

Paskaita 5

K-artimiausių kaimynų
metodas (K-nearest neighbors)

Šis metodas tinka tiek klasifikavimui,
tiek regresijai skaitdamsis atžymėjimams

Pirmas klausimas - koks modelis

(parametrinis ar neparametrinis) yra
naudojamas atlikti sprendimus priėmimus?

Atsakymas: labai naudingas, kas
siekiant sėkmingai modelio apmokymo
etapą (~~ta~~ - apmokymo atlikti
nereikia (inginių algoritmus), o
modelis yra visą ^{apmokymui} ~~testavimui~~ šiek tiek
duomenų arba?!

Sąžada žerū deū galės - ar taip?

Taij modelis sudarytas reikšdama fik
apmokymu ~~testavimo~~ duomenų sutvarkymą, kad
būtų galima efektyviai atlikti skaičia-
nimus (duomenų struktūra) ir jų
pastovų atnaujinimus, bei nereikalingų /
netikslių duomenų pašalinimą (pvz.
beveik dubliuojančių duomenų šalinimą)

KAK metodo algoritmas

- Turime neįėjusius duomenų X .
Reikia rasti juos atitinkančią kreivę.
rijam reikšmę (regresijos uždavinys)
arba nustatyti, keliam klasei x
priskaus (klasifikavimo uždavinys)
- Seiraudame K - artimiausių testa-
vimo duomenų (X_1, X_2, \dots, X_k) .

Pastaba X - gali būti apibestas p
komponentėmis $X = (x_1, x_2, \dots, x_p)$. 2 dimensija
p. 7

Tai sudaringsda algoritmas delis
(kainuynu pavestka) - priduvalute
skalodyh u molyh algoritmu gretimy
tastu pavestka

• Jada regresijis u davinu skaiciojame
K artimauy kainuynu tastu rektimy vidurky (arba
medicauy), klasifikacijis u davinu -
reudame dazindauku pavuotę klase
(cyklo me balsavimy).

Atstumo skaiciojimas

Zinome dany metriky, keuris galime
naudoti skaiciojant atstumus.

Patartim parenkant metriky atditi
analize ("apmelymy"), kur pavuotame
testuojamy tastu aibe, o likusiu tastu
behome "apmelymo" tashais. Trebuome
tę metriky keuri gerdausku prognozeoja
testuojamy tastu rezultatus

Prisiminti populiariausios metrikos.

1. Vektorių Euklidinis atstumas:

$$X_1 = (x_1^1, x_2^1, \dots, x_p^1)$$

$$X_2 = (x_1^2, x_2^2, \dots, x_p^2)$$

$$EA(X_1, X_2) = \|X_1 - X_2\| = \left(\sum_{j=1}^p (x_j^1 - x_j^2)^2 \right)^{1/2}$$

2. Miesty blokinis atstumas (rekomenduojamas, kai taško skirtingos koordinatės apibūdinami skirtingais pagrindiniais parametrais - amžius, lytis, ūgis ir t.t.)

$$MBA(X_1, X_2) = \sum_{j=1}^p |x_j^1 - x_j^2|$$

Parametro K (keimyno skaičius) parinkimas

Šis parametras irgi rekomenduojamas
parinkti „apvalkymo“ būdu - perrenkame
 K reikšmes iš tam skros aibės, pvz
(1, 2, ..., 6).

Pastaba. Jeigu K yra pakankamai
didelis, o deomenų aibė (apvalkymo
tašky aibė) irgi yra didelė, tai
artimiausias keimyno radimo uždav-
nys tampa pakankamai sudėtingu
(reikalauja didelės apimties skaičiavimų)

- Stochastinė paieška (perrenkame apmo-
kyms tašky atsitiktinų porab).
- Lygdauretijs algoritmas