



About LIQOMICS

- **Academic Spin-off from Cologne:**

Focused on cutting-edge tumor diagnostics through liquid biopsy technology.

Key Service: LymphoVista

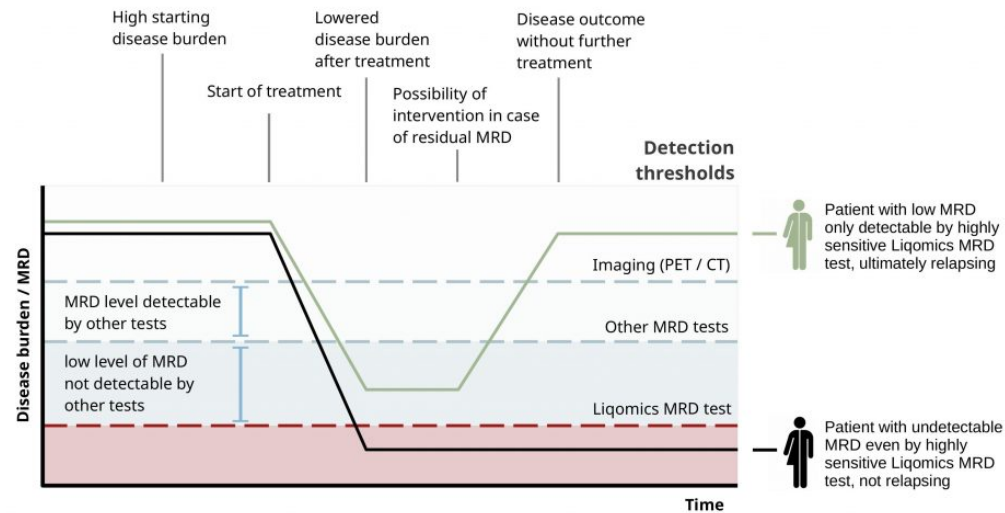
- **What It Does:**

Detects and monitors **lymphomas** and **hematological cancers**.

- **How It Works:**

- **Blood Sample Analysis:** Detects **cell-free DNA** and **circulating tumor DNA**.
- **Disease Monitoring:** Tracks **Minimal Residual Disease (MRD)** to assess treatment response.

Sensitive MRD Testing for Cancer Monitoring



Cell-Free DNA (cfDNA)

What Is cfDNA?

- **Degraded Small DNA Fragments:**
 - Size: **50 - 200 bp**, cleaved by **nucleases**.
-

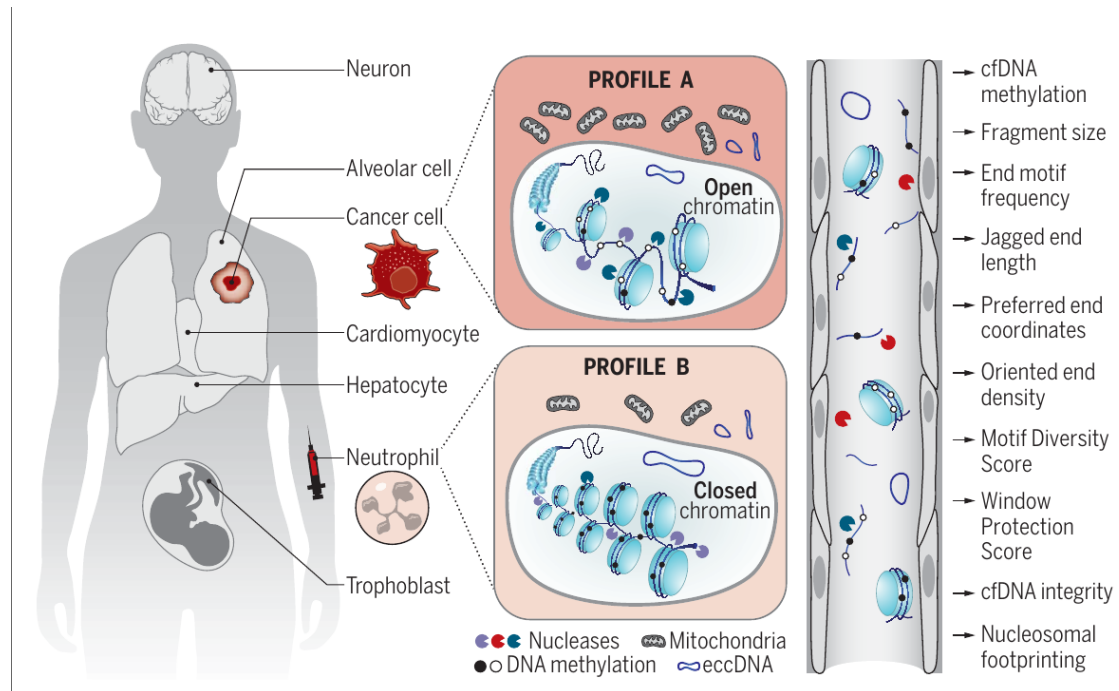
Source of cfDNA:

- **Released by Dead Cells:**
 - Occurs via **apoptosis, necrosis**, or **active secretion** from living cells.
-

Key Molecular Features:

- **Tissue-Specific Signatures:**
 - **Fragment Size:** Reflects tissue-specific fragmentation patterns.
 - **Methylation Status:** Epigenetic modifications characteristic of cell types.
 - **End Motifs:** Specific DNA ends reflecting nuclease activity.

Chromatin Organization and Nuclease Activity Define cfDNA Signatures



Projects Overview

1. Hodgkin Lymphoma Project

- **Data Source:**
 - Cell-free DNA sequencing from **Hodgkin Lymphoma Patients**.
 - Samples collected **before** and **after two cycles of chemotherapy**.
 - Additional **relapse information** included.
 - **Goal:**
 - Build a **Machine Learning Model** to **predict relapse** after treatment.
 - **Result:**
 - **Challenge:** Insufficient sample size for a reliable predictive model.
-

Projects Overview

2. Solid Tumors (Finale_DB Project)

- **Data Source:**
 - Cell-free DNA sequencing from **4 Publications:**
 - **Jiang_2015, Cristiano_2019, Snyder_2016, Sun_2019.**
 - **>800 Samples:** Includes both **healthy** and **cancer** samples.
 - **Goal:**
 - Develop a **Machine Learning Model** to **distinguish cancer from non-cancer** for solid tumors.
-

Note: Detailed results from this project will be presented on the following slides.

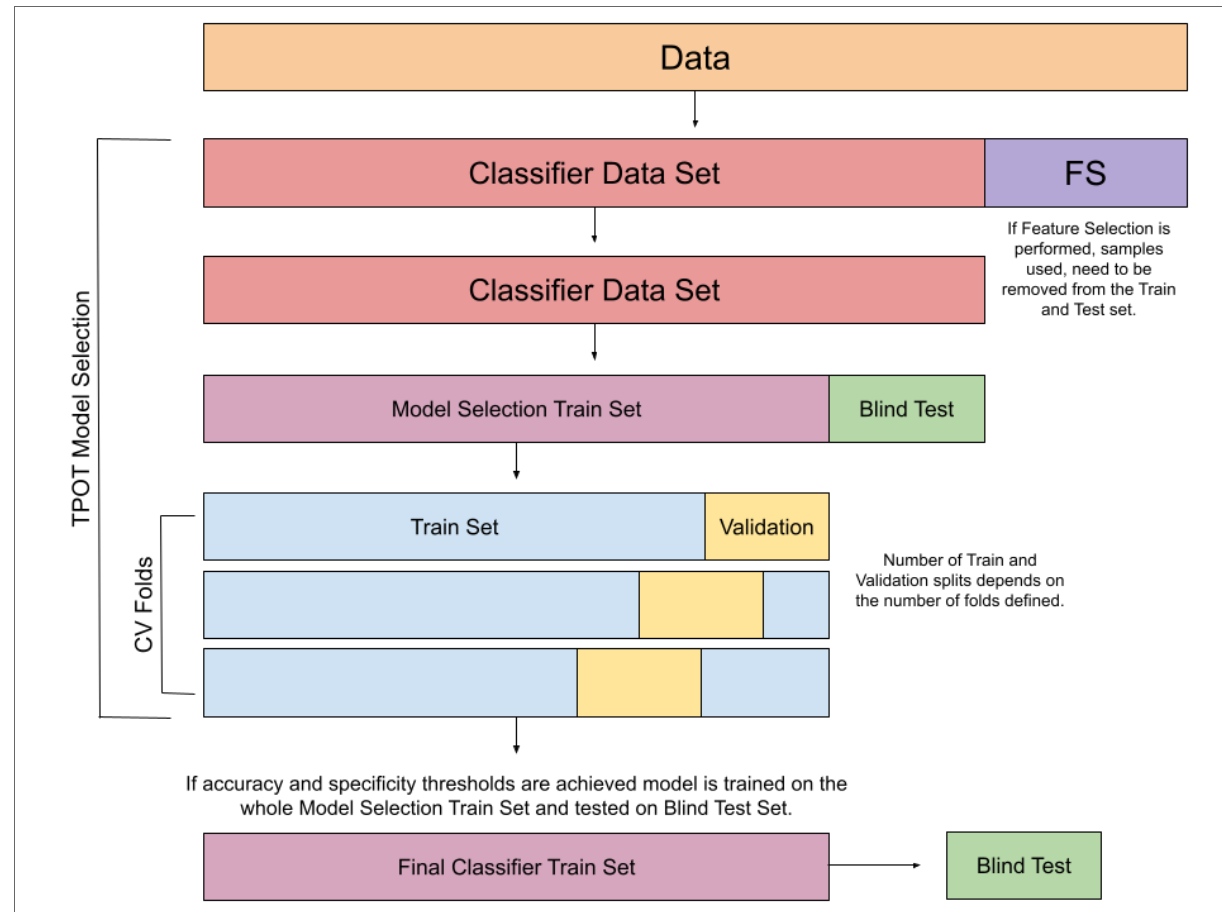
Summary of Data Sources and Samples

Source	Diseases Covered	Total Samples
Cristiano_2019	Bile_duct_cancer	536
	Breast_cancer	
	Colorectal_cancer	
	Duodenal_cancer	
	Gastric_cancer	
	Healthy	
	Lung_cancer	
	Ovarian_cancer	
	Pancreatic_cancer	
Jiang_2015	Cirrhosis	225
	Healthy	
	Hepatitis_B	
	Liver_cancer	

Source	Diseases Covered	Total Samples
Snyder_2016	Bladder_cancer	58
	Breast_cancer	
	Colorectal_cancer	
	Esophageal_cancer	
	Head_and_neck_cancer	
	Healthy	
	Inflammatory_bowel_disease	
	Kidney_cancer	
	Liver_cancer	
	Lung_cancer	
	Ovarian_cancer	
	Pancreatic_cancer	
	Prostate_cancer	
	Skin_cancer	
Sun_2019	Systemic_lupus_erythematosus	29
	Testicular_cancer	
	Uterine_cancer	
Sun_2019	Colorectal_cancer	29
	Liver_transplant	

Pipeline

1. Downsampling
2. Feature Normalization
3. Feature Selection
4. Machine learning (TPOT)



Downsampling Process

Key Steps:

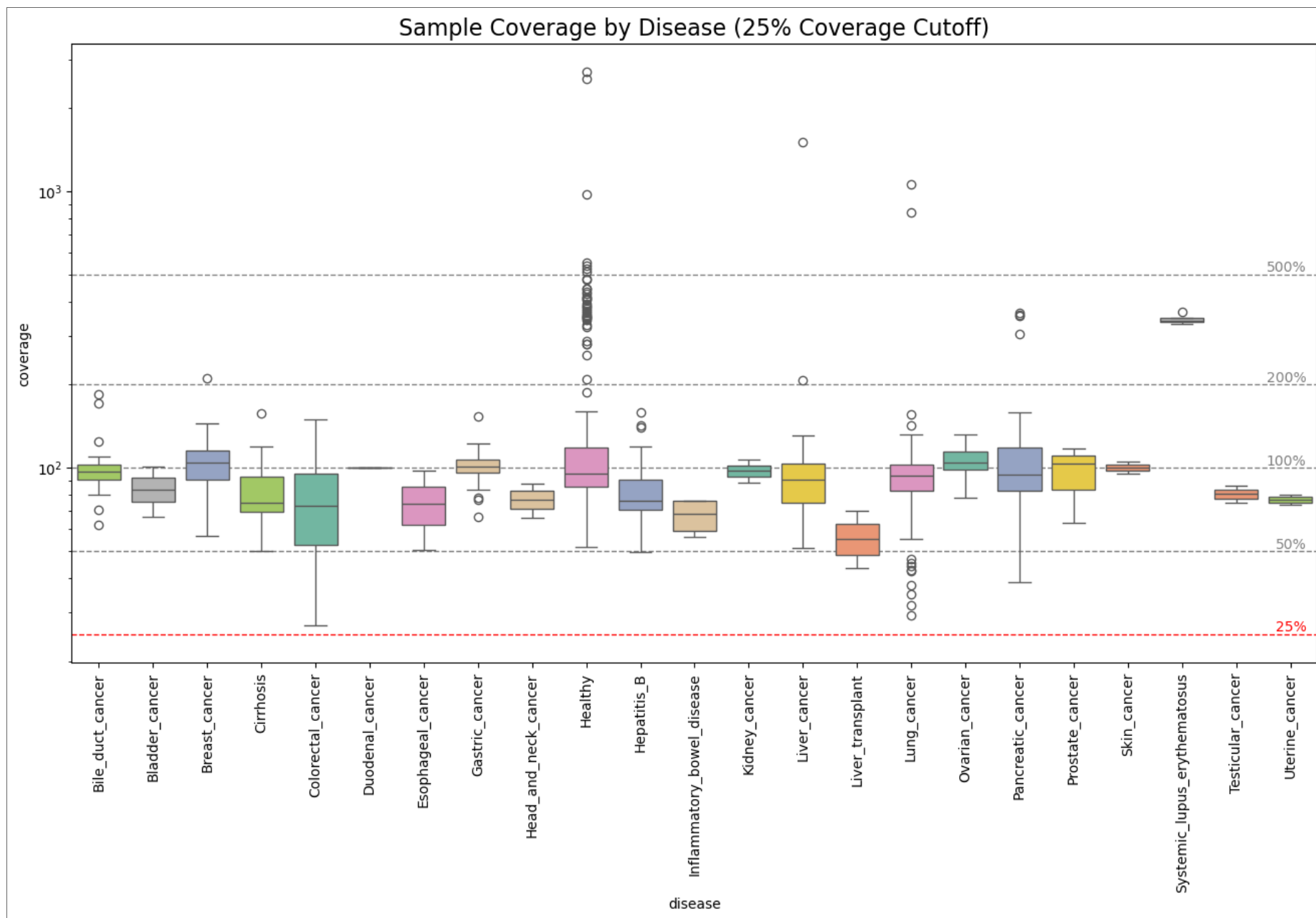
1. Handling Coverage Variability:

- Coverage between samples **varies significantly**, affecting comparability.

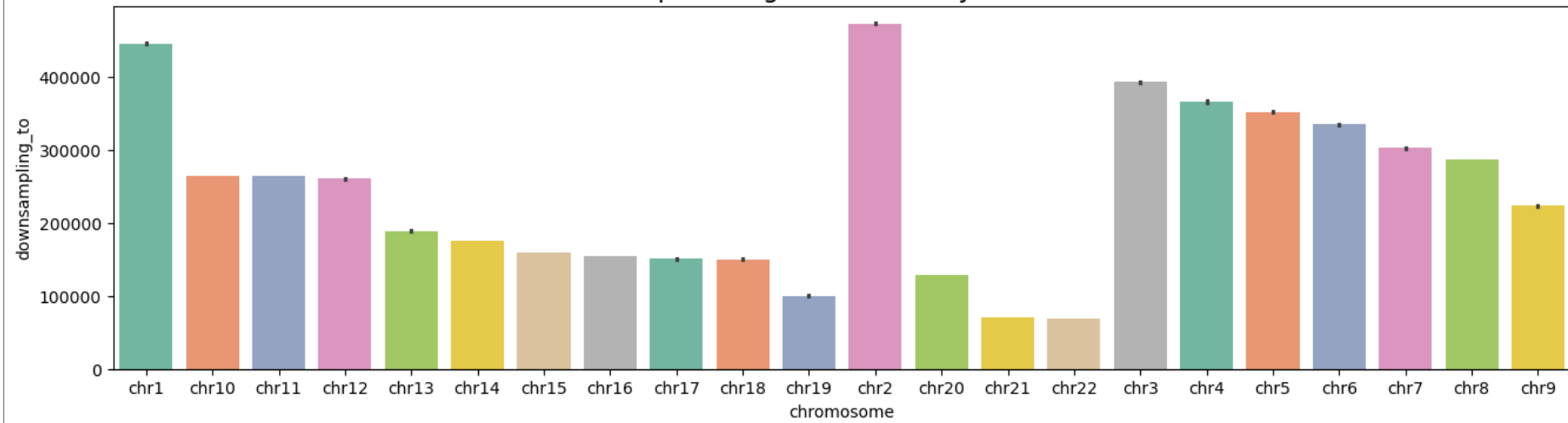
2. Equalization:

- **Downsample all samples** to the same number of **DNA fragments** to ensure consistency in analysis.

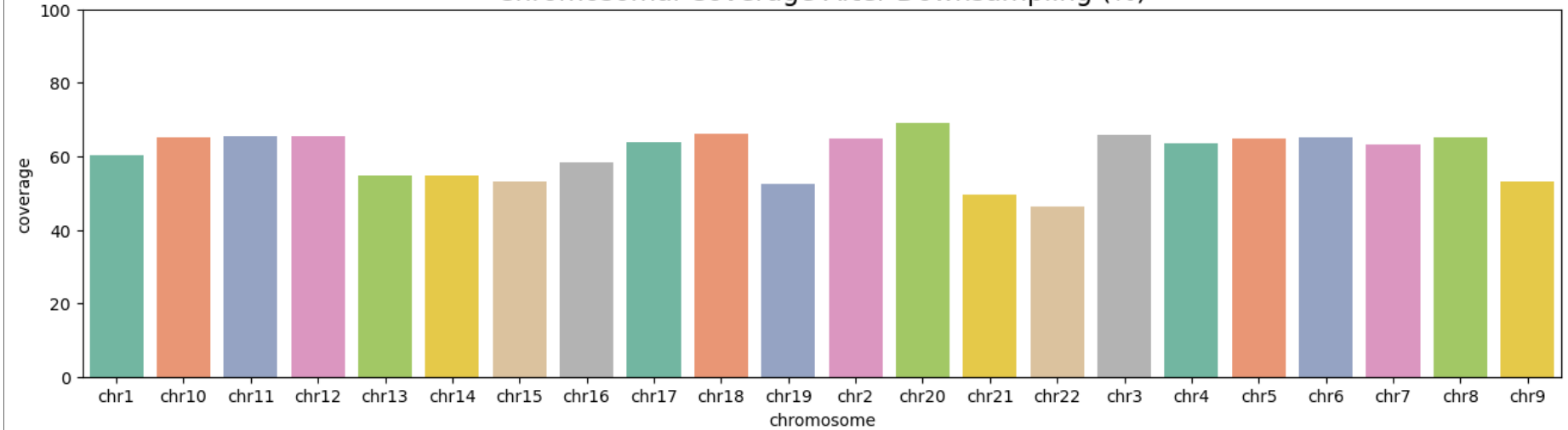
Number of total samples after downsampling: 848



Downsampled Fragment Count by Chromosome



Chromosomal Coverage After Downsampling (%)



Feature Normalization Process

Key Steps:

1. Window Sampling:

- Process **5 million windows** on each chromosome for accurate feature representation.

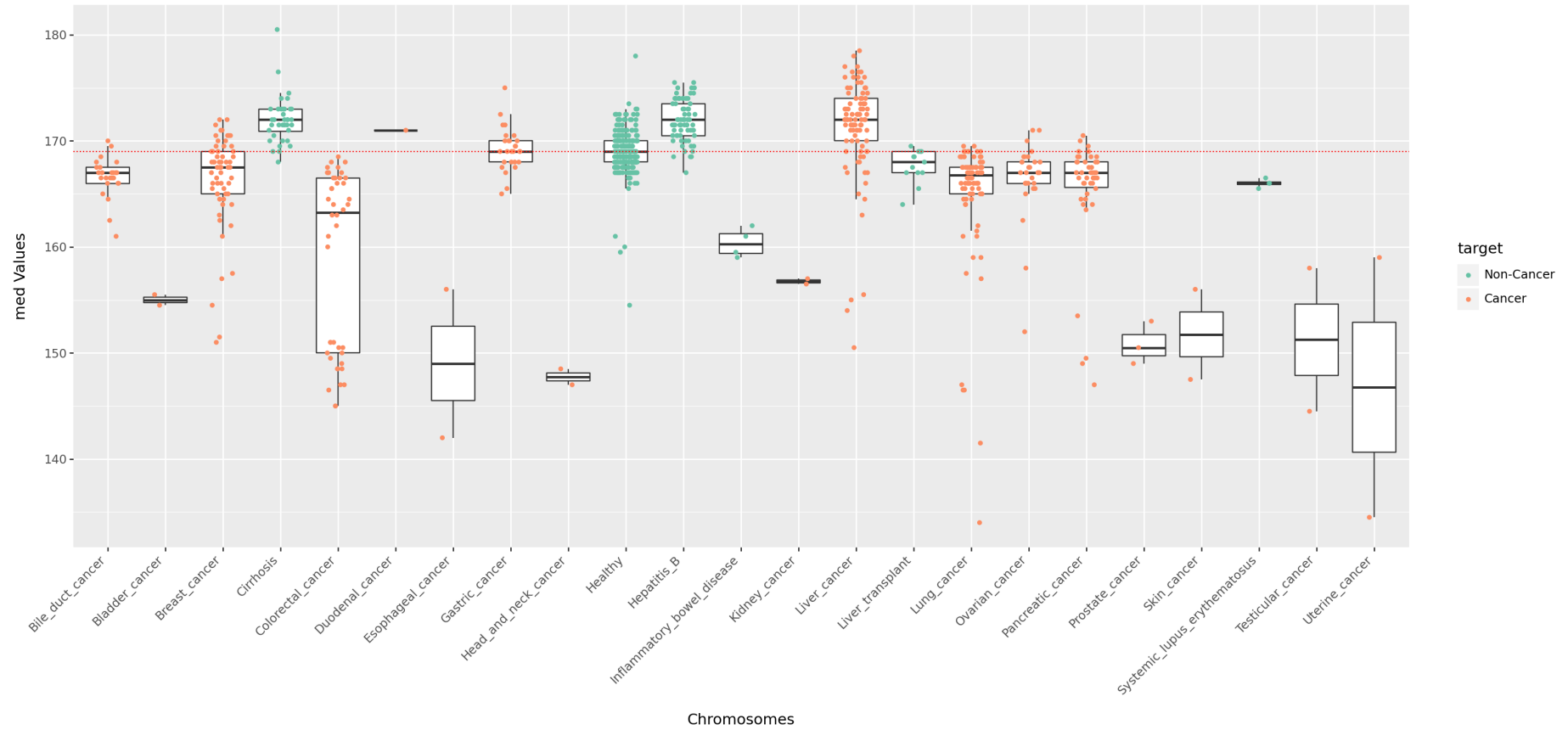
2. Normalization:

- Apply **feature normalization** to prepare data for **machine learning** models.

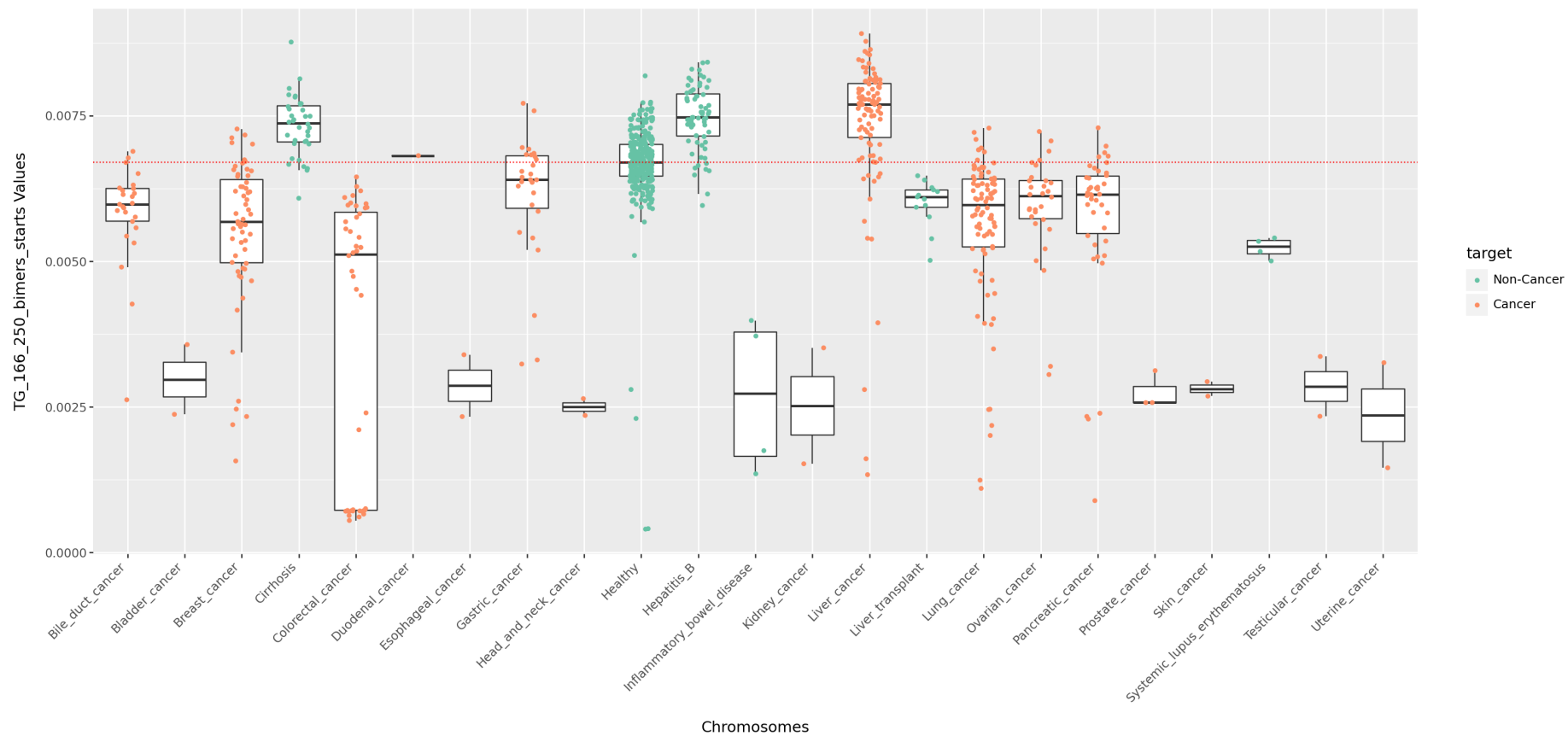
3. Data Readiness:

- Ensure features are **scaled** and **standardized** for better model performance.

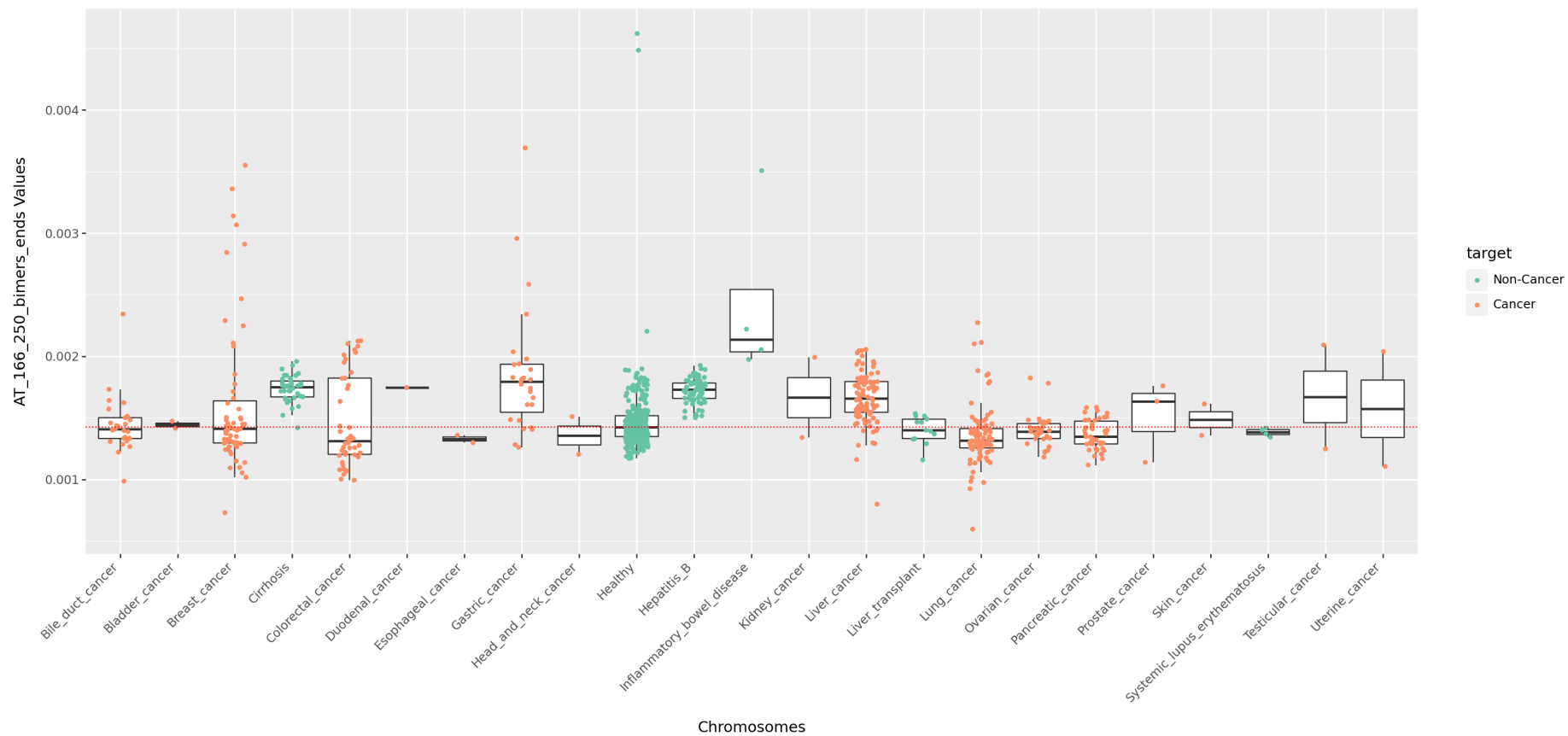
Dot plot of med values for each sample across diseases



Dot plot of TG_166_250_bimers_starts values for each sample across diseases



Dot plot of AT_166_250_bimers_ends values for each sample across diseases



Feature Selection Process

Key Steps:

1. Sample Selection:

- Use **5%** of samples (42 out of 828) for feature selection to avoid overfitting.

2. Significance Testing:

- Apply the **Mann-Whitney U Test** to remove **non-significant features** ($p > 0.05$).

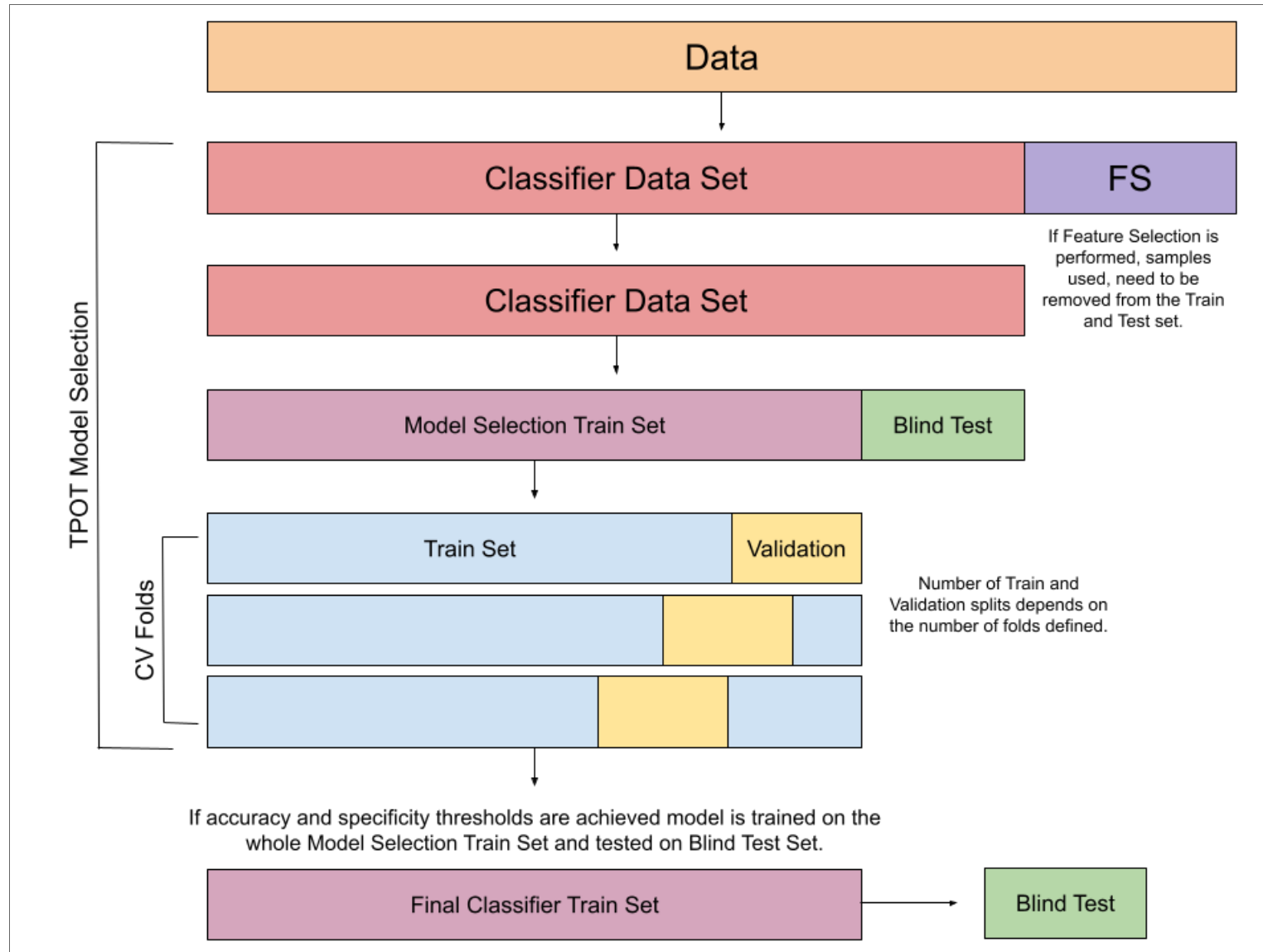
3. Correlation Handling:

- **Group correlated features** (threshold: $|r| > 0.8$) and retain the **most significant feature** from each group.
-

Result:

- **Selected Features: 78 Features** out of **76,801 Features** (0.1% selected)

TPOT



Classifier Performance on Training Data

BEST CLASSIFIERS AFTER TRAINING

Classifier	HP_count	Accuracy	Sensitivity	Specificity	Weighted_Score
BernoulliNB	475	0.581818	0.240984	0.928333	0.652825
ExtraTreesClassifier	55	0.788430	0.660656	0.918333	0.815049
RandomForestClassifier	228	0.814876	0.762295	0.868333	0.825830
XGBClassifier	459	0.806612	0.798361	0.815000	0.808331

SELECTED CLASSIFIERS FOR FULL TRAINING

Classifier	Count
RandomForestClassifier	151
XGBClassifier	8

KEY INSIGHTS:

- **Best Performers:** RandomForestClassifier and XGBClassifier
- **Next Step:** Evaluate the selected classifiers on the **blind test set** for performance validation.

Feature Importance for Training Data

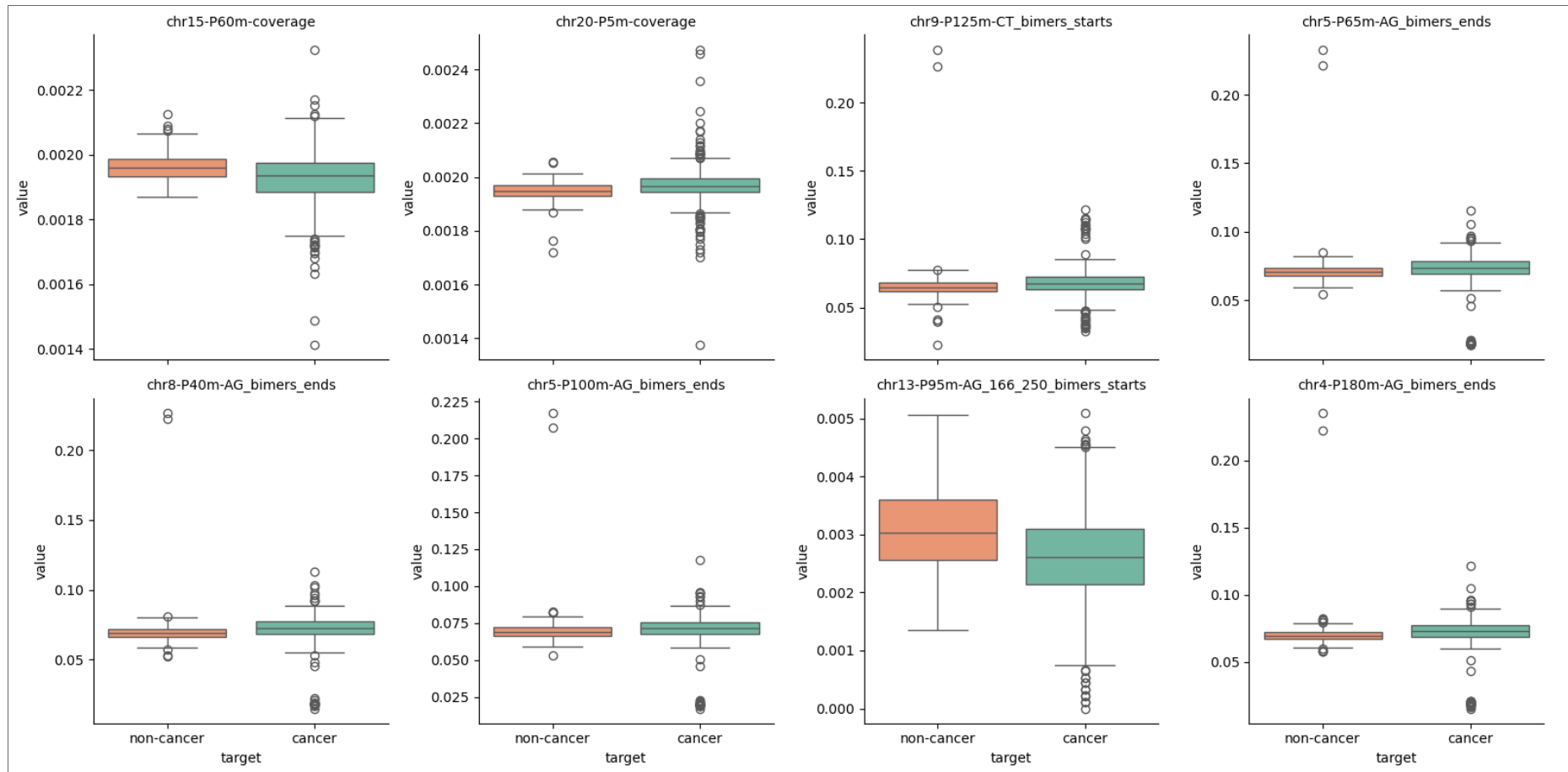
OVERLAP OF TOP 10 IMPORTANT FEATURES BY CLASSIFIER

Feature	RandomForestClassifier_228	XGBClassifier_459
chr15-P60m-coverage	1	1
chr20-P5m-coverage	2	6
chr9-P125m-CT_bimers_starts	3	4
chr5-P65m-AG_bimers_ends	4	2
chr8-P40m-AG_bimers_ends	5	8
chr5-P100m-AG_bimers_ends	7	5
chr13-P95m-AG_166_250_bimers_starts	8	9
chr4-P180m-AG_bimers_ends	10	3

KEY INSIGHTS:

- **Top Feature for Both Models:** chr15-P60m-coverage
- **Different Feature Rankings:** Each classifier relies on slightly different features.
- **Outlook:** Analyze **biological significance** of top-ranked features.

VISUALIZING TOP 10 IMPORTANT FEATURES



Classifier Performance on Blind Test Set

RESULTS OVERVIEW:

The following classifiers achieved perfect performance in predicting blind test set samples.

Classifier	Accuracy	Sensitivity	Specificity	Confusion_Matrix
RandomForestClassifier_228	1.000000	1.000000	1.000000	[[81 0] [0 81]]
XGBClassifier_459	1.000000	1.000000	1.000000	[[81 0] [0 81]]

IMPORTANT NOTE:

- **Testing on a larger sample size is needed** to obtain more **realistic accuracy values**.

Internship Takeaways

What I Learned During My Internship

TECHNICAL & ANALYTICAL SKILLS:

- **Python Development:** Advanced proficiency in Python coding and pandas for data analysis.
 - **Linux Environment:** Comfortable working in Linux-based systems for data processing.
 - **NGS Pipeline Development:** Experience in next-generation sequencing (NGS) workflow design and implementation.
 - **Fragmentomics & Feature Engineering:** Expertise in fragmentomics, feature extraction, and feature selection.
 - **Machine Learning & Model Development:** Applied ML techniques using TPOT and scikit-learn for predictive modeling.
-

Internship Takeaways

What I Learned During My Internship

BUSINESS & PROFESSIONAL INSIGHTS:

- **Start-Up Environment:** Gained insights into start-up dynamics, project management, and business development in a biotech setting.
-

Takeaway:

This internship has significantly enhanced both my **technical** and **professional** skills, preparing me for future challenges in **computational biology**, **bioinformatics**, and **machine learning**.