

Building and Evaluating Predictive Models – Part 1



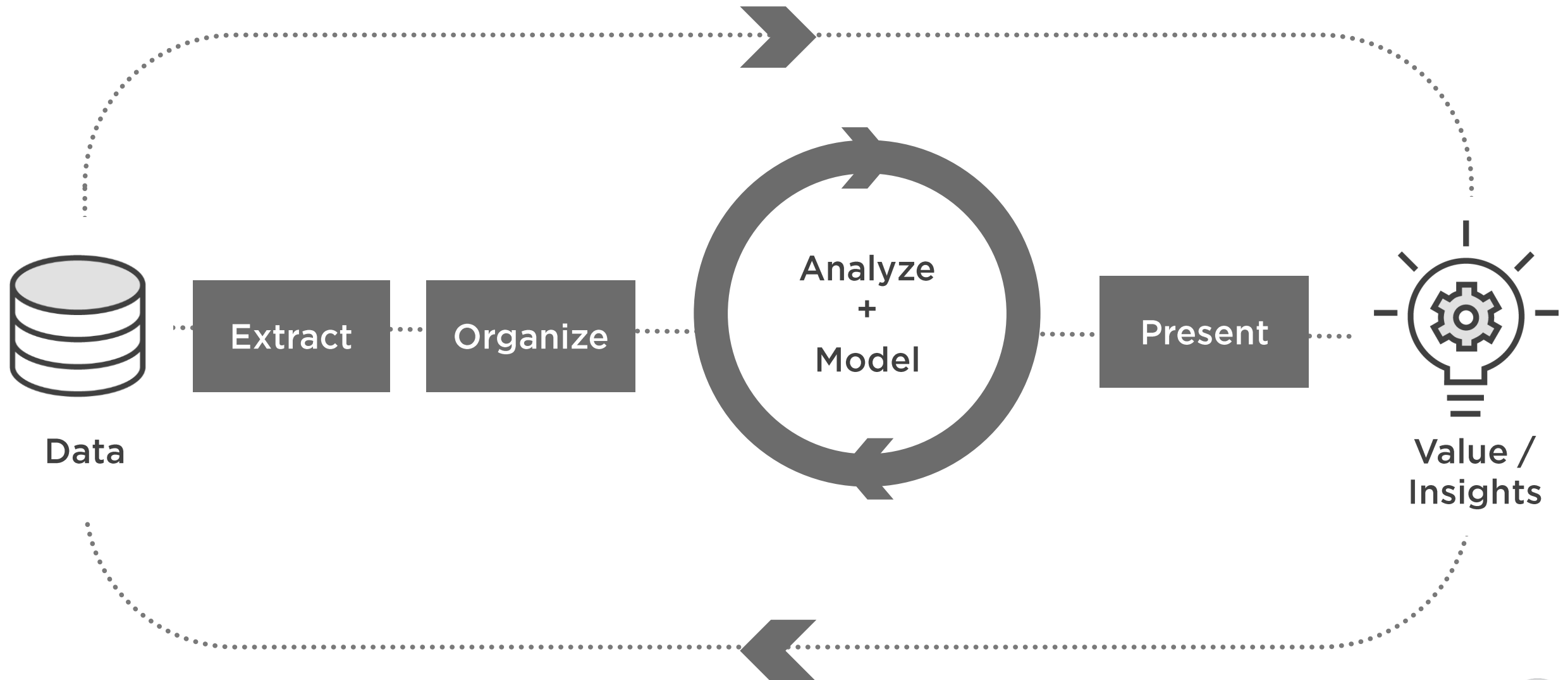
Abhishek Kumar

AUTHOR

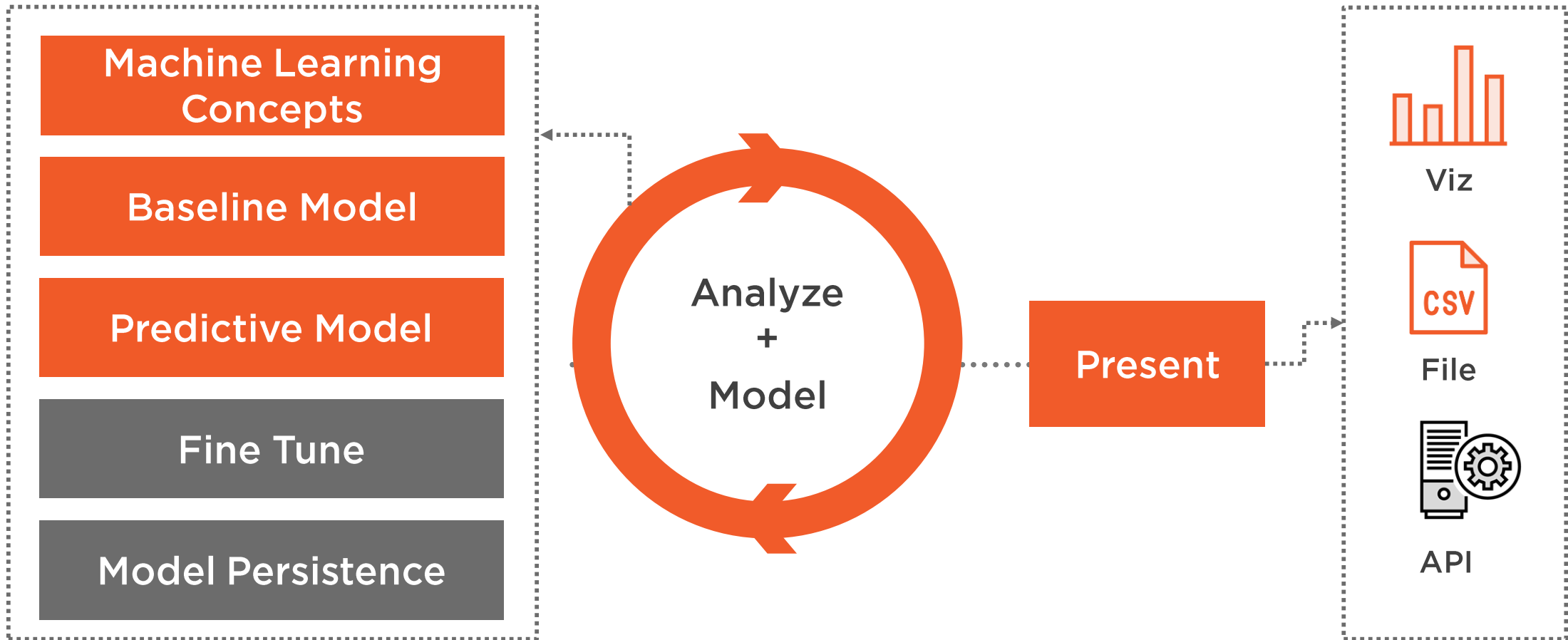
@meabhishekkumar



Data Science Project Cycle



Data Science Project Cycle



Overview (Concepts)

Machine learning basics

Titanic disaster challenge

Classifier

Metrics

Baseline model

Logistic regression model



Overview (Tools)

Python

- Numpy
- Pandas
- Scikit-Learn



Machine Learning Basics



Machine Learning

Learning from data or examples





INBOX (68)

DRAFT

SENT MAIL

SPAM (221)

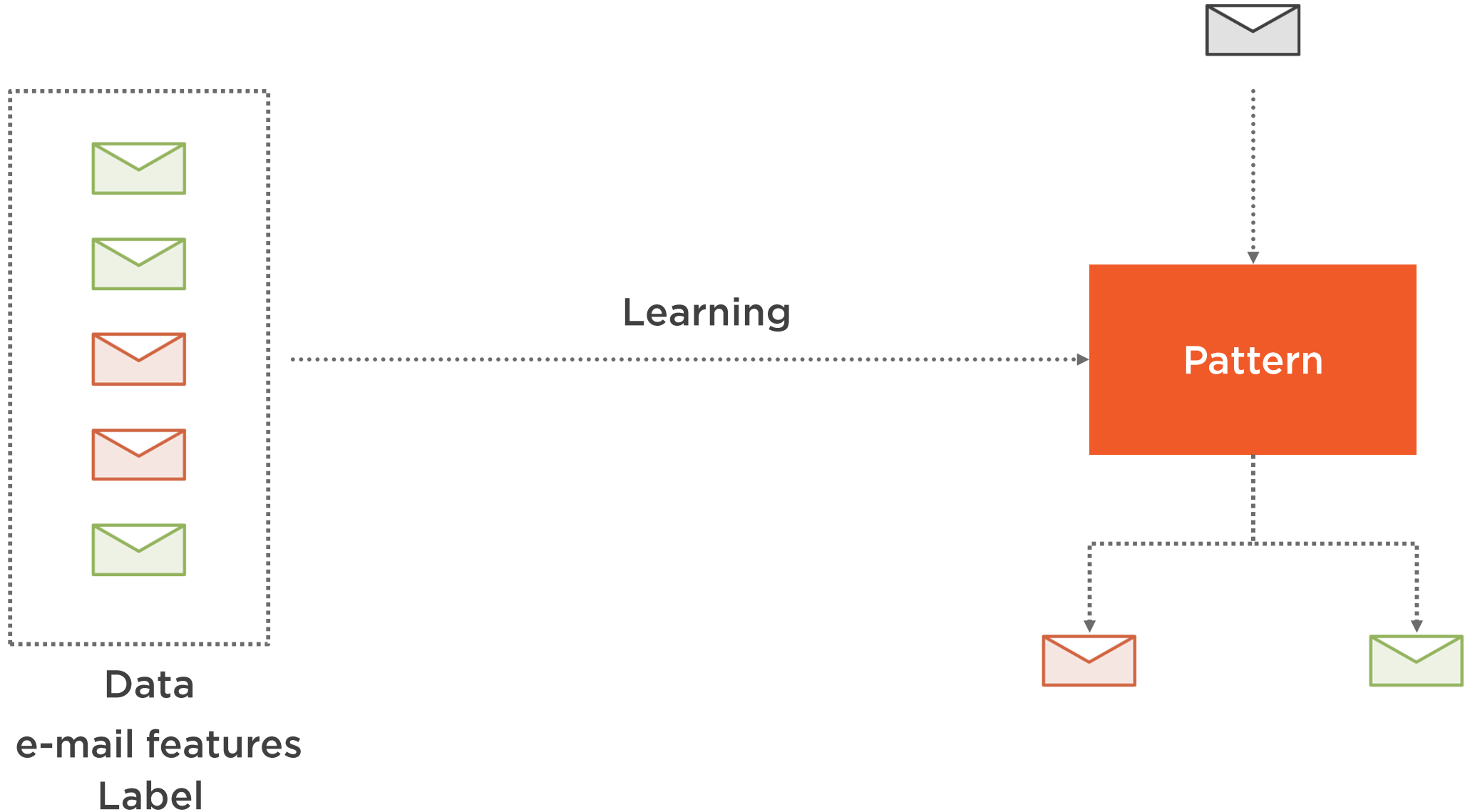
TRASH



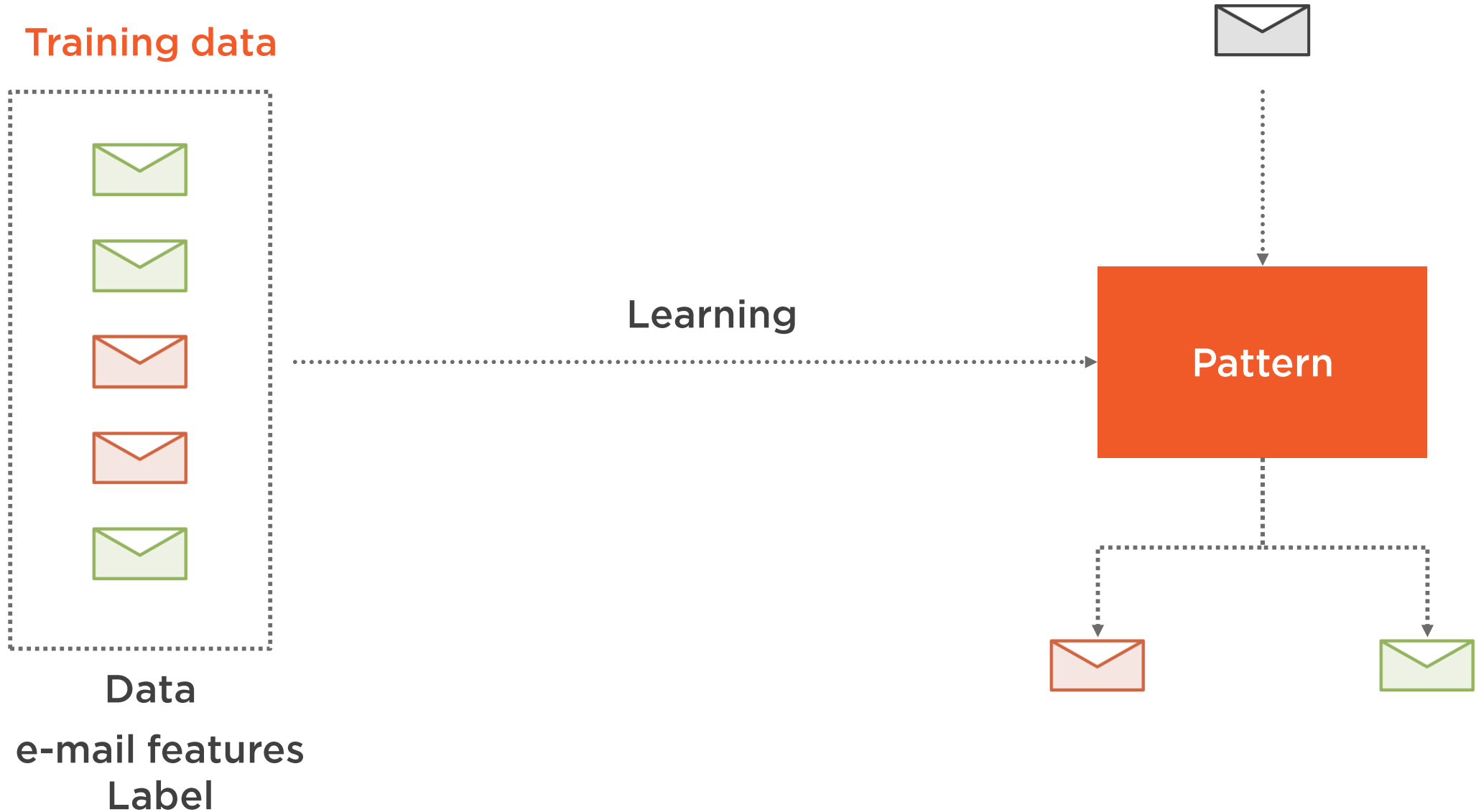
Spam
Detection



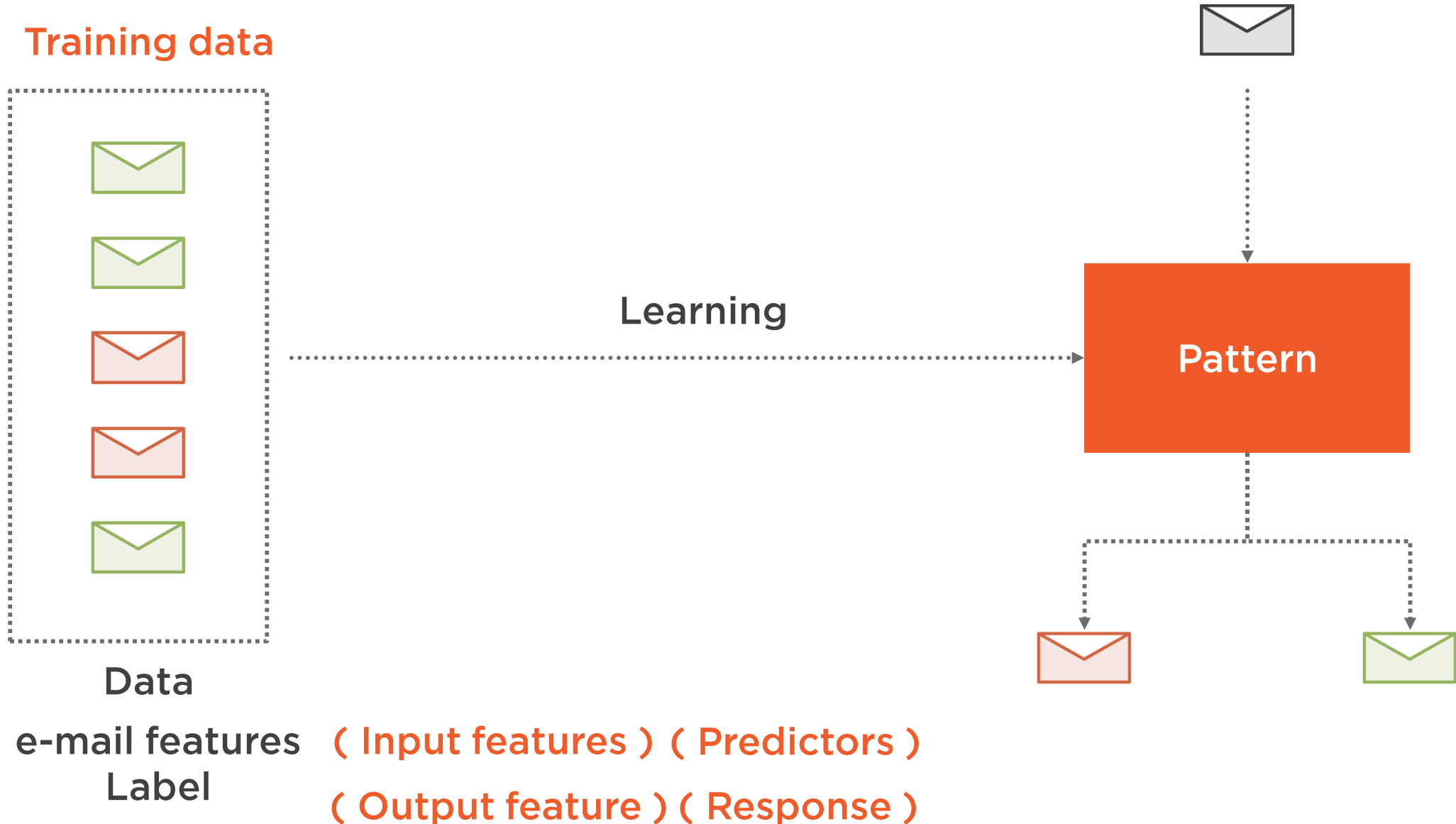
Spam Detection



Training Data



Input and Output Feature

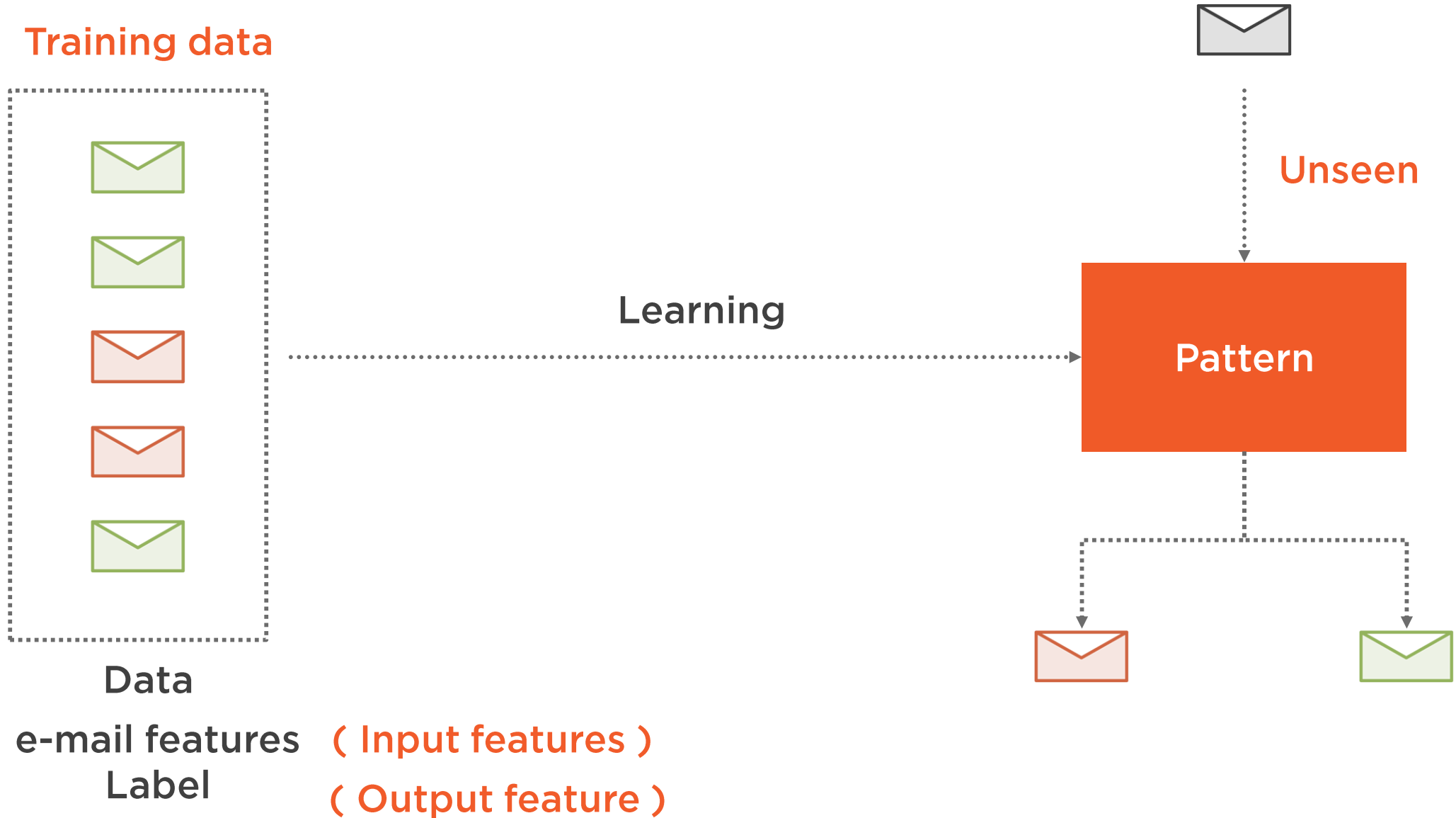


Representation

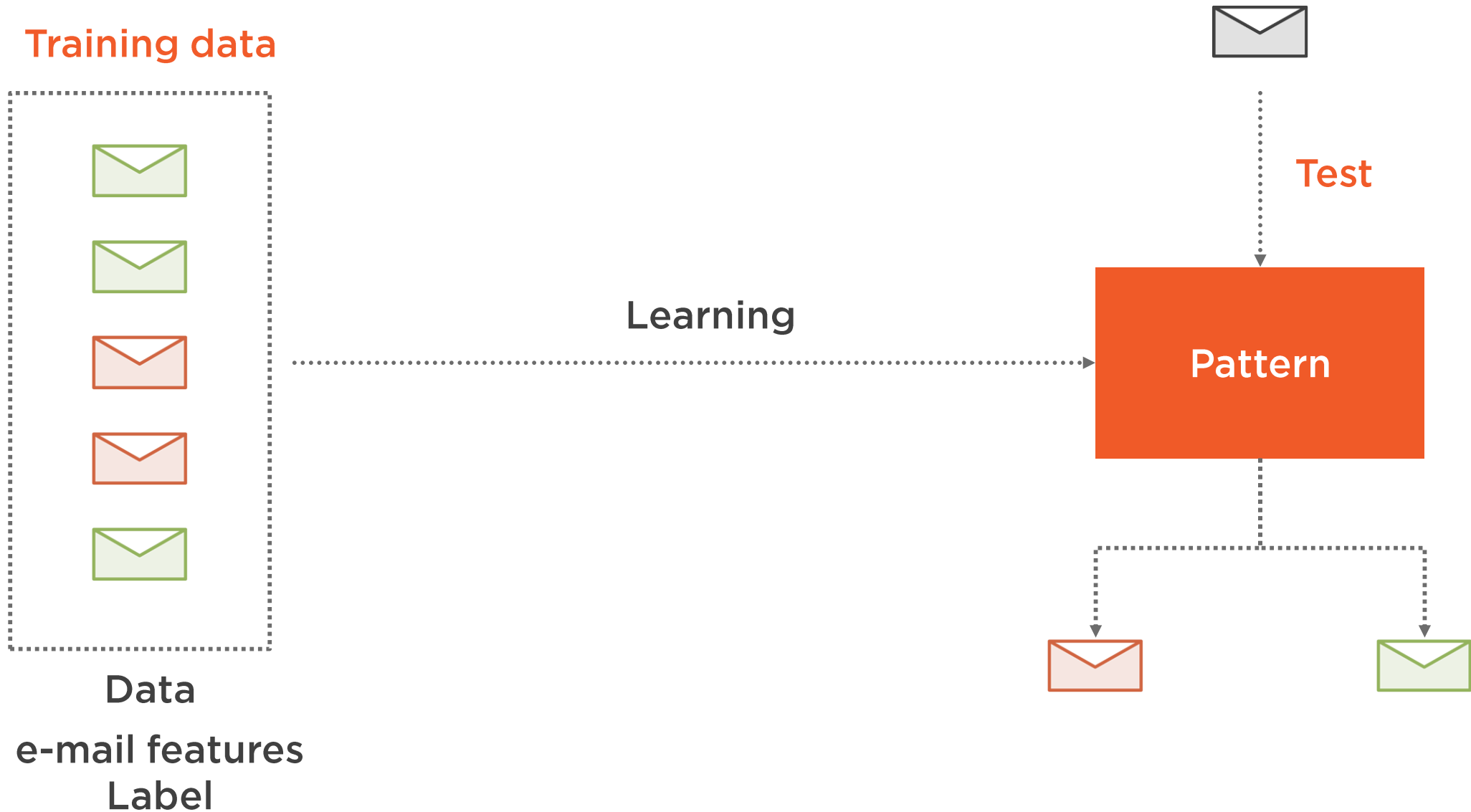
Input				Output
Observation	Sender url	Text	..	Label



Generalization



Generalization



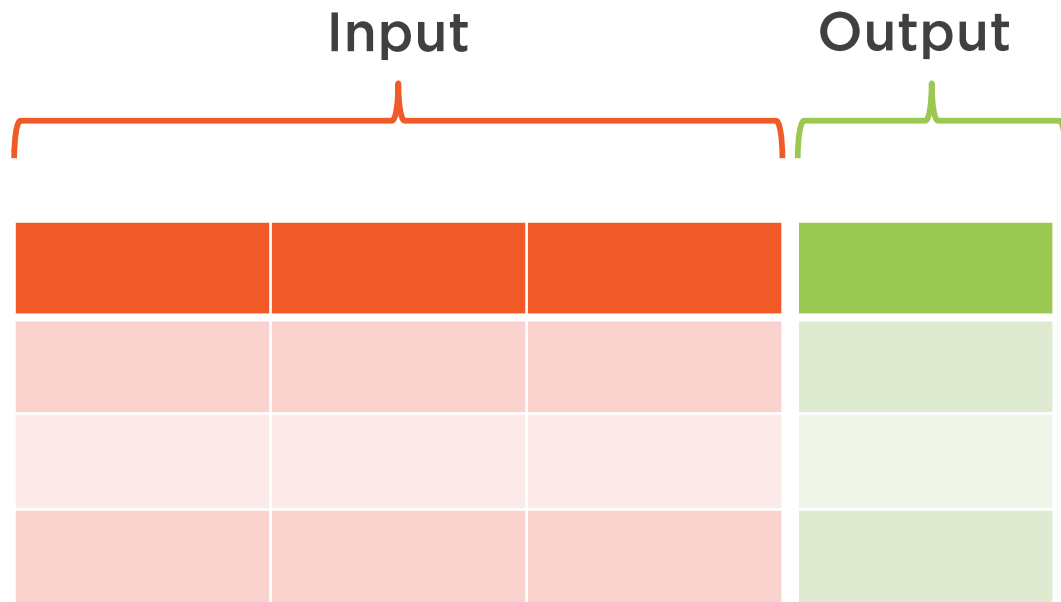
Test Data

					Predicted Output	
Input						
Observation	Sender url	Text	Label	

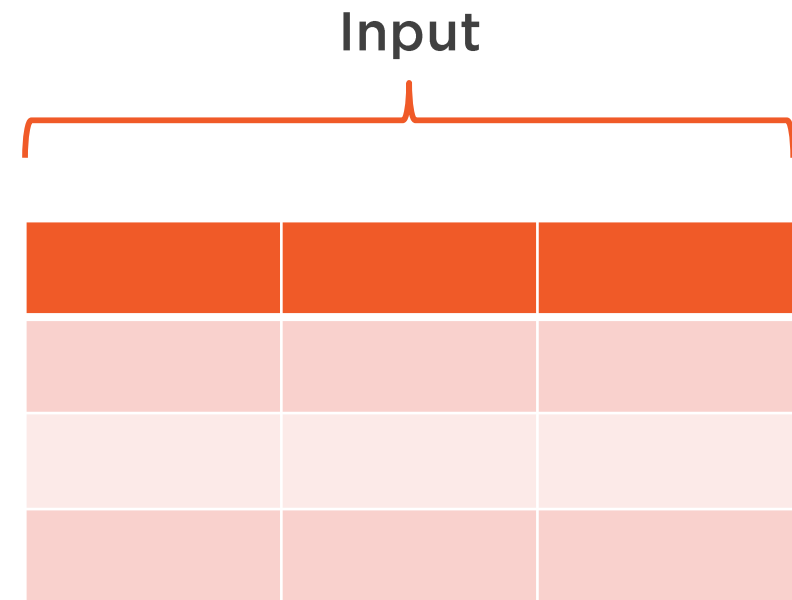


Supervised Learning

Training Data



Test Data



Titanic Disaster

Training Data

Input			Output
Age	Fare	...	Survived

Test Data

Input		
Age	Fare	...

Class : Survived (1)

Class : Not Survived (0)



Classification

Training Data

Input			Output
Age	Fare	...	Survived

Test Data

Input		
Age	Fare	...

Output : Discrete labels

Class : Survived (1)

Class : Not Survived (0)



Regression

Training Data

Input			Output
Height	Width	..	Mileage

Test Data

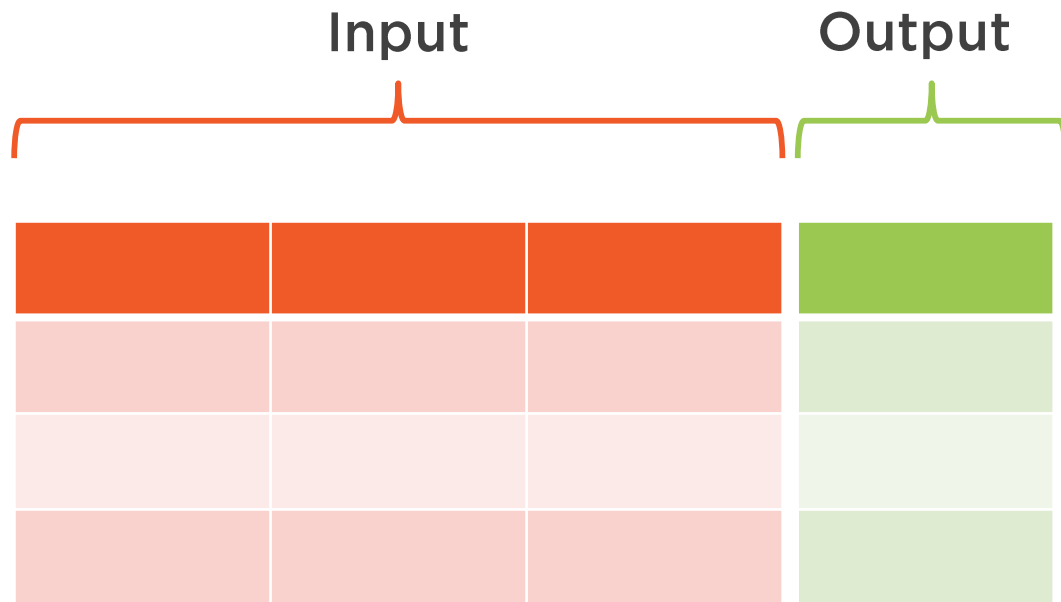
Input		
Height	Width	...

Output : Continuous values

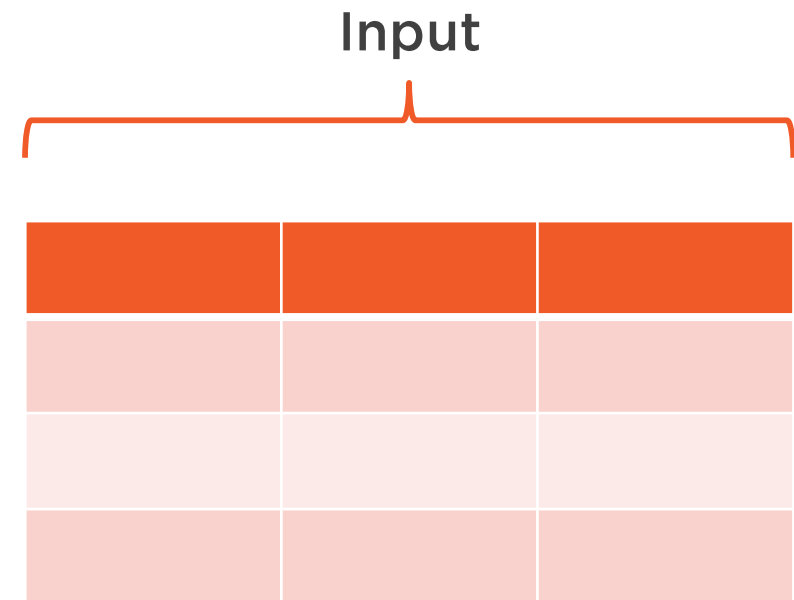


Unsupervised Learning

Training Data



Test Data



Customer Segmentation

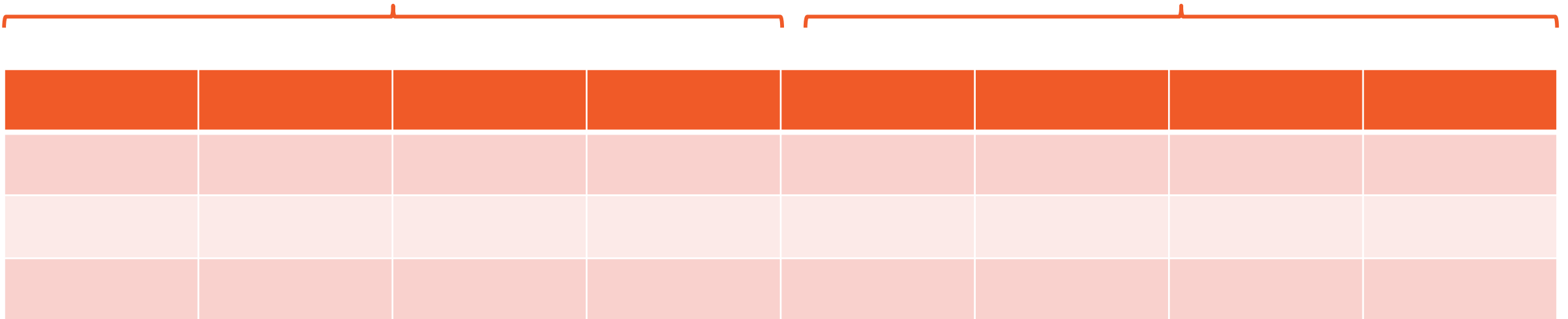


Customer Segmentation

Training Data

Demographic features

Previous purchase history



Demographic features				Previous purchase history			



Customer Segmentation



Clustering



Titanic Disaster Challenge



- Both input and output in training data
- Supervised learning problem
- Classification task
- Binary classification

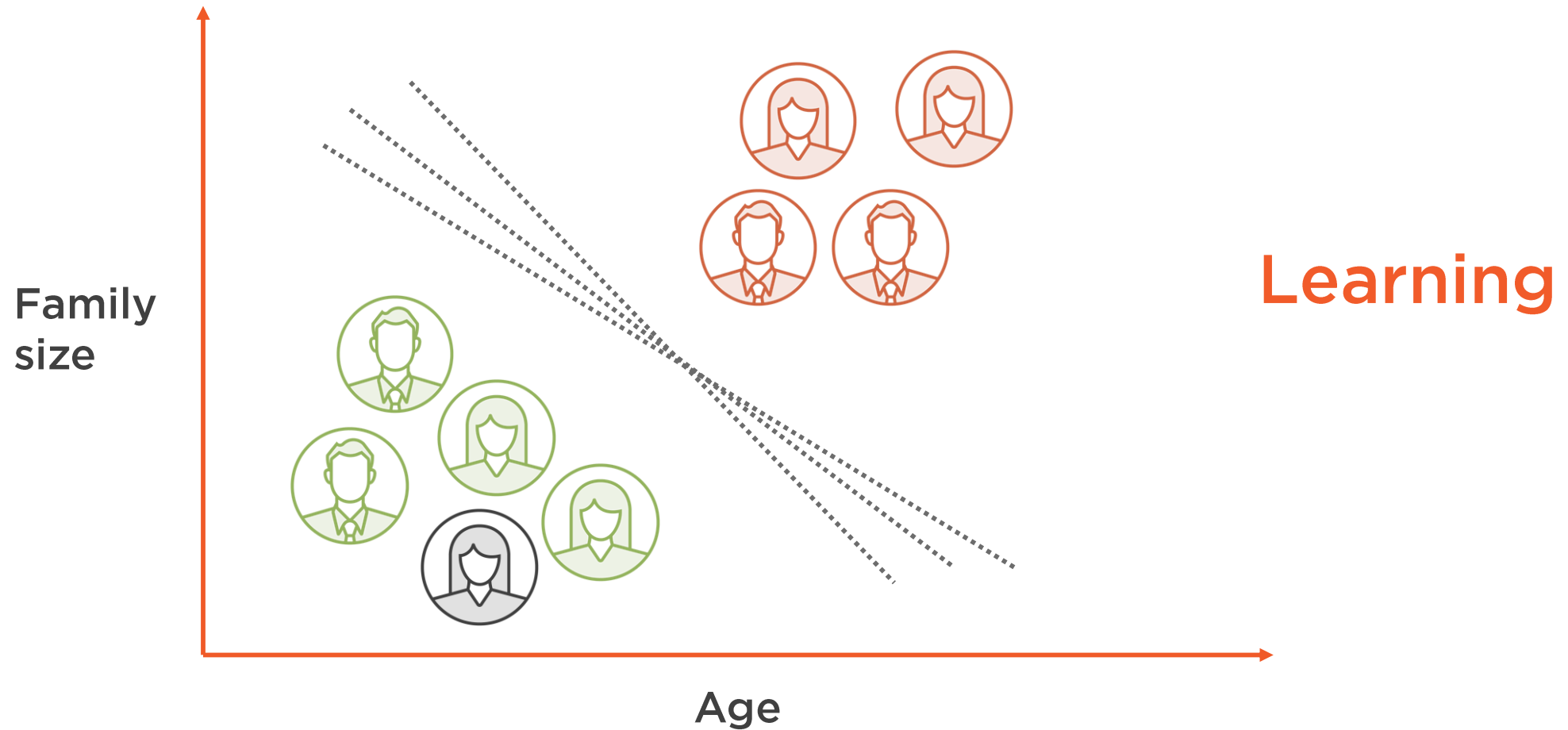
Titanic Disaster Challenge



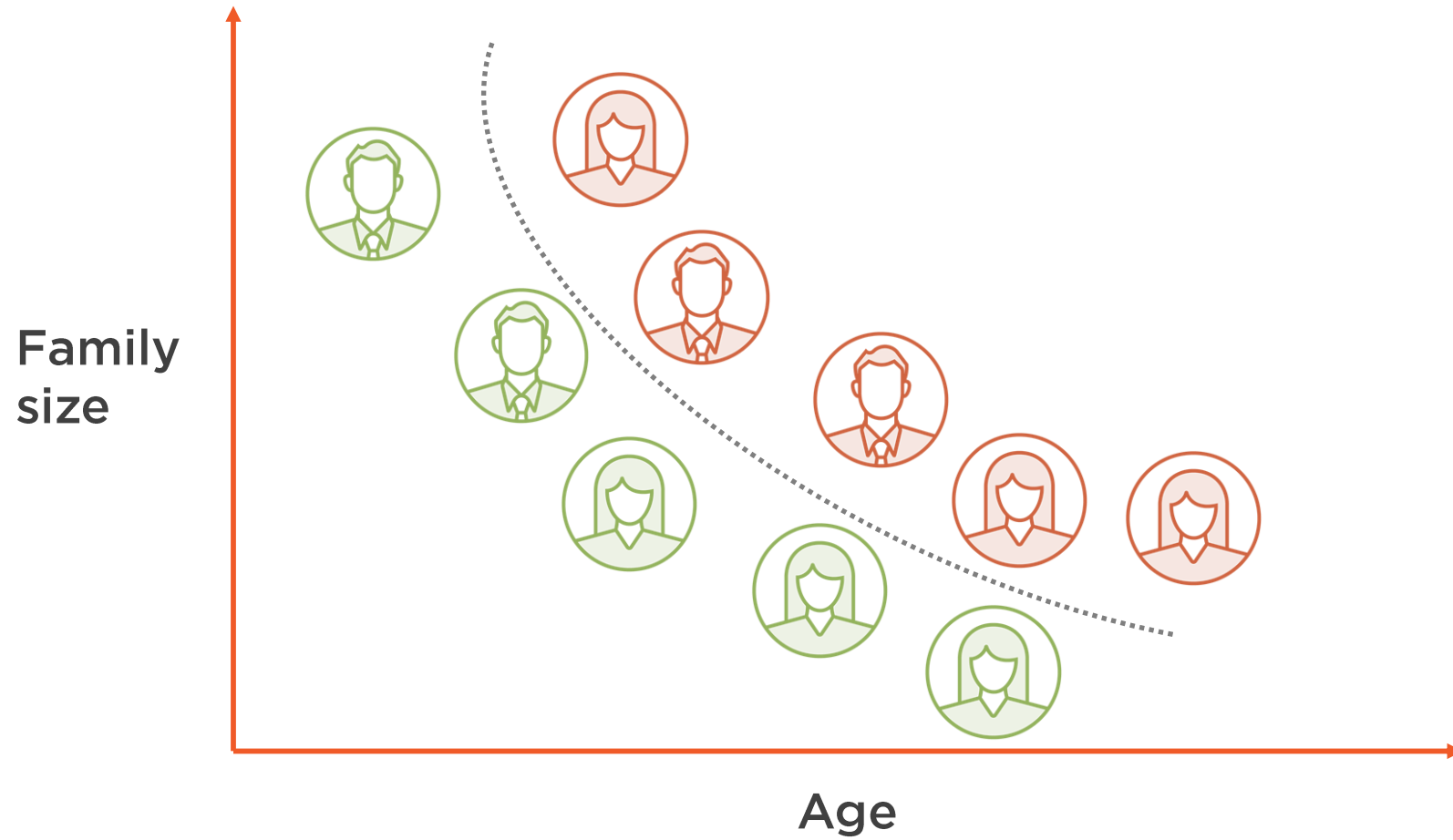
Titanic Disaster Challenge



Classifier

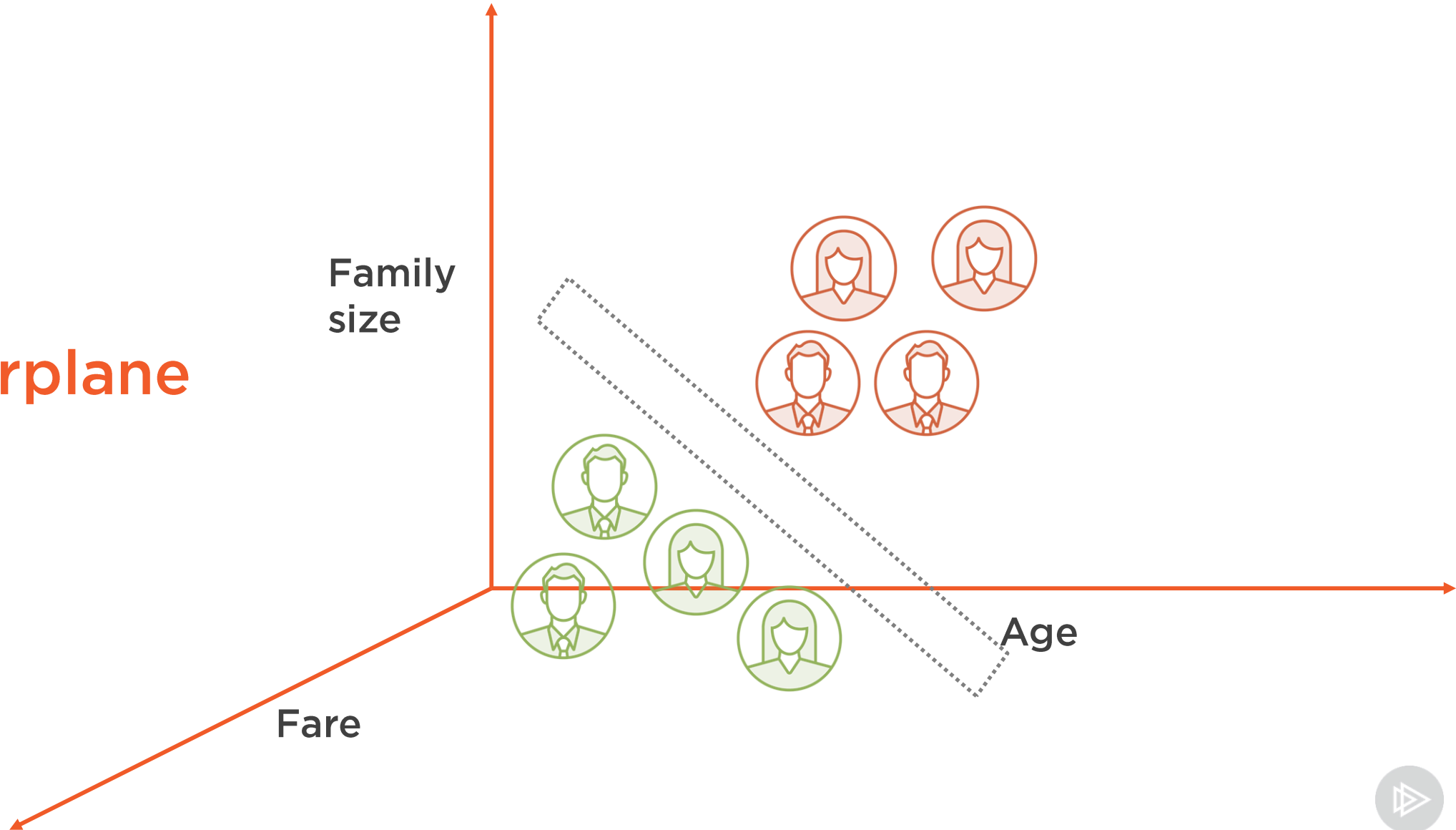


Classifier

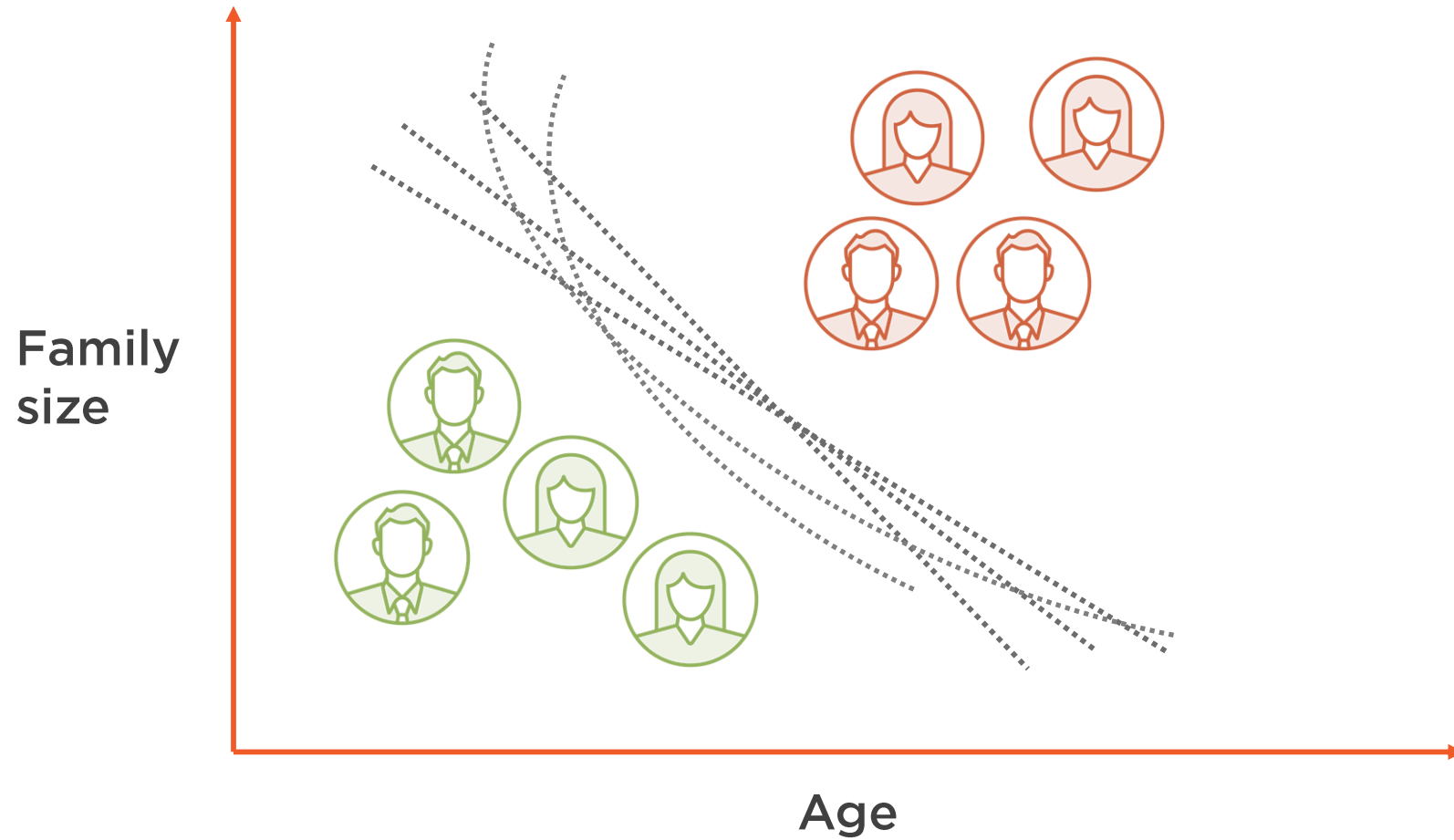


Classifier

Hyperplane



Classifier



Classifiers

Logistic regression

Support vector machine

Neural networks

Random forest



“If you can’t measure it, you can’t improve it.”

Peter Drucker



Performance Metrics

Accuracy

Precision

Recall



Accuracy

Id	F ₁	...	F _n
1			
2			
3			
4			
5			
6			
7			
8			
9			
10			

Test data

Output	Output
1	1
0	0
1	0
1	1
0	1
1	1
1	0
0	1
1	1
1	1

Predicted
output

Actual
output

1

2

3

4

5

6

Accuracy = Correct count / Total Count

Accuracy = 6 / 10 = 0.6 (60%)



Precision

Confusion Matrix

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

	Predicted Negative	Predicted Positive
Actual Negative	20	40
Actual Positive	30	60

Precision

What fraction of positive predictions are correct?

$$\frac{TP}{\text{Total Positive Predictions}}$$

$$\frac{TP}{TP + FP}$$

$$\text{Precision} = 60 / (60 + 40) = 0.6$$



Recall

Confusion Matrix

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

	Predicted Negative	Predicted Positive
Actual Negative	20	40
Actual Positive	30	60

Recall

What fraction of positive cases you predicted correctly?

$$\frac{TP}{\text{Total Positive Cases}}$$

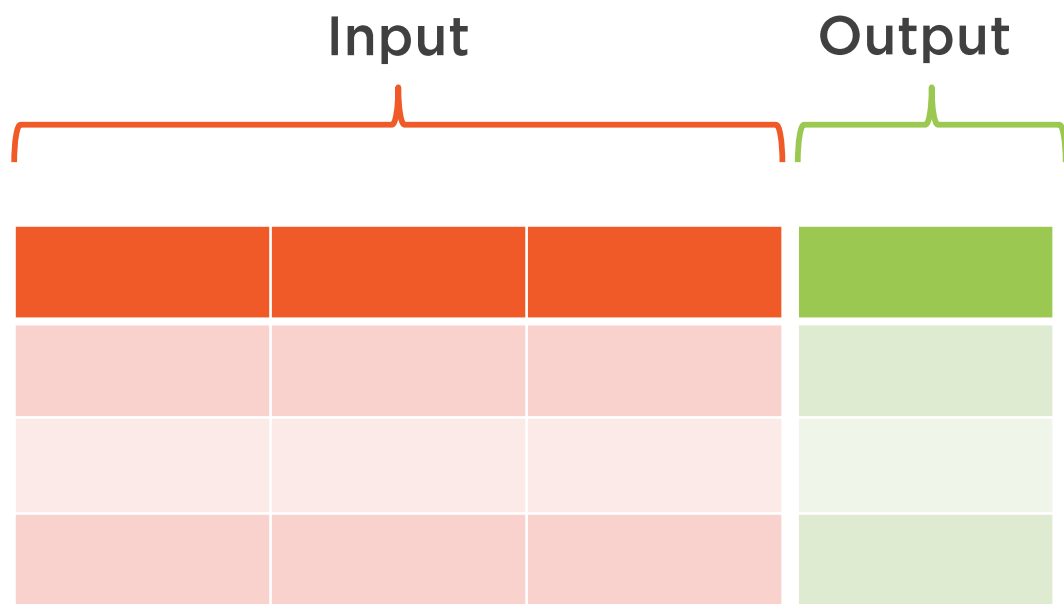
$$\frac{TP}{TP + FN}$$

$$\text{Recall} = 60 / (60 + 30) = 0.67$$

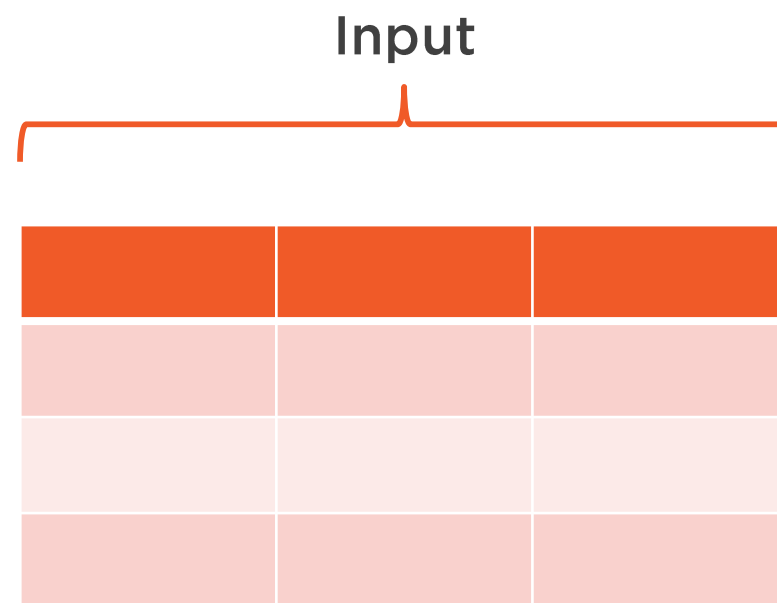


Classifier Evaluation

Training Data



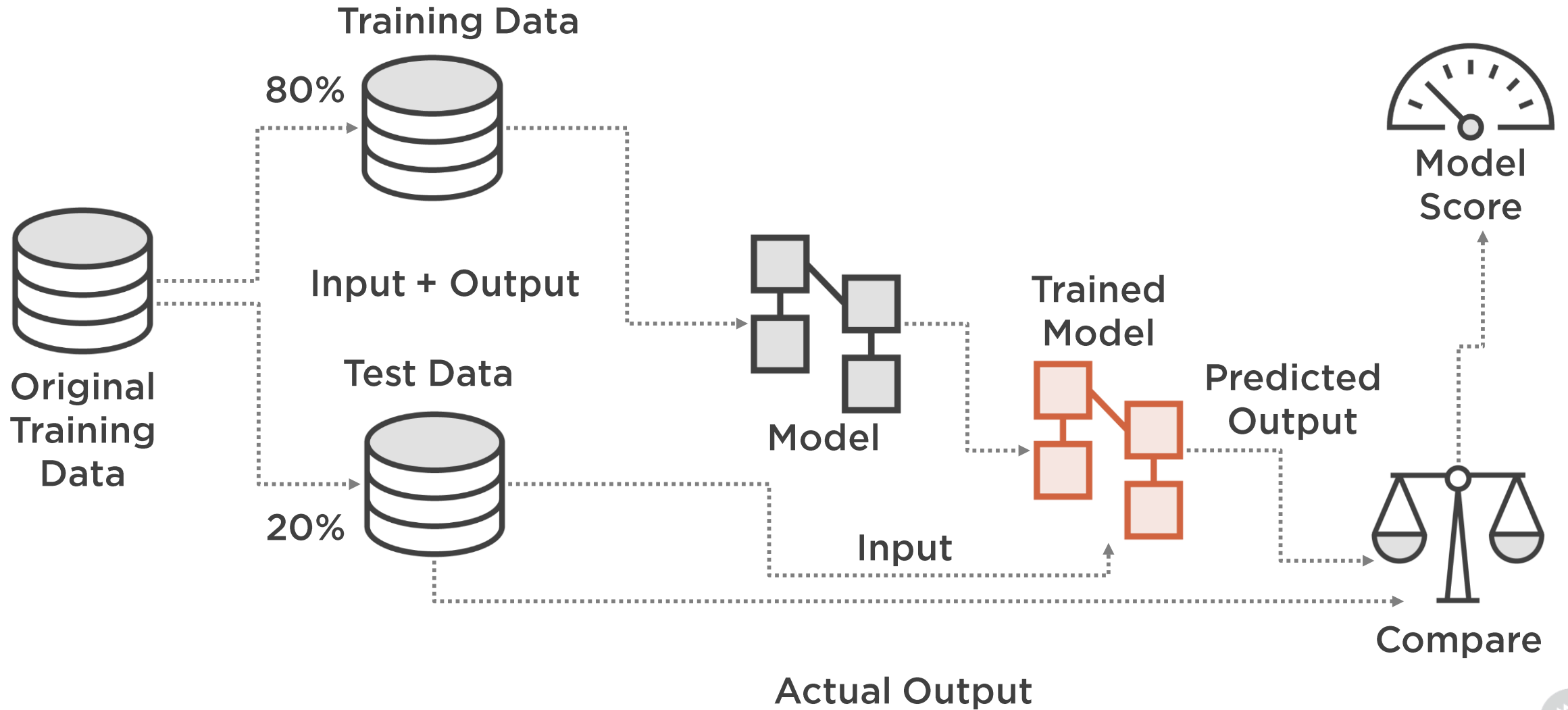
Final Test Data



Train Test Split



Classifier Evaluation



Baseline Model



First Step



Don't Skip



Compare



Baseline Model for Classification

Class	Count
1	60
0	40

Baseline model = Class 1

Baseline model accuracy = $60 / (60 + 40) = 0.6$

- Output majority class
- Predictive model should have better performance than baseline



Demo



Preparing data for machine learning model



Demo



Building and evaluating baseline model



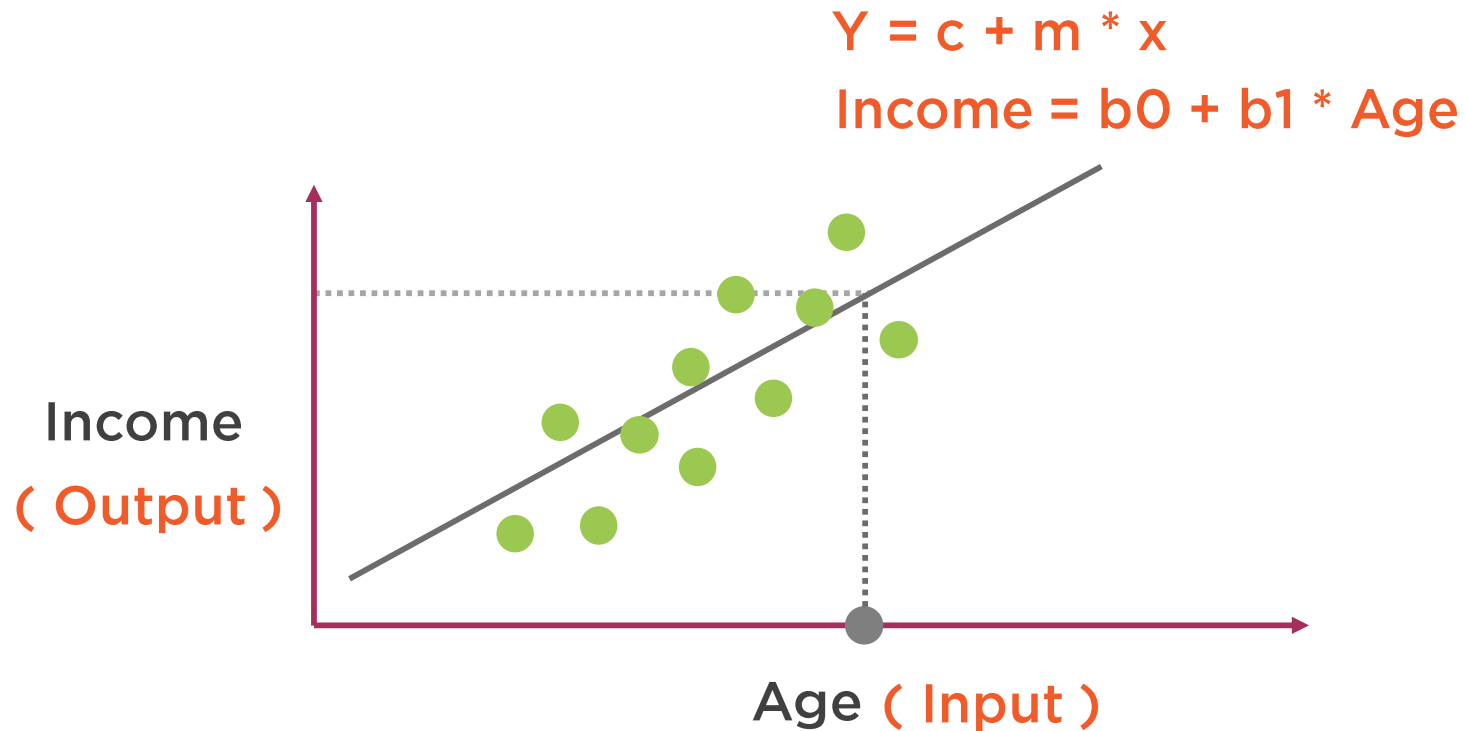
Demo



Making first Kaggle submission

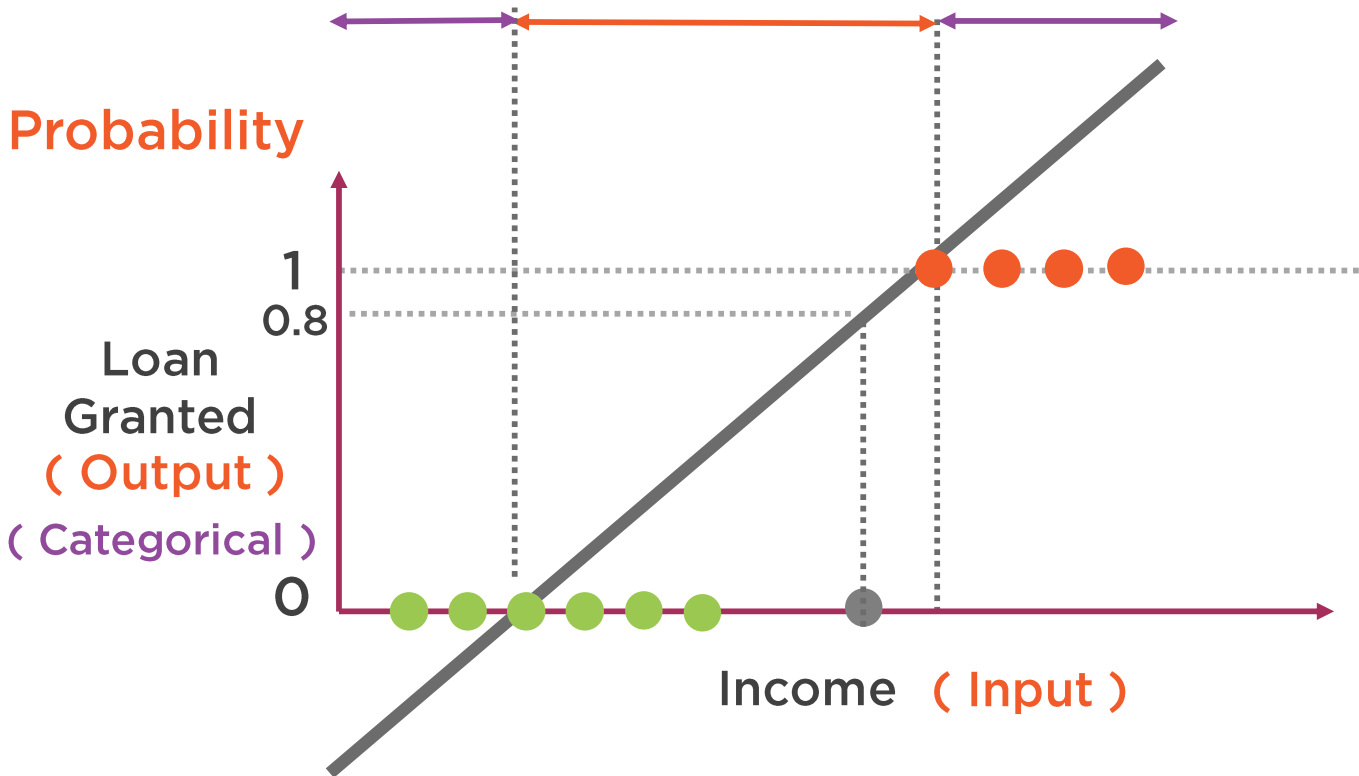


Linear Regression Model



- Supervised learning problem
- Regression task
- Model coefficients : b_0 , b_1

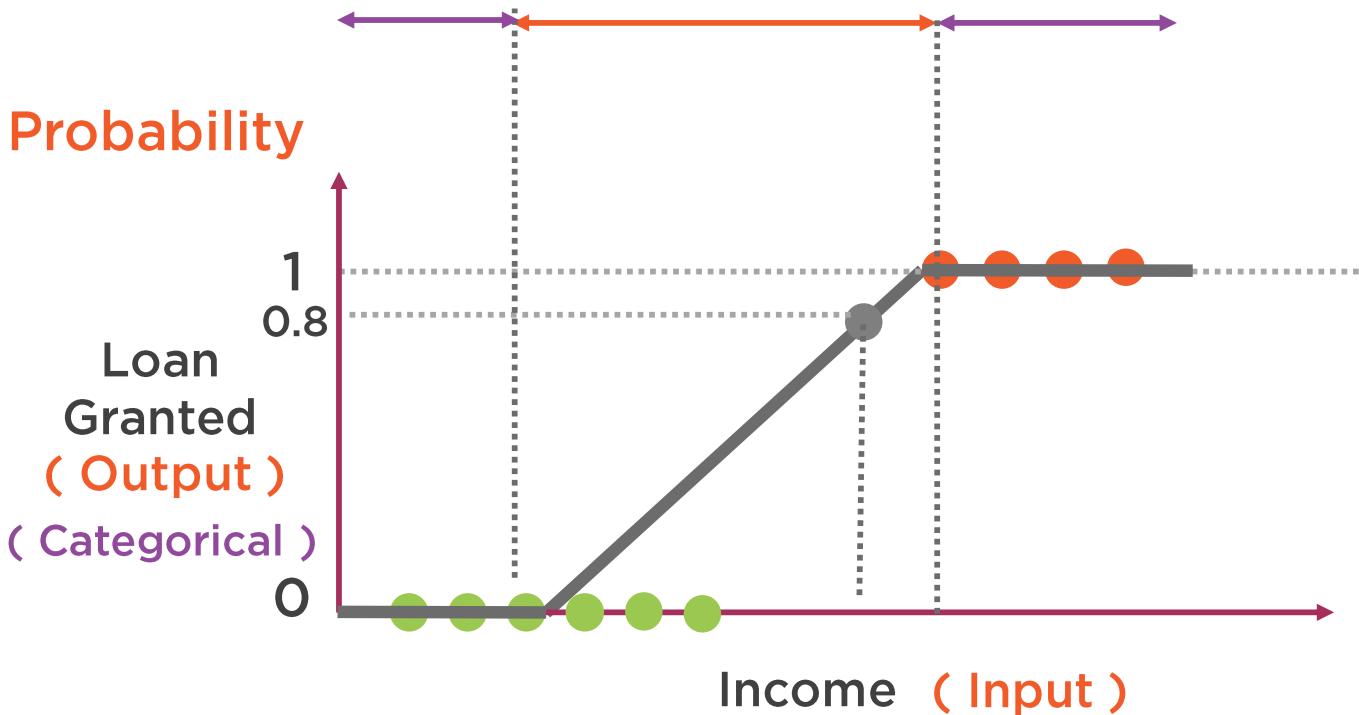
Logistic Regression Model



- Supervised learning problem
- Classification task
- Binary classification



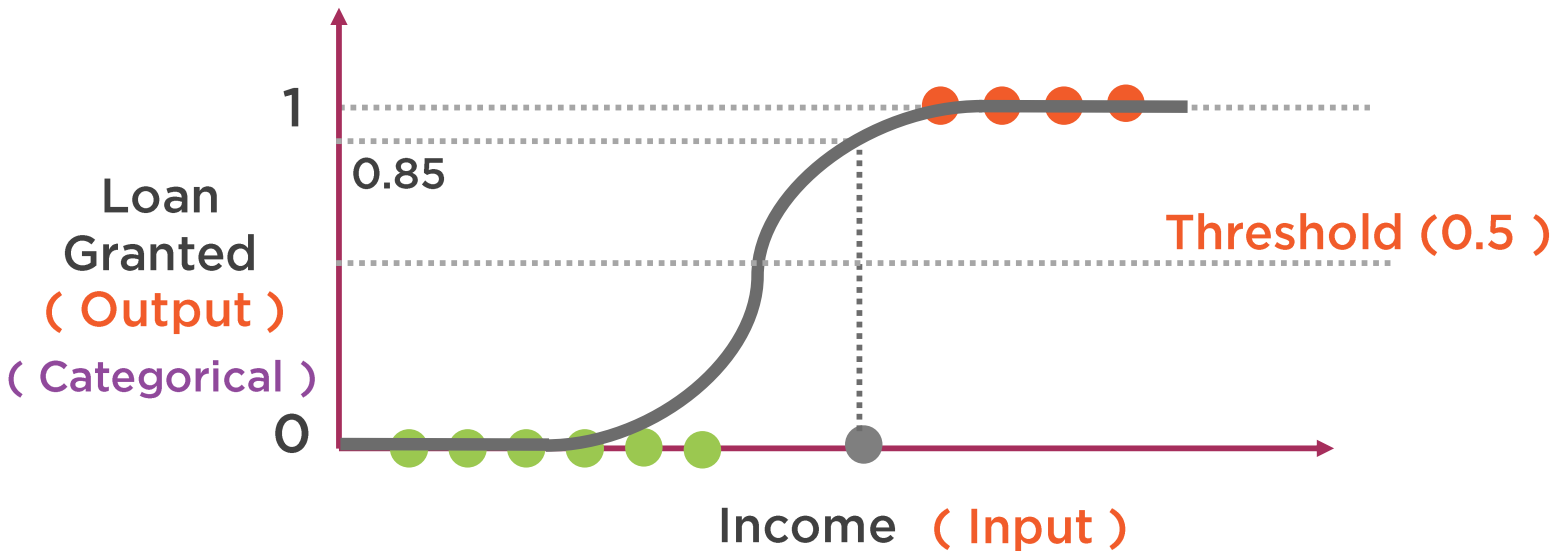
Logistic Regression Model



- Supervised learning problem
- Classification task
- Binary classification

Logistic Regression Model

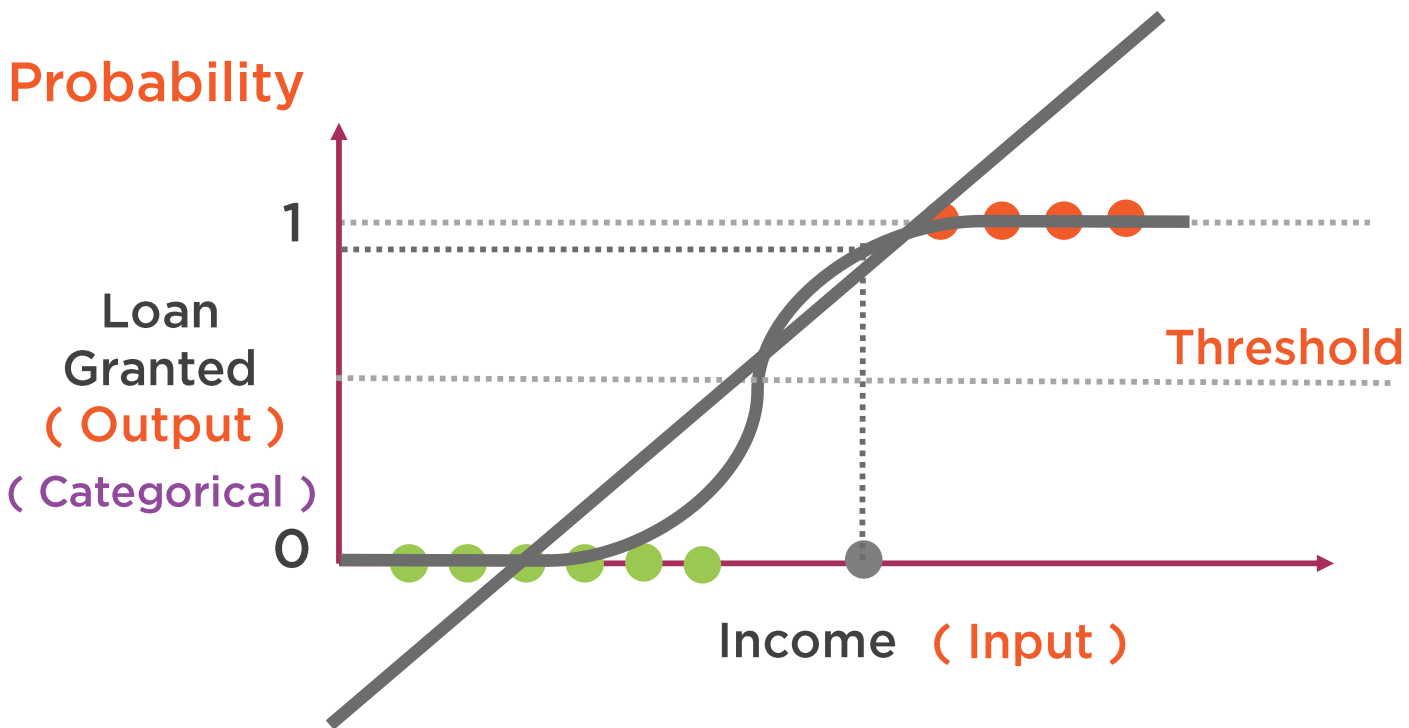
Probability



- Supervised learning problem
- Classification task
- Binary classification
- Sigmoidal curve



Logistic Regression Model



$$\text{Granted} = b_0 + b_1 * \text{income} \quad 1$$

$$P = \frac{1}{1 + e^{-\text{Granted}}} \quad 2$$

$P > \text{threshold} : \text{Class 1}$

$P \leq \text{threshold} : \text{Class 0}$



Demo



Building logistic regression using Scikit-Learn



Demo



Making second Kaggle submission



Summary



Machine learning foundation

Baseline model

Logistic regression model

