

Machine Learning

DA222

Soumitra Samanta
soumitra.samanta@gm.rkmvu.ac.in
Office: PB405

ML Class schedule

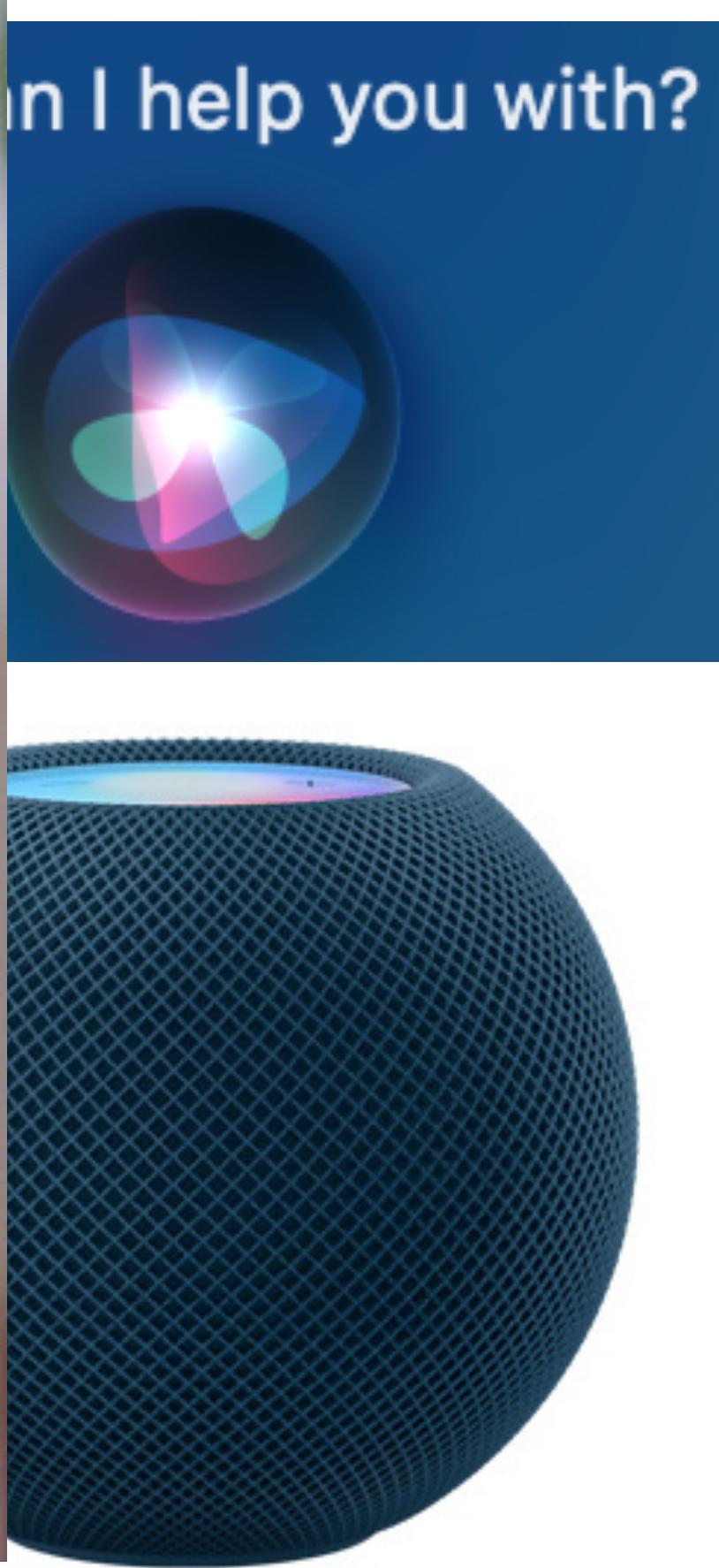
- Monday - 10:30 AM - 12:30 PM (Bhaskara lab)
- Saturday/Wednesday - 10:30 AM - 12:30 PM/12:30 PM - 2:30 PM (Bhaskara lab)
- Programming lab
- TA:
 - ▶ Rajdeep Mondal (PhD student)
 - Room no.- PB412
 - ▶ Suvajit Patra (PhD student)
 - Room no.- PB413



What is machine learning ?

Some applications and motivation

- Virtual assis-

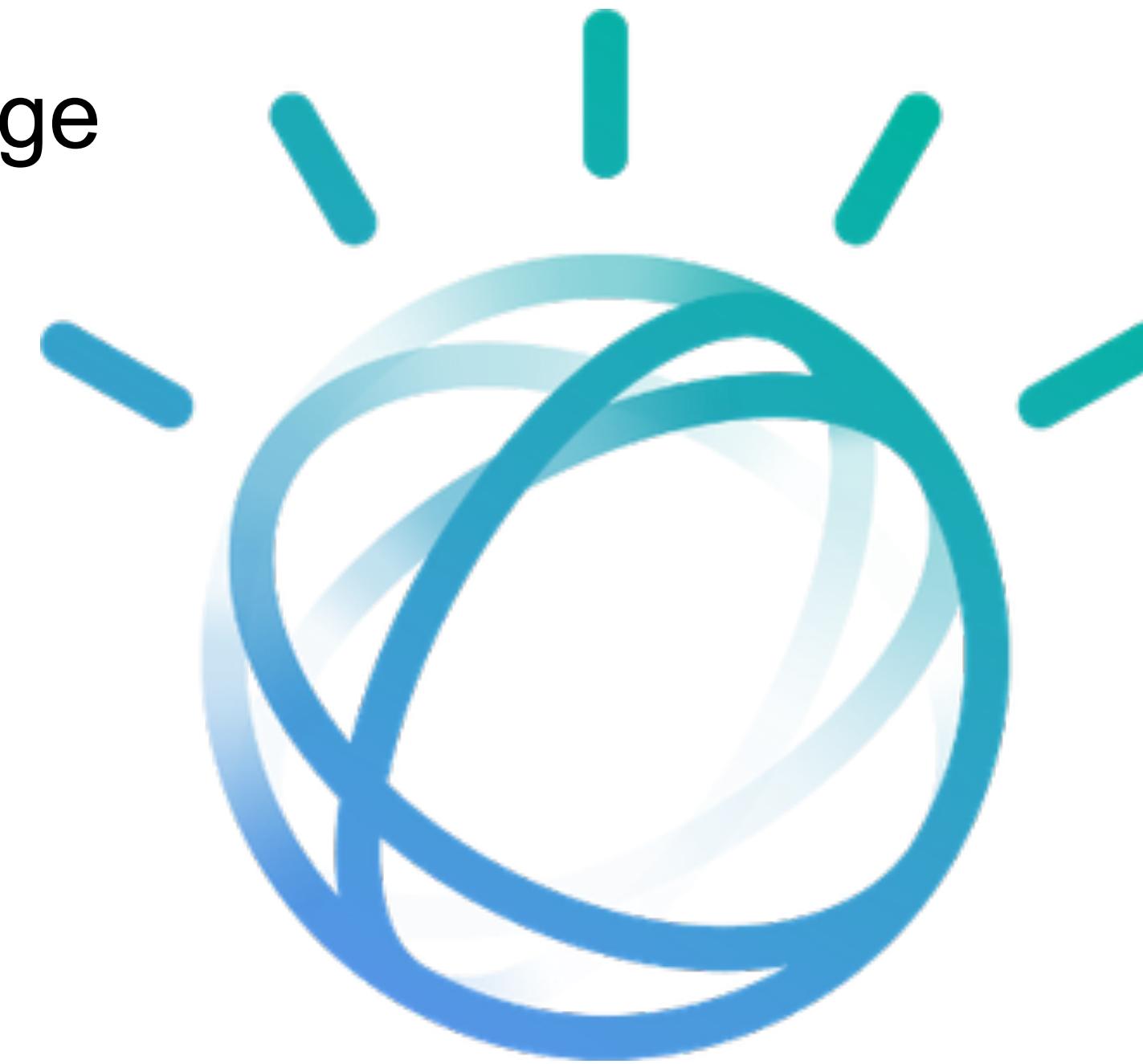


Images: Wikipedia, Amazon, Apple and Google

Video: YouTube

Language processing

- Question-answering system in natural language
 - ▶ IBM Watson
 - Won first prize in **Jeopardy**, 2011



Some applications and motivation (Cont.)

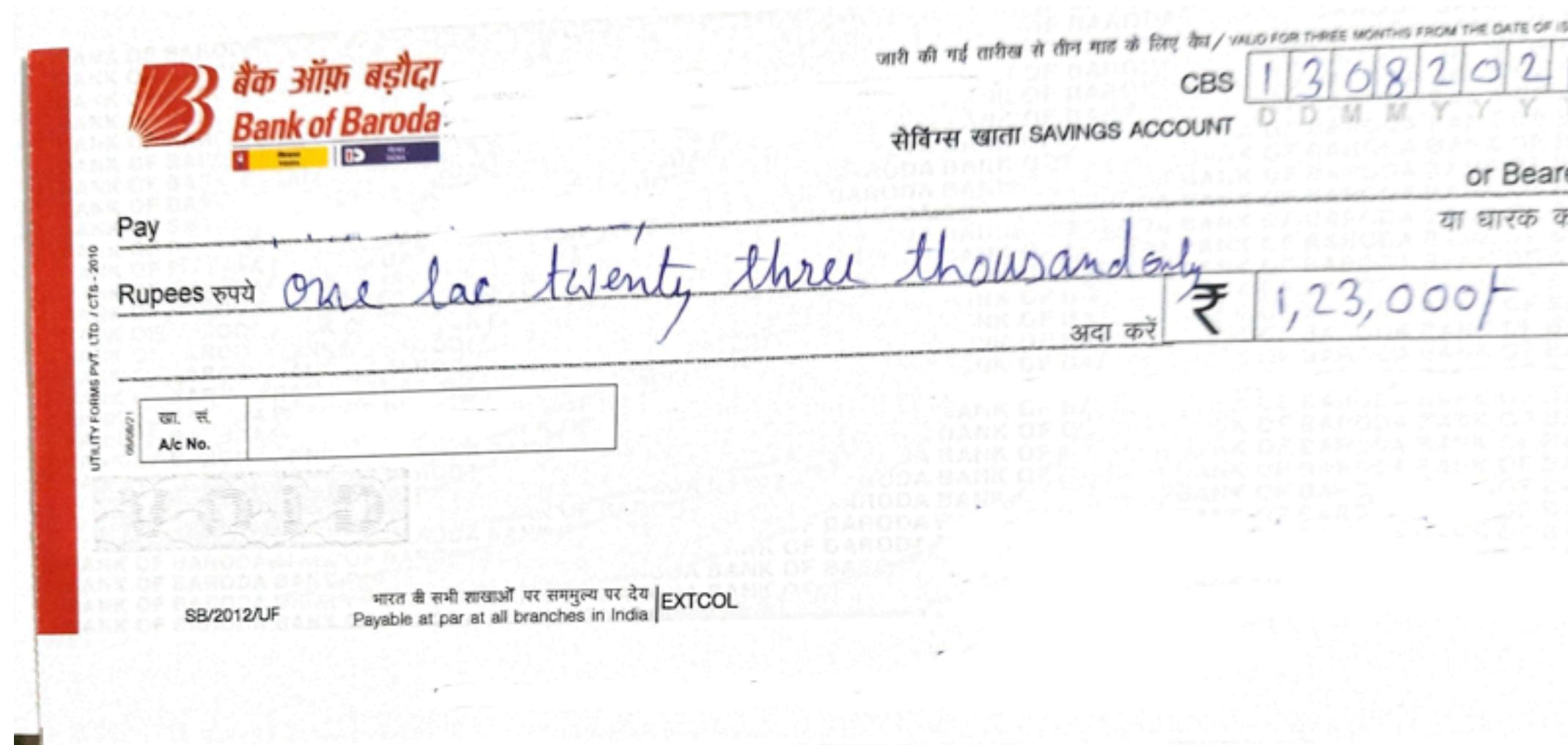


Hand-written character recognition

0,1,2,3,4,5,...

Person-1

Hello everyone, today is our first ML class and we will start from a motivating example.



0.1.2.3.4.5,...

Person-2

Hello everyone, today is our first ML class and we will start from a motivating example.

Hand-written digit recognition

- MNIST digit dataset¹
 - Image: 28×28



THE MNIST DATABASE of handwritten digits

[Yann LeCun](#), Courant Institute, NYU
[Corinna Cortes](#), Google Labs, New York
[Christopher J.C. Burges](#), Microsoft Research, Redmond

Please refrain from accessing these files from automated scripts with high frequency. Make copies!

The MNIST database of handwritten digits, available from this page, has a training set of 60,000 examples, and a test set of 10,000 examples. It is a subset of a larger set available from NIST. The digits have been size-normalized and centered in a fixed-size image.

It is a good database for people who want to try learning techniques and pattern recognition methods on real-world data while spending minimal efforts on preprocessing and formatting.

Four files are available on this site:

[train-images-idx3-ubyte.gz](#): training set images (9912422 bytes)
[train-labels-idx1-ubyte.gz](#): training set labels (28881 bytes)
[t10k-images-idx3-ubyte.gz](#): test set images (1648877 bytes)
[t10k-labels-idx1-ubyte.gz](#): test set labels (4542 bytes)

please note that your browser may uncompress these files without telling you. If the files you downloaded have a larger size than the above, they have been uncompressed by your browser. Simply rename them to remove the .gz extension. Some people have asked me "my application can't open your image files". These files are not in any standard image format. You have to write your own (very simple) program to read them. The file format is described at the bottom of this page.

The original black and white (bilevel) images from NIST were size normalized to fit in a 20x20 pixel box while preserving their aspect ratio. The resulting images contain grey levels as a result of the anti-aliasing technique used by the normalization algorithm. The images were centered in a 28x28 image by computing the center of mass of the pixels, and translating the image so as to position this point at the center of the 28x28 field.

- Extended MNIST²

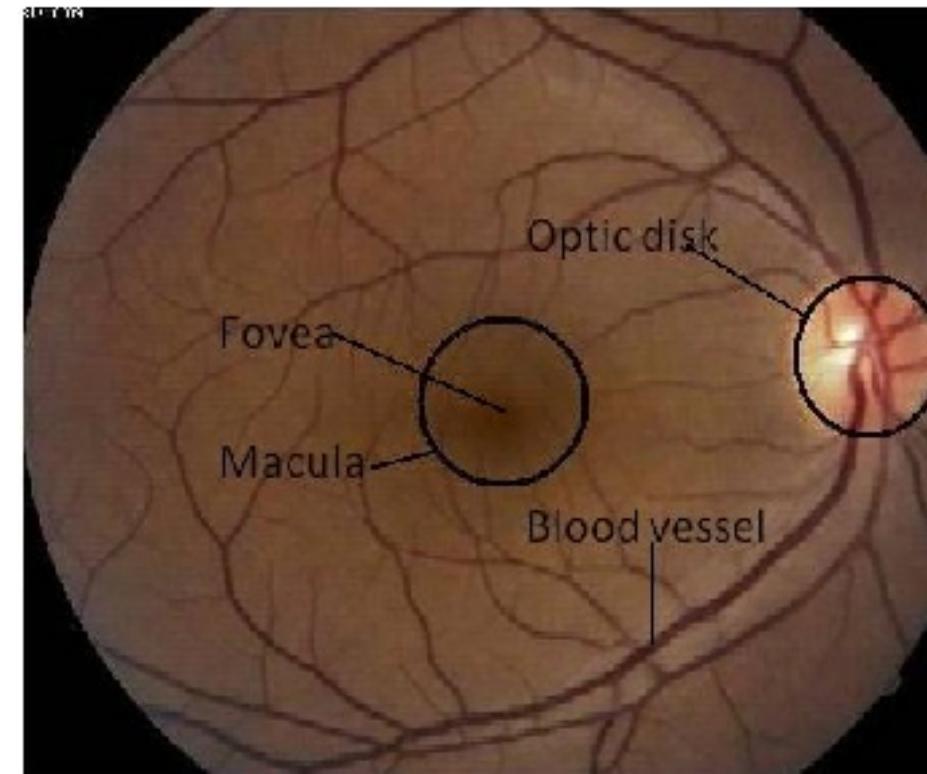
¹<http://yann.lecun.com/exdb/mnist/>

²<https://www.nist.gov/itl/products-and-services/emnist-dataset>

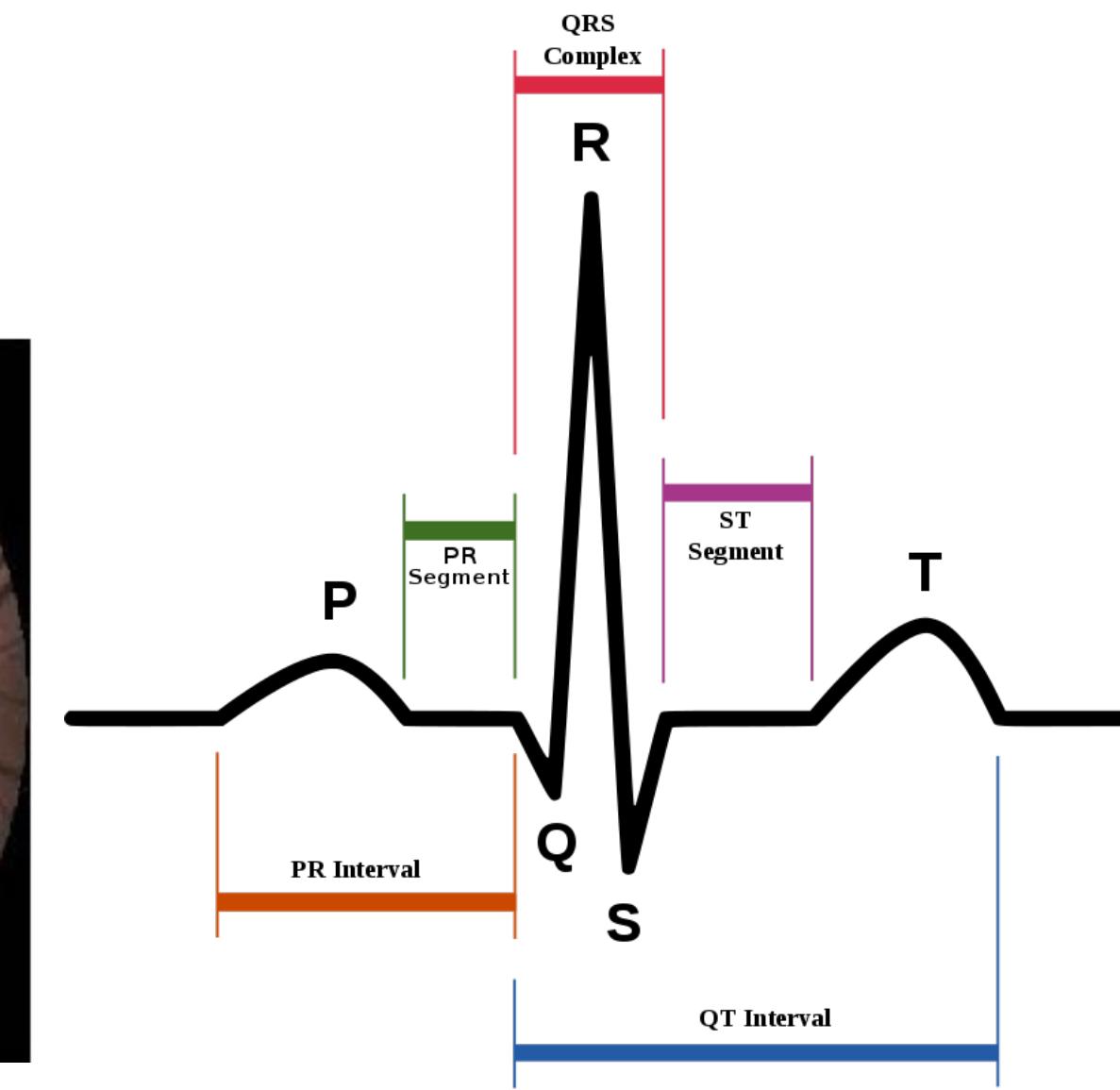
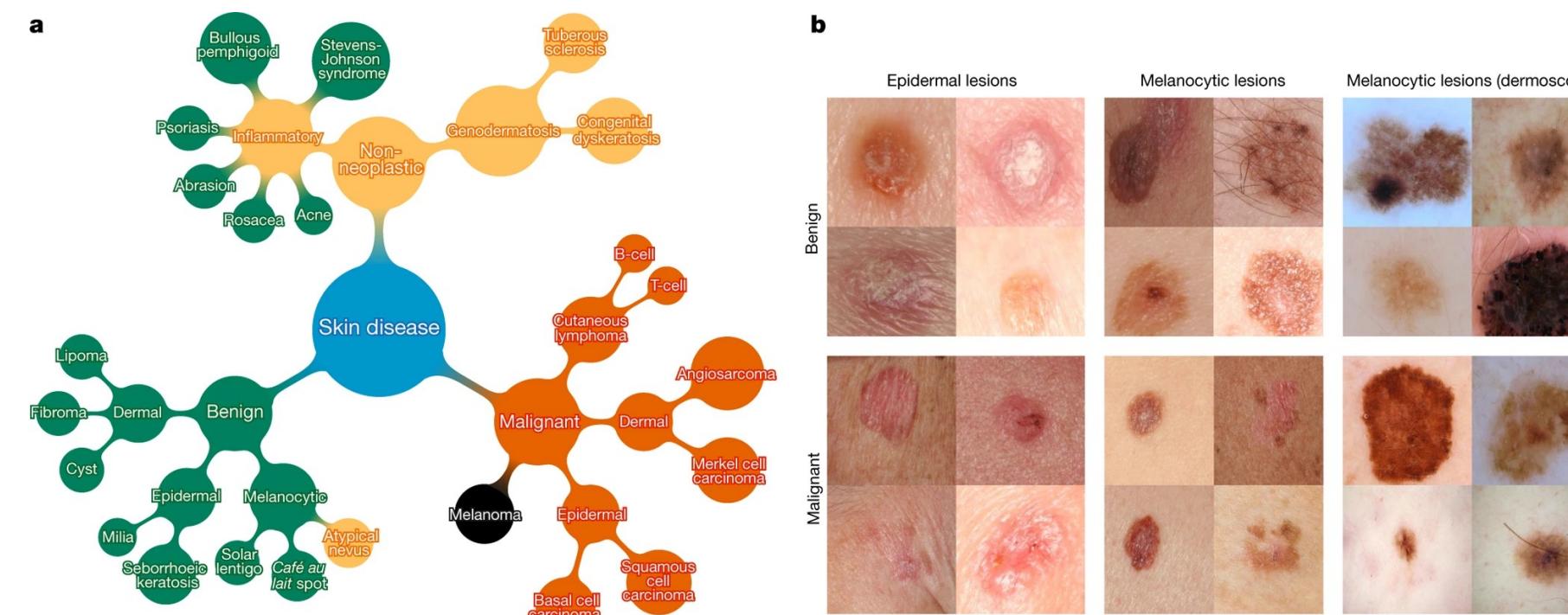
<https://cs.stanford.edu/people/karpathy/convnetjs/demo/mnist.html>

Computer aided diagnosis

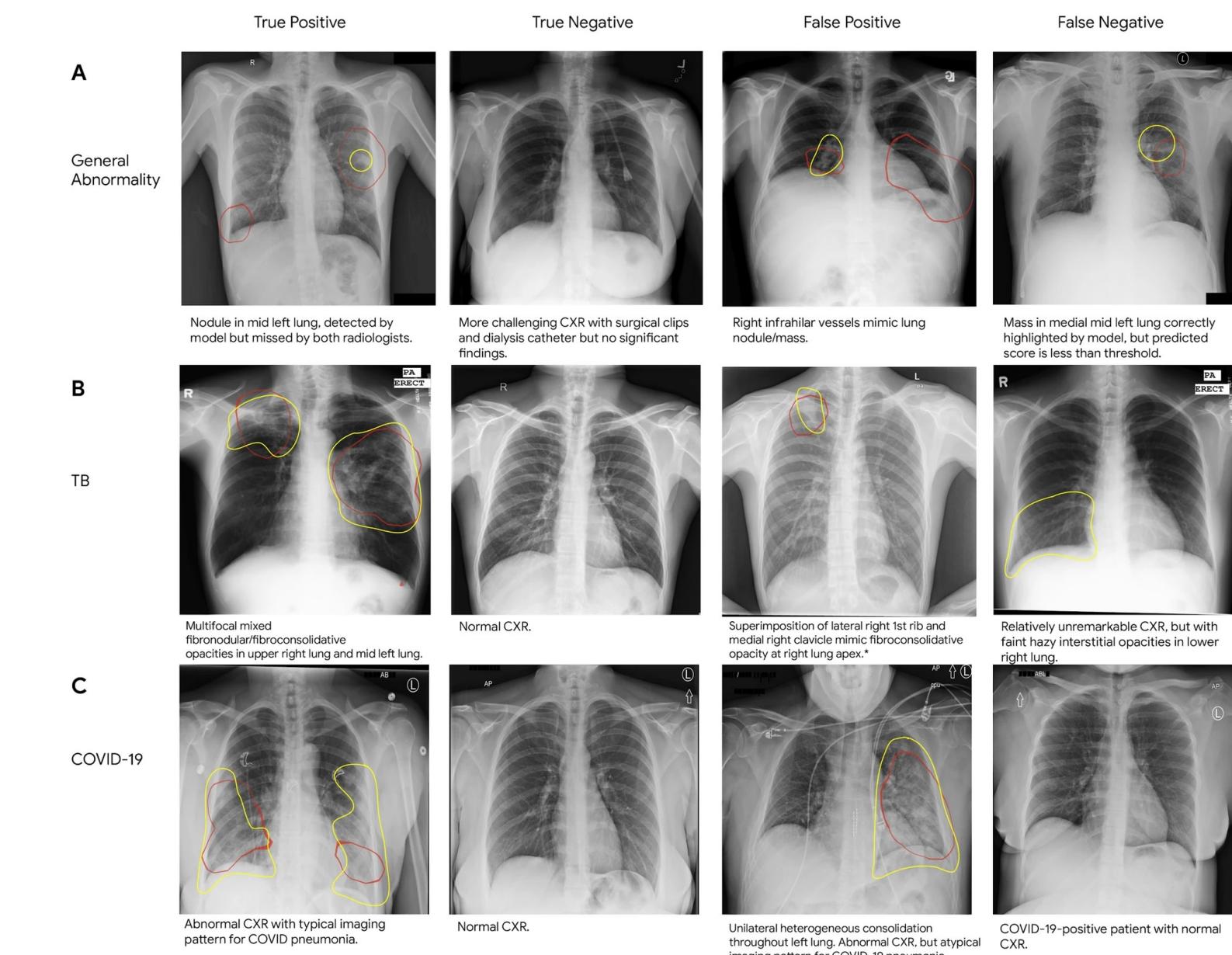
- Automation understanding of biosignals
 - ▶ EEG, ECG, EGG etc.



- Diabetic retinopathy
- Skin cancer detection

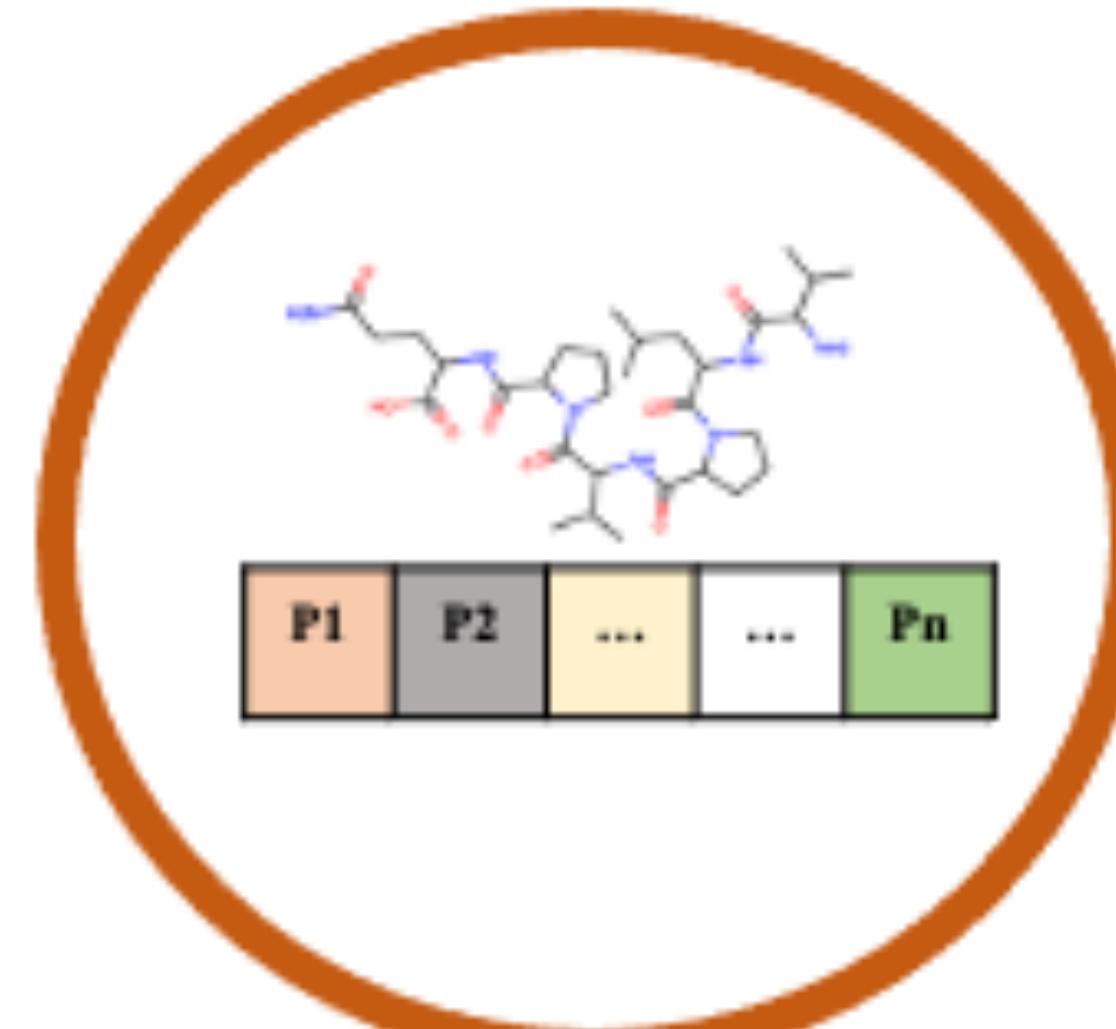


- Abnormal chest x-ray detection
- Designed to assist (**not replace**) physicians

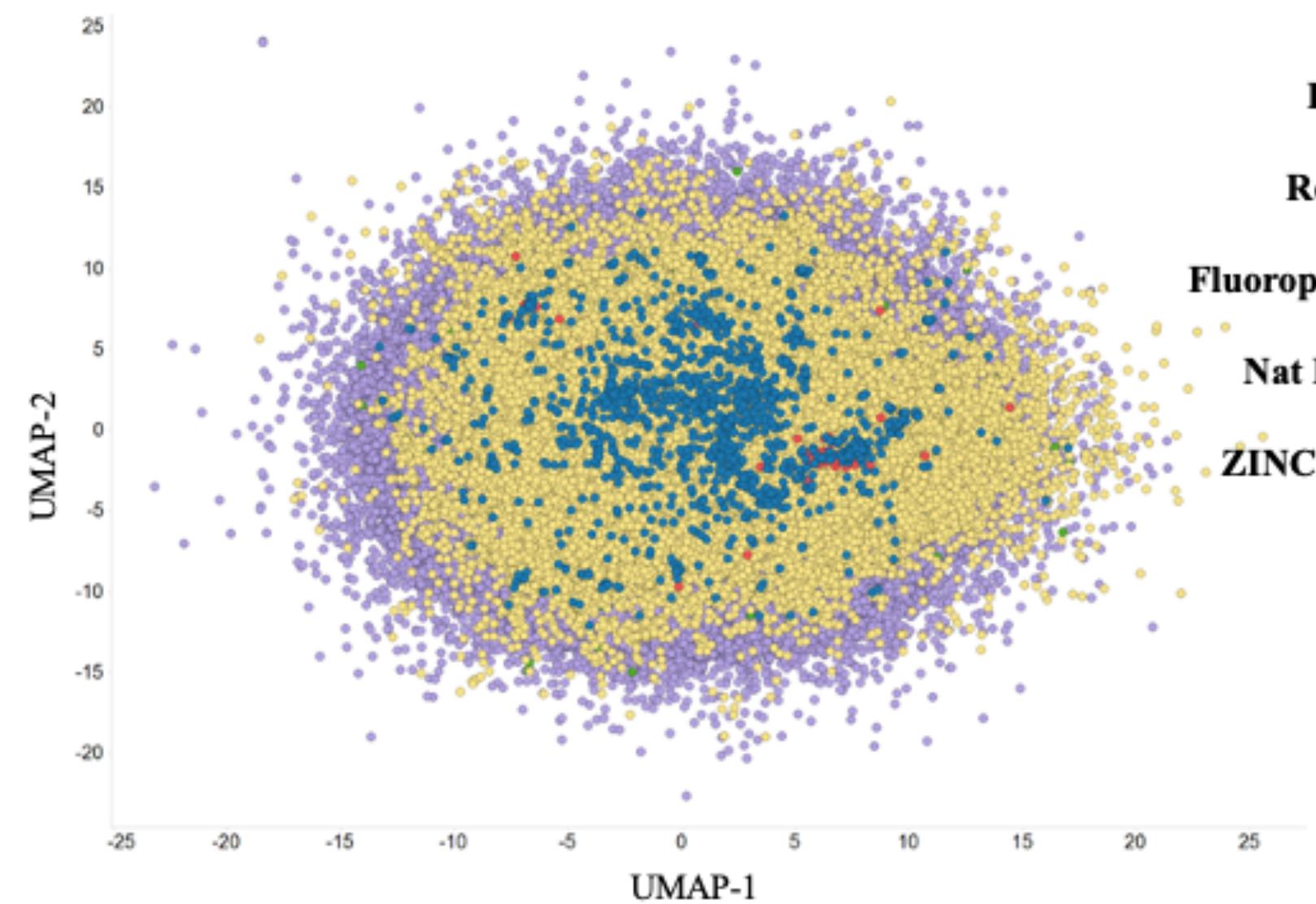


Drug discovery

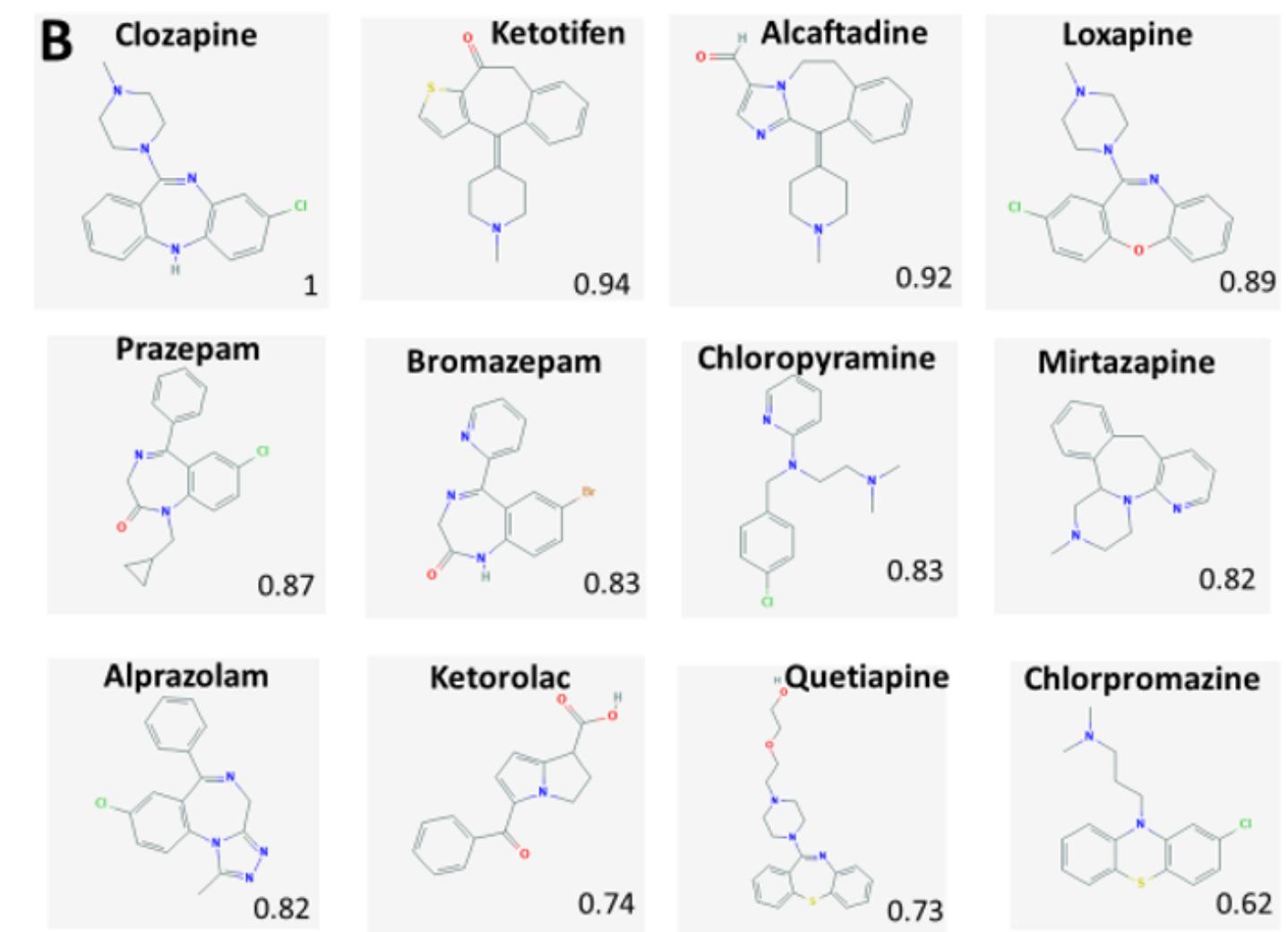
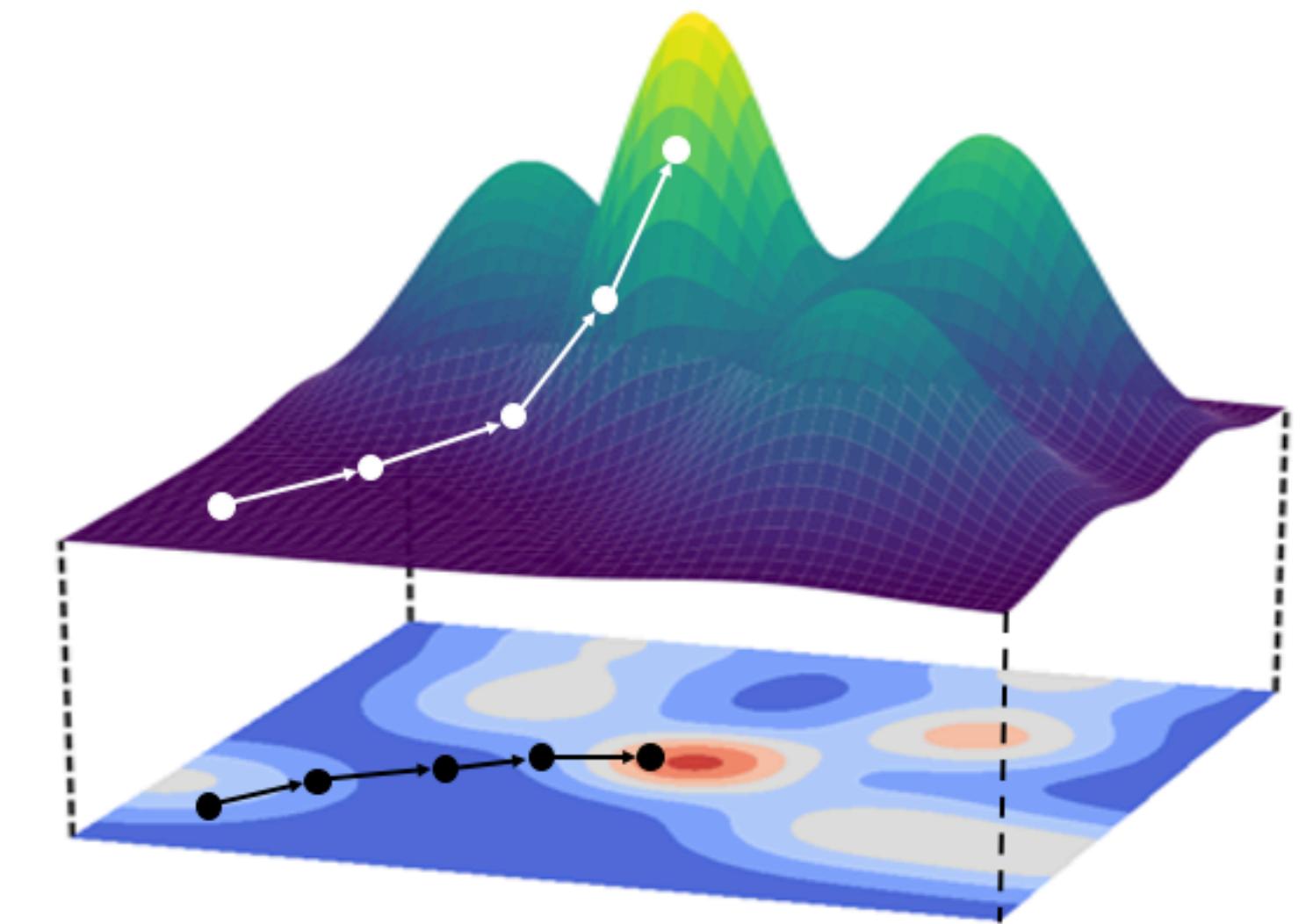
- Search (drug) molecule with desired properties



Nearest molecules (new) search?

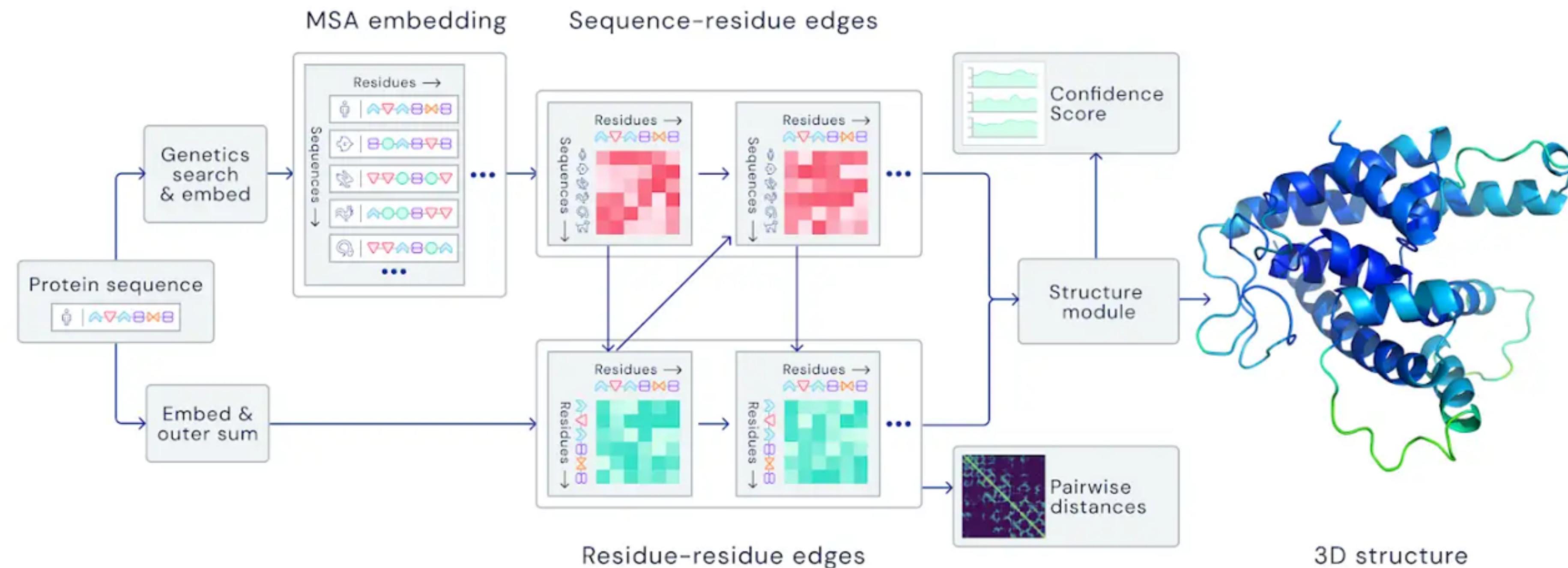


Drugs
Recon2
Fluorophores
Nat Prods
ZINC (6M)

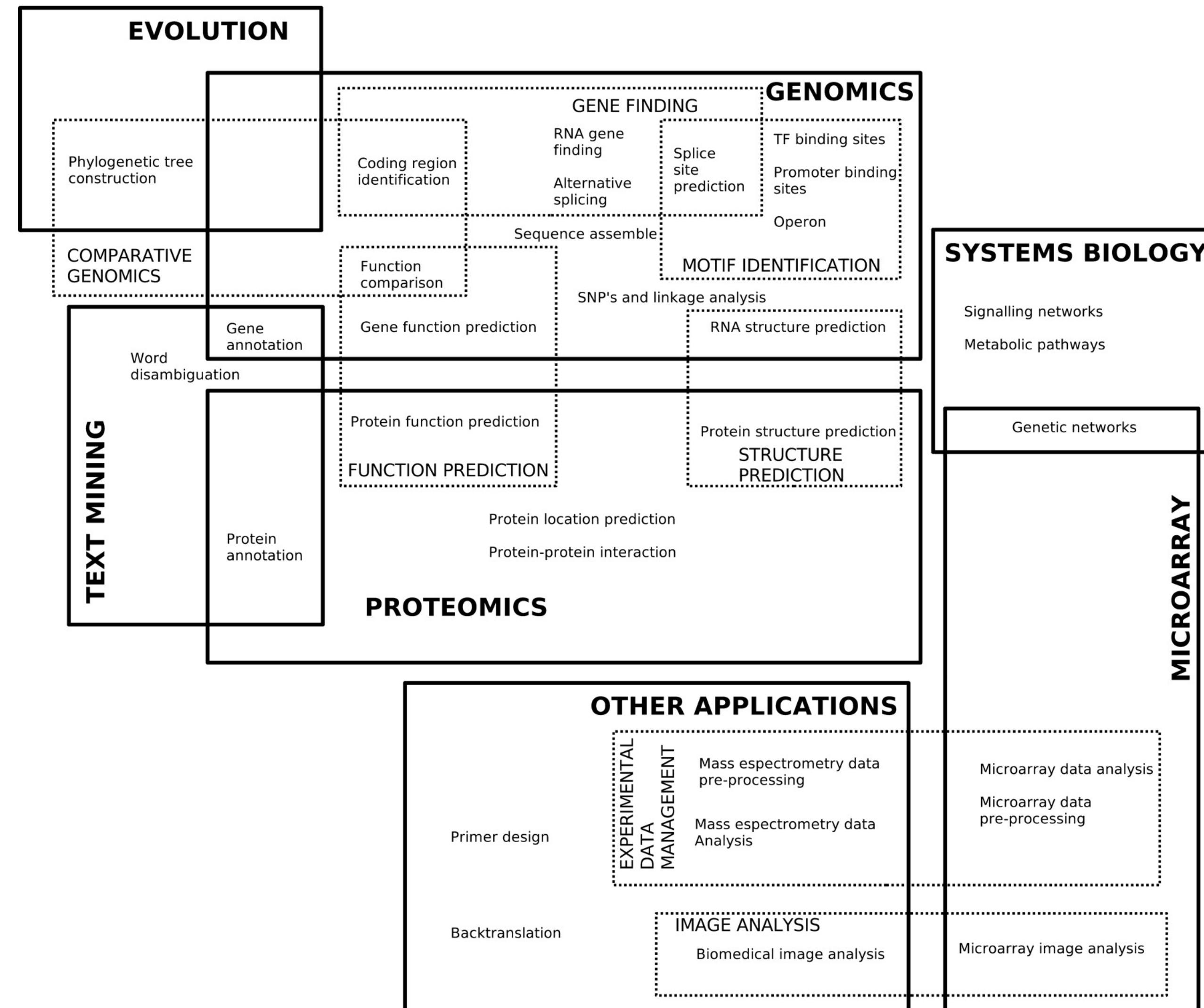


Protein-folding

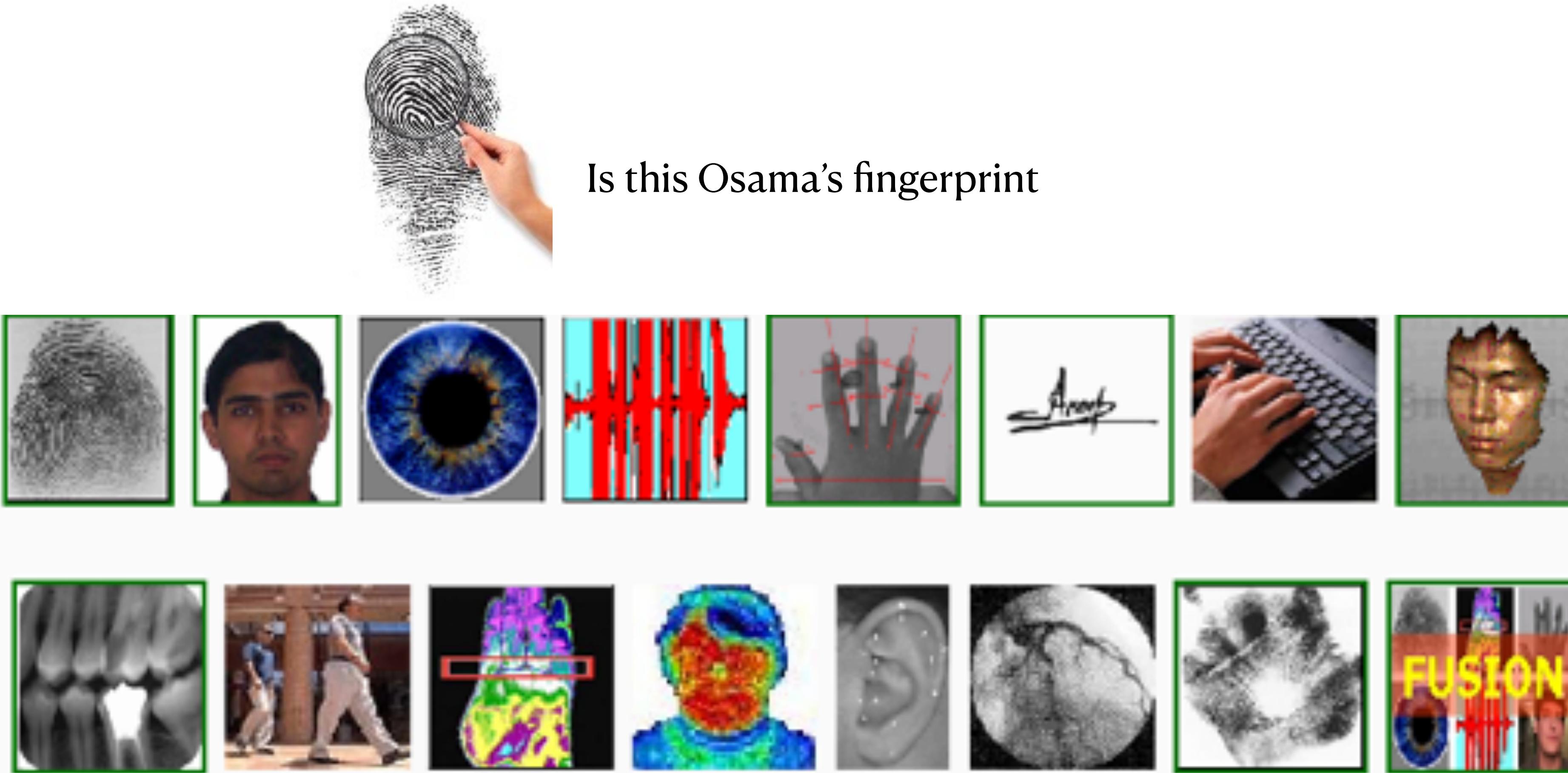
- Prediction of protein structure
 - ▶ CASP: <https://predictioncenter.org/casp15/index.cgi>
 - ▶ AlphaFold2 (DeepMind)- CASP-14, 2020



Bioinformatics

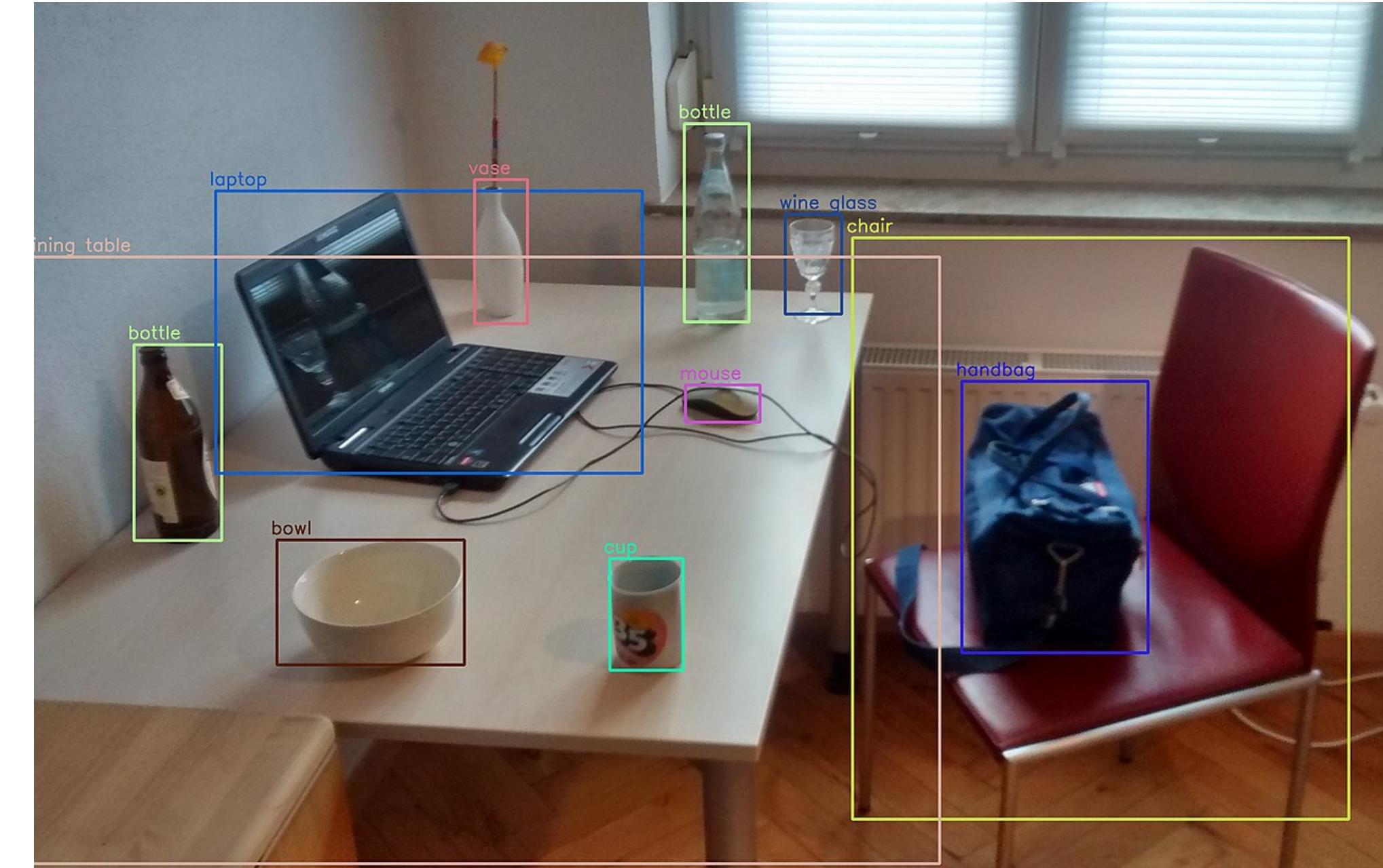
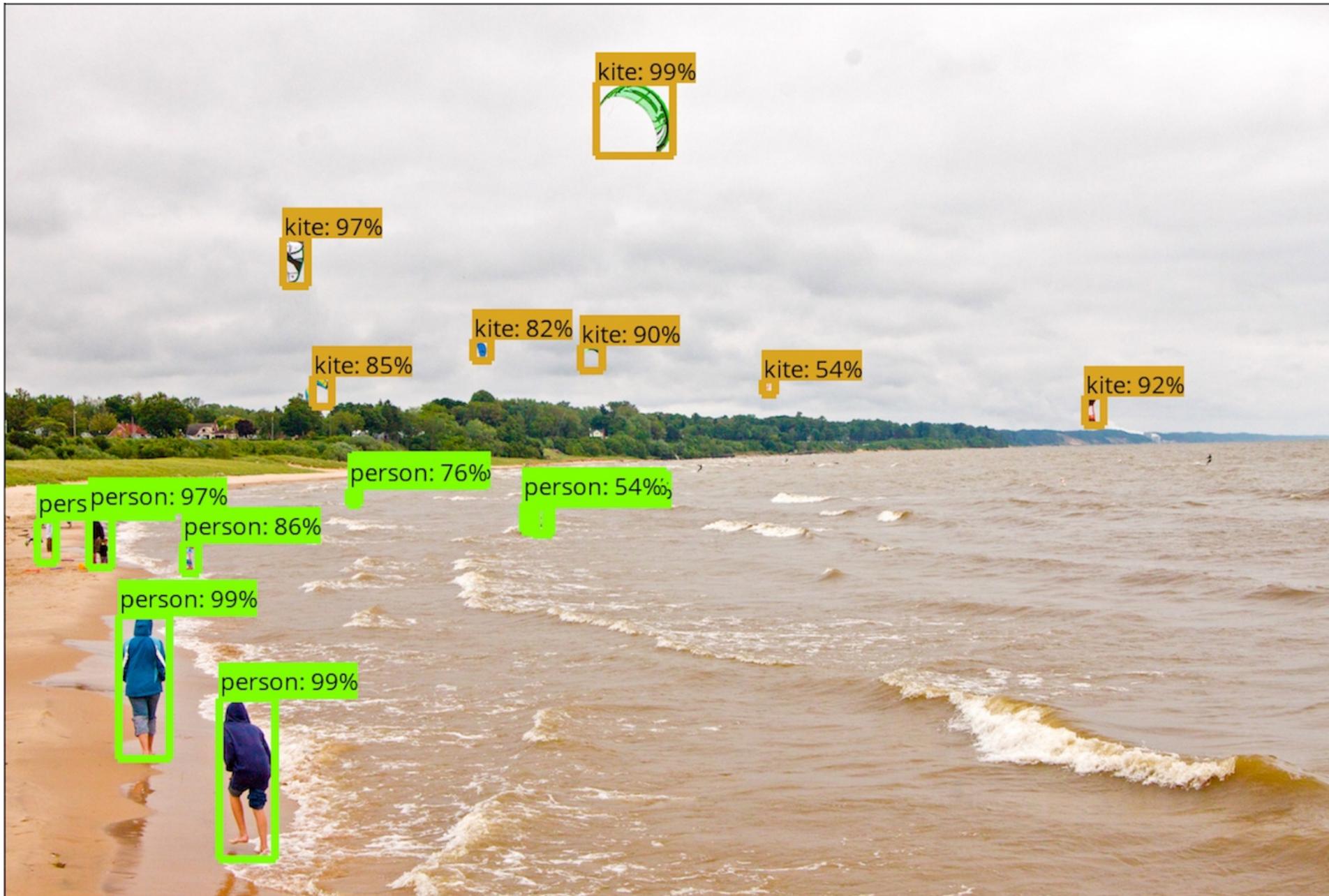


Biometric recognition



Multimodal biometric

Object detection



Started in 2005,
4-object class



Started in 2010, 1000-object class

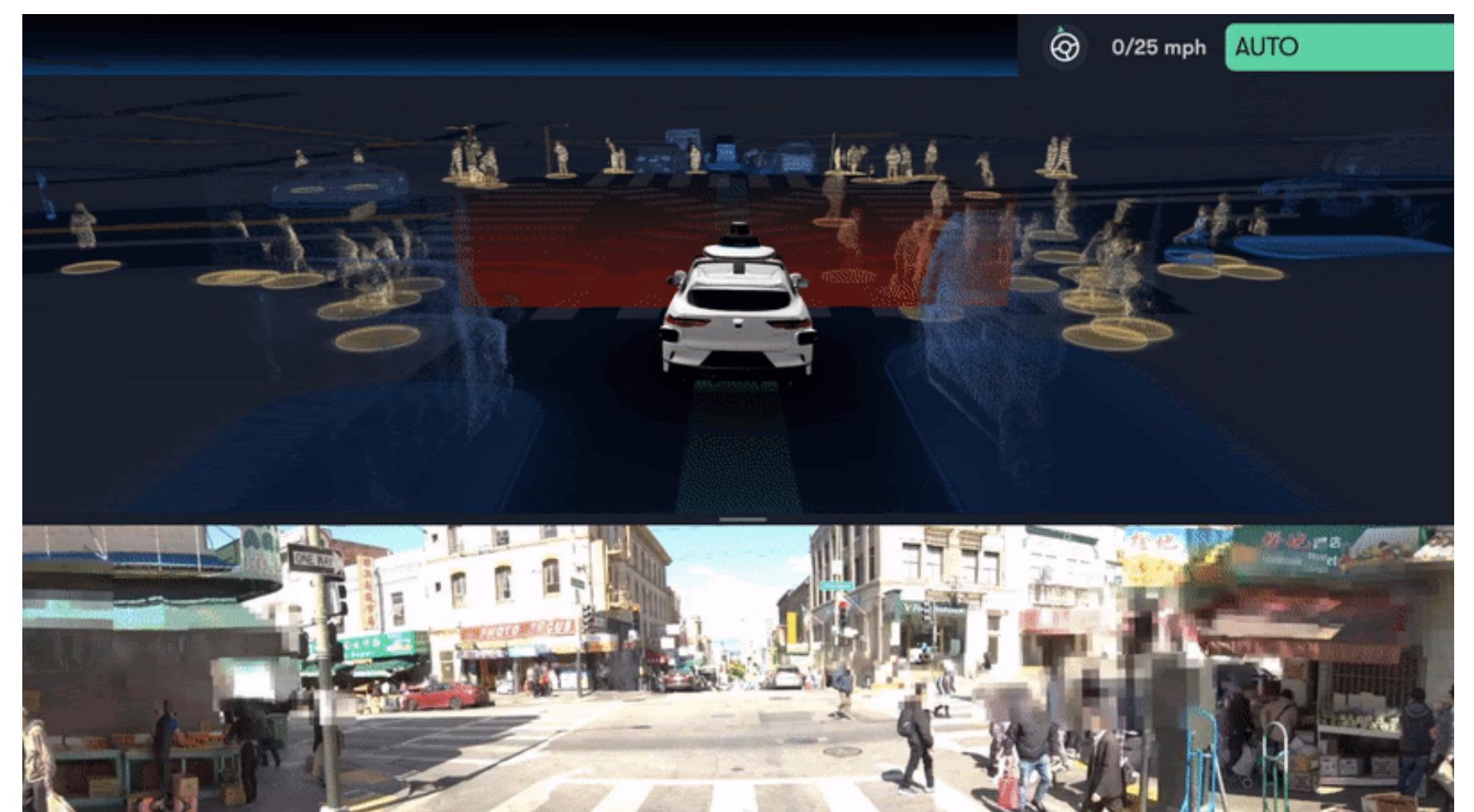


Self driving car

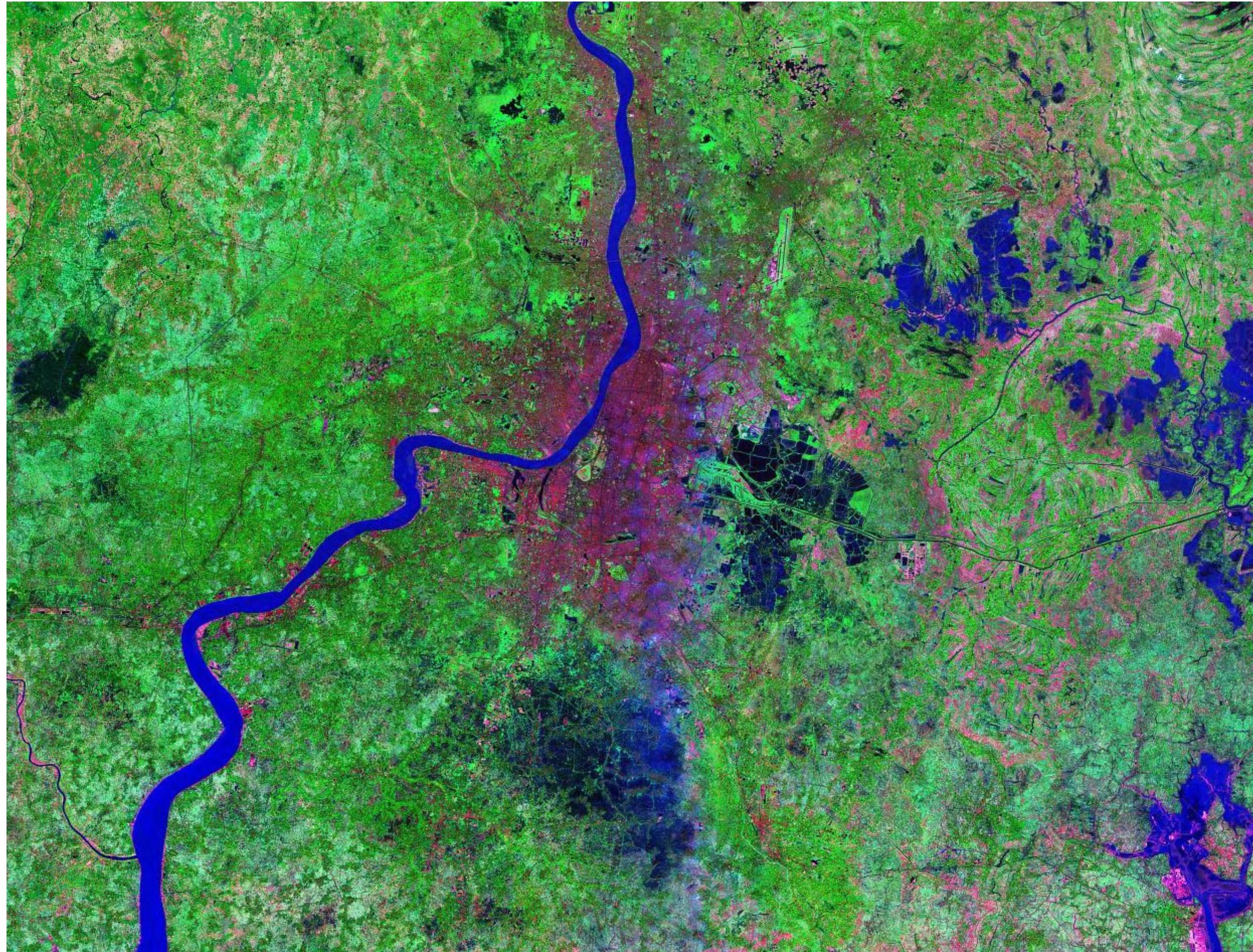
- DARPA grand challenge
 - ▶ 2005- Stanley



- Every car company jump into self driving car
 - ▶ Tata, BMW, GM, Nissan, Mercedes-Benz, etc.
- Big tech company
 - ▶ Google (Waymo), Apple, Uber, Tesla etc.



Satellite image analysis



Where are water body, concrete, and green area?

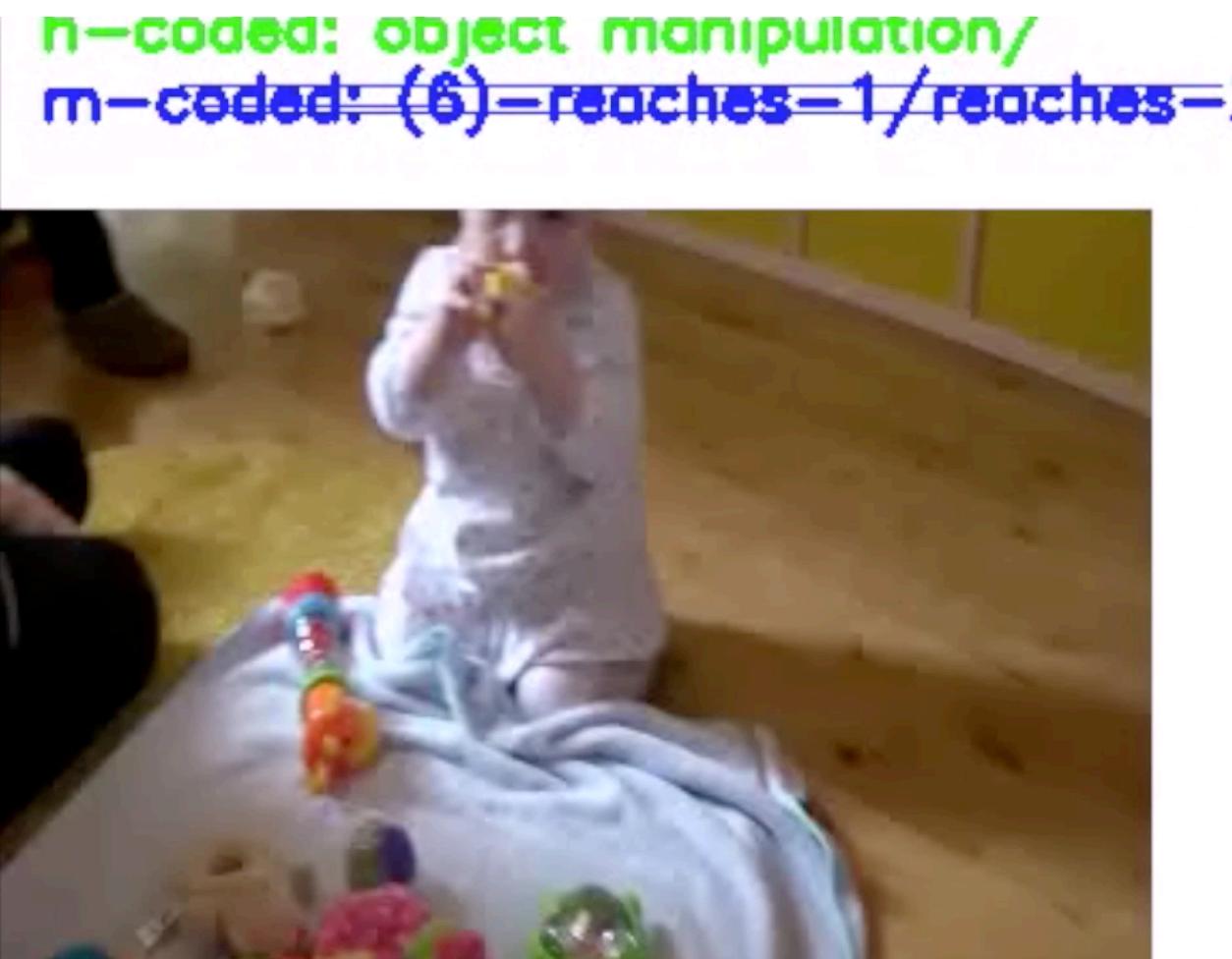
AlphaGo

- Computer program defeat professional human Go player
 - ▶ Beat
 - ▶ Beat world no 1 Go player Ke Jio, 2017



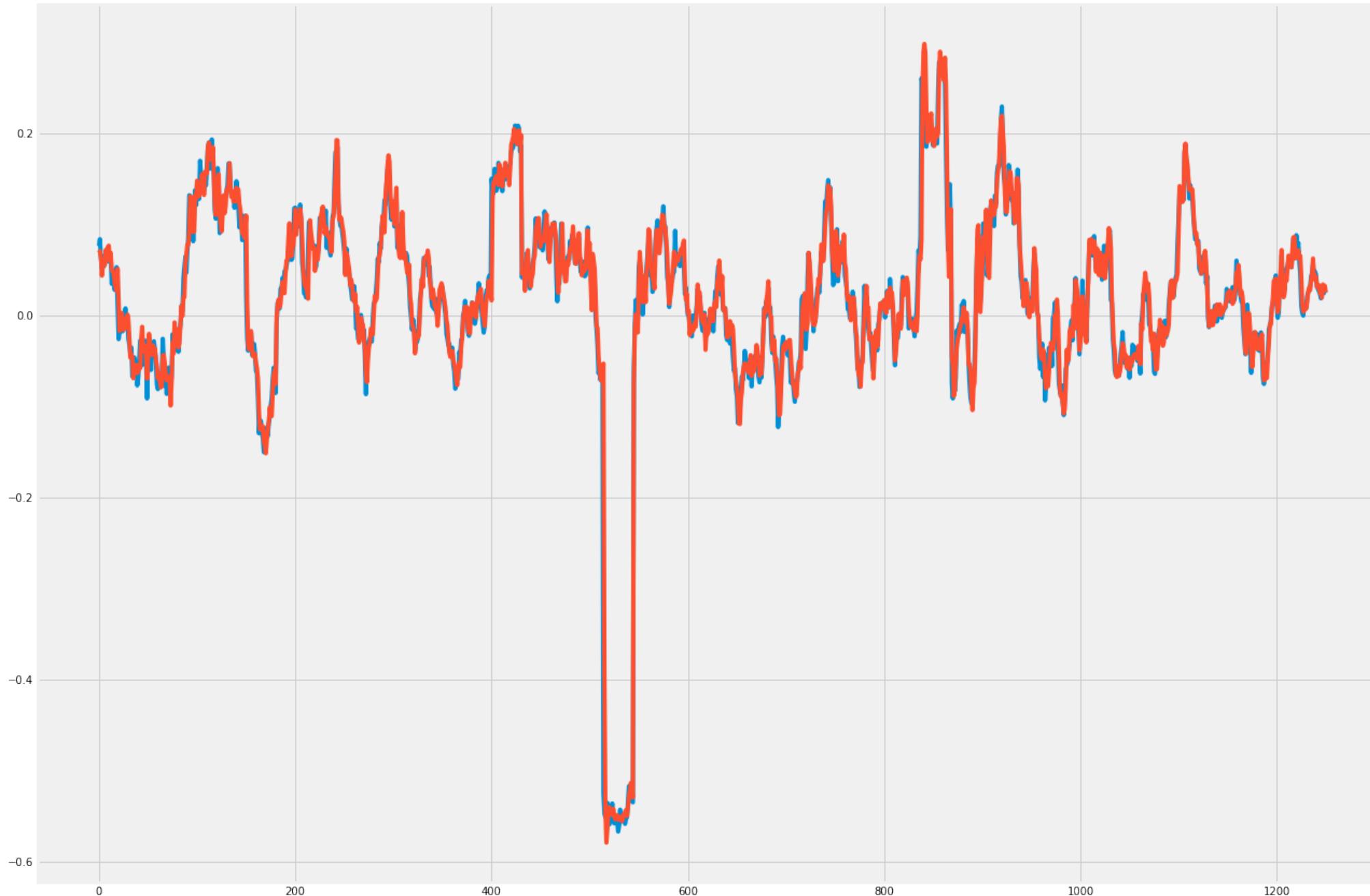
Image: DeepMind
Video: YouTube

Human action detection



Financial forecasting

- Stock market prediction

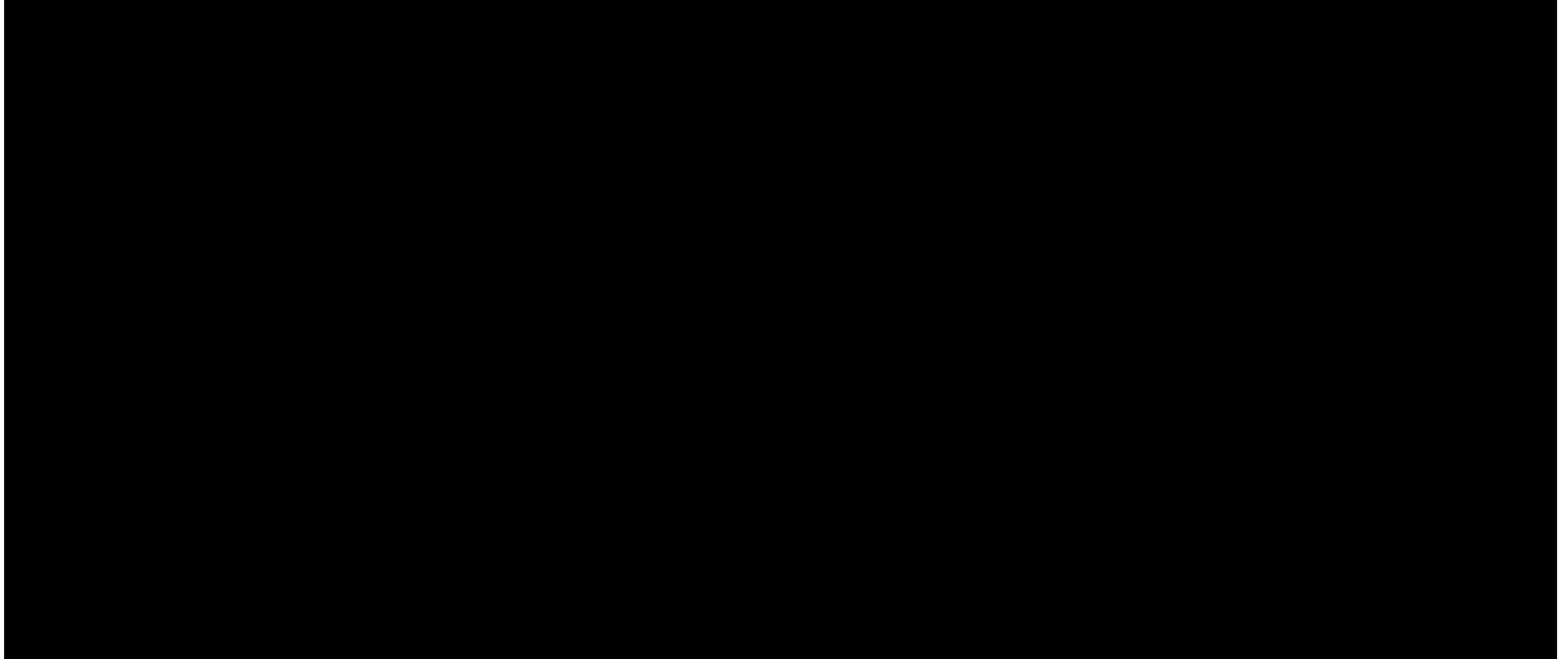


- Unwanted bank transaction (fraud) detection

Application of ML

- Applications
 - ▶ Image processing
 - ▶ Computer vision
 - ▶ Speech recognition
 - ▶ Natural language processing
 - ▶ Optical character recognition
 - ▶ Financial forecast
 - ▶ Medical diagnosis
 - ▶ ...

What is ML?



Video: <https://www.youtube.com/watch?v=QHH3jSeDBLo>

What is machine learning?

- Extract (learn) meaningful information from the examples (training) and answer the query for unseen examples
- Specific task
- How can we allow computer to extract meaningful information from the examples?
 - ▶ Design general rules
 - ▶ Development of algorithms

Inference

deductive vs inductive

- Deductive
 - ▶ Process of reasoning from one or more statements (premises) to reach a logical conclusion
- Example-1:
 - ▶ Premises-1: All RKMVERI students are excellent
 - ▶ Premises-2: Amol study in the RKMVERI
 - ▶ Conclusion: Amol is an excellent student
- In general:
 - ▶ If $A \implies B$ and $B \implies C$
 - ▶ Then $A \implies C$
- If premises are correct, then the conclusion is certain

Inference

deductive vs inductive

- **Inductive**
 - ▶ Method of reasoning in which a body of observation is synthesised to come up with a general principle
- Example-1:
 - ▶ Very often, we drop lots of things
 - ▶ All the times, the things fall downwards, but not upwards.
 - ▶ We can conclude that, if we drop things, likely they always fall downwards
- The truth of the conclusion is probable, based on the evidence given so far
- So, **conclusion is not certain!**
 - ▶ Example?

What is machine learning?

- Extract (learn) meaningful information from the examples (training) and answer the query for unseen examples
- Specific task
- How can we allow computer to extract meaningful information from the examples?
 - ▶ Design general rules
 - ▶ Development of algorithms
- Automate the process of **inductive inference**

Course logistics

Course page in xlms

xlms Home Dashboard My courses

SS Edit mode

The screenshot shows the course page for "Machine Learning 24". The top navigation bar includes links for "Course", "Settings", "Participants", "Grades", "Reports", and "More". The main content area is divided into sections: "General" and "Introduction to Machine Learning". The "General" section contains a forum icon labeled "FORUM Announcements" and a file icon labeled "FILE Suggested reading" with a "Mark as done" button. The "Introduction to Machine Learning" section contains a file icon labeled "FILE Introduction to Machine Learning" with a "Mark as done" button. A question mark icon is located in the bottom right corner of the page.

Machine Learning 24

Course Settings Participants Grades Reports More

▼ General

FORUM Announcements

FILE Suggested reading

Mark as done

▼ Introduction to Machine Learning

FILE Introduction to Machine Learning

Mark as done

?

<https://xlms.rkmvu.ac.in/course/view.php?id=99>

Prerequisites

- Mathematics
 - ▶ Linear Algebra, Multivariate Calculus, Basis Optimisation and Basic probability
 - ▶ No worries, we will touch some background when we need
- Computer programming: Any one from C/C++/**Python (recommended for the class project and assignments)**/MATLAB/Octave
- Basic concept in Algorithms and Data Structure

Linear Algebra

- Vector space
 - ▶ Definition, Basis, Dimension, Eigen value and Eigen vectors etc.
- Matrix
 - ▶ Addition, Multiplication, Trace, Inverse etc.
 - ▶ Positive definite matrices, Singular value decomposition etc.

Multivariate Calculus

- Derivative, Partial derivative, Taylor series expansion, Chain rules etc.

Basic Optimisation

- Convex set, Convex hull, Convex function
- Gradient of a function, Hessian
- Constrained and Unconstrained optimisation problem, Optimality condition

Probability

- Definition, Random variables, Distribution function and their different variants
- Conditional probability, Independence, Expectation, Variance, Moments, Entropy
- Law of large numbers, Central limit theorem

Evaluation

- Mid-term: 15%
- End-term: 50%
- Assignments/class test: 15%
- Project: 20%

Assignments

- Mostly implementation of different ML algorithms
 - ▶ Python (preferable)/C/C++/MATLAB/Octave
- There might be one-two theoretical assignments
- Submission deadline is strict and extra/buffer two days after the deadline for the entire semester
 - ▶ We will consider 11.59PM as our day end

Projects

- Can be done in a group (max **two** students)
- Define your own project
- Submit a one page project proposal- deadline **24-02-2024?**
- Finished the work within the time-line
- Report submission
 - ▶ Submission deadline: **1st May, 2024**
 - ▶ We will consider 11:59PM as our day end
- Final presentation
 - ▶ **20 min** (divided into group members)/**poster presentation**
 - ▶ **20th May, 2024**

Previous year's projects

Project title

Predicting The Stock Market Using Contemporary Current Affairs

Assessing the Feasibility of Diagnosis of Pneumonia using Chest X-Ray Images

Brain Tumor Classification Based on its Presence & Position Using MRI Images

Self Supervised learning in image classification: Introducing Barlow Twins

Suspicious Activity Detection Using Surveillance Footage

WORD-LEVEL INDIAN SIGN LANGUAGE RECOGNITION FROM VIDEO

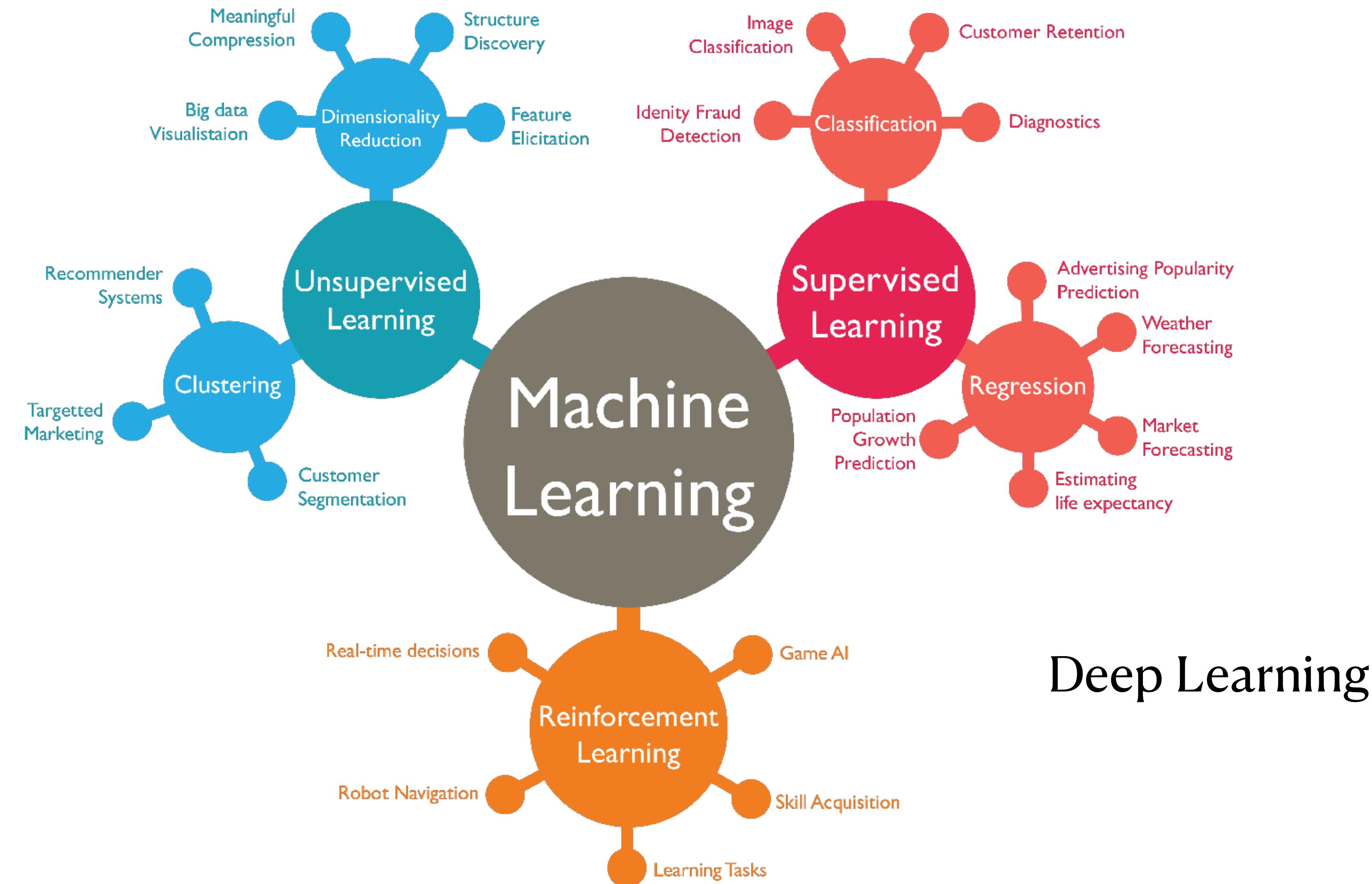
Audio Based Bird Species Classification

Anomaly Detection

FACIAL EMOTION RECOGNITION

Abstract Text Summarization

Course syllabus



Academic ethics

- Your grade should reflect your own work
- Copying or paraphrasing someone's work (code included), or permitting your own work to be copied or paraphrased, even if only in part, is **strictly forbidden**, and will result in an automatic grade of **zero** for the entire assignment or exam in which the copying or paraphrasing was done.
- So, **ask yourself** before copying from others
- If you are going to have trouble completing an assignment, talk to the instructor and TA before due date

Data

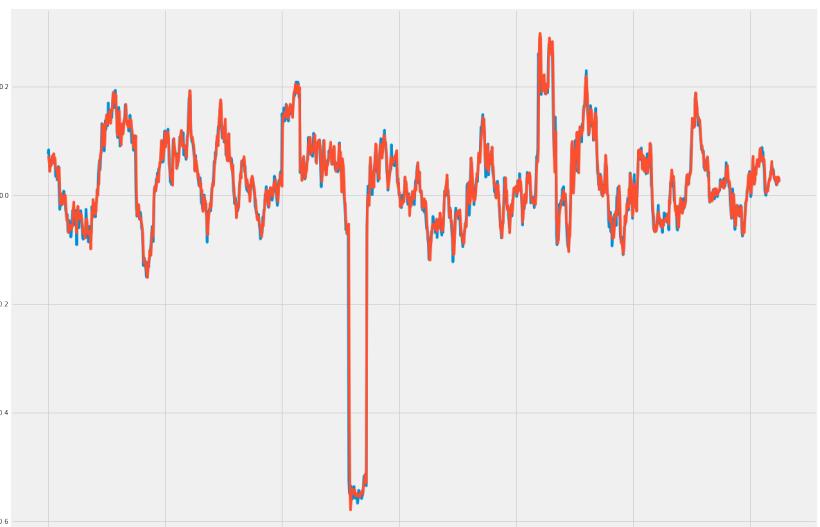
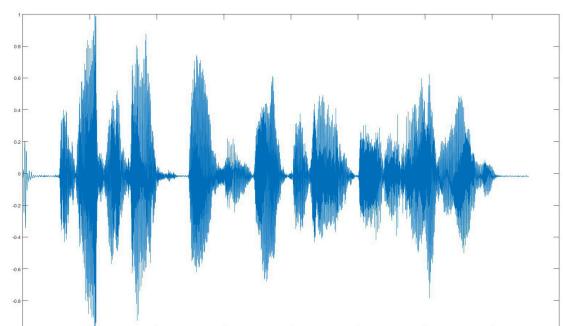
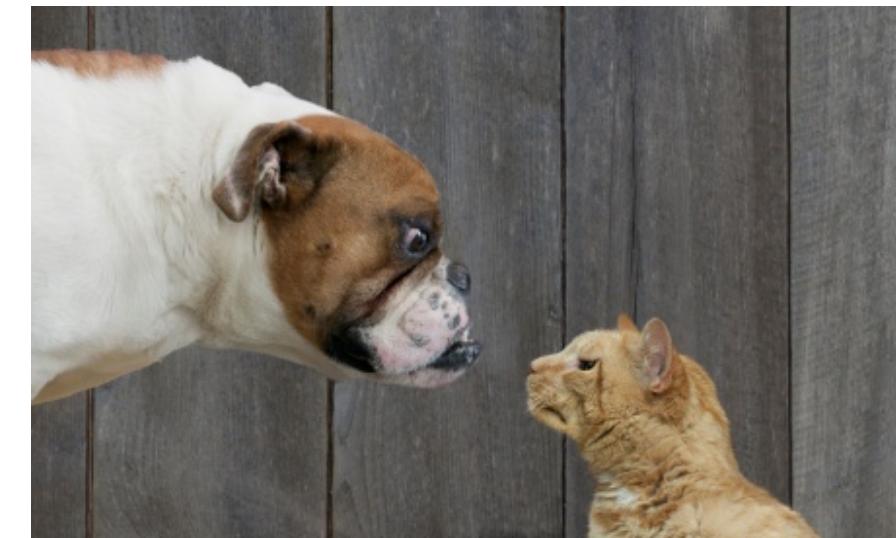
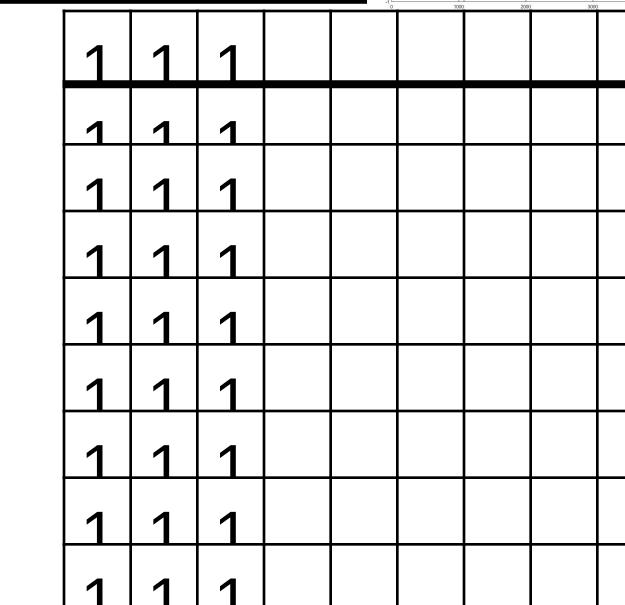
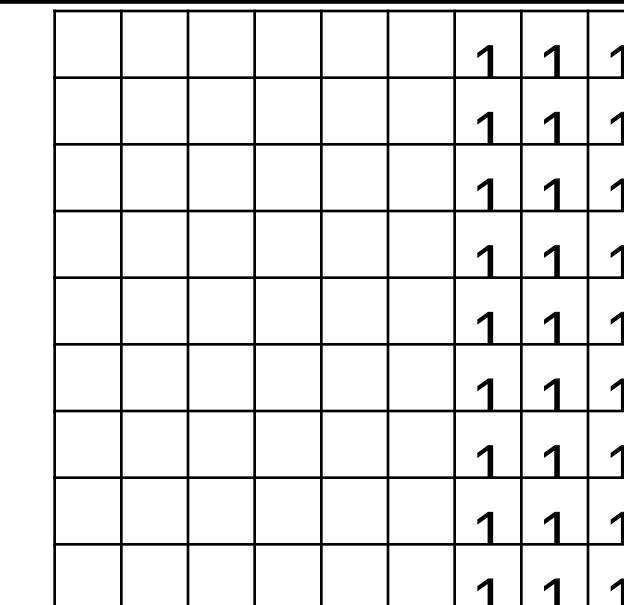
Recap

- What is ML?
- Some applications of ML
- Deductive vs Inductive inference

Data

- Different types of data
 - ▶ Text, Audio, Images, Videos, etc.
- Representation
 - ▶ Example:
 - Hand-written digit recognition
 - Cat vs Dog
 - Stock (particular) price prediction
 - ▶ Feature
 - Any distinct aspect, quality or characteristic
 - Ex. numeric (height)
 - Combination of d features is d -dims column vector, we will call it as **feature vector**
 - d -dims space defined by the feature vectors is called **feature space**
 - Data are represented as a **point** in the feature space

Hello everyone, today is our 1st ML class and we will start from a motivating example.

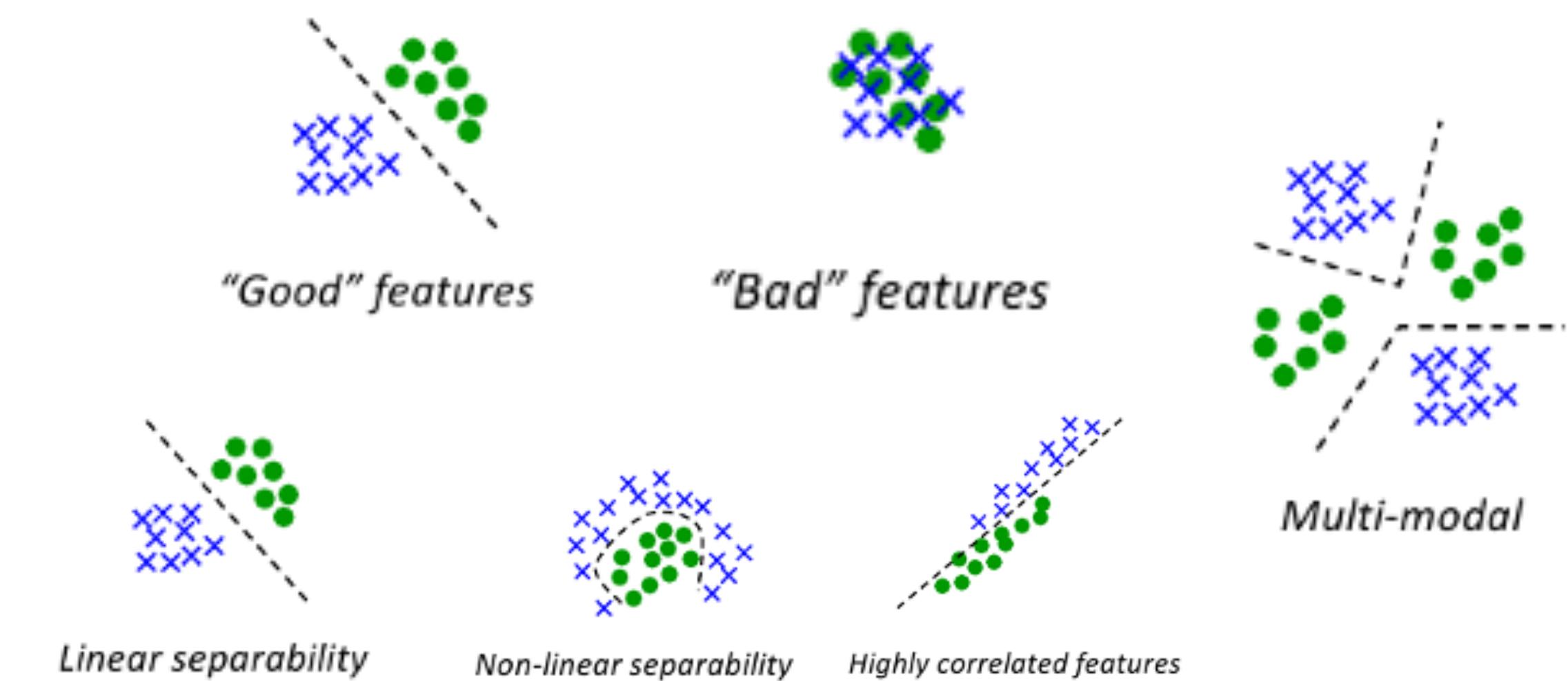
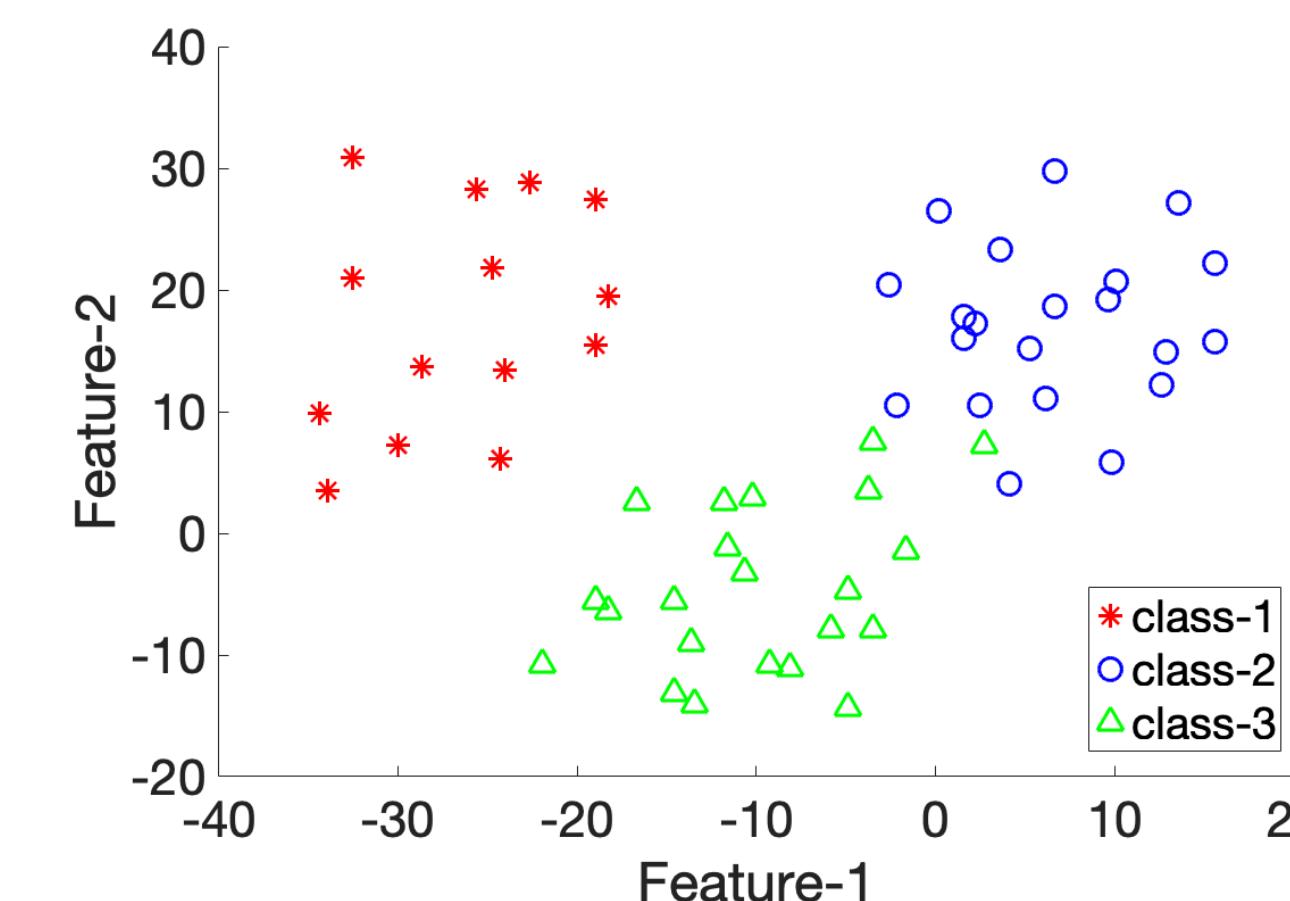


Images: Kaggle, Stackexchange
Video: LuCiD

Data (cont.)

- Feature
 - ▶ Any distinct aspect, quality or characteristic
Ex. numeric (height)
 - ▶ Combination of d features is d -dims column vector, we will call it as **feature vector**
 - ▶ d -dims space defined by the feature vectors is called **feature space**
 - ▶ Data are represented as a **point** in the feature space (\mathcal{X})

$$X = \begin{bmatrix} x^1 \\ x^2 \\ \vdots \\ x^d \end{bmatrix}$$



Data (cont.)

- Syntax:
 - ▶ Data: $(X_i, Y_i); i = 1, \dots, n$
 - Where X_i is the i^{th} data representation (**features**) and $X_i \in R^d$
 - Y_i is the i^{th} data label (class) and $Y_i \in R$
 - ▶ Let f be a model/algorithm for a particular task
 - ▶ \bar{Y}_i is the model/algorithm output: $\bar{Y}_i := f(X_i)$
 - ▶ How do you evaluate your model/algorithm?
 - Loss: $\ell(X_i, Y_i, \bar{Y}_i)$

Data similarity

- How do we compare two objects?



-



=

?

1

-

5

=

?

Data similarity (cont.)

- Why we need comparison?

0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
2	2	2	2	2	2	2	2	2	2	2	2	2	2	2	2	2	2	2	2
3	3	3	3	3	3	3	3	3	3	3	3	3	3	3	3	3	3	3	3
4	4	4	4	4	4	4	4	4	4	4	4	4	4	4	4	4	4	4	4
5	5	5	5	5	5	5	5	5	5	5	5	5	5	5	5	5	5	5	5
6	6	6	6	6	6	6	6	6	6	6	6	6	6	6	6	6	6	6	6
7	7	7	7	7	7	7	7	7	7	7	7	7	7	7	7	7	7	7	7
8	8	8	8	8	8	8	8	8	8	8	8	8	8	8	8	8	8	8	8
9	9	9	9	9	9	9	9	9	9	9	9	9	9	9	9	9	9	9	9

1
?
5
0
?

Similarity measure

- Distance between two points $X_1 \in R^d$ and $X_2 \in R^d$: $\mathcal{D}(X_1, X_2)$?
 - ▶ $\mathcal{D}(X_1, X_2)$ considered a metric if it satisfies the following properties:
 - $\mathcal{D}(X_1, X_2) \geq 0$, $\mathcal{D}(X_1, X_2) = 0$ iff $X_1 = X_2$
 - $\mathcal{D}(X_1, X_2) = \mathcal{D}(X_2, X_1)$
 - $\mathcal{D}(X_1, X_3) \leq \mathcal{D}(X_1, X_2) + \mathcal{D}(X_2, X_3)$
- Commonly used metrics are:
 - ▶ **Minkowski** distance/metric of order p :

- $$- \mathcal{D}^p(X_1, X_2) = \left\{ \sum_{i=1}^d |x_1^i - x_2^i|^p \right\}^{\frac{1}{p}}$$

- ▶ **Manhattan** or city-block distance

- $$- \mathcal{D}^1(X_1, X_2) = \sum_{i=1}^d |x_1^i - x_2^i|$$

Similarity measure (cont.)

- Commonly used metrics are:
 - Minkowski distance/metric of order p :
 - $\mathcal{D}^p(X_1, X_2) = \left\{ \sum_{i=1}^d |x_1^i - x_2^i|^p \right\}^{\frac{1}{p}}$
 - Manhattan or city-block distance:
 - $\mathcal{D}^1(X_1, X_2) = \sum_{i=1}^d |x_1^i - x_2^i|$
 - Euclidean distance:
 - $\mathcal{D}^2(X_1, X_2) = \left\{ \sum_{i=1}^d |x_1^i - x_2^i|^2 \right\}^{\frac{1}{2}}$
 - Chebyshev distance:
 - $\mathcal{D}^\infty(X_1, X_2) = \max_i \{ |x_1^i - x_2^i| \}$

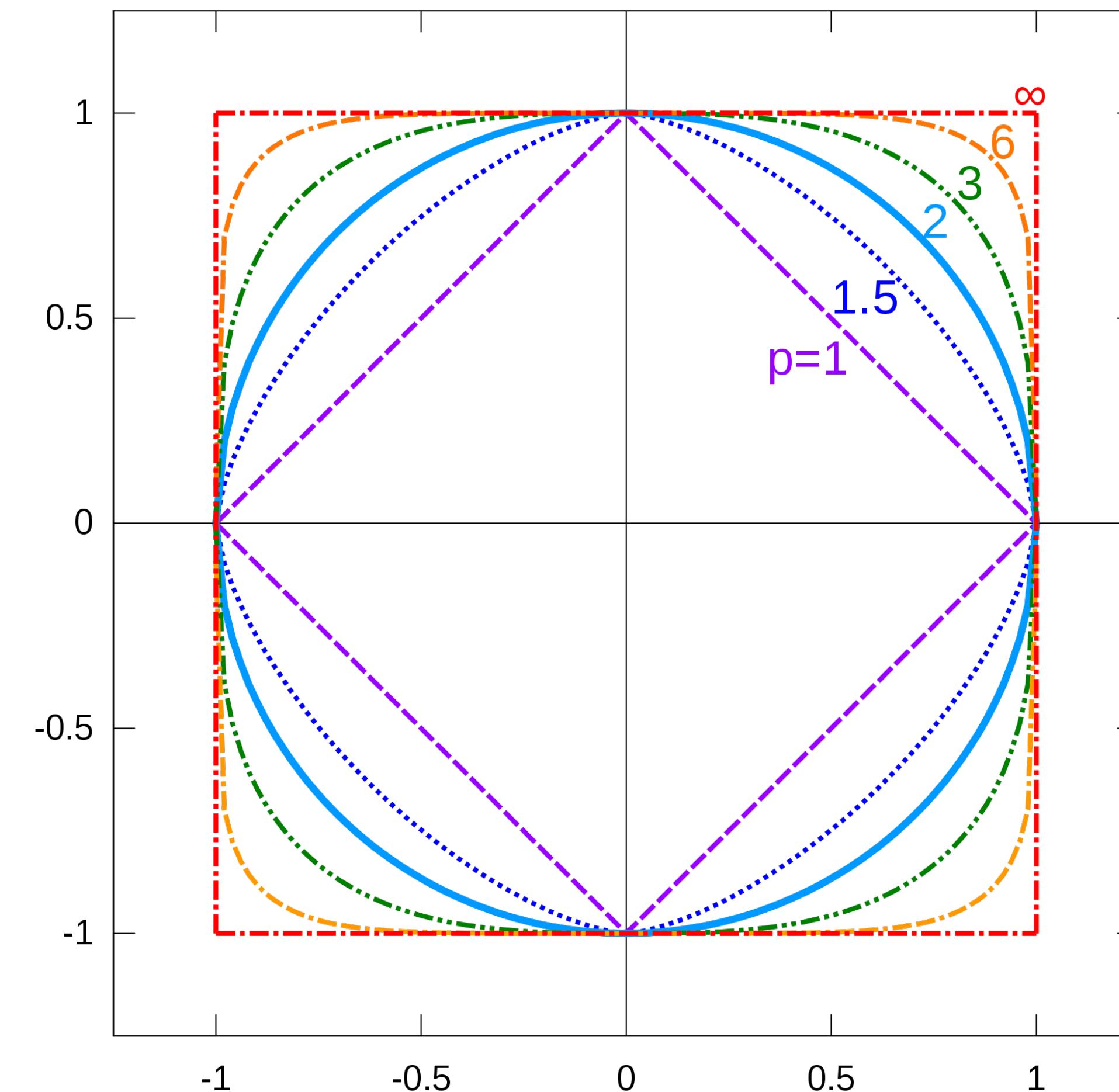
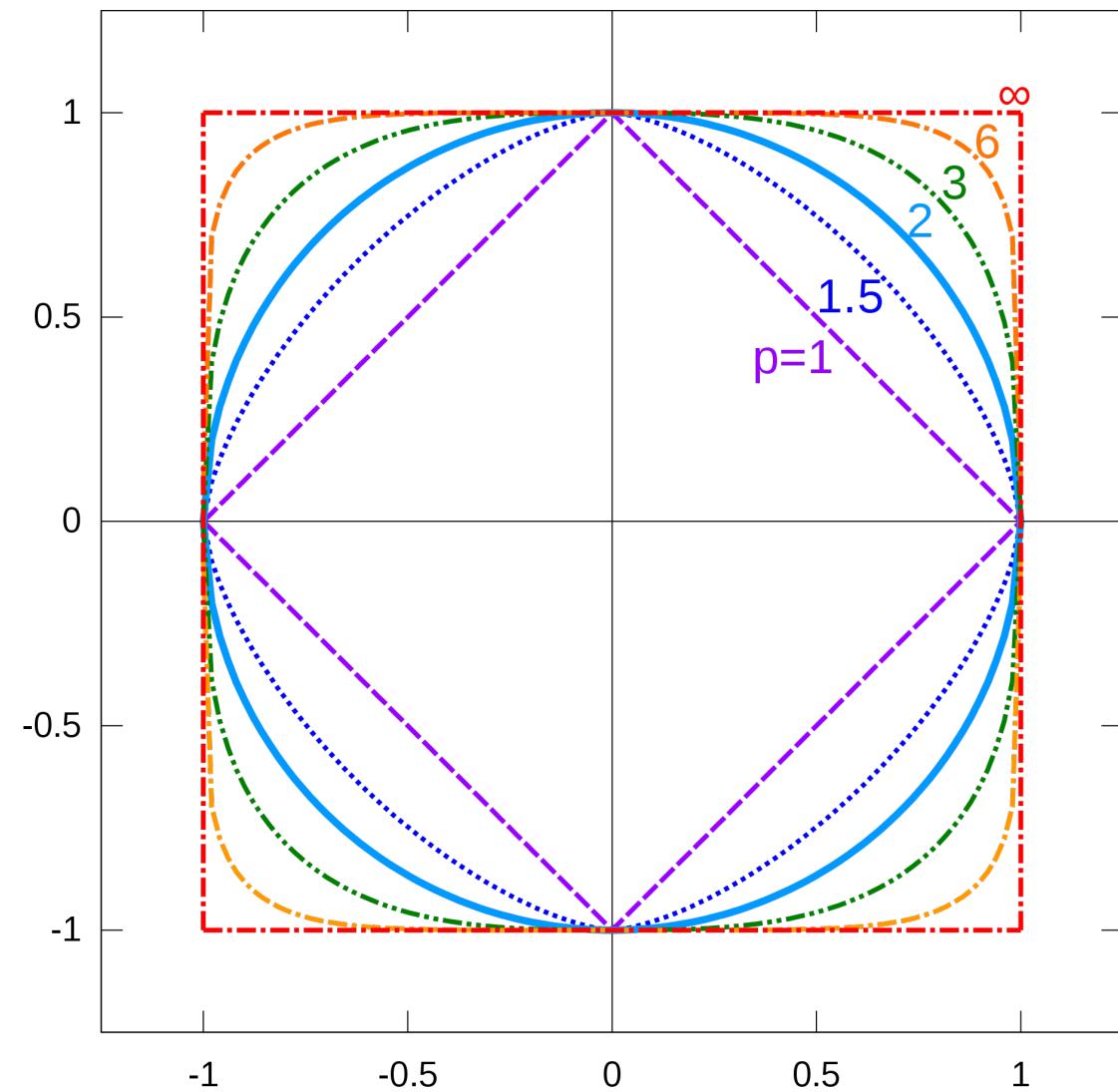


Image: Wikipedia

Homework-1

- Can you draw this type of figure in 2D & 3D?
 - ▶ Use **Minkowski** distance for $p = -2, -1, 0, 1, 1.5, 2, 3, 6, \infty$
 - ▶ Draw all the points within the the intervals has unit distance from the origin for all p 's:
 - 2D: $x \in [-1, 1]$ and $y \in [-1, 1]$
 - 3D: $x \in [-1, 1]; y \in [-1, 1]$ and $z \in [-1, 1]$



Which data to evaluate the model?

- Data partition/division:
 - ▶ Training, Validation and Testing
 - Training- 50%
 - Validation - 20%
 - Testing - 30%

Which data to evaluate the model?

- Model:
 - ▶ Let f be a model/algorithm for a particular task
 - Example: Regression $f: X \rightarrow \mathbf{R}$
 - $f(X) = w^0 + w^T X$

Model parameters: $w^0, w = \begin{bmatrix} w^1 \\ w^2 \\ \vdots \\ w^d \end{bmatrix}$

Which data to evaluate the model?

- Data partition/division:
 - ▶ Training, Validation and Testing
 - Training- 50%
 - Validation - 20%
 - Testing - 30%

Homework-2

- Create a random dataset in \mathbf{R}^{100} of size 50000 with random class labels from {1,2,3,4}. Now partition the data into the following subsets:
 - ▶ Training- 50%
 - ▶ Validation - 20%
 - ▶ Testing - 30%
 - ▶ Plot (bar) the frequency of each class label for each subset.

How do you evaluate your model/algorithm?

- How do you evaluate your model/algorithm?
 - ▶ Accuracy (%) - classification
 - ▶ Mean square error (MSE) - regression
 - ▶ ...

Classification

Recap

- What is ML?
- Some applications of ML
- Deductive vs Inductive inference
- Different types of data and their representation
- Data similarity: different types of distance metrics

Example

- Given digit images:

1
?

K-nearest neighbour classifier

- Let X_1, X_2, \dots, X_n be given feature (observation) vectors, $X_i \in R^d$
- Let Y_i denote the class Lebel of X_i
- Let the number of classes be C
 - ▶ $Y_i \in \{1,2,\dots,C\}$
- Let Y_i 's are known $\forall i = 1,2,\dots,n$
- Let X be a vector for which we don't know its class Lebel

k-nearest neighbour classifier (cont.)

- Let k be a (+)ve integer
- Find k - nearest neighbour of X among X_1, X_2, \dots, X_n
- Let k_i of these nearest neighbours belong to i^{th} class for each $i = 1, 2, \dots, C$
 - ▶
$$\sum_{i=1}^C k_i = k$$
- Put X in the i^{th} class if $k_i > k_j, \forall i \neq j$

K-nearest neighbour classifier (cont.)

- Remark
 - ▶ When $k = 1$, the rule is known as nearest neighbour classifies
 - ▶ There is no universally acceptable way of choosing the value of k
 - ▶ The value of k depends on data point dispersion not only depend on the number of data points
 - ▶ For two different values of k , we may get different results

Classifier (model) evaluation

- Data partition
 - ▶ Training
 - ▶ Validation
 - ▶ Testing
- Model error/loss [$\bar{Y}_i := f(X_i)$]:

- ▶ $\ell(X_i, Y_i, \bar{Y}_i) := \begin{cases} 0 & \text{if } \bar{Y}_i = Y_i \\ 1 & \text{otherwise} \end{cases}$

- ▶ $E[f(X, Y)] = \frac{1}{n} \sum_{i=1}^n \ell(X_i, Y_i, \bar{Y}_i)$

- Pointwise **0 – 1**-loss

Recap

- What is ML?
- Some applications of ML
- Deductive vs Inductive inference
- Different types of data and their representation
- Data similarity: different types of distance metrics
- kNN rule classifier
- Classifier evaluation
 - ▶ data partition

Assignment-1

- Implement kNN classifier and test on MNIST digit data
 - ▶ Download the dataset from here: <http://yann.lecun.com/exdb/mnist/>
 - Strictly follow their data partition
 - There is no validation set!
 - Make your own validation set from the training set (20%)
 - ▶ Use different similarity metrics ($p = 1, 2$ and ∞) and ($k = 1, 3, \dots, 25$) calculate the classifier errors
 - ▶ Plot (3-D) the classification errors/accuracy for different p 's and k 's
- Submission deadline: 19-02-2024
- knn_mnist_data_your_name_version_no.ipynb

kNN algorithm details

- For algorithm:
 - ▶ Duda, Hart, Stork; Section 4.5
- Theory:
 - ▶ Devroye, Gyorfi and Lugosi; A probabilistic Theory of Pattern Recognition;
Chapter-5