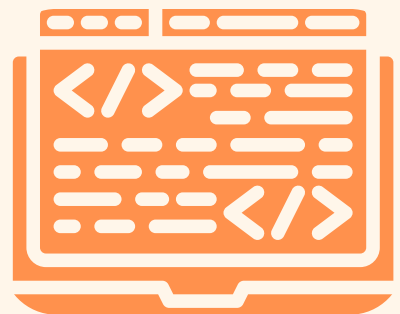


◇ IRIS AND REUTERS ◇

# DS TOOLS LAB WORK



**Supervised by  
Dr. Mohamed  
Eng. Sohaila**



## TEAM NAMES AND IDS

NAME	ID
MAHMOUD ESSAM	20221460231
ABDELRAHMAN ASHRAF	20221374041
ABDULLAH HUSSIEN	20221427861
FARES MUHAMMED	20221461330
ZYAD ASHRAF	20221374025



# IRIS DATA

THIS REPORT DELVES INTO A COMPREHENSIVE ANALYSIS OF THE IRIS DATASET, A WELL-KNOWN DATASET IN THE FIELD OF MACHINE LEARNING. THE DATASET ENCOMPASSES MEASUREMENTS OF SEPAL LENGTH, SEPAL WIDTH, PETAL LENGTH, AND PETAL WIDTH FOR THREE DISTINCT SPECIES OF IRIS FLOWERS: SETOSA, VERSICOLOR, AND VIRGINICA.

## DATA ATTRIBUTES

### SEPAL LENGTH (CM)

THIS COLUMN REPRESENTS THE LENGTH OF THE IRIS FLOWER'S SEPAL (THE OUTERMOST PART OF THE FLOWER).

IT IS MEASURED IN CENTIMETERS.

### SEPAL WIDTH (CM)

THIS COLUMN REPRESENTS THE WIDTH OF THE IRIS FLOWER'S SEPAL.

IT IS MEASURED IN CENTIMETERS.

### PETAL LENGTH (CM)

THIS COLUMN REPRESENTS THE LENGTH OF THE IRIS FLOWER'S PETAL (THE INNER PART OF THE FLOWER).

IT IS MEASURED IN CENTIMETERS.

### PETAL WIDTH (CM)

THIS COLUMN REPRESENTS THE WIDTH OF THE IRIS FLOWER'S PETAL.

IT IS MEASURED IN CENTIMETERS.

### SPECIES

THIS COLUMN REPRESENTS THE SPECIES OF THE IRIS FLOWER.

IT IS A CATEGORICAL VARIABLE WITH THREE POSSIBLE VALUES: 'SETOSA', 'VERSICOLOR', AND 'VIRGINICA'.

THE SPECIES IS THE TARGET VARIABLE THAT WE OFTEN WANT TO PREDICT IN MACHINE LEARNING TASKS.

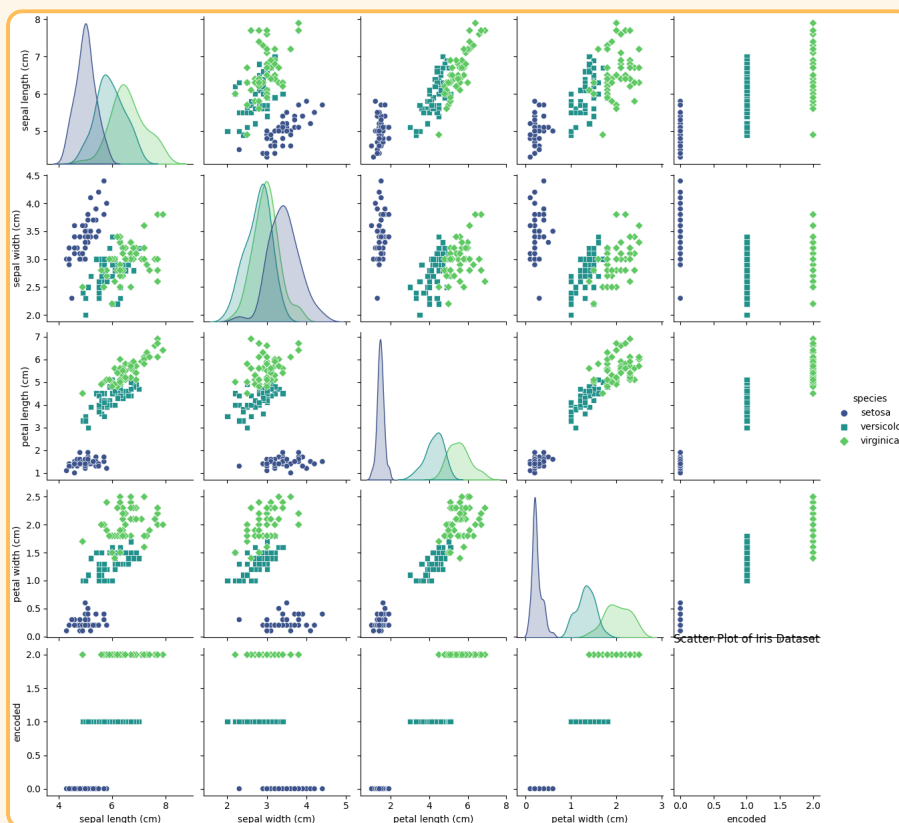
# WHAT HAPPENED IN PREPROCESS?

**A. DUPLICATED ENTRIES:** WE CHECKED FOR THE SAME DATA APPEARING MORE THAN ONCE AND GOT RID OF IT.

**B. NO MISSING STUFF:** WE LOOKED FOR MISSING INFORMATION, AND THERE WASN'T ANY.

**C. CHECK OF OUTLIERS** AND INTERVAL OF DATA

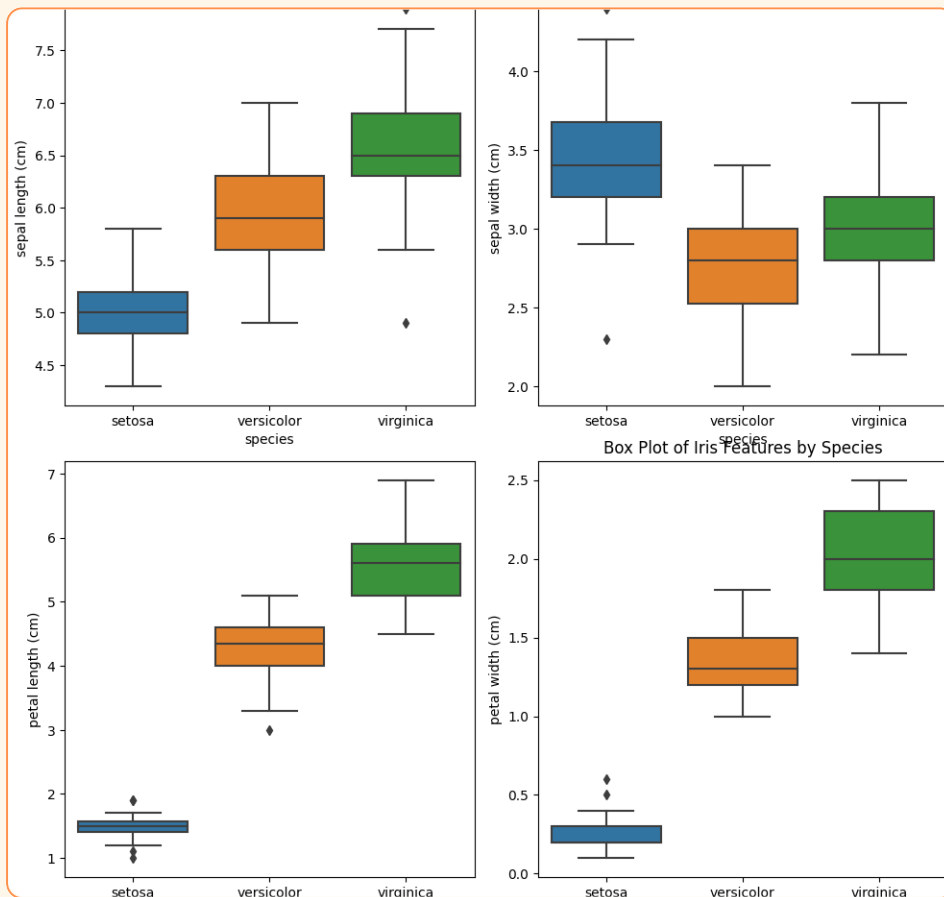
## WHAT ABOUT VISUALIZATION?



### SCATTER PLOT

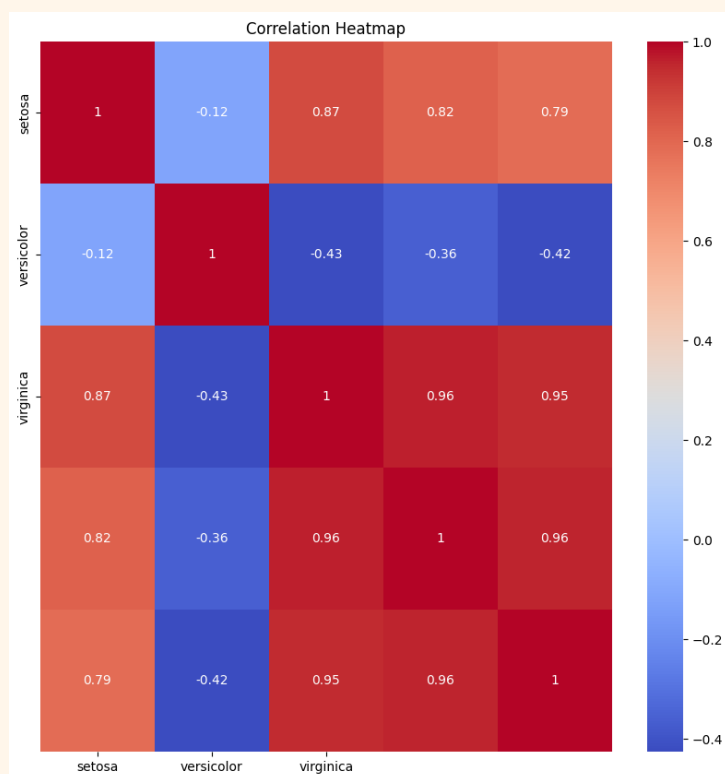
A SCATTER PLOT WAS CONSTRUCTED TO VISUALLY REPRESENT THE DISTRIBUTION OF IRIS FLOWERS BASED ON THEIR SEPAL AND PETAL MEASUREMENTS.

EACH SPECIES IS DIFFERENTIATED BY A DISTINCT MARKER, PROVIDING AN IMMEDIATE VISUAL UNDERSTANDING OF THE DATASET'S STRUCTURE.



### BOX PLOT

**BOX PLOTS WERE UTILISED TO PRESENT A SUMMARY OF THE DISTRIBUTION OF SEPAL AND PETAL MEASUREMENTS WITHIN EACH SPECIES. THIS AIDS IN IDENTIFYING POTENTIAL VARIATIONS AND OUTLIERS ACROSS DIFFERENT FEATURES.**



### CORRELATION HEATMAP

**A HEATMAP OF THE CORRELATION MATRIX WAS GENERATED TO ELUCIDATE THE RELATIONSHIPS BETWEEN SEPAL AND PETAL MEASUREMENTS. THIS VISUALISATION HELPS IN UNDERSTANDING THE INTERDEPENDENCIES AMONG DIFFERENT FEATURES.**

# WHAT MACHINE MODEL USED IN THE IRIS?

## DECISION TREE CLASSIFIER

A DECISION TREE CLASSIFIER WAS IMPLEMENTED WITH HYPERPARAMETER TUNING USING GRID SEARCH. THE BEST MODEL WAS IDENTIFIED BASED ON ITS PERFORMANCE METRICS IN ACCURATELY CLASSIFYING IRIS SPECIES.

## NEURAL NETWORK

A NEURAL NETWORK MODEL WAS CONSTRUCTED USING THE TENSORFLOW KERAS LIBRARY. THE MODEL ARCHITECTURE INVOLVES THREE LAYERS WITH RELU ACTIVATION IN HIDDEN LAYERS AND SOFTMAX ACTIVATION IN THE OUTPUT LAYER, FACILITATING MULTI-CLASS CLASSIFICATION.

# USER INTERFACE FOR PREDICTION

A NEURAL NETWORK MODEL WAS CONSTRUCTED USING THE TENSORFLOW KERAS LIBRARY. THE MODEL ARCHITECTURE INVOLVES THREE LAYERS WITH “RELU” ACTIVATION IN HIDDEN LAYERS AND SOFTMAX ACTIVATION IN THE OUTPUT LAYER, FACILITATING MULTI-CLASS CLASSIFICATION.

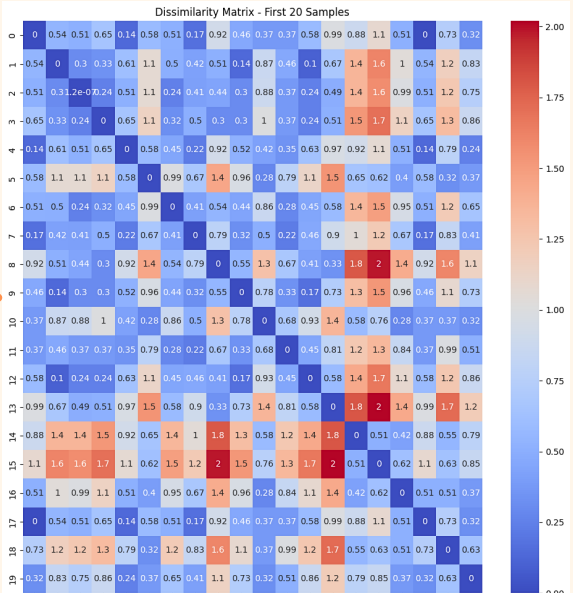
# DISSIMILARITY ANALYSIS

## DISSIMILARITY MATRIX

A DISSIMILARITY MATRIX WAS COMPUTED USING PAIRWISE EUCLIDEAN DISTANCES BETWEEN THE FIRST 20 SAMPLES OF THE IRIS DATASET. THE RESULTING MATRIX WAS VISUALISED THROUGH A HEATMAP, OFFERING INSIGHTS INTO THE DISSIMILARITY RELATIONSHIPS AMONG THE INITIAL DATASET ENTRIES.

# EXPORT TO CSV

THE DISSIMILARITY MATRIX WAS EXPORTED TO A CSV FILE, FACILITATING FURTHER ANALYSIS AND PROVIDING A TANGIBLE RECORD OF DISSIMILARITY VALUES.



# REUTERS

THE DATASET UNDER CONSIDERATION ORIGINATES FROM REUTERS A RENOWNED INTERNATIONAL NEWS AGENCY CELEBRATED FOR ITS EXTENSIVE AND CREDIBLE NEWS COVERAGE. COMPRISING A DIVERSE ARRAY OF NEWS ARTICLES, THIS DATASET SERVES AS A VALUABLE RESOURCE FOR NATURAL LANGUAGE PROCESSING (NLP) AND TEXT CLASSIFICATION TASKS WITHIN THE CONTEXT OF MACHINE LEARNING.

## DATA ATTRIBUTES

### TOPICS

THE 'TOPICS' COLUMN CATEGORISES NEWS ARTICLES INTO DISTINCT TOPICS OR SUBJECTS, PROVIDING A CATEGORICAL LABEL FOR EACH PIECE. THIS COLUMN ACTS AS THE TARGET VARIABLE FOR THE MACHINE LEARNING MODEL, INDICATING THE PRIMARY FOCUS OR THEME OF THE ARTICLES.

### BODY

THE 'BODY' COLUMN ENCOMPASSES THE PRIMARY TEXTUAL CONTENT OF THE NEWS ARTICLES. THROUGH PREPROCESSING AND ANALYSIS, THIS TEXT BECOMES THE FOCAL POINT FOR FEATURE EXTRACTION AND MODEL TRAINING.

WE NEED TO MAKE PERSON IF ADD SENTENCE, THEN THE RESULT WHICH TOPIC SUPPOSE THIS SENTENCE TO BE IN WHICH TOPIC

## WHAT HAPPENED IN PREPROCESS?

THE PROJECT BEGINS BY LOADING THE REUTERS DATASET FROM A CSV FILE, WHERE COLUMNS LIKE 'TOPICS' AND 'BODY' ARE IDENTIFIED.

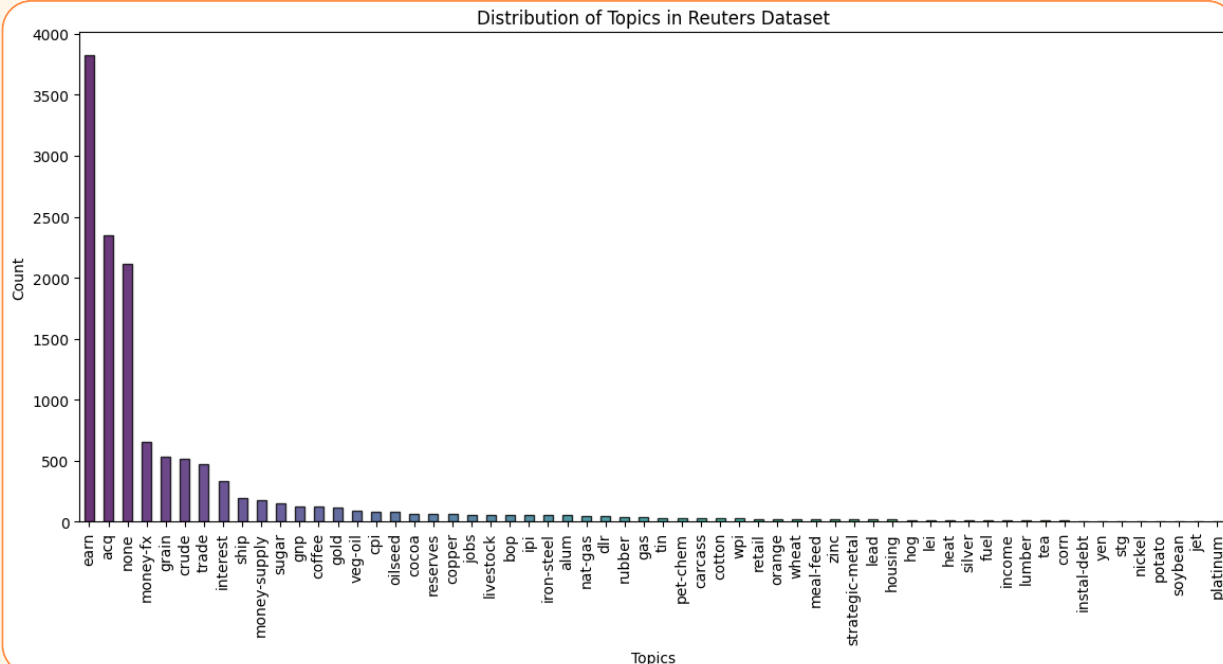
BASIC STATISTICAL INSIGHTS ARE PROVIDED, AND MISSING VALUES IN THE 'BODY' COLUMN ARE HANDLED THROUGH FILLING WITH EMPTY STRINGS.

TEXT PREPROCESSING INVOLVES CONVERTING THE 'BODY' TEXT TO LOWERCASE FOR CONSISTENCY.

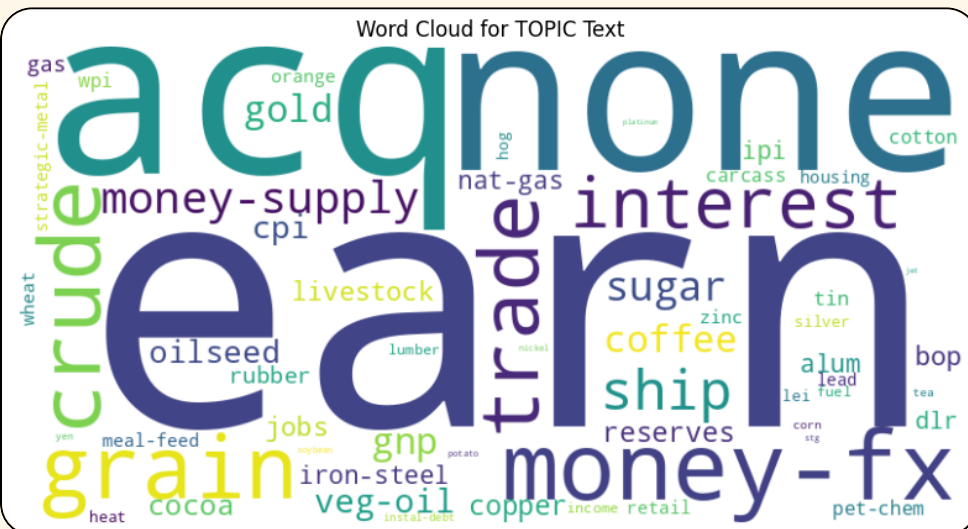


## WHAT ABOUT VISUALIZATION?

**THE CLASS DISTRIBUTION OF TOPICS IS PRESENTED USING A BAR CHART, FILTERING OUT TOPICS WITH COUNTS LESS THAN 3 FOR CLARITY.**

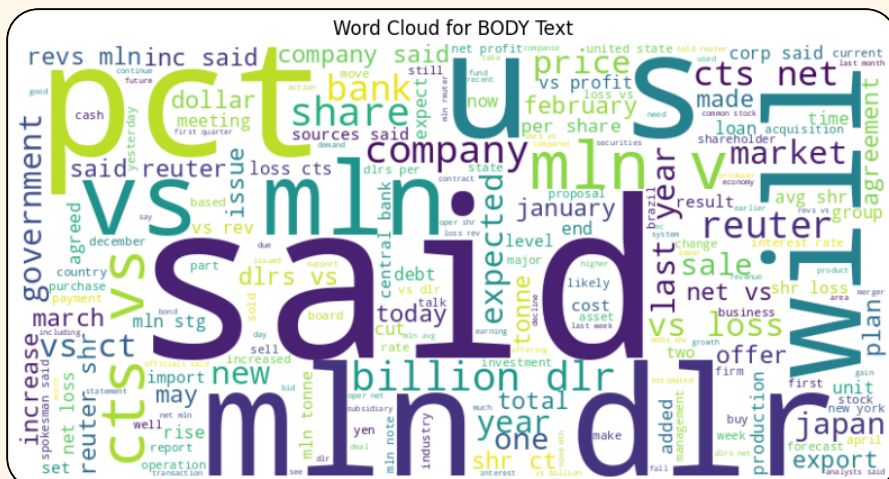


**WORD CLOUDS ARE GENERATED FOR BOTH 'TOPIC' AND 'BODY' TEXTS, OFFERING A VISUAL REPRESENTATION OF THE MOST FREQUENT WORDS IN EACH CATEGORY.**



**IT AIDS IN IDENTIFYING AND  
EXPLORING THE FREQUENT  
TOPICS COVERED IN THE NEWS  
ARTICLES, ALLOWING FOR A  
QUICK AND INTUITIVE  
UNDERSTANDING.**

**THE CODE COMBINES NEWS ARTICLE CONTENT. IT MAKES A WORD CLOUD TO SHOW FREQUENT WORDS VISUALLY. THIS HELPS QUICKLY SEE COMMON TERMS IN THE DATA. THE VISUALISATION IS SIMPLE AND AIDS CONTENT EXPLORATION.**





# WHAT MACHINE MODEL USED IN THE REUTERS?

## PREPARING THE DATA (ENCODING BODY)

THE CODE BEGINS BY PREPARING THE TEXT DATA FROM THE REUTERS DATASET. IT LIMITS THE NUMBER OF WORDS TO CONSIDER AND CONVERTS THE TEXT INTO NUMERICAL SEQUENCES, THIS WILL EASE THE PROCESS OF MAKING TEXT TO NUMBER

## ENCODING TOPICS

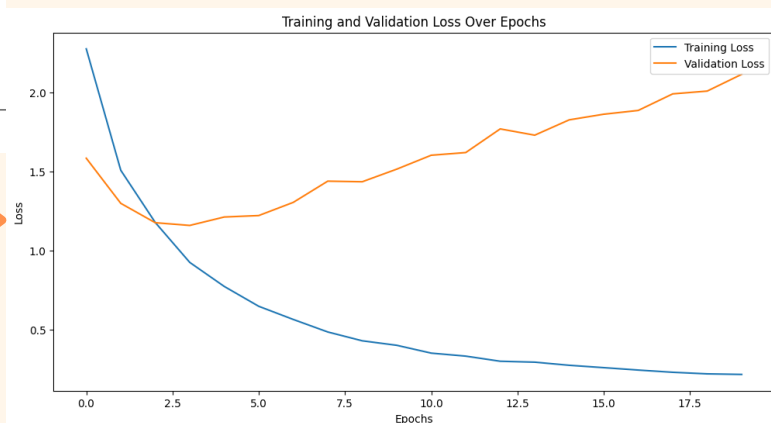
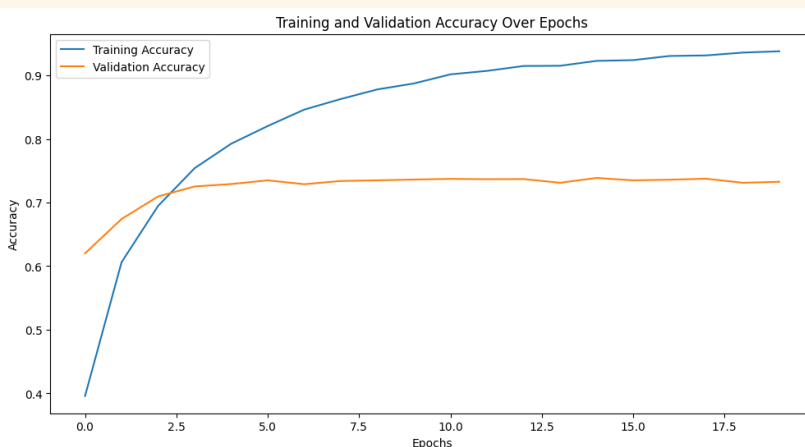
IT ENCODES THE TOPIC LABELS INTO NUMERICAL VALUES SO THAT THE MODEL CAN UNDERSTAND AND LEARN FROM THEM.

## BUILDING THE MODEL

THE NEURAL NETWORK MODEL IS CREATED.  
IT INVOLVES LAYERS THAT HELP THE MODEL UNDERSTAND THE RELATIONSHIPS WITHIN THE DATA. THE MODEL IS DESIGNED TO LEARN PATTERNS AND CONNECTIONS IN THE TEXT.

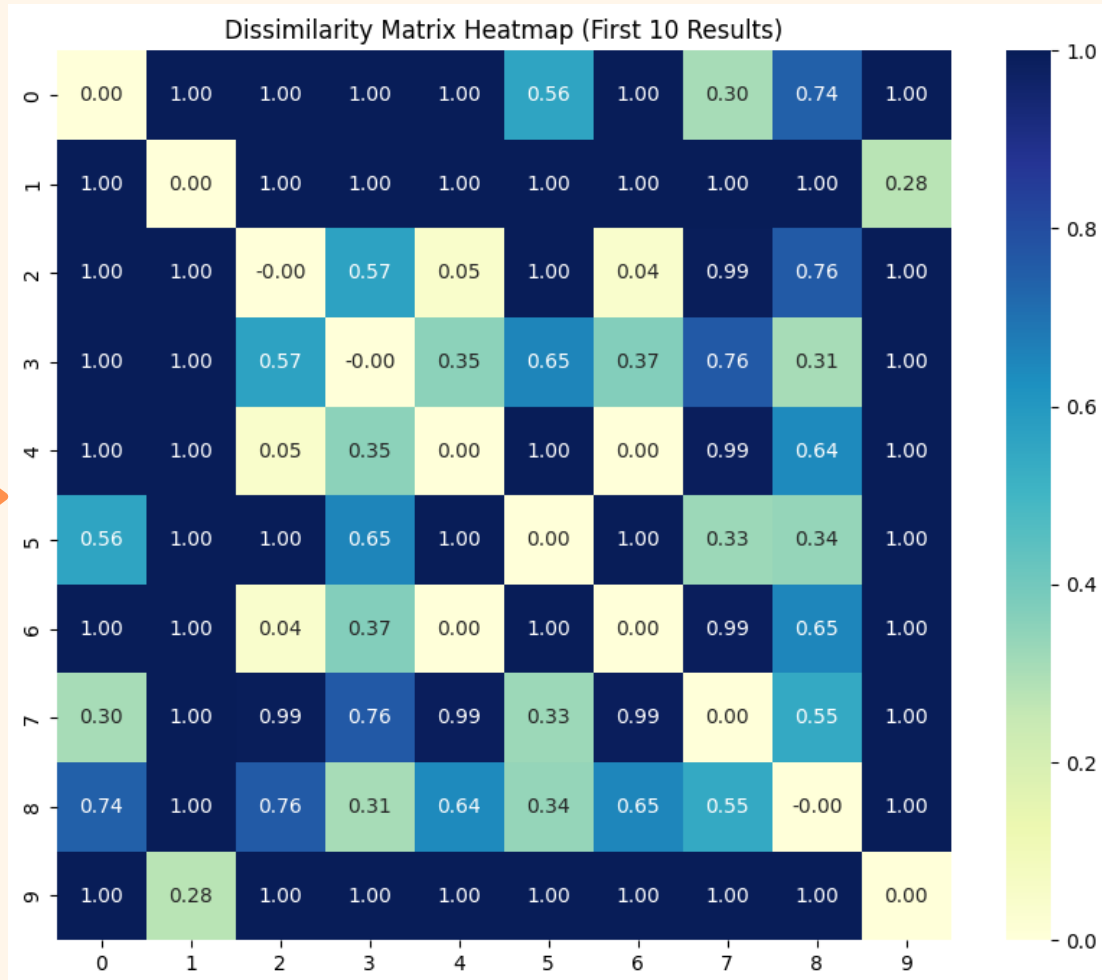
## VISUALIZING TRAINING PROGRESS

GRAPHS ARE CREATED TO SHOW HOW WELL THE MODEL IS LEARNING OVER TIME. THIS HELPS TO UNDERSTAND IF THE MODEL IS IMPROVING OR IF THERE ARE AREAS THAT NEED ADJUSTMENT.



# ANALYZING MODEL PREDICTIONS AND SIMILARITY

ASSUMING THE MODEL IS TRAINED, THE CODE AIMS TO UNDERSTAND HOW SIMILAR OR DISSIMILAR THE PREDICTIONS ARE FOR THE TEST DATA. IT USES A CONCEPT CALLED "COSINE SIMILARITY," A MATHEMATICAL MEASURE TO ASSESS THE SIMILARITY BETWEEN TWO NON-ZERO VECTORS.



## INTERACTIVE PREDICTION

AN INTERACTIVE FUNCTION IS CREATED TO ALLOW USERS TO INPUT NEW TEXT AND SEE THE MODEL'S PREDICTION FOR THE TOPIC OF THAT TEXT.

Enter Text:

1/1 [=====] - ETA: 0s  
1/1 [=====] - 0s 21ms/step  
Predicted Topic: earn

**THATS' ALL!**