



An Introduction to Deep Learning

1

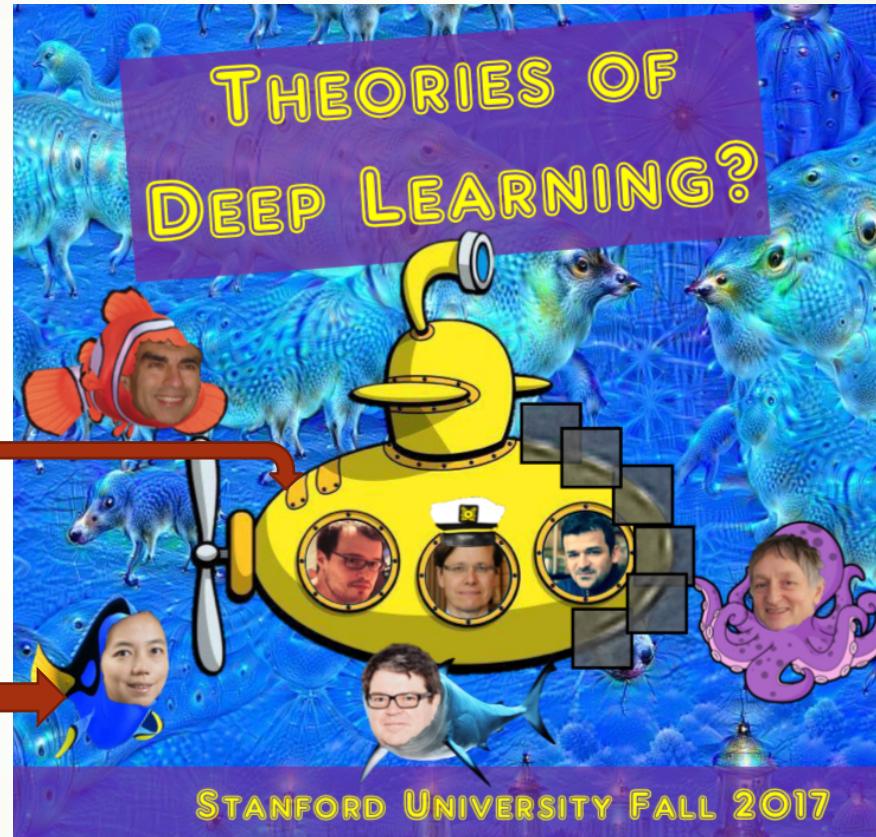
Yuan YAO

HKUST

Acknowledgement

<https://stats385.github.io/>

<http://cs231n.github.io/>



A following-up course at HKUST: <https://deeplearning-math.github.io/>



Some reference books

- ▶ Deep Learning with Python, Manning Publications 2017
 - ▶ by François Chollet
 - ▶ https://www.manning.com/books/deep-learning-with-python?a_aid=keras&a_bid=76564dff
- ▶ Deep Learning, MIT Press 2016
 - ▶ By Ian Goodfellow, Yoshua Bengio, and Aaron Courville,
 - ▶ <http://www.deeplearningbook.org/>
- ▶ Many other public resources

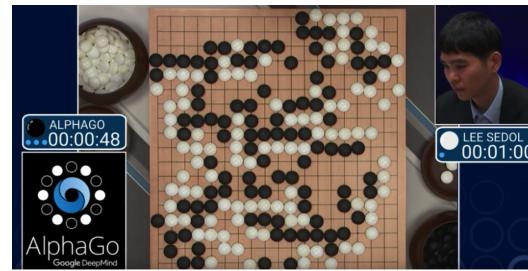


The Tsunami of Deep Learning

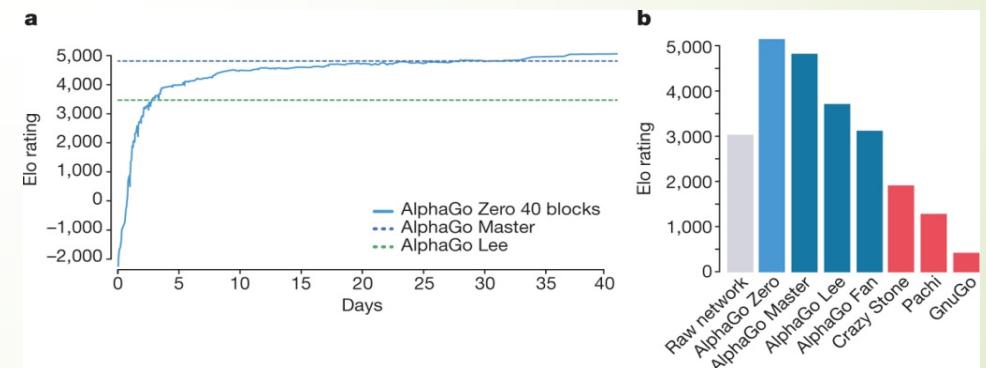
Reaching Human Performance Level in Games



Deep Blue in 1997



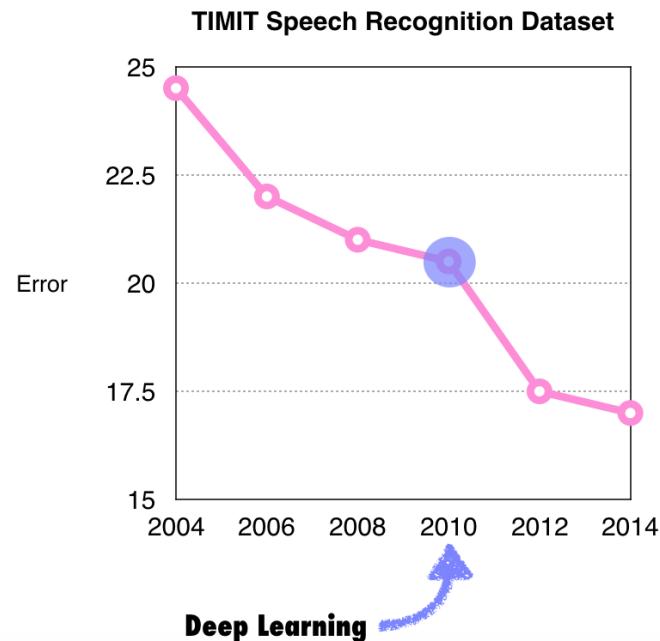
AlphaGo "LEE" 2016



AlphaGo "ZERO" D Silver et al. *Nature* **550**, 354–359 (2017) doi:10.1038/nature24270

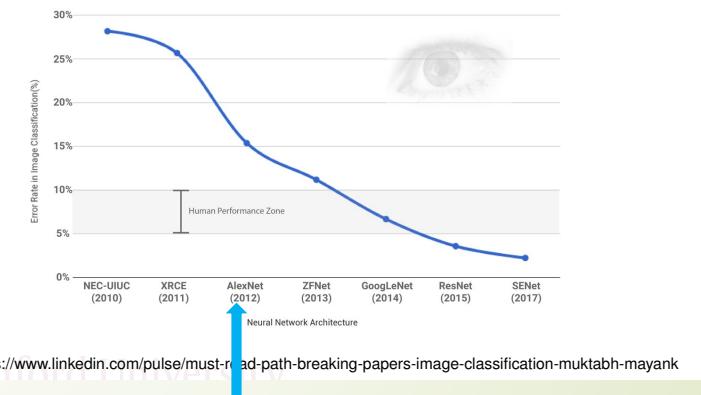
Speech Recognition & Computer Vision

Speech Recognition: TIMIT



Computer Vision: ImageNet

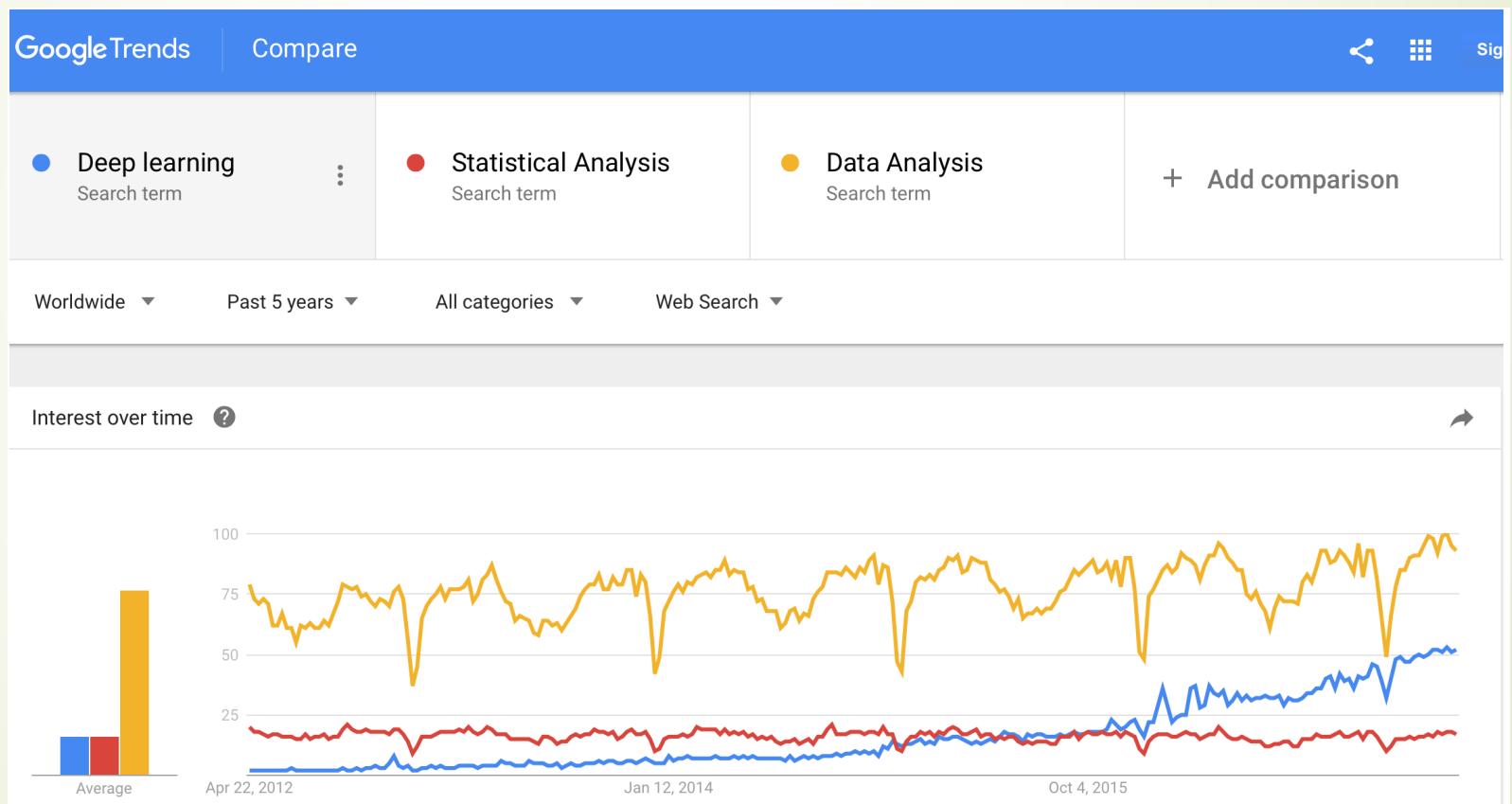
- ImageNet (subset):
 - 1.2 million training images
 - 100,000 test images
 - 1000 classes
- ImageNet large-scale visual recognition Challenge



Deep Learning

Growth of Deep Learning

'Deep Learning' is coined by Hinton et al. in their Restricted Boltzman Machine paper, Science 2006, not yet popular until championing ImageNet competitions.



"We're at the beginning of a new day...
This is the beginning of the AI revolution."
— Jensen Huang, GTC Taiwan 2017



兩股力量驅動電腦的未來

深度學習點亮人工智慧紀元。

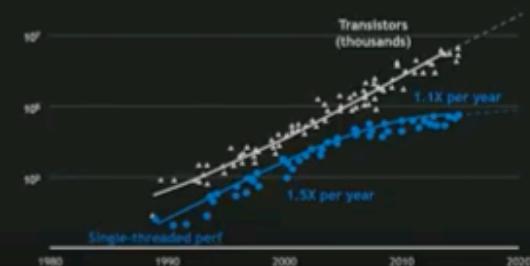
受到人腦的啟發，深度神經網路具備上億的類神經連結，藉由巨量資料來學習，這仰賴極大量的運算。

同時，摩爾定律已到了尾聲 - CPU已不可能再擴張成長。

程式設計人員無法創造出可以更有效率發現更多指令級並行性的CPU架構。

電晶體持續每年增長50%，但是CPU效能僅能成長10%。

TWO FORCES DRIVING THE FUTURE OF COMPUTING



New Moore's Laws

CS231n attendance

 **Andrej Karpathy**  @karpathy 

Came to visit first class of [@cs231n](#) at Stanford. 2015: 150 students, 2016: 350, this year: 750. [#aiinterestsingularity](#)



12:11 PM - 4 Apr 2017

155 Retweets 623 Likes 

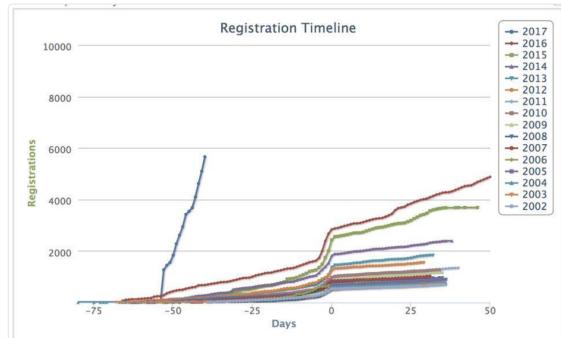
19  155  623 

 **michael_nielsen** @michael_nielsen · Apr 4
Replying to @karpathy @cs231n
Faster than Moore's Law. At this rate - doubling each year - in 24 years everyone on Earth will be enrolled :-)

NIPS registrations

 **Alex Lebrun** @lxbrun 

Deep learning hype in one picture (NIPS conference registrations, 2002 through 2017) [#nips2017](#)



Registration Timeline

Registrations

Days

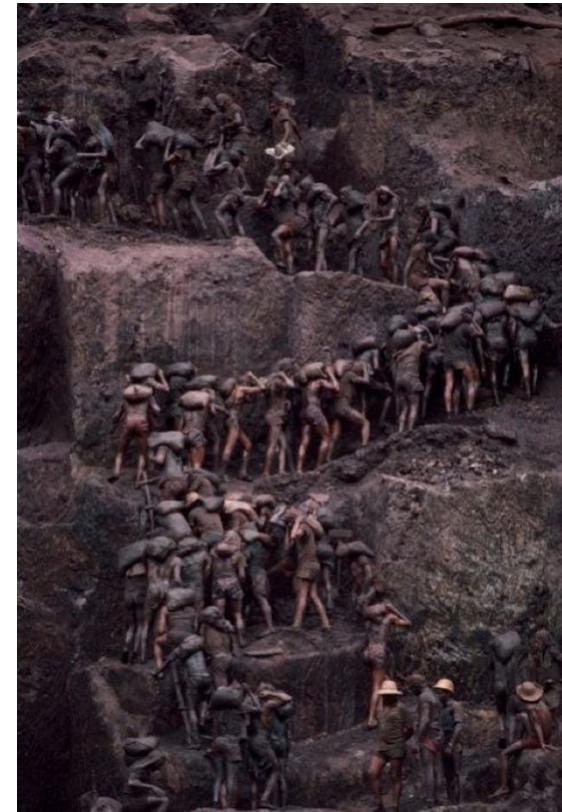
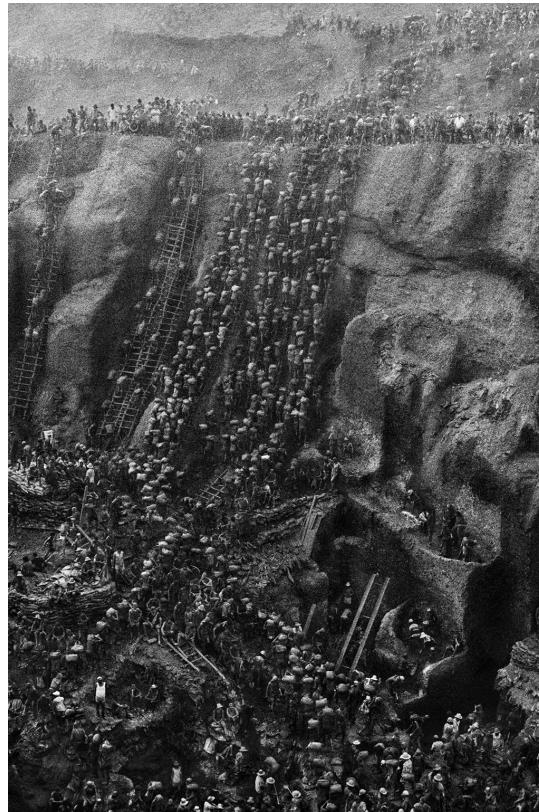
Legend: 2017, 2016, 2015, 2014, 2013, 2012, 2011, 2010, 2009, 2008, 2007, 2006, 2005, 2004, 2003, 2002

8:20 AM - 15 Sep 2017

758 Retweets 1,005 Likes 

20  758  1.0K 

Crowdcomputing: new comers raising the competition record



Some Cold Water: Tesla Autopilot Misclassifies Truck as Billboard



Problem: Why? How can you trust a blackbox?

Deep Learning may be fragile in generalization against noise!

$$x + .007 \times \text{sign}(\nabla_x J(\theta, x, y)) = x + \epsilon \text{sign}(\nabla_x J(\theta, x, y))$$

x
“panda”
57.7% confidence

$\text{sign}(\nabla_x J(\theta, x, y))$
“nematode”
8.2% confidence

$x + \epsilon \text{sign}(\nabla_x J(\theta, x, y))$
“gibbon”
99.3 % confidence

[Goodfellow et al., 2014]

- Small but malicious perturbations can result in severe misclassification
- Malicious examples generalize across different architectures
- What is source of instability?
- Can we robustify network?

Kaggle survey: Top Data Science Methods

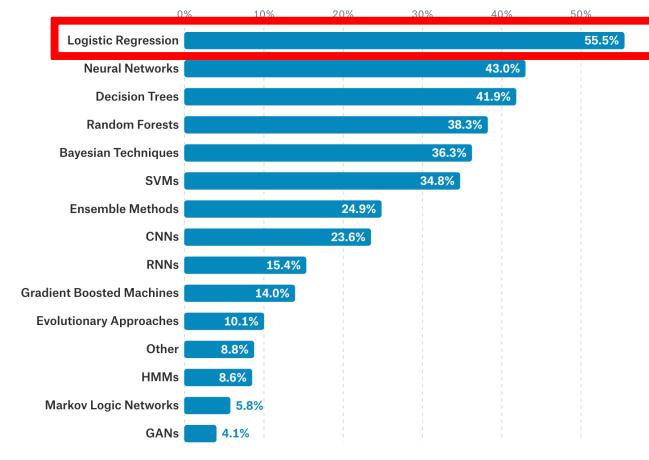
<https://www.kaggle.com/surveys/2017>

Academic

What data science methods are used at work?

Logistic regression is the most commonly reported data science method used at work for all industries except [Military and Security](#) where Neural Networks are used slightly more frequently.

(Company Size ↓) (Academic ↓) (Job Title ↓)



1,201 responses

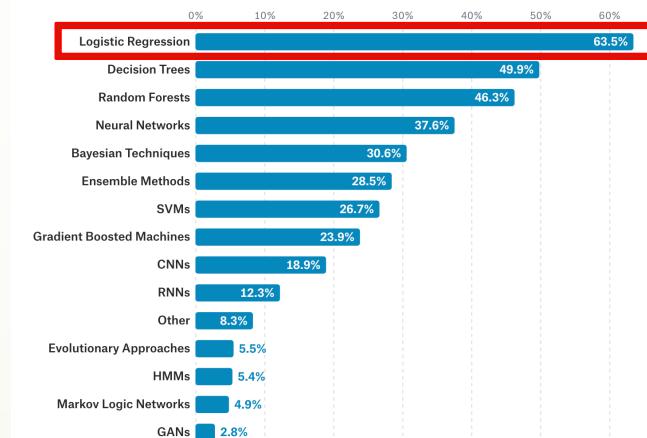
[View code in Kaggle Kernels](#)

Industry

What data science methods are used at work?

Logistic regression is the most commonly reported data science method used at work for all industries except [Military and Security](#) where Neural Networks are used slightly more frequently.

(Company Size ↓) (Industry ↓) (Job Title ↓)



7,301 responses

[View code in Kaggle Kernels](#)

What type of data is used at work?

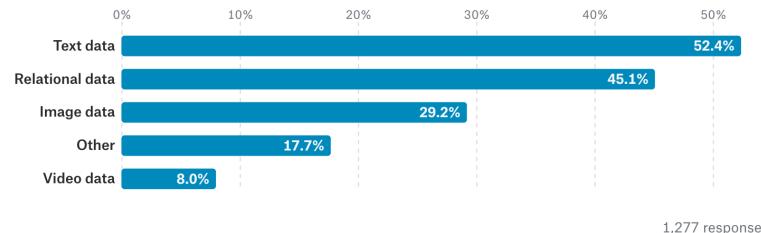
<https://www.kaggle.com/surveys/2017>

Academic

What type of data is used at work?

Relational data is the most commonly reported type of data used at work for all industries except for [Academia](#) and the [Military and Security](#) industry where text data's used more.

Company Size Academic Job Title

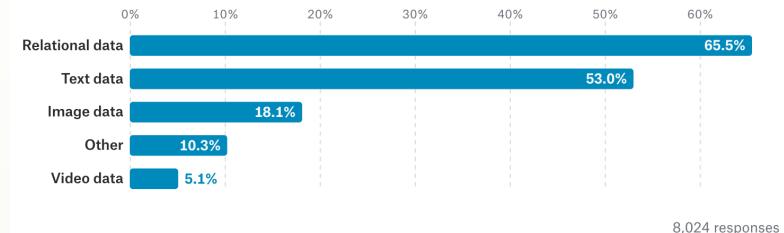


Industry

What type of data is used at work?

Relational data is the most commonly reported type of data used at work for all industries except for [Academia](#) and the [Military and Security](#) industry where text data's used more.

Company Size Industry Job Title



What's wrong with deep learning?

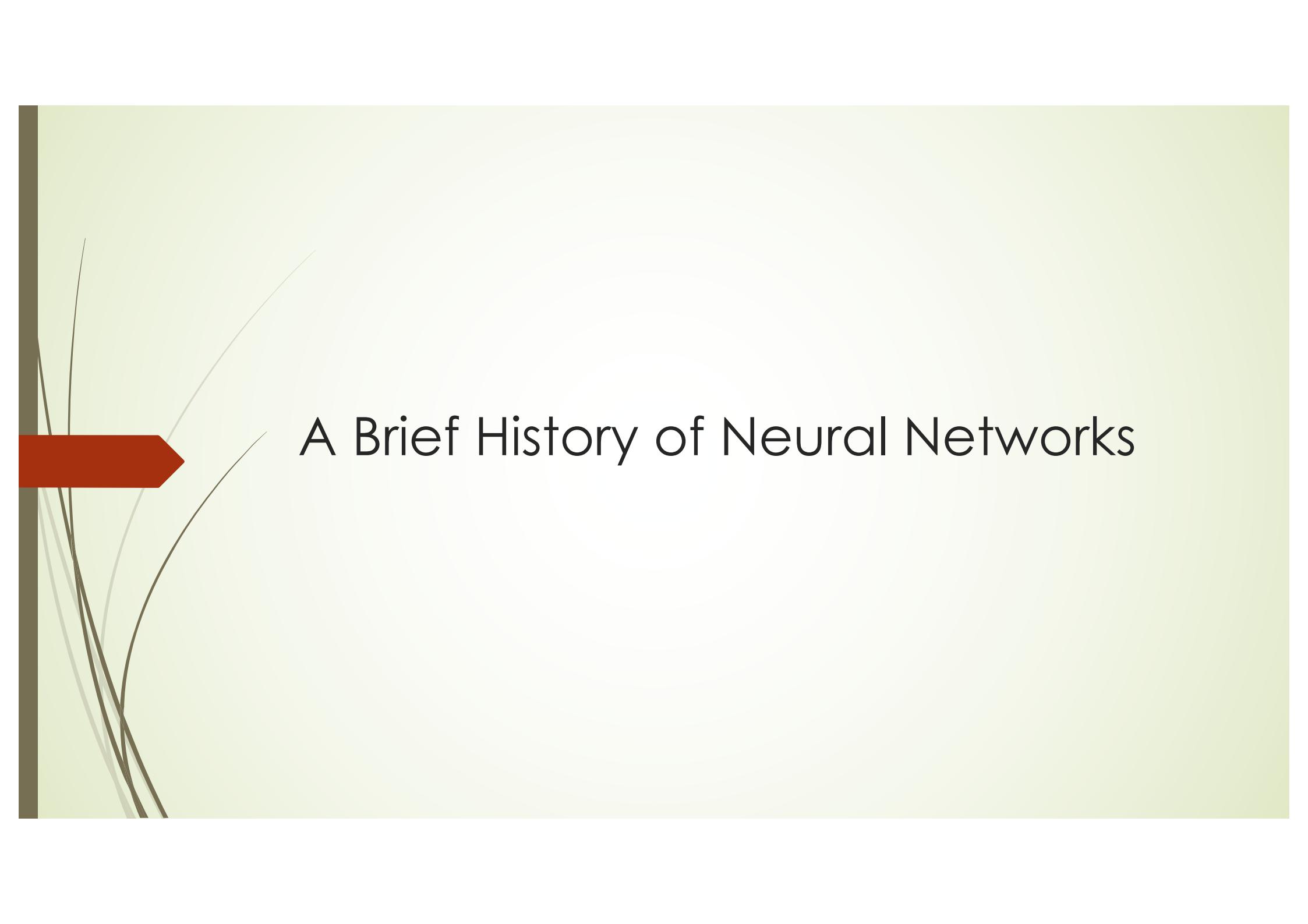
Ali Rahimi NIPS'17: Machine (deep) Learning has become **alchemy**.
<https://www.youtube.com/watch?v=ORHFOnaEzPc>

Yann LeCun CVPR'15, invited talk: **What's wrong with deep learning?**
One important piece: **missing some theory!**

<http://techtalks.tv/talks/whats-wrong-with-deep-learning/61639/>

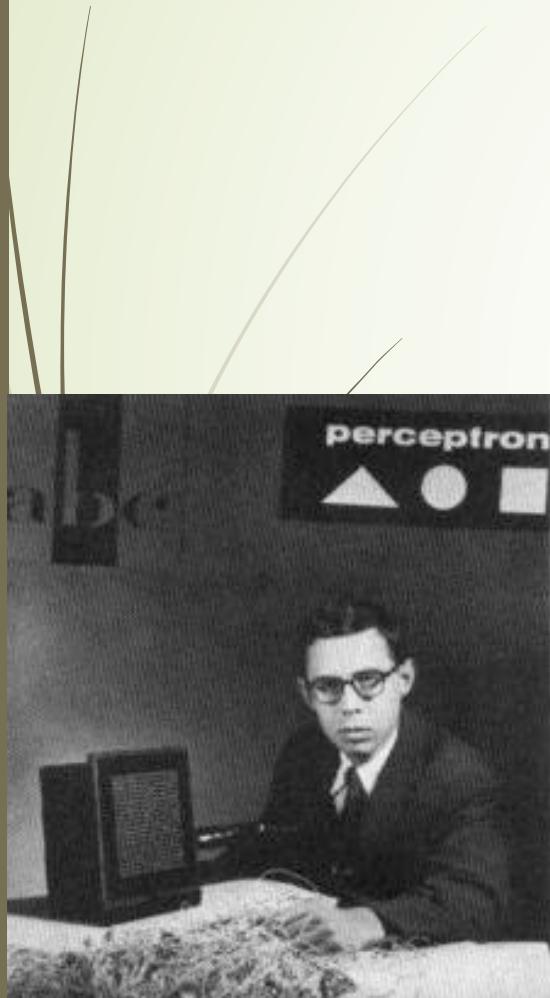


We have to know what we are doing...

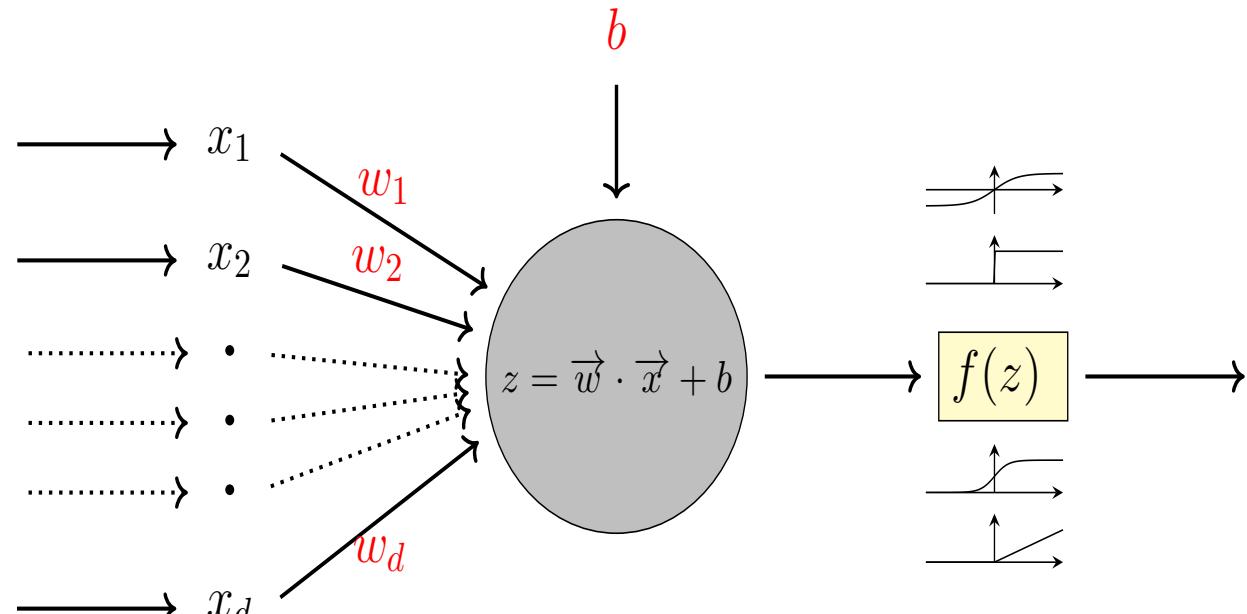


A Brief History of Neural Networks

Perceptron: single-layer



- Invented by Frank Rosenblatt (1957)



The Perceptron Algorithm

$$\ell(w) = - \sum_{i \in \mathcal{M}_w} y_i \langle w, \mathbf{x}_i \rangle, \quad \mathcal{M}_w = \{i : y_i \langle \mathbf{x}_i, w \rangle < 0, y_i \in \{-1, 1\}\}.$$

The Perceptron Algorithm is a *Stochastic Gradient Descent* method
([Robbins–Monro 1950](#); [Kiefer–Wolfowitz 1951](#)) :

$$\begin{aligned} w_{t+1} &= w_t - \eta_t \nabla_i \ell(w) \\ &= \begin{cases} w_t - \eta_t y_i \mathbf{x}_i, & \text{if } y_i w_t^T \mathbf{x}_i < 0, \\ w_t, & \text{otherwise.} \end{cases} \end{aligned}$$

Finite Stop of Perceptron for Separable Data

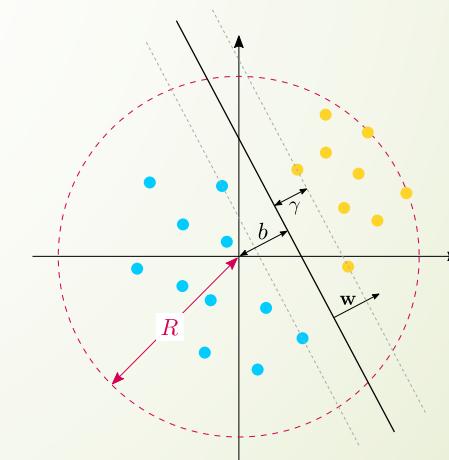
The perceptron convergence theorem was proved by Block (1962) and Novikoff (1962). The following version is based on that in Cristianini and Shawe-Taylor (2000).

Theorem 1 (Block, Novikoff). *Let the training set $S = \{(\mathbf{x}_1, t_1), \dots, (\mathbf{x}_n, t_n)\}$ be contained in a sphere of radius R about the origin. Assume the dataset to be linearly separable, and let \mathbf{w}_{opt} , $\|\mathbf{w}_{\text{opt}}\| = 1$, define the hyperplane separating the samples, having functional margin $\gamma > 0$. We initialise the normal vector as $\mathbf{w}_0 = \mathbf{0}$. The number of updates, k , of the perceptron algorithms is then bounded by*

$$k \leq \left(\frac{2R}{\gamma} \right)^2. \quad (10)$$

Input ball: $R = \max_i \|\mathbf{x}_i\|$.

Margin: $\gamma := \min_i y_i f(\mathbf{x}_i)$





Proof.

Proof. Though the proof can be done using the augmented normal vector and samples defined in the beginning, the notation will be a lot easier if we introduce a different augmentation: $\hat{\mathbf{w}} = (\mathbf{w}^\top, b/R)^\top = (w_1, \dots, w_D, b/R)^\top$ and $\hat{\mathbf{x}} = (\mathbf{x}^\top, R)^\top = (x_1, \dots, x_D, R)^\top$.

Proof (continued, growth of $\|\hat{\mathbf{w}}_k\|$)

We first derive an upper bound on how fast the normal vector grows. As the hyperplane is unchanged if we multiply $\hat{\mathbf{w}}$ by a constant, we can set $\eta = 1$ without loss of generality. Let $\hat{\mathbf{w}}_{k+1}$ be the updated (augmented) normal vector after the k th error has been observed.

$$\|\hat{\mathbf{w}}_{k+1}\|^2 = (\hat{\mathbf{w}}_k + t_i \hat{\mathbf{x}}_i)^T (\hat{\mathbf{w}}_k + t_i \hat{\mathbf{x}}_i) \quad (11)$$

$$= \hat{\mathbf{w}}_k^T \hat{\mathbf{w}}_k + \hat{\mathbf{x}}_i^T \hat{\mathbf{x}}_i + 2t_i \hat{\mathbf{w}}_k^T \hat{\mathbf{x}}_i \quad (12)$$

$$= \|\hat{\mathbf{w}}_k\|^2 + \|\hat{\mathbf{x}}_i\|^2 + 2t_i \hat{\mathbf{w}}_k^T \hat{\mathbf{x}}_i. \quad (13)$$

Since an update was triggered, we know that $t_i \hat{\mathbf{w}}_k^T \hat{\mathbf{x}}_i \leq 0$, thus

$$\|\hat{\mathbf{w}}_k\|^2 + \|\hat{\mathbf{x}}_i\|^2 + 2t_i \hat{\mathbf{w}}_k^T \hat{\mathbf{x}}_i \leq \|\hat{\mathbf{w}}_k\|^2 + \|\hat{\mathbf{x}}_i\|^2 \quad (14)$$

$$= \|\hat{\mathbf{w}}_k\|^2 + (\|\mathbf{x}_i\|^2 + R^2) \quad (15)$$

$$\leq \|\hat{\mathbf{w}}_k\|^2 + 2R^2. \quad (16)$$

This implies that $\|\hat{\mathbf{w}}_k\|^2 \leq 2kR^2$, thus

$$\|\hat{\mathbf{w}}_{k+1}\|^2 \leq 2(k+1)R^2. \quad (17)$$

Proof (continued, projection on $\hat{\mathbf{w}}_{\text{opt}}$)

We then proceed to show how the inner product between an update of the normal vector and $\hat{\mathbf{w}}_{\text{opt}}$ increase with each update:

$$\hat{\mathbf{w}}_{\text{opt}}^T \hat{\mathbf{w}}_{k+1} = \hat{\mathbf{w}}_{\text{opt}}^T \hat{\mathbf{w}}_k + t_i \hat{\mathbf{w}}_{\text{opt}}^T \hat{\mathbf{x}}_i \quad (18)$$

$$\geq \hat{\mathbf{w}}_{\text{opt}}^T \hat{\mathbf{w}}_k + \gamma \quad (19)$$

$$\geq (k+1)\gamma, \quad (20)$$

since $\hat{\mathbf{w}}_{\text{opt}}^T \hat{\mathbf{w}}_k \geq k\gamma$. We therefore have

$$k^2\gamma^2 \leq (\hat{\mathbf{w}}_{\text{opt}}^T \hat{\mathbf{w}}_k)^2 \leq \|\hat{\mathbf{w}}_{\text{opt}}\|^2 \|\hat{\mathbf{w}}_k\|^2 \leq 2kR^2 \|\hat{\mathbf{w}}_{\text{opt}}\|^2, \quad (21)$$

where we have made use of the Cauchy-Schwarz inequality. As $k^2\gamma^2$ grows faster than $2kR^2$, Eq. (21) can hold if and only if

$$k \leq 2\|\hat{\mathbf{w}}_{\text{opt}}\|^2 \frac{R^2}{\gamma^2}. \quad (22)$$

Proof (continued, combined bounds)

As $b \leq R$, we can rewrite the norm of the normal vector:

$$\|\hat{\mathbf{w}}_{\text{opt}}\|^2 = \|\mathbf{w}_{\text{opt}}\|^2 + \frac{b^2}{R^2} \leq \|\mathbf{w}_{\text{opt}}\|^2 + 1 = 2. \quad (23)$$

The bound on k now becomes

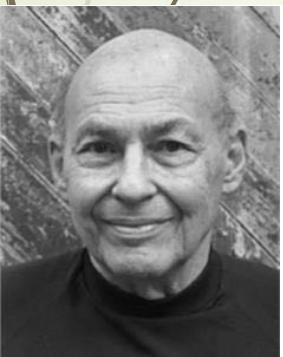
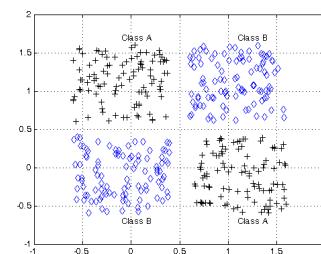
$$k \leq 4 \frac{R^2}{\gamma^2} = \left(\frac{2R}{\gamma} \right)^2, \quad (24)$$

which therefore bounds the number of updates necessary to find the separating hyperplane. \square

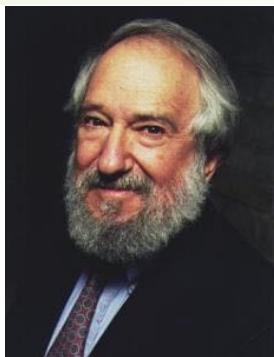
Locality or Sparsity of Computation

Minsky and Papert, 1969

Perceptron can't do **XOR** classification
Perceptron needs infinite global
information to compute **connectivity**

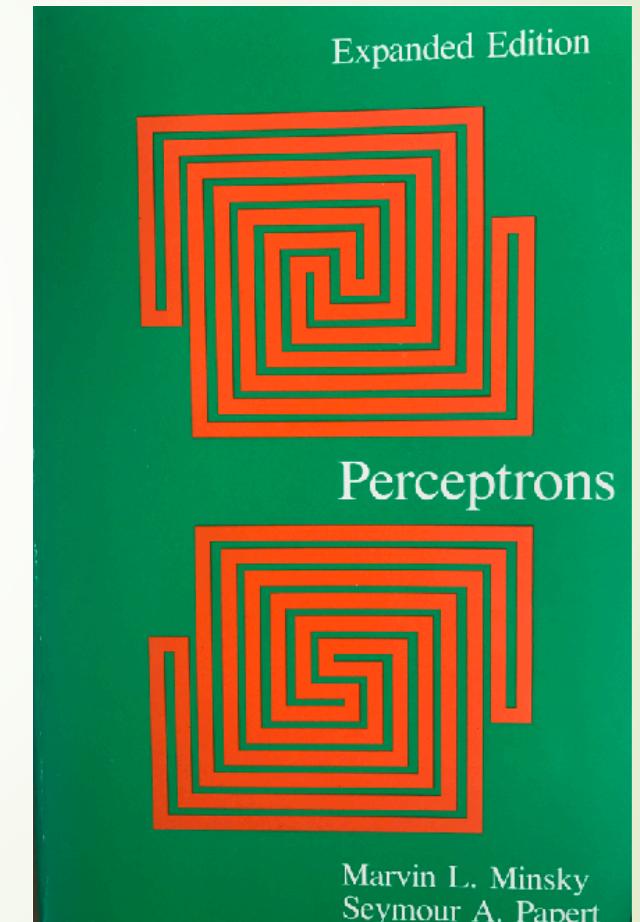


Marvin Minsky



Seymour Papert

Locality or Sparsity is important:
Locality in time?
Locality in space?



Convolutional Neural Networks: shift invariances and locality

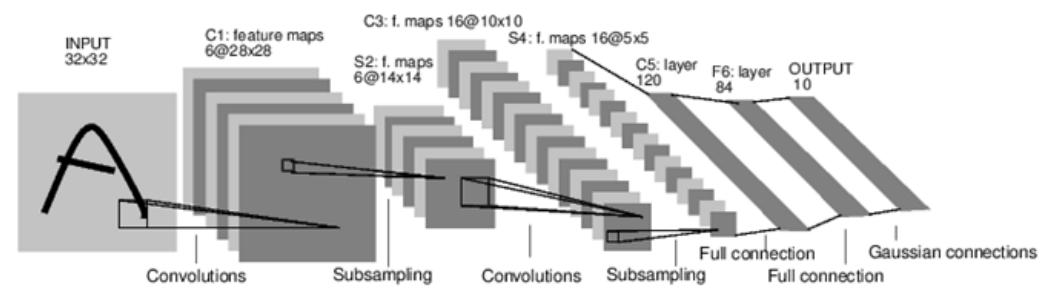
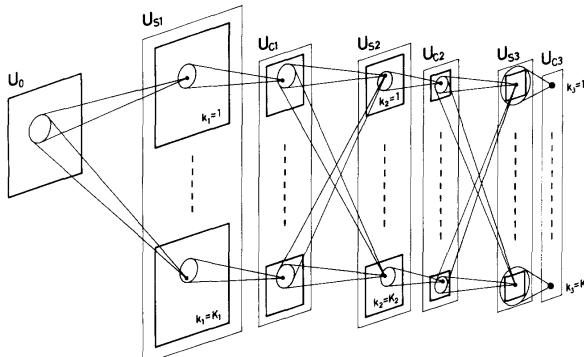
- Can be traced to *Neocognitron* of Kunihiko Fukushima (1979)
- Yann LeCun combined convolutional neural networks with back propagation (1989)
- Imposes **shift invariance** and **locality** on the weights
- Forward pass remains similar
- Backpropagation slightly changes – need to sum over the gradients from all spatial positions

Biol. Cybernetics 36, 193–202 (1980)

Neocognitron: A Self-organizing Neural Network Model for a Mechanism of Pattern Recognition Unaffected by Shift in Position

Kunihiko Fukushima

NHK Broadcasting Science Research Laboratories, Kinuta, Setagaya, Tokyo, Japan



Multilayer Perceptrons (MLP) and Back-Propagation (BP) Algorithms

Rumelhart, Hinton, Williams (1986)

Learning representations by back-propagating errors, Nature, 323(9): 533-536

BP algorithms as **stochastic gradient descent** algorithms (**Robbins–Monro 1950; Kiefer–Wolfowitz 1951**) with Chain rules of Gradient maps

MLP classifies **XOR**, but the global hurdle on topology (connectivity) computation still exists

NATURE VOL. 323 5 OCTOBER 1986 LETTERS TO NATURE 533

Learning representations by back-propagating errors

David E. Rumelhart*, Geoffrey E. Hinton† & Ronald J. Williams*

* Institute for Cognitive Science, C-015, University of California, San Diego, La Jolla, California 92093, USA
† Department of Computer Science, Carnegie-Mellon University, Pittsburgh, Pennsylvania 15213, USA

We describe a new learning procedure, back-propagation, for networks of neurone-like units. The procedure repeatedly adjusts the weights of the connections in the network so as to minimize a measure of the difference between the actual output of the net and the desired output vector. As a result of the weight adjustments, internal 'hidden' units which are not part of the input or output come to represent important features of the task domain, and the regularities in the task are captured by the interactions of these units. The ability to create useful new features distinguishes back-propagation from earlier, simpler methods such as the perceptron-convergence procedure.¹

There have been many attempts to design self-organizing neural networks. The aim is to find a powerful synaptic modification rule that will allow an arbitrarily connected neural network to develop an internal structure that is appropriate for a particular task domain. The task is given by a given desired state vector of the output units for each state vector of the input units. If the input units are directly connected to the output units it is relatively easy to find learning rules that iteratively adjust the relative strengths of the connections so as to progressively reduce the difference between the actual and desired output vectors.² Learning becomes more interesting but

The simplest form of the learning procedure is for layered networks which have a layer of input units at the bottom; any number of intermediate layers of units, and a single layer at the top. Connections within a layer or from higher to lower layers are forbidden, but connections can skip intermediate layers. An input vector is presented to the network by setting the states of the input units. Then the states of the units in each layer are determined by applying equations (1) and (2) to the connections from the previous layer. All units within a layer set their states set in parallel. They then pass their states sequentially, starting at the bottom and working upwards until the states of the output units are determined.

The total input, x_j , to unit j is a linear function of the outputs, y_i , of the units that are connected to j and of the weights, w_{ij} , on these connections

$$x_j = \sum_i y_i w_{ij} \quad (1)$$

Units can be given biases by introducing an extra input to each unit which always has a value of 1. The weight on this extra input is called the bias and is equivalent to a threshold of the opposite sign. It can be treated just like the other weights.

A unit has a real-valued output, y_j , which is a non-linear function of its total input

$$y_j = \frac{1}{1 + e^{-x_j}} \quad (2)$$

¹ To whom correspondence should be addressed

BP Algorithm: Forward Pass

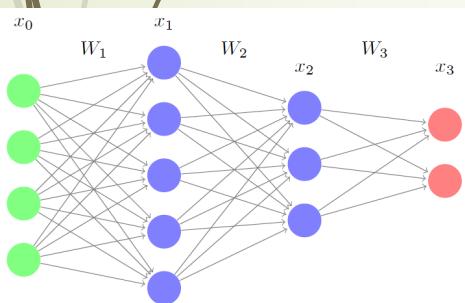
- Cascade of repeated [linear operation followed by coordinatewise nonlinearity]'s
- Nonlinearities: sigmoid, hyperbolic tangent, (recently) ReLU.

Algorithm 1 Forward pass

Input: x_0

Output: x_L

```
1: for  $\ell = 1$  to  $L$  do
2:    $x_\ell = f_\ell(W_\ell x_{\ell-1} + b_\ell)$ 
3: end for
```



BP algorithm = Gradient Descent Method

- Training examples $\{x_0^i\}_{i=1}^n$ and labels $\{y^i\}_{i=1}^n$
- Output of the network $\{x_L^i\}_{i=1}^m$
- Objective

$$J(\{W_l\}, \{b_l\}) = \frac{1}{n} \sum_{i=1}^n \frac{1}{2} \|y^i - x_L^i\|_2^2 \quad (1)$$

Other losses include cross-entropy, logistic loss, exponential loss, etc.

- Gradient descent

$$W_l = W_l - \eta \frac{\partial J}{\partial W_l}$$
$$b_l = b_l - \eta \frac{\partial J}{\partial b_l}$$

In practice: use Stochastic Gradient Descent (SGD)

Derivation of BP: Lagrangian Multiplier

LeCun et al. 1988

Given n training examples $(I_i, y_i) \equiv (\text{input}, \text{target})$ and L layers

- Constrained optimization

$$\min_{W,x} \quad \sum_{i=1}^n \|x_i(L) - y_i\|_2$$

$$\text{subject to} \quad x_i(\ell) = f_\ell \left[W_\ell x_i(\ell-1) \right], \\ i = 1, \dots, n, \quad \ell = 1, \dots, L, \quad x_i(0) = I_i$$

- Lagrangian formulation (Unconstrained)

$$\min_{W,x,B} \mathcal{L}(W, x, B)$$

$$\mathcal{L}(W, x, B) = \sum_{i=1}^n \left\{ \|x_i(L) - y_i\|_2^2 + \sum_{\ell=1}^L B_i(\ell)^T \left(x_i(\ell) - f_\ell \left[W_\ell x_i(\ell-1) \right] \right) \right\}$$

<http://yann.lecun.com/exdb/publis/pdf/lecun-88.pdf>

back-propagation – derivation

- $\frac{\partial \mathcal{L}}{\partial B}$

Forward pass

$$x_i(\ell) = f_\ell \left[\underbrace{W_\ell x_i(\ell-1)}_{A_i(\ell)} \right] \quad \ell = 1, \dots, L, \quad i = 1, \dots, n$$

- $\frac{\partial \mathcal{L}}{\partial x}, z_\ell = [\nabla f_\ell] B(\ell)$

Backward (adjoint) pass

$$z(L) = 2\nabla f_L \left[A_i(L) \right] (y_i - x_i(L))$$

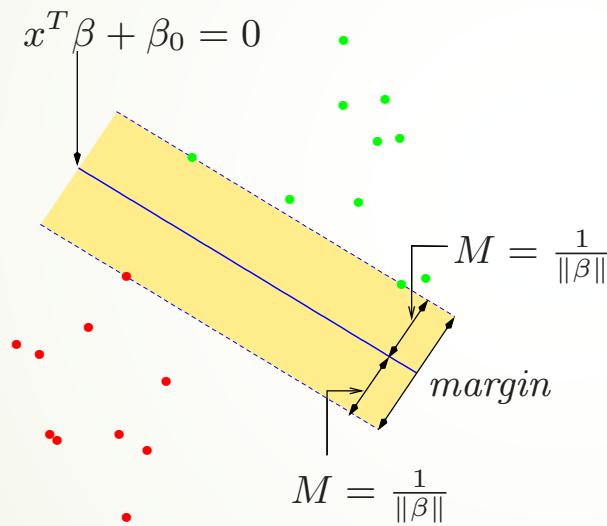
$$z_i(\ell) = \nabla f_\ell \left[A_i(\ell) \right] W_{\ell+1}^T z_i(\ell+1) \quad \ell = 0, \dots, L-1$$

- $W \leftarrow W + \lambda \frac{\partial \mathcal{L}}{\partial W}$

Weight update

$$W_\ell \leftarrow W_\ell + \lambda \sum_{i=1}^n z_i(\ell) x_i^T(\ell-1)$$

Max-Margin Classifier (SVM)

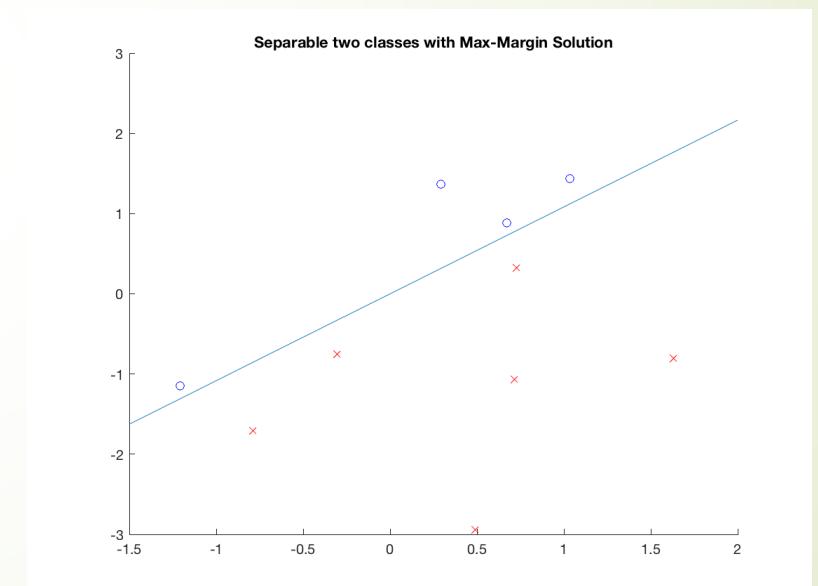


Vladimir Vapnik, 1994

Convex optimization + Reproducing Kernel Hilbert Spaces (Grace Wahba etc.)

$$\text{minimize}_{\beta_0, \beta_1, \dots, \beta_p} \|\beta\|^2 := \sum_j \beta_j^2$$

subject to $y_i(\beta_0 + \beta_1 x_{i1} + \dots + \beta_p x_{ip}) \geq 1$ for all i



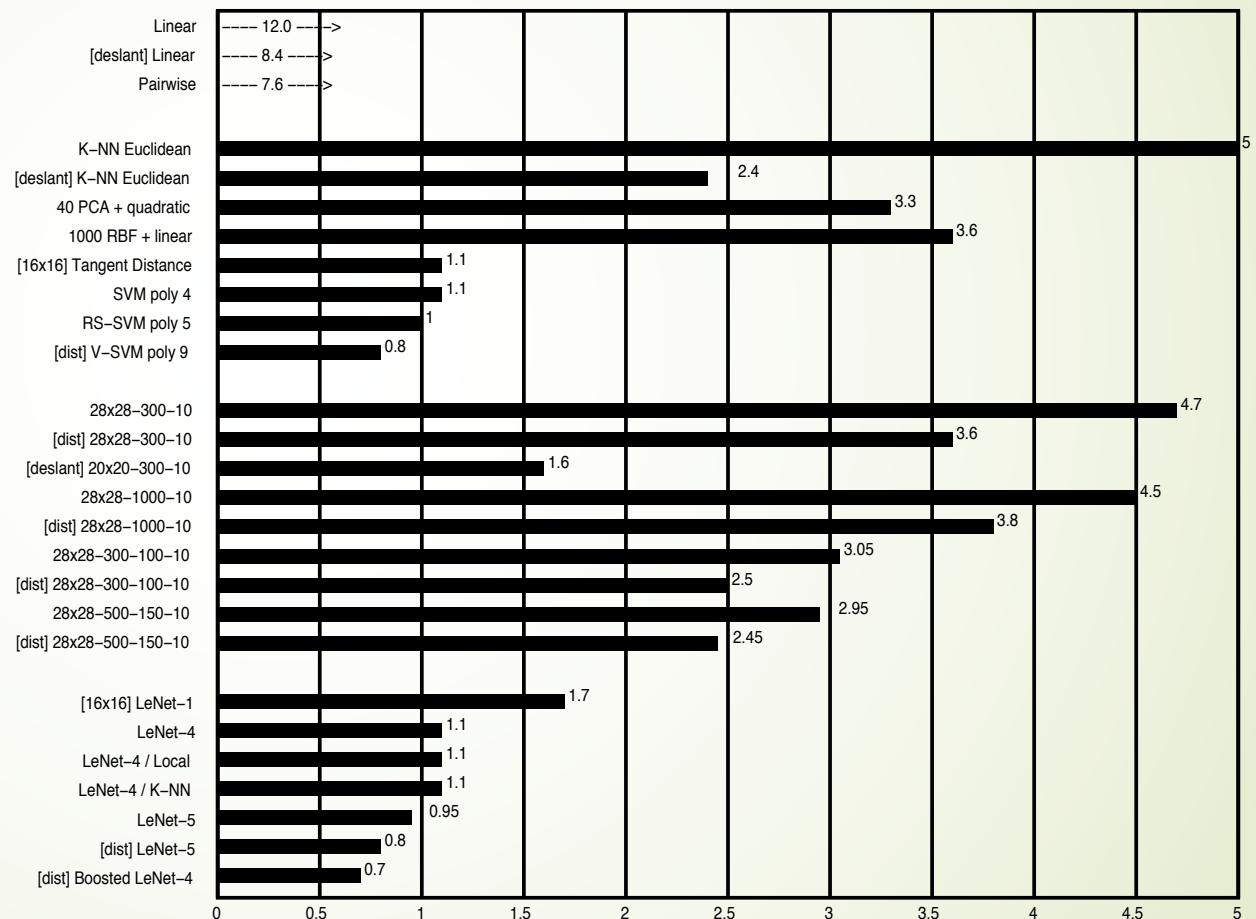
MNIST Dataset Test Error: SVM vs. CNN

LeCun et al. 1998



Simple SVM performs as well as Multilayer Convolutional Neural Networks which need careful tuning (LeNets)

Second dark era for NN:
2000s

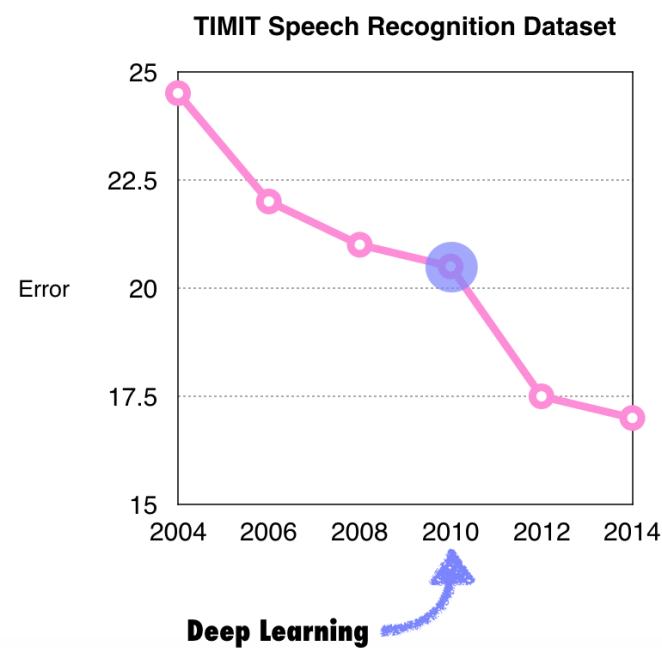


2000-2010: The Era of SVM, Boosting, ... as nights of Neural Networks



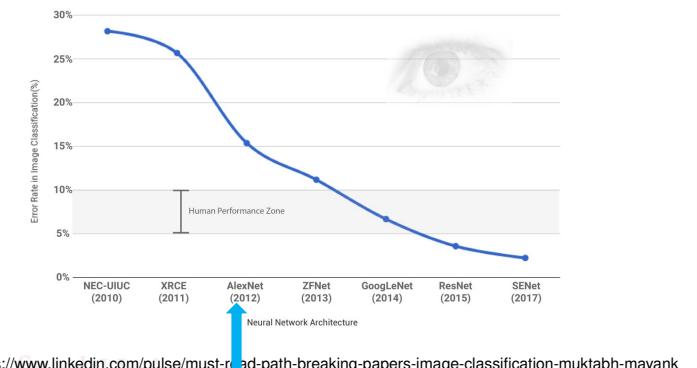
Around the year of 2012...

Speech Recognition: TIMIT



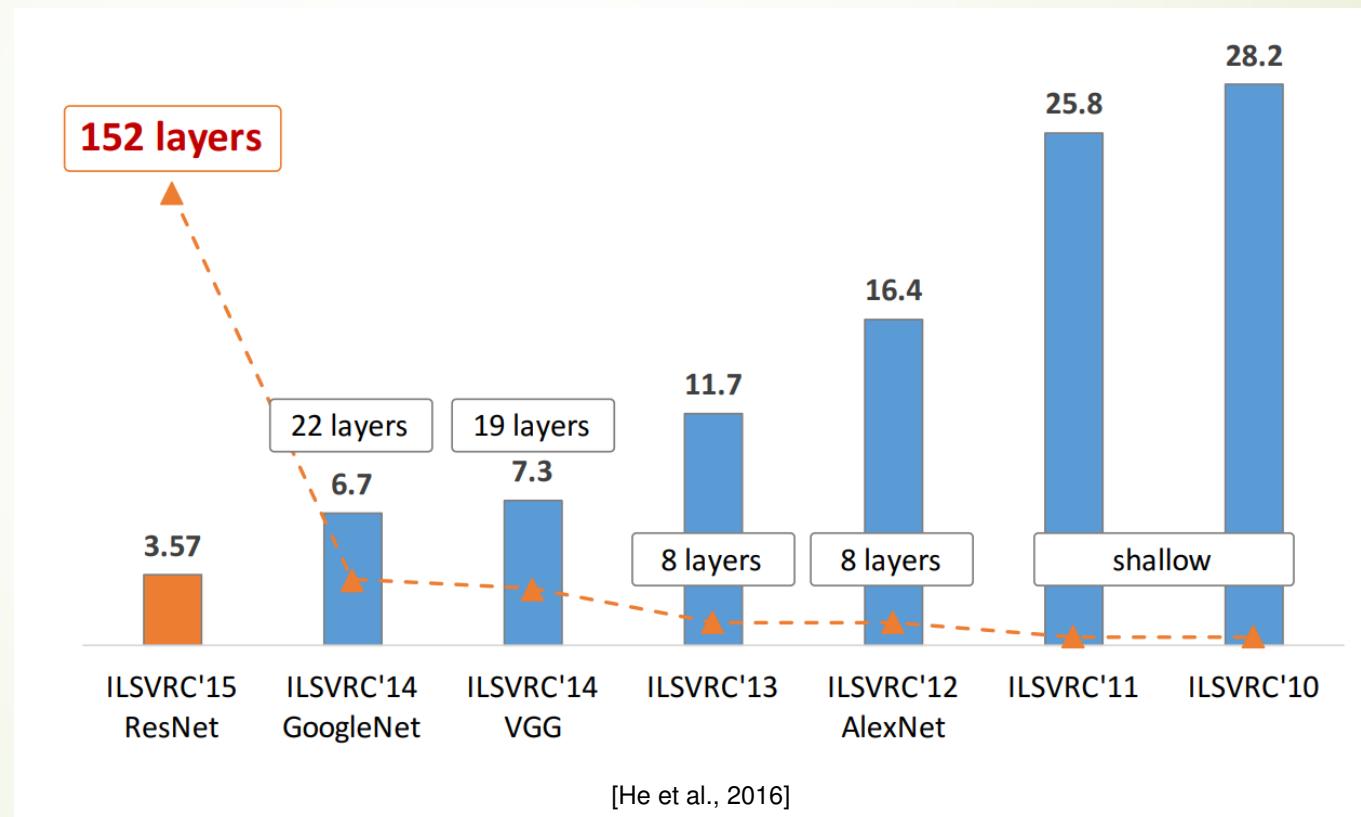
Computer Vision: ImageNet

- ImageNet (subset):
 - 1.2 million training images
 - 100,000 test images
 - 1000 classes
- ImageNet large-scale visual recognition Challenge



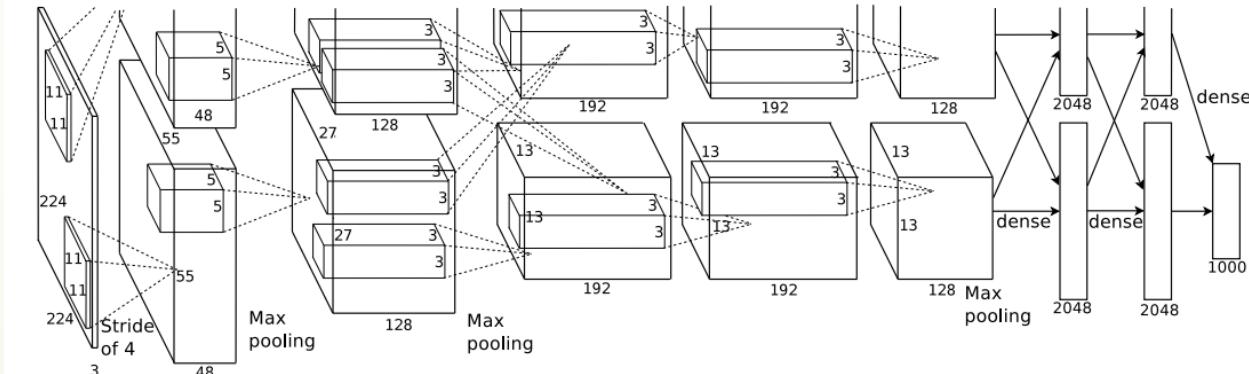
Deep Learning

Depth as function of year



AlexNet (2012): Architecture

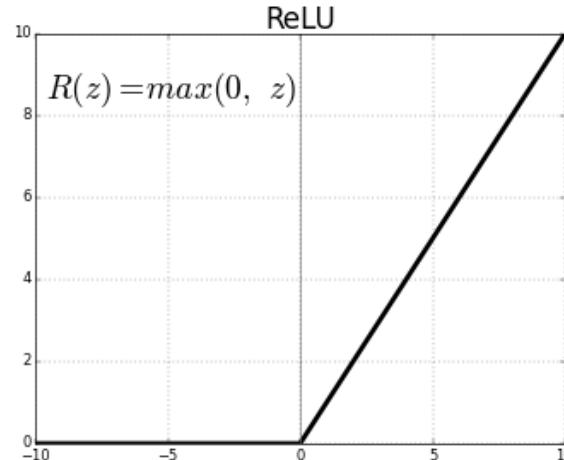
- 8 layers: first 5 convolutional, rest fully connected
- ReLU nonlinearity
- Local response normalization
- Max-pooling
- Dropout



Source: [Krizhevsky et al., 2012]

AlexNet (2012): ReLU

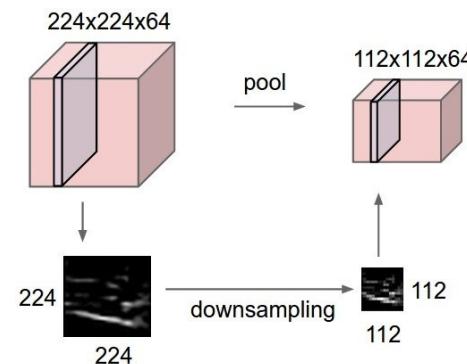
- Non-saturating function and therefore faster convergence when compared to other nonlinearities
- Problem of dying neurons



Source: https://ml4a.github.io/ml4a/neural_networks/

AlexNet (2012): Max Pooling

- Chooses maximal entry in every non-overlapping window of size 2×2 , for example



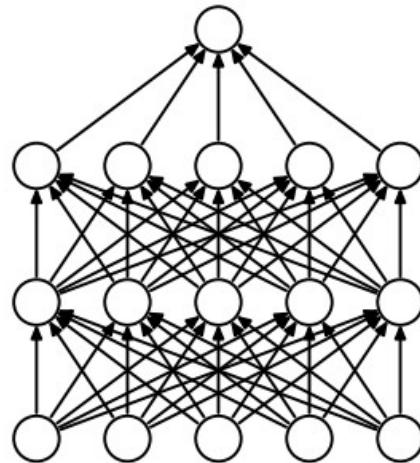
12	20	30	0
8	12	2	0
34	70	37	4
112	100	25	12

$\xrightarrow{2 \times 2 \text{ Max-Pool}}$

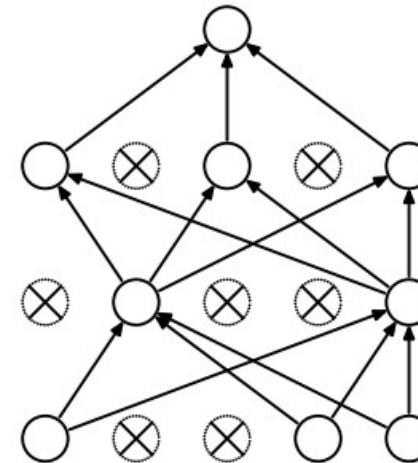
20	30
112	37

Source: Stanford's CS231n github

AlexNet (2012): Dropout



(a) Standard Neural Net



(b) After applying dropout.

Source: [Srivastava et al., 2014]

- Zero every neuron with probability $1 - p$
- At test time, multiply every neuron by p

AlexNet (2012): Training

- Stochastic gradient descent
- Mini-batches
- Momentum
- Weight decay (ℓ_2 prior on the weights)

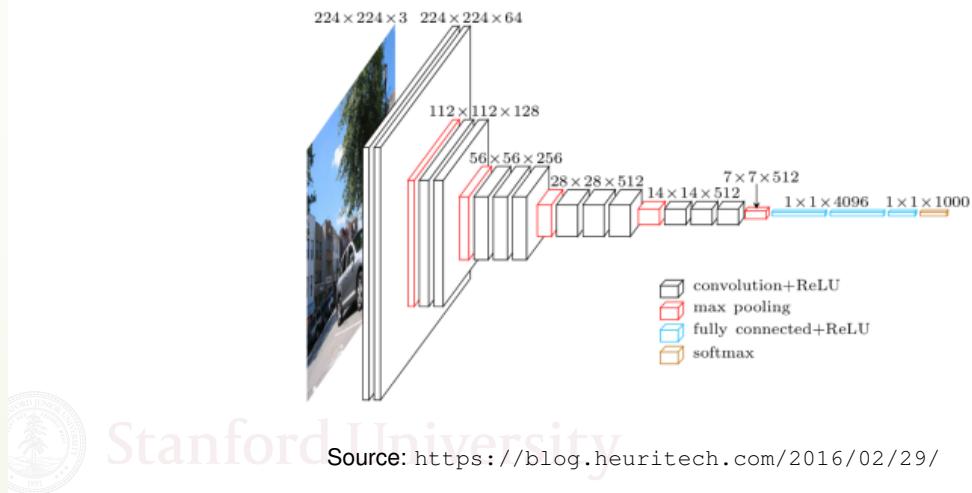


Filters trained in the first layer

Source: [Krizhevsky et al., 2012]

VGG (2014) [Simonyan-Zisserman'14]

- Deeper than AlexNet: 11-19 layers versus 8
- No local response normalization
- Number of filters multiplied by two every few layers
- Spatial extent of filters 3×3 in all layers
- Instead of 7×7 filters, use three layers of 3×3 filters
 - Gain intermediate nonlinearity
 - Impose a regularization on the 7×7 filters

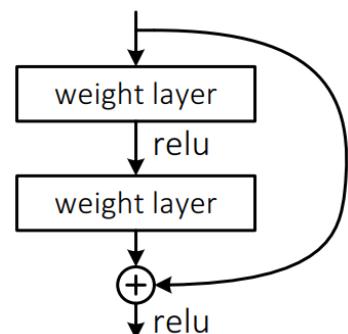


Stanford University

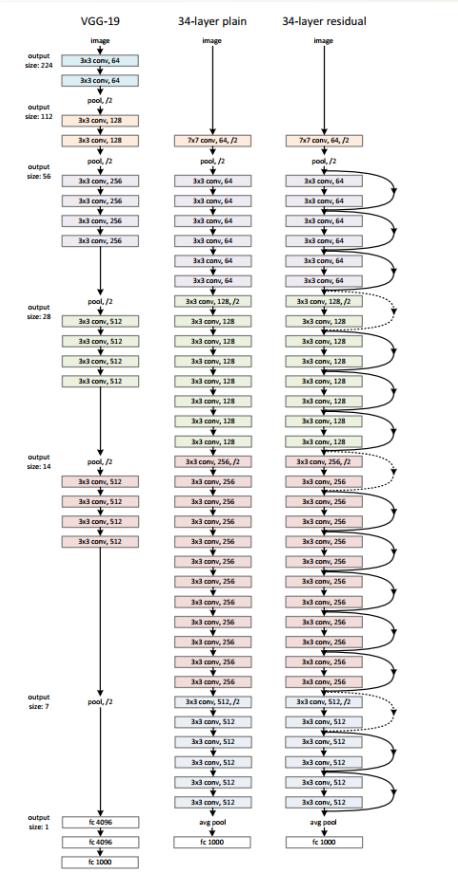
Source: <https://blog.heuritech.com/2016/02/29/>

ResNet (2015) [HGRS-15]

- Solves problem by adding skip connections
- Very deep: 152 layers
- No dropout
- Stride
- Batch normalization



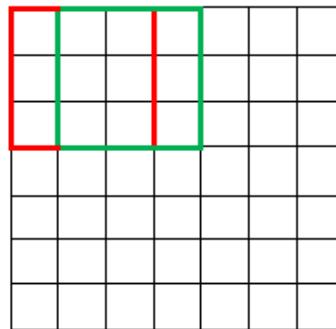
Source: Deep Residual Learning for Image Recognition



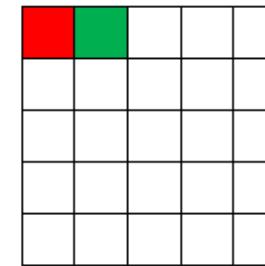
Stride

Stride 1

7 x 7 Input Volume

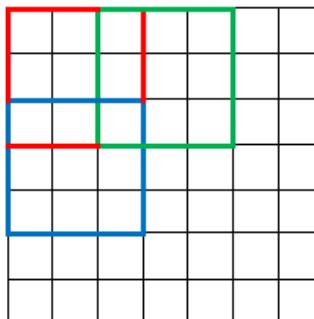


5 x 5 Output Volume

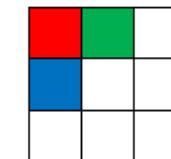


Stride 2

7 x 7 Input Volume



3 x 3 Output Volume



Batch Normalization

Algorithm 2 Batch normalization [Ioffe and Szegedy, 2015]

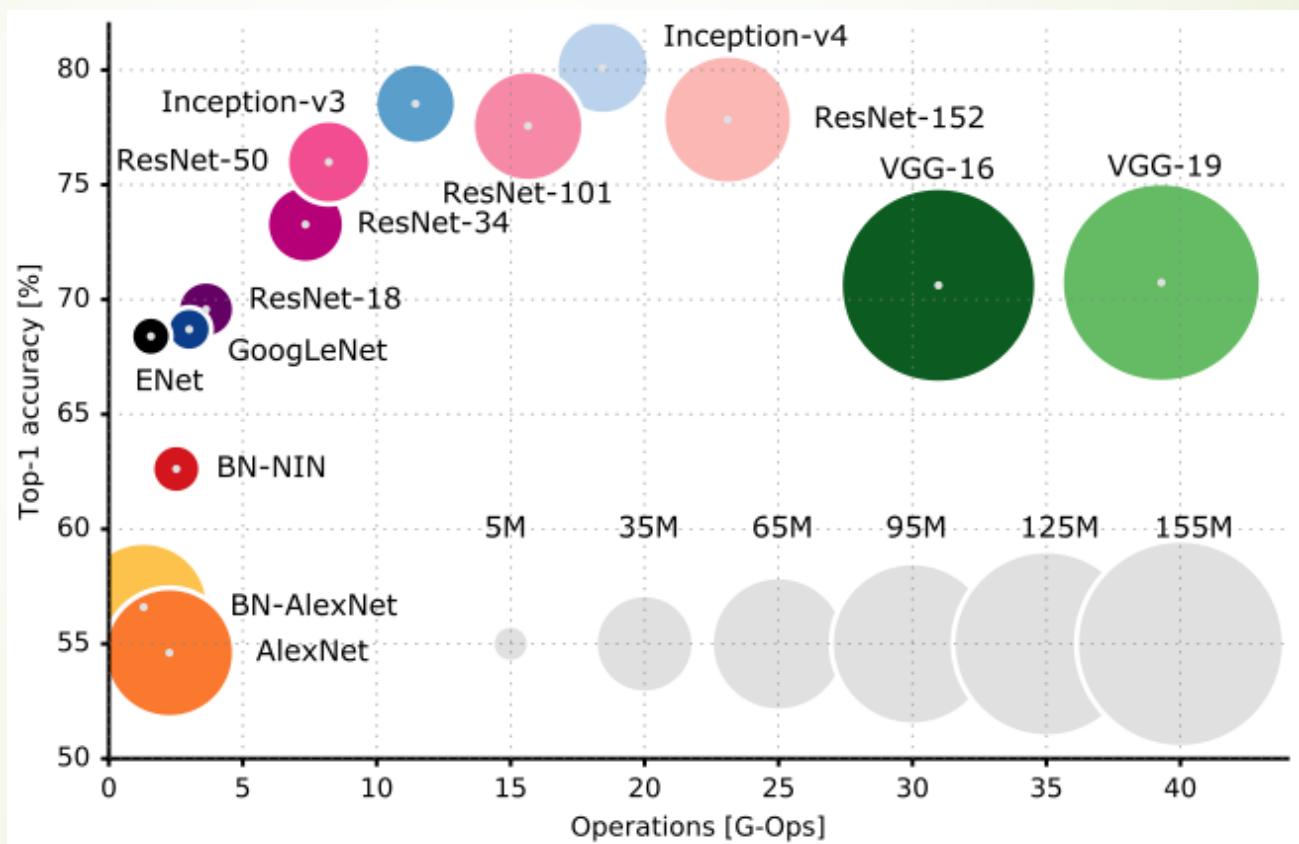
Input: Values of x over minibatch $x_1 \dots x_B$, where x is a certain channel in a certain feature vector

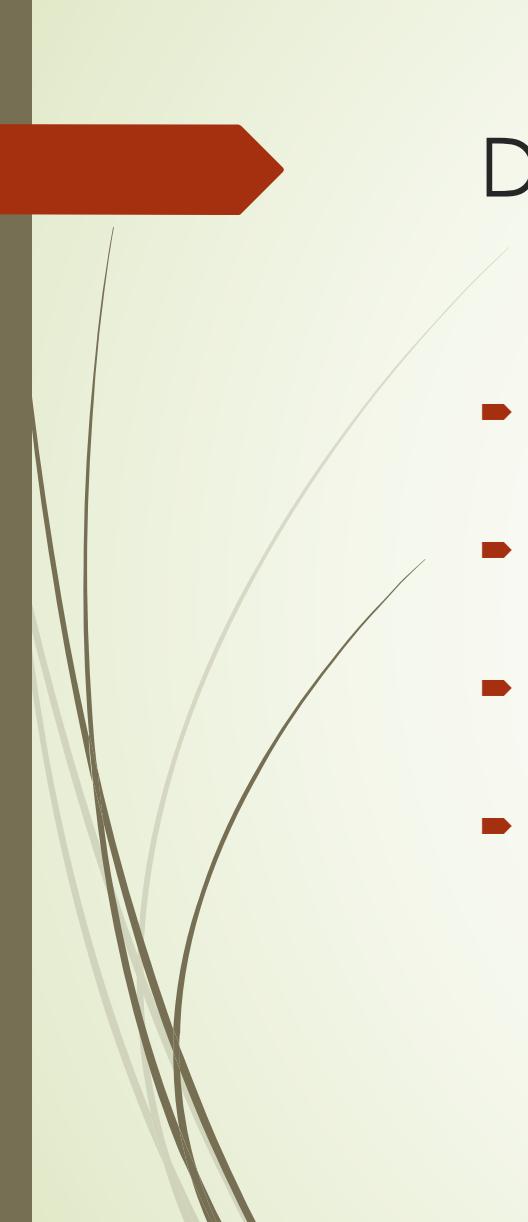
Output: Normalized, scaled and shifted values $y_1 \dots y_B$

- 1: $\mu = \frac{1}{B} \sum_{b=1}^B x_b$
 - 2: $\sigma^2 = \frac{1}{B} \sum_{b=1}^B (x_b - \mu)^2$
 - 3: $\hat{x}_b = \frac{x_b - \mu}{\sqrt{\sigma^2 + \epsilon}}$
 - 4: $y_b = \gamma \hat{x}_b + \beta$
-

- Accelerates training and makes initialization less sensitive
- Zero mean and unit variance feature vectors

Complexity vs. Accuracy of Different Networks





Deep Learning Softwares

- ▶ Pytorch (developed by Yann LeCun and Facebook):
 - ▶ <http://pytorch.org/tutorials/>
- ▶ Tensorflow (developed by Google based on Caffe)
 - ▶ <https://www.tensorflow.org/tutorials/>
- ▶ Theano (developed by Yoshua Bengio)
 - ▶ <http://deeplearning.net/software/theano/tutorial/>
- ▶ Keras (based on Tensorflow or Pytorch)
 - ▶ https://www.manning.com/books/deep-learning-with-python?a_aid=keras&a_bid=76564dff