

train_model.py

```
1 import pandas as pd
2 import numpy as np
3 from sklearn.ensemble import RandomForestRegressor
4 from sklearn.preprocessing import LabelEncoder, OneHotEncoder
5 from sklearn.compose import ColumnTransformer
6 from sklearn.pipeline import Pipeline
7 from sklearn.model_selection import train_test_split
8 import joblib
9 import warnings
10 warnings.filterwarnings('ignore')
11
12 # Load dataset (replace with your CSV path)
13 df = pd.read_csv('digital_nomad_salaries.csv')
14
15 # Data Cleaning
16 df = df.drop(['user_id', 'timestamp'], axis=1) # Remove non-predictive columns
17 df['company_remote'] = df['company_remote'].map({'Y': 1, 'N': 0})
18 df['nomad_visa'] = df['nomad_visa'].map({'Y': 1, 'N': 0})
19
20 # Feature Engineering
21 df['city'] = df['location'].apply(lambda x: x.split(',')[0].strip())
22 df['country'] = df['location'].apply(lambda x: x.split(',')[1].strip())
23
24 # Preprocessing pipeline
25 categorical_features = ['job_role', 'city', 'country']
26 numeric_features = ['productivity', 'burnout_level', 'company_remote', 'nomad_visa']
27
28 preprocessor = ColumnTransformer(
29     transformers=[
30         ('cat', OneHotEncoder(handle_unknown='ignore'), categorical_features),
31         ('num', 'passthrough', numeric_features)
32     ])
33
34 # Model Pipeline
35 model = Pipeline([
36     ('preprocessor', preprocessor),
37     ('regressor', RandomForestRegressor(n_estimators=100, random_state=42))
38 ])
39
40 # Train-Test Split
41 X = df.drop('salary_usd', axis=1)
42 y = df['salary_usd']
43 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
44
45 # Train model
46 model.fit(X_train, y_train)
47
48 # Save model and preprocessing pipeline
49 joblib.dump(model, 'nomad_salary_model.joblib')
50
51 print(f"Model trained and saved. R2 Score: {model.score(X_test, y_test):.2f}")
52
```