

Introduction To Artificial Intelligence And Machine Learning.

Adithya Nair

August 1, 2024

Contents

1	Iris Data Classification	2
2	Overview	2
2.1	Pre-Processing	2
2.1.1	Handling Missing Values (Imputation)	2
2.1.2	Normalization	4
2.1.3	Sampling	6
2.1.4	Binning	7
2.2	TODO Supervised Learning	8
2.3	TODO Unsupervised Learning	8
2.4	Reinforcement Learning	8
2.5	Steps In Implementing An AI Model.	8
2.5.1	Problem identification	8
2.5.2	Data Curation	9
2.5.3	2.1	9
2.5.4	Selection of AI models based on the data	9
2.5.5	Training and tuning the model - A train/test split or a train/validation/testing split.	9
2.5.6	Testing the developed model	9
2.5.7	Analysis of the results	10
2.5.8	Re-iterate as needed	10
2.5.9	Deploy model.	10
3	Data Imbalance	10
3.0.1	TODO Find what vintage means in churn prediction.	10
3.1	One Hot Encoding	10

1 Iris Data Classification

```
import pandas as pd
import numpy as np
from sklearn.model_selection import train_test_split

irisdata = pd.read_csv('iris.csv')

test, train = train_test_split(irisdata, train_size=0.8, test_size=0.2)

print(np.size(test))
print(np.size(train))
print(irisdata.describe())
```

2 Overview

2.1 Pre-Processing

2.1.1 Handling Missing Values (Imputation)

When the no. of missing values in a feature or on a whole in a dataset, is beyond a certain percentage. It might lead to wrong interpretations and might misguide the ML models. Hence it is essential to handle the missing values.

1. CREATING A DATAFRAME

```
import pandas as pd
import numpy as np

# Load the Titanic dataset
df = pd.read_csv('titanic.csv')

# Display the first few rows of the dataset
print("First few rows of the dataset:")
print(df.head())
```

This dataset is not complete, Cabin and Age have values that are unfilled. We can verify this here.

```
# Identify missing values
print("\nMissing values in each column:")
print(df.isnull().sum())
```

2. There are two main methods in dealing with missing values.

- (a) Dropping rows with missing values.
- (b) Filling the empty missing values with zeros.

```
# Method 1: Drop rows with missing values
df_dropped = df.dropna()
print("\n METHOD 1 Shape of dataset after dropping rows with missing values:", df_

# Method 2: Fill missing values with a specific value (e.g., 0)
df_filled_zeros = df.fillna(0)
print("\nMETHOD 2 Missing values filled with 0:")
print(df_filled_zeros.isnull().sum())
```

This isn't exactly ideal. Deleting the rows loses too much of the dataset, and filling with zeros does not work here when that might affect the correctness of the prediction. So here we replace the values with the mean for numerical values and mode for categorical values.

- (a) **TODO** Look into other methods of imputation

```
# Method 3: Fill missing values with the mean (for numerical columns)
df['Age'].fillna(df['Age'].mean(), inplace=True)
print("\nMETHOD 3 Missing values in 'Age' column after filling with mean:")
print(df['Age'].isnull().sum())

# Method 4: Fill missing values with the most frequent value (mode)
df['Embarked'].fillna(df['Embarked'].mode()[0], inplace=True)
print("\nMETHOD 4 Missing values in 'Embarked' column after filling with mode")
print(df['Embarked'].isnull().sum())
```

3. Forward fill and Backward Fill There are two better ways to fill the rows.

- Forward Fill - It iterates down the given data, and fills in missing values with the last value it saw.

- Backward Fill - it iterates up the given data, and fills in missing values with the last value it saw.

```
# Method 5: Forward fill method
df_ffill = df.fillna(method='ffill')
print("\nMethod 5 Missing values handled using forward fill method:")
print(df_ffill.isnull().sum())

# Method 6: Backward fill method
df_bfill = df.fillna(method='bfill')
print("\nMethod 6 Missing values handled using backward fill method:")
print(df_bfill.isnull().sum())
print("*****")
```

2.1.2 Normalization

Used for multiple numerical features in the dataset, which belong to different ranges. It would make sense to normalize the data to a particular range.

Machine learning models tend to give a higher weightage to numerical attributes which have a larger value.

The solution is to normalize. Normalization reduces a given numerical feature into a range that is easier to manage as well as equate with other numerical features.

1. Types Of Normalization

- MinMaxScaler - all data points are brought to the range [0, 1]
- Z-score - Data points are converted in such a way that the mean becomes 0 and the standard deviation is 1.
- LogScaler
- DecimalScaler - divides the number by a power of 10 until it is lesser than 1.

(a) NORMALISING A SET OF VALUES USING MIN MAX NORMALIZATION

```
import numpy as np
from sklearn.preprocessing import MinMaxScaler
```

```
# Example usage:
```

```
data = np.array([2, 5, 8, 11, 14]).reshape(-1, 1) # Reshape to 2D array for
```

```

# Initialize the MinMaxScaler
scaler = MinMaxScaler()

# Apply Min-Max normalization
normalized_data = scaler.fit_transform(data)

# Flatten the normalized data to 1D array
normalized_data = normalized_data.flatten()

print(normalized_data)

```

(b) NORMALISING A SET OF VALUES USING Z-SCORE NORMALIZATION

```

import numpy as np
from sklearn.preprocessing import StandardScaler

# Example usage:
data = np.array([2, 5, 8, 11, 14]).reshape(-1, 1) # Reshape to 2D array for

# Initialize the StandardScaler
scaler = StandardScaler()

# Apply Z-score normalization
normalized_data = scaler.fit_transform(data)

# Flatten the normalized data to 1D array
normalized_data = normalized_data.flatten()

print(normalized_data)

```

(c) NORMALIZING CERTAIN COLUMNS IN THE DATAFRAME

```

# Initialize the MinMaxScaler
from sklearn.preprocessing import MinMaxScaler
scaler = MinMaxScaler()

# List of columns to be normalized
columns_to_normalize = ['Age', 'Fare']

```

```

# Apply Min-Max normalization
df[columns_to_normalize] = scaler.fit_transform(df[columns_to_normalize])

print("\nDataFrame after Min-Max normalization:")
print(df)

```

2.1.3 Sampling

1. RANDOM SAMPLING Random sampling is used for when the dataset is hella large.

```

import random

# Sample data
population = list(range(1, 101)) # Population from 1 to 100
sample_size = 10 # Size of the sample

# Simple random sampling
sample = random.sample(population, sample_size)
print("Simple Random Sample:", sample)

```

2. STRATIFIED SAMPLING

```

import random

# Sample data with strata
strata_data = {
    'stratum1': [1, 2, 3, 4, 5],
    'stratum2': [6, 7, 8, 9, 10],
}

# Sample size per stratum
sample_size_per_stratum = 2

# Stratified sampling
sample = []
for stratum, data in strata_data.items():
    stratum_sample = random.sample(data, sample_size_per_stratum)
    sample.extend(stratum_sample)

print("Stratified Sample:", sample)

```

3. Systematic Sampling

```
# Sample data
data = list(range(1, 101)) # Data from 1 to 100
n = 5 # Every nth data point to be included in the sample

# Systematic sampling
sample = data[::n]
print("Systematic Sample:", sample)

import random

# Sample data with clusters
clusters = {
    'cluster1': [1, 2, 3],
    'cluster2': [4, 5, 6],
    'cluster3': [7, 8, 9],
}

# Number of clusters to sample
clusters_to_sample = 2

# Cluster sampling
selected_clusters = random.sample(list(clusters.keys()), clusters_to_sample)
print("chosen clusters ", selected_clusters)
sample = []
for cluster in selected_clusters:
    sample.extend(clusters[cluster])

print("Cluster Sample:", sample)
```

2.1.4 Binning

```
import pandas as pd

df = pd.read_csv('bollywood.csv')
budget_bins = [0, 10, 20, float('inf')] # Define your budget bins
budget_labels = ['Low Budget', 'Medium Budget', 'High Budget'] # Labels for the bins
df['BudgetBin'] = pd.cut(df['Budget'], bins=budget_bins, labels=budget_labels)
print(df.head(10))
```

```

collection_bins = [0, 20, 40, 60, float('inf')] # Define your collection bins
collection_labels = ['Low Collection', 'Medium Collection', 'High Collection', 'Very H

df['CollectionBin'] = pd.cut(df['BoxOfficeCollection'], bins=collection_bins, labels=c
df.head(10)

import matplotlib.pyplot as plt
budget_bin_counts = df['BudgetBin'].value_counts()
# Plot the data as a bar chart
plt.figure(figsize=(8, 6))
budget_bin_counts.plot(kind='bar', color='skyblue')
plt.title('Number of Movies in Each Budget Bin')
plt.xlabel('Budget Bin')
plt.ylabel('Number of Movies')
plt.xticks(rotation=45) # Rotate x-axis labels for better readability
plt.tight_layout()

```

2.2 TODO Supervised Learning

2.3 TODO Unsupervised Learning

2.4 Reinforcement Learning

This is a method used in game-based systems. It maps:

- A set of states
- A set of actions
- A set of rewards

And tries to take actions, to achieve a goal to get the reward. It receives the reward, when it achieves the goal, and receives a penalty upon failure.

These models maximise the cumulative reward.

2.5 Steps In Implementing An AI Model.

2.5.1 Problem identification

This is done by researching

- Experts in the field
- Personal experience

- Literature survey
- Data curation

2.5.2 Data Curation

- Data collection in person
- Public repos
- Private repos
- Simulated data
- Synthetic data

2.5.3 2.1

2.5.4 Selection of AI models based on the data

- Figure out whether the problem is a regression or a classification problem.
- Figure out the computational capacity
- Try various models for best fit.

2.5.5 Training and tuning the model - A train/test split or a train/validation/testing split.

- The data is separated out into training and testing.
- The training subset is passed onto the chosen AI model.
- Validation is done because it prevents overfitting.
- The model should generalize.

2.5.6 Testing the developed model

- Choose evaluation metrics based on the model.
 - Regression can involve MSPE, MSAE, R^2
- Test the data.

2.5.7 Analysis of the results

2.5.8 Re-iterate as needed

2.5.9 Deploy model.

3 Data Imbalance

We're doing churn prediction, this term means that it predicts how likely a customer is to not buy the product.

3.0.1 TODO Find what vintage means in churn prediction.

3.1 One Hot Encoding

This is used when we have categorical values spread into boolean values for their own category. If a given object is of a certain category, then the column of that category is true instead of giving it a numerical categorical value. This is better than using one column as a categorical value.