



Name: Shayan Umar

Designation: AI intern

Department: Artificial Intelligence

Date: 7-18-2025

Table of Contents

1. Theoretical Understanding:	3
1.1 Working with Numerical Data:	3
1.1.1 Normalization:	3
1.2 Working with Categorical Data:	12
1.3 Datasets, Generalization, and Overfitting:	14
2. Project Work:	20
2.1 Churn Prediction	20
2.2 Heart Disease Prediction	21
3. Challenges and Approaches:	23
3.1 Churn Prediction	23
3.2 Heart Disease Prediction	23
4. Results and Analysis:	25
4.1 Churn Prediction	25
4.2 Heart Disease Prediction	28
5. Conclusion:	30
6. References:	31

1. Theoretical Understanding:

1.1 Working with Numerical Data:

Numerical data includes values like temperature, weight, and counts that can be added and ordered. Some numbers, like postal codes, look numeric but represent categories and are treated as **categorical data**.

Model ingests an array of floating-point values [6.3, 67.4] known as a feature vector and not the data directly.

Working with Numerical data should include steps such as visualizing your data and statistically evaluating the data.

- Visualization (scatter, hist, box, barchart) plots.
- Statistical evaluation (mean, median, mode, variance, std).

1.1.1 Normalization:

After examining your data through statistical and visualization techniques, you should transform your data in ways that will help your model train more effectively. The goal of **normalization** is to transform features to be on a similar scale. For example, consider the following two features:

- Feature X spans the range 154 to 24,917,482.
- Feature Y spans the range 5 to 22.

These two features span very different ranges. Normalization might manipulate X and Y so that they span a similar range, perhaps 0 to 1.

Normalization provides the following benefits:

- Helps models *converge more quickly* during training. When different features have different ranges, gradient descent can "bounce" and slow convergence. That said, more advanced optimizers like Adagrad and Adam protect against this problem by changing the effective learning rate over time.

- Helps models *infer better predictions*. When different features have different ranges, the resulting model might make somewhat less useful predictions.
 - Helps *avoid the "NaN trap"* when feature values are very high. NaN is an abbreviation for *not a number*. When a value in a model exceeds the floating-point precision limit, the system sets the value to NaN instead of a number. When one number in the model becomes a NaN, other numbers in the model also eventually become a NaN.
 - Helps the model *learn appropriate weights* for each feature. Without feature scaling, the model pays too much attention to features with wide ranges and not enough attention to features with narrow ranges.
 - If you normalize a feature during training, you must also normalize that feature when making predictions.
-

Linear scaling

Linear scaling (more commonly shortened to just **scaling** or **MinMaxScaling**) means converting floating-point values from their natural range into a standard range—usually 0 to 1 or -1 to +1.

$$x' = (x - x_{min}) / (x_{max} - x_{min})$$

Linear scaling is a good choice when all of the following conditions are met:

- The lower and upper bounds of your data don't change much over time.
 - The feature contains few or no outliers, and those outliers aren't extreme.
 - The feature is approximately uniformly distributed across its range. That is, a histogram would show roughly even bars for most values.
-

Z-Score Scaling (Standardization)

Z-score scaling is a method to standardize data so it has:

- **Mean = 0**
- **Standard deviation = 1**

$$x' = (x - \mu) / \sigma$$

Where:

- x' = z-score
- x = original value
- μ = mean of the feature
- σ = standard deviation of the feature

Z-score tells you **how far a value is from the average** in terms of standard deviations.

For example:

- A z-score of +2 means the value is **2 standard deviations above the mean**
- A z-score of -1 means it's **1 standard deviation below the mean**
- Centers and scales the data
- Works well for algorithms like **Logistic Regression, SVM, and PCA**

In a classic normal distribution:

- At least 68.27% of data has a Z-score between -1.0 and +1.0.
- At least 95.45% of data has a Z-score between -2.0 and +2.0.
- At least 99.73% of data has a Z-score between -3.0 and +3.0.
- At least 99.994% of data has a Z-score between -4.0 and +4.0.

Log scaling

Log scaling computes the logarithm of the raw value. In theory, the logarithm could be any base; in practice, log scaling usually calculates the natural logarithm (\ln).

Log scaling (or **log transformation**) is used to **compress large values** and **spread out small ones**, especially when data is **highly skewed** or has **large outliers**.

$$x' = \ln(x)$$

Log scaling is helpful when the data conforms to a *power law* distribution. Casually speaking, a power law distribution looks as follows:

- Low values of X have very high values of Y.
- As the values of X increase, the values of Y quickly decrease. Consequently, high values of X have very low values of Y.

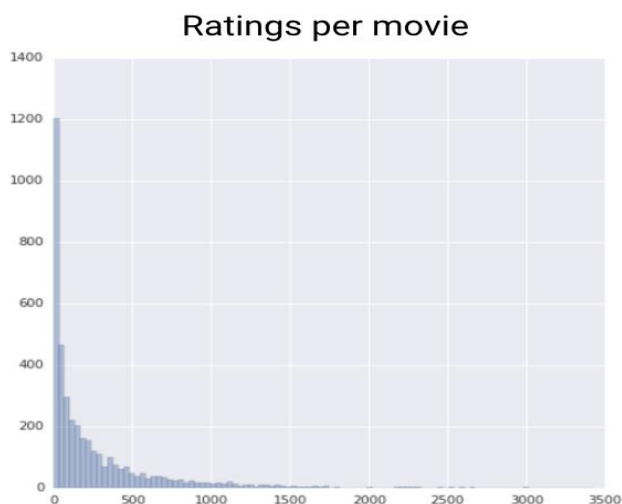
Movie ratings are a good example of a power law distribution. In the following figure, notice:

- A few movies have lots of user ratings. (Low values of X have high values of Y.)
- Most movies have very few user ratings. (High values of X have low values of Y.)

Log scaling changes the distribution, which helps train a model that will make better predictions.

When to Use It?

- Data is **right-skewed**
- You want to **reduce the impact of large values or outliers**
- Common in **financial data**, **population data**, etc



Clipping

Clipping means limiting the values of a feature to stay within a specific range.

Imagine your data has a few values that are **way too high or too low** compared to the rest. Instead of removing them, **clipping** just **cuts them off** at a certain threshold.

Example:

If your valid range is 0–100, and a value is 150, it gets clipped to 100.

Why Use Clipping?

- Reduce the impact of **extreme outliers**
- Keep the values within a **realistic or meaningful range**
- Useful when you **don't want to remove outliers** but **don't want them to dominate** either

You can also clip values after applying other forms of normalization. For example, suppose you use Z-score scaling, but a few outliers have absolute values far greater than 3. In this case, you could:

- Clip Z-scores greater than 3 to become exactly 3.
- Clip Z-scores less than -3 to become exactly -3.

Normalization technique	Formula	When to use
Linear scaling	$x' = \frac{x - x_{min}}{x_{max} - x_{min}}$	When the feature is mostly uniformly distributed across range. Flat-shaped
Z-score scaling	$x' = \frac{x - \mu}{\sigma}$	When the feature is normally distributed (peak close to mean). Bell-shaped
Log scaling	$x' = \log(x)$	When the feature distribution is heavy skewed on at least either side of tail. Heavy Tail-shaped
Clipping	If $x > max$, set $x' = max$ If $x < min$, set $x' = min$	When the feature contains extreme outliers.

Binning:

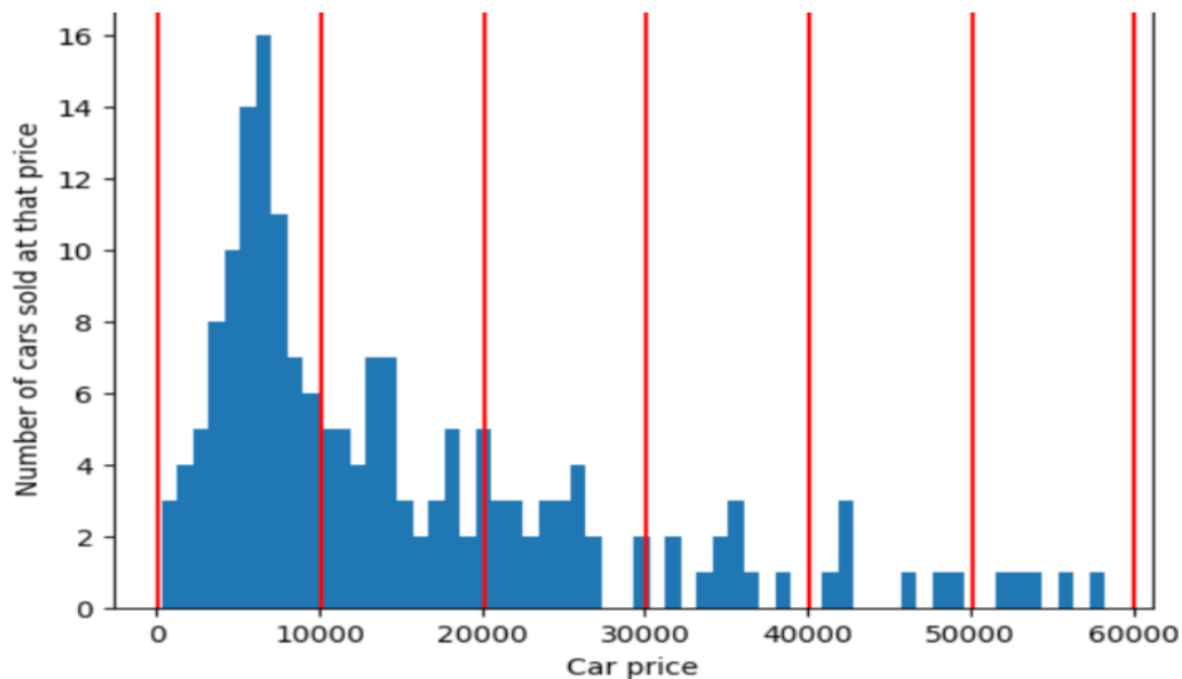
Binning (also called **bucketing**) is a **feature engineering** technique that groups different numerical subranges into *bins* or *buckets*. In many cases, binning turns numerical data into categorical data. For example, consider a **feature** named X whose lowest value is 15 and highest value is 425. Using binning, you could represent X with the following five bins:

- Bin 1: 15 to 34
- Bin 2: 35 to 117
- Bin 3: 118 to 279
- Bin 4: 280 to 392
- Bin 5: 393 to 425

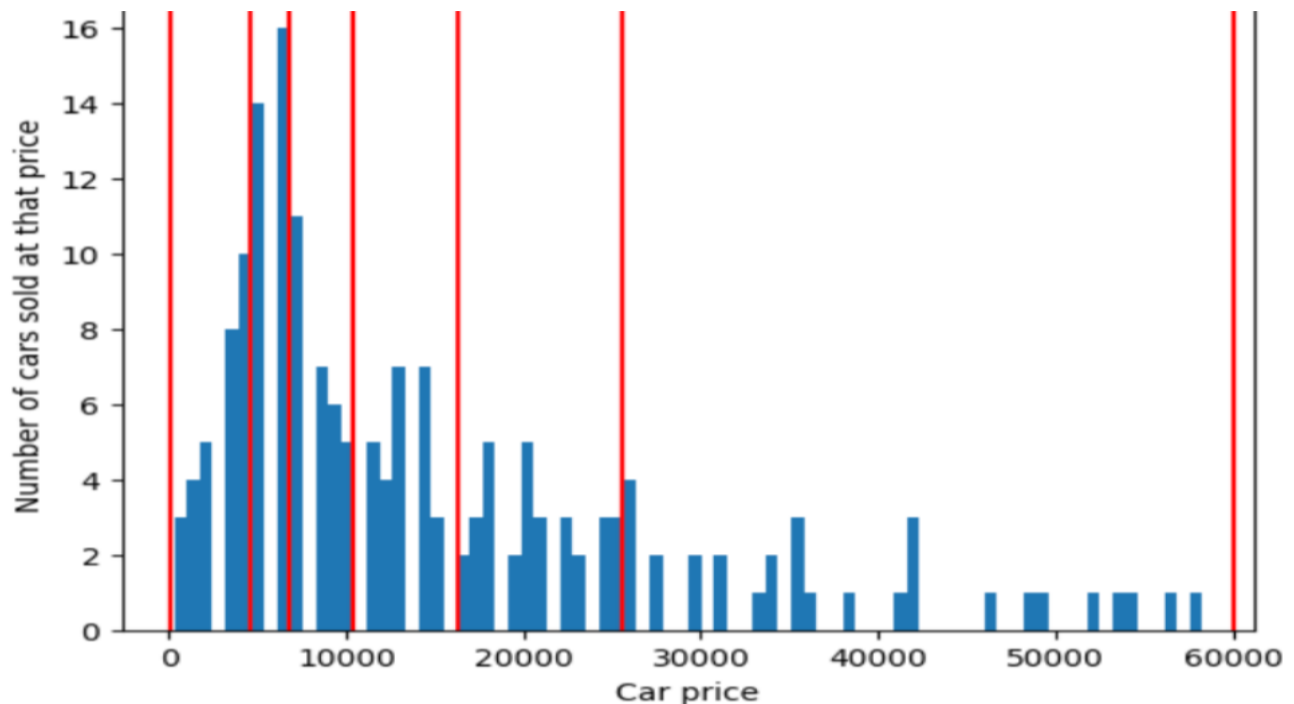
Quantile Bucketing

Quantile bucketing creates bucketing boundaries such that the number of examples in each bucket is exactly or nearly equal. Quantile bucketing mostly hides the outliers.

To illustrate the problem that quantile bucketing solves, consider the equally spaced buckets shown in the following figure, where each of the ten buckets represents a span of exactly 10,000 dollars. Notice that the bucket from 0 to 10,000 contains dozens of examples but the bucket from 50,000 to 60,000 contains only 5 examples. Consequently, the model has enough examples to train on the 0 to 10,000 bucket but not enough examples to train on for the 50,000 to 60,000 bucket.



In contrast, the following figure uses quantile bucketing to divide car prices into bins with approximately the same number of examples in each bucket. Notice that some of the bins encompass a narrow price span while others encompass a very wide price span.



Scrubbing (Data Cleaning)

Scrubbing refers to the process of **cleaning and correcting raw data** to make it usable for analysis or modeling.

Think of scrubbing like **washing dirty data** — removing the mess so your model sees clean and meaningful inputs.

Common Scrubbing Tasks:

- Removing or filling in **missing values**
- Fixing **typos** or **inconsistent formatting**
- Removing **duplicates**
- Correcting **wrong data entries** (e.g., age = -5)
- Filtering **irrelevant or junk data** (like spam rows)

Why is Scrubbing Important?

- Dirty data can lead to **wrong results**
 - Ensures your model learns from **accurate and consistent information**
 - A critical first step in any **data preprocessing pipeline**
-

Synthetic Features

Synthetic features are **new features** you create by combining or transforming existing ones to help your model perform better.

Examples:

- $\text{BMI} = \text{weight} / \text{height}^2$
- $\text{Age} = 2025 - \text{birth_year}$
- $\text{Area} = \text{length} \times \text{width}$

These are **manually created** based on **domain knowledge**.

Polynomial Transformation

Polynomial transformation is a **specific way to automatically create synthetic features** by raising features to powers and multiplying them together.

Example:

For a single feature x , a degree-2 polynomial transform will create:

- x (original)
- x^2 (squared)

With two features x and y , you'll get:

- $x, y, x^2, y^2, x*y$

So you're creating **new features** like x^2 and $x*y$ — which are also **synthetic features**, just **automatically generated**.

Synthetic Features Any new feature created from existing data

Polynomial Transform A method of creating synthetic features using powers & products

Polynomial features are a type of synthetic feature, but not all synthetic features are polynomial.

1.2 Working with Categorical Data:

True numerical data can be meaningfully multiplied, like house area impacting price — e.g., a 200 m² house being worth roughly twice as much as a 100 m² house.

However, not all integers should be treated as numerical. Some, like postal codes, are better represented as categorical data, since their numeric values don't reflect real-world relationships. Treating them numerically may mislead the model; categorically, the model can handle each postal code independently.

Encoding:

Encoding means converting categorical or other data to numerical vectors that a model can train on. This conversion is necessary because models can only train on floating-point values; models can't train on strings such as "dog" or "maple".

One Hot Encoding:

For **categorical features** with a **small number of categories** (low-dimensional), it's best to use **vocabulary encoding**, where each unique value is treated as a separate feature. This allows the model to learn a **specific weight** for each category.

For example, in a feature like `car_color`, the model might learn that:

- Red cars tend to be more valuable than green ones.

However, simply **converting strings to integers** (e.g., red = 1, blue = 2, orange = 3) and feeding those into a model can be **misleading**. The model would interpret these numbers as having **mathematical meaning** — like thinking "orange (3)" is **three times more important** than "red (1)", which is **not true**.

To avoid this, the categorical data must be **encoded properly** (e.g., one-hot or embeddings) so the model doesn't falsely treat categories as numeric quantities

Sparse Feature:

A **sparse feature** is a feature (column in your dataset) that contains **mostly zeros or missing values**.

Sparse features take up a **lot of space** but **contain little info**.

They need **special handling**:

- Use **sparse matrices** to save memory.
 - Use models that can handle sparsity well (like **tree-based models**, or **linear models with sparse input**).
-

Outliers in Categorical Data

Rare categories in categorical data (like unusual car colors) are treated as **outliers**. Instead of giving each a separate category, they're grouped into a single "**out-of-vocabulary (OOV)**" or "**other**" category. The model then learns **one weight** for all these rare values, improving simplicity and generalization.

When the number of categories is high, one-hot encoding is usually a bad choice. Embeddings are usually a much better choice. Embeddings substantially reduce the number of dimensions, which benefits models in two important ways:

- The model typically trains faster.
- The built model typically infers predictions more quickly. That is, the model has lower latency.

Hashing (also called the *hashing trick*) is a less common way to reduce the number of dimensions.

Common Issues with Categorical Data:

Categorical data is often labeled by **humans** or **ML models**. Labels given by humans are called **gold labels** and are usually preferred for training due to higher quality.

However, human-labeled data isn't always perfect — it can include **errors, bias, or malicious inputs**. It's important to **check for quality issues** before training.

Different human raters may label the same data differently. This variability is measured using **inter-rater agreement**, which helps assess the **consistency** of labels by comparing multiple raters' decisions.

Feature Crosses

Feature crosses combine two or more **categorical or bucketed features** by taking their **Cartesian product**. This helps **linear models** capture **nonlinear patterns** and **interactions** between features — similar to what polynomial transforms do for numerical features.

When to Use:

- Use when **domain knowledge** suggests certain feature combinations matter.
- Can be **automatically learned** in deep models (e.g., neural networks).
- But crossing **sparse features** can lead to **very large and sparse** data — e.g., crossing $100 \times 200 = 20,000$ features.

1.3 Datasets, Generalization, and Overfitting:

Understanding Datasets for Machine Learning

What is a Dataset?

- A dataset is a collection of **examples**, often stored in **tables** (e.g., CSVs, spreadsheets).
- Each **row** represents an example, and each **column** is a feature or label.

Types of Data:

- **Numerical** and **categorical** (commonly used in ML)
- **Text, multimedia** (images, audio, videos)
- **Embeddings** and **outputs** from other ML systems

Quantity of Data:

- Good models typically require **more examples than parameters**.

- **Large datasets with fewer features** often perform better than small datasets with many features.
- **Small datasets** can still work well if fine-tuning a pretrained model.

Data Quality and Reliability:

- A **high-quality dataset** helps the model perform its task effectively.
- **Reliability** means the data is trustworthy and consistent.

Common issues that affect quality:

- **Missing values** (e.g., blank entries)
 - **Duplicates** (e.g., repeated records)
 - **Incorrect feature values** (e.g., typos, sensor errors)
 - **Wrong labels** (e.g., misclassified data)
 - **Corrupt data segments** (e.g., due to system errors)
-

Direct vs. Proxy Labels

Direct labels are exactly what the model is trying to predict and should be used when available. For example, a `bicycle_owner` column is a direct label for predicting ownership.

Proxy labels are indirect substitutes used when direct labels are unavailable or not easily usable. Their usefulness depends on how strongly they correlate with the actual prediction target. They are a compromise and should only be used if no better option exists.

Human-Generated Data

Human-labeled data can handle complex or subjective tasks and ensures clear labeling standards. However, it is often expensive and subject to human error, which may require multiple raters per example.

To ensure quality, define clear labeling instructions, assess the needed skill level of raters, and validate their work by comparing with your own labeled examples. If inconsistencies arise, revise instructions and re-label.

Handling Imbalanced Datasets

Imbalanced datasets often lack enough minority class examples, causing models to train mostly on the majority class. This can lead to poor learning of the minority class, especially when batch sizes result in little to no positive examples.

When to Worry:

- If the model performs well on the original dataset, you may not need to rebalance.
 - Start with the original data, use it as a **baseline**, and only apply rebalancing if performance is poor.
-

Techniques to Handle Imbalance

1. Downsampling

- Reduce the number of **majority class** examples during training.
- Example: If the original ratio is 1:200, downsampling by a factor of 10 changes it to ~1:20.

2. Upweighting

- Increase the **importance (weight)** of downsampled examples during loss calculation.
- If you downsample by a factor of 10, assign a **weight of 10** to those examples.
- This maintains influence on the model without overwhelming it with raw quantity.

Note: These **example weights** are not model weights (like w_1 , w_2); they modify how much each example contributes to the loss function.

Why Upweight After Downsampling?

Though counterintuitive, upweighting the majority class after downsampling helps reduce **prediction bias** — it keeps the model's output distribution more in line with the dataset's true label distribution.

Rebalancing Ratios

There's no fixed rule for how much to downsample or upweight. It depends on:

- **Batch size**
- **Imbalance ratio**
- **Dataset size**

Tip: Each batch should include **multiple minority class examples**.

For example, if the imbalance ratio is 100:1, choose a batch size of **at least 500** to increase the chance of including minority class samples in every batch.

Training, Validation, and Test Sets

Why Split the Dataset?

To fairly evaluate model performance, you must test on data that the model **hasn't seen during training**. Splitting ensures more **realistic and trustworthy evaluation**.

Three-Way Split:

1. Training Set

- Used to fit the model and adjust weights/parameters.

2. Validation Set

- Used during model development to tune hyperparameters and evaluate improvements.
- Helps decide when to stop training or which version of the model is best.

3. Test Set

- Used **only once** after finalizing the model to assess real-world performance.
 - Acts as the final evaluation before deployment.
-

Best Practices:

- Avoid overusing the same validation/test sets, as repeated use reduces their effectiveness.

- Periodically **refresh** your validation/test sets with **new, unseen data** to maintain evaluation integrity.
 - Ensure **no duplicates** exist between training and validation/test sets. Duplicates can lead to **inflated performance** and **misleading results**.
-

A Good Test/Validation Set Should Be:

- **Large enough** for statistical reliability.
 - **Representative** of the entire dataset's characteristics.
 - **Similar to real-world data** the model will face in production.
 - **Free from duplicates** found in the training data.
-

1. Generalization

- **Definition:** Generalization refers to a model's ability to perform well on **unseen/test data**, not just the data it was trained on.
 - **Importance:** A good model should not memorize the training data—it should learn the underlying patterns and apply them to new data effectively.
 - **Goal:** Achieve high accuracy (or other metrics) on both training and test sets.
-

2. Overfitting

- **Definition:** Overfitting happens when the model learns **both the actual patterns and the noise** in the training data.
- **Symptoms:**
 - Very high training accuracy.
 - Much lower validation/test accuracy.
- **Cause:** Too complex model (e.g., too many features, too many tree depths, or high-degree polynomial).
- **Consequence:** Poor generalization to new data.

3. Underfitting

- **Definition:** Underfitting occurs when the model is **too simple** to capture the patterns in the data.
 - **Symptoms:**
 - Low accuracy on both training and test data.
 - **Cause:** Model lacks capacity (e.g., missing key features, linear model on non-linear data).
 - **Consequence:** Poor performance across all data.
-

4. L1 Regularization (Lasso)

- **Mechanism:** Adds a penalty equal to the **absolute value** of model coefficients to the loss function:

$$\text{Loss} = \text{Error} + \lambda \sum |w_i|$$

- **Effect:**
 - Can shrink some coefficients to **zero** → performs **feature selection**.
 - Helps in reducing overfitting and simplifying the model.
 - **Use Case:** When you suspect many irrelevant features and want a sparse model.
-

5. L2 Regularization (Ridge)

- **Mechanism:** Adds a penalty equal to the **squared value** of model coefficients:

$$\text{Loss} = \text{Error} + \lambda \sum w_i^2$$

- **Effect:**
 - Shrinks coefficients towards zero but does **not eliminate** them.
 - Reduces the effect of collinearity and **stabilizes** the model.
 - **Use Case:** When all features are expected to contribute, but you want to **reduce model complexity** and overfitting.
-

2. Project Work:

2.1 Churn Prediction

Objective

Predict customer churn (whether a customer will leave a service) based on demographic and transactional features.

Methodology

1. Data Preprocessing

- **Handled Missing Values:**
 - gender: Filled with mode (most frequent category).
 - dependents: Assumed missing = 0 (no dependents).
 - occupation: Created an "Unknown" category.
 - city: Filled with mode (most frequent city).
 - days_since_last_transaction: Used KNN imputation.
- **Feature Engineering:**
 - One-hot encoded occupation (nominal categorical variable).
 - Log-transformed skewed numerical features (vintage, days_since_last_transaction).
 - Standardized numerical features using StandardScaler.
- **Class Imbalance Handling:**
 - Downsampled the majority class (non-churners) to balance the dataset.

2. Model Training

- **Logistic Regression:**
 - Used class weights ({0:1, 1:3}) to handle imbalance.
 - Trained on scaled features.
- **Random Forest:**
 - Optimized threshold for better precision-recall trade-off.

3. Evaluation Metrics

- **Confusion Matrix:**
 - Evaluated true positives, false positives, true negatives, false negatives.
- **Recall:**
 - Measures how many actual churners were correctly identified (critical for retention strategies).
- **ROC-AUC Score:**
 - Assessed model's ability to distinguish between churners and non-churners

Metric	Logistic Regression	Random Forest (Optimized)
Accuracy	0.62	0.72
Precision	0.68	0.72
Recall	0.76	0.86
ROC-AUC	0.82	0.85

2.2 Heart Disease Prediction

Objective

Predict the presence of heart disease in patients based on clinical and demographic features.

Methodology

1. Data Preprocessing

- **Missing Values:** No missing values found (simplified preprocessing).
- **Feature Scaling:** Applied StandardScaler to normalize numerical features.
- **Exploratory Data Analysis (EDA):**
 - Identified top correlated features:
 - cp (chest pain type, +0.43)
 - thalach (max heart rate, +0.42)

- slope (ST segment slope, +0.35)

2. Model Training

- **Logistic Regression:** Baseline model with interpretable coefficients.
- **KNN:** Hyperparameter-tuned for optimal neighbors (k=7).
- **Random Forest:** Ensemble method for non-linear relationships.

3. Evaluation Metrics

- **Accuracy:** Overall correctness of predictions.
- **Precision-Recall Trade-off:** Critical for medical diagnosis (avoiding false negatives).
- **Feature Importance:** cp (chest pain) had the highest impact.

Results

Model	Accuracy	Precision	Recall	F1-Score
Logistic Regression	0.85	0.87	0.84	0.85
KNN (k=7)	0.91	0.88	0.87	0.87
Random Forest	0.83	0.86	0.85	0.85

3. Challenges and Approaches:

3.1 Churn Prediction

1. Class Imbalance

- *Problem:* Far more non-churners than churners.
- *Solution:* Downsampled majority class + used class weights in Logistic Regression.

2. Missing Data

- *Problem:* Gaps in dependents, occupation, days_since_last_transaction.
- *Solution:*
 - Filled categorical columns with mode.
 - Used KNN imputation for numerical gaps.

3. Feature Engineering

- *Problem:* Skewed numerical features (vintage, days_since_last_transaction).
- *Solution:* Applied **log transformation** + standardization.

4. Model Selection

- *Problem:* Logistic Regression had low recall.
- *Solution:* Switched to **Random Forest** + optimized threshold for better recall.

3.2 Heart Disease Prediction

1: Model Selection

- **Problem:** Logistic Regression had lower accuracy (82%) compared to KNN.
- **Solution:** Used **KNN with hyperparameter tuning** (k=7) for better performance.

2: Feature Interpretation

- **Problem:** Non-linear relationships (e.g., thalach vs. target).
- **Solution:** Leveraged **Random Forest** for feature importance and non-linear patterns.

3: Medical Implications

- **Problem:** High recall was critical (avoiding false negatives in disease prediction).
- **Solution:** Prioritized **recall metrics** during model evaluation.

4: Reproducibility

- **Problem:** Needed consistent scaling for deployment.
- **Solution:** Saved the StandardScaler object using joblib.

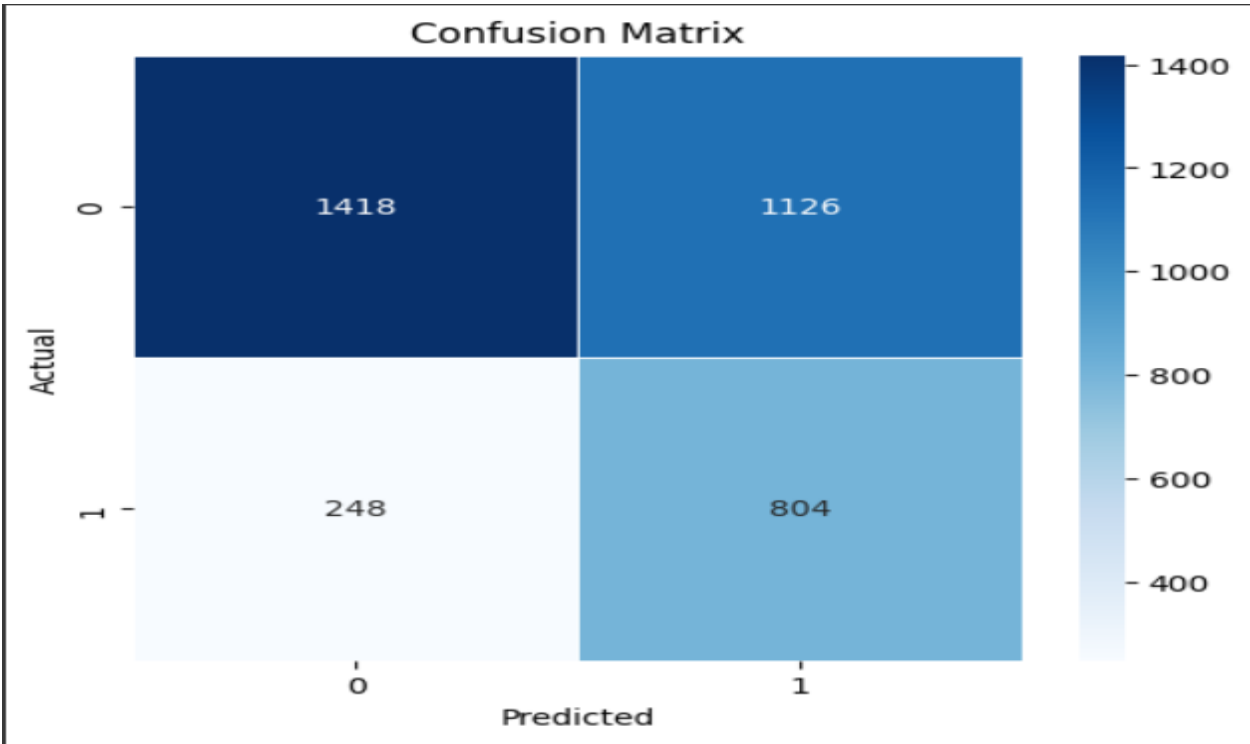
4. Results and Analysis:

4.1 Churn Prediction

Classification Report:

	precision	recall	f1-score	support
0	0.85	0.56	0.67	2544
1	0.42	0.76	0.54	1052
accuracy			0.62	3596
macro avg	0.63	0.66	0.61	3596
weighted avg	0.72	0.62	0.63	3596

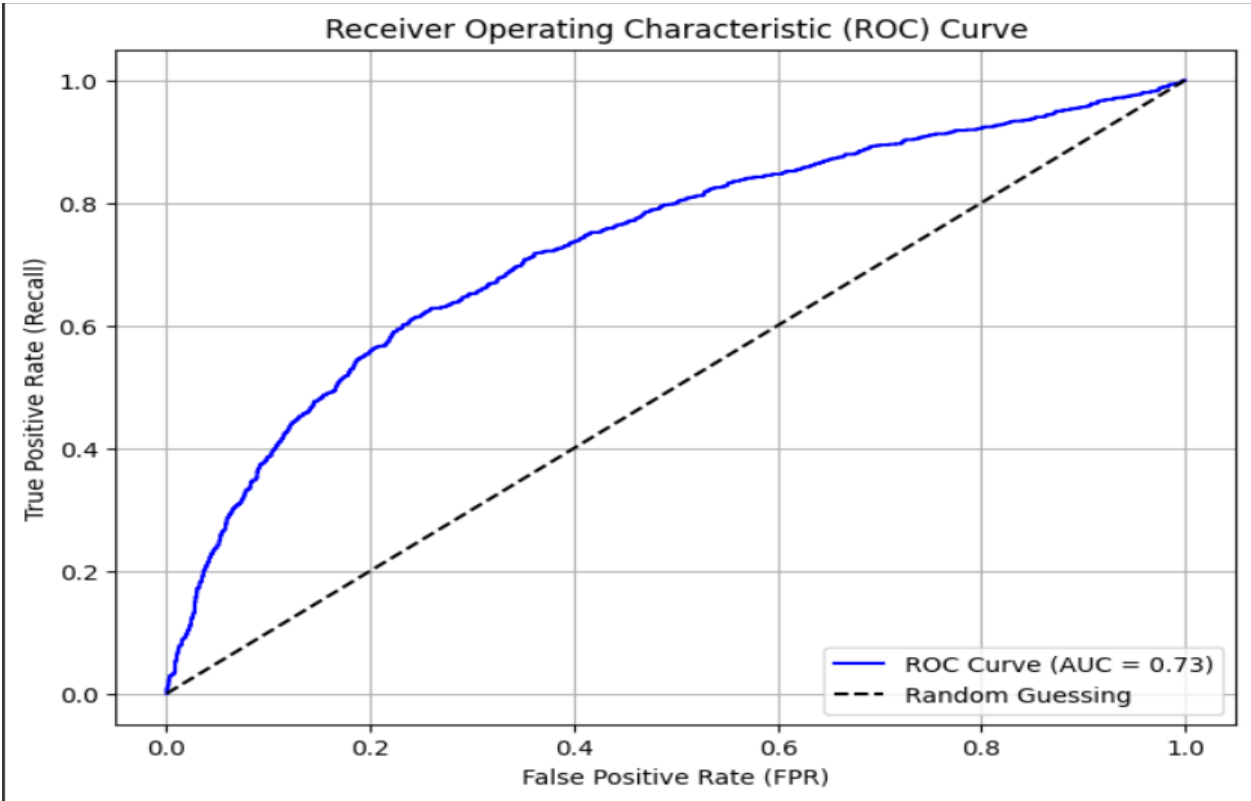
Confusion Matrix:



Accuracy: 0.617908787541713

AUC_SCORE: 0.7301994404189682

ROC Graph:

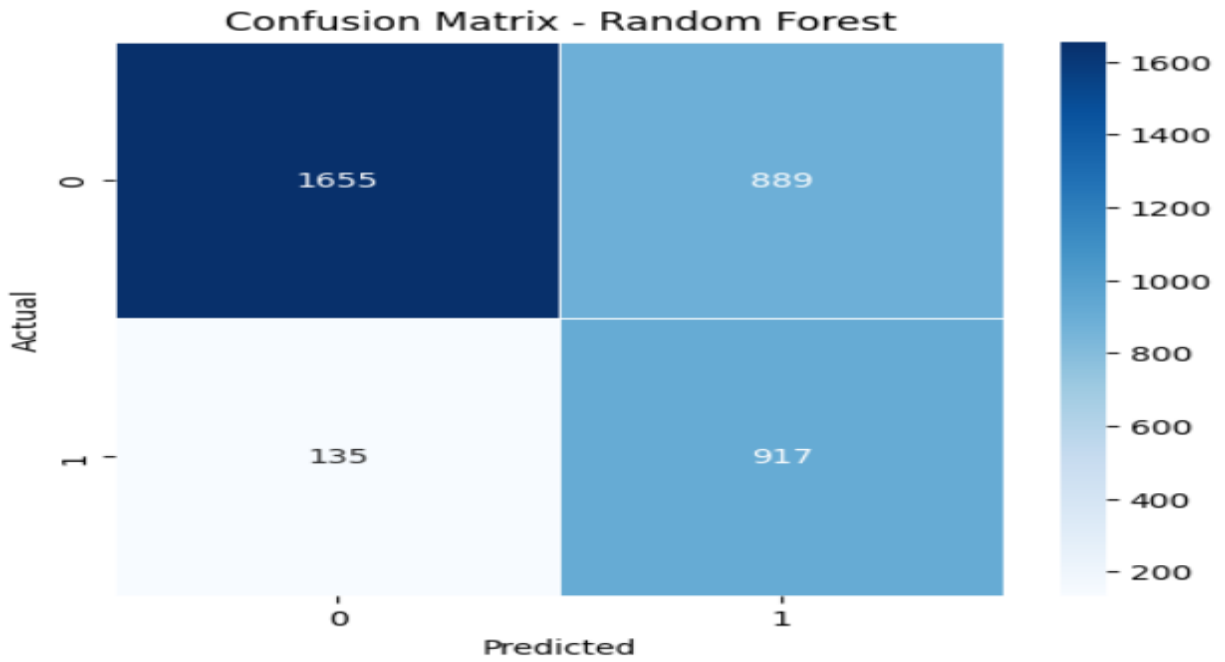


Random Forest:

Classification Report:

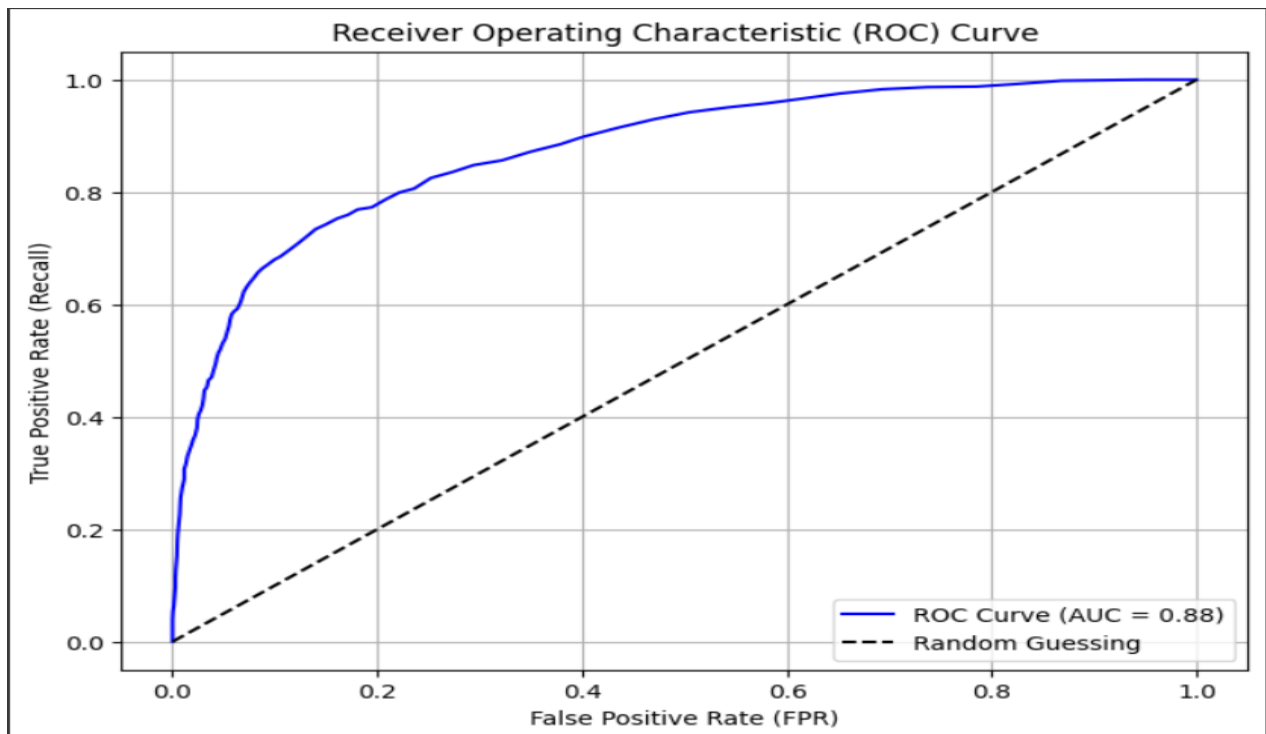
	precision	recall	f1-score	support
0	0.92	0.65	0.76	2544
1	0.51	0.87	0.64	1052
accuracy			0.72	3596
macro avg	0.72	0.76	0.70	3596
weighted avg	0.80	0.72	0.73	3596

Confusion Matrix:



AUC-SCORE: 87

ROC Curve:



4.2 Heart Disease Prediction

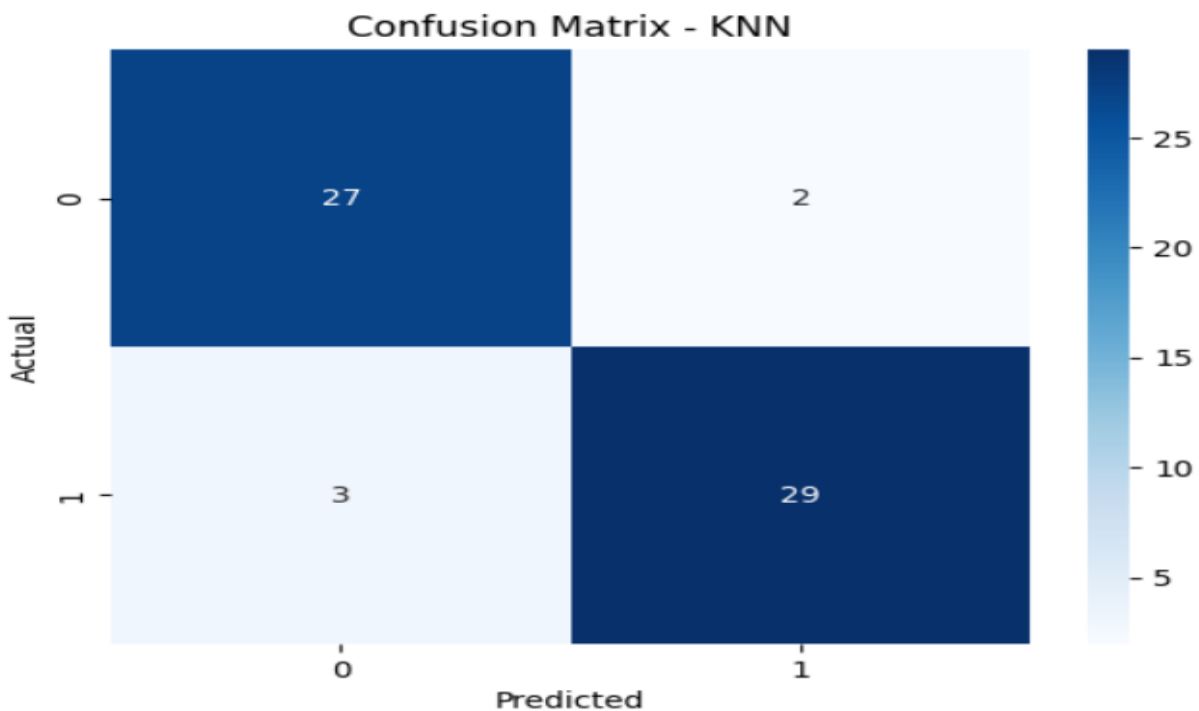
For Heart Disease Prediction I trained 3 machine learning models (Logistic Regression, Random Forest and KNeighborsClassifier) among these 3 KNN performed the best.

Accuracy and Classification Report:

```
Accuracy: 0.9180327868852459
```

Classification Report:					
	precision	recall	f1-score	support	
0	0.90	0.93	0.92	29	
1	0.94	0.91	0.92	32	
accuracy			0.92	61	
macro avg	0.92	0.92	0.92	61	
weighted avg	0.92	0.92	0.92	61	

Confusion Matrix:

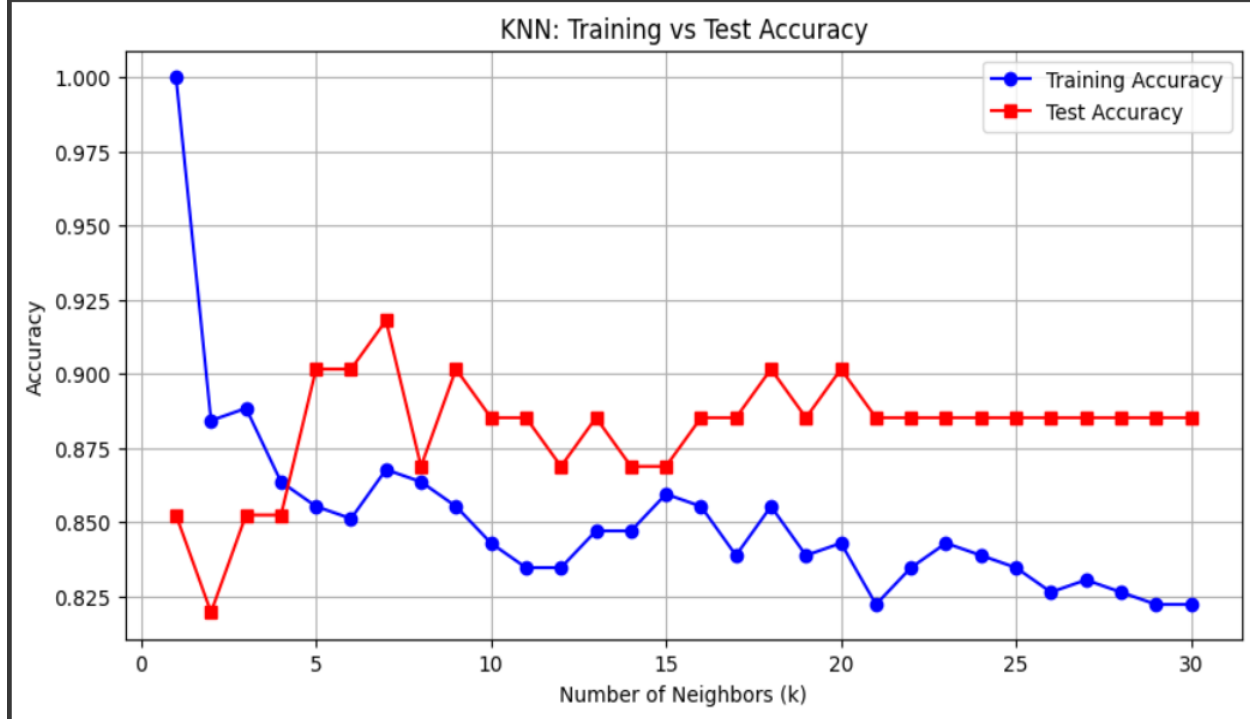


Best Result of KNN after Hyperparameters Tuning:

7. HYPERPARAMETER TUNING FOR KNN:

Optimal number of neighbors: 7

Best test accuracy: 0.9180



5. Conclusion:

Task 1: Churn Prediction

Key Takeaways:

- **Best Model:** Optimized Random Forest achieved **72% accuracy** with high recall (86%) to catch more churners.
- **Top Features:** days_since_last_transaction, current_balance, and occupation type were most influential.
- **Business Impact:** High recall helps retain at-risk customers proactively.

Future Improvements:

1. **Feature Engineering:** Add transaction frequency and customer tenure trends.
2. **Model Tuning:** Test **XGBoost** with class weights for better precision-recall balance.
3. **Real-Time Use:** Deploy as an API to flag high-risk customers in banking apps.
4. **Feedback Loop:** Incorporate customer intervention outcomes to refine predictions.

Task 2: Heart Disease Prediction

Key Takeaways:

- **Best Model:** KNN (k=7) achieved **91% accuracy**; Logistic Regression offered interpretability.
- **Top Features:** Chest pain type (cp), max heart rate (thalach), and ST slope were clinically relevant.
- **Medical Value:** High recall (87%) reduces missed diagnoses.

Future Improvements:

1. **Advanced Models:** Experiment with **LightGBM** and **SHAP** for explainability.
2. **Data Enrichment:** Include blood pressure trends and stress test results.
3. **Clinical Integration:** Build a dashboard for doctors with risk scores and key flags.
4. **Global Validation:** Test model on diverse ethnic groups to ensure robustness.

6. References:

- <https://scikit-learn.org/stable/>
- <https://developers.google.com/machine-learning/crash-course/>
- <https://www.youtube.com/>
- <https://www.youtube.com/@campusx-official>