

DATA SCIENCE 2

March 27, 2022

Data Science 2

Faculty of Mathematics and Physics

AGENDA



Neural Networks Introduction

Gradient Based Optimization

Activation functions

Regularization

Convolutional neural networks

Convolutional neural network case studies

Classical approaches to NLP

Sequence models

Word embeddings

NEURAL NETWORKS INTRODUCTION

TARAN

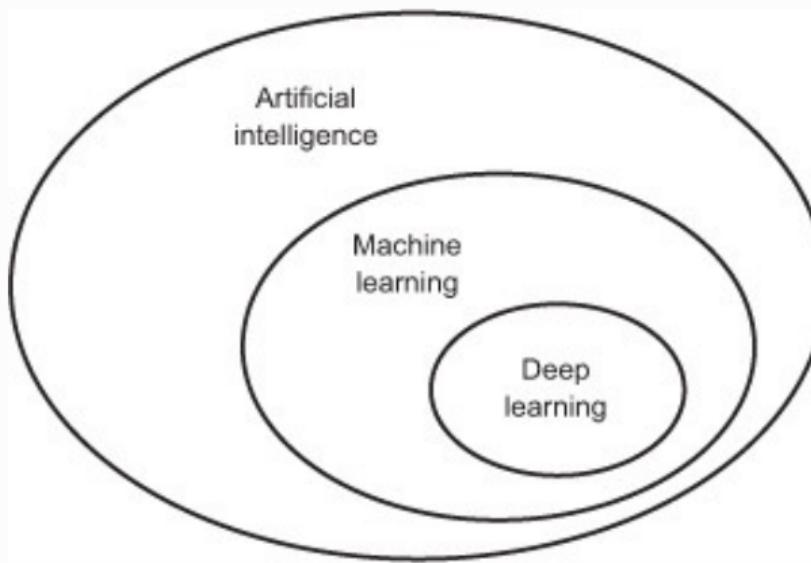


ADVISORY IN DATA & ANALYTICS

INTRODUCTION



- ▶ What is artificial intelligence, machine learning and deep learning?
- ▶ And how do they relate?



ARTIFICIAL INTELLIGENCE



Definition (Artificial intelligence)

Artificial intelligence is an effort to automate intellectual tasks normally performed by humans.

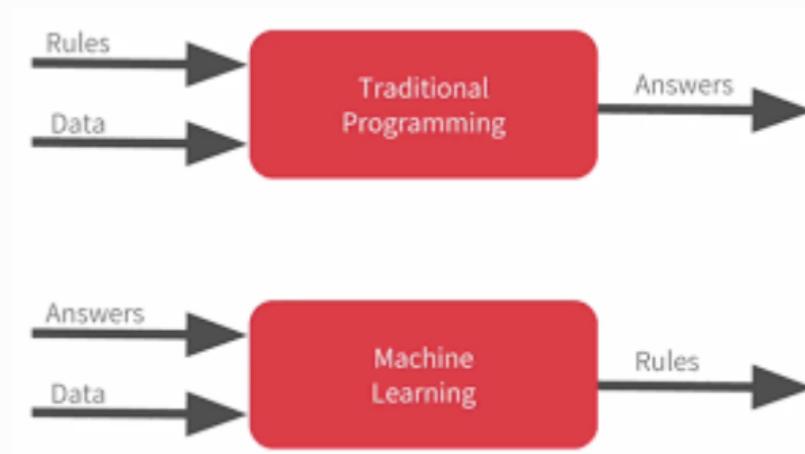
- ▶ AI is a superset of machine learning.
- ▶ AI includes also algorithms, that are not based on any learning - you can learn computer to play chess by defining set of rules to be obeyed.
- ▶ AI was born around 1950s.
- ▶ **Symbolic AI** - at the beginning it was believed that human level AI can be achieved by sufficiently large set of rules.

MACHINE LEARNING



Definition (Machine Learning)

The use and development of computer systems that are able to learn and adapt without following explicit instructions, by using algorithms and statistical models to analyse and draw inferences from patterns in data.



DEEP LEARNING



Definition (Deep learning)

Deep learning is a class of machine learning algorithms that uses multiple layers to progressively extract higher-level features from the raw input. For example, in image processing, lower layers may identify edges, while higher layers may identify the concepts relevant to a human such as digits or letters or faces.

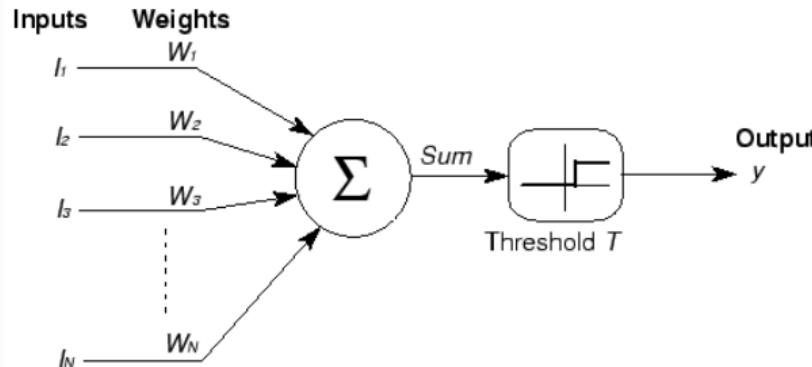
Deep learning achievements:

- ▶ Near-human level image classification
- ▶ Near-human level speech recognition
- ▶ Near-human level handwriting transcription
- ▶ Improved machine translation
- ▶ Improved text-to-speech conversion
- ▶ Near-human level autonomous driving
- ▶ Improved ad targeting
- ▶ Improved search results on the web
- ▶ Superhuman Go playing

PERCEPTRON



In 1958, Frank Rosenblatt introduced an idea of perceptron and called it Mark I Perceptron.



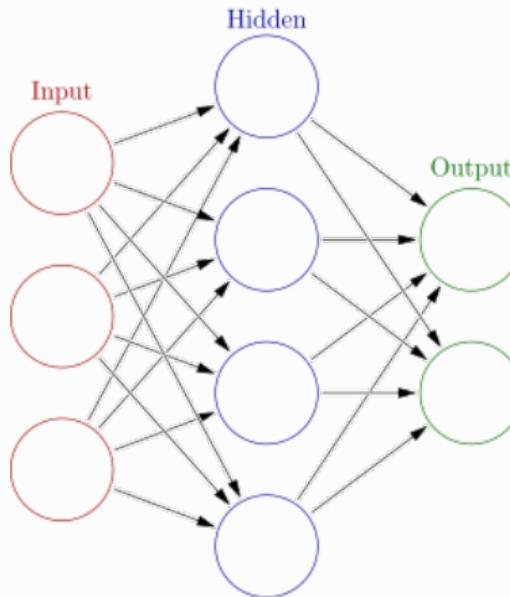
Weights of Mark I Perceptron are 'learnt' through successively passed inputs, while minimizing the difference between desired and actual output.

<https://playground.tensorflow.org/>

MULTI-LAYER PERCEPTRON



Mark I Perceptron is not capable of solving XOR problem. To solve XOR problem, we need to use multi-layer perceptron (MLP) neural network.





INPUT DATA

Input data for neural networks are usually passed in form of multi-dimensional Numpy arrays, also called **tensors**:

- ▶ Scalars (0D tensors)
- ▶ Vectors (1D tensors)
- ▶ Matrices (2D tensors)
- ▶ 3D and higher dimensional tensors

Tensor is defined by tree key attributes:

- ▶ Number of dimensions
- ▶ Shape
- ▶ Data type

For tensors, dimension is often called an *axis*. First axis is usually allocated for sample axis. In neural network models data are not processed all at once, but in batches.

GRADIENT BASED OPTIMIZATION

TARAN



ADVISORY IN DATA & ANALYTICS

GRADIENT BASED OPTIMIZATION



Each neuron transforms input data using following equation:

$$\text{out} = f \left(\sum_{j=1}^N w_j x_j + b \right) \quad (1)$$

where

- ▶ f ... activation function
- ▶ w_i ... weight assigned to a connection coming from i -th input.
- ▶ x_i ... value of i -th input
- ▶ b ... bias

In multilayer layout weights are represented by weight matrix \mathbf{W} and biases by bias vector \mathbf{b} .

BACKPROPAGATION

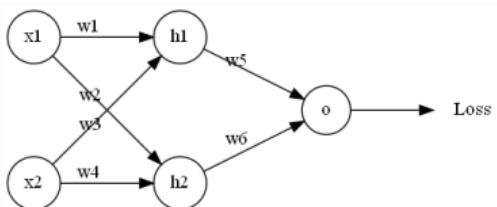


Backpropagation allows to efficiently update weights in the network.

Backpropagation is based on chain rule:

$$\frac{\partial f(g(x))}{\partial x} = \frac{\partial f(x)}{\partial g(x)} \cdot \frac{\partial g(x)}{\partial x} \quad (2)$$

Lets have following network:



To update weight w_1 we calculate:

$$\frac{\partial L}{\partial w_1} = \frac{\partial L}{\partial o} \cdot \frac{\partial o}{\partial h_1} \cdot \frac{\partial h_1}{\partial w_1} \quad (3)$$

TRAINING NEURAL NETWORK



Algorithm (NN training)

Repeat until stopping criterion is met:

1. Sample a minibatch of n observations from training sample.
2. Feed forward step - propagate each observation through neural net and calculate outputs
3. Calculate derivatives of loss function with respect to trainable parameters
4. Backpropagation step - update trainable parameters based on gradient averaged over batch

$$\theta^{(t+1)} = \theta^{(t)} - \alpha \cdot \frac{1}{n} \sum_{i=1}^n \nabla_{\theta^{(t)}} L(\text{output}_i, \text{label}_i) \quad (4)$$

- ▶ **Minibatch** - Subsample to be used for one update of weights (forward and backward pass)
- ▶ **Epoch** - All the training data was used exactly once to update weights

GRADIENT DESCENT

SGD WITH MOMENTUM



Algorithm (SGD with momentum)

Repeat until stopping criteria is met:

- Randomly select m samples from training data.
- $\mathbf{g}_t = \frac{1}{m} \nabla_{\theta} \sum_i L(f(\mathbf{x}^{(i)}; \theta_{t-1}), y^{(i)})$
- $\mathbf{v}_t = \beta \mathbf{v}_{t-1} - \alpha \mathbf{g}_t$
- $\theta_t = \theta_{t-1} + \mathbf{v}_t$

GRADIENT DESCENT

SGD WITH MOMENTUM



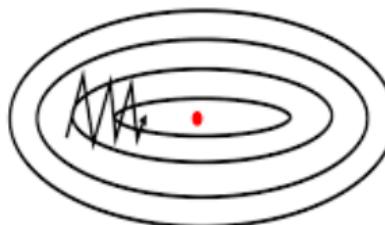
Algorithm (SGD with momentum)

Repeat until stopping criteria is met:

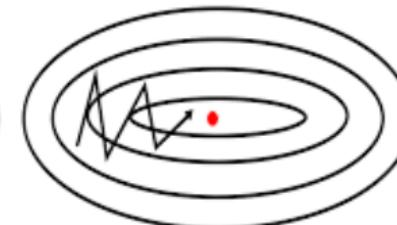
- Randomly select m samples from training data.
- $\mathbf{g}_t = \frac{1}{m} \nabla_{\theta} \sum_i L(f(\mathbf{x}^{(i)}; \theta_{t-1}), y^{(i)})$
- $\mathbf{v}_t = \beta \mathbf{v}_{t-1} - \alpha \mathbf{g}_t$
- $\theta_t = \theta_{t-1} + \mathbf{v}_t$

► Good value for β is 0.9

SGD without momentum



SGD with momentum



GRADIENT DESCENT

NESTEROV ACCELERATED GRADIENT DESCENT



Algorithm (Nesterov accelerated GD)

Repeat until stopping criteria is met:

- ▶ $\theta_{t_0} = \theta_{t-1} + \beta v_{t-1}$
- Randomly select m samples from training data.
- $g_t = \frac{1}{m} \nabla_{\theta} \sum_i L(f(\mathbf{x}^{(i)}; \theta_{t_0}), y^{(i)})$
- $v_t = \beta v_{t-1} - \alpha g_t$
- $\theta_t = \theta_{t_0} - \alpha g_t$

GRADIENT DESCENT

NESTEROV ACCELERATED GRADIENT DESCENT



Algorithm (Nesterov accelerated GD)

Repeat until stopping criteria is met:

- ▶ $\theta_{t_0} = \theta_{t-1} + \beta v_{t-1}$
 - Randomly select m samples from training data.
 - $g_t = \frac{1}{m} \nabla_{\theta} \sum_i L(f(\mathbf{x}^{(i)}; \theta_{t_0}), y^{(i)})$
 - $v_t = \beta v_{t-1} - \alpha g_t$
 - $\theta_t = \theta_{t_0} - \alpha g_t$
-
- ▶ We know we will use momentum, so we can calculate gradient after applying momentum, rather than before.

GRADIENT DESCENT

ADAGRAD



Algorithm (Adagrad)

Repeat until stopping criteria is met:

- Randomly select m samples from training data.
- $\mathbf{g}_t = \frac{1}{m} \nabla_{\theta} \sum_i L(f(\mathbf{x}^{(i)}; \boldsymbol{\theta}_{t-1}), y^{(i)})$
- $\mathbf{r}_t = \mathbf{r}_{t-1} + \mathbf{g}_t^2$
- $\boldsymbol{\theta}_t = \boldsymbol{\theta}_{t-1} - \frac{\alpha}{\sqrt{\mathbf{r}_t + \varepsilon}} \mathbf{g}_t$

GRADIENT DESCENT

ADAGRAD



Algorithm (Adagrad)

Repeat until stopping criteria is met:

- Randomly select m samples from training data.
- $\mathbf{g}_t = \frac{1}{m} \nabla_{\theta} \sum_i L(f(\mathbf{x}^{(i)}; \boldsymbol{\theta}_{t-1}), y^{(i)})$
- $\mathbf{r}_t = \mathbf{r}_{t-1} + \mathbf{g}_t^2$
- $\boldsymbol{\theta}_t = \boldsymbol{\theta}_{t-1} - \frac{\alpha}{\sqrt{\mathbf{r}_t + \varepsilon}} \mathbf{g}_t$

- ▶ Increase learning rate for parameter related to infrequent features
- ▶ Decrease learning rate for parameter related to frequent features
- ▶ Usually $\varepsilon = 10^{-8}$

GRADIENT DESCENT

RMSProp



Algorithm (RMSProp)

Repeat until stopping criteria is met:

- Randomly select m samples from training data.
- $\mathbf{g}_t = \frac{1}{m} \nabla_{\theta} \sum_i L(f(\mathbf{x}^{(i)}; \boldsymbol{\theta}_{t-1}), y^{(i)})$
- $\mathbf{r}_t = \beta \mathbf{r}_{t-1} + (1 - \beta) \mathbf{g}_t^2$
- $\boldsymbol{\theta}_t = \boldsymbol{\theta}_{t-1} - \frac{\alpha}{\sqrt{\mathbf{r}_t + \varepsilon}} \mathbf{g}_t$

GRADIENT DESCENT

RMSProp



Algorithm (RMSProp)

Repeat until stopping criteria is met:

- Randomly select m samples from training data.
- $\mathbf{g}_t = \frac{1}{m} \nabla_{\theta} \sum_i L(f(\mathbf{x}^{(i)}; \boldsymbol{\theta}_{t-1}), y^{(i)})$
- $\mathbf{r}_t = \beta \mathbf{r}_{t-1} + (1 - \beta) \mathbf{g}_t^2$
- $\boldsymbol{\theta}_t = \boldsymbol{\theta}_{t-1} - \frac{\alpha}{\sqrt{\mathbf{r}_t + \varepsilon}} \mathbf{g}_t$

- ▶ Trying to fix monotonically decreasing learning rate in Adagrad
- ▶ Usually $\varepsilon = 10^{-8}$
- ▶ Usually $\beta = 0.9$

GRADIENT DESCENT

ADAM



Adam = Adaptive Moment Estimation

Algorithm (Adam)

Set $s_0 = 0, r_0 = 0$

Repeat until stopping criteria is met:

- Randomly select m samples from training data.
- $\mathbf{g}_t = \frac{1}{m} \nabla_{\theta} \sum_i L(f(\mathbf{x}^{(i)}; \theta_{t-1}), y^{(i)})$
- $\mathbf{s}_t = \beta_1 \mathbf{s}_{t-1} + (1 + \beta_1) \mathbf{g}_t$
- $\mathbf{r}_t = \beta_2 \mathbf{r}_{t-1} + (1 + \beta_2) \mathbf{g}_t^2$
- $\hat{\mathbf{s}}_t = \frac{\mathbf{s}_t}{1 - \beta_1^t}$
- $\hat{\mathbf{r}}_t = \frac{\mathbf{r}_t}{1 - \beta_2^t}$
- $\theta_t = \theta_{t-1} - \frac{\alpha}{\sqrt{\hat{\mathbf{r}}_t + \varepsilon}} \hat{\mathbf{s}}_t$

GRADIENT DESCENT

ADAM

Adam = Adaptive Moment Estimation

Algorithm (Adam)

Set $s_0 = 0, r_0 = 0$

Repeat until stopping criteria is met:

- Randomly select m samples from training data.
- $\mathbf{g}_t = \frac{1}{m} \nabla_{\theta} \sum_i L(f(\mathbf{x}^{(i)}; \theta_{t-1}), y^{(i)})$
- $\mathbf{s}_t = \beta_1 \mathbf{s}_{t-1} + (1 + \beta_1) \mathbf{g}_t$
- $\mathbf{r}_t = \beta_2 \mathbf{r}_{t-1} + (1 + \beta_2) \mathbf{g}_t^2$
- $\hat{\mathbf{s}}_t = \frac{\mathbf{s}_t}{1 - \beta_1^t}$
- $\hat{\mathbf{r}}_t = \frac{\mathbf{r}_t}{1 - \beta_2^t}$
- $\theta_t = \theta_{t-1} - \frac{\alpha}{\sqrt{\hat{\mathbf{r}}_t + \varepsilon}} \hat{\mathbf{s}}_t$

- ▶ Default values: $\beta_1 = 0.9, \beta_2 = 0.999, \alpha = 0.001$
- ▶ Usually $\varepsilon = 10^{-8}$

GRADIENT DESCENT

ALGORITHMS OVERVIEW



- ▶ SGD with Nesterov Momentum
 - momentum
- ▶ AdaGrad (Adaptive Gradient Algorithm)
 - adaptive learning rate
- ▶ RMSProp (Root Mean Square Propagation)
 - momentum
 - adaptive learning rate
- ▶ Adam (Adaptive Moment Estimation)
 - momentum
 - adaptive learning rate

<https://awesomedl.com/project/Jaewan-Yun/optimizer-visualization>

<https://keras.io/api/optimizers/>

LEARNING RATE SCHEDULES



Even for algorithms with adaptive learning rates, further fine-tuning of learning rate can improve performance. Decreasing learning rate can be performed each batch/epoch/several epochs.

- ▶ Exponential decay
$$\eta_t = \eta_0 \cdot c^t$$
- ▶ Polynomial decay
 - Inverse time decay
$$\eta_t = \eta_0 \cdot \frac{1}{t}$$
 - Inverse-square decay
$$\eta_t = \eta_0 \cdot \frac{1}{\sqrt{t}}$$
- ▶ Cosine decay, restarts, ...

https://www.tensorflow.org/api_docs/python/tf/keras/optimizers/schedules/

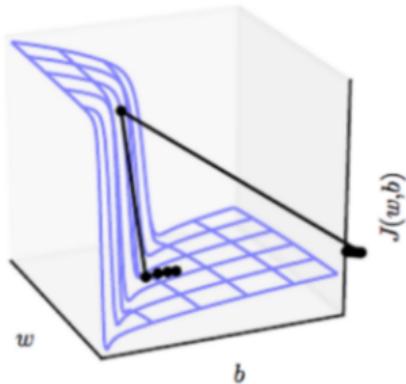
GRADIENT CLIPPING



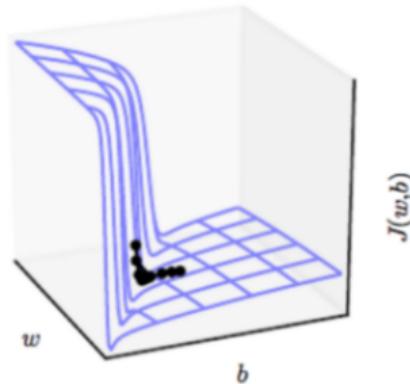
Algorithm (Gradient clipping)

- $\mathbf{g} = \nabla_{\theta} L(f(\mathbf{x}; \theta), y)$
- if $\|\mathbf{g}\| \geq \text{threshold}$ then $\mathbf{g} = \text{threshold}$ end if

Without clipping



With clipping



ACTIVATION FUNCTIONS

TARAN



ADVISORY IN DATA & ANALYTICS

ACTIVATION FUNCTIONS



Each neuron transforms input data using following equation:

$$\text{out} = \mathbf{f} \left(\sum_{j=1}^N w_j x_i + b \right) \quad (5)$$

- ▶ Activation function adds non-linearity at each neuron
- ▶ It helps to control neuron output values range
- ▶ Desired features:
 - Differentiability (required)
 - Know analytical solution of first derivative
 - Low computational expense

https://www.tensorflow.org/api_docs/python/tf/keras/activations

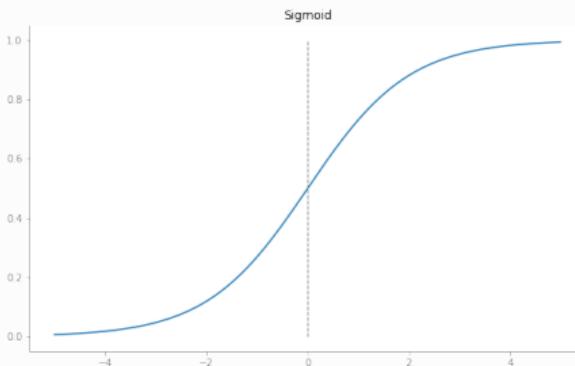
ACTIVATION FUNCTIONS

SIGMOID



Definition:

$$\sigma(x) = \frac{1}{1 + e^{-x}} \quad (6)$$



Derivative:

$$\frac{d\sigma(x)}{dx} = \sigma(x)(1 - \sigma(x)) \quad (7)$$

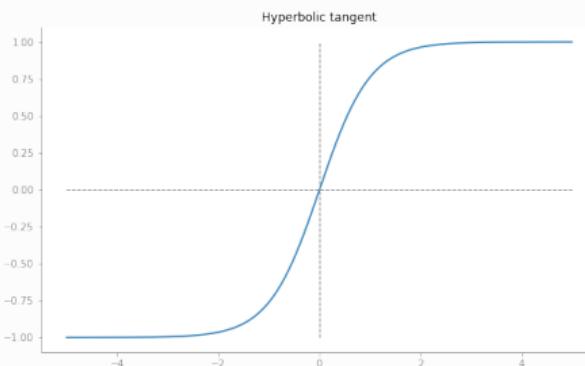
ACTIVATION FUNCTIONS

HYPERBOLIC TANGENT



Definition:

$$\tanh(x) = \frac{e^{2x} - 1}{e^{2x} + 1} \quad (8)$$



Derivative:

$$\frac{d\tanh(x)}{dx} = 1 - \tanh^2(x) \quad (9)$$

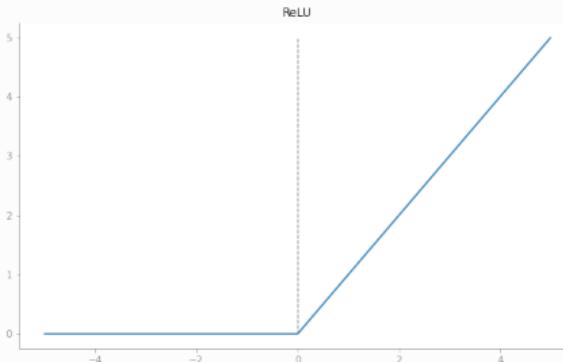
ACTIVATION FUNCTIONS

RECTIFIED LINEAR UNIT



Definition:

$$\text{ReLU}(x) = \max(0, x) \quad (10)$$



Derivative:

$$\frac{d\text{ReLU}(x)}{dx} = \begin{cases} 1 & \text{if } x > 0 \\ NaN & \text{if } x = 0 \\ 0 & \text{if } x < 0 \end{cases} \quad (11)$$

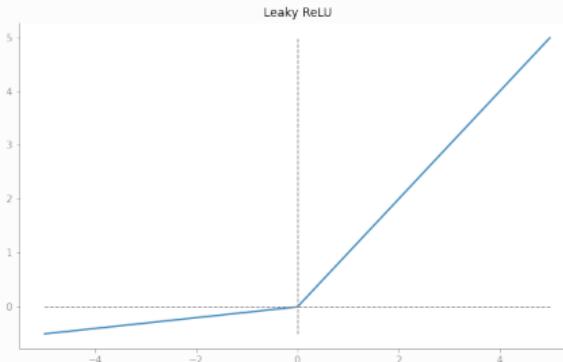
ACTIVATION FUNCTIONS

LEAKY RELU



Definition:

$$\text{LeakyReLU}(x) = \max(\alpha x, x) \quad (12)$$



Derivative:

$$\frac{d\text{LeakyReLU}(x)}{dx} = \begin{cases} 1 & \text{if } x > 0 \\ NaN & \text{if } x = 0 \\ \alpha & \text{if } x < 0 \end{cases} \quad (13)$$

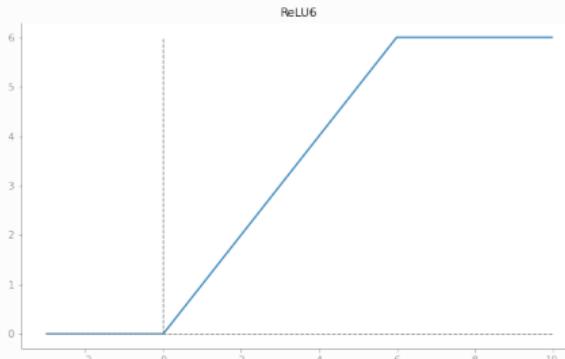
ACTIVATION FUNCTIONS

ReLU6



Definition:

$$\text{ReLU6}(x) = \min(\max(0, x), 6) \quad (14)$$



Derivative:

$$\frac{d\text{ReLU6}(x)}{dx} = \begin{cases} 1 & \text{if } x \in (0, 6) \\ NaN & \text{if } x = 0 \text{ or } x = 6 \\ 0 & \text{if } x < 0 \text{ or } x > 6 \end{cases} \quad (15)$$

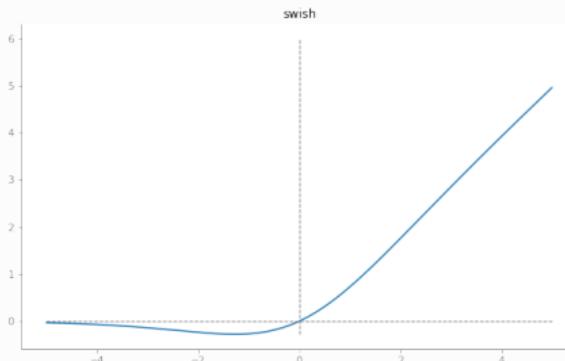
ACTIVATION FUNCTIONS

SWISH



Definition:

$$\text{swish}(x) = \beta x \sigma(x) \quad (16)$$



Derivative:

$$\frac{d\text{swish}(x)}{dx} = \text{swish}(x) + \sigma(x)(\beta - \text{swish}(x)) \quad (17)$$

ACTIVATION FUNCTIONS

ACTIVATION ON OUTPUT LAYER



Common activations on output layer:

- ▶ None - linear regression if there are no hidden layers
- ▶ $\sigma(x)$ - binary classification

$$\sigma(x) = \frac{1}{1 + e^{-x}} \quad (18)$$

- ▶ softmax - multi-class classification

$$\text{softmax}(\mathbf{x})_i = \frac{e^{x_i}}{\sum_j e^{x_j}} \quad (19)$$

VANISHING AND EXPLODING GRADIENT



- ▶ Vanishing gradient is a consequence of chain rule used in backpropagation algorithm.
- ▶ Activation functions such as sigmoid or hyperbolic tangent has derivative in interval $(0,1)$
- ▶ With inappropriate set up derivatives of weights in former layers might **vanish** because of multiplying values lower then 1 in magnitude many times
- ▶ If, for instance, the initial loss is big, the gradient might **explode**

VANISHING AND EXPLODING GRADIENT



- ▶ Vanishing gradient is a consequence of chain rule used in backpropagation algorithm.
- ▶ Activation functions such as sigmoid or hyperbolic tangent has derivative in interval $(0,1)$
- ▶ With inappropriate set up derivatives of weights in former layers might **vanish** because of multiplying values lower than 1 in magnitude many times
- ▶ If, for instance, the initial loss is big, the gradient might **explode**

Symptoms of vanishing/exploding gradient:

Vanishing gradient

- Model improves very slowly
- Weights closer to output layer changes more dynamically

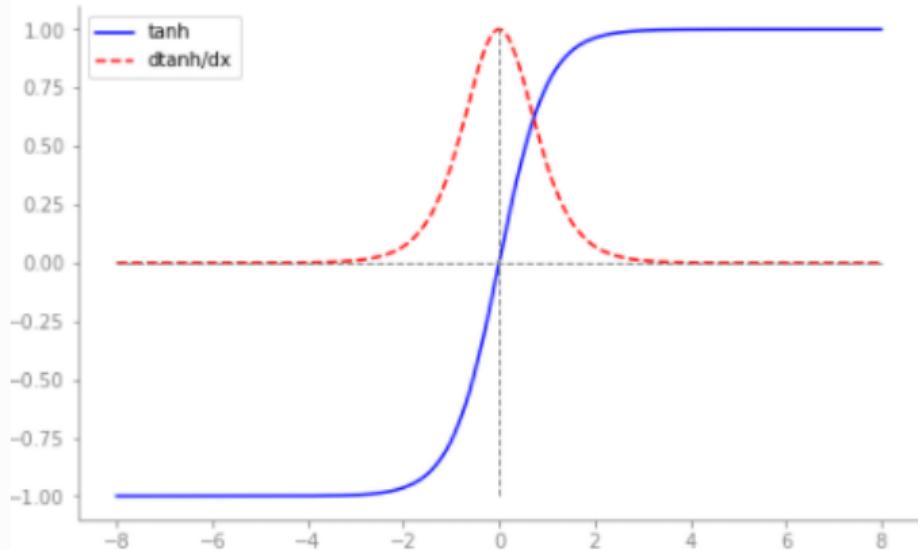
Exploding gradient

- Large changes in loss
- Model loss is NaN
- Model weights grow exponentially
- Model weights are NaN

SATURATING NON-LINEARITIES



With large input to neuron, derivative of saturating activation functions is close to 0



WEIGHTS INITIALIZATION



- ▶ In general practice biases are initialized with 0.
- ▶ Initialization of weights by 0 is not a good idea, since such weights will never break the symmetry inside one layer.
- ▶ Glorot initialization is trying to fix the variance of activations as well as gradients across layers:

$$W_{\text{Glorot}} \sim U \left[-\sqrt{\frac{6}{m+n}}, \sqrt{\frac{6}{m+n}} \right] \quad (20)$$

where n is number of neurons in previous layer and m is a number of neurons in next layer.

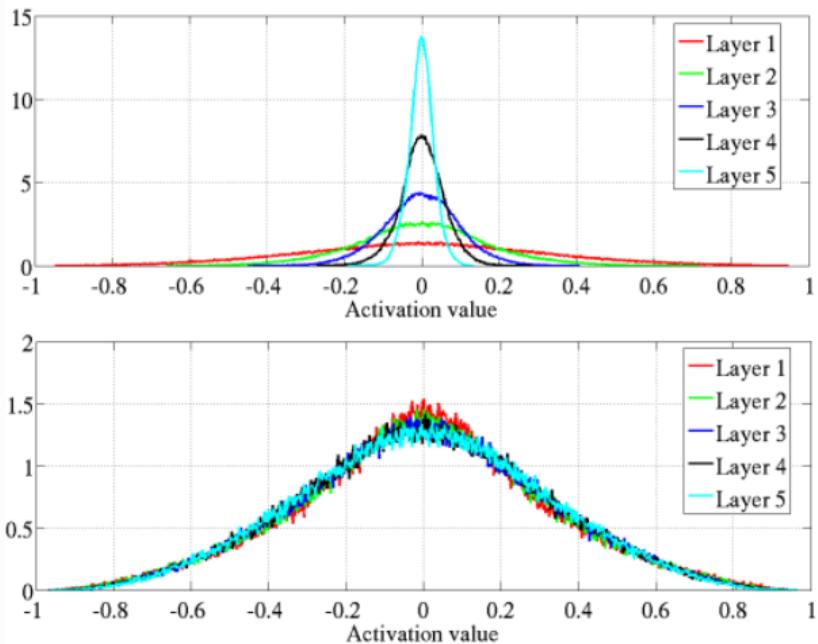
- ▶ For ReLU activation, HeUniform might be more suitable:

$$W_{\text{He}} \sim U \left[-\sqrt{\frac{6}{n}}, \sqrt{\frac{6}{n}} \right] \quad (21)$$

where n is a number of input units.

https://www.tensorflow.org/api_docs/python/tf/keras/initializers/

WEIGHTS INITIALIZATION



Source: <http://proceedings.mlr.press/v9/glorot10a/glorot10a.pdf>

REGULARIZATION

TARAN



ADVISORY IN DATA & ANALYTICS

REGULARIZATION

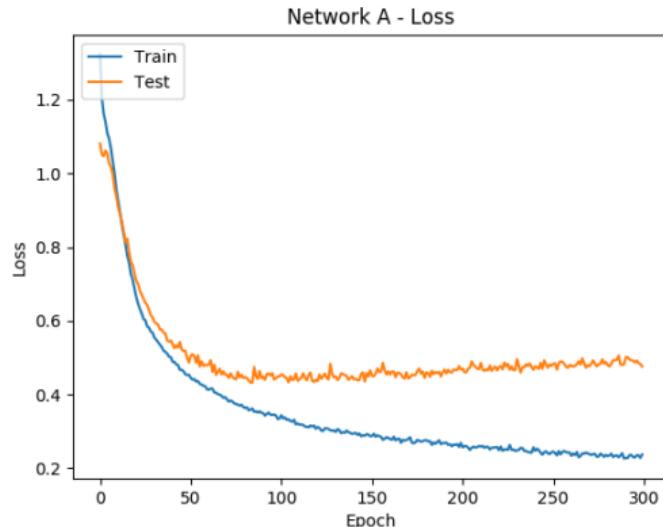


Regularization allows to control generalization error of model.

- ▶ Early stopping
- ▶ L2, L1 regularization
- ▶ Dataset augmentation
- ▶ Ensembling
- ▶ Dropout
- ▶ Label smoothing

REGULARIZATION

EARLY STOPPING



- ▶ Loss on train set will always decrease
- ▶ Stop training once loss on test set stop decreasing

REGULARIZATION

L2, L1 REGULARIZATION

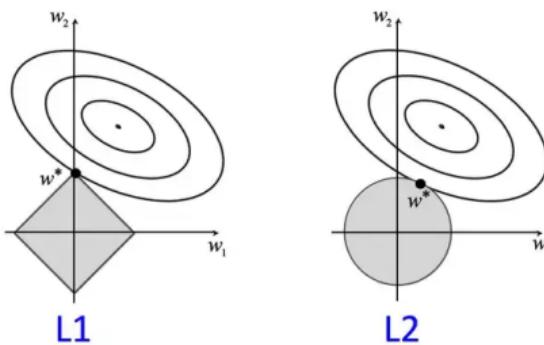
Penalization of large weights. Distribute the model strength amongst more ways in neural net, rather than emphasize one specific way.

- ▶ L2 regularization

$$L_{L2}(\theta; \mathbf{X}) = L(\theta; \mathbf{X}) + \lambda \|\theta\|_2^2 \quad (22)$$

- ▶ L1 regularization

$$L_{L1}(\theta; \mathbf{X}) = L(\theta; \mathbf{X}) + \lambda \|\theta\|_1 \quad (23)$$



REGULARIZATION

DATA AUGMENTATION



Modifying input data can force model to focus on general patterns rather than focusing on details, which can lead to overfitting.

- ▶ Image processing
 - Translation
 - Horizontal flips
 - Scaling
 - Rotation
 - Colour adjustment
 - Mixup - two images are combined using their weighted combination
- ▶ Speech recognition
 - Additional noise
 - Frequency change

https://www.tensorflow.org/api_docs/python/tf/keras/preprocessing

REGULARIZATION

ENSEMBLING



Model ensembling is a technique to reduce generalization error by combining several models. Usually, model outputs are averaged.

Possible set ups:

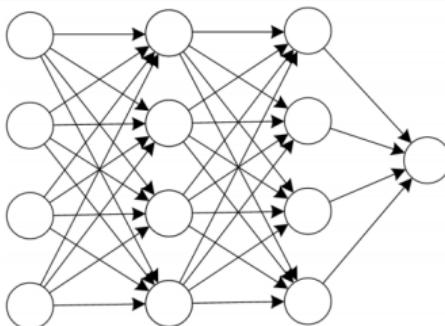
- ▶ Generate random training samples using sampling with replacement.
- ▶ Use different random weight initialization.
- ▶ Average model obtained at different time steps (one hour of training, one day of training, ...)

Ensembling often has high computational requirements.

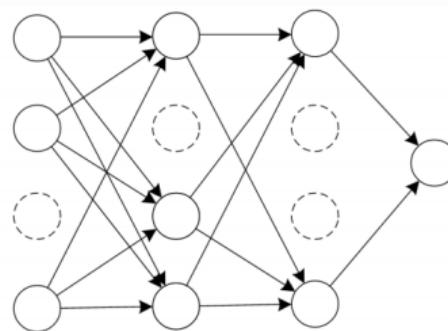
REGULARIZATION

DROPOUT

- ▶ Dropout is applied to a layer.
- ▶ Each neuron of the layer is randomly excluded with certain probability in each batch.
- ▶ Technically, dropout is performed by setting value of the neuron to 0.



(a) Standard Neural Network



(b) Network after Dropout



REGULARIZATION

DROPOUT

During inference, no neurons are dropped out. This means that neuron receives signal from more neurons than during training. Scaling should be performed:

- ▶ Reduce the activation signal during inference by the factor of $1 - p$, where p is a probability of dropout from previous layer.
- ▶ Enhance the activation signal during training by the factor of $\frac{1}{1-p}$

REGULARIZATION

LABEL SMOOTHING

In multi-class classification problem, when using softmax, the model is never satisfied with achieved results. If model prediction for the true category is 0.999, it is still trying to improve.

To overcome this, true label distribution is modified from $\mathbf{1}_{true}$ to

$$(1 - \alpha)\mathbf{1}_{true} + \alpha \frac{\mathbf{1}}{\text{number of classes}}$$



CONVOLUTIONAL NEURAL NETWORKS

TARAN



ADVISORY IN DATA & ANALYTICS

CONVOLUTIONAL NEURAL NETWORKS

MOTIVATION



Convolutional neural networks bring three main advantages:

CONVOLUTIONAL NEURAL NETWORKS

MOTIVATION



Convolutional neural networks bring three main advantages:

- ▶ Local interactions
 - Neighbouring pixels in picture might share an information rather than those far from each other.
 - Neighbouring words in a sentence might have a special meaning.

CONVOLUTIONAL NEURAL NETWORKS

MOTIVATION



Convolutional neural networks bring three main advantages:

- ▶ Local interactions
 - Neighbouring pixels in picture might share an information rather than those far from each other.
 - Neighbouring words in a sentence might have a special meaning.
- ▶ Parameter sharing
 - This concept allow to reduce number of trainable parameters.

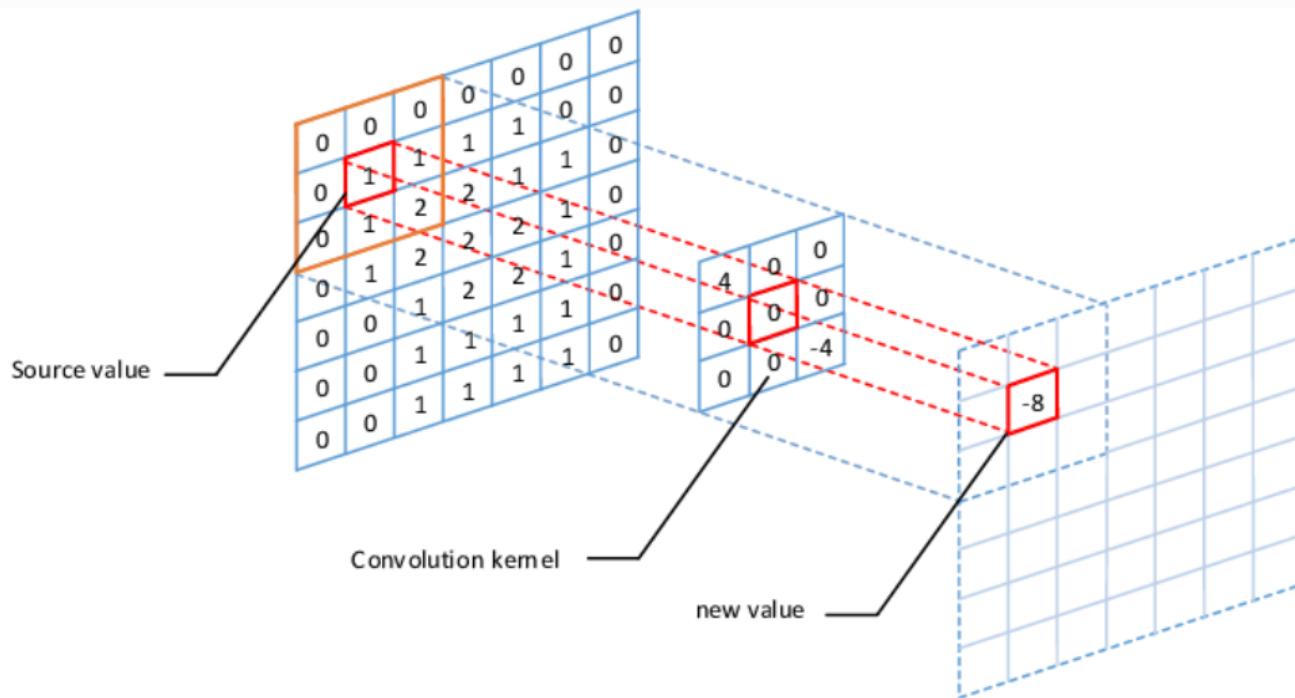
CONVOLUTIONAL NEURAL NETWORKS

MOTIVATION



Convolutional neural networks bring three main advantages:

- ▶ Local interactions
 - Neighbouring pixels in picture might share an information rather than those far from each other.
 - Neighbouring words in a sentence might have a special meaning.
- ▶ Parameter sharing
 - This concept allow to reduce number of trainable parameters.
- ▶ Shift invariance
 - A pattern should be recognized independently on its position.



CONVOLUTIONAL NEURAL NETWORKS

CONVOLUTION OPERATION



Convolution of two functions is defined as:

$$(f * g)(t) = \int_{-\infty}^{\infty} f(\tau)g(t - \tau)d\tau \quad (24)$$

CONVOLUTIONAL NEURAL NETWORKS

CONVOLUTION OPERATION

Convolution of two functions is defined as:

$$(f * g)(t) = \int_{-\infty}^{\infty} f(\tau)g(t - \tau)d\tau \quad (24)$$

In 2D it takes form:

$$(f * g)(t, u) = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} f(\tau, \xi)g(t - \tau, u - \xi)d\tau d\xi \quad (25)$$

CONVOLUTIONAL NEURAL NETWORKS

CONVOLUTION OPERATION

Convolution of two functions is defined as:

$$(f * g)(t) = \int_{-\infty}^{\infty} f(\tau)g(t - \tau)d\tau \quad (24)$$

In 2D it takes form:

$$(f * g)(t, u) = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} f(\tau, \xi)g(t - \tau, u - \xi)d\tau d\xi \quad (25)$$

In discrete domain we have:

$$(\mathbf{w} * \mathbf{x})_t = \sum_i w_i x_{t-i} \quad (26)$$

Or for 2D case:

$$(\mathbf{K} * \mathbf{I})_{i,j} = \sum_{m,n} \mathbf{K}_{m,n} \mathbf{I}_{i-m, j-n} \quad (27)$$

CONVOLUTIONAL NEURAL NETWORKS

CROSS-CORRELATION



CNN actually uses cross-correlation instead of convolution:

$$(\mathbf{K} \star \mathbf{I})_{i,j} = \sum_{m,n} \mathbf{K}_{m,n} \mathbf{I}_{i+m, j+n} \quad (28)$$

CONVOLUTIONAL NEURAL NETWORKS

CROSS-CORRELATION

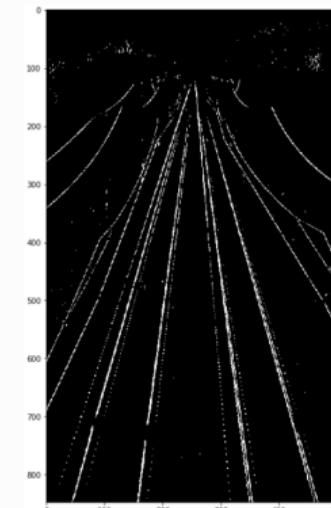
CNN actually uses cross-correlation instead of convolution:

$$(\mathbf{K} * \mathbf{I})_{i,j} = \sum_{m,n} \mathbf{K}_{m,n} \mathbf{I}_{i+m, j+n} \quad (28)$$

Cross-correlation can be used to detect edges in image. Edge is a space with big differences in pixel values.



-1	0	+1
-2	0	+2
-1	0	+1



CONVOLUTIONAL NEURAL NETWORKS

CHANNELS

Matrix \mathbf{K} is called **kernel** or **filter**.

Coloured image consists of three **channels** - RGB. The weights inside the kernel are different for each channel. Thus, we have

$$(\mathbf{K} * \mathbf{I})_{i,j} = \sum_{m,n,c} \mathbf{K}_{m,n,c} I_{i+m,j+n,c} \quad (29)$$

Also we want to use more than one kernels - one kernel can detect vertical lines, another horizontal lines and yet another can detect rounded objects. Neural net will decide on its own what shapes are relevant, but we need to provide kernels, whose weights can be trained.

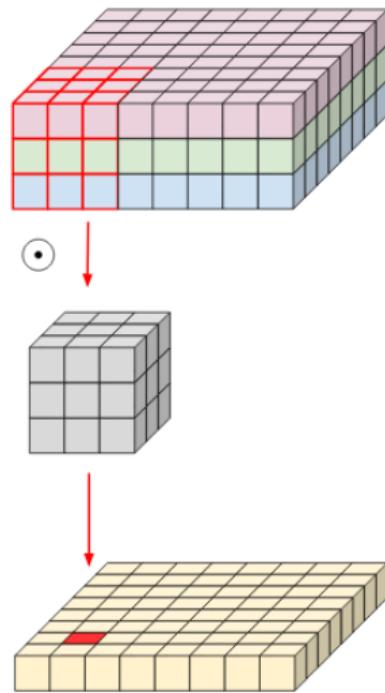
$$(\mathbf{K} * \mathbf{I})_{i,j,o} = \sum_{m,n,c} \mathbf{K}_{m,n,c,o} I_{i+m,j+n,c} \quad (30)$$

Kernel is four dimensional tensor with dimensions:

- ▶ W ... kernel width
- ▶ H ... kernel height
- ▶ C ... input channels
- ▶ F ... output channels (number of filters we want to train)

CONVOLUTIONAL NEURAL NETWORKS

PROCESSING MULTIPLE INPUT CHANNELS



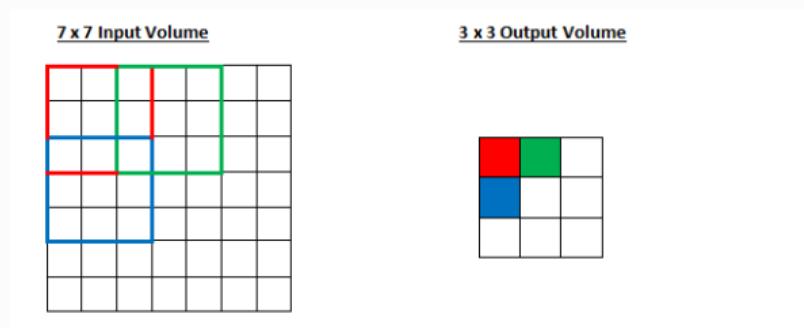
CONVOLUTIONAL NEURAL NETWORKS

STRIDE



Stride defines kernel step when computing cross-correlation.

$$(\mathbf{K} * \mathbf{I})_{i,j,o} = \sum_{m,n,c} \mathbf{K}_{m,n,c,o} \mathbf{I}_{i+S+m, j+S+n, c} \quad (31)$$



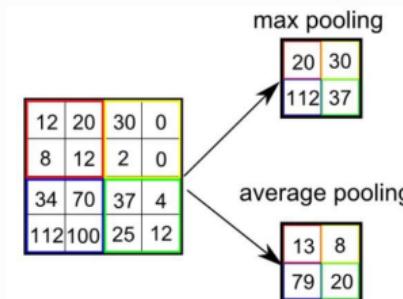
Applying stride reduces size of the output. If $S = 2$, then output is half the size of input.

CONVOLUTIONAL NEURAL NETWORKS

POOLING

Pooling is an operation that reduces width and height of the image. It is fixed operation - it does not include trainable parameters.

- ▶ Max pooling
- ▶ Average pooling



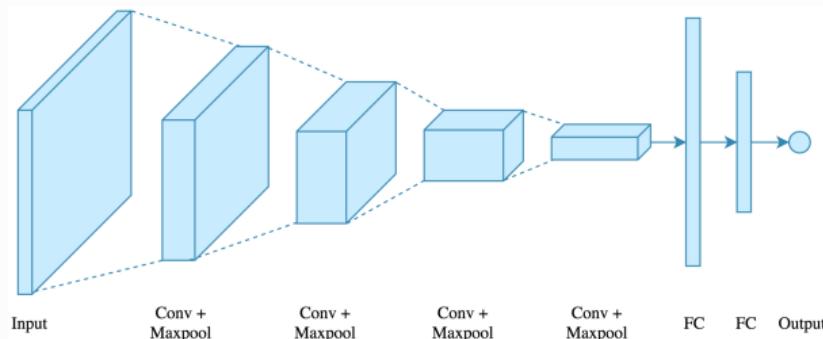
Reducing size of the image allows to detect high level patterns in the picture

CONVOLUTIONAL NEURAL NETWORKS

CNN ARCHITECTURE



- ▶ Double number of channels when performing pooling.
- ▶ Some architectures do not use fully connected layers at the end of the network.





CONVOLUTIONAL NEURAL NETWORKS

BATCH NORMALIZATION

Let us first assume multi layer perceptron neural network. As the weights of the network updates a neuron inside the layer receives different distribution of inputs. This effect is called **internal covariate shift**.

Batch normalization is trying to normalize inputs at each neuron (by normalizing outputs of previous layer). Since we are using minibatch SGD, the inputs are normalized within batches as well.

Empirical mean and variance of the batch B of size m is given by

$$\mu_B = \frac{1}{m} \sum_{i=1}^m x_i \quad (32)$$
$$\sigma_B^2 = \frac{1}{m} \sum_{i=1}^m (x_i - \mu_B)^2$$

CONVOLUTIONAL NEURAL NETWORKS

BATCH NORMALIZATION

Batch normalization layer takes the output of each neuron of previous layer $h_i^{(k)}$ (i denotes a layer and k denotes a neuron) and transforms it in two steps:

1. Normalization

$$\hat{h}_i^{(k)} = \frac{h_i^{(k)} - \mu_B^{(k)}}{\sqrt{\sigma_B^{(k)2} + \varepsilon}} \quad (33)$$

2. Shift and scale

$$x_i^k = \gamma^k \hat{h}_i^{(k)} + \beta^{(k)} \quad (34)$$

γ^k and $\beta^{(k)}$ are trainable parameters.

During inference μ and σ are calculated as a population statistics.

For convolutional networks we don't want to break its properties, mainly the shift invariance. Therefore we perform batch normalization across minibatch and spacial/temporal dimensions.

CONVOLUTIONAL NEURAL NETWORK CASE STUDIES

TARAN



ADVISORY IN DATA & ANALYTICS

CASE STUDIES

NETWORKS OVERVIEW



Motivation:

- ▶ Existing networks can help you to select proper architecture and hyperparameters

CASE STUDIES

NETWORKS OVERVIEW



Motivation:

- ▶ Existing networks can help you to select proper architecture and hyperparameters

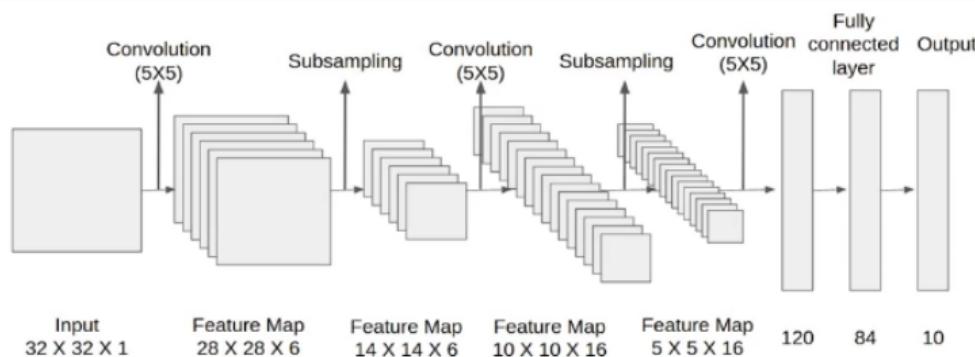
Classic networks:

- ▶ LeNet-5
- ▶ AlexNet
- ▶ VGG-16
- ▶ ResNet
- ▶ Inception

CASE STUDIES

LENET-5

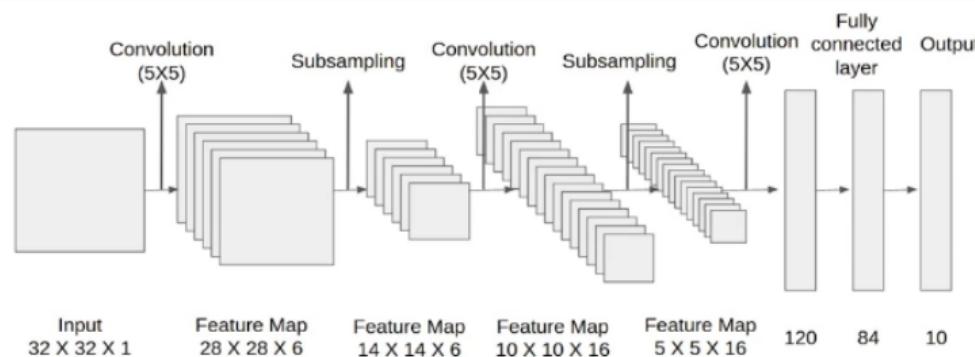
Developed in 1989 by Yann LeCun to classify hand-written digits.



CASE STUDIES

LENET-5

Developed in 1989 by Yann LeCun to classify hand-written digits.



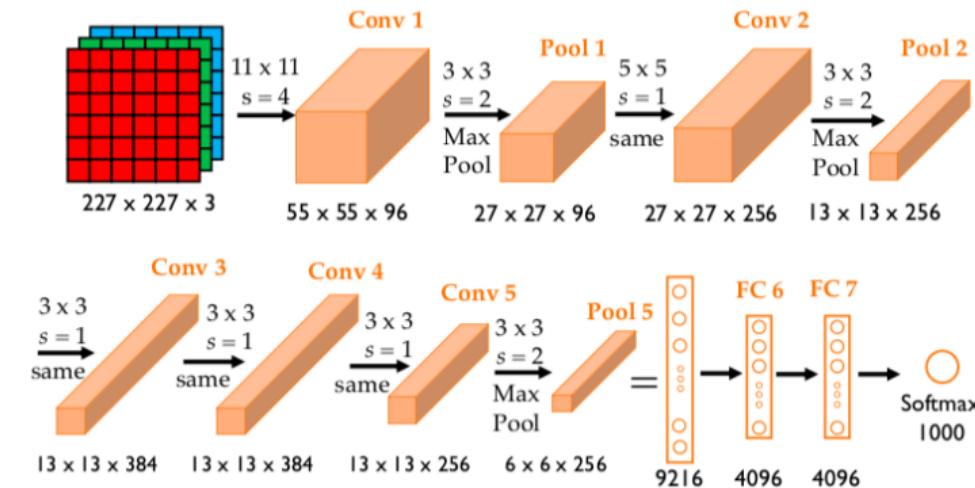
- About 60K trainable parameters

paper: Gradient-Based Learning Applied to Document Recognition (Yann LeCun, Léon Bottou, Yoshua Bengio, Patrick Haffner)

CASE STUDIES

ALEXNET

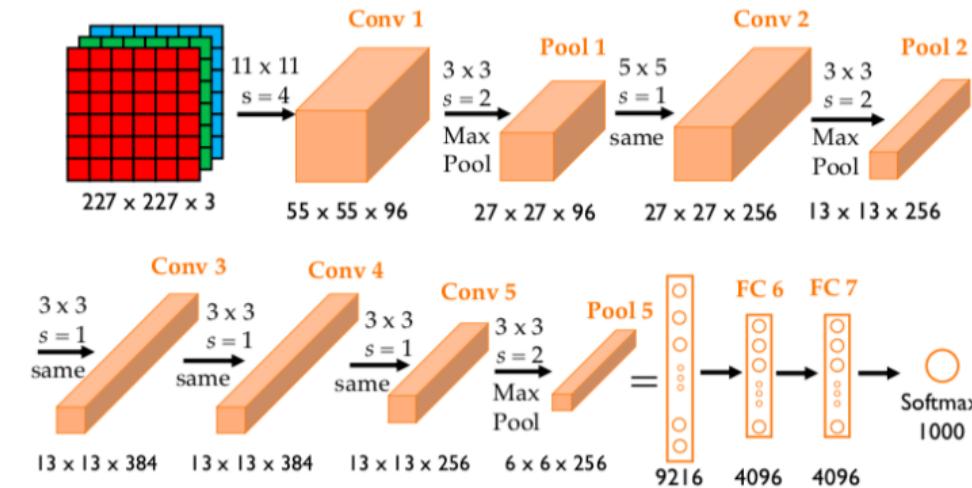
Developed in 2012 to compete in the ImageNet Large Scale Visual Recognition Challenge (image classification of 1000 categories).



CASE STUDIES

ALEXNET

Developed in 2012 to compete in the ImageNet Large Scale Visual Recognition Challenge (image classification of 1000 categories).



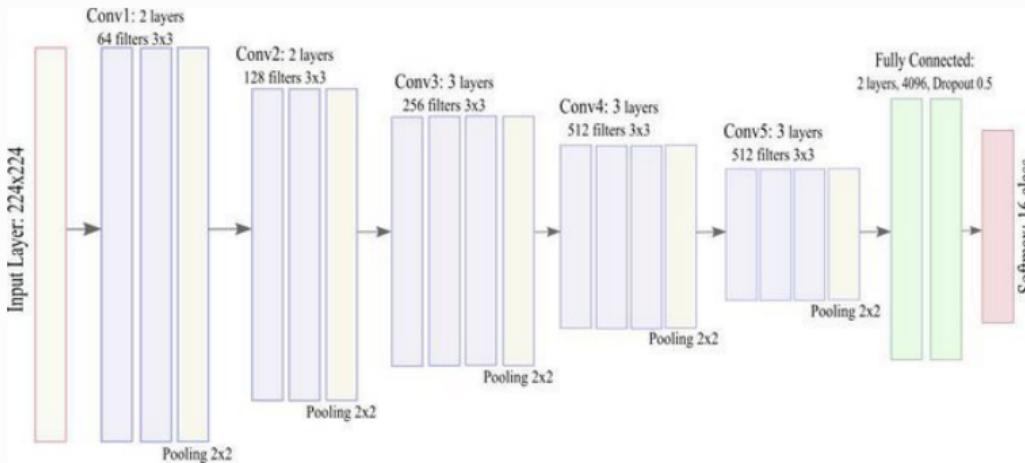
- About 60M trainable parameters

paper: ImageNet Classification with Deep Convolutional Neural Networks (Alex Krizhevsky, Ilya Sutskever, Geoffrey E. Hinton)

CASE STUDIES

VGG-16

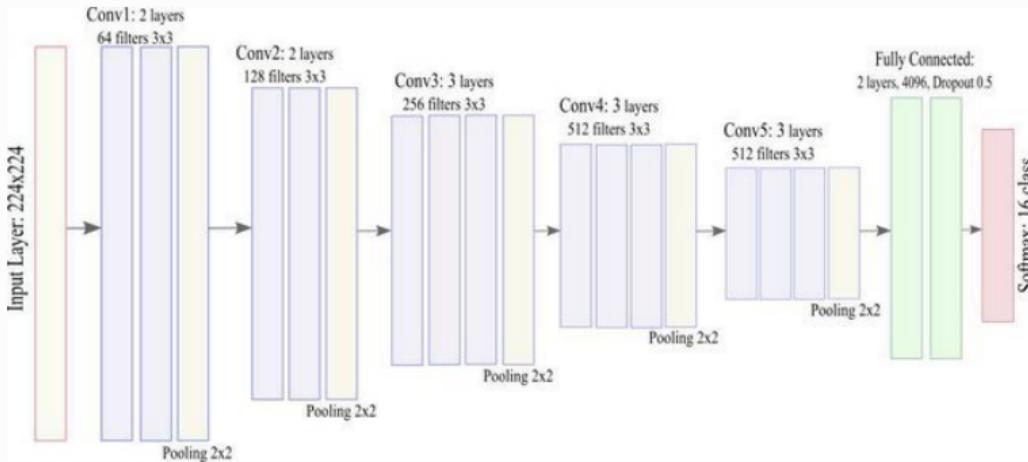
Developed in 2014.



CASE STUDIES

VGG-16

Developed in 2014.



- About 138M trainable parameters

paper: Very deep convolutional networks for large-scale image recognition (Karen Simonyan, Andrew Zisserman)

CASE STUDIES

RESNET



Residual neural networks utilizes skip connections.

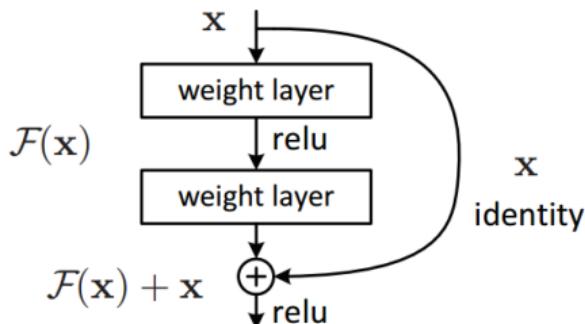
- ▶ Residual networks allows to train deep networks.
- ▶ Skip connections help with vanishing gradient problem.

CASE STUDIES

RESNET

Residual neural networks utilizes skip connections.

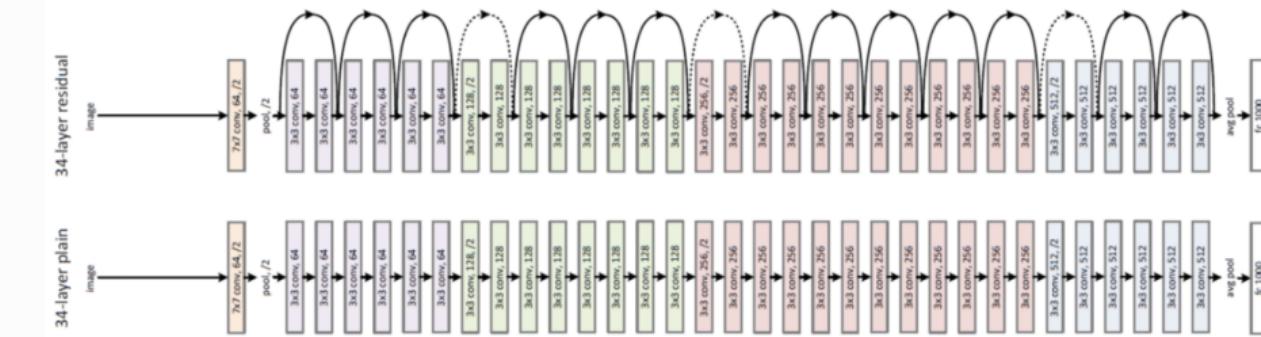
- ▶ Residual networks allows to train deep networks.
- ▶ Skip connections help with vanishing gradient problem.



paper: Deep Residual Learning for Image Recognition (Kaiming He, Xiangyu Zhang, Shaoqing Ren, Jian Sun)

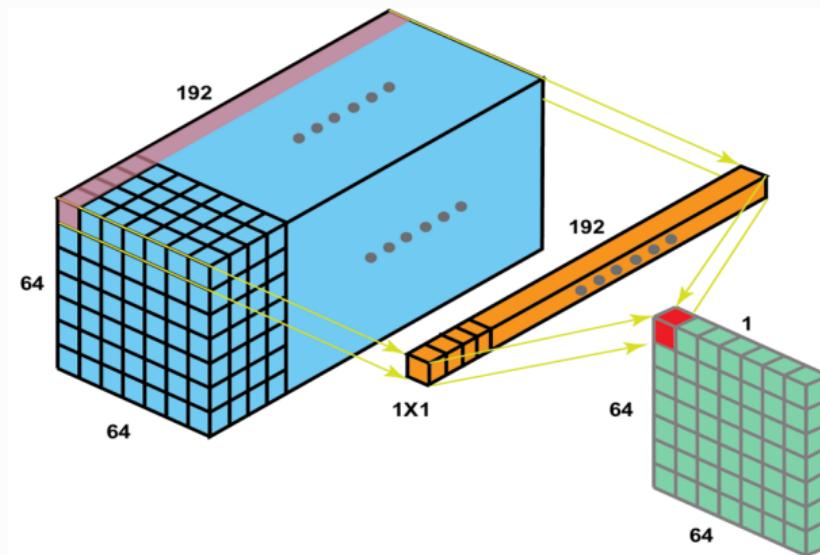
CASE STUDIES

RESIDUAL CONNECTIONS



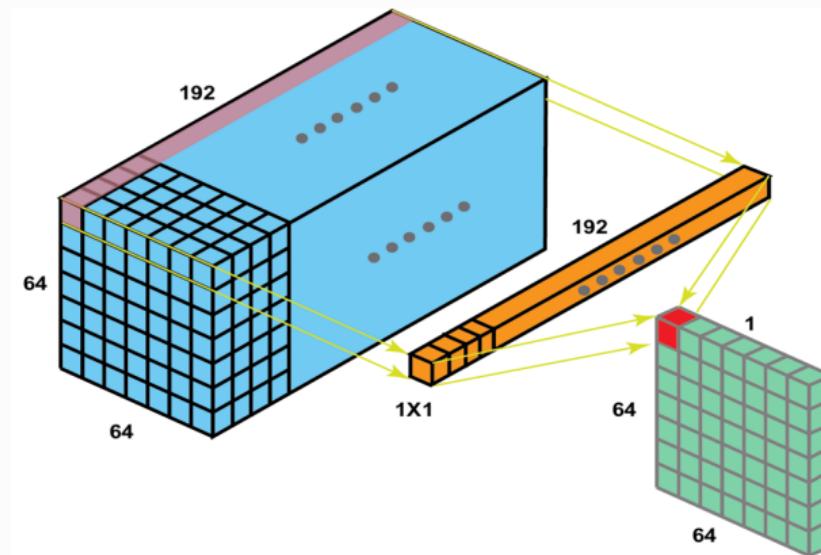
CASE STUDIES

1X1 CONVOLUTION



CASE STUDIES

1X1 CONVOLUTION



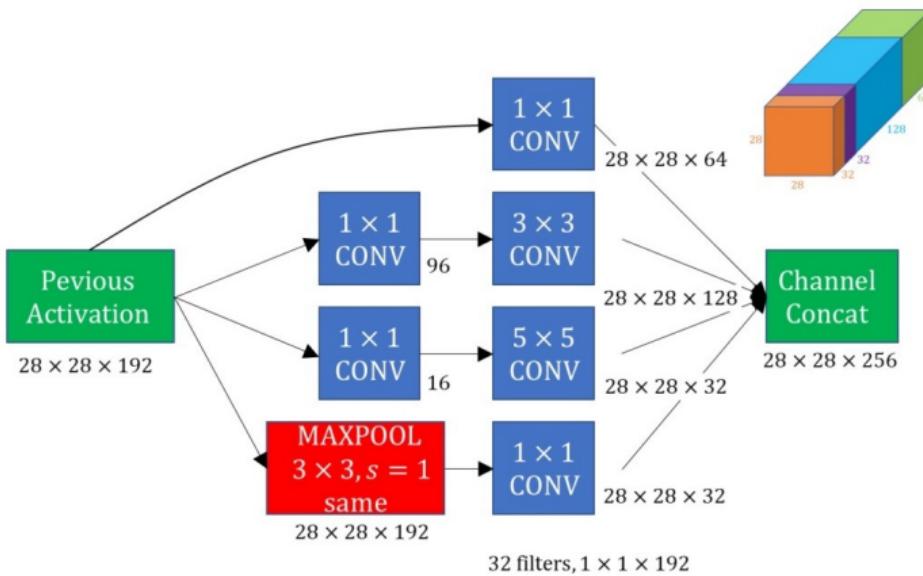
- ▶ Can be used to shrink channels
- ▶ Can add additional non-linearity (keeping channel dimensionality)
- ▶ Can reduce computational costs

CASE STUDIES

INCEPTION MODULE



An example of an Inception module

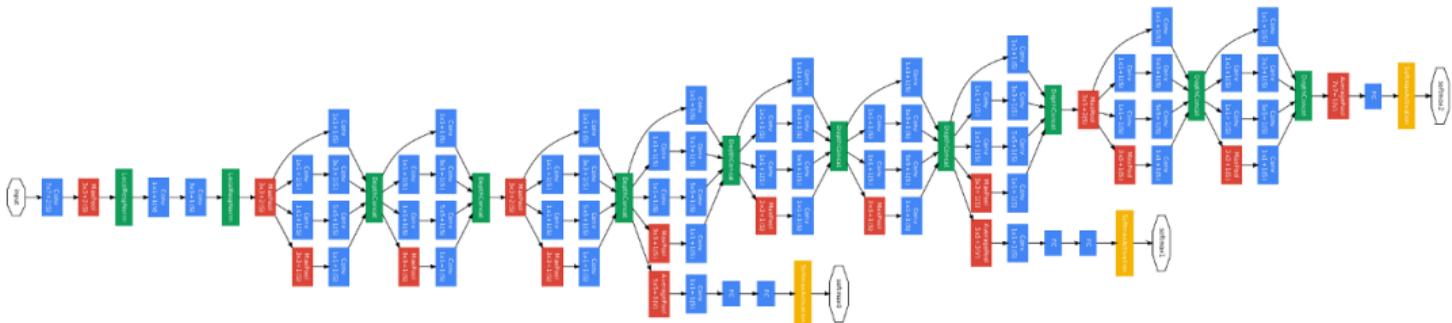


CASE STUDIES

INCEPTION NETWORK



Developed in 2014.



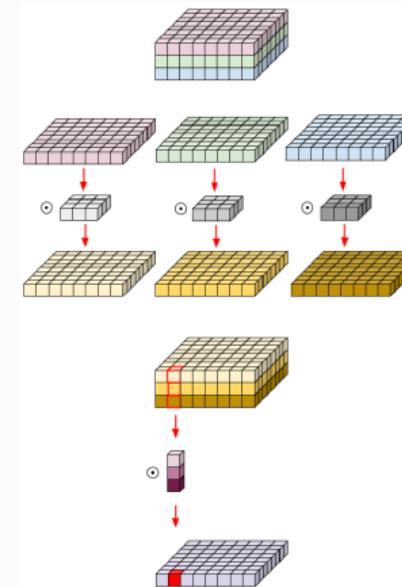
paper: Going Deeper with Convolutions (Szegedy et al.)

CASE STUDIES

DEPTH-WISE SEPARABLE CONVOLUTION

This convolution originated from the idea that depth and spatial dimension of a filter can be separated. For instance Sobel filter for edge detection can be written as vector multiplication.

$$\begin{bmatrix} 1 \\ 2 \\ 1 \end{bmatrix} * \begin{bmatrix} -1 & 0 & 1 \end{bmatrix} = \begin{bmatrix} -1 & 0 & 1 \\ -2 & 0 & 2 \\ -1 & 0 & 1 \end{bmatrix} \quad (35)$$

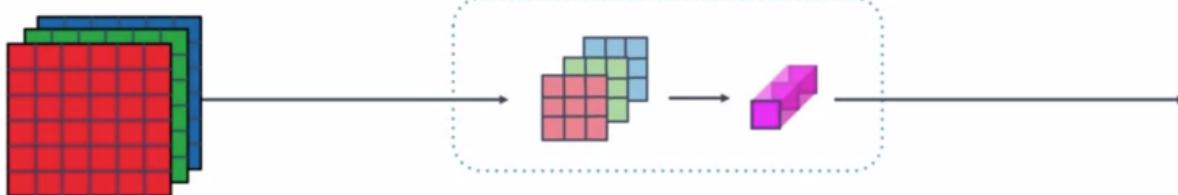


CASE STUDIES

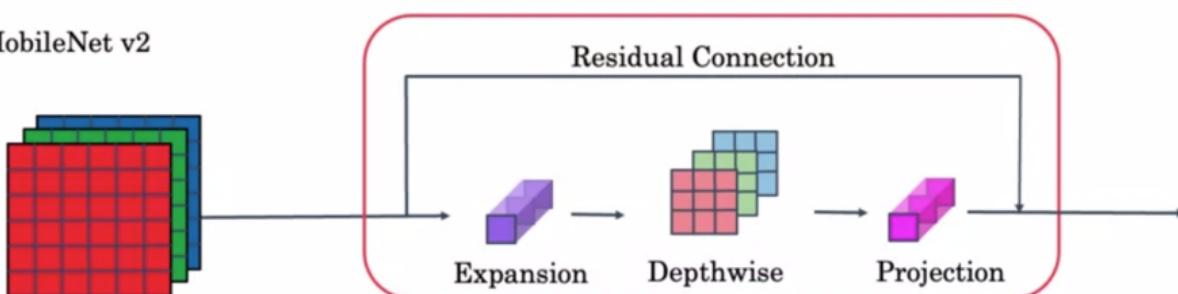
MOBILENET



MobileNet v1



MobileNet v2



CONVOLUTIONAL NEURAL NETWORKS

TRANSFER LEARNING



It might be useful to utilize network trained on different dataset.

- ▶ First layers of pre-trained network capture low level features such as horizontal lines and vertical lines.
- ▶ Pre-trained network harvested large dataset.
- ▶ Using pre-trained network can save lot of computational costs.

CONVOLUTIONAL NEURAL NETWORKS

TRANSFER LEARNING



It might be useful to utilize network trained on different dataset.

- ▶ First layers of pre-trained network capture low level features such as horizontal lines and vertical lines.
- ▶ Pre-trained network harvested large dataset.
- ▶ Using pre-trained network can save lot of computational costs.

We can:

- ▶ Retrain only classification layer (low amount of data).

CONVOLUTIONAL NEURAL NETWORKS

TRANSFER LEARNING



It might be useful to utilize network trained on different dataset.

- ▶ First layers of pre-trained network capture low level features such as horizontal lines and vertical lines.
- ▶ Pre-trained network harvested large dataset.
- ▶ Using pre-trained network can save lot of computational costs.

We can:

- ▶ Retrain only classification layer (low amount of data).
- ▶ Unfreeze last n layers of pre-trained network (larger dataset).

CONVOLUTIONAL NEURAL NETWORKS

TRANSFER LEARNING



It might be useful to utilize network trained on different dataset.

- ▶ First layers of pre-trained network capture low level features such as horizontal lines and vertical lines.
- ▶ Pre-trained network harvested large dataset.
- ▶ Using pre-trained network can save lot of computational costs.

We can:

- ▶ Retrain only classification layer (low amount of data).
- ▶ Unfreeze last n layers of pre-trained network (larger dataset).
- ▶ Unfreeze all layers and keep only architecture (very large dataset).

CLASSICAL APPROACHES TO NLP

TARAN



ADVISORY IN DATA & ANALYTICS



OVERVIEW

Terminology:

- ▶ NLP = Natural Language Processing
- ▶ Corpus - Set of texts (documents) to be used for training
 - Set of emails
 - ...
- ▶ Document - A text unit to be processed
 - Single email
 - ...

Since lot of methods are based on word frequency analysis, data preprocessing is of essence:

- ▶ Stemming/lemmatisation
- ▶ Collocation detection
- ▶ Synonyms replacement

Other techniques used in NLP:

- ▶ TF-IDF (term frequency - inverse document frequency) - for words encoding
- ▶ SVD (singular value decomposition) - for dimensionality reduction

TEXT PRE-PROCESSING

STEMMING



Stemming is a process of reducing inflected words to their word stem Example:

- ▶ fishing -> fish
- ▶ fished -> fish
- ▶ fisher -> fish

TEXT PRE-PROCESSING

STEMMING



Stemming is a process of reducing inflected words to their word stem Example:

- ▶ fishing -> fish
- ▶ fished -> fish
- ▶ fisher -> fish
- ▶ argued -> argu
- ▶ arguing -> argu



TEXT PRE-PROCESSING

STEMMING

Stemming is a process of reducing inflected words to their word stem Example:

- ▶ fishing -> fish
- ▶ fished -> fish
- ▶ fisher -> fish
- ▶ argued -> argu
- ▶ arguing -> argu

How to do stemming:

- a/ Use predefined mapping
 - ▶ simple, fast
 - ▶ new and unfamiliar words are not handled
- b/ Automated mapping
 - ▶ suffix stripping - remove "ed", "ing", "ly"
 - ▶ works poorly for exceptional relations such as "run" and "ran"

TEXT PRE-PROCESSING

LEMMASTICATION



lemma = canonical form

Lemmatisation replaces inflected words with their canonical form.

Lemmatisation uses vocabulary and morphological analysis:

- ▶ Lemmatisation uses part-of-speech (POS) tagging (recognition of nouns, verbs and other types)
- ▶ Lemmatisation uses vocabulary to determine word lemma
- ▶ Lemma for "saw" will be either "saw" or "see" - depending on the word was a noun or a verb

TEXT PRE-PROCESSING

COLLOCATIONS

Collocation refers to a group of two or more words that usually go together.

Examples:

- ▶ post office
- ▶ fast food

Collocation can be detected automatically using normalized point-wise mutual information (NPMI):

$$\begin{aligned} PMI &= \log \left[\frac{p(x,y)}{p(x)p(y)} \right] \\ NPMI &= \frac{PMI}{-\log p(x,y)} = \frac{\log[p(x)p(y)]}{\log p(x,y)} - 1 \end{aligned} \tag{36}$$

WORDS ENCODING

TF-IDF



TF-IDF = term frequency - inverse document frequency

- ▶ term frequency - frequency of a word in a document

$$TF(t, d) = \frac{\text{frequency of } t \text{ in } d}{\text{number of words in } d} \quad (37)$$

- ▶ document frequency - measures document importance in a corpus

$$DF(t) = \text{number of documents with presence of } t \quad (38)$$

- ▶ inverse document frequency - measures informativeness of t

$$IDF(t) = \log \frac{N}{DF(t) + 1} \quad (39)$$

- ▶ tf-idf

$$TF-IDF(t, d) = \frac{TF(t, d)}{IDF(t)} \quad (40)$$

TEXT PRE-PROCESSING

OTHER TIPS



- ▶ stop words - most common words with low informative value. Remove them.
 - the
 - is
 - to
 - ...
- ▶ n grams - sequence of contiguous words
- ▶ To improve performance synonyms can be replaced

DIMENSIONALITY REDUCTION

SVD



Let's assume we have TF-IDF matrix A with m rows that stand for terms and n columns that stands for documents. Rank of the matrix is $r \leq \min(m, n)$.

Singular value decomposition of A is given by:

$$A = U\Sigma V^T \tag{41}$$

where U is an $m \times r$ orthogonal matrix, Σ is $r \times r$ diagonal matrix singular values on diagonal and V is $r \times n$ orthogonal matrix.

- Singular values in Σ are greater than zero
- Largest singular value is in upper left corner of Σ

DIMENSIONALITY REDUCTION

SVD



SVD can be viewed as sum of rank one matrices:

$$A = \sum_{i=1}^r \Sigma_{ii} U_{:,i} V_i^T \quad (42)$$

where $U_{:,i}$ is i-th column of U and V_i^T is i-th row of V .

We can approximate A using only $k < r$ strongest singular values:

$$A \approx A_k = \sum_{i=1}^k \Sigma_{ii} U_{:,i} V_i^T = U_k \Sigma_k V_k^T \quad (43)$$

To project a document d represented by m-dimensional vector to k-dimensional space, we can use:

$$\hat{d} = U_k^T d \quad (44)$$



SEQUENCE MODELS

TARAN



ADVISORY IN DATA & ANALYTICS

SEQUENCE MODELS

EXAMPLES



Examples of sequence data:

- ▶ Speech recognition

- ▶ Music generation

- ▶ Sentiment classification

The food was awful. \mapsto ●○○○○

- ▶ Machine translation

La nourriture était horrible. \mapsto *The food was awful.*

- ▶ Video recognition

- ▶ Name entity recognition

Yesterday, Harry met Hermione. \mapsto *Yesterday, Harry met Hermione.*

SEQUENCE MODELS

NOTATION



We will use following notation:

- x ... input sequence
- $x^{<i>}$... i-th element of sequence x
- T_x ... total number of elements in x
- y ... output sequence
- $y^{<i>}$... i-th element of sequence y
- T_y ... total number of elements in y
- $X^{(i)<t>}$... t-th element of i-th training sequence
- $Y^{(i)<t>}$... t-th element of output sequence for i-th training sequence
- $T_x^{(i)}$... total number of elements in i-th training sequence
- $T_y^{(i)}$... total number of elements in output sequence for i-th training sequence

SEQUENCE MODELS

VOCABULARY



How do we represent words in a sentence?

Vocabulary is enumerated list of known words.

- ▶ It can include punctuation.
- ▶ *< UNKNOWN >* can be added to represent word not included in vocabulary.

Each word can be transformed using one-hot encoding.

SEQUENCE MODELS

FORWARD PROPAGATION

Usually:

$$a^{<0>} = \mathbf{0} \quad (45)$$

After first layer:

$$\begin{aligned} a^{<1>} &= g_1 (W_{aa}a^{<0>} + W_{ax}x^{<1>} + b_a) \\ \hat{y}^{<1>} &= g_2 (W_{ya}a^{<1>} + b_y) \end{aligned} \quad (46)$$

Note: In W_{ax} , first subscript (a) means that W_{at} is used to calculate a -like quantity, while second subscript (x) means that W_{ax} will be multiplied by x

Simplified notation:

$$\begin{aligned} W_a &= [W_{aa} \mid W_{ax}] \\ W_y &= W_{ya} \end{aligned} \quad (47)$$

Thus:

$$\begin{aligned} a^{<t>} &= g_1 (W_a[a^{<t-1>} \mid x^{<t>}] + b_a) \\ \hat{y}^{<t>} &= g_2 (W_ya^{<t>} + b_y) \end{aligned} \quad (48)$$

SEQUENCE MODELS

RNN TYPES



Different types of RNN:

(a) **one-to-one**

This is standard neural network

(b) **many-to-many ($T_x = T_y$)**

example: name entity recognition

(c) **many-to-many ($T_x \neq T_y$)**

example: machine translation

(d) **many-to-one**

example: semantic classification

SEQUENCE MODELS

LANGUAGE MODEL



Let's assume speech recognition system:

Most crimes are committed at **night**.

Most crimes are committed at **knight**.

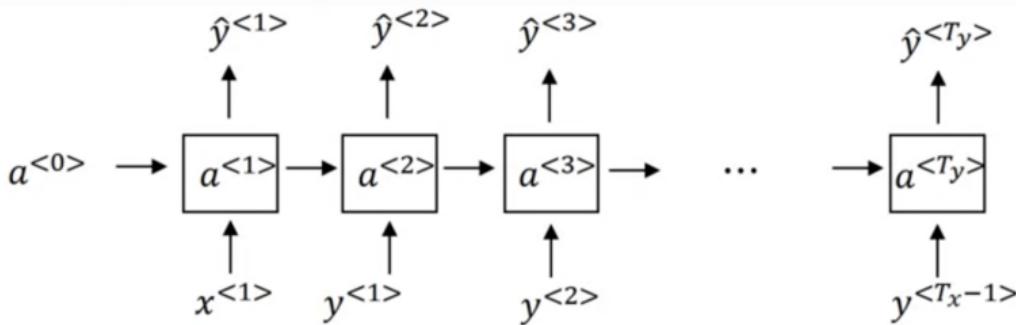
Both sentences have the same pronunciation.

Language model assigns probability to each word given preceding words in a sentence.

$$\mathbb{P} [\textit{night} \mid \textit{Most, crimes, are, committed, at}] > \mathbb{P} [\textit{knight} \mid \textit{Most, crimes, are, committed, at}] \quad (49)$$

SEQUENCE MODELS

LANGUAGE MODEL

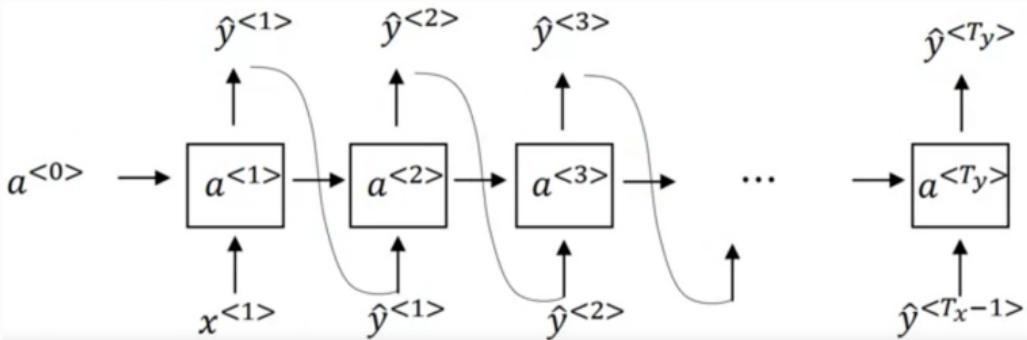


- $\hat{y}^{<t>}$ predicts probability of each word in vocabulary given previous words.
- $x^{<1>} = a^{<0>} = \mathbf{0}$
- loss function:

$$\begin{aligned} l^{<t>}(y_i^{<t>}, \hat{y}_i^{<t>}) &= - \sum_i y_i^{<t>} \log \hat{y}_i^{<t>} \\ L &= \sum_t l^{<t>}(y_i^{<t>}, \hat{y}_i^{<t>}) \end{aligned} \tag{50}$$

SEQUENCE MODELS

SAMPLING NOVEL SEQUENCE



- Stop generating sequence if $<EOS>$ or after fixed number of iteration
- $x^{<1>} = a^{<0>} = \mathbf{0}$
- $\hat{y}^{<t>}$ is drawn from distribution given by words probabilities
- Character level language model can be trained

SEQUENCE MODELS

VANISHING GRADIENT PROBLEM



The **cat**, which ate ..., **was** full.

The **cats**, which ate ..., **were** full.

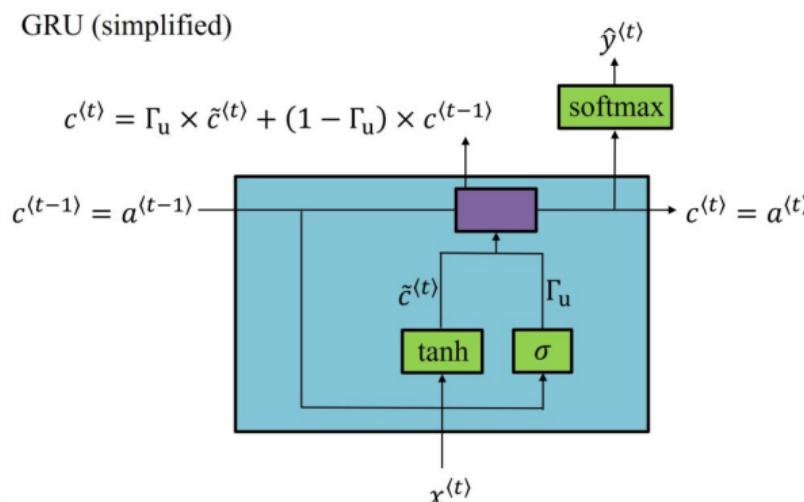
Sometimes probability of a word depends on very distant positions.

- ▶ Output $\hat{y}^{<t>}$ is strongly influenced by close inputs ($x^{<t>} , x^{<t-1>}$).
- ▶ It is not easy to propagate gradient from the end of sequence to the beginning.
- ▶ Exploding gradient can be dealt with by gradient clipping.
- ▶ What about vanishing gradient?

SEQUENCE MODELS

GATED RECURRENT UNIT

$$\begin{aligned}
 c^{<t-1>} &= a^{<t-1>} \\
 \tilde{c}^{<t>} &= \tanh(W_c[c^{<t-1>}, x^{<t>}] + b_c) \\
 \Gamma_u &= \sigma(W_u[c^{<t-1>}, x^{<t>}] + b_u) \\
 c^{<t>} &= \Gamma_u \times \tilde{c}^{<t>} + (1 - \Gamma_u) \times c^{<t-1>}
 \end{aligned} \tag{51}$$



SEQUENCE MODELS

GATED RECURRENT UNIT



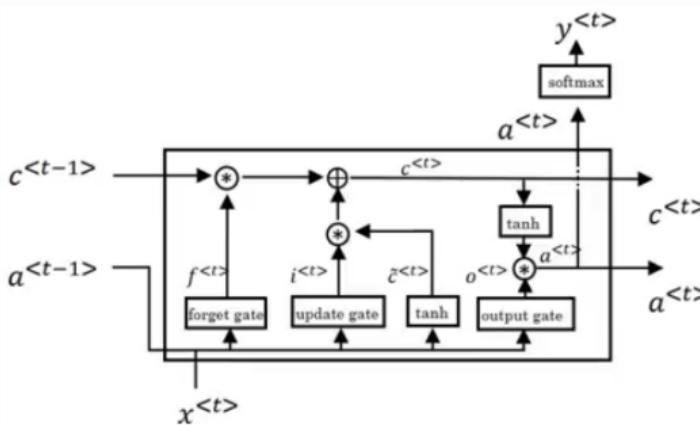
Full GRU:

$$\begin{aligned} c^{<t-1>} &= q^{<t-1>} \\ \tilde{c}^{<t>} &= \tanh(W_c[\Gamma_r \times c^{<t-1>}, x^{<t>}] + b_c) \\ \Gamma_u &= \sigma(W_u[c^{<t-1>}, x^{<t>}] + b_u) \\ \Gamma_r &= \sigma(W_r[c^{<t-1>}, x^{<t>}] + b_r) \\ c^{<t>} &= \Gamma_u \times \tilde{c}^{<t>} + (1 - \Gamma_u) \times c^{<t-1>} \end{aligned} \tag{52}$$

SEQUENCE MODELS

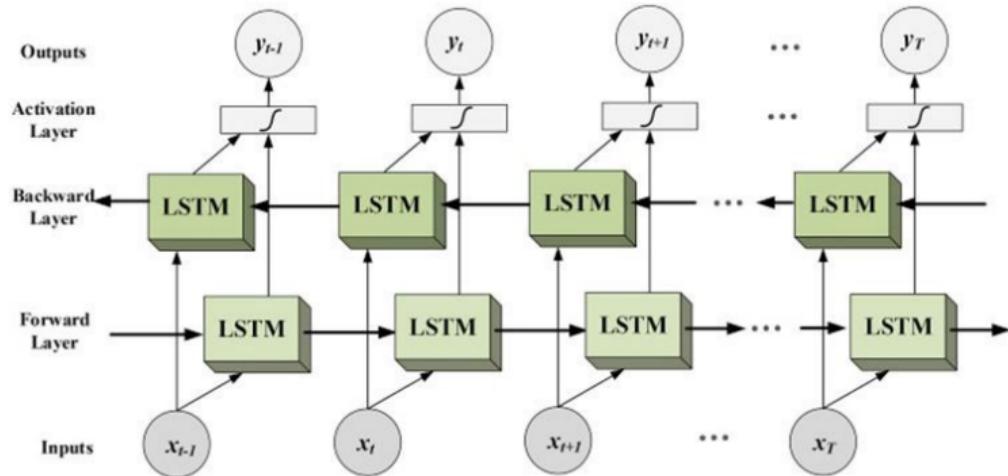
LONG SHORT TERM MEMORY

$$\begin{aligned}
 \tilde{c}^{<t>} &= \tanh(W_c[a^{<t-1>}, x^{<t>}] + b_c) \\
 \Gamma_u &= \sigma(W_u[a^{<t-1>}, x^{<t>}] + b_u) \\
 \Gamma_f &= \sigma(W_f[a^{<t-1>}, x^{<t>}] + b_f) \\
 \Gamma_o &= \sigma(W_o[a^{<t-1>}, x^{<t>}] + b_o) \\
 c^{<t>} &= \Gamma_u \times \tilde{c}^{<t>} + \Gamma_f \times c^{<t-1>} \\
 a^{<t>} &= \Gamma_o \times \tanh(c^{<t>})
 \end{aligned} \tag{53}$$



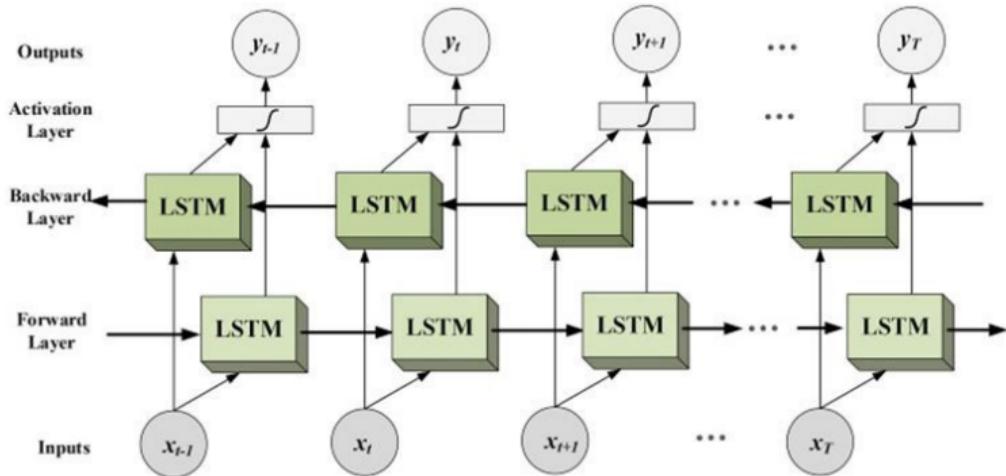
SEQUENCE MODELS

BIDIRECTIONAL RNN



SEQUENCE MODELS

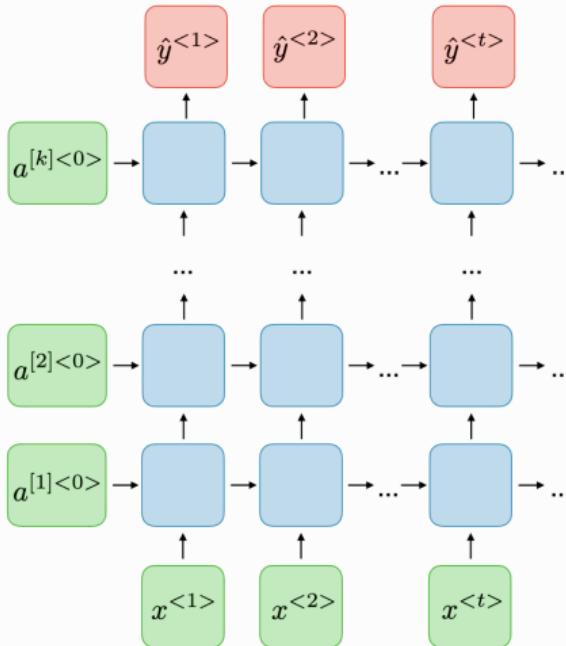
BIDIRECTIONAL RNN



- + To predict $\hat{y}^{}$ the whole sequence is used, not just sequence up to $x^{}$
- To calculate prediction, you need to wait until the whole sequence is finished

SEQUENCE MODELS

DEEP RNN



- ▶ Each layer has one weight matrix assigned
- ▶ Because of time dimension, three layer network might be already quite large
- ▶ Deep RNN can be build using RNN, GRU, LSTM or BRNN layers

WORD EMBEDDINGS

TARAN



ADVISORY IN DATA & ANALYTICS

WORD EMBEDDINGS

WORD REPRESENTATION



So far we were using one-hot encoding for words representation.

Man (5391)	Woman (9853)	King (4914)	Queen (7157)	Water (9532)	Beer (982)
0	0				
⋮					
1	⋮				
		⋮	⋮	⋮	⋮
⋮	1				
		⋮			
0	0				

WORD EMBEDDINGS

WORD REPRESENTATION



Instead one-hot, we might want to use feature representation.

	Man (5391)	Woman (9853)	King (4914)	Queen (7157)	Water (9532)	Beer (982)
Gender	-0.97	0.98	-0.95	0.93	0.02	0.01
Royal	0.04	-0.02	0.96	0.95	0.01	-0.04
Beverage	-0.03	0.01	0.02	-0.04	0.97	0.98
:						

- ▶ Features in embeddings don't have nice meaning like in the table above
- ▶ Vectors for similar things (water and beer) are similar in feature space

WORD EMBEDDINGS

t-SNE



t-SNE = t-distributed stochastic neighbour embedding

Let us have a set of N high-dimensional objects $\mathbf{x}_1, \dots, \mathbf{x}_N$.

$$p_{j|i} = \frac{\frac{\exp(-\|\mathbf{x}_i - \mathbf{x}_j\|^2)}{2\sigma_i^2}}{\sum_{k \neq i} \frac{\exp(-\|\mathbf{x}_i - \mathbf{x}_k\|^2)}{2\sigma_i^2}} \quad \forall i \neq j \quad (54)$$

$$p_{i|i} = 0$$

"The similarity of datapoint x_j to datapoint x_i is the conditional probability $p_{j|i}$ that x_i would pick x_j as its neighbour if neighbours were picked in proportion to their probability density under a Gaussian centred at x_i ."

$$p_{ij} = \frac{p_{j|i} + p_{i|j}}{2N} \quad (55)$$

$$p_{ii} = 0$$

WORD EMBEDDINGS

t-SNE



We want to learn d -dimensional representation $\mathbf{y}_1, \dots, \mathbf{y}_N$ of $\mathbf{x}_1, \dots, \mathbf{x}_N$.

$$q_{ij} = \frac{(1 + \|\mathbf{y}_i - \mathbf{y}_j\|^2)^{-1}}{\sum_k \sum_{l \neq k} (1 + \|\mathbf{y}_k - \mathbf{y}_l\|^2)^{-1}} \quad (56)$$

$$q_{ii} = 0$$

Locations of $\mathbf{y}_1, \dots, \mathbf{y}_N$ are given by minimization of distance between distributions P and Q .

$$KL(P||Q) = \sum_{i \neq j} p_{ij} \log \frac{p_{ij}}{q_{ij}} \quad (57)$$

WORD EMBEDDINGS

SIMILARITY FUNCTIONS



What if, for a given word w , we want to find a most similar word in embedding space:

$$\operatorname{argmax}_{w_i} \text{sim}(e_w, e_{w_i}) \quad (58)$$

Similarity functions can be:

- Cosine similarity

$$\text{sim}(u, v) = \frac{u^T v}{\|u\|_2 \|v\|_2} \quad (59)$$

- Negative Euclidean distance

$$\text{sim}(u, v) = -\|u - v\|^2 \quad (60)$$



WORD EMBEDDINGS

LEARNING EMBEDDINGS

Let's o_i be the one-hot encoding of i -th word in vocabulary.

$$o_i = \begin{bmatrix} 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \end{bmatrix} \leftarrow i\text{-th position} \quad (61)$$

Embedded vector for i -th word is given by embedding matrix E .

$$e_i = E \cdot o_i = \begin{bmatrix} e_{11} & \dots & e_{1m} \\ \vdots & \ddots & \vdots \\ e_{n1} & \dots & e_{nm} \end{bmatrix} \cdot \begin{bmatrix} 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \end{bmatrix} \quad (62)$$

For vocabulary size 10000 and 300-dimensional embedding space, $E \in \mathbb{R}^{300 \times 10000}$

WORD EMBEDDINGS

LEARNING EMBEDDINGS



Let's have a language model.

I want a glass of orange juice.

Based on "I want a glass of orange", we want to predict "juice".

WORD EMBEDDINGS

LEARNING EMBEDDINGS



Let's have a language model.

I want a glass of orange juice.

Based on "I want a glass of orange", we want to predict "juice".

Model:

- ▶ Start with on-hot encoding
- ▶ Apply embedding
- ▶ Add output layer with softmax activation (output layer size is given by vocabulary size)

Embedding matrix is learned through model training (E consists of trainable parameters).



WORD EMBEDDINGS

LEARNING EMBEDDINGS

Let's have a language model.

I want a glass of orange juice.

Based on "I want a glass of orange", we want to predict "juice".

Model:

- ▶ Start with on-hot encoding
- ▶ Apply embedding
- ▶ Add output layer with softmax activation (output layer size is given by vocabulary size)

Embedding matrix is learned through model training (E consists of trainable parameters).

Words used for predicting target word are called context.

- ▶ n words left from target word
- ▶ n words left + n words right from target word
- ▶ 1 word left
- ▶ nearby one word (skip-gram)

WORD EMBEDDINGS

LEARNING EMBEDDINGS

I want a glass of orange juice.

context	target
glass	orange
glass	of
glass	want
:	:

$$\mathbb{P}[t|c] = \frac{e^{\theta_t^T e_c}}{\sum_{j=1}^{10000} e^{\theta_j^T e_c}} \quad (63)$$

Loss function:

$$L(\hat{y}, y) = - \sum_{i=1}^{10000} y_i \log \hat{y}_i \quad (64)$$

WORD EMBEDDINGS

WORD2VEC



Problem with previous approach is the computational cost.

We will use negative sampling to create a training sample:

context	target	y
orange	juice	1
orange	football	0
orange	paper	0
:	:	0
orange	of	0

For one positive context-target pair we added K negative pairs.

The output layer will still be of size of vocabulary size, but each neuron will calculate sigmoid instead of softmax.

$$\mathbb{P}[y = 1 | c, t] = \sigma(\theta_t^T e_c) \quad (65)$$

WORD EMBEDDINGS

WORD2VEC



How to choose negative samples:

- ✗ We can draw negative samples based on distribution over word frequencies
Leads to high presence of words such as "the", "of", "a"
- ✗ We can sample words uniformly
This is very non-representative
- ✓ Use following distribution (discovered empirically)

$$\mathbb{P}[w_i] = \frac{f(w_i)^{3/4}}{\sum_j f(w_j)^{3/4}} \quad (66)$$

WORD EMBEDDINGS

GloVe

GloVe = global vector for word representation

GloVe quantifies how often pair of words can be found close together.

$$X_{ij} = \text{number of times } j \text{ appears in context of } i \quad (67)$$

The model uses following loss function:

$$L = \sum_i \sum_j f(X_{ij})(\theta_i^T e_j + b_i + b_j - \log X_{ij})^2 \quad (68)$$

where

- ▶ $f(X_{ij}) = 0$ if $X_{ij} = 0$ and " $0 \cdot \log 0 = 0$ "
- ▶ $f(X_{ij})$ stands for weights that reflects words frequency

Final embeddings are calculated as:

$$e_i^{final} = \frac{e_i + \theta_i}{2} \quad (69)$$