

Logistic Regression Model description

Logistic regression is a supervised machine learning algorithm used for binary classification problems, where the goal is to predict whether an instance belongs to one of two classes (e.g., "yes" or "no", "true" or "false", "spam" or "not spam"). Unlike linear regression, which predicts a continuous numerical value, logistic regression predicts the probability of an instance belonging to a particular class.

The logistic regression model can be expressed as:

$$P(y=1|x) = 1 / (1 + e^{(-z)})$$

Where:

- $P(y=1|x)$ is the probability that the dependent variable y is 1 (i.e., the instance belongs to the positive class), given the independent variable(s) x .
- $z = \alpha + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_n x_n$ is the linear combination of the independent variables, where:
 - α is the y-intercept (the value of the log-odds when all x 's are 0)
 - $\beta_1, \beta_2, \dots, \beta_n$ are the coefficients that represent the change in the log-odds for a unit change in the corresponding x

The logistic function ($1 / (1 + e^{(-z)})$) maps the linear combination z to a value between 0 and 1, which can be interpreted as the probability of the instance belonging to the positive class.

The key steps in building a logistic regression model are:

1. Data Preparation: Collect and preprocess the data, handling missing values, outliers, and feature engineering as needed.
2. Model Training: Fit the logistic regression model to the training data using techniques like Maximum Likelihood Estimation (MLE) to estimate the model parameters (α and β 's).
3. Model Evaluation: Assess the model's performance using metrics like accuracy, precision, recall, F1-score, and Area Under the Receiver Operating Characteristic (ROC-AUC) curve.
4. Model Interpretation: Analyze the model coefficients, odds ratios, and p-values to understand the relationship between the independent and dependent variables.
5. Model Deployment: Use the trained model to make predictions on new, unseen data.

Logistic regression is a powerful and widely-used technique for binary classification problems due to its simplicity, interpretability, and the availability of well-developed statistical methods for inference and diagnostics.

Creating a Dataset

Create a dataset of 50 records with the features "number of packets", "traffic volume", "repair cost", and "source ID". Generate random data with some relationship between the features and the target variable.

```
import numpy as np
import pandas as pd

# Generate random data
np.random.seed(42)
num_records = 50

# Generate feature values
num_packets = np.random.randint(100, 1000, num_records)
traffic_volume = np.random.randint(1000, 10000, num_records)
repair_cost = 2 * num_packets + 500 + np.random.normal(0, 50, num_records)
source_id = np.random.randint(1, 6, num_records)

# Generate target variable (binary classification)
target = (repair_cost > 1500).astype(int)

# Create the dataset
data = {'number of packets': num_packets,
        'traffic volume': traffic_volume,
        'repair cost': repair_cost,
        'source ID': source_id,
        'target': target}
df = pd.DataFrame(data)

# Print the dataset
print(df)
```

The output will be a pandas DataFrame with 50 records and the following features:

	number of packets	traffic volume	repair cost	source ID	target
0	444	6415	1373	1	1
1	627	7824	1753	3	1
2	675	5118	1850	1	1
3	199	5762	898	5	0
4	395	1413	1290	3	1
5	659	6186	1818	2	1
6	789	9885	2078	1	1
7	289	7648	1077	4	0
8	680	7701	1861	2	1

9	387	7630	1275	3	1
10	624	5713	1744	5	1
11	849	6330	2149	2	1
12	113	2265	726	4	0
13	910	9853	2270	1	1
14	713	2609	1923	3	1
15	595	5623	1699	5	1
16	873	7919	2181	2	1
17	623	2757	1740	4	1
18	465	7516	1464	1	1
19	134	3496	756	3	0
20	286	1793	1068	5	0
21	279	9366	1053	2	0
22	335	4684	1183	4	0
23	112	7947	720	1	0
24	660	6650	1819	3	1
25	393	2808	1285	5	1
26	786	3145	2073	2	1
27	467	5305	1469	4	1
28	756	4427	2027	1	1
29	789	9448	2078	3	1
30	199	5197	898	5	0
31	842	6568	2130	2	1
32	655	1331	1806	4	1
33	473	9831	1480	1	1
34	619	2157	1735	3	1
35	144	3602	771	5	0
36	408	2570	1327	2	1
37	870	4231	2177	4	1
38	677	7945	1855	1	1
39	192	7032	888	3	0
40	524	3954	1563	5	1
41	418	8499	1356	2	1
42	761	7868	2032	4	1
43	366	4613	1232	1	0
44	692	1980	1886	3	1
45	870	6649	2177	5	1
46	412	3893	1339	2	1
47	239	8646	975	4	0
48	370	6551	1241	1	1
49	789	9026	2078	3	1

\\

```
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
from sklearn.linear_model import LogisticRegression
from sklearn.model_selection import train_test_split
from sklearn.metrics import accuracy_score, precision_score, recall_score, f1_score, roc_auc_score

# Load the data
df = pd.DataFrame(data)

# Explore the data
print(df.head())
print(df.describe())

# Prepare the data
X = df[['number of packets', 'traffic volume', 'source ID']].values
y = df['target'].values

# Split the data into training and testing sets
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)

# Fit the logistic regression model
model = LogisticRegression()
model.fit(X_train, y_train)

# Make predictions
y_pred = model.predict(X_test)

# Evaluate the model
accuracy = accuracy_score(y_test, y_pred)
precision = precision_score(y_test, y_pred)
recall = recall_score(y_test, y_pred)
f1 = f1_score(y_test, y_pred)
roc_auc = roc_auc_score(y_test, y_pred)

print('Accuracy:', accuracy)
print('Precision:', precision)
print('Recall:', recall)
print('F1-score:', f1)
print('ROC-AUC:', roc_auc)

# Plot the ROC curve
from sklearn.metrics import roc_curve, auc
```

```
fpr, tpr, thresholds = roc_curve(y_test, model.predict_proba(X_test)[:, 1])
roc_auc = auc(fpr, tpr)

plt.figure(figsize=(8, 6))
plt.plot(fpr, tpr, color='darkorange', lw=2, label='ROC curve (area = %0.2f)' % roc_auc)
plt.plot([0, 1], [0, 1], color='navy', lw=2, linestyle='--')
plt.xlim([0.0, 1.0])
plt.ylim([0.0, 1.05])
plt.xlabel('False Positive Rate')
plt.ylabel('True Positive Rate')
plt.title('Receiver Operating Characteristic (ROC) Curve')
plt.legend(loc="lower right")
plt.show()
```

In this example, we first load the dataset we created earlier. We then prepare the data by selecting the features "number of packets", "traffic volume", and "source ID" as the independent variables (X) and the "target" variable as the dependent variable (y).

Split the data into training and testing sets using the `train_test_split` function from scikit-learn. We then fit the logistic regression model using the `LogisticRegression` class from scikit-learn and make predictions on the test set.

Evaluate the model's performance by calculating the accuracy, precision, recall, F1-score, and ROC-AUC score. We also plot the ROC curve to visualize the model's performance.