

Umělá inteligence a strojové učení

Michal Hradiš: Neural networks

Where to get more information

Ian Goodfellow and Yoshua Bengio and Aaron Courville: **Deep Learning**, MIT Press, 2016. <http://www.deeplearningbook.org/>

Andrew Ng and Kian Katanforoosh: **CS230 Deep Learning**, Stanford, <https://cs230.stanford.edu/>

CS231n: **Convolutional Neural Networks for Visual Recognition**, Stanford, <http://cs231n.stanford.edu/index.html>

Andrew Ng, [CS229 Machine Learning](https://www.youtube.com/watch?v=UzxYIbK2c7E&list=PLEBC422EC5973B4D8), Stanford, <https://www.youtube.com/watch?v=UzxYIbK2c7E&list=PLEBC422EC5973B4D8>

Python Notebooks

Basic Pytorch -

<https://colab.research.google.com/drive/1ytNsahwdl4FJDkXmaw7CCQQJqQDiAxJt?usp=sharing>

Image classification -

<https://colab.research.google.com/drive/1oh-iicfuclPJRudYZVIUWRRLtJiMdM4k?usp=sharing#scrollTo=5PRaK3DSPzQw>

<https://colab.research.google.com/drive/1a8iZeo-pmrRDAH8ck0TrEmaulMnaKNC#scrollTo=7k37fwroNnN>

Pytorch tutorials - <https://pytorch.org/tutorials/>



AI PLAYS BASKETBALL

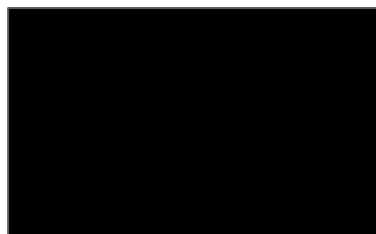
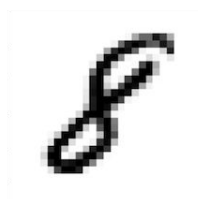
HUMAN DIGITIZATION





What are neural networks?

- Complex and flexible functions (and computationally expensive)
- Can learn to approximate a desired function from training data
- Can accept large range of input types and produce large range of output types
- Can be “customized” to solve very different tasks.

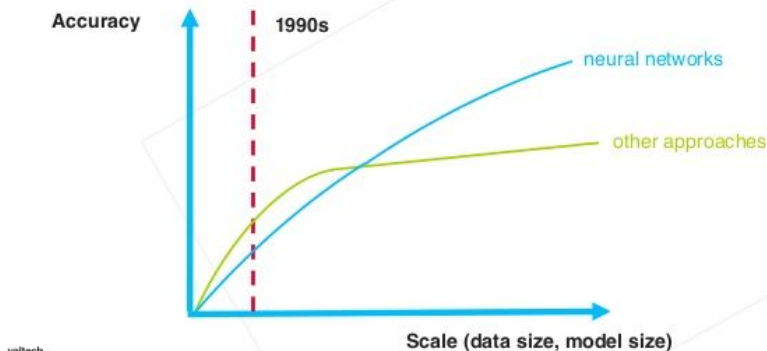


8

Why?

- Scales with computation
- Scales with data
- General framework
- State of the art
- Massive development
- Mature tools

More Data + Bigger Models



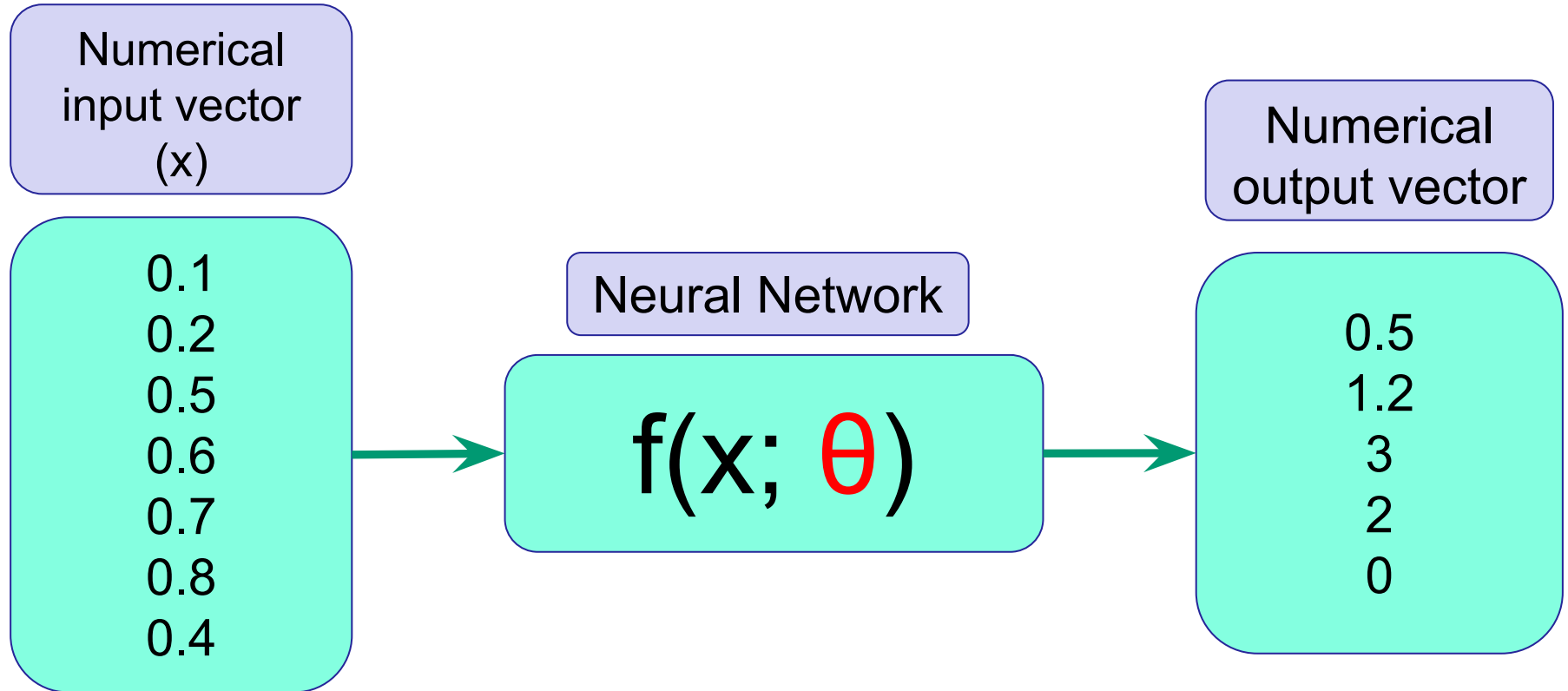
valtech.

<https://www.scribd.com/document/355752799/Jeff-Dean-s-Lecture-for-YC-AI>

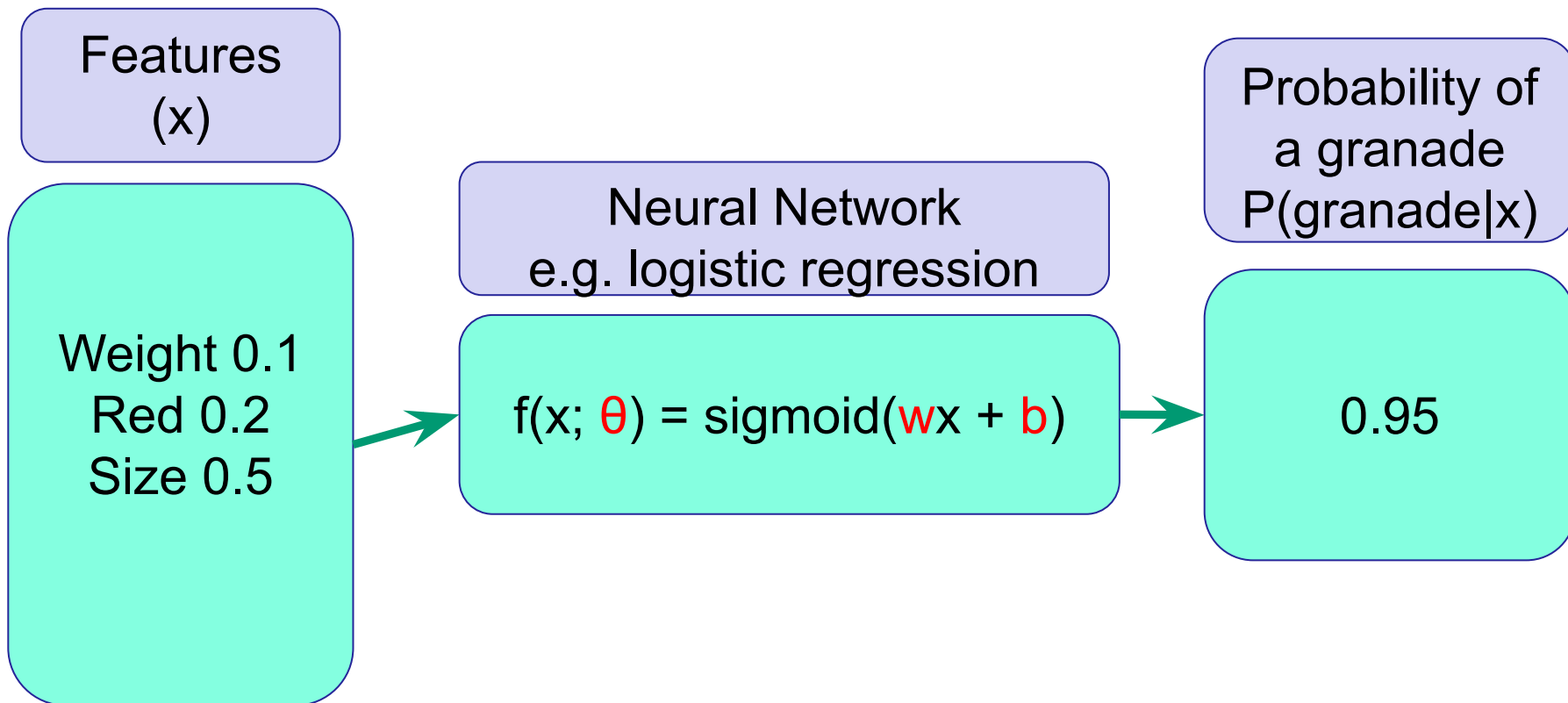
Goals

- Understand basics of machine learning with neural networks
- Understand the flexibility of neural networks
- Have an idea what tools are available
- Complex models and their applications
- Basic intuitions when solving practical tasks
 - Data
 - Generalization
 - Regularization
 - What can go wrong
 - Transfer learning

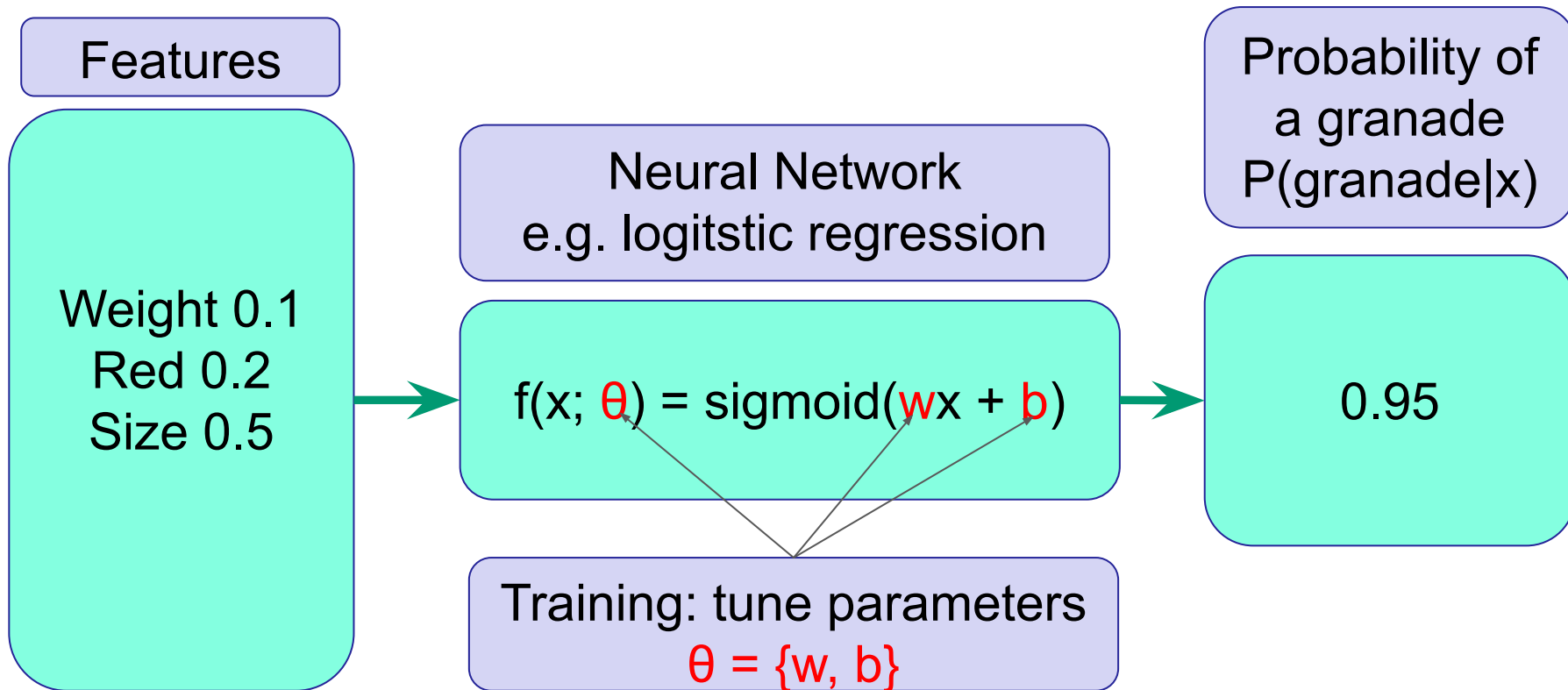
What is a neural network? (in this lecture)



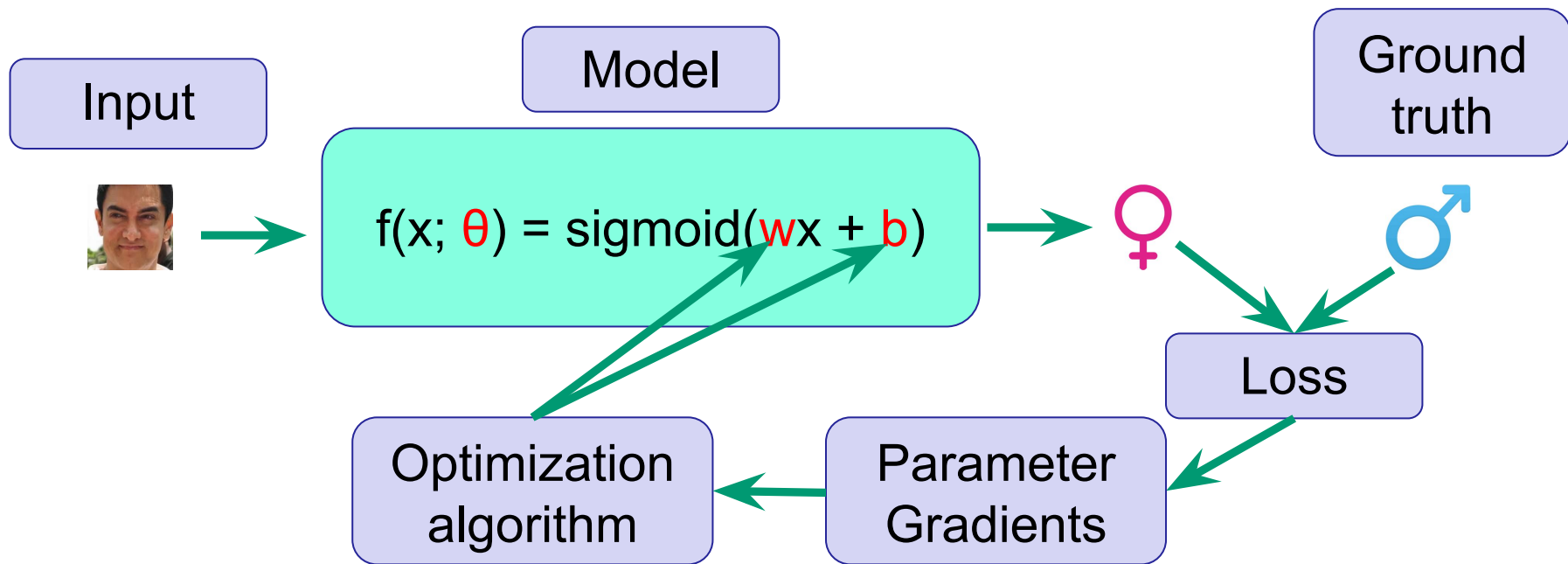
Binary classification network



Binary classification network



Training



Supervised learning components

Training data - inputs and desired network outputs

$$D = \left\{ \left(\text{img}_1, \text{♂} \right), \left(\text{img}_2, \text{♀} \right), \left(\text{img}_3, \text{♂} \right), \left(\text{img}_4, \text{♀} \right) \right\}$$

Model with parameters (network)

$$f(x; \theta) = \text{sigmoid}(wx + b)$$

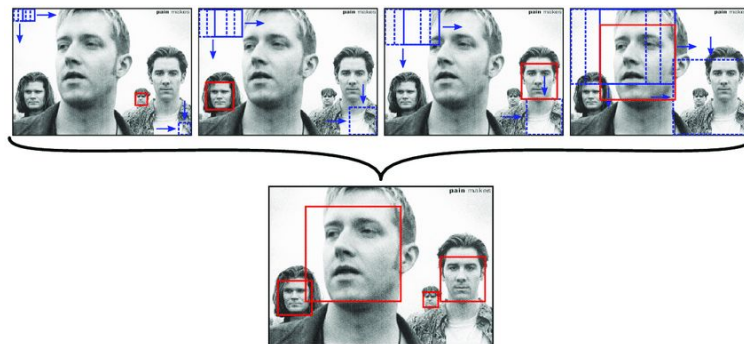
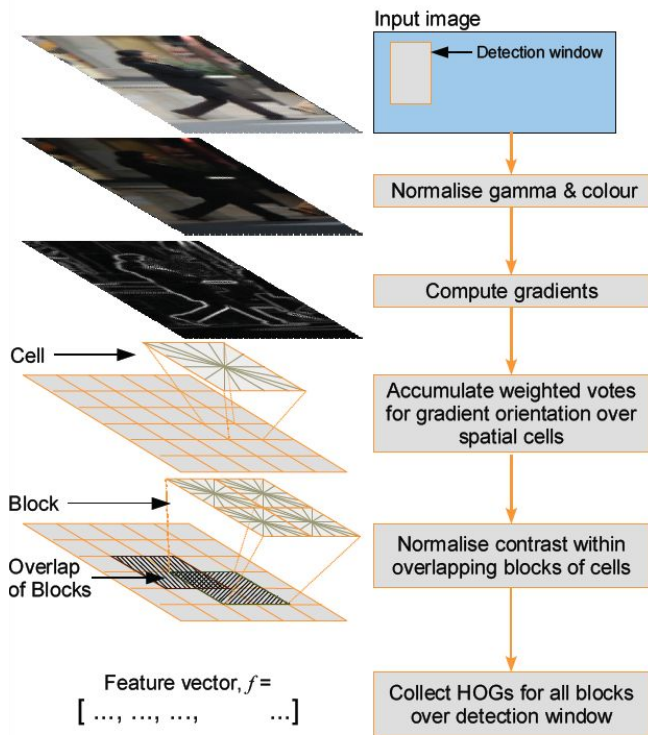
Loss function (e.g. cross entropy)

$$\text{loss}(x, y) = \begin{cases} \log(f(x; \theta)) & \text{if } y = \text{male} \\ \log(1 - f(x; \theta)) & \text{if } y = \text{female} \end{cases}$$

Objective function

$$J(D, \theta) = \sum_{D=\{(x_i, y_i), \dots\}} \text{loss}(f(x_i, \theta), y_i)$$

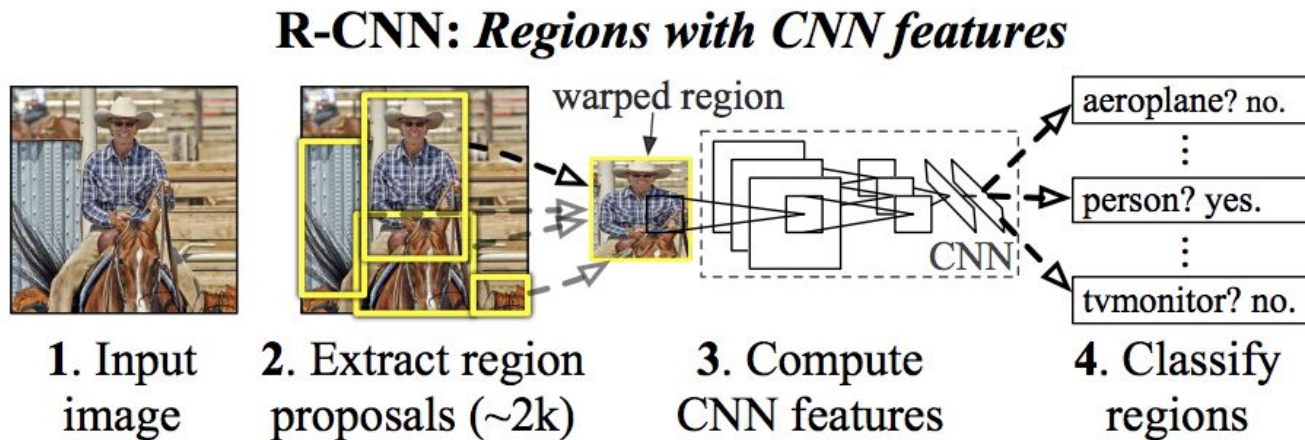
“Shallow” learning



$$f([\text{features}]; \theta) = \text{sigmoid}(w[\text{features}] + b)$$

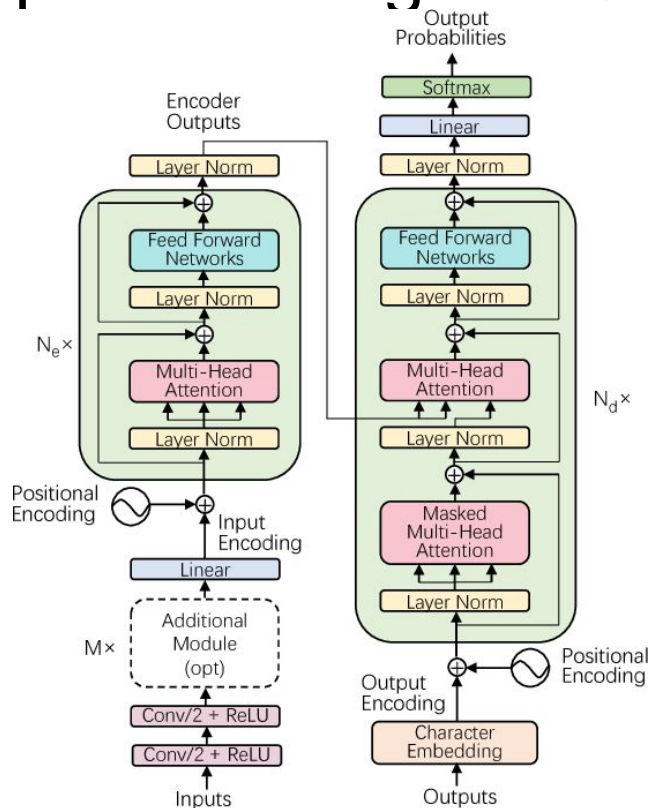
Deep learning

- No hand coded feature extraction
- Our domain knowledge goes into network architecture, regularization, loss, data augmentation, ...

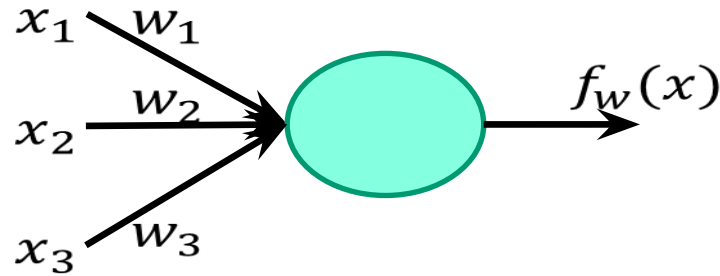


End to end speech recognition/translation

Honza šel na jahody



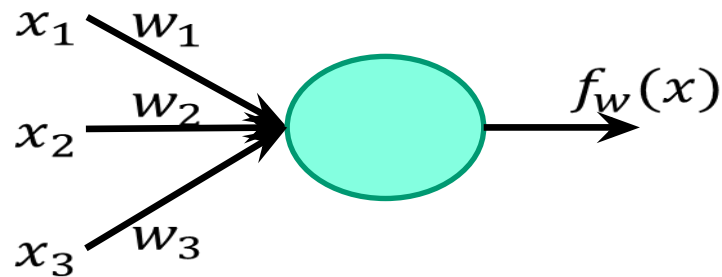
Neural network components (neurons)



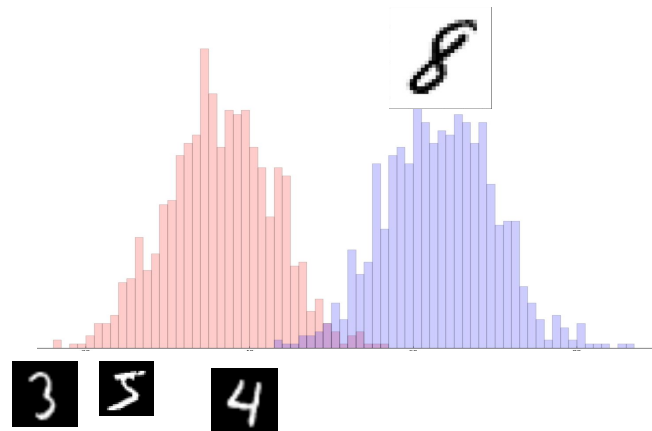
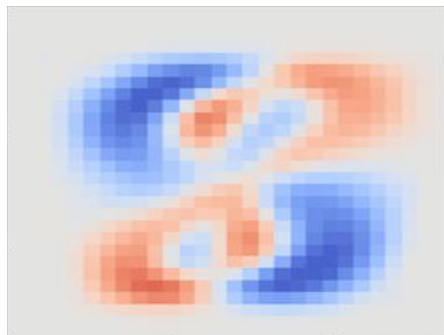
$$\begin{aligned} f(x; w) &= w_1x_1 + w_2x_2 + w_3x_3 \\ &= \sum_i w_i x_i \end{aligned}$$

Neurons as templates

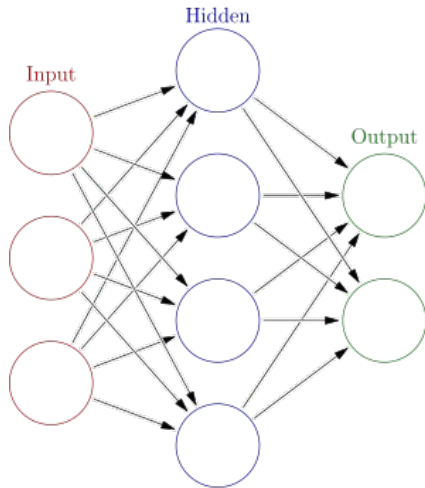
Neurons have higher activations if the input “matches” their weights.



*

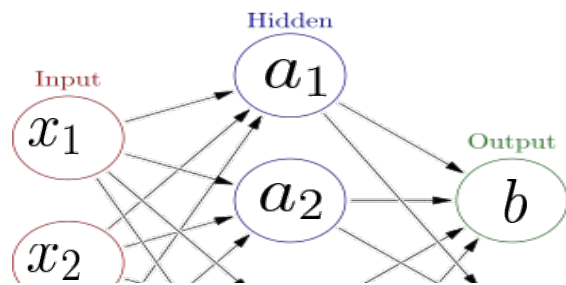


Composition of linear functions?



$$f(\mathbf{x}; \mathbf{w}) = \sum_i w_i x_i$$

Composition of linear functions?



$$f(\mathbf{x}; \mathbf{w}) = \sum_i w_i x_i$$

$$a_1 = 1x_1 + 2x_2$$

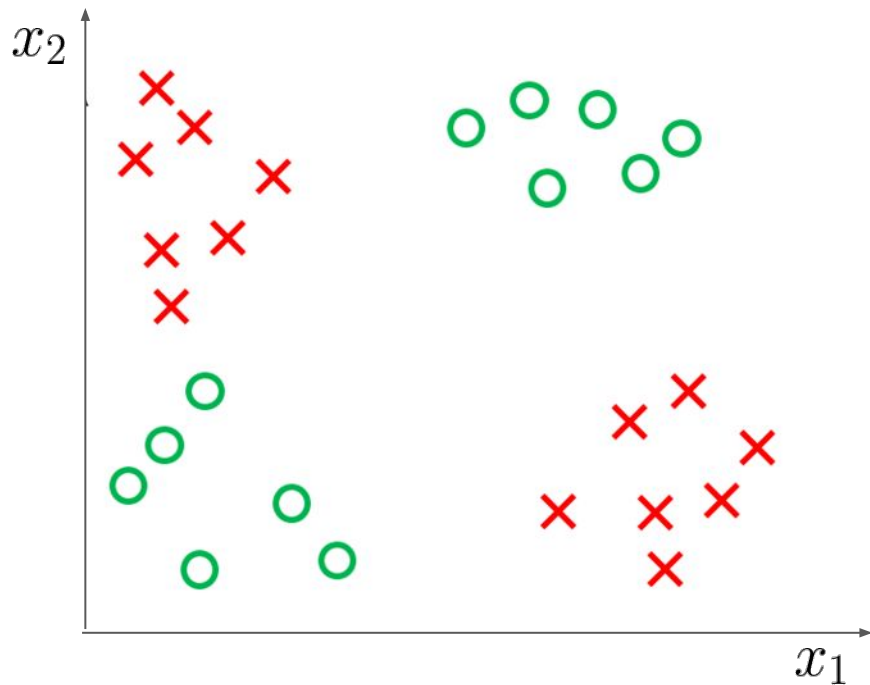
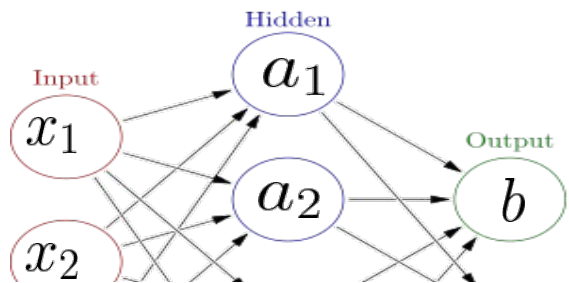
$$a_2 = -1x_1 + 1x_2$$

$$b = 4a_1 + 1a_2$$

$$= 4(1x_1 + 2x_2) + 1(-1x_1 + 1x_2)$$

$$= \underline{3x_1 + 3x_2}$$

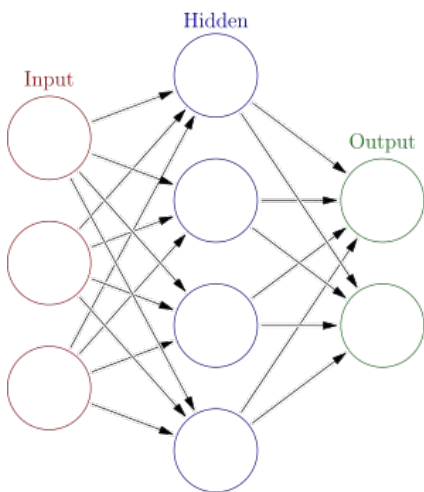
Linear separability and composition of functions



Non-linearity (more at https://en.wikipedia.org/wiki/Activation_function)

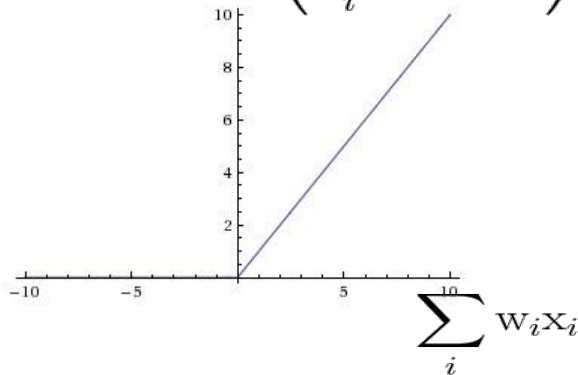
Non-linearities mostly process each activation independently.

One example is ReLU (rectified linear unit).

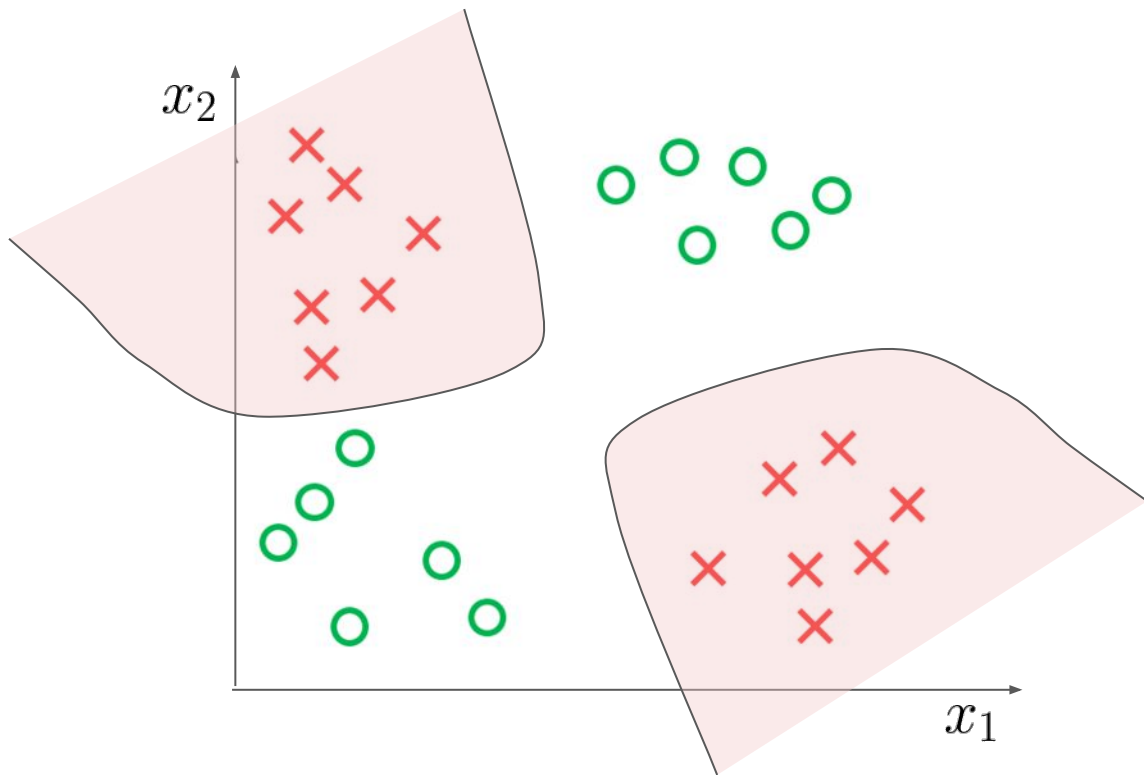
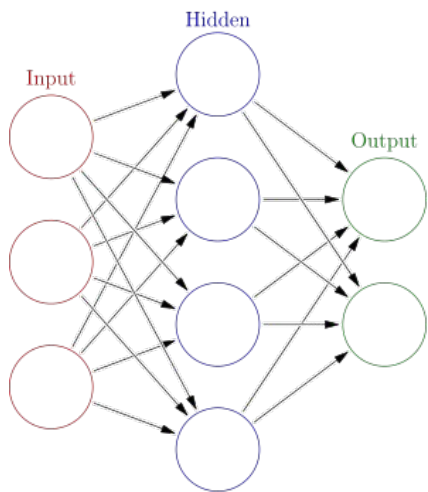


$$f(\mathbf{x}; \mathbf{w}) = \sigma \left(\sum_i w_i x_i \right)$$

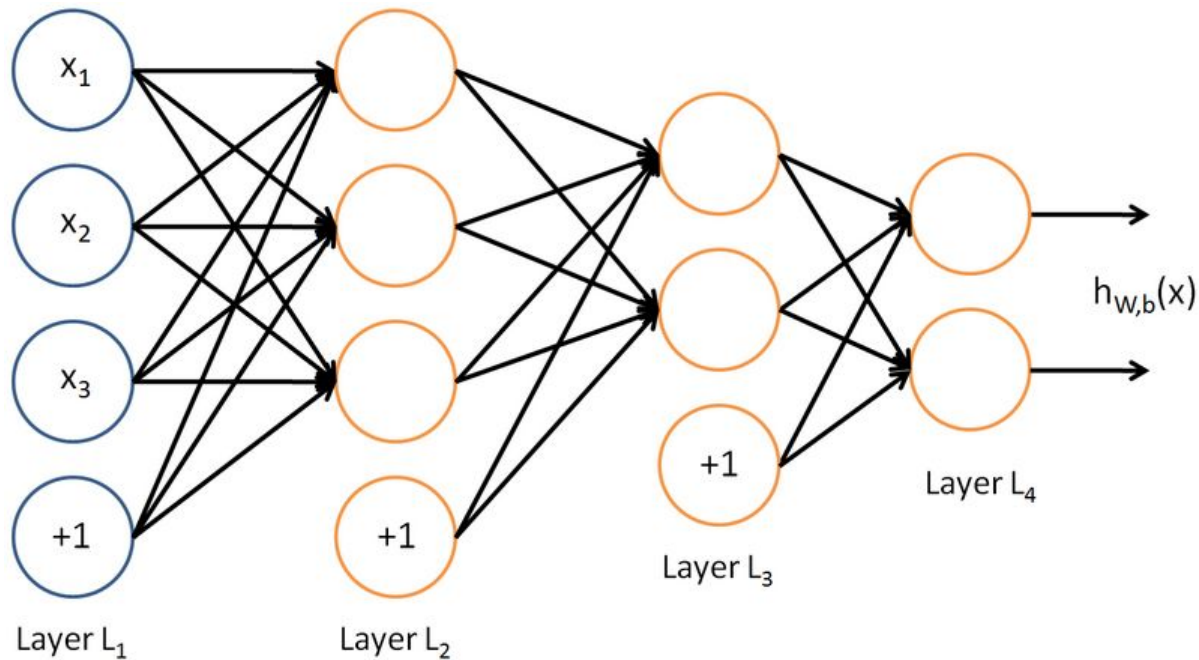
$$f(\mathbf{x}; \mathbf{w}) = \max \left(\sum_i w_i x_i, 0 \right)$$



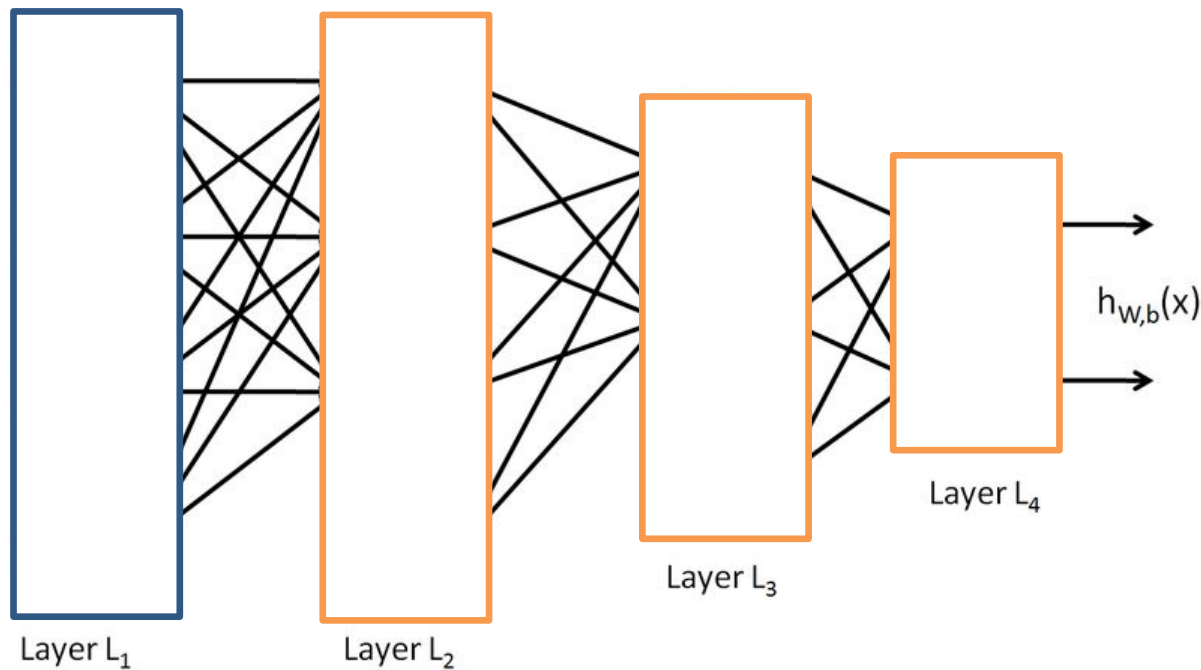
Linear separability and composition of functions



From neurons to layers

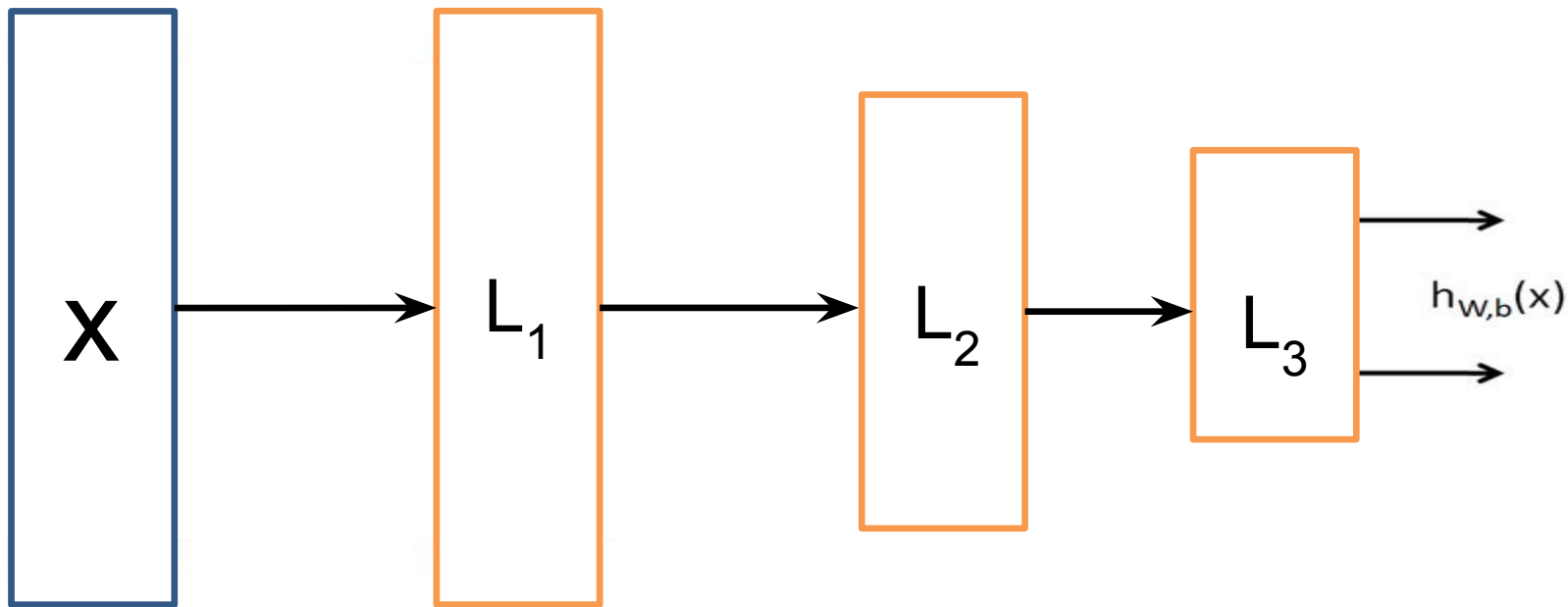


From neurons to layers



Layers

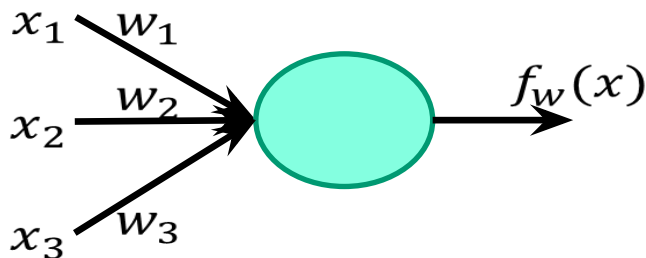
Often, we don't visualize the layers, but their inputs and outputs (activations).



From neurons to matrices

We would like to define layer of neurons as single mathematical operation (as a linear layer).

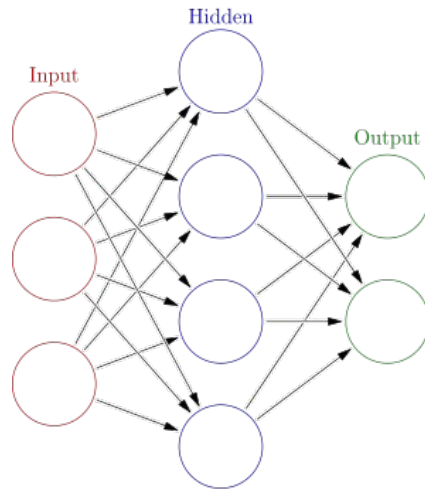
Solution: Linear algebra



$$f(\mathbf{x}; \mathbf{w}) = w_1x_1 + w_2x_2 + w_3x_3 \quad \longrightarrow \quad [1 \ 0 \ 3] \times \begin{bmatrix} 2 \\ -1 \\ -1 \end{bmatrix} = -1$$
$$= \sum_i w_i x_i$$

Linear layer

$$[1 \ 0 \ 3] \times \begin{bmatrix} 2 \\ -1 \\ -1 \end{bmatrix} = -1$$



$$[1 \ 0 \ 3] \times \begin{bmatrix} 2 \\ -1 \\ -1 \end{bmatrix} \begin{bmatrix} 0 \\ 2 \\ 0 \end{bmatrix} \begin{bmatrix} -1 \\ 0 \\ 0 \end{bmatrix} \begin{bmatrix} 1 \\ 1 \\ 0 \end{bmatrix} = -1, 0, -1, -1$$

Linear layer

$$[1 \ 0 \ 3] \times \begin{bmatrix} 2 \\ -1 \\ -1 \end{bmatrix} \begin{bmatrix} 0 \\ 2 \\ 0 \end{bmatrix} \begin{bmatrix} -1 \\ 0 \\ 0 \end{bmatrix} \begin{bmatrix} 1 \\ 1 \\ 0 \end{bmatrix} = -1, 0, -1, -1$$



$$[1 \ 0 \ 3] \times \begin{bmatrix} 2 & 0 & -1 & 1 \\ -1 & 2 & 0 & 1 \\ -1 & 0 & 0 & 0 \end{bmatrix} = [-1 \ 0 \ -1 \ -1]$$

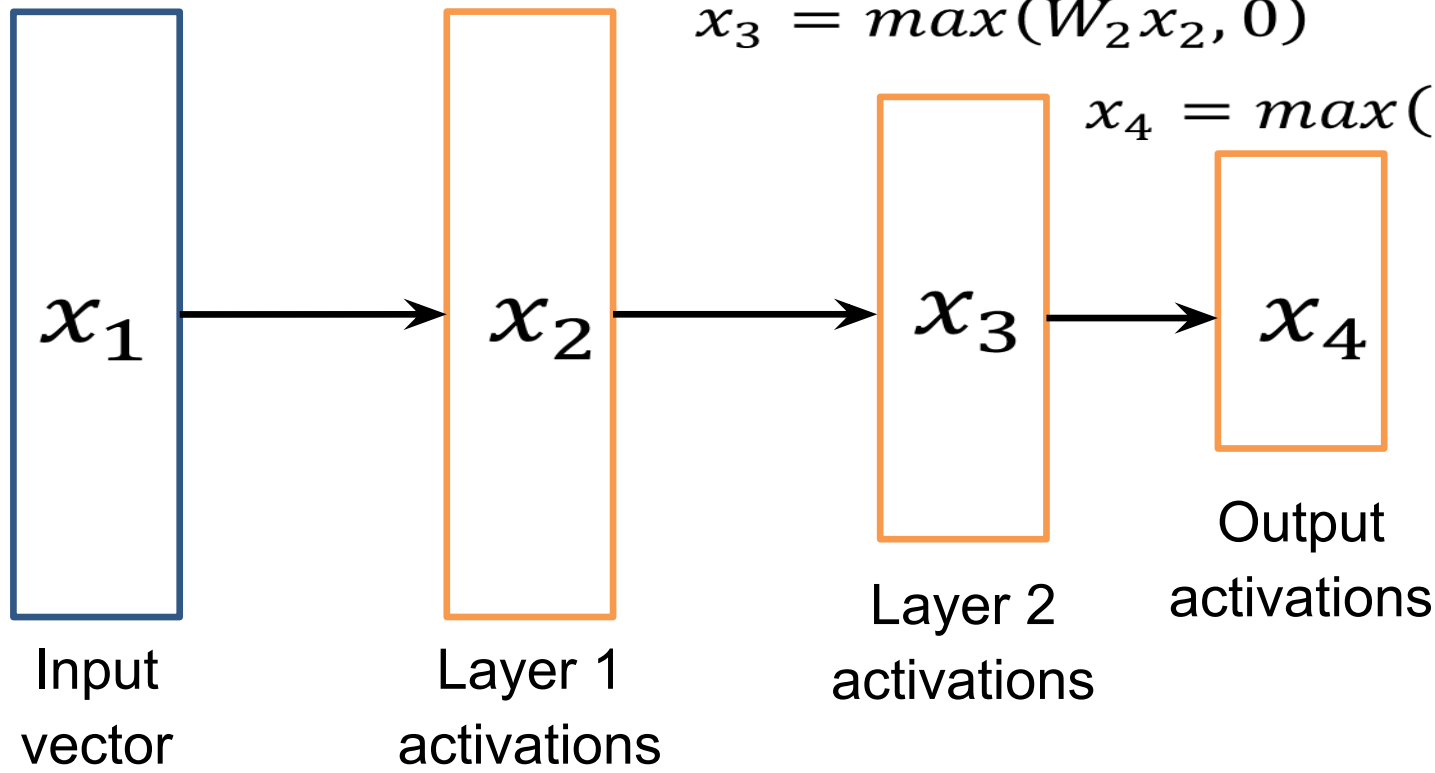
$$f(\mathbf{x}; \mathbf{W}) = \mathbf{W}\mathbf{x}$$

Neural network in code

$$x_2 = \max(W_1 x_1, 0)$$

$$x_3 = \max(W_2 x_2, 0)$$

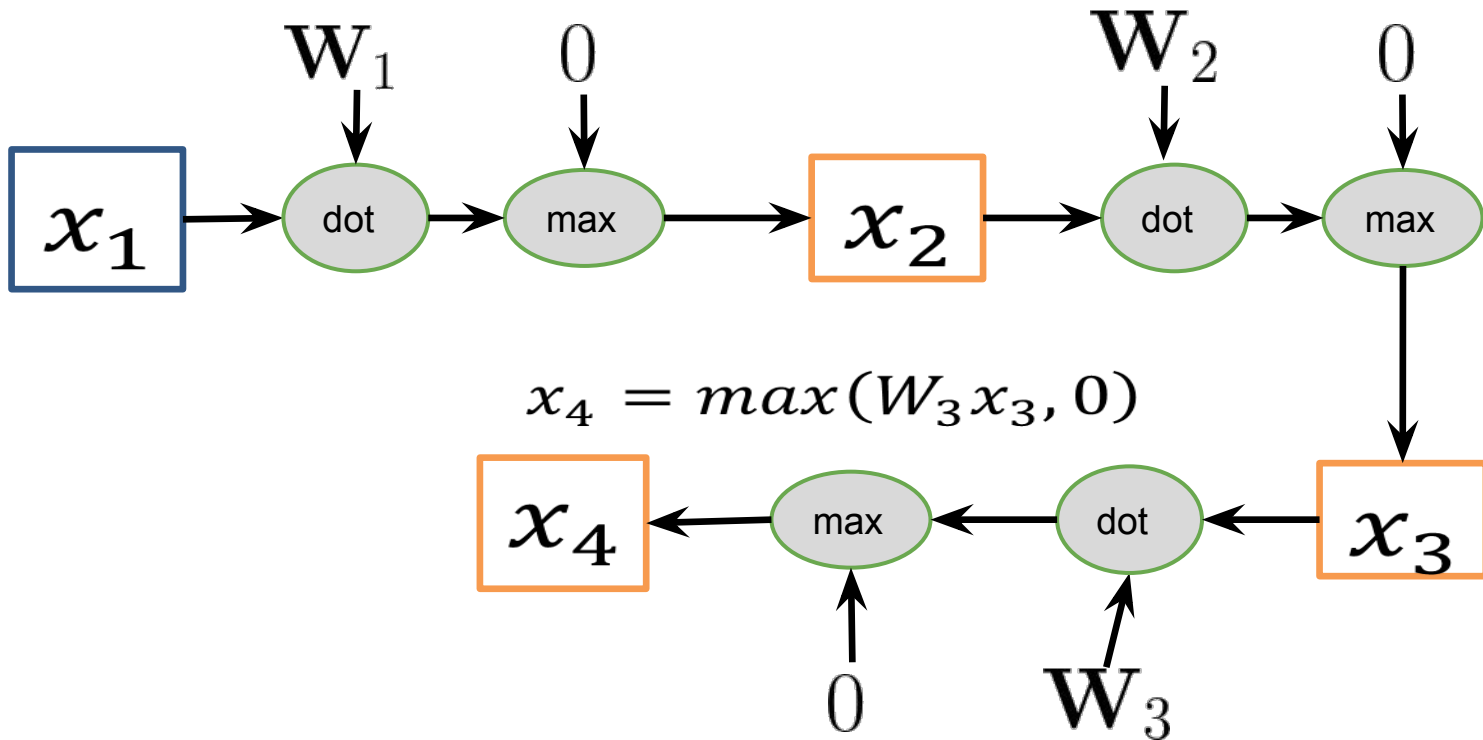
$$x_4 = \max(W_3 x_3, 0)$$



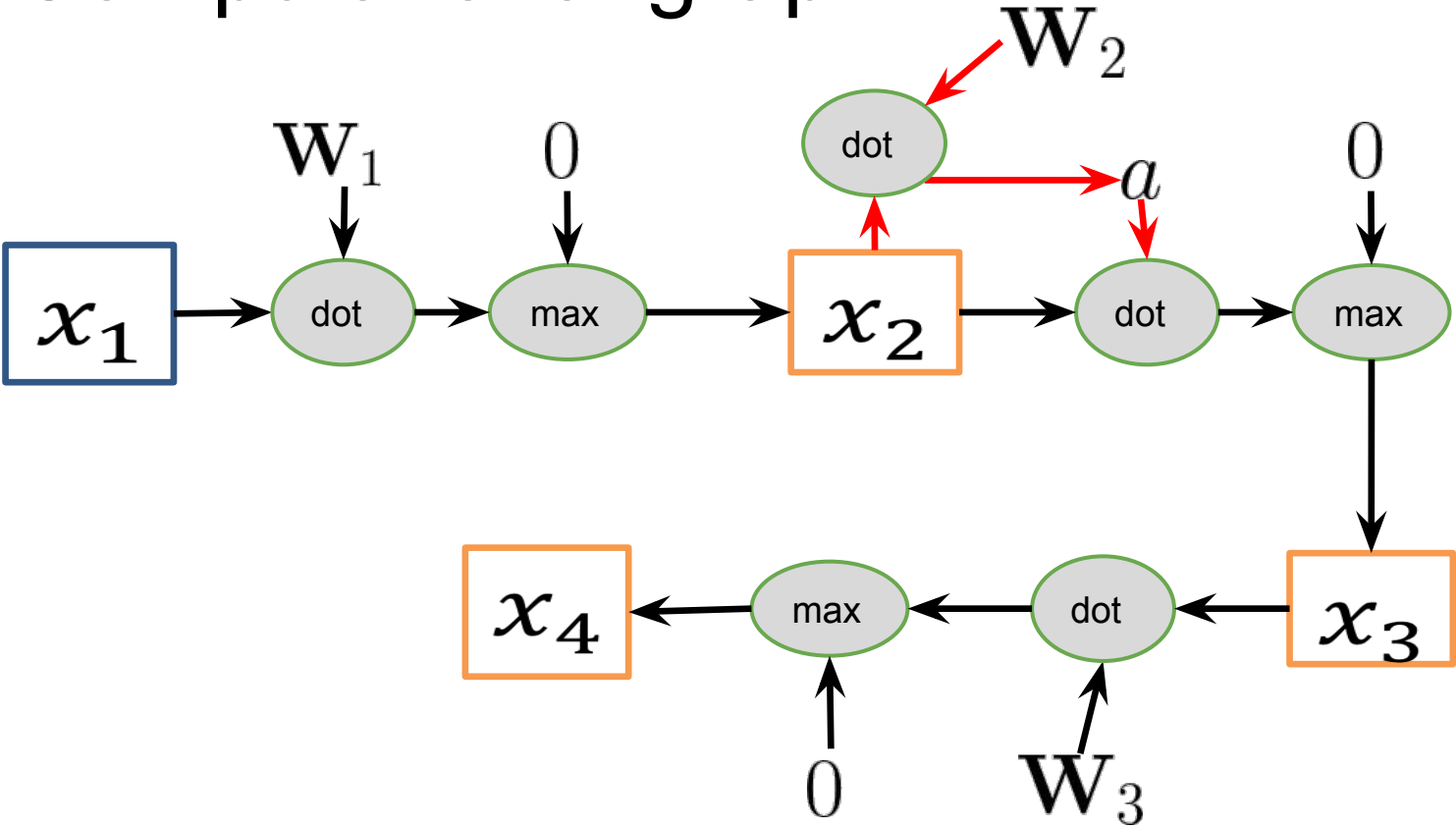
Computational graph

$$x_2 = \max(W_1 x_1, 0)$$

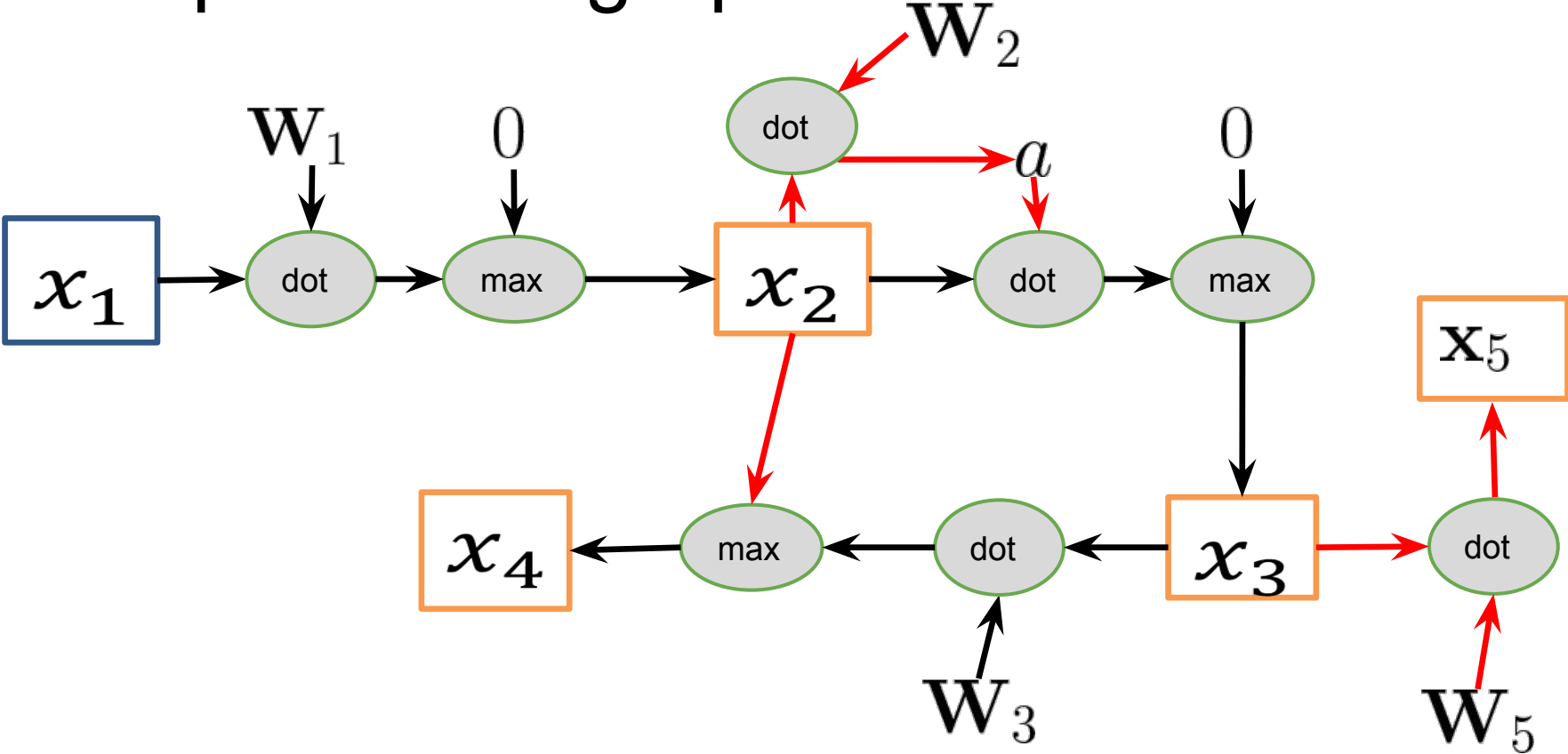
$$x_3 = \max(W_2 x_2, 0)$$



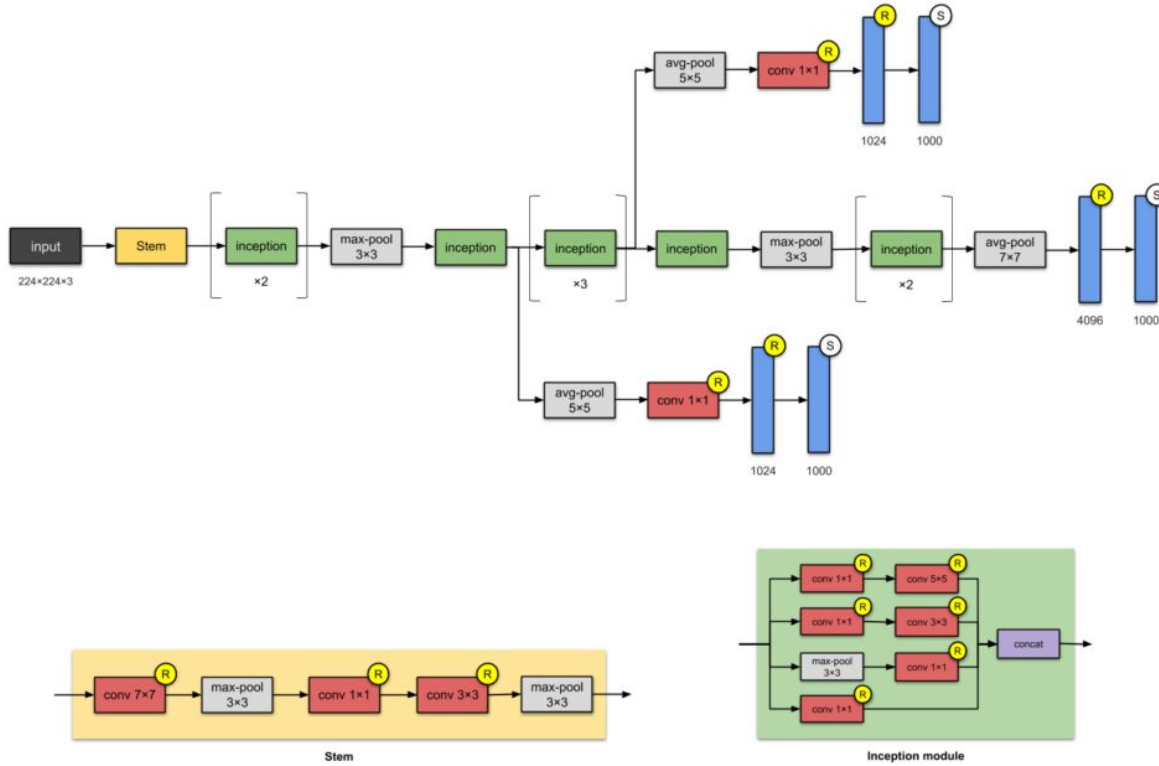
Computational graph



Computational graph



Real world example: Inception-V4



Source: raimibinkarim, Illustrated: 10 CNN Architectures

<https://makerspace.singapore.org/2019/09/illustrated-10-cnn-architectures/>

Loss functions

- During training, loss function measures how network output differs from the desired output.
- It is important to use the correct loss function for the task.
- It is important to use correct activation function for the tasks.
- Multiple loss/activation pairs can usually be used for a single task, but choice affects what the network learns.
(You change the tasks definition.)

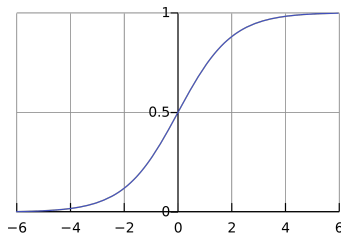
Binary classification

Goal: get $p(c | \text{data})$

Network output: $p(c=1|\text{data})$

Binary problem: $p(c=2|\text{data}) = 1 - p(c=1|\text{data})$

Activation function: sigmoid



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

What loss function?

Binary classification (same as in logistic regression)

Goal: get $p(c | \text{data})$

Network output: $p(c=0|\text{data})$

Binary problem: $p(c=1|\text{data}) = 1 - p(c=0|\text{data})$

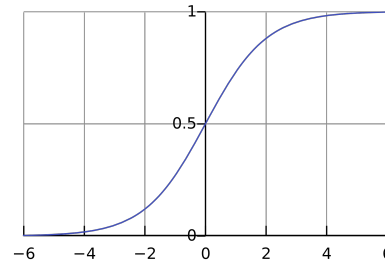
Activation function: sigmoid

Task maximize training set likelihood:

$$P(t|D) = \prod_D P(t|x_i) = \prod_D f(x_i)^{t_i} (1-f(x_i))^{1-t_i}$$

Minimize binary cross-entropy loss (log-likelihood):

$$J(f, D) = \sum_{i \in D} -\ln(P(t|x_i)) = \sum_{i \in D} t_i \ln(f(x_i)) + (1 - t_i) \ln(1 - f(x_i))$$



Binary classification - SVM loss

Goal:

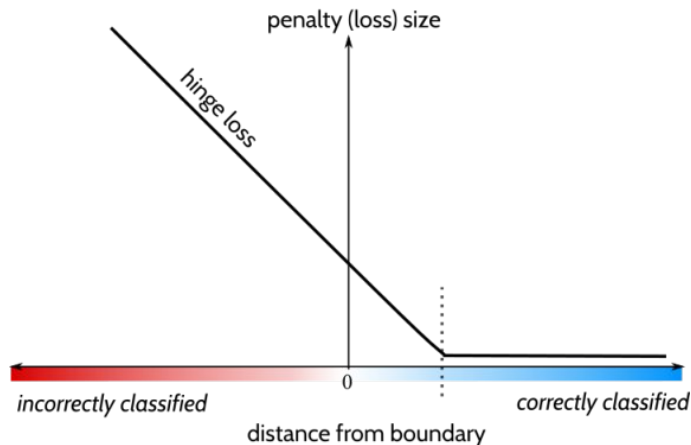
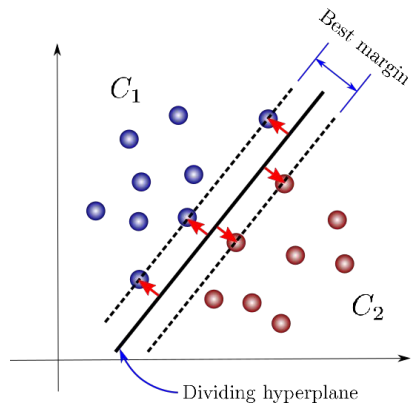
$$f(\mathbf{x}) > 1 \quad \text{if } \mathbf{c} = 1$$

$$f(\mathbf{x}) < -1 \quad \text{if } \mathbf{c} = -1$$

Activation function: linear (none)

$$J(f, D) = \sum_{i \in D} \max(0, 1 - c_i f(\mathbf{x}_i))$$

Loss function: Hinge loss



Multi-class classification

Goal: get $p(c \mid \text{data})$ with number of classes > 2



0.1	Michal
0.0	Martin
0.5	Petra
0.3	Katka
0.1	Honza
0.0	Bětka

Activation function: soft-max

$$\sigma_i(\mathbf{z}) = \frac{e^{z_i}}{\sum_j e^{z_j}}$$

Multi label classification

Goal: get $p(c | \text{data})$ for multiple binary problems



0.1	Asian
1.0	Male
0.9	Black hair
0.3	Actor
0.2	Child
1.0	Photographed outdoors

Activation function: sigmoid

Loss: Binary cross-entropy

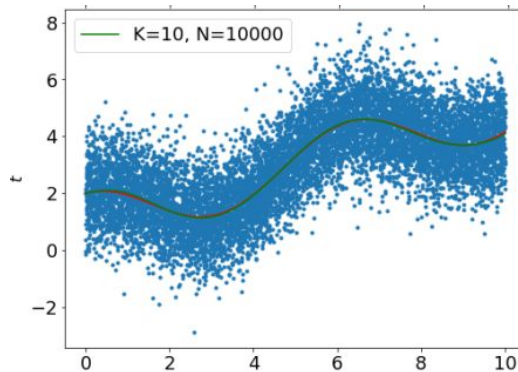
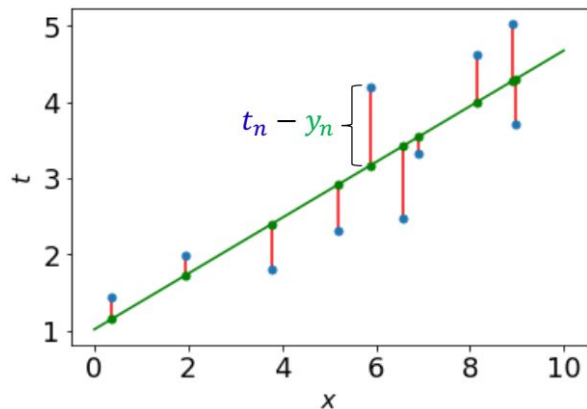
Regression

Goal: get $E[p(\text{value} \mid \text{data})]$ + gaussian noise in data

Activation: Linear (none)

Loss: Mean squared error

$$\sum_{i \in D} \frac{1}{2} (f(x_i) - t_i)^2$$



Regression

Goal: get $E[p(\text{value} \mid \text{data})]$ + lot of non-gaussian outliers

Activation: Linear (none)

Loss: Sum of Absolute Difference

$$\sum_{i \in D} |f(x_i) - t_i|$$

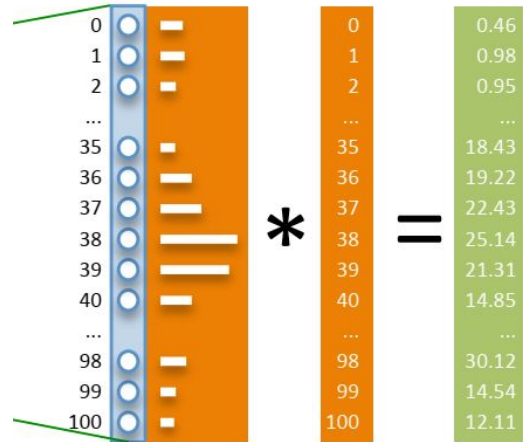
Regression - with discretization

Goal: get $p(\text{value} \mid \text{data})$

Activation: Softmax

Loss: Cross-entropy

5. Prediction



Softmax expected value

$$\Sigma = 38.4 \text{ years}$$

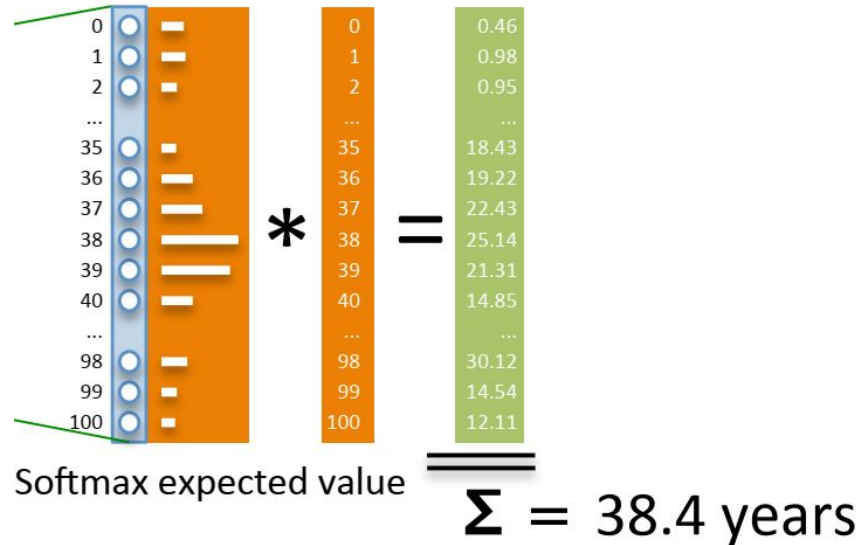
Regression - with discretization

Goal: get $E(p(\text{value} \mid \text{data}))$

Activation: Softmax
(with expectation sum)

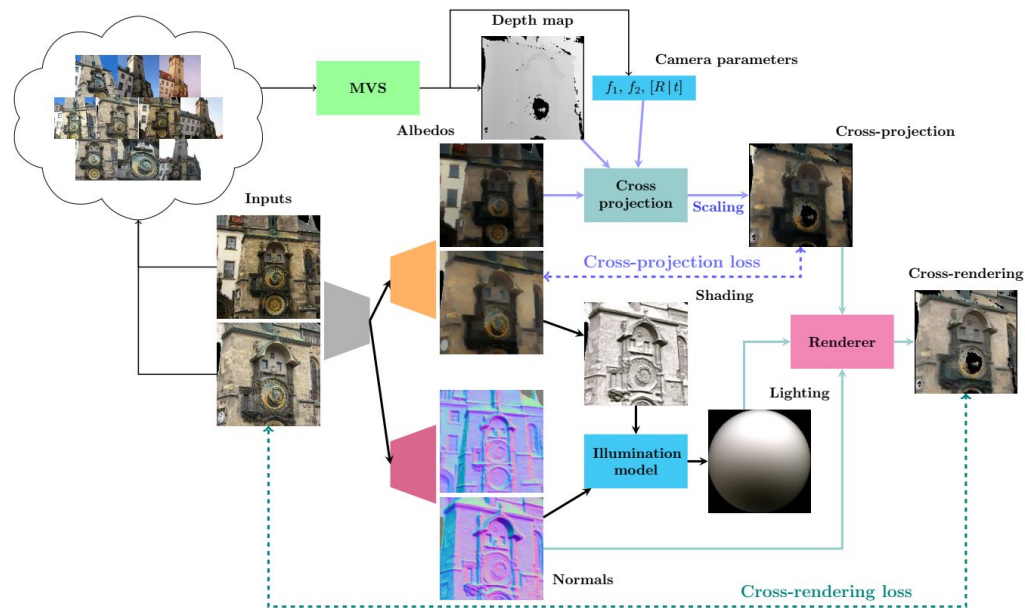
Loss: Mean squared
error?

5. Prediction



Weird loss functions

We train our network to minimise: $\ell = w_1 \ell_{\text{appearance}} + w_2 \ell_{\text{NM}} + w_3 \ell_{\text{albedo}} + w_4 \ell_{\text{cross-rend}} + w_5 \ell_{\text{albedo-smooth}} + w_6 \ell_{\text{albedo-pseudoSup}}$.



Training of neural networks

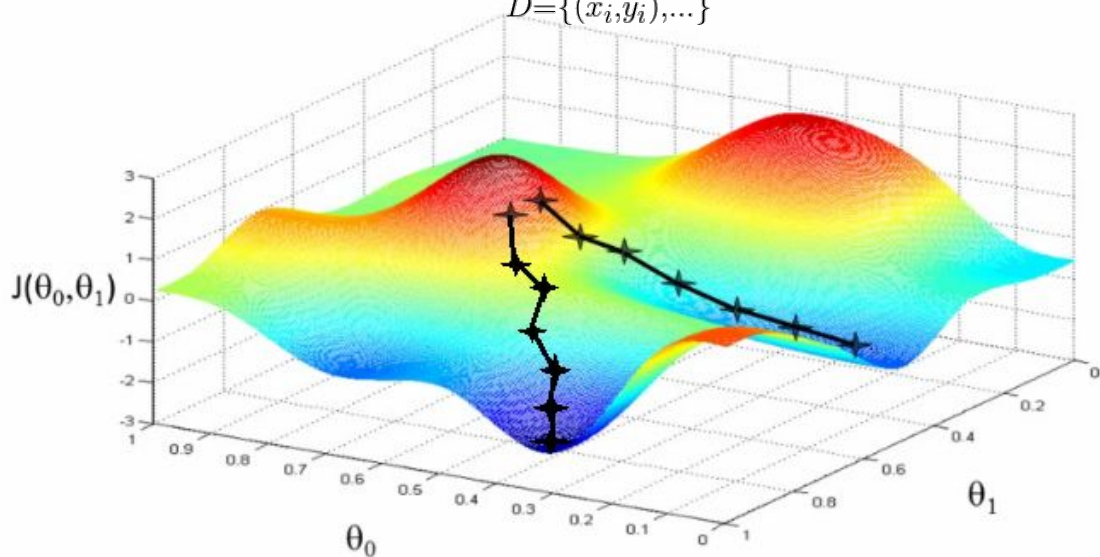
$$J(D, \theta) = \sum_{D=\{(x_i, y_i), \dots\}} \text{loss}(f(x_i, \theta), y_i)$$

$$\theta^* = \text{arg min}_{\theta} J(\mathbf{D}, \theta)$$

Gradient Descent

$$\theta_j \leftarrow \theta_j - \alpha \frac{\partial J(D, \theta)}{\partial \theta_j}$$

$$J(D, \theta) = \sum_{D=\{(x_i, y_i), \dots\}} \text{loss}(f(x_i, \theta), y_i)$$

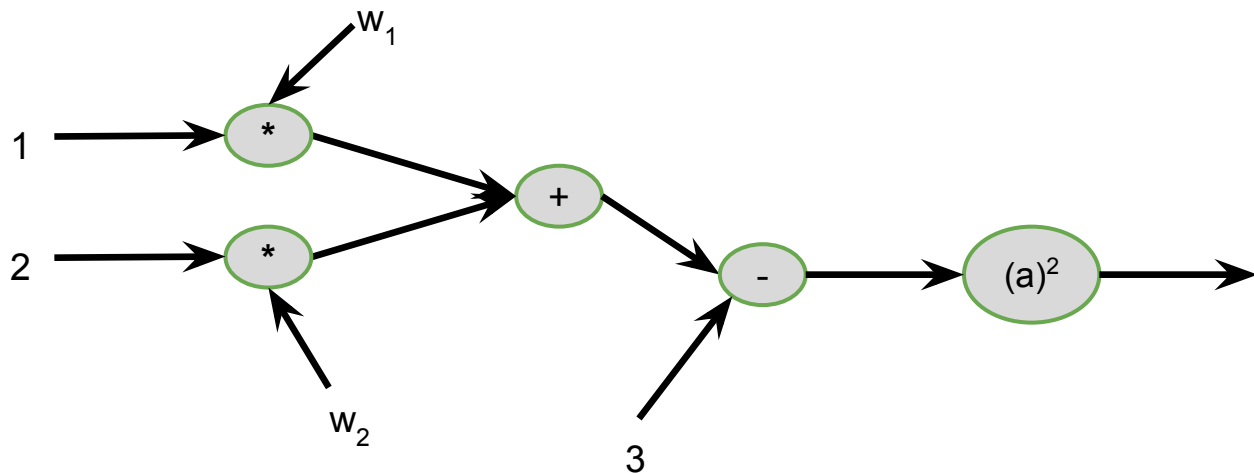


Gradients for training - chain rule

Need $\nabla_w J = \left[\frac{\partial J}{\partial w_1}, \frac{\partial J}{\partial w_2} \right]$

$$\frac{\partial (w_1 + 2w_2 - 3)^2}{\partial w_1} = \frac{\partial (a)^2}{\partial a} \frac{\partial (b - 3)}{\partial b} \frac{\partial (c + 2w_2)}{\partial c} \frac{\partial (1 * w_1)}{\partial w_1}$$

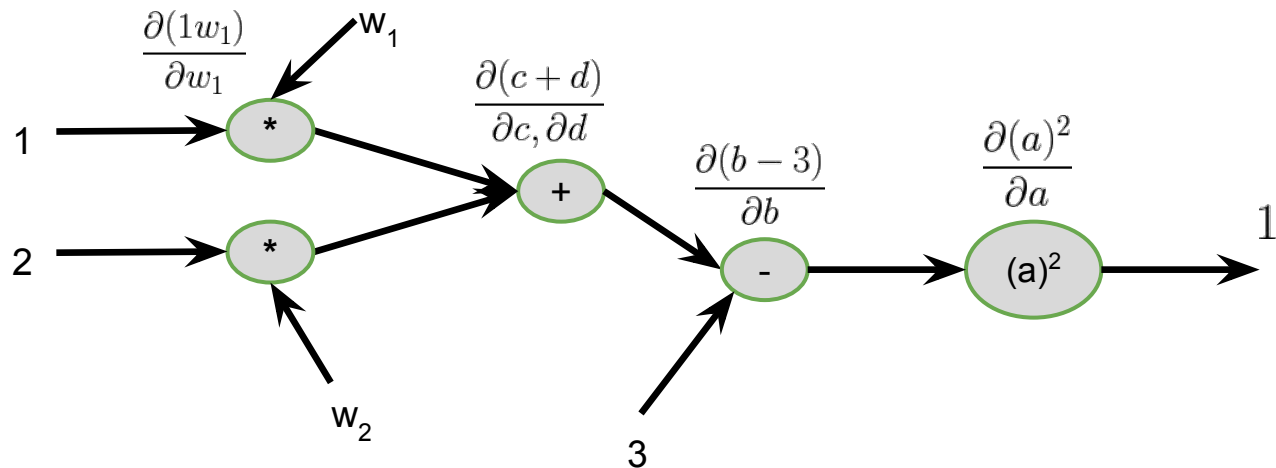
$$\frac{\partial (w_1 + 2w_2 - 3)^2}{\partial w_2} = \frac{\partial (a)^2}{\partial a} \frac{\partial (b - 3)}{\partial b} \frac{\partial (w_1 + c)}{\partial c} \frac{\partial (2 * w_2)}{\partial w_2}$$



Gradient for training

Need $\nabla_w J = \left[\frac{\partial J}{\partial w_1}, \frac{\partial J}{\partial w_2} \right]$

$$\frac{\partial(1w_1 + 2w_2 - 3)^2}{\partial w_1} = \frac{\frac{\partial(a)^2}{\partial a} \frac{\partial(b-3)}{\partial b}}{\frac{\partial(a)^2}{\partial a} \frac{\partial(b-3)}{\partial b}} \frac{\partial(c+2w_2)}{\partial c} \frac{\partial(1 * w_1)}{\partial w_1}$$
$$\frac{\partial(1w_1 + 2w_2 - 3)^2}{\partial w_2} = \frac{\frac{\partial(a)^2}{\partial a} \frac{\partial(b-3)}{\partial b}}{\frac{\partial(a)^2}{\partial a} \frac{\partial(b-3)}{\partial b}} \frac{\partial(1w_1 + c)}{\partial c} \frac{\partial(2 * w_2)}{\partial w_2}$$



Full training

- Forward pass

$$x_2 = \max(W_1 x_1, 0)$$

$$x_3 = \max(W_2 x_2, 0)$$

$$x_4 = \max(W_3 x_3, 0)$$

- Compute loss

$$loss = (x_4 - 1)^2$$

- Backward pass

- Compute derivatives of objective fce. with respect to net. parameters using chain rule / backpropagation

$$\theta_j \leftarrow \theta_j - \alpha \frac{\partial J(D, \theta)}{\partial \theta_j}$$

- Update parameters

Pytorch introduction / basics of tensors

<https://colab.research.google.com/drive/1ytNsahwdl4FJDkXmaw7CCQQJqQDiAxJt?usp=sharing>

The same as Numpy with

- Automatic differentiation

- GPU support

- Neural network layers

- Loss functions

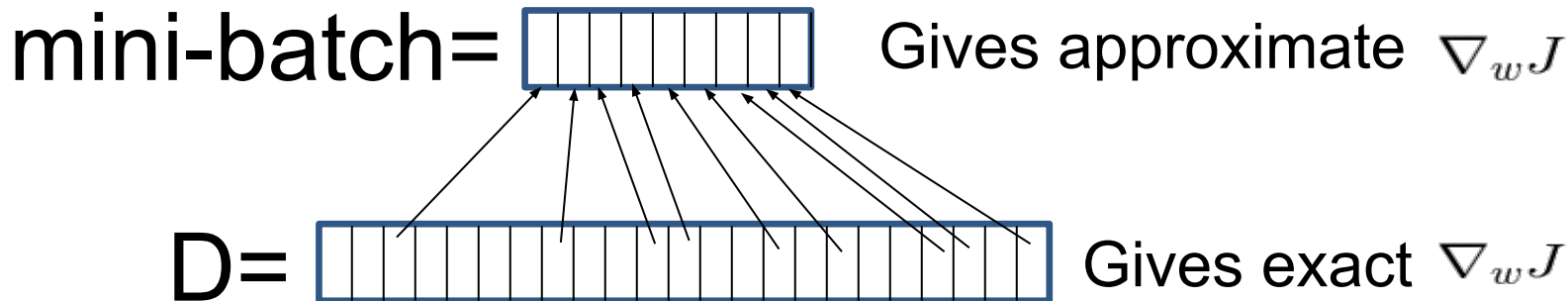
- Optimization algorithms for neural networks

- Datasets, data augmentation, predefined networks, pretrained network, ...

Stochastic/mini-batch gradient descent

- Large datasets -> too much computation per network update -> inefficient training
- Estimate gradient on a small random part of the training set (mini-batch) and update network
- Stochastic Gradient Descent

$$J(D, \theta) = \sum_{D=\{(x_i, y_i), \dots\}} \text{loss}(f(x_i, \theta), y_i)$$



Learning rate and adaptive algorithms

- High learning rate -> divergence
- Low learning rate -> slow training
- For SGD, optimal learning rate depends on
 - Network
 - Objective, consider gradients for $J()$ and $10000*J()$
- You have to try different learning rates

$$\theta_j \leftarrow \theta_j - \alpha \frac{\partial J(D, \theta)}{\partial \theta_j}$$

- Solution:
 - Adaptive GD algorithms which normalize learning rate
 - Adam, AdaDelta, ...

Train good model and guess real-world performance

- Training set
 - Train models on this set
 - In domain data, similar problems, synthetic, generated, augmented
- Validation set
 - Evaluate models, tune hyperparameters, choose alg. and architecture
 - Same (iid) as training set?
- Test set
 - Estimate real-world performance
 - Data ideally from real world
- Loss is usually different than the target metric
 - Loss can improve while “accuracy” gets worse

Review

- What is a neural network?
- What are “trainable/learnable” parameters?
- What is loss function?
- What is gradient?
- GD and SGD?
- What is a training set in supervised learning?
- Shallow learning vs deep learning?
- Linear neuron and matrix multiplication.
- Non-linear activations.
- Layers.
- Computational graph.

Review

- Basic loss functions for classification and regression.
- What is training - objective function minimization.
- Backward pass and gradient descent rule.
- SGD and its benefits.
- What are train set, validation set and test set?

Neural networks for structured data



Local cues



Local cues



Image representations



Layers for images?

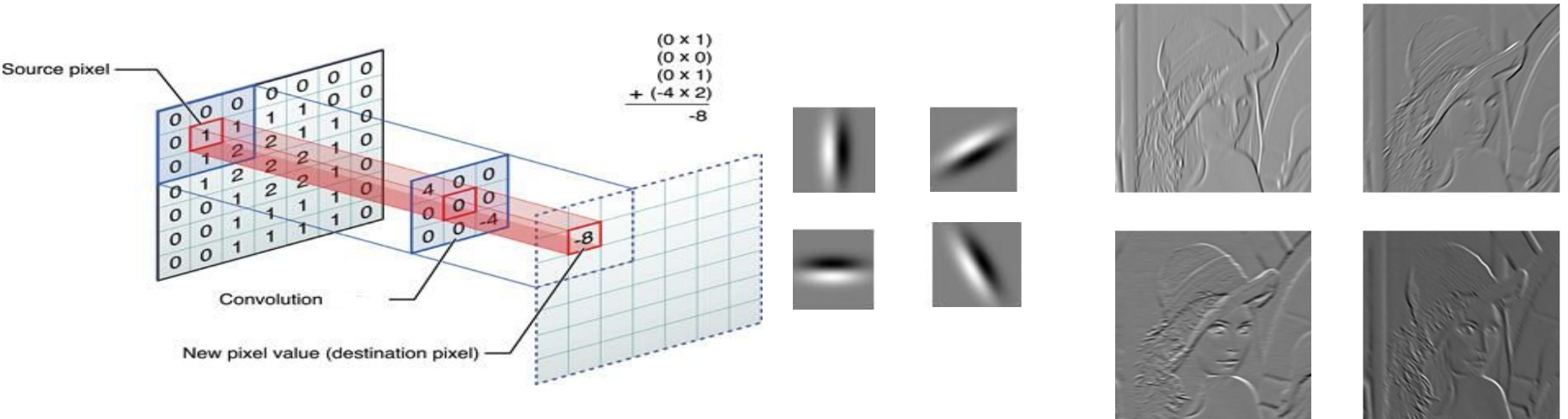
- Which mathematical operation fits?
 - Is local
 - Processes each part of input the same way
 - Is linear
 - Is the building block of neural networks for images, sound, ...

History of Convolutional Networks

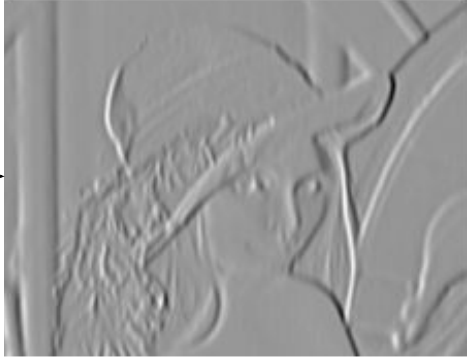
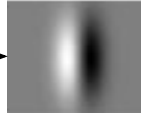
- **1980** – Kuniyuki Fukushima – **Neocognitron: A Self-organizing Neural Network Model for a Mechanism of Pattern Recognition Unaffected by Shift in Position**
- **1986** – Hinton, Rumelhart, Williams, McClelland - **Backpropagation**
- **1998** – LeCun et al.: Gradient-Based Learning Applied to Document Recognition
- **2012** – Krizhevsky et al.: ImageNet Classification with Deep Convolutional Neural Networks
- **2013-today** – explosion of applications

Convolution

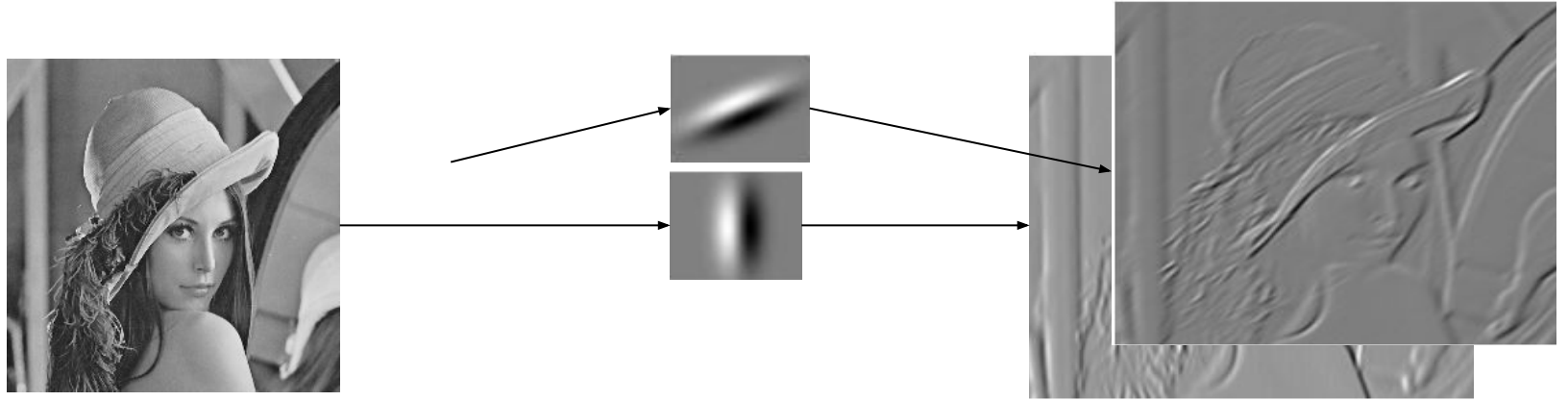
- Linear operation
 - (matrix multiplication on small neighborhood)
- Output: Map of local similarity of input and conv. kernel



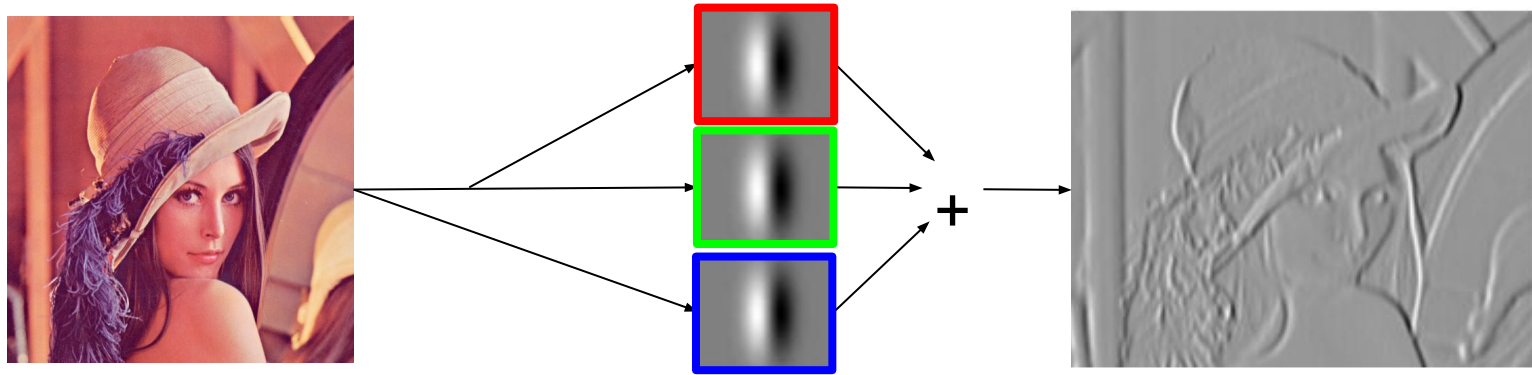
Convolutional layer



Convolutional layer

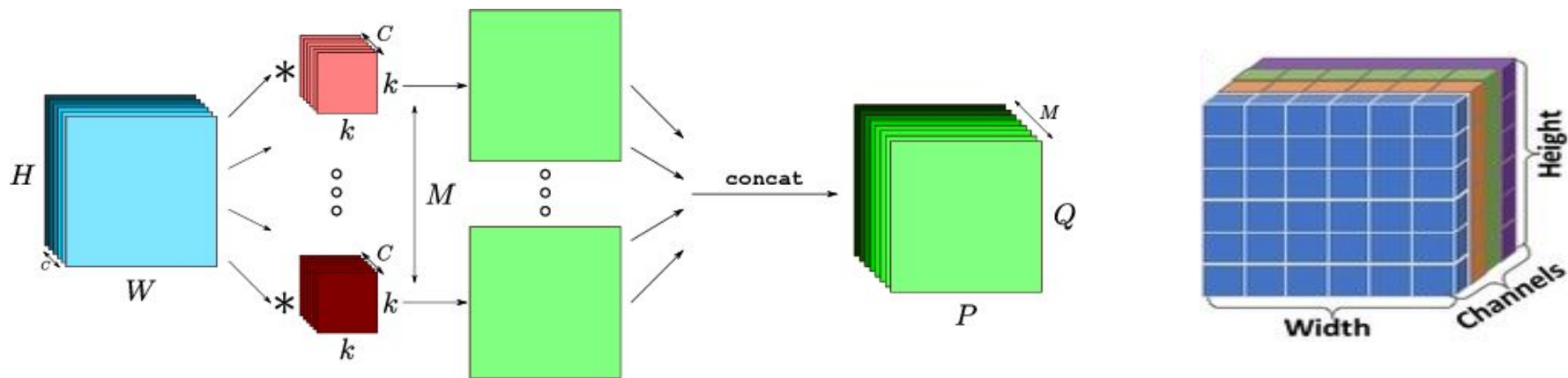


Convolutional layer



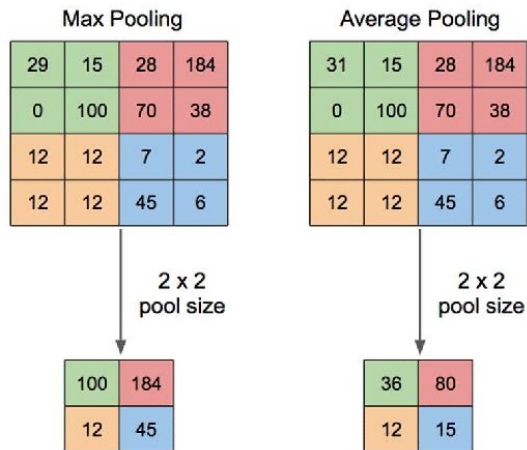
Convolution

- Channels \rightarrow activation map

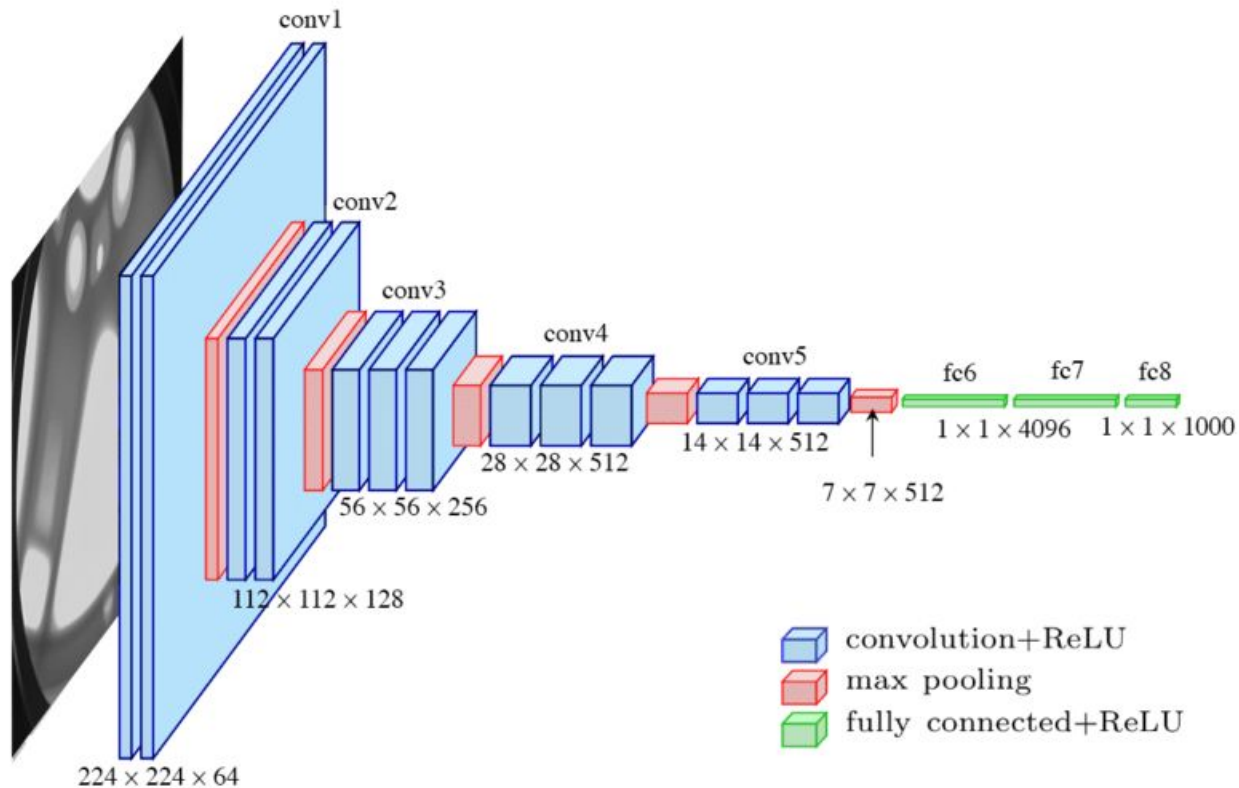


Pooling layer

- Downsampling
 - Local pixels are aggregated into a single pixel
 - Spatial resolution is reduced, channels stay the same
 - Reduction by AVERAGE or MAX functions.
 - MAX-pooling directly encourages translation invariance

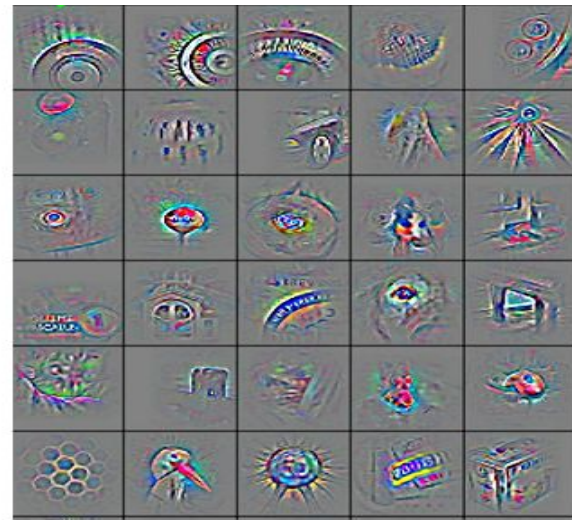
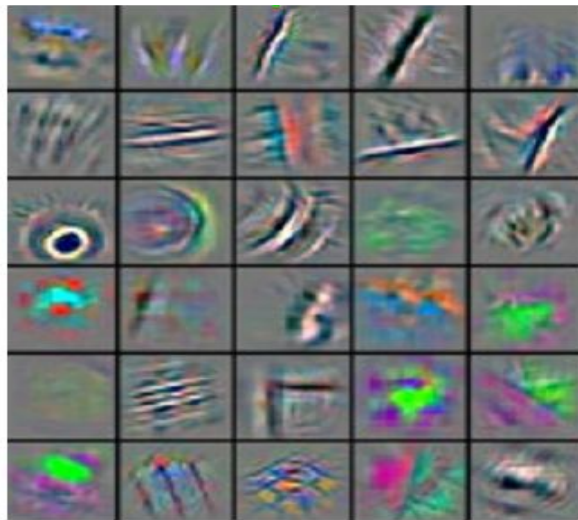
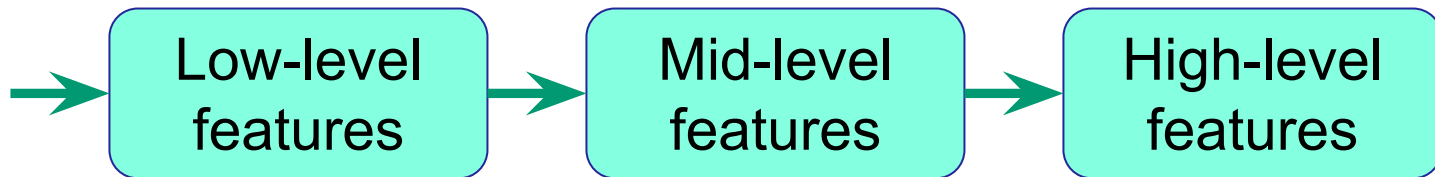


Network architecture (VGG net)



Hierarchy of features

(and human visual system analogy)



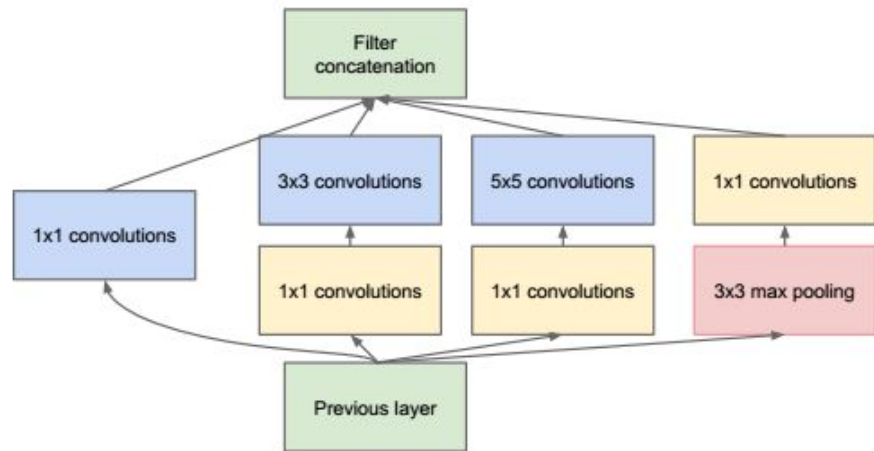
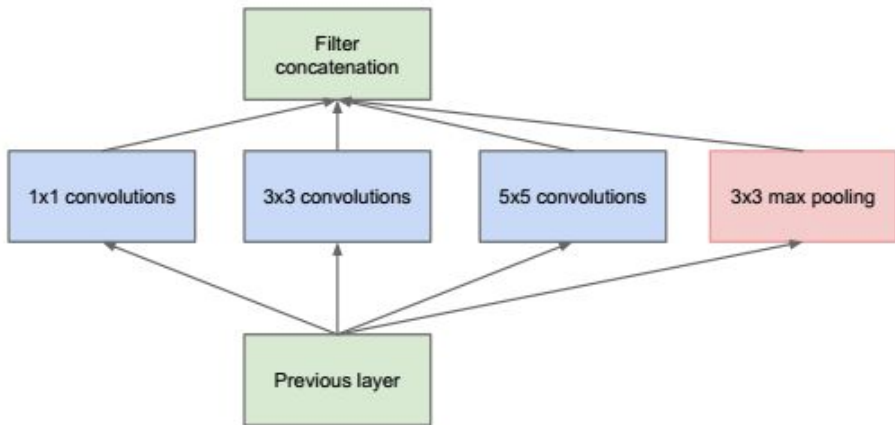
Very Deep Convolutional Networks

ConvNet Configuration					
A	A-LRN	B	C	D	E
11 weight layers	11 weight layers	13 weight layers	16 weight layers	16 weight layers	19 weight layers
input (224 × 224 RGB image)					
conv3-64	conv3-64 LRN	conv3-64 conv3-64	conv3-64 conv3-64	conv3-64 conv3-64	conv3-64 conv3-64
maxpool					
conv3-128	conv3-128	conv3-128 conv3-128	conv3-128 conv3-128	conv3-128 conv3-128	conv3-128 conv3-128
maxpool					
conv3-256 conv3-256	conv3-256 conv3-256	conv3-256 conv3-256	conv3-256 conv3-256 conv1-256	conv3-256 conv3-256 conv3-256	conv3-256 conv3-256 conv3-256 conv3-256
maxpool					
conv3-512 conv3-512	conv3-512 conv3-512	conv3-512 conv3-512	conv3-512 conv3-512 conv1-512	conv3-512 conv3-512 conv3-512	conv3-512 conv3-512 conv3-512 conv3-512
maxpool					
conv3-512 conv3-512	conv3-512 conv3-512	conv3-512 conv3-512	conv3-512 conv3-512 conv1-512	conv3-512 conv3-512 conv3-512	conv3-512 conv3-512 conv3-512 conv3-512
maxpool					
FC-4096					
FC-4096					
FC-1000					
soft-max					

Karen Simonyan, Andrew Zisserman: Very Deep Convolutional Networks for Large-Scale Image Recognition,

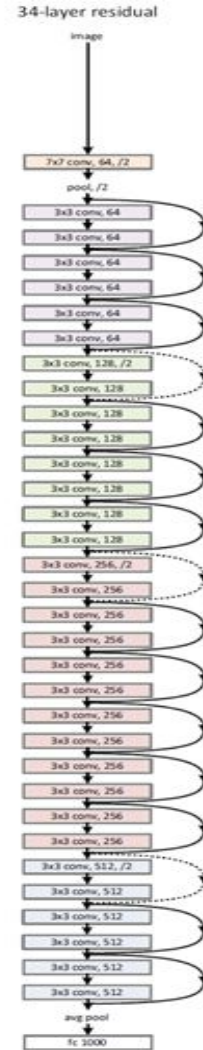
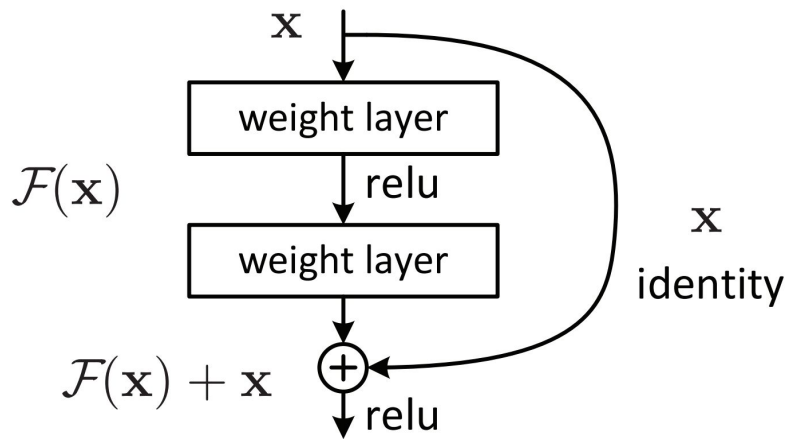
GoogLeNet - Inception Layer

- Winner of ILSVRC 2014

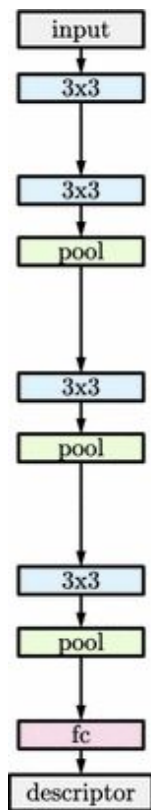


Szegedy et al.: Going deeper with convolutions

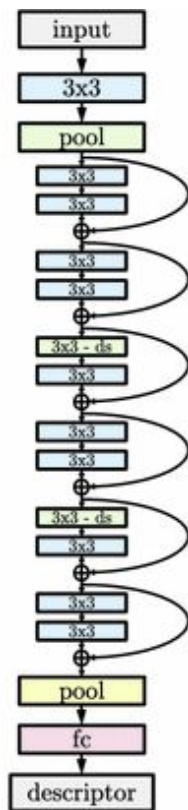
Residual networks



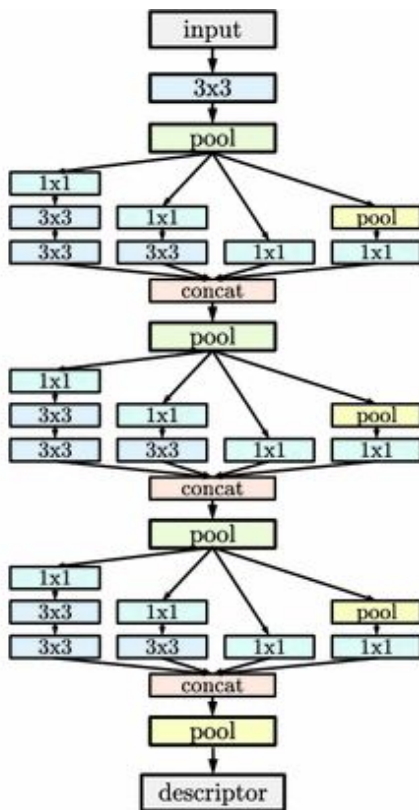
Combinations



(a)



(b)



(c)

Generalization

- **Good model structure**
- Multi-task learning
- Pre-training
- Domain adaptation
- More data - augmentation, synthetic data
- Model size

Regularization and overfitting

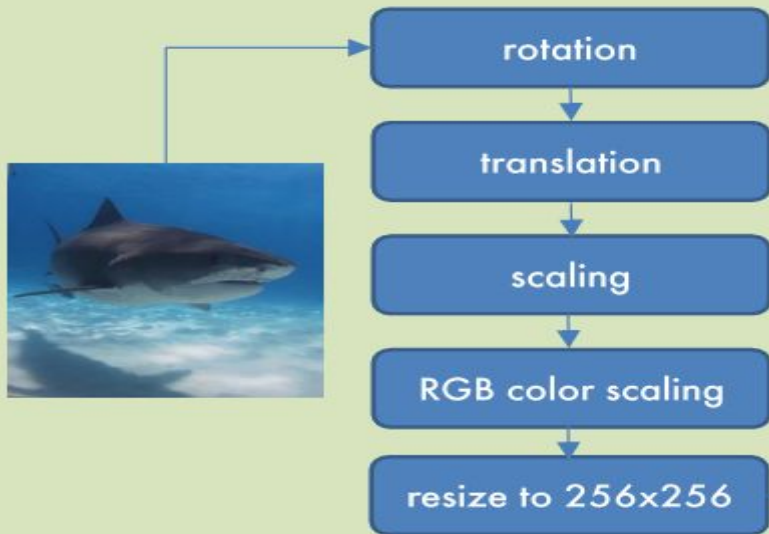
- Deep convolutional networks overfit surprisingly late
- Early stopping (when training)
- Weight regularization (weight decay)
 - does not help much (with ReLU)
- Activation regularization

$$J(D, \theta) = \left(\sum_D \text{loss}(f_\theta(x_i), y_i) \right) + \lambda \|\theta\|^2 + \beta \|\text{activation}\|^2$$

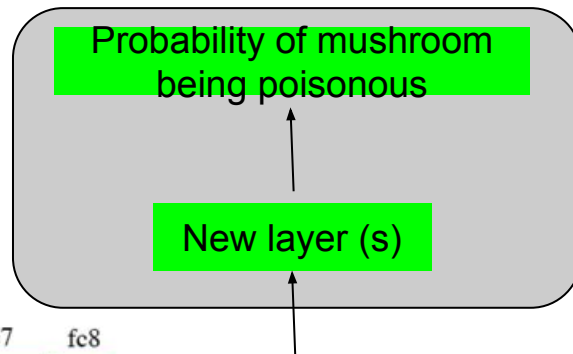
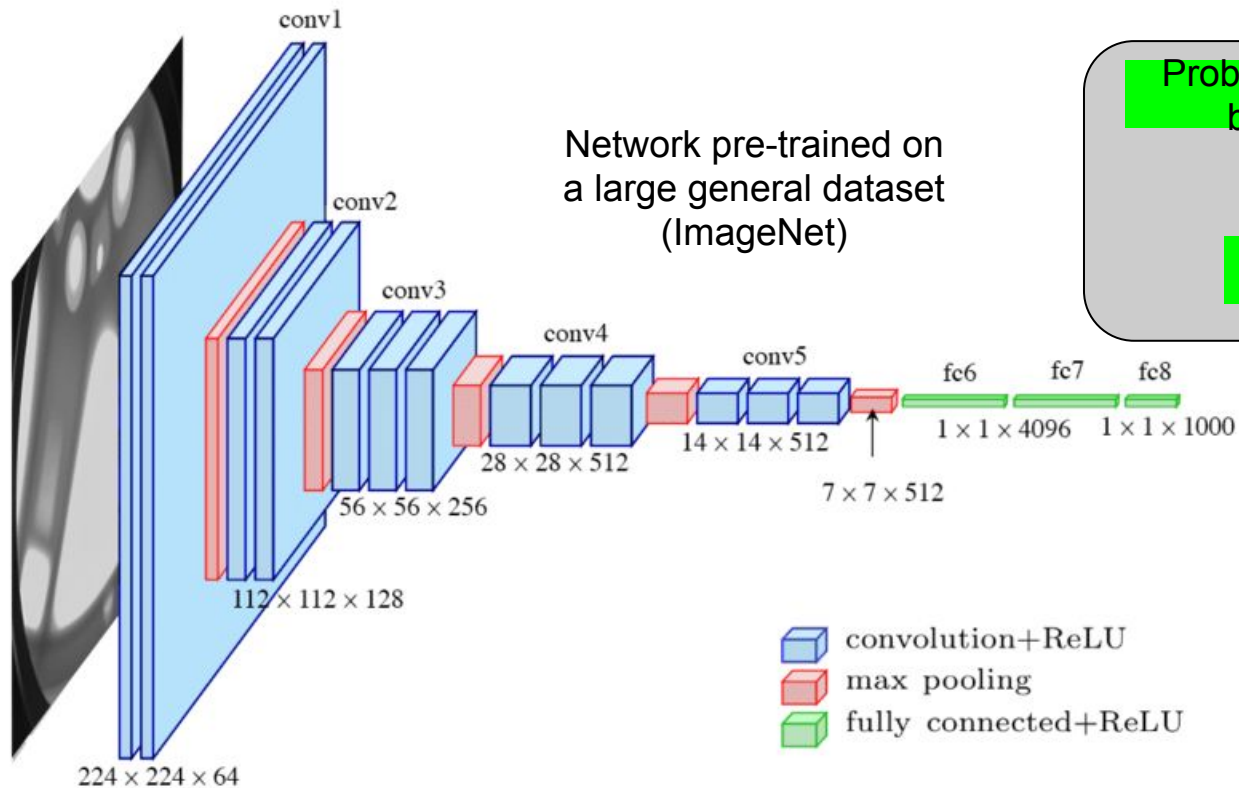
- **Dropout** helps a lot
 - Sets random activations to 0

Data augmentation

Image transformations



Finetuning



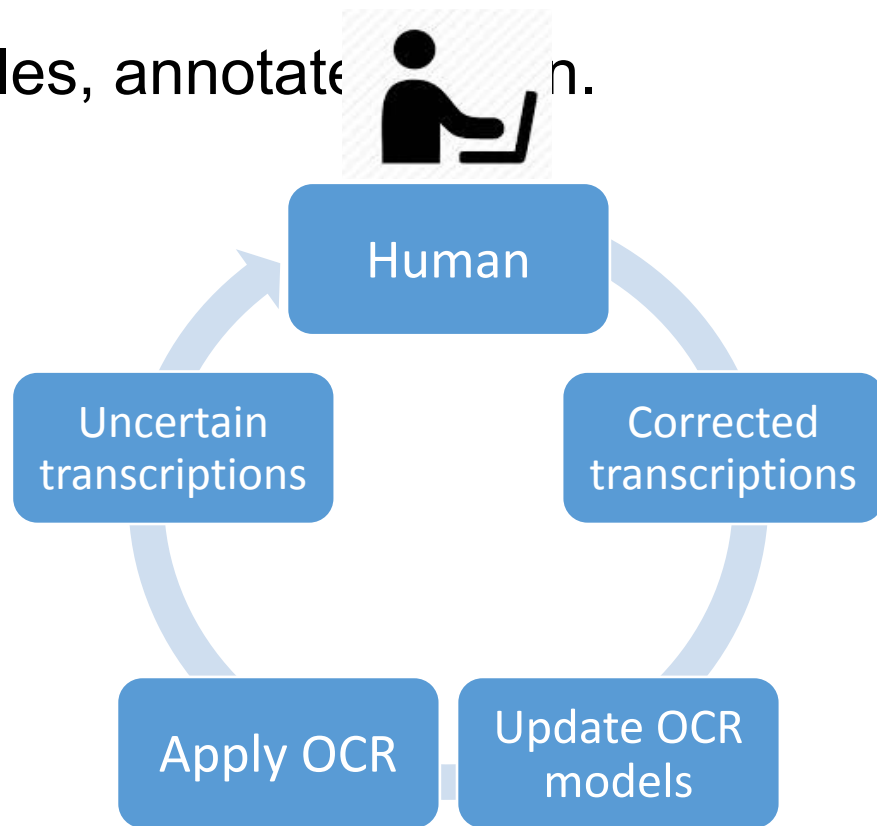
Active learning

Make annotations effective.

Use model to select hard examples, annotate  n.

Use prediction confidence.

f horrors and desolation.	29s to 31s.
f horrors and desolation.	29s to 31s
en He gouwen woonach	ИМОТЬ СЕ ИМОТЬ СЕ
kse-tiende	поставнувало, this Evening Tuesday.
kse-tiende	поставнувало this Evening Tuesday
niets	are the working weet men too veel AMINER will
niets	are the working weet men too veel AMINER will



Weird learning

Semi supervised learning - annotated and unannotated data

Active learning - pick good data to annotate

Self-supervised learning - automatically generate labels (mostly pre-training)

Transfer learning - network for one task is fine-tuned on another task

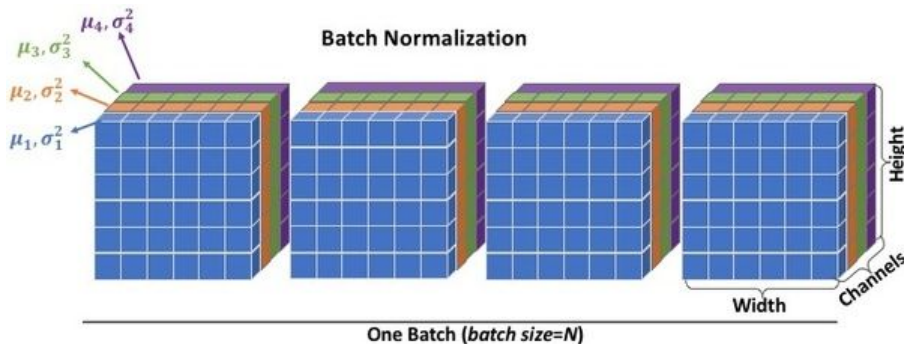
Knowledge distillation - usually a good large model generates “labels” for smaller one

Domain adaptation - the task remains the same, but data changes (day/night)

Self training - model generates labels and is trained on them (semi-supervised)

Federated learning - distributed learning on private data

Batch normalization



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

Parameters to be learned: γ, β

Output: $\{y_i = \text{BN}_{\gamma, \beta}(x_i)\}$

$$\mu_{\mathcal{B}} \leftarrow \frac{1}{m} \sum_{i=1}^m x_i \quad // \text{ mini-batch mean}$$

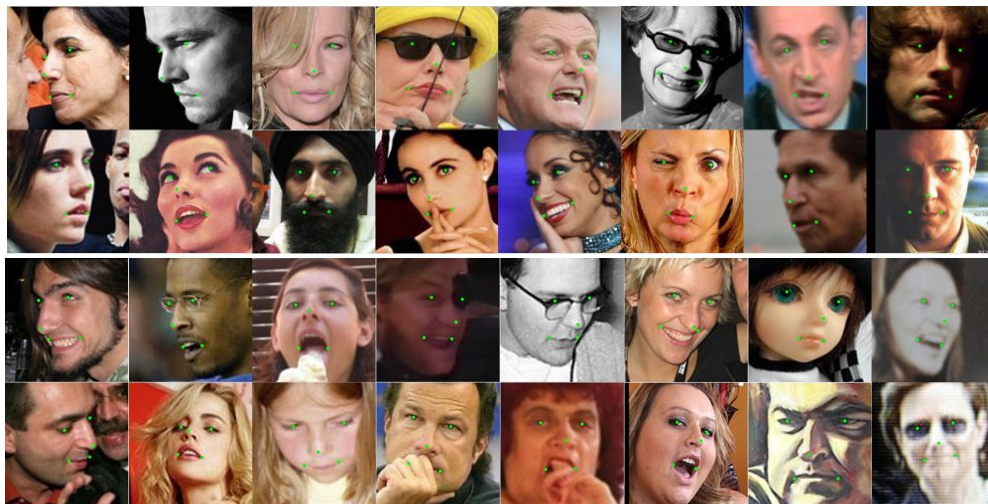
$$\sigma_{\mathcal{B}}^2 \leftarrow \frac{1}{m} \sum_{i=1}^m (x_i - \mu_{\mathcal{B}})^2 \quad // \text{ mini-batch variance}$$

$$\hat{x}_i \leftarrow \frac{x_i - \mu_{\mathcal{B}}}{\sqrt{\sigma_{\mathcal{B}}^2 + \epsilon}} \quad // \text{ normalize}$$

$$y_i \leftarrow \gamma \hat{x}_i + \beta \equiv \text{BN}_{\gamma, \beta}(x_i) \quad // \text{ scale and shift}$$

Facial alignment

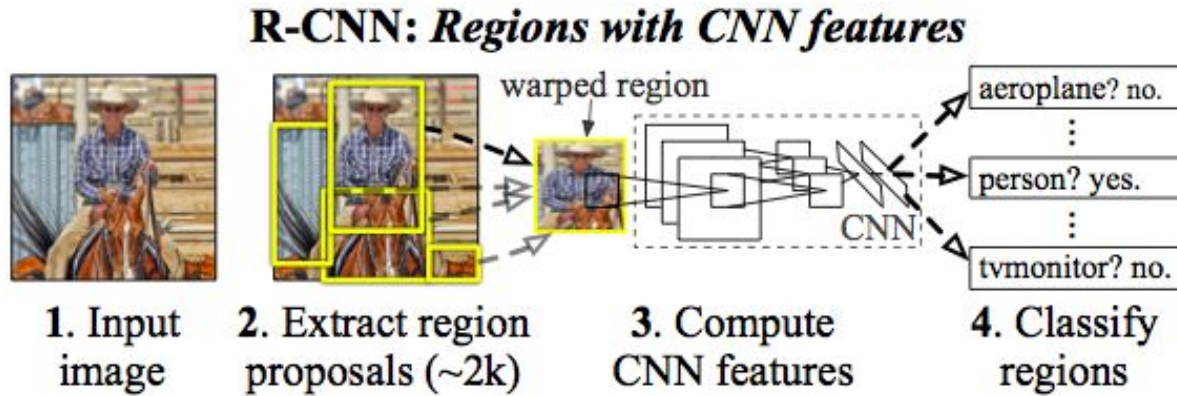
- Input - Facial crop
- Output: x,y coordinates \rightarrow regression loss (e.g. MSE)



Sun et al.: Deep Convolutional Network Cascade for Facial Point Detection, CVPR 2013

Object detection

- R-CNN

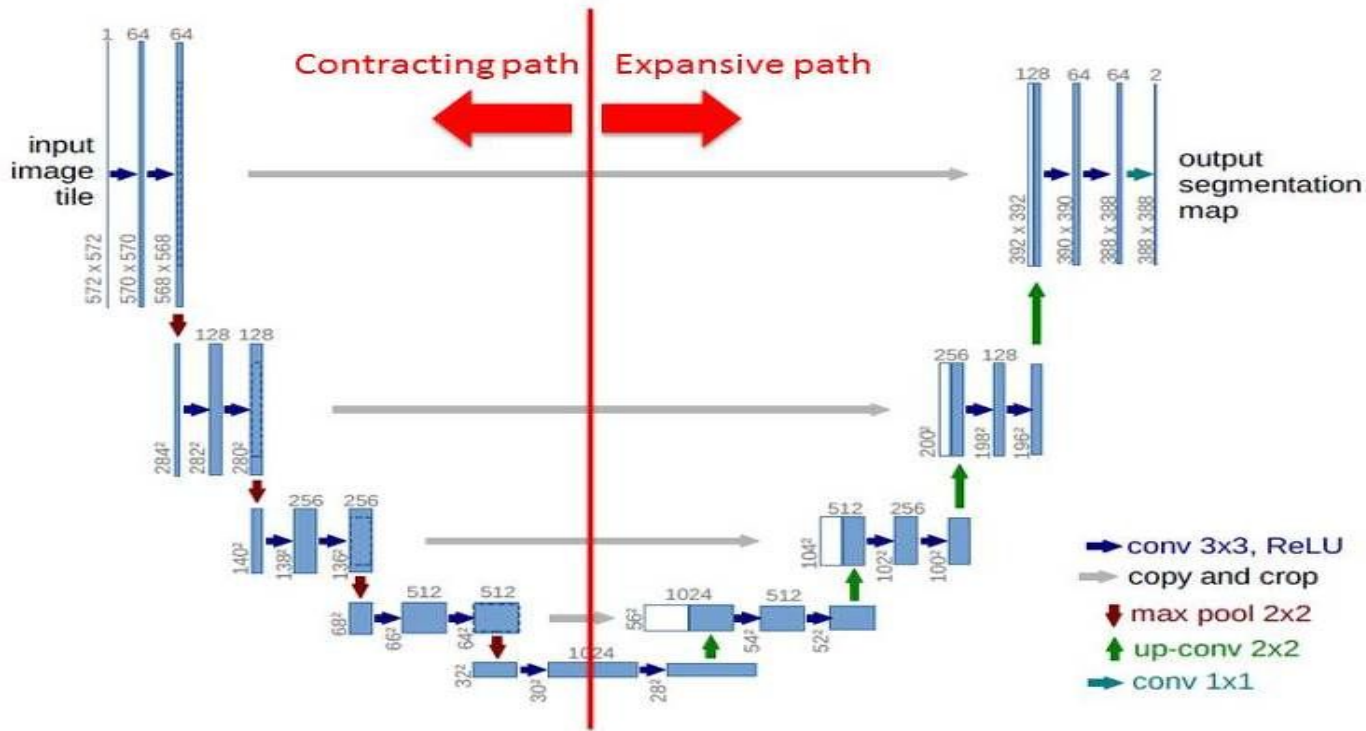


Girshick, R., Donahue, J., Darrell, T., & Malik, J., "Rich feature hierarchies for accurate object detection and semantic segmentation"

Fully Convolutional Networks: FCN

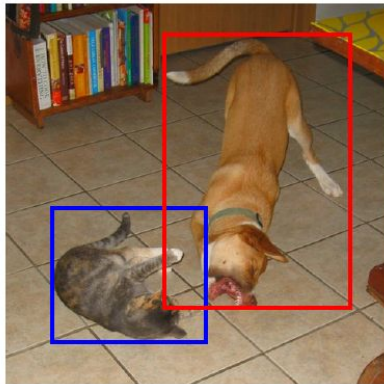


Network Architecture

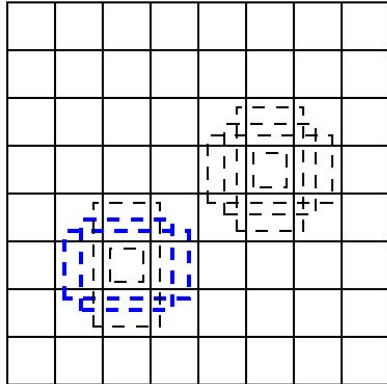


Convolutional object detection

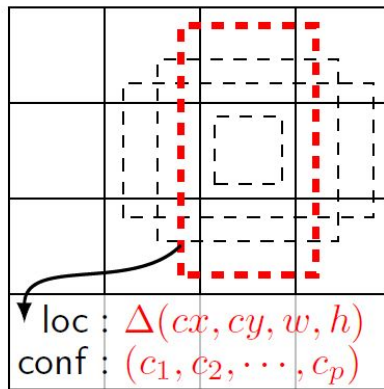
- Output per pixel:
 - For each anchor (base bounding box)
 - Object class (e.g. softmax)



(a) Image with GT boxes



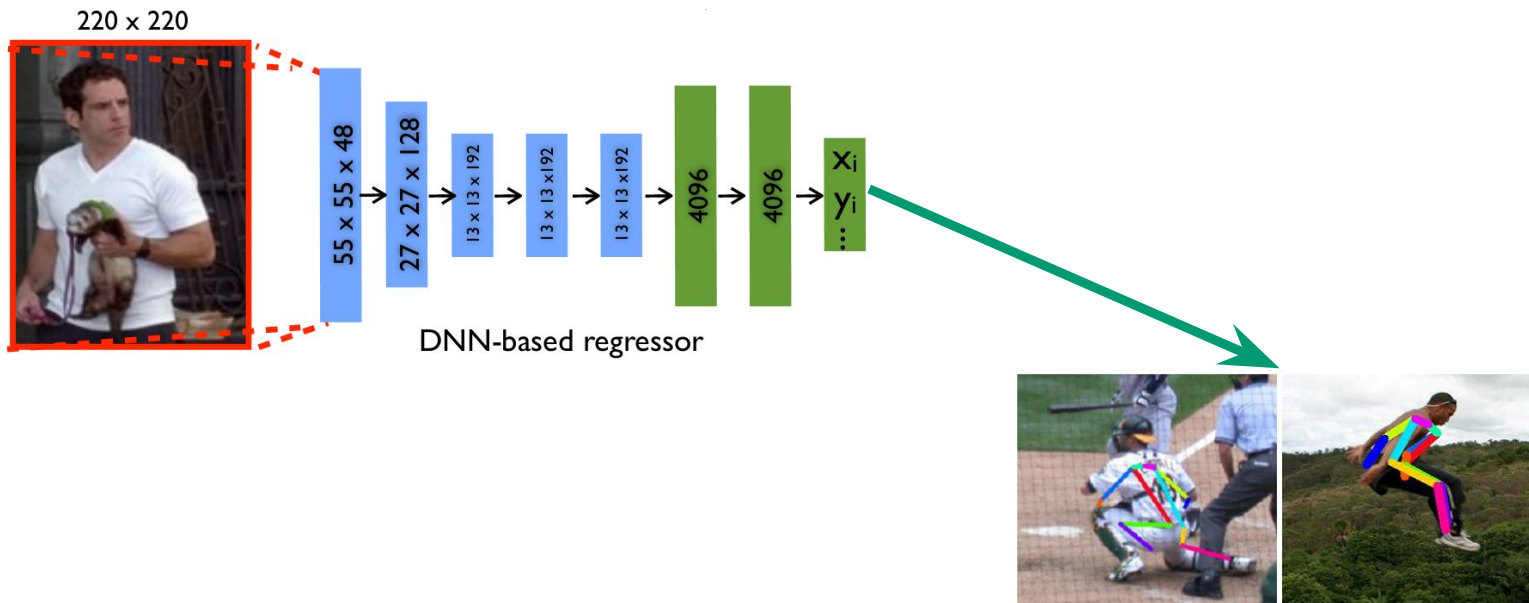
(b) 8×8 feature map



(c) 4×4 feature map

position x, y, w, h)

Pose estimation



Toshev et al.: DeepPose: Human Pose Estimation via Deep Neural Networks. CVPR 2014