

```

1 # Author: Sai Sri Surya Bandaru
2 # CSCE 5222 - Project 2
3 # Pancreatic Cancer Survival Analysis
4 # April 2024

1 # Imports
2 import pandas as pd
3 import numpy as np
4 import matplotlib.pyplot as plt
5 import seaborn as sns
6 from pathlib import Path
7 from sklearn.model_selection import train_test_split
8 from sklearn.metrics import mean_absolute_error, mean_squared_error
9 from sklearn.preprocessing import StandardScaler, OneHotEncoder
10 from sklearn.compose import ColumnTransformer
11 from sklearn.pipeline import Pipeline
12 from sklearn.linear_model import ElasticNet
13 from sklearn.svm import SVR
14 from sklearn.tree import DecisionTreeRegressor
15 from lightgbm import LGBMRegressor
16 from xgboost import XGBRegressor
17 import joblib, os

1 # Setup paths
2
3 # Get the project root directory
4 BASE_DIR = Path(os.getcwd()).parent
5 DATA_DIR = BASE_DIR / "data"
6 MODEL_DIR = BASE_DIR / "models"
7
8 # Print paths to verify
9 print(f"Base directory: {BASE_DIR}")
10 print(f>Data directory: {DATA_DIR}")
11 print(f"Model directory: {MODEL_DIR}")
12
13 # Verify directories exist
14 DATA_DIR.mkdir(parents=True, exist_ok=True)
15 MODEL_DIR.mkdir(parents=True, exist_ok=True)
16
17 # Load data function
18 def load_data(use_vae=False):
19     file = "vae_augmented_data.csv" if use_vae else "synthetic_data.csv"
20     return pd.read_csv(DATA_DIR / file)
21
22 # Build preprocessing pipeline
23 def build_preprocessor(df):
24     num_cols = df.select_dtypes(include=["float64", "int64"]).columns.difference(["SurvivalMonths"])
25     cat_cols = df.select_dtypes(include=["object"]).columns
26     preproc = ColumnTransformer(
27         transformers=[
28             ("num", StandardScaler(), num_cols),
29             ("cat", OneHotEncoder(handle_unknown="ignore"), cat_cols),
30         ]
31     )
32     return preproc

```

Base directory: c:\Users\vckkr\OneDrive - Chaitanya UNT\OneDrive - UNT System\Documents\3rd-Party-Assignments\2025\04-April\CSCE 5222
Data directory: c:\Users\vckkr\OneDrive - Chaitanya UNT\OneDrive - UNT System\Documents\3rd-Party-Assignments\2025\04-April\CSCE 5222\data
Model directory: c:\Users\vckkr\OneDrive - Chaitanya UNT\OneDrive - UNT System\Documents\3rd-Party-Assignments\2025\04-April\CSCE 5222\models

```

1 # Plot feature distributions
2 def plot_feature_distributions(df, title):
3     num_cols = df.select_dtypes(include=["float64", "int64"]).columns.difference(["SurvivalMonths"])
4     n_cols = 3
5     n_rows = (len(num_cols) + n_cols - 1) // n_cols
6
7     plt.figure(figsize=(15, 5*n_rows))
8     for i, col in enumerate(num_cols, 1):
9         plt.subplot(n_rows, n_cols, i)
10         sns.histplot(data=df, x=col, kde=True)
11         plt.title(f"{col} Distribution")
12     plt.suptitle(title)
13     plt.tight_layout()
14     plt.show()
15
16 # Plot survival by stage
17 def plot_survival_by_stage(df, title):
18     plt.figure(figsize=(10, 6))

```

```

19 sns.boxplot(data=df, x="Stage", y="SurvivalMonths")
20 plt.title(title)
21 plt.xlabel("Cancer Stage")
22 plt.ylabel("Survival Months")
23 plt.show()
24
25 # Plot correlation heatmap
26 def plot_correlation(df, title):
27     plt.figure(figsize=(12, 8))
28     corr = df.select_dtypes(include=["float64", "int64"]).corr()
29     sns.heatmap(corr, annot=True, cmap="coolwarm", center=0)
30     plt.title(title)
31     plt.tight_layout()
32     plt.show()

1 # Train models and get scores
2 def train_models(df, use_vae=False):
3     X = df.drop(columns="SurvivalMonths")
4     y = df["SurvivalMonths"]
5
6     Xtr, Xte, ytr, yte = train_test_split(
7         X, y, test_size=0.2, random_state=42, stratify=df["Stage"]
8     )
9
10    pre = build_preprocessor(df)
11
12    # Define models
13    models = {
14        "ElasticNet": ElasticNet(alpha=0.1, l1_ratio=0.5, random_state=42),
15        "SVR": SVR(C=10, gamma="scale"),
16        "DecisionTree": DecisionTreeRegressor(max_depth=5, random_state=42),
17        "LightGBM": LGBMRegressor(n_estimators=200, learning_rate=0.05),
18        "XGBoost": XGBRegressor(n_estimators=300, learning_rate=0.05, max_depth=4),
19    }
20
21    scores = {}
22    for name, est in models.items():
23        pipe = Pipeline([("prep", pre), ("model", est)])
24        pipe.fit(Xtr, ytr)
25
26        y_pred = pipe.predict(Xte)
27        scores[name] = {
28            "MAE": mean_absolute_error(yte, y_pred),
29            "MSE": mean_squared_error(yte, y_pred),
30        }
31
32    # Save model
33    out_dir = MODEL_DIR / "vae" if use_vae else MODEL_DIR
34    out_dir.mkdir(parents=True, exist_ok=True)
35    joblib.dump(pipe, out_dir / f"{name}_model.pkl")
36
37    return pd.DataFrame(scores).T

1 # First, check if we need to generate the data
2 from src.synthetic_data import generate_original_data, generate_synthetic_patients
3
4 # Generate data if it doesn't exist
5 if not (DATA_DIR / "synthetic_data.csv").exists():
6     print("Generating original dataset...")
7     original_df = generate_original_data(1000)
8     original_df.to_csv(DATA_DIR / "synthetic_data.csv", index=False)
9     print("Original dataset saved")
10
11 if not (DATA_DIR / "vae_augmented_data.csv").exists():
12     print("Generating VAE-augmented dataset...")
13     vae_df = generate_synthetic_patients(original_df, 5000)
14     vae_df.to_csv(DATA_DIR / "vae_augmented_data.csv", index=False)
15     print("VAE-augmented dataset saved")
16
17 # Now load the data
18 print("Loading datasets...")
19 original_df = load_data(use_vae=False)
20 vae_df = load_data(use_vae=True)
21
22 # Run analysis
23 print("\nTraining on original data...")
24 original_scores = train_models(original_df, use_vae=False)
25 print("\nOriginal data scores:")
26 print(original_scores)
27
28 print("\nTraining on VAE data...")

```

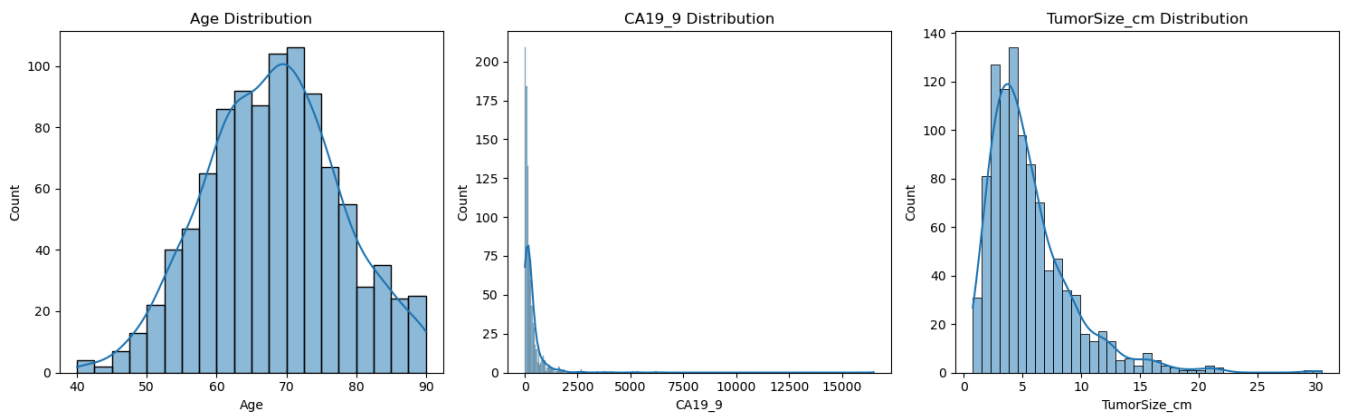
```
29 vae_scores = train_models(vae_df, use_vae=True)
30 print("\nVAE data scores:")
31 print(vae_scores)
```

Generating original dataset...

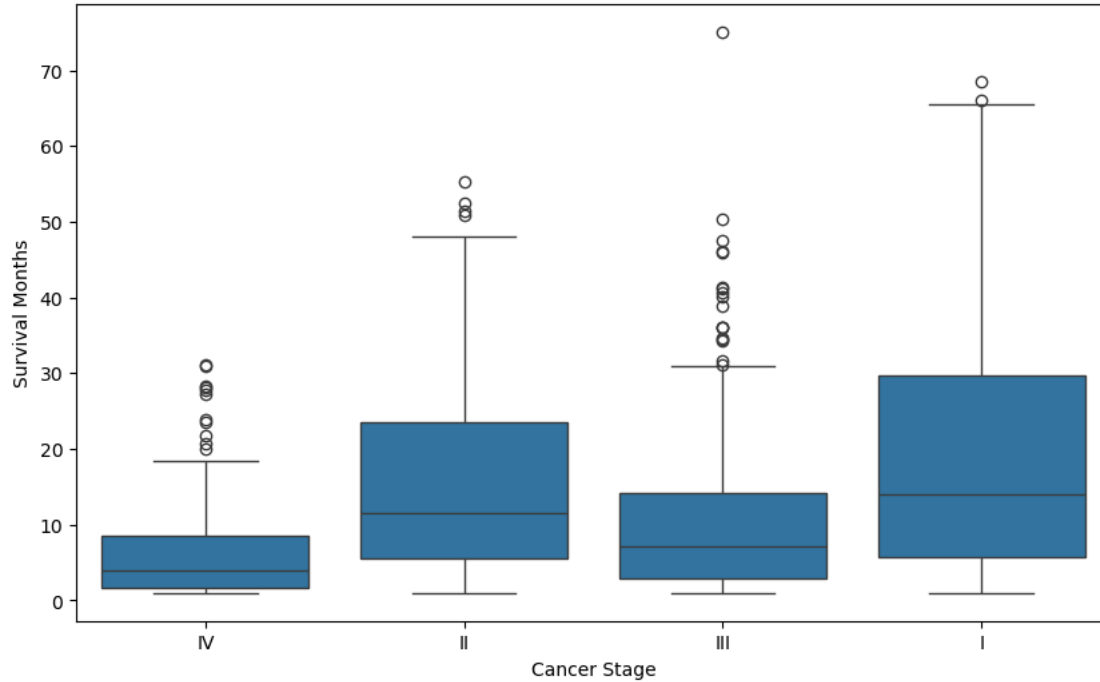
```
1 # Run visualizations for original data
2 print("Original Data Visualizations:")
3 plot_feature_distributions(original_df, "Original Data Feature Distributions")
4 plot_survival_by_stage(original_df, "Survival by Stage (Original Data)")
5 plot_correlation(original_df, "Feature Correlations (Original Data)")
6
7 # Run visualizations for VAE data
8 print("\nVAE-Augmented Data Visualizations:")
9 plot_feature_distributions(vae_df, "VAE Data Feature Distributions")
10 plot_survival_by_stage(vae_df, "Survival by Stage (VAE Data)")
11 plot_correlation(vae_df, "Feature Correlations (VAE Data)")
```

Original Data Visualizations:

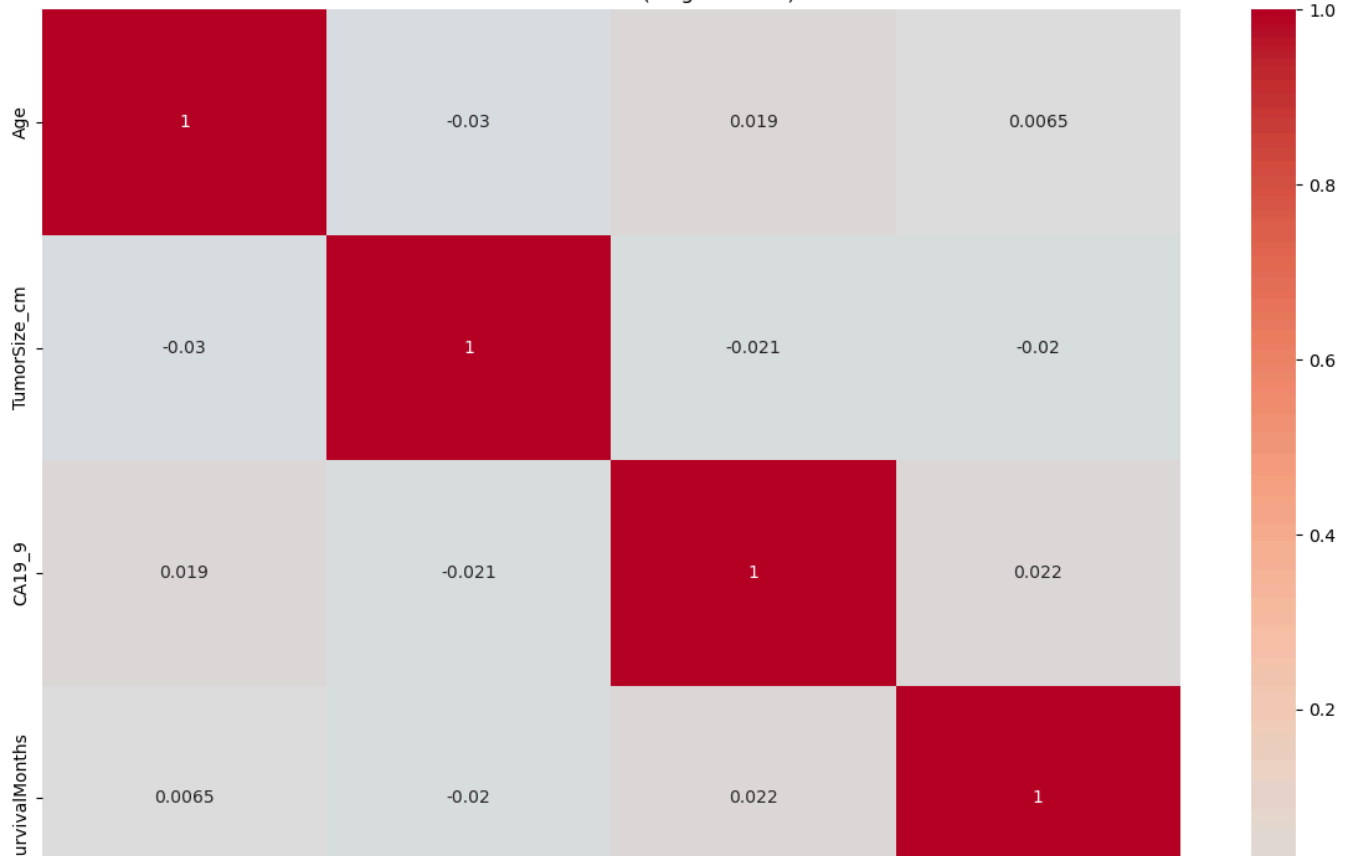
Original Data Feature Distributions



Survival by Stage (Original Data)

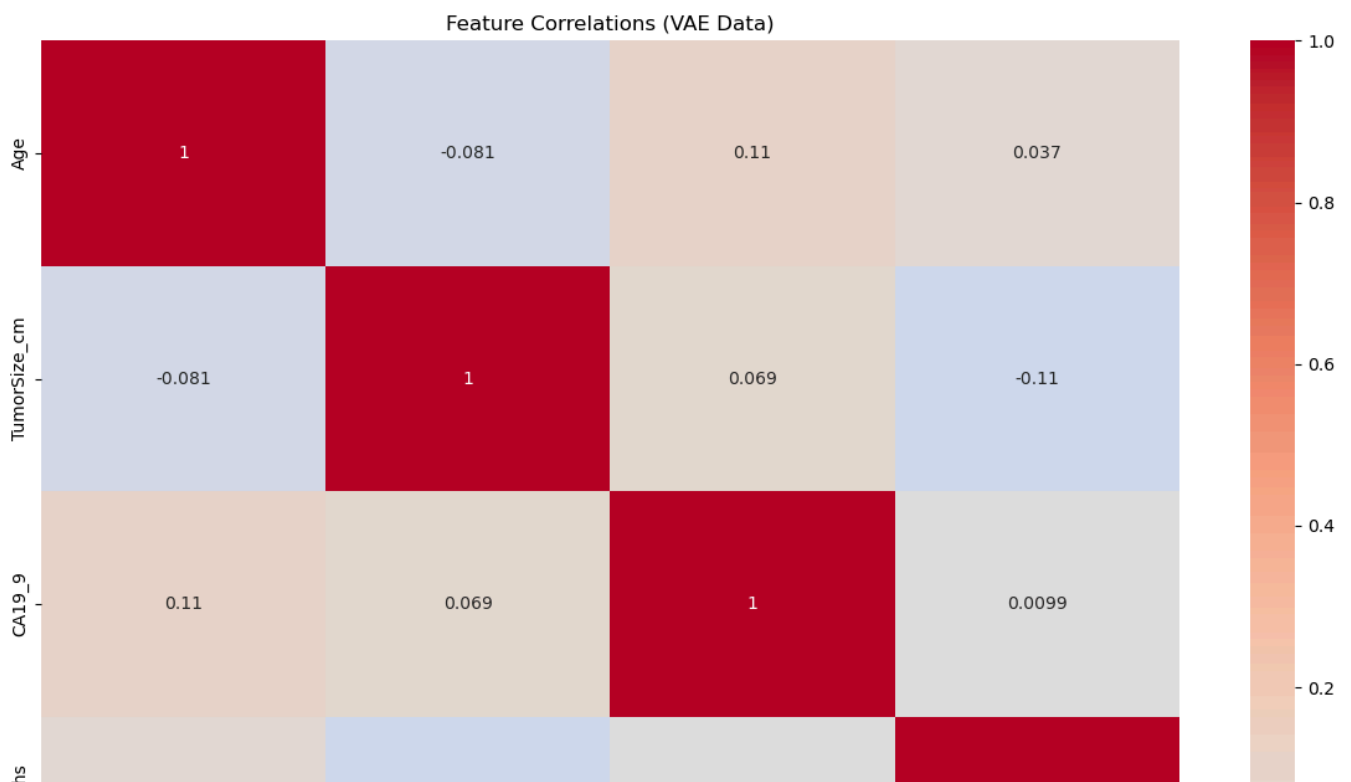
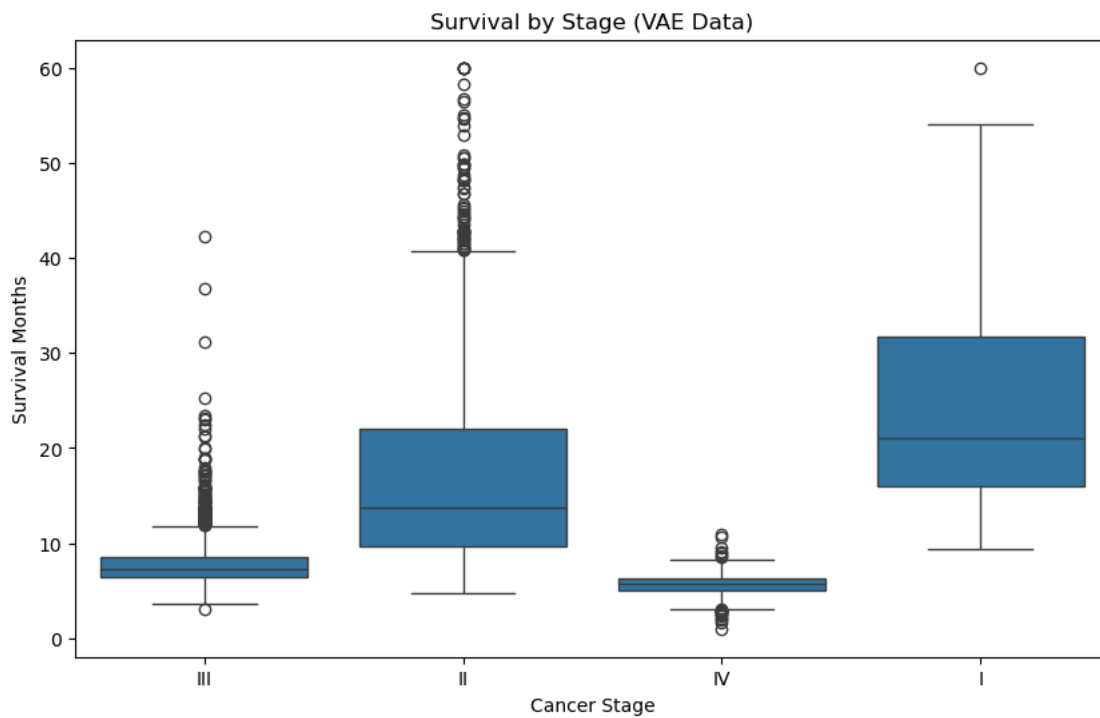
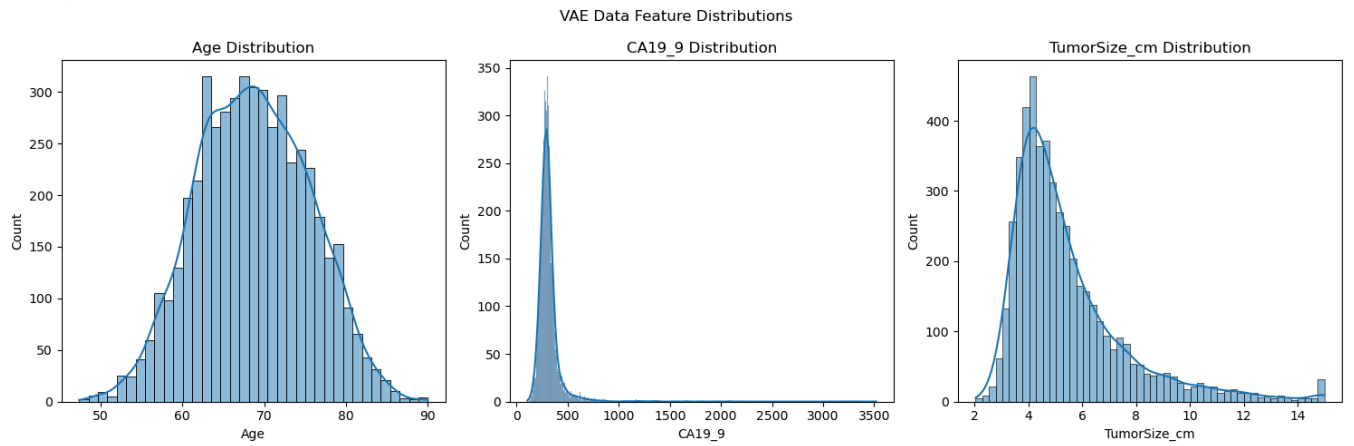


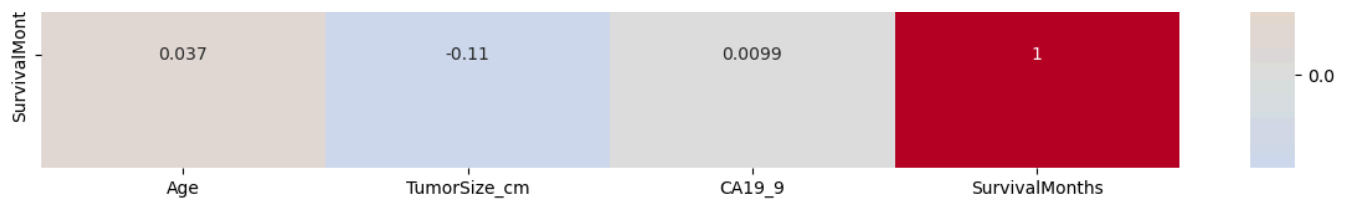
Feature Correlations (Original Data)





VAE-Augmented Data Visualizations:

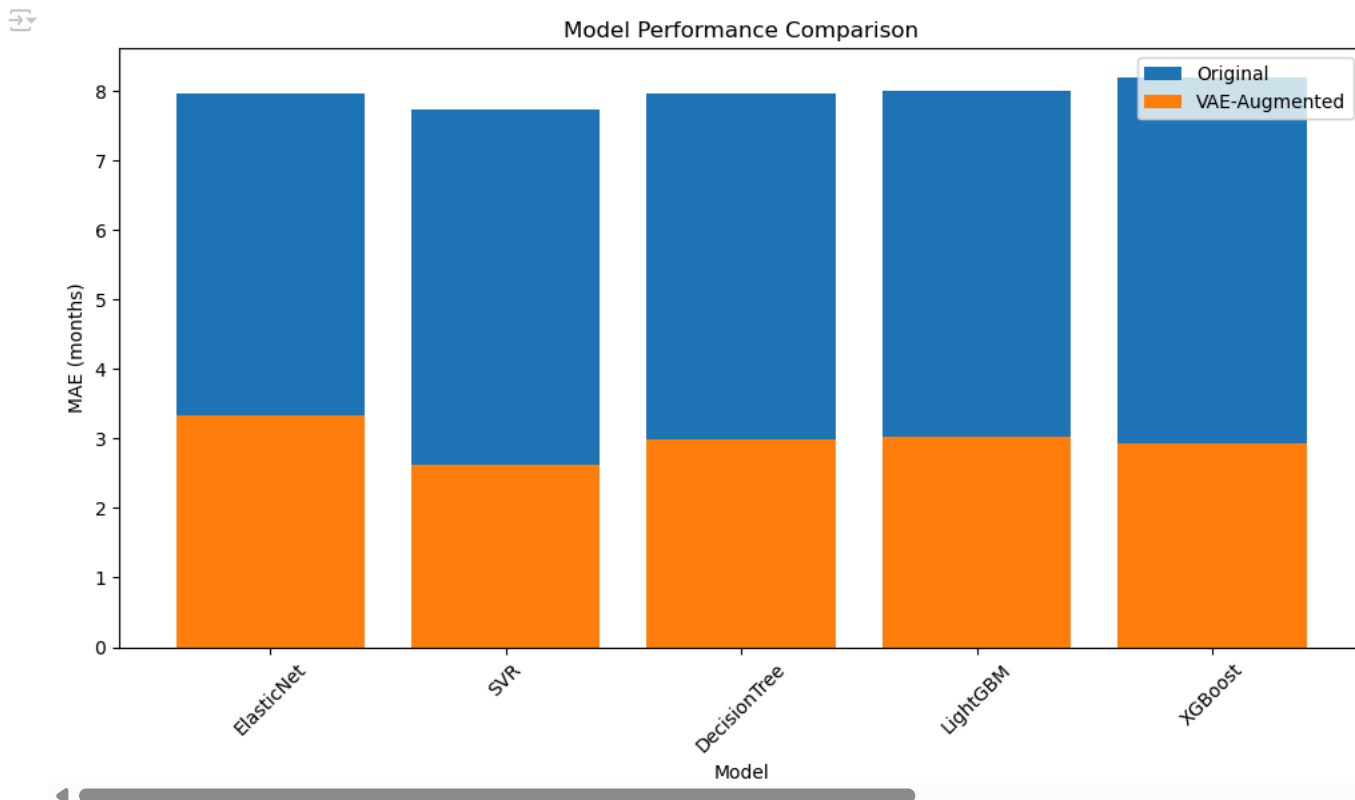




```

1 # Plot results
2 plt.figure(figsize=(10, 6))
3 plt.bar(original_scores.index, original_scores["MAE"], label="Original")
4 plt.bar(vae_scores.index, vae_scores["MAE"], label="VAE-Augmented")
5 plt.title("Model Performance Comparison")
6 plt.xlabel("Model")
7 plt.ylabel("MAE (months)")
8 plt.legend()
9 plt.xticks(rotation=45)
10 plt.tight_layout()
11 plt.show()

```



```

1 # Detailed evaluation metrics
2 def print_detailed_metrics(scores_df, title):
3     print(f"\n{title}")
4     print("-" * 50)
5     print("Mean Absolute Error (MAE):")
6     print(scores_df["MAE"].sort_values())
7     print("\nMean Squared Error (MSE):")
8     print(scores_df["MSE"].sort_values())
9     print("\nBest Model by MAE:", scores_df["MAE"].idxmin())
10    print("Best Model by MSE:", scores_df["MSE"].idxmin())
11
12 # Print metrics for both datasets
13 print_detailed_metrics(original_scores, "Original Data Model Performance")
14 print_detailed_metrics(vae_scores, "VAE-Augmented Data Model Performance")
15
16 # Calculate improvement percentages
17 mae_improvement = ((original_scores["MAE"] - vae_scores["MAE"]) / original_scores["MAE"] * 100)
18 mse_improvement = ((original_scores["MSE"] - vae_scores["MSE"]) / original_scores["MSE"] * 100)
19
20 print("\nImprovement with VAE-Augmented Data:")
21 print("-" * 50)
22 print("MAE Improvement (%):")
23 print(mae_improvement.sort_values(ascending=False))
24 print("\nMSE Improvement (%):")
25 print(mse_improvement.sort_values(ascending=False))

```

```

Original Data Model Performance
-----
Mean Absolute Error (MAE):
SVR          7.738537
ElasticNet   7.964600
DecisionTree 7.966413
LightGBM     8.009767
XGBoost      8.202568
Name: MAE, dtype: float64

Mean Squared Error (MSE):

```

```
ElasticNet      118.423290
DecisionTree    119.455227
LightGBM        121.171024
XGBoost         129.234579
SVR             131.156999
Name: MSE, dtype: float64
```

```
Best Model by MAE: SVR
Best Model by MSE: ElasticNet
```

VAE-Augmented Data Model Performance

Mean Absolute Error (MAE):

```
SVR             2.629438
XGBoost         2.931231
DecisionTree    2.989129
LightGBM        3.023969
ElasticNet      3.331492
Name: MAE, dtype: float64
```

Mean Squared Error (MSE):

```
XGBoost         29.087186
LightGBM        31.146668
DecisionTree    31.189531
SVR             31.434395
ElasticNet      33.328017
Name: MSE, dtype: float64
```

```
Best Model by MAE: SVR
Best Model by MSE: XGBoost
```

Improvement with VAE-Augmented Data:

MAE Improvement (%):

```
SVR             66.021506
XGBoost         64.264477
DecisionTree    62.478362
LightGBM        62.246479
ElasticNet      58.171258
Name: MAE, dtype: float64
```

MSE Improvement (%):

```
XGBoost         77.492722
SVR             76.033003
LightGBM        74.295283
DecisionTree    73.890192
```

```
1 # Plot model performance improvement
2 plt.figure(figsize=(12, 6))
3 plt.subplot(1, 2, 1)
4 mae_improvement.sort_values().plot(kind="barh", color="skyblue")
5 plt.title("MAE Improvement with VAE (%)")
6 plt.xlabel("Improvement Percentage")
7
8 plt.subplot(1, 2, 2)
9 mse_improvement.sort_values().plot(kind="barh", color="lightgreen")
10 plt.title("MSE Improvement with VAE (%)")
11 plt.xlabel("Improvement Percentage")
12
13 plt.tight_layout()
14 plt.show()
```




```

1 # Compare best models with and without VAE
2 print("Model Performance Comparison")
3 print("=" * 50)
4
5 # Find best models
6 best_original = original_scores["MAE"].idxmin()
7 best_vae = vae_scores["MAE"].idxmin()
8
9 print(f"\nBest Model (Original Data): {best_original}")
10 print(f"MAE: {original_scores.loc[best_original, 'MAE']:.2f} months")
11 print(f"MSE: {original_scores.loc[best_original, 'MSE']:.2f}")
12
13 print(f"\nBest Model (VAE Data): {best_vae}")
14 print(f"MAE: {vae_scores.loc[best_vae, 'MAE']:.2f} months")
15 print(f"MSE: {vae_scores.loc[best_vae, 'MSE']:.2f}")
16
17 # Calculate improvement
18 mae_improvement = ((original_scores.loc[best_original, "MAE"] - vae_scores.loc[best_vae, "MAE"])
19                    / original_scores.loc[best_original, "MAE"] * 100)
20
21 print(f"\nImprovement with VAE: {mae_improvement:.1f}% reduction in MAE")
22
23 # Visual comparison
24 plt.figure(figsize=(10, 6))
25 models = [best_original, best_vae]
26 mae_values = [original_scores.loc[best_original, "MAE"], vae_scores.loc[best_vae, "MAE"]]
27 colors = ['skyblue', 'lightgreen']
28 labels = ['Original Data', 'VAE-Augmented']
29
30 plt.bar(models, mae_values, color=colors, label=labels)
31 plt.title("Best Model Performance Comparison")
32 plt.ylabel("MAE (months)")
33 plt.grid(axis='y', linestyle='--', alpha=0.7)
34
35 # Add value labels on top of bars
36 for i, v in enumerate(mae_values):
37     plt.text(i, v, f'{v:.2f}', ha='center', va='bottom')
38
39 plt.legend()

```