

Felüggelt tanulás 2. - regressió \approx modellezés, ha

$$\boxed{X} \rightarrow \boxed{y} \text{ vagy } \boxed{Y} \quad y_i \in \mathbb{R}$$

- modell alkalmazása, ahol csak X ismert
- modell alkalmazása, ahol csak y ismert (kalibráció)
- modell értelmezése

lineáris, nemlineáris \rightarrow adott fv.

\rightarrow nehezebben átlátható pl. neurális háló

Sokszor próbálgatás, de hivatalosan:

- 1) X, y azonosítása; függvénykapcsolat hivatalosítása: ma ma termtud. háttér vagy feltehetően z^2 ; hiba tulajdonságai: ϵ_i homa vagy heteroskedasztikus, ismert e , $cov(\epsilon_i, \epsilon_j)$; hiba X és y közötti kapcsolat; becslési kritérium pl. $\sum |y_i - \hat{y}_i|$, $\sum (y_i - \hat{y}_i)^2$ - MLE, MPE extrapoláció-interpoláció.
- 2) paraméterezés
- 3) modell validálása

Legismertebb eset: egyszerűsített - lineáris regressió $X \rightarrow y$
least squares regression

X determinisztikus (hidamentes)

hiba csak y -ban

$$E(\epsilon_i) = 0, \epsilon_i \text{ ismert}, cov(\epsilon_i, \epsilon_j) = 0 \text{ } i \neq j \text{-re}$$

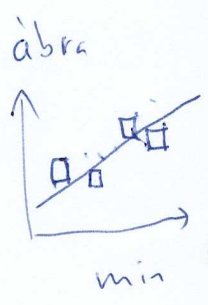
cél fu: $Q = \sum_{i=1}^N (y_i - f(x_i, p))^2 \cdot w_i$

x $y_i \leftarrow \text{mérés}$

$x_i - y_i$ N ismert pár

$w_i = 1/\sigma_i^2$

\uparrow
súly



~~f(x)~~ $y = ax + b$

\downarrow

súlyozatlanra

$Q(a, b) = \sum_{i=1}^N (y_i - ax_i - b)^2 \rightarrow \frac{\partial Q(a, b)}{\partial a} = 0$

parciális = 0

$\frac{\partial Q(a, b)}{\partial b} = 0$

normál
egyenletek

\bar{x}, \bar{y} átlagok

$a = \frac{\sum (x_i y_i - \bar{x} \bar{y})}{\sum (x_i^2 - \bar{x}^2)}$ $b = \bar{y} - a \bar{x}$

Miért négyzet? MLE-ben norm. eloszlás hiányában, σ_y^2 növekedés

val. sűrűség fu. $\rightarrow P_i(x | y_i) = \frac{1}{\sqrt{2\pi} \sigma_y} \cdot \exp\left(-\frac{(y_i - f(x_i, p))^2}{2\sigma_y^2}\right)$

$L = \prod_{i=1}^N P_i(x_i | y_i) = \dots \leftarrow$ ezt maximálni akarjuk

maximálni $\rightarrow L' = \ln L = \dots - \frac{1}{2} \sum \left(\frac{(y_i - f(x_i, p))^2}{\sigma_y^2} \right)$

$\underbrace{\hspace{10em}}$
minimalizálni

ezekkel f torzítatlan $E(f) = y$
kontinuus $W \rightarrow \infty$ PV

hatásos (minimális σ_p -t kapunk)

Többváltozós eset MLR

$$y = p_1 x_1 + p_2 x_2 + p_3 x_3 + \dots + p_m x_m = \underline{X}^T \cdot \underline{p}$$

$$\text{vagy } \hat{y} = \underline{X} \underline{p}$$

$$X = \begin{pmatrix} 1 & x_{11} & \dots & x_{1m} \\ \vdots & \vdots & & \vdots \\ 1 & x_{n1} & \dots & x_{nm} \end{pmatrix}$$

$$N \geq M + 1$$

konstanstható

$N = M + 1$ egyértelmű megoldás, lin. egyenletrendszer $\det(X) \neq 0$
 $\exists y_i \neq 0$

Valóságban $y \approx \underline{X} \underline{p}$ hiba miatt

$$Q = \sum_{i=1}^N (\hat{y}_i - y_i)^2 = (\hat{\underline{y}} - \underline{y})^T (\hat{\underline{y}} - \underline{y})$$

\Rightarrow szil. hiperszil. illeszkedés, legkisebb négyzet

memorizációs levezetés: $\underline{X} \underline{p} \approx \underline{y}$

$$\underline{X}^T \underline{X} \underline{p} = \underline{X}^T \underline{y}$$

$$\underbrace{(\underline{X}^T \underline{X})^{-1}}_E (\underline{X}^T \underline{X}) \underline{p} = (\underline{X}^T \underline{X})^{-1} \underline{X}^T \underline{y}$$

sírgatótt: $\underline{p} = (\underline{X}^T \underline{W} \underline{X})^{-1} \underline{X}^T \underline{W} \underline{y}$ $\underline{w} = \begin{pmatrix} w_1 & 0 \\ 0 & w_n \end{pmatrix}$ $w_i = \frac{1}{\sigma_{y_i}^2}$

Kalibráció: célmodell felhasználása \underline{x} -kevesire, ha csak \underline{y} ismert
 \underline{y} helyett \underline{Y} (több változósú, pl. spektrumcsúcsok, pl. spektrumcsúcsok)

\underline{X} komponens koncentráció, \underline{S} kalibrációs mátrix (reg. koeff. értékek)

modell: $\underline{Y} = \underline{X} \underline{S} + \underline{E}$ 1.) \underline{S} helyett $\hat{\underline{S}} = (\underline{X}^T \underline{X})^{-1} \underline{X}^T \underline{Y}$ kalib. ~~szó~~

2.) $\hat{\underline{X}} = (\hat{\underline{S}} \hat{\underline{S}}^T)^{-1} \hat{\underline{S}} \underline{y}$ alkalmazás
 \uparrow \leftarrow mért
kiseb koncentráció

Miért nem meg sokszor az MLR² kollinearitás

$$p = (X^T X)^{-1} X^T y$$

kovariancia
korrelációs mátrix: $C = \frac{1}{N-1} X_c^T X_c$

← kettős nulla sor
kétváltozós singuláris ~~det(C)~~ $\det(X^T X^{-1}) \approx 0$ $R = \frac{1}{N-1} X_{st}^T X^T$

kétváltozós rögzített cosinusok az objektumok terében

Hogyan kezeljük:

Változó szelekció: - használjuk csak egyet ha gyorsan meg
- átlagoljuk...

Ridge regressió

$p = (X^T X)^{-1} X^T y$ helyett $p = (X^T X + \lambda E)^{-1} X^T y$

λ függvényében mit látunk (ábrákent ell.)

MPE-nél láttuk: $P(p)$ -re előrelátást feltételez,

$$\frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{p-p}{2\sigma^2}}$$

λE -t általában a fixitáshoz regularizációként

nevezik, pl. Tikhonov reg: $(A^T A + \Gamma \Gamma)^{-1}$
 $\Gamma, H \lambda E$

Főkomponens regressió, de ~~ho~~ (PCR)

$$\square = \underbrace{\square + \square + \dots}_{\text{csak a fontosok}}$$

csak a fontosok, kémia: vizsvármérés kell!

PLS - parciális legkisebb négyzetek módszere

25

PCR-ben főkomponensek csak X alapján

PLS X és y vagy X és Y alapján
(PLS1) (PLS, PLS2)

látenis változók (\approx főkomponensek) X és Y terében is

- maximális kovarianciájú pár kiválasztása X-Y-ek közül
- illetékes erre, maradék számolása

eredmény ugyanúgy lineáris, legkisebb négyzetek lineáris regresszió, pluszok spektrum értékei közül csak a kellő

$$X = T \cdot P^T + \overset{\text{hibe, elhanyagolható}}{E} \leftrightarrow Y = UQ^T + \overset{\text{hibe, elhanyagolható}}{F}$$

score loading \uparrow \uparrow hib, elhanyagolható

jól korreláló párok kiválasztása...

MLP további variációi:

polinomiális $(1 \ x \ x^2 \ x^3)$

függvények $y = c_0 + c_1 f_1(x) + c_2 f_2(x) + \dots$

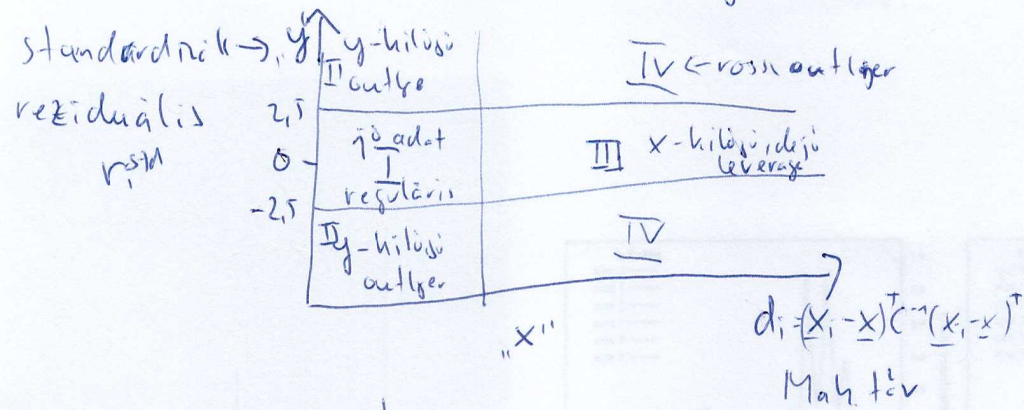
általánosított regr: Mahalanobis-terben

GLS

$$p = (X^T C^{-1} X)^{-1} X^T C^{-1} y$$

\uparrow
súly helyett

Hibás adatokra - robusztus regresszió



- ↓
- úgy dolgozni, hogy csak I és III maradjon meg
 - pár hibás adat eldobása.
 - más becslő fn. (négyzetes til nagy hibánál)

X-y-ban is hibás/hiba vált: (Total, Denominator, ...)

$$x_i \rightarrow \sigma_{x_i}^2 \quad y_i \rightarrow \sigma_{y_i}^2 \quad y = ax + b$$

$$Q(a, b) = \sum_{i=1}^n (y_i - ax_i - b)^2 \cdot w_i = \sum_{i=1}^n \frac{(y_i - ax_i - b)^2}{\sigma_{y_i}^2 + a^2 \sigma_{x_i}^2}$$

$$\text{hossz} \frac{1}{\text{var}(y_i - ax_i - b)} = \frac{1}{\text{var}(y_i) + a^2 \text{var}(x_i)}$$

↑
 igazából ez a jó, de még csak megtört a leírásban

Nonlinear regression ismert függvényalakból:

R7

~ paraméterek becslése

$$f(X, \beta) \quad Q = (y - f(X, \beta))^T W (y - f(X, \beta))$$

\uparrow
 $\begin{pmatrix} f_1(X_1, \beta) \\ f_2(X_2, \beta) \end{pmatrix}$

megoldás J -vel Jacobi mátrix

$$J = \begin{pmatrix} \frac{\partial f(X_1, \beta)}{\partial \beta_1} & \dots & \frac{\partial f(X_1, \beta)}{\partial \beta_m} \\ \vdots & & \vdots \\ \frac{\partial f(X_n, \beta)}{\partial \beta_1} & \dots & \frac{\partial f(X_n, \beta)}{\partial \beta_m} \end{pmatrix}$$

Gauss-Newton

$$\underline{p}^{h+1} = \underline{p}^h + \underbrace{(J^T W J)^{-1}}_{\text{vagy}} J^T W (y - f(X, \underline{p}^h))$$

- Marquardt $(J^T W J)^{-1} \rightarrow (J^T W J + \lambda E)^{-1}$

iteratív, mert nem lineáris, konvergencia nem biztos

Modellzés - elemi kötelesség: modell validálása

adatok felosztása:

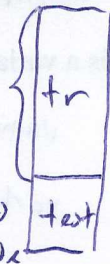
(OECD guidance)
or QSAR

de novo általános
séma

tanító (training)
validáló (validation)
test (test)

MLR

PLS/ANN/...



végső modellnél
használt adatok

← hiperparaméterekhez
PLS komponensek
ANN...

csak ellenőrzés
+ cross validation

+ ez is más lehetőség: val. is mékelt val.



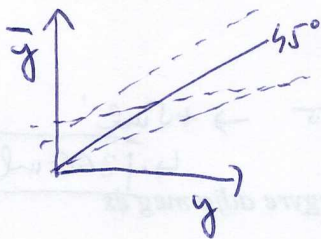
mindig egy validáló

Validációs paraméterek számításánál alapja négyzetösszegek

$$TSS = \sum_{i=1}^{N_{tr}} (y_i - \bar{y})^2 \quad RSS = \sum_{i=1}^{N_{tr}} (y_i - \hat{y}_i)^2 \quad MSS = \sum_{i=1}^{N_{tr}} (\hat{y}_i - \bar{y})^2$$

~~PRESS~~ = "teljes" lineáris regressziónál $TSS = MSS + RSS$, $\bar{y} = \frac{1}{N}$
máskül nem! ← ez okoz kavargást, mely a rengeteg
mutató!

Pearson r - korrelációs e.h. nem igazán jó, csak MLR



← r = 1 mindre
r = -1

r² végképp. de MLR-re jó,
mert $y = \hat{y}$...

átlagos négyzetes eltérés görbe
(root mean square error)

$$RMSE = \sqrt{\frac{RSS}{N}} \quad \left(\begin{array}{l} s = \sqrt{\frac{RSS}{N-1}} \\ 0 \leq RMSE \end{array} \right)$$

determinációs együttható
coefficient of determination

$$R^2 = 1 - \frac{RSS}{TSS} \quad R^2 \leq 1$$

R² = 0 = átlag ugyanolyan jó!

(MLR-nél: $0 \leq R^2 \leq 1$ csatitt)
= $\frac{MSS}{TSS}$ ← modell által
megmagyarázott

Mindseteket lehet mindegyik halmaza

"belső" validáció (modell paramétereinek-hiperparamétereinek meghatározásához használt adatokra)

- illeszkedés jóság (goodness of fit)

RMSE, R^2 ... (Faz MLR-nél, ... korreláció...)

- robusztusság: a használt adatokra mennyire értékes következtetésre (vagy bootstrapping)

$$PRESS = \sum (y_i - \hat{y}_{i(i)})^2$$

↖ ö nem ves részt a modellben épp akkor

$$Q^2 = R^2_{val} = 1 - \frac{PRESS}{TSS}$$

$$RMSE_{val} = \sqrt{\frac{PRESS}{n}}$$

Leave-one-out CV

LOO

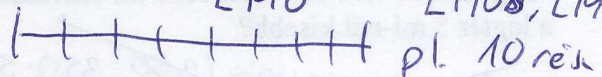
mindenki egyszer marad ki

↘
N modell építése

Leave-many-out CV

LMO

L100, L10, 10%



90%-ra modell (→ 10% használat, mindenki egyszer marad ki)

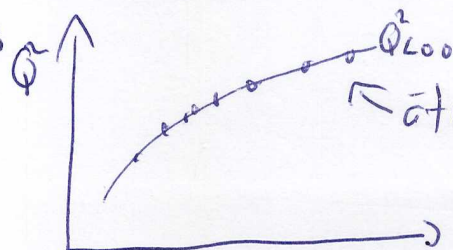
N/10 modell építés

$$-\infty < Q^2 \leq 1$$

$$\text{MLR-nél } -\infty < Q^2 \leq R^2 \leq 1$$

van-e különbség Q^2_{LOO} és Q^2_{LMO} között? Nincs

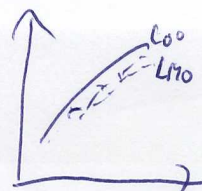
átshálázható Q^2



$$N_{tr, val} = N_{tr} - 1$$

↖ ötleg Q^2 -ek ráhítozódásuk alá

anélkül



"külső" validáció - predikciójoglat

teszt halmaz - ahol se használjuk a modell építésénél

$$RMSE_{\text{test}} (RMSE_{\text{pred}}) = \sqrt{\frac{\sum_{\text{test}} (y_i - \hat{y}_i)^2}{N_{\text{TSS}}}}$$

$$R^2_{\text{test}} = Q_{F2}^2 = 1 - \frac{RSS_{\text{test}}}{TSS_{\text{test}}}$$

Rengegely más mutató van: - pl. négyzetösszeg helyett

variánciával: pl. $(TSS/n-1) \rightarrow$ adjusted R^2 -ek, Q_{FA}^2 , Q_{F3}^2

↑
 sokszor nem nemlétezik, de pl. ezzel
 jobban mérhető, hogy egyszerűs vagy bonyolultabb modellt
 használjunk

- absz. értékes ~~\sum~~ $\rightarrow |y_i - \hat{y}_i|$

- zavarosságok \rightarrow néha félreértés... concordance correlation coefficient (Lin)
 sínessék... R_{011}
 OECD-t én szeretem.

- adott helyen fontos: pl. szabványban: - bármelyik pontnál
 max $\pm 15\%$ eltérés

- $0.95 R_{011}$ $0.99 R_{013}$ kritérium

- visszatérés (visszatérési mérték
 anyagmennyiségre)

pl. $\pm 5\%$

Bias - variance tradeoff \approx kompromissum

Levegethető egyáltalán módon ez a felbontás

$$E(\text{tesztbármely hibája}) = E(\text{variancia}) + E(\text{zajhiba}) + E(\text{torzításhiba})$$

noise

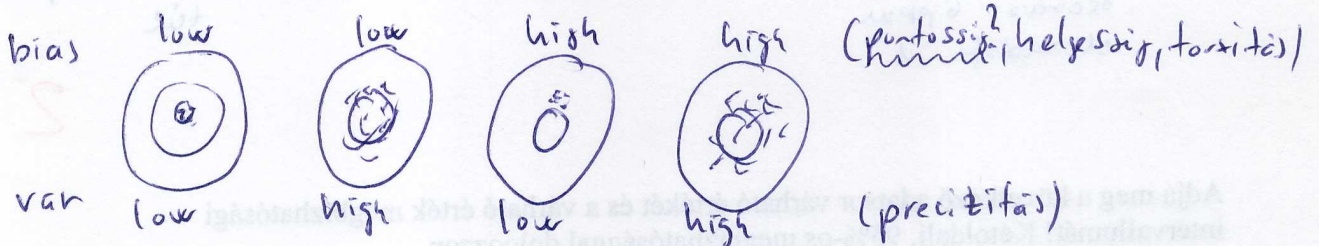
Variancia forrása: mi van, ha más training halmozáson dolgozunk.

Mean square speciális a modellünk (függvények) erre a halmozásra? \leftarrow paraméterek

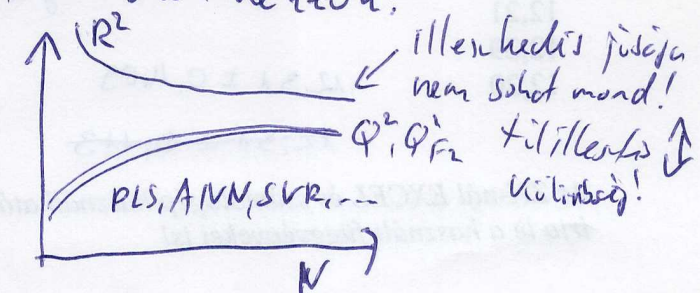
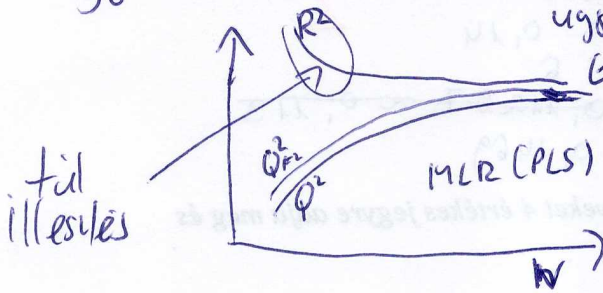
Bias: Mean square tökéletes a modellünk, mert hiányos, még véletlen nagy trainingre is? \leftarrow modellalkotás

Zaj-hiba: Adatok belső zaja és hibája. ~~Ita~~ Prediktorok tökéletlenség-sűrűsége, rosszul mérhető. \leftarrow modellezés előtt

Analitikai kémiai analógia:

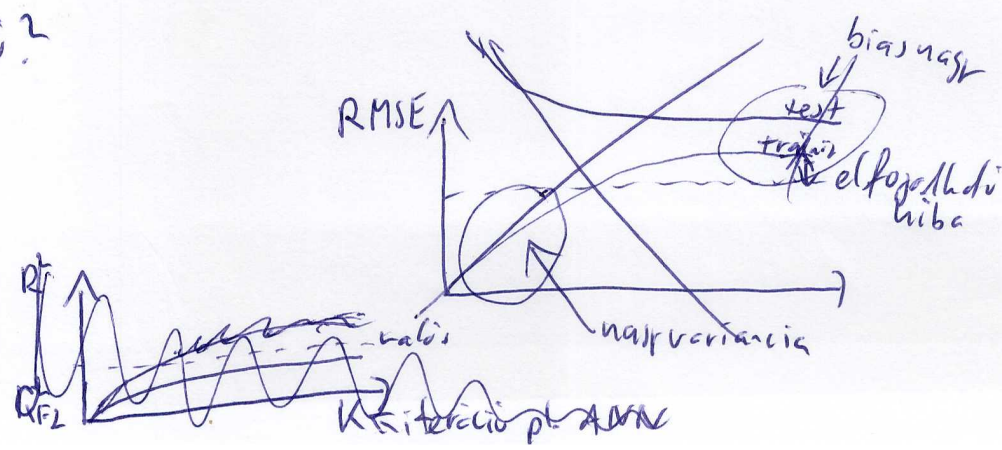


Föl követhető, ha különböző mintasámokban nézzük: (fajti) ugrásoda tart



Mi, hogyan javítható?

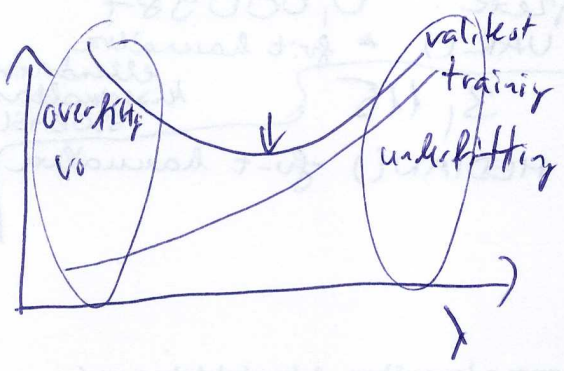
- 1) Variancia
 - nagy mintasám
 - early stop



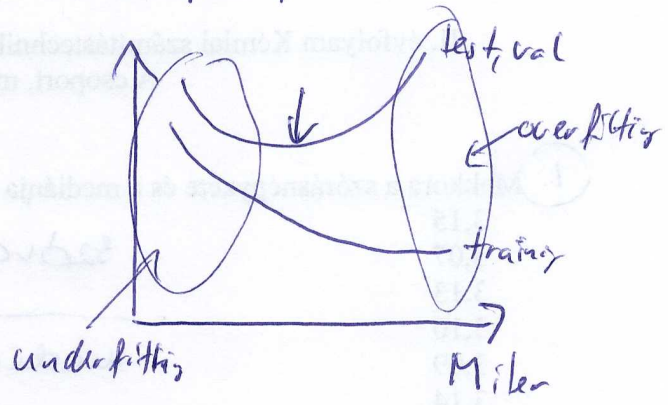
- regularizáció

pl. Ridge, MPE

$$\underline{w} = (X^T X + \lambda E)^{-1} X^T y$$



early stop:



- bagging-konvergens modell

több tanítóhalmaz variánsra illesztési → az eredmény ezek összege

CART → random forest

módszer: bootstrap

Münchhausen bárány...

visztafevéses hivitel $N \rightarrow N$ elemű új minták, sketelésel
 $T_f \approx$ populáció (pedirmiti) ↑ minta

- pl. átlagok számolásá beöle:

nem az összes, csak mondjuk 100-1000

átlagok eloszlása \downarrow 95% konfidenciaintervallum

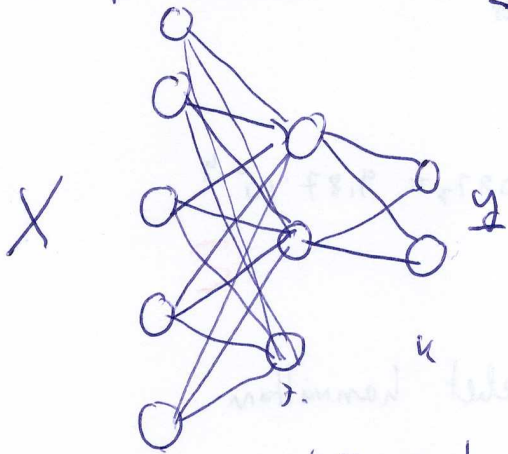
- itt $N_{tr} \rightarrow 100$ db új N_{test} ismétléssel, 100 modell keverése

- mutatónk: out of bag $Q^2 \approx$ robusztusság Konvergens, bagging

Mesterséges idegháló - neurális hálózatok.

Alapmodell:

add. súlyok w_{jk}



egyetleneggy neurónjában (j -dik)

$$\rightarrow \sum w_{ij} x_{ij} + b \rightarrow f(\) \rightarrow$$

n . bemenő i -dik változója

input réteg
rejtett réteg
output

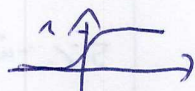
f típusai (ha $f(x) = x$ lineáris)

$$\tanh(u) = \frac{e^u - e^{-u}}{e^u + e^{-u}}$$

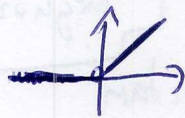


tangens hiperbolikus

$$\text{sigmoid}(u) = \frac{1}{1 + e^{-u}}$$



relu $u < 0$ 0
 $u = 0$ M/A
 $u > 0$ u



elvileg nem szakadós f -k leírhatók egy réteggel, de sok \exists !

használat:

- ismert X, y párok \Rightarrow rengeteg súly optimalizálása, hogy $\hat{y} \approx y \quad \sum (\hat{y} - y)^2$ minimumra

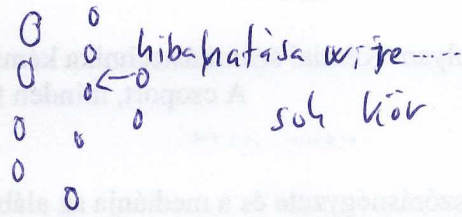
- új X -kre y -ok kiadomlásra

$y \in \mathbb{N}$ osztályozás

$y \in \mathbb{R}$ regresszió

tanulás:

backpropagation alj.



rengetek súly, már itt is $|x| + |y| + |x|k + k = Nw$

N független adat, ritkán így hogy $N > 0,1 Nw$

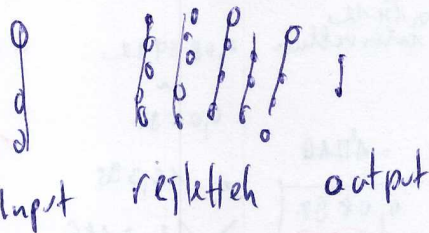
backpropagation \rightarrow alapvetően lokális minimumok
talál meg, helyette + stochasztikus
megoldások (SGD, ADAM)
+ early stop + bagging } min. var.

Mélytanulás, deep learning

Sok réteg:

nagyon összetett dolgokat
tud leírni, de

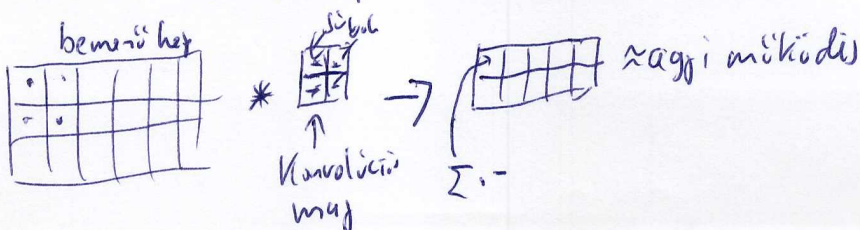
rengeteg optimalizációs súly



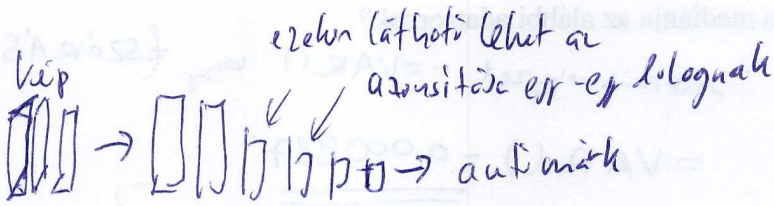
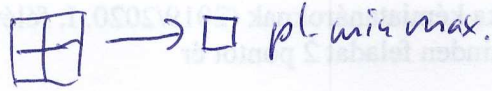
Konvolúciós neurális háló képfeldolgozásra:

kép mérete: $w \times H \times 3$ (szín, zöld, kék)

rengeteg köppont \rightarrow rengeteg súly

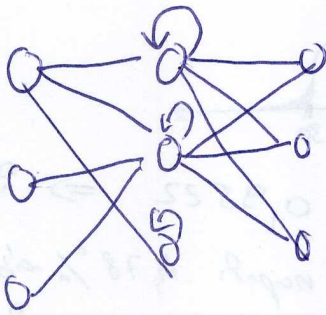


+ pooling

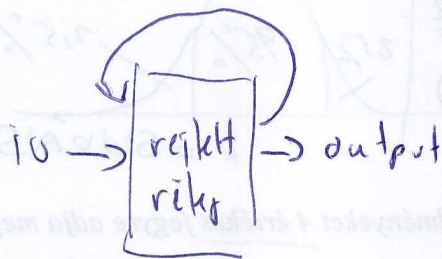


További háló, pl.:

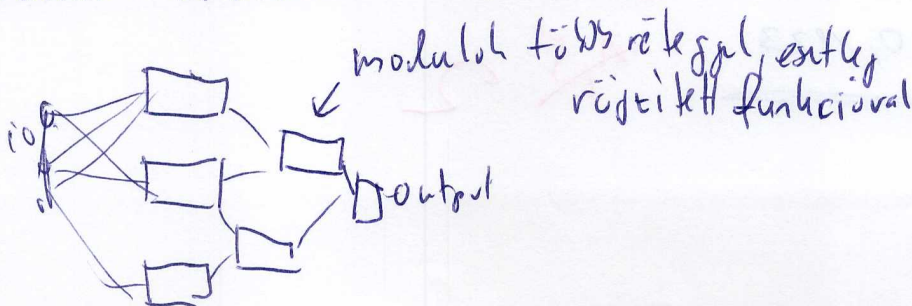
rekurrens háló



long-short term memory



moduláris rendszer



Modellzés Gauss eloszlásokkal:

Gauss-1

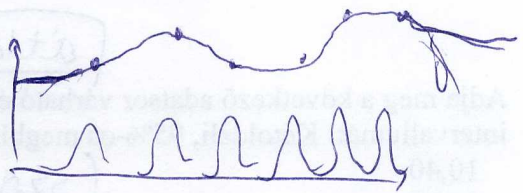
Körváltás K normál módser \rightarrow spam, notspam

távolság függvénye vételeivel:

$$\hat{y} = \sum_{i=1}^{N_{tr}} \frac{K_{in}(x-x_i) y_i}{\sum_{i=1}^{N_{tr}} K_{in}(x-x_i)}$$

↑
súly fu. pl Gauss

lásd hivatott ábra:



problémák többszámítások:

különböző σ -kell (Mah. transz?)

esetlegesen

1 pontot össze N_{tr} -kell (kernel trick)

↓
záró formula)

el(öválogatás)

"Gaussian Process..."

kernel mátrix ~~hat~~ mátrix lin. regnél.
használni lehet
nem lin. ker.

$$\hat{y} = X(X^T X)^{-1} X^T y$$

↑
csak MLR-re