

[70240413 Statistical Machine Learning, Spring, 2018]

Deep Learning **(deep neural nets)**

Jun Zhu

`dcszj@mail.tsinghua.edu.cn`

`http://bigml.cs.tsinghua.edu.cn/~jun`

State Key Lab of Intelligent Technology & Systems

Tsinghua University

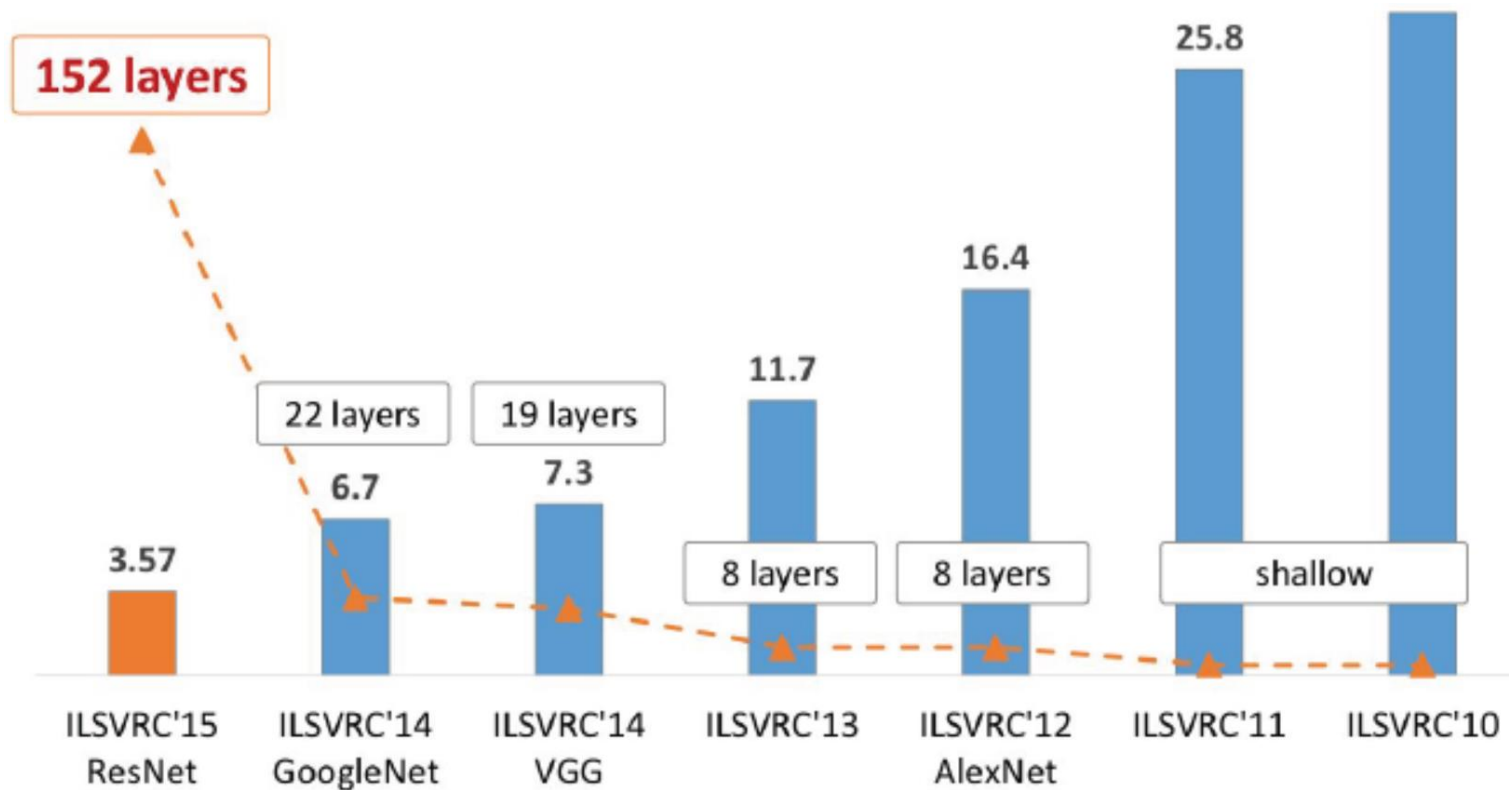
March 27, 2018

Why going deep?

- ◆ Data are often high-dimensional.
- ◆ There is a huge amount of **structure** in the data, but the structure is too complicated to be represented by a simple model.
- ◆ Insufficient depth can require more **computational elements** than architectures whose depth matches the task.
- ◆ Deep nets provide simpler but more descriptive models of many problems.

Resolution in Image Classification

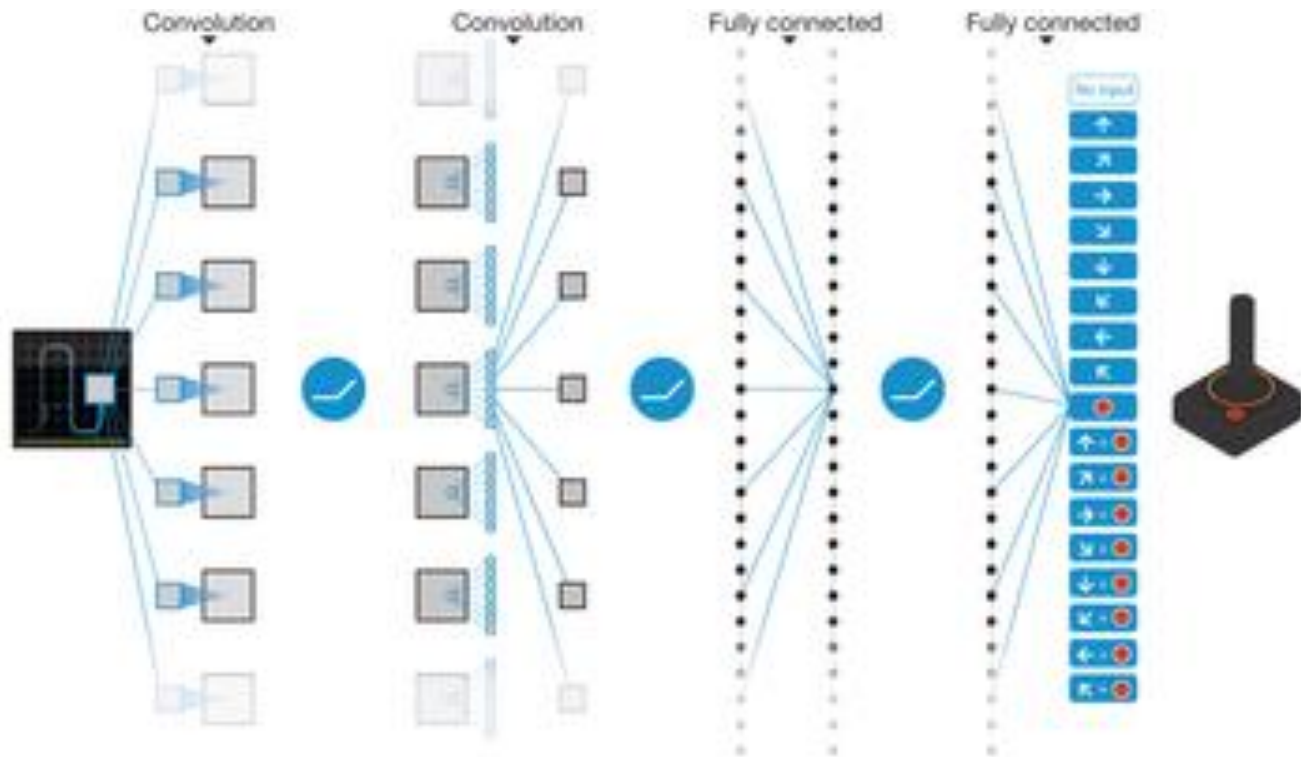
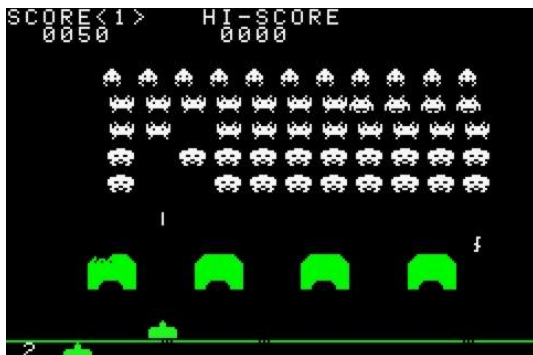
◆ ImageNet Large-Scale Visual Recognition Challenge (ILSVRC)



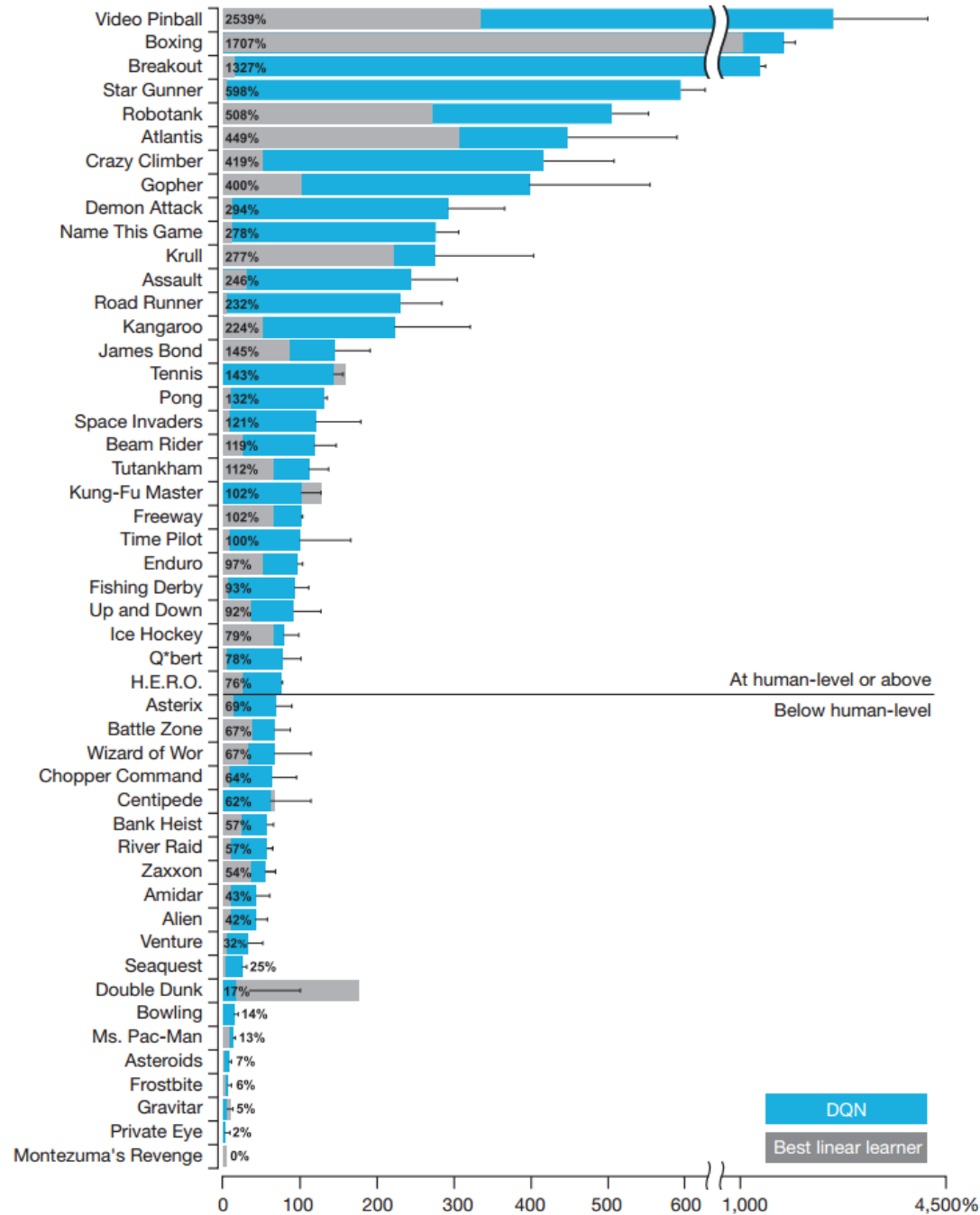


Human-Level Control via Deep RL

- ◆ Deep Q-network with human-level performance on Atari games

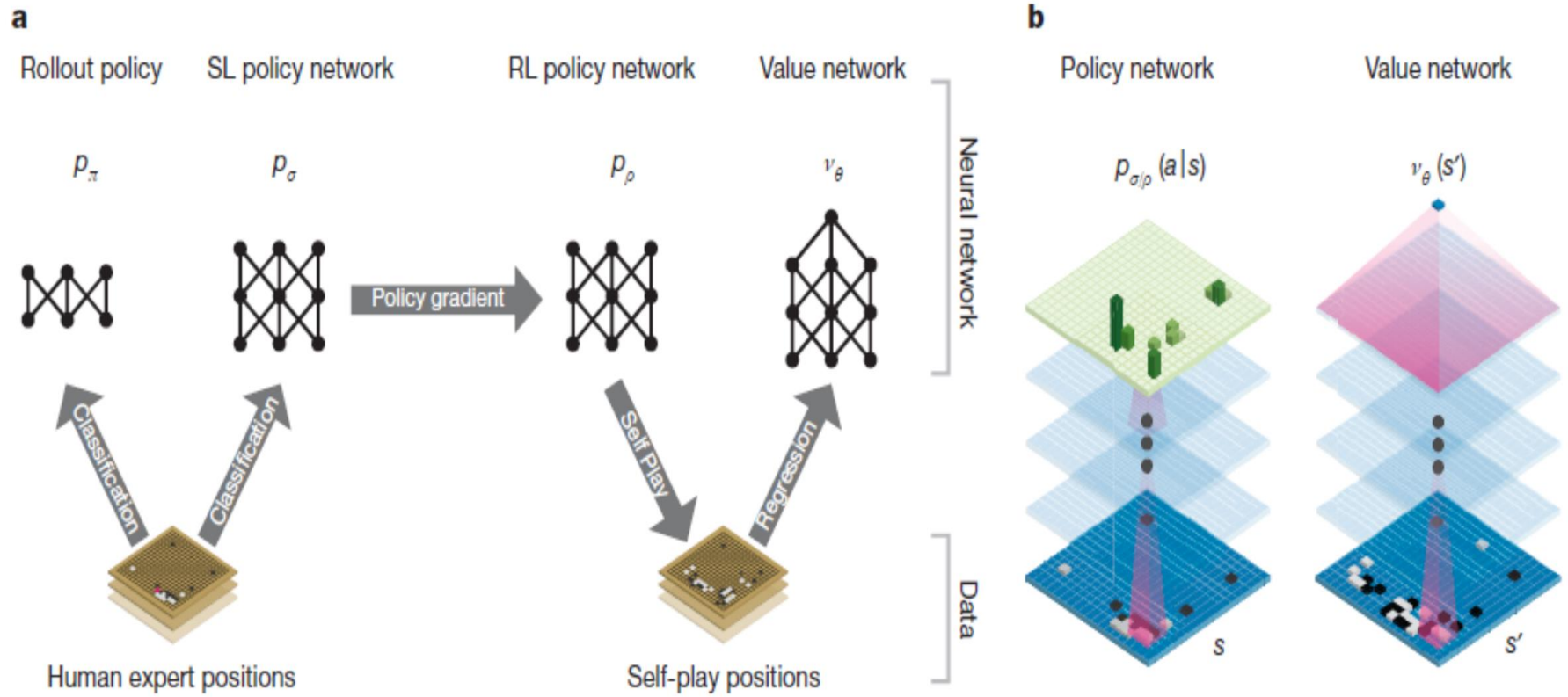


[Mnih et al., Nature 518, 529–533, 2015]



AlphaGo

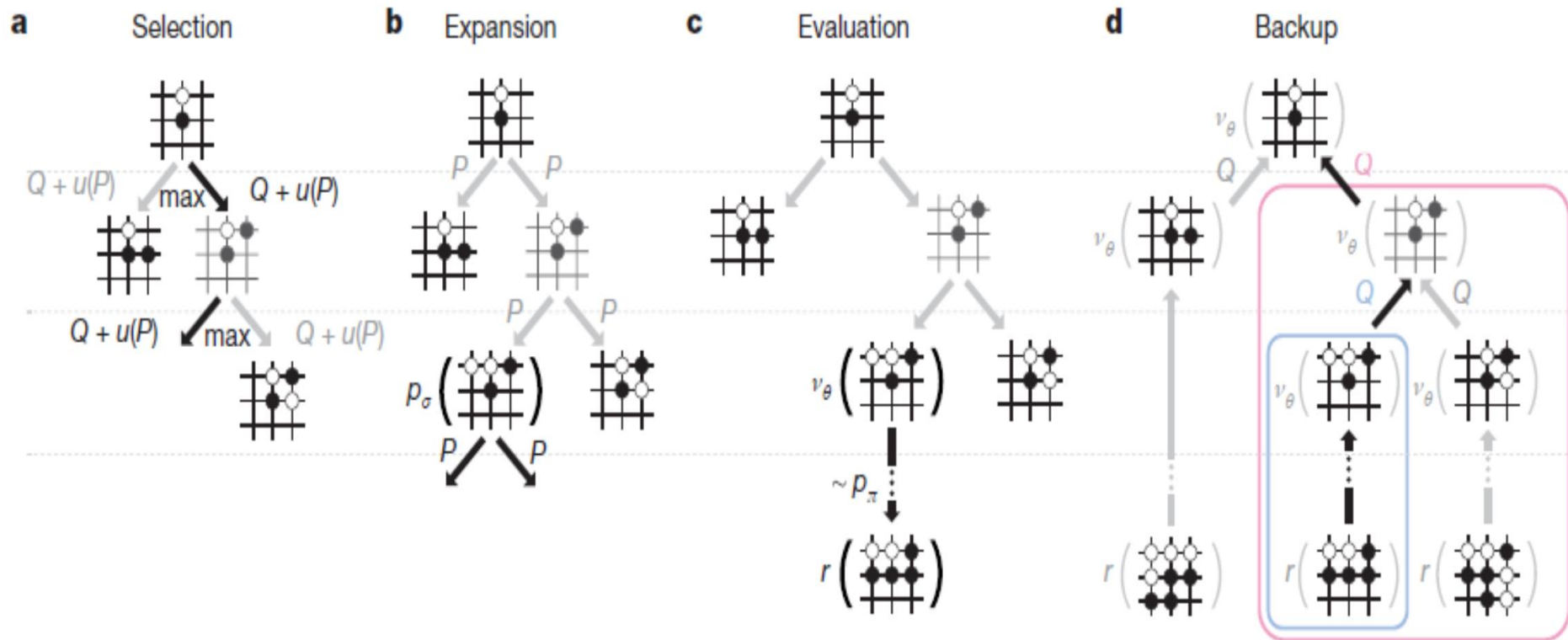
◆ Neural network training pipeline and architecture





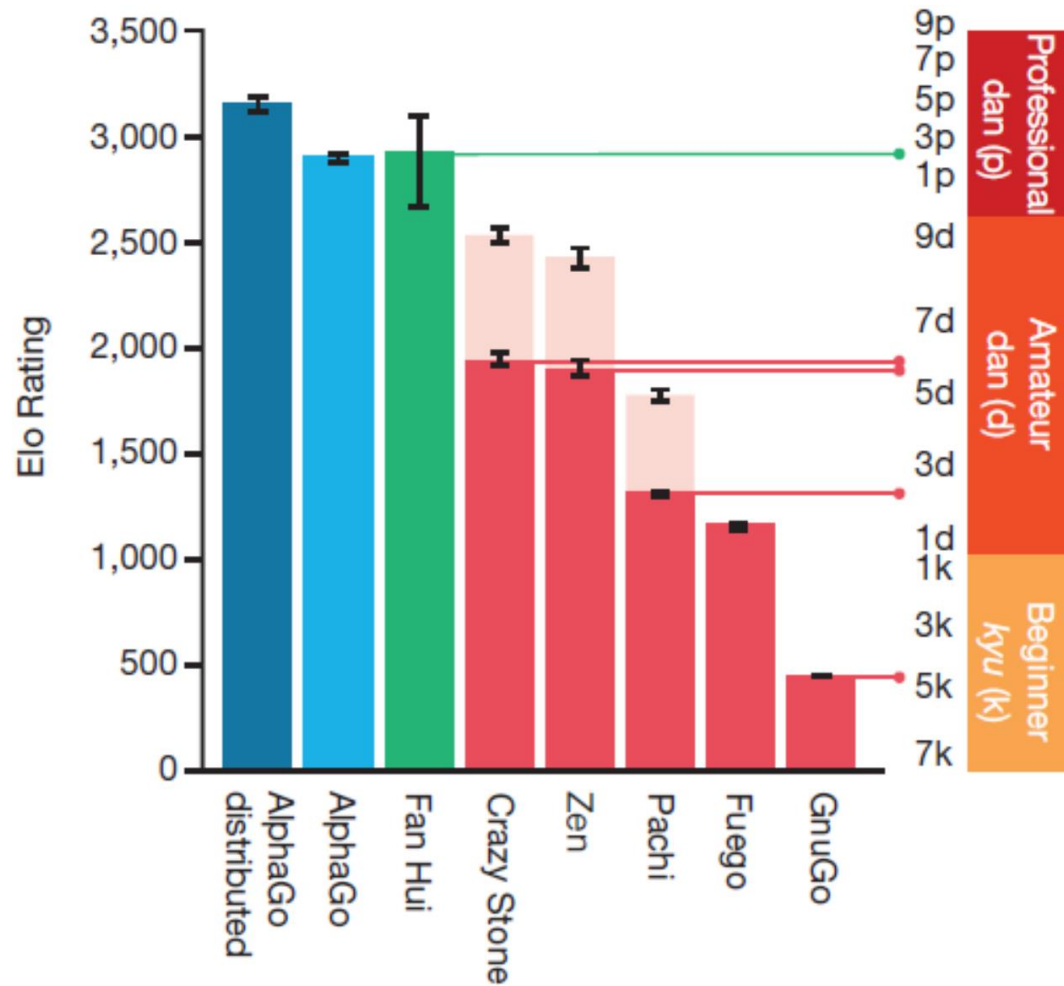
AlphaGo

◆ Monte Carlo tree search

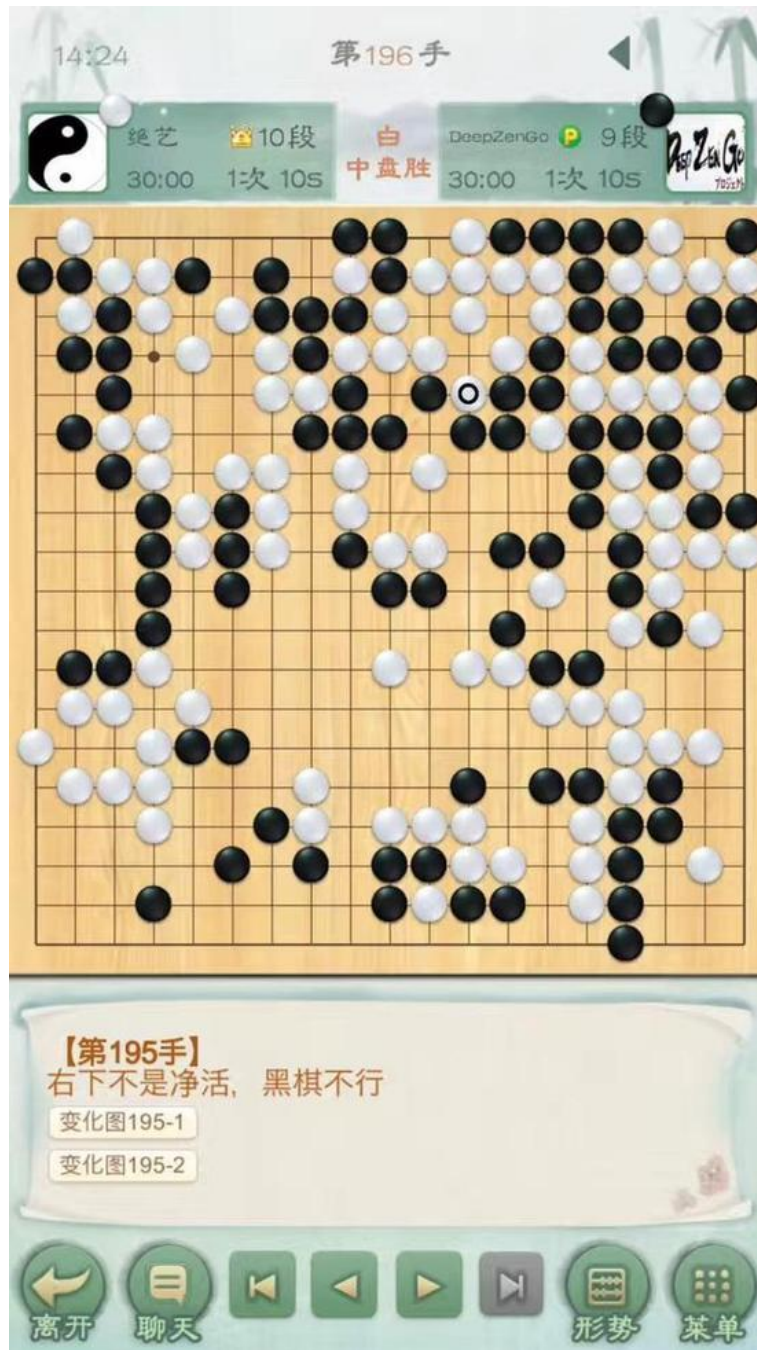




AlphaGo

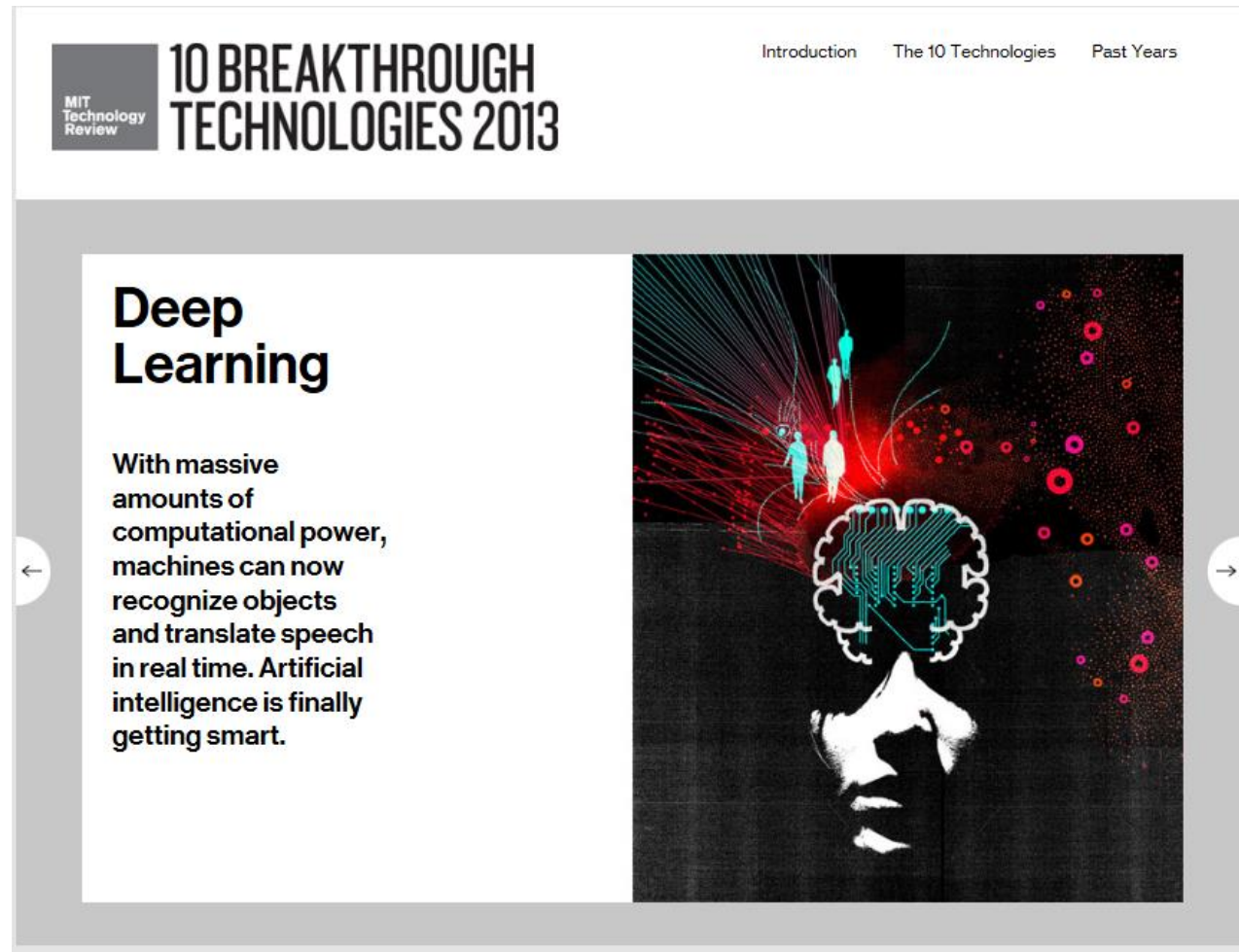


◆ Tencent FineArt





MIT 10 Breakthrough Tech 2013



<http://www.technologyreview.com/featuredstory/513696/deep-learning/>

Deep Learning in industry



Driverless car



Face identification



Speech recognition



Web search

...



...



History of neural networks



Pitts



McCulloch



Rosenblatt



Minsky



Papert



Ackley

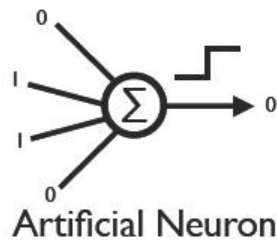


Hinton

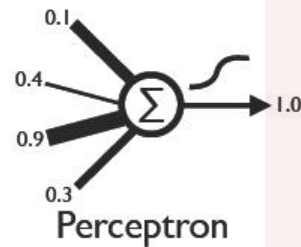


Sejnowski

1943



1960

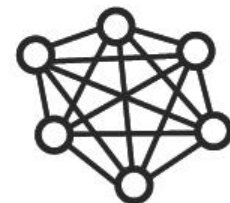


1969



Perceptrons

1985



Boltzmann Machine

History of neural networks



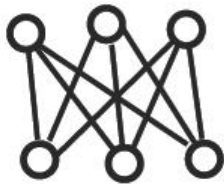
Smolensky



Hinton

Hinton et al.

1986



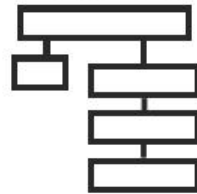
Harmoniums
(Restricted Boltzmann Machine)

2002



Contrastive
Divergence

2006

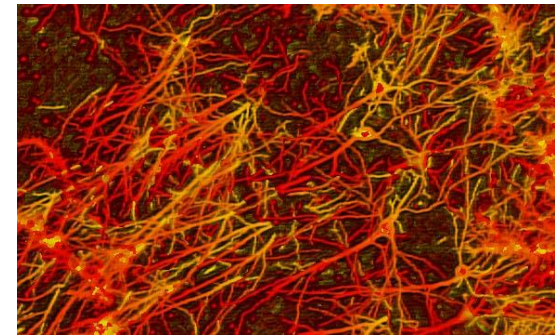
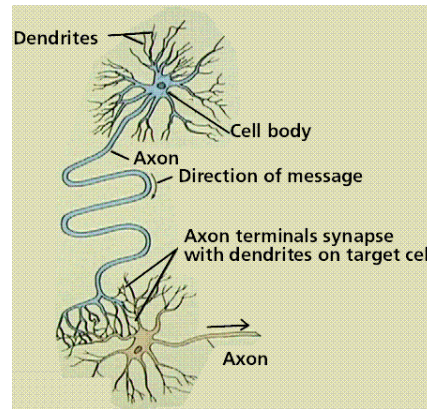
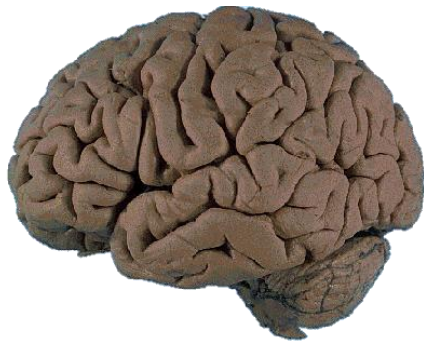


Deep Belief
Networks

Deep Learning Models



How the human brain learns?



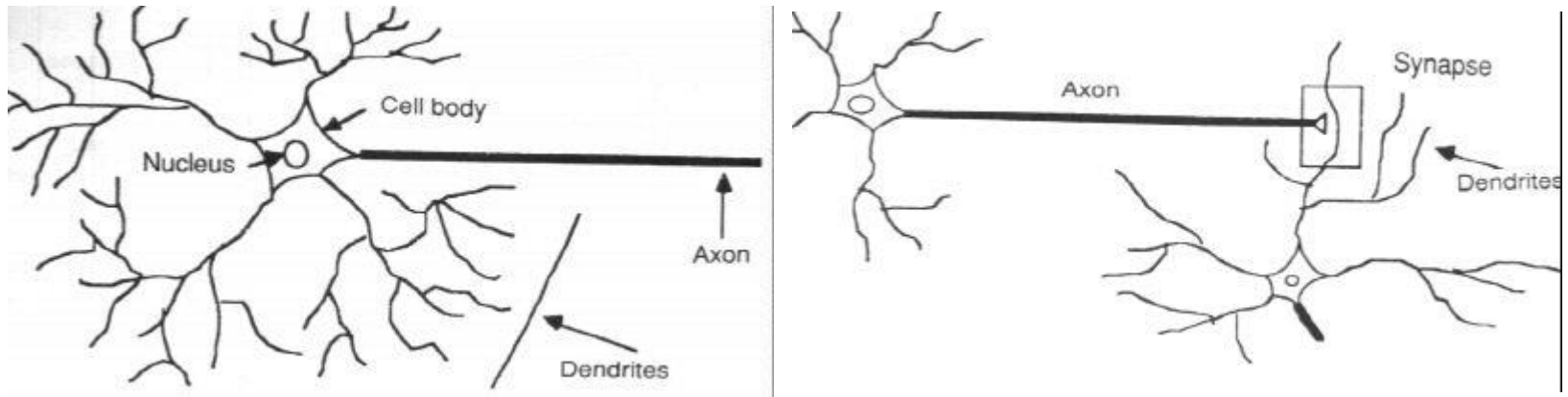
The business end of this is made of lots of these joined in networks like this

Much of our own “computations” are performed in/by this network

Learning occurs by changing the effectiveness of the synapses so that the influence of one neuron on another changes

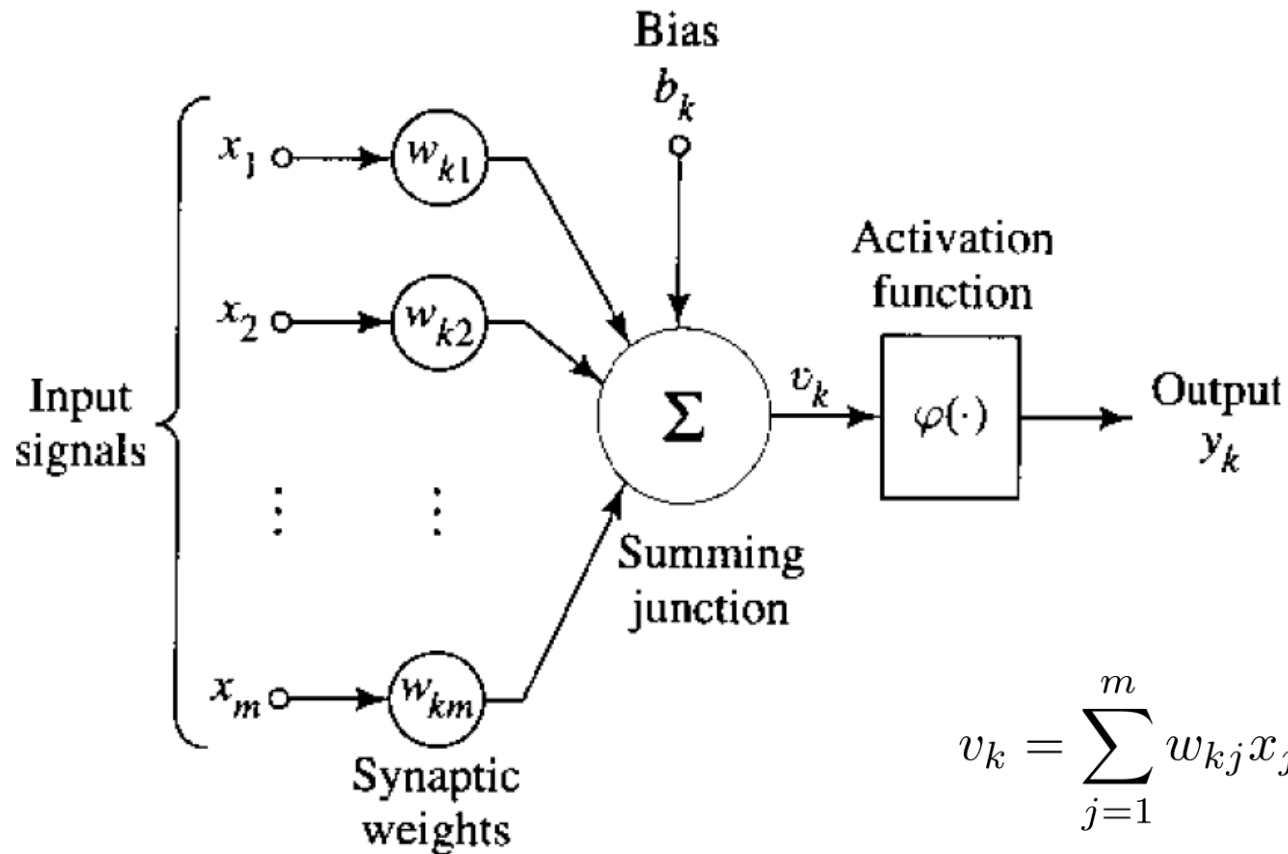
How the human brain learns?

◆ A typical neuron



◆ Learning occurs by changing the effectiveness of the synapses so that the influence of one neuron on another changes

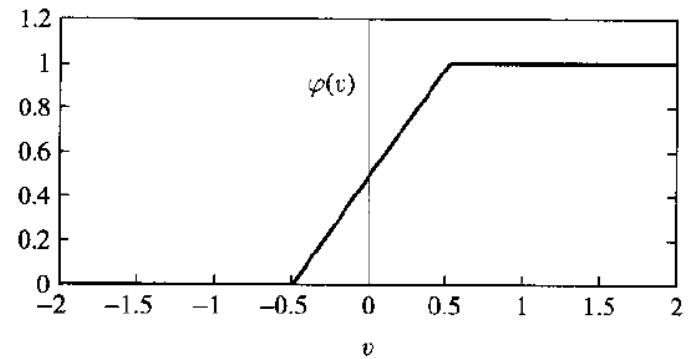
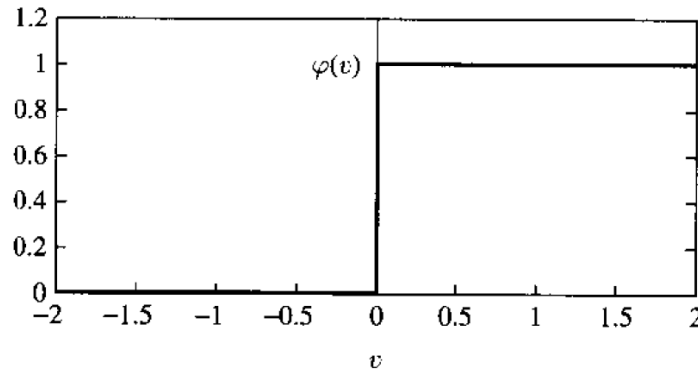
Model of a neuron



$$y_k = \psi(v_k)$$

Activation function

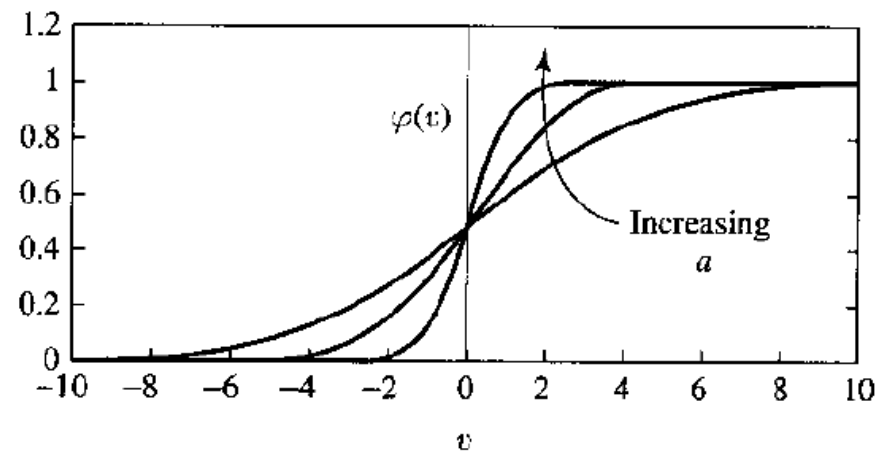
◆ Threshold function & piecewise linear function:



◆ Sigmoid function

$$\psi_{\alpha}(v) = \frac{1}{1 + \exp(-\alpha v)}$$

$a \rightarrow \infty$: step function



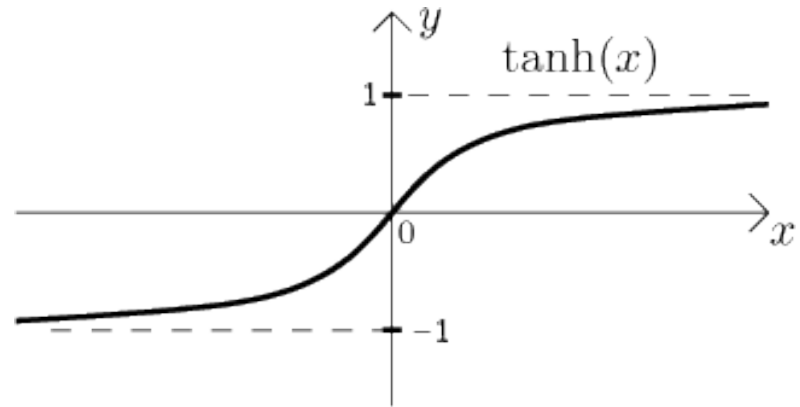
Activation function with negative values

◆ Threshold function & piecewise linear function:

$$\text{sgn}(x) = \begin{cases} 1 & \text{if } x \geq 0 \\ -1 & \text{if } x < 0 \end{cases}$$

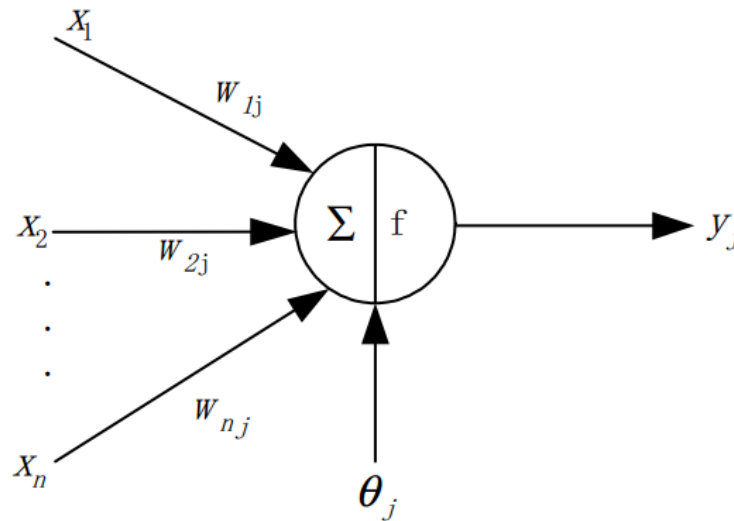
◆ Hyperbolic tangent function

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



McCulloch & Pitts's Artificial Neuron

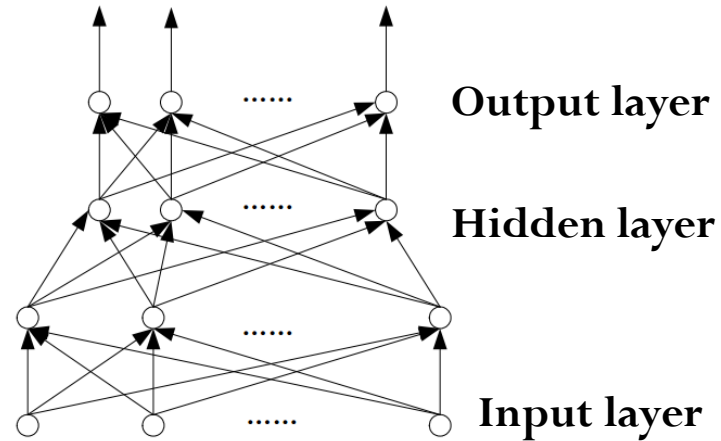
- ◆ The first model of artificial neurons in 1943
 - Activation function: a threshold function



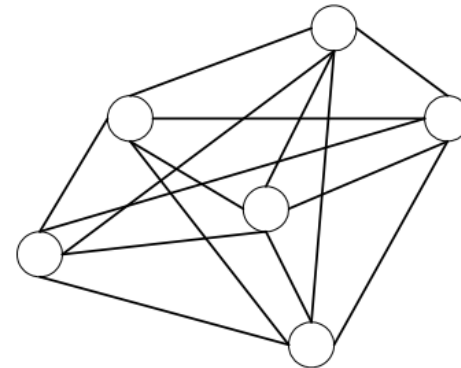
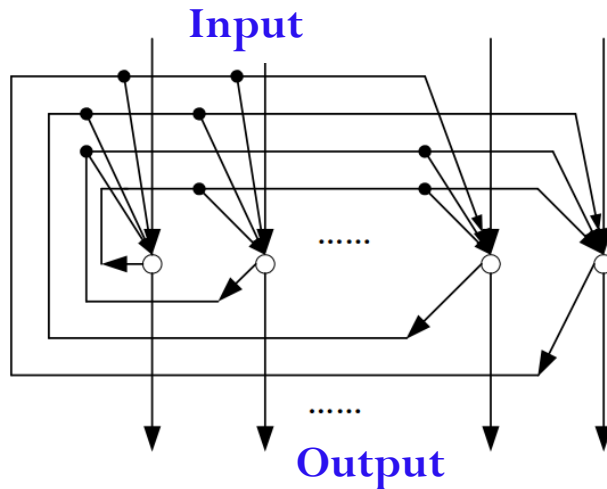
$$y_j = \text{sgn} \left(\sum_i w_{ij} x_i - \theta_j \right)$$

Network Architecture

◆ Feedforward networks

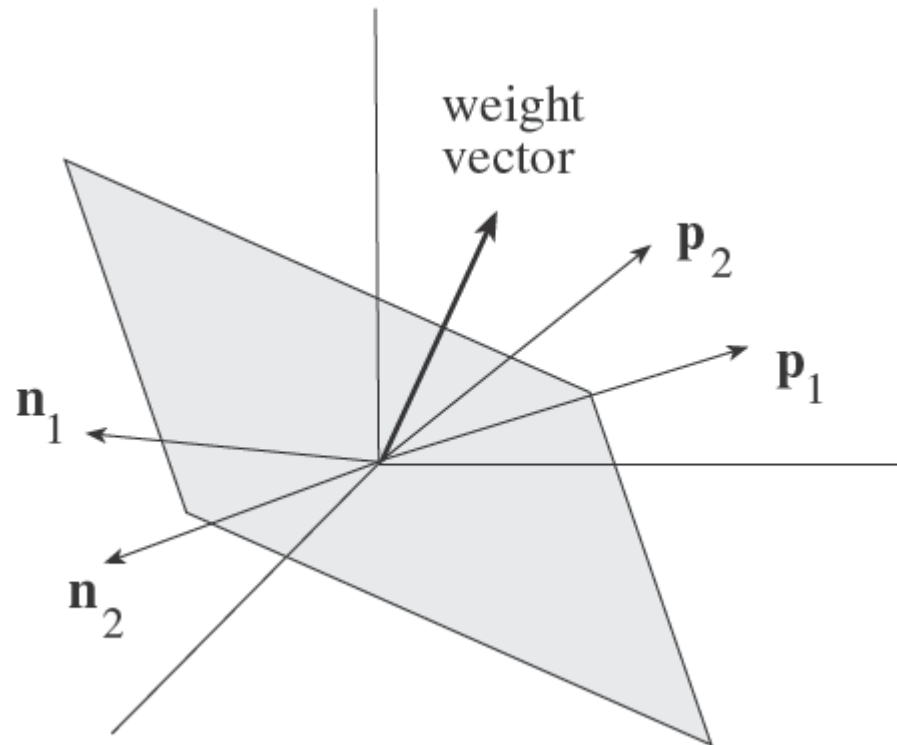


◆ Recurrent networks



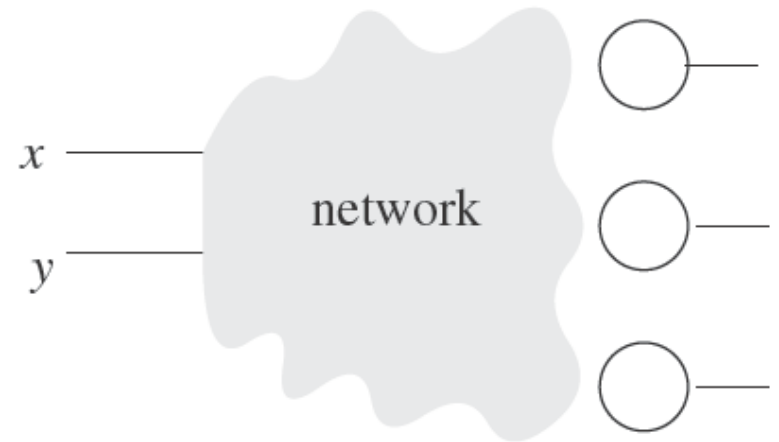
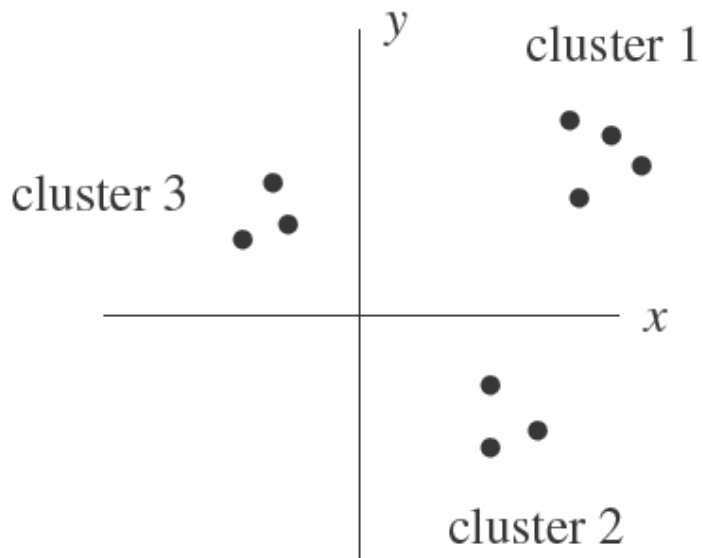
Learning Paradigms

- ◆ Supervised Learning (learning with a teacher)
 - For example, classification: learns a separation plane



Learning Paradigms

- ◆ Unsupervised learning (learning without a teacher)
 - Example: clustering

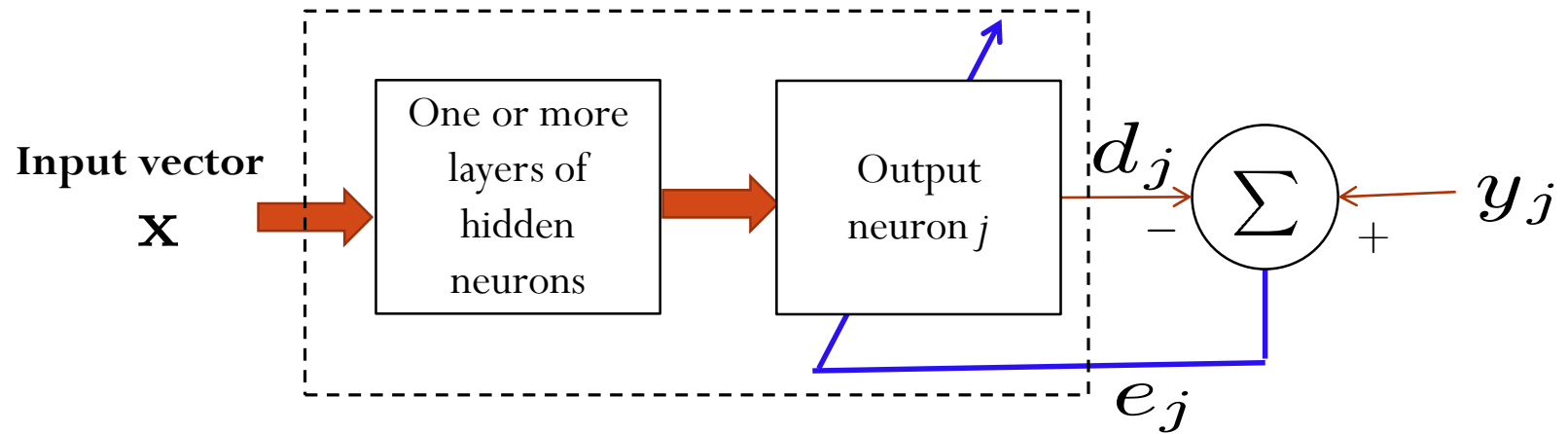


Learning Rules

- ◆ Error-correction learning
- ◆ Competitive learning
- ◆ Hebbian learning
- ◆ Boltzmann learning
- ◆ Memory-based learning
 - Nearest neighbor, radial-basis function network

Error-correction learning

◆ The generic paradigm:



□ Error signal:

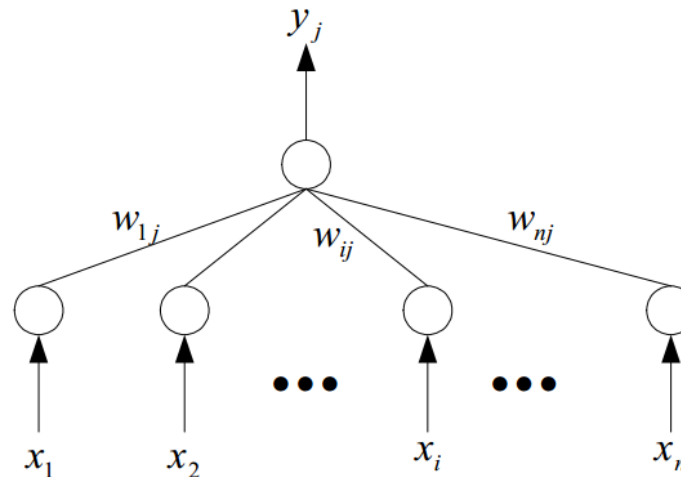
$$e_j = y_j - d_j$$

□ Learning objective:

$$\min_{\mathbf{w}} R(\mathbf{w}; \mathbf{x}) := \frac{1}{2} \sum_j e_j^2$$

Example: Perceptron

- ◆ One-layer feedforward network based on error-correction learning (no hidden layer):



- Current output (at iteration t):

$$d_j = (\mathbf{w}_t^j)^\top \mathbf{x}$$

- Update rule (*exercise?*):

$$\mathbf{w}_{t+1}^j = \mathbf{w}_t^j + \eta(y_j - d_j)\mathbf{x}$$

Perceptron for classification

- ◆ Consider a single output neuron

- ◆ Binary labels:

$$y \in \{+1, -1\}$$

- ◆ Output function:

$$d = \text{sgn} \left(\mathbf{w}_t^\top \mathbf{x} \right)$$

- ◆ Apply the error-correction learning rule, we get ... (next slide)



Perceptron for Classification

◆ Set $\mathbf{w}_1 = 0$ and $t=1$; scale all examples to have length 1
(doesn't affect which side of the plane they are on)

◆ Given example \mathbf{x} , predict positive *iff*

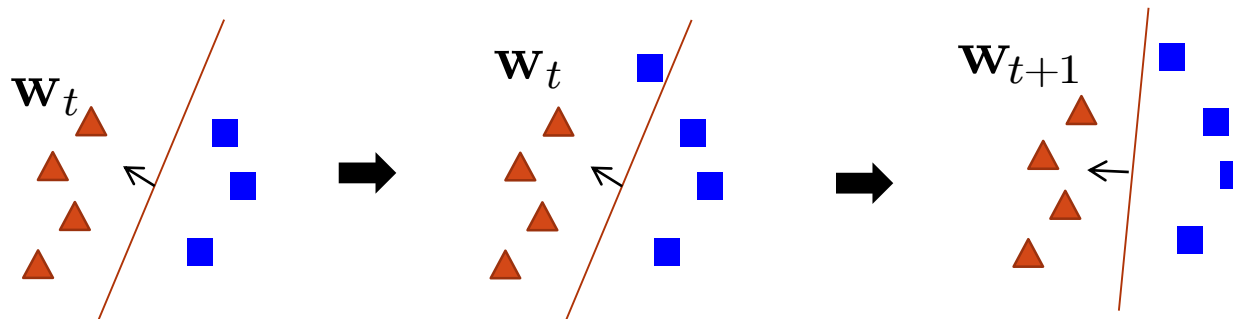
$$\mathbf{w}_t^\top \mathbf{x} > 0$$

◆ If a mistake, update as follows

□ Mistake on positive: $\mathbf{w}_{t+1} \leftarrow \mathbf{w}_t + \eta_t \mathbf{x}$

□ Mistake on negative: $\mathbf{w}_{t+1} \leftarrow \mathbf{w}_t - \eta_t \mathbf{x}$

$t \leftarrow t + 1$



Convergence Theorem

- ◆ For linearly separable case, the perceptron algorithm will converge in a finite number of steps



Mistake Bound

◆ Theorem:

- Let \mathcal{S} be a sequence of labeled examples consistent with a linear threshold function $\mathbf{w}_*^\top \mathbf{x} > 0$, where \mathbf{w}_* is a unit-length vector.
- The number of mistakes made by the online Perceptron algorithm is at most $(1/\gamma)^2$, where

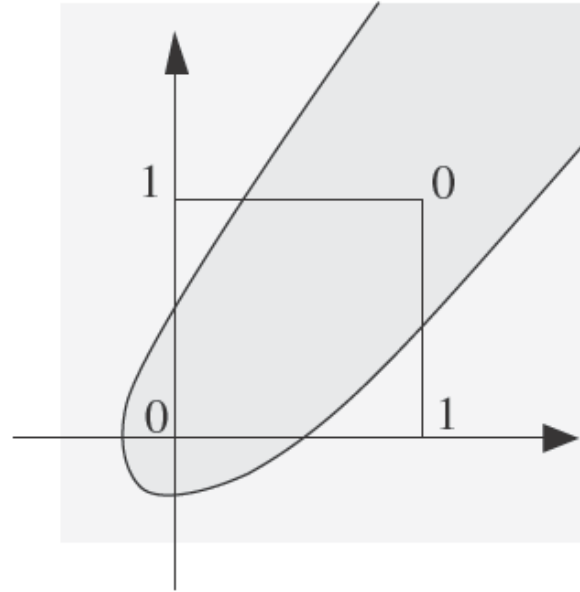
$$\gamma = \min_{\mathbf{x} \in \mathcal{S}} \frac{|\mathbf{w}_*^\top \mathbf{x}|}{\|\mathbf{x}\|}$$

- i.e.: if we scale examples to have length 1, then γ is the minimum distance of any example to the plane $\mathbf{w}_*^\top \mathbf{x} = 0$
- γ is often called the “margin” of \mathbf{w}_* ; the quantity $\frac{\mathbf{w}_*^\top \mathbf{x}}{\|\mathbf{x}\|}$ is the cosine of the angle between \mathbf{x} and \mathbf{w}_*

Deep Nets

- ◆ Deep neural networks
 - Multi-layer Perceptron
 - CNN
 - Deep recurrent nets
- ◆ Deep generative models
 - Auto-encoder
 - RBM
 - Deep belief nets

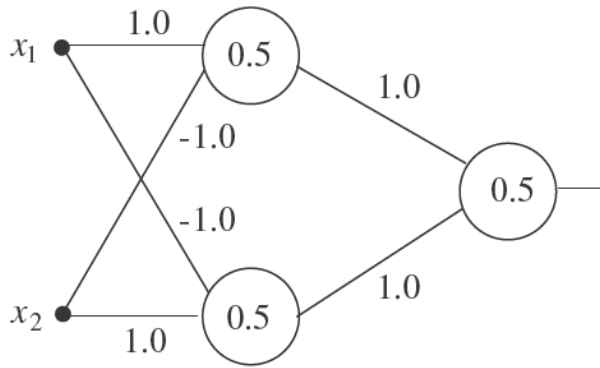
XOR Problem



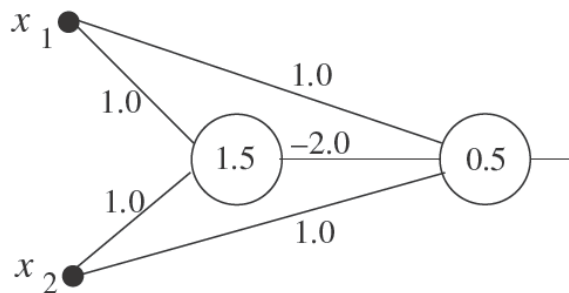
◆ Single-layer perceptron can't solve the problem

XOR Problem

- ◆ A network with 1-layer of 2 neurons works for XOR:
 - threshold activation function



- Many alternative networks exist (not layered)



Multilayer Perceptrons

- ◆ Computational limitations of single-layer Perceptron by Minsky & Papert (1969)

- ◆ Multilayer Perceptrons:
 - Multilayer feedforward networks with an error-correction learning algorithm, known as error *back-propagation*

 - A generalization of single-layer perceptron to allow nonlinearity



Backpropagation

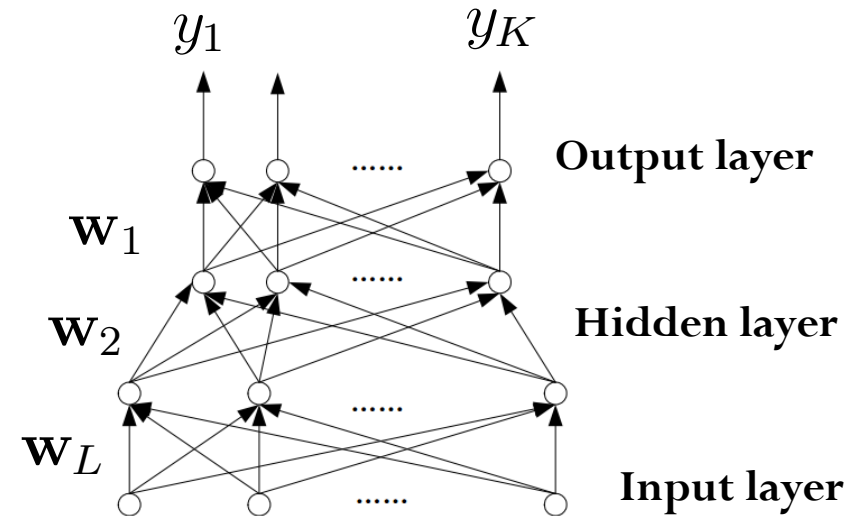
◆ Learning as loss minimization

$$\mathbf{w}^* = \arg \min_{\mathbf{w}} \frac{1}{2} \sum_j e_j^2(\mathbf{x})$$

$$e_j = y_j - d_j$$

◆ Learning with gradient descent

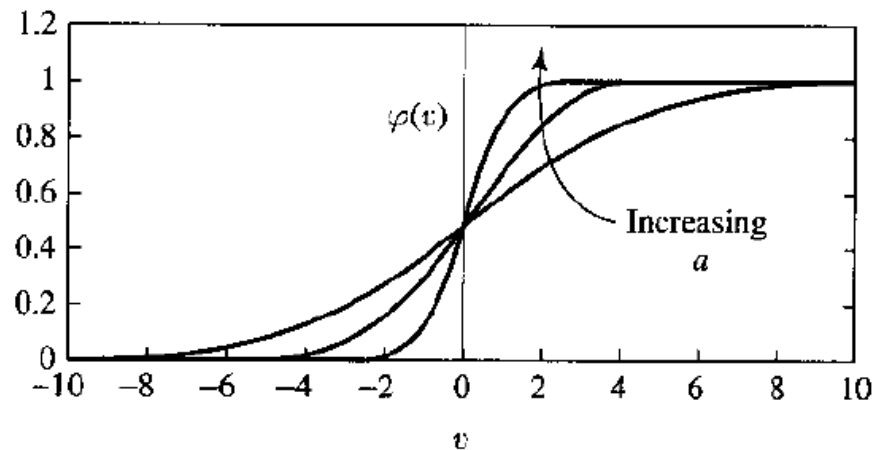
$$\mathbf{w}_{t+1} = \mathbf{w}_t - \lambda_t \nabla R(\mathbf{w}; \mathcal{D})$$



Backpropagation

- ◆ Step function in perceptrons is non-differentiable
- ◆ Differentiable activation functions are needed to calculate gradients, e.g., sigmoid:

$$\psi_{\alpha}(v) = \frac{1}{1 + \exp(-\alpha v)}$$



Backpropagation

◆ Derivative of a sigmoid function ($\alpha = 1$)

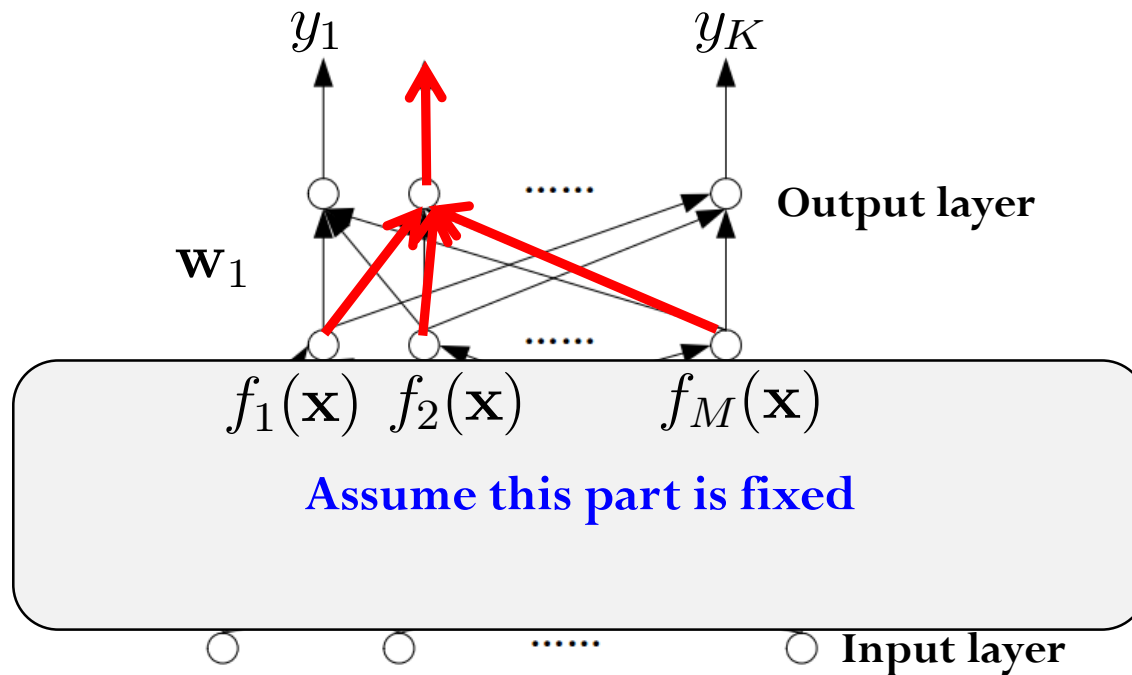
$$\nabla_v \psi(v) = \frac{e^{-v}}{(1 + e^{-v})^2} = \psi(v)(1 - \psi(v))$$

- Note about the small scale of the gradient
- Gradient vanishing issue

◆ Many other activation functions examined

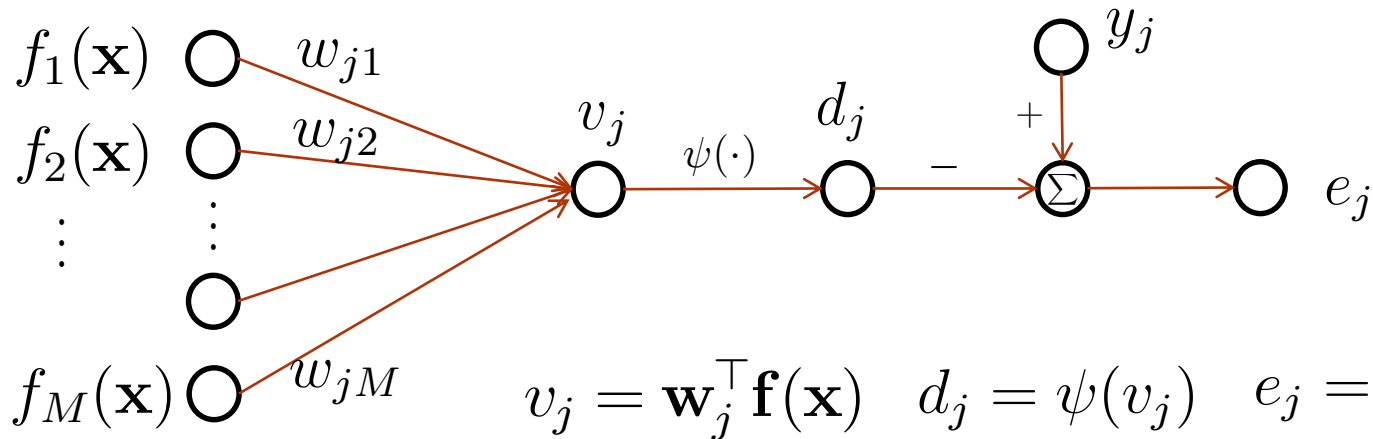
Gradient computation at output layer

- ◆ Output neurons are separate:



Gradient computation at output layer

◆ Signal flow:



$$R_j = \frac{1}{2} e_j^2$$

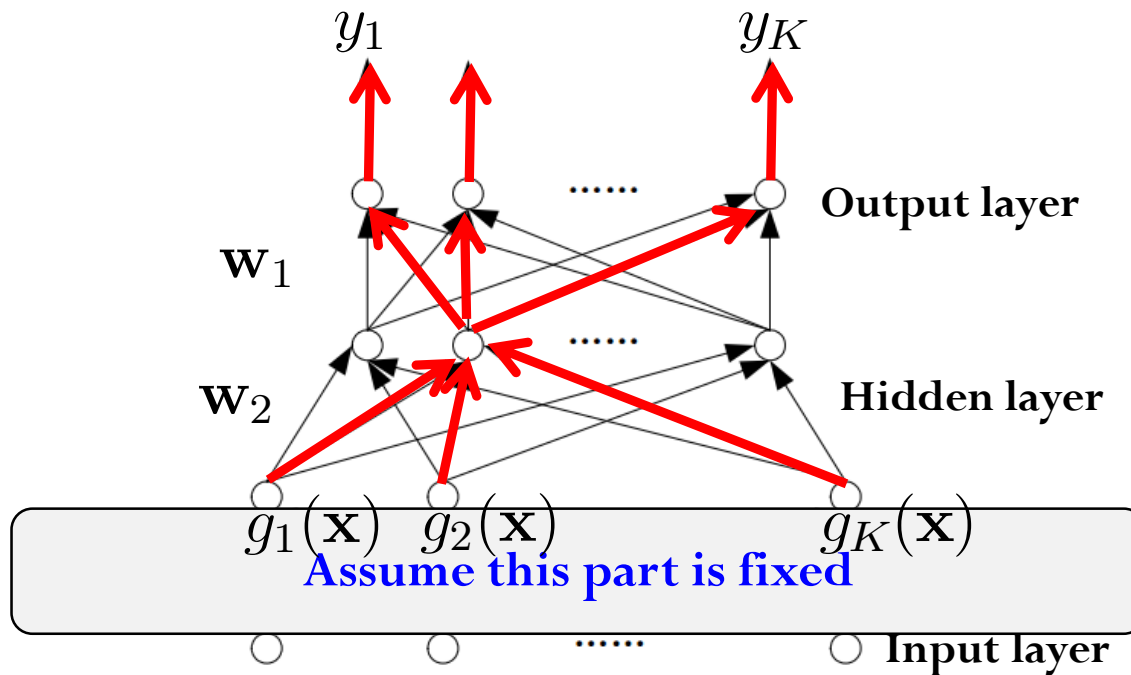
$$R = \frac{1}{2} \sum_j e_j^2$$

$$\begin{aligned}
 \nabla_{w_{ji}} R &= \frac{\partial R_j}{\partial e_j} \frac{\partial e_j}{\partial d_j} \frac{\partial d_j}{\partial v_j} \frac{\partial v_j}{\partial w_{ji}} \\
 &= e_j \cdot (-1) \cdot \psi'(v_j) \cdot f_i(\mathbf{x}) \\
 &= -\boxed{e_j \psi'(v_j)} f_i(\mathbf{x})
 \end{aligned}$$

Local gradient: $\delta_j = -\frac{\partial R}{\partial v_j}$

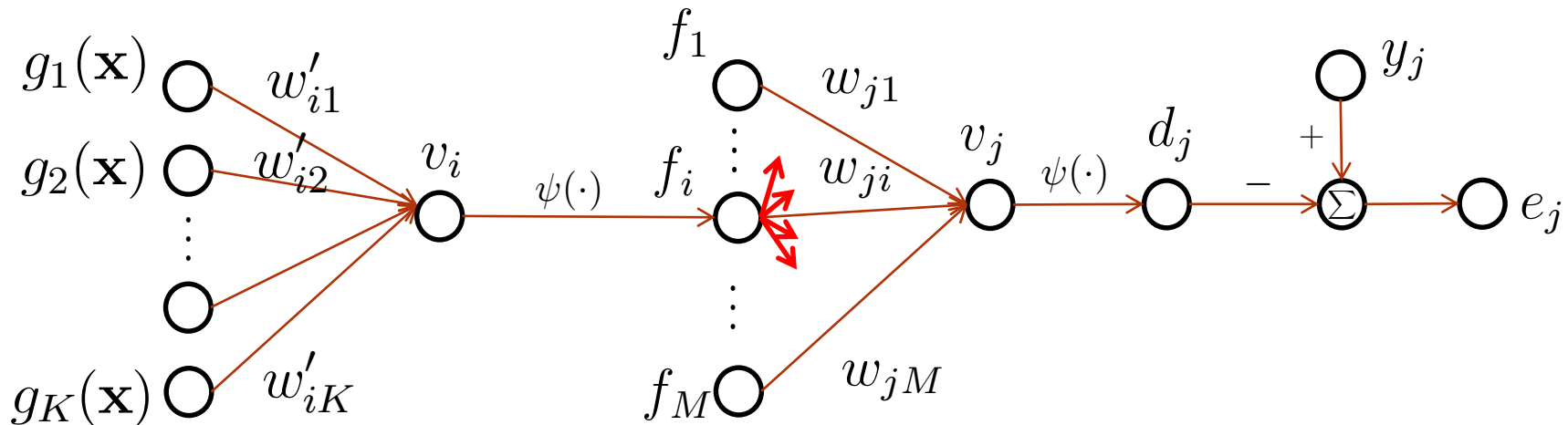
Gradient computation at hidden layer

- ◆ Output neurons are NOT separate:





Gradient computation at hidden layer



$$v_i = (\mathbf{w}'_i)^\top \mathbf{g} \quad f_i = \psi(v_i) \quad v_j = \mathbf{w}_j^\top \mathbf{f} \quad d_j = \psi(v_j) \quad e_j = y_j - d_j$$

$$\begin{aligned} R_j &= \frac{1}{2} e_j^2 \\ R &= \frac{1}{2} \sum_j e_j^2 \\ \nabla_{w'_{ik}} R &= \sum_j \frac{\partial R_j}{\partial e_j} \frac{\partial e_j}{\partial d_j} \frac{\partial d_j}{\partial v_j} \frac{\partial v_j}{\partial f_i} \frac{\partial f_i}{\partial v_i} \frac{\partial v_i}{\partial w'_{ik}} \\ &= - \sum_j e_j \psi'(v_j) w_{ji} \psi'(v_i) g_k(\mathbf{x}) \\ &= - \sum_j \delta_j w_{ji} \psi'(v_i) g_k(\mathbf{x}) \end{aligned}$$

Local gradient: $\delta_i = - \frac{\partial R}{\partial v_i}$

Back-propagation formula

◆ The update rule of **local gradients**:

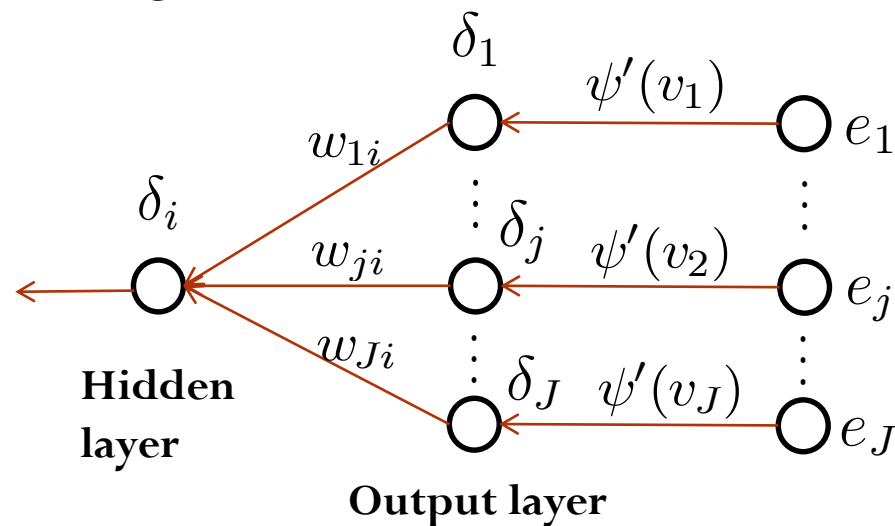
□ for hidden neuron i :

$$\delta_i = \psi'(v_i) \sum_j \delta_j w_{ji}$$

↑

Only depends on the activation function at hidden neuron i

◆ Flow of error signal:



Back-propagation formula

◆ The update rule of weights:

□ Output neuron:

$$\Delta w_{ji} = \lambda \cdot \delta_j \cdot f_i(\mathbf{x})$$

□ Hidden neuron:

$$\Delta w'_{ik} = \lambda \cdot \delta_i \cdot g_k(\mathbf{x})$$

$$\begin{pmatrix} \text{Weight} \\ \text{correction} \\ \Delta w_{ji} \end{pmatrix} = \begin{pmatrix} \text{learning} \\ \text{rate} \\ \lambda \end{pmatrix} \cdot \begin{pmatrix} \text{local} \\ \text{gradient} \\ \delta_j \end{pmatrix} \cdot \begin{pmatrix} \text{input signal} \\ \text{of neuron } j \\ v_i \end{pmatrix}$$

Two Passes of Computation

◆ Forward pass

- Weights fixed
- Start at the first hidden layer
- Compute the output of each neuron
- End at output layer

◆ Backward pass

- Start at the output layer
- Pass error signal backward through the network
- Compute local gradients

Stopping Criterion

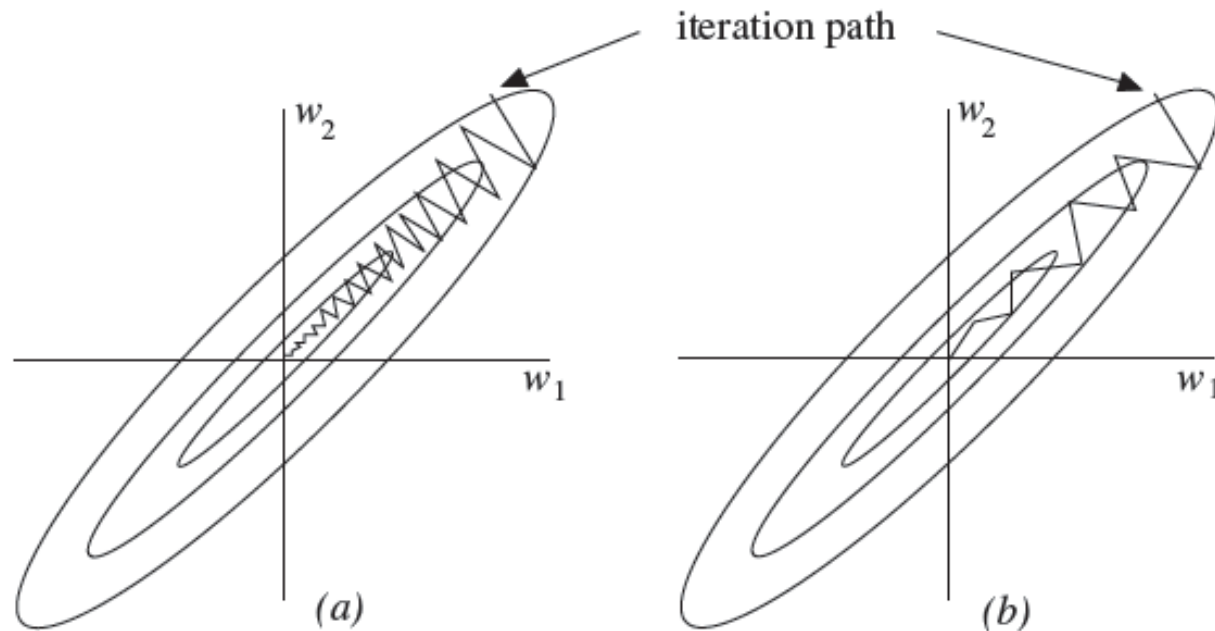
- ◆ No general rules

- ◆ Some reasonable heuristics:
 - The norm of gradient is small enough
 - The number of iterations is larger than a threshold
 - The training error is stable
 - ...

Improve Backpropagation

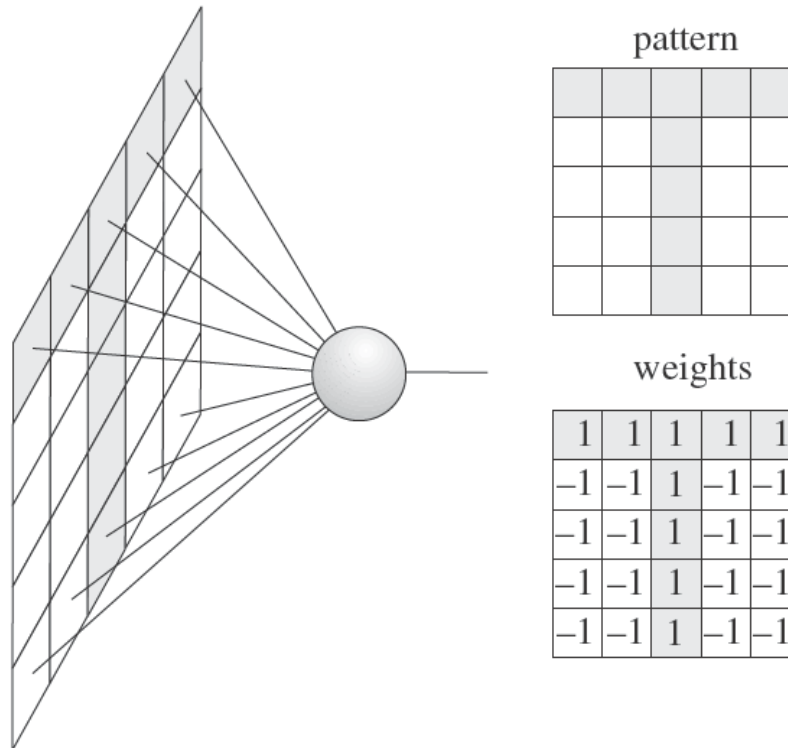
- ◆ Many methods exist to improve backpropagation
- ◆ E.g., backpropagation with momentum

$$\Delta w_{ij}^t = -\lambda \frac{\partial R}{\partial w_{ij}} + \alpha \Delta w_{ij}^{t-1}$$



Neurons as Feature Extractor

- ◆ Compute the similarity of a pattern to the ideal pattern of a neuron
- ◆ Threshold is the minimal similarity required for a pattern
- ◆ Reversely, it visualizes the connections of a neuron



Vanishing gradient problem

◆ The gradient can decrease exponentially during back-prop

◆ Solutions:

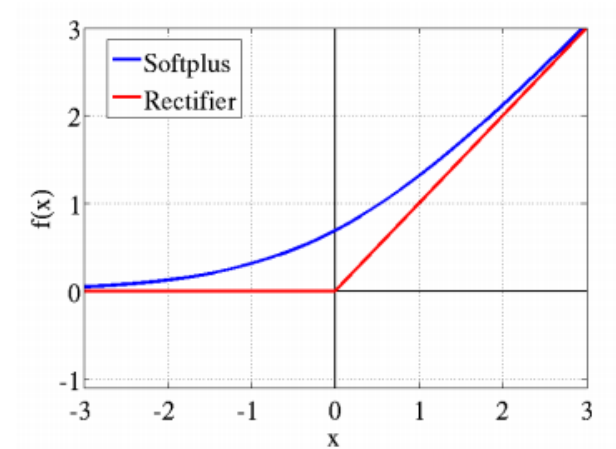
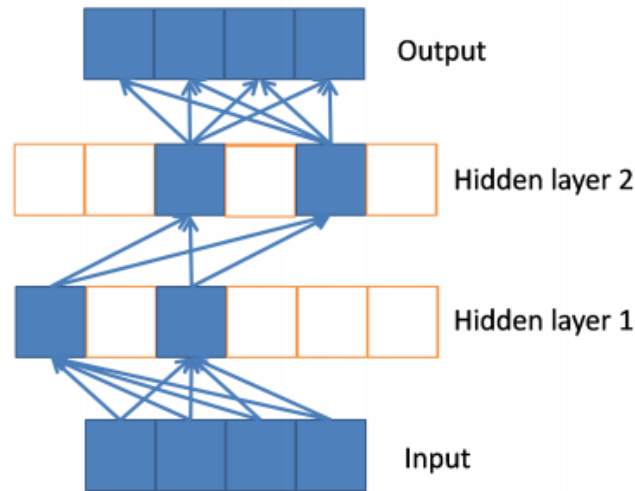
- Pre-training + fine tuning
- Rectifier neurons (sparse gradients)

◆ Ref:

- Gradient flow in recurrent nets: the difficulty of learning long-term dependencies. Hochreiter, Bengio, & Frasconi, 2001

Deep Rectifier Nets

- ◆ Sparse representations without gradient vanishing

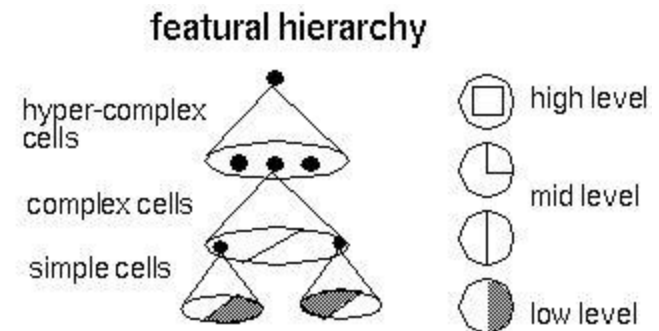
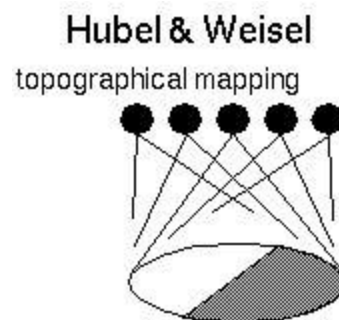
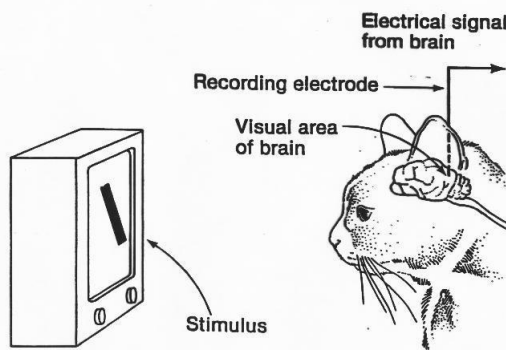


- Non-linearity comes from the path selection
 - Only a subset of neurons are active for a given input
- Can be seen as a model with an exponential number of linear models that share weights

[Deep sparse rectifier neural networks. Glorot, Bordes, & Bengio, 2011]

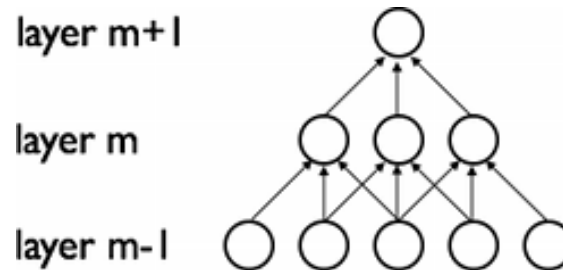
CNN

- ◆ Hubel and Wiesel's study on animal's visual cortex:
 - ❑ Cells that are sensitive to small sub-regions of the visual field, called a *receptive field*
 - ❑ Simple cells respond maximally to specific edge-like patterns within their receptive field. Complex cells have larger receptive fields and are locally invariant to the exact position of the pattern.

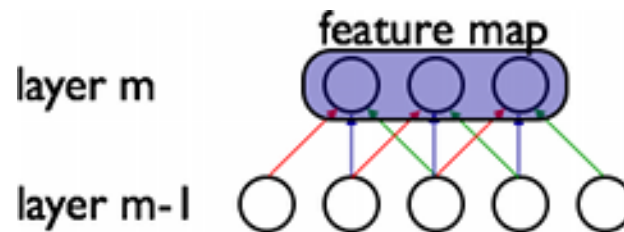


Convolutional Neural Networks

- ◆ Sparse local connections (**spatially contiguous receptive fields**)



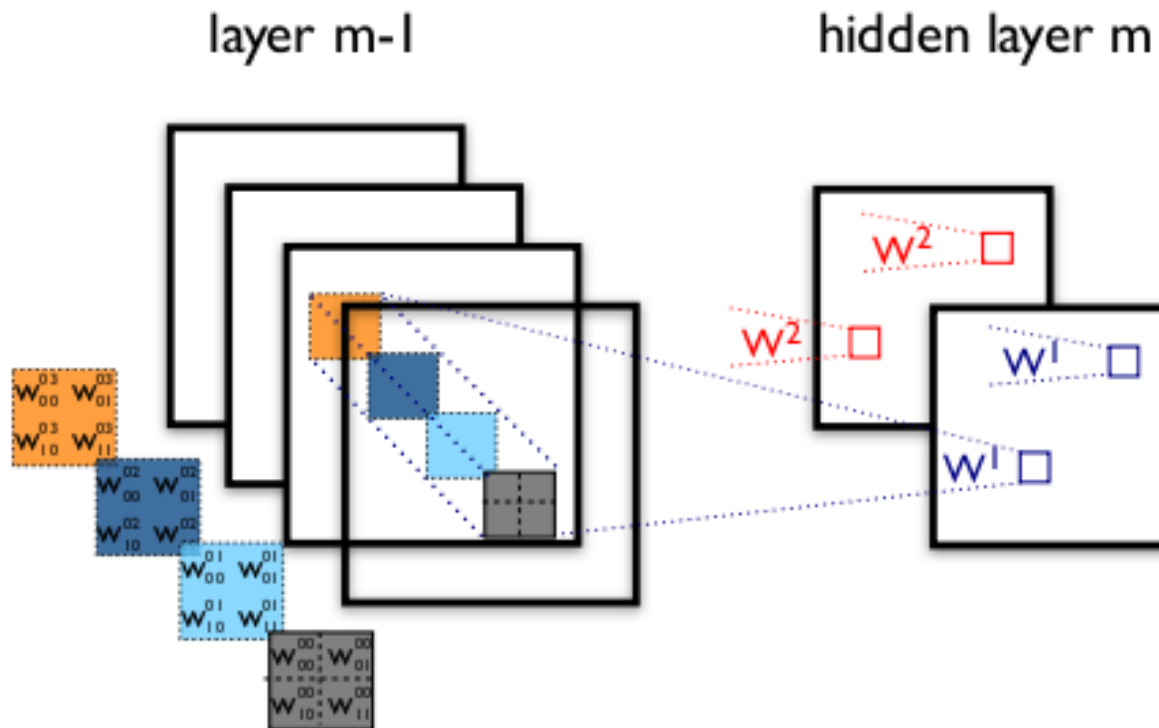
- ◆ Shared weights: each filter is replicated across the entire visual field, forming a feature map





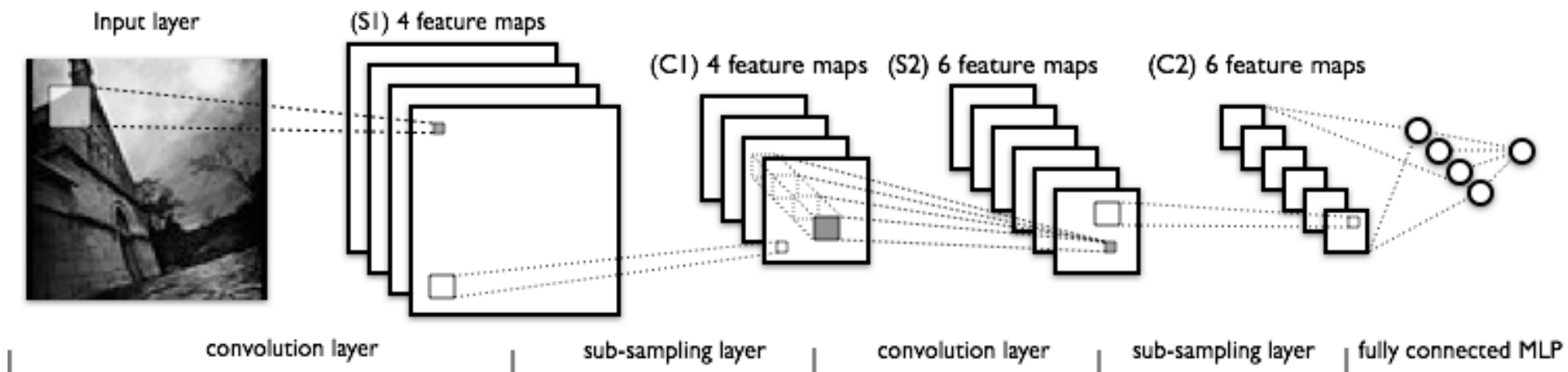
CNN

- ◆ Each layer has multiple feature maps



CNN

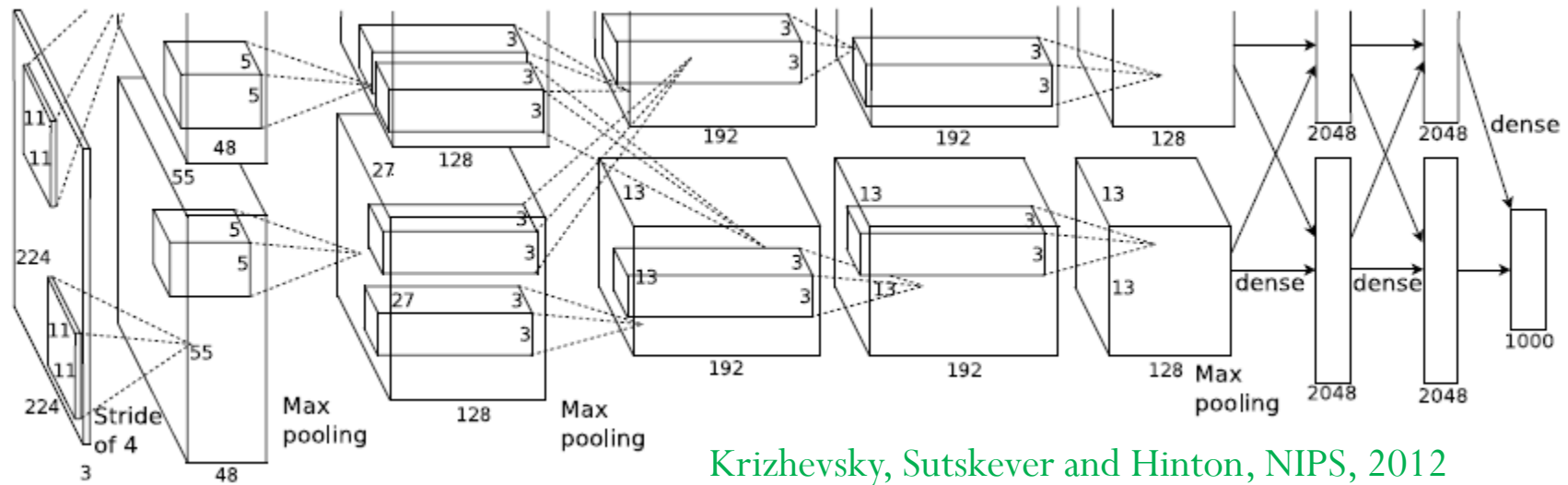
◆ The full model



◆ *Max-pooling*, a form of non-linear down-sampling.

- Max-pooling partitions the input image into a set of non-overlapping rectangles and, for each such sub-region, outputs the maximum value.

Example: CNN for image classification



- ◆ Network dimension: 150,528(input)-253,440-186,624-64,896-64,896-43,264-4096-4096-1000(output)
 - In total: 60 million parameters
 - Task: classify 1.2 million high-resolution images in the ImageNet LSVRC-2010 contest into the 1000 different classes
 - Results: state-of-the-art accuracy on ImageNet

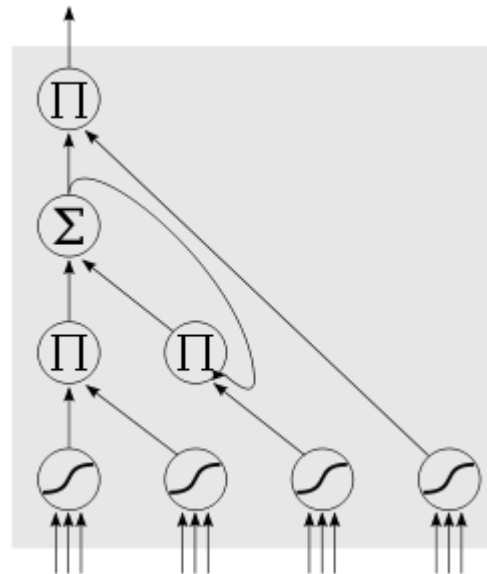
Issues with CNN

- ◆ Computing the activations of a single convolutional filter is much more expensive than with traditional MLPs

- ◆ Many tuning parameters
 - # of filters:
 - Model complexity issue (overfitting vs underfitting)
 - Filter shape:
 - the right level of “granularity” in order to create abstractions at the proper scale, given a particular dataset
 - Usually 5x5 for MNIST at 1st layer
 - Max-pooling shape:
 - typical: 2x2; maybe 4x4 for large images

Long Short-Term Memory

- ◆ A RNN architecture without gradient vanishing issue
- ◆ A RNN with LSTM blocks
 - Each block is a “smart” network, determining when to remember, when to continue to remember or forget, and when to output



Discussions

Challenges of DL

◆ Learning

- Backpropagation is slow and prone to gradient vanishing
- Issues with non-convex optimization in high-dimensions

◆ Overfitting

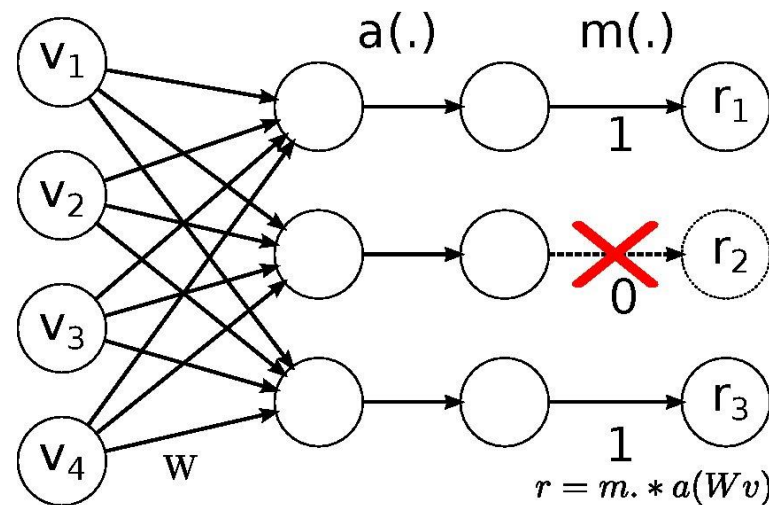
- Big models are lacking of statistical information to fit

◆ Interpretation

- Deep nets are often used as black-box tools for learning and inference

Overfitting in DL

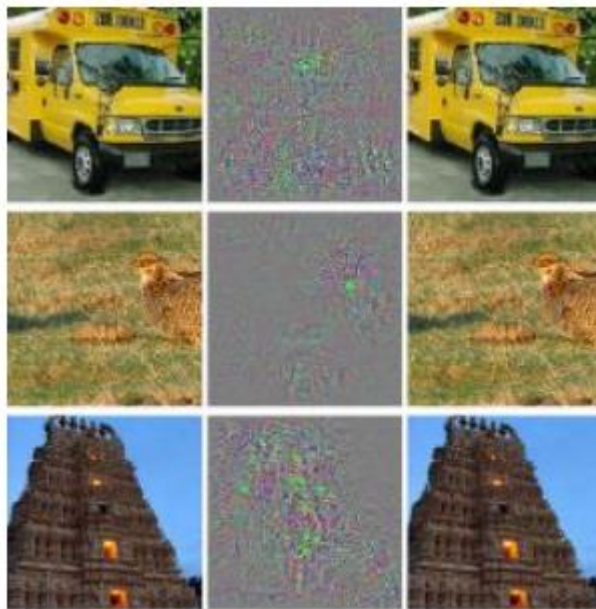
- ◆ Increasing research attention, e.g., dropout training (Hinton, 2012)



- ◆ More theoretical understanding and extensions
 - MCF (van der Maaten et al., 2013); Logistic-loss (Wager et al., 2013); Dropout SVM (Chen, Zhu et al., 2014)

Some counter-intuitive properties

- ◆ Stability w.r.t small perturbations to inputs
 - Imperceptible non-random perturbation can arbitrarily change the prediction (**adversarial examples exist!**)



(a)

10x of
differences



(b)



Cat: 98.96%



Transferability

◆ Cross-model transferability

Target model with
unknown weights,
machine learning
algorithm, training
set; maybe non-
differentiable

Train your
own model

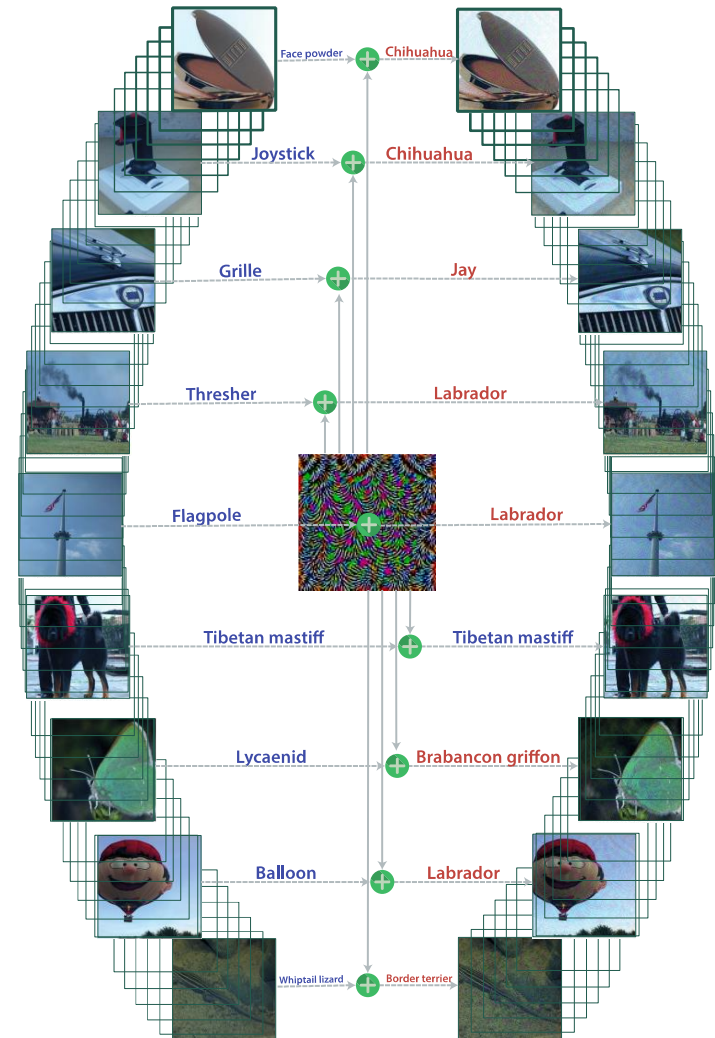
Substitute model
mimicking target
model with known,
differentiable function

Adversarial
examples

Adversarial crafting
against substitute

Deploy adversarial
examples against the
target; transferability
property results in them
succeeding

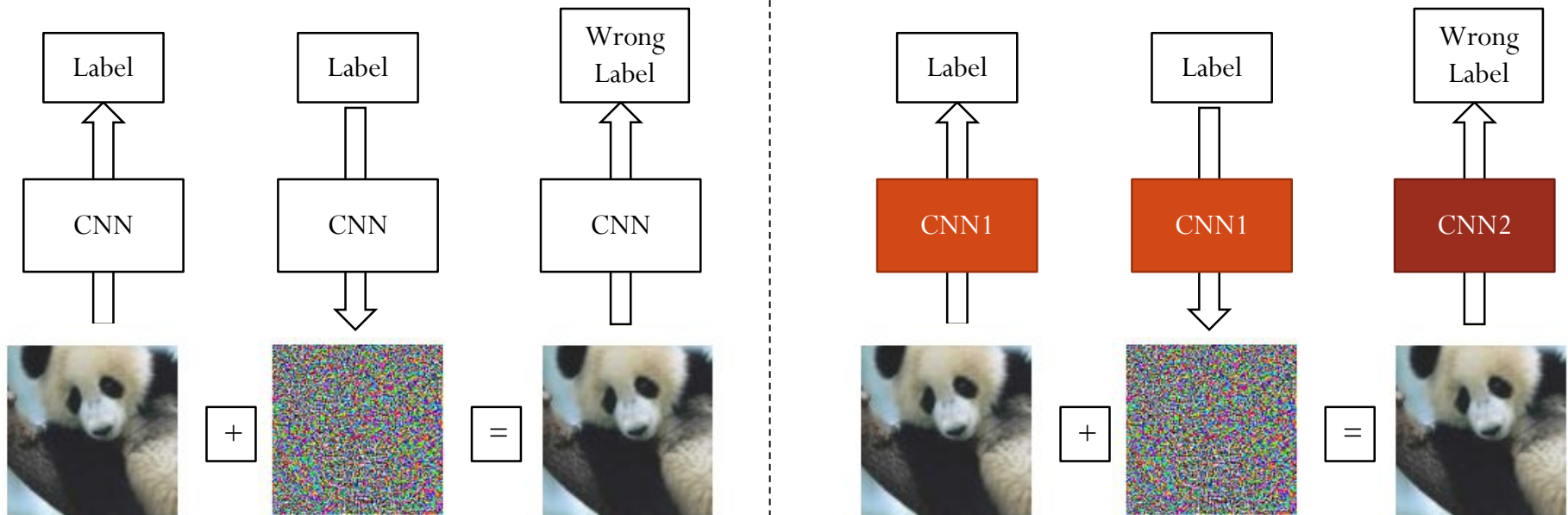
◆ Cross-data transferability





Adversarial Attack & Defense

◆ White-box v.s black-box attack



A game between attack and defense!

Adversarial Training

$$L = \frac{1}{(m - k) + \lambda k} \left(\sum_i L(x_i, y_i) + \lambda \sum_j L(x_j^*, y_j) \right)$$

		Clean	$\epsilon = 2$	$\epsilon = 4$	$\epsilon = 8$	$\epsilon = 16$
Baseline (standard training)	top 1	78.4%	30.8%	27.2%	27.2%	29.5%
	top 5	94.0%	60.0%	55.6%	55.1%	57.2%
Adv. training	top 1	77.6%	73.5%	74.0%	74.5%	73.9%
	top 5	93.8%	91.7%	91.9%	92.0%	91.4%
Deeper model (standard training)	top 1	78.7%	33.5%	30.0%	30.0%	31.6%
	top 5	94.4%	63.3%	58.9%	58.1%	59.5%
Deeper model (Adv. training)	top 1	78.1%	75.4%	75.7%	75.6%	74.4%
	top 5	94.1%	92.6%	92.7%	92.5%	91.6%

◆ Kurakin et al 2017: ADVERSARIAL MACHINE LEARNING AT SCALE

NIPS Challenge

- ◆ Google Brain organized the NIPS 2017 contest on Adversarial Attacks and Defenses
- ◆ Three sub-tasks
 - Un-targeted adversarial attack (mislead the network)
 - Targeted adversarial attack (mislead to a pre-defined category)
 - Defense against adversarial attacks
- ◆ Run all attacks against all defenses to evaluate how each of the attacks performs against each of the defenses
- ◆ 100+ teams in each task
- ◆ We won the 1st places in all three tasks
 - A large-margin from the 2nd in target attack and defense
 - Two papers at CVPR 2018

Criticisms of DL

- ◆ Just a buzzword, or largely a rebranding of neural networks
- ◆ Lack of theory
 - gradient descent has been understood for a while
 - DL is often used as black-box
- ◆ DL is only part of the larger challenge of building intelligent machines, still lacking of:
 - causal relationships
 - logic inferences
 - integrating abstract knowledge



Will DL make other ML methods obsolete?

Quora 2014/12/23

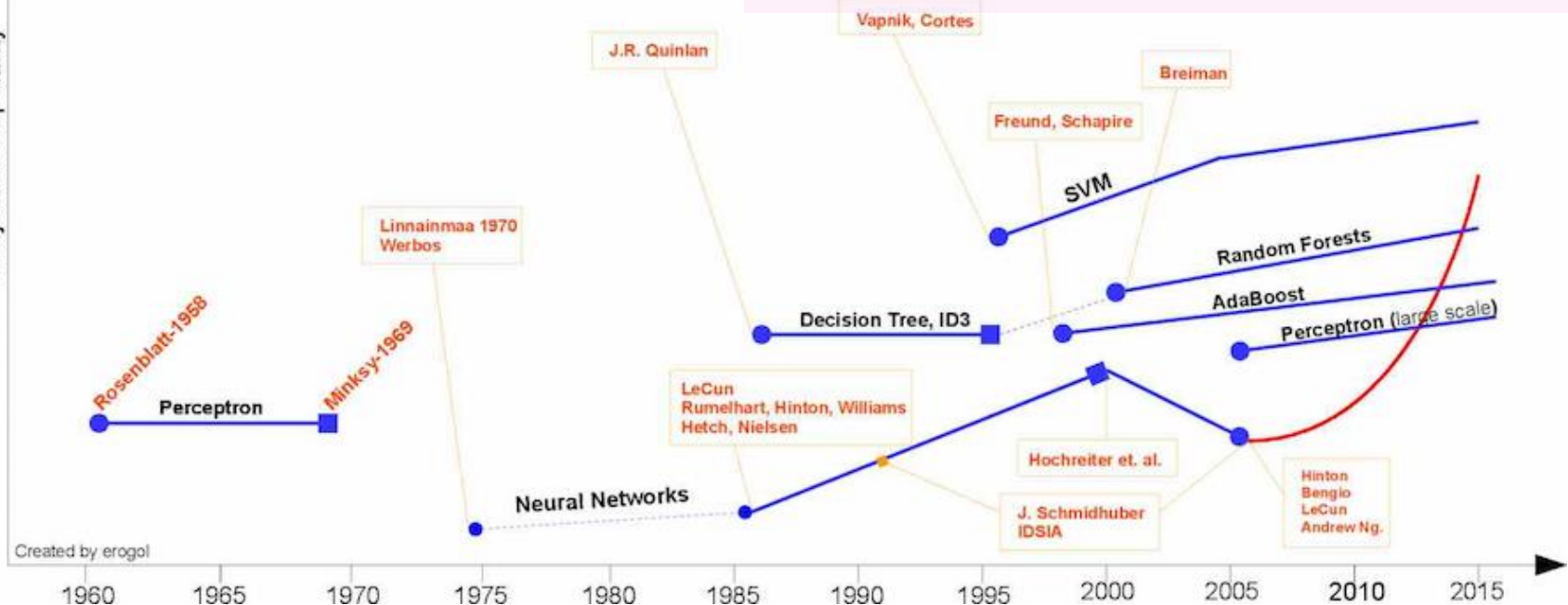
Yes (2 post, 113 upvotes)

- best predictive power when data sufficient
- DL is far from saturated
- Google et al invests on DL, it is the “richest” AI topic

No (10 posts, 284 upvotes)

- simpler algorithms are just fine in many cases
- methods with domain knowledge works better
- DL is feature learning, needs other methods to work
- DL is not that well developed, a lot of work to be done using more traditional methods
- No free lunch
- a lot like how ANN was viewed in the late 80s

Subjective Popularity





What are people saying?

◆ Yann LeCun:

- “AI has gone from failure to failure, with bits of progress. This could be another leapfrog”

◆ Jitendra Malik:

- in the long term, deep learning may not win the day; ... “Over time people will decide what works best in different domains.”
- “Neural nets were always a delicate art to manage. There is some black magic involved”

◆ Andrew Ng:

- “Deep learning happens to have the property that if you feed it more data it gets better and better,”
- “Deep-learning algorithms aren't the only ones like that, but they're arguably the best — certainly the easiest. That's why it has huge promise for the future.”

Thank You!