

# Computational Data Analysis

## Machine Learning

**Yao Xie, Ph.D.**

*Associate Professor*

Harold R. and Mary Anne Nash Early Career Professor  
H. Milton Stewart School of Industrial and Systems  
Engineering

Introduction to Neural Networks

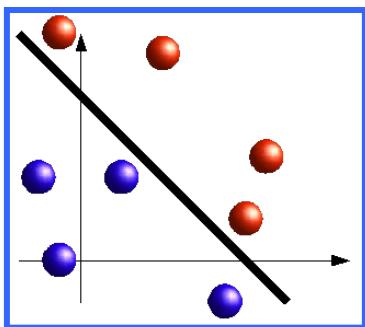


# Main approaches to design classifiers

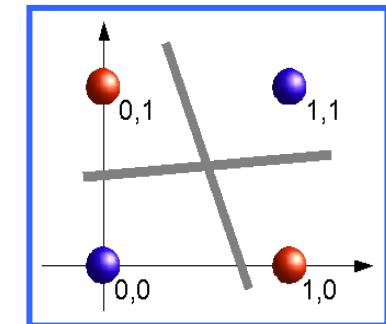
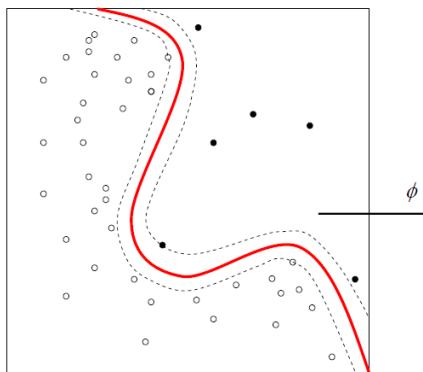
- Bayes rule, use simplifying assumption for  $p(x|y = 1)$ 
  - Assume  $p(x|y = 1)$  is Gaussian
  - Assume  $p(x|y = 1)$  is fully factorized
- Use geometric intuitions
  - k-nearest neighbor classifier
  - Support vector machine
- Directly go for the decision boundary  $h(x) = -\ln \frac{q_i(x)}{q_j(x)}$ 
  - Logistic regression
  - Neural networks

# Learning nonlinear decision boundary

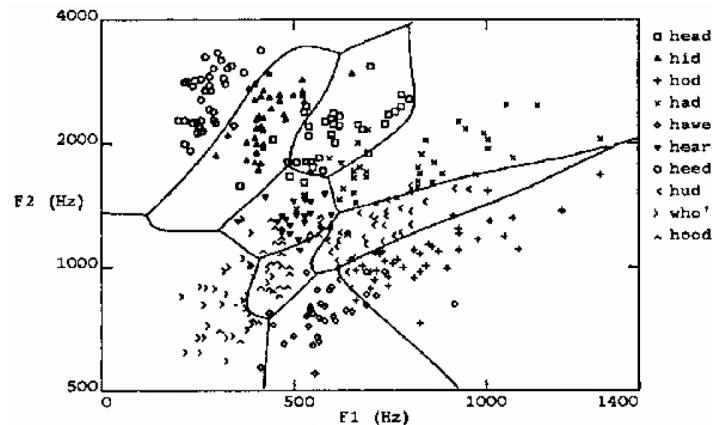
Linearly separable



Nonlinearly separable



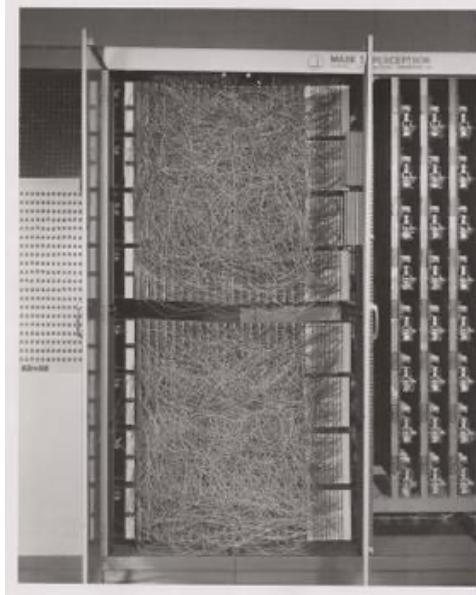
The XOR gate



Speech recognition

# Early artificial intelligence

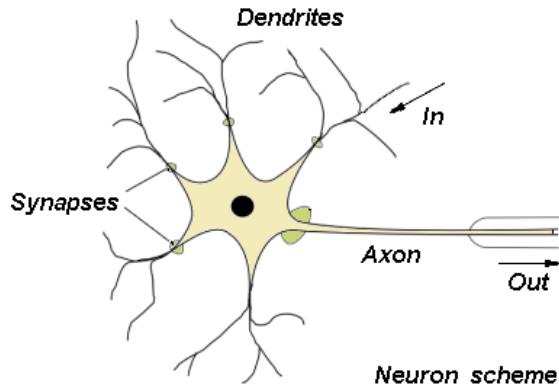
- The perceptron algorithm was first invented in 1958 at the Cornell Aeronautical Lab
- This machine was designed for image recognition: it had an array of 400 photocells, randomly connected to the "neurons".



Mark I Perceptron machine, the first implementation of the perceptron algorithm. It was connected to a camera with  $20 \times 20$  cadmium sulfide photocells to make a 400-pixel image. The main visible feature is a patch panel that set different combinations of input features. To the right, arrays of potentiometers that implemented the adaptive weights.

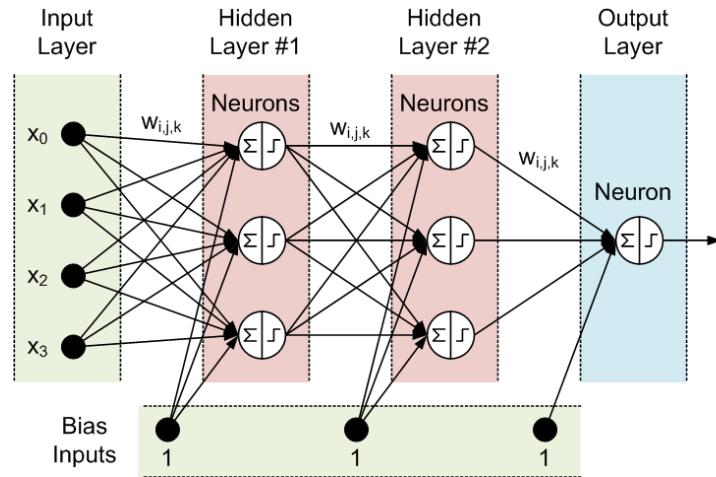
# Perceptron

- From biological neurons to artificial neurons (perceptron)

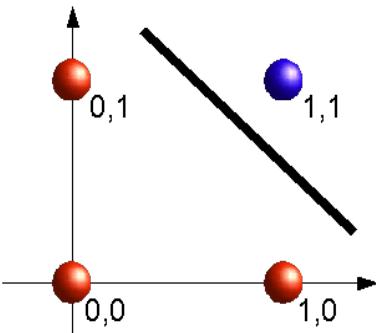


Human brain:  
Number of  
neurons  $\sim 10^{10}$

- Artificial neural networks (ANN)
  - Many neuron-like threshold switching units
  - Weighted interconnections among units at layers (feed-forward)



# Perceptron



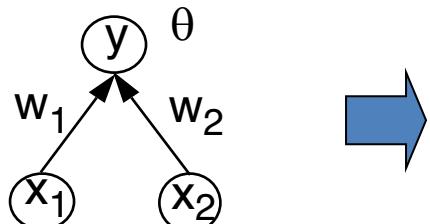
NAND

x	y	Z (color)
0	0	0
0	1	0
1	0	0
1	1	1

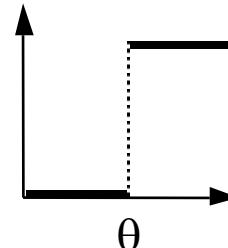
Threshold  $\theta = 0.5$

some possible values for  $w_1$  and  $w_2$

$w_1$	$w_2$
0.20	0.35
0.20	0.40
0.25	0.30
0.40	0.20

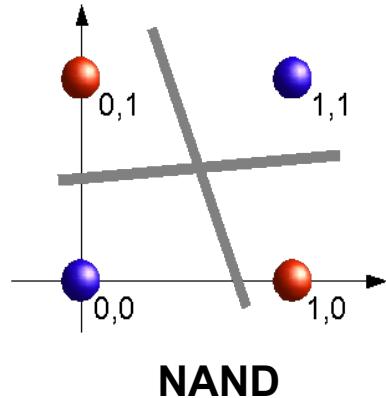


$$f(x_1w_1 + x_2w_2) = y$$
$$f(0w_1 + 0w_2) = 1$$
$$f(0w_1 + 1w_2) = 1$$
$$f(1w_1 + 0w_2) = 1$$
$$f(1w_1 + 1w_2) = 0$$

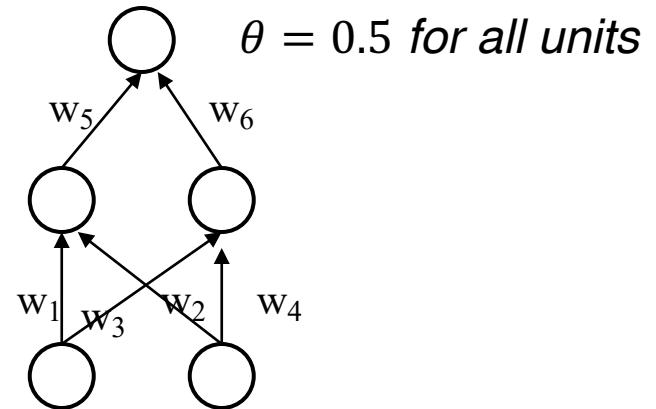


$$f(u) = \begin{cases} 1, & \text{for } u > \theta \\ 0, & \text{for } u \leq \theta \end{cases}$$

# Connecting perceptron can give more complex decision boundary



x	y	Z (color)
0	0	0
0	1	1
1	0	1
1	1	0



a possible set of values for  
 $(w_1, w_2, w_3, w_4, w_5, w_6)$ :  
 $(0.6, -0.6, -0.7, 0.8, 1, 1)$

# Expressive capability of ANN

- Boolean functions:
  - Every Boolean function can be represented by network with single hidden layer
  - But might require exponential (in number of inputs) hidden units
- Continuous functions:
  - Every bounded continuous function can be approximated with arbitrary small error, by network with one hidden layer [Cybenko 1989; Hornik et al 1989]
  - Any function can be approximated to arbitrary accuracy by a network with two hidden layers [Cybenko 1988].
  - Many recent developments for more general activation functions (e.g. ReLU)

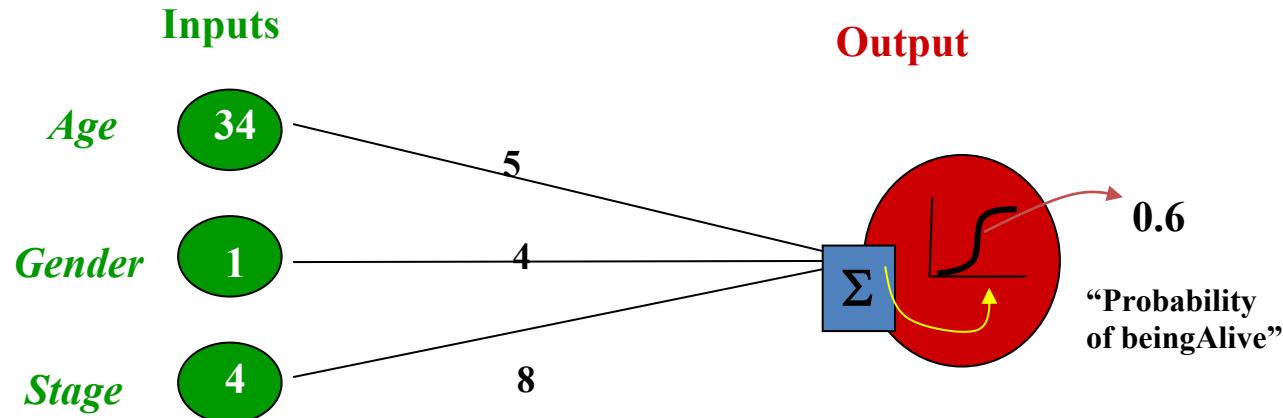
# Logistic regression can be viewed as a “single neuron”

Independent variable = input variable

Dependent variable = output variable

Coefficients = “weights” + “bias”

**Logistic Regression Model (the sigmoid unit) is a perceptron**



*Independent variables*

$x_1, x_2, x_3$

*Coefficients*

$a, b, c$

*Dependent variable*

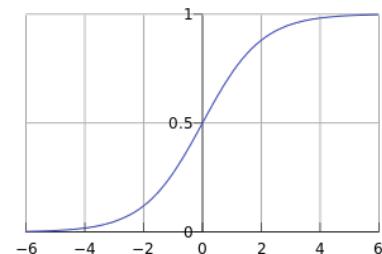
$p$  *Prediction*

# Connection to logistic regression

- Let the original feature be denoted as  $u$  ( $n$  dimensional vector)
- Note that we can define the augmented feature as  $x = [1; u^T]^T$  (a  $n + 1$  dimensional vector)
- The weight vector  $\tilde{w} = [b, w^T]^T$
- Assume that the posterior distribution  $p(y = 1|x)$  take a particular form

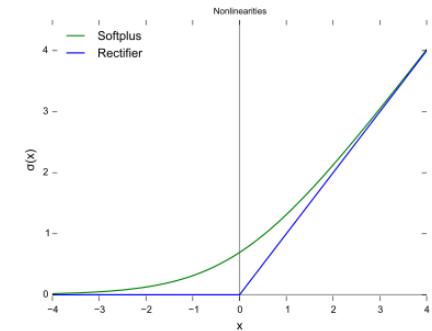
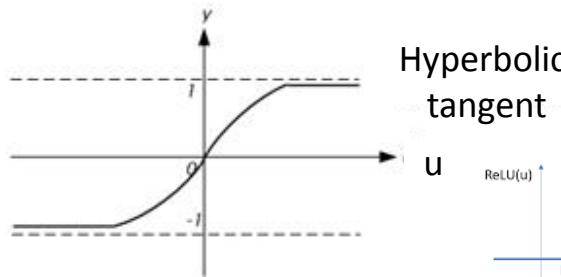
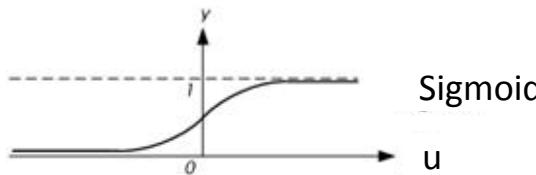
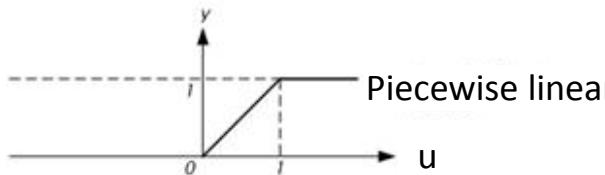
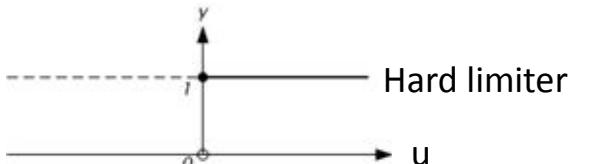
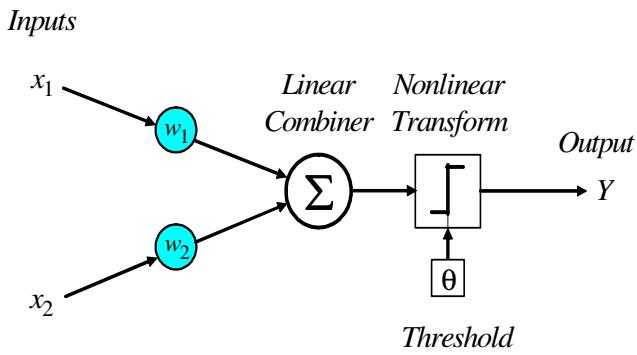
$$p(y = 1|x, \tilde{w}) = \frac{1}{1 + \exp(-\tilde{w}^T x)} = \sigma(\tilde{w}^T x)$$

- Logistic function (or sigmoid function)  $\sigma(u) = \frac{1}{1+\exp(-u)}$
- Logistic regression can be viewed as one neuron



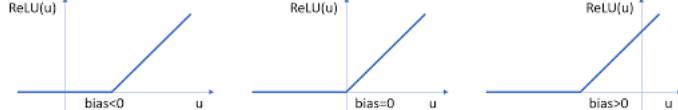
# Other nonlinear neurons

- Use different nonlinear transformations  $f(u)$
- Before that, perform weighted combination of inputs  $u = w^T x$



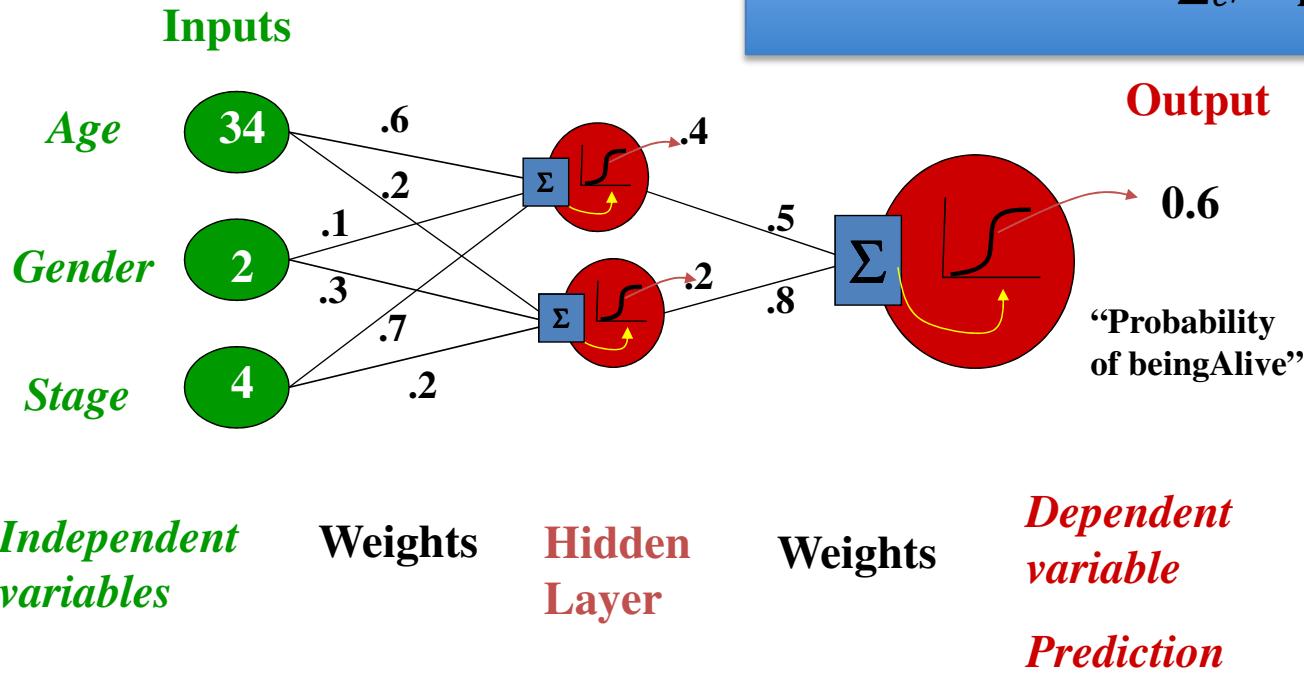
Rectifier linear unit (ReLU)  
 $f(u) = \max\{u, 0\}$

With bias  
 $f(u; \theta) = \max\{u + \theta, 0\}$



# Deep Neural Network Model

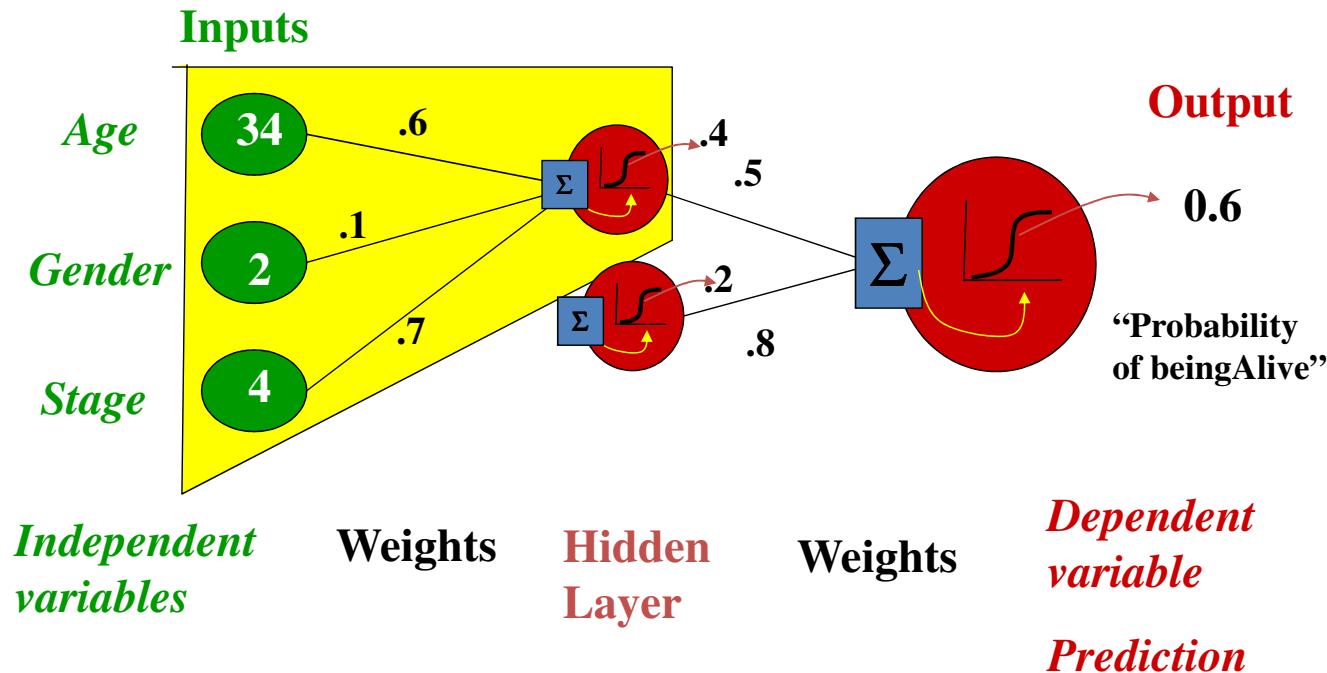
Multiple layers of interconnected perceptrons



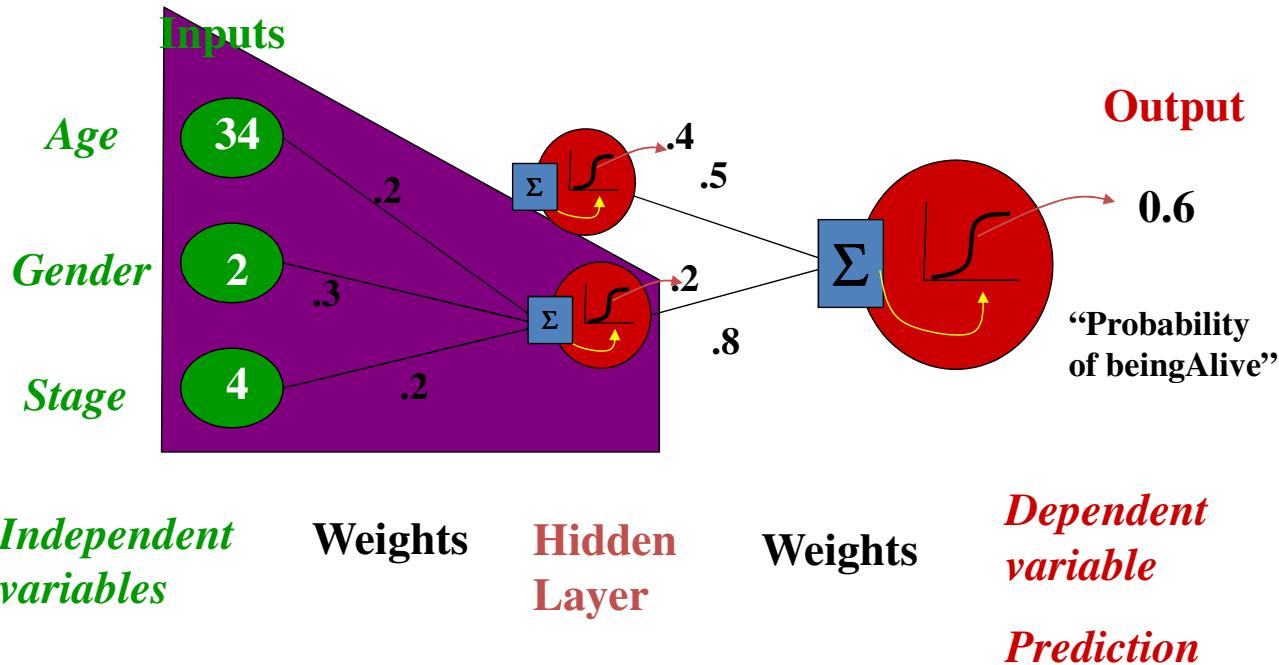
The last layer can be multi-class classification such as soft-max function

$$P(y^i = c|x^i) = \frac{\exp(\theta_c^\top x^i)}{\sum_{c'} \exp(\theta_{c'}^\top x^i)}$$

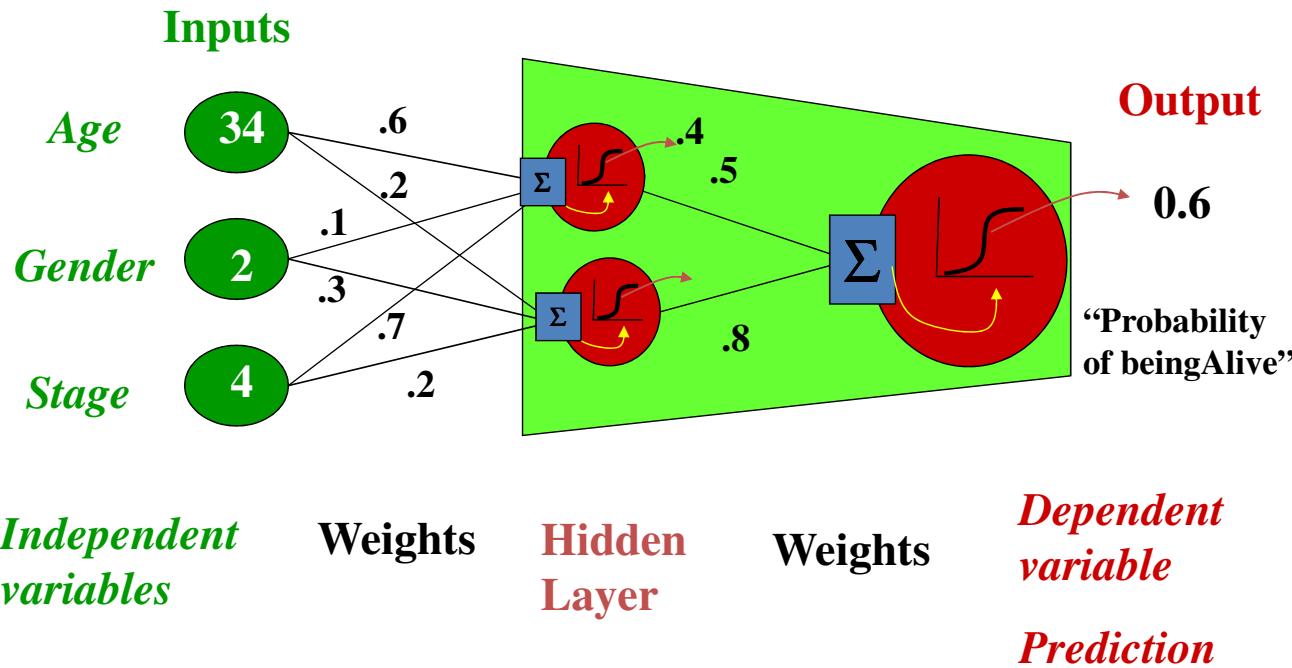
# “Combined logistic models”



# “Combined logistic models”

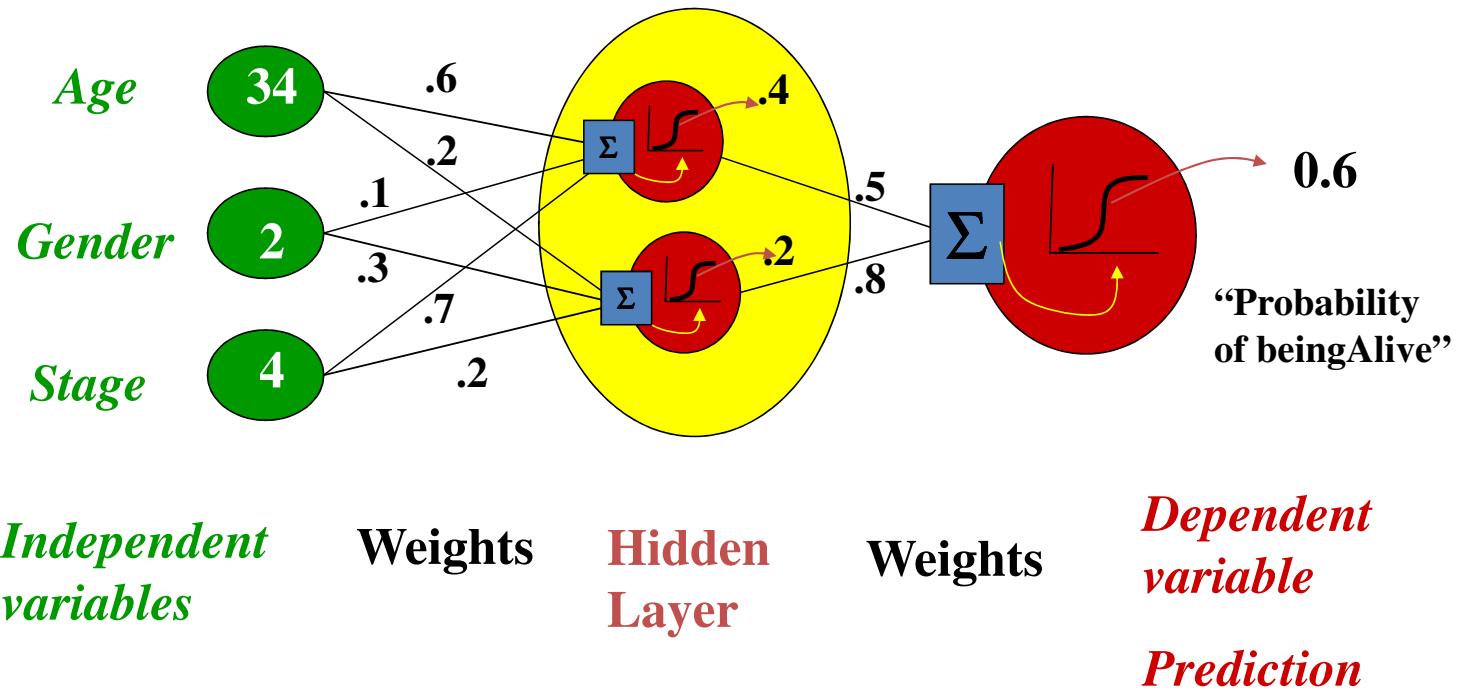


# “Combined logistic models”



# Overall structure

Turn weights and thresholds such that the neural networks will give desired prediction results given input.



# Learning parameters in logistic regression

- Find  $\theta$ , such that the conditional likelihood of the labels is maximized

$$\max_{\theta} l(w) := \log \prod_{i=1}^m P(y^i | x^i, w)$$

- Good news:  $l(w)$  is concave function of  $w$ , and there is a single global optimum.
- Bad new: no closed form solution (resort to numerical method)

# Training/fitting neural networks: Backpropagation

Find  $w, \alpha, \beta$ , such that the value of the output unit are close to the actual labels (0 or 1)

$$\min_{w, \alpha, \beta} l(w, \alpha, \beta) := \sum_i^n \left( y^i - \sigma(w^\top \mathbf{z}^i) \right)^2$$

where (use sigmoid function  $\sigma$ )

$$\begin{aligned} z_1^i &= \sigma(\alpha^\top x^i) \\ z_2^i &= \sigma(\beta^\top x^i) \end{aligned}$$

Non-convex objective function

Use gradient decent to find a local optimum

# The gradient of $l(w, \alpha, \beta)$

$$l(w, \alpha, \beta) := \sum_{i=1}^n \left( y^i - \sigma(w^\top \mathbf{z}^i) \right)^2$$

where  $z_1^i = \sigma(\alpha^\top x^i)$ ,  $z_2^i = \sigma(\beta^\top x^i)$

Let  $u^i = w^\top \mathbf{z}^i$

$$\frac{\partial \sigma(u)}{\partial u} = \sigma(u)(1 - \sigma(u))$$

Gradient

$$\frac{\partial l(w, \alpha, \beta)}{\partial w} = - \sum_i 2 \left( y^i - \sigma(u^i) \right) \sigma(u^i) \left( 1 - \sigma(u^i) \right) \mathbf{z}^i$$

$z^i$  is computed using  $\alpha$  and  $\beta$  from previous iteration

$$z_1^i = \sigma(\alpha^\top x^i), z_2^i = \sigma(\beta^\top x^i)$$

# The gradient with respect to $\alpha, \beta$

Use chain rule of derivatives (backpropagation)

The "error" at layer k depends on the error at the next layer k+1

Note  $z_1^i = \sigma(\alpha^\top x^i)$ ,  $z_2^i = \sigma(\beta^\top x^i)$

Let  $v^i = \alpha^\top x^i$

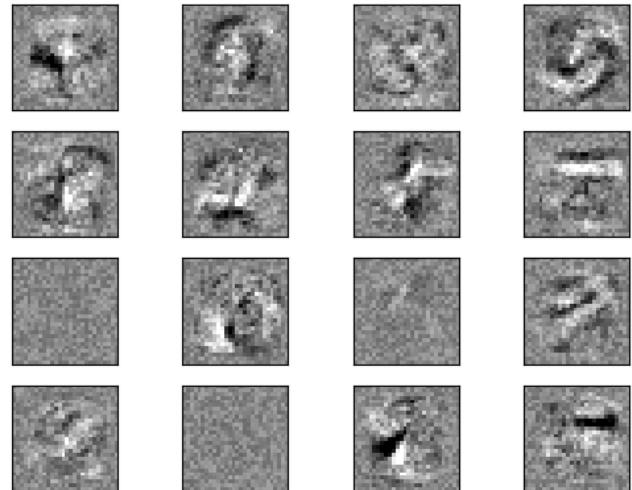
$$\frac{\partial \sigma(u)}{\partial u} = \sigma(u)(1 - \sigma(u))$$

Gradient

$$\begin{aligned}\frac{\partial l(w, \alpha, \beta)}{\partial \alpha} &= \frac{\partial l(w, \alpha, \beta)}{\partial z_1^i} \frac{\partial z_1^i}{\partial \alpha} \\ &= - \sum_i 2 \left( y^i - \sigma(u^i) \right) \sigma(u^i) \left( 1 - \sigma(u^i) \right) w_1 \sigma(v^i) \left( 1 - \sigma(v^i) \right) x^i\end{aligned}$$

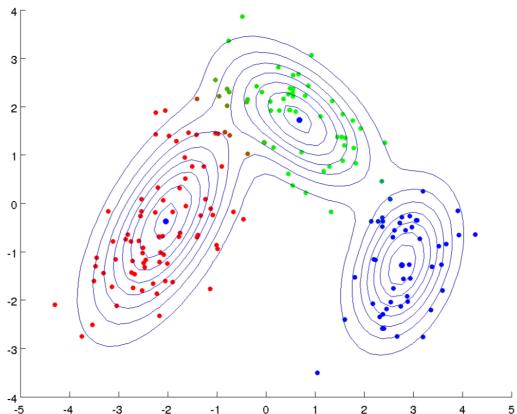
# Demo: hand-written digits

- The input data consists of 28x28 pixel handwritten digits, leading to 784 features in the dataset.
- One-hidden layer
- Therefore the first layer weight matrix have the shape (784, hidden\_layer\_size)
- We can therefore visualize weight for each hidden node as a 28x28 pixel image.



# Demo: Classification of wine data

- Data: chemical analysis of wines grown in the same region in Italy but derived from three different cultivars. 13 features, three types of wines.
- Use multi-class classification in the last layer



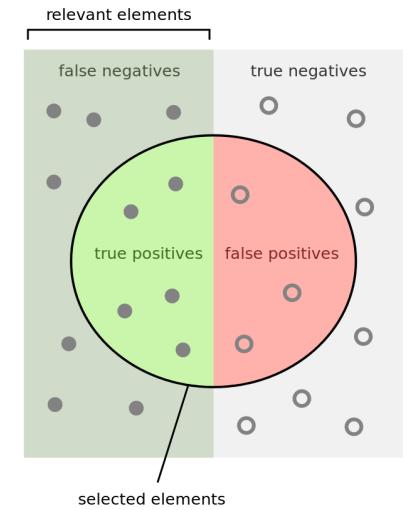
# Classification performance measures

- Mis-classification error:  
 $\# \text{ mis-classified samples} / \# \text{ total}$

- Confusion matrix  
a sample of 13 pictures, 8 of cats and 5 of dogs



		Actual class	
		Cat	Dog
Predicted class	Cat	5	2
	Dog	3	3



How many selected items are relevant?

$$\text{Precision} = \frac{\text{green}}{\text{red} + \text{green}}$$

How many relevant items are selected?

$$\text{Recall} = \frac{\text{green}}{\text{green} + \text{white}}$$

Source: wikipedia

# Classification performance measures

- Precision and recall

Precision =

$$\# \text{ retrieved relevant instances} / \# \text{ total retrieved}$$

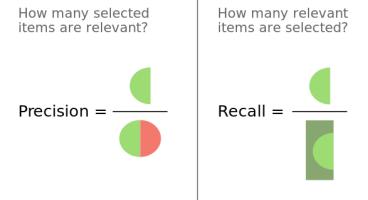
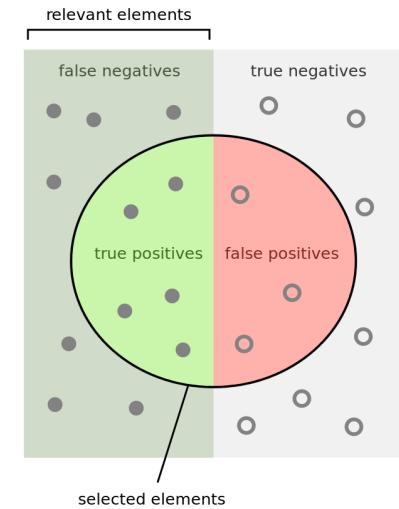
Recall =

$$\# \text{ retrieved relevant instances} / \# \text{ total relevant}$$

$$\text{F1 score} = \frac{2}{\text{precision}^{-1} + \text{recall}^{-1}} \in [0,1]$$

- Using the previous example

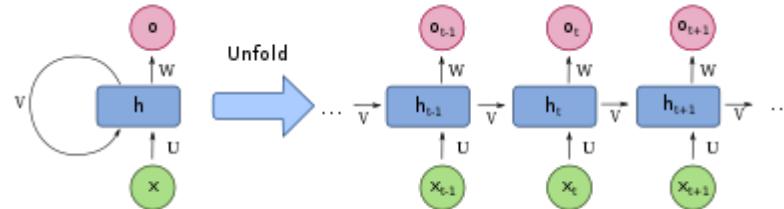
	Precision	Recall	F1-score
Cat	0.714	0.625	0.666
Dog	0.500	0.600	0.545



Source: wikipedia

# Some common neural network structures

- Full connected (all nodes in the  $k$ -th layer are connected to all nodes in the  $(k+1)$ -th layer (not practical))
- Convolutional neural networks (nodes in one layer are only “locally” connected to nodes in the next layer), two types of layers: convolution and max pooling
- Recursive neural network (RNN): good for sequential / time series data; neural network structure “repeats” itself; attention mechanism, transformer

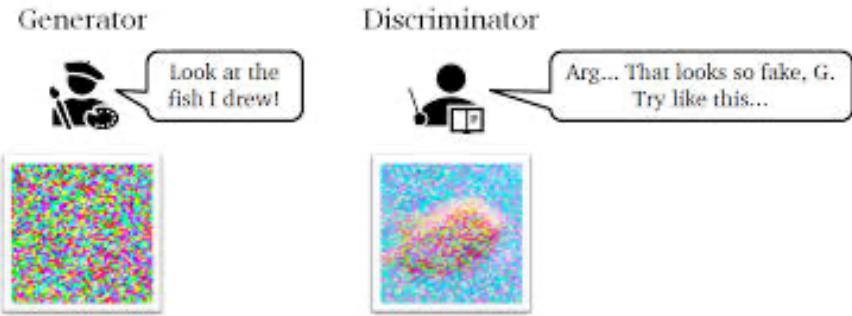


- Residual neural networks: jump over some layers

# GAN

- Generative adversarial neural networks (GAN) (2014) – generative model and better discriminator

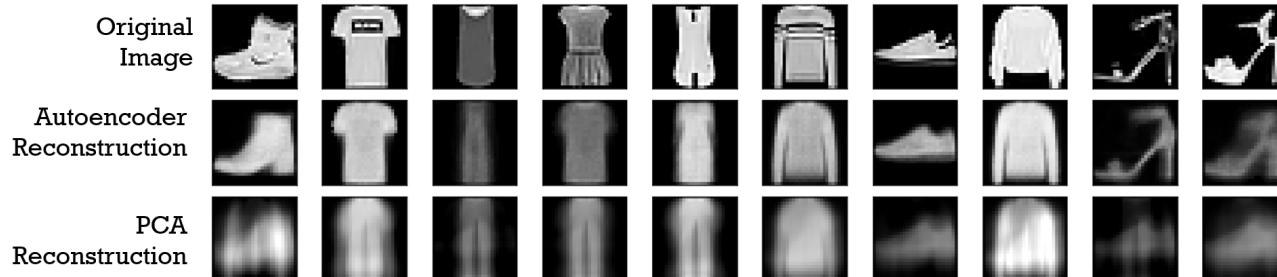
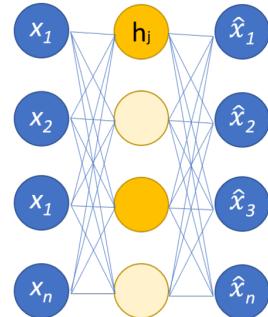
$$\min_G \max_D V(D, G) = \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}(\mathbf{x})} [\log D(\mathbf{x})] + \mathbb{E}_{\mathbf{z} \sim p_{\mathbf{z}}(\mathbf{z})} [\log(1 - D(G(\mathbf{z})))].$$



An image generated by a [StyleGAN](#) that looks deceptively like a photograph of a real person. This image was generated by a StyleGAN based on an analysis of portraits.

# Auto-encoder

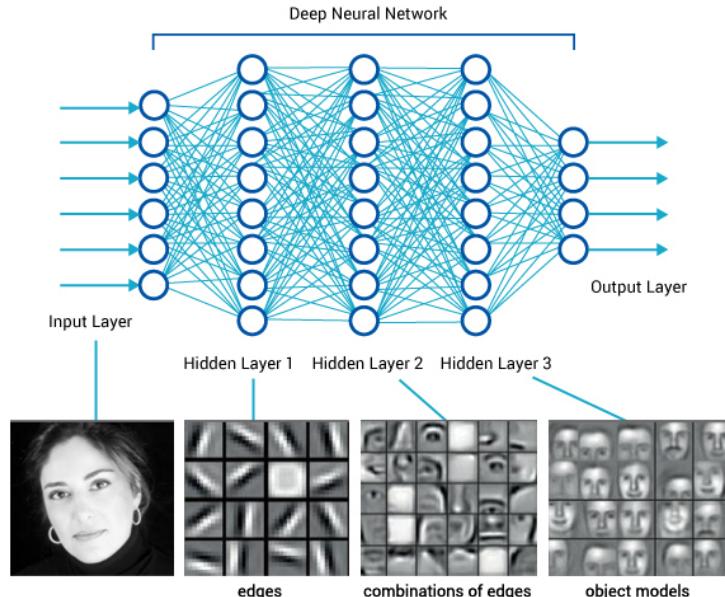
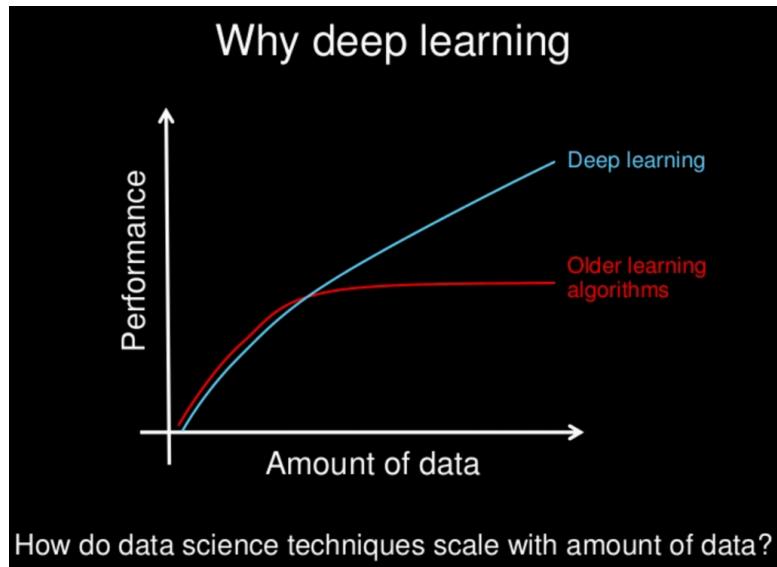
- Suppose we only have unlabeled training samples
- Autoencoder is a neural networks for unsupervised learning: setting the target value to input
- Network learn to represent data itself
- Variational autoencoder (VAE) gives generative model
- If linear activations are used, or only a single sigmoid hidden layer, then the optimal solution to an autoencoder is strongly related to PCA



Reconstruction of 28x28pixel images by an Autoencoder with a code size of two (two-units hidden layer) and the reconstruction from the first two Principal Components of PCA. Images come from the [Fashion MNIST dataset](#).<sup>[31]</sup>

# Why deep learning

- Deep learning (deep neural networks) revolutionary



# ImageNet

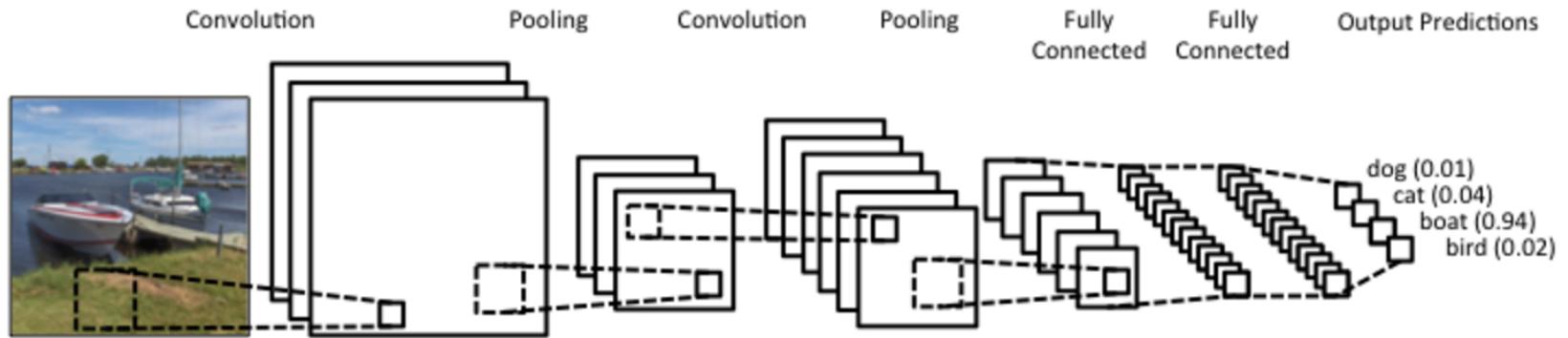
Image classification with 1.3M color images and 1000 classes

Need large scale nonparametric methods



<https://papers.nips.cc/paper/4824-imagenet-classification-with-deep-convolutional-neural-networks.pdf>

# Convolutional Neural Network (CNN)



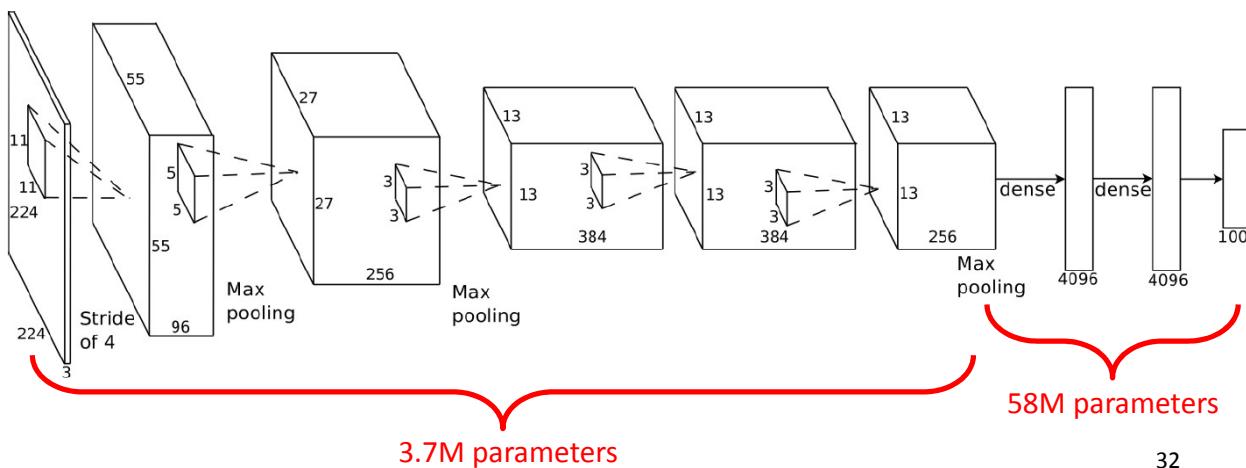
# Convolutional Neural Network (CNN)

8 layer convolution neural network [Krizhevsky12], Achieved state-of-the-art result (beating the second place by 10%)

First 5 layers: convolution + max pooling

**Next 2 layers:** fully connected nonlinear neurons

**Last layer:** multiclass logistic regression

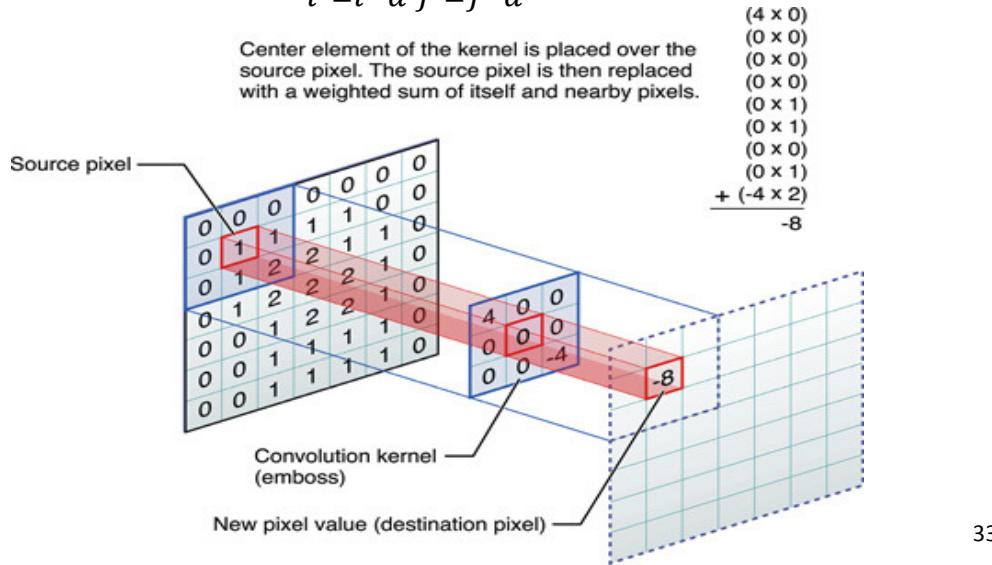


# Convolution layer

- Source image  $I$ , destination image  $O$ , convolving image with a kernel  $K$  of size  $2d+1$

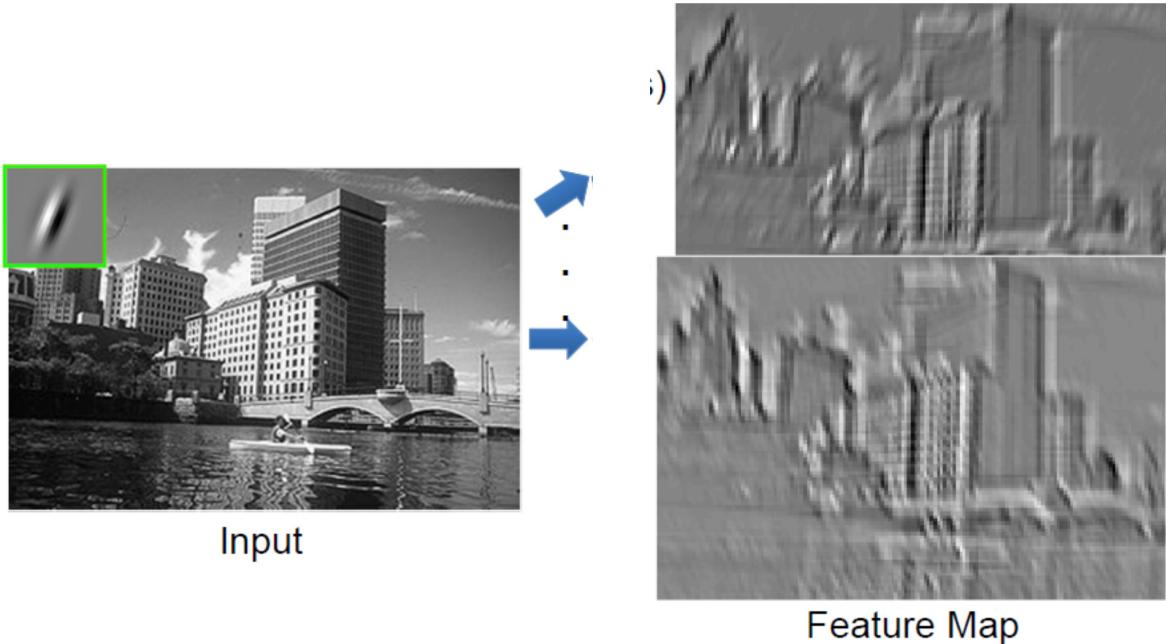
$$O(i,j) = \sum_{i'=i-d}^{i+d} \sum_{j'=j-d}^{j+d} I(i',j')K(i',j')$$

Center element of the kernel is placed over the source pixel. The source pixel is then replaced with a weighted sum of itself and nearby pixels.



# Convolutional feature

Convolution extract local image features

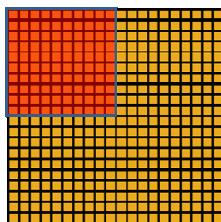


<https://i.imgur.com/TqTJCrz.png>

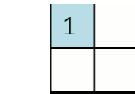
# Convolutional Neural Network

Spatial pooling

- Sum or max (usually max)
- Overlapping or non-overlapping regions
- Invariant to small translations



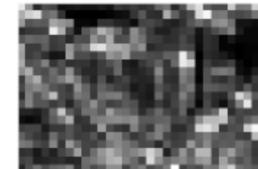
Convolved  
feature



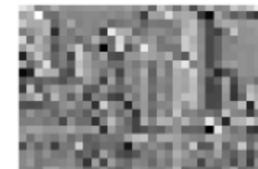
Pooled  
feature



Max

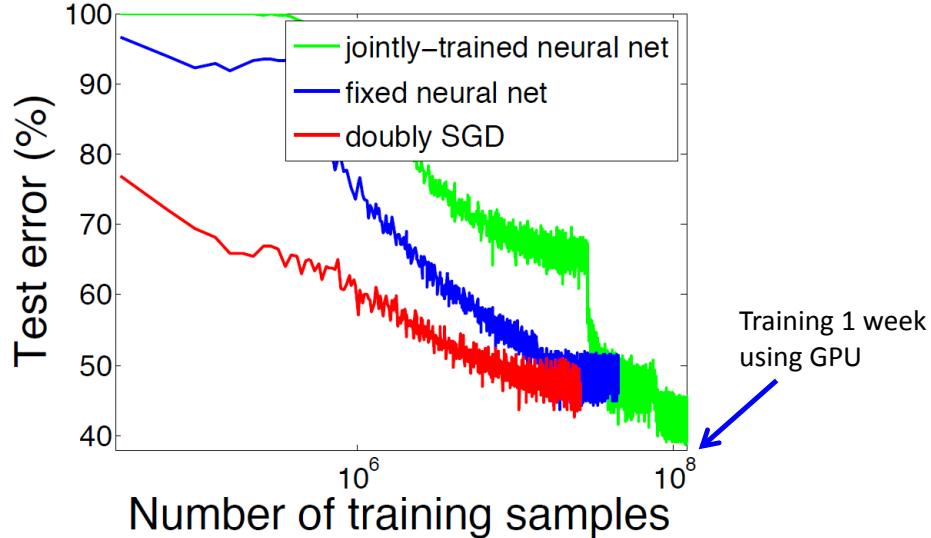


Sum



# ImageNet

Millions of training points, 1000 class classification task



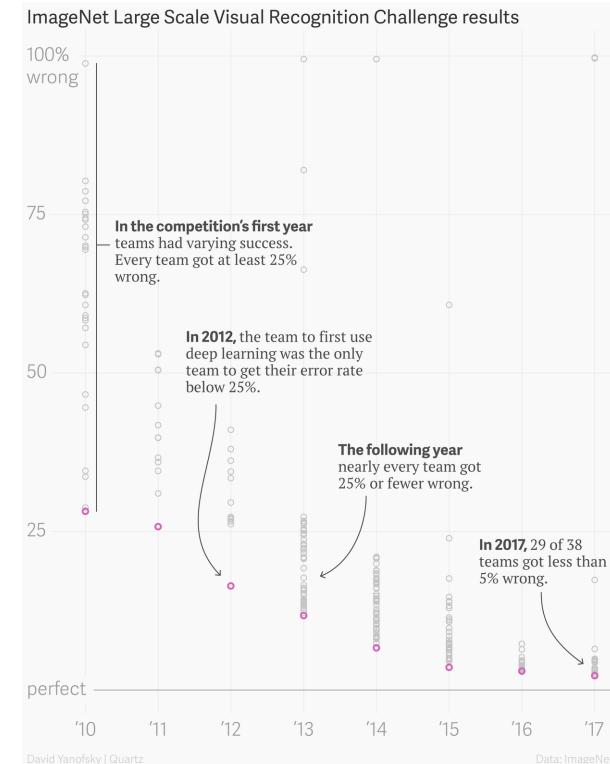
[https://github.com/zixu1986/Doubly\\_Stochastic\\_Gradients](https://github.com/zixu1986/Doubly_Stochastic_Gradients)

3^n

# ImageNet: Demonstrate power of deep learning

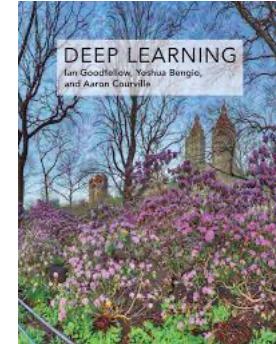
On 30 September 2012, a convolutional neural network (CNN) called AlexNet achieved a top-5 error of **15.3%** in the ImageNet 2012 Challenge, more than 10.8 percentage points lower than that of the runner up. This was made feasible due to the use of Graphics processing units (GPUs) during training, an essential ingredient of the deep learning revolution.

According to The Economist, "Suddenly people started to pay attention, not just within the AI community but across the technology industry as a whole."



# Ending words

- Deep learning is the “hot-topic” in “hot-topics”
- Rapidly evolving
- Many research and engineering efforts, hugely successful in various domain: image, speech etc.
- Many remaining questions: interpretability (black-box), reliability (e.g., healthcare, manufacturing)
- Additional references
- NeurIPS 2015 tutorial on deep learning:  
<https://aiatadams.files.wordpress.com/2016/02/bengio-lecun-20151207-deep-learning-tutorial-nips.pdf>
- NeurIPS 2016 tutorial on GAN: <https://media.nips.cc/Conferences/2016/Slides/6202-Slides.pdf>
- Deep learning book: <https://www.deeplearningbook.org>



2015

