# INTRODUCTION
# TO DEEP LEARNING
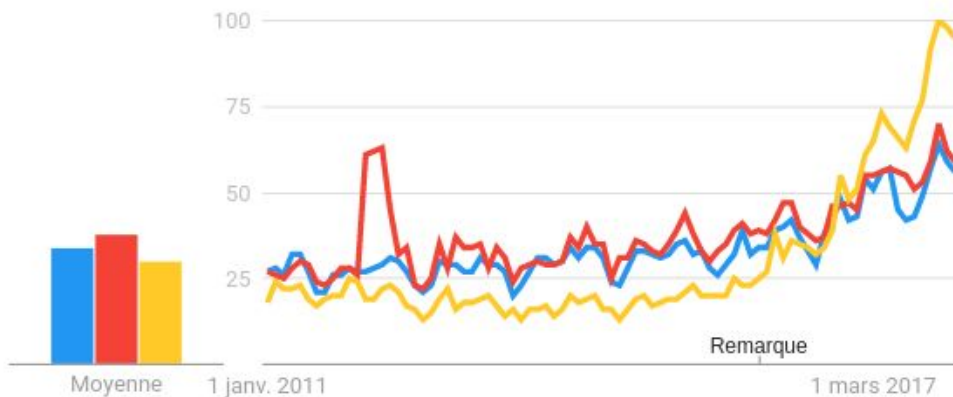
## Mouhidine SEIV
### FOUNDER - RIMINDER

**Deep LEARNING A PRACTICAL COURSE**
ECOLE POLYTECHNIQUE, 05/04/2018

# History of Learning Systems



Évolution de l'intérêt pour cette recherche — Google Trends

- Apprentissage supervisé
- Apprentissage non supervisé
- Apprentissage par renforcement

Moyenne — 1 janv. 2011 — Remarque — 1 mars 2017

Dans tous les pays. 01/01/2011 – 29/01/2018. Recherche sur le Web.

# Program & Course Logistics

- **Course 1 :** (05-04-18)
    - Introduction to Deep Learning - Mouhidine SEIV (Riminder)
- **Course 2 :** (12-04-18)
    - Deep Learning in Computer Vision - Slim FRIKHA (Riminder)
- **Course 3 :** (19-04-18)
    - Deep Learning in NLP - Paul COURSAUX  (Riminder)
- **Course 4 :** (26-04-18)
    - Efficient Methods and Compression for Deep Learning - INVITED GUEST
- **Course 5:** (03-05-18)
    - Introduction to Deep Learning Frameworks - INVITED GUEST
- **Course 6:** (10-05-18)
    - Deployment in Production and Parallel Computing - INVITED GUEST

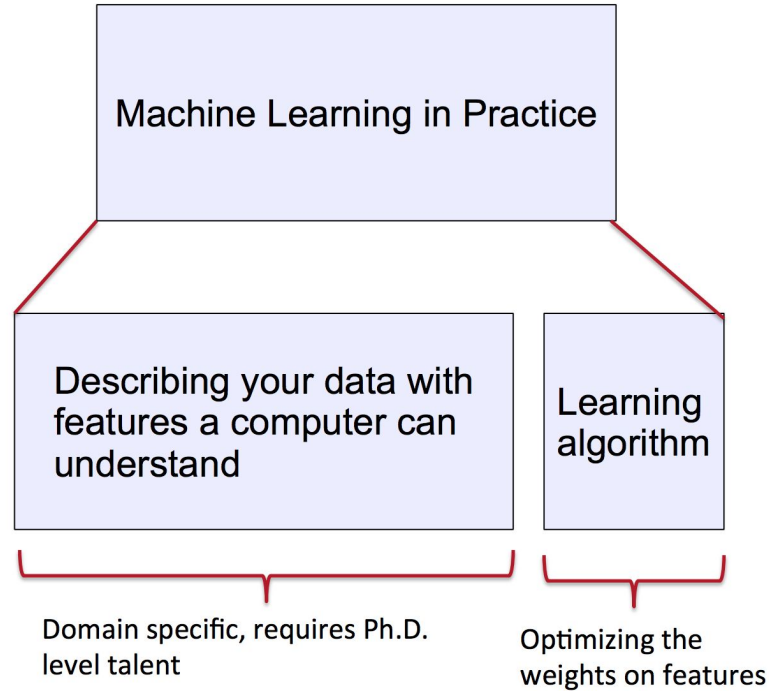**Location: Ecole Polytechnique from 6:30 pm to 7:30pm**

**https://github.com/riminder**

# Prerequisites

- Proficiency in python

- Linear algebra

- Basic probability and statistics

- Basic machine learning ( cost functions, derivatives , gradient methods)

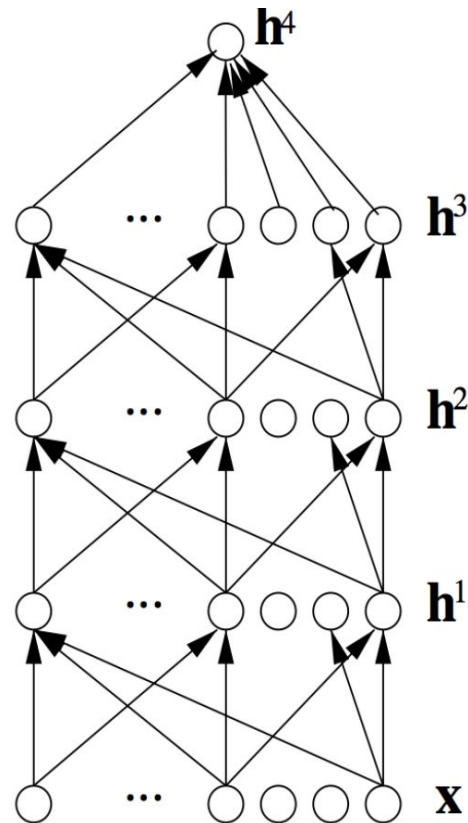# Machine Learning vs Deep Learning

# Reasons for Exploring Deep Learning

- In 2006 deep learning techniques started outperforming other Machine learning techniques. Why now?

- DL techniques benefit more from a lot of data
- Faster machines and multicore CPU/GPU help DL
- New models, algorithms, ideas

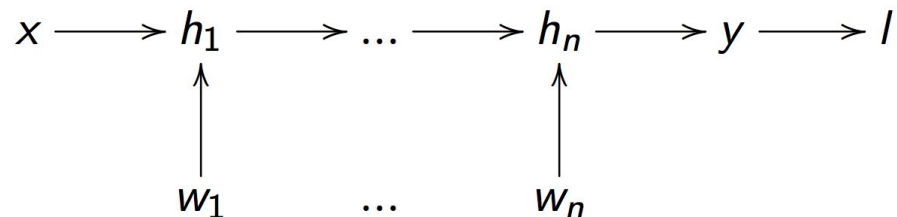→ **Improved performance (first in speech and vision, then NLP)**

# What is Deep Learning?

- Representation learning attempts to automatically learn good features or representations

- Deep learning algorithms attempt to learn (multiple levels of) representation and an outputs
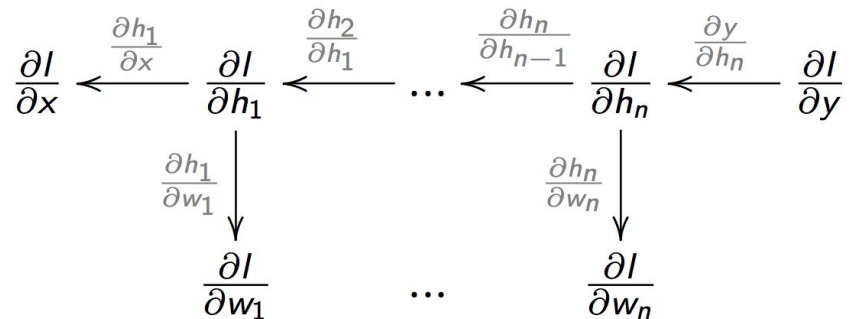
- From "raw" inputs $x$ (e.g. words)

# Deep Representations

- A deep representation is a composition of many functions

$$x \longrightarrow h_1 \longrightarrow \ldots \longrightarrow h_n \longrightarrow y \longrightarrow l$$

$$w_1 \qquad \ldots \qquad w_n$$

- Its gradient can be backpropagated by the chain rule

$$\frac{\partial l}{\partial x} \xleftarrow{\frac{\partial h_1}{\partial x}} \frac{\partial l}{\partial h_1} \xleftarrow{\frac{\partial h_2}{\partial h_1}} \ldots \xleftarrow{\frac{\partial h_n}{\partial h_{n-1}}} \frac{\partial l}{\partial h_n} \xleftarrow{\frac{\partial y}{\partial h_n}} \frac{\partial l}{\partial y}$$

$$\frac{\partial h_1}{\partial w_1} \Big\downarrow \qquad\qquad \frac{\partial h_n}{\partial w_n} \Big\downarrow$$

$$\frac{\partial l}{\partial w_1} \qquad \ldots \qquad \frac{\partial l}{\partial w_n}$$

# Deep Neural Network

- A **deep neural network** is typically composed of:
    - Linear transformations

$$h_{k+1} = Wh_k$$

- Non-linear activation functions

$$h_{k+2} = f(h_{k+1})$$

- A loss function on the output, e.g.
    - Mean-squared error

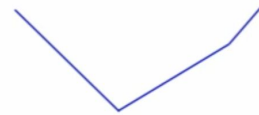$$I = ||y^* - y||^2$$

- Log likelihood

$$I = \log \mathbb{P}[y^*]$$
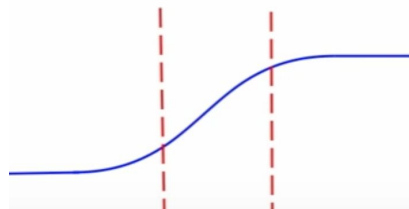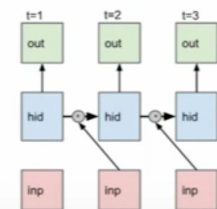
Non-linear activation functions



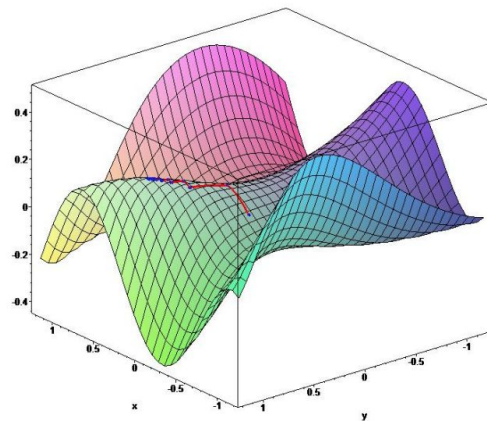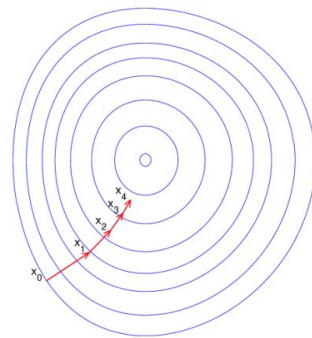Rectified linear unit

Maxout

Carefully tuned sigmoid

LSTM

# Training Neural Networks

- Sample gradient of expected loss L(w) = E [l]

$$\frac{\partial l}{\partial \mathbf{w}} \sim \mathbb{E}\left[\frac{\partial l}{\partial \mathbf{w}}\right] = \frac{\partial L(\mathbf{w})}{\partial \mathbf{w}}$$
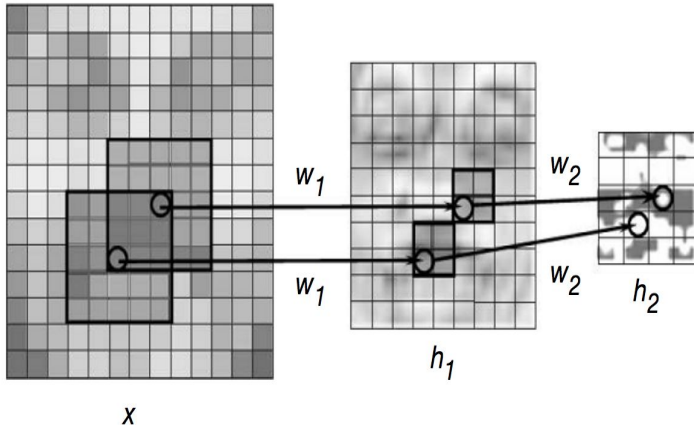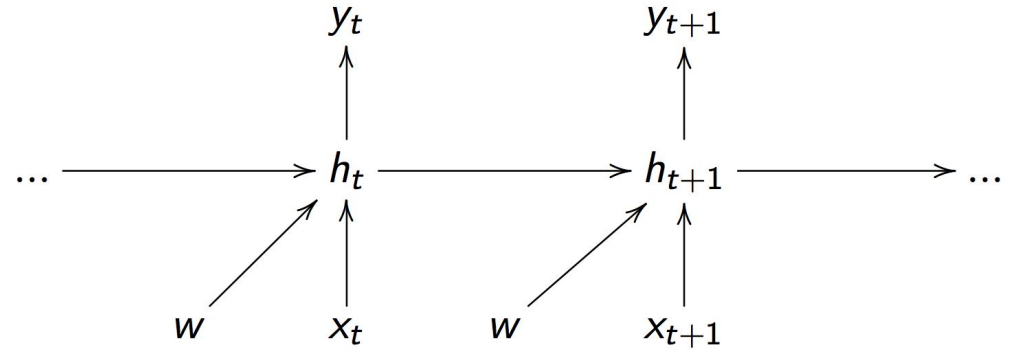
- Adjust **w** down the sampled gradient

$$\Delta w \propto \frac{\partial l}{\partial \mathbf{w}}$$

# Weight Sharing

- Recurrent neural network
  shares weights between time-steps

$$y_t \qquad y_{t+1}$$

$$\ldots \longrightarrow h_t \longrightarrow h_{t+1} \longrightarrow \ldots$$

$$w \quad x_t \qquad w \quad x_{t+1}$$

- Convolutional neural network
  shares weights between local regions
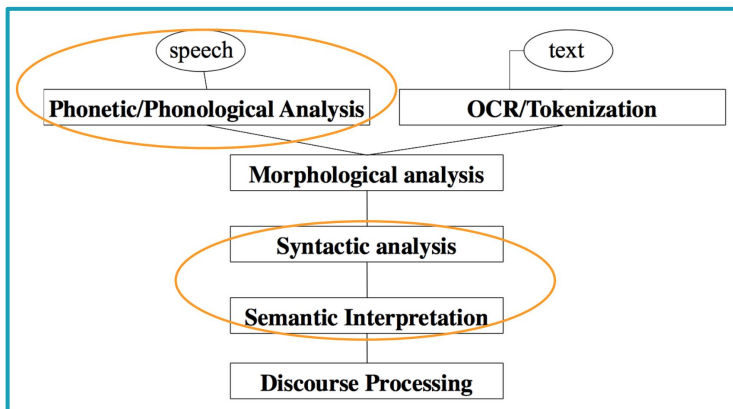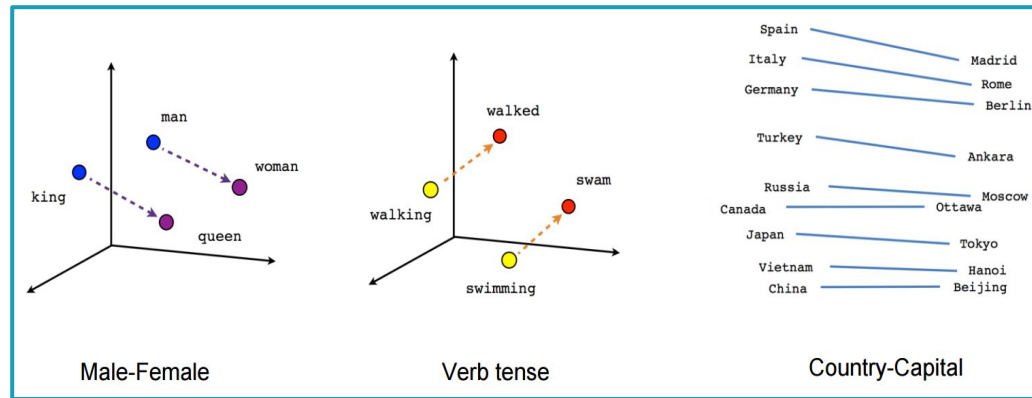
$w_1$

$w_1$

$w_2$

$w_2$

$h_2$

$h_1$

$x$

# What is NLP?

- Natural language processing is a field at the intersection of
    - Computer science
    - Artificial intelligence
    - and linguistics.

- Goal: for computers to process or "understand" natural language inorder to perform tasks that are useful, e.g.
    - Question Answering

- Fully understanding and representing the meaning of language (or even defining it) is an illusive goal.
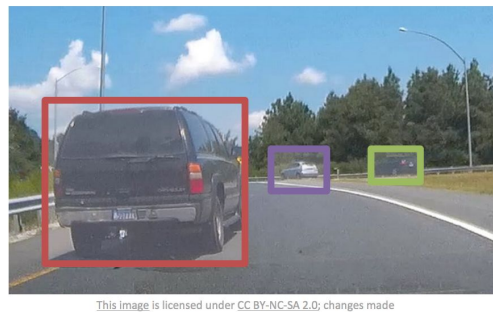- Perfect language understanding is AI-complete
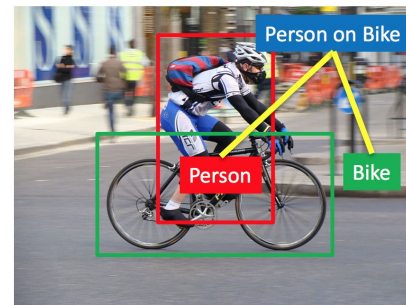
# Deep Learning for NLP



BEFORE Deep Learning



After Deep Learning

# What is Computer Vision?

- Most deep learning groups have (until 2 years ago) focused on computer vision

- Break through paper: ImageNet Classification with Deep Convolutional Neural Networks by Krizhevsky et al. 2012


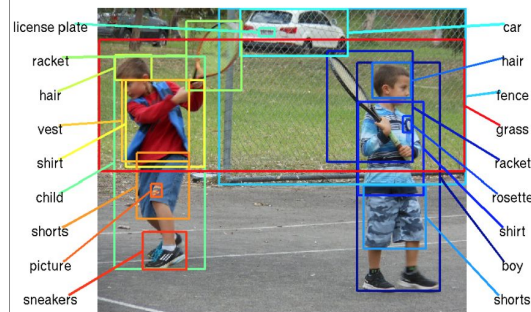This image is licensed under CC BY-NC-SA 2.0; changes made


Person on Bike
Person
Bike
This image is licensed under CC BY-SA 3.0; changes made


flamingo    cock    ruffed grouse    quail    partridge    …

Egyptian cat    Persian cat    Siamese cat    tabby    lynx    …


Person
Hammer
This image is licensed under CC BY-SA 2.0; changes made


license plate, racket, hair, vest, shirt, child, shorts, picture, sneakers
car, hair, fence, grass, racket, rosette, shirt, boy, shorts
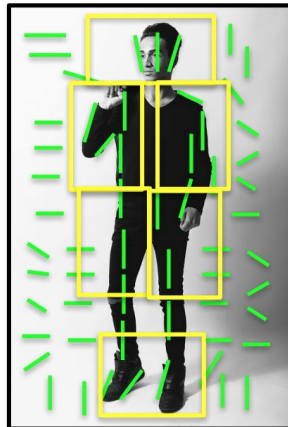
**https://github.com/riminder**

# Deep Learning for Computer Vision



BEFORE Deep Learning

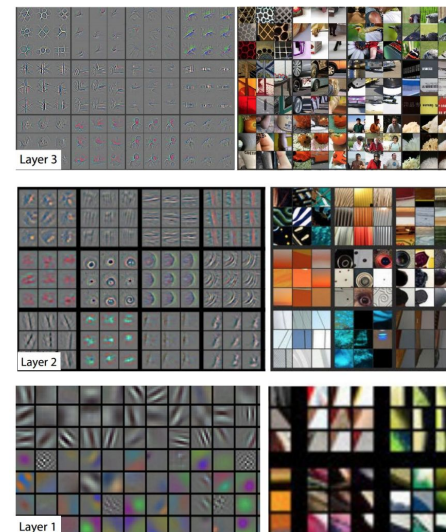After Deep Learning

Histogram of Gradients (HoG)
Dalal & Triggs, 2005

Deformable Part Model
Felzenswalb, McAllester, Ramanan, 2009

Zeiler and Fergus (2013)
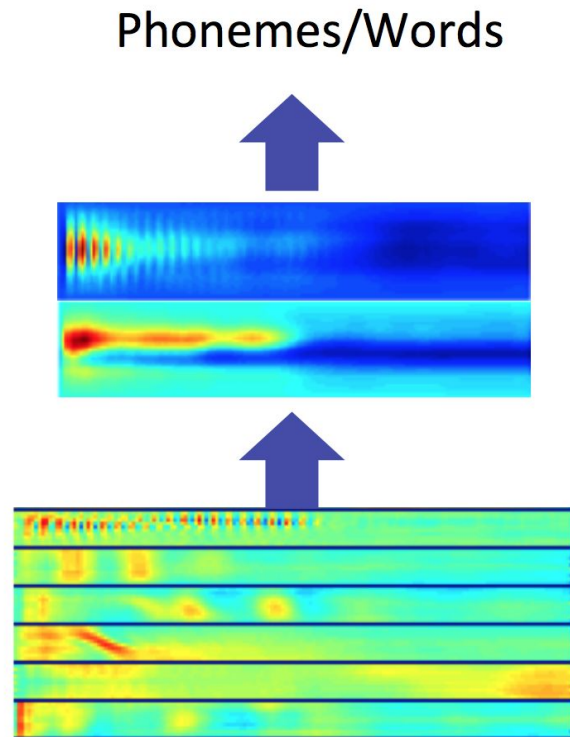
https://github.com/riminder

# Deep Learning for Speech

- The first breakthrough results of "deep learning" on large datasets happened in speech recognition

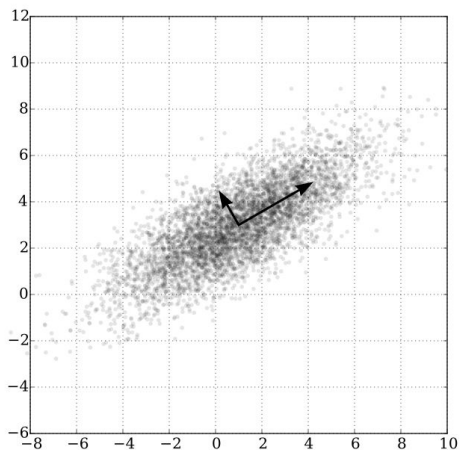- Context-Dependent Pre-trained Deep Neural Networks for Large Vocabulary Speech Recognition Dahl al. (2010)

| Acoustic model | Recog \ WER | RT03S FSH | Hub5 SWB |
|---|---|---|---|
| Traditional features | 1-pass –adapt | **27.4** | **23.6** |
| Deep Learning | 1-pass –adapt | **18.5** (−33%) | **16.1** (−32%) |

Phonemes/Words

# Visualization

- Singular Value Decomposition (SVD)

- Principal component analysis (PCA)





$\hat{X}$ is the best rank $k$ approximation to $X$, in terms of least squares.
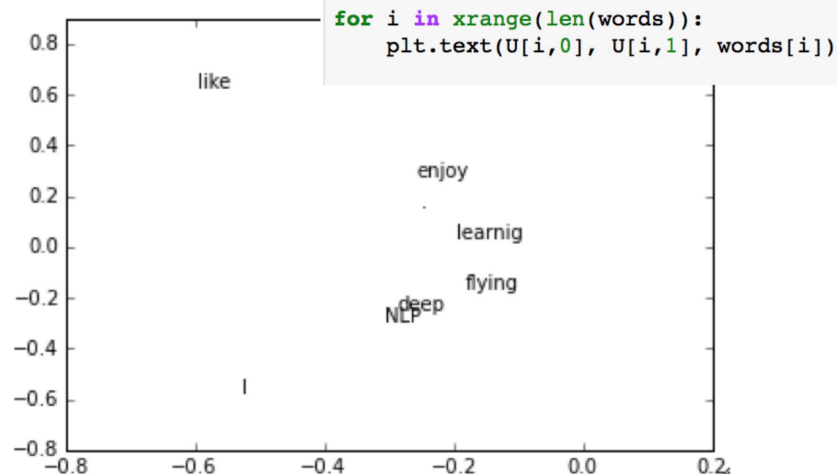
- t-distributed stochastic neighbor embedding (T-SNE)
  https://distill.pub/2016/misread-tsne/

# Singular Value Decomposition (SVD)

**Corpus:**

- I like deep learning. I like NLP. I enjoy flying.

```python
import numpy as np
la = np.linalg
words = ["I", "like", "enjoy",
         "deep","learnig","NLP","flying","."]
X = np.array([[0,2,1,0,0,0,0,0],
              [2,0,0,1,0,1,0,0],
              [1,0,0,0,0,0,1,0],
              [0,1,0,0,1,0,0,0],
              [0,0,1,0,0,0,0,1],
              [0,1,0,0,0,0,0,1],
              [0,0,1,0,0,0,0,1],
              [0,0,0,0,1,1,1,0]])

U, s, Vh = la.svd(X, full_matrices=False)
```



- Printing first two columns of U corresponding to the 2 biggest singular values.

# Trouble with Learning systems



Adversarial example
98% Toaster



Bias example
95% Unqualified

# GENDER STEREOTYPING : WORD EMBEDDINGS

**HE**

homemaker

nurse

receptionist

librarian

**SOFTBALL**

pitcher

bookkeeper

registered nurse

waitress

**SHE**

maestro

skipper

protege

philosopher

**FOOTBALL**

footballer

businessman

maestro

cleric

# GENDER STEREOTYPING : REPRESENTATION

**Formulation:**

- Let B be the gender subspace
- "Hard bias correction": For a gender neutral word w , set
  $w \leftarrow w - w_B \; / \; ||w - w_B||$
- "Soft Bias correction": For word matrix W, gender neutral words N find a linear transformation T

$$\min_T ||(TW)^T \, TW - W^T \, W||_F^2 + \lambda ||(TN)^T \, TB - ||_F^2$$

**Results:**

- Reduced gender stereotyping in the word embeddings
- Performance on downstream tasks still almost the same

Man Is To Computer Programmer as Woman is to Homemaker? Debiasing Word Embeddings. T Bolukbasi, K-W Chang, J Zou, V Saligrama, A Kalai. Arxiv 2016

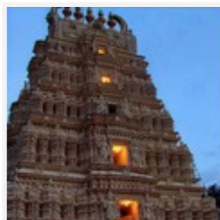# Adversarial examples



(a) Image

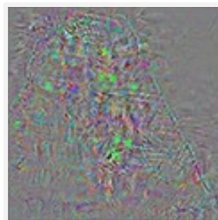(b) Prediction

(c) Adversarial Example
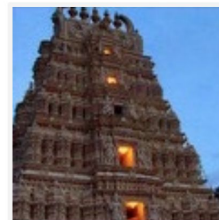
(d) Prediction

Original image
Temple (97%)

Perturbations

Adversarial example
Ostrich (98%)

# Generating Adversarial Examples

The Fast Gradient Sign Method

$$J(\tilde{x}, \theta) \approx J(x, \theta) + (\tilde{x} - x)^\top \nabla_x J(x).$$

Maximize

$$J(x, \theta) + (\tilde{x} - x)^\top \nabla_x J(x)$$

subject to

$$||\tilde{x} - x||_\infty \leq \epsilon$$

$$\Rightarrow \tilde{x} = x + \epsilon \operatorname{sign} \left( \nabla_x J(x) \right).$$

Modern deep nets are very piecewise linear

Rectified linear unit

Maxout

Carefully tuned sigmoid

LSTM

t=1   t=2   t=3

out   out   out

hid   hid   hid

inp   inp   inp
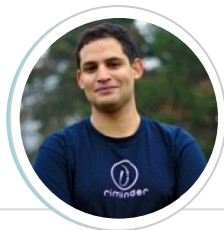
riminder

# POINT OF **CONTACT**



Mouhidine SEIV

CEO & Founder @Riminder

mouhidine.seiv@riminder.net