

L D S S A

SLU13 - Bias-Variance trade-off & Model Selection

December 3rd, 2023



1. Introduction

Context

SLU01 to SLU06



Data



SLU07 to SLU12



Models



SLU13 to SLU17

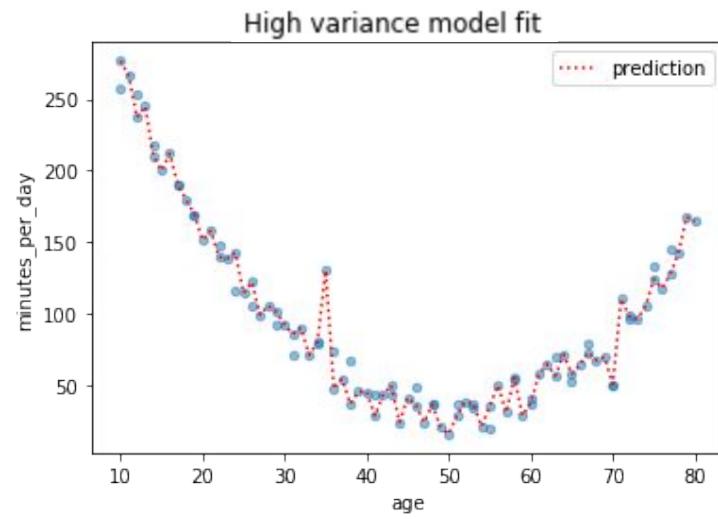
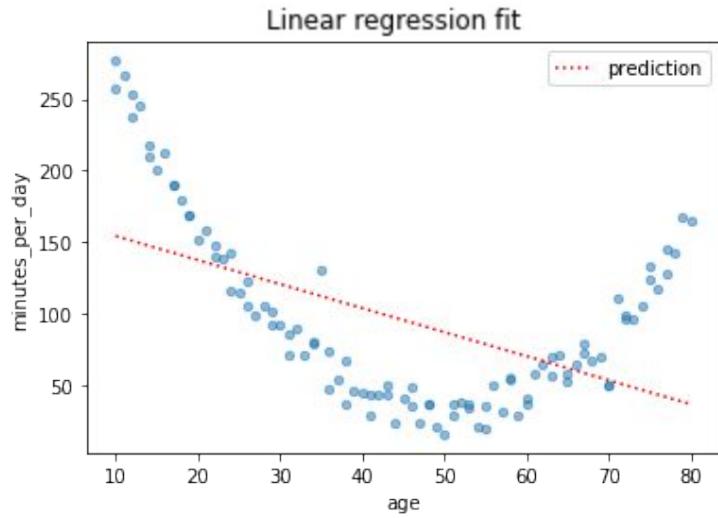


Model Improvement



Motivation

Which model is better?



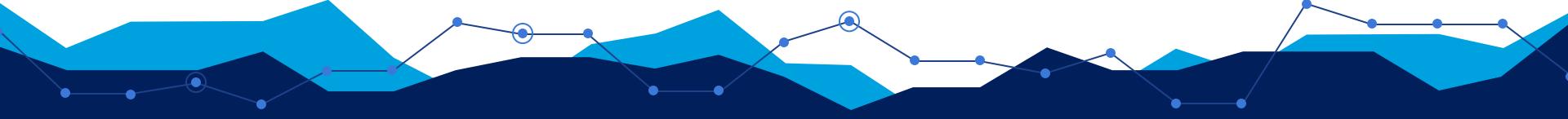
Motivation

Answer: They are both bad!



In this SLU we will try to give you the tools to answer the following questions:

- How do we evaluate our models?
 - (Hint: it's more than just maximizing R^2 , accuracy, etc.)
- How do we make sure our models generalize well?



Overview

- **Bias-variance trade-off**
 - a. Overfitting and underfitting
- **Model Selection**
 - a. Train-test split
 - b. Train-validation-test split
 - c. Cross-validation
- **Learning curves**





2. Topic Explanation

Bias-variance Tradeoff



Model Selection



Learning Curves



Bias-variance Tradeoff

- What is our main goal?
 - To approximate the one true, general target function
- Bias-variance decomposition is a way of analyzing the generalization error

Bias:

- *Always learning the same wrong thing*
- Results from simplistic assumptions and a lack of model adaptability



Variance:

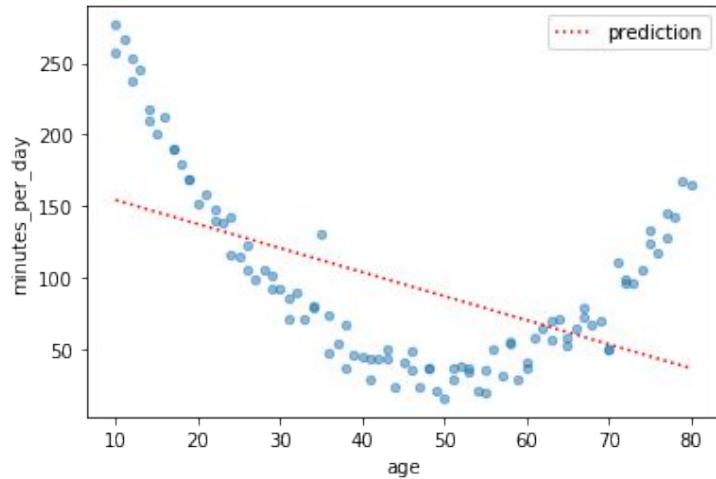
- *Learning random things*
- Results from excessive flexibility, a lack of underlying structure to anchor predictions



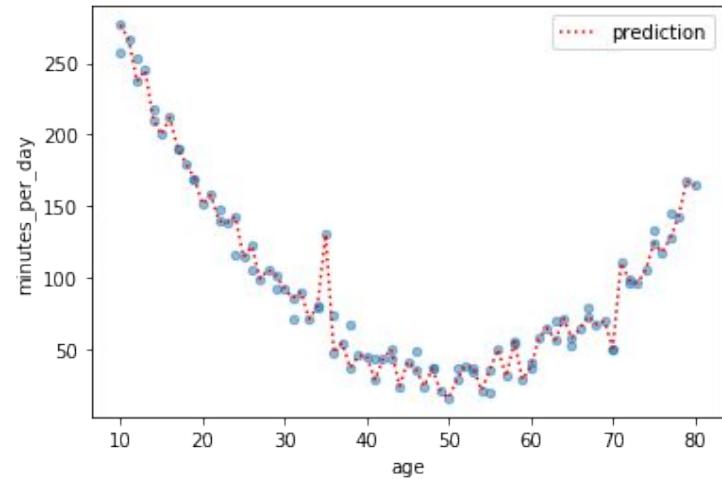
Bias-variance Tradeoff

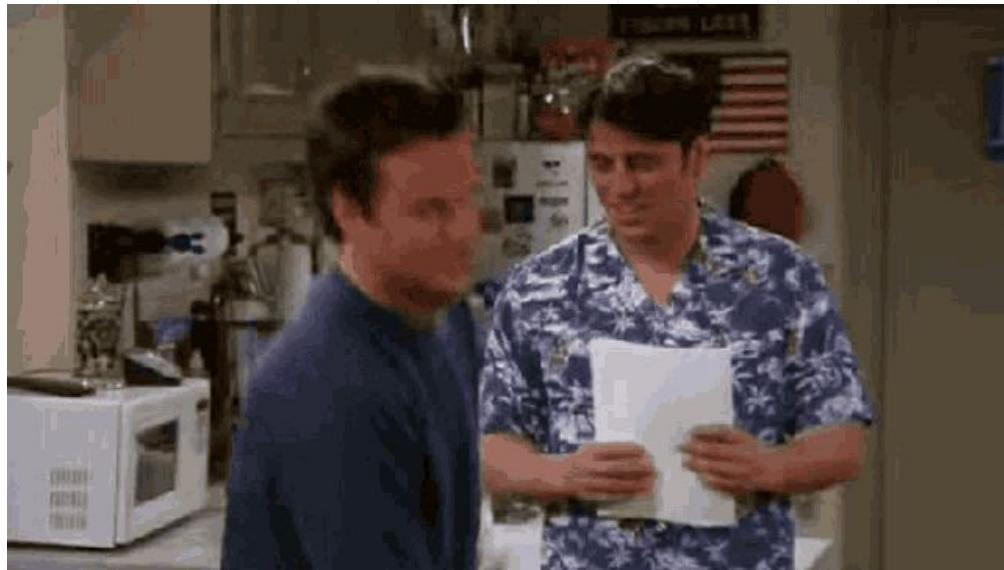
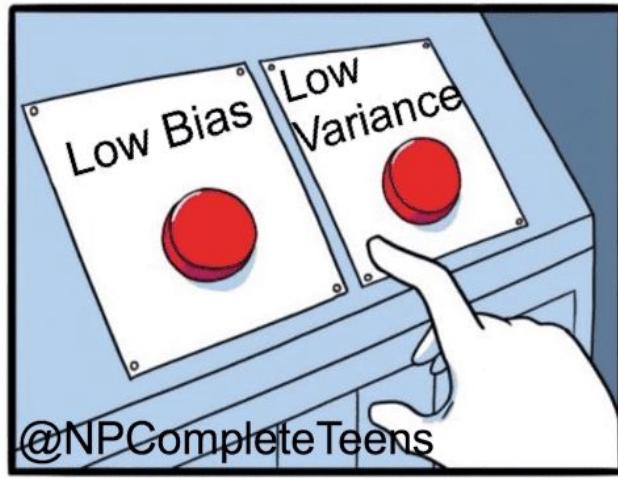
- Regression example:

High bias == Underfitting



High variance == Overfitting

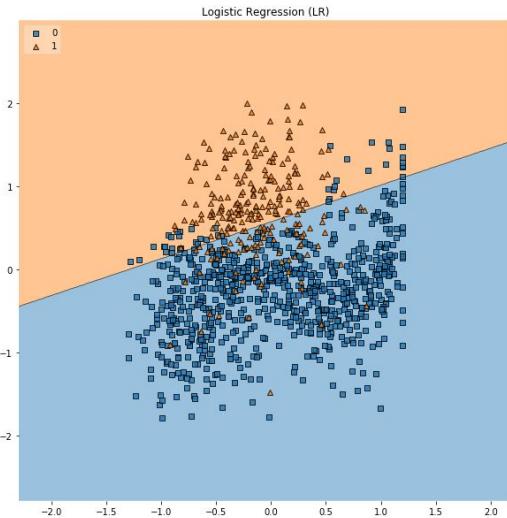




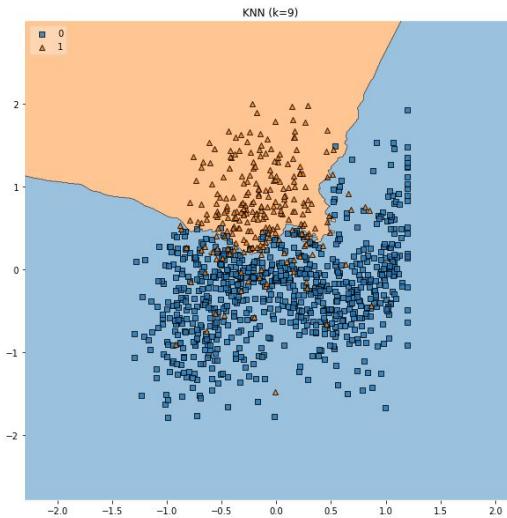
Bias-variance Tradeoff

- Classification example:

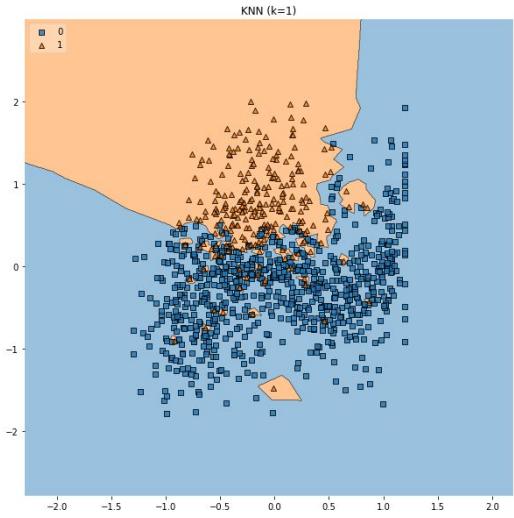
High bias (*underfitting*)



Just right

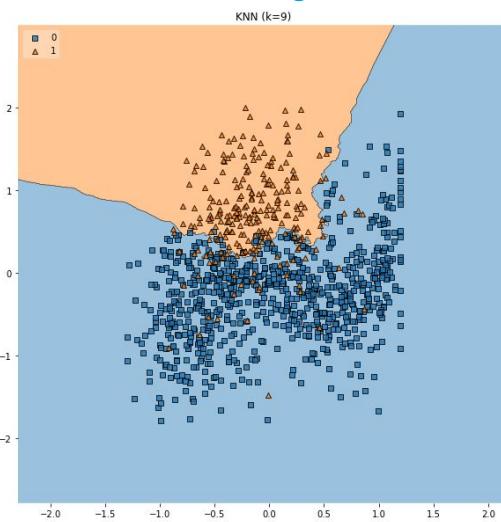


High variance (*overfitting*)

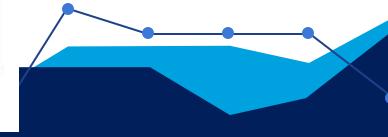
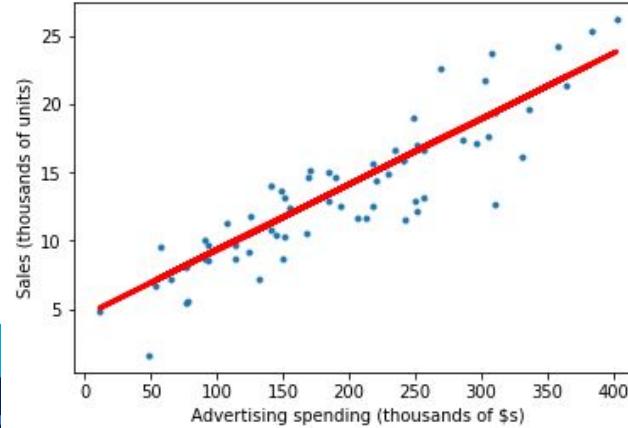


Bias-variance Tradeoff

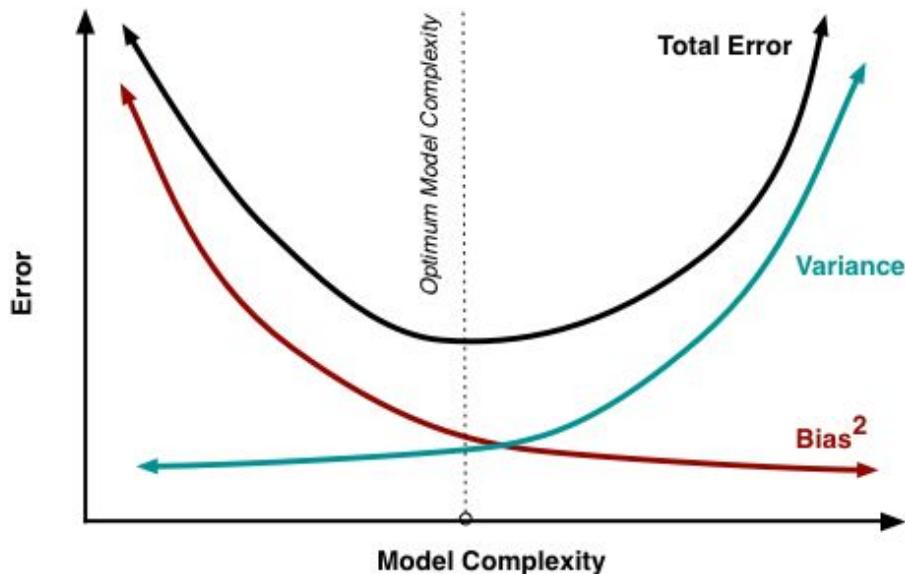
Just right



- But the “just right” model still isn’t perfect!
- There will always be some **irreducible error**:
 - Error that cannot be eliminated by building good models
 - Noise and outliers will exist



Bias-variance Tradeoff

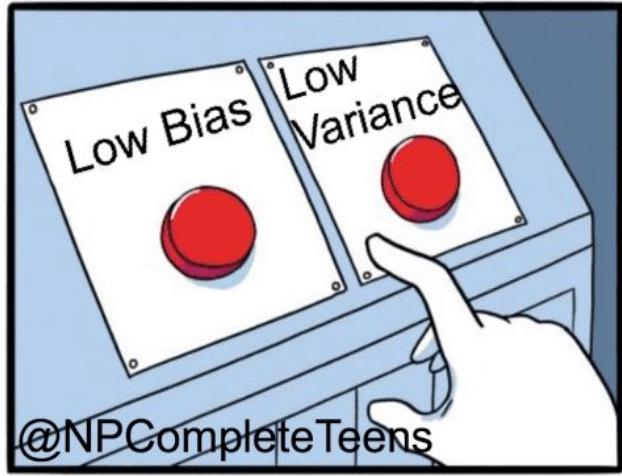


In theory, we reach the right level of complexity when the increase in bias is equivalent to the reduction in variance:

$$\frac{dBias}{dComplexity} = -\frac{dVariance}{dComplexity}$$

But there is no analytical way to find this optimum...

So we need to test our hypothesis!



Bias-variance Tradeoff



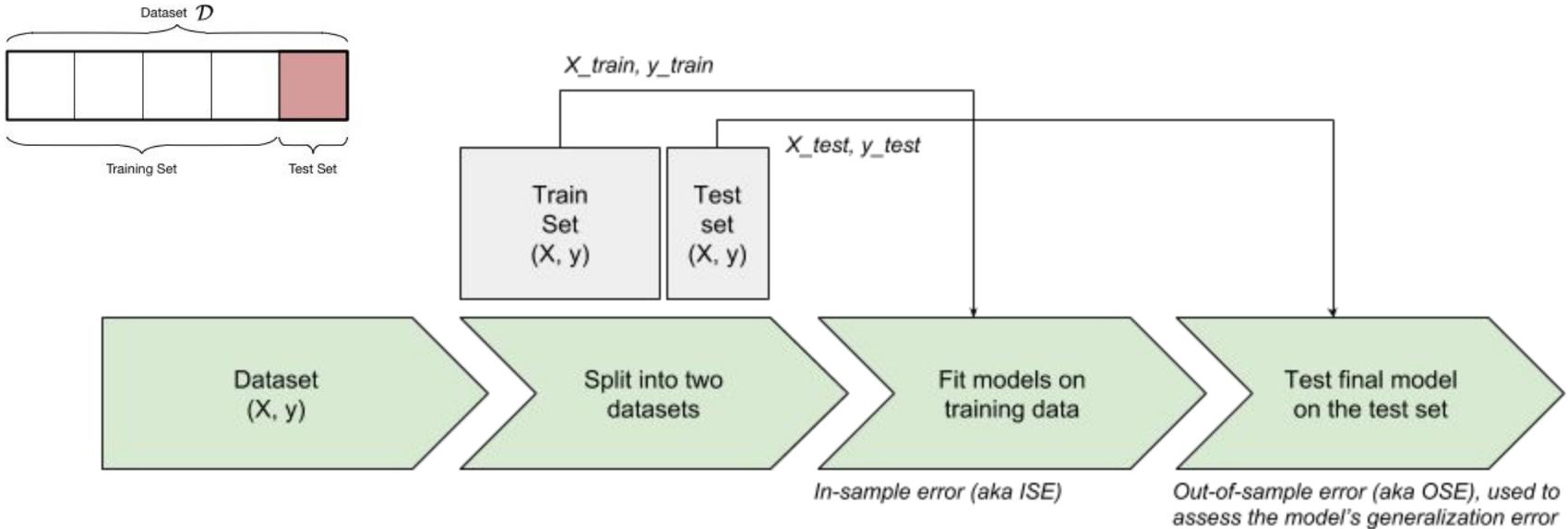
Model Selection



Learning Curves

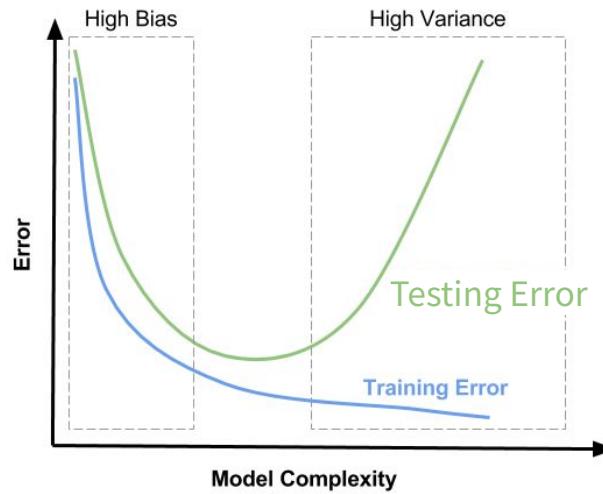


Train-Test Split

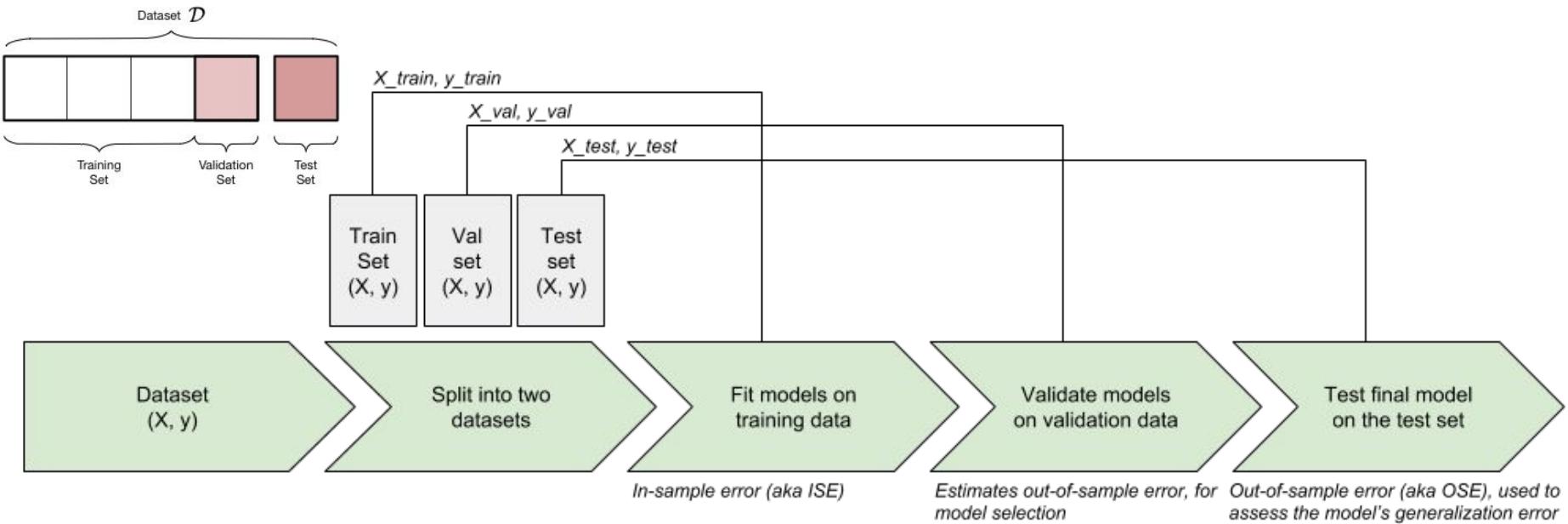


Train and Test Error

- Training error or in-sample error (ISE) measures how well our trained model performs on the training data
- Testing error or out-of-sample error (OSE) is how well our model performs on previously unseen data



Train-Validation-Test Split



Train-Validation-Test Split

Why should we have a validation set? Isn't it the same as the test set?

Validation set can help us decide make decisions regarding hyper-parameters.

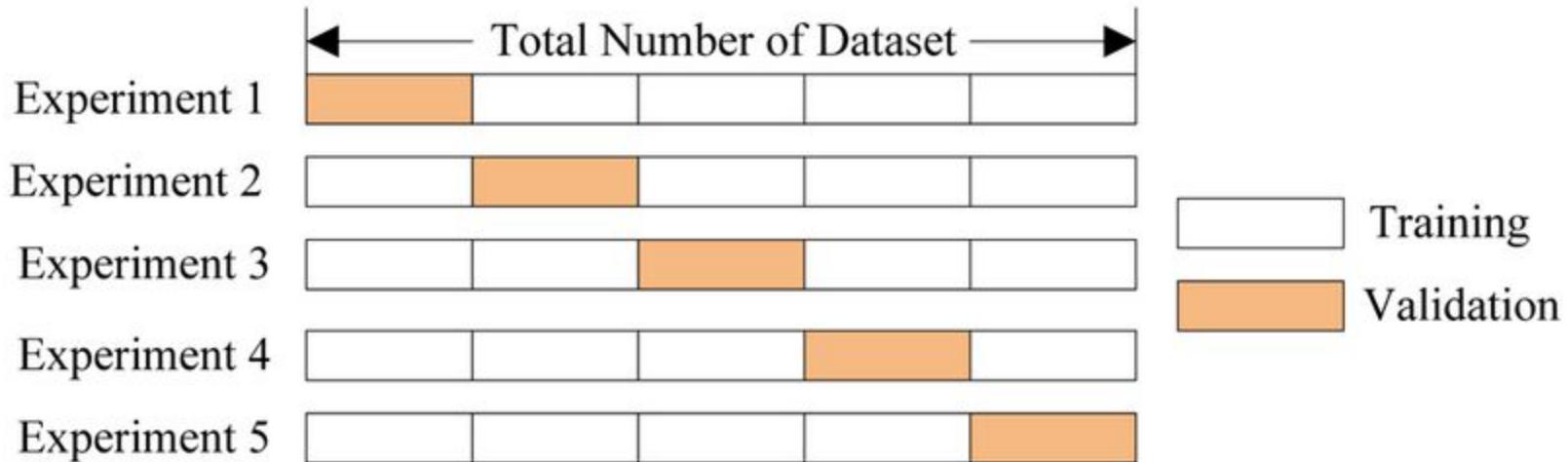
- Example: How many polynomial features to consider in our model?

$$\hat{y} = \beta_0 + \beta_1 x + \beta_2 x^2 + \beta_3 x^3 + \beta_4 x^4 + \beta_5 x^5$$

- Other examples:
 - Which learning rate should we use in iterative processes
 - How many iterations should we use to optimize our model
 - What should be the size of the Neural Network



Cross-Validation



- Measure the mean and variance of the **scores across each “fold”**
- Used to approximate the OSE **when we don’t have a lot of data**

Summary of Evaluation Methods

- **Train-test split**
 - Most basic way to evaluate OSE
- **Train-validation-test split**
 - Needs sufficient data
 - Essential for hyperparameter tuning
- **Cross-validation**
 - Useful when you don't have a lot of data to approximate OSE
 - Measure the mean and variance of the scores across each fold



Bias-variance Tradeoff



Model Selection



Learning Curves



How much data should we use?

- The obvious starting point about which data to use is... "why not just all of it?"



	PassengerId	Survived	Pclass	Name	Sex	Age	SibSp	Parch	Ticket	Fare	Cabin	Embarked
0	1	0	3	Braund, Mr. Owen Harris	male	22.0	1	0	A/5 21171	7.2500	NaN	S
1	2	1	1	Cumings, Mrs. John Bradley (Florence Briggs Th... Heikkinen, Miss. Laina	female	38.0	1	0	PC 17599 STON/O2. 3101282	71.2833 7.9250	C85 NaN	C S
2	3	1	3	Futrelle, Mrs. Jacques Heath (Lily May Peel)	female	35.0	1	0	113803	53.1000	C123	S
3	4	1	1	Allen, Mr. William Henry	male	35.0	0	0	373450	8.0500	NaN	S
4	5	0	3	Moran, Mr. James	male	NaN	0	0	330877	8.4583	NaN	Q
5	6	0	3	McCarthy, Mr. Timothy J	male	54.0	0	0	17463	51.8625	E46	S
6	7	0	1	Palsson, Master. Gosta Leonard	male	2.0	3	1	349909	21.0750	NaN	S
7	8	0	3	Johnson, Mrs. Oscar W (Elisabeth Vilhelmina Berg)	female	27.0	0	2	347742	11.1333	NaN	S
8	9	1	3	Nasser, Mrs. Nicholas (Adele Achem)	female	14.0	1	0	237736	30.0708	NaN	C
9	10	1	2	Sandstrom, Miss. Marguerite Rut	female	4.0	1	1	PP 9549	16.7000	G6	S
10	11	1	3	Bonnell, Miss. Elizabeth	female	58.0	0	0	113783	26.5500	C103	S
11	12	1	1	Saundercock, Mr. William Henry	male	20.0	0	0	A/5. 2151	8.0500	NaN	S
12	13	0	3	Andersson, Mr. Anders Johan	male	39.0	1	5	347082	31.2750	NaN	S
13	14	0	3	Vestrom, Miss. Hulda Amanda Adolfina	female	14.0	0	0	350406	7.8542	NaN	S
14	15	0	3									





slow



Expensive



Insufficient



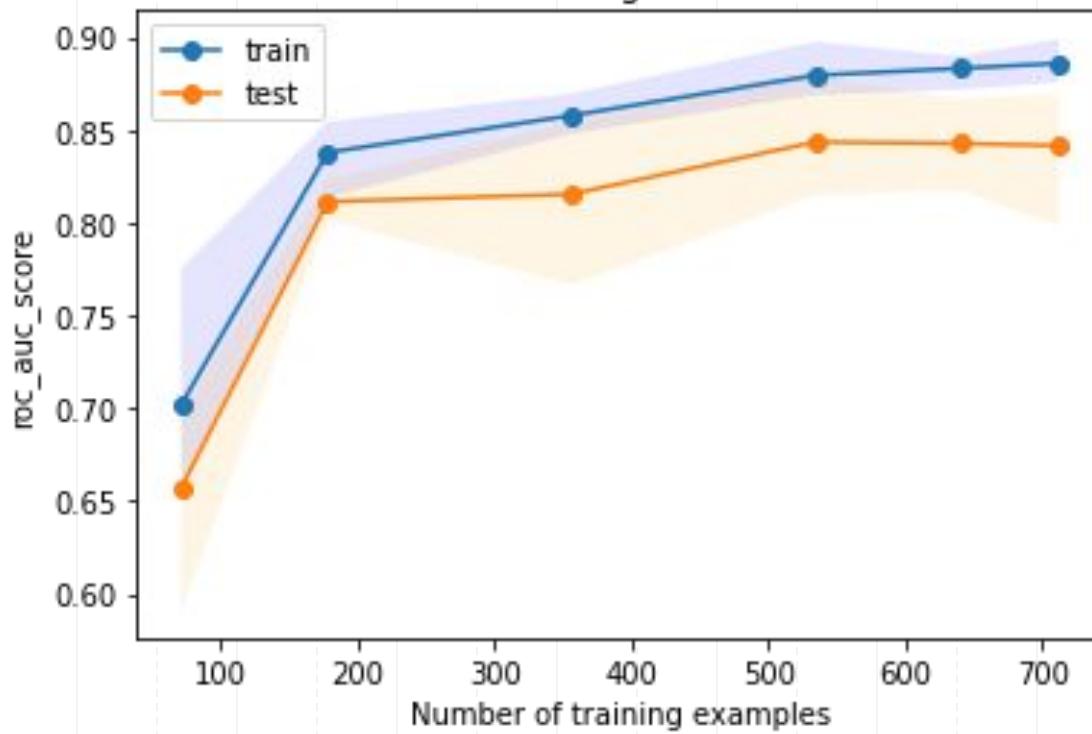
Biased



How much data should we use?

- The obvious starting point about which data to use is... "why not just all of it?"
 - Using too much data makes modelling unnecessarily **slow**
 - The data you have for training may have sample **bias**
 - Data can be **expensive** to come by - Would more data really help improving the model?
 - All the data may **not** be **enough**
- To understand how our model is learning and how much data it needs:
 - `sklearn.model_selection.learning_curve`.

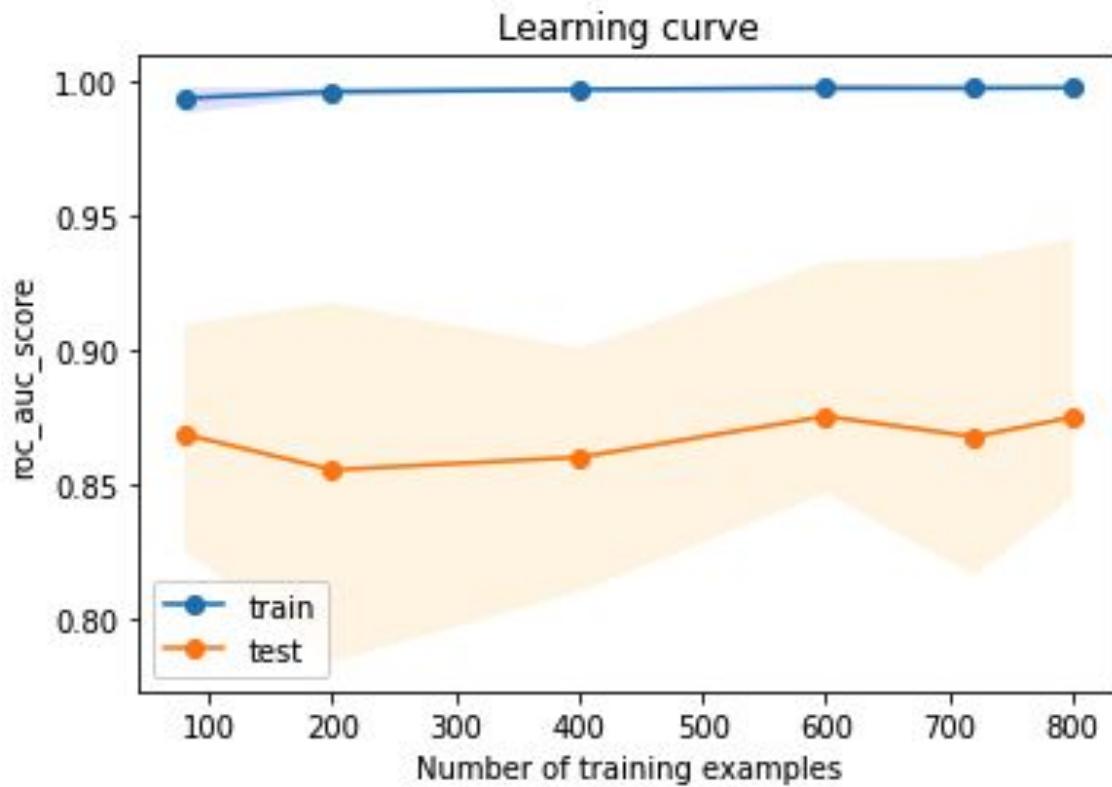
Learning curve

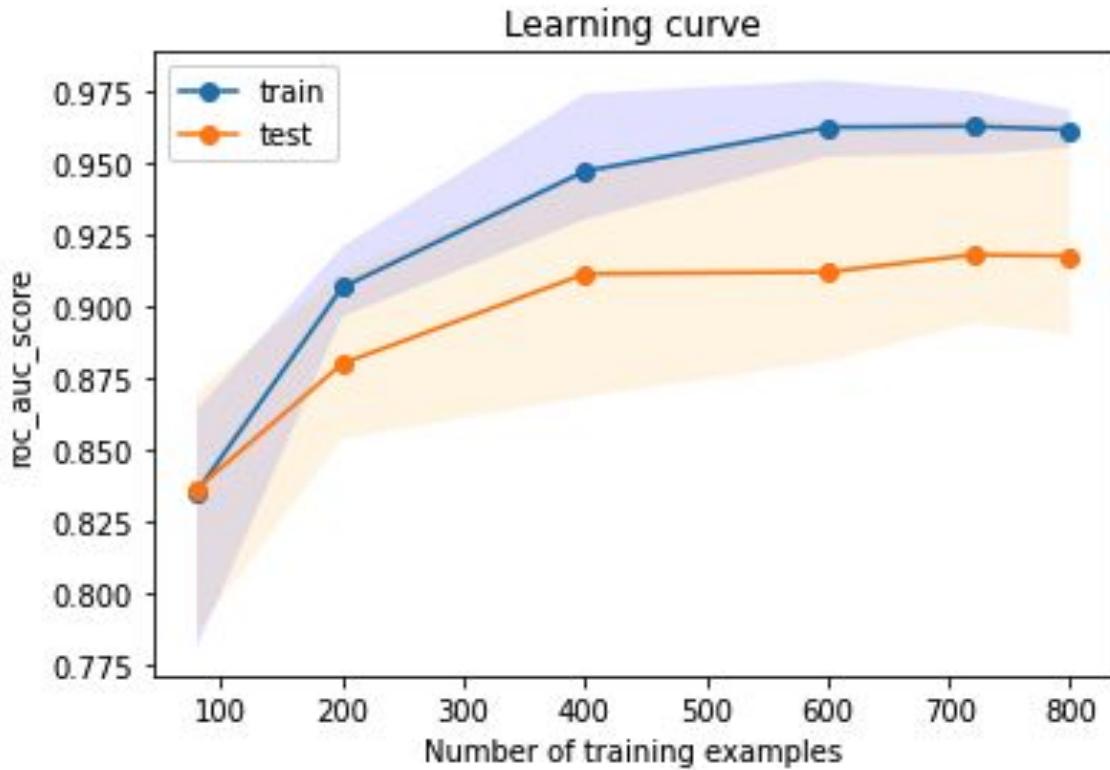


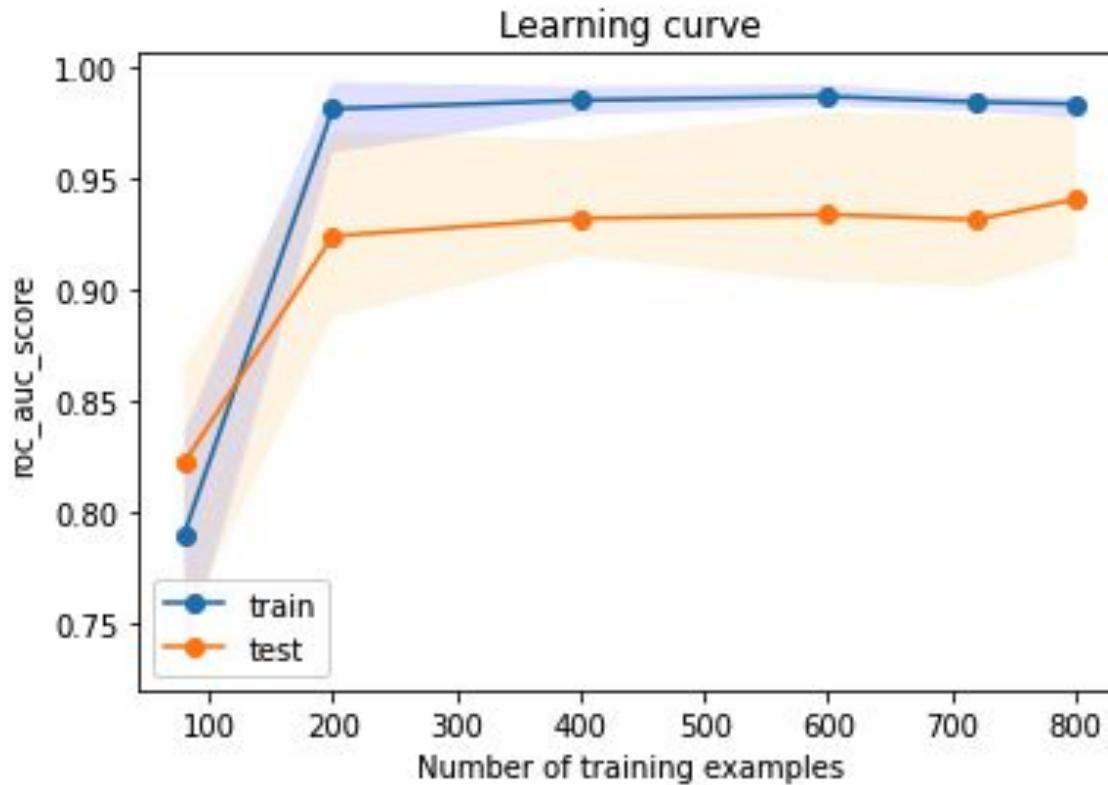
How much data should we use?

- Learning curves show you how the performance changes with dataset size.
- Performs cross-validation on a metric you provide
- Model variance can also be seen in the learning curves









3. Recap

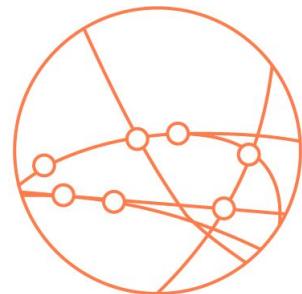
Recap

- There is a tradeoff between **bias and variance**
 - Bias = underfitting = model is always wrong
 - Variance = overfitting = model “memorizes” training data and does not generalize
 - Need a balance between the two
- We can estimate the generalization power of our models by **splitting our data**
 - Train or fit models to training set, evaluate on test set
 - We want both train and test error to be low
- Plot learning curves
 - Don’t use data just because it exists
 - Your model needs your brain more than it needs your CPU



4. Q&A





L D S S A

SLU14 - Model complexity and Overfitting



1. Introduction







PassengerId			Survived		Pclass		Name				Sex	Age	Class	Embarked				
			0	1	2	3	Braund, Mr. Owen Harris	Brickell, Mrs. Florence Briggs Th...	Heldman, Miss. Lula	Homey, Mr. John Bradley	Lucas, Mrs. (Lily May Peel)	McCartry, Mr. Timothy J.	Allen, Mr. William Henry	Moor, Mr. James	Watson, Mr. Charles W.	Ticket	Fare	Cabin
1	1	0	0	1	1	3	Mr. Owen Harris	Miss Florence Briggs Th...	Miss Lula Heldman	Mr. John Bradley Homey	Mrs. (Lily May Peel) Lucas	Mr. Timothy J. McCartry	Mr. William Henry Allen	Mr. James Moor	Mr. Charles W. Watson	3101282	7.2500	NAN
2	2	1	1	1	1	3	Cumings, Mrs. John Bradley	Brickell, Mrs. Florence Briggs Th...	Heldman, Miss. Lula	Homey, Mr. John Bradley	Lucas, Mrs. (Lily May Peel)	McCartry, Mr. Timothy J.	Allen, Mr. William Henry	Moor, Mr. James	Watson, Mr. Charles W.	3101282	71.2500	C85
3	3	1	1	1	1	3	Brickell, Mrs. Florence Briggs Th...	Heldman, Miss. Lula	Homey, Mr. John Bradley	Lucas, Mrs. (Lily May Peel)	McCartry, Mr. Timothy J.	Allen, Mr. William Henry	Moor, Mr. James	Watson, Mr. Charles W.	3101282	7.2500	NAN	
4	4	1	1	1	1	3	Homey, Mr. John Bradley	Lucas, Mrs. (Lily May Peel)	McCartry, Mr. Timothy J.	Allen, Mr. William Henry	Moor, Mr. James	Watson, Mr. Charles W.	3101282	7.2500	C123			
5	5	1	1	1	1	3	Lucas, Mrs. (Lily May Peel)	McCartry, Mr. Timothy J.	Allen, Mr. William Henry	Moor, Mr. James	Watson, Mr. Charles W.	3101282	7.2500	NAN				
6	6	0	0	0	0	3	Futrelle, Mrs. Jacques Heath	Allen, Mr. William Henry	McCarthy, Mr. Timothy J.	Moor, Mr. James	Watson, Mr. Charles W.	3101282	7.2500	NAN				
7	7	0	0	0	0	3	Allen, Mr. William Henry	McCarthy, Mr. Timothy J.	Moor, Mr. James	Watson, Mr. Charles W.	3101282	7.2500	NAN					
8	8	0	0	0	0	3	McCarthy, Mr. Timothy J.	Moor, Mr. James	Watson, Mr. Charles W.	3101282	7.2500	NAN						
9	9	0	0	0	0	3	Moor, Mr. James	Watson, Mr. Charles W.	3101282	7.2500	NAN							
10	10	1	1	1	1	3	Watson, Mr. Charles W.	3101282	7.2500	NAN								
11	11	1	1	2	1	3	Johnson, Mrs. Oscar W. (Elsie Homem Bengt)	Sandstrom, Miss. Margaretta Rut	Donovan, Miss. Elizabeth	Saundercock, Mr. William Henry	Anderson, Mr. Anders John	Westrom, Mrs. Hulda Amanda Adelina	3101282	7.2500	NAN			
12	12	1	1	3	1	3	Elsie Homem Bengt	Margaretta Rut	Elizabeth	William Henry	Anders John	Hulda Amanda Adelina	3101282	7.2500	NAN			
13	13	1	1	1	1	3	Margaretta Rut	Elizabeth	William Henry	Anders John	Hulda Amanda Adelina	3101282	7.2500	NAN				
14	14	0	0	0	0	3	Elizabeth	William Henry	Anders John	Hulda Amanda Adelina	3101282	7.2500	NAN					
15	15	0	0	0	0	3	William Henry	Anders John	Hulda Amanda Adelina	3101282	7.2500	NAN						



```
1 pd.get_dummies(df['Name']).head(3)
```

	Abbing, Mr. Anthony	Abbott, Mr. Rossmore Edward	Abbott, Mrs. Stanton (Rosa Hunt)	Abelson, Mr. Samuel (Hannah Wizosky)	Abelson, Mrs. Samuel (Leah Rosen)	Adahl, Mr. Mauritz Nils Martin	Adams, Mr. John (Johanna Persdotter Larsson)	Ahlin, Mrs. Johan ("Mrs Harbeck")	Aks, Mrs. Sam (Leah Rosen)	Albimona, Mr. Nassef Cassem	Yrois, Miss. Henriette ("Mrs Harbeck")	Zabour, Miss. Hilene	Zabour, Miss. Thamine	Zimmerman, Mr. Leo	I
0	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0
1	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0
2	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0

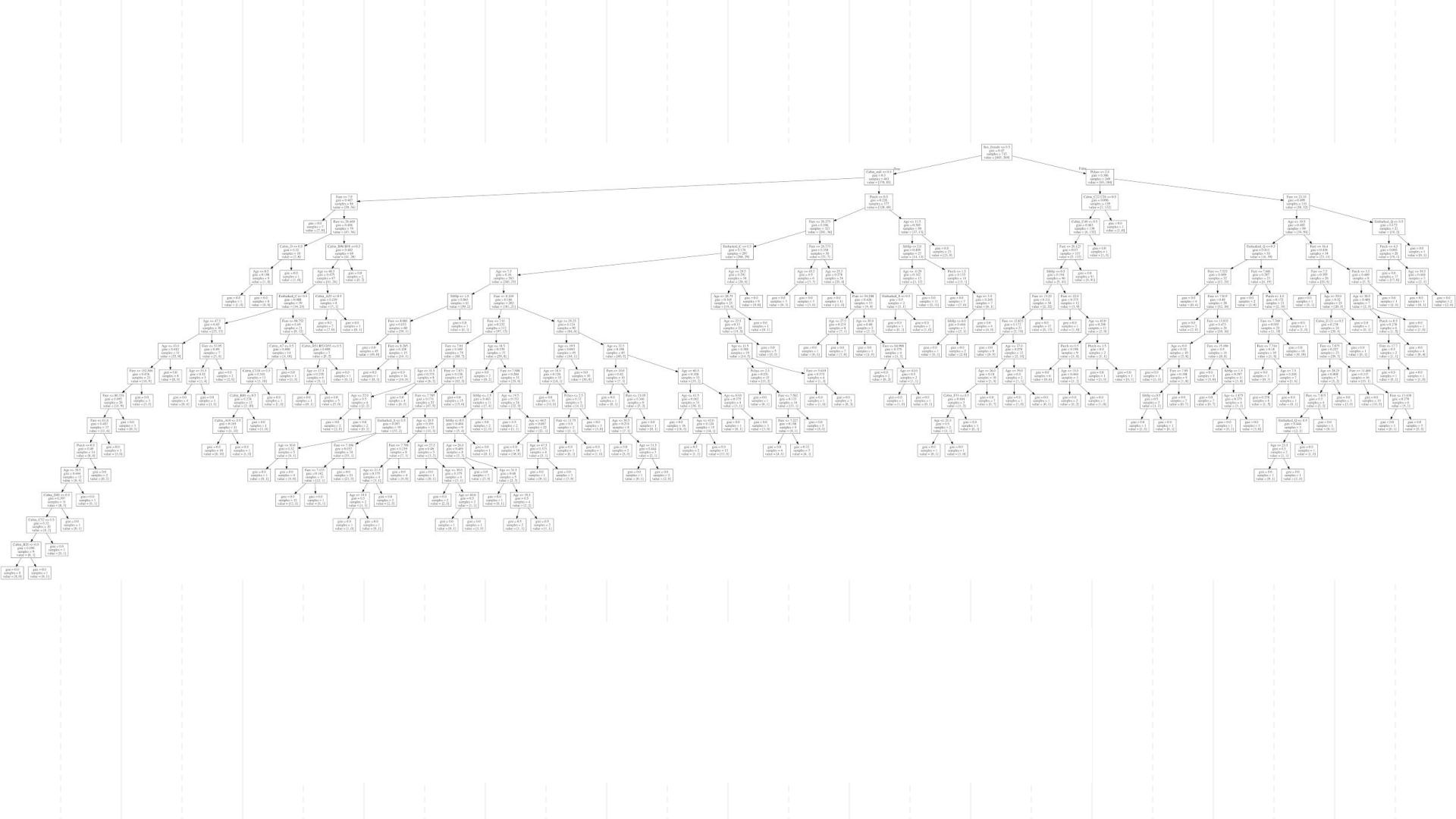
Pasengerid	Survived	Pclass	Name	Sex	Age	SibSp	Parch	Ticket	Fare	Cabin	Embarked
1	1	1	Braund, Mr. Owen Harris	male	30.0	1	0	312378	7.2500	NAN	S
2	2	1	Cummings, Mrs. John Bradley (Florence Th...)	female	30.0	1	0	PC 17599	71.3333	C85	C
3	3	1	Eaton, Mrs. J. Frank (Florence Th...	female	26.0	0	0	3101282	7.8500	NAN	S
4	4	1	Gowrisankaran, Miss. Lakshmi	female	30.0	0	0	113803	53.1000	C123	S
5	5	1	Futrelle, Mrs. Jacques Heath (Lily May Peel)	female	30.0	0	0	373400	6.0000	NAN	S
6	6	0	Allen, Mr. William Henry	male	35.0	0	0	330877	8.4500	NAN	S
7	7	0	McCarthy, Mr. Timothy J.	male	35.0	0	0	174023	51.8625	NAN	Q
8	8	0	Moor, Mr. James	male	35.0	0	0	349809	21.0250	E46	S
9	9	0	McGowan, Mrs. Timothy J.	female	2.0	0	0	347740	11.1333	NAN	S
10	10	1	Johnson, Mrs. Oscar W. (Elisabeth von Breisen)	female	27.0	0	0	237736	32.3708	NAN	S
11	11	1	Palsson, Master. Gustaf Leonardsson	male	30.0	0	0	PP 9542	16.7000	C	S
12	12	1	Nasser, Mrs. Nadezhda (Krstina Berg)	female	30.0	0	0	113783	26.5500	0103	S
12	13	1	Sandström, Miss. Margaretta Rut	female	14.0	0	0	347086	31.3750	NAN	S
14	14	0	Bonnell, Miss. Elizabeth	female	4.0	0	0	350406	7.8542	NAN	S
15	0	3	Southercock, Mr. Charles W.	male	58.0	0	0				
			Anderson, Mr. Anders John	male	20.0	0	0				
			Westrom, Miss. Hulda Amanda Adelina	female	30.0	1	0				
					14.0	0	0				

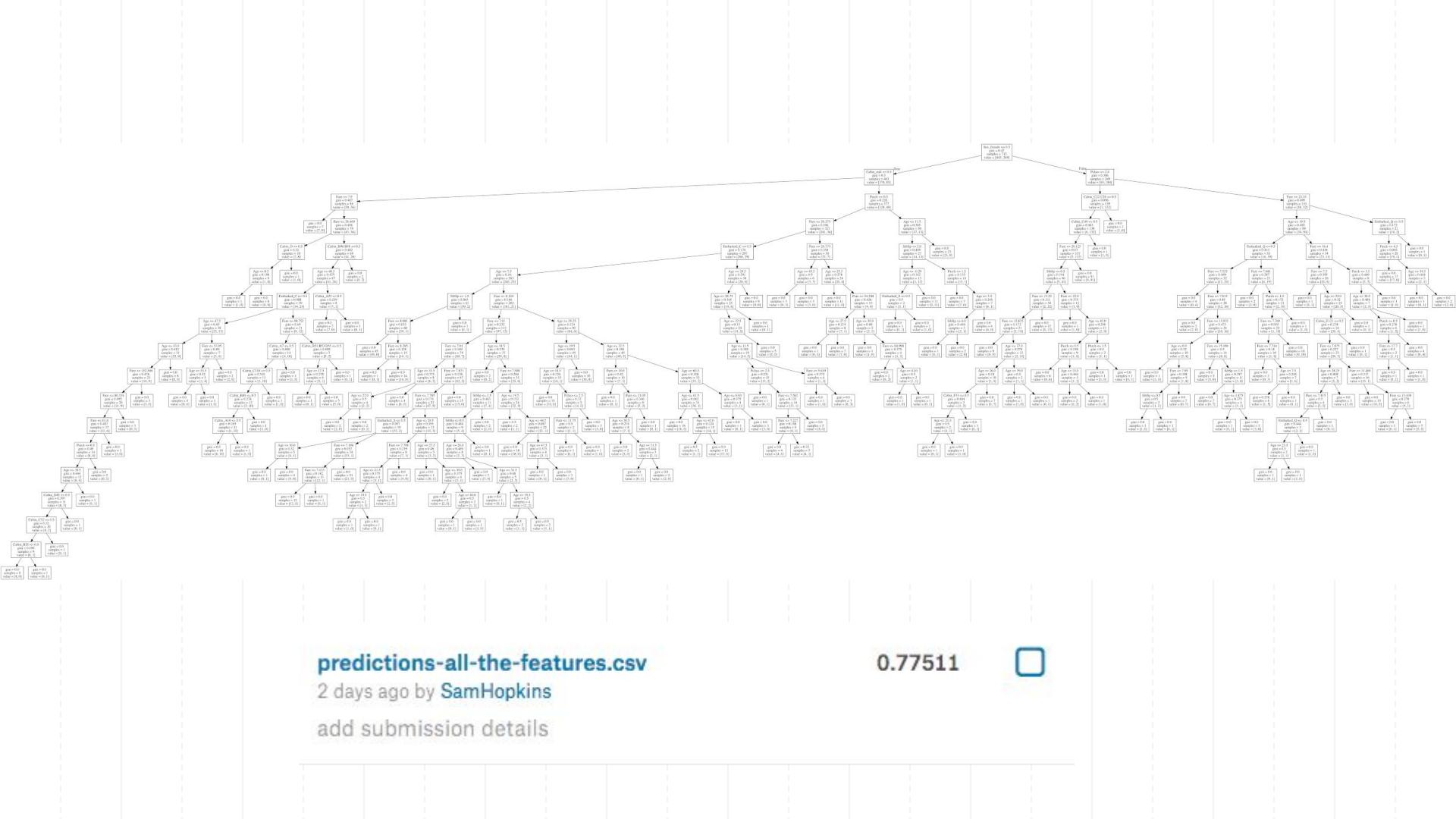


```
1 pd.get_dummies(df[ 'Name' ]).head(3)
```

Abbing, Mr. Anthony	Abbott, Mr. Rossmore Edward	Abbott, Mrs. Stanton (Rosa Hunt)	Abelson, Mr. Samuel	Abelson, Mrs. (Hannah Wizosky)	Adahl, Mr. Mauritz Nils Martin	Adams, Mr. John	Ahlin, Mrs. Johan (Johanna Persdotter Larsson)	Aks, Mrs. Sam (Leah Rosen)	Albimona, Mr. Nassef Cassem	... Henriette ("Mrs Harbeck")	Yrois, Miss. Zabour, Miss. Hileni	Zabour, Miss. Thamine	Zabour, Miss. Thamine	Zimmerman, Mr. Leo
0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
1	0	0	0	0	0	0	0	0	0	0	0	0	0	0
2	0	0	0	0	0	0	0	0	0	0	0	0	0	0







predictions-all-the-features.csv

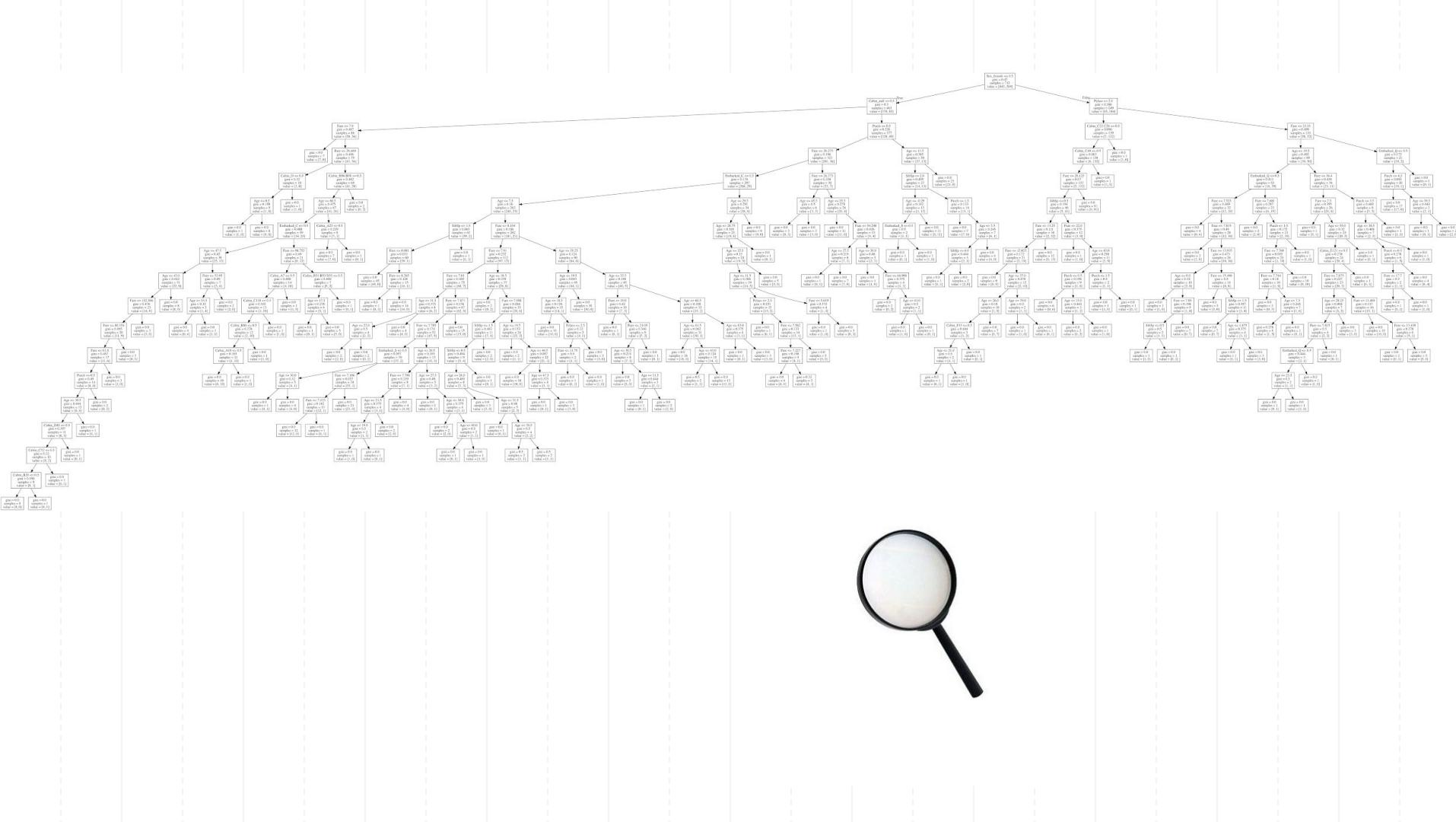
2 days ago by **SamHopkins**

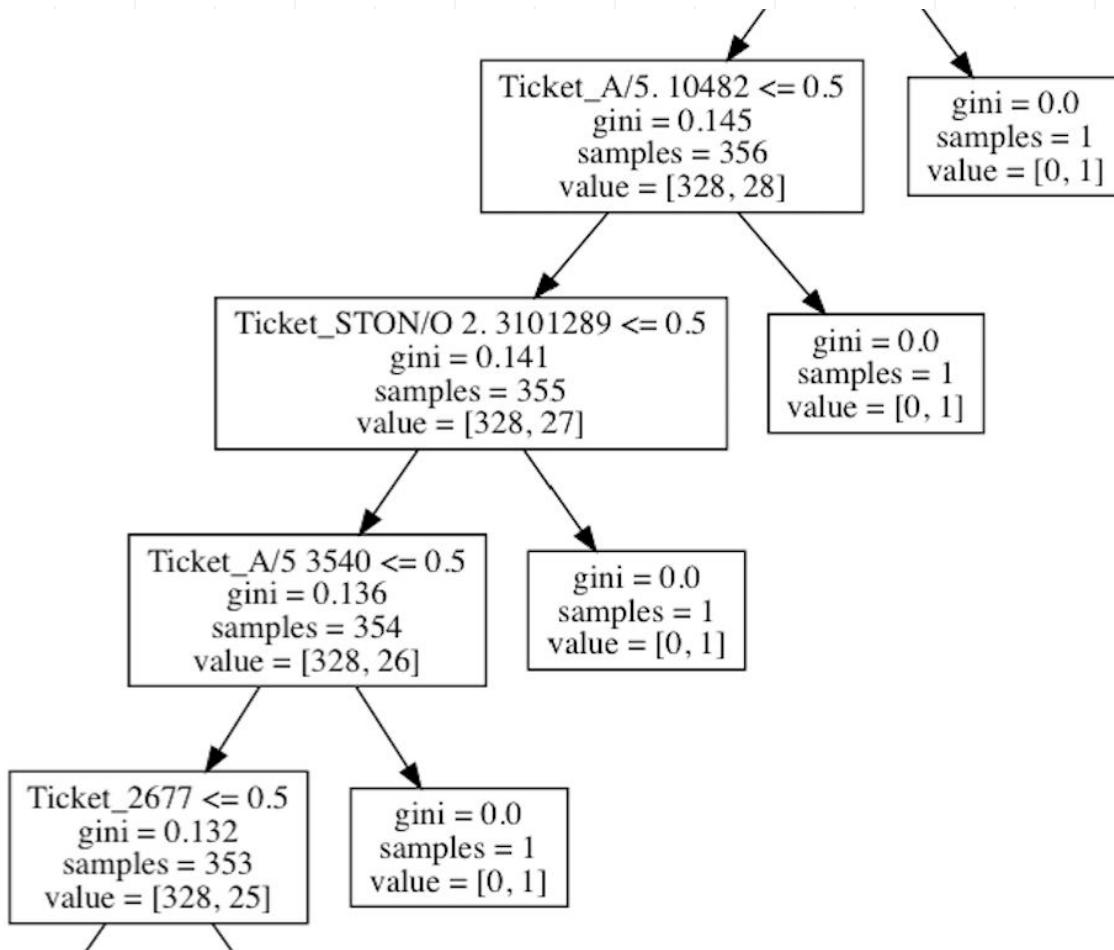
add submission details

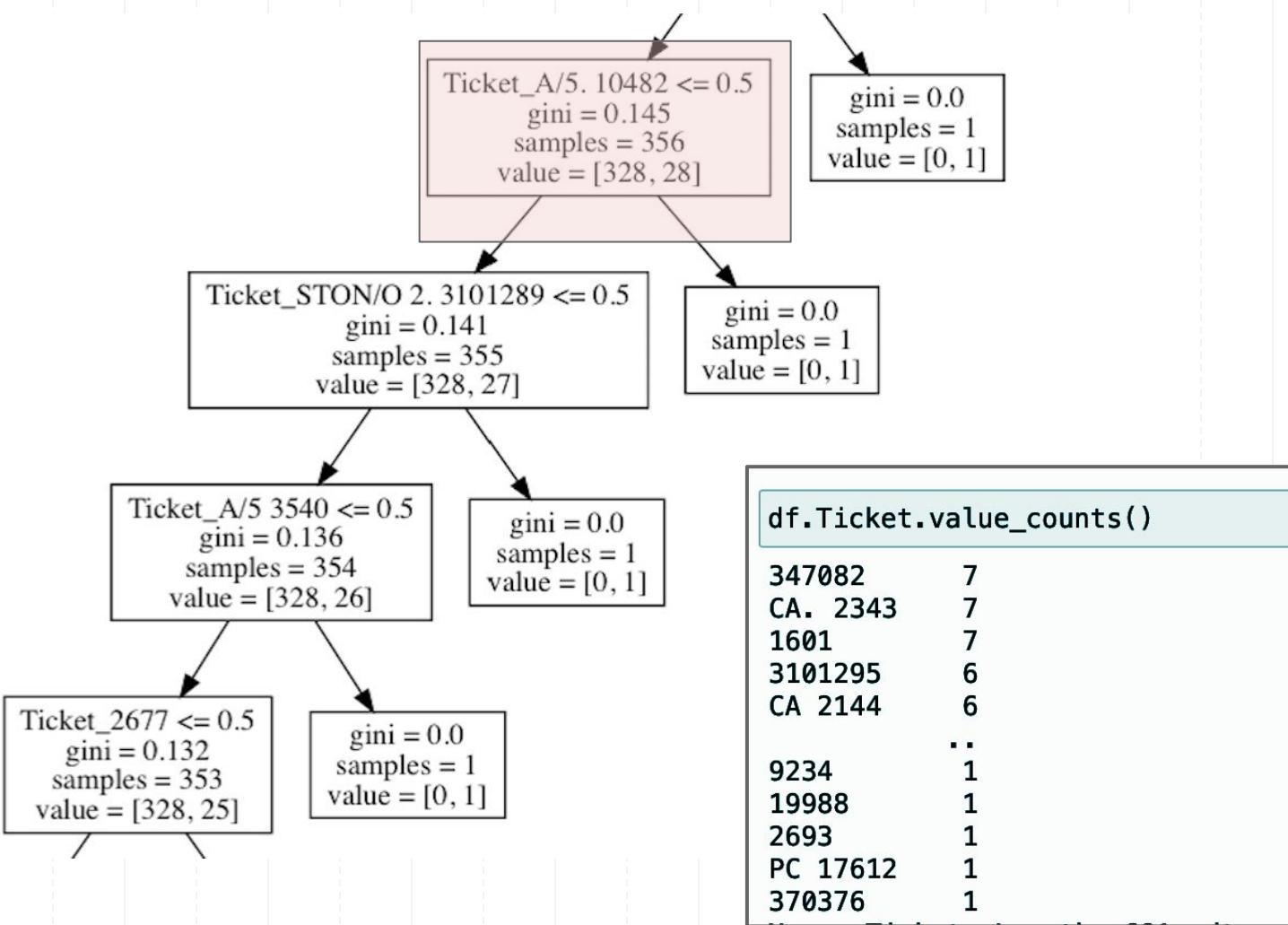
0.77511













WHOOPSIES...

Motivation

Using data blindly is not a wise approach...



Motivation

What issues can arise?

- Unexpected behavior
- Overfitting (poor future performance)
- Lack of interpretability
- Unnecessary model complexity (long training time)

Solution!

There are many, here we are going to look into:

- Classic method of:

‘Let’s look into the features and make sense of them.’

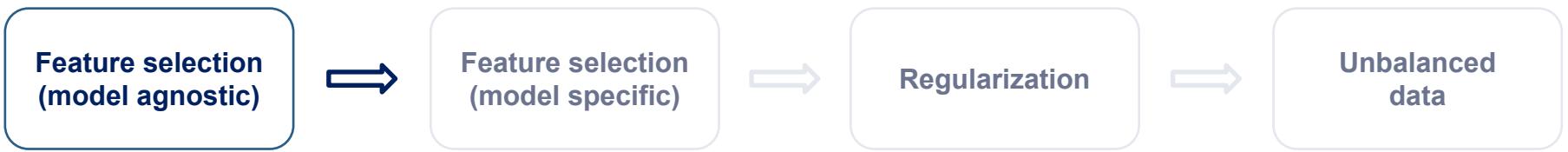
- **Regularisation** of Linear Models

Overview

- 1. Feature Selection**
 - a. Model agnostic
 - b. Model specific
- 2. Regularization**
 - a. Intuition and use-cases
 - b. Types of regularization
- 3. Unbalanced data**
 - a. Recognizing it
 - b. Avoiding some traps
 - c. Oversampling and undersampling

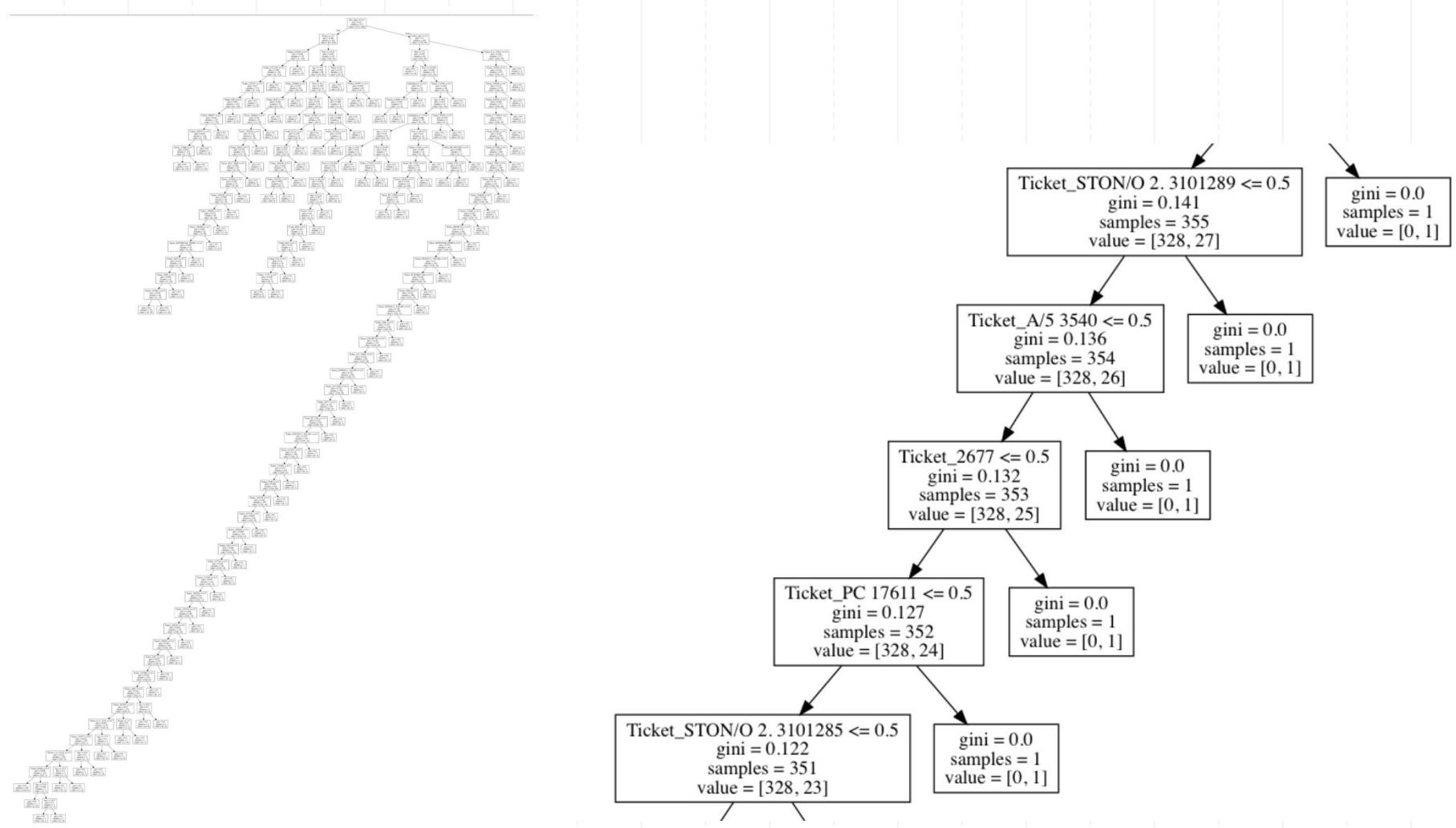


2. Topic Explanation



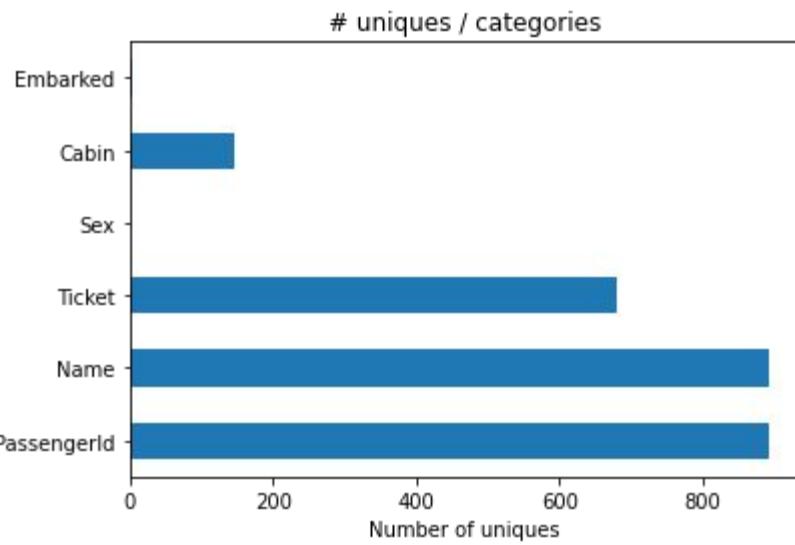
Exploring feature information

In Decision Trees, **looking at the first node** in the tree will most likely tell us what is one of the most important features to define the data



PassengerId	Survived	Pclass	Name	Sex	Age	SibSp	Parch	Ticket	Fare	Cabin	Embarked
160	0	3	Sage, Master. Thomas Henry	male	NaN	8	2	CA. 2343	69.55	NaN	S
181	0	3	Sage, Miss. Constance Gladys	female	NaN	8	2	CA. 2343	69.55	NaN	S
202	0	3	Sage, Mr. Frederick	male	NaN	8	2	CA. 2343	69.55	NaN	S
325	0	3	Sage, Mr. George John Jr	male	NaN	8	2	CA. 2343	69.55	NaN	S
793	0	3	Sage, Miss. Stella Anna	female	NaN	8	2	CA. 2343	69.55	NaN	S
847	0	3	Sage, Mr. Douglas Bullen	male	NaN	8	2	CA. 2343	69.55	NaN	S
864	0	3	Sage, Miss. Dorothy Edith "Dolly"	female	NaN	8	2	CA. 2343	69.55	NaN	S

Exploring feature information



- By doing some exploratory analysis we can see indeed that the Name, Ticket and PassengerId have too many values
- They won't add much information

Exploring feature information

You have a World model in your mind. **Use it.**

**Common
sense**

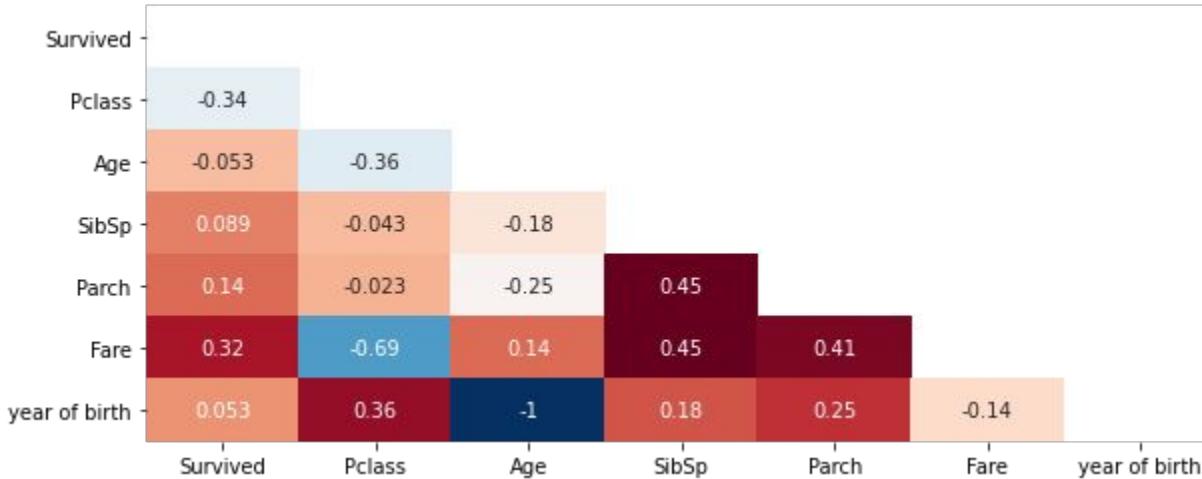
**Tree
visualizations**

Correlation

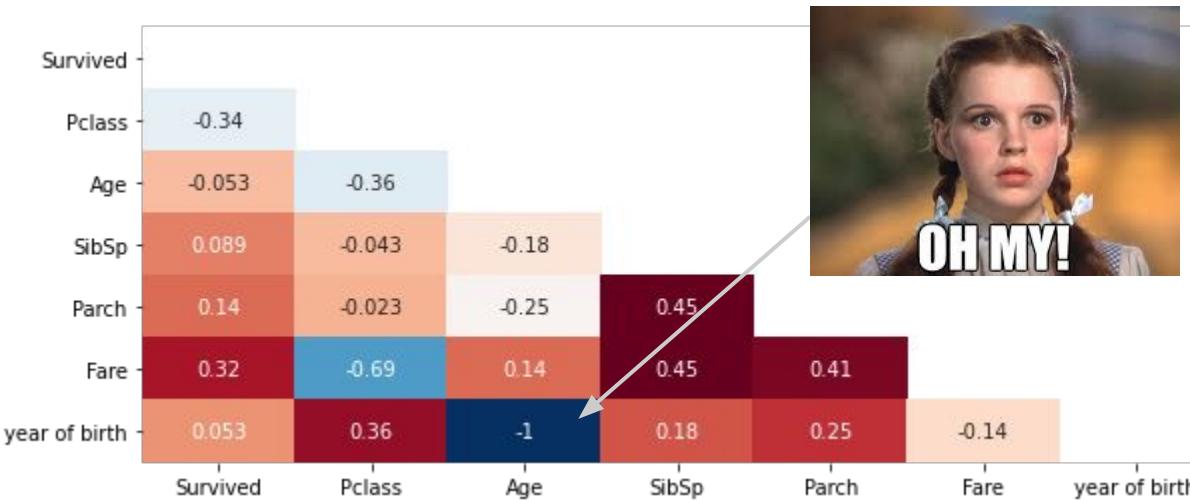
**Mutual
information**

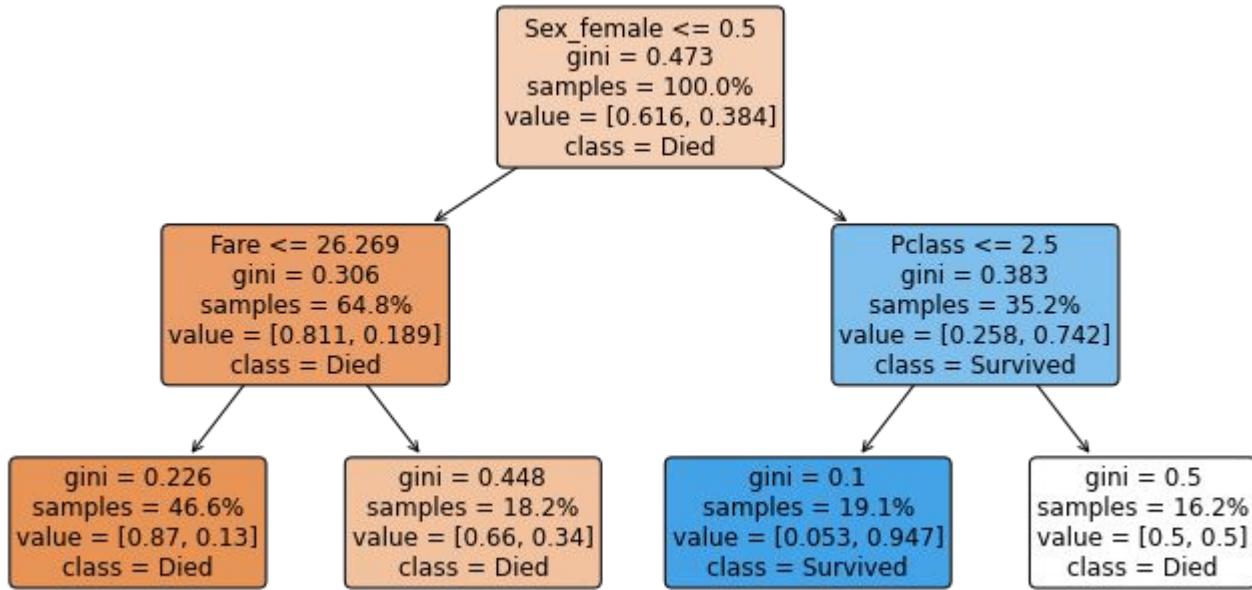


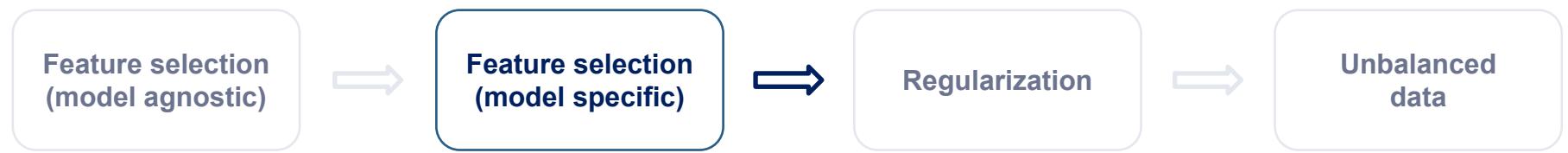
Correlation Matrix



Correlation Matrix

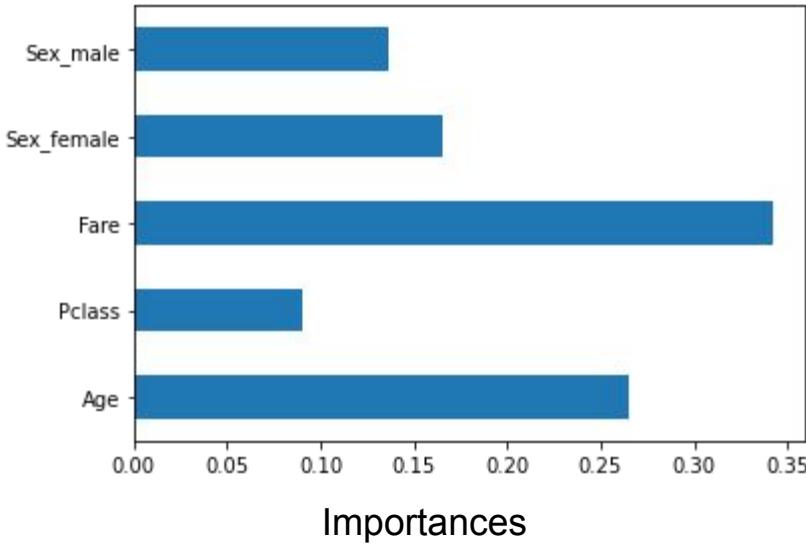


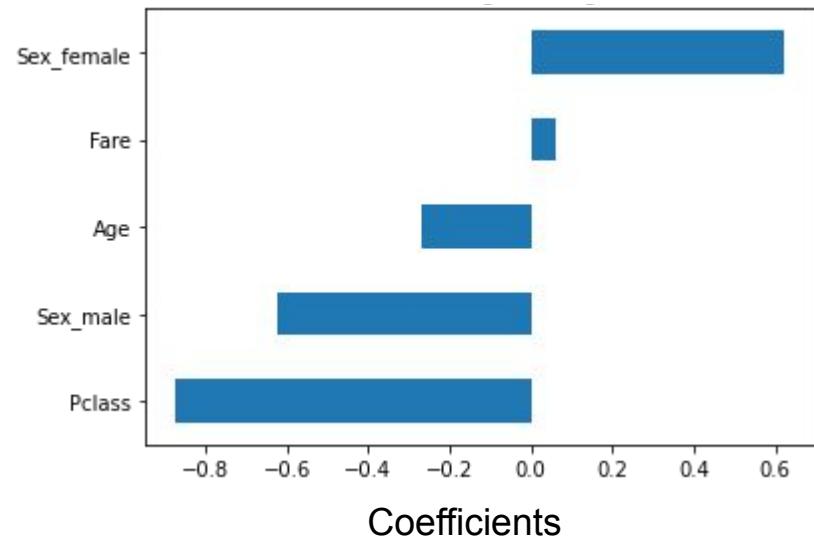
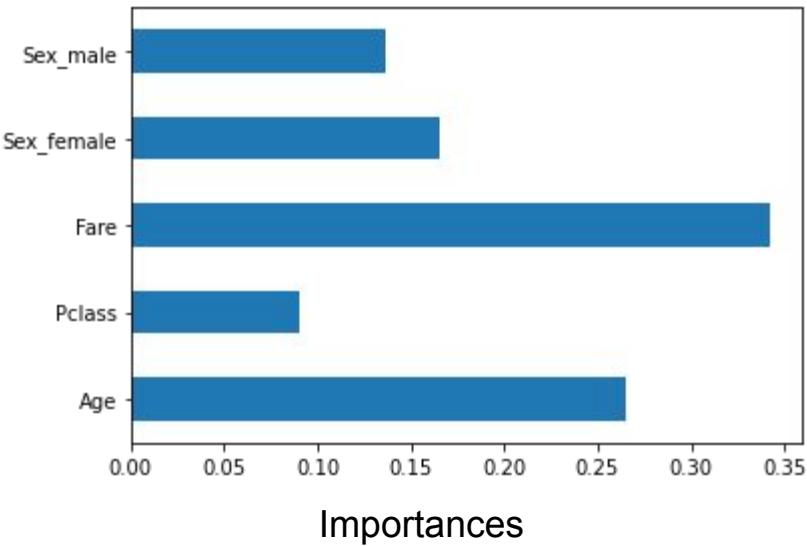
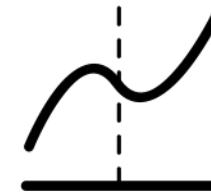


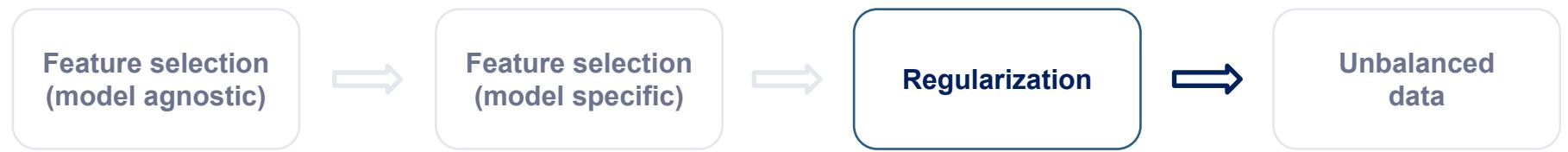


Feature importance in models

- One type of feature selection is that of looking at how important a feature is to a model that was already trained
- You can do this in models through the sets of weights it associates with each feature







Regularization

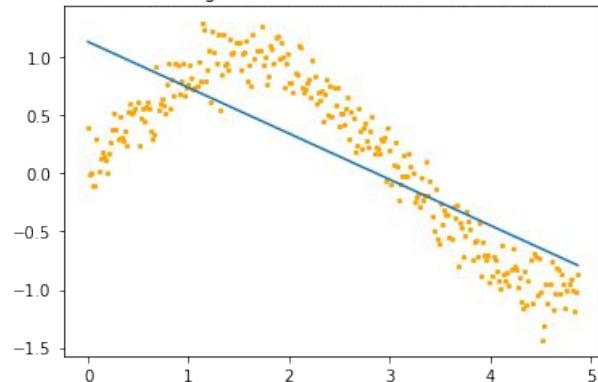
- When you have a dataset that is (a) small, (b) has way more features than observations, or (c) is sparse data, regularization can help a lot
- To force the model to use fewer features and still get similar results
- Regularization strength is a parameter you must set (i.e. not automatically optimized). It is a hyperparameter



Regularization Intuition

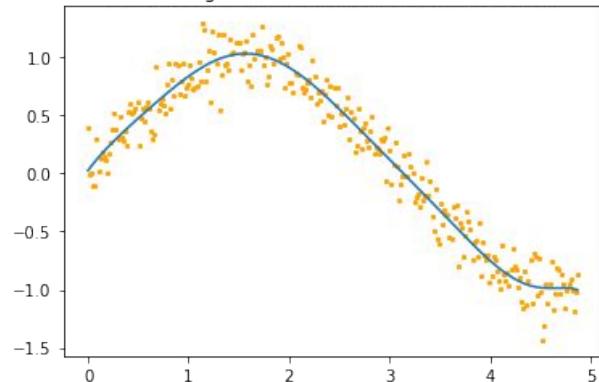
1-degree polynomial

Linear Regression ($R^2: 0.6139021052492375$)



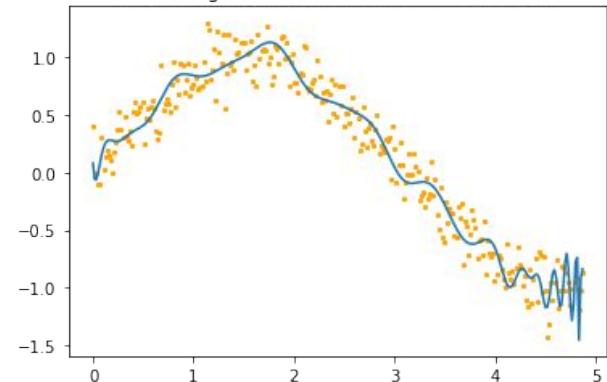
9-degree polynomial

Linear Regression ($R^2: 0.9571738246309704$)



200-degree polynomial

Linear Regression ($R^2: 0.9466925793209918$)



Too many features??
Can we pick the useful ones
automatically?

L_2 Loss Function

We need a new loss that penalizes the growth of useless feature coefficients...

$$L = \frac{1}{N} \sum_{n=1}^N (y_n - \hat{y}_n)^2 + \boxed{\lambda_2 \|\beta\|_2^2}$$

The bigger this number is,
the bigger the cost of
growing the magnitude of
coefficients

$$= \frac{1}{N} \sum_{n=1}^N (y_n - \hat{y}_n)^2 + \boxed{\lambda_2 \sum_{k=1}^K \beta_k^2}$$

This is called L_2 loss

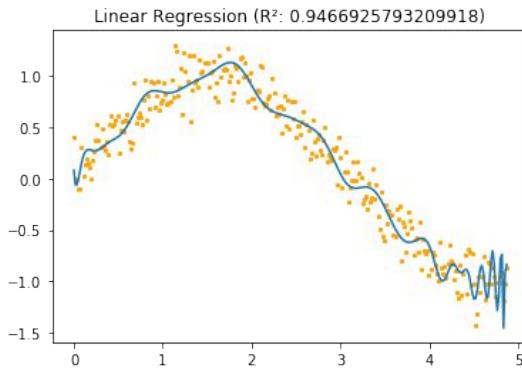
Low: good fit
High: bad fit

Low: simple model
High: complex model

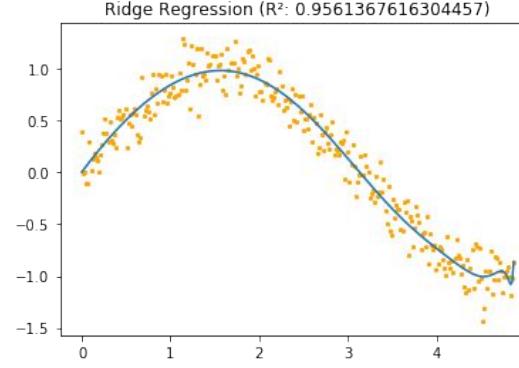
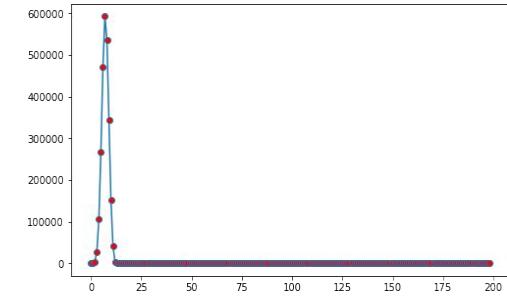


Ridge Regression

Linear
Regression:
200-degree
polynomial



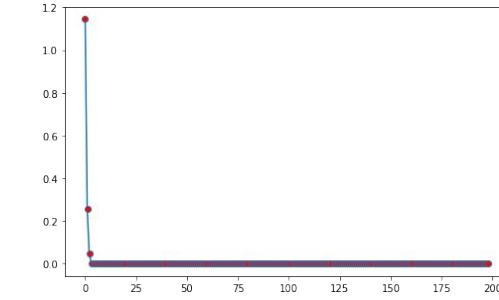
Plot of feature
coefficients



Ridge
Regression:
200-degree
polynomial

Uses L_2 loss

Plot of feature
coefficients



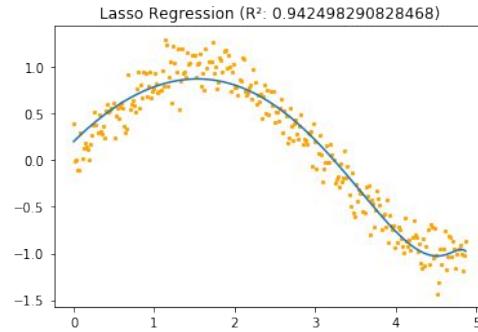
- The model is less overfit
- Most features are zero

Lasso and Elastic Net Regression

Lasso:

$$J = \frac{1}{N} \sum_{n=1}^N (y_n - \hat{y}_n)^2 + \lambda_1 \sum_{k=1}^K |\beta_k|$$

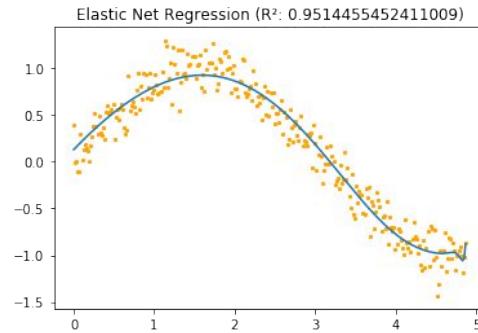
This is called L₁ loss

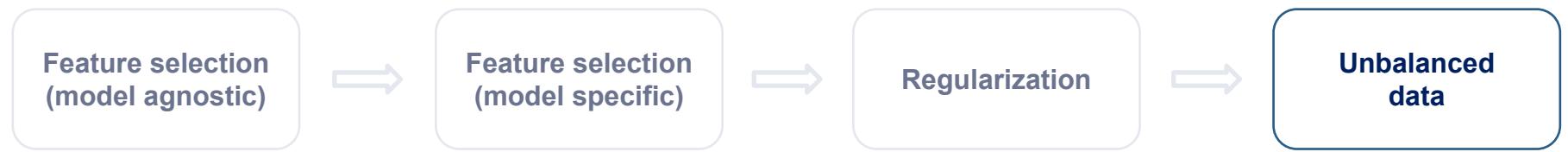


Elastic Net:

$$J = \frac{1}{N} \sum_{n=1}^N (y_n - \hat{y}_n)^2 + \lambda_1 \sum_{k=1}^K |\beta_k| + \lambda_2 \sum_{k=1}^K \beta_k^2$$

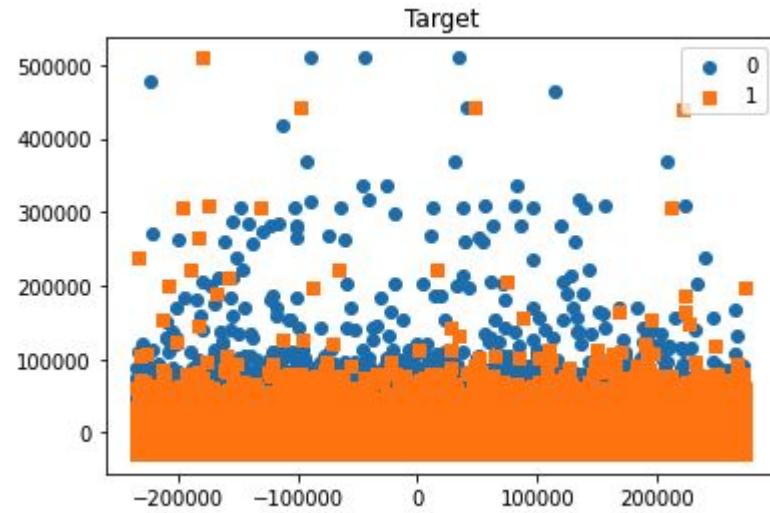
L₁ + L₂ loss





Unbalanced data

- A lot of models work better when the data is balanced
- Traps:
 - Splits
 - Metrics
 - Training



Unbalanced data splits

- The first thing you should be careful with are splits

```
X_left, X_test, y_left, y_test = train_test_split(  
    X, y, test_size=0.2, random_state=42  
)  
  
X_train, X_val, y_train, y_val = train_test_split(  
    X_left, y_left, test_size=0.2/0.8, random_state=42  
)
```

Unbalanced data splits

- The first thing you should be careful with are splits

```
X_left, X_test, y_left, y_test = train_test_split(  
    X, y, test_size=0.2, random_state=42  
)  
  
X_train, X_val, y_train, y_val = train_test_split(  
    X_left, y_left, test_size=0.2/0.8, random_state=42  
)
```

```
Size of train: 120  
Number of class 1 examples in train: 16 (13.33 %)  
Size of validation: 40  
Number of class 1 examples in validation: 6 (15.0 %)  
Size of test: 40  
Number of class 1 examples in test: 8 (20.0 %)
```

- It's easy to get a poor split that impacts your full process if you are not paying attention

Unbalanced data

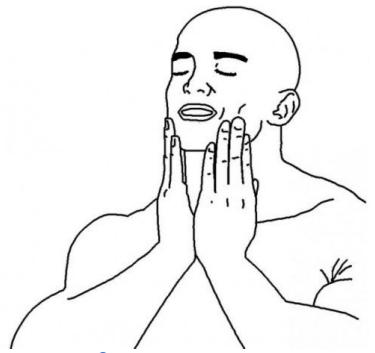
- The first thing you should be careful with are splits

```
X_left, X_test, y_left, y_test = train_test_split(  
    X, y, test_size=0.2, random_state=1234, stratify=y  
)  
  
X_train, X_val, y_train, y_val = train_test_split(  
    X_left, y_left, test_size=0.2/0.8, random_state=1234, stratify=y_left  
)
```

Unbalanced data splits

- The first thing you should be careful with are splits

```
X_left, X_test, y_left, y_test = train_test_split(  
    X, y, test_size=0.2, random_state=1234, stratify=y  
)  
  
X_train, X_val, y_train, y_val = train_test_split(  
    X_left, y_left, test_size=0.2/0.8, random_state=1234, stratify=y_left  
)  
Size of train: 120  
Number of class 1 examples in train: 18 (15.0 %)  
Size of validation: 40  
Number of class 1 examples in validation: 6 (15.0 %)  
Size of test: 40  
Number of class 1 examples in test: 6 (15.0 %)
```



Metrics

Which one is better?

Accuracy: 87.50%

Majority class
Precision: 0.91

Rcall: 0.94

F1-Score: 0.93

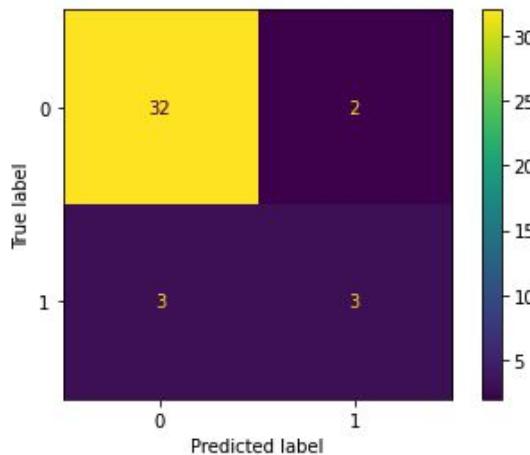
=====

Minority class
Precision: 0.60

Recall: 0.50

F1-Score: 0.55

Validation



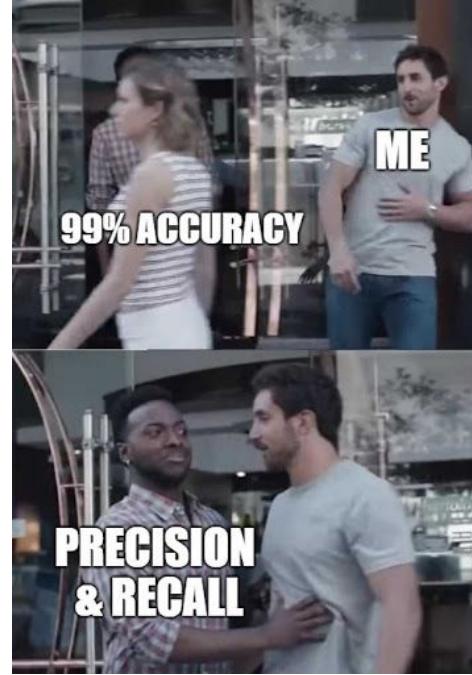
Metrics

Which one is better?

Accuracy: 87.50%

Majority class
Precision: 0.91
Rcall: 0.94
F1-Score: 0.93
=====

Minority class
Precision: 0.60
Recall: 0.50
F1-Score: 0.55

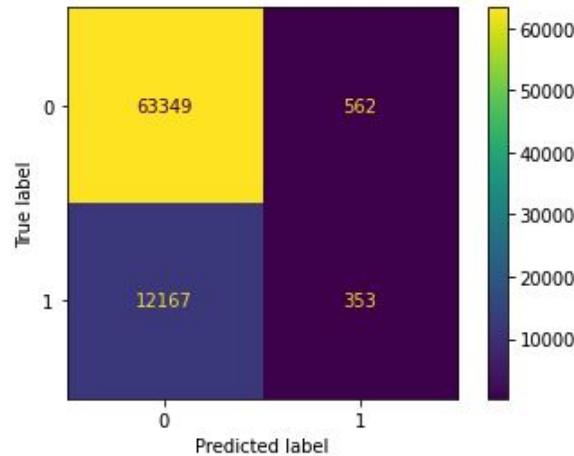


Training with unbalanced data

Majority class
Precision: 0.84
Recall: 0.99
F1-Score: 0.91

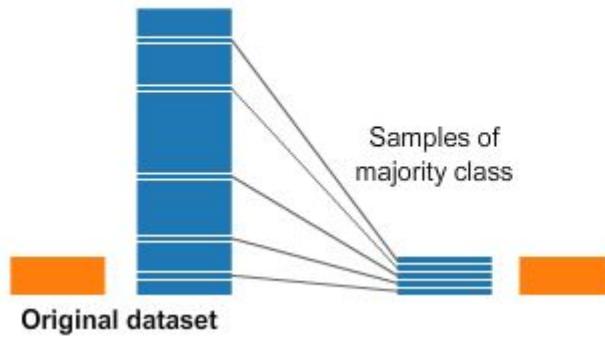
=====

Minority class
Precision: 0.39
Recall: 0.03
F1-Score: 0.05

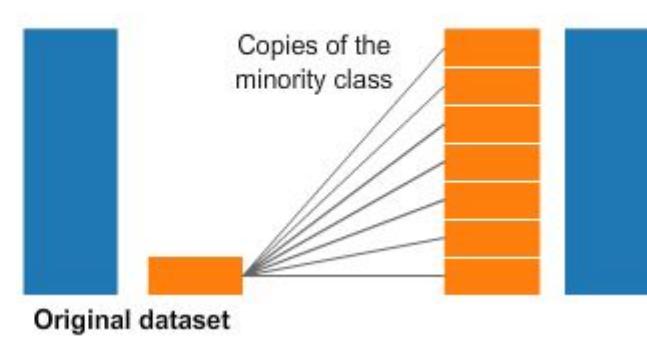


Resampling

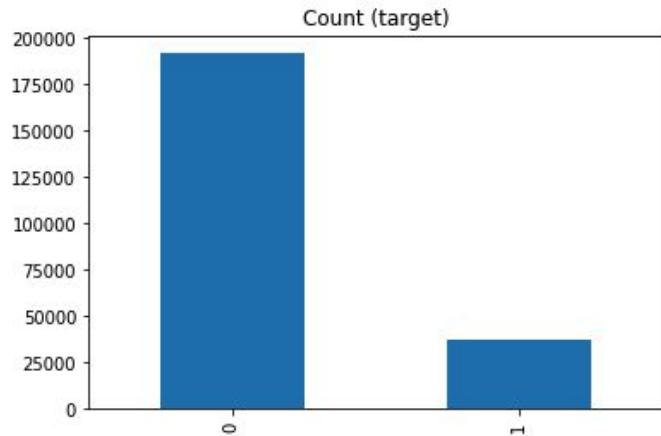
Undersampling



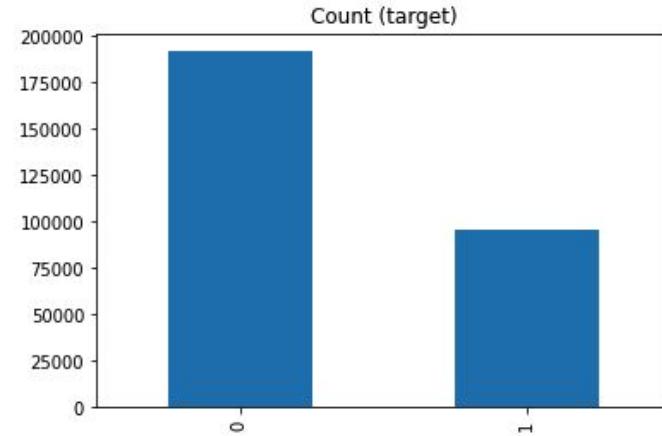
Oversampling



Resampling - oversampling

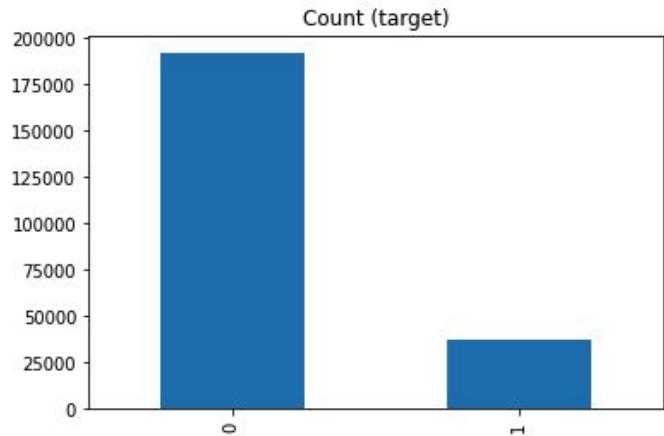


0.5 ratio
→

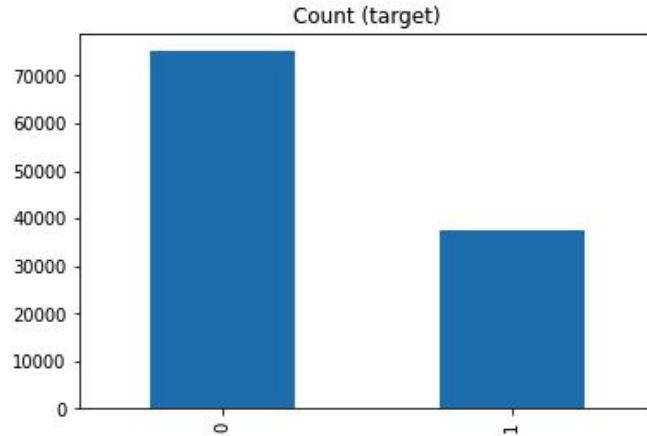


```
ros_sampler = RandomOverSampler(sampling_strategy=.5, random_state=42)  
X_over, y_over = ros_sampler.fit_resample(X_train, y_train)
```

Resampling - undersampling



0.5 ratio
→

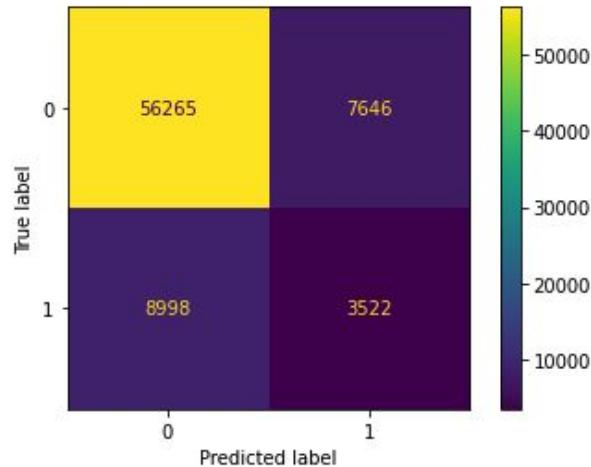


```
rus_sampler = RandomUnderSampler(sampling_strategy=.5, random_state=42)  
X_under, y_under = rus_sampler.fit_resample(X_train, y_train)
```

Training with resampling

Trade-off

Majority class
Precision: 0.86
Recall: 0.88
F1-Score: 0.87
=====
Minority class
Precision: 0.32
Recall: 0.28
F1-Score: 0.30

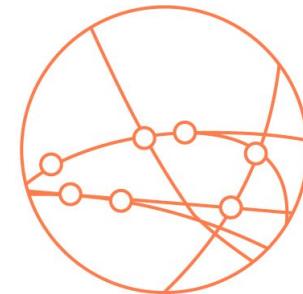


3. Recap

Recap

- 🔎 Not all features are equally important
 - Understand each feature
 - Remove suspect features
 - Examine the features independently
 - Plot a few trees!
 - When modelling, check out the feature importances
- 🛡️ Regularization can help to combat overfitting
- ⚖️ Analyze your data for imbalance
 - Be careful with the data stratification and metric chosen
 - Use resampling techniques to train your model
 - Do not apply the **resampling** to the full dataset to avoid biasing your evaluation

4. Q&A



LDSSA

SLU 15 - **HYPER** Parameter Tuning



1. Introduction

Motivation

- Models can be *tweaked* and modified to better adapt to our specific problem.
- This “fine tuning” can be done in a systematic manner.



Overview

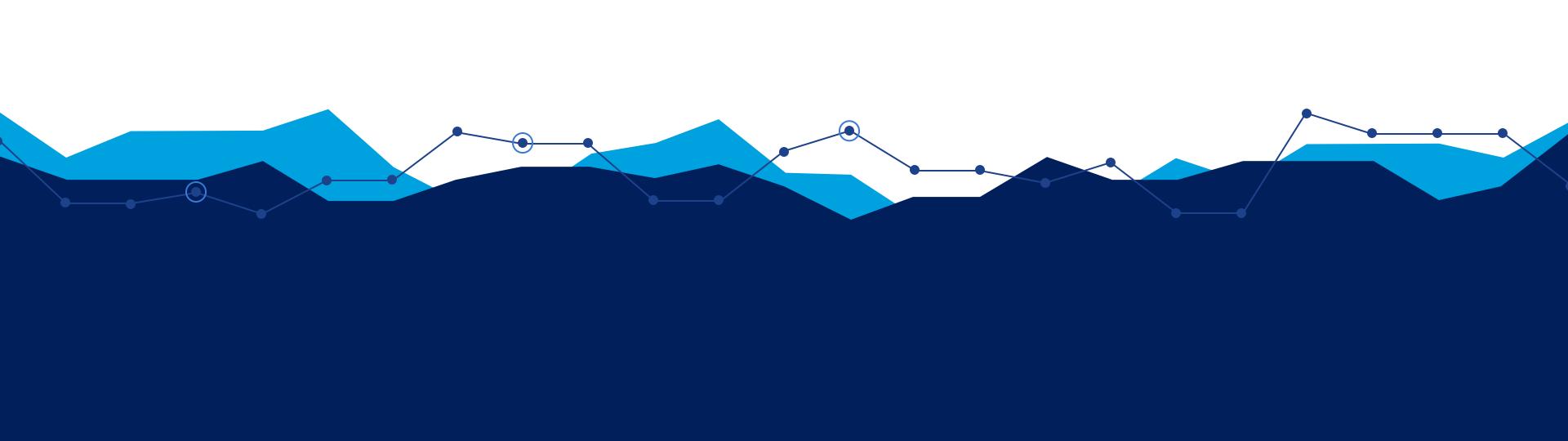
- Short recap
- Hyperparameter definition
- Hyperparameter search
- Model selection

Short Recap

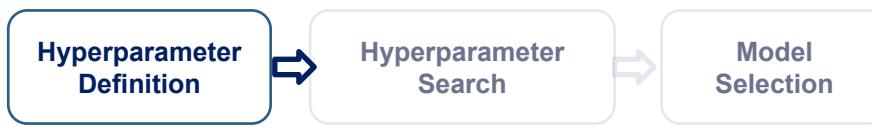
1. Problem Definition
2. Data Preparation
3. Feature Extraction
4. Model training
5. Model testing

Can we do better?





2. Topic Explanation



So... what are **HYPER** parameters?



Hyperparameters Definition

HYPER

parameters

From the greek word hyper: over, above

Parameters that are above? Yes!

The parameters that define the structure of our estimators

...and set how the estimator learns.

Logistic Regression Example

$$\ln\left(\frac{P}{1-P}\right) = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_k X_k$$

parameters

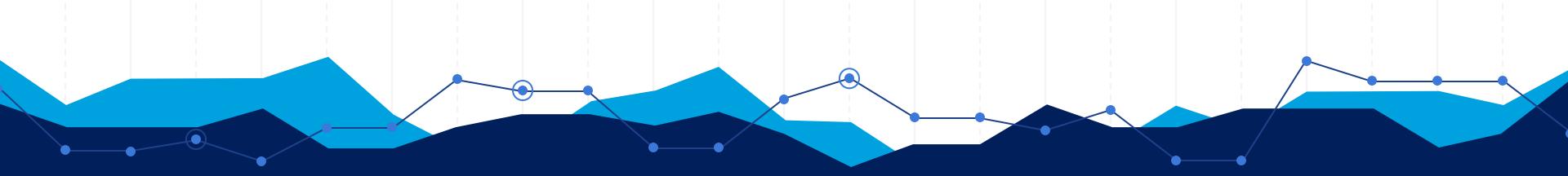
HYPER parameters

```
sklearn.linear_model.LogisticRegression(penalty='l2', dual=False,  
tol=0.0001, C=1.0, fit_intercept=True, intercept_scaling=1,  
class_weight=None, solver='warn', max_iter=100, warm_start=False)
```

Cheatsheet

Estimator	Hyperparameters
Logistic Regression	penalty type (l1 or l2), C (regularization parameter)
SVMs	kernel, C (regularization parameter)
Tree Ensembles	n_estimators, max_depth
KNN	n_neighbors, weights





Don't despair yet!

HYPER parameter tuning can be automated!





Hyperparameter search

General Hyper Parameter tuning algorithm

1. Split **train data** into **train** and **validation** subsets (K-folds);
2. Select hyperparameter values;
3. Instantiate estimator with selected values;
4. **Train estimator on training subset** ;
5. **Score estimator on the validation subset**;
6. Repeat 2-5 for **all desired hyperparameter combinations**.

Goal: Select estimator with the best score on the validation metric.

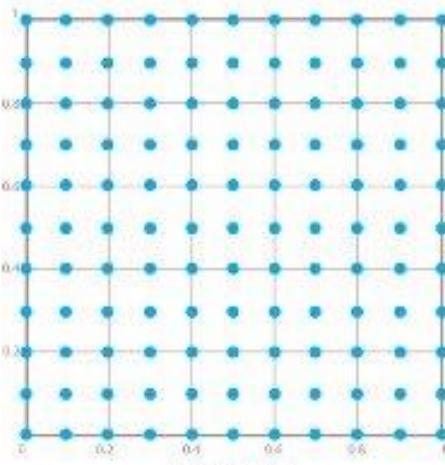


Hyperparameter search (2)

HYPERParameter Search Space- Defines the region of hyperparameters we want to test

Grid Search

([*sklearn.model_selection.GridSearchCV*](#))

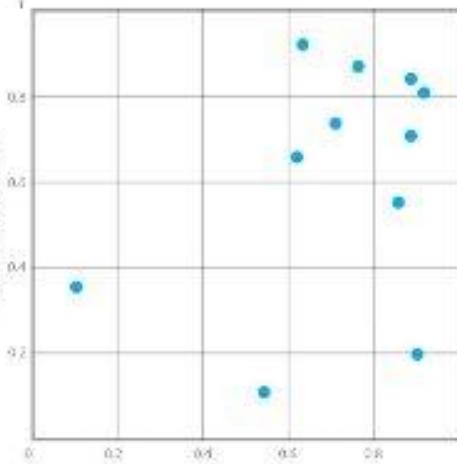


Searches the space linearly
from start to end.

Should not be used when
training takes considerable
time.

Random Search

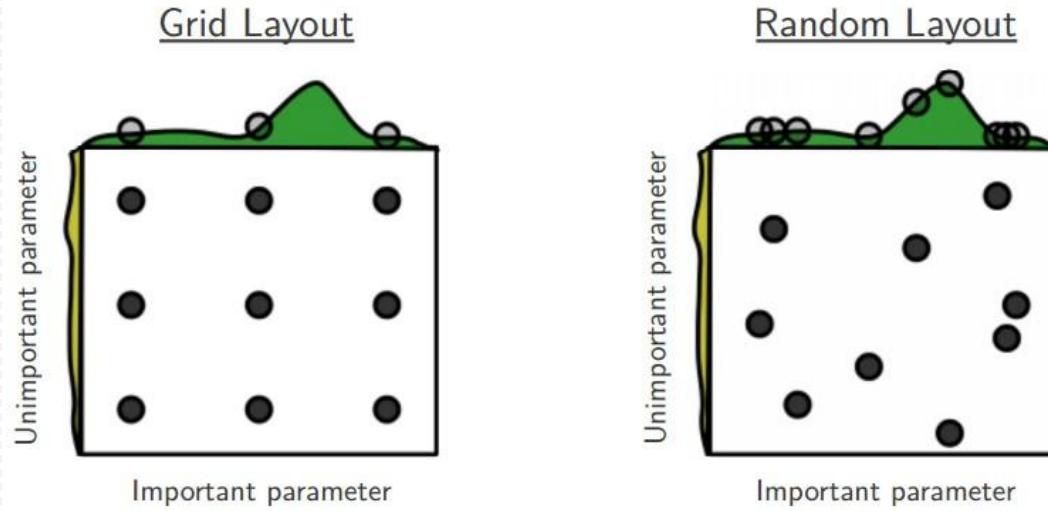
([sklearn.model_selection.RandomizedSearchCV](#))



Searches the space randomly.

Covers a wide range of parameters with a low number of samples

Grid vs Random Search





Model Selection

Now that we have several models assessed...

We just need to choose the best model!

... Right?!



Model Selection

There are other aspects besides the score that we must consider when selecting models:

- Training Time;
- Prediction Time;
- Interpretability.

3. Recap



Recap

- Hyperparameters define the structure of our estimators
 - Different from parameters, which are computed by the model
- Hyperparameter search to select best hyperparameters:
 - Grid search
 - Random search
- Model selection



Happy **HYPER** parameter tuning



4. Q&A