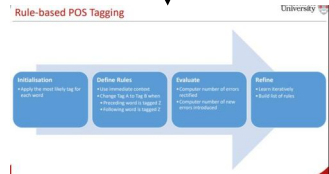


361 4



How to setup data for ML

Clustering and Classification

Split into Training and Testing Data

1. Training Process
2. Prediction
3. Evaluation

Also sometimes has validation data
(Data tested during testing process)

Has no effect on model

Supervised vs Unsupervised

Classification

Assigning categories to groups

Aim is to reduce dissimilarity in a cluster

Maps Feature Vector to Label

K Nearest Neighbour = Classify by closest Neighbours

Clustering

Identifying meaningful patterns/clusters/groups in observed data points

Hierarchical vs K-Means

K-Means = Optimisation Problem
Aim to minimise the within-cluster sum of squares (Variance)

Choose K by:
- Prior knowledge of data space
- Use elbow method
- Run hierarchical clustering on subset of data

K = Amount of Categories

- Efficient
- Number of Clusters can influence "wrong" results
- Non deterministic

```
1 K = Choose
2 while (true) {
3   For 1 to K: Assign point to closest centroid
4   For 1 to K: Compute new centroid by avg points in each cluster
5   If Centroids don't change
6     Stop
7 }
8
9 }
```

Hierarchical =
1. Initialisation
2. Merge iteratively
3. Result

- Deterministic
- Inefficient
- Multilevel Representation
- Flexible

Cluster Created by Splitting Data on Different Levels

Linkage = Distance
Single = 1 to 1
Complete = Furthest
Average = Average