```python
1 import pandas as pd
2 import numpy as np
3 df = pd.read_excel("df.xlsx")
4 df.head()
```

Show hidden output

Next steps:  [ Generate code with df ]  [ ◉ View recommended plots ]  [ New interactive sheet ]
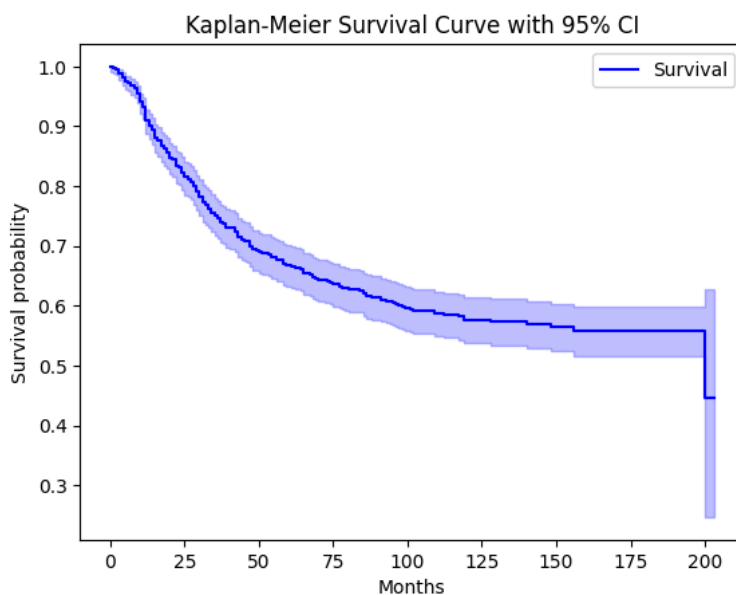
```python
1 !pip install lifelines
```

Show hidden output

```python
1 import pandas as pd
2 from lifelines import KaplanMeierFitter
3 import matplotlib.pyplot as plt
4 # Columns in your df: 'surviva_months' for time, 'status' for event (1=event, 0=censored)
5 T = df['surviva_months']
6 E = df['status']
7
8 # Initialize KM fitter
9 kmf = KaplanMeierFitter()
10
11 # Fit the data
12 kmf.fit(T, event_observed=E, label='Survival')
13
14 # Plot KM curve with 95% confidence interval
15 ax = kmf.plot(ci_show=True, color='blue')
16 ax.set_xlabel('Months')
17 ax.set_ylabel('Survival probability')
18 ax.set_title('Kaplan-Meier Survival Curve with 95% CI')
19 plt.show()
```



```python
1
2 # Define time points in months
3 time_points = [12, 36, 60]  # 1, 3, 5 years
4 # Create binary labels: 1 if patient died before the time point, 0 otherwise
5 for t in time_points:
6     df[f"status_{t}m"] = np.where((df["status"] == 1) & (df["surviva_months"] <= t), 1, 0)
7
```

```python
1 # Drop original survival info from features
2 X = df.drop(columns=["surviva_months", "status", "status_12m", "status_36m", "status_60m"])
3
4 # One-hot encode categorical variables
5 X_enc = pd.get_dummies(X, drop_first=True)
6
7 # Optional: scale numerical columns
8 from sklearn.preprocessing import StandardScaler
9
10 num_cols = ["age", "year_of_diagnosis", "tumor_size"]
11 scaler = StandardScaler()
12 X_enc[num_cols] = scaler.fit_transform(X_enc[num_cols])
```

```
13 import numpy as np
14
15 # Convert to numeric type
16 X_enc_num = X_enc.astype(float)
17
18 # Check for inf values
19 np.isinf(X_enc_num).sum()  # should be 0
20
21 # Optionally replace any remaining NaN or inf
22 X_enc_num = np.nan_to_num(X_enc_num, nan=0.0, posinf=1e10, neginf=-1e10)
23
```

```
1 from sklearn.model_selection import train_test_split
2
3 # Example: predicting 1-year survival
4 y = df["status_60m"]  # change to status_36m or status_60m for other times
5
6 X_train, X_test, y_train, y_test = train_test_split(
7     X_enc, y, test_size=0.2, random_state=42, stratify=y
8 )
9
```
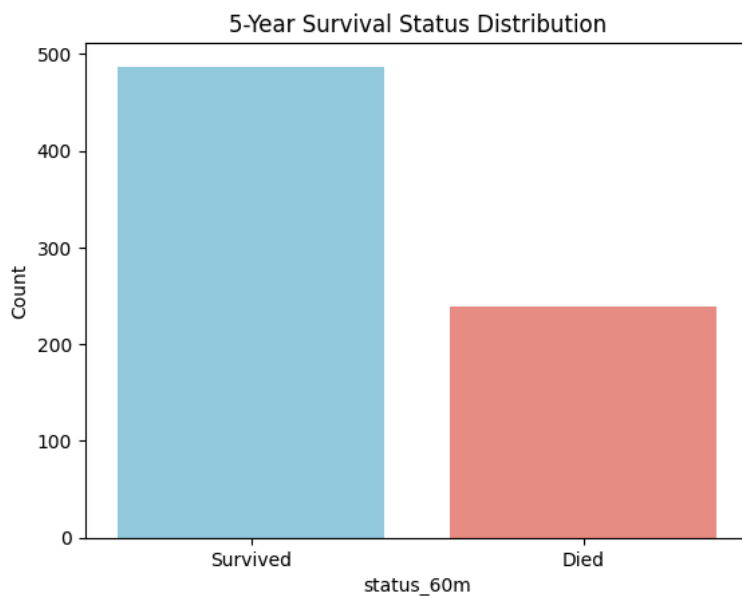
```
1 import matplotlib.pyplot as plt
2 import seaborn as sns
3
4 sns.countplot(x=y, palette=['skyblue','salmon'])
5 plt.xticks([0,1], ['Survived', 'Died'])
6 plt.ylabel("Count")
7 plt.title("5-Year Survival Status Distribution")
8 plt.show()
```

⤺  /tmp/ipython-input-1970634809.py:4: FutureWarning:

Passing `palette` without assigning `hue` is deprecated and will be removed in v0.14.0. Assign the `x` variable to `hue`

    sns.countplot(x=y, palette=['skyblue','salmon'])



```
1 from sklearn.utils import resample
2
3 # Combine X and y temporarily
4 train_data = X_train.copy()
5 train_data['target'] = y_train
6
7 # Separate majority and minority classes
8 majority = train_data[train_data.target == 0]
9 minority = train_data[train_data.target == 1]
10
11 # Downsample majority
12 majority_downsampled = resample(
13     majority,
14     replace=False,    # sample without replacement
15     n_samples=len(minority), # match minority count
16     random_state=42
17 )
18
```

```
19 # Combine minority and downsampled majority
20 train_balanced = pd.concat([majority_downsampled, minority])
21
22 # Shuffle
23 train_balanced = train_balanced.sample(frac=1, random_state=42)
24
25 # Separate X and y again
26 X_train_balanced = train_balanced.drop('target', axis=1)
27 y_train_balanced = train_balanced['target']
28
29
```

```
 1 import numpy as np
 2 import pandas as pd
 3
 4 # Assuming X_train_balanced is a DataFrame or array
 5 # For a pandas DataFrame
 6 print("NaNs in X_train_balanced:\n", X_train_balanced.isna().sum())
 7 print("Any NaNs in X_train_balanced:", X_train_balanced.isna().any().any())
 8
 9 print("NaNs in y_train_balanced:", pd.Series(y_train_balanced).isna().sum())
10
11
```

```
NaNs in X_train_balanced:
 year_of_diagnosis                                                       0
age                                                                      0
tumor_size                                                               0
sex_Male                                                                 0
race_Other (American Indian/AK Native, Asian/Pacific Islander)           0
race_White                                                               0
marital_status_yes                                                       0
AJCC_Stage_IIB                                                           0
AJCC_Stage_III                                                           0
AJCC_Stage_IV                                                            0
AJCC_Stage_IVA                                                           0
AJCC_Stage_IVB                                                           0
AJCC_T_T2                                                                0
AJCC_T_T2b                                                               0
AJCC_T_T3                                                                0
AJCC_N_N1                                                                0
AJCC_M_M1                                                                0
AJCC_M_M1a                                                               0
AJCC_M_M1b                                                               0
surgery_yes                                                              0
radiation_yes                                                            0
chemotherapy_yes                                                         0
primary_site_C40.0–Long bones: upper limb, scapula, and associated joints   0
primary_site_C40.1–Short bones of upper limb and associated joints       0
primary_site_C40.2–Long bones of lower limb and associated joints        0
primary_site_C40.3–Short bones of lower limb and associated joints       0
primary_site_C40.8–Overlap of bones, joints, and art. cartilage of limbs 0
primary_site_C41.0–Bones of skull and face and associated joints         0
primary_site_C41.1–Mandible                                              0
primary_site_C41.2–Vertebral column                                      0
primary_site_C41.3–Rib, sternum, clavicle and associated joints          0
primary_site_C41.4–Pelvic bones, sacrum, coccyx and associated joints    0
primary_site_C41.8–Overlap bones, joints, and art. cartilage             0
stage_Localized                                                          0
stage_Regional                                                           0
dtype: int64
Any NaNs in X_train_balanced: False
NaNs in y_train_balanced: 0
```

```
 1 from sklearn.ensemble import RandomForestClassifier
 2 from sklearn.linear_model import LogisticRegression
 3 from sklearn.tree import DecisionTreeClassifier
 4 from sklearn.svm import SVC
 5 from sklearn.metrics import accuracy_score, roc_auc_score, confusion_matrix
 6
 7 # Drop any row with NaN in features or target
 8 X_enc_clean = X_enc.dropna()
 9 y_clean = y.loc[X_enc_clean.index]  # make sure y matches
10
11
12 rf_model = RandomForestClassifier(n_estimators=200, random_state=42)
13 lr_model = LogisticRegression(max_iter=1000)
14 dt_model = DecisionTreeClassifier(random_state=42)
15 svm_model = SVC(probability=True, random_state=42)
16
17 rf_model.fit(X_train_balanced, y_train_balanced)
18 dt_model.fit(X_train_balanced, y_train_balanced)
19 svm_model.fit(X_train_balanced, y_train_balanced)
20 lr_model.fit(X_train_balanced, y_train_balanced)
```

```
21
22
23
24
25
```

```
▼    LogisticRegression            ⓘ ⓘ
LogisticRegression(max_iter=1000)
```

```python
1  from sklearn.metrics import (
2      accuracy_score,
3      precision_score,
4      recall_score,
5      f1_score,
6      roc_auc_score,
7      confusion_matrix
8  )
9
10 # Example: evaluate a fitted model on the test set
11 def evaluate_model(model, X_test, y_test):
12     y_pred = model.predict(X_test)
13     y_proba = model.predict_proba(X_test)[:, 1]  # probability of class 1
14
15     # Confusion matrix
16     tn, fp, fn, tp = confusion_matrix(y_test, y_pred).ravel()
17
18     results = {
19         "Accuracy": accuracy_score(y_test, y_pred),
20         "Precision (PPV)": precision_score(y_test, y_pred),
21         "Recall (Sensitivity)": recall_score(y_test, y_pred),
22         "Specificity": tn / (tn + fp + 1e-10),
23         "F1-score": f1_score(y_test, y_pred),
24         "ROC-AUC": roc_auc_score(y_test, y_proba)
25     }
26     return results
27
28 # Example usage
29 eval_rf = evaluate_model(svm_model, X_test, y_test)
30 print(eval_rf)
31
```

```
{'Accuracy': 0.815068493150685, 'Precision (PPV)': 0.28125, 'Recall (Sensitivity)': 0.6923076923076923, 'Specificity': r
```

```python
1  from imblearn.over_sampling import SMOTE
2  from sklearn.model_selection import train_test_split
3
4  # Split first
5  X_train, X_test, y_train, y_test = train_test_split(
6      X_enc, y, test_size=0.2, random_state=42, stratify=y
7  )
8
9  # Apply SMOTE on training data only
10 smote = SMOTE(random_state=42)
11 X_train_smote, y_train_smote = smote.fit_resample(X_train, y_train)
12
13 print("Original class distribution:\n", y_train.value_counts())
14 print("After SMOTE:\n", y_train_smote.value_counts())
15 rf_model.fit(X_train_smote, y_train_smote)
16 dt_model.fit(X_train_smote, y_train_smote)
17 svm_model.fit(X_train_smote, y_train_smote)
18 lr_model.fit(X_train_smote, y_train_smote)
19 eval_rf = evaluate_model(svm_model, X_test, y_test)
20 print(eval_rf)
21
22
```

```
Original class distribution:
 status_12m
0    528
1     52
Name: count, dtype: int64
After SMOTE:
 status_12m
0    528
1    528
Name: count, dtype: int64
{'Accuracy': 0.8767123287671232, 'Precision (PPV)': 0.2727272727272727, 'Recall (Sensitivity)': 0.23076923076923078, 'Sp
```

```python
1  from sklearn.ensemble import RandomForestClassifier
```

```
 2 from sklearn.svm import SVC
 3 from sklearn.linear_model import LogisticRegression
 4 from sklearn.tree import DecisionTreeClassifier
 5
 6 # Example: Random Forest with class weights
 7 rf_model = RandomForestClassifier(class_weight='balanced', random_state=42)
 8 rf_model.fit(X_train, y_train)
 9
10 # Example: SVM with class weights
11 svm_model = SVC(class_weight='balanced', probability=True, random_state=42)
12 svm_model.fit(X_train, y_train)
13
14 # Logistic Regression
15 lr_model = LogisticRegression(class_weight='balanced', max_iter=1000, random_state=42)
16 lr_model.fit(X_train, y_train)
17
18 # Decision Tree
19 dt_model = DecisionTreeClassifier(class_weight='balanced', random_state=42)
20 dt_model.fit(X_train, y_train)
21
```

```
▼                    DecisionTreeClassifier                      ⓘ ?
DecisionTreeClassifier(class_weight='balanced', random_state=42)
```

```
1 eval_rf = evaluate_model(svm_model, X_test, y_test)
2 print(eval_rf)
3
```

```
{'Accuracy': 0.8493150684931506, 'Precision (PPV)': 0.2631578947368421, 'Recall (Sensitivity)': 0.38461538461538464, 'Sp
```

```
 1 from sklearn.model_selection import cross_validate
 2 from imblearn.pipeline import Pipeline
 3 from imblearn.over_sampling import SMOTE
 4 from sklearn.ensemble import RandomForestClassifier
 5
 6 pipeline = Pipeline([
 7     ('smote', SMOTE(random_state=42)),
 8     ('dt', DecisionTreeClassifier(random_state=42))
 9 ])
10
11
12 scoring = ['accuracy', 'f1', 'roc_auc', 'precision_macro', 'recall_macro']
13
14 cv_results = cross_validate(
15     pipeline,
16     X_train,
17     y_train,
18     cv=5,
19     scoring=scoring
20 )
21
22 print("Mean accuracy:", cv_results['test_accuracy'].mean())
23 print("Mean precision:", cv_results['test_precision_macro'].mean())
24 print("Mean F1:", cv_results['test_f1'].mean())
25 print("Mean ROC-AUC:", cv_results['test_roc_auc'].mean())
26 print("Mean recall:", cv_results['test_recall_macro'].mean())
27
```

```
Mean accuracy: 0.646551724137931
Mean precision: 0.6148795942846352
Mean F1: 0.5079709000161989
Mean ROC-AUC: 0.6234195288932471
Mean recall: 0.6234195288932471
```

```
 1 from sklearn.model_selection import cross_validate
 2 from imblearn.under_sampling import RandomUnderSampler
 3 from imblearn.pipeline import Pipeline
 4 from sklearn.ensemble import RandomForestClassifier
 5
 6 # Define a pipeline with downsampling + classifier
 7 pipeline = Pipeline([
 8     ('undersample', RandomUnderSampler(random_state=42)),
 9     ('dt', DecisionTreeClassifier(random_state=42))
10 ])
11
12 # Metrics to evaluate
13 scoring = ['accuracy', 'f1', 'roc_auc', 'precision_macro', 'recall_macro']
14
15 # Run cross-validation
```

```
16 cv_results = cross_validate(
17     pipeline,
18     X_train,
19     y_train,
20     cv=5,
21     scoring=scoring
22 )
23
24 # Print mean metrics
25 print("Mean accuracy:", cv_results['test_accuracy'].mean())
26 print("Mean precision:", cv_results['test_precision_macro'].mean())
27 print("Mean ROC-AUC:", cv_results['test_roc_auc'].mean())
28 print("Mean recall:", cv_results['test_recall_macro'].mean())
29 print("Mean F1:", cv_results['test_f1'].mean())
30
```
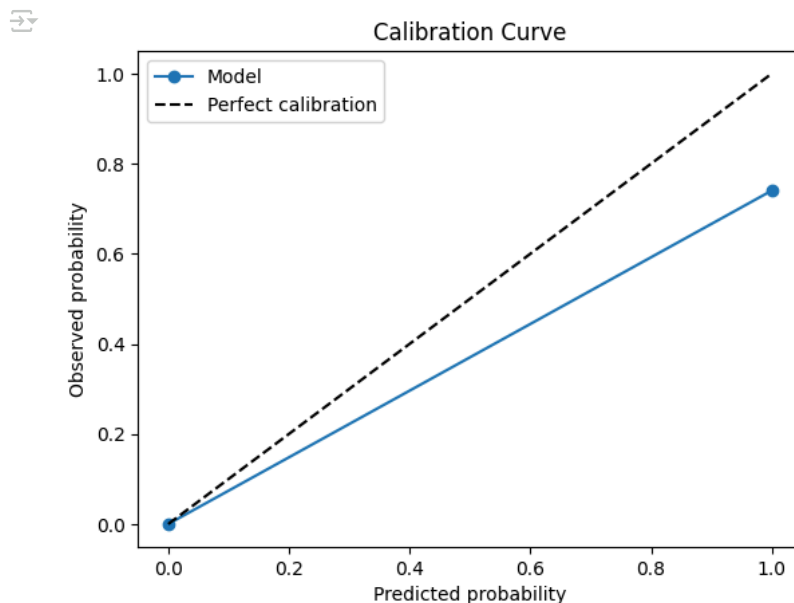
```
Mean accuracy: 0.5844827586206897
Mean precision: 0.5818526168425362
Mean ROC-AUC: 0.5916925180083075
Mean recall: 0.5916925180083075
Mean F1: 0.4928808265441928
```

```
 1 from sklearn.calibration import calibration_curve
 2
 3 y_prob_full = pipeline.fit(X_train, y_train).predict_proba(X_train)[:, 1]
 4 prob_true, prob_pred = calibration_curve(y_train, y_prob_full, n_bins=10)
 5
 6 plt.plot(prob_pred, prob_true, marker='o', label='Model')
 7 plt.plot([0,1],[0,1], 'k--', label='Perfect calibration')
 8 plt.xlabel('Predicted probability')
 9 plt.ylabel('Observed probability')
10 plt.title('Calibration Curve')
11 plt.legend()
12 plt.show()
13
14
```



```
 1 import numpy as np
 2 import matplotlib.pyplot as plt
 3
 4 def net_benefit(y_true, y_prob, thresholds):
 5     nb = []
 6     for t in thresholds:
 7         preds = (y_prob >= t).astype(int)
 8         tp = ((preds==1) & (y_true==1)).sum()
 9         fp = ((preds==1) & (y_true==0)).sum()
10         n = len(y_true)
11         nb.append((tp/n) - (fp/n)*(t/(1-t)))
12     return nb
13
14 thresholds = np.linspace(0.01, 0.99, 99)
15 y_prob = pipeline.predict_proba(X_test)[:,1]
16
17 nb_model = net_benefit(y_test.values, y_prob, thresholds)
18 nb_all = [sum(y_test==1)/len(y_test)]*len(thresholds)  # treat-all strategy
19 nb_none = [0]*len(thresholds)  # treat-none strategy
20
```

```
21 plt.plot(thresholds, nb_model, label='Model')
22 plt.plot(thresholds, nb_all, '--', label='Treat all')
23 plt.plot(thresholds, nb_none, '--', label='Treat none')
24 plt.xlabel('Threshold Probability')
25 plt.ylabel('Net Benefit')
26 plt.title('Decision Curve Analysis')
27 plt.legend()
28 plt.show()
29
30
```