

# Chaptr -10 Features

- Topics already covered in class:  
categorical, Ordinal & Quantitative
- Feature Transformation Techniques:  
Thresholding & Discretization.

Feature transformations that remove scale  
of quantitative features:  $\rightarrow$  Discrete

- ① Thresholding (choose a feature value as a split point)
- ② Discretization (Transforms quantitative feature to ordinal where each ordinal value is referred to as bin)

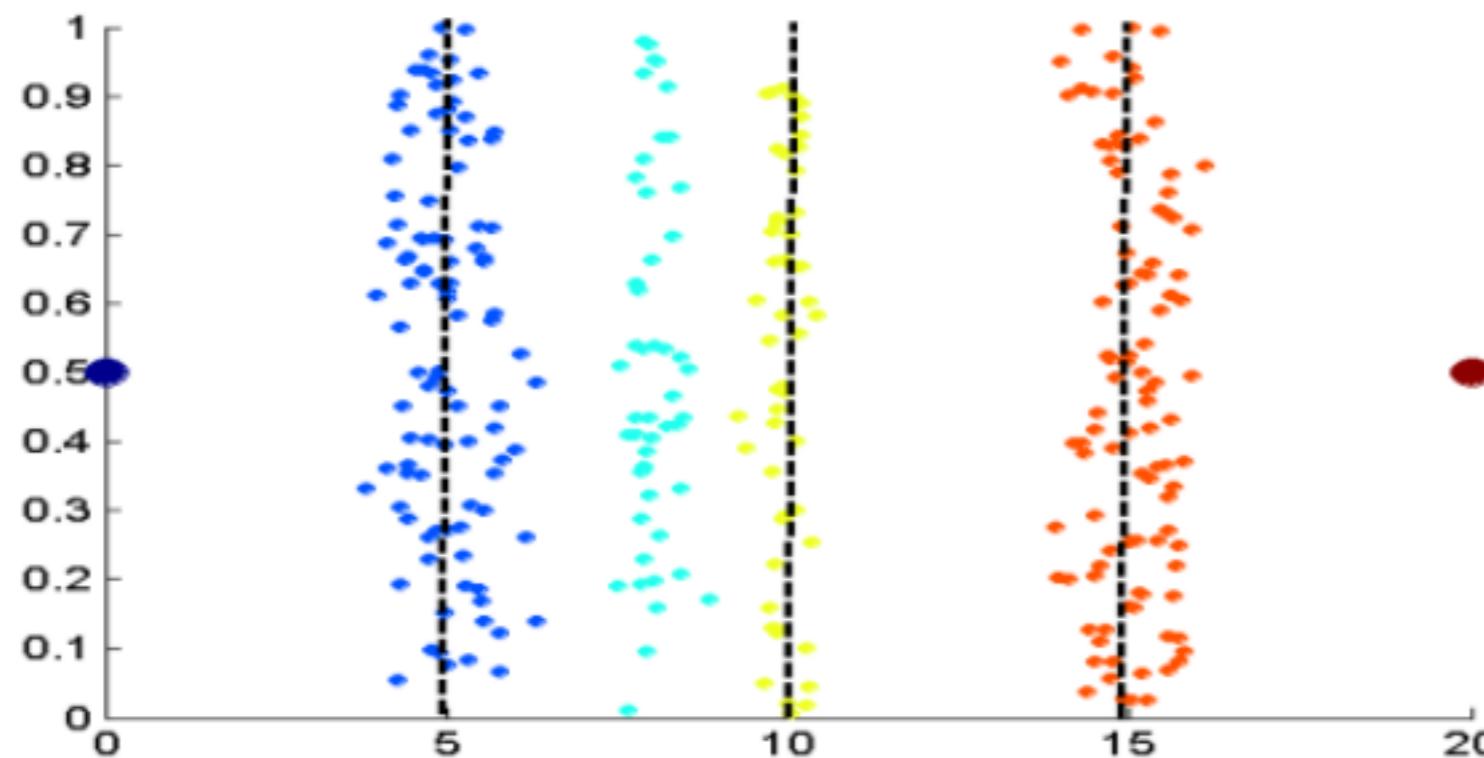
< Quantitative  $\rightarrow$  Discrete >  
(continuous)

## Unsupervised Discretization

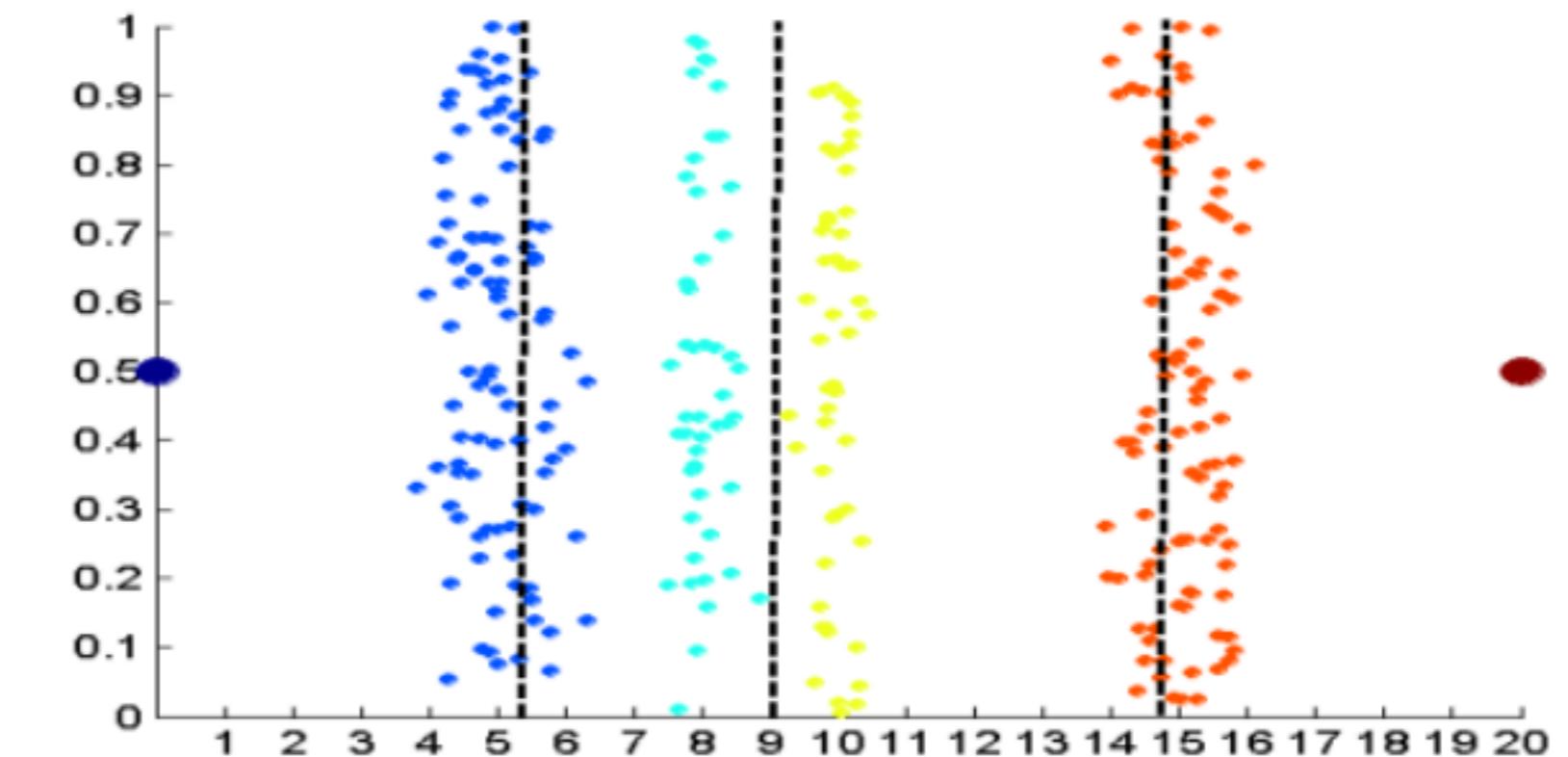
## Unsupervised Discretization

136

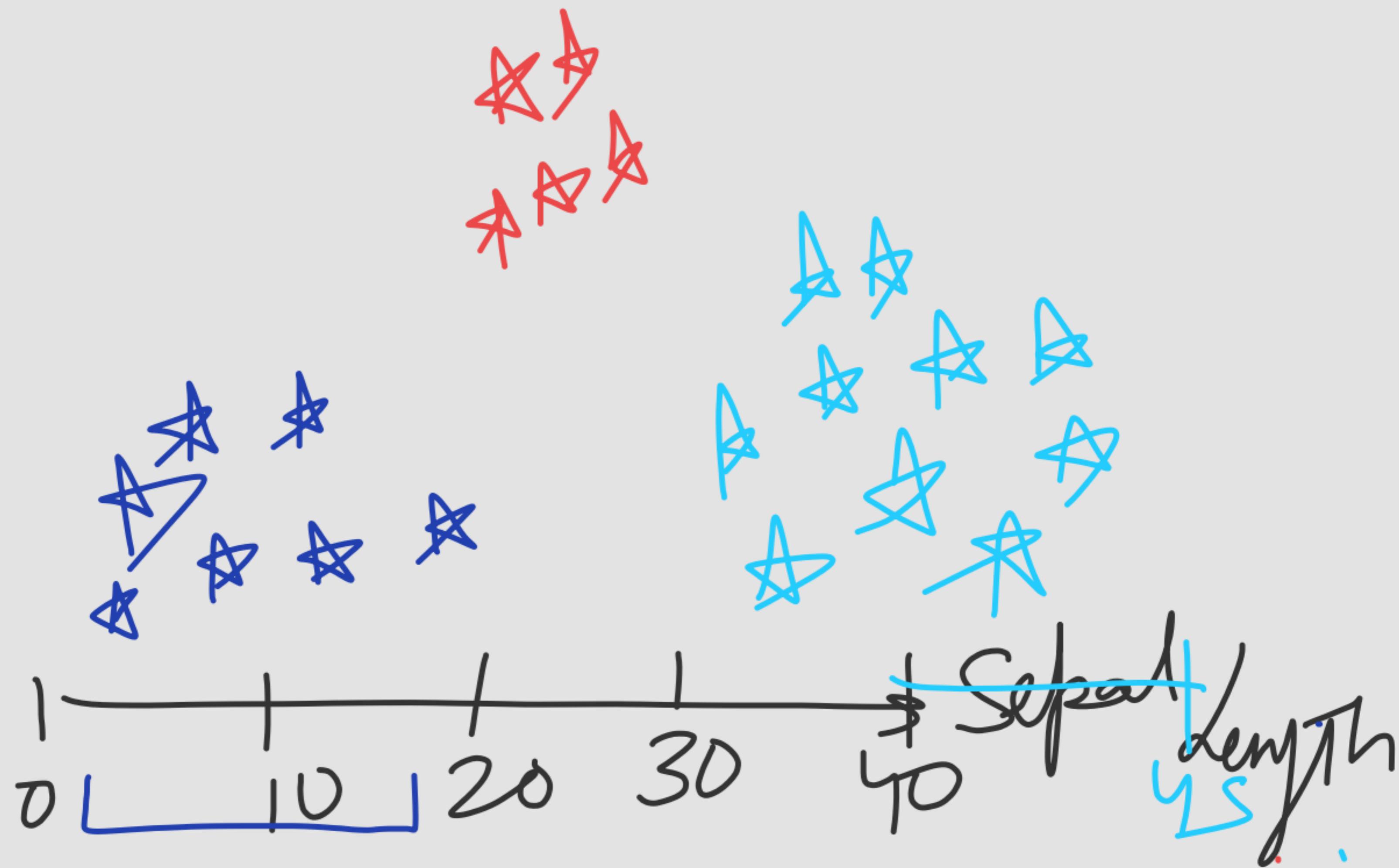
137



Equal interval width approach used to obtain 4 values.



Equal frequency approach used to obtain 4 values.

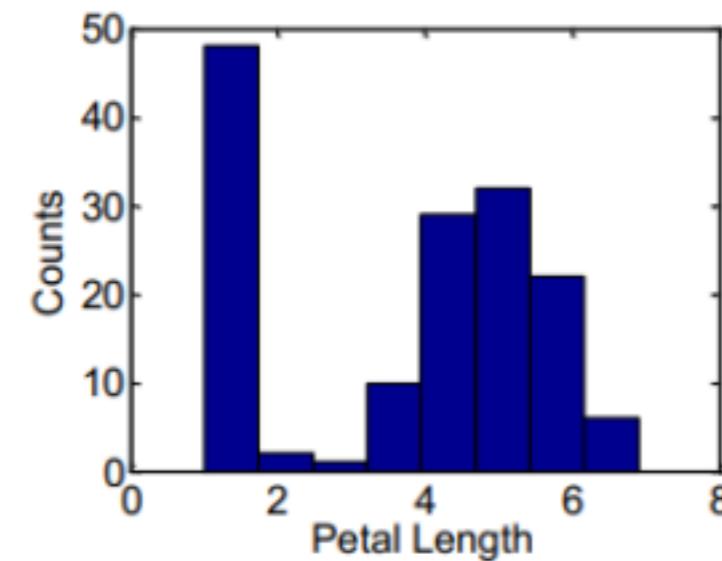


# Supervised Discretization

## Discretization: Iris Example ...

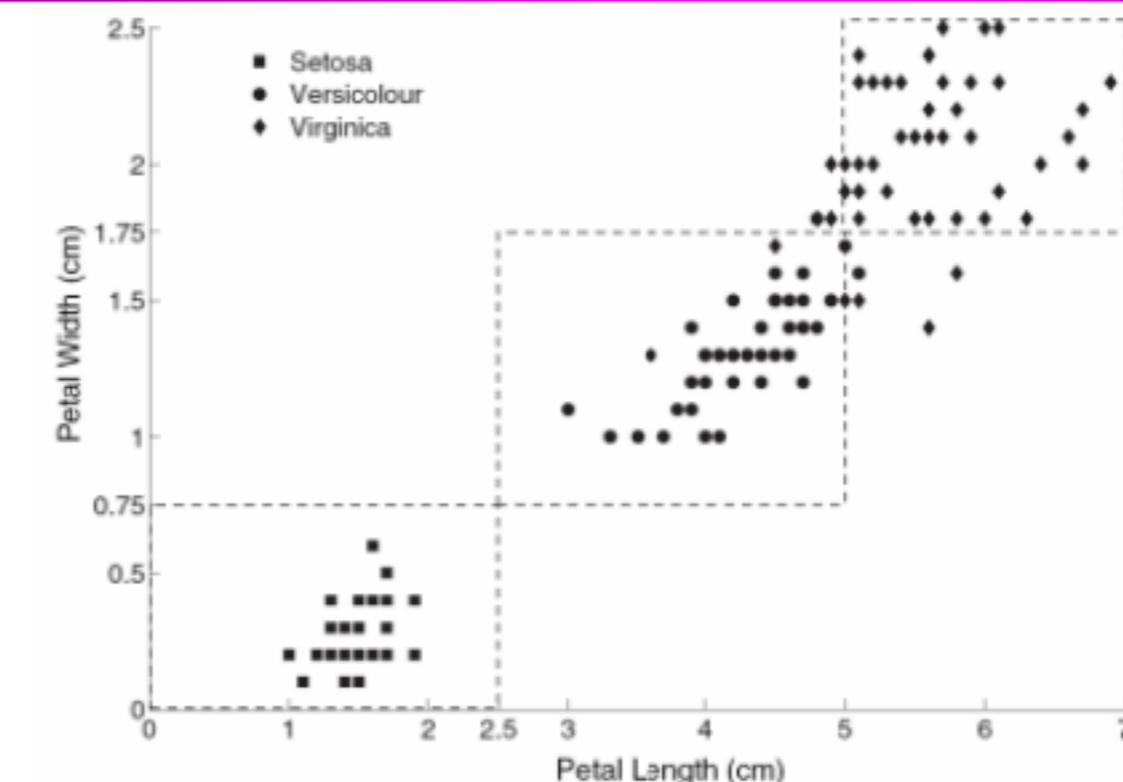
141

- How can we tell what the best discretization is?
  - Unsupervised discretization: find breaks in the data values
    - ◆ Example: Petal Length



## Discretization: Iris Example

142



Petal width low or petal length low implies Setosa.  
Petal width medium or petal length medium implies Versicolour.  
Petal width high or petal length high implies Virginica.

Feature Transformations that deal to  
Ordinal or Categorical features

Quantitative  
① Unsupervised

② Supervised

Normal + One-Hot Encoding  
Categorical  
Normalization → Ordinal  
Calibration

Categorical feature

→ Compute posterior prob. of  
class based on feature  
value

→ Adv. in Naive Bayes  
(No additional training on features)  
Calibrated

Ordinal feature

↓  
Discretize

↓  
Apply Calibration  
Similar to Categorical feature

(1) ~~object~~ Diabetes(y)

Transformed feature  
S.t. it encodes class info

1	Ob	1	
2	Ob	1	
3	Ob	1	
4	Ob	1	
5	Ob	1	
6	Ob	1	
7	Ob	1	
8	Ob	1	
9	Ob	1	
10	Ob	1	
11	Ob	1	
12	Ob	1	
13	Ob	1	
14	Ob	1	
15	Ob	1	
16	Ob	1	
17	Ob	1	
18	Ob	1	
19	Ob	1	
20	Ob	1	
21	Nob	0	
22	Nob	0	
23	Nob	0	
24	Nob	0	
25	Nob	0	
26	Nob	0	
27	Nob	0	
28	Nob	0	
29	Nob	0	
30	Nob	0	

$P(x=Obes / y=1)$

$P(y=1 / \text{obese}) = \frac{6}{9} \times \frac{1}{n} = \frac{6}{66} = \frac{1}{11}$

obes	Diabetes(y)
Ob	1
Ob	0
Ob	0
No b	1
No b	0

Calibrated feature

$$P\left(\frac{D}{Ob}\right) = \frac{6}{9} = \frac{2}{3}$$

$$P\left(\frac{D}{No b}\right) = \frac{1}{4} = \frac{1}{4}$$

# Feature Construction

- ↳ Take Cartesian product of original features
- ① Remove bias of Naive Bayes that have to be treated as independent
- ② Newly added feature may ↑ info gain  
(decision tree)

$x_1$ .Shape  $\in \{\text{Circle}, \text{Triangle}, \text{Square}\}$

$x_2$ .Color  $\in \{\text{Red}, \text{Green}, \text{Blue}\}$

Cartesian Product?

$x_1$	$x_2$
Circle	Red

$\frac{x_1 x_2}{\text{Red Circle}}$

For quantitative features

→ Take arithmetic or polynomial combination

Eg: Marks1 | Marks2 - . - | MarksN | — Pass

→ Dividing mass by  
Volume to get density.

# Curse of Dimensionality

↳ As the no. of dimensions ↑ so, it becomes difficult to learn a model

Classification  
→ hinders models learning capability

Clustering  
→ distance not meaningful

- ① Feature Selection
- ② Feature Reduction

# Feature Selection Approaches

Another way to reduce dimensionality of data- Use subset of features

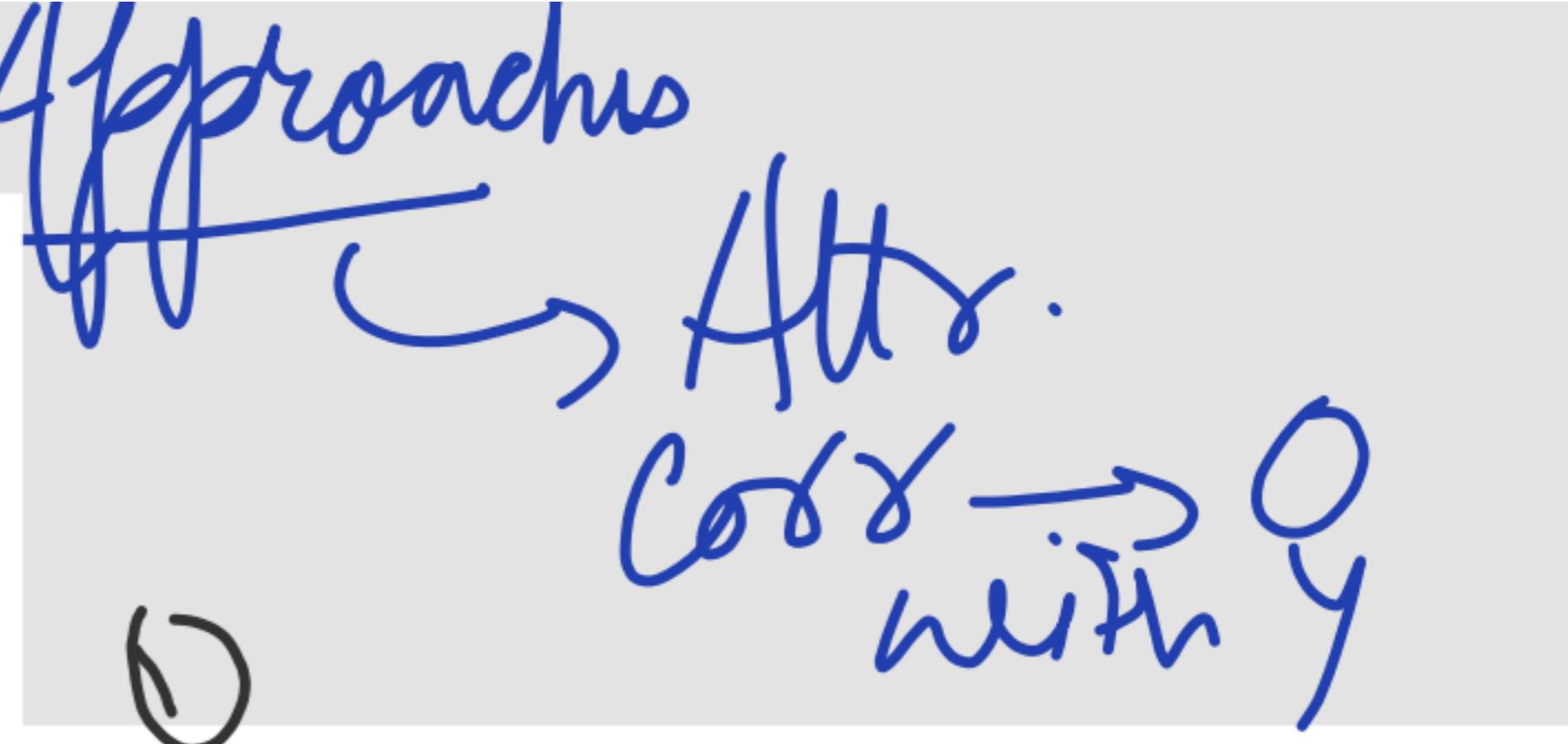
## Redundant features

- Duplicate much or all of the information contained in one or more other attributes
- Example: purchase price of a product and the amount of sales tax paid

## Irrelevant features

- Contain no information that is useful for the data mining task at hand
- Example: students' ID is often irrelevant to the task of predicting students' GPA

Many techniques developed, especially for classification

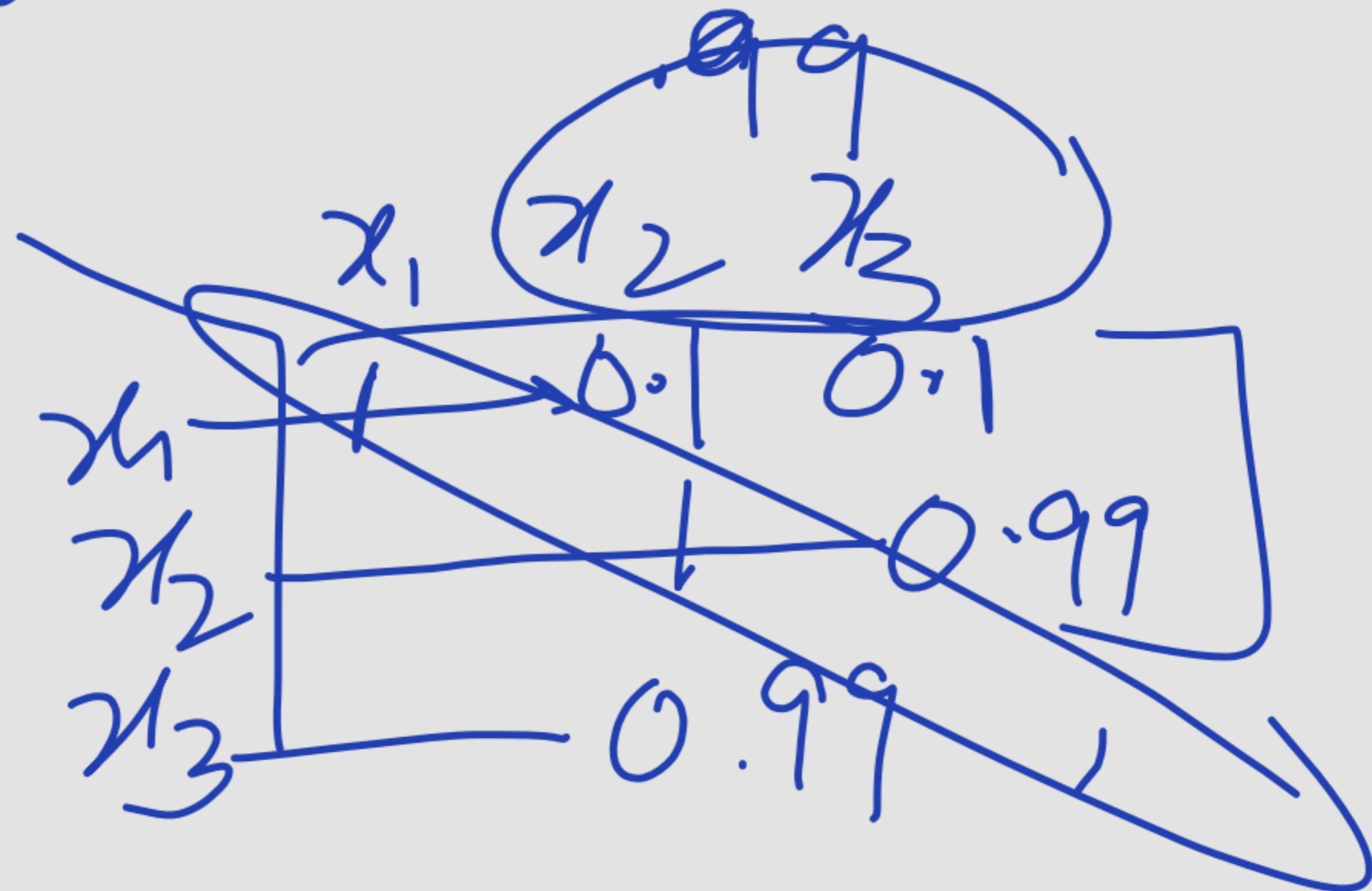


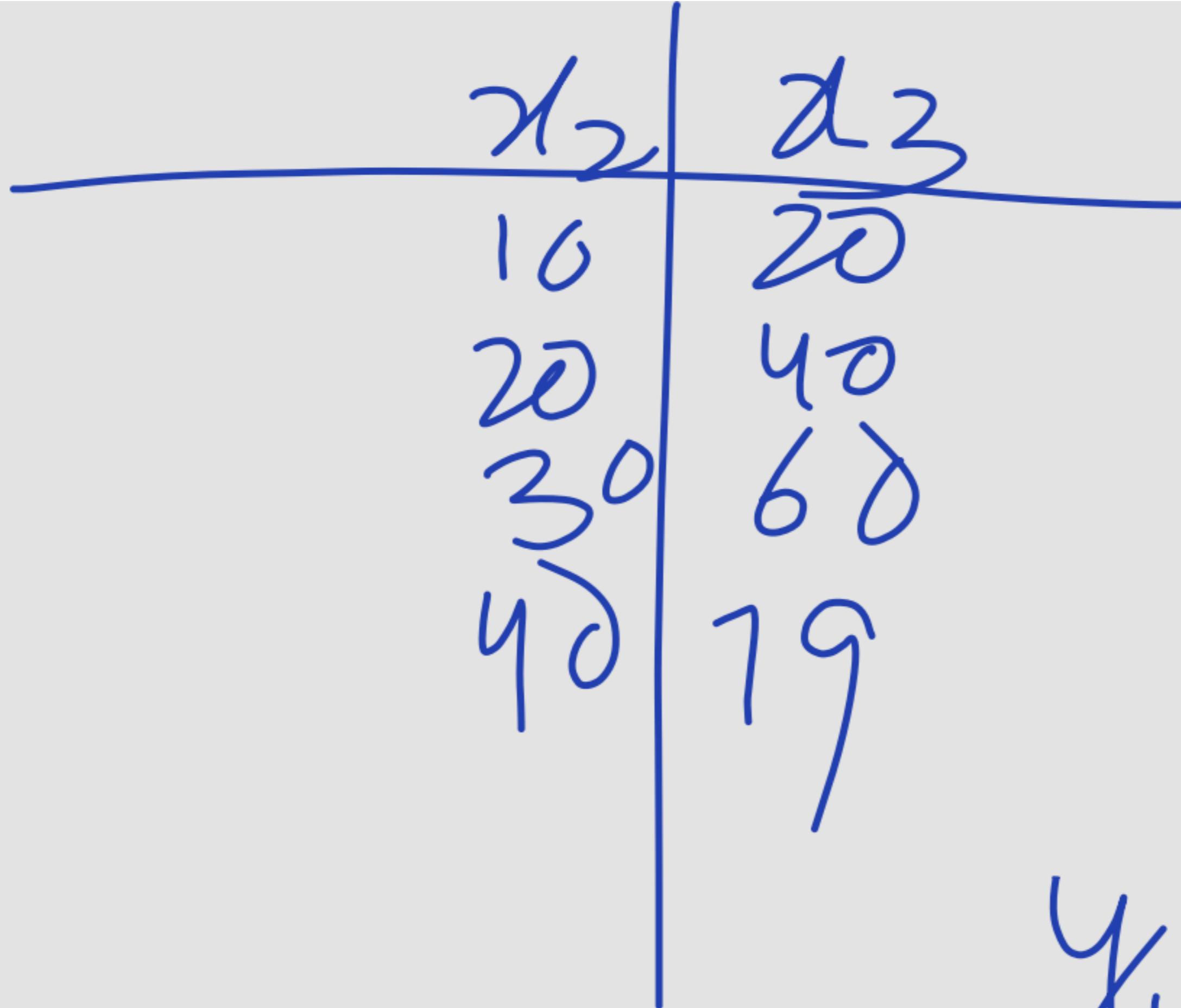
**Embedded approaches** Feature selection occurs naturally as part of the data mining algorithm. Specifically, during the operation of the data mining algorithm, the algorithm itself decides which attributes to use and which to ignore. Algorithms for building decision tree classifiers, which are discussed in Chapter 4, often operate in this manner.

**Filter approaches** Features are selected before the data mining algorithm is run, using some approach that is independent of the data mining task. For example, we might select sets of attributes whose pairwise correlation is as low as possible.

**Wrapper approaches** These methods use the target data mining algorithm as a black box to find the best subset of attributes, in a way similar to that of the ideal algorithm described above, but typically without enumerating all possible subsets.

O Pairwise correlation





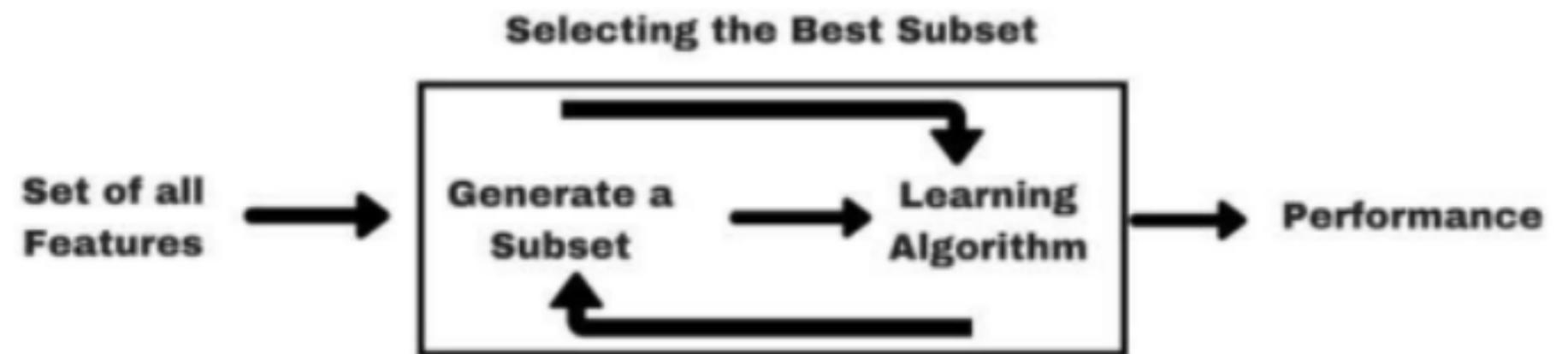
$$y = 2x_2$$

$$y = x_3$$

## 2. Filter Methods



## 3. Wrapper Methods



In wrapper methods, we try to use a subset of features and train a model using them. Based on the inferences that we draw from the previous model, we decide to add or remove features from your subset. The problem is essentially reduced to a search problem. These methods are usually computationally very expensive.

# Feature Reduction Approach: PCA

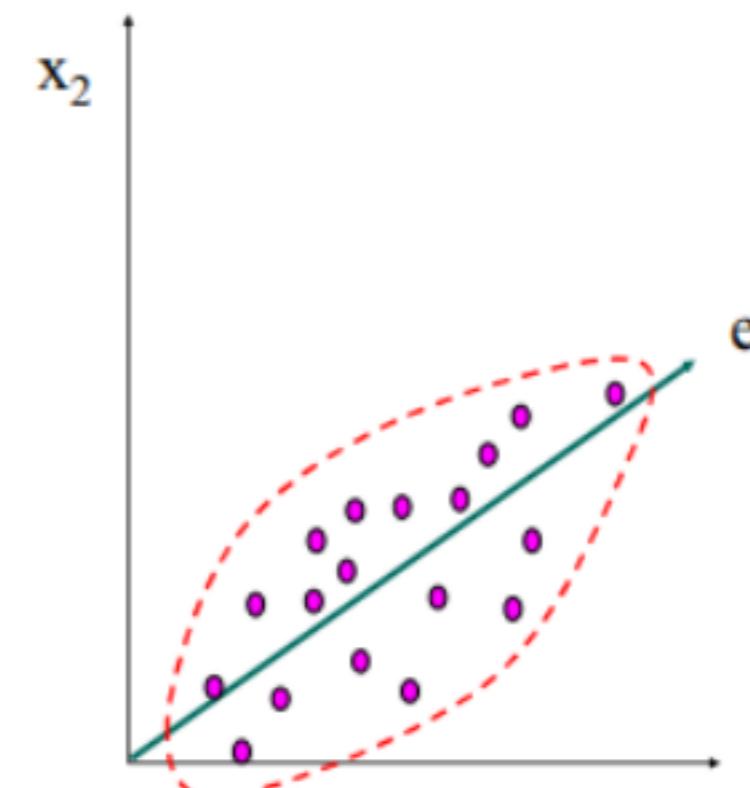
## Dimensionality Reduction: PCA

### □ Purpose:

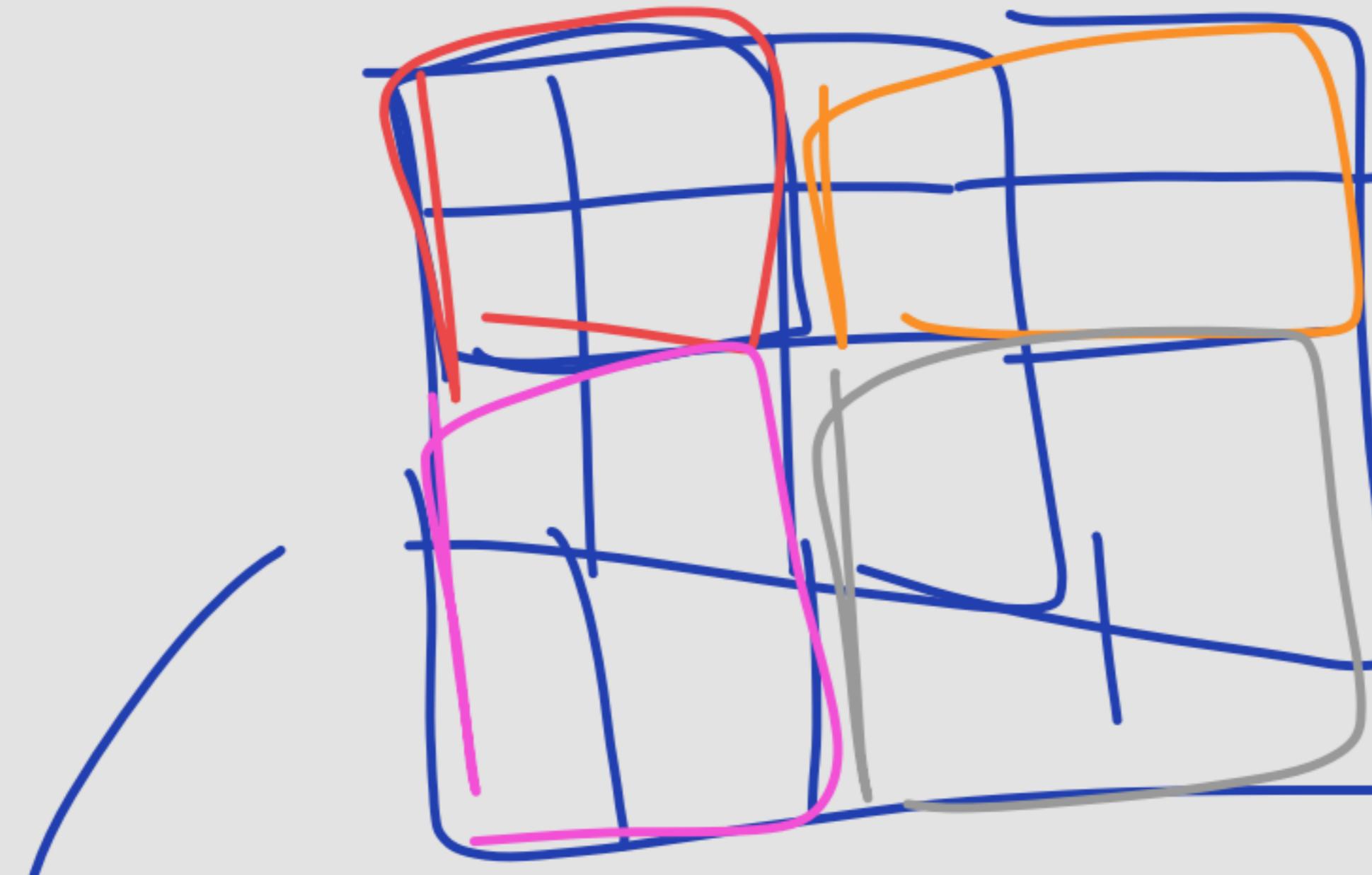
- Avoid curse of dimensionality
- Reduce amount of time and memory required by data mining algorithms
- Allow data to be more easily visualized
- May help to eliminate irrelevant features or reduce noise

### □ Techniques

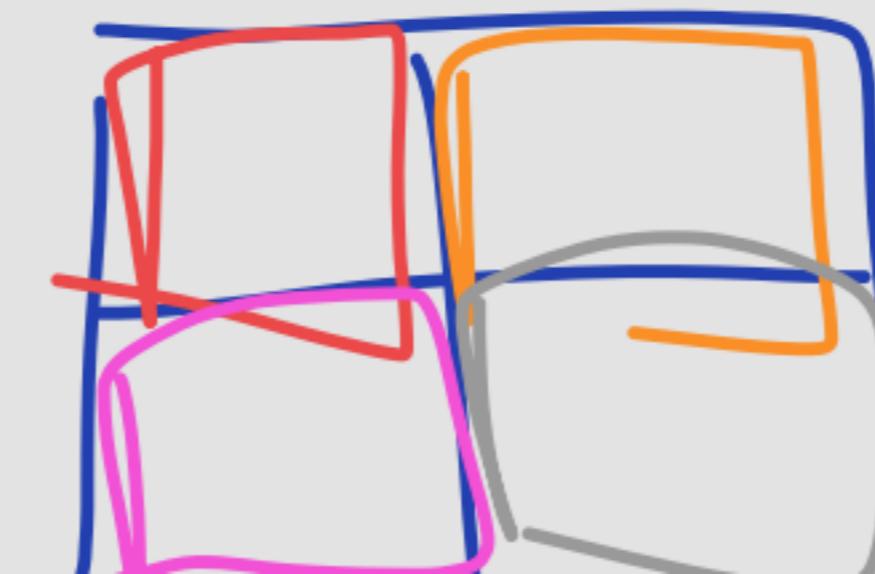
- Principal Components Analysis (PCA)
- Singular Value Decomposition
- Others: supervised and non-linear techniques



→ May be applied on  $256 \times 256$  image to yield  $64 \times 64$  image.

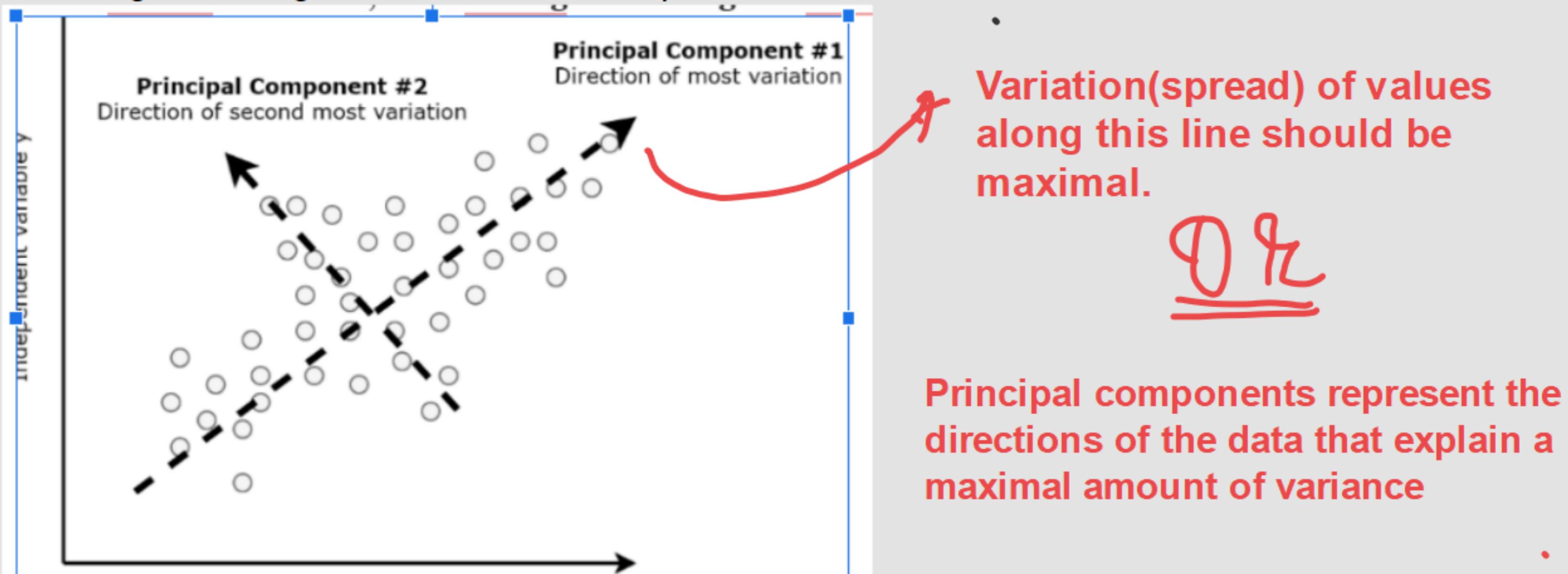


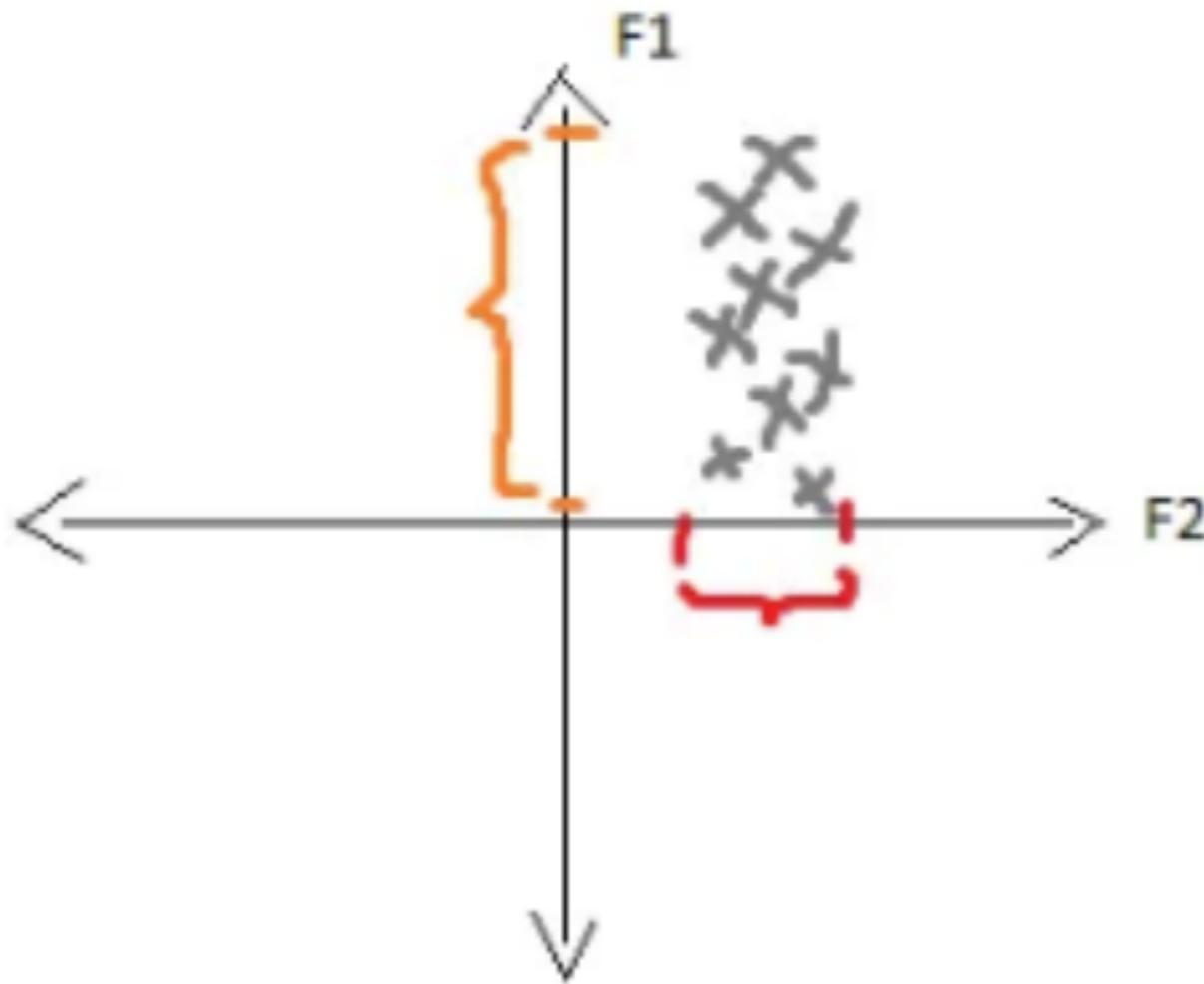
$4 \times 4$



$2 \times 2$

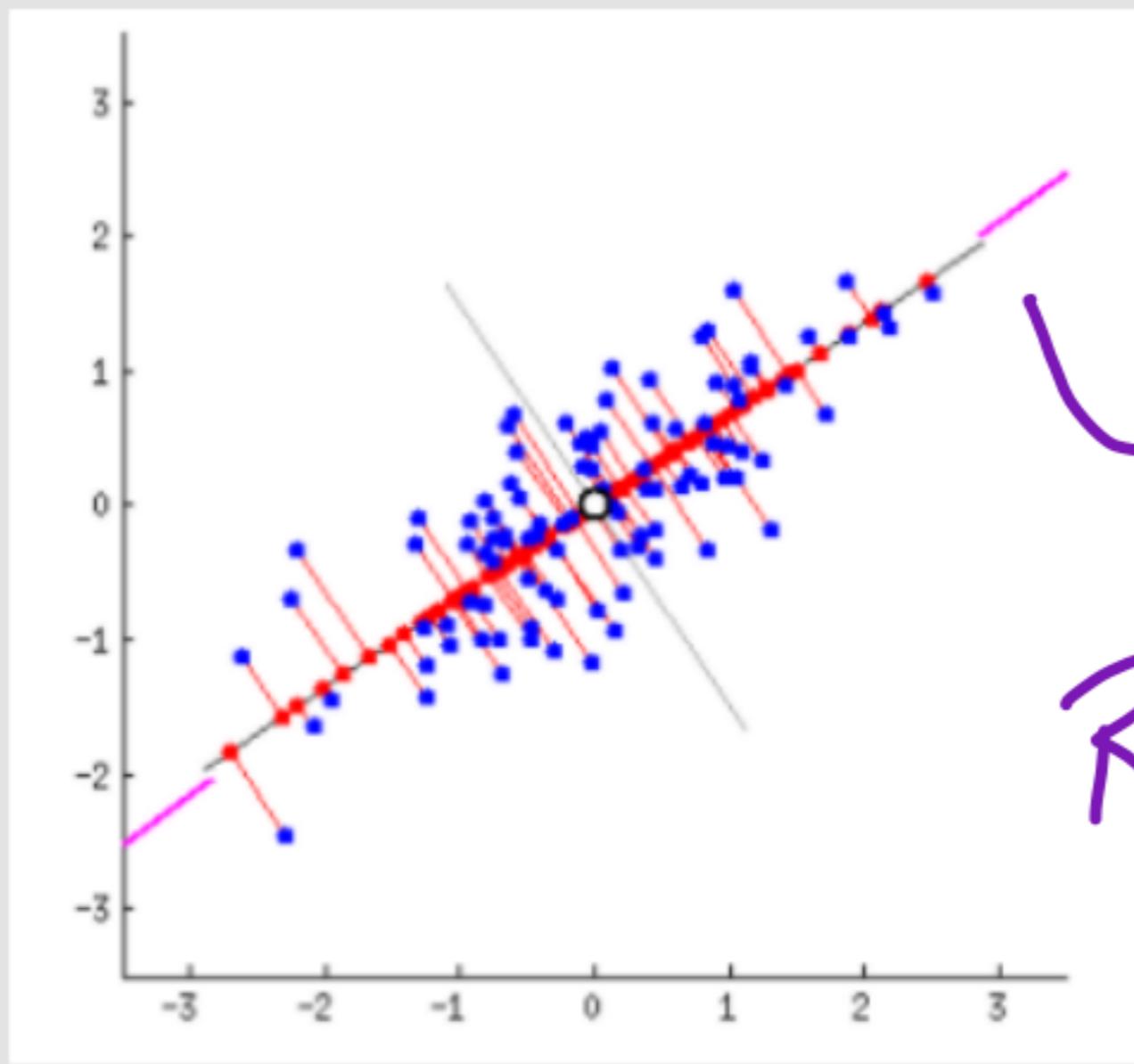
1. Principal components are new variables that are constructed as linear combinations or mixtures of the initial variables.
2. These combinations are done in such a way that the new variables (i.e., principal components) are uncorrelated and most of the information within the initial variables is squeezed or compressed into the first components.
3. So, the idea is 10-dimensional data gives you 10 principal components, but PCA tries to put maximum possible information in the first component, then maximum remaining information in the second and so on, until having something like shown in the scree plot below.





In above image 'x' represents the data points plotted in 2D plot. F1, F2 are the features in dataset(X). Orange selection shows the spread of data points on F1. Red selection shows the spread of data points on F2.

We can observe that spread/variance of data points on F1 is more than



Minimal  
Reconstruction  
Error!

## Summarizing the PCA approach

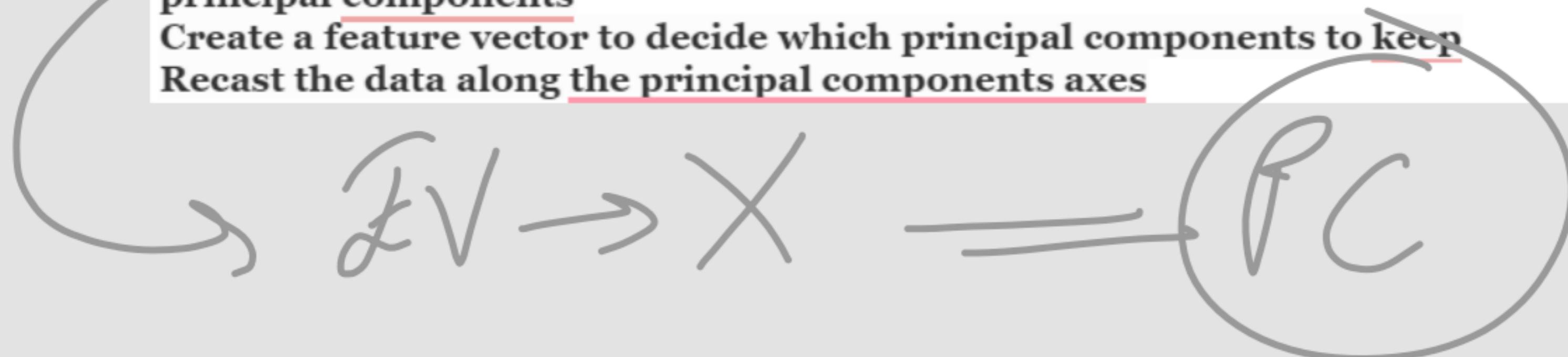
Listed below are the 6 general steps for performing a principal component analysis, which we will investigate in the following sections.

1. Take the whole dataset consisting of  $d$ -dimensional samples ignoring the class labels ✓
2. Compute the  $d$ -dimensional mean vector (i.e., the means for every dimension of the whole dataset)
3. Compute the scatter matrix (alternatively, the covariance matrix) of the whole data set ✓
4. Compute eigenvectors ( $e_1, e_2, \dots, e_d$ ) and corresponding eigenvalues ( $\lambda_1, \lambda_2, \dots, \lambda_d$ ) of covariance matrix ✓
5. Sort the eigenvectors by decreasing eigenvalues and choose  $k$  eigenvectors with the largest eigenvalues to form a  $d \times k$  dimensional matrix  $W$  (where every column represents an eigenvector) ➤ length or magnitude of eigen vector
6. Use this  $d \times k$  eigenvector matrix to transform the samples onto the new subspace. This can be summarized by the mathematical equation:  $y = W^T \times x$  (where  $x$  is a  $d \times 1$ -dimensional vector representing one sample, and  $y$  is the transformed  $k \times 1$ -dimensional sample in the new subspace.)

Just go through blue text

PCA we are trying to find the axes (EIGEN VECTORS) with maximum variances where the data is most spread

Standardize the range of continuous initial variables  
Compute the covariance matrix to identify correlations  
Compute the eigenvectors and eigenvalues of the covariance matrix to identify the principal components  
Create a feature vector to decide which principal components to keep  
Recast the data along the principal components axes



Eg: Binary attribute (Gender)

## KIND OF FEATURES

(Discrete) Qualitative

Categorical or Nominal

→ no scale or ordering

→ Eg: Eye Color

or Gender, Marital Status

→ Measure of Central Tendency:

Mode, Median

distinct values

→ Can only be compared

for equality  $=, \neq$

Feature Transformation

→ Improving utility of a feature by Scoring, Changing, Or adding info

(Discrete) Ordinal

→ Features with Ordering, but no Scale.

→ Ordinal features have Ordered

Set of discrete values as

Eg: Grade or House No. from best to worst

Measure of Central Tendency:

Mod, Median

Measures of Dispersion:

Quantiles

Discrete Values that provide enough info.

to Order Objects

(Continuous)

Quantitative

→ Features with a numerical Scale

Eg: Age, Price

→ Measure of CT: Mean

→ Measure of DT: range, Ig range

Variance Std dev

→ Continuous attr on which arithmetic

operators can be applied

Feature Transformation

→ Improving utility of a feature by Scoring, Changing, Or adding info

Category 1: Converting feature of One type to feature of another type

Category 2: Changing Scale of Quantitative features or adding a scale to Ordinal, categorical and boolean features

Applications: Linear Regression

Type of features req: Numerical / Quantitative

① I/P: Quantitative

② I/P f: Ordinal

Normalization  
(adapt the scale in unsupervised manner)

→ Change to numerical values such that Order is preserved

Calibration

↳ Supervised Feature transformation.

→ adding a meaningful scale carrying class info to arbitrary feature

(Good  
Better  
Best  
None)

Shows that 2 records are linked  
∴ they have common value for this

= # of Unique/Distinct feature values

(i.e. one for each value of Categorical feature)

\* Variance for

Marks

50  
60  
70

$$\frac{1}{3} \left( (10)^2 + (0)^2 + (10)^2 \right)$$

$$\sqrt{\frac{200}{3}} = \sqrt{66.67}$$

10

Deviation from mean

Do you know ?? There are 3 types of Statistics? ~~target~~

**1) Statistics of Central Tendency**

- 1.1 Mean or average Value
- 1.2 Median : Middle Value
- 1.3 Mode : Most Frequent Value (categorical)
- 1.4 & 5. Most Frequent grade (middle value)

\* Median is preferred to mean for skewed distributions.

That attempts to概括 the whole dataset

**2) Statistics of Dispersion**

- 2.1 Variance (Avg. squared deviation from mean)
- 2.2 Std Dev (Square root of Var.) (Both same, SD measured on same scale, as feature i.e. kg as compare to kg)
- 2.3 Range : Diff between max. to min. value
- 2.4 Percentile : p-th percentile is the value st. p% of the instances fall below it.

If multiple of 25, it is called quartile. If it is multiple of 10, it is called decade.

**3) Shape Statistics → Skewness ( $m_3/n^3$ ) :  $\frac{1}{n-1} \sum (x_i - \mu)^3$** : +ve : Distribution is right skewed, Right tail is longer than left tail -ve : Left skewed Distribution

**Skewness ( $m_4/n^4$ ) :  $\frac{1}{n-1} \sum (x_i - \mu)^4$** : +ve value : Distribution is sharply peaked than normal dist.

Q1 Diff b/w Arithmetic Mean :  $\frac{a+b}{2}$   
 Harmonic Mean :  $(\frac{1}{a})^{-1} + (\frac{1}{b})^{-1}/2 = \frac{2ab}{a+b}$   
 Geometric Mean :  $\sqrt{ab}$

extent to which distribution is spread/stretched  
 extent of deviation or spread or Temp  
 or Temp → Cold mean or few  
 Hot Central point

Q2 To Discrete

Application : Classification / Decision Tree  
 Say, Type of feature req: Categorical

Or We want target value to be a category  
 say CGPA given

Feature: Age	18	19	20	21	22
Class: CGPA	Pass	Pass	Pass	Fail	Fail

① I/P : Categorical ✓

② I/P : Ordinal  $\rightarrow$  UnOrdering  $\rightarrow$  Converts Ordinal feature to Category & Discardly

③ I/P : Quantitative Eg: Mapping [4, 5, 6] : [Short, Medium, Tall]

Techniques to convert Ordinal / Quantitative feature to Categorical

① Thresholding Single Threshold  $\rightarrow$  Choose a feature value as split point

$f: x_i \rightarrow \{Class1, Class2\}$

Unsupervised thresholding for any train instance  $i$ , if  $f(x_i^{(1)}) > u$  class 0 else class 1

For check the impact on Target Class

Most sensible thresholds are mean & median

Eg. Marks  $\rightarrow$  Grade

To Score / Mean / Median / Tail

Knee point

Pass / Fail

Ex: Marks  $\rightarrow$  Grade

To Score / Mean / Median / Tail

2 To Ordinal  
 Discretization  
 ↳ Transforms Quantitative feature to Ordinal feature

→ Each Ordinal Value is referred to as bin.  
 and Corresponds to an interval of the Original quantitative feature.

### UnSupervised Discretization

#### 2.1 Equal-frequency Discretization

Prereq → Need to know no. of bins

→ Choose the bins so that each bin has approx. equal no. of instances

#### 2.2 Equal Width Discretization

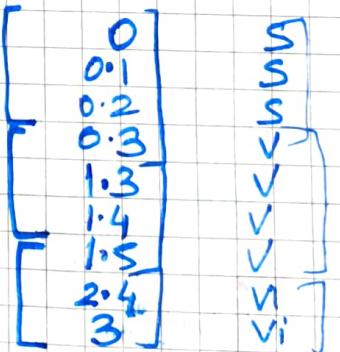
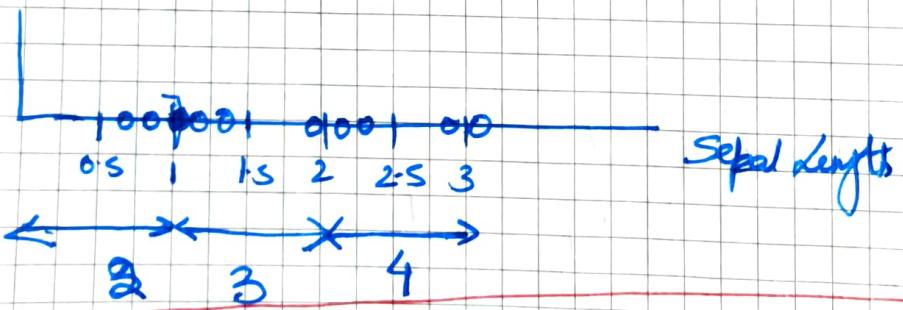
→ Choose bin boundaries so that each interval has the same width.  
 ↳ Can be determined by dividing the feature range by no. of bins.

### Supervised Discretization

(Top Down & Bottom Up)  
 Divide & Agglomerate

→ Progressively split bins  
 ↳ cluster

Start with each instance as a bin  
 Successively merge bins



\* Feature Transformations that make Scale of quantitative features

- ① Thresholding
- ② Discretization

\* Supervised Feature transformation carrying class info to a meaningful scale

- ① Feature Calibration