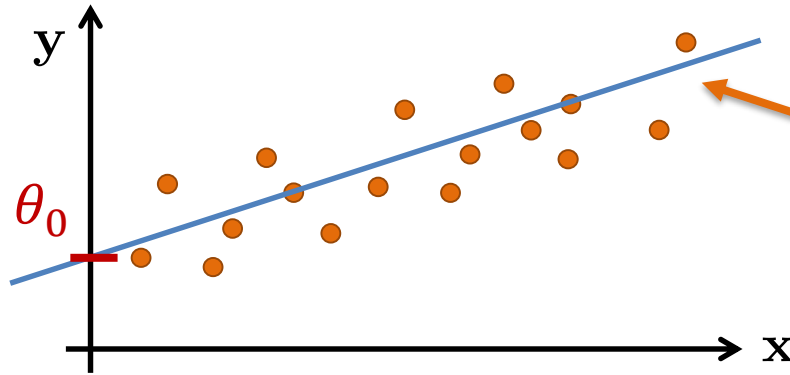# Introduction to Neural Networks

# Lecture 2 Recap

# Linear Regression

= a **supervised learning** method to find a **linear model** of the form

$$\hat{y}_i = \theta_0 + \sum_{j=1}^{d} x_{ij}\theta_j = \theta_0 + x_{i1}\theta_1 + x_{i2}\theta_2 + \cdots + x_{id}\theta_d$$



Goal: find a model that explains a target y given the input x

# Logistic Regression

- Loss function

$$\mathcal{L}(y_i, \widehat{y}_i) = -y_i \cdot \log \widehat{y}_i + (1 - y_i) \cdot \log[1 - \widehat{y}_i])$$
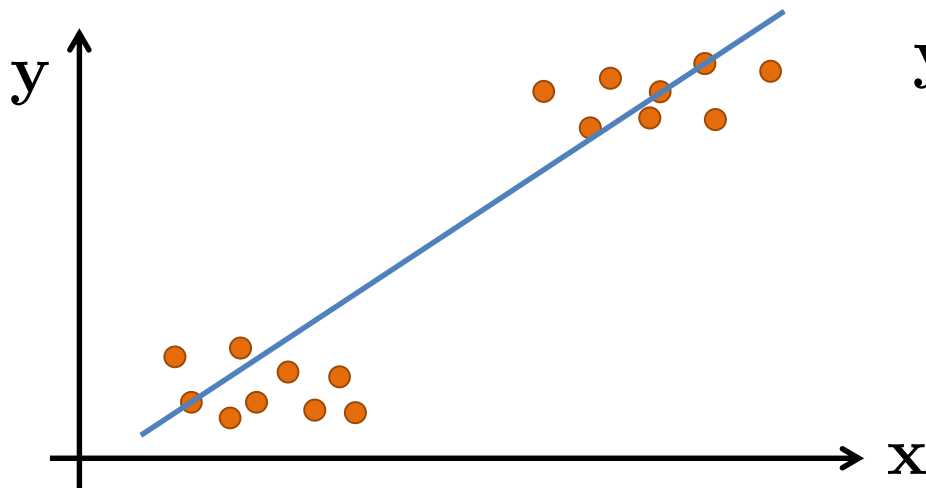
- Cost function

$$\mathcal{C}(\boldsymbol{\theta}) = -\sum_{i=1}^{n}(y_i \cdot \log \widehat{y}_i + (1 - y_i) \cdot \log[1 - \widehat{y}_i])$$

Minimization

$$\widehat{y}_i = \sigma(x_i \boldsymbol{\theta})$$

predictions

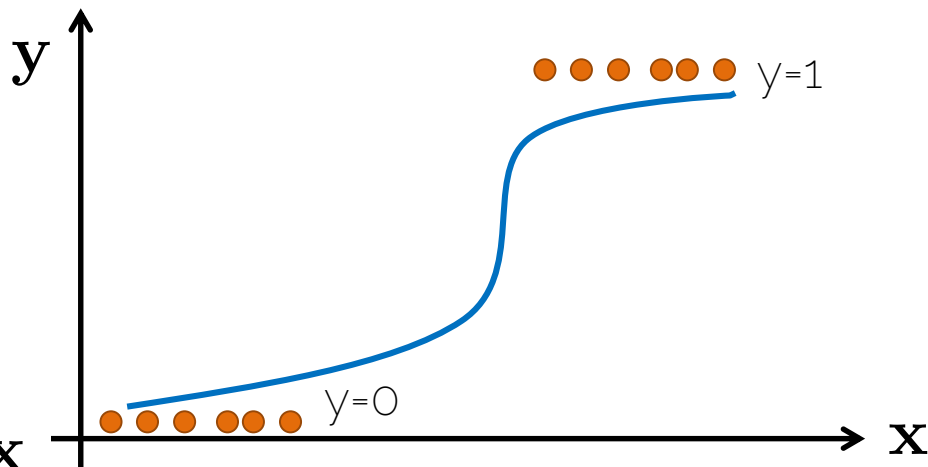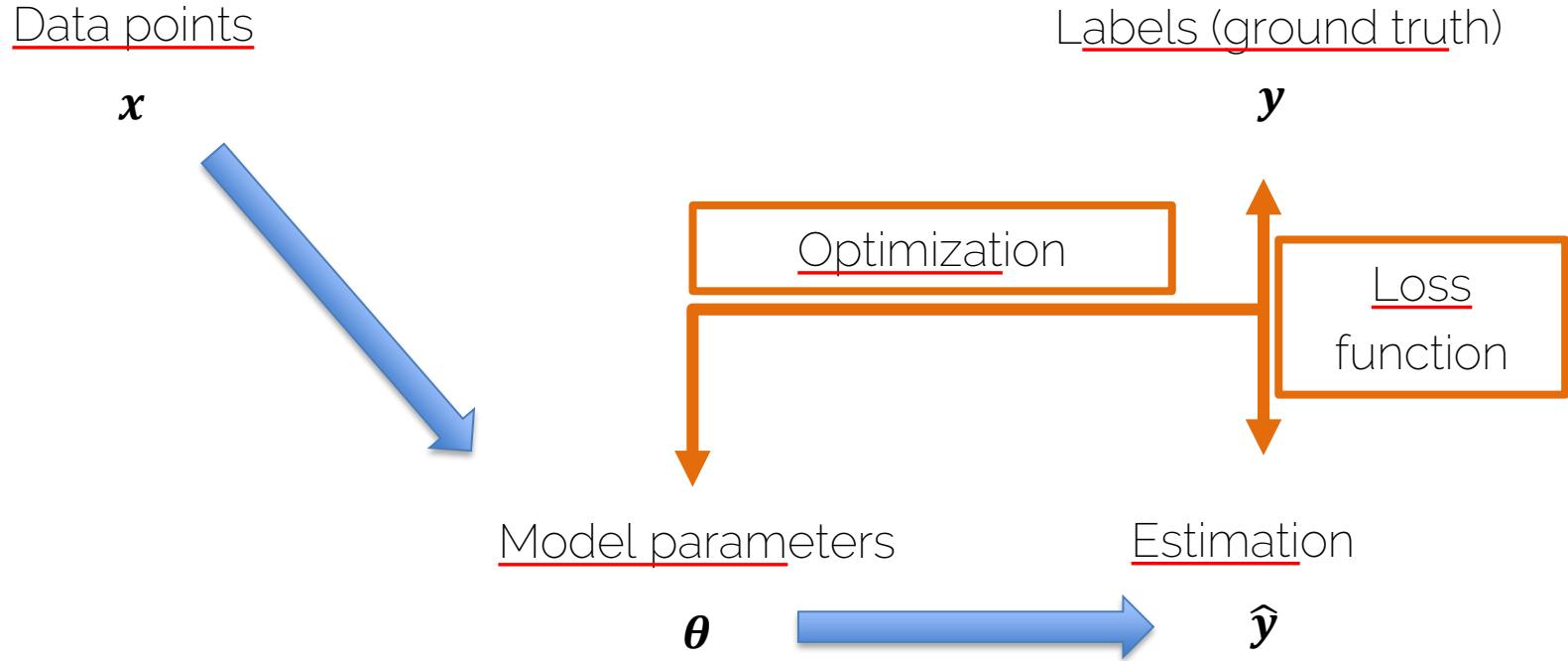Understand

Predictions can exceed the range
of the training samples
  → in the case of classification
[0;1] this becomes a real issue

Predictions are guaranteed
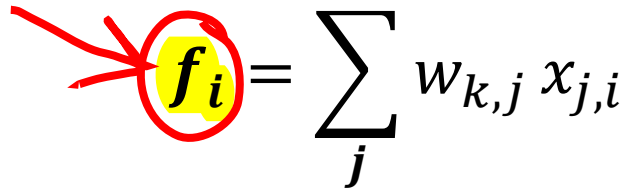to be within [0;1]

# How to obtain the Model?

Data points

$\boldsymbol{x}$

Labels (ground truth)

$\boldsymbol{y}$

Optimization

Loss function

Model parameters

$\boldsymbol{\theta}$

Estimation

$\widehat{\boldsymbol{y}}$

# Linear Score Functions

- Linear score function as seen in linear regression

$$f_i = \sum_j w_{k,j}\, x_{j,i}$$

$$f = W x \qquad \text{(Matrix Notation)}$$
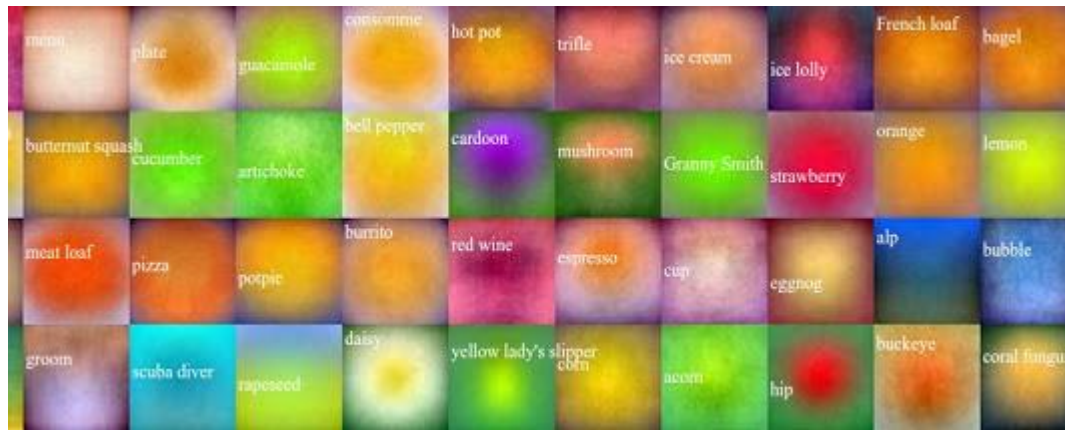
✳ Optimizing in linear score is L2 optim

↳ optimal soln : mean

# Linear Score Functions on Images

- Linear score function $f = Wx$

representation of the average image – so to say.



On CIFAR-10



On ImageNet

object centered in the middle
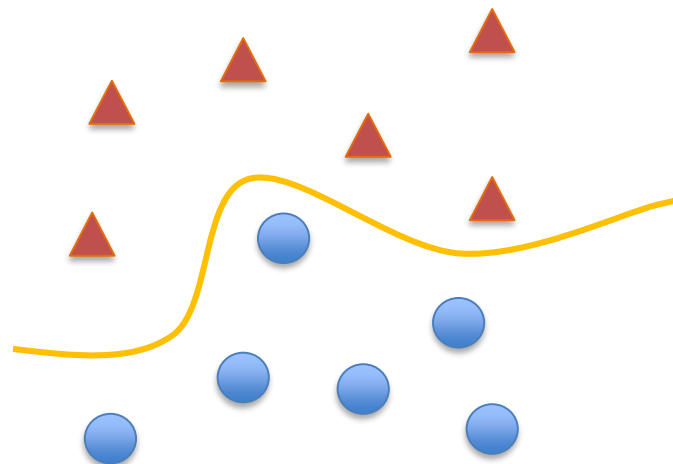
↑

* Simple dataset might be described by linear fn.

Source:: Li/Karpathy/Johnson

# Linear Score Functions?

not powerful enough to rep complex fn

Main Problem

Logistic Regression

Linear Separation Impossible!

# Linear Score Functions?

- Can we make linear regression better?
  - <u>Multiply</u> with another weight matrix $\boldsymbol{W_2}$

$$\hat{\boldsymbol{f}} = \boldsymbol{W_2} \cdot \boldsymbol{f}$$
$$\hat{\boldsymbol{f}} = \boldsymbol{W_2} \cdot \boldsymbol{W} \cdot \boldsymbol{x}$$

- Operation is <u>still linear</u>.

$$\widehat{\boldsymbol{W}} = \boldsymbol{W_2} \cdot \boldsymbol{W}$$
$$\hat{\boldsymbol{f}} = \widehat{\boldsymbol{W}} \, \boldsymbol{x}$$
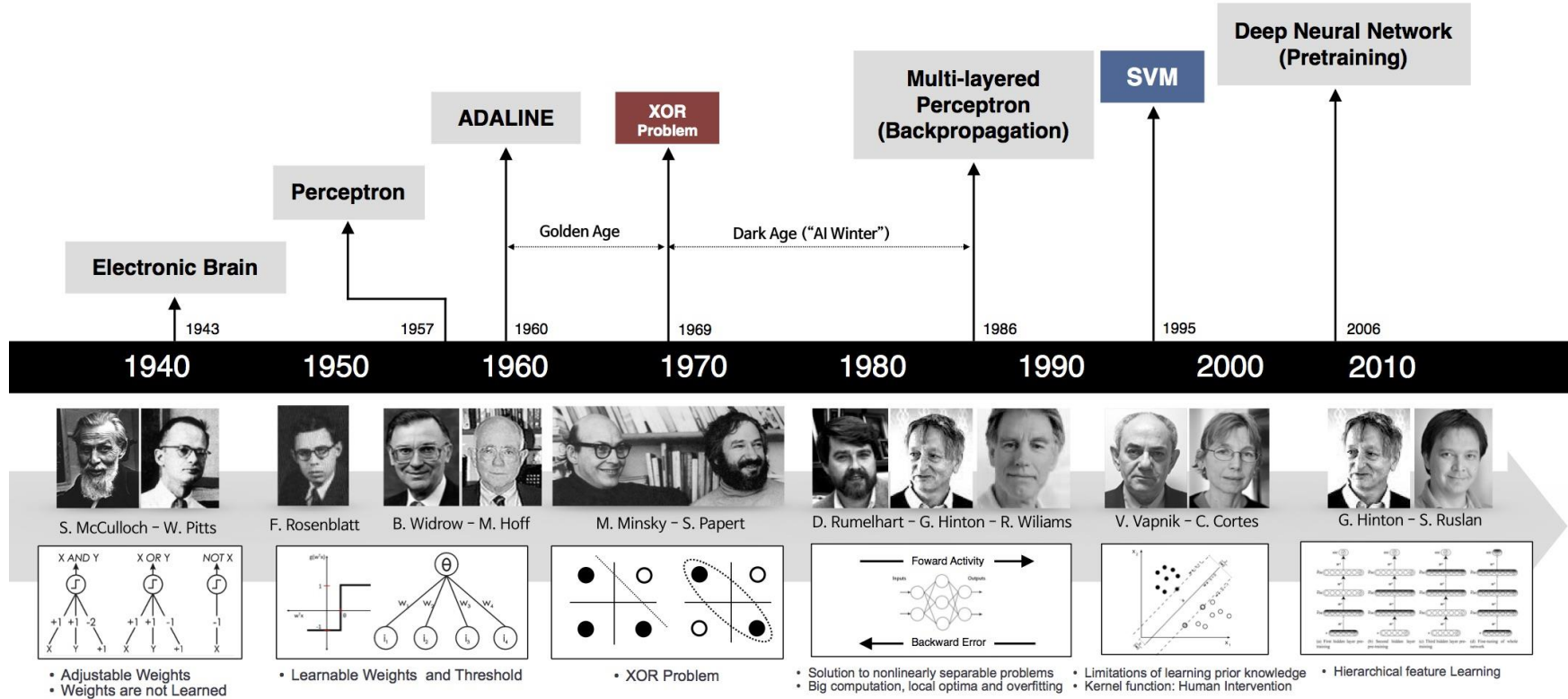
- Solution → <u>add non-linearity</u>!!

# Neural Network

- Linear score function $f = Wx$

- Neural network is a nesting of 'functions'
  - 2-layers: $f = W_2 \max(0, W_1 x)$        *non-linearity*
  - 3-layers: $f = W_3 \max(0, W_2 \max(0, W_1 x))$
  - 4-layers: $f = W_4 \tanh(W_3, \max(0, W_2 \max(0, W_1 x)))$
  - 5-layers: $f = W_5 \sigma(W_4 \tanh(W_3, \max(0, W_2 \max(0, W_1 x))))$
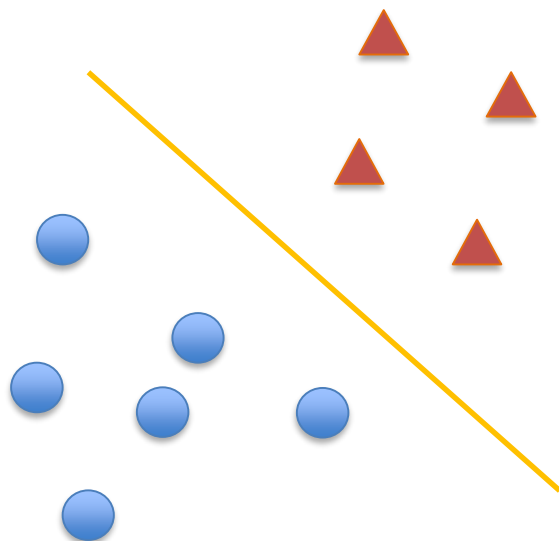  - … up to hundreds of layers

# Introduction to Neural Networks

# History of Neural Networks



Source: http://beamlab.org/deeplearning/2017/02/23/deep_learning_101_part1.html

# Neural Network

Logistic Regression

Neural Networks

# Neural Network

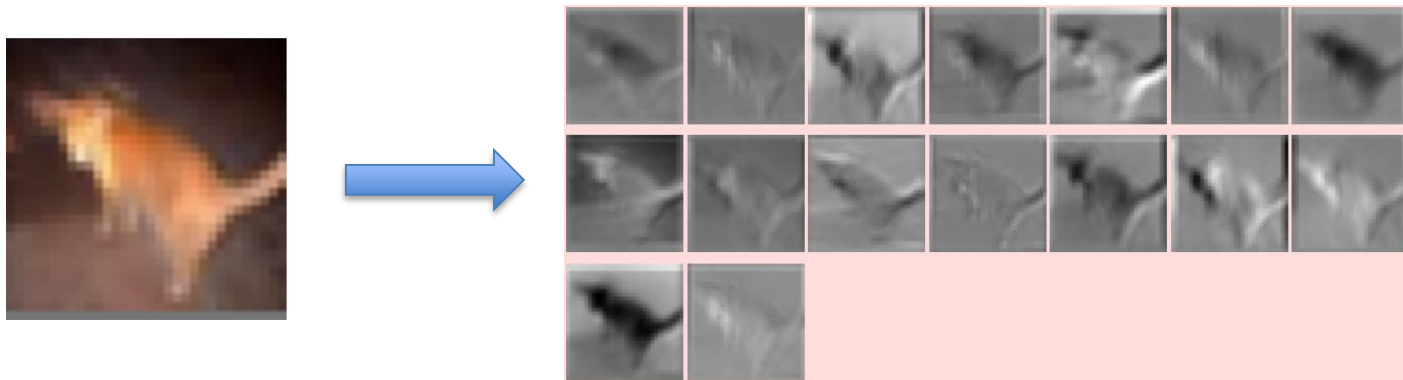- Non-linear score function $f = \ldots (\max(\mathbf{0}, \boldsymbol{W_1 x}))$
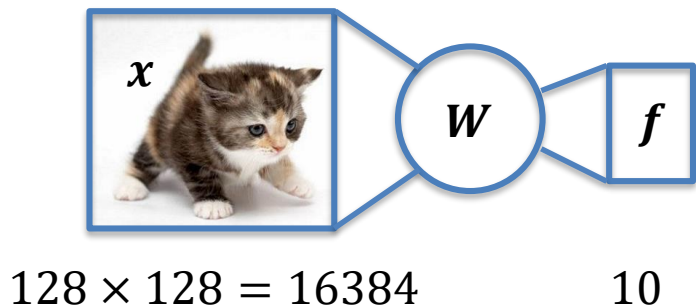
On CIFAR-10

Visualizing activations of first layer.

Source: ConvNetJS

# Neural Network

1-layer network: $f = Wx$

2-layer network: $f = W_2 \max(0, W_1 x)$



$128 \times 128 = 16384$

10

$128 \times 128 = 16384$

1000

10

Why is this structure useful?

# Neural Network

2-layer network: $f = W_2 \max(0, W_1 x)$



$128 \times 128 = 16384$   1000   10

Input Layer   Hidden Layer   Output Layer

input layer

hidden layer

output layer

# Net of Artificial Neurons

dim

$f(W_{0,0}x + b_{0,0})$

$f(W_{0,1}x + b_{0,1})$

$f(W_{0,2}x + b_{0,2})$

$f(W_{0,3}x + b_{0,3})$

$x_1$

$x_2$

$x_3$

$f(W_{1,0}x + b_{1,0})$

$f(W_{1,1}x + b_{1,1})$

$f(W_{1,2}x + b_{1,2})$

$f(W_{2,0}x + b_{2,0})$

$W_{\ell,m,n}$

layer

neuron
in layer

weight
in neuron

NB Do not forget bias
when counting weights

# Neural Network



input layer · hidden layer 1 · hidden layer 2 · hidden layer 3 · output layer

Source: https://towardsdatascience.com/training-deep-neural-networks-9fdb1964b964

*Important*

Sigmoid: $\sigma(x) = \dfrac{1}{(1+e^{-x})}$

Leaky ReLU: $\max(0.1x, x)$

tanh: $\tanh(x)$

Parametric ReLU: $\max(\alpha x, x)$

Maxout $\max(w_1^T x + b_1, w_2^T x + b_2)$

ReLU: $\max(0, x)$

ELU $f(x) = \begin{cases} x & \text{if } x > 0 \\ \alpha(e^x - 1) & \text{if } x \leq 0 \end{cases}$

# Neural Network

$$f = W_3 \cdot (W_2 \cdot (W_1 \cdot x)))$$

*this is still linear*

Why activation functions?

Simply concatenating linear layers would be so much cheaper…

# Neural Network

Why organize a neural network into layers?

# Biological Neurons



impulses carried
toward cell body

branches
of axon

dendrites

nucleus

axon

axon
terminals

cell body

impulses carried
away from cell body

Credit: Stanford CS 231n

# Biological Neurons



totally irrelevant

# Artificial Neural Networks vs Brain



Artificial neural networks are **inspired** by the brain,
but not even close in terms of complexity!
The comparison is great for the media and news articles however… ☺

# Artificial Neural Network

# Neural Network

- Summary → *Supervised*
  - Given a dataset with ground truth training pairs $[x_i; y_i]$,

*idea* →  Find optimal weights $W$ using stochastic gradient descent, such that the loss function is minimized
  - Compute gradients with backpropagation (use batch-mode; more later)
  - Iterate many times over training set (SGD; more later)

# Computational Graphs

# Computational Graphs

- Directional graph

- Matrix operations are represented as compute nodes.

- Vertex nodes are variables or operators like +, -, *, /, log(), exp() ...

- Directional edges show flow of inputs to vertices

# Computational Graphs

- $f(x, y, z) = (x + y) \cdot z$

# Evaluation: Forward Pass

- $f(x, y, z) = (x + y) \cdot z$    Initialization $x = 1, y = -3, z = 4$

# Computational Graphs

- Why discuss compute graphs?

- Neural networks have complicated architectures

$$f = W_5 \sigma(W_4 \tanh(W_3, \max(0, W_2 \max(0, W_1 x))))$$

- Lot of matrix operations!

- Represent NN as computational graphs!

# Computational Graphs

A neural network can be represented as a computational graph...

- it has compute nodes (operations)

- it has edges that connect nodes (data flow)

- it is directional

- it can be organized into 'layers'

# Computational Graphs



$$z_k^{(2)} = \sum_i x_i w_{ik}^{(2)} + b_i^{(2)}$$

$$a_k^{(2)} = f(z_k^{(2)})$$

$$z_k^{(3)} = \sum_i a_i^{(2)} w_{ik}^{(3)} + b_i^{(3)}$$

...

# Computational Graphs

- From a set of neurons to a <u>S</u>tructured <u>C</u>ompute <u>P</u>ipeline



Convolution
AvgPool
MaxPool
Concat
Dropout
Fully connected
Softmax

[Szegedy et al.,CVPR'15] Going Deeper with Convolutions

# Computational Graphs

*Understand*

- The computation of Neural Network has further meanings:
  - The multiplication of $W_i$ and $x$: encode input information
  - The activation function: select the key features



Source; https://www.zybuluo.com/liuhui0803/note/981434

# Computational Graphs

- The computations of Neural Networks have further meanings:
  - The convolutional layers: extract useful features with shared weights



kernel

input

output

# Computational Graphs

- The computations of Neural Networks have further meanings:
  - The convolutional layers: extract useful features with shared weights



Source: https://www.zybuluo.com/liuhui0803/note/981434

# Loss Functions

# What's Next?

Inputs → Neural Network → Outputs ↔ Targets

Are these reasonably close?

**NB**

We need a way to describe how close the network's outputs (= predictions) are to the targets!

# What's Next?

Idea: calculate a '<mark>distance</mark>' between prediction and target!



bad prediction

- prediction
- large distance!
- target

good prediction

- prediction
- small distance!
- target

# Loss Functions

- A function to **measure the goodness** of the predictions (or equivalently, the network's performance)

Intuitively, ...

- – a **large loss** indicates **bad predictions**/performance (→ performance needs to be improved by training the model)

- → the choice of the loss function depends on the concrete problem or the distribution of the target variable

# Regression Loss

- **L1 Loss:**

$$L(\boldsymbol{y}, \widehat{\boldsymbol{y}}; \boldsymbol{\theta}) = \frac{1}{n} \sum_{i}^{n} ||y_i - \widehat{y_i}||_1$$

✳ We normalize
to avoid scaling
issues

- **MSE Loss:**

L2 Loss

$$L(\boldsymbol{y}, \widehat{\boldsymbol{y}}; \boldsymbol{\theta}) = \frac{1}{n} \sum_{i}^{n} ||y_i - \widehat{y_i}||_2^2$$

highly sensitive
to outliers

# Binary Cross Entropy

- Loss function for binary (yes/no) classification

$$L(\boldsymbol{y}, \widehat{\boldsymbol{y}}; \boldsymbol{\theta}) = -\sum_{i=1}^{n}(y_i \cdot \log \widehat{y}_i + (1 - y_i) \cdot \log[1 - \widehat{y}_i])$$

$n \rightarrow$ training samples



Yes! (0.8)

No! (0.2)

The network predicts
the probability of the input
belonging to the "yes" class!

# Cross Entropy

= loss function for multi-class classification

$$L(\boldsymbol{y}, \widehat{\boldsymbol{y}}; \boldsymbol{\theta}) = -\sum_{i=1}^{n} \sum_{k=1}^{k} (y_{ik} \cdot \log \hat{y}_{ik})$$

This <u>generalizes the binary</u> case from the slide before!

dog (0.1)

rabbit (0.2)

duck (0.7)

...

Ⓚ → # Classes

# More General Case

- Ground truth: $y$
- Prediction: $\hat{y}$
- Loss function: $L(y, \hat{y})$
- Motivation:  *it's all about optimization*
  - minimize the loss <=> find better predictions
  - predictions are generated by the NN
  - find better predictions <=> find better NN

# Initially



Prediction

Targets

Bad prediction!

Loss

$t_1$

Training time

# During Training...



Prediction

Targets

Bad prediction!

Loss

$t_2$

Training time

# During Training...



Prediction

Targets

Bad prediction!

Loss

t₃

Training time

# Training Curve

Loss

→ almost the only way to asses a model



Training time

_Understand_

# How to Find a Better NN?

Plotting loss curves against model parameters



Loss
$L(\boldsymbol{y}, f_{\boldsymbol{\theta}}(\boldsymbol{x}))$

$\boldsymbol{\theta}$

Parameters $\boldsymbol{\theta}$

# How to Find a Better NN?

- Loss function: $L(\boldsymbol{y}, \widehat{\boldsymbol{y}}) = L(\boldsymbol{y}, f_{\boldsymbol{\theta}}(\boldsymbol{x}))$
- Neural Network: $f_{\boldsymbol{\theta}}(\boldsymbol{x})$
- Goal:
  - minimize the loss w. r. t. $\boldsymbol{\theta}$

➡️ Optimization! We train compute graphs with some optimization techniques!

└ for this lec, only gradient_based

# How to Find a Better NN?

- Minimize: $L\big(\boldsymbol{y}, f_{\boldsymbol{\theta}}(\boldsymbol{x})\big)$ w.r.t. $\boldsymbol{\theta}$

- In the context of NN, we use gradient-based optimization

Loss
$L(\boldsymbol{y}, f_{\boldsymbol{\theta}}(\boldsymbol{x}))$

$\boldsymbol{\theta}$

Loss function

Gradient Descent

$\boldsymbol{\theta}$

# How to Find a Better NN?

- Minimize: $L\big(\boldsymbol{y}, f_{\boldsymbol{\theta}}(\boldsymbol{x})\big)$ w.r.t. $\boldsymbol{\theta}$



$\boldsymbol{\nabla}_{\boldsymbol{\theta}} L(\boldsymbol{y}, f_{\boldsymbol{\theta}}(\boldsymbol{x}))$

$\boldsymbol{\theta}$

Loss
$L(\boldsymbol{y}, f_{\boldsymbol{\theta}}(\boldsymbol{x}))$

$\boldsymbol{\theta}$

Learning rate

$\boldsymbol{\theta} = \boldsymbol{\theta} - \alpha\, \boldsymbol{\nabla}_{\boldsymbol{\theta}} L(\boldsymbol{y}, f_{\boldsymbol{\theta}}(\boldsymbol{x}))$

$\boldsymbol{\theta}^* = \arg\min L\big(\boldsymbol{y}, f_{\boldsymbol{\theta}}(\boldsymbol{x})\big)$

# How to Find a Better NN?

- Given inputs $\boldsymbol{x}$ and targets $\boldsymbol{y}$

- Given one layer NN with no activation function

$$f_{\boldsymbol{\theta}}(\boldsymbol{x}) = \boldsymbol{W}\boldsymbol{x}, \qquad \boldsymbol{\theta} = \boldsymbol{W}$$

Later $\boldsymbol{\theta} = \{\boldsymbol{W}, \boldsymbol{b}\}$

- Given MSE Loss: $L(\boldsymbol{y}, \widehat{\boldsymbol{y}}; \boldsymbol{\theta}) = \frac{1}{n}\sum_{i}^{n} ||y_i - \widehat{y}_i||_2^2$

# How to Find a Better NN?

- Given inputs $\boldsymbol{x}$ and targets $\boldsymbol{y}$

- Given one layer NN with no activation function

- Given MSE Loss: $L(\boldsymbol{y}, \widehat{\boldsymbol{y}}; \boldsymbol{\theta}) = \frac{1}{n} \sum_{i}^{n} ||y_i - \boldsymbol{W} \cdot x_i||_2^2$



Gradient flow

Backprop

Multiply

# How to Find a Better NN?

- Given inputs $\boldsymbol{x}$ and targets $\boldsymbol{y}$

- Given a one layer NN with no activation function
$$f_{\boldsymbol{\theta}}(\boldsymbol{x}) = \boldsymbol{W}\boldsymbol{x}\,, \qquad \boldsymbol{\theta} = \boldsymbol{W}$$

- Given MSE Loss: $L(\boldsymbol{y}, \widehat{\boldsymbol{y}}; \boldsymbol{\theta}) = \frac{1}{n}\sum_{i}^{n}||\boldsymbol{W}\cdot\boldsymbol{x_i} - \boldsymbol{y_i}||_2^2$

- $\nabla_{\boldsymbol{\theta}}L\big(\boldsymbol{y}, f_{\boldsymbol{\theta}}(\boldsymbol{x})\big) = \frac{1}{n}\sum_{i}^{n}(\boldsymbol{W}\cdot\boldsymbol{x_i} - \boldsymbol{y_i})\cdot\boldsymbol{x_i^T}$

# How to Find a Better NN?

- Given inputs $\boldsymbol{x}$ and targets $\boldsymbol{y}$
- Given a multi-layer NN with many activations

$$\boldsymbol{f} = \boldsymbol{W_5}\sigma(\boldsymbol{W_4}\tanh(\boldsymbol{W_3}, \max(\boldsymbol{0}, \boldsymbol{W_2}\max(\boldsymbol{0}, \boldsymbol{W_1}\boldsymbol{x}))))$$

- Gradient descent for $L\big(\boldsymbol{y}, f_{\boldsymbol{\theta}}(\boldsymbol{x})\big)$ w. r. t. $\boldsymbol{\theta}$
  - Need to propagate gradients from end to first layer ($\boldsymbol{W_1}$).

# How to Find a Better NN?

- Given inputs $x$ and targets $y$
- Given multi-layer NN with many activations

# How to Find a Better NN?

- Given inputs $\boldsymbol{x}$ and targets $\boldsymbol{y}$

- Given multilayer layer NN with many activations

$$f = W_5 \sigma(W_4 \tanh(W_3, \max(0, W_2 \max(0, W_1 x))))$$

- Gradient descent solution for $L\big(\boldsymbol{y}, f_{\boldsymbol{\theta}}(\boldsymbol{x})\big)$ w. r. t. $\boldsymbol{\theta}$

  - Need to propagate gradients from end to first layer ($W_1$)

  Backpropagation: Use chain rule to compute gradients

  - Compute graphs come in handy!

# How to Find a Better NN?

- <mark>Why gradient descent?</mark>
  - <u>Easy to compute</u> using compute graphs
- Other methods include
  - Newtons method
  - L-BFGS
  - Adaptive moments
  - Conjugate gradient

_— efficient with large networks_

# Summary

- Neural Networks are computational graphs
- Goal: for a given train set, find optimal weights

- Optimization is done using gradient-based solvers
  - Many options (more in the next lectures)

- Gradients are computed via backpropagation
  - Nice because can easily modularize complex functions

# Next Lectures

- Next Lecture:
  - Backpropagation and optimization of Neural Networks

- **Check for updates** on website/moodle regarding exercises

# See you next week ☺

# Further Reading

- Optimization:
  - http://cs231n.github.io/optimization-1/
  - http://www.deeplearningbook.org/contents/optimization.html


- General concepts:
  - Pattern Recognition and Machine Learning – C. Bishop
  - http://www.deeplearningbook.org/