



TEXAS A&M UNIVERSITY  
Engineering

# **ECEN 758 Data Mining and Analysis: Lecture 2, Data and Attributes I**

---

Joshua Peeples, Ph.D.

Assistant Professor

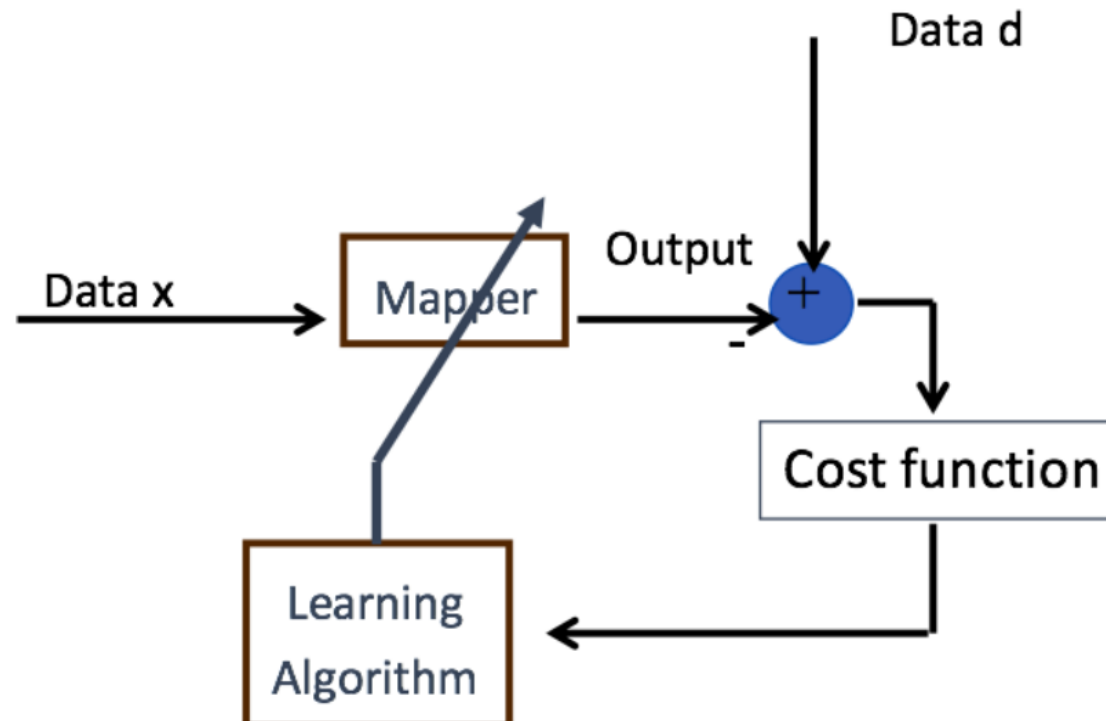
Department of Electrical and Computer Engineering

- Assignment #1 will be released next Wednesday (08/28)
  - Due Friday, 09/06
- Updated exam dates
  - 10/14 (Exam 1) and 11/25 (Exam 2)
- Tentative assignment dates

Assignment #	Released	Due
1	08/28	09/06
2	09/18	09/27
3	10/02	10/11
4	11/06	11/15

- Guest lecture dates
  - October 9<sup>th</sup>: Dr. John Gottula, Director of Crop Science, Agriculture, SAS
  - November 13<sup>th</sup>: Dr. Zigfried Hampel-Arias, Research Scientist, Remote Sensing and Data Science, Los Alamos National Laboratory
- Travel planned for November 18<sup>th</sup> – 21<sup>st</sup>
  - Lecture slides available for November 18<sup>th</sup>
  - No class November 20<sup>th</sup> (Review for exam and work on class project)
- Last day of class: December 2<sup>nd</sup> (No final exam 😊)

- Course objectives and Syllabus material
- Introduction to Data Mining and Machine Learning



- Data and attributes
  - **Numerical**
  - Categorical
- Reading: ZM Chapters 2



# What can we represent with data?

# Data Representations



- Numeric measurements, observations, settings, counts, time intervals, etc. (binary, integer, fixed-point, floating point)
- Text (characters, words, strings, documents)
- Signals (continuous numeric values)
- Time Series (sequence of discrete-time data points often from sensors, communication signals)
- Image and Video (pixel data, series of image data, voxel data, point-clouds)





# Data Types



# Data Types We Will Use

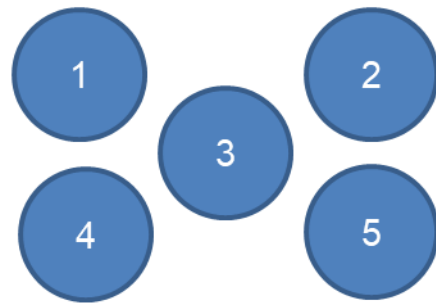


- Data used in Data Mining is generally of two types: Numeric Data and Categorical Data
- Numeric – quantitative, measurable; values are numbers. e.g. 0, 42, 3.1415,  $1.602 \times 10^{-19}$
- Categorical – qualitative, recognizable; values are restricted to the possible values in a category and can be represented by a text value or a number. e.g., Tuesday, Medium Rare, Hawaii

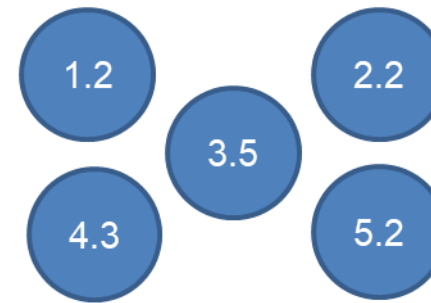
# Types of Numeric Data



- Discrete – variables can take on only specific values over an interval (e.g., counting numbers, integers)
- Continuous – variable can take on any value over an interval (e.g., real values)



Discrete



Continuous



# Univariate Analysis

- Focused on single attribute (e.g., feature)
- Data represented as matrix, **D**
- Each row is a sample and column is an attribute
- $X$  is a random variable
- Each  $x_i$  is independent and identically distributed (iid)

$$D = \begin{pmatrix} X \\ x_1 \\ x_2 \\ \vdots \\ x_n \end{pmatrix}$$

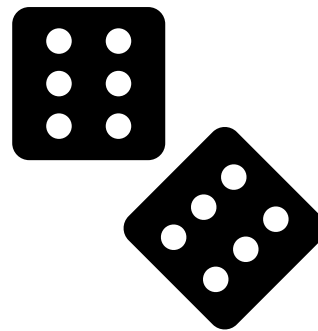


# Probability Review

# Sample Space and Probability



- Sample space( $\Omega$ ,  $S$ ) : Set of all possible outcomes of an experiment  
Example: Throwing a die,  $S = \{1, 2, 3, 4, 5, 6\}$



# Sample Space and Probability

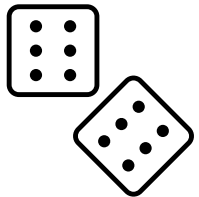


- Sample space( $\Omega$ ,  $S$ ) : Set of all possible outcomes of an experiment

Example: Throwing a die,  $S = \{1, 2, 3, 4, 5, 6\}$

- Probability law: Assigns a non-negative number to each element of the sample space

Example:  $P(1) = P(2) = \dots = P(6) = 1/6$



	1	2	3	4	5	6
Counts	100	100	100	100	100	100
Probability	1/6	1/6	1/6	1/6	1/6	1/6

- Sample space( $\Omega, S$ ) : Set of all possible outcomes of an experiment  
Example: Throwing a die,  $S = \{1, 2, 3, 4, 5, 6\}$
- Probability law: Assigns a non-negative number to each element of the sample space  
Example:  $P(1) = P(2) = \dots = P(6) = 1/6$
- Event ( $E$ ): A subset of sample space  
Example: Getting an even number,  $A = \{2, 4, 6\}$   
 $P(A) = P(2) + P(4) + P(6) = 1/2$



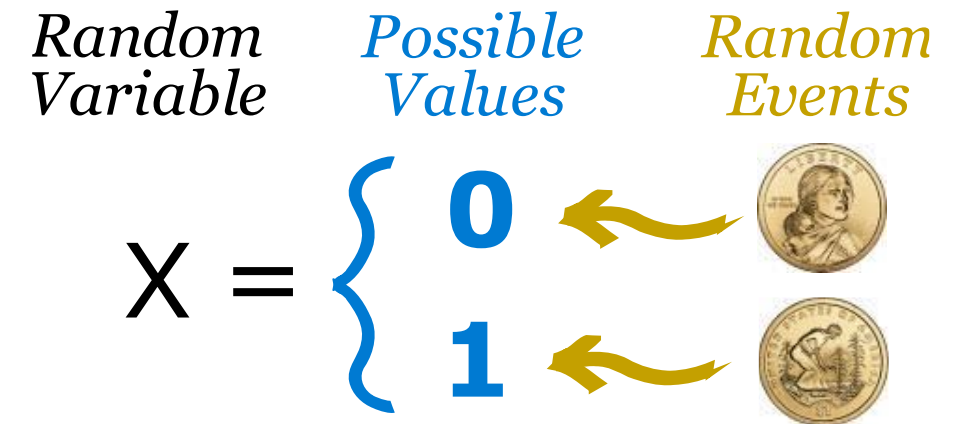


# Random Variables and Probability Distributions

# Random Variables



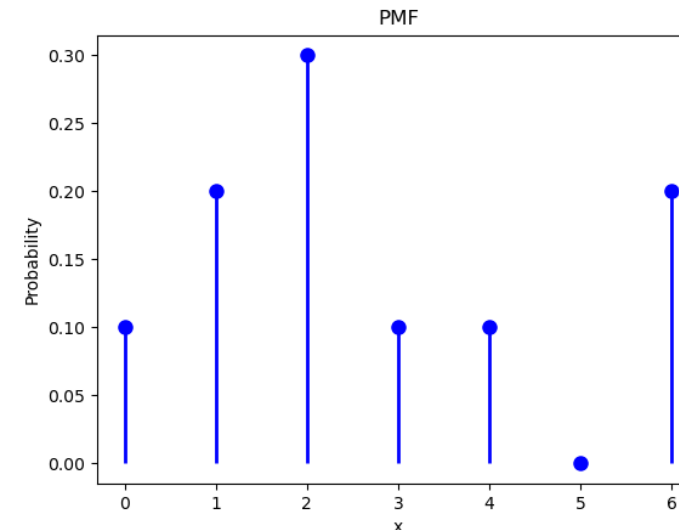
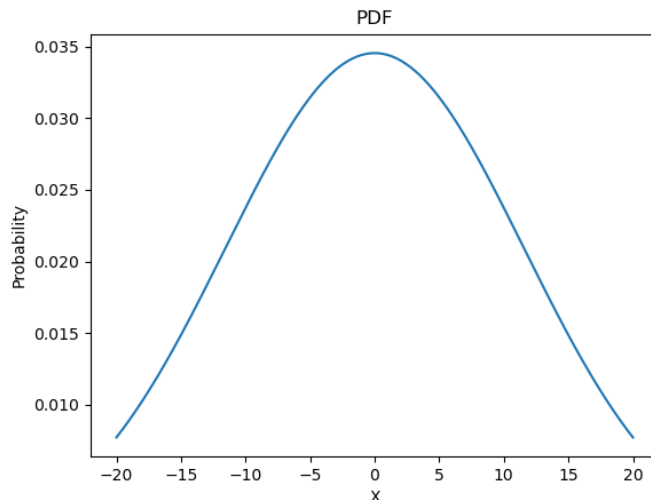
- Random variables: A mapping from the sample space to some range
  - Weather:  $W$  in {Sunny, Rainy, Foggy}
  - Temperature:  $T$  in {Hot, Cold};  $T$  in  $[-50, +50]$
- Random variables can be discrete or continuous
- Random variables can take finitely many values or infinitely many values



# Probability Density Function



- Probability density function (PDF) is an assignment of probability to each possible value of the continuous random variable (RV)
- Probability mass function (PMF) is used for discrete RVs



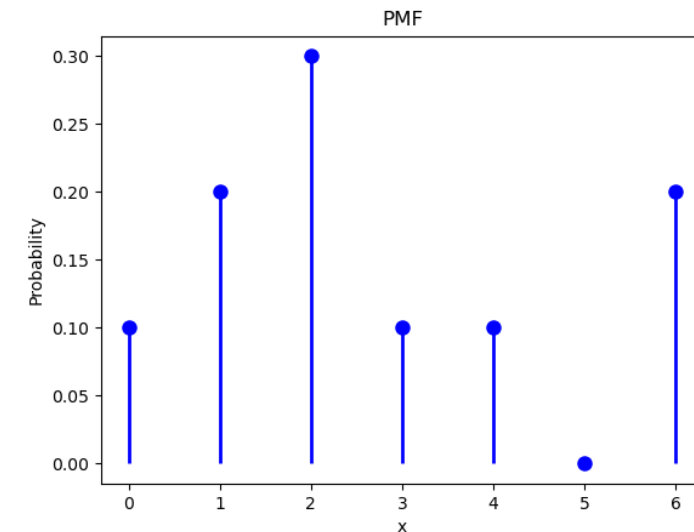
# Probability Mass Function



TEXAS A&M UNIVERSITY  
Engineering

- Probability mass function (PMFs) is used for discrete RVs
- Empirical PMFs assign equal probability to each point

$$\hat{f}(x) = P(X = x) = \frac{1}{n} \sum_{i=1}^n I(x_i = x)$$



- Empirical Cumulative Distribution Function (CDF) is the probability that data points ( $n$ ) in the sample are less than or equal to  $x$
- What is the relationship between CDFs and PDFs?

$$\hat{F}(x) = \frac{1}{n} \sum_{i=1}^n I(x_i \leq x)$$

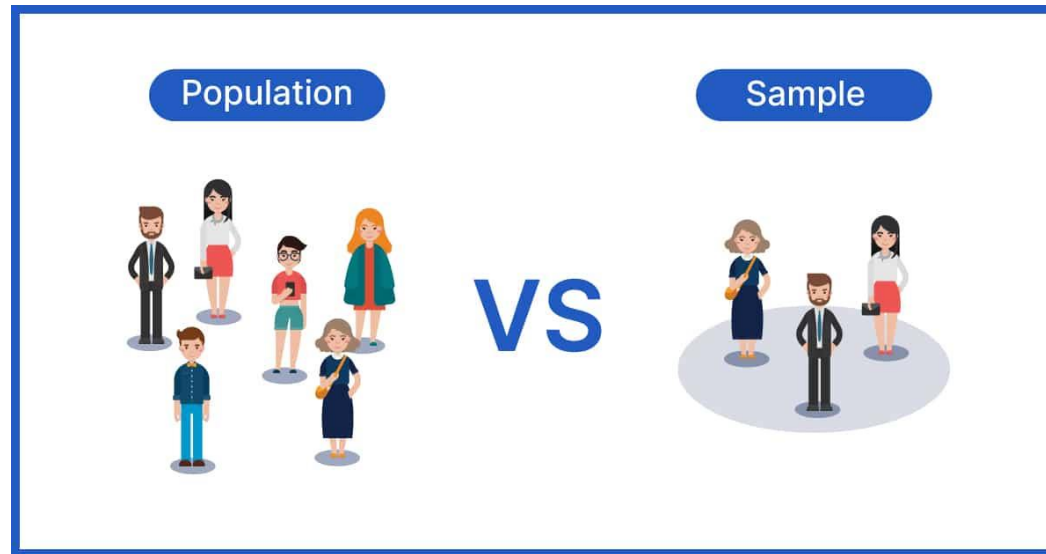


# Statistical Measures

# Populations vs Samples



- Population – set of all data in an area of interest
- Sample – subset of a population



- Measures of central tendency
  - Mean, Median, Mode
- Measures of dispersion
  - Range, Interquartile Range, Variance, Standard Deviation



- **Measures of central tendency**
  - Mean, Median, Mode
- **Measures of dispersion**
  - Range, Interquartile Range, Variance, Standard Deviation

- Arithmetic average of values of  $X$
- Also known as expected value
- $f(x)$  is PMF (discrete) or PDF (continuous)

$$\mu = E[X] = \sum_x x \cdot f(x)$$

Discrete

$$\mu = E[X] = \int_{-\infty}^{\infty} x \cdot f(x) dx$$

Continuous

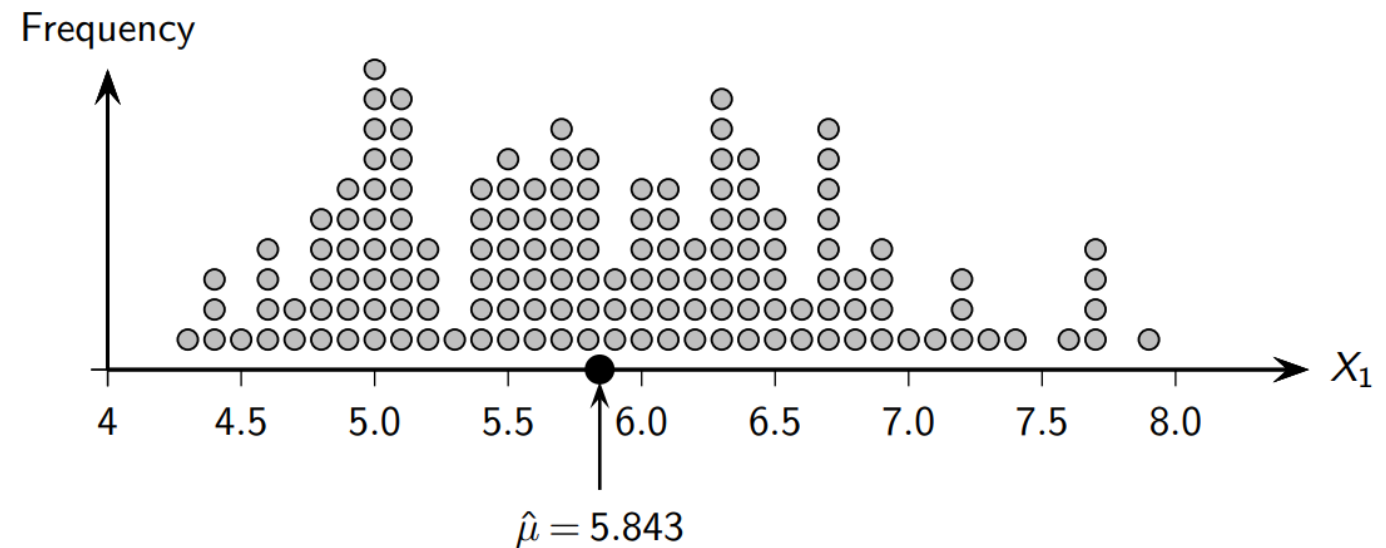
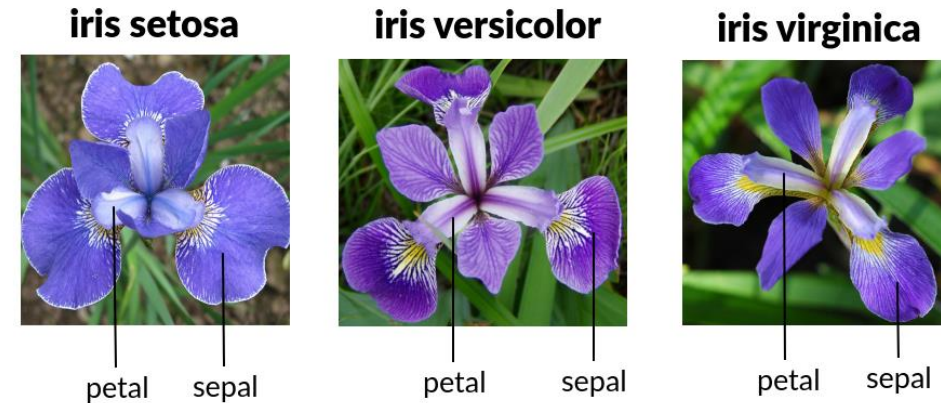
- Statistic defined as average value of  $x_i$
- Estimator for unknown mean value,  $\mu$ , of  $X$
- **Unbiased** estimator for population mean
- Not a **robust** statistic

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

$$E[\hat{\mu}] = E\left[\frac{1}{n} \sum_{i=1}^n x_i\right] = \frac{1}{n} \sum_{i=1}^n E[x_i] = \frac{1}{n} \sum_{i=1}^n \mu = \mu$$

# Sample Mean: Iris Sepal Length

- Iris dataset
  - 150 samples
  - 3 classes
  - 4 attributes
    - Petal and Sepal length and width



- “Middle-most” value
- Half of the values of  $X$  are less and half of the values more than median
- Can use CDF or inverse CDF to find median
- Robust statistic

$$P(X \leq m) \geq \frac{1}{2} \text{ and } P(X \geq m) \geq \frac{1}{2}$$

Median,  $m$

$$F(m) = 0.5 \text{ or } m = F^{-1}(0.5)$$

Median,  $m$  (CDF)

$$\hat{F}(m) = 0.5 \text{ or } m = \hat{F}^{-1}(0.5)$$

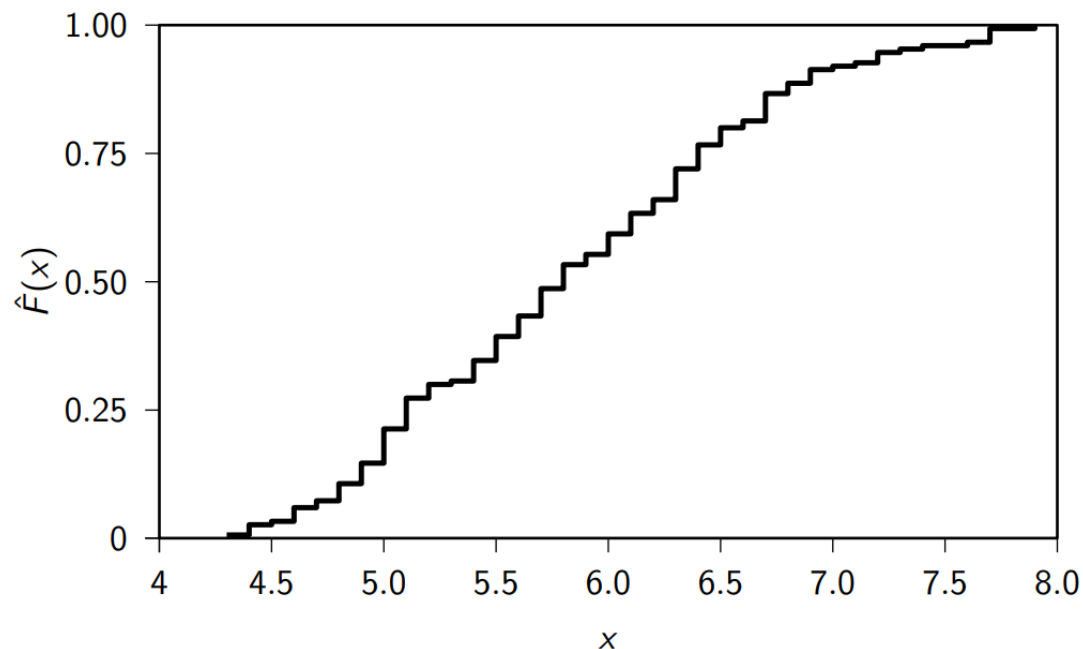
Sample Median,  $\hat{m}$  (CDF)

# Empirical CDF and Inverse CDF

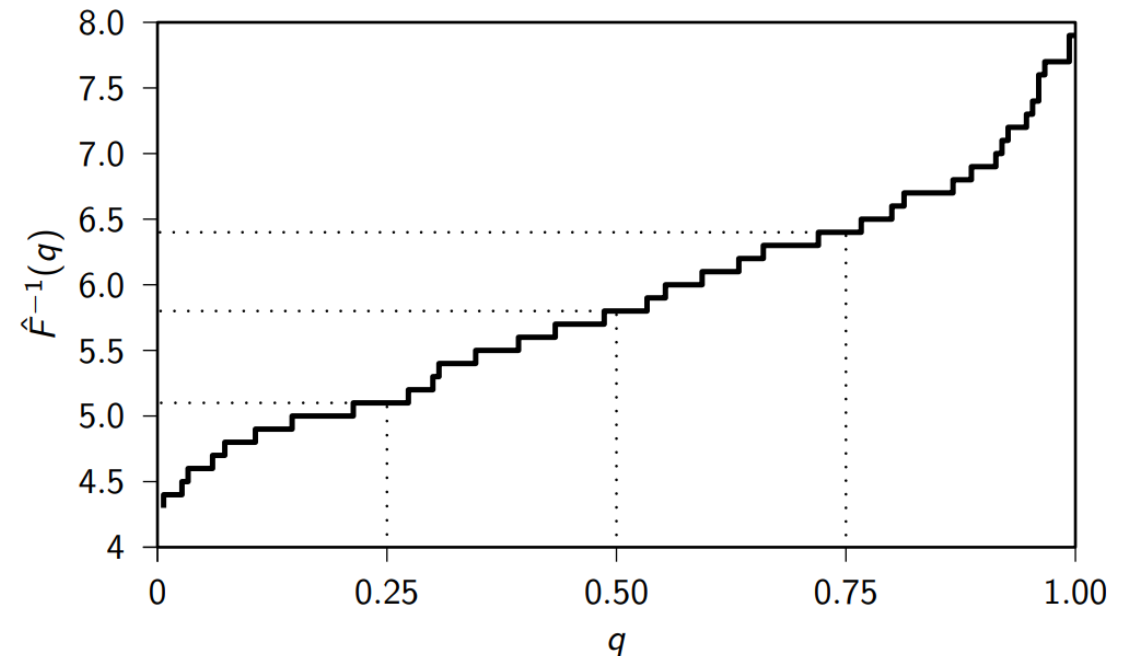


- Sepal length
- Median is 5.8

$$\hat{F}(5.8) = 0.5$$



$$5.8 = \hat{F}^{-1}(0.5)$$



- Value at which PMF or PDF attains maximum value
- Sample mode computed from empirical PMF

$$\text{mode}(X) = \arg \max_x \hat{f}(x)$$

- Measures of central tendency
  - Mean, Median, Mode
- **Measures of dispersion**
  - Range, Interquartile Range, Variance, Standard Deviation



- Difference between the maximum and minimum values of  $X$
- Robust statistic?

$$r = \max\{X\} - \min\{X\}$$

Range

$$\hat{r} = \max_{i=1}^n \{x_i\} - \min_{i=1}^n \{x_i\}$$

Sample range

# Interquartile Range (IQR)



- Difference between the 75<sup>th</sup> and 25<sup>th</sup> percentiles of  $X$
- More robust than range

$$IQR = F^{-1}(0.75) - F^{-1}(0.25)$$

IQR

$$\widehat{IQR} = \hat{F}^{-1}(0.75) - \hat{F}^{-1}(0.25)$$

Sample IQR

- Measure of how much values deviate from the expected value of  $X$

$$\sigma^2 = \text{var}(X) = E[(X - \mu)^2] = \begin{cases} \sum_x (x - \mu)^2 f(x) & \text{if } X \text{ is discrete} \\ \int_{-\infty}^{\infty} (x - \mu)^2 f(x) dx & \text{if } X \text{ is continuous} \end{cases}$$

- Standard deviation is positive square root of variance

$$\hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \hat{\mu})^2$$

Sample variance

$$\hat{\sigma} = \sqrt{\frac{1}{n} \sum_{i=1}^n (x_i - \hat{\mu})^2}$$

Sample standard deviation

# Variance of Sample Mean



- Sample mean is statistic
- Larger values of  $n$  will decrease variation of sample mean from mean

$$E[\hat{\mu}] = \mu$$

$$\text{var}(\hat{\mu}) = E[(\hat{\mu} - \mu)^2] = \frac{\sigma^2}{n}$$

# Variance of Sample Variance



- Sample variance is statistic
- Biased estimator of true population variance
- Larger values of  $n$  will reduce bias of estimator

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

$$E[\hat{\sigma}^2] \rightarrow \sigma^2 \quad \text{as } n \rightarrow \infty$$



# Bivariate Analysis

- Focused on two attributes (e.g., feature)
- Data represented as matrix, **D**
- Each row is a sample and column is an attribute
- **X** is a random variable
- Each  $\mathbf{x}_i$  is independent and identically distributed (iid)

$$\mathbf{D} = \begin{pmatrix} X_1 & X_2 \\ x_{11} & x_{12} \\ x_{21} & x_{22} \\ \vdots & \vdots \\ x_{n1} & x_{n2} \end{pmatrix}$$



- Expected value of the vector random variable ( $\mathbf{X}$ )

$$\boldsymbol{\mu} = E[\mathbf{X}] = E \left[ \begin{pmatrix} X_1 \\ X_2 \end{pmatrix} \right] = \begin{pmatrix} E[X_1] \\ E[X_2] \end{pmatrix} = \begin{pmatrix} \mu_1 \\ \mu_2 \end{pmatrix}$$

Bivariate mean

$$\hat{\boldsymbol{\mu}} = \sum_{\mathbf{x}} \mathbf{x} \hat{f}(\mathbf{x}) = \sum_{\mathbf{x}} \mathbf{x} \left( \frac{1}{n} \sum_{i=1}^n I(\mathbf{x}_i = \mathbf{x}) \right) = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i$$

Sample bivariate mean

- Measure of association or linear dependence between two attributes ( $X_1$  and  $X_2$ )
- Can use to check for independence

$$\begin{aligned}\sigma_{12} &= E[(X_1 - \mu_1)(X_2 - \mu_2)] \\ &= E[X_1 X_2] - E[X_1]E[X_2]\end{aligned}$$

Covariance

$$E[X_1 X_2] = E[X_1] \cdot E[X_2]$$

Independence

$$\hat{\sigma}_{12} = \frac{1}{n} \sum_{i=1}^n (x_{i1} - \hat{\mu}_1)(x_{i2} - \hat{\mu}_2)$$

Sample covariance

# Interpreting Covariance



$\text{cov}(X, Y) > 0 \rightarrow$  X and Y are positively correlated

$\text{cov}(X, Y) < 0 \rightarrow$  X and Y are inversely correlated

$\text{cov}(X, Y) = 0 \rightarrow$  X and Y are independent

- *Covariance values are not constrained, and can be from -infinity to +infinity*
- *Covariance is a measure of the directional relationship between variables*

- Standardized covariance between two attributes ( $X_1$  and  $X_2$ )
- Bounded between -1 (negatively correlated) and 1 (positively correlated)

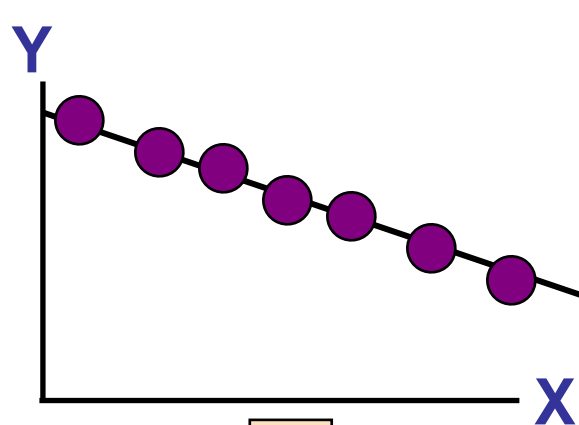
$$\rho_{12} = \frac{\sigma_{12}}{\sigma_1 \sigma_2} = \frac{\sigma_{12}}{\sqrt{\sigma_1^2 \sigma_2^2}}$$

Correlation

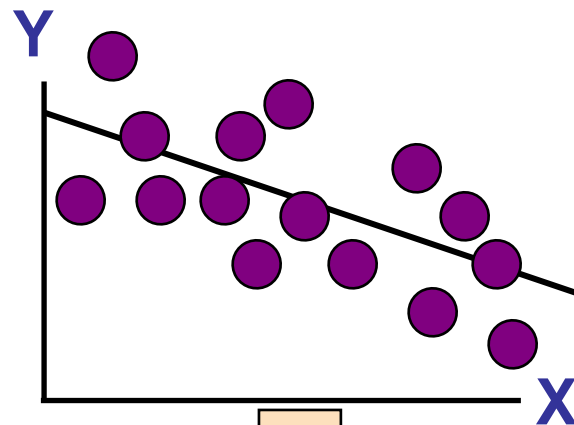
$$\hat{\rho}_{12} = \frac{\hat{\sigma}_{12}}{\hat{\sigma}_1 \hat{\sigma}_2} = \frac{\sum_{i=1}^n (x_{i1} - \hat{\mu}_1)(x_{i2} - \hat{\mu}_2)}{\sqrt{\sum_{i=1}^n (x_{i1} - \hat{\mu}_1)^2 \sum_{i=1}^n (x_{i2} - \hat{\mu}_2)^2}}$$

Sample correlation

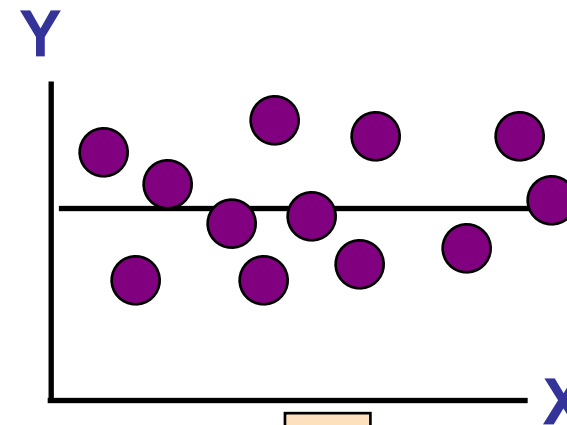
# Scatter Plots of Data with Various Correlation Coefficients



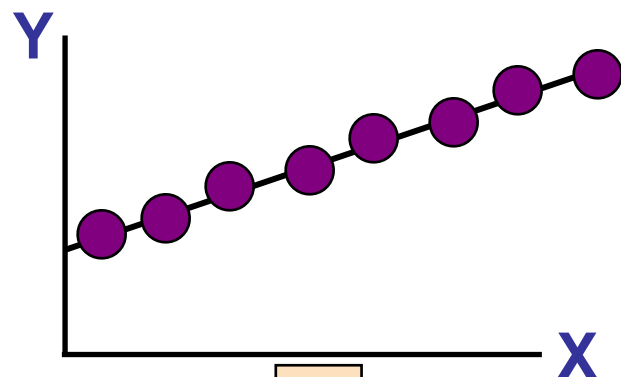
A



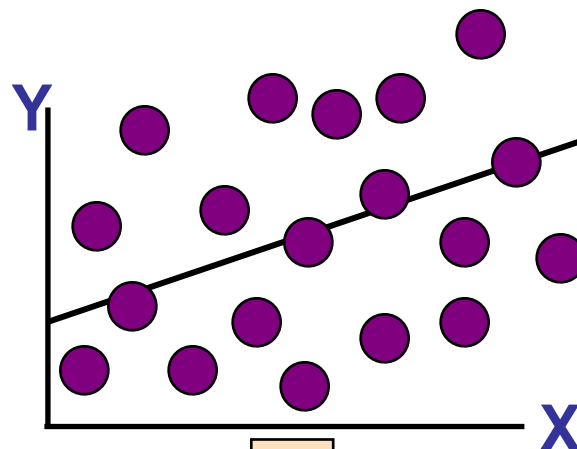
B



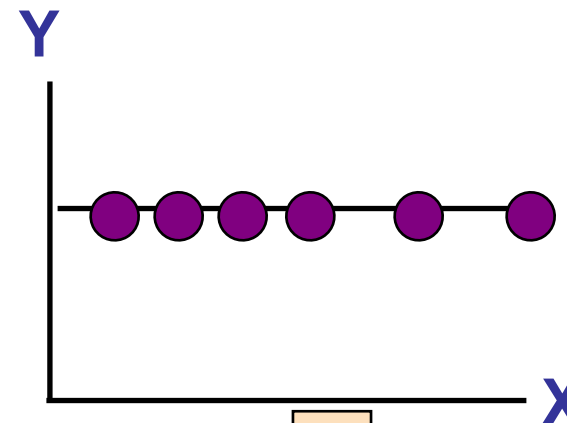
C



D

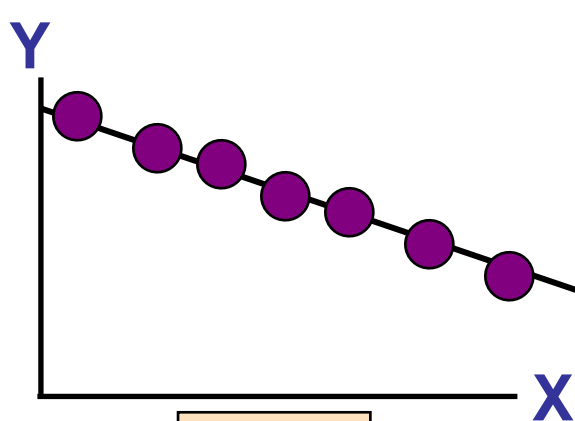


E

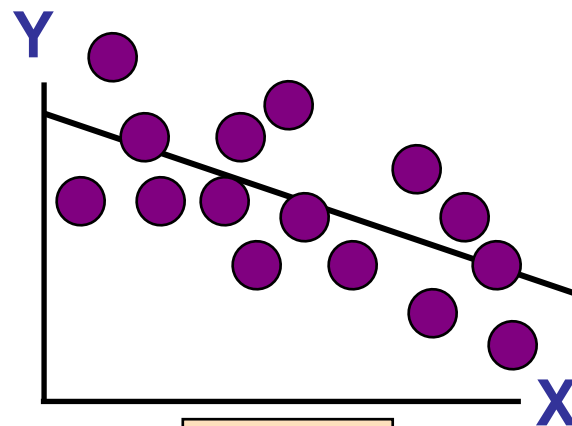


F

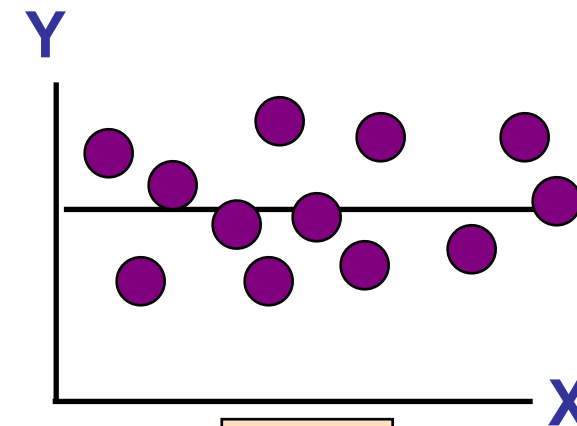
# Scatter Plots of Data with Various Correlation Coefficients



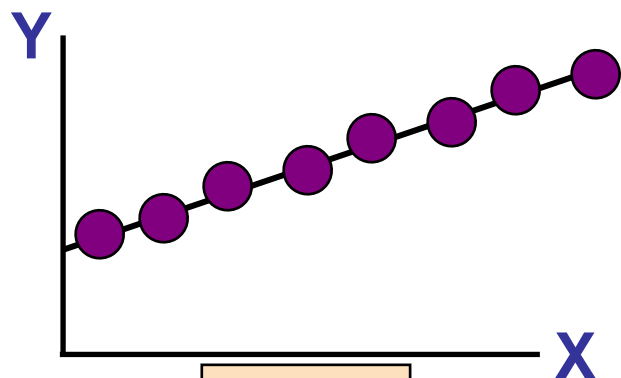
$$r = -1$$



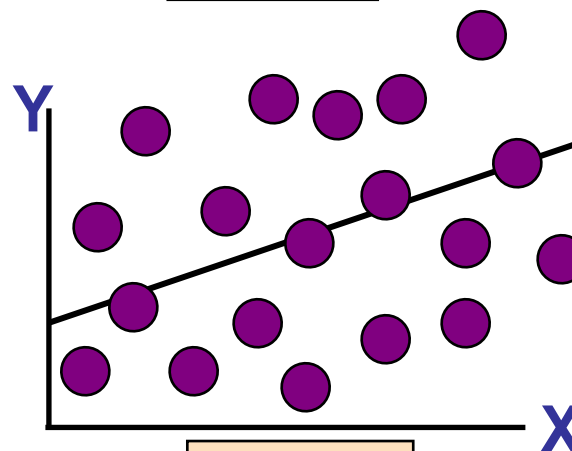
$$r = -.6$$



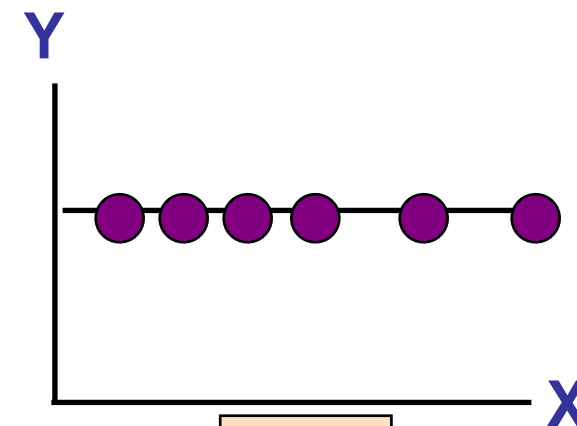
$$r = 0$$



$$r = +1$$



$$r = +.3$$



$$r = 0$$

- Cosine of the angle between two centered attribute vectors

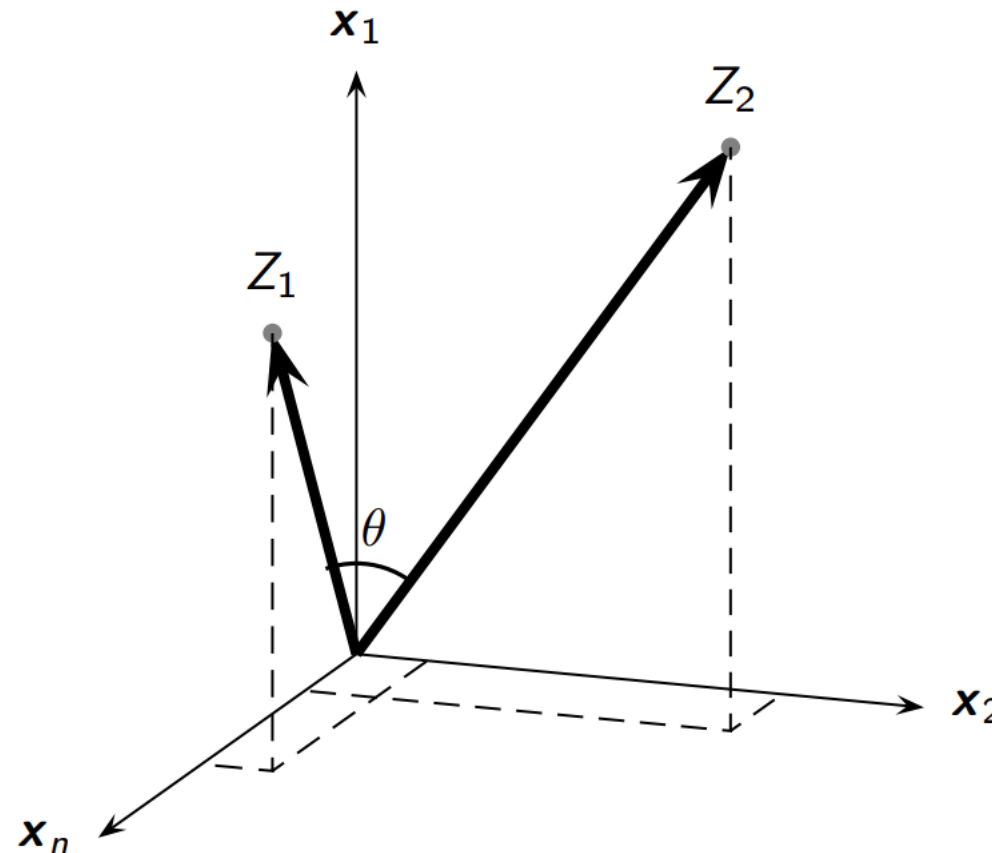
$$\bar{X}_1 = X_1 - 1 \cdot \hat{\mu}_1 = \begin{pmatrix} x_{11} - \hat{\mu}_1 \\ x_{21} - \hat{\mu}_1 \\ \vdots \\ x_{n1} - \hat{\mu}_1 \end{pmatrix} \quad \bar{X}_2 = X_2 - 1 \cdot \hat{\mu}_2 = \begin{pmatrix} x_{12} - \hat{\mu}_2 \\ x_{22} - \hat{\mu}_2 \\ \vdots \\ x_{n2} - \hat{\mu}_2 \end{pmatrix}$$

$$\hat{\rho}_{12} = \frac{\bar{X}_1^T \bar{X}_2}{\sqrt{\bar{X}_1^T \bar{X}_1} \sqrt{\bar{X}_2^T \bar{X}_2}} = \frac{\bar{X}_1^T \bar{X}_2}{\|\bar{X}_1\| \|\bar{X}_2\|} = \left( \frac{\bar{X}_1}{\|\bar{X}_1\|} \right)^T \left( \frac{\bar{X}_2}{\|\bar{X}_2\|} \right) = \cos \theta$$

# Correlation (Geometric Interpretation)



- Cosine of the angle between two centered attribute vectors





- Summary of variance-covariance information
- Symmetric matrix
- Total variance is trace of covariance matrix

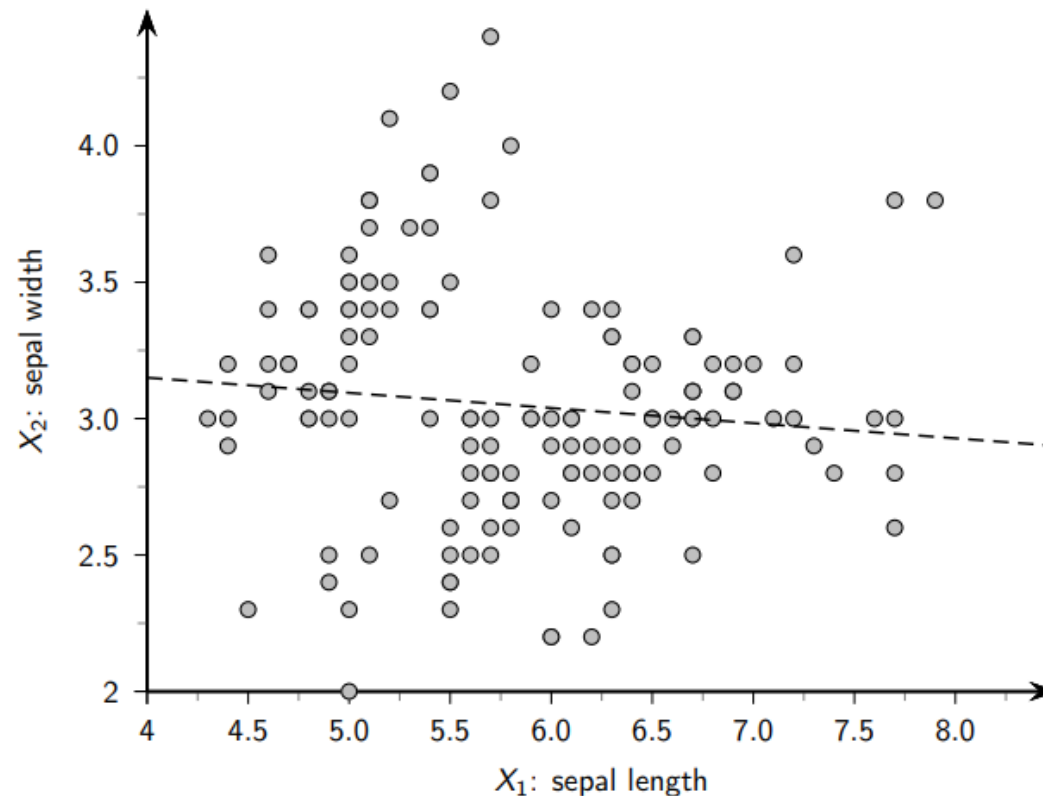
$$\begin{aligned}\Sigma &= E[(\mathbf{X} - \boldsymbol{\mu})(\mathbf{X} - \boldsymbol{\mu})^T] \\ &= \begin{pmatrix} \sigma_1^2 & \sigma_{12} \\ \sigma_{21} & \sigma_2^2 \end{pmatrix}\end{aligned}$$

$$\text{var}(\mathbf{D}) = \text{tr}(\Sigma) = \sigma_1^2 + \sigma_2^2$$

# Correlation and Covariance: Iris Dataset



TEXAS A&M UNIVERSITY  
Engineering



The sample mean is

$$\hat{\mu} = \begin{pmatrix} 5.843 \\ 3.054 \end{pmatrix}$$

The sample covariance matrix is

$$\hat{\Sigma} = \begin{pmatrix} 0.681 & -0.039 \\ -0.039 & 0.187 \end{pmatrix}$$

The sample correlation is

$$\hat{\rho}_{12} = \frac{-0.039}{\sqrt{0.681 \cdot 0.187}} = -0.109$$



# Multivariate Analysis

- Focused on  $d$  attributes (e.g., feature)
- Data represented as matrix,  $\mathbf{D}$
- Each row is a sample and column is an attribute
- $X$  is a random variable
- Each  $x_i$  is independent and identically distributed (iid)

$$\mathbf{D} = \begin{pmatrix} X_1 & X_2 & \cdots & X_d \\ x_{11} & x_{12} & \cdots & x_{1d} \\ x_{21} & x_{22} & \cdots & x_{2d} \\ \vdots & \vdots & \ddots & \vdots \\ x_{n1} & x_{n2} & \cdots & x_{nd} \end{pmatrix}$$

- Expected value of the vector random variable ( $\mathbf{X}$ )

$$\boldsymbol{\mu} = E[\mathbf{X}] = (\mu_1 \quad \mu_2 \quad \cdots \quad \mu_d)^T$$

Mean Vector

$$\hat{\boldsymbol{\mu}} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i$$

Sample mean vector

- Symmetric
- Positive, semi-definite (PSD)
- Total variance is sum of diagonal (trace)

$$\Sigma = \begin{pmatrix} \sigma_1^2 & \sigma_{12} & \cdots & \sigma_{1d} \\ \sigma_{21} & \sigma_2^2 & \cdots & \sigma_{2d} \\ \cdots & \cdots & \cdots & \cdots \\ \sigma_{d1} & \sigma_{d2} & \cdots & \sigma_d^2 \end{pmatrix} \quad \hat{\Sigma} = \begin{pmatrix} \hat{\sigma}_1^2 & \hat{\sigma}_{12} & \cdots & \hat{\sigma}_{1d} \\ \hat{\sigma}_{21} & \hat{\sigma}_2^2 & \cdots & \hat{\sigma}_{2d} \\ \cdots & \cdots & \cdots & \cdots \\ \hat{\sigma}_{d1} & \hat{\sigma}_{d2} & \cdots & \hat{\sigma}_d^2 \end{pmatrix}$$

$$tr(\Sigma) = \sigma_1^2 + \sigma_2^2 + \cdots + \sigma_d^2$$

- Pairwise inner or dot product of centered attribute vectors normalized by sample size
- Sum of rank-one matrices calculated as outer product of each centered point

$$\hat{\Sigma} = \frac{1}{n} (\bar{D}^T \bar{D}) = \frac{1}{n} \begin{pmatrix} \bar{X}_1^T \bar{X}_1 & \bar{X}_1^T \bar{X}_2 & \cdots & \bar{X}_1^T \bar{X}_d \\ \bar{X}_2^T \bar{X}_1 & \bar{X}_2^T \bar{X}_2 & \cdots & \bar{X}_2^T \bar{X}_d \\ \vdots & \vdots & \ddots & \vdots \\ \bar{X}_d^T \bar{X}_1 & \bar{X}_d^T \bar{X}_2 & \cdots & \bar{X}_d^T \bar{X}_d \end{pmatrix}$$

$$\hat{\Sigma} = \frac{1}{n} \sum_{i=1}^n \bar{\mathbf{x}}_i \cdot \bar{\mathbf{x}}_i^T$$



# Data Normalization



# Min-Max or Range Normalization



- Each sample scaled by the sample range
- Features normalized between 0 and 1

$$x'_i = \frac{x_i - \min_i \{x_i\}}{\hat{r}} = \frac{x_i - \min_i \{x_i\}}{\max_i \{x_i\} - \min_i \{x_i\}}$$

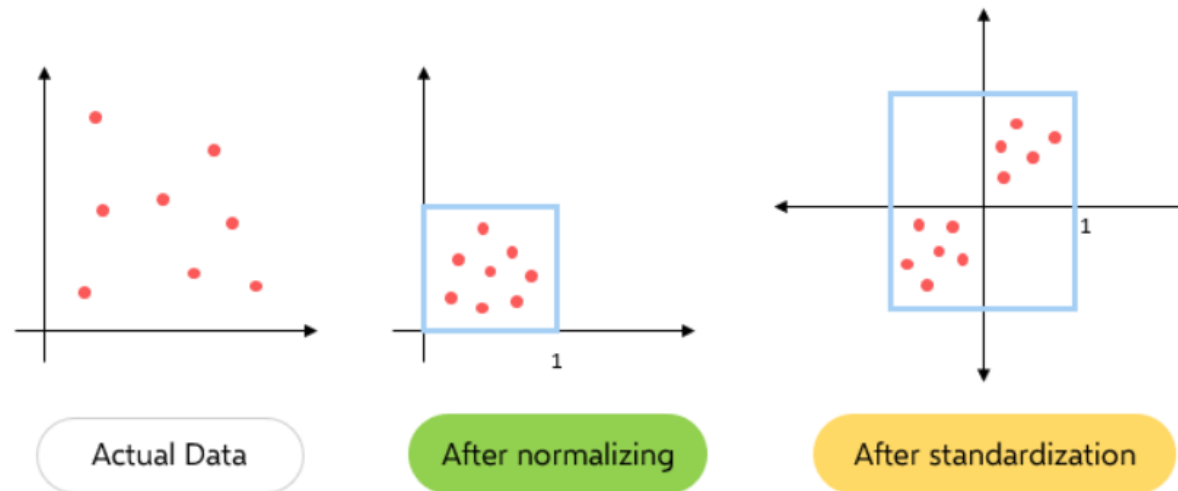
- Z-normalization
- Scales data to be centered (zero mean) and unit variance

$$x'_i = \frac{x_i - \hat{\mu}}{\hat{\sigma}}$$

# Range vs Standard Score Normalization



- Range normalization can be useful where distribution of the data is unknown and in algorithms that do not make assumptions of distribution of the data.
- Standardization is well suited to data that is characterized by a Normal (aka Gaussian) distribution. Its application is not just restricted to such data. Standardization is more robust to outliers.



# Next class



TEXAS A&M UNIVERSITY  
Engineering

- Data and attributes
  - Numerical
    - Normal distribution
  - Categorical



TEXAS A&M UNIVERSITY  
Engineering

**Thank You! Questions?  
Joshua Peeples, Ph.D.**

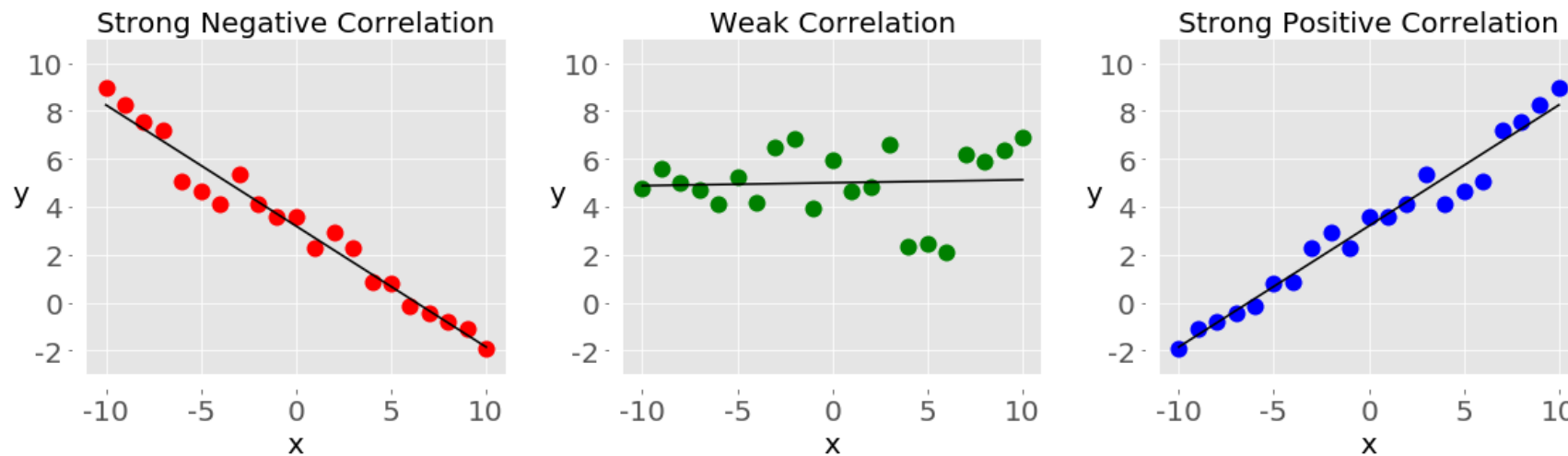
**<https://www.joshpeeples.com/>**  
**[jpeeples@tamu.edu](mailto:jpeeples@tamu.edu)**



TEXAS A&M UNIVERSITY  
Engineering

# Supplemental Slides

- Python [statistics](#) is a built-in Python library for descriptive statistics. You can use it if your datasets are not too large or if you can't rely on importing other libraries.



- [NumPy](#) is a library for numerical computing, optimized for working with single- and multi-dimensional arrays. Its primary type is the array type called [ndarray](#). This library contains many [routines](#) for statistical analysis.
- [SciPy](#) is a library for scientific computing based on NumPy. It offers additional functionality compared to NumPy, including [scipy.stats](#) for statistical analysis.

1	2	3	4
5	6	7	8
9	10	11	12
13	14	15	16

Row one, columns two to four

```
>>> arr[1, 2:4]
array([7, 8])
```

1	2	3	4
5	6	7	8
9	10	11	12
13	14	15	16

All rows in column one

```
>>> arr[:, 1]
array([2, 6, 10, 14])
```

1	2	3	4
5	6	7	8
9	10	11	12
13	14	15	16

All rows after row two,  
all columns after column two

```
>>> arr[2:, 2:]
array([[11, 12],
       [15, 16]])
```

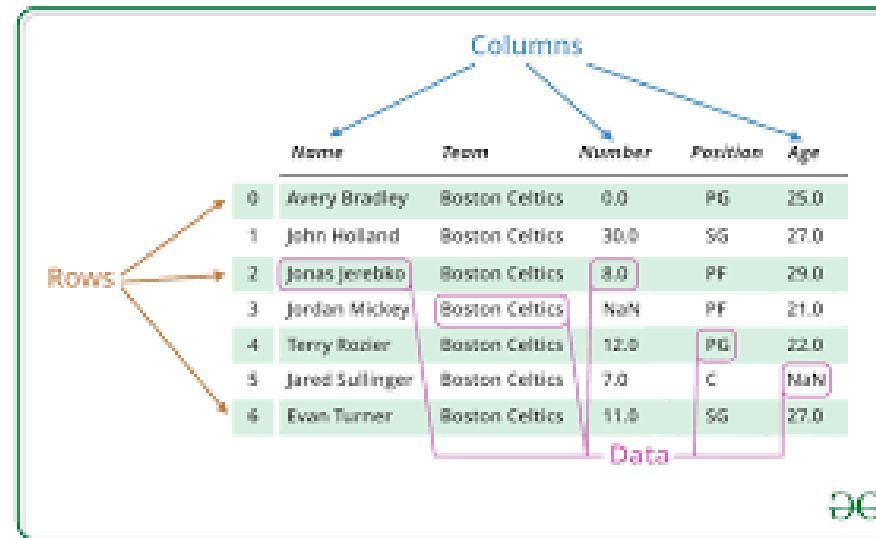
1	2	3	4
5	6	7	8
9	10	11	12
13	14	15	16

Every other row after row one,  
every other column

```
>>> arr[1::2, ::2]
array([[5, 7],
       [13, 15]])
```



- [Pandas](#) is a library for numerical computing based on NumPy. It excels in handling labeled one-dimensional (1D) data with [Series](#) objects and two-dimensional (2D) data with [DataFrame](#) objects.



The diagram illustrates a Pandas DataFrame structure. It features a table with 7 rows and 5 columns. The columns are labeled 'Name', 'Team', 'Number', 'Position', and 'Age'. The rows are indexed from 0 to 6. Annotations include: 'Columns' with arrows pointing to the column headers; 'Rows' with arrows pointing to the row indices; 'Data' with a box around the data cells; and a 'Data' label with an arrow pointing to the data cells. The table data is as follows:

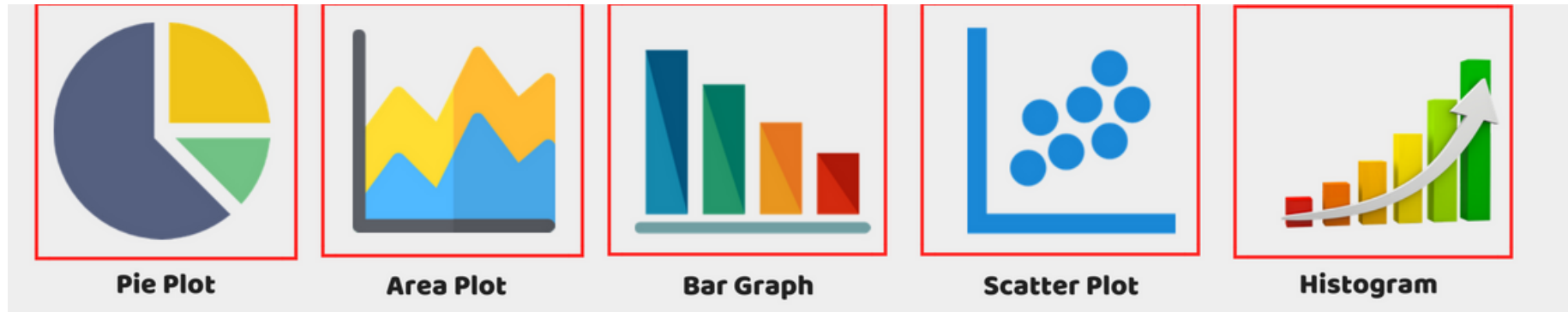
	Name	Team	Number	Position	Age
0	Avery Bradley	Boston Celtics	0.0	PG	25.0
1	John Holland	Boston Celtics	30.0	SG	27.0
2	Jonas Jerebko	Boston Celtics	8.0	PF	29.0
3	Jordan Mickey	Boston Celtics	NaN	PF	21.0
4	Terry Rozier	Boston Celtics	12.0	PG	22.0
5	Jared Sullinger	Boston Celtics	7.0	C	NaN
6	Evan Turner	Boston Celtics	11.0	SG	27.0

# Tools for Statistics



TEXAS A&M UNIVERSITY  
Engineering

- [Matplotlib](#) is a library for data visualization. It works well in combination with NumPy, SciPy, and Pandas.



- [Seaborn](#) pair-plots give us a good way to view correlations between pairs of variables (features):

