

AI 2024 Online Summer Internship

Titanic Project

Lecture 1

Treating Titanic Passenger Survival
Prediction Problem as a Machine
Learning Problem using Train-Test
Split Approach



Lecture Outline

- **Titanic Passenger Survival Prediction Problem**
- **Steps – Treating Titanic Passenger Survival Prediction Problem as a Machine Learning Problem using Train-Test Split Approach**



Titanic Passenger Survival Prediction Problem

SLIDE

Royal Mail Ship (RMS) Titanic – Brief Overview

- RMS Titanic was a **British passenger liner** operated by the White Star Line that sank in the North Atlantic Ocean in the early morning hours of **15 April 1912**, after striking an **iceberg** during her maiden voyage from Southampton to New York City
- Of the estimated **2,224 passengers** and **crew** aboard, **more than 1,500 died**, making the sinking one of modern history's **deadliest** peacetime commercial marine disasters
- RMS Titanic was the **largest ship** afloat at the time she entered service¹

SLIDE

RMS Titanic – Main Features

- **Name of Ship**
 - Royal Mail Ship (RMS) Titanic
- **Manufactured**
 - May 31, 1911
- **Size**
 - Length = 882 feet 9 inches (269.06 m)
 - Width = 92.5 feet (28.2 m)
 - Height = 104 feet (32 m)
- **First Journey**
 - April 10, 1912
- **Source and Destination**
 - Southampton, England to New York, USA
- **Total Passengers On Board**
 - 2,208
- **Passengers Survived**
 - 705
- **Passengers Died**
 - 1,503

¹

Last visited: 07-09-2020

SLIDE

RMS Titanic



Figure 01: Royal Mail Ship (RMS) Titanic []

SLIDE

Lecture Focus

- The **main focus** of this Lecture is **developing a Predictive System** which can automatically predict whether a Passenger in RMS Titanic **Survived or Not**

SLIDE

Titanic Passenger Survival Prediction System

- **Real-world World**
 - **Titanic Passenger Survival**
- **Treated as**
 - **Supervised Machine Learning Problem**
- **Note**
 - **Titanic Passenger Survival Prediction Problem is treated as a**

- **Binary Classification Problem** because the
 - The main aim is to distinguish between Two Classes
 - Class 01 = Survived
 - Class 02 = Not Survived
- Goal
 - Learn an Input-Output Function
 - i.e. Learn from Input to predict the Output

SLIDE

Titanic Passenger Survival Prediction System – Task

- Given
 - A Passenger (Represented as Set of Attributes)
- Task
 - Automatically predict whether the Passenger Survived or Not

SLIDE

Titanic Passenger Survival Prediction System – Input and Output

- Input
 - A Passenger
- Output
 - Survived / Not Survived

SLIDE

Note

- In Kaggle Titanic Dataset, a Passenger is represented with many Attributes
- Kaggle Titanic Dataset
 - URL: <http://kaggle.com/c/titanic>
- For simplicity and to explain things more clearly
 - In this Lecture, we have represented a Passenger with Four Attributes

SLIDE

Titanic Passenger Survival Prediction System – Input Attributes

- In this Lecture, a Passenger is represented with the following Four Attributes
- Attribute 01 – PClass



- Possible Value 01 = First
- Possible Value 02 = Second
- Possible Value 03 = Third
- Attribute 02 – Gender
 - Possible Value 01 = Female
 - Possible Value 02 = Male
- Attribute 03 – Sibling
 - Possible Value 01 = Zero
 - Possible Value 02 = One
 - Possible Value 03 = Two
 - Possible Value 04 = Three
- Attribute 04 – Embarked
 - Possible Value 01 = Cherbourg
 - Possible Value 02 = Southampton
 - Possible Value 03 = Queenstown

SLIDE

Titanic Passenger Survival Prediction System – Output Attributes

- In Titanic Dataset, there is One Output Attribute
 - Attribute 01 – Survived
 - Possible Value 01 = Yes
 - Possible Value 02 = No

SLIDE

Titanic Passenger Survival Prediction System – Summary (Input and Output)

- The following Table summarizes the Input and Output Attributes for Titanic Dataset

Attribute No.	Attribute Names	Possible Values	Data Types
1	PClass	First, Second, Third	Categorical
2	Gender	Male, Female	Categorical



3	Sibling	Zero, One, Two, Three	Categorical
4	Embarked	Southampton, Queenstown, Cherbourg	Categorical
5	Survived	Yes, No	Categorical

Table 01: Attributes of Dataset



Steps – Treating Titanic Passenger Survival Prediction Problem as a Supervised Machine Learning Problem using Train-Test Split Approach

SLIDE

Titanic Passenger Survival Prediction Problem

- Task
 - Develop a Titanic Passenger Survival Prediction System to Predict the Survival of a Passenger
- Input
 - Four Attributes

1. PClass
2. Gender
3. Sibling
4. Embarked

- Output
 - One Attribute

1. Survived

- Treated as a
 - Supervised Machine Learning Problem
- Goal
 - Learn an Input-Output Function
 - i.e. Learn from Input to predict the Output

SLIDE

Titanic Passenger Survival Prediction System is a Classification Problem

- Titanic Passenger Survival Prediction System is a Classification Problem because
 - Output is Categorical



SLIDE

Titanic Passenger Survival Prediction Problem – **Input** and **Output**

- **Input**
 - **Categorical**
- **Output**
 - **Categorical**

SLIDE

Project Focus

Titanic Passenger Survival Prediction System

SLIDE

Steps – Treating Titanic Passenger Survival Prediction Problem as a **Classification Problem**

- In Sha Allah (انشاء الله), I will follow the following steps to **treat** the Titanic Passenger Survival Prediction Problem as a **Classification Problem**
 - Step 1: Decide the **Learning Settings**
 - Step 2: **Obtain** Sample Data
 - Step 3: **Understand and Pre-process** Sample Data
 - Step 4: **Represent** Sample Data in **Machine Understandable Format**
 - Step 5: Select **Suitable** Machine Learning Algorithms
 - Step 6: **Split** Sample Data into **Training Data** and **Testing Data**
 - Step 7: Select **Suitable** Evaluation Measure(s)
 - Step 8: Execute First Two Phases of Machine Learning Cycle
 - Training Phase
 - Testing Phase
 - Step 9: **Analyze** Results

If (Results are **Good**)
Then
Move to the Next Step
Else
Go to Step 1



- Step 10: Execute 3rd and 4th Phases of Machine Learning Cycle
 - Application Phase
 - Feedback Phase
- Step 11: Based on Feedback
 - Go to Step 1 and Repeat all the Steps

Step 1: Decide the Learning Setting

SLIDE

Step 1: Decide the Learning Setting

- In Sha Allah (إِن شاء اللَّهُ), I will treat the Titanic Passenger Survival Prediction Problem as a
 - Supervised Machine Learning Problem
- Since Output is Categorical, it will be treated as a
 - Classification Problem

Step 2: Obtain Sample Data

SLIDE

Step 2: Obtain Sample Data

- Since I am Treating Titanic Passenger Survival Prediction Problem as a Supervised Machine Learning Problem, I will need
 - Annotated Data
- For more accurate learning, I need
 1. Large amount of Annotated Data
 2. High-quality Annotated Data
 3. Balanced Data
- Note
 - For simplicity, In Sha Allah (إِن شاء اللَّهُ) I will use a toy Corpus / Dataset of 100 instances

SLIDE

Step 2: Obtain Sample Data Cont...



- Two Main Choices to Obtain Data
 1. Use an Existing Corpus
 2. Develop Your Corpus
- The Dataset to use is a subset of Kaggle Titanic Dataset
 - Corpus / Dataset
 - Dataset Link:
 - Paper Link:
 -
 - Paper Reference:
 - Cicoria, S., Sherlock, J., Muniswamaiah, M., and Clarke, L., 2014, Classification of titanic passenger data and chances of surviving the disaster, In Proceedings of Student-Faculty Research Day, CSIS (pp. 1-6)

SLIDE

Obtain Sample Data Cont...

- Total Instances in Sample Data = 100
 - Survived = 50
 - Not Survived = 50

SLIDE

Sample Data

- We obtained a Sample Data of 100 instances
 - See sample-data.csv File in Supporting Material
- The following Table shows the Sample Data

Instance No.	Input					Output Survived
	PClass	Gender	Sibling	Embarked		
x ₁	Third	Male	One	Southampton		No
x ₂	Second	Female	Zero	Southampton		Yes
x ₃	Third	Male	Zero	Southampton		No
x ₄	Third	Female	Three	Southampton		Yes
x ₅	Third	Male	Zero	Queenstown		No
x ₆	First	Female	Three	Southampton		Yes
x ₇	Third	Male	Zero	Southampton		No



X8	Third	Male	Zero	Southampton	Yes
X9	First	Male	Zero	Southampton	No
X10	Second	Male	Zero	Southampton	Yes
X11	Third	Male	One	Queenstown	No
X12	First	Male	Zero	Cherbourg	Yes
X13	First	Male	Zero	Cherbourg	No
X14	Second	Female	Zero	Southampton	Yes
X15	Second	Male	Zero	Southampton	No
X16	Third	Male	One	Cherbourg	Yes
X17	Third	Male	Two	Cherbourg	No
X18	Third	Male	Zero	Southampton	Yes
X19	Third	Male	Zero	Southampton	No
X20	Third	Female	One	Cherbourg	Yes
X21	Third	Male	Zero	Cherbourg	No
X22	First	Female	Zero	Southampton	Yes
X23	First	Male	Zero	Cherbourg	No
X24	Second	Female	Zero	Southampton	Yes
X25	Third	Female	One	Southampton	No
X26	Third	Female	Zero	Southampton	Yes
X27	Third	Male	Zero	Southampton	No
X28	Third	Male	Zero	Southampton	Yes
X29	Third	Male	Zero	Southampton	No
X30	Third	Female	One	Queenstown	Yes
X31	Third	Female	Two	Southampton	No
X32	First	Female	Zero	Cherbourg	Yes
X33	Third	Male	Two	Southampton	No
X34	First	Female	Zero	Southampton	Yes
X35	First	Male	One	Cherbourg	No
X36	First	Female	Zero	Southampton	Yes
X37	First	Male	One	Southampton	No
X38	Third	Female	One	Queenstown	Yes
X39	Third	Female	One	Southampton	No
X40	Second	Female	One	Southampton	Yes
X41	Second	Female	One	Southampton	No
X42	Third	Female	Zero	Queenstown	Yes

X43	Third	Male	Zero	Cherbourg	No
X44	First	Female	Zero	Cherbourg	Yes
X45	Third	Male	Three	Southampton	No
X46	Third	Female	One	Southampton	Yes
X47	Third	Male	Zero	Southampton	No
X48	Third	Male	Zero	Southampton	Yes
X49	Third	Male	One	Southampton	No
X50	Second	Female	Zero	Southampton	Yes
X51	First	Male	Three	Southampton	No
X52	Third	Male	Zero	Southampton	Yes
X53	Third	Female	One	Southampton	No
X54	Second	Female	Zero	Southampton	Yes
X55	Third	Female	Three	Southampton	No
X56	Third	Female	Zero	Southampton	Yes
X57	Third	Female	Zero	Southampton	No
X58	Third	Female	One	Southampton	Yes
X59	Third	Male	Zero	Cherbourg	No
X60	Third	Male	Zero	Southampton	Yes
X61	Second	Male	Zero	Southampton	No
X62	First	Female	One	Southampton	Yes
X63	Third	Male	Three	Queenstown	No
X64	Third	Female	Zero	Queenstown	Yes
X65	Second	Female	Zero	Southampton	No
X66	Third	Female	Zero	Cherbourg	Yes
X67	Third	Male	Zero	Southampton	No
X68	Second	Male	Zero	Southampton	Yes
X69	Third	Male	Zero	Southampton	No
X70	First	Male	Zero	Southampton	Yes
X71	Third	Male	Zero	Southampton	No
X72	Second	Male	Two	Southampton	Yes
X73	Third	Male	Zero	Cherbourg	No
X74	Third	Female	Zero	Southampton	Yes
X75	Third	Female	Zero	Southampton	No
X76	First	Female	One	Southampton	Yes
X77	Third	Male	One	Southampton	No

X78	Second	Female	Zero	Southampton	Yes
X79	Third	Male	Zero	Southampton	No
X80	Third	Female	Two	Southampton	Yes
X81	Third	Male	Two	Southampton	No
X82	Third	Female	Zero	Southampton	Yes
X83	Third	Male	Zero	Southampton	No
X84	First	Male	One	Southampton	Yes
X85	Second	Male	Zero	Southampton	No
X86	Third	Female	One	Queenstown	Yes
X87	Second	Male	Zero	Southampton	No
X88	First	Female	One	Southampton	Yes
X89	Third	Male	Zero	Southampton	No
X90	Second	Female	Zero	Southampton	Yes
X91	Second	Male	Zero	Southampton	No
X92	Second	Female	Zero	Southampton	Yes
X93	Third	Male	Zero	Cherbourg	No
X94	First	Male	One	Southampton	Yes
X95	First	Male	One	Southampton	No
X96	Third	Male	One	Cherbourg	Yes
X97	Third	Male	One	Queenstown	No
X98	Third	Male	One	Southampton	Yes
X99	Second	Male	One	Southampton	No
X100	Third	Male	One	Southampton	Yes

Step 03: Understand and Pre-process Data

SLIDE

Step 3: Understand and Pre-process Sample Data

- Understanding Data
 - The **Sample Data** contains **Five Attributes**
 - PClass
 - Gender
 - Sibling
 - Embarked



- Survived
- Separating Input from Output
 - Input comprises of Four Attributes
 - PClass
 - Gender
 - Sibling
 - Embarked
 - The Output comprises of a Single Attribute
 - Survived
- Pre-processing Data
 - Corpus is already pre-processed
 - Therefore, no pre-processing is needed 😊

Step 04: Represent Data in Machine Understandable Format

SLIDE

Step 4: Represent Sample Data in Machine Understandable Format

- Feature-based Classification Algorithms (implemented in Scikit-learn) can understand data in
 - Attribute-Value Pair
 - Values of Attributes / Features must be Numeric
- Problem
 - Our Sample Data is not in Attribute-Value Pair form
 - We need to transform our Sample Data into Machine Understandable Format
- Solution
 - There are many approaches to transform Sample Data into Machine Understandable Format

SLIDE

Feature Extraction

- Features are already extracted
 - Therefore, we will skip the Feature Extraction Step 😊



SLIDE

Important Note

- In this Lecture, we are using **Scikit-learn** implementation of the **Support Vector Classifier** Machine Learning Algorithm
- **Scikit-learn can only understand Data in Numerical Representation**
 - Therefore, we will need to **Convert the Categorical Values to Numerical Values**

SLIDE

Transforming Sample Data in Machine Understandable Format

- In our Sample Data
 - **Input is Categorical**
 - **Output is Categorical**
- Considering **Input** (PClass, Gender, Sibling, Embarked) and **Output** (Survived), we will need to
 - **Transform Input (Categorical) into Numerical Representation**
 - **Transform Output (Categorical) into Numerical Representation**

SLIDE

Converting Output into Numerical Representation

- A Two-Step Process
 - Step 01: Define an **Encoding Scheme**
 - Step 02: Use Encoding Scheme defined in Step 01, to **convert Categorical Output Values to Numerical Output Values for all instances** in the Sample Data

SLIDE

Converting Output into Numerical Representation Cont...

- Step 01: Define an Encoding Scheme
- **Encoding Scheme for Survived Attribute**
 - No = 0
 - Yes = 1

SLIDE

Converting Output into Numerical Representation Cont...



- Step 02: Use Encoding Scheme defined in Step 01, to convert Categorical Output Values to Numerical Output Values for all instances in the Sample Data
- The Table below shows Sample Data after Encoding Categorical Output Values to Numerical Output Values
 - See **sample-data-encoded-output.csv** File in Supporting Material

Instance No.	Input				Output
	PClass	Gender	Sibling	Embarked	
X ₁	Third	Male	One	Southampton	0
X ₂	Second	Female	Zero	Southampton	1
X ₃	Third	Male	Zero	Southampton	0
X ₄	Third	Female	Three	Southampton	1
X ₅	Third	Male	Zero	Queenstown	0
X ₆	First	Female	Three	Southampton	1
X ₇	Third	Male	Zero	Southampton	0
X ₈	Third	Male	Zero	Southampton	1
X ₉	First	Male	Zero	Southampton	0
X ₁₀	Second	Male	Zero	Southampton	1
X ₁₁	Third	Male	One	Queenstown	0
X ₁₂	First	Male	Zero	Cherbourg	1
X ₁₃	First	Male	Zero	Cherbourg	0
X ₁₄	Second	Female	Zero	Southampton	1
X ₁₅	Second	Male	Zero	Southampton	0
X ₁₆	Third	Male	One	Cherbourg	1
X ₁₇	Third	Male	Two	Cherbourg	0
X ₁₈	Third	Male	Zero	Southampton	1
X ₁₉	Third	Male	Zero	Southampton	0
X ₂₀	Third	Female	One	Cherbourg	1
X ₂₁	Third	Male	Zero	Cherbourg	0
X ₂₂	First	Female	Zero	Southampton	1
X ₂₃	First	Male	Zero	Cherbourg	0
X ₂₄	Second	Female	Zero	Southampton	1



X25	Third	Female	One	Southampton	0
X26	Third	Female	Zero	Southampton	1
X27	Third	Male	Zero	Southampton	0
X28	Third	Male	Zero	Southampton	1
X29	Third	Male	Zero	Southampton	0
X30	Third	Female	One	Queenstown	1
X31	Third	Female	Two	Southampton	0
X32	First	Female	Zero	Cherbourg	1
X33	Third	Male	Two	Southampton	0
X34	First	Female	Zero	Southampton	1
X35	First	Male	One	Cherbourg	0
X36	First	Female	Zero	Southampton	1
X37	First	Male	One	Southampton	0
X38	Third	Female	One	Queenstown	1
X39	Third	Female	One	Southampton	0
X40	Second	Female	One	Southampton	1
X41	Second	Female	One	Southampton	0
X42	Third	Female	Zero	Queenstown	1
X43	Third	Male	Zero	Cherbourg	0
X44	First	Female	Zero	Cherbourg	1
X45	Third	Male	Three	Southampton	0
X46	Third	Female	One	Southampton	1
X47	Third	Male	Zero	Southampton	0
X48	Third	Male	Zero	Southampton	1
X49	Third	Male	One	Southampton	0
X50	Second	Female	Zero	Southampton	1
X51	First	Male	Three	Southampton	0
X52	Third	Male	Zero	Southampton	1
X53	Third	Female	One	Southampton	0
X54	Second	Female	Zero	Southampton	1
X55	Third	Female	Three	Southampton	0
X56	Third	Female	Zero	Southampton	1
X57	Third	Female	Zero	Southampton	0
X58	Third	Female	One	Southampton	1
X59	Third	Male	Zero	Cherbourg	0



X ₆₀	Third	Male	Zero	Southampton	1
X ₆₁	Second	Male	Zero	Southampton	0
X ₆₂	First	Female	One	Southampton	1
X ₆₃	Third	Male	Three	Queenstown	0
X ₆₄	Third	Female	Zero	Queenstown	1
X ₆₅	Second	Female	Zero	Southampton	0
X ₆₆	Third	Female	Zero	Cherbourg	1
X ₆₇	Third	Male	Zero	Southampton	0
X ₆₈	Second	Male	Zero	Southampton	1
X ₆₉	Third	Male	Zero	Southampton	0
X ₇₀	First	Male	Zero	Southampton	1
X ₇₁	Third	Male	Zero	Southampton	0
X ₇₂	Second	Male	Two	Southampton	1
X ₇₃	Third	Male	Zero	Cherbourg	0
X ₇₄	Third	Female	Zero	Southampton	1
X ₇₅	Third	Female	Zero	Southampton	0
X ₇₆	First	Female	One	Southampton	1
X ₇₇	Third	Male	One	Southampton	0
X ₇₈	Second	Female	Zero	Southampton	1
X ₇₉	Third	Male	Zero	Southampton	0
X ₈₀	Third	Female	Two	Southampton	1
X ₈₁	Third	Male	Two	Southampton	0
X ₈₂	Third	Female	Zero	Southampton	1
X ₈₃	Third	Male	Zero	Southampton	0
X ₈₄	First	Male	One	Southampton	1
X ₈₅	Second	Male	Zero	Southampton	0
X ₈₆	Third	Female	One	Queenstown	1
X ₈₇	Second	Male	Zero	Southampton	0
X ₈₈	First	Female	One	Southampton	1
X ₈₉	Third	Male	Zero	Southampton	0
X ₉₀	Second	Female	Zero	Southampton	1
X ₉₁	Second	Male	Zero	Southampton	0
X ₉₂	Second	Female	Zero	Southampton	1
X ₉₃	Third	Male	Zero	Cherbourg	0
X ₉₄	First	Male	One	Southampton	1



X95	First	Male	One	Southampton	0
X96	Third	Male	One	Cherbourg	1
X97	Third	Male	One	Queenstown	0
X98	Third	Male	One	Southampton	1
X99	Second	Male	One	Southampton	0
X100	Third	Male	One	Southampton	1

SLIDE**Note**

- Alhamdulillah (الحمد لله), Output is **transformed** into Numerical Representation
- In Sha Allah (إذ شاء الله), in the next Slides, I will try to explain how to transform Input into Numerical Representation

SLIDE**Converting Input into Numerical Representation**

- Step 01: Define an **Encoding Scheme**
- Step 02: Use Encoding Scheme defined in Step 01, to **convert Categorical Input Values to Numerical Input Values for all instances** in the Sample Data

SLIDE**Converting Input into Numerical Representation Cont...**

- Step 01: Define an Encoding Scheme
- **Encoding Scheme for PClass Attribute**
 - First = 0
 - Second = 1
 - Third = 2
- **Encoding Scheme for Gender Attribute**
 - Female = 0
 - Male = 1
- **Encoding Scheme for Sibling Attribute**
 - One = 0
 - Three = 1
 - Two = 2
 - Zero = 3



- **Encoding Scheme for Embarked Attribute**

- Cherbourg = 0
- Queenstown = 1
- Southampton = 2

SLIDE

Converting Input into Numerical Representation Cont...

- Step 02: Use Encoding Scheme defined in Step 01, to convert Categorical Input Values to Numerical Input Values for all instances in the Sample Data
- The Table below shows Sample Data after Encoding Categorical Input Values to Numerical Input Values
 - See **sample-data-encoded.csv** File in Supporting Material

Instance No.	Input				Output
	PClass	Gender	Sibling	Embarked	
X1	2	1	0	2	0
X2	1	0	3	2	1
X3	2	1	3	2	0
X4	2	0	1	2	1
X5	2	1	3	1	0
X6	0	0	1	2	1
X7	2	1	3	2	0
X8	2	1	3	2	1
X9	0	1	3	2	0
X10	1	1	3	2	1
X11	2	1	0	1	0
X12	0	1	3	0	1
X13	0	1	3	0	0
X14	1	0	3	2	1
X15	1	1	3	2	0
X16	2	1	0	0	1
X17	2	1	2	0	0
X18	2	1	3	2	1



X19	2	1	3	2	0
X20	2	0	0	0	1
X21	2	1	3	0	0
X22	0	0	3	2	1
X23	0	1	3	0	0
X24	1	0	3	2	1
X25	2	0	0	2	0
X26	2	0	3	2	1
X27	2	1	3	2	0
X28	2	1	3	2	1
X29	2	1	3	2	0
X30	2	0	0	1	1
X31	2	0	2	2	0
X32	0	0	3	0	1
X33	2	1	2	2	0
X34	0	0	3	2	1
X35	0	1	0	0	0
X36	0	0	3	2	1
X37	0	1	0	2	0
X38	2	0	0	1	1
X39	2	0	0	2	0
X40	1	0	0	2	1
X41	1	0	0	2	0
X42	2	0	3	1	1
X43	2	1	3	0	0
X44	0	0	3	0	1
X45	2	1	1	2	0
X46	2	0	0	2	1
X47	2	1	3	2	0
X48	2	1	3	2	1
X49	2	1	0	2	0
X50	1	0	3	2	1
X51	0	1	1	2	0
X52	2	1	3	2	1
X53	2	0	0	2	0

X54	1	0	3	2	1
X55	2	0	1	2	0
X56	2	0	3	2	1
X57	2	0	3	2	0
X58	2	0	0	2	1
X59	2	1	3	0	0
X60	2	1	3	2	1
X61	1	1	3	2	0
X62	0	0	0	2	1
X63	2	1	1	1	0
X64	2	0	3	1	1
X65	1	0	3	2	0
X66	2	0	3	0	1
X67	2	1	3	2	0
X68	1	1	3	2	1
X69	2	1	3	2	0
X70	0	1	3	2	1
X71	2	1	3	2	0
X72	1	1	2	2	1
X73	2	1	3	0	0
X74	2	0	3	2	1
X75	2	0	3	2	0
X76	0	0	0	2	1
X77	2	1	0	2	0
X78	1	0	3	2	1
X79	2	1	3	2	0
X80	2	0	2	2	1
X81	2	1	2	2	0
X82	2	0	3	2	1
X83	2	1	3	2	0
X84	0	1	0	2	1
X85	1	1	3	2	0
X86	2	0	0	1	1
X87	1	1	3	2	0
X88	0	0	0	2	1

X89	2	1	3	2	0
X90	1	0	3	2	1
X91	1	1	3	2	0
X92	1	0	3	2	1
X93	2	1	3	0	0
X94	0	1	0	2	1
X95	0	1	0	2	0
X96	2	1	0	0	1
X97	2	1	0	1	0
X98	2	1	0	2	1
X99	1	1	0	2	0
X100	2	1	0	2	1

SLIDE

Hooooooooooooorrrrrrrrrrraaaaaaaayyyyyyyyyy! 

- Alhamdulillah (الحمد لله), both **Input** and **Output** are **transformed** into **Numerical Representation**

SLIDE

Recap – **Original Sample Data**

Instance No.	Input				Output
	PClass	Gender	Sibling	Embarked	
X1	Third	Male	One	Southampton	No
X2	Second	Female	Zero	Southampton	Yes
X3	Third	Male	Zero	Southampton	No
X4	Third	Female	Three	Southampton	Yes
X5	Third	Male	Zero	Queenstown	No
X6	First	Female	Three	Southampton	Yes
X7	Third	Male	Zero	Southampton	No
X8	Third	Male	Zero	Southampton	Yes
X9	First	Male	Zero	Southampton	No
X10	Second	Male	Zero	Southampton	Yes



X11	Third	Male	One	Queenstown	No
X12	First	Male	Zero	Cherbourg	Yes
X13	First	Male	Zero	Cherbourg	No
X14	Second	Female	Zero	Southampton	Yes
X15	Second	Male	Zero	Southampton	No
X16	Third	Male	One	Cherbourg	Yes
X17	Third	Male	Two	Cherbourg	No
X18	Third	Male	Zero	Southampton	Yes
X19	Third	Male	Zero	Southampton	No
X20	Third	Female	One	Cherbourg	Yes
X21	Third	Male	Zero	Cherbourg	No
X22	First	Female	Zero	Southampton	Yes
X23	First	Male	Zero	Cherbourg	No
X24	Second	Female	Zero	Southampton	Yes
X25	Third	Female	One	Southampton	No
X26	Third	Female	Zero	Southampton	Yes
X27	Third	Male	Zero	Southampton	No
X28	Third	Male	Zero	Southampton	Yes
X29	Third	Male	Zero	Southampton	No
X30	Third	Female	One	Queenstown	Yes
X31	Third	Female	Two	Southampton	No
X32	First	Female	Zero	Cherbourg	Yes
X33	Third	Male	Two	Southampton	No
X34	First	Female	Zero	Southampton	Yes
X35	First	Male	One	Cherbourg	No
X36	First	Female	Zero	Southampton	Yes
X37	First	Male	One	Southampton	No
X38	Third	Female	One	Queenstown	Yes
X39	Third	Female	One	Southampton	No
X40	Second	Female	One	Southampton	Yes
X41	Second	Female	One	Southampton	No
X42	Third	Female	Zero	Queenstown	Yes
X43	Third	Male	Zero	Cherbourg	No
X44	First	Female	Zero	Cherbourg	Yes
X45	Third	Male	Three	Southampton	No

X46	Third	Female	One	Southampton	Yes
X47	Third	Male	Zero	Southampton	No
X48	Third	Male	Zero	Southampton	Yes
X49	Third	Male	One	Southampton	No
X50	Second	Female	Zero	Southampton	Yes
X51	First	Male	Three	Southampton	No
X52	Third	Male	Zero	Southampton	Yes
X53	Third	Female	One	Southampton	No
X54	Second	Female	Zero	Southampton	Yes
X55	Third	Female	Three	Southampton	No
X56	Third	Female	Zero	Southampton	Yes
X57	Third	Female	Zero	Southampton	No
X58	Third	Female	One	Southampton	Yes
X59	Third	Male	Zero	Cherbourg	No
X60	Third	Male	Zero	Southampton	Yes
X61	Second	Male	Zero	Southampton	No
X62	First	Female	One	Southampton	Yes
X63	Third	Male	Three	Queenstown	No
X64	Third	Female	Zero	Queenstown	Yes
X65	Second	Female	Zero	Southampton	No
X66	Third	Female	Zero	Cherbourg	Yes
X67	Third	Male	Zero	Southampton	No
X68	Second	Male	Zero	Southampton	Yes
X69	Third	Male	Zero	Southampton	No
X70	First	Male	Zero	Southampton	Yes
X71	Third	Male	Zero	Southampton	No
X72	Second	Male	Two	Southampton	Yes
X73	Third	Male	Zero	Cherbourg	No
X74	Third	Female	Zero	Southampton	Yes
X75	Third	Female	Zero	Southampton	No
X76	First	Female	One	Southampton	Yes
X77	Third	Male	One	Southampton	No
X78	Second	Female	Zero	Southampton	Yes
X79	Third	Male	Zero	Southampton	No
X80	Third	Female	Two	Southampton	Yes

X ₈₁	Third	Male	Two	Southampton	No
X ₈₂	Third	Female	Zero	Southampton	Yes
X ₈₃	Third	Male	Zero	Southampton	No
X ₈₄	First	Male	One	Southampton	Yes
X ₈₅	Second	Male	Zero	Southampton	No
X ₈₆	Third	Female	One	Queenstown	Yes
X ₈₇	Second	Male	Zero	Southampton	No
X ₈₈	First	Female	One	Southampton	Yes
X ₈₉	Third	Male	Zero	Southampton	No
X ₉₀	Second	Female	Zero	Southampton	Yes
X ₉₁	Second	Male	Zero	Southampton	No
X ₉₂	Second	Female	Zero	Southampton	Yes
X ₉₃	Third	Male	Zero	Cherbourg	No
X ₉₄	First	Male	One	Southampton	Yes
X ₉₅	First	Male	One	Southampton	No
X ₉₆	Third	Male	One	Cherbourg	Yes
X ₉₇	Third	Male	One	Queenstown	No
X ₉₈	Third	Male	One	Southampton	Yes
X ₉₉	Second	Male	One	Southampton	No
X ₁₀₀	Third	Male	One	Southampton	Yes

SLIDE**Recap - Sample Data in Numerical Representation**

Instance No.	Input				Output
	PClass	Gender	Sibling	Embarked	
X ₁	2	1	0	2	0
X ₂	1	0	3	2	1
X ₃	2	1	3	2	0
X ₄	2	0	1	2	1
X ₅	2	1	3	1	0
X ₆	0	0	1	2	1
X ₇	2	1	3	2	0



X8	2	1	3	2	1
X9	0	1	3	2	0
X10	1	1	3	2	1
X11	2	1	0	1	0
X12	0	1	3	0	1
X13	0	1	3	0	0
X14	1	0	3	2	1
X15	1	1	3	2	0
X16	2	1	0	0	1
X17	2	1	2	0	0
X18	2	1	3	2	1
X19	2	1	3	2	0
X20	2	0	0	0	1
X21	2	1	3	0	0
X22	0	0	3	2	1
X23	0	1	3	0	0
X24	1	0	3	2	1
X25	2	0	0	2	0
X26	2	0	3	2	1
X27	2	1	3	2	0
X28	2	1	3	2	1
X29	2	1	3	2	0
X30	2	0	0	1	1
X31	2	0	2	2	0
X32	0	0	3	0	1
X33	2	1	2	2	0
X34	0	0	3	2	1
X35	0	1	0	0	0
X36	0	0	3	2	1
X37	0	1	0	2	0
X38	2	0	0	1	1
X39	2	0	0	2	0
X40	1	0	0	2	1
X41	1	0	0	2	0
X42	2	0	3	1	1

X43	2	1	3	0	0
X44	0	0	3	0	1
X45	2	1	1	2	0
X46	2	0	0	2	1
X47	2	1	3	2	0
X48	2	1	3	2	1
X49	2	1	0	2	0
X50	1	0	3	2	1
X51	0	1	1	2	0
X52	2	1	3	2	1
X53	2	0	0	2	0
X54	1	0	3	2	1
X55	2	0	1	2	0
X56	2	0	3	2	1
X57	2	0	3	2	0
X58	2	0	0	2	1
X59	2	1	3	0	0
X60	2	1	3	2	1
X61	1	1	3	2	0
X62	0	0	0	2	1
X63	2	1	1	1	0
X64	2	0	3	1	1
X65	1	0	3	2	0
X66	2	0	3	0	1
X67	2	1	3	2	0
X68	1	1	3	2	1
X69	2	1	3	2	0
X70	0	1	3	2	1
X71	2	1	3	2	0
X72	1	1	2	2	1
X73	2	1	3	0	0
X74	2	0	3	2	1
X75	2	0	3	2	0
X76	0	0	0	2	1
X77	2	1	0	2	0

X78	1	0	3	2	1
X79	2	1	3	2	0
X80	2	0	2	2	1
X81	2	1	2	2	0
X82	2	0	3	2	1
X83	2	1	3	2	0
X84	0	1	0	2	1
X85	1	1	3	2	0
X86	2	0	0	1	1
X87	1	1	3	2	0
X88	0	0	0	2	1
X89	2	1	3	2	0
X90	1	0	3	2	1
X91	1	1	3	2	0
X92	1	0	3	2	1
X93	2	1	3	0	0
X94	0	1	0	2	1
X95	0	1	0	2	0
X96	2	1	0	0	1
X97	2	1	0	1	0
X98	2	1	0	2	1
X99	1	1	0	2	0
X100	2	1	0	2	1

Step 05: Select Suitable Machine Learning Algorithms

SLIDE

Step 05: Select Suitable Machine Learning Algorithms

- Previous students have shown that Good Starting Points for Classification Problems are
 - Random Forest Classifier
 - Support Vector Classifier
 - Naïve Bayes
 - Gradient Boosting Classifier



SLIDE**Lecture Focus**

- In Sha Allah, in this Lecture, we will use

Support Vector Classifier

Step 06: Split Sample Data into Training Data and Testing Data

SLIDE**Step 6: Split Sample Data into Training and Testing**

- We **Split** the Sample Data using
 - **Train-Test Split Ratio** of
 - 80% - 20%
- **Training Data**
 - Total Instances = 80
 - Survived = 40
 - Not Survived = 40
- **Testing Data**
 - Total Instances = 20
 - Survived = 10
 - Not Survived = 10

SLIDE**Training Data**

- The following Table shows the **Training Data**
 - See **training-data-encoded.csv** File in Supporting Material

Instance No.	Input				Output Survived
	PClass	Gender	Sibling	Embarked	
x1	2	1	0	2	0



X ₂	1	0	3	2	1
X ₃	2	1	3	2	0
X ₄	2	0	1	2	1
X ₅	2	1	3	1	0
X ₆	0	0	1	2	1
X ₇	2	1	3	2	0
X ₈	2	1	3	2	1
X ₉	0	1	3	2	0
X ₁₀	1	1	3	2	1
X ₁₁	2	1	0	1	0
X ₁₂	0	1	3	0	1
X ₁₃	0	1	3	0	0
X ₁₄	1	0	3	2	1
X ₁₅	1	1	3	2	0
X ₁₆	2	1	0	0	1
X ₁₇	2	1	2	0	0
X ₁₈	2	1	3	2	1
X ₁₉	2	1	3	2	0
X ₂₀	2	0	0	0	1
X ₂₁	2	1	3	0	0
X ₂₂	0	0	3	2	1
X ₂₃	0	1	3	0	0
X ₂₄	1	0	3	2	1
X ₂₅	2	0	0	2	0
X ₂₆	2	0	3	2	1
X ₂₇	2	1	3	2	0
X ₂₈	2	1	3	2	1
X ₂₉	2	1	3	2	0
X ₃₀	2	0	0	1	1
X ₃₁	2	0	2	2	0
X ₃₂	0	0	3	0	1
X ₃₃	2	1	2	2	0
X ₃₄	0	0	3	2	1
X ₃₅	0	1	0	0	0
X ₃₆	0	0	3	2	1

X37	0	1	0	2	0
X38	2	0	0	1	1
X39	2	0	0	2	0
X40	1	0	0	2	1
X41	1	0	0	2	0
X42	2	0	3	1	1
X43	2	1	3	0	0
X44	0	0	3	0	1
X45	2	1	1	2	0
X46	2	0	0	2	1
X47	2	1	3	2	0
X48	2	1	3	2	1
X49	2	1	0	2	0
X50	1	0	3	2	1
X51	0	1	1	2	0
X52	2	1	3	2	1
X53	2	0	0	2	0
X54	1	0	3	2	1
X55	2	0	1	2	0
X56	2	0	3	2	1
X57	2	0	3	2	0
X58	2	0	0	2	1
X59	2	1	3	0	0
X60	2	1	3	2	1
X61	1	1	3	2	0
X62	0	0	0	2	1
X63	2	1	1	1	0
X64	2	0	3	1	1
X65	1	0	3	2	0
X66	2	0	3	0	1
X67	2	1	3	2	0
X68	1	1	3	2	1
X69	2	1	3	2	0
X70	0	1	3	2	1
X71	2	1	3	2	0



X ₇₂	1	1	2	2	1
X ₇₃	2	1	3	0	0
X ₇₄	2	0	3	2	1
X ₇₅	2	0	3	2	0
X ₇₆	0	0	0	2	1
X ₇₇	2	1	0	2	0
X ₇₈	1	0	3	2	1
X ₇₉	2	1	3	2	0
X ₈₀	2	0	2	2	1

SLIDE**Testing Data**

- The following Table shows the **Testing Data**
 - See **testing-data-encoded.csv** File in Supporting Material

Instance No.	Input				Output Survived
	PClass	Gender	Sibling	Embarked	
X ₁	2	1	2	2	0
X ₂	2	0	3	2	1
X ₃	2	1	3	2	0
X ₄	0	1	0	2	1
X ₅	1	1	3	2	0
X ₆	2	0	0	1	1
X ₇	1	1	3	2	0
X ₈	0	0	0	2	1
X ₉	2	1	3	2	0
X ₁₀	1	0	3	2	1
X ₁₁	1	1	3	2	0
X ₁₂	1	0	3	2	1
X ₁₃	2	1	3	0	0
X ₁₄	0	1	0	2	1
X ₁₅	0	1	0	2	0
X ₁₆	2	1	0	0	1



X ₁₇	2	1	0	1	0
X ₁₈	2	1	0	2	1
X ₁₉	1	1	0	2	0
X ₂₀	2	1	0	2	1

Step 07: Select Suitable Evaluation Measure(s)

SLIDE**Step 07: Select Suitable Evaluation Measure(s)**

- I will use the **Accuracy** Evaluation Measure to evaluate the performance of the Model
- Accuracy
 - Accuracy is defined as the proportion of correctly classified Test Instances

$$\text{Accuracy} = \frac{\text{Number of Correctly Classified Test Instances}}{\text{Total Number of Test Instances}}$$

- Note
 - Error = 1 - Accuracy

Step 08: Execute First Two Phases of Machine Learning Cycle

SLIDE**Step 8: Execute First Two Phases of Machine Learning Cycle**

- Recall the Equation

$$\text{Data} = \text{Model} + \text{Error}$$

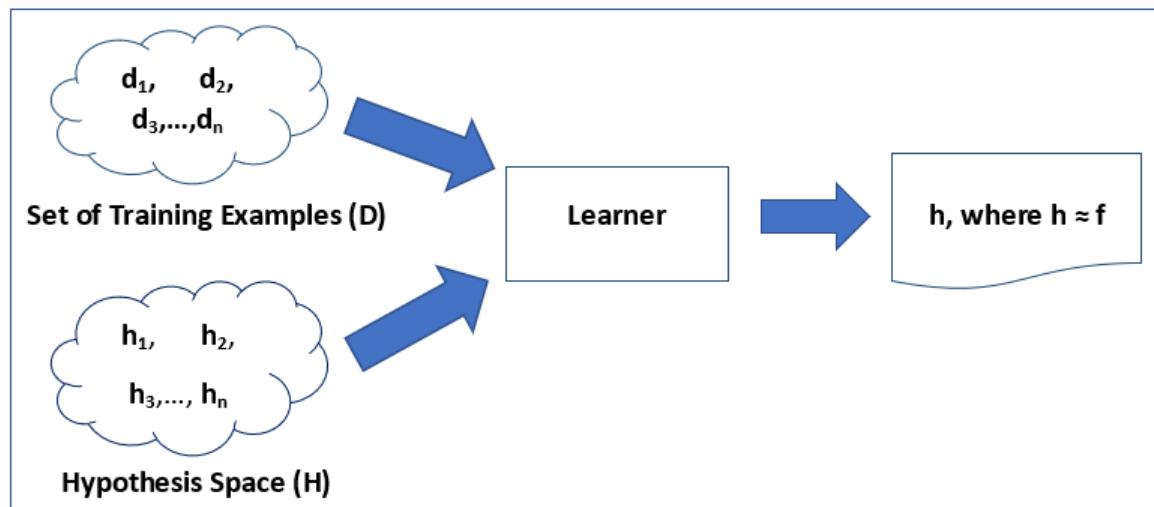


- Training Phase
 - Use Training Data to **build the Model**
- Testing Phase
 - Use Testing Data to **evaluate the performance of the Model**
- Note that we aim to
 - **Learn an Input-Output Function**

SLIDE

General Settings - Learning Input-Output Function

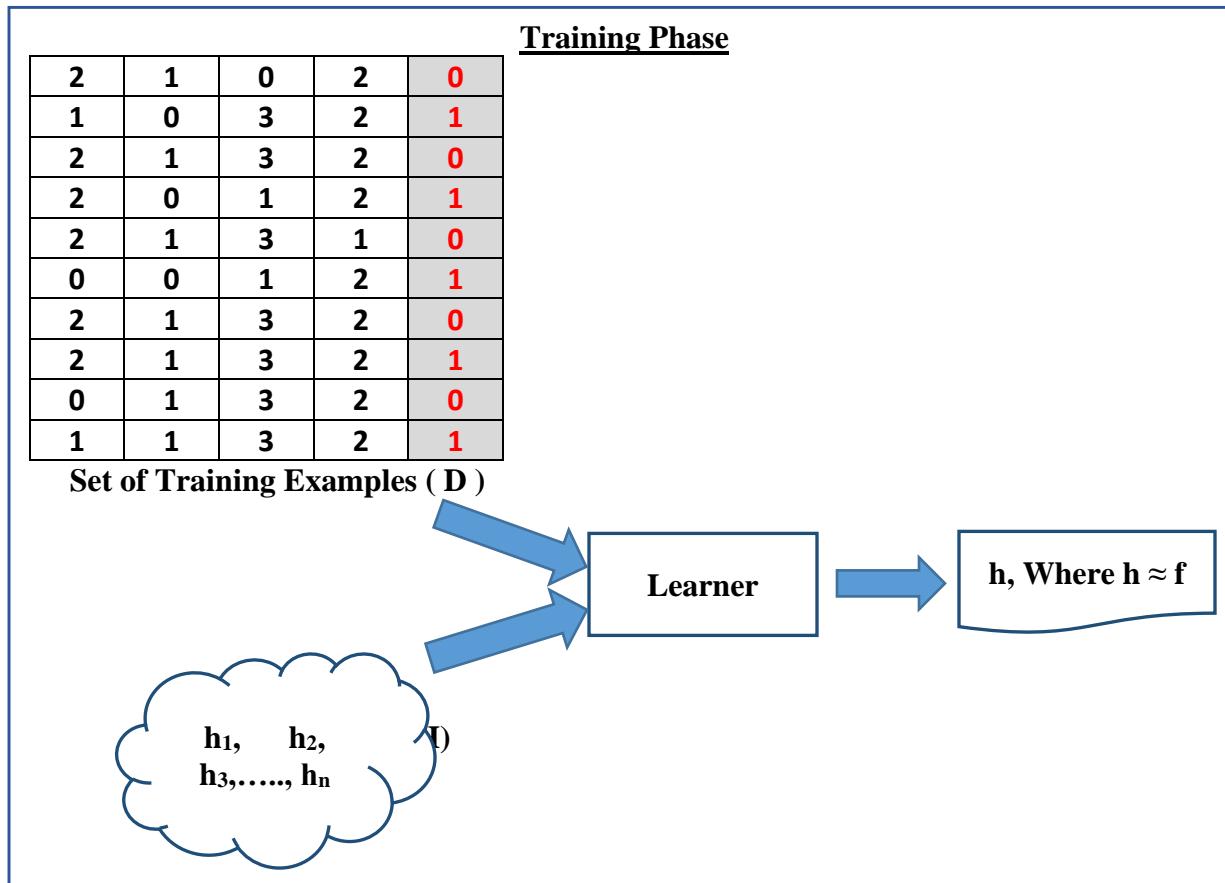
- Recall – Our **goal** is to
 - **Learn an Input-Output Function**



SLIDE



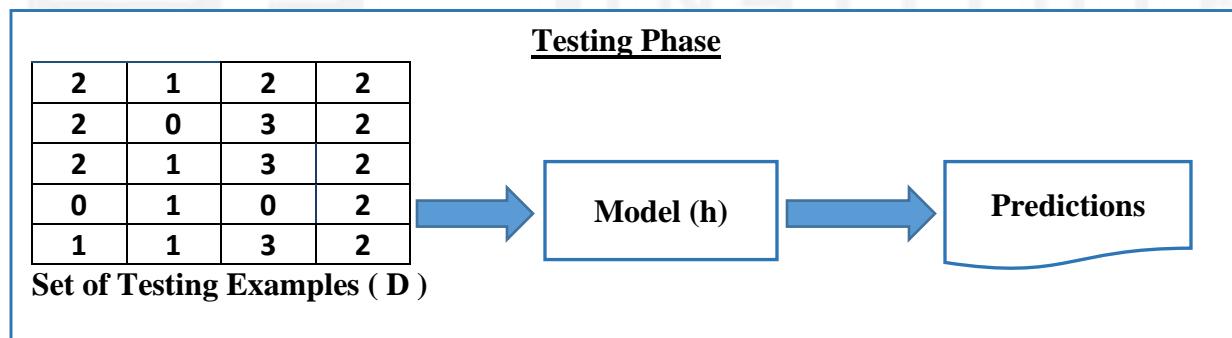
Training Phase



SLIDE

Testing Phase

- Apply Model on the Testing Data



SLIDE



Testing Phase Cont...

- The following Table shows the Predictions Returned by the Model (h)
 - See **model-predictions.csv** File in Supporting Material

Instance No.	Input				Output	
	PClass	Gender	Siblings	Embarked	Survived	Predictions
X ₁	2	1	2	2	0	0
X ₂	2	0	3	2	1	1
X ₃	2	1	3	2	0	0
X ₄	0	1	0	2	1	0
X ₅	1	1	3	2	0	0
X ₆	2	0	0	1	1	1
X ₇	1	1	3	2	0	0
X ₈	0	0	0	2	1	1
X ₉	2	1	3	2	0	0
X ₁₀	1	0	3	2	1	1
X ₁₁	1	1	3	2	0	0
X ₁₂	1	0	3	2	1	1
X ₁₃	2	1	3	0	0	0
X ₁₄	0	1	0	2	1	0
X ₁₅	0	1	0	2	0	0
X ₁₆	2	1	0	0	1	1
X ₁₇	2	1	0	1	0	0
X ₁₈	2	1	0	2	1	0
X ₁₉	1	1	0	2	0	0
X ₂₀	2	1	0	2	1	0

SLIDE

Testing Phase, Continue

- Calculating Accuracy
 - To calculate Accuracy, we will compare



■ Actual Values with Predicted Values

- Note

- To explain calculations more clearly, I have converted Numerical Predicted Values to Categorical Predicted Values

Instance No.	Input				Output		
	PClass	Gender	Sibling	Embarked	Actual Values	Predicted Values	Score
X1	Third	Male	Two	Southampton	No	No	1
X2	Third	Female	Zero	Southampton	Yes	Yes	1
X3	Third	Male	Zero	Southampton	No	No	1
X4	First	Male	One	Southampton	Yes	No	0
X5	Second	Male	Zero	Southampton	No	No	1
X6	Third	Female	One	Queenstown	Yes	Yes	1
X7	Second	Male	Zero	Southampton	No	No	1
X8	First	Female	One	Southampton	Yes	Yes	1
X9	Third	Male	Zero	Southampton	No	No	1
X10	Second	Female	Zero	Southampton	Yes	Yes	1
X11	Second	Male	Zero	Southampton	No	No	1
X12	Second	Female	Zero	Southampton	Yes	Yes	1
X13	Third	Male	Zero	Cherbourg	No	No	1
X14	First	Male	One	Southampton	Yes	No	0
X15	First	Male	One	Southampton	No	No	1
X16	Third	Male	One	Cherbourg	Yes	Yes	1
X17	Third	Male	One	Queenstown	No	No	1
X18	Third	Male	One	Southampton	Yes	No	0
X19	Second	Male	One	Southampton	No	No	1
X20	Third	Male	One	Southampton	Yes	No	0

Accuracy = $\frac{16}{20} = 0.8$



Step 09: Analyze Results

SLIDE

Step 9: Analyze Results

- The assumption for this Example
 - Here, I am **assuming** that the Model
 - **performed well on large Test Data** and we can apply it in the real-world 😊

Step 10: Execute 3rd and 4th Phases of Machine Learning Cycle

SLIDE

Step 10: Execute 3rd and 4th Phases of Machine Learning Cycle

- Application Phase
 - Model is **deployed** in **Real-world** to make **predictions** on **Real-time Data**
- Steps – Make Predictions on Real-time Data
 - Step 1: Take Input from User
 - Step 2: Convert **User Input** into **Feature Vector**
 - **The same** as **Feature Vectors** of Sample Data
 - Step 3: **Apply** Model on the **Feature Vector** of the unseen instance
 - Step 4: Return **Prediction** to the User

SLIDE

Example – Making Predictions on Real-time Data

- Step 1: Take Input from User
 - User Input

Please enter **PClass**: First

Please enter your **Gender**: Female



Please enter your **Sibling**: Two

Please enter **Embarked**: Cherbourg

- Step 2: Convert **User Input** into **Feature Vector**
 - Feature Vector

<First, Female, Two, Cherbourg>

- Feature Vector **after Label Encoding**
 - Exactly same as **Label Encoded Feature Vectors of Sample Data**
 - **Label Encoded Feature Vector**
 - <0, 0, 2, 0>
- Step 3: **Apply Model** on the **Label Encoded Feature Vector** of unseen instance
 - Model (h) is **applied** on <0, 0, 2, 0>
- Step 4: Return **Prediction** to the User
 - 1 (Yes)

SLIDE Application Phase

Application Phase



SLIDE Feedback Phase



- A Two-Step Process
- Step 1: After **some time**, take Feedback from
 - **Domain Experts** and **Users** on **deployed** Titanic Passenger Survival Prediction System
- Step 2: Make a **List of Possible Improvements** based on Feedback receive

Step 11: Improve Titanic Passenger Survival Prediction System based on Feedback

SLIDE

Step 11: Improve Titanic Passenger Survival Prediction System based on Feedback

- Go to Step 1 and **improve** the Titanic Passenger Survival Prediction System based on
 - **List of Possible Improvements** made in Step 10



TODO and Your Turn

SLIDE TODO

- Task
 - Consider the **Heart Disease Classification Problem**. The main aim is to **predict** whether a patient has Heart Disease or Not (i.e. Binary Classification Problem)?
 - Heart Disease Dataset Link
 - URL:
 - For simplicity, I have taken a sample of **100 instances** from the **Original Heart Disease Dataset**
 - See **heart-disease-sample-data.csv** File in Supporting Material
- Note
 - Your **answer** should be
 - Well Justified
- Question
 - Write down the **Input** and **Output** of the **Heart Disease Classification Problem**?
 - Follow the Steps mentioned in this Lecture and show
 - How will you treat the **Heart Disease Classification Problem** as a **Supervised Machine Learning Problem** using Train-Test Split Approach?

SLIDE Your Turn

- Task
 - Select a Problem (similar to the one given in TODO) and **answer the questions** given below
- Note
 - Your **answer** should be
 - Well Justified
- Questions



- Write **Input** and **Output** for the selected **Machine Learning Problem**?
- Follow the Steps mentioned in this Lecture and show
 - How will you treat the selected **Machine Learning Problem** as a **Supervised Machine Learning Problem** using Train-Test Split Approach?



Lecture Summary

SLIDE

Lecture Summary

- Titanic Passenger Survival Prediction System – **Task**
 - Given
 - A Passenger (Represented as Set of Attributes)
 - Task
 - Automatically predict whether the Passenger Survived or Not
- Titanic Passenger Survival Prediction System – **Input and Output**
 - Input
 - A Passenger
 - Output
 - Survived / Not Survived
- The Problem of Titanic Passenger Survival Prediction is treated as a
 - Supervised Machine Learning Task
- The main goal of Titanic Passenger Survival Prediction System is to
 - Learn an Input-Output Function
 - i.e. Learn from Input to predict the Output
- Learning Input-Output Function – General Settings
 - Input to Learner
 - Set of Training Examples (D)
 - Set of Hypothesis (a.k.a. Hypothesis Space (H))
 - Job of Learner
 - The main job of a Learner is to search the Hypothesis Space (H) using the Set of Training Examples (D) to find out a Hypothesis (h) from Hypothesis Space (H), which best fits the Set of Training Examples (D)
 - Output of Learner
 - A Learner outputs a Hypothesis (h) from Hypothesis Space (H), which best fits the Set of Training Examples (D)
- Steps to treat the Titanic Passenger Survival Prediction System Problem as a Classification Problem



- Step 01: Decide the Learning Settings
- Step 02: Obtain Sample Data
- Step 03: Understand and Pre-process Sample Data
- Step 04: Represent Sample Data in Machine Understandable Format
- Step 05: Select Suitable Machine Learning Algorithms
- Step 06: Split Sample Data into Training Data and Testing Data
- Step 07: Select Suitable Evaluation Measure(s)
- Step 08: Execute First Two Phases of Machine Learning Cycle
 - Training Phase
 - Testing Phase
- Step 09: Analyze Results

If (Results are Good)
Then
Move to the Next Step
Else
Go to Step 01

- Step 10: Execute 3rd and 4th Phases of Machine Learning Cycle
 - Application Phase
 - Feedback Phase
- Step 11: Based on Feedback
 - Go to Step 01 and Repeat all the Steps