

Pokročilé metódy analýzy dát 1

úvod do hlbokého učenia

Peter Bednár

Prerekvizity

Prerekvizity 1

Označenia

- $\sum_{i=1}^N a_i$ – suma, $\prod_{i=1}^N a_i$ – súčin
- $\mathbf{x} = (x_1, x_2, \dots, x_N)$ – vektor reálnych čísel, \mathbf{A} – matica (riadky x stĺpce), $\mathbf{x}^T, \mathbf{A}^T$ – transponovanie vektora alebo matice
- $\mathbf{x}^T \mathbf{y} = \sum_{i=1}^N x_i y_i$ – skalárny súčin
- N – počet príkladov, M – počet atribútov, K – počet tried pri klasifikácii, (\mathbf{x}_i, y_i) – trénovacie príklady

Vlastnosti logaritmu

- $\log(\prod_{i=1}^N a_i) = \sum_{i=1}^N \log(a_i)$, $\log(a^b) = b \log(a)$

Prerekvizity 2

Numerická optimalizácia

- Pre zadanú funkciu $f(\mathbf{x})$ hľadáme také hodnoty parametrov \mathbf{x} , pre ktoré nadobúda funkcia maximálnu, alebo minimálnu hodnotu
 - $\operatorname{argmax}_{\mathbf{x}} f(\mathbf{x})$ – maximalizácia
 - $\operatorname{argmin}_{\mathbf{x}} f(\mathbf{x})$ – minimalizácia
- Metódy sú najčastejšie navrhnuté pre minimalizáciu, pre maximalizovanie platí:

$$\operatorname{argmax}_{\mathbf{x}} f(\mathbf{x}) = \operatorname{argmin}_{\mathbf{x}} -f(\mathbf{x})$$

Prerekvizity 3

- Ak vynásobíme funkciu kladným číslom a , optimálne parametre sa nezmenia:

$$\operatorname{argmax}_x af(\mathbf{x}) = \operatorname{argmax}_x f(\mathbf{x})$$

$$\operatorname{argmin}_x af(\mathbf{x}) = \operatorname{argmin}_x f(\mathbf{x})$$

- Podobne, ak transformujeme funkciu monotónne rastúcou funkciou, napr. logaritmom:

$$\operatorname{argmax}_x \log(f(\mathbf{x})) = \operatorname{argmax}_x f(\mathbf{x})$$

Prerekvizity 4

Derivácia a gradient

- Gradient funkcie $f(\mathbf{x})$, $\mathbf{x} = (x_1, x_2, \dots, x_M)$ je vektor parciálnych derivácií $\nabla f(\mathbf{x}) = \left(\frac{\partial f(\mathbf{x})}{\partial x_1}, \frac{\partial f(\mathbf{x})}{\partial x_2}, \dots, \frac{\partial f(\mathbf{x})}{\partial x_M} \right)$, tzn. gradient je vektorová funkcia a vieme ho vypočítať pre daný bod \mathbf{x}

Príklad

$$f(\mathbf{x}) = x_1^2 + 3x_2^2 + 2x_1x_2$$

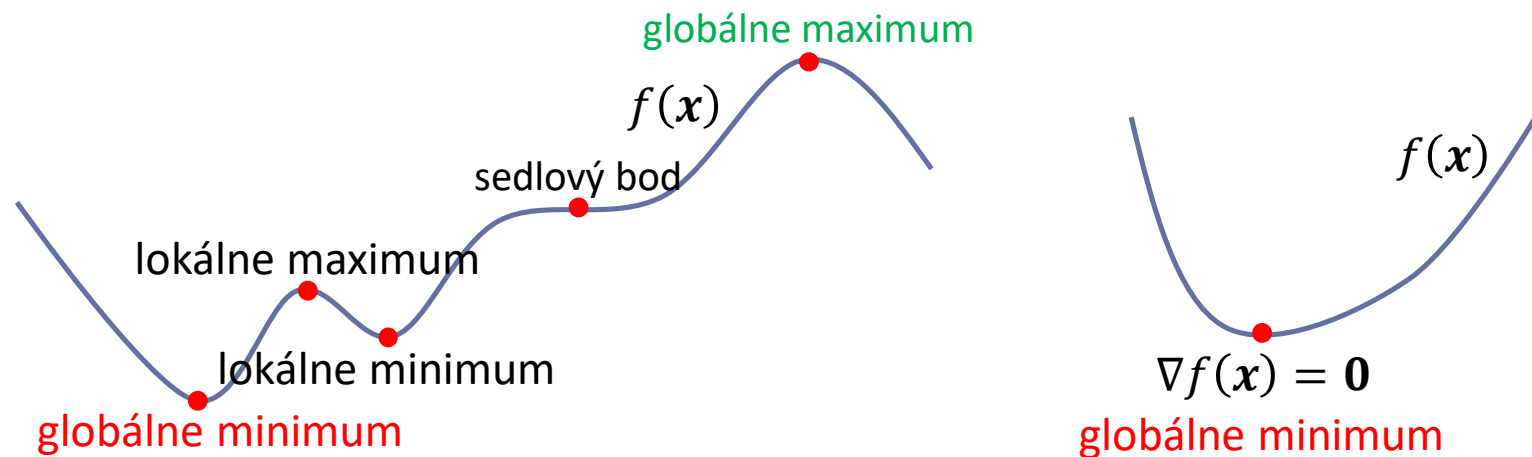
$$\nabla f(\mathbf{x}) = (2x_1 + 2x_2, 6x_2 + 2x_1)$$

Po dosadení gradient v bode $\mathbf{x} = [1, 2]$

$$\nabla f([1, 2]) = (5, 10)$$

Prerekvizity 5

- Body, pre ktoré sa gradient rovná nulovému vektoru $\nabla f(\mathbf{x}) = \mathbf{0}$ sa označujú ako stacionárne – globálne minimum je jeden z nich



- Zvláštny význam pri optimalizácii majú konvexné funkcie, pre ktoré platí, že stacionárny bod zároveň určuje globálne minimum

Prerekvizity 6

Podmienená pravdepodobnosť

- $P(Y|X)$ udáva pravdepodobnosť, že nastane jav Y za predpokladu, že nastal jav X
 - $P(Y|X) = P(Y,X) / P(X)$
- Podmienená pravdepodobnosť $P(Y=y | X_1=x_1, X_2=x_2, \dots, X_M=x_M)$ má význam pre predpovedanie hodnoty cieľového atribútu $Y=y$ na základe známych hodnôt vstupných atribútov $X_1 = x_1, X_2 = x_2, \dots, X_M = x_M$
 - Pre nový prípad zvolíme cieľovú hodnotu, pre ktorú je podmienená pravdepodobnosť maximálna

Prerekvizity 7

Príklad

- Chceme predpovedať, či študent príde na cvičenie ráno keď prší

Študent	Počasiе	Rozvrh	Počet prípadov
príde	slnečno	ráno	65
príde	slnečno	po obede	60
príde	prší	ráno	20
príde	prší	po obede	45
nepríde	slnečno	ráno	10
nepríde	slnečno	po obede	5
nepríde	prší	ráno	30
nepríde	prší	po obede	15
celkovo			250

Prerekvizity 8

$$\begin{aligned} & P(\text{\textcolor{green}{Študent = príde}} \mid \text{Počasie = prší, Rozvrh = ráno}) \\ &= \frac{P(\text{\textcolor{green}{Študent = príde}, Počasie = prší, Rozvrh = ráno})}{P(\text{Počasie = prší, Rozvrh = ráno})} \\ &= \frac{20/250}{(20 + 30)/250} = 0.4 \end{aligned}$$

$$\begin{aligned} & P(\text{\textcolor{red}{Študent = nepríde}} \mid \text{Počasie = prší, Rozvrh = ráno}) \\ &= 1 - P(\text{\textcolor{green}{Študent = príde}} \mid \text{Počasie = prší, Rozvrh = ráno}) \\ &= 0.6 \end{aligned}$$

- Predpovedáme, že študent nepríde s pravdepodobnosťou 0.6

Učenie prediktívnych modelov

Učenie prediktívnych modelov (1)

- Pri predikcii sa snažíme na základe známych atribútov $\mathbf{X} = (X_1, X_2, \dots, X_M)$ predpovedať hodnotu cieľového atribútu Y
- Predpokladáme, že medzi atribútmi \mathbf{X} a Y existuje nejaká závislosť, ktorú však vo všeobecnosti nepoznáme, k dispozícii máme iba tzv. trénovaciu množinu príkladov $(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_N, y_N)$ pre ktorú poznáme hodnoty vstupných atribútov \mathbf{x} a skutočnú hodnotu cieľového atribútu y

Učenie prediktívnych modelov (2)

- **Úloha:** na základe trénovacej množiny príkladov nájsť funkciu $f(x)$, ktorá pre zadané hodnoty vstupných atribútov x vypočíta čo najlepší odhad skutočnej hodnoty cieľového atribútu y vo všeobecnosti pre všetky nové prípady

Príklad regresie

- Máme realitnú kanceláriu a na základe vstupných atribútov:
 - X_1 – veľkosť úžitkovej plochy v m^2
 - X_2 – počet izieb
 - X_3 – navrhovaná cena
- Chceme predpovedať, za ako dlho predáme pre nového zákazníka jeho nehnuteľnosť na trhu:
 - Y – počet dní od zverejnenia ponuky do predaja
- Máme historické dáta napr. za posledný rok

Chybová funkcia

- Pre danú úlohu sa vyhodnocuje zhoda medzi skutočnou hodnotou y a vypočítanou predikciou podľa danej **chybovej funkcie** $L(y, f(\mathbf{x}))$
- Väčšinou sa rovná 0 ak sa y rovná $f(\mathbf{x})$ a rastie s narastajúcim rozdielom medzi y a $f(\mathbf{x})$
- **Príklady:**
- 1/0 chyba (klasifikačná chyba):
$$L(y, f(\mathbf{x})) = 0 \text{ ak } y = f(\mathbf{x}), \text{ inak } 1$$
- Kvadratická chyba:
$$L(y, f(\mathbf{x})) = \frac{1}{2} (y - f(\mathbf{x}))^2$$

Ako navrhnuť algoritmus učenia?

- **Základný princíp:**

- Nevieme spočítať chybu predikcie vo všeobecnosti, pretože nepoznáme cieľovú hodnotu pre všetky prípady
- Preto sa pokúsime nájsť funkciu, ktorá bude mať minimálnu chybu aspoň na tréningových dátach (\mathbf{x}_1, y_1) , (\mathbf{x}_2, y_2) , ..., (\mathbf{x}_N, y_N)
- tzn. zovšeobecnenú chybu sme nahradili tréningovou - tzv. empirickou chybou:

$$J(f) = \frac{1}{N} \sum_{i=1}^N L(y_i, f(\mathbf{x}_i))$$

Algoritmus učenia = Chybová funkcia + Model + Optimalizačná metóda

1. Zvolíme si **model**:

- Matematicky je to množina funkcií, z ktorých algoritmus vyberá výslednú funkciu $f(x)$

Učenie je teda optimalizačná úloha – z modelu vyberieme funkciu, pre ktorú bude chyba na tréningových dátach minimálna

2. Zvolíme si **optimalizačnú metódu**

- Napr. ak je model definovaný ako množina funkcií daného tvaru, ktoré majú okrem vstupných atribútov množinu parametrov $W = \{W_1, W_2, \dots\}$, zvolíme si metódu, ktorá vypočíta optimálne parametre w , pri ktorých je chyba funkcie $f(x, w)$ minimálna

Príklad – lineárna regresia (1)

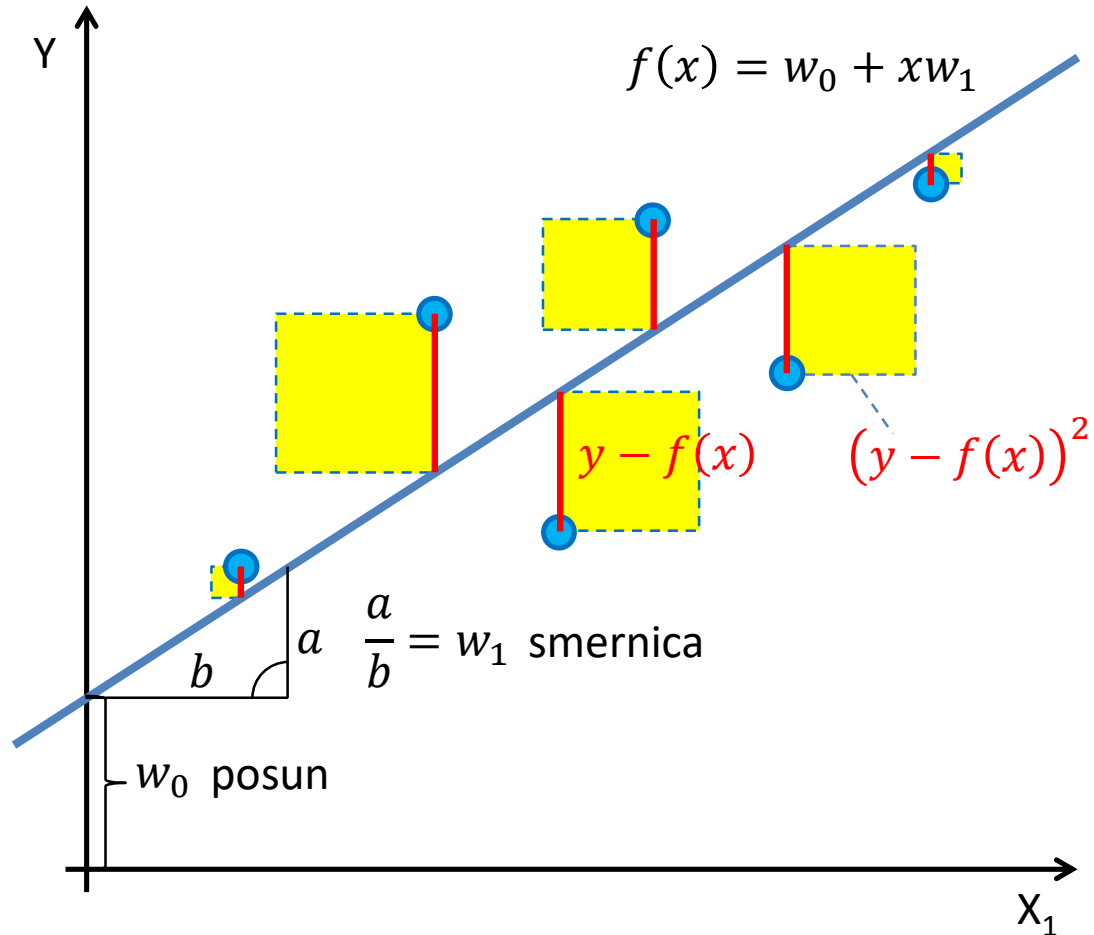
- Úlohou je regresia - cieľový atribút Y je číselný
- Minimalizujeme kvadratickú chybovú funkciu

$$J(f) = \frac{1}{N} \sum_{i=1}^N \frac{1}{2} (y_i - f(\mathbf{x}_i))^2$$

- **Lineárny model** – funkcie majú tvar lineárnej funkcie s parametrami $\mathbf{w} = (w_0, w_1, \dots, w_M)$

$$f(\mathbf{x}, \mathbf{w}) = w_0 + x_1 w_1 + \dots + x_M w_M$$

Príklad – lineárna regresia (2)



Príklad – lineárna regresia (3)

- Optimalizujeme parametre w , hodnoty trénovacích príkladov x , y sú konštanty
- Pri lineárnej regresii vieme priamo vypočítať optimálne parametre aj analyticky:
 1. Zderivujeme $J(f)$ podľa w_0, w_1, \dots, w_M a dostaneme $M+1$ funkcií
 2. Položíme derivácie = 0 a dostaneme sústavu $M+1$ rovníc, z ktorých si vypočítame optimálne parametre
- Algoritmus učenia lineárnej regresie = algoritmus riešenia sústavy rovníc (nie však vo všeobecnosti, pozor na lineárne závislé atribúty!)

Príklad – logistická regresia (1)

- Úlohou je klasifikácia do dvoch tried $Y = \{0, 1\}$
- Minimalizujeme 1/0 chybovú funkciu:
$$L(y, f(\mathbf{x})) = 0 \text{ ak } y = f(\mathbf{x}), \text{ inak } 1$$
- Trénovacia chyba $J(f)$ teda udáva pravdepodobnosť chybnéj klasifikácie
- Problém je, že 1/0 chybová funkcia sa nedá dobre derivovať (nedá sa priamo numericky optimalizovať)
- Preto budeme minimalizovať jej aproximáciu

Príklad – logistická regresia (2)

- Model definujeme pomocou podmienených pravdepodobností
 - $P(Y = 1 | X = x)$ – pravdepodobnosť, že má byť príklad s hodnotami vstupných atribútov x zaradený do triedy 1
 - $P(Y = 0 | X = x) = 1 - P(Y = 1 | X = x)$ – pravdepodobnosť, že má byť zaradený do triedy 0
- 1/0 chybová funkcia bude minimálna, ak zaradíme príklad do najpravdepodobnejšej triedy, tzn. pre dve triedy $\{0,1\}$ zaradíme príklad do triedy 1 ak $P(Y = 1 | X = x) > 0.5$, inak do triedy 0

Príklad – logistická regresia (3)

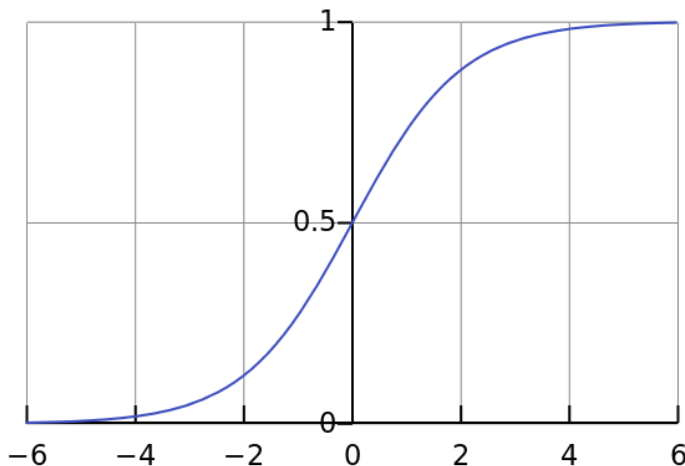
- Pravdepodobnosť $P(Y = 1 | X = \mathbf{x})$ (a teda ani $P(Y = 0 | X = \mathbf{x})$) nepoznáme
- Model logistickej regresie si navrhujeme ako funkciu $p(\mathbf{x})$, pomocou ktorej je možné priamo odhadnúť pravdepodobnosť zaradenia do triedy 1, $p_1(\mathbf{x}) \approx P(Y = 1 | X = \mathbf{x})$, pričom chceme, aby bol tento odhad založený na lineárnej funkcii

$$f(\mathbf{x}, \mathbf{w}) = w_0 + x_1 w_1 + \dots + x_M w_M$$

- Pravdepodobnosť pre triedu 0 potom môžeme odhadnúť ako $p_0(\mathbf{x}) = 1 - p_1(\mathbf{x})$

Príklad – logistická regresia (4)

- Keďže lineárna funkcia môže nadobúdať ľubovoľnú reálnu hodnotu
 - Pre odhad pravdepodobnosti potrebujeme hodnotu predikcie $f(\mathbf{x}, \mathbf{w})$ previesť do intervalu $(0,1)$
 - Pri logistickej regresii sa na transformovanie predikcie používa **logistická funkcia**



$$g(t) = \frac{1}{1+e^{-t}}$$

Príklad – logistická regresia (5)

- Odhad pravdepodobnosti pre triedu 1 dostaneme po dosadení lineárnej funkcie do logistickej

$$p_1(\mathbf{x}) = g(f(\mathbf{x}, \mathbf{w})) = g(w_0 + x_1 w_1 + \dots + x_M w_M)$$

- Ako vypočítame parametre w_0, w_1, \dots, w_M ?
 - Keďže $p_1(\mathbf{x})$ je odhad pravdepodobnosti, môžeme použiť metódu maximálnej vierohodnosti (*max. likelihood*)
1. Vyjadríme si pravdepodobnosť pre celú trénovaciu množinu (tzv. funkciu vierohodnosti $l(\mathbf{w})$) ako súčin pravdepodobností jednotlivých príkladov, tzn. pre príklady z triedy 1 dosadíme priamo $p_1(\mathbf{x})$ a pre príklady z triedy 0 dosadíme $1 - p_1(\mathbf{x})$
 2. Zvolíme také parametre, pre ktoré je funkcia vierohodnosti maximálna

Príklad – logistická regresia (6)

- Keďže $y_i = \{0, 1\}$, celkovú funkciu vierohodnosti môžeme skrátene vyjadriť ako súčin:

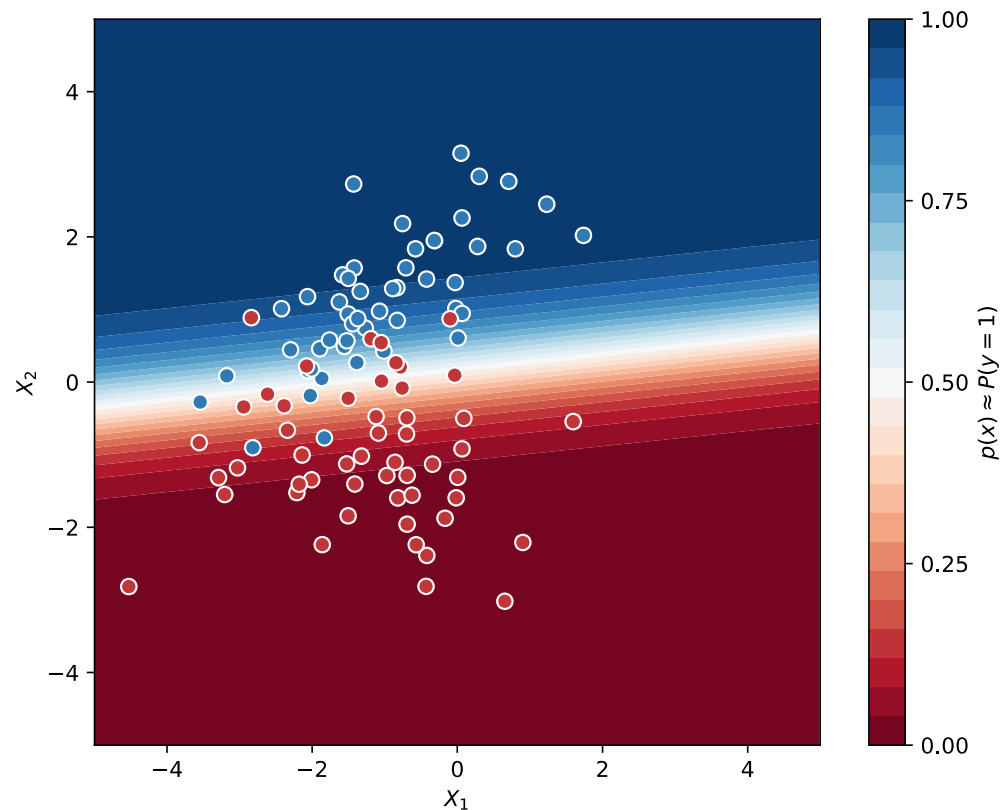
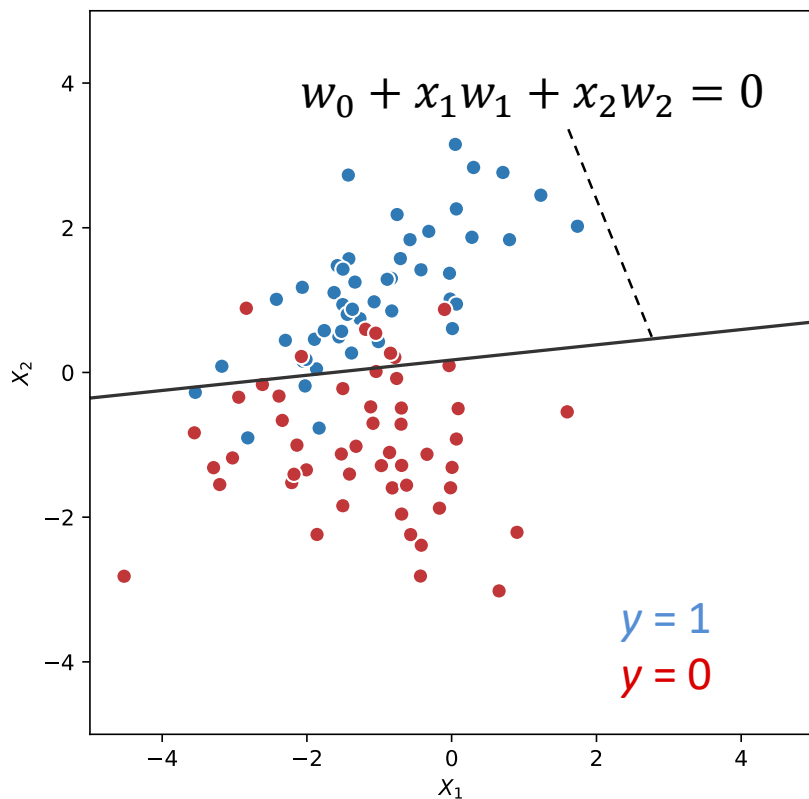
$$l(\mathbf{w}) = \prod_{i=1}^N p_1(\mathbf{x}_i)^{y_i} (1 - p_1(\mathbf{x}_i))^{1-y_i}$$

- Ak $y_i = 1$ $p_1(\mathbf{x}_i)^1 (1 - p_1(\mathbf{x}_i))^0 = p_1(\mathbf{x}_i)$
- Ak $y_i = 0$ $p_1(\mathbf{x}_i)^0 (1 - p_1(\mathbf{x}_i))^1 = 1 - p_1(\mathbf{x}_i) = p_0(\mathbf{x}_i)$
- Parametre \mathbf{w} môžeme vypočítať ľubovoľnou metódou ktorá maximalizuje funkciu vierohodnosti $l(\mathbf{w})$

Príklad – logistická regresia (7)

- Pre prahovú hodnotu 0.5 sa dá funkcia $p_1(x)$ zobrazit' aj geometricky ako deliaca (hyper)plocha oddeľujúca príklady jednotlivých tried
- Z priebehu logistickej funkcie platí, že $g(t) = 0.5$ ak $t = 0$, tzn.
$$g(w_0 + x_1 w_1 + \dots + x_M w_M) = 0.5$$
 ak
$$w_0 + x_1 w_1 + \dots + x_M w_M = 0$$
- Pre logistickú regresiu je deliaca hyperplocha lineárna (priamka pre X_1, X_2 , rovina pre X_1, X_2, X_3, \dots hyperrovina vo všeobecnosti) - lineárny klasifikátor

Príklad – logistická regresia (8)



Zobrazenie deliacej priamky modelu logistickej regresie a kontúr funkcie $p_1(x)$ - odhadu pravdepodobnosti pre triedu 1

Príklad – logistická regresia (9)

- Akú chybovú funkciu minimalizuje logistická regresia?
 - V $l(\mathbf{w})$ máme súčin členov pre jednotlivé príklady, pre chybovú funkciu potrebujeme súčet
 - Chybovú funkciu minimalizujeme, nie maximalizujeme
- 1. Súčin pravdepodobností v $l(\mathbf{w})$ pre jednotlivé príklady sa dá nahradiť za sumu pomocou logaritmu
 - $[\log(a \cdot b \cdot c \dots) = \log(a) + \log(b) + \log(c) + \dots]$
 - Funkcia $\log(l(\mathbf{w}))$ má maximum v rovnakom bode ako $l(\mathbf{w})$, takže sa optimálne riešenie \mathbf{w} nezmení
- 2. Pre minimalizovanie zmeníme znamienko funkcie na
 - $-\log(l(\mathbf{w}))$

Príklad – logistická regresia (10)

- Tzn. logistická regresia minimalizuje zápornú logaritmickú funkciu vierohodnosti $-\log(l(\mathbf{w}))$ – *negative log likelihood*
- Pre jeden príklad má po úprave výsledná chybová funkcia tvar:

$$L(y_i, p(\mathbf{x}_i)) = -[y_i \log(p_1(\mathbf{x}_i)) + (1 - y_i) \log(1 - p_1(\mathbf{x}_i))]$$

– $[\log(a^b) = b \cdot \log(a)]$

- Táto funkcia sa označuje aj ako **krížová entropia**
 - Pri jej minimalizovaní na celej tréningovej množine sa snažíme čo najviac zvýšiť odhad pravdepodobnosti $p_1(\mathbf{x}_i)$ pre príklady z triedy 1 a $1 - p_1(\mathbf{x}_i) = p_0(\mathbf{x}_i)$ pre príklady z triedy 0

Príklad – logistická regresia (11)

- Krížovú entropiu je možné zovšeobecniť pre ďalšie modely a pre klasifikáciu do K tried

$$L(y_i, f(\mathbf{x}_i)) = - \sum_{k=1}^K y_{i,k} \log(p_k(\mathbf{x}_i))$$

- Kde $y_{i,k} = 1$ ak príklad i patrí do triedy k , inak 0 a $p_k(\mathbf{x}_i)$ je odhad pravdepodobnosti modelu pre triedu k
- Ak môže byť príklad zaradený iba do jednej z K tried, pre odhad pravdepodobnosti musí platiť $0 \leq p_k(\mathbf{x}_i) \leq 1$,
 $\sum_{k=1}^K p_k(\mathbf{x}_i) = 1$

Príklad – logistická regresia (12)

- **Príklad**
- Priradenie triedy a predikciu môžeme považovať za dve pravdepodobnostné distribúcie pričom chceme dosiahnuť, aby sa čo najviac podobali pre všetky príklady.

$y_{i,1}$	0	0.001	$p_1(x_i)$
$y_{i,2}$	0	0.0015	$p_2(x_i)$
	
$y_{i,k}$	1	0.78	$p_k(x_i)$
	
$y_{i,K}$	0	0	$p_K(x_i)$

- Pri učení sa snažíme čo najviac zvýšiť odhad pre správnu triedu k (resp. znížiť odhad pre všetky ostatné triedy, keďže platí ohraničenie $\sum_{k=1}^K p_k(x_i) = 1$)

Aký algoritmus vybrať? (1)

- Predpokladajme, že závislosť medzi \mathbf{X} a Y najpresnejšie popisuje funkcia $F^*(\mathbf{x})$, tzn. funkcia, ktorá má minimálnu zovšeobecnenú chybu na všetkých možných prípadoch
- $F^*(\mathbf{x})$ je teda naša cieľová funkcia, ktorú hľadáme
- Ak si zvolíme model s obmedzenými funkciami (napr. lineárnymi), aj najlepšia možná funkcia $f^*(\mathbf{x})$, ktorá sa v ňom bude nachádzať bude veľmi odlišná od cieľovej funkcie $F^*(\mathbf{x})$ a teda bude mať veľkú chybu vo všeobecnosti
- **Nedokážeme presne aproximovať zložité závislosti jednoduchou funkciou**

Aký algoritmus vybrať? (2)

- Naopak, ak zvolíme veľkú množinu komplexných funkcií, je pravdepodobnejšie že sa medzi nimi bude nachádzať aj $F^*(\mathbf{x})$
 - Alebo sa aspoň nebude najlepšia funkcia modelu $f^*(\mathbf{x})$ veľmi odlišovať od $F^*(\mathbf{x})$ (tzn. zníži sa chyba aproximácie)
- Na druhej strane sa však zvýši variancia učenia a môže dochádzať k preučeniu

Variancia učenia a preučenie (1)

- Máme iba obmedzený počet tréningových príkladov
- Ak zvolíme model s komplexnými funkciami - pre danú tréningovú množinu môže existovať veľa navzájom odlišných funkcií, ktoré budú mať minimálnu chybu na tréningových príkladoch
 - Niektoré z týchto funkcií sa však môžu výrazne odlišovať od cieľovej funkcie $F^*(x)$
- Ak algoritmus zvolí funkciu, ktorá síce má minimálnu chybu na tréningových dátach, ale veľkú chybu vo všeobecnosti pre ostatné prípady, dochádza k tzv. **preučeniu**

Variancia učenia a preučenie (2)

- Možnosti preučenia sa nemôžeme nikdy principiálne vyhnúť s úplnou istotou!
- Pre danú konečnú tréningovú množinu sa ju však môžeme pokúsiť znížiť **regulovaním učenia** – tzn. zvolíme si jednoduchší, viac ohraničený model

Koľko tréningových príkladov potrebujeme? (1)

- Predpokladajme, že môžeme vybrať náhodne N tréningových príkladov zo všetkých možných prípadov:
 - Vždy keď si zvolíme iné tréningové príklady, algoritmus môže nájsť inú funkciu $f(x)$ s minimálnou tréningovou chybou
- Ak máme komplexný model – aj malá zmena tréningových príkladov môže spôsobiť výber veľmi odlišnej funkcie $f(x)$ (tzn. veľkú variáciu učenia) – a teda narastá aj možnosť preučenia

Koľko tréningových príkladov potrebujeme? (2)

- Ak budeme pridávať ďalšie tréningové príklady:
 - Chybu bude algoritmus minimalizovať na väčšom počte príkladov, tzn. funkcia $f(x)$ by sa mala približovať k cieľovej funkcii $F^*(x)$ a chyba aproximácie by sa mala znížiť
 - Viac príkladov viac obmedzí výber rôznych funkcií $f(x)$ s minimálnou tréningovou chybou, tzn. variancia učenia sa zníži (a teda sa aj zníži možnosť preučenia)

Zovšeobecnená chyba funkcie $f(\mathbf{x})$ - zhrnutie

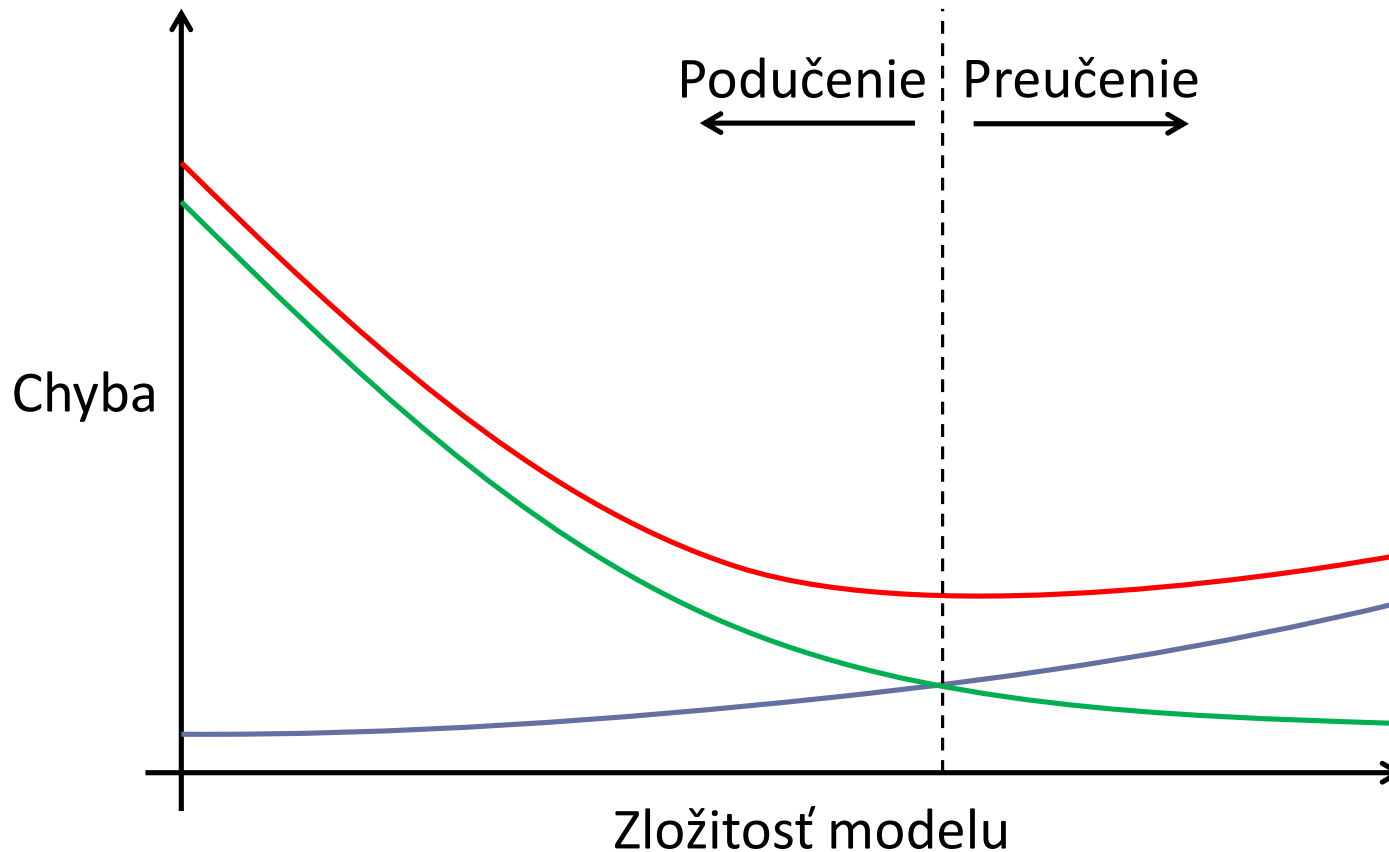
- Je chyba určená na všetkých možných prípadoch
- Môžeme ju vyjadriť rozdielom medzi naučenou funkciou $f(\mathbf{x})$ a neznámou cieľovou funkciou $F^*(\mathbf{x})$
- Nedá sa priamo určiť lebo nepoznáme $F^*(\mathbf{x}) = y$, ale vieme že je ovplyvnená dvoma protichodnými zložkami:
- **Zovšeobecnená chyba = Chyba aproximácie (Bias) + Chyba variancie učenia**

Koľko trénovacích príkladov potrebujeme? (3)

- Zovšeobecnená chyba = Chyba aproximácie (Bias) + Chyba variácie učenia
- Pre malú trénovaciu množinu môže prevládať vplyv variácie učenia a skôr môže dôjsť k preučeniu – preto uprednostníme jednoduchší model
- Pri narastajúcom počte príkladov sa bude variancia zmenšovať, takže zvolíme zložitejší model, aby sme čo najviac minimalizovali chybu aproximácie

Chyba = Bias + Variancia

Zovšeobecnená chyba = Chyba aproximácie (Bias) + Chyba variacie



Aký dobrý je náš model? (1)

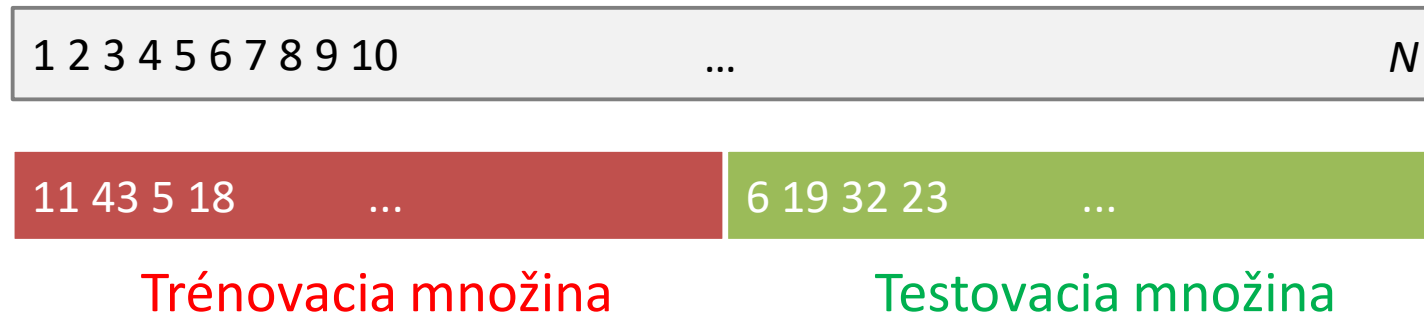
- Nikdy nebudeme vedieť úplne vyhodnotiť zovšeobecnenú chybu
 - A teda nevieme ani presne určiť, ako sa odlišuje niektorá funkcia modelu $f(\mathbf{x})$ od $F^*(\mathbf{x})$
- Podľa **štatistickej teórie učenia** však vieme s veľkou istotou určiť, s akou chybou sa môže funkcia $f(\mathbf{x})$ naučená na N tréningových príkladoch maximálne odlišovať od $f^*(\mathbf{x})$ – tzn. od najpresnejšej funkcie patriacej do modelu
 - Tento odhad je platný nezávisle od $F^*(\mathbf{x})$ a od vybranej tréningovej množiny

Aký dobrý je náš model? (2)

- Pre daný model však nemusíme mať istotu, že cieľová funkcia $F^*(x)$ určite patrí do modelu, resp. že sa $f^*(x)$ od nej výrazne neodlišuje
- Navyše ak chceme spresniť odhad chyby podľa štatistickej teórie, musíme uvažovať veľký počet tréningových príkladov
- A tak isto nás väčšinou zaujíma presnosť algoritmu na danom probléme, nie na všetkých možných problémoch
- Preto sa pri praktickom riešení používa iná metóda na odhad zovšeobecnej chyby

Aký dobrý je náš model? (3)

- Náhodne si rozdelíme dáta na **trénovaciú** a **testovaciú** množinu, na trénovacej množine naučíme funkciu $f(x)$ a vypočítame jej chybu na testovacej množine



Aký dobrý je náš model? (4)

- Testovacia chyba je odhadom zovšeobecnenej chyby
- Pri inom náhodnom rozdelení sa môže testovacia chyba zmeniť:
 - Prejaví sa variancia učenia, pretože sme zmenili trénovacie dáta
 - Prejaví sa variancia náhodného výberu testovacích príkladov
- Tzn., aby sme mali čo najpresnejší odhad, musíme mať dostatočný počet príkladov aj v testovacej aj v trénovacej množine

Aký dobrý je náš model? (5)

- Krížová validácia

- Učenie na väčšej trénovacej množine – zníži sa vplyv variance učenia
- Testovanie vyhodnotené na všetkých dátach – zníži sa vplyv variance testovacej množiny

1	2	3	4	5	6	7	8	9	10	...	N								
11	43	5	...	6	1	23	...	7	33	15	...	8	27	53	...	18	20	2	...
11	43	5	...	6	1	23	...	7	33	15	...	8	27	53	...	18	20	2	...
...																			
11	43	5	...	6	1	23	...	7	33	15	...	8	27	53	...	18	20	2	...

Učenie prediktívnych modelov – zhrnutie (1)

- Úloha je na základe trénovacej množiny nájsť funkciu, pomocou ktorej vieme čo najpresnejšie predpovedať hodnotu cieľového atribútu podľa vstupných atribútov vo všeobecnosti pre všetky prípady
- Algoritmus učenia = Chybová funkcia + Model + Optimalizačná metóda
- Algoritmus učenia zvolí funkciu zo zvolenej množiny funkcií – modelu tak aby sa minimalizovala chyba na trénovacích príkladoch

Učenie prediktívnych modelov – zhrnutie (2)

- Zovšeobecnená chyba = chyba aproximácie + chyba variancie učenia
 - Príliš jednoduchý model – veľká chyba aproximácie – podučenie
 - Príliš zložitý model – veľká varianca učenia – preučenie
- Zovšeobecnenú chybu vieme len odhadnúť:
 - Testovacia chyba na nezávislej náhodne vybranej testovacej množine – pozor na varianciu učenia a varianciu výberu testovacích príkladov
 - Krížová validácia

Lineárne modely na nelineárnych dátach

- Väčšina pozorovaných závislostí je nelineárna
- Transformácia pomocou bázových funkcií
 - Jednoduchý spôsob ako rozšíriť lineárny model na nelineárne dáta:
 - Zvolíme si nelineárne (tzv. bázové) funkcie B_1, B_2, \dots, B_K , pomocou ktorých transformuje pôvodné vstupné atribúty
 - Na transformovaných dátach naučíme lineárny model, keď ho však aplikujeme na pôvodné dáta, vieme reprezentovať aj nelineárne závislosti

Lineárne modely na nelineárnych dátach

- **Príklad:**

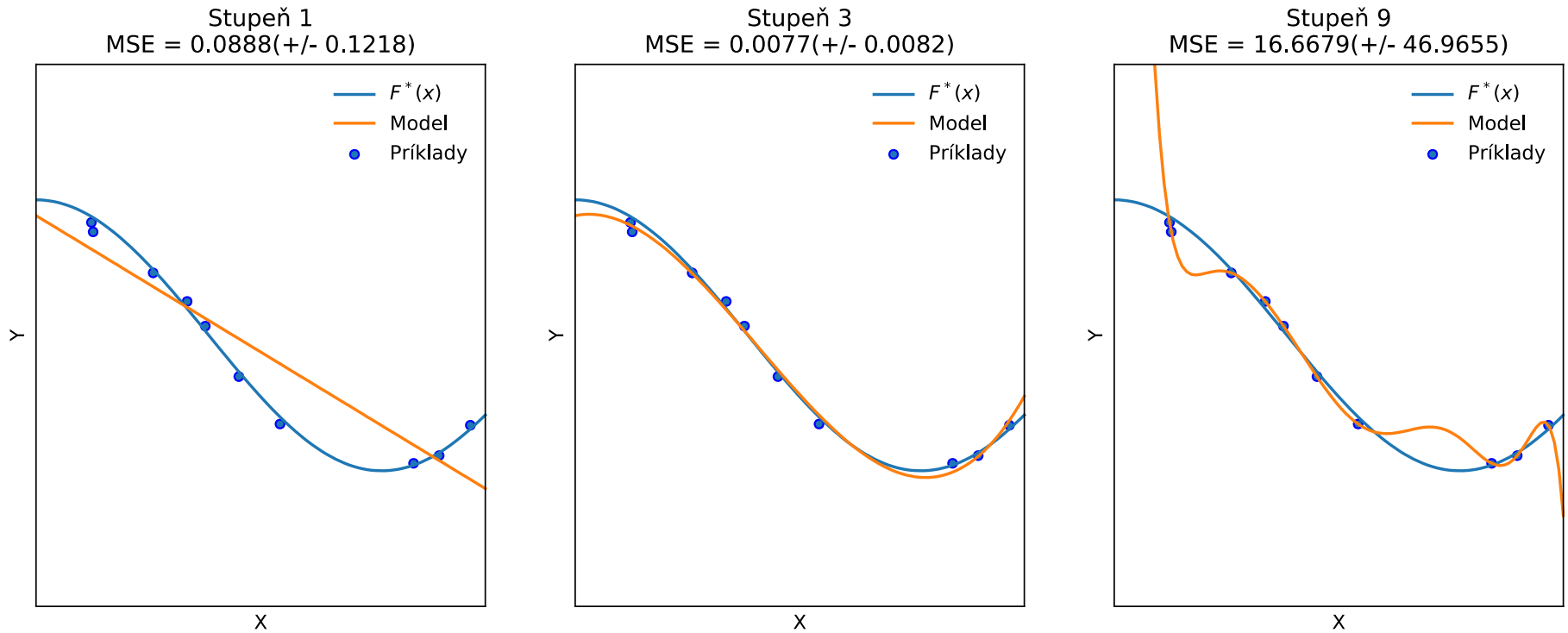
- Polynomiálna regresia

- Ako transformačné funkcie použijeme členy polynómu zvoleného stupňa, tzn. napr. pre dva atribúty X_1 , X_2 a stupeň polynómu 2 budú transformačné funkcie $B_1 = X_1$, $B_2 = X_2$, $B_3 = X_1X_2$, $B_4 = X_1^2$, $B_5 = X_2^2$
- Výsledný model bude mať tvar:

$$f(\mathbf{x}, \mathbf{w}) = w_0 + x_1w_1 + x_2w_2 + x_1x_2w_3 + x_1^2w_4 + x_2^2w_5$$

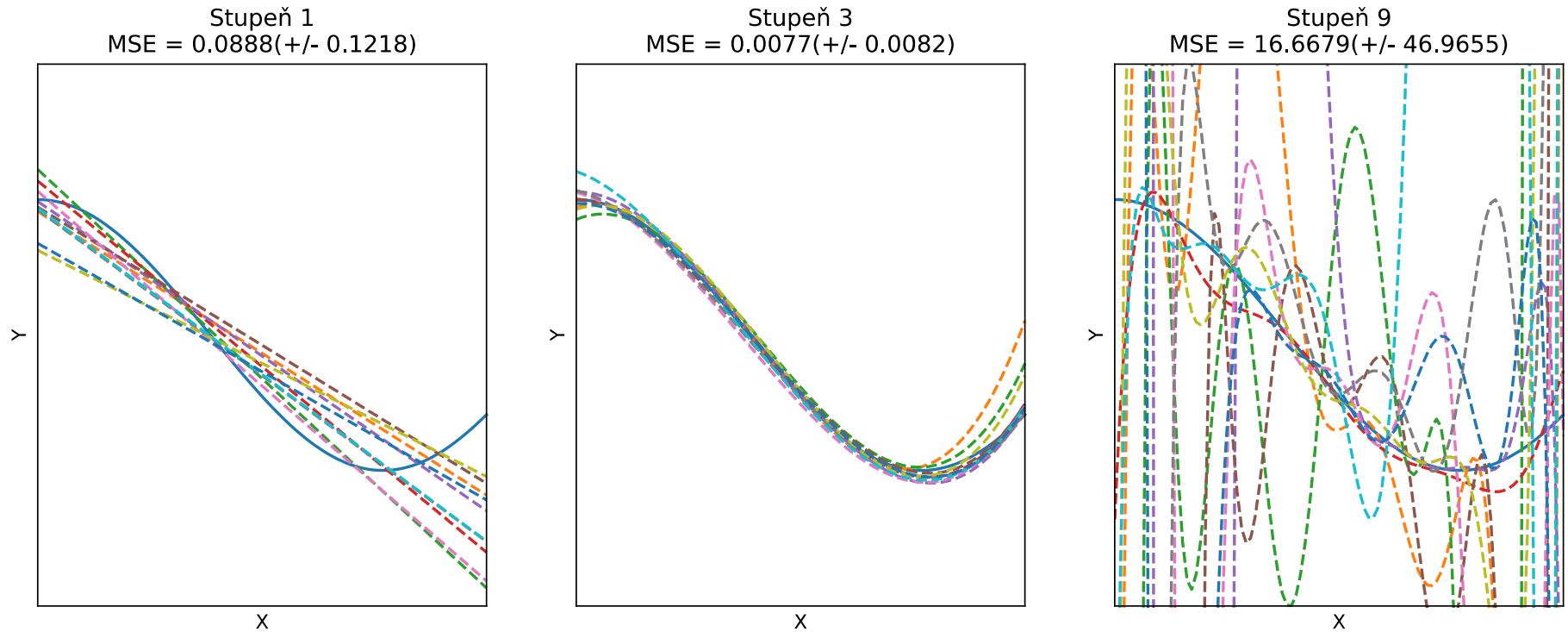
- Okrem členov polynómu môžeme použiť aj iné nelineárne funkcie, napr.: x_i/x_j , $\log(x_i)$, ...

Polynomiálna regresia – príklad (1)



MSE je priemerná kvadratická chyba pri 10-násobnej krížovej validácii,
štandardná odchýlka chyby udáva odhad variancie učenia

Polynomiálna regresia – príklad (2)



Priebeh 10 modelov naučených na náhodne vybraných tréningových množinách s 10 príkladmi. Pri jednoduchom lineárnom modeli (stupeň 1) došlo k podučeniu, pri polynóme 9. stupňa k výraznému preučeniu