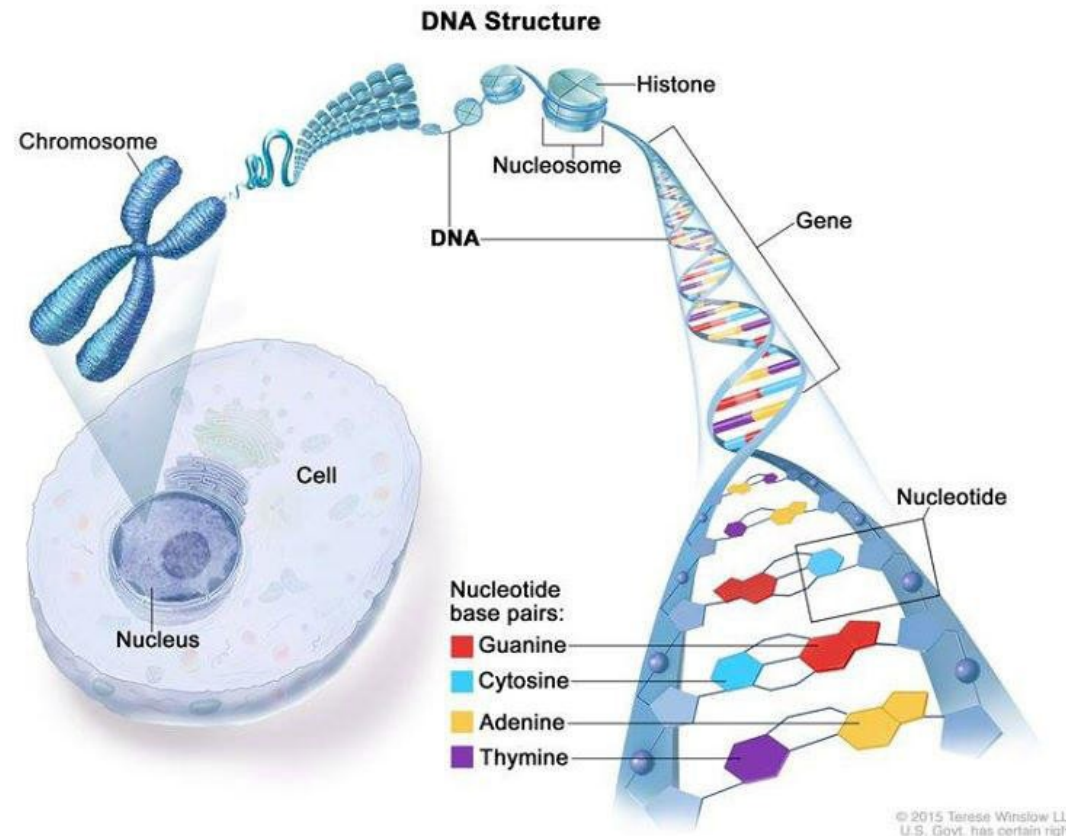# AI for Precision Oncology

# Project Aim

- The aim of my project is to create a **Machine Learning (ML) model** that can **predict** a patient's Immunotherapy **response** using transcriptomic data.
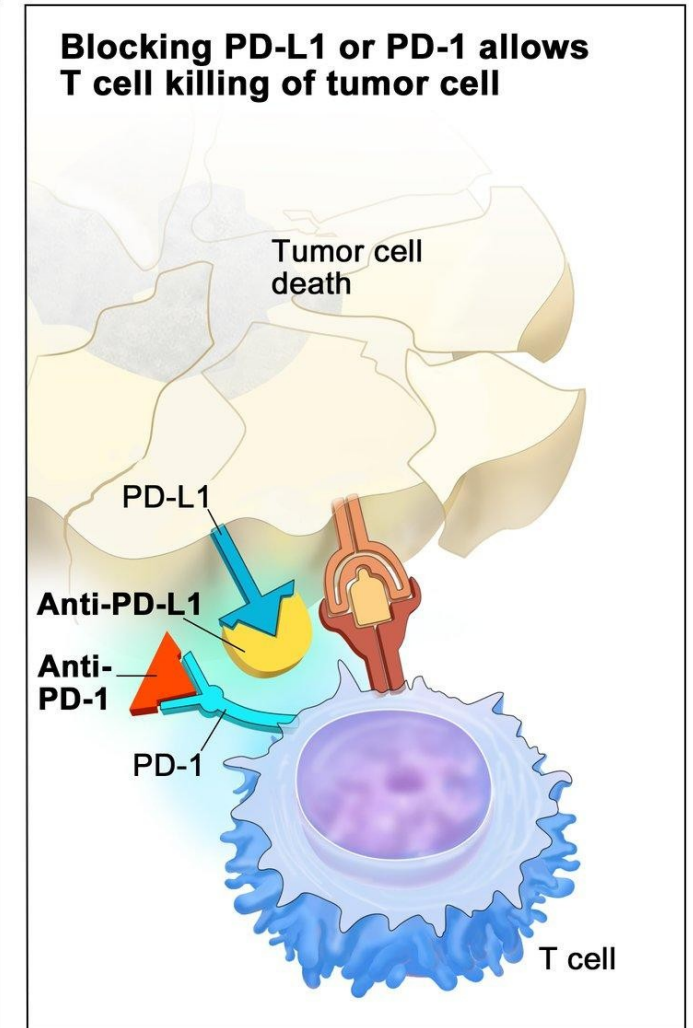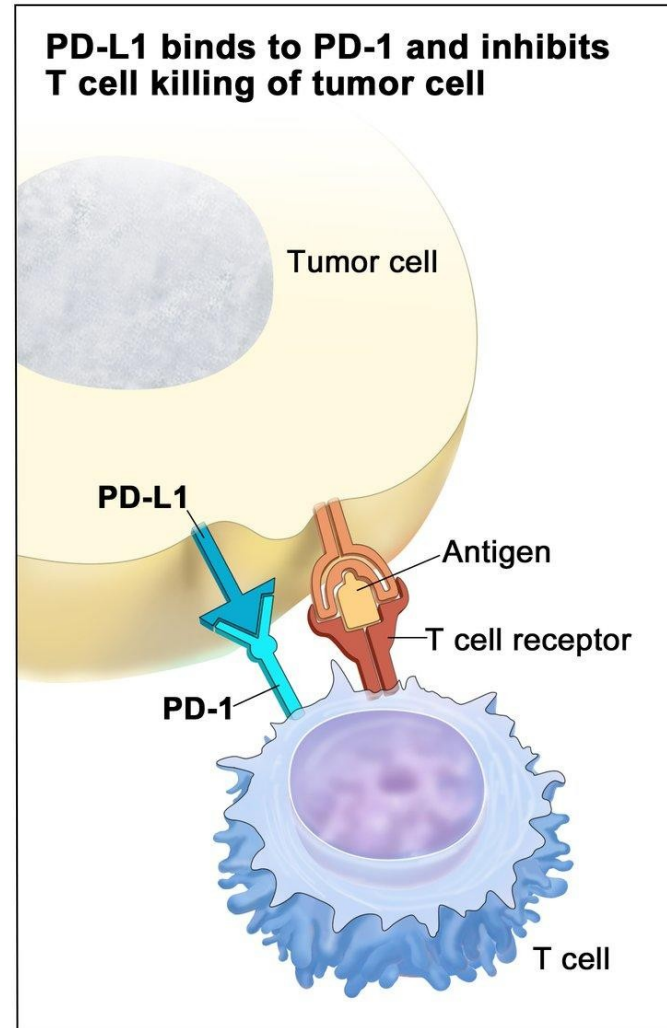
# Cancer

What is cancer?



DNA Structure

Chromosome

Histone

Nucleosome

DNA

Gene

Cell

Nucleus

Nucleotide

Nucleotide base pairs:
- Guanine
- Cytosine
- Adenine
- Thymine

© 2015 Terese Winslow LLC
U.S. Govt. has certain rights

# Cancer Immunotherapy

Immunotherapy leverages the body's immune system to combat cancer.

Immune Checkpoint Inhibitors

- Anti-PD1
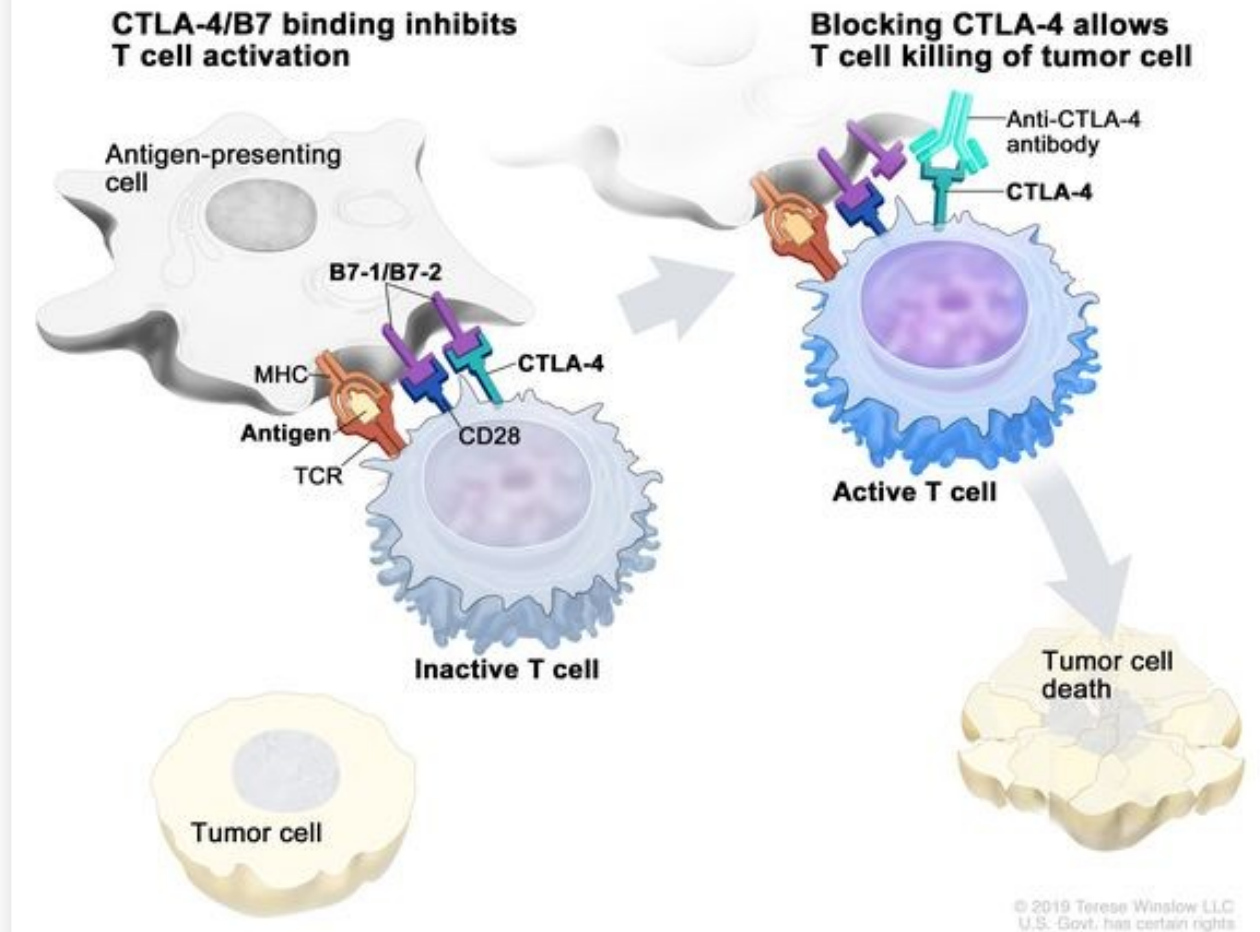  - Pembrolizumab (Keytruda)
  - Nivolumab (Opdivo)



PD-L1 binds to PD-1 and inhibits T cell killing of tumor cell

Tumor cell

PD-L1

Antigen

T cell receptor

PD-1

T cell

Blocking PD-L1 or PD-1 allows T cell killing of tumor cell

Tumor cell death

PD-L1

Anti-PD-L1

Anti-PD-1

PD-1

T cell

[1]

# Cancer Immunotherapy

- Anti-CTLA-4
  - Tremelimumab
  - ipilimumab



CTLA-4/B7 binding inhibits T cell activation

Blocking CTLA-4 allows T cell killing of tumor cell

Antigen-presenting cell

B7-1/B7-2

MHC

Antigen

TCR

CTLA-4

CD28

Inactive T cell

Tumor cell

Anti-CTLA-4 antibody

CTLA-4

Active T cell

Tumor cell death

© 2019 Terese Winslow LLC
U.S. Govt. has certain rights

[1]

# The problem

20-40% of patients [2]. On average, price patient annually is $150,000 in the USA

Side effects of Immunotherapy can be; diarrhoea, fatigue, nausea or even an autoimmune response.

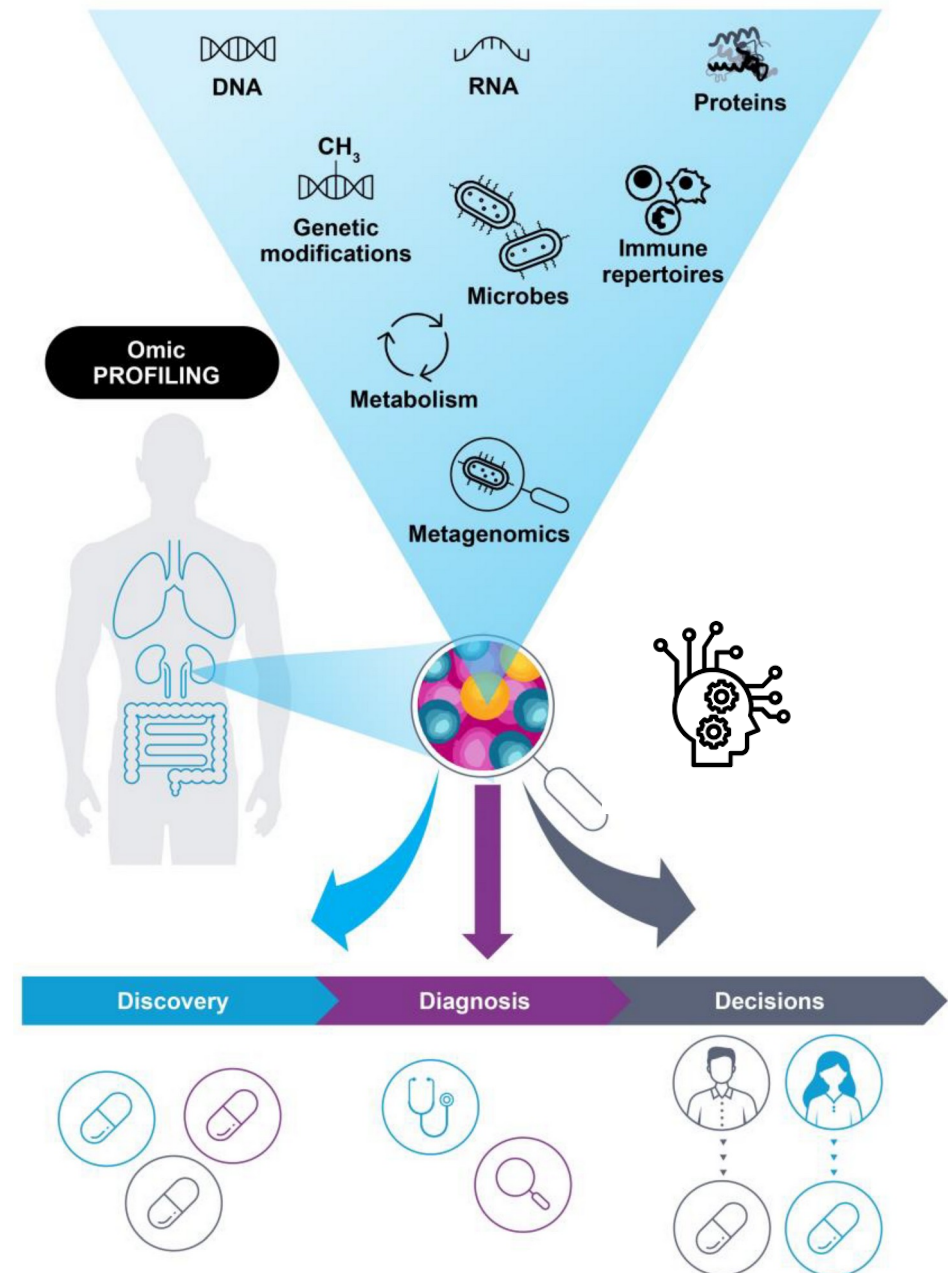Time waste and disease progression.

# The Current Solution

- **Precision Oncology**: Utilises genetic profiling and biomarkers to customise treatments for individuals.

- **Current Biomarker in Clinical Practice**: Tumour mutational burden, PD-L1 expression, etc.

- **Limitations**:
  - Costly
  - Results different across labs

# The Proposed Solution

Combine **ML with precision oncology** and train ML models on pre-treatment gene expression data from next-generation sequencing to **predict patient outcomes**.



[3]

# Transcriptomics

# Biomarker

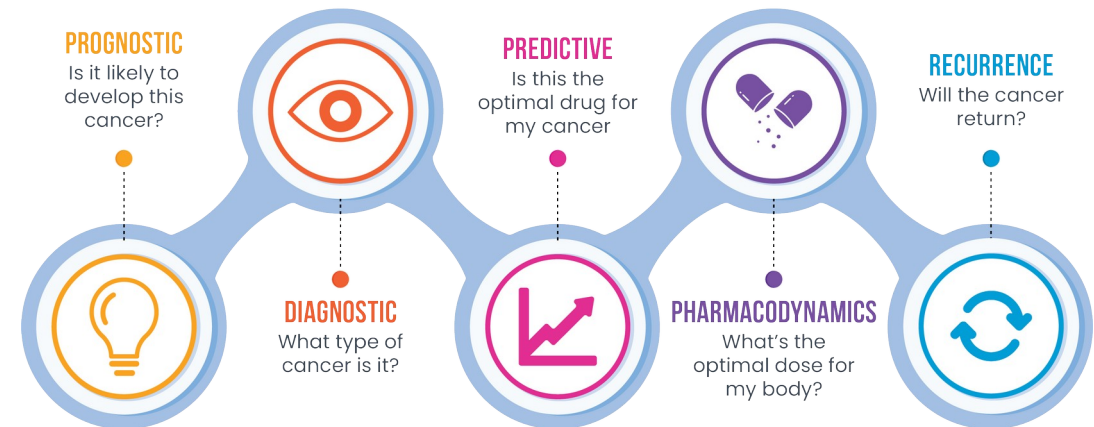## Biomarker – definition

World Stroke Academy

- A molecular, biological, or physical characteristic that indicates a specific physiologic state. It is used in clinical practice to identify risk for disease, diagnose disease and its severity, guide intervention strategies, and monitor patient responses to therapy

*Biomarkers Definitions Working Group. Clin Pharmacol Ther 2008*

'Biomarkers- methodological evaluation' by Ass Prof Ana Catarina Fonseca, University of Lisbon

## TYPES OF BIOMARKERS

**PROGNOSTIC**
Is it likely to develop this cancer?

**PREDICTIVE**
Is this the optimal drug for my cancer

**RECURRENCE**
Will the cancer return?

**DIAGNOSTIC**
What type of cancer is it?

**PHARMACODYNAMICS**
What's the optimal dose for my body?

© MEDITRIAL

# Machine Learning Pipeline

Data Collection → Pre-processing → Model Selection and Training → Model Evaluation → Model Deployment

# Machine Learning Pipeline: Data

- **Transcriptomic** data from pre-treatment tumour biopsies and blood samples.

- Collected through **extensive literature reviews**, **public databases**, and **clinical trials.**

- The pre-treatment tumour provides a baseline of the patient's disease before treatment.

# Biopsy Data

| Characteristics | I01 (N = 119) | I14 (N = 56) | I09 (N = 41) | I15 (N = 28) |
|---|---|---|---|---|
| Studies cohort | Liu D, 2019 | Riaz N, 2017 | Glide TN, 2019 | Hugo W, 2016 |
| Data source | dbGaP (accession number phs000452.v3.p1) | GSE91061 | ENA: PRJEB23709 | GSE78220 |
| Cancer type | Melanoma: 121 (100%) | Melanoma: 56 (100%) | Melanoma: 41 (100%) | Melanoma: 28 (100%) |
| Anti-PD1 received for Melanoma Cancer | Nivolumab or Pembrolizumab (not identified in each patient) | Nivolumab: 56 (100%) | Nivolumab: 9 (22%)  Pembrolizumab: 32 (78%) | Pembrolizumab: 28 (100%) |
| Number of genes in common | 7440 | 7440 | 7440 | 7440 |
| Drug response (RECIST) | | | | |
| CR   (4) | 16 (13%) | 3 (5%) | 4 (10%) | 5 (18%) |
| PR   (3) | 31 (26%) | 8 (14%) | 15 (36%) | 10 (36%) |
| SD   (2) | 16 (13%) | 19 (34%) | 6 (15%) | 0 (0%) |
| PD   (1) | 57 (47%) | 26 (46%) | 16 (39%) | 13 (46%) |
| Drug response | | | | |
| Responder | 47 (39%) | 11 (20%) | 19 (46%) | 15 (54%) |

# Machine Learning Pipeline: Data

# Machine Learning Pipeline: Data Preprocessing

# Z - score Normalization

- The purpose of Z-score normalization is to scale and centre the data such that it has a **mean of zero** and a **standard deviation of one**.

- This makes it easier to compare data that are on different scales or have different units.

**Formula:**

The formula to calculate the Z-score for a given data point is:
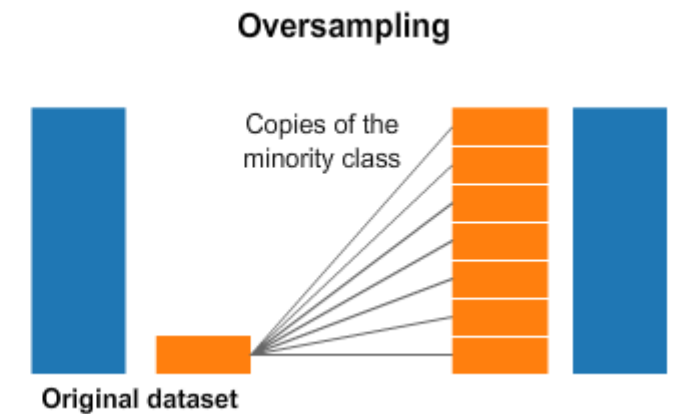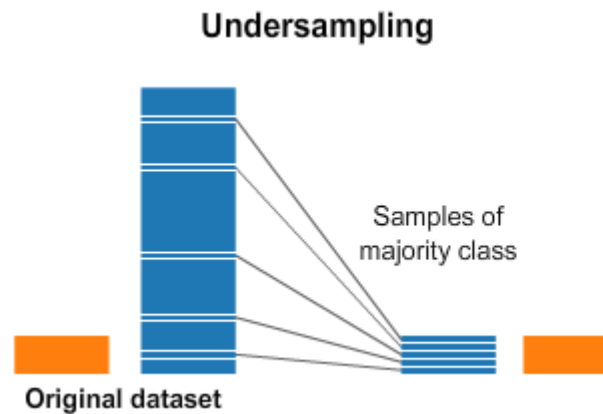
$$Z = \frac{(X - \mu)}{\sigma}$$

where:

- $Z$ is the Z-score.
- $X$ is the original data point.
- $\mu$ is the mean of the dataset.
- $\sigma$ is the standard deviation of the dataset.

# **Handling Imbalanced Datasets**

- Re-sampling techniques
  - Oversampling
  - Under-sampling



[8]

# Machine Learning Pipeline: Model Training

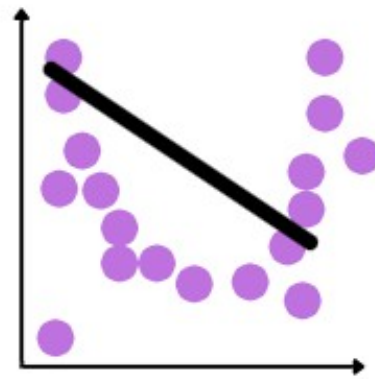# Machine Learning Pipeline: Model Training and Evaluation

Nested k Fold Cross Validation
- Hyperparameter tuning
- Feature selection

Evaluation Metrics
- Model selection
- **MCC** – Matthews Correlation Coefficient
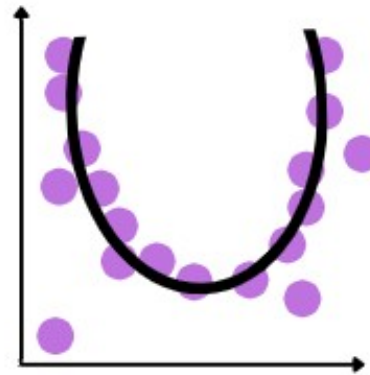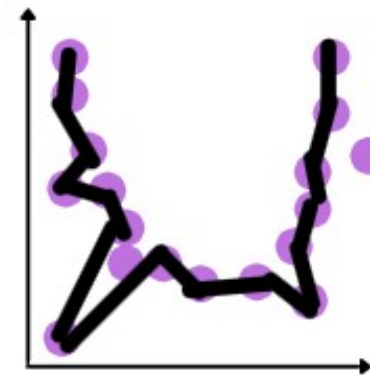- **ROC** – Receiver Operating Characteristic



Outer resampling

Estimate performance

Use tuned parameters

Inner resampling

Tune parameters

Training set outer resampling

Test set outer resampling

Training set inner resampling

Test set inner resampling

[3]

# Hyperparameter tuning



Underfitting
(model is too simple)

Optimal

Overfitting
(model is too complex and captures even noise in the data)

[8]

# Feature Selection
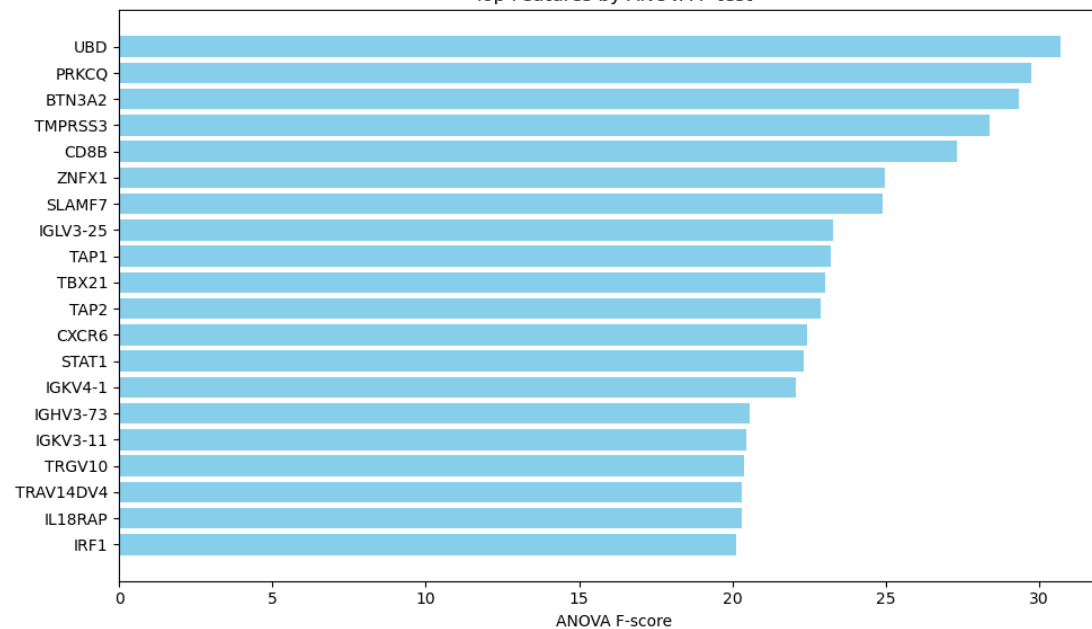
- ANOVA
- Lasso, L1 regularisation
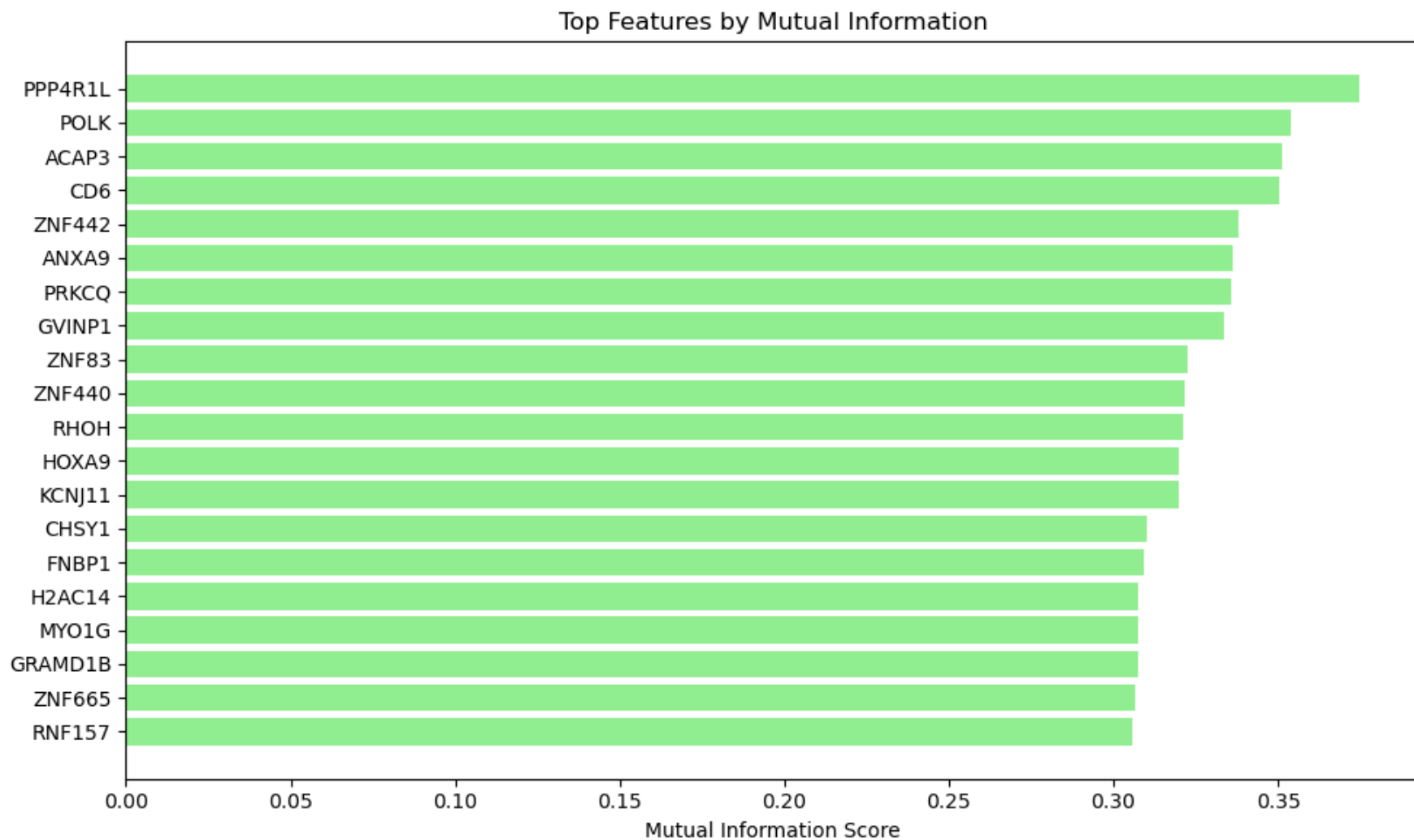- Mutual Info

# ANOVA



ANOVA testing

# Mutual Info
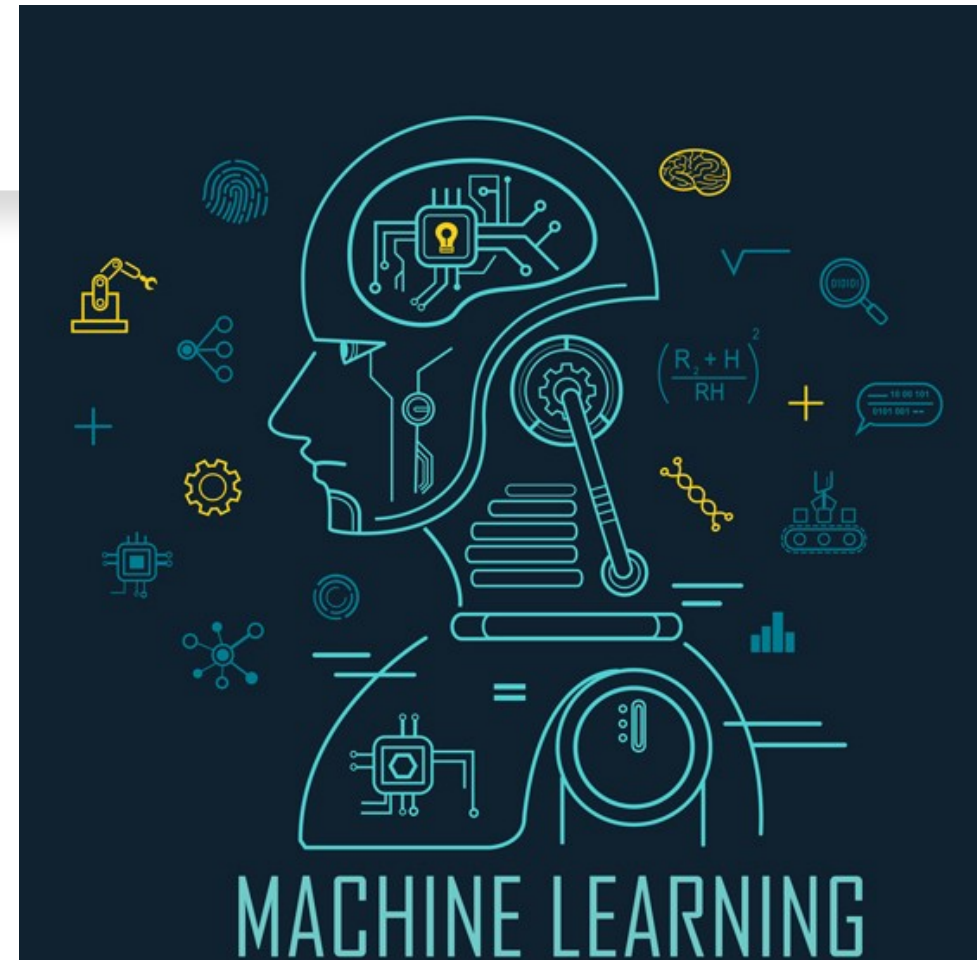
Top Features by Mutual Information

# Lasso, L1 regularisation
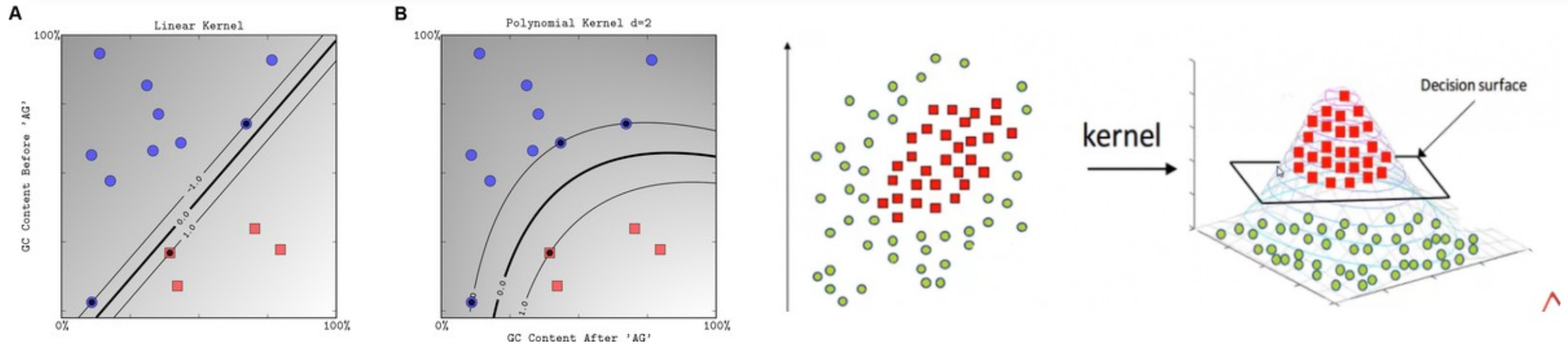


Top Features by Lasso

# Algorithms

- Support Vector Machine
- Logistic Regression
- Extreme Gradient Boost
- Random Forest
- Classification and Regression Trees



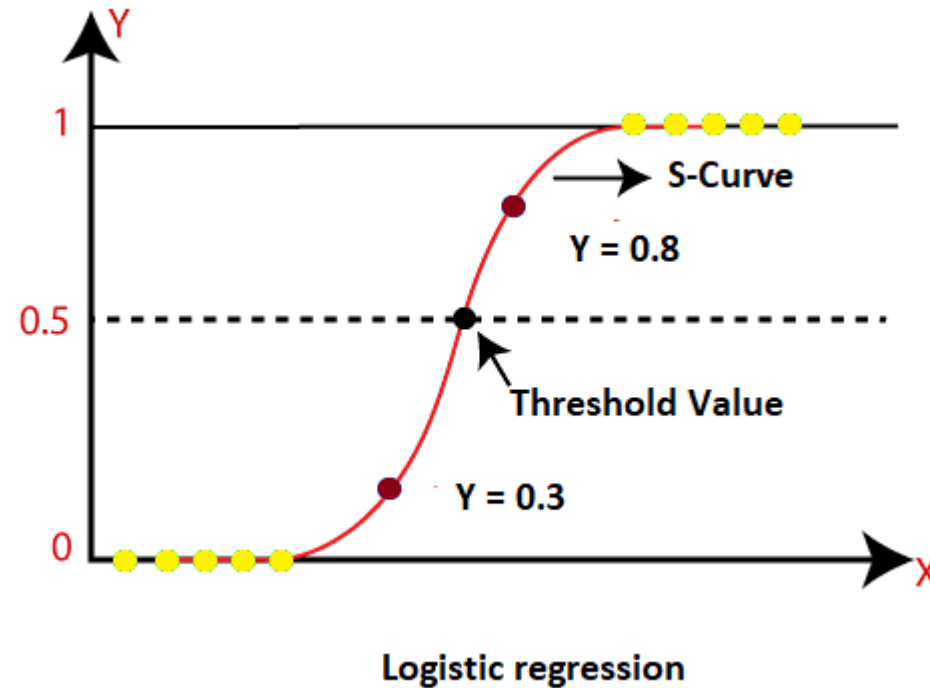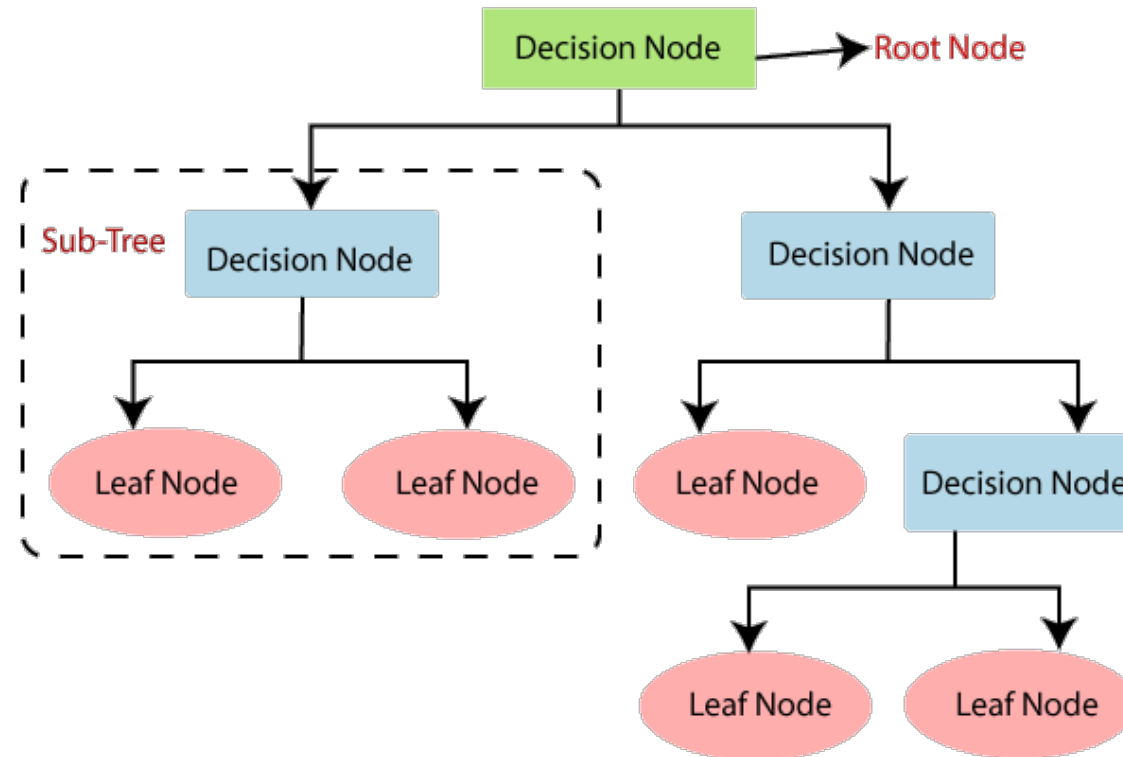[8]

# Support Vector Machine
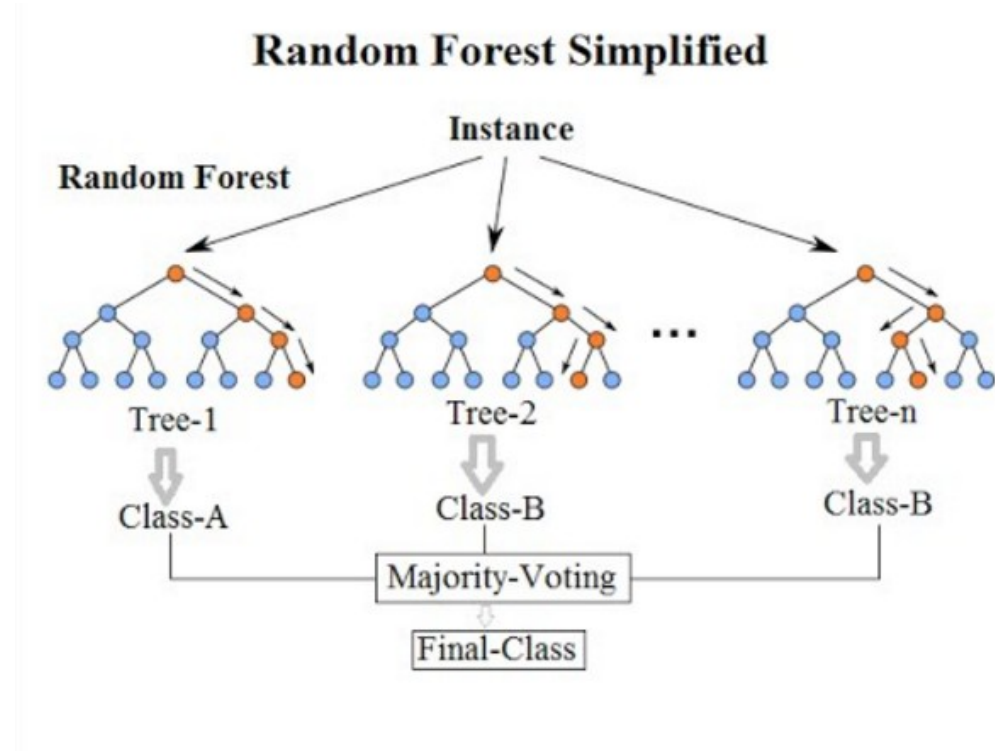


- Hyper Parameters: C, kernel

[8]

# Logistic Regression



Logistic regression

- Hyper Parameters: C, penalty

[8]

# Classification and Regression Trees



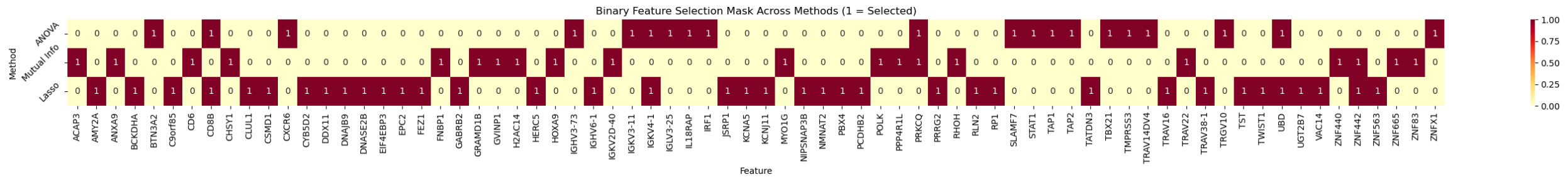- Hyper Parameters:  max depth, min sample split

[8]

# Random Forest



Random Forest Simplified

- Hyper Parameters: n estimator, max depth

[8]
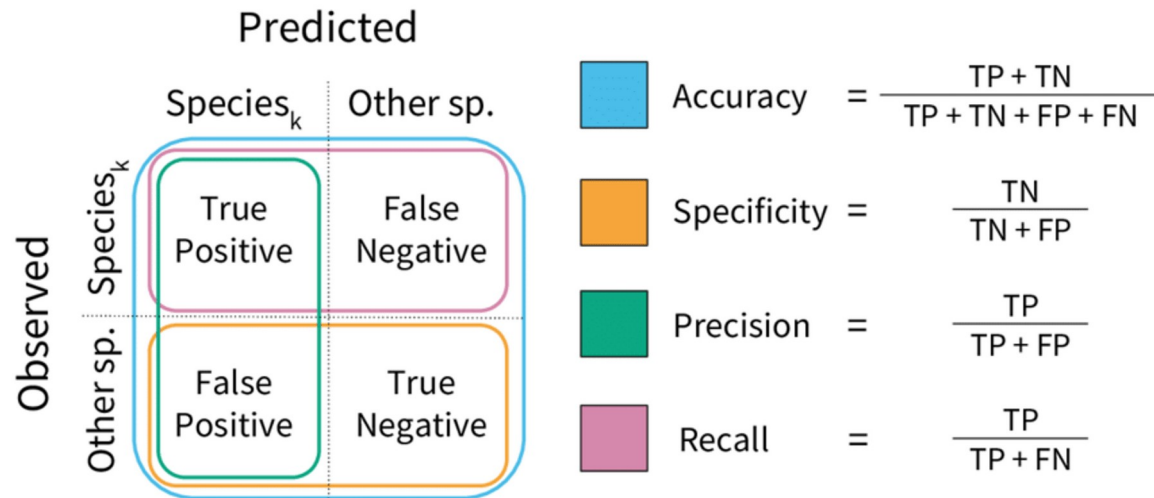
# Extreme Gradient Boost



$$\hat{y} = \sum_{k=1}^{n} f_k(x)$$

Result

- Hyper Parameters: n estimators, learning rate, max depth

[8]

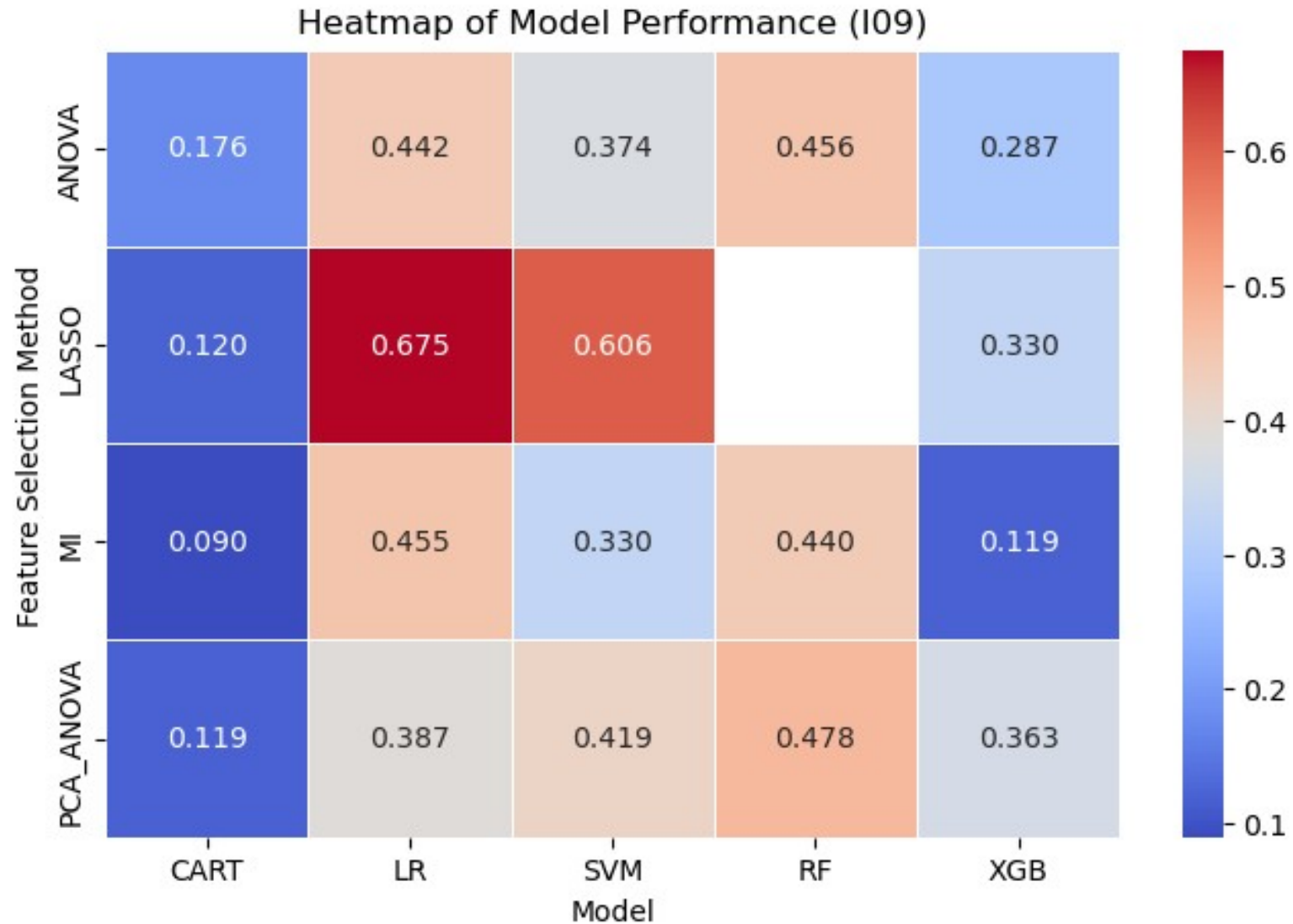# Binary Feature Selection Mask Across Methods (1 = Selected)



Binary Feature Selection Mask Across Methods (1 = Selected)

# Evaluation Metrics



[8]

# I09 Dataset: Results (MCC)



Heatmap of Model Performance (I09)

# Test on I15

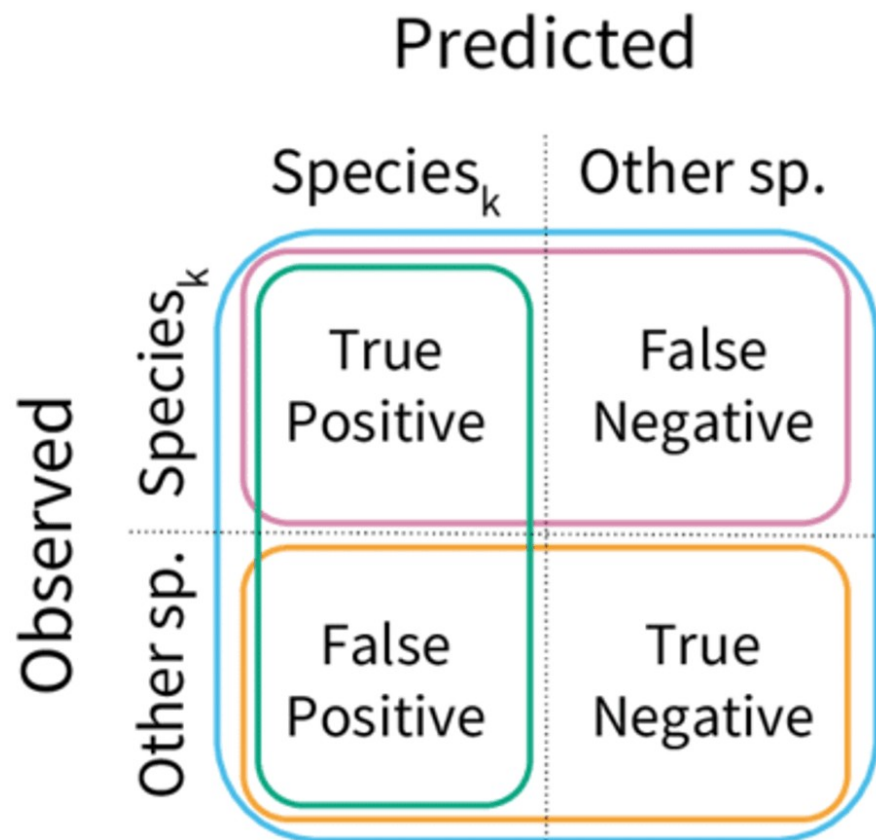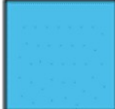| | CV MCC | Test MCC |
|---|---|---|
| LR+LASSO | 0.690228055 | -0.13092 |
| SVM+LASSO | 0.671850832 | 0.005 |

# THANK YOU! QUESTIONS

# References

1. National Cancer Institute, Winslow T. Immune Checkpoint Inhibitor (PD-1). NCI Visuals Online. June 8, 2016. Available from: https://visualsonline.cancer.gov/details.cfm?imageid=10396

2. Wolchok JD, et al. Overall survival with combined nivolumab and ipilimumab in advanced melanoma. N Engl J Med. 2017;377(14):1345-1356.

3. Bischl B, Lang M, Kotthoff L, Schiffner J, Richter J, Studerus E, Casalicchio G, Jones ZM. mlr: Machine Learning in R. J Mach Learn Res. 2016;17(170):1-5.

4. Ogunleye AZ, Piyawajanusorn C, Gonçalves A, Ghislat G, Ballester PJ. Interpretable machine learning models to predict the resistance of breast cancer patients to doxorubicin from their microRNA profiles. Adv Sci (Weinh). 2022 Aug;9(24). doi: 10.1002/advs.202201501. Epub 2022 Jul 3. PMID: 35785523; PMCID: PMC9403644.

5. Ogunleye A, Piyawajanusorn C, Ghislat G, Ballester PJ. Large-Scale Machine Learning Analysis Reveals DNA Methylation and Gene Expression Response Signatures for Gemcitabine-Treated Pancreatic Cancer. Health Data Sci. 2024;4:0108. doi: 10.34133/hds.0108. PMID: 38486621. Available from: https://pubmed.ncbi.nlm.nih.gov/38486621/

6. Tsamardinos I, Lagani V, Papagregoriou G, Tsagris M, Borboudakis G, Zenker M, et al. A perspective on automated machine learning in bioinformatics. Brief Bioinform. 2022;23(1). doi: 10.1093/bib/bbac059. Available from: https://pubmed.ncbi.nlm.nih.gov/35382509/

7. Tanner G. Kernel PCA Explained. ML Explained. January 2, 2022. Available from: https://ml-explained.com/blog/kernel-pca-explained

8. Jappinen R. Resampling strategies for imbalanced datasets. Kaggle. Available from: https://www.kaggle.com/code/rafjaa/resampling-strategies-for-imbalanced-datasets

# AI in Precision Oncology
# Case Studies



Model performance

DOX
0.56 MCC   0.80 ROC-AUC

Doxorubicin response in breast cancer patients:
Decision tree combining 4 isomiRs

DOI: 10.1002/advs.202201501

[4]

GEMZAR
0.44 MCC   0.79 ROC-AUC

Gemcitabine response in pancreatic cancer patients:
Random forest combining 4 mRNAs

DOI: 10.34133/hds.0108

[5]

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

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

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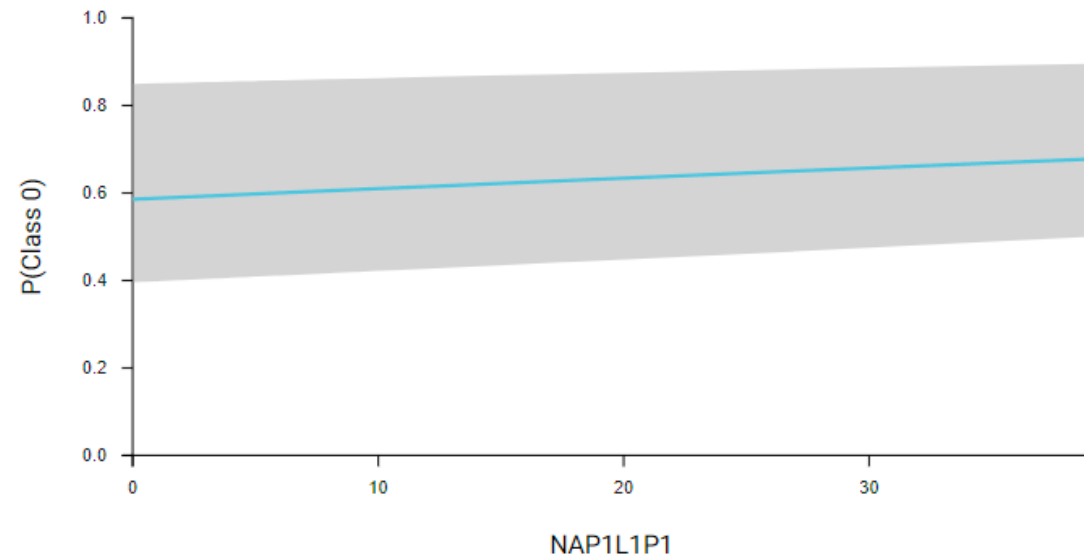
$$\text{Recall} = \frac{TP}{TP + FN}$$

$$MCC = \frac{(TP \times TN) - (FP \times FN)}{\sqrt{(TP + FP)(TP + FN)(TN + FP)(TN + FN)}}$$

# I01 Dataset Feature ICE Plot (Individual Conditional Expectation Plot)

# What is transcriptomics

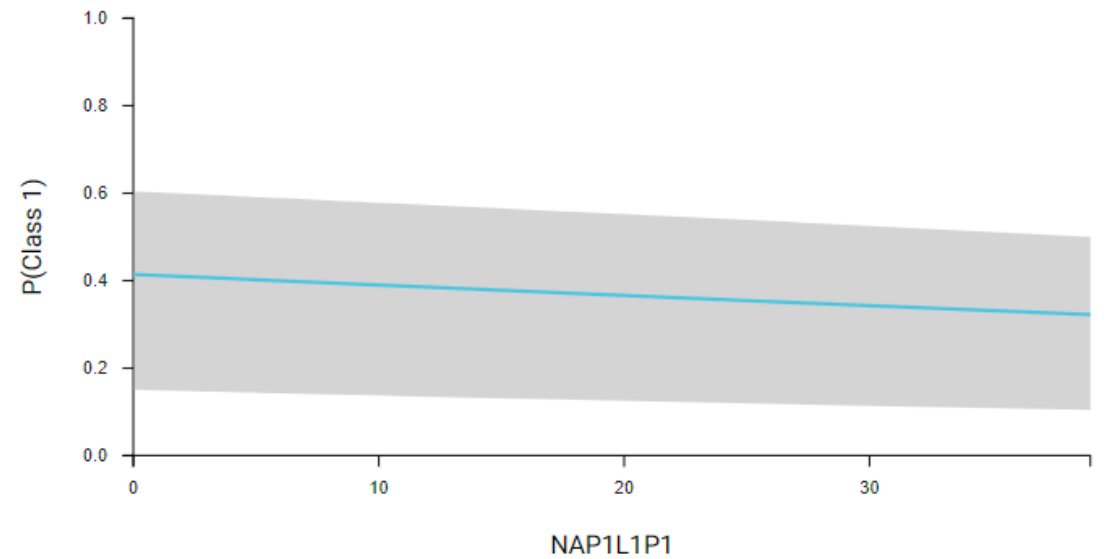**Transcriptomics** is the study of the complete set of RNA transcripts produced by the genome in a specific cell, tissue, or organism.

**TPM – Transcripts per Kilobase Million**

Normalize for gene length first, and then normalize for sequencing depth second