

k NN (k -Nearest Neighbors)

Boris Velichkov

Machine Learning

```
graph TD; ML[Machine Learning] --> SL[Supervised Learning]; ML --> UL[Unsupervised Learning]; SL --> C[Classification]; SL --> R[Regression]; UL --> Cl[Clustering]
```

Supervised Learning

Unsupervised Learning

Classification

Regression

Clustering

Classification

Dataset		Attribute 1	Attribute 2	Attribute N	Class	
		X1						„Yes“
		.						„No“
		.						„Yes“
		.						„Yes“
Xm							„Yes“	

$X_t = (a_1, a_2, \dots, a_N)$

$\text{Class}(X_t) = ?$

Regression

Dataset		Attribute 1	Attribute 2	Attribute N	Class	
		X1						1.5
		.						0.7
		.						1.8
		.						
		Xm						1.2

$X_t = (a_1, a_2, \dots, a_N)$

$\text{Class}(X_t) = ?$

Clustering

		Attribute 1	Attribute 2	Attribute N
Dataset	x1					
	.					
	.					
	.					
	xm					

$\forall x_i \in D, x_i = (a_1, a_2, \dots, a_N)$

$D = (x_1, x_2, \dots, x_m)$

Clusters (D) = ?

Global Learning vs Local Learning

- **Global Learning:** Learning from all instances in the dataset.
 - Naïve Bayes Classifier
- **Local Learning:** Learning from some of the instances in the dataset.
 - k NN

** Local & Global Learning is different from Local & Global Search!*

Instance-Based Learning vs Model-Based Learning

- **Instance-Based Learning:** Use the entire dataset as the model.
 - k NN
- **Model-Based Learning:** Use the training data to create a model that has parameters learned from the training data.
 - Naïve Bayes Classifier

Lazy Learning vs Eager Learning

- Lazy vs. eager learning
 - Lazy learning (e.g., instance-based learning): Simply stores training data (or only minor processing) and waits until it is given a test tuple
 - Eager learning (eg. Decision trees, SVM, NN): Given a set of training set, constructs a classification model before receiving new (e.g., test) data to classify
- Lazy: less time in training but more time in predicting
- Accuracy
 - Lazy method effectively uses a richer hypothesis space since it uses many local linear functions to form its implicit global approximation to the target function
 - Eager: must commit to a single hypothesis that covers the entire instance space

Lazy Learner: Instance-Based Methods

- Instance-based learning:
 - Store training examples and delay the processing (“lazy evaluation”) until a new instance must be classified
- Typical approaches
 - k-nearest neighbor approach
 - Instances represented as points in a Euclidean space.
 - Locally weighted regression
 - Constructs local approximation
 - Case-Based Reasoning (CBR)

kNN

- All instances correspond to points in the n-D space
- The nearest neighbor are defined in terms of Euclidean distance, $\text{dist}(X1, X2)$
- Target function could be discrete- or real-value
- For discrete-valued, k-NN returns the most common value among the k training examples nearest to the test instance.

kNN

- k-NN for real-valued prediction for a given unknown tuple
 - Returns the mean values of the k nearest neighbors
- Distance-weighted nearest neighbor algorithm
 - Weight the contribution of each of the k neighbors according to their distance to the query x_q
 - Give greater weight to closer neighbors
- Robust to noisy data by averaging k-nearest neighbors
- Curse of dimensionality: distance between neighbors could be dominated by irrelevant

kNN

- Kernel estimation
 - k-nearest neighbor

K = 3







