## COMP 204: Python programming for life sciences Introduction to machine learning

Yue Li

イロン イロン イヨン イヨン 三日

1/46

## Outline

#### Introduction to machine learning

Supervised learning Classification Regression

Evaluation of machine learning algorithms

Unsupervised learning Clustering Latent topic models

Reinforcement learning

Python scikit-learn module

# What is machine learning?

Machine learning (ML) is a subset of artificial intelligence (AI). However, the line between ML and AI is becoming blurred. What AI are:

More generic on building automated intelligent agents that can both learn from external environment and also "think internally" (not well defined)

What ML are:

- tools that allow us to perform tasks that are hard to solve by traditional programming
- ways to represent complex structures in the massive amount data by simple interpretable patterns
- data-driven (as opposed to rule-based), leading to novel scientific discoveries

What ML are *not*:

- Computational neuroscience (use math to study brains)
- Rule-based system involving human-crafted rules
- Terminator (neither is AI)

Problems that cannot be solved by typical programs 1

- Typically, programs are written by human in exact sequential order (what we have seen so far)
- However, it is hard to write programs that solve tasks like recognizing faces or speech:
  - We don't know how to program it because we don't know how our brains work
  - Even if we had a good idea how our brain do it, the program will be horrendously complicated



### Problems that cannot be solved by typical programs 2

Detecting disease-causing mutations

- We don't know how to program it because we don't fully understand the functions of our genome
- We have very limited understanding of the physiology underlying most of the complex phenotypes (e.g. Alzheimer's disease, cancers) and how they interact with the environments (e.g., nutrition, exposed to radiation, neighbourhoods)
- There are unknown causal factors that we may not even observe or not yet have a way to measure them (e.g., uncharacterized pathways)



### Problems that cannot be solved by typical programs 3

Building a clinical recommender system:

- We have limited knowledge about how various diseases and symptoms are related to each other (despite ICD taxonomy and other efforts)
- How do we predict future disease onset based on current limited amount of health record for most people?
- How do we handle missing data that are common in healthcare (e.g., unordered lab tests)?



## Introduction to machine learning: a data-driven approach

What we do in machine learning:

- We collect lots of data (e.g., genotype of over 1 million genetic variants and phenotypes for large population cohort)
- We develop a machine learning algorithm that take these individual data as "examples" or "training data" and automatically produces a program that does the job
- The program produced by the learning algorithm may look very different from a typical hand-written program. It may contain millions of numbers (e.g., one decimal number per mutation indicating their impact on the phenotype)
- If we succeed, the program will work on new data ("testing data") that are not seen before (e.g., predicting risk of Alzheimer's disease using genotype of a new individual)
- In practice, we need two essential components:
  - 1. An objective function that quantitatively characterizes the problem (e.g., sum of errors in predicting diseases)

# 'Traditional' programming vs. machine learning

#### Traditional programming

- Program is written first independent of the data
- Program is applied to data to produce an output
- The program does not adapt to the data: it remains the same throughout its execution

#### Machine learning

 Program (or parameters of the program) adjusts itself automatically to **fit** the data

> End result is a program that is trained to achieve a given task



## Types of learning tasks

#### Supervised learning:

- Given examples of inputs (e.g., genotype) and corresponding desired outputs (e.g., disease), predict outputs on future unseen inputs, e.g., classification, regression, time series prediction
- Often the connotation of machine learning (people often ask how accurate is your model?)
- Unsupervised learning
  - Create a new representation of the input, e.g., form clusters, extract latent continuous features, compression
  - This is the new frontier of machine learning because most big datasets do not come with labels

#### Reinforcement learning

- Learn action to maximize payoff (e.g., robotics, self-driving vehicle)
- An important research area but not the focus of this class

A classic example in machine learning: recognizing hand-written digits What makes a "2"?

00011(1112 2222222333 344445555 467777888 888194999

#### First, how do we represent the data such as images: matrix

••••••	
••••••	
••••••	
••••••	
@@	
G	
00000000b	
*	
9 	
10	

We represent each  $28\times28$  image as a matrix or a two-dimensional list (see Lecture28.ipynb) Each entry in the 2D list is either 0 ('.') or number that is greater than 0 ('@')

```
from mnist import MNIST
import random
mndata = MNIST('mnist')
images, labels =
 \rightarrow mndata.load_training()
print(type(images)) # <class 'list'>
# display the first `2' in the data
for index,digit in enumerate(labels):
     if digit == 2:
         break
print(mndata.display(images[index]))
```

## Outline

Introduction to machine learning

Supervised learning Classification Regression

Evaluation of machine learning algorithms

Unsupervised learning Clustering Latent topic models

Reinforcement learning

Python scikit-learn module

## Classification

- Outputs or labels are categorical (1-of-N) (e.g., 10 digits)
- Inputs: feature (e.g., flattened 28 × 28 image = 784 input features)
- Goal: select the correct class for the new inputs (e.g., correctly classify the new image into one of the 10 digits)

How to represent labels? Use "one-hot encoding": create a vector as long as the number of categories we have, and set exactly one of the positions in the vector to 1 and the rest to 0.

# This is the one-hot version of: [5, 0, 4, 1, 9]

#### Classifying MNIST hand-written image with neural network

Classification of flattened  $28 \times 28$  input image (784 features) X into one of the 10 categories (digits) Y:



Simple math behind neural network (not required to understand in this class):

$$h_{j}^{(1)} = \sum_{i=1}^{784} w_{i}^{(1)} x_{i}; \qquad h_{j}^{(m)} = \sum_{i=1}^{H_{m}} w_{i}^{(m)} h_{i}^{(m-1)}$$

$$\hat{y}_{k} = \frac{\exp(-\sum_{i=1}^{M=128} w_{k} a_{i})}{\sum_{k'=1}^{10} \exp(-\sum_{i=1}^{M=128} w_{k'} a_{i})}; \quad CE = -y_{k} \log \hat{y}_{k} - (1 - y_{k}) \log(1 - \hat{y}_{k})$$

### Modern neural network computing infrainstructure

- Neural network also known as "Deep learning" models are typically implemented in a number of different libraries, including: Theano, Tensorflow, or Torch.
- These libraries typically have the following features:
  - Automatic or symbolic differentiation on computational graphs for backpropagation
  - Compilation for CUDA (GPU), enabling speedups due to highly optimized and parallel implementations of core NN functions
- High-level neural network libraries wraps on top of Theano and Tensorflow simplifies neural net creation (Keras, Lasagene)
- More powerful GPUs + Easily accessible cloud GPU computing (e.g., Amazon Web Services (AWS)) → More sophisticated and powerful models!









### Deep learning on classifying MNIST digits (8 lines of code)

```
import tensorflow as tf
1
   mnist = tf.keras.datasets.mnist
2
3
    (x_train, y_train),(x_test, y_test) = mnist load_data()
4
   x_train, x_test = x_train / 255.0, x_test / 255.0
\mathbf{5}
6
   model = tf.keras.models.Sequential([
7
      tf.keras.layers.Flatten(input_shape=(28, 28)),
8
      tf.keras.layers.Dense(128, activation='relu'),
9
      tf.keras.layers.Dropout(0.2),
10
      tf.keras.layers.Dense(10, activation='softmax')
11
12
   ])
13
   model.compile(optimizer='adam',
14
                   loss='sparse_categorical_crossentropy',
15
                  metrics=['accuracy'])
16
17
   model.fit(x_train, y_train, epochs=5)
18
   model.evaluate(x_test, y_test)
                                             (*ロト *個 * * ヨト * ヨト - ヨ
19
```

16 / 46

## Outline

Introduction to machine learning

Supervised learning Classification Regression

Evaluation of machine learning algorithms

Unsupervised learning Clustering Latent topic models

Reinforcement learning

Python scikit-learn module

### Regression

- Outputs are continuous
- Inputs: feature (continuous or discrete)
- Goal: predict outputs accurately for new inputs
- A toy example (see Lecture28.ipynb):
  - Output (y): a continuously measured variable
  - Input (x): another continuously measured variable
  - Task: fitting a line as y = wx + b by estimating w and b



18/46

Example: predict transcription factor binding affinity

Transcription factors (TF) are proteins that bind to specific regions of the genome to regulate nearby gene expression



Binding site

Goal: predict the TF binding affinity based on the DNA sequence Input: representing DNA sequence as 2D matrix:



Matrix representation of DNA sequence (darker = stronger)

## Convolutional neural network (CNN)

Applying 4 bp sequence filter along the DNA matrix:



Yellow = high activity; blue = low activity

### Predicting transcription factor binding affinity



### in-silico prediction of mutation impact [Alipanahi 2015]



≣ •⁄> ৭.ে 22 / 46

## Outline

Introduction to machine learning

Supervised learning Classification Regression

#### Evaluation of machine learning algorithms

Unsupervised learning Clustering Latent topic models

Reinforcement learning

Python scikit-learn module

### Evaluating machine learning algorithms

- How can we get an *unbiased* estimate of the accuracy for a learned model?
- Goal: Estimate accuracy of predictor on examples it has not seen as part of its training.

#### Training data vs Testing data

- split available data into training and testing datasets
- create a learned model from the training data
- measure accuracy of trained model by applying it to the testing data



24 / 46

## Outline

Introduction to machine learning

Supervised learning Classification Regression

Evaluation of machine learning algorithms

Unsupervised learning Clustering Latent topic models

Reinforcement learning

Python scikit-learn module

## Clustering

- Inputs: continuous or categorical data
- Goal: group data cases into a finite number of clusters so that within each cluster all cases have very similar
- One of the simplest algorithm is called: k-means clustering
- following slides are based on Bishop 2006 textbook: Pattern recognition and machine learning

Setting:

- data:  $\{\mathbf{x}_1, \ldots, \mathbf{x}_N\}$
- goal: partition the data into K clusters
- objective function in K-means:

$$J = \sum_{n=1}^{N} \sum_{k=1}^{K} r_{nk} ||\mathbf{x}_{n} - \boldsymbol{\mu}_{k}||^{2}$$
(1)

Algorithm:

- 1. initialize K cluster centers  $\mu_1, \ldots, \mu_K$
- 2. assign each point  $\mathbf{x}_n$  to the closest center k:

$$r_{nk} = egin{cases} 1 & ext{if } k = rgmin_j || \mathbf{x}_n - \boldsymbol{\mu}_j ||^2 \ & j \ 0 & ext{if otherwise} \end{cases}$$

3. update cluster centers:

$$\boldsymbol{\mu}_k = \frac{\sum_n r_{nk} \mathbf{x}_k}{\sum_n r_{nk}}$$

4. repeat 2 & 3 until convergence (i.e. little change from J (1))

Pattern recognition and machine learning (Bishop, 2006)  $\langle \Box \rangle \langle \overline{\Box} \rangle \langle \overline{\Box} \rangle \langle \overline{\Xi} \rangle \langle \overline{\Xi} \rangle \langle \overline{\Xi} \rangle \langle \overline{\Xi} \rangle$ 



















K-means clustering (K=2) convergence



◆□▶ ◆□▶ ◆三▶ ◆三▶ ○○ のへで

#### Clustering mutations based on epigenomic annotations



## Outline

Introduction to machine learning

Supervised learning Classification Regression

Evaluation of machine learning algorithms

Unsupervised learning Clustering Latent topic models

Reinforcement learning

Python scikit-learn module

### Latent topic models (from healthcare data)



40 / 46

# Latent Dirichlet Allocation (LDA) (Blei et al, JMLR 2003)



#### Grouping words by their topics (Blei et al, JMLR 2003) <u>"Arts"</u> "Budgets" "Children" "Education" NEW MILLION CHILDREN SCHOOL

1412144	MILLION	OILLIDIGIN	DOHOOL
FILM	TAX	WOMEN	STUDENTS
SHOW	PROGRAM	PEOPLE	SCHOOLS
MUSIC	BUDGET	CHILD	EDUCATION
MOVIE	BILLION	YEARS	TEACHERS
PLAY	FEDERAL	FAMILIES	HIGH
MUSICAL	YEAR	WORK	PUBLIC
BEST	SPENDING	PARENTS	TEACHER
ACTOR	NEW	SAYS	BENNETT
FIRST	STATE	FAMILY	MANIGAT
YORK	PLAN	WELFARE	NAMPHY
OPERA	MONEY	MEN	STATE
THEATER	PROGRAMS	PERCENT	PRESIDENT
ACTRESS	GOVERNMENT	CARE	ELEMENTARY
LOVE	CONGRESS	LIFE	HAITI

The William Randolph Hearst Foundation will give \$1.25 million to Lincoln Center, Metropolitan Opera Co., New York Philharmonic and Juilliard School. "Our board felt that we had a real opportunity to make a mark on the future of the performing arts with these grants an act every bit as important as our traditional areas of support in health, medical research, education and the social services," Hearst Foundation President Randolph A. Hearst said Monday in announcing the grants. Lincoln Center's share will be \$200,000 for its new building, which will house young artists and provide new public facilities. The Metropolitan Opera Co. and New York Philharmonic will receive \$400,000 each. The Juilliard School, where music and the performing arts are taught, will get \$250,000. The Hearst Foundation, a leading supporter of the Lincoln Center Consolidated Corporate Fund, will make its usual annual \$100,000 donation, too.

## Outline

Introduction to machine learning

Supervised learning Classification Regression

Evaluation of machine learning algorithms

Unsupervised learning Clustering Latent topic models

Reinforcement learning

Python scikit-learn module

### Reinforcement learning

Reinforcement learning (RL) learn action to maximize payoff (e.g., robotics, self-driving vehicle) Approach:

- We start with a robot that doesn't know how to walk, but moves its "muscles" randomly.
- The goal of the robot it to reach a certain destination (e.g. its "mother" at the other end of the room).
- When it reaches its mother, it gets a reward (satisfaction).
- Over time, it realizes that certain actions seem to lead to better rewards (reaching destination faster).

It slowly learns to adjust its behavior to maximize its reward A fun demo using RL:

https://www.youtube.com/watch?v=gn4nRCC9TwQ

## Outline

Introduction to machine learning

Supervised learning Classification Regression

Evaluation of machine learning algorithms

Unsupervised learning Clustering Latent topic models

Reinforcement learning

Python scikit-learn module

#### Python's scikit-learn module

Over the next 3 lectures

- we're going to perform some basic machine learning
- using Python's scikit-learn module

scikit-learn API:

http://scikit-learn.org/stable/modules/classes.html

scikit-learn tutorials:

http://scikit-learn.org/stable/