

Name: Salvador, Louis Abraham	Date: 10-04-25
Section CIT 401A IT31S1	Instructor: Richard Kho

I. BUSINESS UNDERSTANDING Richard

In this CRISP-DM exercise, your task is to build a classification model that achieves the highest possible accuracy by experimenting with different combinations of data preparation techniques. You'll be expected to apply sanitation, feature selection, and feature extraction methods to improve model performance and compare their effects systematically. As you move through the phases — from Business Understanding to Modeling—consider how each preprocessing choice influences the final outcome. Your goal is not just to build a model, but to demonstrate the value of thoughtful data preparation in predictive analytics.

Reflection (10 pts):

- a. Based on your understanding of model creation combined with data preparation task(s), do you see any benefit in performing this task? Please elaborate your answer. (10 pts)
- if the model doesn't have enough data it will be noisy, incomplete or unorganized, then the algorithm will perform poorly. Applying data preparation methods like feature selection, feature extraction will the model can learn patterns more effectively
 - Example if the dataset have many missing values and attributes. ignoring these issues, the model may give inaccurate predictions. if we handle missing data properly, encode the categories and remove the irrelevant attributes, the quality of the dataset improve and the model as well

II. DATA UNDERSTANDING

Extract the dataset from our Git repository located in the '20251004-Lab-Wk11' directory, and perform your own assessment and evaluation of the file's contents.

Reflection (40 pts):

- 1) Based on your assessment, what do you think the dataset is about? (5 pts)
- Each row links an order with metadata: Channel(inbound/outbound), issues category/sub-category, order/product details, agent/manager handling the case, timestamps for reporting/responding, and a customer satisfaction score(CSAT). the goal is to predict the CSAT from the interaction and order features.
- 2) What are the details of the dataset:
- a) List the field names along with their corresponding data types. (15 pts)
- Unique id is a string or identifier (object) that is unique to each record.
 - channel_name — string or category (object) – channel of contact (inbound, outbound, etc.).
 - category — string / categorical (object) — the top-level issue category (for example, Order, Product Query, Return).
 - Sub-category — text or categorical (object) — indicates a more specific issue kind.
 - There are around 28,756 non-null customer remarks, with approximately 66% of them missing. Good for NLP, but scarce.

Name: Salvador, Louis Abraham	Date: 10-04-25
Section CIT 401A IT31S1	Instructor: Richard Kho

- Order_id — string / identifier (object) — numerous missing (about 67,675 non-null).
- The order_date_time string/timestamp (object) has a considerable number of missing timestamps (about 17,214 non-null values). Needs to be parsed to datetime.
- Issue_reported at — string / timestamp (object) — the time when the customer reported the problem; this should be converted to datetime.
- issue_responded — string / timestamp (object) — the response timestamp, which should be converted to datetime.
- Survey_response_Date — string / date (object) — the date the CSAT was provided; parseable to date.
- Customer_City — string/categorical (object) — numerous missing (about 17,079 non-null).
- Product_category: string / categorical (object) — many lacking.
- Item_price — numeric (float) — item price; there are many missing values and some formatting errors in the raw file.
- The connected_handling_time numeric (float) variable is highly sparse, with only a few non-null values compared to the earlier observation of about 242.
- Agent_name — string or categorical (high cardinality) – thousands of distinct names.
- Supervisor — string / categorical (with high cardinality).
- Manager — string / categorical (with high cardinality).
- Tenure Bucket — string / category — agent tenure bracket (for example, 0-30, >90, on-the-job training).
- Agent Shift — string / category — Morning, Evening, Night.
- CSAT Score — integer or numeric (goal) — customer satisfaction score (1–5). This is numeric in the raw data, but it should be considered as the class label (nominal) for classification purposes.

b) Indicate the number of records. (2 pts)

Total rows (instances): 85,907

Total attributes (columns): 20

c) Identify the attribute in the dataset that is suitable to serve as the class label. (3 pts)

- CSAT Score is a good choice for the class name because it symbolizes the outcome (customer satisfaction) that we wish to predict.
- Practical note: Because CSAT is extremely skewed (most ratings = 5), it is frequently beneficial to convert it to a binary goal (e.g., Satisfied vs Unsatisfied) or utilize class-balancing procedures before modeling.

3) Provide a thorough analysis of the dataset in terms of its readiness for model evaluation.

Do you observe any anomalies in the dataset? (15 pts)

Name: Salvador, Louis Abraham	Date: 10-04-25
Section CIT 401A IT31S1	Instructor: Richard Kho

1. Severe class imbalance (CSAT).

- Most buyers give the highest rating (5). In modeling studies, the classifier frequently predicts class 5, whereas other classes have zero true positives => Weka shows? For precision/F-measurement. This is the most significant modeling issue.

2. There is a lot of missing data in critical fields.

- order_date_time and Customer_City contain a substantial fraction of missing data (about 17k non-null out of 85k).
- The Customer Remarks section is generally empty (~28k non-null vs. ~57k absent).
- There are several missing values in item_price and product_category.
- connected_handling_time has significantly few non-null rows. These render certain features useless or necessitate active imputation.

3. Small sample sizes after sampling.

- When balanced incorrectly (or with aggressive sampling), you may wind up with a small training set (e.g., 20 instances), resulting in unstable metrics and... For certain classes. Be sure to keep enough samples per class.

4. Possible mismatched labels or whitespace

- If CSAT values contain stray whitespace or are typed inconsistently (for example, "5"), Weka may treat them as separate categories — check and normalize.

III. DATA PREPARATION

Your task is to import the dataset into a MySQL DB table.

COL 1	COL 2	COL 3	COL 4	COL 5	COL 6	COL 7	COL 8	COL 9	COL 10	COL 11	COL 12	COL 13	COL 14	COL 15
Drag to reorder. Click to mark/favorite. Double-click to copy column name.			Sub-category	Customer Remarks	Order_id	order_date_time	Issue_reported at	Issue_responded	Survey_response_Date	Customer_City	Product_category	Item_price	connected_handling_time	Agent_name
7e9ae164-6a8b-4521-a2d4-5817c9ff13f	Outcall	Product Queries	Life Insurance		c27c9bb4-fa36-4140-9f1f-21009254ffdb		01/08/2023 11:13	01/08/2023 11:47	01-Aug-23					Richard Buchanan
b07ec1b0-f376-43b6-86df-ec03da3b2e16	Outcall	Product Queries	Product Specific Information		d406b0c7-ce17-4654-b9dc-f08d421254bd		01/08/2023 12:52	01/08/2023 12:54	01-Aug-23					Vicki Collins
200814dd-27c7-4149-ba2b-bd3af3092880	Inbound	Order Related	Installation/demo		c273368d-b961-44cb-beaf-62d6fd6c00d5		01/08/2023 20:16	01/08/2023 20:38	01-Aug-23					Duane Norman
eb0d3e53-c1ca-42d3-8486-e42c8d622135	Inbound	Returns	Reverse Pickup Enquiry		5aed0059-55a4-4ec6-bb54-97942092020a		01/08/2023 20:56	01/08/2023 21:16	01-Aug-23					Patrick Flores
ba903143-1e54-406c-b969-46c52f92e5df	Inbound	Cancellation	Not Needed		e8bed5a9-6933-4aff-9dc6-cccfaf7dcd59		01/08/2023 10:30	01/08/2023 10:32	01-Aug-23					Christopher Sanchez
1cfe5b9-6112-44fc-8f3b-892196137a62	Email	Returns	Fraudulent User		a2938961-2833-45f1-83d6-678d9555c603		01/08/2023 15:13	01/08/2023 18:39	01-Aug-23					Desiree Newton
11a3ff08-1d6b-4806-b198-c60b5934c9bc	Outcall	Product Queries	Product Specific Information		bfc6562b-9a2f-4cca-aa79-fd4e2952f901		01/08/2023 15:31	01/08/2023 23:52	01-Aug-23					Shannon Hicks
372b51a5-fa19-4a31-a4b8-a21de117d75e	Inbound	Returns	Exchange / Replacement	Very good	88537e0b-5ffa-43f9-bbe2-fe57a0f4e4ae		01/08/2023 16:17	01/08/2023 16:23	01-Aug-23					Laura Smith

Reflection (40 pts):

Name: Salvador, Louis Abraham	Date: 10-04-25
Section CIT 401A IT31S1	Instructor: Richard Kho

1) Provide a detailed description of the process you performed on the dataset during the data preparation stage? Was this necessary based on your observations from the “Data Understanding” stage? (30 pts)

- Preprocessing steps were necessary based on my observations in this Data Understanding stage.
- The dataset contained unnecessary attributes, missing or imbalanced values, and inconsistent class distributions that would negatively effect the model performance. Cleaning, Balancing and simplifying the data was it made to be ready for model evaluations and ensure the classifiers to produced valid and interpretable results

2) Upload your MySQL DB table dump to your own Git branch. (10 pts)

This branch is 11 commits ahead of, 10 commits behind main . Contribute		
Name	Last commit message	Last commit date
..		
.gitkeep	Create .gitkeep folder	2 minutes ago
20251004__dataset__wk11_1.sql	Add files via upload	1 minute ago

Name: Salvador, Louis Abraham	Date: 10-04-25
Section CIT 401A IT31S1	Instructor: Richard Kho

IV. MODELING

Weka Explorer

Preprocess **Classify** Cluster Associate Select attributes Visualize

Open file... Open URL... Open DB... Generate... Undo Edit... Save...

Filter: Choose **StringToNominal -R first-last** Apply Stop

Current relation: Relation: wk11_dataset-weka.filters.unsupervised.attribute.StringToNominal Attributes: 20 Sum of weights: 10 Instances: 10

Attributes

All None Invert Pattern

No.	Name
2	<input type="checkbox"/> channel_name
3	<input type="checkbox"/> category
4	<input type="checkbox"/> sub_category
5	<input type="checkbox"/> customer_remarks
6	<input type="checkbox"/> order_id
7	<input type="checkbox"/> order_date_time
8	<input type="checkbox"/> issue_reported_at
9	<input type="checkbox"/> issue_responded
10	<input type="checkbox"/> survey_response_date
11	<input type="checkbox"/> customer_city
12	<input type="checkbox"/> product_category
13	<input type="checkbox"/> item_price
14	<input type="checkbox"/> connected_handling_time
15	<input type="checkbox"/> agent_name
16	<input type="checkbox"/> supervisor
17	<input type="checkbox"/> manager
18	<input type="checkbox"/> tenure_bucket
19	<input type="checkbox"/> agent_shift
20	<input type="checkbox"/> csat_score

Remove

Selected attribute

Name: csat_score Missing: 0 (0%) Distinct: 2 Type: Nominal Unique: 0 (0%)

No.	Label	Count	Weight
1	Yes	5	5
2	No	5	5

Class: csat_score (Nom) Visualize All

ABFF-Viewer - C:\Users\spide\Downloads\wk11_clean.aff

1: channel_name	2: category	3: sub_category	4: customer_remarks	5: order_id	6: order_date_time	7: issue_reported_at	8: issue_responded	9: survey_response_date	10: customer_city	11: product_category	12: item_price	13: connected_handling_time	14: agent_name	15: supervisor	16: manager	17: tenure_bucket	18: agent_shift	19: csat_score
mail	Electronics	Mobile	Late delivery	ORD123	2024-05-01 10:00	2024-05-01 09:30	2024-05-02 11:00	2024-05-03	Manila	Smartphone	599	00:12:34	John Doe	Jane Smith	Michael Tan	1-2 years	Morning	Yes
hat	Fashion	Shoes	Missing size	ORD124	2024-05-02 14:10	2024-05-02 13:55	2024-05-03 15:20	2024-05-04	Cebu	Footwear	799	00:10:18	Alice Lee	David Cruz	Michael Tan	1 year	Evening	No
home	Appliances	Refrigerator	Broken handle	ORD125	2024-05-03 09:15	2024-05-03 09:00	2024-05-04 10:00	2024-05-05	Davao	Home Appliance	1599	00:09:45	Carlos Reyes	Jane Smith	Michael Tan	3-4 years	Morning	Yes
mail	Electronics	Laptop	Slow response	ORD126	2024-05-04 11:25	2024-05-04 10:55	2024-05-05 12:30	2024-05-06	Manila	Computer	35999	00:15:10	John Doe	Jane Smith	Michael Tan	1-2 years	Evening	No
hat	Beauty	Skincare	Excellent service	ORD127	2024-05-05 15:00	2024-05-05 14:30	2024-05-06 16:00	2024-05-07	Cebu	Cosmetics	499	00:08:22	Alice Lee	David Cruz	Michael Tan	1 year	Morning	Yes
home	Fashion	Accessories	Late update	ORD128	2024-05-06 09:30	2024-05-06 09:00	2024-05-07 10:15	2024-05-08	Manila	Jewelry	999	00:11:40	Carlos Reyes	Jane Smith	Michael Tan	3-4 years	Evening	No
mail	Electronics	Tablet	Great assistance	ORD129	2024-05-07 13:20	2024-05-07 13:00	2024-05-08 14:30	2024-05-09	Davao	Gadget	12999	00:07:12	John Doe	Jane Smith	Michael Tan	1-2 years	Morning	Yes
hat	Fashion	Bag	Incorrect color	ORD130	2024-05-08 10:45	2024-05-08 10:20	2024-05-09 11:50	2024-05-10	Cebu	Handbag	899	00:10:55	Alice Lee	David Cruz	Michael Tan	1 year	Evening	No
home	Appliances	Microwave	Satisfied customer	ORD131	2024-05-09 08:30	2024-05-09 08:00	2024-05-10 09:45	2024-05-11	Manila	Kitchen Appliance	2999	00:09:00	Carlos Reyes	Jane Smith	Michael Tan	3-4 years	Morning	Yes
mail	Electronics	Camera	Faulty lens	ORD132	2024-05-10 14:00	2024-05-10 13:30	2024-05-11 15:20	2024-05-12	Davao	Photography	10999	00:13:25	John Doe	Jane Smith	Michael Tan	1-2 years	Evening	No

Reflection (40 pts):

1) Using your top three (3) classification algorithms, complete the table below with your recorded measurements (30 pts):

Classifier	IBK	J48	NaiveBayes
Correctly Classified Instances	6(60%)	10(100%)	8(80%)
Incorrectly	4(40%)	0	2(20%)

Name: Salvador, Louis Abraham	Date: 10-04-25
Section CIT 401A IT31S1	Instructor: Richard Kho

Classified Instances			
Kappa statistics	0.2	1	0.6
Precision (Class: NO)	0.667	1	0.800
Precision (Class: YES)	0.571	1	0.800
Recall (Class: NO)	0.400	1	0.800
Recall (Class: YES)	0.800	1	0.800
F-Measures (Class: NO)	0.500	1	0.800
F-Measures (Class: YES)	0.667	1	0.800

2) Upload the model that achieves the highest measurement to your Git branch. (10 pts)

fba_it31s1 / Lab-10-04-25Week11 / Lab-10-04-25Week11 / WekaExplorer_Model /			Add file	...
Lanzel-123 Add files via upload			65c9ac8 · now	History
This branch is 15 commits ahead of, 10 commits behind main.			Contribute	
Name	Last commit message	Last commit date		
..				
.gitkeep	Create .gitkeep	10 minutes ago		
DecisionTable.model	Add files via upload	now		
J48.model	Add files via upload	now		
NaiveBayes.model	Add files via upload	now		

NOTE:

To build your model, follow these steps:

- In the “WEKA Explorer” window, under the “Result list” section, select the recorded result with the highest evaluation and right-click.
- From the options, select “Save model”.
- In the “Save” window, choose the destination where you want to save the model file (with the .model extension)