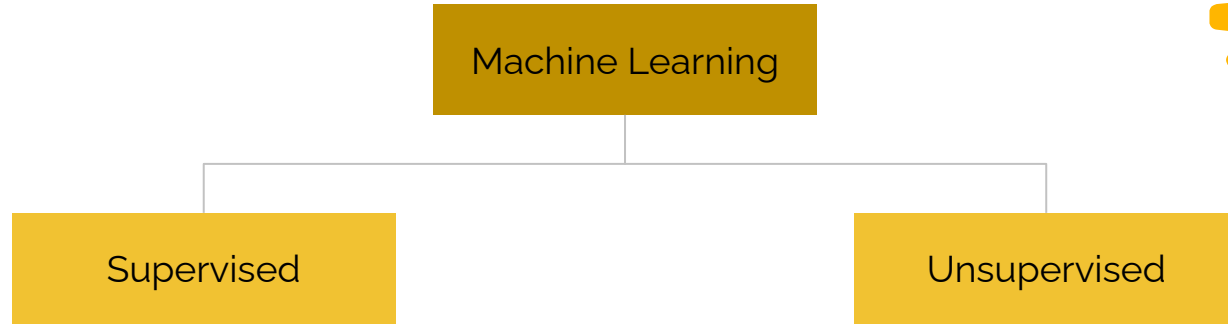
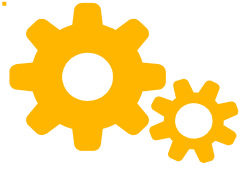




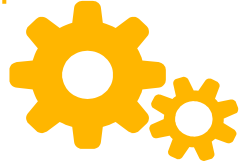
Machine Learning I

Random Forests

August 8th 2020

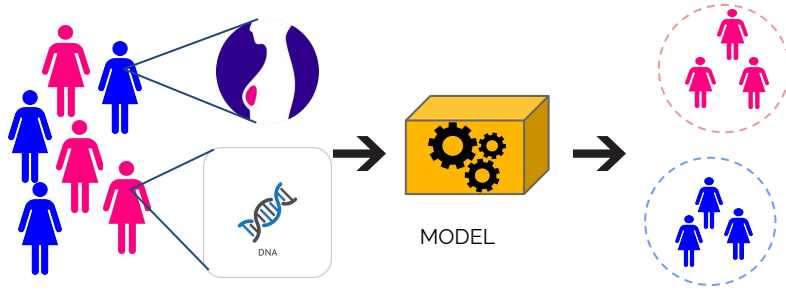


Machine Learning



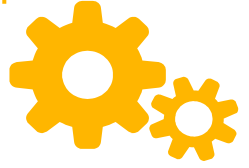
Supervised

Unsupervised



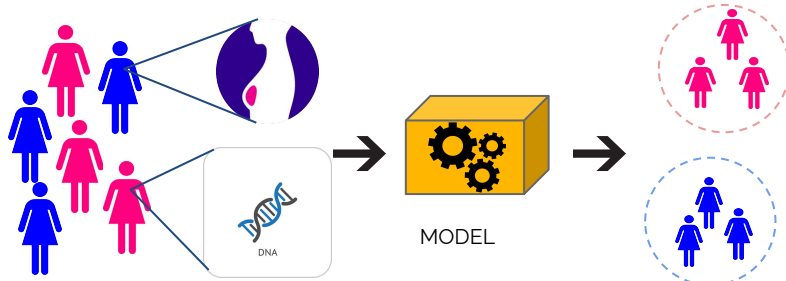
Labelled Data

Machine Learning

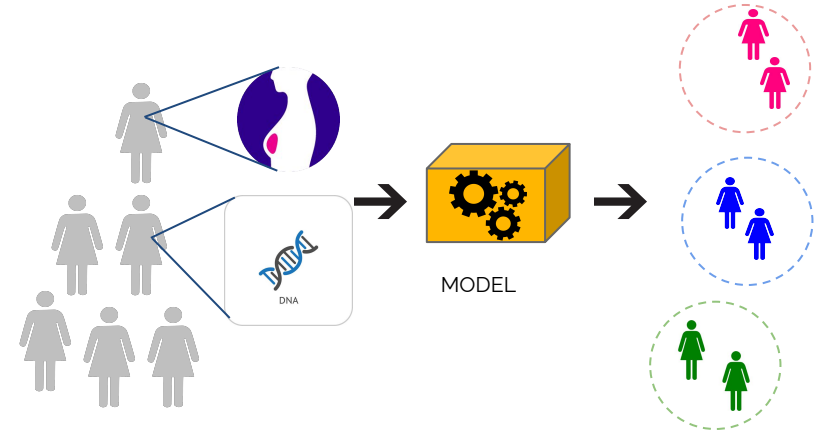


Supervised

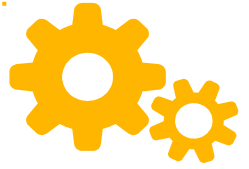
Unsupervised



Labelled Data



Unlabelled Data



Machine Learning

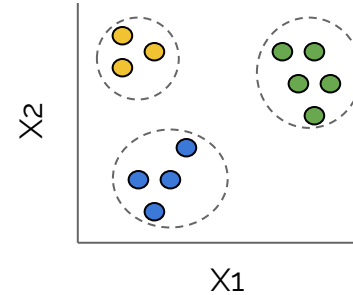
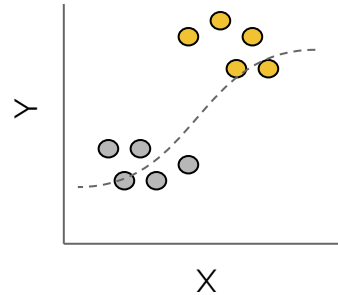
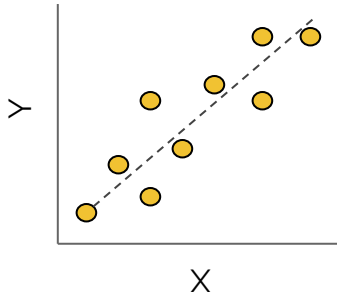
Supervised

Unsupervised

Regression

Classification

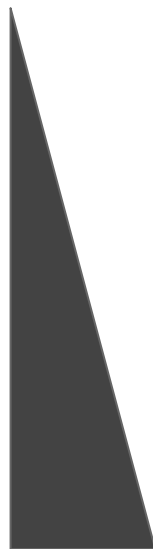
Clustering





Type	Name
Linear	Linear Regression
	Logistic Regression
Survival	Cox Proportional Hazard
Tree-based	Random Forest
	Gradient Boosting
Neural Network	Neural Networks

Performance



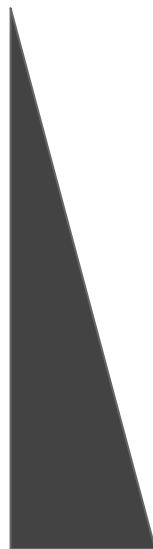
Interpretability



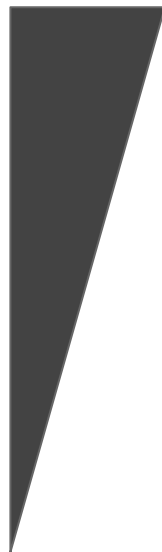


Type	Name
Linear	Linear Regression
	Logistic Regression
Survival	Cox Proportional Hazard
Tree-based	Random Forest
	Gradient Boosting
Neural Network	Neural Networks

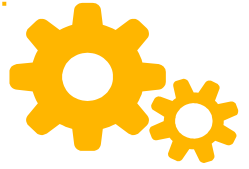
Performance



Interpretability



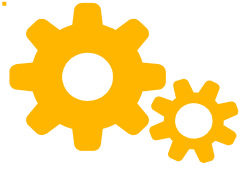
Let's Build an Orange Classifier!



ORANGE



GRAPEFRUIT



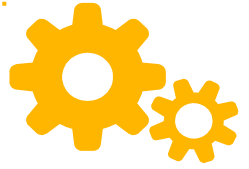
ORANGE

=



GRAPEFRUIT

citrus?

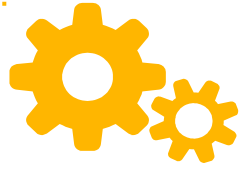


ORANGE



GRAPEFRUIT

sweet?

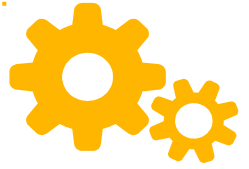


ORANGE



GRAPEFRUIT

weight?



ORANGE



ORANGE



GRAPEFRUIT



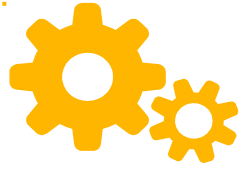
ORANGE



GRAPEFRUIT



GRAPEFRUIT



Collect Data!

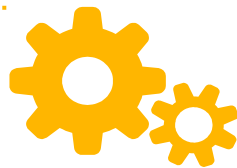


ORANGE



GRAPEFRUIT

fruit	citrus
1	yes
2	yes
3	yes
4	yes
5	yes
6	yes



Collect Data!

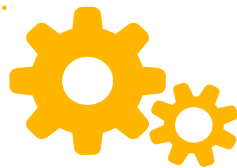


ORANGE



GRAPEFRUIT

fruit	citrus	sugar
1	yes	10
2	yes	11
3	yes	7
4	yes	10
5	yes	6
6	yes	5



Collect Data!

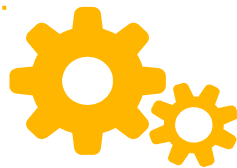


ORANGE



GRAPEFRUIT

fruit	citrus	sugar	weight
1	yes	10	130
2	yes	11	115
3	yes	7	120
4	yes	10	200
5	yes	6	190
6	yes	5	123



Collect Data!

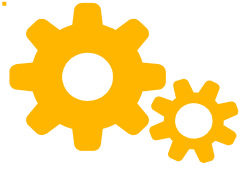


ORANGE



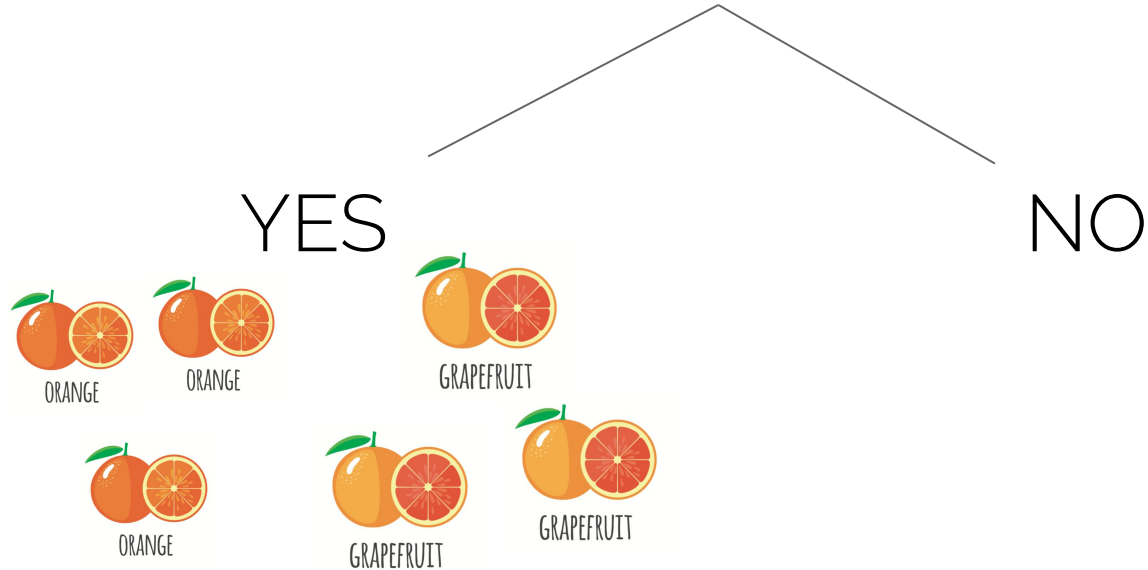
GRAPEFRUIT

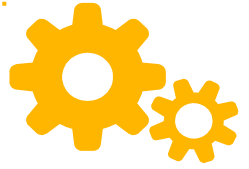
fruit	citrus	sugar	weight	orange
1	yes	10	130	1
2	yes	11	115	1
3	yes	7	120	1
4	yes	10	200	0
5	yes	6	190	0
6	yes	5	123	0



A Bad Question

Is it a citrus fruit?

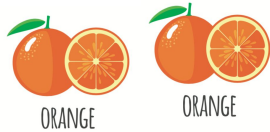




A Better Question

Is sugar $\geq 10g$

YES



ORANGE

ORANGE



GRAPEFRUIT

NO



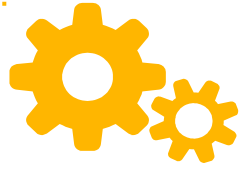
ORANGE



GRAPEFRUIT



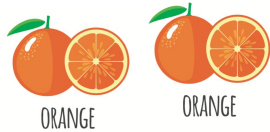
GRAPEFRUIT



A Much Better Question

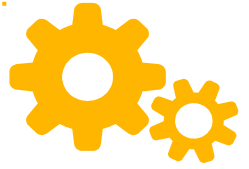
Is weight $< 150\text{g}$

YES



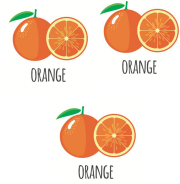
NO





Stack the Questions

Is weight $< 150\text{g}$

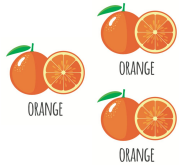


YES

NO



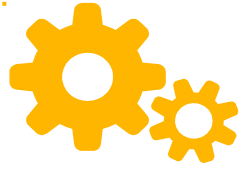
Is sugar $\geq 7\text{g}$



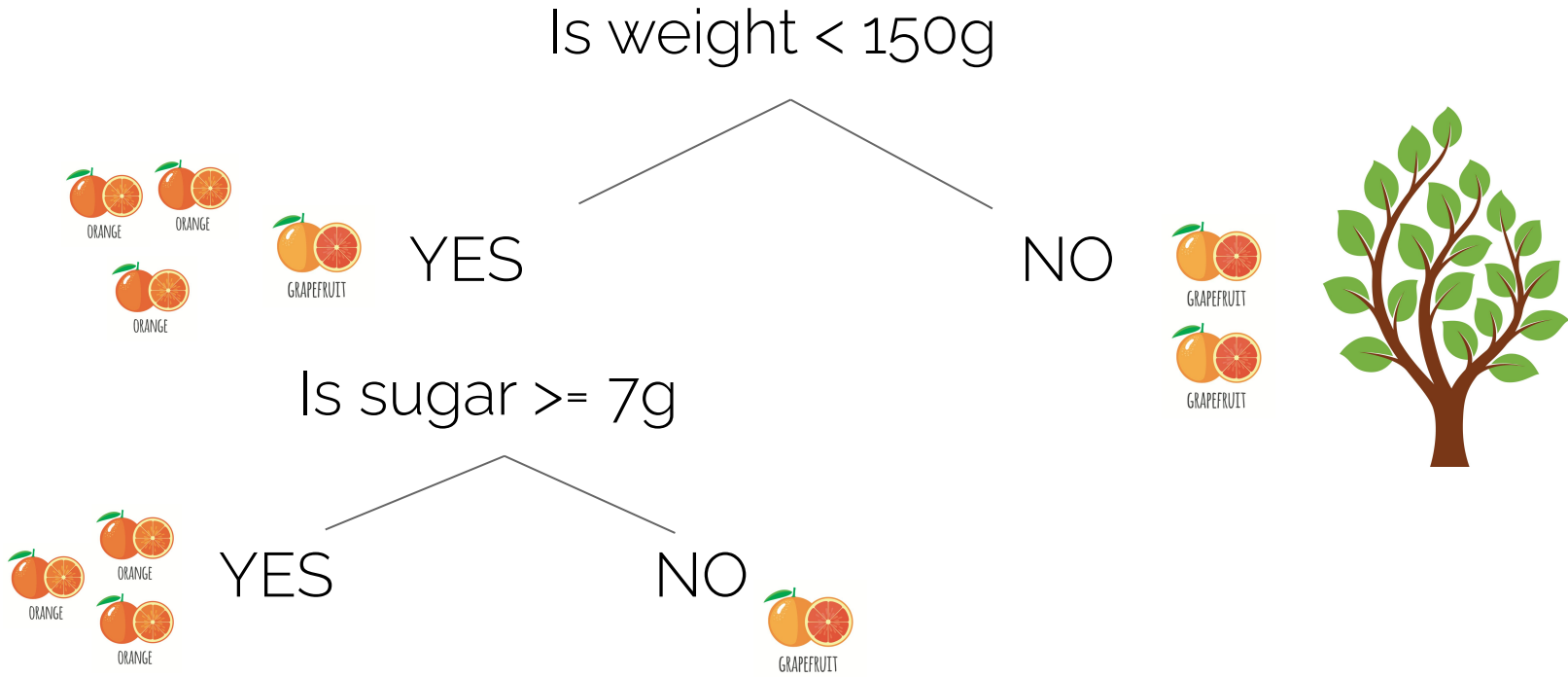
YES

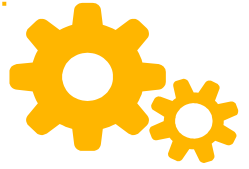
NO



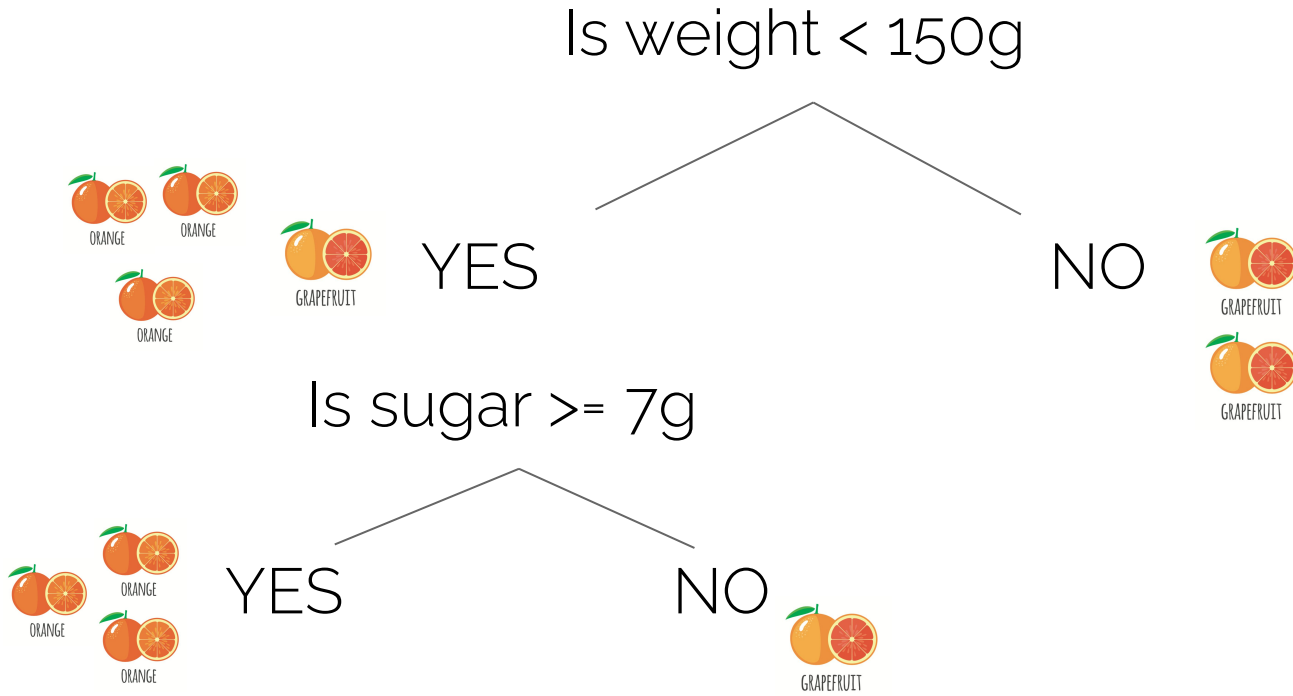


Decision Tree



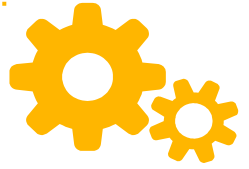


Decision Tree



Which questions should you ask first?

Randomness



ORANGE



ORANGE



GRAPEFRUIT



ORANGE



GRAPEFRUIT



GRAPEFRUIT



ORANGE



ORANGE



ORANGE



ORANGE

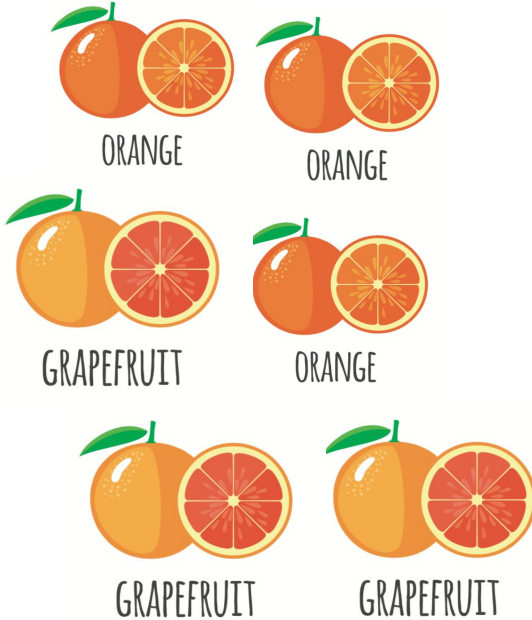
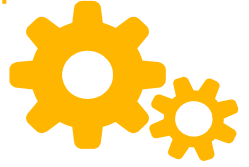


ORANGE

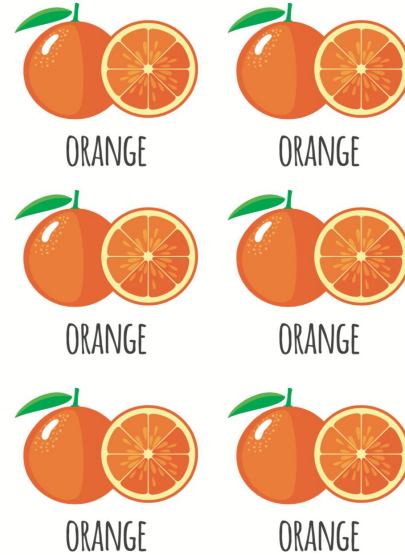


ORANGE

Entropy

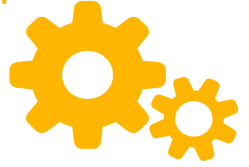


Random! ($e=1$)



Pure! ($e=0$)

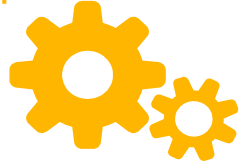
Entropy



$$E(S) = \sum_{i=1}^c -p_i \log_2 p_i$$

A measure of randomness!

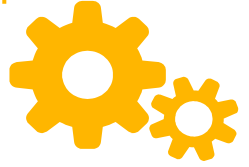
Entropy



$$E(S) = \sum_{i=1}^c -p_i \log_2 p_i$$

Go through each class and add

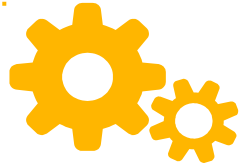
Entropy



$$E(S) = \sum_{i=1}^c -p_i \log_2 p_i$$

Proportion of fruits in each class * log₂
proportion of fruits in each class

Entropy

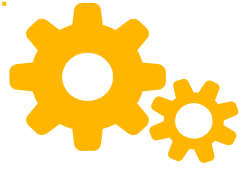


$$E(S) = \sum_{i=1}^c -p_i \log_2 p_i$$

`(-prop(oranges)*log2(prop(oranges))`

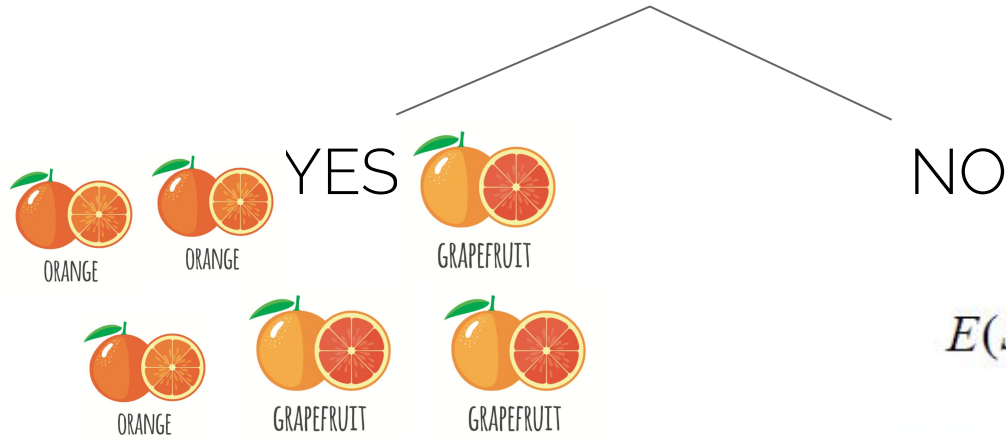
`+`

`(-prop(grapefruits)*log2(prop(grapefruits))`



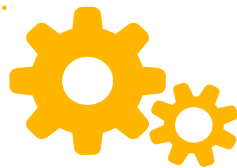
Decision Trees

Is it a citrus fruit?



$$E(S) = \sum_{i=1}^c -p_i \log_2 p_i$$

$$(3/6) \times -\log_2(3/6) + (3/6) \times -\log_2(3/6) = 1$$



Decision Trees

Is sugar $\geq 10g$

$$E(S) = \sum_{i=1}^c -p_i \log_2 p_i$$



YES

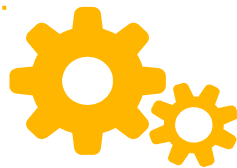


$$\begin{aligned} & (\frac{2}{3})x - \log_2(\frac{2}{3}) + (\frac{1}{3})x - \log_2(\frac{1}{3}) \\ & = 0.92 \end{aligned}$$

NO

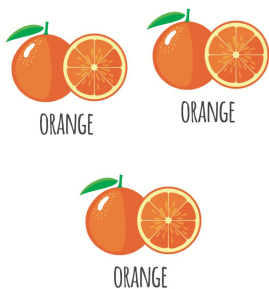


$$\begin{aligned} & (\frac{2}{3})x - \log_2(\frac{2}{3}) + (\frac{1}{3})x - \log_2(\frac{1}{3}) \\ & = 0.92 \end{aligned}$$



Decision Trees

Is weight < 150g



YES

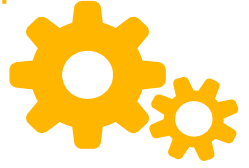


$$\left(\frac{3}{4}\right)x - \log_2\left(\frac{3}{4}\right) + \left(\frac{1}{4}\right)x - \log_2\left(\frac{1}{4}\right) \\ = 0.81$$

NO

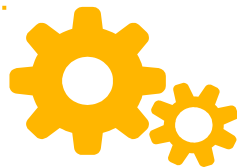


$$\left(\frac{2}{2}\right)x - \log_2\left(\frac{2}{2}\right) \\ = 0$$

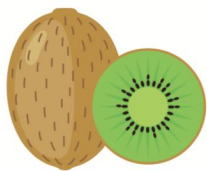


Weight Sugar Citrus

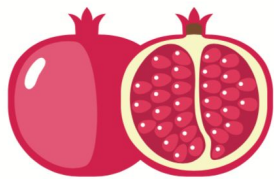




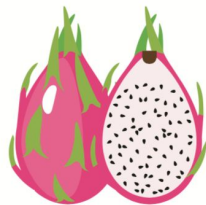
Larger dataset?



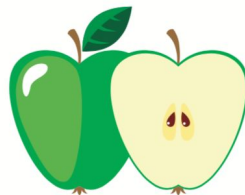
KIWI



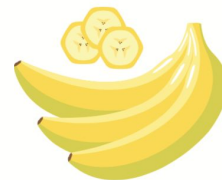
POMEGRANATE



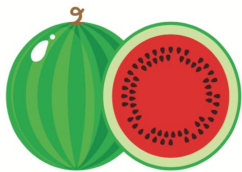
DRAGON FRUIT



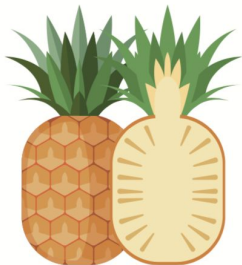
APPLE



BANANA



WATERMELON



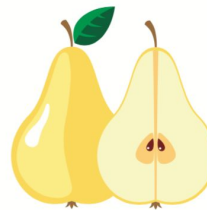
PINEAPPLE



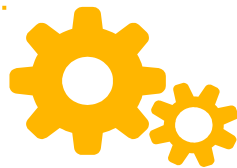
FIGS



ORANGE



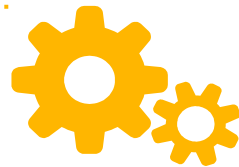
PEAR



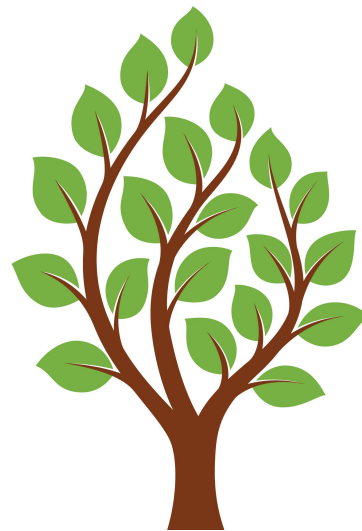
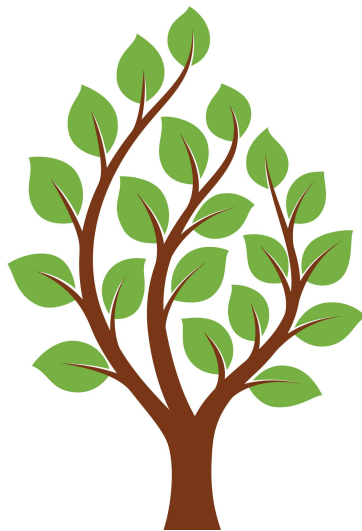
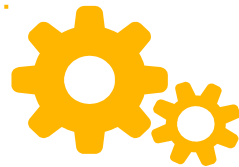
Many more questions!

fruit	citrus	sugar	weight	Col in	Col out	pit	seeds	round	soft	orange
1	yes	10	130	yellow	yellow	yes	no	yes	yes	0
2	yes	11	115	orange	orange	no	yes	no	no	1
3	yes	7	120	red	green	no	yes	no	yes	0
4	yes	10	200	yellow	yellow	no	no	no	no	0
5	yes	6	190	white	yellow	no	yes	yes	no	0
6	yes	5	123	green	green	no	no	no	no	0

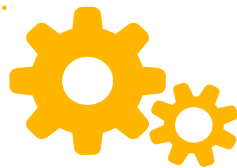
Many more questions!



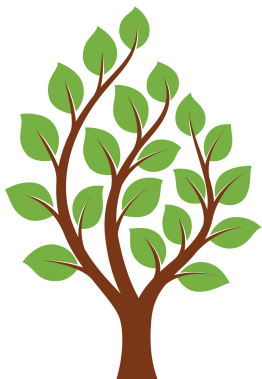
Multiple trees!



Random Forest



Many more columns!

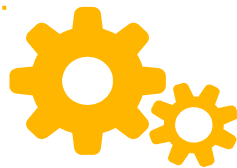


Tree #1

weight
130
115
120
200
190
123

pit	seeds
yes	no
no	yes
no	yes
no	no
no	yes
no	no

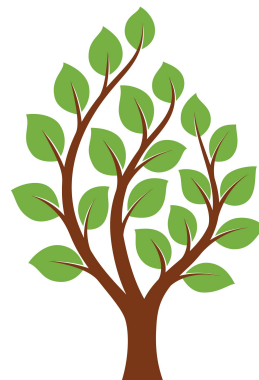
orange
0
1
0
0
0
0



Many more columns!

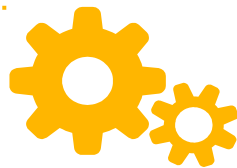
citrus	sugar
yes	10
yes	11
yes	7
yes	10
yes	6
yes	5

Col out
yellow
orange
green
yellow
yellow
green



Tree #2

orange
0
1
0
0
0
0



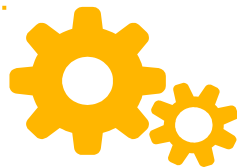
Many more columns!



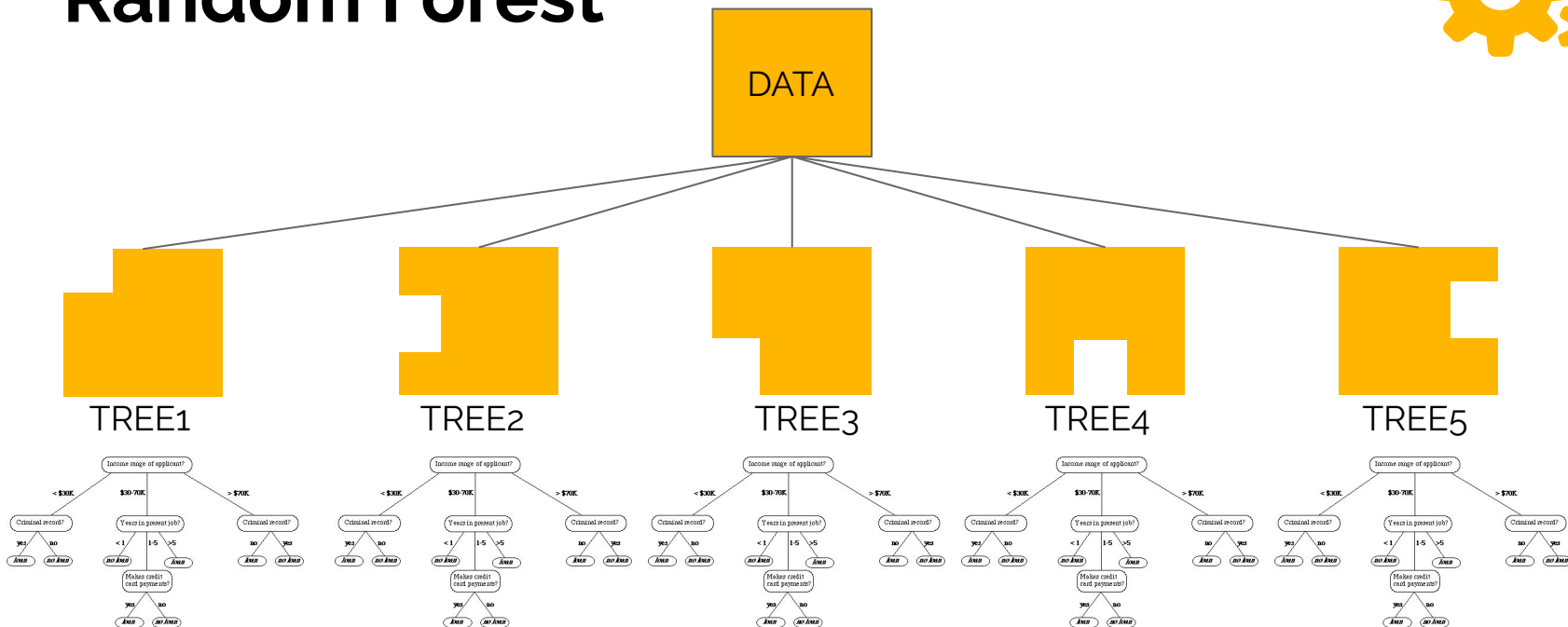
Tree #3

Col in
yellow
orange
red
yellow
white
green

round	soft	orange
yes	yes	0
no	no	1
no	yes	0
no	no	0
yes	no	0
no	no	0

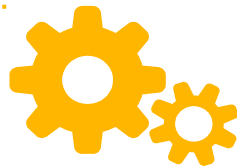


Random Forest

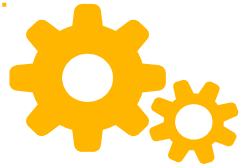


Wisdom of the crowds!

Random Forest

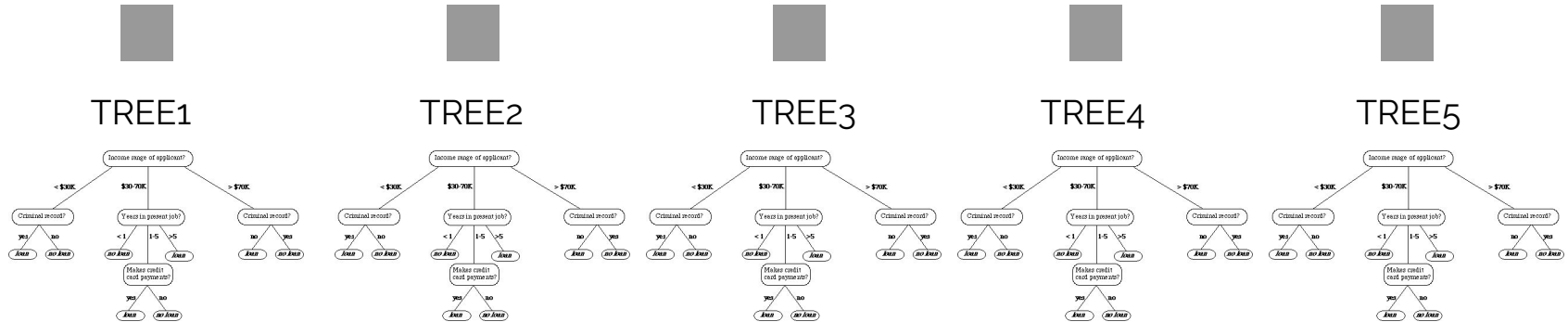


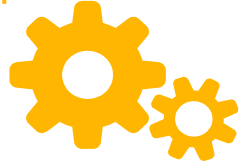
New Fruit



Random Forest

 New Fruit

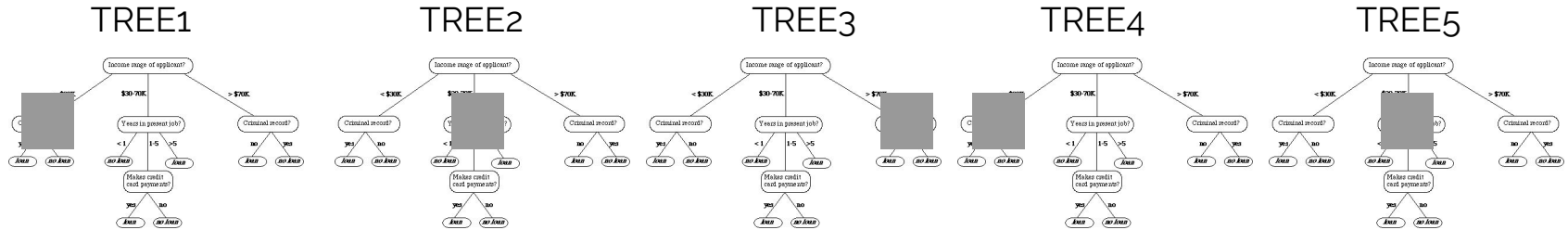




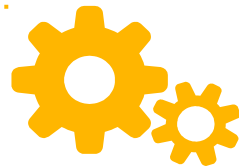
Random Forest



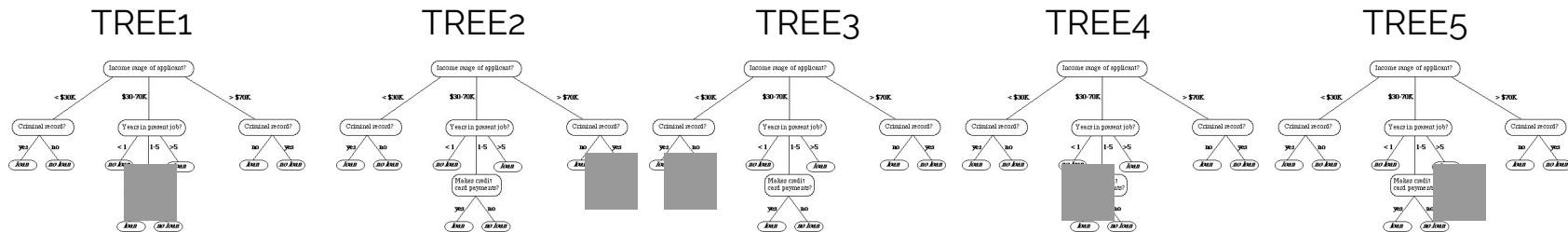
New Fruit



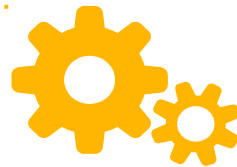
Random Forest



New Fruit

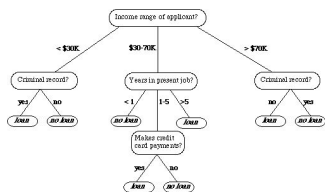


Random Forest

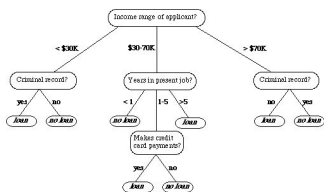


New Fruit

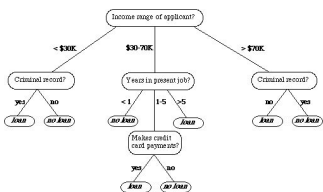
TREE1



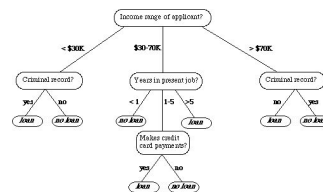
TREE2



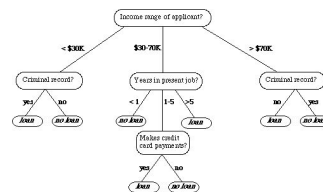
TREE3

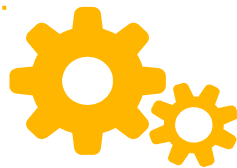


TREE4



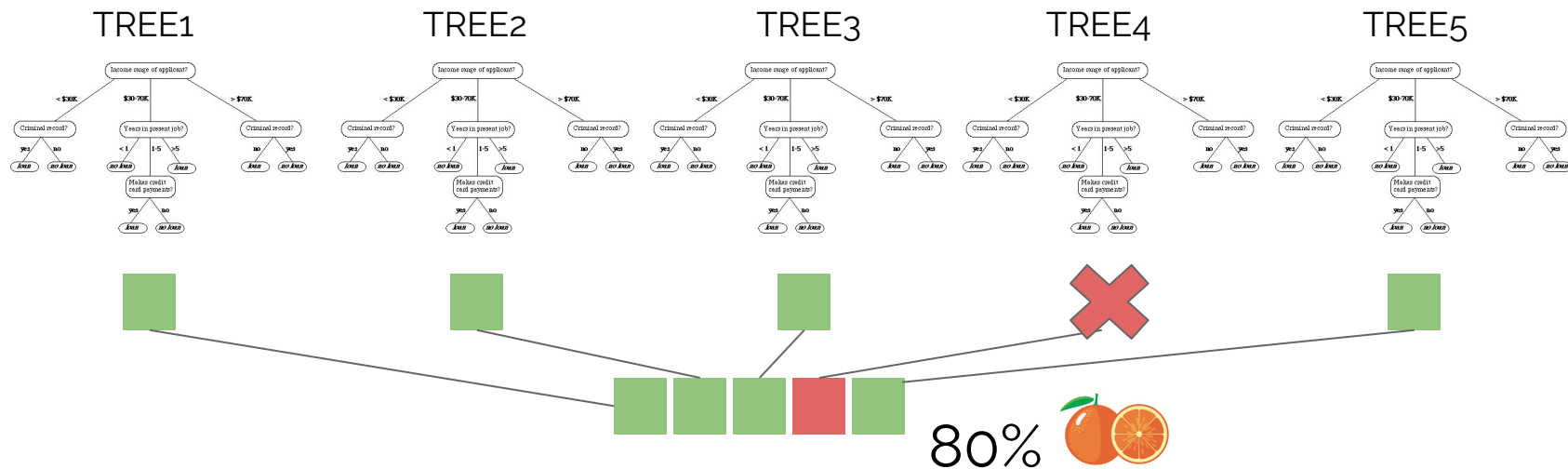
TREE5

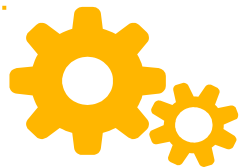




Random Forest

 New Client





Random Forest Parameters

max_depth	n_estimators
Maximum number of questions asked for each branch	Number of trees to grow.



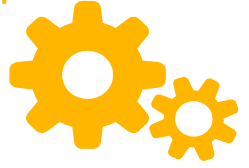


Random Forest Implementation

Python Exercise



Model Training



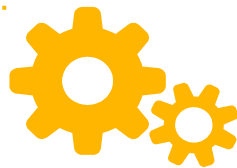
Model Performance

1

$Y, X_1, X_2, X_3, X_4, \dots, X_n$

SAMPLES





Model Performance

1

SAMPLES

$Y, X_1, X_2, X_3, X_4, \dots, X_n$



2

$Y, X_1, X_2, X_3, X_4, \dots, X_n$

TRAIN

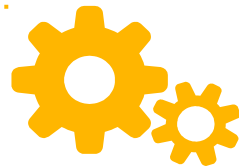
70%

TEST

$Y, X_1, X_2, X_3, X_4, \dots, X_n$

30%

Model Performance



1

SAMPLES

$Y, X_1, X_2, X_3, X_4, \dots, X_n$



2

TRAIN

$Y, X_1, X_2, X_3, X_4, \dots, X_n$

70%

TEST

$Y, X_1, X_2, X_3, X_4, \dots, X_n$

30%



3

TRAIN

$Y_p = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \epsilon$



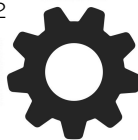
BUILD/TUNE MODEL

Parameter Tuning



Model

Parameter 2



Parameter 4

Parameter 1



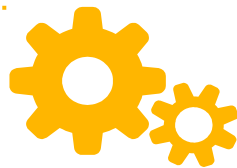
Parameter 3



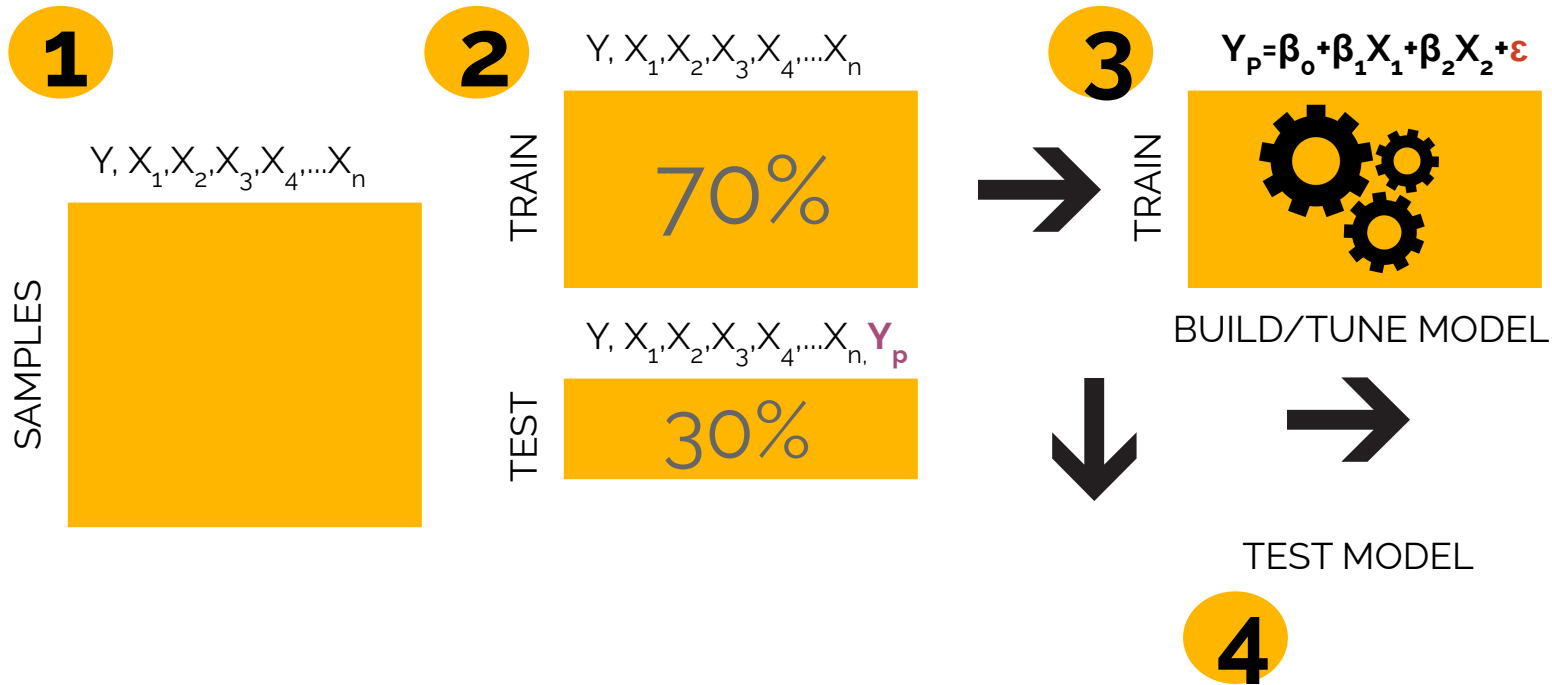
Maximize Performance

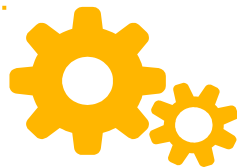


Classification Assessment



Model Performance





Model Performance

1

SAMPLES

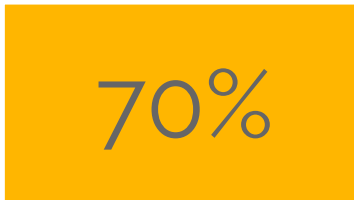
$Y, X_1, X_2, X_3, X_4, \dots, X_n$



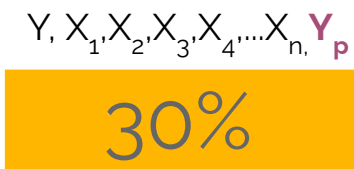
2

$Y, X_1, X_2, X_3, X_4, \dots, X_n$

TRAIN



TEST



Y_p
YES NO

Y	Y_p	
	YES	NO
YES	30	2
NO	3	30

5

ASSESS PERFORMANCE

3

$Y_p = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \epsilon$

TRAIN

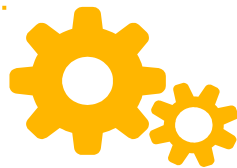


BUILD/TUNE MODEL



TEST MODEL

4



Confusion Matrix

		Predicted Class	
		No	Yes
Observed Class	No	TN	FP
	Yes	FN	TP

TN True Negative
FP False Positive
FN False Negative
TP True Positive

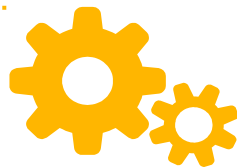
Model Performance

$$\text{Accuracy} = (TN+TP)/(TN+FP+FN+TP)$$

$$\text{Precision} = TP/(FP+TP)$$

$$\text{Sensitivity} = TP/(TP+FN)$$

$$\text{Specificity} = TN/(TN+FP)$$



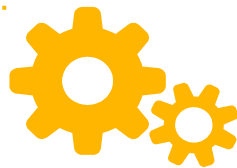
Random Forest

3- Split to Train and Test

```
#split the data to 70% train and 30% test
x_train,x_test,y_train,y_test = train_test_split(X,Y,test_size=0.3,random_state=42)

print(x_train.shape)
print(x_test.shape)
print(y_train.shape)
print(y_test.shape)
```

```
(398, 30)
(171, 30)
(398,)
(171,)
```



Random Forest

4- Train your model: Random Forest

```
rf_model = RandomForestClassifier(max_depth=3,n_estimators=15)
rf_model.fit(x_train, y_train)
rf_model.score(x_train,y_train)
```

```
#define the model
#fit the model (train)
#predict on new observations
```

```
#what is the accuracy of this model?
```

```
0.9849246231155779
```