

## Paskaita 3

### Klasifikācijas uzdevings

Nāpzinēdme nēp ī pagrūndīnīz MM  
uzdāvīnīz - dēsmēnīz klasifikācīnīz.

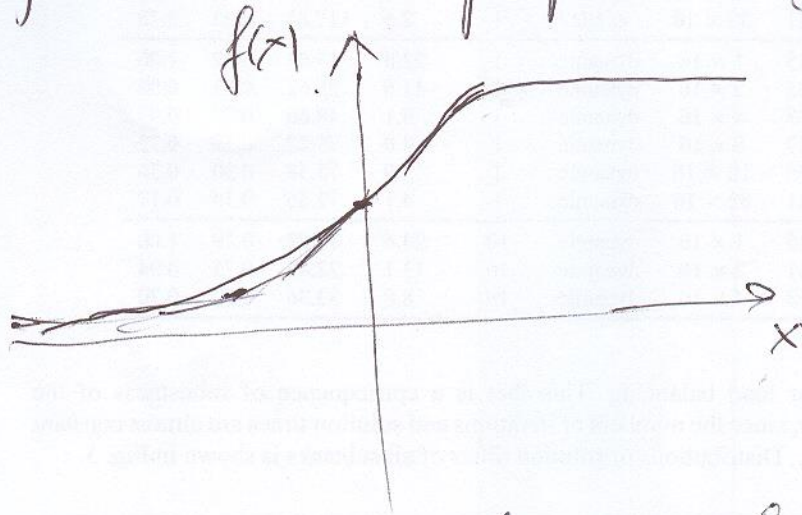
Tāi mokyvosi su mokytojīz algoritmīz  
srotīz ī nārdojāmī parametrvīdāi  
mētdāi fīkslō fūncījīs  $f(x)$  aprōksī-  
māvīnīz.

Tvrdmē dēsmēnīz aībz ī rēkīz  
jūsē gāvtī atsākyms ī klācīnīz  
"tāp / nē", "nārdonē / jūdās",  
(bināriēz uzdāvīdāi) ar āfāsūtī  
"suo / katē / vīstā / vovē" -  $m$ -gālvīnīz  
klasifikācīnīz uzdevīnīz.

Prīmīnsīnē, kad ankscībz nāpzinēdme  
prognozācīnīz uzdevīnīz, kvi  $f(x)$  grā  
solydī  $f$ -jā, dvertīnāntī par srotī, kānāz  
temperatūrīz ī par. Sūspāzīnōmē su fīkslīnīz  
ī rēpīzīnīz mētdāi



Vietoj tiesis naudojame S-formos kreivę, kuri darina netobumų modeliuojant gyvųjų populiacijų dinamiką



Statistinėje aculėje daroma netobumų modelis - sigmoidinė (sigmoid) funkcija

$$S(v) = \frac{1}{1 + e^{-v}}$$

$v$  - nepriklausomo kintamojo reikšmė.

Tada klasifikaicineje modeliuojame gyvųjų tikimybę

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



Taiigi  $p(x) \geq 0.5$  - šai klasifikatorui patvirtina požymį.

### Logistinės regresijos aproksimacija

$$p^x(x) = \frac{e^{\beta_0 + \beta_1 x}}{1 + e^{\beta_0 + \beta_1 x}}$$

Parametrai  $\beta_0, \beta_1$  apibrėžia tiesinę priklausomybę, bet ji transformuojama panaudojant netiesinę funkciją.

Taiigi prognozės rezultatai jau nėra tiesiškai priklausomi nuo  $x$ .



Logistiskās regresijas skaidrojums

$$p(x) = \frac{e^{\beta_0 + \beta_1 x}}{1 + e^{\beta_0 + \beta_1 x}}$$

$$1 - p(x) = 1 - \frac{e^{\beta_0 + \beta_1 x}}{1 + e^{\beta_0 + \beta_1 x}} = \frac{1}{1 + e^{\beta_0 + \beta_1 x}}$$

$$\frac{p(x)}{1 - p(x)} = e^{\beta_0 + \beta_1 x}$$

$$\ln \left( \frac{p(x)}{1 - p(x)} \right) = \beta_0 + \beta_1 x$$

Samais  
odds

↑  
atitulkā lineāras regresijas  
formulā.

$$p = 0.8$$

$$\text{odds} = \frac{0.8}{1 - 0.8} = 4.$$

Koef.  $\beta_0, \beta_1$  var radām  
panaudojam izteiktības  
formulas (ar eksponentu)



Norėdami iš eksperimentinių duomenų surasti logistinę regresijos funkciją naudojame stochastinių gradientinių nusileidimo metodus.

Kaip parodydžiu toliau toliau eksperimentiniuose duomeniuose

$$(x_{1i}, x_{2i}, y_i) \quad i = 1, \dots, m$$

$$y_i = \begin{cases} 0 & \text{binarinė} \\ 1 & \text{klasifikacija} \end{cases}$$

Apriksimacija:

$$y = \frac{1}{1 + \exp(-(\beta_0 + \beta_1 x_1 + \beta_2 x_2))}$$

Jei  $y(x_1, x_2) < 0.5$ , tai turime atvejį 0

Jei  $y(x_1, x_2) > 0.5$ , tai turime atvejį 1.



Imkime paklaides  $f$ -js

$$F(\beta_0, \beta_1, \beta_2; x_1, x_2) = \left( y - \frac{1}{1 + \exp(-\beta_0 - \beta_1 x_1 - \beta_2 x_2)} \right)^2$$

Tad skaitļosim gradienta komponentis

$$\frac{\partial F}{\partial \beta_0} = 2 \left( y - \frac{1}{1 + \exp(\quad)} \right) (-1) \times$$

$$\times \frac{\partial}{\partial \beta_0} \left[ 1 + \exp(-\beta_0 - \beta_1 x_1 - \beta_2 x_2) \right]^{-1}$$

$$= -2 \left( y - \underbrace{\frac{1}{1 + \exp(\quad)}}_{p(\beta_0, \beta_1, \beta_2; x)} \right) \underbrace{\left( 1 + \exp(\quad) \right)^{-2} \cdot \exp(\quad)}_{p \cdot (1-p)}$$

$$\frac{\partial F}{\partial \beta_0} = -2 \left( y - p(\beta_0, \beta_1, \beta_2; x) \right) \cdot p(\cdot) \cdot (1-p(\cdot))$$



Analogiškai apskaičiuojame

$$\frac{\partial F}{\partial \beta_k} = -2(y - p(t)) \cdot p(t) \cdot (1 - p(t)) \cdot x_k,$$

$$k = 1, 2.$$

Gauname iteracinį algoritmą

$$\beta_0^{n+1} = \beta_0^n - \gamma \cdot \frac{\partial F}{\partial \beta_0}(\beta_0^n, \beta_1^n, \beta_2^n; X_1^n, X_2^n)$$

$$\beta_k^{n+1} = \beta_k^n - \gamma \cdot \frac{\partial F}{\partial \beta_k}(\beta_0^n, \beta_1^n, \beta_2^n; X_1^n, X_2^n) \cdot X_k^n,$$
$$k = 1, 2.$$

Tiesinė Diskriminantinė Analizė  
(Linear Discriminant Analysis).

Ankstesniųjų metodų leidžia klasifikuoti  
modelius, kurių surejime tik du  
požymius (klases).

LDA klasifikatorius funkcinė  
tada, kai turime daugiau nei  
2 klases.

LDA modelių atvaizdavimas  
vel yra nesudėtingas. Jis  
priklauso nuo turimų duomenų  
statistinių sąvybių.

Jeigu turime tik vienas įverties  
parametras ( $x$ ), kelių duomenis  
reikia klasifikuoti, tai kiekvienai  
klasei apibrėžiame du parametrus  
(characteristics)

- vidurkis
- dispersija.



Tevgu šķirne daļiņu gēnētiskās  
kontaminācijai, šai skaitļojamē

- vidurķī
- kovariācijas matricā ( šai gra  
daļiņu matricā dispersijas apbēdēt-  
numā ).

Turēdami šīs statistiskās informācijas,  
jā izstrādā  $LDA$  līnija un  
prognozējamā, bērnu klases  
problēma ( $x$ ) - atlikuma klasē  
faktors



Sprendžiant klasifikavimo uždavinį  
LDA metodu darome prielaidas:

- visi kintamieji ( $x$ ) yra pasiskirstę pagal normalųjį pasiskirstymą disus.
- visų kintamųjų ( $x$ ) (jei jis yra kelis) variacijos (dispersijos) yra vienodos.

Todėl patartum duomenis patikrinti ir sufrilyti dar iš LDA algoritmo panaudojimo.

### Algoritmas.

1. ~~Skai~~ Apibrėžiam duomenis skaidinti į kelias grupes priklausomas klasei  $k$  duomenų skaičius  $n_k$ .



$$a_k = \frac{1}{n_k} \sum_{i=1}^{n_k} x_i, \quad x_i - \text{priekšams}$$

klasei k.

2. Skaidrojams visos arbei dispersija (reālamis priekšams, kas visos klasēs deronem, variācijas yca vienoda)

$$\sigma^2 = \frac{1}{n-K} \sum_{i=1}^n (x_i - a_k)^2,$$

čā K yca skirtingų klasių skaidris, is statistikas keruo surime šokį nepasvirulitą dispersijos artrunles formulē



Labar gvertinsime šikimybė, kad nauji duomenys ( $x$ ) priklauso klasei  $k$ . Apšvedusiojame šikimybė  $\neq$  klasei, ir pasirenkame  $\hat{y}$ , kuriam šikimybė yra didžiausia.

Naudosime Bayes teoremą:

$$P(Y=k | X=x) = \frac{P(k) \cdot P(x|k)}{\sum_{j=1}^K P(j) \cdot P(x|j)}$$

- $P(Y=k | X=x)$  yra šikimybė, kad kintamojo  $X$  reikšmė  $x$  apibūdinti klase  $k$ .
- $P(k)$  apibūdinti šikimybė, kad apmokymui skirtose duomenyse  $k$ -klasei šikimybė (dažnis) buvo  $P(k)$ .



eksperimentine

•  $P(x|k)$  yra tikimybi, kad  $x$  priklauso klasei  $k$ .

Įs apskaičiuojame reikšmenis apmokėjimo dėsno menius ir priklauso, kad  $(x)$  yra pasiskirstę pagal normalizę dėsni.

Gauname diskriminantinę funkciją

$$D_k(x) = x \frac{a_k}{\sigma^2} - \frac{a_k^2}{2\sigma^2} + \ln(P(k))$$

Labas jau galime atlikti klasifikavimo sprendimus (žr. Lab.3.txt dėsno menys failų)