

# 732G12 Data Mining

## Föreläsning 4

---

Johan Alenlöv

IDA, Linköping University, Sweden

- K-närmaste grannar
- Bayesianska klassificerare
- Ensemblemetoder
  - Bagging
  - Boosting
  - Random forest

**Idé** basera predikation på de  $K$  datapunkter som är närmast.

Ger en icke-parametrisk metod för klassificering och regression.

Problem: Vad är närmast?

Vi behöver något som talar om för oss hur nära två datapunkter är.  
Finns många alternativ som man kan välja, som ger olika resultat.

## Euklidiskt avstånd

$$d(\mathbf{x}, \mathbf{y}) = \sqrt{\sum_{k=1}^n (x_k - y_k)^2}$$

## Manhattan avstånd

$$d(\mathbf{x}, \mathbf{y}) = \sum_{k=1}^n |x_k - y_k|$$

1. Låt  $k$  vara ditt valda antal grannar och  $D$  din träningsdata.

2. För varje testdata  $z = (\mathbf{x}', y')$   $\in D$ :

2.1 Beräkna  $d(\mathbf{x}, \mathbf{x}')$  (avståndet mellan  $z$  och all träningsdata)

2.2 Välj  $D_z \subseteq D$ , de  $k$  närmaste träningsdatan till  $z$

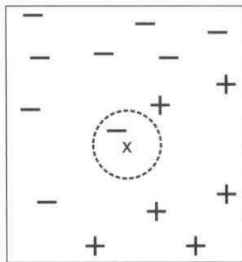
2.3 Låt  $y' = \arg \max_v \sum_{(\mathbf{x}_i, y_i) \in D_z} \mathbf{1}_{v=y_i}$

2.3 är majoritetsvalet. Kan också vikta detta värde med avståndet:

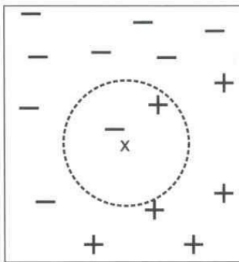
$$2.3 \quad y' = \arg \max_v \sum_{(\mathbf{x}_i, y_i) \in D_z} w_i \mathbf{1}_{v=y_i}.$$

För regression används medelvärde alternativt viktat medelvärde.

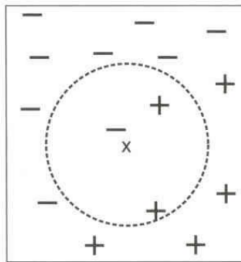
# K-närmaste grannar



(a) 1-nearest neighbor



(b) 2-nearest neighbor



(c) 3-nearest neighbor

- Målet med modellen är att prediktera nya observationer.
- Påverkas stort av olika skalor.
- Långsam anpassning.
- Känslig mot brus.
- Val av  $K$  har stor betydelse!
  - Litet  $K$  ger överanpassning.
  - Stort  $K$  ger underanpassning.
  - Korsvalidering kan användas för att bestämma  $K$ .
- Producerar godtyckligt utformade beslutsgränser.
- Problem i högre dimensioner.

Att direkt modellera en icke-deterministisk funktion kan vara mycket svårt.

Exempel:

- (diet, träning)  $\rightarrow$  (hjärtinfarkt) är svårt
- (diet, träning)  $\rightarrow \mathbb{P}(\text{hjärtinfarkt})$  lättare

Använd Bayes sats för att hjälpa till i modelleringen

$$\mathbb{P}(Y | \mathbf{X}) = \frac{\mathbb{P}(\mathbf{X} | Y)}{\mathbb{P}(\mathbf{X})} \cdot \mathbb{P}(Y) \propto \mathbb{P}(\mathbf{X} | Y) \cdot \mathbb{P}(Y)$$

$$\text{posterior} = \frac{\text{likelihood}}{\text{evidence}} \cdot \text{prior} \propto \text{likelihood} \cdot \text{prior}$$