

# 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  $z$

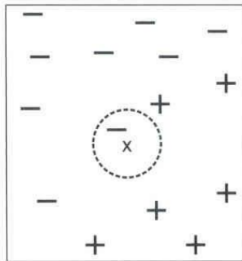
2.3 Låt  $y' = \arg \max_v \sum_{(\mathbf{x}_i, y_i) \in D_z} \mathbf{1}_{v=y_i}$

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

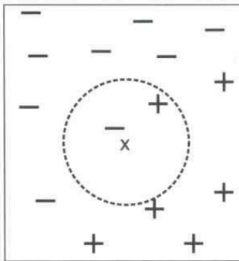
$$2.3 \quad y' = \arg \max_v \sum_{(\mathbf{x}_i, y_i) \in D_z} w_i \mathbf{1}_{v=y_i}.$$

För regression används medelvärde alternativt viktat medelvärde.

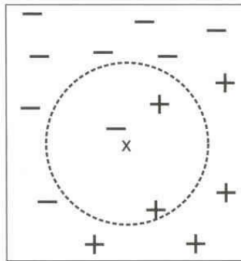
# K-närmaste grannar



(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

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

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



# Kategoriska attribut

$\mathbb{P}(Y = y)$  är andelen datapunkter med klass  $y$ .

$\mathbb{P}(X_i = x_i \mid Y = y)$  andelen datapunkter med attribut  $x_i$  av datapunkterna med klass  $y$ .

	binary	categorical	continuous	class
Tid	Home Owner	Marital Status	Annual Income	Defaulted Borrower
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
6	No	Married	60K	No
7	Yes	Divorced	220K	No
8	No	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes

För kontinuerliga attribut finns olika tillvägagångssätt.

- Diskretisera data i olika kategorier.
  - För få intervall gör att man missar information.
  - För många intervall kan ge intervall utan observationer.
- Anta en sannolikhetsfördelning för variabeln och skatta parametrarna från träningsdatan.
  - Normalfördelningen är vanlig.
  - Conjugate prior.

Träningfasen:

Skatta sannolikheten  $\mathbb{P}(Y | \mathbf{X})$  för alla möjliga  $\mathbf{X}$  och  $y$ .

Klassificeringsfas:

Givet  $\mathbf{X}'$  skatta klass  $Y' = \max_Y \mathbf{P}(Y | \mathbf{X}')$ .

Modellantagande:

$$\mathbb{P}(\mathbf{x} | Y) = \prod_i \mathbb{P}(X_i | Y),$$

alla  $X_i$  är oberoende av varandra. Vi kan då faktorisera likelihooden över  $\mathbf{x}$ .

Använder vi detta får vi en sannolikhet

$$\mathbb{P}(Y | \mathbf{x}) = \prod_i \mathbb{P}(X_i | Y) \mathbb{Y},$$

det räcker med att skatta sannolikheten för varje  $X_i$ . Detta ger oss en enklare modell som går att skatta.

- Metoden är robust mot isolerade bruspunkter.
- Metoden är robust mot irrelevanta attribut.
- Lätt att skatta.
- Korrelerade attribut kan väsentligt försämra prestandan.
  - Behöver en mer komplex modell för att hantera.
  - Simultan sannolikhetsfördelning för likelihooden.

# Ensemblemetoder

