

# Machine Learning on HPC: Supervised Learning

Scalable Computational Intelligence Group  
Texas Advanced Computing Center  
University of Texas at Austin

Juliana Duncan  
June. 2023

# Workshop Series Overview

- There are five workshops this week:
  - Workshop 1: July 15, 2024
    - Introduction on computing cluster environment, Python for machine learning and containers
  - **Workshop 2: July 16, 2024**
    - Common supervised machine learning methods with examples and hands-on exercises
  - Workshop 3: July 17, 2024
    - Common unsupervised machine learning methods with examples and hands-on exercises
  - Workshop 4: July 18, 2024
    - Deep learning concepts, common neural network structure and frameworks with examples and hands-on exercises.
  - Workshop 5: July 19, 2024
    - Advanced ML topics
      - Introduction to Reinforcement Learning and LLMs

# Agenda of Today

- 9:00 - 11:00
  - Machine Learning Basics
  - Common Supervised Learning Methods
- 11:00-12:00
  - Introduction to Hands-on Exercises
- 1:00 - 2:30
  - More Common Supervised Learning Methods
- 2:30 - 4:00
  - Hands-on Exercises

# Outline

- Machine Learning (ML) Basics
  - Define ML
  - Categories of ML
  - ML Workflow
- Common Supervised Learning Techniques
  - K Nearest Neighbors
  - Linear Regression
  - Regularization
  - Logistic Regression
  - Support vector machines
  - Probabilistic Models
  - Decision Tree Based Models

# What is Machine Learning?

- A formal definition from T. Mitchell (1997). *Machine Learning book*:

“A computer program is said to learn from experience E with respect to some class of tasks T and performance measure P, if its performance at the tasks improves with the experiences.”

--- (Mitchell 1997)

HockeyIsLife	MayField	ShangriLaNF
Holiday Springs	Melanie	SmartyPants
Holstein	MicrosoftSansSerif	SnowShoe
<b>Hornswoggled</b>	MissyBT	SoftAutumn
Impress	MISTER SIRLOIN	SconeHenge
<b>Incised</b>	MisterEarl	SunnyDale
initial	Mistral	Susie'sHand
InkRunway	NoteThis	SuzanneQuilts
Jester	<b>NOTMARYKATE</b>	SweetheartScript
Jokerman	NSimSum	TropicalScript
Journalingland	Nupcial	TrottingKings
Juice	<b>OLITOCOMPOLI</b>	TYPEWRNG
<b>JUNIORSTAR</b>	Palatino Linotype	TypeUpright
Kids Scrawl	Papyrus	Unnamed Mbdy
KOMIKANDY	Perpetua	Verdana
Kristen	Perismon	VinerHand
LA Headlights	PeritScript	Vinque
LP CoffeeHouse	Puskin	Vivaldi
Leftovers	QuirkyWiggly	Vrinda
Liarah	RobotTeacher	VarDisney
LooseScript	Rodford	Wenceslas
LoveLetter	Ryan	Wenceslas
Malahini Cuban	SA Amy	whatHappened
Mary Jane Antique	Scrawny Kids	AMasterpiece

Recognize hand written words within images

- Task: Recognize hand written words within images
- Performance: percent of words classified correctly
- Experience: database of classified hand written words

# What is Machine Learning?

- A formal definition from T. Mitchell (1997). *Machine Learning book*:

“A computer program is said to learn from experience E with respect to some class of tasks T and performance measure P, if its performance at the tasks improves with the experiences.”

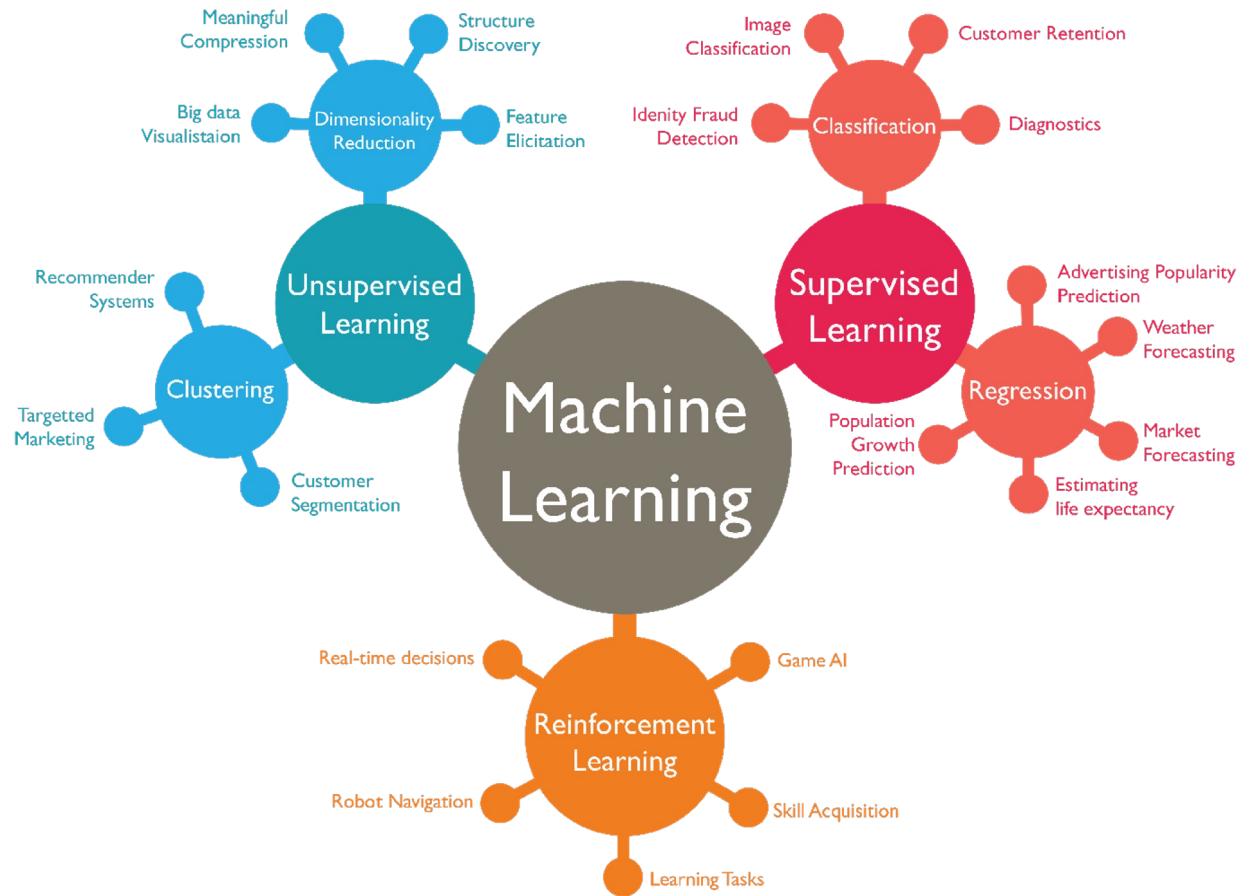
--- (Mitchell 1997)

Next let's take a deeper dive into:

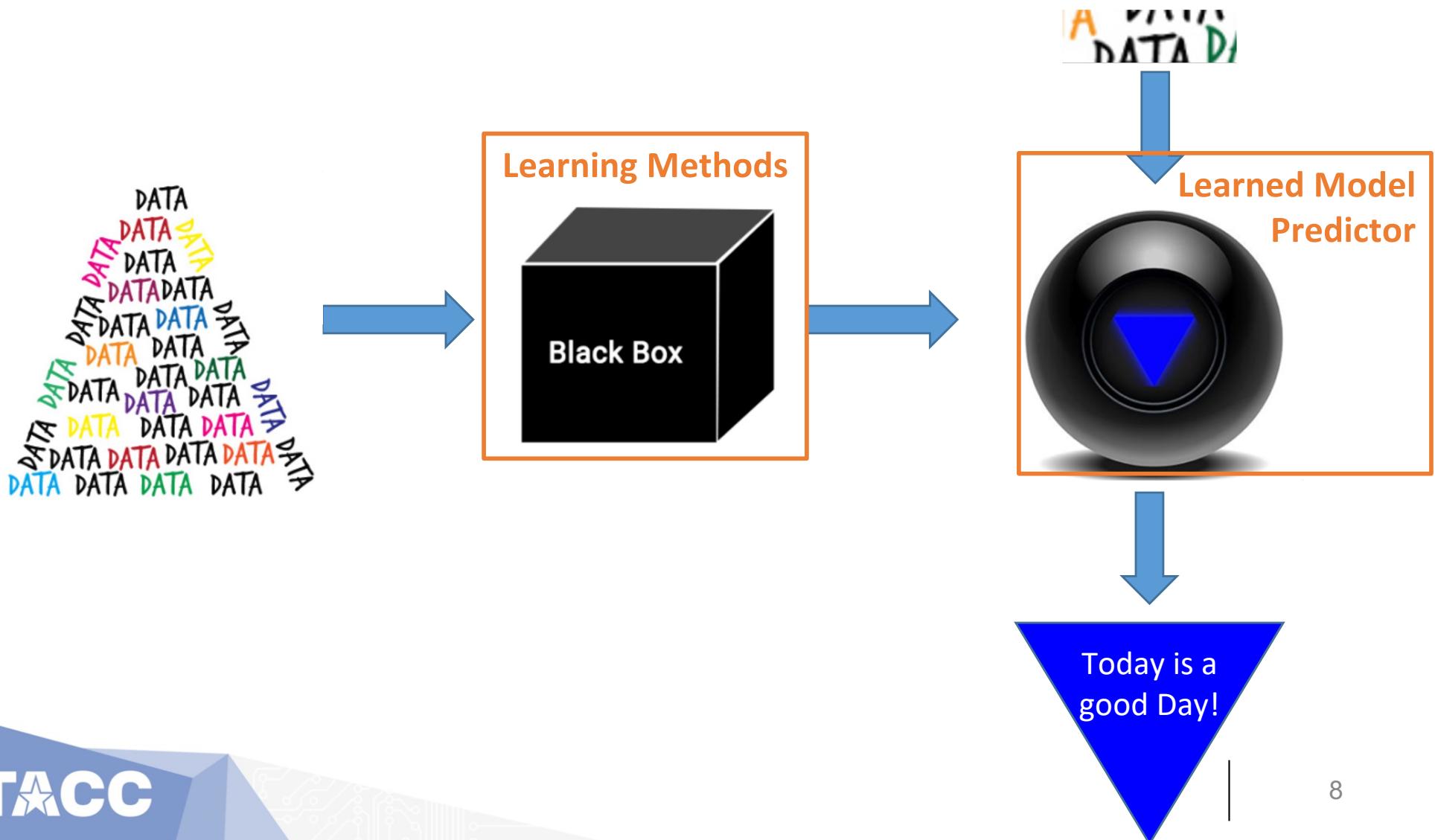
- Common ML tasks
- How do we gain experience and workflow for learning from experience
- Process for evaluating performance

# Common Task for “Learning”

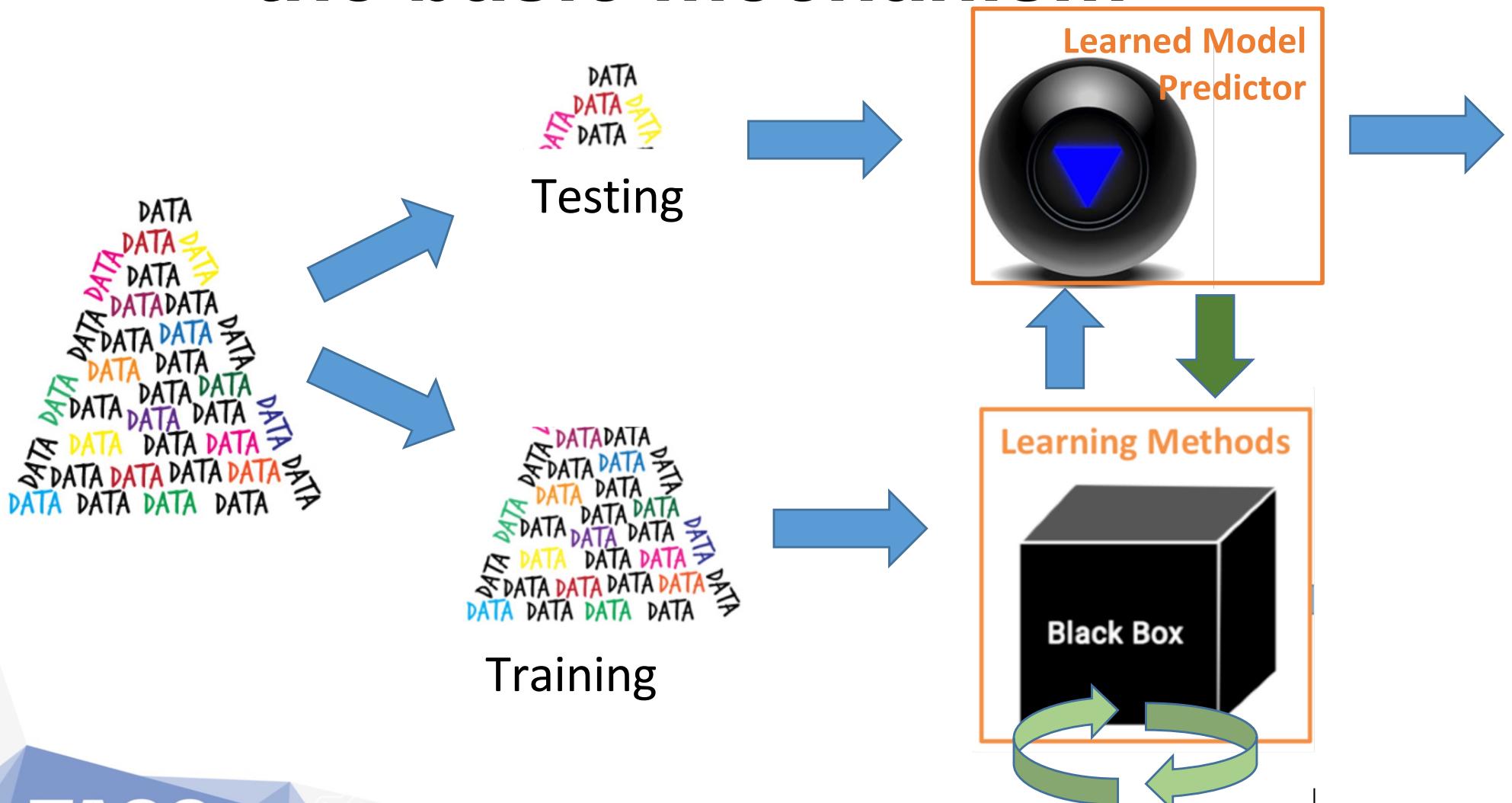
- Supervised Learning
  - Data has Labels
- Unsupervised Learning
  - Data has NO Labels
- Reinforcement Learning
  - Learn from making mistakes



# Learn from Past Experiences



# How to Evaluate Performance of Machine Learning? --- the basic mechanism



# Training and Testing

- **Training**
  - The process for making a model.
  - Training data is the data “observed” by the learning method
- **Testing**
  - Process for evaluating the performance of the model
  - Testing Data is data NOT observed by the learning model
- **80/20 split**
  - 80% available data are randomly selected for training
  - The rest 20% are used for testing.
- **Prediction**
  - Apply model to data not in either training or testing data set.
  - Assume the input data and its prediction are from the same process producing the training data.

# Performance Metrics

- Various Performance Metrics will be discussed throughout the this talk for Supervised Learning (Regression and Classification)
- Simple Example with Classification
  - Accuracy


$$\frac{\text{number of data points classified correctly}}{\text{total number of datapoints}}$$

# Example: tax evasion detection

Historical data records

Tid	Data attribute/ Features				Class Label Evade
	Refund	Marital Status	Taxable Income		
1	Yes	Single	125K	No	
2	No	Married	100K	No	
3	No	Single	70K	No	
4	Yes	Married	120K	No	
5	No	Divorced	95K	Yes	
6	No	Married	60K	No	
7	Yes	Divorced	220K	No	
8	No	Single	85K	Yes	
9	No	Married	75K	No	
10	No	Single	90K	Yes	

Record to be inferreded

Refund	Marital Status	Taxable Income	Evade
No	Married	85k	?

# Outline

- Machine Learning (ML) Basics
  - Define ML
  - Categories of ML
  - ML Workflow
- Common Supervised Learning Techniques
  - K Nearest Neighbors
  - Linear Regression
  - Regularization
  - Logistic Regression
  - Support vector machines
  - Probabilistic Models
  - Decision Tree Based Models

# **Classification: K Nearest Neighbor**

# Nearest Neighbor Classification

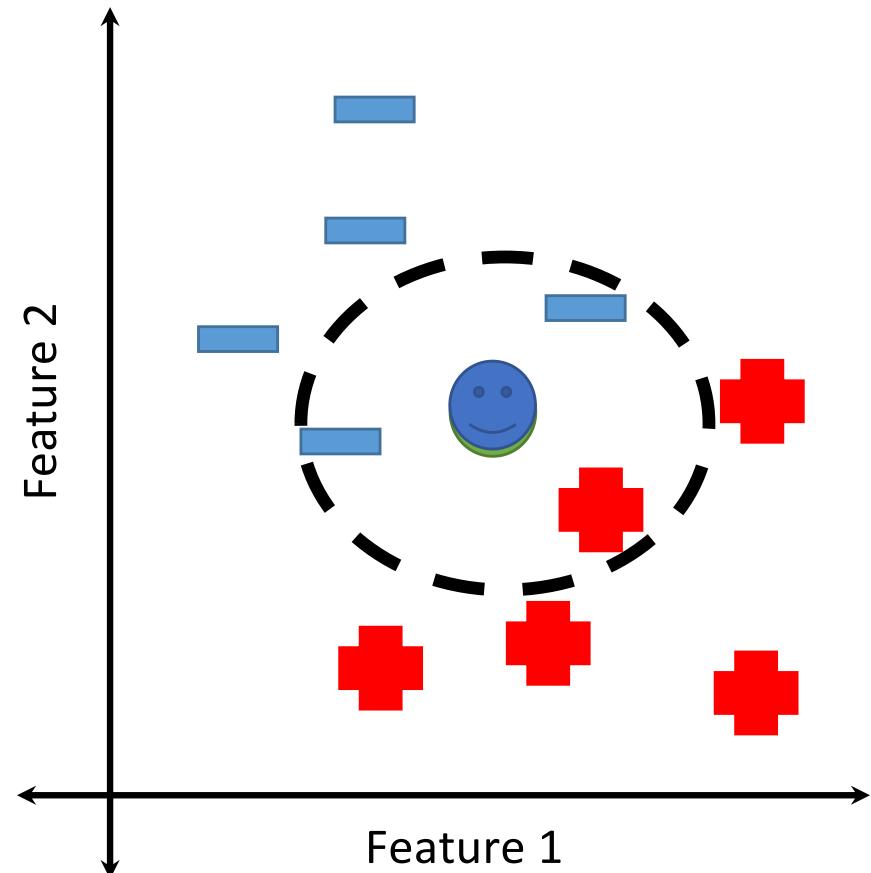
- The basic idea:

If two data objects are similar (close), they are likely from the same class.

- General workflow

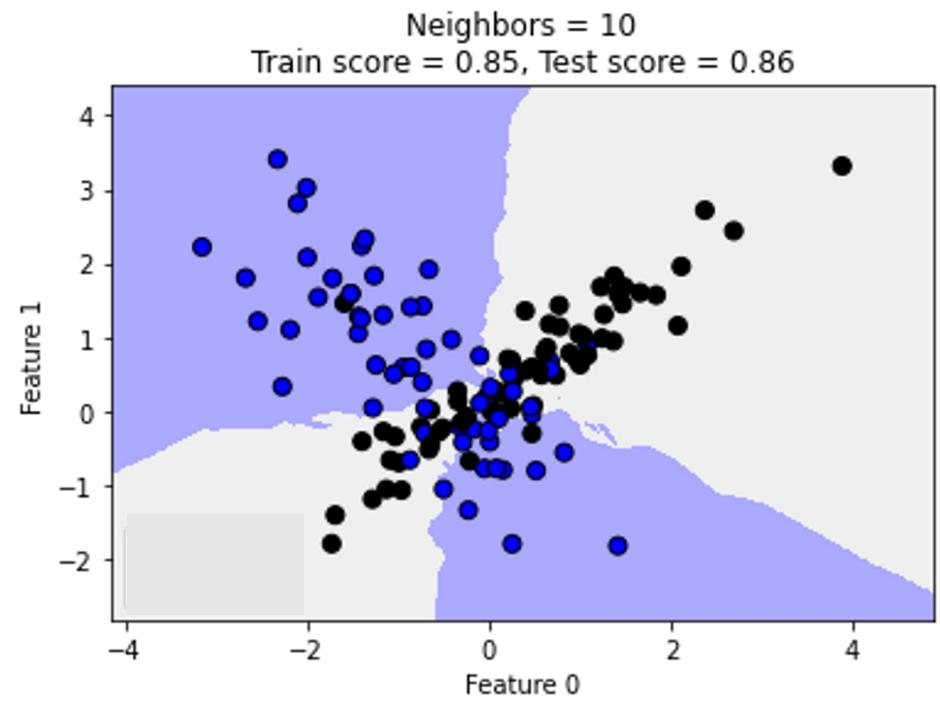
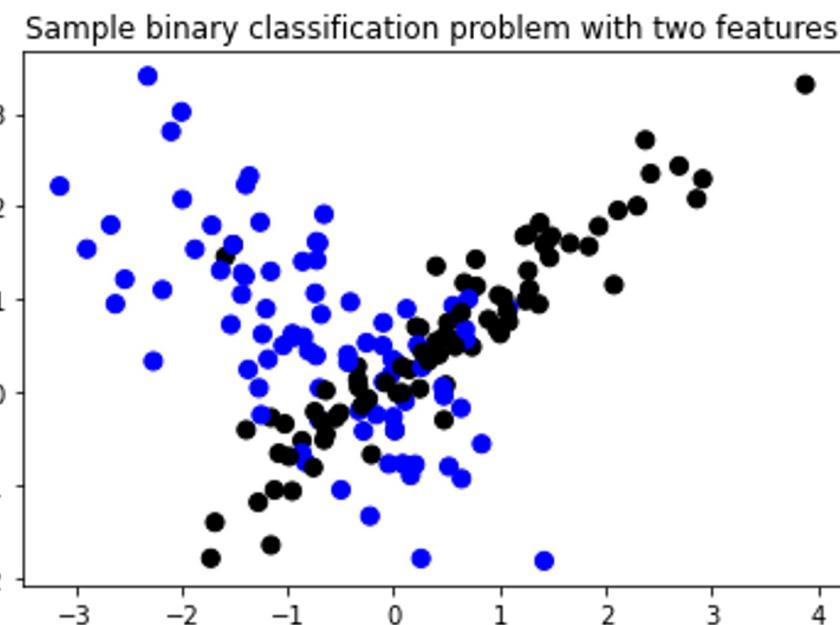
1. Start with a set of data object with class labels
2. Given a distance (similarity) measure of comparing data, compute distance from *unseen data point* to all data in training set
3. Retrieve (k) nearest neighbor for the unknown data
4. Assign class label based on the retrieved record.

- K is a hyperparameter
- Hyperparameters are parameters users select that will impact the performance of ML models



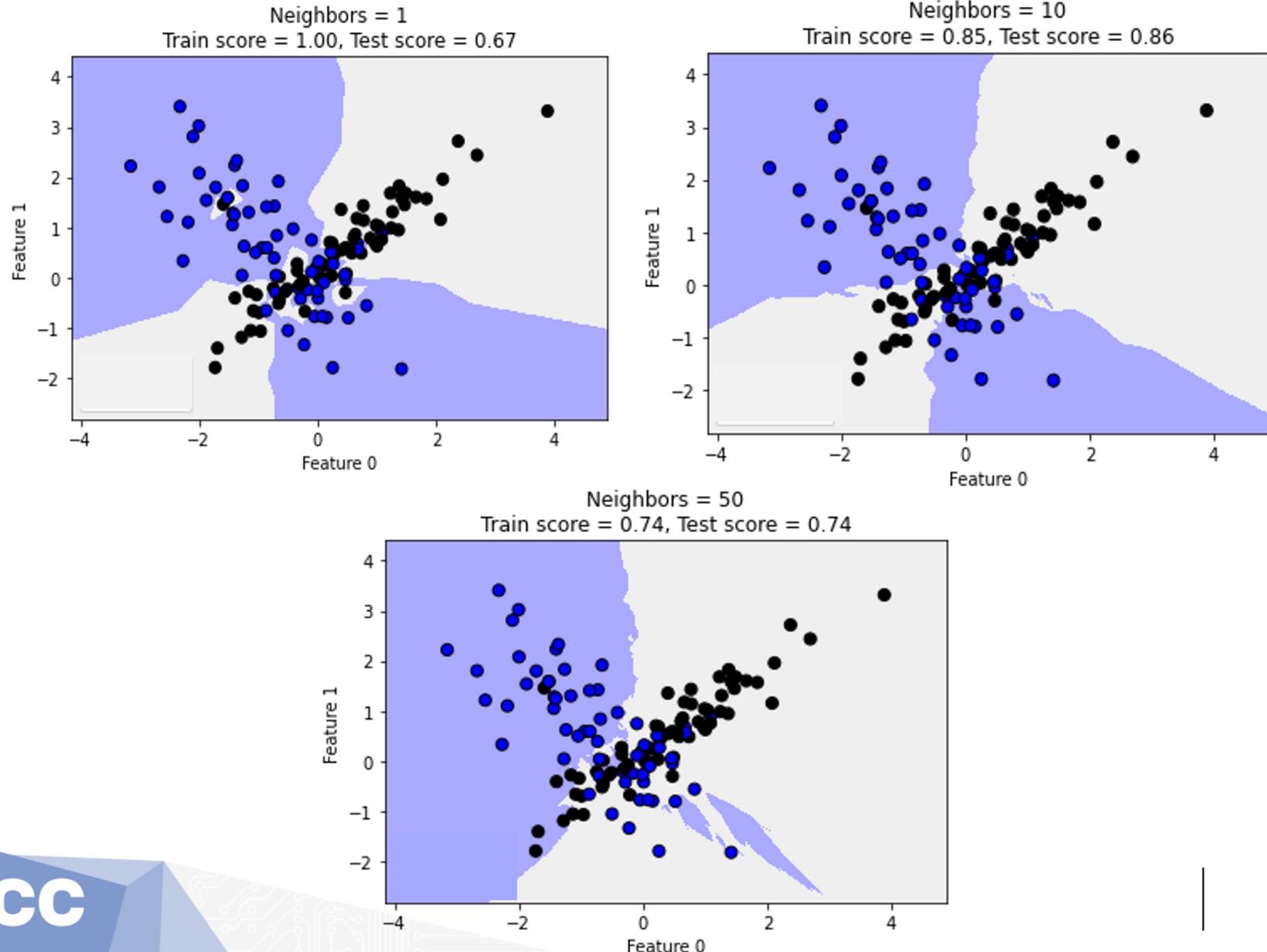
# Decision Boundary for Nearest Neighbor Classification

- A simple binary classification case



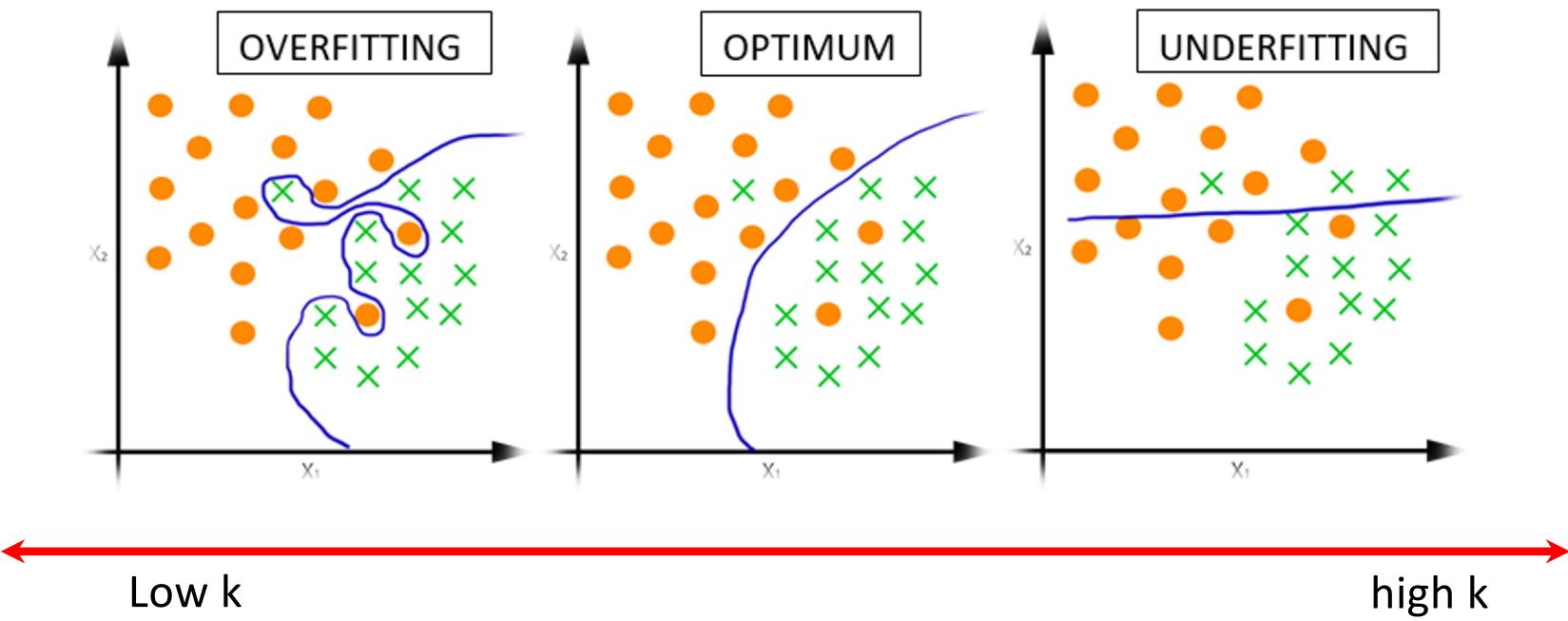
# Choice of K

- What changes?



# Underfitting vs Overfitting

- Underfitting,
  - model cannot reflect all the relations from training data
  - Model performs poorly on training and testing data. “failed to learn”
- Overfitting,
  - Model may lead to poor generalization on new testing data.
  - Performs very well on training data but poorly on testing data
- Hyperparameters impact under or over fitting



# Classification: Additional Evaluation Metrics

Let's say we built a model predicting whether a patient has cancer:

*Confusion Matrix*

	Predicted - Cancer	Predicted – NOT Cancer
Actual - Cancer	0	30
Actual – NOT Cancer	0	270

This is a bad model, but what is the accuracy ...

$$\text{Accuracy} = 270/300 = 90\%$$

Accuracy alone should not be used to evaluate a model's quality when the data is imbalanced.

# Confusion Matrix, Recall, and Precision

- True or False refers to whether we classified data point correctly
- Positive or Negative refers to the **predicted** class label

		Predicted - Positive	Predicted – Negative
Actual - Positive	True Positive (TP)	False Negative (FN)	
	False Positive (FP)	True Negative (TN)	

$$\text{Recall} = \frac{TP}{TP + FN}$$

$$\text{Precision} = \frac{TP}{TP + FP}$$

# Classification: Additional Evaluation Metrics

Let's say we built a model predicting whether a patient has cancer:

<i>Confusion Matrix</i>		
	Predicted - Cancer	Predicted – NOT Cancer
Actual - Cancer	0	30
Actual – NOT Cancer	0	270

An Alternative to this metric is Recall

Recall = Predicted-Cancer / Data Points that are Actually Cancer = 0%

# Classification: Additional Evaluation Metrics

Let's say we built a model predicting whether an email is spam:

<i>Confusion Matrix</i>		
	Predicted - Spam	Predicted – NOT Spam
Actual - Spam	30	90
Actual – NOT Spam	30	270

Precision may make more sense since we do not want to accuse legit emails of being spam (or both recall and precision!)

Precision = Predicted-Spam / Data Points that are Predicted Spam = 50%

# **Regression:** **Linear Regression**

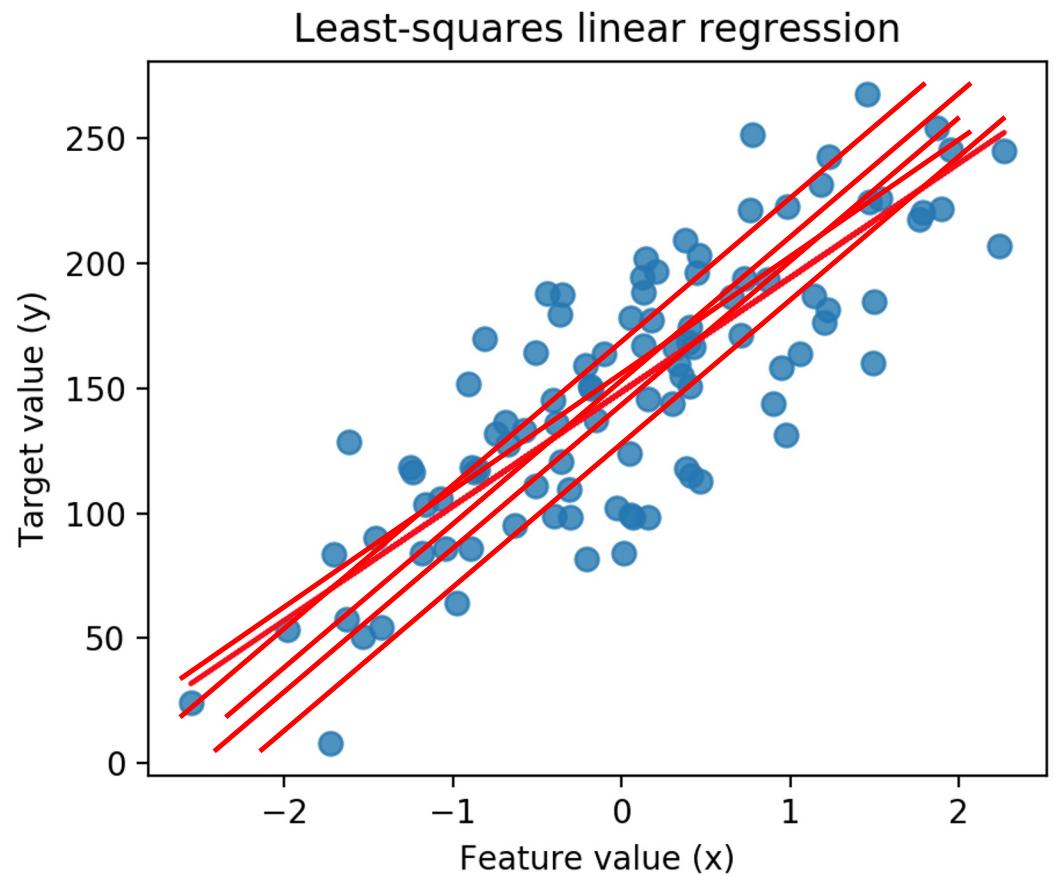
# Regression Model

- The independent variables  $X$
- The dependent variables  $Y$   
outcome, target, or criterion variable
- The unknown parameters  $\theta$
- Predict/estimate  $Y$  with  $f(X, \theta)$

Years experience	Salary
1.1	41,000
1.7	52,000
3	53,000
4.2	70,000

# Simple Linear Regression

- Only one independent variable ( $x$ ) and one dependent variable ( $y$ )
- Assuming label and feature are connected through a function e.g.  $y=ax+b$
- $\theta = \langle a, b \rangle$
- $f(x, \theta) = ax + b$



# Common Goal of the Regression

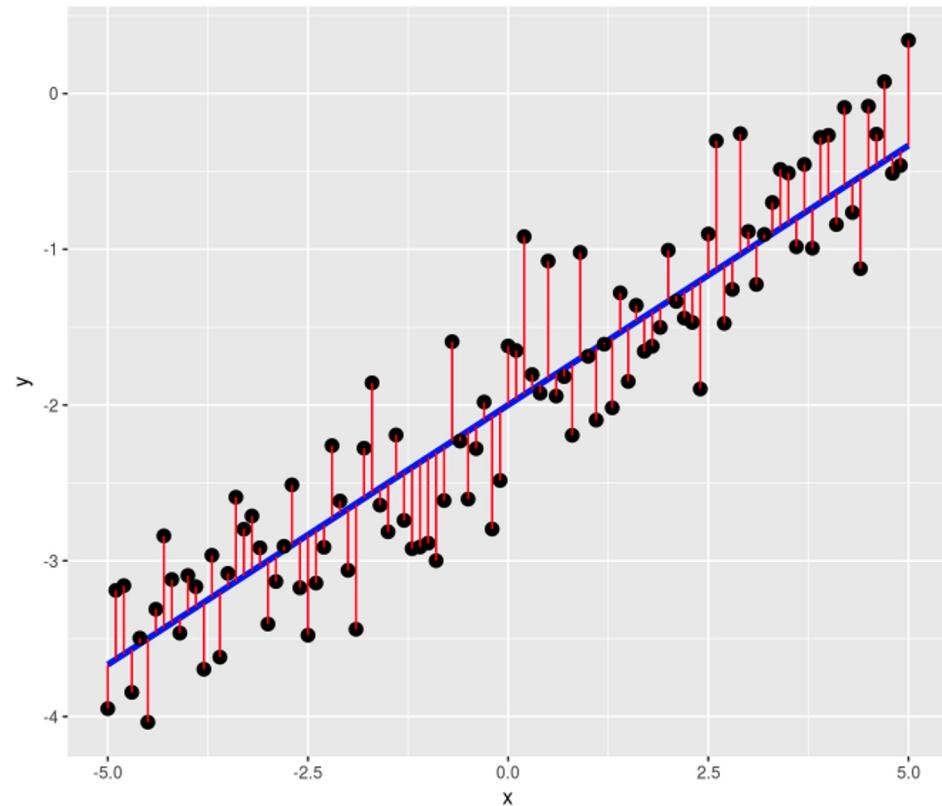
- Learn the parameters by minimize the cost or objective function

- Least absolute deviations: to minimize

$$\sum_{i=1}^N |f(x_i, \theta) - y_i|$$

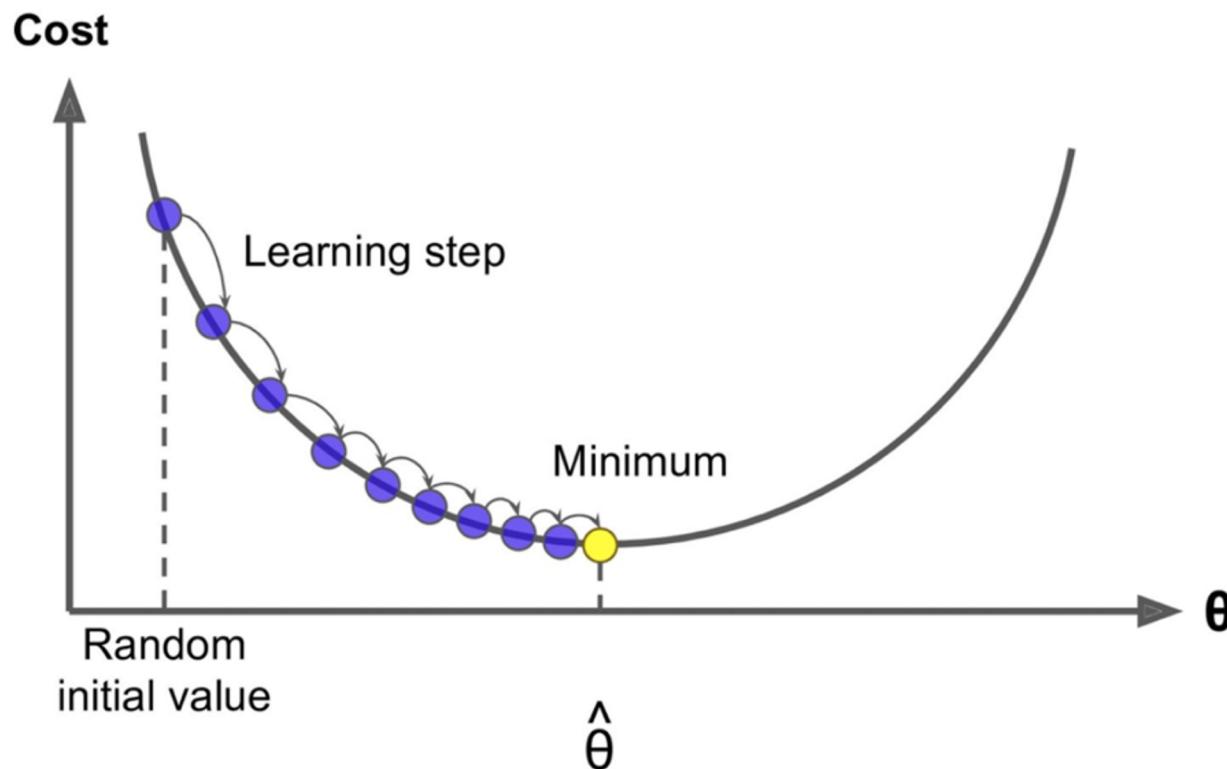
- Ordinary least squares (OLS): to minimize

$$\sum_{i=1}^N (f(x_i, \theta) - y_i)^2$$



# Training Models by Minimize “errors”

- A basic idea in many ML is to minimize loss function
- Numerical optimization methods like gradient descent are commonly used to find optimal values for  $\theta$

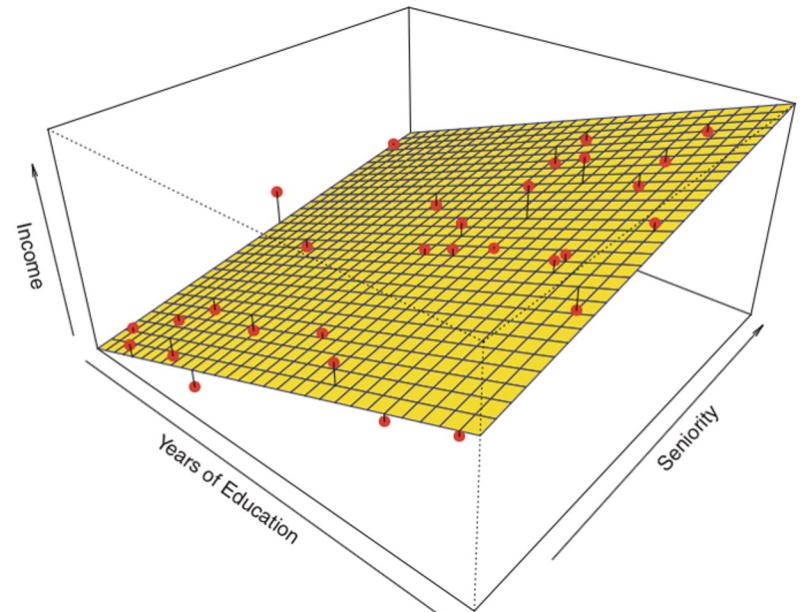


# Multiple Linear Regression

- More than one independent variable
- $x$  is a vector of features,  $x \in \mathbb{R}^N$
- $\theta$  is a vector of weights,  $\theta \in \mathbb{R}^N$
- $\theta = \langle b, a_0, a_1, \dots a_N \rangle$

$$f(x, \theta) = b + a_0x_0 + a_1x_1 + \dots + a_Nx_N$$

Years of Education	Years experience	Salary
13	1.1	41,000
17	1.7	52,000
19	3	53,000
22	4.2	70,000



# Common Metrics of Performance

- Mean Squared Error

$$\frac{1}{N} \sum_{i=1}^N (f(x_i, \theta) - y_i)^2$$

- Mean Absolute Error

$$\frac{1}{N} \sum_{i=1}^N |f(x_i, \theta) - y_i|$$

# Common Metrics of Performance

- From OLS, Residue sum of squares

$$SS_{\text{res}} = \sum_i (y_i - f_i)^2$$

- Total sum of square

$$SS_{\text{tot}} = \sum_i (y_i - \bar{y})^2.$$

- Coefficient of determination, aka,  $R^2$

- The proportion of the variance in the dependent variable that is predictable from the independent variable(s)

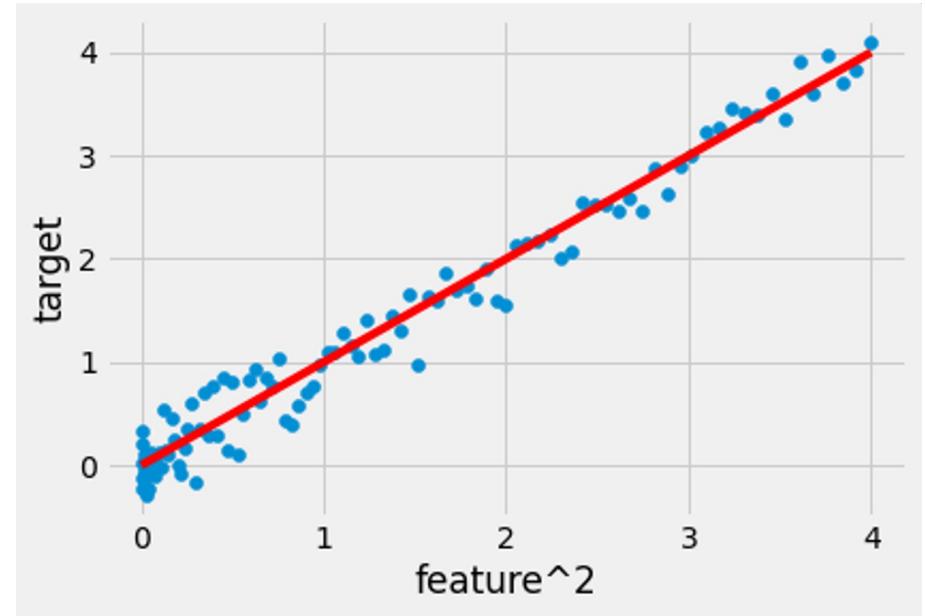
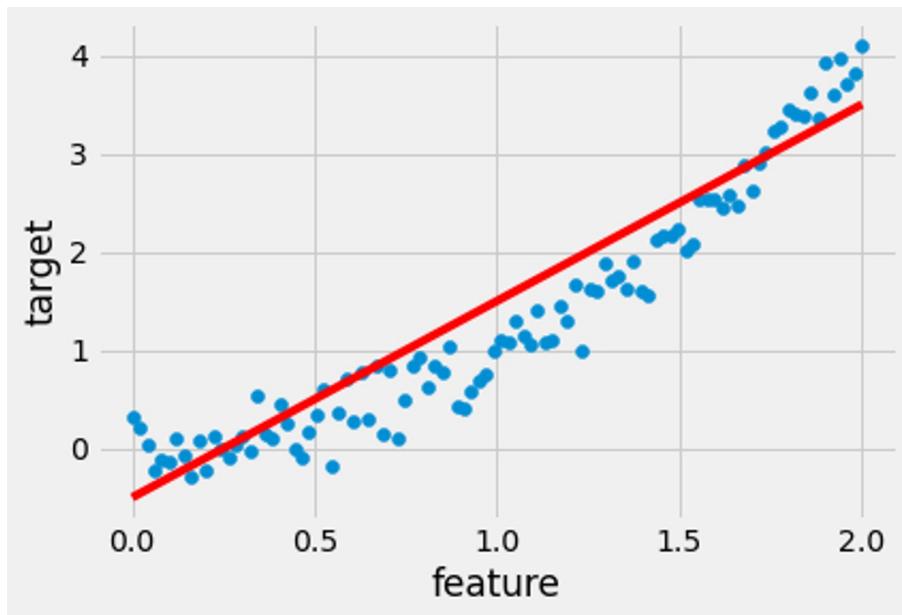
$$R^2 \equiv 1 - \frac{SS_{\text{res}}}{SS_{\text{tot}}}.$$

# Practical Considerations of Linear Regression

- Pros
  - Linear regression is simple
  - It is very interpretable
    - Slopes,  $a_i$ , can be defined as
    - Holding all else constant, for 1 unit increase in feature  $x_i$  the target variable increases by  $a_i$
- Cons
  - not all features in your data have linear relationship with the target
  - need to transform features such that a linear relationship exists to get good results

# Feature Engineering

You may need to change your features so that there is a linear relationship between the features and target

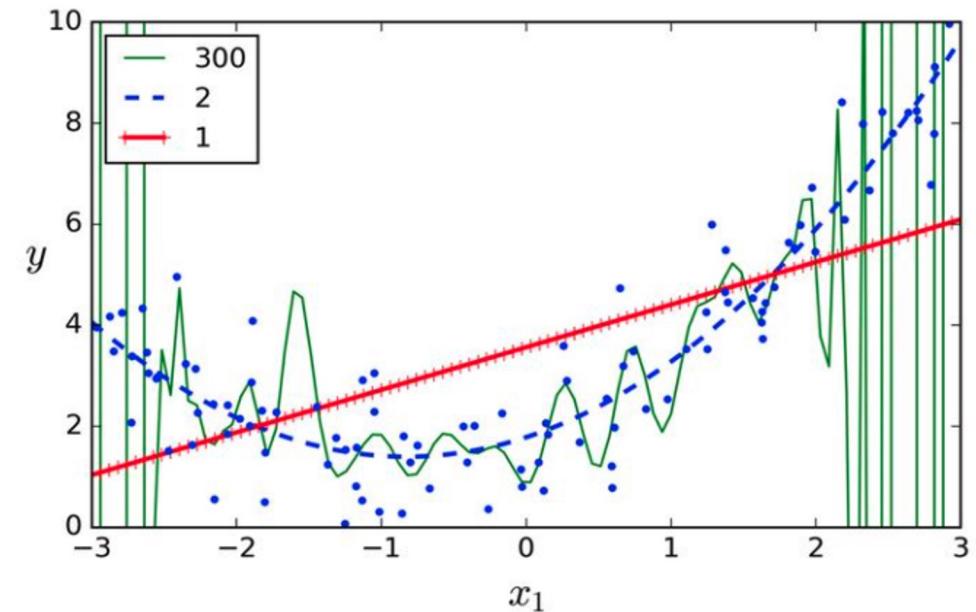


# Polynomial regression

- a special case of **multiple linear regression**

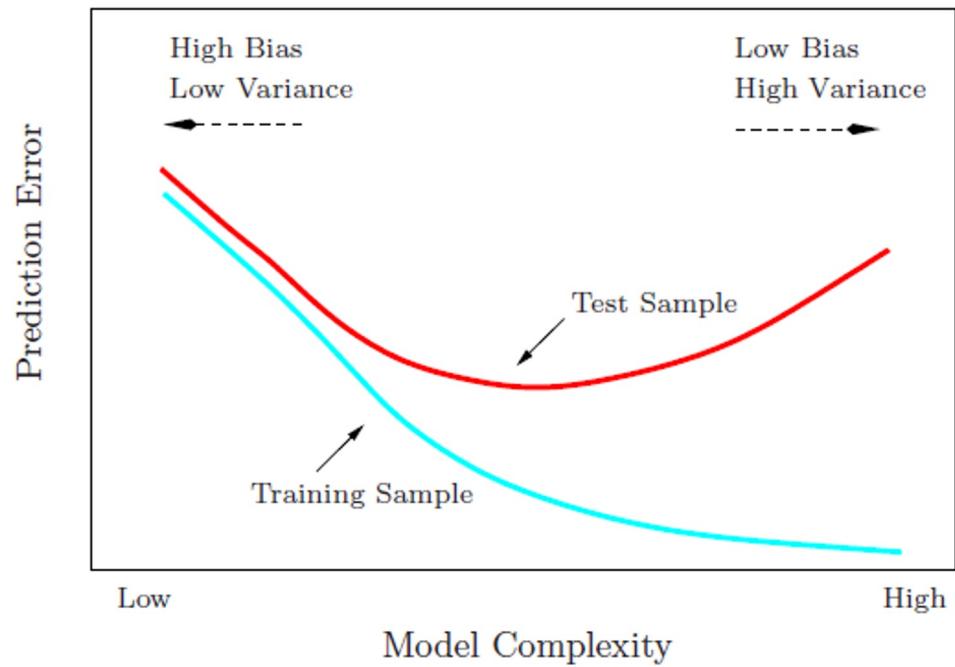
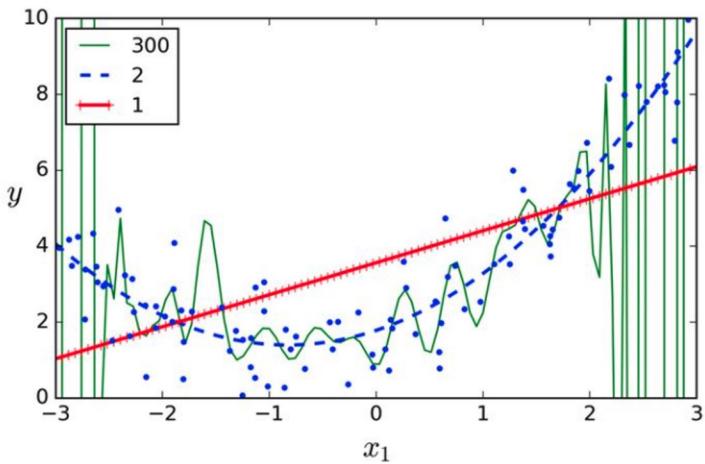
$$y = a_0 + a_1 x + a_2 x^2 + a_3 x^3 + \cdots + a_n x^n$$

- Which is better?
  - Red line: linear regression
  - Blue line: 2-degree polynomial regression
  - Green line: 300-degree polynomial regression



# Underfitting and Overfitting

- Model Complexity  $\approx$  number of parameters in a model
- High degree polynomial or using too many features can lead to over fitting
- How do we know what the appropriate polynomial degree is?



# Ridge Regression

- One way to combat this overfitting is to use **regularization**. An example is Ridge regression.
- Alters cost or objective function to avoid overfitting
- OLS tries to minimize

$$\sum_{i=1}^N (f(x_i, \theta) - y_i)^2$$

- Ridge regression adds a regularization term

$$\sum_{i=1}^N (f(x_i, \theta) - y_i)^2 + \alpha \sum_{k=1}^p a_k^2$$

- Where N is the number of training points, p is the number of features, and  $a_i$  is the slopes associated with feature i
- $\alpha$  – regularization hyperparameter – is a positive number, when  $\alpha = 0$  ridge regression is effectively the same as OLS.

# Standardizing Data

- When using Ridge regression (or the similar Lasso regression) it is important to standardize your data.
- Why? Let's say we are predicting MPG of vehicles given the features
  - weight of vehicle in lbs
  - number of cylinders
- Which beta coefficients will be larger?
  - even though weight is probably the most important, we may see that it has a smaller slope associated with due to the large magnitudes of weight of the vehicles.
- In Ridge regression, we want all our features beta coefficients to have a fair comparison
  - Solution: Standardize data!

# Standardizing Data

There are other ML algorithms where it is important to standardize data. For example,

- KNN
- Any time using a numerical optimizer of a loss function; speeds up gradient descent-like algorithms
- Any time adding regularization term to loss function

# Hands on Exercise:sklearn

- Free open source machine learning library for Python programming language
- Has implementations of all ML algorithms covered today



# Sklearn

```
1 # import needed objects
2 from sklearn.neighbors import KNeighborsClassifier
3 from sklearn.model_selection import train_test_split
4
5 # get data
6 data = get_data(....)
7 X = data.features
8 y = data.targets
9
10 # test train split
11 X_train, X_test, y_train, y_test = train_test_split(X, y, train_size=0.6, random_state = 0)
12
13 # instantiate then fit model with training data
14 knn = KNeighborsClassifier(n_neighbors = 10)
15 knn.fit(X_train, y_train)
16
17 # evaluate model
18 test_accuracy = knn.score(X_test, y_test)
19 train_accuracy = knn.score(X_train, y_train)
20
21 # make predictions |
22 y_pred=knn.predict(X_test)
```

# Sklearn

```
1 # import needed objects
2 from sklearn.neighbors import KNeighborsClassifier
3 from sklearn.model_selection import train_test_split
4
5 # get data
6 data = get_data(....)
7 X = data.features
8 y = data.targets
9
10 # test train split
11 X_train, X_test, y_train, y_test = train_test_split(X, y, train_size=0.6, random_state = 0)
12
13 # instantiate then fit model with training data
14 knn = KNeighborsClassifier(n_neighbors = 10)
15 knn.fit(X_train, y_train)
16
17 # evaluate model
18 test_accuracy = knn.score(X_test, y_test)
19 train_accuracy = knn.score(X_train, y_train)
20
21 # make predictions |
22 y_pred=knn.predict(X_test)
```

# Sklearn

```
1 # import needed objects
2 from sklearn.neighbors import KNeighborsClassifier
3 from sklearn.model_selection import train_test_split
4
5 # get data
6 data = get_data(....)
7 X = data.features
8 y = data.targets
9
10 # test train split
11 X_train, X_test, y_train, y_test = train_test_split(X, y, train_size=0.6, random_state = 0)
12
13 # instantiate then fit model with training data
14 knn = KNeighborsClassifier(n_neighbors = 10)
15 knn.fit(X_train, y_train)
16
17 # evaluate model
18 test_accuracy = knn.score(X_test, y_test)
19 train_accuracy = knn.score(X_train, y_train)
20
21 # make predictions |
22 y_pred=knn.predict(X_test)
```

# Sklearn

```
1 # import needed objects
2 from sklearn.neighbors import KNeighborsClassifier
3 from sklearn.model_selection import train_test_split
4
5 # get data
6 data = get_data(....)
7 X = data.features
8 y = data.targets
9
10 # test train split
11 X_train, X_test, y_train, y_test = train_test_split(X, y, train_size=0.6, random_state = 0)
12
13 # instantiate then fit model with training data
14 knn = KNeighborsClassifier(n_neighbors = 10)
15 knn.fit(X_train, y_train)
16
17 # evaluate model
18 test_accuracy = knn.score(X_test, y_test)
19 train_accuracy = knn.score(X_train, y_train)
20
21 # make predictions |
22 y_pred=knn.predict(X_test)
```

# Sklearn

```
1 # import needed objects
2 from sklearn.neighbors import KNeighborsClassifier
3 from sklearn.model_selection import train_test_split
4
5 # get data
6 data = get_data(....)
7 X = data.features
8 y = data.targets
9
10 # test train split
11 X_train, X_test, y_train, y_test = train_test_split(X, y, train_size=0.6, random_state = 0)
12
13 # instantiate then fit model with training data
14 knn = KNeighborsClassifier(n_neighbors = 10)
15 knn.fit(X_train, y_train)
16
17 # evaluate model
18 test_accuracy = knn.score(X_test, y_test)
19 train_accuracy = knn.score(X_train, y_train)
20
21 # make predictions |
22 y_pred=knn.predict(X_test)
```

# Sklearn

```
1 # import needed objects
2 from sklearn.neighbors import KNeighborsClassifier
3 from sklearn.model_selection import train_test_split
4
5 # get data
6 data = get_data(....)
7 X = data.features
8 y = data.targets
9
10 # test train split
11 X_train, X_test, y_train, y_test = train_test_split(X, y, train_size=0.6, random_state = 0)
12
13 # instantiate then fit model with training data
14 knn = KNeighborsClassifier(n_neighbors = 10)
15 knn.fit(X_train, y_train)
16
17 # evaluate model
18 test_accuracy = knn.score(X_test, y_test)
19 train_accuracy = knn.score(X_train, y_train)
20
21 # make predictions |
22 y_pred=knn.predict(X_test)
```

# Jupyter notebook

- TACC Analysis Portal:

Go to <https://vis.tacc.utexas.edu>; Login with your training account credentials

**TACC** | Analysis Portal [User Guide](#)

**jrduncan** [Log Out](#)

### Submit New Job

System: ---

Application: Select System

Project: Select System

Queue: Select System

Nodes: 1 Tasks: 1

---

### Options

Job Name: 20 characters max

Time Limit: H:M:S (default 2:0:0)

Reservation: reservation name

VNC Desktop Resolution: WIDTHxHEIGHT

[Submit](#) [Utilities](#)

### System Status

System	Status	Utilization	Job Count
Frontera	Open	99%	Running: 312 Queued: 1251
Lonestar6	Open	69%	Running: 135 Queued: 98
Longhorn	Open	74%	Running: 30 Queued: 40
Maverick2	Open	16%	Running: 4 Queued: 7
Stampede2	Open	96%	Running: 830 Queued: 736

### Past Jobs

JNB-Frontera	03/18/2022	<a href="#">Details</a>	<a href="#">Resubmit</a>
JNB-Frontera	03/18/2022	<a href="#">Details</a>	<a href="#">Resubmit</a>
JNB-Frontera	03/18/2022	<a href="#">Details</a>	<a href="#">Resubmit</a>
JNB-Frontera	03/18/2022	<a href="#">Details</a>	<a href="#">Resubmit</a>
JNB-Frontera	02/21/2022	<a href="#">Details</a>	<a href="#">Resubmit</a>

TACC

46

# Jupyter notebook

- TACC Visualization Portal:  
Select the following options

**Submit New Job**

System	Frontera
Application	Jupyter notebook
Project	Frontera training
Queue	small

Nodes: 1 Tasks: 1

---

**Options**

Job Name	20 characters max
Time Limit	H:M:S (default 2:0:0)
Reservation	ML_Institute_Tue
VNC Desktop Resolution	WIDTHxHEIGHT

**TACC**

Submit Utilities

# Jupyter notebook

## TAP Job Status

**Job:** Jupyter notebook on Frontera (4175197, 2022-03-21T17:28-05:00)

**Status:** RUNNING

**Start:** March 21, 2022, 5:28 p.m.

**End:** March 21, 2022, 5:33 p.m.

**Refresh:** in 873 seconds

### Message:

```
TAP: Your session is running at https://frontera.tacc.utexas.edu:60752  
/?token=9cbad0f26752e7dd14fcf090d6a30b6ec5c15c63ed7d9e2b626f214712fb8b4d
```

[Connect](#)

[End Job](#)

[Show Output](#)

[Back to Jobs](#)

# Morning Exercise

- You can find today's hand on exercises in `ml_institute_summer_24/day2`
- In this directory is a jupyter notebook which will use KNN, linear regression, and the algorithms that will be introduced this afternoon in the notebooks called **SupervisedLearning.ipynb**
- This notebook contains several optional exercises

# Hands-on Exercise / Break

# Outline

- Machine Learning (ML) Basics
  - Define ML
  - Categories of ML
  - ML Workflow
- Common Supervised Learning Techniques
  - K Nearest Neighbors
  - Linear Regression
  - Regularization
  - Logistic Regression
  - Support vector machines
  - Probabilistic Models
  - Decision Tree Based Models

# Logistic Regression

# (Binary) Logistic Regression

- Used for classification problems with a (Binary) dependent variable
- "y" can only take two values (usually 0 and 1)
  - Positive/Negative:
  - Pass/fail
  - Healthy/sick
- One or more independent variables x
- Predict the probability of a binary response

# How Logistic Regression Works

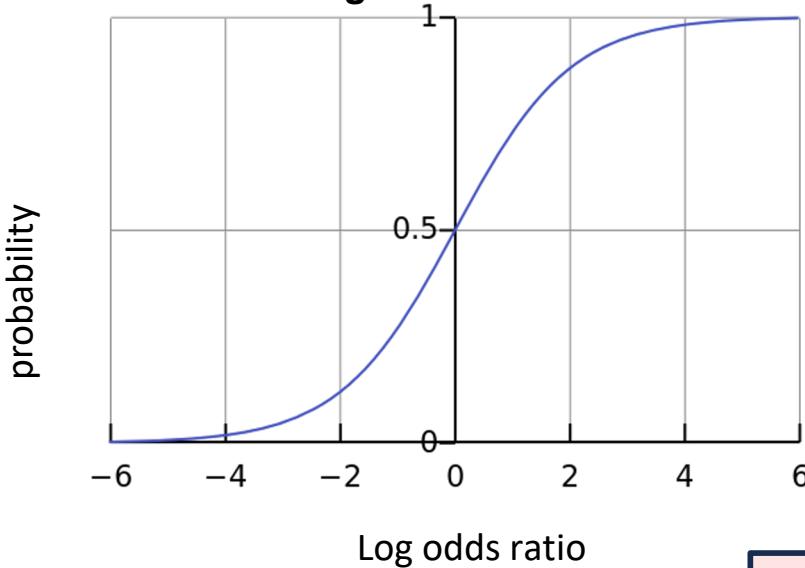
Linear Regression for classification?


$$pr(y = f(x, \theta) = b + a_0x_0 + a_1x_1 + \dots + a_Nx_N)$$

Instead linear model predicts the log odds ratio

$$t = \ln\left(\frac{p}{1-p}\right) = f(x, \theta) = b + a_0x_0 + a_1x_1 + \dots + a_Nx_N$$

**Logistic Function**



**Logistic Function**

Converts log odds ratio to probabilities

$$\sigma(t) = \frac{1}{1 + e^{-t}}$$

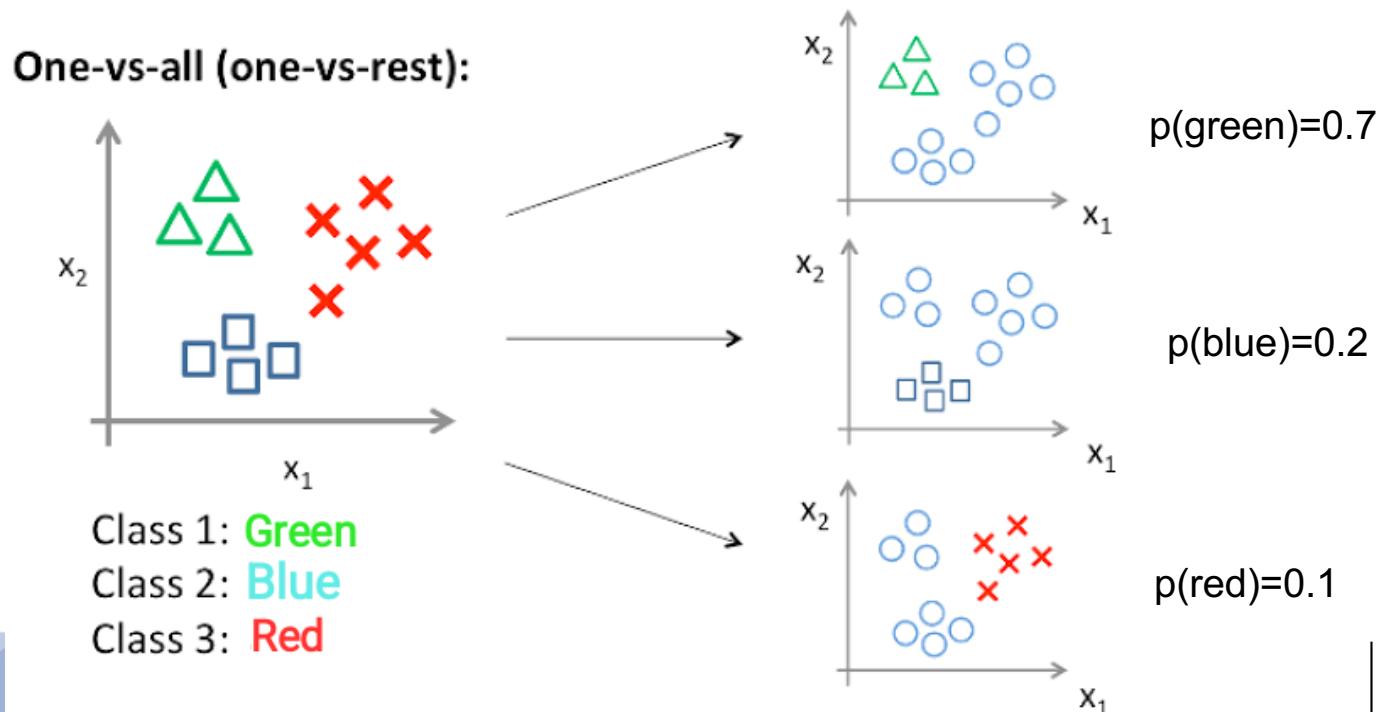
**Loss Function for Classification**

$$loss = \sum_{i=1}^n y_i \log(p_i) + (1 - y_i) \log(1 - p_i)$$

$$p(x, \theta) = \sigma(b + a_0x_0 + a_1x_1 + \dots + a_Nx_N)$$

# Multiclass Logistic Regression

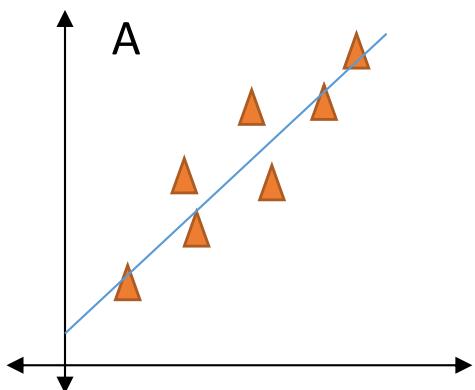
- Logistic regression can be applied to multiclass classification problems (more than 2 categories (i.e. green, blue, red))
- **One-v-rest:**
  - Build binary logistic regression for each class label vs all other labels; this produces probability of that label
  - Selects highest confidence model



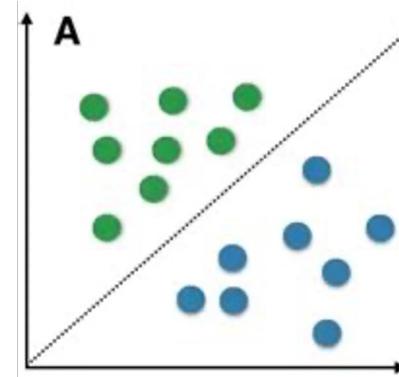
# Practical Usage

- Linear regression and Logistic regression create linear models (A below)

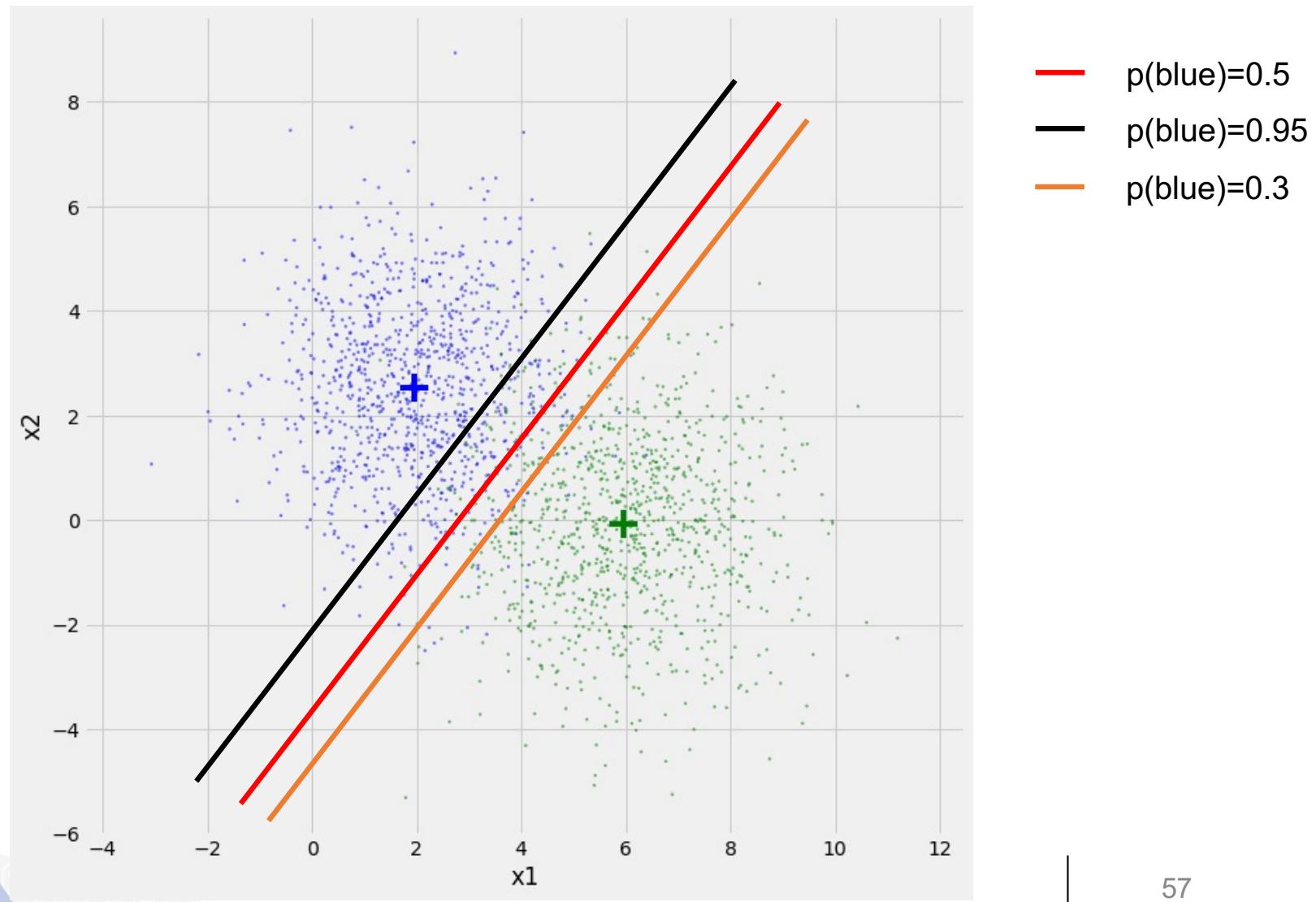
*Regression*



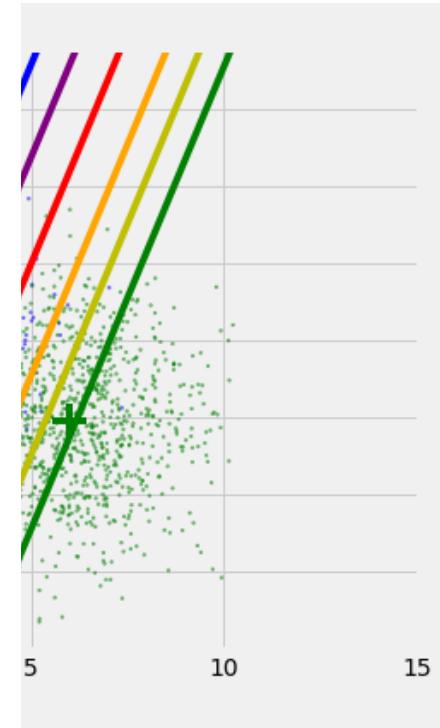
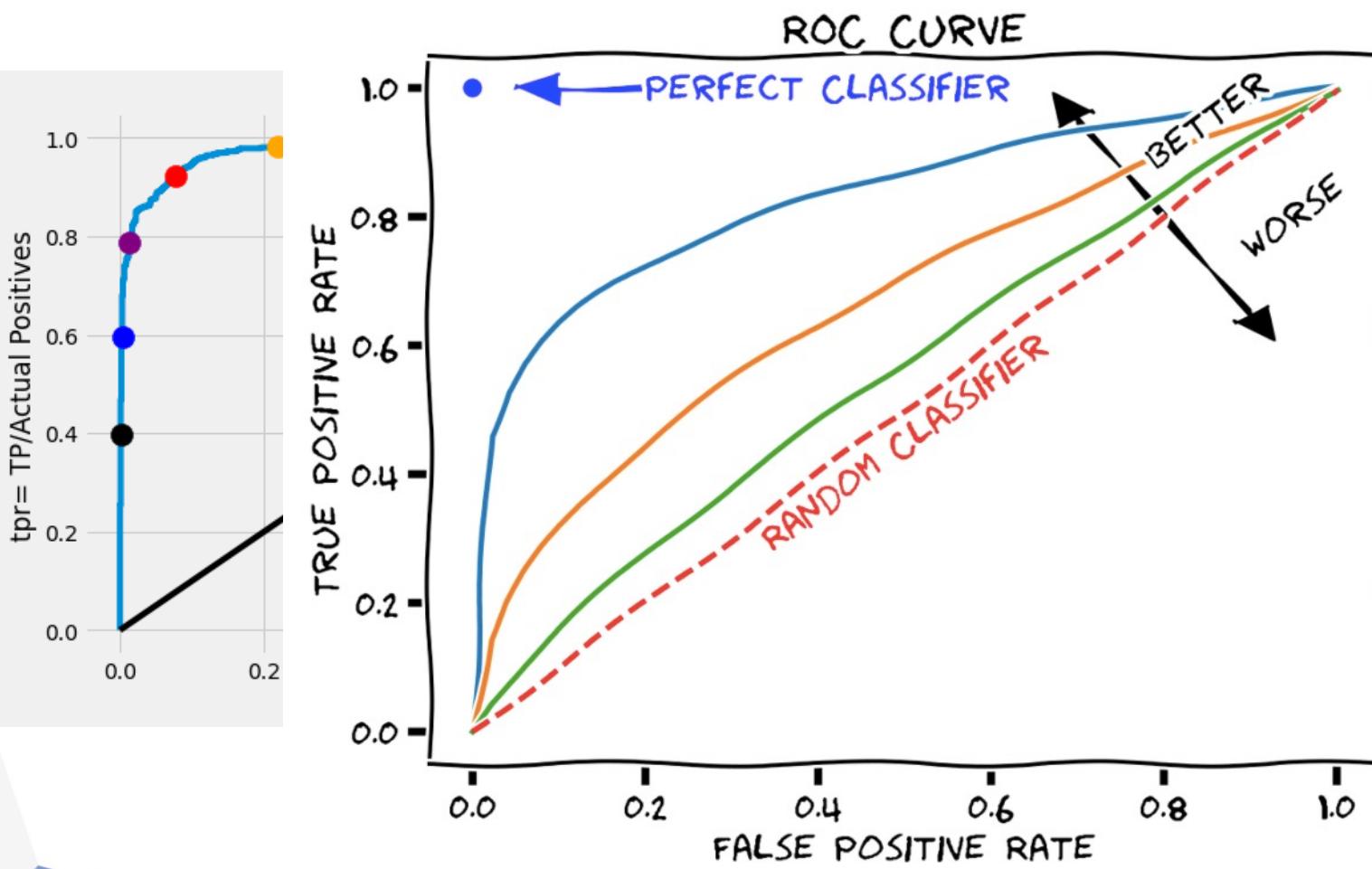
*Classification*



# Thresholding and ROC curves



# Model Evaluation: ROC curves



# Pros and Cons of Linear Models

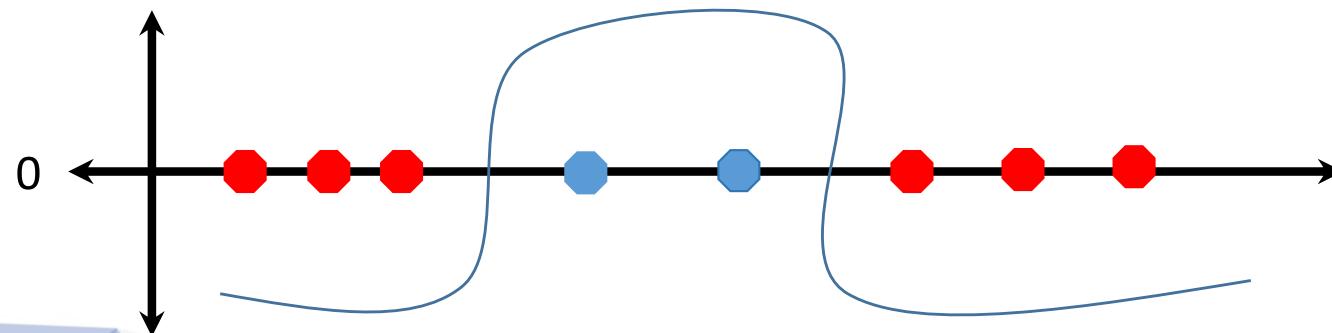
- Pros
  - Simple
  - Interpretable
    - Holding all else constant, for 1 unit increase in feature  $x_i$  the response variable  $y$  increases by  $a_i$
  - If linear relationship exist, use linear models
- Cons
  - Very Difficult to engineer linear relationships

# Outline

- Machine Learning (ML) Basics
  - Define ML
  - Categories of ML
  - ML Workflow
- Common Supervised Learning Techniques
  - K Nearest Neighbors
  - Linear Regression
  - Logistic Regression
  - Support vector machines
  - Probabilistic Models
  - Decision Tree Based Models

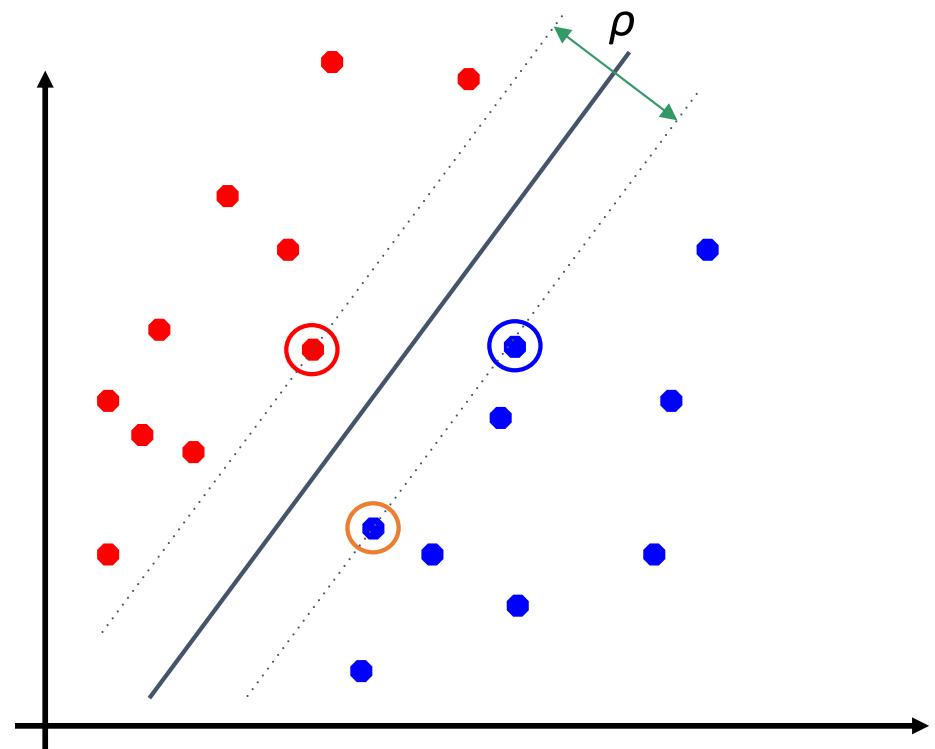
# Support Vector Machine (SVM)

- Used for classification
- Given labeled training data, SVM generates an optimal hyperplane to separate classes
- Creates a decision function



# Binary Classification with Linear Separator

- A line is used as the hyperplane separating two classes
- The closets data to the hyperplane are the **support vectors**
- Hyperplane is selected such that the margin is as big as possible
- A line is bad if it passes too close to some points (noise sensitive)



# Linear Support Vector Machine

**Decision Function**

$$\mathbf{w}^T \mathbf{x} + b$$

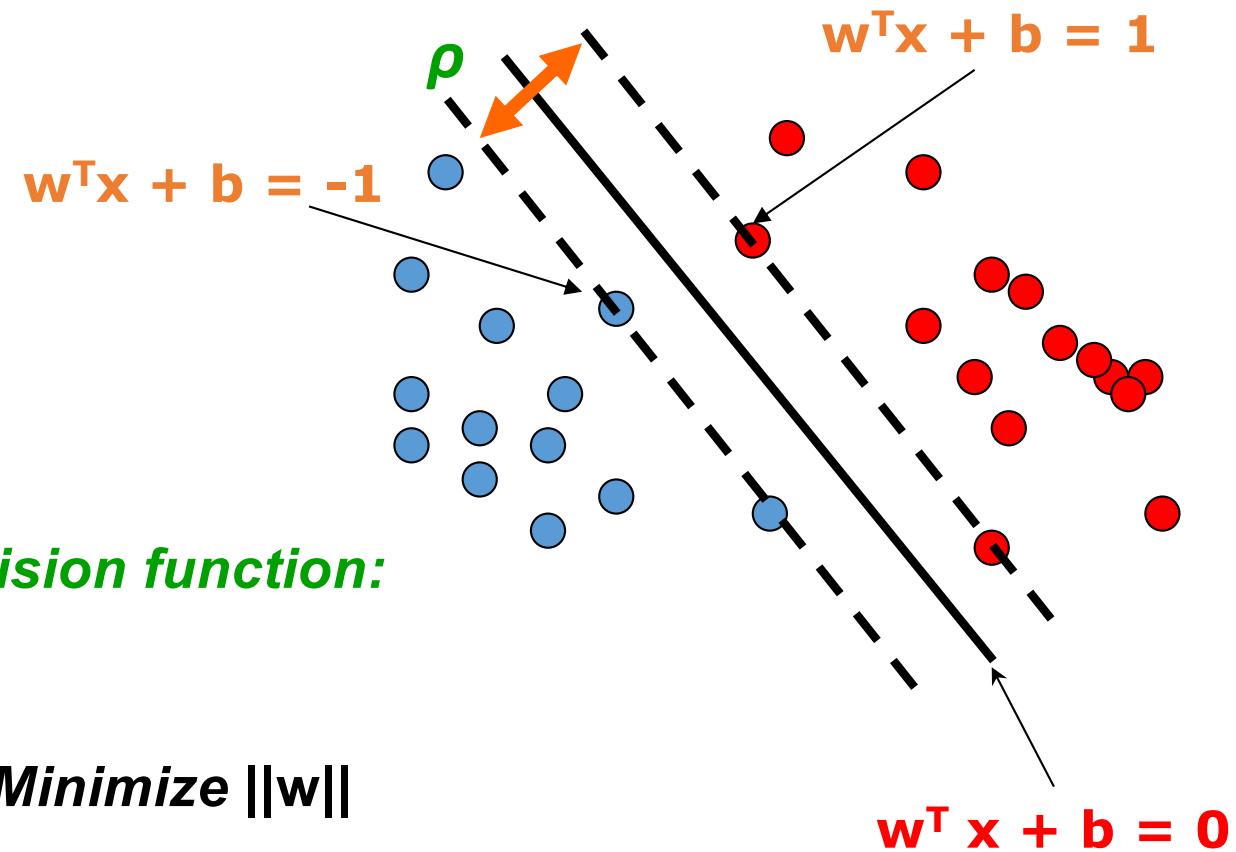
**Hyperplane**

$$\mathbf{w}^T \mathbf{x} + b = 0$$

*Given constraints to decision function:*

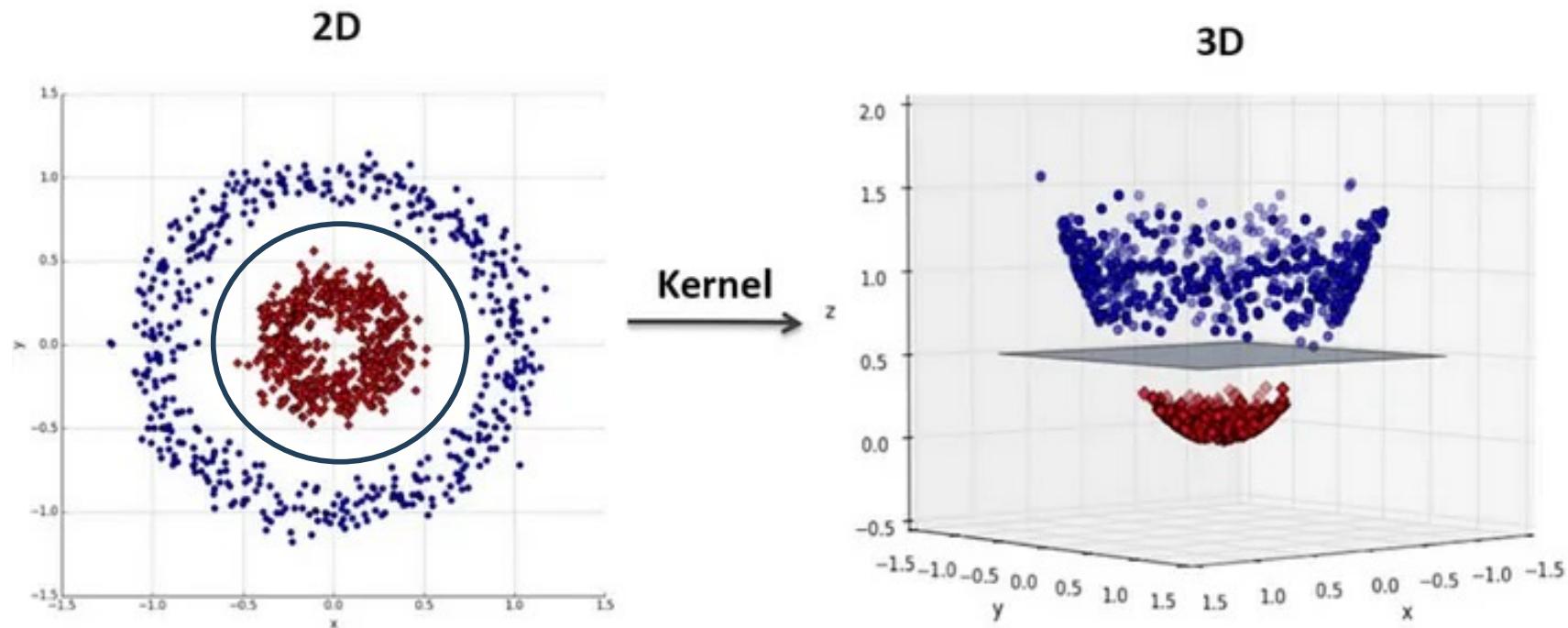
$$\rho = 2/\|\mathbf{w}\|_2$$

Maximize the margin  $\rho$ , Minimize  $\|\mathbf{w}\|$



# Non-linear SVM

Kernels functions: Mapping the original feature space to some higher-dimensional feature space where the training set is separable



You will try out the rbf kernel and a polynomial kernel in the hand's on exercise.

# Support Vector Machines

- Pros
  - Performs well in high dimensional spaces
    - # features > # training points
  - Can find non linear relationships by using kernel trick
- Cons
  - Struggles with noisy data
  - Not as interpretable as logistic regression
    - No probabilities
    - No slopes

# Naïve Bayes Classifier

# Naïve Bayes Classifier

- Up until this point we have mostly made predictions by computing distances
  - KNN – distance to nearest neighbors
  - Linear Regression – Residuals
  - SVM – maximum margin
- Naïve Bayes is a **probabilistic model**
  - Similar to Logistic Regression, Naïve Bayes estimates probabilities of outputs

# Naïve Bayes Classifier

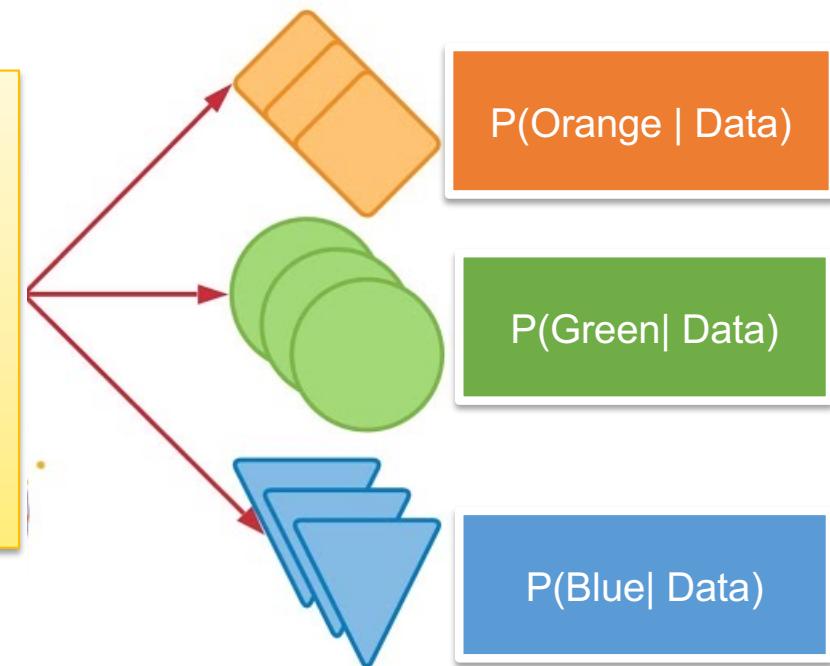
**Labelled Data**



Bayes Theorem

$$P(C|D) = \frac{P(D|C)P(C)}{P(D)}$$

*Estimate probabilities of class given training data*



# Bayes Theorem for Inference

Let's use C (Class) and D (Data) instead of A, B

$$P(C|D) = \frac{P(D|C)P(C)}{P(D)}$$

Prior probability of  $C$

Likelihood of  $D$  given  $C$

Posterior probability of  $C$  given  $D$

Evidence for  $D$

The diagram illustrates the components of Bayes' Theorem. At the top right is a light blue box labeled "Prior probability of C". A blue line points from this box down to the term "P(C)" in the numerator of the equation. Below the equation, a yellow box contains the term "Likelihood of D given C". A green line points from this box up to the term "P(D|C)" in the numerator. To the left of the equation, a yellow box contains the term "Posterior probability of C given D". A yellow line points from this box down to the term "P(C|D)" in the equation. At the bottom right is a grey box labeled "Evidence for D", with a black line pointing up to the term "P(D)" in the denominator.

Bayes Theorem: update probability of class based on the data

# Example Calculation

Refund	Marital Status	Evade
Yes	Single	No
No	Married	No
No	Single	No
Yes	Married	No
No	Divorced	Yes
No	Married	No
Yes	Divorced	No
No	Single	Yes
No	Married	No
No	Single	Yes

Compute  $P(C|D)$  where

Classes (C):

1. Yes – Evade Taxes
2. No – Do not Evade

# How to Estimate Probabilities from Data?

Refund	Marital Status	Evade
Yes	Single	No
No	Married	No
No	Single	No
Yes	Married	No
No	Divorced	Yes
No	Married	No
Yes	Divorced	No
No	Single	Yes
No	Married	No
No	Single	Yes

$$P(C|D) = \frac{P(D|C)P(C)}{P(D)}$$

We do not need to compute this.

# How to Estimate Probabilities from Data?

Refund	Marital Status	Evade
Yes	Single	No
No	Married	No
No	Single	No
Yes	Married	No
No	Divorced	Yes
No	Married	No
Yes	Divorced	No
No	Single	Yes
No	Married	No
No	Single	Yes

$$P(C|D) = \frac{P(D|C)P(C)}{P(D)}$$

- $P(C) = N_c/N$
- $N_c$  is total number of data points of class C
- N is total number of data points

e.g.,  $P(\text{No}) = 7/10$ ,  
 $P(\text{Yes}) = 3/10$

# How to Estimate Probabilities from Data?

Refund	Marital Status	Evade
Yes	Single	No
No	Married	No
No	Single	No
Yes	Married	No
No	Divorced	Yes
No	Married	No
Yes	Divorced	No
No	Single	Yes
No	Married	No
No	Single	Yes

$$P(C|D) = \frac{P(D|C)P(C)}{P(D)}$$

- **Assumption 1**
- Assumes Features are independent
- This allows us to treat each column independently and multiply computed probabilities

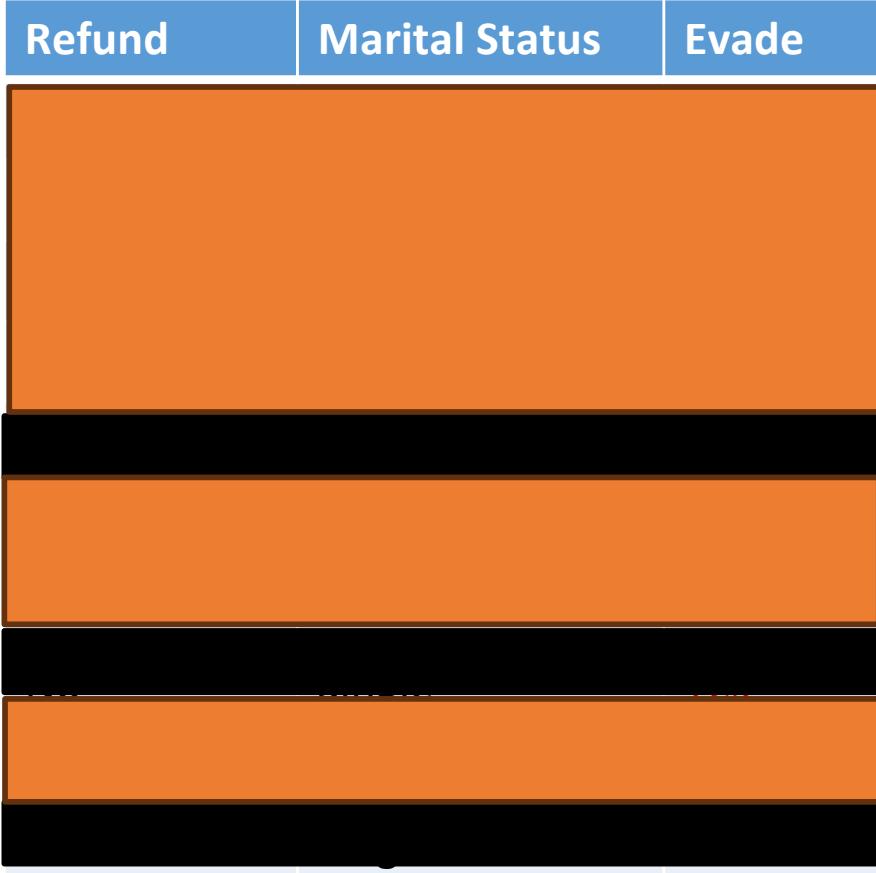
Given a Test Record:

$$X = (Refund = No, Married)$$

Compute:

$$P(D|No) = P(Refund = No|No)P(Married|No)$$

# How to Estimate Probabilities from Data?



$$P(C|D) = \frac{P(D|C)P(C)}{P(D)}$$

- Example of computing Likelihood
- $P(\text{Status}=\text{Married}|\text{Evade}=\text{No})$ 
  - only consider training data where class is no
  - compute the percent of data points where marital status = Married
- $P(\text{Status}=\text{Married}|\text{No}) = 4/7$
- $P(\text{Refund}=\text{Yes}|\text{Yes})=0$

# How to Estimate Probabilities from Data?

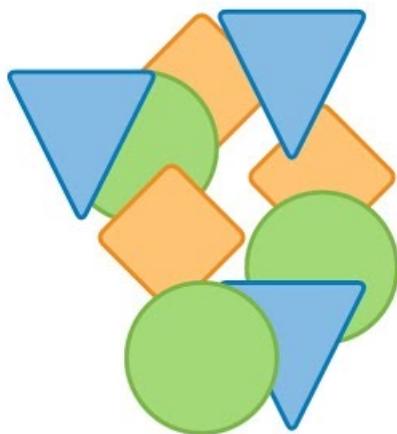
Refund	Marital Status	Evade
Yes	Single	No
No	Married	No
No	Single	No
Yes	Married	No
No	Divorced	Yes
No	Married	No
Yes	Divorced	No
No	Single	Yes
No	Married	No
No	Single	Yes

$$P(C|D) = \frac{P(D|C)P(C)}{P(D)}$$

- **Assumption 2:**
- Assume multinomial probability distribution for each feature
- There are different types of Naïve Bayes depending on the assumed underlying distribution per feature
- Multinomial, Gaussian, etc.

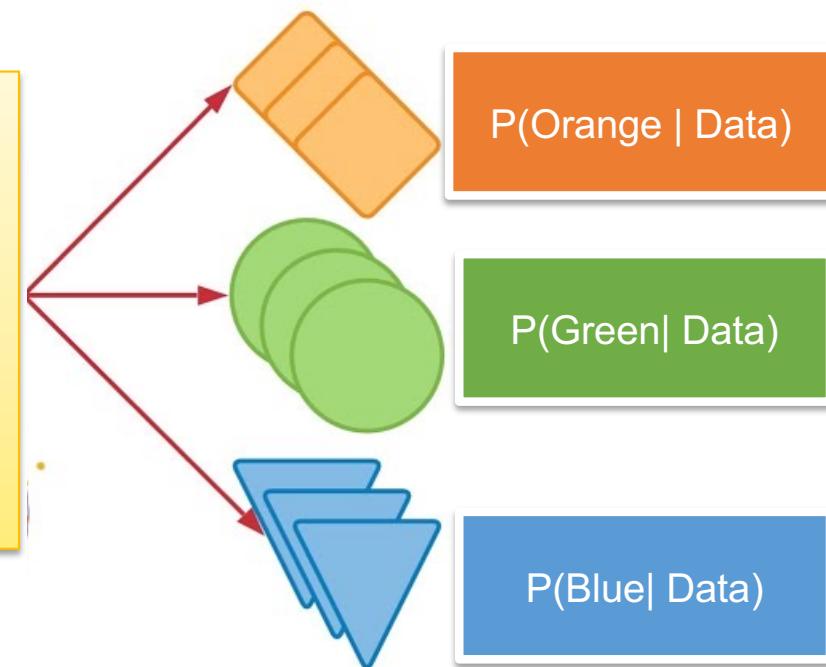
# Naïve Bayes Classifier

**Labelled Data**



Bayes Theorem

$$P(C|D) = \frac{P(D|C)P(C)}{P(D)}$$



***Highest posterior  
computed is the predicted  
class***

# Naïve Bayes (Assumptions)

Why is Naïve Bayes Naïve?

- Assumes features independent from one another
- Assumes variables come from a specific distribution

Despite assumptions, naïve bayes can still make good predictions in some use cases and is efficient at making predictions

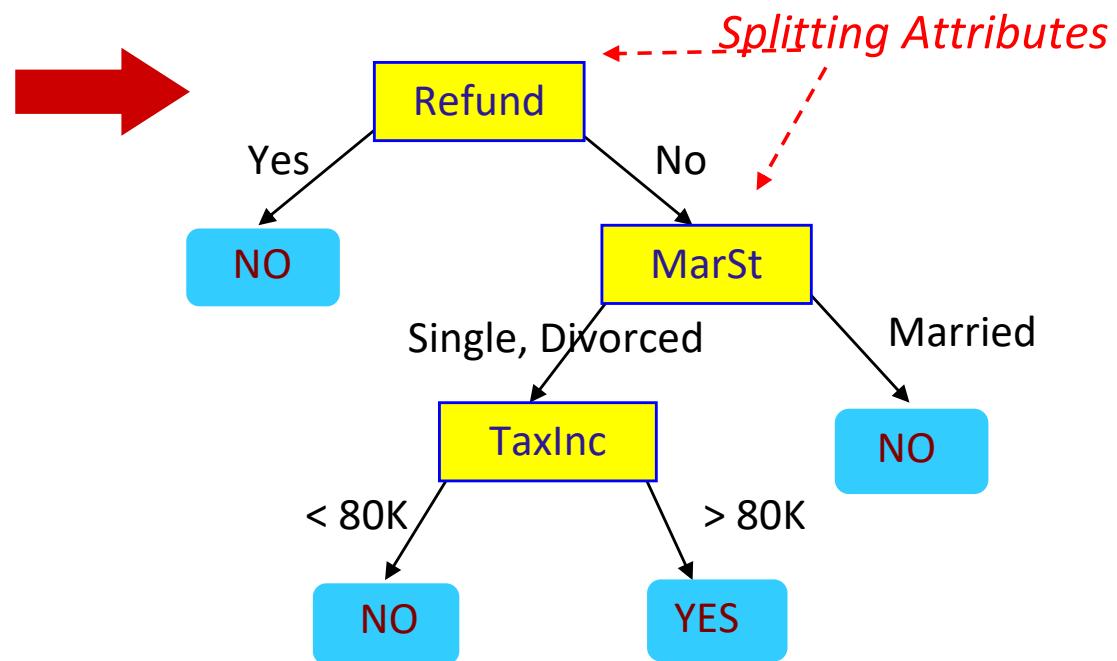
- Text classification

# Tree-based Classifier

# Decision Tree

Tid	Refund	Marital Status	Taxable Income	Cheat	class
1	Yes	Single	125K	No	categorical
2	No	Married	100K	No	categorical
3	No	Single	70K	No	continuous
4	Yes	Married	120K	No	continuous
5	No	Divorced	95K	Yes	continuous
6	No	Married	60K	No	continuous
7	Yes	Divorced	220K	No	continuous
8	No	Single	85K	Yes	continuous
9	No	Married	75K	No	continuous
10	No	Single	90K	Yes	continuous

- A decision tree is a (binary) tree with each branch having a left and right node.
- Each **branch node** specifies a feature to split along with information on which values to pass to the next node or leaf.
- **Leaf nodes** contain the predicted class label and are terminal.



Training Data

Model: Decision Tree

# Decision Trees

Tid	Refund	Marital Status	Taxable Income	Cheat
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
6	No	Married	60K	No
7	Yes	Divorced	220K	No
8	No	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes

When building a decision tree:



1. Considering all possible binary splits based on a single feature
2. Calculate the information gain for all possible splits
3. Pick the best split
4. Repeat with children nodes until pure node is created or some other stopping criteria

Class 1: Don't cheat (default class)

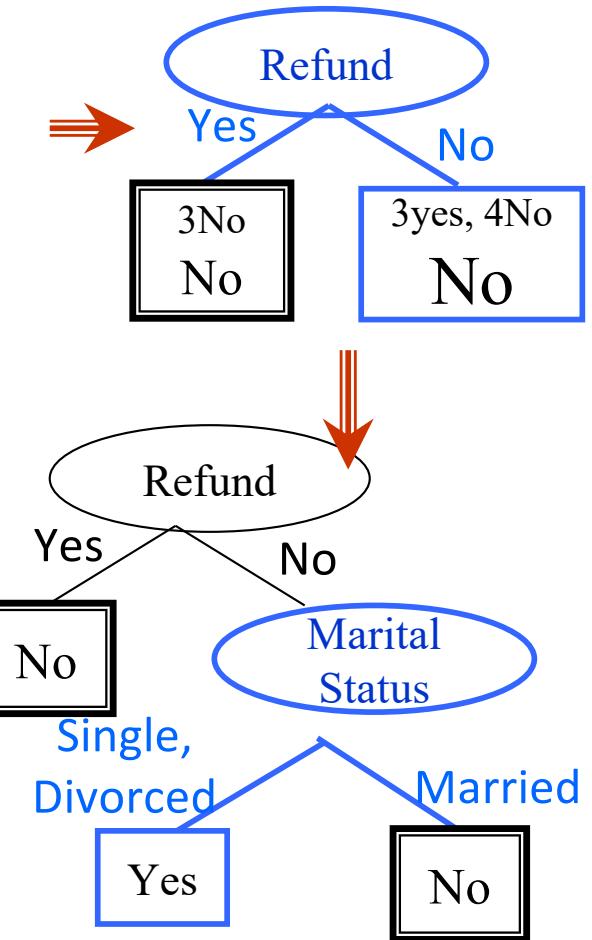
Class 2: Cheat

# Decision Trees

Tid	Refund	Marital Status	Taxable Income	Cheat
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
6	No	Married	60K	No
7	Yes	Divorced	220K	No
8	No	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes

When building a decision tree:

1. Considering all possible binary splits based on a single feature
2. Calculate the information gain for all possible splits
3. Pick the best split
4. Repeat with children nodes until pure node is created or some other stopping criteria



Class 1: Don't cheat (default class)

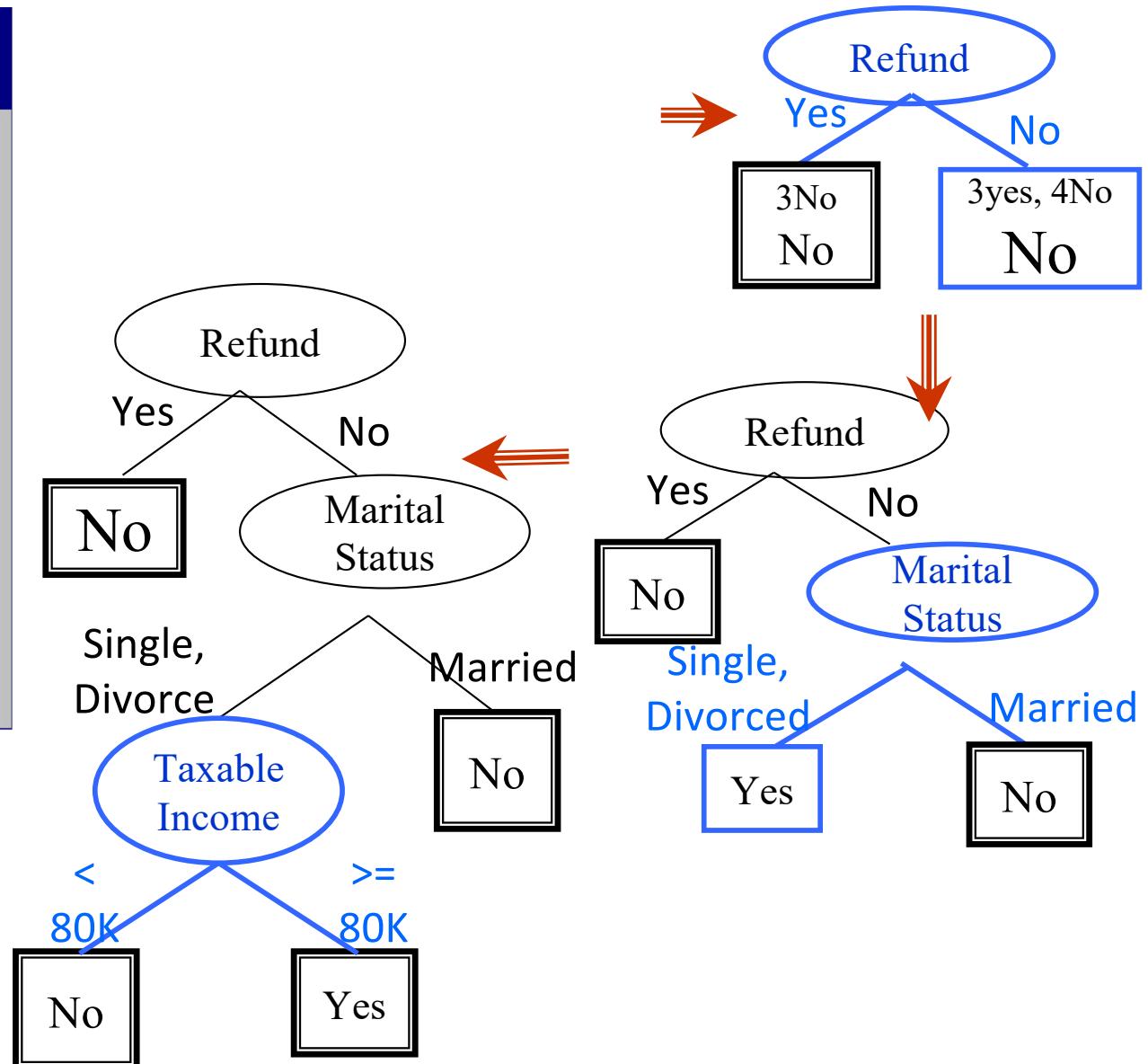
Class 2: Cheat

# Decision Trees

Tid	Refund	Marital Status	Taxable Income	Cheat
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
6	No	Married	60K	No
7	Yes	Divorced	220K	No
8	No	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes

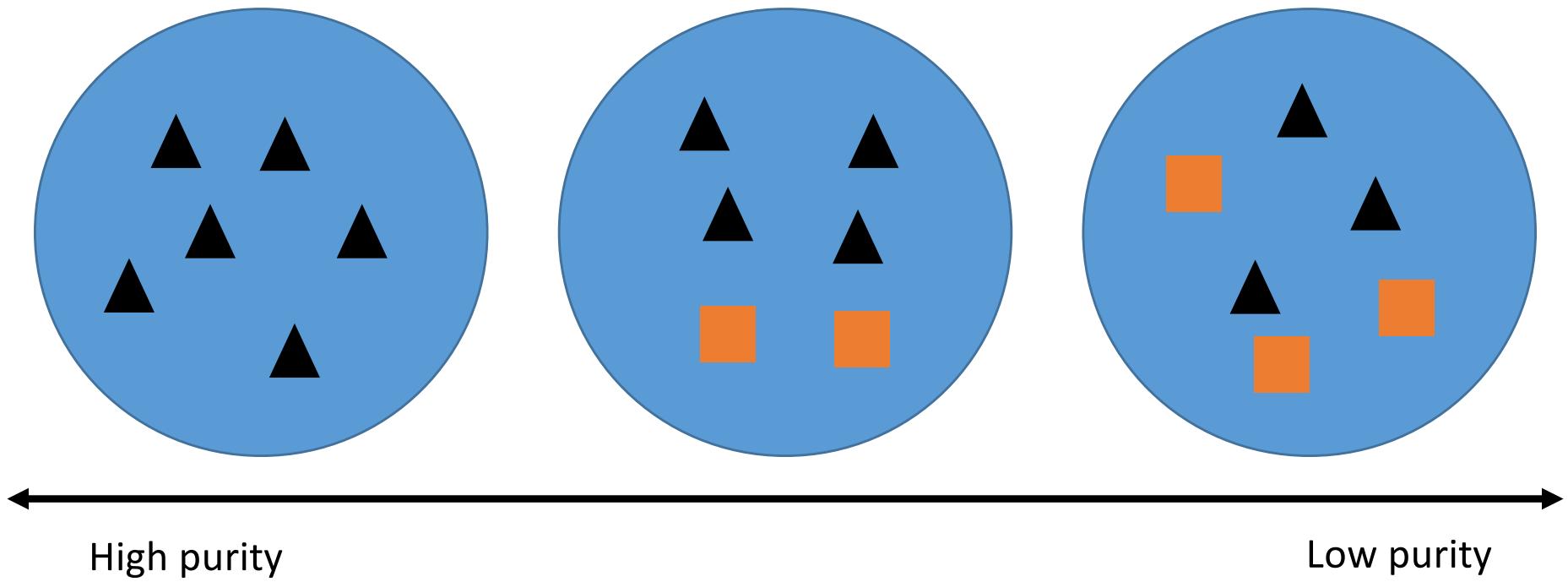
Class 1: Don't cheat (default class)

Class 2: Cheat



# Which split is the best?

-- measure impurity



$$P(\text{black}) = 1$$

$$P(\text{orange}) = 0$$

$$P(\text{black}) = 0.5$$

$$P(\text{orange}) = 0.5$$

# Measuring Node Impurity

- Let  $p(j | t)$  is the relative frequency of class j at node t).
  - Entropy

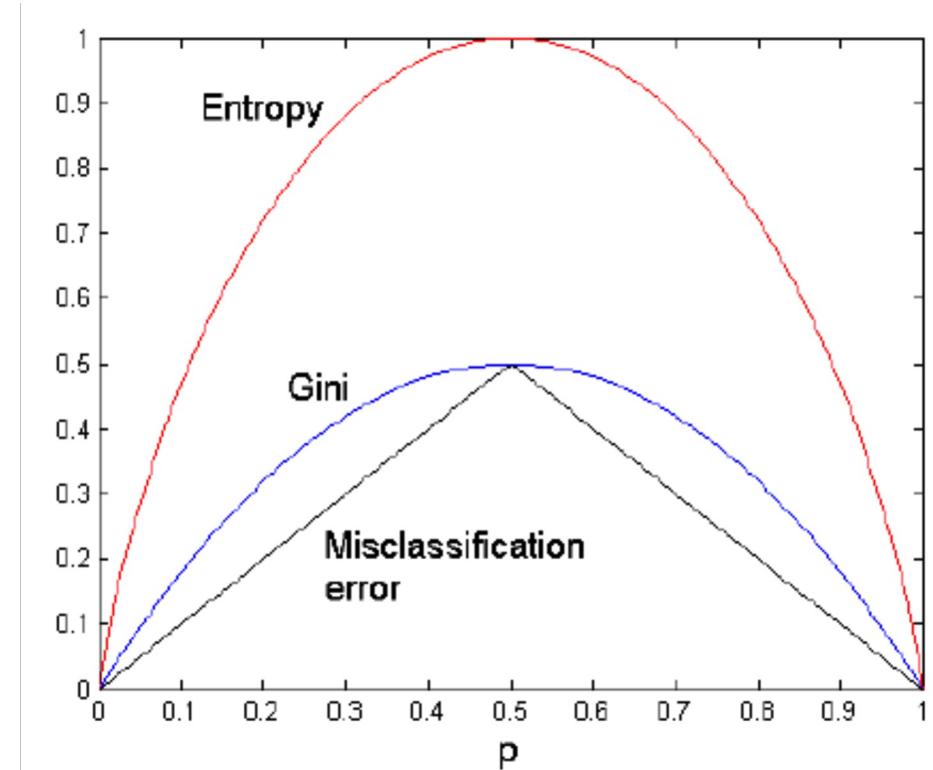
$$\text{Entropy}(t) = -\sum_j p(j | t) \log p(j | t)$$

- Gini Index

$$GINI(t) = 1 - \sum_j [p(j | t)]^2$$

- Misclassification Error

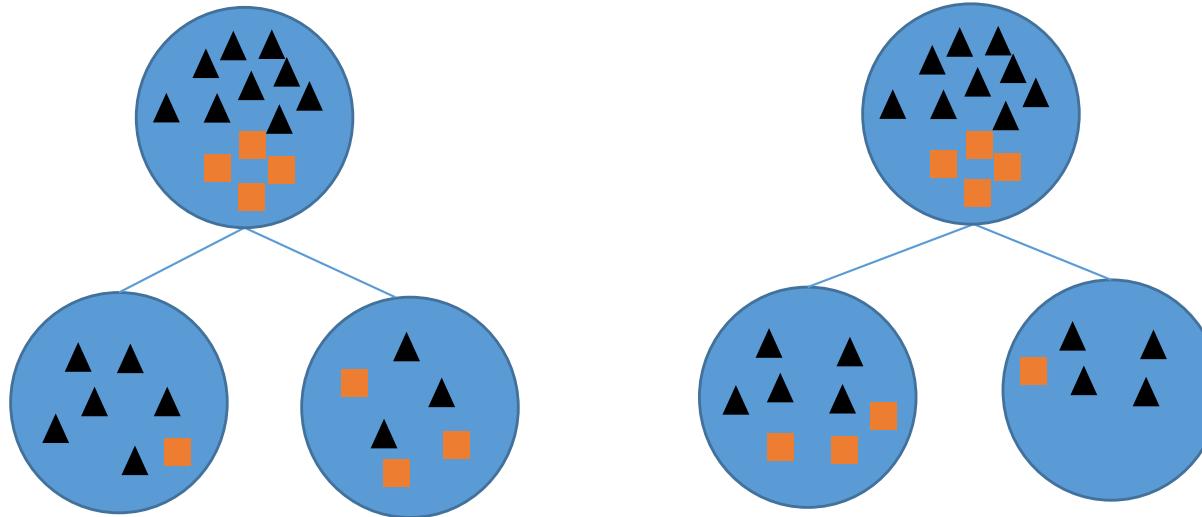
$$\text{Error}(t) = 1 - \max_i P(i | t)$$



# Which split is the best?

## -- information gain

Which split is better?



1. Compute purity of all nodes
2. Compare the purity of the parent node to the weighted average of the children

$$IG(P, C) = H(P) - \frac{|C_1|}{|P|}H(C_1) - \frac{|C_2|}{|P|}H(C_2)$$

Where P is the parent node

C<sub>i</sub> is one of the two children nodes

H is the impurity metric

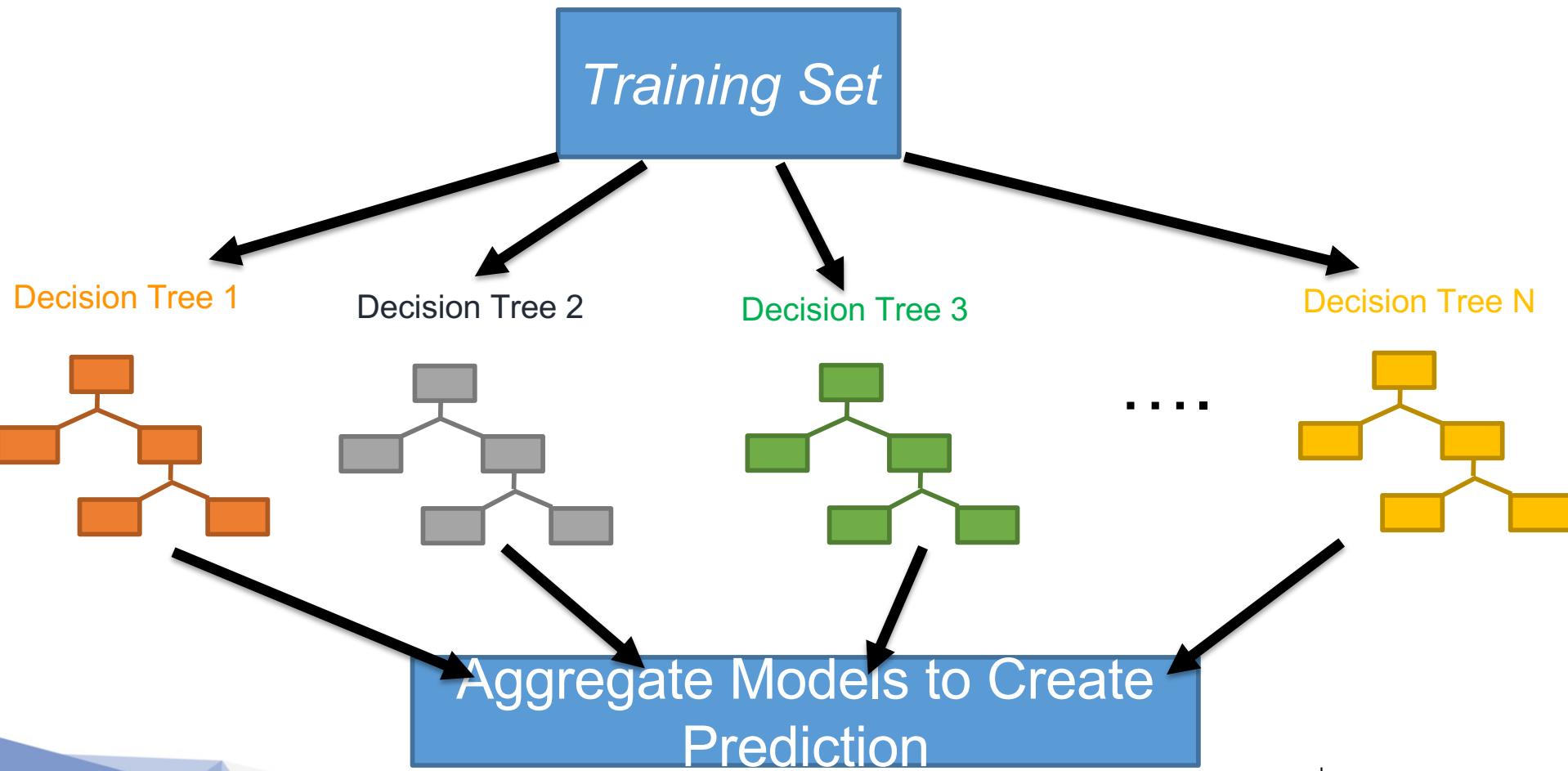
|| refers to the data points in that node

# Decision Tree Based Classification

- Advantages:
  - Inexpensive to construct
  - Extremely fast at classifying unknown records
  - Easy to interpret for small-sized trees
  - Accuracy is comparable to other classification techniques for many simple data sets
- Less effective when
  - High background noise.
  - Large scale data
  - Data with high dimension

# Ensemble Methods

procedure that combines the outputs of many “weak” models to produce a “strong” model



# Random Forest

A collection of decorrelated decision trees

## Building stage

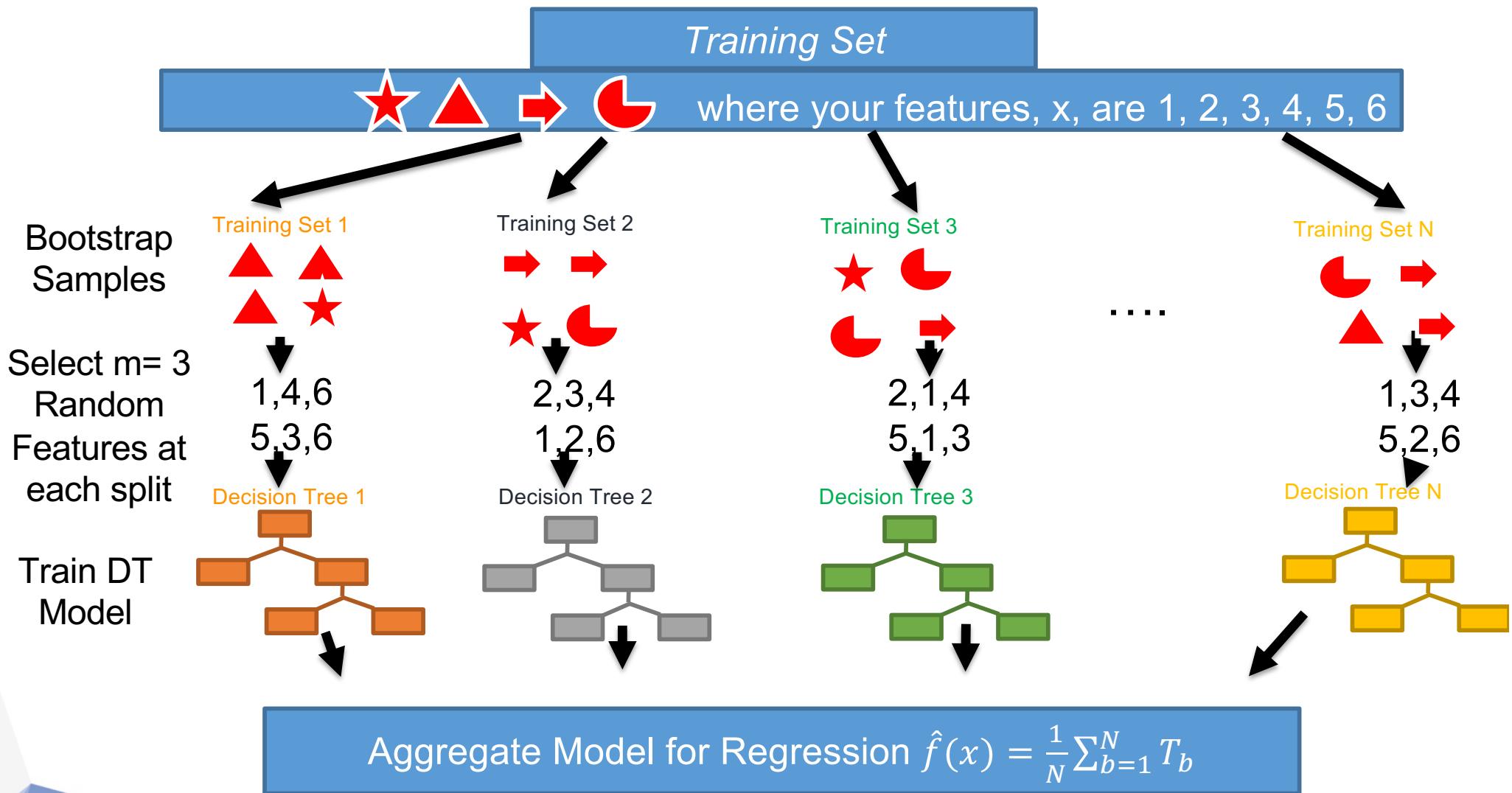
- Randomly sample data **with replacement**
- Only a subset of attributes are chosen to make decision at each node of the tree
- Select best split from selected subset of attributes.

## Classifying

All trees are used to make a decision.

The final class label is determined statistically e.g. majority rule.

# Random Forest



# Feature Importance

- When building random forest models we can get a sense for what features were most important to making predictions
- For example
  - Gini importance
    - Total decrease in node impurity for a particular feature weighted by the probability of reaching a particular node
  - Permutation importance
    - Quantify how much the performance of your model decrease if we randomly permute a feature.
- While the ability to quantify feature importance is great, it is not as interpretable as for example slopes in linear regression.

# Random Forest

## Advantage

Accurate

Good for high dimensional data

Reduce the need of feature detection and selection.

Give estimate on important features.

## Disadvantage

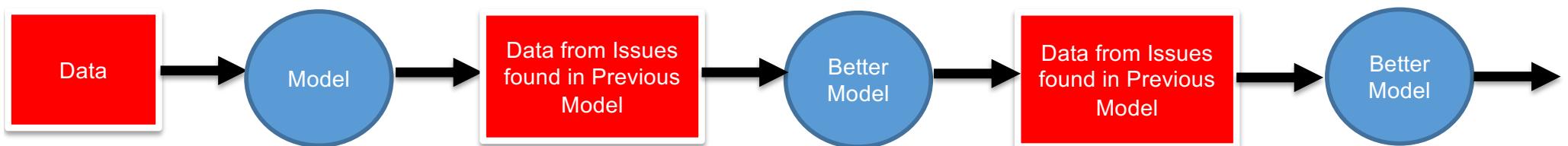
More computational requirement

Less effective with noisy training data.

The result is an statistical ensemble and might be hard to interpret.

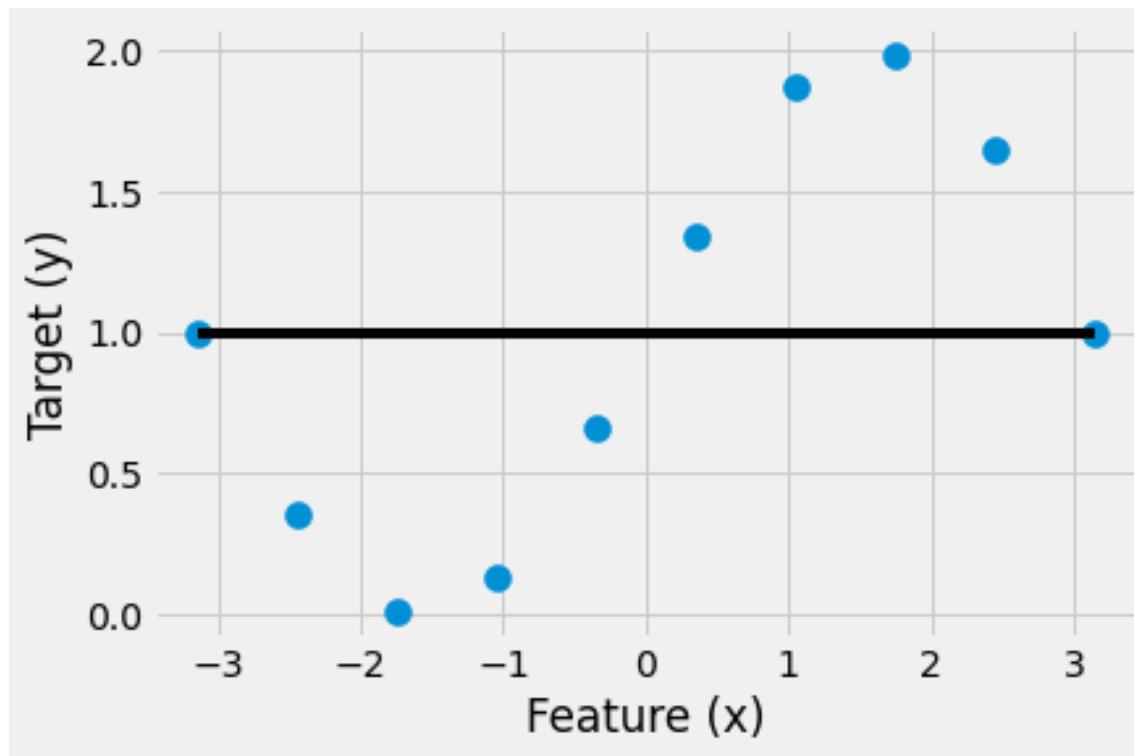
# Gradient Boosting

**Different Approach:** sequentially add weak models to repeatedly modified versions of the data to improve model



# Gradient Boosting

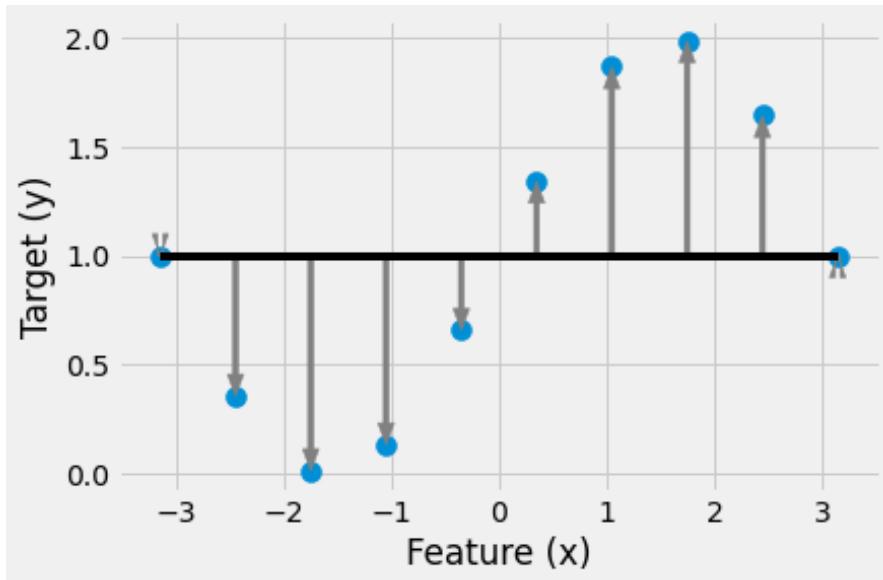
Start by creating a weak model; e.g. mean regressor



$$f_0(x) = \bar{y}$$

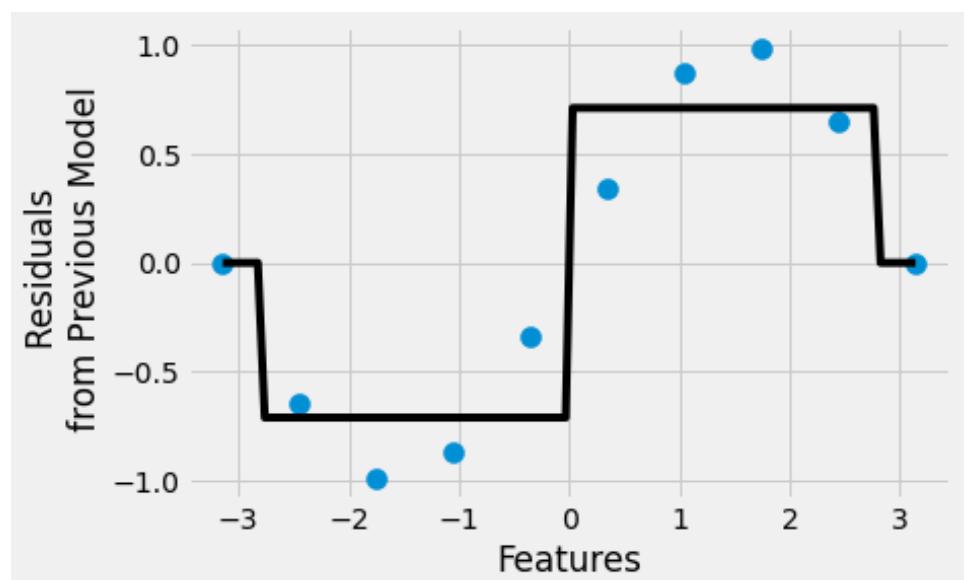
# Gradient Boosting

Quantify errors by computing residuals



$$y_{true} - f_0(x)$$

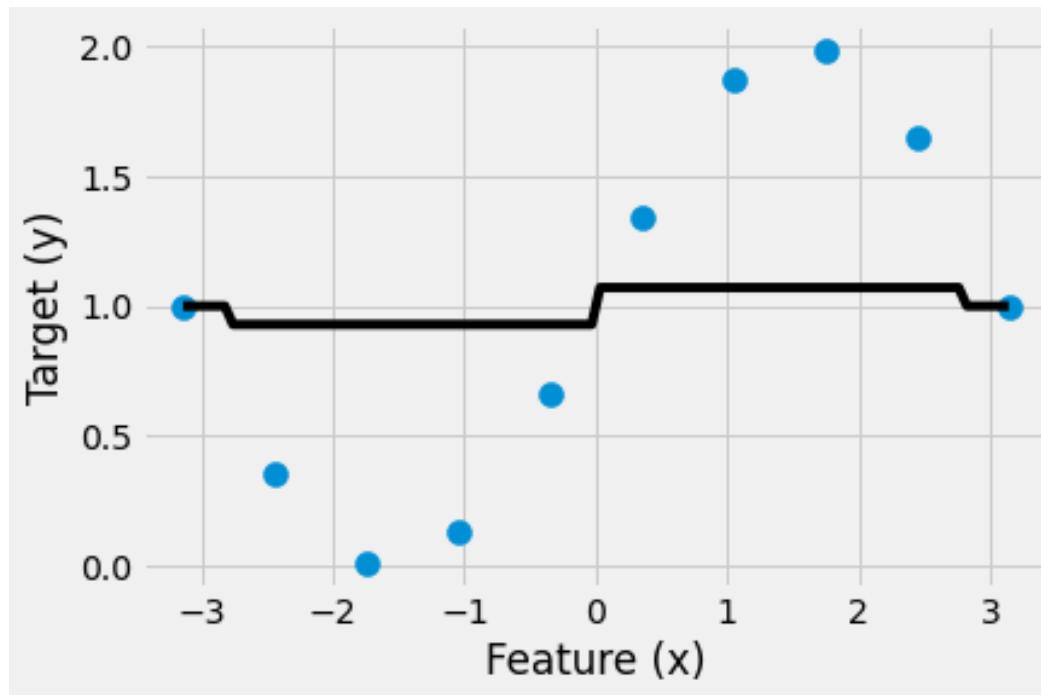
Build a decision tree to learn the residual



$$h_1(x)$$

# Gradient Boosting

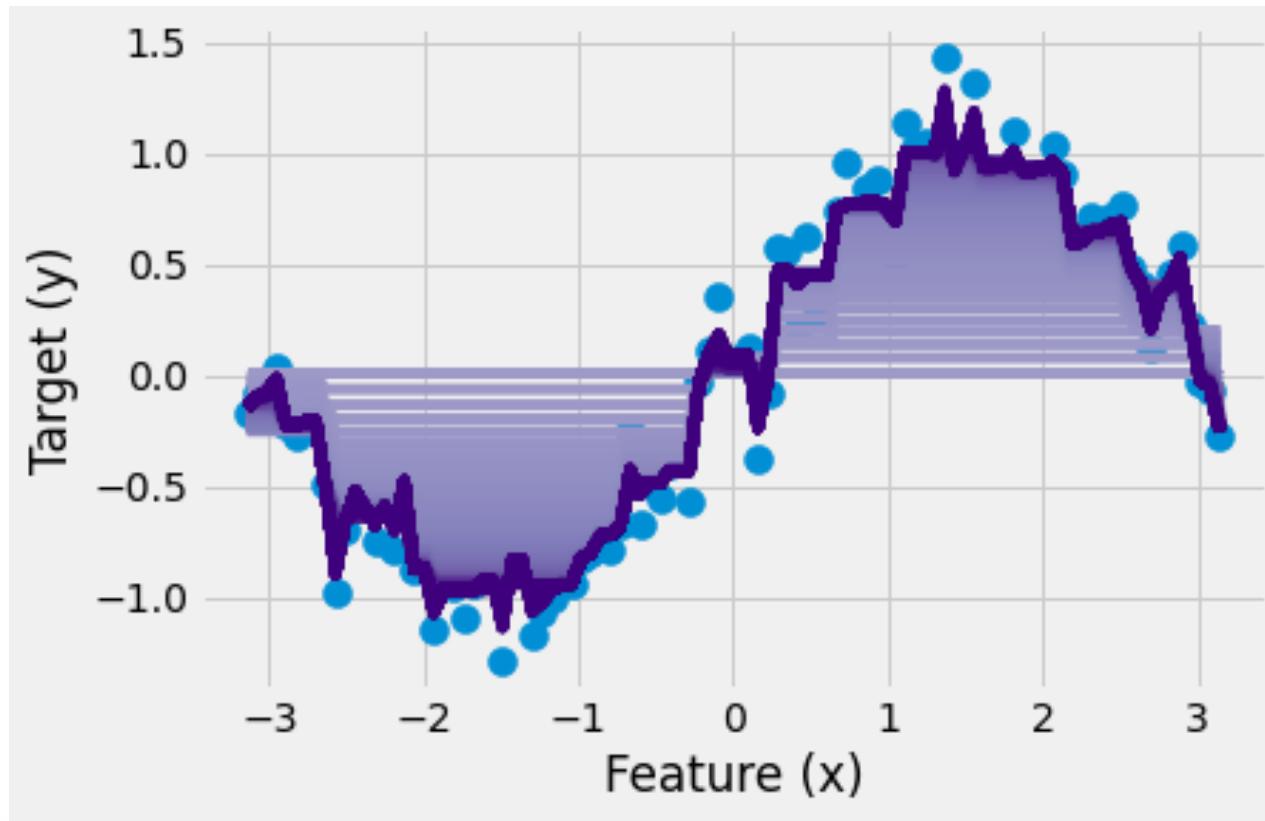
Add decision tree approximating residuals to mean regressor



$$f_1(x) = f_0(x) + \lambda h_1(x)$$

# Gradient Boosting

Repeat



# Gradient Boosting vs Random Forest

- Both are powerful models
- Random Forest tends to get better out of the box performance and requires less hyperparameter tuning
- If tuned correctly, Gradient Boosting tends to match or out perform random forest.

# Summary

- Supervised learning methods model the connections between training data with their known outcome value/label.
- By distance
  - E.g. knn, regression, SVM
  - Predictions are based on position in the space.
  - Data can be projected into new spaces
- By probability distribution
  - E.g. Naïve Bayes Classifier
  - A statistical distribution is “learned”.
  - Predictions are based on likelihood of observation.
- By feature space
  - E.g. decision trees, Random forest.
  - Data are divided according to their feature spaces.
  - Predictions are based on group of features that shared.

# Summary

- Machine Learning (ML) Basics
  - Define ML, Categories of ML, ML Workflow
- Common Supervised Learning Techniques
  - K Nearest Neighbors, Linear Regression , Logistic Regression, Support vector machines , Probabilistic Models , Decision Tree Based Models
- Evaluation Metrics
  - Classification: Confusion Matrix, Precision, Recall, ROC curve
  - Regression: MSE, MAE, R<sup>2</sup>
- Preparing Data for Training
  - Standardizing Data
  - Feature Engineering / Selection

# Machine Learning on HPC: Supervised Learning Exercises and Case Study

# Jupyter notebook

- TACC Analysis Portal:

Go to <https://vis.tacc.utexas.edu>; Login with your training account credentials

**TACC** | Analysis Portal [User Guide](#)

**jrduncan** [Log Out](#)

### Submit New Job

System: ---

Application: Select System

Project: Select System

Queue: Select System

Nodes: 1 Tasks: 1

---

### Options

Job Name: 20 characters max

Time Limit: H:M:S (default 2:0:0)

Reservation: reservation name

VNC Desktop Resolution: WIDTHxHEIGHT

[Submit](#) [Utilities](#)

### System Status

System	Status	Utilization	Job Count
Frontera	Open	99%	Running: 312 Queued: 1251
Lonestar6	Open	69%	Running: 135 Queued: 98
Longhorn	Open	74%	Running: 30 Queued: 40
Maverick2	Open	16%	Running: 4 Queued: 7
Stampede2	Open	96%	Running: 830 Queued: 736

### Past Jobs

JNB-Frontera	03/18/2022	<a href="#">Details</a>	<a href="#">Resubmit</a>
JNB-Frontera	03/18/2022	<a href="#">Details</a>	<a href="#">Resubmit</a>
JNB-Frontera	03/18/2022	<a href="#">Details</a>	<a href="#">Resubmit</a>
JNB-Frontera	03/18/2022	<a href="#">Details</a>	<a href="#">Resubmit</a>
JNB-Frontera	02/21/2022	<a href="#">Details</a>	<a href="#">Resubmit</a>

TACC

101

# Jupyter notebook

- TACC Visualization Portal:  
Select the following options

**Submit New Job**

System	Frontera
Application	Jupyter notebook
Project	Frontera training
Queue	small

Nodes: 1 Tasks: 1

---

**Options**

Job Name	20 characters max
Time Limit	H:M:S (default 2:0:0)
Reservation	ML_Institute_Tue
VNC Desktop Resolution	WIDTHxHEIGHT

**TACC**

Submit Utilities

# Jupyter notebook



Analysis Portal User Guide

jrduncan

Log Out

## TAP Job Status

**Job:** Jupyter notebook on Frontera (4175197, 2022-03-21T17:28-05:00)

**Status:** RUNNING

**Start:** March 21, 2022, 5:28 p.m.

**End:** March 21, 2022, 5:33 p.m.

**Refresh:** in 873 seconds

### Message:

```
TAP: Your session is running at https://frontera.tacc.utexas.edu:60752  
/?token=9cbad0f26752e7dd14fcf090d6a30b6ec5c15c63ed7d9e2b626f214712fb8b4d
```

Connect

End Job

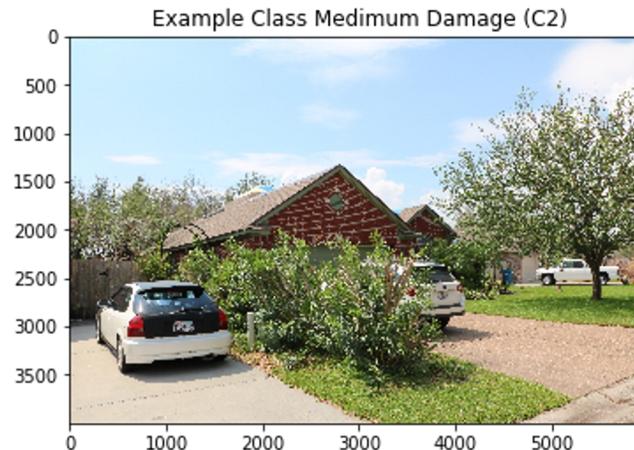
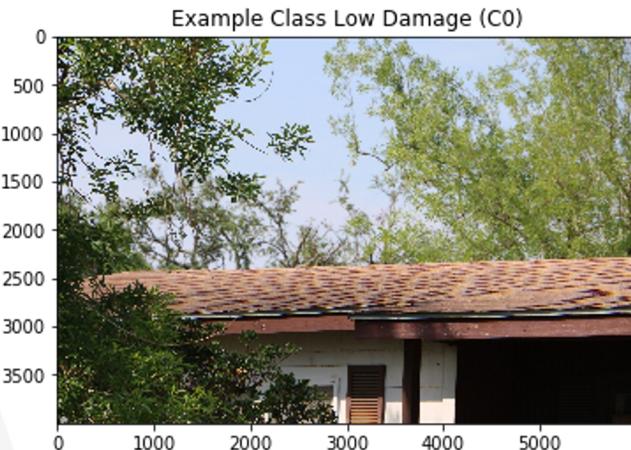
Show Output

Back to Jobs

# Exercise



## Image Classification with Hurricane Harvey Dataset

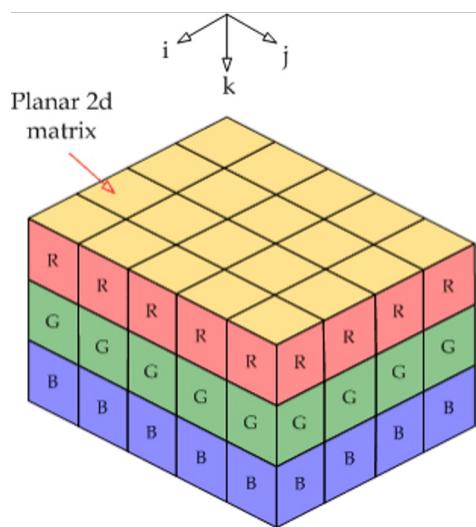


# Image Processing

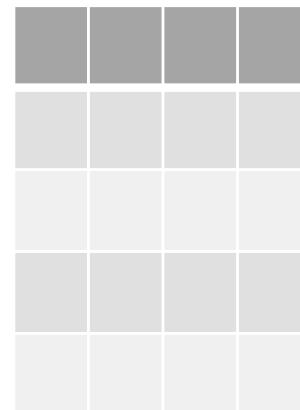
- Process of manipulating images
- Common image processing tasks:
  - Gray-scaling images
  - Resizing images
  - Reshaping data
- Image processing is used to for feature selection in supervised learning

# Image Processing

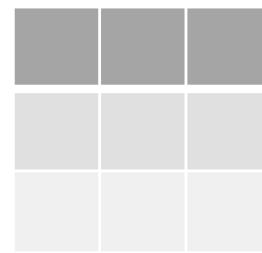
Color Image



Grayscale



Resize



Vectorize



# Today: Supervised Learning

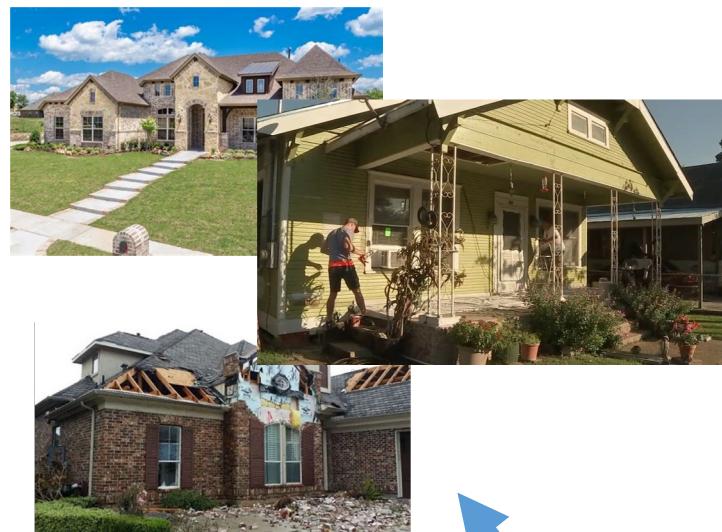
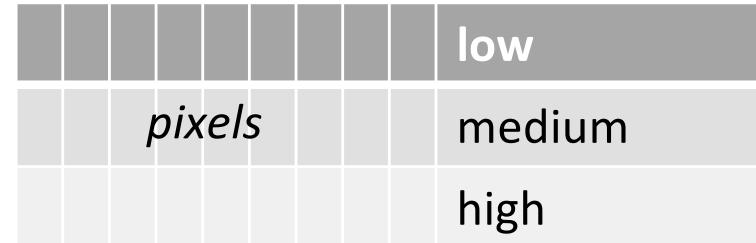


Image  
processing

Vectorized  
Data



Train Model



Evaluate  
Model

# Building and Evaluating Classification Models

- Take vectorized images from pipeline and use them to build supervised learning models
  - Two different datasets:
    - City streets versus Forest
    - Design Safe dataset
  - Decision tree, Logistic Regression, Naïve Bayes, Support Vector Machine, Random Forest
- Model Evaluation
  - Ok, so I built some models, now what?
  - Showcases taking a deeper dive into the models you built so that you can learn how to improve performance

# Instructions

In this exercise you will run two notebooks which build and evaluate image classification models

- Run the **city\_vs\_forest\_classifier.ipynb**

This exercise will build models to distinguish between images of city landscapes vs forest landscapes

- Run the **designsafe\_classifier.ipynb**

This exercise will build models using the DesignSafe Dataset you were introduced to yesterday. We will build classification models which predict whether the damage level is low, medium or high.