

## Paskaita 2

### Iradas apie tiesinės regresijos metodus

Šioje paskaitoje nagrinėsime gerai žinomą iš statistikos kurso tiesinės regresijos metodus, bet papildomai dėmesį skirsime tiesų naujienis klasifikavimui, kuriuos nagrinėjame MM metode. Mūsų domūs tokie

~~temai~~ punktai:

- a) duomenų skaidymo strategija
- b) apmokyto fikslumo vertinimas
- c) statistikos metodai ir algoritmai, kurie esmingai naudojami MM technologijoje

1. Nagrinėjame klasterinę (standartinę) situaciją, kai efektyviai naudojami MM.

Turime  $m$  eksperimentinių duomenų

$$(x_i, y_i), \quad i = 1, \dots, m,$$

kurie gauti „matuojant“ funkcijos  $f(x)$  reikšmes,  $x \in X$  - yra nepriklausomo kintamojo reikšmių aibė (jūs yra determinuotos, nesusitiktinės),

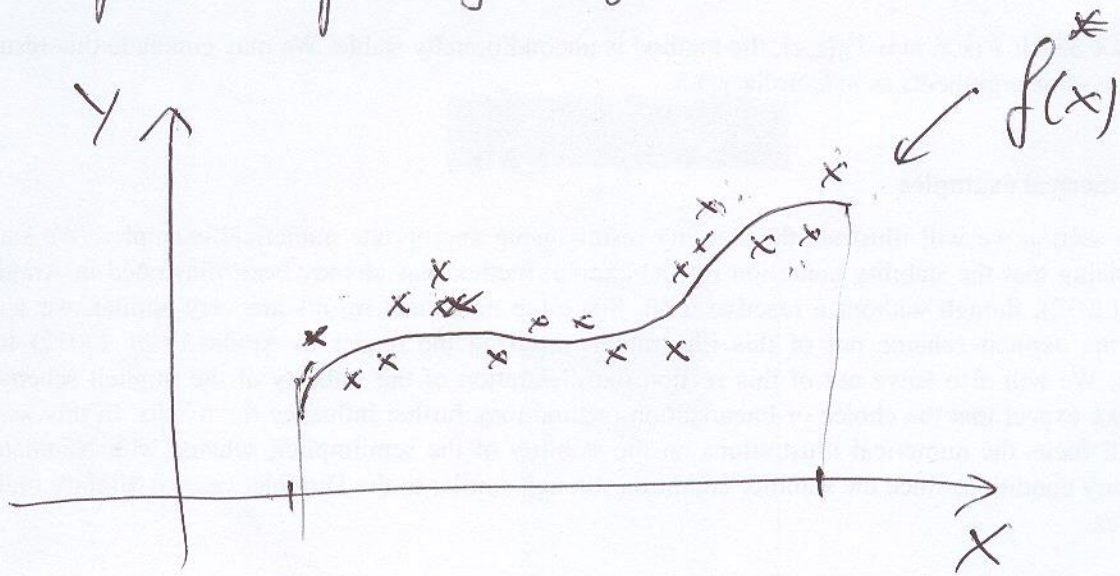
$y \in Y$  yra priklausomo kintamojo reikšmės (jūs yra atsitiktinės, paveiktos „triukšmo“).

Tikslas yra rasti  $\hat{f}$  funkciją  $f(x)$ , kuri geriausiai nusako ryšį tarp duomenų eksperimentinių duomenų

$$Y = f(X) + \varepsilon$$

$\varepsilon$  yra matavimų paklaidai, apie jos tikimybinį pasiskirstymą turime tam tikrą informaciją (matavimo prietaisų charakteristikos).

Tipinis pavyzdys:



Pakleidos apė pakleidas

1.  $\varepsilon_i$  yra normaliai pasiskirstę dydžiai
2. Pakleidy vidurkiu lygūs nuliui  $E\varepsilon_i = 0$
3. ~~Dispersijos~~ yra Pakleidy dispersijos  
yra lygūs  $D\varepsilon_i = \sigma^2$
4. Visos atsitiktinės pakleidos yra nepriklausomos

Toliau esnėje šioi pasakitos neoklaidosje  
pirmuoliuė paprastosi tiesinė regresij  
analizė,

Darome prielaidą (hipotezę), kad eksperimentiniai duomenys gali būti aprašomi tokia modeliu:

$$y_i = \alpha_0 + \alpha_1 x_i + \varepsilon_i$$

Čia  $\alpha_0$  ir  $\alpha_1$  yra parametrai kurių reikšmės yra nežinomos (jas reikia lininiu būdu parasti).

Dėl to, kad  $\varepsilon_i$  yra atsitiktiniai dydžiai, tai net ir turint tą pačią  $x_i$  reikšmę kelis kartus:  $x_i = x_j = \tilde{x}$ ,

$y_i$  ir  $y_j$  reikšmės gali būti skirtingos dėl matavimo paklaidos poveikio  $\varepsilon_i \neq \varepsilon_j$ . Parametrai  $\alpha_0, \alpha_1$  radiniam formuluojame optimizavimo uždavinį

$$\arg \min_{\alpha_0, \alpha_1} \sum_{i=1}^m \left( y_i - (\alpha_0 + \alpha_1 x_i) \right)^2 \quad (1)$$

$J(\alpha_0, \alpha_1)$

(1) uzdevums sprendzot (parametrus  $\hat{\alpha}_0, \hat{\alpha}_1$ ) randamu izteikšimui būdus (2x2 - tiesnesis līdņu sistēmas sprendis)

$$\frac{\partial J}{\partial \alpha_0} = -2 \sum_{i=1}^m (y_i - \alpha_0 - \alpha_1 x_i) = 0$$

$$\frac{\partial J}{\partial \alpha_1} = -2 \sum_{i=1}^m (y_i - \alpha_0 - \alpha_1 x_i) x_i = 0$$

$$\begin{cases} m \hat{\alpha}_0 + \left( \sum_{i=1}^m x_i \right) \hat{\alpha}_1 = \sum_{i=1}^m y_i \\ \left( \sum_{i=1}^m x_i \right) \hat{\alpha}_0 + \left( \sum_{i=1}^m x_i^2 \right) \hat{\alpha}_1 = \sum_{i=1}^m x_i y_i \end{cases}$$

$$\hat{\alpha}_0 = \bar{y} - \bar{x} \hat{\alpha}_1,$$

$$\hat{\alpha}_1 = \frac{\sum_{i=1}^m (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^m (x_i - \bar{x})^2}$$

Mūsu šikslas jau naudoti  
sudaryto funkciję

$$f(x) = d_0 + d_1 x,$$

kai prognozuosime  $Y$  reikšmes naujuose  
tastuose. Tada svarbiausias teorinis  
klausimas - kokia bus šios  
prognozių paklaida

$$\text{error}(x) = |f^*(x) - f(x)|,$$

Kadangi  $f^*(x)$  nežinome, tai šerime  
sudaryti aposteriorinius įverčius,  
leidžiančius ~~prognozuoti~~ gauti informaciją  
apie apmokymo metu gautos funkcijos  
 $f(x)$  šikslumą.

Tam naudojame - Duomenų skaidymą.

## Daomenu skaidyru strategia 1.

Visus dromus eksperimentuuis dromenis dalijame  $2$  dvi ~~ai~~ dalis.

1. Apmokyumi skirti dromeny.

$$(x_i, y_i), \quad i=1, \dots, k$$

2. Testavimui skirti dromeny.

$$(x_i, y_i), \quad i=k+1, \dots, m.$$

Tikslinga sivas arba parinkti, remiantis proporcija 70% : 30% arba 60% : 40%

### Darbo eige.

1. Apmokyimo procese sudarome  $f_j$   
 $f(x)$  (per naudjame tiesinę  
regrasj)  $f = f(x, \alpha)$ .

2. Ivertinami gautosis  $f$ -jos  $f(x)$   
tikslumų  $f = f(x, \alpha)$ .

$$\hat{\epsilon} = \max_{k+1 \leq i \leq m} |y_i - f(x_i)|.$$

Jeigu toks tikslumas yra nepa-  
kaniamas, tai renkames kitą  
funkcijos  $f$  klase - pvz. sudarome  
parabolę  $f(x, \alpha) = \alpha_0 + \alpha_1 x + \alpha_2 x^2$   
arba spline'ą  
 $S(x, \alpha)$ .

Vėl naudojame apmokymui skirtus  
duomenis pirmame etape ir  
testavimui skirtus duomenis, kai  
įvertiname paklaidą



Sabam svarbu susipažinti su  
netrivialia (ir tiek tiek netiketa)  
išvada.

- Pasirenkdamas ypatas  $f^*(x)$  aproksi-  
mavimui funkcijas galime gauti  
efektyviam realizuojamam algoritmui  
(tiek apmokymu stadijoje, tiek ai  
modeliavimo/prognozavimo metu.)

Taigi  $f(x)$  yra paprasta (~~tai~~ pr.  
tiesine  $f(x)$ ), tai jos realizacijos  
kostas yra mažas. Tačiau tokios  
 $f$ -jos tikslumas gali būti nepakan-  
kamas, neefektyviam panaudojame  
turimus eksperimentinius duomenis.  
Tokios situacijos vadiname - nepakan-  
kamu apmokymu.

Apriņai šai gali signalizēt ar pakļaudes apmēkyno metu

$$\tilde{E} = \max_{1 \leq i \leq k} |y_i - f(x_i)|$$

Vēlvar, šai patvrtina ir ~~to~~ fikslumo fikslumas neudojant testavimni skrtus dnoemni

$$\tilde{e} = \max_{k+1 \leq i \leq m} |y_i - f(x_i)|$$

(Dorinai  $\tilde{e} \approx \tilde{E}$  ~~to~~ arba net  $\tilde{e} > \tilde{E}$ ).

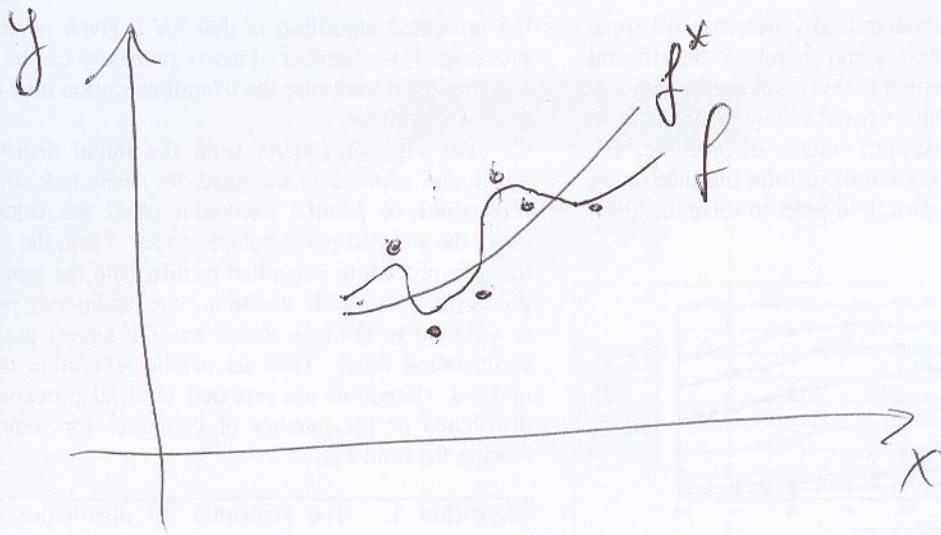
• Pasvinkdceem specialioi struktūros  $f$ -jās (splaineis, racionolis trupnes, Pade apromuwojā) - apmēkyno metu sudarome  $f(x)$   $f$ -jā, kevi lotai gerai apromuwojā apmēkynim skrtus dnoemni, t. y.  $\tilde{e} \ll 1$  (l. mātā pakluda)

Tačiau atlikę testavimą su antrąja duomenų aibe gyvename, kuol paklaida  $\tilde{\epsilon}$  yra didelė.

Tokią situaciją vadiname persukio apmokymo procesu metu.

Tai nesuletinga paaukinti - apmokymo metu per daug tiksliai aproksimuojam eksperimentiniai duomenys  $y_i$ , kurie patys yra nefiksūs (jei išmatuojami su paklaida  $\epsilon$ ).

Reguliarizavimo principas - apmokymo paklaida turi būti tos pačios eilės, kaip ir duomenų nustatymo paklaida.



## Modelio vertinimas

Sukonstruoti modelį (pvz. tiesines regresijos modelį), jo tinkamumą pirmiausia galime įvertinti naudojant žinomus statistinius metodus

(Fischerio ir Stjuarto kriterijus).  
(statistinės vertės charakteristikos)

Apibūdinti tokias dydžius: (n-iesiems skaičiams).

$$SST = \sum_{i=1}^n (y_i - \bar{y})^2$$

(nenormalizuota dispersija).

Normalizacija  $\frac{1}{n}$   
arba patiksl.  $\frac{1}{n-1}$

Regresijos  $f(x)$  ~~to~~ skirtumas nuo vidurkio kvadratinis reikšmė

$$SSR = \sum_{i=1}^n (\hat{y}_i - \bar{y})^2$$

Regresijos paklaidų normos kvadratas

$$SSE = \sum_{i=1}^n e_i^2 = \sum_{i=1}^n SST - SSR$$

### Determinacijos koeficientas

$$r^2 = \frac{SSR}{SST}$$

Jei  $r^2 < 0.25$ , tai reiškia, kad duomenys gali būti apytiksliai fituoti regresijos kreive (tikse)

nėra patikima (nėra tikslaus prognozės)

# Gautojo modelio sąryšis aptardum (santrauka)

1. Modelio sudėtingumui didėjant apmėlymo paklaida mažėja

$$MSE = \sum_{i=1}^n (\hat{y}_i - y_i)^2$$

(modelio reikšmės vs artimėnis eksperimentuotoms duomenims - bet neįmanoma, kad  $y_i$  yra matuojami su matavimo paklaida  $\tau$ ).

Svarbi charakteristika - modelio variacija. Ji nusako, kaip pasikeičia modelio prognozė, pakętes apmėlymo duomenų intę. Labai fikslūs (low variance) modeliams variacija yra didelę. (The Variance).

Antroji charakteristika - modelio  
 poslinkis. Jis parodo, kiek gerai  
 duotasis modelis gali aproksimuoti  
 fikslę  $f$  - jį  $f^*(x)$ . ~~Repus~~ Daskiau-  
 siai fikslėn, suoleitų modeliai (   
 laukitus modeliai) teoriskai gali  
 fikslėn aproksimuoti ~~u testu~~  
 apuskyms duomenis ir fikslę  $f$  - jį  
 $f^*(x)$ . Problem - surisė fiklė su  
 paklauda smatstus duomenis, šidel  
 per fikslėn aproksimuoti apuskyms  
 duomenys gali blogai aproksimuoti  
 $f$  - jį  $f^*(x)$ , o štai lems, kad netikslia  
 prognozuoime y reikšmes nežina-  
 vidus x.