

Introduction to Modeling in Python

Five College DataFest 2019
Quentin Dupupet (MassMutual DSDP)

What is a model?

- A **model** is a mathematical function of form $y = f(X)$
 - y is some quantity of interest that we want to predict
 - X is a collection of individual observations and some information about them
 - f is our model
- Statistical models - derive insight about a real-world event while quantifying uncertainty about predictions
- Machine learning models - output the best prediction possible without regard for interpretability
- “All models are wrong, but some are useful”
 - A model is a simplified representation of a potentially complex real-world phenomenon

How does a model work?

- Models “learn” by associating a set of various observations of different feature inputs (e.g. weight, age) with a set of corresponding outputs (e.g. height) in a way such that an arbitrary **cost function** is minimized
 - This combination is known as the “**training set**”
- Modeling is typically one of the last steps in a data science pipeline
 - Make sure data is cleaned and relevant features are selected or processed beforehand
 - Typical feature engineering tasks include removal of missing values, normalization/standardization, or one-hot encoding depending on model type
- The process of modeling includes many potentially subjective steps that need to be justified by the modeler, so be careful!

Model performance

- Models are ultimately judged on their ability to generalize to never-before seen data (frequently known as the “**test set**”)
 - A train/test partition of 80%/20% is often used in practice
- A fundamental concept in modeling is the **bias-variance** tradeoff
 - **Bias** - error introduced in the specification of the form of the model
 - **Variance** - error introduced by inherent variability in the data itself
- Too much bias indicates the model has **underfit** to the data, while too much variance indicates the model is **overfit** to the set of observations it has seen
- The best models balance both!

Train/test code

```
from sklearn.datasets import load_boston

data = load_boston()
df = pd.DataFrame(data.data, columns = data.feature_names)
df['target'] = data.target
df.head()
```

	CRIM	ZN	INDUS	CHAS	NOX	RM	AGE	DIS	RAD	TAX	PTRATIO	B	LSTAT	target
0	0.00632	18.0	2.31	0.0	0.538	6.575	65.2	4.0900	1.0	296.0	15.3	396.90	4.98	24.0
1	0.02731	0.0	7.07	0.0	0.469	6.421	78.9	4.9671	2.0	242.0	17.8	396.90	9.14	21.6
2	0.02729	0.0	7.07	0.0	0.469	7.185	61.1	4.9671	2.0	242.0	17.8	392.83	4.03	34.7
3	0.03237	0.0	2.18	0.0	0.458	6.998	45.8	6.0622	3.0	222.0	18.7	394.63	2.94	33.4
4	0.06905	0.0	2.18	0.0	0.458	7.147	54.2	6.0622	3.0	222.0	18.7	396.90	5.33	36.2

Train/test code

```
from sklearn.model_selection import train_test_split

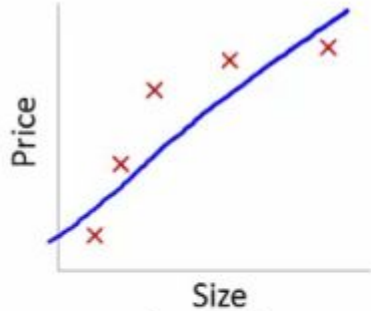
X = data.data
y = data.target

X_train, X_test, y_train, y_test = train_test_split(X, y, test_size = 0.25)

print(X_train.shape, y_train.shape, X_test.shape, y_test.shape)
```

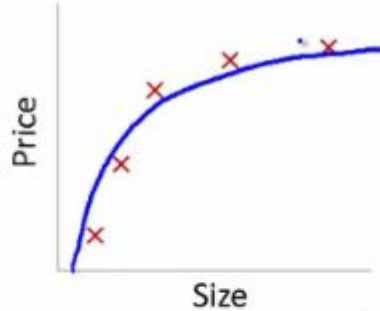
```
(379, 13) (379,) (127, 13) (127,)
```

Bias-variance tradeoff



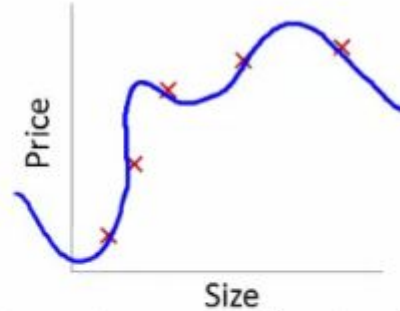
$$\theta_0 + \theta_1 x$$

High bias
(underfit)



$$\theta_0 + \theta_1 x + \theta_2 x^2$$

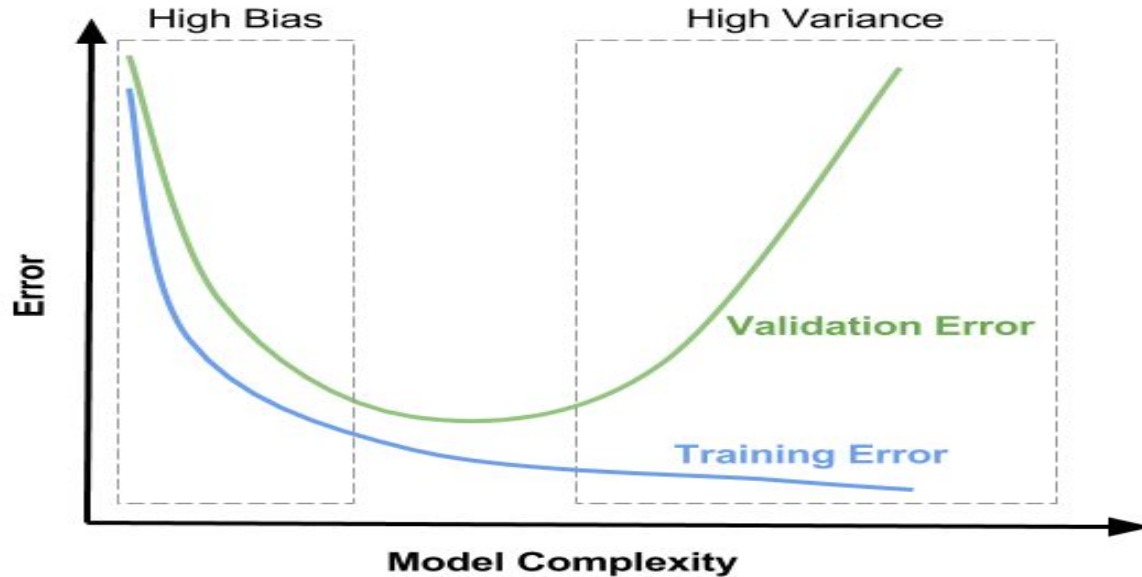
“Just right”



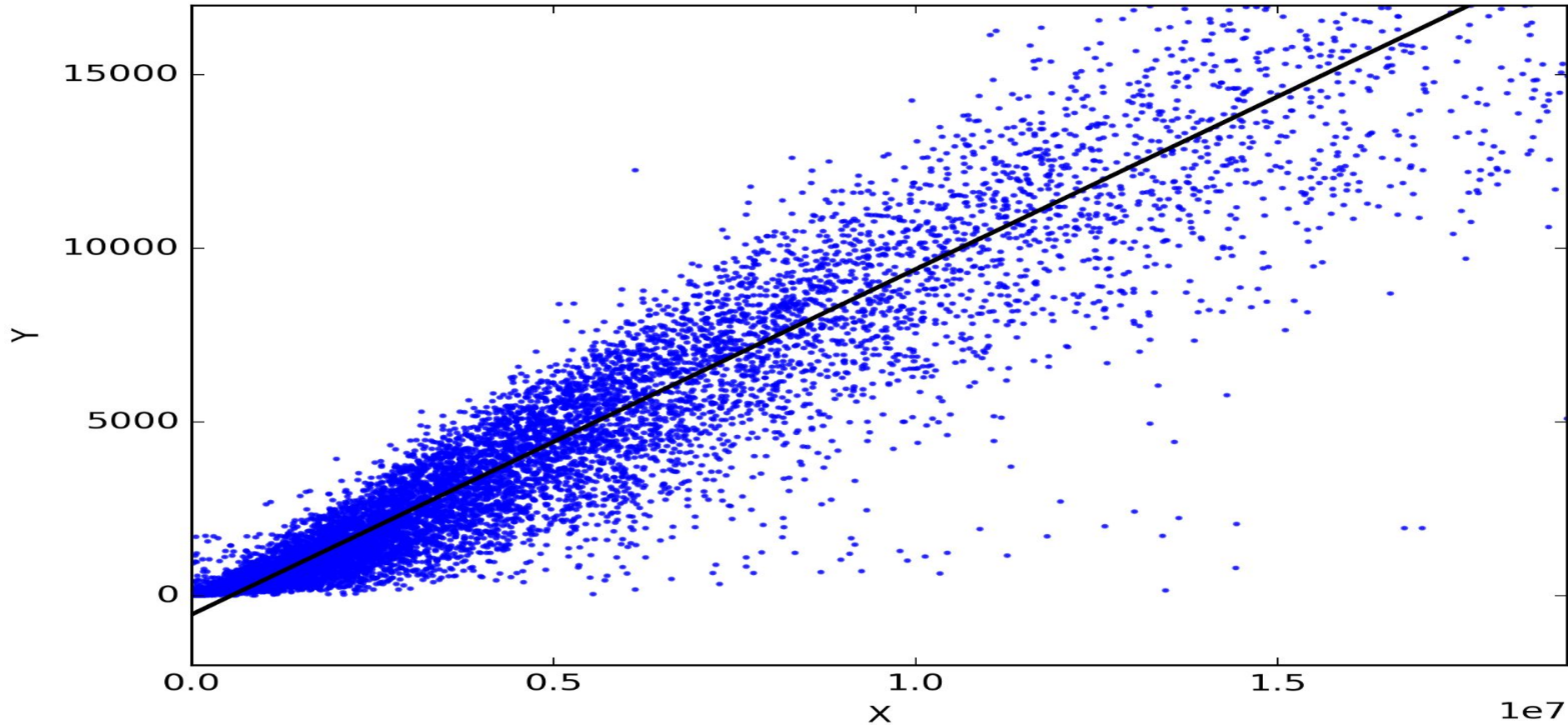
$$\theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$$

High variance
(overfit)

Bias-variance tradeoff



Regression

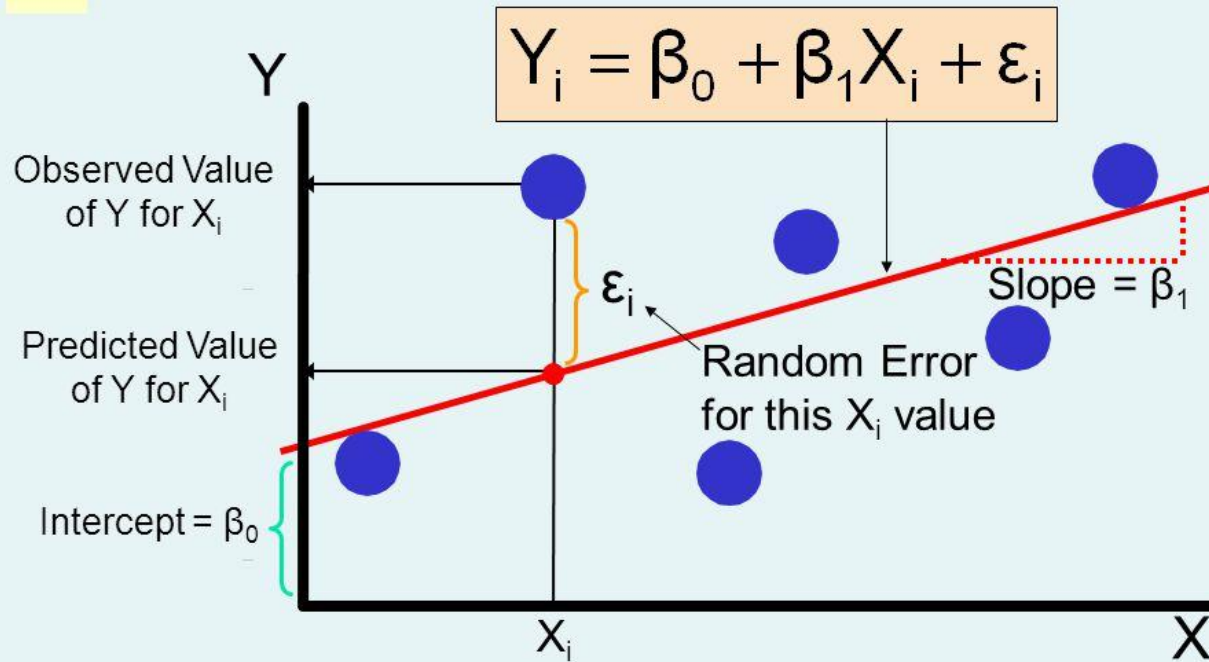


Regression background

- “**Regression**” is the act of predicting a continuous quantity (e.g. the price of something)
- An example is **linear regression**, which fits an $(N-1)$ -dimensional hyperplane to an $M \times N$ numeric matrix of observations such that distance between the hyperplane and all observations is minimized
 - A typically used cost function is **mean-squared error** (MSE) e.g. the average of the sum of $(y - y_{pred})^2$ across all M observations
- Linear regression assumes normally distributed error terms, and is therefore a **parametric** model
 - Other assumptions include independence and even variance of errors, a true underlying linear relationship, and uncorrelated input variables

Simple Linear Regression Model

DCOVA
(continued)



Linear regression code - fitting the model

```
import numpy as np
from sklearn.linear_model import LinearRegression

lin_reg = LinearRegression()
lin_reg.fit(X_train, y_train)

print(
    np.round(lin_reg.coef_, 1),
    '\n',
    np.round(lin_reg.intercept_, 1)
)
```

```
[-0.1  0.   -0.2  0.   -0.   0.1 -0.3 -0.2 -0.   0.1 -0.2 -0.3 -0. ]
3.5
```

$$y = -0.1(CRIM) - 0.2(INDUS) + 0.1(RM) - 0.3(AGE) - 0.2(DIS) + 0.1(TAX) - 0.2(PTRATIO) - 0.3(B) + 3.5$$

Linear regression code - evaluation

```
from sklearn.metrics import mean_squared_error
y_pred = lin_reg.predict(X_test)

print('R^2 Score:', lin_reg.score(X_test, y_test))
print('Mean squared error: ', mean_squared_error(y_test, y_pred))
```

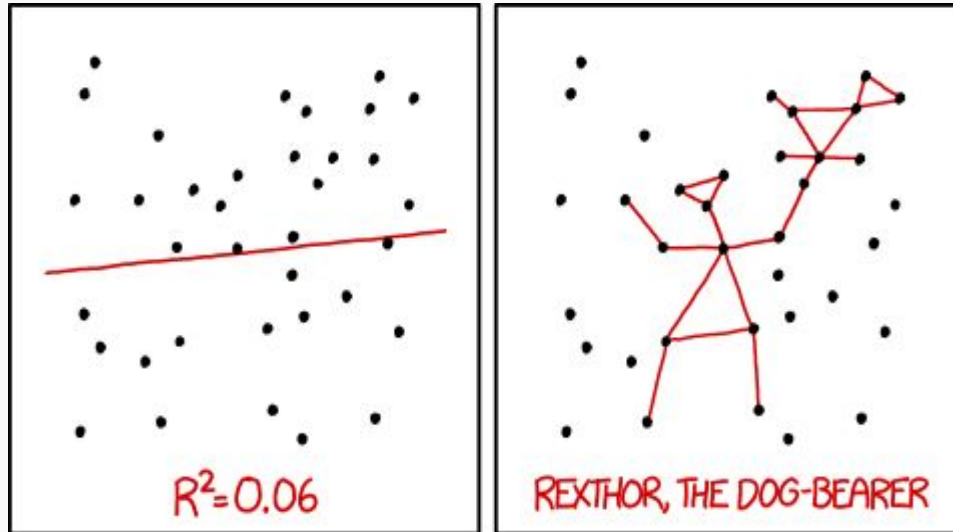
R^2 Score: 0.8850825151531964

Mean squared error: 0.06843974653098528

The **score()** method uses R^2 , or the coefficient of determination, which measures how much of the total variation in the data is explained successfully by the model.

Values range from 0 to 1

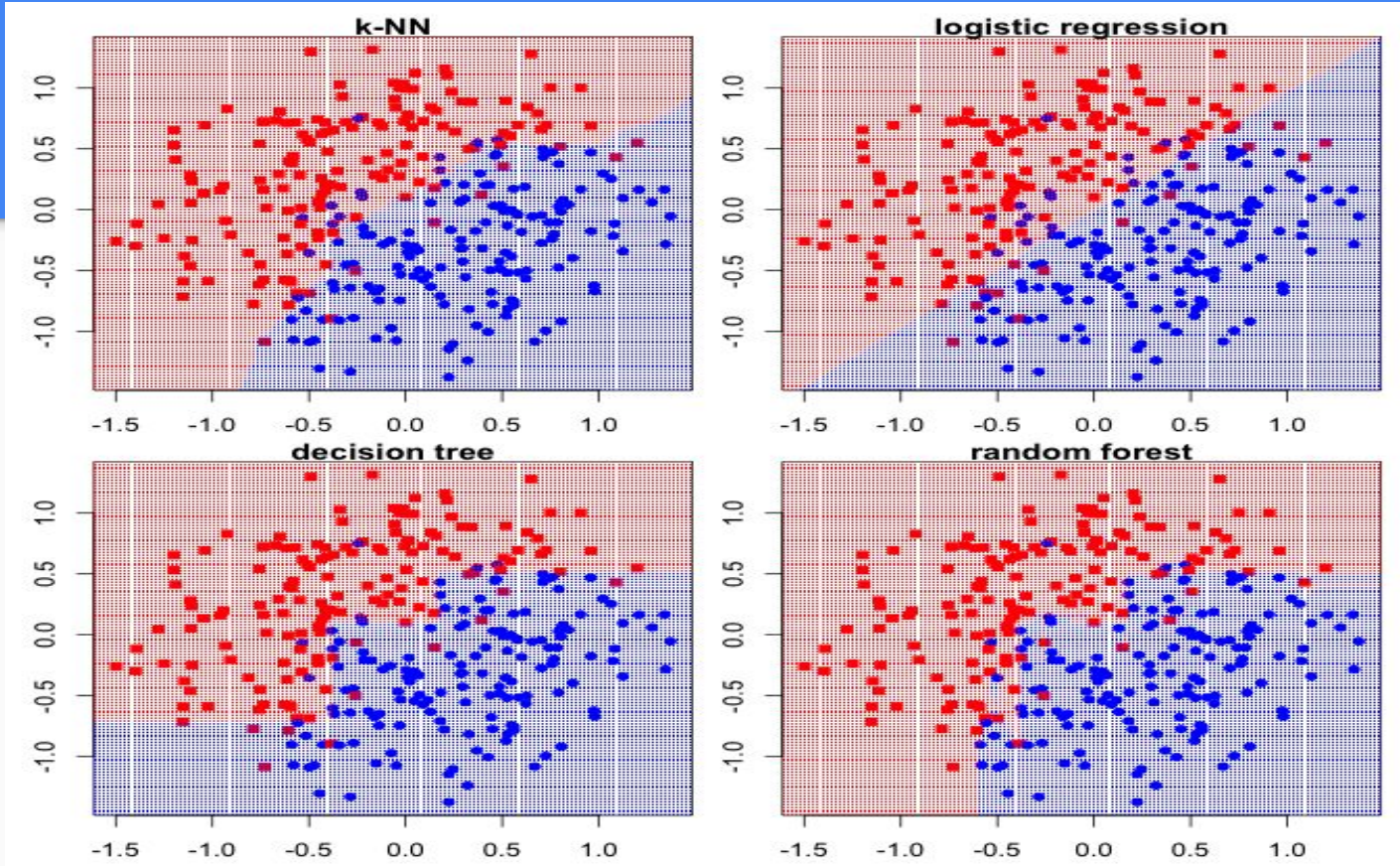
Relevant XKCD



I DON'T TRUST LINEAR REGRESSIONS WHEN IT'S HARDER TO GUESS THE DIRECTION OF THE CORRELATION FROM THE SCATTER PLOT THAN TO FIND NEW CONSTELLATIONS ON IT.

https://www.explainxkcd.com/wiki/index.php/1725:_Linear_Regression

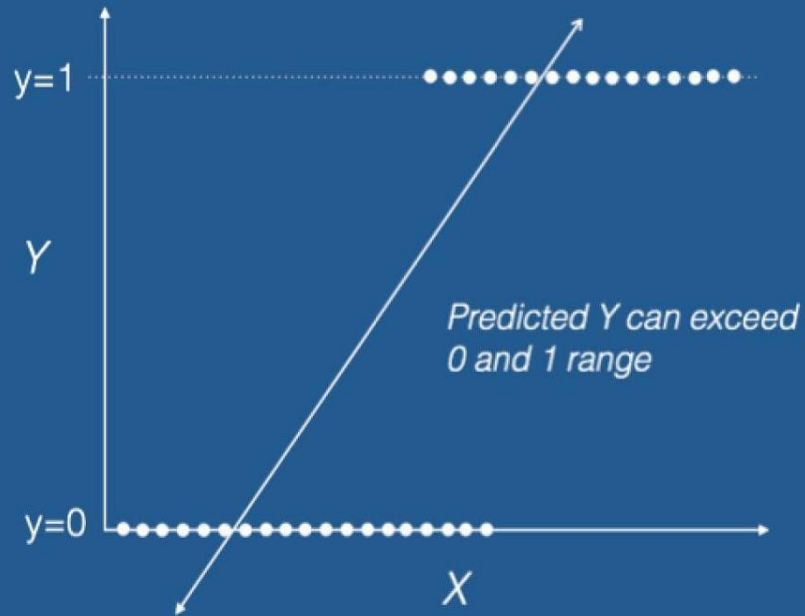
Classification



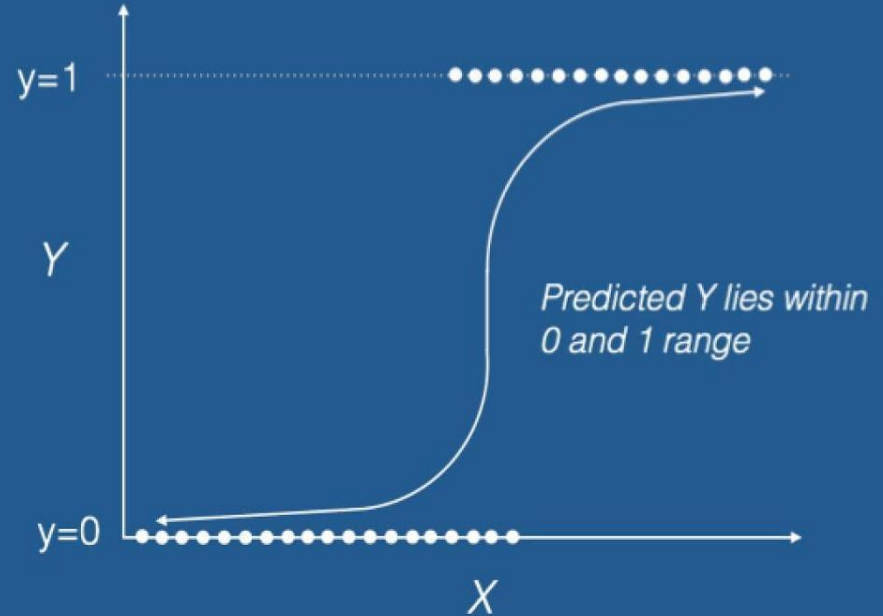
Classification background

- In **classification**, our interest is a discrete class rather than a continuous quantity (e.g. a label indicating whether a sale was made or not), and our goal is to find a **decision boundary** that accurately separates classes
- A common classification model is **logistic regression**, which outputs a probability of each observation belonging to a class
 - Like linear regression, it consists of a linear set of coefficients which are fed into a link function ensuring they output a real number between 0 and 1
- Classifiers are generally evaluated through metrics like **accuracy**, **precision**, and **recall**
 - The right metric to optimize will depend on your modeling situation

Linear Regression



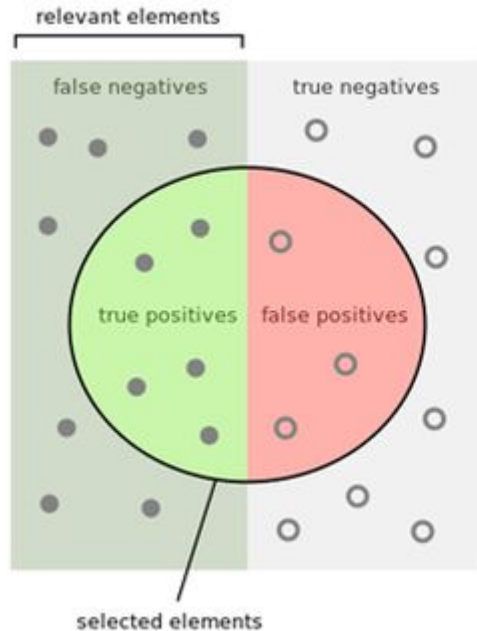
Logistic Regression



Classification evaluation - confusion matrix

		True class		Measures
		Positive	Negative	
Predicted class	Positive	True positive TP	False positive FP	Positive predictive value (PPV) $\frac{TP}{TP+FP}$
	Negative	False negative FN	True negative TN	Negative predictive value (NPV) $\frac{TN}{FN+TN}$
Measures		Sensitivity $\frac{TP}{TP+FN}$	Specificity $\frac{TN}{FP+TN}$	Accuracy $\frac{TP+TN}{TP+FP+FN+TN}$

Classification evaluation - precision/recall



How many selected items are relevant?

$$\text{Precision} = \frac{\text{true positives}}{\text{true positives} + \text{false positives}}$$

How many relevant items are selected?

$$\text{Recall} = \frac{\text{true positives}}{\text{true positives} + \text{false negatives}}$$

Classification Code

```
from sklearn.datasets import load_wine

wine = load_wine()
wine_df = pd.DataFrame(wine.data, columns = wine.feature_names)
wine_df['target'] = wine.target
wine_df.head()
```

resium	total_phenols	flavanoids	nonflavanoid_phenols	proanthocyanins	color_intensity	hue	od280/od315_of_diluted_wines	proline	target
127.0	2.80	3.06	0.28	2.29	5.64	1.04	3.92	1065.0	0
100.0	2.65	2.76	0.26	1.28	4.38	1.05	3.40	1050.0	0
101.0	2.80	3.24	0.30	2.81	5.68	1.03	3.17	1185.0	0
113.0	3.85	3.49	0.24	2.18	7.80	0.86	3.45	1480.0	0
118.0	2.80	2.69	0.39	1.82	4.32	1.04	2.93	735.0	0

```
X = wine.data
y = wine.target
```

```
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size = 0.25, stratify = y)
print(X_train.shape, y_train.shape, X_test.shape, y_test.shape)
```

```
(133, 13) (133,) (45, 13) (45,)
```

Classification Code - fitting the model

```
from sklearn.linear_model import LogisticRegression
```

```
log_reg = LogisticRegression()
```

```
log_reg.fit(X_train, y_train)
```

```
print(  
    np.round(log_reg.coef_, 2),  
    '\n',  
    np.round(log_reg.intercept_)  
)
```

```
[[-0.57  0.78  1.03 -0.58 -0.02 -0.04  1.23  0.05 -0.44 -0.   -0.15  0.82  
  0.02]  
[ 0.68 -1.04 -0.7   0.26  0.01  0.32  0.52  0.21  0.45 -1.87  0.73  0.15  
 -0.01]  
[-0.08  0.72 -0.08  0.06  0.01 -0.61 -1.49 -0.01 -0.64  1.33 -0.41 -1.26  
 -0.   ]]  
[-0.  0. -0.]
```

Classification code - evaluation

```
from sklearn.metrics import confusion_matrix, classification_report

print('Accuracy: ', log_reg.score(X_test, y_test))
y_pred = log_reg.predict(X_test)

conf_matrix = confusion_matrix(y_test, y_pred)
print('\nConfusion Matrix:\n',
      pd.DataFrame(
          confusion_matrix(y_test, y_pred),
          columns = ['pred: 0', 'pred: 1', 'pred: 2'],
          index = ['true: 0', 'true: 1', 'true: 2']),
      '\n\nClassification Report:\n',
      classification_report(y_test, y_pred))
```

Accuracy: 0.9555555555555556

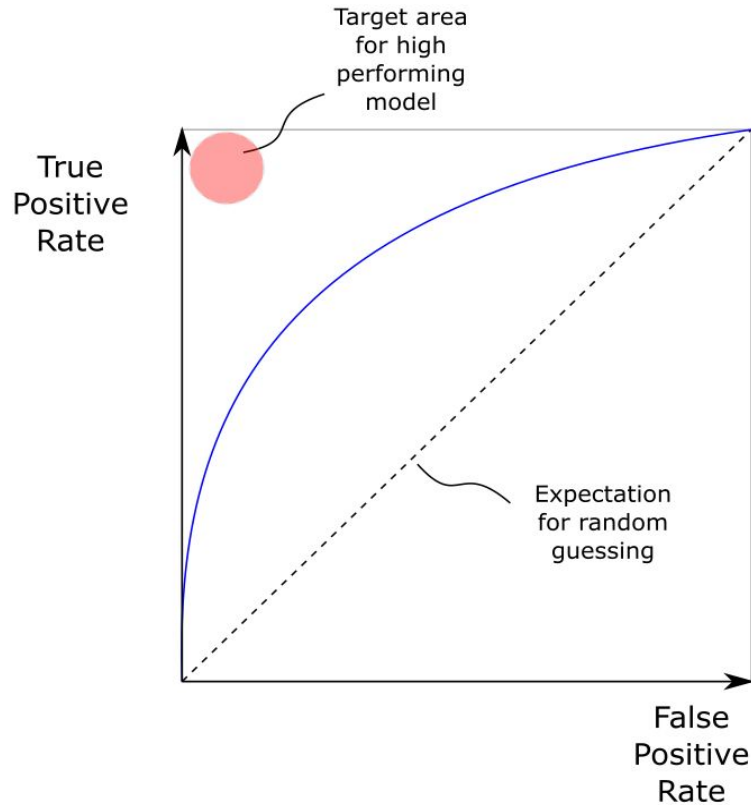
Confusion Matrix:

	pred: 0	pred: 1	pred: 2
true: 0	15	0	0
true: 1	0	16	2
true: 2	0	0	12

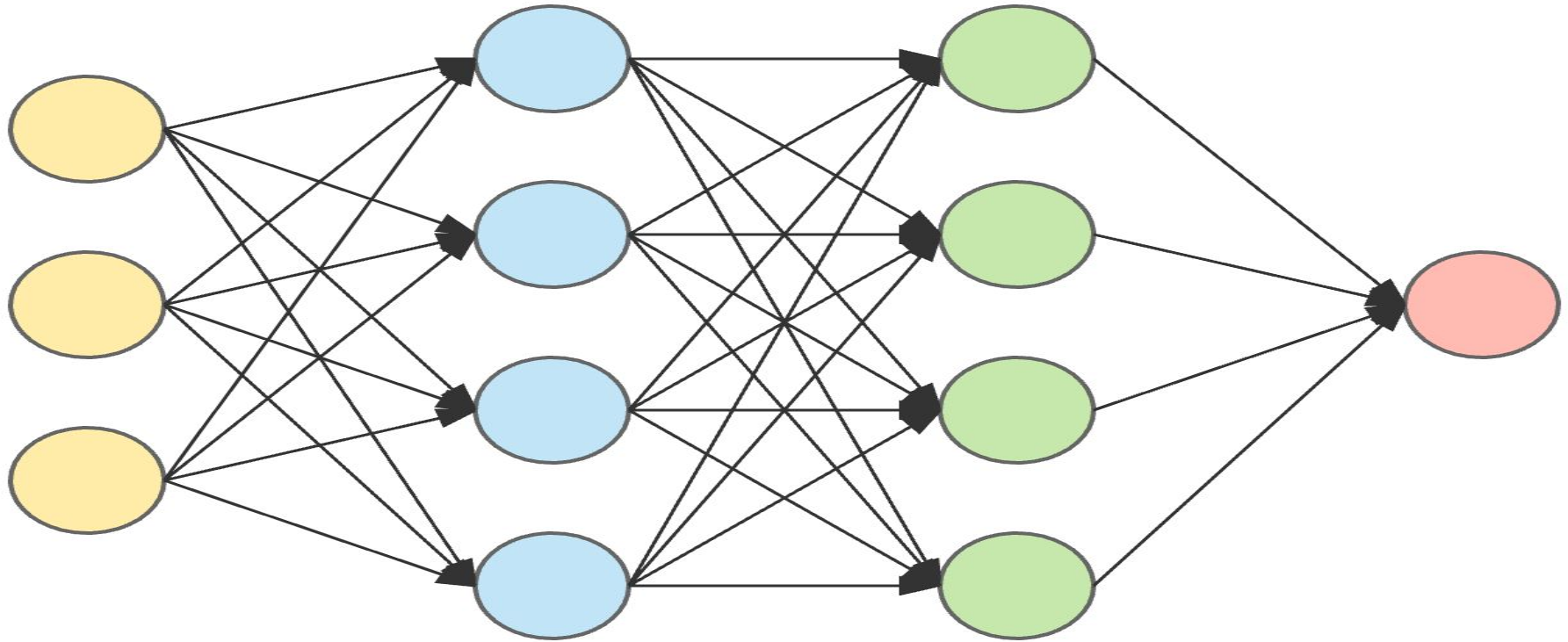
Classification Report:

	precision	recall	f1-score	support
0	1.00	1.00	1.00	15
1	1.00	0.89	0.94	18
2	0.86	1.00	0.92	12
avg / total	0.96	0.96	0.96	45

ROC curve



Advanced Methods



input layer

hidden layer 1

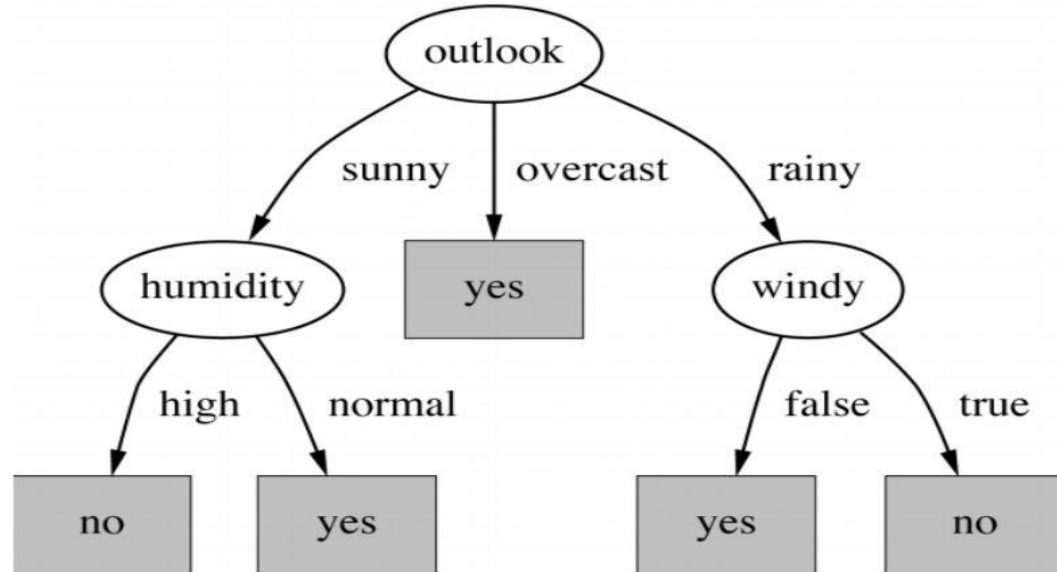
hidden layer 2

output layer

Decision tree background

- Some advanced methods can perform both regression and classification
- **Decision trees** attempt to split observations into distinct groupings based on feature values such that bin purity is maximized, ultimately creating a binary tree in the process
 - The predicted class will then be the maximally-represented class in the leaf nodes for classification, or an average of continuous values for a regression
- Decision trees are prone to overfitting, but the number of splits (e.g. tree depth) can be specified beforehand by the modeler
- **Non-parametric**, meaning it makes no assumption of the underlying form of the data

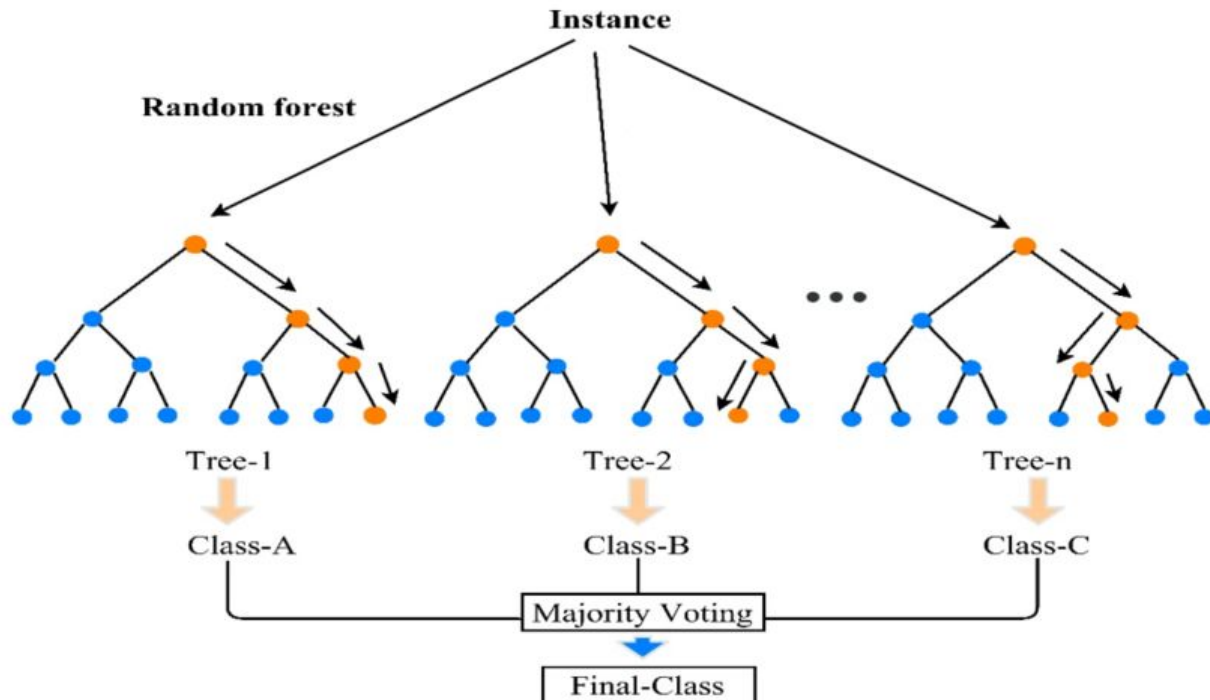
Decision tree (playing tennis)



Random forest background

- In order to reduce variance and chance of overfitting, a ***random forest*** is an average of the results of many decision trees splitting on random subsets of features
- Random forests are typically a ***robust*** model working well in many situations (including high-dimensional data), but they may lack interpretability as a result
- Other forms of tree-based learning include ***boosting*** and ***bagging*** (beyond the scope of this presentation, but worth looking into for those interested!)

Random forest example



https://www.researchgate.net/figure/Random-Forests-Naive-Bayes-NB-NB-approaches-are-a-family-of-simple-probabilistic_fig2_326722598

Random forest code example: Boston

```
from sklearn.ensemble import RandomForestRegressor

rf_reg = RandomForestRegressor(max_depth = 4)
rf_reg.fit(X_train, y_train)
y_pred_rf = rf_reg.predict(X_test)
print(rf_reg.feature_importances_)
print('\nR^2 Score:', rf_reg.score(X_test, y_test))
print('Mean squared error: ', mean_squared_error(y_test, y_pred_rf))
```

```
[0.06543473 0.00104719 0.00340049 0.          0.00376146 0.
 0.38850392 0.          0.00343326 0.08661245 0.          0.22357216
 0.22423434]
```

```
R^2 Score: 0.9059701492537314
Mean squared error: 0.056
```

Additional types of models

- We have only covered a small subset of modeling approaches that exist!
- Regression
 - Regularized linear regression (LASSO, Ridge, Elastic Net), polynomial regression/splines
- Classification
 - Naive Bayes, linear/quadratic discriminant analysis
- Regression/classification
 - Support vector machines, K-nearest neighbors, XGBoost
- Other approaches
 - Survival analysis, time series, structural equation modeling

Other types of modeling

- Unsupervised learning - infer y from the data (no explicit optimization criteria)
 - E.g. k -means clustering (split data in k distinct groupings, useful for market segmentation)
 - Also encompasses dimensionality reduction techniques like principal component analysis
- Deep learning - send input through a nested series of “hidden” activation functions (or layers) before reaching the output
 - E.g. neural networks (require lots of data, often used for language or vision related tasks)
- Reinforcement learning - choose the best actions in a system where many are possible in order to optimize some reward function over time
 - E.g. q-learning (useful for robotics and training AI to play video games)

Additional resources

- Linear regression:
https://scikit-learn.org/stable/modules/generated/sklearn.linear_model.LinearRegression.html
- Logistic regression:
https://scikit-learn.org/stable/modules/generated/sklearn.linear_model.LogisticRegression.html
- More methods and explanations:
https://scikit-learn.org/stable/tutorial/statistical_inference/supervised_learning.html

Exercises

Link to respective Jupyter Notebook / RMD files here