

A04 – Predikcia kvality vína, lineárna regresia pomocou L^1 , L^∞

Piatí proti optimalizácii

Tomáš Antal, Erik Božík, Róbert Kendereš,

Teo Pazera, Andrej Špitalský

2DAV

Január 2024

Predstavenie projektu – lineárna regresia

- ▶ lineárna regresia – predikcia závislej premennej $y \in \mathbb{R}^n$ pomocou nezávislých $x_1, \dots, x_k \in \mathbb{R}^n$

$$\min ||y - \hat{y}||$$

$$\hat{y} = \beta_0 + \beta_1 x_1 + \dots + \beta_k x_k$$

atribúty	x_1	x_2	\dots	x_k	y
pozorovanie 1	1	0.84	\dots	121	4.25
\vdots	\vdots	\vdots	\vdots	\vdots	\vdots
pozorovanie n	4	0.12	\dots	117	5.68

- ▶ vyjadriteľné ako úloha lineárneho programovania – L^1, L^∞

Predstavenie projektu – obsah

- ▶ formulácia LP úloh a dokázanie optimality
- ▶ implementácia v Python-e a predikcia kvality vína
- ▶ počítanie a interpretácia R^2 koeficientu
- ▶ implementácia všeobecnej triedy na počítanie L^1 a L^∞ lineárnej regresie
- ▶ minimalizácia váženej sumy noriem

Formulácia úloh lineárneho programovania

Úloha

Nájsť koeficienty $\beta_0, \beta_1, \dots, \beta_k$ tak, aby predikovaný vektor

$$\hat{y} = \beta_0 + \beta_1 x_1 + \dots + \beta_k x_k \quad (1)$$

bol čo najbližšie k výstupu y , kde y označuje závislú premennú a $x_1, x_2, \dots, x_k \in \mathbb{R}^n$ označujú nezávislé premenné. Túto vzdialenosť $|y - \hat{y}|$ sme minimalizovali l_1 a l_∞ normami

Minimalizovanie l_1 normy

Chceme minimalizovať normu $\|y - \hat{y}\|_1$

označíme:

$$\begin{aligned} A &:= (1_n, x_1, \dots, x_k) \\ \beta &:= (\beta_0, \beta_1, \dots, \beta_k)^T \end{aligned} \tag{2}$$

Problém prevedieme do tvaru:

$$\begin{aligned} \min \quad & c^T x \\ & Ax \geq b \end{aligned}$$

Zavedieme nový vektor $t \in \mathbb{R}^n$, ktorým ohraničíme $y - A\beta$

Minimalizovanie l_1 normy ako úloha lineárneho programovania:

$$\begin{aligned} \min \quad & \left(0_{k+1}^T \mid 1_n^T \right) \begin{pmatrix} \beta \\ t \end{pmatrix} \\ & \begin{pmatrix} A & \mathbb{I}_n \\ -A & \mathbb{I}_n \end{pmatrix} \begin{pmatrix} \beta \\ t \end{pmatrix} \geq \begin{pmatrix} y \\ -y \end{pmatrix} \\ & \beta \in \mathbb{R}^{k+1}, \quad t \geq 0_n \end{aligned}$$

Minimalizovanie l_∞ normy

Chceme minimalizovať normu $\|y - \hat{y}\|_\infty$

Zavedieme skalárnu premennú $\gamma \in \mathbb{R}$, prevedieme na úlohu LP

$$-\gamma \mathbf{1}_n \leq y - A\beta \leq \gamma \mathbf{1}_n$$

Pomocou značenia z (2), výsledná úloha:

$$\begin{aligned} \min \quad & \left(0_{k+1}^T \mid \mathbf{1} \right) \begin{pmatrix} \beta \\ \gamma \end{pmatrix} \\ & \left(\begin{array}{c|c} A & \mathbf{1}_n \\ \hline -A & \mathbf{1}_n \end{array} \right) \begin{pmatrix} \beta \\ \gamma \end{pmatrix} \geq \begin{pmatrix} y \\ -y \end{pmatrix} \\ & \beta \in \mathbb{R}^{k+1}, \gamma \geq 0 \end{aligned}$$

B

...

Predikcia kvality vína

dáta o víne

- ▶ množstvo dážďa v zime
- ▶ priemerná teplota počas zretia vína
- ▶ množstvo dažďa počas zberu
- ▶ vek vína
- ▶ populácia Francúzska
- ▶ cena

Orley Ashenfelter



Výsledky predikcie

L^1

- ▶ + vplyv teplota počas zretia
- ▶ + vplyv veku vína
- ▶ - vplyv dážď počas zberu
- ▶ + vplyv dážď počas zimy
- ▶ - vplyv populácie Francúzska

L^∞

- ▶ rovnaké poradie ako L^1
- ▶ ale - vplyv veku vína

D

...

Nadstavba

```
from models.models import L1Model, LInfModel
```

- ▶ Zovšeobecnenie problému
- ▶ Voľnosť dimenzionality
- ▶ Vstupný vektor y a matica X
- ▶ Hodnoty β výstupom

```
# inicializacia
model1 = L1Model(y, X)
model2 = LInfModel(y, X)

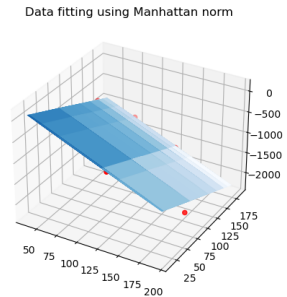
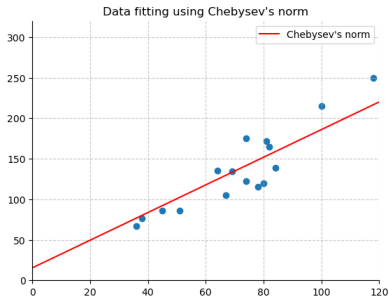
# riesenie
beta1 = model1.solve()
beta2 = model2.solve()
```

Nadstavba

- ▶ Hodnota R^2
- ▶ Vizualizácia pre 2D a 3D

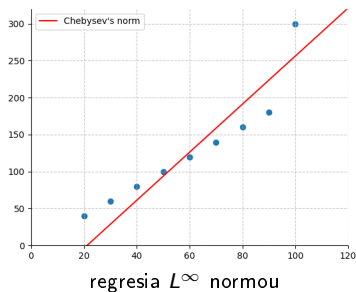
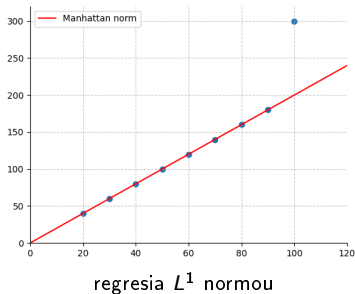
```
model.r2()
```

```
model.visualize()
```



Porovnanie L^1 a L^∞ lineárnej regresie

- ▶ L^1 – veľmi dobre zachytáva lineárny vzťah, môže viesť k *overfittingu*
- ▶ L^∞ – príliš ovplyvňovaná outliermi



Minimalizácia váženého súčtu noriem

- ▶ redukcia *overfittingu* L^1 regresie váženým súčtom s L^∞ normou

$$\min \omega \|y - \hat{y}\|_1 + (1 - \omega) \|y - \hat{y}\|_\infty, \quad \omega \in [0; 1]$$

- ▶ stále implementovateľné ako úloha lineárneho programovania

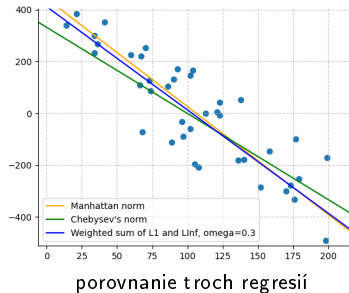
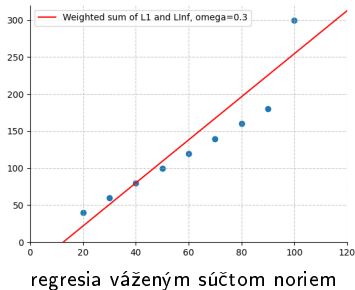
$$\min \left(0_{k+1}^T \mid \omega 1_n^T \mid (1 - \omega) \right) \begin{pmatrix} \beta \\ t \\ \gamma \end{pmatrix}, \quad \omega \in [0; 1]$$

$$\left(\begin{array}{c|c|c} A & \mathbb{I}_n & 0_n \\ \hline -A & \mathbb{I}_n & 0_n \\ \hline A & 0_{n \times n} & 1_n \\ \hline -A & 0_{n \times n} & 1_n \end{array} \right) \begin{pmatrix} \beta \\ t \\ \gamma \end{pmatrix} \geq \begin{pmatrix} y \\ -y \\ y \\ -y \end{pmatrix}$$

$$\beta \in \mathbb{R}^{k+1}, \quad t \geq 0_n, \quad \gamma \geq 0$$

Minimalizácia váženého súčtu noriem

► implementované ako `WeightedL1LInfModel`



Zhrnutie

- ▶ formulácia lineárnej regresie ako úlohy LP
- ▶ predikcia ceny vín
- ▶ jednoduchý framework na počítanie lineárnej regresie pomocou L^1 a L^∞ normami, resp. ich váženej sumy

Ďalšie kroky

- ▶ analýza časovej zložitosti, napr. voči najmenším štvorcám
- ▶ porovnanie vhodnosti jednotlivých prístupov podľa vstupných dát