

## Data Cleaning, EDA, and Baseline Model

```
In [1]: # Imports
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LinearRegression
from sklearn.metrics import mean_squared_error, r2_score
from sklearn.preprocessing import StandardScaler
from sklearn.impute import SimpleImputer
from sklearn.pipeline import Pipeline
from sklearn.compose import ColumnTransformer
import warnings
warnings.filterwarnings('ignore')
```

```
In [2]: # Step 1: Import required libraries
import pandas as pd
```

```
In [8]: # Step 3: Load the parquet files
train_df = pd.read_parquet(train_file)
test_df = pd.read_parquet(test_file)

print("Train shape:", train_df.shape)
print("Test shape:", test_df.shape)

train_df.head()

Train shape: (415, 73)
```

© 1996 by J. B. Lippincott Company

url	call_transcript	VWAP	exchangeCountry	securityType	CIK	name	securityID	incorporationCountry	exchangeName	..
ticker										
Fair Isaac										

5 rows x 73 columns

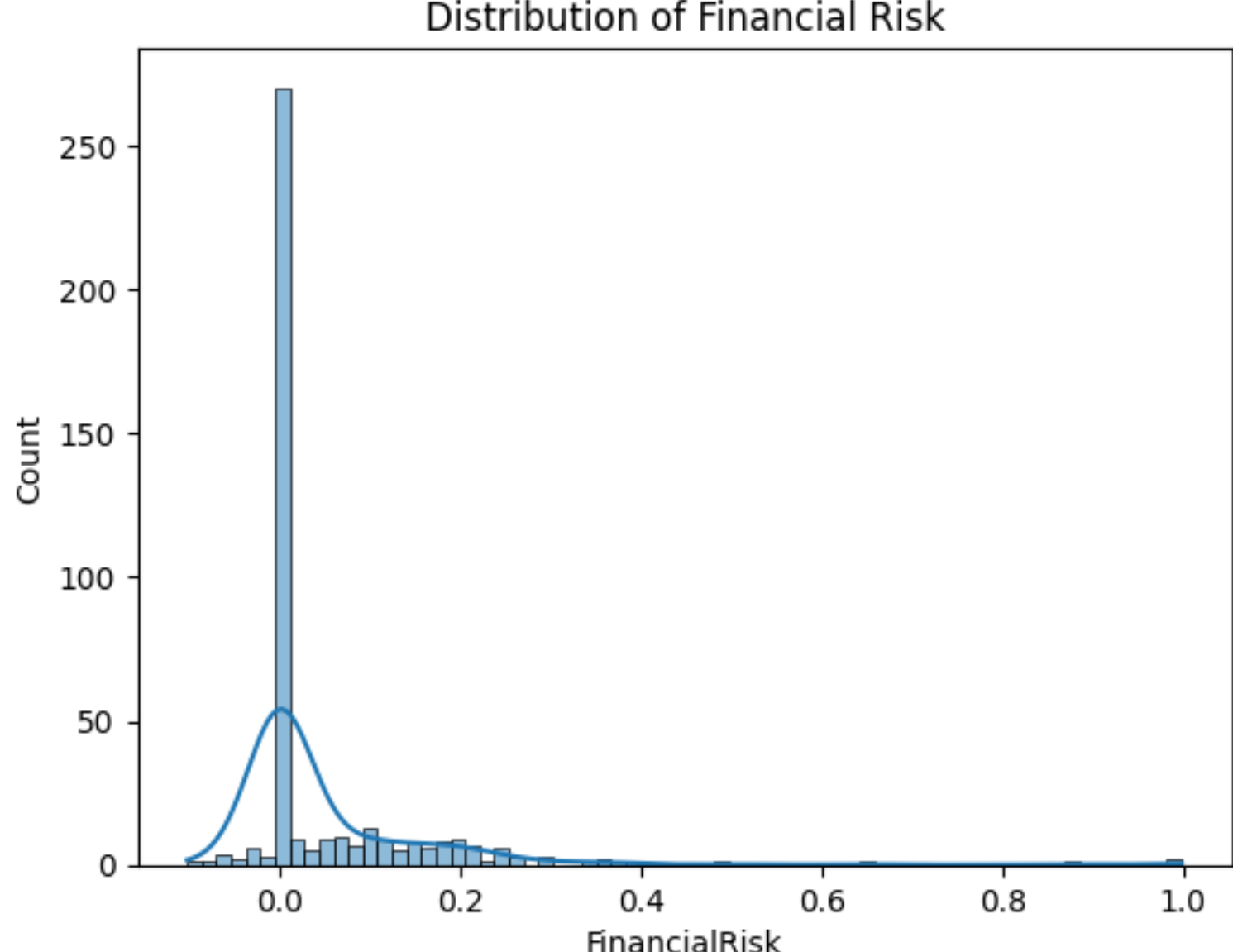
```
[9]: # Drop irrelevant features
    irrelevant_cols = ['url', 'exchangeCountry', 'CIK', 'securityID',
                      'incorporationCountry', 'exchangeName', 'exchangeID', 'businessDescription']
    df.drop(columns=irrelevant_cols, inplace=True, errors='ignore')

    # Drop columns with more than 20% missing values
    missing_threshold = 0.2
    df.dropna(axis=1, thresh=(1 - missing_threshold) * len(df), inplace=True)

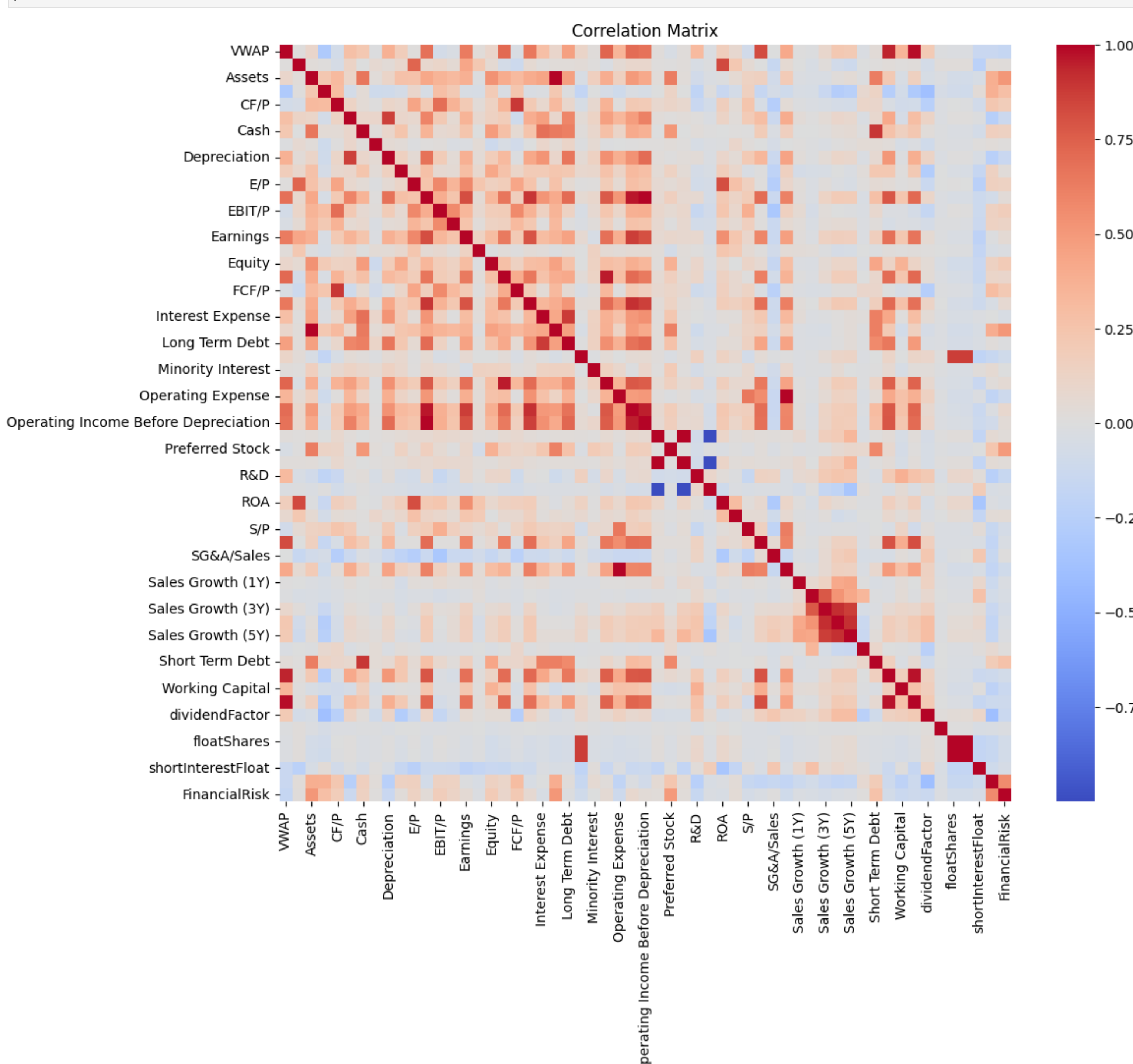
    # Separate numerical and categorical columns
    num_cols = df.select_dtypes(include=['float64', 'int64']).columns.tolist()
    cat_cols = df.select_dtypes(include=['object']).columns.tolist()
```

```
In [10]: # Check distribution of target
sns.histplot(df['FinancialRisk'], kde=True)
plt.title('Distribution of Financial Risk')
plt.show()

# Log-transform target
df['FinancialRisk_log'] = np.log1p(df['FinancialRisk'])
```



```
[11]: # Correlation matrix
plt.figure(figsize=(12, 10))
corr = df[num_cols].corr()
sns.heatmap(corr, cmap='coolwarm', center=0)
plt.title('Correlation Matrix')
```



```
[31]: # Baseline Model: Linear Regression
selected_features = ['Assets', 'Cash', 'Debt/Equity', 'B/P', 'SG&A/Sales', 'Long Liabilities']
print('Columns in DataFrame:', df.columns.tolist())
missing_features = [col for col in selected_features if col not in df.columns]
if missing_features:
    print('The following features are missing from the data:', missing_features)
else:
    print('All selected features are present in the data.')
if 'FinancialRisk_log' not in df.columns:
    print('FinancialRisk_log column is missing. Please make sure the previous cell has been executed.')
else:
    print('FinancialRisk_log column exists.')
print('Missing value count for selected features:')
print(df[selected_features].isnull().sum())
# Convert selected features to numeric type
for col in selected_features:
    df[col] = pd.to_numeric(df[col], errors='coerce')
# Remove all rows with missing values in selected features or target
df_clean = df.dropna(subset=selected_features + ['FinancialRisk_Log']).copy()
print('Shape after dropping missing:', df_clean.shape)
# Continue with modeling
X = df_clean[selected_features]
y = df_clean['FinancialRisk_Log']
from sklearn.model_selection import train_test_split
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
from sklearn.impute import SimpleImputer
from sklearn.preprocessing import StandardScaler
from sklearn.pipeline import Pipeline
from sklearn.compose import ColumnTransformer
numeric_transformer = Pipeline(steps=[
    ('imputer', SimpleImputer(strategy='median')),
    ('scaler', StandardScaler())
])
preprocessor = ColumnTransformer(
    transformers=[
        ('num', numeric_transformer, X.columns)
    ])
from sklearn.linear_model import LinearRegression
model = Pipeline(steps=[('preprocessor', preprocessor),
    ('regressor', LinearRegression())])
model.fit(X_train, y_train)
y_pred = model.predict(X_test)
from sklearn.metrics import mean_squared_error, r2_score
r2 = r2_score(y_test, y_pred)
rmse = mean_squared_error(y_test, y_pred) ** 0.5
print(f'R^2 Score: {r2:.2f}')
print(f'RMSE: {rmse:.2f}')
```

eciation', 'Dividend', 'E/P', 'EBIT', 'EBIT/P', 'EBIT/TEV', 'Earnings', 'Earnings Growth (1Y)', 'Equity', 'FCF', 'FCF/P', 'Income Tax', 'Interest Expense', 'Long Liability', 'Long Term Debt', 'Market Cap', 'Minority Interest', 'Operating Cash Flow', 'Operating Expense', 'Operating Income', 'Operating Income Before Depreciation', 'Operating Margin', 'Preferred Stock', 'Profit Margin', 'R&D', 'R&D/Sales', 'ROA', 'ROE', 'S/P', 'SG&A', 'SG&A/Sales', 'Sales', 'Sales Growth (1Y)', 'Sales Growth (2 Y)', 'Sales Growth (3Y)', 'Sales Growth (4Y)', 'Sales Growth (5Y)', 'Sales Variability', 'Short Term Debt', 'TEV', 'Working Capital', 'close', 'dividendFactor', 'fiscalYr', 'floatShares', 'outstandingShares', 'shortInterestFloat', 'FinancialSector', 'FinancialRisk', 'FinancialRisk\_log']

FinancialRisk log column exists

Financial Risk Log - 2014

Missing value count to

Assets 0

ASSOCIATES	
Associates	

Cash	0
------	---

Debt/Equity 18

B/P 0

SCCA/Coloc 7

30Q/A/ 3a LES 7

Long Liabilities 0

```
dtype: int64
```

Shape after dropping m

D3.60000000 0.00

R<sup>2</sup> Score: -0.08

RMSF: 0.08

**TABLE 1**