Adobe Workshop on Machine Learning

# How many of you

- Are you eager to work in a real project with Machine Learning (ML)?
- Have you used machine learning on your hobby projects?
- Are involved in projects that use ML technology implementation or related activities?
- Feel you understand the underlying ML paradigm?

# About me

- A Machine Learning enthusiast and expert;
- 10 years of experience in software development;
- Associate Professor at UPB, ETTI;
- 8 year of experience in research:

3

- PhD and postdoc research in the field of computer vision and multimedia retrieval;
- Over 40 publications in international journals, conferences and workshops;
- Participation at various international competitions;
- Part of the organizing team for various conferences and competitions;
- Technical reviewer for various journals, conferences and workshops.

#### More info on http://ionut.mironica.ro





# PART 1



• What is and what types of ML exist?

- How does is work?
- How can be used?
- How we can combine different ML algorithms?
- Ethics in ML ?
- Applications



- Types of Machine Learning
- Feature extraction
- Classical Machine Learning algorithms
- Deep learning algorithms
- Tools for Machine Learning

# **Artificial Intelligence**

Artificial Intelligence (AI) is the science of making the things smart and it can be defined as:

### "Human Intelligence exhibited by machines"

A broad term of getting computers to perform human tasks. The scope of AI is disputed and constantly changing over the time.



[Images from: http://bryanrussell.org and http://m.el-dosuky.com]

# **Machine Learning**

Machine Learning (ML) can be defined as:

### "An approach to achieve artificial intelligence through systems that can learn from experience to find patterns in that data"

ML involves teaching computer by examples and learning patterns and not programming specific rules. And these patterns can be found in **data**.

# **Deep Learning**

9

The subset of machine learning composed of algorithms that permit software to train itself to perform tasks, like speech and image recognition, by exposing multilayered neural networks to vast amounts of data.



[Image from: http://www.global-engage.com/life-science/deep-learning-in-digital-pathology]

# **Machine Learning**



# **Classical programming vs Machine learning**



- Features are used to train a ML system.
- Represents properties of the things that you are trying to learn about

#### Weather forecasting

- humidity
- pressure
- temperature





• sepal / petal width / height

[Image from: www.collinsdictionary.com]

### WEATHER FORECASTING

- humidity
- pressure

Choosing the useful features can have a high impact on the system's classification accuracy.



### WEATHER FORECASTING

- humidity
- pressure

Choosing the useful features is not a trivial task

There are systems that may have millions of features, so the visualization of the feature space is almost impossible



Adding more dimensions may help on improving the performance, allowing the algorithm to better separate the classes

More dimensions can generate

#### "CURSE OF DIMENSIONALITY":

A phenomena that occurs when the dimensionality of the data increases, the sparsity of the data increases.



[Image from: https://www.futuristspeaker.com]



#### **UNSUPERVISED LEARNING**

- Training data doesn't contain output;
- The basic idea is to find templates and patterns in data to be automatically highlighted.





Unsupervised Learning

17



#### **UNSUPERVISED LEARNING**



#### **UNSUPERVISED LEARNING**



- Similarity a subjective concept;
- Hard to define at human level.



[Images from: www.cesar.com and www.purina.com]

#### **SEMI-SUPERVISED LEARNING**

 Semi-supervised learning is a class of supervised learning tasks and techniques that also make use of unlabeled data for training – typically a small amount of labeled data with a large amount of unlabeled data.





#### **REINFORCEMENT LEARNING**

- Instead of optimizing a cost function is maximizing a reward function
- Learning by trial and error throw reward and punishment;
- The program learns by playing the game millions on times





# **Machine learning algorithms**

#### SUPERVISED LEARNING



#### What is the best algorithm?

## No free lunch theorem



#### [Image from: https://cynicalbabblings.wordpress.com]



**UNDERFITTING** A model is underfitted when it is too simple with regards to the data that is trying to model.



This is because the model is unable to capture the relationship between the input examples (called X) and the target values (called Y).

**OVERFITTING** The model is overfitted when on the training data the model performs well but it does not have reasonable results on the evaluation data.



This is because the model is memorizing the data used during the training sample and it is unable to generalize to unseen examples.





- Increase the complexity / performance of your features
- Remove unnecessary features



- Add more data
- Reduce the model complexity
- Add techniques against overfitting:
  - regularizations, normalizations, noise
- Remove unnecessary features

# **A classical Machine Learning flow**



## **Databases used on training the model**

Center for Naval Analyses to minimize the plane losses during next missions





# **A classical Machine Learning flow**

#### **COMPONENTS**

- Feature extraction Today
- Machine Learning algorithms:
  - Classical approaches Course 2
  - Deep Learning Course 3
- Performance measurement Course 4
- •Hands-on Course 5

### **FEATURE EXTRACTION**

....

The algorithms we used are very standard for Kagglers. [...] We spent most of our efforts in feature engineering.

- Xavier Conort, on "Q&A with Xavier Conort" on winning the Flight Quest challenge on Kaggle

Better features means flexibility.

Better features means simpler models.

Better features means better results.





**KEYPOINTS** 


••••••

#### **KEYPOINTS**

General algorithm

(1) detects the regions that contains the keypoints;

(2) for each keypoint we need to define a neighborhood and extract a descriptor;

(3) compute the distance between the keypoints extracted from the object template and the test dataset.

#### **KEYPOINTS**



#### **KEYPOINTS**

**Objectives:** 

- Repetitivity:
  - invariant to translation, rotation, scale change;
  - invariant la oclusions;
  - Invariant to ilumination changes.
- Precise localization;
- Relevant content.

#### **KEYPOINTS**

#### **Dense extraction**



#### **Selective extraction**



Which is the best algorithm?



- Robust to oclusions; •
- More interest points. ٠

- More errors (one keypoint can be part of • more objects);
- Description is more precise. •



### **KEYPOINTS** Algoritmul SIFT

Cum se pot detecta locații care sunt invariante la schimbări de scală ale imaginii?

Soluție: utilizarea funcției de diferențe de gausiene (DoG).

• Pentru o imagine I(x,y), fie o reprezentare liniară a acesteia:

$$L(x, y, \sigma) = G(x, y, \sigma) * I(x, y)$$
$$G(x, y, \sigma) = \frac{1}{2\pi\sigma^2} e^{-(x^2 + y^2)/2\sigma^2}$$

 Se caută valorile de minim/maxim local pentru o serie de diferențe de gausiene:

$$D(x, y, \sigma) = (G(x, y, k\sigma) - G(x, y, \sigma)) * I(x, y)$$
$$= L(x, y, k\sigma) - L(x, y, \sigma)$$

## **Feature extraction** FILTRU GAUSSIAN - REAMINTIRE

 $\sigma$  = 1, W = 5

0.002969020.013306210.021938230.013306210.002969020.013306210.05963430.098320330.05963430.013306210.021938230.098320330.162102820.098320330.021938230.013306210.05963430.098320330.05963430.013306210.002969020.013306210.021938230.013306210.00296902



Algoritmul SIFT - Detecția punctelor de extrem



Algoritmul SIFT - Detecția punctelor de extrem



DoG

Dacă X este cel mai mare sau cel mai mic dintre toți vecinii, atunci X este denumit punct de interes (keypoint).

DoG

#### Algoritmul SIFT - Detecția punctelor de extrem

(a) puncte de interes înainte de aplicarea pragului;(b) puncte de interes după aplicarea pragului;



În funcție de selecția pragurilor anterioare vom avea un număr mai mare / mai mic de puncte de interes.



(a) Top 100 keypoints.(b) Top 200 keypoints

- 1. Se preia o fereastră de 16 x16 în jurul punctului de interes.
- 2. Se împarte în 4x4 celule.
- 3. Se calculează o histogramă de gradienți pe 8 orientări pentru fiecare celulă.



16 histograme x 8 orientări = 128 trăsături

#### PCA-SIFT

- Algoritmul este identic cu SIFT, mai puțin pasul 4 (calculul descriptorului);
- În loc să se utilizeze histogramele ponderate din algoritmul SIFT clasic, se calculează gradienții (verticali / orizontali) locali pe o suprafață de 41x41 pixeli (2 pixeli reprezintă conturul);
- Se concatenează toate valorile într-un singur vector: 2 x 39 x 39 = 3042 elemente.



#### PCA-SIFT

 Se reduce dimensionalitatea vectorului folosind Analiza Componentelor Principale (Principal Component Analysis - PCA) - de ex: de la 3042 la 64 / 36 / 20;

$$\begin{array}{c} PCA \\ \hline & & \\ \hline & & \\ \hline & & \\ N \times 1 \end{array} \qquad A_{K \times N} I_{N \times 1} = I'_{K \times 1}$$

- GLOH reprezintă un descriptor similar cu PCA-SIFT, singura diferență fiind calculul descriptorului final:
- Se împarte spațiul din jurul punctului de interes în coordonate log-polare;
- Lungime descriptor: (2 x 8 + 1) \* 16 = 272;
- Se aplică PCA și se rețin 128 elemente.



- Proprietăți:
  - Acuratețe mai mare de clasificare;
  - Viteză mai mică;
  - Mai rezistent la zgomot.

[Mikolajczyk & Schmid '05]

Algoritmul **SURF** (Speeded Up Robust Features) îmbunătățește viteza de calcul prin:

(1) utilizarea matricei de aproximare Hessiană

$$H(x,\sigma) = \begin{bmatrix} L_{xx}(\sigma) & L_{xy}(\sigma) \\ L_{yx}(\sigma) & L_{yy}(\sigma) \end{bmatrix}$$

(2) a imaginii integrale în calculul descriptorului.

Ce este o imagine integrală?

### **III. Descriptori locali** UTILIZAREA PUNCTELOR DE INTERES SURF – Imaginea integrală

Imaginea integrală  $I_{\Sigma}(x,y)$  a unei imagini I(x, y) reprezintă suma tuturor pixelilor din I(x,y) dintr-o regiune dreptunghiulară.



$$I_{\Sigma}(x, y) = \sum_{i=0}^{i \le x} \sum_{j=0}^{j \le y} I(i, j)$$

Prin utilizarea imaginii integrale este nevoie de doar patru valori pentru calculul sumei pixelilor dintr-o suprafată dreptunghiulară

$$S = A - B - C + D$$

#### SURF – Imaginea integrală



$$S = A - B - C + D$$

#### IVOM – dr.ing. Ionuţ Mironică

### **III. Descriptori locali** UTILIZAREA PUNCTELOR DE INTERES SURF – Imaginea integrală



$$S = A - B - C + D$$

S = 64-32-32+16 = 16

#### IVOM – dr.ing. Ionuț Mironică

### **III. Descriptori locali** UTILIZAREA PUNCTELOR DE INTERES SURF – Aproximarea matricei Hessiene

Pentru  $\sigma$ =1.2 (9x9), funcția LoG poate fi aproximată din:



IVOM – dr.ing. Ionuț Mironică

#### SURF – Aproximarea matricei Hessiene



Acestea pot fi foarte ușor calculate utilizând principiul de imagine integrală.

#### SURF – Aproximarea matricei Hessiene

Similar cu algoritmul SIFT, se aplică filtrul se calculează la mai multe octave (3 sau 4 în funcție de implementare):



#### SURF – Aproximarea matricei Hessiene

Similar cu algoritmul SIFT, se aplică filtrul se calculează la mai multe scale (3 sau 4 în funcție de implementare):



06.05.2019

IVOM – dr.ing. Ionuţ Mironică



#### IVOM – dr.ing. Ionuț Mironică



[https://github.com/imironica/IVOM-Demo/tree/master/IVOM\_Demo/KeypointsDetector]

06.05.2019



[https://github.com/imironica/IVOM-Demo/tree/master/IVOM\_Demo/KeypointsDetector]



[https://github.com/imironica/IVOM-Demo/tree/master/IVOM\_Demo/KeypointsDetector]

## **III. Descriptori locali** AGREGAREA PUNCTELOR DE INTERES

Căutarea cu puncte de interes

Algoritmul de căutare poate fi unul destul de consumator de resurse de calcul.

Posibilă soluție:

• Agregarea punctelor de interes într-un descriptor global.



- Extragere trăsături (puncte de interes / descriptori de mișcare)

- Învățarea unui dicționar (vizual / de mișcare)

- Cuantizarea trăsăturilor prin utilizarea vocabularului

- Reprezentarea cadrelor prin utilizarea frecvenței de apariție cuvintelor

[Slide-uri adaptate din prezentările lui Rob Fergus, Svetlana Lazebnik și Noah Snavely]

- Extragere trăsături (puncte de interes / descriptori de mișcare)







- Învățarea unui dicționar (vizual / de mișcare)



- K-means
- Clusterizare ierarhică
- Gaussian Mixture Model
- Arbori aleatori

- Învățarea unui dicționar vizual



- Cat de mare trebuie să fie un dicționar?
  - Prea mic: numărul de cuvinte vizuale nu vor fi reprezentative pentru toate conceptele
  - Prea mare: supra-învățare (overfitting)

Există metode care propun de la câteva sute de cuvinte la sute de mii de cuvinte.

- Reprezentarea cadrelor prin utilizarea frecvenței de apariție cuvintelor






## III. Descriptori locali

## **MODELUL BAG-OF-WORDS**

Descriptori - Bag of Words

• Au rezultate bune atunci când obiectele sunt asemănătoare



• Dar ce facem cu scaunele?



## III. Descriptori locali MODELUL BAG-OF-WORDS

## Dezavantaje model "Bag of Words"

- nu există nici o metodă riguroasă de reprezentare a distribuției spațiale dintre anumite perechi de cuvinte.
- există multe cuvinte care nu sunt relevante
- procesul de cuantizare a cuvintelor generează zgomot de cuantizare.
- costul computațional crește foarte mult odată cu dimensiunea vocabularului de cuvinte.

•

III. Descriptori locali MODELUL FISHER KERNEL FK vs BoW

Bag of Words conține apartenența fiecărui punct proeminent către un element al unui dicționar (histogramă de cuvinte)

Rezultat: D = [0;0;0;1];



Dimensiune: K

## III. Descriptori locali MODELUL FISHER KERNEL

## FK vs BoW

**Fisher Kernels** 

Calculează probabilitățile de apartenență la un cuvânt din dicționar

Rezultat: D = [0.3;0.1;0.1;0.5];

calculează gradientul mediei și
 a varianței probabilităților de
 apartenență la un cuvânt din dicționar.



Dimensiune: 2\*D\*K

## III. Descriptori locali MODELUL FISHER KERNEL Îmbunătățiri FK

Normalizare L2

- elimină erorile ce apar din diferența de scală a obiectelor Normalizare de putere (Power Normalization)

- eliminarea efectului de matrice rară (majoritatea elementelor din FK au valori foarte mici)

 $f(z) = \operatorname{sign}(z)|z|^{\alpha}$ 

Aplicare Piramide Spațiale

- utilizează informația geometrică a obiectelor



III. Descriptori locali MODELUL FISHER KERNEL Piramide spațiale



IVOM – dr.ing. Ionuț Mironică

## **Modelul Fisher kernel**

## Arhitectura reprezentării "Fisher kernel"



IVOM – dr.ing. Ionuț Mironică

**VIDEO CLASSIFICATION** Visual You Tube Sign in - color Q Upload - texture GUIDE NEW - shape Dalida & Alain Delon - Paroles, - keypoints paroles 60 pascalocool 6,066 views Motion Norah Jones - Feels Like Home by kritikospa12 422,131 views Andre Rieu - Romantic Paradise **Audio** Part 1 - TORTONA by Carmine Ciampa 246 780 views - music ☆ • □ □ [ VIENNE - speech by epagneule 431,412 views 6.50 Fanfan - Les valses de Vienne Maria Nazionale - Ciao, ciao - sounds seaofdeath · 45 videos 449,233 v flaavia78 **P** Subscribe (171 📫 1,453 🛛 🏺 2 028 676 views Moldovenii au talent 📹 Like About Share Add to մա - Pu **Textual information** Uploaded on Oct 5, 2010 Artist No description available. François Feldman Celine Dion Et Garou - Sous Le Vent (Live) mickaeldeconincl 25.972 views [www.youtube.com] 'Brindisi from La Traviata' Show more

#### **Visual features**

Histograms of Oriented Gradients



**Convolutional Neural Networks** 



#### Audio features



[Tzanetakis al., Ismir, 2011]

#### **Motion features**

3D-HoG / 3D-HOF



[J. Uijlings et al., IJMIR, 2014]

#### **VIDEO CLASSIFICATION** Fast automatic indexing of video databases; $t_{stop} = -2.5$ Random I = -5.3forests I = -3.9l = -2.4 > $v_{\mu,i} = \frac{1}{T\sqrt{P(x_i)}} \sum_{t=1}^{T} \frac{(x_t - \mu_i)}{\sigma_i}$ $v_{\sigma,i} = \frac{1}{T\sqrt{2P(x_i)}} \sum_{t=1}^{T} \left[ \frac{(x_t - \mu_i)^2}{\sigma_i^2} \right]$ **Nonlinear Extraction of multimodal** Generate a video **Compute VLAD SVM** dictionary descriptors features classifiers

[Mironică et al., MTAP, 2016]

#### Fuziunea trăsăturilor – "Late Fusion"



[Mironică et al., CBMI 2013, IEEE/ACM]

#### IVOM – dr.ing. Ionuț Mironică

## **TEXT CLASSIFICATION**

#### **Preprocessing steps**

- Remove the stop words
- Word stemming
   "swimmers", "swimming" -> swim
   "stems", "stemmer", "stemming", "stemmed" -> "stem".
- Extract the dictionary

a	been	get
about	before	getti
after	being	go
again	between	goes
age	but	goin
all	by	gon
almost	came	got
also	can	gotte
am	cannot	had
an	come	has
and	could	ha

english money library build loss address trigger corner manifest build loss pc layout pc software archieve investig Costserver security came to monitor RAM survey light situation message parkusername modul list procedure thelephone task salary lanreview letter virus almost submit page group wish computer file exista access

## [http://wiki.com]

## **TEXT CLASSIFICATION**

#### **Bag of Words**

(1) John likes to watch movies. Mary likes movies too.(2) John also likes to watch football games.

Bag of Words

BoW1 = {"John":1,"likes":2,"to":1,"watch":1,"movies":2,"Mary":1,"too":1}; BoW2 = {"John":1,"also":1,"likes":1,"to":1,"watch":1,"football":1,"games":1};

(1) [1, 2, 1, 1, 2, 1, 1, 0, 0, 0]
(2) [1, 1, 1, 1, 0, 0, 0, 1, 1, 1]

## **TEXT CLASSIFICATION**

**Term Frequency – Inverse Document Frequency** 

• A popular weighting schema that normalize the word frequency:

$$tfidf(w) = tf.\log(\frac{N}{df(w)})$$
• The word

• The word is more important if it appears several times in a target document

• The word is more important if it appears in less documents



## **TEXT CLASSIFICATION**

#### **Word Embeddings**





Male-Female

Verb tense

**Country-Capital** 

[GloVe (Pennington et al., 2014)]

## **STOCK MARKET PREDICTION**

- Values of stocks (open, close, max, min)
- Bid/Ask volume
- Bid/Ask volume misbalance
- Price volatility
- Aggregated parameters on the last hours / days / weeks
- Correlations
- News & Tweets & Facebook (sentiment analysis)



[http://wiki.com]

## **CREDIT DEFAULT PREDICTION**

- Credit history
- Purpose of the credit
- Savings account/bonds
- Personal status and sex
- Other debtors / guarantors
- Present residence since
- Age in years
- Housing (rent / own / for free)
- Job
- Residency



"I'd like to borrow enough money so that I'll be completely out of debt!"

[http://jantoo.com]

# **Feature engineering**

.....

## **CREDIT DATASET**

Status of existing account	Credit history	Country	Age	Salary	Telephone
0 <= < 200 EUR	no credits taken	Germany	22	1000	none
>= 2000 EUR	all credits paid back duly	France	38	500	none
0 <= < 200 EUR	delay in paying off in the past		28		yes
200 EUR <= < 2000 EUR		Thailand	55	100000	yes
>= 2000 EUR	no credits taken		55	1000	

# Feature engineering – missing data

Status of existing account	Credit history	Country	Age	Salary	Telephone
0 <= < 200 EUR	no credits taken	Germany	22	1000	none
>= 2000 EUR	all credits paid back duly	France	38	500	none
0 <= < 200 EUR	delay in paying off in the past	?	28	?	yes
200 EUR <= < 2000 EUR	?	Thailand	55	100000	yes
>= 2000 EUR	no credits taken	?	55	1000	?

# **Feature engineering – missing data**

Alternative 1: remove these lines

Alternative 2: dummy methods

Use aggregation operators for numeric categories:

• Average, Max, Min, Median value

**Categorical features:** 

- Most used value
- Other category

Alternative 3: A value estimated by another predictive model.



Label Encoder: It is used to transform non-numerical labels to numerical labels. Numerical labels are always between **0** and **n\_classes-1**.



**Dummy Encoder:** Dummy coding is a commonly used method for converting a categorical input variable into continuous variable. 'Dummy', as the name suggests is a duplicate variable which represents one level of a categorical variable. Presence of a level is represented by 1 and absence is represented by 0.

Credit history	no credits taken	all credits paid back duly	delay in paying off in the past
no credits taken	1	0	0
all credits paid back duly	 0	1	0
delay in paying off in the past	0	0	1
all credits paid back duly	0	1	0
no credits taken	1	0	0

Country
Germany
France
Australia
Thailand
US

Germany	France		US
1	0	250 columns	0
0	1		0
0	0		0
0	0		0
0	0		1



# Feature engineering – different scales

## **Feature normalization**

Age	Salary
22	1000
38	500
28	1000000
55	100000
55	1000

[18 - 80] [0 - 1000000]

Age	Salary
0.22	0.001
0.38	0.0005
0.28	1
0.55	0.1
0.55	0.001

#### Robust scaler

$$x' = \frac{x_i - Q_1(x)}{Q_3(x) - Q_1(x)}$$

Rescaling

$$x' = rac{x-\min(x)}{\max(x)-\min(x)}$$

# Mean normalization $x' = rac{x - ext{mean}(x)}{ ext{max}(x) - ext{min}(x)}$

#### **Standardization**

$$x'=rac{x-ar{x}}{\sigma}$$

Scaling to unit length

$$x'=rac{x}{||x||}$$

# Feature engineering – different scales

## **Feature normalization**









$$x' = rac{x-\min(x)}{\max(x)-\min(x)}$$



#### **Robust scaler**

$$x' = \frac{x_i - Q_1(x)}{Q_3(x) - Q_1(x)}$$



# Feature engineering – imbalanced data



# Feature engineering – imbalanced data

- Oversampling (may generate overfitting)
- Subsampling (may generate underfitting)
- SMOTE

... But what if there is a majority sample Nearby?

Majority sample



# Feature engineering – unnecessary features





Distribution of features values with target = 1 and target = 0

[G. Preda, I. Mironică AAD Conference 2018]

## Feature engineering – most important features





Distribution of features values with target = 1 and target = 0

[G. Preda, I. Mironică AAD Conference 2018]

## **Session 2**

. . . . . . . . . . . . . . . . . .

## CLASSICAL

## **MACHINE LEARNING ALGORITHMS**

## **Nearest neighbor**



K-NN computes the distance between the item that need to be classified to the most close objects from the database.

The item it will be classified according to the class of the neighbors.

# **Nearest neighbor**


## **Nearest neighbor**

#### **Parameters to consider**

• Number of neighbors

If *k* is too small, sensitive to noise points If *k* is too large, neighborhood may include points from other classes

• The way of computing the metric



#### **ADVANTAGES**

Small amount of data to train your model

Training is very fast (almost zero) **DISADVANTAGES** 

You need to keep in memory all the data

Slow: you need to compare your feature with all the stored database

Highly depends on the metric used

You cannot use it for complex algorithms

## **Linear classification**

#### **FISHER METHOD**

LDA assumes that the conditional probability density functions p(x|y=0) and p(x|y=1) are both normally distributed with mean and covariance.



### **Linear classification**

#### LINEAR PERCEPTRON METHOD

$$h_{\theta}(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 \qquad h(x) = \sum_{i=0}^n \theta_i x_i = \theta^T x,$$

Use Gradient Descent to update the weights:

- Choose some initial weights
- Define cost function

$$J(\theta) = \frac{1}{2} \sum_{i=1}^{m} (h_{\theta}(x^{(i)}) - y^{(i)})^2.$$

- Update weights until convergence

$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta).$$

$$\frac{\partial}{\partial \theta_j} J(\theta) = \frac{\partial}{\partial \theta_j} \frac{1}{2} \left( h_\theta(x) - y \right)^2$$
$$= 2 \cdot \frac{1}{2} \left( h_\theta(x) - y \right) \cdot \frac{\partial}{\partial \theta_j} (h_\theta(x) - y)$$
$$= \left( h_\theta(x) - y \right) \cdot \frac{\partial}{\partial \theta_j} \left( \sum_{i=0}^n \theta_i x_i - y \right)$$
$$= \left( h_\theta(x) - y \right) x_j$$

## **Logistic regression**

Logistic regression's output represents the probability of an independent variable to belong to a certain category

$$h_{\theta}(x) = g(\theta^{T}x) = \frac{1}{1 + e^{-\theta^{T}x}},$$

$$g(z) = \frac{1}{1 + e^{-z}} \quad g'(z) = \frac{d}{dz} \frac{1}{1 + e^{-z}} = \frac{1}{(1 + e^{-z})^{2}} (e^{-z}) = \frac{1}{(1 + e^{-z})^{2}} \cdot \left(1 - \frac{1}{(1 + e^{-z})}\right) = g(z)(1 - g(z)).$$

$$J(\theta) = -\left[\frac{1}{m}\sum_{i=1}^{m} y^{(i)}\log h_{\theta}(x^{(i)}) + (1-y^{(i)})\log(1-h_{\theta}(x^{(i)}))\right]$$

Use Gradient Descent to update the weights.



Logistics functions

## **Logistic regression**





SVM constructs a hyperplane that is used to separate different classes.



- The hyperplane is chosen where the distance between distinct categories is maximized.
- Hard-margin vs Soft-margin
- Function that needs to be minimized:

$$egin{bmatrix} rac{1}{n}\sum_{i=1}^n \max\left(0,1-y_i(ec{w}\cdotec{x}_i-b)
ight) \end{bmatrix}+\lambda \|ec{w}\|^2 \ ext{Hinge loss function} \ ext{Regularization} \end{cases}$$



$$\begin{bmatrix} \frac{1}{n} \sum_{i=1}^{n} \max\left(0, 1 - y_i(\vec{w} \cdot \vec{x}_i - b)\right) \end{bmatrix} + \lambda \|\vec{w}\|^2$$
  
Hinge loss function Regularization

Partial derivates:

$$\frac{\delta}{\delta w_k} \lambda \parallel w \parallel^2 = 2\lambda w_k$$

$$\frac{\delta}{\delta w_k} \left( 1 - y_i \langle x_i, w \rangle \right)_+ = \begin{cases} 0, & \text{if } y_i \langle x_i, w \rangle \ge 1\\ -y_i x_{ik}, & \text{else} \end{cases}$$

 $egin{aligned} & ext{Gradient Update} - ext{No misclassification} \ & w = w - lpha \cdot (2\lambda w) \ & ext{Gradient Update} - ext{Misclassification} \ & w = w + lpha \cdot (y_i \cdot x_i - 2\lambda w) \end{aligned}$ 



# **Nonlinear Support vector machines**

General idea: the original feature space can always be mapped to some higher-dimensional feature space where the training set is separable:



## **Nonlinear Support vector machines**



Map a multidimentional function  $\varphi(x) = (x, x^2)$ 



$$\varphi(x) \cdot \varphi(y) = (x, x^2) \cdot (y, y^2) = xy + x^2 y^2$$
$$K(x, y) = xy + x^2 y^2$$



#### **PARAMETERS TO CONSIDER**

- C The way of computing the hyperplane:
  - if it is two tight (C has a higher value) the model will overfit
  - if it is too general the model will underfit



Gamma and the kernel type

#### **ADVANTAGES**

Small amount of data to train your model

Training may be very fast on liner model

Can be used for complex algorithms

#### **DISADVANTAGES**

Training for nonlinear approaches can be slow

The data should be normalized

### **Decision trees**



**Training Data** 

## **Decision trees**





#### PARAMETERS TO CONSIDER

**Criterion** - the function to measure the quality of a split (Gini vs Entropy)

#### Max depth of the tree:

- high values the model will overfit
- small values the model will underfit



#### **ADVANTAGES**

Small amount of data to train your model

Training is very fast

Provide a dendogram that is able to show the importance of each feature

The data do not need to be normalized

#### DISADVANTAGES

It is very easy to overfit

# **Bagging algorithms - Random forests**



the

# **Bagging algorithms - Random forests**



### **Extremely Random forests**



# **Boosting algorithms - Gradient boosted trees**



In boosting – every tree is learning from other which in turn boosts the learning. Trees are learned sequentially (slow).

### Adaboost



Train a set of weak hypotheses: h1, ...., hT.

The combined hypothesis H is a **weighted** majority vote of the T weak hypotheses.

Each hypothesis  $h_t$  has a weight  $\alpha_t$ .

$$H(x) = sign(\sum_{t=1}^{T} \alpha_t h_t(x))$$





Training set: 10 points (represented by plus or minus)

Original Status: Equal Weights for all training samples

### Adaboost



#### Round 1:

Three "plus" points are not correctly classified;

They are given higher weights.





Three "minuses" points are not correctly classified; They are given higher weights.





#### Round 3:

One "minus" and two "plus" points are not correctly classified; They are given higher weights.



137



**Final Classifier:** integrate the three "weak" classifiers and obtain a final strong classifier.





### **Gradient boosted trees**



## **Gradient boosted trees vs Random forests**

- RF are much easier to tune than GBM
- RF are harder to overfit than GBM
- XGBoost is an optimized version of GBM (reduce the overfitting)



XGBoost has gained a lot of popularity recently and has been used in the most winning Kaggle competition models. It is powerful tool to have in your Data Science Toolbox.

XGBoost can work with Trees as well as Linear Models. I recommend reading the XGBoost documentation for further parameter tuning options.

XGBoost is a recent, most preferred and powerful gradient boosting method. Instead of making hard Yes and No Decision at the Leaf Nodes, XGBoost assigns positive and negative values to every decision made. All Trees are weak learners and provide a decisions slightly better than a random guess. But collectively averaged out, XGBoost performs really well.

## **Gradient boosted trees and Random forests**

#### **PARAMETERS TO CONSIDER**

**Number of trees** – high number of trees – the model will have better results / the training classification will last longer

**Criterion** - the function to measure the quality of a split (Gini vs Entropy)

#### Max depth of the tree:

- high values the model will overfit
- small values the model will underfit

#### Learning rate (on Gradient Boosting algorithms):

Learning rate parameter shrinks the contribution of each tree. Lowering the value of learning rate increases the number of trees in the ensemble. Be vary that increasing the number of estimators to a large value may overfit the model.



Naive Bayes classifiers are a family of simple probabilistic classifiers based on applying Bayes' theorem with strong (naive) independence assumptions between the features.





Detect the sales of ice-cream by using the weather forecasting [Image from www.criminisi.com]




## **Comparison between different methods**



## **Comparison between different methods**



# **Deep Learning vs classical Machine Learning**

#### DO WE NEED HUNDREDS OF CLASSIFIERS TO SOLVE REAL WORLD CLASSIFICATION PROBLEMS?

- Performed experiments with 179 classifiers arising from 17 families on 121 datasets (from UCI repository);
- They concluded that the classifiers most likely to be the bests are the random forest and SVM

### **Deep Learning**

....

# **Deep Learning vs classical Machine Learning**

#### SUPERVISED LEARNING



classical Machine Learning?

# **Deep Learning vs classical Machine Learning**





[Images from: https://appliedgo.net/perceptron, wiki.org and quora.com]

#### Deep neural network





- Good to find the best separating plane between two classes.
- Complicated structure, with many parameters and several hyper-parameters, non – trivial to tune.
- Prone to overfitting.

[Image from: http://quora.com]

### **Neural networks history**











#### TWO LAYERS OF NEURONS



#### **ACTIVATION FUNCTIONS**

**Sigmoid Function** 

**Tanh Function** 



$$A = \frac{1}{1+e^{-x}}$$

- Sigmoid outputs are not zero-centered
- Sigmoids saturate and kill gradients.

$$f(x) = tanh(x) = \frac{2}{1+e^{-2x}} - 1$$

### **ACTIVATION FUNCTIONS**

#### **RELU Function**

- has accelerated convergence
- it does not saturate
- involve inexpensive operations (a simple threshold)
- ReLU units can die during the training (the neuron will never activate)
  more sensible to learning rate parameter



 $f(x)=x^+=\max(0,x)$  ,

#### **ACTIVATION FUNCTIONS**



Some people report success with this form of activation function, but the results are not always consistent.

#### **ACTIVATION FUNCTIONS**



#### ACTIVATION FUNCTIONS

#### Softmax

The idea of softmax is to define a new type of output layer for our neural networks  $z_{\mu}^{L}$ 

$$a_j^L = rac{e^{z_j^L}}{\sum_k e^{z_k^L}},$$

**SWISS KNIFE** "What neuron type should I use?"

- Use the ReLU non-linearity, be careful with your learning rates and possibly monitor the fraction of "dead" units in a network.
- If this concerns you, give Leaky ReLU or Maxout a try. Never use sigmoid. Try tanh, but expect it to work worse than ReLU/Maxout.
- Use Softmax for the final layer

### L1 / L2 NORMALIZATION

• Penalize large weights, and tend to make the network to prefer small weights.

L1 normalization 
$$S = \sum_{i=1}^{n} |y_i - f(x_i)|.$$

L2 normalization 
$$S = \sum_{i=1}^{n} (y_i - f(x_i))^2$$

#### DROPOUT

- Penalize the neurons.
- Unlike L1 and L2 regularization, dropout doesn't rely on modifying the cost function. Instead, in dropout we modify the network itself.





#### DROPOUT

- When teams up, if everyone expect the partner will do the work, nothing will be done finally.
- However, if you know your partner will dropout, you will do better.
- When testing, no one dropout actually, so obtaining good results eventually.
- You will need more neurons to train your model.

### **BATCH NORMALIZATION**

• Normalize distribution of each input feature in each layer across each minibatch to N(0, 1) learn the scale and shift

Input: Values of x over a mini-batch: $\mathcal{B} = \{x_{1m}\};$ Parameters to be learned: $\gamma, \beta$ Output: $\{y_i = BN_{\gamma,\beta}(x_i)\}$	
$\mu_{\mathcal{B}} \leftarrow \frac{1}{m} \sum_{i=1}^{m} x_i$	// mini-batch mean
$\sigma_{\mathcal{B}}^2 \leftarrow \frac{1}{m} \sum_{i=1}^m (x_i - \mu_{\mathcal{B}})^2$	// mini-batch variance
$\widehat{x}_i \leftarrow \frac{x_i - \mu_{\mathcal{B}}}{\sqrt{\sigma_{\mathcal{B}}^2 + \epsilon}}$	// normalize
$y_i \leftarrow \gamma \widehat{x}_i + \beta \equiv BN_{\gamma,\beta}(x_i)$	// scale and shift

Algorithm 1: Batch Normalizing Transform, applied to activation *x* over a mini-batch.

[loffe et al.]

#### **DROPOUT ENSEMBLE**



Train a bunch of networks with different structures

#### **DROPOUT ENSEMBLE**



### **ARTIFICIALLY EXPANDING THE TRAINING DATA**















[image source: http://wiki.com]

#### **ARTIFICIALLY EXPANDING THE TRAINING DATA**



#### **ARTIFICIALLY EXPANDING THE TRAINING DATA**





[Goodfellow et al.]

#### **Gradient Descent**

- To find argminx f(x):
- Start with x0
- For t=1....

xt+1 = xt + alpha f'(x t)
 where alpha is smaller than 1

- Simple and often quite effective on ML tasks
- Often very scalable
- Only applies to smooth functions (differentiable)
- Might find a local minimum, rather than a global one





How to solve this problem?

**LEARNING RATE** 

175



**Gradient Descent** 

**Stochastic Gradient Descent** 

**Batch Gradient Descent** 

Momentum / RMSProp / Adam etc



#### Momentum

Instead of using only the gradient of the current step to guide the search, momentum also accumulates the gradient of the past steps to determine the direction to go:

$$egin{aligned} v_{dw} &= eta \cdot v_{dw} + (1-eta) \cdot dw \ W &= W - lpha \cdot v_{dw} \end{aligned}$$



#### **RMSProp**

RMSprop is developed by Prof. Geoffrey Hinton in his neural nets class.

Instead of letting all of the gradients accumulate for momentum, it only accumulates gradients in a fixed window.

$$v_{dw} = eta \cdot v_{dw} + (1-eta) \cdot dw^2$$

$$W = W - lpha \cdot rac{dw}{\sqrt{v_{dw}} + \epsilon} \, .$$

#### Adam

Adam combines RMSProp with Momentum.

$$egin{aligned} & v_{dw} = eta \cdot v_{dw} + (1-eta) \cdot dw \ & v_{dw} = eta \cdot v_{dw} + (1-eta) \cdot dw^2 \ & W = W - lpha \cdot rac{dw}{\sqrt{v_{dw}} + \epsilon} \; v_{dw}_{_{ ext{M}}} \end{aligned}$$

## **Training Neural Architectures**

#### **FEED FORWARD**

Feed Forward (FF)



#### Deep Feed Forward (DFF)



Why we are going to deeper architectures?
#### **FEED FORWARD**

Structure	Types of Decision Regions	Exclusive-OR Problem	
Single-Layer Q	Half Plane Bounded By	A B	
	Hyper plane	B A	
Two-Layer Q	Convex Open Or	A B	
	Closed Regions	BA	
Three-Layer	Arbitrary (Complexity Limited by No.	A B	
	of Nodes)	BA	

### **FEED FORWARD**

Ex: applying deep learning to image recognition



#### **FEED FORWARD**

Little or no invariance to shifting, scaling, and other forms of distortion



### **FEED FORWARD**

Little or no invariance to shifting, scaling, and other forms of distortion







• scaling, and other forms of distortion

### **FEED FORWARD**

- the topology of the input data is completely ignored
- work with raw data.





#### **GENERAL SCHEMA**



[image source: http://wiki.com]

[Zeiler and Fergus 2013]

### CONVOLUTION





Local image neighborhood mask

Modified image data

*New Value* = 3 + 1 + 4 = 8

Other parameters: Stride & Padding



Image: 5x5x3
Padding: 1
Stride: 2
Number of convolutional kernels: 2
Size of convolutional kernels (3x3)

Image: 3x3x6



Filter W1 $(3x3x3)$				
-1	-1	-1		
1	1	-1		
1	1	0		
w1 (	:,:	:,1]		
1	-1	0		
1	0	0		
-1	0	0		
w1[:,:,2]				
-1	1	0		
-1	1	1		
1	1	-1		

Output Volume (3x3x2)

<u>o[:</u>,:,0]

3

2

0[:,:,1]

5

-3 3 -3

3

3 5

4

6

0 5

Bias b1 (1x1x1) b1[:,:,0] 0

[animation from http://cs231n.github.io/convolutional-networks]







Filter W1 (3x3x3)					
w1[:,:,0]					
0	1	1			
-1	1	1			
1	1	0			
w1[	:,:	,1]			
1	1	0			
0	1	0			
-1	0	-1			
w1[:,:,2]					
1	0	1			
-1	1	1			
-1	1	1			

Output Volume (3x3x2)					
o[:,:,0]					
0	-3	-4			
5	Δ	7			
-5	-4	-/			
-1	-8	-5			
0[:	<i>,</i> :,	ŢŢ			
6	9	3			
0	0	2			
9	8	3			
11	7	1			
••	1	1			

Bias b1 (1x1x1) b1[:,:,0] 0



Output Volume (3x3x2)				
o[:,:,0]				
2	1	0		
2	5	2		
3	-5	-6		
<u>0[</u> :	,:,	1]		
3	-5	-5		
-1	-6	-4		
-3	-6	-5		



Out	put V	Volu	me (3x3x2)		
0[:	,:,	0]			
3	3	5			
7	3	1			
0	2	4			
o[:,:,1]					
2	5	6			
0	0	5			

-3 3 -3







### SUBSAMPLING LAYER

Feature map



	Sing	gle d	epth	slice	
×	1	1	2	4	may pool with 2x2 filtors
	5	6	7	8	and stride 2
	3	2	1	0	
	1	2	3	4	
					•

This reduces the number of inputs to the next layer of feature extraction, thus allowing us to have many more different feature maps.

8

4

6

3

- reduce the effect of noises and shift or distortion
- reduce the overfitting

### **FLATTEN LAYER**



Flattening is the process of converting all the resultant 2 dimensional arrays into a single long continuous linear vector.

### **CONCLUSIONS ON PARAMETERS**

#### Size of the convolutional kernels

- Smaller filters collect as much local information as possible, bigger filters represent more global, high-level and representative information.
- It is a common standard to use odd numbers

#### Padding

 Padding is generally used to add columns and rows of zeroes to keep the spatial sizes constant after convolution, doing this might improve performance as it retains the information at the borders

### **CONCLUSIONS ON PARAMETERS**

Stride

• It is used to throw away the duplicate information.

#### Number of channels

- It is the equal to the number of color channels for the input but in later stages is equal to the number of filters we use for the convolution operation.
- The more the number of channels, more the number of filters used, more are the features learnt, and more is the chances to over-fit and vice-versa.

### **CONCLUSIONS ON PARAMETERS**

- Always start by using smaller filters is to collect as much local information as possible, and then gradually increase the filter width to reduce the generated feature space width to represent more global, high-level and representative information
- The number of channels should be low in the beginning such that it detects lowlevel features which are combined to form many complex shapes (by increasing the number of channels) which help distinguish between classes
- Keep adding layers until you over-fit.
- Always use classic networks like LeNet, AlexNet, VGG-16, VGG-19 etc

- **LENET 5** C1, C3, C5: Convolutional layer.
  - $5 \times 5$ : Convolution matrix.
  - S2, S4: Subsampling layer.
  - Subsampling by factor 2.
  - F6: Fully connected layer.



#### **LENET 5**



- About 187,000 connections
- About 14,000 trainable weights

#### **LENET 5**







#### **ALEX NET**



- It consisted 11x11 5x5,3x3 convolutions
- max pooling
- dropout
- data augmentation
- ReLU activations
- SGD with momentum

AlexNet architecture (May look weird because there are two different "streams". This is because the training process was so computationally expensive that they had to split the training onto 2 GPUs)

# Introduced the concept of dropout and ReLU Winner on LSVRC2012

[http://vision.stanford.edu/teaching/cs231b\_spring1415/slides/alexnet\_t ugce\_kyunghee.pdf]



### **GOOGLE LENET (INCEPTION ARCHITECTURE)**



(a) Inception module, naïve version

#### Added more layers - 2015

### **GOOGLE LENET (INCEPTION ARCHITECTURE) – V2/3 IMPROVEMENTS**



**Factorize 5x5** convolution **to two 3x3** convolution operations to improve computational speed. Although this may seem counterintuitive, a 5x5 convolution is **2.78 times more expensive** than a 3x3 convolution.



### **GOOGLE LENET (INCEPTION ARCHITECTURE)**



Added more layers - 2015



2016

[https://towardsdatascience.com]

214

### RESNET

- Full ResNet architecture:
- Stack residual blocks
- Every residual block has two 3x3 conv layers
- Periodically, double # of filters and downsample F(x) spatially using stride 2 (/2 in each dimension)
- Additional conv layer at the beginning
- No FC layers at the end (only FC 1000 to output classes)



[He et al. 2015]

#### DENSENET

The idea here is that if connecting a skip connection from the previous layer improves performance, why not connect every layer to every other layer?
## **Training Neural architectures**

#### **CAPSULE NETWORKS**







[https://towardsdatascience.com/animated-rnn-lstm-and-gru-ef124d06cf45]



**PROBLEM OF LONG-TERM DEPENDENCIES** 

E.g. 1: The grass is green

E.g. 2: I am Romanian. I love ... (200 words)... My mother tongue is Romanian.

LONG-SHORT TERM MEMORY (LSTM) / GRU





[https://towardsdatascience.com/animated-rnn-lstm-and-gru-ef124d06cf45]

LSTM



https://medium.com/mlreview/understanding-lstm-and-itsdiagrams-37e2f46f1714



[https://towardsdatascience.com/animated-rnn-lstm-and-gru-ef124d06cf45]

*Transfer Learning (TL):* The ability of a system to recognize and apply knowledge and skills learned in previous tasks to novel tasks (in new domains)

**ASSUMPTION:** the source domain and target domain data use exactly the same features and labels.

**MOTIVATION:** Although the source domain data can not be reused directly, there are some parts of the data that can still be reused by re-weighting.



#### **VERSION 1**



#### **VERSION 2**



Replace the last layers and retrain only the last layers of the network

#### **VERSION 3**



Replace the last layers and retrain all the network

# IMAGE CLASSIFICATION, OBJECT DETECTION, LOCALIZATION, ACTION RECOGNITION, SCENE UNDERSTANDING



#### **PEDESTRIAN DETECTION, TRAFFIC SIGN RECOGNITION**



# IMAGE CLASSIFICATION, OBJECT DETECTION, LOCALIZATION, ACTION RECOGNITION, SCENE UNDERSTANDING



Sliding windows 2 Window size bigger than 1 [Image source: https://towardsdatascience.com/]

#### **IMAGE CLASSIFICATION, OBJECT DETECTION, LOCALIZATION**





Figure 2: Two examples of our selective search showing the necessity of different scales scales. On the right we necessarily find the objects at different scales as the girl is contain

#### **R-CNN:** Regions with CNN features



[Images source: https://blog.athelas.com/a-brief-history-of-cnns-in-image-segmentation-from-r-cnn-to-mask-r-cnn-34ea83205de4/]

# IMAGE CLASSIFICATION, OBJECT DETECTION, LOCALIZATION, ACTION RECOGNITION, SCENE UNDERSTANDING







Yolo: You Only Look Once

#### **DEEPMIND'S ALPHAGO**



#### **ROAD FINDER**





[Vlad Mnih, ICML 2012]

#### **ARTIFICIALLY EXPANDING THE TRAINING DATA**





[Goodfellow et al.]



[Image source O'Reilly]

#### **COMMON PROBLEMS WITH GENERATIVE ADVERSIAL NETWORKS**

- The loss of the discriminator is going very rapidly to zero and the generator is not able to fool the discriminator
  - Take a look on the gradients of the generator vs discriminator (at the begining the discriminator should learn faster)
- Mode Collapse refers to the scenario when the Generator produces the same (or almost same) images every time and is able to successfully fool the discriminator.

#### **GENERATIVE ADVERSIAL NETWORKS TIPS AND TRICKS**

#### Normalize the inputs

- normalize the images between -1 and 1
- Tanh as the last layer of the generator output

#### **Generative sampling**

- don't use uniform distribution
- sample from a gaussian distribution





#### **GENERATIVE ADVERSIAL NETWORKS TIPS AND TRICKS**

• Construct different mini-batches for real and fake, i.e. each mini-batch needs to contain only all real images or all generated images.

• When batchnorm is not an option use instance normalization (for each sample, subtract mean and divide by standard deviation).



#### **GENERATIVE ADVERSIAL NETWORKS TIPS AND TRICKS**

- the stability of the GAN game suffers if you have sparse gradients
- LeakyReLU = good (in both G and D)
- for Downsampling, use: Average Pooling, Conv2d + stride
- for Upsampling, use: PixelShuffle, ConvTranspose2d + stride
- avoid Sparse Gradients: ReLU, MaxPool



#### **GENERATIVE ADVERSIAL NETWORKS TIPS AND TRICKS**

#### Use Soft and Noisy Labels

• *Label Smoothing*, i.e. if you have two target labels: Real=1 and Fake=0, then for each incoming sample, if it is real, then replace the label with a random number between 0.7 and 1.2, and if it is a fake sample, replace it with 0.0 and 0.3 (for example).

• Salimans et. al. 2016

• make the labels the noisy for the discriminator: occasionally flip the labels when training the discriminator

#### **GENERATIVE ADVERSIAL NETWORKS TIPS AND TRICKS**

- Use Adam optimizer
  - See Radford et. al. 2015
- Use SGD for discriminator and ADAM for generator
- Add some artificial noise to inputs to D (Arjovsky et. al., Huszar, 2016)
  - <u>http://www.inference.vc/instance-noise-a-trick-for-stabilising-gan-training/</u>
  - <u>https://openreview.net/forum?id=Hk4\_qw5xe</u>
- Adding gaussian noise to every layer of generator (Zhao et. al. EBGAN)
  - Improved GANs: OpenAI code also has it (commented out)



#### **GAN - MUSIC GENERATION**

http://www.asimovinstitute.org/analyzing-deep-learning-tools-music/

#### **GAN – TEXT GENERATION**

#### **AUTOMATIC TEXT GENERATION**

Given "my name is ivan and i live in stockholm" generated:

my name is ivan and i live in stockholm and i can do it and i wanna get closer so call me all alone i said its all i act so good and i wanna know the time i could stay in the world, youre always the only one that i wanna be...

https://towardsdatascience.com/yet-another-textgeneration-project-5cfb59b26255 [I, Donald J. Trump, president of the United ]States. Were going to have you to say what happens to try to make America great again. We're going to tell you that we have to win. We owe \$100.

[Despite the constant negative press cov]erage. I was going to be great. But if they're going to stop the good guy who has a big night.

[Despite the constant negative press cov]erage. Looking forward to it. The fundraisers are totally promised by the Central Park South...

[Despite the constant negative press cov]erage lines to Trump Tower. That's not the relationships. It doesn't have to be a disaster to Iraq.

[Despite the constant negative press cov]erage car millions of people and spot. We are going to come back. We have to go to golf.

#### HARRY POTTER: WRITTEN BY ARTIFICIAL INTELLIGENCE

"The Malfoys!" said Hermione.

Harry was watching him. He looked like Madame Maxime. When she strode up the wrong staircase to visit himself.

"I'm afraid I've definitely been suspended from power, no chance — indeed?" said Snape. He put his head back behind them and read groups as they crossed a corner and fluttered down onto their ink lamp, and picked up his spoon. The doorbell rang. It was a lot cleaner down in London.

Hermione yelled. The party must be thrown by Krum, of course.

Harry collected fingers once more, with Malfoy. "Why, didn't she never tell me. ..." She vanished. And then, Ron, Harry noticed, was nearly right.

"Now, be off," said Sirius, "I can't trace a new voice."

He punished Uncle Vernon, so loudly that she could barely use it with Viktor Krum, these faces in one side of their previous year. Mr. Dursley again was this champion. He was about to give the Remembrall there.... every time he .....

#### **AUTOMATIC COLORIZATION OF BLACK AND WHITE IMAGES**



#### **AUTO-ENCONDERS**



250

#### **AUTO-ENCONDERS**



#### **AUTO-ENCONDERS – IMAGE SEGMENTATIONS**



[image from https://www.frontiersin.org]
### **Applications of Neural Networks**

#### **CODE GENERATION**



Original

Prediction

https://blog.floydhub.com/turning-design-mockups-into-code-with-deep-learning

### **Applications of Neural Networks**

#### **OSCAR PREDICTION 2018**



**Big ML** was able to predict all the winners on 2018 Oscars using:



- details about the film (duration, budget, genre, etc.),
- IMDB ratings
- nominations in previous awards (Golden Globes, BAFTA, Screen Actors Guild and Critics Choice) from 2000 to 2017 used in last year's predictions.

### Let's make things happen!

### Languages in Machine Learning







### **Machine learning frameworks**

#### **PYTHON LIBRARIES**









### **Machine learning frameworks**



NLTK TextBlob



Standford NLP





### **Deep learning frameworks**



### **Cloud platforms**







#### **Google** Cloud Platform

### **GPU programming**

GPU may provide more than 10x performance and 5x energy efficiency





#### GPU Accelerator Optimized for Parallel Tasks



#### **TESLA K80**

**Specs:** 24GB GDDR5 CUDA 4992 processors

Launch price: \$5000



#### [images and content from www.nvidia.com]

#### <u>GTX 1080</u>

#### Specs

VRAM: 8 GB

Memory bandwidth: 320 GBs/second Processing power: 2560 cores @ 1733 MHz



#### GTX 1080 TI Specs VRAM: 11 GB Memory bandwidth: 484 GBs/sec Processing power: 3584 cores @ 1582 MHz



[images from www.nvidia.com]

#### GTX 1060 / 1050 TI (6 / 4 GB)

#### **Specs**

Processing power: 1280 cores @ 1708 MHz (~ 2,19 M CUDA Core Clocks)



- I work with data sets > 250GB: K80 or multiple GTX 1080 Ti.
- I have little money: GTX 1060 (6GB).
- I have almost no money: GTX 1050 Ti (4GB).
- I do Kaggle: GTX 1060 (6GB) for any "normal" competition, or GTX 1080 Ti for "deep learning competitions".
- I am a competitive computer vision researcher: NVIDIA Titan Xp or at least few GTX 1080 Ti.



## **Udemy COURSERC U D A C I T Y**

# kaggle Medium



### root:~\$Siraj Raval

Artificial Intelligence Education New Video Every Week!

### Measuring performance in Machine Learning

#### **TWO CLASSES**

	Actual			
Predicted		TRUE	FALSE	
	TRUE	True positive (TP)	False pozitive (FP)	
	FALSE	False negative (FN)	True negative (TN)	

Precision = 
$$\frac{TP}{TP+FP}$$
 Recall =  $\frac{TP}{TP+FN}$ 

Are these parameters OK for every measure?

What is more important: precision or recall?

Accuracy =  $\frac{TP+TN}{TP+TN+FN+FP}$ 

#### **CONFUSION MATRIX**



#### Ex: fraud detection



Precision = 3%

Mean = 
$$(P+R)/2 = 51.5\%$$

F1 score =  $\frac{2*Precision*Recall}{Precision+Recall}$ 

#### **ROC CURVES**

ROC Curves are used to see how well your classifier can separate positive and negative examples and to identify the best threshold for separating them.

To be able to use the ROC curve, your classifier has to be *ranking* - that is, it should be able to rank examples such that the ones with higher rank are more likely to be positive.



[Image from habrastorage.org]

#### **PRECISION / RECALL CURVES**

- If we retrieve more document, we improve recall (if return all docs, R=1)
- if we retrieve fewer documents, we improve precision, but reduce recall
- so there's a trade-off between them

 The are under the curve is also known as Mean Average Precision



#### **MODEL EVALUATION TECHNIQUES**



- Training set consists of records with known class labels used to build a classification model
- Evaluation set is used to measure the performance of the system
- A labeled test set of previously unseen data records is used to evaluate the quality of the model.

#### **K-FOLDS (CROSS VALIDATION)**



### **Combining more algorithms together**

- Ensemble methods are based around the hypothesis that an aggregated decision from multiple experts can be superior to a decision from a single system.
- More accurate ONLY if the individual classifiers disagree.

### **Combining more algorithms together**

#### **Early Fusion**



### **Combining more algorithms together**

#### **Late Fusion**



#### Hard voting vs Soft Voting

### **Ethics in Machine Learning**



- Big Data Increases Inequality and Threatens Democracy, pointed out that predictive analytics based on algorithms tend to punish the poor, using algorithmic hiring practices, insurance and credit risk analysis, bespoke offers and advertisement.
- Prevent using discriminatory information: sex, religion, race, ideology
- How many jobs will disappear because of AI?
- A Hippocratic oath for AI developers?
- Is the accuracy principle equal to equity principle?

#### Conclusion



[image from http://www.quotehd.com]

### Conclusion

• No free lunch: machine learning algorithms are tools with advantages and disadvantages.

• First try the simple classifiers. If they don't map well on the application then other more complicated models can be used.

• It is better to have simple classifiers and more intelligent data than complex classifiers and simple data.

• Use complex classifiers when we have diverse and/or large data set.



## PART 2

## Hands-on time

## Exercise 1:

GitHub

https://github.com/imironica/keras-without-a-phd

282

### **Demo details**

#### **MNIST DIGIT DATASET**

- The MNIST database of handwritten digits, available from NIST
- Training set: 60,000 image of digits of size 28x28



• Test set: 10,000 images

### **Demo details**

#### **MNIST DIGIT DATASET** a г $\bigcirc$ Ъ $\mathcal{O}$ Ó Ô I Ζ г $\mathcal{O}$ $\bigcirc$ г a $\odot$ a $\mathcal{O}$ Ø Ó а ð а O Ø C 0: D a

[more details on http://yann.lecun.com/exdb/mnist/]

## Hands-on time

User / password: deeplearning

cd /home/deeplearning/work/keras-without-a-phd

sudo mkdir keras-without-a-phd cd keras-without-a-phd git clone <u>https://github.com/imironica/keras-without-a-phd.git</u>

sudo python3 0.1\_visualizeDataset.py -v 1

## Hands-on time



GitHub

https://github.com/imironica/traffic-signs-keras

#### **Demo details**

.....

#### **TRAFFIC SIGN RECOGNITION**



#### **Demo details**

#### **TRAFFIC SIGN RECOGNITION**






## **THANK YOU**

Ionuț Mironică







http://ionut.mironica.ro CitHub imironica

**G** ionut.mironica