# Deep Learning

## COMP 527
## Danushka Bollegala
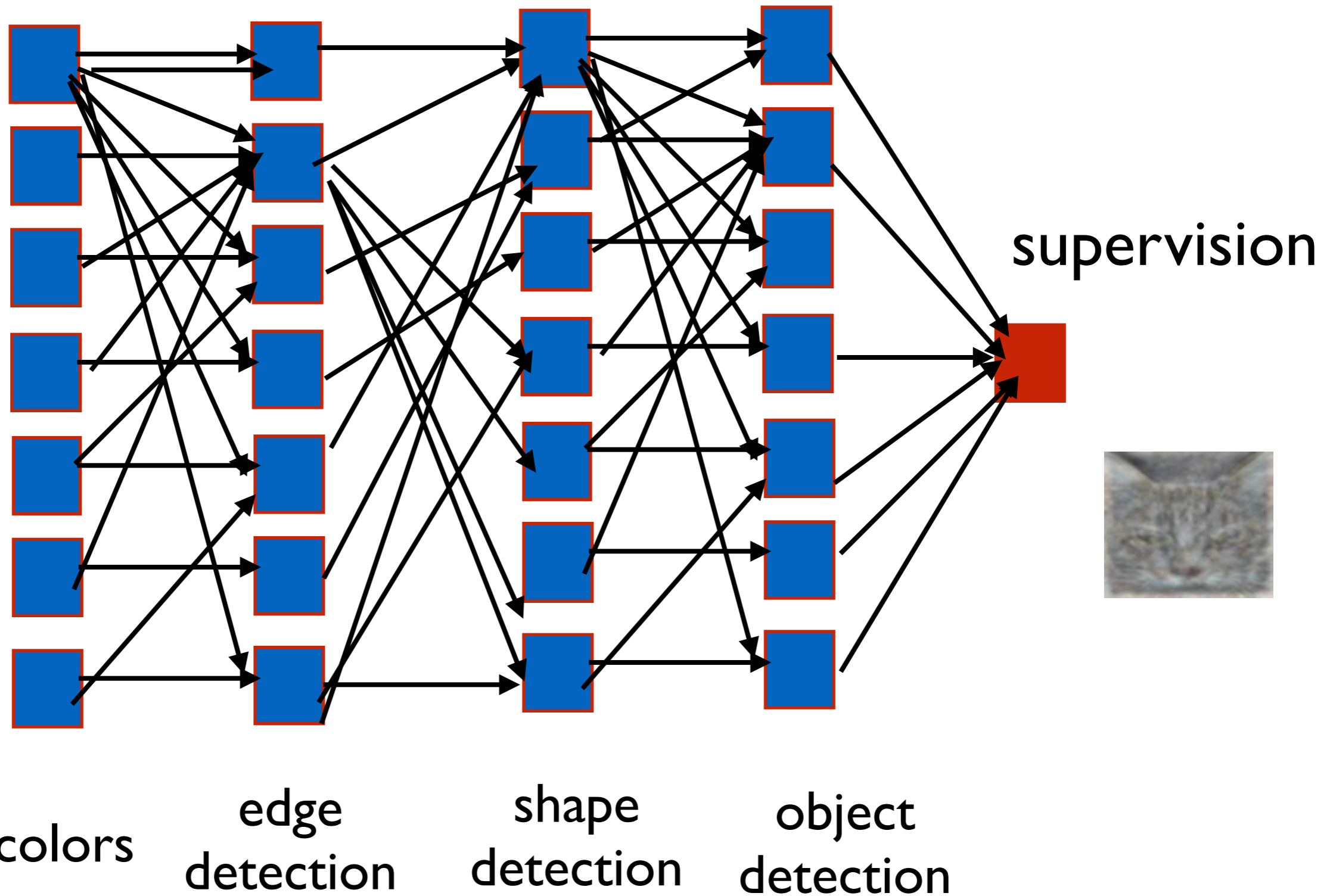
# Deep Learning

- So far, in all our machine learning, we designed features by ourselves. But can we do this automatically?

  - feature learning from data

- How can we combine different types of features and decide the useful combinations as part of the learning process?

- One solution

  - kernels

    - Considers only fixed, limited, and specific combinations. (eg. polynomial kernel considers only pairwise combinations)

- Another solution

  - Multi-layer perceptrons

    - Overfitting, difficulty to train, time consuming

- Can we train deep models (with many hidden layers) efficiently and without over-fitting?

  - This is the central problem considered in *Deep learning*

# Deep Learning at a glance



supervision

pixel colors     edge detection     shape detection     object detection

# Big news



Yan LeCun (Facebook)

The New York Ti...
Monday, June 25, 2012    Last Update: 11:50 PM ET
...e Hired to Make AI a Reality    WEEKS

**EXCLUSIVE**

# Facebook, Google in 'Deep Learning' Arms Race

Yann LeCun, an NYU artificial intelligence researcher who now works for Facebook. *Photo: Josh Valcarcel/WIRED*

**WIRED**

**NEWS BULLETIN**

Yoshua Bengio
Univ. of Toronto

Geoff Hinton (Google)

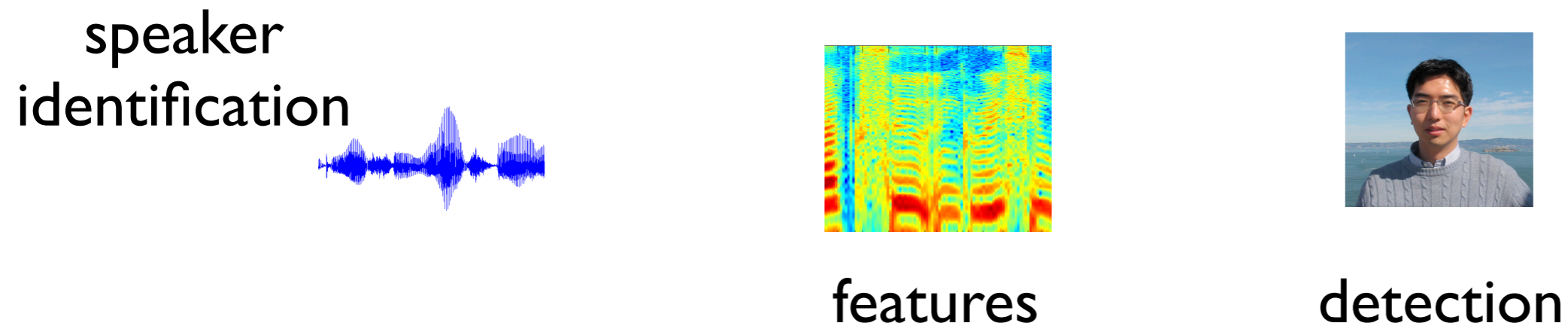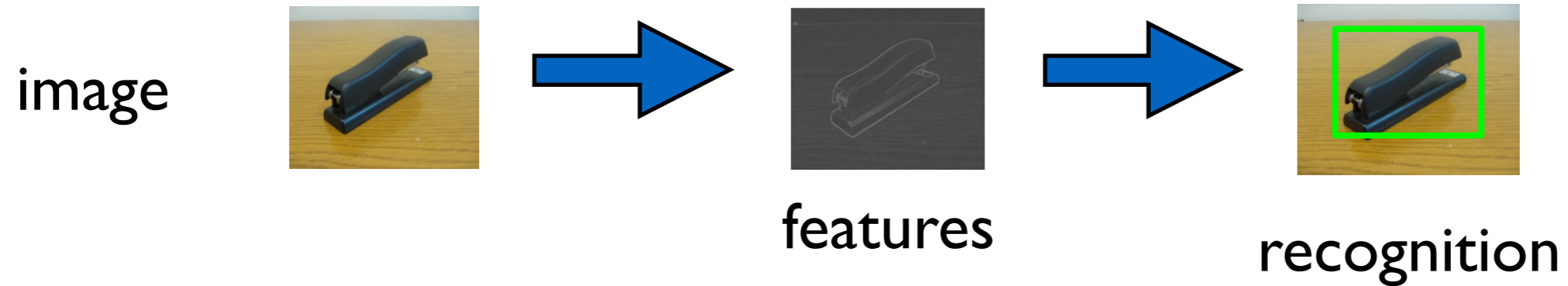# Google Beat Facebook for DeepMind

# Google Acquires Artificial Intelligence Startup DeepMind For More Than $500M

slide credit: Bengio KDD'14    4

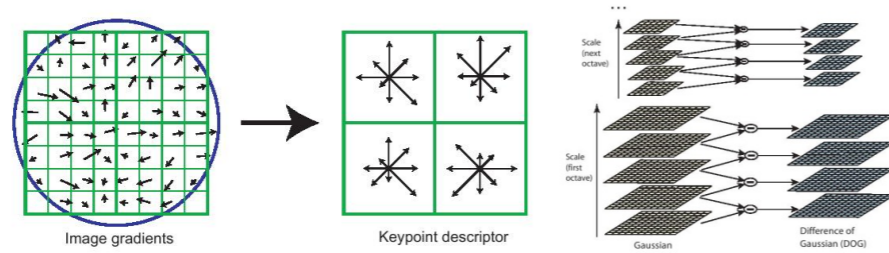# Applications of DL

- Image Recognition

    - ILVSCR 14, 15: outperformed human level

- NLP

    - Machine translation, text similarity, sentiment analysis

- Voice

    - Voice recognition

- Robotics

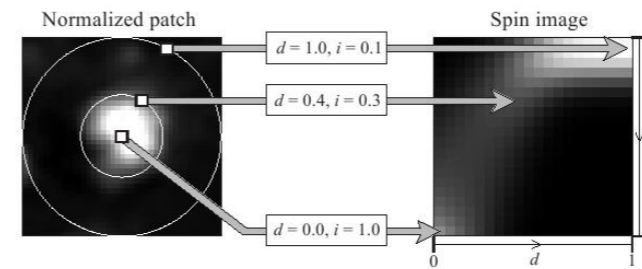    - DeepMind, Computer Games, Reinforcement Learning

# Object Recognition



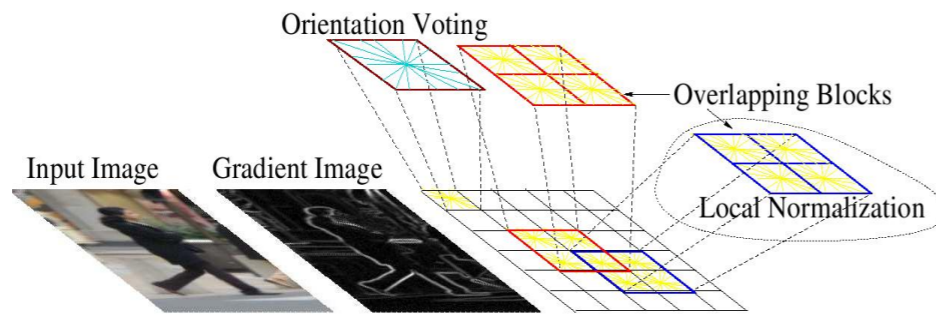input → features → classifier
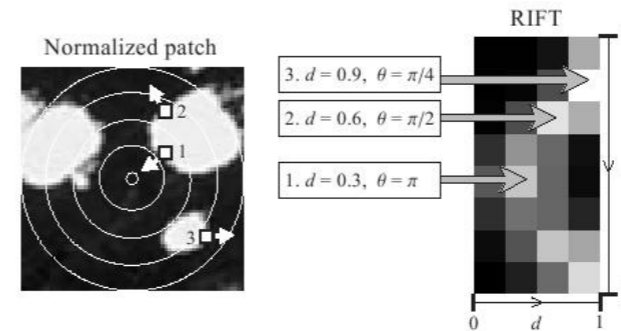
image
speaker identification
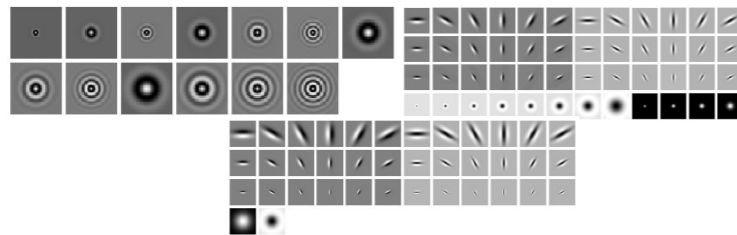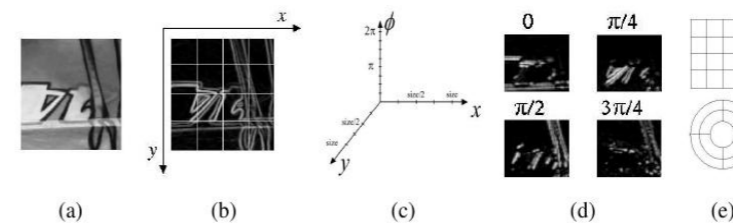features        detection

6

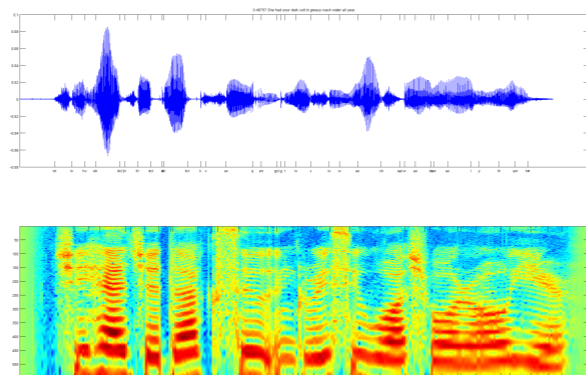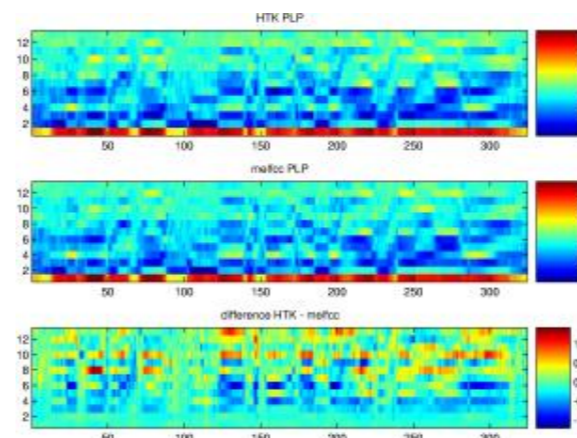# Image Features



SIFT



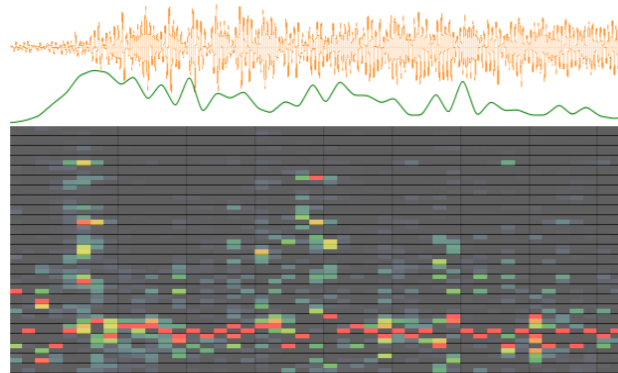Spin image



HoG



RIFT



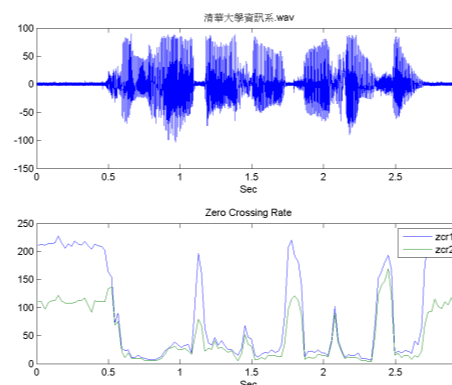Textons



GLOH

Slide Credit: Honglak Lee

# Voice Features



Spectrogram



MFCC



Flux



ZCR



Rolloff

# Image labeling



Krizhevsky+ NIPS'12

- demo：http://deeplearning.cs.toronto.edu/

# Similar image retrieval



Krizhevsky+ NIPS'12

## 20 layer NN by Google!

# Detecting roads from satellite images

# Word Analogy Detection

- How to learn representations for words?

- *you shall know a word by the company it keeps* — Firth

- Distributional Hypothesis

  - We can predict the meaning of a word by looking into its local context

- Can we learn vector representations for words such that we can accurately predict its neigbours in a sentence?

  - word2vec (skip-gram model) Mikolov+13

  - GloVe (Global Vector Prediction) Pennington+14

- v(king) - v(man) + v(woman) = v(queen)

# Brief History

- Neural networks were around even back in 1950s

- It was shown that you cannot learn non-linearly separable data using single layer neural networks (ca. Perceptron)

  - Marvin Minsky [1960]

- First NN winter

# Perceptron (revision)

- Perceptron is a single layer neural network

# Perceptron (revision)



$x_1$  $x_2$  $x_3$  ...  $x_n$

$w_1$

$w_2$

$w_n$

$w_n$

Logistic function

$s = x_1 w_1 + x_2 w_2 + \ldots + x_n w_n$
if s > 0:
    return 1
else:
    return 0

16

# linear separability



In 2D space linear separability means you can separate the two classes by a straight line.  ax + by +c = 0

# Non-linear separable case

XOR (exclusive OR)

|  | x=0 | x=1 |
| --- | --- | --- |
| y=0 | 0 | 1 |
| y=1 | 1 | 0 |

# Multi-layer Neural Networks

- XOR is a very common logical operation

  - Perceptron being unable to handle this common case was seen as a show stopper for neural networks

- We could get over this issue by using multiple hidden layers.

  - But there was no algorithm to learn the weights

  - Until, error backpropagation (Rummelhart+86) was proposed

- However, deep neural networks are likely to overfit and training them was time consuming (no GPUs back then!)

  - Second NN winter.

# Advantages of DL

- Deep learning can learn the features useful for a particular task automatically

  - Can use **unlabeled** data to learn the features

- Can learn distributed representations

# Local vs. Distributed Representations

• Clustering, Nearest Neighbors, RBF SVM, local density estimators

• RBMs, Factor models, PCA, Sparse Coding, Deep models

> • Parameters for each region.
> • # of regions is linear with # of parameters.



Learned prototypes

C1=1
C2=1
C3=1

C1=1
C2=1
C3=0

C1=0
C2=1
C3=0

C1=1
C2=0
C3=0

C1=0
C2=1
C3=1

C1=0
C2=0
C3=0

C1=0
C2=0
C3=1

C1    C2    C3

Bengio, 2009, Foundations and Trends in Machine Learning

# The breakthrough!

- It was shown that by learning two layers at a time, and then stacking those to create a deep neural network was an effective method for overfitting.

- greedy layer-wise training

  - *A Fast Learning Algorithm for Deep Belief Nets*, Hinton et al., Neural Computing, 2006.

# unsupervised pre-training



Even more abstract
features

More abstract
features

features

input

# supervised post-training

# Deep Learning Methods

- Two main techniques exist

- Autoencoders  (AE)

  - A non-probabilistic method

  - Easier to implement

  - Theoretical analysis is difficult (although some work has been done lately)

- Restricted Boltzman Machine (RBM)

  - A probabilistic method

  - Under certain conditions it could be shown that both RBMs and AEs are optimizing the same objective

# Autoencoder

$y$     **Feature Representations**

$y = f(Wx + b)$

**Encoder**

**Decoder**

$z = f(W^{T}y + b')$

$x$     **Input**

Reconstruction error $= ||x-z||^2$

# Autoencoder

Hidden layer

$y_1$  $y_2$  $y_3$  +1

$x_1$  $x_2$  $x_3$  +1

Input

# Autoencoder



hidden layer

$y_1$ $y_2$ $y_3$ +1

encoder

$b_3$

$W_{32}$

input layer

$w_{ij} : x_j \longrightarrow y_i$

# Details

- By using the transpose matrix W' for the decoder where W is the encoder matrix, we can reduce the number of parameters in the model. (less likely to overfit)

- b and b' are respectively encoding and decoding bias terms.

- Bias terms can be incorporated as features into the autoencoder by setting a feature that is always ON.

- The non-linear function f is performing some elementwise non-linear operation on each element of a vector.

- Without non-linearity, autoencoders are equivalent to PCA.

# Training procedure

1. encode the input x using the encoder

    1. Calculate Wx + b,   and insert it in f

2. Let the output of (1) be y. Insert y into the decoder and reconstruct the input

    1. Calculate $W^T y + b'$,   and insert it in f

3. Let the output of (2) be z. Compare z and x.

    1. The loss function to be used is problem specific. For real values a popular loss function is the squared loss. For binary values use the cross-entropy error function.

4. Adjust the parameters (W, b, b') such that the loss computed in (3) is minimized.

    1. Compute the partial derivative of the loss w.r.t. each parameter and apply the stochastic gradient descent method.

# Stochastic Gradient Descent (SGD)

- We have already seen SGD in Perceptron, logistic regression and multi-layer neural networks.

- Move in the opposite direction of the gradient of the loss

$$x^{(t+1)} = x^{(t)} - \eta \frac{\partial E(\boldsymbol{x}, \boldsymbol{z})}{\partial x}\big|_{t=t}$$

tan(θ)=∂y/∂x

θ

gradient is given by the derivative

# Example



$$y_1 = f(x_1 W_{11} + x_2 W_{12} + x_3 W_{13} + b_1)$$
$$y_2 = f(x_1 W_{21} + x_2 W_{22} + x_3 W_{23} + b_2)$$
$$y_3 = f(x_1 W_{31} + x_2 W_{32} + x_3 W_{33} + b_3)$$

$$f(t) = \frac{1}{1 + \exp(-t)}$$

# In matrix form…

$$\mathbf{W} = \left( \begin{array}{ccc} w_{11} & w_{12} & w_{13} \\ w_{21} & w_{22} & w_{23} \\ w_{31} & w_{32} & w_{33} \end{array} \right) \qquad \boldsymbol{b} = \left( \begin{array}{c} b_1 \\ b_2 \\ b_3 \end{array} \right) \qquad \boldsymbol{x} = \left( \begin{array}{c} x_1 \\ x_2 \\ x_3 \end{array} \right) \qquad \boldsymbol{y} = \left( \begin{array}{c} y_1 \\ y_2 \\ y_3 \end{array} \right)$$

$$\boldsymbol{z} = \left( \begin{array}{c} z_1 \\ z_2 \\ z_3 \end{array} \right) \qquad \boldsymbol{b}' = \left( \begin{array}{c} b_1' \\ b_2' \\ b_3' \end{array} \right)$$

$$\boldsymbol{y} = f(\mathbf{W}\boldsymbol{x} + \boldsymbol{b})$$

33

# Decoder becomes…

$$z_1 = f(y_1 W_{11} + y_2 W_{21} + y_3 W_{31} + b_1')$$

$$z_2 = f(y_1 W_{12} + y_2 W_{22} + y_3 W_{32} + b_2')$$

$$z_3 = f(y_1 W_{13} + y_2 W_{23} + y_3 W_{33} + b_3')$$

$$\boldsymbol{z} = f(\mathbf{W}^\top \boldsymbol{y} + \boldsymbol{b}')$$

squared loss $\quad ||\boldsymbol{x} - \boldsymbol{z}||^2 = \sum_{k=1}^{d}(x_k - z_k)^2$

# Parameter update

Let us consider the update of w₁₂

$$\frac{\partial \, \|\boldsymbol{x} - \boldsymbol{z}\|^2}{\partial w_{12}} = -2(\boldsymbol{x} - \boldsymbol{z})\frac{\partial z_2}{\partial w_{12}}$$

$$\frac{\partial z_2}{\partial w_{ij}} = \frac{\partial}{\partial w_{ij}}f(w_{12j}y_1 + w_{22}y_2 + w_{32}y_3 + b_2')$$

Here

$$t = w_{12j}y_1 + w_{22}y_2 + w_{32}y_3 + b_2'$$

$$\frac{\partial f(t)}{\partial w_{12}} = \frac{\partial f(t)}{\partial t}\frac{\partial t}{\partial w_{12}} = \frac{\exp(-t)}{(1 + \exp(-t))^2}y_1 = \sigma(t)(1 - \sigma(t))y_1$$

$$\frac{\partial \, \|\boldsymbol{x} - \boldsymbol{z}\|^2}{\partial w_{12}} = -2(\boldsymbol{x} - \boldsymbol{z})\sigma(t)(1 - \sigma(t))y_1$$

$$w_{12}^{(n+1)} = w_{12}^{(n)} + 2\eta(\boldsymbol{x} - \boldsymbol{z})\sigma(t)(1 - \sigma(t))y_1$$

# References

- Deep Learning Tutorial

    - http://deeplearning.net/tutorial/gettingstarted.html

    - git clone git://github.com/lisa-lab/DeepLearningTutorials.git

- You need Theano

    - http://deeplearning.net/software/theano/install.html

- Dependencies Python >= 2.6, g++, python-dev, NumPy, SciPy, BLAS)

- Can be installed via sudo apt-get install in Debian/ Ubuntu or for Mac OSX brew/macports.