

Zpracování digitalizovaného obrazu (ZDO) - Klasifikace

Přehled klasifikačních metod

Ing. Zdeněk Krňoul, Ph.D.

Katedra Kybernetiky
Fakulta aplikovaných věd
Západočeská univerzita v Plzni



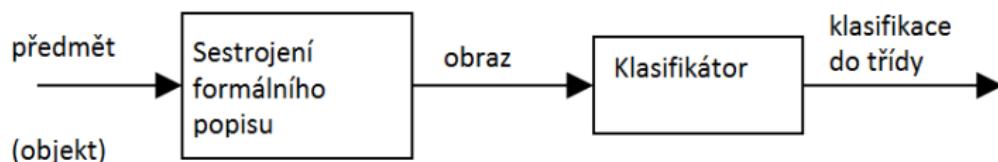
Obsah:

- ▶ klasické přístupy
- ▶ DNN (CNN)



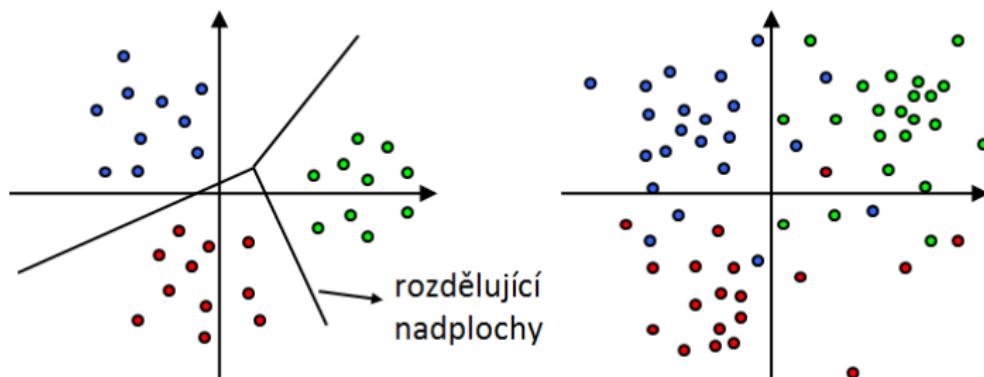
Klasifikace

- ▶ Rozpoznávání (klasifikace, angl. Pattern recognition) – zařazování předmětů do tříd
- ▶ Klasifikátor nerozeznává objekty, nýbrž jejich obrazy (popisy)
- ▶ PŘÍZNAKOVÉ ROZPOZNÁVÁNÍ
- ▶ STRUKTURÁLNÍ (SYNTAKTICKÉ) ROZPOZNÁVÁNÍ

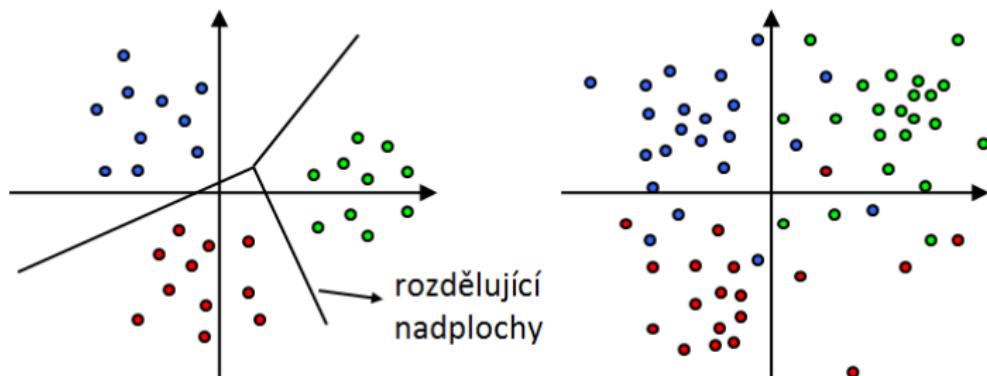


PŘÍZNAKOVÉ ROZPOZNÁVÁNÍ

- ▶ obrazy jsou charakterizovány vektorem, jehož souřadnice tvoří hodnoty jednotlivých příznaků.
- ▶ množina všech možných obrazů vytváří n-rozměrný obrazový prostor.
- ▶ při vhodném výběru příznaků je podobnost předmětů v každé třídě vyjádřena geometrickou blízkostí jejich obrazů



- pokud lze obrazy jednotlivých tříd (různých) od sebe oddělit rozdělující nadplochou → (mluvíme o separabilních množinách obrazů),
- úloha klasifikace je relativně jednoduchá a lze očekávat bezchybnou klasifikaci
- ve valné většině případů však množiny obrazů nejsou stoprocentně separabilní a část předmětů bude vždy chybně klasifikována, viz obrázek vpravo



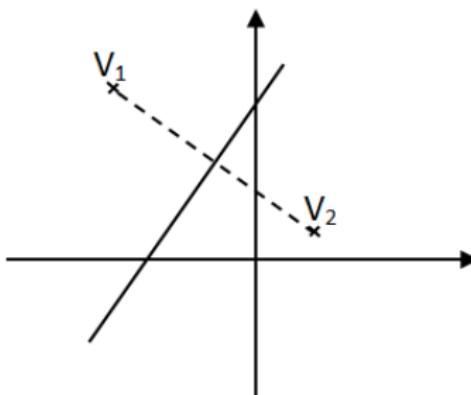
Učení s učitelem - Supervised Learning

- ▶ Základem je trénovací množina vzorových obrazů
- ▶ + u každého je uvedeno zařazení do správné třídy.
- ▶ na základě této trénovací množiny je pak určena reprezentace tříd,
- ▶ např. pro metodu **minimální vzdálenosti**



Metoda minimální vzdálenosti

- ▶ metodou minimální vzdálenosti jsou vypočteny centroidy (průměry) vzorových obrazů pro jednotlivé třídy.
- ▶ každá třída je reprezentována jedním vzorovým obrazem
- ▶ tento obraz lze vypočítat např. průměrem všech vzorových obrazů dané třídy



Klasifikátor KNN (K-Nearest Neighbours)

- ▶ each image is matched with all images in training data
- ▶ top K with minimum distances are selected
- ▶ majority class of those top K is predicted as output class of the image
- ▶ various distance metrics can be used like
 - ▶ L1 distance (sum of absolute distance) $E = \sum_{i=1}^N |y_i - f(x_i)|$,
 - ▶ L2 distance (sum of squares) $E = \sum_{i=1}^N (y_i - f(x_i))^2$,
 - ▶ etc.



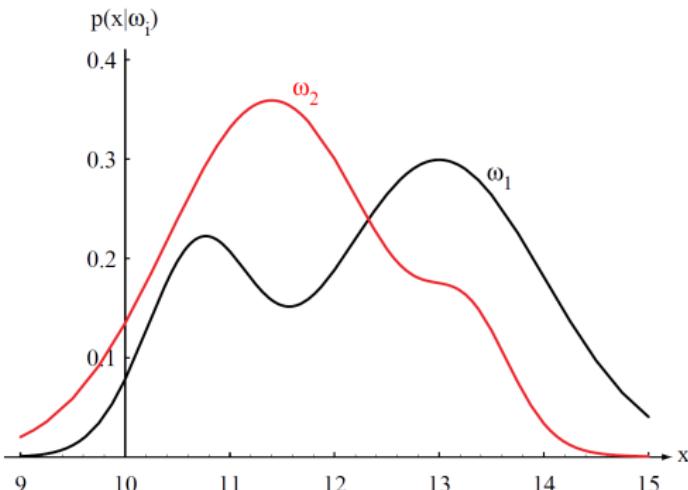
Bayesův klasifikátor

- ▶ základní statistický přístup ke klasifikaci založený na Bayesovu teorému (Thomas Bayes 1701–1761)
- ▶ patří do množiny tzv. pravděpodobnostních klasifikátorů, tedy vyčíslení kompromisu mezi různými klasifikačními rozhodnutími za využití pravděpodobnosti
- ▶ princip je znám od 50-tých a 60-tých let minulého století, kdy byl používán na úlohu vyhledávání informací v textu (information retrieval), kategorizaci textu, vyhodnocení diagnoz v lékařství aj.
- ▶ stále je považován za alternativu k více sofistikovanějším klasifikačním metodám, jako je např. SVM (support vector machine)

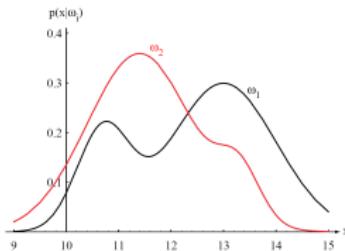


- ▶ ze statistického pohledu Bayesův klasifikátor minimalizuje pravděpodobnost chybné klasifikace,
- ▶ princip založen na podmíněné distribuci $P(Y|X)$), kde Y je konečná množina možných tříd a X je vstup (měřené příznaky).
- ▶ potom *rozhodnutí* pro nejlepší třídu pro námi naměřené X je dáno $\hat{y} = \operatorname{argmax}_y P(Y = y|X)$
- ▶ trénování Bayesova klasifikátoru ... často metoda maximální věrohodnosti (maximum likelihood, MLE)
- ▶ pokud jsou proměnné X nezávislé, pak je potřeba určit celkem **málo parametrů**, např. x se řídí normálním (Gausovským) rozdelením → statistický popis jednotlivých tříd je střední hodnota a rozptyl (není potřeba plná kovarianční matice)





- ▶ třídy ω_1 a ω_2
- ▶ $p(x|\omega_j)$... značí podmíněné závislosti vzorku,
- ▶ x ... spojité náhodná veličina, jejíž distribuce závisí na *přirozeném stavu*
- ▶ dále zavedeme apriorní pravděpodobnost, která předurčuje příslušnost vzorku k dané přídě ... tedy pro třídu j označíme jako $P(\omega_j)$,



Trénování:

- ▶ pro dvě třídy hledáme $p(x|\omega_1)$ a $p(x|\omega_2)$ pro podmíněnou hustotu pravděpodobnosti třídy ω_1 resp. ω_2
- ▶ dále dvě apriorní pravděpodobnost $P(\omega_1)$, která předurčuje, že vzorek patří do třídy ω_1 , resp. $P(\omega_2)$ pro třídu ω_2 ... neexistují žádné další třídy, proto $P(\omega_1) + P(\omega_2) = 1$
- ▶ rozdíl mezi $p(x|\omega_1)$ a $p(x|\omega_2)$ popisuje rozdíl mezi jednotlivými třídami

Sdružená hustota pravděpodobnosti $p(\omega_j, x)$

- ▶ sdružená hustota pravděpodobnosti popisuje, že vzorek patří do třídy ω_j a má hodnotu příznaku x
- ▶ může být zapsáno dvěma způsoby:
$$p(\omega_j, x) = P(\omega_j|x)p(x) = p(x|\omega_j)P(\omega_j),$$
- ▶ přeuspořádáním tohoto vztahu dostaneme Bayesův vztah:

$$P(\omega_j|x) = \frac{p(x|\omega_j)P(\omega_j)}{p(x)}, \quad (1)$$

kde $p(x) = \sum_j p(x|\omega_j)P(\omega_j)$

- ▶ jde o Bayesův vztah, který může být vyjádřen slovně jako:

$$\text{posterior} = \frac{\text{likelihood} \times \text{prior}}{\text{evidence}} \quad (2)$$



Učení bez učitele - shluková analýza

- ▶ u trénovací množiny není udána informace o příslušnosti obrazů k třídám
- ▶ snahou je rozdělit obrazy do k tříd tak, aby byla minimalizována hodnota kritéria optimality
- ▶ často globální minimum nelze z výpočetních důvodů nalézt → přijatelné lokální minimum
- ▶ Metody lze rozdělit na:
 1. hierarchické – vytvářejí shlukovací strom:
 - aglomerativní – vycházíme od jednotlivých obrazů a postupně spojujeme menší shluky do větších
 - divizní – vycházíme od celé trénovací množiny jako jednoho shluku a postupně dělíme větší shluky na menší
 2. nehierarchické – různé iterační metody – např. MacQueenův algoritmus (k-means)



STRUKTURÁLNÍ (SYNTAKTICKÉ) METODY

- ▶ Syntaktický popis
 - ▶ je vhodný tam, kde potřebujeme zachytit strukturu objektů
 - ▶ nebo kde pro jejich složitost chceme využít strukturu pro rozpoznávání
- ▶ Syntaktický popis objektu je hierarchická struktura jeho elementárních vlastností
- ▶ elementární vlastnosti se zde nazývají **primitiva**
- ▶ Obraz je pak reprezentován **řetězcem primitiv**
- ▶ Množina všech primitiv bývá nazývána **abecedou**
- ▶ Množina všech řetězců, pomocí nichž lze charakterizovat obrazy jedné třídy, se nazývá **jazyk popisu**
- ▶ Jazyk je generován nějakou **gramatikou**.
- ▶ Gramatika je soubor pravidel, pomocí nichž lze ze symbolů abecedy vytvářet řetězce, charakterizující možné tvary objektů.



STRUKTURÁLNÍ METODY - Trénování:

1. volba primitiv
2. konstrukce gramatiky generující řetězce, reprezentující možné tvary objektu
 - ▶ na základě znalosti úlohy, zkušenosti
 - ▶ případně z trénovací množiny
 - ▶ tento krok je převážně prováděn ručně
 - ▶ úloha automatické inference (odvozování) gramatik je jen velmi obtížně řešitelná



STRUKTURÁLNÍ METODY - Rozpoznávání:

1. příslušnost neznámého obrazu do dané třídy testujeme procesem tzv. syntaktické analýzy.
2. snažíme se pomocí gramatiky charakterizující danou třídu vygenerovat neznámý řetězec
3. např. konečné automaty



DNN (CNN)

- ▶ jedna z technik machine learning pro zpracování obrazu
- ▶ více vrstev - nelineální přístup jak pro extrakci příznaků, tak i jejich zpracování
- ▶ obrázek reprezentován jako vektor
- ▶ často jako feedforward NN ale i recurrent NN (LSTM language modeling),
- ▶ trénování backpropagation algoritmem s gradient descent technikou

Rozdíly oproti standardním přístupům:

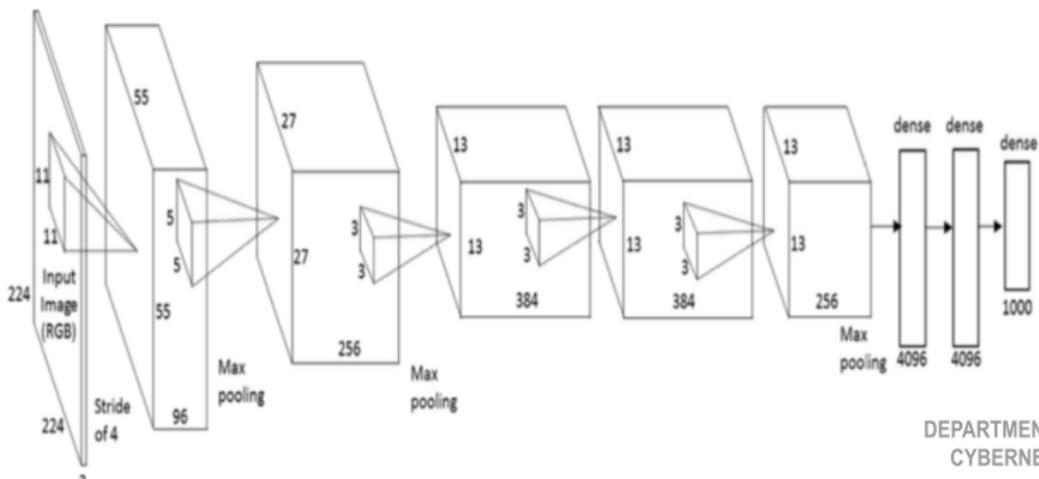
- ▶ trénovatelné feature extractory
- ▶ end-to-end systém, tedy včetně klasifikace (regrese)
- ▶ feature extractory ve více úrovních



Klasická CNN je kombinací následujících vrstev:

1. Convolution Layer
2. Pooling Layer
3. Fully Connected Layer

- ▶ první konvoluční vrstvy se starají o detekci hran
- ▶ Následující vrstvy kombinují informace do jednodušších tvarů
- ▶ následuje modelování menších kousků informace
- ▶ předposlední vrstvy (FC) kombinují (vážená kombinace) do finální predikce



Vlastnosti:

- ▶ konvoluční jádro (ve více vrstvách)
- ▶ maxpooling - maximální hodnota v oblasti (okně) ... dělá pod-vzorkování výběrem pouze významných odezv
- ▶ aktivační funkce (často RELu - rychlejší konvergence SGD a implementačně jednodušší)
- ▶ Dropout (regularizace přetrénování) - náhodné vymaskování některých neuronů v FC vrstvách
- ▶ inicializace vah



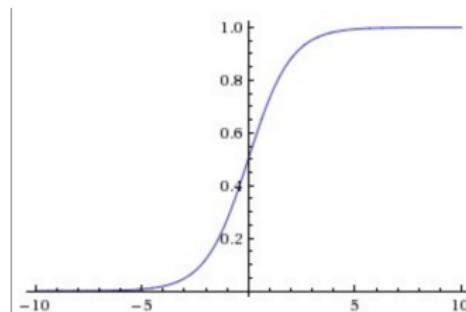
Inicializace vah

- ▶ Gaussian Random Variables
 - ▶ $\mu = 0$ a $\sigma = 0.01 \dots 10^{-5}$
- ▶ Xavier Initialization
 - ▶ inicializace Gaussem pro větší sítě způsobuje, že výstupy neuronů budou blízké nule
 - ▶ variance gausovská distribuce závislá na počtu vstupů neuronu (velikosti předchozí vrstvy)
 - ▶ doporučuje se $\text{var}(w) = \frac{1}{n_{in}}$ (Caffe framework) nebo i
$$\text{var}(w) = \frac{2}{n_{in} + n_{out}}$$



Aktivační funkce

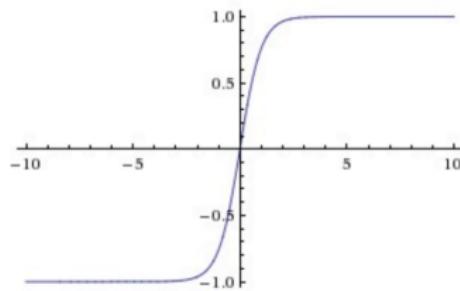
$$\text{Sigmoid Function } \sigma(x) = \frac{1}{(1+e^{-x})}$$



- ▶ existují problémy v souvislosti CNN
- ▶ saturované neurony potlačují gradient (viz -5 +5) kde gradient je skoro nula
- ▶ při zpětné propagaci se gradient na tomto neuronu "ztratí"
- ▶ náročný výpočet exp funkce
- ▶ výstup (0, 1) není centrován na nulu ... nevýhoda jako vstup do další vrstvy
- ▶ není proto využívána pro CNN

Aktivační funkce

tanh activation - tvar hyperbolické tangent funkce

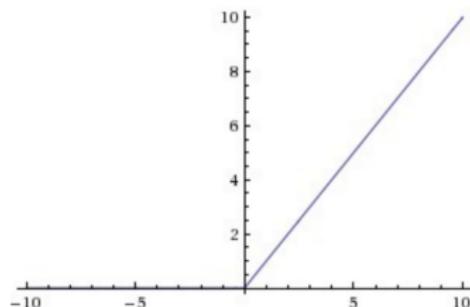


- ▶ stále problém ... saturované neurony potlačují gradient (viz -5 +5) kde gradient je skoro nula
- ▶ při zpětné propagaci se gradient na tomto neuronu "ztratí"
- ▶ náročný výpočet



Aktivační funkce

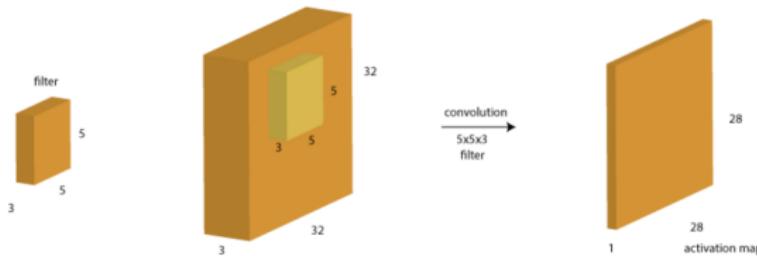
ReLU (Rectified Linear Unit) $f(x) = \max(0, x)$



- ▶ nejčastěji používaná aktivační funkce pro CNN
- ▶ rychlá konvergence trénování
- ▶ implementačně jednodušší
- ▶ gradient je stále potlačován pro záporné x

Konvoluční vrstva - Convolution Layer

Pro vstupní obrázek a daný konvoluční filtr máme výstup jako odezvu (aktivační mapu)



Pro více filtrů (zde na ukázku 10) pak získáváme tvar ve formě 10 odezv:



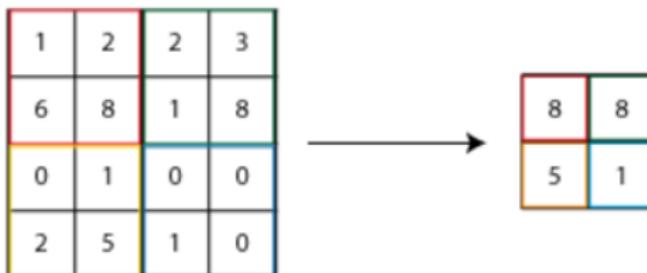
Vlastnosti:

- ▶ posun okénka (Stride) je definován pro danou konvoluční vrstvu
- ▶ hodnota posunu z principu 1 nebo více, závisí také i na zvolené velikosti filtru (okénka)
- ▶ ovlivňuje velikost následující vrstvy
- ▶ "Zero-padding" vyplnění nulou za okrajem obrázku (výpočet konvoluce)
- ▶ počet parametrů $V * K + K$
- ▶ V velikost (objem) filtru např. $5 \times 5 \times 3$
- ▶ K je počet filtrů



Pooling vrstva - Pooling Layer

- ▶ Když použijeme zero padding a stride jedna, pak velikost výstupu konvoluční vrstvy, tedy aktivační mapa(y), mají stejnou velikost
- ▶ musí se provést redukce velikosti, např. odezva 4×4 je redukována na 2×2 metodou max-pooling
- ▶ možné i jiné techniky např. average pooling (GoogLeNet)



Plně spojená vrstva - Fully Connected Layer

- ▶ jsou použity ke konci CNN
- ▶ každý neuron trénováním získává danou funkci důležitou pro rozhodnutí
- ▶ při trénování aplikován Dropout
- ▶ poslední vrstva má velikost rovnou počtu klasifikačních tříd, nebo velikosti regresního vektoru
- ▶ FC vrstvy nemusí být v architektuře použity



DNN as a machine learning model

- ▶ in the case of classification or regression we try to:
 - ▶ such values of ω_i that will minimize a chosen criterion
 - ▶ criterium usually incorporates information from teacher: t

Classification criteria:

- ▶ Binary cross entropy :

$$E_k = - \sum_i t_i \log o_i + (1 - t_i) \log(1 - o_i) \quad (3)$$

- ▶ Categorical cross entropy:

$$E_k = - \sum_i t_{k,i} \log o_{k,i} \quad (4)$$

(with soft-max layer)

Regression criteria: Mean squared/absolute error

$$E_k = \sum_i (t_i - o_i)^2 \quad (5)$$



Trénování

- ▶ často použit algoritmus Gradient Descent

$$\omega_{n+1} = \omega_n - \gamma_n \nabla E(\omega_n) \quad (6)$$

- ▶ kde $E(\omega_n)$ je zvolené kritérium (nelineární funkce) a ω_n parametry v kroku n ,
- ▶ γ_n Learning rate ... u GD metod bohužel neznáme velikost kroku optimalizace

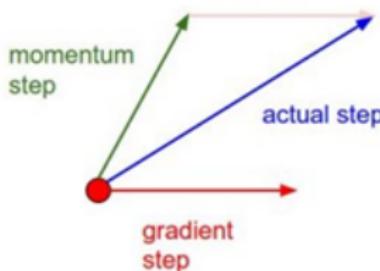


Realizace:

- ▶ omezená velikost paměti
- ▶ mini-batch (Stochastic Gradient Descent - SGD metoda)
- ▶ náhodně rozdělená data, update postupně na balíčcích
- ▶ rychlosť konvergencie - momentum SGD

$$\Delta \omega_n = \gamma_n \nabla E(\omega_n) + \alpha \Delta \omega_{n-1} \quad (7)$$

- ▶ α momentum = váha jak *cumulative gradient direction* bude při optimalizaci směrován

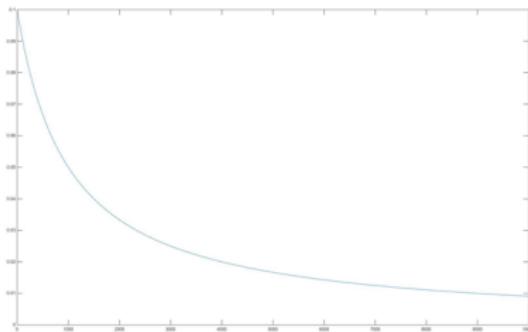


Learning rate decay:

- ▶ změna learning rate během trénování
- ▶ **Step decay** skokové snižování po určitém počtu trénovacích epoch
- ▶ **Exponential decay**
- ▶ **1/t decay**

1/t decay

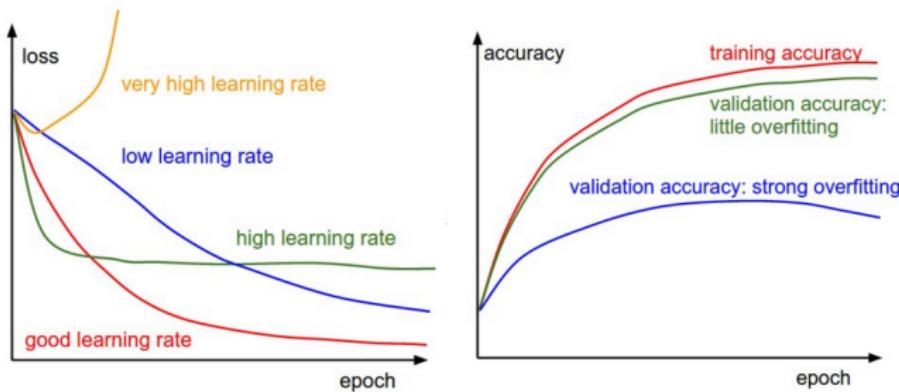
* $\gamma_0 = 0.1$
* $k = 0.001$



Per-parameter adaptive LR methods:

- ▶ adaptivně mění Learning rate pro jednotlivé parametry
- ▶ **Adagrad**
- ▶ **RMSprop**
- ▶ **Adam**



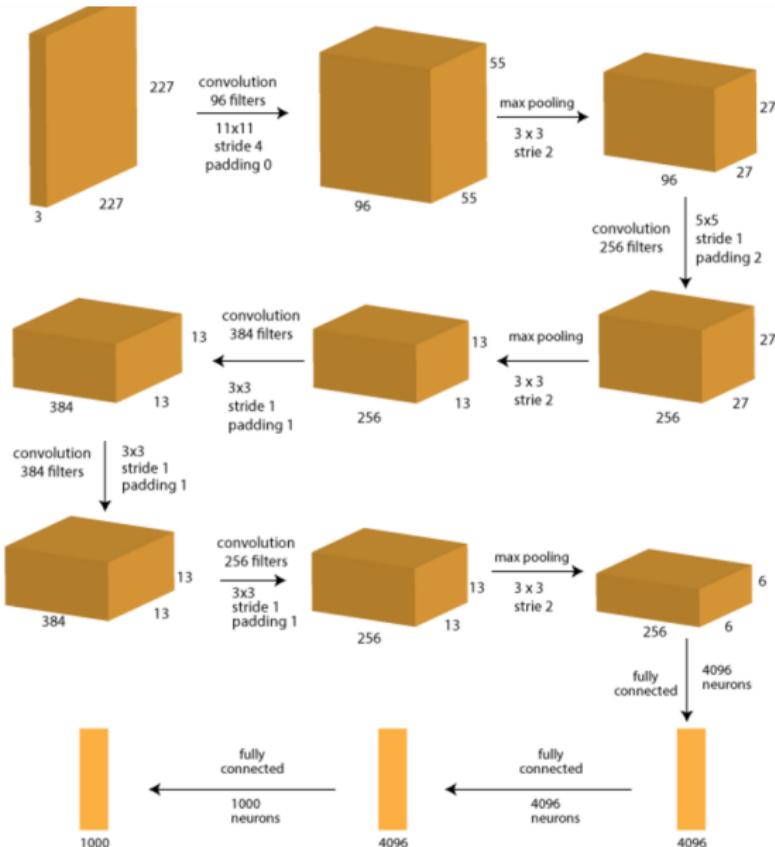


- ▶ zde vidíme efekt rozdílného Learning rate
- ▶ rozdíl mezi trénovací a validační přesností udává míru přetrénování

AlexNet (2012)

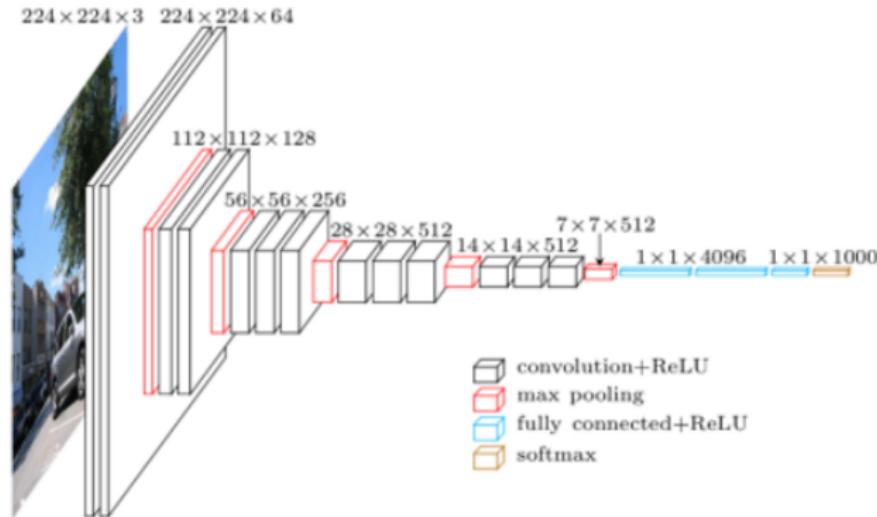
- ▶ vítěz IMAGENET Challenge 2012
- ▶ 1000 tříd, 1.3 mil. trénovacích obrázků
- ▶ nebyl zcela zvládnut problém trénování na GPU
- ▶ v dané době získali významně lepších výsledků než jiné přístupy
- ▶ označena jako první z "moderních" CNN
- ▶ hloubka 11
- ▶ možné použít jako inicializace pro podobnou úlohu + dotrénování





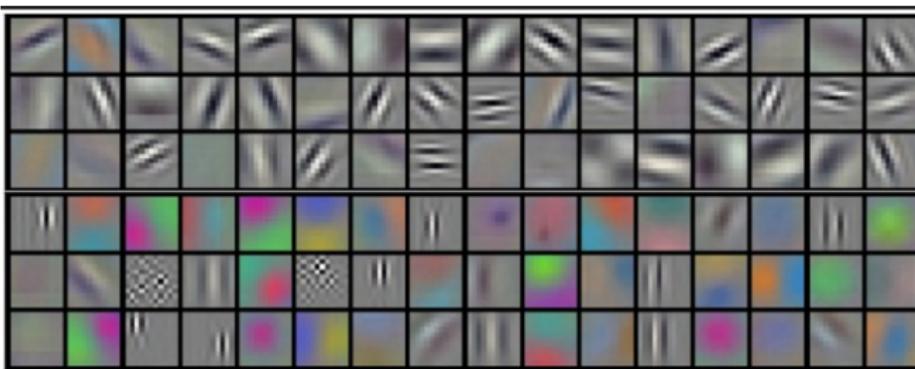
VGGNet (2014)

- ▶ malé jen 3×3 kernely, stride 1
- ▶ augmentace dat
- ▶ hloubka 19



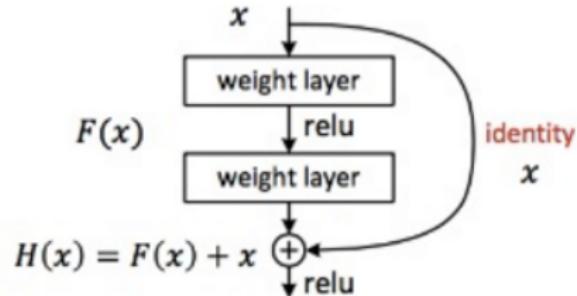
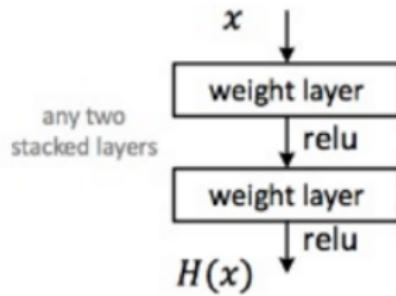
GoogLeNet (2014)

- ▶ vítěz IMAGENET Challenge 2014
- ▶ jen kernely
- ▶ vzniká otázka ztráty gradientu pro hlubší sítě
(BackPropagation)
- ▶ problémy s FC vrstvami (tendence k přetrénování)
- ▶ hloubka 22



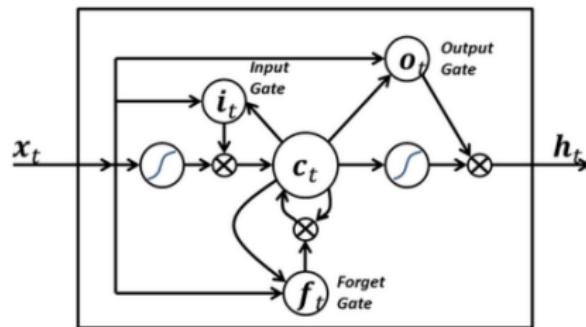
ResNet (2014)

- ▶ vítěz IMAGENET Challenge 2015
- ▶ inspirováno VGG, ale hloubka **152**
- ▶ aktivační funkce P-RELU (adaptivní RELU - více parametrů)
- ▶ vzniká problém ztráty gradientu pro hlubší sítě (BackPropagation) ... cca už od hloubky 30
- ▶ vymyšleny a použity tzv. **shortcut connection**



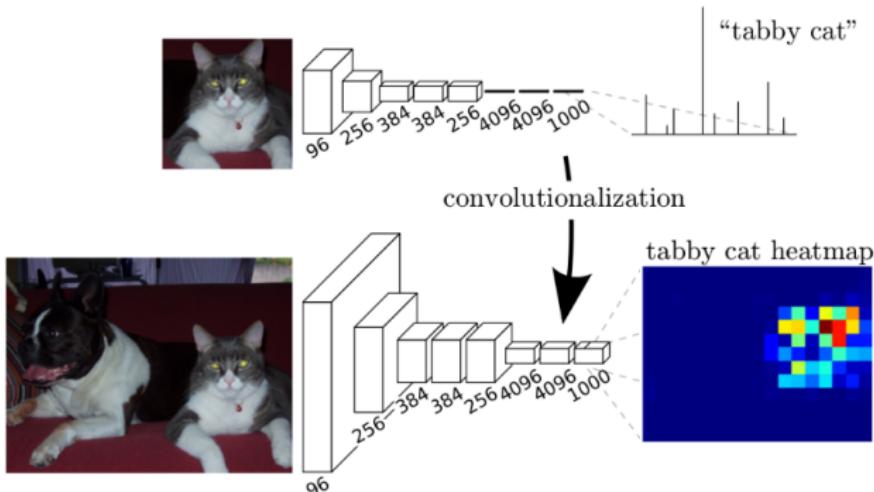
Long short-term memory (LSTM) 1997

- ▶ RNN architektura, jako alternativa k HMM
- ▶ v současnosti velký úspěch v aplikacích v ASR nebo AVSR (audio + obraz rtů), Handwriting recognition, Robot control aj.
- ▶ řeší problém RNN se ztrátou gradientu s rostoucí dobou mezi důležitými událostmi
- ▶ na množině trénovacích sekvenčních



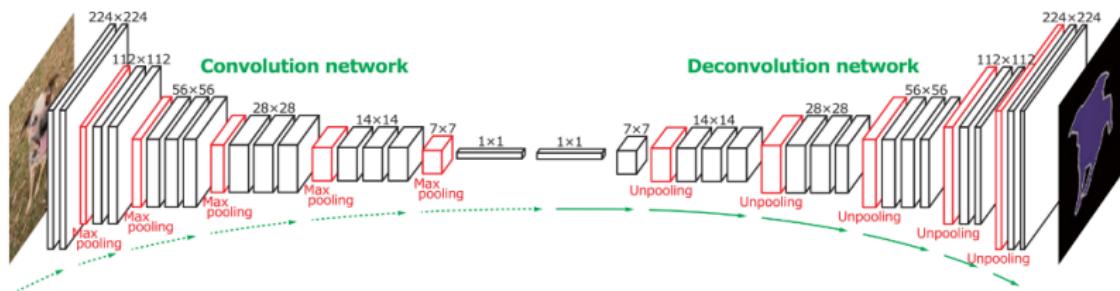
FCN - Fully Convolutional Networks:

- ▶ FCN má pouze konvoluční vrstvy
- ▶ nejsou FC vrstvy
- ▶ přebírá úlohu segmentace - výstupem je heat mapa
- ▶ použitelná na libovolnou velikost vstupního obrázku



Deconvolution Network:

- ▶ místo full connected vrstev jsou použity de-konvoluční a un-pooling vrstvy
- ▶ dekonvoluce časo-prostorových feature map
- ▶ produkuje per-pixel predikci (labels)
- ▶ může např. vycházet AlexNet, VGG aj. architektur

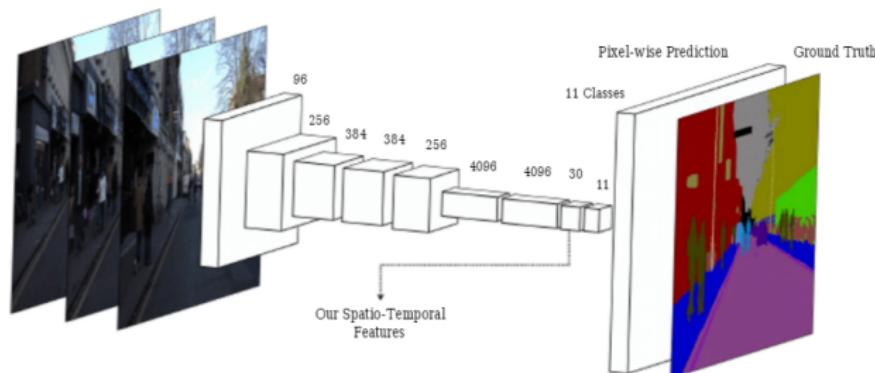


1

¹Learning Deconvolution Network for Semantic Segmentation

Spatio-Temporal FCN for Semantic Video Segmentation:

- ▶ vychází z AlexNet architektury
- ▶ použity video data LSTM verze rekurze
- ▶ dekonvoluce časo-prostорových feature map ... produkuje per-pixel predikci



2

²STFCN: Spatio-Temporal FCN for Semantic Video Segmentation

DEPARTMENT OF
CYBERNETICS



AARSHAY JAIN, Deep Learning for Computer Vision – Introduction to Convolution Neural Networks

