

CY 7790

Special Topics in Security and Privacy:
Machine Learning Security and
Privacy
Fall 2021

Alina Oprea
Associate Professor
Khoury College of Computer Science

September 20 2021

Logistical Issues

- HW1 will be released today
 - It will be due on Friday, October 1
 - You will receive a Gradescope invitation
- Will release the presentation schedule later today
- Several resources on Piazza on reading papers
 - Paper template

Outline: Review of ML and Deep Learning

- Ensemble learning
 - Bagging, boosting
- Feed-forward neural networks
- Convolutional networks
 - Common architectures
- Regularization
- Backpropagation
- Comparing classifiers

Ensemble Learning

Consider a set of classifiers h_1, \dots, h_L


Idea: construct a classifier $H(\mathbf{x})$ that combines the individual decisions of h_1, \dots, h_L

- e.g., could have the member classifiers vote, or
- e.g., could use different members for different regions of the instance space

Successful ensembles require **diversity**

- Classifiers should make different mistakes
- Can have different types of base learners

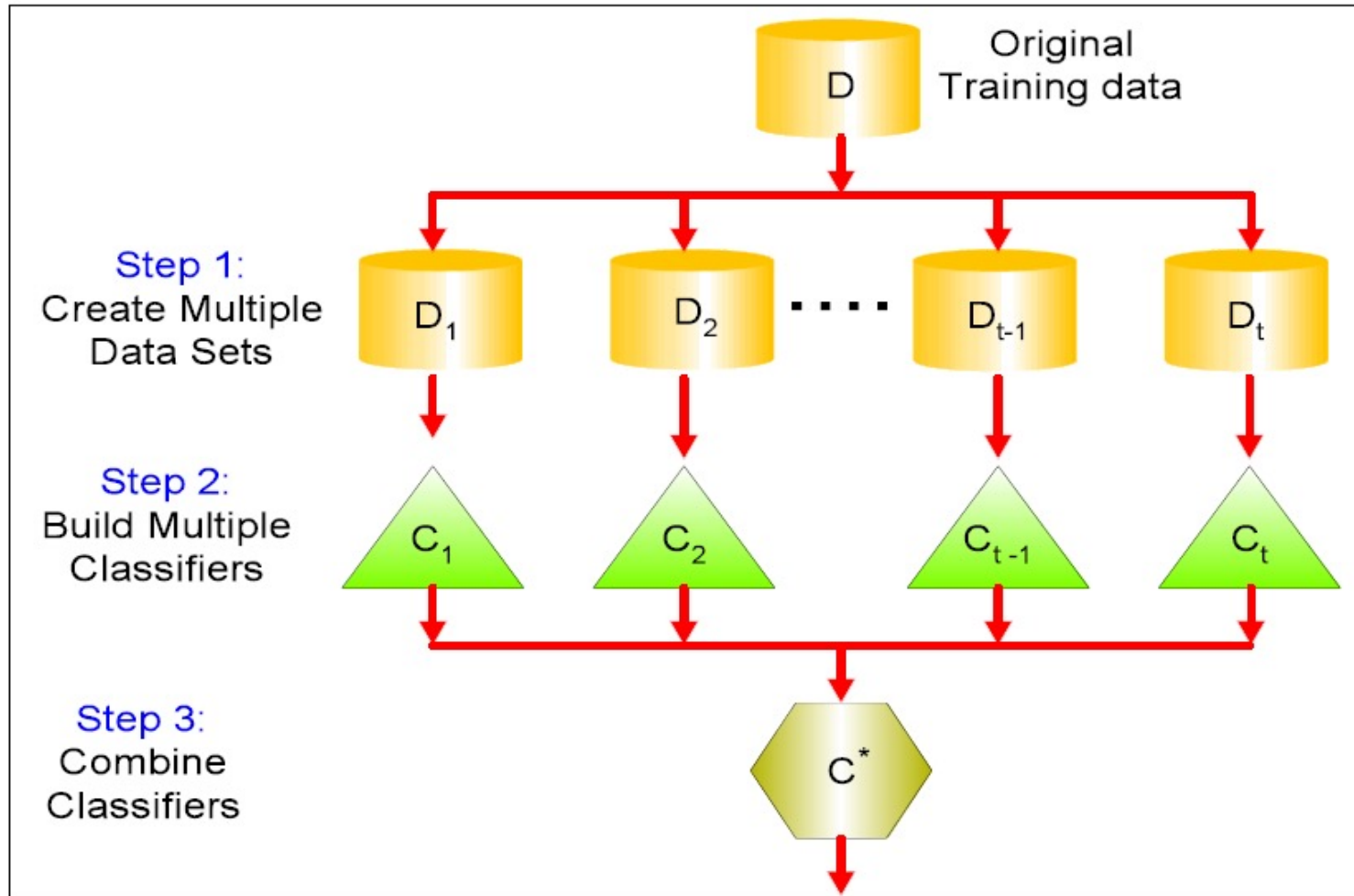
How to Achieve Diversity

- 
- Avoid overfitting
 - Vary the training data
 - Features are noisy
 - Vary the set of features

Two main ensemble learning methods

- Bagging PARALLEL RANDOM FOREST
- Boosting SEQUENTIAL Ada Boost
Gradient Boosting

Bagging



Majority Votes / Avg. for regression

Random Forest Algorithm

1. For $b = 1$ to B :

- (a) Draw a **bootstrap sample** \mathbf{Z}^* of size N from the training data.
- (b) Grow a random-forest tree T_b to the bootstrapped data, by recursively repeating the following steps for each terminal node of the tree, until the minimum node size n_{min} is reached.
 - i. Select **m variables at random** from the p variables.
 - ii. Pick the best variable/split-point among the m .
 - iii. Split the node into two daughter nodes.

$$m \approx \sqrt{d}$$

2. Output the ensemble of trees $\{T_b\}_1^B$.

To make a prediction at a new point x :

Regression: $\hat{f}_{\text{rf}}^B(x) = \frac{1}{B} \sum_{b=1}^B T_b(x).$

Classification: Let $\hat{C}_b(x)$ be the class prediction of the b th random-forest tree. Then $\hat{C}_{\text{rf}}^B(x) = \text{majority vote } \{\hat{C}_b(x)\}_1^B$.

Overview of AdaBoost

$$Loss = \frac{1}{N} \sum_{i=1}^N L(x_i, y_i)$$

$$Loss_W = \frac{1}{N} \sum_{i=1}^N w_i L(x_i, y_i)$$

$$G(x) = \text{sign}\left(\sum_{i=1}^T \beta_i h_i(x)\right)$$

Better classifiers will get higher weights

- Mis-classified examples get higher weights
- Correct examples get lower weights

Uniform weights

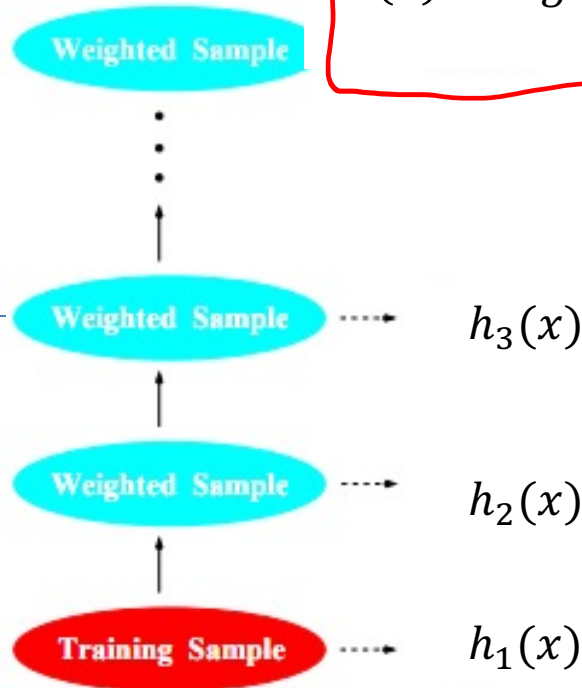


FIGURE 10.1. Schematic of AdaBoost. Classifiers are trained on weighted versions of the dataset, and then combined to produce a final prediction.

Bagging vs Boosting

Bagging

vs.

Boosting

Resamples data points

Weight of each classifier is the same

Only variance reduction

Applicable to complex models with low bias, high variance

E.g. DECISION TREES

Reweights data points (modifies their distribution)

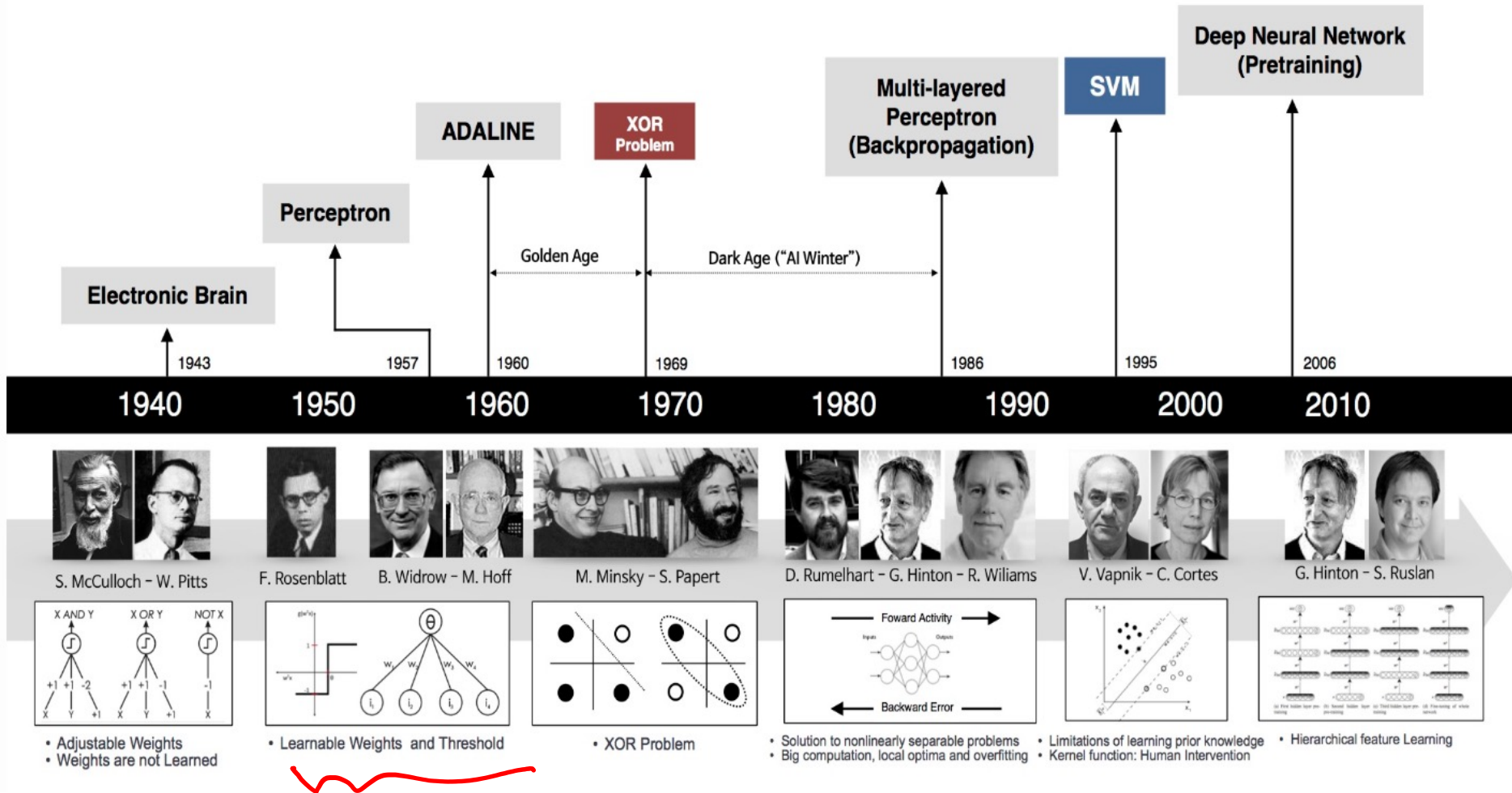
Weight is dependent on classifier's accuracy

Both bias and variance reduced – learning rule becomes more complex with iterations

Applicable to weak models with high bias, low variance

LINEAR MODELS / DECISION STUMP

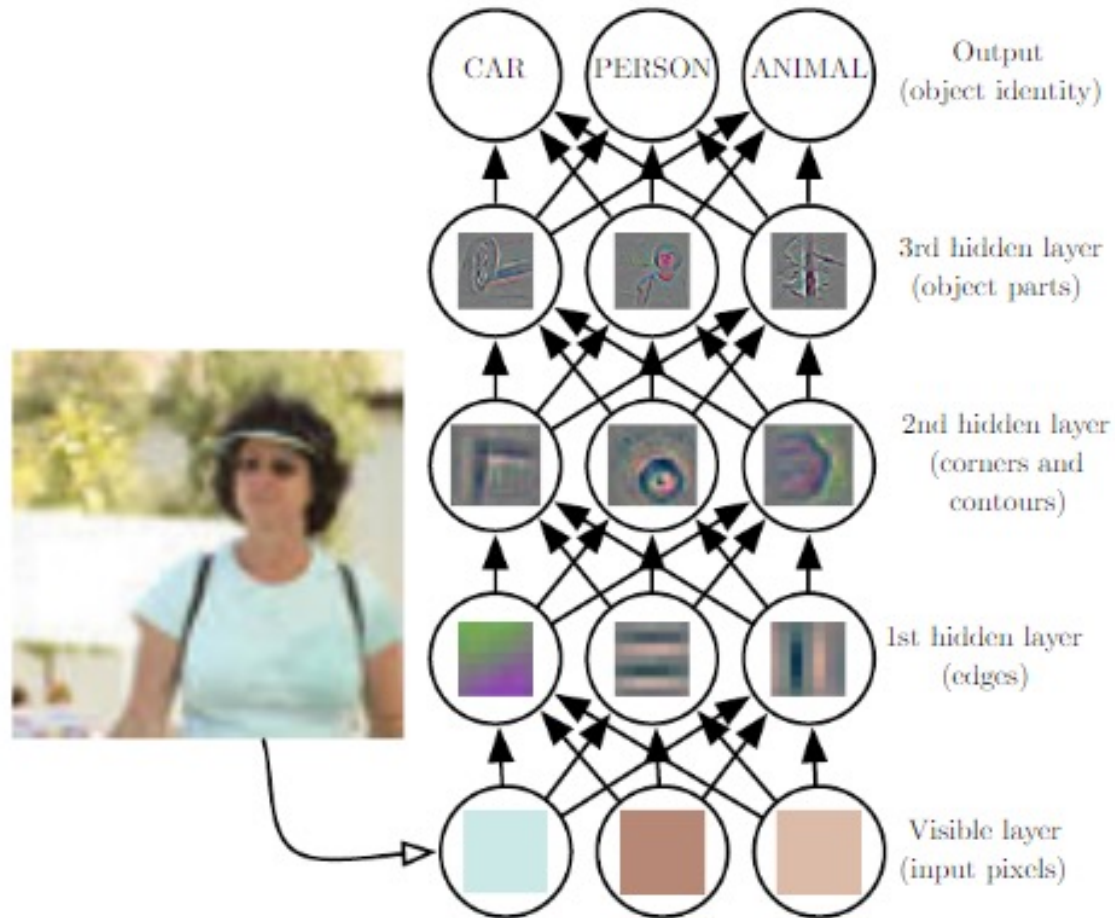
Timeline



Deep Learning References

- Chapter 10 from Introduction to Statistical Learning
 - <https://www.statlearning.com/>
- Deep Learning books
 - <https://d2l.ai/> (D2L)
 - <https://www.deeplearningbook.org/> (advanced)
- Stanford notes on deep learning
 - http://cs229.stanford.edu/summer2020/cs229-notes-deep_learning.pdf
- Stanford tutorial on training Multi-Layer Neural Networks
 - <http://ufldl.stanford.edu/tutorial/supervised/MultiLayerNeuralNetworks/>
- Notes on backpropagation by Andrew Ng
 - <http://cs229.stanford.edu/notes-spring2019/backprop.pdf>

Learning Representations



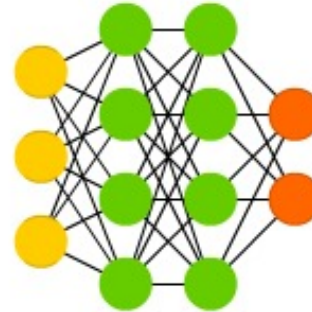
Deep Learning addresses the problem of learning hierarchical representations

Neural Network Architectures

Feed-Forward Networks

- Neurons from each layer connect to neurons from next layer

Deep Feed Forward (DFF)

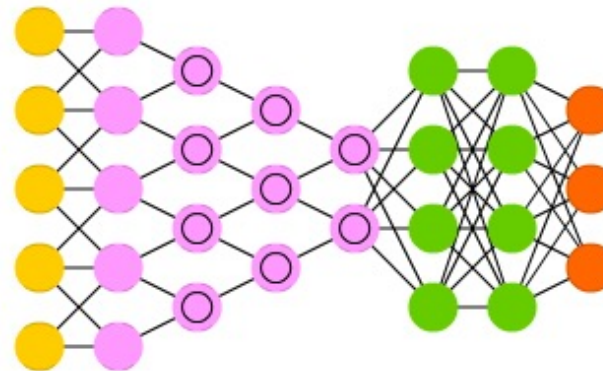


MLP

Convolutional Networks

- Includes convolution layer for feature reduction
- Learns hierarchical representations

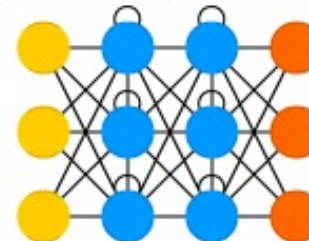
Deep Convolutional Network (DCN)



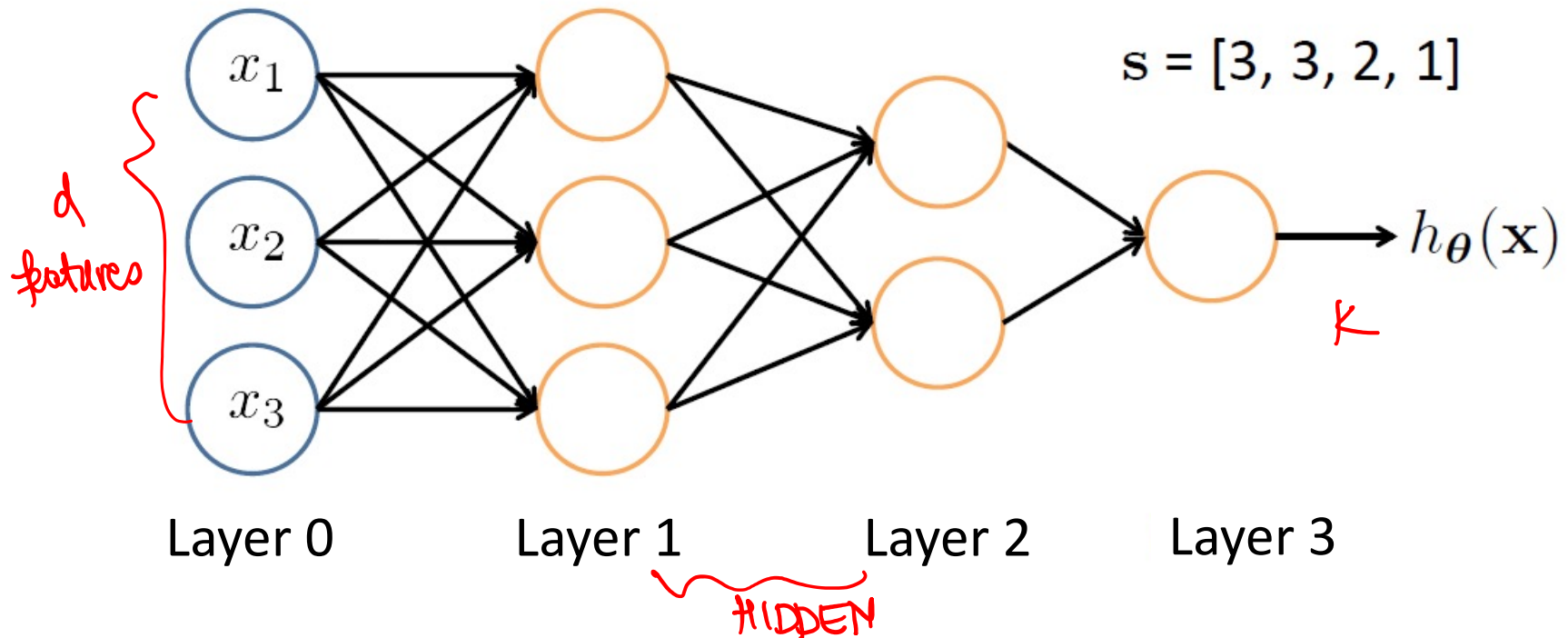
Recurrent Networks

- Keep hidden state
- Have cycles in computational graph

Recurrent Neural Network (RNN)



Feed-Forward Networks

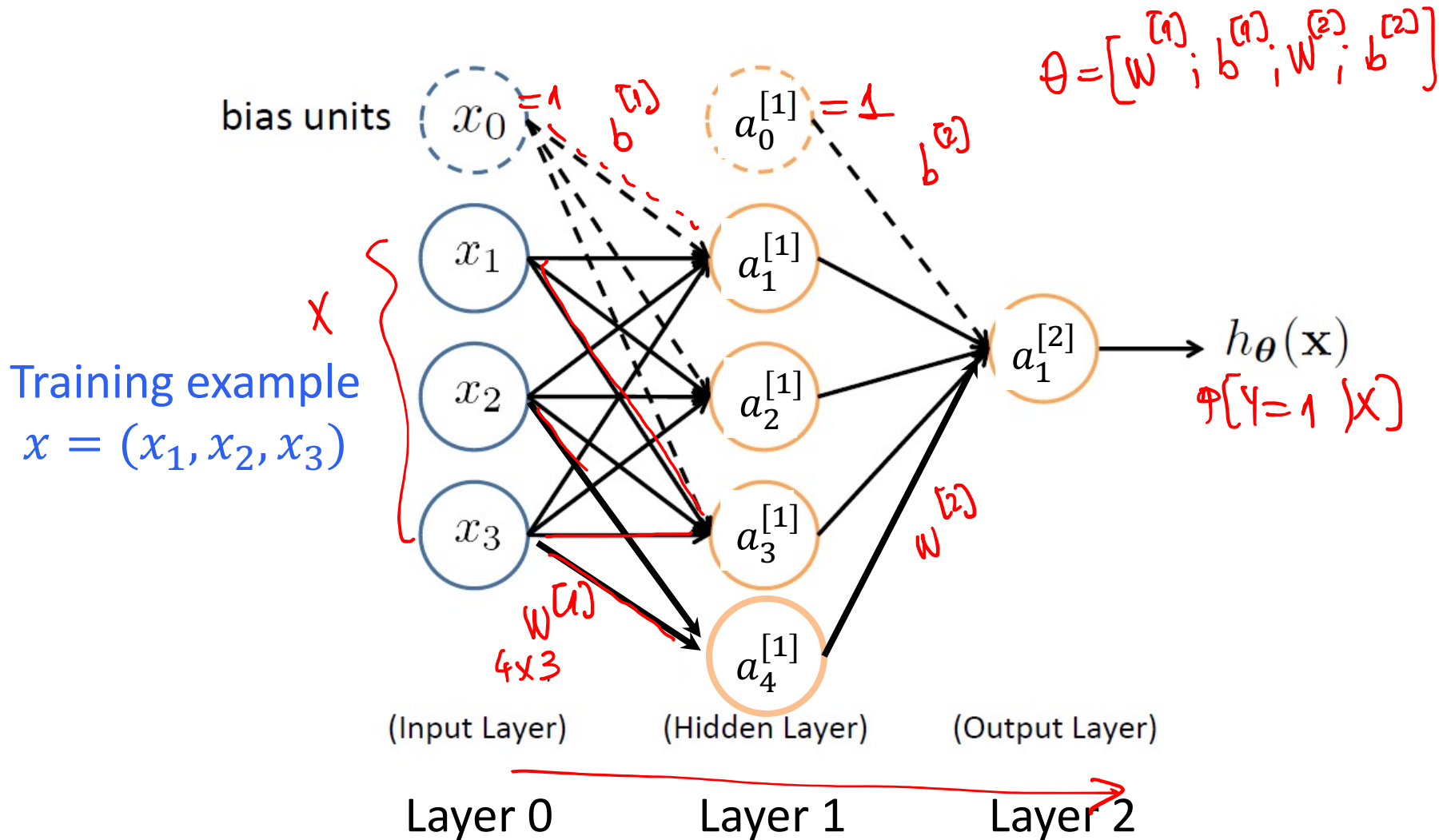


L denotes the number of layers

$s \in \mathbb{N}^{+L}$ contains the numbers of nodes at each layer

- Not counting bias units
- Typically, $s_0 = d$ (# input features) and $s_{L-1} = K$ (# classes)

Feed-Forward Neural Network



Layer Operations

$$z_1^{[1]} = W_1^{[1]} x + b_1^{[1]} \quad \text{and} \quad a_1^{[1]} = g(z_1^{[1]})$$

$$\vdots$$
$$\vdots$$
$$\vdots$$

$$z_4^{[1]} = W_4^{[1]} x + b_4^{[1]} \quad \text{and} \quad a_4^{[1]} = g(z_4^{[1]})$$

$$\underbrace{\begin{bmatrix} z_1^{[1]} \\ \vdots \\ \vdots \\ z_4^{[1]} \end{bmatrix}}_{z^{[1]} \in \mathbb{R}^{4 \times 1}} = \underbrace{\begin{bmatrix} - & W_1^{[1]} & - \\ - & W_2^{[1]} & - \\ & \vdots & \\ - & W_4^{[1]} & - \end{bmatrix}}_{W^{[1]} \in \mathbb{R}^{4 \times 3}} \underbrace{\begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}}_{x \in \mathbb{R}^{3 \times 1}} + \underbrace{\begin{bmatrix} b_1^{[1]} \\ b_2^{[1]} \\ \vdots \\ b_4^{[1]} \end{bmatrix}}_{b^{[1]} \in \mathbb{R}^{4 \times 1}}$$

$$z^{[1]} = W^{[1]}x + b^{[1]}$$

Linear

ACTIVATION

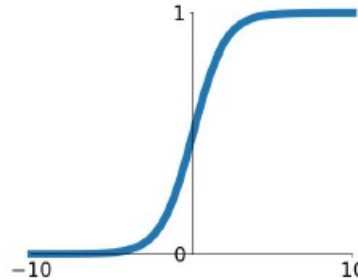
$$a^{[1]} = g(z^{[1]})$$

Non-Linear

Activation Functions

Sigmoid

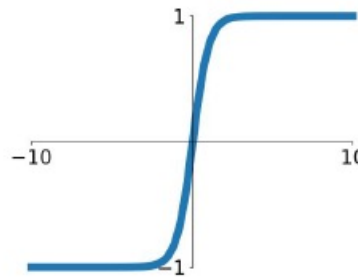
$$\sigma(x) = \frac{1}{1+e^{-x}}$$



Binary
Classification

tanh

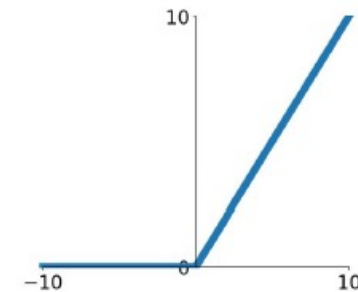
$$\tanh(x)$$



Regression

ReLU

$$\max(0, x)$$

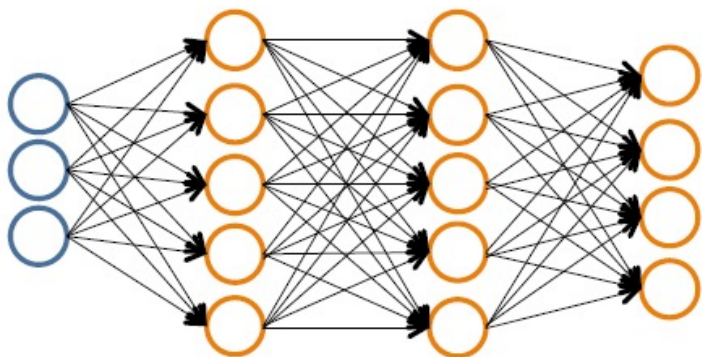


Intermediary
layers

SOFTMAX



Neural Network Classification



Given:

$$\{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_n, y_n)\}$$

$\mathbf{s} \in \mathbb{N}^{+L}$ contains # nodes at each layer

– $s_0 = d$ (# features)

Binary classification

$y = 0$ or 1

1 output unit ($s_{L-1} = 1$)

Sigmoid

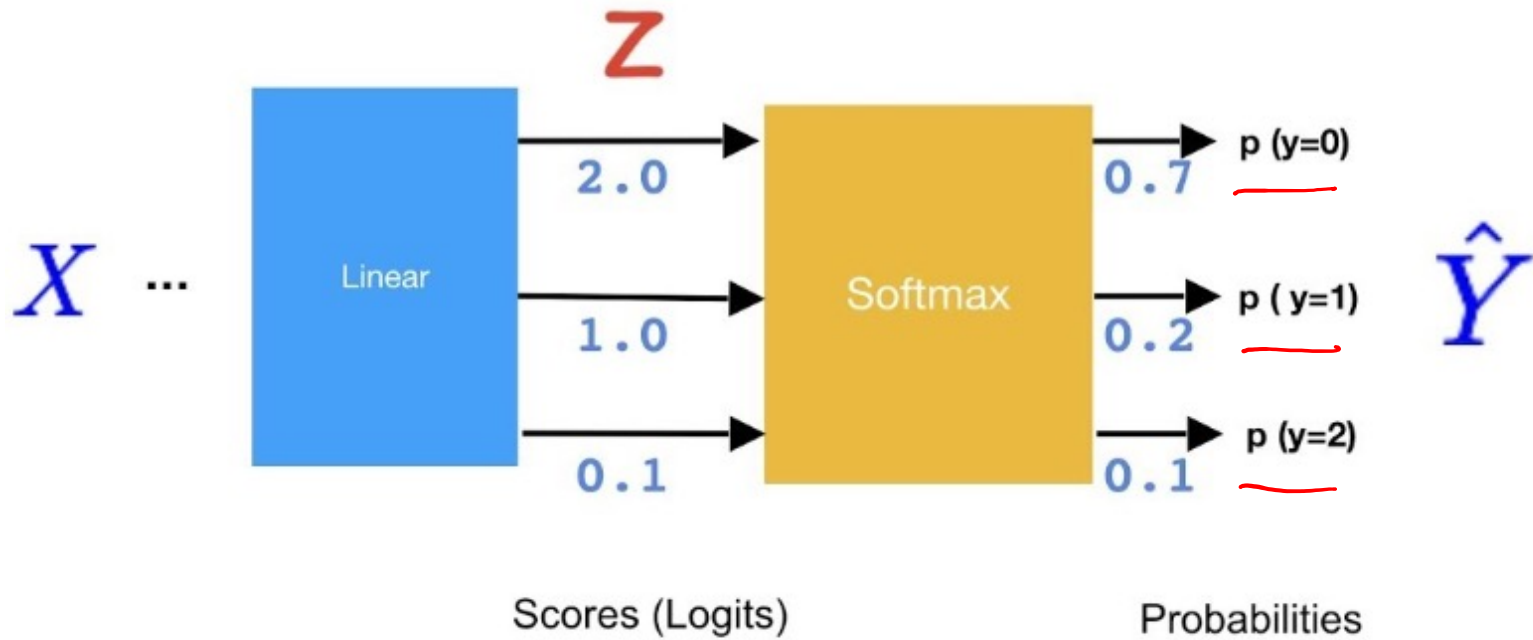
Multi-class classification (K classes)

$\mathbf{y} \in \mathbb{R}^K$ e.g. $\begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}$, $\begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix}$, $\begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \end{bmatrix}$, $\begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \end{bmatrix}$
pedestrian car motorcycle truck

K output units ($s_{L-1} = K$)

Softmax

Softmax classifier



$$Z = [W]X + [b] \quad \text{LINEAR}$$

$$y_j = \frac{e^{z_j}}{\sum_{k=1}^K e^{z_k}} \in [0, 1]$$

$$\sum_{j=1}^K y_j = 1$$

$$\frac{z_j}{\sum_{k=1}^K z_k}$$

Convolutional Nets

- Particular type of Feed-Forward Neural Nets
 - Invented by [LeCun 89]
- Applicable to data with natural grid topology
 - Time series
 - Images
- Use convolutions on at least one layer
 - Convolution is a linear operation that uses local information
 - Also use pooling operation
 - Used for dimensionality reduction and learning hierarchical feature representations

Convolutional Nets

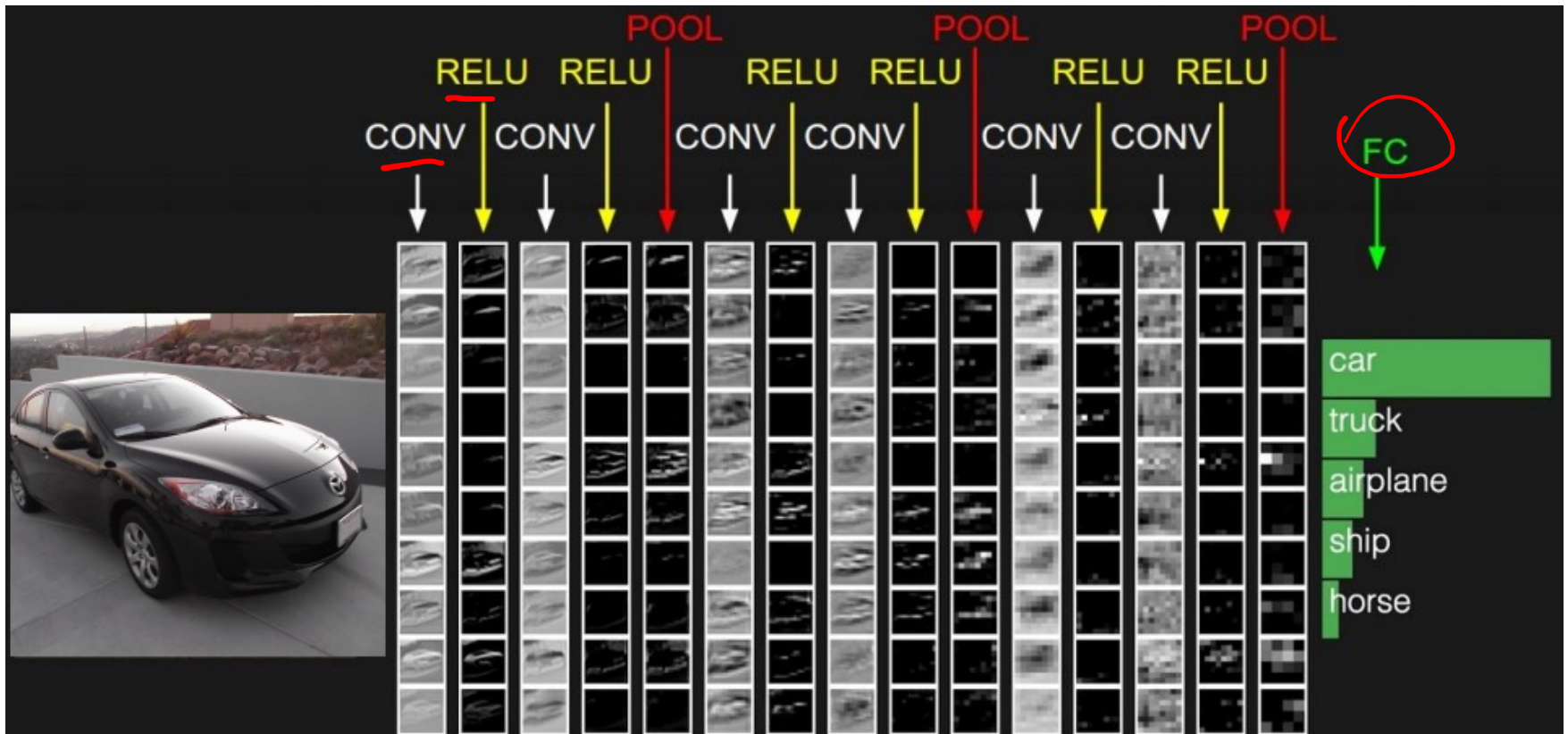
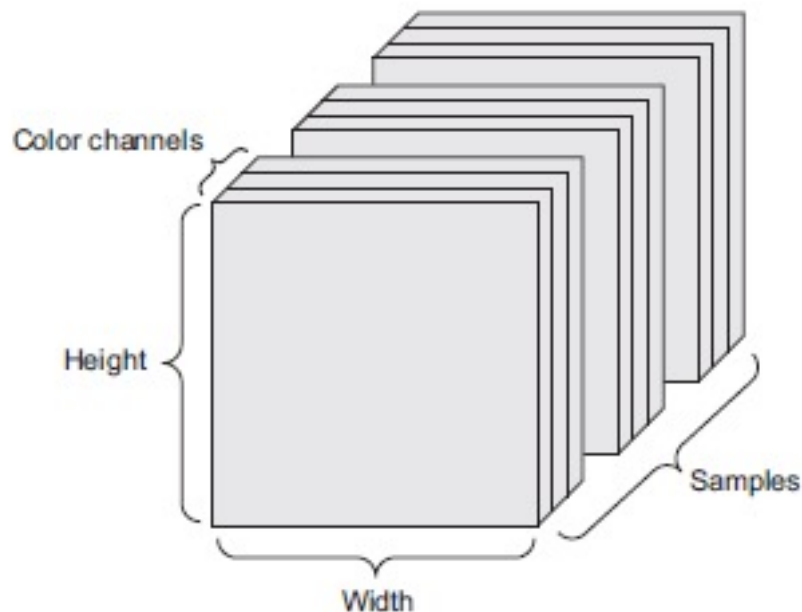


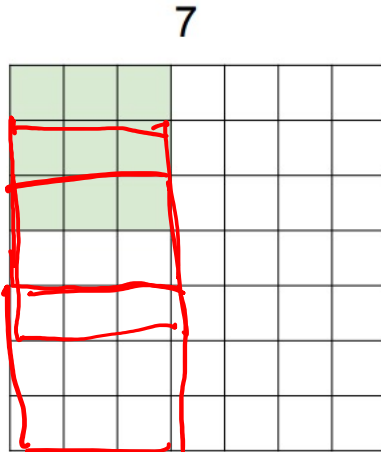
Image Representation

- Image is 3D “tensor”: height, width, color channel (RGB)
- Black-and-white images are 2D matrices: height, width
 - Each value is pixel intensity

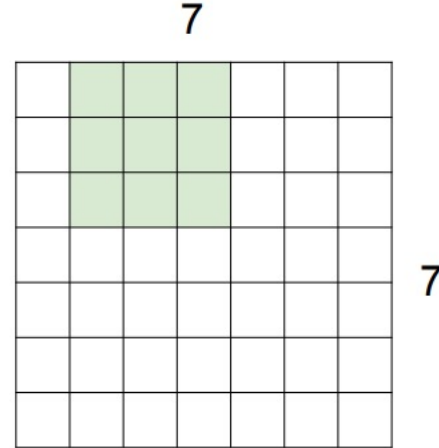


Convolutions

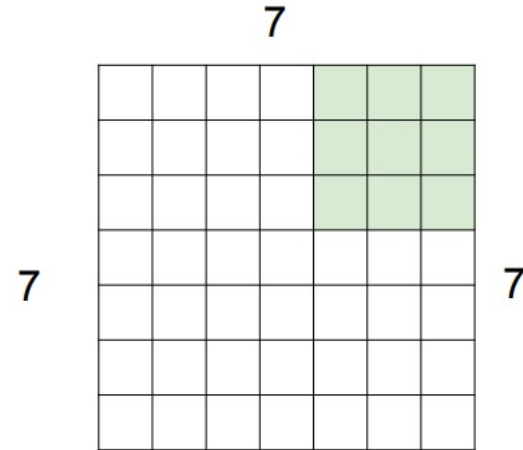
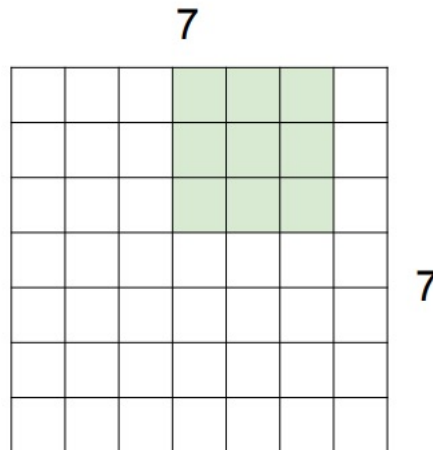
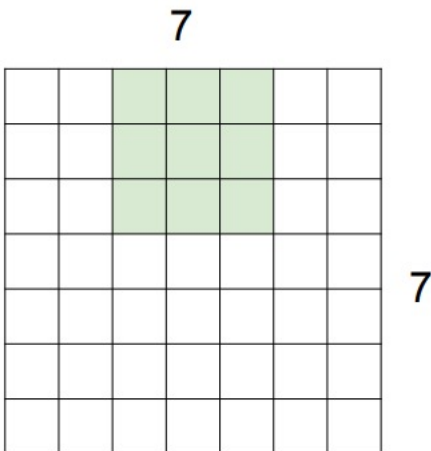
A closer look at spatial dimensions:



7x7 input (spatially)
assume 3x3 filter



=) 5x5



Example

$$[0, 1, 3, 4] \cdot [0, 1, 2, 3] = 19$$

0	1	2
3	4	5
6	7	8

Input

*

0	1
2	3

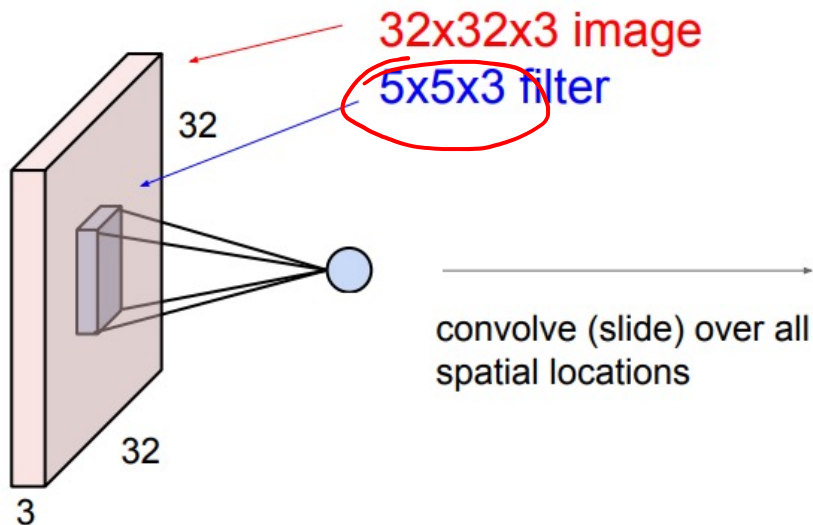
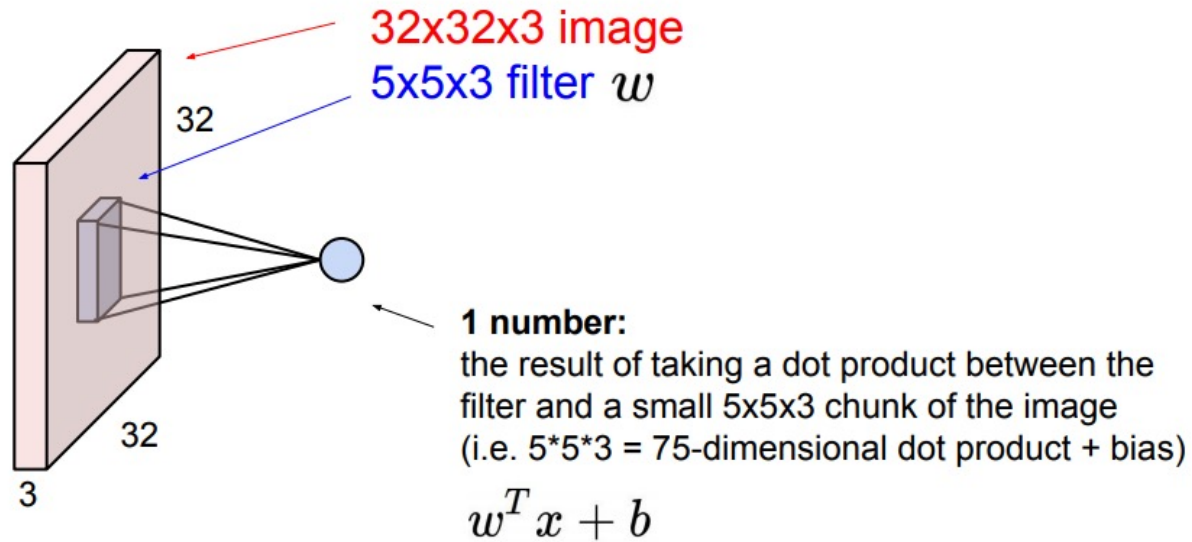
Filter

=

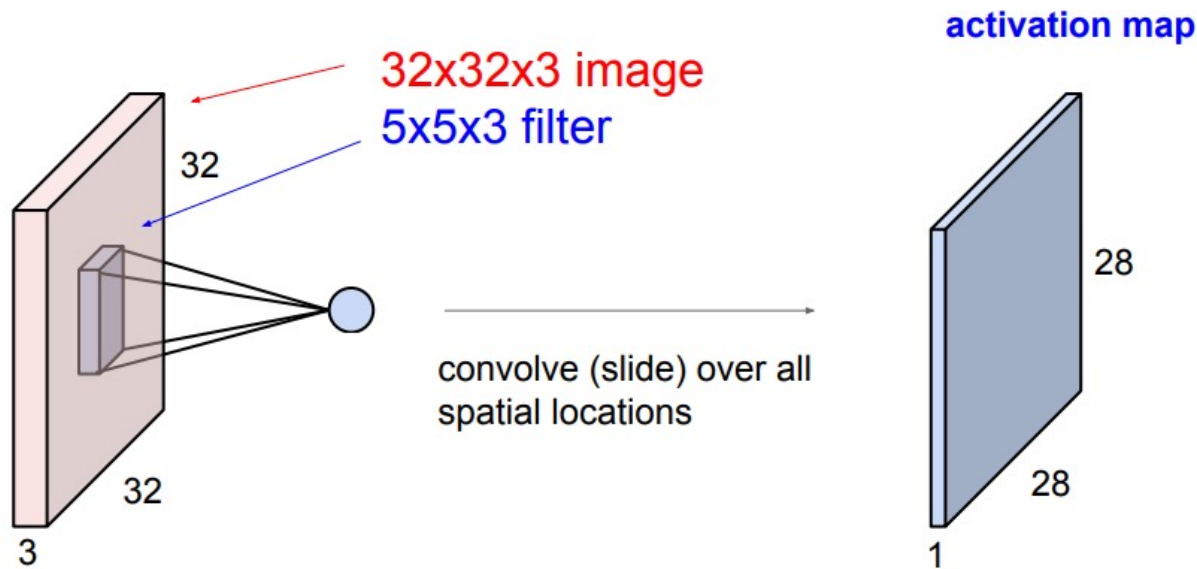
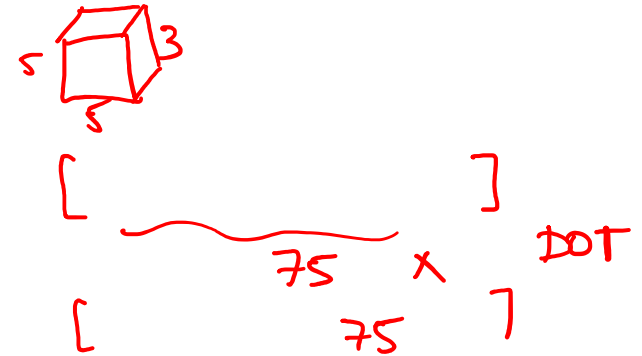
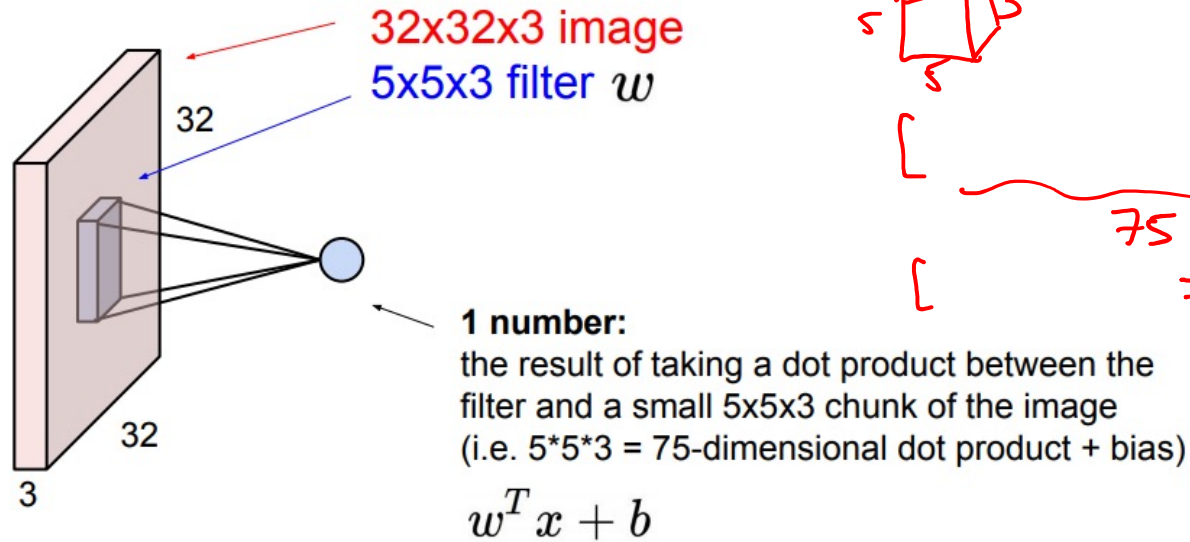
19	25
37	43

Output

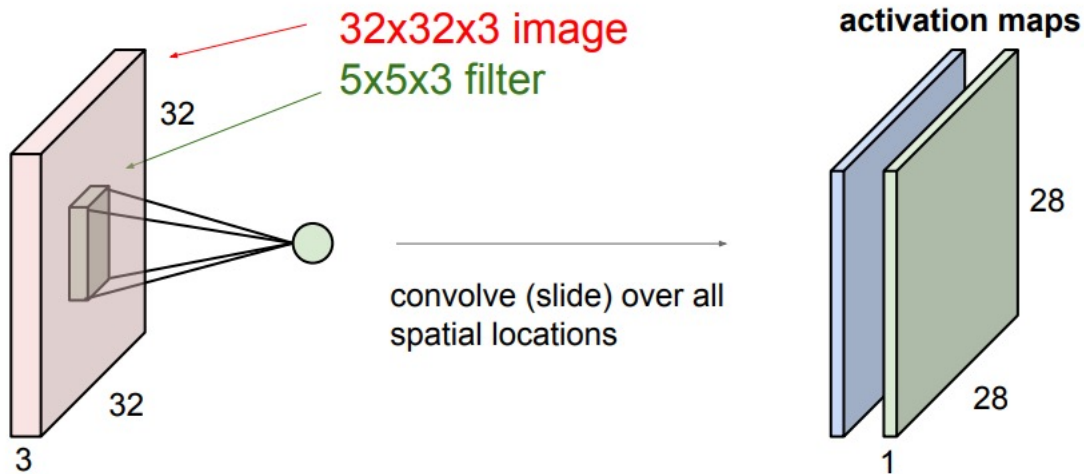
Convolution Layer



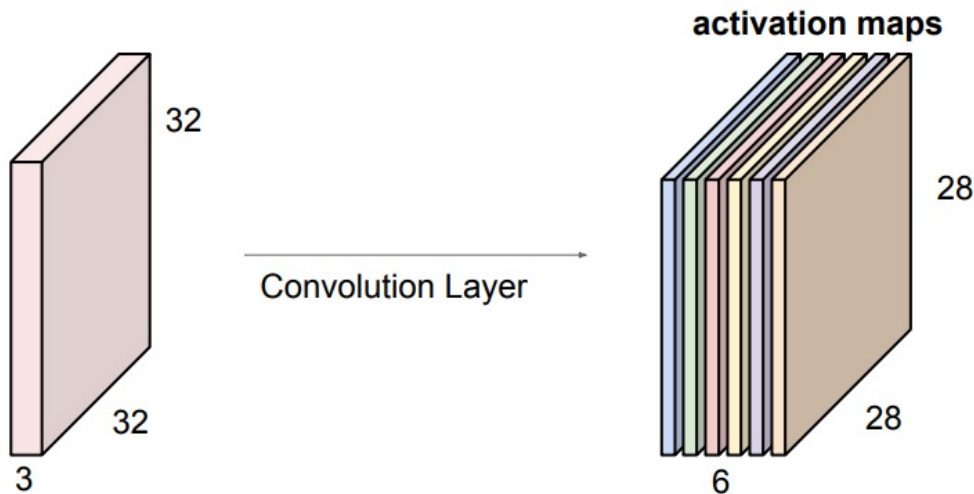
Convolution Layer



Convolution Layer

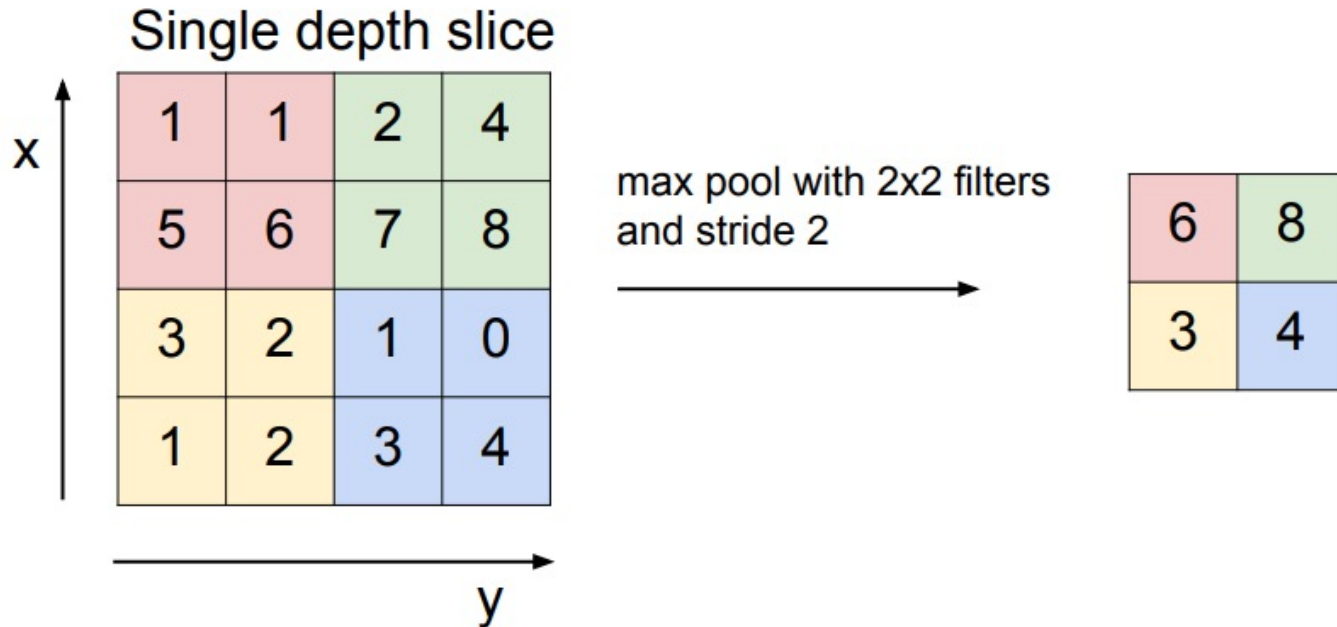


Second, green filter



6 filters

Max Pooling

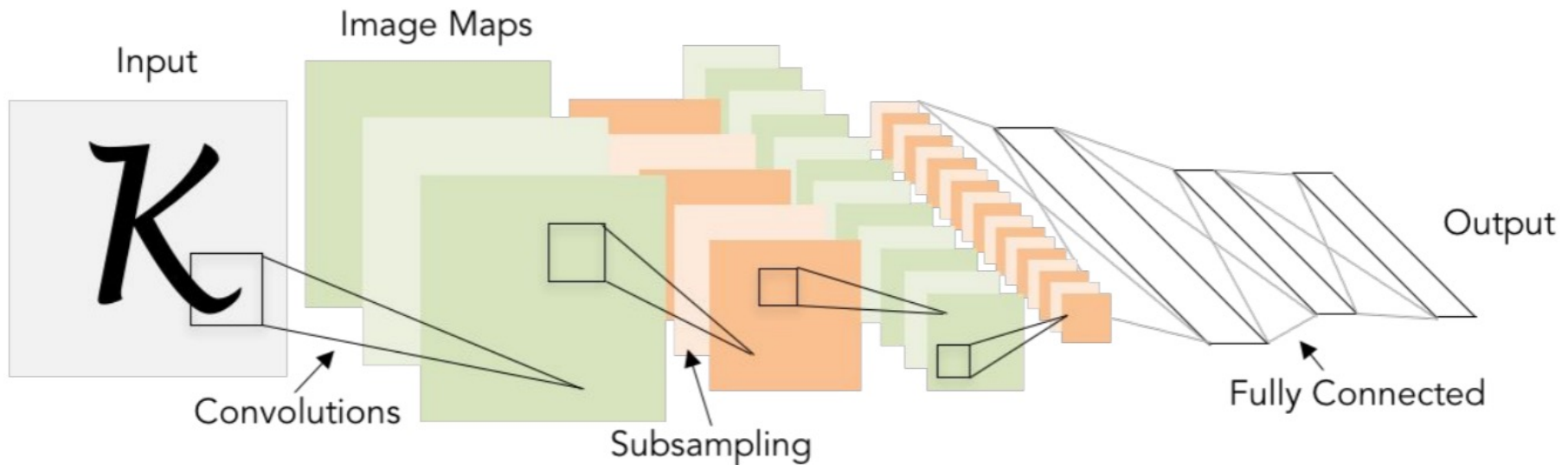


- Accepts a volume of size $W_1 \times H_1 \times D_1$
- Requires three hyperparameters:
 - their spatial extent F ,
 - the stride S ,
- Produces a volume of size $W_2 \times H_2 \times D_2$ where:
 - $W_2 = (W_1 - F) / S + 1$
 - $H_2 = (H_1 - F) / S + 1$
 - $D_2 = D_1$
- Introduces zero parameters since it computes a fixed function of the input
- Note that it is not common to use zero-padding for Pooling layers

LeNet 5

MNIST

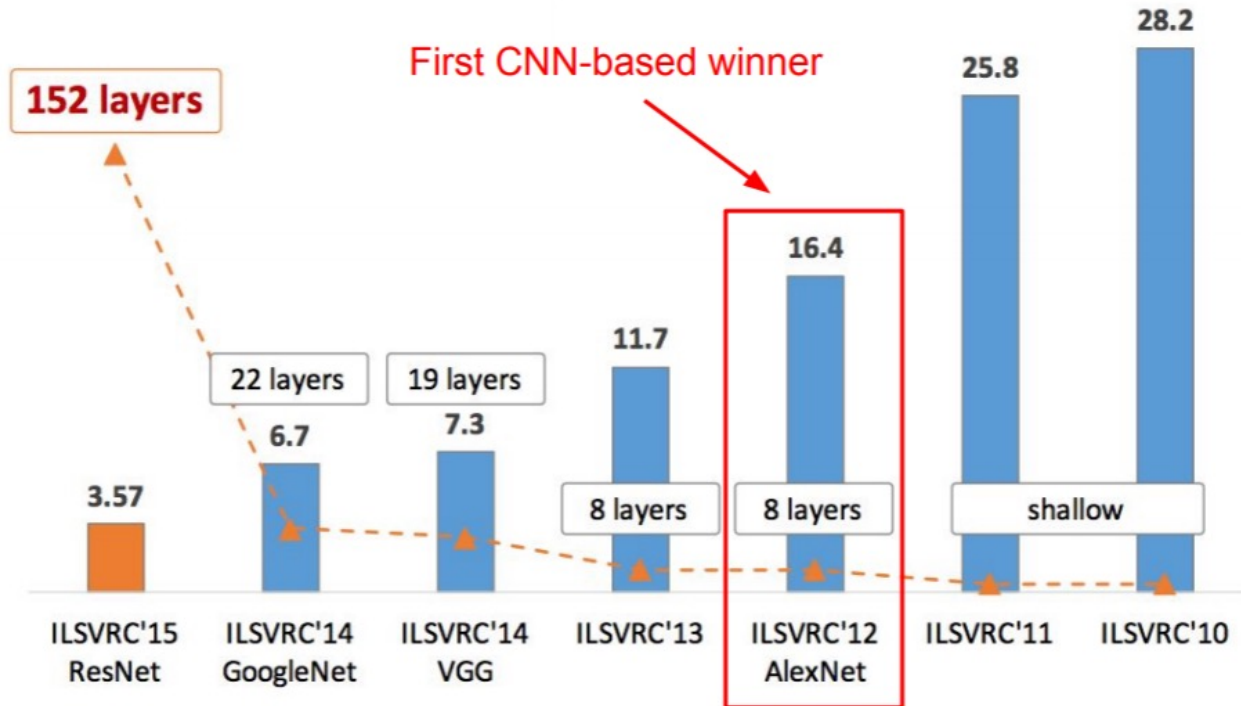
[LeCun et al., 1998]



Conv filters were 5x5, applied at stride 1
Subsampling (Pooling) layers were 2x2 applied at stride 2
i.e. architecture is [CONV-POOL-CONV-POOL-FC-FC]

History

ImageNet Large Scale Visual Recognition Challenge (ILSVRC) winners



VGGNet

Case Study: VGGNet

[Simonyan and Zisserman, 2014]

Small filters, Deeper networks

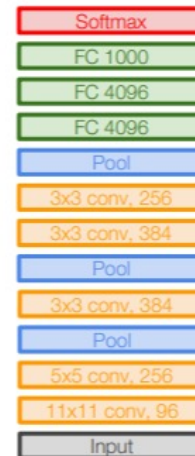
8 layers (AlexNet)

-> 16 - 19 layers (VGG16Net)

Only 3x3 CONV stride 1, pad 1
and 2x2 MAX POOL stride 2

11.7% top 5 error in ILSVRC'13
(ZFNet)

-> 7.3% top 5 error in ILSVRC'14



AlexNet



VGG16

VGG19

138 million
parameters

How to train Neural Networks?

- Backpropagation algorithm
- David Rumelhart, Geoffrey Hinton, Ronald Williams. "Learning representations by back-propagating errors". Nature. 323 (6088): 533–536. 1986
- Applicable to both FFNN and CNN
- Extension of Gradient Descent to multi-layer neural networks

Training Neural Networks

- Training data $x_1, y_1, \dots, x_N, y_N$ BINARY, $y_i \in \{0, 1\}$
- One training example $x_i = (x_{i1}, \dots, x_{id})$, label y_i
- One forward pass through the network
 - Compute prediction $\hat{y}_i = h(x_i)$
- Loss function for one example
 - $L(\hat{y}, y) = -[(1 - y) \log(1 - \hat{y}) + y \log \hat{y}]$

Cross-entropy loss

- Loss function for training data

$$- J(W, b) = \frac{1}{N} \sum_{i=1}^N L(\hat{y}_i, y_i)$$

GD for Neural Networks

$$\theta = [W^{[\ell]} ; b^{[\ell]}], \ell = 1, \dots, N$$

- Initialization

- For all layers ℓ
 - Initialize $W^{[\ell]}, b^{[\ell]}$

$$\theta \leftarrow \theta - \alpha \frac{\nabla_{\theta} J(\theta)}{|\nabla \theta|}$$

FORWARD



- Backpropagation

- Fix learning rate α

- For all layers ℓ (starting backwards) = $L, \dots, 1$

← MEASURE ERROR
BACKWARD

- $W^{[\ell]} = W^{[\ell]} - \alpha \sum_{i=1}^N \frac{\partial L(\hat{y}_i, y_i)}{\partial W^{[\ell]}}$
 - $b^{[\ell]} = b^{[\ell]} - \alpha \sum_{i=1}^N \frac{\partial L(\hat{y}_i, y_i)}{\partial b^{[\ell]}}$

GD for Neural Networks

BATCH GD

- Initialization

- For all layers ℓ
 - Set $W^{[\ell]}, b^{[\ell]}$ at random

- Backpropagation

- Fix learning rate α
- Repeat
 - For all layers ℓ (starting backwards)

ALL TRAINING DATA

$$\bullet W^{[\ell]} = W^{[\ell]} - \alpha \sum_{i=1}^N \frac{\partial L(\hat{y}_i, y_i)}{\partial W^{[\ell]}}$$

$$\bullet b^{[\ell]} = b^{[\ell]} - \alpha \sum_{i=1}^N \frac{\partial L(\hat{y}_i, y_i)}{\partial b^{[\ell]}}$$

This is
expensive!

Stochastic Gradient Descent (SGD)

- Initialization

- For all layers ℓ
 - Set $W^{[\ell]}, b^{[\ell]}$ at random

- Backpropagation

- Fix learning rate α
- Repeat
 - For all layers ℓ (starting backwards)

- For all training examples x_i, y_i

$$W^{[\ell]} = W^{[\ell]} - \alpha \frac{\partial L(\hat{y}_i, y_i)}{\partial W^{[\ell]}}$$

$$b^{[\ell]} = b^{[\ell]} - \alpha \frac{\partial L(\hat{y}_i, y_i)}{\partial b^{[\ell]}}$$

Incremental
version of GD

Mini-batch Gradient Descent

- Initialization

$$\min_{\Theta} L(\Theta)$$

- For all layers ℓ
 - Set $W^{[\ell]}, b^{[\ell]}$ at random

- Backpropagation

- Fix learning rate α
- Repeat

- For all layers ℓ (starting backwards)
 - For all batches b of size B with training examples x_{ib}, y_{ib}

$$\left\{ \begin{aligned} W^{[\ell]} &= W^{[\ell]} - \alpha \sum_{i=1}^B \frac{\partial L(\hat{y}_{ib}, y_{ib})}{\partial W^{[\ell]}} \\ b^{[\ell]} &= b^{[\ell]} - \alpha \sum_{i=1}^B \frac{\partial L(\hat{y}_{ib}, y_{ib})}{\partial b^{[\ell]}} \end{aligned} \right.$$

SHUFFLE POINTS
RANDOMLY

Gradient Descent Variants

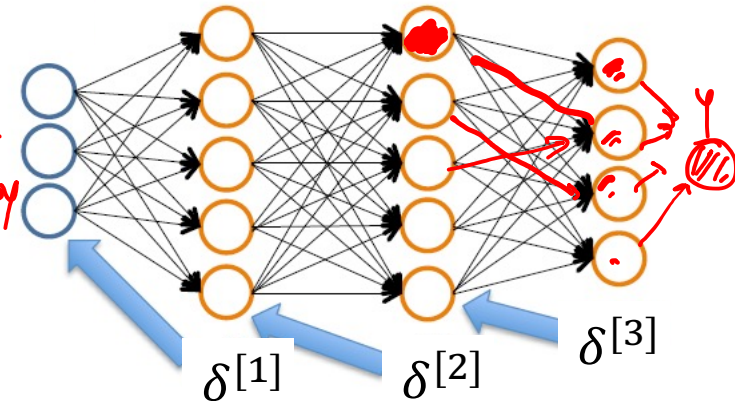


Backpropagation

Let $\delta_j^{(l)}$ = "error" of node j in layer l

$$L(y, \hat{y}) = -[(1 - y) \log(1 - \hat{y}) + y \log \hat{y}]$$

CROSS-ENTROPY



Definitions

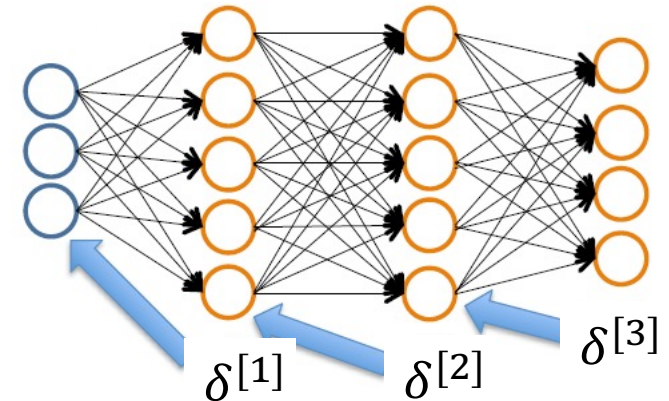
$$\begin{cases} - z^{[\ell]} = W^{[\ell]} a^{[\ell-1]} + b^{[\ell]}, a^{[\ell]} = g(z^{[\ell]}) \\ - \delta^{[\ell]} = \frac{\partial L(\hat{y}, y)}{\partial z^{[\ell]}}; \text{Output } \hat{y} = a^{[L]} = g(z^{[L]}) \end{cases}$$

$$\begin{aligned} 1) \text{ LAYER } L: \delta^{[L]} &= \frac{\partial L(\hat{y}, \hat{y})}{\partial z^{[L]}} = \frac{\partial L(\hat{y}, \hat{y})}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial z^{[L]}} = \frac{\partial L(\hat{y}, \hat{y})}{\partial \hat{y}} \cdot g'(z^{[L]}) \\ 2) \text{ layer } l < L: \delta^{[l]} &= \frac{\partial L(\hat{y}, \hat{y})}{\partial z^{[l]}} = \underbrace{\frac{\partial L(\hat{y}, \hat{y})}{\partial z^{[l+1]}}}_{\delta^{[l+1]}} \cdot \underbrace{\frac{\partial z^{[l+1]}}{\partial a^{[l]}}}_{W^{[l+1]}} \cdot \frac{\partial a^{[l]}}{\partial z^{[l]}} \\ &= \delta^{[l+1]} \cdot W^{[l+1]} \cdot g'(z^{[l]}) \end{aligned}$$

Backpropagation

Let $\delta_j^{(l)}$ = “error” of node j in layer l

$$L(y, \hat{y}) = -[(1 - y) \log(1 - \hat{y}) + y \log \hat{y}]$$



Definitions

$$- z^{[\ell]} = W^{[\ell]} a^{[\ell-1]} + b^{[\ell]}, a^{[\ell]} = g(z^{[\ell]})$$

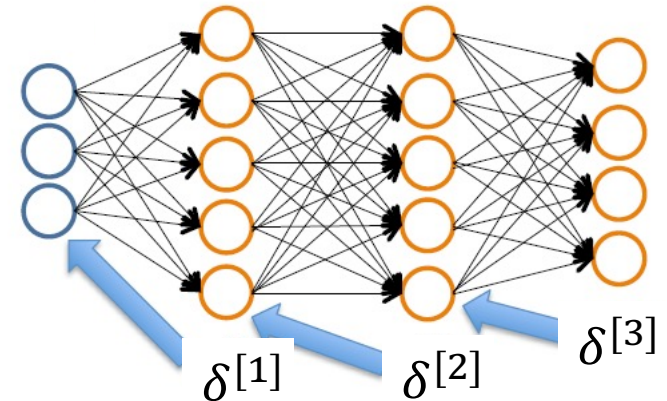
$$- \delta^{[\ell]} = \frac{\partial L(\hat{y}, y)}{\partial z^{[\ell]}}; \text{Output } \hat{y} = a^{[L]} = g(z^{[L]})$$

$$\left\{ \begin{aligned} \frac{\partial L(y, \hat{y})}{\partial w^{[\ell]}} &= \frac{\partial L(y, \hat{y})}{\partial z^{[\ell]}} \cdot \frac{\partial z^{[\ell]}}{\partial w^{[\ell]}} = \delta^{[\ell]} \cdot a^{[\ell-1]T} \\ \frac{\partial L(y, \hat{y})}{\partial b^{[\ell]}} &= \frac{\partial L(y, \hat{y})}{\partial z^{[\ell]}} \cdot \frac{\partial z^{[\ell]}}{\partial b^{[\ell]}} = \delta^{[\ell]} \end{aligned} \right.$$

Backpropagation

Let $\delta_j^{(l)}$ = “error” of node j in layer l

$$L(y, \hat{y}) = -[(1 - y) \log(1 - \hat{y}) + y \log \hat{y}]$$



Definitions

- $z^{[\ell]} = W^{[\ell]} a^{[\ell-1]} + b^{[\ell]}, a^{[\ell]} = g(z^{[\ell]})$
- $\delta^{[\ell]} = \frac{\partial L(\hat{y}, y)}{\partial z^{[\ell]}}$; Output $\hat{y} = a^{[L]} = g(z^{[L]})$

1. For last layer L : $\delta^{[L]} = \frac{\partial L(\hat{y}, y)}{\partial z^{[L]}} = \frac{\partial L(\hat{y}, y)}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial z^{[L]}} = \frac{\partial L(\hat{y}, y)}{\partial \hat{y}} g'(z^{[L]})$
2. For layer ℓ : $\delta^{[\ell]} = \frac{\partial L(\hat{y}, y)}{\partial z^{[\ell]}} = \frac{\partial L(\hat{y}, y)}{\partial z^{[\ell+1]}} \frac{\partial z^{[\ell+1]}}{\partial a^{[\ell]}} \frac{\partial a^{[\ell]}}{\partial z^{[\ell]}} = \delta^{[\ell+1]} W^{[\ell+1]} g'(z^{[\ell]})$
3. Compute parameter gradients

- $\frac{\partial L(\hat{y}, y)}{\partial W^{[\ell]}} = \frac{\partial L(\hat{y}, y)}{\partial z^{[\ell]}} \frac{\partial z^{[\ell]}}{\partial W^{[\ell]}} = \delta^{[\ell]} a^{[\ell-1]T}$
- $\frac{\partial L(\hat{y}, y)}{\partial b^{[\ell]}} = \frac{\partial L(\hat{y}, y)}{\partial z^{[\ell]}} \frac{\partial z^{[\ell]}}{\partial b^{[\ell]}} = \delta^{[\ell]}$

Training NN with Backpropagation

Given training set $(x_1, y_1), \dots, (x_N, y_N)$

Initialize all parameters $W^{[\ell]}, b^{[\ell]}$ randomly, for all layers ℓ

Loop

Set $\Delta_{ij}^{[l]} = 0$, for all layers l and indices i, j

EPOCH

For each training instance (x_k, y_k) :

Compute $a^{[1]}, a^{[2]}, \dots, a^{[L]}$ via forward propagation

Compute errors $\delta^{[L]}, \delta^{[L-1]}, \dots, \delta^{[1]}$

Aggregate gradients $\Delta_{ij}^{[l]} = \Delta_{ij}^{[l]} + a_j^{[l-1]} \delta_i^{[l]}$

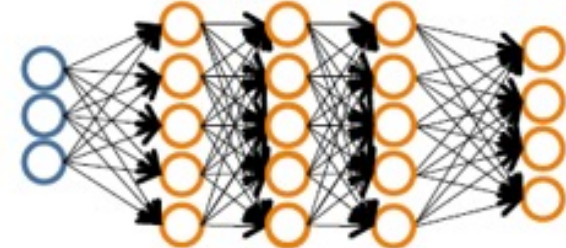
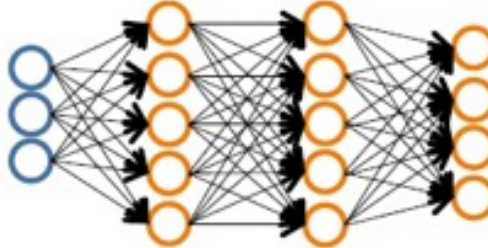
MINI-BATCH

Update weights via gradient step

- $W_{ij}^{[\ell]} = W_{ij}^{[\ell]} - \alpha \Delta_{ij}^{[\ell]}$
- Similarly for $b_{ij}^{[\ell]}$

Until weights converge or maximum number of epochs is reached

Overfitting



- The larger the network, the higher the capacity (more model parameters)
- **But also more prone to overfitting!**

Regularization

$$L(W) = \underbrace{\frac{1}{N} \sum_{i=1}^N L_i(f(x_i, W), y_i)}_{\text{Data loss}} + \underbrace{\lambda R(W)}_{\text{Regularization}}$$

Data loss: Model predictions should match training data

Regularization: Prevent the model from doing *too* well on training data

λ = regularization strength (hyperparameter)

L2 regularization: $R(W) = \sum_k \sum_l W_{k,l}^2$

L1 regularization: $R(W) = \sum_k \sum_l |W_{k,l}|$

Elastic net (L1 + L2): $R(W) = \sum_k \sum_l \beta W_{k,l}^2 + |W_{k,l}|$

RIDGE | L2

Weight decay

- When computing gradients of loss function, regularization term needs to be taken into account

Dropout

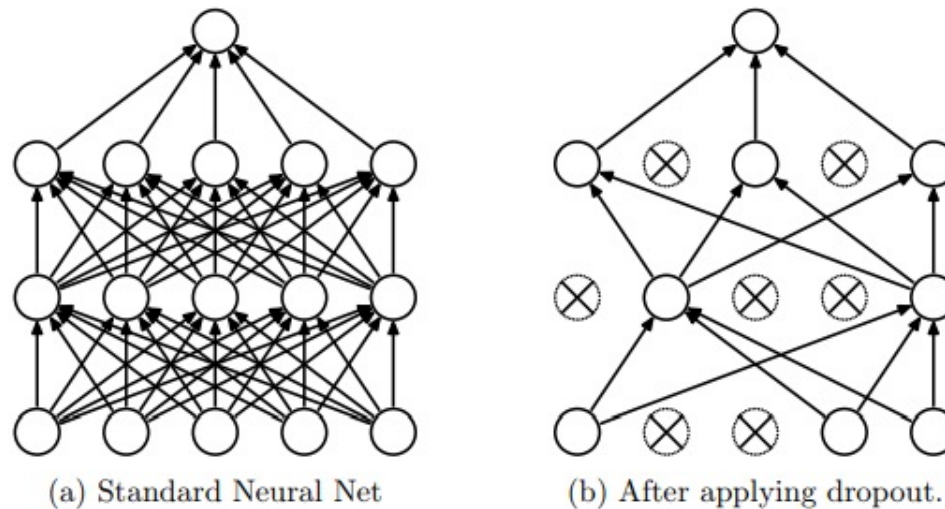


Figure 1: Dropout Neural Net Model. **Left:** A standard neural net with 2 hidden layers. **Right:** An example of a thinned net produced by applying dropout to the network on the left. Crossed units have been dropped.

- At training time, sample a sub-network per epoch (batch) and learn weights
 - Keep each neuron with probability p
- At testing time, all neurons are there, but multiply weight by a factor of p
- Srivastava et al. Dropout: A Simple Way to Prevent Neural Networks from Overfitting. Journal of Machine Learning Research 15, 2014

Comparing classifiers

Algorithm	Interpretable	Model size	Predictive accuracy	Training time	Testing time
Logistic regression	High	Small	Lower	Low	Low
kNN	Medium	Large	Lower	No training	High
LDA	Medium	Small	Lower	Low	Low
Decision trees	High	Medium	Lower	Medium	Low
Ensembles	Low	Large	High	High	High
Naïve Bayes	Medium	Small	Lower	Medium	Low
SVM	Medium	Small	Lower	Medium	Low
Neural Networks	Low	Large	High	High	Low

Summary

- Deep Learning performs well for datasets with certain structure: images, text, speech
 - Learns hierarchical feature representations
- Multiple types of NN architectures
 - Feed-forward neural networks
 - Convolutional neural networks (applicable to images)
 - Recurrent neural networks, Transformers (text)
- Training most NN architectures via backpropagation
- Regularization (weight decay, dropout) to avoid overfitting