

Predicción de precio en propiedades de Airbnb

Mario Emilio Jiménez Vizcaíno
A01173359@itesm.mx
Ingeniería en Tecnologías Computacionales
Tecnológico de Monterrey

Alberto Dávila Almaguer
A01196987@itesm.mx
Ingeniería en Innovación y Desarrollo
Tecnológico de Monterrey

Franco Daniel Pérez Reyes
A00822080@itesm.mx
Ingeniería en Mecatrónica
Tecnológico de Monterrey

Carlos Andrés Luna Leyva
A00826148@itesm.mx
Ingeniería en Tecnologías Computacionales
Tecnológico de Monterrey

ABSTRACT

En los últimos años han surgido plataformas en las cuales anfitriones pueden publicitar y contratar el arriendo de sus propiedades con sus huéspedes; anfitriones y huéspedes pueden valorarse mutuamente, como referencia para futuros usuarios. Una de las primeras, y la más exitosa en Estados Unidos es AirBnB, fundada en el año 2008. En este proyecto se utilizarán técnicas de aprendizaje automático con el fin de predecir y optimizar una serie de métricas sobre propiedades, utilizando un dataset de más de 70 mil propiedades listadas desde noviembre de 2008 hasta octubre 2017 en esta plataforma.

1 INTRODUCCIÓN

El precio de cualquier producto o servicio se ve afectado por variables como oferta, demanda, inflación en el país, etc. Para bienes inmuebles esto es mucho más complejo, ya que entran en consideración factores como la zona donde se encuentra la propiedad, el espacio total que abarca, el número de habitaciones, etc. Cada una de estas variables puede o no afectar a las demás y juntas caracterizan un sistema digno de ser analizado.

2 CONCEPTOS PREVIOS

- Conocimientos básicos de estadística
- Programación básica en el lenguaje Python
- Conocimientos sobre librerías como `scikit-learn`, `pandas` y `numpy`

3 METODOLOGÍA

TODO

4 RESULTADOS

TODO

5 CONCLUSIONES

TODO

REFERENCES

A CÓDIGO DE EJECUCIÓN DE LOS REGRESORES

```

1 import pathlib
2 import numpy as np
3 import pandas as pd
4 from sklearn.ensemble import GradientBoostingRegressor, RandomForestRegressor, VotingRegressor
5 from sklearn.linear_model import LinearRegression
6 from sklearn.metrics import mean_squared_error
7 from sklearn.model_selection import train_test_split
8 from sklearn.neighbors import KNeighborsRegressor
9
10 df_base = pd.read_csv('airbnb_clean.csv').set_index("id")
11 df_ammenities = pd.read_csv('airbnb_amenities_clean.csv').set_index("id")
12 df = pd.concat([df_base, df_ammenities], axis=1)
13
14 X = df.drop(columns=["log_price", "host_since", "host_has_profile_pic", "city",
15                    "host_identity_verified", "property_type", "room_type"])
16 Y = df["log_price"]
17
18 X["cleaning_fee"] = X["cleaning_fee"].astype(int)
19 X["instant_bookable"] = X["instant_bookable"].astype(int)
20 X[df_ammenities.columns] = X[df_ammenities.columns].astype(int)
21
22 rf_mses = []
23 gb_mses = []
24 lr_mses = []
25 kn_mses = []
26 vr_mses = []
27
28 for i in range(30):
29     print(f"Calculating scores for iteration {i+1}")
30     X_train, X_test, y_train, y_test = train_test_split(
31         X, Y, test_size=0.2, random_state=i)
32
33     random_forest = RandomForestRegressor(n_estimators=100, n_jobs=-1)
34     random_forest.fit(X_train, y_train)
35     rf_pred = random_forest.predict(X_test)
36     rf_mses.append(mean_squared_error(y_test, rf_pred))
37
38     gradient_boosting = GradientBoostingRegressor(n_estimators=150)
39     gradient_boosting.fit(X_train, y_train)
40     gb_pred = gradient_boosting.predict(X_test)
41     gb_mses.append(mean_squared_error(y_test, gb_pred))
42
43     linear_reg = LinearRegression(fit_intercept=True)
44     linear_reg.fit(X_train, y_train)
45     lr_pred = linear_reg.predict(X_test)
46     lr_mses.append(mean_squared_error(y_test, lr_pred))
47
48     k_neighbors = KNeighborsRegressor(n_neighbors=30, n_jobs=-1)
49     k_neighbors.fit(X_train, y_train)
50     kn_pred = k_neighbors.predict(X_test)
51     kn_mses.append(mean_squared_error(y_test, kn_pred))
52
53     voting_reg = VotingRegressor([("RF", random_forest), ("GB", gradient_boosting),
54                                ("KN", k_neighbors), ("LN", linear_reg)], n_jobs=-1)
55     voting_reg.fit(X_train, y_train)
56     vr_pred = voting_reg.predict(X_test)
57     vr_mses.append(mean_squared_error(y_test, vr_pred))
58
59     print(f"{rf_mses[i]:2.5f} {gb_mses[i]:2.5f}", end=" ")
60     print(f"{lr_mses[i]:2.5f} {kn_mses[i]:2.5f} {vr_mses[i]:2.5f}")

```

B SALIDA DE EJECUCIÓN DE LOS REGRESORES

```
1 Calculating scores for iteration 1
2 0.17583 0.18862 0.29061 0.27962 0.19436
3 Calculating scores for iteration 2
4 0.17881 0.18835 0.28611 0.27449 0.19290
5 Calculating scores for iteration 3
6 0.17301 0.18280 0.27537 0.26925 0.18600
7 Calculating scores for iteration 4
8 0.17623 0.18881 0.28657 0.28166 0.19390
9 Calculating scores for iteration 5
10 0.17483 0.18737 0.28715 0.27676 0.19200
11 Calculating scores for iteration 6
12 0.17366 0.18453 0.28149 0.27784 0.19029
13 Calculating scores for iteration 7
14 0.17639 0.18611 0.28709 0.27331 0.19206
15 Calculating scores for iteration 8
16 0.17819 0.18926 0.29059 0.28183 0.19590
17 Calculating scores for iteration 9
18 0.17487 0.18631 0.28020 0.27405 0.18994
19 Calculating scores for iteration 10
20 0.17268 0.18663 0.28243 0.27379 0.18992
21 Calculating scores for iteration 11
22 0.17527 0.18710 0.28620 0.27405 0.19206
23 Calculating scores for iteration 12
24 0.17455 0.18893 0.28500 0.27107 0.19063
25 Calculating scores for iteration 13
26 0.17531 0.18802 0.28308 0.27220 0.19208
27 Calculating scores for iteration 14
28 0.17485 0.18806 0.28952 0.27267 0.19275
29 Calculating scores for iteration 15
30 0.17623 0.18632 0.28253 0.27538 0.19105
31 Calculating scores for iteration 16
32 0.17335 0.18552 0.27961 0.27192 0.18910
33 Calculating scores for iteration 17
34 0.17291 0.18649 0.28193 0.27358 0.18973
35 Calculating scores for iteration 18
36 0.17437 0.18559 0.28341 0.27621 0.19148
37 Calculating scores for iteration 19
38 0.17343 0.18511 0.28152 0.27081 0.18863
39 Calculating scores for iteration 20
40 0.17776 0.18946 0.28510 0.27851 0.19392
41 Calculating scores for iteration 21
42 0.18095 0.19456 0.29053 0.27785 0.19696
43 Calculating scores for iteration 22
44 0.17118 0.18228 0.27875 0.26795 0.18694
45 Calculating scores for iteration 23
46 0.17113 0.18675 0.28304 0.27240 0.19027
47 Calculating scores for iteration 24
48 0.17124 0.18415 0.28199 0.27108 0.18764
49 Calculating scores for iteration 25
50 0.17631 0.18797 0.28490 0.27225 0.19187
51 Calculating scores for iteration 26
52 0.17590 0.18882 0.28670 0.27830 0.19293
53 Calculating scores for iteration 27
54 0.17297 0.18495 0.27996 0.26918 0.18770
55 Calculating scores for iteration 28
56 0.17261 0.18690 0.28336 0.27616 0.19117
57 Calculating scores for iteration 29
58 0.17671 0.19054 0.28790 0.27424 0.19418
59 Calculating scores for iteration 30
60 0.17552 0.18500 0.28215 0.27992 0.19099
```

C CÓDIGO DE PRUEBA DE WILCOXON

```

1  from pathlib import Path
2  import numpy as np
3  import pandas as pd
4  from scipy.stats import ranksums
5
6  mses_dict = {}
7
8  for mses in Path('mses').iterdir():
9      if mses.is_file():
10         mses_dict[mses.stem] = np.loadtxt(mses)
11
12  if len(mses_dict) == 0:
13      print("No MSEs files found")
14      print("Run `generate_mses.py` file first")
15
16  matrix_greater = []
17  matrix_less = []
18
19  for mses1 in mses_dict.values():
20      row_greater = []
21      row_less = []
22
23      for mses2 in mses_dict.values():
24          row_greater.append(
25              ranksums(mses1, mses2, alternative="greater").pvalue)
26          row_less.append(ranksums(mses1, mses2, alternative="less").pvalue)
27
28      matrix_greater.append(row_greater)
29      matrix_less.append(row_less)
30
31  df_greater = pd.DataFrame(
32      matrix_greater, index=mses_dict.keys(), columns=mses_dict.keys())
33  df_less = pd.DataFrame(
34      matrix_less, index=mses_dict.keys(), columns=mses_dict.keys())
35
36  pd.options.display.float_format = "{:.14f}".format
37  print("Wilcoxon rank sums test with a greater hypothesis")
38  print(df_greater)
39  print()
40
41  print("Wilcoxon rank sums test with a less hypothesis")
42  print(df_less)

```

D SALIDA DE EJECUCIÓN DE PRUEBA DE WILCOXON

D.1 Wilcoxon rank sums test with a greater hypothesis

	RandomForest	GradientBoosting	LinearRegression	KNeighbors	VotingRegressor
RandomForest	0.5	0.99999999998564	0.99999999998564	0.99999999998564	0.99999999998564
GradientBoosting	0.00000000001436	0.5	0.99999999998564	0.99999999998564	0.99999991000660
LinearRegression	0.00000000001436	0.00000000001436	0.5	0.00000000022015	0.00000000001436
KNeighbors	0.00000000001436	0.00000000001436	0.9999999977985	0.5	0.00000000001436
VotingRegressor	0.00000000001436	0.00000008999340	0.99999999998564	0.99999999998564	0.5

D.2 Wilcoxon rank sums test with a less hypothesis

	RandomForest	GradientBoosting	LinearRegression	KNeighbors	VotingRegressor
RandomForest	0.5	0.00000000001436	0.00000000001436	0.00000000001436	0.00000000001436
GradientBoosting	0.99999999998564	0.5	0.00000000001436	0.00000000001436	0.00000008999340
LinearRegression	0.99999999998564	0.99999999998564	0.5	0.9999999977985	0.99999999998564
KNeighbors	0.99999999998564	0.99999999998564	0.00000000022015	0.5	0.99999999998564
VotingRegressor	0.99999999998564	0.99999991000660	0.00000000001436	0.00000000001436	0.5