



## Classifier

ตัวเรียนรู้กระตือรือร้น (Eager Learner) ทำการประมาณค่าฟังก์ชันเป้าหมาย (target function) ณ ขณะเวลาเรียนรู้ด้วยชุดตัวอย่างสอน ค่าประมาณฟังก์ชันเป้าหมายที่ได้เป็นแบบ **Global Approximation** ซึ่งเป็นค่าประมาณที่ใช้จำแนกประเภทตัวอย่างทั้งหมดในปริภูมิตัวอย่าง รวมทั้งตัวอย่างค้นถามในอนาคต ได้แก่ Decision Tree, Naïve Bayes, Neural Networks

ตัวเรียนรู้เกียจคร้าน (Lazy Learner) จะรีรอไม่ดำเนินการเรียนรู้ด้วยชุดตัวอย่างสอน จนกว่าจะมีตัวอย่างใหม่ที่ต้องการจำแนกประเภทเข้ามา โดยจะสร้างค่าประมาณฟังก์ชันเป้าหมายแบบท้องถิ่น (local) ซึ่งค่าประมาณจะแตกต่างกันไปขึ้นอยู่กับค่าตัวอย่างละแวกใกล้เคียง (neighbors) กับตัวอย่างค้นถาม  $X_q$  (query instance)

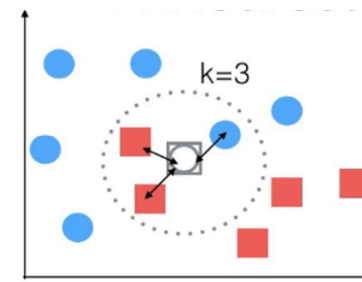
2110773-6 2/2567

2

## Instance-based/ Memory-based Learning

- memorizes the training observations for classifying the unseen test data, instead of performing explicit generalization.
- also termed as *lazy learning*, as it does not learn anything during the training phase like decision tree, neural networks, and so on
- Instead, it starts working only during the evaluation phase to compare the querying observations with nearest training observations, consuming significant time in comparison.
- the technique is thus not efficient on big data
- also, performance does deteriorate when the number of variables is high due to the *curse of dimensionality*.

## Non-Parametric Classifier: K-Nearest Neighbors



- Simple
- Non-Parametric means no assumption for underlying data distribution

2110773-6 2/2567

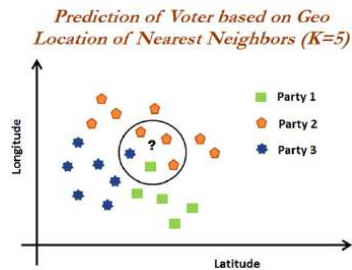
3

2110773-6 2/2567

4

## k-NN Voter Example

predict the party for which voter will vote based on their neighborhood, precisely geolocation (latitude and longitude)



- The querying voter will vote for **Party 2**. As within the vicinity, one neighbor voted for Party 1 and the other voter voted for Party 3. But three voters voted for Party 2.
- kNN solves any given **classification** problem with **Majority vote**.
- While **regression** problems are solved by taking **mean** of its neighbors within the given circle or vicinity or k-value.

2110773-6 2/2567

5

## Two parts of Classification/ Regression issues:

- First, how to measure “close”

- Minkowski distance

$$d(i, j) = \sqrt[q]{(|x_{i1} - x_{j1}|^q + |x_{i2} - x_{j2}|^q + \dots + |x_{ip} - x_{jp}|^q)}$$

- Manhattan distance (q=1)
  - Euclidean distance (q=2) is the most common one

- Second, how many close observations (K) neighbors

2110773-6 2/2567

6

## Circumventing Scaling Issues

- Similarity measured with Distance function can be dominated by those large-scaled attributes.
- Data Standardization/ Normalization; for example, Min-Max Normalization, Scale each feature to zero mean and unit variance.

$$x_{i, \text{std}} = \frac{x_i - \mu_i}{\sigma_i}$$

Note: A scale is an ordered set of values, continuous or discrete, or a set of categories to which an attribute is mapped. Type of scale: Ratio numeric, Interval numeric, Nominal, Ordinal.

2110773-6 2/2567

7

## k-NN Algorithm

Training Algorithm:

- For each training example  $\langle x, f(x) \rangle$ , add the example to the list *training\_examples*

Classification Algorithm:

- Given a query instance  $x_q$  to be classified,
    - Let  $x_1 \dots x_k$  denote  $k$  instances from *training\_examples* nearest to  $x_q$
    - CASE discrete-valued target function: RETURN
 
$$\hat{f}(x_q) \leftarrow \operatorname{argmax}_{v \in V} \sum_{i=1}^k \delta(v, f(x_i))$$

where  $\delta(a, b) = 1$  if  $a=b$  and  $\delta(a, b) = 0$  otherwise.
    - CASE real-valued target function: RETURN
 
$$\hat{f}(x_q) \leftarrow \frac{\sum_{i=1}^k f(x_i)}{k}$$

2110773-6 2/2567

8

# Weighted k-NN Algorithm

- Regression

$$\hat{f}(x_q) \leftarrow \frac{\sum_{i=1}^k w_i f(x_i)}{\sum_{i=1}^k w_i}$$

Intuitively, give more weight to the nearby points, and less weight to the farther away points. The simple function used is the inverse distance function.

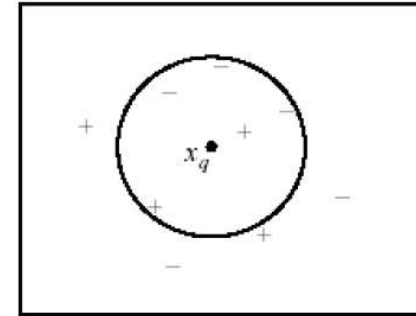
$$w_i \equiv \frac{1}{d(x_q, x_i)^2}$$

- Classification

$$\hat{f}(x_q) \leftarrow \operatorname{argmax}_{v \in V} \sum_{i=1}^k w_i \delta(v, f(x_i))$$

where  $\delta(a,b) = 1$  if  $a=b$  and  $\delta(a,b) = 0$  otherwise.

## ประเด็นค่า k

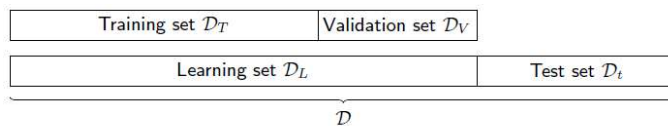


$K=1, \hat{f}(x_q) =$

$K=5, \hat{f}(x_q) =$

## Choosing k

- Goal is **generalization**: pick k (called a hyper-parameter) that performs best<sup>1</sup> on unseen (future) data.
- Split the dataset D:

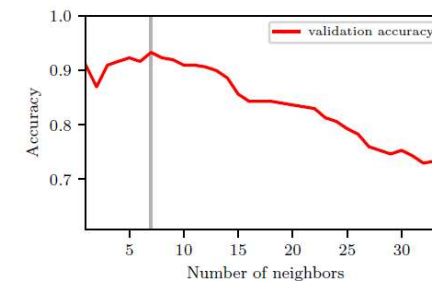


Hyper-parameter tuning procedure

- Learn the model using the training set
- Evaluate performance with different k on the *validation set* picking the best k
- Report final performance on the test set.<sup>2</sup>

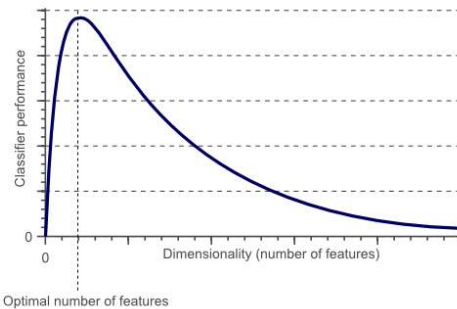
<sup>1</sup>In terms of some predefined metric, i.e. accuracy

## Using validation set to choose k



We choose  $k = 7$ .

## Curse of Dimensionality



**Figure 1.** As the dimensionality increases, the classifier's performance increases until the optimal number of features is reached. Further increasing the dimensionality without increasing the number of training samples results in a decrease in classifier performance.

2110773-6 2/2567 <https://www.visiondummy.com/2014/04/curse-dimensionality-affect-classification/>

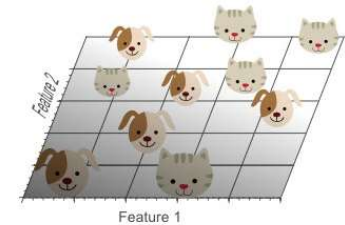
13

## Curse of Dimensionality and Overfitting

- Given 10 training images of cats and dogs
- Start by a single feature (Fig.2)
- Add another feature (Fig.3)



**Figure 2.** A single feature does not result in a perfect separation of our training data.



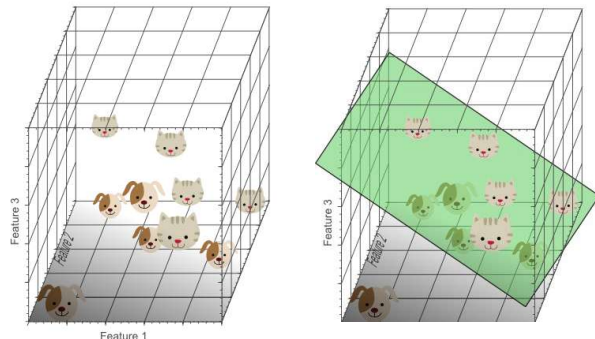
**Figure 3.** Adding a second feature still does not result in a linearly separable classification problem: No single line can separate all cats from all dogs in this example.

2110773-6 2/2567

14

## Curse of Dimensionality and Overfitting

- Add a third feature, yielding a three-dimensional feature space:



**Figure 4.** Adding a third feature results in a linearly separable classification. A plane exists that perfectly separates dogs from cats.

**Figure 5.** The more feature we use, the higher the likelihood that we can successfully separate the classes perfectly.

2110773-6 2/2567

15

## Curse of Dimensionality and Overfitting

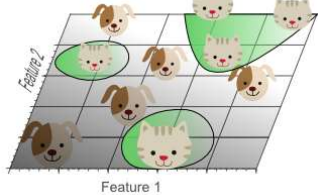
- Keep adding features, the dimensionality of the feature space grows, and becomes sparser and sparser.
- In the 1D case (figure 2), 10 training instances covered the complete 1D feature space, the width of which was 5-unit intervals. Therefore, in the 1D case, the sample density was  $10/5=2$  samples/interval.
- In the 2D case however (figure 3), we still had 10 training instances, which now cover a 2D feature space with an area of  $5 \times 5 = 25$ -unit squares. Here, the sample density was  $10/25 = 0.4$  samples/interval.
- Finally, in the 3D case, the 10 samples had to cover a feature space volume of  $5 \times 5 \times 5 = 125$ -unit cubes. Therefore, the sample density was  $10/125 = 0.08$  samples/interval.
- N has to grow exponentially with the number of features.
- By using less features, the *curse of dimensionality* was avoided such that the classifier did not *overfit* the training data.

2110773-6 2/2567

16

## Curse of Dimensionality and Overfitting

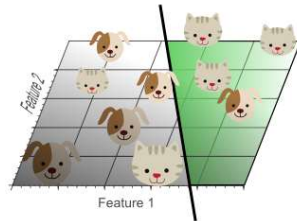
- Fig. 6 shows 3D classification results, projected onto a 2D feature space. Whereas the data was linearly separable in the 3D space, this is not the case in a lower dimensional feature space.



**Figure 6.** Using too many features results in overfitting. The classifier starts learning exceptions that are specific to the training data (by coincidence) and do not generalize well when new data is encountered.

2110773-6 2/2567

- Fig. 7 shows the result of a linear classifier that has been trained using only 2 features instead of 3.



**Figure 7.** Although the training data is not classified perfectly, this classifier achieves better results on unseen data than the one from figure 5.

17

## Issues of kNN

- Distance measures** → Normalization/ Standardization
- Expense of searching for nearest neighbors** as we need to compare with the entire of training data stored in the repository
  - Solution: use tree-based search structures, such as *k-d tree* search for efficient (approximate) NN on high-dimensional data
- Sensitive to Curse of Dimensionality**
  - Remove irrelevant attributes
  - Suggest large datasets
  - Suggest the number of attributes not more than 20

2110773-6 2/2567

18

## Pros & Cons & Constraints

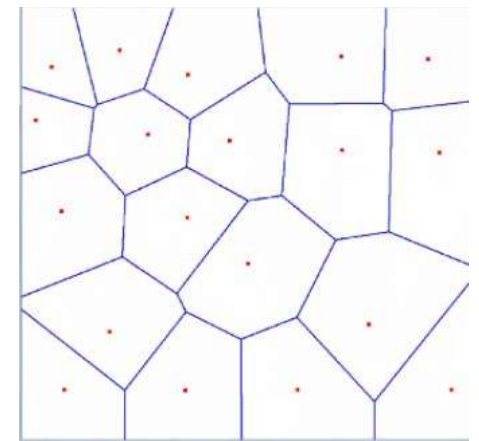
- Pros:** high accuracy, insensitive to outlier (not using all points), no assumption about data distribution
- Cons:** computationally expensive when classifying; susceptible to curse of dimensionality
- Constraints:** all features need to be standardized before fitting the model

2110773-6 2/2567

19

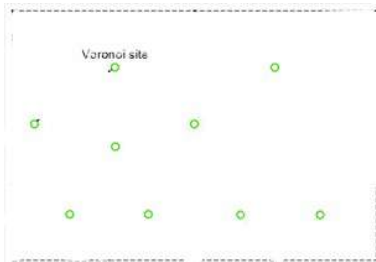
## Voronoi Diagram 1

- method decomposes a set of objects in a spatial space to a set of polygonal partitions.

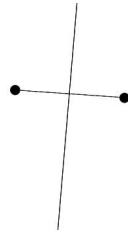


20

## Voronoi Diagram 2



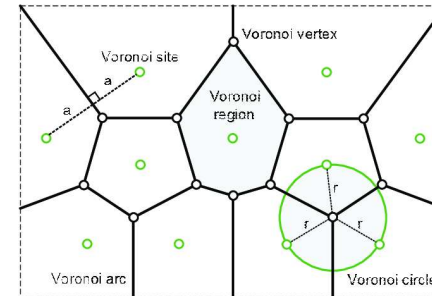
The points  $p_1, \dots, p_n$  are called **Voronoi sites**. The Voronoi diagram for two sites  $p_i$  and  $p_j$  can be easily constructed by drawing the perpendicular bisector of line segment  $\overline{p_i p_j}$ .



2110773-6 2/2567

21

## Voronoi Diagram 3



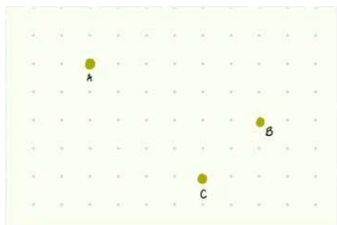
- Given a set of  $S$  points  $p_1, \dots, p_n$  in the plane, a Voronoi diagram divides the plane into  $n$  Voronoi regions with the following properties:
- Each point  $p_i$  lies in exactly one region.
- If a point  $q \notin S$  lies in the same region as  $p_i$ , then the Euclidean distance from  $p_i$  to  $q$  will be shorter than the Euclidean distance from  $p_j$  to  $q$ , where  $p_j$  is any other point in  $S$ .

2110773-6 2/2567

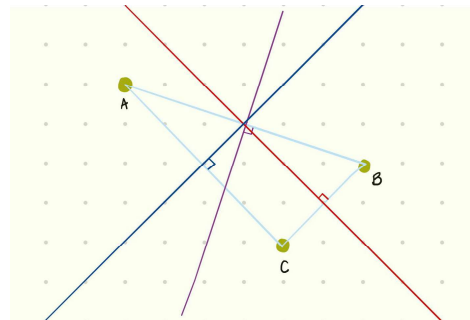
22

## Example

- Given Voronoi sites,  $p_i = A, B, C$



- Half plane,  $H(p_i p_j)$



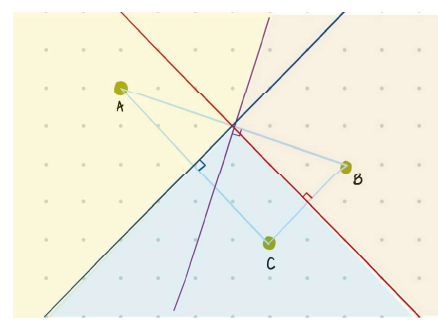
2110773-6 2/2567

23

## Example cont.

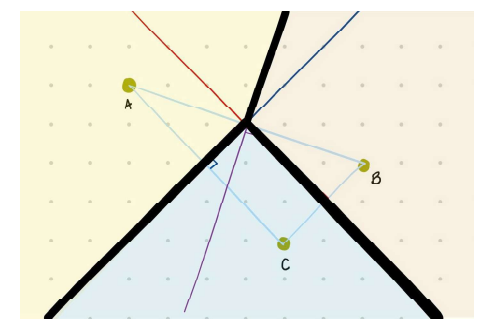
- Voronoi regions,  $V(p_i)$

$$V(p_i) = H(p_i p_1) \cap H(p_i p_2) \cap \dots \cap H(p_i p_n)$$



2110773-6 2/2567

- Voronoi edges** = line segments form the boundaries of Voronoi regions
- Voronoi vertex** = intersection of adjacent edges.



24