

Tre delar:

1. Maskininlärning:

- grundläggande begrepp.
- ingredienser (optimering, validering).
- neurala nätverk.
- etiskt perspektiv.

2. Bayesiansk inferens:

- formellt perspektiv på inlärning.
- tolkning av sannolikhet.
- Bayesiansk linjär regression.

3. Stokastiska processer:

- Markovkedjor.
- MCMC (Markov Chain Monte Carlo) för Bayesiansk inferens.

Vetenskaplig modellering:

$$y_j = y_j(x_j) \quad \begin{cases} y_j & \text{beroende (respons-) variabel/-ler.} \\ x_j & \text{oberoende (prediktor-) variabel/-ler.} \end{cases}$$

Introducera modell M (approximation).

$$y_j \approx M(\Theta; x_j), \quad \Theta \text{ modellparameter/-rar.} \quad (\text{reella})$$

y_j, x_j, Θ kan vara olika typer av tal. Oftast betraktas $\bullet \in \mathbb{R}^p$.

Kategorisering:

- Deterministisk modell: (input + fixa parametrar \rightarrow output)
 - Explicit: Välddefinierade matematiska funktioner.
 - Linjär.
 - Icke-linjär.
 - Implicit: T.ex. lösning av different. eller egenvärdesproblem.
- Stokastisk modell: innehåller fluktuationer etc.

Vetenskapliga modeller är oftast parametriska.

T.ex. $x(t) = A \sin(\omega t + \phi)$.

- parametrarna har betydelse i sig själva.
- bygger oftast på en fysikalisk insikt/hypotes.

Alt. konstruera en flexibel klass av modeller och en inlärningsalgoritm för att hitta samband.
Kräver oftast en stor mängd observationer (träningsdata).
⇒ Maskininlärning (ML) (icke-parametriska modeller).

Regressionsanalys:

Att hitta samband, oftast med statistiska metoder.
Mest vanligt förekommande i ML är optimering.

Antag modell $M(\theta; x)$, N_d observationer $\{y_1, \dots, y_{N_d}\}$,
motsv. oberoende variabel $\{x_1, \dots, x_{N_d}\}$.

Tillsammans utgör detta en datamängd ①.

Definiera kostnadsfunktion $C(\theta)$ som "mäter" hur väl modellen beskriver data.

T.ex. $C(\theta) = \frac{1}{N_d} \sum_{i=1}^{N_d} \frac{[y_i - M(\theta; x_i)]^2}{\sigma_i^2}$, σ_i skalfaktor s.a C dim.löst.

Hitta $\theta^* = \underset{\theta \in \mathbb{R}^p}{\operatorname{argmin}} C(\theta)$.

Detta är ett optimeringsproblem.

De vanligaste optimeringsmetoderna bygger på gradientstegsoptimering (gradient descent).

Betrakta $C(\theta)$ som en flervariabelfunktion.

Grundalgoritm:

1. Starta från slumpmässigt vald punkt θ_0 .
2. Vid θ_n :
 - evaluera gradienten $\nabla_{\theta} C(\theta)|_{\theta=\theta_n} \equiv \nabla C_n$.
 - välj "inlärningshastighet" η_n
3. $\theta_{n+1} = \theta_n - \eta_n \nabla C_n$ för fixt antal iterationer eller tills gradienten i θ_k är liten.

Linjära modeller:

$$y \in \mathbb{R}, x \in \mathbb{R}, \Theta \in \mathbb{R}^p.$$

Linjära modeller innebär linjärt parameterberoende.

$$M(\Theta; x) = \sum_{j=0}^{p-1} \Theta_j \phi_j(x), \quad \phi_j(x) \text{ beror ej på } \Theta.$$

Ex.

• polynommodell: $M(\Theta; x) = \sum_{j=0}^{p-1} \Theta_j x^j = \Theta_0 + \Theta_1 x + \Theta_2 x^2 + \dots + \Theta_{p-1} x^{p-1}.$

• trigonometrisk modell: $M(\Theta; x) = A_0 + \sum_{n=1}^p (A_n \cos nx + B_n \sin nx),$

$$\Theta = \{A_0, \dots, A_p, B_0, \dots, B_p\}.$$

Θ_0 kallas intercept (skärningspunkt) eller bias.

Regressionsanalys med linjära modeller:

Data $\mathbb{D} = [y_1, \dots, y_{N_d}]^T$ (kolumnvektor).

Mot varje y_i finns en oberoende variabel x_i .

$$y_i \approx M(\Theta; x_i) = \sum_{j=0}^{p-1} \Theta_j \phi_j(x_i).$$

Introducera designmatrisen $X = \begin{bmatrix} \phi_0(x_1) & \dots & \phi_{p-1}(x_1) \\ \vdots & & \vdots \\ \phi_0(x_{N_d}) & \dots & \phi_{p-1}(x_{N_d}) \end{bmatrix}$ ($N_d \times p$)

med parametervektorn $\Theta = [\Theta_0, \dots, \Theta_{p-1}]^T$

$$\Rightarrow \mathbb{D} = \begin{bmatrix} y_1 \\ \vdots \\ y_{N_d} \end{bmatrix} \approx X\Theta.$$

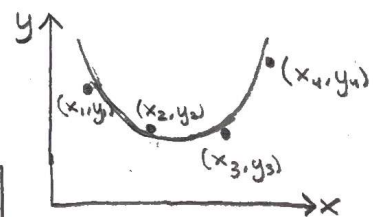
Introducera residualer $\epsilon_i = y_i - M_i$ ($M_i = M(\Theta; x_i)$)

$$\Rightarrow \mathbb{D} = X\Theta + \mathbb{E}, \quad \mathbb{E} = [\epsilon_1, \dots, \epsilon_{N_d}]^T.$$

Ex.

$$M(\Theta; x) = \Theta_0 + \Theta_1 x + \Theta_2 x^2, \quad \mathbb{D} = [y_1, y_2, y_3, y_4]^T.$$

$$X = \begin{bmatrix} 1 & x_1 & x_1^2 \\ 1 & x_2 & x_2^2 \\ 1 & x_3 & x_3^2 \\ 1 & x_4 & x_4^2 \end{bmatrix} \Rightarrow \mathbb{D} = \begin{bmatrix} y_1 \\ y_2 \\ y_3 \\ y_4 \end{bmatrix} = \begin{bmatrix} 1 & x_1 & x_1^2 \\ 1 & x_2 & x_2^2 \\ 1 & x_3 & x_3^2 \\ 1 & x_4 & x_4^2 \end{bmatrix} \begin{bmatrix} \Theta_0 \\ \Theta_1 \\ \Theta_2 \end{bmatrix} + \begin{bmatrix} \epsilon_1 \\ \epsilon_2 \\ \epsilon_3 \\ \epsilon_4 \end{bmatrix}$$



Hitta Θ som minimerar $|\mathbb{E}|^2$. (vanligaste sättet att minimera \mathbb{E})

Formulera som kostnadsfunktion $C(\Theta) = \sum_{i=1}^{N_d} \epsilon_i^2 = \sum_{i=1}^{N_d} (y_i - M_i)^2 = |\epsilon|^2$
 $= \{M_i = (\mathbb{X}\Theta)_i\} = (\mathbb{I} - \mathbb{X}\Theta)^T (\mathbb{I} - \mathbb{X}\Theta) = \epsilon^T \epsilon$.

Söker $\Theta^* = \underset{\Theta \in \mathbb{R}^p}{\operatorname{argmin}} (\mathbb{I} - \mathbb{X}\Theta)^T (\mathbb{I} - \mathbb{X}\Theta)$ (optimeringsproblem).

$C(\Theta)$ konvex funktion av $\Theta \Rightarrow$ en extrempunkt (minimipunkt).

$$\nabla_{\Theta} C(\Theta) = \begin{bmatrix} \frac{\partial C}{\partial \theta_0} \\ \vdots \\ \frac{\partial C}{\partial \theta_{p-1}} \end{bmatrix} = \left\{ \frac{\partial C}{\partial \theta_j} = \frac{\partial}{\partial \theta_j} \left(\sum_{i=1}^{N_d} (y_i - \theta_0 \phi_0(x_i) - \dots - \theta_{p-1} \phi_{p-1}(x_i))^2 \right) \right\} =$$

$$= -2 \sum_{i=1}^{N_d} \phi_j(x_i) (y_i - \theta_0 \phi_0(x_i) - \dots - \theta_{p-1} \phi_{p-1}(x_i)) = -2 \mathbb{X}^T (\mathbb{I} - \mathbb{X}\Theta).$$

Söker Θ^* s.a. $\nabla_{\Theta} C(\Theta)|_{\Theta=\Theta^*} = 0$.

$$\Rightarrow -2 \mathbb{X}^T \mathbb{I} + 2 \mathbb{X}^T \mathbb{X} \Theta^* = 0 \Rightarrow \mathbb{X}^T \mathbb{I} = \mathbb{X}^T \mathbb{X} \Theta^* \text{ (normalekvationen).}$$

Hitta $(\mathbb{X}^T \mathbb{X})^{-1}$ s.a. $\underline{\underline{\Theta^* = (\mathbb{X}^T \mathbb{X})^{-1} \mathbb{X}^T \mathbb{I}}}$.

$\mathbb{X}^T (p \times N_d)$, $\mathbb{X} (N_d \times p)$, $\mathbb{I} (N_d \times 1) \Rightarrow \Theta^* (p \times 1)$ (dim. ok).

Regressionsresidualerna $\epsilon^* = \mathbb{I} - \mathbb{X}\Theta^* = 0$ vid "perfekt" anpassad regression (osannolikt).

Linjär regression inom fysik:

- Basfunktionerna $\phi_j(x)$ är fysikaliskt motiverade.
- y och x är fysikaliska storheter (med enheter).
- Designar och utför experiment för att samla in data \mathbb{I} .
- Modellparametrarna är intressanta i sig.
- Kan ha kunskaper om parametrarna från andra fysikaliska samband.
- Kan ofta estimeras storlek på fel.

Linjär regression inom ML:

- En form av övervakad (supervised) maskininlärning.
- Ofta kombineras oberoende variabler och basfunktion till prediktor. T.ex. $x, x^2, \frac{1}{x}, 1 \Rightarrow \mathbb{X} = [1, x, x^2, x^{-1/2}]$.

- Oftast konstateras en lämplig kostnadsfunktion och man prövar sig fram med olika prediktorer.
- Parametrarna är oftast inte intressanta i sig själva. Ibland skrivs modellen $\hat{y}(x)$ (beror ej på parametrar Θ).
- Begränsad möjlighet att påverka mängden träningsdata.
- Svårt att estimeras "generaliseringsfel" (d.v.s. då modellen används för förutsägelser).

Maskininlärning:

Fokus på övervakad (supervised) inlärning.

Tre ingredienser:

1. Data.
2. Matematisk modell.
3. Inlärningsalgoritm.

Kategorisering av ML efter typ av uppgift:

- i) Regressionsanalys: -- hitta samband (funktionsliknande) mellan prediktor- och responsvariabler.
- ii) Klassificering: - kategorisera.
- tränas med kategoriserad data.
- iii) Klusteralgoritmer: - hitta mönster i data.
- dela in i olika kluster.
- kan ibland utvärderas utan fördefinierade kategorier.
- iv) Dimensionsreduktion: - identifiera de viktigaste prediktorerna.

Fokus på i) och ii).

(Numeriskt exempel på linjär regression i kap. 8.)

Klassificering:Binär klassificering:

(0 eller 1)

Ex. Två prediktorer $\mathbf{x} = [x_1, x_2]$, respons $y \in \{0, 1\}$,

träningsdata $\mathcal{T} = \{\mathbf{x}_i, y_i\}_{i=1}^{N_d}$

Skapa en "hård" klassificerare (hard classifier) genom att göra en linjär modell

$$z = w_0 + w_1 x_1 + w_2 x_2, \quad \mathbf{w} = [w_0, w_1, w_2]^T \text{ parametrar, } z \in \mathbb{R}.$$

Projicera på $\{0, 1\}$ genom "perceptron"

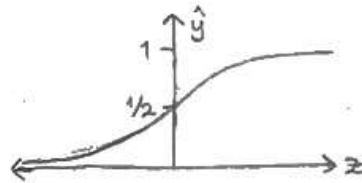
$$\hat{y} = \hat{y}(z) = \frac{\text{sign}(z) + 1}{2} = \begin{cases} 1 & \text{om } \text{sign}(z) = 1 \text{ (} z > 0 \text{)} \\ 0 & \text{om } \text{sign}(z) = -1 \text{ (} z < 0 \text{)} \end{cases}$$

Dvs. $\hat{y}_i = \hat{y}(z_i)$, $z_i = w_0 + w_1 x_{1,i} + w_2 x_{2,i}$, eller

$$\hat{y}_j = \hat{y}(z), \quad z = Xw, \quad X = \begin{bmatrix} 1 & x_{1,1} & x_{2,1} \\ \vdots & \vdots & \vdots \\ 1 & x_{1,N_d} & x_{2,N_d} \end{bmatrix}.$$

Skapa en "mjuk" klassificerare genom sigmoidfunktionen

$$\hat{y} = \hat{y}(z) = \frac{e^z}{1+e^z} = \frac{1}{1+e^{-z}}.$$

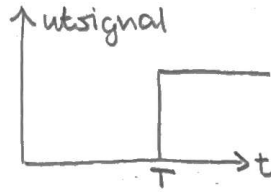
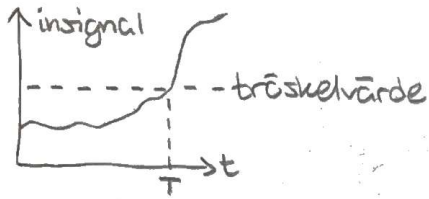


Tolka \hat{y} som sannolikhet:

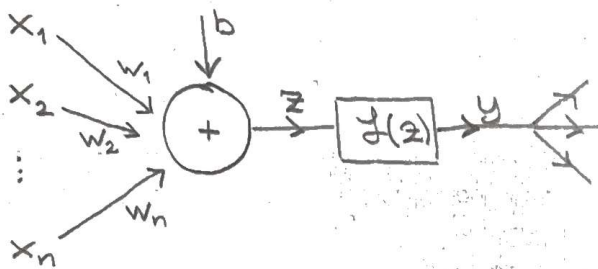
$$\hat{y} = P(y(x)=1), \quad 1-\hat{y} = P(y(x)=0).$$

Artificiella neurala nätverk:

Byggstenarna för neurala nätverk kallas neuroner.



Matematisk modell av neuron:



$$z = \sum_{j=1}^n w_j x_j + b, \quad y = f(z) \text{ aktiveringsfunktion.}$$

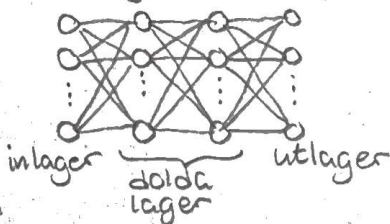
Ex.

$$f(z) = \frac{\text{sign}(z) + 1}{2}, \quad f(z) = \frac{1}{1 + e^{-z}} \text{ (sigmoid).}$$

Ett neuralt nätverk (NN) består av flera ihopkopplade neuroner (noder) som kan delas in i olika typer (beroende på uppgift).

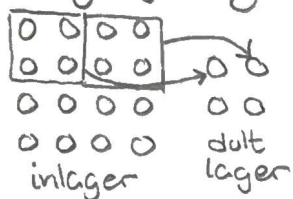
1. Framåtriktade NN (FFNN):

- Många olika användningsområden inkl. djupa NN (djupa → med dolda lager mellan in- och utsignal).
- Oftast övervakad inlärning.
- Oftast fullt anslutna → alla noder i ett lager skickar ut signaler till alla noder i nästa lager.



2. Neurala konvolutionsnätverk (CNN):

- Ej fullt anslutna.
- Utnyttjar geometrisk struktur.



- Används framförallt vid bildigenkänning.

3. Återkopplade NN (RNN):

- Ej enbart framåtriktade.
- Sekventiell information (ordning spelar roll).
- Används framförallt vid text och tal.

4. Autoencoders:

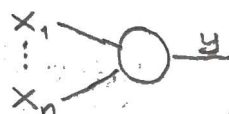
- För övervakad inlärning (eller Boltzmann machines).

FFNN definieras av sin:

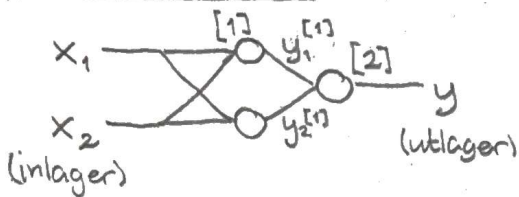
- Arkitektur.
- Aktiveringsfunktion(er).
- Inlärningsalgoritm.

Arkitektur:

- Antal lager (L) (inkl. de med aktiveringsfunktioner
→ inkl. utlagret, exkl. inlagret).
- Antal noder per lager N_l , $l=1, \dots, L$.
- Signalpropageringen beskrivs bäst med matris-vektor-operationer.
- Grafiskt förenklad bild av noder:



Litet nätverk:



$$y_1^{[1]} = f(z_1^{[1]}), \quad z_1^{[1]} = w_{11}^{[1]} x_1 + w_{12}^{[1]} x_2 + b_1^{[1]}$$

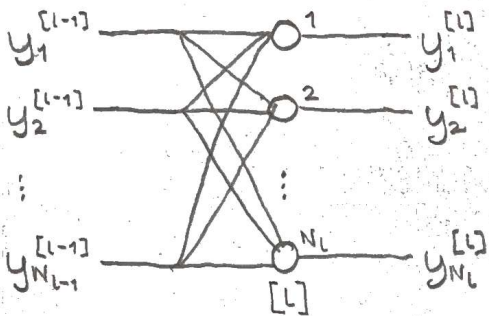
Med matris-vektor-operationer:

$$\begin{aligned} \mathbf{x} &= [x_1, x_2] \\ \mathbf{y}^{[1]} &= [y_1^{[1]}, y_2^{[1]}] \\ \mathbf{w}^{[1]} &= \begin{bmatrix} w_{11}^{[1]} & w_{21}^{[1]} \\ w_{12}^{[1]} & w_{22}^{[1]} \end{bmatrix} \\ \mathbf{b}^{[1]} &= [b_1^{[1]}, b_2^{[1]}] \end{aligned} \quad \Rightarrow \quad \begin{cases} y_j^{[1]} = f(z^{[1]}), \\ z^{[1]} = \mathbf{x} \cdot \mathbf{w}^{[1]} + \mathbf{b}^{[1]} \\ \text{Utsignal } y = y^{[2]} = f(z^{[2]}), \\ z^{[2]} = \mathbf{y}^{[1]} \cdot \mathbf{w}^{[2]} + \mathbf{b}^{[2]}, \\ \mathbf{w}^{[2]} = \begin{bmatrix} w_1^{[2]} \\ w_2^{[2]} \end{bmatrix} \end{cases}$$

Beteckning: $w_{12}^{[1]}$ ← lager
nod ↑ insignal

Djupt nätverk:

Betrakta lager L:



$$\begin{aligned} y_j^{[L]} &= f(z^{[L]}), \\ z^{[L]} &= \mathbf{y}^{[l-1]} \cdot \mathbf{w}^{[L]} + \mathbf{b}^{[L]} \end{aligned}$$

$[1 \times N_L] \quad [1 \times N_{L-1}] \quad [N_{L-1} \times N_L] \quad [1 \times N_L]$

⇒ # parametrar i lager L är $N_L(N_{L-1} + 1)$.

Ex. Hur många parametrar har ett FFNN som ska klassificera bilder (8x8 pixlar) av handskrivna siffror 0-9? Det har 10 dolda lager med 20 noder var. Utlagret ger 10 utsignaler y_j , $j=0, \dots, 9$, där $y_j = IP(t=j) \in [0, 1]$ (mjuk klassificerare).

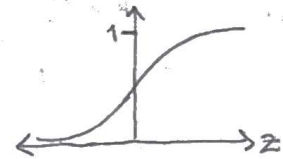
$$\left. \begin{aligned} \text{Dolt lager 1: } & 20(64+1) = 1300 \\ & 2-10: 20(20+1) = 420 \\ \text{Utlager: } & 10(20+1) = 210 \end{aligned} \right\} \Rightarrow \underline{\underline{5200 \text{ parametrar.}}}$$

Aktiveringsfunktioner:

a) Sigmoid: $f(z) = \frac{1}{1+e^{-z}} \in [0, 1]$, $f(z) = \tanh(z) \in [-1, 1]$.

• Nollskild $\forall z$.

• Små gradienter \Rightarrow Långsam inlärning.

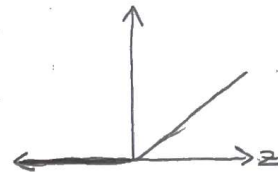


b) Modifierad linjär enhet (ReLU):

$$f(z) = \max(z, 0)$$

• Ovikliga noder är "tysta".

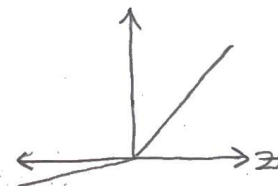
• $f'(z)$ är enkel, $f'(z) = 0$ då $z < 0$.



c) Läckande ReLU:

$$f(z) = \max(0, z) + \alpha \min(0, z), \quad \alpha \ll 1.$$

• Inga "tysta" noder.

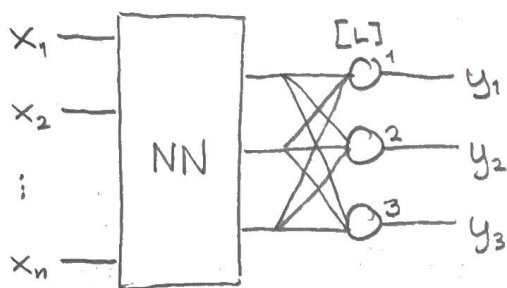


d) ELU: (exp.)

Notera att aktiveringsfunktionen i utlagret väljs efter systemet.

Ex. NN ska förutsäga position i \mathbb{R}^3 .

$$y_j = (y_1, y_2, y_3), \quad y_j \in \mathbb{R}, \quad \text{ex.v. } y_j = y_j^{[L]} = z_j^{[L]}$$



Ex. NN ska förutsäga fyra diskreta sannolikheter.

$$y_j = P(t=j), \quad \sum_{j=1}^4 y_j = 1.$$

Låt $\tilde{y}_j = y_j^{[L]} = \sigma(z_j^{[L]})$ (sigmoid) och transformera enligt

$$y_j = \frac{\tilde{y}_j}{\sum_{j=1}^4 \tilde{y}_j}.$$

Inlärningsalgoritmer:

a) Val av kostnadsfunktion:

Jämför modellens utdata med träningsdata (ex.v. medelkvadratfel (MSE)).

$C(y, \hat{y})$, y data, \hat{y} modell utdata.

b) Optimeringsalgoritmer:

Optimera vikterna för att minimera kostnadsfunktion genom någon form av gradientstegsoptimering.

$$\frac{\partial C}{\partial w_{ij}^{[l]}} \quad \forall l, i, j \quad (\text{bakåtpropagering}).$$

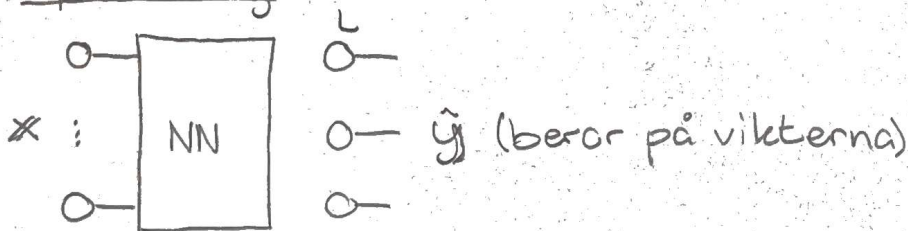
c) Uppdelning av data:

- i) träningsdata.
- ii) valideringsdata. (se modellvalidering)
- iii) testdata.

d) Inlärningsprocedur:

Hur många iterationer med gradientstegsoptimering?
Variera arkitektur, aktiveringsfunktion?

Optimering:



Kostnadsfunktion $C(y, \hat{y})$.

Justera vikterna iterativt mha grad.stegsopt:

$$w_{ij}^{[l]} \leftarrow w_{ij}^{[l]} - \eta \delta_{ij}^{[l]}, \quad \delta_{ij}^{[l]} = \frac{\partial C}{\partial w_{ij}^{[l]}}$$

(iteration) (iteration)
 $n+1$ n

Pss för bias:

$$b_i \leftarrow b_i - \eta \delta_i$$

Notera att C beror på \hat{y}_j .

$$\hat{y}_j = y_j^{[L]} = f(z^{[L]}) = f(y^{[L-1]} w^{[L]} + b^{[L]}) = f(f(z^{[L-1]} w^{[L]} + b^{[L]})) = \dots,$$

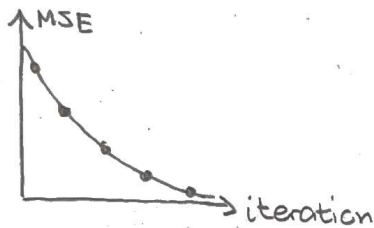
dvs alla gradienter kan beräknas med bakåtpropagering.

Inläring:

Träningsdata $x_t = \begin{bmatrix} x_{t1} \\ \vdots \\ x_{tN_t} \end{bmatrix}$, $y_t = \begin{bmatrix} y_{t1} \\ \vdots \\ y_{tN_t} \end{bmatrix}$. $x_t \rightarrow \boxed{\text{NN}} \rightarrow y_t$

Vill optimera nätverket s.a. $\text{MSE}(y_j, \hat{y}_j)$ blir så liten som möjligt.

Vikterna justeras iterativt.



(iterationer mäts ofta i epoker
(1 epok = all träningsdata har använts
en gång))

Modellvalidering:

ML-modeller har ofta stor flexibilitet \Rightarrow risk för överanpassning mot träningsdata och opålitliga prediktioner.

Använd linjär regression med polynommodell som prototypex.

$$\hat{y}(x) = \sum_{i=0}^p \theta_i x^i \quad (\text{ökande } p \Rightarrow \text{ökande komplexitet}).$$

(se ex. på Canvas: Files/demo/demo-ModelValidation-skeleton.ipynb)

Överanpassad modell:

Har anpassats till brus lika mycket som verklig respons.

Underanpassad modell:

Ignorerar variabilitet i data som faktiskt motsvarar verklig respons.

Detta kan vara svårt att upptäcka visuellt. Behöver andra metoder.

Kännetecken:

1. Stora värden på modellens parametrar vid överanpassning.
2. Valideringsdata kan användas för att detektera både över- och underanpassning.

Stora parametervärden kan undvikas genom regularisering.

Istället för att optimera $\Theta^* = \operatorname{argmin}_{\Theta \in \mathbb{R}^{p+1}} C(\Theta)$, optimeras en modifierad kostnadsfunkt., ex.v. $\Theta \in \mathbb{R}^{p+1}$

$$\Theta_{\lambda,2}^* = \operatorname{argmin}_{\Theta \in \mathbb{R}^{p+1}} C_{\lambda,2}(\Theta), \quad C_{\lambda,2}(\Theta) = C(\Theta) + \lambda \|\Theta\|_2^2.$$

Detta kallas Ridge-regularisering.

(2-norm av parametervektorn)

Lasso-regularisering: $C_{\lambda,1}(\Theta) = C(\Theta) + \lambda \|\Theta\|_1$.

Notera: $\|\Theta\|_2^2 = \sum_{i=1}^{p+1} \theta_i^2$, $\|\Theta\|_1 = \sum_{i=1}^{p+1} |\theta_i|$.

λ är en hyperparameter.

Validering:

Bygger på att dela upp tillgänglig data.

Data $\begin{cases} \rightarrow \text{Valideringsdata } (D_{\text{val}}) \\ \rightarrow \text{Träningsdata } (D_{\text{train}}) \end{cases}$

Optimera modellen mot $D_{\text{train}} \Rightarrow$ minimera $C(y_{\text{train}}, \hat{y}_{\text{train}})$ \Rightarrow
 $\Rightarrow \Theta^*$ tränad modell \Rightarrow Beräkna \hat{y}_{val} , jämför med y_{val} ,
där $D_{\text{val}} = (x_{\text{val}}, y_{\text{val}})$.

Jämförelsen sker mha en fel-funktion. Ex. för regression

$$E(y, \hat{y}) = \frac{1}{N_{\text{val}}} \sum_{i=1}^{N_{\text{val}}} (y_i - \hat{y}(x_i))^2$$

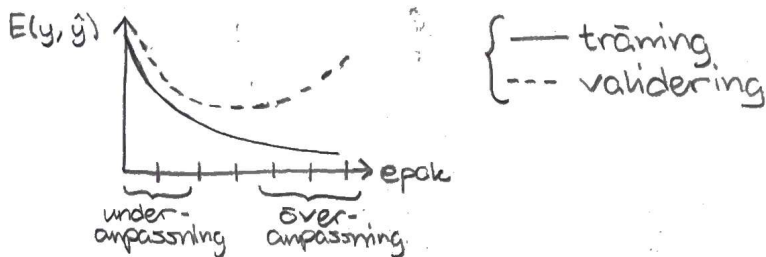
Kan vara (men måste inte) identisk med kostnads-funk.
Ex. för klassificering

$$E(y, \hat{y}) = \frac{1}{N_{\text{val}}} \sum_{i=1}^{N_{\text{val}}} I(y_i \neq \hat{y}_i), \quad I(y_i \neq \hat{y}_i) = \begin{cases} 0 & \text{om } y_i = \hat{y}_i \\ 1 & \text{om } y_i \neq \hat{y}_i \end{cases}$$

Två exempel på användning av valideringsdata:

a) Inlärningskurva:

Vid iterativ träning av ett NN ex.v.



b) Optimering av hyperparametrar:

Betrakta ex.v. ordningen p hos en polynommodell som hyperparameter.

Testa flera val $\{p_i\}$. För varje val optimeras modellen mot träningsdata och jämförs mot valideringsdata. Valet av p som minimerar felet mot valideringsdata är då det "bästa" valet.

Rättviseaspekten vid träning och validering av ML-modeller:

- Kostnadsfunktion används vid träning (optimering) av en ML-modell mot given träningsdata.
- Fel används för att validera modellen mot valideringsdata.

Huvudsakligen betraktas tekniska aspekter, men val av fel-/kostnadsfunktion tillsammans med vilken data som är tillgänglig kan påverka hur rättvis modellen är.

Betrakta binär klassificering, vilken kan utvärderas med felklassificeringsrelevansen.

$$E_{\text{val}} = \frac{1}{N_{\text{val}}} \sum_{i=1}^{N_{\text{val}}} (1 - \delta_{y_i, \hat{y}(x_i)}),$$

(måldata) (ML-modellens förutsägelse)

Notera $y_i \in \{0, 1\}$, $\hat{y}(x_i) \in \{0, 1\}$.

Sanningsmatrisen:

	$y=0$	$y=1$	Σ	$\left\{ \begin{array}{l} T - \text{true} \\ F - \text{false} \\ P - \text{positive} \\ N - \text{negative} \end{array} \right.$
$\hat{y}=0$	TN	FN	N^*	
$\hat{y}=1$	FP	TP	P^*	
Σ	N	P	N_{val}	

FP kallas ibland för typ-1 fel.
FN -||- typ-2 fel.

Beroende på sammanhang kan endera av dessa två feltyper ge allvarliga konsekvenser.

Modellens träffsäkerhet:

$$\frac{TN+TP}{N_{\text{val}}} = 1 - \text{fel.klass.frekv.} = 1 - \frac{FN+FP}{N_{\text{val}}}.$$

Andel falskt positiva: $\frac{FP}{N}$.

-||- negativa: $\frac{FN}{P}$.

Modellens precision: $\frac{TP}{(TP+FP)} = \frac{TP}{P^*}$.

Notera att det finns en risk att något/flera av dessa mått blir olika för olika grupper inom populationen.

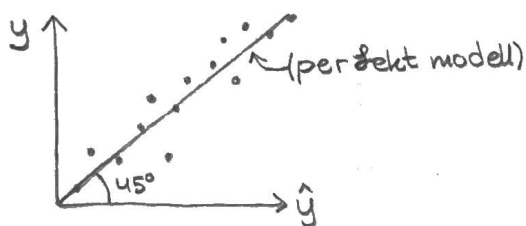
Detta kan ske pga:

- snedvriden data.
- val av kostnadsfunktion.

Svårt att upptäcka. Medvetenhet är viktigt. Presentera och utvärdera modellen utifrån olika rättvisaspekter.

Ex. Compasalgoritmen.

Ett sätt att undersöka fel för en regressionsmodell (motsv. sanningsmatis) är paritetsplot.



Begränsningar i träningsdata:

Målet är att hitta $y \approx \hat{y}(x)$.
(ML-modell som tränas på träningsdata)

Bristfällig data \Rightarrow bristfällig modell.

Fokuserar på snedvriden data (biased data). Ex.v.

• Systematisk snedvridning:

- En viss typ av data favoriseras (är vanligare) medan annan ignoreras.

• Automatisk snedvridning:

- Data samlas in genom automatiserad process, ex.v. cookies,

• Urvals-snedvriden data:

- Urvalet representerar inte populationen.
- Generellt vanligaste typen.

• Språklig snedvridning:

Det finns etiska principer och riktlinjer för pålitlig AI.

Det finns 7 krav:

1. Mänskligt agentskap och tillsyn.
2. Teknisk robusthet och säkerhet.
- ...
4. Transparens.
- ...
7. Ansvarsskyldighet.

Repetition:

Sannolikhetsmått: $IP: A \rightarrow [0, 1]$.
händelse i utfallsrum S

Betingad sannolikhet: $IP(B|A) = \frac{IP(A \cap B)}{IP(A)}$. (*)
smitt
B givet A

A, B oberoende händelser $\Rightarrow IP(B|A) = IP(B) \stackrel{(*)}{\Rightarrow} IP(A \cap B) = IP(A)IP(B)$.

Marginal sannolikhet:

B_1, \dots, B_n unika och spänner upp utfallsrummet, $IP(B_i) > 0 \forall i \in [1, n]$

$$\Rightarrow IP(A) = \sum_{i=1}^n IP(A|B_i)IP(B_i) \stackrel{(*)}{=} \sum_{i=1}^n IP(A \cap B_i)$$

Slumpvariabel: $X: S \rightarrow \mathbb{R}$.

Fokuserar på utfallen $x \in \mathbb{R}$.

Fördelningsfunktion: $P: \mathbb{R} \rightarrow [0, 1]$ som uppfyller $IP(X \leq x) = P(x)$.

(Notera att $IP(X=x) = 0$ då X kontinuerlig.)

Ibland skrivs $P_X(x)$.

Täthetsfunktion: $p: \mathbb{R} \rightarrow [0, \infty)$ s.a. $P(x) = \int_{-\infty}^x p(\tilde{x}) d\tilde{x}$.

Notera att

- $p(x) \geq 0$.
- Enheter $[P] = 1$, $[p] = [x^{-1}]$.

$$\int_{-\infty}^{\infty} p(\tilde{x}) d\tilde{x} = 1.$$

$$\bullet IP(a \leq X \leq b) = IP(b) - IP(a) = \int_a^b p(\tilde{x}) d\tilde{x}.$$

($IP(x \leq X \leq x+dx) = p(x)dx$ för infinitesimalt intervall.)

Marginal täthetsfunktion:

Betrakta den bivarata fördelningsfunktionen

$$P(x, y) = IP(X \leq x, Y \leq y) = \int_{-\infty}^x \int_{-\infty}^y p(\tilde{x}, \tilde{y}) d\tilde{y} d\tilde{x}.$$

$$p(\tilde{x}) = \int_{-\infty}^{\infty} p(\tilde{x}, \tilde{y}) d\tilde{y}.$$

Vad är sannolikhet?

Kan man mäta $P(A)$? Ja, om man kan göra ett stort antal observationer N och tolka den relativa frekvensen

$$\frac{N_A}{N} \approx P(A), \quad N_A = \# \text{ gånger som } A \text{ observerades.}$$

Detta kan endast tillämpas på slumpvariabler.

Bayesiansk sannolikhet:

Tolkar sannolikhet som en grad av tro. Kan tillämpas på alla intressanta vetenskapliga frågeställningar.

$$\text{Produktregeln: } P(A \cap B | I) \stackrel{(*)}{=} P(A | B, I) P(B | I) = P(B | A, I) P(A | I).$$

$$\Rightarrow P(A | B, I) = \frac{P(B | A, I) P(A | I)}{P(B | I)}. \quad (\text{Bayes sats}) \quad \uparrow \text{ tillgänglig information}$$

- $P(A | I)$ - A priorisannolikhet.
- $P(A | B, I)$ - A posteriorisannolikhet.

Ex.

A = Higgsbosonens massa är mellan $[m_H, m_H + dm_H]$.

$$\Rightarrow P(A | I) = p(m_H | I) dm_H.$$

B = Observerar datapunkt i $[D, D + dD]$.

$$\Rightarrow P(B | I) = p(D | I) dD.$$

$$P(A | B, I) = p(m_H | D \cap I) dm_H.$$

$$P(B | A, I) = p(D | m_H \cap I) dD.$$

$$\begin{aligned} \text{Bayes sats} \Rightarrow p(m_H | D, I) dm_H &= \frac{p(D | m_H, I) dD p(m_H | I) dm_H}{p(D | I) dD} \\ &= \frac{p(D | m_H, I) p(m_H | I)}{p(D | I)}, \end{aligned}$$

där $p(D | m_H, I)$ är trolighetsfunktion (likelihood);

$p(D | I)$ är marginell trolighet.

Inferens:

Slutsats som dras utifrån kontext.

Deduktiv inferens:

A, B, \dots = premisser som betraktas som sanna.

H = slutsats.

$P(H|A, B, \dots) = \{0, 1\}$ ($0 = H$ falskt, $1 = H$ sant).

Induktiv inferens:

Premisserna påverkar slutsatsen, men tillåter inte att avgöra sannolikhetshalten definitivt. (vanligaste inferensen)

Statistisk inferens används för att kvantifiera styrkan hos induktiv inferens via sannolikheter.

Data och modeller:

Betrakta en observation D av ett fysikaliskt fenomen. En modell $M(\theta)$ beskriver detta fenomen.

$$D = M(\theta) + \delta D + \delta M. \quad (*)$$

\uparrow exp. osäkerhet (stumpvariabel) \uparrow modellosäkerhet (motverkar överanpassning)

Datan lär oss om modellens parametrar. Formaliseras med Bayes teorem enligt

$$P(\theta | D, I) = \frac{P(D | \theta, I) p(\theta | I)}{p(D | I)}$$

$$(*) \Rightarrow p(D | \theta, I) \equiv \mathcal{L}(\theta).$$

Att inferera posteriorifördelningen $p(\theta | D, I)$ kallas för parameterestimering.

Förutsägelser:

Antag att modellen $M(\theta)$ förutsäger ett annat fysikaliskt fenomen, för vilket det ej ännu finns data. Kallar detta framtida data F .

$D \Rightarrow$ lärdom om $M \Rightarrow$ förutsäg F , d.v.s.

$$p(F | D, I) = \int_{\Omega} p(F, \theta | D, I) d\theta = \int_{\Omega} \underbrace{p(F | \theta, D, I)}_{\text{modellens förutsägelse}} \underbrace{p(\theta | D, I)}_{\text{a posteriori-fördelning}} d\theta.$$

\uparrow modell ingår här \uparrow marginalisering

Detta kallas för en a-posteriori förutsägelse.

Dessa integraler utförs i praktiken av Monte Carlo-tekniker

Bayesiansk linjär regression:

Betrakta modellen $M(\Theta; x) = \sum_{j=0}^{P-1} \Theta_j \phi_j(x)$.

Anpassas till data $D = \begin{pmatrix} y_1 \\ \vdots \\ y_{N_d} \end{pmatrix}$ med designmatrix $X = \begin{pmatrix} \phi_0(x_1) & \dots & \phi_{p-1}(x_1) \\ \vdots & & \vdots \\ \phi_0(x_{N_d}) & \dots & \phi_{p-1}(x_{N_d}) \end{pmatrix}$

$$D = X\Theta + \epsilon,$$

där $\Theta = \begin{pmatrix} \Theta_0 \\ \vdots \\ \Theta_{P-1} \end{pmatrix}$, $\epsilon = \begin{pmatrix} \epsilon_1 \\ \vdots \\ \epsilon_{N_d} \end{pmatrix}$ ← residualer
parametrar

Minimering av $\|\epsilon\|_2^2$ leder till normalekvationen med lösningen

$$\Theta^* = (X^T X)^{-1} X^T D.$$

Gör istället Bayesiansk inferens.

$$p(\Theta | D, I) = \frac{p(D | \Theta, I) p(\Theta | I)}{p(D | I)}.$$

Trolighetsfunktion:

Under antagandet att residualerna är oberoende, är

$$p(\epsilon | I) = \prod_{i=1}^{N_d} p(\epsilon_i | I).$$

Under antagandet att $E[\epsilon_i] = 0$, $\text{Var}[\epsilon_i] = \sigma_i^2$, samt att

trolighetsfunktionen beskrivs av en normalfördelning, är

$$p(\epsilon_i | I) = \frac{1}{\sqrt{2\pi\sigma_i^2}} \exp\left(-\frac{1}{2} \frac{\epsilon_i^2}{\sigma_i^2}\right).$$

$$D_i = \underbrace{(X\Theta)_i}_{\equiv M_i} + \epsilon_i. \quad (**)$$

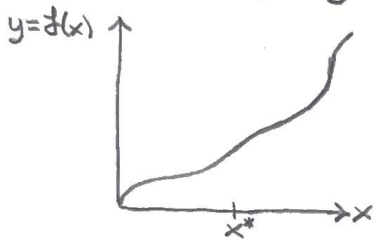
Osäkerhet i ϵ_i leder till osäkerhet i D_i , men hur?

Felpropagering, variabelbyte:

Betrakta slumpvariabel X , $x \stackrel{\text{"fördelat enligt"}}{\sim} p_X(x | I)$.

Betrakta $Y = \phi(X)$ (unik avbildning).

Vad blir $p_Y(y|I)$?



Betrakta godtycklig punkt x^* och litet intervall δx .

$$\Rightarrow P(x^* - \frac{\delta x}{2} \leq x \leq x^* + \frac{\delta x}{2}) = p_X(x^*|I) \delta x.$$

Avbildningen ger $y^* = f(x^*)$.

Konsivering av sannolikhet innebär att det måste finnas ett intervall δy s.a.

$$p_X(x^*|I) \delta x = p_Y(y^*|I) \delta y.$$

Detta ska gälla för alla punkter. Låt $\delta x, \delta y \rightarrow 0$.

$$\Rightarrow p_X(x|I) = p_Y(y=y(x)|I) \left| \frac{dy}{dx} \right| \text{ eller } p_Y(y|I) = p_X(x=x(y)|I) \left| \frac{dx}{dy} \right|.$$

Betrakta (**). ^{Jacobian} Ovanstående ger

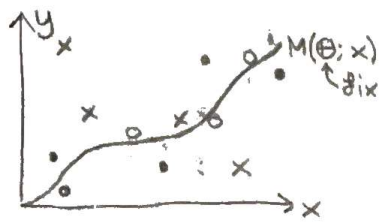
$$p(D_i | \theta, I) = p(\varepsilon_i(D_i) | \theta, I) \cdot 1 = \frac{1}{\sqrt{2\pi\sigma_i^2}} \exp\left(-\frac{1}{2} \frac{(D_i - M_i)^2}{\sigma_i^2}\right).$$

Bayesianisk linjär regression: (forts.)

$$M(\theta; x) = \sum_{j=0}^{p-1} \theta_j \phi_j(x) \text{ anpassas till data } D_{\text{obs}} = \begin{pmatrix} y_1 \\ \vdots \\ y_{N_d} \end{pmatrix}.$$

Bayes teorem ger $p(\theta | D_{\text{obs}}, I) \propto p(D_{\text{obs}} | \theta, I) p(\theta | I)$.

Trolighetsfunktion - beskriver sannolikheten att observera en uppsättning data givet en specifik datagenererande process (i detta fall $M(\theta; x)$).



$\begin{cases} \bullet & D_1 \text{ (observerade datapunkter).} \\ \circ & D_2 \\ \times & D_3 \end{cases}$

$$p(D_2 | \theta, I) > p(D_1 | \theta, I) \gg p(D_3 | \theta, I).$$

Med den observerade datan vill man hitta den modell som maximerar trolighetsfunktionen.

$$p(D_{\text{obs}} | \theta, I) \equiv L(\theta).$$

Betrakta felmodell för residualerna ϵ som oberoende slumpvariabler.

$$p(\epsilon_i | I) = N(0, \sigma_i^2) \text{ (normalförd.)}$$

$$\Rightarrow p(D_i | \theta, I) = \frac{1}{\sqrt{2\pi\sigma_i^2}} \exp\left(-\frac{(D_i - M_i)^2}{2\sigma_i^2}\right), \quad M_i = M(\theta; x_i).$$

Oberoende fel ger

$$p(D | \theta, I) = \prod_{i=1}^{N_d} p(D_i | \theta, I) = \frac{1}{(2\pi)^{N_d/2}} \frac{1}{|\Sigma_\epsilon|^{1/2}} \exp\left(-\frac{1}{2}(D - X\theta)^T \Sigma_\epsilon^{-1} (D - X\theta)\right),$$

$$\Sigma_\epsilon = \begin{pmatrix} \sigma_1^2 & & 0 \\ & \ddots & \\ 0 & & \sigma_{N_d}^2 \end{pmatrix} \text{ kovariansmatris, } X \text{ designmatris.}$$

I specialfallet då $\sigma_i^2 = \sigma_\epsilon^2 \forall i$, gäller

$$p(D | \theta, I) = \frac{1}{(2\pi\sigma_\epsilon^2)^{N_d/2}} \exp\left(-\frac{1}{2} \frac{(D - X\theta)^T (D - X\theta)}{\sigma_\epsilon^2}\right).$$

För en linjär modell kan man skriva om den normalfördelade troligheten till något som är proportionellt mot

$$\exp\left(-\frac{1}{2}(\theta - \theta^*)^T \Sigma_\theta^{-1} (\theta - \theta^*)\right), \quad \theta^* = (X^T X)^{-1} X^T D, \quad \Sigma_\theta^{-1} = \frac{X^T X}{\sigma_\epsilon^2}.$$

• Det ser ut som en normalfördelning för θ , men är ej normaliserad.

- Moden (maximum) är lika med Θ^* , dvs lösningen till normal ekvationen. \Rightarrow linjär regression motsvarar "Maximum likelihood estimate" (MLE).

A priori-fördelning:

Betrakta två möjligheter:

$$a) p(\Theta | I_a) = \begin{cases} \left(\frac{1}{\Delta\Theta}\right)^{Np} & \text{om } \Theta_i \in \left[-\frac{\Delta\Theta}{2}, \frac{\Delta\Theta}{2}\right] \forall i. \\ 0 & \text{annars.} \end{cases}$$

dvs likförmig fördelning med bredd $\Delta\Theta$.

$$b) p(\Theta | I_b) = \frac{1}{(2\pi\sigma_\Theta^2)^{Np/2}} \exp\left(-\frac{\Theta^T \Theta}{\sigma_\Theta^2}\right).$$

dvs en i.i.d normalfördelning för p parametrar.

A posteriori-fördelning:

Betrakta analys av a):

$$p(\Theta | D_{\text{obs}}, I_a) \propto \begin{cases} \exp\left(-\frac{1}{2}(\Theta - \Theta^*)^T \Sigma_\Theta^{-1} (\Theta - \Theta^*)\right) & \text{om } \Theta_i \in \left[-\frac{\Delta\Theta}{2}, \frac{\Delta\Theta}{2}\right] \forall i. \\ 0 & \text{annars.} \end{cases}$$

- Givet $\Theta_i^* \in \left[-\frac{\Delta\Theta}{2}, \frac{\Delta\Theta}{2}\right]$, har denna a-posteriori-fördelning samma mod som trolighetsfunktionen.
- Kan tolkas som en pdf för Θ .

Betrakta analys av b):

$$p(\Theta | D_{\text{obs}}, I_b) = \left\{ \begin{array}{l} \text{produkt av två} \\ \text{normalförd.} \end{array} \right\} \propto \exp\left(-\frac{1}{2}(\Theta - \tilde{\Theta})^T \tilde{\Sigma}^{-1} (\Theta - \tilde{\Theta})\right),$$

$$\tilde{\Theta} = \tilde{\Sigma}_\Theta \Sigma_\Theta^{-1} \Theta^*, \quad \tilde{\Sigma}^{-1} = \Sigma_\Theta^{-1} + \frac{1}{\sigma_\Theta^2} \mathbb{I}_{p \times p}.$$

- Inferensen ger en normalfördelning
- Jämfört med a priori:

$$\mathbb{E}[\Theta]: 0 \xrightarrow{\text{data}} \tilde{\Theta}.$$

$$\text{Var}[\Theta]: \sigma_\Theta^2 \mathbb{I}_{p \times p} \xrightarrow{\text{data}} \tilde{\Sigma}_\Theta.$$

Posteriorifördelningar och kostnadsfunktioner:

Betrakta modell $M(\theta; x)$ och data $D = \begin{pmatrix} y_1 \\ \vdots \\ y_{N_d} \end{pmatrix}$.

$$D_i = y_i, \quad M_i = M(\theta; x_i).$$

Förutsätt (oändligt stor) likförmig apriorifördelning och normalfördelad tröghetsfunktion.

A posteriori:

$$p(\theta | D, I) \propto \exp\left(-\frac{1}{2} \sum_{i=1}^{N_d} \frac{(D_i - M_i)^2}{\sigma_\varepsilon^2}\right).$$

Moden kallas MAP (maximum a posteriori). (vanlig kostnadsfunktion)

$$\theta^* = \operatorname{argmax}_{\theta \in \mathbb{R}^p} p(\theta | D, I) = \operatorname{argmax}_{\theta \in \mathbb{R}^p} \log p = \operatorname{argmin}_{\theta \in \mathbb{R}^p} \sum_{i=1}^{N_d} (D_i - M_i)^2.$$

Antag att apriorifördelningen är:

$$p(\theta | I) \propto \exp\left(-\frac{1}{2} \sum_{i=1}^p \frac{\theta_i^2}{\sigma_\theta^2}\right).$$

Vilket ger a posteriorifördelningen:

$$p(\theta | D, I) \propto \exp\left(-\frac{1}{2} \left(\sum_{i=1}^{N_d} \frac{(D_i - M_i)^2}{\sigma_\varepsilon^2} + \sum_{i=1}^p \frac{\theta_i^2}{\sigma_\theta^2} \right)\right).$$

MAP-punkten för denna är:

$$\theta^* = \operatorname{argmin}_{\theta \in \mathbb{R}^p} \left(\sum_{i=1}^{N_d} (D_i - M_i)^2 + \frac{\sigma_\varepsilon^2}{\sigma_\theta^2} \sum_{i=1}^p \theta_i^2 \right).$$

(motsvarar Ridgeregularisering)

Modellering med osäkerhet:Felpropagering:

Betrakta parametrisk modell. Dataanalys för att estimeras parametrarna har genomförts.

Propagera osäkerheter till prediktioner:

Tre metoder:

1. Marginalisering.
2. Variabeltransformation.
3. Approximativ felpropagering.

Marginalisering:

Ex. 1 Modell med två parametrar θ, ϕ , varav bara θ är av fysikaliskt intresse.

Inferens ger $p(\theta, \phi | D, I)$. (nuisance parameter)

$$p(\theta | D, I) = \int p(\theta, \phi | D, I) d\phi \quad (\text{marginalisering}).$$

Ex. 2 $z = f(x, y)$. Givet $p_{x,y}(x, y | D, I)$, vad är $p_z(z | D, I)$?

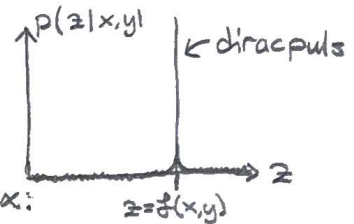
Använd marginalisering och produktregeln.

$$p(z | D, I) = \int p(x, y, z | D, I) dy dx = \int p(z | x, y, D, I) \cdot p(x, y | D, I) dy dx.$$

Eftersom x, y ger z , är

$$p(z | x, y, D, I) = \delta(z - f(x, y))$$

$$\Rightarrow p(z | D, I) = \int \delta(z - f(x, y)) p(x, y | D, I) dy dx.$$



Kan evaluera integral mha deltafunktioner.

(se t.ex. 17.1 ($z = x + y$))

Variabeltransformation:

Betrakta $\mathbf{x} = \{x_j\}_{j=1}^M$, $\mathbf{y} = \{y_j\}_{j=1}^M$, $x_j = x_j(\mathbf{y})$, $y_j = y_j(\mathbf{x})$.

Om $p_{\mathbf{y}}(\mathbf{y} | I)$ är känd, vad är $p_{\mathbf{x}}(\mathbf{x} | I)$?

Konservering av sannolikhet ger

$$p_{\mathbf{x}}(\mathbf{x} | I) = p_{\mathbf{y}}(\mathbf{y}(\mathbf{x}) | I) \left| \frac{\partial(y_1, \dots, y_M)}{\partial(x_1, \dots, x_M)} \right| \quad \leftarrow (\text{Jacobian})$$

T.ex. för $y = y(x)$:

$$p_x(x | I) = p_y(y(x) | I) \left| \frac{dy}{dx} \right|.$$

Approximativ felpropagering:

Betrakta $Z = f(X, Y, \dots)$,

$$E[X] = x_0, \quad E[Y] = y_0, \dots$$

$$\text{Var}[X] = \sigma_x^2, \quad \text{Var}[Y] = \sigma_y^2, \dots$$

Vad är $E[z]$, $\text{Var}[z]$?

• Ofta gäller $E[z] = z_0 = f(x_0, y_0, \dots)$ åtminstone approximativt.

För $\text{Var}[z]$:

• Differentiera sambandet runt (x_0, y_0, \dots) , dvs Taylorutveckla runt (x_0, y_0, \dots) och behåll upp t.o.m. linjära termer i $(x-x_0), (y-y_0), \dots$

$$\begin{cases} \delta X = X - x_0 \\ \delta Y = Y - y_0 \\ \delta z = z - z_0 \dots \end{cases}$$

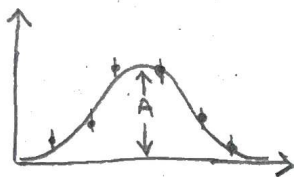
• Kvadrera detta uttryck och beräkna väntevärdet. Notera att

$$\begin{aligned} E[\delta z^2] &= E[(z - E[z])^2] = \sigma_z^2 \\ E[\delta X^2] &= \dots = \sigma_x^2 \\ E[\delta Y^2] &= \dots = \sigma_y^2 \dots \end{aligned}$$

Var försiktig! Antagligen förs termer $E[\delta X \delta Y] = E[(X - E[X])(Y - E[Y])] = \text{Cov}[X, Y]$, alltså påverkar det om X, Y, \dots är beroende av varandra.

Ex. 17.4

Mät ett spektrum och estimeramplitud A .



$$\begin{aligned} E[A] &= A_0 \\ \text{Var}[A] &= \sigma_A^2 \end{aligned}$$

Givet att $A = |f|^2$, f komplex strukturfaktor.

Uppskatta $|f| = f_0 + \sigma_f^2$.

$$A_0 = f_0^2 \Rightarrow f_0 = \sqrt{A_0}$$

Differentiera uttrycket.

$$\delta A = 2f \delta f$$

Taylorutveckla runt $f = f_0$.

$$A = f_0^2 + \left. \frac{dA}{df} \right|_{f=f_0} (f - f_0) + O((f - f_0)^2) = f_0^2 + 2f_0(f - f_0) + O((f - f_0)^2)$$

Definiera $\delta A = A - A_0$, $\delta f = f - f_0 \Rightarrow \delta A = 2f_0 \delta f$.

$$\Rightarrow \Delta A^2 = 4z_0^2 \Delta z^2$$

$$E[\Delta A^2] = E[4z_0^2 \Delta z^2] \Rightarrow \sigma_A^2 = 4z_0^2 E[\Delta z^2] = 4z_0^2 \sigma_z^2.$$

$$\sigma_z = \frac{\sigma_A}{2\sqrt{A_0}}, \text{ dvs } z = \sqrt{A_0} \pm \frac{\sigma_A}{2\sqrt{A_0}}.$$

Använd även pdf:er och variabeltransformation:

$$p_z(z | D, I) = p_A(A(z) | D, I) \left| \frac{dA}{dz} \right| = 2z p_A(A(z) | D, I).$$

Konkret om

$$p_A(A | D, I) = \frac{1}{\sqrt{2\pi}\sigma_A} \exp\left(-\frac{1}{2} \frac{(A-A_0)^2}{\sigma_A^2}\right) \Rightarrow p_z(z | D, I) = \frac{2z}{\sqrt{2\pi}\sigma_A} \exp\left(-\frac{1}{2} \frac{(z^2 - z_0^2)^2}{\sigma_A^2}\right).$$

(Se figurer i kompendiet med vad som kan gå fel med approximationen)

Stokastiska processer:

Deterministiska processer - givet resultat för given input.

Stokastiska processer - beskriver fenomen som kännetecknas av osäkerheter och slumpmässiga förlopp.

Ex. Tillväxtmodell.

Antalet bakterier $y(t)$ växer deterministiskt med tillväxtfaktor λ

$$\Rightarrow \frac{dy}{dt} = \lambda y(t) \Rightarrow \{y(0) = y_0\} \Rightarrow y(t) = y_0 e^{\lambda t}$$

Detta beskriver inte det stokastiska förloppet.

$\lambda \Delta t$ = sannolikheten att en bakterie delar på sig under ett kort tidsintervall.

Def. Stokastisk process:

[En uppsättning indexerade slumpvariabler X_t med] gemensamt utfallsrum S .

Diskreta processer: $t \in \{0, 1, 2, \dots\}$.

Kontinuerliga processer: $t \in \mathbb{R}, t \in [0, \infty), \dots$

Olika sekvenser av en given process kommer ge olika utfall. Processen beskrivs med (betingade) sannolikheter

$$P_{X_0}(x_0), P_{X_1|X_0}(x_1|x_0), \dots, P_{X_n|X_0, \dots, X_{n-1}}(x_n|x_0, \dots, x_{n-1}).$$

Ex. Slumpvandring (diskret, endimensionell).

$$S_n = a + \sum_{i=1}^n X_i$$

↑ position efter n steg ↙ startposition

← slumpvariabel, $S = \{-1, 1\}$.

$\{X_t: t \in \mathbb{N}\}$ är den intressanta stokastiska processen.

Notera att $P_{S_n|S_0, \dots, S_{n-1}}(s_n|s_0, \dots, s_{n-1}) = P_{S_n|S_{n-1}}(s_n|s_{n-1}) =$

$$= \begin{cases} 0,5 & \text{för } s_n = s_{n-1} + 1. \\ 0,5 & \text{för } s_n = s_{n-1} - 1. \\ 0 & \text{för alla andra } s_n. \end{cases}$$

Dvs att ett utfall beror endast av det senaste utfallet.
Detta kallas Markovkedja.

För kontinuerliga utfallsrum: "

$$P_{z_n | z_0, \dots, z_{n-1}}(x_n | x_0, \dots, x_{n-1}) = P_{z_n | z_{n-1}}(x_n | x_{n-1}).$$

Stationära Markovkedjor:

En Markovkedja är stationär om de betingade sannolikheterna ej beror på tidsindexet, dvs för diskreta utfallsrum

$$P_{z_{n+1} | z_n}(j | i) = P_{z_1 | z_0}(j | i) = T(i, j) \quad \forall n, i, j.$$

T kallas övergångsmatris ($|S| \times |S|$).

För kontinuerliga utfallsrum:

$$P_{z_{n+1} | z_n}(x_j | x_i) = P_{z_1 | z_0}(x_j | x_i) = T(x_i, x_j),$$

där T är en (kontinuerlig) övergångsdensitet.

Ibland skrivs T som $T(x_j \leftarrow x_i)$, T_{ij} .

Gränsdistributioner:

Betrakta diskret utfallsrum.

Frekvensfunktionen blir en vektor π med element

$$\pi_j = P_z(j), \quad j \in S.$$

En vektor π är gränsdistributionen till Markovkedjan som beskrivs av övergångsmatrisen T om

$$1) \pi_j \geq 0 \quad \forall j, j \in S.$$

$$2) \sum_j \pi_j = 1.$$

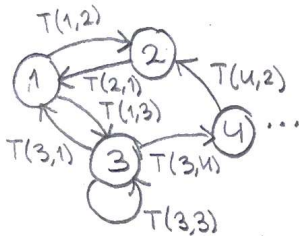
$$3) \pi_j = \lim_{n \rightarrow \infty} T^n(i, j) \quad \forall i, j \in S.$$

dvs $\pi = \lim_{n \rightarrow \infty} \alpha T^n$, α godtycklig startdistribution.

Markk kedjor (forts.):

Betrakta diskret stationär Markk kedja med diskret utfallsrum. Detta beskrivs entydigt med

- 1) $p_{z_0}(i) = P(Z_0 = i)$ (trelvensfunktion för Z_0).
- 2) $T(i, j) = P_{z_{n+1}|z_n}(j|i)$ (övergångsmatrix).



Inför $\pi^{(0)} = (\pi_1^{(0)}, \dots, \pi_{|\mathcal{S}|}^{(0)})$, $\pi_i^{(0)} = p_{z_0}(i)$.
↖ # möjliga utfall

Vad blir π^1 ?

$$\pi_i^{(1)} = p_{z_1}(i) = \sum_j p_{z_0, z_1}(j, i) = \sum_j p_{z_1|z_0}(i|j) p_{z_0}(j) = \sum_j T(j, i) \pi_j^{(0)} \Rightarrow$$

$$\Rightarrow \pi^{(1)} = \pi^{(0)} T \Rightarrow \pi^{(n)} = \pi^{(0)} T^n$$

Under vissa förutsättningar existerar en gränsfördelning π så att

$$\pi = \lim_{n \rightarrow \infty} \pi^{(n)} = \pi T, \quad \alpha = \pi^{(0)} \text{ godtycklig startfördelning.}$$

Stationär fördelning:

Vektorn π är en stationär fördelning till Markk kedjan (som beskrivs av T) om

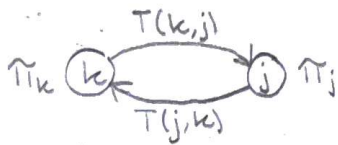
- 1) $\pi_j \geq 0 \quad \forall j, j \in \mathcal{S}$.
- 2) $\sum_j \pi_j = 1$.
- 3) $\pi_j = \sum_i \pi_i T(i, j) \quad \forall j \quad (\pi = \pi T)$.

En gränsfördelning är alltid en stationär fördelning, men det motsatta gäller inte alltid.

Reversibla kedjor:

En stationär Markk kedja (övergångsmatrix T) med stationär fördelning π är reversibel om

$$\pi_k T(k, j) = \pi_j T(j, k) \quad \forall j, k \in S.$$



Detta beskriver en detaljerad balans av förflyttad sannolikhetsvekt. Sambandet innebär också att

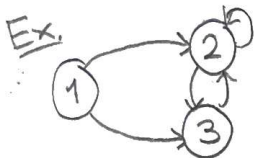
$$IP_{z_{n-1}, z_n}(k, j) = IP_{z_{n-1}, z_n}(j, k),$$

dvs invariant under tidsreversion.

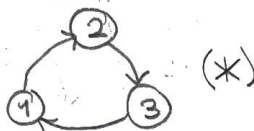
Denna egenskap garanterar att processen är stationär.

Existensvillkor för gränsdistribution:

- 1) Icke-reducerbar: Alla tillstånd kan nå oavsett startpunkt.
- 2) Positivt återkommande: Sannolikheten att återkomma till ett tillstånd någon gång är $> 0 \quad \forall i$.
- 3) Icke-periodisk: Sekvenser av tillstånd upprepas inte genom hela kedjan.



reducerbar
ej positivt återkommande
ej periodisk



ej reducerbar
positivt återkommande
periodisk

Studera (*):

- Har en stationär fördelning $\pi = (\frac{1}{3}, \frac{1}{3}, \frac{1}{3})$.

- Men detta är ej en gränsfördelning, ty exv

$$\pi^{(0)} = (1, 0, 0), \quad \pi^{(1)} = (0, 1, 0), \dots$$

Repetition: (Markovkedjor)

Betrakta stationär Markovkedja med diskret utfallsrum.

$$P_{Z_{n+1}|Z_n}(j|i) = T(i, j) \quad \forall i, j \in S \text{ ober. av } n.$$

$$\pi_j^{(n)} = P_{Z_n}(j).$$

Vid jämvikt:

$$\pi = \pi T, \quad \pi \text{ kedjans stationärfördelning.}$$

Ett starkare villkor som garanterar jämvikt är detaljerad balans

$$\pi_i T(i, j) = \pi_j T(j, i).$$

I vissa fall är stationärfördelningen även kedjans gränsfördelning, dvs kedjan konvergerar mot π oavsett startfördelning α .

$$\pi = \lim_{n \rightarrow \infty} \alpha T^n.$$

Tillräckliga (men ej nödvändiga) villkor för existens av gränsfördelning:

- Irreducibel.
- Positivt återkommande.
- Aperiodisk.

Metropolisdesign:

Man kan konstruera en Markovkedja (diskret utfallsrum) så att den har en specifik gränsfördelning π .

Bevis:

Introducera en stegförslagsfunktion:

$$S(i, j) = P_{Z_{n+1}|Z_n}(\text{föreslå position } j | i).$$

• Måste vara en stokastisk matris, annars valfri.

$$\sum_j S(i, j) = 1 \quad \forall i, \quad S(i, j) \in [0, 1].$$

Introducera acceptansfunktion:

$$A(i, j) = P_{Z_{n+1}|Z_n}(\text{acceptera } j | i), \quad A(i, j) \in [0, 1].$$

A och S är oberoende för $i \neq j$, vilket ger

$$T(i, j) = P_{Z_{n+1}|Z_n}(j|i) = S(i, j)A(i, j), \quad i \neq j.$$

Detaljerad balans

$$\pi_i S(i, j) A(i, j) = \pi_j S(j, i) A(j, i) \Rightarrow \frac{A(i, j)}{A(j, i)} = \frac{\pi_j S(j, i)}{\pi_i S(i, j)},$$

vilket uppfylls om

$$A(i, j) = \min\left(1, \frac{\pi_j S(j, i)}{\pi_i S(i, j)}\right).$$

Diagonalelementen fås genom normalisering

$$\begin{aligned} T(i, i) &= 1 - \sum_{j \neq i} T(i, j) = \sum_j S(i, j) - \sum_{j \neq i} S(i, j) A(i, j) = \\ &= S(i, i) + \sum_{j \neq i} S(i, j) [1 - A(i, j)]. \end{aligned}$$

sannolikheten att inte acceptera ett steg från $i \rightarrow j$.

Algoritm:

1. Välj $S(i, j)$.
2. Beräkna $A(i, j) \forall i \neq j \Rightarrow T(i, j) \forall i \neq j$.
3. Beräkna $T(i, i) \forall i$.

\Rightarrow Markovkedja med stationärfördelning π .

För att vara säker på att π också är en gränsfördelning kan man bekräfta att T har tillräckliga egenskaper, dvs aperiodisk, irreducibel, positivt återkommande.

Kontinuerliga, multivariata utfallsrum. Skillnader från tidigare:

Utfall: $i \rightarrow \Theta \in \mathbb{R}^p$.

Fördelningar: $\pi_i = p(i) \rightarrow p(\Theta)$.

Sannolikheter: $P_{Z_n}(i) = \pi_i \rightarrow P(Z_n \in \Delta\Theta) = \int_{\Delta\Theta} p(\Theta') d\Theta'$.

För stationära Markovkedjor:

$T(i, j) \rightarrow T(\Theta, \Theta') = P_{Z_{n+1}|Z_n}(\Theta'|\Theta)$ (ober. av n).

uppfyller $\int_{\mathbb{R}^p} T(\Theta, \Theta') d\Theta' = 1 \quad \forall \Theta \in \mathbb{R}^p$.

Stationär Markovkedja i jämvikt:

$$p(\theta) = \int_{\mathbb{R}^p} p(\theta') T(\theta', \theta) d\theta'$$

Starkare villkor med detaljerad balans

$$p(\theta) T(\theta, \theta') = p(\theta') T(\theta', \theta).$$

Metropolis(-Hastings) algoritmen:

$$T(\theta, \theta') = J(\theta, \theta') A(\theta, \theta').$$

Stegförslagsfunktionen måste vara en pdf:

$$P_{z_{n+1}|z_n}(\text{föreslå } \theta' | \theta).$$

T.ex. likförmig i en liten volym nära θ ,

$$J(\theta, \theta') = \begin{cases} \frac{1}{\Delta \theta^p} & \text{då } |\theta'_i - \theta_i| \leq \frac{\Delta \theta}{2} \quad \forall i \in \{1, \dots, p\}. \\ 0 & \text{annars.} \end{cases}$$

eller normalfördelad

$$J(\theta, \theta') = \mathcal{N}(\theta' | \theta, \Sigma).$$

Acceptansfunktionen

$$A(\theta, \theta') = \min\left(1, \underbrace{\frac{p(\theta')}{p(\theta)} \frac{J(\theta, \theta')}{J(\theta', \theta)}}_{\equiv r \text{ (Metropoliskvoten)}}\right).$$

För symmetriska J fås

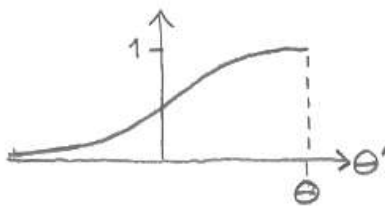
$$r = \frac{p(\theta')}{p(\theta)}.$$

Stickprov från en pdf:

En (pseudo-)slumptalsgenerator kan producera

$$\{u_i\}_{i=1}^N, u_i \sim U([0, 1]).$$

Hur produceras slumptal $\{\theta_i\}_{i=1}^N, \theta_i \sim N(0, 1)$?



$$P(\theta) = \int_{-\infty}^{\theta} p(\theta') d\theta'$$

$$P(\theta) = u_i \Rightarrow \theta_i = P^{-1}(u_i) \sim p(\theta).$$

Fungerar för analytiska, icke-inverterbara fördelningar.

Man kan använda MCMC-algoritmer för att dra stickprov.

Metropolisalgoritmen:

1. Initialisera samplingen genom att dra θ_0 ur en startfördelning.
2. Dra stickprov från slumpvariabel X_{i+1} enligt följande:
 - a) Givet utfallet θ_i för X_i , föreslå ϕ genom dragning ur $S(\theta_i, \phi)$.
 - b) Beräkna $r = \frac{p(\phi)}{p(\theta_i)}$ (Metropolis - om S symmetrisk),
 $r = \frac{p(\phi)S(\phi, \theta_i)}{p(\theta_i)S(\theta_i, \phi)}$ (Metropolis-Hastings - om S ej symmetrisk)
 - c) Om $r \geq 1$: $\theta_{i+1} \leftarrow \phi$.
 Om $r < 1$: Dra $u \sim U([0, 1])$, om
 - $u \leq r$: $\theta_{i+1} \leftarrow \phi$.
 - $u > r$: $\theta_{i+1} \leftarrow \theta_i$.

Fortsätt tills kedjan har konvergerat (svårt) för då kommer alla nya utfall vara stickprov ur $p(\theta)$. Notera att man bara evaluerar $p(\theta)$ punktvís (i b). Algoritmen fungerar i höga dimensioner, men ineffektiv. (Mer sofistikerade algoritmer finns.)

Bayesiansk inferens:

Utmaningar:

1. Utforska a posteriori-fördelningar $p(\theta | D, I)$.

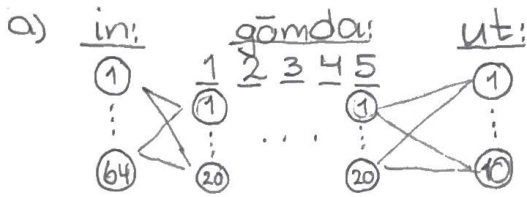
- Var finns moden?
- Finns det flera moder?
- Vad är sannolikhetsmassan koncentrerad?
- Parameterkorrelationer?
- Hur ser marginalfördelningarna ut?

2. A posteriori-förutsägelser.

- Evaluera modell $g(\theta)$ för $\theta \sim p(\theta | D, I)$.
- Förväntansvärde $E_p[g(\theta)] = \int g(\theta) p(\theta | D, I) d\theta$
- Fördelning av modellprediktioner (PPD). $\{g(\theta): \theta \sim p(\theta | D, I)\}$.

Samtliga frågor kan besvaras med stickprov $\{\theta_i\}_{i=1}^N$,
 $\theta_i \sim p(\theta | D, I)$.

2301-2 | (Läs uppgiftsformulering själv)



Utsignal från lager 1:

$$y_j^1 = f(z^1), \quad y_i^1 = f(z_i^1) = \frac{1}{1 + e^{-z_i^1}}$$

y_j^1, z^1 har storleken (1×20) .

Aktiveringarna ges av

$$z^1 = y_j^{\text{in}} W^1 + W^1_0$$

$\begin{matrix} \uparrow & \uparrow & \uparrow \\ \text{(insignaler)} & \text{(vikter)} & \text{(biasvikter)} \\ (1 \times N_{\text{in}}) & (N_{\text{in}} \times N_1) & (1 \times N_1) \end{matrix}$

b) Det finns vikter i alla gömda lager samt i utlagret. Varje nod (i) i ett lager (j) har $N_{j-1} + 1$ vikter.

Antal parametrar per lager:

(vikter för varje nod i tidigare lagret + bias)

$$20 \cdot \frac{1}{64+1} + \frac{2}{20} + \frac{3}{20} + \frac{4}{20} + \frac{5}{20} + \frac{\text{ut}}{10} \cdot (20+1) = \underline{\underline{3190 \text{ vikter}}}$$

c) Mjuk klassificerare, 10 klasser

$$\Rightarrow \text{utsignal } \hat{y}_i = P(t=i), \quad i \in \{0, \dots, 9\}, \quad \sum_{i=0}^9 \hat{y}_i = 1.$$

Använd t.ex. sigmoidfunktion i utlagret

$$y_i^{\text{ut}} = \frac{1}{1 + e^{-z_i^{\text{ut}}}} \in [0, 1].$$

Normalisera

$$\hat{y}_i = \frac{y_i^{\text{ut}}}{\sum_{i=0}^9 y_i^{\text{ut}}}$$

vilket ger en diskret sannolikhetsfördelning \hat{y} .

Bevis 1 | Prediktion $\hat{y}^0 \equiv \hat{y}(x^0)$ jämförs med framtida observation $y^0 \equiv y(x^0)$. $y(x) = f(x) + \epsilon$, $f(x)$ underliggande funktionen, ϵ fel. $E[\epsilon] = 0$, $\text{Var}(\epsilon) = \sigma_\epsilon^2$.

Visa $E[(y^0 - \hat{y}^0)^2] = (\beta^0 - E[\hat{y}^0])^2 + \text{Var}[\hat{y}^0] + \sigma_\varepsilon^2$,
där $\beta^0 = \beta(x^0)$.

(Bias-variance theorem)

$$y^0 = \beta^0 + \varepsilon^0, \quad E[y^0] = \beta^0, \quad \text{Var}[y^0] = \sigma_\varepsilon^2. \quad \text{(I)}$$

$$\text{Söker } E[(y^0 - \hat{y}^0)^2] = E[(\beta^0 + \varepsilon^0 - \hat{y}^0)^2] = E[(\beta^0 + \varepsilon^0 - E[\hat{y}^0] + E[\hat{y}^0] - \hat{y}^0)^2].$$

$$\text{(I)}^2 \cdot E[(\beta^0 - E[\hat{y}^0] + \varepsilon^0)^2] = E[(\beta^0 - E[\hat{y}^0])^2] + E[\varepsilon^{02}] + 2E[(\beta^0 - E[\hat{y}^0])\varepsilon^0] =$$

$$= \left\{ \begin{array}{l} \text{Notera } \cdot \text{Var}[\varepsilon^0] = E[(\varepsilon^0 - E[\varepsilon^0])^2] = E[\varepsilon^{02}] = \sigma_\varepsilon^2 \\ \cdot 2(\beta^0 - E[\hat{y}^0])E[\varepsilon^0] = 0 \end{array} \right\} = E[(\beta^0 - E[\hat{y}^0])^2] + \sigma_\varepsilon^2.$$

$$\text{(II)}^2 \cdot E[(\hat{y}^0 - E[\hat{y}^0])^2] = \text{Var}[\hat{y}^0]. \quad (\text{modellens prediktionsvarians})$$

$$\text{(I)(II)} \cdot 2E[(\beta^0 - E[\hat{y}^0] + \varepsilon^0)(E[\hat{y}^0] - \hat{y}^0)]$$

$$E[(\cdot)] = \dots = 0$$

Bevis 3 | Betrakta $D = X\theta + \varepsilon$, $p(\varepsilon_i | I) = \mathcal{N}(0, \sigma_\varepsilon^2) = \frac{1}{\sqrt{2\pi}\sigma_\varepsilon} \exp\left[-\frac{1}{2} \frac{\varepsilon_i^2}{\sigma_\varepsilon^2}\right]$

a) Betrakta observation D_i , jämför $M_i = (X\theta)_i$, $D_i = M_i + \varepsilon_i$.

$$p(D_i | \theta, I) = \{\text{var. byte}\} = p_\varepsilon(\varepsilon_i = D_i - M_i | \theta, I) \left| \frac{d\varepsilon_i}{dD_i} \right| =$$

$$= \frac{1}{\sqrt{2\pi}\sigma_\varepsilon} \exp\left[-\frac{1}{2} \frac{(D_i - M_i)^2}{\sigma_\varepsilon^2}\right].$$

Felen oberoende \Rightarrow

$$\Rightarrow p(D | \theta, I) = \prod_{i=1}^{N_d} p(D_i | \theta, I) = \left(\frac{1}{\sqrt{2\pi}\sigma_\varepsilon}\right)^{N_d/2} \exp\left[-\frac{1}{2\sigma_\varepsilon^2} \sum_{i=1}^{N_d} (D_i - M_i)^2\right] =$$

$$= \underbrace{\left(\frac{1}{\sqrt{2\pi}\sigma_\varepsilon}\right)^{N_d/2} \exp\left[-\frac{1}{2} \frac{(D - X\theta)^T (D - X\theta)}{\sigma_\varepsilon^2}\right]}_{\equiv -L(\theta)}. \quad \square$$

b) Taylorutveckla $L(\theta)$ runt $\theta^* = (X^T X)^{-1} X^T D$.

$$L(\theta) = L(\theta^*) + \sum_{j=0}^{N_p-1} \frac{\partial L}{\partial \theta_j} \Big|_{\theta=\theta^*} (\theta_j - \theta_j^*) + \frac{1}{2} \sum_{j=0}^{N_p-1} \left(\sum_{i=0}^{N_p-1} \frac{\partial^2 L}{\partial \theta_i \partial \theta_j} \Big|_{\theta=\theta^*} \right)$$

$$(\theta_i - \theta_i^*)(\theta_j - \theta_j^*) + \dots = L(\theta^*) + \nabla L \Big|_{\theta=\theta^*} (\theta - \theta^*) + \dots$$

$$+\frac{1}{2}(\Theta - \Theta^*)^T H|_{\Theta = \Theta^*} (\Theta - \Theta^*) + \dots$$

∇L och H ges i uppgiftsledtråden.

$$\nabla L|_{\Theta = \Theta^*} = \frac{X^T D - X^T X \Theta^*}{\sigma_\epsilon^2} = \{D = X \Theta^*\} = 0.$$

$$H = \frac{X^T X}{\sigma_\epsilon^2} \text{ oberoende av } \Theta \Rightarrow \frac{\partial^{(n)} L}{\partial \theta^i} = 0 \text{ för } n \geq 3.$$

$$\text{Alltså, } L(\Theta) = L(\Theta^*) + \frac{1}{2}(\Theta - \Theta^*)^T \underbrace{\frac{X^T X}{\sigma_\epsilon^2}}_{\equiv \Sigma_\Theta^{-1}} (\Theta - \Theta^*).$$

$$p(D|\Theta, I) = \underbrace{\left(\frac{1}{2\pi\sigma_\epsilon^2}\right)^{N_0/2} \exp[-L(\Theta^*)]}_{= p(D|\Theta^*, I)} \exp\left[-\frac{1}{2}(\Theta - \Theta^*)^T \Sigma_\Theta^{-1} (\Theta - \Theta^*)\right]. \quad \square$$

(Läs uppgiftsbeskrivningarna själv)