# Deep Learning on Azure

## Deep Learning Crash Course

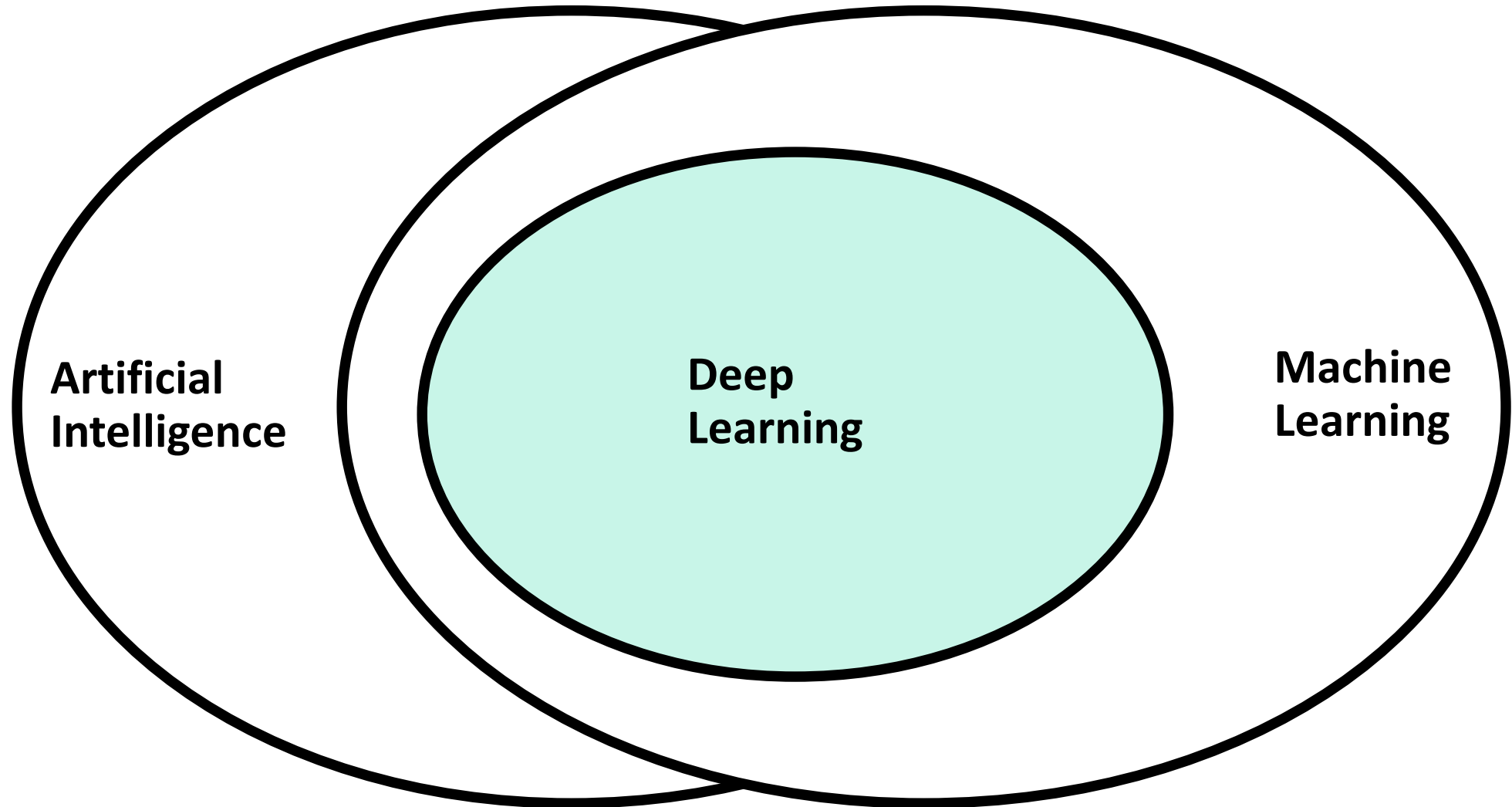Jarek Kazmierczak
MTC Silicon Valley

# Agenda

## Morning session:

- Introductions
- Deep Learning Crash Course
- Overview of Microsoft Cognitive Toolkit
- Hands-on labs:
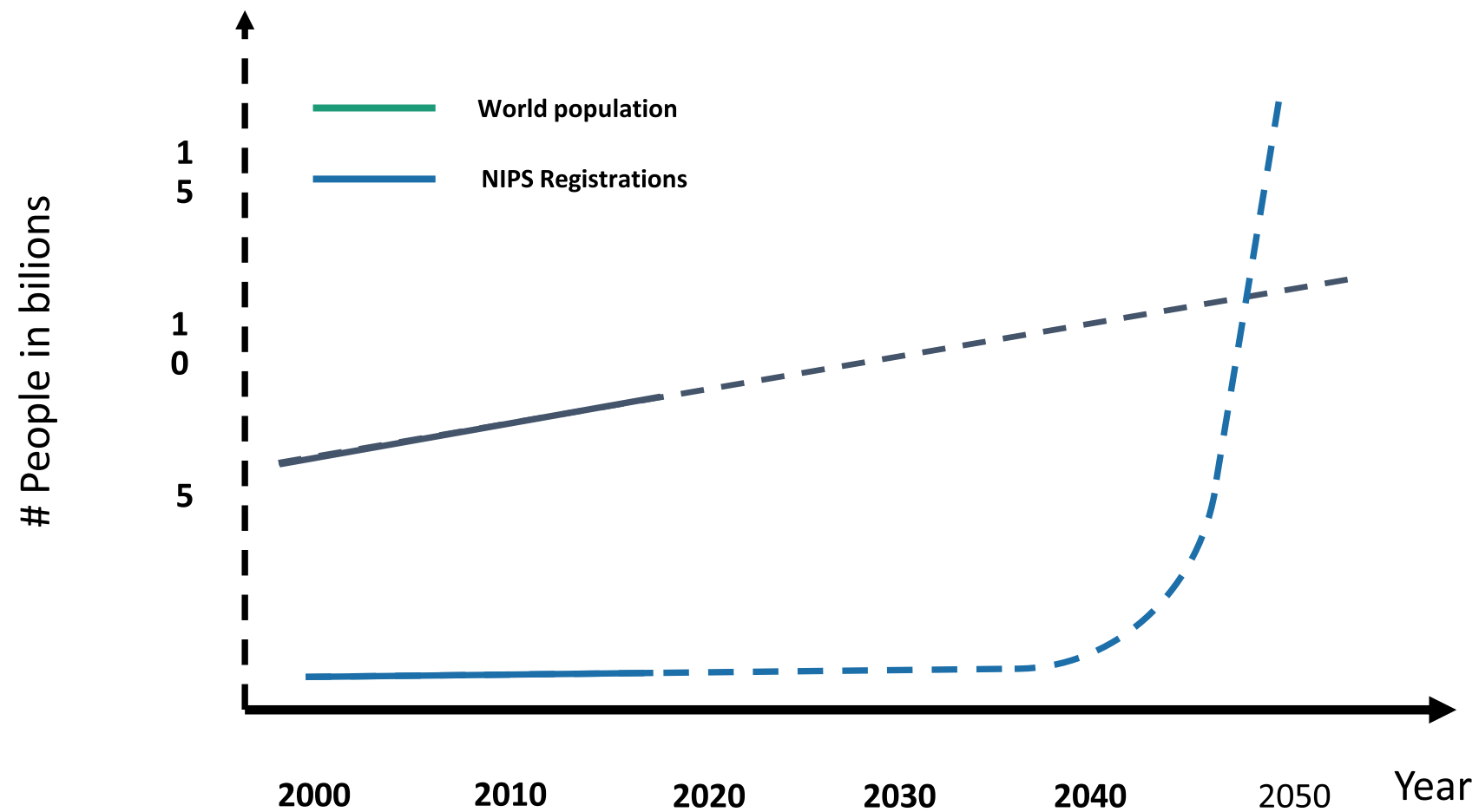  - Lab 1 - Multiclass Logistic Regression – Computer Vision

## Afternoon session:

- Hands-on labs:
  - Lab 2 - Fully Connected Neural Network – Computer Vision
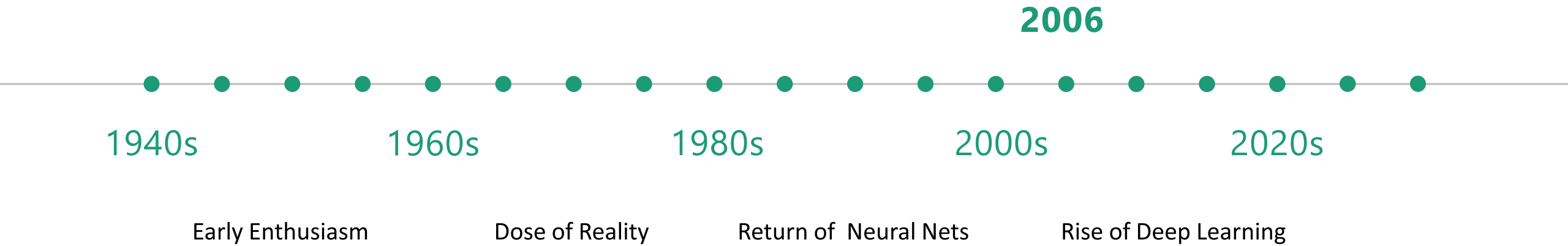  - Lab 3 – LSTM – Time series with IoT data

# Artificial Intelligence, Machine Learning and Deep Learning

# Have we been there before?

**2006**

1940s     1960s     1980s     2000s     2020s

Early Enthusiasm     Dose of Reality     Return of  Neural Nets     Rise of Deep Learning

# The Deep Learning Triumv[ei]rate

*LeCun*: "You have to realize that deep learning .... is really a conspiracy between Geoff Hinton and myself and Yoshua Bengio"



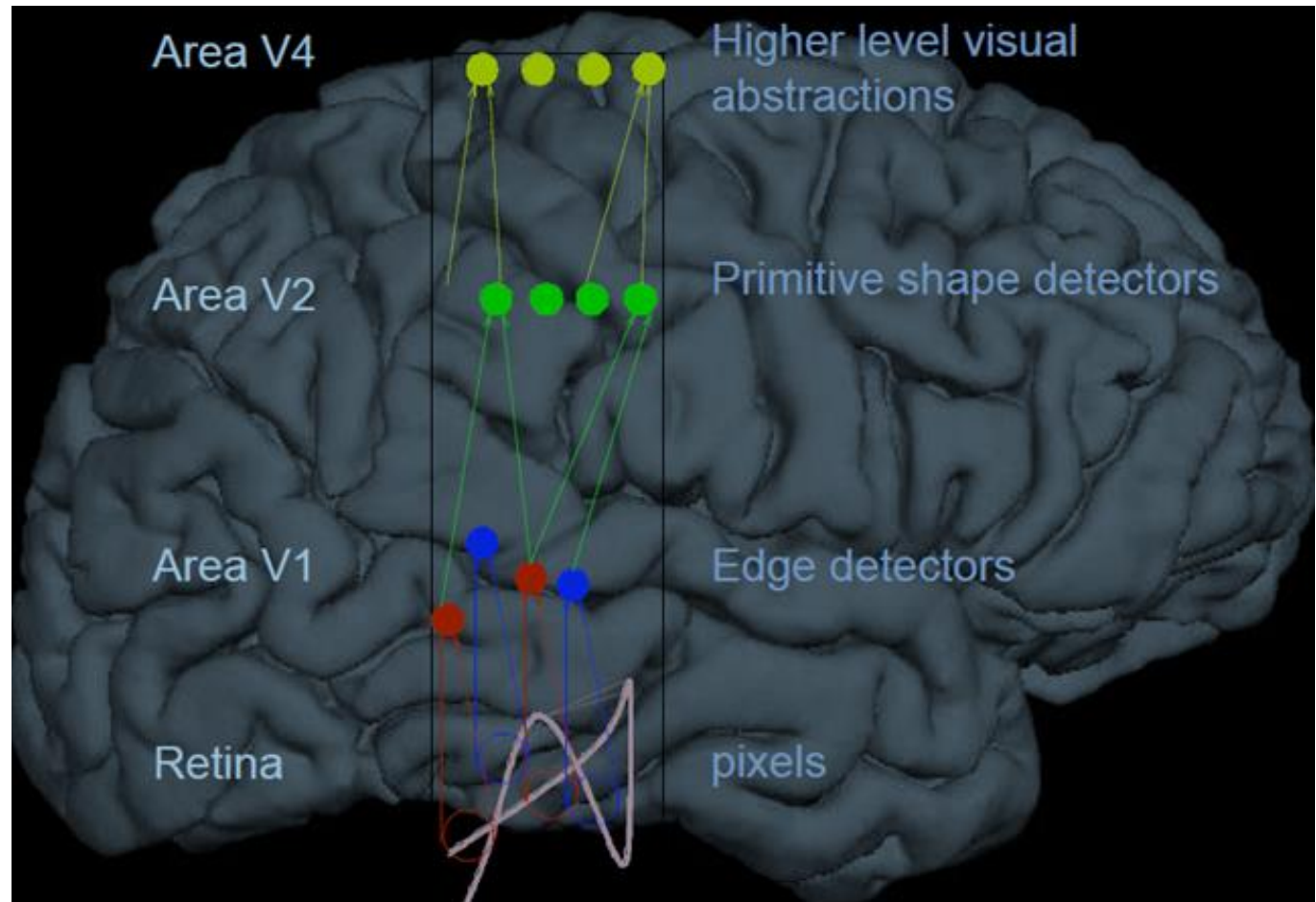Geoff Hinton          Yann LeCun          Yoshua Bengio

Latin  :  Trium  -  {ver,vir}  -  ate

English  :  Three  -  {truth, men}  -  official

# Why go deep?

- Deep learning algorithms attempt multiple levels of representation of increasing complexity/abstraction

- Brains have a deep architecture

- Deep Learning has been successful in tasks that have been a challenge for "traditional ML"
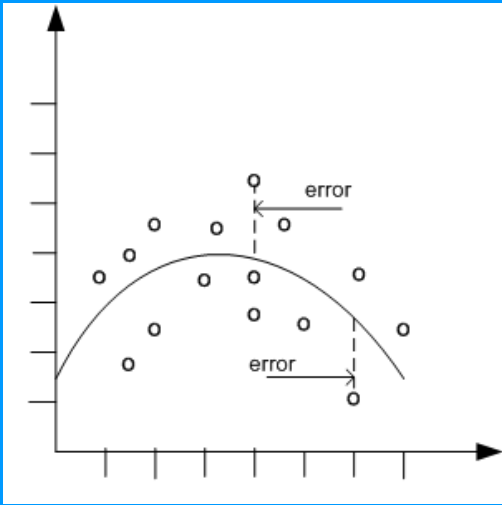
# What is Machine Learning?

Tom Mitchell (1998). Well-posed Learning Problem: A computer program is said to *learn* from experience E with respect to some task T and some performance measure P, if its performance on T, as measured by P, improves with experience E.
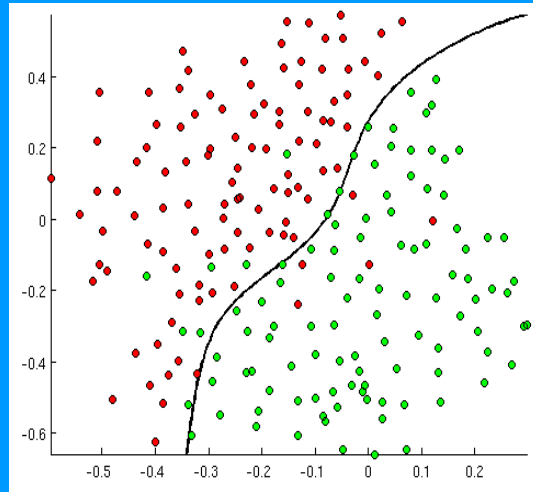
# The Task, *T*



**Predict a number**
Regression

**Predict a Class**
Classification

**Find groups/patterns**
Clustering

**Find unusual items**
Anomaly Detection

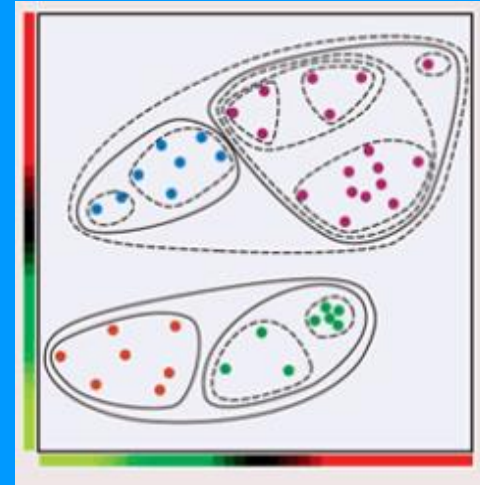What will referral fee revenue be in Q3?

What will assets under management be in 2017?

What is propensity of customer to purchase a variable annuity?

Probability the customer will churn to competitor?

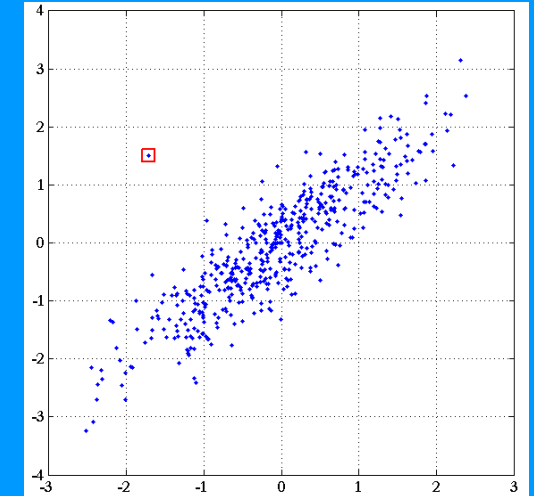Find similar customers for a new "wealth builder" segment.

What is the profile of customers with 3+ products?

Identify fraudulent expense report filings.

# Traditional ML Vs DL

Traditional ML requires manual feature extraction/engineering

Deep learning can automatically learn features in data

Feature extraction for unstructured data is very difficult

Deep learning is largely a "black box" technique, updating learned weights at each layer
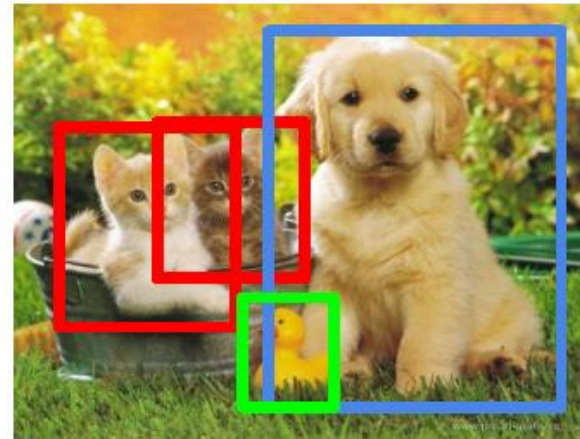
# Computer vision tasks



| Classification | Classification + Localization | Object Detection | Instance Segmentation |

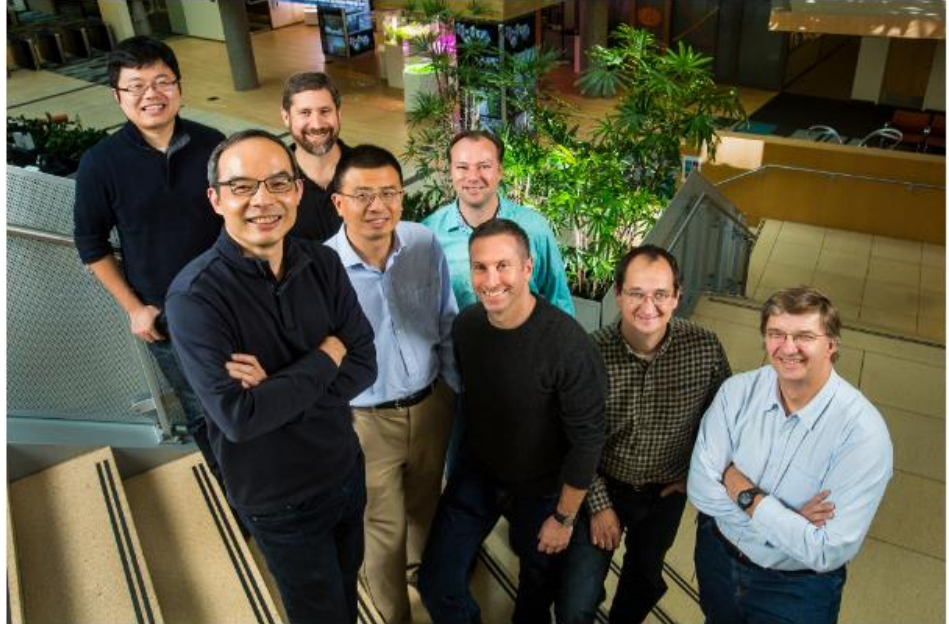CAT · CAT · CAT, DOG, DUCK · CAT, DOG, DUCK

Single object · Multiple objects

# Speech recognition tasks

- Microsoft 2016 research system for conversational speech recognition
- 5.9% word-error rate
- enabled by CNTK's multi-server scalability

[W. Xiong, J. Droppo, X. Huang, F. Seide, M. Seltzer, A. Stolcke, D. Yu, G. Zweig: "Achieving Human Parity in Conversational Speech Recognition," https://arxiv.org/abs/1610.05256]



Historic Achievement: Microsoft researchers reach human parity in conversational speech recognition

Microsoft researchers from the Speech & Dialog research group include, from back left, Wayne Xiong, Geoffrey Zweig, Xuedong Huang, Dong Yu, Frank Seide, Mike Seltzer, Jasha Droppo and Andreas Stolcke. (Photo by Dan DeLong)
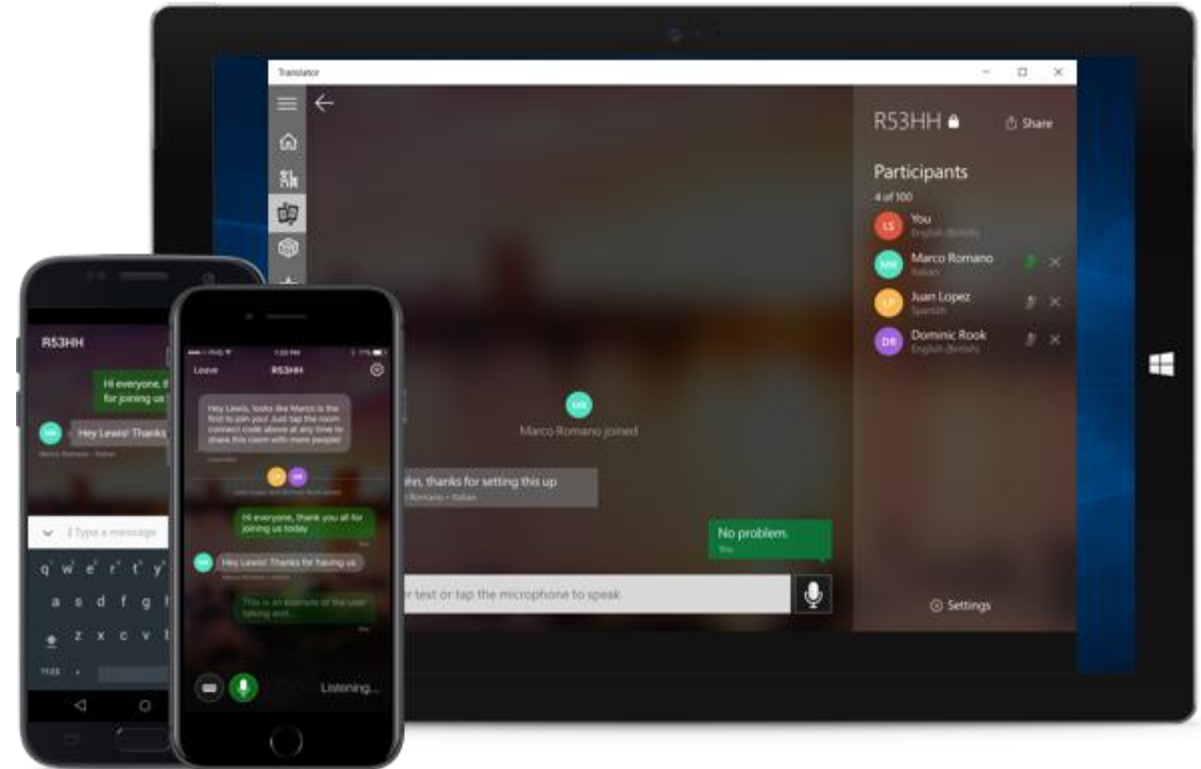
Posted October 18, 2016

By **Allison Linn**

Microsoft has made a major breakthrough in speech recognition, creating a technology that recognizes the words in a conversation as well as a person does.

# Language translation and understanding tasks

- Microsoft Translator live is an in-person, multi-device translation service for two or more participants, speech or text.

- Start a conversation, share the code and break the language barrier.

- 9 speech languages and 60 text languages.

- Apps on iOS, Android, Windows UWP and web. API to manage the conversation.

Training examples

Training labels

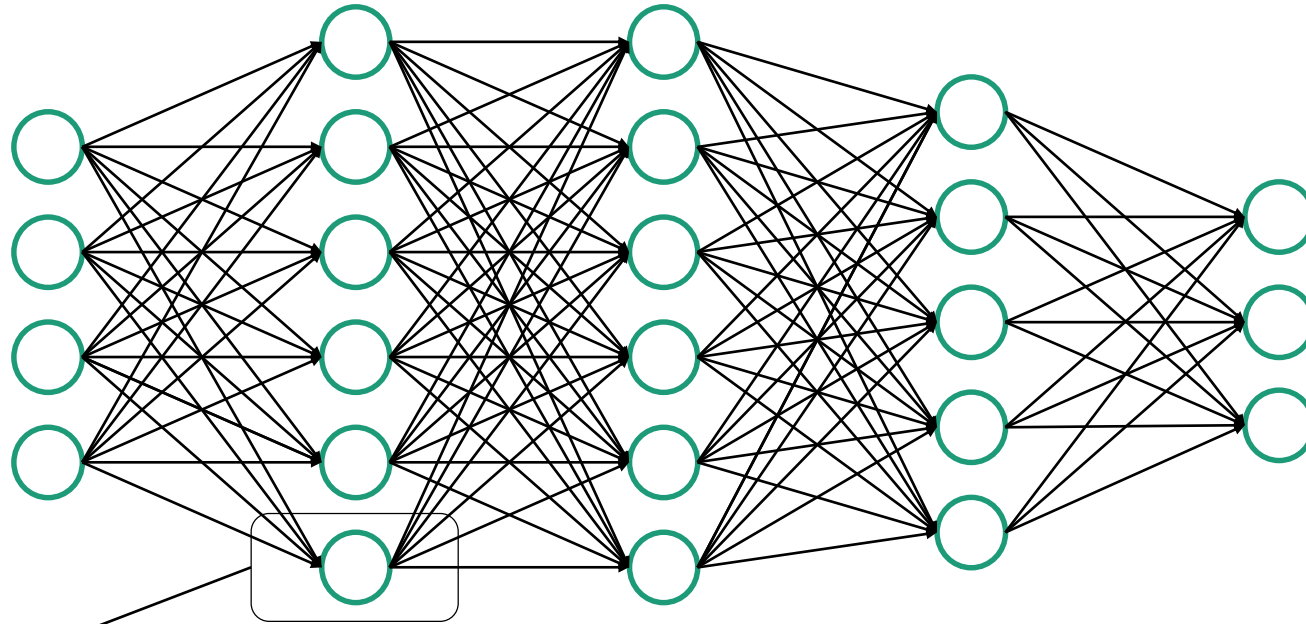| 1 | 1 | 5 | 4 | 3 |
| 7 | 5 | 3 | 5 | 3 |
| 5 | 5 | 9 | 0 | 6 |
| 3 | 5 | 2 | 0 | 0 |

Accurate digit classifier

2

Machine learning system

# Simplified learning model

# Artificial Neural Networks - ANNs

Input layer    Hidden layer    Hidden layer    Hidden layer    Output layer

$$a_j = \sigma\left(\sum_{i=1}^{n} w_i x_i + b_j\right)$$

**Activation functions**

0.5

0.0

**Sigmoid**

0.0

0.0

**ReLU**

# Neurons can connect in various ways …

"Dense"

"Sparse"

"Feedback loops"

**F**ully **C**onnected **N**eural **N**etworks

**C**onvolutional **N**eural **N**etworks

**R**ecurrent **N**eural **N**etworks

# … and can be arranged in layers



Convolutions

Fully connected

# Revolution of depth

## ResNet: 152 layers, and 1001 layers later on

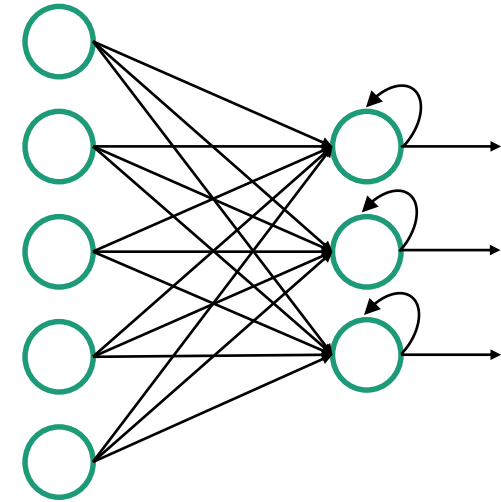*MSRA's ResNet won the 1st places in ImageNet classification, ImageNet detection, ImageNet localization, COCO detection, and COCO segmentation in ILSVRC &COCO competitions 2015*



ImageNet classifcation top-5 error (%)

| | | | | | |
|---|---|---|---|---|---|
| 28.2 | 25.8 | 16.4 | 11.7 | 6.7 | 3.5 |
| 2010 | 2011 | 2012 | 2013 | 2014 | 2015 |
| NEC | Xerox | AlexNet (8 layers) | Clarifi (8 layers) | GoogleNet (22 layers) | ResNet (152 layers) |

Human Performance



7x7 conv, 64, /2, pool/2
1x1 conv, 64
3x3 conv, 64
1x1 conv, 256
1x1 conv, 64
3x3 conv, 64
1x1 conv, 256
1x1 conv, 64
3x3 conv, 64
1x1 conv, 256
1x2 conv, 128, /2
3x3 conv, 128
1x1 conv, 512
1x1 conv, 128
3x3 conv, 128
1x1 conv, 512
1x1 conv, 128
3x3 conv, 128
1x1 conv, 512
1x1 conv, 128
3x3 conv, 128
1x1 conv, 512
1x1 conv, 128
3x3 conv, 128
1x1 conv, 512
1x1 conv, 128
3x3 conv, 128
1x1 conv, 512

# ANN Learning



Prepare **Training Data**
$(\mathbb{x}_1, y_1), \ldots, (\mathbb{x}_N, y_N)$

Define a **Network Model**

Define **Loss Function**

and

Run **Optimization Procedure**

Trained Network

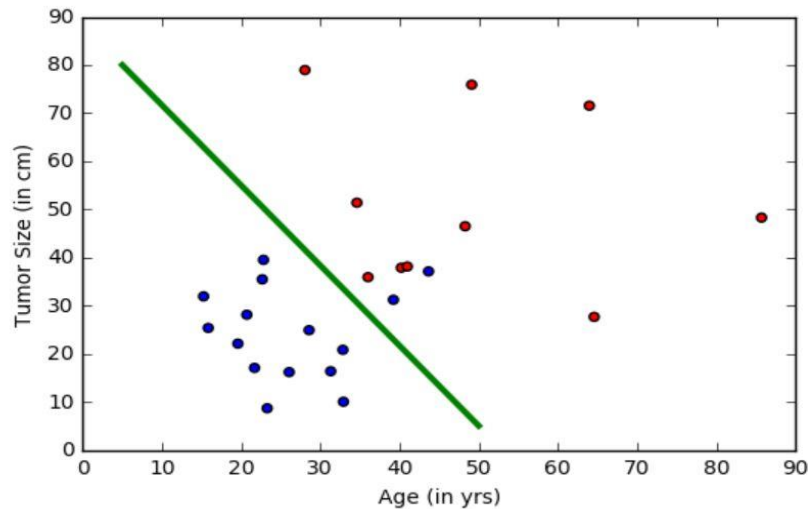# Two class logistic regression – aka. single neuron ANN



Training Data: $\quad D = \{(x_1, y_1), (x_2, y_2), \dots, (x_N, y_N)\}$

Hypothesis Space: $\quad h(x) = \dfrac{1}{1 + e^{-(w^T x + b)}}$

Loss Function: $\quad E(w) = -\displaystyle\sum_{j=1}^{N} \left(y^{(j)} log\left(h(x^{(j)})\right) + \left(1 - y^{(j)}\right) log\left(1 - h(x^{(j)})\right)\right)$

Optimization Procedure: $\quad \textbf{\textit{Gradient Descent}}$
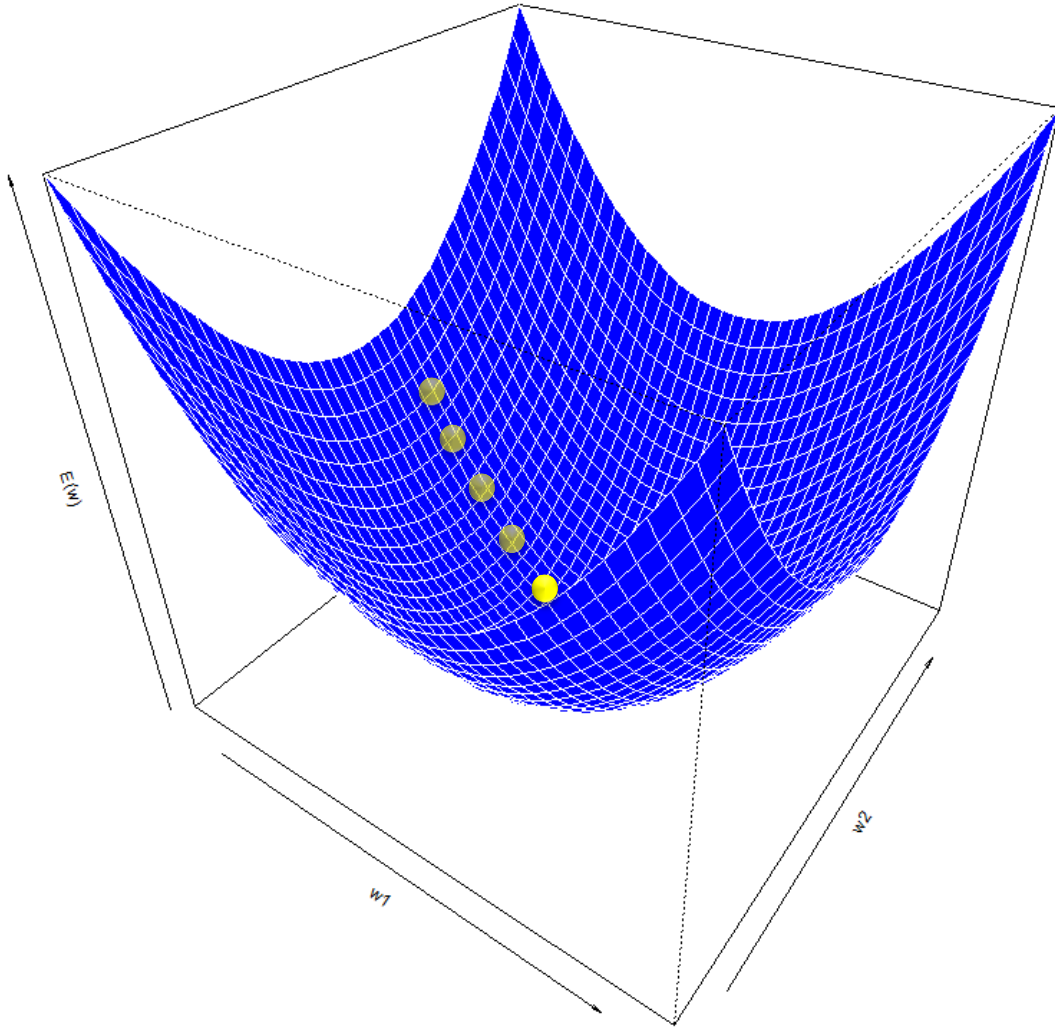
# Gradient descent (Augustin Cauchy, 1847)

# Gradient descent in two dimensional parameter space



$$E(w) = -\sum_{j=1}^{N}(y^{(j)}log(h(x^{(j)}) + (1 - y^{(j)})log(1 - h(x^{(j)})))$$

$$\nabla J(w) \equiv \left[\frac{\partial J}{\partial w_1}, \frac{\partial J}{\partial w_2}\right]$$

$$w = w - \eta \nabla J(w)$$

# Stochastic Gradient Descent and Backpropagation

**Loss function**

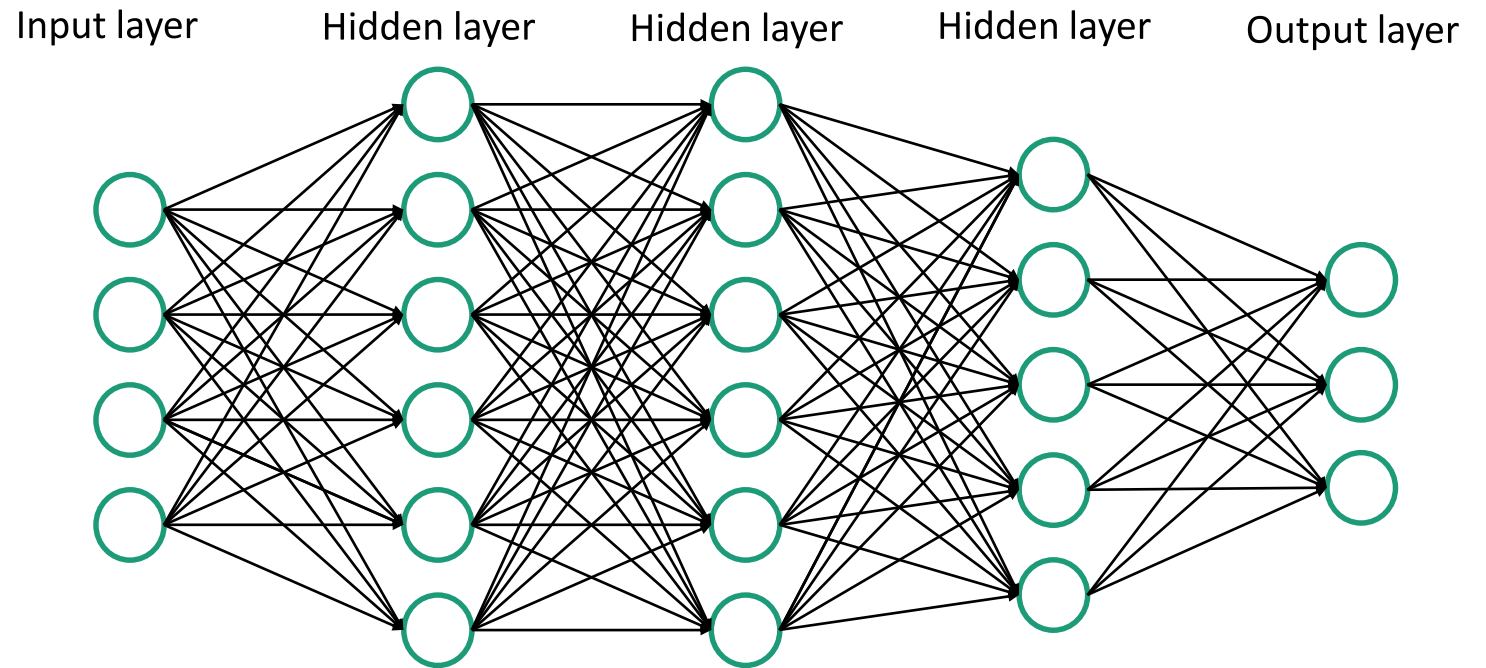$$\min_{\boldsymbol{w}} \sum_{i=1}^{N} f(x_i, y_i; w)$$

**Stochastic Gradient Descent (SGD)**

$$g(w_t) = \nabla f(x_i, y_i; w_t)$$

$$w_{t+1} = w_t - \eta_t g(w_t)$$



Input layer    Hidden layer    Hidden layer    Hidden layer    Output layer

# Stochastic Gradient Descent

Initialize learning rate $\boldsymbol{\eta}$

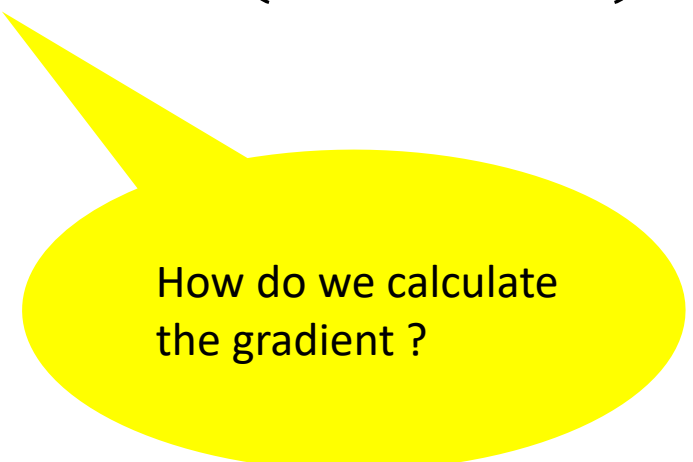Initialize parameter vector $\boldsymbol{w}$

**while** stopping criterion not met **do**

Sample a minibatch of $\boldsymbol{m}$ examples from the training set

Compute gradient estimates: $\widehat{\boldsymbol{g}} \leftarrow \frac{1}{\boldsymbol{m}} \boldsymbol{\nabla_w} \sum_{i=1}^{m} \boldsymbol{E}\left(\boldsymbol{w}, \boldsymbol{x}^{(i)}, \boldsymbol{y}^{(i)}\right)$

Apply update: $\boldsymbol{w} \leftarrow \boldsymbol{w} - \eta\widehat{\boldsymbol{g}}$

**end**

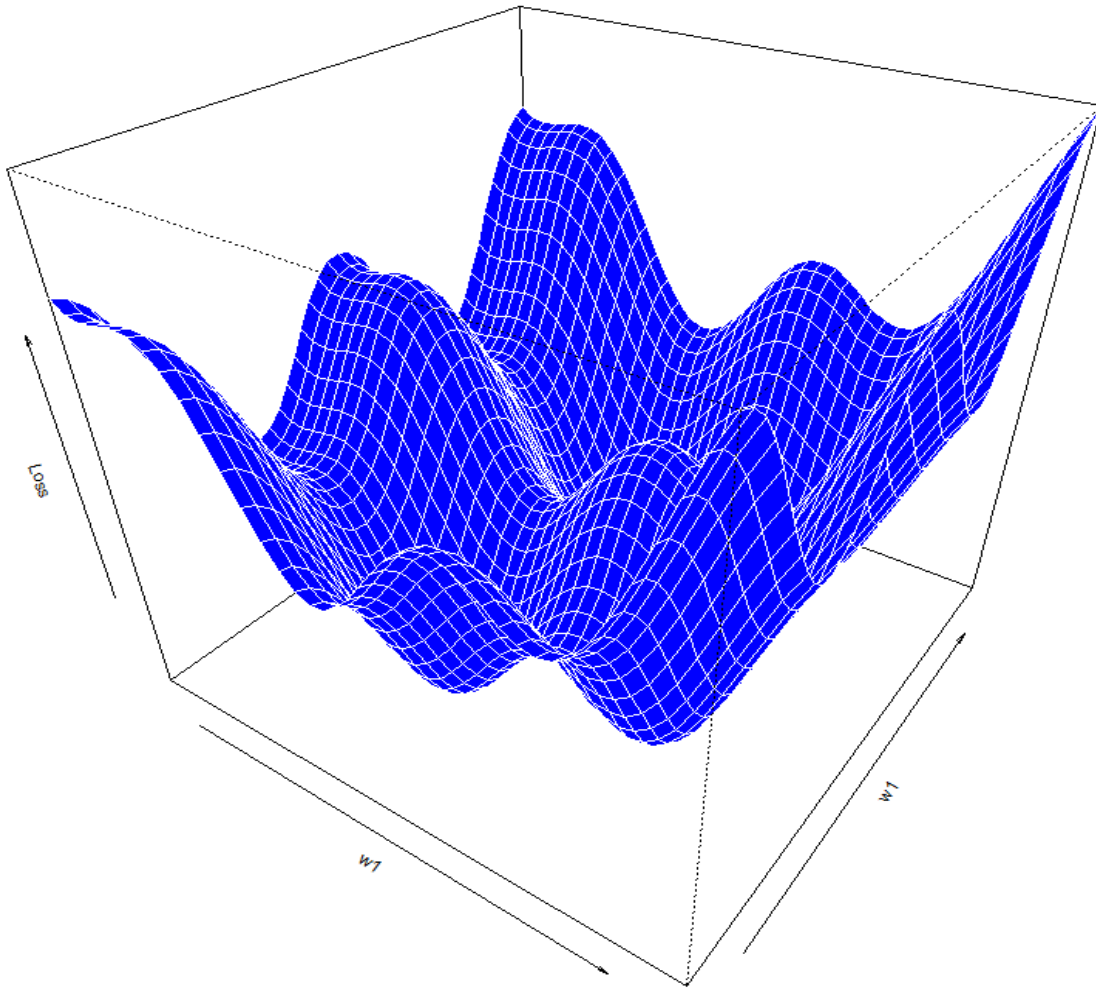How do we calculate the gradient ?

# The back propagation algorithm

1. Input $x$:
   Set the activations for the input layer $l = 1$

2. Feedforward:
   For each layer $l = 2,3, \dots, L$ compute:
   $$z^l = w^l a^{l-1} + b^l \text{ and } a^l = \sigma(z^l)$$

3. Output error $\delta^L$:
   Compute the vector:
   $$\delta^L = \nabla_a E \odot \sigma'(z^L)$$

4. Backpropagate the error:
   For each layer $l = L - 1, L - 2, \dots, 2$
   $$\delta^l = \left( (w^{l+1})^T \delta^{l+1} \right) \odot \sigma'(z^l)$$

5. Calculate gradient:
   $$\frac{\partial E}{w_{jk}^l} = a_k^{l-1} \delta_j^l \text{ and } \frac{\partial E}{b_j^l} = \delta_j^l$$

# Loss functions and Deep Learning

- ML literature and toolkits use different names for the same concept:
  - Objective function, loss function, cost function, error function, **criterion**
- In Deep Learning we often optimize a criterion which is not the same as the performance measure *P* we care about
  - Since optimizing the direct performance measure may be intractable,
  - One optimizes a **surrogate loss function**
- An example, common loss function for multinomial classification is **cross-entropy**

$$E(w) = -\frac{1}{N}\sum_{i=1}^{N}\sum_{j=1}^{C} y_j \, log(p_j)$$

# Deep Learning reality - non-convex loss functions



Imagine that in 1,000,000 dimensional space

# Deep Learning optimization procedures

- ## Stochastic Gradient Descent – SDG

- SDG with momentum

- AdaGrid

- RMSProp

- Adam
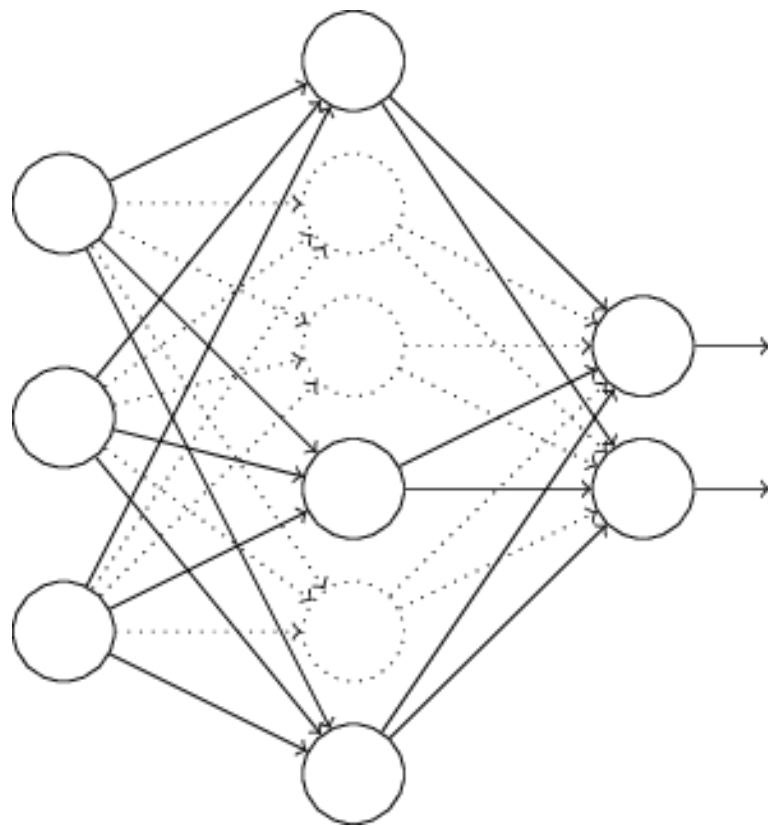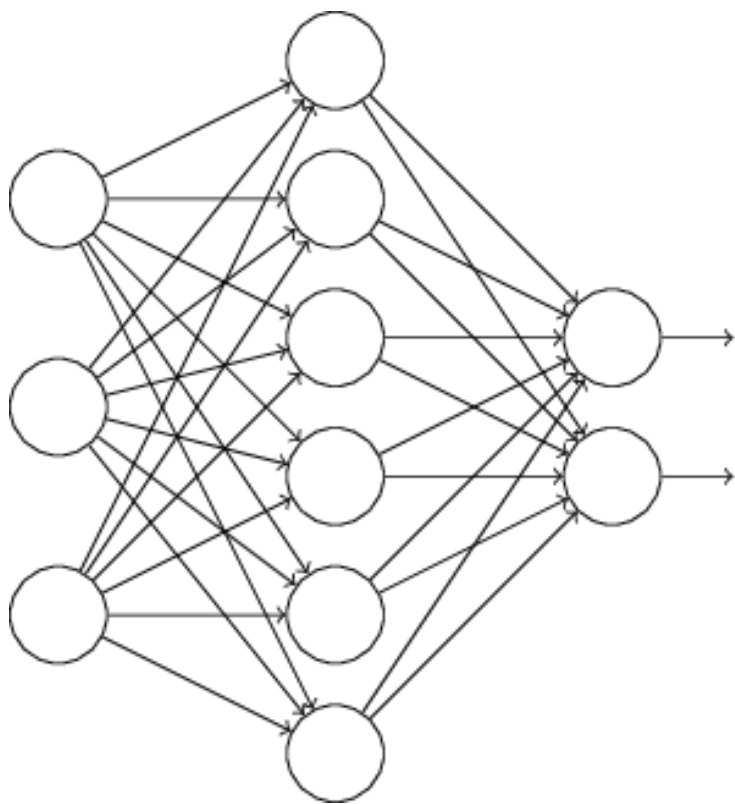
- Newton's Method

- Conjugate Gradients

- BFGS

# Preventing overfitting in Neural Networks

- Regularization
  - L2 regularization $E(w) = E(w)_0 + \frac{\lambda}{2n}\sum_w w^2$
  - L1 regularization $E(w) = E(w)_0 + \frac{\lambda}{n}\sum_w |w|$
- Early stopping
- Data augmentation



- Dropout

# Dropout

# Descending the rugged terrain of a multidimensional world can be a scary adventure

… You will be disoriented by countless saddle points

… You may get lost in vast, sparse, plateaus

… You risk falling down deep cliffs

… You could be injured by exploding gradients

# Sources and acknowledgements

- Frank Seide (2017). Training Deep Models Like Microsoft Product Groups.
- Christopher Bishop (2006). Pattern Recognition and Machine Learning.
- Ian Goodfellow, Yoshua Bengio, Aaron Courville (2016). Deep Learning.
- Michael Nielsen (2017). Neural Networks and Deep Learning.
- Yaser Abu-Mostafa et al. (2012) Learning From Data.
- Trevor Hastie et al. (2008). The Elements of Statistical Learning.
- Tom M. Mitchell (1997). Machine Learning.