

# Mašinsko učenje - metoda najstrmijeg spusta, generalizacija i regularizacija

Tatjana Jakšić Krüger

[tatjana@turing.mi.sanu.ac.rs](mailto:tatjana@turing.mi.sanu.ac.rs)

- Linearni model i estimatori.
- Metoda maksimalne verodostojnosti.
- Uvod u metodu najstrmijeg spusta.

## Cilj za danas



- Nastavak za metodu najstrmijeg spusta.
- Izbor modela i generalizacija.
- Regularizacija.

# Osnovna terminologija

## □ Notacija:

$x$  ulaz, svojstvo (*eng* feature, atribut).

$y$  izlaz, ciljna promenljiva.

$n$  broj atributa/svojstava.

$m$  broj primera/instanci.

$(x^{(i)}, y^{(i)})$   $i$ -ti primer iz skupa za obučavanje.

$\{x^{(i)}, y^{(i)}\}_{i=1}^m$  skup za obučavanje (trening skup).

$h(x) = y$  hipoteza (model).

$\epsilon_i | e_i$  pravi|empirijski reziduali.

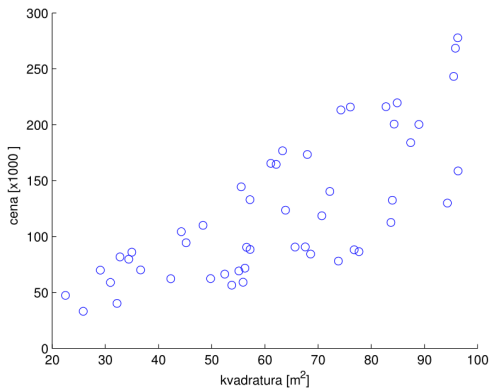
## □ Ulaz može biti:

□ vektor  $\vec{x} = \mathbf{x} = [x_1 \ x_2 \ \dots \ x_n]^T$ .

□ matrica dimenzije  $m \times n$ .

$$X = \begin{bmatrix} x_{11} & x_{12} & \dots & x_{1n} \\ x_{21} & x_{22} & \dots & x_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ x_{m1} & x_{m2} & \dots & x_{mn} \end{bmatrix}$$

# Problem predviđanja - model regresije



Zadatak: predvidi cenu na osnovu date kvadrature.

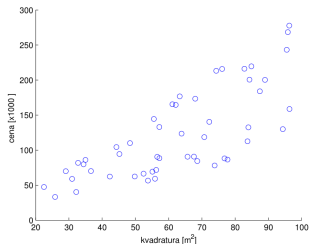
Izvor grafika:

[http://automatika.etf.bg.ac.rs/images/FAJLOVI\\_srpski/predmeti/master\\_studije/MU/01%20Linearna%20regresija.pdf](http://automatika.etf.bg.ac.rs/images/FAJLOVI_srpski/predmeti/master_studije/MU/01%20Linearna%20regresija.pdf)

## Hipoteza u $\mathbb{R}^2$

- Ulazna promenljiva  $x$  je skalarna veličina.
- Ciljna promenljiva  $y$  je skalarna veličina.
- Parametri linearnog modela su  $\theta_0$  (odsečak na  $y$ -osi) i  $\theta_1$  (koeficijent pravca).
- Funkcija hipoteze (model) našeg predviđanja je oblika:

$$h(x) = \theta_0 + \theta_1 x.$$



Zadatak: predvidi cenu na osnovu date kvadrature.

## Hipoteza u $\mathbb{R}^3$

- Ulazne promenljive  $x_1$  i  $x_2$  definišu vektor:

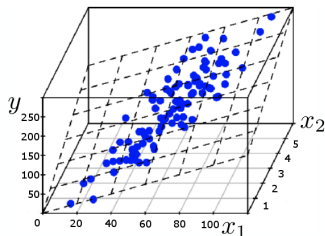
$$\mathbf{x} = [x_0 \ x_1 \ x_2]^T, \quad x_0 = 1$$

- Ciljna promenljiva  $y$  je skalarna veličina.
- Parametri linearnog modela su  $\theta_0, \theta_1, \theta_2$ .
- Hipoteza (model) našeg predviđanja je:

$$\begin{aligned} h(\mathbf{x}; \boldsymbol{\theta}, \theta_0) &= \theta_0 + \theta_1 x_1 + \theta_2 x_2 \\ &= \sum_{i=0}^3 \theta_i x_i, \end{aligned}$$

pri čemu parametri modela čine vektor

$$\boldsymbol{\theta} = [\theta_0 \ \theta_1 \ \theta_2]^T.$$



Zadatak: predvidi cenu na osnovu kvadrature ( $x_1$ ) i broja spavaćih soba ( $x_2$ ).

## Postavka problema u $\mathbb{R}^{n+1}$

- Vektor svojstava je:

$$\mathbf{x} = [x_0 \ x_1 \ x_2 \ \dots \ x_n]^T, \quad x_0 = 1.$$

- Vektor parametara je:

$$\boldsymbol{\theta} = [\theta_0 \ \theta_1 \ \theta_2 \ \dots \ \theta_n]^T.$$

- Usvajamo linearnu hipotezu:

$$\begin{aligned} h(\mathbf{x}; \boldsymbol{\theta}, \theta_0) &= h_{\boldsymbol{\theta}}(\mathbf{x}) = \theta_0 + \theta_1 x_1 + \dots + \theta_n x_n. \\ &= \boldsymbol{\theta}^T \mathbf{x}. \end{aligned}$$

- Kako biramo  $\boldsymbol{\theta}$ ?

$$h_{\boldsymbol{\theta}}(\mathbf{x}) \approx y,$$

tj. takvo da je rizik najmanji:

$$\min_{\boldsymbol{\theta}} \frac{1}{2} \sum_{i=1}^m (h_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}) - y^{(i)})^2 = \min_{\boldsymbol{\theta}} J(\boldsymbol{\theta}).$$



- 1 Metoda najstrmijeg spusta (*eng.* gradient descent).
- 2 "Šaržni" gradijentni spust (*eng.* batch gradient descent).
- 3 Stohastički gradijentni spust.
- 4 Lokalno ponderisana linearna regresija.
- 5 Njutnova metoda.

## Parcijalni izvod u $\mathbb{R}^2$

- Neka je  $f = f(x, y)$ .
- Parcijalni izvod po  $x$  u nekoj tački  $(x_0, y_0)$ :

$$\frac{\partial f(x_0, y_0)}{\partial x} = \lim_{l \rightarrow 0} \frac{f(x_0 + l, y_0) - f(x_0, y_0)}{l} = \partial_x f(x_0, y_0)$$

- Parcijalni izvod po  $y$  u nekoj tački  $(x_0, y_0)$ :

$$\frac{\partial f(x_0, y_0)}{\partial y} = \lim_{l \rightarrow 0} \frac{f(x_0, y_0 + l) - f(x_0, y_0)}{l} = \partial_y f(x_0, y_0)$$

- Ukoliko  $\partial_x f(x_0, y_0)$  postoji, funkcija  $f$  je diferencijabilna po  $x$  u tački  $(x_0, y_0)$ .

## Parijalni izvod i gradijent u $\mathbb{R}^d$

- Neka je  $f : \mathbb{R}^d \rightarrow \mathbb{R}$ .
- Neka je  $f$  definisana i diferencijabilna u okolini  $U_a$  tačke  $\vec{a} = (a_1, a_2, \dots, a_d)$ .
- *Gradijent* je vektor  $\nabla f$  čije su koordinate:

$$\nabla f(a_1, \dots, a_d) = \left( \frac{\partial f}{\partial x_1}(a_1, \dots, a_d), \dots, \frac{\partial f}{\partial x_n}(a_1, \dots, a_d) \right)$$

$$\nabla f(\mathbf{a}) = \left( \frac{\partial f}{\partial x_1}(\mathbf{a}), \dots, \frac{\partial f}{\partial x_n}(\mathbf{a}) \right).$$

- Gradijent u tački  $\mathbf{a} \in \mathbb{R}^d$  predstavlja vektor u čijem smeru funkcija  $f$  najbrže raste/opada u okolini tačke  $\mathbf{a}$ .
- Gradijent se koristi za određivanje maksimalne/minimalne vrednosti  $f$  u okolini tačke  $\mathbf{a}$ .

## Metoda najstrmijeg spusta

---

**Input:**  $f$ ,  $\mathbf{a}_0$ : analitički izraz sa  $f$  i polazna tačka

**Output:**  $\nabla f(\mathbf{a}_k)$

```
1 repeat
2   | Izračunaj  $\nabla f(\mathbf{a}_k)$ ;
3   | Pomeri se ka tački  $\mathbf{a}_{k+1}$  u smeru gradijenta  $-\nabla f(\mathbf{a}_k)$ ;
4   |  $k = k + 1$ ;
5 until  $|f(\mathbf{a}_k - f'(\mathbf{a}_{k-1}))| \leq \epsilon |f(\mathbf{a}_{k-1})|$ ;
```

---

- Nedostaje nam korak koji kaže za koliko moramo da se pomeramo.
- Parametar  $\lambda$  pomaže da naučimo dužinu pomeraja iz  $\mathbf{a}_k$  u  $\mathbf{a}_{k+1}$ :

$$\mathbf{a}_{k+1} = \mathbf{a}_k - \lambda \nabla f(\mathbf{a}_k).$$

- Izbor vrednosti za  $\lambda$  može biti određen:
  - teorijski:  $\sum_{i=0}^{\infty} \lambda_i = \infty$  i  $\sum_{i=0}^{\infty} \lambda_i^2 < \infty$ .
  - empirijski:  $\lambda = 0.01$ .

## Gradijentna metoda za određivanje $\theta_k$

- Neka je  $\theta_k$   $k$ -ta komponenta vektora  $\theta$  vektor u  $\mathbb{R}^{n+1}$  prostoru.
- Primenom gradijentnog spusta  $\theta_k$  ažuriramo:

$$\theta_k := \theta_k - \lambda \frac{\partial J(\theta)}{\partial \theta_k}.$$

- Parcijalni izvod od  $J(\theta)$  po  $\theta_k$ :

$$\begin{aligned} \frac{\partial J(\theta)}{\partial \theta_k} &= \frac{\partial}{\partial \theta_k} \left( \frac{1}{2} \sum_{i=1}^m (h_{\theta}(\mathbf{x}^{(i)}) - y^{(i)})^2 \right) \\ &= \frac{\partial}{\partial \theta_k} \left( \frac{1}{2} \sum_{i=1}^m \left( \sum_{j=0}^n \theta_j x_j^{(i)} - y^{(i)} \right)^2 \right) \\ &= \frac{\partial}{\partial \theta_k} \left( \frac{1}{2} \sum_{i=1}^m (\theta_k x_k^{(i)} - y^{(i)})^2 \right) \\ &= \sum_{i=1}^m (\theta_k x_k^{(i)} - y^{(i)}) x_k^{(i)}. \end{aligned} \tag{1}$$

## Grafik funkcije $J(\theta)$



prikazani gif je moguće pokrenuti u Adobe Reader pdf editoru.

## Minimizacija nad celim skupom za obučavanje za izračunavanje $\theta$

- Imamo matricu svojstava  $X$ , dimenzije  $m \times n$ :

$$X = \begin{bmatrix} x_1^{(1)} & x_2^{(1)} & \dots & x_n^{(1)} \\ x_1^{(2)} & x_2^{(2)} & \dots & x_n^{(2)} \\ \vdots & \vdots & \ddots & \vdots \\ x_1^{(m)} & x_2^{(m)} & \dots & x_n^{(m)} \end{bmatrix}$$

- Neka je  $\mathbf{y}$   $m$ -dimenzioni vektor koji sadrži ciljne vrednosti (labele) iz skupa za obučavanje:

$$\mathbf{y} = \begin{bmatrix} y^{(1)} \\ y^{(2)} \\ \vdots \\ y^{(m)} \end{bmatrix} = [y^{(1)} \ y^{(2)} \ \dots \ y^{(m)}]^T$$

- Koristićemo sledeću osobinu:

$$h_{\theta}(\mathbf{x}^{(i)}) = \sum_{j=0}^n \theta_j x_j^{(i)} \iff h_{\theta}(\mathbf{x}^{(i)}) = (\mathbf{x}^{(i)})^T (\boldsymbol{\theta})$$

## Minimizacija nad celim skupom za obučavanje

$$\mathbf{X}\boldsymbol{\theta} - \mathbf{y} = \begin{bmatrix} (\mathbf{x}^{(1)})^T \boldsymbol{\theta} \\ \vdots \\ (\mathbf{x}^{(m)})^T \boldsymbol{\theta} \end{bmatrix} - \begin{bmatrix} y^{(1)} \\ \vdots \\ y^{(m)} \end{bmatrix} = \begin{bmatrix} (\mathbf{x}^{(1)})^T \boldsymbol{\theta} - y^{(1)} \\ \vdots \\ (\mathbf{x}^{(m)})^T \boldsymbol{\theta} - y^{(m)} \end{bmatrix}$$

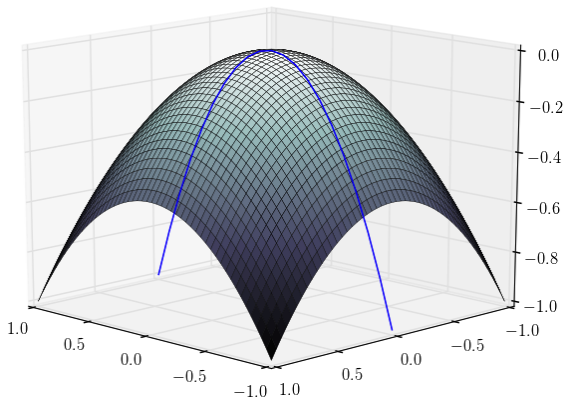
□ Na osnovu osobine  $\mathbf{z}$  važi  $\mathbf{z}^T \mathbf{z} = \sum_i z_i^2$  dobijamo.

$$\begin{aligned} \frac{1}{2}(\mathbf{X}\boldsymbol{\theta} - \mathbf{y})^T(\mathbf{X}\boldsymbol{\theta} - \mathbf{y}) &= \frac{1}{2} \sum_{i=1}^m ((\mathbf{x}^{(i)})^T \boldsymbol{\theta} - y^{(i)})^2 \\ &= \sum_{i=1}^m ((h(\mathbf{x}^{(i)}) - y^{(i)})^2 \\ &= J(\boldsymbol{\theta}). \end{aligned} \tag{2}$$

□ Cilj je minimizacija ovako zapisane funkcije cilja  $J(\boldsymbol{\theta})$ .



## Primeri kvadratne površi bez minimalne vrednosti



Konkavna površ

## Globalni minimum u jednom koraku

□ Dimenzija gradijenta je iste dimenzije kao i vektor  $\theta$ .

□ Pomoćne formule:

$$(AB)^T = B^T A^T$$

$$\text{tr}A = \text{tr}A^T$$

$$\nabla_A \text{tr}ABA^T C = CAB + C^T AB^T$$

□ dovode do sledećeg:

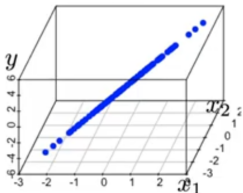
$$\nabla J(\theta) = \nabla \frac{1}{2} (X\theta - y)^T (X\theta - y) = \dots = X^T X\theta - X^T y. \quad (3)$$

□ Globalni optimum nalazimo kada je  $\nabla J(\theta) = 0$  tj:

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

Kada formula za  $\theta$  ne radi?

- Ponekada ne postoji jedinstvena hiperravan.

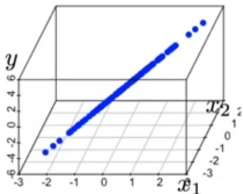


Podaci čine liniju.

- Ponekada usled šuma deluje da imamo dobro rešenje, ali ono nije značajno bolje od drugih.
- Visoko korelisani atributi.
- Preveliki broj atributa.

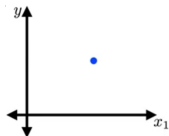
Kada formula za  $\theta$  ne radi?

- Ponekada ne postoji jedinstvena hiperravan.



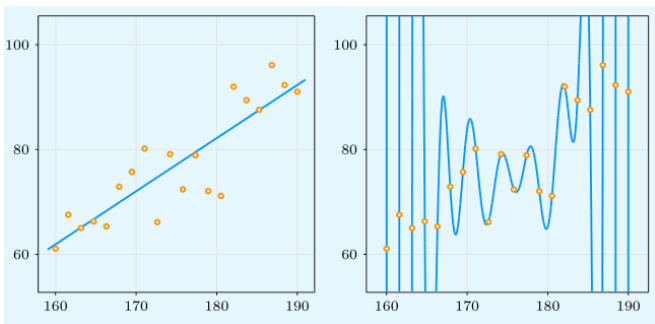
Podaci čine liniju.

- Ponekada usled šuma deluje da imamo dobro rešenje, ali ono nije značajno bolje od drugih.
- Visoko korelisani atributi.
- Preveliki broj atributa.



## Kada minimizacija ne radi? Preprilagođavanje. Generalizacija.

- Slika levo polinom prvog stepena. Slika desno je polinom višeg stepena. Deluje da druga slika bolje fituje jer je greška nula.
- Preprilagođavanje (*eng.* overfitting) je fenomen koji se javlja kada sa složnošću modela raste nepouzdanost.



- *Generalizacija* je mera koliko dobro će istrenirani model predvideti nove rezultate.

## Složenost modela

- Složenost odgovara na pitanje koliko je dobro imati veliku ili premalu složenost.
- Moguće je da je složenost manja od instinske funkcije - *eng.* underfitting.



## Regularizacija.

- Analiza atributa, potraga za najuticajnijim svojstvima.
- Regularizacija - smanjivanje složenosti modela.
- Poznati načini regularizacije kod linearnih modela je penalizovanje:
  - Lasso regression.
  - Ridge regression.
- Ridž regresija penalizuje na taj način da neki parametri dobijaju vrednosti bliske nuli.
- Koristi se  $l_2$  norma (euklidsko rastojanje)

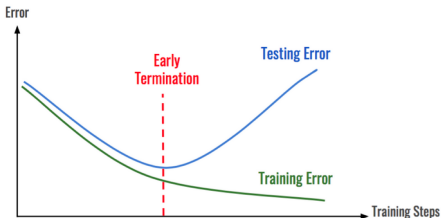
$$J_{ridge}(\theta) = \frac{1}{2} \sum_{i=1}^m ((x^{(i)})^T - y^{(i)})^2 + \lambda \|\theta\|^2, \quad \lambda > 0.$$

- Ukoliko želimo da neke komponente vektora budu nula, koristimo Lasso regresiju koja se zasniva na  $l_1$ -meri.

$$J_{lasso}(\theta) = \frac{1}{2} \sum_{i=1}^m ((x^{(i)})^T - y^{(i)})^2 + \lambda \sum_{i=1}^p |\theta_i|.$$

## "Early stopping" regularizacija

- Regularizacija u vidu ranog zaustavljanja (*eng.* early stopping).
- Obučavanje se zaustavlja pre nego što se javi problem sa prilagođavanjem.



- Neophodno je balansiranje između:
  - 1 Složenosti modela (hipoteze).
  - 2 Veličine skupa za obučavanje.
  - 3 Greške generalizacije nad skupom za testiranje.



- Predavanja prof. Predraga Tadića,  
<http://automatika.etf.bg.ac.rs/sr/13m051mu>.

U nastavku...



- Podela skupa na 3 dela (obučavanje, validacija, testiranjej).
- Cross-validation.



Hvala na pažnji.

Sada idemo na prikaz  
komandi.

Da li imate pitanja?