

Prof.dr.sc. Bojana Dalbello Bašić

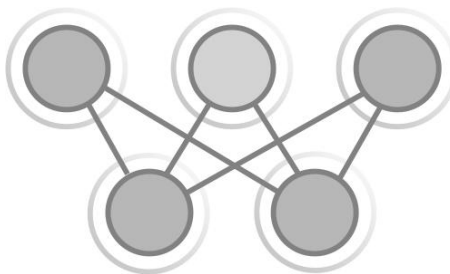
Fakultet elektrotehnike i računarstva
Zavod za elektroniku, mikroelektroniku, računalne i inteligentne sustave

www.zemris.fer.hr/~bojana
bojana.dalbello@fer.hr

Strojno učenje

Bayesova teorija odlučivanja i parametarske metode

Željko Juretić

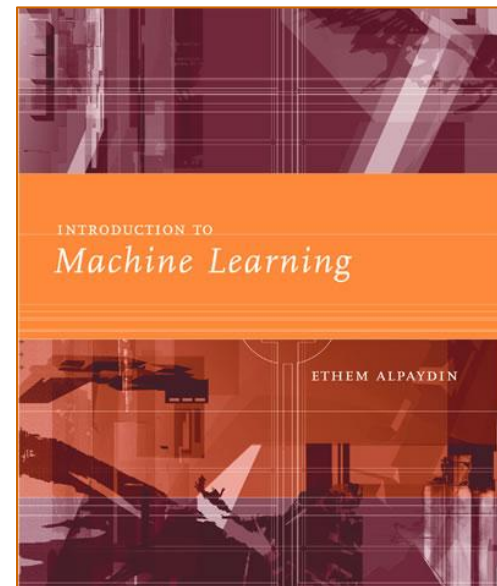


Bayesova teorija odlučivanja (Chapter 3)

- Uvod
- Vjerojatnost i zaključivanje
- Klasifikacija
- Rizici i gubitci
- Funkcije odluke
- Teorija korisnosti
- Vrijednost informacije

Parametarske metode (Chapter 4)

- Uvod
- Kriterij najveće izglednosti (*Maximum Likelihood Estimation*)
- Pristranost i varijanca
- Bayesov estimator
- Parametarska klasifikacija
- Regresija
- Ugađanje složenosti modela: dvojba pristranost/varijanca
- Postupci odabira modela



Bayesova teorija odlučivanja

Ethem Alpaydin,
Introduction to Machine Learning:
Chapter 3: *Bayesian Decision Theory*

T. Bayes.



Thomas Bayes, (1702.-1761.?)

- Donošenje odluka pri nesigurnosti ima dugu povijest: Zvijezde, kristalne kugle, vidovnjaci, ...
- Teorija vjerojatnosti je stara tek nekoliko stotina godina, a nastala je kao pokušaj analiziranja igara na sreću: Gerolamo Cardano, Pierre de Fermat, Blaise Pascal, (17. stoljeće).
- U današnje vrijeme spoj statistike i računarske znanosti omogućuje nam zaključivanje na temelju dostupnih podataka (*data mining, automatic classification, itd.*)

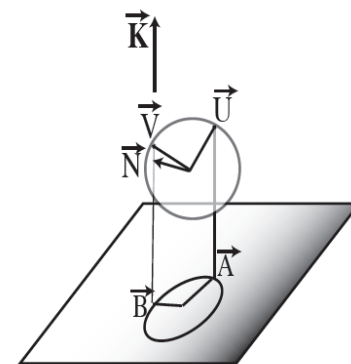


- Vidljive i nevidljive varijable (engl. *observable and unobservable*)
- Rezultat pokusa bacanja novčića: $\in \{pismo, glava\}$
- Slučajna varijabla: $X \in \{1, 0\}$
- Bernoullijev pokus: $P\{X = 1\} = p_0^x (1 - p_0)^{(1-x)}$
- Uzorak: $X = \{x^t\}_{t=1}^N$

- Procjena: $p_0 = \frac{\#\{pisama\}}{\#\{bacanja\}} = \frac{\sum_t x^t}{N}$

- Predviđanje ishoda sljedećeg bacanja:

$$\begin{cases} pismo, & \text{ako } p_0 > \frac{1}{2} \\ glava, & \text{inače} \end{cases}$$



Simulacija



- Analiza kreditne sposobnosti građana
- Ulazi su visina prihoda i uštedjevina građana, a izlaz je odluka da li dodjela kredita predstavlja nizak ili visok rizik
- Ulaz: $\mathbf{x} = [x_1, x_2]^T$ Izlaz: $C \in \{0, 1\}$

- Predviđanje:
$$\begin{cases} C = 1, \text{ ako } P(C = 1 | x_1, x_2) > 0.5 \\ C = 0, \text{ inače} \end{cases}$$

ili ekvivalentno

$$\begin{cases} C = 1, \text{ ako } P(C = 1 | x_1, x_2) > P(C = 0 | x_1, x_2) \\ C = 0, \text{ inače} \end{cases}$$

Bayesovo pravilo

A priori vjerojatnost

Izglednost (likelihood)

A posteriori vjerojatnost

$$P(C | \mathbf{x}) = \frac{P(C) p(\mathbf{x} | C)}{p(\mathbf{x})}$$

Činjenica (evidence)

Vrijedi ako su hipoteze C_i međusobno isključive, a zbroj njihovih vjerojatnosti iznosi 1

$$P(C = 0) + P(C = 1) = 1$$

$$p(\mathbf{x}) = p(\mathbf{x} | C = 1)P(C = 1) + p(\mathbf{x} | C = 0)P(C = 0)$$

$$P(C = 0 | \mathbf{x}) + P(C = 1 | \mathbf{x}) = 1$$

Zadatak



Bayesovo pravilo: $K > 2$ razreda

$$\begin{aligned} P(C_i | \mathbf{x}) &= \frac{p(\mathbf{x} | C_i) P(C_i)}{p(\mathbf{x})} \\ &= \frac{p(\mathbf{x} | C_i) P(C_i)}{\sum_{k=1}^K p(\mathbf{x} | C_k) P(C_k)} \end{aligned}$$

Pri čemu mora vrijediti: $P(C_i) \geq 0$ i $\sum_{i=1}^K P(C_i) = 1$

Odabiremo: C_i ako $P(C_i | \mathbf{x}) = \max_k P(C_k | \mathbf{x})$

- Akcija α_i predstavlja odluku da ulaz pridružimo razredu C_i
- Gubitak nastao poduzimanjem akcije α_i ako ulaz pripada razredu C_k označit ćemo sa: λ_{ik}
- Očekivani rizik (engl. *Risk*):

$$R(\alpha_i | \mathbf{x}) = \sum_{k=1}^K \lambda_{ik} P(C_k | \mathbf{x})$$

Odabiremo: α_i ako $R(\alpha_i | \mathbf{x}) = \min_k R(\alpha_k | \mathbf{x})$

Rizici i gubici: 0/1 gubitak

$$\lambda_{ik} = \begin{cases} 0, & \text{ako } i = k \\ 1, & \text{ako } i \neq k \end{cases}$$

$$\begin{aligned} R(\alpha_i | \mathbf{x}) &= \sum_{k=1}^K \lambda_{ik} P(C_k | \mathbf{x}) \\ &= \sum_{k \neq i} P(C_k | \mathbf{x}) \\ &= 1 - P(C_i | \mathbf{x}) \end{aligned}$$

Ako želimo minimalni rizik, odabrat ćemo najvjerojatniji slučaj

Rizici i gubici: odbacivanje

$$\lambda_{ik} = \begin{cases} 0, & \text{ako } i = k \\ \lambda, & \text{ako } i = K + 1, \ 0 < \lambda < 1 \\ 1, & \text{inače} \end{cases}$$

$$R(\alpha_{K+1} | \mathbf{x}) = \sum_{k=1}^K \lambda P(C_k | \mathbf{x}) = \lambda \underbrace{\sum_{k=1}^K P(C_k | \mathbf{x})}_{=1} = \lambda$$

$$R(\alpha_i | \mathbf{x}) = \sum_{k \neq i} P(C_k | \mathbf{x}) = 1 - P(C_i | \mathbf{x})$$

Odabiremo:

$$\begin{cases} C_i, & \text{ako } P(C_i | \mathbf{x}) > P(C_k | \mathbf{x}) \forall k \neq i \text{ i } P(C_i | \mathbf{x}) > 1 - \lambda \\ \text{odbaci,} & \text{inače} \end{cases}$$

- Što je s rubnim slučajevima kada je $\lambda = 0$, odnosno $\lambda \geq 1$?

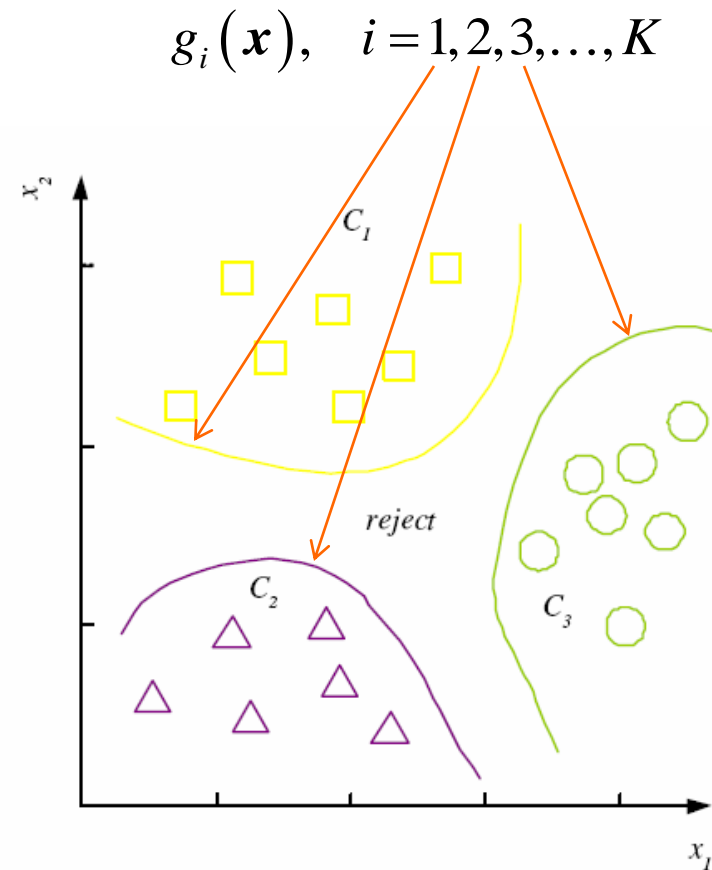
Odabiremo:

$$C_i \text{ ako } g_i(\mathbf{x}) = \max_k g_k(\mathbf{x})$$

$$g_i(\mathbf{x}) = \begin{cases} -R(\alpha_i | \mathbf{x}) \\ P(C_i | \mathbf{x}) \\ p(\mathbf{x} | C_i)P(C_i) \end{cases}$$

K područja odluke R_1, \dots, R_K

$$R_i = \left\{ \mathbf{x} \mid g_i(\mathbf{x}) = \max_k g_k(\mathbf{x}) \right\}$$



- Vjerojatnost stanja S_k uz poznat vektor značajki \mathbf{x} : $P(S_k | \mathbf{x})$
- Korisnost akcije α_i kada je stvarno stanje S_k : U_{ik}
- Očekivana korisnost (engl. *Expected Utility*):

$$EU(\alpha_i | \mathbf{x}) = \sum_k U_{ik} P(S_k | \mathbf{x})$$

Odabiremo: α_i ako $EU(\alpha_i | \mathbf{x}) = \max_j EU(\alpha_j | \mathbf{x})$

- Očekivana korisnost ako koristimo samo vektor značajki \mathbf{x} :

$$EU(\mathbf{x}) = \max_i \sum_k U_{ik} P(S_k | \mathbf{x})$$

- Očekivana korisnost ako uz vektor značajki \mathbf{x} koristimo novu značajku z :

$$EU(\mathbf{x}, z) = \max_i \sum_k U_{ik} P(S_k | \mathbf{x}, z)$$

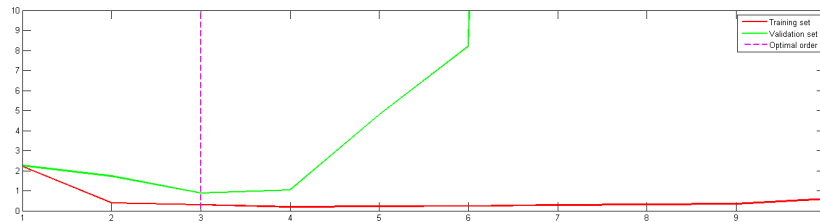
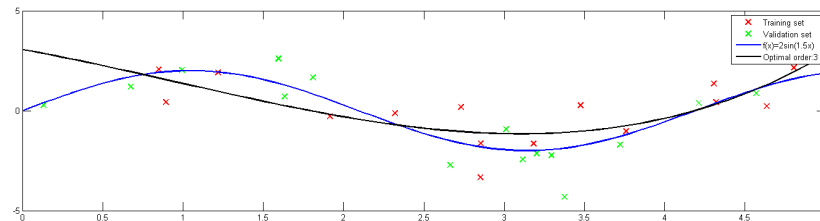
- Kažemo da je značajka z korisna ako vrijedi:

$$EU(\mathbf{x}, z) > EU(\mathbf{x})$$

Parametarske metode

Ethem Alpaydin,
Introduction to Machine Learning:

Chapter 4: *Parametric Methods*



- Statistika je bilo koja vrijednost izračunata iz danog uzorka.
- Kod statističkog zaključivanja odluke donosimo koristeći informacije koje su nam dostupne iz uzorka.
- Parametarski pristup podrazumijeva da uzorak dolazi iz određene distribucije kojoj se podvrgavaju i još neviđeni podaci.
- Prednost parametarskih metoda je da podatke možemo opisati malim brojem parametara.
- Dovoljna statistika (engl. *Sufficient statistics*) je skup parametara koje nužno moramo poznavati kako bi u potpunosti opisali distribuciju, npr: srednja vrijednost i varijanca čine dovoljnu statistiku Gaussove distribucije $\mathcal{N}(\mu, \sigma^2)$.
- Parametre distribucije procjenjujemo iz danog uzorka te na temelju procijenjene distribucije donosimo potrebne zaključke.

Parametarska procjena

- Smatramo da su značajke iz uzorka nezavisne i identično distribuirane (engl. *Independent identically distributed – iid*) $X = \{x^t\}_{t=1}^N$
- Pretpostavljamo da su x^t instance dobivene iz neke poznate familije gustoće vjerojatnosti čija je dovoljna statistika skup parametara θ .

$$x^t \sim p(x|\Theta)$$

- Želimo pronaći θ koja bi uzorkovanje x^t iz distribucije $p(x|\theta)$ učinilo najvjerojatnijim što je moguće.
- Dovoljan skup parametara θ procjenjujemo na temelju danog uzorka.

Kriterij najveće izglednosti

- Kriterij najveće izglednosti (engl. *Maximum Likelihood Estimation* – *MLE*) je metoda kojom nalazimo dovoljan skup parametara θ koja bi uzorkovanje x^t iz distribucije $p(x|\theta)$ učinilo najvjerojatnijim što je moguće.

- ***Izglednost** (engl. *Likelihood*) od θ uz dani uzorak X :

$$l(\theta | X) \equiv p(X | \theta) = \prod_{t=1}^N p(x^t | \theta)$$

Zašto logaritmiramo?

- **Logaritamska izglednost** (engl. *Log likelihood*):

$$L(\theta | X) = \log l(\theta | X) = \sum_{t=1}^N \log p(x^t | \theta)$$

- **Procjenitelj najveće izglednosti** (engl. *Maximum likelihood estimator* – *MLE*):

$$\theta^* = \arg \max_{\theta} L(\theta | X)$$

*U hrvatskoj literaturi koristi se još i naziv *vjerodostojnost*

MLE – Bernoullijeva razdioba

- **Bernoullijeva razdioba**
- Dva stanja, neuspjeh/uspjeh, $x \in \{0,1\}$

$$p(x) = p_0^x (1 - p_0)^{(1-x)}$$

$$\begin{aligned} L(p_0 | X) &= \log l(p_0 | X) \\ &= \log \prod_{t=1}^N p_0^{(x^t)} (1 - p_0)^{(1-x^t)} \\ &= \sum_{t=1}^N x^t \log p_0 + \left(N - \sum_{t=1}^N x^t \right) \log(1 - p_0) \end{aligned}$$

- **MLE:** $\hat{p} = \frac{dL}{dp} = 0 \quad \Rightarrow \quad \hat{p} = \frac{\sum_t x^t}{N}$

Dokaz



MLE – Polinomna razdioba

- **Polinomna razdioba**

- Generalizacija Bernoullijeve razdiobe, $K > 2$ stanja, $x \in \{0, 1\}$

$$p(x_1, x_2, \dots, x_K) = \prod_{i=1}^K p_i^{x_i}$$

- Radimo N nezavisnih pokusa s ishodima $X = \{x^t\}_{t=1}^N$ pri čemu vrijedi:

$$x_i^t = \begin{cases} 1, & \text{ako je ishod pokusa } t \text{ stanje } i \\ 0, & \text{inače} \end{cases} \quad \sum_i x_i^t = 1$$

- Polinomnu razdiobu možemo promatrati kao K odvojenih Bernoullijevih pokusa stoga vrijedi:

- **MLE:** $\hat{p} = \frac{\sum_t x^t}{N}$

MLE – Gaussova (normalna) razdioba

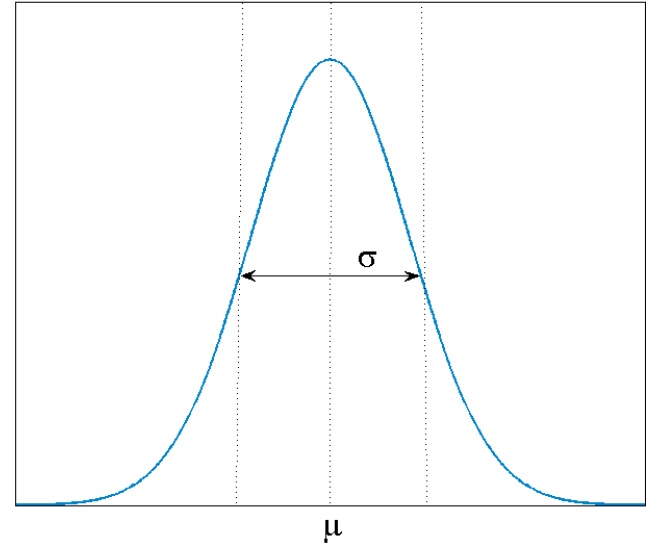
- **Gaussova razdioba**

$$p(x) = \mathcal{N}(\mu, \sigma^2)$$
$$= \frac{1}{\sqrt{2\pi}\sigma} \exp\left[-\frac{(x-\mu)^2}{2\sigma^2}\right]$$

- **MLE** za μ i σ^2 :

$$m = \frac{\sum_t x^t}{N} \quad \sigma^2 = \frac{\sum_t (x^t - m)^2}{N}$$

- **Konvencija:** Malim grčkim slovima označavamo parametre populacije, a malim slovima latinice njihove procjenitelje dobivene iz uzoraka.



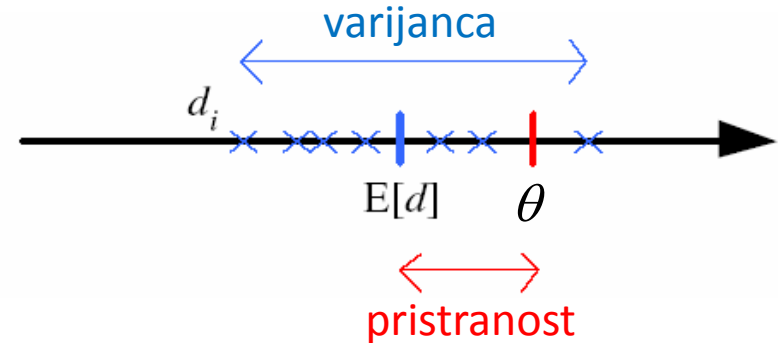
Dokaz



Priistranost i varijanca

- Nepoznati parametar θ
- Procjenitelj parametra θ na uzorku X_i je $d_i = d(X_i)$
- **Priistranost:** $b_\theta(d) = E[d] - \theta$
- **Varijanca:** $E[(d - E[d])^2]$
- **Srednja kvadratna pogreška:**

$$\begin{aligned} r(d, \theta) &= E[(d - \theta)^2] \\ &= (E[d] - \theta)^2 + E[(d - E[d])^2] \\ &= \text{Priistranost}^2 + \text{Varijanca} \end{aligned}$$



$$\text{Varijanca} = \sigma^2$$

- Ponekad i prije promatranja uzoraka imamo *a priori* informacije o mogućem intervalu vrijednosti koje neki parametar može poprimiti. Takva informacija je izuzetno korisna i uvijek ju treba iskoristiti, posebno u slučajevima kad je skup uzoraka malen.
- PRIMJER: Ekspert iz domene problema nam je rekao da se vrijednost parametra λ u 90% slučajeva nalazi u intervalu $[5, 9]$, (simetrični interval oko 7). Također znamo da traženi parametar ima normalnu razdiobu.

$$P\left\{-1.64 < \frac{\theta - \mu}{\sigma} < 1.64\right\} = 0.9$$

$$P\{\mu - 1.64\sigma < \theta < \mu + 1.64\sigma\} = 0.9$$

- Zaključujemo: $p(\theta) \propto \mathcal{N}\left(7, (2/1.64)^2\right)$

- Promatramo parametar θ kao slučajnu varijablu s pripadnom gustoćom vjerojatnosti $p(\theta)$
- **Bayesovo pravilo:**

$$p(\theta | X) = \frac{p(X | \theta) p(\theta)}{p(X)} = \frac{p(X | \theta) p(\theta)}{\int p(X | \theta') p(\theta') d\theta'}$$

- **Procjena gustoće vjerojatnosti kod uzorka x :**

$$\begin{aligned} p(x | X) &= \int p(x, \theta | X) d\theta \\ &= \int p(x | \theta, X) p(\theta | X) d\theta \\ &= \int p(x | \theta) p(\theta | X) d\theta \end{aligned}$$

Bayesov estimator

- Promatramo parametar θ kao slučajnu varijablu s pripadnom gustoćom vjerojatnosti $p(\theta)$
- Maximum a Posteriori (MAP):** $\theta_{MAP} = \arg \max_{\theta} p(\theta | X)$
$$= \arg \max_{\theta} \frac{p(X | \theta) p(\theta)}{p(X)}$$

Konstantan izraz, ne zavisi od skupa parametara, θ

$$= \arg \max_{\theta} p(X | \theta) p(\theta)$$
- Maximum Likelihood (ML):** $\theta_{ML} = \arg \max_{\theta} p(X | \theta)$
- Bayesov estimator:** $\theta_{Bayes} = E[\theta | X] = \int \theta p(\theta | X) d\theta$

Bayesov estimator - primjer

- $x^t \square \mathcal{N}(\theta, \sigma_0^2)$ i $\theta \square \mathcal{N}(\mu, \sigma^2)$

- $\theta_{MAP} = \theta_{Bayes} =$

$$E[\theta | X] = \frac{N/\sigma_0^2}{N/\sigma_0^2 + 1/\sigma^2} m + \frac{1/\sigma^2}{N/\sigma_0^2 + 1/\sigma^2} \mu$$

Ako je je $p(\theta | X)$ gustoća normalne razdiobe tada vrijedi $\theta_{MAP} = \theta_{Bayes}$

- U ovom slučaju Bayesov estimator je težinski prosjek od apriorne srednje vrijednosti, μ i srednje vrijednosti procjenjene iz uzoraka, m .
- Što možemo zaključiti iz dobivenog izraza:
 - 1) kad se povećava broj uzoraka, N ?
 - 2) kad je varijanca σ^2 mala?
- Kakvo je fizikalno tumačenje toga?

Parametarska klasifikacija

$$g_i(x) = p(x | C_i) P(C_i)$$

ili ekvivalentno

$$g_i(x) = \log p(x | C_i) + \log P(C_i)$$

$$p(x | C_i) = \frac{1}{\sqrt{2\pi}\sigma_i} \exp\left[-\frac{(x - \mu_i)^2}{2\sigma_i^2}\right]$$

$$g_i(x) = -\frac{1}{2} \log 2\pi - \log \sigma_i - \frac{(x - \mu_i)^2}{2\sigma_i^2} + \log P(C_i)$$

Parametarska klasifikacija

- Neka nam je dan skup primjera za učenje: $X = \{x^t, r^t\}_{t=1}^N$

$$x \in \square \quad r_i^t = \begin{cases} 1, & \text{ako } x^t \in C_i \\ 0, & \text{ako } x^t \in C_j, j \neq i \end{cases}$$

- ML procjenitelji su:

$$\hat{P}(C_i) = \frac{\sum_t r_i^t}{N} \quad m_i = \frac{\sum_t x^t r_i^t}{\sum_t r_i^t} \quad s_i^2 = \frac{\sum_t (x^t - m_i)^2 r_i^t}{\sum_t r_i^t}$$

- Funkcija odluke sada postaje:

$$g_i(x) = -\frac{1}{2} \log 2\pi - \log s_i - \frac{(x - m_i)^2}{2s_i^2} + \log \hat{P}(C_i)$$

Parametarska klasifikacija

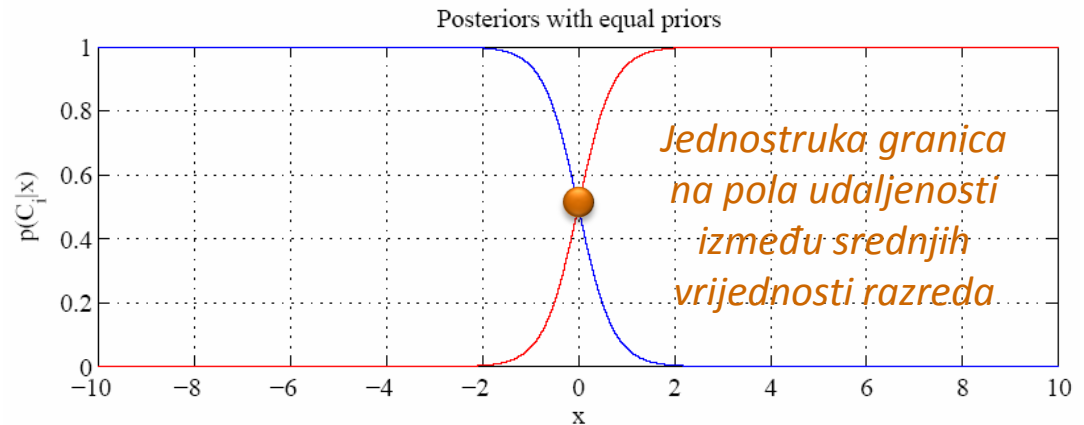
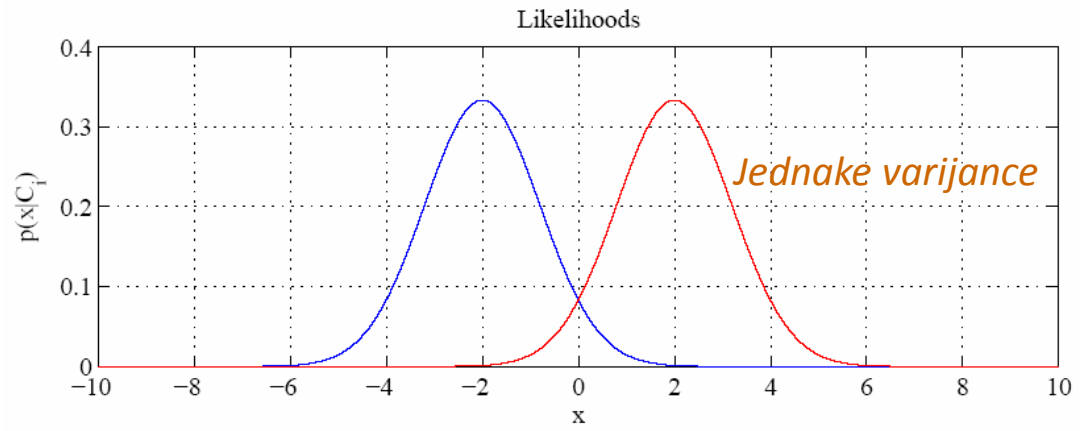
Funkcija izglednosti $p(x | C_i)$
i posteriori vjerojatnost
 $p(C_i | x)$ za slučaj kada su:

- *a priori* vjerojatnosti za dva razreda jednake
- varijance su jednake

$$g_1(x) = g_2(x)$$

$$(x - m_1)^2 = (x - m_2)^2$$

$$x = \frac{m_1 + m_2}{2}$$



$$g_i(x) = -\frac{1}{2} \log 2\pi - \log s_i - \frac{(x - m_i)^2}{2s_i^2} + \log \hat{P}(C_i)$$

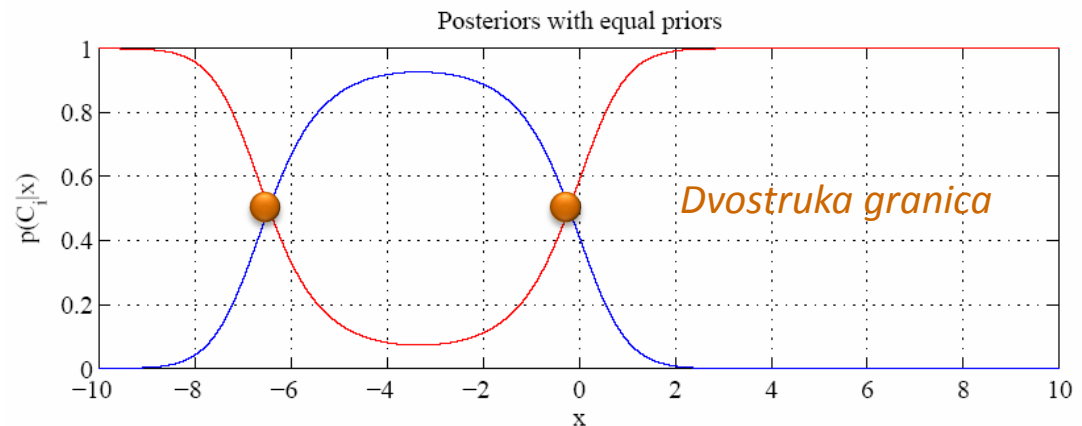
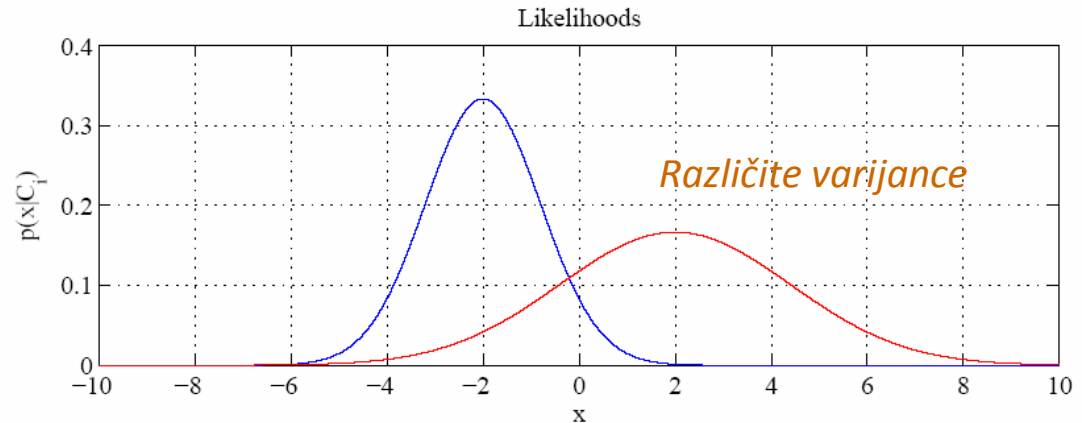
Parametarska klasifikacija

Funkcija izglednosti $p(x | C_i)$
i posteriori vjerojatnost
 $p(C_i | x)$ za slučaj kada su:

- *a priori* vjerojatnosti za dva razreda jednake
- varijance su različite

Što se događa ako su
a priori vjerojatnosti
različite?

- Prag odluke se pomiče prema srednjoj vrijednosti manje izglednijeg razreda.



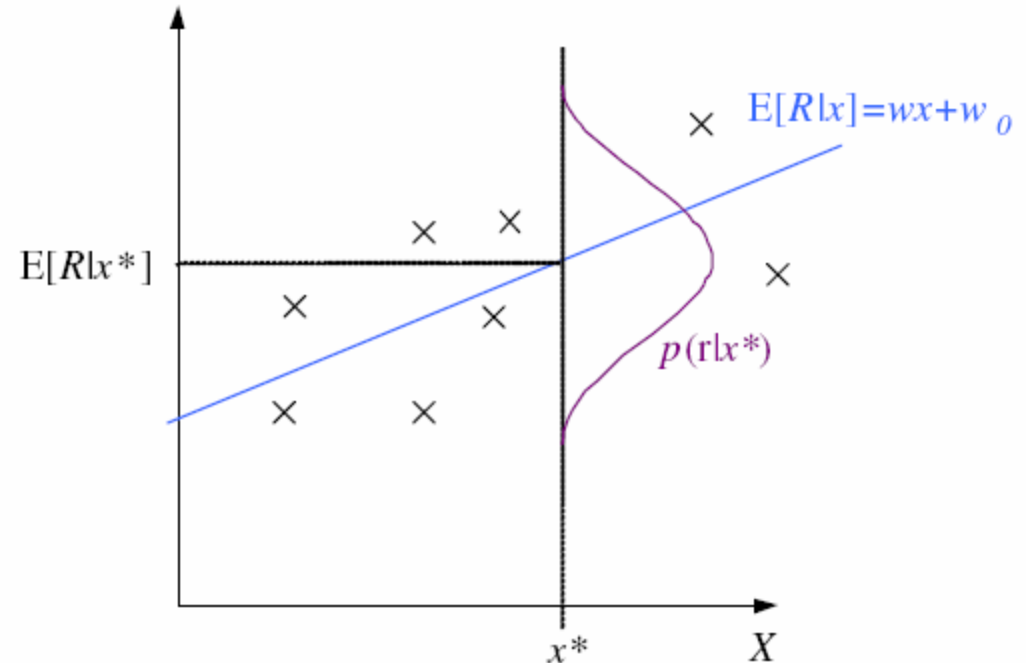
$$g_i(x) = -\frac{1}{2} \log 2\pi - \log s_i - \frac{(x - m_i)^2}{2s_i^2} + \log \hat{P}(C_i)$$

$$r = f(x) + \varepsilon$$

procjenitelj: $g(x|\theta)$

$$\varepsilon \square \mathcal{N}(0, \sigma^2)$$

$$p(r|x) \square \mathcal{N}(g(x|\theta), \sigma^2)$$



$$L(\theta | X) = \log \prod_{t=1}^N p(x^t, r^t)$$

$$= \log \prod_{t=1}^N p(r^t | x^t) + \log \prod_{t=1}^N p(x^t)$$

$$p(x, r) = p(r|x)p(x)$$

Srednja kvadratna pogreška

$$L(\theta | X) = \log \prod_{t=1}^N p(r^t | x^t) + \log \prod_{t=1}^N p(x^t)$$

$$\begin{aligned} L(\theta | X) &= \log \prod_{t=1}^N \frac{1}{\sqrt{2\pi}\sigma} \exp \left[-\frac{[r^t - g(x^t | \theta)]^2}{2\sigma^2} \right] \\ &= \log \left(\frac{1}{\sqrt{2\pi}\sigma} \right)^N \exp \left[-\frac{1}{2\sigma^2} \sum_{t=1}^N [r^t - g(x^t | \theta)]^2 \right] \\ &= -N \log(\sqrt{2\pi}\sigma) - \frac{1}{\sigma^2} \frac{1}{2} \sum_{t=1}^N [r^t - g(x^t | \theta)]^2 \end{aligned}$$

$$E(\theta | X) = \frac{1}{2} \sum_{t=1}^N [r^t - g(x^t | \theta)]^2$$

$$g(x^t | w_1, w_0) = w_1 x^t + w_0$$

$$E(\theta | X) = \frac{1}{2} \sum_{t=1}^N \left[r^t - g(x^t | \theta) \right]^2$$

$$\frac{\partial E(\theta | X)}{\partial w_0} : \sum_t r^t = N w_0 + w_1 \sum_t x^t$$

$$\frac{\partial E(\theta | X)}{\partial w_1} : \sum_t r^t x^t = w_0 \sum_t x^t + w_1 \sum_t (x^t)^2$$

$$\mathbf{A} = \begin{bmatrix} N & \sum_t x^t \\ \sum_t x^t & \sum_t (x^t)^2 \end{bmatrix} \quad \mathbf{w} = \begin{bmatrix} w_0 \\ w_1 \end{bmatrix} \quad \mathbf{y} = \begin{bmatrix} \sum_t r^t \\ \sum_t r^t x^t \end{bmatrix}$$

- Zapis u matričnom obliku: $\mathbf{w} = \mathbf{A}^{-1} \mathbf{y}$

- Generalizacija linearne regresije:

$$g(x^t | w_k, \dots, w_2, w_1, w_0) = w_k (x^t)^k + \dots + w_2 (x^t)^2 + w_1 x^t + w_0$$

$$\mathbf{D} = \begin{bmatrix} 1 & x^1 & (x^1)^2 & \dots & (x^1)^k \\ 1 & x^2 & (x^2)^2 & \dots & (x^2)^k \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x^N & (x^N)^2 & \dots & (x^N)^k \end{bmatrix} \quad \mathbf{r} = \begin{bmatrix} r^1 \\ r^2 \\ \vdots \\ r^N \end{bmatrix}$$

- Zapis u matričnom obliku: $\mathbf{w} = (\mathbf{D}^T \mathbf{D})^{-1} \mathbf{D}^T \mathbf{r}$

- Srednja kvadratna pogreška:

$$E(\theta | X) = \frac{1}{2} \sum_{t=1}^N [r^t - g(x^t | \theta)]^2$$

- Relativna kvadratna pogreška:

$$E(\theta | X) = \frac{\sum_{t=1}^N [r^t - g(x^t | \theta)]^2}{\sum_{t=1}^N [r^t - \bar{r}]^2}$$

- Apsolutna pogreška:

$$E(\theta | X) = \sum_{t=1}^N |r^t - g(x^t | \theta)|$$

Ugađanje složenosti modela

- Očekivana kvadratna pogreška:

$$E\left[\left(r - g(x)\right)^2 \mid x\right] = \underbrace{E\left[\left(r - E[r \mid x]\right)^2 \mid x\right]}_{\text{šum}} + \underbrace{\left(E[r \mid x] - g(x)\right)^2}_{\text{kvadratna pogreška}}$$

Šum ne ovisi o izboru estimatora.
To je dio pogreške kojeg nikad ne
možemo ukloniti

Kvadratna pogreška ovisi o
izboru procjenitelju i o skupu
primjera za učenje

- Očekivana vrijednost (*prosjek nad svim uzorcima X*):

$$E_x\left[\left(E[r \mid x] - g(x)\right)^2\right] = \underbrace{\left(E[r \mid x] - E_x[g(x)]\right)^2}_{\text{pristranost}} + \underbrace{E_x\left[\left(g(x) - E_x[g(x)]\right)^2\right]}_{\text{varijanca}}$$

Procjena pristranosti i varijance

- M skupova primjra za učenje, $X_i = \{x_i^t, r_i^t\}, i = 1, \dots, M$ koristimo kako bi podesili regresijske funkcije u obliku polinoma $g_i(x), i = 1, \dots, M$.

$$\text{Pristranost}^2(g) = \frac{1}{N} \sum_t \left[\bar{g}(x^t) - f(x^t) \right]^2$$

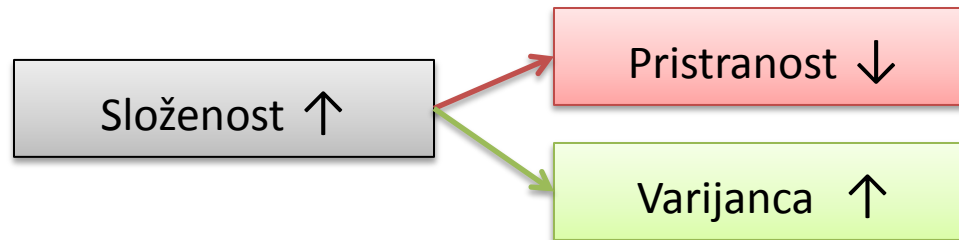
$$\text{Varijanca}(g) = \frac{1}{NM} \sum_t \sum_i \left[g_i(x^t) - \bar{g}(x^t) \right]^2$$

- Pri čemu $\bar{g}(x)$ računamo kao:

$$\bar{g}(x) = \frac{1}{M} \sum_i g_i(x)$$

Dvojba pristranost/varijanca

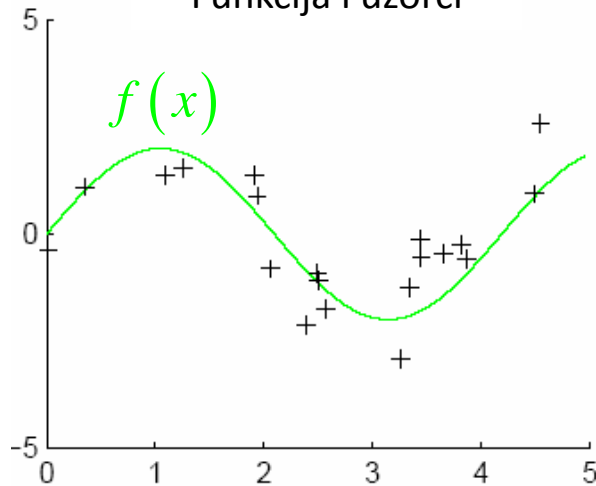
- Dvojba pristranost/varijanca, (engl. *Bias/Variance dilemma*): (Geman et al., 1992)
- Kako povećavamo složenost modela:
 - pristranost se smanjuje (bolje poklapanje s podacima)
 - varijanca se povećava (poklapanje više varira)



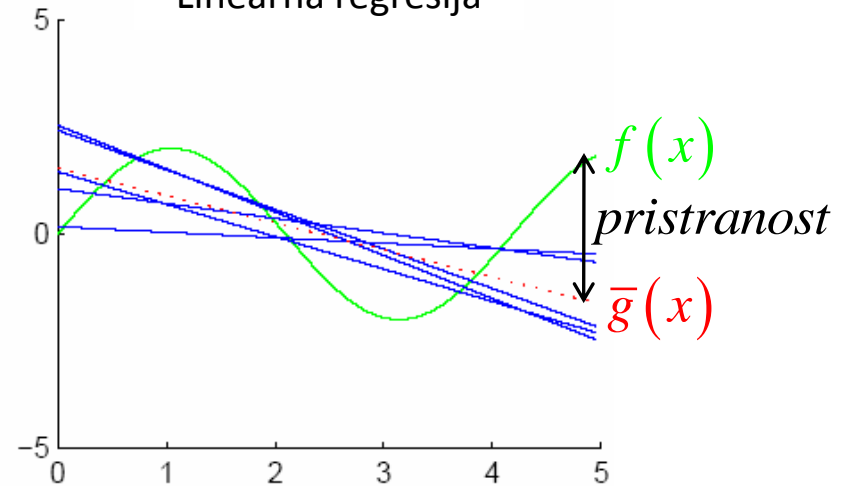
- PRIMJERI:
 - $g_i(x) = 2 = \text{const.}$ nema varijancu, ali ima veliku pristranost.
 - $g_i(x) = \sum_t r_i^t / N$ ima manju pristranost, ali ima i varijancu

Dvojba pristranost/varijanca

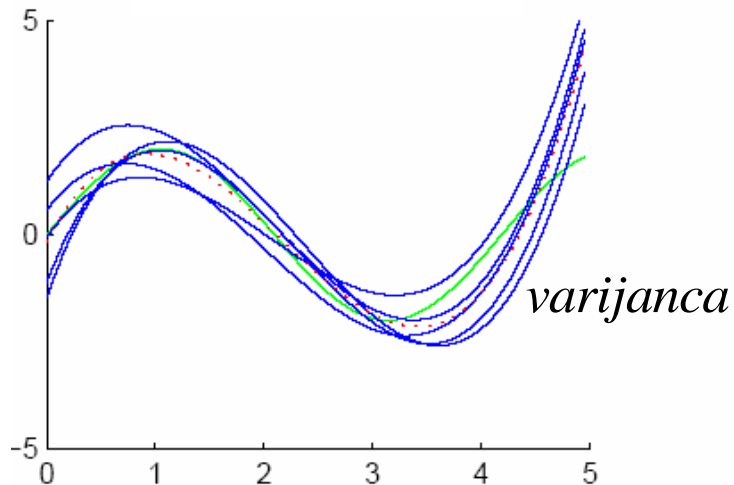
Funkcija i uzorci



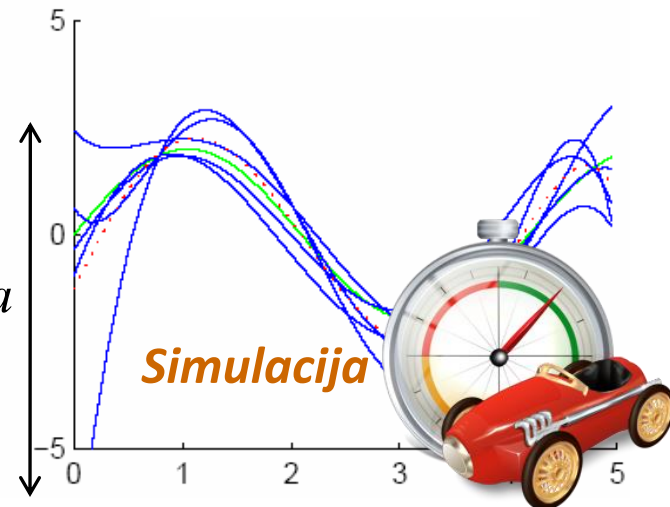
Linearna regresija



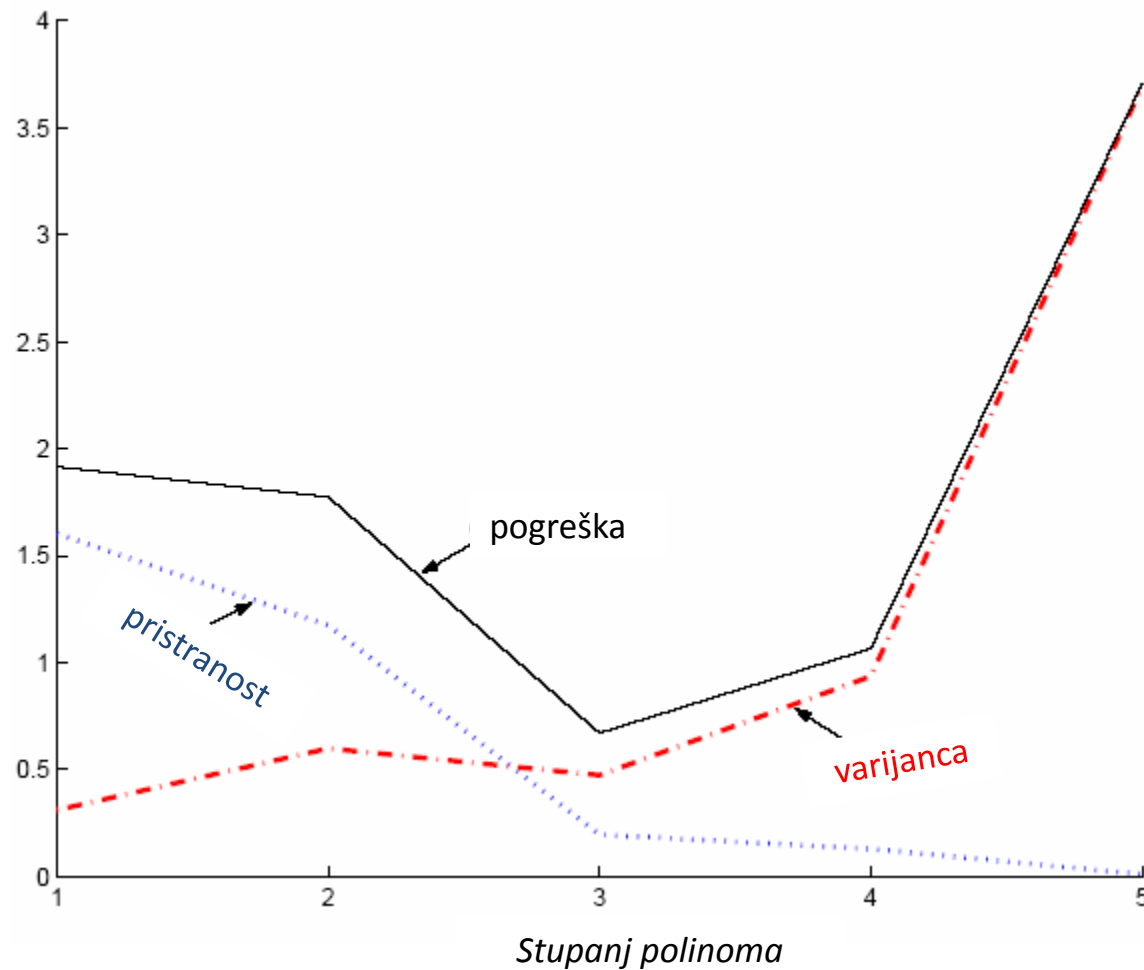
Polinomi 3. reda



Polinomi 5. reda



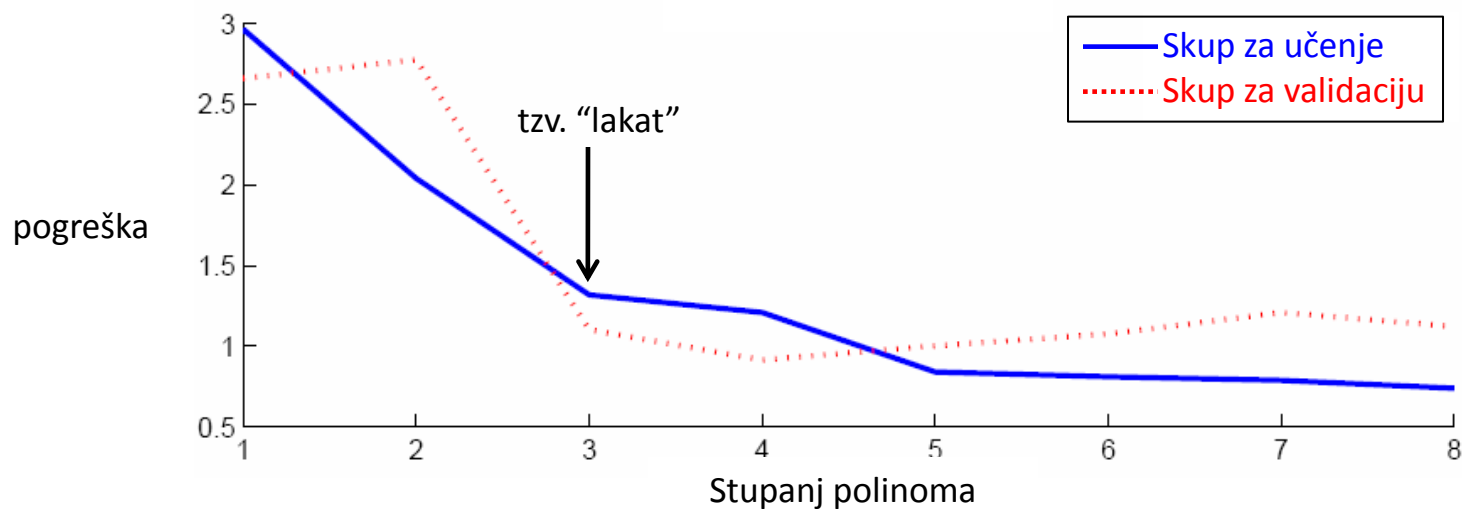
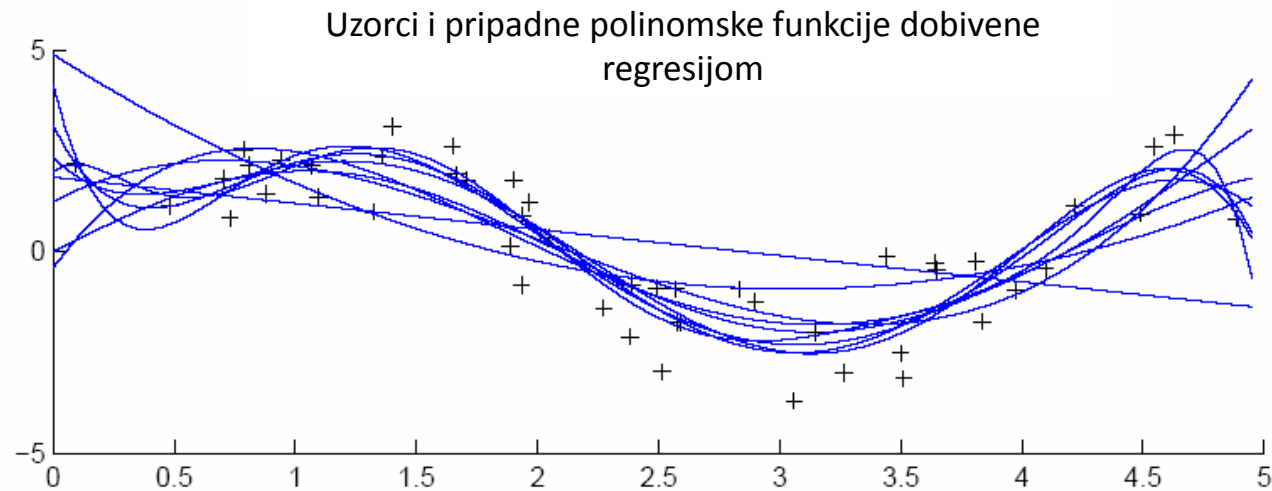
Polinomska regresija



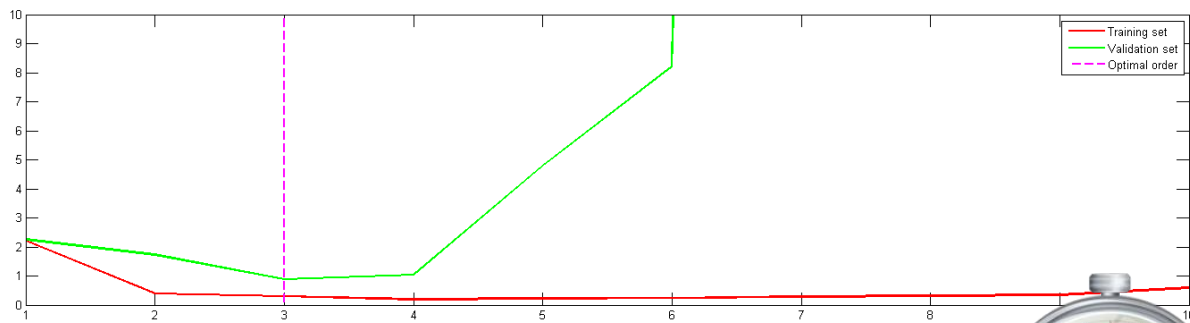
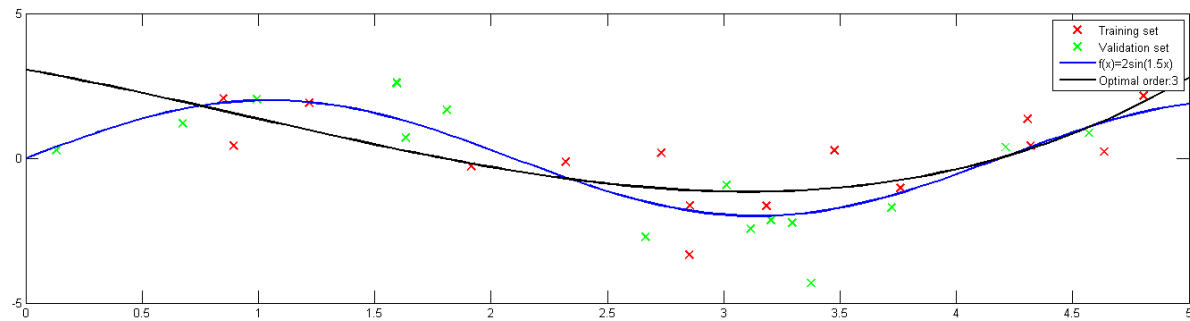
Postupci odabira modela

- **Unakrsna validacija (engl. *Cross-validation*):** ne možemo izračunati pristranost i varijancu za model, ali možemo izračunati ukupnu pogrešku. Uzorke podijelimo na skup za učenje i skup za validaciju. Uvježbamo modele različitih složenosti i ispitujemo njihove pogreške na skupu za validaciju.
- Dok povećavamo složenost modela, pogreška na skupu za učenje se smanjuje. Pogreška na skupu za validaciju se smanjuje do određene razine složenosti. Pri toj razini pogreška se zaustavlja smanjivati, ili se beznačajno smanjuje. Ako je u podacima šum jako izražen, pogreška se može i povećavati. Optimalnu složenost modela na grafu prepoznamo prema karakterističnom obliku lakta (engl. *elbow*).

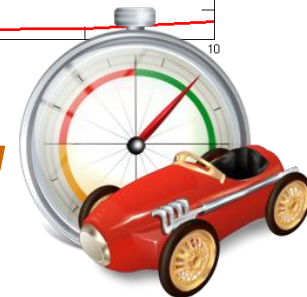
Postupci odabira modela



Postupci odabira modela



Simulacija



- **Regularizacija**, (Breiman, 1998): Koristimo uvećanu (engl. *augmented*) funkciju pogreške:

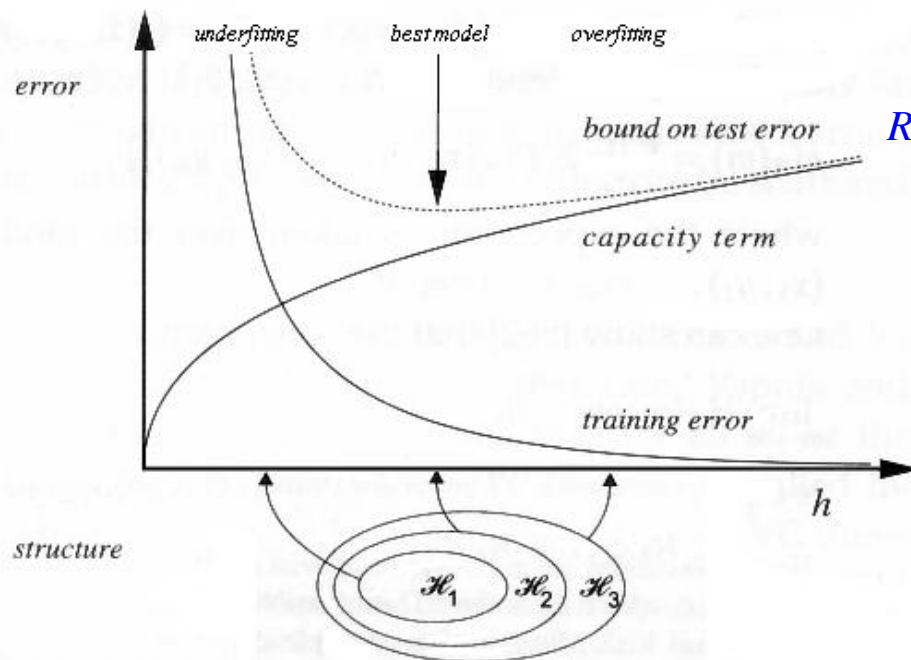
$$E' = \text{pogreška nad podacima} + \lambda \cdot \text{složenost modela}$$

Drugi član izraza “kažnjava” složene modele koji imaju veliku varijancu. Parametar λ predstavlja težinu kazne.

Postupci odabira modela

- **Minimizacija strukturnog rizika** (Vapnik, 1995):

Predstavlja opći model kontrole složenosti modela i omogućuje nam balansiranje između složenosti samog modela (u nekim slučajevima se može predstaviti VC dimenzijom) i kvalitete “poklapanja” modela s uzorcima iz skupa za učenje. Vrsta regularizacije.



$$R(\alpha) \leq R_{emp}(\alpha) + \sqrt{\frac{h \log\left(\frac{2m}{h} + 1\right) - \log\left(\frac{\eta}{4}\right)}{m}}$$

Empirijski rizik
(skup za učenje)

Složenost skupa
model

Minimizacija strukturnog rizika

▪ PRIMJER:

1. Koristeći apriorno znanje iz domene odabiremo klasu funkcija kao što su: polinomi stupnja n , neuronske mreže sa n skrivenih slojeva i slično.
2. Podijelimo klasu funkcija u hijerarhijski ugnježdene podskupove poredane po rastućoj složenosti. Npr. polinomi rastućeg stupnja.
3. Minimiziramo empirijski rizik na svakom podskupu (odabir optimalnih parametara u pogledu empirijskog rizika).
4. Odabiremo model iz niza sortiranih modela za kojeg je zbroj empirijskog rizika i funkcije ovisne o VC dimenziji minimalan.

- **Dužina minimalnog opisa (engl. *Minimum description length - MDL*),** (Rissanen, 1978): promotrimo maksimum a posteriori hipotezu u svijetlu teorije informacije:

$$\theta_{MAP} = \arg \max_{\theta} p(X | \theta) p(\theta)$$

$$\theta_{MAP} = \arg \max_{\theta} (\log_2 P(X | \theta) + \log_2 P(\theta))$$

$$\theta_{MAP} = \arg \min_{\theta} \left(\underbrace{-\log_2 P(X | \theta) - \log_2 P(\theta)}_{\text{entropija}} \right)$$

- MDL nam pruža način odabira složenosti modela. Možemo odabrati između jednostavnijeg modela koji radi određenu pogrešku i složenijeg modela koji savršeno klasificira primjere iz skupa za učenje.

Postupci odabira modela

- **Bayesov odabir modela:** Koristi se kada unaprijed postoji znanje o prikladnom razredu aproksimacijskih funkcija.

$$p(model | podatci) = \frac{p(podatci | model) p(model)}{p(podatci)}$$

- Unakrsna validacija se razlikuje od svih ostalih postupaka odabira modela jer ne pretpostavlja nikakve *a priori* pretpostavke o modelu. Unakrsna validacija je najbolji pristup odabiru modela ako je skup za validaciju dovoljno velik. Ostali modeli postaju korisni kada je skup uzoraka malen.