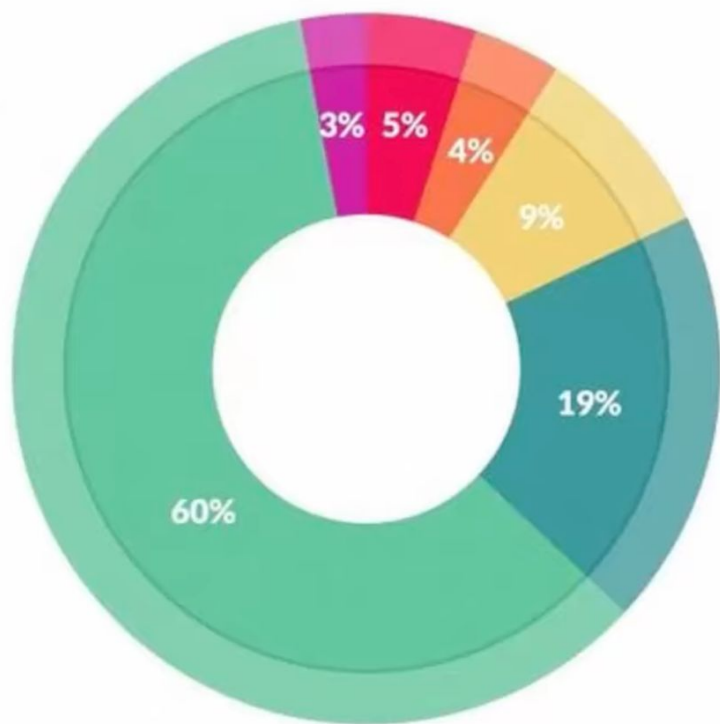# MIDS W207
# Applied Machine Learning

## Summer 2022

## Week 3

What data scientists spend the most time doing

- Building training sets: 3%
- Cleaning and organizing data: 60%
- Collecting data sets; 19%
- Mining data for patterns: 9%
- Refining algorithms: 4%
- Other: 5%

# Prediction



|   | Size | Beds | Baths | Zip | Price |
|---|------|------|-------|-----|-------|
| | 1100 | 1 | 1 | 64576 | 1.29 |
| | 1900 | 3 | 1.5 | 78321 | 2.14 |
| | 2800 | 3 | 3 | 98712 | 3.10 |
| | 3400 | 4 | 3.5 | 25721 | 3.75 |

Features · Label

Rows

Columns

Two classes in coordinate system

Two classes in polar coordinates

Feature engineering

Tangled

Transparent

Selection & Cleaning

Preprocessing

Feature Engineering

Machine Learning

Interpretation & Evaluation

Raw Data

Target Data

Preprocessed Data

Transformed Data

Patterns & Rules

Knowledge

# Missing Values

| | col1 | col2 | col3 | col4 | col5 |
|---|---|---|---|---|---|
| **0** | 2 | 5.0 | 3.0 | 6 | NaN |
| **1** | 9 | NaN | 9.0 | 0 | 7.0 |
| **2** | 19 | 17.0 | NaN | 9 | NaN |

**mean()** →

| | col1 | col2 | col3 | col4 | col5 |
|---|---|---|---|---|---|
| **0** | 2.0 | 5.0 | 3.0 | 6.0 | 7.0 |
| **1** | 9.0 | 11.0 | 9.0 | 0.0 | 7.0 |
| **2** | 19.0 | 17.0 | 6.0 | 9.0 | 7.0 |

# Transforming Features

# Scaling

# Bucketing

```
#Numerical Binning Example

Value        Bin
0-30    ->   Low
31-70   ->   Mid
71-100  ->   High


#Categorical Binning Example

Value        Bin
Spain   ->   Europe
Italy   ->   Europe
Chile   ->   South America
Brazil  ->   South America
```
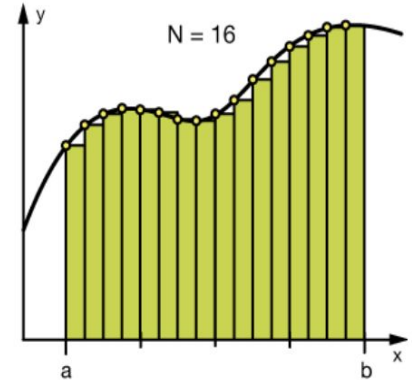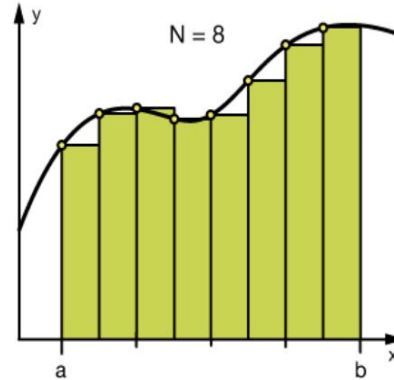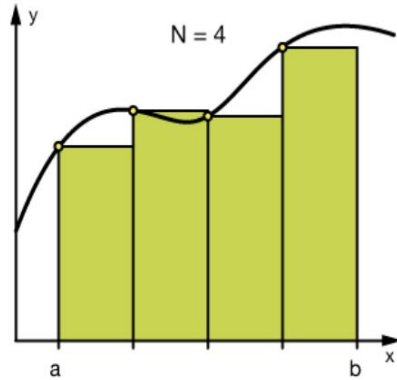
# Encoding

## Label Encoding

| Food Name | Categorical # | Calories |
|-----------|---------------|----------|
| Apple | 1 | 95 |
| Chicken | 2 | 231 |
| Broccoli | 3 | 50 |

$\rightarrow$

## One Hot Encoding

| Apple | Chicken | Broccoli | Calories |
|-------|---------|----------|----------|
| 1 | 0 | 0 | 95 |
| 0 | 1 | 0 | 231 |
| 0 | 0 | 1 | 50 |

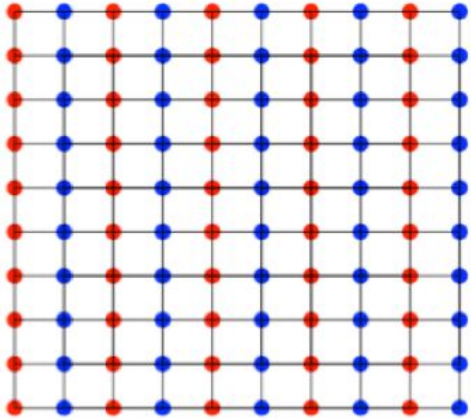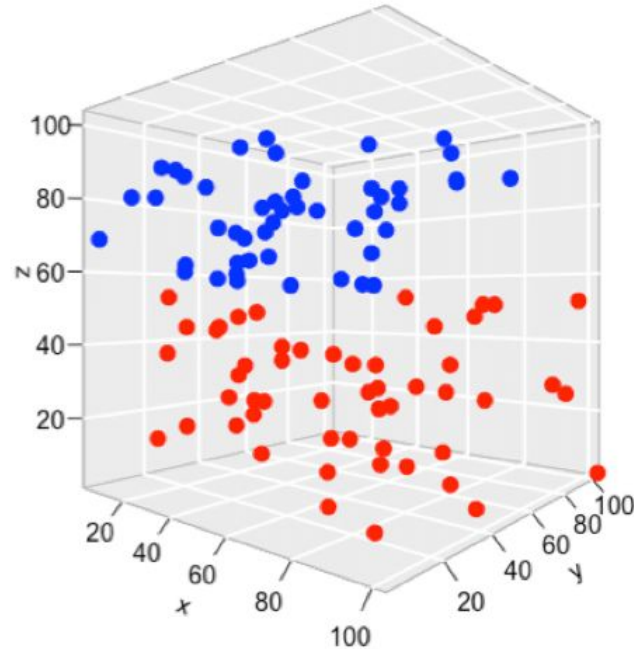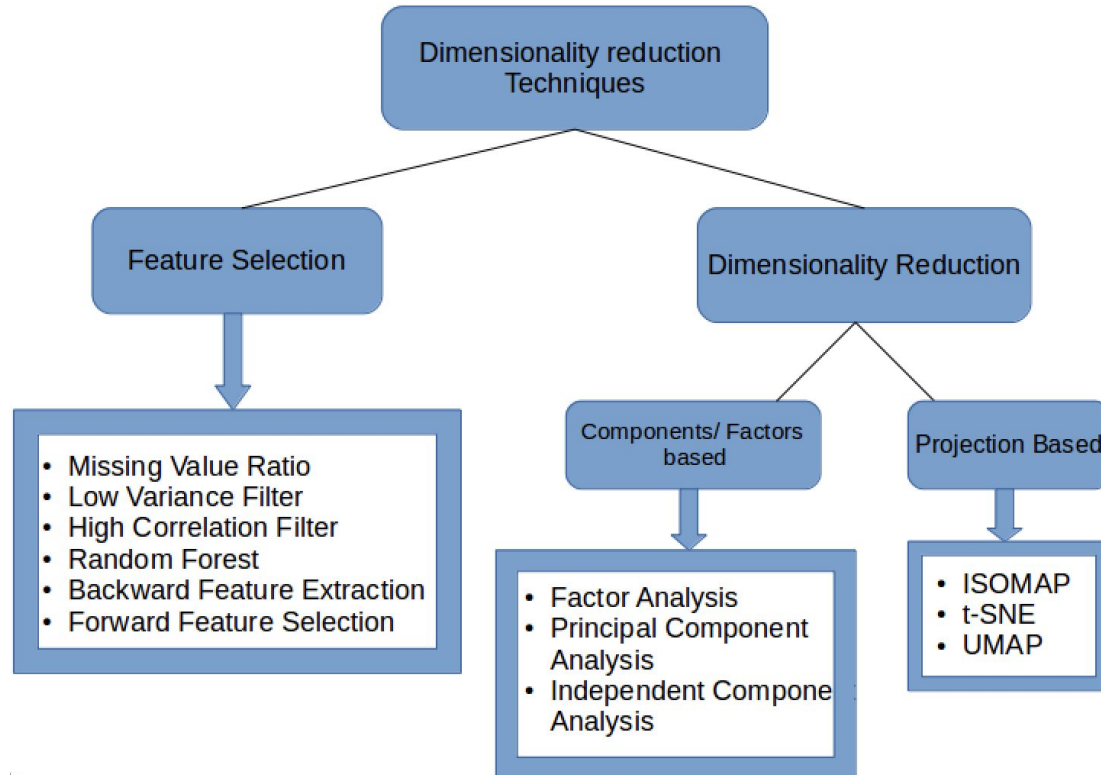| Feature Selection | Feature Extraction | Feature Engineering | Feature Learning |
|---|---|---|---|
| Subsetting the features | Creating new features when we could **NOT** have used raw features | Creating new features when we could have used raw features | Constructing features automatically |
| Ex: Using correlation with the dependent variable | Ex: from images to RGB values. Automatic methods such as PCA | Ex: Creating a new dummy variable for working days | Ex: Supervised neural networks, Independent component analysis |

# Curse of Dimensionality



(A) 1-D

(B) 2-D

(C) 3-D

# Dimensionality Reduction

## Numerical

- **Standardization**

$$Z = \frac{X - \mu}{\sigma}$$

- **Normalization**

$$X_{normalized} = \frac{(x - x_{minimum})}{(x_{minimum} - x_{minimum})}$$

- **Bucketing**

| Age<18 | 19<=Age<30 | 30<=Age<40 | Age>=40 |
|--------|------------|------------|---------|

## Categorical

- **One-hot encoding**



- **TF-IDF**
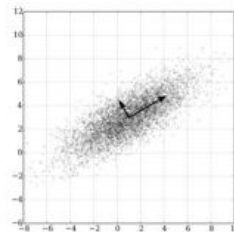
$$w_{i,j} = tf_{i,j} \times \log\left(\frac{N}{df_i}\right)$$

$tf_{i,j}$ = number of occurrences of $i$ in $j$
$df_i$ = number of documents containing $i$
$N$ = total number of documents

- **Word embeddings**



## Dimensionality Reduction

- **Principal component analysis (PCA)**



- **t-SNE**