# Linear Regression-1

Prof. Asim Tewari
IIT Bombay

# Supervised Machine Learning

- Our goal in supervised machine learning is to extract a relationship from data (ordered pairs of *(y,x)* )

The real relation is

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

$\epsilon$ is noise with zero mean.

What we get from learning from data is

$$\hat{y} = h(x)$$

# Regression vs Classification

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

- The task of <u>classification</u> differs from <u>regression</u> in that we assign a <u>discrete number</u> of classes (nominal scale or ordinal scale), instead of assigning it a <u>continuous value</u> (interval or ratio scale).

- If $y$ is in interval or ratio scale, then it is regression
- If $y$ is in Nominal or ordinal (?) scale, then it is classification

# 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     $\omega \in R^d$

$\uparrow$

Parameters or the weights

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

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

# Regression

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

$$\downarrow \quad \downarrow$$
$$\omega_0, \omega_1$$

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

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

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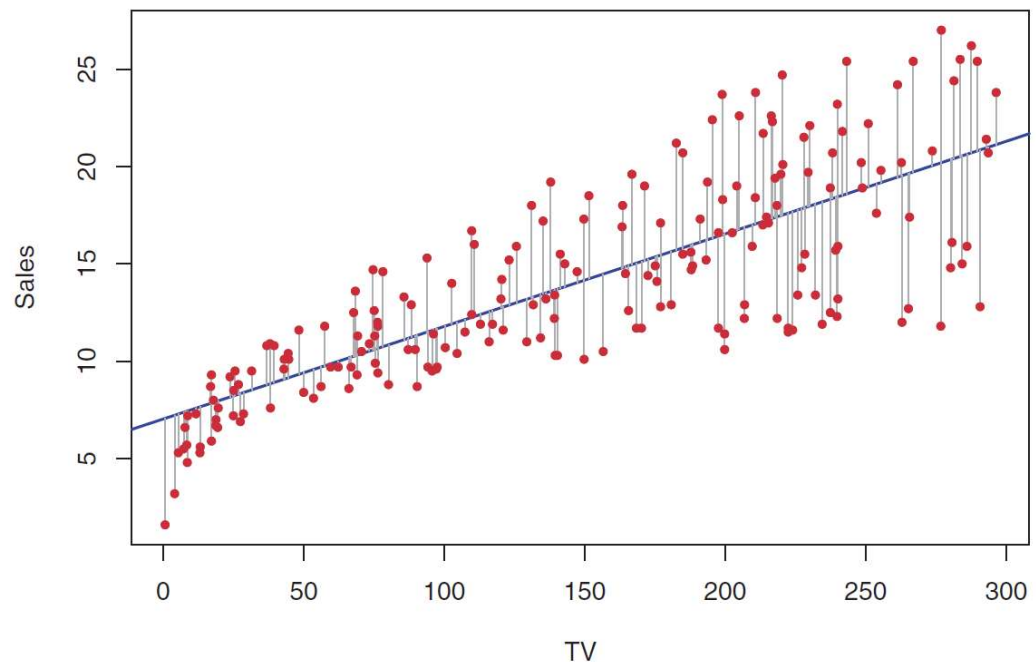
# Linear Regression

$$y = f(x) + \epsilon \quad \underset{\text{Error/Noise}}{\uparrow}$$

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

- 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

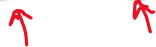$$\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x$$

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

error (zero mean)

$$x_i = \begin{bmatrix} x_i^1 \\ x_i^2 \\ \vdots \\ x_i^p \end{bmatrix}, \quad x_i \in R^p$$

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

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

# 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

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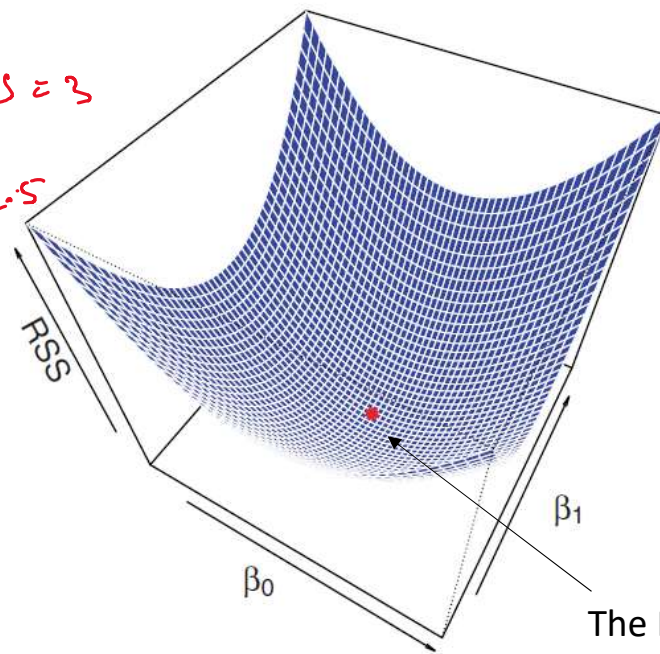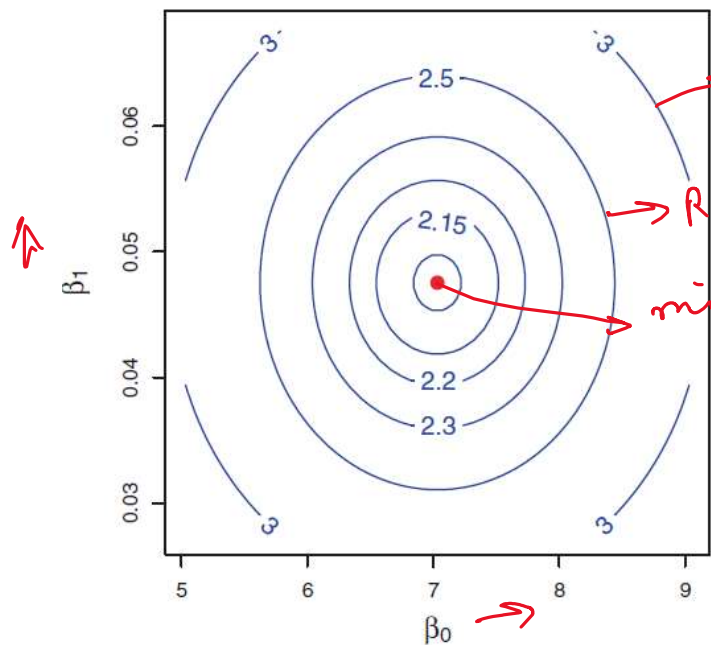
By minimizing RSS we can find

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

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

$$e_i = y_i - \hat{y}_i$$
$$\hookrightarrow \hat{\beta}_0 + \hat{\beta}_1 x_i$$

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

# Linear Regression

$$RSS = f\left(\hat{\beta}_0, \hat{\beta}_1\right) ;$$



RSS = 3

RSS = 2.5

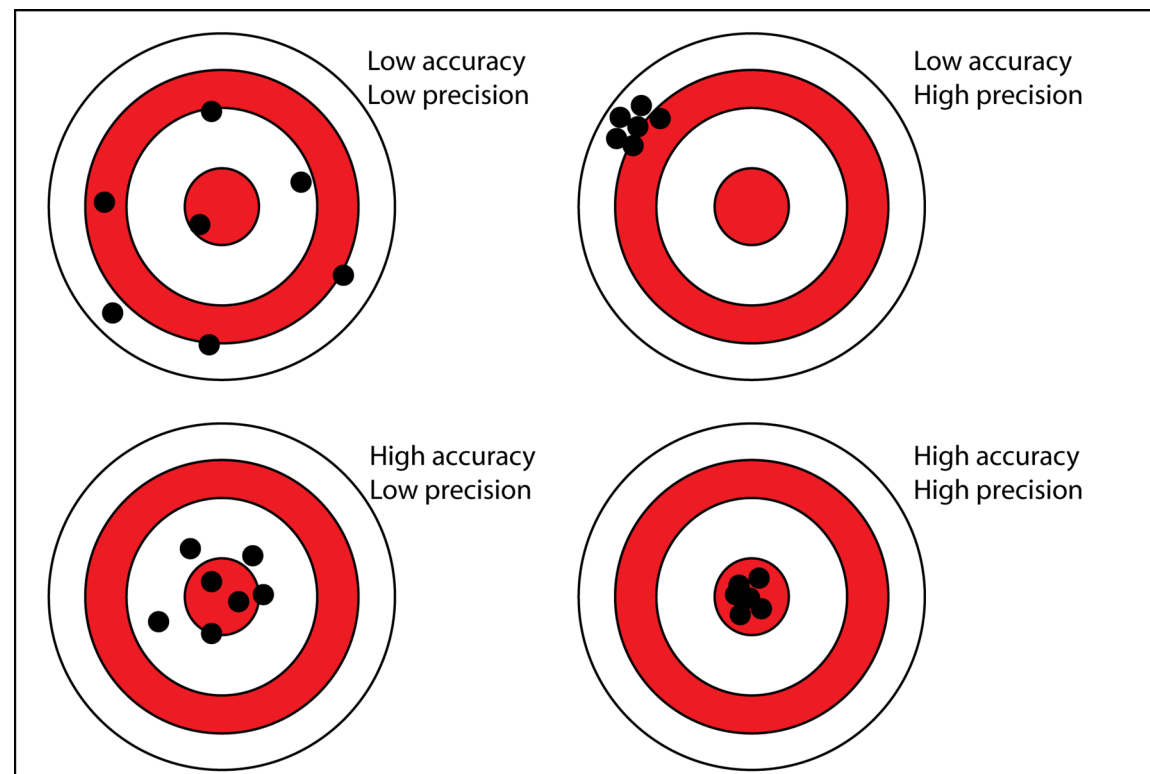min RSS

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

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

$$\hat{Y} = \hat{\beta_0} + \hat{\beta_1} X$$

- True relation between X and Y

  Y=f(x) + €

  $\uparrow$ Mean zero

random error Independent of x

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

1) How close is $\hat{\beta_0}$ to $\beta_0$ and $\hat{\beta_1}$ to $\beta_1$ ?

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

$\hat{\beta_0}$ 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 $\beta_1$

# 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

$$y_1 \quad x_1$$
$$y_2 \quad x_2$$
$$y_3 \quad x_3$$

$$\vdots \quad \vdots$$

$$y_n \quad x_n$$
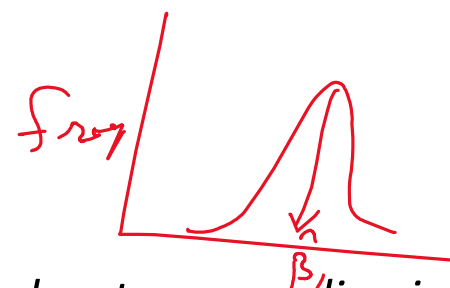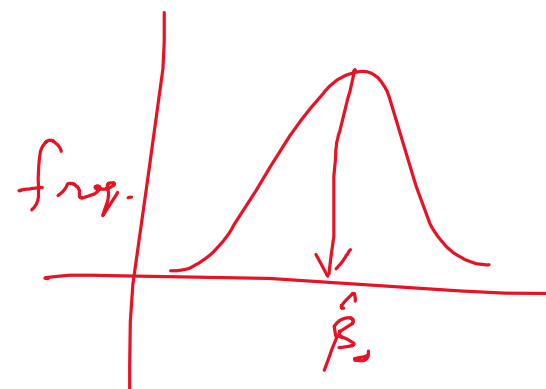
$\hat{\beta}_0, \hat{\beta}_1$

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

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

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

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

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

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

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

Expected value of $g(x)$

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

# 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.) $\text{Var}(X) \geq 0$

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

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

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

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

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

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

# Variance and Covariance

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

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

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

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

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

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

# Variance of mean

Mean of $n$ r.v.s is $\frac{1}{n}\sum_{i=1}^{n} x_i$

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

$$Var\left(\frac{1}{n}\sum_{i=1}^{n} x_i\right) = \frac{1}{n^2}\sum_{i=1}^{n} var(x_i) \quad \Bigg\} \text{ Assuming } x_i\text{'s are independent}$$

$$= \frac{1}{n^2}\left(n\,\sigma^2\right) = \frac{\sigma^2}{n} \quad \Bigg\} \text{ Assuming } x_i\text{'s are iid}$$

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

# Population and sample

Population : Population of size $N$ 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. $y_1, y_2, \cdots y_n$

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

Sample variance : ?

# Sample mean

Expected value of Sample mean:

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

$$= \frac{1}{n}\sum_{i=1}^{n} \mu = \mu \qquad \therefore E(\bar{y}) = \mu$$

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

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

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

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

# Sample Variance

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

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

$j=i$ $\qquad$ $j\neq i$ $\qquad\qquad$ $k\neq j$ $\qquad$ $k=j$

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

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

# 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

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

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

$\therefore$ Sample variance $S^2$ is defined as

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