

<pre># import library to project import numpy as np import pandas as pd import matplotlib.pyplot as plt import seaborn as sns  from sklearn.preprocessing import StandardScaler from sklearn.model_selection import train_test_split from sklearn.linear_model import LogisticRegression from sklearn.metrics import classification_report, confusion_matrix, f1_score, ConfusionMatrixDisplay from sklearn.model_selection import cross_val_score  from sklearn.svm import SVC from sklearn.tree import DecisionTreeClassifier from sklearn.neighbors import KNeighborsClassifier from sklearn.linear_model import LogisticRegression from sklearn.ensemble import RandomForestClassifier from sklearn.metrics import accuracy_score from sklearn.metrics import recall_score, precision_score import pickle from sklearn.metrics import accuracy_score, precision_score, recall_score, f1_score from imblearn.under_sampling import RandomUnderSampler</pre>	<p>احنا هنا عملنا تصدير للمكتبات الي بنستخدمه في الكود</p>
<pre>from google.colab import files uploaded = files.upload()</pre>	<p>هنا بنحمل الداتا من علي الجهاز علي السيرفر</p>
<pre># information the lung cancer # show the first 5 row  datacan = pd.read_csv("Lung Cancer.csv") print(datacan.head())</pre>	<p>هنا بنقر الملف ونعمل اظهار لاول خمس صفوف للملف</p>
<pre># Display summary for columns print(datacan.describe())</pre>	<p>عملنا دا عشان نعرف ملخص لداتا الموجوده في الملف</p>
<pre># Check the conversion results print(datacan.info())</pre>	<p>عشان نعرف الصفوف الي موجود ونوع البيانات الموجود واذا فيه حاجه فاضيه ولا لا</p>
<pre># checking number of rows print(datacan.shape)</pre>	<p>لمعرفه عدد الصفوف ولا الاعمده</p>

<pre># finding number of unique print(datacan.nunique())</pre>	<p>عشان نعرف عدد الاشياء المميزه في الصف</p>
<pre># Plot histograms for all columns datacan.hist(figsize=(20, 14), bins= 15) plt.show()</pre>	<p>عشان نعمل اشكال بيانيه لكل لصف</p>

<pre># change str to int number #to make easily more to make compersion between them #change female = 0 , male = 1 datacan['gender'] = datacan['gender'].map({'Male': 1, 'Female': 0})  #change cancer_stage in stage to number datacan['cancer_stage'] = datacan['cancer_stage'].map({'Stage I' : 1, 'Stage II' : 2, 'Stage III' : 3 , 'Stage IV' : 4, 'Stage V' : 5 , 'Stage VI' : 6})  #change smoking_status in stage to number datacan['smoking_status'] = datacan['smoking_status'].map({"Never Smoked": 0, "Passive Smoker": 1,"Former Smoker": 2,"Current Smoker": 3})  #change family_history in stage to number datacan['family_history'] = datacan['family_history'].map({'Yes' : 1 , 'No' : 0 })  ##change family_history in stage to number datacan['treatment_type'] = datacan['treatment_type'].map({"Chemotherapy": 0, "Surgerly": 1, "Radiation": 2,"Combined": 3 })</pre>	<p>هنا حولنا اي حاجه نص في العواميد دا عشان نعرف نعدل او نعمل حساب في الكود عشان تبقي اسهل في العمليات وميعملش خط</p>
<pre># Checking for null values (if any) print(datacan.isna().any()) print(datacan.isnull().sum())</pre>	<p>بنشوف ان فيه حاجه فاضيه ولا لا ولو فيه يظهر عددهم</p>
<pre>#show boxplot for col in datacan: plt.figure(figsize=(6, 2)) sns.boxplot(data=datacan, x=col) plt.title(f"Outlier Check: {col}") plt.tight_layout() plt.show()</pre>	<p>احنا استخدمنا الشكل دا عشان يظهر الاشكال عشان برضوا نعرف الاوت ليرس</p>

<pre># Country distribution pie chart country_counts = datacan['country'].value_counts().head(10)  plt.figure(figsize=(12, 8)) plt.pie(country_counts.values, labels=country_counts.index, autopct='%1.1f%%') plt.title('Top 10 Countries by Lung Cancer Patient Distribution') plt.axis('equal') plt.show()  print("Top 10 countries:") print(country_counts)</pre>	<p>هنا عملنا شكل بياني للعمود الدول لمعرفة اكثر 10 دول مصابه</p>
<pre># delete the information str and i dont need that #impact the acurcye the programe for xy in ['end_treatment_date', 'diagnosis_date', 'country', 'id']:     if xy in datacan.columns:         datacan.drop(columns=xy, inplace=True)</pre>	<p>هنا حذفنا الاشياء التي تؤثر ولا احتاجها في العمليات</p>
<pre># Family history vs cancer stage heatmap family_stage_crosstab = pd.crosstab(datacan['family_history'], datacan['cancer_stage'],  plt.figure(figsize=(10, 6)) sns.heatmap(family_stage_crosstab, annot=True, cmap='YlOrRd', fmt='.1f',             cbar_kws={'label': 'Percentage'}) plt.title('Cancer Stage Distribution by Family History') plt.xlabel('Cancer Stage') plt.ylabel('Family History') plt.show()</pre>	<p>عملنا شكل بياني عشان نعرف اذا في لو فيه حد عندوا كان حد عندوا تاريخ او لا نعرف مستوي درجه الاصابه</p>
<pre>#check target variable balance  sns.countplot( x='survived',data=datacan) plt.title('Survival Count') plt.show()</pre>	<p>عشان نعرف عدد الناس الناجين</p>
<pre># correlation matrix correlation_matrix = datacan.corr() plt.figure(figsize=(12, 8)) sns.heatmap(correlation_matrix, annot=True, fmt=".2f", cmap='coolwarm', square=True) plt.title('Correlation Matrix') plt.show()</pre>	<p>عشان نعرف العلاقه بين العوامل والصوف</p>

<pre>columns_to_exclude = ['survived'] for column in datacan.select_dtypes(include=['float64', 'int64']).columns:     if column not in columns_to_exclude:         Q1 = datacan[column].quantile(0.25)         Q3 = datacan[column].quantile(0.75)         IQR = Q3 - Q1         lower_bound = Q1 - 1.5 * IQR         upper_bound = Q3 + 1.5 * IQR         datacan = datacan[(datacan[column] &gt;= lower_bound) &amp; (datacan[column] &lt;= upper_bound)]</pre>	<p>عملنا الكود دا عشان نحذف الاوت لير التي توثر علي عمليات التدريب</p>
<pre>print(datacan['gender'].value_counts()) print(datacan.groupby('gender')['survived'].value_counts())  # Survival rate by gender print(datacan.groupby('gender')['survived'].mean())</pre>	<p>عشان نعرف العلاقة بين بين النوع هياثر علي احتماليه النجاه</p>
<pre>#Split and scale data X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2,                                                     random_state=42)  scaler = StandardScaler() X_train_scaled = scaler.fit_transform(X_train) X_test_scaled = scaler.transform(X_test)</pre>	<p>عملنا دا عشان نفصل الداتا شويه للتدريب وشويه عشان نعمل بيهم للتحقق من دقة النموذج</p>
<pre>#Train and evaluate base model  logistic_model = LogisticRegression() logistic_model.fit(X_train_scaled, y_train) y_pred = logistic_model.predict(X_test_scaled)  # Accuracy print(f"Logistic Regression Accuracy: {accuracy_score(y_test, y_pred):.2%}")</pre>	<p>عشان ندرب علي Logistic Regression ونشوف دقة التنبؤ بتاع النموذج</p>
<pre># Confusion Matrix conf_matrix = confusion_matrix(y_test, y_pred) conf_matrix_df = pd.DataFrame(conf_matrix,                               index=['Actual Negative', 'Actual Positive'],                               columns=['Predicted Negative', 'Predicted Positive']) print("Confusion Matrix:\n", conf_matrix_df)  # Classification Report print("Classification Report:\n", classification_report(y_test, y_pred))</pre>	<p>عشان نعرف precision, recall, f1-score, support, Predicted Negative, Predicted Positive, Actual Negative, Actual Positive</p>

```
#train model
models = {
    'Logistic Regression': LogisticRegression(),
    'SVM': SVC(class_weight='balanced'),
    'Decision Tree': DecisionTreeClassifier(),
    'Random Forest': RandomForestClassifier(),
    'KNN': KNeighborsClassifier()
}

results = []
f1_scores = {}
```

هنا عملنا الكود عشان نعمل  
تدريب لاداتا  
ونحفظها

```
for name, model in models.items():
    model.fit(X_train_scaled, y_train)
    y_pred = model.predict(X_test_scaled)

    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)

    results.append({
        'Model': name,
        'Accuracy': acc,
        'Precision': prec,
        'Recall': rec,
        'F1 Score': f1
    })
    f1_scores[name] = f1
```

عشان نعمل تدريب لنماذج  
كلها مع بعض  
عشان نعرف افضل طريقه  
نستخدمها نمشي بيها

```
# Visualize performance
results_df = pd.DataFrame(results)
print("\nModel Comparison Results:")
print(results_df)

plt.figure(figsize=(10, 6))
for metric in ['Accuracy', 'Precision', 'Recall', 'F1 Score']:
    plt.plot(results_df['Model'], results_df[metric], marker='o', label=metric)
plt.title('Model Performance Comparison')
plt.ylabel('Score')
plt.ylim(0, 1.1)
plt.legend()
plt.grid(True)
plt.show()
```

عملنا الكود عشان نظهر  
الفرق بين طرق التدريب  
علي شكل رسم بياني

<pre># F1 Bar Chart plt.figure(figsize=(8, 5)) plt.bar(f1_scores.keys(), f1_scores.values(), color='skyblue') plt.title("F1-score Comparison between Models") plt.xlabel("Model") plt.ylabel("F1-score") plt.ylim(0, 1) plt.grid(axis='y') for i, (model, score) in enumerate(f1_scores.items()):     plt.text(i, score + 0.01, f"{score:.2f}", ha='center', fontweight='bold') plt.tight_layout() plt.show()</pre>	<p>عملنا رسم لمقارنة نماذج بناءً على F1-score</p>
<pre># Cross-Validation  print("\nCross-Validation Scores:") for name, model in models.items():     scores = cross_val_score(model, X, y, cv=5)     print(f"{name}: Mean CV Accuracy = {scores.mean():.2f}")</pre>	<p>هنا يعمل التدريب أكثر من مره وبعدين المتوسط عشان نعرف اقرب حاجه للدقه النموذج للحقيقه</p>
<pre>rus = RandomUnderSampler(random_state=42) X_resampled, y_resampled = rus.fit_resample(X_train, y_train)  scaler = StandardScaler() X_resampled_scaled = scaler.fit_transform(X_resampled) X_test_scaled = scaler.transform(X_test)</pre>	<p>عشان نحل مشكله الاصفار التي طلعت في التدريب في حاله الناجين مطلعهم صفر</p>
<pre>cm = confusion_matrix(y_test, y_pred) sns.heatmap(cm, annot=True, fmt='d', cmap='Blues') plt.title('Confusion Matrix ') plt.xlabel('Predicted') plt.ylabel('Actual') plt.show()  # for name, model in models.items():     with open(f'{name.replace(" ", "_").lower()}_model.pkl', 'wb') as f:         pickle.dump(model, f)</pre>	<p>استخدمنا الكود توضيح أداء النموذج في التنبؤ بالفئات الفعلية ويقوم بحفظها</p>

<pre># Save trained models  for name, model in models.items():     with open(f'{name.replace(" ", "_").lower()}_model.pkl', 'wb') as f:         pickle.dump(model, f)</pre>	يقوم عمليات التدريب التي تم تنفيذها في الكود
<pre># save the new values after that remove the text #datacan.to_csv('Lung Cancer.csv')</pre>	يقوم بحفظ اي تعديل في الملف