## Introduction to Machine Learning

### Augustin Cosse.



Summer 2021

May 24, 2021

▲□▶ ▲□▶ ▲ 三▶ ▲ 三▶ 三三 - のへぐ

## Schedule

- Class and labs: Monday/Tuesday, 3pm-6.30pm + Thursday 3-5pm (CEST).
- Recitation (Mandatory): Probably Tuesday 4pm 6.30pm.
- Office hour : Thursday 5 6pm (CEST), to be confirmed.
- Location: Zoom !
- Combination between programming sessions (python) and lectures

・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・

- Final Exam: Midterm: 30%, Final : 30%
- Assignements throughout the semester: 30%
- Independent project: 10%

## Course organization

- Notes + Sample exam questions can be found on the course webpage
- Sample exam questions = help you with the study but not comprehensive
- If a section of the notes is not covered in class, you don't have to study it for the exam

・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・

• See http://www.augustincosse.com/teaching for details (Scroll down and select "NYU Paris ML Summer School")

#### Introduction to Machine Learning (NYU Paris, Spring 2020)



Tentative schedule:

Legend: Lab sessions are in green, Homeworks are in red (right side of the table),

#### dates related to the project are in orange.

Week #	date	Торіс
Week 1	02/03, 02/07	General Intro + reminders on proba and inference. Part I, Part II Lab 1
		Part I : supervised Learning

## **Reference Books**

**Springer Series in Statistics** 

Trevor Hastie Robert Tibshirani Jerome Friedman

#### The Elements of Statistical Learning

Data Mining, Inference, and Prediction

Second Edition

Deringer





# My objective: give you the tools to start working on your own ideas later



## And connect you to the ML sphere in Paris

ΜΛΤΗ & ΙΛ

**LUNDI 9 MARS 2020** DE 9H À 18H

#### CONFÉRENCES DE

Francis Bach (Inria) Emmanuel Candès (Stanford University) Igor Carron (LightOn) Aurélie Jean (In Silico Veritas) Michael Jordan (Berkeley) Kathryn Hess (EPFL) Yan Le Cur (Facebook)

Table ronde animée par Stéphane Mallat (Collège de France)

## Artificial Intelligence and the Future / Demain, l'intelligence artificielle

Conférence du 22 novembre 2018

#### Demis Hassabis



HACKATHON – 2EME INTELLIGENCE ARTIF

hisis

En partenariat avec

**N**I▲

SESAM

## En partenariat avec

28-29 février 2020 Université Paris-Dauphine

old Dist Till Chief.

Etudiants en Finance et IT, jeunes diplômés

Inscription & informations : www.qminitiative.org

9 & 10 mars 2020

Palais des Congrès • Paris

## Class + programming sessions

- The objective is to provide you with all the theory/material needed to tackle some practical problems
- Then you will be given the opportunity to deal with those problems through Jupyter notebooks
- Many exercises rely on Scikit learn (free software machine learning library for Python, see https://scikit-learn.org/stable/) which features multiple learning algorithms.

・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・

## General material (theory), see webpage

- Reminders in Stats/Probability, Inference.
- Supervised learning
  - Classification (Logistic regression, LDA)
  - Regression (Linear regression, regularization)
  - Neural Networks (+ implementation/Keras/TensorFlow)
  - Support Vector Machines
  - ...
- Unsupervised learning
  - Clustering (K-means, K-medoid, EM)
  - PCA, ICA
  - Manifold Learning
- Advanced topics (Adversarial and Reinforcement Learning)

What type of data are we going to use

- Today data can be found everywhere
- Some major sources of data available online include
  - Kaggle (https://www.kaggle.com/)
  - UCI Machine Learning Repository (https://archive.ics.uci.edu/ml/datasets.html)
  - ENS Challenge Data (https://challengedata.ens.fr)
  - Many possibilities to get data directly from python (e.g. pandas, yahoo finance, scikit-learn...)

・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・



## How do the programming sessions work?

- Two key components : Github and Jupyter notebooks
- Most efficient way to teach ML today = notebooks
  - Part of the code is given, remaining part has to be provided by the students
  - Many online examples that can be used directly or modified to generate new programming exercises
- When dealing with Neural nets (including GANs) so far I use TensorFlow/Keras (high-level API to build and train deep learning models)
- Exercises will be given through Github, see https://github.com/acosse/Intro2ML\_Summer2021 (will be updated soon, Please fork the repo)

#### 2.3 Convolutional Neural Networks for image classification

(adapted from R.Hon)

In this example we will learn how to train and use convolutional neural network for image recognition. Start by logging on Kaggle and download the MNIST Data : https://www.kaggle.com/c/digit-recognizer/data

Then store the content of the csv files using pandas in a train and a test variables

```
In [ ]: import numpy as np # linear algebra
    import pandas as pd # data processing, CSV file I/O (e.g. pd.read csv)
```

```
# put your code here
```

In this exercice, we will use Keras, which is a high level library which runs on top of Tensor Flow and enable fast experimentation with deep neural networks

⇔ Code ① Issues 0 11 Pull requests 0

In [ ]: from keras.models import Sequential from keras.layers import Dense, Dropout, Flatten from keras.layers.convolutional import Conv2D, acosse / IntroMLFall2018 from keras.utils import np\_utils from keras.optimizers import RMSprop from keras.callbacks import ReduceLROnPlateau from keras.preprocessing.image import ImageDat import matplotlib.pvplot as plt website for the Fall 2018 NYU Class Introduction to Machine Learning from sklearn.model selection import train test Manage topics

Using the tools form the pandas library, look at your data, W the first column of each pandas data array contain ?

#### 2.3.1 Displaying the images

Display the first few images from the MNIST dataset using im-

In [ ]: # put your code here



(e) 168 commits O releases # 1 environment \$\$ 1 contributor Branch: master - New pull request Create new file Upload files Find file Clone or download acosse Add files via upload Latest commit 6ba19f4 on 6 Dec 2018 III Assignement2 Add files via upload 2 months ago III Assignements 3 months ago ExamplesNotes 6 months ago III Exercise1 Rename Exercice1.jpvnb to Exercise1.jpvnb 2 months ago Exercise2 Renaming Ex2 2 months ago Exercise3 Add files via -----ProjectsProposals Add files III Readings Add files in TestLab GitHub README md Set theme \_config.yml Update \_c

III Projects 0 III Wiki



## Notebook

## A group project

- A key aspect of the course is the project
- The idea of the project is to give you the opportunity to develop your own ideas (insofar as possible)
- A suggestion of possible topics can be found online. However you are free to come up with new ideas

• View this as a first step towards a startup. If the project works well, you can perhaps extend it later !

## A group project

- You choose one subject in group and start working on it throughout the semester
- The project has to (1) be related to at least one algorithm studied in class and (2) has to exhibit some programming aspects.
- You are free to come during OH with questions on the project
- You then get to present your work at the end of the semester

・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・

## A group project

Some examples from last year:

- Training Tic Tac Toe
- Training a neural network to win at a pong game
- Using Natural Language processing to detect and remove aggressive comments from Twitter
- Detecting health conditions from medical data
- Use generative adversarial networks to generate cartoons

・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・

• ...



◆□▶ ◆□▶ ◆ ≧▶ ◆ ≧▶ ○ ≧ ○ � � �

## Table of contents

▲□▶ ▲圖▶ ▲ 臣▶ ▲ 臣▶ ● 臣 ● 9 Q @

## Machine learning today



▲ロト▲御ト▲臣ト▲臣ト 臣 めるの

## Machine learning today



## BBC NEWS

# Artificial intelligence-created medicine to be used on humans for first time

By Jane Wakefield Technology reporter

# This Al tool helps identifies breast cancer with 90% accuracy rate

The drug was much quicker to market than ones developed in more traditi

A drug molecule "invented" by artificial intellig human trials in a world first for machine learni



The AI tool identifies malignant lesions (red) and benign lesions (green).

Image: NYU School of Medicine

▲□▶ ▲□▶ ▲□▶ ▲□▶ □ のQで



# Autonomous truck carts 40,000 lb of butter coast to coast in the US

## Bloomberg Businessweek

Tesla's Autopilot Could Save the Lives of Millions, But It Will Kill Some People First



UDACITY



#### **Artificial Intelligence**

### **Machine Learning**

#### **Deep Learning**

Multilayered (deep) Neural Networks + vast amount of data General set of algorithms enabling machines to improve performance when being exposed to more and more data

< D >

Programs able to learn and reason, mimicking human intelligence

## Some achievements

Artificial Intelligence

**IBM DeepBlue** 

### **Machine Learning**

Netflix recommendation

**Deep Learning** 

**Google RankBrain** 

Gmail smart reply Tesia autopilot Skype Translator

**Facebook photo tagging** 

commendation

Email spam filter

Google pageRank

Facebook NewsFeed

▲□▶ ▲ 同

Google Alien, AlphaZero, AlphaGo

Wolfram

Cyc

A = 
A = 
A

Vol. LIX. No. 236.]

MIND

[October, 1950

#### A OUARTERLY REVIEW OF

#### PSVCHOLOGY AND PHILOSOPHY

#### I.-COMPUTING MACHINERY AND INTELLIGENCE

#### By A. M. TURING

#### 1 The Imitation Game

I PROPOSE to consider the question, 'Can machines think ?' This should begin with definitions of the meaning of the terms 'machine 'and 'think '. The definitions might be framed so as to reflect so far as possible the normal use of the words, but this attitude is dangerous. If the meaning of the words 'machine' and ' think ' are to be found by examining how they are commonly used it is difficult to escape the conclusion that the meaning and the answer to the question, ' Can machines think ?' is to be sought in a statistical survey such as a Gallup poll. But this is absurd. Instead of attempting such a definition I shall replace the question by another, which is closely related to it and is expressed in relatively unambiguous words.

The new form of the problem can be described in terms of a game which we call the 'imitation game'. It is played with three people, a man (A), a woman (B), and an interrogator (C) who may be of either sex. The interrogator stays in a room apart from the other two. The object of the game for the interrogator is to determine which of the other two is the man and which is the woman. He knows them by labels X and Y, and at the end of the game he says either 'X is A and Y is B ' or 'X is B and Y is A'. The interrogator is allowed to put questions to A and B thus:

C: Will X please tell me the length of his or her hair ? Now suppose X is actually A, then A must answer. It is A's 98 433



#### Lovelace essentially envisioned what computers were capable of before computers really existed.



## McKinsey Global Institute

McKinsey Global Institute

## Notes from the AI frontier: Modeling the impact of AI on the world economy

September 2018 | Discussion Paper

Figure 2 – Expected gains from AI in the different regions of the world by 2030

## Economic impacts of artificial intelligence (AI)

**pwc** 

**trillion** from AI while the outer circles are proportional to % GDP the 2030 regional GDP according to PwC projections

Dark circles represent the expected gains derived

Source: The macroeconomic impact of artificial intelligence, PwC, 2018.

## **PwC's Global Artificial Intelligence Study: Sizing**

## the prize

## Technology, communications and entertainment

### Healthcare

#### Three areas with the biggest AI potential

- Supporting diagnosis in areas such as detecting small variations from the baseline in patients' health data or comparison with similar patients.
- Early identification of potential pandemics and tracking incidence of the disease to help prevent and contain its spread.
- Imaging diagnostics (radiology, pathology).

#### Three areas with the biggest AI potential

- Media archiving and search bringing together diffuse content for recommendation.
- Customised content creation (marketing, film, music, etc.).
- Personalised marketing and advertising.

### Automotive

#### Three areas with the biggest AI potential

- Autonomous fleets for ride sharing.
- Semi-autonomous features such as driver assist.

▲□▶ ▲□▶ ▲□▶ ▲□▶ ▲□ ● ● ●

• Engine monitoring and predictive, autonomous maintenance.

## Nine charts that really bring home just how fast Al is growing

Artificial intelligence is booming in Europe, Chi but it's still a very male industry.

by Will Knight December 12, 2018

1. Al is being commercialized at a dizzying pace.





MIT Technology Review





(日) э

#### 2. The focal points are China and the US, but also Europe.

Much has been made of China's rising AI provess (see "China's AI avakening") and its growing rivalry with the US. As the data shows, Europe is also a huge hub of AI activity. But it seems that three main centers of power are emerging.



#### 5. Artificial intelligence is a political issue.



Mentions of artificial intelligence and machine learning in the US Congress (above) and the UK Parliament (below) have exploded in the past few years. This reflects a growing awareness of the technology's economic and strategic importance (see "Canada and France propose an international panel on AI").



▲□▶ ▲□▶ ▲ □▶ ▲ □ ▶ ▲ □ ● ● ● ●

#### 4. The state of the art is improving fast.



The report includes several measures of technical progress, including the accuracy of object recognition in images, measured against average human performance (top), and the accuracy of machine translations of news articles, measured using a score assigned by human judges (bottom). These don't mean that the field is getting closer to developing a human-level AI, but they show how key techniques have been honed in recent years. English > German



German > English





## **10 Breakthrough Technologies** 2017

hese technologies all have staying power. They will affect the economy and our politics, improve medicine, or influence our culture. Some are unfolding now; others will take a decade or more to develop. But you should know about all of them right now.

### **Dueling Neural Networks**



ILLUSTRATION BY DEREK BRAHNEY | DIAGRAM COURTESY OF MICHAEL NIELSEN, "NEURAL NETWORKS AND DEEP LEARNING", DETERMINATION PRESS, 2015



#### **Reinforcement Learning**

By experimenting, computers are figuring out how to do things that no programmer could teach them.

人口 医水黄 医水黄 医水黄素 化甘油

Availability: 1 to 2 years

by Will Knight

## The Rise of the 'Unicorns'



## The PayPal Mafia's Golden Touch

#### Meet the Dons who run Silicon Valley – the PayPal Mafia.

Despite the name, this 'outfi' is an entirely legit group of tech entrepreneurs - the founders and former employees of PayPal. The 'Maffai' have continued to work and invest together since the company's 2002 sale to eBay. Between them, they've founded, funded or led some of the world's biggest tech Imms.

Iousehold names



Bubble size indicates total value of all investments in a company

・ロト ・ 国 ト ・ ヨ ト ・ ヨ ト

э

COMPANY	MOST RECENT VALUATION	YEARLY INVESTMENT IN MILLIONS	EED ANGEL	WISC. GROWTH	TOTAL AMOUNT RAISED
Uber Uber Technologies' app connects drivers to riders in 55 countries, with pricing that varies with demand.	\$40 billion	'05 '06 '07 '08	*0.2 1.3 50 258	up to 1,000 3,200	\$6.3 billion
Airbnb A website that allows users to book apartments, homes and other temporary accommodations.	\$20 billion		\$0.6 7.2 112 200	up 475 to 1,000	\$1.8 billion
Snapchat Its messaging app lets users take and share images and videos that disappear once they are viewed.	\$15 billion		\$0.5 ÷	485 200	\$809 million
Dropbox A service that lets users upload and share photos, documents and videos with specific people.	\$10 billion	\$1.2 6	250	350 500	\$1.1 billion
Palantir Palantir Technologies specializes in software used in counterterrorism and the financial industry.	\$15 billion	unk. unk. seed V.C. \$35	90 120 56 197	444	\$942 million

S



gofundme



## Renaissance

## Al in Finance

## DE Shaw & Co



## Rich Formula: Math And Computer Wizards Now Billionaires Thanks To Quant Trading Secrets



Nathan Vardi Forbes Staff Hedge Funds & Private Equity Following the money trail

By then, quants and others in various industries had embraced Simons' machine learning techniques. Indeed, he and his colleagues at Renaissance had anticipated the transformation in decision-making that's sweeping almost every business and walk of life. More companies and individuals are accepting and embracing models that continuously learn from their successes and failures. As investor Matthew Granade has noted, Amazon, Tencent, Netflix, and others that rely on dynamic, ever-changing models are emerging dominant. The more data that's fed into machines, the smarter they're supposed to become. Just as Jim Simons predicted.

#### Markets

## Two Sigma Hires Google Scientist Mike Schuster for Al Expansion

60 44

(-55.90)

March of the machines

(-60.01)

598.71

The stockmarket is now run by computers, algorithms and passive managers

Such a development raises questions about the function of markets, how companies are governed and financial stability

(日)
# **Q** Palantir





・ロト ・四ト ・ヨト ・ヨ





# THE TIMES

▲□▶ ▲□▶ ▲□▶ ▲□▶ ■ ●の00

# Palantir, the tech spooks who found bin Laden

# MANAGE DATA LIKE CODE

Our platforms use versioning technology so you can manage data like software engineers manage code.

#### Tech Policy Jan 27

# 40 groups have called for a US moratorium on facial recognition technology



# It's too late to ban face recognition – here's what we need instead

The Secretive Company That Might End Privacy as We Know It



A little-known start-up helps law enforcement match photos of unknown people to their online images — and "might lead to a dystopian future or something," a backer says.



#### In The News



### Facial-recognition tech is a win for the little guy

"For all the hysteria over facialrecognition technology, the breathtakingly fast arrest of Larry Griffin II for prompting a major subway panic shows that the tech is a boon to law enforcement."

#### The New Hork Times THE WALL STREET JOURNAL.

#### Video Games and Online Chats Are 'Hunting Grounds' for Sexual Predators

"And so law enforcement officials from across the state took over a building near the Jersey Shore last year, and started chatting under assumed identities as children. In less than a week, they arrested 24 people."

#### Have No Fear of Facial Recognition

"In another disturbing case, with explicit video, someone sexually harassed and extorted young girls. Police used facial recognition to identify 14 of 22 victims, who were carefully interviewed. They eventually identified the offender."

▲ロ ▶ ▲周 ▶ ▲ 国 ▶ ▲ 国 ▶ ● の Q @

### Available now for Law Enforcement

**Request Access** 

# Stop Searching. Start Solving.

Clearview provides clients with its proprietary technology, database and investigative tools on a subscription basis. A Licensed User's subscription includes:

✓ Unlimited Use of CV's Proprietary Research System for its Licensed Users.

✓ Unlimited Access to CV's Proprietary Image Database for its Licensed Users.

✓ Each Licensed User Account Includes iPhone/Android CV Application

Each Licensed User Account Includes Lap/Desktop Versions of CV Program
Help-Desk Support

Annual 12-month Subscription Rate: 5 Seats: \$10,000 10 Seats: \$15,000 20 Seats: \$25,000 50 Seats: \$50,000 12 Seats: \$00,000 12 Seats: \$25,00,000 10 Initial License Unlimited Users: Nerrotiated Flat Fee

For More Information: Jessica Medeiros Garrison (e) Jessica@clearview.ai (c) 205.568.4371

Time is law enforcement's most valuable resource. Cleanlew puts the world's most advanced facial-recognition technology and largest image database into their hands, allowing them to turn a photograph into a solid lead in an instant.

GClearview Ai

Clearview Ai, Inc. 15 West 72nd St. Suite 23-S. New York, NY 10023 Our office spent 12 man hours over a month's time trying to identify a theft suspect. We ran the picture through Clearview and identified the suspect in seconds.



If we had Clearview at the time when the report came in, we would not only have identified the suspect sooner, but also would have prevented other thefts that the suspect committed before we arrested him.

John Hodger



World's best facial-recognition technology combined with the world's largest database of headshots.

Real-time Results.



World-Class Accuracy. Control of the set of





Detective Sgt. Nick Ferrara in Gainesville, Fla., said he had used Clearview's app to identify dozens of suspects. Charlotte Kesl for The New York Times





The patterns on the goods in this shop are designed to trigger Automated License Plate Readers, injecting junk data in to the systems used by the State and its contractors to monitor and track civilians and their locations.

(日) (四) (日) (日) (日)





Artificial intelligence and the core technology of machine learning are likely to be a general purpose technology of the scale and scope similar to electricity or even the steam engine before it - fundamentally revolutionizing many, many sectors of the economy.

> ERIK BRYNJOLFSSON DIRECTOR OF THE MIT INITIATIVE ON THE DIGITAL ECONOMY

> > ・ロト ・ 理 ト ・ ヨ ト ・ ヨ ト

-

### Perhaps as important as steam engines back in the days..

Erik Brynjolfsson, ICLR 2018

FIGURE 1.2 What Bent the Curve of Human History? The Industrial Revolution.



= 900

A D > A P > A B > A B >

## Back to the industrial revolution

from Erik Brynjolfsson, ICLR 2018

- What can history tell us?
- Steam engine was classified by Bresnahan et Trajtenberg (1996) as belonging to the so-called General Purpose Technologies

- Those technologies are characterized by 3 features:
  - Pervasive
  - Able to be improved over time
  - Able to spawn complemetary innovations
- Does that remind you of something ?

## Back to the industrial revolution

from Erik Brynjolfsson, ICLR 2018

- Technology is not neutral
- You can do a small pox vaccine...

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 のへぐ

## Back to the industrial revolution

from Erik Brynjolfsson, ICLR 2018

- Technology is not neutral
- You can do a small pox vaccine...But you can also create a nuclear weapon

▲□▶ ▲□▶ ▲□▶ ▲□▶ ■ ●の00

• In fact let's compare..

# **IS 'Progress' Good for Humanity?**

Rethinking the narrative of economic development, with sustainability in mind

JEREMY CARADONNA SEP 9, 2014



Rage against the machine: Luddites smashing a loom. (CHRIS SUNDE / W



— Colin Stretch, general counsel for Facebook, Sean Edgett, acting general counsel for Twitter, Richard Salgado, director of law enforcement and information security at Google, testify before the Senate Judiciary Committee's hearing on "Extremist Content and Russian Disinformation Online: Working with Tech to Find Solution's on Capitol Hill in Washington DC on Oct. 31, 2017. Sharen Teev. / PA

Mr. Colin Stretc

Mr. Sean Edgett

Mr. Richard Salead

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 のへで



# The challenges of driving a yellow cab in the age of Uber

By John Crudele

September 18, 2017 | 10:35p

BERATC.COM/CAR

UBER



UBER

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 三臣 - のへ⊙





◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 のへぐ

Population of horses in the US during industrialization









# Cambridge Analytica

# Facebook Bans Deepfake Videos That Could Sway Voters, But Is It Enough?





John Brandon Contributor ① Social Media John Brandon covers social media trends. @jmbrandonbb

# Deep fake' imagery manipulation poses threat to society not just military, US General warns



Reference

**Our Result** 



#### Snow and Ice Pose a Vexing Obstacle for Self-Driving Cars

Most testing of autonomous vehicles until now has been in sunny, dry climates. That will have to change before the technology will be useful everywhere.



Self-driving cars will need to identify





AI

Kathleen Walch Contributor COGNITIVE WORLD Contributor Group ③

▲□▶ ▲□▶ ▲□▶ ▲□▶ □ のQで

The Canadian driving data set includes lane markings and vehicles that are covered with snow. COURTESY OF JEFF HILNBRAND

Horses, Equine Law And The Future Of The Autonomous Vehicle Legal Framework

DEADLY CRASH WITH SELF-DRIVING UBER

ARIZON 11:01 64<mark>0NA</mark>



TEMPE

••• 15

Rahul Razdan Contributor ③ Transportation I focus on the impact of autonomous systems on society.

# Universities told to be Al's ethical watchdogs



Industry and government alone cannot oversee new technology, experts say



# Engineers will be crucial to Al-driven economies

Just as the AI revolution calls for more computer scientists, engineers will be needed to develop next-generation AI hardware, says Bashir M. AI-Hashimi



#### **CIO JOURNAL**

# What Machine Learning Can and Cannot Do

Jul 27, 2018 1:56 pm ET



A doctor examines a magnetic resonance image of a human brain during a Beijing neuroimaging competition between human doctors and AI, June 30, 2018. PHOTO: MARK SCHIEFELBEIN / ASSOCIATED PRESS

# What Machine Learning can and cannot do

- We have seen many achievements (essentially in vision, language)..
- In particular Supervised learning has been quite successful
- But there are still plenty of tasks that computers still cannot handle (see Lex Friedman, MIT Sloan lecture)

▲□▶ ▲□▶ ▲□▶ ▲□▶ ■ ●の00

- Awareness of self
- Emotion
- Imagination
- Morality
- Consciousness
- high level reasoning

Still many tasks that machines cannot do

from Erik Brynjolfsson, ICLR 2018 keynote.



# Immediate Challenges



# Immediate Challenges

(Lex Friedman, MIT Sloan)

- Occlusions
- Sensor spoofing (camera, Lidar)
- adversarial noise
- Risk quantification
- Data is costly ⇒ Ideally, we would like to move on to more unsupervised learning.

- ロ ト - 4 回 ト - 4 □

# MIT Technology Review

Feature p. 42 A 3-D Printer That Really Matters

Feature p. 78 Cancer Cures For a Lucky Few

Feature p. 28 Time to Consider Geoengineering?

### **Mysterious Machines**





Artificial Intelligence / Machine Learning

# The Dark Secret at the Heart of AI

No one really knows how the most advanced algorithms do what they do. That could be a problem.

by Will Knight

"We can build these models, but we don't know how they work."



# How well can we get along with machines that are unpredictable and inscrutable?

Just as many aspects of human behavior are impossible to explain in detail, perhaps it won't be possible for AI to explain everything it does. "Even if somebody can give you a reasonablesounding explanation [for his or her actions], it probably is incomplete, and the same could very well be true for AI," says Clune, of the University of Wyoming. "It might just be part of the nature of intelligence that only part of it is exposed to rational explanation. Some of it is just instinctual, or subconscious, or inscrutable."

#### Deep neural networks are coming to your phone



Already, mathematical models are being used to help determine who makes parole, who's approved for a loan, and who gets hired for a job. If you could get access to these mathematical models, it would be possible to understand their reasoning. But banks, the military, employers, and others are now turning their attention to more complex machine-learning approaches that could make automated decision-making altogether inscrutable.

> "Whether it's an investment decision, a medical decision, or maybe a military decision, you don't want to just rely on a 'black box' method."

> > ▲□▶ ▲□▶ ▲□▶ ▲□▶ ■ ●の00

## Machine Learning is not new..

• General principle of machine learning is very simple



 One of the reason for the renewed excitement is Massive parallelism through Graphical Processing Units ⇒ Essential for neural network training on massive databases (think of imageNet, GoogleNet,..)



▲□▶ ▲□▶ ▲□▶ ▲□▶ □ のQで

### Some new architectures are coming



The general picture: Supervised vs Unsupervised

• Supervised learning tries to understand the relation between data  $\mathcal{D} = \{x_i\}_{i=1}^N$  and the associated labels (targets)  $\{y_i\}$  based on a subset of samples for which the labels are known. The goal is thus to find a mapping between the points  $x_j \in \mathcal{D}$  and their targets  $y_i$ . During training all labels are known.

Ex: Handwriting recognition. Data = images from MNIST, labels,knowledge = actual numbers displayed

• Unsupervised learning tries to find structure within a given (unlabeled) dataset

Ex: customized advertising (cluster users in groups in order to send specific advertising to each group)

# The big picture: Supervised vs Unsupervised

- Semi-supervised Learning (SSL): small amount of labeled data with large amount of unlabeled data. Other examples of partial supervision can involve constraints on the prototypes x<sub>i</sub> (such as requiring subgroups of prototypes to have the same target)
- Reinforcement: The machine looks for suitable actions, in a given situation, in order to maximize a "reward" (e.g. Neural Network backgammon playing : board position + dice value ⇒ move + reward)

・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・



◆□▶ ◆圖▶ ◆臣▶ ◆臣▶ 三臣 - 釣��

#### Semi-supervised: How unlabelled data can help?





◆□▶ ◆□▶ ◆臣▶ ◆臣▶ ○臣 ○○へ⊙

# Supervised learning

Supervised learning can be subdivided into two main classes of problems:

Regression problems (approx. a mapping f(x<sub>μ</sub>) from inputs x<sub>μ</sub> to continuous output variables y<sub>μ</sub>)

*Ex.* 
$$\hat{\beta} = \min_{\beta} \sum_{i=1}^{N} (y_i - \langle x_i, \beta \rangle)^2$$
 (residual SS)

Classification problems (approx. a mapping f(x<sub>μ</sub>) from inputs x<sub>μ</sub> to discrete output variables y<sub>μ</sub>)

$$E_{X}. \quad p(\mathcal{C}_{k}|\boldsymbol{x}) = \frac{p(\boldsymbol{x}|\mathcal{C}_{k})p(\mathcal{C}_{k})}{p(\boldsymbol{x})} = \frac{N_{\delta,K}}{N_{\delta}}, \quad (\text{Nearest Neighbors})$$


## Examples of supervised learning algorithms

Many possible regression/classification models (some models can be used in both frameworks)

▲□▶ ▲□▶ ▲□▶ ▲□▶ ▲□ ● ● ●

- Linear Regression
- Neural networks
- Support vector machines (SVMs)
- K-nearest neighbors
- Tree based models
- ...

#### Supervised Learning: What are we going to learn?

# SVMs for face recognition





















Logistic regression, neural nets for handwritten digits recognition

#### Neural nets pong training





イロト 不得 トイヨト イヨト 3

#### Parametric vs non-parametric

- Fixed number of parameters = parametric
  - +: faster to use
  - -: stronger assumptions regarding data distribution.

#### Ex. linear regression

• Complexity of the model grows with size of training dataset = non-parametric

▲□▶ ▲□▶ ▲□▶ ▲□▶ ▲□ ● ● ●

- +: more flexible regarding data.
- -: often computationally intractable for large datasets

#### Ex. KNN

#### Parametric vs non-parametric

#### • From Hastie, Tibshirani, Friedman



FIGURE 2.1. A classification example in two dimensions. The classes are as a binary variable (BLUE = 0, ORMGE = 1), and then fit by linear regre. The line is the decision boundary defined by  $x^T \beta = 0.5$ . The orange shadel 4 denotes that part of input space classified as ORANCE, while the blue regrees lossified as BLUE.

#### 15-Nearest Neighbor Classifier



FIGURE 2.2. The same classification example in two dimensions as in Figure 2.1. The classes are coded as a binary variable (BLUE = 0, ORNOE = 1) and then fit by 15-nearest-neighbor averaging as in (2.8). The predicted class is hence chosen by majority vote amongst the 15-nearest neighbors.

・ロト ・ 日 ト ・ 日 ト ・ 日 ト

э

#### Unsupervised learning

In unsupervised learning, we are only given prototypes  $x_\mu$  and we want to extract pattern from the data.

Examples of unsupervised learning approaches include

- Clustering (K-means, K-medoid)
- Manifold learning/Latent variable models
- Factor Analysis, Principal and Independent Component Analysis (PCA/ICA)

・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・

• Gaussian mixture models (GMM)



▲□▶ ▲□▶ ▲臣▶ ★臣▶ = 臣 = のへで

## Image segmentation through clustering



#### You Can Trick Self-Driving Cars by Defacing Street Signs





What are we going to learn? Combined supervised and unsupervised

Samir Bhatt, Bhaskar Trivedi, Ankur Devani, Hemang Bhimani (eInfochips)

## Short reminder on Linear Algebra: motivation

Why are Linear algebra and Differential Calculus useful in Machine Learning?

- 1. Linear algebra is useful in both supervised and supervised learning to derive the expression of the model's parameters given the training data.
  - Several models in unsupervised learning rely on matrix factorization. Singular value decomposition, for example can be used for dimensionality reduction (PCA, ICA)
  - When looking for an optimal model, given a particular dataset *D*, we will often have to find the solution to multivariate systems of linear equations. Linear algebra will be useful to understand how we can compute the solution efficiently when those systems are overdetermined or even ill conditioned.
- 2. The tools from linear algebra will also be useful when will need concepts from multivariate statistics

(also see the notes of Z. Kolter)

 Vector products. Given two vectors x and y one can define the quantity x<sup>T</sup>y which is called the inner product or scalar product of x and y

$$\langle \boldsymbol{x}, \boldsymbol{y} \rangle = \boldsymbol{x}^T \boldsymbol{y} \in \mathbb{R} = [x_1, x_2, \dots, x_n] \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} = \sum_{i=1}^n x_i y_i.$$

- Norms. A norm is any function N : ℝ<sup>n</sup> → ℝ that satisfies the following properties
  - 1. For all  $x \in \mathbb{R}^n$ ,  $N(x) \ge 0$
  - 2. N(x) = 0 if and only if x = 0
  - 3. For all  $x \in \mathbb{R}^n$ ,  $\alpha \in \mathbb{R}$ ,  $N(\alpha x) = |\alpha|N(x)$
  - 4. For all  $x, y \in \mathbb{R}^n$ ,  $N(x + y) \le N(x) + N(y)$

(also see the notes of Z. Kolter)

• Common examples of norms include the  $\ell_2$  norm  $\|x\|_2 = \sqrt{\sum_{i=1}^n x_i^2}$ , the  $\ell_1$  norm  $\|x\|_1 = \sum_{i=1}^n |x_i|$  and the  $\ell_\infty$  norm  $\|x\|_\infty = \max_i |x_i|$ .

・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・

• More generally, the familly of  $\ell_p$  norms is defined as  $||x||_p = (\sum_{i=1}^n |x_i|^p)^{1/p}$ 

(also see the notes of Z. Kolter)

- Similar ideas hold for matrices. For two matrices A and B, one can define an inner product ⟨A, B⟩ as ⟨A, B⟩ = ∑<sub>i=1</sub><sup>n</sup> ∑<sub>j=1</sub><sup>n</sup> A<sub>ij</sub>B<sub>ij</sub>.
- Just as we define the  $\ell_p$  norms for vectors, we can define spectral *p*-norms on matrices (which are jsut  $\ell_p$  on the vector of singular values)

$$\|A\|_{1} = \sum_{i=1}^{n} \sigma_{i}(A)$$
$$\|A\|_{2} = \sqrt{\sum_{i=1}^{n} \sigma_{i}^{2}(A)}$$
$$\|A\|_{\inf} = \sup_{1 \le i \le n} \sigma_{i}(A)$$

◆□ ▶ ◆□ ▶ ◆ 三 ▶ ◆ 三 ● ● ● ●

(also see the notes of Z. Kolter)

- Linear independence. A set of vectors is said to be linearly independent if no vector can be represented as a linear combination of the remaining set of vectors.
- The rank of a matrix A is the largest subset of linearly independent columns of A.
- The inverse of a square matrix  $A \in \mathbb{R}^{n \times n}$  is the unique matrix  $A^{-1}$  such that  $A^{-1}A = I = AA^{-1}$
- Two vectors  $x, y \in \mathbb{R}^d$  are orthogonal if  $\langle x, y \rangle = 0$ . A vector x is normalized if ||x|| = 1
- A square matrix  $U \in \mathbb{R}^{n \times n}$  is orthogonal if  $U^T U = UU^T = I$ .

#### Eigenvalues and eigenvectors

(also see the notes of Z. Kolter)

• Given a square matrix  $A \in \mathbb{R}^{n \times n}$ , we say that  $\lambda \in \mathbb{C}$  is an eigenvalue of A and  $\mathbf{x} \in \mathbb{C}^n$  is its corresponding eigenvector if

$$A\mathbf{x} = \lambda \mathbf{x}, \quad \mathbf{x} \neq \mathbf{0}$$

Given a symmetric matrix A ∈ S<sup>n×n</sup>, all the eigenvalues of A are real. Moreover, the eigenvectors of A are orthonormal. If we store the eigenvectors in U and the eigenvalues in the diagonal matrix Λ, A can read as

$$A = U \Lambda U^*$$

• From this, for any vector  $x \in \mathbb{R}^n$ , we thus have  $x = \sum_i t_i u_i + \overline{x}$  where  $\langle \overline{x}, u_i \rangle = 0$  for all  $u_i$  which gives  $x^T A x = \sum_i \lambda_i \left( \sum_j t_j u_j^T \right) u_i u_i^T \left( \sum_j t_j u_j \right) = \sum_i |t_i|^2 \lambda_i$ 

#### Eigenvalues and eigenvectors

• In particular, when  $||x||^2 = 1 = \sum_i t_i^2 ||u_i||^2 = \sum_i t_i$ you see that the quantity

$$x^{T}Ax = \sum_{i} \lambda_{i} \left( \sum_{j} t_{j} u_{j}^{T} \right) u_{i} u_{i}^{T} \left( \sum_{j} t_{j} u_{j} \right) = \sum_{i} |t_{i}|^{2} \lambda_{i}$$

is maximized when  $x = t_j u_j$  where  $u_j$  is the eigenvector corresponding to the largest eigenvalue  $\lambda_{max}$ .

• The solution to the maximization problem

$$\max_{x \in \mathbb{R}^n} x^T A x \text{ subject to } ||x||_2^2 = 1$$

is thus the eigenvector corresponding to the largest eigenvalue

 We will use this later when discussing latent variable models in unsupervised learning.

#### (Matrix) differential calculus

Given a function f : ℝ<sup>m×n</sup> → ℝ, that takes as input the matrix of unknowns X the gradient of f(X) with respect to X is defined as fhte matrix

$$\nabla_{\boldsymbol{X}} f(\boldsymbol{X}) \in \mathbb{R}^{m \times n} = \begin{bmatrix} \frac{\partial f(\boldsymbol{X})}{\partial X_{11}} & \frac{\partial f(\boldsymbol{X})}{\partial X_{12}} & \cdots & \frac{\partial f(\boldsymbol{X})}{\partial X_{1n}} \\ \frac{\partial f(\boldsymbol{X})}{\partial X_{21}} & \frac{\partial f(\boldsymbol{X})}{\partial X_{22}} & \cdots & \frac{\partial f(\boldsymbol{X})}{\partial X_{2n}} \\ \vdots & \ddots & \vdots \\ \frac{\partial f(\boldsymbol{X})}{\partial X_{m1}} & \frac{\partial f(\boldsymbol{X})}{\partial X_{m2}} & \cdots & \frac{\partial f(\boldsymbol{X})}{\partial X_{mn}} \end{bmatrix}$$

• When X is a vector, we recover

$$\nabla_{\boldsymbol{x}} f(\boldsymbol{x}) = \begin{bmatrix} \frac{\partial f(\boldsymbol{x})}{\partial x_1} \\ \frac{\partial f(\boldsymbol{x})}{\partial x_2} \\ \frac{\partial f(\boldsymbol{x})}{\partial x_n} \end{bmatrix}$$

▲□▶ ▲□▶ ▲□▶ ▲□▶ ▲□ ● ● ●

## (Matrix) differential calculus

• When looking for the optimal model in terms of empirical risk function, we will often have to find the value x that minimizes the sum of squred residuals

$$\min_{x} \|Ax - b\|_{2}^{2} = x^{T} A^{T} A x - 2b^{T} A x + b^{T} b$$

• The solution in this case can be computed in closed form by first taking the derivative with respect to *x*,

$$\nabla_{\mathbf{x}}(\mathbf{x}^{\mathsf{T}}\mathbf{A}^{\mathsf{T}}\mathbf{A}\mathbf{x} - 2\mathbf{b}^{\mathsf{T}}\mathbf{A}\mathbf{x} + \mathbf{b}^{\mathsf{T}}\mathbf{b}) = 2\mathbf{A}^{\mathsf{T}}\mathbf{A}\mathbf{x} - 2\mathbf{A}^{\mathsf{T}}\mathbf{b}$$

Setting the derivative to 0, we get

$$x = (A^T A)^{-1} A^T b$$

・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・

provided that  $A^T A$  is invertible.

## (Matrix) differential calculus

- When deriving a closed form solution is not possible, the model (A, b) can be learned through gradient descent
- For any given function f(x), the gradient gives a vector pointing in the direction of greatest increase of the function. Hence -∇f(x) is a vector pointing in the direction of greatest decrease of f(x)
- Starting from an initial guess x<sup>(0)</sup> for the solution of min<sub>x</sub> f(x), one can thus define an iterative procedure (known as gradient descent) as

$$\mathbf{x}^{(t+1)} \leftarrow \mathbf{x}^{(t)} - \alpha \nabla f(\mathbf{x}^{(t)})$$

 $\alpha$ , which represents the step size, is also known as the learning rate. When  $\mathbf{x} \in \mathbb{R}^n$  is a vector, recall that we have

$$\nabla f(\mathbf{x}) = \left[\frac{\partial f}{\partial x_1}, \dots, \frac{\partial f}{\partial x_n}\right]$$

#### Gradient descent



#### Reminders on Statistics and probability

- Why using stats/proba?
- Machine Learning relies on complex distributions (cancerous cells, possible moves in Go, Existing sign roads, possible evolutions of stocks,...)
- Only a few samples are usually available
- ⇒ We need a way to measure how well those samples are representing the underlying (unknown) distribution

#### Why is that important?

#### Self-Driving Uber Car Kills Pedestrian in Arizona, Where Robots Roam



A woman crossing Mill Avenue at its intersection with Curry Road in Tempe, Ariz. on Monday. A pedestrian was struck and killed by a self-driving Uber vehicle at the intersection a night earlier. Catilin O'Hara for The New York Times

## Reminders (I)

#### (Discrete sets of events)

- Sum rule  $p(X) = \sum_{Y} p(X|Y)$
- Product rule p(X, Y) = p(X|Y)p(Y)
- Bayes theorem

$$p(Y|X) = \frac{p(X|Y)p(Y)}{p(X)}$$

(continuous sets of events)

• density p(x),

$$p(x \in [a, b]) = \int_a^b p(x) dx, \quad p(x) = \int p(x, y) dy$$

▲□▶ ▲□▶ ▲□▶ ▲□▶ ▲□ ● ● ●

## Reminders (II)

• Cumulative distribution Function (CDF)  $F(z) = \int_{-\infty}^{z} p(x) dx$ 

- Expectation  $\mathbb{E}[x] = \int xp(x)dx$ ,  $\mathbb{E}[x] = \sum_i x_i p(x_i)$
- Conditional expectation  $\mathbb{E}_x f(x|y) = \sum_x f(x)p(x|y)$
- Variance Var[x]  $\equiv \mathbb{E}\left\{(x \mathbb{E}x)^2\right\}$
- Covariance  $Cov[x, y] \equiv \mathbb{E} \{ (x \mathbb{E}x)(y \mathbb{E}y) \}$

## Reminders (III) A few important distributions

• The gaussian distribution

$$\mathcal{N}(x|\mu,\sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2\sigma^2}(x-\mu)^2\right)$$

- Uniform distribution:  $P(y) = \frac{1}{|b-a|}, y \in [a, b]$
- $\chi^2$  distribution:  $\chi^2 \sim \sum_{i=1}^{N} Z_i^2$  with  $Z_i$  independent standard normal RV.

## Reminders (IV) A few important distributions

• Binary variables: Bernoulli and Rademacher,

$$Bern(x|\mu) = \mu^{x}(1-\mu)^{1-x}, \quad x = \begin{cases} 1 \\ 0 \end{cases}, 0 \le \mu \le 1$$
  
Rademacher:  $\varepsilon(x) = \begin{cases} (1/2), & x = +1 \\ (1/2), & x = -1 \\ 0, & \text{otherwise} \end{cases}$ 

### The exponential family

- Many of the distributions we have discussed are part of a general family called The exponential family
- The exponential family has many interesting properties
  - It is the only family of distribution with finite-sized sufficient statistics (see next slides)

- It is the only family with known conjugate priors
- It is at the core of generalized linear models
- it is at the core of variational inference
- We will come back to these notions later

#### The exponential family

 A pdf p(x|θ) is said to be in the exponential family for x = (x<sub>1</sub>,..., x<sub>m</sub>) and θ ⊆ ℝ<sup>d</sup> if

$$p(\mathbf{x}|\boldsymbol{\theta}) = \frac{1}{Z(\boldsymbol{\theta})} h(\mathbf{x}) \exp(\boldsymbol{\theta}^{T} \phi(\mathbf{x}))$$
$$= h(\mathbf{x}) \exp(\boldsymbol{\theta}^{T} \phi(\mathbf{x}) - A(\boldsymbol{\theta}))$$

•  $Z(\theta)$  and  $A(\theta)$  are defined as

$$Z(\theta) = \int_{\mathcal{X}^m} h(\mathbf{x}) \exp[\theta^T \phi(\mathbf{x})] \, d\mathbf{x}$$
$$A(\theta) = \log(Z(\theta))$$

•  $Z(\theta)$  is called the partition function,  $\theta$  are the mutual parameters,  $\phi(x) \in \mathbb{R}^d$  is the vector of sufficient statistics,  $A(\theta)$  is the log partition function or cumulant function.

## The exponential family

- Two examples
  - Bernoulli

$$Ber(x|\mu) = \mu^{x}(1-\mu)^{1-x} = \exp(x\log(\mu) + (1-x)\log(1-\mu))$$
  
=  $\exp(\phi(x)^{T}\theta)$ 

◆□▶ ◆□▶ ◆ 臣▶ ◆ 臣▶ ○ 臣 ○ の Q @

- Univariate Gaussian
- The Uniform distribution does not belong to the exponential family

#### Parameter inference: What does it mean?



Background vector created by freepik - www.freepik.com

(日) (四) (日) (日) (日)

Parameter/model inference: Bayesian vs frequentist

- Several models in ML are special instance of a more general idea called model inference
- Inference can be used in both supervised (learn new labels from training labels) and unsupervised (learn parameters from distribution) frameworks

・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・

• Two main approaches: frequentist and Bayesian.

Parameter/model inference: Bayesian vs frequentist

- Bayesian statistics.
  - Considers the (distribution) parameters as random
  - Relies heavily on the posterior distribution  $p(\theta | D)$
  - dominated statistical practice before 20<sup>th</sup> century
  - Ex: MAP  $\underset{\theta}{\operatorname{argmax}} P(\mathcal{D}|\theta)P(\theta)$
- Frequentist statistics (a.k.a classical stat.)
  - Parameters θ viewed as fixed, sample D as random (Randomness in the data affects the posterior)
  - Relies on the likelihood or some other function of the data

- dominated statistical practice during 20<sup>th</sup> century
- Ex. MLE :  $\underset{\theta}{\operatorname{argmax}} P(\mathcal{D}|\theta)$

Bayesian statistics: Some vocabulary

- We saw Bayesian inference relies on the posterior  $p( heta | \mathcal{D})$
- The posterior reads from the Bayes rule as

$$p(\theta|\mathcal{D}) = \frac{p(\mathcal{D}|\theta)p(\theta)}{p(\mathcal{D})} \\ = \frac{p(\mathcal{D}|\theta)p(\theta)}{\int p(\mathcal{D}|\theta)p(\theta)d\theta}$$

- p(θ) is called the prior, p(D|θ) is called the likelihood function and Z = p(D) is the normalizing constant (independent of θ)
- Given a set of patterns  $(\mathbf{x}_{\mu}, \mathbf{y}_{\mu})$ , probabilistic classifiers are usually of two types:
  - Generative (learn model for  $p(\mathbf{x}, \mathbf{y}|\theta)$ )
  - Discriminative (learn model for  $p(y|x, \theta)$ )

#### Bayesian statistics: Some vocabulary

- An example of discriminative classifier : Logistic regression
  - Here we take µ(x) = sigm(w<sup>T</sup>x) and define the classifier as a Bernoulli distribution

$$p(y|\boldsymbol{x}, \boldsymbol{w}) = Ber(y|\mu(\boldsymbol{x}))$$

- Good when the output is binary
- An example of generative classifier :
  - relies on the assumption that the features (hidden variables) are independent

$$p(\mathbf{x}|y=c, \boldsymbol{\theta}) = \prod_{j=1}^{D} p(x_j|y=c, \theta_{jc})$$

- $\theta_{j,c}$  is the parameters of the distribution of class c for  $j^{th}$  entry in the D-dimensional pattern vector  $\mathbf{x} \in \{1, \dots, K\}^{D}$ .
- We will study those models in further detail when discussing classifiers.

#### Bayesian statistics

- In Bayesian statistics, randomness is most often used to encode uncertainty
- The posterior  $p(\theta|\mathcal{D})$  summarizes all we know on the parameters
- Bayesian inference is not always the right choice because of the following

▲□▶ ▲□▶ ▲□▶ ▲□▶ ▲□ ● ● ●

- The Mode is not a typical point in the distribution
- MAP estimator depends on re-parametrization

Bayesian statistics: Drawbacks and solutions

- A solution to the first part is to use a more robust loss function  $\ell(\hat{\theta},\theta)$
- A solution to the second part is to replace the MAP with the following estimator (when available)

$$\hat{\theta} = \underset{\theta}{\operatorname{argmax}} p(\mathcal{D}|\theta) p(\theta) |I(\theta)|^{-1/2}$$
(1)

where  $I(\theta)$  is the Fischer information matrix
#### Fischer information matrix

• For a generative model  $p(\mathbf{x}|\theta)$ , we let  $g(\theta, \mathbf{x})$  denote the Fischer score

$$g( heta, \mathbf{x}) = 
abla_{ heta} \log(p(\mathbf{x}| heta))$$

• the Fischer Kernel is the defined as

$$k(\pmb{x},\pmb{x}') = g(\pmb{ heta},\pmb{x})^T \pmb{F}^{-1} g(\pmb{ heta},\pmb{x}')$$

• The matrix  ${\pmb F}$  is called the Fischer matrix and defined as

$$m{F} = \mathbb{E}_{m{x}}\left\{g(m{ heta},m{x})g(m{ heta},m{x})^{T}
ight\}$$

• Note that it is often computed empirically as

$$\boldsymbol{F} \approx rac{1}{N} \sum_{n=1}^{N} g(\theta, \boldsymbol{x}) g(\theta, \boldsymbol{x})^{T}$$

## Occam's razor and Model selection

- Only looking for the best model often leads to overfitting (we will see that later in more details)
- Bayesian framework offers and alternative called Bayesian model selection
- For a family of models, we can express the posterior

$$p(m|\mathcal{D}) = \frac{p(\mathcal{D}|m)p(m)}{\sum_{m \in \mathcal{M}} p(m|\mathcal{D})}$$
$$\propto p(\mathcal{D}|m)p(m)$$

where  $p(\mathcal{D}|m) = \int p(\mathcal{D}|\theta)p(\theta|m)d\theta$  is called the marginal likelihood, integrated likelihood or evidence

# Occam's razor

• Integrating the parameters heta such as in

$$p(\mathcal{D}|m) = \int p(\mathcal{D}|\theta)p(\theta|m)d\theta$$

acts as a natural regularization and prevents overfitting when solving for  $\max_m p(m|D)$ . This idea is known as Bayesian Occam's razor

- The evidence p(D|m) can be understood as the probability to generate a particular dataset from a family of model (all values of the parameters included).
- When the family of models is too simple, or too complex, this probability will be low.

# Bayesian decision theory

- How do we resolve the lack of robustness of Bayesian estimators vis a vis the distribution (recall the bimodal distribution)?
- Statistical decision theory can be viewed as a game against nature.
- Nature has a parameter value in mind and gives us a sample
- We then have to guess what the value of the parameter is by choosing an action *a*
- As an additional piece of information, we also get a feedback from a loss function L(y, a) which tells us how compatible our action is w.r.t Nature's hidden state.

#### Bayesian decision theory

• The goal of the game is to determine the optimal decision procedure,

$$\underset{a \in \mathcal{A}}{\operatorname{argmin}} \mathbb{E} \left\{ L(y, a) \right\}$$

• In economics L(y, a) = U(y, a) and leads to the Maximum utility principle which is considered as rational behavior

$$\delta(x) = \operatorname*{argmax}_{a \in \mathcal{A}} \mathbb{E} \left\{ U(y, a) \right\}$$

 In the Bayesian framework, we want to minimize the loss over the models compatible with the observations {x<sub>μ</sub>}

$$\delta(\mathbf{x}) = \underset{\mathbf{a} \in \mathcal{A}}{\operatorname{argmin}} \mathbb{E}_{p(\theta | \{\mathbf{x}_{\mu}\})} \left\{ L(\theta, \mathbf{a}) \right\} = \sum_{\theta \in \Theta} L(\theta, \mathbf{a}) p(\theta | \{\mathbf{x}_{\mu}\}_{\mu})$$

・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・

# Bayesian decision theory (continued)

• The MAP is equivalent to minimizing a 0/1 loss

$$L(\hat{\theta}, \theta) = \mathbb{1}_{\theta \neq \hat{\theta}} = \begin{cases} 0 & \text{if } \hat{\theta} \neq \theta \\ 1 & \text{if } \hat{\theta} = \theta. \end{cases}$$

we then have

$$\mathbb{E}L(\hat{\theta},\theta) = p(\theta \neq \hat{\theta} | \{\boldsymbol{x}_{\mu}\}_{\mu}) = 1 - p(\hat{\theta} = \theta | \{\boldsymbol{x}_{\mu}\}_{\mu})$$

which is maximized for  $\hat{\theta}=\theta$  with  $\theta$  taken as

$$heta^*(\{m{x}_\mu\}_\mu) = \operatorname*{argmax}_{\hat{ heta}} p( heta|\{m{x}_\mu\}_\mu)$$

• The 0/1 loss means that each time you have an outcome different from your estimator, you are maximally penalized

What other losses can we choose?

• Is it possible to take more robust losses?

• 
$$\ell_2$$
 loss,  $L(\hat{\theta}, \theta) = |\hat{\theta} - \theta|^2$  gives posterior mean  
 $\mathbb{E}\left\{(\hat{\theta} - \theta)^2 | \mathbf{x}_{\mu}\right\} = \mathbb{E}[\theta^2 | \mathbf{x}_{\mu}] - 2\hat{\theta}\mathbb{E}[\theta | \mathbf{x}_{\mu}] + \hat{\theta}^2$ 

- setting derivative to 0,  $\partial_{\hat{ heta}} \mathbb{E}\{\hat{ heta}|m{x}_{\mu}\}=$  0, we get

$$-2\mathbb{E} \left\{ \theta | \mathbf{x}_{\mu} 
ight\} + 2\hat{ heta} = 0$$
 $\hat{ heta} = \int heta p( heta | \mathbf{x}_{\mu}) \ d heta$ 

## What other losses can we choose? (continued)

- Is it possible to take more robust losses?
- $\ell_1$  loss,  $L(\hat{ heta}, heta) = |\hat{ heta} heta|$  gives posterior median
- The value  $\hat{\theta}$  such that

$$oldsymbol{
ho}( heta < \hat{ heta} | oldsymbol{x}_{\mu}) = oldsymbol{
ho}( heta \geq \hat{ heta} | oldsymbol{x}_{\mu}) = 1/2$$

・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・

# What other losses can we choose? (continued)

- Now assume θ̂ defines the value of some hidden variable y
   (e.g. the class of a point x<sub>μ</sub> defined by a gaussian mixture θ̂).
- Finding the optimal parameters (or equivalently estimate the hidden state) can be done by considering the error

$$egin{aligned} & \mathcal{L}_{m{g}}( heta, \hat{ heta}) = \mathbb{E}_{(m{x}_{\mu}, y_{\mu}) \sim m{p}(m{x}_{\mu}, y_{\mu} | heta)} \left\{ \ell( heta, f_{\hat{ heta}}) 
ight\} \ &= \sum_{m{x}_{\mu}} \sum_{y_{\mu}} \ell(y_{\mu}, f_{\hat{ heta}}(x_{\mu})) m{p}(x_{\mu}, y_{\mu} | m{ heta}) \end{aligned}$$

• The Bayesian approach then minimizes the posterior expected loss

$$\underset{\hat{\theta}}{\operatorname{argmin}} \int p(\theta | \mathcal{D}) L_g(\theta, \hat{\theta}) \ d\theta$$

・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・

• Note that here the model is fixed and we want to learn the parameters (>< model selection)

# How to pick up the priors?

- The controversial aspect of Bayesian statistics are the priors
- The main argument of Bayesians is that we most often know something about the world
- When it is possible, it makes things easier to pick up a prior from the same family as the likelihood function

- Imagine that we have access to a set of obervations x<sub>i</sub> ∈ ℝ<sup>n</sup> and we can reasonably assume those samples are drawn independently from gaussian distributions.
- Because the observations are i.i.d, we can write the expression for the probability of oberving the  $x_i$  given the common  $\mu$  and  $\sigma^2$ ,

$$p(x|\mu,\sigma^2) = \prod_{n=1}^{N} \mathcal{N}(x_n|\mu,\sigma^2)$$
(2)

・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・

• A reasonably good idea to derive estimates for  $\mu$  and  $\sigma$  is then to maximize this likelihood

• Since the log is a monotonically increasing function,

$$\underset{\mu,\sigma^{2}}{\operatorname{argmax}} \quad p(x|\mu,\sigma^{2}) = \prod_{n=1}^{N} \mathcal{N}(x_{n}|\mu,\sigma^{2})$$

is equivalent to maximizing the log likelihood function

$$\begin{aligned} \operatorname*{argmax}_{\mu,\sigma^2} \quad \log\left(p(x|\mu,\sigma^2)\right) &= -\frac{1}{2\sigma^2}\sum_{n=1}^N (x_n - \mu)^2 \\ &\quad -\frac{N}{2}\log(\sigma^2) - \frac{N}{2}\log(2\pi). \end{aligned}$$

▲□▶ ▲□▶ ▲ 三▶ ▲ 三▶ 三 のへぐ

• Maximizing the log likelihood function with respect to  $\mu$  first and then  $\sigma^2$  gives the maximum likelihood estimators

$$\hat{\mu}_{ML} = \frac{1}{N} \sum_{n=1}^{N} x_n$$
$$\hat{\sigma}_{ML}^2 = \frac{1}{N} \sum_{n=1}^{N} (x_n - \hat{\mu}_{ML})^2$$

• Those two estimates are functions of the data set  $x_1, \ldots, x_N$ 

・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・

• Remember the ML estimators

$$\hat{\mu}_{ML} = \frac{1}{N} \sum_{n=1}^{N} x_n$$
$$\hat{\sigma}_{ML} = \frac{1}{N} \sum_{n=1}^{N} (x_n - \hat{\mu}_{ML})^2$$

 Now take the expectation of those estimators with respect to the known distribution,

$$\mathbb{E}\hat{\mu}_{ML} = \mu$$
$$\mathbb{E}\hat{\sigma}_{ML}^2 = \left(\frac{N-1}{N}\right)\sigma^2$$

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 のへぐ

• On average, the MLE will get you the right mean, but will underestimate the variance

$$\mathbb{E}\hat{\mu}_{ML} = \mu$$
$$\mathbb{E}\hat{\sigma}_{ML}^2 = \left(\frac{N-1}{N}\right)\sigma^2$$

- This problem is called bias and is related to the problem of overfitting
- In fact it turns out that a better estimator for  $\sigma^2$  is given by

$$\hat{\sigma}^2 = \frac{1}{N-1} \sum_{n=1}^{N} (x_n - \hat{\mu}_{ML})^2$$