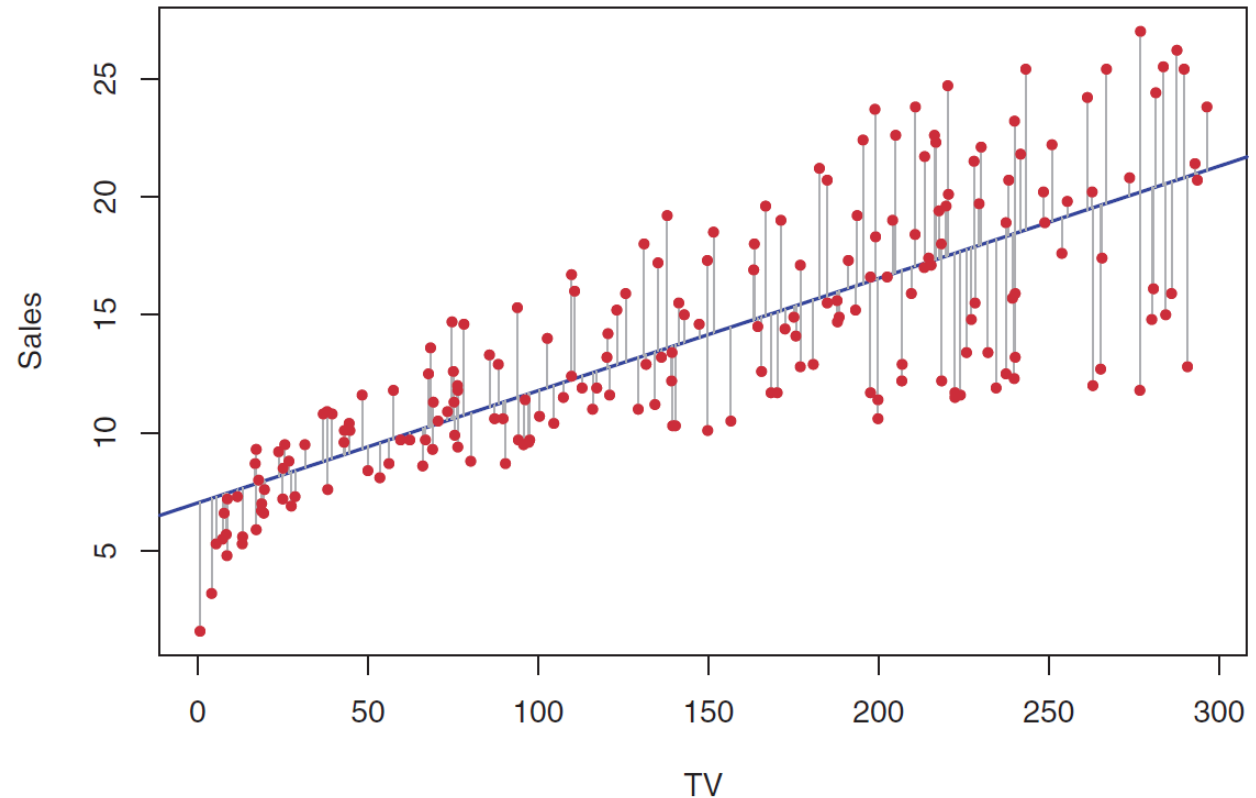
# Linear Regression

$$Y = f(X) + \epsilon \quad \leftarrow \text{Error/Noise}$$

$$Y = \beta_0 + \beta_1 X + \epsilon$$

- It assumes a linear relation between input x and output   y

$$Y \approx \hat{\beta}_0 + \hat{\beta}_1 X$$

Approximately modeled $\hat{\beta}0$ and $\hat{\beta}1$ are unknown coefficients or parameters which are estimated from training data.

# Linear Regression

$$\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x$$

$$y = \beta_0 + \beta_1 x + \epsilon$$

$\epsilon$ error (zero mean)

$$x_i = \begin{bmatrix} x_{i1} \\ x_{i2} \\ \vdots \\ x_{ip} \end{bmatrix}, \quad x_i \in R^p$$

Prediction (estimate) of Y

**Estimating the coefficients:**     $P = 1$

- Training data (**x1**, **y1**), (**x2**, **y2**)... and (**xn**, **yn**)

- n data pair

- we need $\hat{\beta}_0$ and $\hat{\beta}1$ ^ such that the linear model fits the data well

- measure of data fits the data well or closeness?

- One possible closeness measure is least square criterion

$$e_i = y_i - \hat{y}_i$$     - represents ith residual

# Linear Regression

*Residual sum of squares* (RSS)

$$\text{RSS} = e_1^2 + e_2^2 + \cdots + e_n^2$$

Minimize RSS is least square criterion

$$\text{RSS} = (y_1 - \hat{\beta}_0 - \hat{\beta}_1 x_1)^2 + (y_2 - \hat{\beta}_0 - \hat{\beta}_1 x_2)^2 + \ldots + (y_n - \hat{\beta}_0 - \hat{\beta}_1 x_n)^2$$
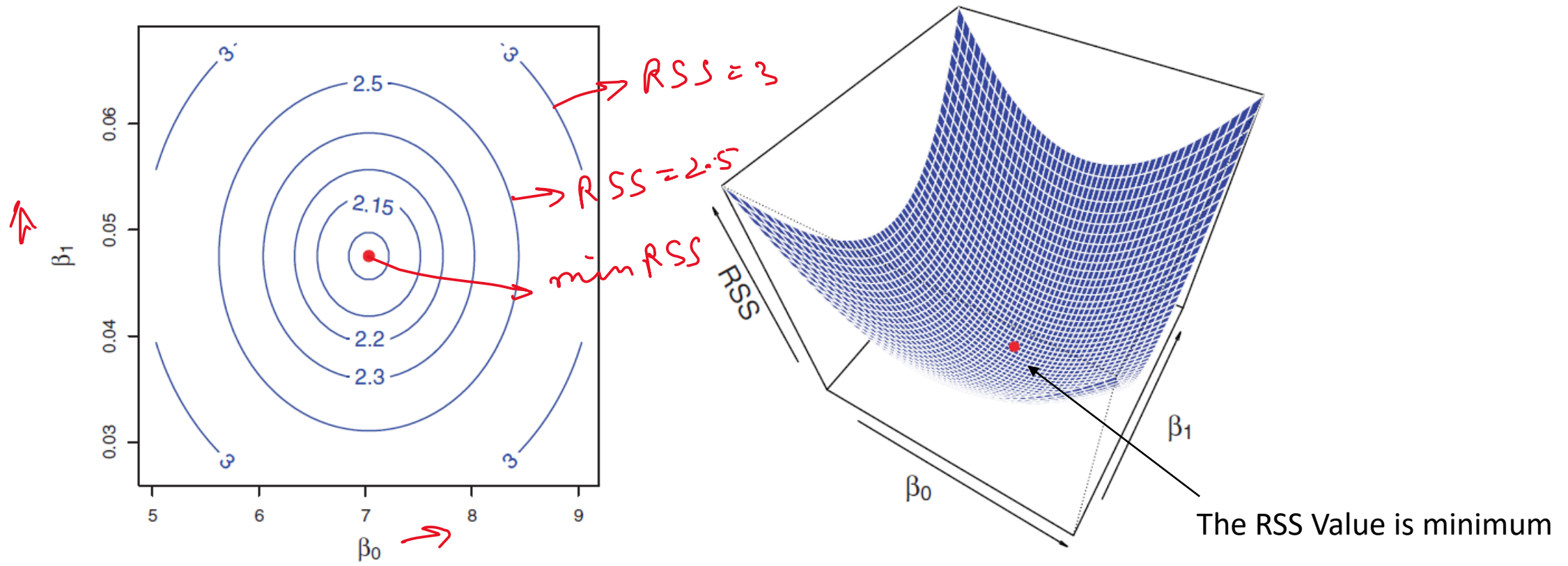
By minimizing RSS we can find

$$\hat{\beta}_1 = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2},$$

$$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}, \qquad \text{where } \bar{y} \equiv \frac{1}{n}\sum_{i=1}^n y_i \text{ and } \bar{x} \equiv \frac{1}{n}\sum_{i=1}^n x_i$$
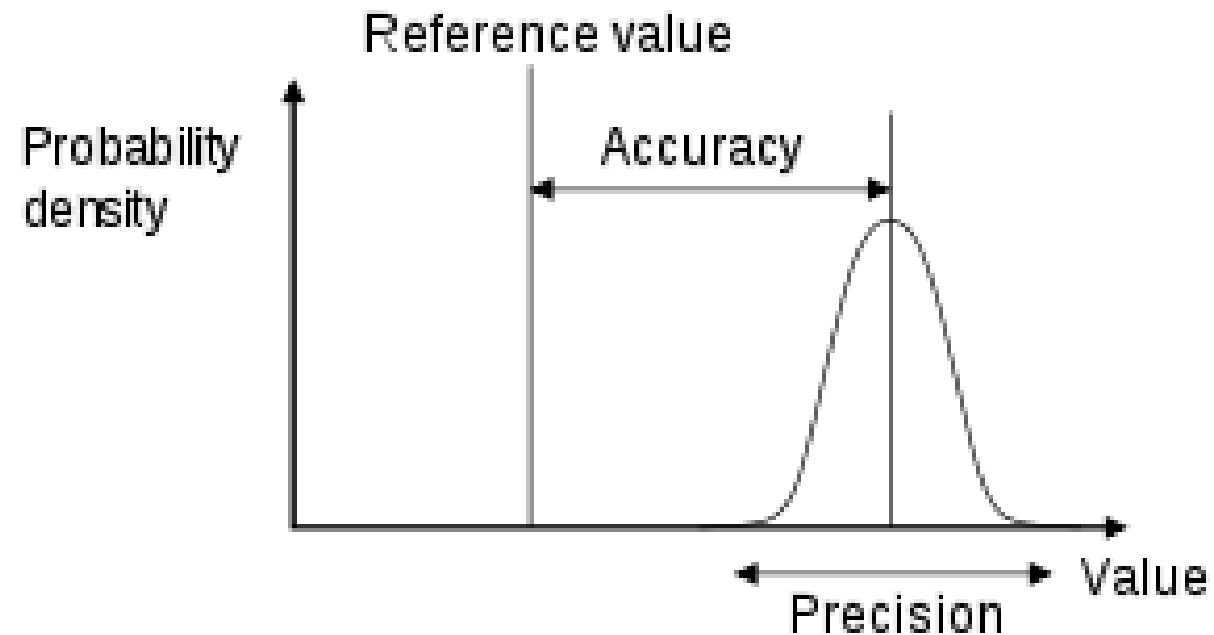
$$e_i = y_i - \hat{y}_i$$
$$\hookrightarrow \hat{\beta}_0 + \hat{\beta}_1 x_i$$

$$\frac{\partial \text{RSS}}{\partial \hat{\beta}_1} = 0, \quad \frac{\partial \text{RSS}}{\partial \hat{\beta}_0} = 0$$

# Linear Regression

$$RSS = f(\hat{\beta_0}, \hat{\beta_1}) \; ;$$



$\rightarrow RSS = 3$

$\rightarrow RSS = 2.5$
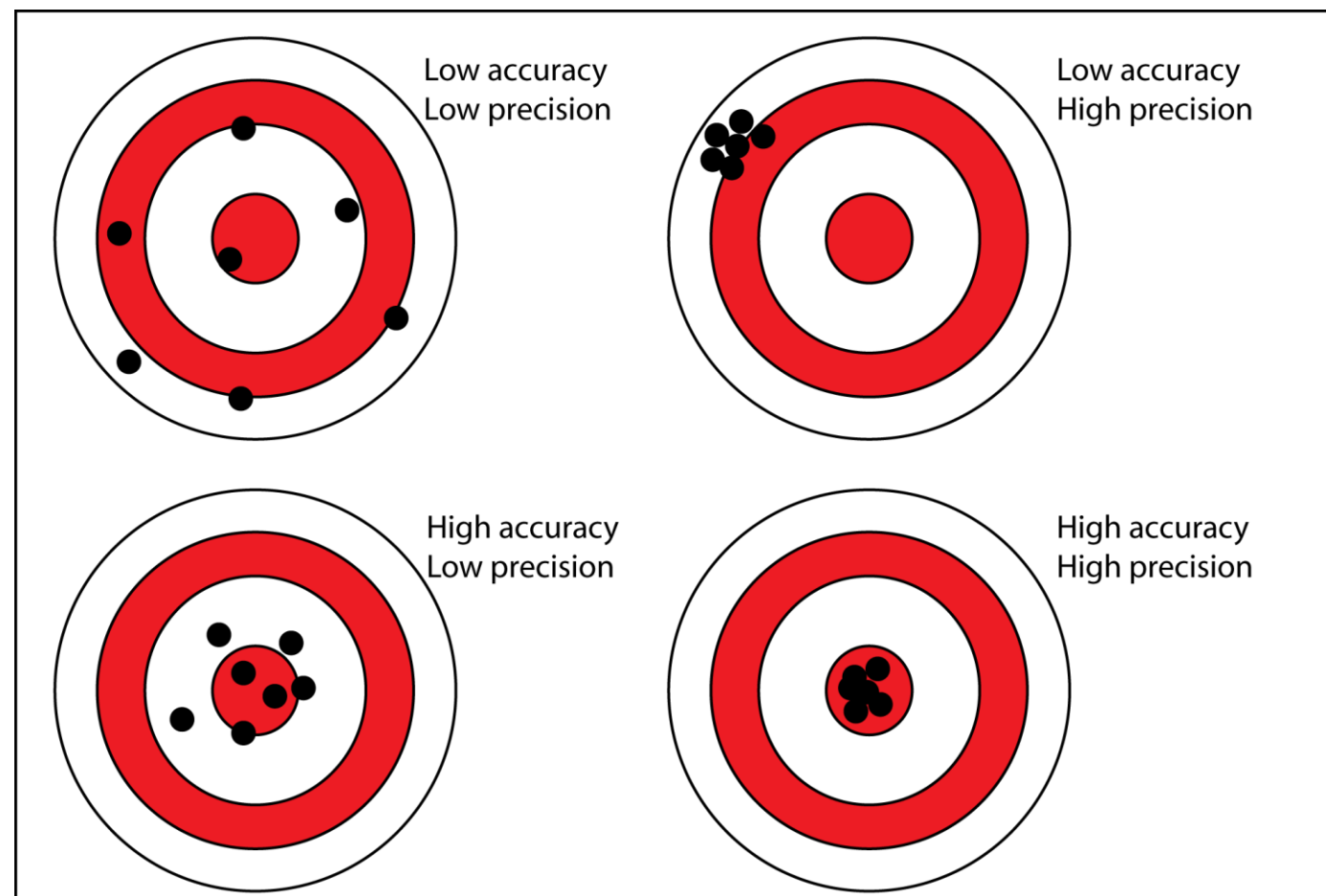
min RSS

The RSS Value is minimum

# Accuracy and Precision

- Accuracy refers to the closeness of a measured value to a standard or known value. Accuracy is a description of systematic errors, a measure of statistical bias.

- Precision refers to the closeness of two or more measurements to each other. Precision is a description of random errors, a measure of statistical variability.

# Accuracy Vs. Precision

# Linear Regression

- True relation between X and Y

  Y=f(x) + €

  ↑ Mean zero

random error Independent of x

If f(x) is linear then $Y = \beta_0 + \beta_1 X + \epsilon.$

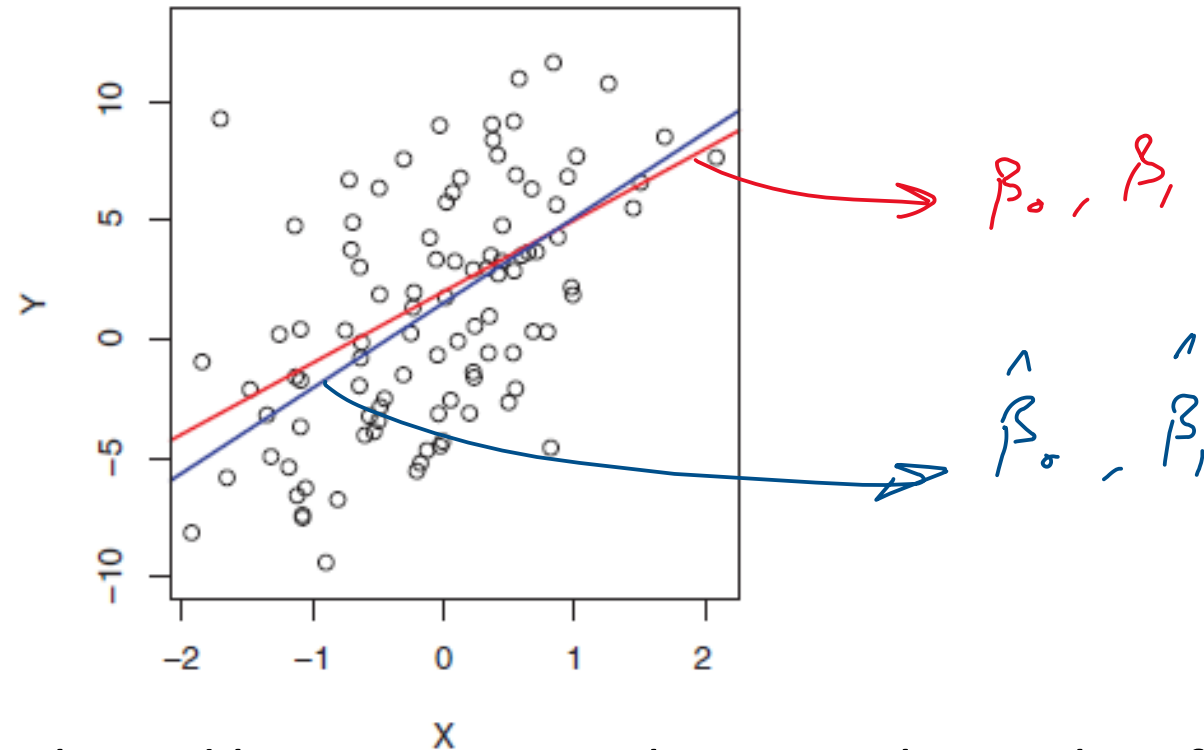$\hat{\beta}_0$ and $\hat{\beta}_1$ are analogous to estimation of population mean from sample mean

I.e $\hat{\beta}_0$ an $\hat{\beta}_1$ are unbiased estimate of true $\beta_0$ and $\beta_1$

$$Y = \beta_0 + \beta_1 X + e$$

$$\hat{Y} = \hat{\beta}_0 + \hat{\beta}_1 X$$

1.) How close is $\hat{\beta}_0$ to $\beta_0$ and $\hat{\beta}_1$ to $\beta_1$?

2.) How close is $\hat{y}$ to $y$?

# Linear Regression
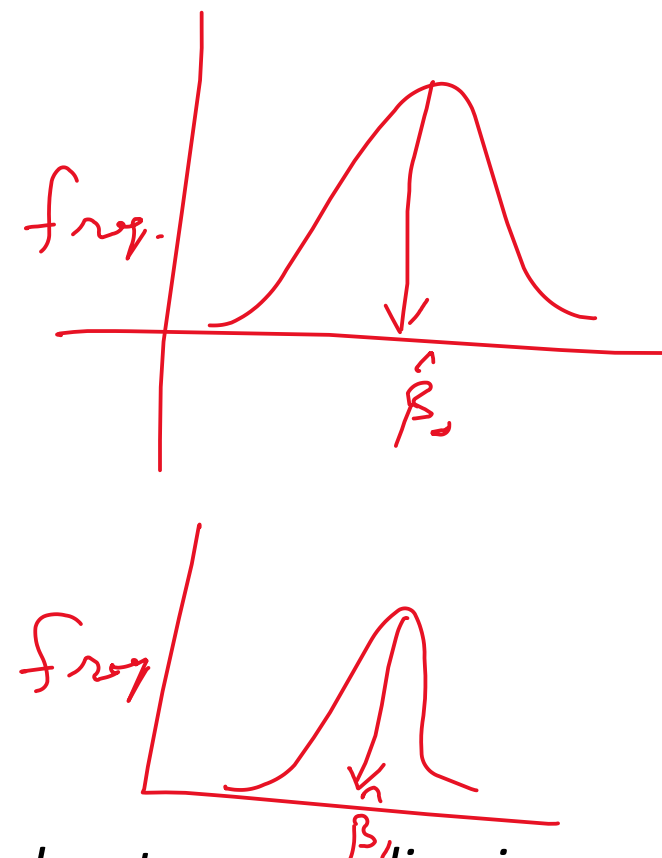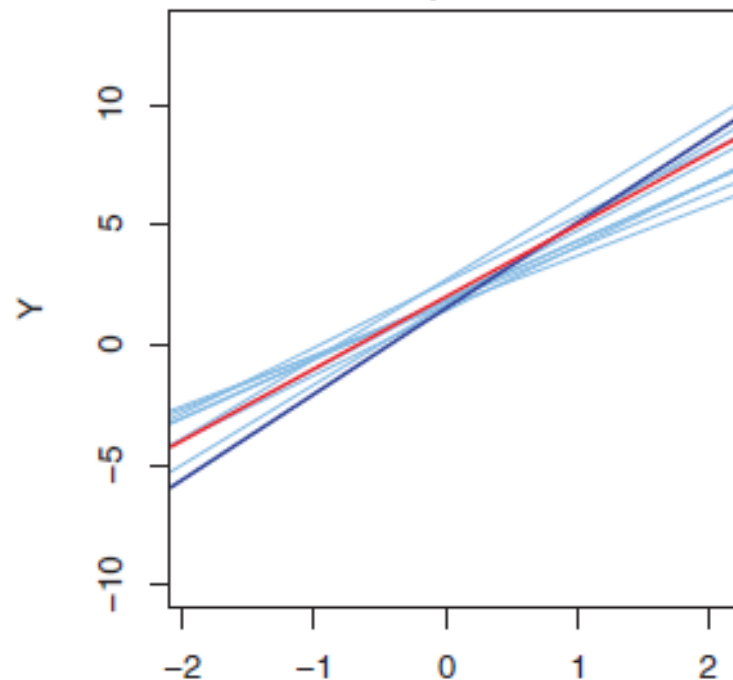


$\beta_0, \beta_1$

$\hat{\beta}_0, \hat{\beta}_1$

*A simulated data set.* Left: *The red line represents the true relationship, f(X) = 2+3X, which is known as the population regression line. The blue line is the least squares line; it is the least squares estimate for f(X) based on the observed data, shown in black*

# Linear Regression

$$\hat{\beta}_0, \hat{\beta}_1$$

$$\hat{\beta}_0, \hat{\beta}_1$$

$$\hat{\beta}_0, \hat{\beta}_1$$

$y_1$   $x_1$

$y_2$   $x_2$

$y_3$   $x_3$

$y_n$   $x_n$

freq.

$\hat{\beta}_0$

freq.

$\hat{\beta}_1$



*The population regression line is again shown in red, and the least squares line in dark blue. In light blue, ten least squares lines are shown, each computed on the basis of a <u>separate random set of observations</u>. Each least squares line is different, but on average, the least squares lines are quite close to the population regression line.*

# Concepts of Population and Sample

- Mean

- Variance

- Covariance


- Population and Sample

- Population mean and variance

- Sample mean and variance

# Mean and Variance

Expected value of $g(x)$

$$E(g(x)) = \int_{-\infty}^{\infty} g(x) f(x) dx$$

$X$ is a random variable with P.d.f. $f(x)$

Mean value of $X$ is $\mu \equiv \int_{-\infty}^{\infty} x f(x) dx = E(x)$

Variance of $X$ is $\sigma^2 = \int_{-\infty}^{\infty} (x-\mu)^2 f(x) dx$

$$= E\left[(X - E(x))^2\right) = E\left[x^2 - 2x E(x) + \left[E(x)\right]^2\right]$$

$$= E(x^2) - 2E(x)E(x) + \left[E(x)\right]^2 = E(x^2) - \left[E(x)\right]^2$$

$$\Rightarrow Var(X) = E(X^2) - \left[E(x)\right]^2$$

# Variance and Covariance

$$Cov(X,Y) = E\left([X-E(X)][Y-E(Y)]\right)$$

$$= E(XY) - 2E(X)E(Y) + E(X)E(Y)$$

$$= E(XY) - E(X)E(Y)$$

$$\therefore Var(X) = E(X^2) - [E(X)]^2$$

$$\text{and } Cov(X,Y) = E(X\cdot Y) - E(X)E(Y)$$

# Variance and Covariance

1.) $\text{Var}(X) \geqslant 0$

2.) $P(X=a)=1 \iff \text{Var}(X)=0$

3.) $\text{Var}(X+a) = \text{Var}(X)$

4.) $\text{Var}(aX) = a^2 \text{Var}(X)$

5.) $\text{Var}(X) = \text{Cov}(X,X)$

6.) $\text{Var}(aX+bY) = a^2 \text{Var}(X) + b^2 \text{Var}(Y) + 2ab\,\text{Cov}(X,Y)$

# Variance of sum of random variables

$$\text{Var}\left[\sum_{i=1}^{N} X_i\right] = \sum_{i,j}^{N} \text{Cov}(X_i, X_j)$$

$$= \sum_{i=1}^{N} \text{Var}(X_i) + \sum_{i \neq j}^{N} \text{Cov}(X_i, X_j)$$

$$\text{Var}\left[\sum_{i=1}^{N} a_i X_i\right] = \sum_{i,j}^{N} a_i a_j \text{Cov}(X_i, X_j)$$

$$= \sum_{i=1}^{N} a_i^2 \text{Var}(X_i) + \sum_{i \neq j}^{N} a_i a_j \text{Cov}(X_i, X_j)$$

$$= \sum_{i=1}^{N} a_i^2 \text{Var}(X_i) + \sum_{1 \leq i < j \leq N} a_i a_j \text{Cov}(X_i, X_j)$$

# Variance and Covariance

$$\text{If } Cov(X_i, X_j) = 0 \quad \forall \; i \neq j$$

$$\Rightarrow \quad X_i, X_j \text{ are } \quad \text{uncorrelated}$$

For N independent r.v. $X_1, X_2 \ldots \ldots X_N$

$$Var\left[\sum_{i=1}^{N} X_i\right] = \sum_{i=1}^{N} Var(X_i)$$

If all the N r.v. have the same variance $\sigma^2$

then $\quad Var\left[\sum_{i=1}^{N} X_i\right] = \sum_{i=1}^{N} Var(X_i) = N\sigma^2$

# Variance of mean

Mean of $n$ r.V.s is $\frac{1}{n}\sum\limits_{i=1}^{n} X_i$

Variance of mean of $n$ r.v.s would be

$$\text{Var}\left(\frac{1}{n}\sum_{i=1}^{n} X_i\right) = \frac{1}{n^2}\sum_{i=1}^{n}\text{var}(X_i) \quad \left\}\begin{array}{l}\text{Assuming } X_i\text{'s}\\\text{are independent}\end{array}\right.$$

$$= \frac{1}{n^2}\left(n\sigma^2\right) = \frac{\sigma^2}{n} \quad \left\}\begin{array}{l}\text{Assuming } X_i\text{'s}\\\text{are i i d}\end{array}\right.$$

$$\therefore \text{Var}\left(\frac{1}{n}\sum_{i=1}^{n} X_i\right) = \frac{\sigma^2}{n} \quad \text{if } X_i\text{'s are i i d}$$

# Population and sample

Population : Population of singe $N$ with value $X_i$

Population mean : $\mu \equiv \dfrac{1}{N} \sum\limits_{i=1}^{N} X_i$ $\Big\}$ $\mu = E(X)$

Population variance : $\sigma^2 \equiv \dfrac{1}{N} \sum\limits_{i=1}^{N} (X_i - \mu)^2$ $\Big\}$ $\sigma^2 = E[X - E(X)]^2$

Sample : Take $n$ random values (with replacement) from the population. $y_1, y_2, \ldots \ldots y_n$

Sample mean : $\bar{y} = \dfrac{1}{n} \sum\limits_{i=1}^{n} y_i$

Sample variance : ?

# Sample mean

Expected value of Sample mean :

$$E(\bar{y}) = E\left[\frac{1}{n}\sum_{i=1}^{n} y_i\right] = \frac{1}{n}\sum_{i=1}^{n} E(y_i)$$

$$= \frac{1}{n}\sum_{i=1}^{n} \mu \qquad = \mu$$

$$\therefore \; E(\bar{y}) = \mu$$

Expected value of sample mean is equal to population mean.

$\Rightarrow$ Sample mean is an unbiased estimator of the population mean.

# Sample Variance

How to define sample variance so that it is an unbiased estimator of the population variance.

Let squared deviation be defined as

$$\sigma_y^2 = \frac{1}{n} \sum_{i=1}^{n} \left( y_i - \bar{y} \right)^2$$

then $E(\sigma_y^2) = E\left( \frac{1}{n} \sum_{i=1}^{n} \left( y_i - \bar{y} \right)^2 \right)$

$$= E\left( \frac{1}{n} \sum_{i=1}^{n} \left( y_i - \frac{1}{n} \sum_{j=1}^{n} y_j \right)^2 \right)$$

# Sample Variance

$$\Rightarrow E(\sigma_y^2) = E\left[\frac{1}{n} \sum_{i=1}^{n} \left(y_i - \frac{1}{n} \sum_{j=1}^{n} y_j\right)^2\right]$$

$$= \frac{1}{n} \sum_{i=1}^{n} E\left(y_i^2 - \frac{2}{n} y_i \sum_{j=1}^{n} y_j + \frac{1}{n^2} \sum_{j=1}^{n} y_j \sum_{k=1}^{n} y_k\right)$$

$$= \frac{1}{n} \sum_{i=1}^{n} \left(1 - \frac{2}{n}\right) E(y_i^2) - \frac{2}{n} \sum_{j \neq i} E(y_i y_j) + \frac{1}{n^2} \sum_{j=1}^{n} \sum_{k \neq j} E(y_i y_k) + \frac{1}{n^2} \sum_{j=1}^{n} E(y_j^2)$$

$$\Rightarrow E(\sigma_y^2) = \frac{1}{n} \sum_{i=1}^{n} \left[\left(1 - \frac{2}{n}\right)(\sigma^2 + \mu^2) - \frac{2}{n}(n-1)\mu^2 + \frac{1}{n^2} n(n-1)\mu^2 + \frac{n}{n^2}(\sigma^2 + \mu^2)\right]$$

# Sample Variance

$$\Rightarrow E\left(\sigma_g^2\right) = \frac{1}{n} \sum \left(\frac{n-1}{n}\right) \sigma^2$$

$$= \frac{1}{n} \, n \left[\frac{n-1}{n}\right] \sigma^2 = \left(\frac{n-1}{n}\right) \sigma^2$$

$$\therefore E\left(\sigma_g^2\right) = \frac{n-1}{n} \, \sigma^2$$

$\therefore$ $\sigma_g^2$ estimates the population variance with a

bias of $\frac{n-1}{n}$ factor

# Sample Variance

∴ if we define sample variance $S^2 = \frac{n}{n-1} \sigma_y^2$

then it will be an embiased estimator of the population variance.

∴ Sample variance $S^2$ is defined as

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