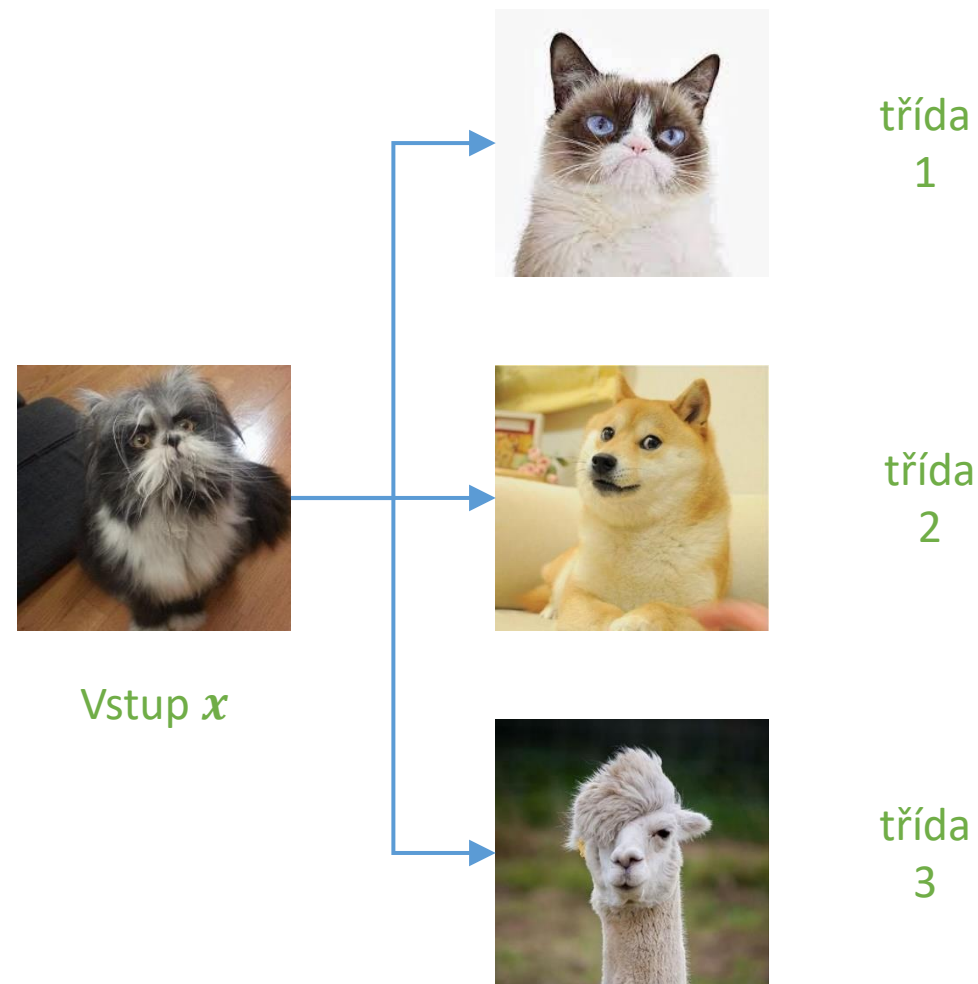


Počítačové vidění

Hluboké neuronové sítě, klasifikace obrázků

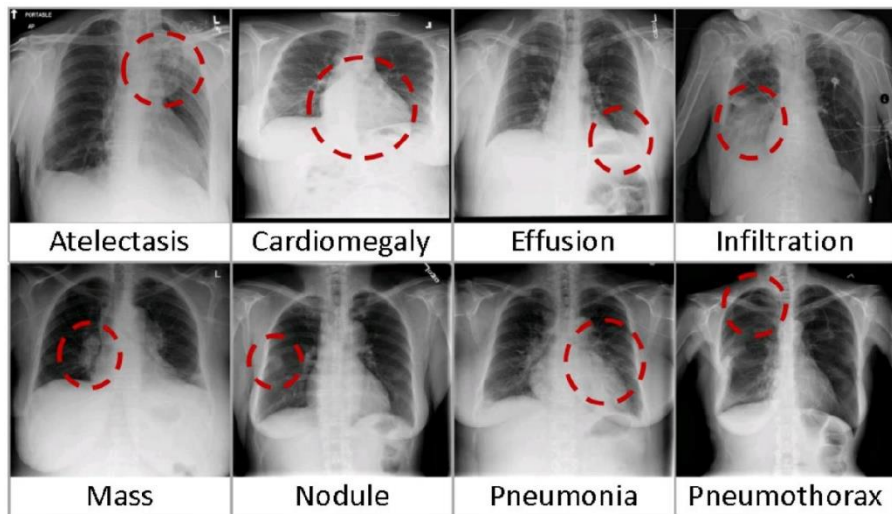
Úloha klasifikace obrázků

- Vstupem x_n je RGB obrázek
- Úkolem zařadit x_n do jedné ze tří předdefinovaných kategorií (tříd):
 1. “kočka”
 2. “pes”
 3. “alpaka”
- Počet tříd označíme jako K
- Výstupem bude jedno z následujících:
 - skóre jednotlivých tříd s_n ,
 - pravděpodobnost jednotlivých tříd p_n ,
 - nebo celé číslo $\hat{y}_n \in \{1, \dots, K\}$



Proč klasifikace?

Klasifikace je zajímavá a užitečná sama o sobě



Amanita_fulva1



Lepista_saeva48



Coprinopsis_atramentaria35



Morchella_esculenta33



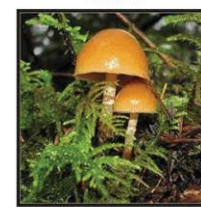
Amanita_fulva11



Auricularia_auricula-judae22



Pleurocybella_porrigena13



Galerina_sulciiceps19

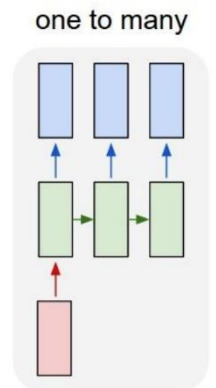


Amanita_arocheae37

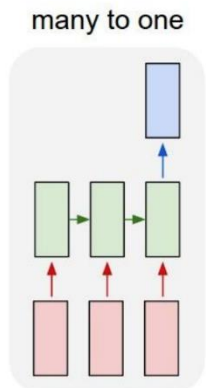


Clavulinaceae15

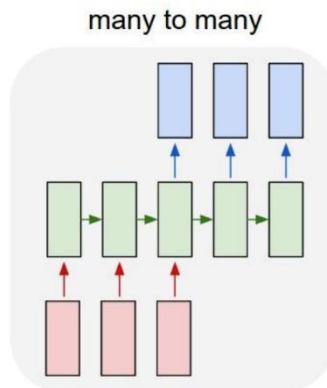
zároveň ale tvoří základ mnoha dalších aplikací



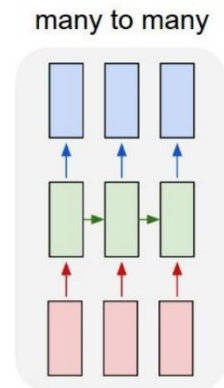
tagování obrázku,
generování hudby



analýza sentimentu,
klasifikace videa



strojový překlad,
rozpznávání řeči



rozpoznávání fonémů, jazykový
model



Návrh a trénování lineárního klasifikátoru

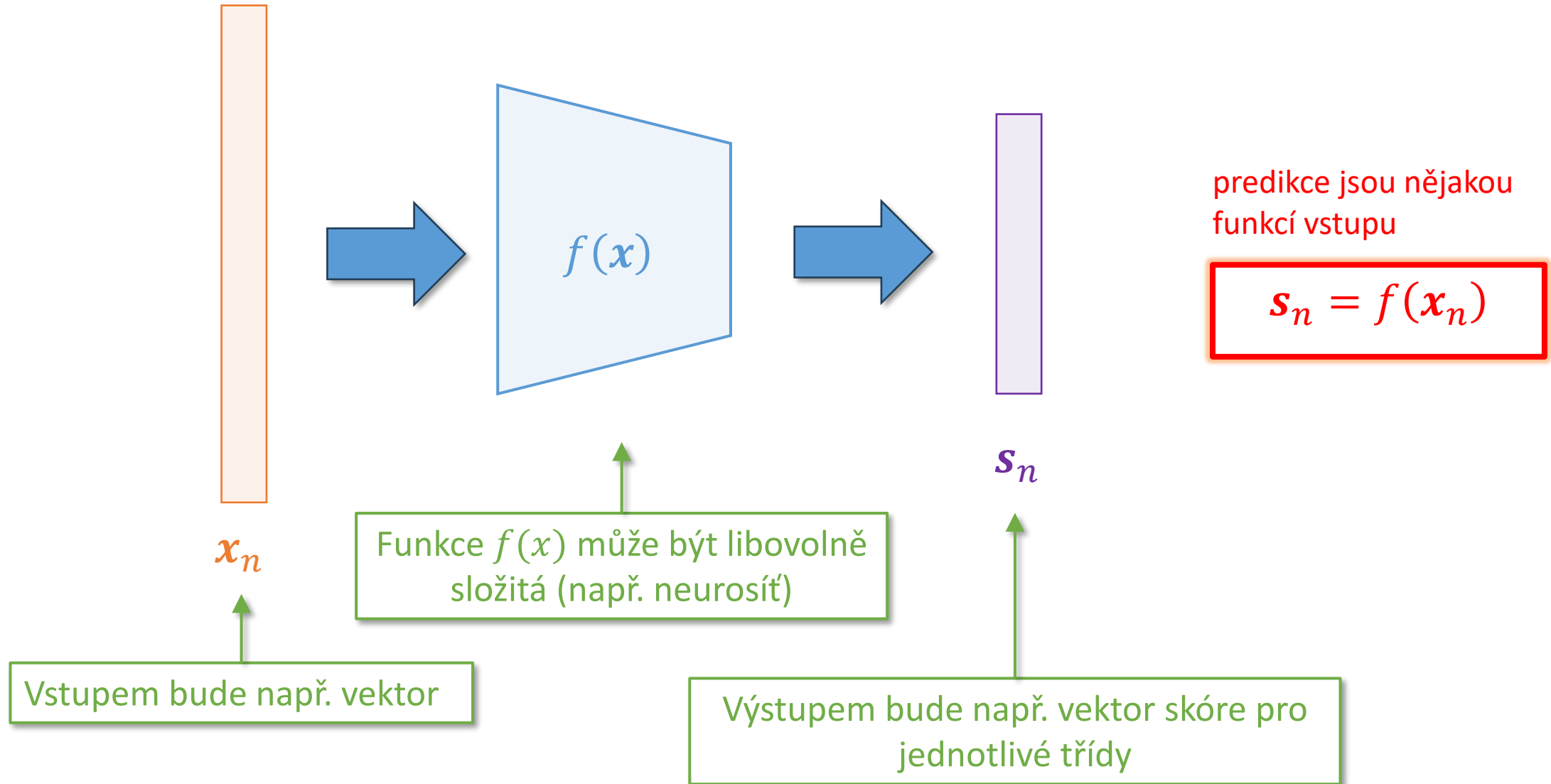
1. Navrhujeme diskriminativní klasifikační funkci s upravitelnými parametry
2. Kvantifikujeme její úspěšnost klasifikace nějakým kritériem
3. Nastavíme parametry klasifikátoru tak, abychom optimalizovali zvolené kritérium

Návrh a trénování lineárního klasifikátoru

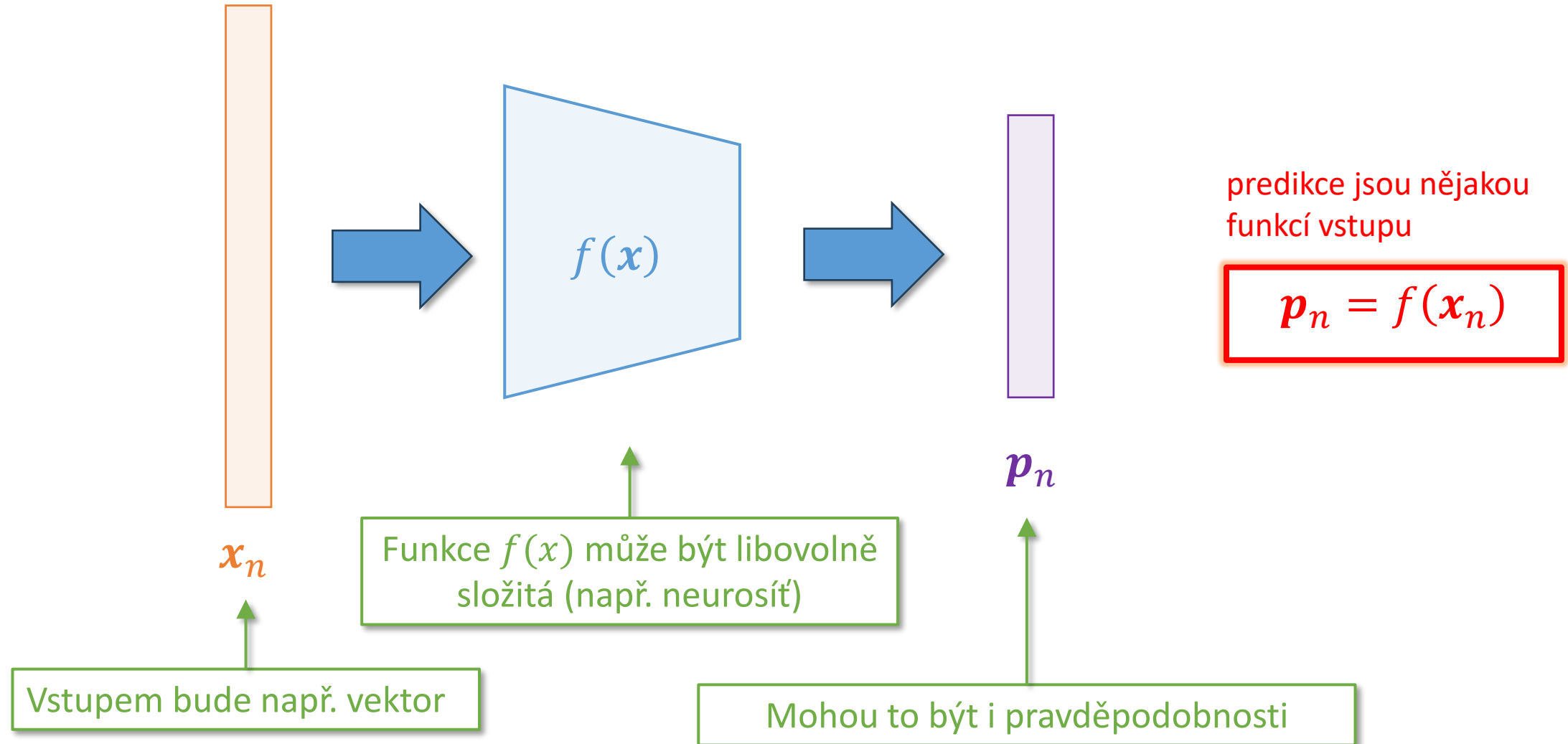
1. Navrhujeme diskriminativní klasifikační funkci s upravitelnými parametry
2. Kvantifikujeme její úspěšnost klasifikace nějakým kritériem
3. Nastavíme parametry klasifikátoru tak, abychom optimalizovali zvolené kritérium

Diskriminativní klasifikace

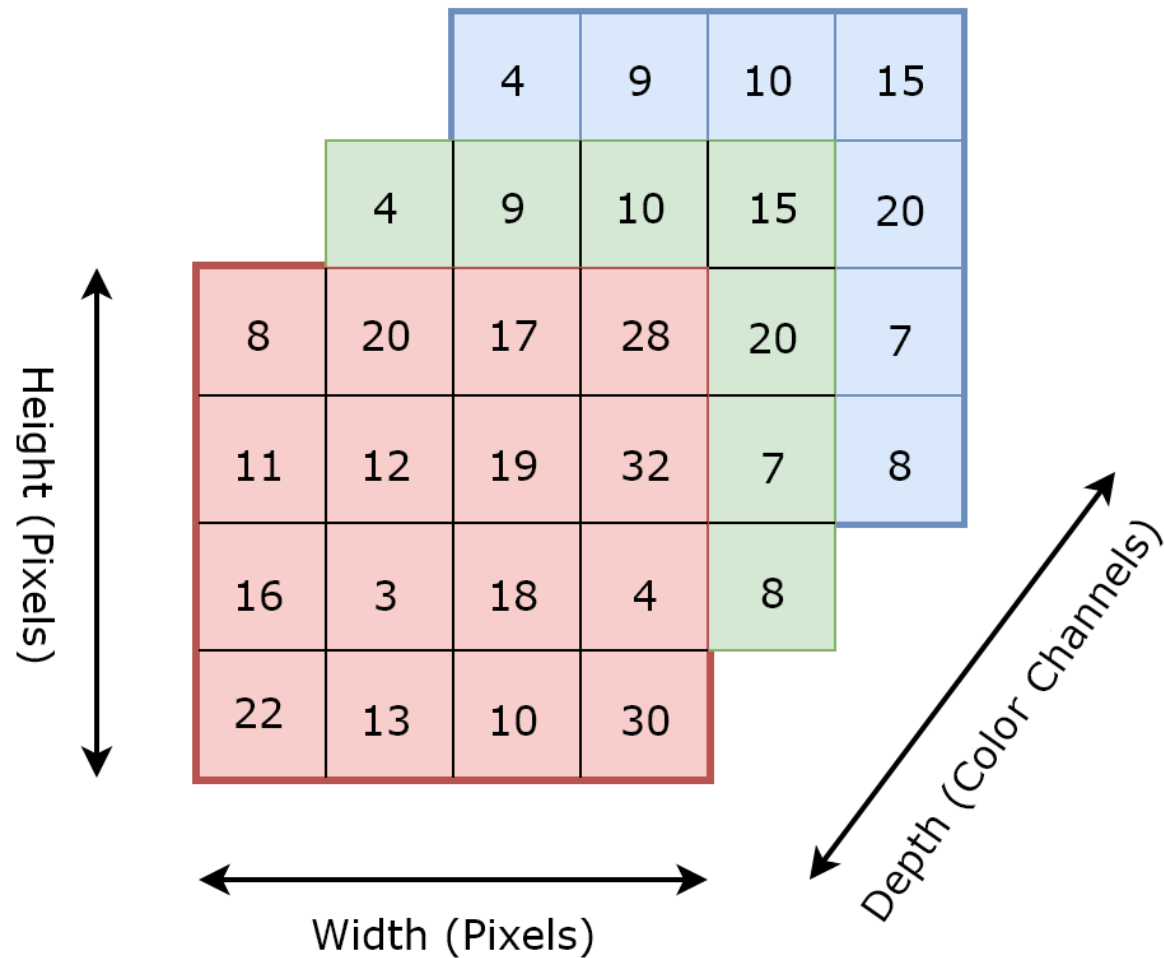
Diskriminativní klasifikátor



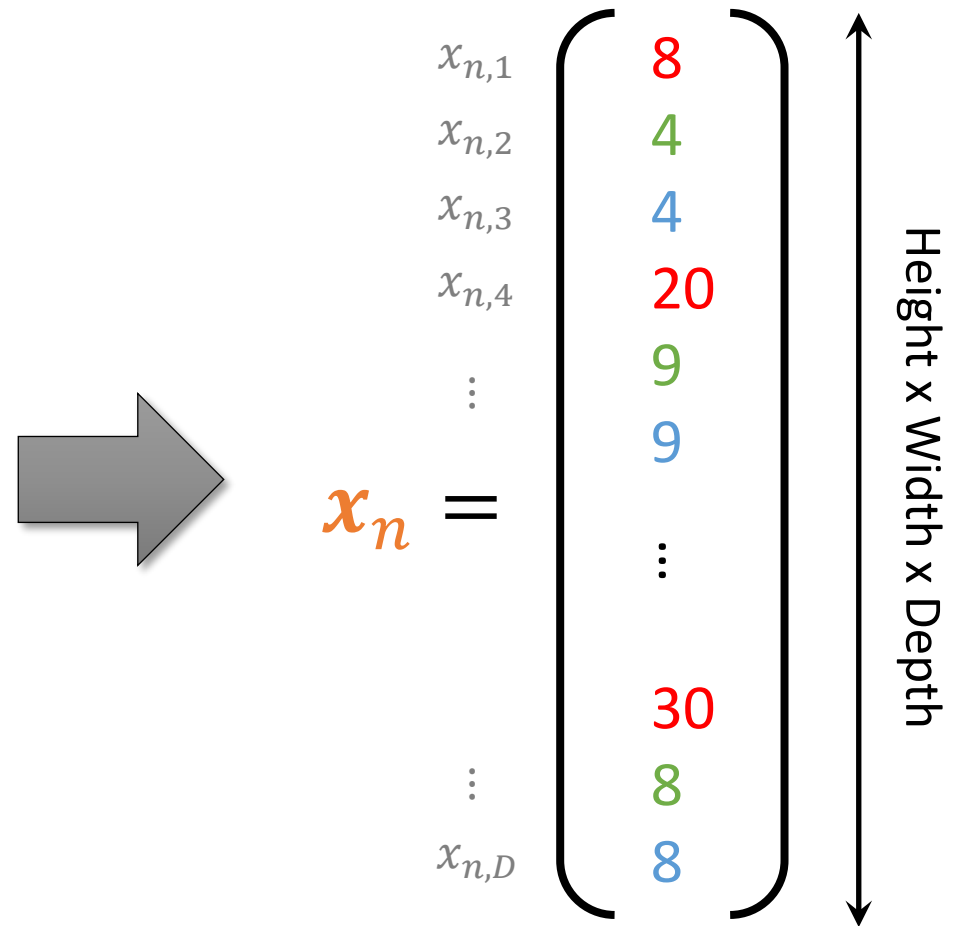
Diskriminativní klasifikátor



Reprezentace RGB obrázku jako vektoru



Tensor tvaru (výška, šířka, hloubka) (HWC formát)



Vektor délky D = výška x šířka x hloubka

Lineární klasifikátor

- Lineární model předpokládá afinní* vztah mezi skóre třídy s_n a vstupem x_n

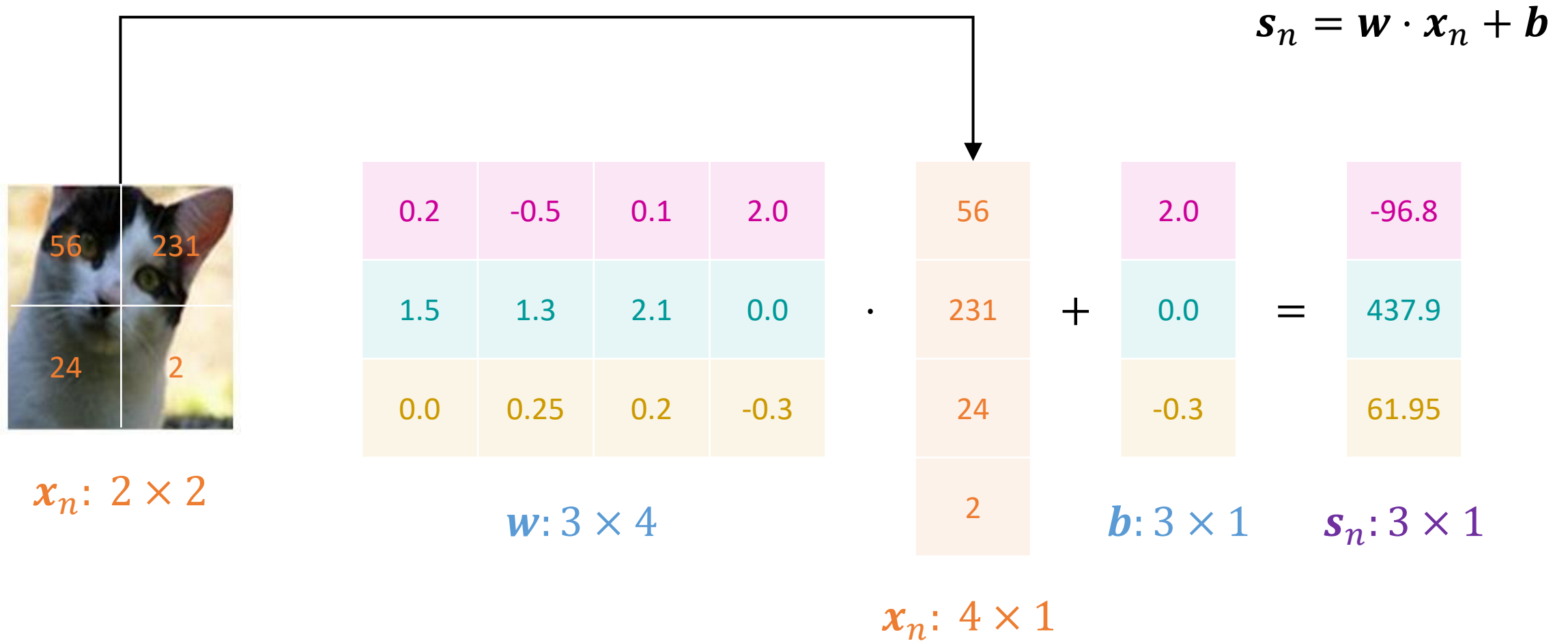
$$s_n = w \cdot x_n + b$$

kde

- x_n je sloupcový vektor o rozměru D
- w je matice vah klasifikátoru s rozměry $K \times D$
- b je sloupcový vektor biasů klasifikátoru s rozměrem K

} parametry klasifikátoru

Lineární predikce skóre: příklad



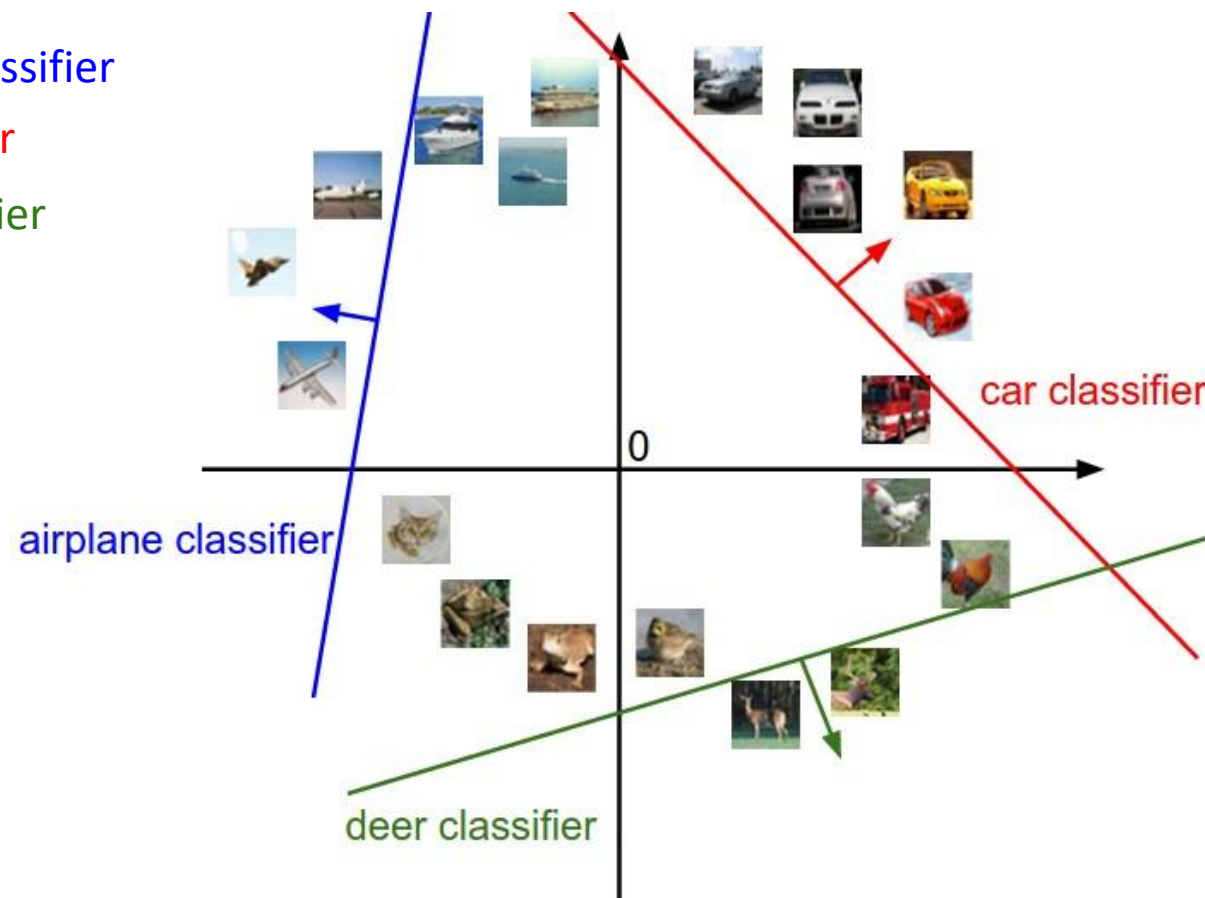
příklad: <https://cs231n.github.io/linear-classify/>

https://web.eecs.umich.edu/~justincj/slides/eecs498/WI2022/598_WI2022_lecture03.pdf

Geometrická interpretace

$$\mathbf{w} = \begin{bmatrix} \mathbf{w}_{1,1} & \mathbf{w}_{1,2} & \dots & \mathbf{w}_{1,D} \\ \mathbf{w}_{2,1} & \mathbf{w}_{2,2} & \dots & \mathbf{w}_{2,D} \\ \mathbf{w}_{3,1} & \mathbf{w}_{3,2} & \dots & \mathbf{w}_{3,D} \end{bmatrix} \begin{array}{l} \text{airplane classifier} \\ \text{car classifier} \\ \text{deer classifier} \end{array}$$

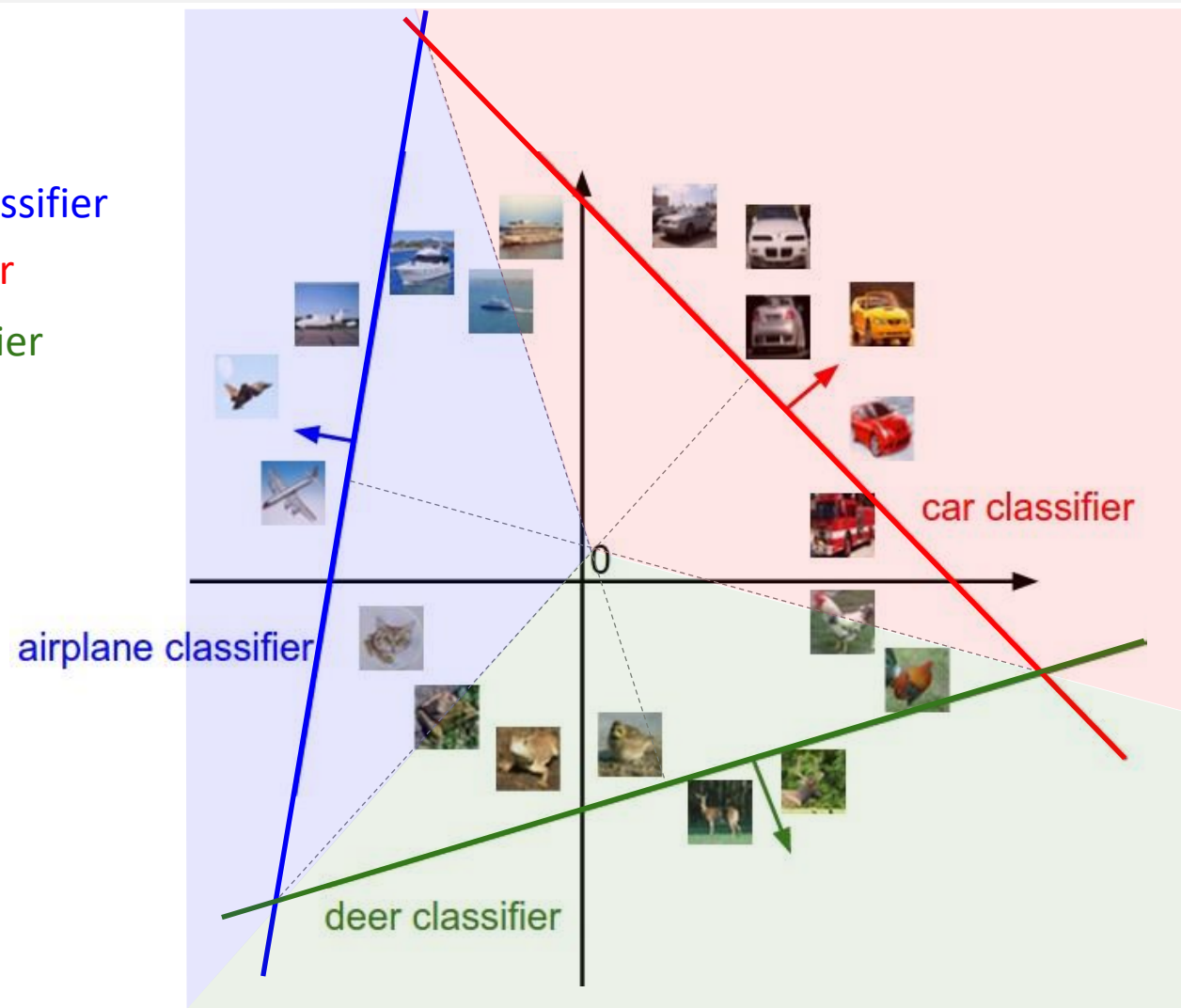
Každý řádek matice \mathbf{w} je binární klasifikátor diskriminující třídu k od ostatních



Geometrická interpretace

$$\mathbf{w} = \begin{bmatrix} \mathbf{w}_{1,1} & \mathbf{w}_{1,2} & \dots & \mathbf{w}_{1,D} \\ \mathbf{w}_{2,1} & \mathbf{w}_{2,2} & \dots & \mathbf{w}_{2,D} \\ \mathbf{w}_{3,1} & \mathbf{w}_{3,2} & \dots & \mathbf{w}_{3,D} \end{bmatrix} \begin{array}{l} \text{airplane classifier} \\ \text{car classifier} \\ \text{deer classifier} \end{array}$$

Každý řádek matice \mathbf{w} je binární klasifikátor diskriminující třídu k od ostatních



Návrh a trénování lineárního klasifikátoru

1. Navrhne diskriminativní klasifikační funkci s upravitelnými parametry
2. Kvantifikujeme její úspěšnost klasifikace nějakým kritériem
3. Nastavíme parametry klasifikátoru tak, abychom optimalizovali zvolené kritérium

Klasifikační kritérium

Softmax cross entropy

Klasifikační kritérium

- Klasifikátor predikuje $\hat{y}_n \in \{1, \dots, K\}$

$$\hat{y}_n = \operatorname{argmax}_k \mathbf{s}_n$$

- Pro každý obrázek \mathbf{x}_n přitom známe správnou odpověď $y_n \in \{1, \dots, K\}$ (target)
- Porovnáním \hat{y}_n vs y_n (nebo vs \mathbf{s}_n) můžeme vyčíslit, jak dobrá/špatná predikce je
- Celkem máme **trénovací dataset** \mathbf{X} s N obrázky a tedy i páry (\mathbf{x}_n, y_n)
- Nakonec tedy můžeme spočítat tzv. loss

$$L(\mathbf{X}) = \frac{1}{N} \sum_{n=1}^N L_n(\mathbf{s}_n, y_n)$$

Multiclass cross entropy

- Logistická regrese definuje loss jako tzv. **křížovou entropii**

$$l_n = - \sum_{k=1}^K p_{n,k} \cdot \log(\hat{p}_{n,k})$$

Pro jeden obrázek

- kde

$$\mathbf{p}_n = [p_{n,1}, \dots, p_{n,K}]^T \quad \dots \text{cílové rozdělení (ground truth / target)}$$

$$\hat{\mathbf{p}}_n = [\hat{p}_{n,1}, \dots, \hat{p}_{n,K}]^T \quad \dots \text{výstupní pravd. (predikce) klasifikátoru}$$

jsou vektory, na které nahlížíme jako na diskrétní pravděpodobnostní rozdělení

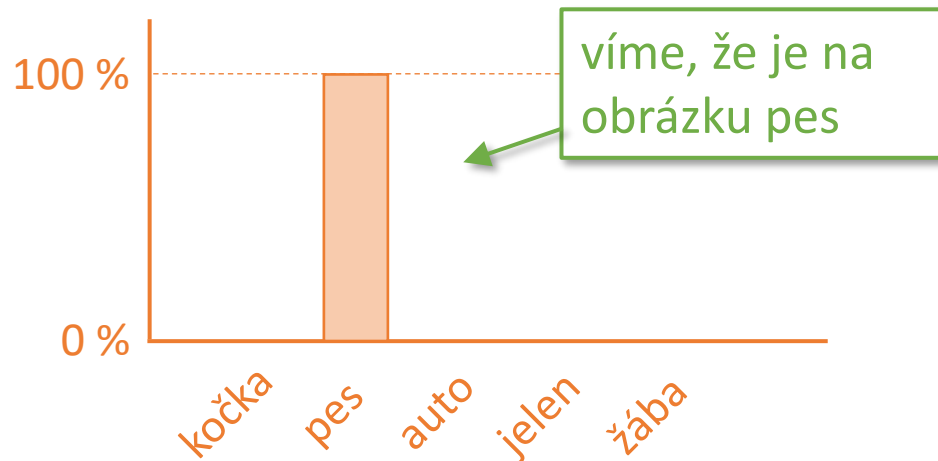
→ cross entropy = rozdíl mezi dvěma pravděpodobnostními rozděleními

Multiclass cross entropy

$$l_n = - \sum_{k=1}^K p_{n,k} \cdot \log(\hat{p}_{n,k})$$

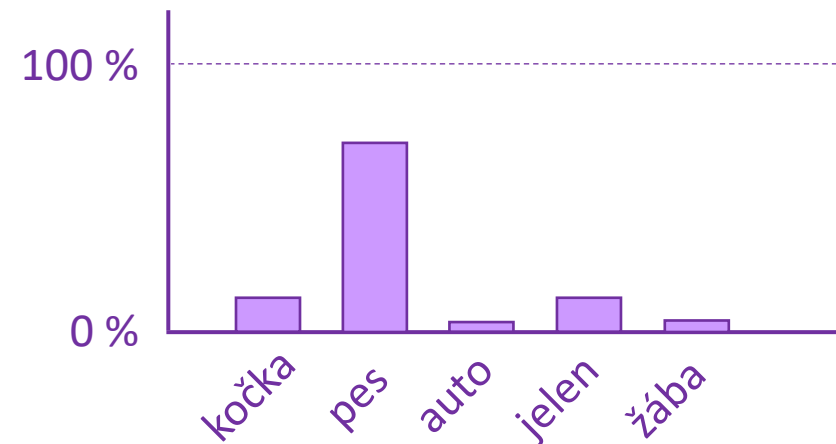
$$\mathbf{p}_n = [p_{n,1}, \dots, p_{n,K}]^T$$

cílové rozdělení (ground truth / target)



$$\hat{\mathbf{p}}_n = [\hat{p}_{n,1}, \dots, \hat{p}_{n,K}]^T$$

výstup (predikce) klasifikátoru



Převod číselného označení třídy na rozdělení: one hot encoding

- Značka y_n pro každý obrázek je celé číslo, tj. $y_n \in \{1, \dots, K\}$
- Pokud počet tříd $K = 5 \rightarrow$ požadované rozdělení je

$$y_n = 2 \quad \Rightarrow \quad \mathbf{p}_n = [0, 1, 0, 0, 0]^\top$$

$$y_n = 5 \quad \Rightarrow \quad \mathbf{p}_n = [0, 0, 0, 0, 1]^\top$$

Převod výstupních skóre modelu na rozdělení: softmax

- Normalizuje vektor skóre \mathbf{s}_n tak, že výstup lze interpretovat jako pravděpodobnosti
- Pravděpodobnost, že na obrázku \mathbf{x}_n je objekt třídy k definuje jako

$$\hat{p}_{n,k} = P(\text{třída } k | \mathbf{x}_n) = \frac{e^{s_{n,k}}}{\sum_{i=1}^K e^{s_{n,i}}}$$

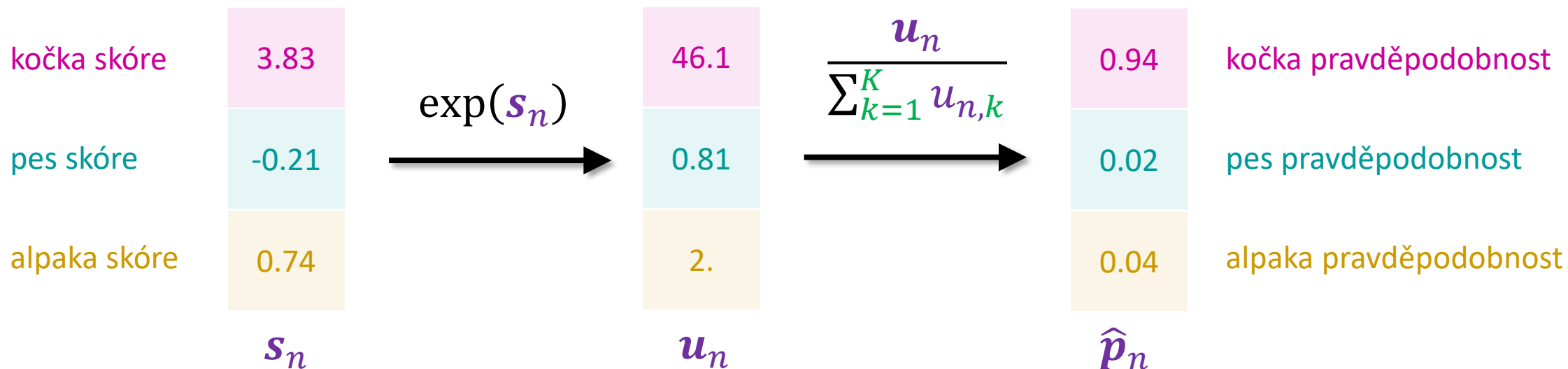
- Výstupem K -dimezionální vektor $\hat{\mathbf{p}}_n$ pravděpodobností jednotlivých tříd

$$\hat{\mathbf{p}}_n = [\hat{p}_{n,1}, \dots, \hat{p}_{n,K}]^T, \quad 0 \leq \hat{p}_{n,k} \leq 1, \quad \sum_{k=1}^K \hat{p}_{n,k} = 1$$

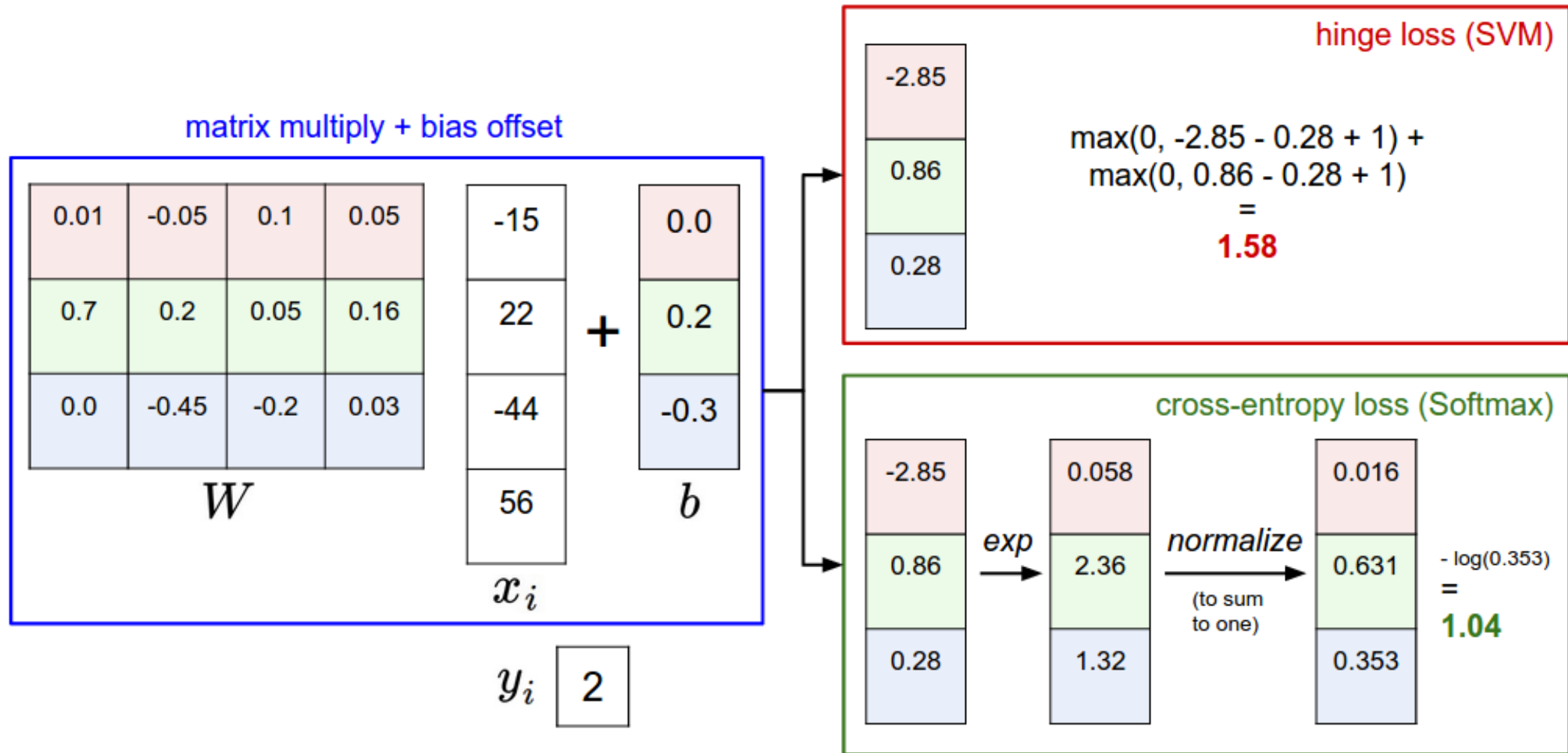
- Chová se jako “měkké” maximum: exponenciováním se zvýrazní rozdíly (nejvyšší hodnota vynikne), až teprve pak se normalizuje (ostatní jsou staženy k nule)

Softmax: příklad

$$\hat{p}_{n,k} = \frac{e^{s_{n,k}}}{\sum_{i=1}^K e^{s_{n,i}}}$$



Lineární model a klasifikační loss



Návrh a trénování lineárního klasifikátoru

1. Navrhujeme diskriminativní klasifikační funkci s upravitelnými parametry
2. Kvantifikujeme její úspěšnost klasifikace nějakým kritériem
3. Budeme upravovat parametry a poznamenávat si výsledek (hodnotu kritéria)

Celkový loss

- Chybovost klasifikátoru (loss)

$$L(\mathbf{X}) = \frac{1}{N} \sum_{n=1}^N L_n(\mathbf{s}_n, \mathbf{y}_n)$$

- Je funkce, která kromě dat závisí na parametrech klasifikátoru, protože

$$\mathbf{s}_n = \mathbf{w} \cdot \mathbf{x}_n + \mathbf{b}$$

- Celkově tedy vyhodnocujeme funkci $L(\mathbf{X}, \boldsymbol{\theta})$ závislou na datech \mathbf{X} a parametrech $\boldsymbol{\theta}$

Strojové učení je optimalizace


$$\theta^* = \underset{\theta}{\operatorname{argmin}} L(\mathbf{X}, \theta)$$

Minimalizace funkce metodou největšího spádu

Minimalizace funkce

- Mějme jednorozměrnou funkci

$$f(x): \mathbb{R} \rightarrow \mathbb{R}$$

tj. vstup i výstup jsou reálná čísla (skaláry)

- Chceme najít bod x^* , kde funkce f nabývá minimální hodnoty, tj.

$$x^* = \operatorname{argmin}_x f(x)$$

Metoda největšího spádu (Gradient Descent, GD)

- Metoda největšího spádu využívá derivaci pro zjištění, kterým směrem funkce roste
 1. V aktuální pozici x spočítá derivaci $\frac{df(x)}{dx}$ (exaktně nebo numericky)
 2. Vylepší aktuální odhad x posunutím se ve směru opačném ke směru růstu, formálně

$$x' := x - \gamma \cdot \frac{df(x)}{dx}$$

Výsledná velikost **rozdílu** mezi původním odhadem a novým odhadem je ovlivněna

- a) velikostí kroku γ , kterou se výsledek derivace škáluje
 - b) absolutní hodnotou derivace $\frac{df(x)}{dx}$, tj. strmostí $f(x)$ v aktuálním bodě x
3. S novým odhadem x' se uvedený postup opakuje
- Na začátku vyžaduje nějaký počáteční odhad x

Příklad GD: funkce $f(x) = x^4 + x + 2$

- Mějme např. funkci

$$f(x) = x^4 + x + 2$$

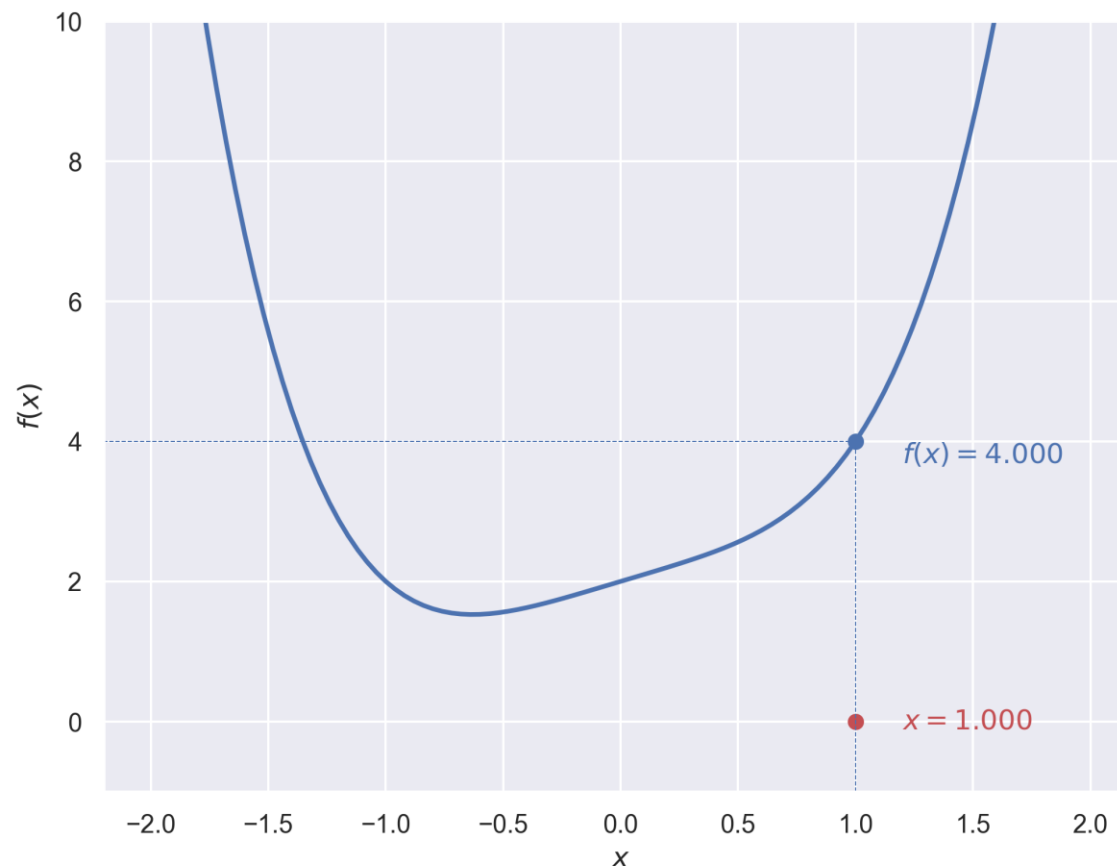
- Hledáme pozici x^* jejího minima, tj.

$$x^* = \operatorname{argmin}_x f(x)$$

- Jako počáteční odhad minima zkusíme

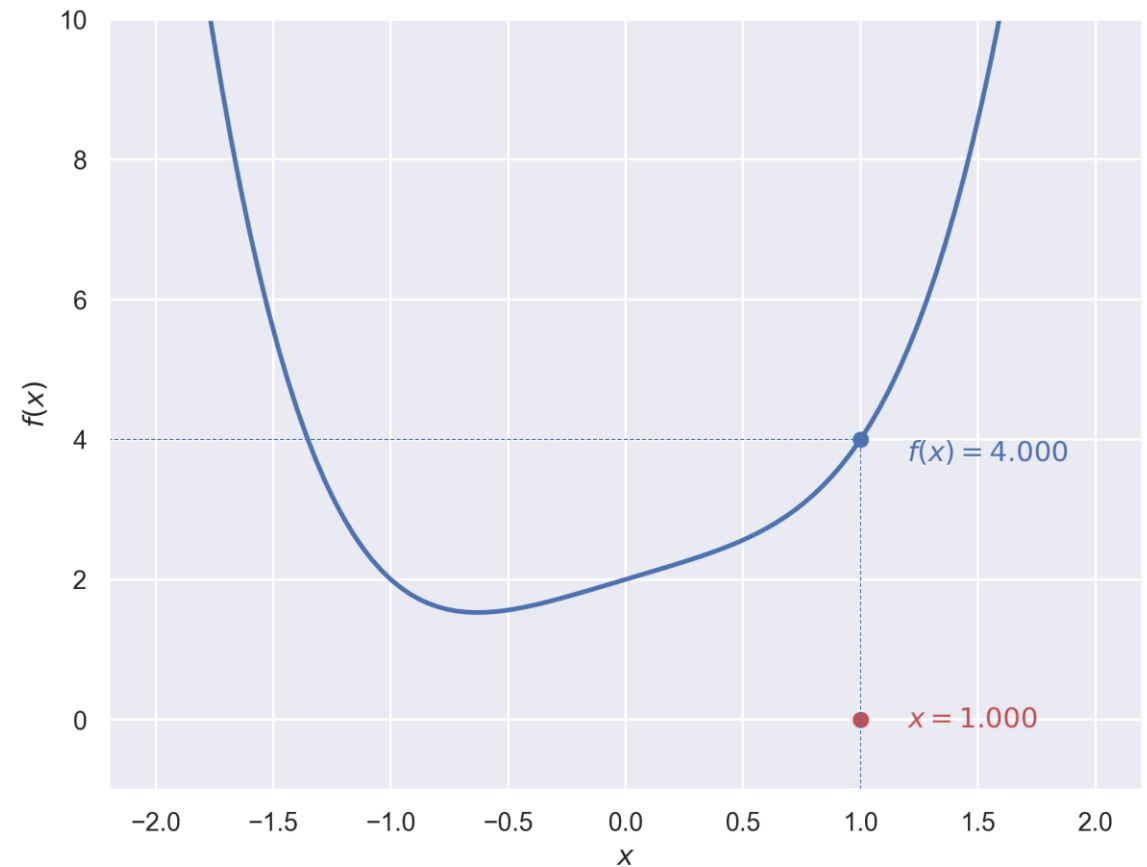
$$x^{(0)} = 1$$

- Nastavíme $\gamma = 0.2$ a $h = 0.001$



Příklad GD: funkce $f(x) = x^4 + x + 2$

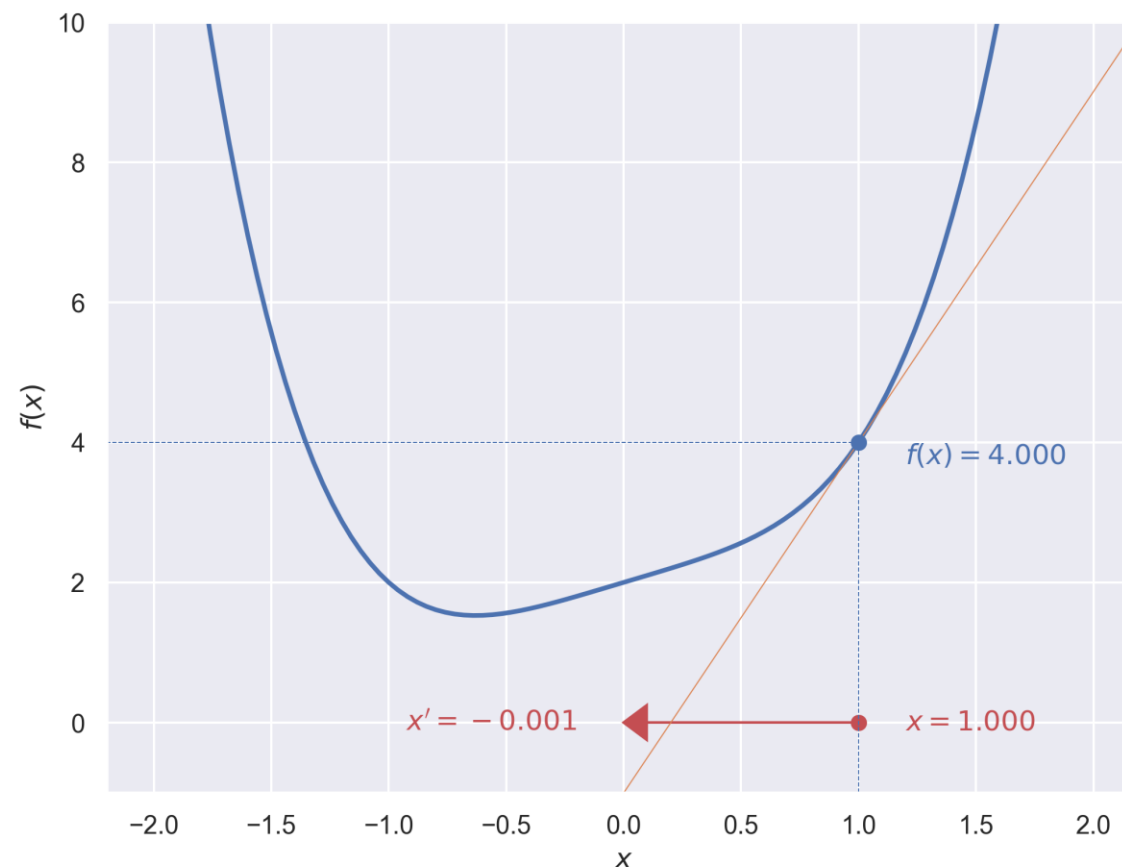
t	x	$f(x)$	$\frac{df}{dx}$



Příklad GD: funkce $f(x) = x^4 + x + 2$

t	x	$f(x)$	$\frac{df}{dx}$
0	1.000	4.000	5.006

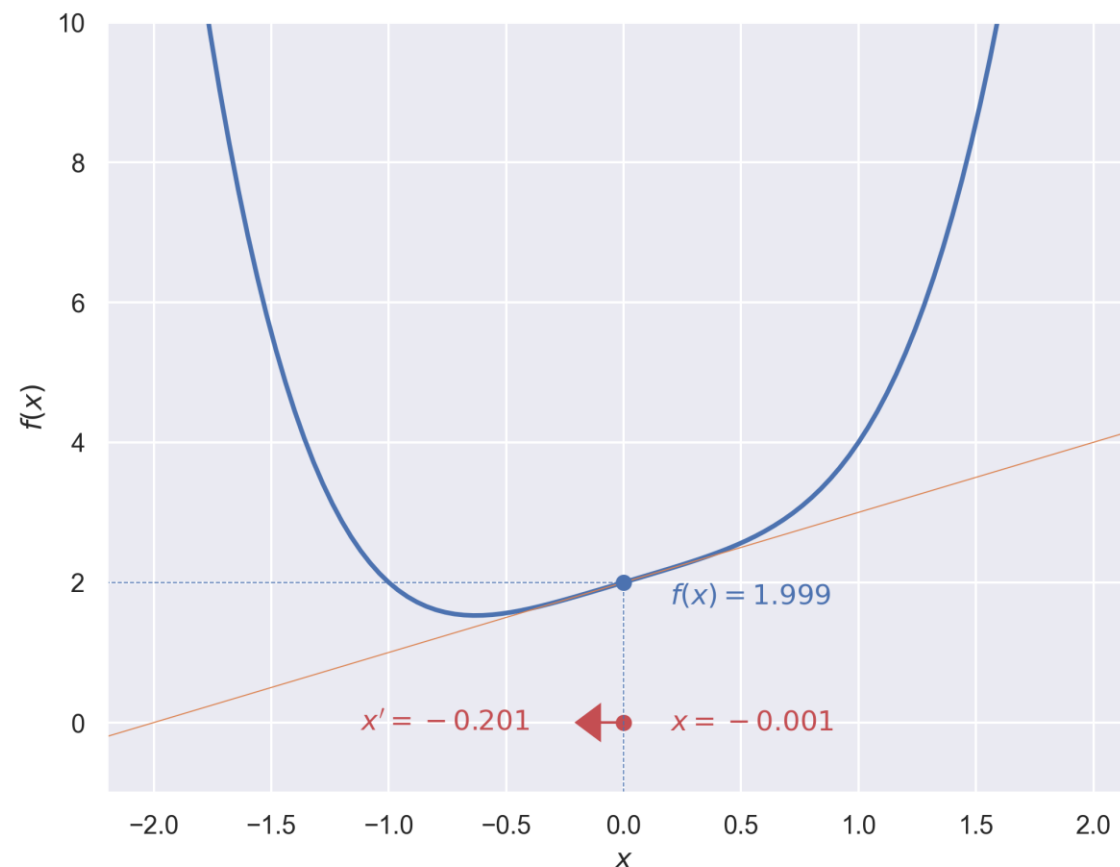
$$\begin{aligned}x &= 1.000 \\f(x) &= 1.000^4 + 1.000 + 2 = 4.000 \\ \frac{df(x)}{dx} &= \frac{f(1.000 + 0.001) - f(1.000)}{0.001} \\ &= \frac{4.005 - 4.000}{0.001} \\ &= 5.006 \\x' &= 1.000 - 0.2 \cdot 5.006 \\ &= -0.001\end{aligned}$$



Příklad GD: funkce $f(x) = x^4 + x + 2$

t	x	$f(x)$	$\frac{df}{dx}$
0	1.000	4.000	5.006
1	-0.001	1.999	1.000

$$\begin{aligned}x &= -0.001 \\f(x) &= (-0.001)^4 + (-0.001) + 2 = 1.999 \\ \frac{df(x)}{dx} &= \frac{f(-0.001 + 0.001) - f(-0.001)}{0.001} \\ &= \frac{2.000 - 1.999}{0.001} \\ &= 1.000 \\ x' &= -0.001 - 0.2 \cdot 1.000 \\ &= -0.201\end{aligned}$$



Příklad GD: funkce $f(x) = x^4 + x + 2$

t	x	$f(x)$	$\frac{df}{dx}$
0	1.000	4.000	5.006
1	-0.001	1.999	1.00
2	-0.201	1.800	0.968

$$x = -0.201$$

$$f(x) = (-0.201)^4 + (-0.201) + 2 = 1.800$$

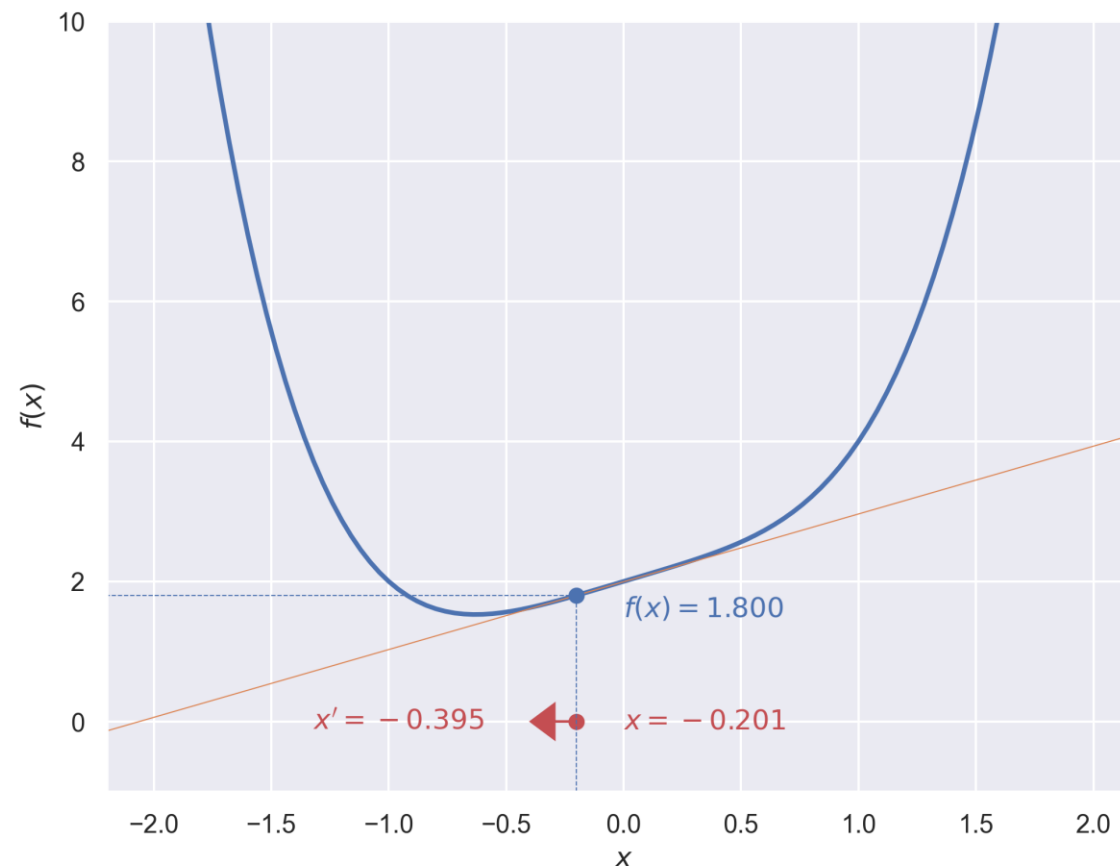
$$\frac{df(x)}{dx} = \frac{f(-0.201 + 0.001) - f(-0.201)}{0.001}$$

$$= \frac{1.801 - 1.800}{0.001}$$

$$= 0.968$$

$$x' = -0.201 - 0.2 \cdot 0.968$$

$$= -0.395$$



Příklad GD: funkce $f(x) = x^4 + x + 2$

t	x	$f(x)$	$\frac{df}{dx}$
0	1.000	4.000	5.006
1	-0.001	1.999	1.00
2	-0.201	1.800	0.968
3	-0.395	1.630	0.755

$$x = -0.395$$

$$f(x) = (-0.395)^4 + (-0.395) + 2 = 1.630$$

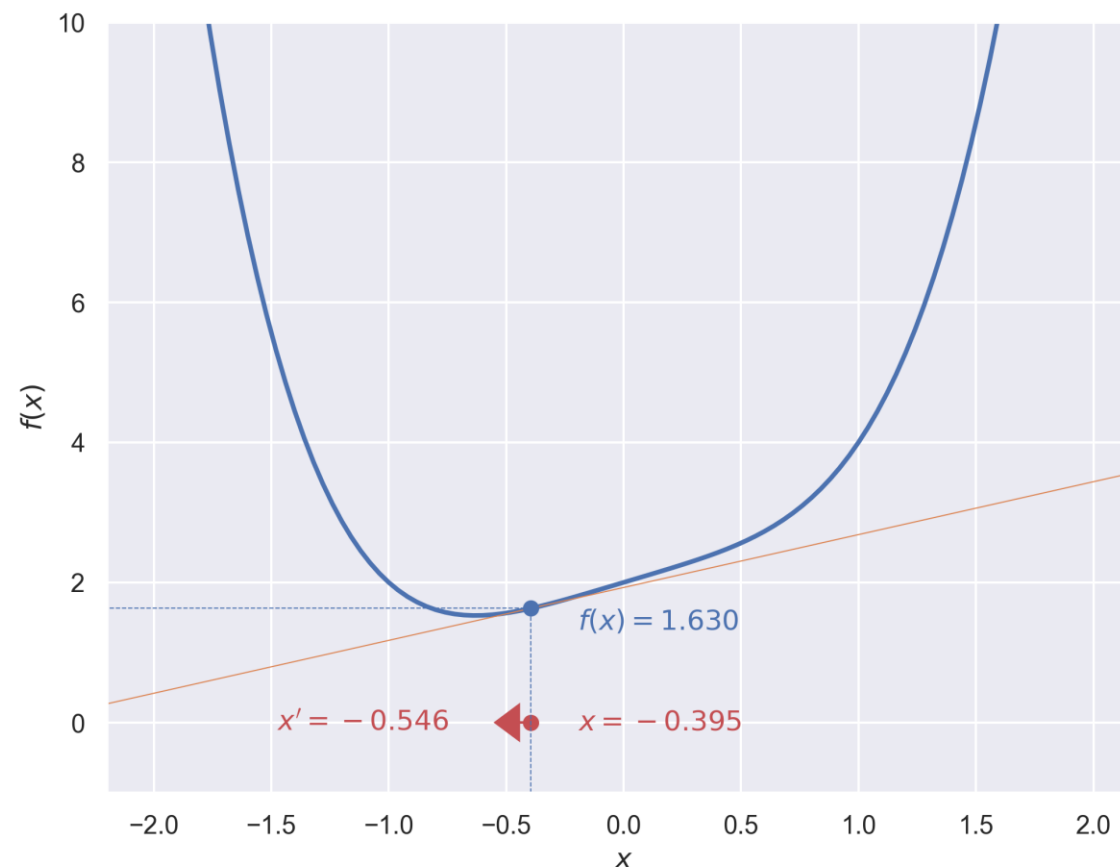
$$\frac{df(x)}{dx} = \frac{f(-0.395 + 0.001) - f(-0.395)}{0.001}$$

$$= \frac{1.630 - 1.630}{0.001}$$

$$= 0.755$$

$$x' = -0.395 - 0.2 \cdot 0.755$$

$$= -0.546$$



Příklad GD: funkce $f(x) = x^4 + x + 2$

t	x	$f(x)$	df/dx
0	1.000	4.000	5.006
1	-0.001	1.999	1.00
2	-0.201	1.800	0.968
3	-0.395	1.630	0.755
4	-0.546	1.543	0.352

$$x = -0.546$$

$$f(x) = (-0.546)^4 + (-0.546) + 2 = 1.543$$

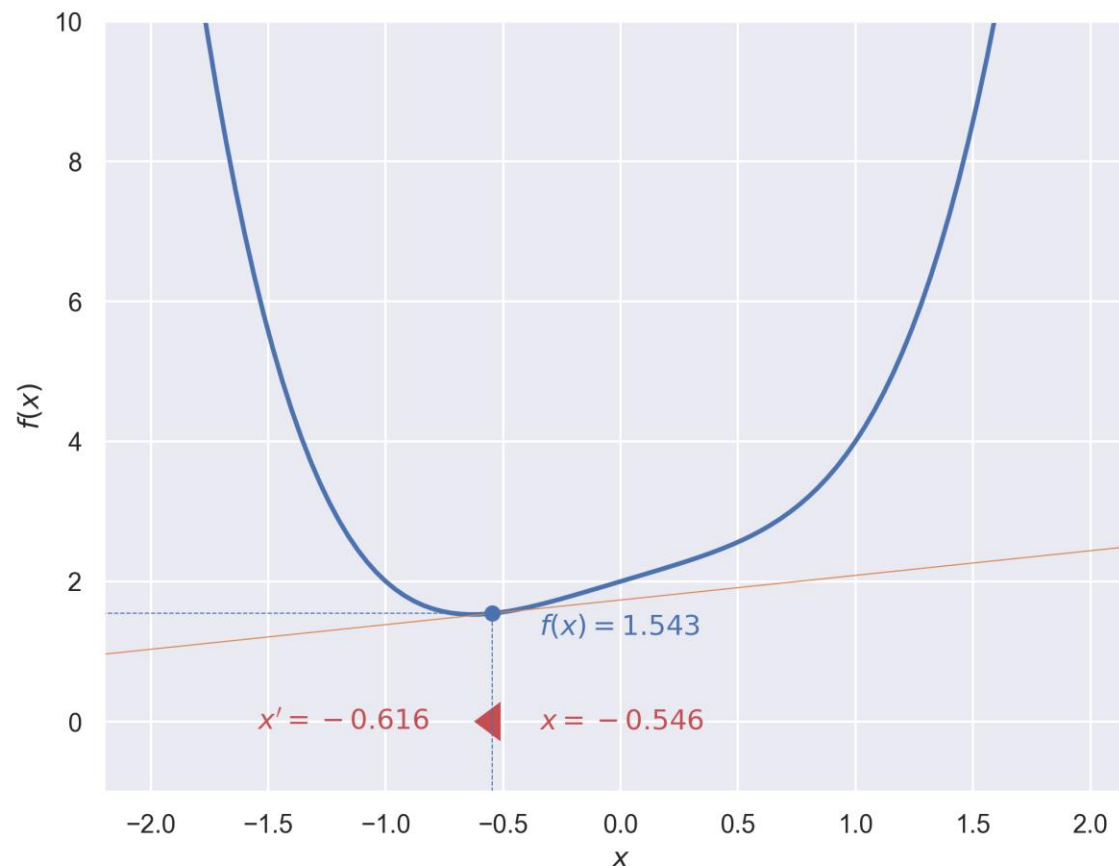
$$\frac{df(x)}{dx} = \frac{f(-0.546 + 0.001) - f(-0.546)}{0.001}$$

$$= \frac{1.543 - 1.543}{0.001}$$

$$= 0.352$$

$$x' = -0.546 - 0.2 \cdot 0.352$$

$$= -0.616$$



Příklad GD: funkce $f(x) = x^4 + x + 2$

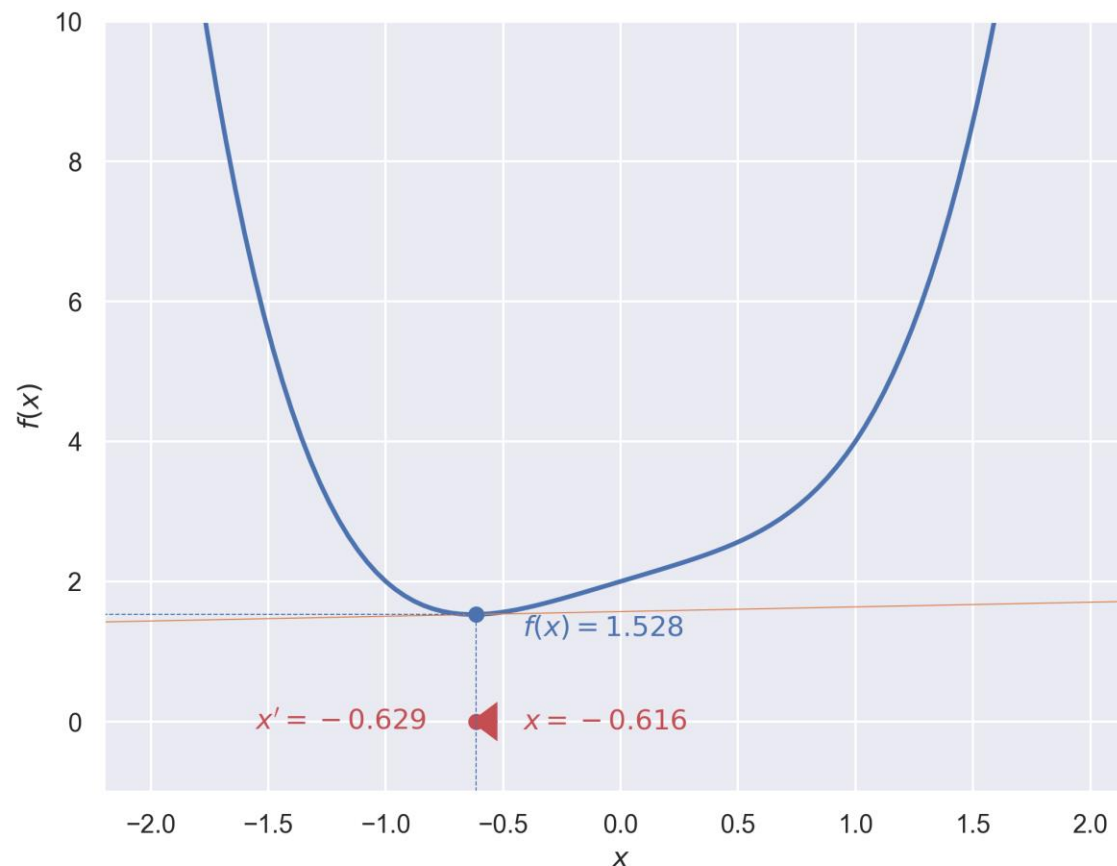
t	x	$f(x)$	df/dx
0	1.000	4.000	5.006
1	-0.001	1.999	1.00
2	-0.201	1.800	0.968
3	-0.395	1.630	0.755
4	-0.546	1.543	0.352
5	-0.616	1.528	0.067

$$x = -0.616$$

$$f(x) = (-0.616)^4 + (-0.616) + 2 = 1.528$$

$$\begin{aligned}\frac{df(x)}{dx} &= \frac{f(-0.616 + 0.001) - f(-0.616)}{0.001} \\ &= \frac{1.528 - 1.528}{0.001} \\ &= 0.067\end{aligned}$$

$$\begin{aligned}x' &= -0.616 - 0.2 \cdot 0.067 \\ &= -0.629\end{aligned}$$



Příklad GD: funkce $f(x) = x^4 + x + 2$

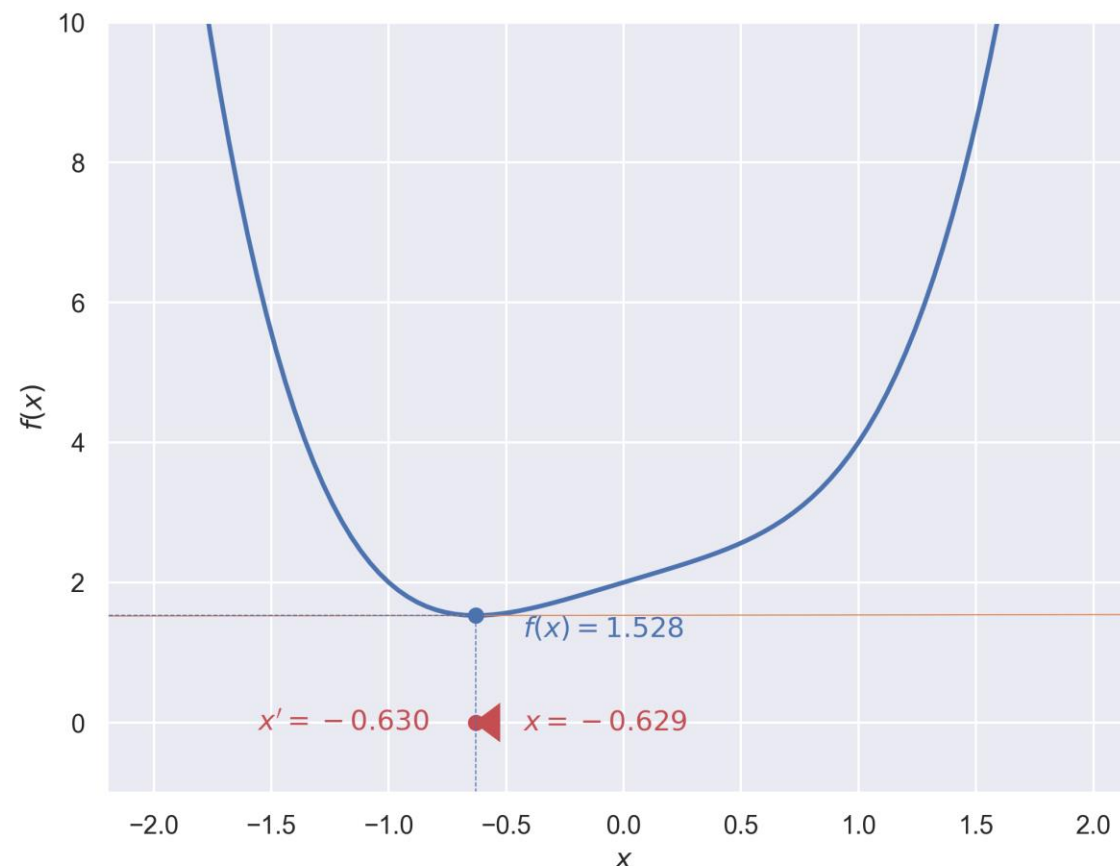
t	x	$f(x)$	df/dx
0	1.000	4.000	5.006
1	-0.001	1.999	1.00
2	-0.201	1.800	0.968
3	-0.395	1.630	0.755
4	-0.546	1.543	0.352
5	-0.616	1.528	0.067
6	-0.629	1.528	0.005

$$x = -0.629$$

$$f(x) = (-0.629)^4 + (-0.629) + 2 = 1.528$$

$$\begin{aligned}\frac{df(x)}{dx} &= \frac{f(-0.629 + 0.001) - f(-0.629)}{0.001} \\ &= \frac{1.528 - 1.528}{0.001} \\ &= 0.005\end{aligned}$$

$$\begin{aligned}x' &= -0.629 - 0.2 \cdot 0.005 \\ &= -0.630\end{aligned}$$



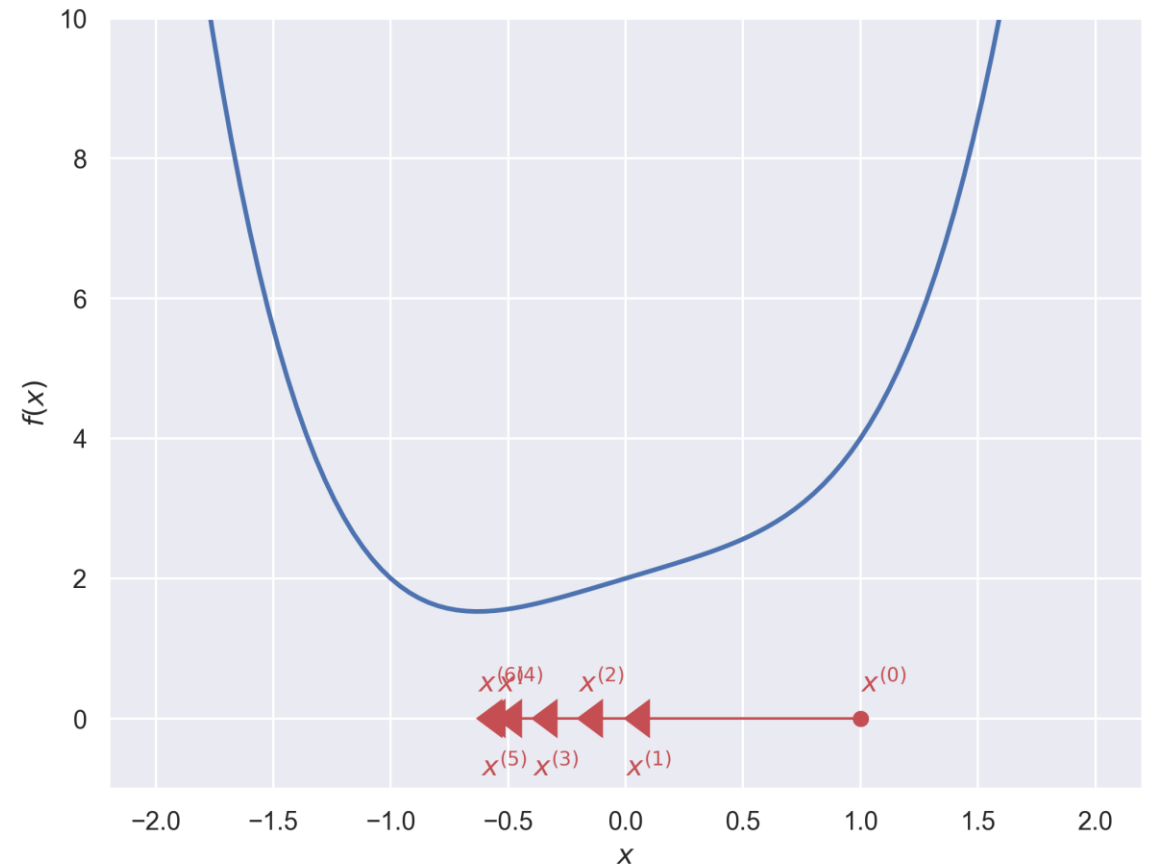
Optimální řešení dle
<https://www.wolframalpha.com/>

$$x^* = \frac{-1}{2^{2/3}} = 0.62996$$

Příklad GD: funkce $f(x) = x^4 + x + 2$

t	x	$f(x)$	df/dx
0	1.000	4.000	5.006
1	-0.001	1.999	1.00
2	-0.201	1.800	0.968
3	-0.395	1.630	0.755
4	-0.546	1.543	0.352
5	-0.616	1.528	0.067
6	-0.629	1.528	0.005

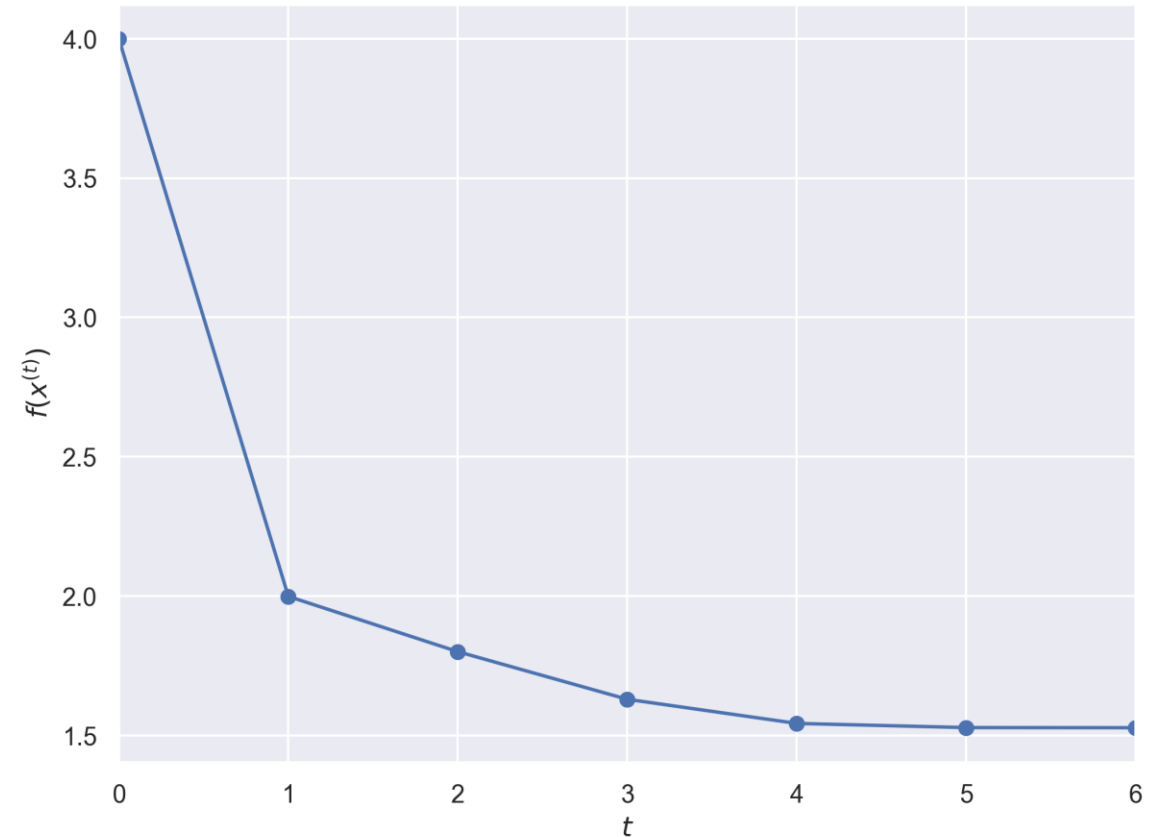
Vizualizace optimálního odhadu v čase



Příklad GD: funkce $f(x) = x^4 + x + 2$

t	x	$f(x)$	df/dx
0	1.000	4.000	5.006
1	-0.001	1.999	1.00
2	-0.201	1.800	0.968
3	-0.395	1.630	0.755
4	-0.546	1.543	0.352
5	-0.616	1.528	0.067
6	-0.629	1.528	0.005

Průběh lossu (optimalizované funkce) v čase



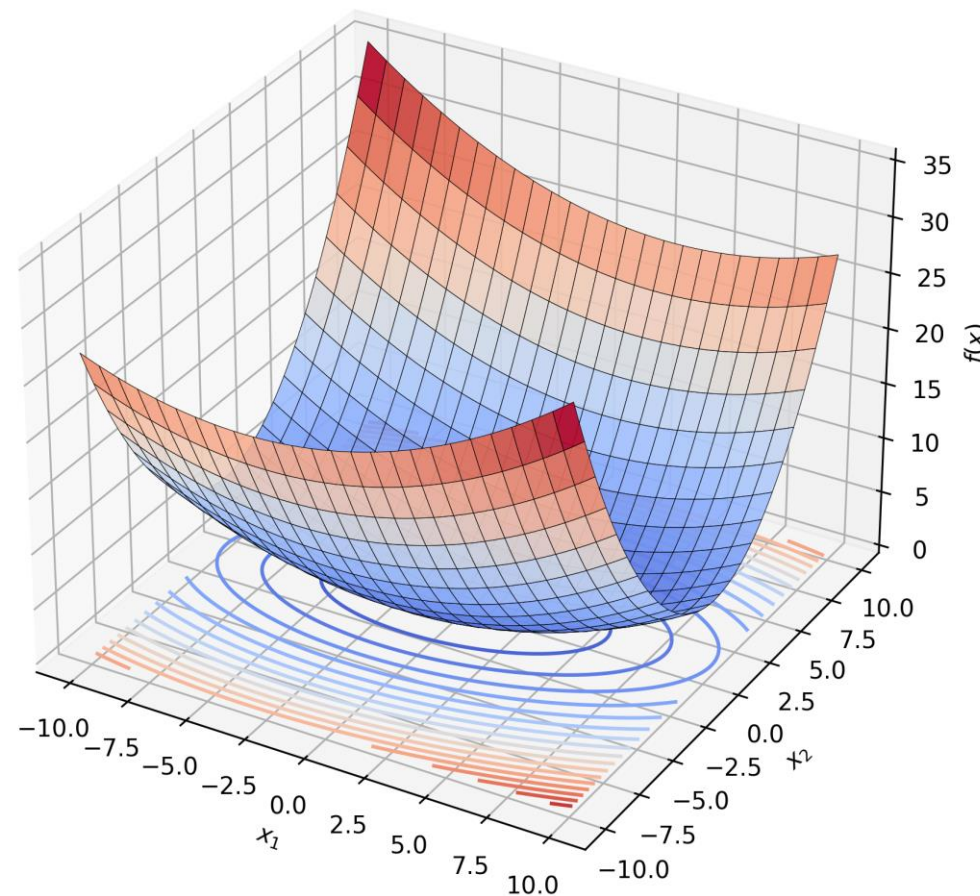
Funkce: n D vstup, 1D výstup

- Trénovaný model má obvykle více než jeden parametr
- Optimalizovaná loss funkce má tedy vstup \mathbf{x} jako vektor obecně s rozměrem D
- Vrací přitom skalární hodnotu (chybovost modelu na datasetu jako jediné číslo)
- Potřebujeme tedy minimalizovat

$$\mathbf{x}^* = \underset{\mathbf{x}}{\operatorname{argmin}} f(\mathbf{x})$$

kde

$$f(\mathbf{x}): \mathbb{R}^D \rightarrow \mathbb{R}$$
$$\mathbf{x} \in \mathbb{R}^D$$



Derivace funkce: n D vstup, 1D výstup

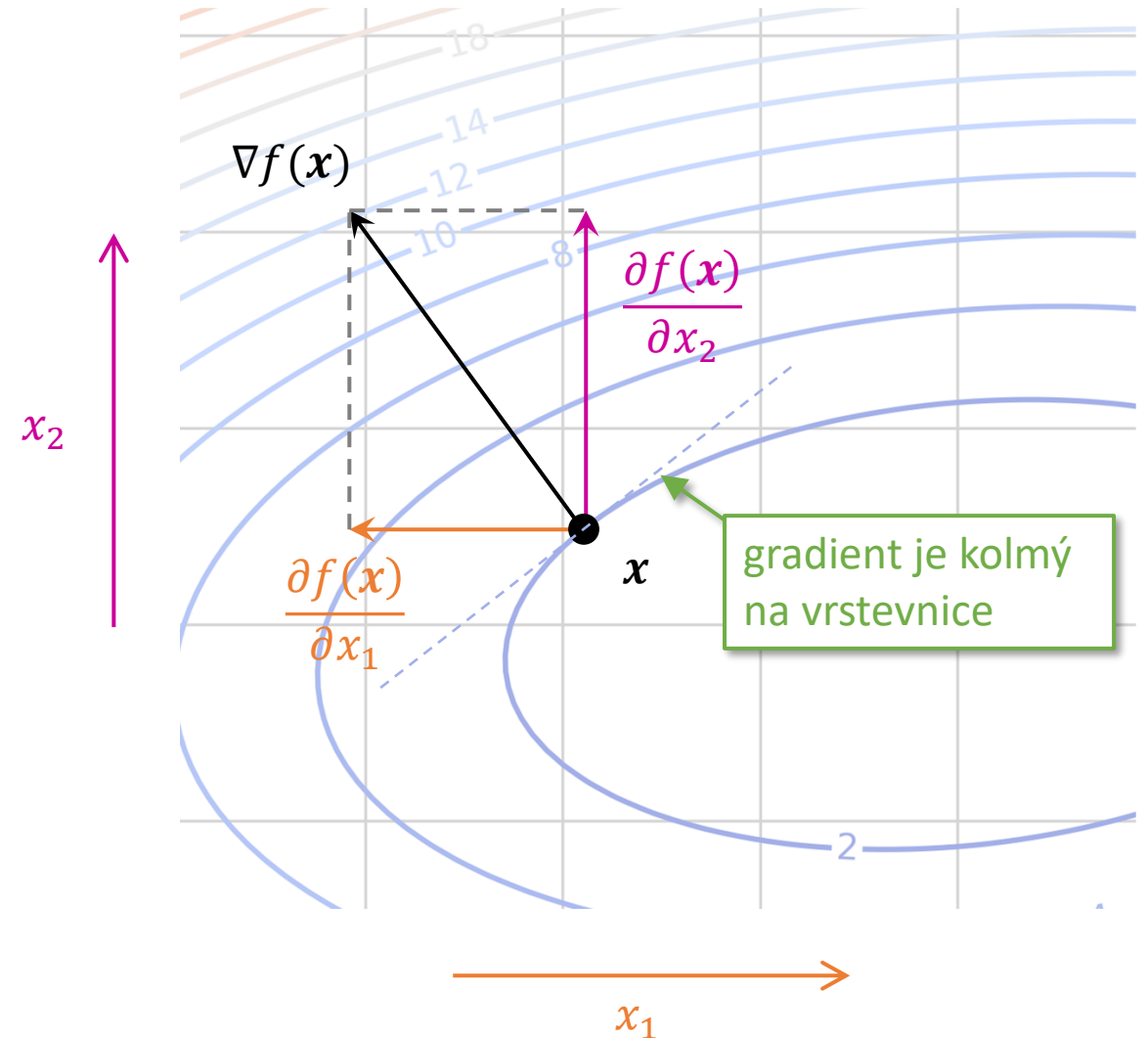
- U funkcí $f(\mathbf{x}): \mathbb{R}^D \rightarrow \mathbb{R}$, tj.

$$f(\mathbf{x}) = f(x_1, \dots, x_D)$$

tedy funkcí s D vstupy a 1 výstupem
můžeme derivovat vzhledem ke každému
z jednotlivých $x_d \rightarrow$ **parciální derivace**

- Uspořádání všech D parciálních derivací
do vektoru se nazývá **gradient**:

$$\nabla f(\mathbf{x}) = \frac{\partial f(\mathbf{x})}{\partial \mathbf{x}} = \left[\frac{\partial f(\mathbf{x})}{\partial x_1}, \dots, \frac{\partial f(\mathbf{x})}{\partial x_D} \right]^T$$



Metoda největšího spádu (Gradient Descent, GD)

- Metoda největšího spádu se pro vícerozměrné funkce principiálně nemění
- Jediný rozdíl je, že derivaci nahrazuje gradient a pravidlo je vektorové:

$$\mathbf{x}' := \mathbf{x} - \gamma \cdot \nabla f(\mathbf{x})$$

kde

původní odhad $\mathbf{x} \in \mathbb{R}^D$ je vektor s rozměrem D

nový odhad $\mathbf{x}' \in \mathbb{R}^D$ je vektor s rozměrem D

gradient $\nabla f(\mathbf{x}) \in \mathbb{R}^D$ je vektor s rozměrem D

krok učení (learning rate) $\gamma \in \mathbb{R}$ je skalár

Metoda největšího spádu (Gradient Descent, GD)

- Metoda největšího spádu se pro vícerozměrné funkce principiálně nemění
- Jediný rozdíl je, že derivaci nahrazuje gradient a pravidlo je vektorové:

$$\begin{bmatrix} x'_1 \\ \vdots \\ x'_D \end{bmatrix} := \begin{bmatrix} x_1 \\ \vdots \\ x_D \end{bmatrix} - \gamma \cdot \begin{bmatrix} \frac{\partial f(\mathbf{x})}{\partial x_1} \\ \vdots \\ \frac{\partial f(\mathbf{x})}{\partial x_D} \end{bmatrix}$$

kde

původní odhad $\mathbf{x} \in \mathbb{R}^D$ je vektor s rozměrem D

nový odhad $\mathbf{x}' \in \mathbb{R}^D$ je vektor s rozměrem D

gradient $\nabla f(\mathbf{x}) \in \mathbb{R}^D$ je vektor s rozměrem D

krok učení (learning rate) $\gamma \in \mathbb{R}$ je skalár

Příklad GD: lineární softmax cross entropy na CIFAR-10

- Optimalizovaná funkce má pro lineární klasifikátor se softmaxem a křížovou entropií formu

$$L(\mathbf{w}_{1,1}, \dots, \mathbf{w}_{K,D}, b_1, \dots, b_K) = -\frac{1}{N} \sum_{n=1}^N \log \frac{\exp(\mathbf{w}_{y_n, \cdot} \cdot \mathbf{x}_n + b_{y_n})}{\sum_{k=1}^K \exp(\mathbf{w}_{k, \cdot} \cdot \mathbf{x}_n + b_k)}$$

- Chceme nalézt bod $\boldsymbol{\theta}^* = [w_{1,1}^*, \dots, w_{K,D}^*, b_1^*, \dots, b_K^*]^\top$ takový, ve kterém $L(\boldsymbol{\theta})$ nabývá minimální hodnoty, tj.

$$\boldsymbol{\theta}^* = \underset{\boldsymbol{\theta}}{\operatorname{argmin}} L(\boldsymbol{\theta})$$

- Data tedy považujeme za konstantu
- Jako počáteční odhad minima zkusíme

$$\boldsymbol{\theta}^{(0)} \sim \mathcal{N}(0, 0.001)$$

inicializujeme na náhodné hodnoty z normálního (gaussovského) rozdělení s nulovým průměrem a std. odch. 0.001

Gradient lineárního softmaxu a křížové entropie analyticky

Optimalizovaná funkce (loss) je

$$L(\mathbf{w}, \mathbf{b}) = -\frac{1}{N} \sum_{n=1}^N \log \frac{\exp(\mathbf{w}_{y_n, :} \cdot \mathbf{x}_n + b_{y_n})}{\sum_{k=1}^K \exp(\mathbf{w}_{k, :} \cdot \mathbf{x}_n + b_k)} \hat{p}_{n, y_n}$$

Potřebujeme

$$\frac{\partial L(\mathbf{w}, \mathbf{b})}{\partial \mathbf{w}} = ?$$



$$\frac{\partial L(\mathbf{w}, \mathbf{b})}{\partial \mathbf{w}} = \frac{1}{N} \sum_{n=1}^N (\hat{\mathbf{p}}_n - \mathbf{p}_n) \cdot \mathbf{x}_n^T$$

$K \times D$ $(K \times 1) \cdot (1 \times D)$

$$\frac{\partial L(\mathbf{w}, \mathbf{b})}{\partial \mathbf{b}} = ?$$



$$\frac{\partial L(\mathbf{w}, \mathbf{b})}{\partial \mathbf{b}} = \frac{1}{N} \sum_{n=1}^N (\hat{\mathbf{p}}_n - \mathbf{p}_n)$$

$K \times 1$ $(K \times 1)$

Gradient descent (GD) pro multiclass logistickou regresi

Inicializujeme:

$$w_{kd} \sim \mathcal{N}(0, 0.001)$$

$$b_k \sim \mathcal{N}(0, 0.001)$$

Opakujeme:

$$\hat{p}_n := \frac{\exp(\mathbf{w}_{y_n} \cdot \mathbf{x}_n + b_{y_n})}{\sum_{k=1}^K \exp(\mathbf{w}_{k} \cdot \mathbf{x}_n + b_k)}$$

$$L(\mathbf{w}, \mathbf{b}) := -\frac{1}{N} \sum_{n=1}^N \log p_{n, y_n}$$

$$\frac{\partial L(\mathbf{w}, \mathbf{b})}{\partial \mathbf{w}} := \frac{1}{N} \sum_{n=1}^N (\hat{p}_n - p_n) \cdot \mathbf{x}_n^\top$$

$$\frac{\partial L(\mathbf{w}, \mathbf{b})}{\partial \mathbf{b}} := \frac{1}{N} \sum_{n=1}^N (\hat{p}_n - p_n)$$

$$\mathbf{w} := \mathbf{w} - \gamma \cdot \frac{\partial L(\mathbf{w}, \mathbf{b})}{\partial \mathbf{w}}$$

$$\mathbf{b} := \mathbf{b} - \gamma \cdot \frac{\partial L(\mathbf{w}, \mathbf{b})}{\partial \mathbf{b}}$$

```
1. w = 1e-3 * np.random.randn(10, 3072)
2. b = 1e-3 * np.random.randn(10)
3. for t in range(200):
4.     l = 0.
5.     dw = np.zeros_like(w) # (K, D)
6.     db = np.zeros_like(b) # (K,)
7.     ids = np.random.permutation(len(X))[:bs]
8.     for n, (xn, yn) in enumerate(zip(X[ids], Y[ids])):
9.         # loss
10.        sn = np.dot(w, xn) + b # (K,)
11.        pn = np.exp(sn) / np.sum(np.exp(sn)) # (K,)
12.        ln = -np.log(pn[yn])
13.
14.        # gradients
15.        dbn = pn.copy() # (K,)
16.        dbn[yn] -= 1
17.        dwn = np.dot(dbn.reshape(-1, 1), xn.reshape(1, -1)) # (K, D)
18.
19.        # accumulate
20.        l += ln
21.        dw += dwn
22.        db += dbn
23.
24.        # average
25.        l /= n + 1
26.        dw /= n + 1
27.        db /= n + 1
28.
29.        # update
30.        w -= lr * dw
31.        b -= lr * db
```

Gradient descent (GD) pro multiclass logistickou regresi

Inicializujeme:

- $\theta = \{W, b\}$ na náhodné hodnoty

Opakujeme:

1. Pro každý vzorek x_n v trénovací sadě x_1, \dots, x_N
 - a. predikujeme pravděpodobnosti \hat{p}_n
 - b. vypočteme dílčí kritérium l_n a akumulujeme k celkovému l
 - c. vypočteme dílčí gradient ∇L_n a akumulujeme k celkovému ∇L
2. updatujeme parametry θ akumulovaným gradientem ∇L s krokem γ

Zastavíme:

- po fixním počtu iterací
- parametry θ se ustálí
- hodnota kritéria $J(\theta)$ již delší dobu neklesá

```
1. w = 1e-3 * np.random.randn(10, 3072)
2. b = 1e-3 * np.random.randn(10)
3. for t in range(200):
4.     l = 0.
5.     dw = np.zeros_like(w) # (K, D)
6.     db = np.zeros_like(b) # (K,)
7.     ids = np.random.permutation(len(X))[:bs]
8.     for n, (xn, yn) in enumerate(zip(X[ids], Y[ids])):
9.         # loss
10.        sn = np.dot(w, xn) + b # (K,)
11.        pn = np.exp(sn) / np.sum(np.exp(sn)) # (K,)
12.        ln = -np.log(pn[yn])
13.
14.        # gradients
15.        dbn = pn.copy() # (K,)
16.        dbn[yn] -= 1
17.        dwn = np.dot(dbn.reshape(-1, 1), xn.reshape(1, -1)) # (K, D)
18.
19.        # accumulate
20.        l += ln
21.        dw += dwn
22.        db += dbn
23.
24.        # average
25.        l /= n + 1
26.        dw /= n + 1
27.        db /= n + 1
28.
29.        # update
30.        w -= lr * dw
31.        b -= lr * db
```

Gradient descent (GD) pro multiclass logistickou regresi

Inicializujeme:

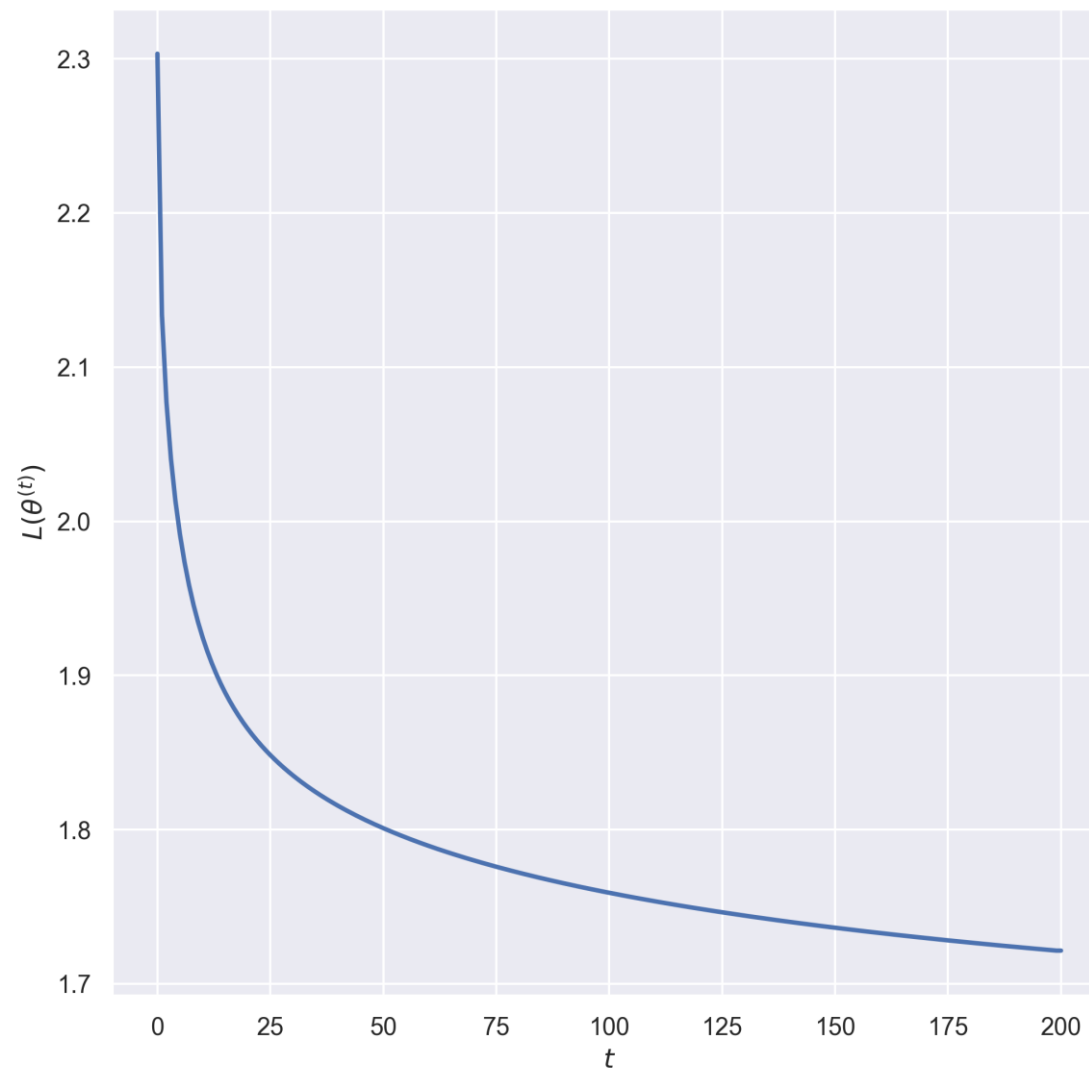
- $\theta = \{W, b\}$ na náhodné hodnoty

Opakujeme:

1. Pro každý vzorek x_n v trénovací sadě x_1, \dots, x_N
 - a. predikujeme pravděpodobnosti \hat{p}_n
 - b. vypočteme dílčí kritérium l_n a akumulujeme k celkovému l
 - c. vypočteme dílčí gradient ∇L_n a akumulujeme k celkovému ∇L
2. updatujeme parametry θ akumulovaným gradientem ∇L s krokem γ

Zastavíme:

- po fixním počtu iterací
- parametry θ se ustálí
- hodnota kritéria $J(\theta)$ již delší dobu neklesá



Gradient descent (GD) pro multiclass logistickou regresi

Inicializujeme:

- $\theta = \{W, b\}$ na náhodné hodnoty

← jak budeme inicializovat?

Opakujeme:

1. Pro každý vzorek x_n v trénovací sadě x_1, \dots, x_N
 - a. predikujeme pravděpodobnosti \hat{p}_n
 - b. vypočteme dílčí kritérium l_n a akumulujeme k celkovému l
 - c. vypočteme dílčí gradient ∇L_n a akumulujeme k celkovému ∇L
2. updatujeme parametry θ akumulovaným gradientem ∇L s krokem γ

hyperparametry

← jaký krok učení?

Zastavíme:

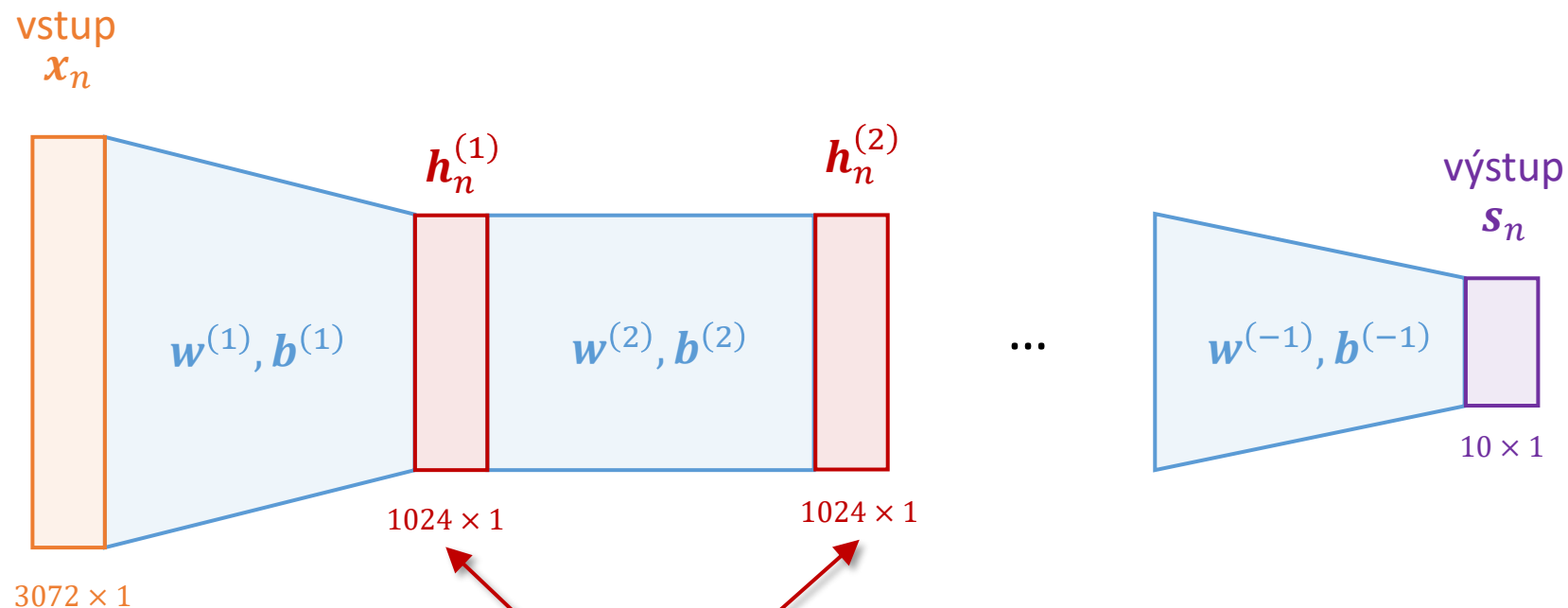
- po fixním počtu iterací
- parametry θ se ustálí
- hodnota kritéria $J(\theta)$ již delší dobu neklesá

← kolik iterací? jak dlouho?

Vícevrstvý perceptron

Vícevrstvý perceptron (Multi-Layer Perceptron, MLP)

- Označuje se také jako dopředná síť (Feed-Forward Network, FFN)
- Opakují se především dva+ typy vrstev: plně propojená a aktivace



Bloky sestávající z plně propojené vrstvy a aktivace se obvykle označují jako vrstvy a jejich výstup jako tzv. skrytá vrstva (h jako hidden)

Manuální výpočet gradientu?

Optimalizovaná funkce (loss) je

$$L(\theta) = -\frac{1}{N} \sum_{n=1}^N \log \frac{\exp(f(\mathbf{x}_n, \theta)_{y_n})}{\sum_{k=1}^K \exp(f(\mathbf{x}_n, \theta)_k)}$$

Co když $f(\mathbf{x}_n, \theta)$ potažmo $L(\theta)$ jsou složité funkce jako např. hluboké neuronové sítě?

Potřebujeme

$$\frac{\partial L(\theta)}{\partial \theta} = ?$$



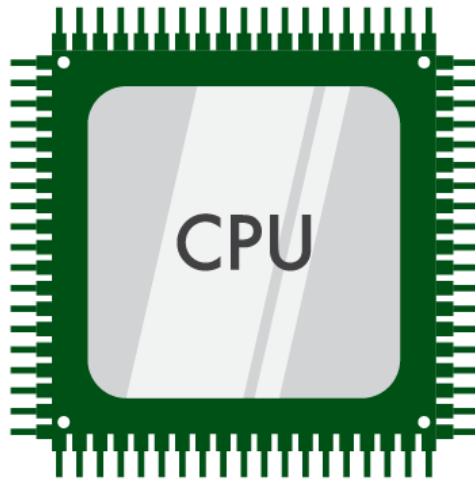
Automatický výpočet gradientu = zpětná propagace

Optimalizovaná funkce (loss) je

$$L(\theta) = -\frac{1}{N} \sum_{n=1}^N \log \frac{\exp(f(\mathbf{x}_n, \theta)_{y_n})}{\sum_{k=1}^K \exp(f(\mathbf{x}_n, \theta)_k)}$$

Potřebujeme

$$\frac{\partial L(\theta)}{\partial \theta} = ?$$



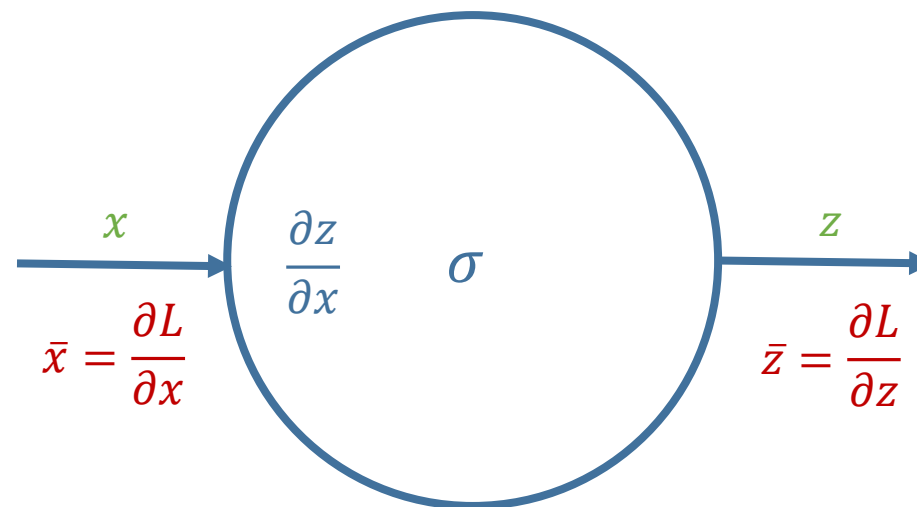
$$\frac{\partial L(\theta)}{\partial \theta} = \dots$$

Funkce jako uzel ve výpočetním grafu

```
class Sigmoid(Function):
```

```
    @staticmethod
    def forward(x: float): # zatím pouze skalary
        z = 1 / (1 - math.exp(-x))
        cache = z,
        return z, cache
```

```
    @staticmethod
    def backward(dz: float, cache: tuple):
        z, = cache
        dx = dz * z * (1 - z) # retizkove pravidlo
        return dx
```



zpětný průchod je řetízkové pravidlo

$$\frac{\partial L}{\partial x} = \frac{\partial L}{\partial z} \cdot \frac{\partial z}{\partial x}$$

Každá funkce, kterou chceme použít jako stavební blok, musí mít definovaný

dopředný průchod

zpětný průchod

Dvouvrstvý perceptron v numpy na 11 řádků

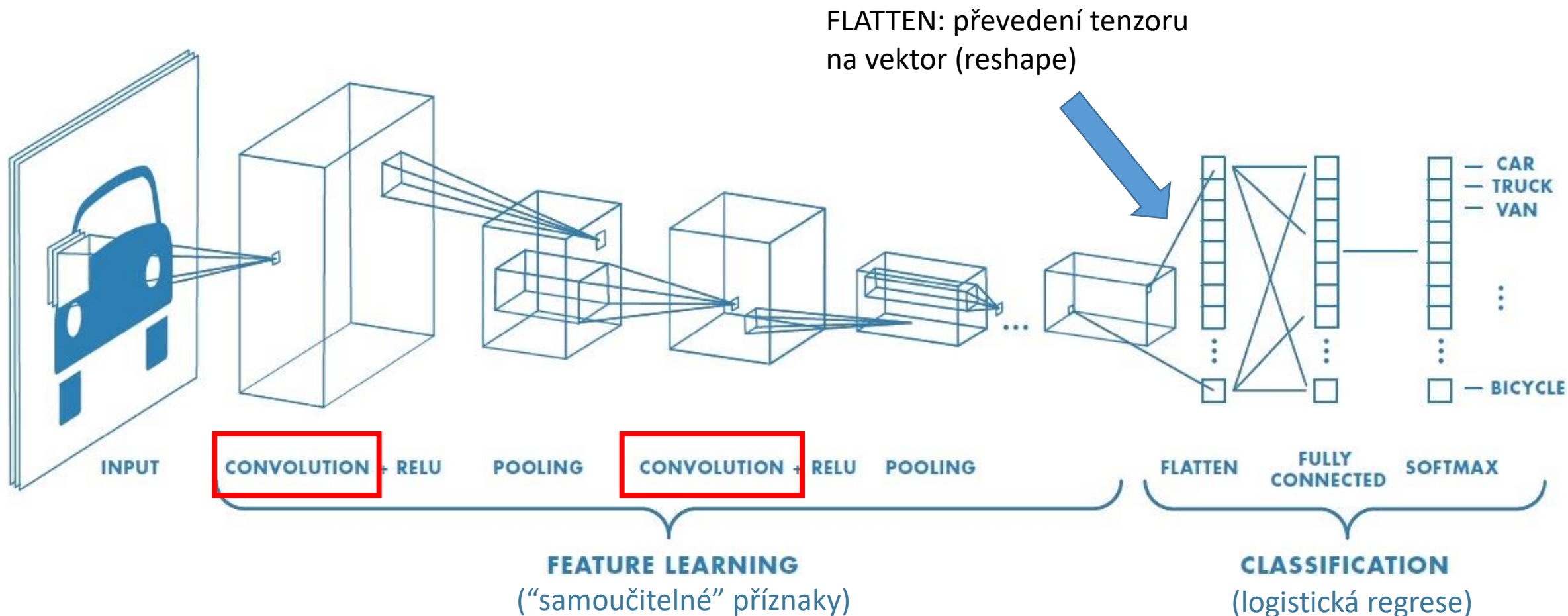
```
X = np.array([ [0,0,1],[0,1,1],[1,0,1],[1,1,1] ])
y = np.array([[0,1,1,0]]).T
syn0 = 2*np.random.random((3,4)) - 1
syn1 = 2*np.random.random((4,1)) - 1
for j in range(60000):
    l1 = 1/(1+np.exp(-(np.dot(X,syn0))))
    l2 = 1/(1+np.exp(-(np.dot(l1,syn1))))
    l2_delta = (y - l2)*(l2*(1-l2))
    l1_delta = l2_delta.dot(syn1.T) * (l1 * (1-l1))
    syn1 += l1.T.dot(l2_delta)
    syn0 += X.T.dot(l1_delta)
```

<http://iamtrask.github.io/2015/07/12/basic-python-network/>

Konvoluční sítě

Konvoluční síť (Convolutional Neural Network, CNN)

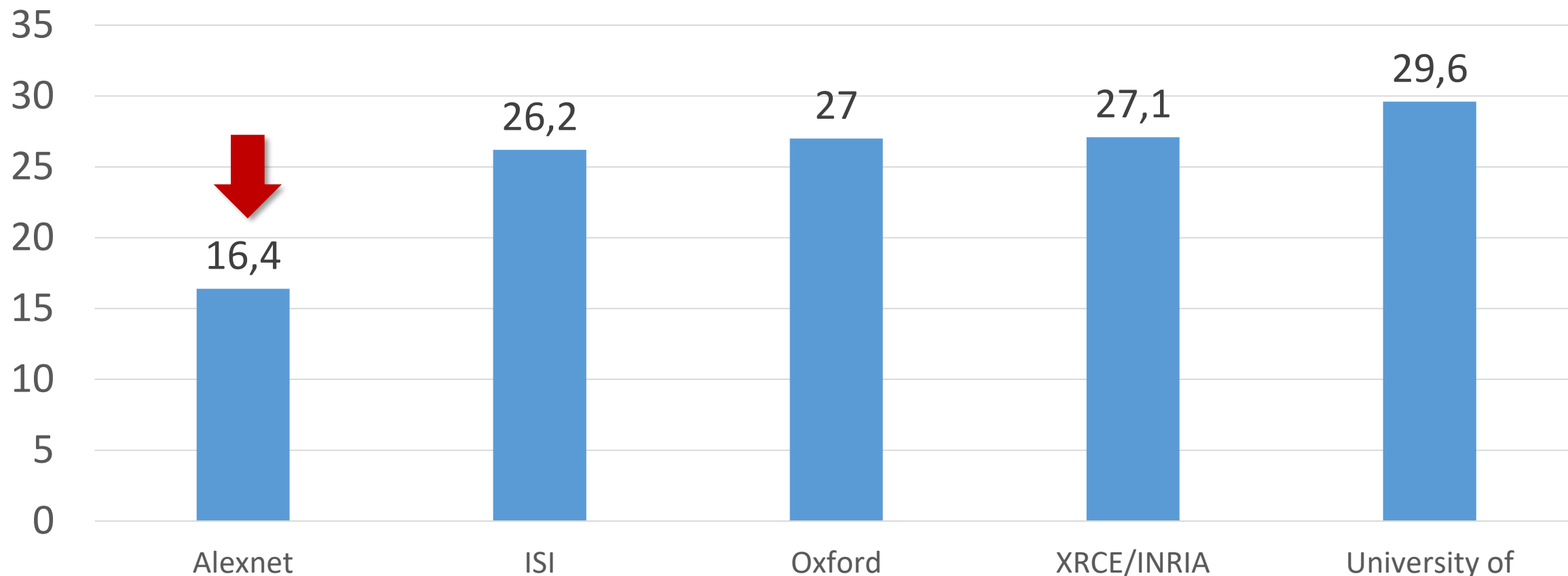
- zadefinováním konvoluce jako bloku v neurosíti nyní můžeme libovolně kombinovat s ostatními vrstvami



obrázek: <https://ch.mathworks.com/solutions/deep-learning/convolutional-neural-network.html>

Rozpoznávání ImageNet: Alexnet (2012)

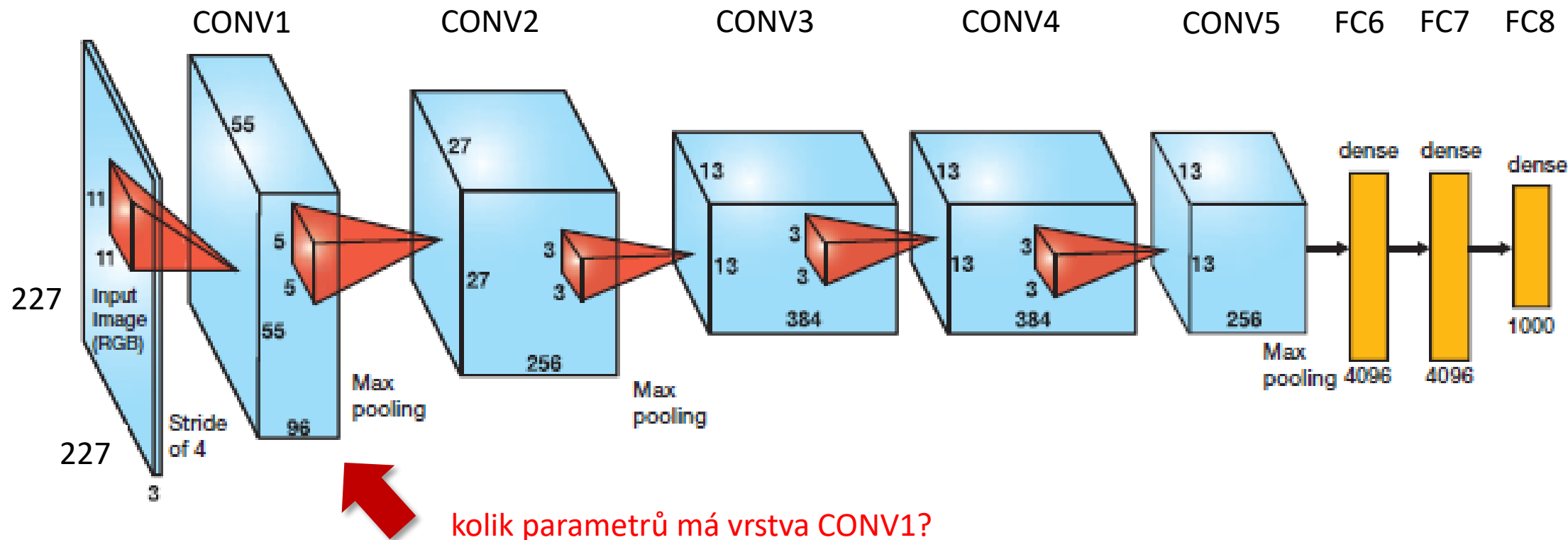
1000 classes, Top-5 error [%]



top-5 ... správná třída musí být
do 5. místa dle skóre z výstupu
klasifikátoru

non-DNN modely

Alexnet (2012)



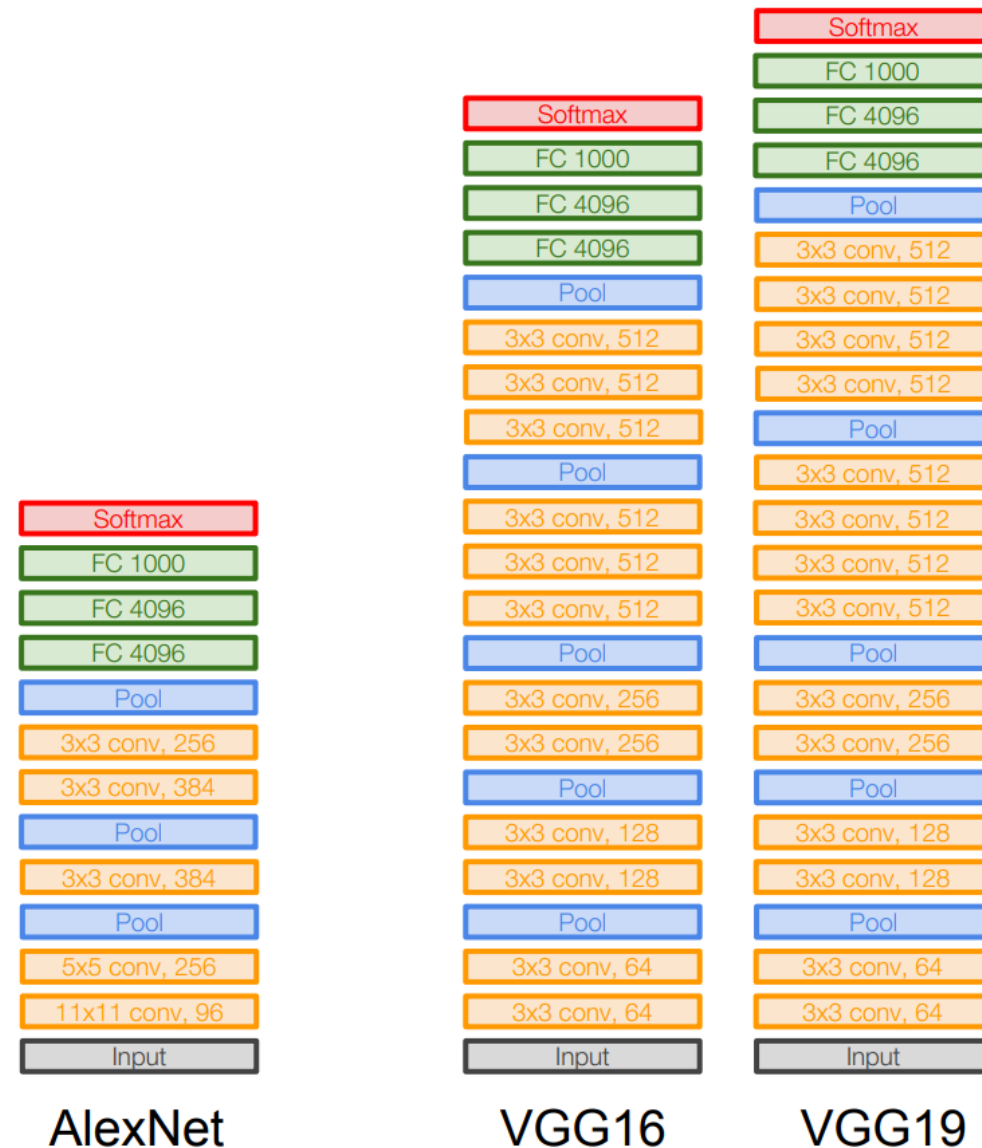
kolik parametrů má vrstva CONV1?

$$11 \cdot 11 \cdot 3 \cdot 96 = 34\,848$$

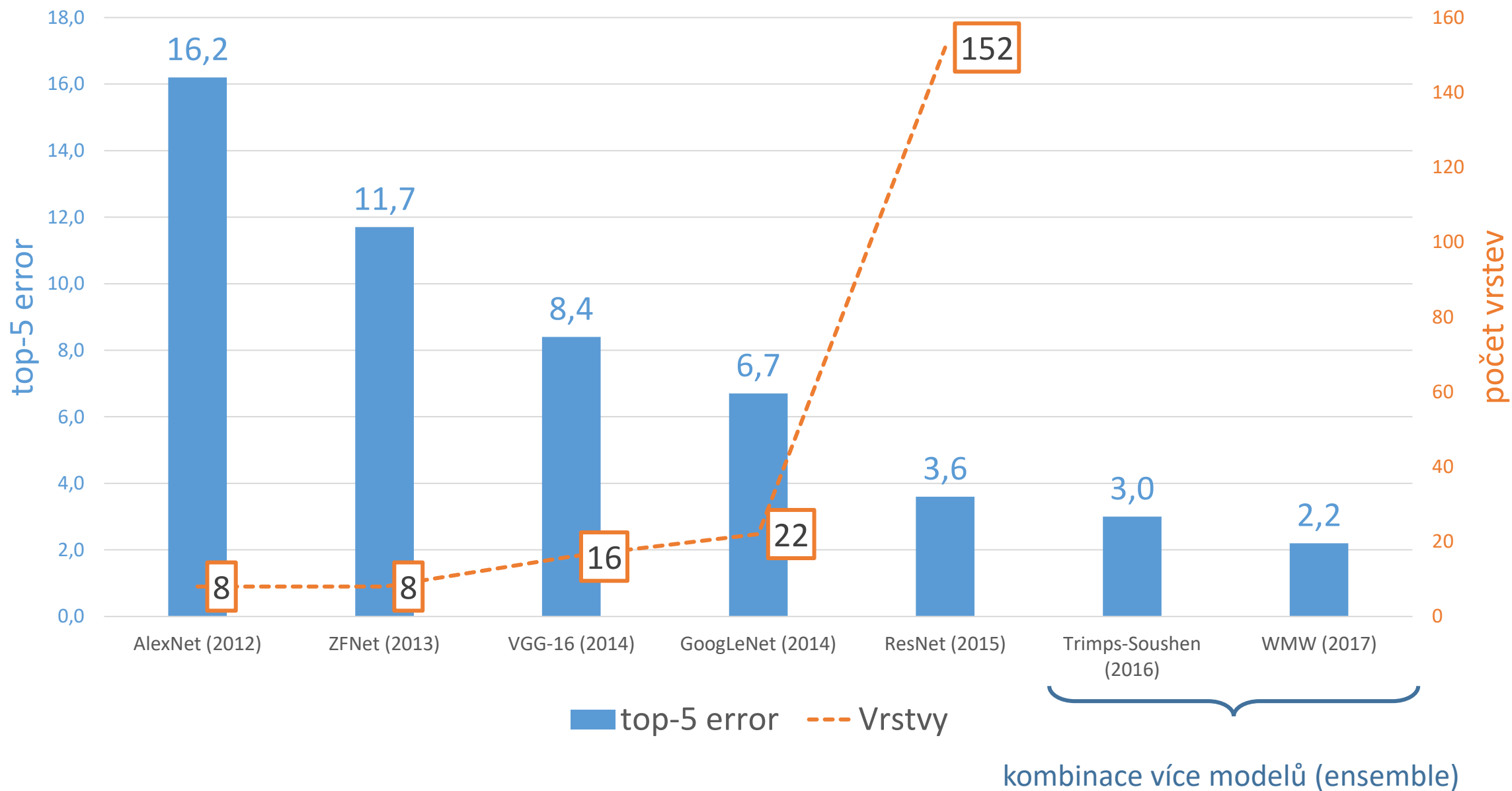
- architektura: CONV-POOL-NORM-CONV-POOL-NORM-CONV-CONV-CONV-FC-FC-FC
- “naškálovaná” LeNet-5

VGG (2014)

- [Simonyan, Zisserman: “Very Deep Convolutional Networks for Large-Scale Image Recognition”](#)
- Druhé místo ImageNet competition 2014
- Mnohem jednodušší architektura než vítěz (GoogLeNet)
- Velmi podobné AlexNet
- Místo 11x11 apod. konvolucí pouze 3x3
- Pouze 2x2 max-pooling
- Žádná lokální normalizace
- 16 a 19 vrstev
- VGG-16: 8.4 % top-5 error



ImageNet klasifikace



ResNet (2015)

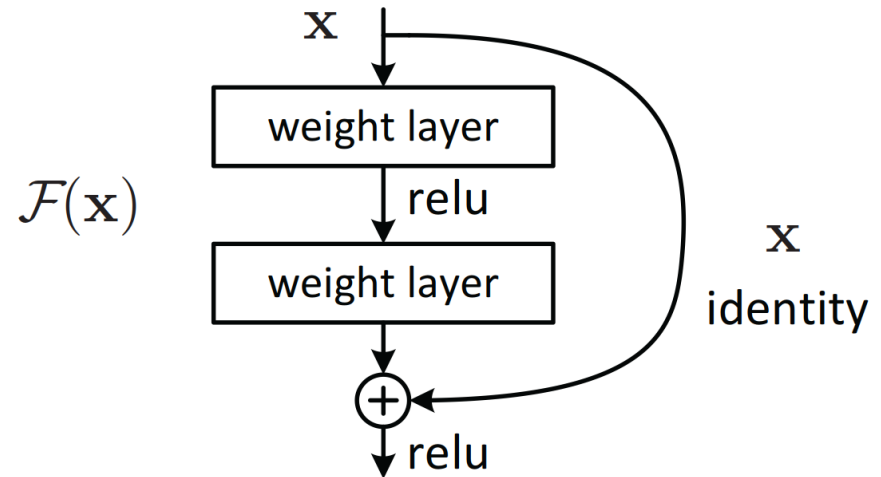
- [He et al.: “Deep Residual Learning for Image Recognition”](#)
- Cílem návrhu být co nejhlubší → 152 vrstev!
- Vítěz ImageNet 2015 ve všech kategoriích
- Vítěz MS COCO challenge
- 3.6 % top-5 error na ImageNet: lepší než člověk (cca 5 %)

Reziduální blok

- Podobně jako inception používá složitější bloky
- Výstup sestává ze součtu konvoluce a přímo mapovaného vstupu (identity)
- Síť se tedy učí pouze rezidua

$$\mathcal{F}(x) = \mathcal{H}(x) - x$$

- “Naučit se nuly je jednodušší než identitu”



$$\mathcal{H}(x) = \mathcal{F}(x) + x$$

Figure 2. Residual learning: a building block.

EfficientNet (2019)

- [Tan, Le: EfficientNet: Rethinking Model Scaling for Convolutional Neural Networks](#)

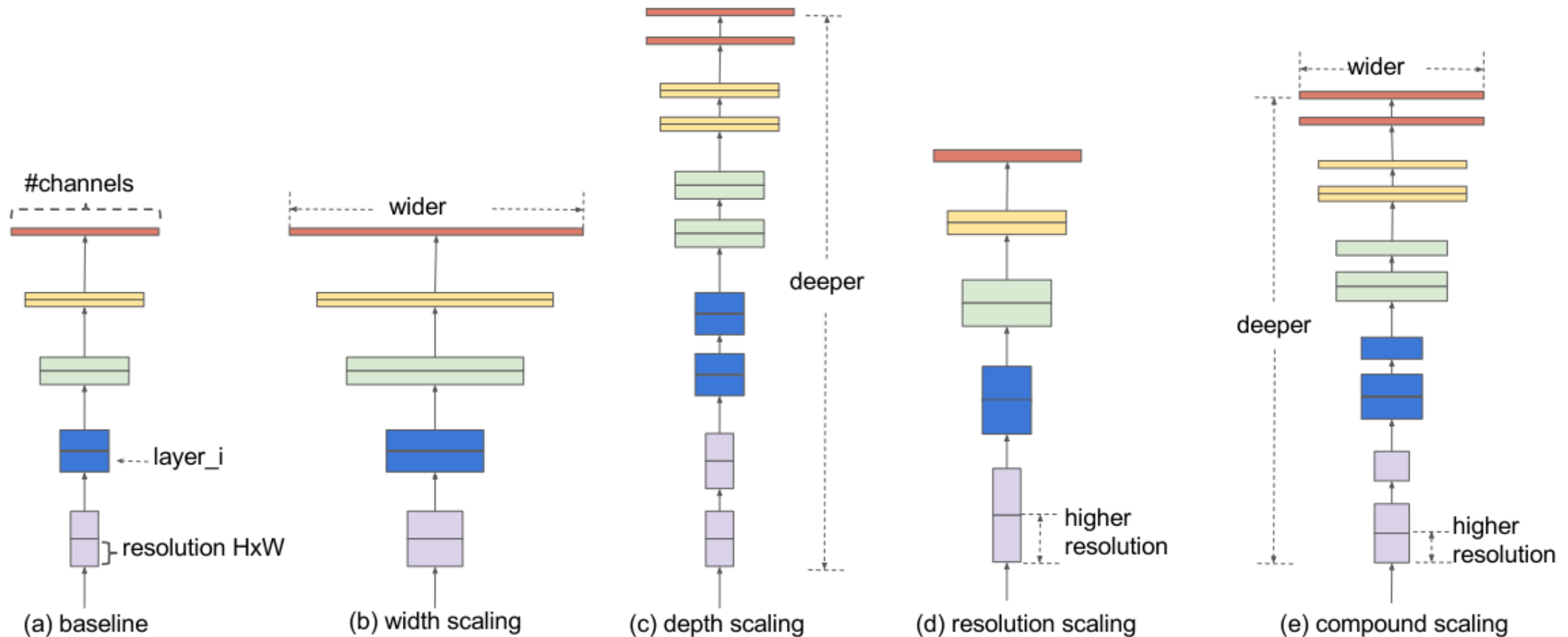
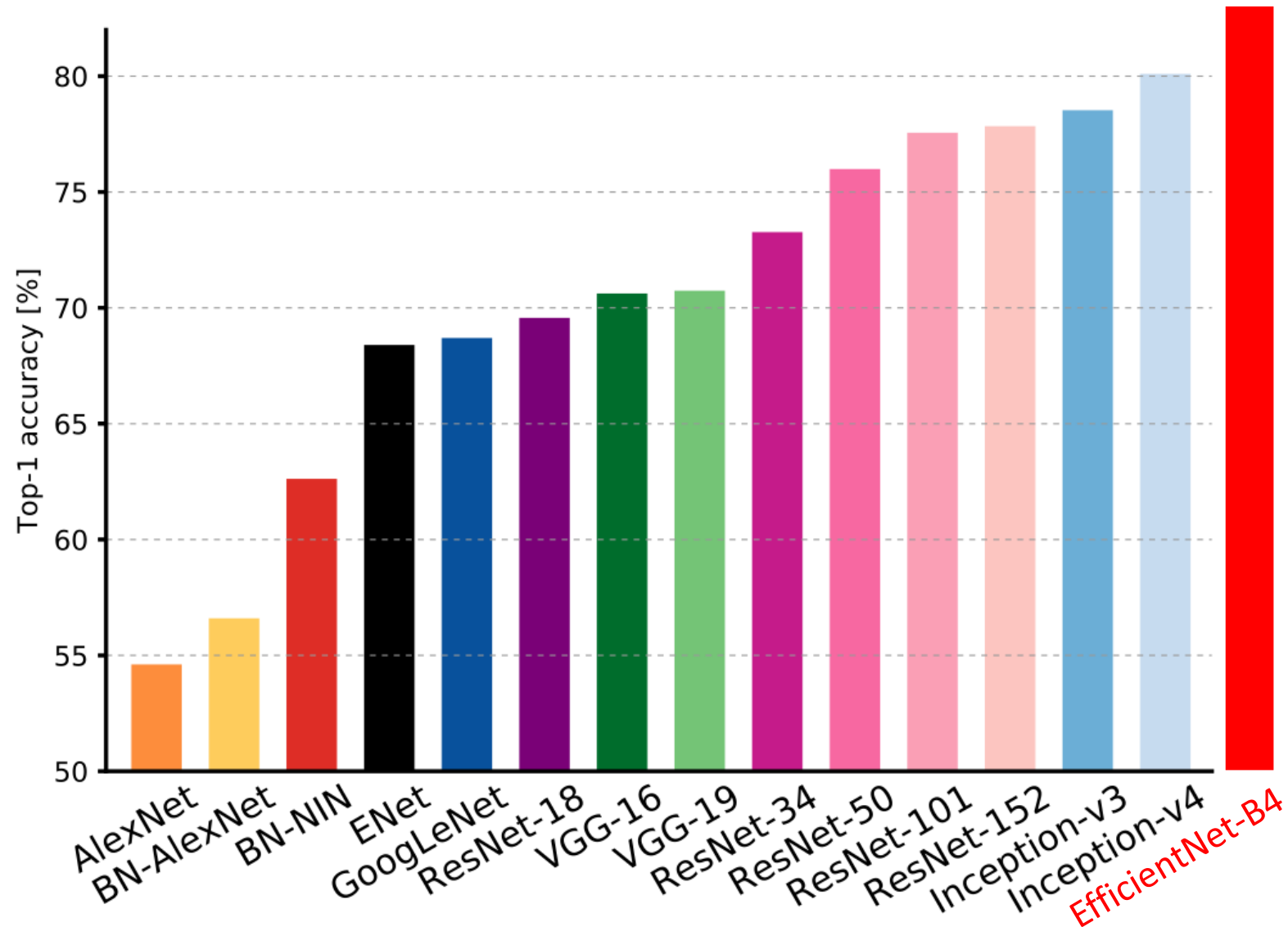


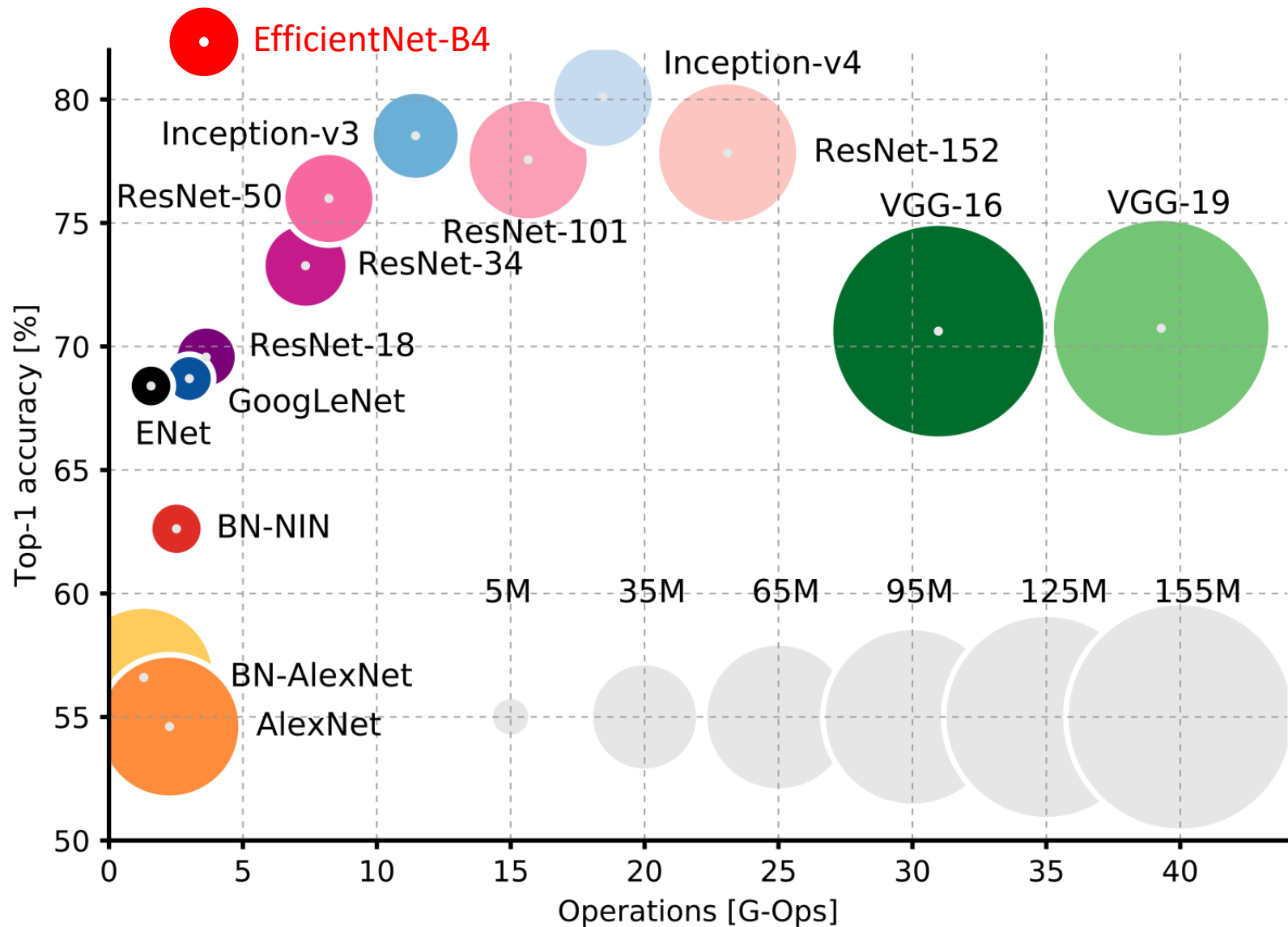
Figure 2. Model Scaling. (a) is a baseline network example; (b)-(d) are conventional scaling that only increases one dimension of network width, depth, or resolution. (e) is our proposed compound scaling method that uniformly scales all three dimensions with a fixed ratio.

Srovnání nejpoužívanějších CNN architektur



obrázek: [Canziani et al.: “An Analysis of Deep Neural Network Models for Practical Applications”](#)

Srovnání nejpoužívanějších CNN architektur



velikost znázorňuje
celkový počet
parametrů

Trénování sítí

Trénování sítí

- Typicky několik kroků
 1. Příprava dat
 2. Inicializace modelu (sítě)
 3. Definice metrik a optimizéru
 4. Trénovací smyčka

Trénování sítí

- Typicky několik kroků

1. Příprava dat
2. Inicializace modelu (sítě)
3. Definice metrik a optimizéru
4. Trénovací smyčka

```
training_set =  
torchvision.datasets.FashionMNIST(  
    './data',  
    train = True,  
    transform = transform,  
    download = True  
)
```

```
validation_set =  
torchvision.datasets.FashionMNIST(  
    './data',  
    train = False,  
    transform = transform,  
    download = True  
)
```

Trénování sítí

- Typicky několik kroků

1. Příprava dat
2. Inicializace modelu (sítě)
3. Definice metrik a optimizéru
4. Trénovací smyčka

```
training_loader =  
torch.utils.data.DataLoader(  
    training_set,  
    batch_size = 4,  
    shuffle = True  
)
```

```
validation_loader =  
torch.utils.data.DataLoader(  
    validation_set,  
    batch_size = 4,  
    shuffle = False  
)
```

Trénování sítí

- Typicky několik kroků

1. Příprava dat

2. Inicializace modelu (sítě)

3. Definice metrik a optimizéru

4. Trénovací smyčka

```
class GarmentClassifier(nn.Module):  
    def __init__(self):  
        super(GarmentClassifier, self).__init__()  
        self.conv1 = nn.Conv2d(1, 6, 5)  
        self.pool = nn.MaxPool2d(2, 2)  
        self.conv2 = nn.Conv2d(6, 16, 5)  
        self.fc1 = nn.Linear(16 * 4 * 4, 120)  
        self.fc2 = nn.Linear(120, 84)  
        self.fc3 = nn.Linear(84, 10)
```

```
    def forward(self, x):  
        x = self.pool(F.relu(self.conv1(x)))  
        x = self.pool(F.relu(self.conv2(x)))  
        x = x.view(-1, 16 * 4 * 4)  
        x = F.relu(self.fc1(x))  
        x = F.relu(self.fc2(x))  
        x = self.fc3(x)  
        return x
```

```
model = GarmentClassifier()
```

Trénování sítí

- Typicky několik kroků

1. Příprava dat
2. Inicializace modelu (sítě)
3. Definice metrik a optimizéru
4. Trénovací smyčka

```
loss_fn = torch.nn.CrossEntropyLoss()
```

```
optimizer = torch.optim.SGD(  
    model.parameters(),  
    lr = 0.001,  
    momentum = 0.9  
)
```


Trénování sítí

- Typicky několik kroků

1. Příprava dat
2. Inicializace modelu (sítě)
3. Definice metrik a optimizéru
4. Trénovací smyčka

```
def train_one_epoch():
    running_loss = 0.
    last_loss = 0.

    for i, data in enumerate(training_loader):
        # forward
        inputs, labels = data
        outputs = model(inputs)
        loss = loss_fn(outputs, labels)

        # backward
        optimizer.zero_grad()
        loss.backward()
        optimizer.step()

        # Gather data and report
        running_loss += loss.item()
        if i % 1000 == 999:
            last_loss = running_loss / 1000 # loss per batch
            ...
    return last_loss
```

Trénování sítí

- Typicky několik kroků

1. Příprava dat
2. Inicializace modelu (sítě)
3. Definice metrik a optimizéru

4. Trénovací smyčka

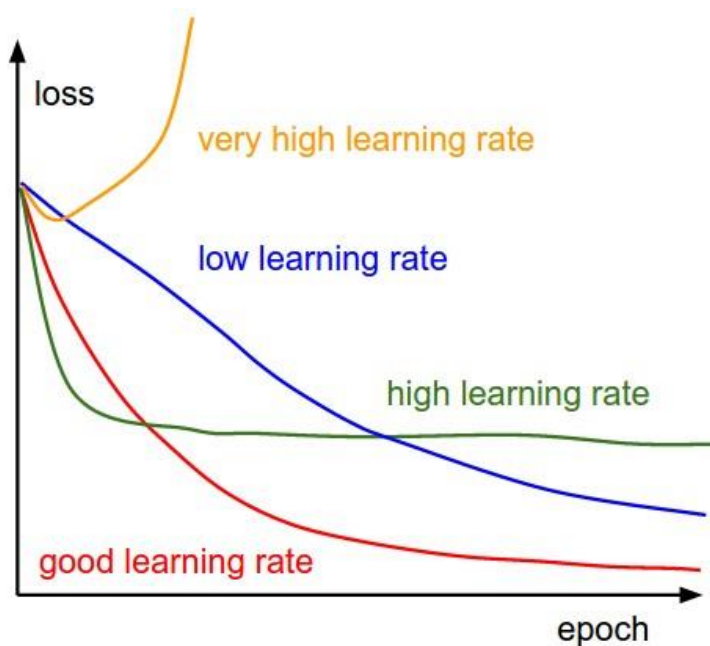
```
for epoch in range(EPOCHS):
    # training
    model.train(True)
    avg_loss = train_one_epoch()

    # validation
    running_vloss = 0.0
    model.eval()
    with torch.no_grad():
        for i, vdata in enumerate(validation_loader):
            vinputs, vlabels = vdata
            voutputs = model(vinputs)
            vloss = loss_fn(voutputs, vlabels)
            running_vloss += vloss

    avg_vloss = running_vloss / (i + 1)
    print('LOSS train {} valid {}'.format(avg_loss, avg_vloss))

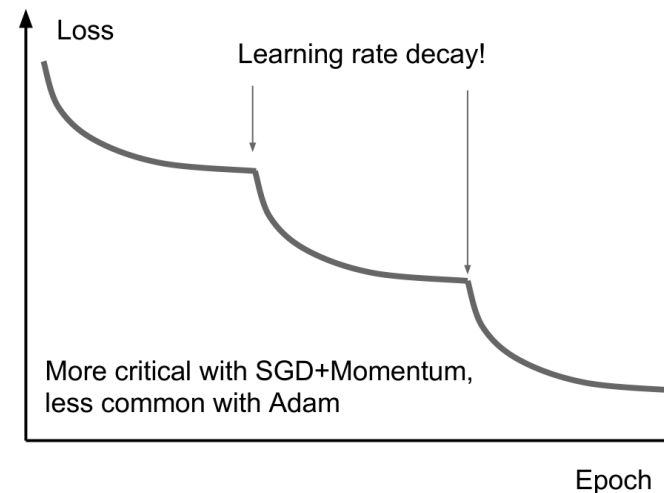
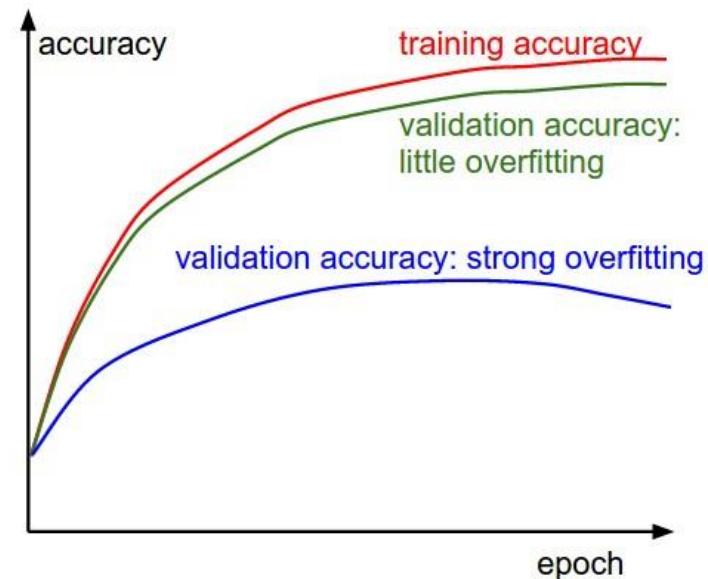
    # Track best performance, and save the model's state
    if avg_vloss < best_vloss:
        best_vloss = avg_vloss
        torch.save(model.state_dict(), model_path)
```

Trénování



- Monitorovat hodnotu lossu a podle toho nastavit lr
- Nebo lze použít automatické hledání lr
- Pokud funguje, zkusit lr decay

obrázky: <https://cs231n.github.io/neural-networks-3/>



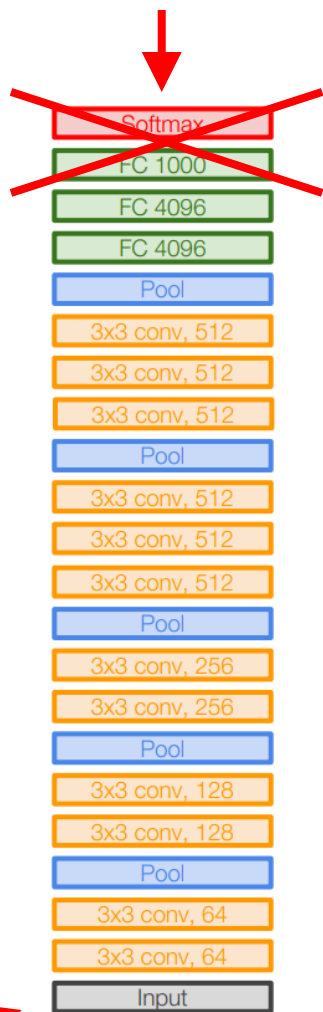
Transfer learning

Trénování konvolučních sítí při málo datech

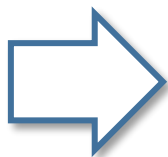
- Popsané architektury mají obvykle miliony parametrů
- Malé datasety na jejich trénování nestačí → výrazný overfit
- I pokud data máme: trénování VGG na ImageNet trvalo autorům 2-3 týdny, a to i s 4x NVIDIA Titan Black GPU
- **Naštěstí lze obejít!**
 1. Můžeme vzít existující již natrénovaný model (např. VGG-16)
 2. Odstraníme poslední klasifikační vrstvu
 3. Nahradíme vlastní

Transfer learning

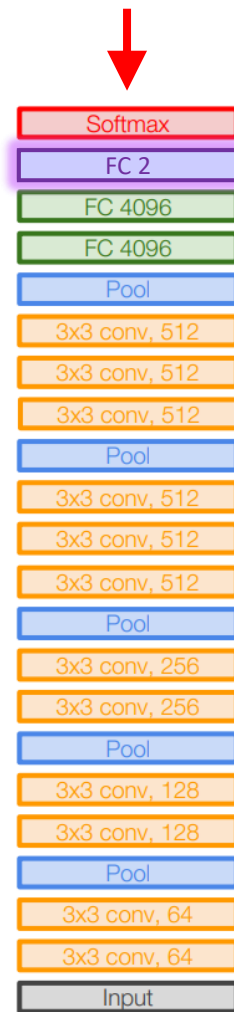
1000 tříd pro ImageNet



poslední lineární vrstvu
tvaru 4096×1000
zahodíme



2 třídy: kočky vs psi



připojíme vlastní lineární vrstvu
tvaru 4096×2 a náhodně
inicializujeme

pokud máme málo dat, je lepší
konvoluční vrstvy netrénovat →
layer freeze

můžeme použít vlastní data

ImageNet data

VGG16

VGG16

Transfer learning v PyTorch

```
model_conv = torchvision.models.resnet18(pretrained=True)
```

```
for param in model_conv.parameters():  
    param.requires_grad = False
```

“zmrazení” vrstev, nebudou se trénovat a zůstávají konst. → síť pouze jako extractor příznaků

```
# Parameters of newly constructed modules have requires_grad=True by default
```

```
num_ftrs = model_conv.fc.in_features  
model_conv.fc = nn.Linear(num_ftrs, 2)
```

poslední vrstvu klasifikující do 1000 ImageNet tříd nahradíme vlastní, která má pouze 2 třídy

```
model_conv = model_conv.to(device)
```

```
criterion = nn.CrossEntropyLoss()
```

jako seznam parametrů pro optimalizaci předáváme pouze poslední lineární vrstvu (pouze pro urychlení)

```
# Observe that only parameters of final layer are being optimized as  
# opposed to before.
```

```
optimizer_conv = optim.SGD(model_conv.fc.parameters(), lr=0.001, momentum=0.9)
```

- poté, co je poslední vrstva natrénovaná, je možné opět uvolnit (“rozmrazit”) i konvoluční vrstvy a model dále zlepšit (fine tuning)

https://pytorch.org/tutorials/beginner/transfer_learning_tutorial.html

Předzpracování dat u konvolučních sítí

- U konvolučních sítí často používané konstantní hodnoty

- Odečtení průměrného pixelu

`out = rgb - mean_pixel`

kde `mean_pixel` je trojice `[r, g, b]`

často lze vidět konkrétní hodnoty `[123.68, 116.779, 103.939]`, které jsou průměrným pixelem na databázi ImageNet

- Méně časté: odečtení průměrného obrázku

`out = rgb - mean_image`

kde `mean_image` je `32x32x3`

- Např. všechny předtrénované modely Pytorch:

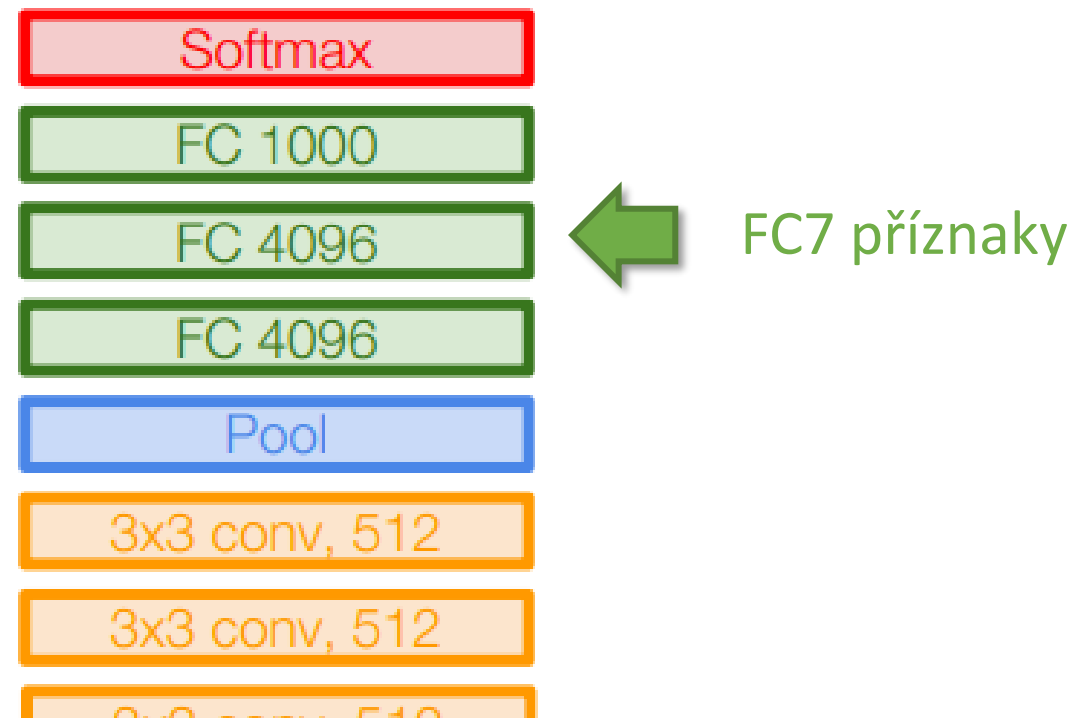
`transforms.Normalize([0.485, 0.456, 0.406], [0.229, 0.224, 0.225])`

odečtení průměru a normalizace standardní odchylky

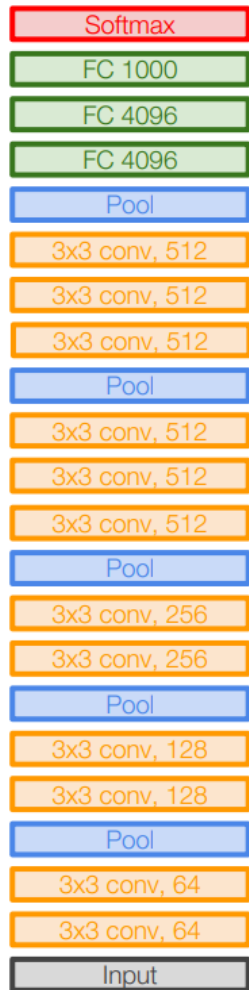
CNN příznaky

- Výstup z posledních lineárních vrstev lze použít např. jako příznaky (tzv. FC7) → CNN jako “feature extractor”
- Např. VGG-16 předposlední vrstva má rozměr 4096
- Nad těmito příznaky je možné natrénovat libovolný klasifikátor, třeba i rozhodovací stromy/lesy, bayesovské klasifikátory, ...
- Lze také využít pro urychlení trénování: celý dataset projet sítí a pro každý obrázek uložit na disk FC7 příznaky
- Během trénování se pak nemusí znovu a znovu provádět dopředný průchod celou sítí, pouze těmi posledními

např. VGG-16:



Transfer learning: shrnutí



specifičtější příznaky
(obličeje, text, ...)

obecné příznaky
(hrany, bloby, ...)

	podobná data	odlišná data
málo dat	trénovat spíše jen poslední vrstvu	problém 😊
hodně dat	fine tune několika vrstev (lze ale i celou síť)	fine tune více vrstev nebo i celé sítě