

# T 991210 Mönstergenkammer TF4

1 a)  $\hat{m}_x = \frac{1}{N} \sum_{n=1}^N x_n$        $\sigma_x^2 = E\{(x - m_x)^2\} = E\{x^2\} - m_x^2$

$$E\{\hat{m}_x\} = \frac{1}{N} \sum_{n=1}^N E\{x_n\} = \frac{1}{N} N m_x = m_x$$

b)  $S_x^2 = E\{(\hat{m}_x - m_x)^2\} = E\{\hat{m}_x^2\} - m_x^2$

$$E\{\hat{m}_x^2\} = \frac{1}{N^2} \sum_{n=1}^N \sum_{m=1}^N E\{x_n x_m\}$$

$$= \frac{1}{N^2} \sum_{n=1}^N E\{x_n^2\} + \frac{1}{N^2} \sum_{n=1}^N \sum_{m \neq n} E\{x_n\} E\{x_m\}$$

oberende

$$= \frac{1}{N^2} \sum_{n=1}^N E\{x_n^2\} + m_x^2 \frac{N(N-1)}{N^2}$$

$$= \frac{1}{N^2} \sum_{n=1}^N \underbrace{E\{x_n^2\}}_{\sigma_x^2 + m_x^2} + m_x^2 - \frac{m_x^2}{N}$$

$$= \frac{\sigma_x^2}{N} + \frac{m_x^2}{N} + m_x^2 - \frac{m_x^2}{N}$$

$$S_x^2 = \frac{\sigma_x^2}{N} + m_x^2 - m_x^2 = \frac{\sigma_x^2}{N} \quad \text{Q.E.D.}$$

c)  $S_x = 0.01 = \frac{\sigma_x}{\sqrt{N}}$

$$N = \frac{\sigma_x^2}{S_x^2}$$

$$\Rightarrow N = \frac{\sigma_x^2}{0.01^2}$$

$$N = 100$$

$$N = 10000 \sigma_x^2$$



3. a) Outliers förkisa: Påverkan <sup>skattning av</sup> medel och kovariansmatris (mer generellt parameter-skattningar).  $\Rightarrow$  Bestulsgränser kan bli helt fel. Detta demonstrerades på lab.

b) Outliers kan uppstå tex p.g.a trasiga (felaktiga) sensor eller på manuell felavläsning.

c) Man kan "filtrera" bort outliers före träning genom att tex utföra PCA och titta på data i 2D/3D och eller genom att införa villkor

$$x_{\min}^{(i)} < x_i < x_{\max}^{(i)}$$

för varje komponent.

4 a) ickeparametriska metoder: "Univerella", kan justeras (via parameter) så att de utför godtyckliga klassificeringar och funktionsapproximationer. Kallas ibland för data-drivna metoder

b) Parametriska metoder: En familj av tänkbara lösningar parametreras och en sökmethd väljer ut den bästa. Dessa metoder ej Univerelle dvs de kan inte justeras till att passa godtyckliga problem.

Eulagerperception: Parametrisk, dess beslutgränser består av hyperplan bestämt av parameter.

Flerlagerperception: Icke-parametrisk, dess beslutgränser kan bli godtyckligt komplicerade genom att öka antalet gömda noder

b) K-NN: Snabb "träning"      Långsam respons  
Konceptuellt enkel att förstå      Mycket minne krävs

Flerlagerperception: Snabb respons      Långsam träning  
Svår att tolka  
Lokala minimer  
Likt minne krävs

5

$$P(\text{klass } i | \bar{x}) = \frac{f(\bar{x} | \text{klass } i) P(\text{klass } i)}{f(\bar{x})} \quad \text{Bayes}$$

$$P(\text{klass } i | \bar{x}) > P(\text{klass } j | \bar{x}) = g_j(\bar{x}) \quad ?$$

$$g_j(\bar{x}) = f(\bar{x} | \text{klass } i) P(\text{klass } i)$$

Beskrivning:

$$\{ \bar{x} \mid g_1(\bar{x}) = g_2(\bar{x}) \}$$

$$\left\{ \bar{x} \mid \frac{e^{-\frac{1}{2}(\bar{x}-\bar{m}_1)^T C_1^{-1}(\bar{x}-\bar{m}_1)} P(\text{klass } 1)}{(2\pi)^{d/2} |C_1|^{1/2}} = \frac{e^{-\frac{1}{2}(\bar{x}-\bar{m}_2)^T C_2^{-1}(\bar{x}-\bar{m}_2)}}{(2\pi)^{d/2} |C_2|^{1/2}} \right\}$$

$$\left\{ \bar{x} \mid -\frac{1}{2}(\bar{x}-\bar{m}_1)^T C_1^{-1}(\bar{x}-\bar{m}_1) + \frac{1}{2}(\bar{x}-\bar{m}_2)^T C_2^{-1}(\bar{x}-\bar{m}_2) = \frac{1}{2} \log |C_1| - \frac{1}{2} \log |C_2| \right\}$$

Detta är en kvadratisk form

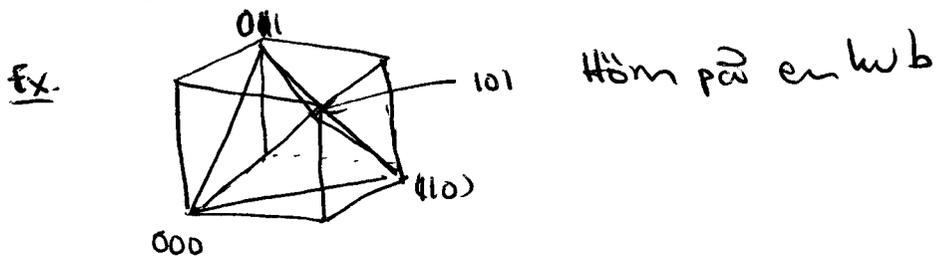
$$\{ \bar{x} \mid \bar{x}^T A \bar{x} + b^T \bar{x} + c = 0 \}$$

välkänd i linjär algebra. Om kvadratiske ytor erhålles beroende på hur A, b, c ser ut.

$C_1 = C_2 \Rightarrow$  kvadratiske termer tar ut varandra  
 $\{ \bar{x} \mid b^T \bar{x} + \tilde{c} = 0 \}$  Hyperplan

6 önsker samma avstånd  
 mellan alla bilstäder

Ledning:



Samma avstånd ( $\sqrt{2}$ ) mellan  
 alla.

Välj tex

|     |   |   |
|-----|---|---|
| 000 | ↔ | A |
| 011 | ↔ | C |
| 101 | ↔ | G |
| 110 | ↔ | T |

Anm.

Man kan tänka sig att flytta denna tetraeder  
 så att den har sina tyngdpunkter i origo men  
 detta är inte nödvändigt. Erhålls direkt  
 med kodningen

|      |   |   |
|------|---|---|
| -11  | ↔ | A |
| -1+1 | ↔ | C |
| +1-1 | ↔ | G |
| +1+1 | ↔ | T |

7

- a) K-MN: 1. För ett givet  $\bar{x}$ , bestäm avstånd till de  $K$  närmaste grannarna i träningsmängden
2. Bestäm antalet medlemmar  $K_c$  från klass  $c$  bland de  $K$  närmaste grannarna
3. Välj klass  $c$  om  $K_c > K_i$   $c \neq i$ .

- b) Korsvalderthy: 1. Dela upp tillgängliga data i tränings och testset på flera olika sätt, exempelvis är att placera ut endast ett testexempel varje gång
2. Testa prestanda m.h.a testmängderna för olika antal  $K$ . Upprepa
3. Beräkna medelprestanda för alla olika antal  $K$  över de testmängder som studerats
4. Välj det  $K$  som ger högst medelprestanda

- c)  $K$  stor  $\Rightarrow$  Liten varians i  $\frac{K_c}{N_c}$
- $\Rightarrow$  Bra estimat men inte lokalt;
- $\Rightarrow$  Globalt eftersom alla punkter i rummet med i skattning
- $K$  liten  $\Rightarrow$  Lokalt estimat men stor varians

8 a) K-NN: Alltör långsam respons  
Måste utföra många beräkningar  
Vd alla jämförelser

b) Flerlagerperceptron: Alltör långsam att  
träna.

c) RBF: Träna gånnde lagret med K-means  
Snabbt och enkelt

Linjär regresson: Mycket snabb träning och  
respons. Gäller även RBF  
där gånnde lagret fixerat.

Linjär regresson absolut snabbast men  
om sambandet olinjärt sä ger (troligen)  
RBF bättre resultat.

9) a) olinearity

$$\varphi(x) = \frac{1}{1+e^{-x}}$$

$$\varphi(x) = \tanh(x)$$

$$\varphi(x) = \begin{cases} 1 & x > 0 \\ 0 & x \leq 0 \end{cases}$$

$$\varphi(x) = e^{-x^2}$$

b) Blas: Gör det möjligt att få beslutningen/ytan som går utenför origo.