



COL333/671: Introduction to AI

Semester I, 2022-23

Learning – II

Rohan Paul

Outline

- Last Class
 - Basics of machine learning
- This Class
 - Neural Networks
- Reference Material
 - Please follow the notes as the primary reference on this topic. Additional reading from AIMA book Ch. 18 (18.2, 18.6 and 18.7) and DL book Ch 6 sections 6.1 – 6.5 (except 6.4).

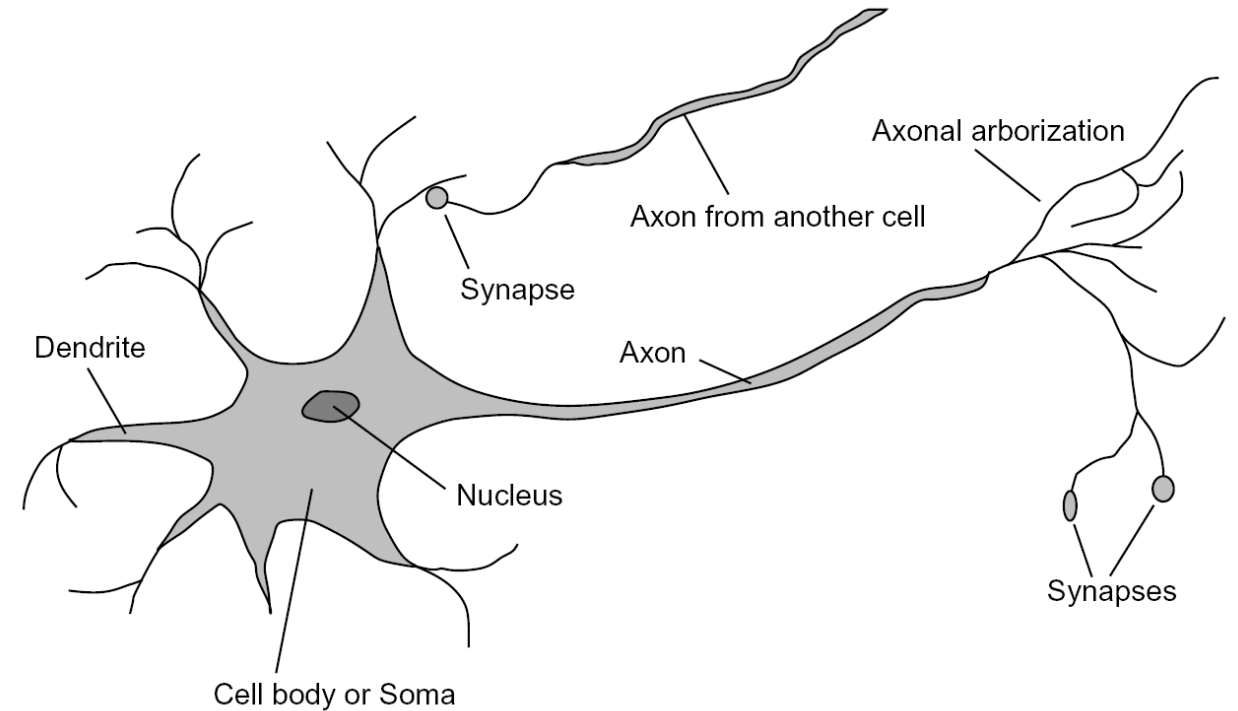
Acknowledgement

These slides are intended for teaching purposes only. Some material has been used/adapted from web sources and from slides by Doina Precup, Dorsa Sadigh, Percy Liang, Mausam, Parag, Emma Brunskill, Alexander Amini, Dan Klein, Anca Dragan, Nicholas Roy and others.

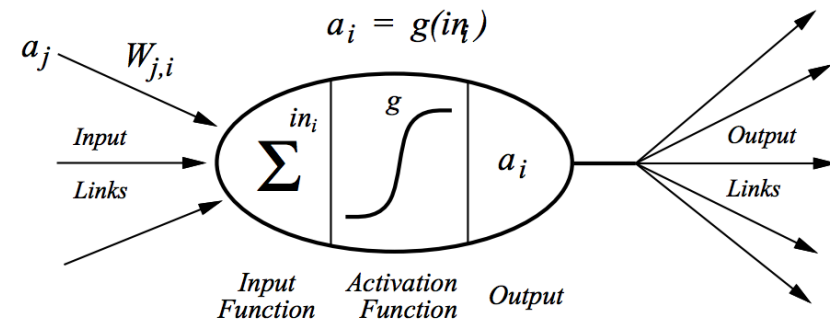
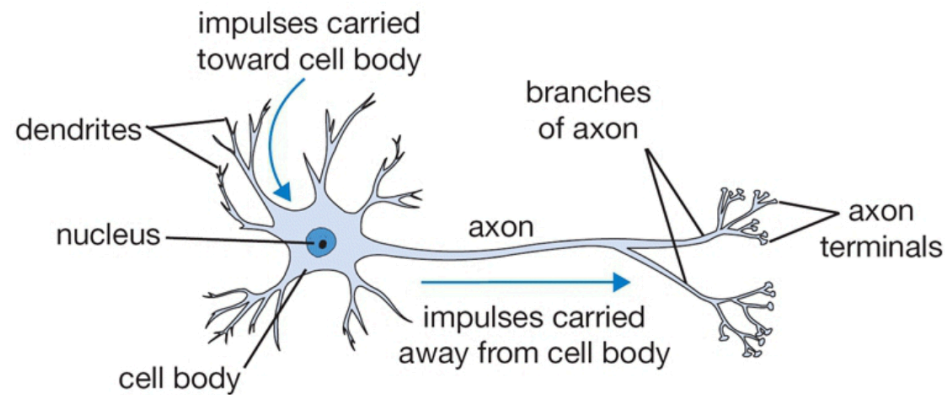
Neuron

A simplified view

- Activations and inhibitions
- Parallelism
- Connected networks



Modeling a Neuron



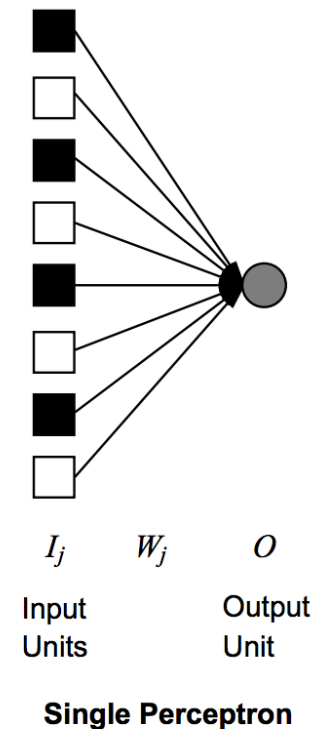
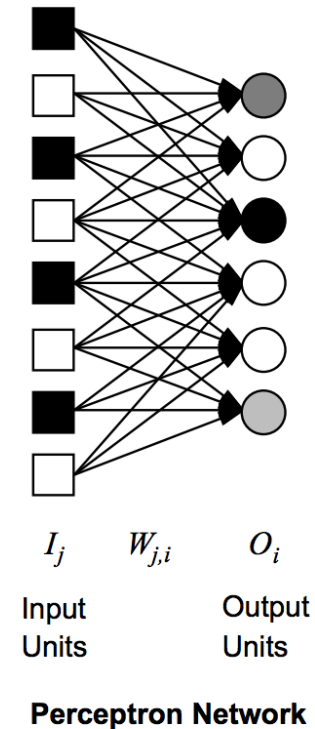
$$a_i = g\left(\sum_j W_{j,i} a_j\right)$$

Main processing unit

- Setup where there is a function that connects the inputs to the output.
- Problem: learning this function that maps the input to the output.

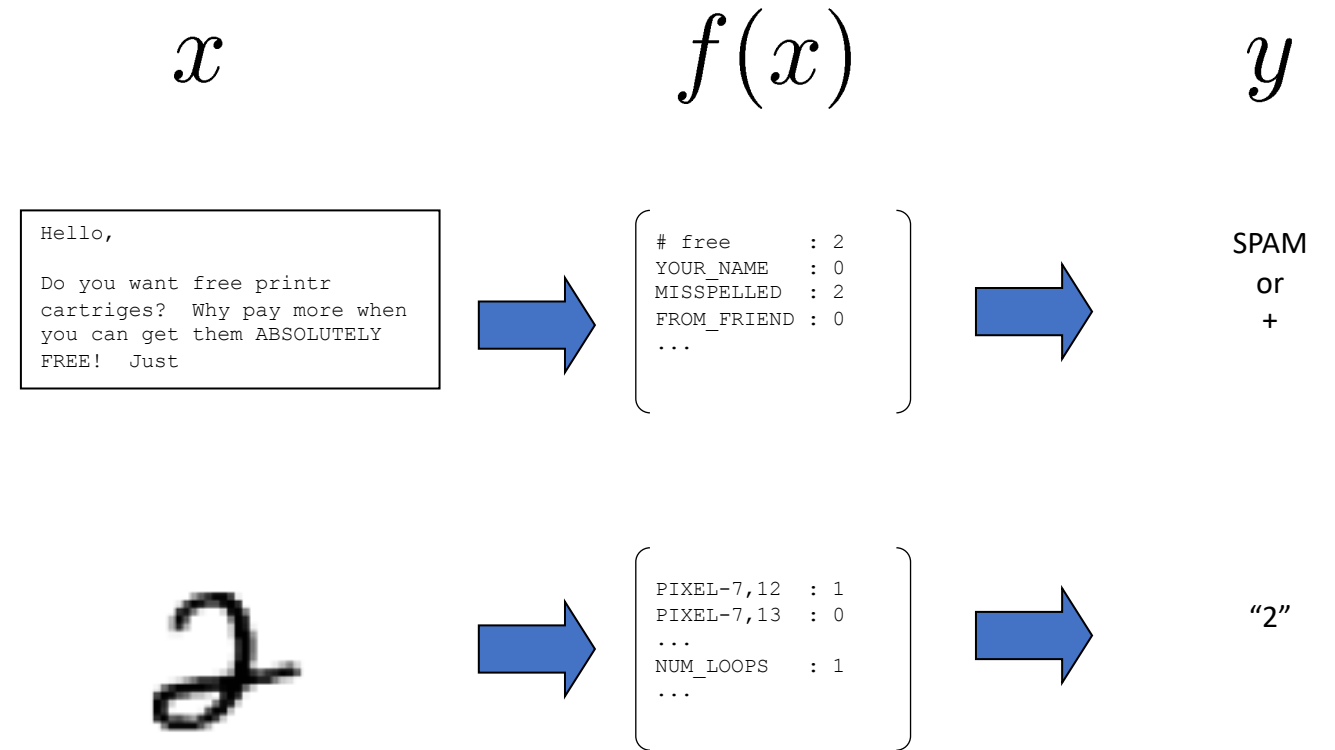
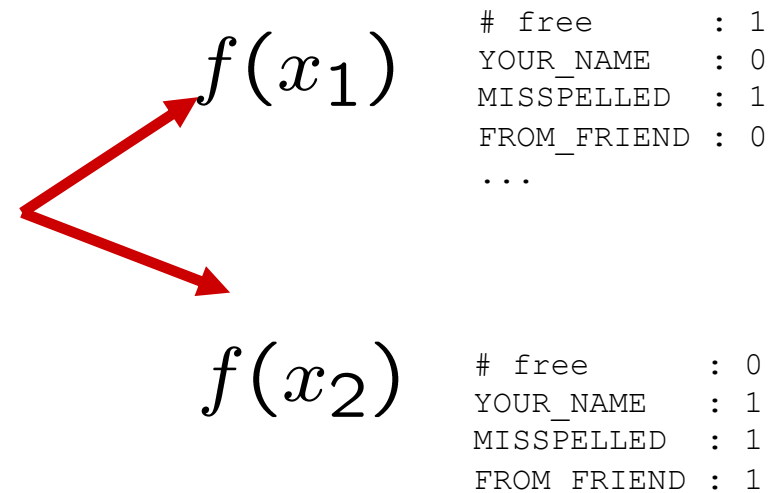
Perceptron

- Introduced in the late 50s
 - Minsky and Papert.
- Perceptron convergence theorem
Rosenblatt 1962:
 - Perceptron will learn to classify any linearly separable set of inputs
- Note: the earlier class talked about model-based classification. Here, we do not build a model. Operate directly on feature weights.



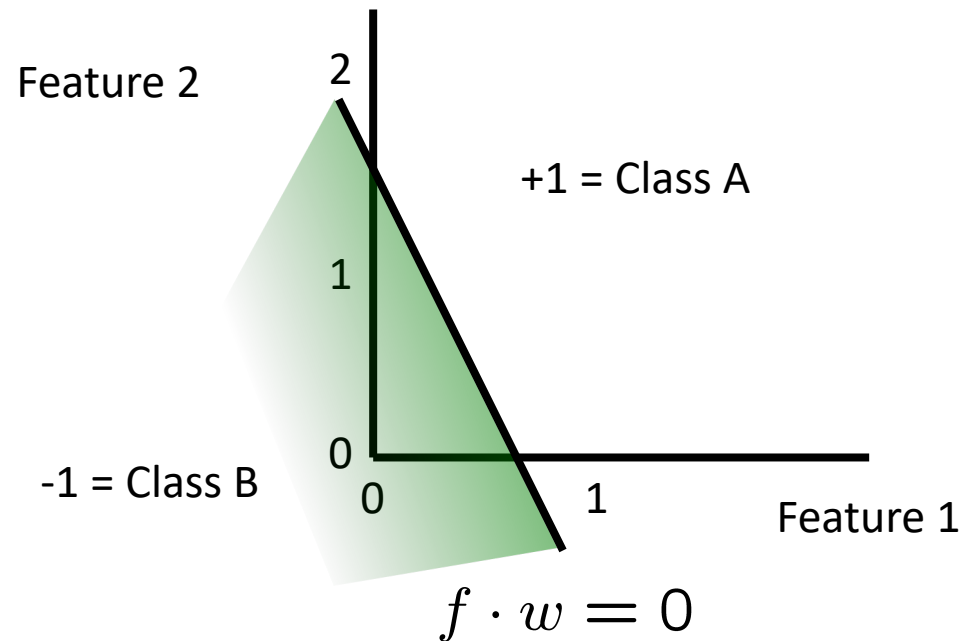
Feature Space

- Extract Features from data
- Learn a model with these features
- Data can be viewed as a point in the feature space.

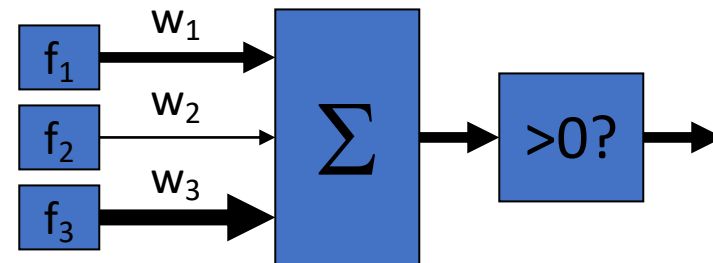


Linear Classification

- A decision boundary is a hyperplane orthogonal to the weight vector.

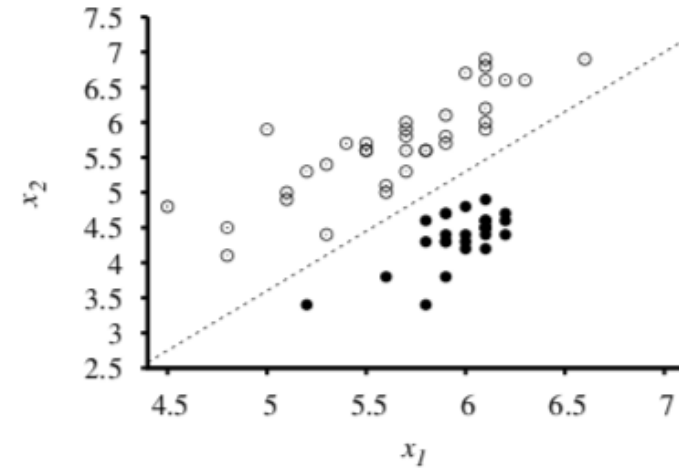


$$\text{activation}_w(x) = \sum_i w_i \cdot f_i(x) = w \cdot f(x)$$



Perceptron

Decision rule (Binary case)



Binary classification task

$h_{\mathbf{w}}(\mathbf{x}) = 1$ if $\mathbf{w} \cdot \mathbf{x} \geq 0$ and 0 otherwise.

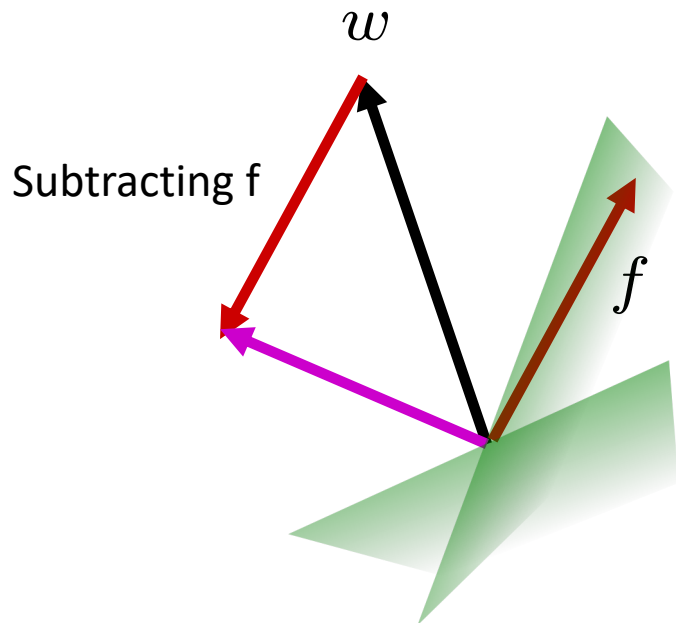
$h_{\mathbf{w}}(\mathbf{x}) = \text{Threshold}(\mathbf{w} \cdot \mathbf{x})$ where $\text{Threshold}(z) = 1$ if $z \geq 0$ and 0 otherwise.

One side of the decision boundary is class A and the other is class B. A threshold is introduced.

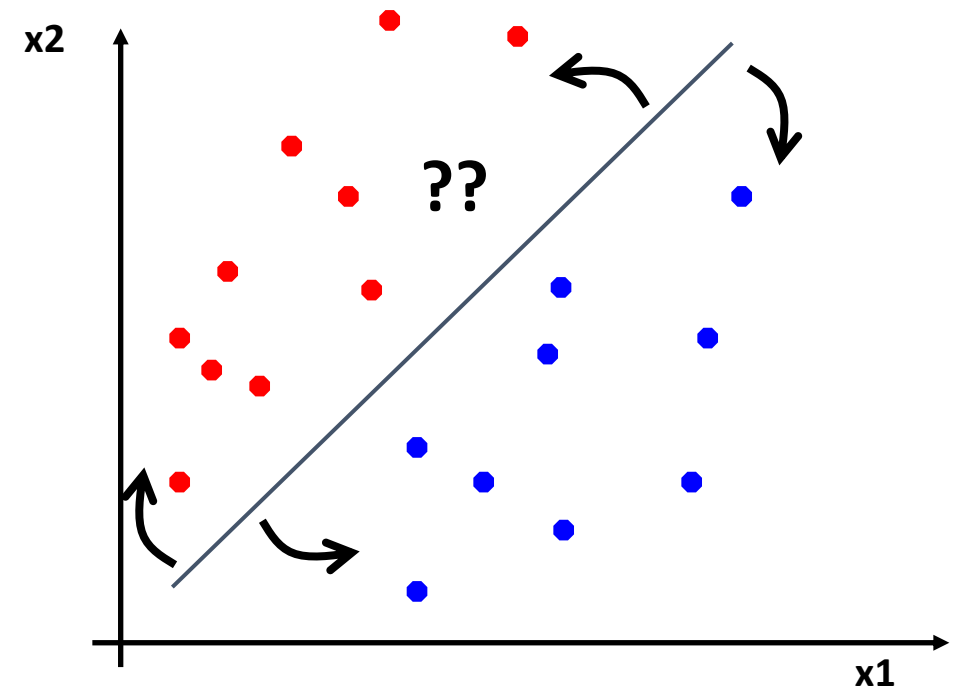
Perceptron

Learning rule

- Classify with the current weights
- If correct no change.
- If the classification is wrong: *adjust* the weight vector by adding or subtracting the feature vector.



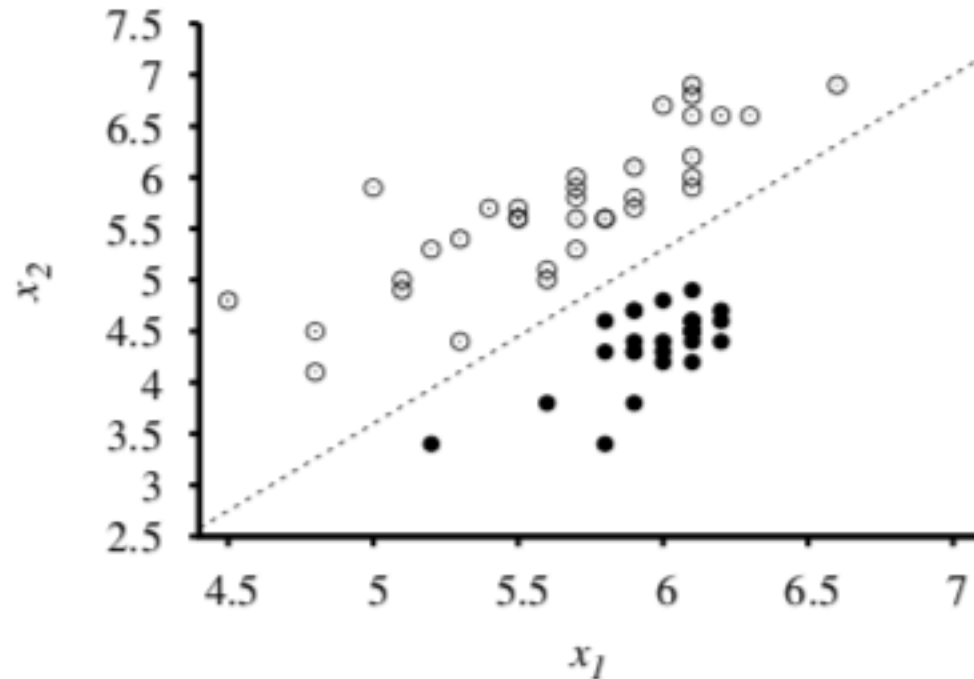
$$w_i \leftarrow w_i + \alpha (y - h_{\mathbf{w}}(\mathbf{X})) \times x_i$$



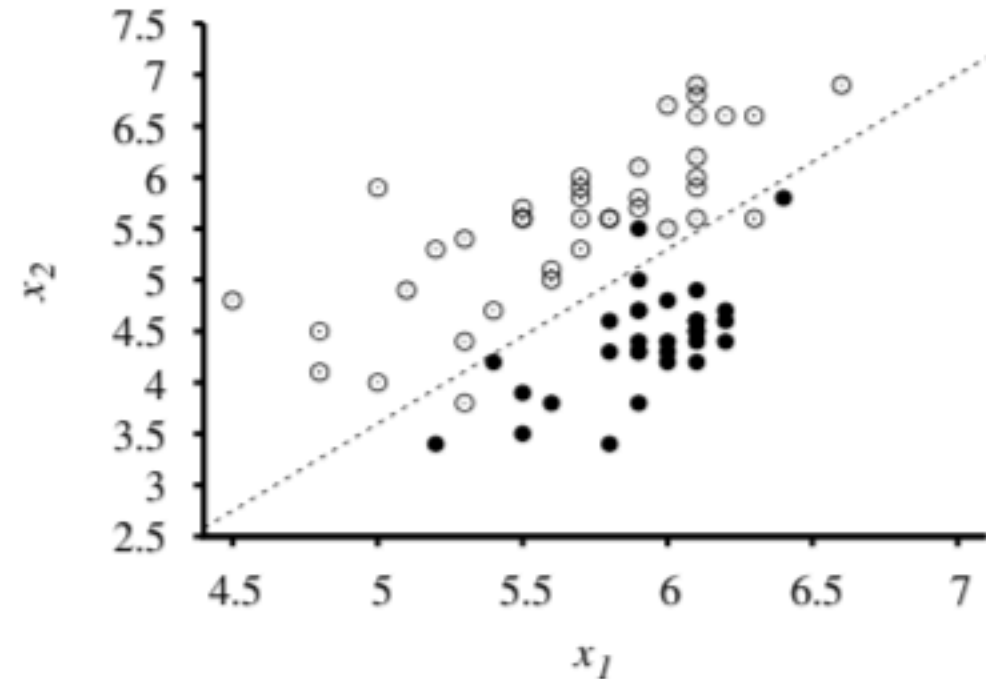
Binary classification task

Perceptron

Case: Linearly-separable

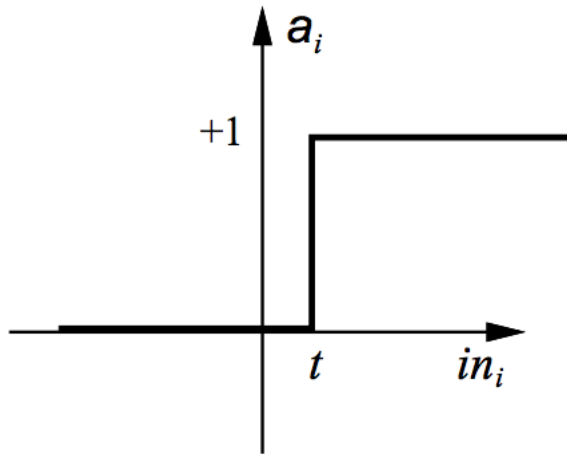


Case: Not Linearly-separable

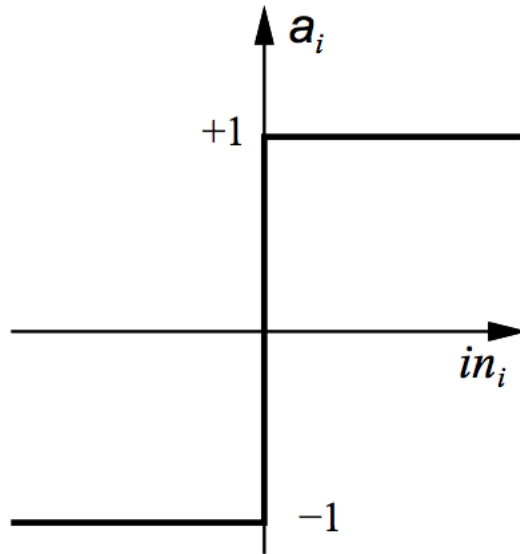


Perceptron learning rule converges to a perfect linear separator when the data points are linearly separable. Problem when there is non-separable data.

Threshold Functions



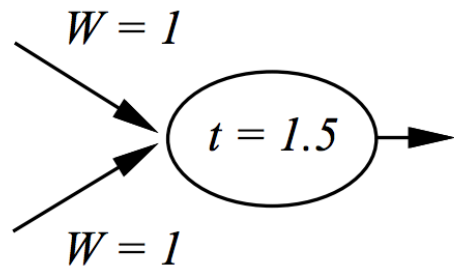
(a) Step function



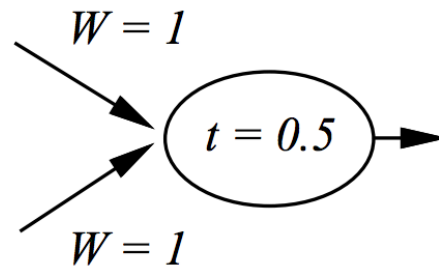
(b) Sign function

- Till now, threshold functions were linear.
- Can we modify the threshold function to handle the non-separable case?
- Can we "soften" the outputs?

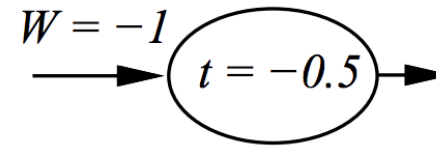
Boolean Functions and Perceptron



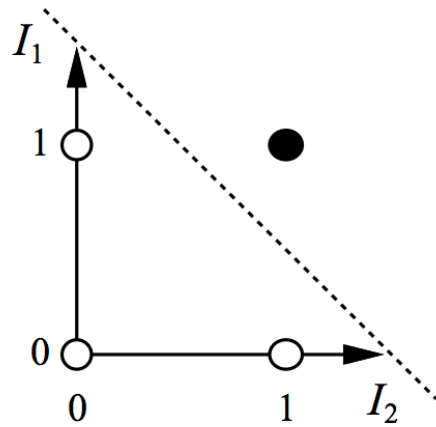
AND



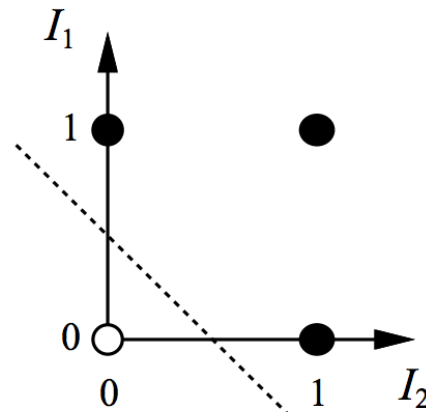
OR



NOT



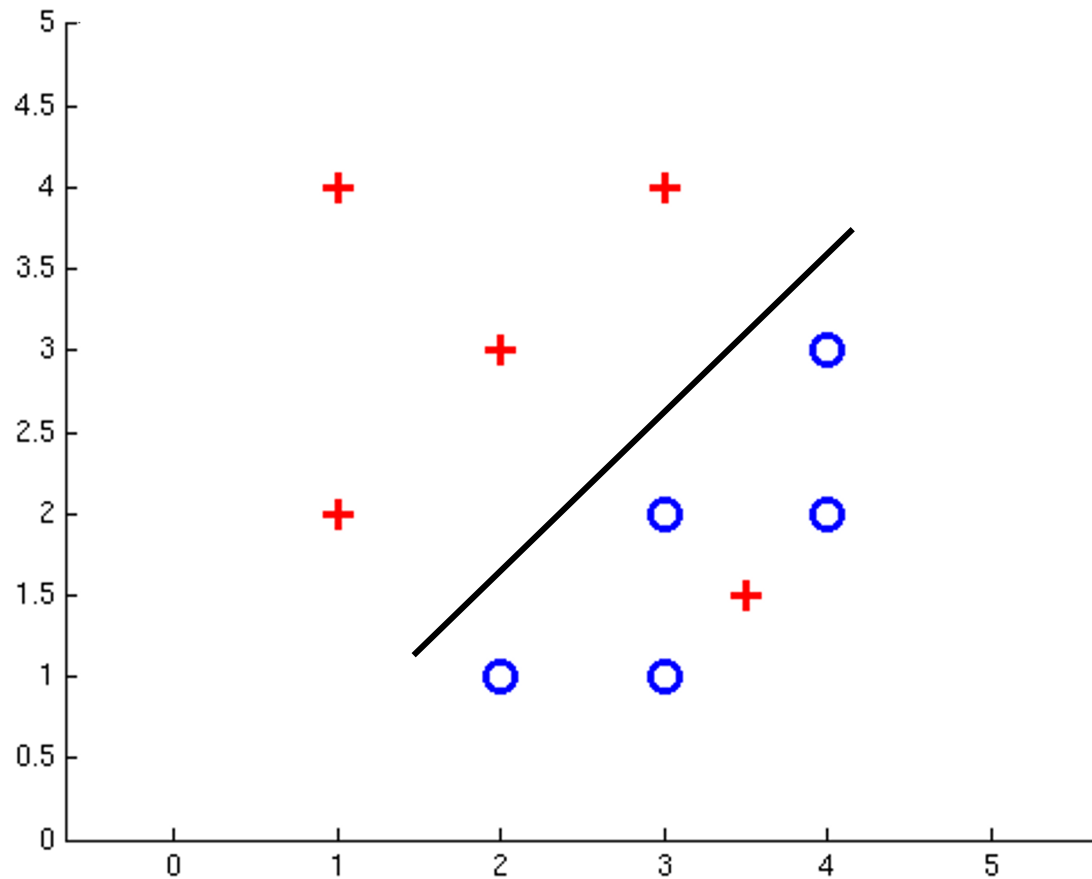
(a) I_1 and I_2



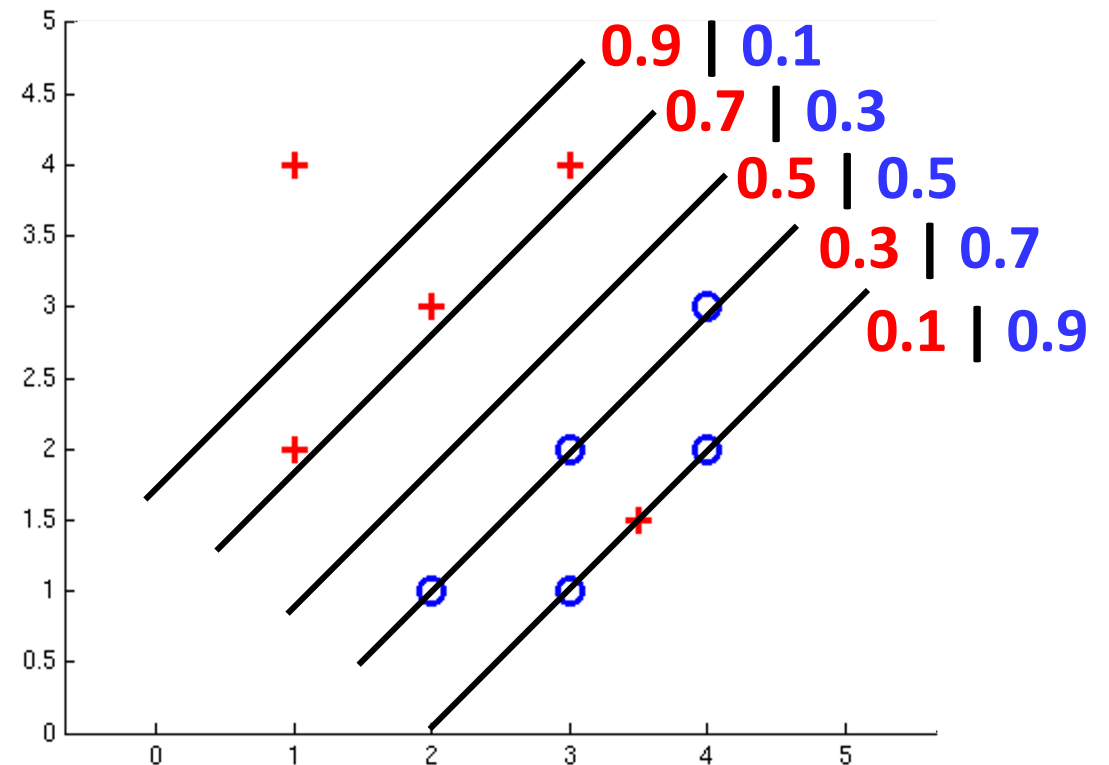
(b) I_1 or I_2

Non-separable case

Deterministic Decisions



Probabilistic Decisions



Logistic Output

- **Logistic Function**

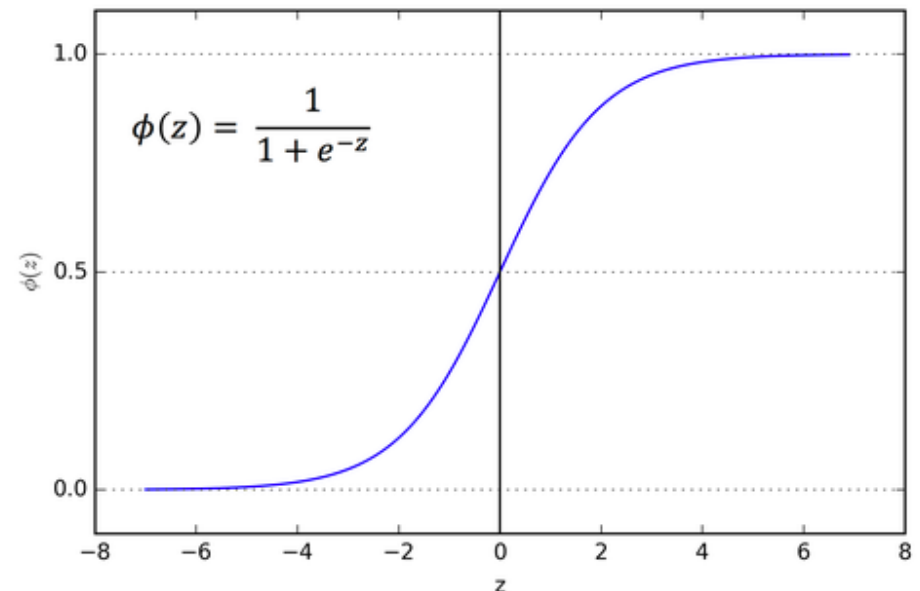
- Very positive values. Probability $\rightarrow 1$
- Very negative values. Probability $\rightarrow 0$.
- Makes the prediction. Converts to a probability
- Softens the decision boundary.

- **Logistic Regression**

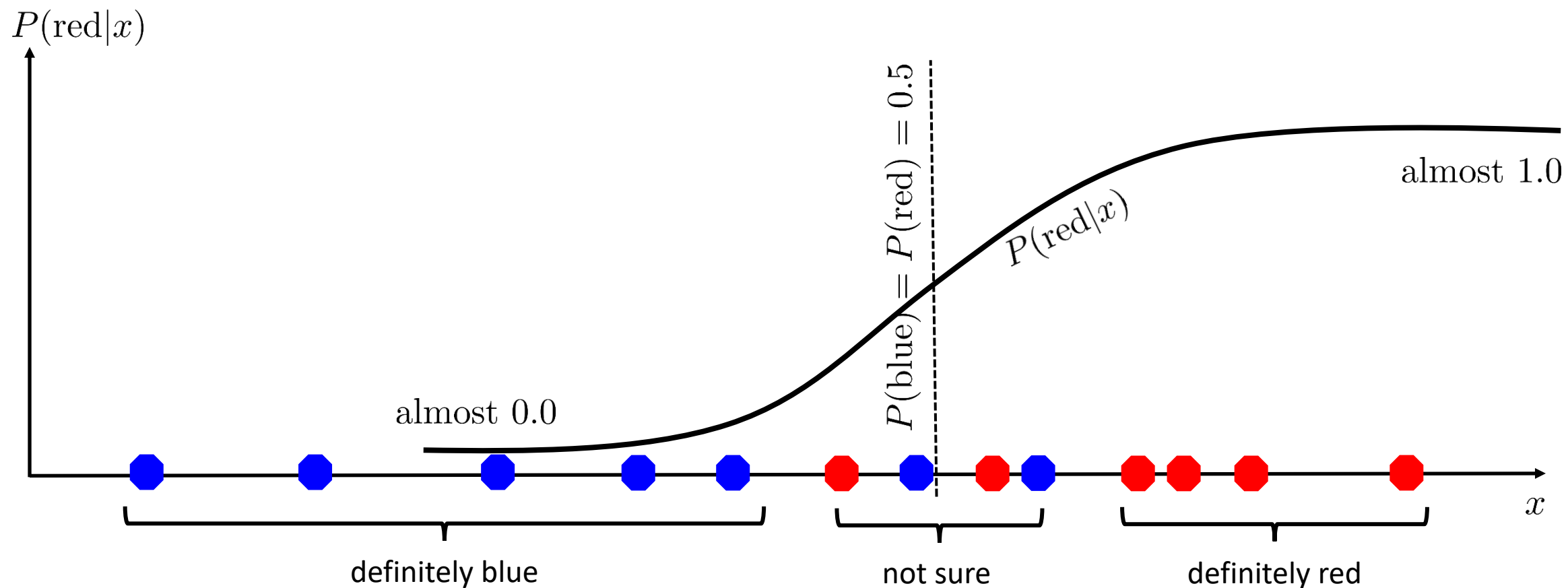
- Fitting the weights of this model to minimize loss on a data set is called logistic regression.

$$\text{Logistic}(z) = \frac{1}{1 + e^{-z}}$$

$$h_{\mathbf{w}}(\mathbf{x}) = \text{Logistic}(\mathbf{w} \cdot \mathbf{x}) = \frac{1}{1 + e^{-\mathbf{w} \cdot \mathbf{x}}}$$



Example (red or blue classes)



$$P(\text{red}|x) = \frac{e^{w_{\text{red}} \cdot x}}{e^{w_{\text{red}} \cdot x} + e^{w_{\text{blue}} \cdot x}}$$

probability increases exponentially as we move away from boundary

Normalizer

Estimating weights using MLE

Logistic Regression

Maximize the log-likelihood

$$\max_w ll(w) = \max_w \sum_i \log P(y^{(i)} | x^{(i)}; w)$$

$$P(y^{(i)} = +1 | x^{(i)}; w) = \frac{1}{1 + e^{-w \cdot f(x^{(i)})}}$$

$$P(y^{(i)} = -1 | x^{(i)}; w) = 1 - \frac{1}{1 + e^{-w \cdot f(x^{(i)})}}$$

Softmax Output

- Multi-class setting
 - A probability distribution over a discrete variable with n possible values.
 - Generalization of the sigmoid function to multiple outputs.
- Output of a classifier
 - Distribution over n different classes. The individual outputs must sum to one.

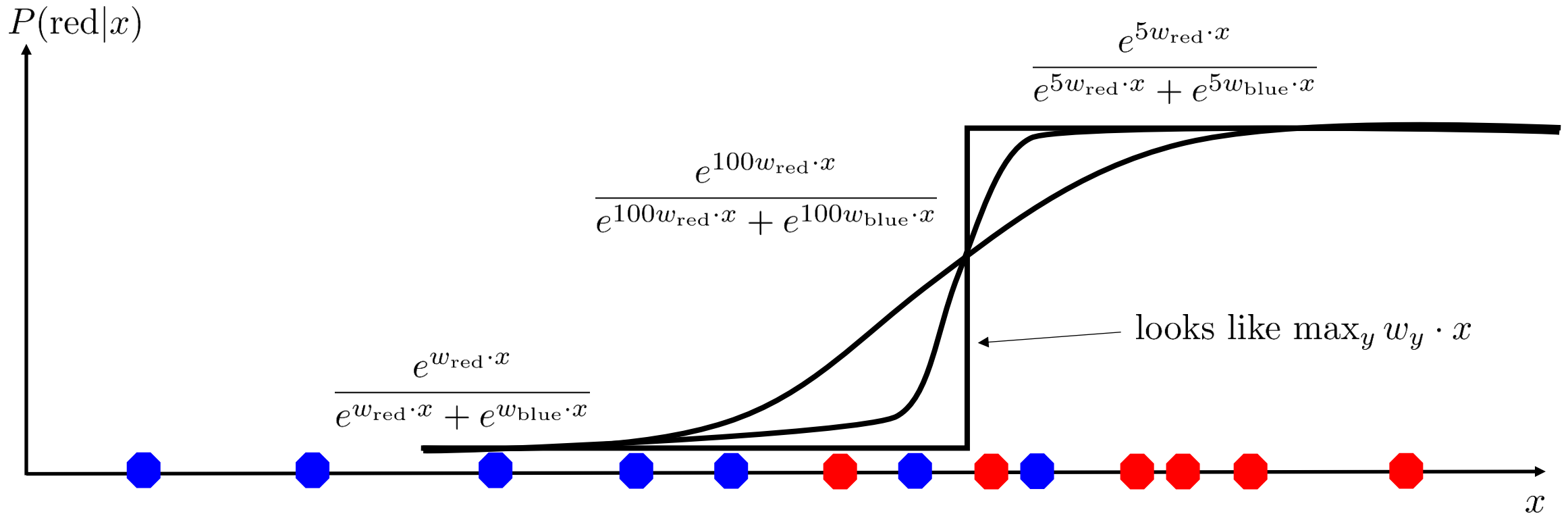
Prediction of the unnormalized probabilities.

$$z_i = \log \tilde{P}(y = i | \mathbf{x})$$

Exponentiate and normalize the values.

$$\text{softmax}(\mathbf{z})_i = \frac{\exp(z_i)}{\sum_j \exp(z_j)}$$

Softmax Example



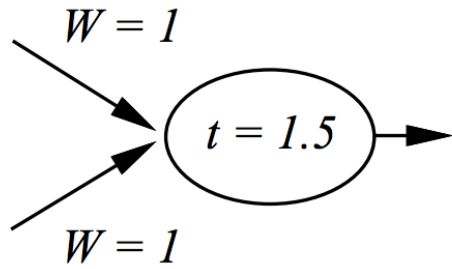
$$P(\text{red}|x) = \frac{e^{w_{\text{red}} \cdot x}}{e^{w_{\text{red}} \cdot x} + e^{w_{\text{blue}} \cdot x}}$$

Multi-class Logistic Regression

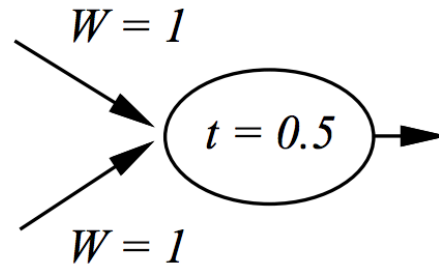
$$\max_w ll(w) = \max_w \sum_i \log P(y^{(i)} | x^{(i)}; w)$$

$$P(y^{(i)} | x^{(i)}; w) = \frac{e^{w_{y^{(i)}} \cdot f(x^{(i)})}}{\sum_y e^{w_y \cdot f(x^{(i)})}}$$

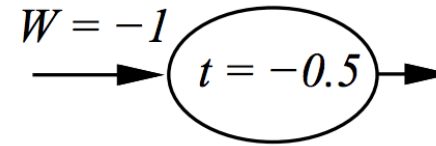
Can a perceptron learn XOR?



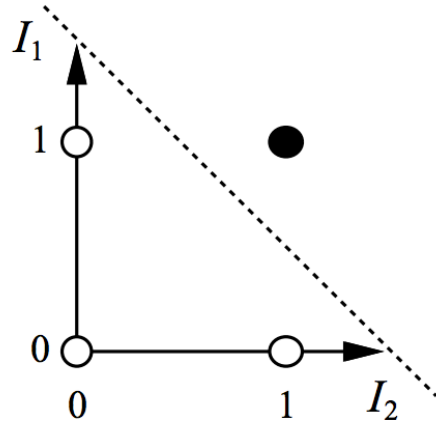
AND



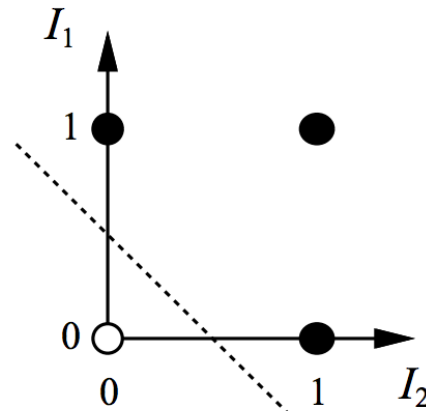
OR



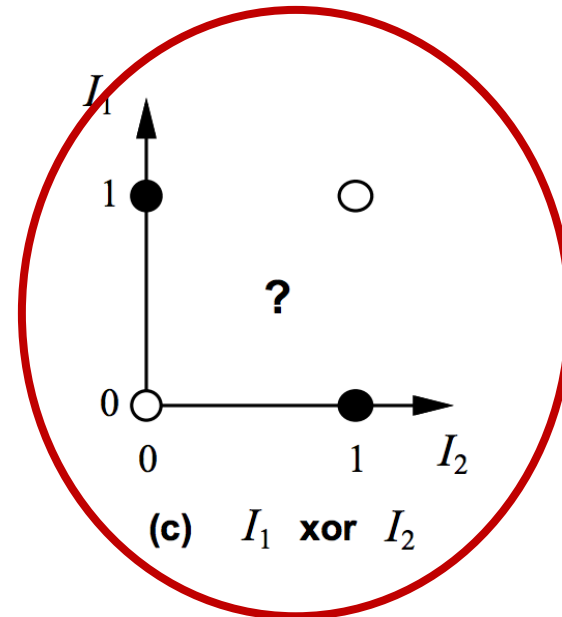
NOT



(a) I_1 and I_2



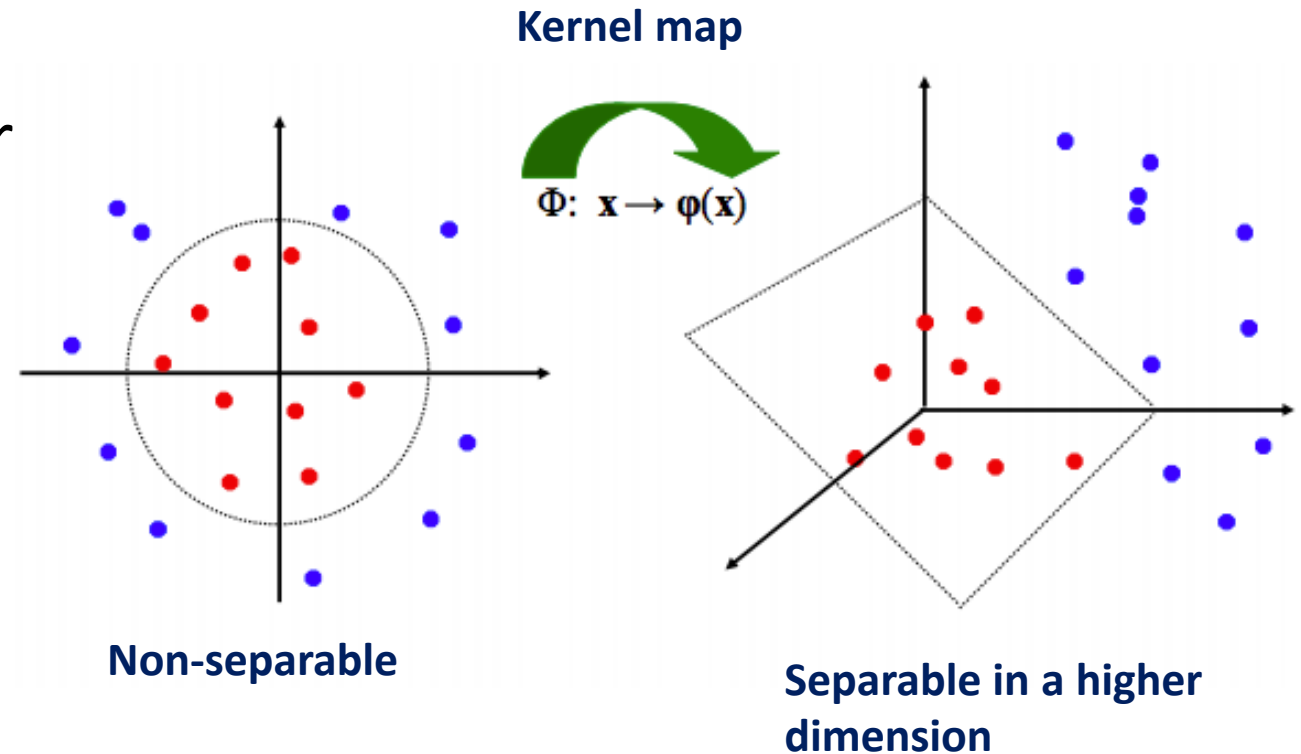
(b) I_1 or I_2



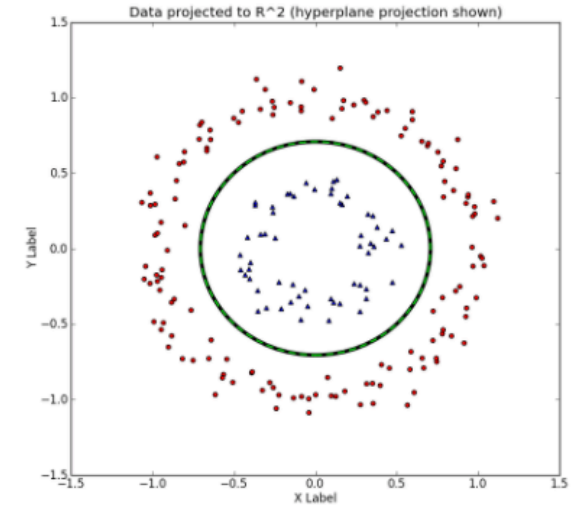
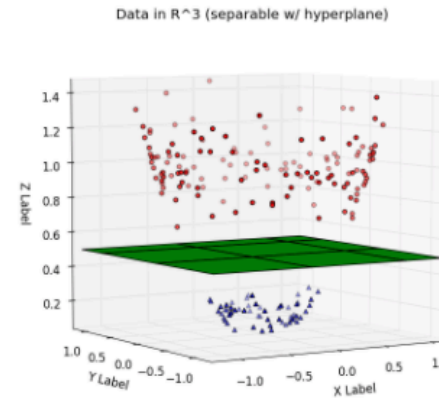
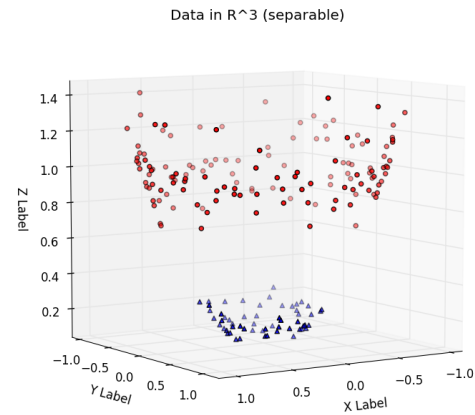
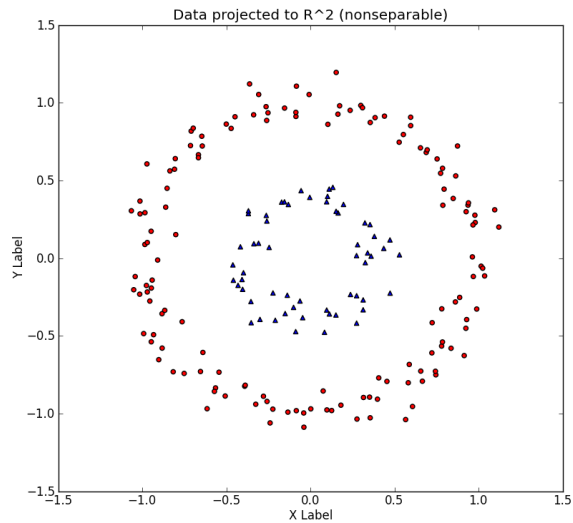
(c) I_1 xor I_2

Non-separability and Non-linear Functions

- The original feature space is mapped to some higher-dimensional feature space where the training set is separable.
- Need a non-linear function to describe the features.
- Applying a non-linear kernel map. Affine transformation.

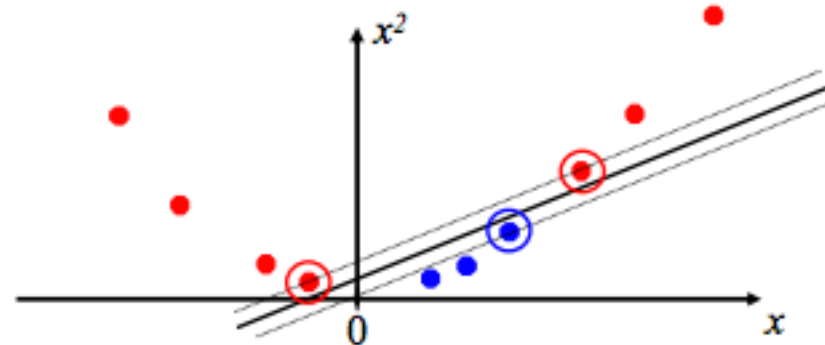
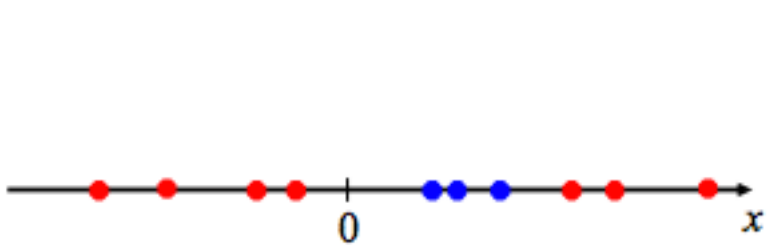


Example: Kernel Map

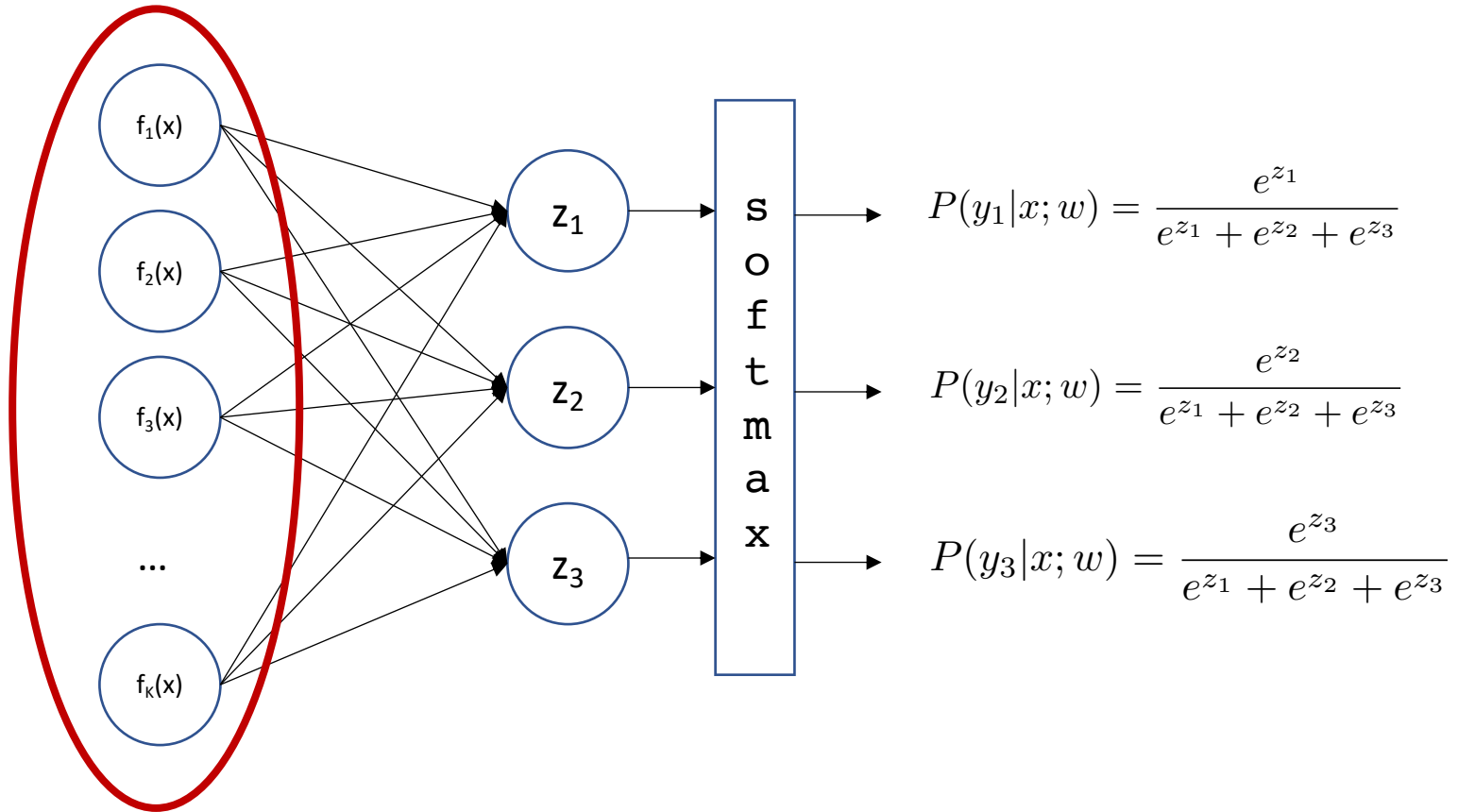


Left) A dataset in \mathbb{R}^2 , not linearly separable. (Right) The same dataset transformed by the transformation: $[x_1, x_2] = [x_1, x_2, x_1^2 + x_2^2]$.

δ: (Left) The decision boundary \vec{w} shown to be linear in \mathbb{R}^3 . (Right) The decision boundary \vec{w} , when transformed back to \mathbb{R}^2 , is nonlinear.

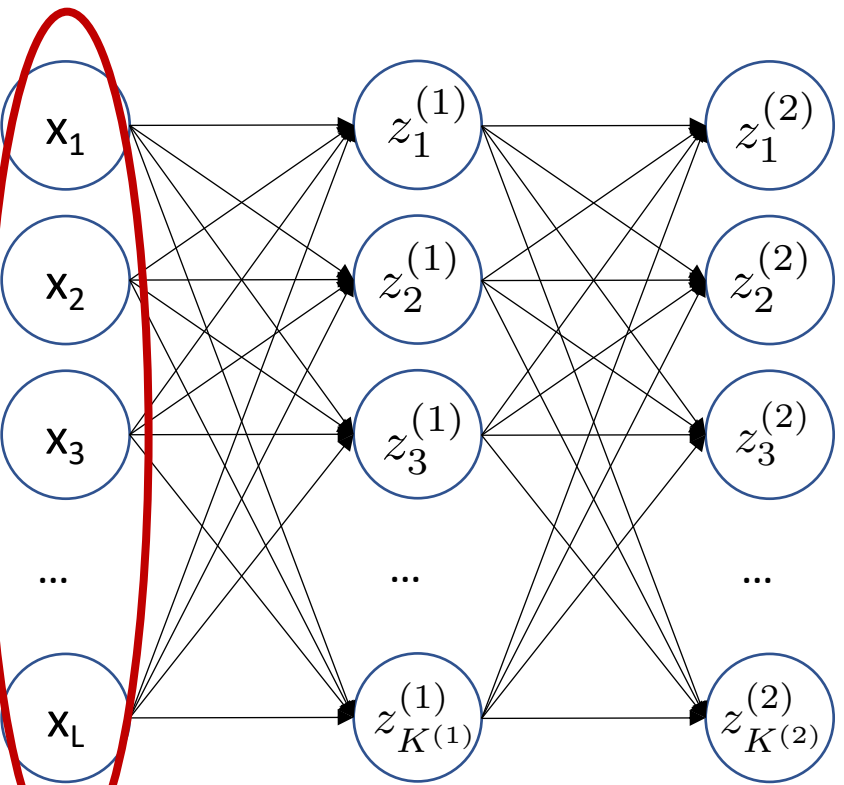


Till now, the features were hand-crafted

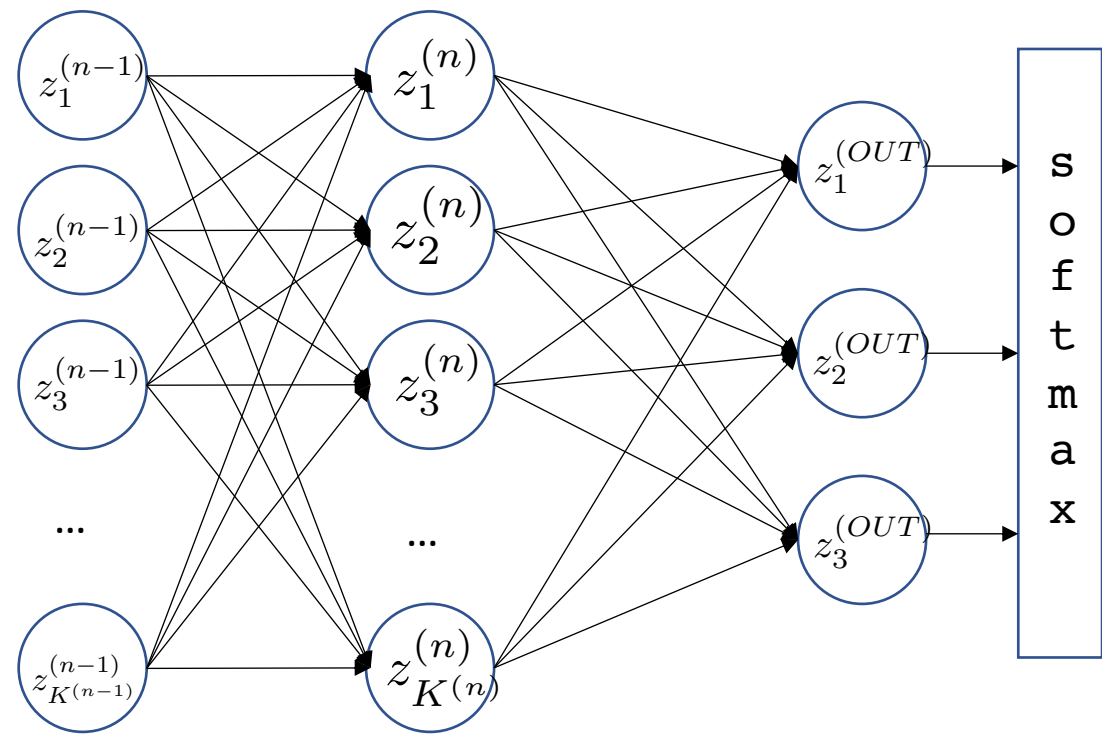


Still, we are designing these features. Can these be acquired in a data-driven manner? Can the parameters controlling these non-linear functions be learned?

Neural networks: learning the features



...



$$P(y_1|x; w) =$$

$$P(y_2|x; w) =$$

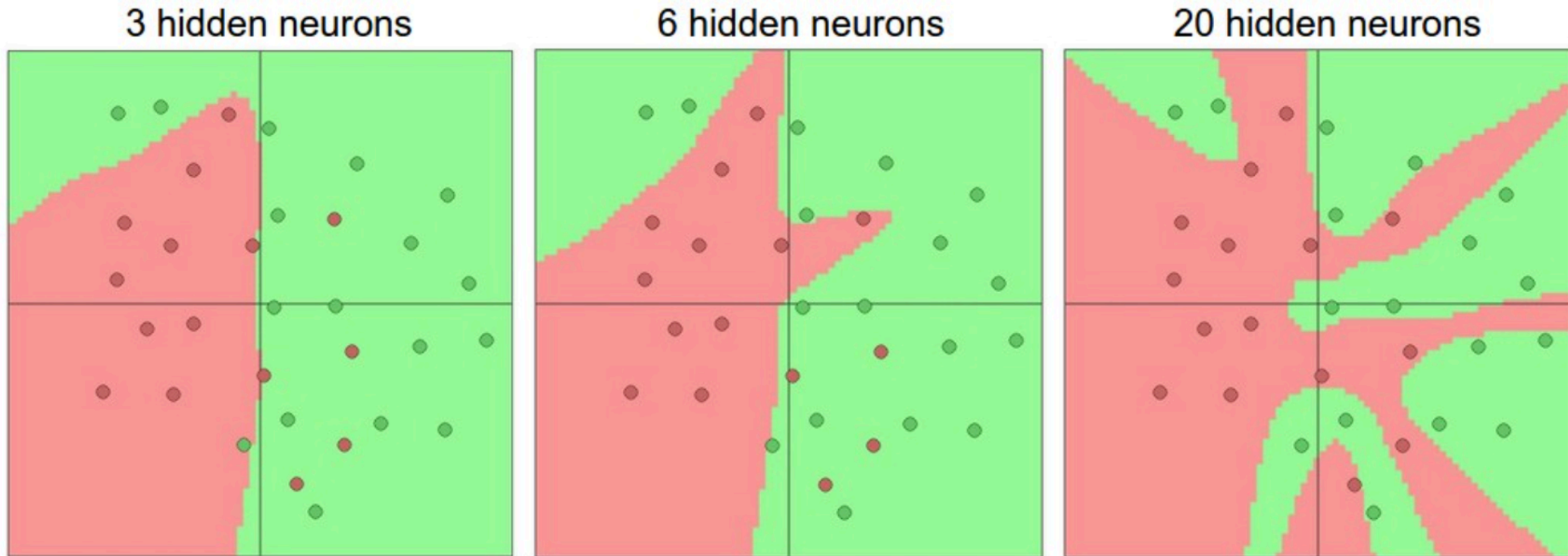
$$P(y_3|x; w) =$$

$$z_i^{(k)} = g\left(\sum_j W_{i,j}^{(k-1,k)} z_j^{(k-1)}\right)$$

g = nonlinear activation function

Structure these models by composing many units. Paradigm is called *deep learning*.

Representation of complex functions

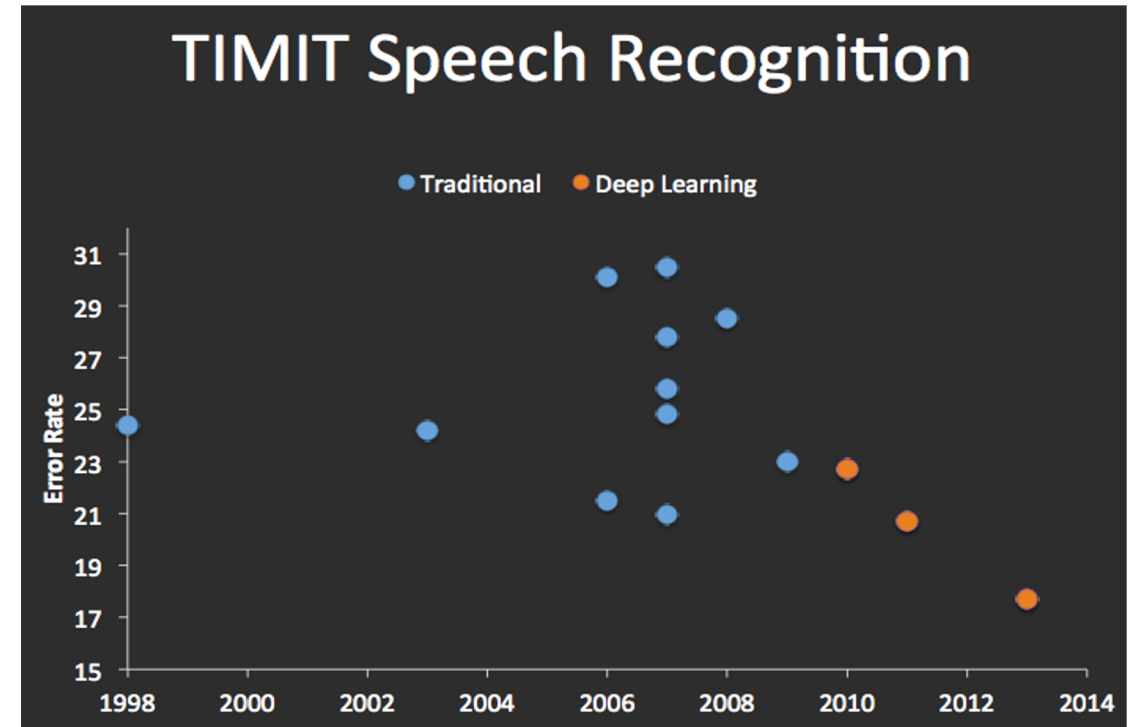
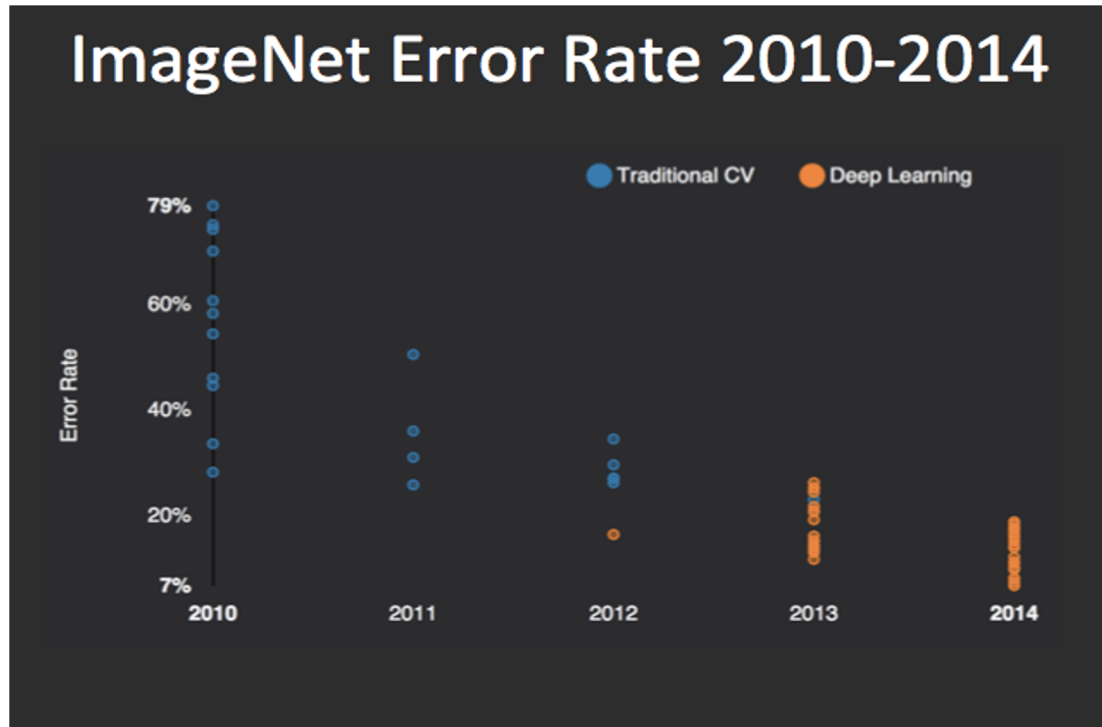


Larger Neural Networks can represent more complicated functions. The data are shown as circles colored by their class, and the decision regions by a trained neural network are shown underneath. You can play with these examples in this [ConvNetsJS demo](#).

Deep Neural Networks

- Last layer
 - Logistic regression
- Several Hidden Layers
 - Computing the features. The features are learned rather than hand-designed.
- Universal function approximation theorem
 - If neural net is large enough
 - Then neural net can represent any continuous mapping from input to output with arbitrary accuracy
 - Note: overfitting is a challenge.
 - In essence, hyper-parametric function approximation.

Neural Networks: Successes



Learning XOR

XOR is not linearly separable.

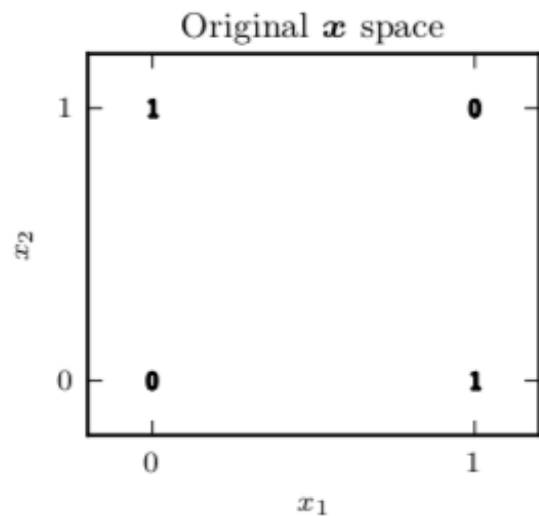


Figure 6.1, left

Rectified Linear Activation

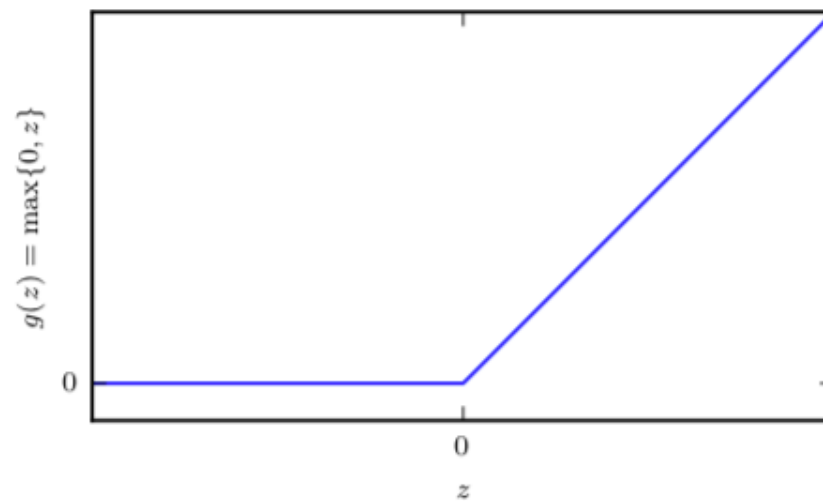


Figure 6.3

Network Diagram

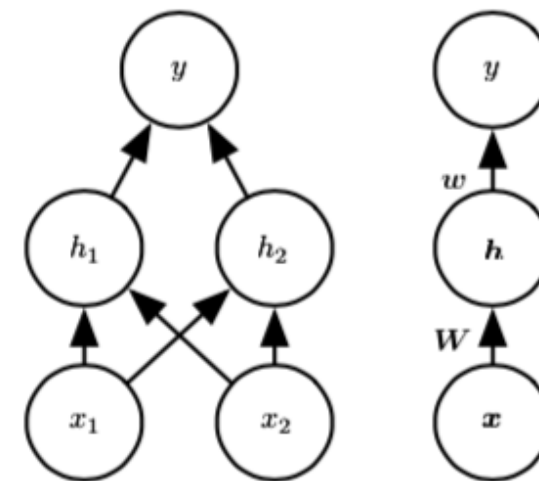


Figure 6.2

More compact representation

Learning XOR

Network Diagram

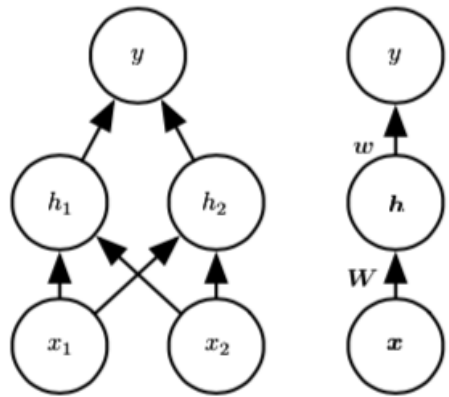


Figure 6.2

Model

$$f(\mathbf{x}; \mathbf{W}, \mathbf{c}, \mathbf{w}, b) = \mathbf{w}^\top \max\{0, \mathbf{W}^\top \mathbf{x} + \mathbf{c}\} + b.$$

$$\mathbf{W} = \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix},$$
$$\mathbf{c} = \begin{bmatrix} 0 \\ -1 \end{bmatrix},$$
$$\mathbf{w} = \begin{bmatrix} 1 \\ -2 \end{bmatrix},$$

XOR is separable in the transformed space

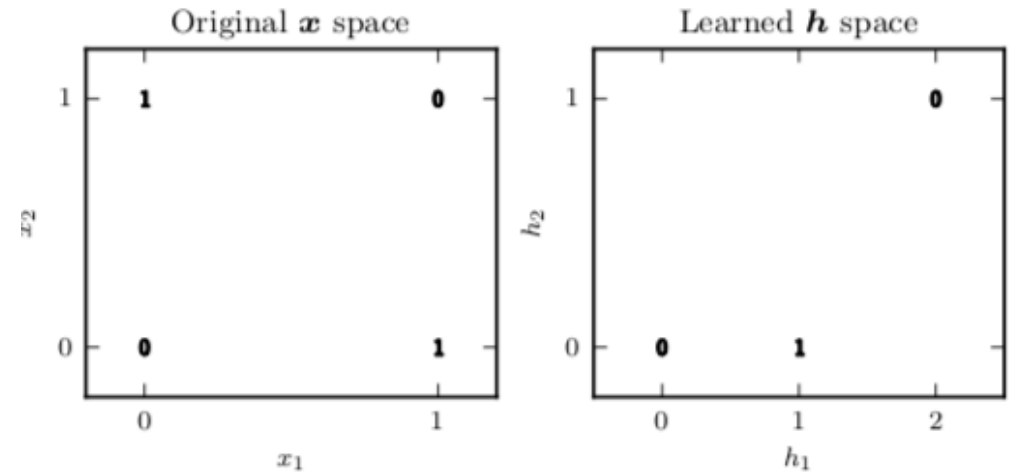
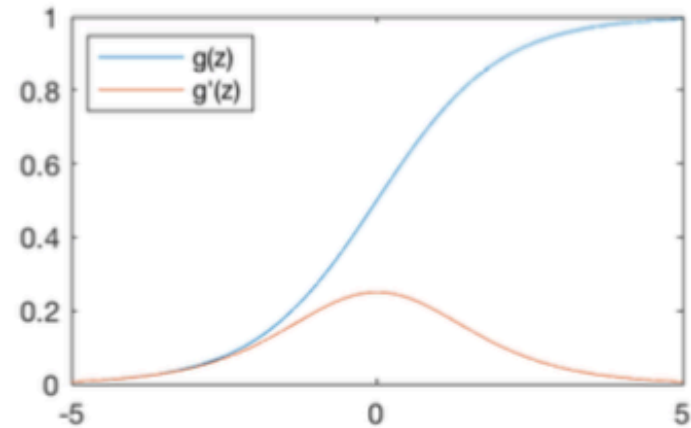


Figure 6.1

Takeaway: Applying ReLU to the output of a linear transformation yields a non-linear transformation. The problem can be solved in the transformed space.

Common Activation Functions

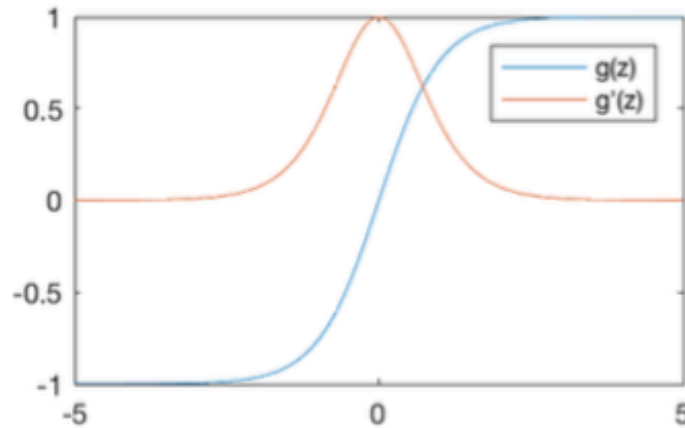
Sigmoid Function



$$g(z) = \frac{1}{1 + e^{-z}}$$

$$g'(z) = g(z)(1 - g(z))$$

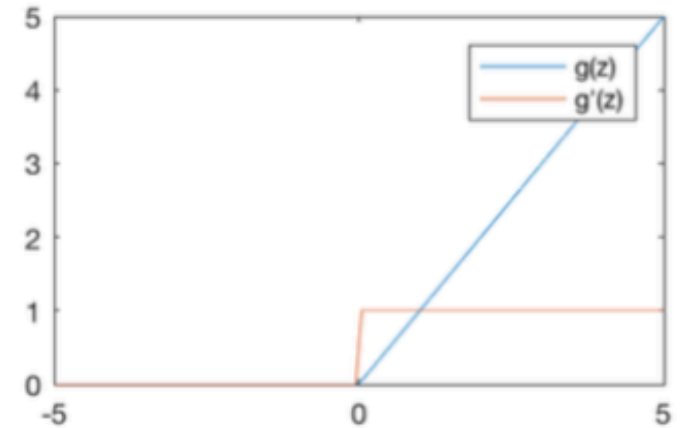
Hyperbolic Tangent



$$g(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$$

$$g'(z) = 1 - g(z)^2$$

Rectified Linear Unit (ReLU)



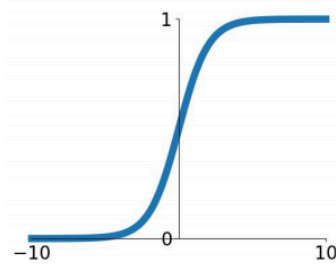
$$g(z) = \max(0, z)$$

$$g'(z) = \begin{cases} 1, & z > 0 \\ 0, & \text{otherwise} \end{cases}$$

Activation Functions

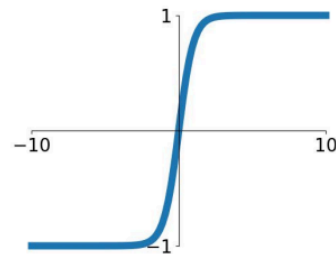
Sigmoid

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



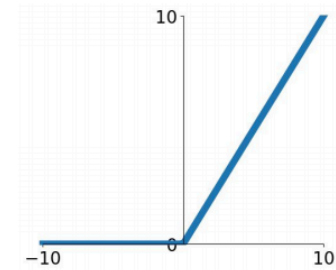
tanh

$$\tanh(x)$$



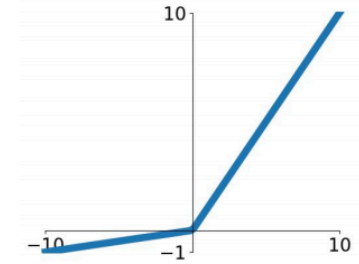
ReLU

$$\max(0, x)$$



Leaky ReLU

$$\max(0.1x, x)$$

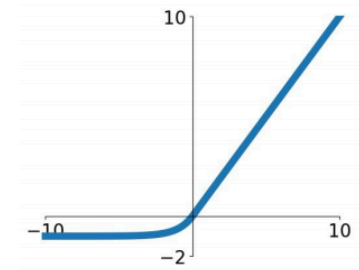


Maxout

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

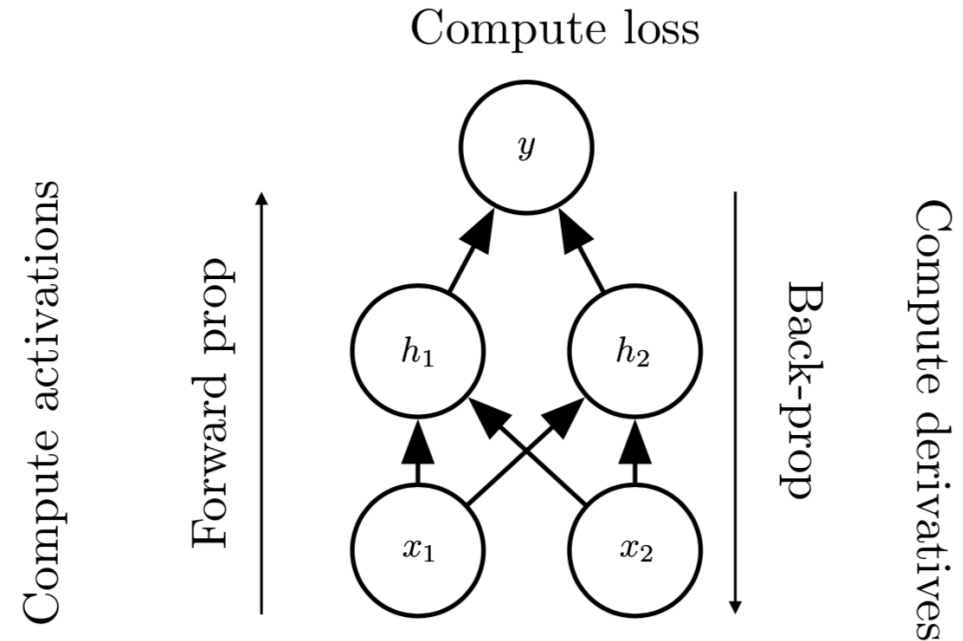
ELU

$$\begin{cases} x & x \geq 0 \\ \alpha(e^x - 1) & x < 0 \end{cases}$$



Backpropagation and Computation Graphs

- Backpropagation
 - In a NN, need a way to optimize the output loss with respect to the inputs.
 - Apply the chain rule to obtain the gradient.



$$\frac{dz}{dx} = \frac{dz}{dy} \frac{dy}{dx}$$

$$\nabla_{\mathbf{x}} z = \left(\frac{\partial \mathbf{y}}{\partial \mathbf{x}} \right)^{\top} \nabla_{\mathbf{y}} z,$$

Backpropagation and Computation Graphs

- Computation Graphs
 - A way to organize the computation in a neural network.
 - Also enables identification and caching of repeated sub-expressions.

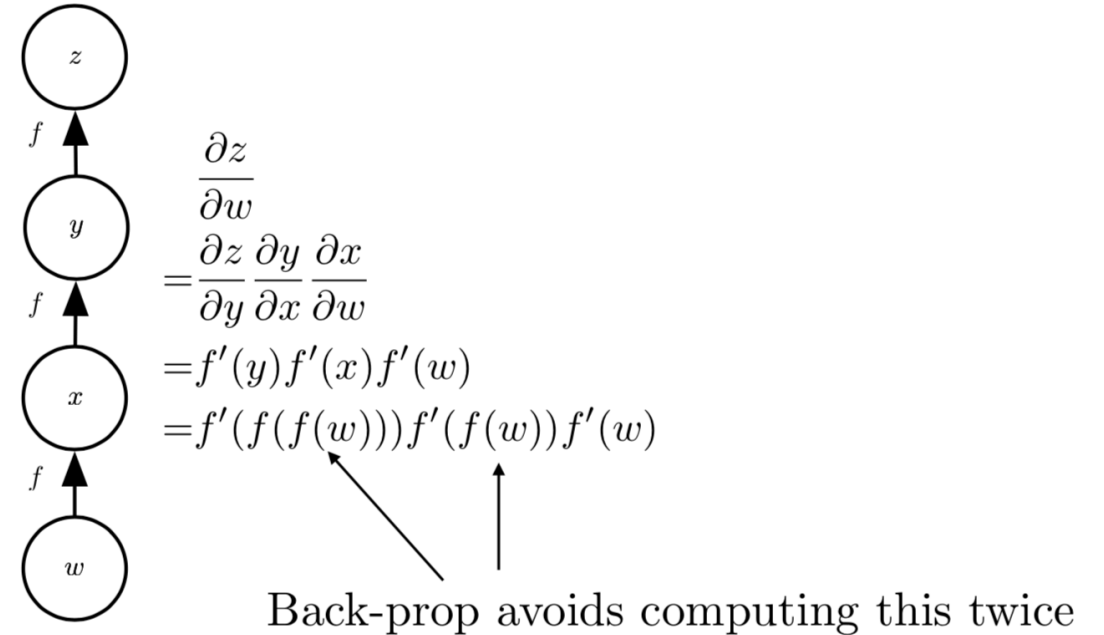


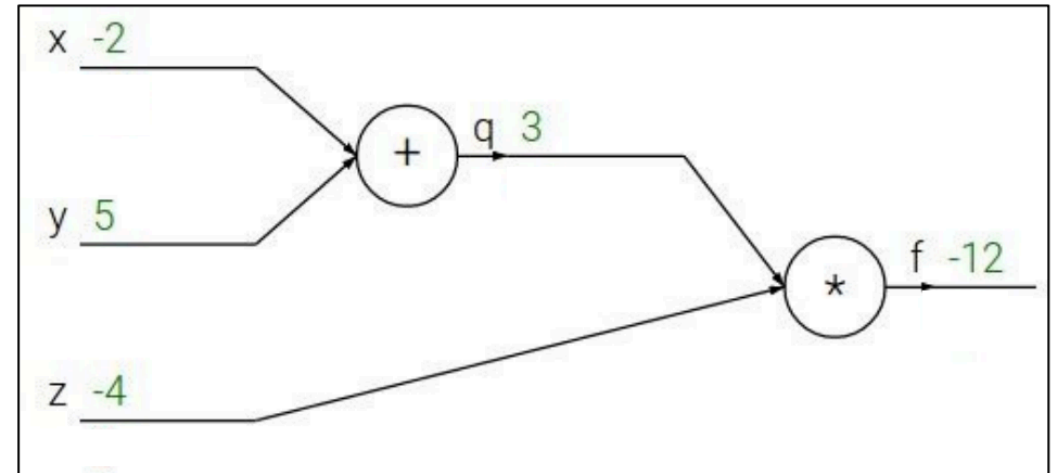
Figure 6.9

Backpropagation: Toy Example

Backpropagation: a simple example

$$f(x, y, z) = (x + y)z$$

e.g. $x = -2, y = 5, z = -4$



Backpropagation: Toy Example

Backpropagation: a simple example

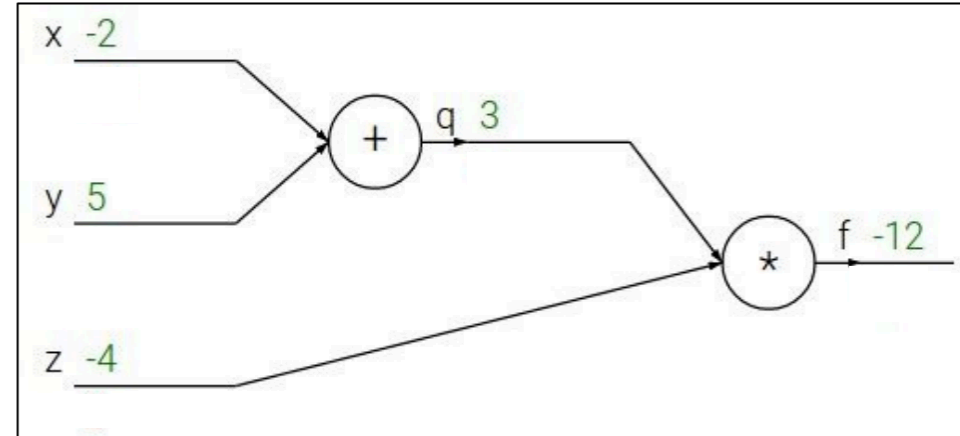
$$f(x, y, z) = (x + y)z$$

e.g. $x = -2, y = 5, z = -4$

$$q = x + y \quad \frac{\partial q}{\partial x} = 1, \frac{\partial q}{\partial y} = 1$$

$$f = qz \quad \frac{\partial f}{\partial q} = z, \frac{\partial f}{\partial z} = q$$

Want: $\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$



Backpropagation: Toy Example

Backpropagation: a simple example

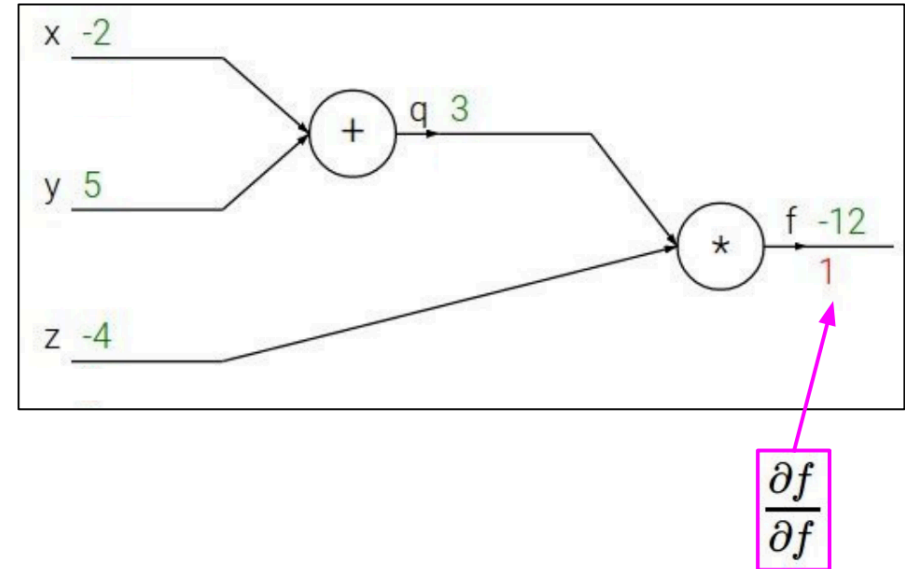
$$f(x, y, z) = (x + y)z$$

e.g. $x = -2, y = 5, z = -4$

$$q = x + y \quad \frac{\partial q}{\partial x} = 1, \frac{\partial q}{\partial y} = 1$$

$$f = qz \quad \frac{\partial f}{\partial q} = z, \frac{\partial f}{\partial z} = q$$

Want: $\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$



Backpropagation: Toy Example

Backpropagation: a simple example

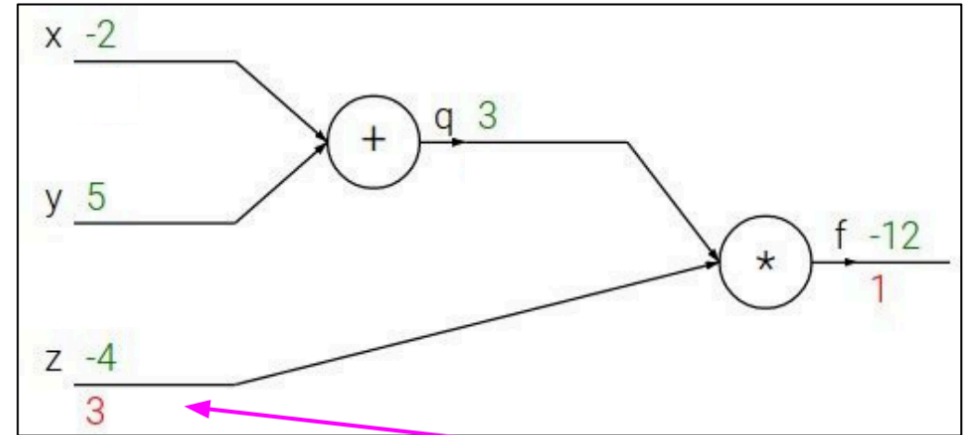
$$f(x, y, z) = (x + y)z$$

e.g. $x = -2, y = 5, z = -4$

$$q = x + y \quad \frac{\partial q}{\partial x} = 1, \frac{\partial q}{\partial y} = 1$$

$$f = qz \quad \frac{\partial f}{\partial q} = z, \frac{\partial f}{\partial z} = q$$

Want: $\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$



$$\frac{\partial f}{\partial z}$$

Backpropagation: Toy Example

Backpropagation: a simple example

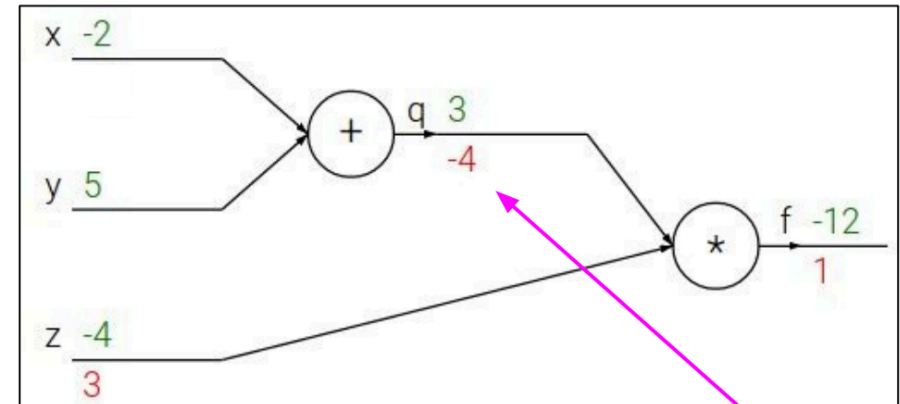
$$f(x, y, z) = (x + y)z$$

e.g. $x = -2, y = 5, z = -4$

$$q = x + y \quad \frac{\partial q}{\partial x} = 1, \frac{\partial q}{\partial y} = 1$$

$$f = qz \quad \frac{\partial f}{\partial q} = z, \frac{\partial f}{\partial z} = q$$

Want: $\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$



$$\frac{\partial f}{\partial q}$$

Backpropagation: Toy Example

Backpropagation: a simple example

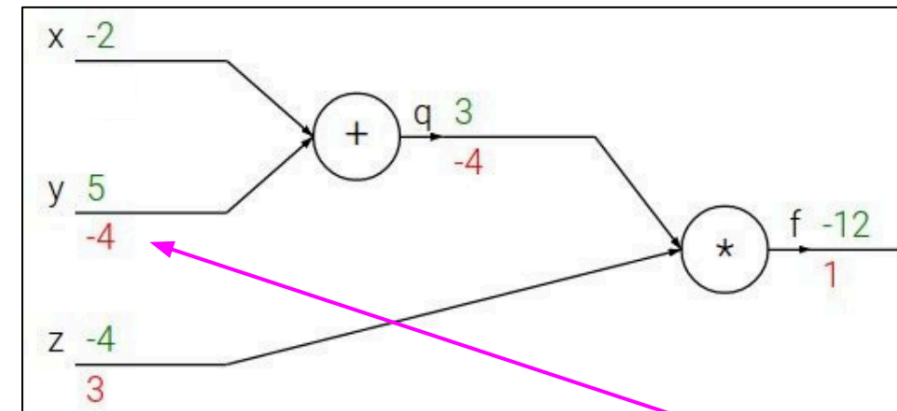
$$f(x, y, z) = (x + y)z$$

e.g. $x = -2, y = 5, z = -4$

$$q = x + y \quad \frac{\partial q}{\partial x} = 1, \frac{\partial q}{\partial y} = 1$$

$$f = qz \quad \frac{\partial f}{\partial q} = z, \frac{\partial f}{\partial z} = q$$

Want: $\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$



Chain rule:

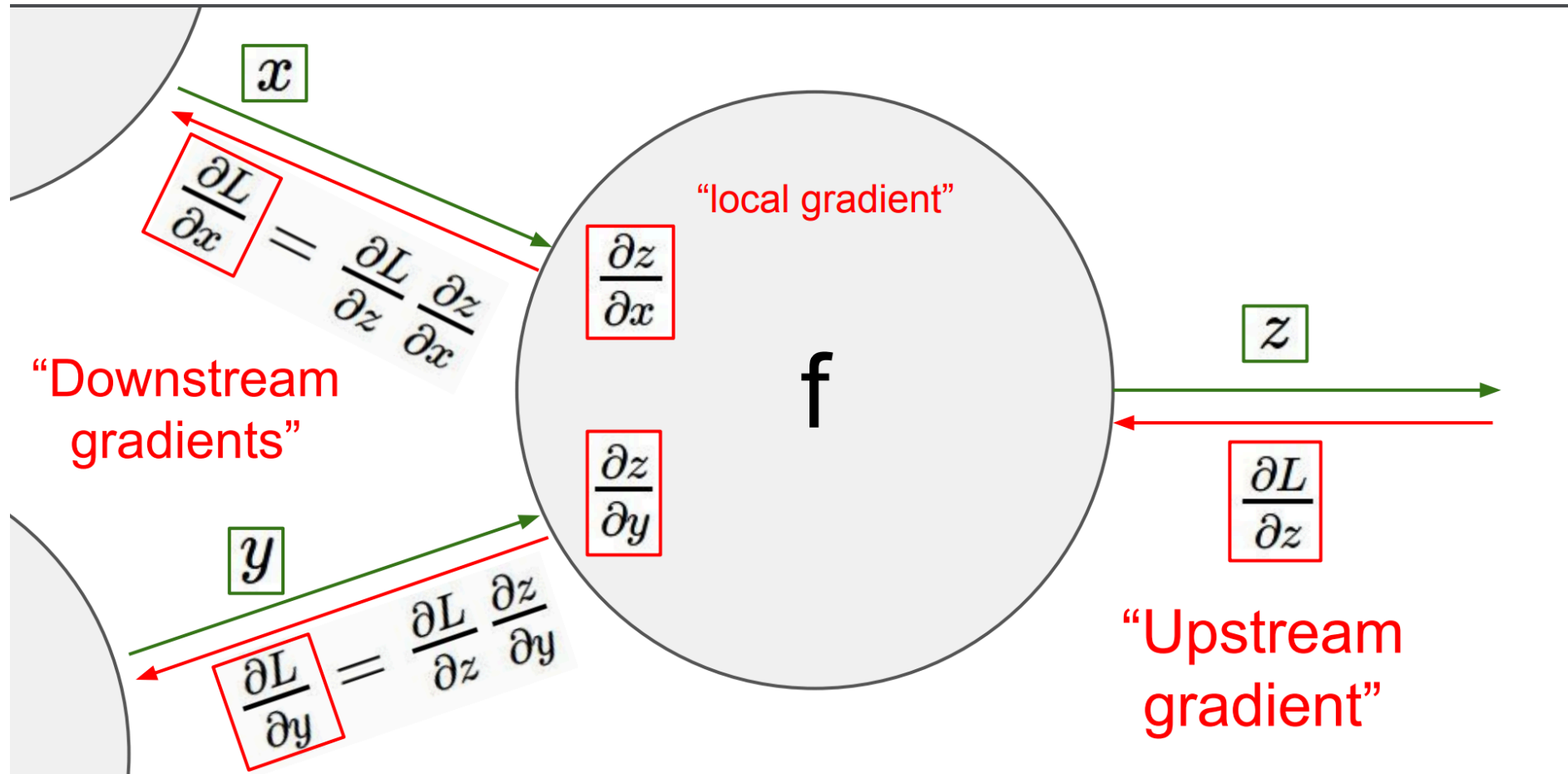
$$\frac{\partial f}{\partial y} = \frac{\partial f}{\partial q} \frac{\partial q}{\partial y}$$

Upstream
gradient

Local
gradient

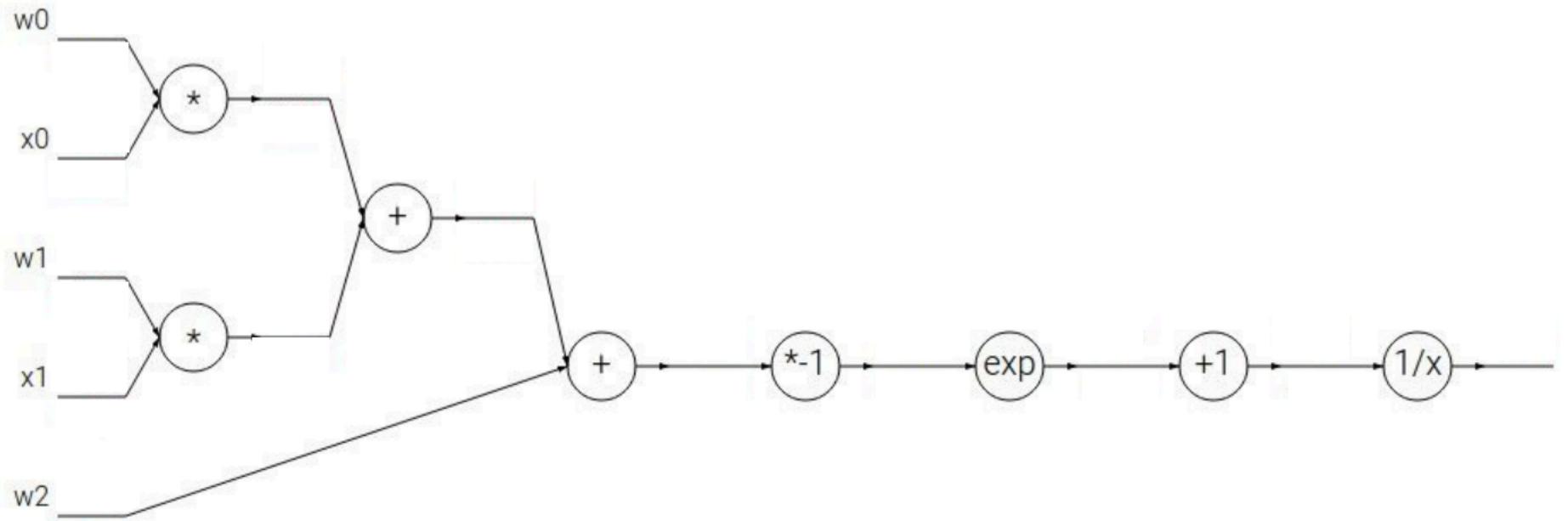
$$\frac{\partial f}{\partial y}$$

Backpropagation



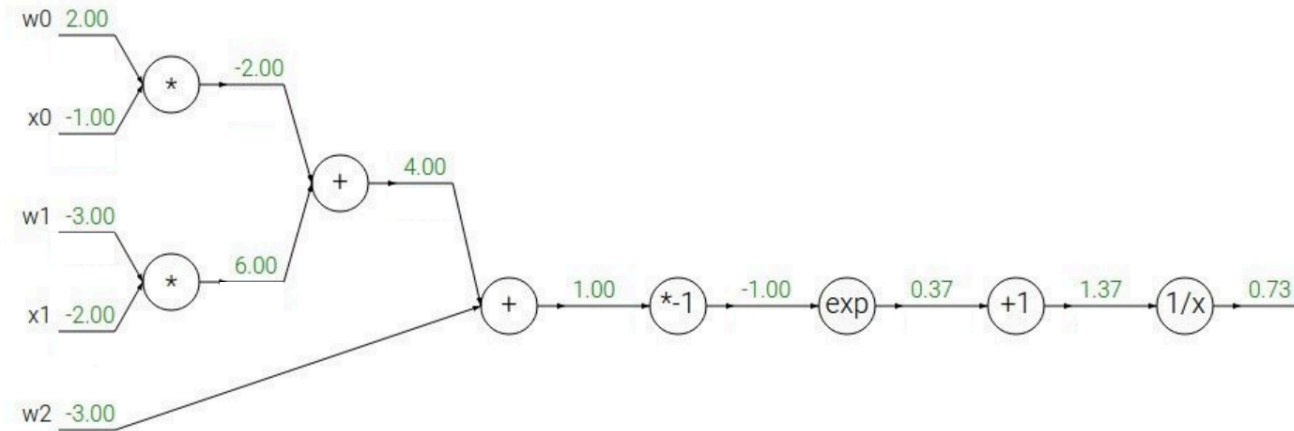
Backpropagation: Example

Another example: $f(w, x) = \frac{1}{1 + e^{-(w_0x_0 + w_1x_1 + w_2)}}$



Backpropagation: Example

Another example: $f(w, x) = \frac{1}{1 + e^{-(w_0x_0 + w_1x_1 + w_2)}}$

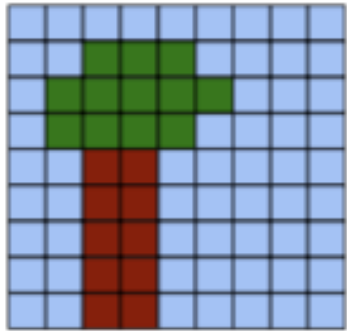


$f(x) = e^x$	→	$\frac{df}{dx} = e^x$		$f(x) = \frac{1}{x}$	→	$\frac{df}{dx} = -1/x^2$
$f_a(x) = ax$	→	$\frac{df}{dx} = a$		$f_c(x) = c + x$	→	$\frac{df}{dx} = 1$

Other Links

- Visualization
 - <http://playground.tensorflow.org>
- Libraries
 - <https://pytorch.org/>
 - <https://www.tensorflow.org/>
 - <https://pypi.org/project/Theano/>
 - <http://pyro.ai/>

Locality and Translational Invariance



A digital image is a 2D grid of pixels.

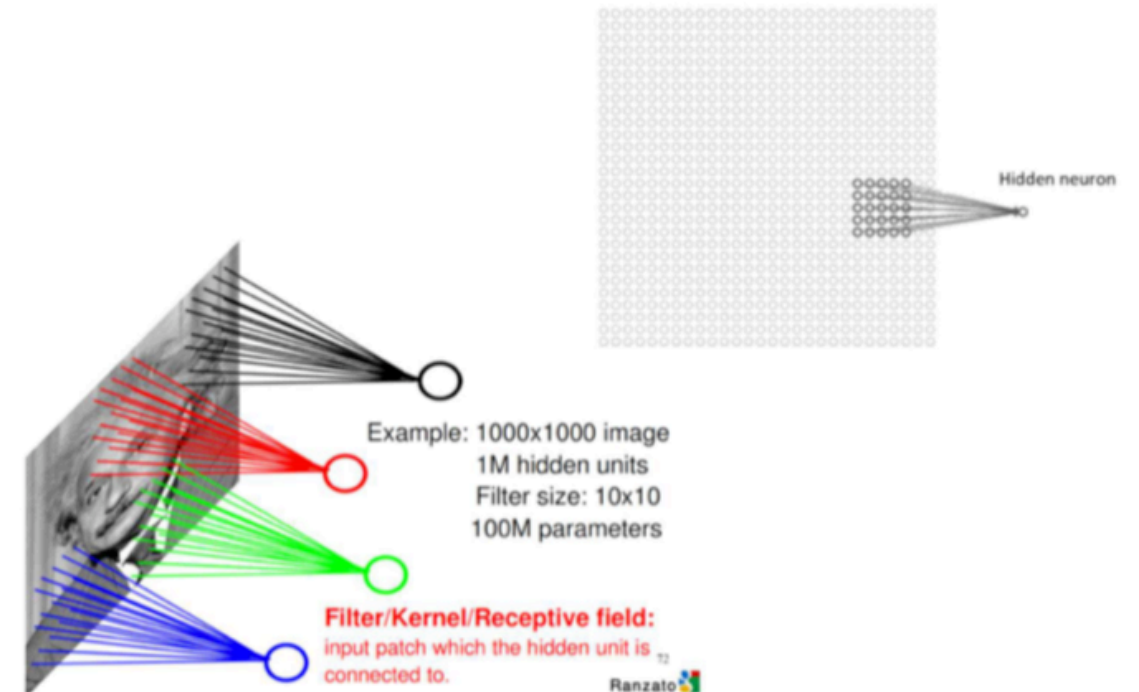
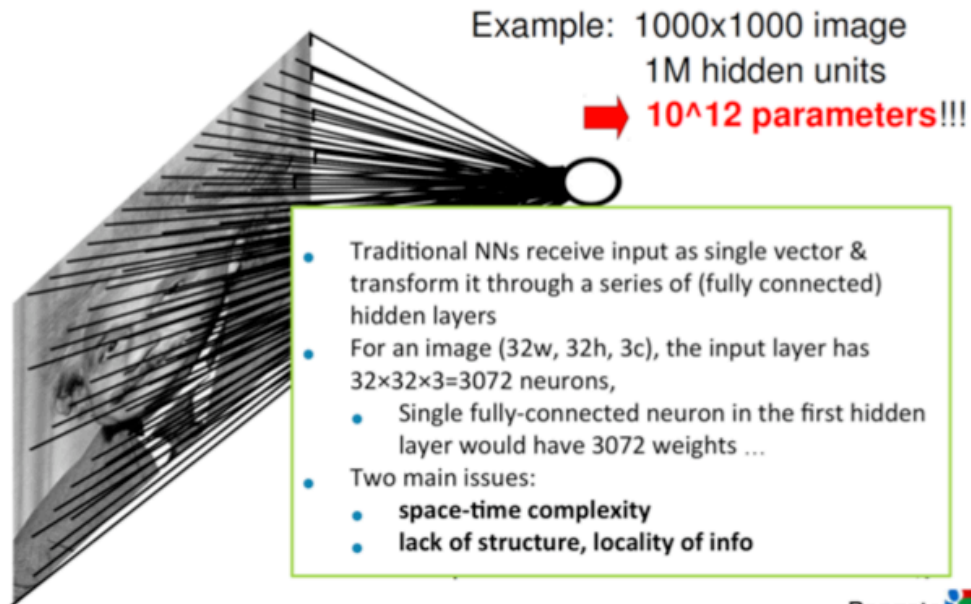
A neural network expects a **vector of numbers** as input.



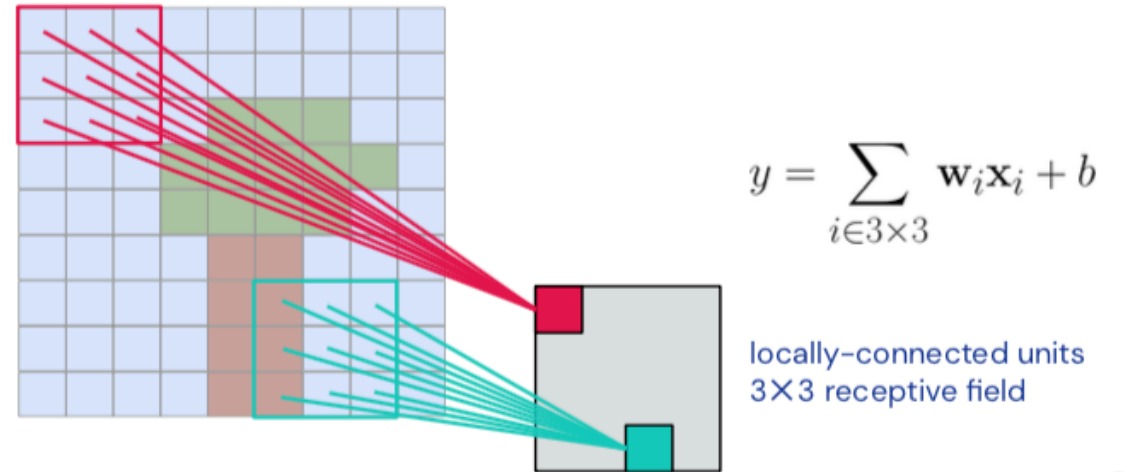
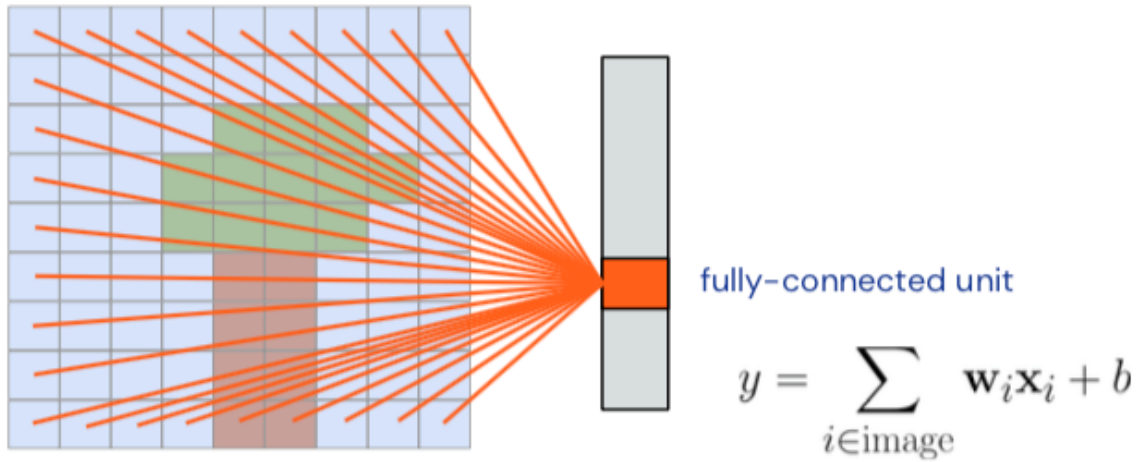
Locality of Information: Receptive Fields

Fully connected network.

Convolutional NN

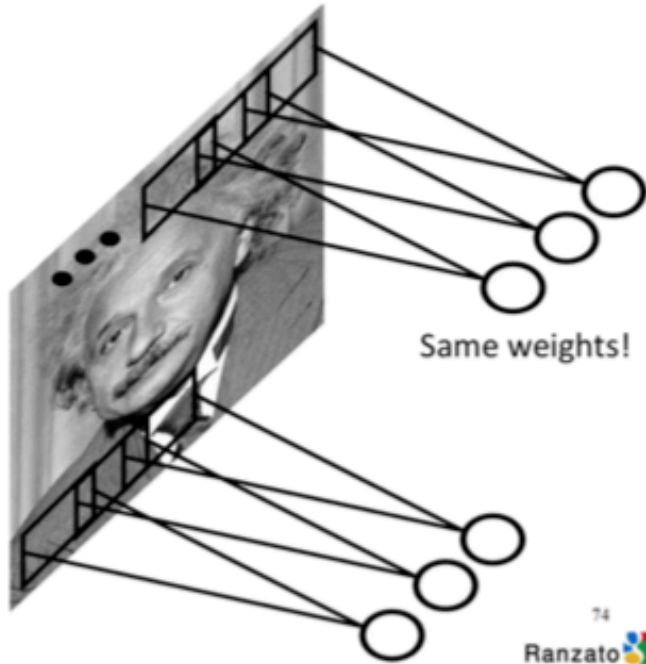


From Globally to Locally Connected



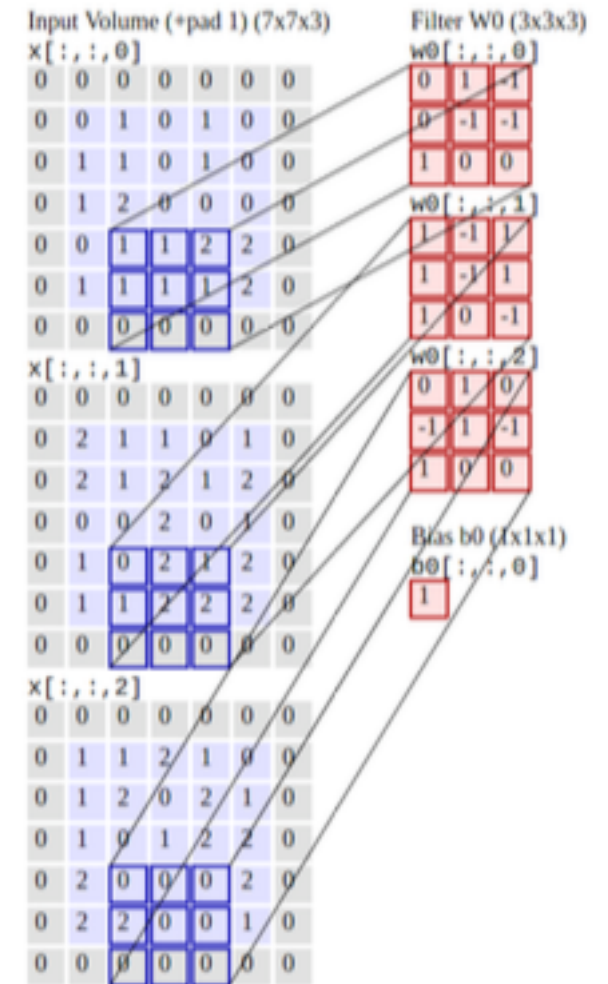
Convolutional NN

Feature Maps



74
Ranzato

- The map from the input layer to the hidden layer is therefore a feature map: all nodes detect the same feature in different parts
- The map is defined by the shared weights and bias
- The shared map is the result of the application of a convolutional filter (defined by weights and bias), also known as convolution with learned kernels

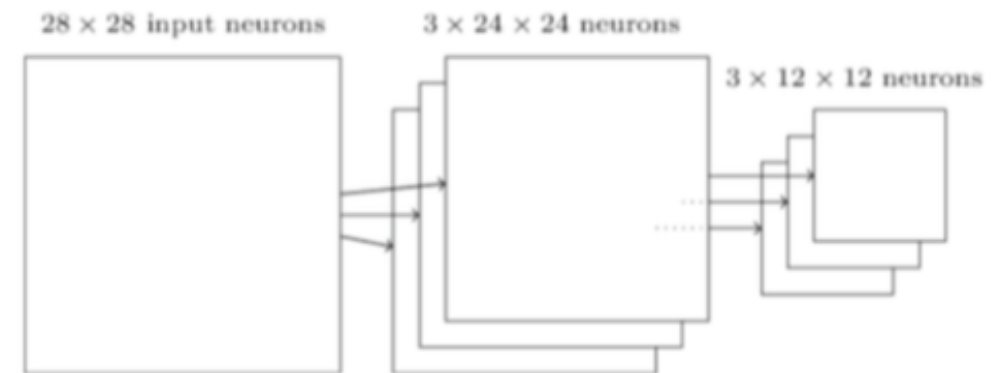
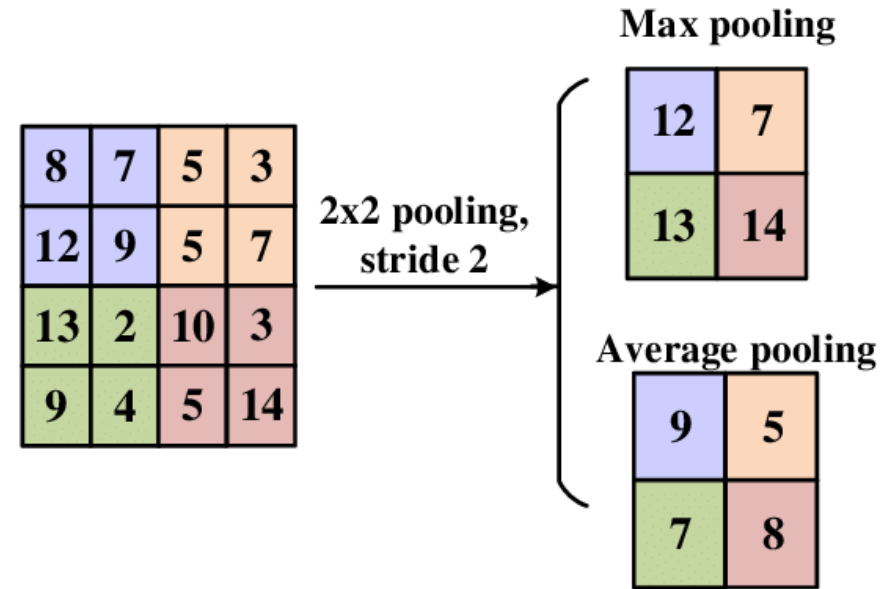


Pooling layers

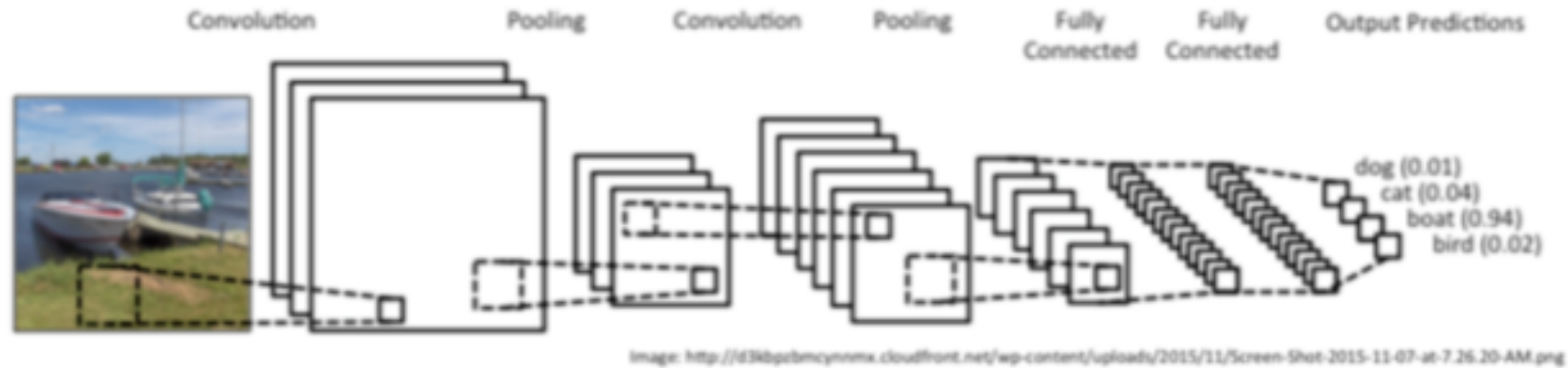
Pooling layers are usually used immediately after convolutional layers.

Pooling layers simplify / subsample / compress the information in the output from convolutional layer

A pooling layer takes each feature map output from the convolutional layer and prepares a condensed feature map



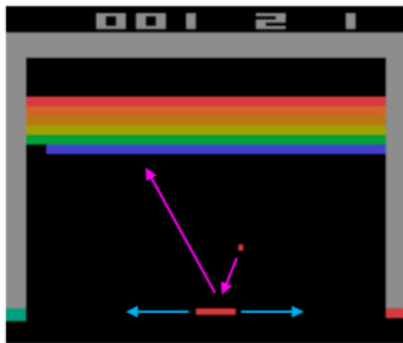
Convolution NN



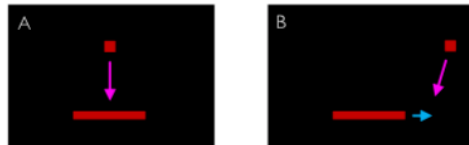
- Consider local structure and common extraction of features
- Not fully connected. Locality of processing
- Weight sharing for parameter reduction
- Learn the parameters of multiple convolutional filter banks
- Compress to extract salient features & favor generalization

Application

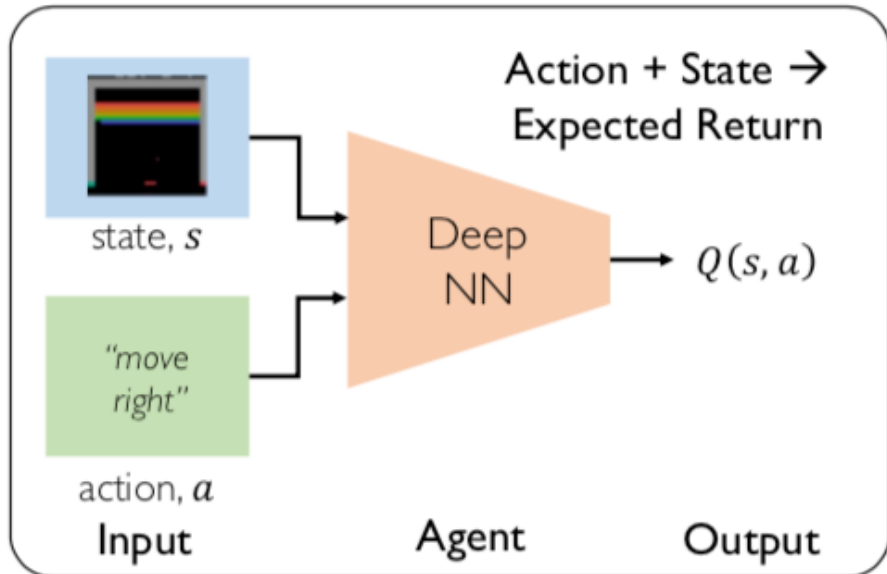
Example: Atari Breakout



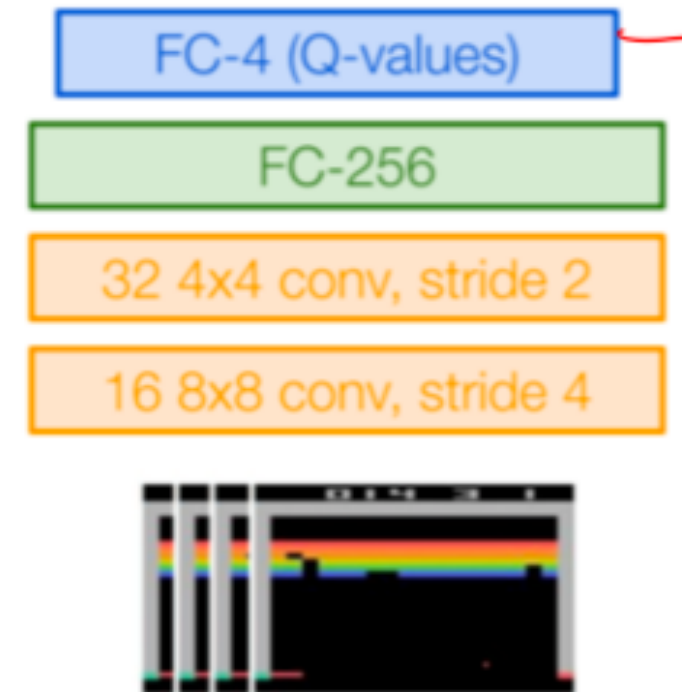
It can be very difficult for humans to accurately estimate Q-values



Which (s, a) pair has a higher Q-value? 🤔

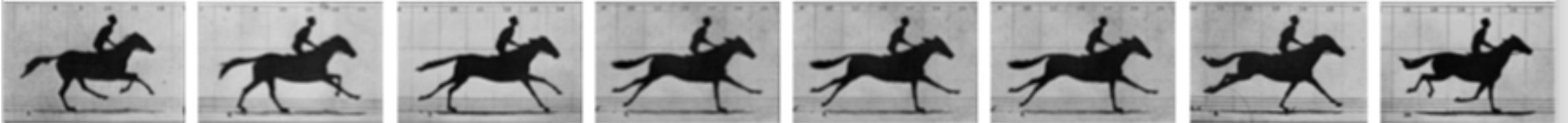


Output: $Q(s, \text{left}), Q(s, \text{right}), \dots$



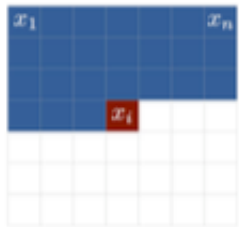
Current state s_t : 84*84*4 stack of last four frames. After RGB \rightarrow grayscale conversion, downsampling and cropping.

Sequences

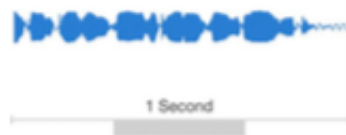


"Sequences really seem to be everywhere! We should learn how to model them. What is the best way to do that? Stay tuned!"

Words, letters



Images



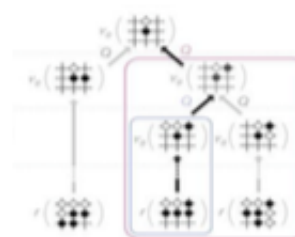
Speech



Videos



Programs



Decision making

Collection of elements where elements can be repeated, order matters and can be of variable or infinite length.

Sequences

	Supervised learning	Sequence modelling
Data	$\{x, y\}_i$	$\{x\}_i$
Model	$y \approx f_\theta(x)$	$p(x) \approx f_\theta(x)$
Loss	$\mathcal{L}(\theta) = \sum_{i=1}^N l(f_\theta(x_i), y_i)$	$\mathcal{L}(\theta) = \sum_{i=1}^N \log p(f_\theta(x_i))$
Optimisation	$\theta^* = \arg \min_{\theta} \mathcal{L}(\theta)$	$\theta^* = \arg \max_{\theta} \mathcal{L}(\theta)$

Modeling the conditional distribution

--

The chain rule

Computing the joint $p(\mathbf{x})$ from conditionals

$$p(\mathbf{x}) = \prod_{t=1}^T p(x_t | x_1, \dots, x_{t-1})$$

Modeling

Modeling **word**

Modeling word **probabilities**

Modeling word probabilities **is**

Modeling word probabilities is **really**

Modeling word probabilities is really **difficult**

$$p(x_1)$$

$$p(x_2 | x_1)$$

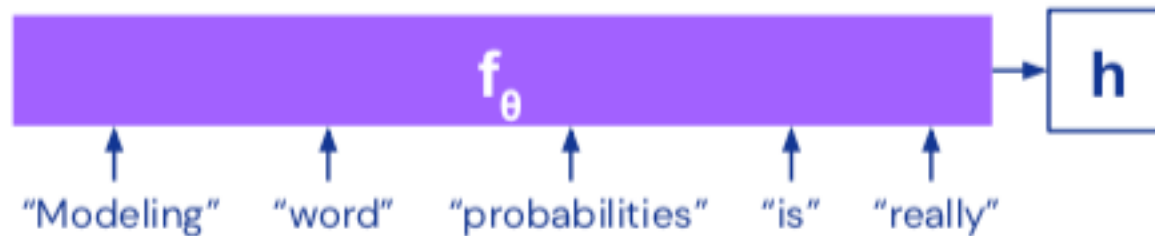
$$p(x_3 | x_2, x_1)$$

$$p(x_4 | x_3, x_2, x_1)$$

$$p(x_5 | x_4, x_3, x_2, x_1)$$

$$p(x_6 | x_5, x_4, x_3, x_2, x_1)$$

Vectorizing the conditional likelihood

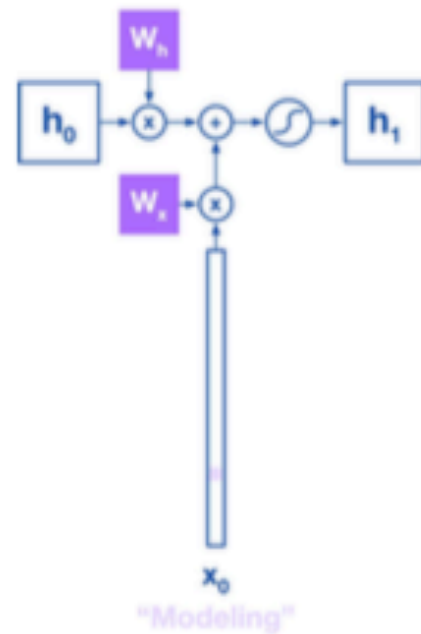


Desirable properties for f_θ :

- Order matters
- Variable length
- Learnable (differentiable)
- Individual changes can have large effects (non-linear/deep)

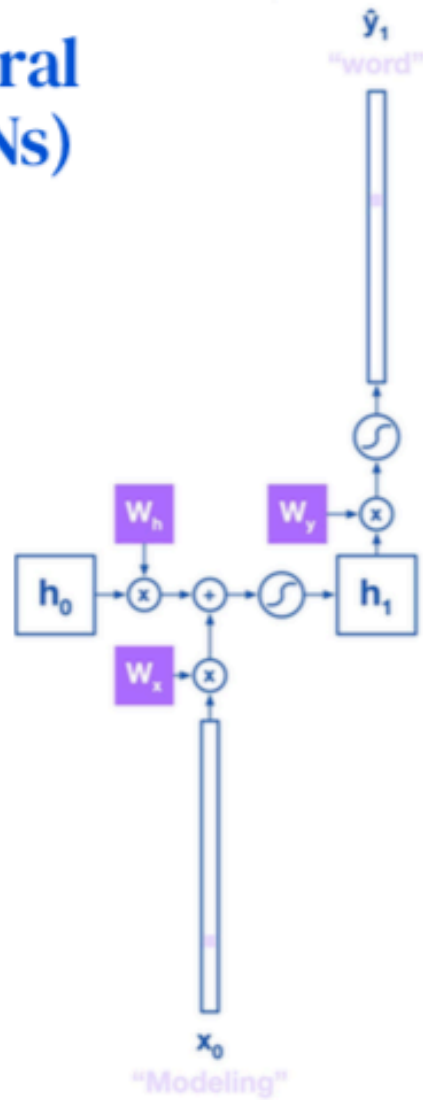
Recurrent Neural Networks (RNNs)

Persistent state variable \mathbf{h} stores information from the context observed so far.



$$\mathbf{h}_t = \tanh(\mathbf{W}_h \mathbf{h}_{t-1} + \mathbf{W}_x \mathbf{x}_t)$$

Recurrent Neural Networks (RNNs)

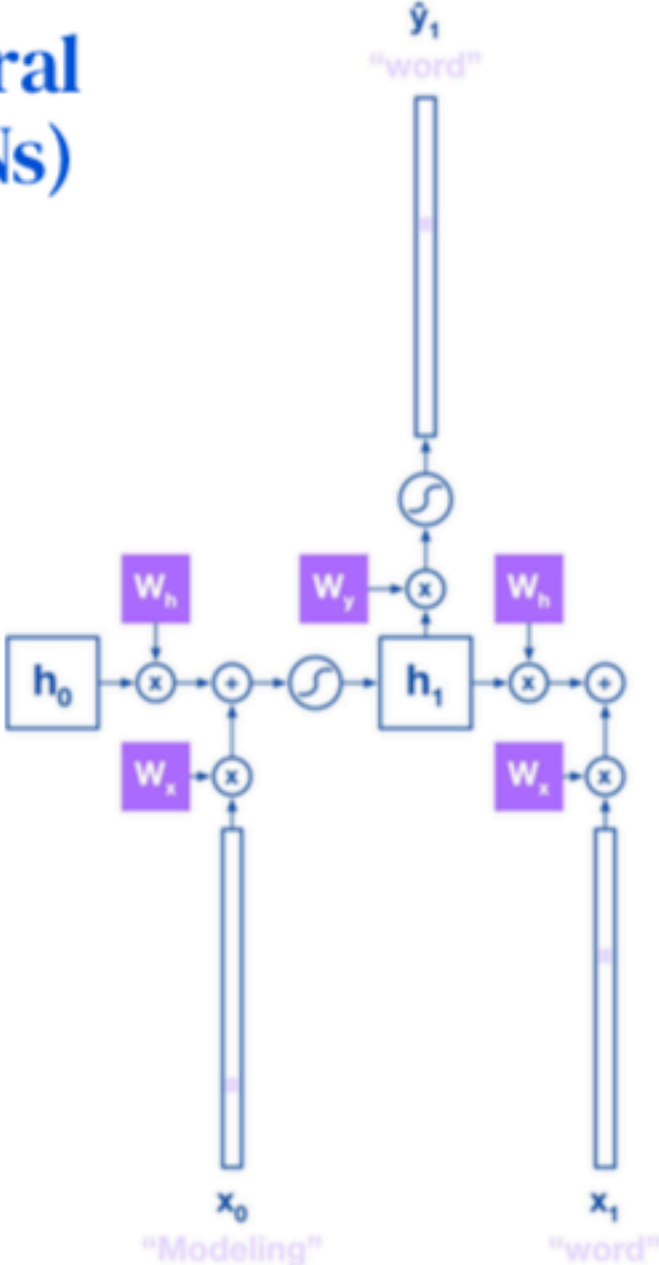


RNNs predict the target \mathbf{y} (the next word) from the state \mathbf{h} .

$$p(\mathbf{y}_{t+1}) = \textit{softmax}(\mathbf{W}_y \mathbf{h}_t)$$

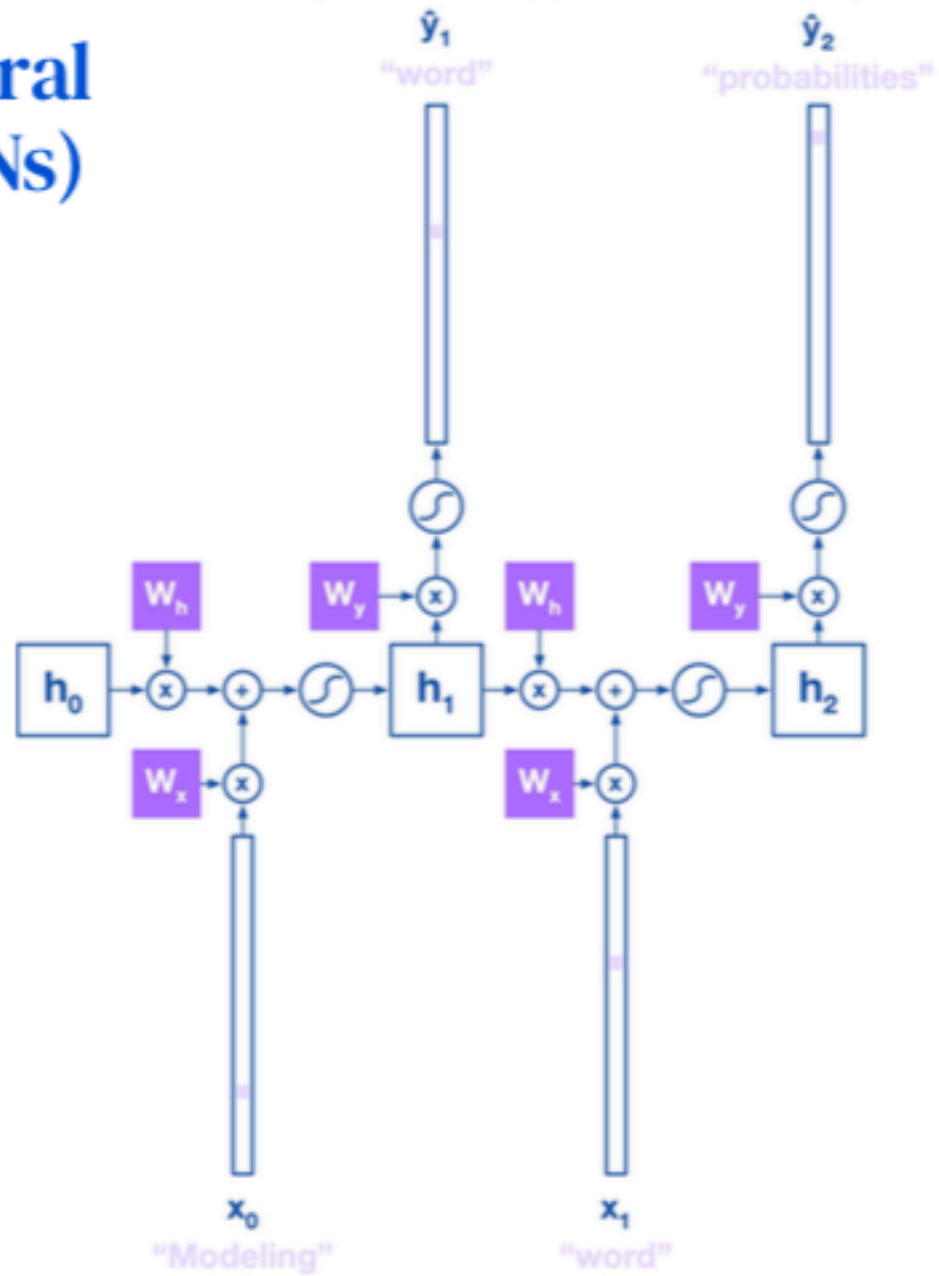
Softmax ensures we obtain a distribution over all possible words.

Recurrent Neural Networks (RNNs)



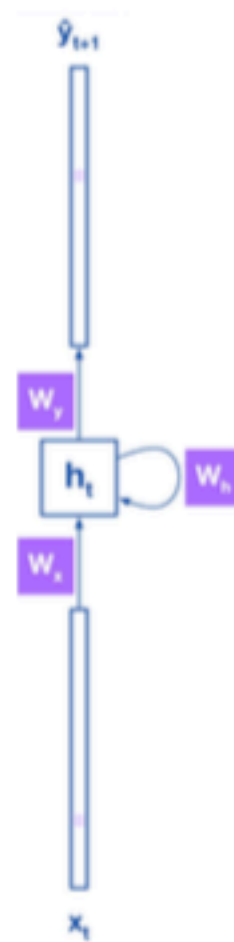
Input next word in sentence x_1

Recurrent Neural Networks (RNNs)

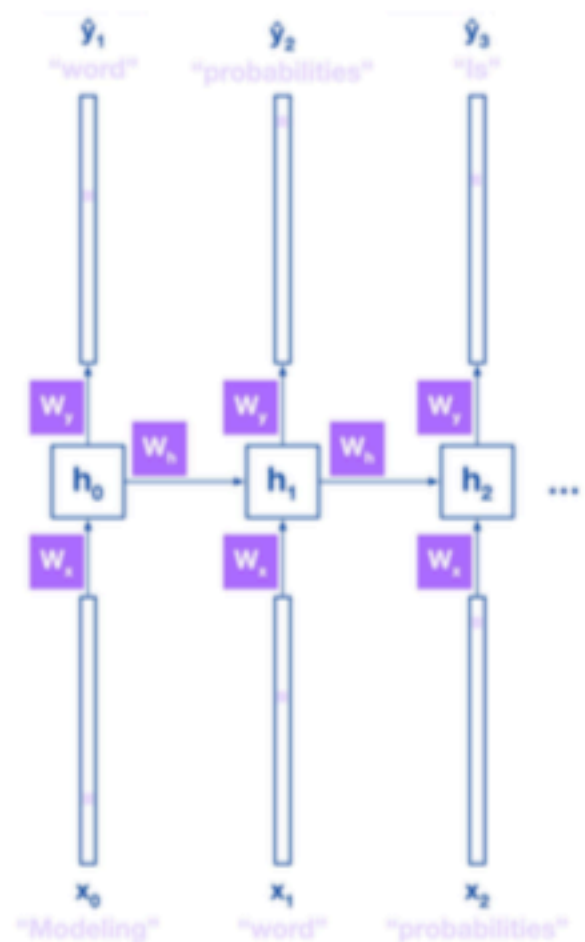


Recurrent Neural Networks (RNNs)

Weights are shared over time steps



RNN



RNN rolled out over time

Loss: Cross Entropy

Next word prediction is essentially a classification task where the number of classes is the size of the vocabulary.

As such we use the cross-entropy loss:

For one word: $\mathcal{L}_{\theta}(\mathbf{y}, \hat{\mathbf{y}})_t = -\mathbf{y}_t \log \hat{\mathbf{y}}_t$

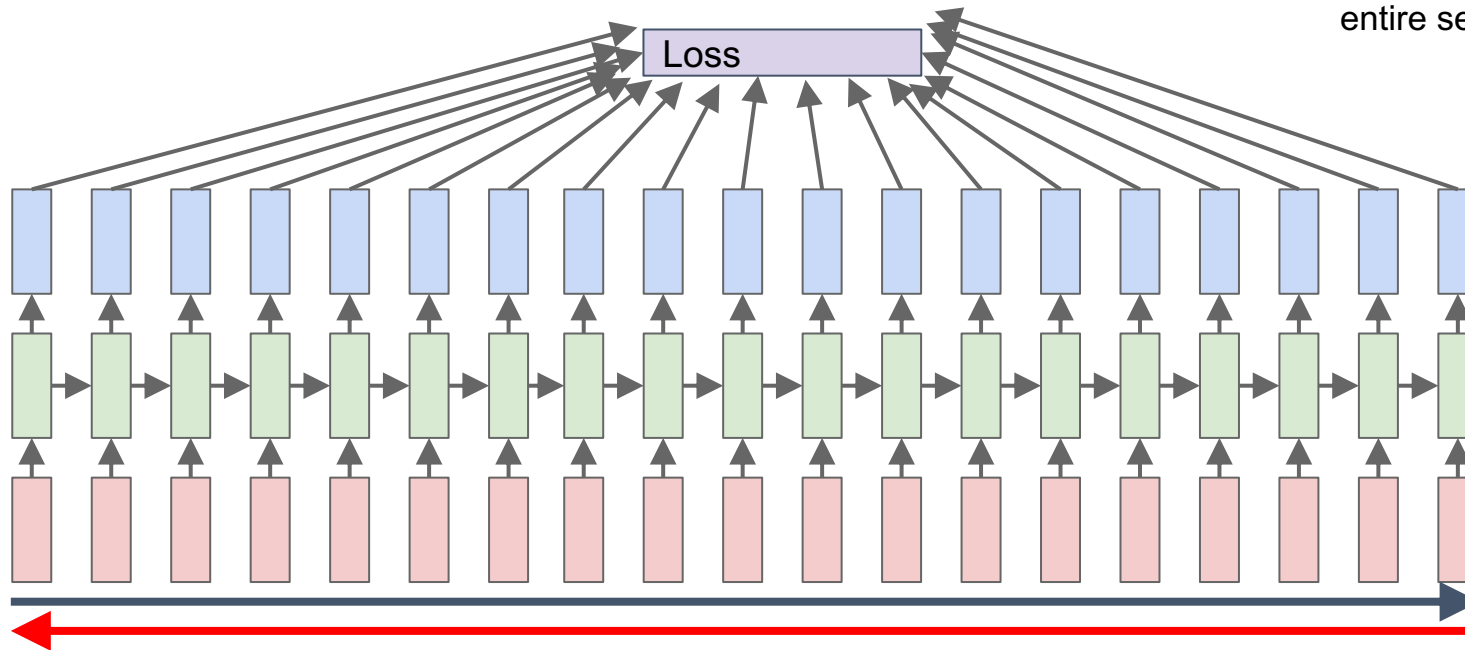
For the sentence: $\mathcal{L}_{\theta}(\mathbf{y}, \hat{\mathbf{y}}) = -\sum_{t=1}^T \mathbf{y}_t \log \hat{\mathbf{y}}_t$

With parameters $\theta = \{\mathbf{W}_y, \mathbf{W}_x, \mathbf{W}_h\}$



Backprop through Time

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

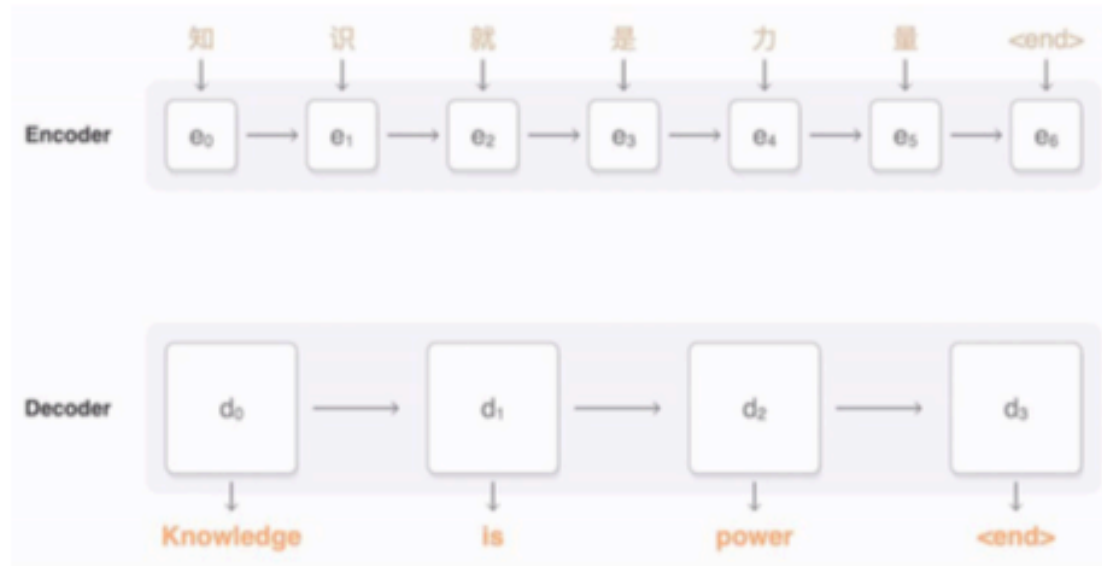


RNNs can have long or short dependencies. When there are long dependencies, gradients have trouble back-propagating through.

Other models such as LSTMs and beyond address that problem.

Applications

Google Neural Machine Translation

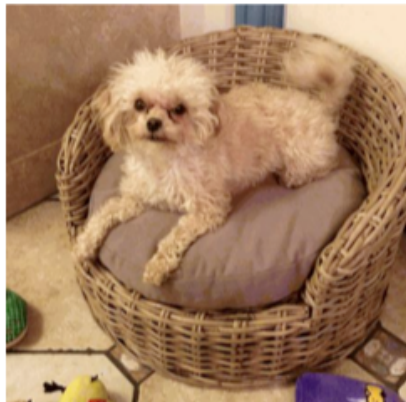
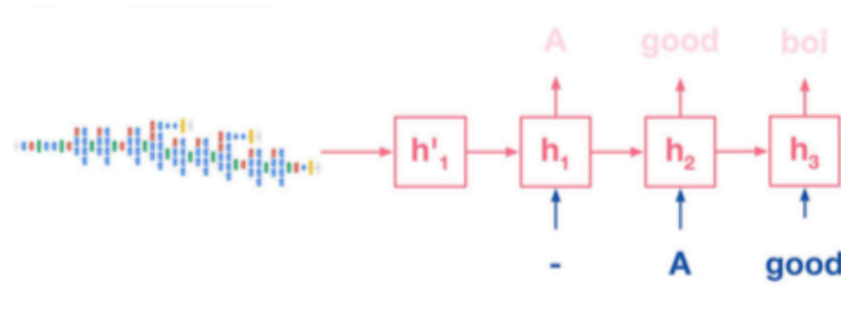


Wu et al, 2016

(Kalchbrenner et al, 2013; Sutskever et al, 2014; Cho et al, 2014; Bhadanau et al, 2014; ...)

Applications

$$p(\text{language}_1 | \text{language}_2) \rightarrow p(\text{language}_1 | \text{image})$$

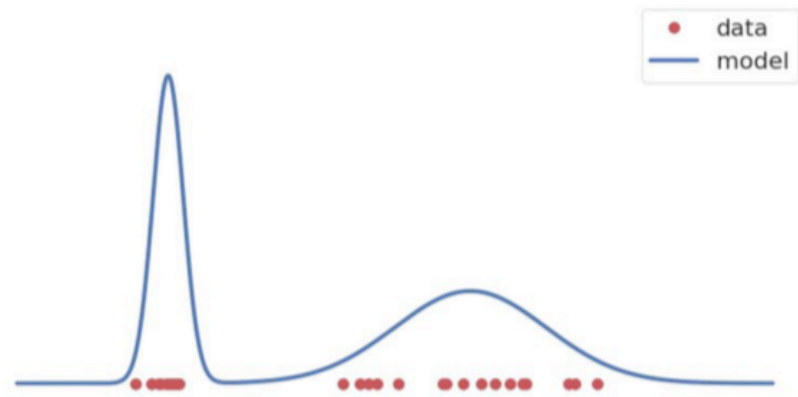


Human: A brown dog
laying in a red wicker
bed.

Best Model: A small dog
is sitting on a chair.

Initial Model: A large
brown dog laying on top
of a couch.

Generative Models



Goal of generative modeling is to learn a model of the true (unknown) underlying data distribution from samples.



Goodfellow, et al. Generative adversarial networks. NIPS (2014)



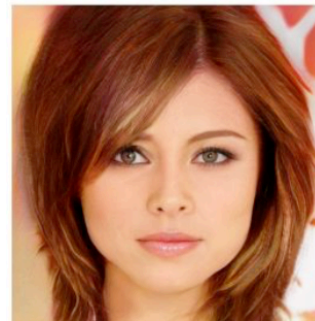
Denton, et al. Deep Generative Image Models using a Laplacian Pyramid of Adversarial Networks. NIPS (2015)



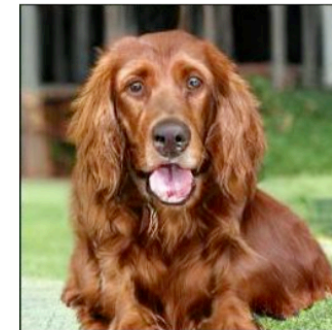
Radford et al. Unsupervised Representation Learning with Deep Convolutional Generative Adversarial Networks. ICLR (2015)



Miyato et al. Spectral normalization for Generative Adversarial Networks. ICLR (2018)



Karras et al. Large Scale GAN Training for High Fidelity Natural Image Synthesis. ICLR (2018)



Brock et al. Large Scale GAN Training for High Fidelity Natural Image Synthesis. ICLR (2019)



Karras et al. A Style-Based Generator Architecture for Generative Adversarial Networks. CVPR (2018)

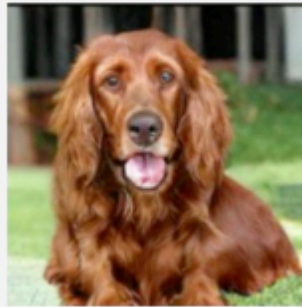
Generative Adversarial Networks

Discriminator

Learns to distinguish between real and generated data.

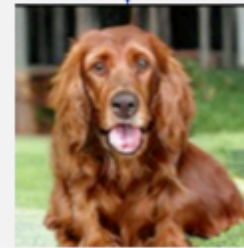


vs



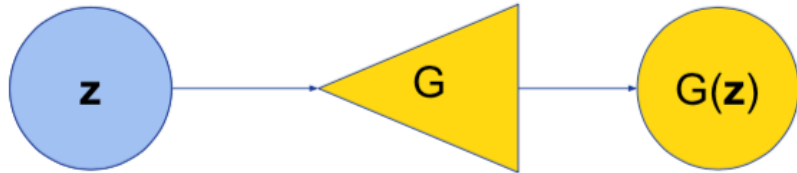
Generator

Learns to generate data to “fool” the discriminator.



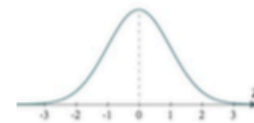
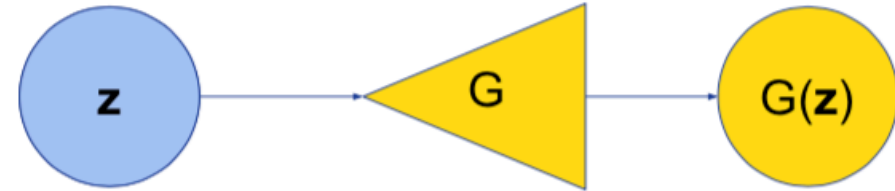
Generative Adversarial Networks

latent ("noise") vector $z \sim P(z)$ generator G: a deep neural network generated data $G(z)$



Generation of image data

latent ("noise") vector $z \sim P(z)$ generator G: a deep neural network generated data $G(z)$

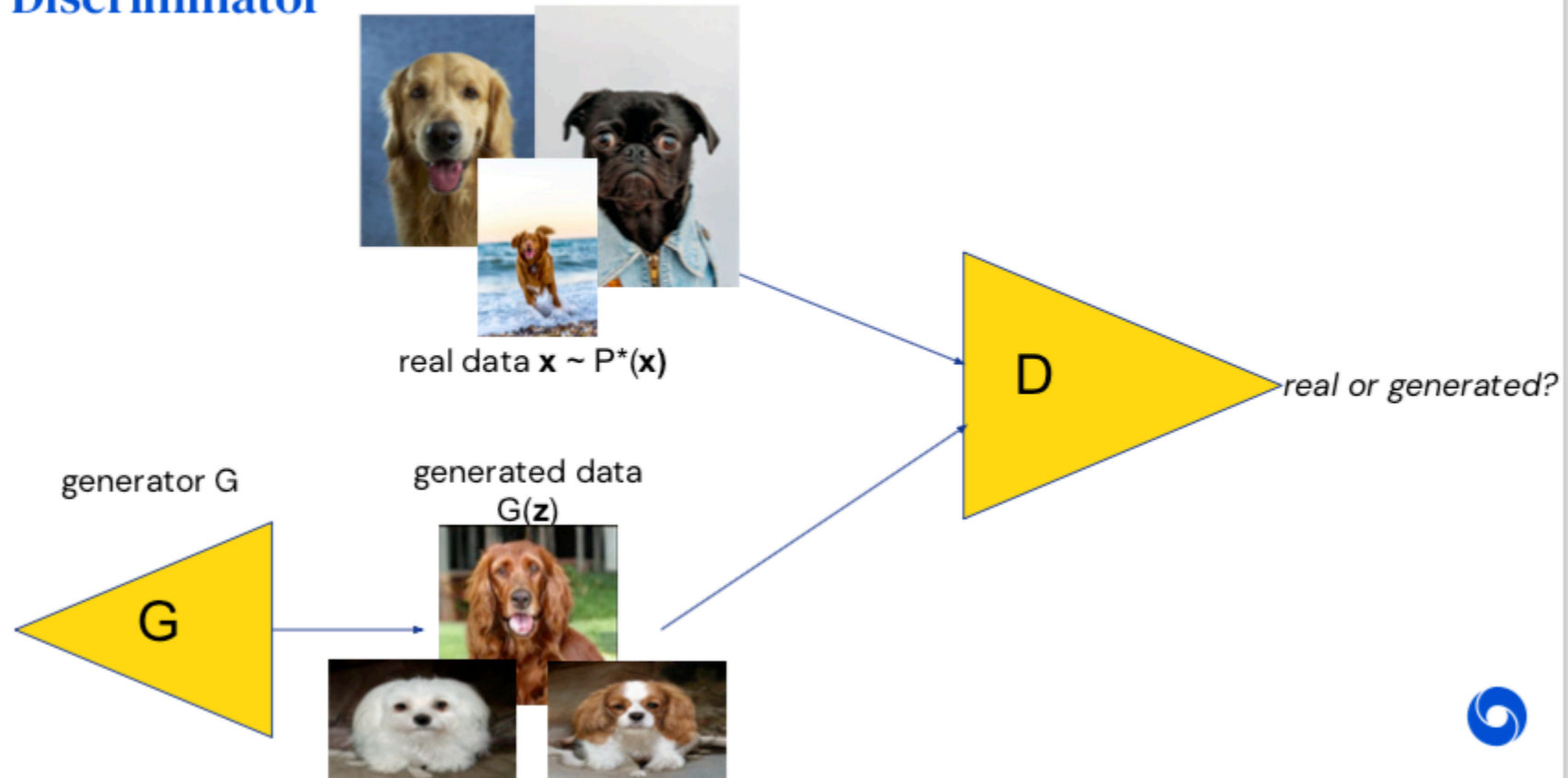


It was the best of times, it was the worst of times, it was the age of wisdom, it was the age of foolishness...

Generation of text data

Generative Adversarial Networks

Discriminator



Generative Adversarial Networks

$$\min_G \max_D V(D, G) = \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}(\mathbf{x})} [\log D(\mathbf{x})] + \mathbb{E}_{\mathbf{z} \sim p_{\mathbf{z}}(\mathbf{z})} [\log(1 - D(G(\mathbf{z})))]$$

log-probability that D correctly predicts real data \mathbf{x} are real

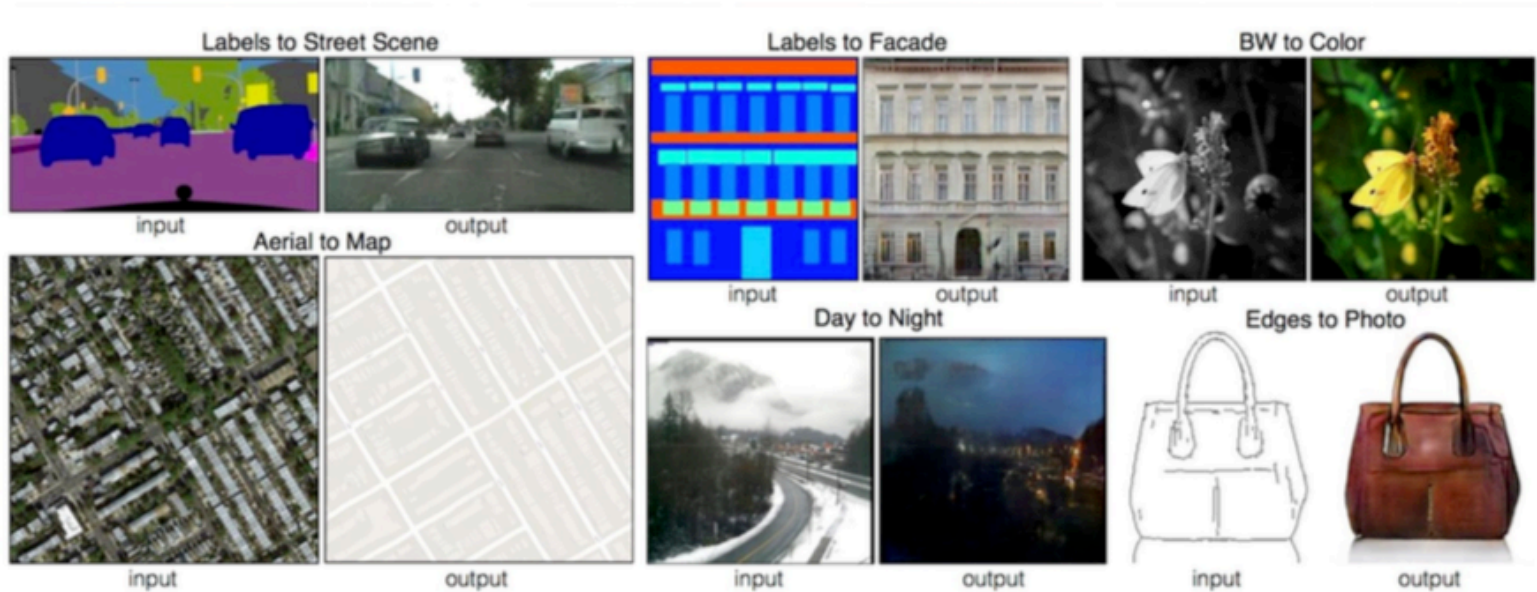
log-probability that D correctly predicts generated data $G(\mathbf{z})$ are generated

discriminator's (D) goal: **maximize** prediction accuracy

generator's (G) goal: **minimize** D's prediction accuracy, by **fooling** D into believing its outputs $G(\mathbf{z})$ are real as often as possible



Applications



Example results on several image-to-image translation problems. In each case we use the same architecture and objective, simply training on different data.

- Train a generator to **translate** between images of two different domains
- Standard GAN objective combined with reconstruction error

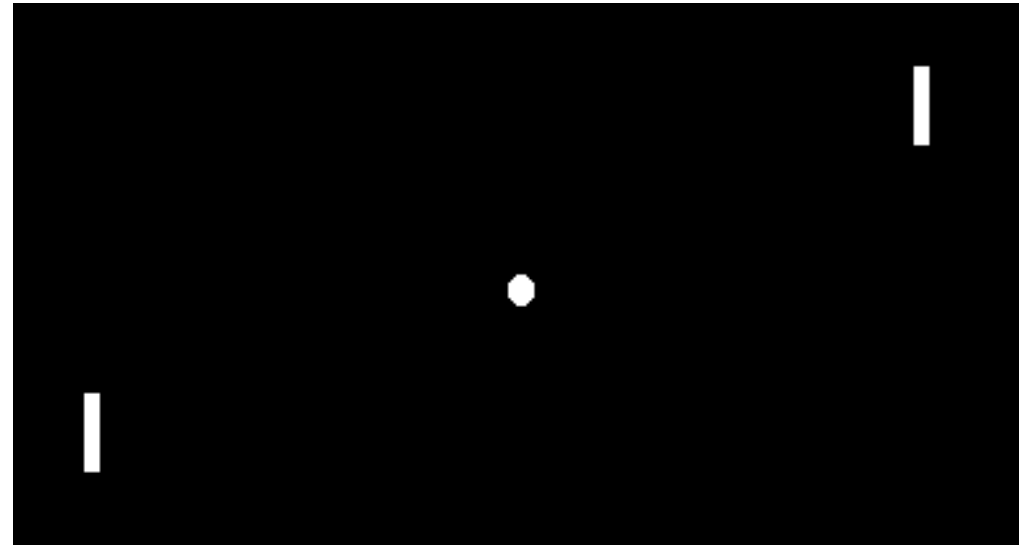
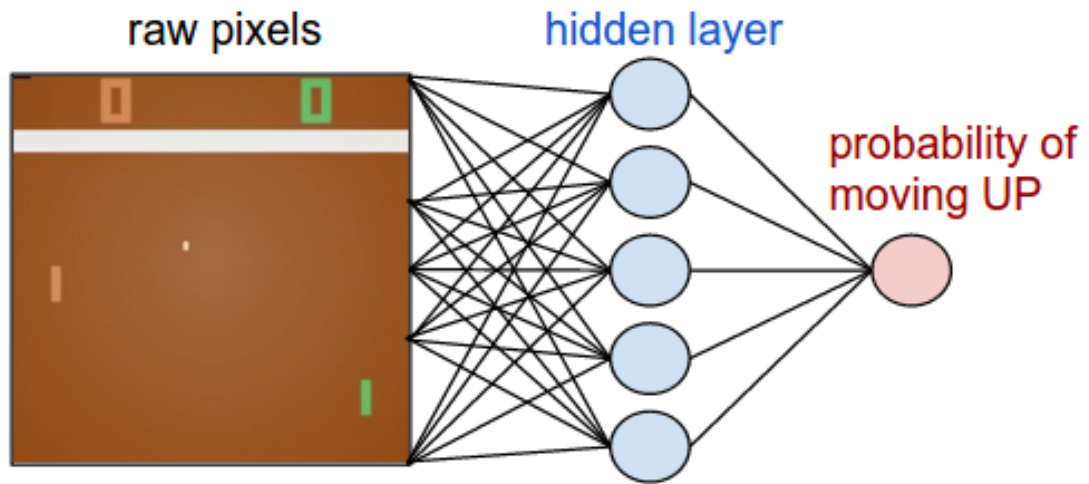
$$\mathcal{L}_{GAN}(G, D) = \mathbb{E}_y[\log D(y)] + \mathbb{E}_{x,z}[\log(1 - D(G(x, z)))].$$

$$\mathcal{L}_{L1}(G) = \mathbb{E}_{x,y,z}[\|y - G(x, z)\|_1].$$

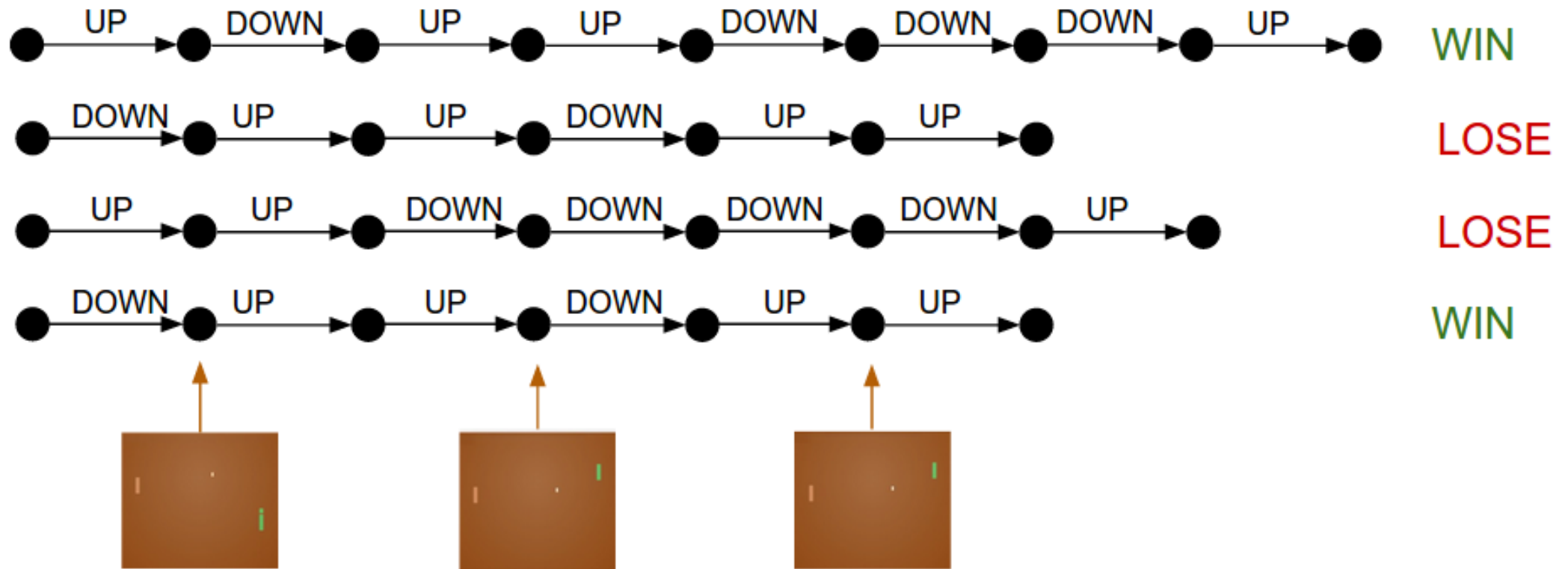
$$G^* = \arg \min_G \max_D \mathcal{L}_{cGAN}(G, D) + \lambda \mathcal{L}_{L1}(G).$$

Pix2Pix (Isola et al.)

Neural Networks in RL



Sequential Task



Parameterized Policy

- Class of policies defined by parameters θ

$$\boxed{\pi_{\theta}(a|s) : \mathcal{S} \rightarrow \mathcal{A}}$$

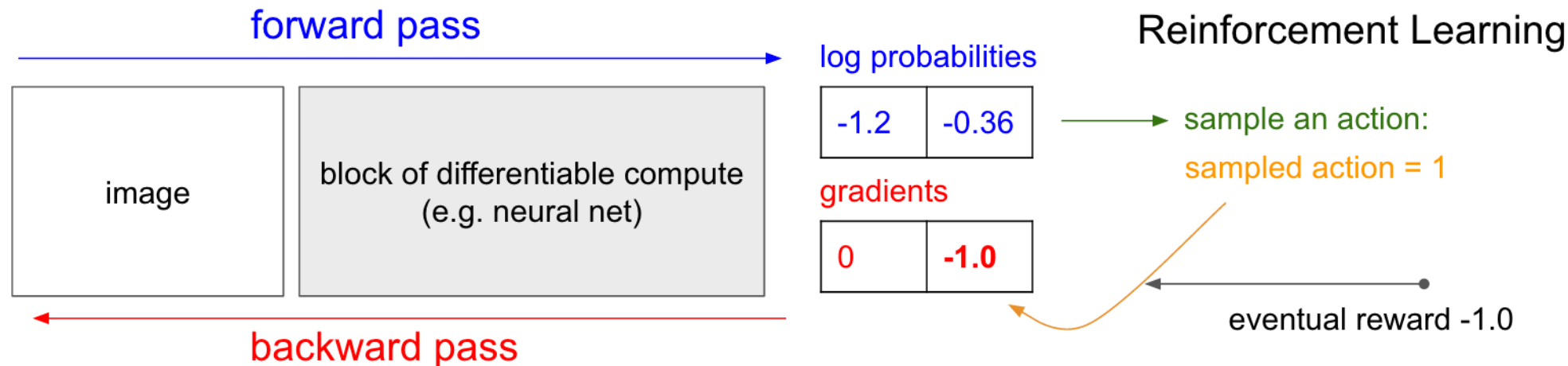
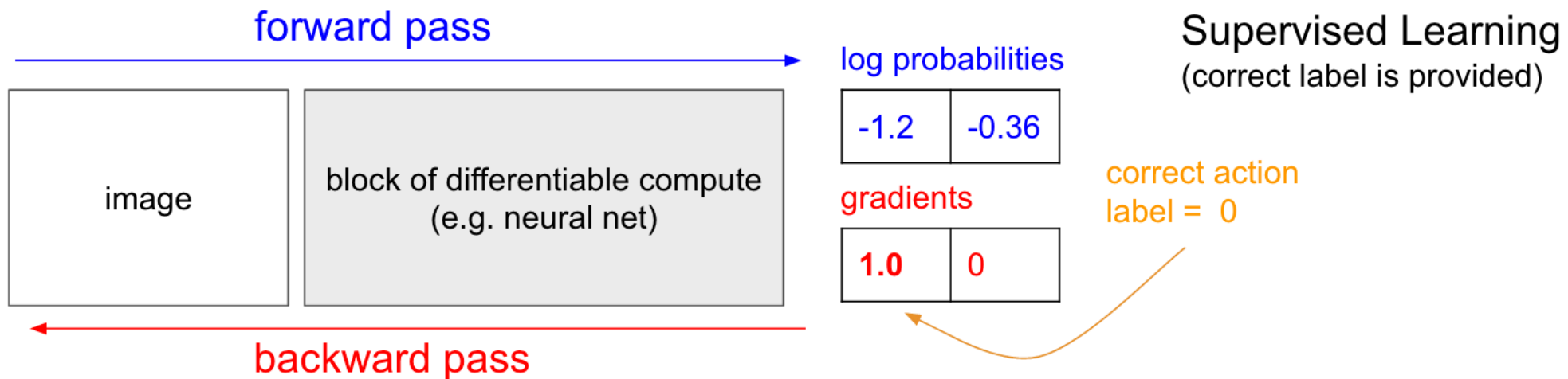
- Eg: θ can be parameters of linear transformation, deep network, etc.

- Want to maximize:
$$J(\pi) = \mathbb{E} \left[\sum_{t=1}^T \mathcal{R}(s_t, a_t) \right]$$

- In other words,

$$\pi^* = \arg \max_{\pi: \mathcal{S} \rightarrow \mathcal{A}} \mathbb{E} \left[\sum_{t=1}^T \mathcal{R}(s_t, a_t) \right] \quad \longrightarrow \quad \theta^* = \arg \max_{\theta} \mathbb{E} \left[\sum_{t=1}^T \mathcal{R}(s_t, a_t) \right]$$

Training signal comes from reward



Gathering Experience

- ◆ Slightly re-writing the notation

Let $\tau = (s_0, a_0, \dots, s_T, a_T)$ denote a trajectory

$$\begin{aligned}\pi_{\theta}(\tau) &= p_{\theta}(\tau) = p_{\theta}(s_0, a_0, \dots, s_T, a_T) \\ &= p(s_0) \prod_{t=0}^{T-1} p_{\theta}(a_t | s_t) \cdot p(s_{t+1} | s_t, a_t)\end{aligned}$$

$$\arg \max_{\theta} \mathbb{E}_{\tau \sim p_{\theta}(\tau)} [\mathcal{R}(\tau)]$$

Gathering Experience

$$\begin{aligned} J(\theta) &= \mathbb{E}_{\tau \sim p_{\theta}(\tau)} [\mathcal{R}(\tau)] \\ &= \mathbb{E}_{a_t \sim \pi(\cdot | s_t), s_{t+1} \sim p(\cdot | s_t, a_t)} \left[\sum_{t=0}^T \mathcal{R}(s_t, a_t) \right] \end{aligned}$$

- ◆ How to gather data?
 - ◆ We already have a policy: π_{θ}
 - ◆ Sample N trajectories $\{\tau_i\}_{i=1}^N$ by acting according to π_{θ}

$$\approx \frac{1}{N} \sum_{i=1}^N \sum_{t=1}^T r(s_t^i, a_t^i)$$

Reinforce Algorithm

◆ Sample trajectories $\tau_i = \{s_1, a_1, \dots, s_T, a_T\}_i$ by acting according to π_θ

◆ Compute policy gradient as

$$\nabla_\theta J(\theta) \approx \frac{1}{N} \sum_i \left[\sum_{t=1}^T \nabla_\theta \log \pi_\theta (a_t^i | s_t^i) \cdot \sum_{t=1}^T \mathcal{R} (s_t^i | a_t^i) \right]$$

◆ Update policy parameters: $\theta \leftarrow \theta + \alpha \nabla_\theta J(\theta)$

