

# Credit Card Default Prediction

Mohammed El-Bathy  
Department of Computer Science and Engineering  
Oakland University  
Rochester, Michigan, USA  
elbathy@oakland.edu

**Abstract—** This project is designed to predict credit card defaults by analyzing a variety of financial and demographic data attributes. The dataset undergoes multiple stages, including data preprocessing, feature engineering, and model training, which culminate in the evaluation and comparison of three machine learning models: Decision Tree, Logistic Regression, and Random Forest. Each model is tested through various evaluation metrics to determine its effectiveness, and the results are carefully documented to support model selection

## I. INTRODUCTION

Credit card default prediction is a critical issue for financial institutions that aim to mitigate risk and minimize losses. By leveraging customer data such as payment history, credit limits, and demographic attributes, we can build predictive models to estimate the likelihood of default. This project explores the use of decision tree classifiers, focusing on preparing the data, handling class imbalance, and evaluating model performance.

## II. DOMAIN KNOWLEDGE

In the domain of credit risk, default refers to the failure to meet the legal obligations of debt repayment. Credit card issuers face significant financial risks due to defaults, so it's crucial to develop predictive models to assess the likelihood of default before issuing credit.

The features used in this analysis cover various aspects of customer behavior and demographics, such as credit limit, payment history, bill amounts, and personal attributes (e.g., age, gender, education, and marital status). Understanding these variables is essential to modeling the risk of default

## III. DATASET ANALYSIS & UNDERSTANDING

All the headings in the main body are numbered (automatically).

### A. Data Characteristics

The dataset contains 30 variables, including both numerical and categorical data. Key variables include:

- **LIMIT\_BAL**: The credit limit assigned to the customer.
- **AGE**: The age of the customer
- **PAY\_0 to PAY\_6**: Repayment status in different months.
- **BILL\_AMT1 to BILL\_AMT6**: Amount of bill statements from the past six months.
- **PAY\_AMT1 to PAY\_AMT6**: Amount paid in the past six months.
- **SEX, EDUCATION, MARRIAGE**: Categorical features indicating demographic details

The target variable is default payment next month, indicating whether the customer defaulted (1) or not (0)

### B. Feature Analysis & Selection

- The key financial variables are credit limit (LIMIT\_BAL), bill amounts (BILL\_AMT1-6), and payment history (PAY\_0-6). These features directly reflect the customer's ability to manage credit and meet payment obligations
- Demographic variables such as AGE, SEX, EDUCATION, and MARRIAGE also provide insights into the customer profile and their associated risk levels

## IV. DATA PREPROCESSING FUNCTIONS

To ensure a smooth workflow, the project begins by importing the required libraries, which helps avoid any errors from missing dependencies. The `import_libraries` function accomplishes this initial step, setting the foundation for the project by confirming that all necessary packages are available. In addition we defined global variables.

```
import gradio as gr
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
from io import BytesIO
from PIL import Image
from sklearn.tree import DecisionTreeClassifier
from sklearn.ensemble import RandomForestClassifier
from sklearn.linear_model import LogisticRegression
from sklearn.model_selection import StratifiedKFold, cross_val_score, GridSearchCV
from sklearn.metrics import confusion_matrix, classification_report, roc_auc_score
from sklearn.preprocessing import StandardScaler
from statsmodels.stats.outliers_influence import variance_inflation_factor
from statsmodels.tools.tools import add_constant
from imblearn.over_sampling import SMOTE
from sklearn import tree
```

```
# Global variables
data = None
encoded_data = None
X_resampled = None
y_resampled = None
dt_model = None
lr_model = None
rf_model = None
lr_results = None
rf_results = None
```

### A. Loading Dataset.

The next step is loading the dataset through the `load_dataset` function. This function reads in the credit card client data, providing access to the information necessary for further processing and analysis.

```
# Function 2: Load the Dataset
def load_dataset():
    global data
    file_path = 'default_of_credit_card_clients.csv'
    data = pd.read_csv(file_path, header=1)
    return data
```

#	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30	31	32	33	34	35	36	37	38	39	40	41	42	43	44	45	46	47	48	49	50	51	52	53	54	55	56	57	58	59	60	61	62	63	64	65	66	67	68	69	70	71	72	73	74	75	76	77	78	79	80	81	82	83	84	85	86	87	88	89	90	91	92	93	94	95	96	97	98	99	100	101	102	103	104	105	106	107	108	109	110	111	112	113	114	115	116	117	118	119	120	121	122	123	124	125	126	127	128	129	130	131	132	133	134	135	136	137	138	139	140	141	142	143	144	145	146	147	148	149	150	151	152	153	154	155	156	157	158	159	160	161	162	163	164	165	166	167	168	169	170	171	172	173	174	175	176	177	178	179	180	181	182	183	184	185	186	187	188	189	190	191	192	193	194	195	196	197	198	199	200	201	202	203	204	205	206	207	208	209	210	211	212	213	214	215	216	217	218	219	220	221	222	223	224	225	226	227	228	229	230	231	232	233	234	235	236	237	238	239	240	241	242	243	244	245	246	247	248	249	250	251	252	253	254	255	256	257	258	259	260	261	262	263	264	265	266	267	268	269	270	271	272	273	274	275	276	277	278	279	280	281	282	283	284	285	286	287	288	289	290	291	292	293	294	295	296	297	298	299	300	301	302	303	304	305	306	307	308	309	310	311	312	313	314	315	316	317	318	319	320	321	322	323	324	325	326	327	328	329	330	331	332	333	334	335	336	337	338	339	340	341	342	343	344	345	346	347	348	349	350	351	352	353	354	355	356	357	358	359	360	361	362	363	364	365	366	367	368	369	370	371	372	373	374	375	376	377	378	379	380	381	382	383	384	385	386	387	388	389	390	391	392	393	394	395	396	397	398	399	400	401	402	403	404	405	406	407	408	409	410	411	412	413	414	415	416	417	418	419	420	421	422	423	424	425	426	427	428	429	430	431	432	433	434	435	436	437	438	439	440	441	442	443	444	445	446	447	448	449	450	451	452	453	454	455	456	457	458	459	460	461	462	463	464	465	466	467	468	469	470	471	472	473	474	475	476	477	478	479	480	481	482	483	484	485	486	487	488	489	490	491	492	493	494	495	496	497	498	499	500	501	502	503	504	505	506	507	508	509	510	511	512	513	514	515	516	517	518	519	520	521	522	523	524	525	526	527	528	529	530	531	532	533	534	535	536	537	538	539	540	541	542	543	544	545	546	547	548	549	550	551	552	553	554	555	556	557	558	559	560	561	562	563	564	565	566	567	568	569	570	571	572	573	574	575	576	577	578	579	580	581	582	583	584	585	586	587	588	589	590	591	592	593	594	595	596	597	598	599	600	601	602	603	604	605	606	607	608	609	610	611	612	613	614	615	616	617	618	619	620	621	622	623	624	625	626	627	628	629	630	631	632	633	634	635	636	637	638	639	640	641	642	643	644	645	646	647	648	649	650	651	652	653	654	655	656	657	658	659	660	661	662	663	664	665	666	667	668	669	670	671	672	673	674	675	676	677	678	679	680	681	682	683	684	685	686	687	688	689	690	691	692	693	694	695	696	697	698	699	700	701	702	703	704	705	706	707	708	709	710	711	712	713	714	715	716	717	718	719	720	721	722	723	724	725	726	727	728	729	730	731	732	733	734	735	736	737	738	739	740	741	742	743	744	745	746	747	748	749	750	751	752	753	754	755	756	757	758	759	760	761	762	763	764	765	766	767	768	769	770	771	772	773	774	775	776	777	778	779	780	781	782	783	784	785	786	787	788	789	790	791	792	793	794	795	796	797	798	799	800	801	802	803	804	805	806	807	808	809	810	811	812	813	814	815	816	817	818	819	820	821	822	823	824	825	826	827	828	829	830	831	832	833	834	835	836	837	838	839	840	841	842	843	844	845	846	847	848	849	850	851	852	853	854	855	856	857	858	859	860	861	862	863	864	865	866	867	868	869	870	871	872	873	874	875	876	877	878	879	880	881	882	883	884	885	886	887	888	889	890	891	892	893	894	895	896	897	898	899	900	901	902	903	904	905	906	907	908	909	910	911	912	913	914	915	916	917	918	919	920	921	922	923	924	925	926	927	928	929	930	931	932	933	934	935	936	937	938	939	940	941	942	943	944	945	946	947	948	949	950	951	952	953	954	955	956	957	958	959	960	961	962	963	964	965	966	967	968	969	970	971	972	973	974	975	976	977	978	979	980	981	982	983	984	985	986	987	988	989	990	991	992	993	994	995	996	997	998	999	1000	1001	1002	1003	1004	1005	1006	1007	1008	1009	1010	1011	1012	1013	1014	1015	1016	1017	1018	1019	1020	1021	1022	1023	1024	1025	1026	1027	1028	1029	1030	1031	1032	1033	1034	1035	1036	1037	1038	1039	1040	1041	1042	1043	1044	1045	1046	1047	1048	1049	1050	1051	1052	1053	1054	1055	1056	1057	1058	1059	1060	1061	1062	1063	1064	1065	1066	1067	1068	1069	1070	1071	1072	1073	1074	1075	1076	1077	1078	1079	1080	1081	1082	1083	1084	1085	1086	1087	1088	1089	1090	1091	1092	1093	1094	1095	1096	1097	1098	1099	1100	1101	1102	1103	1104	1105	1106	1107	1108	1109	1110	1111	1112	1113	1114	1115	1116	1117	1118	1119	1120	1121	1122	1123	1124	1125	1126	1127	1128	1129	1130	1131	1132	1133	1134	1135	1136	1137	1138	1139	1140	1141	1142	1143	1144	1145	1146	1147	1148	1149	1150	1151	1152	1153	1154	1155	1156	1157	1158	1159	1160	1161	1162	1163	1164	1165	1166	1167	1168	1169	1170	1171	1172	1173	1174	1175	1176	1177	1178	1179	1180	1181	1182	1183	1184	1185	1186	1187	1188	1189	1190	1191	1192	1193	1194	1195	1196	1197	1198	1199	1200	1201	1202	1203	1204	1205	1206	1207	1208	1209	1210	1211	1212	1213	1214	1215	1216	1217	1218	1219	1220	1221	1222	1223	1224	1225	1226	1227	1228	1229	1230	1231	1232	1233	1234	1235	1236	1237	1238	1239	1240	1241	1242	1243	1244	1245	1246	1247	1248	1249	1250	1251	1252	1253	1254	1255	1256	1257	1258	1259	1260	1261	1262	1263	1264	1265	1266	1267	1268	1269	1270	1271	1272	1273	1274	1275	1276	1277	1278	1279	1280	1281	1282	1283	1284	1285	1286	1287	1288	1289	1290	1291	1292	1293	1294	1295	1296	1297	1298	1299	1300	1301	1302	1303	1304	1305	1306	1307	1308	1309	1310	1311	1312	1313	1314	1315	1316	1317	1318	1319	1320	1321	1322	1323	1324	1325	1326	1327	1328	1329	1330	1331	1332	1333	1334	1335	1336	1337	1338	1339	1340	1341	1342	1343	1344	1345	1346	1347	1348	1349	1350	1351	1352	1353	1354	1355	1356	1357	1358	1359	1360	1361	1362	1363	1364	1365	1366	1367	1368	1369	1370	1371	1372	1373	1374	1375	1376	1377	1378	1379	1380	1381	1382	1383	1384	1385	1386	1387	1388	1389	1390	1391	1392	1393	1394	1395	1396	1397	1398	1399	1400	1401	1402	1403	1404	1405	1406	1407	1408	1409	1410	1411	1412	1413	1414	1415	1416	1417	1418	1419	1420	1421	1422	1423	1424	1425	1426	1427	1428	1429	1430	1431	1432	1433	1434	1435	1436	1437	1438	1439	1440	1441	1442	1443	1444	1445	1446	1447	1448	1449	1450	1451	1452	1453	1454	1455	1456	1457	1458	1459	1460	1461	1462	1463	1464	1465	1466	1467	1468	1469	1470	1471	1472	1473	1474	1475	1476	1477	1478	1479	1480	1481	1482	1483	1484	1485	1486</
---	---	---	---	---	---	---	---	---	---	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	--------



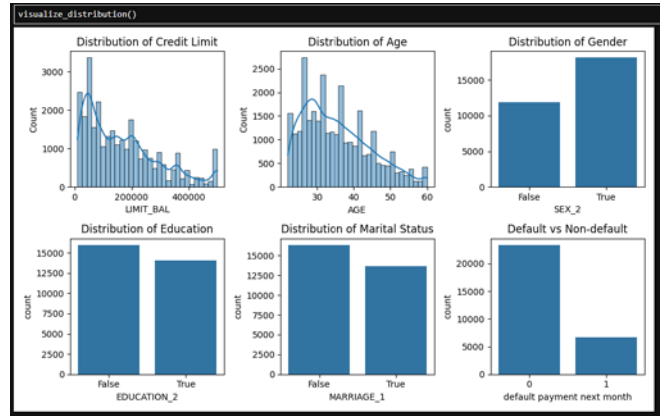
recalculate_vif()		
	Feature	VIF
0	const	5.502704
1	ID	1.012959
2	LIMIT_BAL	1.513392
3	AGE	1.024821
4	PAY_0	2.008695
5	PAY_2	3.212227
6	PAY_3	3.743819
7	PAY_4	4.402964
8	PAY_5	4.838624
9	PAY_6	3.311836
10	BILL_AMT1	13.271621
11	BILL_AMT2	24.583129
12	BILL_AMT3	21.497720
13	BILL_AMT4	21.390585
14	BILL_AMT5	25.037119
15	BILL_AMT6	14.797413
16	PAY_AMT1	1.762232
17	PAY_AMT2	1.885662
18	PAY_AMT3	1.778416
19	PAY_AMT4	1.764788
20	PAY_AMT5	1.799162
21	PAY_AMT6	1.233234
22	default payment next month	1.143918

#### J. Visualize Distribution

To further explore the dataset, visualize distribution generates plots that display the distribution of key features. This visualization provides additional insights into data patterns, helping identify skewness or other anomalies that might necessitate transformation or further preprocessing.

```
# Function 11: Visualize Distribution
def visualize_distribution():
    global data, encoded_data
    if data is not None and encoded_data is not None:
        fig, axes = plt.subplots(2, 3, figsize=(10, 6))
        sns.histplot(data['LIMIT_BAL'], bins=30, kde=True, ax=axes[0, 0]).set_title('Distribution of Credit Limit')
        sns.histplot(data['AGE'], bins=30, kde=True, ax=axes[0, 1]).set_title('Distribution of Age')
        sns.countplot(x='SEX_2', data=encoded_data, ax=axes[0, 2]).set_title('Distribution of Gender')
        sns.countplot(x='EDUCATION_2', data=encoded_data, ax=axes[1, 0]).set_title('Distribution of Education')
        sns.countplot(x='MARRIAGE_1', data=encoded_data, ax=axes[1, 1]).set_title('Distribution of Marital Status')
        sns.countplot(x='default payment next month', data=data, ax=axes[1, 2]).set_title('Default vs Non-default')
        plt.tight_layout()

        buf = BytesIO()
        fig.savefig(buf, format='png')
        buf.seek(0)
        plt.close(fig)
        return Image.open(buf)
    else:
        return "Data not loaded."
```



### V. MODEL TRAINING AND EVALUATION FUNCTIONS

#### A. Decision Tree Model

1) The first machine learning model employed in this project is a Decision Tree. The initialize\_decision\_tree function initializes and trains a Decision Tree classifier on the resampled dataset, establishing a straightforward and interpretable baseline model.

```
# Function 12: Initialize and Train Decision Tree
def initialize_decision_tree():
    global dt_model, X_resampled, y_resampled
    if X_resampled is not None and y_resampled is not None:
        dt_model = DecisionTreeClassifier(random_state=42)
        dt_model.fit(X_resampled, y_resampled)
        return "Decision Tree Classifier initialized and trained."
    else:
        return "SMOTE resampled data not available."
```

```
initialize_decision_tree()
```

```
'Decision Tree Classifier initialized and trained.'
```

2) To evaluate the Decision Tree, cross\_validation\_decision\_tree performs cross-validation, calculating accuracy and ROC-AUC scores across folds. Cross-validation reduces the risk of overfitting and provides a reliable estimate of the model's generalization capabilities.

```
# Function 13: K-Fold Cross-Validation for Decision Tree
def cross_validation_decision_tree():
    global dt_model, X_resampled, y_resampled
    if dt_model is not None and X_resampled is not None:
        kf = StratifiedKFold(n_splits=5, shuffle=True, random_state=42)
        accuracy_scores = cross_val_score(dt_model, X_resampled, y_resampled, cv=kf, scoring='accuracy')
        roc_auc_scores = cross_val_score(dt_model, X_resampled, y_resampled, cv=kf, scoring='roc_auc')
        return f"Mean Accuracy: {accuracy_scores.mean():.4f}, Std Dev: {accuracy_scores.std():.4f} \n" \
            f"Mean ROC-AUC: {roc_auc_scores.mean():.4f}, Std Dev: {roc_auc_scores.std():.4f}"
    else:
        return "Decision Tree or resampled data not available."
```

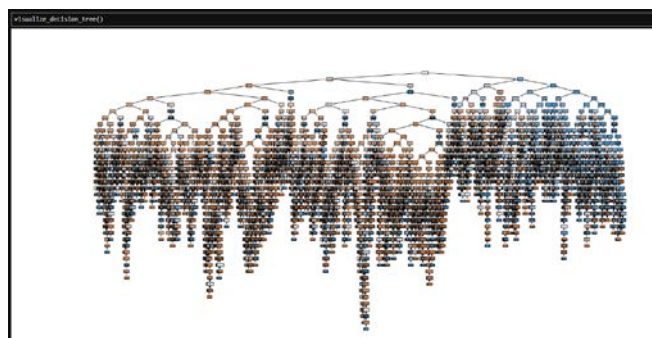
```
cross_validation_decision_tree()
```

```
'Mean Accuracy: 0.8043, Std Dev: 0.0021\nMean ROC-AUC: 0.8043, Std Dev: 0.0021'
```



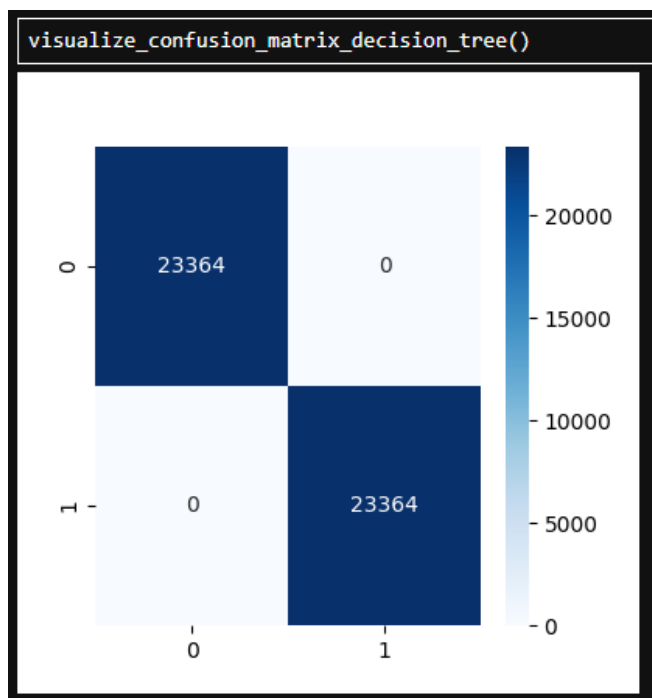
3) To enhance interpretability, `visualize_decision_tree` produces a graphical representation of the tree structure, displaying the decision-making paths and feature splits.

```
# Function 14: Visualize Decision Tree Structure
def visualize_decision_tree():
    global dt_model, X_resampled
    if dt_model is not None and X_resampled is not None:
        plt.figure(figsize=(20, 10))
        tree.plot_tree(dt_model, feature_names=X_resampled.columns, class_names=["No Default", "Default"], filled=True)
        buf = BytesIO()
        plt.savefig(buf, format="png")
        buf.seek(0)
        plt.close()
        return Image.open(buf)
    else:
        return "Decision Tree model not trained."
```



4) The confusion matrix for the Decision Tree, generated by `visualize_confusion_matrix_decision_tree`, offers a detailed view of true versus predicted classifications. This matrix highlights the model's accuracy on each class, helping to identify any biases in prediction.

```
# Function 15: Confusion Matrix Visualization for Decision Tree
def visualize_confusion_matrix_decision_tree():
    global dt_model, X_resampled, y_resampled
    if dt_model is not None and X_resampled is not None:
        plt.figure(figsize=(8, 4))
        sns.heatmap(confusion_matrix(y_resampled, dt_model.predict(X_resampled)), annot=True, fmt='d', cmap='Blues')
        buf = BytesIO()
        plt.savefig(buf, format="png")
        buf.seek(0)
        plt.close()
        return Image.open(buf)
    else:
        return "Decision Tree or resampled data not available."
```



5) Finally, `classification_report_dt` provides a classification report with precision, recall, and F1-score metrics for the Decision Tree. This report allows for a

comprehensive evaluation of the model's performance across each class, especially in imbalanced datasets.

```
# Function 16: Classification Report for Decision Tree
def classification_report_dt():
    global dt_model, X_resampled, y_resampled
    if dt_model is not None and X_resampled is not None:
        y_pred = dt_model.predict(X_resampled)
        report = classification_report(y_resampled, y_pred)
        return report
    else:
        return "Decision Tree or resampled data not available."
```

classification_report_dt()										
	precision	recall	f1-score	support	0	1.00	1.00	1.00	23364	n
1.00	46728	n	macro avg	1.00	1.00	1.00	46728	n	weighted avg	1.00

	precision	recall	f1-score	support
0	1.00	1.00	1.00	23364
1	1.00	1.00	1.00	23364
accuracy			1.00	46728
macro avg	1.00	1.00	1.00	46728
weighted avg	1.00	1.00	1.00	46728

## B. Logistic Regression Model

1) Logistic Regression is the second model used in this analysis. The `initialize_logistic_regression` function initializes and trains the Logistic Regression model on the resampled data, providing a robust, interpretable model that serves as a benchmark.

```
# Function 17: Initialize and Train Logistic Regression
def initialize_logistic_regression():
    global lr_model, X_resampled, y_resampled
    if X_resampled is not None and y_resampled is not None:
        lr_model = LogisticRegression(max_iter=1000, random_state=42)
        lr_model.fit(X_resampled, y_resampled)
        return "Logistic Regression model initialized and trained."
    else:
        return "SMOTE resampled data not available."
```

2) Evaluation of the Logistic Regression model is conducted through `evaluate_logistic_regression`, which generates a classification report and calculates the ROC-AUC score. These metrics offer a thorough view of the model's ability to discriminate between classes and provide a basis for comparing it to other models.

```
# Function 18: Logistic Regression Evaluation
def evaluate_logistic_regression():
    global lr_model, X_resampled, y_resampled
    if lr_model is not None and X_resampled is not None:
        y_pred = lr_model.predict(X_resampled)
        report = classification_report(y_resampled, y_pred)
        auc = roc_auc_score(y_resampled, lr_model.predict_proba(X_resampled)[:, 1])
        return f"Classification Report:\n{report}\nROC-AUC Score: {auc:.4f}"
    else:
        return "Logistic Regression model or resampled data not available."
```

#### Classification Report:

	precision	recall	f1-score	support
0	0.77	0.75	0.76	23364
1	0.76	0.77	0.76	23364
accuracy		0.76	0.76	46728
macro avg	0.76	0.76	0.76	46728
weighted avg	0.76	0.76	0.76	46728

ROC-AUC Score: 0.8431

### C. Random Forest Model

1) The third and final model in this project is a Random Forest, which is both flexible and powerful for handling complex data patterns. The initialize\_random\_forest function initializes and tunes the Random Forest model using GridSearchCV, optimizing hyperparameters based on ROC-AUC scores. This tuning process enhances the model's performance by selecting the best combination of parameters.

```
# Function 19: Initialize and Train Random Forest with Hyperparameter Tuning
def initialize_random_forest():
    global rf_model, X_resampled, y_resampled, rf_results
    if X_resampled is not None and y_resampled is not None:
        rf_model = RandomForestClassifier(random_state=42)
        param_grid = {
            'n_estimators': [50, 100, 200],
            'max_depth': [None, 10, 20, 30],
            'min_samples_split': [2, 5, 10]
        }
        grid_search = GridSearchCV(rf_model, param_grid, cv=5, scoring='roc_auc', n_jobs=-1)
        grid_search.fit(X_resampled, y_resampled)
        rf_model = grid_search.best_estimator_
        rf_results = {
            "Best Parameters": grid_search.best_params_,
            "Best ROC-AUC": grid_search.best_score_
        }
        return f"Random Forest model initialized and trained with tuning. Best ROC-AUC: {rf_results['Best ROC-AUC']:.4f}"
    else:
        return "SMOTE resampled data not available."
```

output

Random Forest model initialized and trained with tuning. Best ROC-AUC: 0.5598

2) The evaluation function evaluate\_random\_forest then generates a classification report and ROC-AUC score, which reveal the model's performance in distinguishing between classes. This assessment is critical for comparing the Random Forest's predictive power against the other models.

```
# Function 20: Random Forest Evaluation
def evaluate_random_forest():
    global rf_model, X_resampled, y_resampled
    if rf_model is not None and X_resampled is not None:
        y_pred = rf_model.predict(X_resampled)
        report = classification_report(y_resampled, y_pred)
        auc = roc_auc_score(y_resampled, rf_model.predict_proba(X_resampled))
        return f"Classification Report:\n{report}\nROC-AUC Score: {auc:.4f}"
    else:
        return "Random Forest model or resampled data not available."
```

output

#### Classification Report:

	precision	recall	f1-score	support
0	1.00	1.00	1.00	23364
1	1.00	1.00	1.00	23364
accuracy		1.00	1.00	46728
macro avg	1.00	1.00	1.00	46728
weighted avg	1.00	1.00	1.00	46728

ROC-AUC Score: 1.0000

## VI. GARDIO USER INTERFACE

```
# Gardio Interface Setup
with gr.Blocks() as credit_default_interface:
    with gr.Tab("Import Libraries"):
        gr.Interface(fn=import_libraries, inputs=[], outputs="text").render()
    with gr.Tab("Load Dataset"):
        gr.Interface(fn=load_dataset, inputs=[], outputs="dataframe").render()
    with gr.Tab("Display Statistics"):
        gr.Interface(fn=display_statistics, inputs=[], outputs="dataframe").render()
    with gr.Tab("Display Correlation"):
        gr.Interface(fn=display_correlation, inputs=[], outputs=gr.Image(type="pil")).render()
    with gr.Tab("Outlier Detection"):
        gr.Interface(fn=outlier_detection, inputs=[], outputs="text").render()
    with gr.Tab("Encoding Categorical Variables"):
        gr.Interface(fn=encoding_categorical, inputs=[], outputs="dataframe").render()
    with gr.Tab("Feature Scaling"):
        gr.Interface(fn=feature_scaling, inputs=[], outputs="dataframe").render()
    with gr.Tab("Apply SMOTE"):
        gr.Interface(fn=apply_smote, inputs=[], outputs="text").render()
    with gr.Tab("Calculate VIF"):
        gr.Interface(fn=calculate_vif, inputs=[], outputs="dataframe").render()
    with gr.Tab("Recalculate VIF"):
        gr.Interface(fn=recalculate_vif, inputs=[], outputs="dataframe").render()
    with gr.Tab("Visualize Distribution"):
        gr.Interface(fn=visualize_distribution, inputs=[], outputs=gr.Image(type="pil")).render()
    with gr.Tab("Initialize and Train Decision Tree"):
        gr.Interface(fn=initialize_decision_tree, inputs=[], outputs="text").render()
    with gr.Tab("Decision Tree Cross-Validation"):
        gr.Interface(fn=cross_validation_decision_tree, inputs=[], outputs="text").render()
    with gr.Tab("Visualize Decision Tree"):
        gr.Interface(fn=visualize_decision_tree, inputs=[], outputs=gr.Image(type="pil")).render()
    with gr.Tab("Decision Tree Classification Report"):
        gr.Interface(fn=classification_report_dt, inputs=[], outputs="text").render()
    with gr.Tab("Initialize and Train Logistic Regression"):
        gr.Interface(fn=initialize_logistic_regression, inputs=[], outputs="text").render()
    with gr.Tab("Evaluate Logistic Regression"):
        gr.Interface(fn=evaluate_logistic_regression, inputs=[], outputs="text").render()
    with gr.Tab("Initialize and Train Random Forest"):
        gr.Interface(fn=initialize_random_forest, inputs=[], outputs="text").render()
    with gr.Tab("Evaluate Random Forest"):
        gr.Interface(fn=evaluate_random_forest, inputs=[], outputs="text").render()

# Launch Interface
credit_default_interface.launch(share=True)
```

## VII. CONCLUSION

This project successfully built a **Credit Card Default Prediction** model using decision trees, providing valuable insights into the factors that contribute to credit default risk. The preprocessing steps, including outlier handling, one-hot encoding, feature scaling, and SMOTE, ensured that the model was well-prepared to handle the data complexities. While the model performed moderately well, future improvements could include hyperparameter tuning and trying other models like **Random Forest** or **Gradient Boosting** for potentially better accuracy.

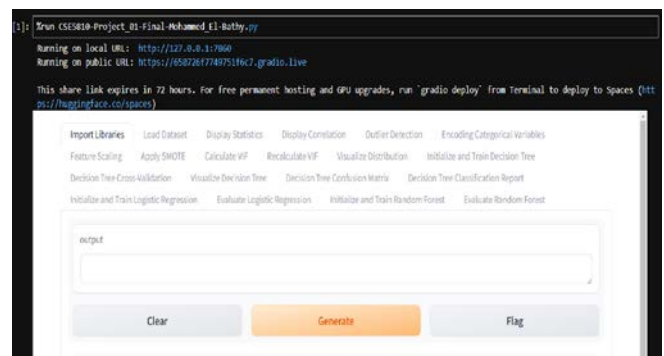
By analyzing the **confusion matrix** and the **classification report**, we can identify areas where the model's performance could be improved, particularly in reducing false positives and false negatives.

This solution is highly interpretable, making it a useful tool for financial institutions to assess and mitigate credit risk.

## APPENDIX A

The Jupyter Notebook, the Python file and this report with the data file are uploaded to **GITHUB** and accessed using the URL <https://github.com/MohammedEl-Bath/creditCardDefaultPrediction>

## STEPS EXECUTION USING GRADIO



## APPENDIX B

The execution steps of the application using Jupiter Lap are illustrated after the references

### REFERENCES

- [1] Breiman, L. (2001). Random forests. Machine learning, 45(1), 5-32.
- [2] Chawla, N. V., et al. (2002). SMOTE: Synthetic Minority Over-sampling Technique. Journal of Artificial Intelligence Research..
- [3] Handling Imbalanced Datasets. Towards Data Science. Retrieved from <https://github.com/MohammedEl-Bathy/CreditCardDefaultPrediction/tree/main>.

