

# Pharmaceutics Meets Machine Learning: Predicting Drug Price by Simplified Approaches



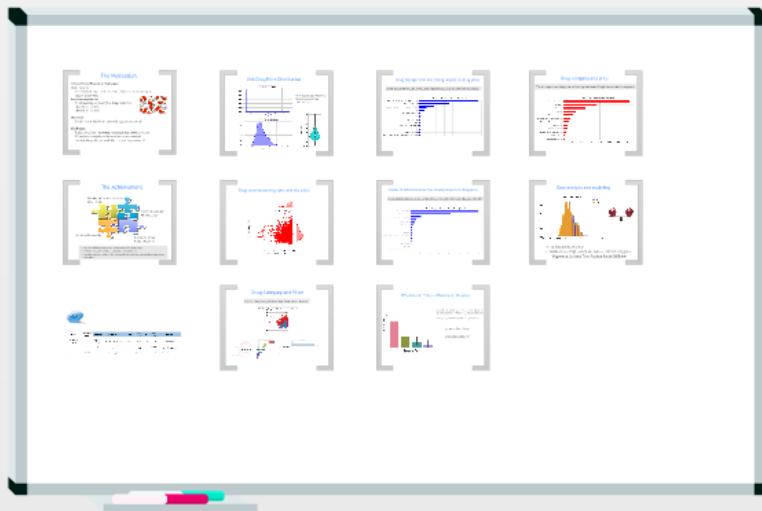
Bo Lian, Choutine Zhou, Robin Gane-McCalla

# Pharmaceutics Meets Machine Learning: Predicting Drug Price by Simplified Approaches



Bo Lian, Choutine Zhou, Robin Gane-McCalla





# The Motivation

## **Life and Social Economic Challenges:**

### High drug price

Prescription drug prices in the United States have been among the highest in the world.

### Enormous expenditure

Total spending on prescription drugs in the U.S

\$425 billion in 2015

\$600 billion in 2021



## **Our Goal:**

Develop a simplified and universal drug prediction model

## **Challenges:**

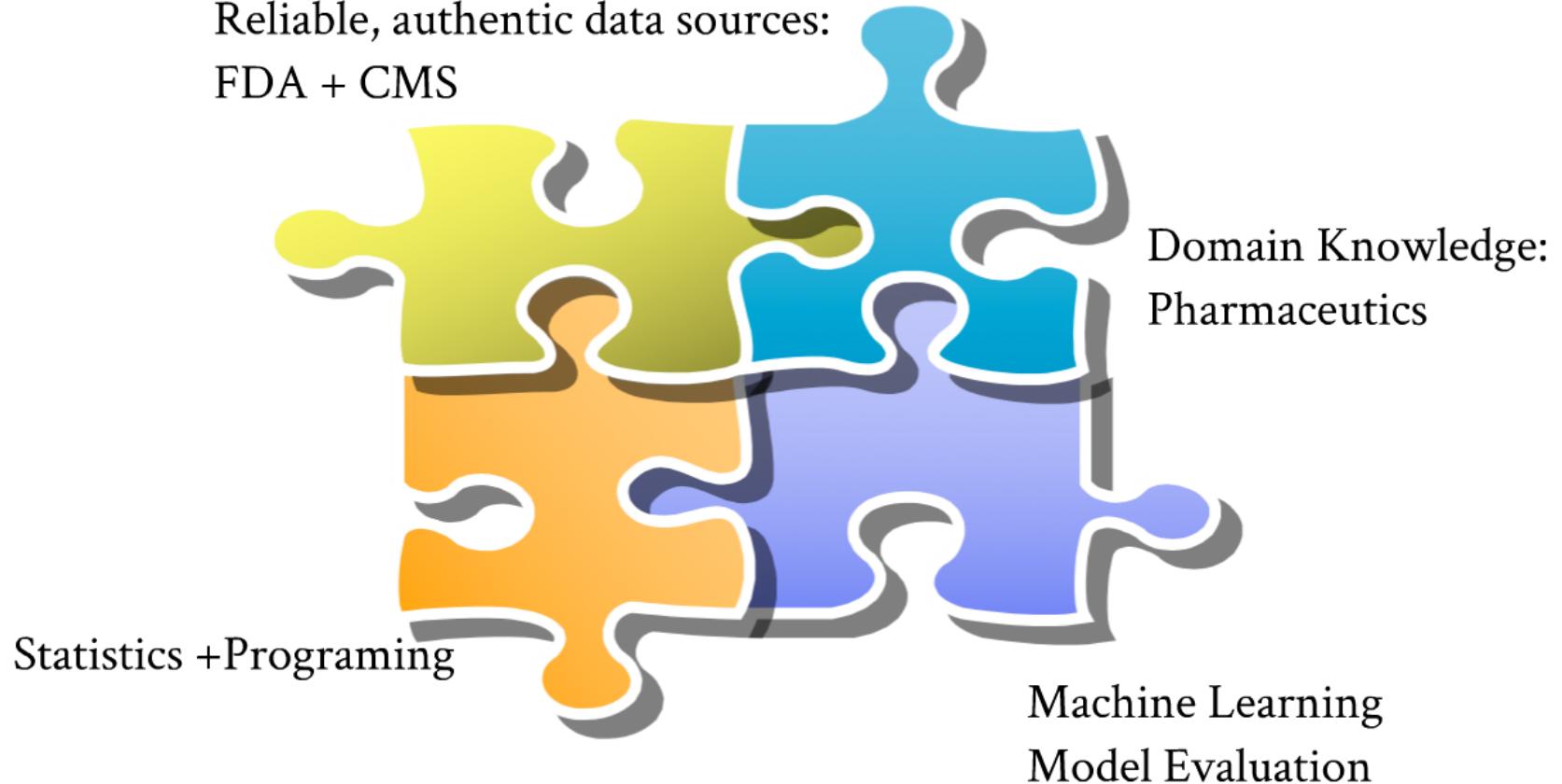
R&D, clinical trial, marketing, manufacturing, distribution cost

Policy/politics/regulation/business/economics involved

No desirable public data available, they cost huge amount!!!

# The Achievement

Reliable, authentic data sources:  
FDA + CMS

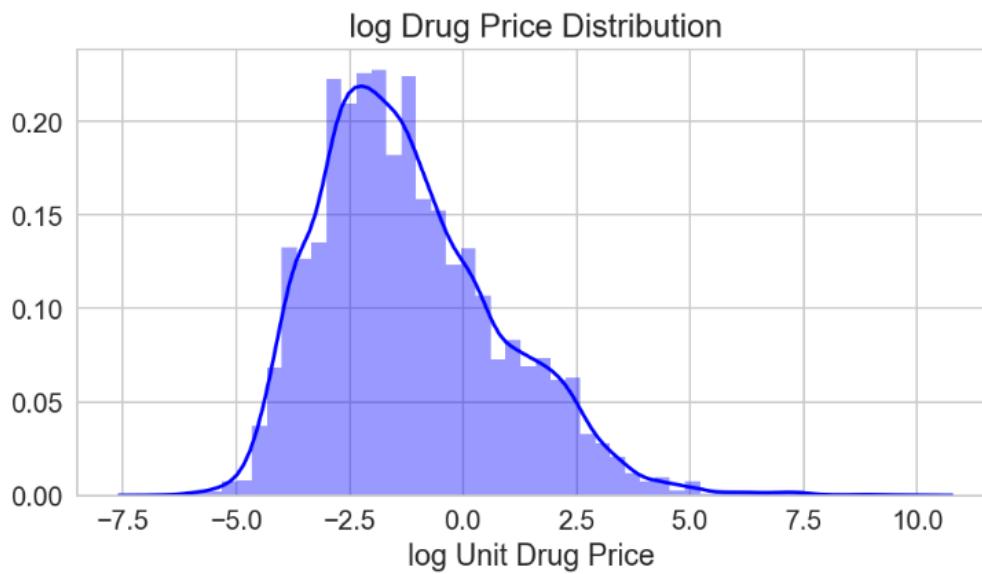
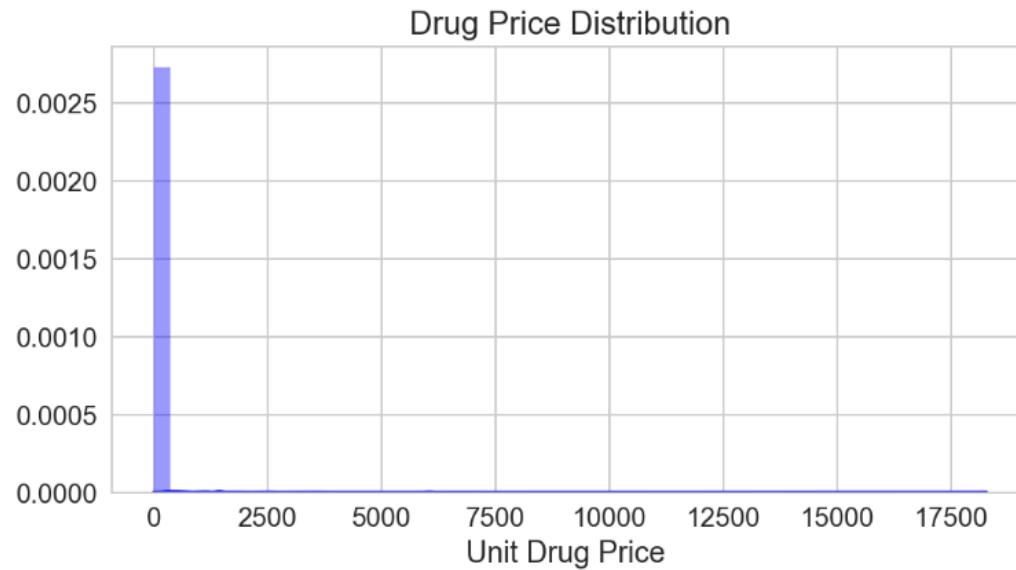


- Fully utilized limited data info and imbalance data distribution
- 20+ Features scaled down to 7 , mostly are drug features
- 4 models built and evaluated for price prediction and price classification, respectively
- Innovation

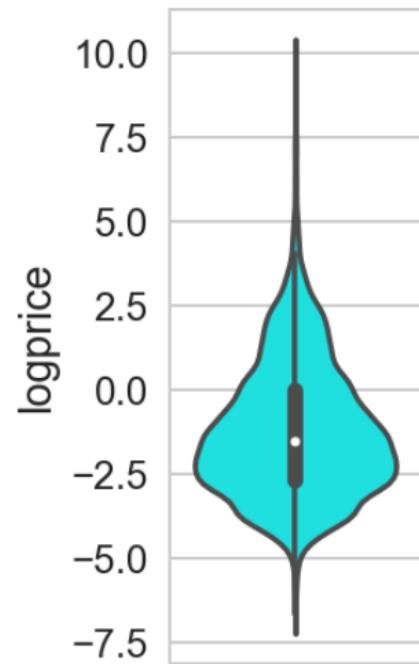


Product	Unit Price (dollar)	Company	Drug Category	Route	Dosage Form	DEA	Starting Date	Indication
Good neighbor pharmacy day time	0.13	Amerisource Bergen	Over-the-counter (OTC)	oral	capsule	N/A	8/2/12	cold and flu
Viagra	59.1	Pfizer	Brand chemical drugs (NDA)	oral	film coated tablet	N/A	3/27/98	Pulmonary hypertension Sexual dysfunction
Neulasta	9070.6	Amgen	Biologics (BLA)	subcutaneous	Injection	N/A	4/1/02	Cancer chemo-therapy

# Unit Drug Price Distribution



Right Skewed, Super Imbalanced  
Most are low price drugs  
Log transformed



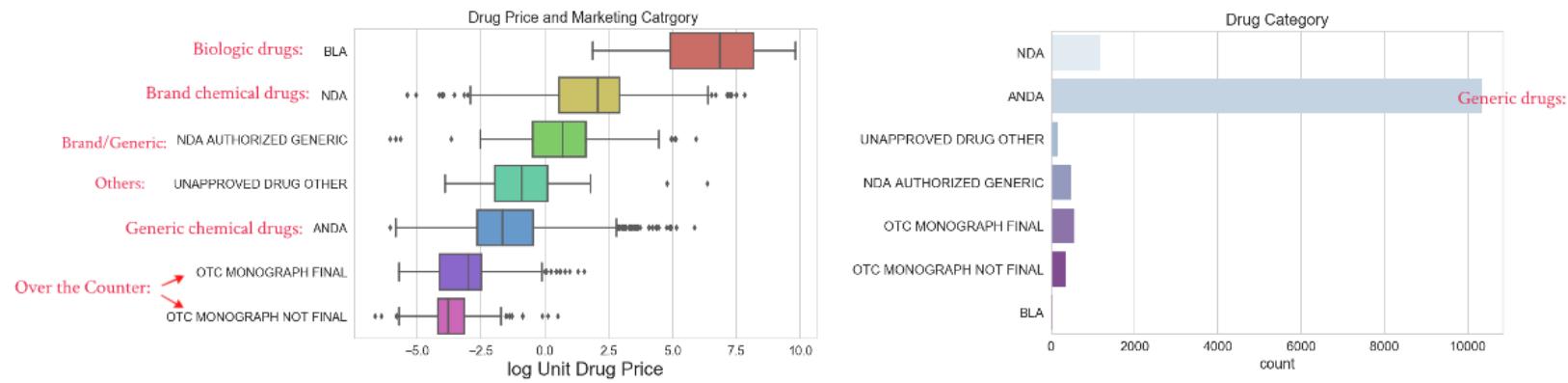
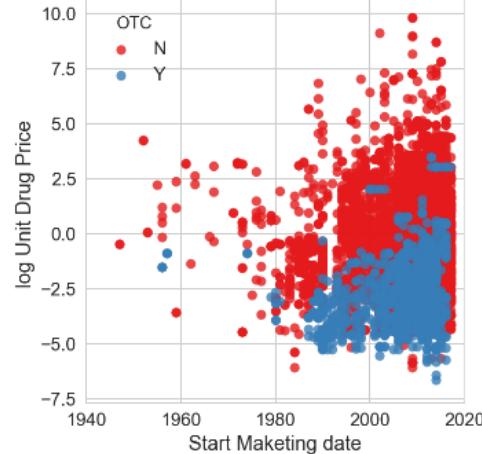
# Drug start-marketing-date and the price



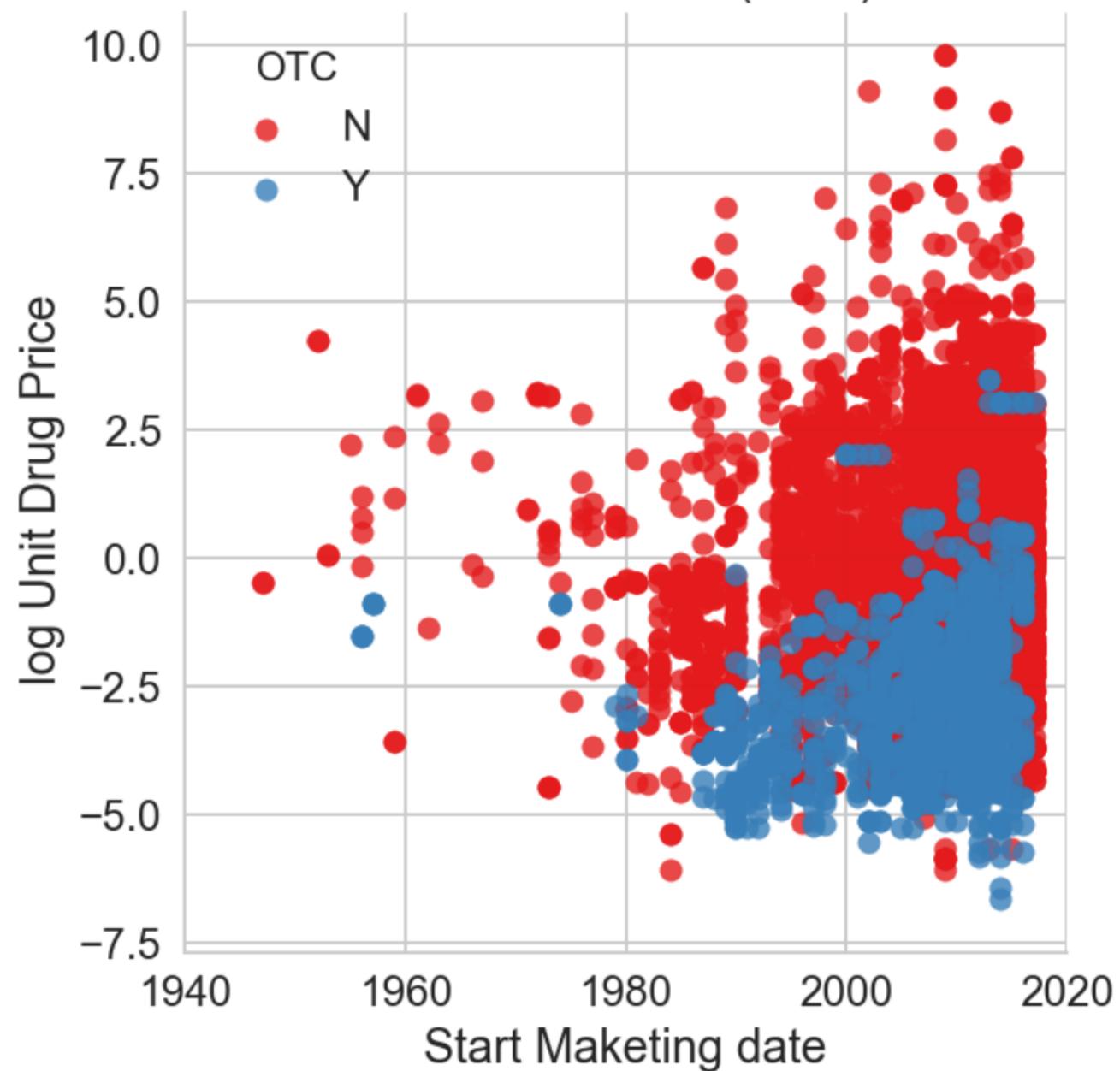
# Drug Category and Price

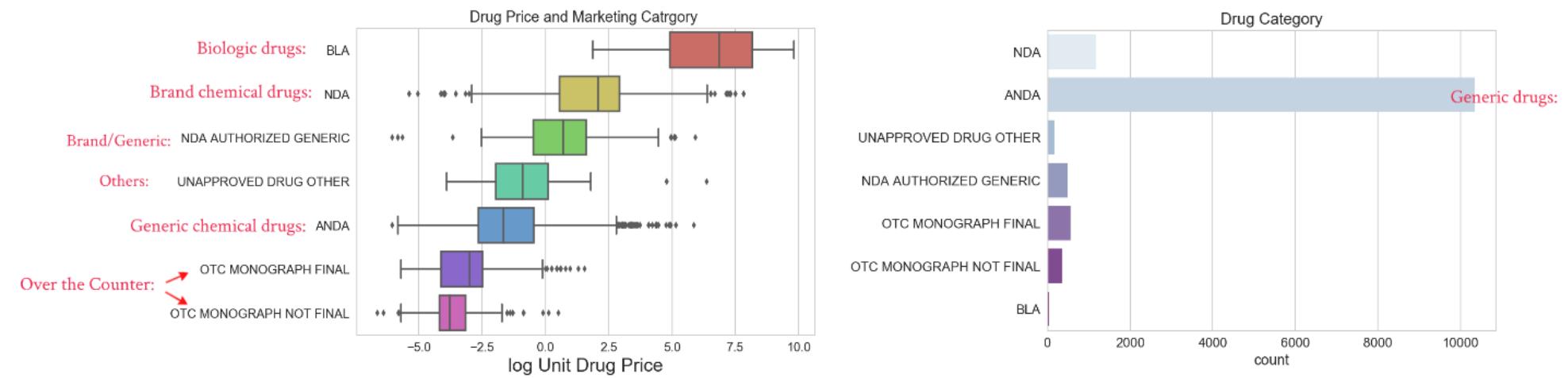
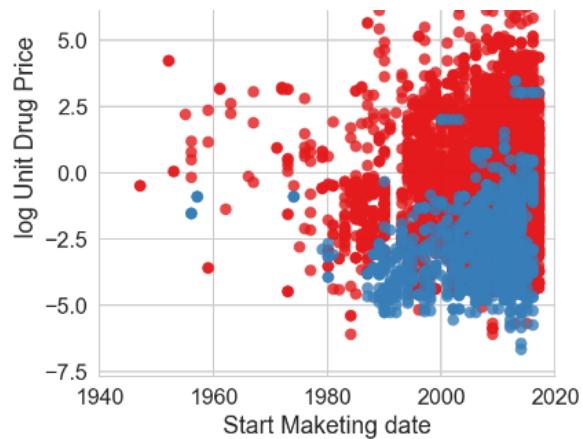
Different drug categories have significant price difference

Price Difference between Over-the-Counter (OTC) and Prescription (Rx) Drugs



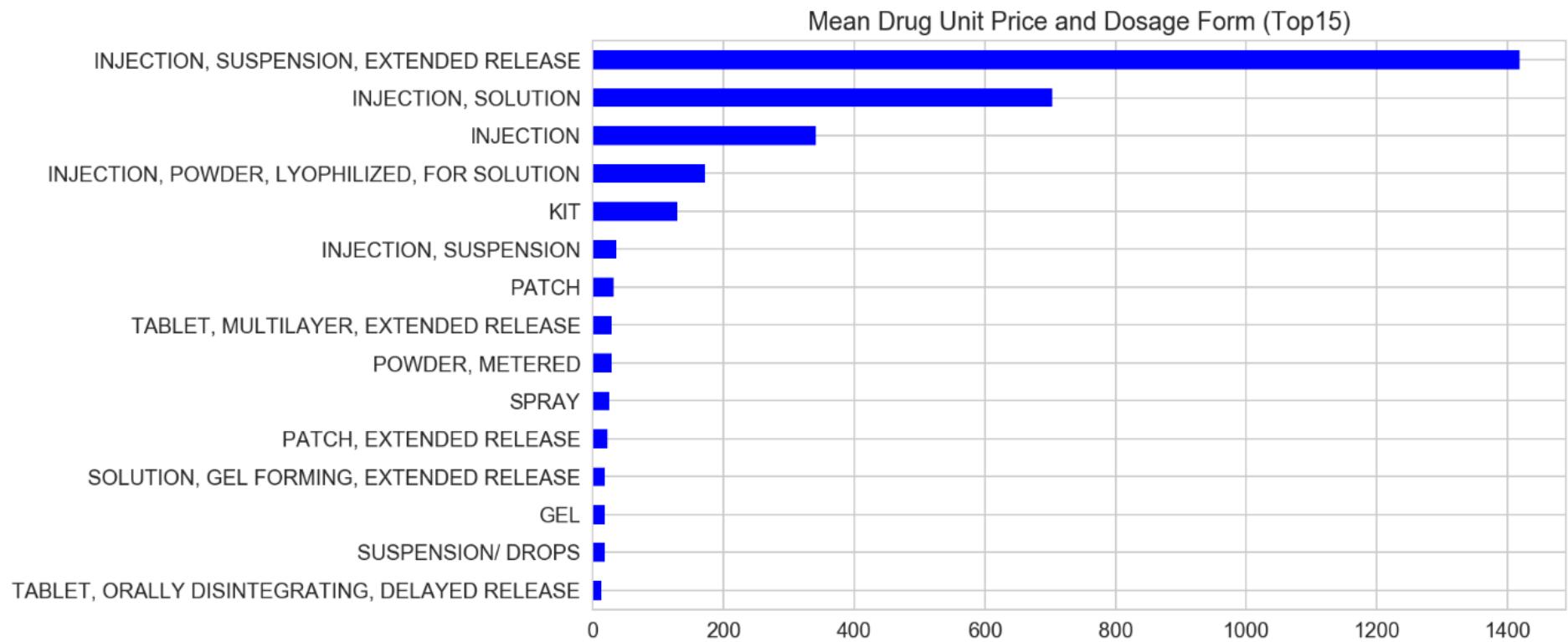
## Price Difference between Over-the-Counter (OTC) and Prescription (Rx) Drugs





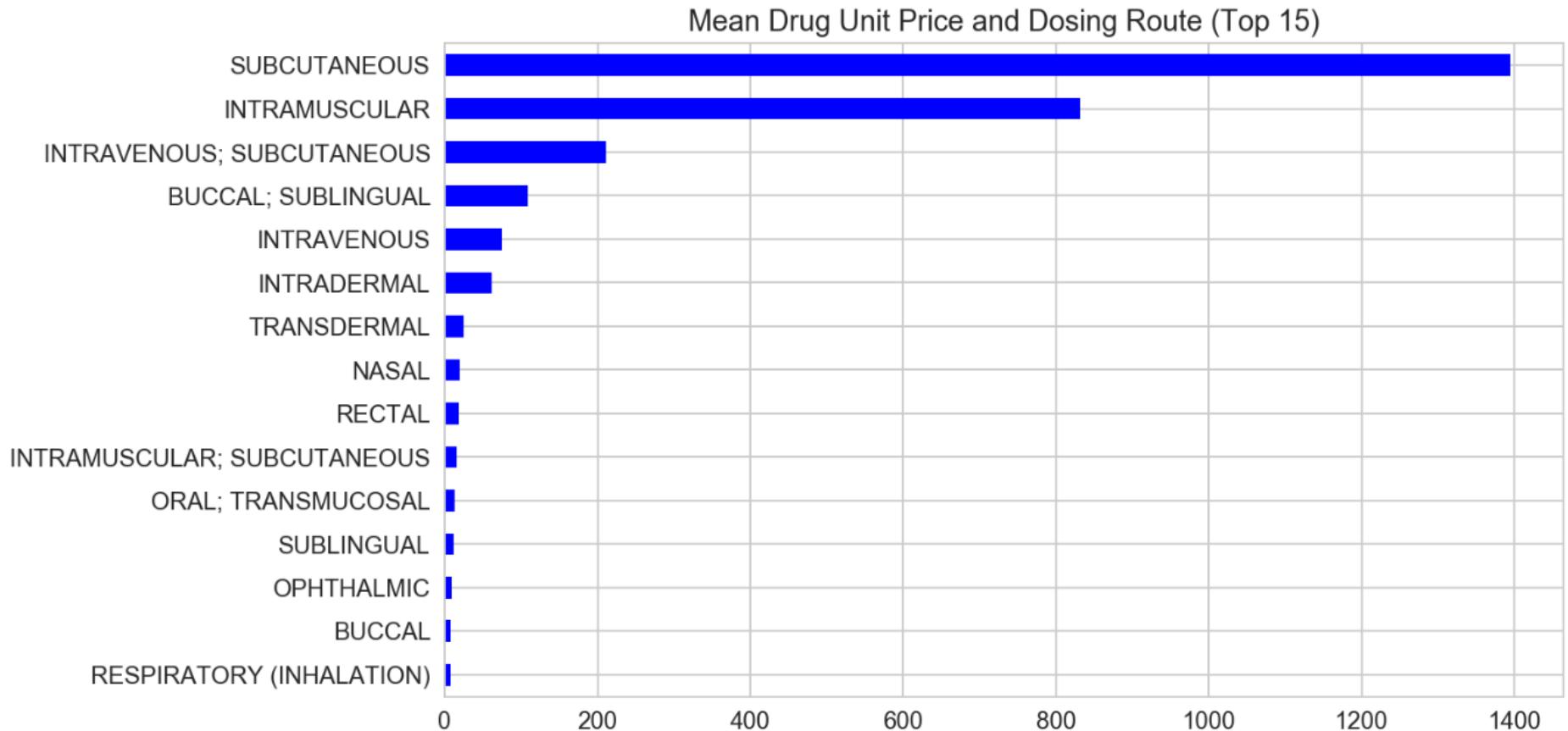
# Drug Dosage form has strong impact on drug price

Liquid dosageforms usually are much more expensive than solid and semi-solid dosageform

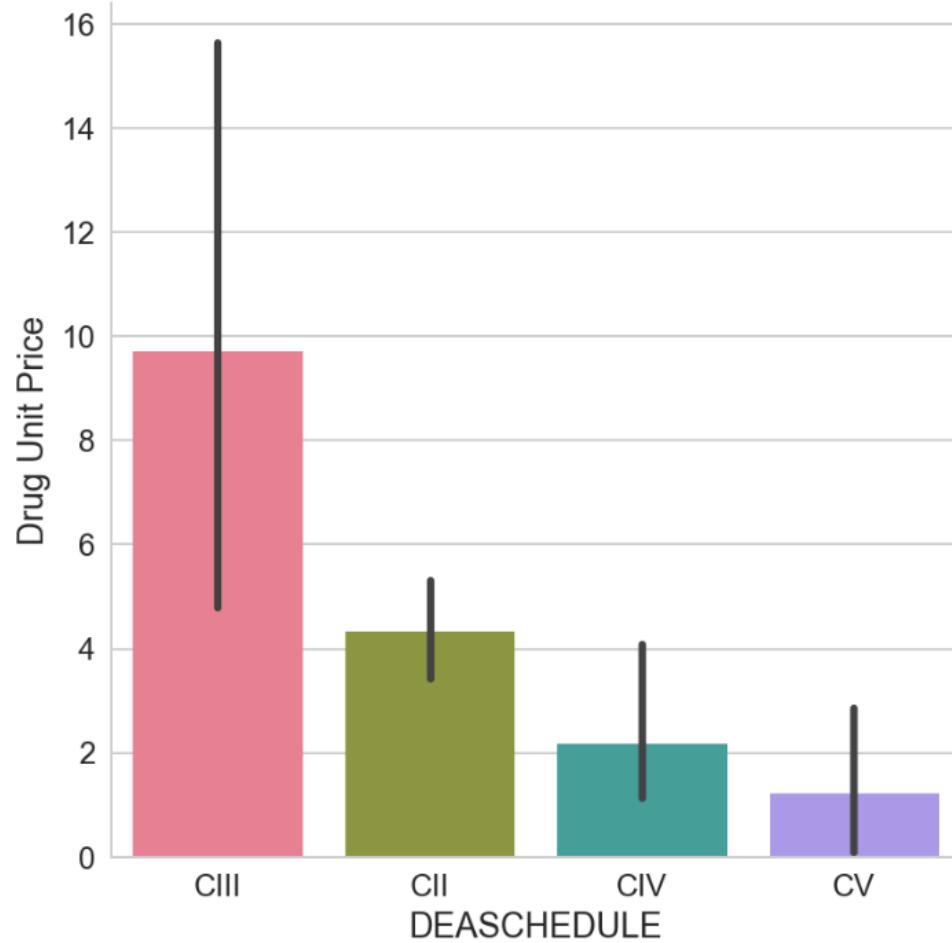


# Routes of Administration has strong impact on drug price

Drugs administrated by injection are mostly more expensive than drugs administrated orally



## DEA class can make a difference on the price



Schedule I : drugs with no currently accepted medical use and a high potential for abuse: (not covered by Medicaid):

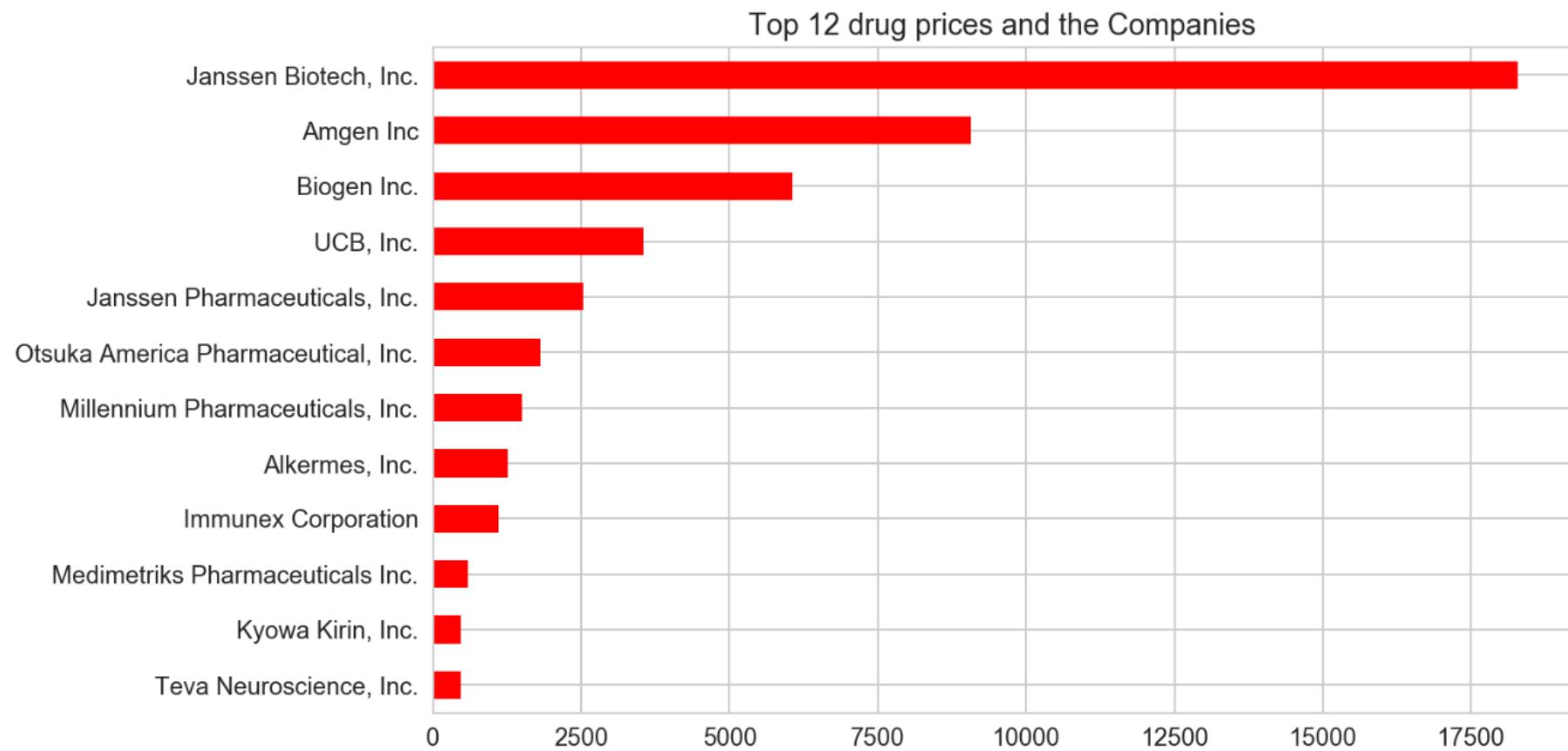
heroin, lysergic acid diethylamide (LSD), marijuana

potential for abuse:

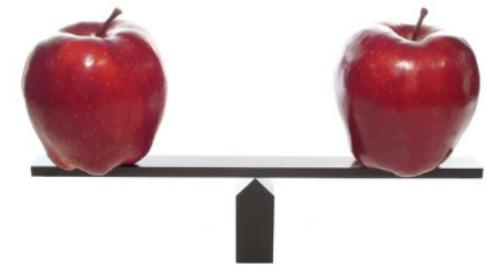
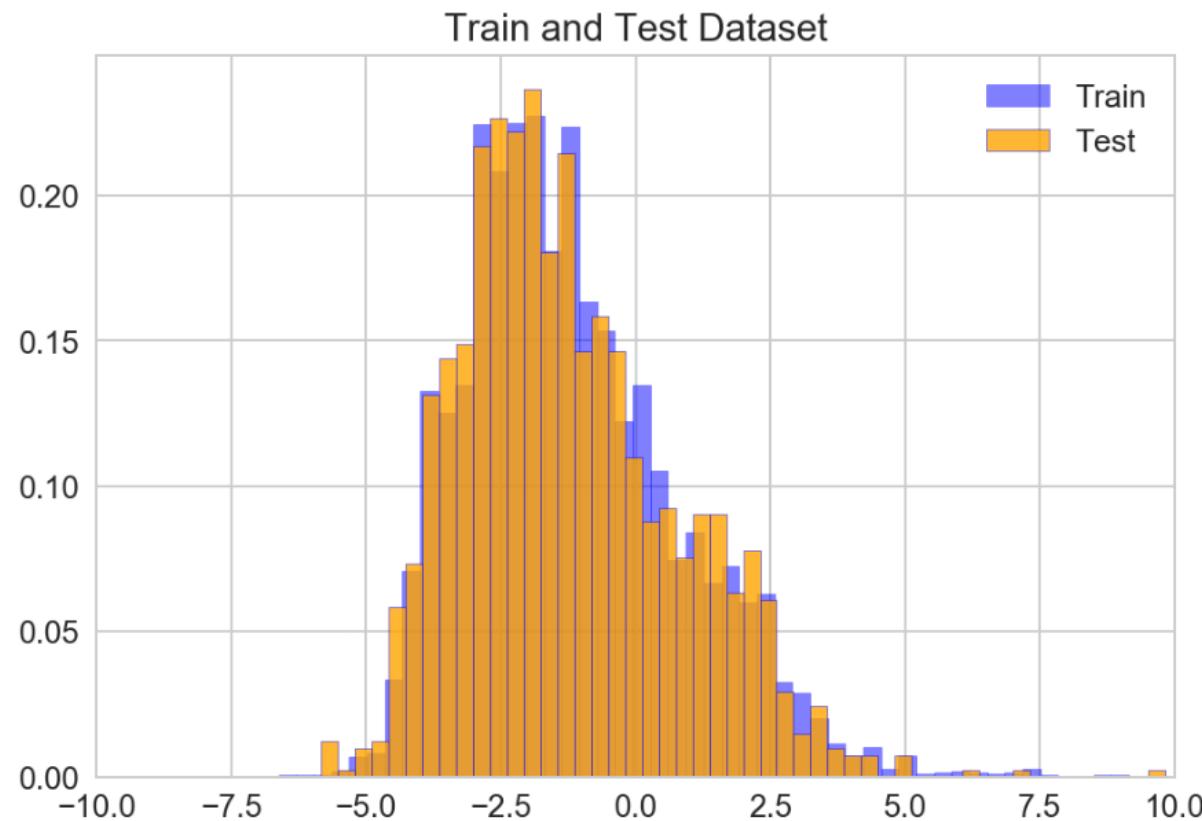
CV < CIV < CIII < CII

# Drug company and price

The top expensive drugs come from top innovative biopharmaceutical companies

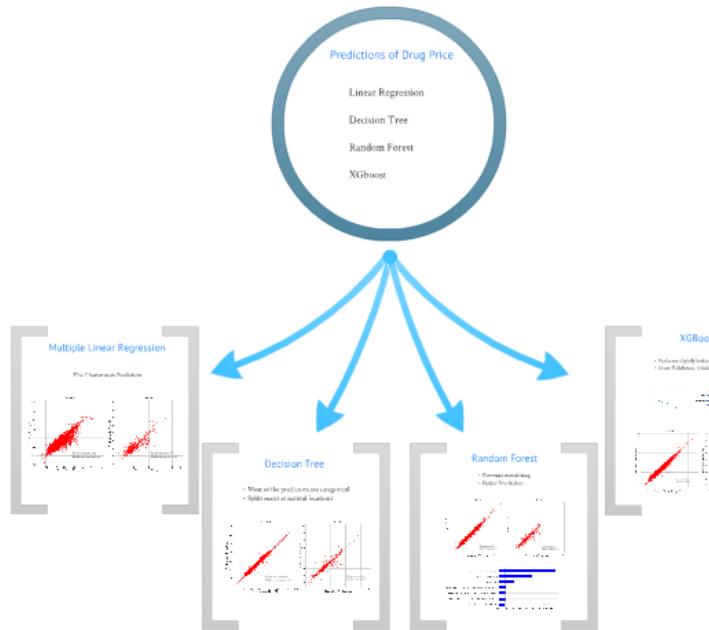


# Data analysis and modeling



- Unit drug price prediction
  - Classification of high ( $>10 \$$ ) and low ( $< 10\$$ ) unit drug price
- Regression, Decision Tree, Random Forest, XGBoost

## Price Prediction



## Price Classification



# Predictions of Drug Price

Linear Regression

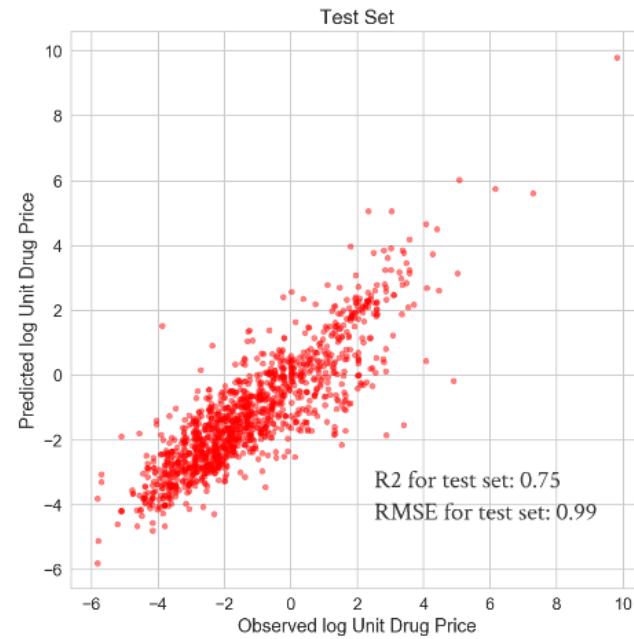
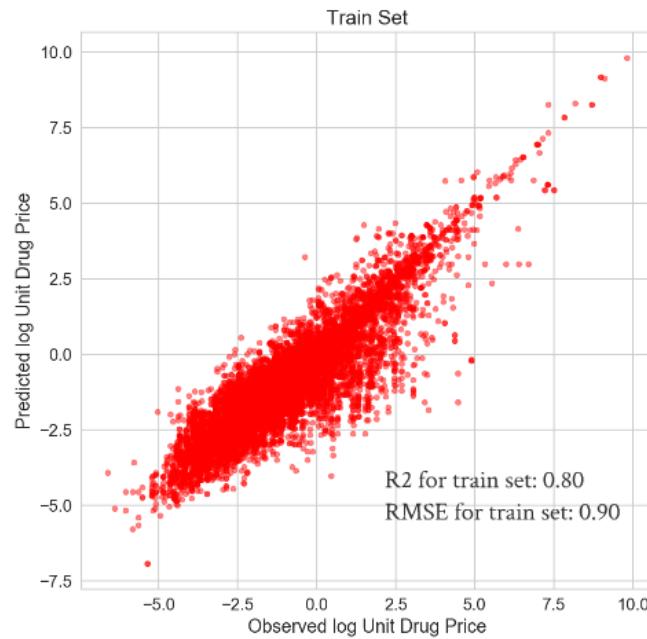
Decision Tree

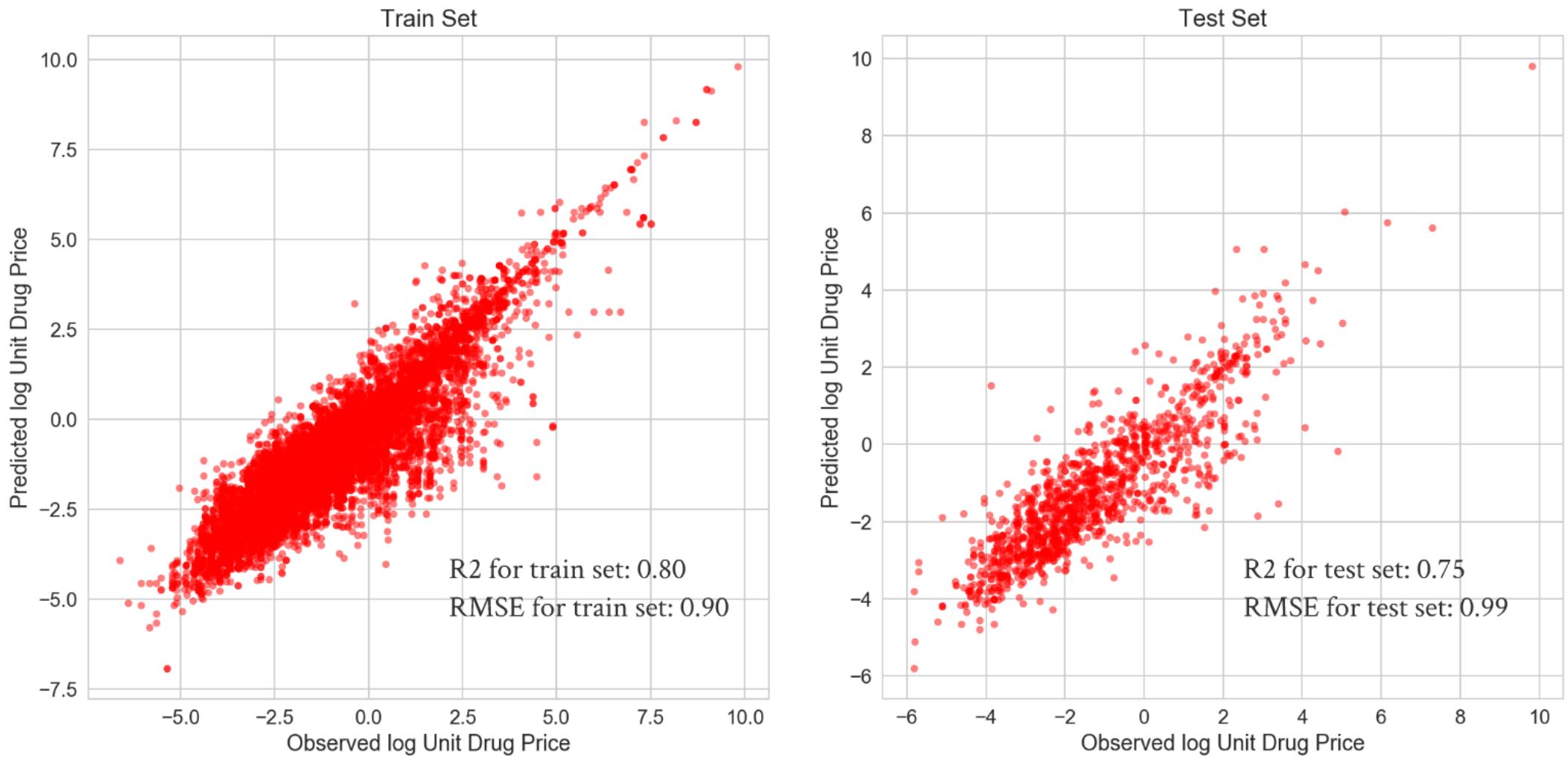
Random Forest

XGboost

# Multiple Linear Regression

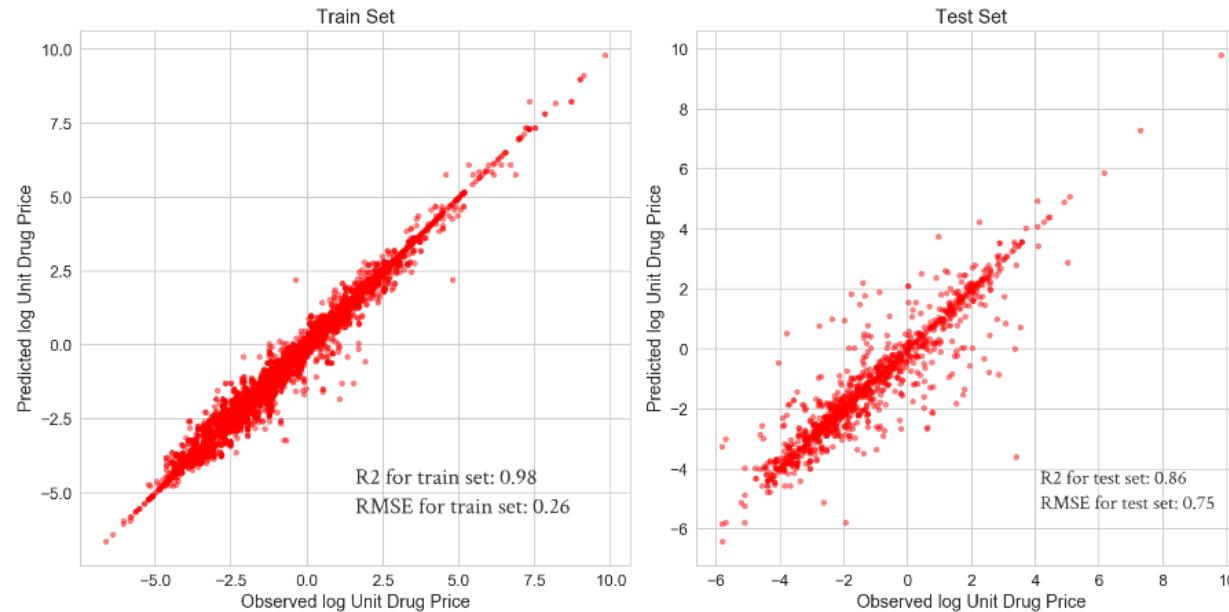
The 7 features as Predictors

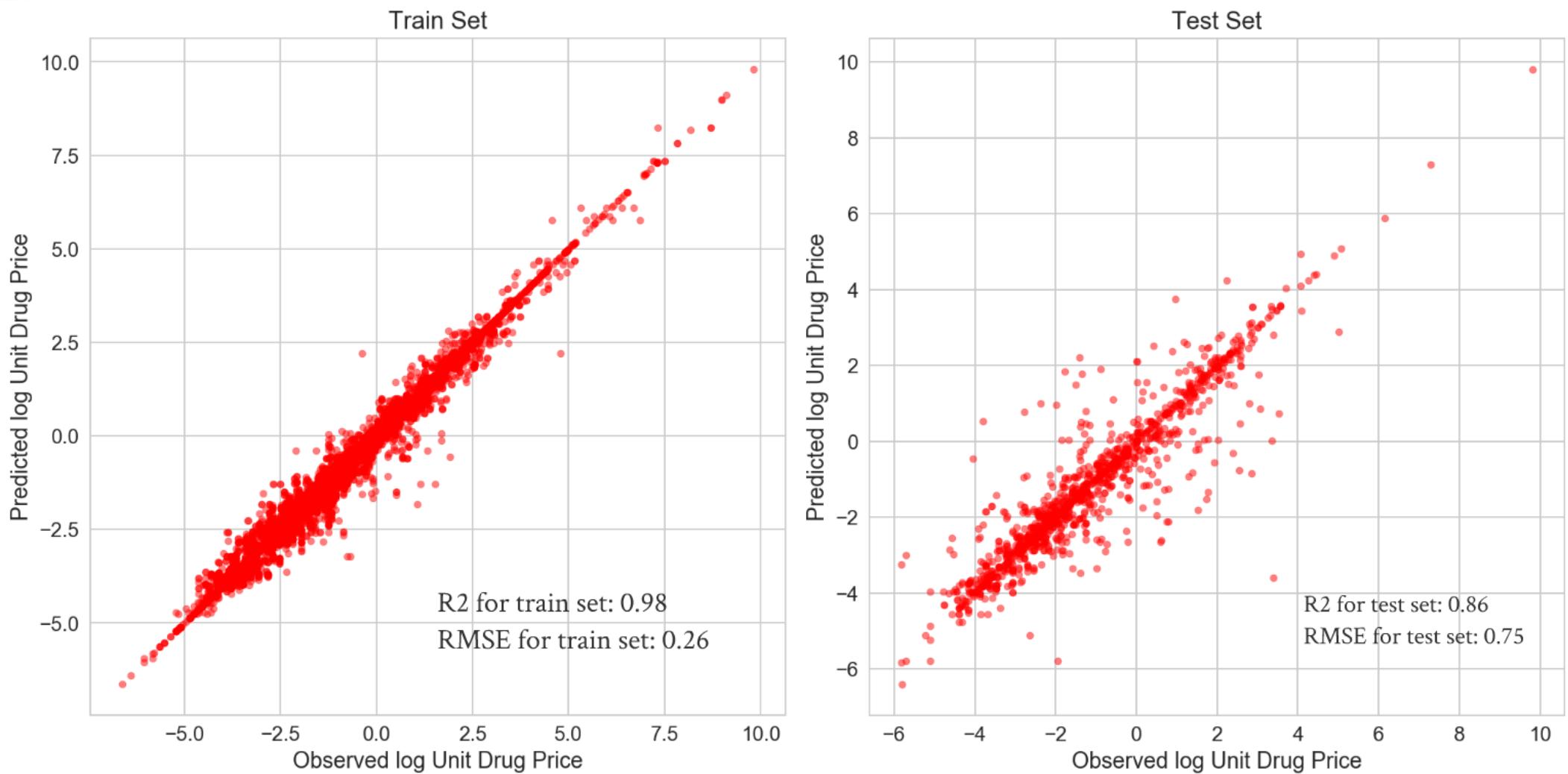




# Decision Tree

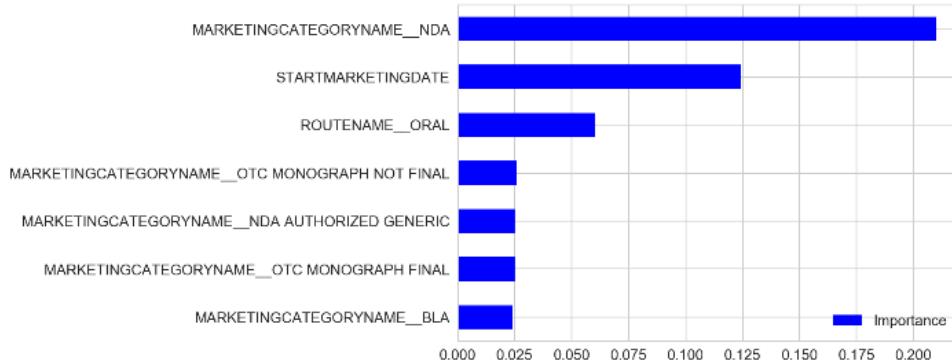
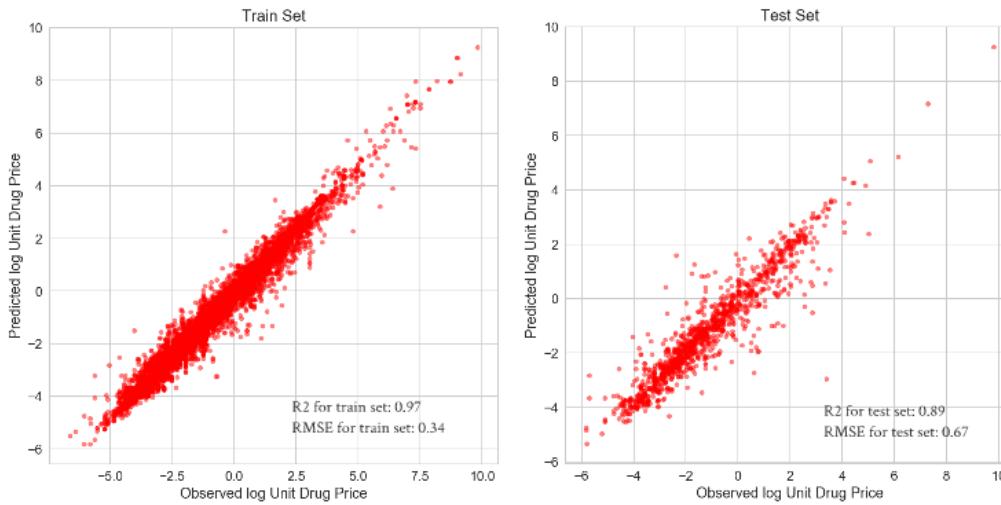
- Most of the predictors are categorical
- Splits occur at natural locations



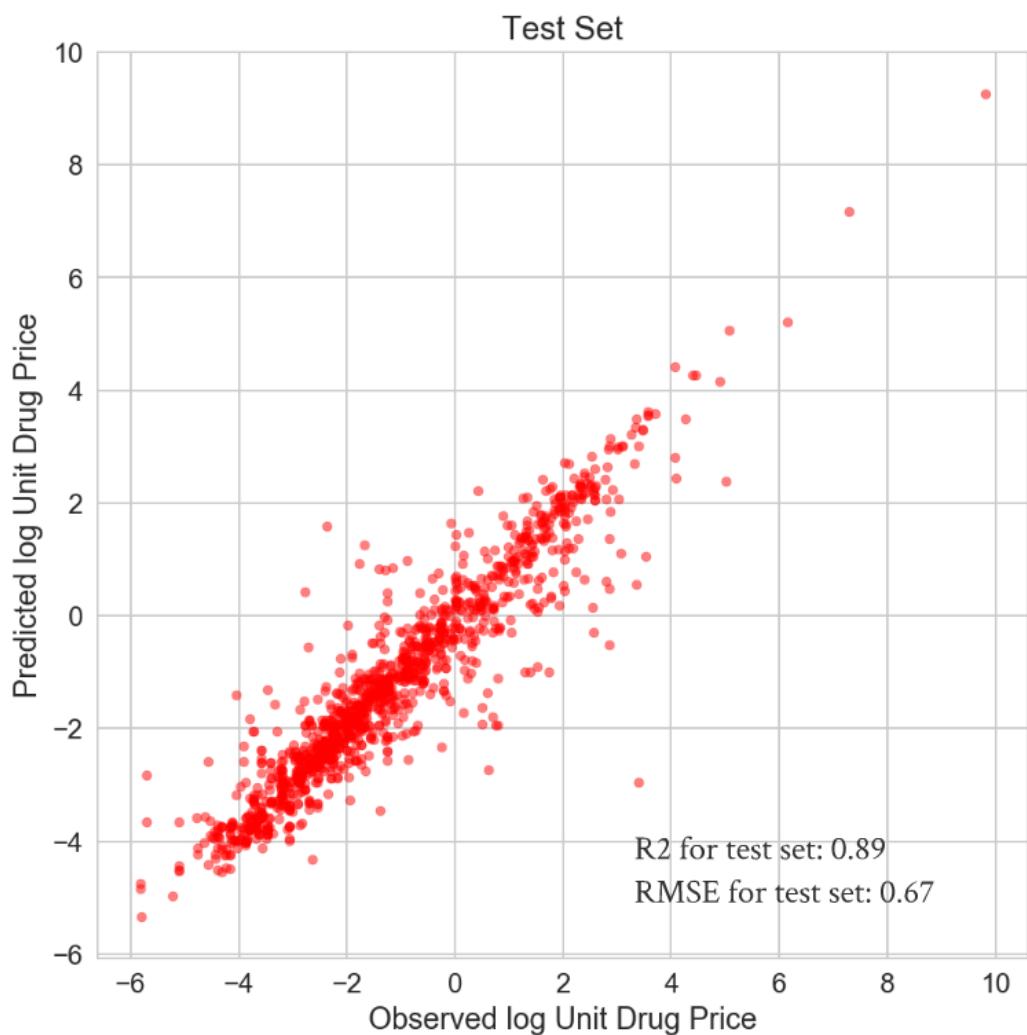
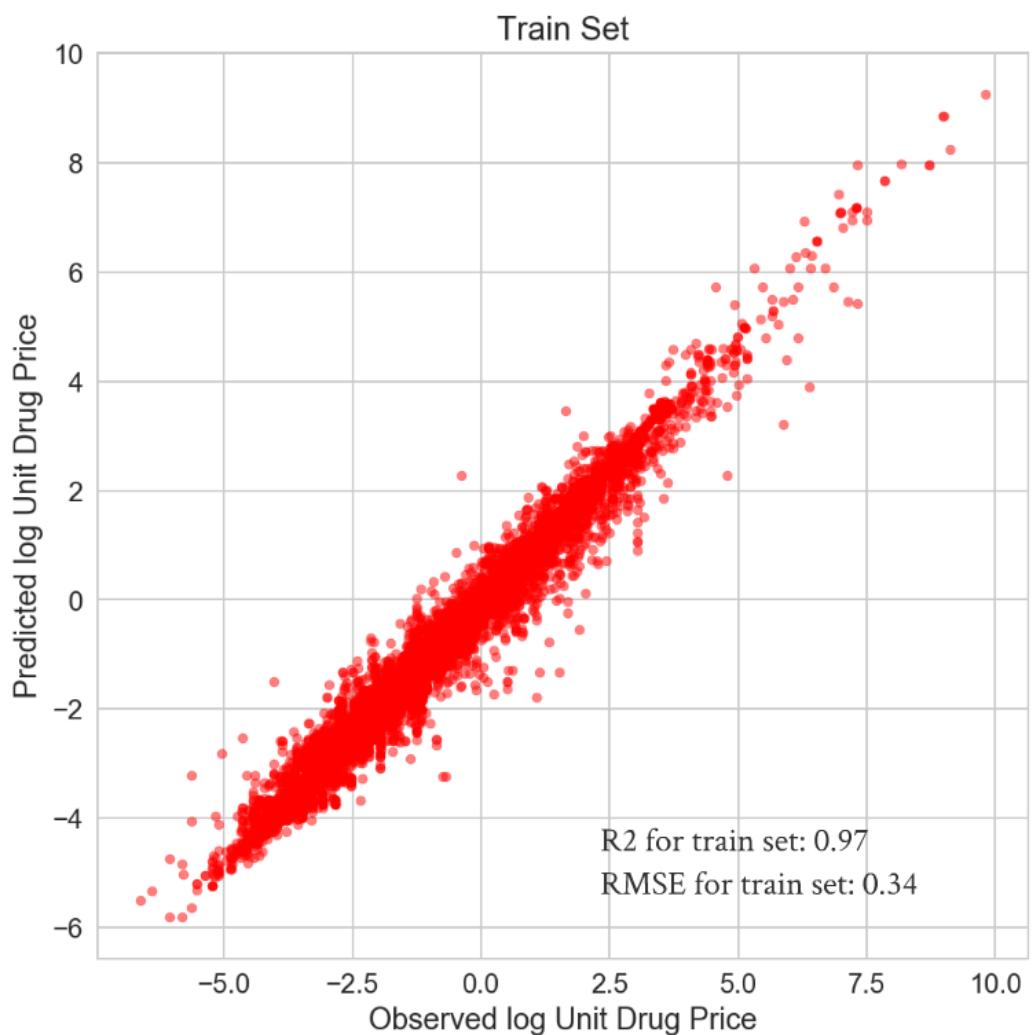


# Random Forest

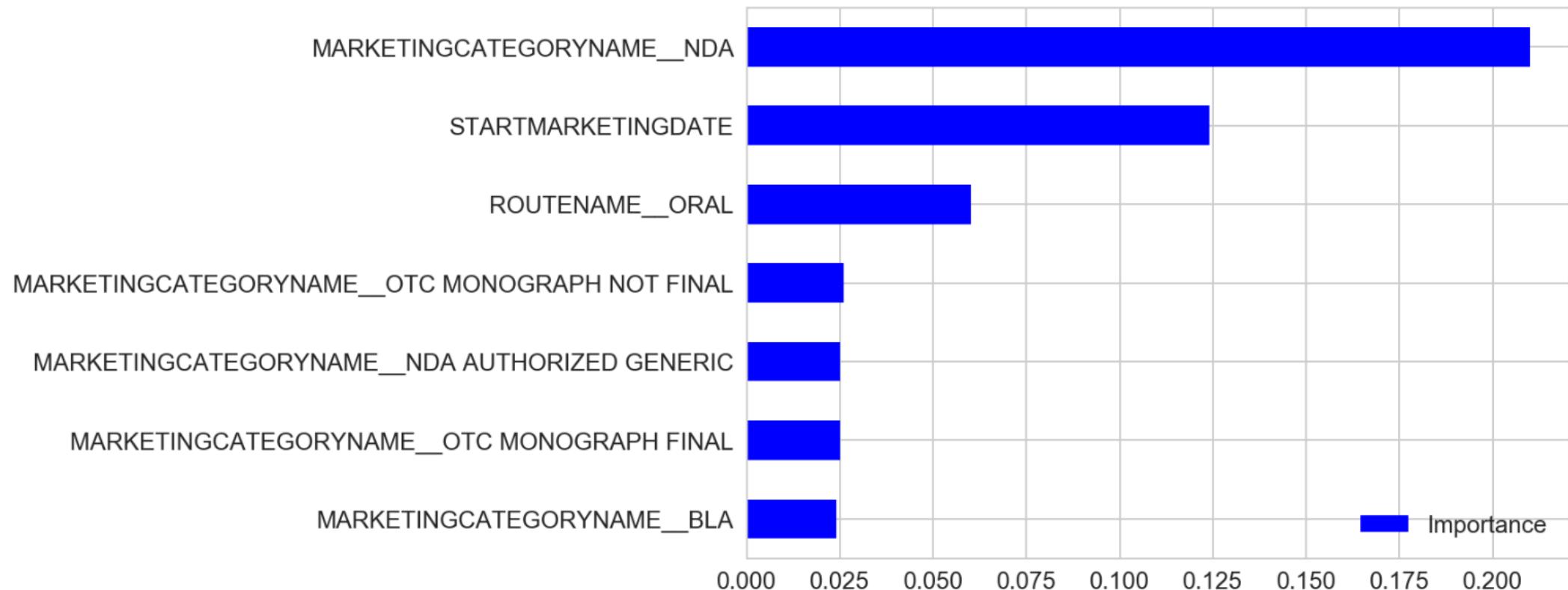
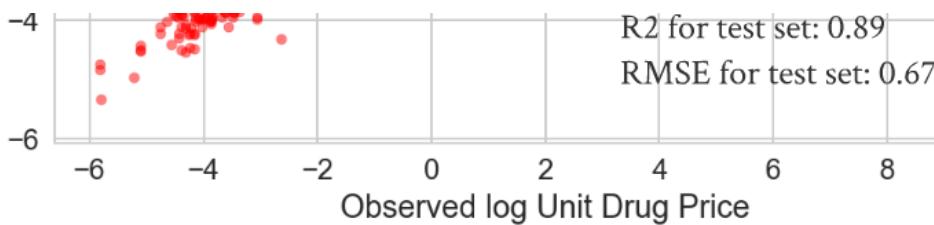
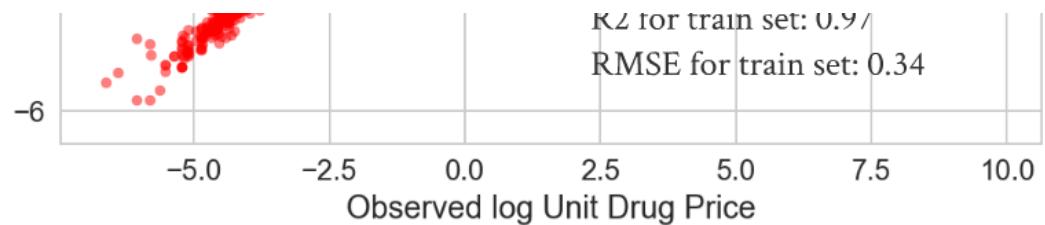
- Prevents overfitting
- Better Prediction



## • DEEPLIFT EXPLANATION

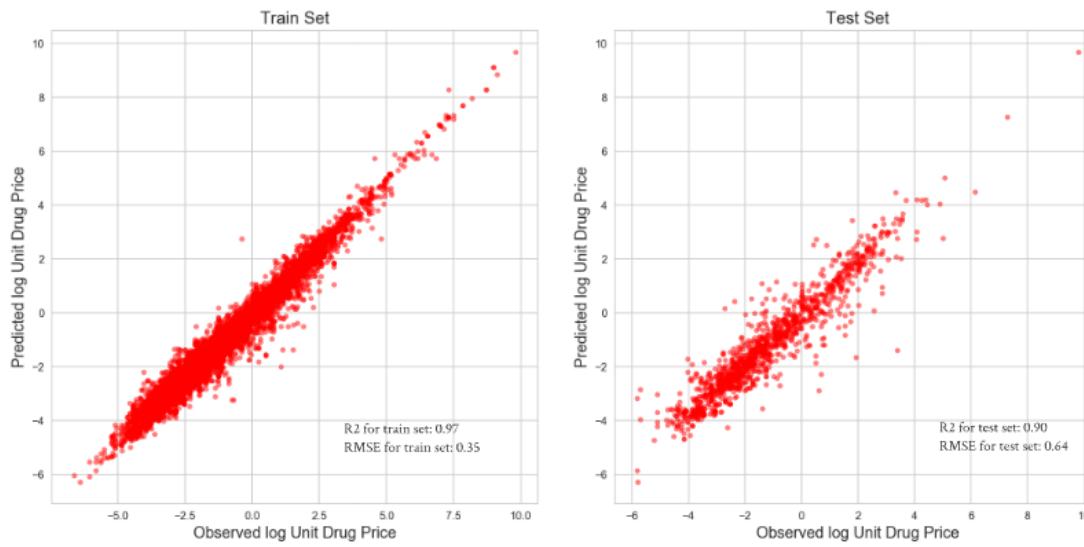
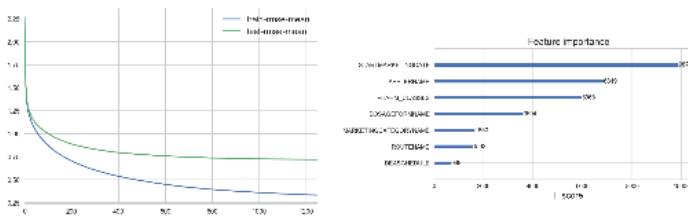


MARKETINGCATEGORYNAME\_NDA

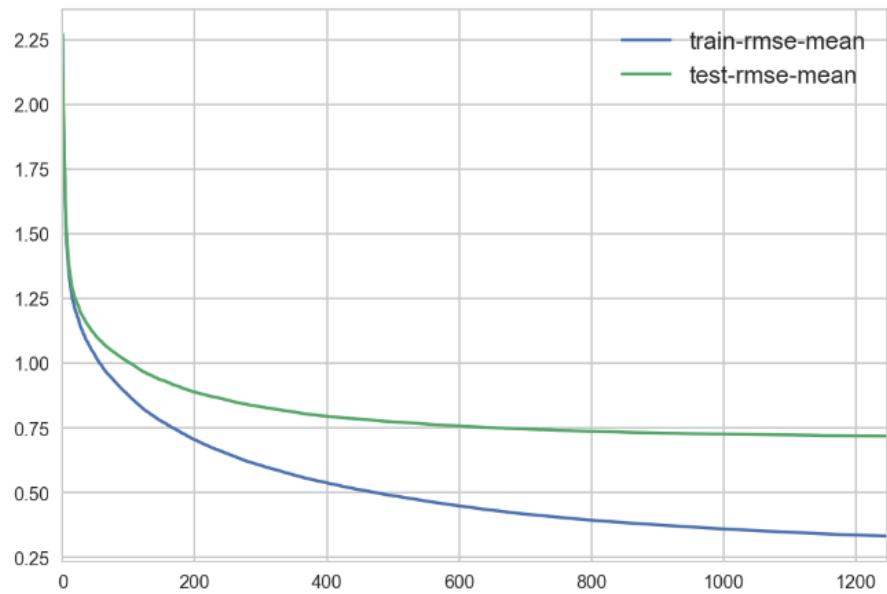


# XGBoost

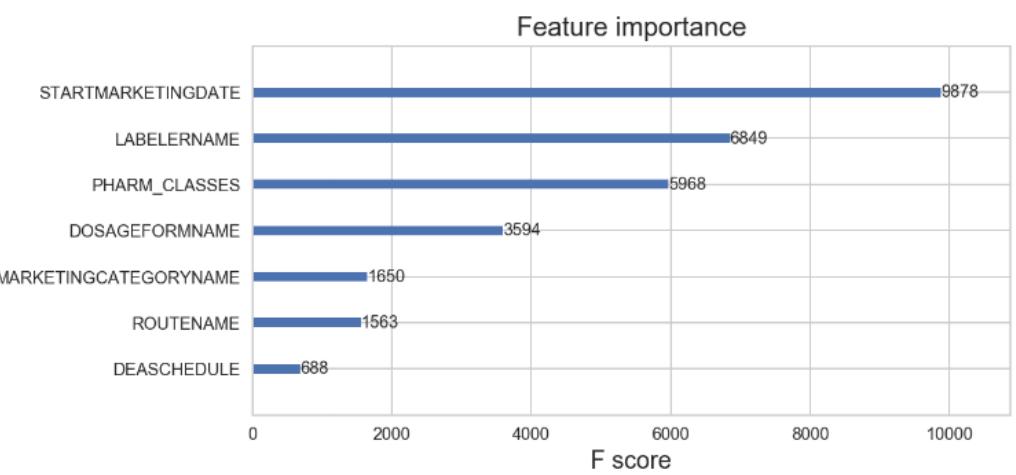
- Performs slightly better than Random Forest
- Cross Validation, 5-fold



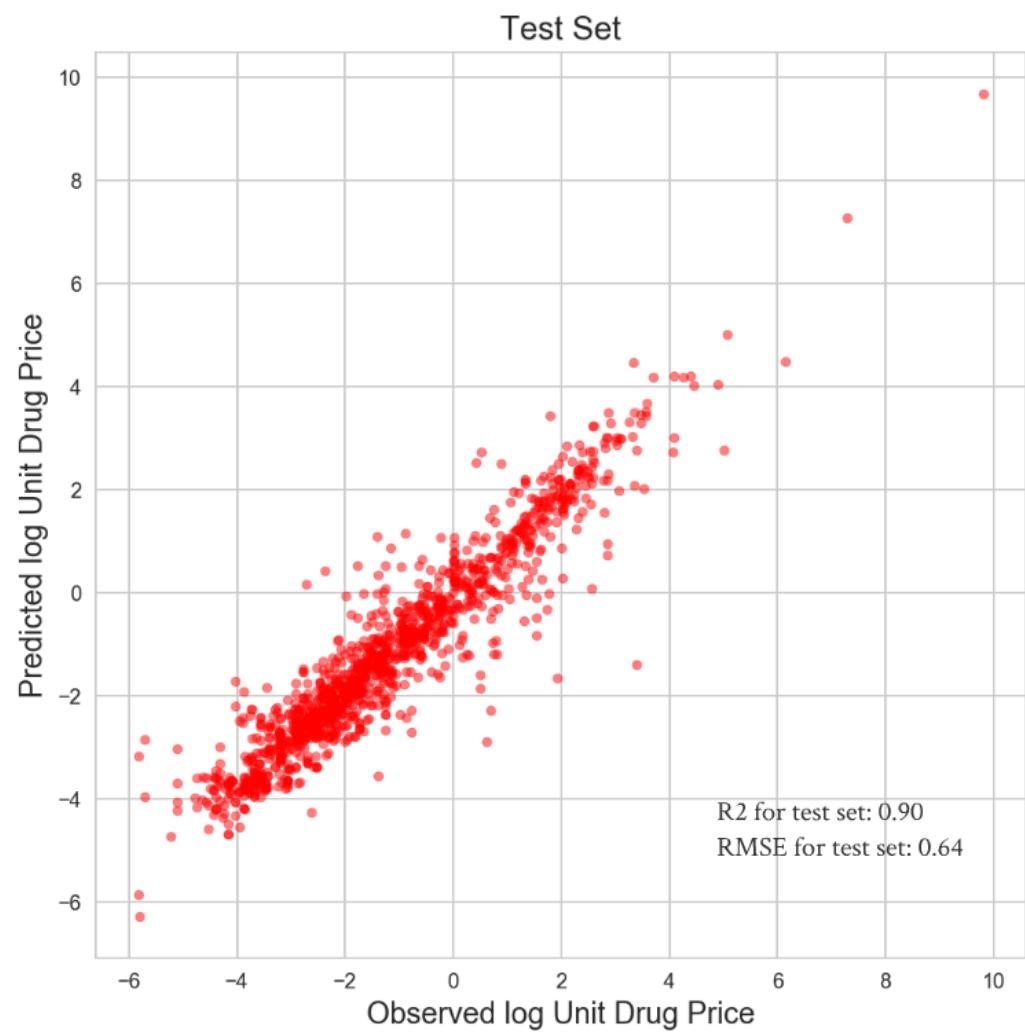
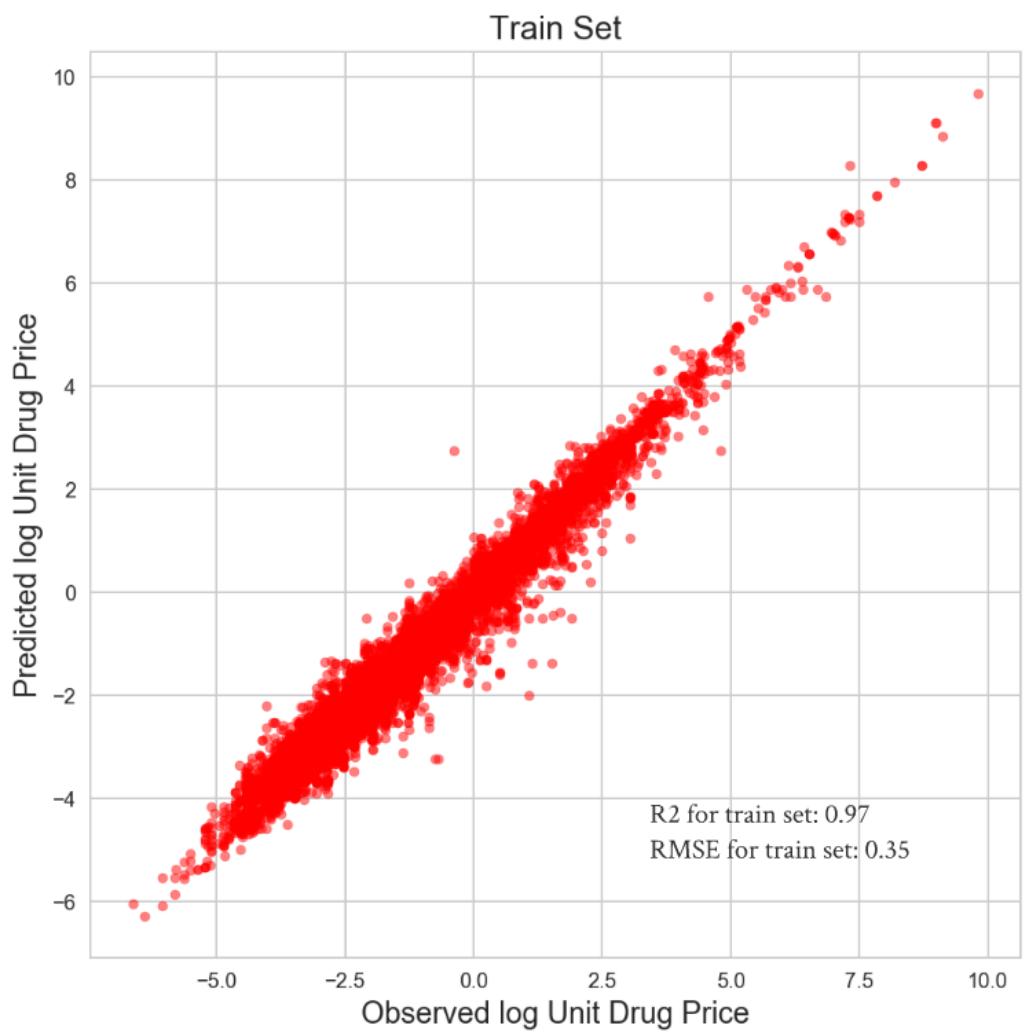
- Cross Validation, 5-fold



Train Set



Test Set



# Drug Price Classification

Logistic Regression

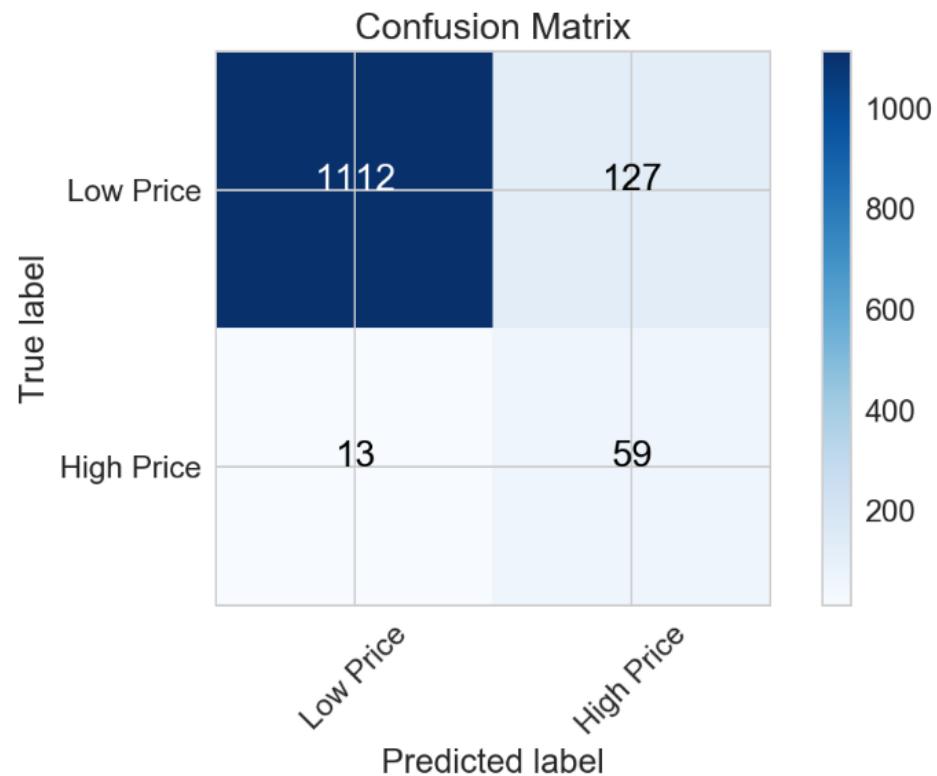
Decision Tree

Random Forest

XGboost

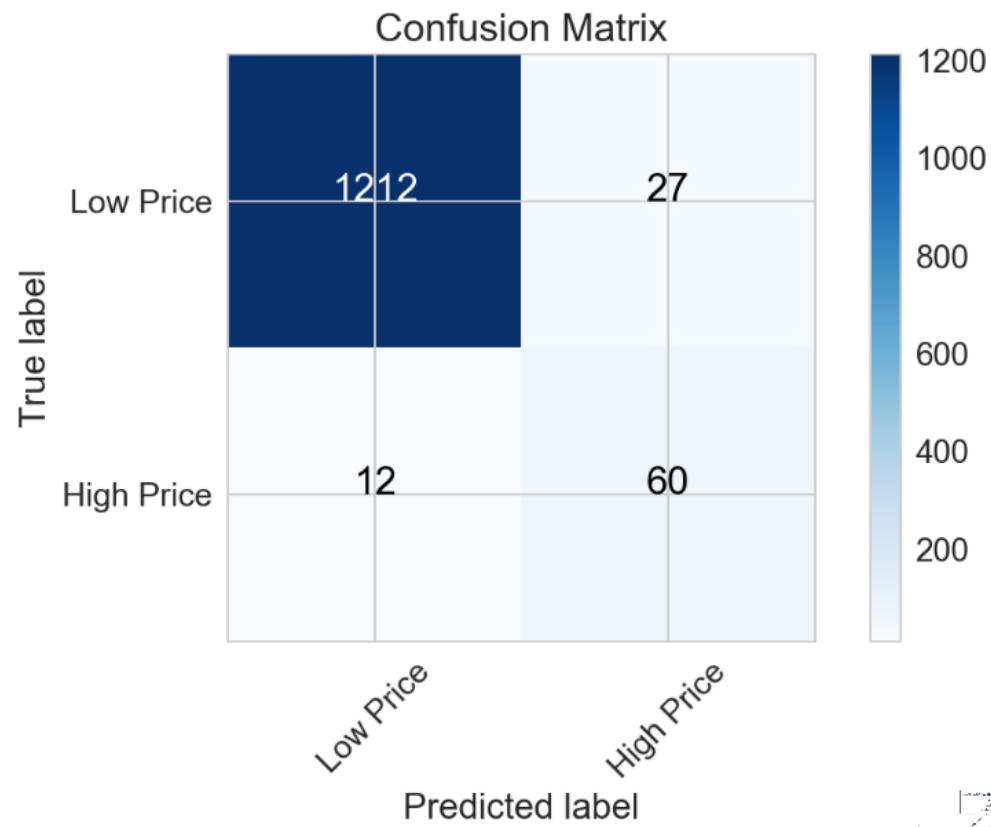
# Logistic Regression

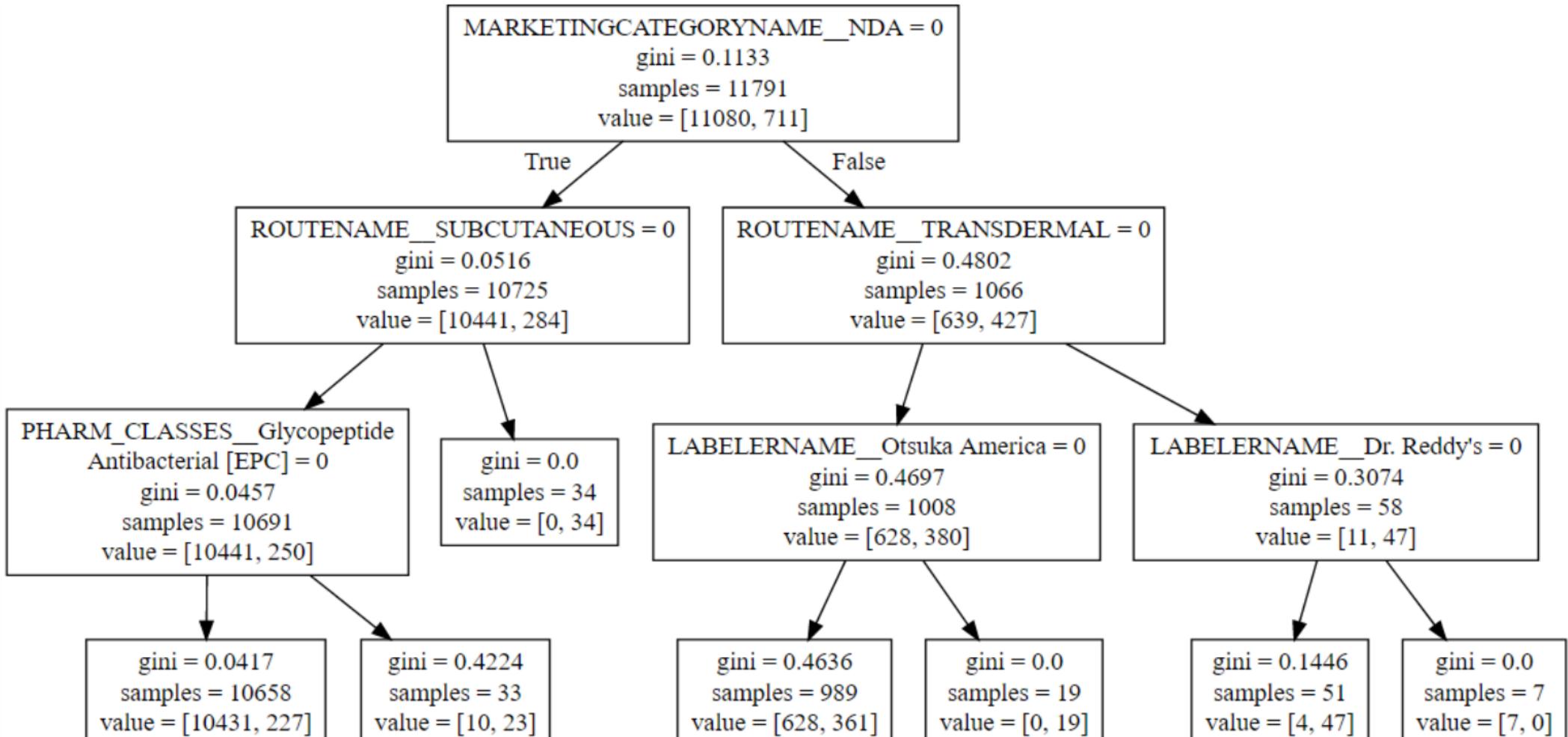
- Poor precision on high price drugs
- Accuracy Score: 0.92
- Low Price Precision: 0.99
- High Price Precision: 0.32
- Stratified sampling of cross validation, 5-fold



# Decision Tree

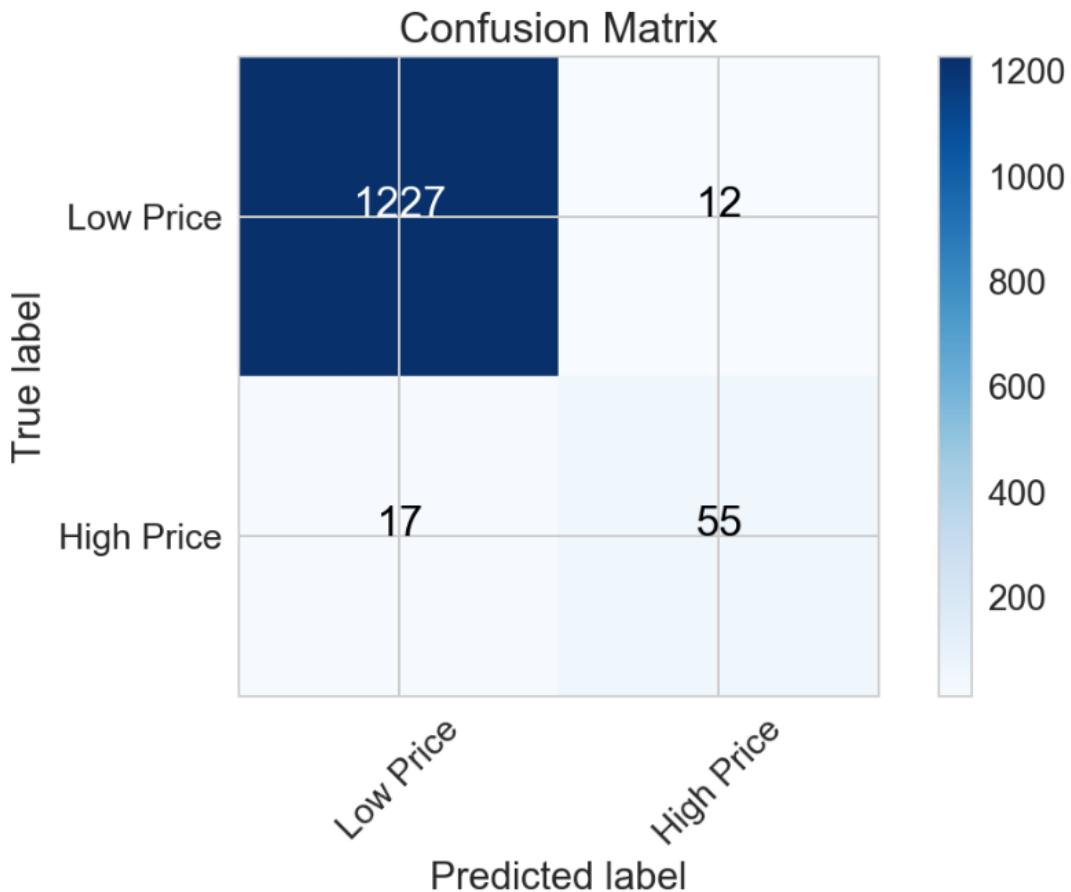
- Improved precision on high price drugs
- Accuracy Score: 0.97
- Low Price Precision: 0.99
- High Price Precision: 0.69
- Stratified sampling of cross validation





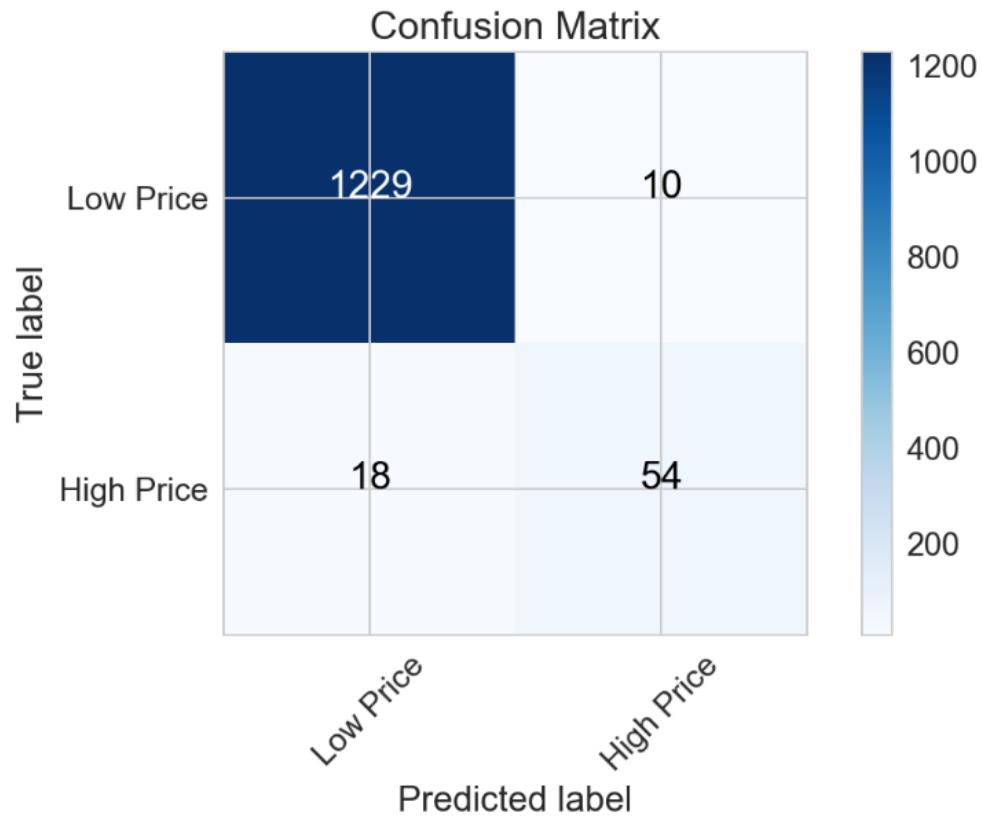
# Random Forest

- Much improved estimation on high price drugs
- Accuracy Score: 0.98
- Low Price Precision: 0.99
- High Price Precision: 0.82



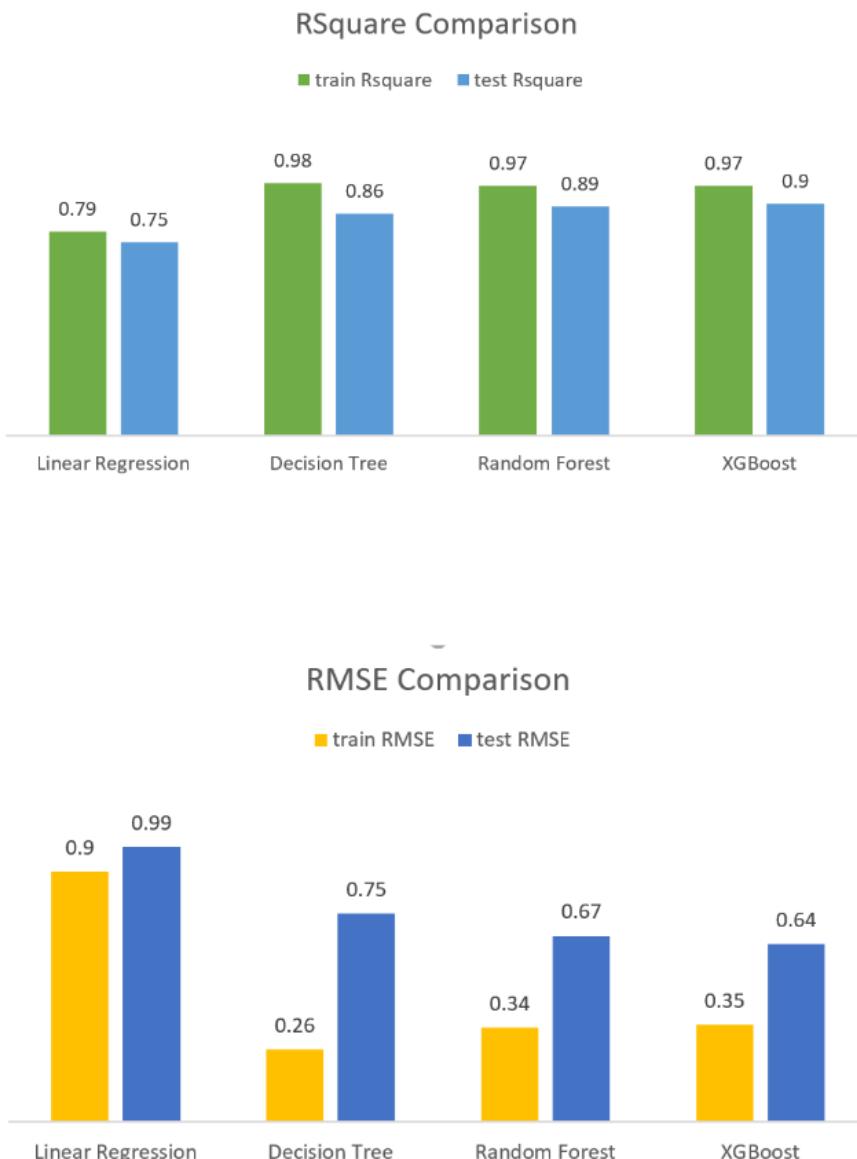
# XGBoost

- Best result among all models
- Accuracy Score: 0.98
- Low Price Precision: 0.99
- High Price Precision: 0.84



# Comparisons: Machine Learning Models

## Price Prediction Models



## Price Classification Models



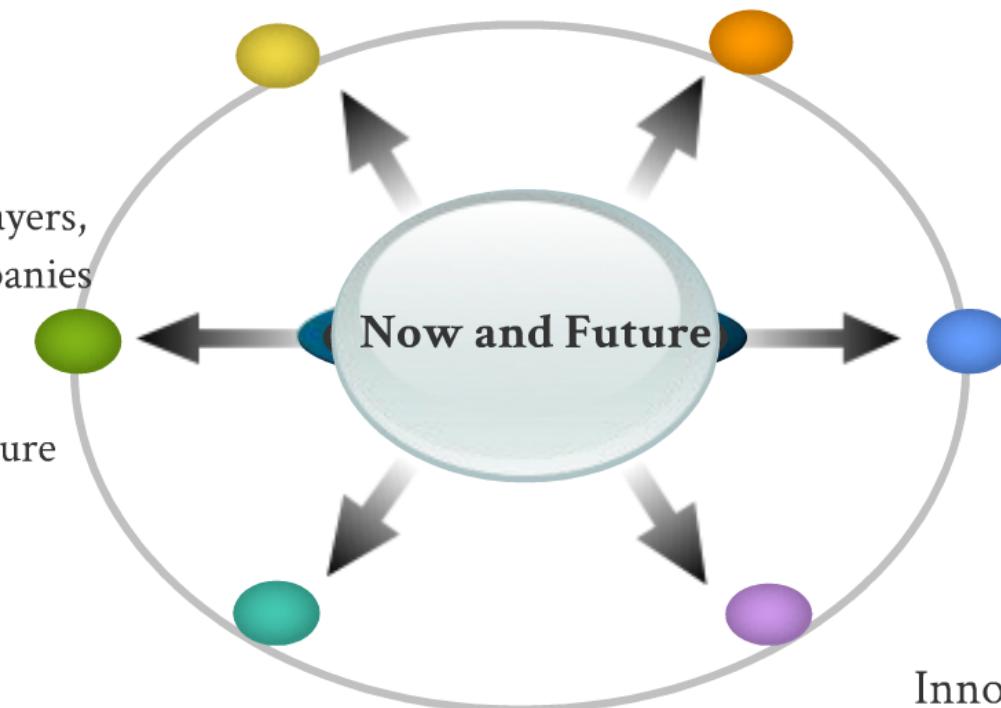
# Conclusions

Prediction and classification models  
Reasonably robust

Regression<Tree<Random Forest<=XGBoost

Significant values:

- patients, prescribers, payers, medical insurance companies
- new drug development evaluation
- target selection for venture capital investment.



Apply to other fields and products

Expand the dataset  
Optimize the model  
supply-demand,  
economics, policy,  
regulatory variables,  
time series analysis

Innovative work  
for publication

# Acknowledgment

The NYCDSA Team

The Fellow Bootcamp Students

The Core Team: Stay Healthy, Live Longer!

# Questions

---

---

---