# Galaxy Distance Estimation Using ML
## Supervised Learning Approach

Samuel Ghalayini

# Table of contents

01 **Introduction**

02 **Scope**

03 **Data**
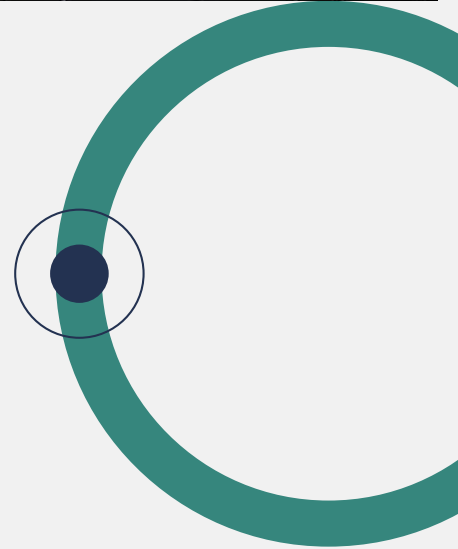
04 **Models**

# Introduction

- Motivation & Background
  - Interest sparked by data-driven astrophysics
  - Predict galaxy distances from SDSS photometry
  - Extend the cosmic distance ladder

# Project Scope



- Spectroscopic redshift & Hubble's Law
  - Redshift z from spectral line shifts
  - Doppler Effect
  - Velocity $v \approx c \times z$
  - Distance $D = v / H_0$
  - Train ML models to predict z

## Data Source

## Data Cleaning

- Sloan Digital Sky Survey 18 via astroquery.sdss
  - Up to 500k rows per redshift bin
  - Stacked into ~2M galaxy entries
  - ugriz magnitudes, uncertainties, redshift

- Key steps
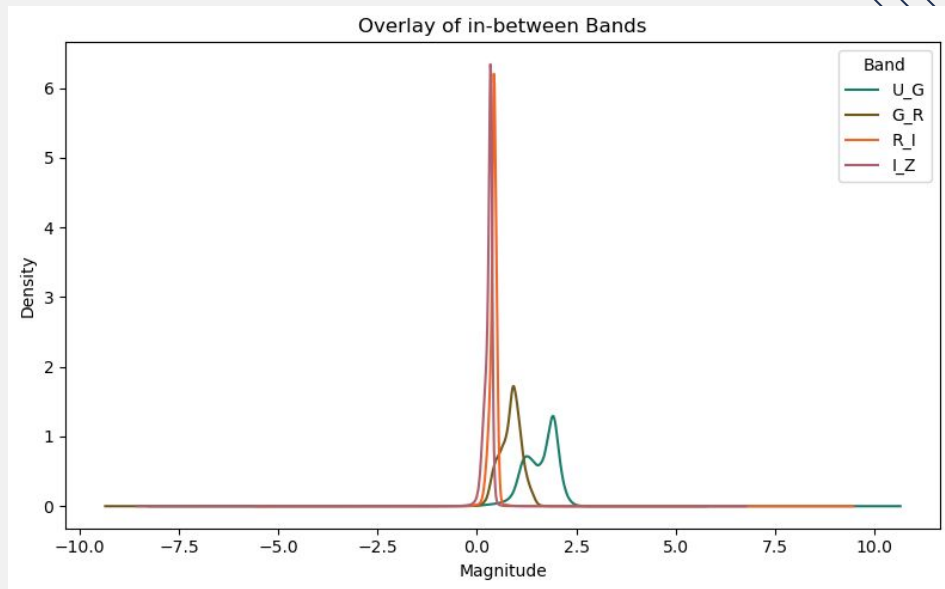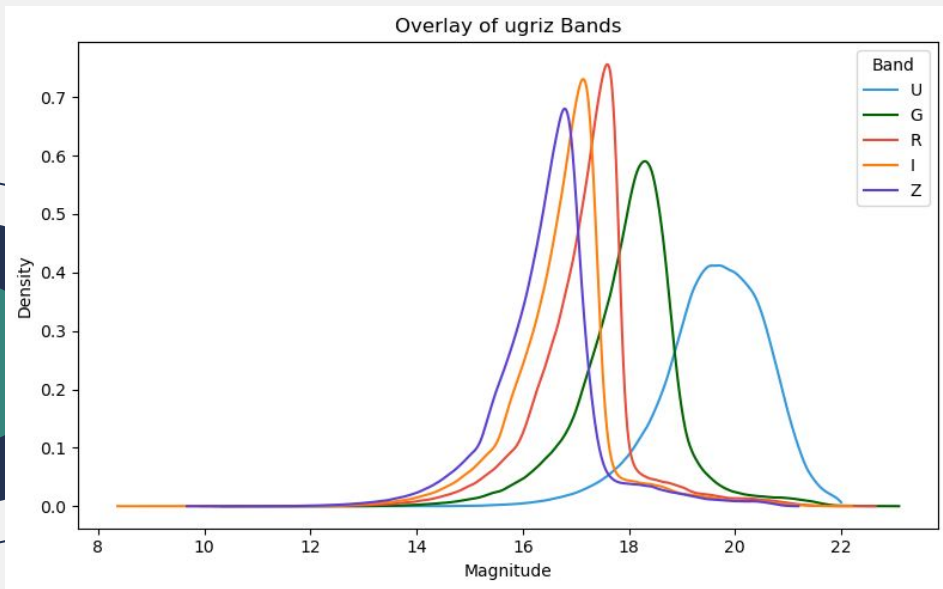  - Replace -9999 with NaN and drop
  - Exclude $z < 0.01$ or $z > 1.0$
  - Require err < 0.2 mag (S/N > 5)
  - Remove fainter than 5σ depth
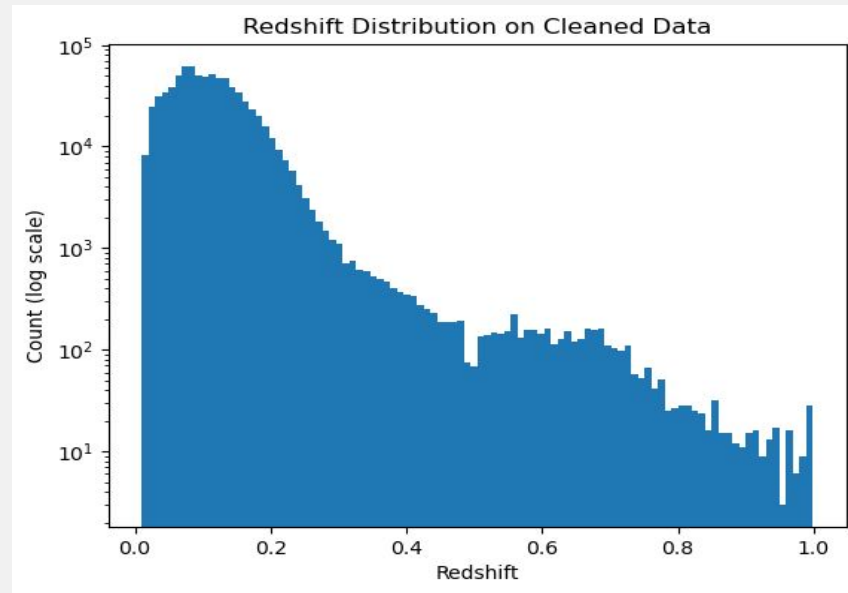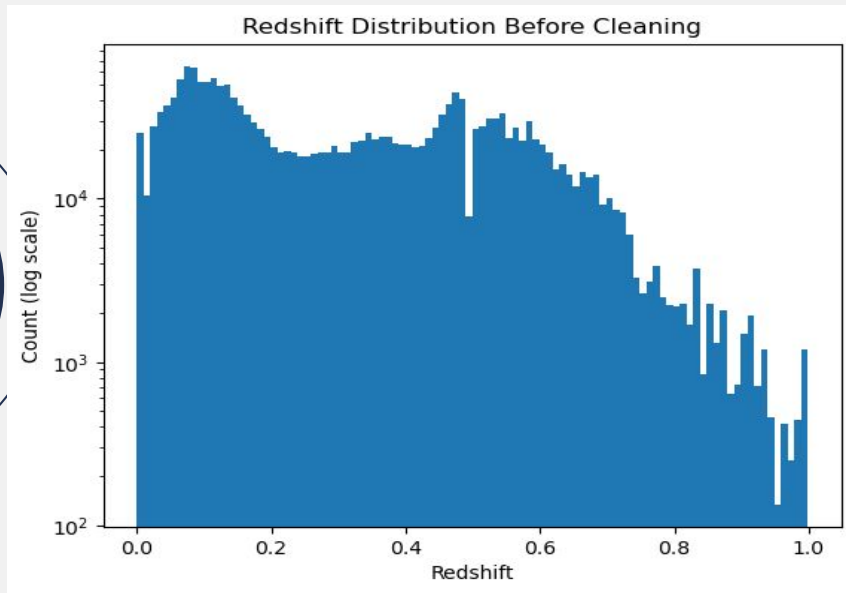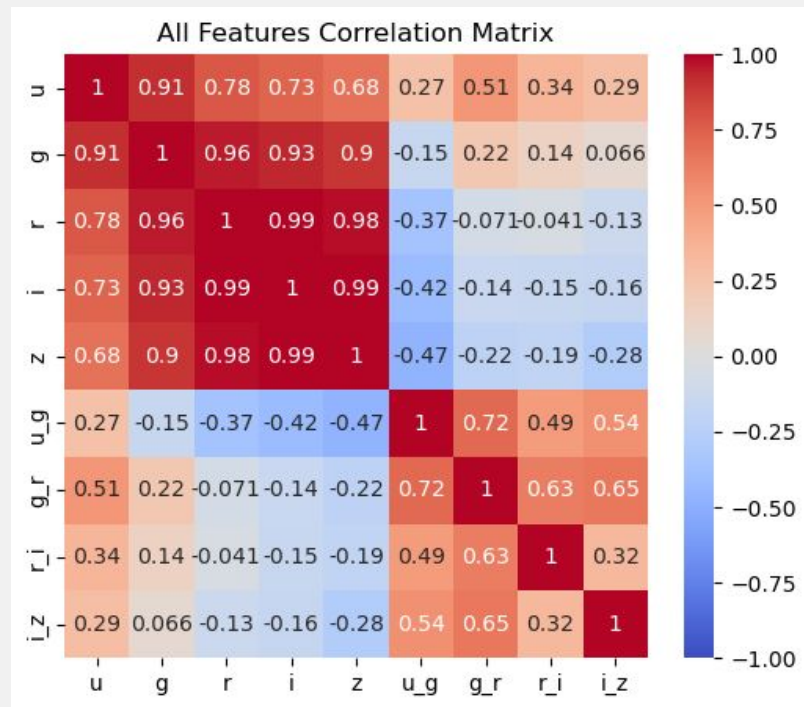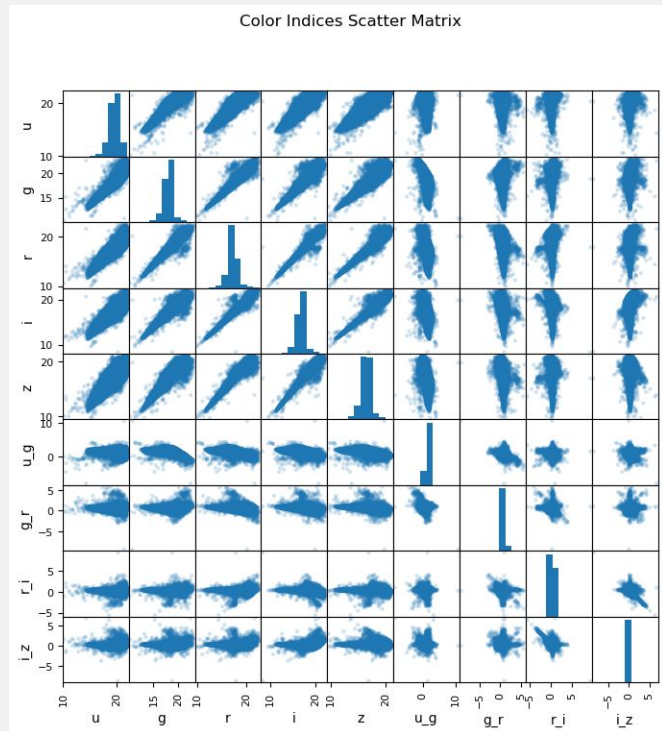  - Create u−g, g−r, r−i, i−z indices

# Exploratory Data Analysis

- Highlights
  - ugriz KDEs show spectral sequence
  - Pairplot: strong band correlations
  - Color–color clouds reveal variance
  - Cleaned redshifts

## Overlay of ugriz Bands

## Overlay of in-between Bands

Color Indices Scatter Matrix
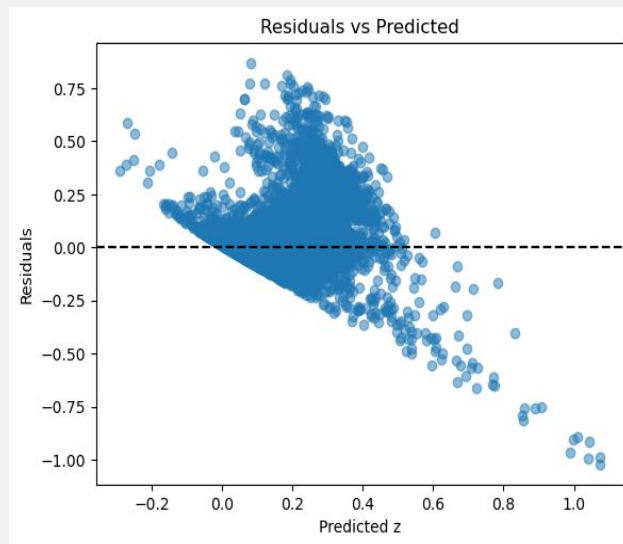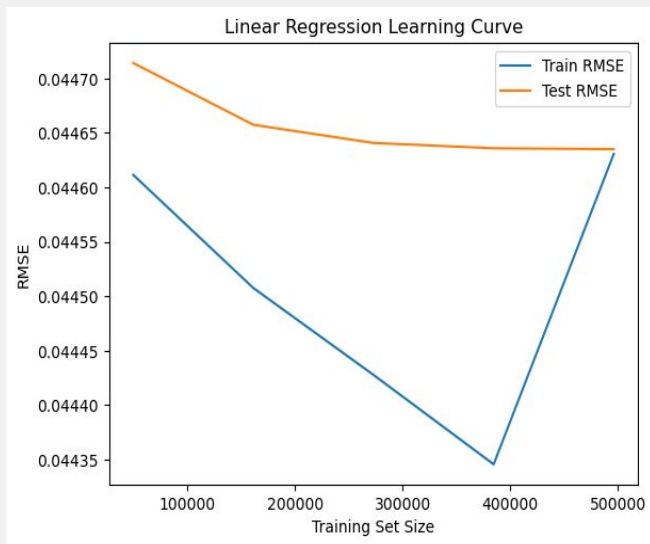

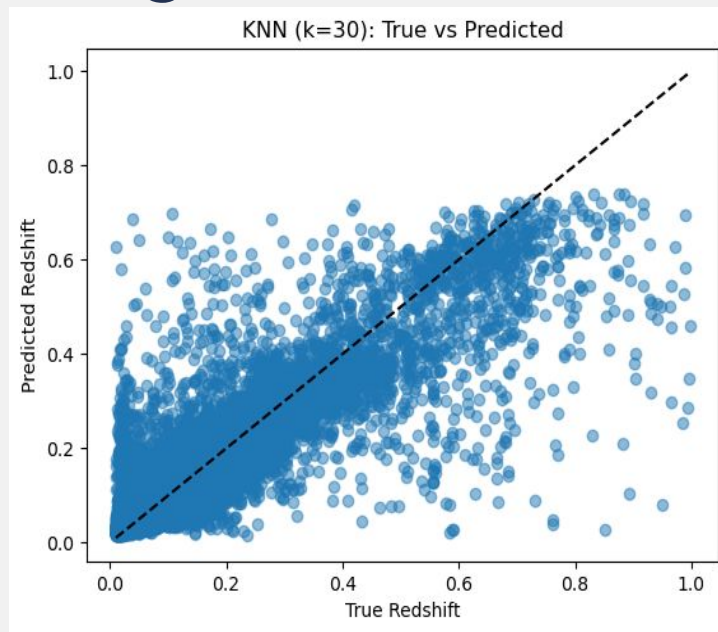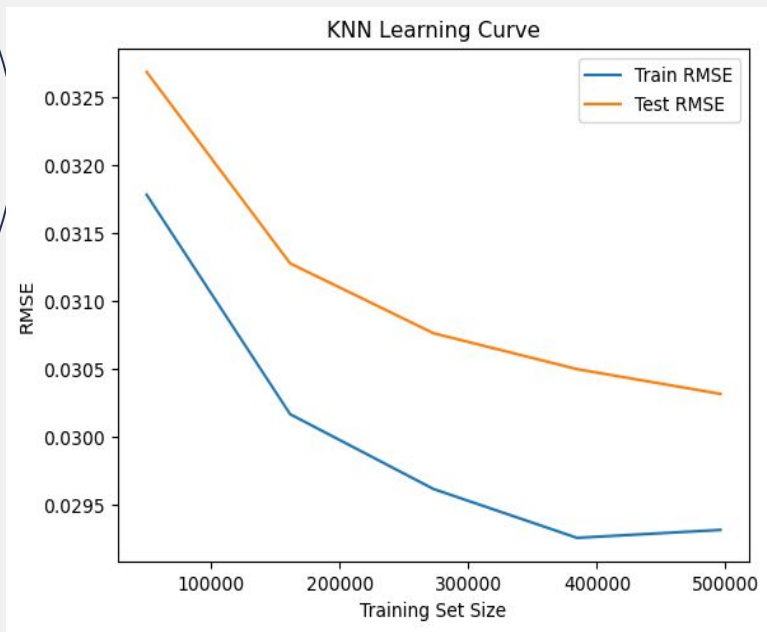All Features Correlation Matrix

# Models

- Algorithms
  - Linear Regression (feature selection)
  - KNN (k=30)
  - Random Forest (RandomizedSearchCV)
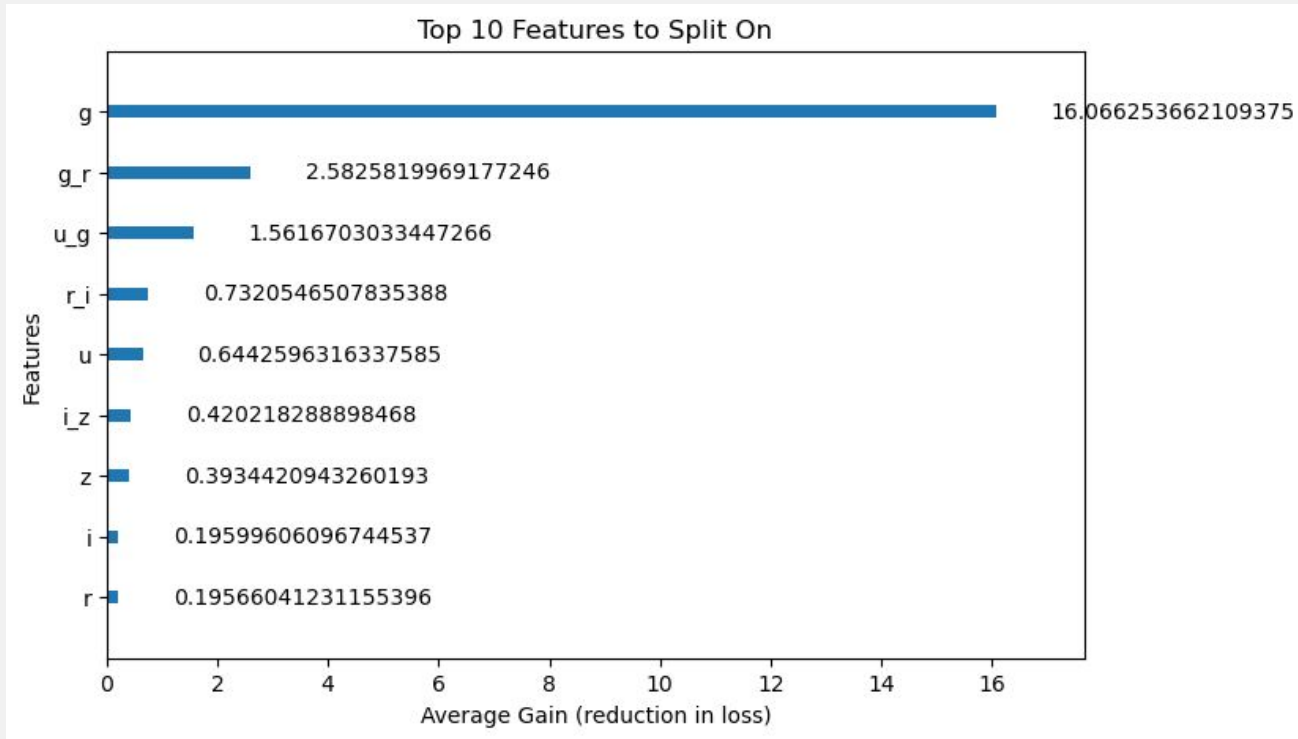  - XGBoost (Boosting)

# Linear Regression

# K-Nearest Neighbors

# Random Forest and Notebook

# XGBoost



Top 10 Features to Split On

# Results & Performance

| Model | MAE | RMSE | R² | Time (s) |
|---|---|---|---|---|
| Linear Regression | 0.0263 | 0.0444 | 0.6079 | 0.0283 |
| KNN (k=24) | 0.0180 | 0.0311 | 0.8076 | 5.1731 |
| Random Forest | 0.0374 | 0.0522 | 0.4565 | 353.6183 |
| XGBoost | 0.0176 | 0.0307 | 0.8126 | 1.9744 |

# Discussion & Next Steps

- Takeaways
  - Data cleaning drives performance
  - XGBoost best balance of speed & accuracy
  - Future: more bands, cloud hyperparameter tuning
- Github: https://github.com/sam-ghala/galaxy_distance