



## Week-4 : Logistic Regression

### Vehicle Cross-Sell Prediction Project Implementation

V Semester - ML and AIUP, Aug-Oct 2024

**Session Date and Time: 6<sup>th</sup> Oct 2024 10:30 AM IST – 12:00 Noon IST**

**Venkateswar Reddy Melachervu**

Visiting Faculty and [CTO, Brillium Technologies](#)

[Department of Computer Science Engineering – AI and ML \(CSM\)](#)

Email: [venkat.reddy.gf@gprec.ac.in](mailto:venkat.reddy.gf@gprec.ac.in)



**G. Pulla Reddy Engineering College (Autonomous)**

G. Pulla Reddy Nagar, Nandyal Road, Kurnool, AP 518007, India

Website: <https://www.gprec.ac.in>

# Disclaimer and Confidentiality Notice

The content of this guest lecture, including all discussions, materials, and demonstrations, code is intended for educational purposes only and reflects the views and opinions of the speaker. While every effort has been made to ensure the accuracy and relevance of the information presented, it should not be considered as legal, financial, or professional advice.

Brillium Technologies retains unrestricted ownership of all information shared during this session. Participants must not record, reproduce, distribute, or disclose any part of the lecture or materials without prior written permission from Brillium Technologies. Unauthorized use or distribution of the content may result in legal action.

Additionally, all trademarks, service marks, and intellectual properties referenced or used in this presentation are the property of their respective owners. No ownership or rights over such third-party content are claimed or implied by the author or Brillium Technologies or by GPREC.

By attending this lecture, you agree to respect the confidentiality of the information shared and refrain from using it in any unauthorized manner. Failure to comply with these terms may result in legal action.

Thank you for your understanding and cooperation.

# Lecture Outline



- Predicting Discrete or Categorical Values
- Discrete Value Predictions
- Logistic Regression Model and Logistic Function
- Prediction Probabilities
- Determining Regression Coefficients
- Model Evaluation Metrics
- RoC and Key Terms
- Vehicle Insurance Cross-Sell Prediction Project Overview
- Project Implementation Steps
- Key Implementation Advanced Concepts

# Predicting Discrete or Categorical Values

- Linear Regression
  - Predicting the house value in Boston housing area
  - Predicting the value of a used e2w etc.
  - Predicting the future price of an equity stock
  - Response or inferred variable  $y$  is **Continuous**

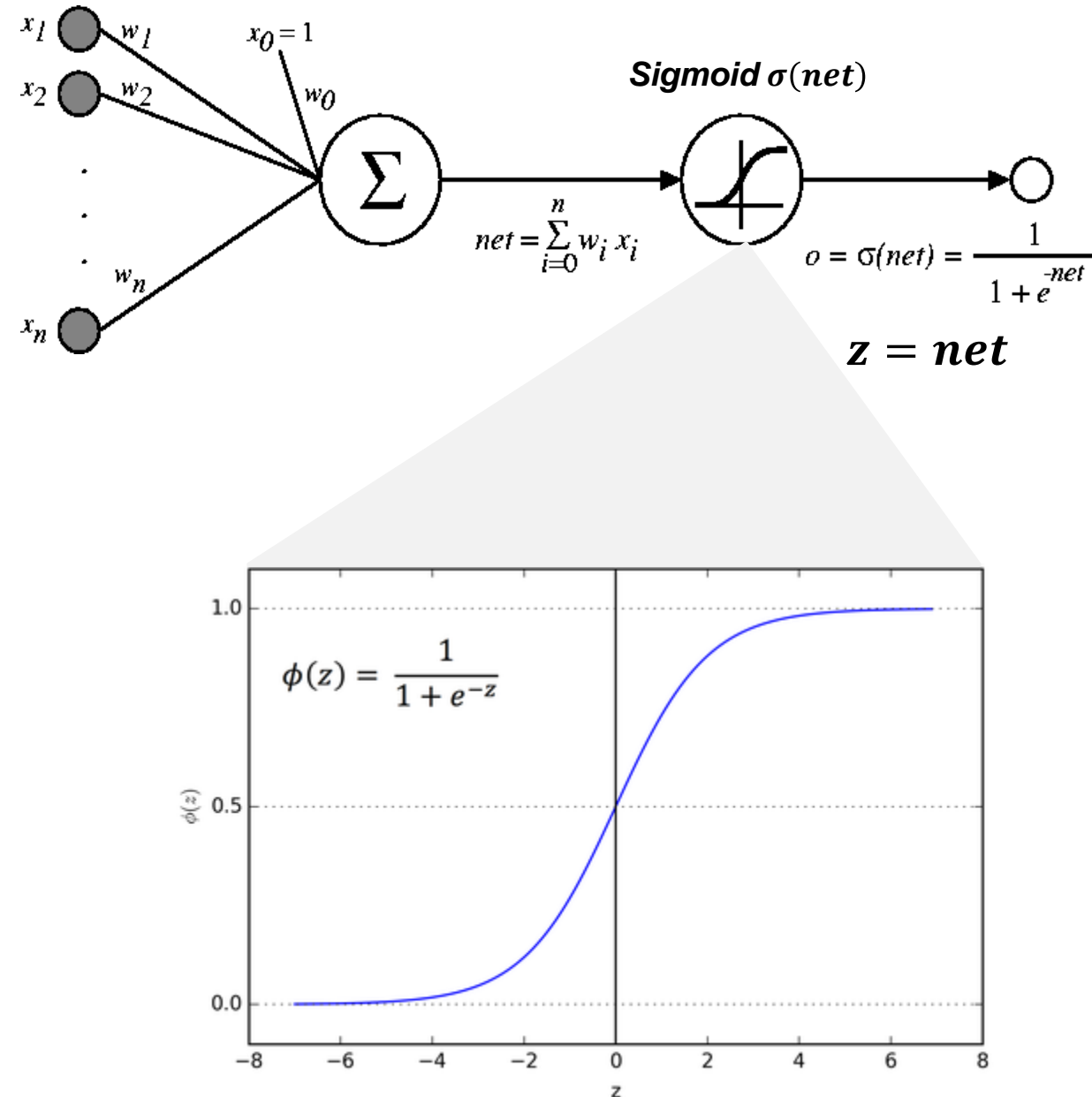
What if the response or inferred variable  $y$  is *discrete or categorical*?

# Discrete Value Predictions

- Two possible outcomes
  - Customer's interest in buying a vehicle insurance of a company -  $y \in \{0, 1\}$ 
    - Response - 1 (Interested in buying - 1)
    - Response - 0 (Not Interested in buying - 0)
  - Received email is a spam or genuine -  $y \in \{0, 1\}$
  - A credit card transaction is a fraudulent transaction or genuine transaction -  $y \in \{0, 1\}$
  - Predicting a person has a disease (diabetes, heart-related etc.) or not based on diagnostic reports and data
  - Yes/No, True/False, Success/Failure
- Multiple possible outcomes
  - Predicting a customer's preferred product category or brand (iPhone, Google Pixel, Samsung Galaxy, OnePlus etc.)
  - Image classification into cat, dog, bird, car, bus etc.
- Is Linear Regression suitable?
  - Predicted values are continuous and unbounded
    - Not between 0 and 1
    - Not bounded
  - Hence - Not suitable
- The outcomes  $\{0, 1\}$  are more like probabilities
- Need a model that maps real-valued numbers to a value in the range  $[0, 1]$  - Binary Outcomes

# Logistic Regression Model and Logistic Function

- A model that transforms the output of a linear equation (**regression**) into a probabilities (**logical or discrete or categorical**) using a special function
- The **Logistic Function** aka **Sigmoid Function** squashes linear equation output to a range between 0 and 1
  - $\sigma(z) = \frac{1}{(1+e^{-z})}$
  - Where
    - $z = w_0x_0 + w_1x_1 + w_2x_2 + w_nx_n$  - linear combination of input features
    - $\sigma(z)$  is the probability that the output prediction belongs to value/class 1
- $\sigma(z)$ :
  - 1 – The input instance belongs to class 1 (Positive class)
  - 0 – The input instance belongs to class 0 (Negative class)
  - Typical boundary condition for classification:
    - $\sigma(z) = \begin{cases} 1 & \text{if } \sigma(z) \geq 0.5 \\ 0 & \text{if } \sigma(z) < 0.5 \end{cases}$
- The Sigmoid function is:
  - Smooth
  - Continuous
  - Differentiable



# Logistic Regression Model and Logistic Function

- Two types based on the outcome (dependent variable)
  - Binary Logistic Regression – two possible outcomes/classes
  - Multinomial Logistic Regression – multiple possible outcomes/classes
- Estimates the discrete probability that a given input instance belongs to a particular class – 1 or 0
- **We will focus on Binary Logistic Algorithm in this programme**

**What is  $\sigma(z)$  value for below cases**

$$1. \lim_{z \rightarrow -\infty} \sigma(z) = \lim_{z \rightarrow -\infty} \frac{1}{(1+e^{-z})} = ?$$

$$2. \lim_{z \rightarrow \infty} \sigma(z) = \lim_{z \rightarrow \infty} \frac{1}{(1+e^{-z})} = ?$$

$$3. \sigma(z) = ? \text{ when } z = 0$$

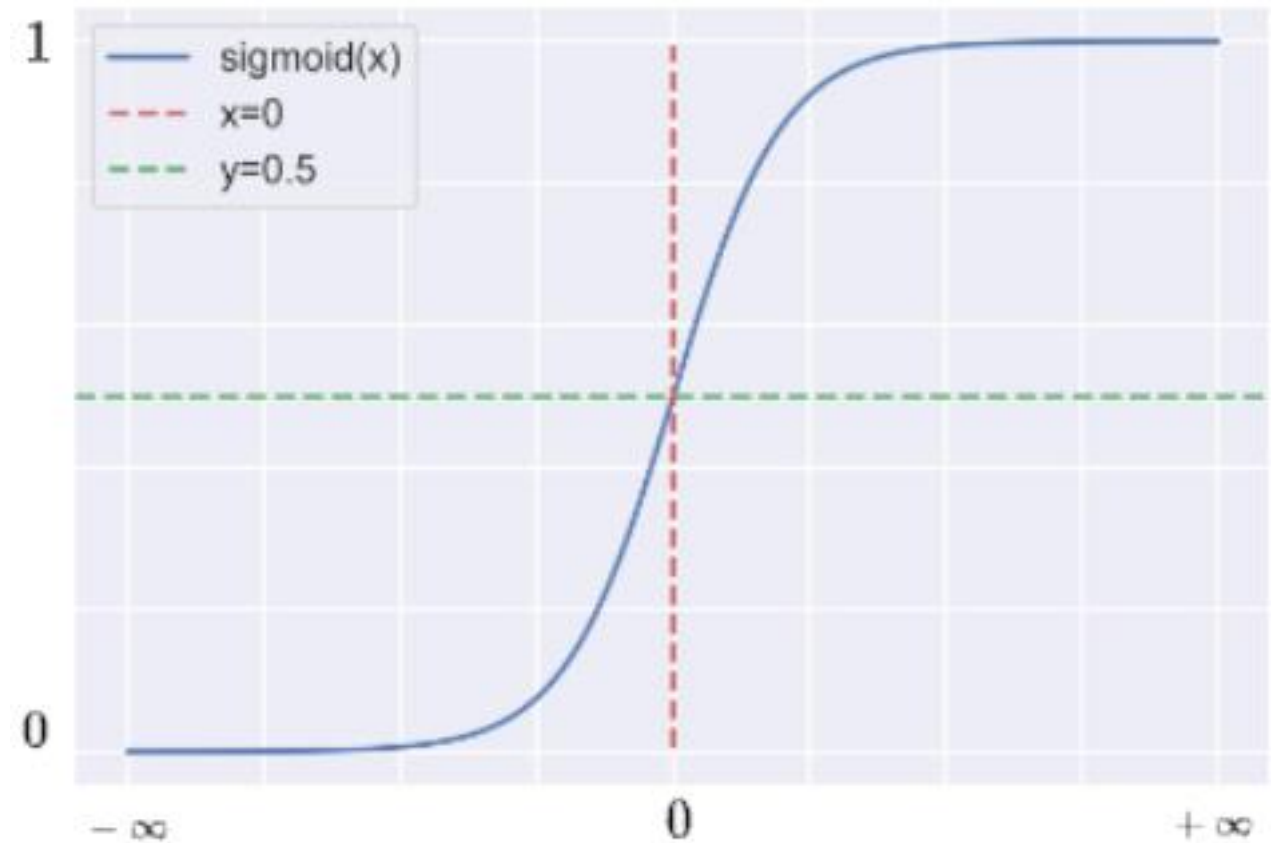
# Logistic Regression Model and Logistic Function

What is  $\sigma(z)$  value for below cases

1.  $\lim_{z \rightarrow -\infty} \sigma(z) = \lim_{z \rightarrow -\infty} \frac{1}{1+e^{-z}} = 0$

2.  $\lim_{z \rightarrow \infty} \sigma(z) = \lim_{z \rightarrow \infty} \frac{1}{1+e^{-z}} = 1$

3.  $\sigma(z) = 0.5$  when  $z = 0$





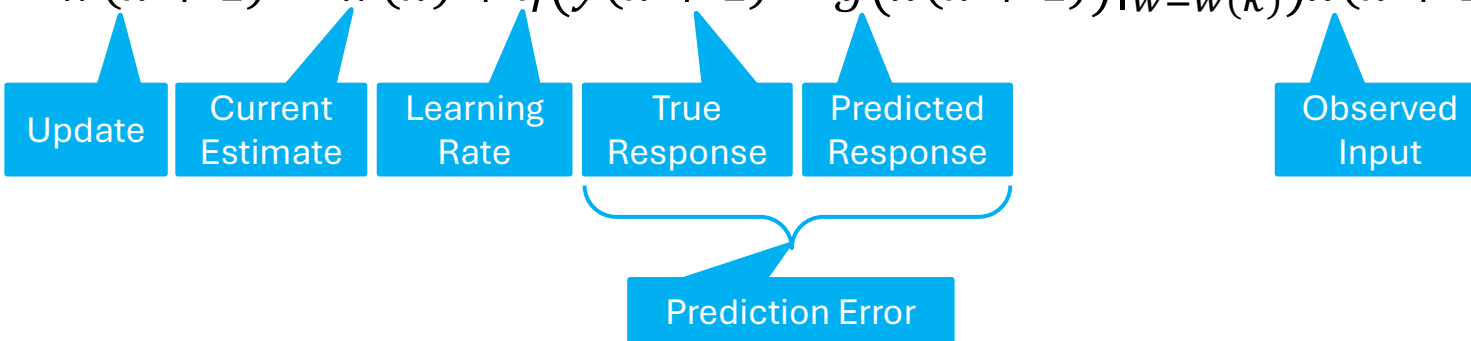
# Prediction Probabilities

- Input regression vector  $\bar{x} = \begin{bmatrix} 1 \\ x_1 \\ x_2 \\ x_3 \\ \vdots \\ x_n \end{bmatrix}$
- $P(y = 1 | \bar{x}) = \frac{1}{(1+e^{-\bar{x}^T \bar{w}})} = g(x)$  where  $\bar{x}^T \bar{w} = w_0x_0 + w_1x_1 + w_2x_2 + w_nx_n$ 
  - $P(y = 1 | \bar{x})$  – probability of response/prediction belongs to value/class 1, given input feature vector  $\bar{x}$
  - $y$  – Response/target value
  - $\bar{x}$  - Input feature vector or regression vector
  - $\bar{w}$  - Regression coefficients or weights vector
- $P(y = 0 | \bar{x}) = 1 - P(y = 1 | \bar{x}) = 1 - \frac{1}{(1+e^{-\bar{x}^T \bar{w}})} = \frac{e^{-\bar{x}^T \bar{w}}}{(1+e^{-\bar{x}^T \bar{w}})} = \mathbf{1} - g(x)$

# Determining Regression Coefficients

- Logistic regression uses **Maximum Likelihood Estimate – MLE** to determine the values of regression coefficients
  - Set of coefficients that maximize the likelihood of observing the actual data under the model
- The likelihood of  $(y(k), \bar{x}(k))$  ( $k^{th}$  observation)  $L(\bar{w}) = \left(g(\bar{x}(k))\right)^{y(k)} \left(1 - g(\bar{x}(k))\right)^{1-y(k)}$
- The joint likelihood of all responses and inputs  $L(\bar{w}) = \prod_{k=1}^M \left(g(\bar{x}(k))\right)^{y(k)} \left(1 - g(\bar{x}(k))\right)^{1-y(k)}$
- Joint log-likelihood  $L(\bar{w}) = \sum_{k=1}^M \left[ y(k) \ln \left(g(\bar{x}(k))\right) + (1 - y(k)) \ln \left(1 - g(\bar{x}(k))\right) \right]$
- To maximize the log-likelihood, **gradient ascent technique** may be employed
- The regression coefficients update rule, during the training, reduces to

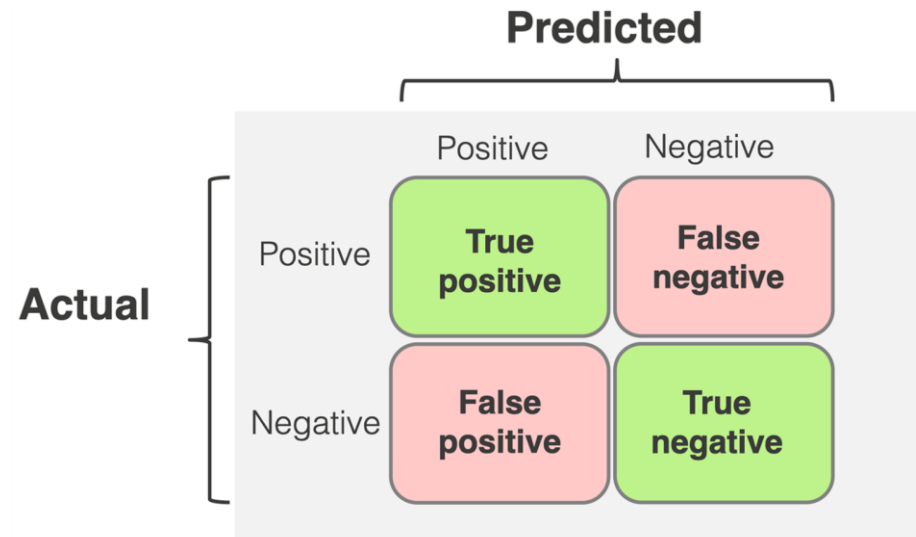
$$\bar{w}(k+1) = \bar{w}(k) + \eta(y(k+1) - g(\bar{x}(k+1))|_{\bar{w}=\bar{w}(k)})\bar{x}(k+1) \text{ where } g(\bar{x}(k+1)) = \frac{1}{1+e^{-(\bar{x}(k+1)^T \bar{w}(k))}}$$



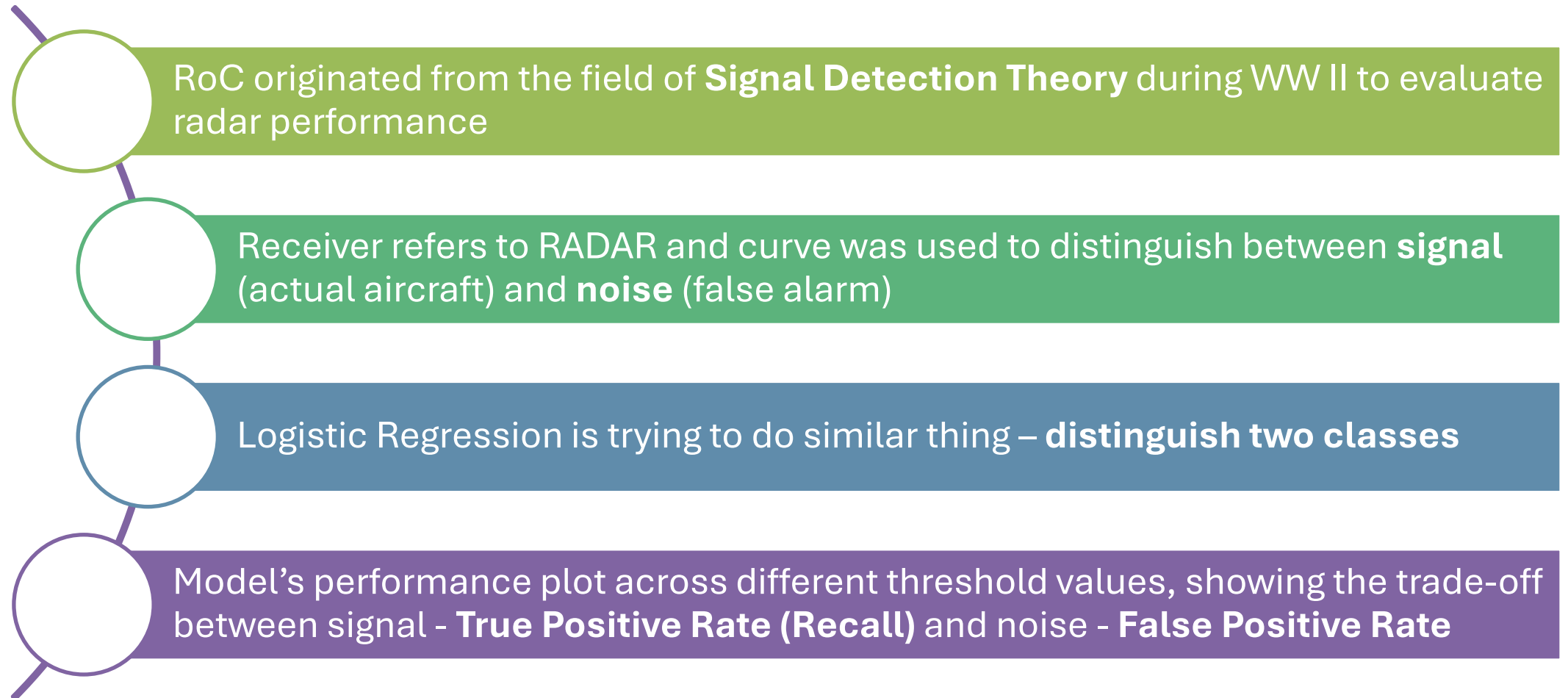
**Also Called Online Learning Algorithm**

# Model Evaluation Metrics

- Accuracy
  - Measures overall correctness
  - However, in imbalanced datasets, accuracy can be misleading.
- Confusion Matrix



# Model Evaluation Metrics - Receiver Operating Characteristic Curve



# Model Evaluation Metrics - RoC Key Terms

## True Positives

- The instances where the model correctly predicts the positive class as the positive class

## False Negatives

- The instances where the model incorrectly predicts the positive class as the negative class

## True Negatives

- The instances where the model correctly predicts the negative class as the negative class

## False Positives

- The instances where the model incorrectly predicts the negative class as the positive class

# Model Evaluation Metrics – TPR, FPR, Specificity, AUC

## True Positive Rate aka Recall aka Sensitivity

- Out of all actual positives, how many did the model correctly predict as positive -  $\frac{TP}{TP+FN}$

## False Positive Rate - FPR

- How many actual negatives are incorrectly identified as positives -  $\frac{FP}{FP+TN}$

## Specificity

- Out of all actual negatives, how many are correctly predict as negatives -  $\frac{TN}{TN+FP}$

## Area Under The Curve

- Single value summarizing the model's performance - A perfect model has an AUC of 1

# Model Evaluation Metrics – F1 Score

- Harmonic mean of precision and recall
  - Harmonic Mean
    - An average calculated by dividing the number of values by the sum of reciprocals of those values – useful to calculate the average of rates or ratios
    - $$HM = \frac{n}{\sum_{i=1}^n \frac{1}{x_i}}$$
- $$F1 = 2 \times \frac{Precision \times Recall}{Precision + Recall}$$
- Precision asks - "Out of the predicted positives, how many were actually positive?"
- Recall asks - "Out of the actual positives, how many did we correctly predict as positive?"
- Significance
  - If either precision or recall is very low, the harmonic mean (and thus the F1 Score) will also be low
  - This reflects that the model isn't performing well overall if one of these metrics is very weak

# Key Takeaways

- Logistic regression is simple but powerful for binary classification problems.
- Proper data preprocessing, handling class imbalance, and model evaluation are crucial for building an effective logistic regression model.
- Real-world project
  - Cross-sell vehicle insurance prediction, which illustrates how logistic regression can be applied to solve classification problems.



# Project Overview - Vehicle Insurance Cross-Sell Prediction

- Project Title
  - Cross-Sell Vehicle Insurance Prediction
- Object and Project Statement
  - A healthcare insurance company that has a sizeable existing customer for its health insurance product needs your help in building a ML to predict whether the existing health insurance customers/policyholders will also be interested in buying vehicle/automobile insurance provided by the company (2,3, and 4-wheelers) – such selling of another product to an existing customer is called cross-selling
  - Vehicle insurance is a mandatory insurance that needs to be bought by each vehicle owner and usually purchased/renewed annually.
  - Building an ML model to predict whether a customer would be interested in Vehicle Insurance is extremely helpful for the company because it can then accordingly plan its communication strategy to reach out to those customers and optimise its business model and increase revenue.

# Project Overview - Vehicle Insurance Cross-Sell Prediction

- Dataset

- The provided dataset - train and test CSVs – contain information about demographics (gender, age, region code type), vehicles (vehicle age, damages), healthcare policy (premium, sourcing channel) etc. The train CSV contains values for the response column whereas the values for this column are blank and need to be predicted and filled in the test CSV file.





Variable	Definition
id	Unique ID for the customer
Gender	Gender of the customer
Age	Age of the customer
Driving_License	0 : Customer does not have DL, 1 : Customer already has DL
Region_Code	Unique code for the region of the customer
Previously_Insured	1 : Customer already has Vehicle Insurance, 0 : Customer doesn't have Vehicle Insurance
Vehicle_Age	Age of the Vehicle
Vehicle_Damage	1 : Customer got his/her vehicle damaged in the past. 0 : Customer didn't get his/her vehicle damaged in the past.
Annual_Premium	The amount customer needs to pay as premium in the year
Policy_Sales_Channel	Anonymized Code for the channel of outreaching to the customer ie. Different Agents, Over Mail, Over Phone, In Person, etc.
Vintage	Number of Days, Customer has been associated with the company
Response	1 : Customer is interested, 0 : Customer is not interested

# Project Overview - Vehicle Insurance Cross-Sell Prediction

- Dataset

- The provided dataset - train and test CSVs – contain information about demographics (gender, age, region code type), vehicles (vehicle age, damages), healthcare policy (premium, sourcing channel) etc.
- **train.csv** contains values the data for training
- **test.csv** contains the data for prediction with response column being blank
- Training Data Set
  - Number of rows in train.csv : 3,81,109
  - Number of rows in test.csv : 1,27,037
  - Number of Features: 10
  - Response Variable (Target): Customer is interest or NOT interested (1 or 0)
- Features
  - Gender, Age, Driving\_License, Region\_Code, Previously\_Insured, Vehicle\_Age, Vehicle\_Damage
  - Annual\_Premium, Policy\_Sales\_Channel, Vintage

# Project Overview - Vehicle Insurance Cross-Sell Prediction

- Usage in Machine Learning:
  - The dataset is used to predict the binary response/interest of the customer in buying another product from the same company.
- Data Source and Published By:
  - The data is published by Anmol Kumar on Kaggle - [Health Insurance Cross Sell Prediction](https://www.kaggle.com/datasets/anmolkumar1993/health-insurance-cross-sell-prediction)   [\(kaggle.com\)](https://www.kaggle.com/datasets/anmolkumar1993/health-insurance-cross-sell-prediction)
- Data Download Link
  - Kaggle - [Health Insurance Cross Sell Prediction](https://www.kaggle.com/datasets/anmolkumar1993/health-insurance-cross-sell-prediction)   [\(kaggle.com\)](https://www.kaggle.com/datasets/anmolkumar1993/health-insurance-cross-sell-prediction)
- Submit the below
  - A new csv file “**AAM-IPL-insurance-prediction-submission.csv**” containing customer ID and predicted response for the interest in buying Vehicle Insurance, to be created for **test.csv** data points
  - PDF of the code, output, graphs etc. in a file titled “**AAM-IPL-Wk-4-LogiReg-Cross-Sell-Vehicle-Insurance-Prediction-Full-Code.pdf**”

# Project Implementation Steps

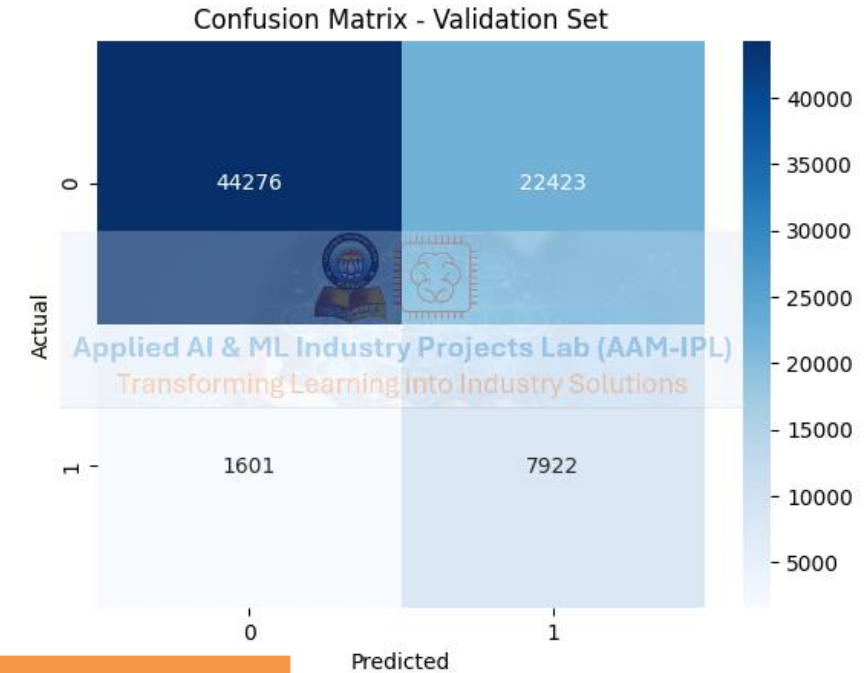
- Load the train.csv and test.csv datasets using pandas
- Print the number of rows in the training and test datasets
- Check for any missing values in the datasets
- Encode categorical features
- Standardise the data and train the model
  - SMOTE technique may be used optionally for handling data imbalance
- Optional hyper parameter tuning can be employed
- Evaluate the model using Metrics – Classification report, confusion matrix etc.
- Plot RoC and Precision Curves
- Plot Feature Importance
- Prediction on Test Data and Create Submission File

# Low Accuracy and Challenges

Validation Accuracy: 68.4815%

Classification Report (Validation Set):

	precision	recall	f1-score	support
0	0.97	0.66	0.79	66699
1	0.26	0.83	0.40	9523
accuracy			0.68	76222
macro avg	0.61	0.75	0.59	76222
weighted avg	0.88	0.68	0.74	76222



Target Variable	Count	Percentage
0	334399	87.8%
1	46400	12.2%

- Class Imbalance – **Bias Towards Majority Class**
- Feature Scaling – **Model Under Performance**
- Improper Regularization – **Non-generalized Model for Unseen Data**

# SMOTE – Synthetic Minority Over-sampling Technique

- In imbalanced datasets, models tend to favour the majority class, leading to poor performance on the minority class
- SMOTE helps by creating new synthetic data points for the minority class, resulting in a more balanced dataset
- Generates synthetic samples for the minority class (Response = 1), balancing the dataset
- Helps the model learn better from the minority class, improving **recall** and **F1-score** for the buyers
- Reduces false negatives by giving more importance to minority class predictions

# Overfitting and Underfitting

- Overfitting

- What is it?
  - The model captures noise and irrelevant details in the training data, leading to poor performance on unseen data
- Characteristic and Impact
  - High accuracy on training data, but low accuracy on test data
  - Model is too complex (e.g., too many parameters or deep trees)
  - Memorizes the training data rather than learning general patterns
- Solution
  - Use regularization, cross-validation, simplify the model, or increase training data

- Underfitting

- What is it?
  - The model is too simple to capture the underlying patterns in the data
- Characteristics
  - Low accuracy on both training and test data
  - Model fails to capture the relationship between features and target
  - Too simple (e.g., shallow trees or insufficient parameters)
- Solution
  - Use a more complex model, add more features, or reduce bias.



# Regularisation in ML Models

- Techniques that help prevent overfitting and improve a model's generalizability
- Adds a penalty to the model's loss function inhibiting the model from finding parameters that over-assign importance to its features
- This results in a small decrease in training accuracy, but a larger increase in generalizability
- Some common regularization techniques
  - L2 regularization
  - L1 regularization
  - Elastic Net etc.

# Feature Scaling

- Standardizes features so that they all have the same weight in the model, improving the convergence of logistic regression
- Ensures that no feature dominates the learning process, improving the model's ability to generalize
- Helps the model converge faster and with more stable predictions

# Hyperparameter Tuning

- Hyperparameters are model parameters that are set before training (e.g., C and solver in logistic regression)
- Tuning involves finding the best combination of hyperparameters to improve model performance
- GridSearchCV systematically tests different combinations of hyperparameters to find the best settings that maximize model performance, using cross-validation for reliability
- Parameters like C and solver can be employed for enhancing the model accuracy
- C – Inverse of regularisation strength
  - Regularization is a technique used to prevent overfitting by penalizing large or complex model coefficients
  - Smaller C values (e.g.,  $C = 0.01$ ) imply stronger regularization, which helps prevent overfitting by penalizing large coefficients
  - Larger C values (e.g.,  $C = 10$ ) reduce the regularization effect, allowing the model more flexibility to fit the data.
- Solver
  - Solver specifies the optimization algorithm used to fit the logistic regression model
  - Common solvers
    - liblinear: Suitable for small datasets and supports both L1 (Lasso) and L2 (Ridge) regularization
    - saga: Efficient for large datasets and also supports both L1 and L2 regularization. It can handle more complex problems like sparse data.

# Feature Importance

- Feature Importance measures the impact each input variable (feature) has on a model's predictions. It shows which features contribute the most to predicting the target variable
- In logistic regression, feature importance is based on the coefficients of the model
  - Positive coefficients: Increase the likelihood of the target class
  - Negative coefficients: Decrease the likelihood of the target class
  - Larger absolute values: Indicate stronger influence on the prediction
- Features like Annual Premium or Vehicle Age may have higher importance in predicting whether a customer will purchase insurance (Response = 1)







**Interested in building a Gen AI application?**

**Reach out to [venkat@brillium.in](mailto:venkat@brillium.in)**



THANK YOU!

**AAML-IPL Brought You in Partnership with:**



**Brillium Technologies**

Sector 7, HSR Layout, Bengaluru 560102, Karnataka, India

Website: [www.brillium.in](http://www.brillium.in) | Email: [connect@brillium.in](mailto:connect@brillium.in)