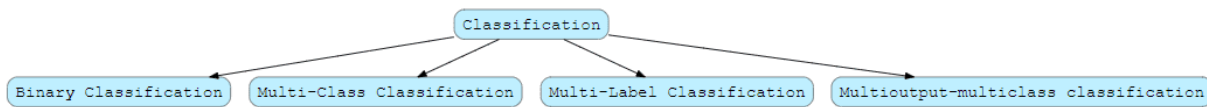
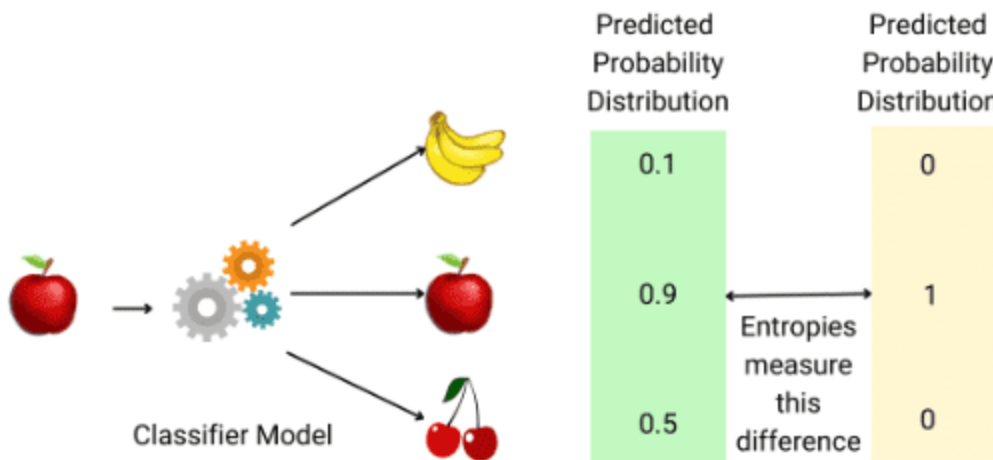


# Types of Classification Problems



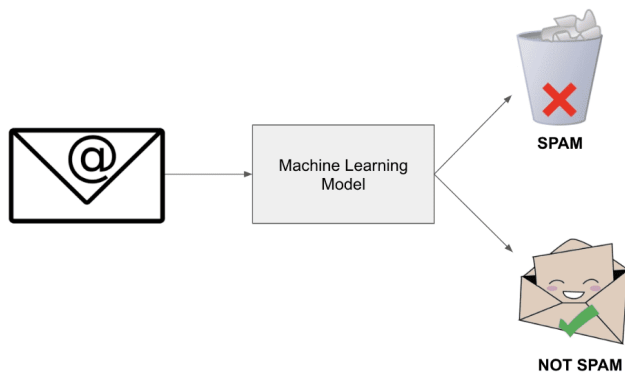
## What is Classification?

- **Classification Problem:** A type of machine learning problem where the goal is to predict the category or class of given data points.
- **Output Variable:** In a classification problem, the output variable is a category or a discrete value.
- **Mapping Function:** The classification predictive modeling approximates a mapping function ( $f$ ) from input variables ( $X$ ) to discrete output variables ( $y$ ).

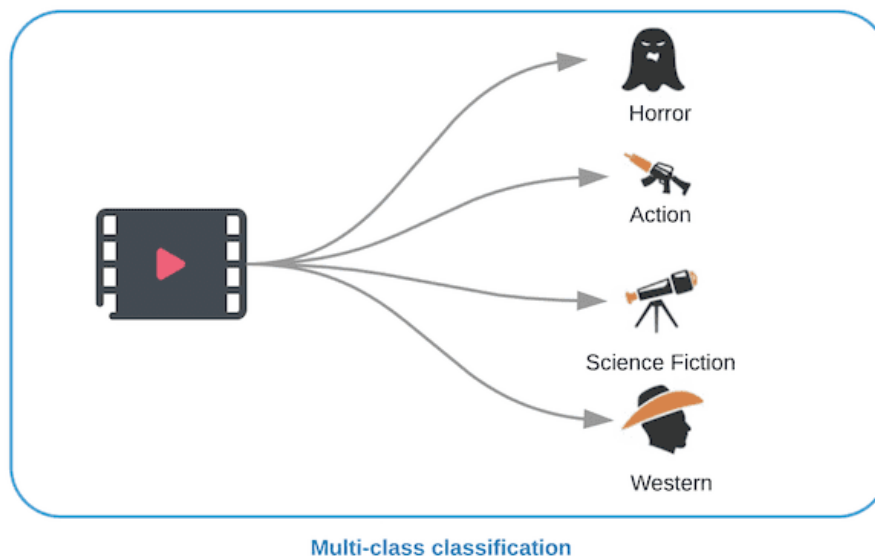


## Classification problems in Machine Learning

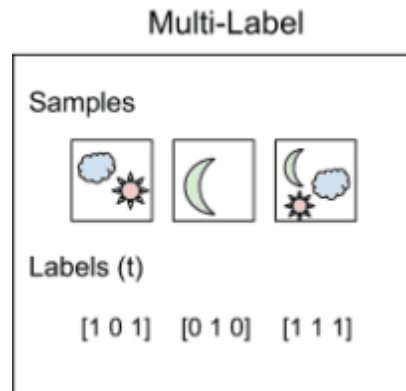
- **Binary Classification:** This is a type of supervised learning where instances are labeled with two classes. For example, a medical test could be used to determine whether a patient has a particular disease (positive) or not (negative)



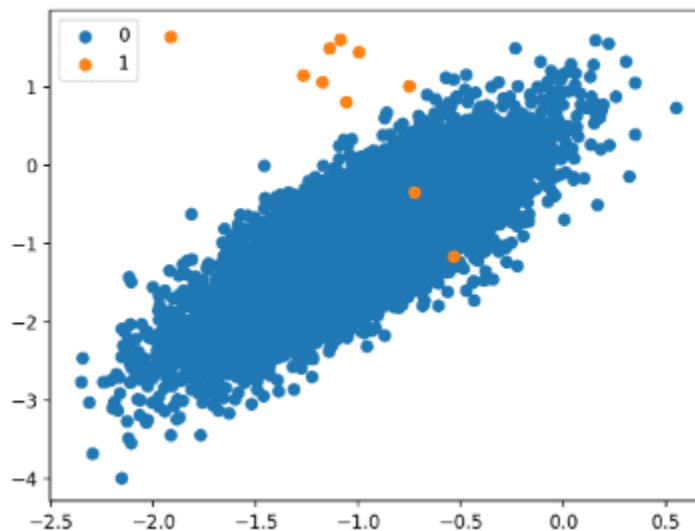
- **Multi-class Classification:** In this type of classification, instances are labeled with more than two classes. For instance, classifying a city's weather into one of the following classes: dry, hot, cold, or humid



- **Multi-label Classification:** This involves predicting more than one label for each instance. For example, a picture might consist of multiple objects like a car, an airplane, and a train, and the task might be to predict all the labels of the image



- Imbalanced Classification:** This refers to classification tasks where the number of instances in each class are unequally distributed. For example, in fraud detection, most of the transactions that happen daily are normal, and fraudulent transactions are rare.



## Strategies for Multi-class and Multi-label Classification

- One-vs-Rest (One-vs-All):** This strategy fits one binary classification model for each class vs. all other classes  
 For example, if we have a multi-class classification problem with classes red, green, and blue, the binary classification can be categorized as

follows:

- Problem one: red vs. green/blue
- Problem two: blue vs. green/red
- Problem three: green vs. blue/red

This means that for each class, a binary classifier is trained against all the other classes. The classifier with the highest confidence in its prediction is typically chosen as the final output.

- **One-vs-One:** This strategy fits one binary classification model for each pair of classes

For instance, given a multi-class classification problem with classes red, blue, and green, this could be divided into three binary classification datasets as follows:

- Problem one: red vs. blue
- Problem two: red vs. green
- Problem three: blue vs. green

Each binary classification model may predict one class label and the model with the most predictions or votes is predicted by the one-vs-one strategy.

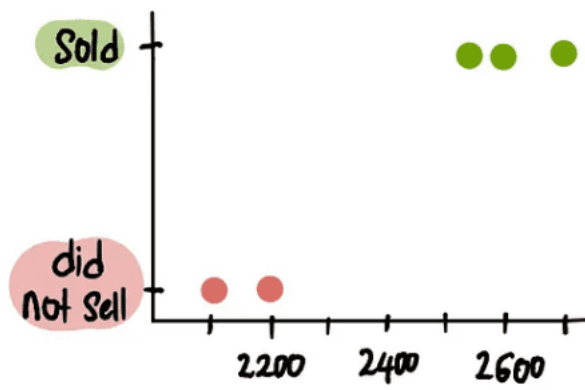
## Logistic Regression

- Logistic regression allows us to estimate the probability of a given instance.
- This probability will range from 0 to 1.
- This probability will be capped by a decision boundary and associated with a discrete positive or negative class.

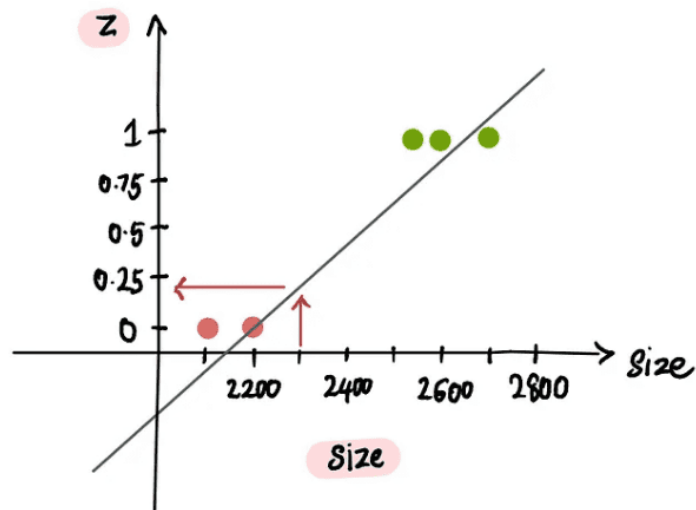
house size	did it sell?
2100 feet <sup>2</sup>	0
2200 feet <sup>2</sup>	0
2550 feet <sup>2</sup>	1
2600 feet <sup>2</sup>	1
2700 feet <sup>2</sup>	1

0 = did not sell  
1 = sold

If this is plot ,



Now let plot as probabilities



So what is required to get this probability ?

predicted probability of the house selling

$$\hat{p} = g(\beta_0 + \beta_1 \times \text{size}) = \frac{1}{1 + e^{-(\beta_0 + \beta_1 \times \text{size})}}$$

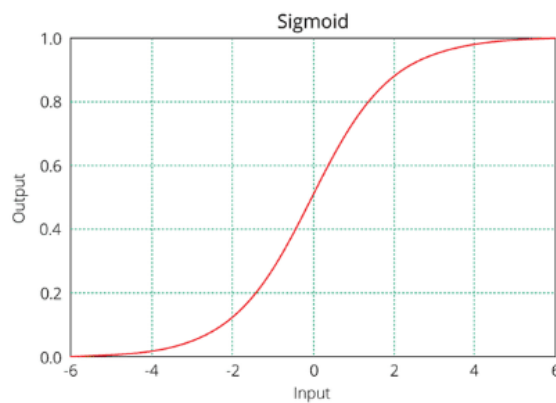
Sigmoid function

- **Odds** : Likelihood or probability of a favorable event happening
- **Log Odds** : Natural logarithm of the odds

## Linear Regression Expression

$$y_{\text{pred}} = w_0 * x_0 + w_1 * x_1 + w_2 * x_2 + w_3 * x_3 + w_4 * x_4 \dots w_n * x_n$$

1.  $e^{\ln \left[ \frac{p}{1-p} \right]} = e^{\beta_0 + \beta_1 \times \text{size}}$
2.  $\frac{p}{1-p} = e^{\beta_0 + \beta_1 \times \text{size}}$
3.  $p = e^{\beta_0 + \beta_1 \times \text{size}} - p \cdot e^{\beta_0 + \beta_1 \times \text{size}}$
4.  $p [1 + e^{\beta_0 + \beta_1 \times \text{size}}] = e^{\beta_0 + \beta_1 \times \text{size}}$
5.  $p = \frac{e^{\beta_0 + \beta_1 \times \text{size}}}{1 + e^{\beta_0 + \beta_1 \times \text{size}}}$
6.  $p = \frac{1}{1 + e^{-(\beta_0 + \beta_1 \times \text{size})}}$



Sigmoid Function

## Inducing the Predictions

$$\text{ClassLabel} = \{ 1 \quad g(y^\wedge) \geq 0.5, \{ 0 \quad g(y^\wedge) < 0.5$$

- Predict a class label of 1 if  $g(y^\wedge) \geq 0.5$ , and predict a class label of 0 if  $g(y^\wedge) < 0.5$ .
- We can also change the threshold of 0.5 to evaluate the model.

## Negative Log Loss or Cross Entropy Error:

$$\text{log loss} = -\frac{1}{n} \sum_{i=1}^n [y \cdot \log(\hat{p}) + (1-y) \cdot \log(1-\hat{p})]$$

*Handwritten notes:*  
- Above  $y \cdot \log(\hat{p})$ : this evaluates to 0 if  $y=0$   
- Above  $(1-y) \cdot \log(1-\hat{p})$ : this evaluates to 0 if  $y=1$

① sum of all the costs

$$\text{cost} = \frac{1}{n} \sum_{i=1}^n \begin{cases} -\log(\hat{p}) & \text{if } y = 1 \\ -\log(1-\hat{p}) & \text{if } y = 0 \end{cases}$$

② divide the sum by  $n$ , to get an average

If  $y = 1$  and we predicted 1 (i.e., 100% probability it sold), there is no penalty. However, if we predicted 0 (i.e., 0% probability it didn't sell), then we get penalized heavily.

- if  $y = 1$  :
- the house actually sold
  - cost =  $-\log(\hat{p})$



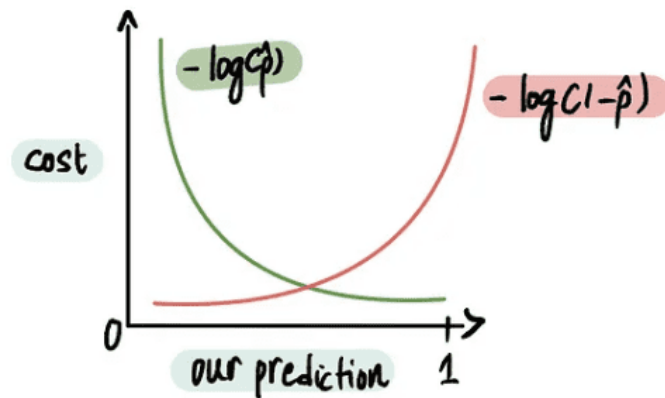
Similarly, if  $y = 0$  and we predicted a high probability of the house selling, we should be penalized heavily, and if we predicted a low probability of the house selling, there should be a lower penalty. The more off we are, the more it costs us.

- if  $y = 0$  :
- the house actually didn't sell
  - cost =  $-\log(1 - \hat{p})$



Combined graph of cost function:





## How to Determine if the Model Is Good or Bad - *Evaluation Metrics*

- The model predicts a class as positive and it was positive (TP).
- The model predicts a class as positive and it was negative (FP).
- The model predicts a class as negative and it was negative (TN).
- The model predicts a class as negative and it was positive (FN).

		Predicted Response		
		$\hat{y} = 1$	$\hat{y} = 0$	
True Response	$y = 1$	True Positive	False Negative	Recall (Sensitivity) $TP/(y=1)$
	$y = 0$	False Positive	True Negative	Specificity $TN/(y=0)$
Prevalence $(y=1)/total$		Precision $TP/(\hat{y} = 1)$	Accuracy $(TP+TN)/total$	

- Accuracy:** Measures the proportion of true results (both true positives and true negatives) among the total number of cases examined.
  - Formula:**  $Accuracy = TP + TN / FP + FN + TP + TN$
- Precision:** Also known as the Positive Predictive Value (PPV), measures the proportion of true positives among all the positive results predicted by the model.
  - Formula:**  $Precision = TP / FP + TP$
- Recall:** Also known as the True Positive Rate (TPR) or Sensitivity, measures the proportion of actual positives that are correctly identified as such.
  - Formula:**  $Recall = TP / FN + TP$
- Specificity:** Also known as the True Negative Rate (TNR), measures the proportion of actual negatives that are correctly identified as such.
  - Formula:**  $Specificity = TN / FP + TN$
- True Positive Rate (TPR),** also known as Sensitivity or Recall, measures the proportion of actual positives that are correctly identified by the model.
  - Formula:**  $TPR = TP / FN + TP$
- False Positive Rate (FPR)** measures the proportion of actual negatives that are incorrectly identified as positives by the model.

- **Formula:**  $FPR = FP / (TN + FP)$
- **True Negative Rate (TNR)**, also known as Specificity, measures the proportion of actual negatives that are correctly identified by the model.
  - **Formula:**  $TNR = TN / (FP + TN)$
- **False Negative Rate (FNR)** measures the proportion of actual positives that are incorrectly identified as negatives by the model.
  - **Formula:**  $FNR = TP / (FN + TP)$
- **F-Measure** F-measure is an attempt to combine both measures into a single value. Specifically, F-measure, also referred to as  $F_1$ , is computed as the weighted harmonic mean of precision and recall,
  - **F1 = 2 \* (precision \* recall) / (precision + recall)**

## Scenario:

Imagine we have a dataset of 1000 emails, and our model has classified them as follows:

- True Positives (TP): 400 emails correctly identified as spam.
- False Positives (FP): 100 emails incorrectly identified as spam (they are actually not spam).
- True Negatives (TN): 450 emails correctly identified as not spam.
- False Negatives (FN): 50 emails incorrectly identified as not spam (they are actually spam).

## Confusion Matrix:

	Predicted: Yes (Spam)	Predicted: No (Not Spam)
Actual: Yes (Spam)	TP = 400	FN = 50
Actual: No (Not Spam)	FP = 100	TN = 450

## Calculations:

### 1. Accuracy:

- $400+50+100+450 / 400+450 = 1000 / 850 = 0.85$  or 85%
- Interpretation: The model correctly classified 85% of the emails.

### 2. Precision (Positive Predictive Value):

- $400+100 / 400 = 500 / 400 = 0.8$  or 80%
- Interpretation: When the model predicts an email is spam, it is correct 80% of the time.

### 3. Recall (True Positive Rate or Sensitivity):

- $400+50 / 400 = 450 / 400 = 0.89$  or 89%
- Interpretation: The model correctly identifies 89% of all actual spam emails.

### 4. Specificity (True Negative Rate):

- $450+100 / 450 = 550 / 450 = 0.82$  or 82%
- Interpretation: The model correctly identifies 82% of all actual non-spam emails.

### 5. False Positive Rate:

- $100+450 / 100 = 550 / 100 = 0.18$  or 18%
- Interpretation: 18% of the non-spam emails were incorrectly classified as spam.

### 6. False Negative Rate:

- $400+50 / 50 = 450 / 50 = 0.11$  or 11%
- Interpretation: 11% of the spam emails were missed by the model and incorrectly classified as not spam.

### ● F1 Score

- Let's calculate the F1 Score:
- $F1 = 2 \times 0.8 \times 0.89 / (0.8 + 0.89)$
- Interpretation: 0.843. This indicates a strong balance between Precision and Recall, suggesting the model is relatively robust in terms of both correctly identifying spam emails and minimizing the number of non-spam emails incorrectly classified as spam

## Interpretations:

- **Accuracy** shows a general effectiveness of the model across both classes but doesn't distinguish between the types of errors.
- **Precision** indicates that when the model predicts an email as spam, it is fairly reliable, though there is still a 20% chance it could be wrong.
- **Recall** shows the model is quite good at catching spam emails, though it misses about 11% of them.
- **Specificity** tells us the model is also reasonably good at recognizing non-spam emails, but it mistakenly flags some as spam.
- **FPR** and **FNR** provide insight into the error rates for each class, indicating areas where the model's performance could be improved.

## AUC-ROC

An **ROC curve (receiver operating characteristic curve)** is a graph showing the performance of a classification model at all classification thresholds. This curve plots two parameters:

- True Positive Rate
- False Positive Rate

**True Positive Rate (TPR)** is a synonym for recall and is therefore defined as follows:

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

**False Positive Rate (FPR)** is defined as follows:

$$FPR = \frac{FP}{FP + TN}$$

