

```
In [13]: from pathlib import Path
import sys, subprocess, warnings

import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns

def pip_import(pkg):
    try:
        return __import__(pkg)
    except ModuleNotFoundError:
        print(f"[INFO] Installing {pkg} ...")
        subprocess.check_call([sys.executable, "-m", "pip", "install",
                                pkg])
        return __import__(pkg)

wordcloud = pip_import("wordcloud")
shap       = pip_import("shap")
sklearn    = pip_import("sklearn")

from wordcloud import WordCloud, STOPWORDS
from sklearn.feature_extraction.text import TfidfVectorizer
from sklearn.linear_model import LogisticRegression
from sklearn.pipeline import make_pipeline
from sklearn.model_selection import train_test_split

warnings.filterwarnings("ignore")
plt.style.use("default")
```

```
In [14]: # 1 Load data
CSV = Path("fake_job_postings.csv")
if not CSV.is_file():
    sys.exit(f"[ERROR] {CSV.name} not found in {Path.cwd()}")

df = pd.read_csv(CSV)
print(f"Loaded {len(df):,} rows | {len(df.columns)} columns")

Loaded 17,880 rows | 18 columns
```

```
In [15]: # 2 Basic cleaning
if df.columns.duplicated().any():
    df = df.loc[:, ~df.columns.duplicated()].copy()

LABEL = "fraudulent"
df[LABEL] = df[LABEL].astype(int)

before = len(df)
df = df.drop_duplicates()
print("Duplicates removed:", before - len(df))

Duplicates removed: 0
```

```
In [16]: # 3 Feature engineering
TEXT_FIELDS = ["title", "company_profile", "description", "requirements"]
for col in TEXT_FIELDS:
    df[f"len_{col}"] = df[col].fillna("").str.len()

BINARY_COLS = [c for c in ["telecommuting", "has_company_logo", "has_company_logo"]]
NUM_FEATURES = [f"len_{c}" for c in TEXT_FIELDS] + BINARY_COLS
```

```
In [17]: # 4 Visuals
# Figure 1: class balance
class_counts = df[LABEL].value_counts().rename({0:"real", 1:"fake"})
class_counts.plot(kind="bar", color=["steelblue", "indianred"], figsize=(8,6))
plt.title("Real vs. Fake Job Postings"); plt.ylabel("Count"); plt.xticks([0,1]);
plt.tight_layout(); plt.savefig("class_balance.png"); plt.close()

# Figure 2: text-length distributions
plt.figure(figsize=(8,6))
for field, color in zip(["description", "requirements"], ["steelblue", "indianred"]):
    sns.kdeplot(df[f"len_{field}"], fill=True, alpha=0.3, label=field, color=color)
plt.xlim(0,5000); plt.xlabel("Character Count"); plt.title("Text-Length Distributions");
plt.legend(); plt.tight_layout(); plt.savefig("text_length_dist.png"); plt.close()

# Figure 3: correlation heat-map
plt.figure(figsize=(8,6))
corr = df[[LABEL] + NUM_FEATURES].corr().round(2)
sns.heatmap(corr, annot=True, cmap="coolwarm", vmin=-1, vmax=1)
plt.title("Correlation Matrix vs. Fraudulent Flag")
plt.tight_layout(); plt.savefig("correlation_heatmap.png"); plt.close()

# Figure 4: word clouds
for cls, fname in [(0,"top_words_real.png"), (1,"top_words_fake.png")]:
    blob = " ".join(df.loc[df[LABEL]==cls, "description"].fillna("").str.split())
    WordCloud(width=800, height=400, background_color="white", stopwords=stopwords).generate(blob).to_image().save(fname)
```

```
In [18]: # 5 Quick model
df[TEXT_COL] = df[TEXT_COL].fillna("")
X_train, X_test, y_train, y_test = train_test_split(df[TEXT_COL], df[LABEL],
                                                    test_size=0.2, random_state=42)
model = make_pipeline(TfidfVectorizer(max_features=10000, ngram_range=(1,2)),
                      LogisticRegression(max_iter=1000, class_weight='balanced'))
model.fit(X_train, y_train)
print("Validation accuracy:", round(model.score(X_test, y_test), 3))

# SHAP explanation for ad in test set
explainer = shap.LinearExplainer(model.named_steps["logisticregression"],
                                model.named_steps["tfidfvectorizer"].get_feature_names_out(),
                                feature_perturbation="interventional")

idx = X_test[y_test==1].index[0]
row_vec = model.named_steps["tfidfvectorizer"].transform([X_test.loc[idx, TEXT_COL]])
shap_values = explainer(row_vec)
plt.figure(figsize=(10,4))
shap.plots.waterfall(shap_values[0], max_display=12, show=False)
plt.title("Why this listing is predicted fake")
plt.tight_layout(); plt.savefig("shap_waterfall.png"); plt.close()

print("EDA + SHAP waterfall complete – images saved.")
```

Validation accuracy: 0.96  
EDA + SHAP waterfall complete – images saved.

In [ ]: