

哈尔滨工业大学（深圳）

大数据导论大作业报告

**题目：基于随机森林回归预测的房价
预测**

姓 名 陈致远

学 号 200111521

报告日期 2022-10-28

一、实验目的

本次实验主要在题目中给出的各种类型的房价数据进行预处理，然后利用机器学习中的一种预测算法对其进行预测，对模型进行评估，在训练好的模型上对其进行预测房价。

二、实验内容分析

①数据预处理

在查看 `train.csv` 文件时，发现“公寓号”一列存在着大量的空缺，认为其对模型没有直接的帮助，可以直接删除。

并且房价跟出售日期没有直接的关系，考虑直接删去，不再保留。

可以认为修建年份对房价有比较大的影响，在查看修建年份中，发现有 3444 条为 0 的扰动数据，将其所在行全部删除，避免对训练模型造成干扰。

查看到地役权几乎全部空白，予以删除。

删除当前税收级别为空白的两行数据。

删除出售价格小于 1000 的行。

考虑到测试集里面土地平方英尺和总平方英尺同样存在着空缺数据，所以将训练集里面的两个指标所在列全部删除。

考虑到邮箱编码可以由所属区域和社区确定，所以删除邮箱编码一列。

删除税收级别为空白的 439 条数据。

②进行相关编码操作， 对社区的编码：

AIRPORT LA GUARDIA	1	DONGAN HILLS-OLD TOWN	6 1	INWOOD	12 1	PELHAM PARKWAY NORTH	181	WHITESTONE	241
ALPHABET CITY	2	DOUGLASTON	6 2	JACKSON HEIGHTS	12 2	PELHAM PARKWAY SOUTH	182	WILLIAMSBRIDGE	242
ANNADALE	3	DOWNTOWN-FULTON FERRY	6 3	JAMAICA	12 3	PLEASANT PLAINS	183	WILLIAMSBURG-CENTRAL	243
ARDEN HEIGHTS	4	DOWNTOWN-FULTON MALL	6 4	JAMAICA BAY	12 4	PORT IVORY	184	WILLIAMSBURG-EAST	244
ARROCHAR	5	DOWNTOWN-METROTECH	6 5	JAMAICA ESTATES	12 5	PORT RICHMOND	185	WILLIAMSBURG-NORTH	245
ARROCHAR-SHORE ACRES	6	DYKER HEIGHTS	6 6	JAMAICA HILLS	12 6	PRINCES BAY	186	WILLIAMSBURG-SOUTH	246
ARVERNE	7	EAST ELMHURST	6 7	JAVITS CENTER	12 7	PROSPECT HEIGHTS	187	WILLOWBROOK	247
ASTORIA	8	EAST NEW YORK	6 8	KENSINGTON	12 8	QUEENS VILLAGE	188	WINDSOR TERRACE	248
BATH BEACH	9	EAST RIVER	6 9	KEW GARDENS	12 9	RED HOOK	189	WOODHAVEN	249

BATHGATE	1 0	EAST TREMONT	7 0	KINGSBRIDGE HTS/UNIV HTS	13 0	REGO PARK	190	WOODLAWN	250
BAY RIDGE	1 1	EAST VILLAGE	7 1	KINGSBRIDGE /JEROME PARK	13 1	RICHMOND HILL	191	WOODROW	251
BAYCHESTER	1 2	ELMHURST	7 2	KIPS BAY	13 2	RICHMONDTOWN	192	WOODSIDE	252
BAYSIDE	1 3	ELTINGVILLE	7 3	LAURELTON	13 3	RICHMONDTOWN-LIGHTS HILL	193	WYCKOFF HEIGHTS	253
BEDFORD PARK/NORWOOD	1 4	EMERSON HILL	7 4	LITTLE ITALY	13 4	RIDGEWOOD	194		
BEDFORD STUYVESANT	1 5	FAR ROCKAWAY	7 5	LITTLE NECK	13 5	RIVERDALE	195		
BEECHHURST	1 6	FASHION	7 6	LIVINGSTON	13 6	ROCKAWAY PARK	196		
BELLE HARBOR	1 7	FIELDSTON	7 7	LONG ISLAND CITY	13 7	ROOSEVELT ISLAND	197		
BELLEROSE	1 8	FINANCIAL	7 8	LOWER EAST SIDE	13 8	ROSEBANK	198		
BELMONT	1 9	FLATBUSH- CENTRAL	7 9	MADISON	13 9	ROSEDALE	199		
BENSONHURST	2 0	FLATBUSH- EAST	8 0	MANHATTAN BEACH	14 0	ROSSVILLE	200		
BERGEN BEACH	2 1	FLATBUSH- LEFFERTS GARDEN	8 1	MANHATTAN VALLEY	14 1	ROSSVILLE-CHARLESTON	201		
BLOOMFIELD	2 2	FLATBUSH- NORTH	8 2	MANOR HEIGHTS	14 2	ROSSVILLE-PORT MOBIL	202		
BOERUM HILL	2 3	FLATIRON	8 3	MARINE PARK	14 3	ROSSVILLE-RICHMOND VALLEY	203		
BOROUGH PARK	2 4	FLATLANDS	8 4	MARINERS HARBOR	14 4	SCHUYLERVILLE/PELHAM BAY	204		
BRIARWOOD	2 5	FLORAL PARK	8 5	MASPETH	14 5	SEAGATE	205		
BRIGHTON BEACH	2 6	FLUSHING MEADOW PARK	8 6	MELROSE/CO NCOURSE	14 6	SHEEPSHEAD BAY	206		
BROAD CHANNEL	2 7	FLUSHING- NORTH	8 7	MIDDLE VILLAGE	14 7	SILVER LAKE	207		
BRONXDALE	2 8	FLUSHING- SOUTH	8 8	MIDLAND BEACH	14 8	SO. JAMAICA-BAISLEY PARK	208		
BROOKLYN HEIGHTS	2 9	FORDHAM	8 9	MIDTOWN CBD	14 9	SOHO	209		

BROWNSVILLE	3 0	FOREST HILLS	9 0	MIDTOWN EAST	15 0	SOUNDVIEW	210
BULLS HEAD	3 1	FORT GREENE	9 1	MIDTOWN WEST	15 1	SOUTH BEACH	211
BUSH TERMINAL	3 2	FRESH KILLS	9 2	MIDWOOD	15 2	SOUTH JAMAICA	212
BUSHWICK	3 3	FRESH MEADOWS	9 3	MILL BASIN	15 3	SOUTH OZONE PARK	213
CAMBRIA HEIGHTS	3 4	GERRITSEN BEACH	9 4	MORNINGSID E HEIGHTS	15 4	SOUTHBRIDGE	214
CANARSIE	3 5	GLEN OAKS	9 5	MORRIS PARK/VAN NEST	15 5	SPRING CREEK	215
CARROLL GARDENS	3 6	GLENDALE	9 6	MORRISANIA /LONGWOOD	15 6	SPRINGFIELD GARDENS	216
CASTLE HILL/UNIONPORT	3 7	GOWANUS	9 7	MOTT HAVEN/PORT MORRIS	15 7	ST. ALBANS	217
CASTLETON CORNERS	3 8	GRAMERCY	9 8	MOUNT HOPE/MOUNT EDEN	15 8	STAPLETON	218
CHELSEA	3 9	GRANT CITY	9 9	MURRAY HILL	15 9	STAPLETON-CLIFTON	219
CHINATOWN	4 0	GRASMERE	1 0 0	NAVY YARD	16 0	SUNNYSIDE	220
CITY ISLAND	4 1	GRAVESEND	1 0 1	NEPONSIT	16 1	SUNSET PARK	221
CITY ISLAND - PELHAM STRIP	4 2	GREAT KILLS	1 0 2	NEW BRIGHTON	16 2	THROGS NECK	222
CIVIC CENTER	4 3	GREAT KILLS- BAY TERRACE	1 0 3	NEW BRIGHTON- ST. GEORGE	16 3	TODT HILL	223
CLINTON	4 4	GREENPOINT	1 0 4	NEW DORP	16 4	TOMPKINSVILLE	224
CLINTON HILL	4 5	GREENWICH VILLAGE- CENTRAL	1 0 5	NEW DORP- BEACH	16 5	TOTTENVILLE	225

CLOVE LAKES	4 6	GREENWICH VILLAGE-WEST	1 0 6	NEW DORP- HEIGHTS	16 6	TRAVIS	226
COBBLE HILL	4 7	GRYMES HILL	1 0 7	NEW SPRINGVILLE	16 7	TRIBECA	227
COBBLE HILL- WEST	4 8	HAMMELS	1 0 8	OAKLAND GARDENS	16 8	UPPER EAST SIDE (59-79)	228
COLLEGE POINT	4 9	HARLEM- CENTRAL	1 0 9	OAKWOOD	16 9	UPPER EAST SIDE (79-96)	229
CONCORD	5 0	HARLEM-EAST	1 1 0	OAKWOOD- BEACH	17 0	UPPER EAST SIDE (96-110)	230
CONCORD-FOX HILLS	5 1	HARLEM- UPPER	1 1 1	OCEAN HILL	17 1	UPPER WEST SIDE (59-79)	231
CONEY ISLAND	5 2	HARLEM-WEST	1 1 2	OCEAN PARKWAY- NORTH	17 2	UPPER WEST SIDE (79-96)	232
CO-OP CITY	5 3	HIGHBRIDGE/M ORRIS HEIGHTS	1 1 3	OCEAN PARKWAY- SOUTH	17 3	UPPER WEST SIDE (96-116)	233
CORONA	5 4	HILLCREST	1 1 4	OLD MILL BASIN	17 4	VAN CORTLANDT PARK	234
COUNTRY CLUB	5 5	HOLLIS	1 1 5	OZONE PARK	17 5	WAKEFIELD	235
CROTONA PARK	5 6	HOLLIS HILLS	1 1 6	PARK SLOPE	17 6	WASHINGTON HEIGHTS LOWER	236
CROWN HEIGHTS	5 7	HOLLISWOOD	1 1 7	PARK SLOPE SOUTH	17 7	WASHINGTON HEIGHTS UPPER	237
CYPRESS HILLS	5 8	HOWARD BEACH	1 1 8	PARKCHESTE R	17 8	WEST NEW BRIGHTON	238
DONGAN HILLS	5 9	HUGUENOT	1 1 9	PELHAM BAY	17 9	WESTCHESTER	239

DONGAN HILLS- COLONY	6 0	HUNTS POINT	1 2 0	PELHAM GARDENS	18 0	WESTERLEIGH	240
-------------------------	--------	-------------	-------------	-------------------	---------	-------------	-----

对建筑类型的编码:

01 ONE FAMILY DWELLINGS	1
02 TWO FAMILY DWELLINGS	2
03 THREE FAMILY DWELLINGS	3
04 TAX CLASS 1 CONDOS	4
05 TAX CLASS 1 VACANT LAND	5
06 TAX CLASS 1 - OTHER	6
07 RENTALS - WALKUP APARTMENTS	7
08 RENTALS - ELEVATOR APARTMENTS	8
09 COOPS - WALKUP APARTMENTS	9
10 COOPS - ELEVATOR APARTMENTS	10
11A CONDO-RENTALS	11
12 CONDOS - WALKUP APARTMENTS	12
13 CONDOS - ELEVATOR APARTMENTS	13
14 RENTALS - 4-10 UNIT	14
15 CONDOS - 2-10 UNIT RESIDENTIAL	15
16 CONDOS - 2-10 UNIT WITH COMMERCIAL UNIT	16
17 CONDO COOPS	17
21 OFFICE BUILDINGS	18
22 STORE BUILDINGS	19
23 LOFT BUILDINGS	20
25 LUXURY HOTELS	21
26 OTHER HOTELS	22
27 FACTORIES	23
28 COMMERCIAL CONDOS	24
29 COMMERCIAL GARAGES	25
30 WAREHOUSES	26
31 COMMERCIAL VACANT LAND	27
32 HOSPITAL AND HEALTH FACILITIES	28
33 EDUCATIONAL FACILITIES	29
34 THEATRES	30
35 INDOOR PUBLIC AND CULTURAL FACILITIES	31
36 OUTDOOR RECREATIONAL FACILITIES	32
37 RELIGIOUS FACILITIES	33
38 ASYLUMS AND HOMES	34
39 TRANSPORTATION FACILITIES	35
40 SELECTED GOVERNMENTAL FACILITIES	36
41 TAX CLASS 4 - OTHER	37

42 CONDO	38
CULTURAL/MEDICAL/EDUCATIONAL/ETC	
43 CONDO OFFICE BUILDINGS	39
44 CONDO PARKING	40
45 CONDO HOTELS	41
46 CONDO STORE BUILDINGS	42
47 CONDO NON-BUSINESS STORAGE	43
48 CONDO TERRACES/GARDENS/CABANAS	44
49 CONDO WAREHOUSES/FACTORY/INDUS	45

建筑类别的编码:

A0	1	D6	31	H8	61	M4	91	RG	121
A1	2	D7	32	HB	62	M9	92	RH	122
A2	3	D8	33	HH	63	N2	93	RK	123
A3	4	D9	34	HR	64	N9	94	RP	124
A4	5	E1	35	HS	65	O1	95	RR	125
A5	6	E2	36	I1	66	O2	96	RS	126
A6	7	E7	37	I3	67	O3	97	RT	127
A7	8	E9	38	I4	68	O4	98	RW	128
A9	9	F1	39	I5	69	O5	99	S0	129
B1	10	F2	40	I6	70	O6	100	S1	130
B2	11	F4	41	I7	71	O7	101	S2	131
B3	12	F5	42	I9	72	O8	102	S3	132
B9	13	F9	43	J1	73	P2	103	S4	133
C0	14	G0	44	J8	74	P5	104	S5	134
C1	15	G1	45	J9	75	P6	105	S9	135
C2	16	G2	46	K1	76	P8	106	T2	136
C3	17	G3	47	K2	77	P9	107	V0	137
C4	18	G4	48	K3	78	Q1	108	V1	138
C5	19	G5	49	K4	79	Q8	109	V3	139
C6	20	G6	50	K5	80	Q9	110	V6	140
C7	21	G7	51	K6	81	R1	111	V9	141
C8	22	G8	52	K7	82	R2	112	W1	142
C9	23	G9	53	K9	83	R3	113	W2	143
CM	24	GU	54	L1	84	R4	114	W3	144
D0	25	GW	55	L3	85	R5	115	W4	145
D1	26	H1	56	L8	86	R6	116	W8	146
D2	27	H2	57	L9	87	R8	117	W9	147
D3	28	H3	58	M1	88	R9	118	Y1	148
D4	29	H4	59	M2	89	RA	119	Y3	149
D5	30	H6	60	M3	90	RB	120	Z0	150
								Z2	151
								Z9	152

③模型训练

本次作业第一次试探使用线性回归模型，但是测试结果的拟合优度 R^2 一直小于 0.1，效果很差，所以决定采用随机森林回归模型，原理如下：

随机森林回归模型由多棵回归树构成，且森林中的每一棵决策树之间没有关联，模型的最终输出由森林中的每一棵决策树共同决定。

随机森林的随机性体现在两个方面：

1、样本的随机性，从训练集中随机抽取一定数量的样本，作为每颗回归树的根节点样本；

2、特征的随机性，在建立每颗回归树时，随机抽取一定数量的候选特征，从中选择最合适的特征作为分裂节点。

算法原理如下：

(a) 从训练样本集 S 中随机的抽取 m 个样本点，得到一个新的 $S_1 \cdots S_n$ 个子训练集；

(b) 用子训练集，训练一个 CART 回归树(决策树)，这里在训练的过程中，对每个节点的切分规则是先从所有特征中随机的选择 k 个特征，然后在从这 k 个特征中选择最优的切分点在做左右子树的划分。(这里的得到决策树都是二叉树)

(c) 通过第二步，可以生成很多个 CART 回归树模型。

(d) 每一个 CART 回归树最终的预测结果为该样本点所到叶节点的均值。

(e) 随机森林最终的预测结果为所有 CART 回归树预测结果的均值。

三、 实验过程及结果

在最开始进行实验时，使用随机森林回归模型进行训练测试得到的 RMSE 虽然比较大，但是查看 R^2 已经达到 0.8 多，所以决定采用此算法，在数据基础上进行多次训练。由于房价数额普遍比较大，所以在查看 MSE、RMSE 两个指标之外，另外新建一个百分比指标，如下：

$$y_error_percentage = \frac{|Y_pred - Y_test|}{Y_test}$$

该指标反映出预测的差错与原始房价对比误差的占比，下面给出训练集上的前部分指标数据：

	y_true	y_pred	y_error	y_error_percentage
0	900000	815800.0	84200.0	0.093556
1	1440000	1762000.0	322000.0	0.223611
2	499500	504600.0	5100.0	0.010210
3	11015000	12573172.9	1558172.9	0.141459
4	160000	253280.0	93280.0	0.583000

测试集上的前部分指标数据：

	y_true	y_pred	y_error	y_error_percentage
0	610000	6.450000e+05	35000.000000	0.057377
1	1185000	1.070858e+06	114141.666667	0.096322
2	960000	1.248000e+06	288000.000000	0.300000
3	110000	2.582667e+05	148266.666667	1.347879
4	438000	4.376570e+05	343.000000	0.000783

可以看出大部分 `y_error_percentage` 都很小，在容忍范围之内。

使用 `train_test_split` 方法，按测试数据集占 0.1 的比例对数据进行分割,计算训练集上的平均绝对误差(MSE)、 R^2 决定系数、均方根误差 RMSE 三个指标，最开始时查看到训练集上的 R^2 决定系数高达 0.845，但是测试集上的决定系数仅有 0.345，合理怀疑是训练集上的脏数据导致。于是每轮对训练集上的 `y_error_percentage` 最大的 100 行数据进行删除，并且每删除一轮就测试在训练集和测试集上的表现效果。

在经过多轮数据筛选、删除敏感值之后，在训练集上得到的三个指标数值如下：

```
{ '平均绝对误差(MSE)': [3461354204184.418], 'R2 决定系数': [0.881], '均方根误差RMSE': [1860471.501]}
```

在测试集上的三个指标如下：

```
{ '平均绝对误差(MSE)': [3983560516191.027], 'R2 决定系数': [0.596], '均方根误差RMSE': [1995885.898]}
```

其中，由于房价的差距比较悬殊，所以这里只关注 R^2 .在训练集上达到 0.881

说明预测效果十分良好，可以接受，在测试集上的 R^2 也达到了 0.596，说明预测精度可以接受（继续删除发现测试集的决定系数开始下降了），在接下来使用这个训练好的模型去预测测试集里面的房价。

对测试集采用与问题一同样的方法进行编码。

查看到测试集的税收级别里面有空白值，使用税收级别里面的中位数 5 进行填充。

下面给出预测结果的部分展示图，详细数据可以参考附件。

支撑材料列表：

数值化后的 `train.xlsx`:处理后的训练集数据

数值化后的 `test.xlsx`:处理和预测后的测试集数据

大数据导论作业.py:使用上面两个数据集进行模型训练的过程

Excel 代码.txt:数据预处理过程中使用到的宏代码

Submit.csv:预测结果表格

四、 实验心得

这次大作业完成之后我对于机器学习的预测算法有了更深层次的理解，对机器学习预测模型的一些指标如 MSE 、 $RMSE$ 、 R^2 有了更为深刻的理解，同时也对于在 python 中利用 numpy、pandas、sklearn 库处理数据更加的熟练掌握，感谢这次大作业提供的机会。