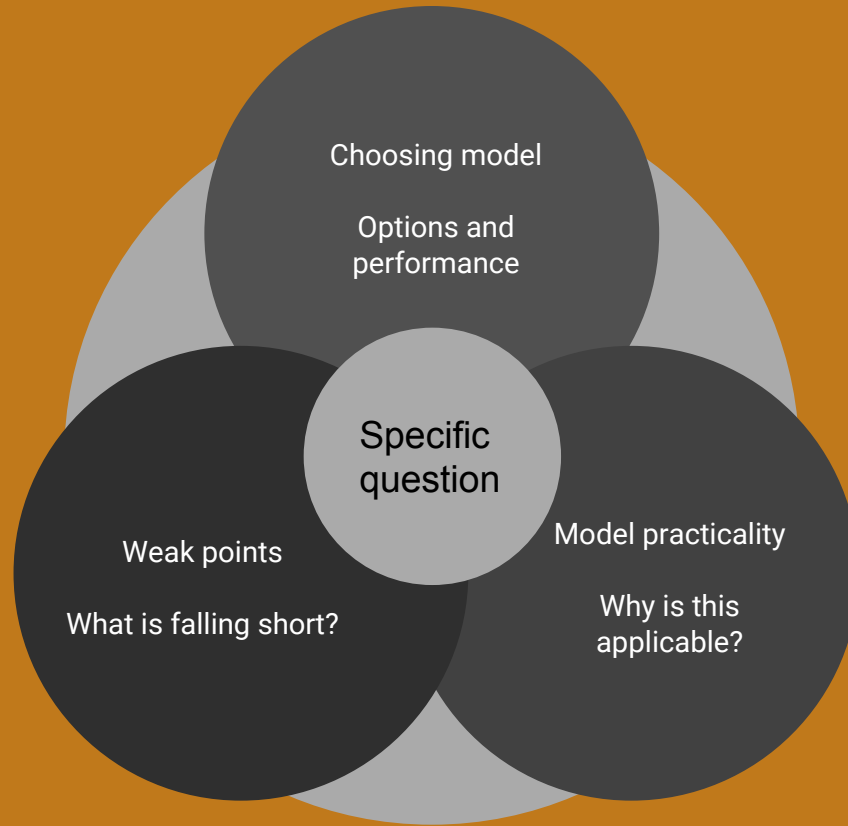


Supervised learning

An end-to-end project

<https://tinyurl.com/y9u4tqm3>







Adoption rates and outcomes

Research question: Are there features that can explain or predict the most likely outcome for animals at the Austin shelter?





Characteristics of the data

Context

The data set comes from an open data initiative in Austin - the data was stored on servers on the Kaggle domain.

Content

The dataset contains shelter outcomes of several types of animals and breeds from 10/1/2013 to the present with a hourly time frequency.

Approach

Outcome of interest: what can we predict?
Animal status outcomes or adoption rates seemed like the best candidate.



High-level process: three steps

Consistency and quality

Data preparation

The data quality is not model-ready. It needs transformations, selection, and filtering.

Initial models

Classification models

The course has covered many different models for classification and regression. We have at least eight different models to attempt.

Model tuning

Increase performance

After finding inherently strong models against this data, we can tune a few of the best - single and ensemble.

```
In [2]: # Reading in the data set.
df = pd.read_csv('../_Datasets/aac_shelter_outcomes.csv')
df2 = df.copy()
df.head(3)
```

Out[2]:

	age_upon_outcome	animal_id	animal_type	breed	color	date_of_birth	datetime	monthyear	name	o
0	2 weeks	A684346	Cat	Domestic Shorthair Mix	Orange Tabby	2014-07-07T00:00:00	2014-07-22T16:04:00	2014-07-22T16:04:00	NaN	P
1	1 year	A666430	Dog	Beagle Mix	White/Brown	2012-11-06T00:00:00	2013-11-07T11:47:00	2013-11-07T11:47:00	Lucy	P
2	1 year	A675708	Dog	Pit Bull	Blue/White	2013-03-31T00:00:00	2014-06-03T14:20:00	2014-06-03T14:20:00	*Johnny	N

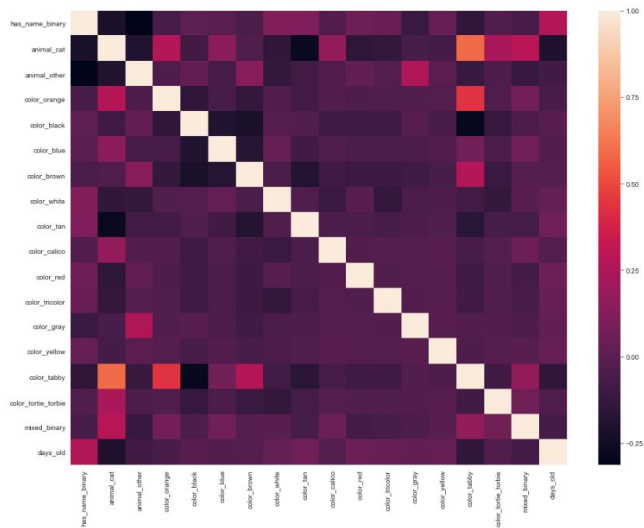
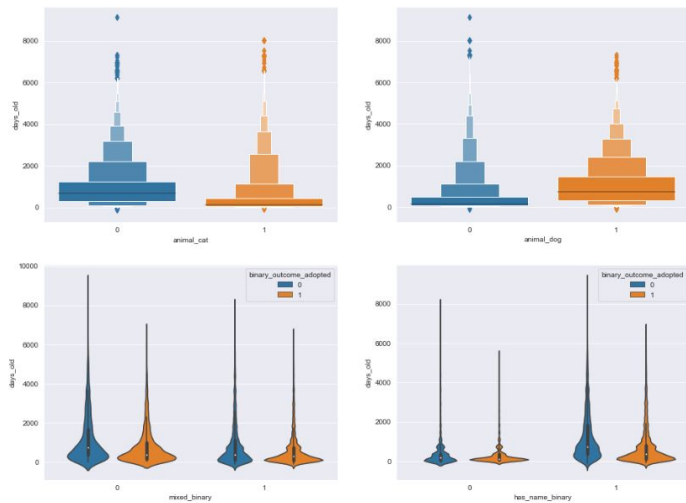
```
In [3]: # Nulls aren't too bad, except for name and outcome_subtype. (I did not originally see the two on outcome_
df.info())

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 78256 entries, 0 to 78255
Data columns (total 12 columns):
age_upon_outcome    78248 non-null object
animal_id           78256 non-null object
animal_type         78256 non-null object
breed               78256 non-null object
color               78256 non-null object
date_of_birth       78256 non-null object
datetime            78256 non-null object
monthyear           78256 non-null object
name                54378 non-null object
outcome_subtype     35963 non-null object
outcome_type        78344 non-null object
sex_upon_outcome    78254 non-null object
dtypes: object(12)
memory usage: 7.2+ MB
```

Data prep: wrangling, cleaning up, and creating features



With so few continuous features, visualizing our data is hard to do.





Testing different models

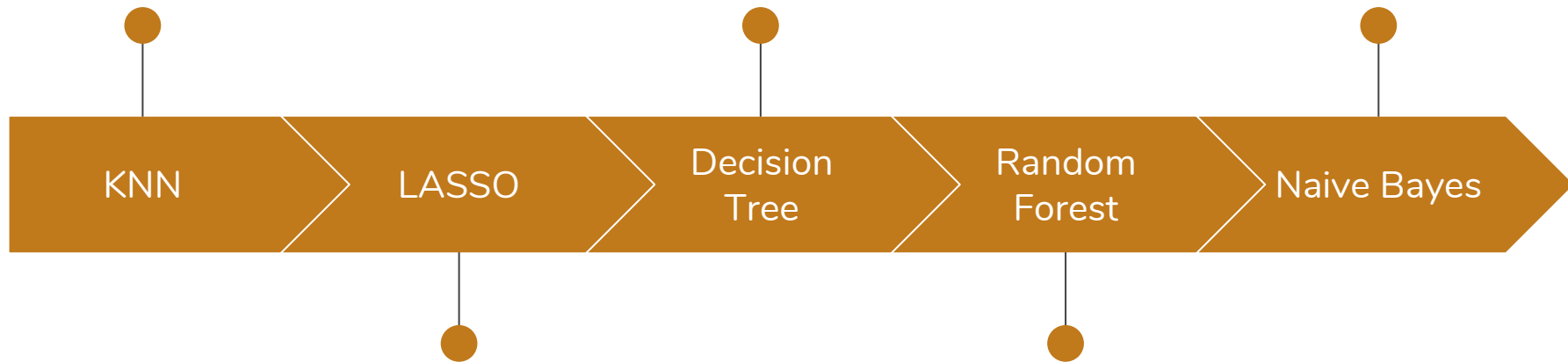
Classification algorithms



Algorithm to identify
and group data points
based on neighbors

Decision process
selecting strong
information gain splits

Gaussian and
Bernoulli for
classification



Logistic regression
approach discarding
weaker features

Black box ensemble
approach of many
trees

0.5-0.79

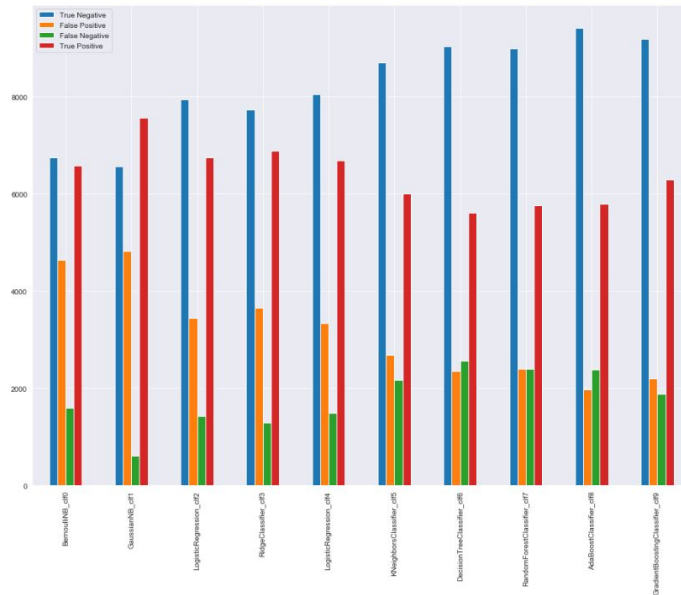
SKLearn metric scores across our models on binary outcomes



Model metrics:

Graphing confusion matrices

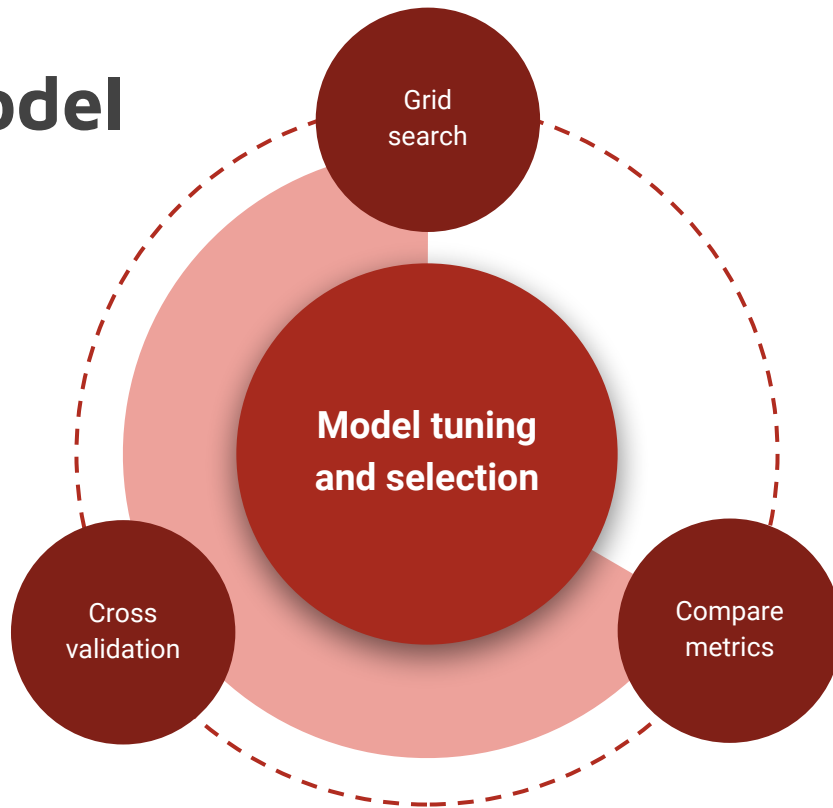
More specifically, how do the
models measure side-by-side?





Improving the model

We started with a handful of classification algorithms. After checking the default parameters and their scores, we gather an idea of what is working. Then, after pushing the model through these steps, we can determine which are working, and if they might improve.





Tuning a few parameters did not make a significant difference - even after trying several: gini, entropy, max_depth, and alpha values. In general, I believe this fine tuning works very well when you have a better understanding of what inputs the model needs. The score improvements sat around one to four percent.

[illegible]



Back to the drawing board

What about a multiclass model?



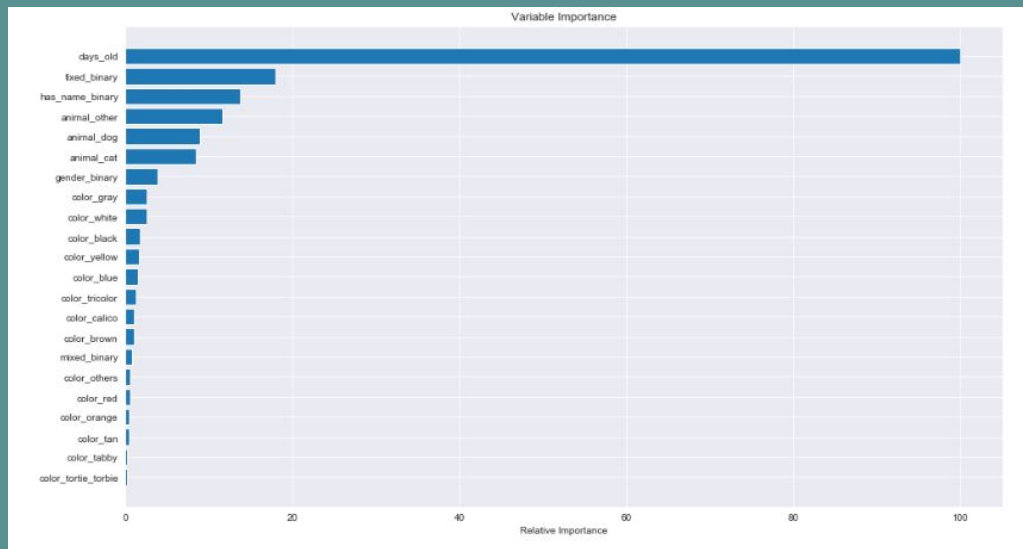


Fixing features: trade offs between retention or tossing

Binary outcomes from many options and some feature discarding may be ill-advised

Hitting more than 0.75 accuracy across a binary outcome is certainly an improvement over random guessing and chance. Our models were able to do this pretty early on - probably due to a healthy-sized, simple, and clean data set.

Revisiting the feature transformation and including some information previously left out was able to improve this at orders of magnitude. Hitting F1 scores of about 0.7 against random-guessing likelihoods of 0.25 is certainly trending in the right direction.



Ensemble model: gradient boosting does provide feature importances



Models: best performers

Weaknesses: Overall, the numbers turned out to be decent. Hitting higher accuracies is always nice, but this is a trade off with potential overfitting. The inability of this data set to intuitively provide dozens of visualization options is probably the weakest part of its modeling process; its absence was noticed.

Pure predictive power: The Gradient boosting classifier seemed to pull some of the best scores off - and it does so against an F1 score. This covers several metrics to indicate performance. It is an ensemble method, though, and this does not tell us much about this question - *why?*

Explanatory power: K-Neighbors classifier does a solid job for a single model. The F1 score indicates it performs across different metrics. Luckily, this model does provide outcome insight - an explanation as to why.



Q & A

Thanks for attending!

