

End-to-end machine learning



Saptarshi Purkayastha, Ph.D.

<https://plhi.lab.Indianapolis.iu.edu>

- Associate Professor, Health Informatics, Data Science at Indiana University Indianapolis
 - Director, Health Informatics, School of informatics & Computing
 - Director of Undergraduate Education & Research, BHI
- Affiliated Scientist, CBMI, Regenstrief Institute
- Faculty Associate, STEM Education Innovation and Research Institute (SEIRI)
- Research funded by:



Atika Paddo
PhD Student



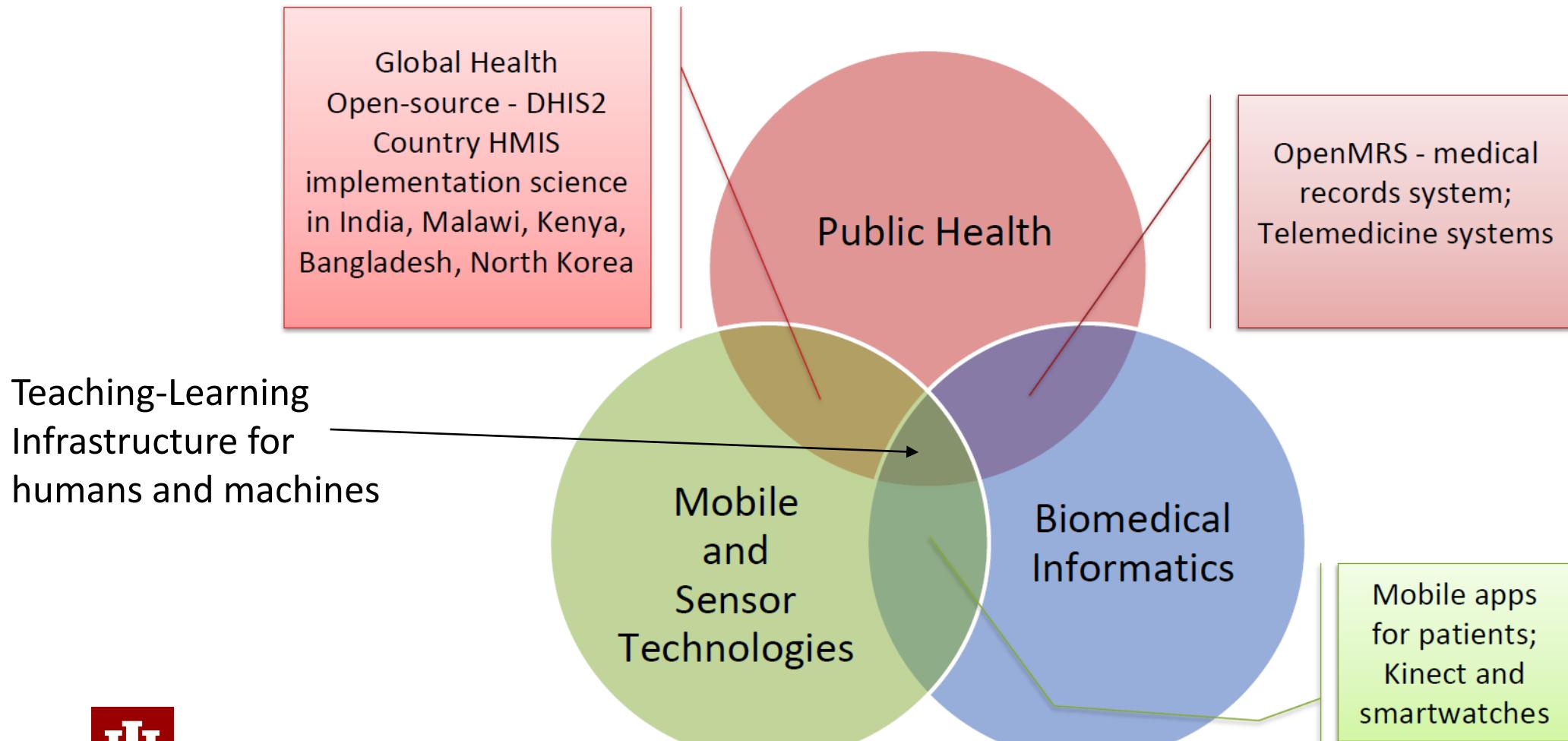
Google
Cloud Platform



Work centered around building teaching-learning infrastructure

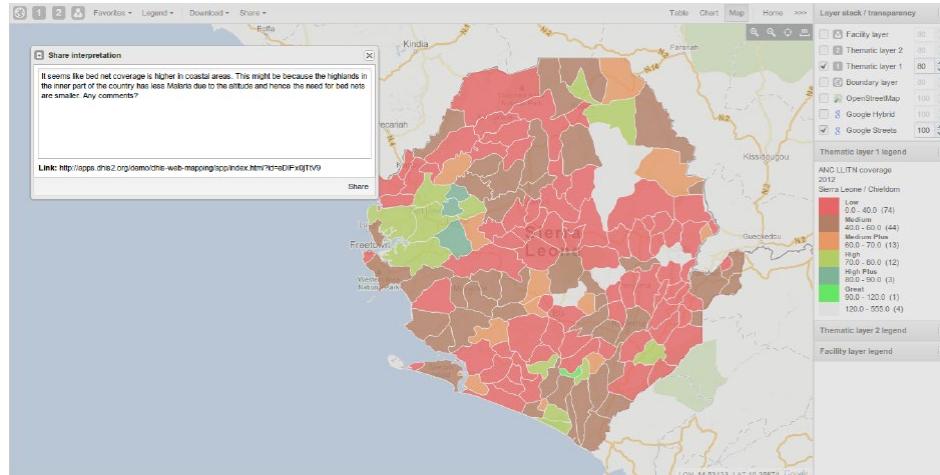
"Information is Care"

A world-wide recognized fact that clinicians provide best treatment to patients, when medical history of the patient is available

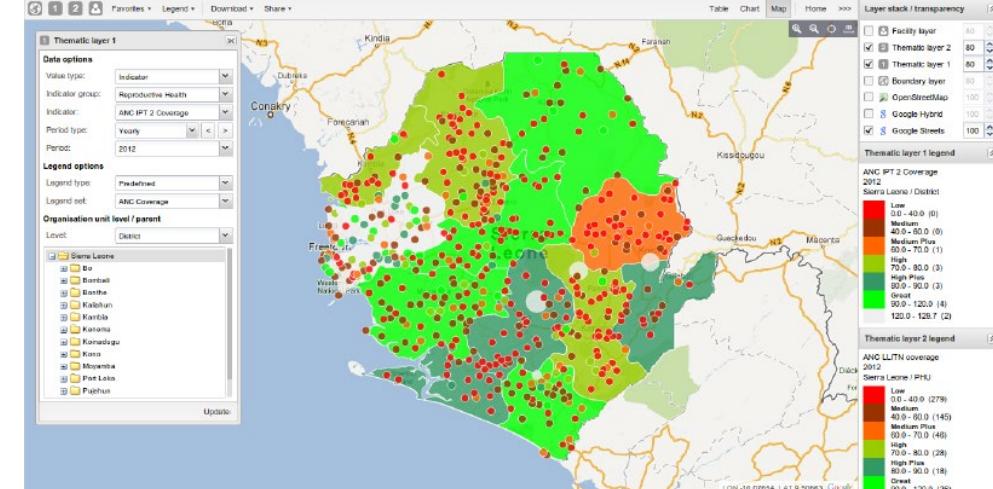


Health Information Infrastructure

- National eHealth Architectures: With WHO SEARO, I developed the first country-wide telemedicine system in Bhutan (2015) based on web-telephony and OpenMRS EHR. This was followed by DHIS2 optimization for patient records in Bangladesh (2015), training and implementation of national HMIS in North Korea (2015). The most recent is the development of the eHealth Architecture for Myanmar (2017).
- Cybersecurity in the cloud for privacy preservation: Development and integration of novel cybersecurity technologies in the NSF first public cloud called JetStream. \$600k NSF-funded project (co-PI)
- Wound care platform: Using DHIS2, we've built a platform for Indiana Center for Regenerative Medicine and Engineering for program monitoring and evaluation about amputations after non-recoverable wounds.

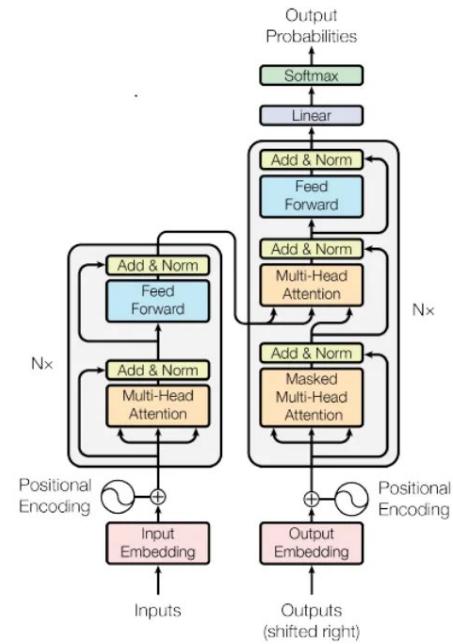


And building a social network to distribute the analysis of the data

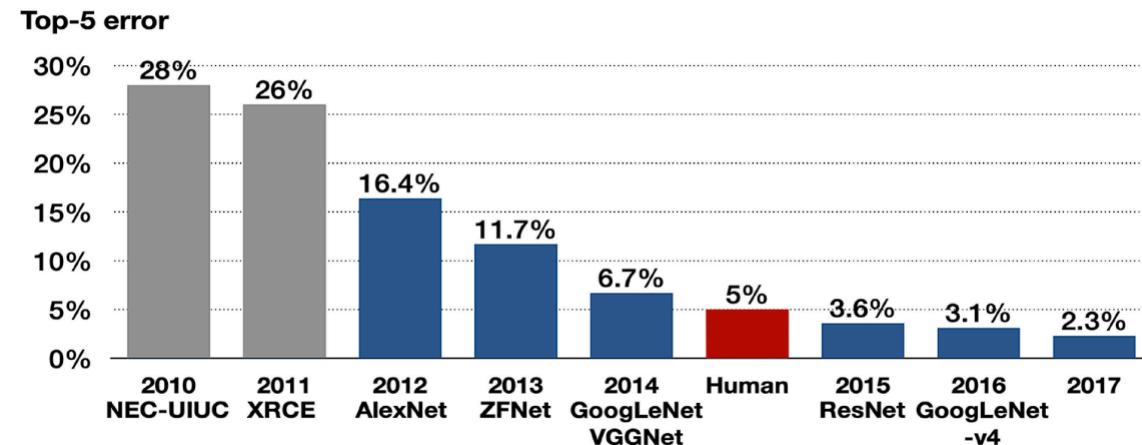


Displaying disease outbreak in Sierra Leone



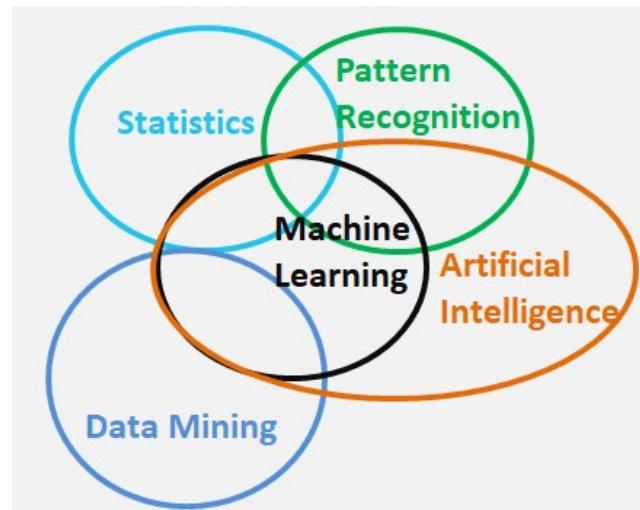


Why need Artificial Intelligence?



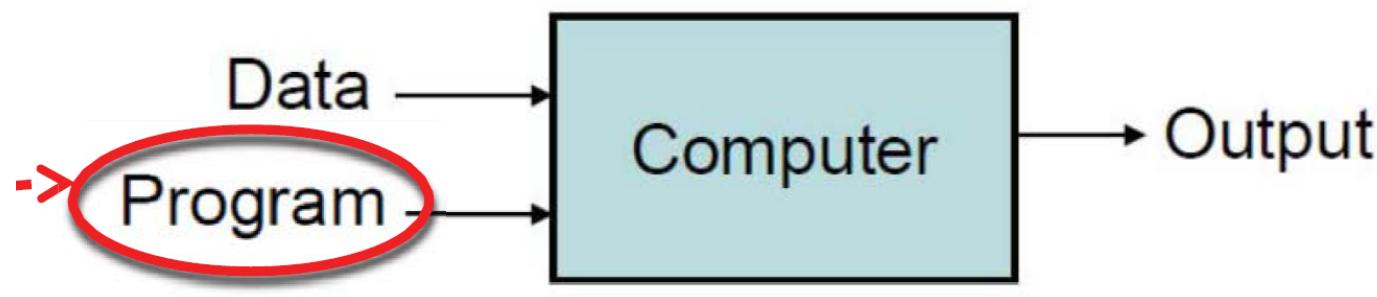
What is machine learning?

- All useful programs "learn" something.
- Suppose you wrote a program to calculate the square root of any number with some acceptable error, say 0.001. Is the program "learning" something?
- Machine learning as a field started with Arthur Samuel's 1959 paper which said, "Field of study that gives computers the ability to learn without being explicitly programmed."
 - Computer pioneer who wrote first self-learning program, which played checkers - learned from "experience".
 - Invented alpha-beta pruning - widely used in decision tree searching

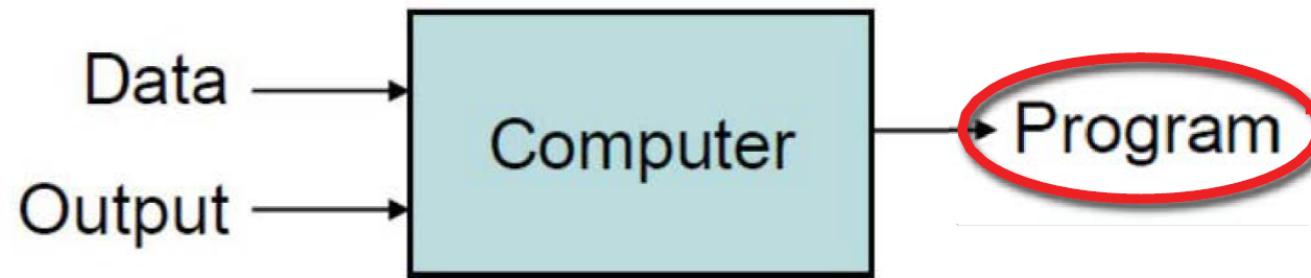


Traditional vs machine learning programs

Traditional Programming



Machine Learning



Square root
finder

Curve fitting by
linear regression

How are things learned?

- Memorization
 - Accumulation of individual facts
- Limited by
 - Time to observe facts
 - Memory to store facts
- Generalization
 - Deduce new facts from old facts
 - Limited by accuracy of deduction process
 - Essentially a predictive activity
 - Assumes that the past predicts the future
- Interested in extending to programs that can infer useful information from **implicit** patterns in data

Declarative knowledge

Imperative knowledge

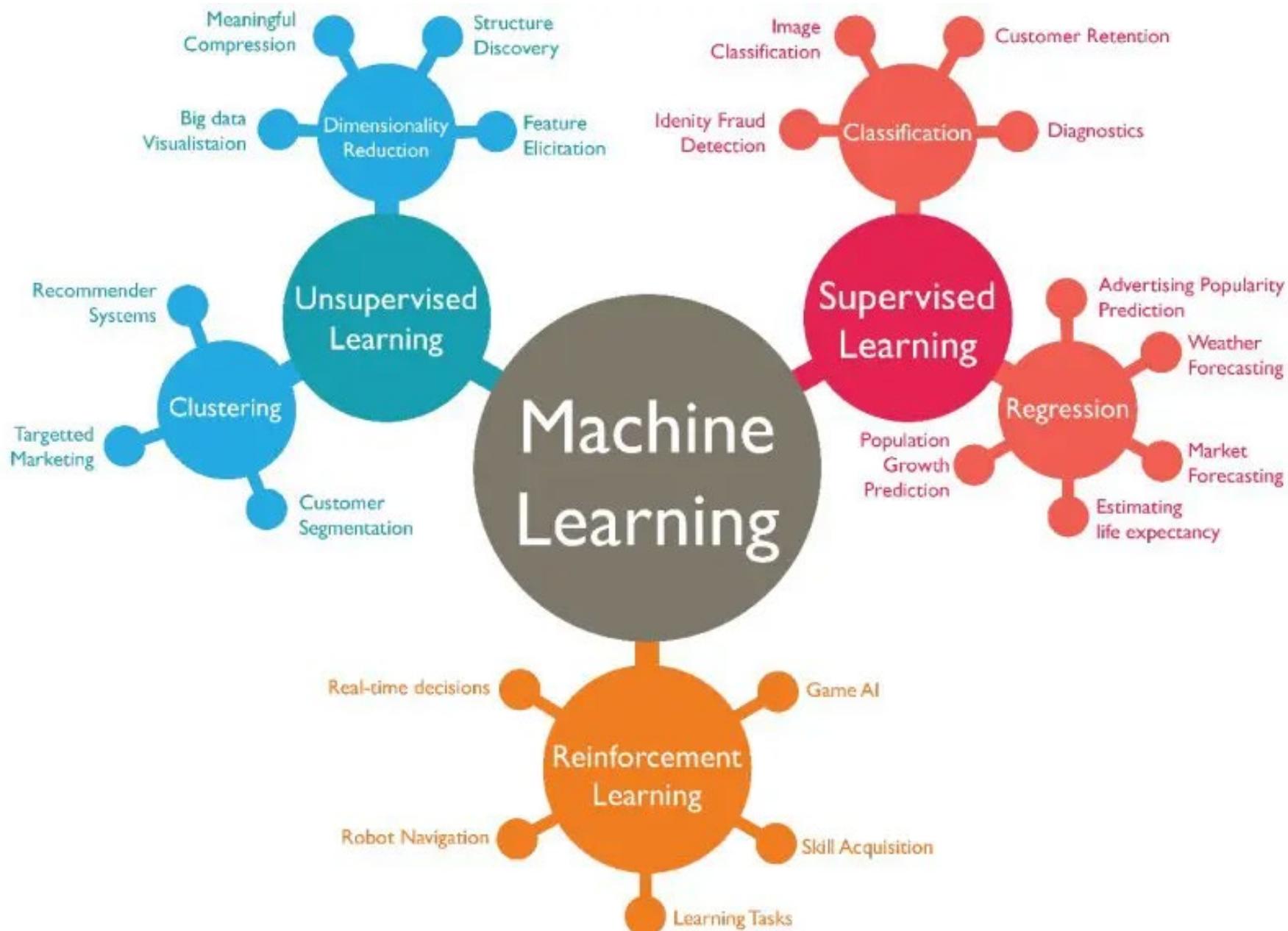
Basic paradigm of machine learning

- Observe set of examples
 - *Training data*
- Infer something about process that generated the data
 - Types of model
 - Create a model
 - *Parameters* to a model
- Use inference to make predictions about previously unseen data
 - *Test data*
- Two paradigms on how to learn:
 - **Supervised**: given a set of feature/label pairs, find a rule that predicts the label associated with a previously unseen input.
 - **Unsupervised**: given a set of feature vectors {without labels}, group them into "natural clusters" {or create labels for the groups}

Spatial deviations relative to mass displacement of spring

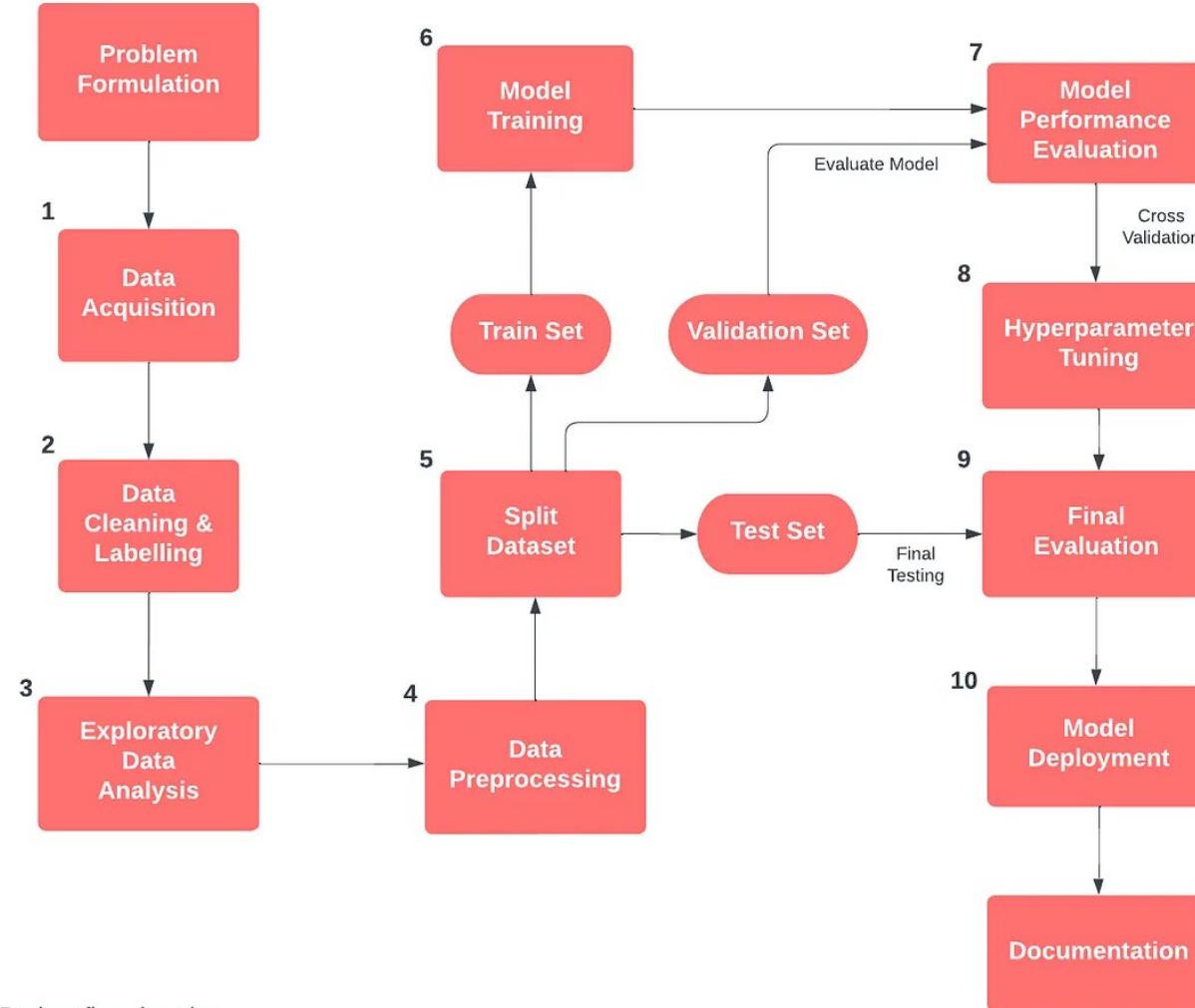
Fit polynomial curve using linear regression

Predict displacements for other weights



What is the end-to-end process for machine learning?

Machine Learning Project Life Cycle



Emory ICU dataset

Problem formulation

- The length of stay (LOS) in MIMIC-IV is typically measured in days, not hours.
- Some important features often used for LOS prediction in MIMIC-IV include:
 1. Patient demographics: Age, Gender, Ethnicity
 2. Admission details: Admission type (emergency, elective, etc.), Admission location, Insurance type
 3. Clinical data: Primary diagnosis, Comorbidities, Severity scores (e.g. SOFA, SAPS II)
 4. Vital signs: Heart rate, Blood pressure, Respiratory rate, Temperature, Oxygen saturation
 5. Laboratory values: Complete blood count, Basic metabolic panel, Liver function tests
 6. Interventions: Mechanical ventilation, Vasopressor use, Dialysis
 7. Medication data
 8. Previous hospital visits/admissions
- The specific set of most important features can vary depending on the particular prediction model and patient population. Some studies have found factors like age, admission type, primary diagnosis, and certain lab values to be particularly predictive.



Data acquisition - Where is my data?

What is parquet?

Tables in parquet format

- Apache Parquet is a free and open-source column-oriented data storage format in the Apache Hadoop ecosystem.
- Pandas library native

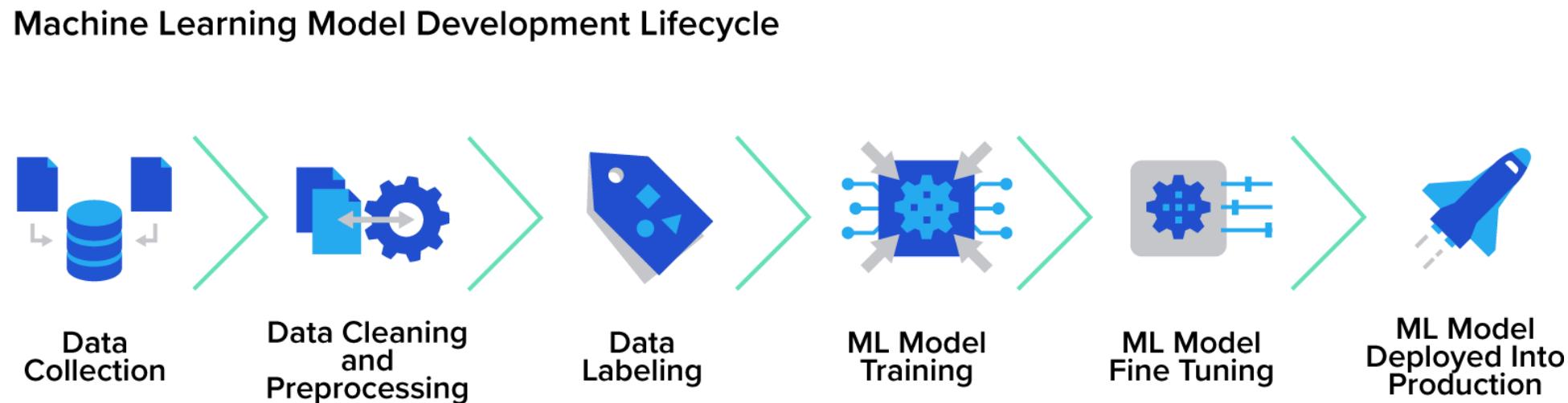
	Column 1	Column 2	Column 3	Column 4	Column 5
	Product	Customer	Country	Date	Sales Amount
Row Group 1	Ball	John Doe	USA	2023-01-01	100
	T-Shirt	John Doe	USA	2023-01-02	200
Row Group 2	Socks	Maria Adams	UK	2023-01-01	300
	Socks	John Doe	Australia	2023-01-02	100
Row Group 3	T-Shirt	Maria Adams	UK	2023-01-02	500
	Socks	John Doe	USA	2023-01-05	200

Spark Format Showdown		File Format		
		CSV	JSON	Parquet
A t t r i b u t e	Columnar	No	No	Yes
	Compressable	Yes	Yes	Yes
	Splittable	Yes*	Yes**	Yes
	Human Readable	Yes	Yes	No
	Nestable	No	Yes	Yes
	Complex Data Structures	No	Yes	Yes
	Default Schema: Named columns	Manual	Automatic (full read)	Automatic (instant)
	Default Schema: Data Types	Manual (full read)	Automatic (full read)	Automatic (instant)



Data cleaning and labeling

- Remove duplicates.
- Data imputation – remove where an estimate is unrealistic, fill where distribution remains the same.
- Merging data

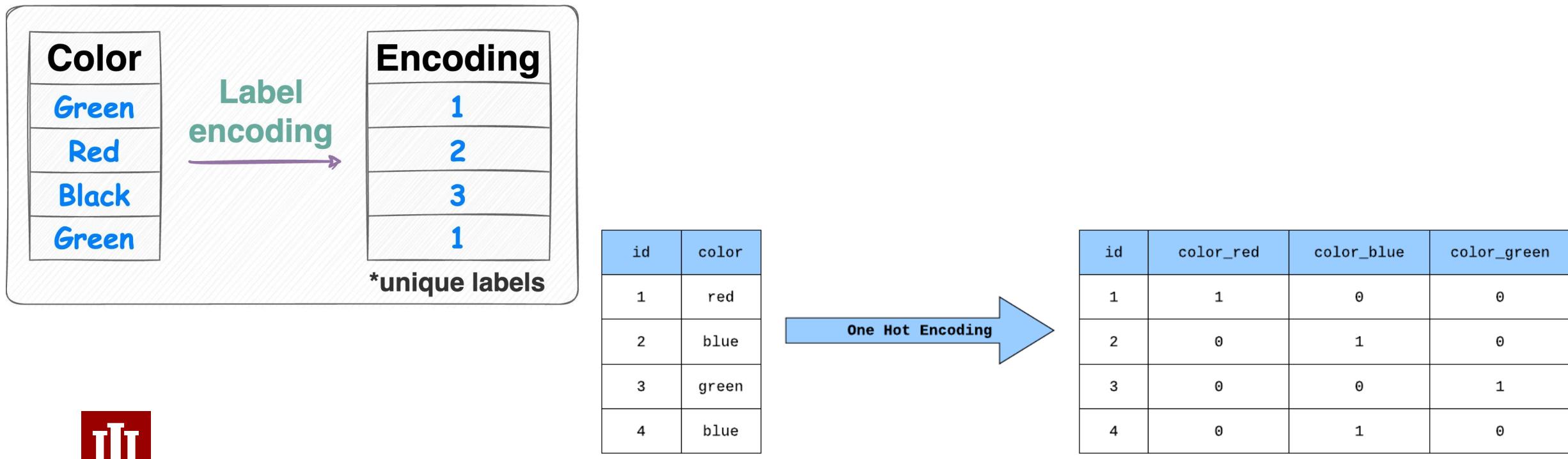


Data from the production system is
collected to further tune the model.



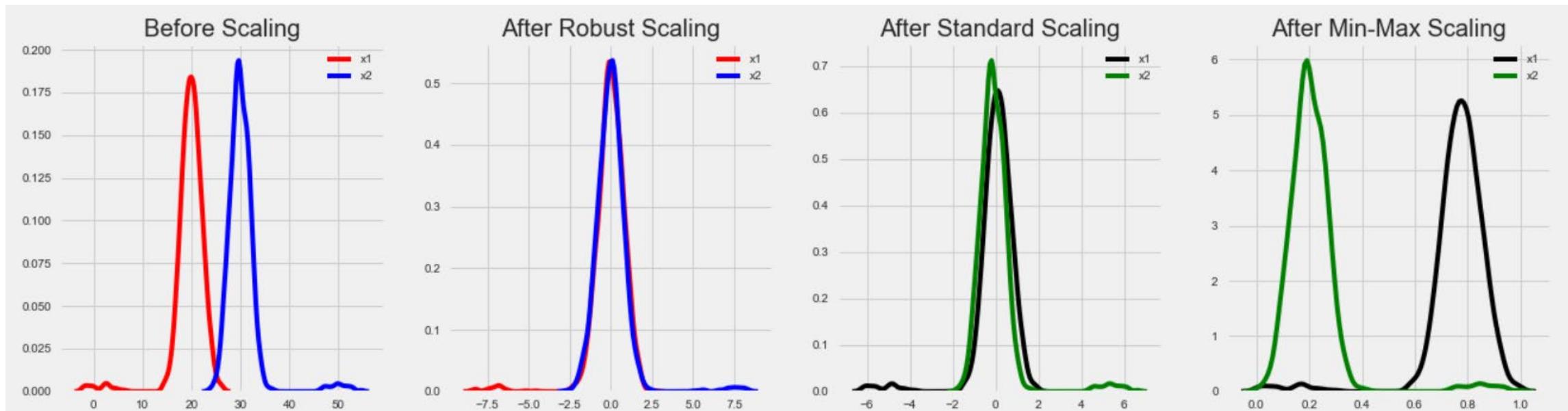
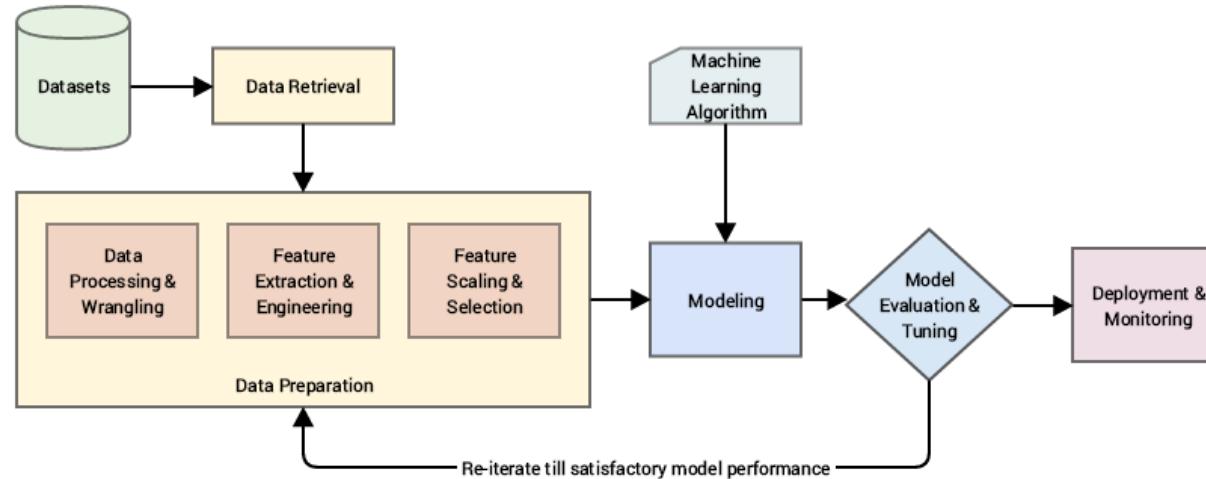
Exploratory data analysis (EDA)

- Distribution of the variables. Different types of visualization.
- Reviewing the Data types.
- Applying specific pre-processing based on continuous or discrete variables.
- Label Encoder vs One Hot Encoder



Feature engineering

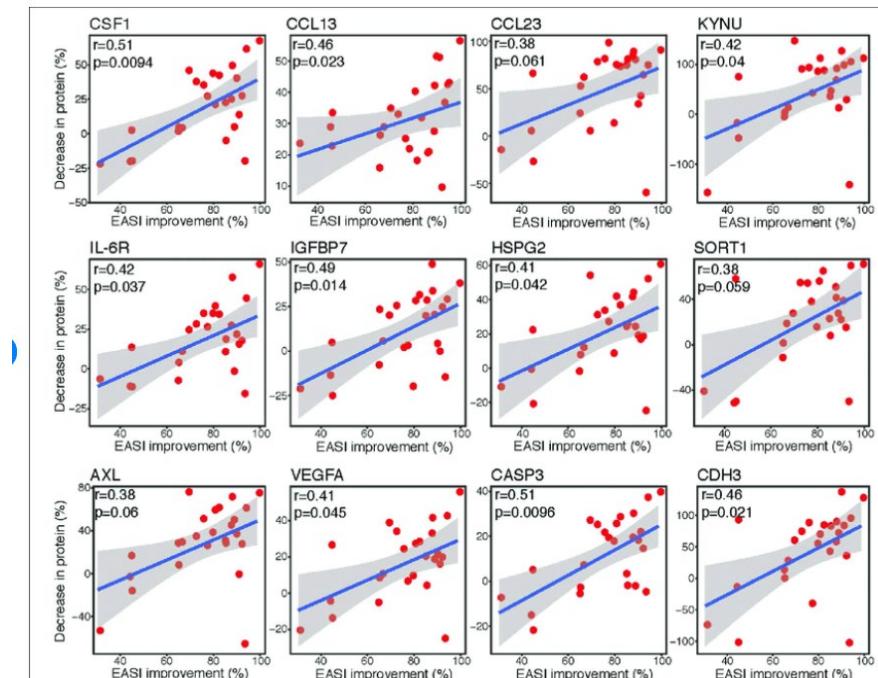
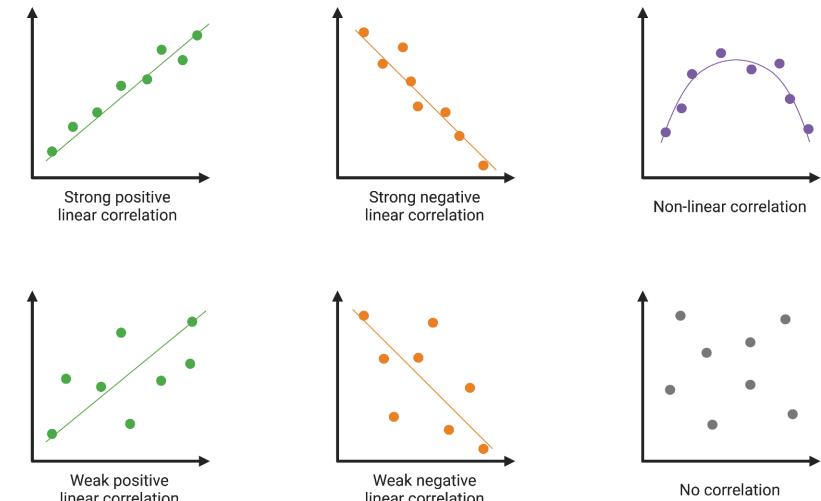
- An iterative process with domain knowledge – knowing what to change and associate.



Correlation

- Keeping one out of strongly correlated features – so that it doesn't affect the model compared to other variables.
- Once you've understood the datatype, choice of correlation method will explain type
- Most people use heatmap, but corr plots are possibility, if you know what it means.

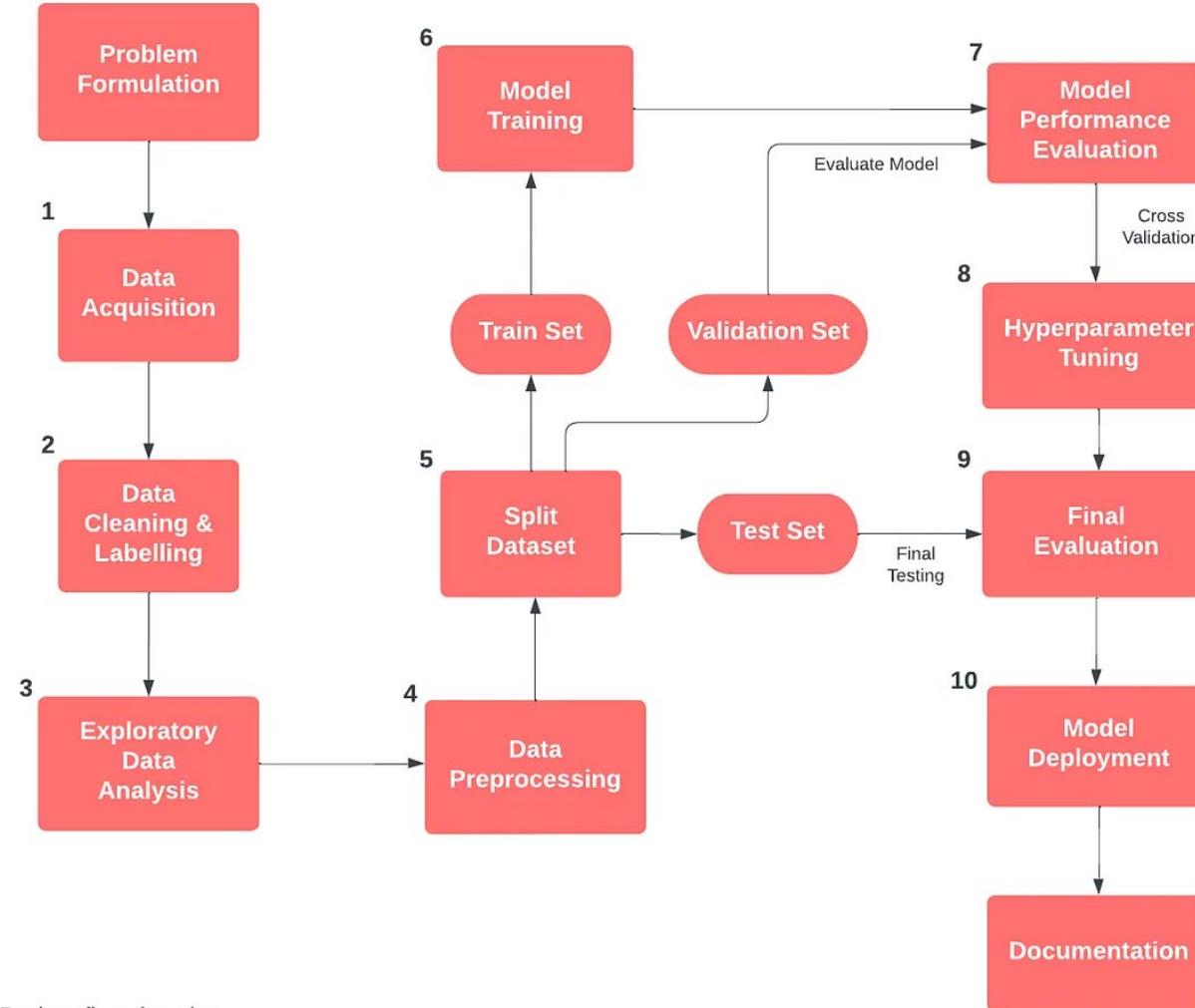
Types of Correlation



Variable Type	Continuous Variable	Ordinal Variable	Categorical Variable
Continuous variable	Pearson coefficient	Spearman's coefficient	Eta coefficient or difference analysis
Ordinal variable	Spearman's coefficient	Spearman's coefficient	Contingency coefficient, Phi (Φ) coefficient, or difference analysis
Categorical variable	Eta coefficient or difference analysis	Contingency coefficient, Phi (Φ) coefficient, or difference analysis	Contingency coefficient or Phi (Φ) coefficient

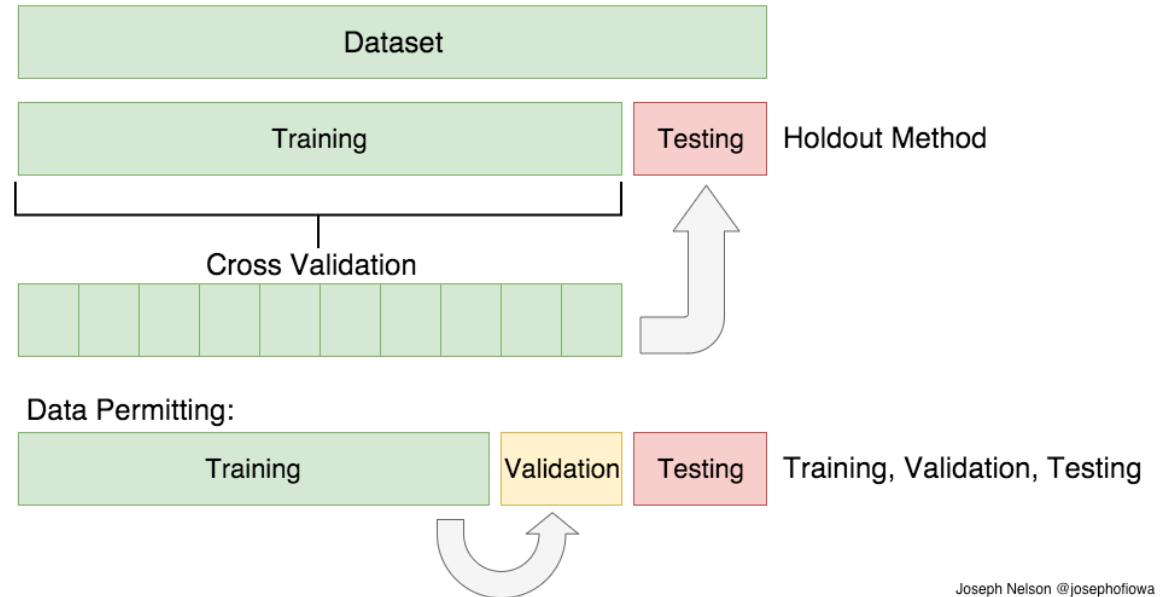
What is the end-to-end process for machine learning?

Machine Learning Project Life Cycle

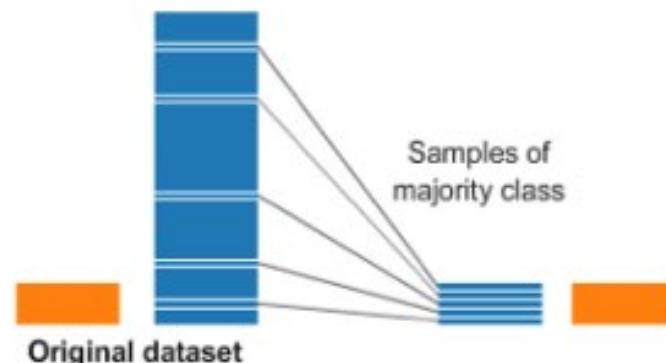


Dataset split

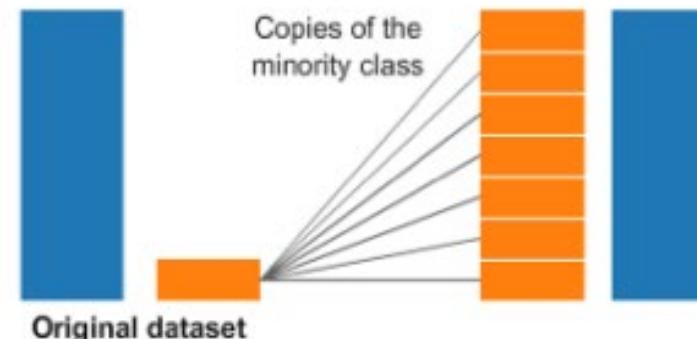
- Hiding data from the model, to validate in the training process and also test once the model is built.



Undersampling



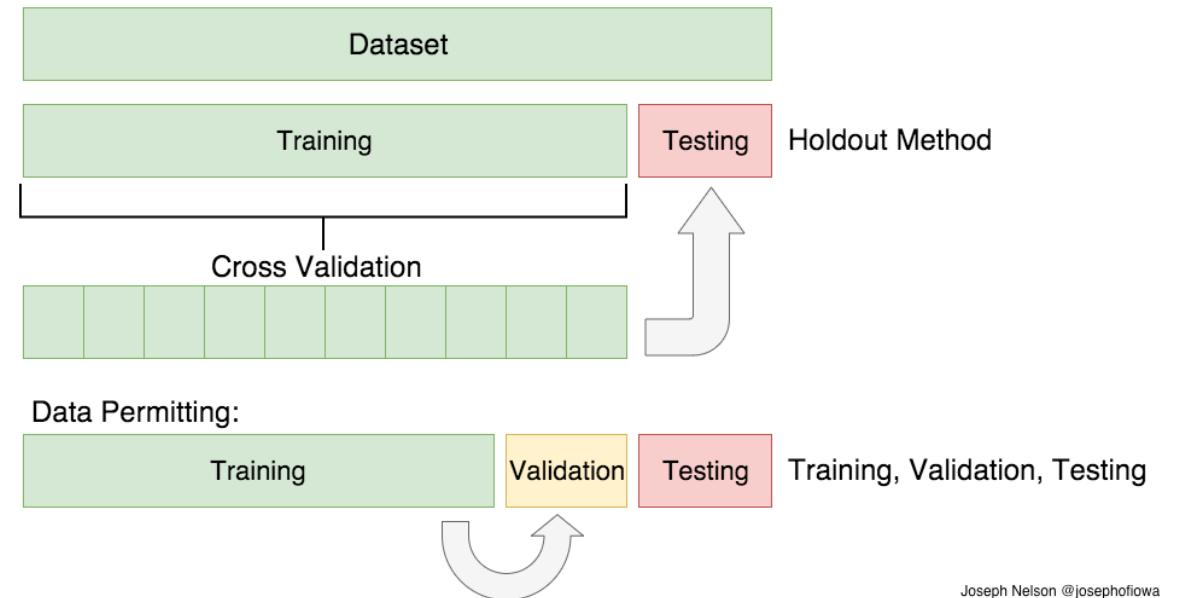
Oversampling



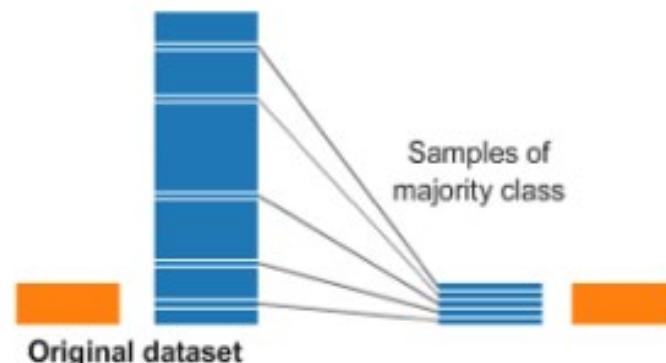
Ψ

Decision Trees

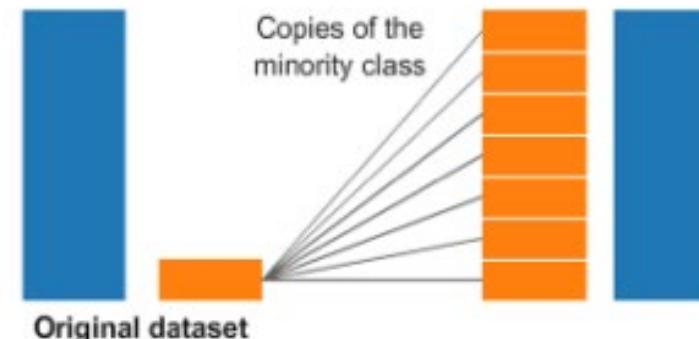
- Hiding data from the model, to validate in the training process and also test once the model is built.



Undersampling

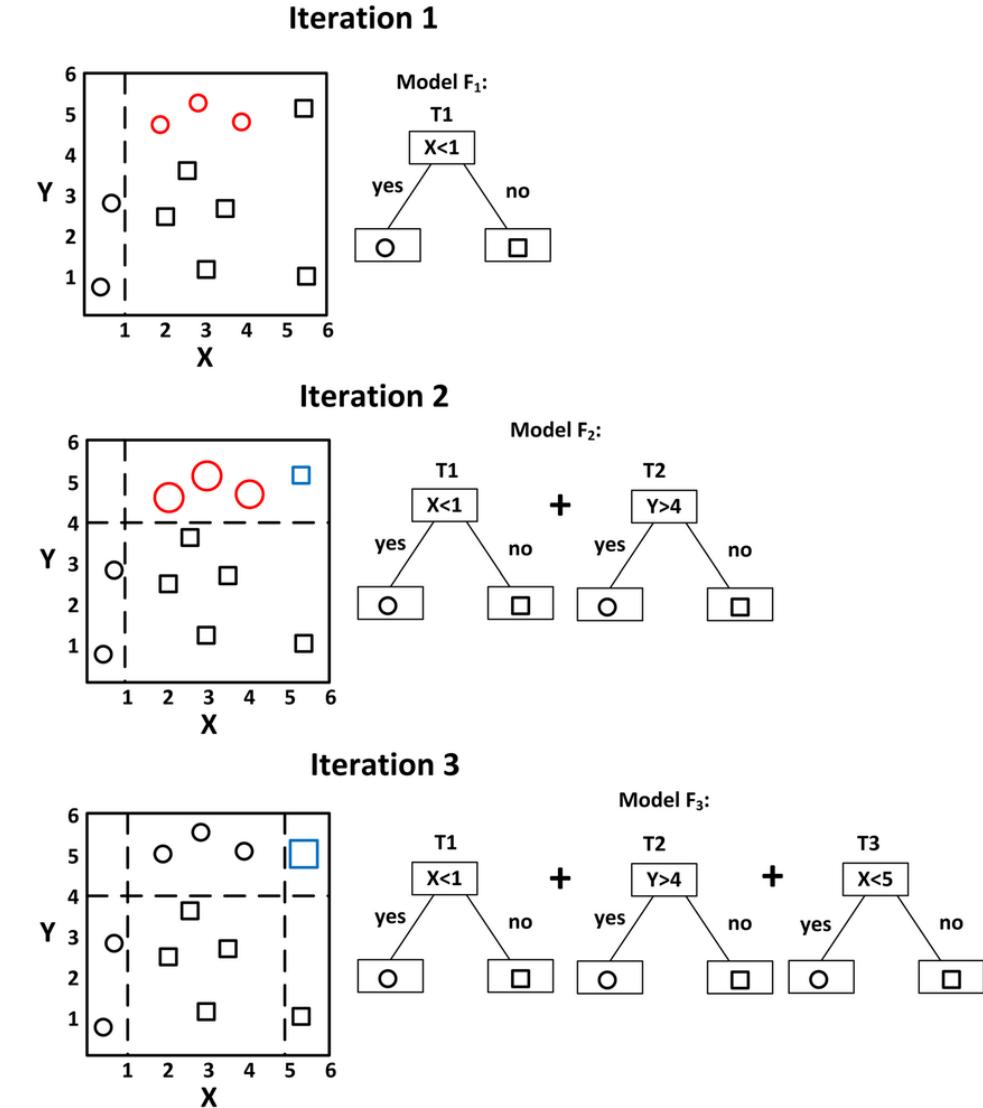
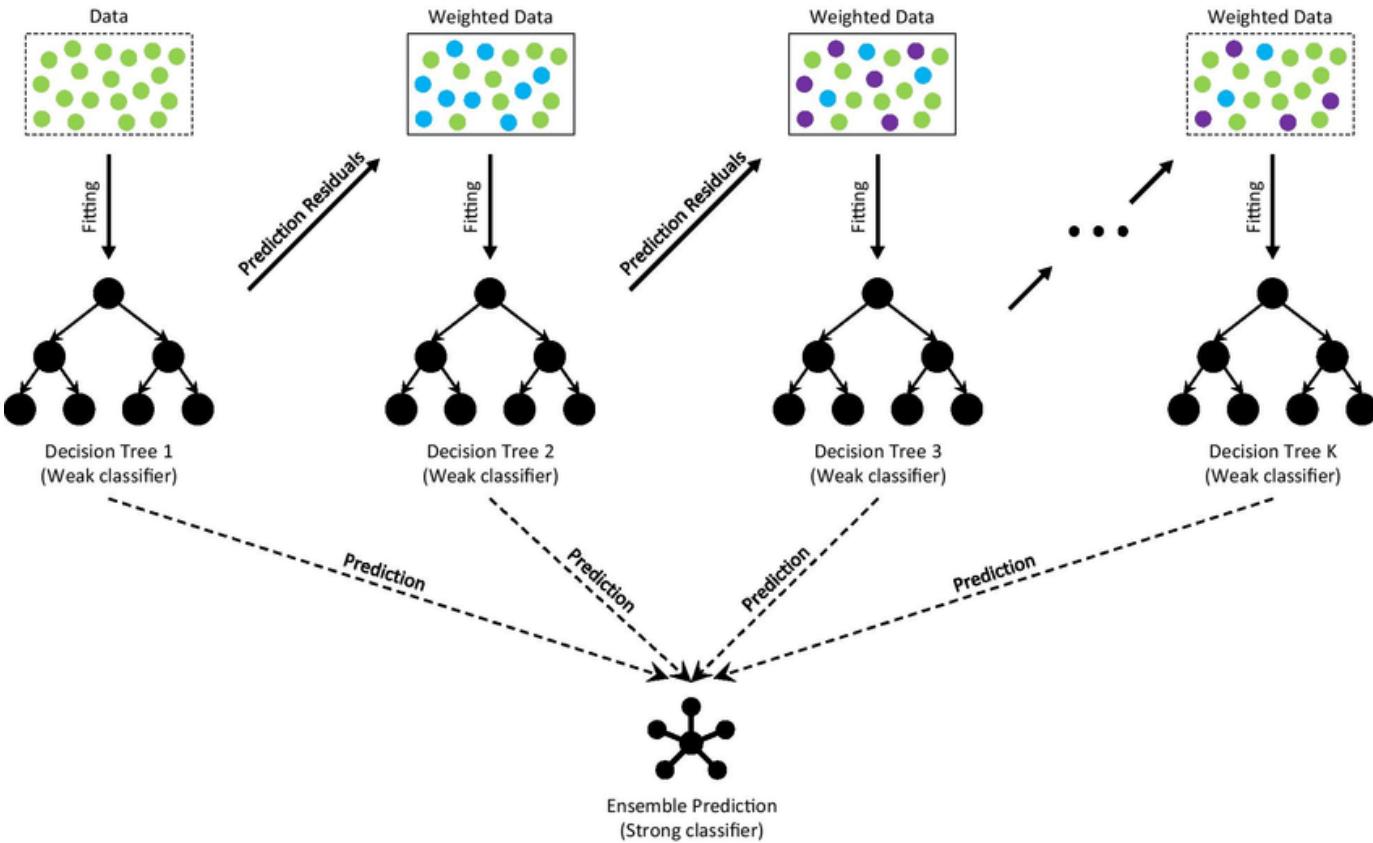


Oversampling

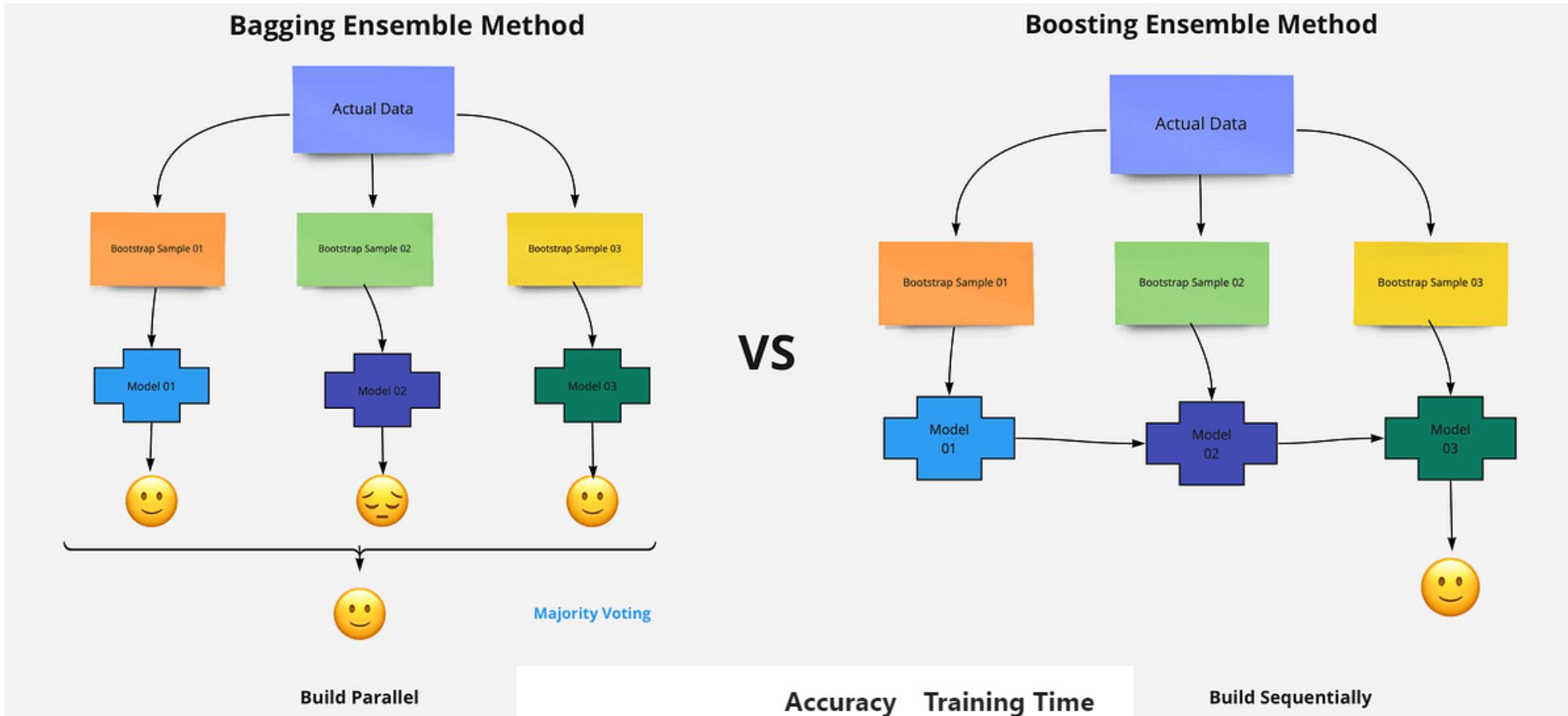


Ψ

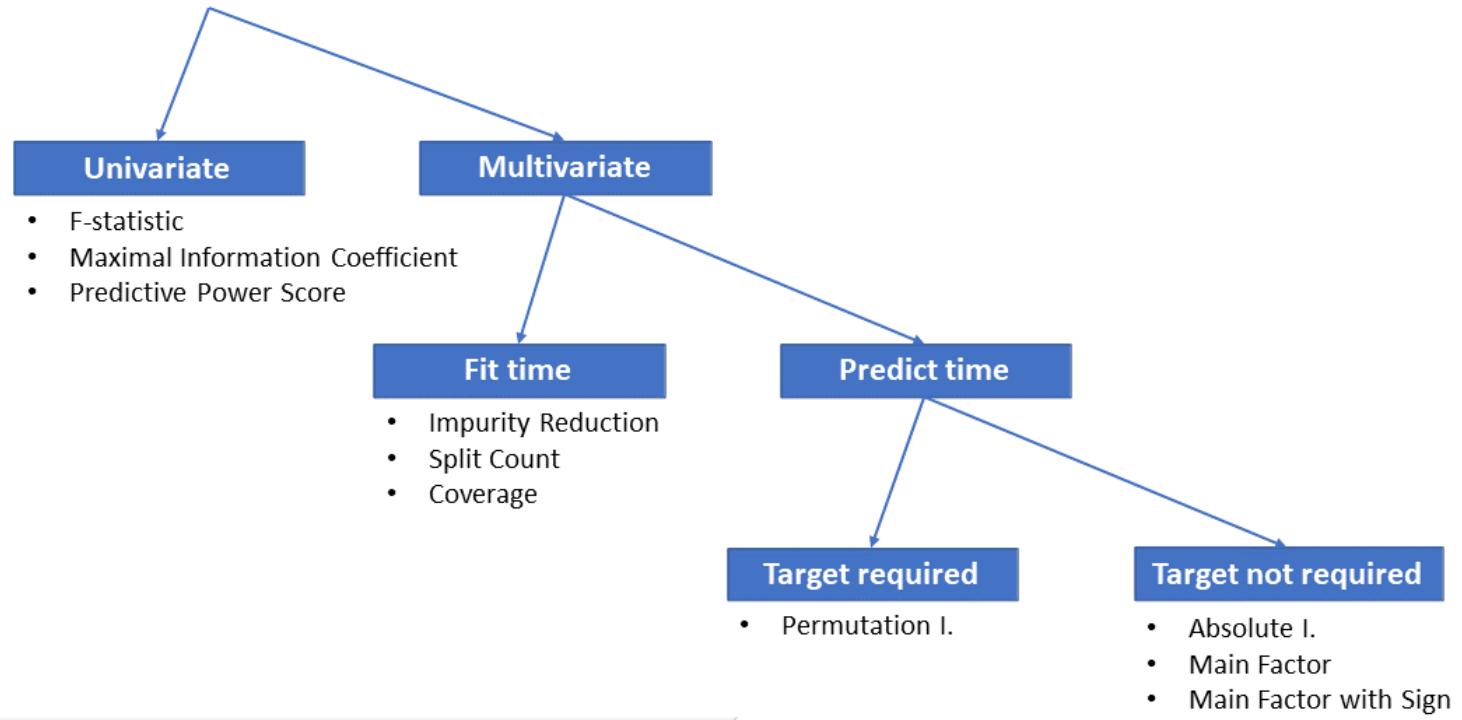
Gradient Boosting methods



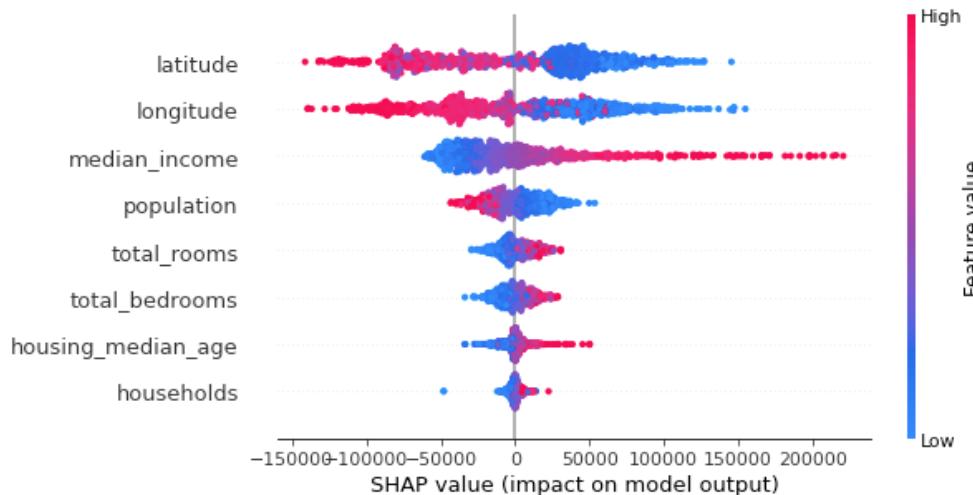
Boosting vs Bagging (Random Forest)



Explaining with feature importance



```
shap.summary_plot(xg_shap_values, X_st)
```



Metrics to evaluate the regression model

$$\begin{aligned} R^2 &= 1 - \frac{\text{Residual variance}}{\text{Total variance}} \\ &= \frac{\text{Total variance} - \text{Residual variance}}{\text{Total variance}} \\ &= \frac{\text{Explained variance}}{\text{Total variance}} \\ &= \text{Fraction of total variance explained} \end{aligned}$$

R	Interpretation
0.00-0.199	Very weak
0.20-0.399	Weak
0.40-0.599	Medium
0.60-0.799	Strong
0.80-1.00	Very strong

number of samples n

$$\sum_{i=1}^n (\text{real value } Y_i - \text{predicted value } \hat{Y}_i)^2$$

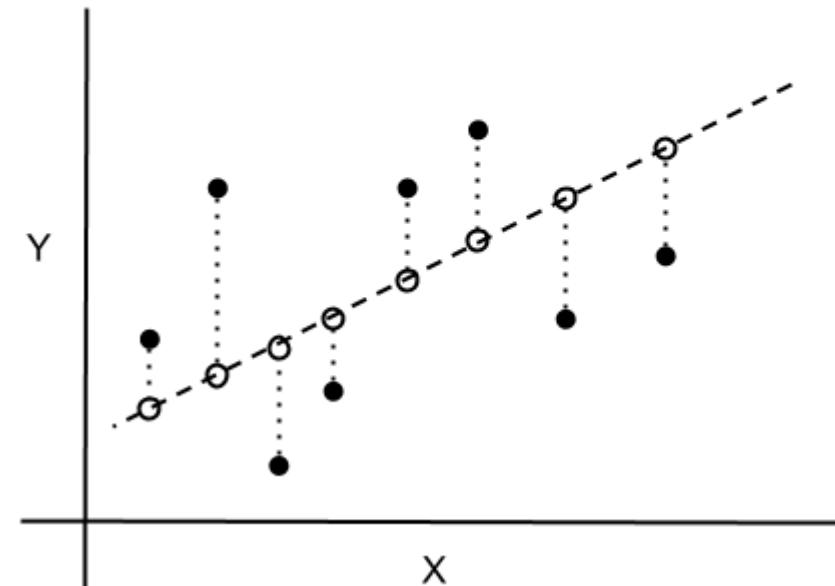
sum of the errors of all samples

● Observed Values

○ Predicted Values

- - - Regression Line

..... Y Scale Difference Between Observed and Predicted Values



Metrics to evaluate the classification model

		True condition	
		Condition positive	Condition negative
Predicted condition	Predicted condition positive	True positive	False positive, Type I error
	Predicted condition negative	False negative, Type II error	True negative
	True positive rate (TPR), Recall, Sensitivity, probability of detection, Power = $\frac{\sum \text{True positive}}{\sum \text{Condition positive}}$	False positive rate (FPR), Fall-out, probability of false alarm = $\frac{\sum \text{False positive}}{\sum \text{Condition negative}}$	
	False negative rate (FNR), Miss rate $= \frac{\sum \text{False negative}}{\sum \text{Condition positive}}$	Specificity (SPC), Selectivity, True negative rate (TNR) = $\frac{\sum \text{True negative}}{\sum \text{Condition negative}}$	

CALCULATE
PRECISION,
RECALL,
F1-SCORE

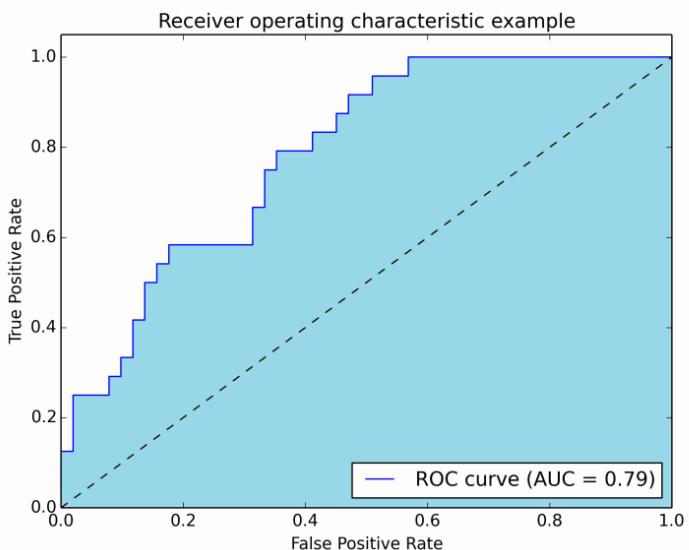
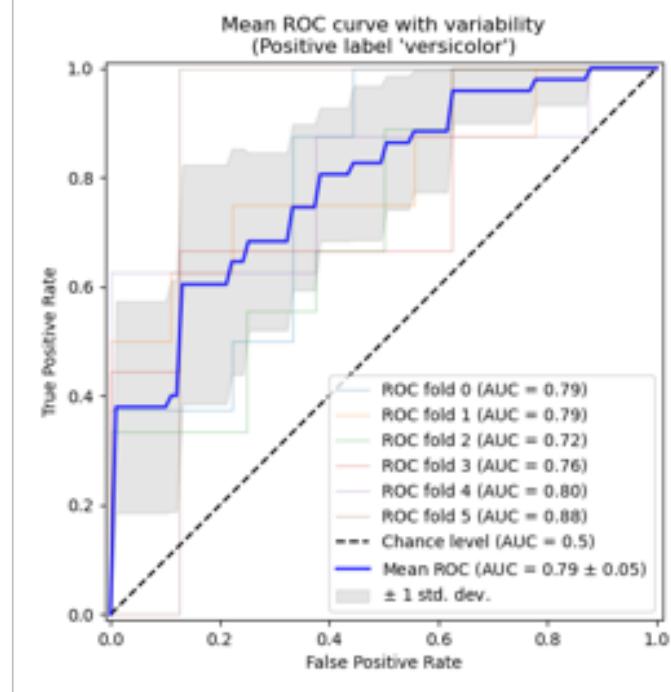
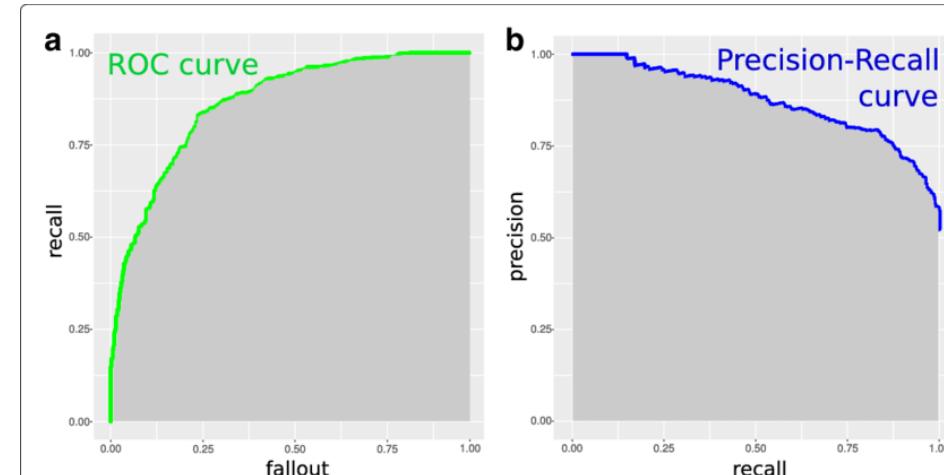
$$\text{precision} = \frac{TP}{TP + FP}$$

$$\text{recall} = \frac{TP}{TP + FN}$$

$$F1 = \frac{2 \times \text{precision} \times \text{recall}}{\text{precision} + \text{recall}}$$

$$\text{accuracy} = \frac{TP + TN}{TP + FN + TN + FP}$$

$$\text{specificity} = \frac{TN}{TN + FP}$$



Thank you for your attention

