

# Applied Machine Learning Workshop

*Come to our workshop to learn how to build Python machine learning systems for your data science projects!*

Thursday Feb. 28

Pupin 420, 8–9pm



Columbia Data Science Society



Connect with Us. Sign up for our News Letter.

<https://cdssatcu.com>



colab

## Basic Introductory Tutorial

[https://colab.research.google.com/notebooks/  
welcome.ipynb](https://colab.research.google.com/notebooks/welcome.ipynb)

*dmlc*  
***XGBoost***



**Sci-Kit Learn:** <https://scikit-learn.org/stable/>

**XGBoost:** <https://xgboost.readthedocs.io/en/latest/>

# Acknowledgements:

The content of this talk is motivated from the two classes, I am taking this semester at Columbia Data Science Institute.

<div> <div>Andreas C. Müller - Associate Research Scientist</div> <div> <div>Courses</div> <div>Personal Website</div> </div> </div>	
<div> <div>COMS W4995</div> <div> <div>Overview</div> <div>Schedule</div> <div>Piazza</div> <div>Courseworks</div> <div>GitHub</div> </div> </div>	
	<div>COMS W4995 Applied Machine Learning Spring 2019 - Schedule</div>
	<div>Press P on slides for presenter notes.</div>

COMS W4995 Applied Machine Learning  
taught by **Andread C. Müller**, core  
contributer of sci-kit learn and author of  
the O'Reilly book "Introduction to  
machine learning with Python", describing  
a practical approach to machine learning  
with python and scikit-learn.

<p align="center"><b>COMS W4212 Machine Learning for Data Science</b></p> <p align="center"><b>Columbia University, Spring 2017</b></p>		
<p>Instructor: John Paisley  Location: 501 Schermerhorn Hall  Time: T/Th 7:40pm - 8:55pm  Office hours: Monday 11am-12pm @ 422 Mudd Building</p>		
TA's:	Ghazal Fazelnia Tianhao Lu Dheeraj Kalmekolan Ashutosh Nanda Avinash Bukkittu Yuhao Zhang Jiefu Ying George Yu Peng Wu	gf2293@columbia.edu tf2710@columbia.edu drk2143@columbia.edu an2655@columbia.edu ab4377@columbia.edu js3044@columbia.edu jy2799@columbia.edu gy2206@columbia.edu pw2393@columbia.edu
		CVN office hours via email (no fixed time)  Tue/Thu 1pm - 2pm @ CS TA room, Mudd 122A (1st floor) Wed 4:30pm - 6:30pm @ CS TA room, Mudd 122A (1st floor) Fri 1:30pm - 3:30pm @ CS TA room, Mudd 122A (1st floor) Tues 10am - 12pm @ CS TA room, Mudd 122A (1st floor) Tues 9pm - 11pm @ CS TA room, Mudd 122A (1st floor) Fri 4pm - 6pm @ CS TA room, Mudd 122A (1st floor) Tues 4pm - 5pm & Wed 12pm-1pm @ CS TA room, Mudd 122A (1st floor) Mon 6:30pm - 8:30pm @ CS TA room, Mudd 122A (1st floor)
<p><b>Synopsis:</b> This course provides an introduction to supervised and unsupervised techniques for machine learning. We will cover both probabilistic and non-probabilistic approaches to machine learning. Focus will be on classification and regression models, clustering methods, matrix factorization and sequential models. Methods covered in class include linear and logistic regression, support vector machines, boosting, K-means clustering, mixture models, expectation-maximization algorithm, hidden Markov models, among others. We will cover algorithmic techniques for optimization, such as gradient and coordinate descent methods, as the need arises.</p>		
<p><b>Prerequisites:</b> Basic linear algebra and calculus, introductory-level courses in probability and statistics. Comfort with a programming language (e.g., Matlab) will be essential for completing the homework assignments. Not open to students who have taken COMS 4771, STATS 4400 or IEOB 4525.</p>		
<p><b>Text:</b> There is no required text for the course. Suggested readings for each class will be given from the textbooks below. These readings are meant to be general pointers and may contain more material than we cover in class.</p>		
<p>T. Hastie, R. Tibshirani and J. Friedman, <i>The Elements of Statistical Learning, Second Edition</i>, Springer. <a href="#">[link]</a>  C. Bishop, <i>Pattern Recognition and Machine Learning</i>, Springer. <a href="#">[link]</a>  H. Daume, <i>A Course in Machine Learning</i>, Draft. <a href="#">[link]</a></p>		
<p><b>Grading:</b> 5 homework assignments (50%), midterm exam (25%), final in-class exam (25%). Each homework assignment will have a programming component that will</p>		

COMS W4721 Machine Learning for Data Science  
taught by **Prof. John Paisley**. He also teaches  
a course on Machine Learning through edX. The  
content of his on-campus course largely overlaps  
with his edX class.

# Let's get started.

EXAMPLES

RECENT

GOOGLE DRIVE


GITHUB


UPLOAD


Enter a GitHub URL or search by organization or user

☐ Include private repos

[https://github.com/AkhilPunia/cdss\\_applied\\_ml\\_2019](https://github.com/AkhilPunia/cdss_applied_ml_2019)




Repository: 



Branch: 

AkhilPunia/cdss\_applied\_ml\_2019

master

Path

 cdss\_aml\_talk.ipynb

NEW PYTHON 3 NOTEBOOK

CANCEL

# Part 1/5

Let's talk about the Data





## House Prices: Advanced Regression Techniques

Predict sales prices and practice feature engineering, RFs, and gradient boosting

4,189 teams · Ongoing

[Overview](#)[Data](#)[Kernels](#)[Discussion](#)[Leaderboard](#)[Rules](#)[Team](#)[My Submissions](#)[Submit Predictions](#)

### Overview

#### Description

#### Evaluation

#### Tutorials

#### Frequently Asked Questions

### Start here if...

You have some experience with R or Python and machine learning basics. This is a perfect competition for data science students who have completed an online course in machine learning and are looking to expand their skill set before trying a featured competition.

### Competition Description



Ask a home buyer to describe their dream house, and they probably won't begin with the height of the basement ceiling or the proximity to an east-west railroad. But this playground competition's dataset proves that much more influences price negotiations than the number of bedrooms or a white-picket fence.

With 79 explanatory variables describing (almost) every aspect of residential homes in Ames, Iowa, this competition challenges you to predict the final price of each home.



# Zillow Prize: Zillow's Home Value Prediction (Zestimate)

Can you improve the algorithm that changed the world of real estate?

**\$1,200,000**

Prize Money



Zillow · 3,779 teams · a year ago

Overview

Data

Kernels

Discussion

Leaderboard

Rules

Team

My Submissions

Late Submission

## Overview

### Description

### Evaluation

### Prizes

### Timeline

### Competition Overview

Zillow's Zestimate home valuation has shaken up the U.S. real estate industry since first released 11 years ago.

A home is often the largest and most expensive purchase a person makes in his or her lifetime. Ensuring homeowners have a trusted way to monitor this asset is incredibly important. The Zestimate was created to give consumers as much information as possible about homes and the housing market, marking the first time consumers had access to this type of home value information at no cost.

"Zestimates" are estimated home values based on 7.5 million statistical and machine learning models that analyze hundreds of data points on each property. And, by continually improving the median margin of error (from 14% at the onset to 5% today), Zillow has since become established as one of the largest, most trusted marketplaces for real estate information in the U.S. and a leading example of impactful machine learning.

**111 Archer Ave,  
New York, NY 10031**  
4 beds • 3 baths • 3,410 sqft

**FOR SALE**  
**\$1,175,000**  
Zestimate: \$1,275,448

EST. MORTGAGE  
**\$4,461/mo**  
[Get pre-qualified](#)

Built in 2009, perfectly blending elegance with functional living space. Excellent floor plan with 3 beds up and 1 on main. Open living, kitchen & dining w/ huge fireplace & sound views. Spacious kitchen w/ slab granite surfaces & center island. Huge master suite with Jacuzzi tub & separate shower. Features: hwd floors, all

CONTACT  
Your Name  
Phone  
Email  
I am interested in NY 10031  
☐ I want to see this property

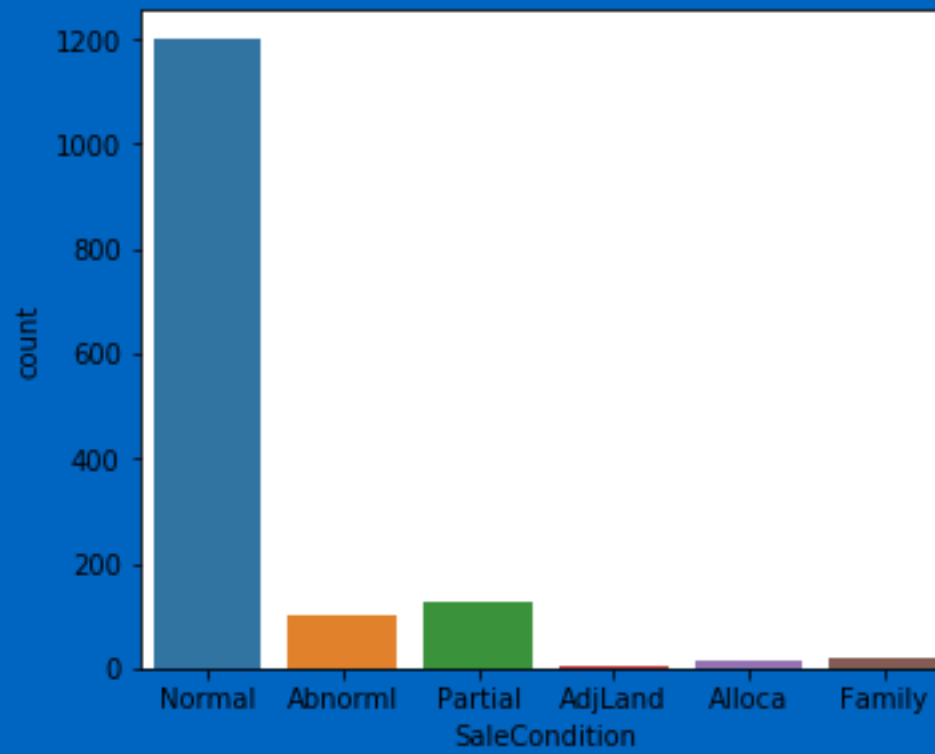
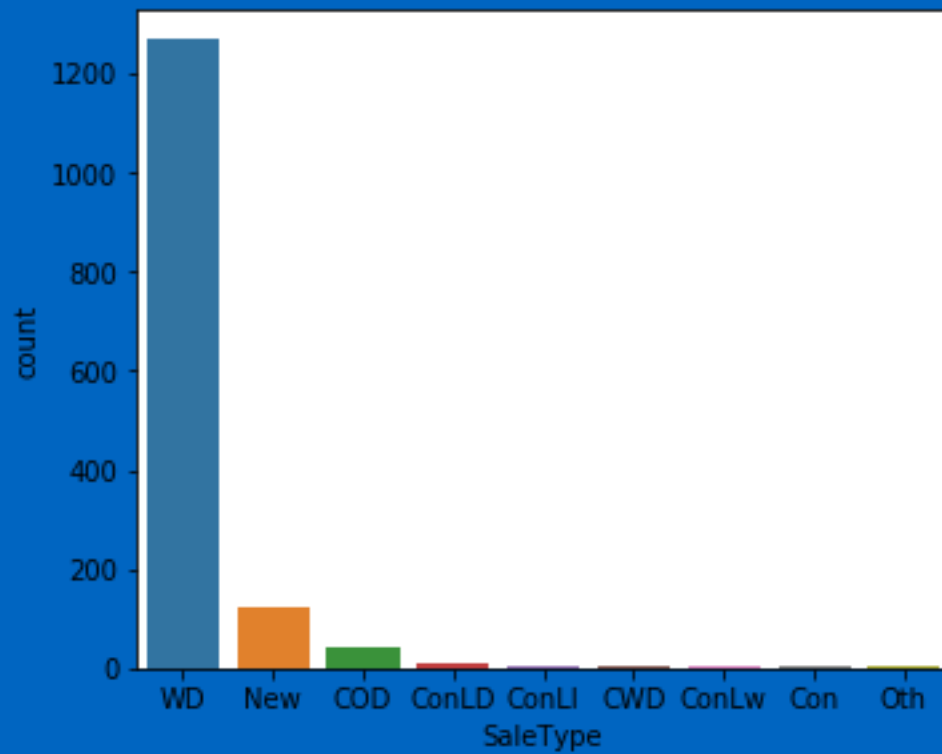
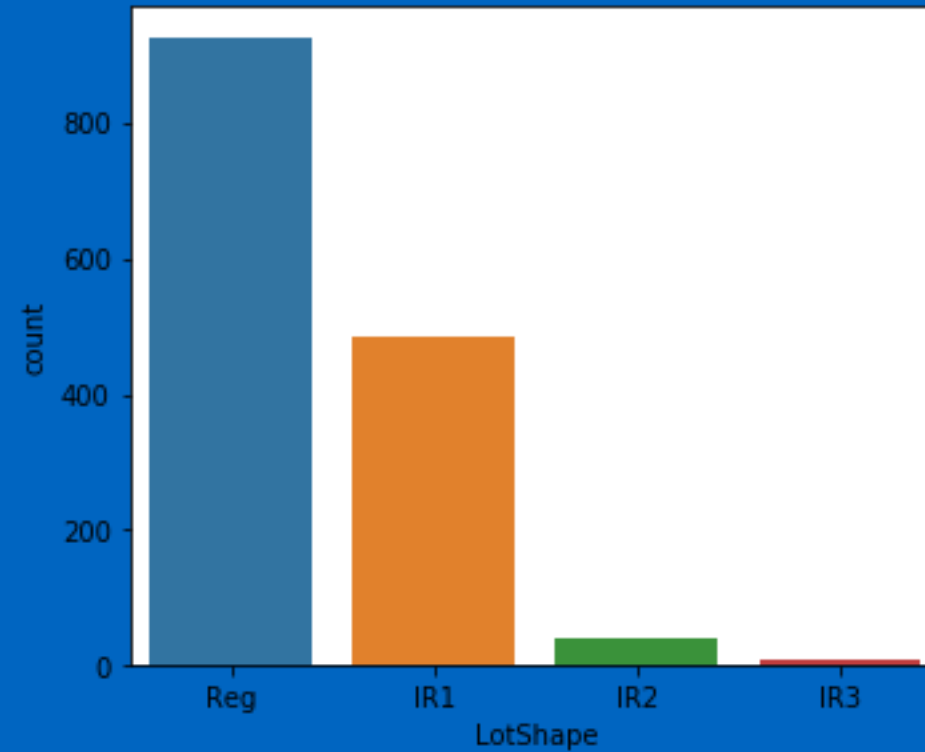
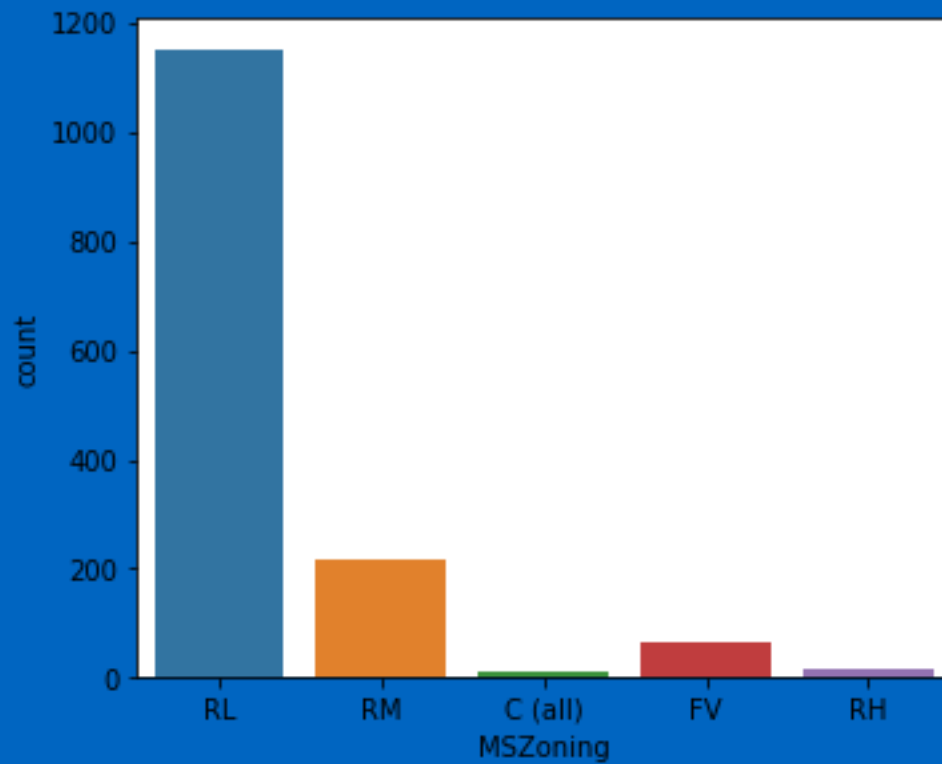
## Inflation-adjusted U.S. home prices, Population, Building costs, and Bond yields (1890–2005)

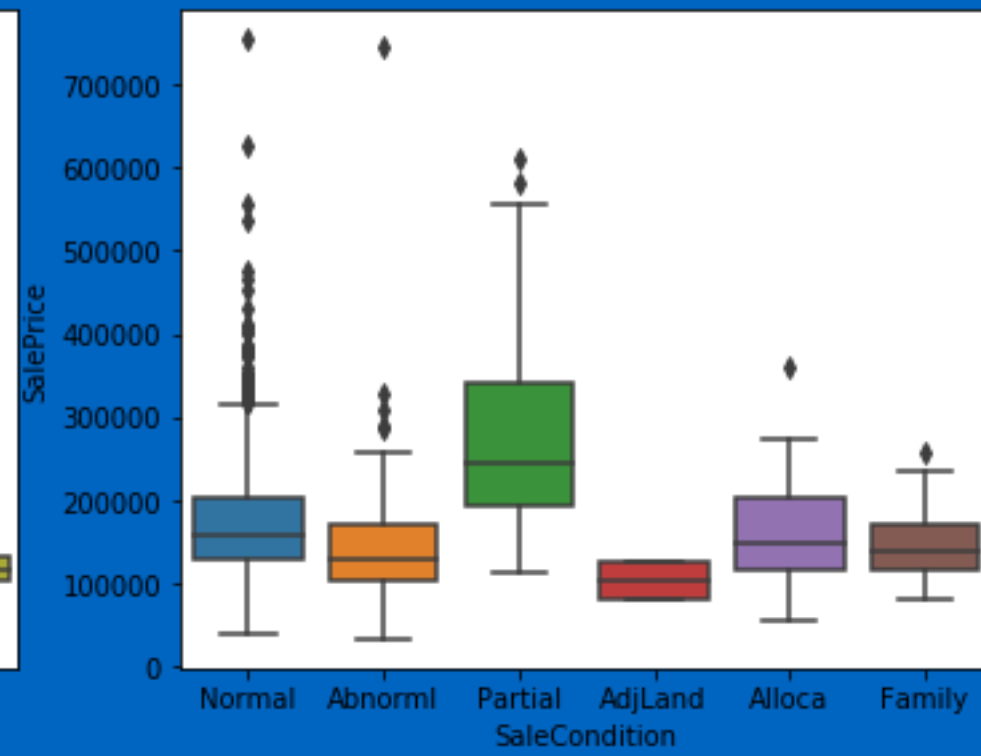
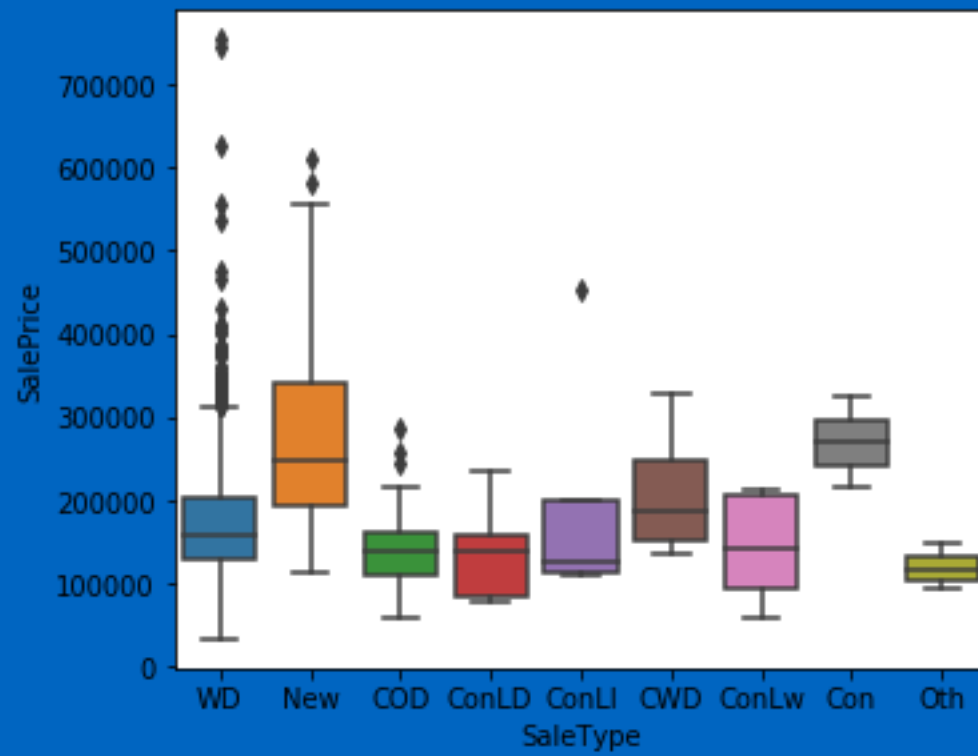
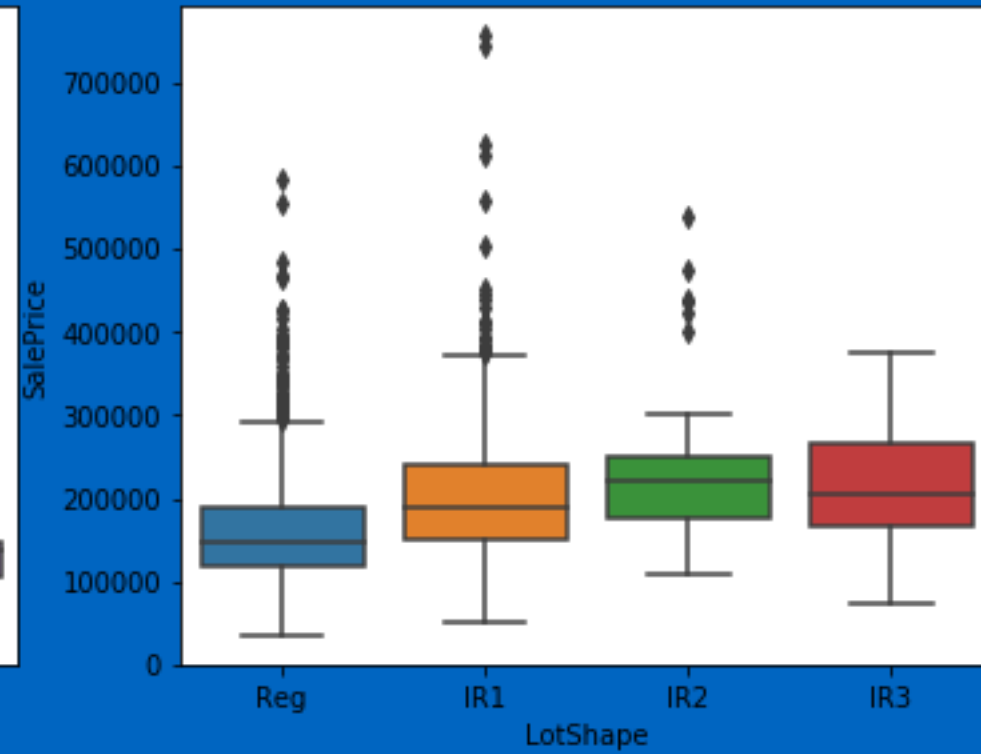
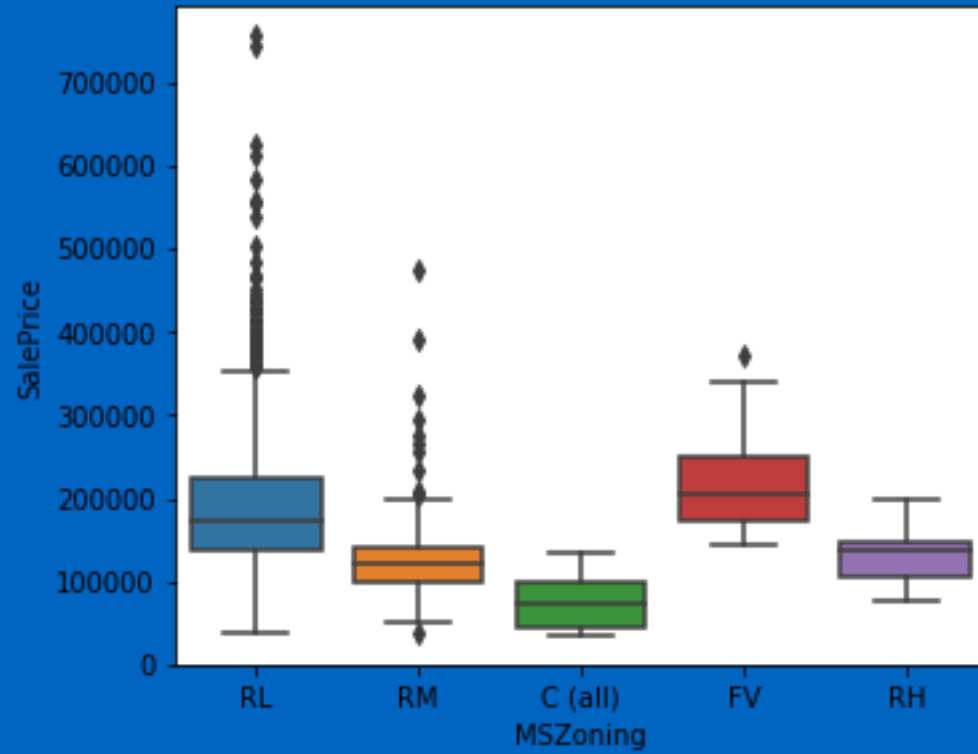


Source: *Irrational Exuberance*, 2d ed. (Fig. 2.1)

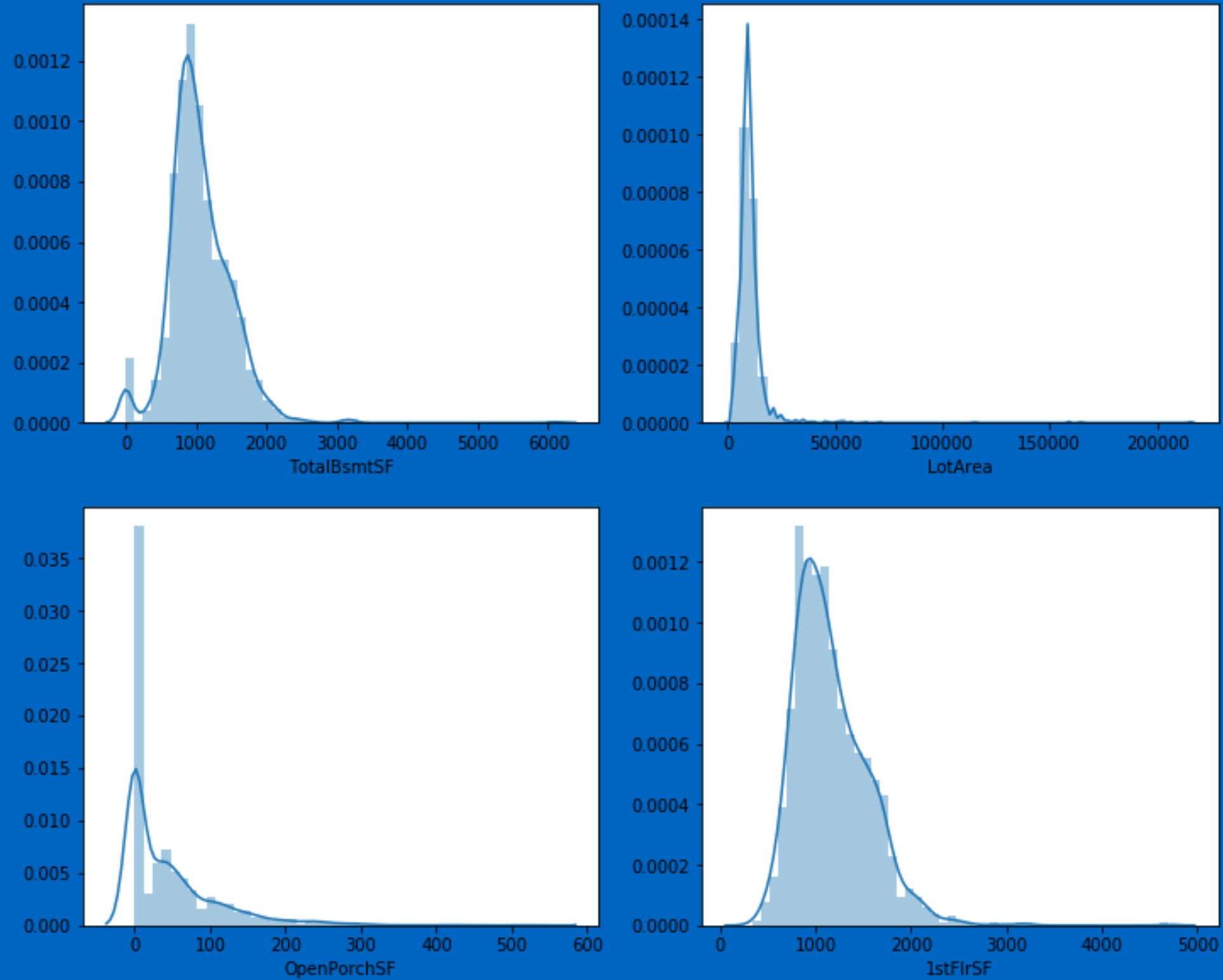
**What's with a feature?**

# Continuous

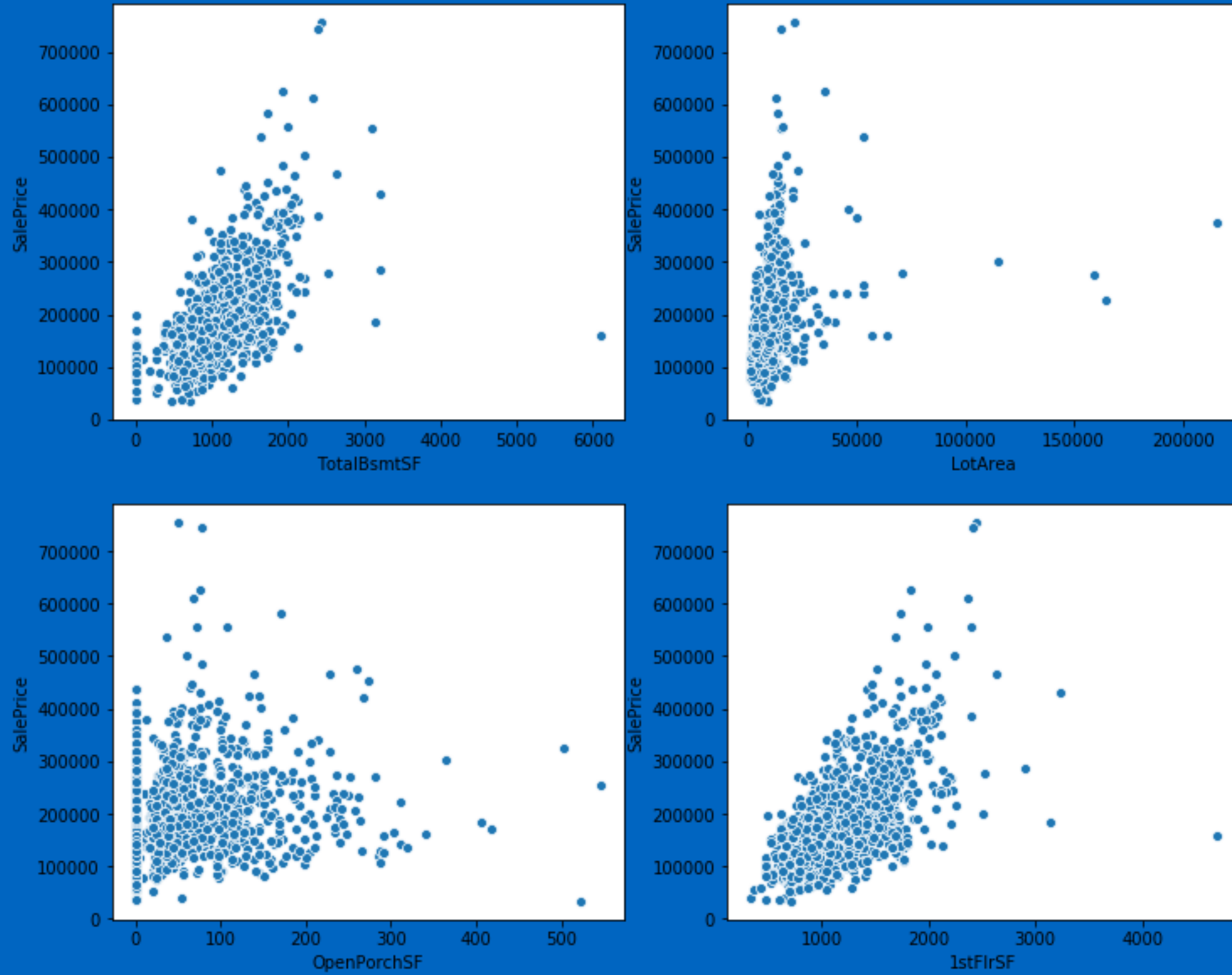




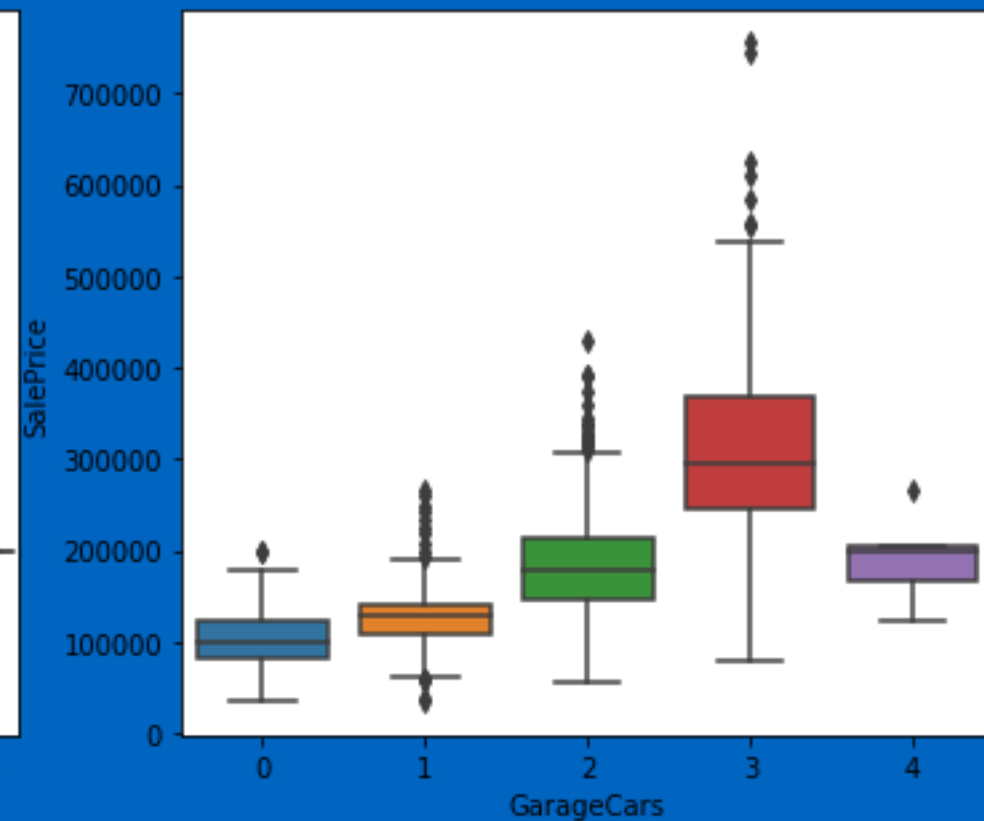
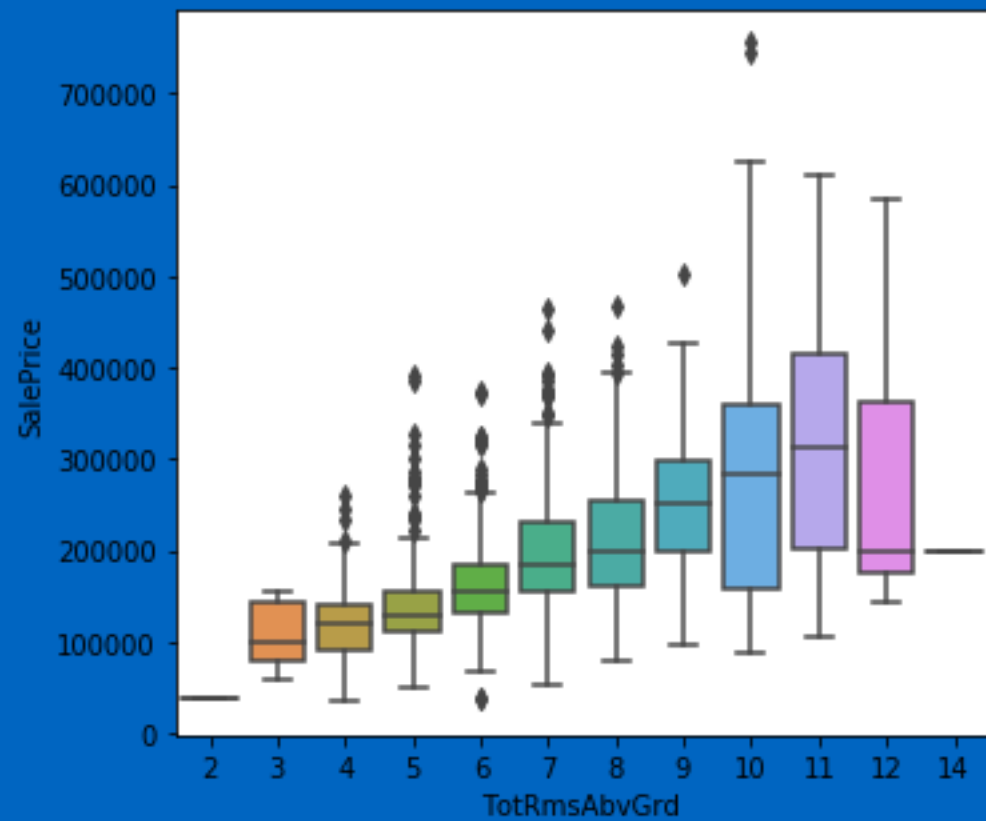
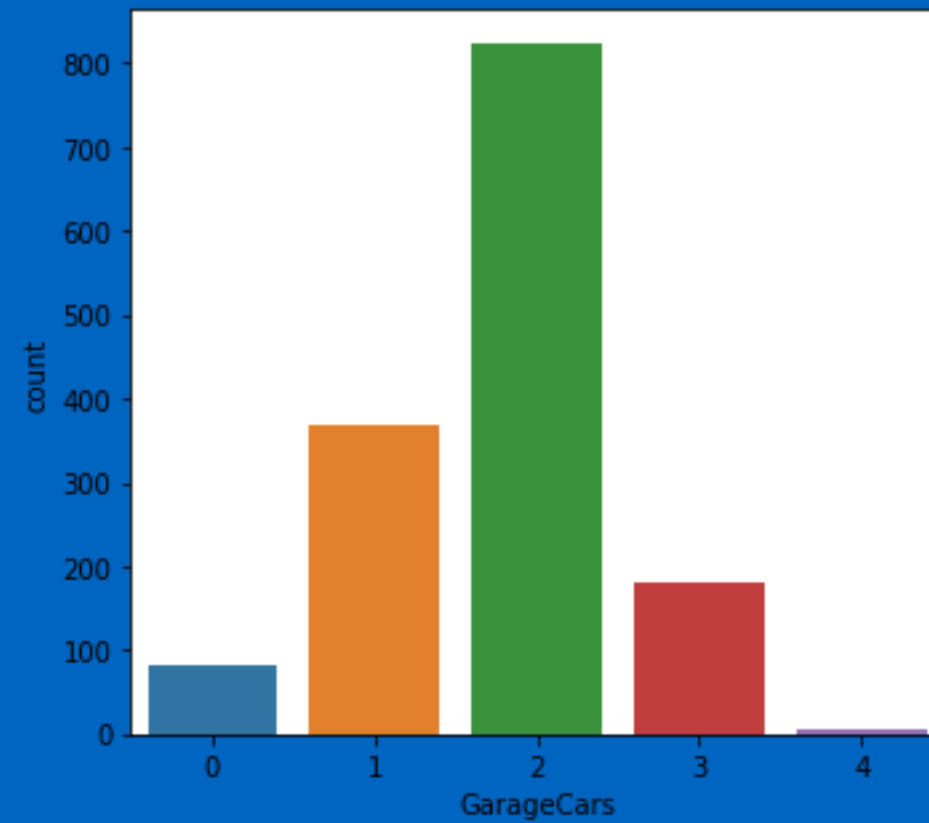
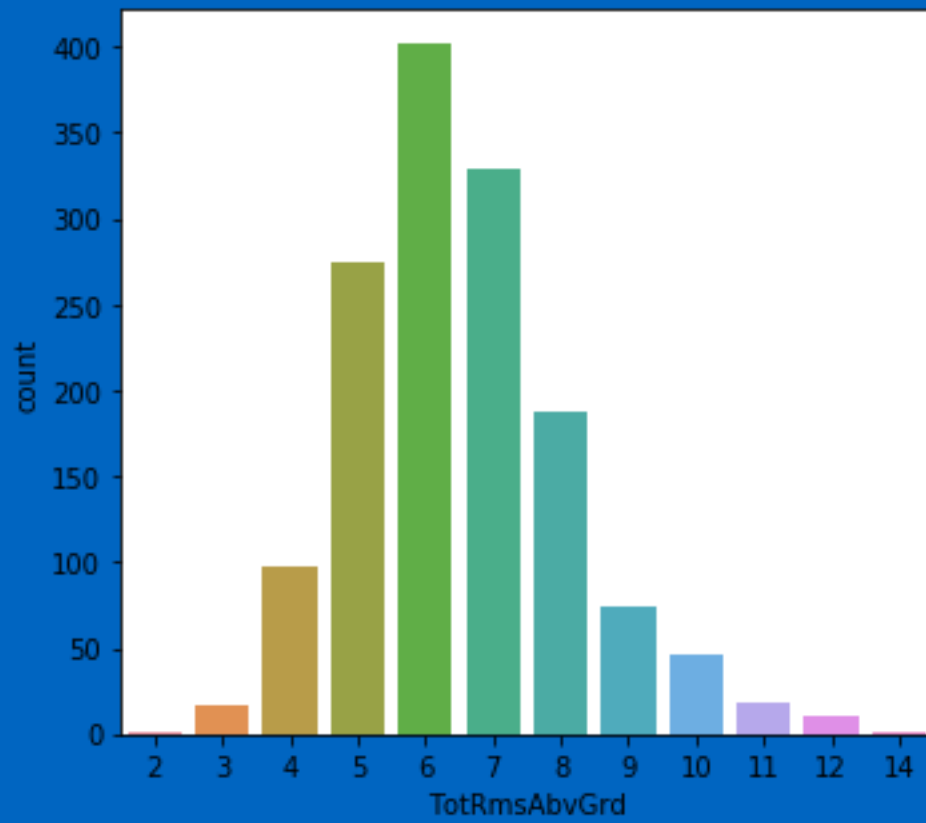
# Categorical



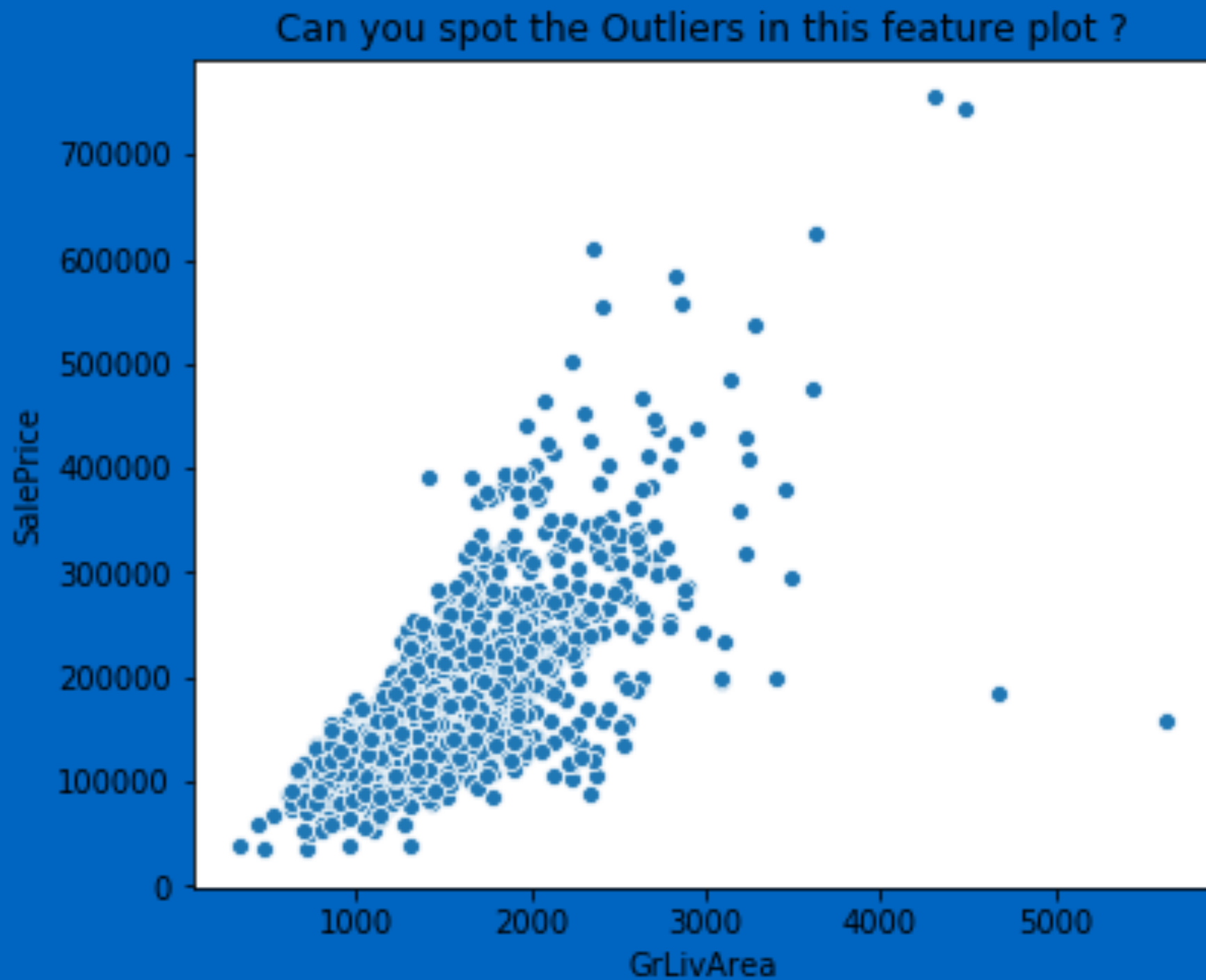




# Ordinal



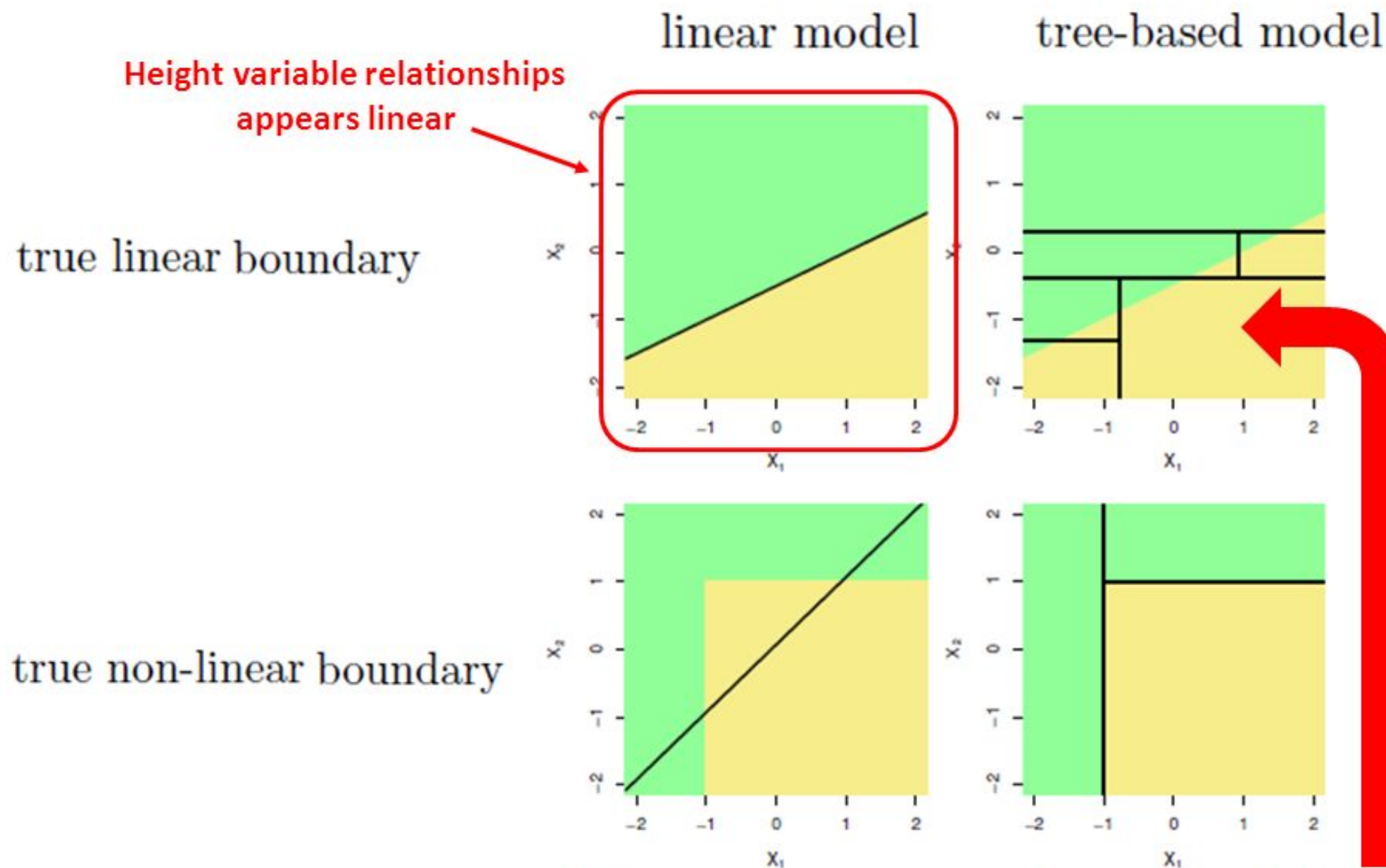
# Outliers



# Part 2/5

What's in a Model ?

# Linear Models vs. Decision Trees

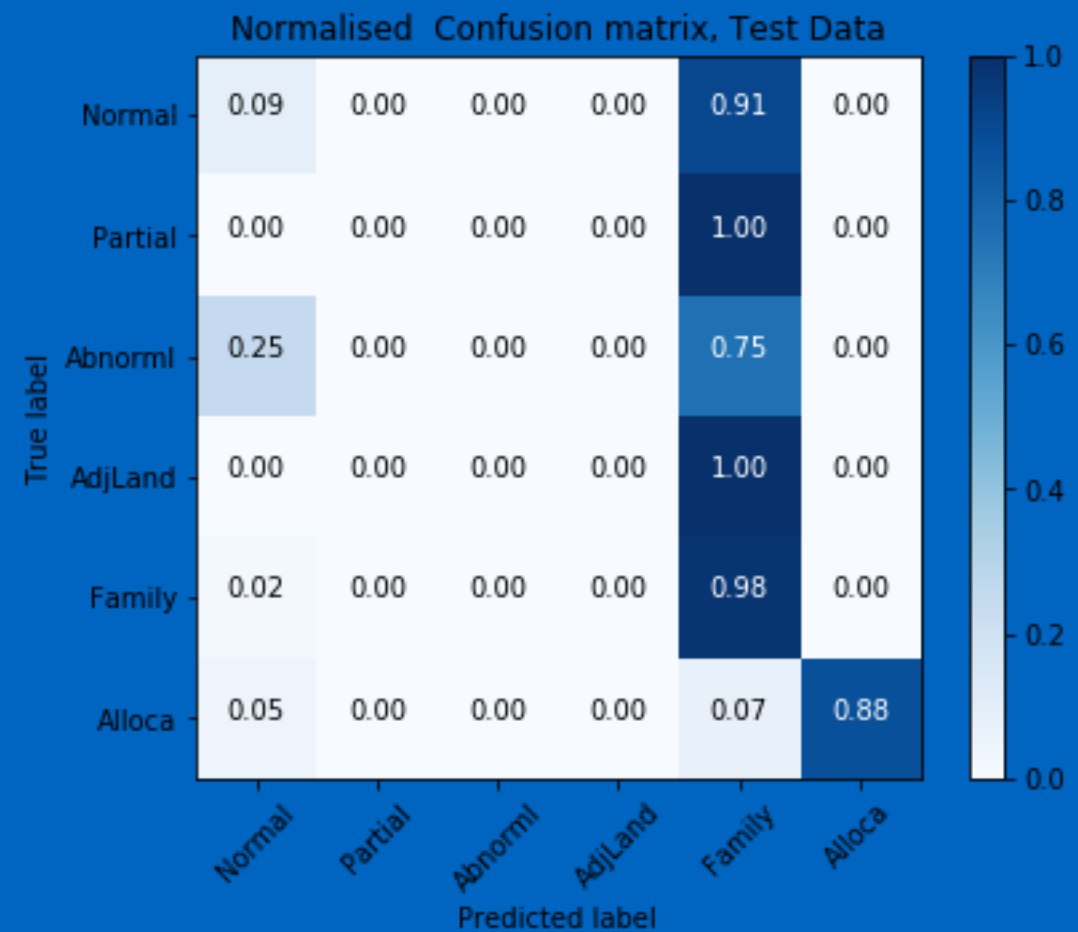
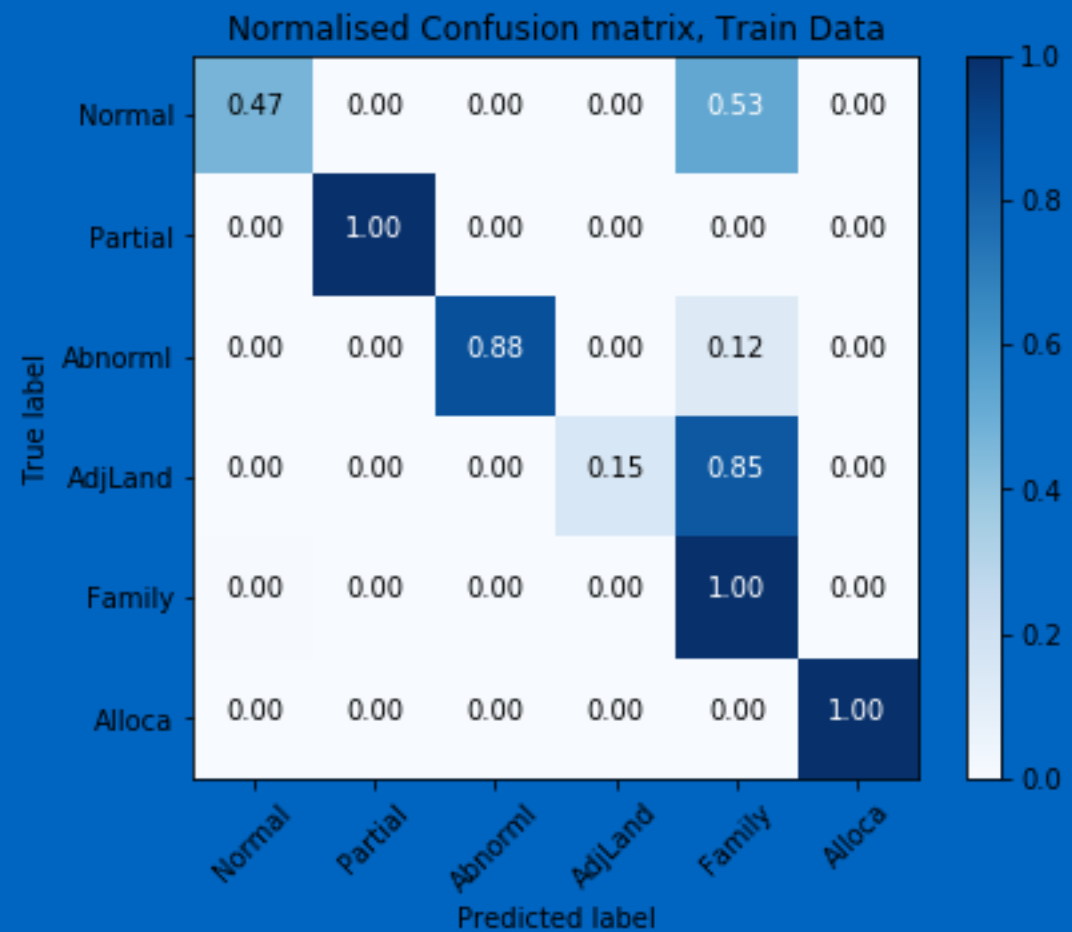
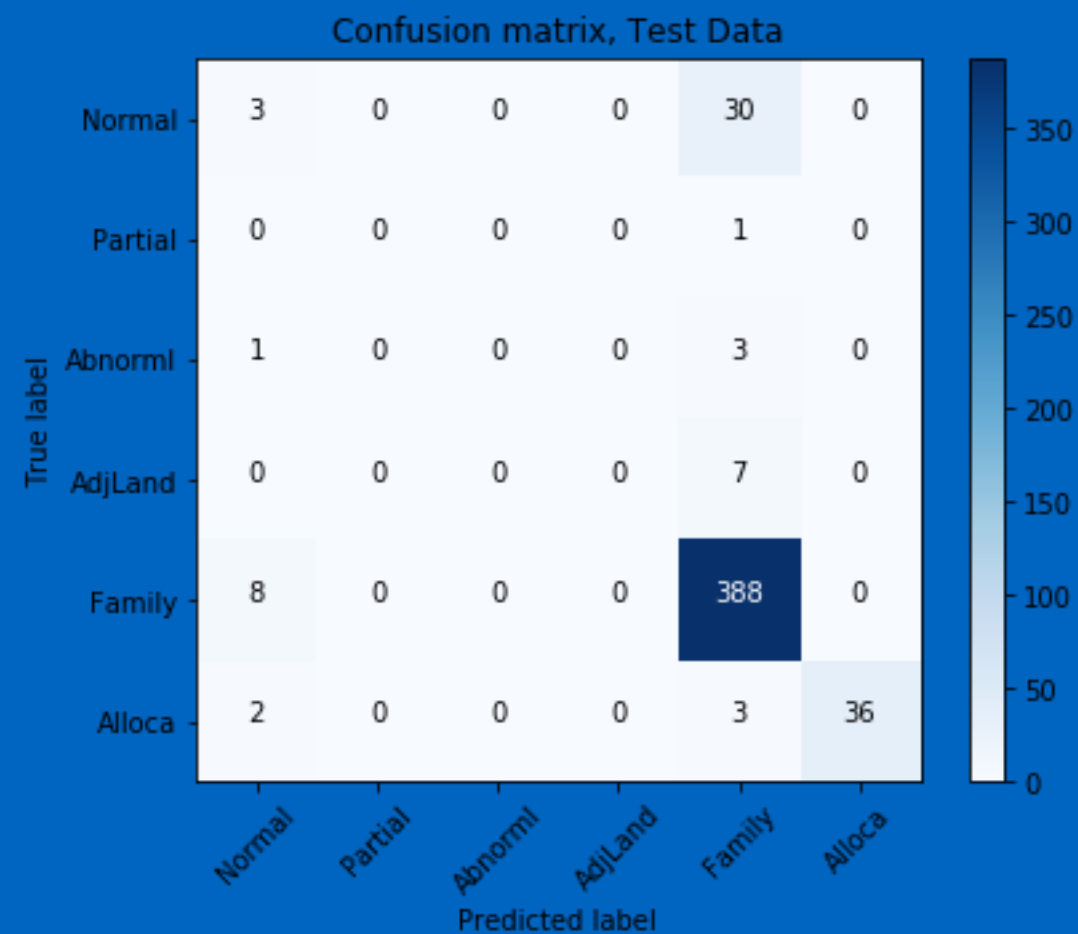
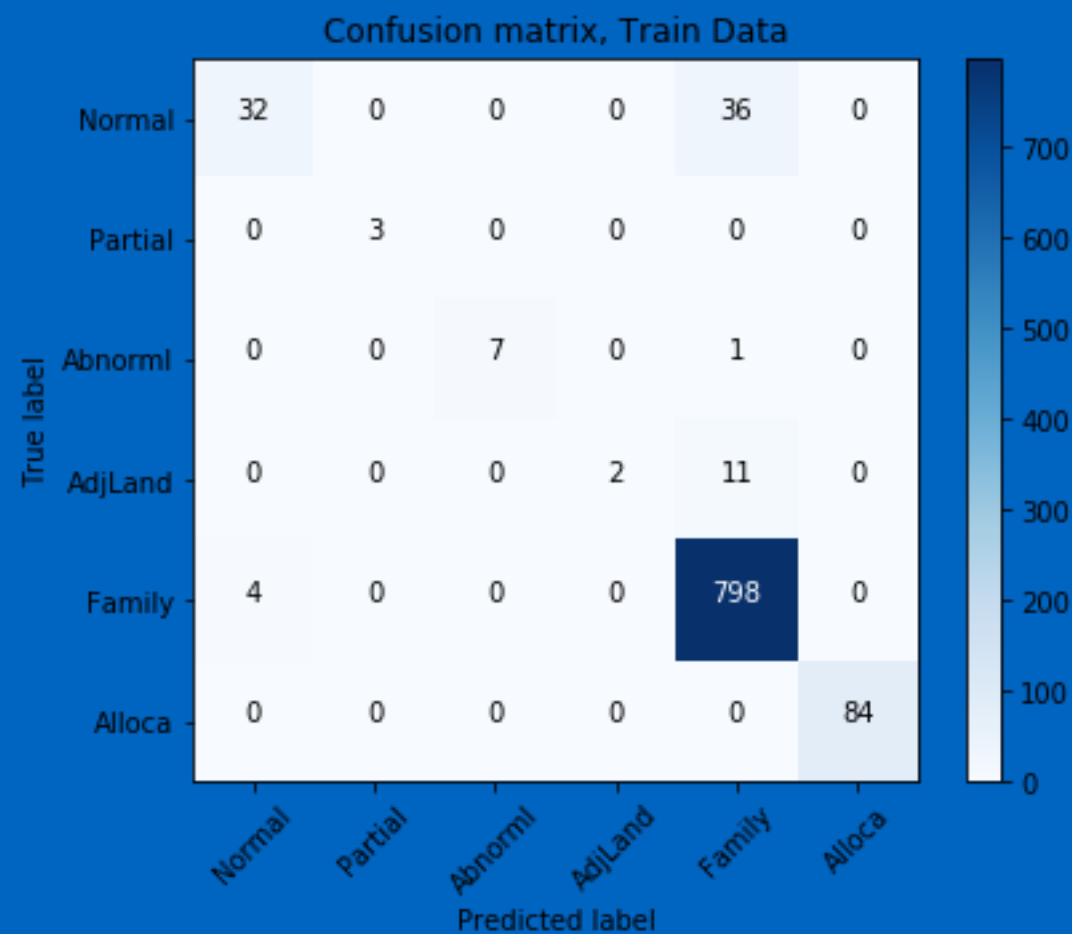


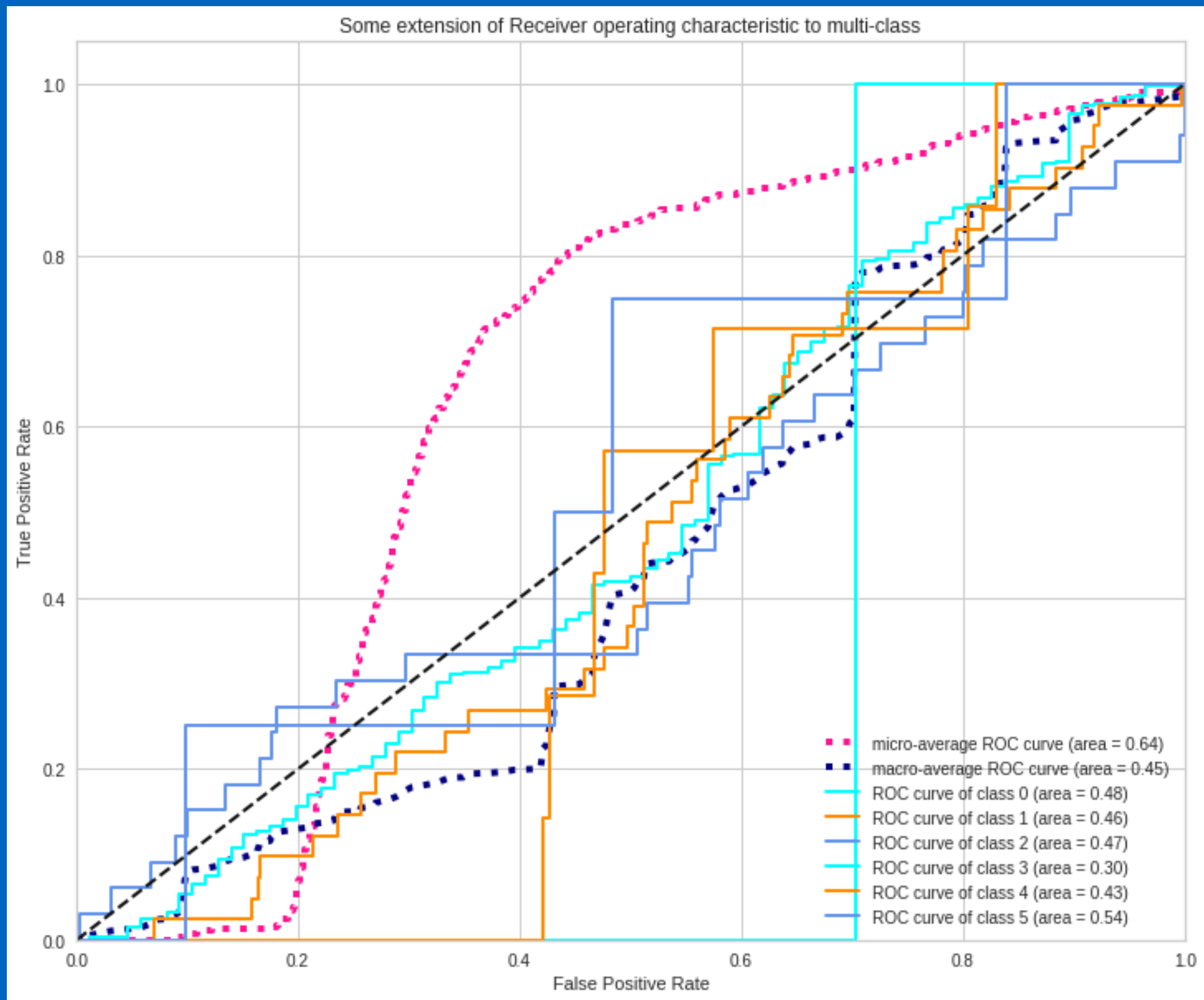
Decision Trees Would Not Appear to be as Helpful

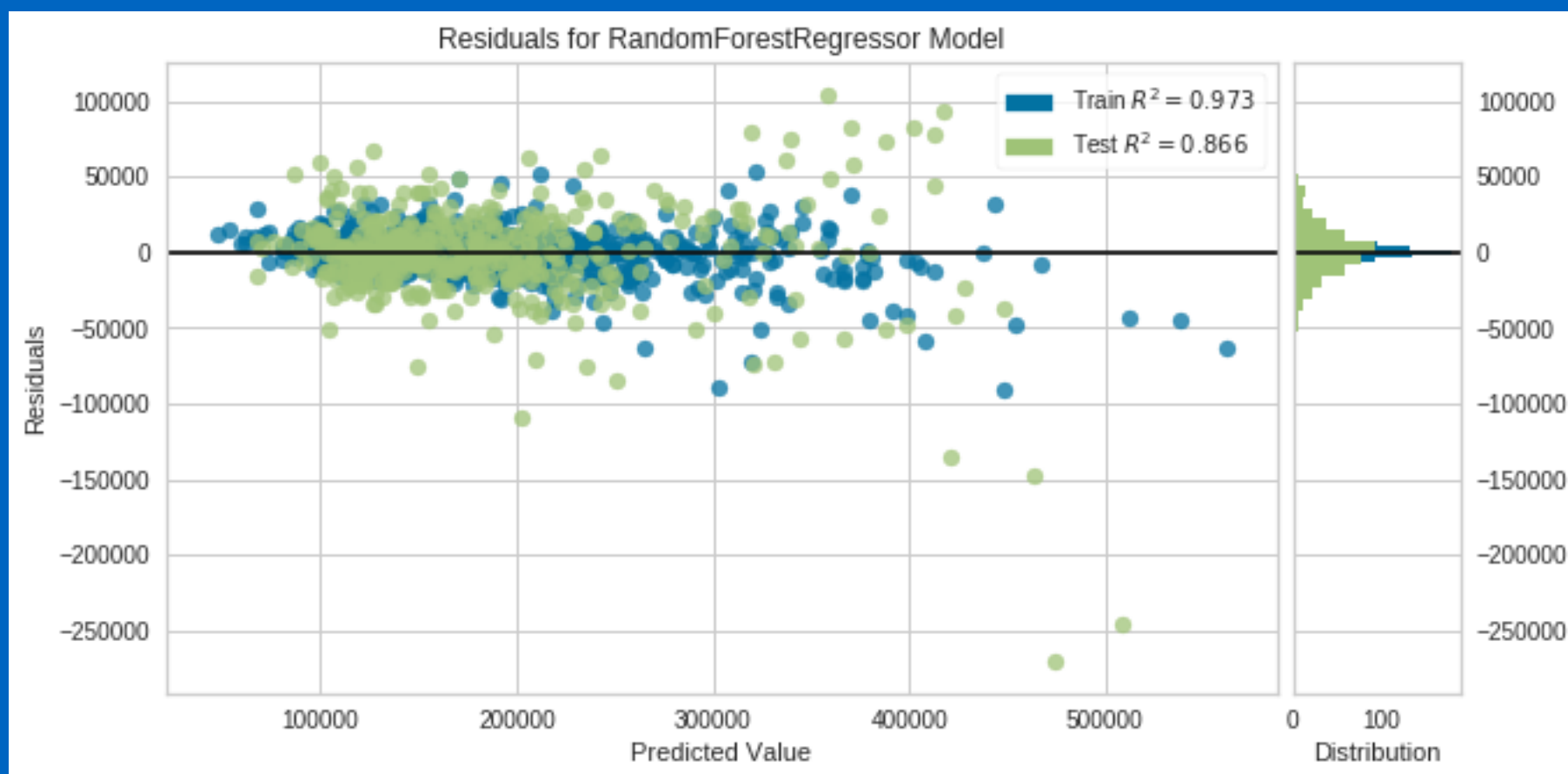
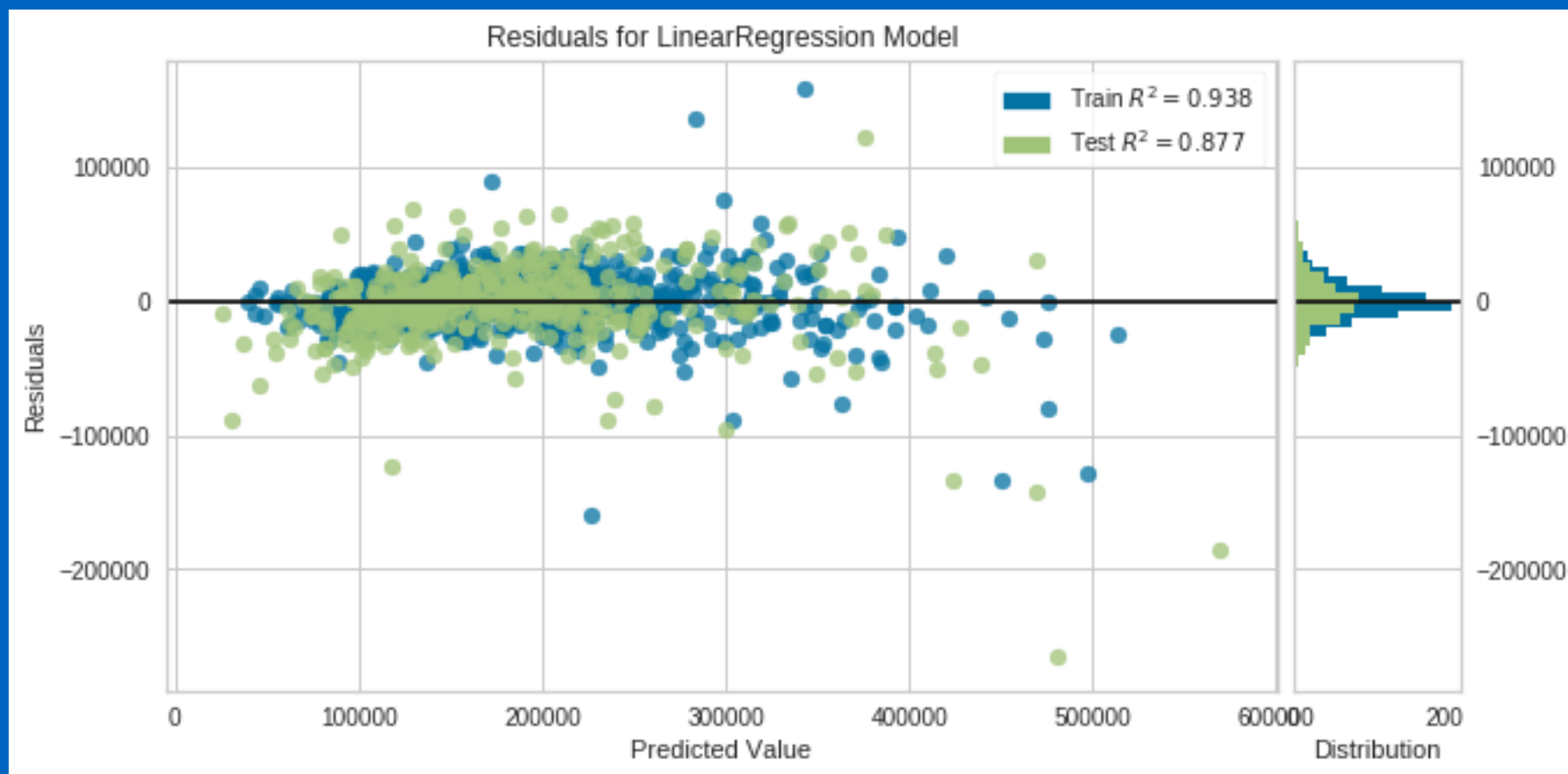
# Part 2/5

## Room for Improvement.



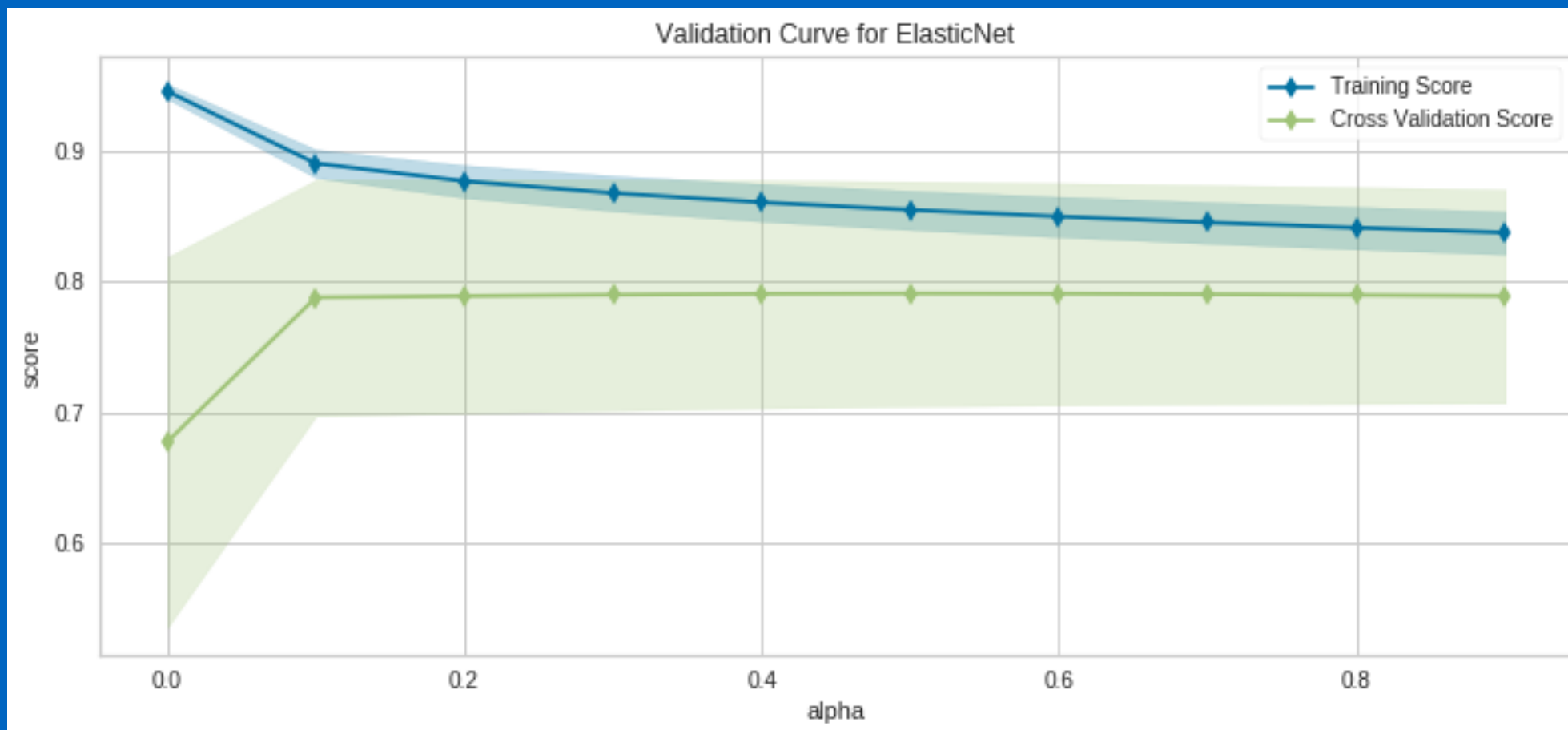
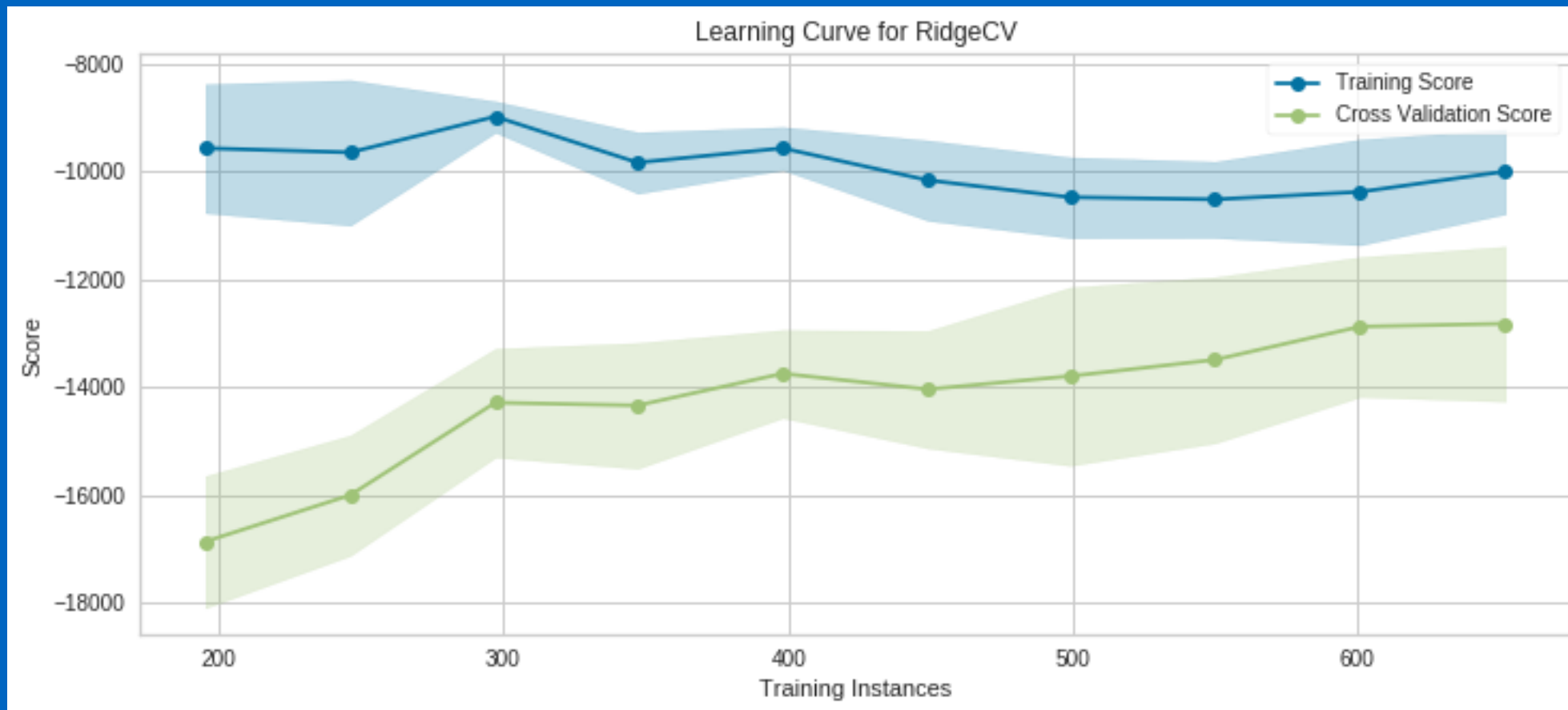






## Part 4/5

**ML is so easy. I am getting 99.5% accuracy. Wait..  
something is wrong !**



## Part 5/5

**Explain me what you learnt !**



