

## K-means Clustering

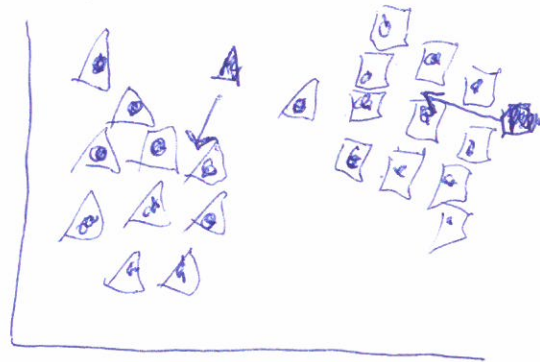
- unsupervised learning
- given features discover clusters ex) Developed / underdeveloped nations  
species from measurements  
~~partis~~ partisanship  
~~Ketipubaw~~ poor / wealthy  
neighborhood boundaries

Given  $x_1, \dots, x_n$  points

Place  $c_1, \dots, c_k$  centroids randomly

Eg

Weight



$k=2$

Items

① For each  $x_i$ :

- find the nearest centroid  $\min d(x_i, c_j)$
- assign  $x_i$  to cluster  $c_j$

② For each  $c_j$ :

- recalculate position of cluster  $j$   
using the mean

$$\bar{\Delta} = \frac{1}{n_{\Delta}} \sum_{i \in \Delta} x_i \quad \bar{\square} = \frac{1}{n_{\square}} \sum_{i \in \square} x_i$$

Repeat ① & ② until no cluster

<u>Al's</u>	<u>Type</u>	<u>±</u>	<u>Best use</u>	<u>Tuning</u>
Neural Networks	supervised	+ very accuracy high prob	- slow - uninterpretable - high variance - complex model - need a lot of data	- Layers - Backpropagation - Nodes - Loss functions

11 men → clusters	both supervised & unsupervised	+ fast + simple + interpretation	- specify clusters - cluster size be "sensible"	Classification text or non-text
-------------------	--------------------------------	--	---	---------------------------------------

Topic model →	unsupervised	+ interpretation + quite useful	- comparison/sim - depends on structure of data	- K - sample method
---------------	--------------	------------------------------------	--	------------------------

- sparse matrix  
topic >

✓

Algo Type

Random Forests

- + high accuracy
- + ~~high accuracy~~
- + good interpretability

Tuning

- Number of trees
- tree depth

Support Vector Machines

Supervised

- + high accuracy
- + fast - scales well
- + works for regression

Kernel

- Sensitivity to outliers
- Complexity
- model tuning

Kernel

- Kernel  $\Rightarrow$  linear  $\Rightarrow$  random basis
- C - regularization parameter
- gamma - regularization

Linear Regression / LASSO

Supervised

- + simple
- + fast
- + interpretation

Kernel

- kernel used for classification
- regularization
- poor performance

Regression

- regularization parameter  $\lambda$

Linear Regression (regularized)

Supervised

- + interpretation

Kernel

- inflexible for complex data

Regression

- regularization parameter  $\lambda$

Algorithm  
 $K-NN$

Type  
Supervised

- |   |                                    |
|---|------------------------------------|
| <u>+</u>                                  | <u>-</u>                           |
| + Very fast                               | - overfitting                      |
| + nonparametric                           | - poor accuracy                    |
| + useful for classification or regression | - Low interpretability             |
|   | - sensitive to irrelevant features |
- ↓  
scale feature

Best Use

Classification  
regression  
w/ low # of features

Tuning Parameters

$K$  - # of points around used in estimation

$\eta$

$\eta \uparrow$  prevents overfitting but may reduce accuracy

Naive Bayes

Supervised

- |                                      |                                    |
|--------------------------------------|------------------------------------|
| + very fast                          | - mostly useful for classification |
| + few assumptions                    | - missing accuracy                 |
| <del>high accuracy</del>             | - TF-IDF                           |
| + highly interpretable               | - pre-processing necessary         |
| + robust to variation of assumptions |                                    |
| + small                              |                                    |

Decision Classification

Laplace Smoothing

Decision trees

Supervised

- |                         |   |
|-------------------------|---|
| + high accuracy         | - Skilled person to avoid overfitting   |
| + high interpretability | - Sensitive to small data perturbations |
|                         | - Computationally intensive             |

Classification

→ pruning / tree depth