# MIMIC-IV Mortality Prediction Tutorial

A step-by-step guide using clinical data and machine learning

# Introduction

- Predictive Analytics in Healthcare
  - Leveraging large-scale EHR data to predict patient outcomes
  - MIMIC-IV: Publicly available ICU dataset with structured and unstructured data

# Dataset Overview

- Key MIMIC-IV Tables:
  - ADMISSIONS
  - PATIENTS
  - ICUSTAYS
  - DIAGNOSES_ICD

# Goal

- Predict in-hospital mortality at the time of ICU admission using EHR data

# Import Libraries

- pandas

- numpy

- scikit-learn

- tensorflow / PyTorch

- transformers (ClinicalBERT)

# Data Loading

- Load data from MIMIC IV BigQuery into DataFrames



Length of Stay, Mortality and Gender



Admissions Data



Diagnoses code categories

# Feature Engineering

— Demographics: age, gender, ethnicity

— Admission info: type, source

— Clinical measurements: lab results, vitals

— Diagnosis codes grouped by ICD-9 categories

```python
# Drop datetime and 'gender' columns for modeling
datetime_cols = ['admittime', 'dischtime', 'deathtime']
# Drop subject_id, hadm_id, datetime columns, and consolidate diagnosis_x and diagnosis_y columns
drop_cols = ['subject_id', 'hadm_id'] + datetime_cols

# Consolidate diagnosis_x and diagnosis_y columns by summing them
diag_categories = [
    'blood', 'circulatory', 'congenital', 'digestive', 'endocrine', 'genitourinary',
    'infectious', 'injury', 'mental', 'misc', 'muscular', 'neoplasms', 'nervous',
    'pregnancy', 'prenatal', 'respiratory', 'skin'
]
for cat in diag_categories:
    x_col = f'{cat}_x'
    y_col = f'{cat}_y'
    if x_col in X_train.columns and y_col in X_train.columns:
        X_train[f'{cat}'] = X_train[x_col].fillna(0) + X_train[y_col].fillna(0)
        X_test[f'{cat}'] = X_test[x_col].fillna(0) + X_test[y_col].fillna(0)

# Drop the original diagnosis_x and diagnosis_y columns
diag_x_cols = [f'{cat}_x' for cat in diag_categories if f'{cat}_x' in X_train.columns]
diag_y_cols = [f'{cat}_y' for cat in diag_categories if f'{cat}_y' in X_train.columns]
X_train_clean = X_train.drop(columns=drop_cols + diag_x_cols + diag_y_cols + ['gender'])
X_test_clean = X_test.drop(columns=['gender'] + diag_x_cols + diag_y_cols + drop_cols)

# Impute missing values with column mean
imputer = SimpleImputer(strategy='mean')
# Ensure dense output in case imputer returns a sparse matrix
X_train_imputed = imputer.fit_transform(X_train_clean)
if hasattr(X_train_imputed, "toarray"): (function) toarray: Unknown
    X_train_imputed = X_train_imputed.toarray()
X_train_clean = pd.DataFrame(X_train_imputed, columns=X_train_clean.columns, index=X_train_clean.index)

X_test_imputed = imputer.transform(X_test_clean)
if hasattr(X_test_imputed, "toarray"):
```



Distribution by Gender



Distribution of Age

# Modeling

- • Train classifiers: Logistic Regression, Random Forest, XGBoost and SVM

- • Fine-tune ClinicalBERT for text classification

- • Split data: 80% train, 20% test

```python
from sklearn.impute import SimpleImputer


# Fine-tuned models with parameters
models = {
    'Logistic Regression': LogisticRegression(max_iter=1000, solver='lbfgs'),
    'Random Forest': RandomForestClassifier(n_estimators=100, max_depth=7),
    'XGBoost': XGBClassifier(use_label_encoder=False, eval_metric='logloss', max_depth=7, learning_rate=0.1),
    'SVM': SVC(probability=True, kernel='rbf', C=1.0)
}

results = {}
for name, model in models.items():
    model.fit(X_train_clean, y_train)
    y_pred = model.predict(X_test_clean)
    results[name] = {'model': model, 'y_pred': y_pred}
    print(f'{name} training and prediction done.')
✓ 18s
```

# Evaluation

— R^2 and RMSE Score Comparison for various models

— Confusion matrix analysis

— Feature Importance

```python
# Evaluate models using R² and RMSE
for name, result in results.items():
    y_pred = result['y_pred']
    r2 = r2_score(y_test, y_pred)
    rmse = np.sqrt(mean_squared_error(y_test, y_pred))
    print(f'{name}: R² = {r2:.3f}, RMSE = {rmse:.3f}')
    results[name]['r2'] = r2
    results[name]['rmse'] = rmse

# Plot comparison of models for R2 and RMSE
model_names = list(results.keys())
r2_scores = [results[m]['r2'] for m in model_names]
rmse_scores = [results[m]['rmse'] for m in model_names]

plt.figure(figsize=(10,4))
plt.subplot(1,2,1)
sns.barplot(x=model_names, y=r2_scores)
plt.title('R² Comparison')
plt.ylabel('R²')

plt.subplot(1,2,2)
sns.barplot(x=model_names, y=rmse_scores)
plt.title('RMSE Comparison')
plt.ylabel('RMSE')
plt.tight_layout()
plt.show()
```
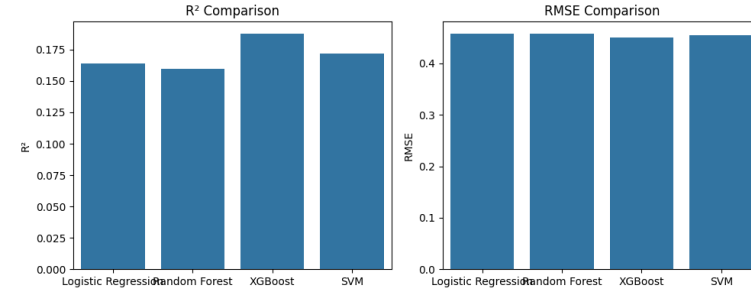


Table 1

| Logistic Regression: $R^2$ = 0.164 | RMSE = 0.457 |
|---|---|
| Random Forest: $R^2$ = 0.160 | RMSE = 0.458 |
| XGBoost: $R^2$ = 0.188 | RMSE = 0.451 |
| SVM: $R^2$ = 0.172 | RMSE = 0.455 |

# Evaluation - Logistics Regression



```python
# Confusion matrix for each model
from sklearn.metrics import confusion_matrix, ConfusionMatrixDisplay
from sklearn.metrics import accuracy_score, precision_score, recall_score, f1_score

for name, result in results.items():
    y_pred = result['y_pred']
    acc = accuracy_score(y_test, y_pred)
    prec = precision_score(y_test, y_pred)
    rec = recall_score(y_test, y_pred)
    f1 = f1_score(y_test, y_pred)
    print(f"{name}: Accuracy={acc:.3f}, Precision={prec:.3f}, Recall={rec:.3f}, F1={f1:.3f}")
    cm = confusion_matrix(y_test, y_pred)
    disp = ConfusionMatrixDisplay(confusion_matrix=cm, display_labels=["No Mortality (0)", "Mortality (1)"])
    disp.plot(cmap='Blues')
    plt.title(f'Confusion Matrix - {name}')
    plt.show()
```

Confusion Matrix - Logistic Regression

|                    | No Mortality (0) | Mortality (1) |
|--------------------|------------------|---------------|
| No Mortality (0)   | 396              | 111           |
| Mortality (1)      | 98               | 395           |

Logistic Regression: Accuracy=0.791, Precision=0.781, Recall=0.801, F1=0.791

# Evaluation - Random Forest

```python
# Confusion matrix for each model
from sklearn.metrics import confusion_matrix, ConfusionMatrixDisplay
from sklearn.metrics import accuracy_score, precision_score, recall_score, f1_score

for name, result in results.items():
    y_pred = result['y_pred']
    acc = accuracy_score(y_test, y_pred)
    prec = precision_score(y_test, y_pred)
    rec = recall_score(y_test, y_pred)
    f1 = f1_score(y_test, y_pred)
    print(f"{name}: Accuracy={acc:.3f}, Precision={prec:.3f}, Recall={rec:.3f}, F1={f1:.3f}")
    cm = confusion_matrix(y_test, y_pred)
    disp = ConfusionMatrixDisplay(confusion_matrix=cm, display_labels=["No Mortality (0)", "Mortality (1)"])
    disp.plot(cmap='Blues')
    plt.title(f'Confusion Matrix - {name}')
    plt.show()
```
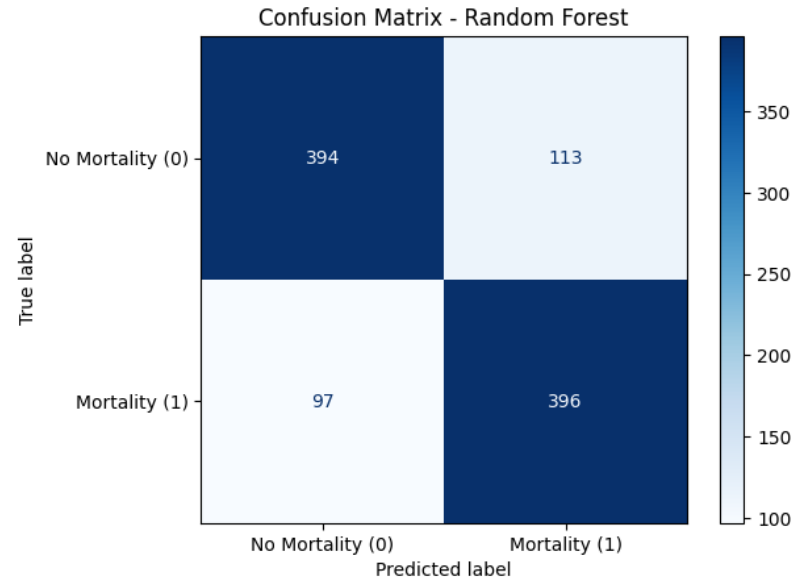


Confusion Matrix - Random Forest

Random Forest: Accuracy=0.790, Precision=0.778, Recall=0.803, F1=0.790

# Evaluation - XGBoost



```python
# Confusion matrix for each model
from sklearn.metrics import confusion_matrix, ConfusionMatrixDisplay
from sklearn.metrics import accuracy_score, precision_score, recall_score, f1_score

for name, result in results.items():
    y_pred = result['y_pred']
    acc = accuracy_score(y_test, y_pred)
    prec = precision_score(y_test, y_pred)
    rec = recall_score(y_test, y_pred)
    f1 = f1_score(y_test, y_pred)
    print(f"{name}: Accuracy={acc:.3f}, Precision={prec:.3f}, Recall={rec:.3f}, F1={f1:.3f}")
    cm = confusion_matrix(y_test, y_pred)
    disp = ConfusionMatrixDisplay(confusion_matrix=cm, display_labels=["No Mortality (0)", "Mortality (1)"])
    disp.plot(cmap='Blues')
    plt.title(f'Confusion Matrix - {name}')
    plt.show()
```
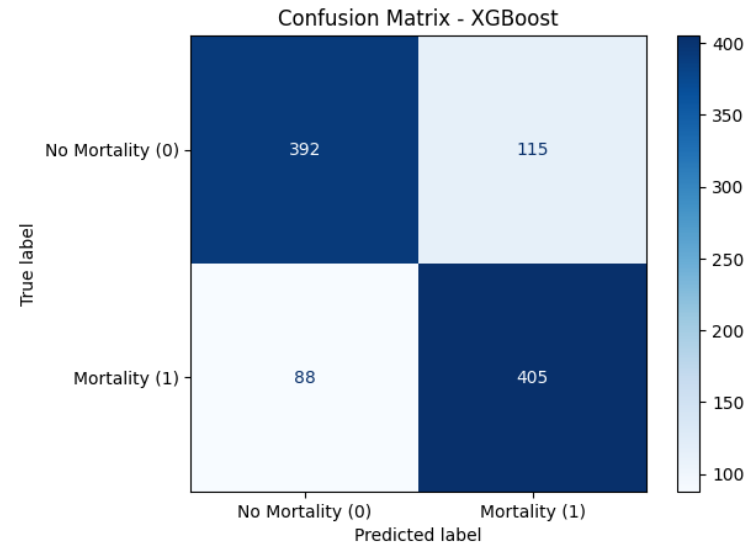
Confusion Matrix - XGBoost

XGBoost: Accuracy=0.797, Precision=0.779, Recall=0.822, F1=0.800

# Evaluation - SVM

```python
# Confusion matrix for each model
from sklearn.metrics import confusion_matrix, ConfusionMatrixDisplay
from sklearn.metrics import accuracy_score, precision_score, recall_score, f1_score

for name, result in results.items():
    y_pred = result['y_pred']
    acc = accuracy_score(y_test, y_pred)
    prec = precision_score(y_test, y_pred)
    rec = recall_score(y_test, y_pred)
    f1 = f1_score(y_test, y_pred)
    print(f"{name}: Accuracy={acc:.3f}, Precision={prec:.3f}, Recall={rec:.3f}, F1={f1:.3f}")
    cm = confusion_matrix(y_test, y_pred)
    disp = ConfusionMatrixDisplay(confusion_matrix=cm, display_labels=["No Mortality (0)", "Mortality (1)"])
    disp.plot(cmap='Blues')
    plt.title(f'Confusion Matrix - {name}')
    plt.show()
```
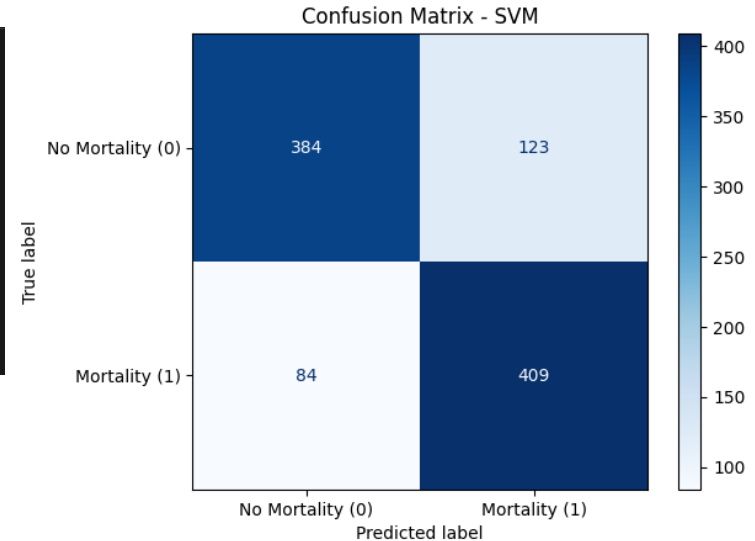


Confusion Matrix - SVM

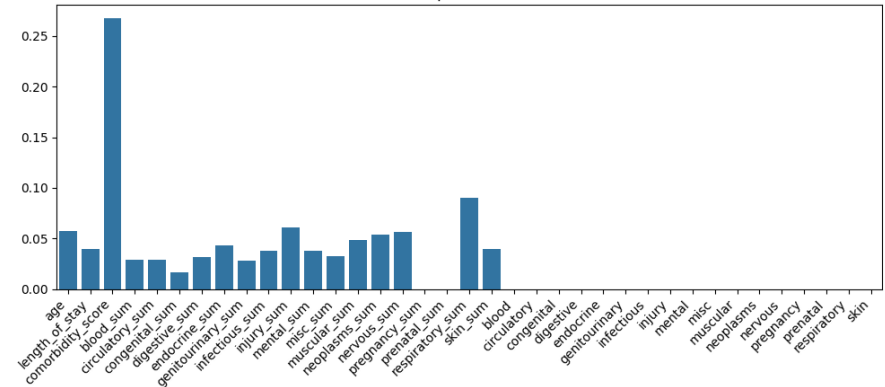SVM: Accuracy=0.793, Precision=0.769, Recall=0.830, F1=0.798

# Feature Importance

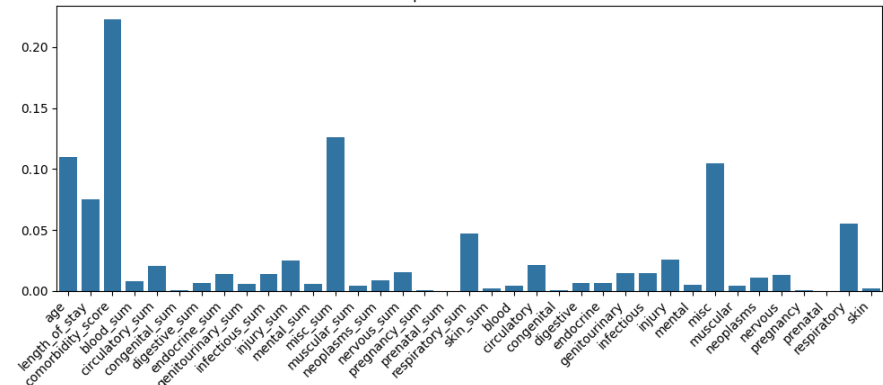- Identify top predictive features influencing mortality risk

```python
# Feature importance for tree-based models
for name in ['Random Forest', 'XGBoost']:
    model = results[name]['model']
    importances = model.feature_importances_
    feature_names = X_train_clean.columns
    plt.figure(figsize=(10, 5))
    sns.barplot(x=feature_names, y=importances)
    plt.title(f'Feature Importance - {name}')
    plt.xticks(rotation=45, ha='right')
    plt.tight_layout()
    plt.show()

# Coefficients for linear model
lr_model = results['Logistic Regression']['model']
if hasattr(lr_model, 'coef_'):
    coef = lr_model.coef_[0]
    feature_names = X_train_clean.columns
    plt.figure(figsize=(8, 5))
    sns.barplot(x=coef, y=feature_names, orient='h')
    plt.title('Feature Coefficients - Logistic Regression')
    plt.xlabel('Coefficient Value')
    plt.ylabel('Feature')
    plt.tight_layout()
    plt.show()
```
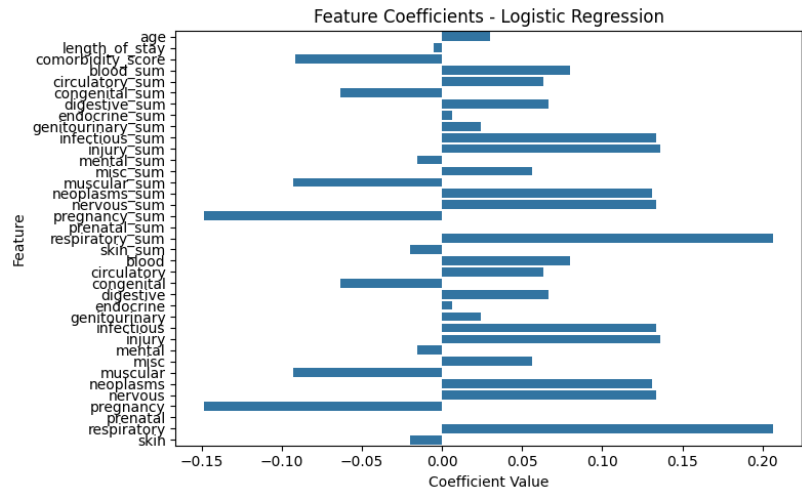```
✓  0.3s
```

# Feature Importance

- Identify top predictive features influencing mortality risk



```python
# Feature importance for tree-based models
for name in ['Random Forest', 'XGBoost']:
    model = results[name]['model']
    importances = model.feature_importances_
    feature_names = X_train_clean.columns
    plt.figure(figsize=(10, 5))
    sns.barplot(x=feature_names, y=importances)
    plt.title(f'Feature Importance – {name}')
    plt.xticks(rotation=45, ha='right')
    plt.tight_layout()
    plt.show()

# Coefficients for linear model
lr_model = results['Logistic Regression']['model']
if hasattr(lr_model, 'coef_'):
    coef = lr_model.coef_[0]
    feature_names = X_train_clean.columns
    plt.figure(figsize=(8, 5))
    sns.barplot(x=coef, y=feature_names, orient='h')
    plt.title('Feature Coefficients – Logistic Regression')
    plt.xlabel('Coefficient Value')
    plt.ylabel('Feature')
    plt.tight_layout()
    plt.show()
```
0.3s

# Conclusion & Next Steps

- With R^2 score of 0.172 and RMSE of 0.455, SVM model provide the best recall (True Positive) for prediction

- Integrate real-time EHR streaming data

- Deploy model in clinical decision support systems