# CIS 519/419
# Applied Machine Learning
## www.seas.upenn.edu/~cis519

Dan Roth

danroth@seas.upenn.edu

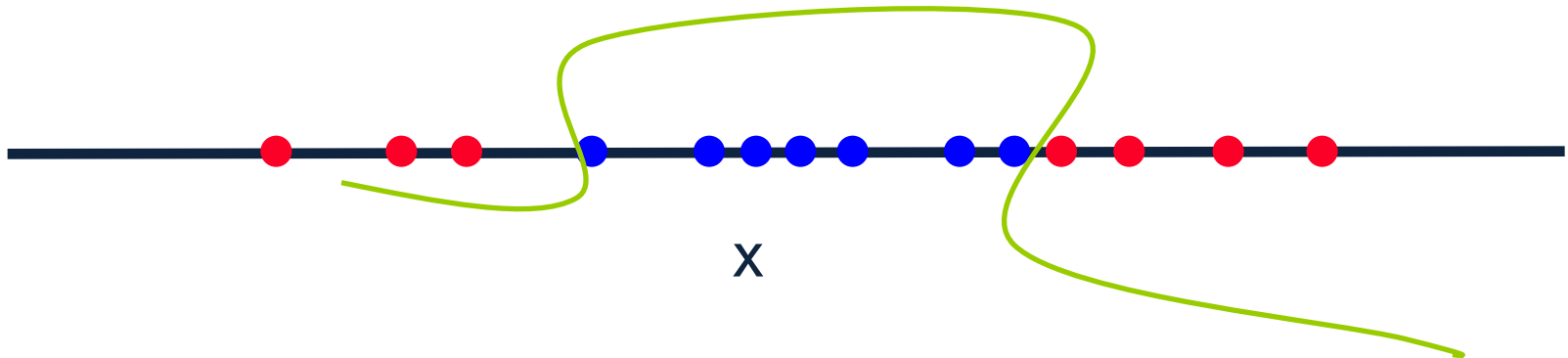http://www.cis.upenn.edu/~danroth/

461C, 3401 Walnut

Lecture given by Daniel Khashabi
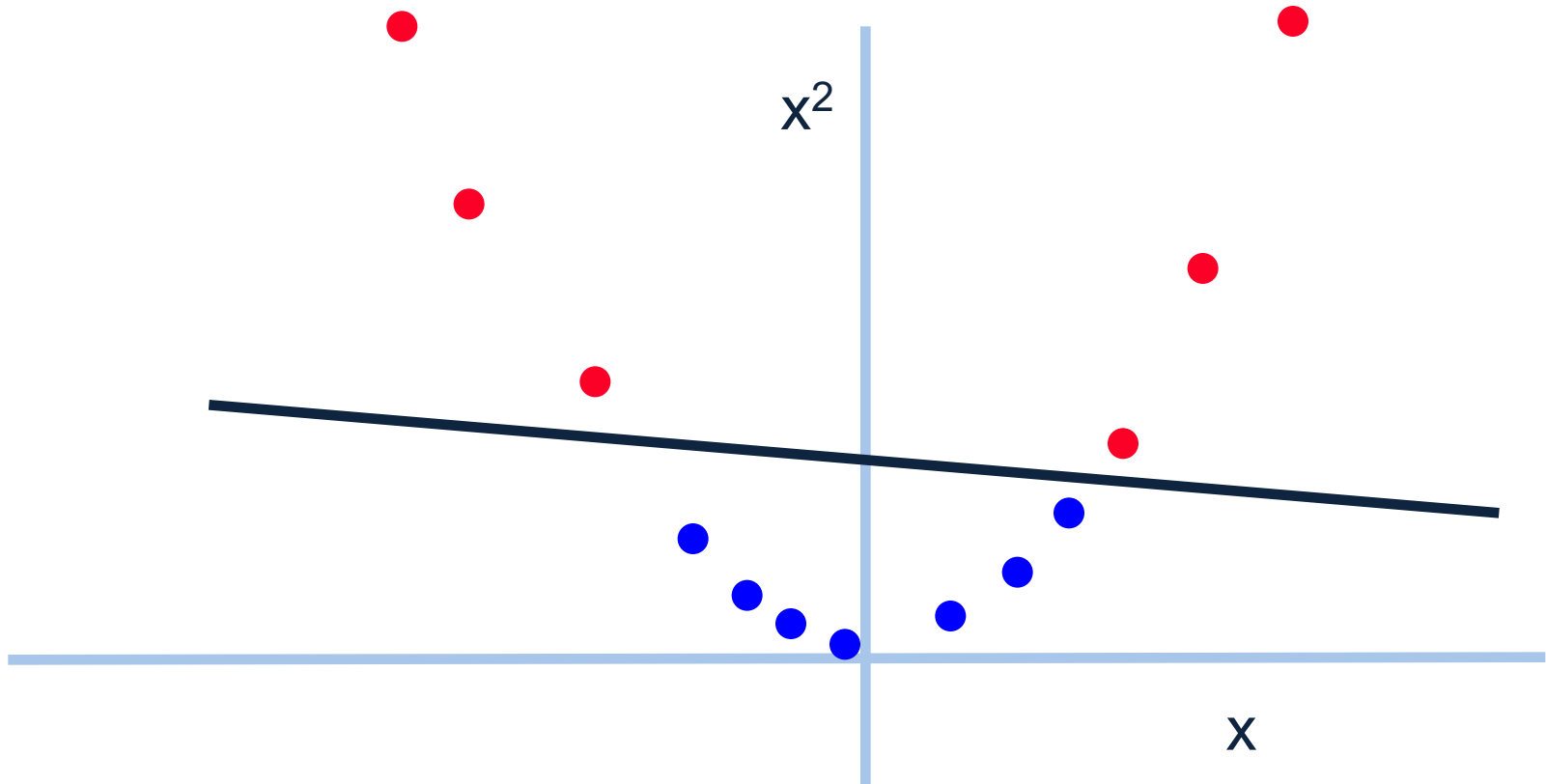
CIS419/519 Spring '18

# Functions Can be Made Linear

- Data are not linearly separable in one dimension
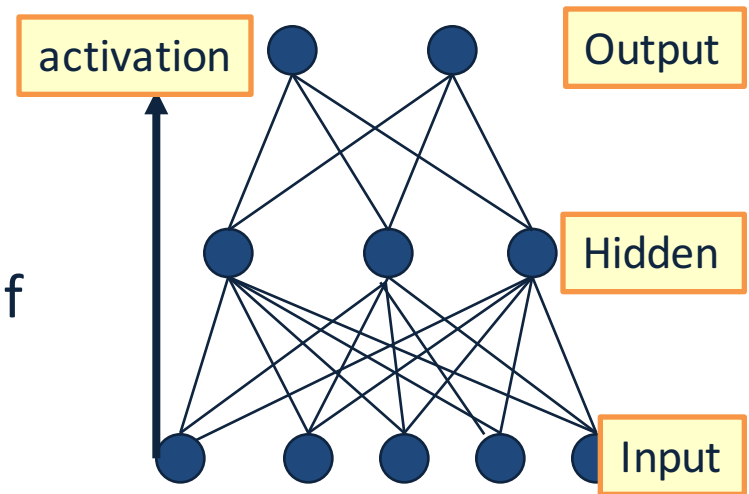- Not separable if you insist on using a specific class of functions

x

# Blown Up Feature Space
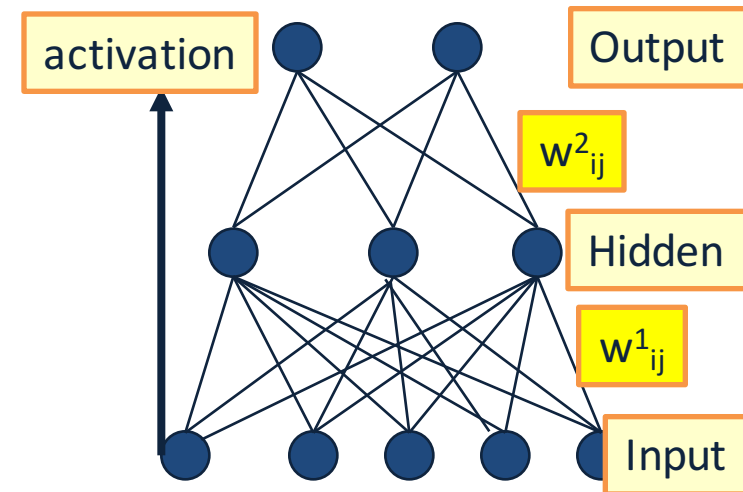
- Data are separable in $\langle x, x^2 \rangle$ space

# Multi-Layer Neural Network

- Multi-layer network were designed to overcome the computational (expressivity) limitation of a single threshold element.

- The idea is to stack several layers of threshold elements, each layer using the output of the previous layer as input.

activation

Output

Hidden

Input

- Multi-layer networks can represent arbitrary functions, but building effective learning methods for such network was [thought to be] difficult.
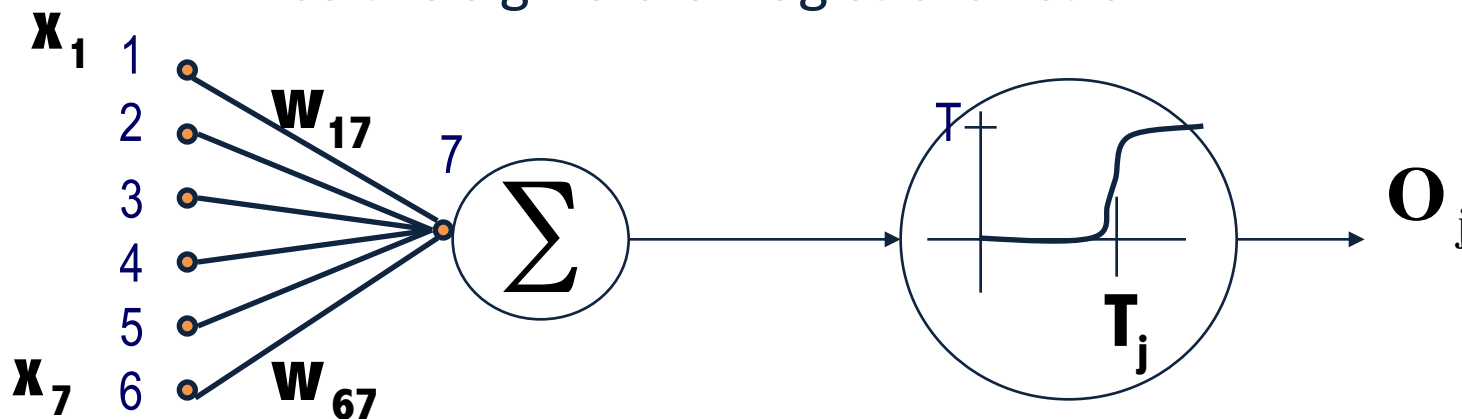
# Basic Units

- Linear Unit: Multiple layers of linear functions $o_j = w \cdot x$ produce linear functions. We want to represent nonlinear functions.

- **Need to do it in a way that facilitates learning**

- Threshold units: $o_j = \text{sgn}(w \cdot x)$ are not differentiable, hence unsuitable for gradient descent.

- The key idea was to notice that the discontinuity of the threshold element can be represents by a smooth non-linear approximation: $o_j = [1+ \exp\{-w \cdot x\}]^{-1}$

- (Rumelhart, Hinton, Williiam, 1986), (Linnainmaa, 1970) , see: http://people.idsia.ch/~juergen/who-invented-backpropagation.html )

activation

Output

$w^2_{ij}$

Hidden

$w^1_{ij}$

Input

# Model Neuron (Logistic)

- Us a non-linear, differentiable output function such as the sigmoid or logistic function



- Net input to a unit is defined as: $net_j = \sum w_{ij} \bullet x_i$
- Output of a unit is defined as:

$$O_j = \frac{1}{1 + e^{-(net_j - T_j)}}$$
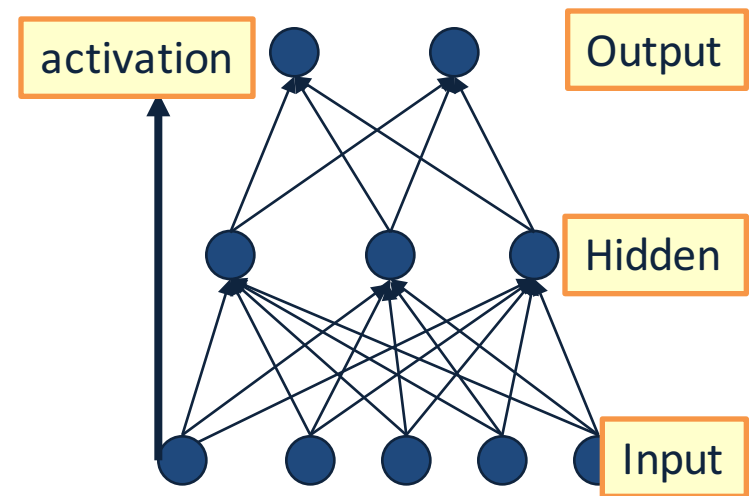
# Neural Networks

- Neural Networks are **functions**: $NN: X \rightarrow Y$
  - where $X = [0,1]^n$, or $\{0,1\}^n$ and $Y = [0,1], \{0,1\}$
  - **Robust** approach to approximating **real-valued**, **discrete-valued** and **vector valued** target functions.
- Among the most effective **general purpose** supervised learning method currently known.
- Effective especially for **complex and hard to interpret input data** such as real-world sensory data, where a lot of supervision is available.
- The **Backpropagation algorithm** for neural networks has been shown successful in many practical problems
  - handwritten character recognition, speech recognition, object recognition, some NLP problems

# Neural Networks

- Neural Networks are **functions**: $NN: X \rightarrow Y$
  - where $X = [0,1]^n$, or $\{0,1\}^n$ and $Y = [0,1], \{0,1\}$
- NN can be used as an approximation of a target classifier
  - In their general form, even with a single hidden layer, NN can approximate any function
  - Algorithms exist that can learn a NN representation from labeled training data  (e.g., Backpropagation).

# Multi-Layer Neural Networks

- Multi-layer network were designed to overcome the computational (**expressivity**) limitation of a single threshold element.

- The idea is to **stack** several layers of threshold elements, each layer using the output of the previous layer as input.

# Motivation for Neural Networks

- Inspired by **biological systems**
  - **But don't take this (as well as any other words in the new on "emergence" of intelligent behavior) seriously;**
- We are currently on rising part of a wave of interest in NN architectures, after a long downtime from the mid-90-ies.
  - Better computer architecture (GPUs, parallelism)
  - A lot more data than before; in many domains, supervision is available.
- Current surge of interest has seen very minimal algorithmic changes
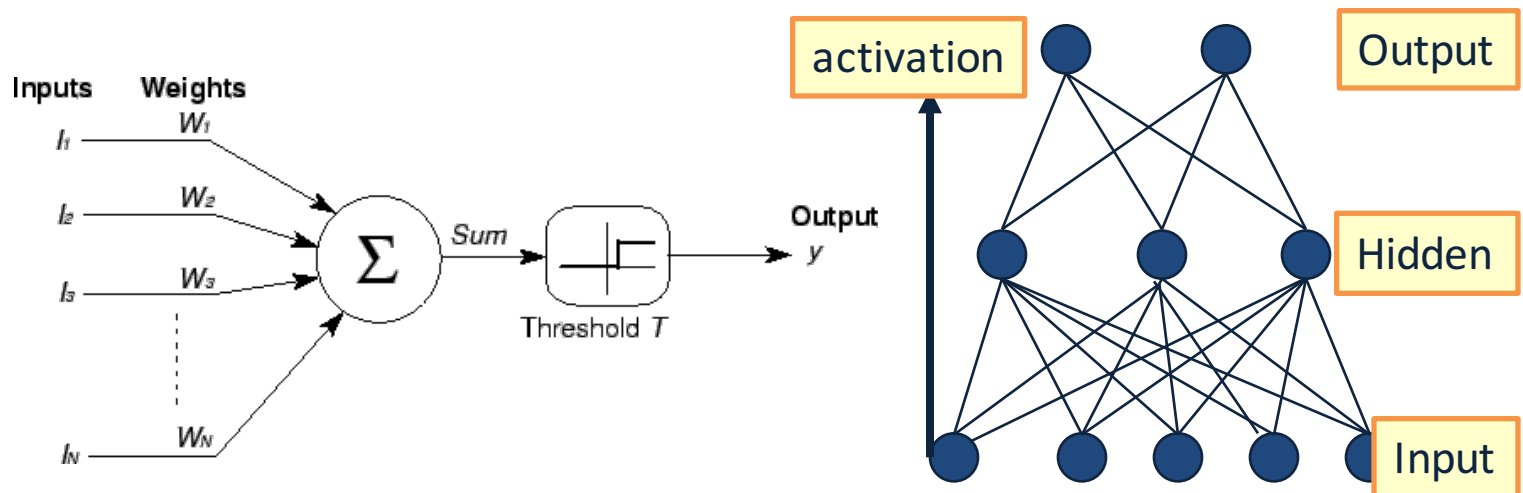
# Motivation for Neural Networks

- Minimal to no algorithmic changes

- One potentially interesting perspective:
  - Before we looked at NN only as function approximators.
  - Now, we look at the intermediate representations generated while learning as meaningful
  - Ideas are being developed on the value of these intermediate representations for transfer learning etc.

- We will present in the next two lectures a few of the basic architectures and learning algorithms, and provide some examples for applications

# Neural Speed Constraints

- Neuron "switching time" is **O(milliseconds)**, compared to nanosecond for transistors.

    - However, biological systems can perform significant cognitive tasks (vision, language understanding) in fractions of a second.

- Even for limited abilities, current AI systems require orders of magnitude more steps.

- Human brain has approximately 10^10 neurons, each connected to 10^4; must explore massive parallelism (but there's more...)
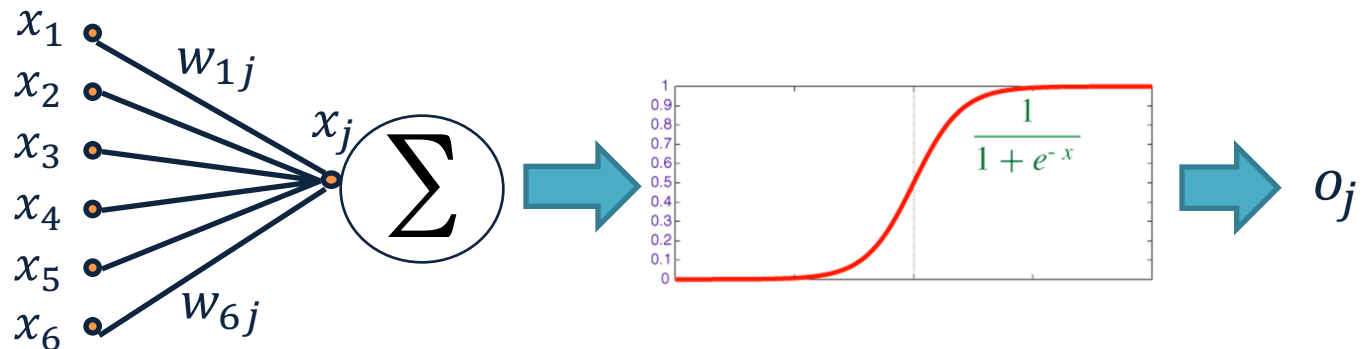
# Basic Unit in Multi-Layer Neural Network

- **Linear Unit**: $o_j = \vec{w}.\vec{x}$ multiple layers of linear functions produce linear functions. We want to represent nonlinear functions.

- **Threshold units:** $o_j = sgn(\vec{w}.\vec{x} - T)$ are not differentiable, hence unsuitable for gradient descent

# Model Neuron (Logistic)

- Neuron is modeled by a unit $j$ connected by weighted links $w_{ij}$ to other units $i$.



- Use a non-linear, differentiable output function such as the sigmoid or logistic function

- Net input to a unit is defined as:
$$\text{net}_j = \sum w_{ij}.x_i$$

- Output of a unit is defined as:
$$o_j = \frac{1}{1 + \exp\big(-(\text{net}_j - T_j)\big)}$$

# History: Neural Computation

- **McCulloch and Pitts (1943)** showed how linear threshold units can be used to compute logical functions

- Can build basic logic gates
  - **AND:** $w_{ij} = T_j/n$
  - **OR:** $w_{ij} = T_j$
  - **NOT:** use negative weight

$$\text{net}_j = \sum w_{ij}.x_i$$
$$o_j = \frac{1}{1 + \exp\left(-(\text{net}_j - T_j)\right)}$$

- Can build arbitrary logic circuits, finite-state machines and computers given these basis gates.

- Can specify any Boolean function using two layer network (w/ negation)
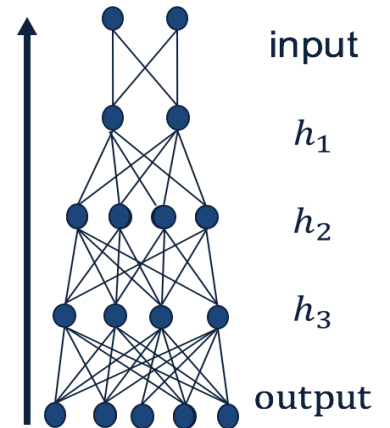  - DNF and CNF are universal representations

# Representational Power

- Any Boolean function can be represented by a **two layer** network (simulate a two layer AND-OR network)

- Any **bounded continuous function** can be approximated with **arbitrary small error** by a **two layer** network.

- Sigmoid functions provide a set of **basis function** from which arbitrary function can be composed.

- **Any function** can be approximated to arbitrary accuracy by a **three layer** network.

# Quiz Time!

- Given a neural network, how can we make predictions?

  - Given input, calculate the output of each layer (starting from the first layer), until you get to the output.

- What is required to fully specify a neural network?

  - The weights.

- Why NN predictions can be quick?

  - Because many of the computations could be parallelized.

- What makes a neural networks **non**-linear approximator?

  - The non-linear units.



input

$h_1$

$h_2$

$h_3$

output

# Training a Neural Net

# Widrow-Hoff Rule

- This incremental update rule provides an approximation to the goal:
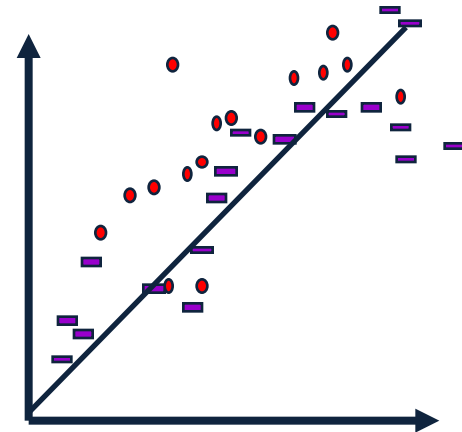  - Find the best linear approximation of the data

$$Err\left(\vec{w}^{(j)}\right) = \frac{1}{2}\sum_{d \in D}(t_d - o_d)^2$$

  - where:

$$o_d = \sum_i w_{ij}.x_i = \vec{w}^{(j)}.\vec{x}$$

  output of linear unit on example $d$
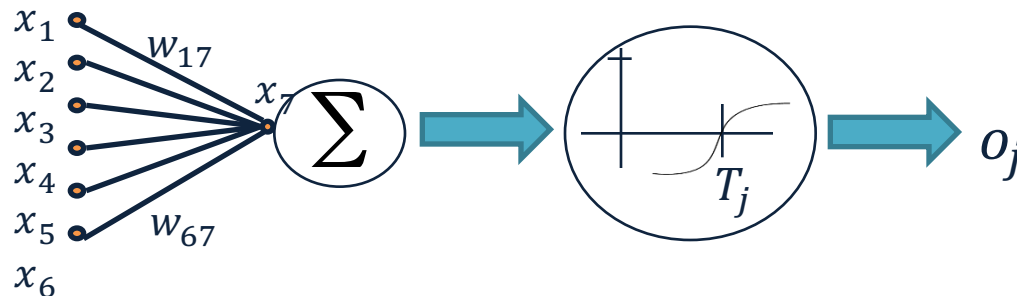  - $t_d$ = Target output for example d

# History: Learning Rules

- Hebb (1949) suggested that if two units are both active (firing) then the weights between them should increase:

$$w_{ij} = w_{ij} + R o_i o_j$$

  - $R$ and is a constant called **the learning rate**
  - Supported by physiological evidence

- Rosenblatt (1959) suggested that when a target output value is provided for a single neuron **with fixed input**, it can **incrementally change weights** and learn to produce the output using the **Perceptron learning rule**.

  - assumes **binary output** units; single linear threshold unit
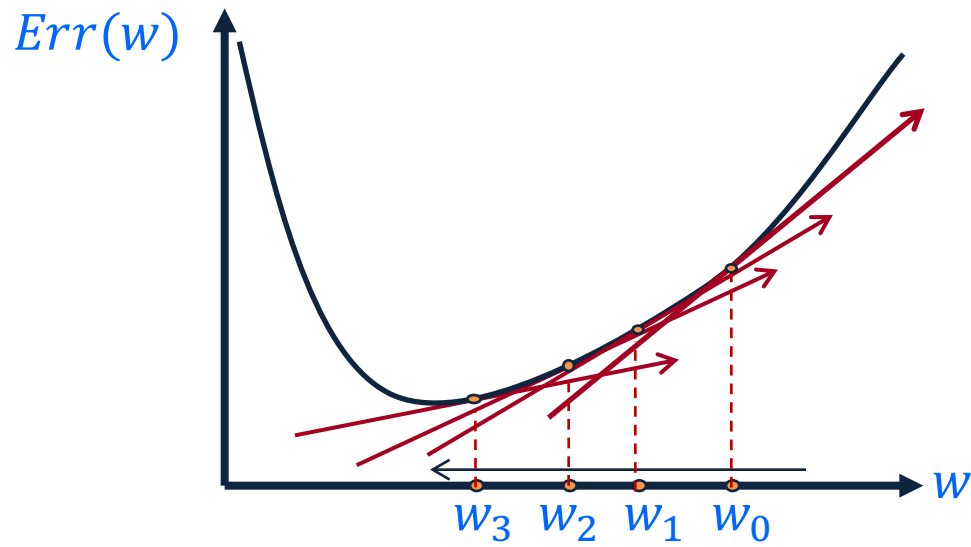  - Led to the Perceptron Algorithm

- See: http://people.idsia.ch/~juergen/who-invented-backpropagation.html

# Perceptron Learning Rule

- Given:
  - the **target** output for the output unit is $t_j$
  - the **input** the neuron sees is $x_i$
  - the **output** it **produces** is $o_j$

- Update weights according to $w_{ij} \leftarrow w_{ij} + R(t_j - o_j)x_i$
  - If output is **correct**, don't change the weights
  - If output is **wrong**, change weights for all inputs which are 1
    - If output is low (0, needs to be 1) increment weights
    - If output is high (1, needs to be 0) decrement weights

# Gradient Descent

- We use gradient descent determine the weight vector that minimizes $Err(\vec{w}^{(j)})$;

- Fixing the set $D$ of examples, $E$ is a function of $\vec{w}^{(j)}$

- At each step, the weight vector is modified in the direction that produces the steepest descent along the error surface.

# Summary: Single Layer Network

- Variety of update rules
  - Multiplicative
  - Additive
- **Batch** and **incremental** algorithms
- Various convergence and efficiency conditions
- There are other ways to learn linear functions
  - Linear Programming (general purpose)
  - Probabilistic Classifiers ( some assumption)
- Key algorithms are driven by gradient descent

# General Stochastic Gradient Algorithms

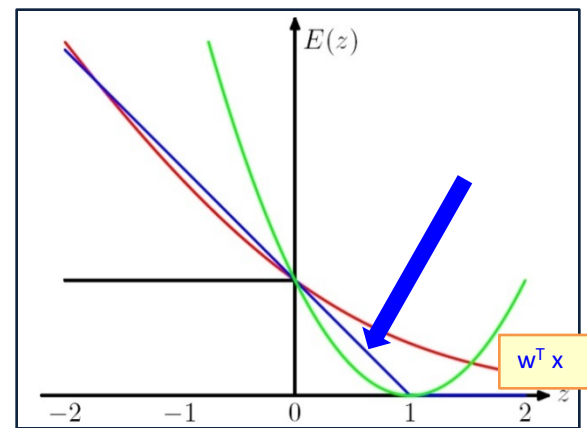Learning rate — gradient — **The loss Q:** a function of x, w and y

$$w_{t+1} = w_t - r_t\, g_w\, Q(x_t, y_t, w_t) = w_t - r_t\, g_t$$

- LMS: $Q((x, y), w) = 1/2\,(y - w^T x)^2$
- leads to the update rule (Also called Widrow's Adaline):

$$w_{t+1} = w_t + r\,(y_t - w_t^T x_t)\, x_t$$

- Here, even though we make binary predictions based on sgn $(w^T x)$ we do not take the sign of the dot-product into account in the loss.

- Another common loss function is:
- Hinge loss:

  $Q((x, y), w) = \max(0, 1 - y\, w^T x)$

- This leads to the perceptron update rule:



- If $y_i\, w_i^T \cdot x_i > 1$   (No mistake, by a margin):   No update
- Otherwise   (Mistake, relative to margin):   $w_{t+1} = w_t + r\, y_t\, x_t$

Here g = -yx
Good to think about the case of Boolean examples

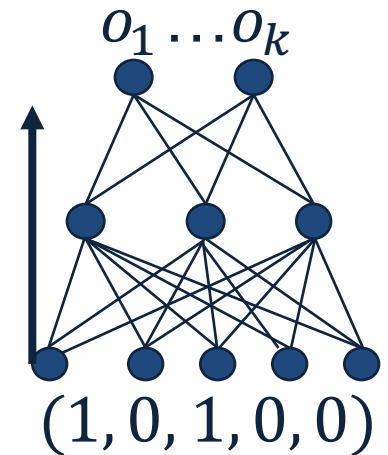# Summary: Single Layer Network

- Variety of update rules
    - Multiplicative
    - Additive
- **Batch** and **incremental** algorithms
- Various convergence and efficiency conditions
- There are other ways to learn linear functions
    - Linear Programming (general purpose)
    - Probabilistic Classifiers ( some assumption)
- Key algorithms are driven by gradient descent
- However, the representational restriction is limiting in many applications

# Backpropagation Learning Rule

- Since there could be multiple output units, we define the **error** as the sum over all the network output units.

$$Err(\vec{w}) = \frac{1}{2}\sum_{d \in D}\sum_{k \in K} \quad (t_{kd} - o_{kd})^2$$

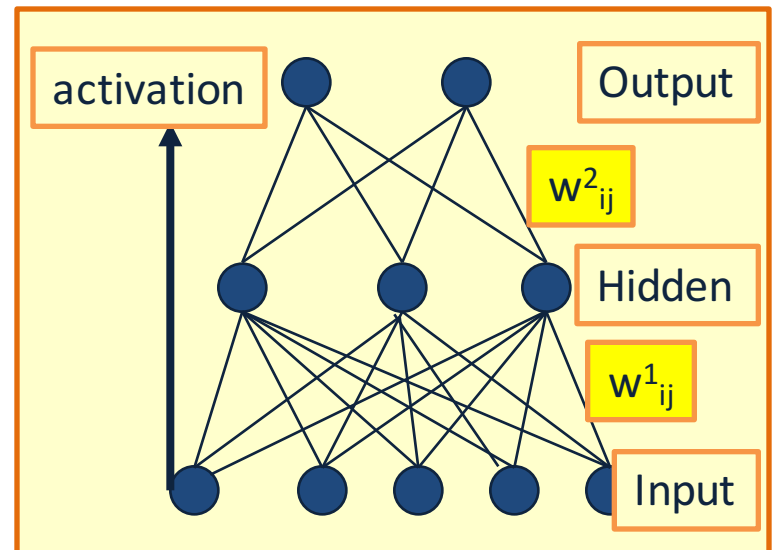  - where $D$ is the set of training examples,
  - $K$ is the set of output units

$$o_1 \dots o_k$$

$$(1, 0, 1, 0, 0)$$

- This is used to derive the (global) learning rule which performs gradient descent in the weight space in an attempt to minimize the error function.

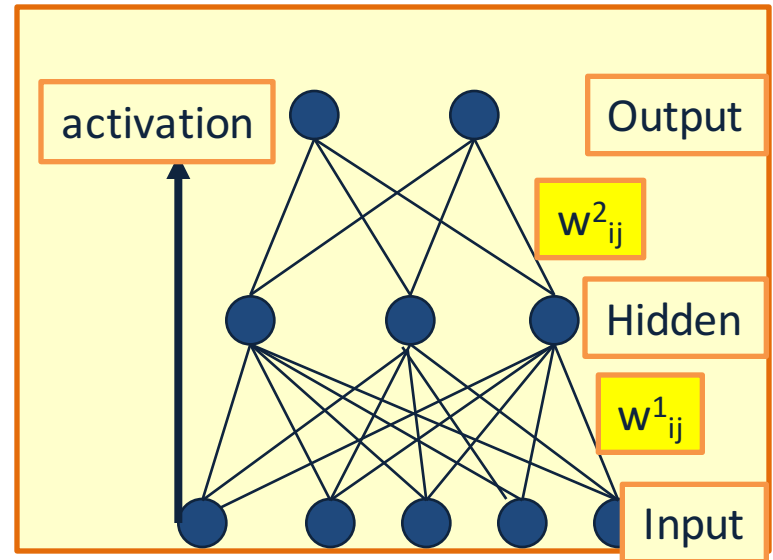$$\Delta w_{ij} = -R\frac{\partial E}{\partial w_{ij}}$$

Function 1

# Learning with a Multi-Layer Perceptron

- It's easy to learn the top layer – it's just a linear unit.

- Given feedback (truth) at the top layer, and the activation at the layer below it, you can use the Perceptron update rule (more generally, gradient descent) to updated these weights.

- The problem is what to do with the other set of weights – we do not get feedback in the intermediate layer(s).

# Learning with a Multi-Layer Perceptron

- The problem is what to do with the other set of weights – we do not get feedback in the intermediate layer(s).
- Solution: If all the activation functions are differentiable, then the **output** of the network is also a differentiable function of the input and weights in the network.
- Define an error function (e.g., sum of squares) that is a differentiable function of the output, i.e. this error function is also a differentiable function of the weights.
- We can then evaluate the derivatives of the error with respect to the weights, and use these derivatives to find weight values that minimize this error function, using gradient descent (or other optimization methods).
- This results in an algorithm called back-propagation.



activation | Output
$w^2_{ij}$
Hidden
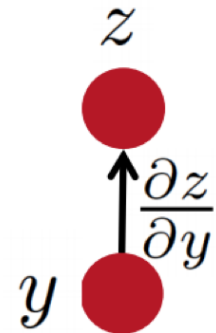$w^1_{ij}$
Input

# Some facts from real analysis

First let's get the notation right:

The arrow shows functional dependnence of z on y
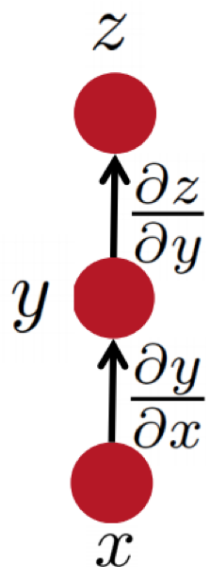    i.e. given y, we can calculate z.
    e.g., for example: z(y) = 2y^2


The derivative of z, with respect to y.

$$z$$

$$\frac{\partial z}{\partial y}$$

$$y$$

# Some facts from real analysis

- Simple chain rule
  - If $z$ is a function of $y$, and $y$ is a function of $x$
    - Then $z$ is a function of $x$, as well.

  - Question: how to find $\dfrac{\partial z}{\partial x}$

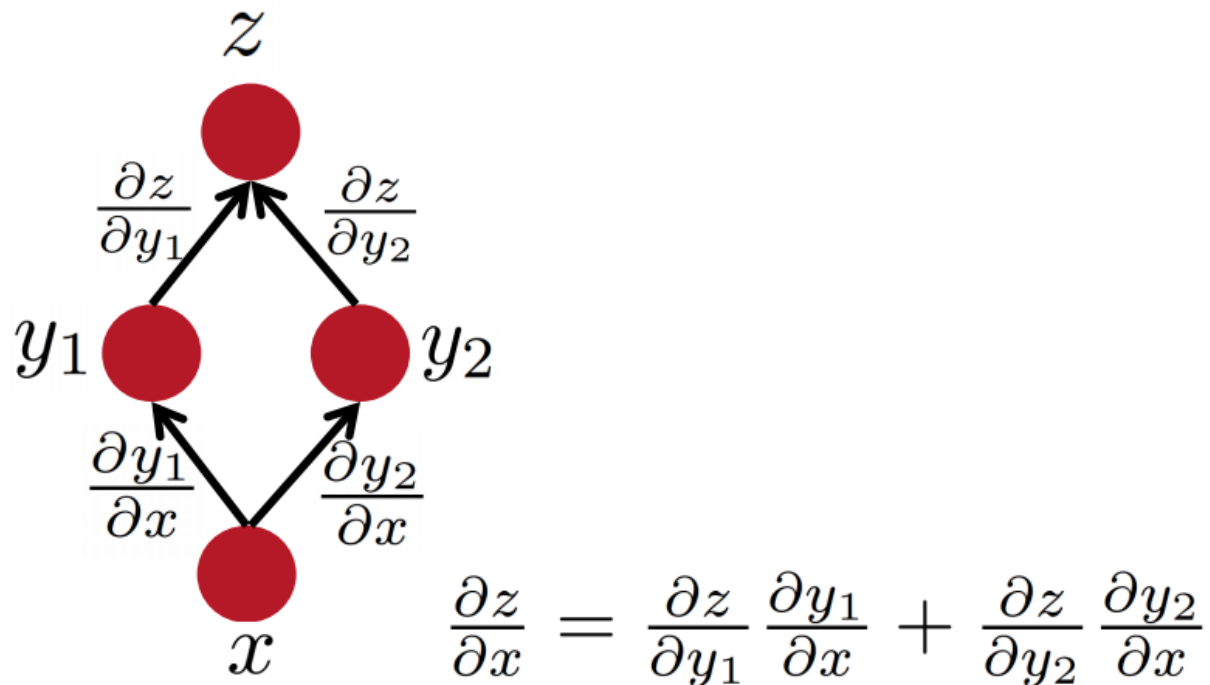$$\frac{\partial z}{\partial x} = \frac{\partial z}{\partial y} \frac{\partial y}{\partial x}$$

We will use these facts to derive the details of the Backpropagation algorithm.

z will be the error (loss) function.
- We need to know how to differentiate z

Intermediate nodes use a logistics function (or another differentiable step function).
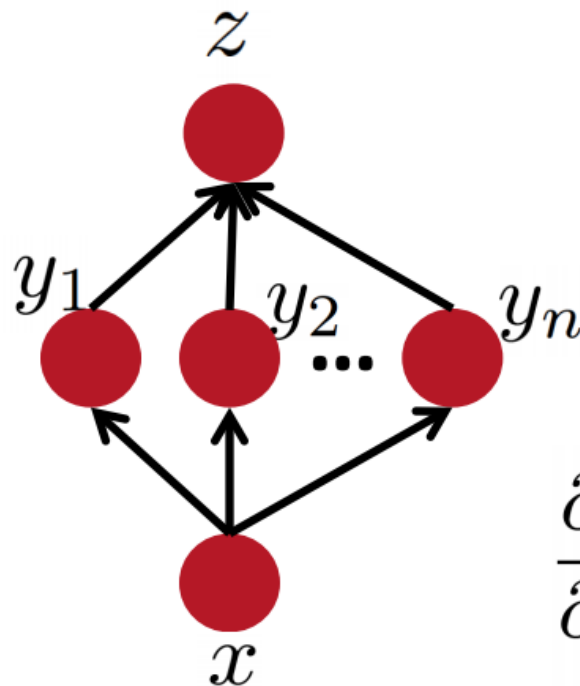- We need to know how to differentiate it.

# Some facts from real analysis

- Multiple path chain rule

$$z$$

$$\frac{\partial z}{\partial y_1} \qquad \frac{\partial z}{\partial y_2}$$

$$y_1 \qquad y_2$$

$$\frac{\partial y_1}{\partial x} \qquad \frac{\partial y_2}{\partial x}$$

$$x$$

$$\frac{\partial z}{\partial x} = \frac{\partial z}{\partial y_1}\frac{\partial y_1}{\partial x} + \frac{\partial z}{\partial y_2}\frac{\partial y_2}{\partial x}$$

Slide Credit: Richard Socher

# Some facts from real analysis

- Multiple path chain rule: general

$$\frac{\partial z}{\partial x} = \sum_{i=1}^{n} \frac{\partial z}{\partial y_i} \frac{\partial y_i}{\partial x}$$

Slide Credit: Richard Socher

# Key Intuitions Required for BP

- **Gradient Descent**
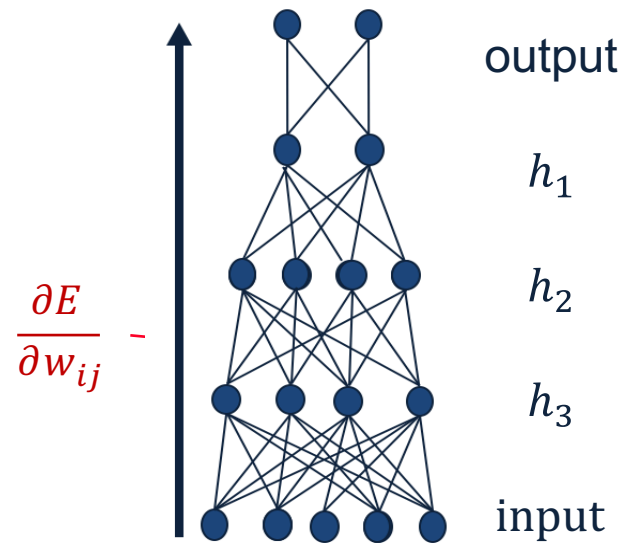  - Change the weights in the direction of gradient to minimize the error function.

$Err(w)$

- **Chain Rule**
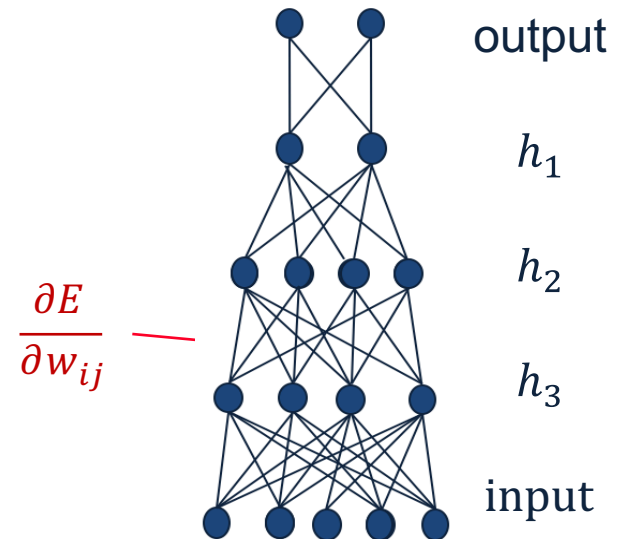  - Use the chain rule to calculate the weights of the intermediate weights

- **Dynamic Programming (Memoization)**
  - Memoize the weight updates to make the updates faster.
  - The "back" part of "backpropagation"

$\frac{\partial E}{\partial w_{ij}}$

output

$h_1$

$h_2$

$h_3$

input

# Backpropagation: the big picture

- Loop over instances:

  1. The forward step
     - Given the input, make predictions layer-by-layer, starting from the first layer)

  2. The backward step
     - Calculate the error in the output
     - Update the weights layer-by-layer, starting from the final layer

output

$h_1$

$h_2$

$\dfrac{\partial E}{\partial w_{ij}}$

$h_3$

input

# Quiz time!

- **What is the purpose of forward step?**
  - To make predictions, given an input.

- **What is the purpose of backward step?**
  - To update the weights, given an output error.

- **Why do we use the chain rule?**
  - To calculate gradient in the intermediate layers.

- **Why backpropagation could be efficient?**
  - Because it can be parallelized.

# Deriving the update rules

# Reminder: Model Neuron (Logistic)

- Neuron is mod███ ██ ████ ██ ███████ ████ ███████ weighted links $w_{ij}$ to oth███

<div style="border:1px solid #ccc; background:#fdf6cc; padding:8px;">

**The parameters so far?**

The set of connective weights: $w_{ij}$

The threshold value: $T_j$

</div>

$x_1$
$x_2$    $w_{17}$
$x_3$    $x_7$   $\Sigma$    $\dfrac{1}{1+e^{-x}}$    $o_j$
$x_4$
$x_5$
$x_6$    $w_{67}$

- Use a non-linear, differentiable output function such as the sigmoid or logistic function

- Net input to a unit is defined as: $\qquad \text{net}_j = \sum w_{ij} . x_i$

- Output of a unit is defined as: $\qquad o_j = \dfrac{1}{1 + \exp\bigl(-(\text{net}_j - T_j)\bigr)}$

# Derivatives

- Function 1 (error):

  - $E = \frac{1}{2}\sum_{k \in K}(t_k - o_k)^2$

  - $\frac{\partial E}{\partial o_i} = -(t_i - o_i)$

- Function 2 (linear gate):

  - $\text{net}_j = \sum w_{ij}.x_i$

  - $\frac{\partial \text{net}_j}{\partial w_{ij}} = x_i$

- Function 3 (differentiable step function):

  - $o_i = \frac{1}{1+\exp\{-(\text{net}_j - T)\}}$

  - $\frac{\partial o_i}{\partial \text{net}_j} = \frac{\exp\{-(\text{net}_j - T)\}}{(1+\exp\{-(\text{net}_j - T)\})^2} = o_i(1 - o_i)$



$o_1 \ldots o_k$

$j$

$w_{ij}$

$i$

# Derivation of Learning Rule

- The weights are updated incrementally;  the error is computed **for each example** and the weight update is then derived.

$$E_d(\vec{w}) = \frac{1}{2} \sum_{k \in K} (t_k - o_k)^2$$

- $w_{ij}$ influences the output only through $\text{net}_j$

$$o_i = \frac{1}{1 + \exp\{-(\text{net}_j - T)\}} \quad \text{and} \quad \text{net}_j = \sum w_{ij} . x_i$$

- Therefore:

$$\frac{\partial E_d}{\partial w_{ij}} = \frac{\partial E_d}{\partial o_j} \frac{\partial o_j}{\partial \text{net}_j} \frac{\partial \text{net}_j}{\partial w_{ij}}$$

$o_1 \dots o_k$

$j$

$w_{ij}$

$i$

# Derivation of Learning Rule (2)

- Weight updates of output units:
  - $w_{ij}$ influences the output only through $\text{net}_j$

- Therefore:

$$\frac{\partial E_d}{\partial w_{ij}} = \frac{\partial E_d}{\partial o_j} \frac{\partial o_j}{\partial \text{net}_j} \frac{\partial \text{net}_j}{\partial w_{ij}}$$

$$= -(t_j - o_j) \, o_j(1 - o_j) \, x_i$$

$$E_d(\vec{w}) = \frac{1}{2} \sum_{k \in K} (t_k - o_k)^2$$

$$\frac{\partial o_j}{\partial \text{net}_j} = o_j(1 - o_j)$$

$$o_j = \frac{1}{1 + \exp\{-(\text{net}_j - T_j)\}}$$

$$\sum w_{ij} . x_i$$

$o_1 \ldots o_k$

$w_{ij}$

# Derivation of Learning Rule (3)

- Weights of output units:
  - $w_{ij}$ is changed by:

$$\Delta w_{ij} = R(t_j - o_j)o_j(1 - o_j)x_i$$
$$= R\delta_j x_i$$

Where we defined:

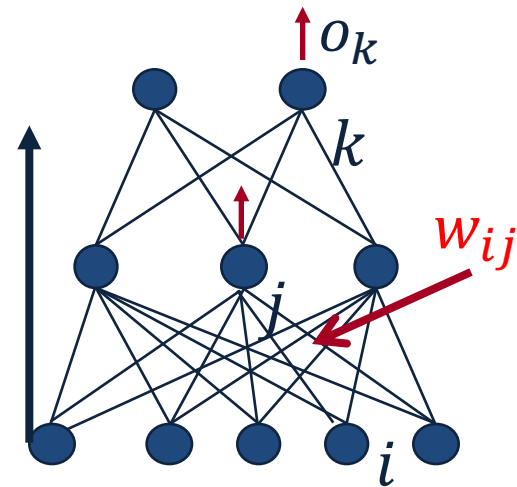$$\delta_j = \frac{\partial E_d}{\partial net_j} = (t_j - o_j)o_j(1 - o_j)$$



$o_j$

$j$

$w_{ij}$

$x_i$  $i$

# Derivation of Learning Rule (4)

- Weights of hidden units:

  - $w_{ij}$ Influences the output only through all the units whose direct input include $j$

$$\frac{\partial E_d}{\partial w_{ij}} = \frac{\partial E_d}{\partial \text{net}_j} \frac{\partial \text{net}_j}{\partial w_{ij}} =$$

$$\boxed{\text{net}_j = \sum w_{ij}.x_i}$$

$$= \frac{\partial E_d}{\partial \text{net}_j} x_i =$$

$$= \sum_{k \in parent(j)} \frac{\partial E_d}{\partial \text{net}_k} \frac{\partial \text{net}_k}{\partial \text{net}_j} x_i$$

$$= \sum_{k \in parent(j)} -\delta_k \frac{\partial \text{net}_k}{\partial \text{net}_j} x_i$$



$o_k$

$k$

$j$

$w_{ij}$

$i$

# Derivation of Learning Rule (5)

- **Weights of hidden units:**

  - $w_{ij}$ influences the output only through all the units whose direct input include $j$

$$\frac{\partial E_d}{\partial w_{ij}} = \sum_{k \in parent(j)} -\delta_k \frac{\partial \text{net}_k}{\partial \text{net}_j} x_i =$$

$$= \sum_{k \in parent(j)} -\delta_k \boxed{\frac{\partial \text{net}_k}{\partial o_j}} \boxed{\frac{\partial o_j}{\partial \text{net}_j}} x_i$$

$$= \sum_{k \in parent(j)} -\delta_k \boxed{w_{jk}} \boxed{o_j(1-o_j)} x_i$$



$o_k$

$k$

$w_{ij}$

$j$

$i$

# Derivation of Learning Rule (6)

- Weights of hidden units:
    - $w_{ij}$ is changed by:

$$\Delta w_{ij} = R \, o_j (1 - o_j) . \left( \sum_{k \in parent(j)} -\delta_k \, w_{jk} \right) x_i$$

$$= R \delta_j x_{ij}$$

- Where

$$\delta_j = o_j (1 - o_j) . \left( \sum_{k \in parent(j)} -\delta_k \, w_{jk} \right)$$

- First determine the error for the output units.
- Then, backpropagate this error layer by layer through the network, changing weights appropriately in each layer.

$o_k$

$k$

$w_{ij}$

$j$

$i$

# The Backpropagation Algorithm

- Create a fully connected three layer network. Initialize weights.
- Until all examples produce the correct output within $\epsilon$ (or other criteria)

For each example in the training set do:

1. Compute the network output for this example
2. Compute the error between the output and target value
$$\delta_k = (t_k - o_k)o_k(1 - o_k)$$
1. For each output unit $k$, compute error term

$$\delta_j = o_j(1 - o_j). \sum_{k \in downstream(j)} -\delta_k \ w_{jk}$$

1. For each hidden unit, compute error term:
$$\Delta w_{ij} = R\delta_j x_i$$
1. Update network weights with $\Delta w_{ij}$

End epoch

# More Hidden Layers

- The same algorithm holds for more hidden layers.

output

$h_1$

$h_2$

$h_3$

input

# Demo time!

- Link: [https://playground.tensorflow.org/](https://playground.tensorflow.org/)

# Comments on Training

- **No guarantee of convergence**; may **oscillate** or reach a local minima.

- In practice, many large networks can be trained on **large amounts of data** for realistic problems.

- **Many epochs** (tens of thousands) may be needed for adequate training. Large data sets may require many hours of CPU

- **Termination criteria**: Number of epochs; Threshold on training set error; No decrease in error; Increased error on a validation set.

- To **avoid local minima**: several trials with different random initial weights with majority or voting techniques

# Over-training Prevention

- Running too many epochs may **over-train** the network and result in over-fitting. (improved result on training, decrease in performance on test set)

- Keep an **hold-out validation** set and test accuracy after every epoch

- Maintain weights for best performing network on the validation set and return it when performance decreases significantly beyond that.

- To avoid losing training data to validation:
  - Use 10-fold cross-validation to determine the average number of epochs that optimizes validation performance
  - Train on the full data set using this many epochs to produce the final results

# Over-fitting prevention

- **Too few hidden units** prevent the system from adequately fitting the data and learning the concept.

- **Using too many hidden units** leads to over-fitting.

- Similar cross-validation method can be used to determine an appropriate number of hidden units. (general)

- Another approach to prevent over-fitting is weight-decay: all weights are multiplied by some fraction in (0,1) after every epoch.

  - Encourages smaller weights and less complex hypothesis

  - Equivalently: change Error function to include a term for the sum of the squares of the weights in the network. (general)

# Dropout training

- Proposed by (Hinton et al, 2012)



- Each time decide whether to delete one hidden unit with some probability $p$

# Dropout training



- Dropout of 50% of the hidden units and 20% of the input units (Hinton et al, 2012)

# Dropout training

- Model averaging effect
  - Among $2^H$ models, with shared parameters
    - *H*: number of units in the network
  - Only a few get trained
  - Much stronger than the known regularizer

- What about the input space?
  - Do the same thing!

# Input-Output Coding

- Appropriate coding of inputs and outputs can make learning problem easier and improve generalization.

- Encode each binary feature as a separate input unit;

- For mu~~~~~~~~~~~~~~~~~~~~t per value r~~~~~~~~~~~~~~~~~~tion in fewer u~~~~~~~~

  - Very~~~~~~~~~~~~~~~~~~~~~~~~~ of the input – real valued, dense representation.

> One way to do it, if you start with a collection of sparsely representation examples, is to use dimensionality reduction methods:
> - Your $m$ examples are represented as a $m \times 10^6$ matrix
> - Multiple it by a random matrix of size $10^6 \times 300$, say.
> - Random matrix: Normal(0,1)
> - New representation: $m \times 300$ dense rows

- For disjoint categorization problem, best to have one output unit for each category rather than encoding N categories into log N bits.

# Hidden Layer Representation

- Weight tuning procedure sets weights that define whatever hidden units representation is most effective at minimizing the error.

- Sometimes Backpropagation will define new hidden layer features that are not explicit in the input representation, but which capture properties of the input instances that are most relevant to learning the target function.

- Trained hidden units can be seen as newly constructed features that re-represent the examples so that they are linearly separable

# Gradient Checks are useful!

- Allow you to know that there are no bugs in your neural network implementation!

> - Implement your gradient
> - Implement a finite difference computation by looping through the parameters of your network, adding and subtracting a small epsilon (~10^-4) and estimate derivatives
>
> $$f'(\theta) \approx \frac{f(\theta^+) - f(\theta^-)}{2\epsilon} \qquad \theta^\pm = \theta \ \pm \epsilon$$
>
> - Compare the two and make sure they are almost the same

# Auto-associative Network

- An auto-associative network trained with 8 inputs, 3 hidden units and 8 output nodes, where the output must reproduce the input.

- When trained with vectors with only one bit on

| INPUT | HIDDEN | | |
|---|---|---|---|
| 1 0 0 0 0 0 0 0 | .89 | .40 | 0.8 |
| 0 1 0 0 0 0 0 0 | .97 | .99 | .71 |
| .... | | | |
| 0 0 0 0 0 0 0 1 | .01 | .11 | .88 |

1 0 0 0 1 0 0 0

1 0 0 0 1 0 0 0

- Learned the standard 3-bit encoding for the 8 bit vectors.

- Illustrates also data compression aspects of learning

# Sparse Auto-encoder

- Encoding: $\qquad y = f(Wx + b)$
- Decoding: $\qquad \widehat{x} = g(W'y + b')$
  - Goal: perfect reconstruction of input vector $x$, by the output $\widehat{x}$
    - Where $\theta = \{W, W'\}$
  - Minimize an error function $l(\widehat{x}, x)$
    - For example:
      $$l(\widehat{x}, x) = \|\widehat{x} - x\|^2$$
  - And regularize it
    $$\min_\theta \sum_x l(\widehat{x}, x) + \sum_i |w_i|$$
- After optimization drop the reconstruction layer and add a new layer

# Stacking Auto-encoder

- Add a new layer, and a reconstruction layer for it.
- And try to tune its parameters such that
- And continue this for each layer

# Beyond supervised learning

- So far what we had was purely **supervised**.
  - Initialize parameters randomly
  - Train in supervised mode typically, using backprop
  - Used in most practical systems (e.g. speech and image recognition)
- Unsupervised, layer-wise pre-training + supervised classifier on top
  - Train each layer unsupervised, one after the other
  - Train a supervised classifier on top, keeping the other layers fixed
  - Good when very few labeled samples are available
- Unsupervised, layer-wise pre-training + supervised fine-tuning
  - Train each layer unsupervised, one after the other
  - Add a classifier layer, and retrain the whole thing supervised
  - Good when label set is poor (e.g. pedestrian detection)

We won't talk about unsupervised pre-training here. But it's good to have this in mind, since it is an active topic of research.

# NN-2

# Recap: Multi-Layer Perceptrons

- Multi-layer network
  - A global approximator
  - Different rules for training it
- The Back-propagation
  - Forward step
  - Back propagation of errors



- Congrats! Now you know the one of the important algorithms in neural networks!
- Today:
  - Convolutional Neural Networks
  - Recurrent Neural Networks

# Receptive Fields

- The **receptive field** of an individual sensory neuron is the particular region of the sensory space (e.g., the body surface, or the retina) in which a stimulus will trigger the firing of that neuron.

  - In the auditory system, receptive fields can correspond to volumes in auditory space

- Designing "proper" receptive fields for the input Neurons is a significant challenge.

- Consider a task with image inputs

  - Receptive fields should give expressive features from the raw input to the system

  - How would you design the receptive fields for this problem?

- A **fully connected layer**:
    - Example:
        - 100x100 images
        - 1000 units in the input
    - Problems:
        - 10^7 edges!
        - Spatial correlations lost!
        - Variables sized inputs.

Input layer

Slide Credit: Marc'Aurelio Ranzato

- Consider a task with image inputs:
- A **locally connected layer**:
  - Example:
    - 100x100 images
    - 1000 units in the input
    - Filter size: 10x10
  - Local correlations preserved!
  - Problems:
    - $10^5$ edges
    - This parameterization is good when input image is registered (e.g., face recognition).
    - Variable sized inputs, again.

Input layer

# Convolutional Layer

- **A solution:**
  - **Filters** to capture different patterns in the input space.
    - **Share** parameters across different locations (assuming input is stationary)
    - **Convolutions** with learned filters
  - Filters will be **learned** during training.
  - The issue of variable-sized inputs will be resolved with a **pooling** layer.

  So what is a convolution?

  Input layer

# Convolution Operator

- Convolution operator: *
  - takes two functions and gives another function
- One dimension:

$$(x * h)(t) = \int x(\tau)h(t - \tau)d\tau$$

$$(x * h)[n] = \sum_m x[m]h[n - m]$$

"Convolution" is very similar to "cross-correlation", except that in convolution one of the functions is flipped.

# Convolution Operator (2)

- Convolution in two dimension:
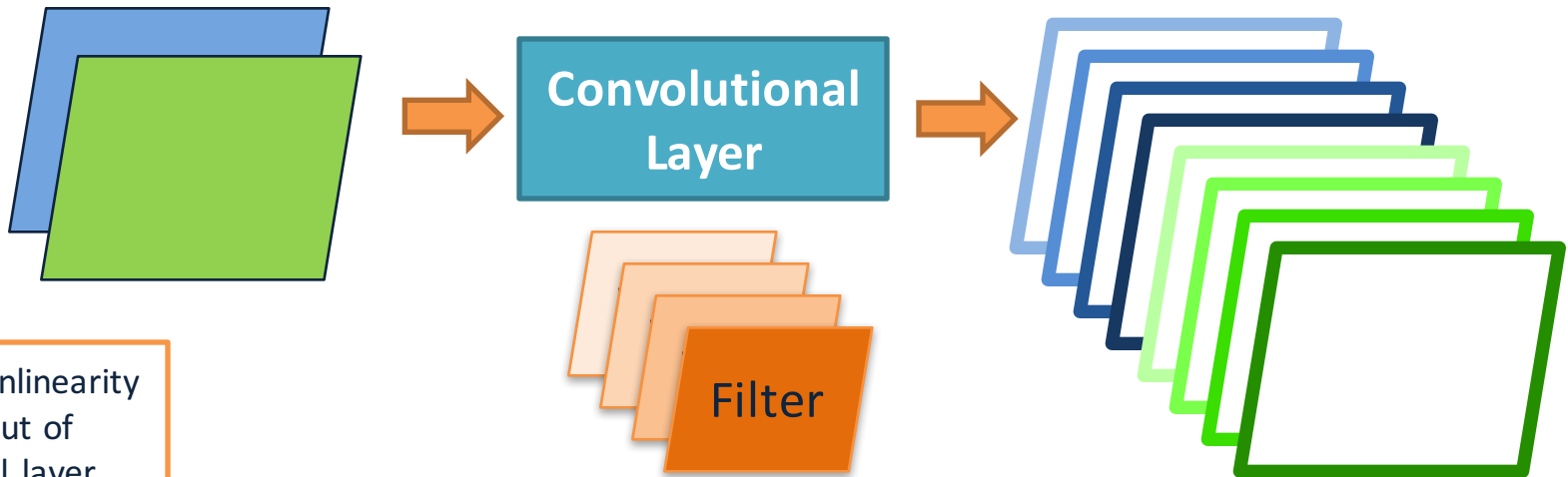  - The same idea: flip one matrix and slide it on the other matrix
  - Example: Sharpen kernel:

|  |  |  |
|---|---|---|
| 0 | -1 | 0 |
| -1 | 5 | -1 |
| 0 | -1 | 0 |



input image

$173 \times 0 + 255 \times -1 + 254 \times 0$

$+ 54 \times -1 + 140 \times 5 + 213 \times -1$

$+ 46 \times 0 + 81 \times -1 + 118 \times 0$

$= 97$

output image

# Convolution Operator (3)

■ Convolution in two dimension:

  ❑ The same idea: flip one matrix and slide it on the other matrix

# Complexity of Convolution

- Complexity of convolution operator is $n\log(n)$, for $n$ inputs.
  - Uses Fast-Fourier-Transform (FFT)
- For two-dimension, each convolution takes $MN\log(MN)$ time, where the size of input is $MN$.



Slide Credit: Marc'Aurelio Ranzato

# Convolutional Layer

- The convolution of the **input (vector/matrix)** with weights **(vector/matrix)** results in a **response vector/matrix**.

- We can have **multiple filters** in each convolutional layer, each producing an output.

- If it is an intermediate layer, it can have **multiple inputs**!

One can add nonlinearity at the output of convolutional layer

**Convolutional Layer**

Filter

# Pooling Layer

- How to handle variable sized inputs?
  - A layer which reduces inputs of different size, to a fixed size.
  - **Pooling**



Slide Credit: Marc'Aurelio Ranzato

# Pooling Layer

- How to handle variable sized inputs?

  - A layer which reduces inputs of different size, to a fixed size.

  - **Pooling**

  - Different variations

    - Max pooling

    $$h_i[n] = \max_{i \in N(n)} \tilde{h}[i]$$

    - Average pooling

    $$h_i[n] = \frac{1}{n} \sum_{i \in N(n)} \tilde{h}[i]$$

    - L2-pooling

    $$h_i[n] = \frac{1}{n} \sqrt{\sum_{i \in N(n)} \tilde{h}^2[i]}$$

    - etc

# Convolutional Nets

- One stage structure:



- Whole system:

# Training a ConvNet

- The same procedure from Back-propagation applies here.
  - Remember in backprop we started from the error terms in the last stage, and passed them back to the previous layers, one by one.

- Back-prop for the pooling layer:
  - Consider, for example, the case of "max" pooling.
  - This layer only routes the gradient to the input that has the highest value in the forward pass.
  - Hence, during the forward pass of a pooling layer it is common to keep track of the index of the max activation (sometimes also called *the switches*) so that gradient routing is efficient during backpropagation.
  - Therefore we have: $\delta = \frac{\partial E_d}{\partial y_i}$

# Training a ConvNet

- Back-prop for the convolutional layer:

We derive the update rules for a 1D convolution, but the idea is the same for bigger dimensions.

$$\tilde{y} = w * x \iff \tilde{y}_i = \sum_{a=0}^{m-1} w_a\, x_{i-a} = \sum_{a=0}^{m-1} w_{i-a}\, x_a \qquad \text{The convolution}$$
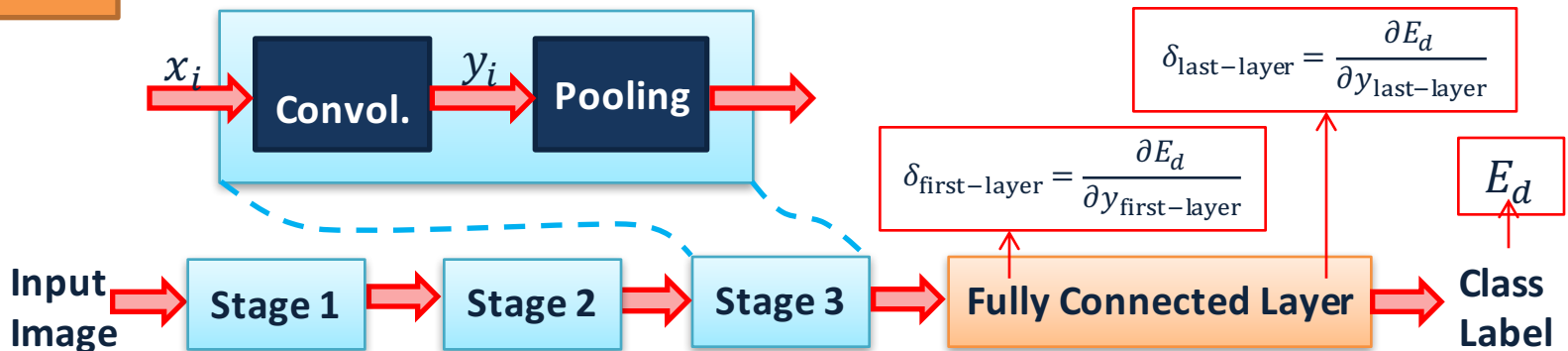
$$y = f(\tilde{y}) \iff y_i = f(\tilde{y}_i) \quad \forall i \qquad \text{A differentiable nonlinearity}$$

$$\frac{\partial E_d}{\partial w_a} = \sum_{i=0}^{m-1} \frac{\partial E_d}{\partial \tilde{y}_i} \frac{\partial \tilde{y}_i}{\partial w_a} = \sum_{i=0}^{m-1} \frac{\partial E_d}{\partial \tilde{y}_i} x_{i-a}$$
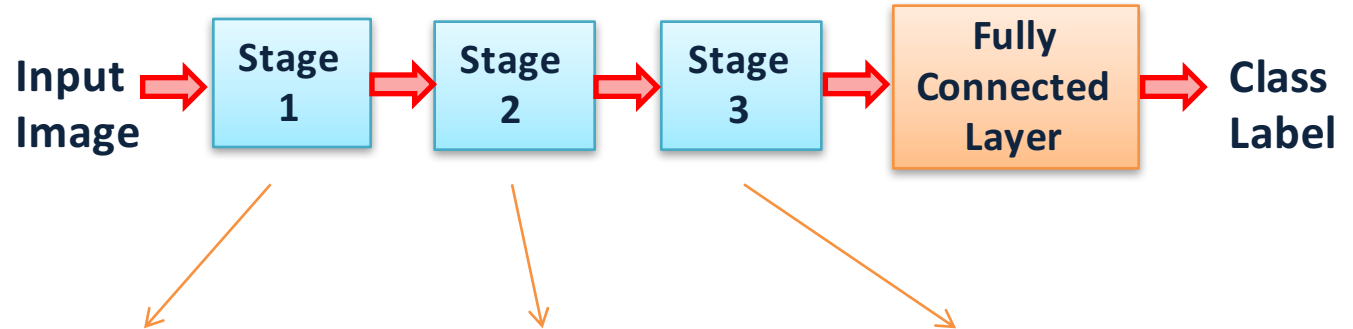
$$\frac{\partial E_d}{\partial \tilde{y}_i} = \frac{\partial E_d}{\partial y_i} \frac{\partial y_i}{\partial \tilde{y}_i} = \frac{\partial E_d}{\partial y_i} f'(\tilde{y}) \qquad \text{Now we have everything in this layer to update the filter}$$

$$\delta = \frac{\partial E_d}{\partial x_a} = \sum_{i=0}^{m-1} \frac{\partial E_d}{\partial \tilde{y}_i} \frac{\partial \tilde{y}_i}{\partial x_a} = \sum_{i=0}^{m-1} \frac{\partial E_d}{\partial \tilde{y}_i} w_{i-a} \qquad \text{We need to pass the gradient to the previous layer}$$

Now we can repeat this for each stage of ConvNet.

$x_i$ → **Convol.** → $y_i$ → **Pooling** →

$$\delta_{\text{last-layer}} = \frac{\partial E_d}{\partial y_{\text{last-layer}}}$$

$$\delta_{\text{first-layer}} = \frac{\partial E_d}{\partial y_{\text{first-layer}}}$$

$E_d$

**Input Image** → **Stage 1** → **Stage 2** → **Stage 3** → **Fully Connected Layer** → **Class Label**

# Convolutional Nets
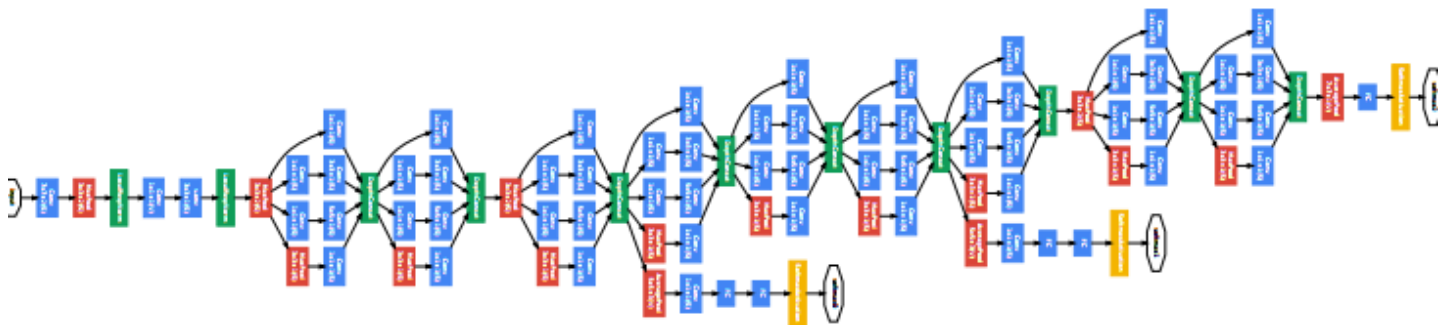


Feature visualization of convolutional net trained on ImageNet from [Zeiler & Fergus 2013]

# Demo (Teachable Machines)

https://teachablemachine.withgoogle.com/
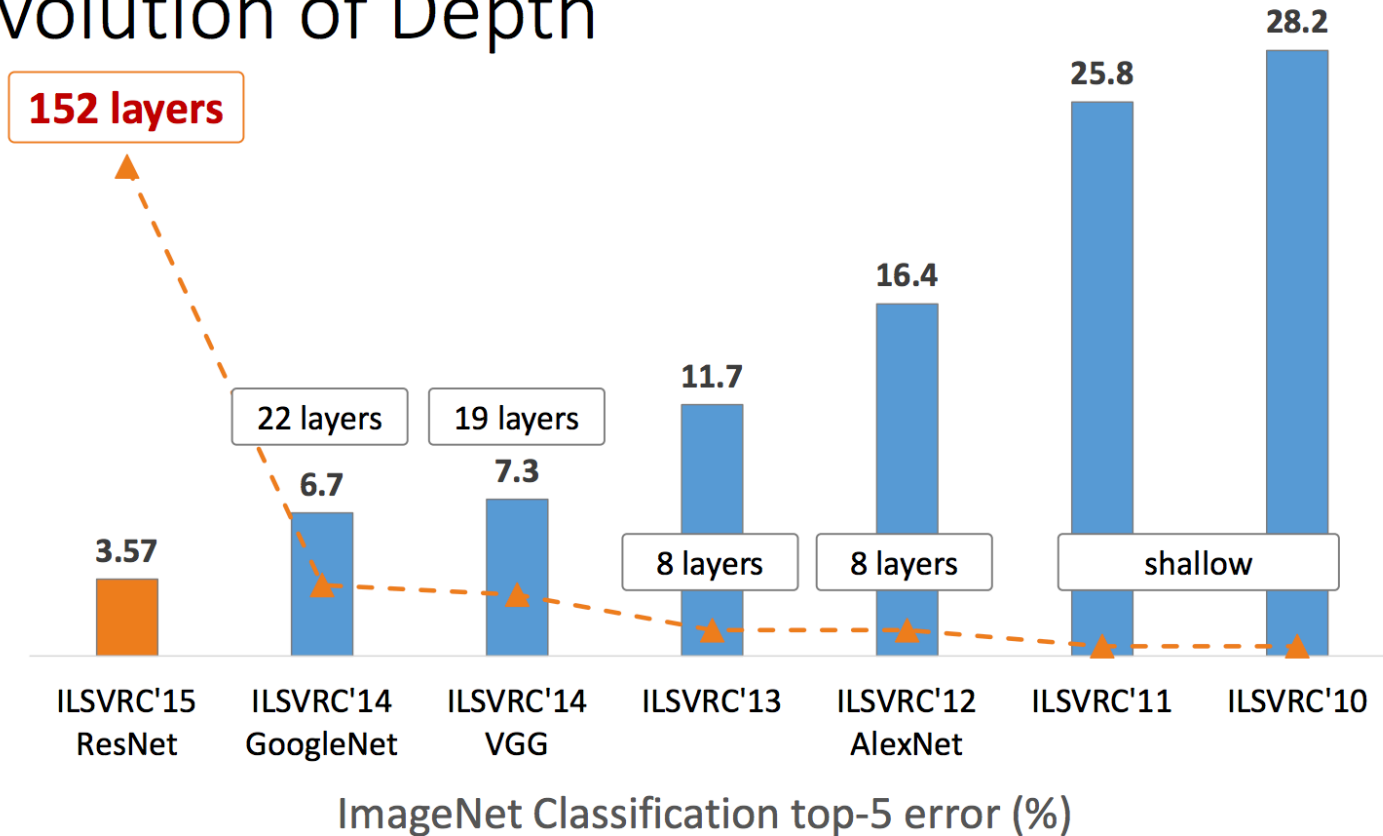
# ConvNet roots

- **Fukushima, 1980s** designed network with same basic structure but did not train by backpropagation.

- The first successful applications of **Convolutional Networks** by Yann LeCun in 1990's (LeNet)
  - Was used to read zip codes, digits, etc.

- Many variants nowadays, but the core idea is the same
  - Example: a system developed in Google (GoogLeNet)
    - Compute different filters
    - Compose one big vector from all of them
    - Layer this iteratively

See more: http://arxiv.org/pdf/1409.4842v1.pdf

# Depth matters

## Revolution of Depth



ImageNet Classification top-5 error (%)

Slide from [Kaiming He 2015]

# Vanishing/exploding gradients

Gradient can become very **small or very large quickly**, and the locality assumption of gradient descent breaks down (Vanishing gradient) [Bengio et al 1994]

- Vanishing gradients are quite prevalent and a serious issue.
- A real example
  - Training a feed-forward network
  - y-axis: sum of the gradient norms
  - Earlier layers have exponentially smaller sum of gradient norms
  - This will make training earlier layers much slower.

# Vanishing/exploding gradients

- In architectures with many layers (*e.g.* >10) the gradients can easily explode or vanish.

- Many methods proposed for reduce the effect of vanishing gradients; although it is still a problem
  - Introduce shorter path between long connections
  - Abandon stochastic gradient descent in favor of a much more sophisticated Hessian-Free (HF) optimization
  - Clip gradients with bigger sizes:

Defnne $g = \frac{\partial E}{\partial W}$

If $\|g\| \geq threshold$ then
$$g \leftarrow \frac{threshold}{\|g\|} g$$

# Practical Tips

- Before large scale experiments, test on a small subset of the data and check the error should go to zero.
  - Overfitting on small training
- Visualize features (feature maps need to be uncorrelated) and have high variance
- Bad training: many hidden units ignore the input and/or exhibit strong correlations.



| GOOD | BAD | BAD | BAD |
| --- | --- | --- | --- |
| | too noisy | too correlated | lack structure |

Figure Credit: Marc'Aurelio Ranzato

# Debugging

- Training diverges:
  - Learning rate may be too large → decrease learning rate
  - BackProp is buggy → numerical gradient checking

- Loss is minimized but accuracy is low
  - Check loss function: Is it appropriate for the task you want to solve? Does it have degenerate solutions?

- NN is underperforming / under-fitting
  - Compute number of parameters → if too small, make network larger

- NN is too slow
  - Compute number of parameters → Use distributed framework, use GPU, make network smaller

> Many of these points apply to many machine learning models, no just neural networks.

# CNN for vector inputs

- Let's study another variant of CNN for language
  - Example: sentence classification (say spam or not spam)
- First step: represent each word with a vector in $\mathbb{R}^d$

This     is     not     a     spam

O O O O O O O    O O O O O O O    O O O O O O O    O O O O O O O    O O O O O O O

Concatenate the vectors

O O O O O O O   O O O O O O O   O O O O O O O   O O O O O O O   O O O O O O O

- Now we can assume that the input to the system is a vector $\mathbb{R}^{dl}$
  - Where the input sentence has length $l$ ($l = 5$ in our example )
  - Each word vector's length $d$ ($d = 7$ in our example )

# Convolutional Layer on vectors

- Think about a single convolutional layer
  - A bunch of **vector** filters
    - Each defined in $\mathbb{R}^{dh}$
      - Where $h$ is the number of the words the filter covers
      - Size of the word vector $d$   ⊙⊙⊙⊙⊙⊙⊙ ⊙⊙⊙⊙⊙⊙⊙
  - Find its (modified) convolution with the input vector

$$c_1 = f(w.x_{1:h}) \quad c_2 = f(w.x_{h+1:2h}) \quad c_3 = f(w.x_{2h+1:3h}) \quad f(w.x_{3h+1:4h})$$

  - Result of the convolution with the filter

$$c = [c_1, \ldots, c_{n-h+1}]$$

  - Convolution with a filter that spans 2 words, is operating on all of the bi-grams (vectors of two consecutive word, concatenated): "this is", "is not", "not a", "a spam".
  - Regardless of whether it is grammatical (not appealing linguistically)

# Convolutional Layer on vectors

Get word vectors for each words

This     is     not     a     spam

○○○○○○○  ○○○○○○○  ○○○○○○○  ○○○○○○○  ○○○○○○○

Concatenate vectors

○○○○○○○ ○○○○○○○ ○○○○○○○ ○○○○○○○ ○○○○○○○

*

Perform convolution with each filter

○○○○○○○○
○○○○○○○ ○○○○○○○
○○○○○○○ ○○○○○○○○
○○○○○○○ ○○○○○○○ ○○○○○○○
○○○○○○○ ○○○○○○○ ○○○○○○○

Filter bank

How are we going to handle the **variable sized** response vectors?
**Pooling!**

#of filters {

○○○○○
○○○○○
○○○○○
○○○○
○○○

Set of response vectors

**#words** - #length of filter + 1

# Convolutional Layer on vectors

This      is      not      a      spam

Get word vectors for each words

Concatenate vectors

*

Perform convolution with each filter

Filter bank

Pooling on filter responses

Some choices for pooling: *k-max, mean*, etc

#of filters

**#words** - #length of filter + 1

- Now we can pass the fixed-sized vector to a logistic unit (softmax), or give it to multi-layer network (last session)

# Recurrent Neural Networks

- Prediction on chain-like input:
  - Example: POS tagging words of a sentence

| This | → | is | → | a | → | sample | → | sentence |

DT    VBZ    DT    NN    NN

  - Issues :
    - Structure in the output:  There is connections between labels
    - Interdependence between elements of the inputs: The final decision is based on an intricate interdependence of the words on each other.
    - Variable size inputs:  e.g. sentences differ in size

- How would you go about solving this task?

# Recurrent Neural Networks

- Infinite uses of finite structure

Input

X0    X1    X2    X3

H0 → H1 → H2 → H3 →

Hidden state representation

Y0    Y1    Y2    Y3

Output

# Recurrent Neural Networks

- A chain RNN:
  - Has a chain-like structure
  - Each input is replaced with its vector representation $x_t$
  - Hidden (memory) unit $h_t$ contain information about previous inputs and previous hidden units $h_{t-1}, h_{t-2}$, etc
    - Computed from the past memory and current word. It summarizes the sentence up to that time.

# Recurrent Neural Networks

- A popular way of formalizing it:
$$h_t = f(W_h h_{t-1} + W_i x_t)$$
  - Where $f$ is a nonlinear, differentiable (why?) function.

- Outputs?
  - Many options; depending on problem and computational resource

# Recurrent Neural Networks

▪ Prediction for $x_t$, with $h_t$:    $y_t = \mathrm{softmax}(W_o h_t)$



▪ Some inherent issues with RNNs:

  ▪ Recurrent neural nets cannot capture phrases without prefix context

  ▪ They often capture too much of last words in final vector

    ▪ A slightly more sophisticated solution: Long Short-Term Memory (LSTM) units

# Recurrent Neural Networks

- **Multi-layer feed-forward NN: <span style="color:red">DAG</span>**

  - Just computes a fixed sequence of

  non-linear learned transformations to convert an input patter into an output pattern

- **Recurrent Neural Network: <span style="color:red">Digraph</span>**

  - Has cycles.

  - Cycle can act as a memory;

  - The hidden state of a recurrent net can carry along  information about a "potentially" unbounded number of previous inputs.

  - They can model sequential data in a much more natural way.

# Equivalence between RNN and Feed-forward NN

- Assume that there is a time delay of 1 in using each connection.

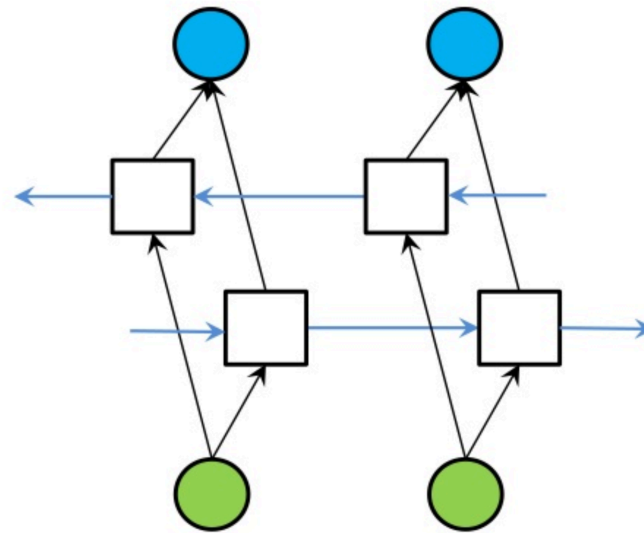- The recurrent net is just a layered net that keeps reusing the same weights.



Slide Credit: Geoff Hinton

# Bi-directional RNN

- One of the issues with RNN:
  - Hidden variables capture only one side context
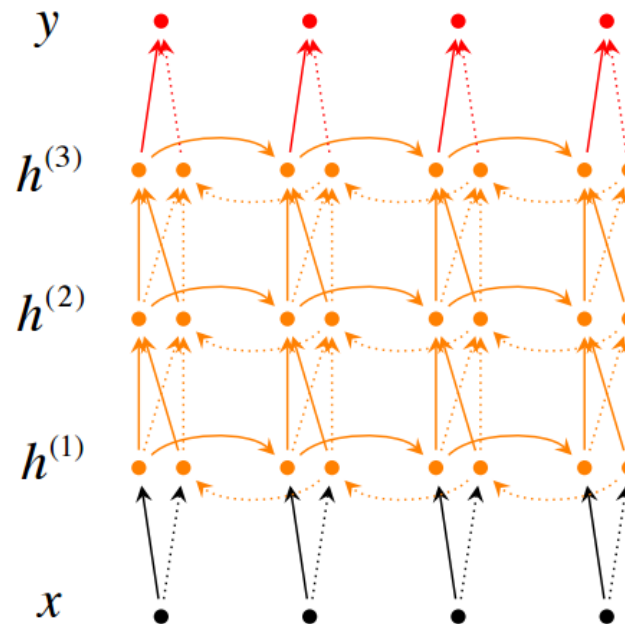- A bi-directional structure



RNN



Bi-directional RNN

# Stack of bi-directional networks

- Use the same idea and make your model further complicated:

# Training RNNs

- How to train such model?

  - Generalize the same ideas from back-propagation

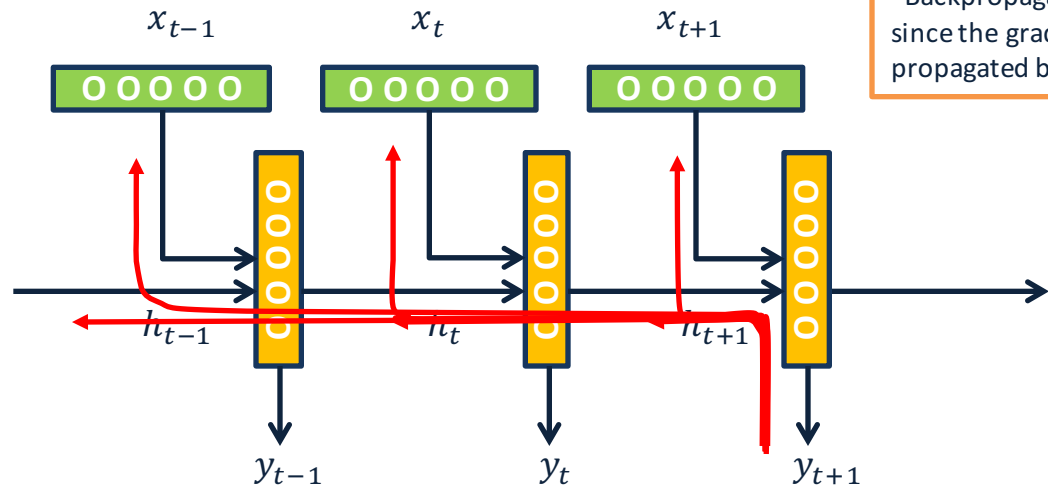- Total output error: $E(\vec{y}, \vec{t}) = \sum_{t=1}^{T} E_t(y_t, t_t)$

Parameters?
$W_o, W_i, W_h$ +
vectors for
input

$$\frac{\partial E}{\partial W} = \sum_{t=1}^{T} \frac{\partial E_t}{\partial W}$$

$$\frac{\partial E_t}{\partial W} = \sum_{t=1}^{T} \frac{\partial E_t}{\partial y_t} \frac{\partial y_t}{\partial h_t} \boxed{\frac{\partial h_t}{\partial h_{t-k}}} \frac{\partial h_{t-k}}{\partial W}$$

Reminder:
$y_t = \text{softmax}(W_o h_t)$
$h_t = f(W_h h_{t-1} + W_i x_t)$

This sometimes is called
"Backpropagation Through Time",
since the gradients are
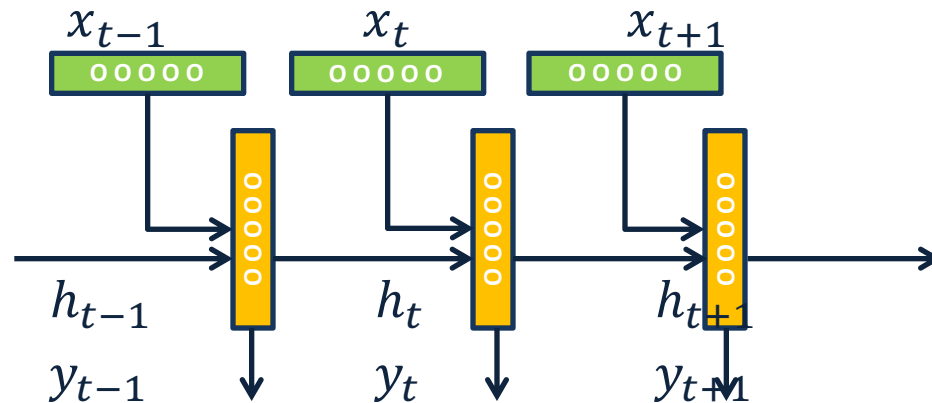propagated back through time.

# Recurrent Neural Network

$$\frac{\partial E}{\partial W} = \sum_{t=1}^{T} \frac{\partial E_t}{\partial y_t} \frac{\partial y_t}{\partial h_t} \boxed{\frac{\partial h_t}{\partial h_{t-k}}} \frac{\partial h_{t-k}}{\partial W}$$

Reminder:
$$y_t = \text{softmax}(W_o h_t)$$
$$h_t = f(W_h h_{t-1} + W_i x_t)$$

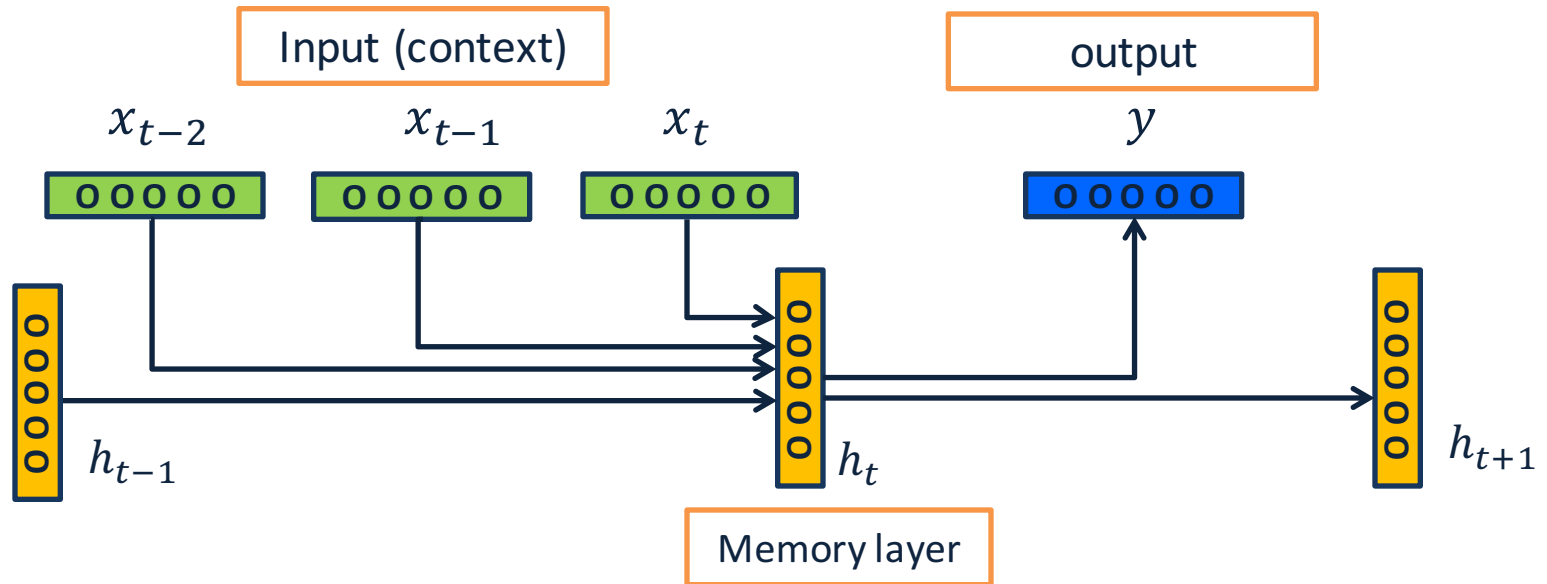$$\frac{\partial h_t}{\partial h_{t-1}} = W_h \text{diag}[f'(W_h h_{t-1} + W_i x_t)]$$

$$\text{diag}[a_1, ..., a_n] = \begin{bmatrix} a_1 & 0 & 0 \\ 0 & \ddots & 0 \\ 0 & 0 & a_n \end{bmatrix}$$

$$\frac{\partial h_t}{\partial h_{t-k}} = \prod_{j=t-k+1}^{t} \frac{\partial h_j}{\partial h_{j-1}} = \prod_{j=t-k+1}^{t} W_h \text{diag}[f'(W_h h_{t-1} + W_i x_t)]$$
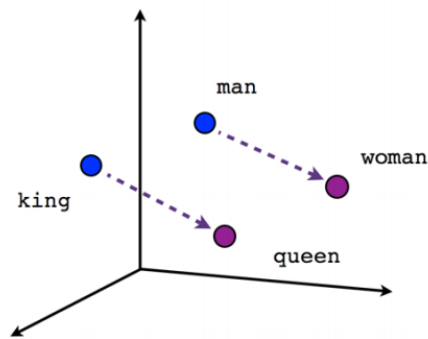
# Unsupervised RNNs

- **What to put here?**
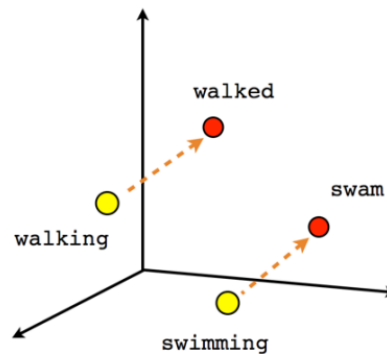  - He was locked up after he _____ .



- **Note that:**
  - This is unsupervised; you can use tons of data to train this.
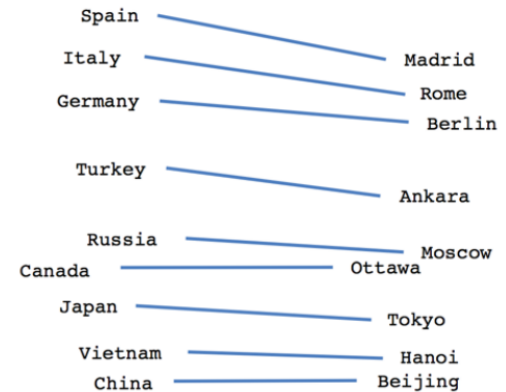  - While training the model, we train the word representations too.

# Unsupervised RNNs

- **This would result in word representations**
  - that convey information about their co-occurrence
  - Or some form of weak "semantic" similarity



Male-Female

Verb tense

Country-Capital

- A big part of progress (past 5-10 years) is partly due to discovering better ways create unsupervised context-sensitive representations