# Predicting Cancer Patient Response to ICI therapy using Machine Learning

## INTRODUCTION

ICI is a type of immunotherapy that treats cancer by blocking immune checkpoint proteins that prevent immune responses from destroying healthy cells. Cancer cells exploit these checkpoints to "hide" from the immune system, and ICI disrupts this evasion, allowing immune cells to attack the tumor. Patient response to ICI varies signficantly and is not consistent and prediciting responses is extremely helpful.

## METHODOLOGY

**Dataset**
Reused the public melanoma ICI dataset GSE120575 from the PRECISE study.

**Preprocessing**
Removed "low-quality" cells and genes expressed in <3% of cells.
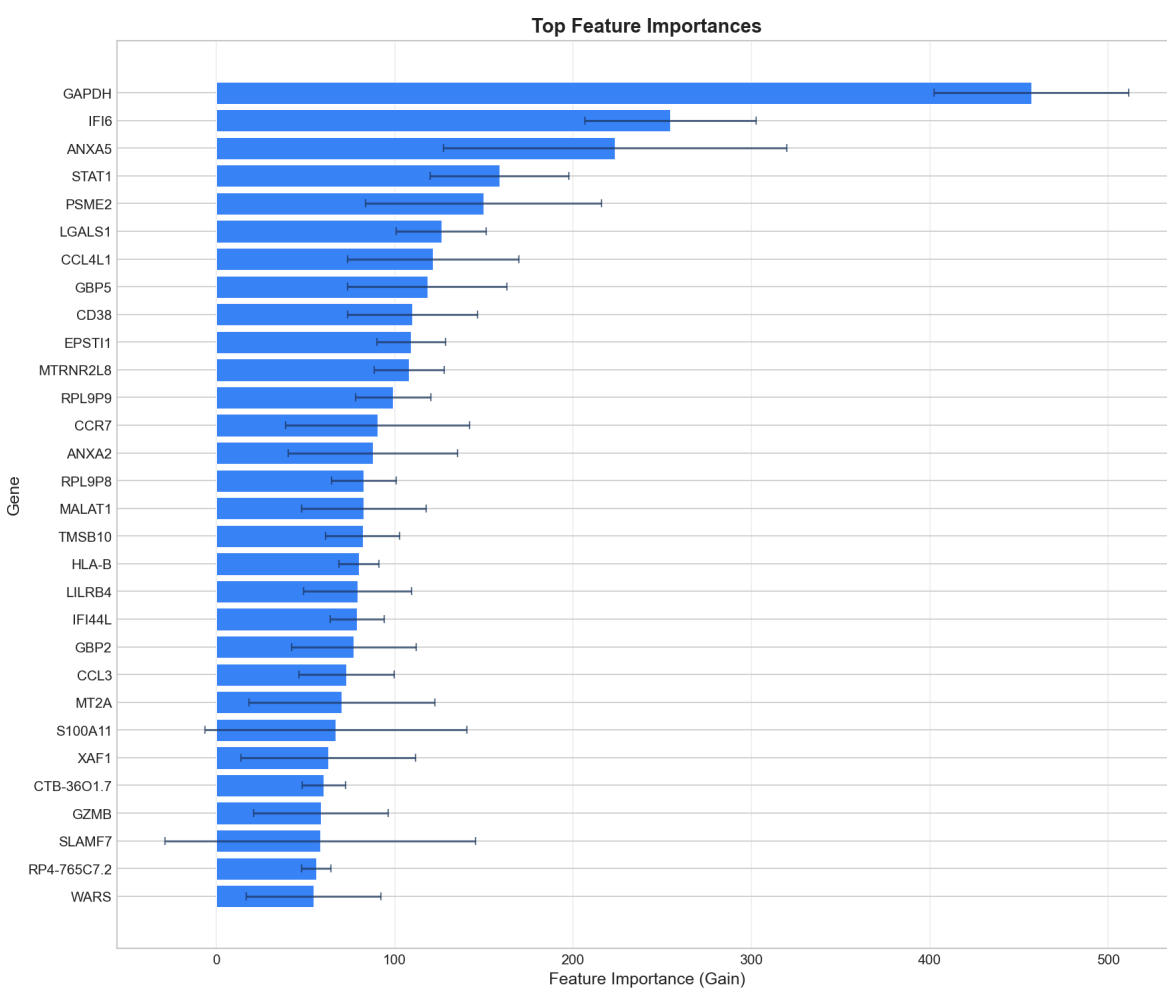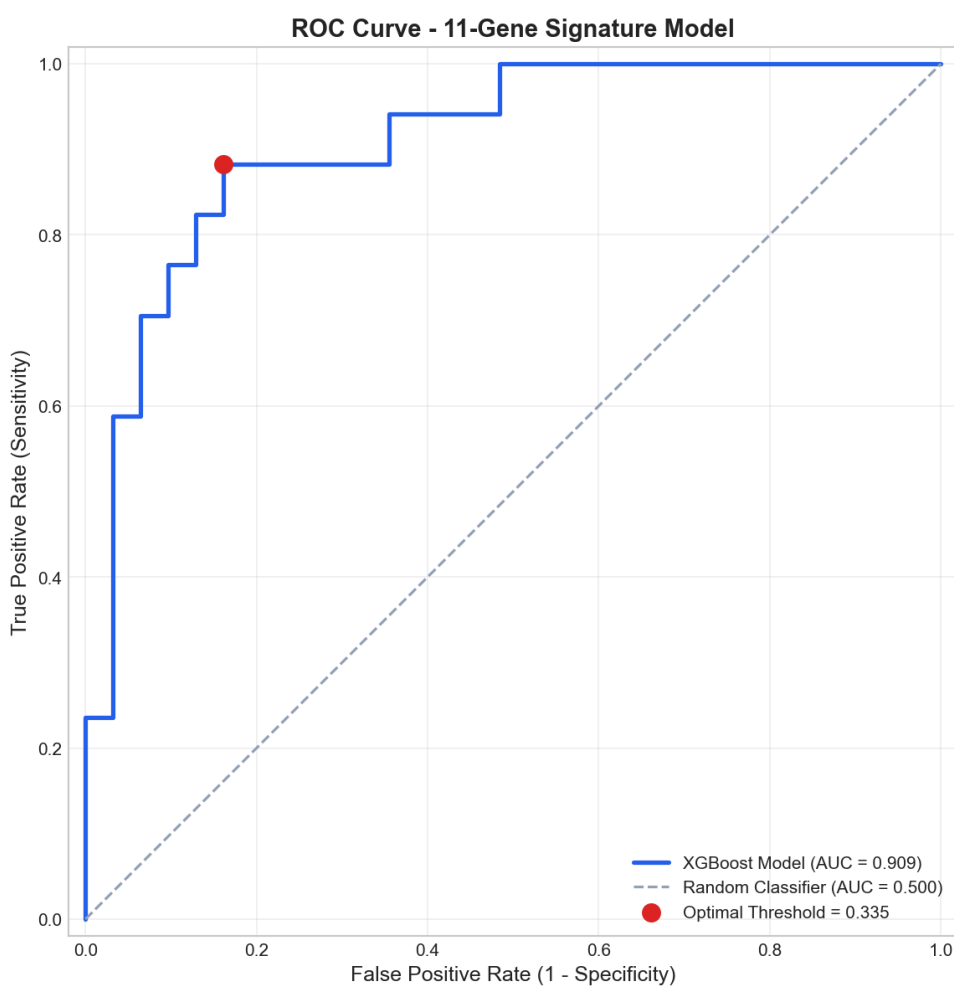
**Training & Analysis**
An XGBoost classifier was trained to predict ICI response using leave-one-patient-out cross-validation

## RESULT

- Baseline model performance: The full-gene XGBoost model (12,785 genes) achieved a patient-level AUC of 0.77 using leave-one-patient-out cross-validation.

- 10 out of an 11-gene signature, uncovered in the original PRECISE paper as predictive of patient response across different cancer types, were present.
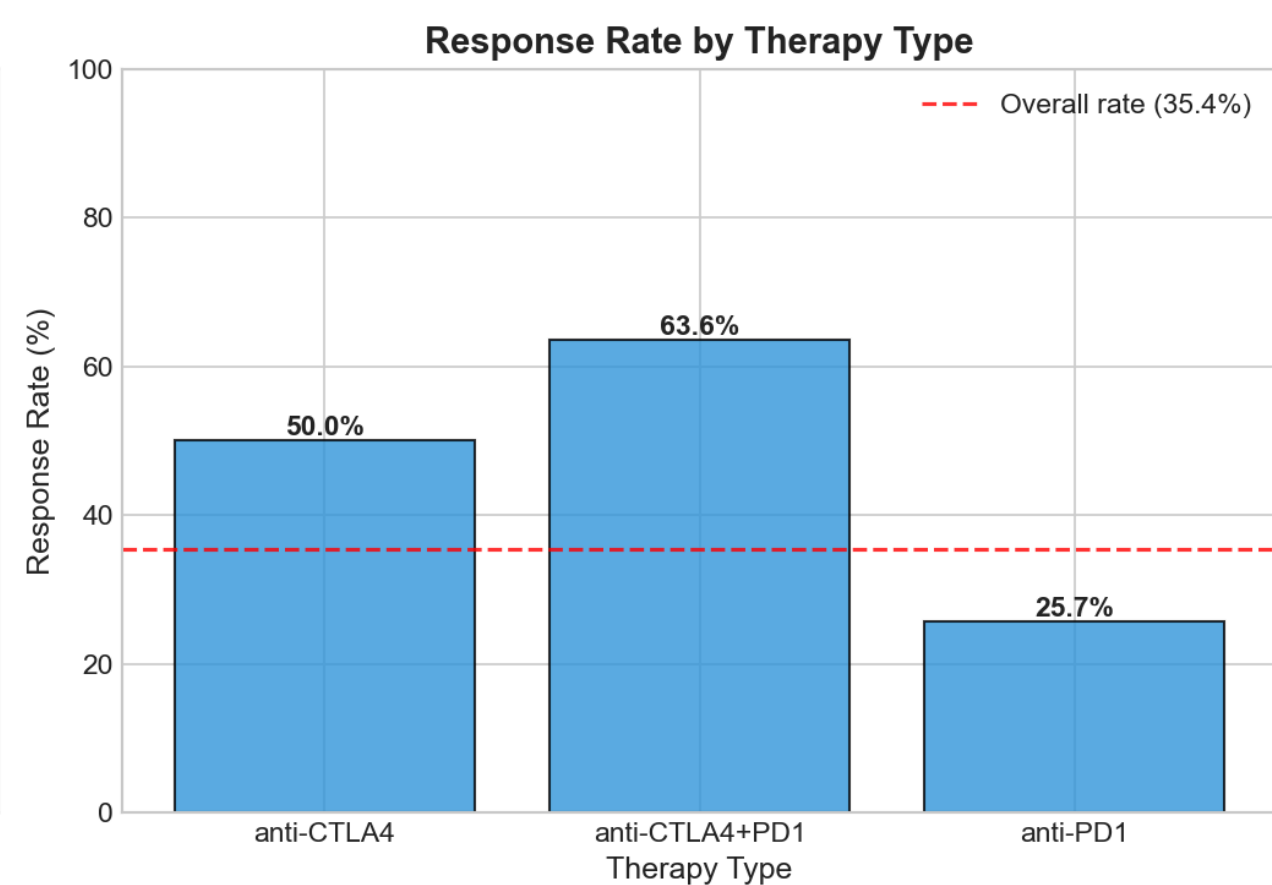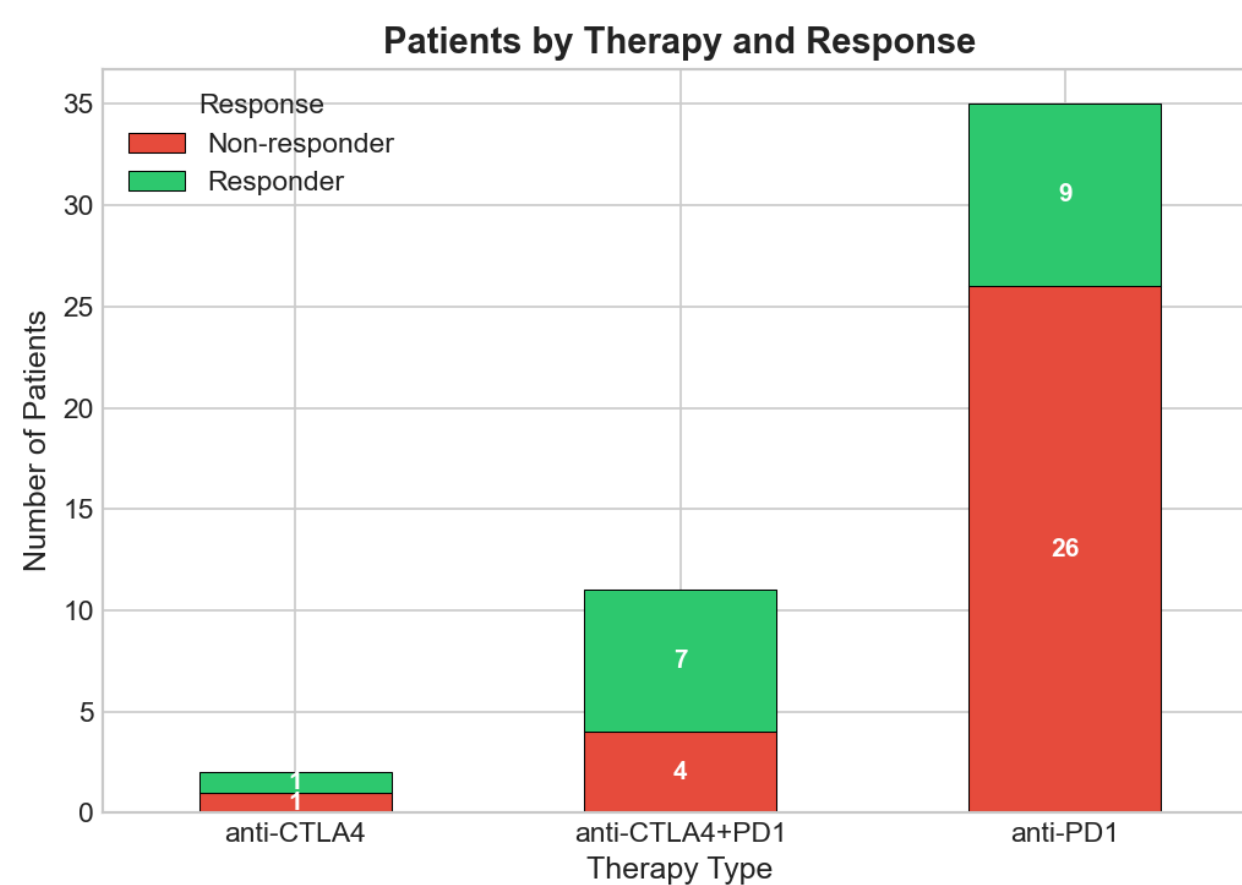
## CONCLUSION

ICI response can be predicted from single-cell immune profiles using gradient-boosted trees, and a small, biologically meaningful gene panel captures most of the predictive signal. Our reimplementation closely aligns with the PRECISE 11-gene signature, reinforcing its robustness across modeling choices.
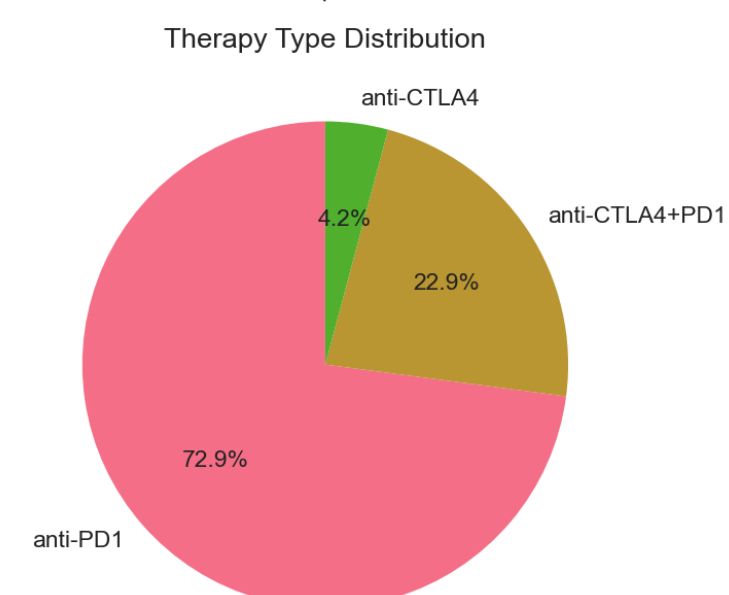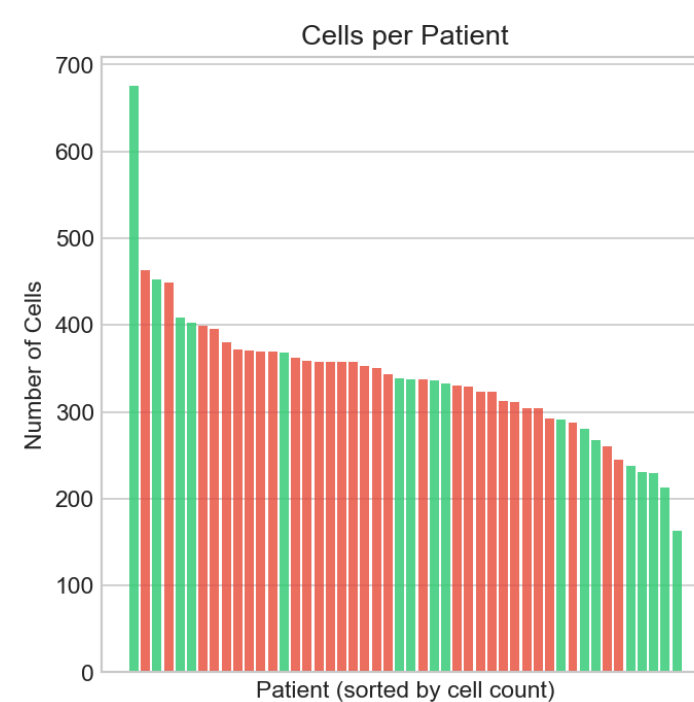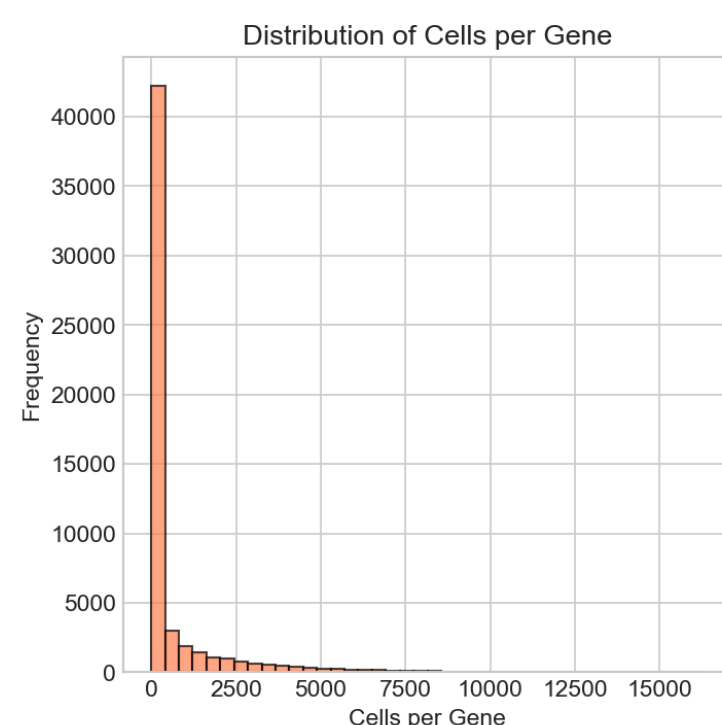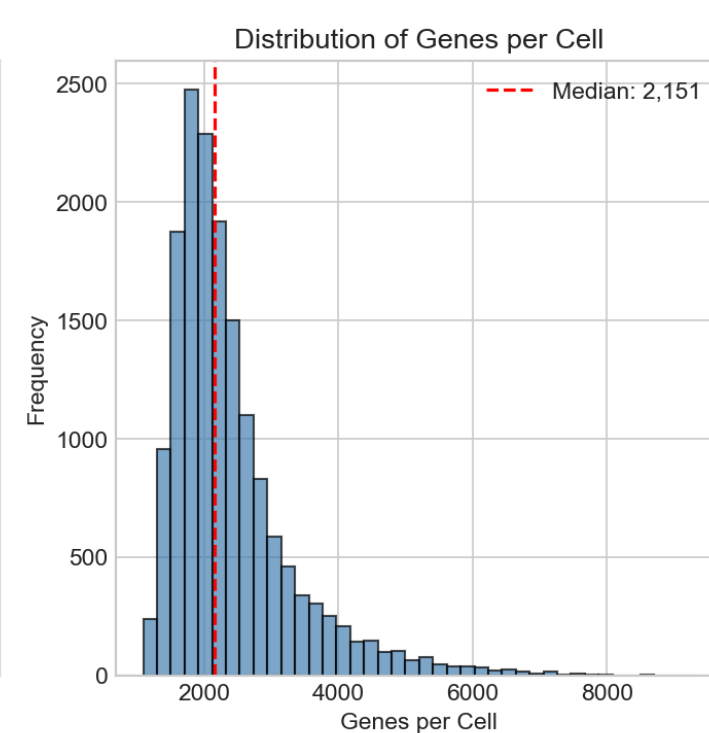


ROC Curve - 11-Gene Signature Model

XGBoost Model (AUC = 0.909)
Random Classifier (AUC = 0.500)
Optimal Threshold = 0.335



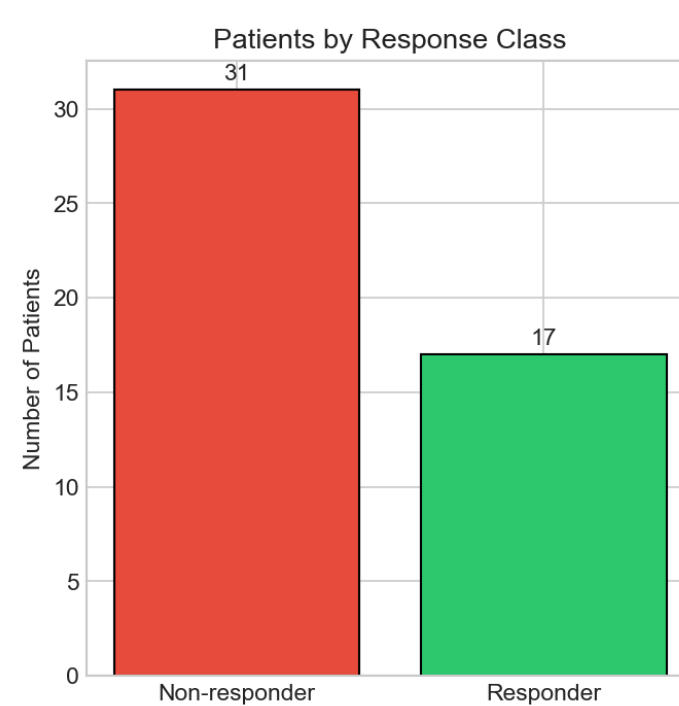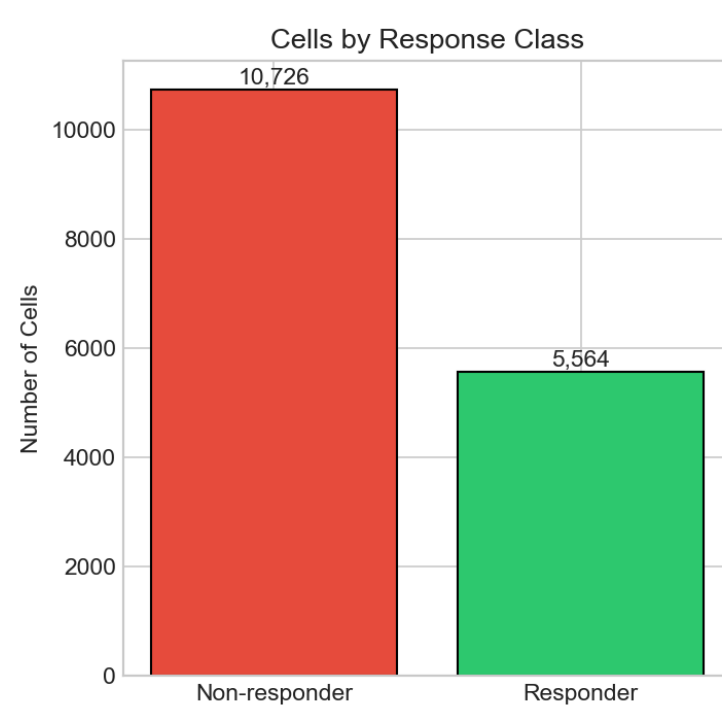Top Feature Importances

# Dataset Details, Labels, & Preprocessing

## DATASET

- 16,290 cells* across 48 unique tumor samples from melanoma patients
- Cells were labeled as Responder or Non-responder

- ICI therapy types were also present in the dataset (i.e. which checkpoint protein was targeted).



Patients by Therapy and Response; Response Rate by Therapy Type



GSE120575 Melanoma Dataset Overview

# XGBoost Model Training and Feature Selection

## BASELINE MODEL

The model was trained on individual cells but evaluated at the patient level using leave-one-patient-out cross-validation, where each patient is held out once as a test set. This full-gene model achieved a patient-level ROC AUC of 0.77, establishing that single-cell immune profiles contain meaningful signal about ICI response, but leaving room for improvement.

## FEATURE SELECTION

### Importance Ranking

XGBoost's gain-based feature importances were aggregated across cross-validation folds to rank all genes by how much they improved the model's splits. Many immune-related genes and members of the published PRECISE 11-gene panel rose to the top of this ranking.

### Top-50 Gene Model

For each training fold, we selected the top 50 most important genes and retrained XGBoost using only this reduced feature set.
This dimensionality reduction shrank the input space by about 250× while preserving the same training procedure and evaluation metric.

### Nested vs. Non-Nested

A nested leave-one-patient-out scheme was used so that gene selection was based only on training patients, avoiding train–test leakage.
A quick non-nested variant gave a much higher AUC but was intentionally treated as a biased estimate rather than a reproducible result.

## PERFORMANCE COMPARISON

- Baseline model reached AUC = 0.77

- After feature selection, the nested top-50-gene model achieved AUC approx. 0.78

- The 11-gene signature identified by original PRECISE paper produced a 0.91 AUC score



Patient Predictions - Feature Selection Model