



భారతీయ సాంకేతిక విజ్ఞాన సంస్థ హైదరాబాద్
भारतीय प्रौद्योगिकी संस्थान हैदराबाद
Indian Institute of Technology Hyderabad

Biostatistics BT2023

Lecture 9 + 10 Correlation and regression

Himanshu Joshi
5 and 8 September 2023

Measure of dispersion

Chebyshev's inequality

The rule is often known as Chebyshev's theorem, tells about the range of standard deviations around the mean, in statistics. In a probability distribution, no more than a certain fraction of values can be more than a certain distance from the mean.

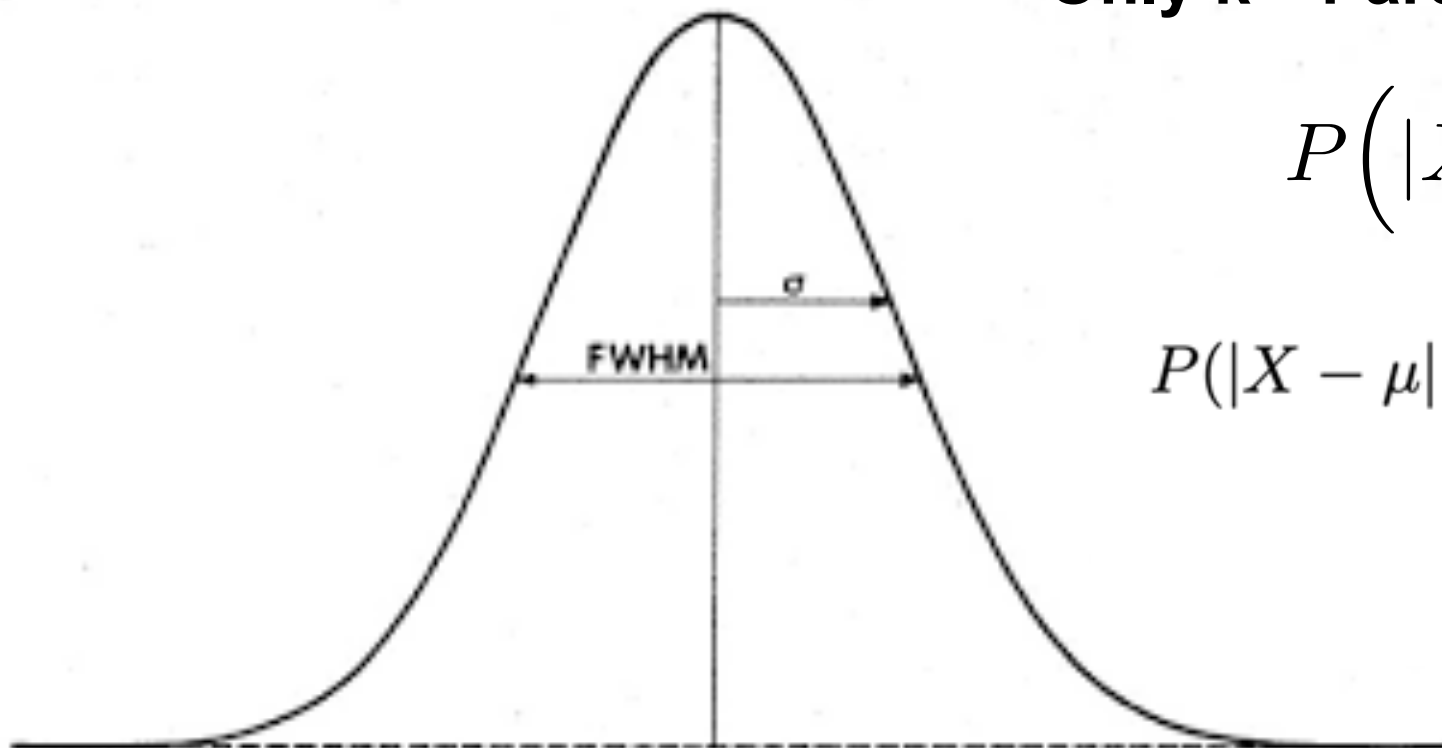
$$P(r) \left(|X - \mu| \geq k \times \sigma \right) \leq \frac{1}{k^2}$$

Only $k > 1$ are interesting because for $k < 1$ it trivial

$$P \left(|X - \mu| < k \times \sigma \right) \geq 1 - \frac{1}{k^2}$$

$$P(|X - \mu| < k\sigma) = P(\mu - k\sigma < X < \mu + k\sigma) \geq 1 - \frac{1}{k^2}$$

$$P(|X - \mu| \geq r) \leq \frac{\text{Var}(X)}{r^2}.$$



Moments of central measure

Skewness

$$\beta_1 = \frac{\mu_3^2}{\mu_2^3}$$

First moment : μ_1 $\frac{\sum (X - \bar{X})}{N}$ $\frac{\sum f(X - \bar{X})}{N}$
Always 0

Second moment : μ_2 $\frac{\sum (X - \bar{X})^2}{N}$ $\frac{\sum f(X - \bar{X})^2}{N}$

Measure of variance

Third moment : μ_3 $\frac{\sum (X - \bar{X})^3}{N}$ $\frac{\sum f(X - \bar{X})^3}{N}$
Measure skewness

Kurtosis

$$\beta_2 = \frac{\mu_4}{\mu_2^2}$$

Forth moment : μ_4 $\frac{\sum (X - \bar{X})^4}{N}$ $\frac{\sum f(X - \bar{X})^4}{N}$
Measure Kurtosis

$\beta > 3$ Leptokurtic

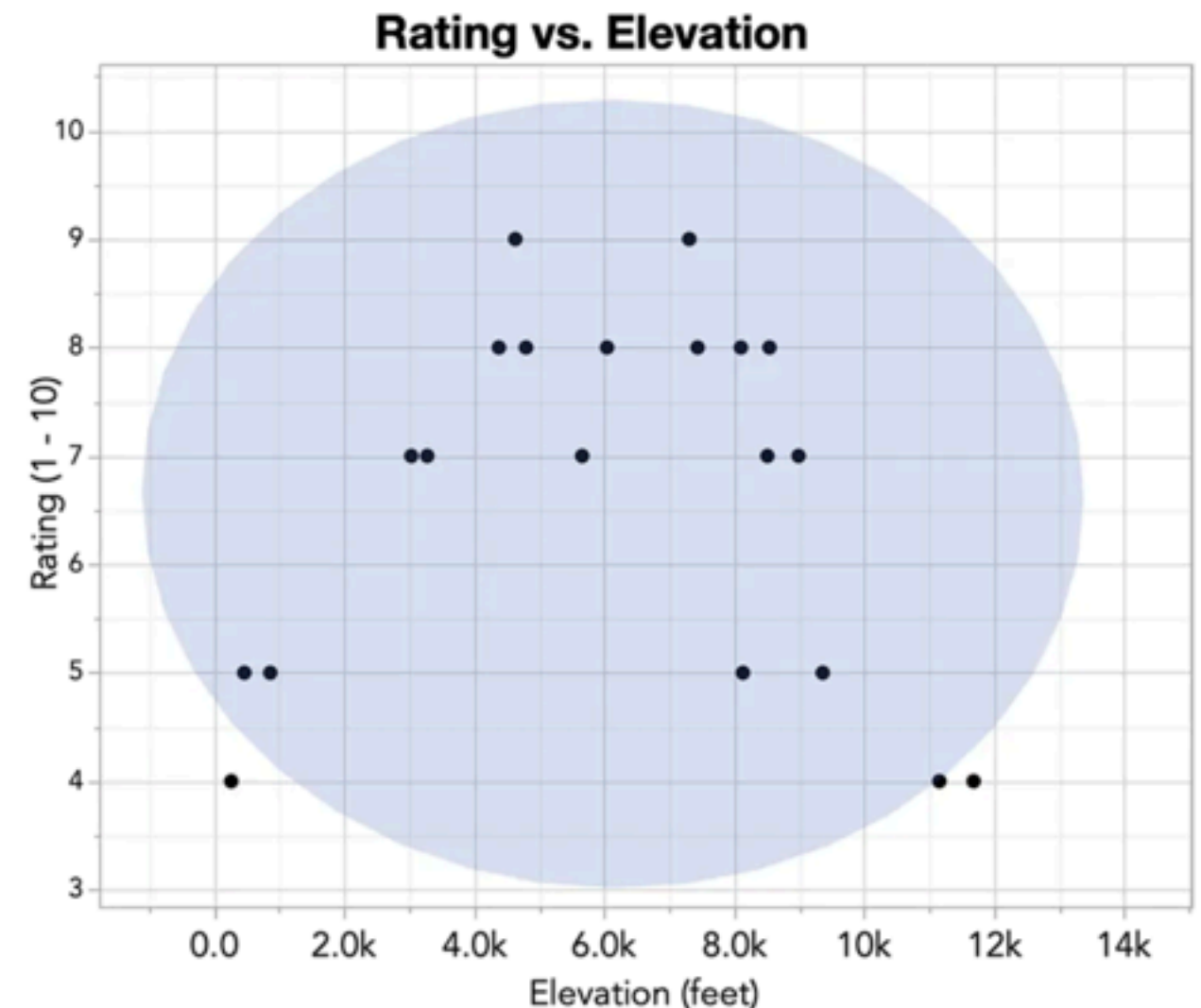
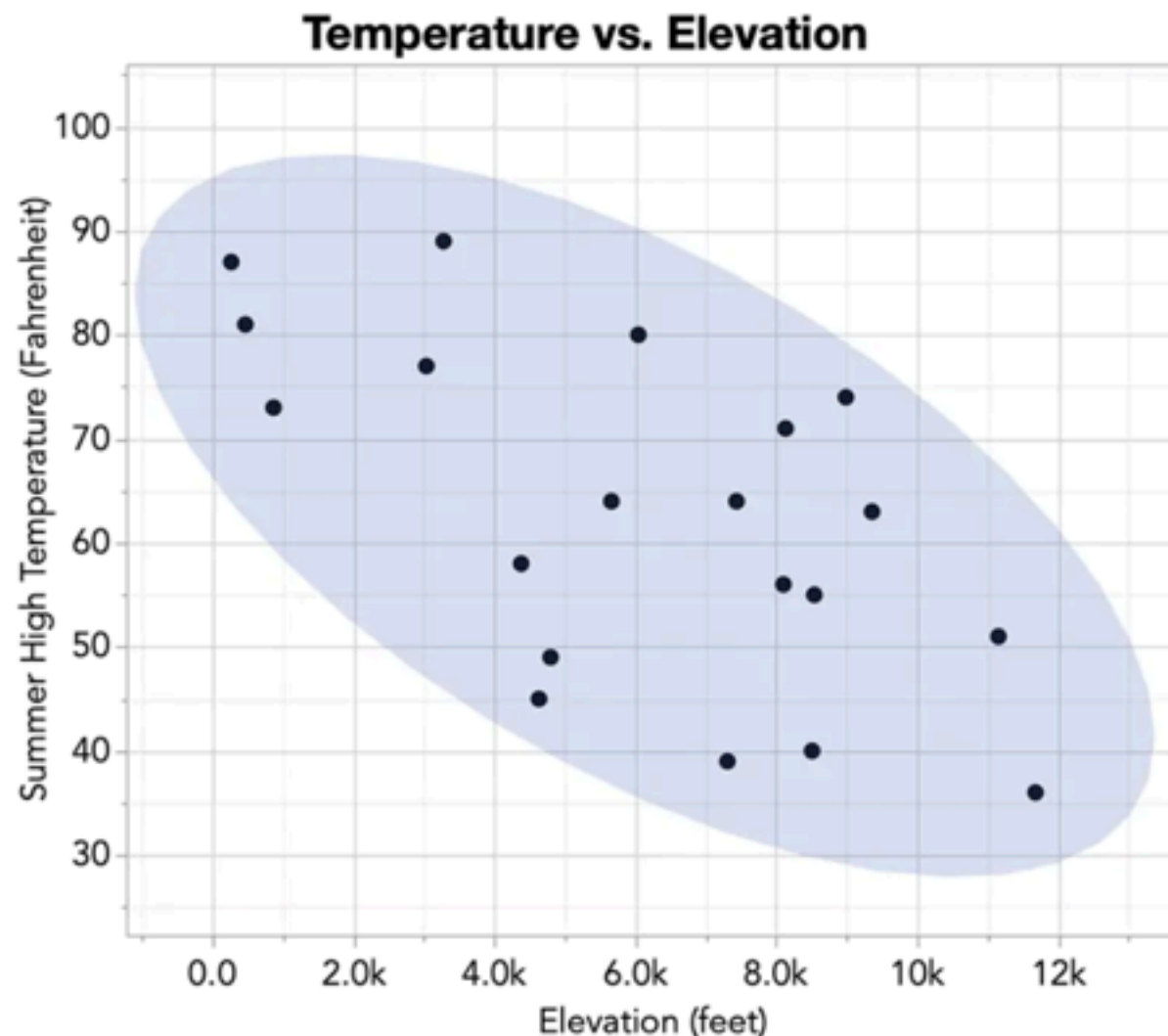
$\beta < 3$ Platykurtic

$\beta = 3$ Mesokurtic

Causality and Correlation

Causality is an influence which leads to another event or production

We are often interested in if two incidents are related to each other, if so how to quantify them. Lets consider the following scatter plot





Covariance

Covariance in two quantities variable X and Y on a given set is given by

$$\begin{aligned} Cov. (X, Y) &= \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}) \\ &= \frac{1}{n} \sum_{i=1}^n x_i y_i - \frac{1}{n} \sum_{i=1}^n (x_i) \frac{1}{n} \sum_{i=1}^n (y_i) \end{aligned}$$

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



Correlation coefficient

$$r = \frac{\sum_{i=1}^n \frac{(x_i - \bar{x})}{\sigma_x} \frac{(y_i - \bar{y})}{\sigma_y}}{n}$$
$$\implies r = \frac{\sigma_{xy}}{\sigma_x \sigma_y}$$

where $\sigma_{xy} = Cov(X, Y)$

Properties of the correlation coefficient

1. It will range from -1 to +1
2. Measures the closeness of the fit

$$dx = x - \bar{x}$$

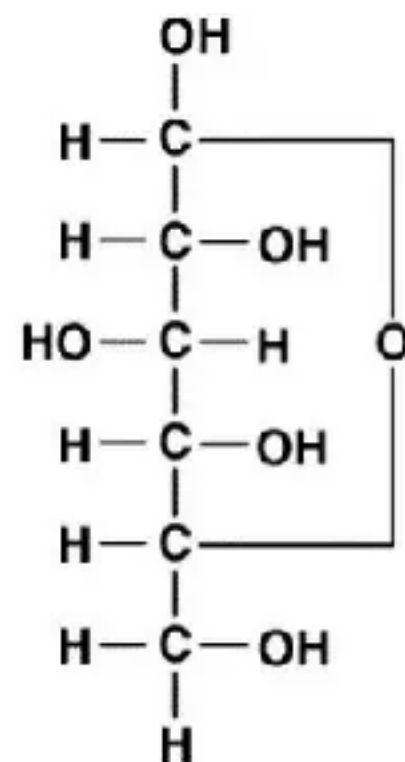
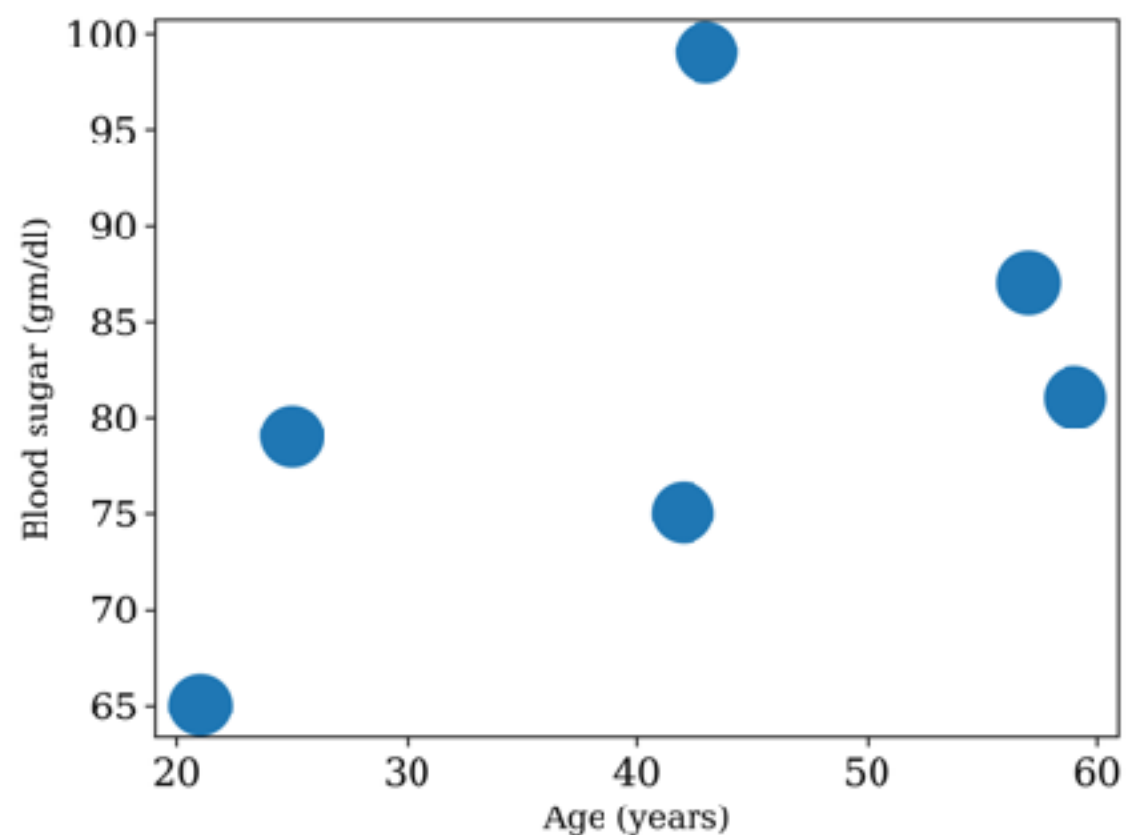
$$dy = y - \bar{y}$$

$$r_{xy} = \frac{\sum_{i=1}^n dx dy - \left(\sum dx \sum dy \right)}{\sqrt{\sum dx^2 - \frac{(\sum dx)^2}{n}}} \times \sqrt{\sum dy^2 - \frac{(\sum dy)^2}{n}}$$

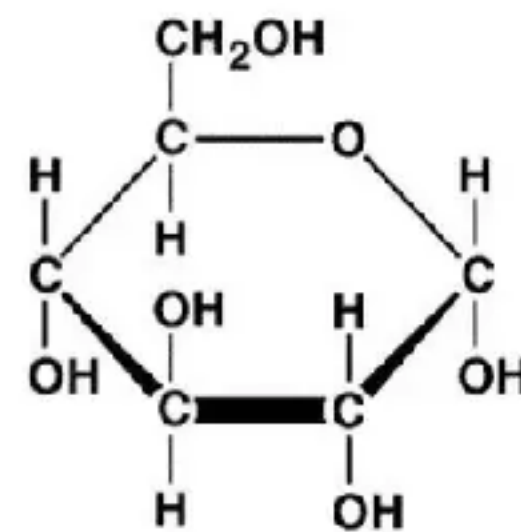
Blood sugar

Units

4.4 to 6.1 mm/liter
 or 82 to 110 mg/deciliter



Glucose



Four types of biomolecules

?

Regression

Regression shows the relationship b/w the average values of two variable. Its is very helpful in finding correlations.

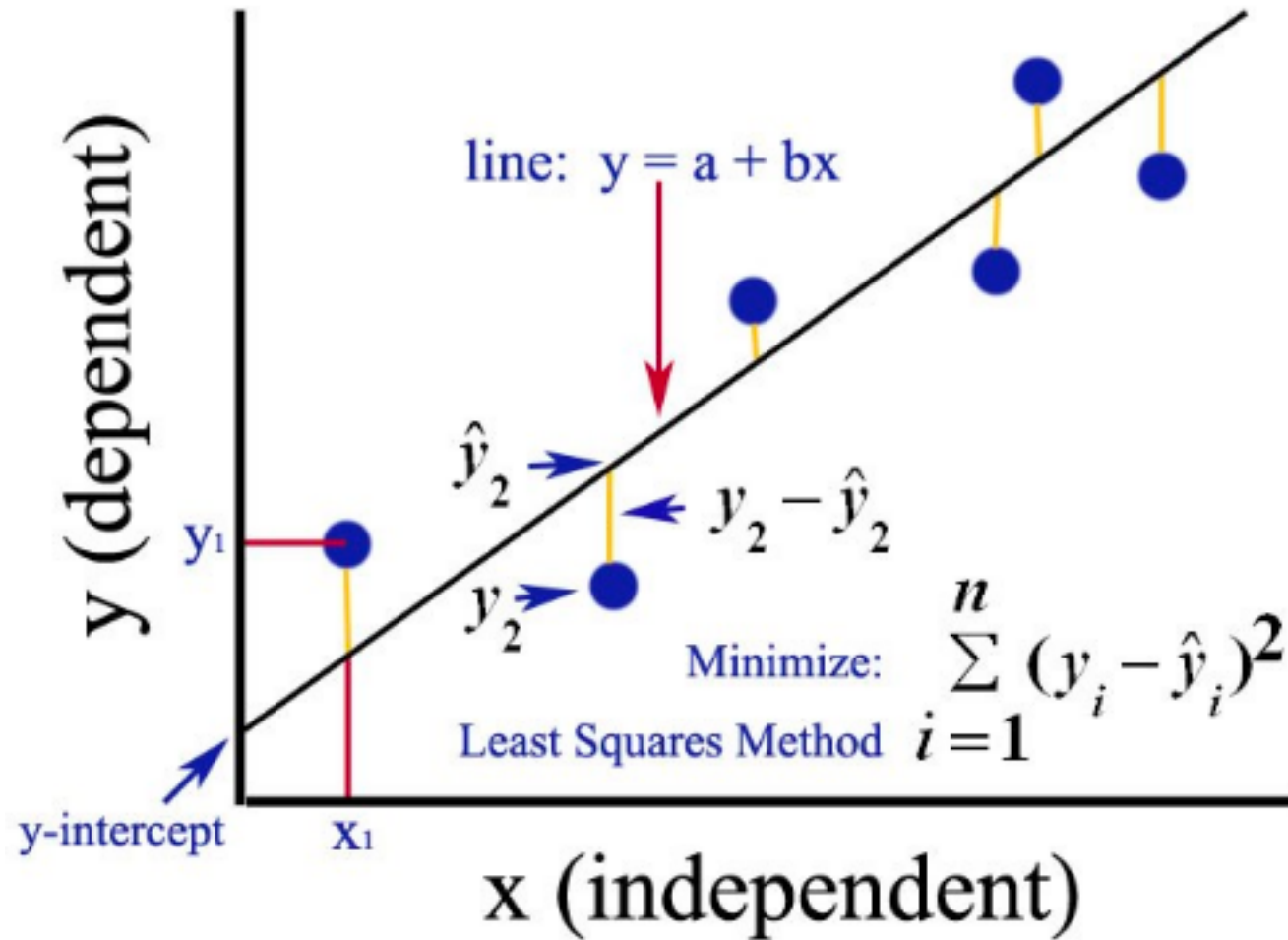
1. **Linear regression** can be solved exactly
2. **Non-linear regression** solved using approximation or iteration

Linear regression by least square method

It is a method is the process of finding the best-fitting curve or line of best fit for a set of data points by reducing the sum



Linear regression by least square method



$$\sum y = na + b \sum x$$
$$\sum xy = a \sum x + b \sum x^2$$

$$b_{yx} = \frac{n \sum xy - \sum x \sum y}{n \sum x^2 - (\sum x)^2}$$



భారతీయ సాంకేతిక విజ్ఞాన సంస్థ హైదరాబాద్
भारतीय प्रौद्योगिकी संस्थान हैदराबाद
Indian Institute of Technology Hyderabad

Covariance Matrix

Principal component analysis



భారతీయ సాంకేతిక విజ్ఞాన సంస్థ హైదరాబాద్
भारतीय प्रौद्योगिकी संस्थान हैदराबाद
Indian Institute of Technology Hyderabad

Next Class

2:30 PM Friday, 12 September 2023