



Repetition DMI, m.m.

- I. Terminologi och Grundproblem
- II. Linjär algebra
- III. Optimering
- IV. Sannolikhetslära
- V. Parameterestimering

2003-08-22

Signaler & System
Uppsala universitet

1



Några begrepp

- Mönstervektor (egenskapsvektor/indata)
 - lista med numeriska värden som beskriver mönstret. Brukar heta \mathbf{x} .

$$\mathbf{x} = \begin{pmatrix} \text{egenskap 1} \\ \text{egenskap 2} \\ \vdots \\ \text{egenskap } d \end{pmatrix}$$

varje egenskap ses som en koordinat i ett d -dimensionellt rum

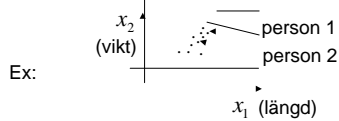
2003-08-22

Signaler & System
Uppsala universitet

2



- Börvärde. Ofta vektor (eng. target vector, el. desired output)
 - den vektor som hör ihop med \mathbf{x} .
 - brukar heta \mathbf{t} , \mathbf{y} eller \mathbf{d} .
- Mönsterrum
 - vi betraktar ett mönster, \mathbf{x} , som en punkt i ett vektorrum med dimension d .



2003-08-22

Signaler & System
Uppsala universitet

3



Grundproblem i DMI

- Klassificering (förk. "klassning")
 - synonym: mönsterigenkänning
 - nästan synonym: detektion
- Regression
 - nästan synonym: funktionsapproximation
- Kodning
 - effektiv representation av data

2003-08-22

Signaler & System
Uppsala universitet

4



För detta krävs ...

- Mängd med träningsexempel (träning mängd)
 - synonymer: designexempel, träningsdata
 - antingen (a) $X = \{\mathbf{x}_1, \mathbf{t}_1, \dots, \mathbf{x}_N, \mathbf{t}_N\}$
 - eller (b) $X = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$
- Inläring
 - synonymer: träning, parameterestimering
 - kan var övervakad (supervised, fall (a)) eller oövervakad (unsupervised, fall (b))

2003-08-22

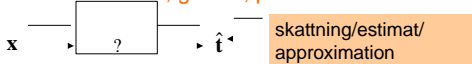
Signaler & System
Uppsala universitet

5



Övervakad inläring

- Träning data $X = \{\mathbf{x}_1, \mathbf{t}_1, \dots, \mathbf{x}_N, \mathbf{t}_N\}$
 - Målet är att, baserat på träningsdata, skapa en funktion som, givet \mathbf{x} , predikterar \mathbf{t} .



Vanlig vid klassificering att låta \mathbf{t} vara en s.k. indikatorfunktion

$$\mathbf{t}_n = \begin{pmatrix} 0 \\ \vdots \\ 1 \\ 0 \\ \vdots \end{pmatrix} \begin{array}{l} 1 \text{ på pos. } i \text{ om} \\ \mathbf{x}_n \in C_i \\ 0 \text{ annars} \end{array}$$

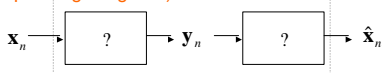
2003-08-22

Signaler & System
Uppsala universitet

6

Oövervakad inlärning

- **Träningsdata** $X = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$
 - Mål: finna (oftast kompakt) representation av \mathbf{x} , dvs *koda* eller *komprimera* i syfte att:
 - Förenkla lösningen av problem som kan beskrivas som övervakad inlärning.
 - Visualisera (Koden \mathbf{y} har typiskt max 3 dimensioner. Kan då plotta mönstren som punkter i spridningsdiagram.)

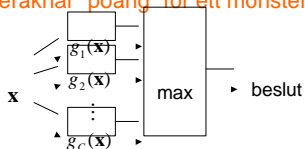


\mathbf{y}_n har lägre dimension än $\hat{\mathbf{x}}_n$

- **Exempel på oövervakad inlärning**
 - Via principalkomponentanalys (PCA). Koden är kontinuerlig.
 - Olika former av klustring. Koden är diskret (heltal)
 - Hierarkisk klustring
 - K-means clustering (KMC)
 - Kombinationer av PCA och klustring

Flera begrepp

- **Diskriminantfunktioner.**
 - förekommer vid klassificering. En funktion per klass (undantag ev. för 2-klassproblem)
 - Beräknar "poäng" för ett mönster, \mathbf{x} .





- **Beslutsområden**
 - ett område som associeras till en viss diskriminantfunktion (klass)
 - beslutsområdet för klass C_i är det område i mönsterrummet där $g_i(\mathbf{x}) > g_j(\mathbf{x})$ för alla $j \neq i$
- **Beslutsgräns/beslutsyta**
 - gräns mellan två beslutsområden
 - punkter på beslutsgränsen definieras av ekv:

$$g_i(\mathbf{x}) = g_j(\mathbf{x})$$

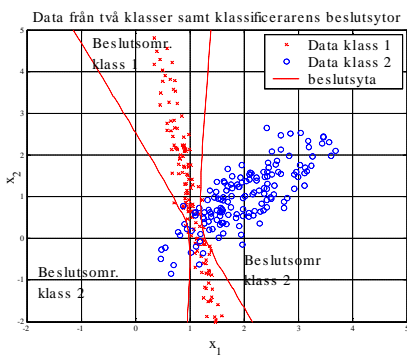
2003-08-22

Signaler & System
Uppsala universitet

10



Ex:



2003-08-22

Signaler & System
Uppsala universitet

11



II. Linjär algebra

- Hur man "läser" matris/vektor-uttryck lättare

Definitioner: $\mathbf{y} = \begin{pmatrix} y_1 \\ \vdots \\ y_m \end{pmatrix}$, $\mathbf{x} = \begin{pmatrix} x_1 \\ \vdots \\ x_d \end{pmatrix}$ och $\mathbf{W} = \begin{pmatrix} w_{11} & \cdots & w_{1d} \\ \vdots & & \vdots \\ w_{m1} & \cdots & w_{md} \end{pmatrix}$


eller $\mathbf{W} = (\mathbf{w}_1 \cdots \mathbf{w}_d)$ där $\mathbf{w}_k = \begin{pmatrix} w_{1k} \\ \vdots \\ w_{mk} \end{pmatrix}$

Dessutom $\mathbf{U} = \begin{pmatrix} \mathbf{u}_1^T \\ \vdots \\ \mathbf{u}_m^T \end{pmatrix}$ radvektorer av dimension d

2003-08-22

Signaler & System
Uppsala universitet

12


Hur man "läser" ...


(a) **Matris/vektorprodukt sett som en viktad summa av kolumnvektorer**

$$\mathbf{y} = \mathbf{W}\mathbf{x} = \sum_{k=1}^d \mathbf{w}_k x_k$$

(b) **Matris/vektorprodukt sett som en serie skalärprodukter**

$$\mathbf{y} = \mathbf{U}\mathbf{x} = \begin{pmatrix} \mathbf{u}_1^T \mathbf{x} \\ \vdots \\ \mathbf{u}_m^T \mathbf{x} \end{pmatrix}$$

2003-08-22 Signaler & System Uppsala universitet 13


Hur man "läser" ...

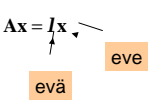
– (c) **Matris/matrixprodukt sett som en serie matrix/vektorprodukter**

Bilda $\mathbf{Y} = (\mathbf{y}_1 \ \dots \ \mathbf{y}_N)$, $\mathbf{X} = (\mathbf{x}_1 \ \dots \ \mathbf{x}_N)$

$$\Rightarrow \mathbf{Y} = \mathbf{W}\mathbf{X} = (\mathbf{W}\mathbf{x}_1 \ \dots \ \mathbf{W}\mathbf{x}_N)$$


Eigenvärde (evä) och egenvektorer (eve)

– Def. $\mathbf{A}\mathbf{x} = \lambda\mathbf{x}$



A är kvadratisk matris

2003-08-22 Signaler & System Uppsala universitet 14


Användbart "verktyg": spåret (trace) av en matris

– def: $\text{tr}(\mathbf{A}) = \sum_k \mathbf{A}_{k,k}$, dvs summan av diagonalelementen

– egenskaper (se utdelat papper)

– Erbjuder ofta ett kompakt skrivsätt för uttryck som innehåller summor

ex: Kriteriefunktion vid linjär regression (OLS)

$$J(\mathbf{W}) = \sum_n \| \mathbf{y}_n - \hat{\mathbf{y}}_n \|^2 = \sum_n (\mathbf{y}_n - \mathbf{W}\mathbf{x}_n)^T (\mathbf{y}_n - \mathbf{W}\mathbf{x}_n)$$

kan skrivas $J(\mathbf{W}) = \text{tr}((\mathbf{Y} - \mathbf{W}\mathbf{X})^T (\mathbf{Y} - \mathbf{W}\mathbf{X}))$

med $\mathbf{Y} = (\mathbf{y}_1 \ \dots \ \mathbf{y}_N)$, $\mathbf{X} = (\mathbf{x}_1 \ \dots \ \mathbf{x}_N)$

2003-08-22 Signaler & System Uppsala universitet 15



III. Optimering

- Problem: minimera $f(\mathbf{x})$ m.a.p \mathbf{x} , dvs vi söker vektorn $\mathbf{x}^* = \arg \min_{\mathbf{x}} f(\mathbf{x})$
- Lokalt minimum definieras av

$$\nabla_{\mathbf{x}} f(\mathbf{x}) = \mathbf{0}, \text{ där } \nabla_{\mathbf{x}} f(\mathbf{x}) = \begin{pmatrix} \frac{\partial f}{\partial x_1} \\ \frac{\partial f}{\partial x_2} \\ \vdots \end{pmatrix} \text{ samt att } \mathbf{x}^T \mathbf{H} \mathbf{x} > 0$$

Hessianen

2003-08-22

Signaler & System
Uppsala universitet

16



Metoder för optimering

- gradientmetoden (steepest descent)
- Newtons metod: Ide' approximera funktionen (lokalt) med en Taylorutveckling upp till kvadratiske termer. Minimera sedan denna kvadratiske funktion.
 - **Fördel:** Mkt snabb konvergens nära ett minimum
 - **Problem:**
 - Kräver att Hessianen är pos. def.
 - om dimensionen för \mathbf{x} är stor, stora beräkningar

2003-08-22

Signaler & System
Uppsala universitet

17



Optimering under bivillkor

- söker $\mathbf{x}^* = \arg \min_{\mathbf{x}} f(\mathbf{x})$, b.v. $c_i(\mathbf{x}) = 0$, för $i = 1, \dots, K$
- Enligt Lagrange: Om \mathbf{x}^* är ett lok. min. till $f(\cdot)$ (och uppfyller bivillkoren) så existerar

L_1, \dots, L_K sådana att \mathbf{x}^* och L_1, \dots, L_K uppfyller ekv.

$$\nabla_{\mathbf{x}} L(\mathbf{x}, L_1, \dots, L_K) = \mathbf{0}$$

$$c_i(\mathbf{x}) = 0, L_i \geq 0 \text{ för alla } i$$

$$L(\mathbf{x}, L_1, \dots, L_K) = f(\mathbf{x}) - \sum_i L_i c_i(\mathbf{x})$$

2003-08-22

Signaler & System
Uppsala universitet

18

IV. Sannolikhetslära

- **sannolikhet** associeras med utsagor, tex
 - A="det regnar i Skövde"
 - I=vår bakgrundsinformation
 - $P(A|I)$ utläses: sannolikhet att utsaga A är sann, givet att I är sann.
 - Egenskap: $0 \leq P(A) \leq 1$
- **Räknerregler för sannolikheter**
 - Produktregeln: $P(A, B|I) = P(A|B, I)P(B|I)$

"givet att"

"och"

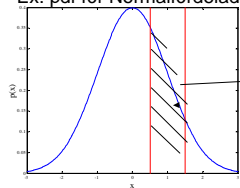
• Räknerregler, forts

- Summaregeln: $P(A \text{ eller } B|I) = P(A|I) + P(B|I) - P(A, B|I)$
- **Konsekvens av produktregeln: Bayes sats**
 - Gäller att $P(A, B|I) = P(B, A|I)$
 - $\Rightarrow P(A|B, I)P(B|I) = P(B|A, I)P(A|I)$
 - $\Rightarrow P(B|A, I) = \frac{P(A|B, I)P(B|I)}{P(A|I)}$
- **Specialfall av summaregeln:**
 - Om A och B är *uteslutande* (kan ej ske samtidigt)
 - $P(A \text{ eller } B|I) = P(A|I) + P(B|I)$, ty $P(A, B|I) = 0$


• Täthetsfördelning (pdf) för skalär, x

- Betecknas $p(x)$. Egenskap $p(x) \geq 0$. En funktion sådan att $P(a \leq x \leq b) = \int_a^b p(x) dx$

Ex: pdf för Normalfördelade x



Arean = $P(0.5 < x < 1.5)$




- **Täthetsfördelning (pdf) för vektor, $\mathbf{x} = (x_1, \dots, x_d)^T$**
 - Betecknas **$p(\mathbf{x})$** . Egenskap: $p(\mathbf{x}) \geq 0$.
 - “betydelse”: $p(\mathbf{x}) = p(x_1, x_2, \dots, x_d)$.
 - Produktregeln gäller, dvs

$$p(x_1, x_2, \dots, x_d) = p(x_1 | x_2, \dots, x_d) p(x_2, \dots, x_d) \quad \text{osv.}$$
 - Blandning av diskreta och kontinuerliga variabler möjligt, tex

$$p(\mathbf{x}, \text{klasshörighet } C_i) = p(\mathbf{x}, C_i) = p(\mathbf{x} | C_i) P(C_i)$$

\uparrow
förkortat skrivsätt

2003-08-22 Signaler & System Uppsala universitet 22




- **Tillämpning av shtslära inom MI**
 - Anta att \mathbf{x} “dras” från olika klasser C_1, \dots, C_c med *a priori-sannolikheter* $P(C_1), \dots, P(C_c)$
 - Olika pdf:er, $p(\mathbf{x} | C_1), \dots, p(\mathbf{x} | C_c)$, för varje klass
 - Minimering av felsannolikhet leder till den optimala beslutsregeln:

välj C_i om $P(C_i | \mathbf{x}) \geq P(C_j | \mathbf{x})$ för alla $j \neq i$
 - *a-posteriorisannolikheterna* $P(C_i | \mathbf{x})$ fås via **Bayes sats**

$$P(C_i | \mathbf{x}) = \frac{p(\mathbf{x} | C_i) P(C_i)}{p(\mathbf{x})}$$

\swarrow
måste vanl. skattas

2003-08-22 Signaler & System Uppsala universitet 23



Några viktiga begrepp:

- **Väntevärde**

$$\mathbf{m} = E[\mathbf{x}] = \int_{\mathbf{x}} \mathbf{x} p(\mathbf{x}) d\mathbf{x} = \int \int \dots \int \mathbf{x} p(\mathbf{x}) dx_d \dots dx_2 dx_1$$
 - kan tolkas som en *tyngdpunkt* för en massfördelning, $p(\mathbf{x})$.
- **Kovariansmatris**

$$\mathbf{C} = E[(\mathbf{x} - \mathbf{m})(\mathbf{x} - \mathbf{m})'] = \int_{\mathbf{x}} (\mathbf{x} - \mathbf{m})(\mathbf{x} - \mathbf{m})' p(\mathbf{x}) d\mathbf{x}$$
 - ger information om fördelningens spridning
 - analyseras m.h.a principalkomponentanalys

2003-08-22 Signaler & System Uppsala universitet 24

• **Betingad pdf, $p(y|x)$**
 – betecknas $p(y|x)$
 – täthetsfördelningen för y då x är låst.

• **Betingat väntevärde**

$$m_{y|x} = E[y|x] = \int_y y p(y|x) dy$$

Tyngdpunkten längs denna linje

Ex: y och x är skalärer.
 De har simultana pdf:en $p(x,y)$

2003-08-22 Signaler & System Uppsala universitet 25

• **Vanligt ex på vektorvärd pdf, normalfördelningen:**

$$p(\mathbf{x}) = \frac{1}{(2\pi)^{d/2} |\mathbf{C}|^{1/2}} e^{-\frac{1}{2}(\mathbf{x}-\mathbf{m})^T \mathbf{C}^{-1}(\mathbf{x}-\mathbf{m})}$$

kovariansmatris väntevärde

d=dimension determinant

2003-08-22 Signaler & System Uppsala universitet 26

V: Parameterskattning

• **Problem:** Vi har data $X = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$ som antas vara "dragna" från pdf $p(\mathbf{x})$ som har känd parametrisk form. Parametrarnas värden är dock okända.

– ex: Vi antar att \mathbf{x} är normalfördelad med väntevärde \mathbf{m} och kovariansmatris \mathbf{C} .

• **Två (besläktade) principer för parameterskattning:**

– Maximum likelihood (ML)
 – Maximum a posteriori (MAP)

2003-08-22 Signaler & System Uppsala universitet 27



Låt \mathbf{q} beteckna en vektor som innehåller alla okända parametrar

- ML: $\hat{\mathbf{q}}_{ML} = \arg \max p(X | \mathbf{q})$
s.k. likelihoodfunktion
- MAP: $\hat{\mathbf{q}}_{MAP} = \arg \max p(\mathbf{q} | X) = \arg \max p(X | \mathbf{q})p(\mathbf{q})$
s.k. "prior"

– Obs: om $p(\mathbf{q})$ är relativt konstant för \mathbf{q} i ett stort område så sammanfaller ML- och MAP-skattningarna

2003-08-22

Signaler & System
Uppsala universitet

28
