



ISDM (INDEPENDENT SKILL DEVELOPMENT MISSION)



REGRESSION MODELS: LINEAR REGRESSION & POLYNOMIAL REGRESSION

CHAPTER 1: INTRODUCTION TO REGRESSION MODELS

1.1 WHAT IS REGRESSION?

REGRESSION IS A STATISTICAL METHOD USED IN MACHINE LEARNING AND DATA SCIENCE TO MODEL RELATIONSHIPS BETWEEN VARIABLES. IT HELPS PREDICT A DEPENDENT VARIABLE (OUTPUT) BASED ON ONE OR MORE INDEPENDENT VARIABLES (INPUTS).

1.2 WHY IS REGRESSION IMPORTANT?

REGRESSION ANALYSIS IS WIDELY USED IN VARIOUS INDUSTRIES: ✓

FINANCE – PREDICTING STOCK PRICES AND FINANCIAL TRENDS.

✓ HEALTHCARE – DIAGNOSING DISEASES BASED ON SYMPTOMS AND TEST RESULTS.

✓ MARKETING – UNDERSTANDING CONSUMER BEHAVIOR AND SALES FORECASTING.

✓ ENGINEERING – ESTIMATING PRODUCT PERFORMANCE BASED ON INPUT FACTORS.

1.3 TYPES OF REGRESSION MODELS

- REGRESSION MODELS CAN BE CATEGORIZED AS:
- ◆ **LINEAR REGRESSION** – MODELS A STRAIGHT-LINE RELATIONSHIP.
 - ◆ **POLYNOMIAL REGRESSION** – CAPTURES NON-LINEAR RELATIONSHIPS.
 - ◆ **LOGISTIC REGRESSION** – USED FOR CLASSIFICATION PROBLEMS.
-

❖ CHAPTER 2: LINEAR REGRESSION

2.1 UNDERSTANDING LINEAR REGRESSION

LINEAR REGRESSION IS THE SIMPLEST FORM OF REGRESSION, WHERE THE RELATIONSHIP BETWEEN INPUT (X) AND OUTPUT (Y) IS MODELED USING A STRAIGHT LINE.

MATHEMATICAL FORMULA:

$$Y = MX + B$$

WHERE:

- YY = DEPENDENT VARIABLE (PREDICTED OUTPUT)
- XX = INDEPENDENT VARIABLE (INPUT)
- MM = SLOPE OF THE LINE
- BB = Y-INTERCEPT (WHERE THE LINE CROSSES THE Y-AXIS)

2.2 EXAMPLE OF LINEAR REGRESSION

A COMPANY WANTS TO PREDICT SALES BASED ON ADVERTISING SPENDING:

ADVERTISING SPEND (X)	SALES (Y)
10,000	20,000

20,000	40,000
30,000	60,000

A LINEAR REGRESSION MODEL WILL IDENTIFY A RELATIONSHIP BETWEEN ADVERTISING SPEND AND SALES, PREDICTING FUTURE SALES BASED ON NEW SPENDING VALUES.

2.3 ASSUMPTIONS OF LINEAR REGRESSION

- ✓ THERE IS A LINEAR RELATIONSHIP BETWEEN X AND Y.
- ✓ DATA HAS MINIMAL NOISE (OUTLIERS ARE RARE).
- ✓ VARIABLES ARE INDEPENDENT OF EACH OTHER.

2.4 IMPLEMENTING LINEAR REGRESSION IN PYTHON

```
IMPORT NUMPY AS NP  
IMPORT PANDAS AS PD  
IMPORT MATPLOTLIB.PYPLOT AS PLT  
FROM SKLEARN.LINEAR_MODEL IMPORT LINEARREGRESSION  
  
# SAMPLE DATA  
X = NP.ARRAY([10, 20, 30, 40, 50]).RESHAPE(-1,1)  
Y = NP.ARRAY([20, 40, 60, 80, 100])  
  
# CREATE AND TRAIN MODEL  
MODEL = LINEARREGRESSION()  
MODEL.FIT(X, Y)
```

```
# PREDICT AND PLOT
```

```
PREDICTIONS = MODEL.PREDICT(X)
```

```
PLT.SCATTER(X, Y, COLOR='BLUE')
```

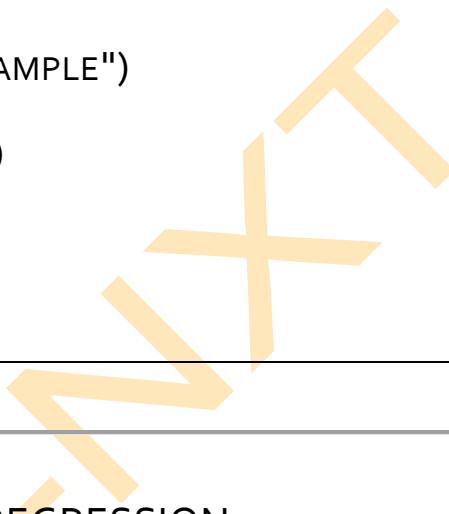
```
PLT.PLOT(X, PREDICTIONS, COLOR='RED')
```

```
PLT.TITLE("LINEAR REGRESSION EXAMPLE")
```

```
PLT.XLABEL("ADVERTISING SPEND")
```

```
PLT.YLABEL("SALES")
```

```
PLT.SHOW()
```



📌 CHAPTER 3: POLYNOMIAL REGRESSION

3.1 WHAT IS POLYNOMIAL REGRESSION?

POLYNOMIAL REGRESSION EXTENDS LINEAR REGRESSION BY MODELING A CURVED RELATIONSHIP BETWEEN VARIABLES. INSTEAD OF A STRAIGHT LINE, IT FITS A PARABOLIC OR HIGHER-ORDER CURVE TO THE DATA.

MATHEMATICAL FORMULA:

$$Y = AX^2 + BX + C$$

WHERE:

- YY = DEPENDENT VARIABLE (PREDICTED OUTPUT)
- XX = INDEPENDENT VARIABLE (INPUT)
- A, B, C = COEFFICIENTS

3.2 WHEN TO USE POLYNOMIAL REGRESSION?

- ✓ WHEN DATA HAS A **NON-LINEAR RELATIONSHIP**.
- ✓ WHEN A STRAIGHT LINE CANNOT ACCURATELY REPRESENT THE DATA TREND.
- ✓ WHEN HIGHER-DEGREE FEATURES IMPROVE MODEL ACCURACY.

3.3 EXAMPLE OF POLYNOMIAL REGRESSION

A COMPANY'S WEBSITE TRAFFIC IS INCREASING IN A NON-LINEAR PATTERN:

AD SPEND (X)	WEBSITE VISITS (Y)
10,000	5,000
20,000	15,000
30,000	35,000
40,000	70,000

A POLYNOMIAL REGRESSION MODEL CAN **CURVE-FIT** THIS DATA BETTER THAN LINEAR REGRESSION.

3.4 IMPLEMENTING POLYNOMIAL REGRESSION IN PYTHON

FROM SKLEARN.PREPROCESSING IMPORT POLYNOMIALFEATURES

FROM SKLEARN.PIPELINE IMPORT MAKE_PIPELINE

SAMPLE DATA

```
X = NP.ARRAY([10, 20, 30, 40, 50]).RESHAPE(-1,1)
```

```
Y = NP.ARRAY([5, 15, 35, 70, 120]) # NON-LINEAR GROWTH
```

```
# CREATE POLYNOMIAL REGRESSION MODEL
```

```
POLY_MODEL =  
MAKE_PIPELINE(POLYNOMIALFEATURES(DEGREE=2),  
LINEARREGRESSION())  
  
POLY_MODEL.FIT(X, Y)
```

```
# PREDICT AND PLOT
```

```
PREDICTIONS = POLY_MODEL.PREDICT(X)  
  
PLT.SCATTER(X, Y, COLOR='BLUE')  
  
PLT.PLOT(X, PREDICTIONS, COLOR='RED')  
  
PLT.TITLE("POLYNOMIAL REGRESSION EXAMPLE")  
  
PLT.XLABEL("ADVERTISING SPEND")  
  
PLT.YLABEL("WEBSITE VISITS")  
  
PLT.SHOW()
```

📌 CHAPTER 4: COMPARISON OF LINEAR & POLYNOMIAL REGRESSION

FEATURE	LINEAR REGRESSION	POLYNOMIAL REGRESSION
RELATIONSHIP	STRAIGHT-LINE (LINEAR)	CURVED (NON-LINEAR)
COMPLEXITY	SIMPLE	MORE COMPLEX

ACCURACY	LIMITED FOR NON-LINEAR DATA	FITS COMPLEX PATTERNS
FORMULA	$Y = MX + B$	$Y = AX^2 + BX + C$

4.1 WHEN TO CHOOSE WHICH MODEL?

- ✓ **LINEAR REGRESSION** – WHEN DATA SHOWS A SIMPLE INCREASING OR DECREASING TREND.
- ✓ **POLYNOMIAL REGRESSION** – WHEN DATA HAS CURVATURE OR NON-LINEARITY.

CHAPTER 5: EXERCISES & ASSIGNMENTS

5.1 MULTIPLE CHOICE QUESTIONS

1. WHAT DOES LINEAR REGRESSION PREDICT?

- (A) CATEGORIES
- (B) NUMERICAL VALUES
- (C) TEXT
- (D) IMAGES

2. WHICH FORMULA REPRESENTS A POLYNOMIAL REGRESSION?

- (A) $Y = MX + B$
- (B) $Y = AX^2 + BX + C$
- (C) $Y = X \times Z$
- (D) $Y = X / Z$

3. WHEN SHOULD YOU USE POLYNOMIAL REGRESSION?

- (A) WHEN DATA HAS A LINEAR PATTERN
- (B) WHEN DATA HAS A CURVED TREND

- (C) WHEN THERE IS NO PATTERN
- (D) ALWAYS

5.2 PRACTICAL ASSIGNMENTS

📌 **TASK 1:** IMPLEMENT LINEAR AND POLYNOMIAL REGRESSION ON A DATASET (E.G., HOUSING PRICES, SALES DATA).

📌 **TASK 2:** COMPARE THE ACCURACY OF BOTH MODELS AND EXPLAIN WHICH ONE WORKS BEST FOR YOUR DATA.

📌 **CHAPTER 6: SUMMARY**

- LINEAR REGRESSION MODELS A STRAIGHT-LINE RELATIONSHIP.**
- POLYNOMIAL REGRESSION CAPTURES CURVED RELATIONSHIPS.**
- POLYNOMIAL MODELS IMPROVE ACCURACY WHEN DATA IS NON-LINEAR.**
- BOTH MODELS HAVE DIFFERENT APPLICATIONS DEPENDING ON DATA TRENDS.**

🌟 **CONCLUSION: THE FUTURE OF REGRESSION MODELS**

AS MACHINE LEARNING ADVANCES, REGRESSION MODELS ARE BEING ENHANCED USING AI AND DEEP LEARNING. HYBRID MODELS COMBINING LINEAR AND NON-LINEAR TECHNIQUES ARE IMPROVING PREDICTIVE ACCURACY IN **FINANCE, HEALTHCARE, AND ENGINEERING.**

CLASSIFICATION MODELS: LOGISTIC REGRESSION, DECISION TREES, SUPPORT VECTOR

INTRODUCTION TO CLASSIFICATION MODELS

CLASSIFICATION IS A SUPERVISED LEARNING TECHNIQUE USED IN **MACHINE LEARNING (ML)** TO CATEGORIZE DATA INTO PREDEFINED CLASSES. IT IS WIDELY APPLIED IN **SPAM DETECTION, FRAUD DETECTION, MEDICAL DIAGNOSIS, SENTIMENT ANALYSIS, AND MORE.**

◆ TYPES OF CLASSIFICATION MODELS

THERE ARE SEVERAL CLASSIFICATION MODELS, BUT IN THIS GUIDE, WE WILL FOCUS ON:

1. LOGISTIC REGRESSION
2. DECISION TREES
3. SUPPORT VECTOR MACHINES (SVMs)

1. LOGISTIC REGRESSION

◆ WHAT IS LOGISTIC REGRESSION?

- LOGISTIC REGRESSION IS A STATISTICAL METHOD USED FOR BINARY CLASSIFICATION (YES/NO, 0/1, TRUE/FALSE).
- IT IS AN EXTENSION OF LINEAR REGRESSION BUT USES A SIGMOID FUNCTION TO PREDICT PROBABILITIES.

◆ THE SIGMOID FUNCTION

THE SIGMOID FUNCTION MAPS INPUT VALUES TO A PROBABILITY RANGE BETWEEN **0 AND 1**:

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

where z is the weighted sum of input features.

◆ WHEN TO USE LOGISTIC REGRESSION?

- WHEN THE TARGET VARIABLE IS **BINARY** (E.G., SPAM/NOT SPAM, FRAUD/NOT FRAUD).
- WHEN YOU NEED **INTERPRETABLE** PROBABILITY OUTPUTS.
- WORKS WELL WHEN FEATURES ARE **LINEARLY SEPARABLE**.

◆ ADVANTAGES

- ✓ SIMPLE AND EASY TO INTERPRET.
- ✓ WORKS WELL WITH SMALL DATASETS.
- ✓ OUTPUTS PROBABILITIES, WHICH HELP IN DECISION-MAKING.

◆ DISADVANTAGES

✗ ASSUMES LINEAR RELATIONSHIP BETWEEN INDEPENDENT VARIABLES AND THE LOG ODDS.

✗ NOT EFFECTIVE WHEN DEALING WITH **NON-LINEAR** RELATIONSHIPS.

◆ EXAMPLE: PREDICTING CUSTOMER CHURN

```
IMPORT PANDAS AS PD
```

```
FROM SKLEARN.MODEL_SELECTION IMPORT TRAIN_TEST_SPLIT  
FROM SKLEARN.LINEAR_MODEL IMPORT LOGISTICREGRESSION  
FROM SKLEARN.METRICS IMPORT ACCURACY_SCORE  
  
# LOAD DATASET  
DF = PD.READ_CSV("CUSTOMER_CHURN.CSV")  
  
# SELECT FEATURES AND TARGET VARIABLE  
X = DF[['MONTHLY_CHARGES', 'TENURE']]  
Y = DF['CHURN'] # 0: NO, 1: YES  
  
# SPLIT DATASET INTO TRAINING AND TESTING SETS  
X_TRAIN, X_TEST, Y_TRAIN, Y_TEST = TRAIN_TEST_SPLIT(X, Y,  
TEST_SIZE=0.2, RANDOM_STATE=42)  
  
# TRAIN LOGISTIC REGRESSION MODEL  
MODEL = LOGISTICREGRESSION()  
MODEL.FIT(X_TRAIN, Y_TRAIN)  
  
# PREDICTIONS  
Y_PRED = MODEL.PREDICT(X_TEST)
```

```
# ACCURACY SCORE
```

```
PRINT("ACCURACY:", ACCURACY_SCORE(Y_TEST, Y_PRED))
```

📌 2. DECISION TREES

◆ WHAT IS A DECISION TREE?

- A TREE-LIKE MODEL WHERE DATA IS SPLIT INTO BRANCHES BASED ON CONDITIONS.
- USES A HIERARCHICAL DECISION-MAKING PROCESS.

◆ HOW IT WORKS?

1. STARTS FROM THE ROOT NODE (ENTIRE DATASET).
2. SPLITS INTO BRANCHES USING CONDITIONS (E.G., IS AGE > 30?).
3. ENDS AT LEAF NODES, WHICH CONTAIN THE FINAL CLASS LABELS.

◆ KEY CONCEPTS

- ENTROPY: MEASURES IMPURITY IN A DATASET.
- GINI INDEX: MEASURES HOW WELL A NODE IS SPLIT.
- PRUNING: REMOVING UNNECESSARY BRANCHES TO PREVENT OVERRFITTING.

◆ WHEN TO USE DECISION TREES?

- WHEN THE DATASET HAS COMPLEX AND NON-LINEAR RELATIONSHIPS.
- IF INTERPRETABILITY IS IMPORTANT (EASY TO VISUALIZE).

◆ ADVANTAGES

- ✓ WORKS FOR BOTH CLASSIFICATION AND REGRESSION PROBLEMS.
- ✓ HANDLES NON-LINEAR RELATIONSHIPS EFFECTIVELY.
- ✓ NO NEED FOR FEATURE SCALING.

◆ DISADVANTAGES

- ✗ OVERRFITTING IF THE TREE IS TOO DEEP.
- ✗ SENSITIVE TO NOISY DATA, LEADING TO INSTABILITY.

◆ EXAMPLE: PREDICTING LOAN APPROVAL

```
FROM SKLEARN.TREE IMPORT DECISIONTREECLASSIFIER
```

```
# TRAIN DECISION TREE CLASSIFIER
```

```
MODEL = DECISIONTREECLASSIFIER(MAX_DEPTH=3,  
RANDOM_STATE=42)
```

```
MODEL.FIT(X_TRAIN, Y_TRAIN)
```

```
# PREDICTIONS
```

```
Y_PRED = MODEL.PREDICT(X_TEST)
```

```
# ACCURACY SCORE
```

```
PRINT("ACCURACY:", ACCURACY_SCORE(Y_TEST, Y_PRED))
```

◆ VISUALIZING THE DECISION TREE

```
FROM SKLEARN IMPORT TREE
```

```
IMPORT MATPLOTLIB.PYPLOT AS PLT
```

```
PLT.FIGURE(FIGSIZE=(12,6))
```

```
TREE.PLOT_TREE(MODEL, FILLED=TRUE,  
FEATURE_NAMES=['MONTHLY_CHARGES', 'TENURE'],  
CLASS_NAMES=['No', 'Yes'])
```

```
PLT.SHOW()
```

📌 3. SUPPORT VECTOR MACHINES (SVMs)

◆ WHAT IS SVM?

- A POWERFUL CLASSIFICATION ALGORITHM THAT SEPARATES DATA USING A HYPERPLANE.
- MAXIMIZES THE MARGIN BETWEEN DIFFERENT CLASSES.
- CAN BE LINEAR OR NON-LINEAR.

◆ KEY CONCEPTS

- HYPERPLANE: A DECISION BOUNDARY THAT SEPARATES DIFFERENT CLASSES.
- SUPPORT VECTORS: DATA POINTS CLOSEST TO THE HYPERPLANE THAT INFLUENCE ITS POSITION.
- KERNEL TRICK: USED WHEN DATA IS NOT LINEARLY SEPARABLE (E.G., RBF KERNEL).

◆ WHEN TO USE SVM?

- WHEN THE DATASET IS HIGH-DIMENSIONAL.

- WORKS WELL WHEN DATA IS **NOT LINEARLY SEPARABLE**.
- EFFECTIVE FOR **SMALL TO MEDIUM-SIZED DATASETS**.

◆ **ADVANTAGES**

- ✓ WORKS WELL WITH **HIGH-DIMENSIONAL DATA**.
- ✓ CAN CLASSIFY **NON-LINEARLY SEPARABLE DATA** (WITH KERNEL TRICK).
- ✓ ROBUST TO OUTLIERS.

◆ **DISADVANTAGES**

- ✗ COMPUTATIONALLY EXPENSIVE FOR **LARGE DATASETS**.
- ✗ CHOOSING THE RIGHT KERNEL REQUIRES TUNING.

◆ **EXAMPLE: CLASSIFYING TUMORS (MALIGNANT OR BENIGN)**

```
FROM SKLEARN.SVM IMPORT SVC  
  
# TRAIN AN SVM CLASSIFIER  
  
SVM_MODEL = SVC(KERNEL='RBF', C=1.0, GAMMA='SCALE')  
  
SVM_MODEL.FIT(X_TRAIN, Y_TRAIN)  
  
# PREDICTIONS  
  
Y_PRED = SVM_MODEL.PREDICT(X_TEST)  
  
# ACCURACY SCORE  
  
PRINT("ACCURACY:", ACCURACY_SCORE(Y_TEST, Y_PRED))
```

◆ VISUALIZING THE SVM DECISION BOUNDARY

```
IMPORT NUMPY AS NP
```

```
IMPORT MATPLOTLIB.PY PLOT AS PLT
```

```
# CREATE MESHGRID
```

```
X_MIN, X_MAX = X_TRAIN['MONTHLY_CHARGES'].MIN() - 1,  
X_TRAIN['MONTHLY_CHARGES'].MAX() + 1
```

```
Y_MIN, Y_MAX = X_TRAIN['TENURE'].MIN() - 1,  
X_TRAIN['TENURE'].MAX() + 1
```

```
XX, YY = NP.MESHGRID(NP.LINSPACE(X_MIN, X_MAX, 100),  
NP.LINSPACE(Y_MIN, Y_MAX, 100))
```

```
# PREDICT ON THE GRID
```

```
Z = SVM_MODEL.PREDICT(NP.C_[XX.RAVEL(), YY.RAVEL()])
```

```
Z = Z.RESHAPE(XX.SHAPe)
```

```
# PLOT DECISION BOUNDARY
```

```
PLT.CONTOURF(XX, YY, Z, ALPHA=0.3)
```

```
PLT.SCATTER(X_TRAIN['MONTHLY_CHARGES'], X_TRAIN['TENURE'],  
C=Y_TRAIN, EDGECOLORS='K')
```

```
PLT.TITLE('SVM DECISION BOUNDARY')
```

```
PLT.XLABEL('MONTHLY CHARGES')
```

```
PLT.YLABEL('TENURE')
```

PLT.SHOW()

SUMMARY

ALGORITHM	BEST FOR	ADVANTAGES	DISADVANTAGES
LOGISTIC REGRESSION	BINARY CLASSIFICATION, INTERPRETABLE RESULTS	SIMPLE, PROBABILITY-BASED, EFFICIENT	ONLY WORKS WITH LINEAR DATA
DECISION TREES	COMPLEX DATASETS, EASY VISUALIZATION	HANDLES NON-LINEARITY, INTERPRETABLE	PRONE TO OVERRFITTING
SVM	HIGH-DIMENSIONAL AND NON-LINEAR DATA	WORKS WELL FOR COMPLEX BOUNDARIES, ROBUST	COMPUTATIONALLY EXPENSIVE FOR LARGE DATASETS

CONCLUSION

- **LOGISTIC REGRESSION** IS GREAT FOR SIMPLE, LINEARLY SEPARABLE PROBLEMS.
- **DECISION TREES** ARE USEFUL WHEN INTERPRETABILITY IS IMPORTANT.

- **SVMs** WORK WELL FOR COMPLEX DECISION BOUNDARIES BUT CAN BE SLOW.

ISDM-NxT

MODEL EVALUATION METRICS: ACCURACY, PRECISION, RECALL, F₁-SCORE

CHAPTER 1: INTRODUCTION TO MODEL EVALUATION METRICS

1.1 WHAT IS MODEL EVALUATION?

MODEL EVALUATION IS THE PROCESS OF ASSESSING HOW WELL A MACHINE LEARNING MODEL PERFORMS ON A DATASET. IT HELPS DETERMINE WHETHER A MODEL IS **ACCURATE, RELIABLE, AND SUITABLE** FOR REAL-WORLD APPLICATIONS.

1.2 WHY IS MODEL EVALUATION IMPORTANT?

- ✓ HELPS AVOID **OVERTFITTING** (WHEN A MODEL MEMORIZES TRAINING DATA BUT FAILS ON NEW DATA).
- ✓ ENSURES THE MODEL **GENERALIZES WELL** TO UNSEEN DATA.
- ✓ HELPS IN SELECTING THE **BEST MODEL** AMONG MULTIPLE MODELS.

1.3 COMMON MODEL EVALUATION METRICS

- ◆ **ACCURACY** – MEASURES THE PROPORTION OF CORRECT PREDICTIONS.
- ◆ **PRECISION** – EVALUATES HOW MANY PREDICTED POSITIVES WERE ACTUALLY CORRECT.
- ◆ **RECALL** – MEASURES HOW WELL THE MODEL DETECTS ACTUAL POSITIVES.
- ◆ **F₁-SCORE** – A BALANCE BETWEEN PRECISION AND RECALL.

CHAPTER 2: ACCURACY METRIC

2.1 WHAT IS ACCURACY?

ACCURACY IS THE MOST BASIC EVALUATION METRIC, MEASURING THE PERCENTAGE OF CORRECTLY CLASSIFIED INSTANCES.

FORMULA FOR ACCURACY:

$$\text{Accuracy} = \frac{\text{Correct Predictions}}{\text{Total Predictions}} \times 100$$

2.2 EXAMPLE OF ACCURACY CALCULATION

CONSIDER A SPAM EMAIL CLASSIFIER PREDICTING WHETHER AN EMAIL IS SPAM OR NOT.

ACTUAL	PREDICTED
SPAM	SPAM
NOT SPAM	NOT SPAM
SPAM	NOT SPAM
NOT SPAM	NOT SPAM

- ◆ **CORRECT PREDICTIONS = 3**
- ◆ **TOTAL PREDICTIONS = 4**

$$\text{Accuracy} = \frac{3}{4} \times 100 = 75\%$$

2.3 WHEN IS ACCURACY USEFUL?

- ✓ WHEN THE DATASET IS BALANCED (EQUAL POSITIVES AND NEGATIVES).

- ✓ IN SIMPLE CLASSIFICATION PROBLEMS WHERE ALL ERRORS HAVE EQUAL CONSEQUENCES.

2.4 WHEN IS ACCURACY MISLEADING?

- 🚫 IF THE DATASET IS IMBALANCED (E.G., DETECTING FRAUD, WHERE FRAUD CASES ARE RARE), ACCURACY MAY BE MISLEADING.
- 🚫 IF PREDICTING ONLY THE MAJORITY CLASS (E.G., ALWAYS PREDICTING "NOT SPAM"), ACCURACY MIGHT BE HIGH BUT THE MODEL IS USELESS.

CHAPTER 3: PRECISION METRIC

3.1 WHAT IS PRECISION?

PRECISION MEASURES THE PROPORTION OF CORRECT POSITIVE PREDICTIONS OUT OF ALL POSITIVE PREDICTIONS MADE BY THE MODEL.

FORMULA FOR PRECISION:

$$\text{Precision} = \frac{\text{True Positives (TP)}}{\text{True Positives (TP)} + \text{False Positives (FP)}}$$

3.2 EXAMPLE OF PRECISION CALCULATION

CONSIDER A MEDICAL TEST FOR DETECTING A DISEASE.

ACTUAL	PREDICTED
DISEASED	DISEASED
HEALTHY	DISEASED
DISEASED	DISEASED

HEALTHY	HEALTHY
---------	---------

- ◆ **TRUE POSITIVES (TP)** = 2 (CORRECTLY DETECTED DISEASED CASES)
- ◆ **FALSE POSITIVES (FP)** = 1 (HEALTHY PERSON MISCLASSIFIED AS DISEASED)

$\text{PRECISION} = \frac{\text{TP}}{\text{TP} + \text{FP}} = \frac{2}{2+1} = \frac{2}{3} = 0.66 \text{ (66\%)}$

3.3 WHEN IS PRECISION IMPORTANT?

- ✓ WHEN FALSE POSITIVES ARE COSTLY (E.G., FRAUD DETECTION, CANCER DIAGNOSIS).
- ✓ IN APPLICATIONS WHERE INCORRECT POSITIVE PREDICTIONS SHOULD BE MINIMIZED (E.G., RECOMMENDING INAPPROPRIATE PRODUCTS TO USERS).

CHAPTER 4: RECALL METRIC

4.1 WHAT IS RECALL?

RECALL (ALSO CALLED SENSITIVITY OR TRUE POSITIVE RATE) MEASURES THE PROPORTION OF ACTUAL POSITIVE CASES THAT WERE CORRECTLY PREDICTED.

FORMULA FOR RECALL:

$$\text{Recall} = \frac{\text{True Positives (TP)}}{\text{True Positives (TP)} + \text{False Negatives (FN)}}$$

4.2 EXAMPLE OF RECALL CALCULATION

USING THE SAME MEDICAL TEST EXAMPLE:

ACTUAL	PREDICTED
DISEASED	DISEASED
HEALTHY	DISEASED
DISEASED	DISEASED
DISEASED	HEALTHY

- ◆ **TRUE POSITIVES (TP)** = 2 (CORRECTLY DETECTED DISEASED CASES)
- ◆ **FALSE NEGATIVES (FN)** = 1 (MISSED DISEASED CASE)

$\text{RECALL} = \frac{\text{TP}}{\text{TP} + \text{FN}} = \frac{2}{2+1} = 0.66 \text{ (66\%)}$

4.3 WHEN IS RECALL IMPORTANT?

- ✓ WHEN FALSE NEGATIVES ARE COSTLY (E.G., DETECTING CANCER, SAFETY ALARMS).
- ✓ IN SCENARIOS WHERE MISSING A POSITIVE CASE HAS SERIOUS CONSEQUENCES (E.G., DETECTING TERRORIST THREATS).

CHAPTER 5: F1-SCORE METRIC

5.1 WHAT IS F1-SCORE?

F1-SCORE IS THE HARMONIC MEAN OF PRECISION AND RECALL. IT BALANCES THE TWO WHEN ONE IS HIGH AND THE OTHER IS LOW.

FORMULA FOR F1-SCORE:

$$F1 = 2 \times \frac{\text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}}$$

5.2 EXAMPLE OF F1-SCORE CALCULATION

USING OUR PREVIOUS PRECISION (66%) AND RECALL (66%):

$$F_1 = 2 \times \frac{\text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}} = 2 \times \frac{0.66 \times 0.66}{0.66 + 0.66} = 0.66 \text{ (66%)}$$

5.3 WHEN IS F_1 -SCORE IMPORTANT?

- ✓ WHEN BOTH FALSE POSITIVES AND FALSE NEGATIVES MATTER (E.G., SPAM DETECTION).
- ✓ IN IMBALANCED DATASETS, WHERE ACCURACY IS MISLEADING.

📌 CHAPTER 6: COMPARISON OF METRICS

METRIC	DEFINITION	WHEN TO USE?
ACCURACY	% OF CORRECT PREDICTIONS	BALANCED DATASETS, GENERAL CLASSIFICATION
PRECISION	% OF PREDICTED POSITIVES THAT ARE CORRECT	WHEN FALSE POSITIVES ARE COSTLY (E.G., FRAUD DETECTION)
RECALL	% OF ACTUAL POSITIVES CORRECTLY IDENTIFIED	WHEN FALSE NEGATIVES ARE COSTLY (E.G., MEDICAL TESTS)
F_1 -SCORE	BALANCE BETWEEN PRECISION AND RECALL	WHEN BOTH FP AND FN ARE IMPORTANT

📌 CHAPTER 7: IMPLEMENTING METRICS IN PYTHON

7.1 USING SCIKIT-LEARN

```
FROM SKLEARN.METRICS IMPORT ACCURACY_SCORE,  
PRECISION_SCORE, RECALL_SCORE, F1_SCORE  
  
# EXAMPLE PREDICTIONS  
  
Y_TRUE = [1, 0, 1, 1, 0, 1, 0, 0, 1, 1] # ACTUAL LABELS  
Y_PRED = [1, 0, 1, 0, 0, 1, 1, 0, 1, 1] # PREDICTED LABELS  
  
# CALCULATE METRICS  
  
ACCURACY = ACCURACY_SCORE(Y_TRUE, Y_PRED)  
PRECISION = PRECISION_SCORE(Y_TRUE, Y_PRED)  
RECALL = RECALL_SCORE(Y_TRUE, Y_PRED)  
F1 = F1_SCORE(Y_TRUE, Y_PRED)  
  
# PRINT RESULTS  
  
PRINT(F"ACCURACY: {ACCURACY:.2F}")  
PRINT(F"PRECISION: {PRECISION:.2F}")  
PRINT(F"RECALL: {RECALL:.2F}")  
PRINT(F"F1-SCORE: {F1:.2F}")
```



CHAPTER 8: EXERCISES & ASSIGNMENTS

8.1 MULTIPLE CHOICE QUESTIONS

1. WHAT DOES PRECISION MEASURE?

- (A) CORRECTLY CLASSIFIED INSTANCES
- (B) CORRECT POSITIVE PREDICTIONS
- (C) MISCLASSIFIED INSTANCES

2. WHICH METRIC IS BEST WHEN MISSING A POSITIVE CASE IS COSTLY?

- (A) ACCURACY
- (B) RECALL
- (C) PRECISION

3. WHAT DOES F₁-SCORE BALANCE?

- (A) ACCURACY AND PRECISION
- (B) PRECISION AND RECALL
- (C) FALSE POSITIVES AND TRUE NEGATIVES

8.2 PRACTICAL ASSIGNMENTS

- 📌 **TASK 1:** IMPLEMENT MODEL EVALUATION METRICS ON A DATASET (E.G., SPAM DETECTION).
- 📌 **TASK 2:** COMPARE MODEL PERFORMANCE USING ACCURACY, PRECISION, RECALL, AND F₁-SCORE.

📌 CHAPTER 9: SUMMARY

- ✓ ACCURACY IS USEFUL FOR BALANCED DATASETS.
- ✓ PRECISION IS IMPORTANT WHEN FALSE POSITIVES ARE COSTLY.
- ✓ RECALL IS NEEDED WHEN FALSE NEGATIVES ARE COSTLY.
- ✓ F₁-SCORE BALANCES PRECISION AND RECALL FOR IMBALANCED DATA.



HYPERPARAMETER TUNING & MODEL OPTIMIZATION

📌 CHAPTER 1: INTRODUCTION TO HYPERPARAMETER TUNING & MODEL OPTIMIZATION

1.1 WHAT IS HYPERPARAMETER TUNING?

HYPERPARAMETER TUNING REFERS TO THE PROCESS OF SELECTING THE BEST HYPERPARAMETERS FOR A MACHINE LEARNING MODEL TO IMPROVE ITS PERFORMANCE. UNLIKE MODEL PARAMETERS (WHICH ARE LEARNED FROM DATA), HYPERPARAMETERS ARE SET BEFORE TRAINING BEGINS.

◆ **EXAMPLE:** IN A DECISION TREE, HYPERPARAMETERS INCLUDE:

- MAXIMUM DEPTH OF THE TREE
- MINIMUM SAMPLES PER LEAF
- SPLITTING CRITERION (GINI IMPURITY VS. ENTROPY)

1.2 WHY IS HYPERPARAMETER TUNING IMPORTANT?

PROPER HYPERPARAMETER TUNING CAN: ✓ IMPROVE MODEL ACCURACY AND GENERALIZATION

- ✓ PREVENT OVERFITTING OR UNDERFITTING
- ✓ OPTIMIZE TRAINING TIME AND COMPUTATIONAL RESOURCES

1.3 DIFFERENCE BETWEEN PARAMETERS & HYPERPARAMETERS

FEATURE	MODEL PARAMETERS	HYPERPARAMETERS
DEFINITION	VALUES LEARNED BY THE MODEL FROM DATA	VALUES SET BEFORE TRAINING
EXAMPLES	WEIGHTS IN NEURAL NETWORKS, COEFFICIENTS IN LINEAR REGRESSION	LEARNING RATE, NUMBER OF LAYERS, NUMBER OF NEURONS
ADJUSTED BY	TRAINING ALGORITHM	MANUALLY OR VIA TUNING TECHNIQUES

📌 CHAPTER 2: COMMON HYPERPARAMETERS IN MACHINE LEARNING

DIFFERENT ML ALGORITHMS HAVE DIFFERENT HYPERPARAMETERS.
HERE ARE SOME COMMON ONES:

2.1 HYPERPARAMETERS IN DECISION TREES

- ✓ MAX_DEPTH – LIMITS TREE DEPTH TO AVOID OVERFITTING
- ✓ MIN_SAMPLES_SPLIT – MINIMUM SAMPLES REQUIRED TO SPLIT A NODE
- ✓ CRITERION – MEASURES SPLIT QUALITY (GINI IMPURITY VS. ENTROPY)

2.2 HYPERPARAMETERS IN RANDOM FOREST

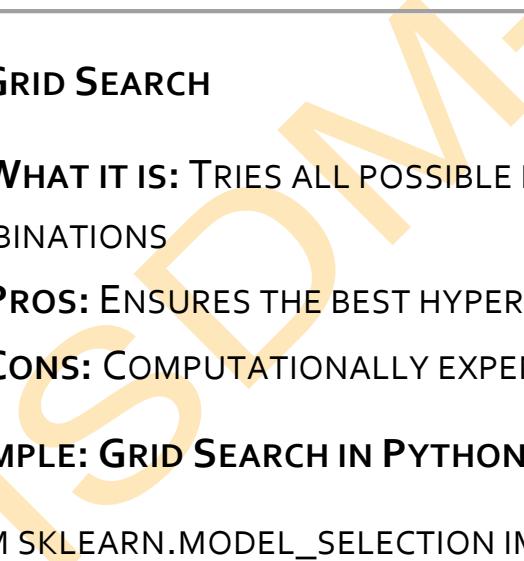
- ✓ N_ESTIMATORS – NUMBER OF DECISION TREES
- ✓ MAX_FEATURES – NUMBER OF FEATURES USED PER TREE
- ✓ BOOTSTRAP – WHETHER SAMPLES ARE DRAWN WITH REPLACEMENT

2.3 HYPERPARAMETERS IN NEURAL NETWORKS

- ✓ LEARNING_RATE – CONTROLS WEIGHT UPDATES DURING TRAINING
 - ✓ BATCH_SIZE – NUMBER OF SAMPLES PER TRAINING ITERATION
 - ✓ EPOCHS – NUMBER OF TIMES THE MODEL SEES THE FULL DATASET
-

CHAPTER 3: HYPERPARAMETER TUNING TECHNIQUES

3.1 MANUAL TUNING

- ◆ **WHAT IT IS:** SELECTING HYPERPARAMETERS MANUALLY BASED ON DOMAIN KNOWLEDGE
 - ◆ **PROS:** SIMPLE FOR SMALL MODELS
 - ◆ **CONS:** TIME-CONSUMING AND INEFFICIENT
- 

3.2 GRID SEARCH

- ◆ **WHAT IT IS:** TRIES ALL POSSIBLE HYPERPARAMETER COMBINATIONS
- ◆ **PROS:** ENSURES THE BEST HYPERPARAMETERS ARE FOUND
- ◆ **CONS:** COMPUTATIONALLY EXPENSIVE

EXAMPLE: GRID SEARCH IN PYTHON

```
FROM SKLEARN.MODEL_SELECTION IMPORT GRIDSEARCHCV  
FROM SKLEARN.ENSEMBLE IMPORT RANDOMFORESTCLASSIFIER
```

```
# DEFINE MODEL  
  
MODEL = RANDOMFORESTCLASSIFIER()
```

```
# DEFINE HYPERPARAMETER GRID  
  
PARAM_GRID = {'N_ESTIMATORS': [50, 100, 200], 'MAX_DEPTH': [5,  
10, None]}
```

```
# PERFORM GRID SEARCH  
  
GRID_SEARCH = GRIDSEARCHCV(MODEL, PARAM_GRID, CV=5)  
GRID_SEARCH.FIT(X_TRAIN, Y_TRAIN)  
  
# GET BEST PARAMETERS  
  
PRINT("BEST PARAMETERS:", GRID_SEARCH.BEST_PARAMS_)
```

3.3 RANDOM SEARCH

- ◆ **WHAT IT IS:** RANDOMLY SELECTS HYPERPARAMETERS INSTEAD OF TRYING ALL COMBINATIONS
- ◆ **PROS:** FASTER THAN GRID SEARCH
- ◆ **CONS:** MAY NOT FIND THE ABSOLUTE BEST PARAMETERS

EXAMPLE: RANDOM SEARCH IN PYTHON

```
FROM SKLEARN.MODEL_SELECTION IMPORT  
RANDOMIZEDSEARCHCV  
  
IMPORT NUMPY AS NP
```

```
# DEFINE HYPERPARAMETER DISTRIBUTIONS
```

```
PARAM_DIST = {'N_ESTIMATORS': np.arange(50, 200, 50),  
'MAX_DEPTH': [5, 10, None]}
```

```
# PERFORM RANDOM SEARCH
```

```
RANDOM_SEARCH = RandomizedSearchCV(MODEL, PARAM_DIST,  
n_iter=5, cv=5)
```

```
RANDOM_SEARCH.FIT(X_TRAIN, Y_TRAIN)
```

```
PRINT("BEST PARAMETERS:", RANDOM_SEARCH.BEST_PARAMS_)
```

3.4 BAYESIAN OPTIMIZATION

- ◆ **WHAT IT IS:** USES PROBABILITY MODELS TO FIND THE BEST HYPERPARAMETERS
- ◆ **PROS:** MORE EFFICIENT THAN GRID AND RANDOM SEARCH
- ◆ **CONS:** REQUIRES ADDITIONAL LIBRARIES LIKE SCIKIT-OPTIMIZE

EXAMPLE: BAYESIAN OPTIMIZATION IN PYTHON

```
FROM SKOPT IMPORT BAYESSEARCHCV
```

```
# PERFORM BAYESIAN OPTIMIZATION
```

```
BAYES_SEARCH = BAYESSEARCHCV(MODEL, PARAM_GRID,  
n_iter=10, cv=5)
```

```
BAYES_SEARCH.FIT(X_TRAIN, Y_TRAIN)
```

```
PRINT("BEST PARAMETERS:", BAYES_SEARCH.BEST_PARAMS_)
```

3.5 AUTOMATED HYPERPARAMETER TUNING (AUTOML)

- ◆ **WHAT IT IS:** USES AI TO OPTIMIZE HYPERPARAMETERS AUTOMATICALLY
- ◆ **PROS:** HANDS-OFF APPROACH FOR COMPLEX MODELS
- ◆ **CONS:** LIMITED CONTROL OVER TUNING PROCESS

POPULAR AUTOML TOOLS ✓ GOOGLE AUTOML

✓ TPOT (TREE-BASED PIPELINE OPTIMIZATION TOOL)

✓ H2O AUTOML

📌 CHAPTER 4: MODEL OPTIMIZATION TECHNIQUES
HYPERPARAMETER TUNING ALONE ISN'T ENOUGH. HERE'S HOW TO FURTHER IMPROVE MODEL PERFORMANCE:

4.1 FEATURE ENGINEERING

- ✓ REMOVING IRRELEVANT FEATURES
- ✓ CREATING NEW INFORMATIVE FEATURES
- ✓ NORMALIZING OR SCALING DATA

4.2 CROSS-VALIDATION

- ✓ SPLITS DATA INTO TRAINING & VALIDATION SETS MULTIPLE TIMES
- ✓ ENSURES MODEL GENERALIZES WELL TO NEW DATA

EXAMPLE OF K-FOLD CROSS-VALIDATION:

```
FROM SKLEARN.MODEL_SELECTION IMPORT CROSS_VAL_SCORE
```

```
SCORES = CROSS_VAL_SCORE(MODEL, X_TRAIN, Y_TRAIN, CV=5)
```

```
PRINT("MEAN ACCURACY:", SCORES.MEAN())
```

4.3 EARLY STOPPING (FOR NEURAL NETWORKS)

✓ STOPS TRAINING WHEN MODEL PERFORMANCE STARTS DEGRADING

✓ PREVENTS OVERFITTING

EXAMPLE IN TENSORFLOW:

```
FROM TENSORFLOW.KERAS.CALLBACKS IMPORT EARLYSTOPPING
```

```
EARLY_STOPPING = EARLYSTOPPING(MONITOR='VAL_LOSS',  
PATIENCE=3)
```

```
MODEL.FIT(X_TRAIN, Y_TRAIN, VALIDATION_DATA=(X_VAL, Y_VAL),  
CALLBACKS=[EARLY_STOPPING])
```

📌 CHAPTER 5: EXERCISES & ASSIGNMENTS

5.1 MULTIPLE CHOICE QUESTIONS

1. WHAT IS HYPERPARAMETER TUNING USED FOR?

- (A) CHANGING MODEL STRUCTURE
- (B) FINDING THE BEST MODEL PARAMETERS 

- (C) INCREASING TRAINING TIME
- (D) REMOVING MISSING VALUES

2. WHICH METHOD EXHAUSTIVELY TRIES ALL HYPERPARAMETER COMBINATIONS?

- (A) GRID SEARCH
- (B) RANDOM SEARCH
- (C) BAYESIAN OPTIMIZATION
- (D) EARLY STOPPING

3. WHICH OPTIMIZATION TECHNIQUE IS USED IN NEURAL NETWORKS TO PREVENT OVERFITTING?

- (A) FEATURE ENGINEERING
- (B) CROSS-VALIDATION
- (C) EARLY STOPPING
- (D) RANDOM SEARCH

5.2 PRACTICAL ASSIGNMENTS

📌 **TASK 1:** PERFORM GRID SEARCH AND RANDOM SEARCH ON A DATASET AND COMPARE RESULTS.

📌 **TASK 2:** APPLY BAYESIAN OPTIMIZATION TO A RANDOM FOREST MODEL AND REPORT THE BEST HYPERPARAMETERS.

📌 **CHAPTER 6: SUMMARY**

HYPERPARAMETER TUNING IMPROVES MODEL PERFORMANCE BY SELECTING THE BEST CONFIGURATIONS.

GRID SEARCH IS EXHAUSTIVE BUT SLOW, WHILE **RANDOM SEARCH** IS FASTER.

- BAYESIAN OPTIMIZATION BALANCES EFFICIENCY AND ACCURACY.**
 - MODEL OPTIMIZATION TECHNIQUES LIKE CROSS-VALIDATION AND EARLY STOPPING PREVENT OVERRFITTING.**
-

 **CONCLUSION: THE FUTURE OF HYPERPARAMETER TUNING**
AUTOMATED TOOLS LIKE **GOOGLE AUTOML** AND **TPOT** ARE
MAKING HYPERPARAMETER TUNING FASTER AND MORE EFFICIENT.
WITH ADVANCEMENTS IN AI-DRIVEN OPTIMIZATION, MODEL
PERFORMANCE WILL CONTINUE TO IMPROVE ACROSS INDUSTRIES.

ISDM-NINJA



IMPLEMENTING SUPERVISED ML MODELS USING SCIKIT-LEARN

📌 CHAPTER 1: INTRODUCTION TO SUPERVISED MACHINE LEARNING

1.1 WHAT IS SUPERVISED LEARNING?

SUPERVISED LEARNING IS A TYPE OF MACHINE LEARNING WHERE A MODEL LEARNS FROM **LABLED DATA**. IT MAPS INPUT VARIABLES (X) TO AN OUTPUT VARIABLE (Y).

1.2 TYPES OF SUPERVISED LEARNING

- ◆ **CLASSIFICATION** – PREDICTING CATEGORIES (E.G., SPAM OR NOT SPAM).
- ◆ **REGRESSION** – PREDICTING CONTINUOUS VALUES (E.G., HOUSE PRICES).

1.3 POPULAR SUPERVISED LEARNING ALGORITHMS

- ✓ **LOGISTIC REGRESSION** – USED FOR BINARY CLASSIFICATION.
- ✓ **DECISION TREES** – RULE-BASED DECISION MAKING.
- ✓ **RANDOM FOREST** – A COLLECTION OF DECISION TREES FOR BETTER ACCURACY.
- ✓ **SUPPORT VECTOR MACHINES (SVMs)** – WORKS WELL FOR CLASSIFICATION.
- ✓ **K-NEAREST NEIGHBORS (KNN)** – FINDS THE NEAREST NEIGHBORS FOR PREDICTION.
- ✓ **NAÏVE BAYES** – WORKS WELL FOR TEXT CLASSIFICATION.

✓ NEURAL NETWORKS – ADVANCED MODEL FOR COMPLEX PATTERNS.

📌 CHAPTER 2: IMPLEMENTING CLASSIFICATION USING SCIKIT-LEARN

2.1 STEPS TO BUILD A CLASSIFICATION MODEL

1. LOAD THE DATASET.
2. SPLIT DATA INTO TRAINING AND TESTING SETS.
3. TRAIN THE MODEL USING A SUPERVISED LEARNING ALGORITHM.
4. MAKE PREDICTIONS.
5. EVALUATE THE MODEL USING METRICS LIKE ACCURACY AND F₁-SCORE.

2.2 EXAMPLE: IMPLEMENTING LOGISTIC REGRESSION

PYTHON

COPY EDIT

IMPORT PANDAS AS PD

FROM SKLEARN.MODEL_SELECTION IMPORT TRAIN_TEST_SPLIT

FROM SKLEARN.LINEAR_MODEL IMPORT LOGISTICREGRESSION

FROM SKLEARN.METRICS IMPORT ACCURACY_SCORE

LOAD DATASET (EXAMPLE: BREAST CANCER DATASET)

FROM SKLEARN.DATASETS IMPORT LOAD_BREAST_CANCER

DATA = LOAD_BREAST_CANCER()

```
X = PD.DATAFRAME(DATA.DATA,  
COLUMNS=DATA.FEATURE_NAMES)  
  
Y = DATA.TARGET  
  
# SPLIT DATASET  
  
X_TRAIN, X_TEST, Y_TRAIN, Y_TEST = TRAIN_TEST_SPLIT(X, Y,  
TEST_SIZE=0.2, RANDOM_STATE=42)  
  
# TRAIN LOGISTIC REGRESSION MODEL  
  
MODEL = LOGISTICREGRESSION(MAX_ITER=1000)  
MODEL.FIT(X_TRAIN, Y_TRAIN)  
  
# MAKE PREDICTIONS  
  
Y_PRED = MODEL.PREDICT(X_TEST)  
  
# EVALUATE ACCURACY  
  
PRINT("ACCURACY:", ACCURACY_SCORE(Y_TEST, Y_PRED))
```

📌 CHAPTER 3: IMPLEMENTING REGRESSION USING SCIKIT-LEARN

3.1 EXAMPLE: IMPLEMENTING LINEAR REGRESSION

PYTHON

COPYEDIT

```
FROM SKLEARN.LINEAR_MODEL IMPORT LINEARREGRESSION  
IMPORT NUMPY AS NP
```

SAMPLE DATA

```
X = NP.ARRAY([1, 2, 3, 4, 5]).RESHAPE(-1, 1)  
Y = NP.ARRAY([2, 4, 5, 4, 5])
```

TRAIN MODEL

```
REGRESSOR = LINEARREGRESSION()  
REGRESSOR.FIT(X, Y)
```

PREDICT NEW VALUES

```
PREDICTIONS = REGRESSOR.PREDICT(X)  
PRINT("PREDICTIONS:", PREDICTIONS)
```

-
- 📌 CHAPTER 4: EXERCISES & ASSIGNMENTS
 - 📌 **TASK 1:** IMPLEMENT A DECISION TREE CLASSIFIER USING SCIKIT-LEARN.
 - 📌 **TASK 2:** TRAIN A RANDOM FOREST MODEL AND COMPARE IT WITH DECISION TREES.
-

📌 CHAPTER 5: SUMMARY

- ✓ **SUPERVISED LEARNING** REQUIRES LABELED DATA.
- ✓ **CLASSIFICATION MODELS** PREDICT CATEGORIES.
- ✓ **REGRESSION MODELS** PREDICT CONTINUOUS VALUES.
- ✓ **SCIKIT-LEARN** MAKES MODEL TRAINING EASY.

🌟 CONCLUSION: THE POWER OF SCIKIT-LEARN

SCIKIT-LEARN IS A POWERFUL TOOL FOR **QUICKLY BUILDING ML MODELS**. WITH MORE DATA AND **HYPERPARAMETER TUNING**, WE CAN IMPROVE ACCURACY AND GENERALIZATION.

ISDM-NXT

  **ASSIGNMENT 1:**
 **BUILD A HOUSE PRICE PREDICTION
MODEL USING LINEAR REGRESSION.**

ISDM-NXT

📌 ⚡ ASSIGNMENT SOLUTION 1: BUILD A HOUSE PRICE PREDICTION MODEL USING LINEAR REGRESSION

🎯 OBJECTIVE

THE GOAL OF THIS ASSIGNMENT IS TO BUILD A **HOUSE PRICE PREDICTION MODEL** USING **LINEAR REGRESSION**. YOU WILL LEARN TO LOAD DATA, CLEAN IT, PERFORM EXPLORATORY DATA ANALYSIS (EDA), BUILD A REGRESSION MODEL, AND EVALUATE ITS PERFORMANCE.

🛠 STEP 1: IMPORT NECESSARY LIBRARIES

WE WILL USE **PYTHON** AND ITS DATA SCIENCE LIBRARIES TO BUILD THE MODEL.

IMPORT NUMPY AS NP

IMPORT PANDAS AS PD

IMPORT MATPLOTLIB.PYPLOT AS PLT

IMPORT SEABORN AS SNS

FROM SKLEARN.MODEL_SELECTION IMPORT TRAIN_TEST_SPLIT

FROM SKLEARN.LINEAR_MODEL IMPORT LINEARREGRESSION

FROM SKLEARN.METRICS IMPORT MEAN_ABSOLUTE_ERROR,
MEAN_SQUARED_ERROR, R2_SCORE

📁 STEP 2: LOAD THE DATASET

WE WILL USE A PUBLICLY AVAILABLE **HOUSE PRICE DATASET** FROM **KAGGLE** OR GENERATE SAMPLE DATA.

```
# LOAD THE DATASET (REPLACE 'HOUSE_PRICES.CSV' WITH YOUR  
ACTUAL DATASET)
```

```
DF = PD.READ_CSV("HOUSE_PRICES.CSV")
```

```
# DISPLAY FIRST 5 ROWS
```

```
DF.HEAD()
```

🔍 STEP 3: EXPLORATORY DATA ANALYSIS (EDA)

3.1 CHECK DATASET INFORMATION

```
# GET DATASET INFO
```

```
DF.INFO()
```

```
# CHECK FOR MISSING VALUES
```

```
DF.ISNULL().SUM()
```

✓ ACTION PLAN:

- IF MISSING VALUES EXIST, FILL THEM USING THE MEDIAN OR DROP THEM.

3.2 VISUALIZING HOUSE PRICES DISTRIBUTION

```
# HISTOGRAM OF HOUSE PRICES
```

```
PLT.FIGURE(figsize=(8,5))
```

```
SNS.HISTPLOT(DF['PRICE'], BINS=30, KDE=TRUE, COLOR='BLUE')

PLT.TITLE("HOUSE PRICE DISTRIBUTION")

PLT.XLABEL("PRICE")

PLT.YLABEL("COUNT")

PLT.SHOW()
```

3.3 IDENTIFY CORRELATIONS

```
# CORRELATION HEATMAP

PLT.FIGURE(figsize=(10,6))

SNS.HEATMAP(DF.CORR(), ANNOT=True, CMAP="COOLWARM")

PLT.TITLE("FEATURE CORRELATION MATRIX")

PLT.SHOW()
```

✓ INSIGHT:

- IDENTIFY THE MOST IMPORTANT FEATURES AFFECTING HOUSE PRICES.

STEP 4: SELECT FEATURES & TARGET VARIABLE

4.1 CHOOSING THE RIGHT FEATURES

WE SELECT IMPORTANT INDEPENDENT VARIABLES (**X**) AFFECTING HOUSE PRICES.

```
# DEFINE INDEPENDENT VARIABLES (X) AND DEPENDENT VARIABLE (Y)
```

```
X = DF[['SQFT_LIVING', 'BEDROOMS', 'BATHROOMS', 'FLOORS',  
'SQFT_LOT']]
```

```
Y = DF['PRICE']
```

STEP 5: SPLIT THE DATA INTO TRAINING & TESTING SETS

WE DIVIDE THE DATASET INTO **TRAINING (80%) AND TESTING (20%)** SETS.

```
# SPLIT INTO TRAINING AND TESTING SETS (80-20 RATIO)
```

```
X_TRAIN, X_TEST, Y_TRAIN, Y_TEST = TRAIN_TEST_SPLIT(X, Y,  
TEST_SIZE=0.2, RANDOM_STATE=42)
```

STEP 6: BUILD AND TRAIN THE LINEAR REGRESSION MODEL

WE TRAIN A **LINEAR REGRESSION MODEL** USING **SCIKIT-LEARN**.

```
# CREATE LINEAR REGRESSION MODEL
```

```
MODEL = LINEARREGRESSION()
```

```
# TRAIN THE MODEL
```

```
MODEL.FIT(X_TRAIN, Y_TRAIN)
```

STEP 7: MAKE PREDICTIONS

NOW, WE USE THE TRAINED MODEL TO PREDICT HOUSE PRICES.

```
# MAKE PREDICTIONS ON THE TEST SET
```

```
Y_PRED = MODEL.PREDICT(X_TEST)
```

✍ STEP 8: EVALUATE MODEL PERFORMANCE

WE EVALUATE THE MODEL USING **MEAN ABSOLUTE ERROR (MAE)**, **MEAN SQUARED ERROR (MSE)**, AND **R² SCORE**.

```
# CALCULATE PERFORMANCE METRICS
```

```
MAE = MEAN_ABSOLUTE_ERROR(Y_TEST, Y_PRED)
```

```
MSE = MEAN_SQUARED_ERROR(Y_TEST, Y_PRED)
```

```
RMSE = NP.SQRT(MSE)
```

```
R2 = R2_SCORE(Y_TEST, Y_PRED)
```

```
# DISPLAY RESULTS
```

```
PRINT(F"MEAN ABSOLUTE ERROR (MAE): {MAE}")
```

```
PRINT(F"MEAN SQUARED ERROR (MSE): {MSE}")
```

```
PRINT(F"ROOT MEAN SQUARED ERROR (RMSE): {RMSE}")
```

```
PRINT(F"R2 SCORE: {R2}")
```

✓ INTERPRETATION:

- **MAE & RMSE** SHOULD BE LOW (LESS ERROR MEANS BETTER PREDICTIONS).
- **R² SCORE** SHOULD BE CLOSE TO 1 (HIGHER MEANS BETTER ACCURACY).

📊 STEP 9: VISUALIZING THE MODEL PREDICTIONS

WE COMPARE ACTUAL VS PREDICTED HOUSE PRICES.

```
PLT.FIGURE(FIGSIZE=(8,5))
```

```
PLT.SCATTER(Y_TEST, Y_PRED, COLOR='BLUE', ALPHA=0.5)
```

```
PLT.PLOT([MIN(Y_TEST), MAX(Y_TEST)], [MIN(Y_TEST),  
MAX(Y_TEST)], COLOR='RED', LINESTYLE='DASHED')
```

```
PLT.TITLE("ACTUAL VS PREDICTED HOUSE PRICES")
```

```
PLT.XLABEL("ACTUAL PRICE")
```

```
PLT.YLABEL("PREDICTED PRICE")
```

```
PLT.SHOW()
```

✓ INSIGHT:

- IF POINTS ALIGN WITH THE RED DASHED LINE, THE PREDICTIONS ARE ACCURATE.
-

📌 STEP 10: CONCLUSION

✓ KEY TAKEAWAYS

- WE LOADED AND EXPLORED THE DATASET.
- WE IDENTIFIED CORRELATIONS AND SELECTED THE BEST FEATURES.
- WE TRAINED A LINEAR REGRESSION MODEL TO PREDICT HOUSE PRICES.
- WE EVALUATED PERFORMANCE USING MAE, MSE, RMSE, AND R² SCORE.

- MODEL PREDICTIONS WERE VISUALIZED TO ASSESS ACCURACY.

NEXT STEPS

- TRY ADDING MORE FEATURES (E.G., LOCATION, YEAR_BUILT).
- EXPERIMENT WITH POLYNOMIAL REGRESSION FOR BETTER PREDICTIONS.
- OPTIMIZE THE MODEL USING FEATURE ENGINEERING AND HYPERPARAMETER TUNING.

ISDM-NXT

📌 ⚡ ASSIGNMENT 2:
🎯 IMPLEMENT A SPAM EMAIL CLASSIFIER
USING LOGISTIC REGRESSION.

ISDM-NxT

ASSIGNMENT SOLUTION 2: IMPLEMENT A SPAM EMAIL CLASSIFIER USING LOGISTIC REGRESSION

OBJECTIVE

THE GOAL OF THIS ASSIGNMENT IS TO BUILD A **SPAM EMAIL CLASSIFIER USING LOGISTIC REGRESSION**. THE CLASSIFIER WILL ANALYZE AN EMAIL'S TEXT AND PREDICT WHETHER IT IS **SPAM (1)** OR **NOT SPAM (0)**.

STEP 1: INSTALL AND IMPORT NECESSARY LIBRARIES

WE WILL USE THE FOLLOWING PYTHON LIBRARIES:

- **PANDAS:** FOR HANDLING DATASETS.
 - **NUMPY:** FOR NUMERICAL OPERATIONS.
 - **SKLEARN:** FOR MACHINE LEARNING AND TEXT PROCESSING.
- ◆ **INSTALL DEPENDENCIES (IF NOT INSTALLED)**

!PIP INSTALL PANDAS NUMPY SCIKIT-LEARN

◆ **IMPORT REQUIRED LIBRARIES**

IMPORT PANDAS AS PD

IMPORT NUMPY AS NP

IMPORT RE

IMPORT STRING

```
FROM SKLEARN.MODEL_SELECTION IMPORT TRAIN_TEST_SPLIT  
FROM SKLEARN.FEATURE_EXTRACTION.TEXT IMPORT  
TFIDFVECTORIZER  
FROM SKLEARN.LINEAR_MODEL IMPORT LOGISTICREGRESSION  
FROM SKLEARN.METRICS IMPORT ACCURACY_SCORE,  
CLASSIFICATION_REPORT, CONFUSION_MATRIX
```

STEP 2: LOAD AND EXPLORE THE DATASET

WE WILL USE THE SPAM EMAIL DATASET, WHICH CONSISTS OF LABELED EMAILS.

◆ LOAD THE DATASET

```
# LOAD THE DATASET (ASSUMING A CSV FILE WITH 'MESSAGE' AND  
'LABEL' COLUMNS)
```

```
DF = PD.READ_CSV('SPAM.CSV', ENCODING='LATIN-1')
```

```
# DISPLAY FIRST FEW ROWS
```

```
PRINT(DF.HEAD())
```

◆ CHECK DATASET INFORMATION

```
DF.INFO()
```

```
DF.ISNULL().SUM() # CHECK FOR MISSING VALUES
```

◆ UNDERSTANDING THE DATA

- **MESSAGE:** THE ACTUAL EMAIL TEXT.

- **LABEL:** THE CLASSIFICATION (SPAM = 1, NOT SPAM = 0).

CONVERT LABELS INTO BINARY VALUES

```
DF = DF[['V1', 'V2']] # SELECT RELEVANT COLUMNS  
  
DF.COLUMNS = ['LABEL', 'MESSAGE'] # RENAME COLUMNS
```

```
# CONVERT LABELS TO NUMERIC VALUES (SPAM = 1, NOT SPAM = 0)
```

```
DF['LABEL'] = DF['LABEL'].MAP({'SPAM': 1, 'HAM': 0})
```

```
# DISPLAY UPDATED DATA
```

```
PRINT(DF.HEAD())
```

✍ STEP 3: PREPROCESS THE DATA

BEFORE TRAINING THE MODEL, WE NEED TO CLEAN AND PREPROCESS THE TEXT.

◆ DEFINE A FUNCTION TO CLEAN TEXT

```
DEF PREPROCESS_TEXT(TEXT):  
  
    TEXT = TEXT.LOWER() # CONVERT TO LOWERCASE  
  
    TEXT = RE.SUB(R'\D+', "", TEXT) # REMOVE NUMBERS  
  
    TEXT = RE.SUB(R'\S+', ' ', TEXT) # REMOVE EXTRA SPACES  
  
    TEXT = TEXT.TRANSLATE(STR.MAKETRANS("", "",  
    STRING.PUNCTUATION)) # REMOVE PUNCTUATION  
  
    RETURN TEXT STRIP()
```

◆ APPLY CLEANING FUNCTION

```
DF['CLEANED_MESSAGE'] =  
DF['MESSAGE'].APPLY(PREPROCESS_TEXT)  
  
PRINT(DF.HEAD()) # DISPLAY CLEANED DATA
```

■ STEP 4: CONVERT TEXT TO NUMERICAL FEATURES

SINCE MACHINE LEARNING MODELS DO NOT UNDERSTAND TEXT, WE NEED TO CONVERT THE EMAILS INTO **NUMERICAL FEATURES** USING **TF-IDF VECTORIZATION**.

◆ CONVERT TEXT TO TF-IDF FEATURES

```
VECTORIZER = TfidfVectorizer(STOP_WORDS='ENGLISH',  
MAX_FEATURES=5000)
```

```
# TRANSFORM TEXT DATA INTO NUMERICAL FORMAT
```

```
X = VECTORIZER.FIT_TRANSFORM(DF['CLEANED_MESSAGE'])
```

```
# TARGET VARIABLE (SPAM OR NOT SPAM)
```

```
Y = DF['LABEL']
```

■ STEP 5: SPLIT THE DATA INTO TRAINING AND TESTING SETS

```
X_TRAIN, X_TEST, Y_TRAIN, Y_TEST = train_test_split(X, Y,  
TEST_SIZE=0.2, RANDOM_STATE=42)
```

```
PRINT("TRAINING DATA SHAPE:", X_TRAIN.SHAPE)
```

```
PRINT("TESTING DATA SHAPE:", X_TEST.SHAPE)
```

STEP 6: TRAIN THE LOGISTIC REGRESSION MODEL

```
# INITIALIZE AND TRAIN THE LOGISTIC REGRESSION MODEL
```

```
MODEL = LOGISTICREGRESSION()
```

```
MODEL.FIT(X_TRAIN, Y_TRAIN)
```

```
PRINT("MODEL TRAINING COMPLETE ")
```

STEP 7: MAKE PREDICTIONS

```
# PREDICT ON THE TEST SET
```

```
Y_PRED = MODEL.PREDICT(X_TEST)
```

STEP 8: EVALUATE THE MODEL

WE EVALUATE THE MODEL USING ACCURACY, CONFUSION MATRIX, AND CLASSIFICATION REPORT.

◆ ACCURACY SCORE

```
ACCURACY = ACCURACY_SCORE(Y_TEST, Y_PRED)
```

```
PRINT(F"MODEL ACCURACY: {ACCURACY:.2F}")
```

◆ CLASSIFICATION REPORT

```
PRINT("CLASSIFICATION REPORT:\n",
CLASSIFICATION_REPORT(Y_TEST, Y_PRED))
```

◆ CONFUSION MATRIX

```
IMPORT SEABORN AS SNS
```

```
IMPORT MATPLOTLIB.PYPLOT AS PLT
```

```
# GENERATE CONFUSION MATRIX
CM = CONFUSION_MATRIX(Y_TEST, Y_PRED)

# VISUALIZING THE CONFUSION MATRIX
PLT.FIGURE(figsize=(5,4))
SNS.HEATMAP(CM, ANNOT=TRUE, FMT="D", CMAP="BLUES",
XTICKLABELS=['NOT SPAM', 'SPAM'], YTICKLABELS=['NOT SPAM',
'SPAM'])
PLT.XLABEL("PREDICTED")
PLT.YLABEL("ACTUAL")
PLT.TITLE("CONFUSION MATRIX")
PLT.SHOW()
```

STEP 9: TEST ON NEW EMAILS

NOW, LET'S TEST THE MODEL WITH SOME SAMPLE EMAILS.

```
DEF PREDICT_SPAM(TEXT):
```

```
TEXT = PREPROCESS_TEXT(TEXT) # PREPROCESS THE TEXT  
  
TEXT_TFIDF = VECTORIZER.TRANSFORM([TEXT]) # CONVERT TO  
TF-IDF  
  
PREDICTION = MODEL.PREDICT(TEXT_TFIDF)[0] # MAKE  
PREDICTION  
  
RETURN "SPAM" IF PREDICTION == 1 ELSE "NOT SPAM"
```

EXAMPLE EMAILS

EMAIL1 = "CONGRATULATIONS! YOU HAVE WON A LOTTERY. CLICK
HERE TO CLAIM YOUR PRIZE."

EMAIL2 = "HI JOHN, CAN WE SCHEDULE A MEETING FOR
TOMORROW?"

```
PRINT("EMAIL 1 PREDICTION:", PREDICT_SPAM(EMAIL1))
```

```
PRINT("EMAIL 2 PREDICTION:", PREDICT_SPAM(EMAIL2))
```

FINAL SUMMARY

STEP	DESCRIPTION
STEP 1	INSTALL AND IMPORT NECESSARY LIBRARIES
STEP 2	LOAD AND EXPLORE THE DATASET
STEP 3	PREPROCESS THE TEXT (CLEANING)

STEP 4	CONVERT TEXT TO NUMERICAL FEATURES USING TF-IDF
STEP 5	SPLIT DATA INTO TRAINING AND TESTING SETS
STEP 6	TRAIN A LOGISTIC REGRESSION MODEL
STEP 7	MAKE PREDICTIONS ON THE TEST SET
STEP 8	EVALUATE THE MODEL USING ACCURACY AND CONFUSION MATRIX
STEP 9	TEST THE MODEL ON NEW EMAIL SAMPLES

★ CONCLUSION

- WE SUCCESSFULLY BUILT A **SPAM EMAIL CLASSIFIER** USING **LOGISTIC REGRESSION**.
- THE MODEL CAN PREDICT WHETHER AN EMAIL IS **SPAM** OR **NOT SPAM** BASED ON ITS TEXT **CONTENT**.
- THE **TF-IDF** TECHNIQUE HELPED CONVERT TEXT INTO NUMERICAL VALUES.
- THIS CLASSIFIER CAN BE FURTHER **IMPROVED** BY USING **ADVANCED MODELS** LIKE **NAIVE BAYES** OR **DEEP LEARNING**.