# Neural Networks exam 2019/20

No. 1 Topic

- a. Local optimization in non-convex cases (reason for non-convexity)
- b. RMSP optimizer
- c. Dropout
- d. ResNet
- e. Gradient ascent

## No. 2 Topic

- a. Weight update strategies
- b. ReLU and dying ReLU problem
- c. LSTM cell
- d. Convolution as a mathematical operation in continuous and discrete cases

## No. 3 Topic

- a. Newton optimization method
- b. Ensembling, bagging
- c. Comparison of loss functions
- d. Machine learning vs traditional programming
- e. Inception

## No. 4 Topic

- a. McCulloch-Pitts model
- b. Parameters of conv filter, stride, padding, etc
- c. Linear classifier, margin of the classifier
- d. Data augmentation
- e. YOLO

## No. 5 Topic

- a. Statistical learning theory
- b. Various activation functions and their properties
- c. Autoencoders
- d. Graph unrolling and parameter sharing in recurrent NN
- e. MobileNet

## No. 6 Topic

- a. Machine learning problem definition
- b. Newton optimizer
- c. Effects and relationship of model capacity and complexity overfitting, underfitting
- d. t-distributed Stochastic Neighbor Embedding
- e. ShuffleNet

## No. 7 Topic

- a. Credit approval problem
- b. Objective functions in neural networks
- c. Nesterov momentum optimizer
- d. Decomposition of convolutional kernels
- e. Alexnet + ILSVRC

## No. 8 Topic

- a. Delta learning rule
- b. Batch normalization
- c. Transposed conv., atrous conv.
- d. Object classification, localization VS object detection, semantic segmentation VS instance segmentation
- e. ResNext

## No. 9 Topic

- a. ADAM optimizer
- b. The softmax function
- c. R-CNN architectures: R-CNN, Fast R-CNN, Faster R-CNN
- d. Supervised VS unsupervised learning
- e. EfficientNet
- No. 10 Topic
  - a. Optimization problem of objective functions of NN
  - b. AdaGrad optimizer
  - c. Input vector normalization
  - d. DeconvNet, U-Net
  - e. Neural style transfer

## No. 11 Topic

- a. Multilayer perceptron
- b. Early stopping
- c. Gradient descent (multidimensional cases as well)
- d. Weight regularization (L1, L2)
- e. Pooling

## No. 12 Topic

- a. Perceptron convergence theorem (no proof)
- b. Momentum optimizer
- c. Properties of CNN: sparsity, parameter sharing, equivariance, invariance to shifting
- d. RNN examples: predicting the next letter, image captioning
- e. Adversarial attacks

## No. 13 Topic

- a. Elementary set separation by a single neuron
- b. Local response normalization
- c. Unpooling
- d. Representations: Blum and Li theorem, construction

## No. 14 Topic

- a. Principal component analysis (PCA)
- b. Back-propagation through time
- c. Stochastic gradient descent optimizer
- d. Object detection problem explained
- e. Effects of filter size on convolution

## No. 15 Topic

- a. Rosenblatt perceptron training algorithm
- b. Back-propagation
- c. Curse of dimensionality
- d. SqueezeNet
- e. Back-propagation and gradient-based optimizers

## No. 1 Topic

a. Local optimization in non-convex cases (reason for non-convexity)



## b. RMSP optimizer

# **RMSP** algorithm

- The RMSProp algorithm (2012) modifies AdaGrad to perform better in the nonconvex setting by changing the gradient accumulation into an exponentially weighted moving average
- In each step AdaGrad reduces the learning rate, therefore after a while it stops entirely!
- AdaGrad shrinks the learning rate according to the entire history of the squared gradient and may have made the learning rate too small before arriving at such a convex structure
- RMSProp uses an exponentially decaying average to discard history from the extreme past so that it can converge rapidly after finding a convex bowl, as if it were an instance of the AdaGrad algorithm initialized within that bowl

RMSP algorithm					
Algorithm The RMSProp algorithm	The closer parts of the history are counted more				
<b>Require:</b> Global learning rate $\epsilon$ , decay rate $\rho$ .					
<b>Require:</b> Initial parameter $\boldsymbol{\theta}$	strongly.				
<b>Require:</b> Small constant $\delta$ , usually $10^{-6}$ , used to stabilize division by small					
numbers.					
Initialize accumulation variables $\boldsymbol{r} = 0$					
while stopping criterion not met do					
Sample a minibatch of m examples from the training set $\{x^{(1)}, \ldots, x^{(m)}\}$ with					
corresponding targets $\boldsymbol{u}^{(i)}$ .					
Compute gradient: $\boldsymbol{q} \leftarrow \frac{1}{2} \nabla_{\boldsymbol{\theta}} \sum_{i} L(f(\boldsymbol{x}^{(i)}; \boldsymbol{\theta}), \boldsymbol{y}^{(i)})$					
Accumulate squared gradient: $\mathbf{r} \leftarrow \rho \mathbf{r} + (1 - \rho) \mathbf{q} \odot \mathbf{q}$					
Compute parameter update: $\Delta \theta = -\frac{\epsilon}{\sqrt{s_{\pm}}} \odot g$ . $(\frac{1}{\sqrt{s_{\pm}}}$ applied element-wise)					
Apply update: $\theta \leftarrow \theta + \Delta \theta$					
end while					

## c. Dropout

neurons



No computational penalty in testing phase

Use it for fully connected layers

temporally deactivated

## Reasoning behind dropout

- Dropout can be considered as averaging of multiple thinned networks ("ensemble")
- Dropout avoids training separate models - Would be very expensive
- Avoids computatinal penalty in the test phase
- But still gets benefits of ensemble methods

Reduces overfitting, because the network is forced to learn the functionality in different configurations using different neural paths.

## Intuitive explanation

Imagine that you have a team of workers and the overall goal is to learn how to erect a building. When each of the workers is overly specialized, if one gets sick or makes a mistake, the whole building will be severely affected. The solution proposed by "dropout" technique is to pick randomly every week some of the workers and send them to business trip. The hope is that the team overall still learns how to build the building and thus would be more resilient to noise or workers being on vacation.



## Data regularization techniques:

Temporal Modification of the net architecture in training phase – Dropout



## d. ResNet

## Deep Residual Networks (ResNets)

- "Deep Residual Learning for Image Recognition". CVPR 2016
- A simple and clean framework of training "very" deep nets

### • State-of-the-art performance for

- Image classification
- Object detection
- Semantic segmentation
- and more...

## Deep Residual Learning



## **Deep Residual Learning**



Kaiming He, Xiangyu Zhang, Shaoqing Ren, & Jian Sun. "Deep Residual Learning for Image Recognition". CVPR 201

, Shaoqing Ren, & Jian Sun. "Deep Residual Learning for

## Deep Residual Learning

• F(x) is a residual mapping w.r.t. identity



# Static samples vs Data signal flow

AlexNet could recognize 1000s of images. ResNet could reach better then human performance.

•	Though human can recognize - Single letters - Single sounds - Single tunes - Single pictures	→	But in real life we handle – Texts – Speech – Music – Movies	Story (temporal analysis of sequential data)

Can feed-forward neural networks (perceptrons, conv. nets) solve these problems?

DATA MEMORY

Naturally, we can extend the data dimension with the time, but this leads to data size and computational load explosion



ResNets had the lowest error rate at most competitions since 2015

1st places in all five main tracks

- ImageNet Classification: "Ultra-deep" 152-layer nets
- ImageNet Detection: 16% better than2nd
- ImageNet Localization: 27% better than2nd
- COCO Detection: 11% better than2nd
- COCO Segmentation: 12% better than2nd

## How do ResNets address these issues?



## https://towardsdatascience.com/introduction-to-resnets-c0a830a288a4

## e. Gradient ascent

# **Gradient Ascent**

Generate a synthetic image that maximizes the response of a neuron.



This image has to be "natural". The response should not depend on pixels and can not have arbitrary values

- Guassian blur on the image
- Clipping image values
- Clipping small gradients to 0

 $I^* = \arg \max_{I} f(I) + R(I)$ Neuron value Natural image regularizer

Intermediate Layers



Classes

Using a network which can learn feature inversion



# **Gradient Ascent**

Finding the maximizing patterns for each kernelImage: Second second

# **Gradient Ascent**



We could search in our database and find typical samples.

It helps, but usually the network is good on this set (train accuracy). We are curious about those images which the network has not seen.

Could we generate and ideal image for a given class?







## Gradient Ascent - activation maximization

We could search in our database and find typical samples.

It helps, but usually the network is good on this set (train accuracy). We are curious about those images which the network has not seen.

Could we generate and ideal image for a given class?



# Neural Style Transfer



An interesting application of the gradient ascent method is neural style transfer Could we use an input image and transform it into the style of an other input image?

https://demos.algorithmia.com/deep-style/



# Neural Style Transfer

Could we use an input image and transform it into the style of an other input image?

Gradient ascent transforms the image according to a loss function.

Can we find a loss function, which would preserve objects and another which preserves features connected to style?



# Neural Style Transfer

Could we use an input image and transform it into the style of an other input image?

Gradient ascent transforms the image according to a loss function.

Can we find a loss function, which would preserve objects and another which preserves features connected to style?



More weight to content loss

 More weight style loss

# Fast Neural Style Transfer



Could we use an input image and transform it into the style of an other input image?

Gradient ascent transforms the image according to a loss function.

Can we find a loss function, which would preserve objects and another which preserves features connected to style?



## No. 2 Topic

## a. Weight update strategies

## Weight update strategy

- Apply all the input vectors in one after the others, selecting them randomly
- Instance update
  - Update the weights after each input
- Batch update
  - Add up the modifications
  - Update the weights with the sum of the modifications,
  - after all the inputs were applied
- Mini batch
  - Select a smaller batch of input vectors, and do with that as
  - in the batch mode



#### Weight update: very simple example Start again: Test with the again with the first vector The result is OK! - Do not modify!!! 20 Test with the again with the second vector W. The result is OK! Do not modify!!! Q Test with the again with the third vector 3 The result is OK! Do not modify!!! Since all input vectors are correctly classified: we are ready 9/17/2019 P-ITEFA-0011 lei

## Formalization of the update rules

- Positive misclassification : ADD
   ε = d<sub>j</sub> y<sub>j</sub> = 1
   w(k+1)=w(k)+x<sub>j</sub>
   w(k+1)=w(k)+x\_j
   w(k)+x\_j
   w(k
- Negative misclassification : SUBTRACT ε = d<sub>j</sub>-y<sub>j</sub> = -1 w(k+1)=w(k)-x<sub>j</sub>
- Correct classification : DO NOTHING
   ε = d<sub>j</sub>·y<sub>j</sub> = 0 w(k+1)=w(k)
- In general:

```
w(k+1)=w(k)+\varepsilon x_i
```



## **Back-propagation**

- Though we showed how to modify the weights with back propagation, its most important value that it can calculate the gradient
- The weight updates can be calculated with different optimization methods, after the gradients are calculated
- Various optimization method can drastically speed up the training (100x, 1000x)

## Sequential back propagation

Adapting the weights of the FFNN (recursive algorithm)

$$w_{ij}^{(l)}(k+1) = w_{ij}^{(l)}(k) + \Delta w_{ij}^{(l)}(k)$$
$$\Delta w_{ii}^{(l)}(k) = ?$$

The weights are modified towards the differential of the error function (delta rule):  $\partial R_{max}$ 

$$\Delta w_{ij}^{(l)} = -\eta \, \frac{\partial \kappa_{emp}}{\partial w_{ij}^{(l)}}$$

The elements of the training set adapted by the FFNN sequentially B = B = (x(x), d)

$$R_{emp} = R_{emp}(y(\mathbf{x}), d)$$

# What are we optimizing here?

 Cost function in quadratic case for one x<sub>i</sub> → d<sub>i</sub> pair:

 $\mathcal{E}_i = \left(\mathbf{d}_i - Net(\mathbf{x}_i, \mathbf{w})\right)^2$ 

- Error surface is in the w space
- Error surface depends on the
- x<sub>i</sub> → d<sub>i</sub> pair
   Moreover, we do not see the entire surface, just

 $\mathcal{E}$  and the gradients  $\frac{\partial \mathcal{E}}{\partial w_{ij}^{(l)}}$ 



a Pro

# Update strategies

- Single vector update approach (instant update) •
  - Weights are updated after each input vector
- Batched update approach
  - All the input vectors are applied
    - this is actually the correct entire error funtion, which is used by the original Gradient Descent Method
  - Updates  $(\Delta w_{ii})$  are calculated for each vector, and averaged
  - Update is done with the averaged values  $(\Delta w_{ij})$  after the entire batch is calculated
- Mini batch approach
  - When the number of inputs are very high (104-106), batch would be ineffective
  - Random selection of m input vectors (m is a few hundred)
  - Updates  $(\Delta w_{ij})$  are calculated for each vector, and averaged
  - Update is done with the averaged values ( $\Delta w_{ii}$ ) after the mini batch is calculated
  - Works efficiently when far away from minimum, but inaccurate close to minimum

10/1/2019 Requires reducing learning rate



- · Activation function
- · Half-wave rectifier
- Not compressing the gradient
- learns much faster
- ReLU types
  - Softmax
  - Leaky ReLU
  - ELU, SELU Relu6

29



*a* is a hyper-parameter: Tuned during training

- Variation of leaky ELU
- Two fixed parameters
  - Not trained, but selected to be fixed
  - $-\lambda$  is the scaling parameter



a is a hyper-parameter to be tuned and  $a \geq 0$  is a constraint.



ReLUs do not require input normalization to prevent them from saturating. However, Local Response Normalization aids generalization.











# Unrolling LSTM network





# Gradient calculation in LSTM Uninterrupted gradient flow! $c_{0} + f_{1} + f_{2} + f_{3} + h_{1} + h_{2} + h_{2} + h_{3} + h_{4} +$

- Though we multiply the memory content with a smaller than 1 number
- And the W matrix is part of the memory update
- But it still preserves the content for longer time
- As it comes from the name: It is a elongated time <u>short term</u> memory

# Achevements with LSTM networks

- Record results in natural language text compression
- Unsegmented connected handwriting recognition
- Natural speech recognition
  - Smart voice assistants
    - Google Translate

.

- Amazon Alexa
- Microsoft Cortana
- Apple Quicktype
- G X Google Translate
  - Hi, I'm Cortana.

ron

from

iom

- 95.1% recognition accuracy on the Switchboard corpus, incorporating a vocabulary of 165,000 words
  - Continuous spontaneous English native speech



# Variants of LSTM II : Joined forget and input ht Input and forget gates has practically the same role Why not to join them?



# d. Convolution as a mathematical operation in continuous and discrete cases

# Convolution

- Convolution is a mathematical operation that
  - does the integral of the product of 2 functions (signals),
  - with one of the signals flipped and shifted
- Mathetmatically:

$$egin{aligned} &(f*g)(t) \stackrel{ ext{def}}{=} \int_{-\infty}^\infty f( au) g(t- au) \, d au \ &= \int_{-\infty}^\infty f(t- au) g( au) \, d au \end{aligned}$$

• Convolution is commutative 10/22/2019



- 1. Flipp g signal
- 2. Slide the flipped *g* over *f*
- Integrate the product in continious space or Multiply and accumulate it in discrete space with each shift

# **Discrete convolution**

• For continiuous:

$$egin{aligned} (f st g)(t) \stackrel{ ext{def}}{=} & \int_{-\infty}^{\infty} f( au) g(t- au) \, d au \ &= & \int_{-\infty}^{\infty} f(t- au) g( au) \, d au \end{aligned}$$

• For discrete functions:

$$egin{aligned} (fst g)[n] &= \sum_{m=-\infty}^\infty f[m]g[n-m] \ &= \sum_{m=-\infty}^\infty f[n-m]g[m] \end{aligned}$$

# Convolution vs correlation II

## As the only difference is kernel flipping...

Why convolution rather that correlation?

- Commutativity, Associativity, Distributivity helps to prove mathematical statements
- Since the network learns its own weights, it is invariant whether that flip is there or not (just a convention)
- In many cases, correlation is implemented even when it is called convolution



## No. 3 Topic

## a. Newton optimization method

Simplest 2<sup>nd</sup> order Gradient descent method: Newton Method

$$f(\mathbf{x}) \approx f(\mathbf{x}_0) + (\mathbf{x} - \mathbf{x}_0)^{\mathrm{T}} \nabla f(\mathbf{x}_0) + \frac{1}{2} (\mathbf{x} - \mathbf{x}_0)^{\mathrm{T}} \mathbf{H}(f(\mathbf{x}_0)) (\mathbf{x} - \mathbf{x}_0)$$

• Replacing  $(\mathbf{x} - \mathbf{x}_0) \rightarrow \Delta \mathbf{x}$  and differentiating it with  $\Delta \mathbf{x}$ , assuming that we can jump to a minima, where:  $\nabla f(\mathbf{x}) \approx 0$ 

$$0 = \frac{\partial}{\partial \Delta \mathbf{x}} \left( \begin{array}{c} f(\mathbf{x}_0) + \Delta \mathbf{x}^{\mathrm{T}} \nabla f(\mathbf{x}_0) + \frac{1}{2} \Delta \mathbf{x}^{\mathrm{T}} \mathbf{H}(f(\mathbf{x}_0)) \Delta \mathbf{x} \\ f(\mathbf{x}_0) - \mathbf{x} \\ \text{Constant} \rightarrow 0 \\ (\Delta x)' \rightarrow 1 \\ (\mathcal{Y}_2(\Delta x)^2)' \rightarrow \Delta x \end{array} \right) = \nabla f(\mathbf{x}_0) + \mathbf{H}(f(\mathbf{x}_0)) \Delta \mathbf{x}$$

Newton optimization:  $\Delta \mathbf{x} = -\mathbf{H}(f(\mathbf{x}_0))^{-1} \nabla f(\mathbf{x}_0) \quad \mathbf{x}(n+1) = \mathbf{x}(n) - \eta \mathbf{H}(f(\mathbf{x}(n)))^{-1} \nabla f(\mathbf{x}(n))$ 

# Properties of Newton optimization method

- When *f* is a positive definite quadratic function, Newton's *method* jumps in a single step to the minimum of the function directly.
- Newton's method can reach the critical point much faster than 1<sup>st</sup> order gradient descent.

Newton optimization:  $\Delta \mathbf{x} = -\mathbf{H}(f(\mathbf{x}_0))^{-1} \nabla f(\mathbf{x}_0) \quad \mathbf{x}(n+1) = \mathbf{x}(n) - \eta \mathbf{H}(f(\mathbf{x}(n)))^{-1} \nabla f(\mathbf{x}(n))$ 

# Newton's algorithm

Algorithm Newton's method with objective  $J(\theta) = \frac{1}{m} \sum_{i=1}^{m} L(f(\boldsymbol{x}^{(i)}; \theta), y^{(i)}).$ Require: Initial parameter  $\theta_0$ Require: Training set of m examples while stopping criterion not met do Compute gradient:  $\boldsymbol{g} \leftarrow \frac{1}{m} \nabla_{\theta} \sum_{i} L(f(\boldsymbol{x}^{(i)}; \theta), \boldsymbol{y}^{(i)})$ Compute Hessian:  $\boldsymbol{H} \leftarrow \frac{1}{m} \nabla_{\theta}^2 \sum_{i} L(f(\boldsymbol{x}^{(i)}; \theta), \boldsymbol{y}^{(i)})$ Compute Hessian inverse:  $\boldsymbol{H}^{-1}$ Compute update:  $\Delta \theta = -\boldsymbol{H}^{-1}\boldsymbol{g}$ Apply update:  $\theta = \theta + \Delta \theta$ end while

Typically not used, due to the computational complexity

Parameter space much higher than first order (where it is already very high)

## b. Ensembling, bagging

Ensemble methods:

- Network duplication
- Bagging
- Dropout

# **Ensemble methods**

- Idea of ensemble methods:
  - Generate multiple copies of your net
    - Same or slightly modified architectures
  - Train them separately
    - Using different subsets of the training sets
    - Different objective functions
    - Different optimization methods
  - The different trained models have independent error characteristics
  - Averaging the results will lead to smaller error
- Requires more computation and memory both in training and inferencing (testing) phase



First learns the upper loop, the second the lower. When both say yes, it is an 8.

## Ensembling: duplicating the network I

- Train two architecturally identical copies of the network on two GPUs
  - Half of the neuron layers are on each GPU
    - GPUs communicate only in certain layers
  - Improvement (as compared with a net with half as many kernels in each convolutional layer trained on one GPU):
    - Top 1 error rate by 1.7%
    - Top 5 error rate by 1.2%

Ensembling: duplicating the network II



## c. Comparison of loss function

## Loss functions

- Loss function determines the training process
  - Tells the net, whether an error is big or small, and penalize accordingly
  - There can be other errors, not just the difference of the output and the desired output
- Most used loss function types:

K

Quadratic, in case of regression

$$R_{emp}(\mathbf{w}) = \frac{1}{K} \sum_{k=1}^{K} \left( d_k - Net(\mathbf{x}_k, \mathbf{w}) \right)^2$$

 Conditional log-likelihood, in case of classification The sum of the negative logarithmic likelihood is minimized

$$C(\mathbf{w}) = -\frac{1}{K} \sum_{k=1}^{K} \left( -logP(\mathbf{y}_k | \mathbf{x}_k, \mathbf{w}) \right)$$

What is the problem with quadratic loss function in classification tasks?

- In case of classification, the convergence can be very slow
- Consider the following very simple case • Check out the example! bias bhttp://neuralnetworksan weight wddeeplearning.com/chap d=0 x=1 3.html **Case A**: Start the learning from w(0)=0.6, b(0)=0.9 cost . Loss function decreases quickly Case A **Case B**: Start the learning from w(0)=2, b(0)=2 • 300 Epoch - Loss function decreases very slowly at the beginning Cost Why is that? Case B - Because the  $\Delta w$  is proportional with the gradient 10/8/2019  $\rightarrow$  Epoch 300

## Calculation of the gradient

• Loss function:

 $L = \frac{1}{2}(d - y)^2$ , where  $y = \sigma(wx + b)$ 

Gradient, using chain rule:

 $\frac{\partial L}{\partial w} = (y - d) \sigma'(wx + b)x = y\sigma'(wx + b)x$ 

- **Case A**: w(0)=0.6, b(0)=0.9, x=1, d=0
  - Slope of the gradient is fine: (wx + b) = 1.5
  - Fast convergence
- **Case B**: w(0)=2, b(0)=2, x=1, d=0

- Slope of the gradient is very small: 
$$(wx + b) = 4$$

Very slow convergence



sigmoid is a large value.

# Introducing Cross Entropy

- Idea: replace the quadratic Loss function with a more appropriate Loss function: Try cross entropy!
- In general:  $C = -\frac{1}{n} \sum_{x} [y \ln a + (1-y) \ln(1-a)]$

• 
$$C(\mathbf{w}) = -\frac{1}{\kappa} \sum_{k=1}^{\kappa} \left( d_k \log P(\mathbf{y}_k | \mathbf{x}_k, \mathbf{w}) + (1 - d_k) \log \left( 1 - P(\mathbf{y}_k | \mathbf{x}_k, \mathbf{w}) \right) \right)$$

- Is it always positive?
  - d<sub>k</sub> is either 0 or 1 (binary classification)
    - $-\,$  Either the first or the second term is zero
  - $P(\mathbf{y}_k | \mathbf{x}_k, \mathbf{w}) = \sigma(\mathbf{w}\mathbf{x}_k + b)$ 
    - The probability is the output of the network
    - $-\,$  Due to the sigmoid, it is between 0 and 1  $\,$
    - Therefore, its logarithm is negative

10/8/2019 http://neuralnetworksanddeeplearning.com/chap3.html

- Is it a good loss function?

Good decision (small loss):

- When d<sub>k</sub> is 0 and  $P(\mathbf{y}_k | \mathbf{x}_k, \mathbf{w})$  is close to 0, than  $-log(1 P(\mathbf{y}_k | \mathbf{x}_k, \mathbf{w})) \sim \mathbf{0}$
- When d<sub>k</sub> is 1 and  $P(\mathbf{y}_k | \mathbf{x}_k, \mathbf{w})$  is close to 1, than  $-log(P(\mathbf{y}_k | \mathbf{x}_k, \mathbf{w})) \sim \mathbf{0}$ Bad decision (large loss):
- When d<sub>k</sub> is 0 and  $P(\mathbf{y}_k | \mathbf{x}_k, \mathbf{w})$  is close to 1, than  $-log(1 P(\mathbf{y}_k | \mathbf{x}_k, \mathbf{w})) \sim \infty$
- When d<sub>k</sub> is 1 and  $P(\mathbf{y}_k | \mathbf{x}_k, \mathbf{w})$  is close to 0, than  $-log(P(\mathbf{y}_k | \mathbf{x}_k, \mathbf{w})) \sim \infty$

## Why is cross entropy good?

$$C(\mathbf{w}) = -\frac{1}{K} \sum_{k=1}^{K} \left( d_k \log P(\mathbf{y}_k | \mathbf{x}_k, \mathbf{w}) + (1 - d_k) \log \left( 1 - P(\mathbf{y}_k | \mathbf{x}_k, \mathbf{w}) \right) \right)$$

– Because its partial derivative does not contain  $\sigma'$ 

$$\frac{\partial C}{\partial w_j} = \frac{1}{K} \sum_{k=1}^{K} x_j (\sigma(\mathbf{wx} + \mathbf{b}) - d)$$

 The gradient is proportional with the value of the sigmoid, and not with its derivative!



## Loss function for softmax: Negative log-likelihood

- $L(\mathbf{y}) = \sum_{k=1}^{K} -\log(\mathbf{y})$
- The negative logarithm of the probability of the correct decision classes are summed up
- It is small, if the confidence of a good decision was high for a certain class
- Large, when the confidence is low
- Partial derivative of a softmax layer with negative log-likelihood:

$$\frac{\partial C}{\partial v_j} = y_j - 1$$



# Probability type loss: Cross Entropy and Softmax

- Mathematically:
   Normalized exponential functions of the units
- Probability distribution of n discrete classes:
   One-of-n classes problems

$$= \frac{e^{v_j}}{\sum_{j=1}^n e^{v_j}}$$
$$v = w^T x$$

ovi

 $y_i = softmax(v)_i$ 

$$0 < y_i < 1$$

- $\sum_{i=1}^{n} y_i = 1$
- Architectural difference:
  - Previously learned activation functions were based on the inputs of one neuron
  - Softmax combines a layer of output neurons

10/22/2019



https://medium.com/deep-learning-demystified/loss-functions-explained-3098e 8ff2b27

## d. Machine learning vs traditional programming

## Conventional approach

- Trivial, or at least analitically solvable tasks
  - Well established mathematical solution exist or at least can be derived
- Example:
  - Finding well defined data constellations in a database
  - Formal verification of the operation is easy

- Machine learning approach
  - Complex underspecified tasks
    - No exact mathematical solution exists, the function to be implemented is not known
- Example:
  - Searching for "strange" data constellations in a database
  - Verification of the operation is difficult

In case of very complex problems, verification of the operation is very difficult. Typically done by exhaustive testing in case of machine learning.

Machine learning



We consider each task as an input-output problem



- 1. Appearance of machine learning methods and frameworks, optimization know-how, new tools for rapid experimentation
- 2. New architectures are available for computation
  - (1980: VIC-20 5kb RAM, MOS 6502 CPU 1.02Mhz)
  - (2018: NVIDIA GeForce GTX 1080, 8GB RAM, 1733 MHz, 2560 cores)
- 3. Vast amount of data is
- available – Billions of labeled images available quasi free



000 (ts an

Typical Machine Learning Types

- Supervised Learning
  - Learning from labeled examples (for which the answer is known)
- Unsupervised Learning
  - Learning from unlabeled examples (for which the answer is unknown)
- Reinforcement Learning

   Learning by trial and feedback, like the "child learning" example

What is intelligence? • The ability to acquire and apply knowledge and skills. Intelligence is the ability to adapt to change "Stephen Hawking" Providing computers the ability to learn without being explicitly programmed:

Involves: programming, Computational statistics, mathematical optimization, image processing, natural language processing etc...



#### - the math of the functionality is known - the known math should be programmed

## At Neural Networks

- the math behind the functionality is unknown - the functionality is "illustrated" with examples

## e. Inception

Inception module:



# Inception

Google, Christian Szegedy

2014 with a top 5 error rate of 6.7%

This can be thought of as a "pooling of features" because we are reducing the depth of the volume, similar to how we reduce the dimensions of height and width with normal maxpooling layers.

Idea:

Not to introduce different size kernels in different layers, but introduce 1x1, 3x3, 5x5 in each layers, and let the Neural Net figure out, what representation is the most useful, and use that!

Parallel multi-scale approach.



# **Rethinking Inception**



# **Rethinking Inception**

Larger (5x5) convolutions were substituted by series of 3x3 convolutions

#### Advantages:

- 1. Reduction of number of parameters,
- 2. Additional non-linearities (RELUs) can be introduced







2D convolution were substituted by two 1D convolutions

AlexNet: 60 million parameters VGGNet :180 million parameters GoogLeNet / Inception-v3: 7 million parameters



33

# **GoogleNet Inception v4**



Inception architecture applied to residual networks





41

## No. 4 Topic

## a. McCulloch-Pitts model

Artificial neuron model, 40's (McCulloch-Pitts, J. von Neumann);

The artificial neuron (McCulloch-Pitts)



- The artificial neuron is an information processing unit that is basic constructing element of an artificial neural network.
- Extracted from the biological model



- Receives input through its synapsis (x<sub>i</sub>)
- Synapsis are weighted (w<sub>i</sub>)
  - if  $w_i > 0$ : amplified input from that source (excitatory input)
  - if  $w_i < 0$ : attenuated input from that source (inhibitory input)
- A b value biases the sum Bias b to enable asymmetric behavior C  $v_{k1}$ Activation A weighted sum is calculated function Activation function shapes the Output Input Σ output signal  $y_k$ Summing  $x_i$ : input vector iunction  $w_{ki}$ : weight coefficient vector of neuron k  $b_k$ : bias value of neuron k Synaptic  $o_k$ : output value of neuron k weights P-ITEEA-0011 2/

## b. Parameters of conv - filter, stride, padding, etc




#### Why use padding:

, , ,	
Simplifies the execution	unpadded f: 3 2 2 -1 -2 -3
code	Code type for boundary 1 🛛 -1 0 2 0 -1
No branches	Code type for boundary 2 -1 0 2 0 -1
Do not have to deal with the different calculation methods at	Code type for central -1 0 2 0 -1
the boundaries	
Same code runs in the	-1 0 2 0 -1
entire array	One code for all the array
	Though it is more multiply-add operation, but as f>>g a branch free simpler code is more efficient

# Parameter number and computational load

- Number of trainable free parameters:
  - k in 1D convolution | size(g)= k
  - $k^2$  in 2D convolution | size(g)=  $k \times k$
- Operation number
  - $k^*n$  for a padded 1D convolution | size(f)= n
  - $k^2 * n^2 = O(n^2)$  for a padded 2D convolution | size(f)=  $n \times n$



# c. Linear classifier, margin of the classifier

# Linear separability



- Today, we assume that the IO sets are linearly separable
- The decision boundary is a hyperplane defined:

 $\mathbf{w}^T \mathbf{x} = \mathbf{0}$ 

- Positive side of the hyperplane is classified: +1 (yes)
- Negative side of the hyperplane is classified : 0 (no).

9/17/2019

P-ITEEA-0011 Lecture 2





#### Maximum Margin:

Define the margin of a linear classifier as the width that the boundary could be increased by before hitting a data point.

## The learning algorithm: Datasets

#### • Training set

Set of input – desired output pairs
Will be used for training

t pairs  $X^+ = \{\mathbf{x} : d = +1\}$  $X^- = \{\mathbf{x} : d = 0\}$ 

Test set

X 2

decision

 $w_0 = 0$ 

12

region for C1

 $W_1X_1$ 

 $+ W_2 X_2 +$ 

- Used, when we have large set of input vectors (not used today)
  Set of input desired output pairs
- Will be used for testing and scoring the result
- We assumed that  $X^+$  and  $X^-$  must be linearly separable

• We are looking for an optimal parameter set:  $X^{+} = \left\{ \mathbf{x} : \mathbf{w}_{opt}^{T} \mathbf{x} > 0 \right\},$  $X^{-} = \left\{ \mathbf{x} : \mathbf{w}_{opt}^{T} \mathbf{x} < 0 \right\}.$ 

# Elementary set separation by a single neuron (2)

decision

P-ITFFA-0011

Lecture 1

- in a 2-D input space, the hyper plane is a straight line.
- Above the line is boundary classified: +1 (C1: yes)
- Below the line is decision classified : 0 (C2: no)region for C2

9/10/2019

## d. Data augmentation

# Regularization and optimization methods 📲

Different methods to increase the loss in the learning phase, but reduce overfitting and increase generalization capabilities

- Local response normalization
- Batch normalization
- Data augmentation (Enriching the data set)
- Early stopping
- Ensemble methods
  - Network duplication
  - Bagging
  - Dropout

Ian Goodfellow: regularization is "any modification we make to the learning algorithm that is intended to reduce the generalization error, but not its training error"

- Data augmentation
- Idea:
  - Increase the generalization capability of the net by enlarging the training set
- Increase the number of the training vector by introducing fake (artificial) input-output pairs
- Typical methods
  - Translating
  - Slight rotation
  - Rescaling
  - Adding noise
  - Flipping
  - Cutting out parts
  - Manipulating with pixel values

Enlarge your Dataset

# Input normalization and Data augmentation I

Images were down-sampled and cropped to 256×256 pixels and normalized

- 1<sup>st</sup> : image translations and horizontal reflections
  - random 224x224 patches + horizontal reflections from the 256x256 images
  - Testing: five 224x224 patches + horizontal reflections → averaging the predictions over the ten patches



1/5/2019

# Data augmentation II

- 2<sup>nd</sup> : change the intensity of RGB channels
  - PCA on the set of RGB pixel values throughout the ImageNet training set
  - To each RGB image pixel  $I_{xy} = [I_{xy}^R, I_{xy}^G, I_{xy}^B]$  following is added



## **U-Net**



Designed for biomedical image processing: cell segmentation
Data augmentation via applying elastic deformations,
Natural since deformation is a common variation of tissue
Smaller dataset is enough

# Alexnet

Alex Krizhevsky, Ilya Sutskever, and Geoffrey Hinton (2012)

Trained whole ImageNet (15 million, 22,000 categories)

Used data augmentation (image translations, horizontal reflections, and patch extractions)

## e. YOLO = you only look once (2016 May)

#### -> Object detection as regression

Classification

Classification + Localization

**Object Detection** 

Instance Segmentation



# YOLO, Detectnet

Models detection as a regression problem:

Divide the image into a grid and each cell can vote for the bounding box position of possible object. (Four output per cell for the corner positions.)

Boxes can have arbitrary sizes

Each cell can proposes a bounding box one category (more layers, more categories per position).

Non-suppression on the boxes

No need for scale search, the image is processed once and objects in different scales can be detected



Resize The Image And bounding boxes to 448 x 448

Redmon, Joseph, et al. "You only look once: Unified, real-time obje

Divide The Image Into a 7 x 7 grid. Assign detection grid cells based on their centers.















confidence scores: reflect how confident is that the box contains an object+how accurate the box is.

conditional class probabilities: conditioned on the grid cell containing an object

$$\Pr(\text{Class}_i | \text{Object}) * \Pr(\text{Object}) * \text{IOU}_{\text{pred}}^{\text{truth}} = \Pr(\text{Class}_i) * \text{IOU}_{\text{pred}}^{\text{truth}}$$

- At test time, multiply the conditional class probabilities and the individual box confidence predictions
- giving class-specific confidence scores for each box
- Showing both the probability of that class appearing in the box and how well the predicted box fits the object



Handles oclusion

## **Pixel level segmentation**



The expected output of the network is not a class, but a map representing the pixels belonging to a certain class.

Creation of a labeled dataset (handmade pixel level mask) is a tedious task

More complex architectures are needed (compared to classification)

Popular architectures (Sharpmask, U-NET ...)



sky tree mod grass water bidg mmtn fg ob). SharpMask: Learning to Refine Object Segments. Pedro O. Pinheiro, Tsung-Yi Lin, Ronan Collobert, Piotr Dollàr (ECCV 2016)



SEMANTIC IMAGE SEGMENTATION WITH DEEP CONVOLUTIONAL NETS AND FULLY CONNECTED CRFS Liang-Chieh Chen et al. ICLR 2015

# No. 5 Topic

- a. Statistical learning theory Statistical learning theory
- Empirical error

$$R_{emp}\left(\mathbf{w}\right) = \frac{1}{K} \sum_{k=1}^{K} \left(d_k - Net\left(\mathbf{x}_k, \mathbf{w}\right)\right)^2$$

• Theoretical error

 $\left\| \mathbf{F}(\mathbf{x}) - \operatorname{Net}(\mathbf{x}, \mathbf{w}) \right\|^{2} = \int \dots \int \left( \mathbf{F}(\mathbf{x}) - \operatorname{Net}(\mathbf{x}, \mathbf{w}) \right)^{2} dx_{1} \dots dx_{N}$ 

- Let us have  $\boldsymbol{x}_k$  random variables subject to uniform distribution

 $\mathbf{x}_{k}$  random variable, where  $d=F(\mathbf{x})$ 

$$\lim_{k \to \infty} = \frac{1}{K} \sum_{k=1}^{K} (d_k - Net(\mathbf{x}_k, \mathbf{w}))^2 = E(d - Net(\mathbf{x}, \mathbf{w}))^2 =$$

$$\int \cdots_{X} \int (F(\mathbf{x}) - Net(\mathbf{x}, \mathbf{w}))^2 p(\mathbf{x}) dx_1 \dots dx_N =$$
Because it is ~ constant due to the uniformity
$$\frac{1}{|X|} \int \cdots_{X} \int (F(\mathbf{x}) - Net(\mathbf{x}, \mathbf{w}))^2 dx_1 \dots dx_N \square$$

$$\int \cdots_{X} \int (F(\mathbf{x}) - Net(\mathbf{x}, \mathbf{w}))^2 dx_1 \dots dx_N$$

# Therefore

$$\lim_{K \to \infty} \mathbf{w}_{\text{opt}} = \mathbf{w}_{\text{opt}}^{(K)}$$

# Where l.i.m. means: lim in mean

$$\lim_{K \to \infty} R_{emp} \left( \mathbf{w} \right) = R_{th} \left( \mathbf{w} \right)$$
$$\lim_{K \to \infty} \frac{1}{K} \sum_{k=1}^{K} \left( d_k - Net \left( \mathbf{x}_k, \mathbf{w} \right) \right)^2 = \int \dots \int \left( F(\mathbf{x}) - Net \left( \mathbf{x}, \mathbf{w} \right) \right)^2 dx_1 \dots dx_N$$

#### Weak learning is satifactory!

Ez mar nem is ide tartozik, de mindu: Learning – in practice

• Learning based on the training set:

$$\tau^{(K)} = \{ (\mathbf{x}_k, d_k); k = 1, ..., K \}$$

- Minimize the empirical error function  $(R_{emp})$  $\mathbf{w}_{opt}^{(K)} : \min_{\mathbf{w}} \frac{1}{K} \sum_{k=1}^{K} (\underline{d_k - Net(\mathbf{x}_k, \mathbf{w})})^2 = \min_{\mathbf{w}} R_{emp}(\mathbf{w})$
- Learning is a multivariate optimization task

# b. Various activation functions and their properties

- Activation function shapes the output signal
- <u>https://missinglink.ai/guides/neural-network-concepts/7-types-neural-network-activati</u> on-functions-right/

# Activation function: $\phi(.)$

- Always a nonlinear function
- Typically it clamps the output (introduces boundaries)
- Monotonic increasing function
- Differentiable
  - Important from theoretical point of view
- Or at least continuous (except in simplified cases)
  - Sophisticated training algorithms require continuity

## Activation function

- Originally step function
- Sigmoid or Tanh or their piece-wise linear approximation is used nowadays
- Sophisticated training algorithms require differentiable or at least continuous functions

Sigmoid (or logistic) function is a widely activation function





Derivative of sigmoid function I  
Sigmoid(x) = 
$$\frac{1}{1 + e^{-x}}$$
  $\frac{d}{dx}S(x) = \frac{d}{dx}\frac{1}{1 + e^{-x}}$   
quotient rule:  
 $\frac{d}{dx}f = \frac{(denominator * \frac{d}{dx}numerator) - (numerator * \frac{d}{dx}denominator)}{denominator^2}$   
 $\frac{d}{dx}S(x) = \frac{(1 + e^{-x})(0) - (1)(-e^{-x})}{(1 + e^{-x})^2}$   
 $\frac{d}{dx}S(x) = \frac{(1 + e^{-x})(0) - (1)(-e^{-x})}{(1 + e^{-x})^2}$   
 $\frac{d}{dx}S(x) = \frac{e^{-x}}{(1 + e^{-x})^2}$  This is the correct result, but it is not in a nice form.  
Derivative of sigmoid function II  
 $\frac{d}{dx}S(x) = \frac{e^{-x}}{(1 + e^{-x})^2}$   $\longrightarrow$   $\frac{d}{dx}S(x) = \frac{1 - 1 + e^{-x}}{(1 + e^{-x})^2}$   
 $\frac{d}{dx}S(x) = \frac{1 + e^{-x}}{(1 + e^{-x})^2} - \frac{1}{(1 + e^{-x})^2}$  reduction  
 $\frac{d}{dx}S(x) = \frac{1}{(1 + e^{-x})} - \frac{1}{(1 + e^{-x})^2}$  Multiply out  
 $\frac{d}{dx}S(x) = \frac{1}{(1 + e^{-x})}(1 - \frac{1}{1 + e^{-x}})$  Sigmoid(x) =  $\frac{1}{1 + e^{-x}}$   
 $\frac{d}{dx}S(x) = S(x)(1 - S(x))$  Much nicer form!  
• Bipolar activation function:  
 $tanh$   
• Continuously differentiable

f(x) 0

-0.5

-1 -3

-2

-1

 $\varphi$ 

1

-1

0

0

1

2

3

- Monotonic
- Useful, when bipolar output • is expected
- Hard approximations:
  - Piece-wise
  - Step-wise

# Activation function II

- Hyperbolic tangent function
  - Continuous
  - Continuously differentiable
  - It is used in the output layer of the fully connected neural network

 $\varphi(x) = \tanh(x)$ 

$$\frac{d}{dx}\varphi(x) = 1 - \tanh^2(x) = (1 - \tanh(x))(1 + \tanh(x))$$

æ



+ other ReLU in No.2 Topic

# Probabilistic decision (n discrete categories)

00

11

22

5

66 77

88

2

 $x_i$ 

- Assume we have annotated input vectors with *n* different classes (MNIST data base)
- Expect a probability distribution on the output layer!
  - $0 \le y_i \le 1$  sigmoid OK! \_

- 
$$\sum_{i=1}^{n} y_i = 1$$
 sigmoid NOT OK! **?**



# Softmax

- Mathematically:
  - Normalized exponential functions of the output units
- Probability distribution of n discrete classes: • - One-of-n classes problems

$$- 0 \le y_i \le 1$$

- $\sum_{i=1}^{n} y_i = 1$ \_
- Architectural difference: .
  - Previously learned activation functions were based on the inputs of one neuron
  - Softmax combines a layer of output neurons

10/8/2019

# **Properties of Softmax**

Generalization of sigmoid function for one-of-n class

Squashes a vector of size n between 0 and 1 Improves the interpretability of the output of a Neural Net

Describes the probability distribution of a certain class

- We may use the word "confidence"
- Winner take all
- exponential function strongly penalize the nonwinners
- Similar to lateral negative feedback in the natural neural systems

$$y_i = softmax(v)_i$$

$$=\frac{e^{v_i}}{\nabla v_i}$$

$$\sum_{j=1}^{n} e^{v_j}$$
$$v = w^T x$$





- $L(\mathbf{y}) = \sum_{k=1}^{K} -\log(y)$
- The negative logarithm of the probability of the correct decision classes are summed up
- It is small, if the confidence of a good decision was high for a certain class
- Large, when the confidence is low
- Partial derivative of a softmax layer with negative log-likelihood:

$$\frac{\partial C}{\partial v_i} = y_j - 1$$



#### c. Autoencoders

- Neural network used for efficient data coding
- Uses the same vector for the input and the output
  - No labelled data set is needed
  - Unsupervised learning
- Two parts
  - Encoder: reduces data dimension
  - Decoder: reconstructs data
  - Middle layer: code

11/19/2019







## d. Graph unrolling and parameter sharing in recurrent NN

Még nézd át az RNN-t külön, mert ez homály





- Parameter sharing
  - Same parameters everywhere in the layer
  - Contribution to the gradient of a weight from many positions
  - Reduces the risk of overfitting
  - Reduces the risk of dying RELU (dying cell)
    - When it happens, an entire feature extractor on a layer is dying



## **EZ STACKOVERFLOW:**

# **Parameter Sharing**

Being able to efficiently process sequences of varying length is not the only advantage of parameter sharing. As you said, you can achieve that with padding. The main purpose of parameter sharing is a reduction of the parameters that the model has to learn. This is the whole purpose of using a RNN.

If you would learn a different network for each time step and feed the output of the first model to the second etc. you would end up with a regular feed-forward network. For a number of 20 time steps, you would have 20 models to learn. In Convolutional Nets, parameters are shared by the Convolutional Filters because when we can assume that there are similar interesting patterns in different regions of the picture (for example a simple edge). This drastically reduces the number of parameters we have to learn. Analogously, in sequence learning we can often assume that there are similar patterns at different time steps. Compare 'Yesterday I ate an apple' and 'I ate an apple yesterday'. These two sentences mean the same, but the 'I ate an apple' part occurs on different time steps. By sharing parameters, you only have to learn what that part means once. Otherwise, you'd have to learn it for every time step, where it could occur in your model.

There is a drawback to sharing the parameters. Because our model applies the same transformation to the input at every time step, it now has to learn a transformation that makes sense for all time steps. So, it has to remember, what word came in which time step, i.e. 'chocolate milk' should not lead to the same hidden and memory state as 'milk chocolate'. But this drawback is small compared to using a large feed-forward network.

## e. MobileNet

# MobileNet

In this arhcitecture feature depths are squeezed before each operation

In a squeezed architecture we will use downscale the 128 feature maps to 16, using a linear combination (1x1 convolution)

After the 3x3 covolutions, we expanded back to 128 layers by 1x1 convolution again

From the linear combination of these elements the new maps are created



Scaling in feature map depths.

## No. 6 Topic

## a. Machine learning problem definition

Nem vágom ez alatt mit ért...

# Machine learning



We consider each task as an input-output problem



Machine learning is an application of artificial intelligence (AI) that provides systems the ability to automatically learn and improve from experience without being explicitly programmed. **Machine learning focuses on the development of computer programs** that can access data and use it learn for themselves.

The process of learning begins with observations or data, such as examples, direct experience, or instruction, in order to look for patterns in data and make better decisions in the future based on the examples that we provide. **The primary aim is to allow the computers learn automatically** without human intervention or assistance and adjust actions accordingly.

## **Problem Definition Framework**

I use a simple framework when defining a new problem to address with machine learning. The framework helps me to quickly understand the elements and motivation for the problem and whether machine learning is suitable or not.

The framework involves answering three questions to varying degrees of thoroughness:

- Step 1: What is the problem? Describe the problem informally and formally and list assumptions and similar problems.
- Step 2: Why does the problem need to be solve? List your motivation for solving the problem, the benefits a solution provides and how the solution will be used.
- Step 3: How would I solve the problem? Describe how the problem would be solved manually to flush domain knowledge.

#### b. Newton optimizer -> See: No. 3 Topic

# c. Effects and relationship of model capacity and complexity - overfitting, underfitting

Do we have to reach the global minimum? **§** 

- Not really
- Global minimum means:
   Overfitting
- Overfitting: The network exactly learned the training vectors
- However, it loses the generalization capabilities





**Overfitting** occurs when a model with high **capacity** fits the noise in the data instead of the (assumed) underlying relationship

# Losing the generalization capabilities!!!

# Network complexity vs. capacity

 In general, the more layers we have, and the more neurons there are, the larger the capacity.



- There is no adequate method to predict the required complexity.
- Even if a network is capable to learn a task, it is not guaranteed that it will.



## How could we create deeper networks?

A deeper network always have the potential to perform better, but training becomes difficult

How could we ensure that additional layers will not decrease accuracy (might even increase it)?



## **Overfitting in Machine Learning**

Overfitting refers to a model that models the training data too well.

Overfitting happens when a model learns the detail and noise in the training data to the extent that it negatively impacts the performance of the model on new data. This means that the noise or random fluctuations in the training data is picked up and learned as concepts by the model. The problem is that these concepts do not apply to new data and negatively impact the models ability to generalize.

Overfitting is more likely with nonparametric and nonlinear models that have more flexibility when learning a target function. As such, many nonparametric machine learning algorithms also include parameters or techniques to limit and constrain how much detail the model learns.

For example, decision trees are a nonparametric machine learning algorithm that is very flexible and is subject to overfitting training data. This problem can be addressed by pruning a tree after it has learned in order to remove some of the detail it has picked up.

## **Underfitting in Machine Learning**

Underfitting refers to a model that can neither model the training data nor generalize to new data.

An underfit machine learning model is not a suitable model and will be obvious as it will have poor performance on the training data.

Underfitting is often not discussed as it is easy to detect given a good performance metric. The remedy is to move on and try alternate machine learning algorithms. Nevertheless, it does provide a good contrast to the problem of overfitting.

#### A Good Fit in Machine Learning

Ideally, you want to select a model at the sweet spot between underfitting and overfitting.

This is the goal, but is very difficult to do in practice.

To understand this goal, we can look at the performance of a machine learning algorithm over time as it is learning a training data. We can plot both the skill on the training data and the skill on a test dataset we have held back from the training process.

Over time, as the algorithm learns, the error for the model on the training data goes down and so does the error on the test dataset. If we train for too long, the performance on the training dataset may continue to decrease because the model is overfitting and learning the irrelevant detail and noise in the training dataset. At the same time the error for the test set starts to rise again as the model's ability to generalize decreases.

The sweet spot is the point just before the error on the test dataset starts to increase where the model has good skill on both the training dataset and the unseen test dataset.

You can perform this experiment with your favorite machine learning algorithms. This is often not useful technique in practice, because by choosing the stopping point for training using the skill on the test dataset it means that the testset is no longer "unseen" or a standalone objective measure. Some knowledge (a lot of useful knowledge) about that data has leaked into the training procedure.

There are two additional techniques you can use to help find the sweet spot in practice: resampling methods and a validation dataset.

#### How To Limit Overfitting

Both overfitting and underfitting can lead to poor model performance. But by far the most common problem in applied machine learning is overfitting.

Overfitting is such a problem because the evaluation of machine learning algorithms on training data is different from the evaluation we actually care the most about, namely how well the algorithm performs on unseen data.

There are two important techniques that you can use when evaluating machine learning algorithms to limit overfitting:

- 1. Use a resampling technique to estimate model accuracy.
- 2. Hold back a validation dataset.

The most popular resampling technique is k-fold cross validation. It allows you to train and test your model k-times on different subsets of training data and build up an estimate of the performance of a machine learning model on unseen data.

A validation dataset is simply a subset of your training data that you hold back from your machine learning algorithms until the very end of your project. After you have selected and tuned your machine learning algorithms on your training dataset you can evaluate the learned models on the validation dataset to get a final objective idea of how the models might perform on unseen data.

Using cross validation is a gold standard in applied machine learning for estimating model accuracy on unseen data. If you have the data, using a validation dataset is also an excellent practice.

# d. t-distributed Stochastic Neighbor Embedding

Unsupervised learning technique:

# t-Distributed Stochastic Neighbor Embedding (t-SNE)

- Introduced by Laurens Van Der Maaten (2008)
- Generates a low dimensional representation of the high dimensional data set iteratively
- · Aims to minimize the divergence between two distributions
  - Pairwise similarity of the points in the higher-dimensional space
  - Pairwise similarity of the points in the lower-dimensional space
- Output: original points mapped to a 2D or a 3D data space
  - similar objects are modeled by nearby points and
  - dissimilar objects are modeled by distant points with high probability
- Unlike PCA, it is stochastic (probabilistic)

## Step 1: Generate the points in the low dimensional data set (2D or 3D)

- random initialization
- First two or three components of PCA



# Step 2: Calculate the pair-wise similarities measures between data pairs (probability measure)







Exponential normalization of the Euclidian distances are needed due to the high dimensionality. (Curse of dimensionality)



$$q_{ij} = rac{(1+||y_i-y_j||^2)^{-1}}{\sum_{k
eq l}(1+||y_k-y_l||^2)^{-1}}$$

The similarity of datapoint  $x_j$  to datapoint  $x_i$  means the conditional probability  $p_{ji}$  that  $x_i$  would pick  $x_j$ as its nearest neighbor.

#### Step 3: Define the cost function

Similarity of data points in High dimension:

$$p_{ij} = \frac{\exp(-||x_i - x_j||^2/2\sigma^2)}{\sum_{k \neq l} \exp(-||x_l - x_k||^2/2\sigma^2)}$$

- Similarity of data points in Low dimension:
- $q_{ij} = \frac{(1+||y_i y_j||^2)^{-1}}{\sum_{k \neq l} (1+||y_k y_l||^2)^{-1}}$ Cost function (called Kullback-Leiber divergence between the two distributions): Dii ----

$$C = KL(P||Q) = \sum_{i} \sum_{j} p_{ij} \log \frac{r_{ij}}{q_{ij}}$$

- Large  $p_{ii}$  modeled by small  $q_{ii} \rightarrow$ Large penalty
- Large  $p_{ji}$  modeled by large  $q_{ji} \rightarrow Small penalty$
- Local similarities are preserved

## Step 4: Minimize the cost function using gradient descent

Gradient has a surprisingly simple form:

$$\frac{\partial C}{\partial y_i} = 4 \sum_{j \neq i} (p_{ij} - q_{ij})(1 + ||y_i - y_j||^2)^{-1}(y_i - y_j)$$

Optimization can be done using momentum method

# **Physical analogy**



- Our map points are all connected with springs in the low dimensional data map
- Stiffness of the springs depends on  $p_{i|i} q_{i|i}$
- Let the system evolve according to the laws of physics
  - If two map points are far apart while the data points are close, they are attracted together
  - If they are nearby while the data points are dissimilar, they are repelled.
- Illustration (live)
  - https://www.oreilly.com/learning/an-illustrated-introduction-tothe-t-sne-algorithm





## e. ShuffleNet

# ShuffleNet





#### Neten találtam:

We introduce an extremely computation-efficient CNN architecture named ShuffleNet, which is designed specially for mobile devices with very limited computing power (e.g., 10-150 MFLOPs). The new architecture utilizes two new operations, pointwise group convolution and channel shuffle, to greatly reduce computation cost while maintaining accuracy. Experiments on ImageNet classification and MS COCO object detection demonstrate the superior performance of ShuffleNet over other structures, e.g. lower top-1 error (absolute 7.8%) than recent MobileNet on ImageNet classification task, under the computation budget of 40 MFLOPs. On an ARM-based mobile device, ShuffleNet achieves ~13x actual speedup over AlexNet while maintaining comparable accuracy.

#### Maga a cikk:

#### https://zpascal.net/cvpr2018/Zhang\_ShuffleNet\_An\_Extremely\_CVPR\_2018\_paper.pdf

...To overcome the side effects brought by group convolutions, we come up with a novel channel **shuffle operation** to help the information flowing across feature channels. Based on the two techniques, we build a highly efficient architecture called **ShuffleNet**. Compared with popular structures like [31, 9, 41], for a given computation complexity budget, our ShuffleNet allows **more feature map channels**, which helps to **encode more information** and is especially critical to the performance of very small networks. We evaluate our models on the challenging ImageNet classification [4, 30] and MS COCO object detection [24] tasks. A series of controlled experiments shows the **effectiveness** of our design principles and the **better performance over other structures.** Compared with the state-of-the-art architecture MobileNet [12], ShuffleNet achieves **superior performance by a significant margin**, e.g. absolute 7.8% lower ImageNet top-1 error at level of 40 MFLOPs. We also examine the speedup on real hardware, i.e. an off-the-shelf ARM-based computing core. The ShuffleNet model achieves ~13× actual speedup over AlexNet [22] while **maintaining comparable accuracy.** 

## No. 7 Topic

# a. Credit approval problem -> na hogy itt mire gondolt a költő???? Neten találtam:

The last decade has seen an important rise of data gathering, especially in the financial sectors. Banks are indeed one of the biggest producers of big data, as a matter of fact no other company than the bank has so much data gathered on its customers. Gathering and analyzing this data is a key feature for decision making, particularly in banking sector. One of the most important and frequent decision banks has to make, is loan approval. The challenge is to know how to build a proactive, powerful, responsible and ethical exploitation of personal data, to make loan applicant proposals more relevant and personalized. Machine learning is a promising solution to deal with this problem. Therefor, in the last years, many algorithms based on machine learning have been proposed to solve loan approval issue. However, these algorithms have not taken into consideration model to deal with loan approval. Our proposed model is based on a deep neural network, and it permits to classify loan applicant as good or bad risk. Experimental results prove that our proposed Real-Time model, based on deep neural network, outperforms typical binary classifiers, in terms of precision recall and accuracy.

#### Indiai bácsi a youtube-on:

https://www.youtube.com/watch?v=HNj8dzw3H\_E

1. It is used in Distributed Systems

2. This can be divided into Temporal Credit Assignment Problem (Credit or blame to Outcome of internal Decisions) and Structural Credit Assignment Problem (Credit or blame to actions of internal decisions).

3. By these two, we can train the learning machine easily

## b. Objective functions in neural networks

# **Optimization**

- Given an Objective function to optimize
  - Also called: Error function, Cost function, Loss function, Criterion
  - Derived from the network topology and the input/output pairs
- Function types:
  - Quadratic, in case of regression (stochastic process)

$$R_{emp}\left(\mathbf{w}\right) = \frac{1}{K} \sum_{k=1}^{K} \left(d_{k} - Net\left(\mathbf{x}_{k}, \mathbf{w}\right)\right)^{2}$$

- · Conditional log-likelihood, in case of classification (classification process)
  - The sum of the negative logarithmic likelihood (probability) is minimized

$$\Theta(\mathbf{w}) = \sum_{k=1}^{n} -logP(\mathbf{y}_{\mathbf{k}}|\mathbf{x}_{\mathbf{k}};\mathbf{w})$$

9/30/2019

# **Optimizations**

- Here we always minimize the objective function
  - Parametric equation
    - x are the variables
    - w are the parameters
- Optimization targets to find the optimal weights

 $\mathbf{w}_{opt} = min f(\mathbf{x}, \mathbf{d}, Net(\mathbf{x}, \mathbf{w}))$ 

goals:

- Acceptable error level
- Acceptable computational time assuming reasonable computational effort

#### As we discussed ...

- Stocastic Gradient Descent (SGD) Method
  - Uses a random subset of the training vectors (mini batches)
    - One update is fast to calcualte
    - The objective function changes stocastically with the minibatch selection
      - More fluctuation in the objective function than in case of Gradient Descent
  - It helps to come out from local minima and saddles
     Decreases the learning rate during the training time to reduce overshoot
  - Still very slow! (Many update steps are needed)
- More advanced optimization methods required!



## ez talán már nem is ide tartozik:

## **Ensemble methods**

- Idea of ensemble methods:
  - Generate multiple copies of your net
    - Same or slightly modified architectures
  - Train them separately
    - Using different subsets of the training sets
    - Different objective functions
    - Different optimization methods
  - The different trained models have independent error characteristics
  - Averaging the results will lead to smaller error
- Requires more computation and memory both in training and inferencing (testing) phase

## c. Nesterov momentum optimizer

# Momentum I

- Introduced in 1964
- Physical analogy
- The idea is to simulate a unity weight mass
- It flows through on the surface of the error function
- Follows Newton's laws of dynamics
- Having *v* velocity
- Momentum correctly traverses the canyon lengthwise, while gradient steps waste time moving back and forth across the narrow axis of the canyon.



# Momentum II: velocity considerations

N/P

The update rule is given by:

$$\boldsymbol{v} \leftarrow \alpha \boldsymbol{v} - \epsilon \nabla_{\boldsymbol{\theta}} \left( \frac{1}{m} \sum_{i=1}^{m} L(\boldsymbol{f}(\boldsymbol{x}^{(i)}; \boldsymbol{\theta}), \boldsymbol{y}^{(i)}) \right),$$
$$\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} + \boldsymbol{v}$$

The velocity  $\boldsymbol{v}$  accumulates the gradient elements  $\nabla_{\boldsymbol{\theta}} \left( \frac{1}{m} \sum_{i=1}^{m} L(\boldsymbol{f}(\boldsymbol{x}^{(i)}; \boldsymbol{\theta}), \boldsymbol{y}^{(i)}) \right)$ . The larger  $\alpha$  is relative to  $\epsilon$ , the more previous gradients affect the current direction.

#### Terminal velocity is applied when it finds descending gradient permanently:



Require: Learning rate  $\boldsymbol{\epsilon}$ , momentum parameter  $\boldsymbol{\alpha}$ . Require: Initial parameter  $\boldsymbol{\theta}$ , initial velocity  $\boldsymbol{v}$ . while stopping criterion not met do Sample a minibatch of m examples from the training set  $\{\boldsymbol{x}^{(1)}, \ldots, \boldsymbol{x}^{(m)}\}$  with corresponding targets  $\boldsymbol{y}^{(i)}$ . Compute gradient estimate:  $\boldsymbol{g} \leftarrow \frac{1}{m} \nabla_{\boldsymbol{\theta}} \sum_{i} L(f(\boldsymbol{x}^{(i)}; \boldsymbol{\theta}), \boldsymbol{y}^{(i)})$ Compute velocity update:  $\boldsymbol{v} \leftarrow \alpha \boldsymbol{v} - \epsilon \boldsymbol{g}$ Apply update:  $\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} + \boldsymbol{v}$ end while

#### Nesterov Momentum update:

- It calculates the gradient not in the current point, but in the next point, and correct the velocity with the gradient over there (look ahead function)
- It does not runs through a minimum, because if there is a hill behind a minimum, than it starts decreasing the speed in time.



9/30/2019

d. Decomposition of convolutional kernels



- Decomposition is not exact in most cases
  - In general case, it approximates the kernels with a limited accuracy only
- Neural nets does not sensitive for inaccurate decomposition
- Decomposition of larger kernels leads to higher savings!
- Wildly used!

## e. Alexnet + ILSVRC

Alexnet: for image classification



 Won the ImageNet Large Scale Visual Recognition (ILVSRC) Challenge in 2012 (ILVSRC2012)





Alex Krizhevsky, Ilya Sutskever, and Geoffrey E. Hinton, "Imagenet classification with deep convolutional neural networks", Advances in neural information processing systems, 2012

16

# ImageNet Large Scale Visual Recognition Challenge I



- ImageNet:
  - 15+ million labeled high-resolution images
  - 22000 categories
- ILSVRC uses a subset of ImageNet:
  - 1000 categories
  - ~1000 images per category
  - 1.2 million training images | 50 000 validation images | 150 000 testing images

# ImageNet Large Scale Visual Recognition Challenge II

- Each image should be classified
  - Probability distribution
- Top 1 error rate:
  - What percentage was wrongly classified as highest probability? (38,9%)
- Top 5 error rate:
  - What percentage was not in the first five? (18.9%)

11/5/2019





# Alexnet



Alex Krizhevsky, Ilya Sutskever, and Geoffrey Hinton (2012)

Trained whole ImageNet (15 million, 22,000 categories)

Used data augmentation (image translations, horizontal reflections, and patch extractions)

Used ReLU for the nonlinearity functions (Decreased training time compared to tanh) - Trained on two GTX 580 GPUs for six days

#### **Dropout layers**

2012 marked the first year where a CNN was used to achieve a top 5 test error rate of 15.4% (next best entry was with error of 26.2%)



## No. 8 Topic

a. Delta learning rule

Sequential back propagation

Adapting the weights of the FFNN (recursive algorithm)

$$w_{ij}^{(l)}(k+1) = w_{ij}^{(l)}(k) + \Delta w_{ij}^{(l)}(k)$$
$$\Delta w_{ij}^{(l)}(k) = ?$$

• The weights are modified towards the differential of the error function (delta rule):  $\partial R$ 

$$\Delta w_{ij}^{(l)} = -\eta \, rac{\partial R_{emp}}{\partial w_{ij}^{(l)}}$$

• The elements of the training set adapted by the FFNN sequentially

$$R_{emp} = R_{emp}(y(\mathbf{x}), d)$$



# b. Batch normalization

- In very deep networks the distribution of the input vectors changes from layer to layer
  - The first layer got normalized input
  - The second layer somewhat shifts and twists on this normalization
  - And it goes on, and the (originally normalized) data propagating trough the layers will be lose its normalized properties (called *"covariance shift"*)
  - This will shift the neuron out of its zero centered position, where the activation function performs well (where the nonlinearity is)
- Solution: normalization on each layers!
- It also introduce a noise (loss function increase), which helps to avoid local minima and avoids overfitting

# **Batch Normalization**

- Done on layer level like softmax
- Training:
  - Done on minibatch level
- Inferencing:
  - Do the normalization with the precalculated parameters of the entire training set
- Batch normalization is differenciable via chain rule
  - Back propagation can be applied for batch normalized layers
- Rewriting the normalization using probability terms:

$$\hat{x}^{(k)} = \frac{x^{(k)} - \mathbf{E}[x^{(k)}]}{\sqrt{Var[x^{(k)}]}}$$

E: the expectation Var: the variance





# Local response norm. vs batch norm.

Both work within one convolutional layer

- Normalization either through the feature maps or within one feature map
- Normalization is done for one input image



- Normalization done for all the pixels in all the feature maps within a layer
   Normalization is done for
- Normalization is done for the entire batch



c. Transposed conv., atrous conv. Upsampling II: "transpose convolution"





67	



- Semantic Image Segmentation
  - U-Net
  - DeConvNet
  - SegNet
    - Resolution controlling
    - Atrous convolutions, sub-pixel image combination



Atrous convolution goes deeper without further reducing resolution



stone!

# Visualizing atrous convolution

11/12/2019

- d. Object classification, localization VS object detection, semantic segmentation VS instance segmentation -> ez kicsit megfoghatatlan a diasor alapján...
- How to do image classification

   Alexnet
   One decision per image (*classification*)

   Detection and Localization is more complex

   Multiple (few) decision per image
   <u>Regressions</u> for localization
   <u>Classification</u> for detection

   Pixel level Segmentation

   Very high number of decisions (*classification*)
   Very high number of decisions (*classification*)
   Very high number of decisions (*classification*)
   Very high number of decisions (*classification*)

# Object detection/localization and classification

- Chicken and egg problem
  - You need to know that it is a bicycle before able to say that both a wheel part and a pipe segment belongs to the same object
  - You need to know that the red box contains an object before you can recognize it. (Cannot recognize a bicycle if you try it from separated parts)
- Our brain does it parallel
  - How neural nets can solve it?
  - Detection by regression?
  - Detection by classification?

## **Neural networks for regression**

Multiple object detection on a single image

Classification is good for a single object (can be extended for k objects – top k candidates)

How could we detect objects in general, when the number of objects is unknow

Classification

Classification + Localization

Object Detection









# **Object classification:**



## Approximation



- When solving engineering task by FFNN we are faced with the following theoretical questions:
- 1. Representation
- What kind of functions can be Approximated by an FFNN?
   Learning
- How to set up the weights to solve a specific task?
- 3. Generalization
  - If only limited knowledge is available about the task which is to be solved, then how the FFNN is going to generalize this knowledge?

#### fullos link:

https://medium.com/analytics-vidhya/image-classification-vs-object-detection-vs-image-segm entation-f36db85fe81

# **Object detection:**

One or Multiple object per image

- Task:
  - Find the objects
  - Identify them with bounding boxes

Area or pixel level one-of-two decision!



Annotated image database:

- Detection (squared objects)
- Segmentation (segmented objects)

## Segmentations:

# Segmentation

- Semantic Segmentation
  - Label each pixel in the image with a category label
  - Don't differentiate Instances, only care about pixels
     Pixel level one-of-n classification!
- Semantic Instance Segmentation
  - Differentiate instances

11/5/2019





Ing

Input Image Semantic Segmentation

Semantic Instance Segmentation

# Semantic Segmentation Idea I: Sliding Window



it pixel-wise! No reuse of shared features between overlapping Patches.

Farabet et al, "Learning Hierarchical Features for Scene Labeling," TPAMI 2013 Pinheiro and Collobert, "Recurrent Convolutional Neural Networks for Scene Labeling",

# Semantic Segmentation Idea II: Fully Convolutional



# Semantic Segmentation Idea III: Fully Convolutional

**Downsampling**: Pooling, strided convolution Design network as a bunch of convolutional layers, with downsampling and upsampling inside the network!

 convolution

 H/8 x W/8

 H/8 x W/8

 H/4 x W/4

 H/16 x W/16

 H/32 x W/32

 tupsampling



**Upsampling:** 

???

61

Predictions: H x W

# Increasing spatial resolution in segmentation I

nonlinearity



Higher resolution layers directly forwarded to transfer finer spatial information Called "Skipping". It skips using the coarser (more downsampled) layers Can be considered of an ensembling of three networks




The model name, ResNeXt, contains Next. It means the next dimension, on top of the <u>ResNet</u>. This next dimension is called the "cardinality" dimension. And ResNeXt becomes the 1st Runner Up of ILSVRC classification task.

#### No. 9 Topic

a. ADAM optimizer

#### ADAM algorithm (2014)

- The name "Adam" derives from the phrase "adaptive moments."
- In the context of the earlier algorithms, it is perhaps best seen as a variant on the combination of RMSProp and momentum with a few important distinctions.
- in Adam, momentum is incorporated directly as an estimate of the first order moment (with exponential weighting) of the gradient.
- Adam includes bias corrections to the estimates of both the firstorder moments (the momentum term) and the (uncentered) second-order moments to account for their initialization at the origin



#### b. The softmax function -> lasd in No. 5 Topic

#### c. R-CNN architectures: R-CNN, Fast R-CNN, Faster R-CNN

#### Detection and Localization

- PASCAL Database and Competion
- R-CNN
  - Region proposal, Classification
    - Support Vector Machine (SVN), Bounding box refinement
- Fast R-CNN
- Faster R-CNN

#### **R-CNN in a Glance**

#### **R-CNN:** Regions with CNN features



- 1. Input image
- 2. Region proposals
- 3. Compute CNN features with warped images
- 4. Classification with Support Vector Machine (SVM)
- 5. Ranking/selecting/merging  $\rightarrow$  detections
- 6. Bounding box regression

# **R-CNN: Region Proposal**

- Requirements:
  - Propose a large number ( up to 2000) of regions (boxes) with different sizes
  - Still much better than exhausting search with multi-scale sliding window (brute force)
  - Boxes should contain all the candidate objects with high probability
- R-CNN works with various Region proposal methods:
  - Objectness
  - Constrained Parametric Min-Cuts for Automatic Object Segmentation
  - Category Independent Object Proposals
  - Randomized Prim
  - <u>Selective Search</u>
- Selective Search is the fastest and provides best regions

# **R-CNN: Selective Search I**

- Graph based segmentation (Felzenszwalb and Huttenlocher method)
  - cannot be used in this form, because one object is covered with multiple segments, moreover regions for occluded objects will not be covered
- Idea: oversegment it and apply scaled similarity based merging



Input image



Segmented image



Oversegmented image

#### Step-by-step merging regions at multiple scales based on similarities Convert regions to boxes Step n merging Original fine scale Step one merging Similarity measures I **Texture Similarity Color Similarity** Generate color histogram of each Texture features: Gaussian segment (descriptor) derivatives at 8 orientations in - 25 bins/ color channels each pixel - Descriptor vector $(c_i^k)$ size: 3x25=75 10 bins/color channels Descriptor vector $(t_i^k)$ size: Calculate histogram similarity for 3x10x8=240 each region pair $s_{color}(r_i, r_j) = \sum \min(c_i^k, c_j^k)$ Each region will have a texture histogram ··· Histogram Calculate histogram similarity for c<sup>k</sup> is the Intersection: 0.66 each region pair 240 similarity histogram value $s_{texture}(r_i, r_j) = \sum_{k=1}^{k} \min(t_i^k, t_j^k)$ for the kth bin in color descriptor t<sup>k</sup><sub>i</sub> is the histogram value for the k<sup>th</sup> bin in texture descriptor Shape Similarity Size Similarity Measures how well Helps merging the smaller sized • two regions are fit objects How close they similar are Since we do bottom up merging, How large is the the small segments will be overlap More merged first, because their size similar similarity score is higher $s_{fill}(r_i, r_j) =$ $= 1 - \frac{size(BB_{ij}) - size(r_i) - size(r_j)}{size(image)}$ $s_{size}(r_i, r_j) = 1 - \frac{size(r_i) + size(r_j)}{size(image)}$ $size(BB_{ij})$ is the size of the bounding box of size(image) is the size of the entire image in $r_i$ and $r_i$ pixels

**R-CNN: Selective Search II** 

#### **Final Similarity**

- Linear combination of the four similarities
- $s_{final}(r_i, r_j) = a_1 s_{color}(r_i, r_j) + a_2 s_{texture}(r_i, r_j) + a_3 s_{shap}(r_i, r_j) + a_4 s_{fill}(r_i, r_j)$

#### List or proposed region

- 1. Initial oversegmentation
- 2. Calculation the similarities
- 3. Merge the similar regions
- 4. The formed regions are added to the region list (this ensures that there will be smaller and larger regions in the list as well)
- 5. Goto 2



# **Proposed regions**

- Few hundreds or few thousand boxes
- Includes all the objects with high probability
- Number of the boxes are much smaller than with brute force method
- C and python functions
   exist



# Computing the features of the regions

after the otherResize (warp) the

•

- regions to the input size of the ConvNet
- Calculate the features of the individual regions



#### **Convolution network**

- Pre-trained AlexNet, later VGGNet
- The decision maker SoftMax layer was cut





Similar to single layer perceptron, but optimized • for maximum margin

#### Why SVM?

- Why not use simple the classification output of the AlexNet?
- During the training, the AlexNet/VGGNet is not trained •
- Only SVM is trained •

•

.

vectors

· The number of category is much smaller

Greedy non-maximum suppression

Designed for 20-200 categories rather than 1000



As many separate

Feature vector of all the other categories plus the background e.g.: No Cat

a ya





The result: Each region is categorized in every image classes.



Ranking, selecting, merging



#### **Bounding Box Regression**

Linear regression model •

11/12/2019

- One per object category
- Input: last feature map cube of the conv net (pool5)
- Output: size and position modification to the bounding box:

(0, 0, 0, 0)

Proposal is good

- dx, dy, dw, dh

Training image regions

Input: Cached feature map cube (pool5)

Regression targets: (dx, dy, dw, dh) (normalized)



(.25, 0, 0, 0) Proposal too far to left

(0, 0, -0.125, 0)

pool5

Proposal too wide

78



**R-CNN Training** 

Step 1: Take a pretrained Convolutional Neural Network (e.g. AlexNet)

Cached region features vectors





#### Step 4: Train one SVM per class to classify region features



W

**Step 5** (bbox regression): For each class, train a linear regression model to map from cached features cubes to offsets/size of the boxes to fix "slightly wrong" position proposals

Training image regions

Cached region feature cube (pool5)

Regression targets (dx, dy, dw, dh) Normalized coordinates



(0, 0, 0, 0)

Proposal is good



(.25, 0, 0, 0) Proposal too far to left

## **R-CNN Problems**



(0, 0, -0.125, 0)

Proposal too

wide

- 1. Slow at test-time: need to run full forward pass of CNN for each region proposal
  - Recalculate the features again-and-again in the overlapping regions
- 2. SVMs and bbox regressors are post-hoc:
  - CNN features not updated in response to SVMs and regressors
- 3. Complex multistage training pipeline
  - Calculate the features for all the regions for all the training image first
  - Then train for SVM and bbox regressor separately





**R-CNN Problem #2**: Post-hoc training: CNN not updated in response to final classifiers and regressors

R-CNN Problem #3: Complex training pipeline

#### Solution:

Just train the whole system end-to-end all at once!

Slide credit: Ross Girschick







ennyi a lényeg: (training time) R-CNN -> Fast R-CNN 84 hours -> 9.5 hours







	R-CNN	Fast R-CNN	Faster R-CNN	
Test time per image (with proposals)	50 seconds	2 seconds	0.2 seconds	
(Speedup)	1x	25x	250x	
mAP (VOC 2007)	66.0	66.9	66.9	

#### d. Supervised VS unsupervised learning

#### Supervised Learning

 Learning from labeled examples (for which the answer is known)

#### Unsupervised Learning

 Learning from unlabeled examples (for which the answer is unknown)

#### • Reinforcement Learning

 Learning by trial and feedback, like the "child learning" example

#### Supervised learning

- We have prior knowledge of the desired output
  - Always have data set with ground truth (like image data sets with labels)
- Typical tasks
  - Classification
  - Regression

#### Unsupervised learning

- No prior knowledge of the desired output
  - Bessived radio sign
    - Received radio signals from deep space
- Typical tasks
  - Clustering
    - Representation learning
    - Density estimation

We wish to learn the inherent structure of (patterns in) our data.

# Use cases for unsupervised learning

- Exploratory analysis of a large data set
  - Clustering by data similarity
  - Enables verifying individual hypothesizes after analyzing the clustered data
- Dimensionality reduction
  - Represents data with less columns
  - Allows to present data with fewer features
  - Selects the relevant features
  - Enables less power consuming data processing, and/or human analysis

#### Unsupervised learning techniques

- Curse of dimensionality
- Principal component analysis (PCA)
- t-Distributed Stochastic Neighbor Embedding (t-SNE)
- Autoencoder



- · Scale the width, the depth, and the resolution uniformly!
- Can be used for any existing architecture, and the efficiency will be significantly better with the same performance
  - EfficientNet-B7 achieves stateof-the-art 84.4% top-1 / 97.1% top-5 accuracy on ImageNet, while being 8.4x smaller (number of parameters) and 6.1x faster on inference than the best existing ConvNet.
- Best performance can be reached by using NN to generate the optimal baseline ConvNet.



#### No. 10 Topic

a. Optimization problem of objective functions of NN - tessék??

https://www.youtube.com/watch?v=AoJQS10Ewn4

Optimization problems are commonly written in the form minimize f(x)Here, f is the objective function

- Objective Function
  - The value you are trying to optimize
  - Minimized or maximized
- Decision Variables
- The values the optimizer can change
- Also called design or manipulated variables

# Learning

- The questions are the following
  - What is the relationship of these optimal weights?

$$\mathbf{w}_{\text{opt}} \stackrel{???}{\Leftrightarrow} \mathbf{w}_{\text{opt}}^{(K)}$$
$$\mathbf{w}_{\text{opt}}^{(K)} : \min_{\mathbf{w}} \frac{1}{K} \sum_{k=1}^{K} (d_k - Net(\mathbf{x}_k, \mathbf{w}))^2$$

 How this new objective function should be minimized as quickly as possible?

### Learning – in practice

• Learning based on the training set:

$$\tau^{(K)} = \{ (\mathbf{x}_k, d_k); k = 1, ..., K \}$$

• Minimize the empirical error function  $(R_{emp})$ 

$$\mathbf{w}_{_{\text{opt}}}^{(K)}:\min_{\mathbf{w}}\frac{1}{K}\sum_{k=1}^{K}\left(d_{k}-Net\left(\mathbf{x}_{k},\mathbf{w}\right)\right)^{2}=\min_{\mathbf{w}}R_{emp}\left(\mathbf{w}\right)$$

• Learning is a multivariate optimization task

#### b. AdaGrad optimizer

- The AdaGrad algorithm (2011) individually adapts the learning rates of all model parameters by scaling them inversely proportional to the square root of the sum of all of their historical squared values
- The parameters with the largest partial derivative of the loss have a correspondingly rapid decrease in their learning rate, while parameters with small partial derivatives have a relatively small decrease in their learning rate
- The net effect is greater progress in the more gently sloped directions of parameter space
- AdaGrad performs well for some but not all deep learning models

# AdaGrad algorithm



Algorithm The AdaGrad algorithm	Remembers the
<b>Require:</b> Global learning rate $\epsilon$	entire history
<b>Require:</b> Initial parameter $\boldsymbol{\theta}$	evenly
<b>Require:</b> Small constant $\delta$ , perhaps $10^{-7}$ , for numeric	al stability
Initialize gradient accumulation variable $\boldsymbol{r}=\boldsymbol{0}$	
while stopping criterion not met $do$	
Sample a minibatch of $m$ examples from the training	g set $\{\boldsymbol{x}^{(1)},\ldots,\boldsymbol{x}^{(m)}\}$ with
corresponding targets $\boldsymbol{y}^{(i)}$ .	
Compute gradient: $\boldsymbol{g} \leftarrow \frac{1}{m} \nabla_{\boldsymbol{\theta}} \sum_{i} L(f(\boldsymbol{x}^{(i)}; \boldsymbol{\theta}), \boldsymbol{y}^{(i)})$	)
Accumulate squared gradient: $\boldsymbol{r} \leftarrow \boldsymbol{r} + \boldsymbol{g} \odot \boldsymbol{g}$	
Compute update: $\Delta \theta \leftarrow -\frac{\epsilon}{\delta + \sqrt{r}} \odot g$ . (Division	and square root applied
element-wise)	
Apply update: $\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} + \Delta \boldsymbol{\theta}$	
end while	

#### **RMSP** algorithm

- The RMSProp algorithm (2012) modifies AdaGrad to perform better in the nonconvex setting by changing the gradient accumulation into an exponentially weighted moving average
- In each step AdaGrad reduces the learning rate, therefore after a while it stops entirely!
- AdaGrad shrinks the learning rate according to the entire history of the squared gradient and may have made the learning rate too small before arriving at such a convex structure
- RMSProp uses an exponentially decaying average to discard history from the extreme past so that it can converge rapidly after finding a convex bowl, as if it were an instance of the AdaGrad algorithm initialized within that bowl

#### Video comparing adaptive and non-adaptive:

Three optimizer types are compared:

- SGD
  - Momentum types
    - Momentum
    - Nesterov AG
- Adaptív
   AdaGrad
  - Aua
  - AdaDelta

RmsProp
Adaptive ones are the fastest

#### c. Input vector normalization

#### Data regularization techniques

- Modification of the input vectors and internal data and internal parameters of the net
- Targeting to perform better in generalization
- Increases the loss during training phase
- Puts the parameters further away from a minimum with an expectation of it will find a deeper minimum
- In many cases these are heuristic methods with mostly experimental and partial mathematical proof

## Input vector normalization

0,45 Input · When the input vector contains high 1589,2 x =vector: 0.00143 and small mean values in different 0,32 vector positions it is usefull to 1423,2 mean:  $\bar{x} =$ normalize them 0.00132 Squeezes the number to the same 0,11 155,2 deviation:  $\sigma =$ range 0.00042 Speeds up the training process  $\frac{x-\bar{x}}{\sigma} = \begin{pmatrix} 1,18\\ 1,06\\ 0,26 \end{pmatrix}$ normalized  $x_{normed} =$ input vector:



40

#### **Input Normalization**

- Different normalization strategies exists for different input types
- · Showing it in two dimension, it shapes the input vector



Once you trained you net with a normalized training set, you have to apply normalization when a previously unseen vector (a new observation) is appled during inference. OK, but how do you know the statistics?

10/8/2019

# Input normalization ezample



#### d. DeconvNet, U-Net

#### Semantic Image Segmentation

- U-Net
- DeConvNet
- SegNet
  - Resolution controlling
  - Atrous convolutions, sub-pixel image combination

#### DeconvNet



- Instance-wise segmentation
- Two-stage training:
  - train on easy examples (cropped bounding boxes centered on a single object) first and
  - then more difficult examples



# Deconvnet: Extreme segmentation I



- Fully symmetrical convolutional network

   All convolution and pooling layers are reversed
  - Two stage training (first side trained for classification first)
- Takes 6 days to train on titan GPU
- Output probability map same size as input



# Deconvnet: Extreme segmentation II



correspondingly d nanual labels

# **U-Net**







- Designed for biomedical image processing: cell segmentation
- Data augmentation via applying elastic deformations,
  - Natural since deformation is a common variation of tissue
  - Smaller dataset is enough

64



#### e. Neural style transfer

An interesting application of the gradient ascent method is neural style transfer Could we use an input image and transform it into the style of an other input image?



Could we use an input image and transform it into the style of an other input image? Gradient ascent transforms the image according to a loss function. Can we find a loss function, which would preserve objects and another which preserves features connected to style?



Style transfer works, but It requires a lot of time, to generate an image. Many forward and backward passes are needed.

We could train a network that learns the result of this iterative transformation, and tries to predict it. Only a single pas is needed.



We have a loss function for content:

Can the same objects be found on both images? Content loss, Perceptual loss: this is a distance between the two embedded image vectors in the last features layers Style loss:

Can the same low level features, edges structures, simple patterns be found on both images Style loss: Distances between lower level representations of the images

Could we use an input image and transform it into the style of an other input image? Gradient ascent transforms the image according to a loss function. Can we find a loss function, which would preserve objects and another which preserves features connected to style?



#### Neural style transfer with Cycle Consistent GANs Input ı Monet Van Gogh Cezanne Ukiyo-e I I T I I I I. н L I I. L ī ı I L I

# Fast Neural Style Transfer

I



Could we use an input image and transform it into the style of an other input image?

Gradient ascent transforms the image according to a loss function.

Can we find a loss function, which would preserve objects and another which preserves features connected to style?



#### No. 11 Topic



Usage of Multilayer Perceptron

Hidder

- Multilayer perceptrons are used for
  - Classification
    - Supervised learning for classification
    - Given inputs and class labels
  - Approximation
    - Approximate an arbitrary function with arbitrary precision

#### b. Early stopping

#### Regularization and optimization methods

Different methods to increase the loss in the learning phase, but reduce overfitting and increase generalization capabilities

- Local response normalization
- Batch normalization
- Data augmentation (Enriching the data set)
- Early stopping
- Early stopping
   Ensemble methods
  - Network duplication
  - Bagging
  - Dropout

lan Goodfellow: regularization is "any modification we make to the learning algorithm that is intended to reduce the generalization error, but not its training error"



 Since the error function is not necessarily monotonic, the optimization goes on, but the suboptima are saved



 We have to start out from one point (say x<sub>1</sub>) and with an iterative method, we need to go towards the minimum

- We do not know where the curve is
- We know the value at  $f(x_1)$
- We know the derivative at x<sub>1</sub>
   f'(x<sub>1</sub>)
- Which way to go?
- Idea: follow the descending gradient!



$$f(x+\varepsilon) \approx f(x) + \varepsilon f'(x)$$

therefore

 $f(x - \varepsilon \operatorname{sign}(f'(x))) \leq f(x)$ 

This technique is called Gradient Descent (Cauchy, 1847).

Optimization goal is to find the f'(x) = 0 position.

(Critical or stationary points)

30/2019.

## **Stationary points**

f(x)

f(x)

f'(x) < 0

Negative gra

 $f(x_1)$ 

 $x_1$ 

x

f'(x) > 0

tangents

tive gradient

f(x)

4

- Local minimum, where f`(x)=0, and f(x) is smaller than all neighboring points
- Local maximum, where f`(x)=0, and f(x) is larger than all neighboring points
- Saddle points, where f'(x)=0, and neither minimum nor maximum



# Local and global minimum



In neural network parameter optimization we usually settle for finding a value of f that is very low, but not necessarily minimal in any formal sense.



# Gradient Descent in multidimensional input case 🚢

Steepest gradient descent iteration

$$\mathbf{x}(n+1) = \mathbf{x}(n) - \varepsilon \nabla f(\mathbf{x}(n))$$

- ε is the learning rate
- Choosing *ɛ*:
  - Small constant
  - Decreases as the iteration goes ahead
  - Line search: checked with several values, and the one selected, where f(x) is the smallest
- Stopping condition of the gradient descent iteration
  - When the gradient is zero or close to zero

#### Egy kis extra ide:

#### Jacobean Matrix

- Partial derivative of a vector  $\rightarrow$  vector function
- Specifically, if we have a function  $\mathbf{f}: \mathfrak{R}^m \to \mathfrak{R}^n$ then the Jacobian matrix  $\mathbf{J} \in \mathfrak{R}^{n \times m}$

of **f** is defined such that: 
$$\mathbf{J}_{i,j} = \frac{\partial}{\partial x_i} f(x_i)$$

$$\mathbf{J} = \begin{bmatrix} \frac{\partial \mathbf{f}}{\partial x_1} & \cdots & \frac{\partial \mathbf{f}}{\partial x_n} \end{bmatrix} = \begin{bmatrix} \frac{\partial f_1}{\partial x_1} & \cdots & \frac{\partial f_1}{\partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial f_m}{\partial x_1} & \cdots & \frac{\partial f_m}{\partial x_n} \end{bmatrix}$$

## 2<sup>nd</sup> order gradient descent method I

- $2^{nd}$  derivative in a specific direction:  $\mathbf{u}^{\mathrm{T}}\mathbf{H}\mathbf{u}$
- Second-order Taylor series approximation to the function f(x) around the current point  $\mathbf{x}_0$  where it

$$(\mathbf{x}) \approx f(\mathbf{x}_0) + (\mathbf{x} - \mathbf{x}_0)^{\mathrm{T}} \mathbf{g} + \frac{1}{2} (\mathbf{x} - \mathbf{x}_0)^{\mathrm{T}} \mathbf{H} (\mathbf{x} - \mathbf{x}_0)$$
 g: gradient at  $\mathbf{x}_0$   
H: Hessian at  $\mathbf{x}_0$ 

stepping towards the largest gradient:

f

$$\mathbf{x}_{0} - \varepsilon \, \mathbf{g} \approx \mathbf{x} \quad \rightarrow \quad \mathbf{x} - \mathbf{x}_{0} \approx -\varepsilon \, \varepsilon \, \mathbf{g}$$
$$f(\mathbf{x}) \approx f(\mathbf{x}_{0} - \varepsilon \, \mathbf{g}) \approx f(\mathbf{x}_{0}) - \varepsilon \, \mathbf{g}^{\mathsf{T}} \mathbf{g} + \frac{1}{2} \varepsilon^{2} \, \mathbf{g}^{\mathsf{T}} \mathbf{H} \mathbf{g}$$

(+itt már folyt a newton opt.)

$$x_2 \qquad \nabla f(\mathbf{x}) \qquad x_1$$

#### 2<sup>nd</sup> derivatives

- 2<sup>nd</sup> derivative determines the curvature of a line in 1D
- In nD, it is described by the Hessian Matrix

$$H(f(x_{i,j})) = \frac{\partial^2}{\partial x_i \partial x_j} f(x) = \frac{\partial^2}{\partial x_j \partial x_i} f(x)$$

• The Hessian is the Jacobian of the gradient.



#### 2<sup>nd</sup> order gradient descent method II

• Analyzing:  $f(\mathbf{x}_0 - \varepsilon \mathbf{g}) \approx f(\mathbf{x}_0) - \varepsilon \mathbf{g}^{\mathsf{T}} \mathbf{g} + \frac{1}{2} \varepsilon^2 \mathbf{g}^{\mathsf{T}} \mathbf{H} \mathbf{g}$ 

Original value Expected Correction due to improvement curvature

- When the third term is too large, the gradient descent step can actually move uphill.
- When it *is zero or negative*, the Taylor series approximation predicts that increasing ε *forever will decrease f* forever.
- In practice, the Taylor series is unlikely to remain accurate for large ε, so one must resort to more heuristic choices of ε in this case.

• When it is positive, solving for  
the optimal step 
$$\mathcal{E}^* = \frac{\mathbf{g}^T \mathbf{g}}{\mathbf{g}^T \mathbf{H} \mathbf{g}}$$

99

#### d. Weight regularization (L1, L2)



- Modification of the input vectors or the internal data composition of the network
  - Input normalization
  - Batch normalization
- Modification of the cost function (involving the weight magnitudes)
  - L1 and L2 regularization
  - (weight penalty)
- Temporal Modification of the net architecture in training phase

- Dropout



#### e. Pooling

Methods of data size reduction:

- Pooling
- Convolution with strides
- Convolution without padding

Data aggregation:

• Stride convolution, pooling

# Pooling

- Pooling summarizes statistically the extracted features from the same location on a feature map
- Mathematically, it is a local function over 1D or 2D data
  - input:
    - Segment of a vector in 1D
    - rectangular neighborhood in 2D
  - Function
    - Statistical (maximum: max-pool)
    - L2 norm
    - Weighted average (weights proportional of the distance of the central element)
- In most cases: stride > 1
  - This leads to downsampling
- Pooling introduces some shift invariancy

х



Convolved	Pooled
feature	feature

s = 10

- Max pooling is the most used pooling in CNN
- Picks the largest value from a neighborhood
- Non-linear
- Statistical filter
- Downsampling depends on the stride



У

Max pooling









after average pooling



(b) Illustration of average pooling drawback

## Architecture of a typical Convolution Neural Network



47



- Data reduction (pooling)
- Combination of the features aggregating information (fully connected layer)
- Decision (fully connected layer with soft-max activation)







• Pooling layers summarize the outputs of neighboring neurons in the same kernel map.



- Top 1 error rate by 0.4%
- Top 5 error rates by 0.3%



#### No. 12 Topic

a. Perceptron convergence theorem (no proof) - nem tudom hol kezdodik a bizonyitas.hahaha

# Perceptron Convergence theorem (1)

Assumptions:

- **w**(0)=0
- the input space is linearly separable, therefore w<sub>o</sub> (stands for w<sub>optimal</sub>) exists:

$$x \in X^+$$
:  $w_1^T x > 0: d = 1$   
 $x \in X^-$ :  $w_1^T x < 0: d = 0$ 

- Let us denote  $\tilde{x} = -x$ 

$$\widetilde{x} \in \widetilde{X}^-$$
:  $w^T \widetilde{x} > 0$ :  $d = 1$ 

For the proof, see also: Simon Haykins: Neural Networks and Learning Machines, Section 1.3: <u>http://dai.fmph.uniba.sk/courses/NN/haykin.neural-networks.3ed.2009.pdf</u>

# Perceptron Convergence theorem (2)

- Idea:
  - During the training, the network will be activated with those input vectors (one after the other), where the decision is wrong, hence non zero adaptation is needed:

$$x(j) \in X^+$$
:  $w^T(j)x(j) < 0, y = 0, d = 1$   
 $x(j) \in \widetilde{X}^-$ :  $w^T(j)x(j) < 0, y = 0, d = 1$ 

- Note: The error function is always positive (  $\varepsilon = 1$  )

## Perceptron Convergence theorem (3)

- According to the learning method:
- $w(n+1)=w(0)+\eta x(0)+\eta x(1)+\eta x(2)+\eta x(3)+...+\eta x(n)$

$$x(j) \in X^+$$
:  $w^T(j)x(j) < 0, y = -1, d = 1$   
 $x(j) \in \widetilde{X}^-$ :  $w^T(j)x(j) < 0, y = -1, d = 1$ 

- The decision boundary will be:

 $\eta w^T x=0$ which means that  $\eta$  is a scaling factor, therefore it can be choosen for any positive number. Let us use  $\eta=1$ , therefore  $\eta\varepsilon=1$  • We will calculate  $||w(n+1)||^2$  in two ways, and give an upper and a lower boundary, and it will turn out that an  $n_{max}$  exists, and beyond that the lower boundary is higher than the upper boundary (squeeze theorem, sandwitch lemma (*közrefogási elv*, *rendőr elv*))

#### lower limit (1)

According to the learning method, the presented input vectors are added up:

$$w(n+1) = w(0) + x(0) + x(1) + \dots + x(n)$$
 w(0)=0

Multiply it with  $w_o^T$  from the left:

$$w_o^T w(n+1) = w_o^T x(0) + w_o^T x(1) + \dots + w_o^T x(n)$$
  

$$0 < \alpha \le w_o^T x(j)$$
 Because each input vector (or its opposite) were selected that way.  

$$0 < \alpha = \min_{x(n) \in \{X^+, \tilde{X}^-\}} w_o^T x(n)$$

 $w_o^T w(n+1) \ge n\alpha$ 

#### lower limit (2)

 $w_o^T w(n+1) \ge n\alpha$ 

We apply Cauchy Schwarty inequality  $\|a\|^2 \|b\|^2 \ge \|a^T b\|^2$ 

$$\|w_0^T\|^2 \|w(n+1)\|^2 \ge \|w_o^T w(n+1)\|^2 \ge n^2 \alpha^2$$

Lower limit:

$$\|w(n+1)\|^2 \ge \frac{n^2 \alpha^2}{\|w_0^T\|^2}$$
 Lower limit proportional with n<sup>2</sup>

# upper limit (1)

ß

Ø

Let us have a different synthetization approach of w(n+1):

$$w(k+1) = w(k) + x(k)$$
 for k= 0 ... r

Squared Euclidian norm:  $||w(k+1)||^2 = ||w(k)||^2 + ||x(k)||^2 + 2w(k)^T x(k)$ 

 $w(k)^T x(k) < 0$  Because each input vector (or its opposite) were selected that way.

$$\|w(k+1)\|^2 \le \|w(k)\|^2 + \|x(k)\|^2$$

$$||w(k+1)||^2 - ||w(k)||^2 \le ||x(k)||^2$$
 for k= 0 ... n

upper limit (2)  

$$\|w(k+1)\|^{2} - \|w(k)\|^{2} \le \|x(k)\|^{2}$$
Telescoping sum:  $\sum_{i=1}^{n} (a_{i+1} - a_{i}) = a_{n+1} - a_{1}$   
Example:  $\sum_{i=1}^{n} (a_{i+1} - a_{i}) = a_{2} - a_{1} + a_{3} - a_{2} + a_{3} - a_{2} + a_{3} - a_{3} + a_{5} - a_{4} + a_{5} - a_{5} + a_{5}$ 

Note that there is a telescoping sum in the left hand side.

$$\|w(n+1)\|^{2} - \|w(0)\|^{2} = 0 \qquad 0 < \beta = \max_{x(k) \in \{X^{+}, \widetilde{X}^{-}\}} \|x(k)\|^{2}$$
$$\|w(0)\|^{2} = 0 \qquad \|w(n+1)\|^{2} \le (n+1)\beta$$
Upper limit linearly proportional with n

#### comparing upper and lower limits $||w(n+1)||^{2}$ $||w(n+1)||^2 \le (n+1)\beta$ $\|w(n+1)\|^2 \ge$ $n^2\alpha$ n

Linear upper limit and squared lower limit cannot grow unlimitedly

n<sub>max</sub> should exist

$$n_{\max} = \frac{\beta \|w_0\|^2}{\alpha^2}$$

# b. Momentum optimizer

# Momentum I

- Introduced in 1964
- Physical analogy
- · The idea is to simulate a unity weight mass
- It flows through on the surface of the error function
- Follows Newton's laws of dynamics
- Having *v* velocity
- Momentum correctly traverses the canyon • lengthwise, while gradient steps waste time moving back and forth across the narrow axis of the canyon.



# **Momentum II: velocity considerations**

The update rule is given by:

$$\boldsymbol{v} \leftarrow \alpha \boldsymbol{v} - \epsilon \nabla_{\boldsymbol{\theta}} \left( \frac{1}{m} \sum_{i=1}^{m} L(\boldsymbol{f}(\boldsymbol{x}^{(i)}; \boldsymbol{\theta}), \boldsymbol{y}^{(i)}) \right),$$
$$\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} + \boldsymbol{v}$$

The velocity  $\boldsymbol{v}$  accumulates the gradient elements  $\nabla_{\boldsymbol{\theta}} \left( \frac{1}{m} \sum_{i=1}^{m} L(\boldsymbol{f}(\boldsymbol{x}^{(i)}; \boldsymbol{\theta}), \boldsymbol{y}^{(i)}) \right)$ . The larger  $\alpha$  is relative to  $\epsilon$ , the more previous gradients affect the current direction.

#### Terminal velocity is applied when it finds descending gradient permanently:

$\epsilon   m{g}  $	
9/30/2019 $1-lpha$ 37	,
Algorithm Stochastic gradient descent (SGD) with momentum	-
<b>Require:</b> Learning rate $\epsilon$ , momentum parameter $\alpha$ .	_
<b>Require:</b> Initial parameter $\boldsymbol{\theta}$ , initial velocity $\boldsymbol{v}$ .	
while stopping criterion not met $\mathbf{do}$	
Sample a minibatch of $m$ examples from the training set $\{x^{(1)}, \ldots, x^{(m)}\}$ wit	h
corresponding targets $\boldsymbol{y}^{(i)}$ .	
Compute gradient estimate: $\boldsymbol{g} \leftarrow \frac{1}{m} \nabla_{\boldsymbol{\theta}} \sum_{i} L(f(\boldsymbol{x}^{(i)}; \boldsymbol{\theta}), \boldsymbol{y}^{(i)})$	
Compute velocity update: $\boldsymbol{v} \leftarrow \alpha \boldsymbol{v} - \epsilon \boldsymbol{g}$	
Apply update: $oldsymbol{ heta} \leftarrow oldsymbol{ heta} + oldsymbol{v}$	
end while	

#### Momentum demo

 What does the parameter of the momentum method means, and how to set them?
 https://distill.pub/2017/momentum/

Why Momentum Really Works



9/30/2019

Algorithms with changing but not adaptive learning rate:

- Stochastic Gradient Descent algorithm
- Momentum algorithm
- Nesterov momentum update

itt tuti van valami hasznos még:

https://mlfromscratch.com/optimizers-explained/#/
c. Properties of CNN: sparsity, parameter sharing, equivariance, invariance
 to shifting -> mi a kutyafule

Properties of Convolutional Neural Networks I:

Sparsity



- A few dozen free parameter describes the operation of a layer
- Receptive field organization similar to natural neural vision systems



#### Parameter sharing

- Same parameters everywhere in the layer
- Contribution to the gradient of a weight from many positions
- Reduces the risk of overfitting
- Reduces the risk of dying RELU (dying cell)
  - When it happens, an entire feature extractor on a layer is dying



- Variable input size
  - The input image is either resized or padded



Equivalent representation:

• Equvariance to translation

• The output shifts with an input shift

- In a fully connected neural network, each input is a dedicated channel for a certain input parameter-therefore the inputs cannot be swapped
  - · Like bank example, one cannot replace the age input with the salary input
- In CNN, the image can be shifted, because the inputs are not dedicated and the features are identified anywhere

Space invariance means here that the functionality of a 2D function is not changing in space. This enables the detection of a certain image feature anywhere on the image.



### Simple RNN Training Example: Predicting the next letter

Example: Character-le Language M	evel Iodel				
Vocabulary: [h,e,l,o]	One-hot encoding				
Example trai sequence: " <b>hello</b> "	ning	input layer 0 0 input chars: "h"	0 1 0 0	0 0 1 0 "I"	0 0 1 0





Example: Character-level Language Model Sampling

Vocabulary: [h,e,l,o]

At test-time sample characters one at a time, feed back to model



### Example: Character-level Language Model Sampling Ba

Vocabulary: [h,e,l,o] Backpropagation can be started using negative log likelihood cost function

At test-time sample characters one at a time, feed back to model





Explain Images with Multimodal Recurrent Neural Networks, Mao et al. Deep Visual-Semantic Alignments for Generating Image Descriptions, Karpathy and Fei-Fei Show and Tell: A Neural Image Caption Generator, Vinyals et al. Long-term Recurrent Convolutional Networks for Visual Recognition and Description, Donahue et al.

Learning a Recurrent Visual Representation for Image Caption Generation, Chen and Zitnick

## Image captioning example Recurrent Neural Network



# **Convolutional Neural Network**

2019-11-25





Alexnet: scored 5 best guesses

test image

image
conv-64
conv-64
maxpool
conv-128
conv-128
maxpool
conv-256
conv-256
maxpool
conv-512
conv-512
maxpool
conv-512
conv-512
maxpool
FC-4096
FC-4096
FG 1090
softwax



#### image conv-64 conv-64 maxpool conv-128 conv-128 maxpool conv-256 conv-256 maxpool conv-512 conv-512 maxpool conv-512 conv-512 maxpool FC-4096 FC-4096



### test image

x0
<STA
RT>







## Image captioning Example: Results



A cat sitting on a suitcase on the floor



Two people walking on the beach with surfboards 2019-11-25



A tennis player in action on the court

branch

A cat is sitting on a tree



A dog is running in the grass with a frisbee



Two giraffes standing in a grassy field



A white teddy bear sitting in the grass



A man riding a dirt bike on a dirt track 40

# Image captioning: Failure cases

A bird is perched on a tree branch



A woman is holding a cat in her hand



A person holding a computer mouse on a desk



beach holding a surfboard



A man in a baseball uniform throwing a ball

e. Adversarial attacks

We have a high number of parameters to be optimized An even higher-dimensional input

The network works well in practice, but can not cover all the possible inputs





One can exploit that there will be regions in the input domain, which were not seen during training

Adversarial noise:

• I have a working well-trained classifier:



What should I add to the input to cause misclassification:

A special, low amplitude additive noise:



The two images are the same for human perception

Knowing a trained network one can identify modifications (which does not happen in real life), which change the network output completely





#### Adversarial noise - does not work in practice



Knowing a trained network one can identify modifications (which does not happen in real life), which change the network output completely

Luckily this low amplitude noise is not robust enough in real life (lens distortion and other additive noises)



### Sticker based adversarial attacks

High intensity noise concentrated on a small region of the input image:

 $C_{d} = N\left(I + \sum_{i=1}^{k} St_{i}(x_{i}, y_{i}, w_{i}, h_{i}) + \sum_{j=1}^{l} St_{j}(x_{j}, y_{j}, w_{j}, h_{j})\right)$ 

Parameters are the positions (x,y) and size (w,h) of the stickers

It was shown that these attacks are **robust** enough to be applied in practical applications

Does this mean that convolutional neural networks can not be used in critical problem in practice anymore?





#### No. 13 Topic

### a. Elementary set separation by a single neuron

Elementary set separation by a single neuron (1)

• Let us use  $\varphi(.)$  step nonlinear function for siplicity:

$$y = \varphi(u) = \frac{sign(u)}{2} + \frac{1}{2} = \begin{cases} 1, \text{ if } u \ge 0\\ 0, \text{ else} \end{cases}$$

The output of the neuron will be binary:



• Neuron with *m* inputs has an *m* dimensional input space

 $\mathbf{w}^T \mathbf{x} = \mathbf{0}$ 

- Neuron makes a linear decision for a 2 class problem
- The decision boundary is a hyperplane defined:





•

•

P-ITEEA-0011

### Implementation of a single logical function by a single



- Furthermore instead of 2D, we can actually come up with the *R* dimensional AND function.
- The weights corresponding to the inputs are all 1 and threshold should be R – 0.5. As a result the actual weights of the neuron are the following:

$$\mathbf{w}^{\mathrm{T}} = (-(R-0.5), 1, ..., 1)$$



- However we cannot implement every logical function by a linear hyper plane.
- Exclusive OR (XOR) cannot be implemented by a single neuron (linearly not separable)



#### b. Local response normalization

- Regularization and normalization methods
  - Local response normalization
  - Data augmentation
  - Early stopping
  - Ensembling

### Regularization and optimization methods

Different methods to increase the loss in the learning phase, but reduce overfitting and increase generalization capabilities — Local response normalization

Batch normalization

- Batch normalization
   Data augmentation
- (Enriching the data set)
- Early stopping
- Ensemble methods
  - Network duplication
     Bagging
  - BaggingDropout

# Local response normalization I

- Implementation of the Lateral inhibition from neurobiology
  - If a neuron starts spiking strongly in a layer it inhibits (suppresses) the of the neighboring cells
  - Winner take all (have a strong decision)
  - Balances the asymmetric responses of neurons in different areas of the layer
- Useful when we are dealing with ReLU neurons
  - Normalizes the unbounded activation of the ReLU neurons
     Avoids concentrating and delivering large values through layers
  - It enhances high spatial frequencies by suppressing the local neighbors of the strongest neuron

11/5/2019

https://towardsdatascience.com/difference-between-localresponse-normalization-and-batch-normalization-272308c034ac lan Goodfellow: regularization is "any modification we make to the learning algorithm that is intended to reduce the generalization error, but not its training error"



### Local response norm. vs batch norm.

Both work within one convolutional layer

- Normalization either through the feature maps or within one feature map
- Normalization is done for one input image



- Normalization done for all the pixels in all the feature maps within a layer
- Normalization is done for the entire batch





c. Unpooling -> lasd No. 11 Topic

d. Representations: Blum and Li theorem, construction

### Representation - Blum and Li theorem

- $F(x) \in L^2$ Theorem:  $\forall \varepsilon > 0, \exists w$ 
  - $\int \cdots_{\mathbf{x}} \int \left( F(\mathbf{x}) Net(\mathbf{x}, \mathbf{w}) \right)^2 \mathbf{d}x, \dots \mathbf{d}x_N < \varepsilon$
- Proof:

  - Using the step functions: S
  - From elementary integral theory it is clear every function can be approximated by appropriate step function sequence



9/24/2019

## Representation – Blum and Li theorem



- The step function can have arbitrary narrow steps
- For example each step could be divided into two sub-steps
- Therefore we can synthetize a function with arbitrary precision



$$F(x) \cong \underbrace{\sum_{i} F(x_i) I(x_i)}_{s(x)}$$

## Representation – Blum and Li construction

- This construction ...
  - ... has no dimensional limits
  - ... has no equidistance restrictions on tiles (partitions)
  - ... can be further fined, and the approximation can be any precise
- 2 dimensional example
  - The tiles are the top of the columns for each approximation cell

## Blum and Li – Limitations



- The size of the FFNN constructed via this method is quite big
- Consider the task on the picture, where there are 1000 by 1000 cell to approximate the function
- General case:
   ~2 Million neurons are needed
- Smoother approximation needs more
- The network architecture is synthetized (constracted), the weights are generated
- We are after to find a less complicated architectures

### No. 14 Topic

### a. Principal component analysis (PCA) Principal component analysis (PCA)

- Technique for dimensionality reduction
- Invented by Karl Pearson (1901)
- Linear coordinate transformation
  - converts a set of observations of possibly correlated variables
  - into a set of values of linearly uncorrelated orthogonal variables called principal components
- Deterministic algorithm

#### PCA algorithm

1. Mean normalization: For every value in the data, subtract its mean dimension value. This makes the average of each dimension zero.



2. Standardization (optional): Do it, if you want to have each of your features the same variance.



3. Covariance matrix: Calculate the covariance matrix

Covariance (formal definition)

- Assume that x are random variable vectors
- We have *n* vectors

Variance(x) = 
$$\frac{1}{n} \sum_{i=1}^{n} (x_i - \bar{x})^2$$
  
=  $\frac{1}{n} \sum_{i=1}^{n} (x_i - \bar{x}) (x_i - \bar{x})$ 

 $Covariance(\mathbf{x}, \mathbf{y}) = \frac{1}{n} \sum_{i=1}^{n} (x_i - \bar{x}) (\mathbf{y}_i - \bar{y})$ 

- Covariance(x, x) = var(x)
- Covariance(x, y) = Covariance(y, x)



x

#### **Covariance matrix** $\begin{bmatrix} cov(x_1, x_1) & cov(x_1, x_2) & \cdots & cov(x_1, x_m) \\ cov(x_2, x_1) & cov(x_2, x_2) & \cdots & cov(x_2, x_m) \end{bmatrix}$ $\cdots$ $cov(x_1, x_m)^{-1}$ **Diagonal elements** $Cov(\Sigma) =$ are variances, i.e. : . 1 . Cov(x, x) = var(x) $\begin{bmatrix} cov(x_m, x_1) & cov(x_m, x_2) & \cdots & cov(x_m, x_m) \end{bmatrix}$ *n* is the number of the vectors *m* is the $Cov(\Sigma) = \frac{1}{n}(X - \overline{X})(X - \overline{X})^T$ ; where X =dimension **Covariance Matrix** $Cov (\Sigma) = \begin{bmatrix} var(x_1, x_1) & cov(x_1, x_2) & \cdots & cov(x_1, x_m) \\ cov(x_2, x_1) & var(x_2, x_2) & \cdots & cov(x_2, x_m) \\ \vdots & \vdots & \vdots & \vdots \\ cov(x_m, x_1) & cov(x_m, x_2) & \cdots & var(x_m, x_m) \end{bmatrix}$ is symmetric commutative 11/19/2019

### PCA algorithm

- 1. Mean normalization: For every value in the data, subtract its mean dimension value. This makes the average of each dimension zero.
- 2. Standardization (optional): Do it, if you want to have each of your features the same variance.
- 3. Covariance matrix: Calculate the covariance matrix
- 4. Eigenvectors and eigenvalues of the covariance matrix
  - Note: Each new axis (PC) is an eigenvector of the data. The standard deviation of the data variance on the new axis is the eigenvalue for that eigenvector.



Principal components will be orthogonal. Uncorrelated, independent!

- 5. Rank eigenvectors by eigenvalues
- 6. Keep top k eigenvectors and stack them to form a feature vector
- 7. Transform data to PCs:
  - New data = feature vectors (transposed) \* original data

From *k* original variables:  $x_1, x_2, \dots, x_k$ :

Produce k new variables:  $y_1, y_2, \dots, y_k$ :  $y_1 = a_{11}x_1 + a_{12}x_2 + \dots + a_{1k}x_k$  $y_2 = a_{21}x_1 + a_{22}x_2 + \dots + a_{2k}x_k$  $y_k$ 's are ... **Principal Components**  $y_k = a_{k1}x_1 + a_{k2}x_2 + \dots + a_{kk}x_k$ 

 $\{a_{11}, a_{12}, ..., a_{1k}\}$  is 1st **Eigenvector** of first principal component {a21,a22,...,a2k} is 2nd Eigenvector of of 2nd principal component

 $\{a_{k1}, a_{k2}, ..., a_{kk}\}$  is kth **Eigenvector** of of kth principal component

# Principal Component Analysis (PCA)



- The idea is to project the data onto a subspace which compresses most of the variance in as little dimensions as possible.
- Each new dimension is a principle component
- The principle components are ordered according to how much variance in the data they capture

Proportion

Example:



We have to choose how many PCs to use from the top

## How many PCs to use?

Calculate the proportion of variance for each feature

- prop. of var. = 
$$\frac{\lambda_i}{\sum_{i=1}^n \lambda_i}$$

- $-\lambda_i$  are the eigen values
- Rich a predefined threshold
- Or find the elbow of the • Scree plot

of variance Scree plot 0.9 0.8 0.7 0.6 Variance Cumulative variance 0.5 0.4 0.3 Scree plot elbow 0.2 0. 010 ŝ PO 4 ŝ ñ 50 D0 5 S ò Principal components

11/19/2019

		England	Wales	Scotland	N Ireland
PCA Example	Cheese	105	103	103	66
I CA Example	Carcass meat	245	227	242	267
<ul> <li>Weekly food</li> </ul>	Other meat	685	803	750	586
consumption of the	Fish	147	160	122	93
four countries	Fats and oils	193	235	184	209
<ul> <li>food types: variables</li> </ul>	Sugars	156	175	147	139
	Fresh potatoes	720	874	566	1033
- countries: observations	Fresh Veg	253	265	171	143
<ul> <li>Clustering the</li> </ul>	Other Veg	488	570	418	355
countries:	Processed potatoes	198	203	220	187
<ul> <li>Needs visualization in</li> </ul>	Processed Veg	360	365	337	334
17 dimension	Fresh fruit	1102	1137	957	674
DCA: roduco	Cereals	1472	1582	1462	1494
PCA. reduce	Beverages	57	73	53	47
dimensionality	Soft drinks	1374	1256	1572	1506
	Alcoholic drinks	375	475	458	135
http://www.sdss.jhu.edu/~szalay/clas	Confectionery	54	64	62	41

s/2016-oldold/SignalProcPCA.pdf 11/19/2019

http:

UK food consumption in 1997 (g/person/week). Source: DEFRA

Projections onto first principal component (1-D space) 0.5 Wal Eng Scot N Ire -0.5 -1 100 PC1 200 300 -200 -100 400 500 0

# **PCA Example**

- From PC1, two clusters ٠ are well separable
- Including PC2, the four • clusters can be well separated







11/19/2019

# **Coefficients of the Principal Components**



Load plot shows the coefficients of the original feature vectors to the principal components

a pa

b. Back-propagation through time -> shit ezt csak most vettem eszre, lent a megoldas

## Learning

- The Rosenblatt algorithm is inapplicable,
  - the error and desired output in the hidden layers of the FFNN is unknown
- Someway the error of the whole network has to be distributed to the internal neurons, in a feedback way



## Sequential back propagation

• Adapting the weights of the FFNN (recursive algorithm)

$$w_{ij}^{(l)}(k+1) = w_{ij}^{(l)}(k) + \Delta w_{ij}^{(l)}(k)$$
  
$$\Delta w_{ij}^{(l)}(k) = ?$$

The weights are modified towards the differential of the error function (delta rule):

$$\Delta w_{ij}^{(l)} = -\eta \, \frac{\partial R_{emp}}{\partial w_{ij}^{(l)}}$$

The elements of the training set adapted by the FFNN sequentially

$$R_{emp} = R_{emp}(y(\mathbf{x}), d)$$



- Though we showed how to modify the weights with back propagation, its most important value that it can calculate the gradient
- The weight updates can be calculated with different optimization methods, after the gradients are calculated
- Various optimization method can drastically speed up the training (100x, 1000x)

Conclusion

- For known functions (according to Blum-Li)
  - One can define a Neural Network architecture
  - And generate the weights
  - That it can represent the known function with arbitrary precision
- For unknown but existing function defined by IO pairs (according to statistic learning)
  - One can find a Neural Network architecture
  - And train the network (optimize the weights)
  - Reach arbitrary precision with high number of IO pairs
  - The trained network will be able to well predict previously unknown IO pairs (generalization)

## **Back propagation**

- We have seen last time how to calculate the gradient in a multilayer fully connected network using back propagation
  - The introduced method was based on gradient descent method
- However, being able to calculate gradient, we might select any of the above methods, which leads to orders of magnitude faster convergence



#### Back propagation can be applied for batch normalized layers



## Simple RNN Training Example: Predicting the next letter



## Back propagation through time

- Assuming that the length of the input vector sequence is limited
- · It became a feedforward neural net
- Possible to apply back propagation
- We need multiple vector sequences to train!





Backpropagation through time

Forward through entire sequence to compute loss, then backward through entire sequence to compute gradient





## Truncated Backpropagation through time

Run forward and backward through chunks of the sequence instead of whole sequence



Carry hidden states forward in time forever, but only backpropagate for some smaller number of steps!



-> en ebbol nem ertettem meg, de azért beraktam a diakat

### c. Stochastic gradient descent optimizer

Stochastic process is a process, where we cannot observe the exact values. In these processes, our observations are always corrupted with some random noise.

- Algorithms with changing but not adaptive learning rate
  - Stochastic Gradient Descent algorithm
  - Momentum algorithm
  - Nesterov momentum update

#### Stochastic Gradient Descent (SGD) algorithm:

- Introduced in 1945
- Gradient Descent method, plus:
- Applying mini batches
- Changing the learning rate during the iteration

## Learning rate at SGD

Sufficient conditions to guarantee convergence of SGD:

$$\sum_{k=1}^{\infty} \epsilon_k = \infty, \quad \text{and} \quad \sum_{k=1}^{\infty} \epsilon_k^2 < \infty. \qquad \begin{array}{c} \epsilon \text{ is the learning} \\ \text{rate, also marked} \\ \text{with } \eta \text{ sometimes} \end{array}$$

• In practice:

10/1/2010

$$\epsilon_k = (1 - \alpha)\epsilon_0 + \alpha\epsilon_\tau \qquad \alpha = \frac{k}{\tau}$$

• After iteration  $\tau$ , it is common to leave  $\varepsilon$  constant

# Stochastic Gradient Descent algorithm



**Algorithm** Stochastic gradient descent (SGD) update at training iteration k

Require: Learning rate  $\epsilon_k$ . Require: Initial parameter  $\boldsymbol{\theta}$ while stopping criterion not met do Sample a minibatch of m examples from the training set  $\{\boldsymbol{x}^{(1)}, \ldots, \boldsymbol{x}^{(m)}\}$  with corresponding targets  $\boldsymbol{y}^{(i)}$ . Compute gradient estimate:  $\hat{\boldsymbol{g}} \leftarrow +\frac{1}{m} \nabla_{\boldsymbol{\theta}} \sum_{i} L(f(\boldsymbol{x}^{(i)}; \boldsymbol{\theta}), \boldsymbol{y}^{(i)})$ Apply update:  $\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} - \epsilon \hat{\boldsymbol{g}}$ end while

#### where: L is the cost function

 $\theta$  is the total set of  $w_{i,j}^{(l)}$  (and all other parameters to optimize)

## **Stochastic Gradient Descent algorithm**

- This very elongated quadratic function resembles a long canyon.
- Gradient descent wastes time repeatedly descending canyon walls, because they are the steepest feature.
- Because the step size is somewhat too large, it has a tendency to overshoot the bottom of the function and thus needs to descend the opposite canyon wall on the next iteration.



# **Training I**

- Stochastic Gradient Descent (with momentum)
   ADAM method was introduced in 2014 only (2 years later)
- Minimizing the negative log-likelihood (cross-entropy) loss function
- With L2 regularization (weight penalty):

$$L(w) = \sum_{i=1}^{N} \sum_{c=1}^{1000} -y_{ic} \log f_c(x_i) + \epsilon ||w||_2^2$$

predicted probability of class c for image x

indicator that example i has label c

14/6/2010

**d. Object detection problem explained** -> itt nem tudom mirol maradtam le, de en ezt se latom a diaban...



138

## Object detection/localization and classification

- Chicken and egg problem
  - You need to know that it is a bicycle before able to say that both a wheel part and a pipe segment belongs to the same object
  - You need to know that the red box contains an object before you can recognize it. (Cannot recognize a bicycle if you try it from separated parts)
- Our brain does it parallel
- How neural nets can solve it?
  - Detection by regression?
  - Detection by classification?

#### **Neural networks for regression**





Multiple object detection on a single image

Classification is good for a single object (can be extended for  ${\sf k}$  objects – top  ${\sf k}$  candidates)

How could we detect objects in general, when the number of objects is unknow



#### Dunno what I'm doin:

#### **Object detection as regression**

RCNN

Single Shot Object Detector (SSD) (2016 March)

You Only Look Once YOLO (2016 May)

### **R-CNN**

Region proposal CNN network

Separate the problem of object detection and calssification

It consists of three modules.

The first generates category-independent region proposals. These proposals define the set of candidate detection avail-able to detector.

The second module is a large convolutional neural network that extracts a fixed-length feature vector from each region.

The third module is a set of class- specific linear SVMs

#### Ezek jók lehetnek:

- <u>https://towardsdatascience.com/5-significant-object-detection-challenges-and-solutio</u> ns-924cb09de9dd
- <u>https://www.coursera.org/lecture/deep-learning-in-computer-vision/object-detection-pr</u>
   <u>oblem-tvLPq</u>
- + a szövege:

Object detection: probably the key problem in computer vision.

The <u>goal</u> of object detection is to detect the presence of object from a certain set of classes, and locate the exact position in the image.

We can informally divide all objects into two big groups: things and stuff. Things are objects of certain size and shape like cars, bicycles, people, animals, planes. We can specify where object is located in image with a bounding box. Stuff is more likely a region of image which correspond to objects like road, or grass, or sky, or water. It is easier to specify the location of a sky by marking the region in an image, not by a bounding box.

Now, we will talk mostly about **detection of things**. To detect a stuff, it is better to use semantic image segmentation methods, which we will discuss during week five. Compared to image classification, output of the detector is structured. Each object is usually marked with a bounding box and class label. Bounding box is described by position of one of the corners, and by width and size of the box. Object position and class are annotated in ground truth data. If only part of this information is annotated, we call it the "weak" annotation. For example, only the presence of object in image can be annotated without a bounding box. There is a lot of research into object detection with weak annotation, but performance of such algorithms is lower compared to algorithms trained to use full annotation. To check whether the detection is correct, we compare the predicted bounding box with ground truth bounding box. The metric is intersection over union or IoU. It is the ratio of area of intersection of predicted in ground truth bounding boxes to the area of the union on these boxes as shown on the slide. Either IoU is larger than the threshold, then the detection is correct. The larger the threshold, the more precisely detector should localize objects. Currently, the threshold is usually set to 0.5. The detector output is a set of detection proposals. Usually, for each proposal the detector also gives a score as a measure of confidence in the detection. So, we can rank all proposals according to the score. Each proposal is considered. If IoU is larger than the threshold, then it is the true positive detection. If IoU is lower, then it's false positive detection. If some ground truth object is not detected, then it is marked as misdetection or false negative. On the whole dataset, you can measure the precision and the recall of the detector. The precision is the ratio between the number of true detections and the number of all detections. The recall is the ratio of number of true detection to the number of objects annotated in the ground truth data. By varying some parameter in the detector, usually the threshold on detection score, you can simultaneously change precision and recall. Then you can plot the precision-recall curve. To compare two detectors correctly, you should compare their precision-recall curves. If one curve is generally higher than the other curve, then the first detector is better than the other. To measure the overall quality of the detector with one number, we compute average precision. It is the mean of 11 points on the curve for recalls from zero to one, by 0.1

intervals. If you have multi-class detector, then you can compute mean average precision by averaging of average precision across classes. On the plot, we can select the working point, the precision, and the recall are best suited for the task at hand. Know that for production algorithm, the point is selected so that precision is closer one by sacrificing the recall. Sometimes, another charting metric are used to measure the quality of the detector. It is the plot of a miss rate with false detections per image? The curve is constructed similar to the precision-recall curve by varying the threshold on the detection score. For the object detection, the creation of ground truth annotation is very important. It has been demonstrated in the recent papers that annotations are usually different across datasets and annotators. **This lead to significant error in training of detector and its evaluation**. A lot of objects can be missed in ground truth data. This is especially true for the small objects or objects with very similar appearance to the target.

For example, if you want to detect faces, then detectors can detect not only real faces but for example prints on clothes or portraits. It is much better to annotate all such objects, and then mark some of them as objects to ignore during training and evaluation. It can also use several annotators for verification of annotation. The higher localization precision you want to obtain, the higher are usually the requirements on annotation precision. It is very important to have strict annotation protocol with clear definition how objects should be marked in difficult cases. For example, in a recent city person data set, each pedestrian should be marked with a line from toe to the head, to the middle point between the legs, and then bounding boxes width to height aspect ratio is placed on top of this annotation. Some objects are more important for us than other, so detection of such object classes as faces, pedestrian, or cars has a lot of practical applications. Thus, a lot of algorithms has been proposed for the detection of specific class of object. Such detectors reach top performance on their classes compared to the multi-class detectors. Practical multi-class detectors repeat only this development of deep learning methods. There are several peripheral data sets for multi-class detection. ImageNet is the first example of such data sets. It has objects of 1000 classes same as classification on ImageNet. Each image has annotation of one class and at least one bounding box. There are 800 training images per class. Detector should produce five guesses per image. Detection on the is correct if at least one guess has correct class and the corresponding bounding box is close to correct.

#### e. Effects of filter size on convolution

### Filter size considerations

- <u>Small</u> field-of-view → accurate <u>localization</u>
- <u>Large</u> field-of-view → context <u>assimilation</u>
- Effective filter size increases (enlarge the field-of-view of filter)

 $\begin{array}{ll} n_o: \ k \times k & \rightarrow & n_a: \left(k + (k-1)(r-1)\right) \times \left(k + (k-1)(r-1)\right) \\ n_o: \ \text{original convolution kernel size} \\ n_a: \ \text{atrous convolution kernel size} \\ r: \ \text{rate} \end{array}$ 

- However, we take into account only the non-zero filter values:
  - Number of filter parameters is the same
  - Number of operations per position is the same







Usage of convolution III : 2D filtering 7x7 Laplacian of Gaussian kernel





0.3

0.4

-0.7

-1.3

-0.7

0.4

0.3

0.2

0.11

-0.3

-0.7

-0.3

0.11

0.2

0.09

0.13

0.11

0.4

0.11

0.13

0.09

0.2

-0.3

-0.7

-0.3

0.11

0.2

0.11





10/22/2019

# Usage of convolution IV : 2D filtering

0.13

0.11

0.4

0.11

0.13

0.09

- Seeking for a known patter
- Large convolution kernel is applied •
- Kernel size is equivalent with the size of the sought pattern •





Scale variant ٠



Filter responeded with a strong white peek in the matching position

oszinten nem tudom mi kene meg ide..
#### No. 15 Topic

#### a. Rosenblatt perceptron training algorithm

#### The learning algorithm: Recursive algorithm

We have to develop a recursive algorithm called learning, which can learn the weight step by step, based on observing

- the <u>(i) input,</u>
- the (ii) weight vector,
- the (iii) desired output, and
- the <u>(iv) actual output</u> of the system.
- This can be described formally as follows:

 $\mathbf{w}(k+1) = \Psi(\mathbf{x}(k), \mathbf{w}(k), d(k), y(k)) \rightarrow \mathbf{w}_{opt}$ 

#### The learning algorithm: Recursive steps

1. Initialization.

Set w(0)=0 or w(0)=rand

- → 2. Activation.
  - Select a  $\mathbf{x}_k \rightarrow \mathbf{d}_k$  pair
  - 3. Computation of actual response  $y(k) = sign(w^{T}(k)x(k))$
  - 4. Adaptation of the weight vector  $\mathbf{w}(k+1) = \Psi(\mathbf{x}(k), \mathbf{w}(k), d(k), y(k))$
  - 5. Continuation

Until all responses of the perceptron are OK

### Learning

- The Rosenblatt algorithm is inapplicable,
  - the error and desired output in the hidden layers of the FFNN is unknown
- Someway the error of the whole network has to be distributed to the internal neurons, in a feedback way



Forward propagation of function signals and back-propagation of errors signals

#### The learning algorithm: Perceptron Learning Algorithm

- In a more ambitious way it can be called intelligent, because
  - perceptron can learn through examples (adapt),
  - even the function parameters are fully hidden.
- Perceptron learning was introduced by Frank Rosenblatt 1958
  - Built a 20x20 image sensor
  - With analog perceptron
  - 400 weights controlled by electromotors



b. Back-propagation -> tok random lasd No. 14 Topic b.

### **Batch Normalization**

- Done on layer level like softmax
- Training:
- Done on minibatch level Inferencing:
  - Do the normalization with the precalculated parameters of the entire training set
- Batch normalization is differenciable via chain rule
  - Back propagation can be applied for batch normalized layers
- Rewriting the normalization using probability terms:





**Input:** Values of x over a mini-batch:  $\mathcal{B} = \{x_{1...m}\};$ 

bias

Parameters to be learned:  $\gamma$ ,  $\beta$ 

## **Back propagation**

- We have seen last time how to calculate the gradient in a multilayer fully connected network using back propagation
  - The introduced method was based on gradient descent method
- However, being able to calculate gradient, we might select any of the above methods, which leads to orders of magnitude faster convergence

#### c. Curse of dimensionality

- What is it? •
  - A name for various problems that arise when analyzing data in high dimensional space.
  - Dimensions = independent features in ML
    - Input vector size (different measurements, or number of pixels in an image)
  - Occurs when d (# dimensions) is large in relation to n (number of samples).
- ٠ Real life examples:
  - Genomics
    - We have ~20k genes, but disease sample sizes are often in the 100s or 1000s.

# So what is this curse?

- Sparse data:
  - When the dimensionality d increases, the volume of the space increases so fast that the available data becomes **sparse**, **i.e.** a few points in a large **space**
  - Many features are not balanced, or are 'rarely occur' sparse features
- Noisy data: More features can lead to increased noise → it is harder to find the true signal
- Less clusters: Neighborhoods with fixed k points are less concentrated as d increases.
- **Complex features**: High dimensional functions tend to have more complex features than low-dimensional functions, and hence harder to estimate

### Data becomes sparse as dimensions increase

• A sample that maps 10% of the 1x1 squares in 2D represent only 1% of the 1x1x1 cubes in 3D



• There is an exponential increase in the search-space

#### Data sample number increase to avoid sparsity

- e.g. 10 observations /dimension
  - 1D: 10 observations
  - 2D: 100 observations
  - 3D: 1000 observations

```
- ...
```

11/19/2019



## Curse of dim - Running complexity

- Many data points (labeled measurements) are needed •
- Complexity (running time) increase with dimension d •
- A lot of methods have at least  $O(n^*d^2)$  complexity, where **n** is the number of samples
- As **d** becomes large, this complexity becomes very costly. • - Compute = \$







S<sub>100</sub>=1

#### Distances in high dimension

- Assume, we have a unit side (2D) square, what we divided to 100 equal small squares
  - Calculate the ratio of the largest distance in a small square and the largest distance of the big square (in 2D)
- Assume, we have a unit side 100D cube, • what we divided to 100 equal small 100D cubes
  - Calculate the ratio Ratio of the largest distance in a small cube and the largest distance of the big cube (in 100D)
  - The average nearest neighbor distance is 95% of the largest distance!!!
  - Euclidian distance becomes meaningless, most two points are "far" from each others

$$D_{100} = \sqrt{100} = 10$$
  $d_{100} = \sqrt{100 * 0.95^2} = 9.5$   
 $R_{100} = \frac{d_{100}}{D_{100}} = 0.95$ 

s<sub>100</sub>=

13

 $\sqrt[100]{\frac{1}{100}} = 0.95$ 







## Curse of dim - Some mathematical (weird) effects



- sube for d=3:  $\frac{(\frac{\tau}{3})\pi r^3}{(2r)^3} \approx \frac{4r^3}{8r^3} \approx 0.5$
- Ratio between the volume of a sphere and a cube for d=3: (2r)<sup>3</sup> ≈ 8r<sup>3</sup>
   When d tends to infinity the volume of the sphere (this ratio) tends to zero

d	3	5	10	20	30	50
ratio	0.52	0.16	0.0025	2.5E-08	2.0E-14	1.5E-28

- Most of the data is in the corner of the cube
  - Thus, Euclidian distance becomes meaningless, most two points are "far" from each others
- Very problematic for methods such as k-NN classification or k-means clustering because most of the neighbors are equidistant

## The nearest neighbor problem in a sphere

- · Assume randomly distributed points in a sphere with a unit diameter
- The median of the nearest neighbors is *l*
- As dimension tends to infinity
  - The median of the nearest neighbors converges to 1



"The Curse of Dimensionality" by Raúl Rojas https://www.inf.fu-berlin.de/inst/agki/rojas\_home/documents/tutorials/dimensionality.pdf

11/19/2019

### How to calculate dimensionality?

G		feature vectors (x)						
bservations (		<b>x</b> <sub>1</sub>	<i>x</i> <sub>2</sub>	<b>X</b> 3	<i>X</i> <sub>4</sub>			
	<i>d</i> <sub>1</sub>	1	2	1	1			
	$d_2$	2	4	3.5	1			
	$d_3$	3	6	17	1			
0								

 How many dimensions does the data intrinsically have here? (How many independent coordinates?)

#### – Two!



x1 = ½ \* x2 (no additional information, correlated, not independent)
x4 is constant (carries no information at all!)



**X**3

## How to avoid the curse?



- Feature selection Choose only a subset of features
- Use algorithms that transform the data into a lower dimensional space (example PCA, t-SNE)
   \*Both methods often result in information loss
- Less is More
  - In many cases the information that is lost by discarding variables is made up for by a more accurate mapping/sampling in the lower-dimensional space



11/19/2019

#### d. SqueezeNet



Figure 2. The SqueezeNet architecture

Ø,

17



Figure 2. The SqueezeNet architecture

#### SqueezeNext



In this arhcitecture depths are squeezed before each operation



#### e. Back-propagation and gradient-based optimizers

Simple gradient based optimizers:

• 1st and 2nd order optimizers

#### Quora - TL;DR:

The backpropagation algorithm is an instruction set for computing the gradient of a multi-variable function.

The Adam optimizer is a specialized gradient-descent algorithm that uses the computed gradient, its statistics, and its historical values to take small steps in its opposite direction inside the input parameter space as a means of minimizing a function. It is used for optimization in neural network training.

In other words, the Adam optimizer would need to use an algorithm like the backpropagation algorithm to first compute the gradient of the function. Then the Adam optimizer would use this computation to perform gradient-descent in a specialized manner.

#### Másik jó oldalka:

https://towardsdatascience.com/types-of-optimization-algorithms-used-in-neural-networks-an d-ways-to-optimize-gradient-95ae5d39529f

Optimization Algorithm falls in 2 major categories -

- First Order Optimization Algorithms These algorithms minimize or maximize a Loss function E(x) using its Gradient values with respect to the parameters. Most widely used First order optimization algorithm is Gradient Descent. The First order derivative tells us whether the function is decreasing or increasing at a particular point. First order Derivative basically give us a line which is Tangential to a point on its Error Surface.
- 2. <u>Second Order Optimization Algorithms</u> Second-order methods use the second order derivative which is also called Hessian to minimize or maximize the Loss function. The Hessian is a Matrix of Second Order Partial Derivatives. Since the second derivative is costly to compute, the second order is not used much. The second order derivative tells us whether the first derivative is increasing or decreasing which hints at the function's curvature. Second Order Derivative provide us with a quadratic surface which touches the curvature of the Error Surface.