

پروژه داده کاوی

استاد راهنما: دکتر محمدرضا فقیهی حبیب آبادی

گرد آورنده: آرش سجادی، شماره دانشجویی: ۴۰۰۴۲۲۰۹۶

دانشکده ریاضی دانشگاه شهید بهشتی

۳۰ خرداد ۱۴۰۱

چکیده:

یکی از مشکلاتی که در هنگام خرید و فروش خانه با آن مواجه هستیم، این است که مرجع خوبی برای قیمت گذاری خانه ها وجود ندارد. اغلب مشاورین املاک با توجه به در نظر گرفتن منفعت خود، اقدام به قیمت گذاری های جانب دارانه در پیشگاه خریدار و فروشنده می کنند. در این پروژه تلاش خواهیم کرد یک سیستم پیشگویی منصفانه از قیمت خانه را در شهر تهران ارائه دهیم. منبع اصلی من در کل این پروژه کتاب داده کاوی برای تحلیل خودکار کسب و کار: مفاهیم، فنون و کاربردهای R [۱] است. به علاوه منبع من برای تولید واژه نامه، پس از واژه نامه منبع اصلی مذکور، واژه نامه رسمی انجمن ریاضی ایران [۲] و انجمن آمار ایران [۳] است.

کلمات کلیدی: داده کاوی، پیشگویی قیمت مسکن

فهرست مطالب

۷	۱-۲	بافت نگار ویژگی ها	۲	۱	مقدمات و معرفی داده ها
۹	۲-۲	نمودار جعبه ای	۲	۱-۱	هدف داده کاوی
۱۲	۳-۲	نمودار حرارتی	۲	۲-۱	تعریف هر متغیر
۱۲	۴-۲	تصویری سازی چند بُعدی	۳	۳-۱	نمونه یک ثبت
۱۴	۵-۲	نمودار با محورهای موازی	۳	۴-۱	خلاصه وضعیت ویژگی ها
۱۵	۶-۲	کاهش بُعد	۴	۵-۱	پیش پردازش داده ها
۱۵	۷-۲	تجزیه و تحلیل مؤلفه اصلی	۴	۱-۵-۱	ویژگی آدرس
۱۶	۸-۲	داده کاهی	۵	۲-۵-۱	مقادیر گم شده
۱۷	۹-۲	اعمال تغییرات نتیجه گیری شده	۵	۳-۵-۱	داده های دور افتاده
۱۸	۳	رگرسیون خطی چندگانه	۵	۴-۵-۱	افراز داده ها
۲۱	۴	مدل درخت تصمیم	۵	۵-۵-۱	نرمالیده کردن متغیرها
۲۳	۵	مدل k نزدیک ترین همسایه	۶	۶-۵-۱	کدگذاری متغیرهای رسته ای
۲۴	۶	مدل شبکه عصبی	۷	۲	تصویری سازی و اکتشاف داده ها
۲۶	۷	انتخاب مدل نهایی			

مراجع

۲۸

واژه نامه انگلیسی به فارسی

۲۹

واژه نامه فارسی به انگلیسی

۳۱

مقدمه

داده‌های من ثبت^۲‌های مشخصات ۳۴۷۹ خانه در شهر تهران است. این داده‌ها از وبسایت کگل [۴] گرفته شده است. شخصی که داده‌ها را بارگذاری کرده است مدعی شده است که این داده‌ها تحت نظر دانشگاه تربیت مدرس با تکیه بر داده‌های وبسایت‌های معاملات مسکن گرفته شده است. تاریخ ثبت داده‌ها همگی مربوط به آذر ماه سال ۱۴۰۰ می‌باشند. هر ثبت در این داده‌ها شامل ۸ ویژگی^۳ است و هدف نهایی پروژه، پیشگویی قیمت مسکن در شهر تهران است.

۱. مقدمات و معرفی داده‌ها

هر ثبت در داده‌های این پروژه مشخص‌کننده ویژگی‌های یک خانه در شهر تهران است که برای فروش در بنگاه‌های معاملاتی آنلاین مکتوب گردیده. در ادامه به ویژگی‌های توصیفگر این ثبت‌ها اشاره خواهم کرد. اما قبل از هر چیز بایستی مشخص شود چرا این پروژه ارزش این را دارد که برای آن وقت صرف کنیم.

۱-۱ هدف داده کاوی

به گزارش باشگاه خبرنگاران جوان [۵] وجود مشاورین املاک در کشور نه تنها موجب کنترل قیمت‌ها نشده است بلکه علیرغم ادعای رئیس اتحادیه مشاوران املاک کشور، در موارد بسیاری به جای تنظیم بازار ملک به ساخت یک حباب مثبت یا منفی دامن زده است. از این رو وجود یک سیستم بی طرف برای قیمت‌گذاری منصفانه مسکن می‌تواند به این آشفتگی بازار مسکن سروسامان ببخشد و در صورتی که کار جدی‌تر با داده‌ی بیشتر و دقیق‌تر روی این قبیل پروژه‌ها صورت بگیرد، می‌توان انتظار داشت این التهاب در بازار مسکن تا حد خوبی کنترل شود. از این رو من تصمیم گرفتم به عنوان پروژه داده‌کاوی روی این موضوع تمرکز کنم.

۲-۱ تعریف هر متغیر

برای درک بهتر داده‌ها اولین قدم تعریف دقیق هر متغیر است. در این بخش هر ۸ ویژگی موجود در ماتریس داده‌ها تشریح خواهم کرد. به علاوه متغیرها از نظر پاسخ یا پیشگو^۴ بودن نیز مشخص

۱. مساحت^۵: متغیر اول داده‌ها، مشخص‌کننده مساحت داده‌ها به وسیله یک عدد صحیح^۶ است. بدیهی است که مساحت، در حالت کلی یک عدد حقیقی^۷ است. (البته بسیاری از فیزیک‌دانان معتقدند که فضا و زمان اساساً گسسته‌اند. پس شاید بهتر باشد که بگویم مساحت عددی گویاست.) در اینجا مساحت به نزدیک‌ترین عدد صحیح گرد شده است.

۲. تعداد اتاق: متغیر دوم، عددی صحیح است که تعداد اتاق‌های هر ثبت را مشخص می‌کند.

۳. پارکینگ: یک متغیر دودویی^۸ نمایانگر تعلق داشتن و یا نداشتن پارکینگ به واحد مسکونی مورد نظر است. البته از نظر من بهتر بود تعداد پارکینگ‌های متعلق به یک واحد مسکونی با عددی صحیح مشخص می‌شد. ولیکن گردآورنده داده‌ها به صورتی که گزارش کرده‌ام اقدام به گردآوری نمود است.

۴. انباری: یک متغیر دودویی نمایانگر تعلق داشتن و یا نداشتن انباری به واحد مسکونی مورد نظر است. البته مانند ویژگی قبلی، از نظر من بهتر بود مساحت انباری متعلق به یک واحد مسکونی توسط یک عدد صحیح یا حقیقی مشخص می‌شد. اما مجدداً گردآورنده داده‌ها به صورتی که گزارش کرده‌ام، داده‌ها را گردآوری نموده‌اند.

۵. آسانسور: متغیری دودویی نمایانگر مجهز بودن و یا نبودن ساختمان شامل واحد مسکونی مورد نظر، به آسانسور است.

۶. آدرس: در هر ثبت آدرس واحد مسکونی توسط یک رشته^۹ از حروف مشخص شده است. بدیهی است که این رشته‌ها توسط یک متغیر رسته‌ای^{۱۰} قابل بیان خواهند بود.

۷. قیمت به تومان: متغیر هفتم، یک عدد صحیح نمایانگر قیمت در نظر گرفته شده برای ملک مسکونی است. البته بهتر بود برای چنین پروژه‌ای یک کارشناس این قیمت‌گذاری را انجام می‌داد. چرا که این قیمت‌ها نمایانگر نظر مالکان واحدهای مسکونی است.

۸. قیمت به دلار: این متغیر به جهت ثبات بیشتر این ارز نسبت به قیمت خانه‌ها، صرفاً برای مقایسه قیمت واحدهای

⁵Area⁶Integer⁷Real Number⁸Binary⁹String¹⁰Categorical¹Data²Record³Feature⁴Predictor

```

7 3rd Qu.: 120
8 Max. : 3600
9 NA's : 4
10 > summary(housePrice["Room"])
11 Room
12 Min. : 0.00
13 1st Qu.: 2.00
14 Median : 2.00
15 Mean : 2.08
16 3rd Qu.: 2.00
17 Max. : 5.00
18
19 > summary(housePrice["Parking"])
20 Parking
21 Min. : 0.0000
22 1st Qu.: 1.0000
23 Median : 1.0000
24 Mean : 0.8479
25 3rd Qu.: 1.0000
26 Max. : 1.0000
27
28 > summary(housePrice["Warehouse"])
29 Warehouse
30 Min. : 0.0000
31 1st Qu.: 1.0000
32 Median : 1.0000
33 Mean : 0.9146
34 3rd Qu.: 1.0000
35 Max. : 1.0000
36
37 > summary(housePrice["Elevator"])
38 Elevator
39 Min. : 0.0000
40 1st Qu.: 1.0000
41 Median : 1.0000
42 Mean : 0.7873
43 3rd Qu.: 1.0000
44 Max. : 1.0000
45
46 > summary(housePrice["Address"])
47 Address
48 Length: 3479
49 Class : character
50 Mode : character
51
52 > summary(housePrice["Price"])
53 Price
54 Min. : 3.600e+06
55 1st Qu.: 1.418e+09
56 Median : 2.900e+09
57 Mean : 5.359e+09
58 3rd Qu.: 6.000e+09
59 Max. : 9.240e+10

```

مسکونی در بازه‌های زمانی مختلف تهیه شده است. در واقع این ستون از ماتریس^{۱۱} داده‌ها اطلاعاتی به داده‌های ما اضافه نمی‌کنند چرا که نرخ دلار آمریکا، به مبلغ ثابت ۳۰ هزار تومان در زمان جمع‌آوری داده‌ها در نظر گرفته شده است. در نتیجه همبستگی^{۱۲} ۱ میان این ستون از ماتریس داده‌ها و ویژگی قبلی که قیمت واحد مسکونی به تومان بود، وجود دارد.

در این پروژه، ستون قیمت به دلار حذف خواهد شد همچنین ستون قیمت به تومان نقش متغیر پاسخ^{۱۳} را بازی خواهد کرد. (از این به بعد این متغیر را به اختصار قیمت می‌نامم) همچنین بقیه متغیرها پیشگو هستند.

۳-۱ نمونه یک ثبت

در این بخش می‌خواهم ثبت اول داده‌ها را به صورت عینی بررسی کنم. در جلد شماره ۱ می‌توانید سه ثبت اول داده‌ها را مشاهده کنید. در این بخش به توصیف اولین ثبت خواهم پرداخت.

Price(USD)	Price	Address	Elevator	Warehouse	Parking	Room	Area
\$ ۶۱,۶۶۶,۶۷	۱,۸۵۰,۰۰۰,۰۰۰ T	Shahran	TRUE	TRUE	TRUE	۱	۶۳
\$ ۶۱,۶۶۶,۶۷	۱,۸۵۰,۰۰۰,۰۰۰ T	Shahran	TRUE	TRUE	TRUE	۱	۶۰
\$ ۱۸,۲۳۲,۲۳	۵۵۰,۰۰۰,۰۰۰ T	Pardis	TRUE	TRUE	TRUE	۲	۷۹

جدول ۱: سه ثبت اول داده‌ها

توصیف ثبت اول ما به این صورت است که یک خانه ۱ خواب با مساحت ۶۳ مترمربع، دارای پارکینگ، انباری و آسانسور در منطقه شهران به مبلغ ۱,۸۵۰,۰۰۰,۰۰۰ تومان معادل ۶۱,۶۶۶,۶۷ دلار آمریکا (در آذر ۱۴۰۰ با حساب هر دلار معادل ۳۰ هزار تومان) در یک سامانه آنلاین آگهی املاک برای فروش عرضه شده است.

۴-۱ خلاصه وضعیت ویژگی‌ها

در هنگام شروع کار با داده‌ها خوب است که خلاصه^{۱۴} ای از شاخص‌های مرکزی مهم آنها بدانیم. من کدها به همراه نتایج حاصل شده از آنها را که با مترجم^{۱۵} زبان برنامه‌نویسی R اجرا شده است را در ادامه قرار خواهم داد.

```

1 > summary(housePrice["Area"])
2 Area
3 Min. : 30
4 1st Qu.: 69
5 Median : 90
6 Mean : 108

```

¹¹Matrix

¹²Correlation

¹³Response Variable

¹⁴Summary

¹⁵Compiler

```

2 housePrice.df <- read.csv("E:/DS/term2/
  faghihi/housePrice.csv")
3 frequency_of_add=housePrice.df %>% count_('
  Address')
4 tmp=frequency_of_add[order(-frequency_of_add
  $n),] #sort frequency_of_add by "freq"
  column
5 other_samples=(tmp)[,1][70:193] #We selected
  all addresses with a frequency of less
  than 9.
6
7 for (i in 1:length(other_samples)){
8   for(j in 1:dim(housePrice.df["Address"])
9     [1]){
10     if(housePrice.df[j,"Address"]==other
11       _samples[i]) {
12       housePrice.df[j,"Address"]="
13       Other"
14     }}
15 rm(frequency_of_add,tmp)
16 frequency_of_add=housePrice.df %>% count_('
17   Address')
18 tmp=frequency_of_add[order(-frequency_of_add
19   $n),]
20
21 for (i in 1:dim(tmp)[1]){
22   for(j in 1:dim(housePrice.df["Address"])
23     [1]){
24     if(housePrice.df[j,"Address"]==tmp[i
25       ,1]) {
26       housePrice.df[j,"Address"]=i
27     }}
28 }
29 as.factor(housePrice.df$Address)

```

با اجرای کد فوق در محیط Rstudio آنچه من به دنبال آن هستم محقق می‌شود. ضمناً همان‌طور که در گزارش‌های پیشین دیدیم ۲۳ ثبت من فاقد آدرس بودند. با اجرای کد تا اینجا این ۲۳ ثبت در یک رسته قرار می‌گیرند. برچسب این دسته در ویژگی آدرس، "۳۸" می‌باشد.

۲-۵-۱ مقادیر گم‌شده

در حال حاضر، فقط ۴ ثبت ما شامل متغیر تهی^{۱۹} است و مقدار گم‌شده^{۲۰} محسوب می‌شود. مساحت این ۴ ثبت‌ها به درستی درج نشده است و با توجه به تعداد کمی که در مقابل کل داده‌ها دارند من تصمیم گرفتم که از آنها چشم‌پوشی کنم و یا به عبارتی، آنها را از ماتریس داده‌ها حذف کنم.

```

60
61 > summary(housePrice["Price.USD."])
62   Price.USD.
63   Min.      : 120
64   1st Qu.: 47275
65   Median : 96667
66   Mean    : 178634
67   3rd Qu.: 200000
68   Max.    : 3080000

```

تابع describe در زبان برنامه‌نویسی R نیز می‌تواند اطلاعات خوبی را از داده‌ها در اختیار ما بگذارد. از این رو در مورد هر ویژگی، بخشی از گزارش خروجی این داده‌ها تحت تابع مذکور را در جدول ۲ قرار داده‌ام.

Variable	n	missing	distinct	Info	Sum	Mean	Gmd
Area	3475	4	239	1		108	89.60
Room	3479	0	6	804.0		08.2	7626.0
Parking	3479	0	2	387.0	2950	8479.0	2579.0
Warehouse	3479	0	2	234.0	3182	9146.0	1562.0
Elevator	3479	0	2	502.0	2739	7873.0	335.0
Address	3456	23	192				
Price	3479	0	934	1		359e+09.5	122e+09.6
Price.USD.	3479	0	932	1		178634	204055

جدول ۲: گزیده‌ای از خروجی تابع describe

۵-۱ پیش پردازش داده‌ها

پیش‌پردازش^{۱۶} داده به مرحله‌ای گفته می‌شود که در آن داده‌ها برای داده‌کاوی آماده می‌شود. لازم به ذکر است که این مراحل جزء مهم‌ترین گام‌ها در داده‌کاوی هستند. در این بخش بایستی ماتریس داده‌ها را نهایی کنیم تا در بدو شروع فصل بعد، تصویری سازی^{۱۷} روی ماتریس نهایی داده‌ها صورت بگیرد.

۱-۵-۱ ویژگی آدرس

یکی از چالش‌هایی که در این پروژه پیش روی من است این است که آدرس‌ها به صورت رشته نوشته شده‌اند. آدرس برخی از ثبت‌ها فراوانی^{۱۸} کمی دارند و ممکن است برای ما مشکل‌ساز باشد. از این رو آدرس‌هایی که کمتر از ۱۰ فراوانی داشته باشند را در دسته سایر قرار می‌دهم. سپس به ترتیب فراوانی (از فراوانی زیاد به کم) یک رشته عددی در چهارچوب یک متغیر رسته‌ای به هر آدرس در هر ثبت نسبت می‌دهیم. این‌گونه می‌توانم جزئیاتی که ممکن است در هنگام کار کردن با رشته‌های حروفی با آن مواجه شوم را از قلم ببندم.

```
1 library(plyr)
```

¹⁶Preprocessing

¹⁷Visualization

¹⁸Frequency

¹⁹Null

²⁰Missing Value

```

1 set.seed(1)
2 train.rows<-sample(rownames(X),dim(X)[1]*
  0.6)
3 valid.rows<-sample(setdiff(rownames(X),train
  .rows),dim(X)[1]*0.3)
4 test.rows<-setdiff(rownames(X),union(train.
  rows,valid.rows))
5 training_set<-X[train.rows,]
6 validation_set<-X[valid.rows,]
7 test_set<-X[test.rows,]

```

در انتها برای گزارش از ماتریس نهایی داده‌ها گزارش زیر را مشاهده بفرمایید.

```

1 > dim(training_set)
2 [1] 2080 7
3 > dim(validation_set)
4 [1] 1040 7
5 > dim(test_set)
6 [1] 354 7

```

۵-۵-۱ نرمالیده کردن متغیرها

در اینجا از روش نرمالیده^{۲۸} کردن داده‌ها با استفاده از Standard Scaler خواهیم پرداخت [۷] با اجرای کد پیش رو ویژگی‌های مساحت، تعداد اتاق و قیمت نرمالیده خواهند شد. ضمناً لازم به توضیح است که متغیرهای پارکینگ، انباری و آسانسور دودویی هستند و نیازی به نرمالیده شدن ندارند. همچنین متغیر آدرس یک متغیر رشته‌ای است و قابل نرمالیده شدن نیست.

این مرحله حتماً بایستی بعد از افراز داده‌ها صورت بگیرد که اطلاعات مجموعه تست از طریق نرمالیده شدن به داخل مجموعه آموزشی اصطلاحاً نشت نکند و به تعبیر علمی، نشت داده^{۲۹} رخ ندهد.

```

1 X_standard_scaler.train<-training_set
2 X_standard_scaler.valid<-validation_set
3 X_standard_scaler.test<-test_set
4 #- - - - -
5 for(i in c(1,2,7)){
6   X_standard_scaler.train[,i]=scale(
7     training_set[i])
8 }
9 #- - - - -
10 for(j in c(1,2,7)){
11   for (i in 1:dim(validation_set)[1]){
12     X_standard_scaler.valid[i,j]=(
13       validation_set[i,j]-mean(training_set[,j]
14       ))/sd(training_set[,j])
15   }
16 }

```

^{۲۸}Normalized

^{۲۹}Data Leakage

```

1 > X=na.omit(housePrice.df)
2 > dim(X)
3 [1] 3475 8

```

۳-۵-۱ داده‌های دورافتاده

داده دورافتاده^{۲۱} (داده پرت هم گفته می‌شود) داده‌هایی است که نسبت به سایر مشاهدات تفاوت قابل ملاحظه‌ای داشته باشد. [۶] در داده‌های من یک داده وجود دارد که با وجود مساحت زیاد، قیمت نسبتاً پایینی دارد. این داده، در سطر ۲۱۶۹ ماتریس داده‌های من وجود دارد. این داده شرح دهنده یک خانه ۳۶۰۰ متری ۲ خواب، بدون پارکینگ انباری و آسانسور در محله شهریار^{۲۲} به مبلغ ۹ میلیارد و ۷۲۰ میلیون تومان است. به عبارتی هر متر این خانه ۲ میلیون و ۷۰۰ هزار تومان قیمت‌گذاری شده است. همه این اطلاعات به ما نشان می‌دهد این ثبت مربوط به یک باغ است در شهریار که احتمالاً کاربری باغ دارد. اما پروژه ما مربوط به خانه‌های مسکونی در تهران است. لذا حذف این داده تحت عنوان داده دورافتاده خیلی غیر قابل انتظار نیست.

```

1 > X[2169,]
2      Area Room Parking Warehouse Elevator
3 2172 3600      2      0      0      0
4
5 Address      Price      Price.USD.
6      58      9.72e+09      324000

```

با اجرای دستور پیش رو داده دورافتاده مذکور در ادامه حذف خواهد شد.

```
1 > X=X[-c(2169),]
```

۴-۵-۱ افراز داده‌ها

در این بخش می‌خواهم داده‌ها را به سه بخش مجموعه آموزشی^{۲۳}، مجموعه اعتبار سنجی^{۲۴} و مجموعه آزمون^{۲۵} افراز^{۲۶} کنیم. میدانیم که بیشترین سهم مربوط به مجموعه آموزشی است لذا ۶۰ درصد داده‌ها را به صورت تصادفی^{۲۷} برای این مجموعه انتخاب می‌کنیم. ۳۰ درصد داده‌ها برای مجموعه اعتبار سنجی و ۱۰ درصد برای مجموعه آزمون به تصادف انتخاب می‌شود. اجرای کد پیشرو این فرایند را برای ما اجرا می‌کند.

^{۲۱}Outlier

^{۲۲}متغیر غیر عددی ۵۸ برای مشخص کردن خانه‌هایی که ویژگی آدرس آنها "شهریار" است انتخاب شده است.

^{۲۳}Training Set

^{۲۴}Validation Set

^{۲۵}Test Set

^{۲۶}Partition

^{۲۷}Random

```

12     }
13 }
14 #- - - - -
15 for(j in c(1,2,7)){
16     for (i in 1:dim(test_set)[1]){
17         X_standard_scaler.test[i,j]=(test_
18             set[i,j]-mean(training_set[,j]))/sd(
19                 training_set[,j])
17     }
18 }
19 }

```

۶-۵-۱ کدگذاری متغیرهای رسته‌ای

با توجه به اینکه تعداد رسته‌های متغیر رسته‌ای من در ویژگی رسته‌ای آدرس عدد ۷۰ است، استفاده از روش‌های مبتنی بر Dummy encoding و One hot encoding منجر به افزایش چشمگیر تعداد ویژگی‌ها می‌شوند. لذا بهتر است در این مورد از روش مرسوم Target encoding بر اساس میانه‌ی متغیر پاسخ استفاده کنیم. [۸] این مرحله نیز مانند بخش قبل حتماً بایستگی پس از بخش افراز داده‌ها انجام شود. دلیل این امر هم مانند آنچه در بخش قبل تشریح کرده‌ام است.

```

1 target_encode <- build_target_encoding( X_
2     standard_scaler.train, cols_to_encode =
3     "Address", target_col = "Price",
4     functions = c("median"))
5 X_standard_scaler.train=target_encode( X_
6     standard_scaler.train, target_encoding =
7     target_encode)[,c(2,3,4,5,6,8,7)]
8 X_standard_scaler.valid=target_encode( X_
9     standard_scaler.valid, target_encoding =
10    target_encode)[,c(2,3,4,5,6,8,7)]
11 X_standard_scaler.test=target_encode( X_
12    standard_scaler.test, target_encoding =
13    target_encode)[,c(2,3,4,5,6,8,7)]

```

با توجه به اینکه در فصل بعدی ممکن است مقادیر عددی منسوب به متغیرهای رسته‌ای من تغییر کند فعلاً در این بخش از وارد کردن این مقادیر در ماتریس اصلی داده‌ها خودداری می‌کنم. در ادامه بعد از نهایی شدن این مقادیر با داده‌های اصلی آنها را جایگزین خواهم کرد.

ضمناً دلیل استفاده من از شاخص مرکزی میانه، شکل ۱۳ است. این شکل به من نشان می‌دهد که میانگین پارامتر خوبی برای کدگذاری این متغیرهای رسته‌ای نخواهد بود. چراکه نمودارهای جعبه‌ای در اکثر موارد نامتقارن هستند. لذا میانگین معیار مناسبی نمی‌تواند باشد.

۲ تصویری سازی و اکتشاف داده‌ها

در فصل تصویری سازی احتیاج داریم که فرایند تصویری سازی روی کل داده‌های ما صورت بگیرد. از این رو در این فصل بیشتر با شیء‌های X و X_scaled ^{۳۰} سر و کار خواهیم داشت.

۱-۲ بافت‌نگار ویژگی‌ها

در این بخش ویژگی‌هایی که قابلیت دارند از آنها بافت‌نگار^{۳۱} رسم کنیم را مورد تحلیل و بررسی قرار می‌دهم. این فرایند به من کمک می‌کند تا در مورد توزیع هر یک از ویژگی‌ها اطلاعات خوبی را به دست آورم.

از آنجایی که بسیاری از ویژگی‌های ما مثل ویژگی مساحت نرمالیده شده‌اند، ممکن است تجسم خوبی از مقیاس آن‌ها نداشته باشیم. از این رو من در چنین مواردی مانند شکل ۱، دو بافت‌نگار ترسیم می‌کنم که به صورت عینی‌تر توزیع ویژگی مربوطه را به تصویر بکشم.

در خصوص بازه‌های بافت‌نگار ویژگی‌ها، حالت‌های مختلفی برای تعیین بازه‌های دسته‌بندی وجود دارد. هرچه طول بازه‌ها بیشتر باشد، نویز ناشی از نمونه‌گیری تصادفی را کم‌تر به نمایش می‌گذارد. از طرف دیگر هرچه طول بازه‌ها کمتر باشد، تخمین بهتری از توزیع می‌توان پیدا کرد.^[۹]

من در این پروژه تلاش کرده‌ام که طول بازه‌های دسته‌بندی را برای هر ویژگی به گونه‌ای تعیین شوند که بهترین ارتباط بصری را با مخاطب برقرار کنند. البته طول و تعداد بازه‌های انتخاب شده لزوماً بهترین نیستند و قضاوت این موضوع به عهده استاد محترم درس خواهد بود.

شکل ۱: همانطور که در گزارش اولیه فصل اول دیدیم، میانگین ویژگی مساحت ۱۰۸ متر و میانه آن ۹۰ متر بوده است. از آنجا که در بافت‌نگارها میانه^{۳۲} به شکل مشهودتری دیده می‌شود، باید بگویم که کاملاً واضح است که می‌توان خانه‌های حدود ۹۰ متر را خانه‌ای با مساحت عادی در نظر گرفت. نکته دوم اینکه در نمودار نرمالیده شده، "صفر" همان میانگین داده‌های نرمالیده است. کاملاً تأثیر داده‌های غیرعادی بزرگ (و نه داده‌های دورافتاده) را در فاصله دادن میانگین از میانه می‌توانیم مشاهده کنیم. ضمناً با اجرای یک کد کوتاه که در ادامه خواهیم آورد متوجه خواهیم شد که بیش از ۵.۶۴ درصد خانه‌هایی که روی آنها مطالعه می‌کنم مساحتی بین ۶۰ تا ۱۲۵ مترمربع دارند. این حقیقت نیز از روی نمودار بافت‌نگار قابل مشاهده و صحنه سنجی است.

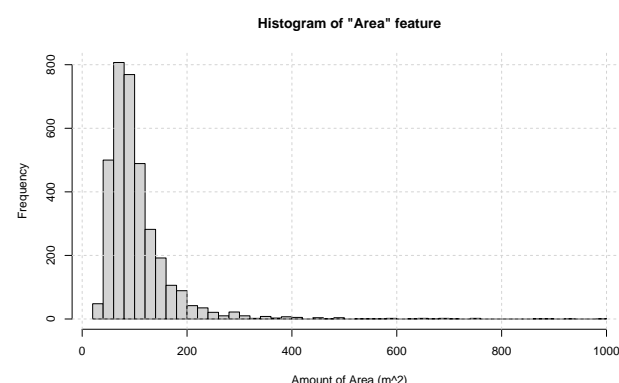
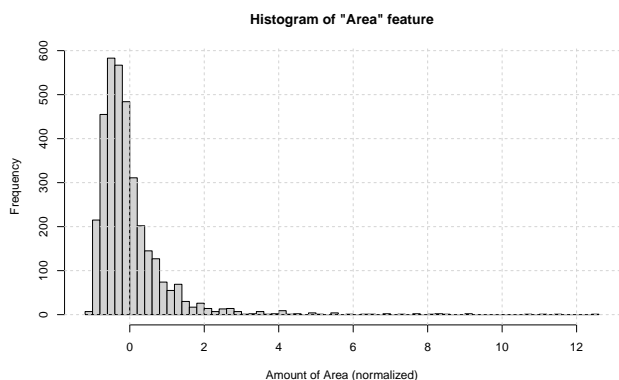
1 > j=0

^{۳۰}متغیرهای عددی غیر دودویی X در X_scaled نرمالیده شده‌اند. این ماتریس از داده‌ها صرفاً برای تصویری سازی ایجاد شده است.

^{۳۱}Histogram

^{۳۲}Mean

^{۳۳}Median



شکل ۱: بافت‌نگار ویژگی مساحت

```
2 > for (i in 1:3474){
3 +   if(X[i,1]<=125 & X[i,1]>=60){j=j+1}
4 + }
5 > print(100*j/3474)
6 [1] 64.50777
```

شکل ۲: همان‌طور که از شکل ۲ پیداست، به نظر می‌رسد در مورد تعداد اتاق‌ها با یک توزیع^{۳۴} تقریباً شبیه به توزیع نرمال^{۳۵} سروکار داریم. بیشترین فراوانی مربوط به خانه‌های دو خواب و کمترین فراوانی برای خانه‌های بدون اتاق خواب است. ضمناً خانه‌های سه خواب فراوانی بیشتری نسبت به خانه‌های یک خواب دارند.

در ادامه بافت‌نگار ویژگی‌های پارکینگ، انباری و آسانسور را خواهید دید که چون متغیرهای دودویی داشتند از نرمالیده کردن آنها صرف‌نظر کردم. از این رو نمودار بافت‌نگار آنها را به صورت خام ارائه خواهیم کرد.

شکل ۳: بیش از ۸۴ درصد خانه‌هایی که مورد مطالعه قرار دادم دارای دست‌کم یک پارکینگ هستند. این حقیقت را می‌توانستیم از خروجی تابع summary در فصل اول نیز دریابیم.

توجه داشته باشید آمارهایی که ارائه می‌کنم فقط مربوط به

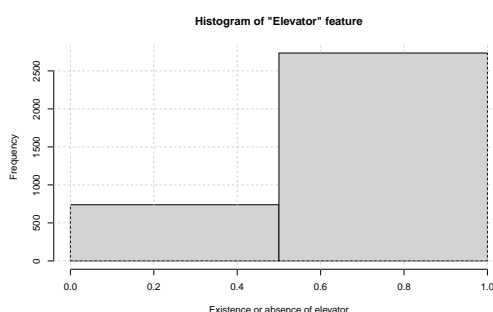
^{۳۴}Distribution

^{۳۵}Normal Distribution



شکل ۴: بافت‌نگار ویژگی انباری

یکی از انتقادات من به این مجموعه داده‌ها این است که ویژگی‌ها به صورت مدیرانه انتخاب نشده‌اند. مثلاً در مورد ویژگی پارکینگ شاید بهتر بود تعداد پارکینگ به جای وجود یا عدم وجود آن ثبت می‌شد. همچنین شاید سال ساخت واحد مسکونی می‌توانست یک ویژگی بسیار مؤثر باشد. به‌رحال من روی این داده‌ها کار می‌کنم و نمی‌توانم تغییری در آن ایجاد کنم ولیکن هر چند بسیار به داده‌ها انتقاد دارم امیدوارم با این ویژگی‌ها بتوانم نتیجه نسبتاً خوبی بگیرم.

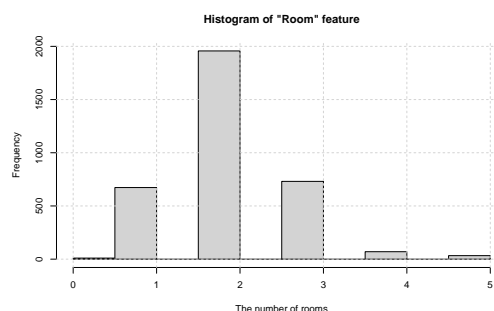
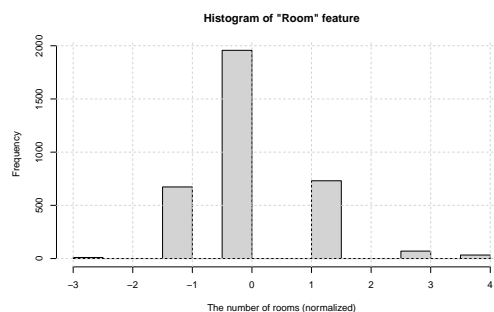


شکل ۵: بافت‌نگار ویژگی آسانسور

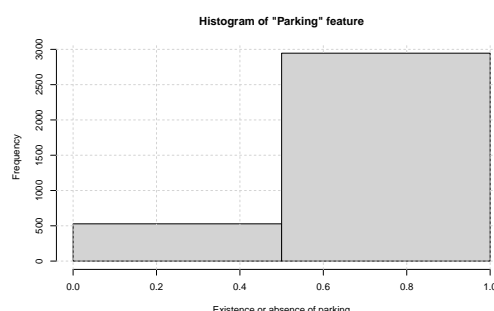
شکل ۵: در شکل ۵ آخرین بافت‌نگار دودویی را در مورد ویژگی آسانسور می‌بینیم که نسبت به دو ویژگی دودویی دیگر کمتر نامتعادل است. لازم به ذکر است که کمتر از ۸۰ درصد ثبت‌های ما دارای آسانسور هستند.

شکل ۶: این نمودار حاوی اطلاعات متغیر پاسخ ماست. چیزی که توجه من را جلب کرد این بود که شباهت نسبی‌ای بین این نمودار و شکل ۱ وجود دارد که می‌تواند بیانگر یک همبستگی میان متغیر پیشگوی مساحت و پاسخ باشد. در کل می‌توانیم به‌وضوح ببینیم هر چه قیمت خانه‌ها بیشتر باشند، فراوانی آنها نیز کمتر است.

به‌علاوه در شکل ۶ دو نمودار مشاهده می‌کنیم. نمودار اول، بافت‌نگار داده‌های نرمالیده است و نمودار دوم بافت‌نگار قیمت



شکل ۲: بافت‌نگار ویژگی تعداد اتاق

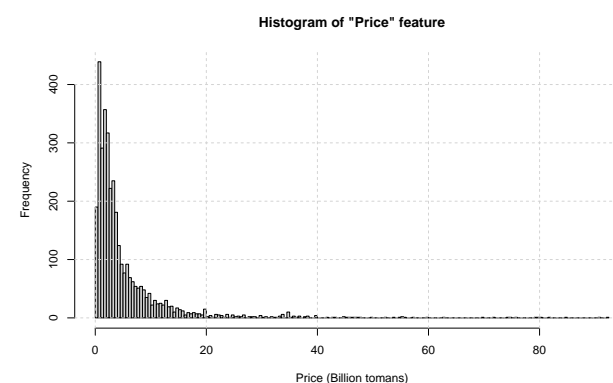
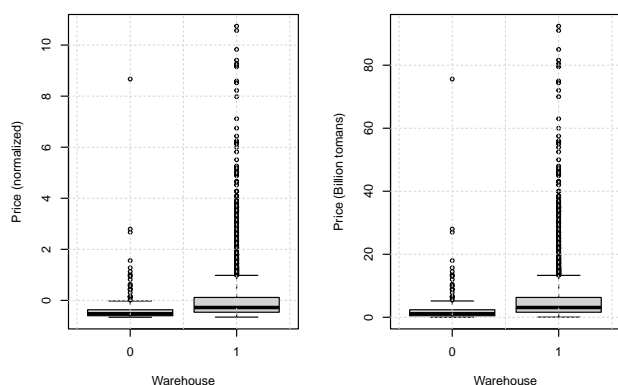
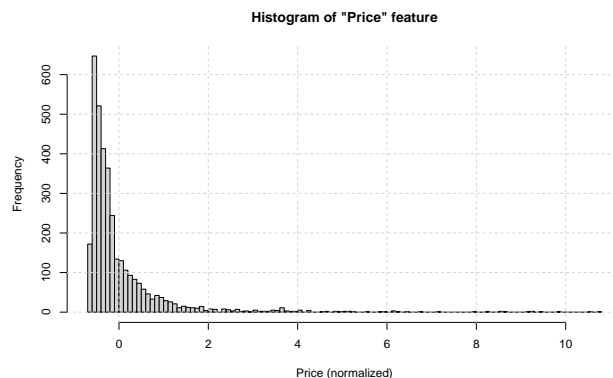
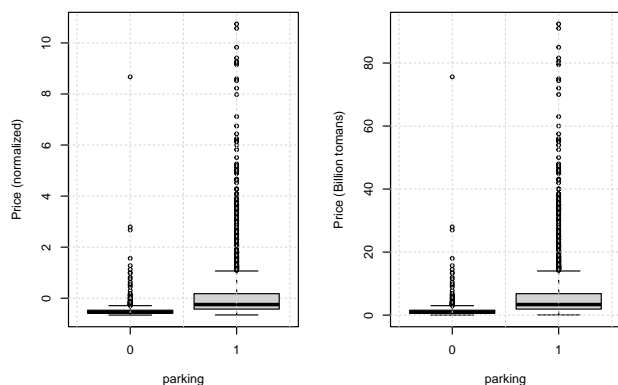


شکل ۳: بافت‌نگار ویژگی پارکینگ

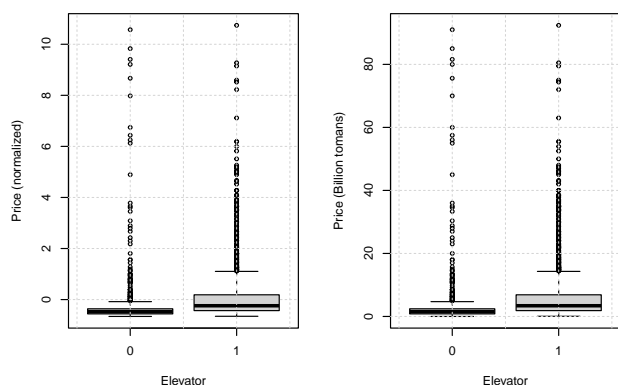
داده‌های این پروژه است. ممکن است برداشت‌های ما از داده‌های این پروژه با آنچه در سطح استان تهران می‌بینیم متفاوت باشد. به‌عنوان مثال بسیاری از خانه‌هایی که قصد فروش ندارند، به‌هیچ‌وجه در این مجموعه داده‌ها نمی‌توانستند قرار بگیرند. پس نمی‌توان ادعا کرد که مثلاً بیش از ۸۴ درصد خانه‌های استان تهران دارای پارکینگ هستند.

شکل ۴: در مورد ویژگی انباری نیز باید بگویم مشابه ویژگی پارکینگ، یک نامتعادل^{۳۶} بودن را در داده‌ها مشاهده می‌کنم. ضمناً بیش از ۹۱ درصد خانه‌هایی که ما آنها را در این پروژه بررسی می‌کنیم دارای انباری هستند. (در مورد مساحت انباری هیچ اطلاعاتی در دست نیست.)

³⁶Imbalance



شکل ۶: بافت‌نگار متغیر پاسخ



شکل ۷: نمودار جعبه‌ای پهلوه‌پهلوه متغیرهای دودویی - قیمت

در مورد انباری و آسانسور نیز صدق می‌کند. چیزی که توجه من را جلب کرد این بود که برخلاف نمودارهای جعبه‌ای مربوط به انباری و پارکینگ، داده‌هایی که نرم‌افزار آنها را دور افتاده تشخیص داده است در نموداری که مربوط به ویژگی آسانسور است به شکل متوازن‌تری بین خانه‌های با آسانسور و بی آسانسور توزیع شده است. حدس من این است که این داده‌های دور افتاده می‌توانند خانه‌های ویلایی گران قیمتی باشند که نیاز به آسانسور ندارند. اما فارغ از این داده‌هایی که تعداد کمی دارند، عمده داده‌ها به ما می‌گویند داشتن هر سه این امکانات می‌تواند تأثیر افزایشی روی

خانه‌ها برحسب میلیارد تومان است.

۲-۲ نمودار جعبه‌ای

در این بخش تلاش خواهیم کرد از نمودار جعبه‌ای^{۳۷} و نمودار جعبه‌ای پهلوه‌پهلوه^{۳۸} برای نشان دادن توزیع ویژگی‌ها نسبت به هم خصوصاً نسبت به متغیر پاسخ، استفاده کنیم. شکل ۷: در ابتدا در شکل ۷ برای متغیرهایی که بیانگر داشتن یا نداشتن امکاناتی در منزل مسکونی است از نمودارهای جعبه‌ای پهلوه‌پهلوه استفاده می‌کنم. این کار به من کمک خواهد کرد تا بینم تا چه حد به طور کلی وجود یا عدم وجود این امکانات از قبیل پارکینگ، انباری و آسانسور روی قیمت یک خانه مسکونی در شهر تهران (بر اساس داده‌های در اختیار ما) می‌تواند اثربخش باشد.

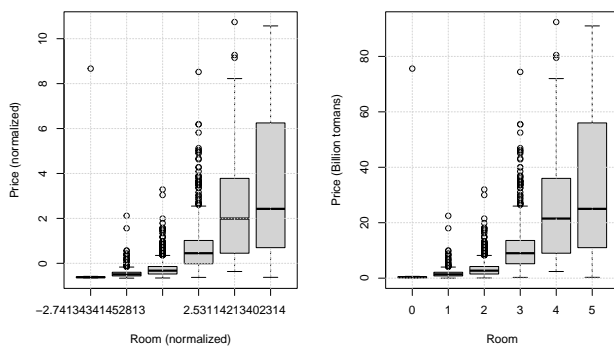
از آنجایی که داده‌ها نرمالیده شده‌اند، ممکن است نمودار داده‌های نرمالیده نتواند مقدار متغیر قیمت را به مخاطب انتقال دهد، لذا تصمیم گرفتم یک ستون به متغیر قیمت برحسب میلیارد تومان اختصاص دهم.

همان‌طور که از شکل ۷ مشخص است قیمت خانه‌هایی که دارای پارکینگ هستند به شکل مشهودی بالاتر است. این موضوع

^{۳۷}Boxplot

^{۳۸}Side-By-Side Boxplot

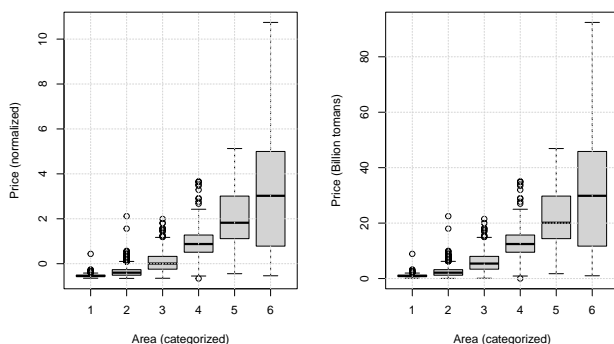
قیمت خانه‌ها داشته باشند.



شکل ۹: نمودار جعبه‌ای پهلوه پهلوه متغیر تعداد اتاق - قیمت

تعداد اتاق‌هایی که در خود جای داده‌اند، خیلی گران هستند، غالباً دو و سه خواب هستند.

یک خانه حدوداً ۸۰ میلیارد تومانی داریم که هیچ اتاق خوابی ندارد. این ثبت را بررسی کردم و متوجه شدم این ثبت مربوط به یک ملک ۶۳۰ متری با قیمت ۷۵ میلیارد و ۶۰۰ میلیون تومان، واقع در منطقه تجریش (در دسته بندی ما در رسته ۱ یا همان سایر قرار دارد) این طور که مشخص است، احتمالاً این ثبت مربوط به یک قطعه زمین بایر در منطقه تجریش است. در مورد حذف آن در صورت لزوم در ادامه تصمیم خواهیم گرفت.



شکل ۱۰: نمودار جعبه‌ای پهلوه پهلوه متغیر رسته‌ای شده مساحت - قیمت

شکل ۱۰: از آنجایی که میان شکل‌های ۸ و ۹ جای خالی نمودار جعبه‌ای پهلوه پهلوه مساحت - قیمت حس می‌شد و از طرفی مساحت متغیری پیوسته بود برای نمایش بهتر، در شکل ۱۰ مساحت خانه‌ها را به ۶ دسته مطابق جدول ۳ طبقه‌بندی کردم. همان طور که پیش‌تر حدس زده بودم رابطه مستقیمی میان متغیر پاسخ و متغیر پیشگوی مساحت وجود دارد. یعنی به طور کلی خانه‌های بزرگ‌تر گران‌تر هستند. البته این حدس خیلی واضحی است اما برداشت ما از تحلیل علمی داده‌ها سندیت دارد، لذا تا

شکل ۸: نمودار جعبه‌ای پهلوه پهلوه متغیر تعداد اتاق - مساحت

شکل ۸: به صورت کلی به من این نتیجه را می‌دهد که تعداد اتاق، با مساحت خانه‌ها رابطه مستقیم دارد. این حقیقت نشان دهنده همبستگی مثبت میان متغیر پیشگو تعداد اتاق و متغیر پیشگو مساحت است. (مقدار دقیق همبستگی در اجرای کد پیش رو قابل مشاهده است) خالی از لطف نیست که اشاره کنم، از آنجایی که حدس می‌زنم به احتمال بسیار قوی متغیر مساحت، با متغیر پاسخ، همبستگی قابل توجهی دارد؛ هر چه همبستگی میان متغیر مساحت و متغیر تعداد اتاق کمتر از ۱ باشد و در عین حال متغیر تعداد اتاق با متغیر پاسخ همبستگی نزدیک‌تری به ۱ داشته باشند، این دو ویژگی در پیشگویی کمک بیشتری به من خواهند کرد.

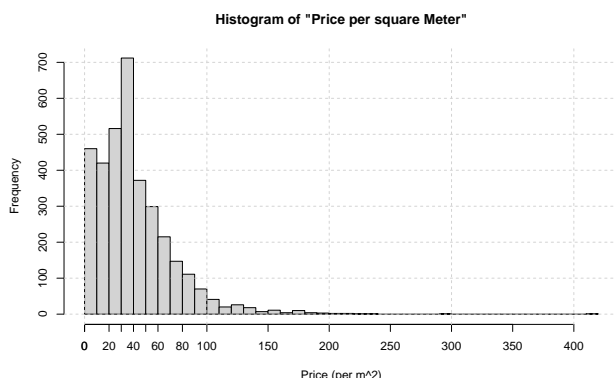
```
1 >cor(X[,1] , X[,2])
2 Room
3 Area 0.6570648
```

همان‌طور که در شکل ۲، یا بافت نگار ویژگی تعداد اتاق قابل ملاحظه است، بیشترین فراوانی مربوط به خانه‌های دو خواب است. از این رو طبیعی است که داده‌های دور افتاده بیشتری را در شکل ۸ برای خانه‌های دو خواب و پس از آن برای خانه‌های سه خواب شاهد باشیم. (تصور من این است که بسیاری از این داده‌های دور افتاده مربوط به خانه‌های ویلایی است)

شکل ۹: اولین نکته‌ای که این شکل به من نشان می‌دهد این است که سه شاخص آماری، شامل میانه، چارک^{۳۹} سوم و صدک^{۴۰} صدم برای متغیر پاسخ یا همان قیمت خانه با افزایش تعداد اتاق خواب‌ها افزایش می‌یابد. نکته دیگری من متوجه آن هستم این است که داده‌های دور افتاده در خانه‌های دو و سه خواب بیشترین سهم را به خود اختصاص داده‌اند. یعنی خانه‌هایی که به نسبت

^{۳۹}Quartile

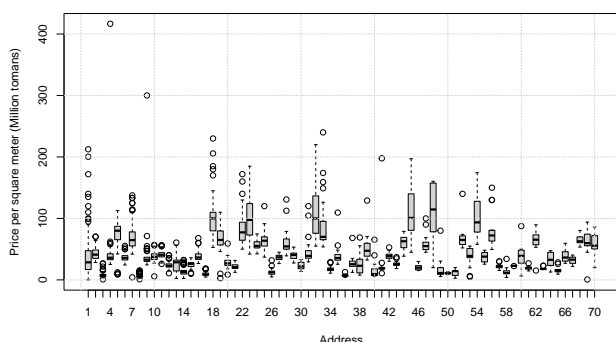
^{۴۰}Percentile



شکل ۱۲: بافت‌نگار توزیع قیمت هر متر خانه در شهر تهران

کردم. اگر این عدد روی ۶۰ تنظیم می‌شد متوجه می‌شدیم که تقریباً تعداد خانه‌های ۱۰ تا ۲۰ میلیون تومان (به ازای واحد مساحت) با خانه‌های ۳۰ تا ۳۵ میلیون تومانی (به ازای واحد مساحت) برابر است. اما تعداد خانه‌های ۰ تا ۱۰ میلیون تومانی بسیار کمتر از خانه‌های ۳۵ تا ۴۰ میلیون تومانی است. لذا قضاوت اینکه قیمت یک متر خانه در تهران چقدر است کار سخت‌تری از اعلام یک عدد یا یک بازه قیمتی است.

آخرین نموداری که در بخش نمودارهای پهلوه‌پهلوه می‌خواهم بررسی کنم متشکل از ۷۰ نمودار جعبه‌ای است. این نمودار به قیمت هر متر خانه در مناطق مختلف تهران می‌پردازد.



شکل ۱۳: نمودار جعبه‌ای پهلوه‌پهلوه متغیر رسته‌ای آدرس- قیمت هر متر مربع خانه

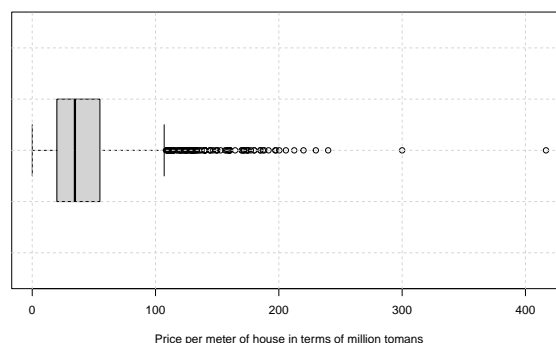
شکل ۱۳: همان‌طور که قبلاً هم اشاره کرده‌ام متغیر پیشگو رسته‌ای آدرس، با توجه به فراوانی آدرس‌ها نام‌گذاری شدند. رسته شماره ۱ نشان دهنده رسته سایر و باقی رسته‌ها بر اساس فراوانی هر یک متناظر با عددی از ۲ تا ۷۰ هستند. این نمودار به من نشان می‌دهد که اولاً بازه قیمتی در مناطقی که گران‌قیمت‌تر هستند بزرگ‌تر است. یعنی در مناطق گران‌قیمت‌تر، خانه‌ها قیمت دقیق‌تری ندارند. لذا این نکته کلیدی را به من نشان می‌دهد که ضرورت

Area ≤ ۵۰	۵۰ < Area ≤ ۱۰۰	۱۰۰ < Area ≤ ۱۵۰	۱۵۰ < Area ≤ ۲۰۰	۲۰۰ < Area ≤ ۲۵۰	۲۵۰ < Area
۱	۲	۳	۴	۵	۶

جدول ۳: طبقه‌بندی ویژگی مساحت برای رسم شکل ۱۰

زمانی که داده‌ها به ما اطلاعاتی را ندهند هیچ سوگیری نخواهم داشت.

Boxplot of one meter of houses price in Tehran



شکل ۱۱: نمودار جعبه‌ای قیمت هر متر مربع خانه در شهر تهران

شکل ۱۱: شکل ۱۰ این ایده را به من داد که پراکندگی^{۴۱} قیمت هر متر خانه را در تهران بررسی کنم. همچنین برای درک بهتر از جدول ۴ برای تفسیر بهتر شکل ۱۱ نیز بهره بردم. برداشت من این است که قیمت هر متر خانه‌های مورد بررسی ما به طور متوسط قیمتی بین ۳۰ تا ۴۰ میلیون تومان دارند. اما خانه‌هایی که قیمت آنها بسیار بالاتر از این اعداد هستند نیز با فراوانی کمتر نیز وجود دارند. از این رو می‌توان پیش‌بینی کرد که احتمالاً بافت‌نگار متناظر با شکل ۱۱، مانند شکل‌های ۱ و ۶ دارای چولگی مثبت^{۴۲} خواهد بود.

Min.	1st Qu.	Median	Mean	3rd Qu.	Max.
۰٫۲۲۵	۲۰٫۰۰	۳۴٫۶۱۵۴	۴۱٫۱۶۰۷	۵۴٫۸۸۶۱	۴۱۶٫۶۶۶۷

جدول ۴: جدول شاخص‌های آماری مرتبط با شکل ۱۱ بر حسب میلیون تومان

تحلیل ما از جدول ۴ و شکل ۱۱ به ما نشان می‌داد که به طور متوسط خانه‌ها بین ۳۰ تا ۴۰ میلیون تومان (به ازای هر متر) قیمت دارند. اما بررسی بافت‌نگار توزیع قیمت هر متر خانه در تهران اطلاعات دیگری نیز به من می‌دهد. در نتیجه بر آن شدم تا برای مقایسه استثنائاً در این قسمت مربوط به نمودارهای جعبه‌ای، یک بافت‌نگار را نیز بررسی کنم
من در این بافت‌نگار تعداد میله‌ها را حدود ۴۰ عدد تنظیم

^{۴۱}Dispersion

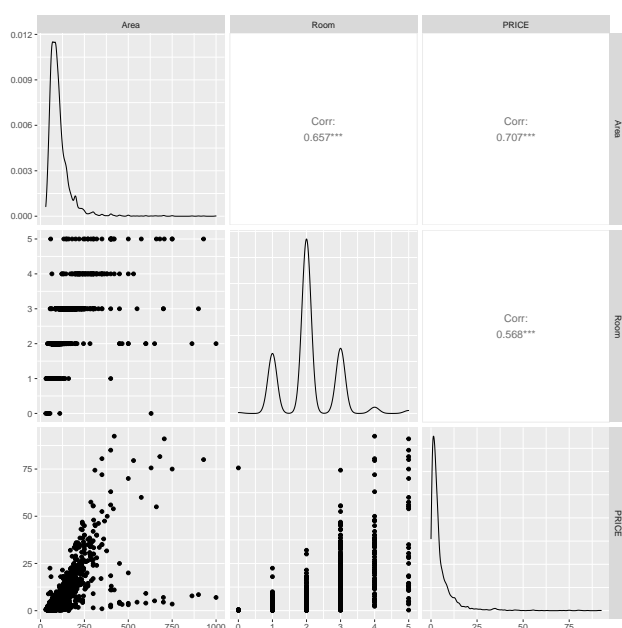
^{۴۲}Positive Skewness

ویژگی‌ای که در این تصویر قابل مشاهده نیست (چون یک مغیر رسته‌ای است) ویژگی آدرس هست. همان‌طور که در شکل ۱۳ قابل مشاهده است یک ویژگی بسیار اثرگذار روی قیمت خانه‌ها آدرس است که در ادامه گزارش به اهمیت آنها اشاره خواهم کرد.

	Area	Room	Parking	Warehouse	Elevator	Price
Area	۱	۰/۶۵۷۱	۰/۱۹۴۷	۰/۷۶۲	۰/۴۴۱۸	۰/۷۰۶۷
Room	۰/۶۵۷۱	۱	۰/۲۷۴۹	۰/۱۲۹۹	۰/۸۱۱۷	۰/۵۶۷۶
Parking	۰/۱۹۴۷	۰/۲۷۴۹	۱	۰/۴۳۰۸	۰/۴۳۰۴	۰/۱۹۰۵
Warehouse	۰/۷۶۲	۰/۱۲۹۹	۰/۴۳۰۸	۱	۰/۲۰۱۶	۰/۱۱۰۰
Elevator	۰/۴۴۲	۰/۸۱۱۷	۰/۴۳۰۴	۰/۲۰۱۶	۱	۰/۱۱۱۷
Price	۰/۷۰۶۷	۰/۵۶۷۶	۰/۱۹۰۵	۰/۱۱۰۰	۰/۱۱۱۷	۱

جدول ۵: جدول همبستگی داده‌ها

برای پیدا کردن دید بهتر روی سه ستون مهم در ماتریس داده‌ها شامل مساحت، تعداد اتاق و قیمت (برحسب میلیارد تومان) شکل پیش رو را که حاوی ماتریس نمودار پراکنش است را نیز ضمیمه نمودارهای قبل کرده‌ام.



شکل ۱۵: ماتریس نمودار پراکنش برای قیمت و دو پیشگوی عددی مساحت و تعداد اتاق

۴-۲ تصویری سازی چند بُعدی

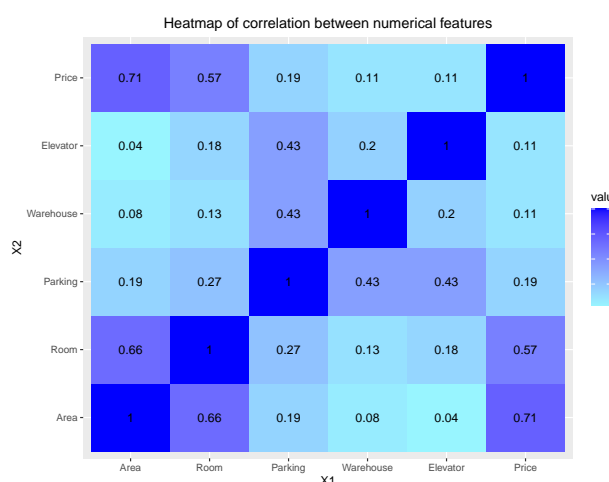
در این بخش تلاش می‌کنم با افزودن متغیرهای رسته‌ای از طریق کدگذاری رنگی و ... جنبه‌های بیشتر از دو بُعد را با استفاده از نمودار پراکنش^{۴۵} داده‌ها به تصویر بکشم.

شکل ۱۶: این تصویر در واقع در شکل ۱۵ (پایین سمت چپ تصویر) یک‌بار رسم شده است اما در این تصویر شاهد

وجود یک سیستم پیش‌گویانه قیمت در این مناطق بیشتر از مناطق ارزان قیمت‌تر است. چرا که آمار به من یک آشفتگی در قیمت گذاری را نشان می‌دهد. ثانیاً بیشترین و کمترین تعداد خانه‌های برای فروش گذاشته شده، در مناطقی هستند که قیمت هر متر نزدیک به میانگین و میانه است. مناطق خیلی گران‌تر بر خلاف ادعای بنگاه‌های معاملاتی دارای بازار داغ‌تر خرید و فروش نیستند. ضمناً با توجه به بررسی‌های من در خود داده‌ها منطقه ۵ داغ‌ترین بازار خرید و فروش املاک را دارد. پیش‌تر یک وب‌سایت خرید و فروش مسکن نیز این ادعا را مطرح کرده بود و دلایل خود را برای این امر ذکر کرده بود. [۱۰]

۳-۲ نمودار حرارتی

نمودار حرارتی^{۴۳} یک نوع تصویرسازی داده‌ها^{۴۴} است که در آن همبستگی بین ویژگی‌ها با یک رنگ نمایش داده می‌شود. ویژگی‌هایی که دارای همبستگی بالایی با یکدیگر هستند با رنگ‌های تیره‌تری نسبت به ویژگی‌هایی که دارای همبستگی کمتری نسبت به هم هستند به نمایش گذاشته می‌شوند. به طور معمول برای نمایش بصری اعداد یک بازه رنگی از روشن تا تیره (از کم به زیاد) در نظر گرفته می‌شود. البته این نمودارها همواره توسط یک ماتریس عددی قابل بیان هستند. اما نمودار حرارتی از نظر بصری قابلیت درک بیشتری دارد. [۱۱]



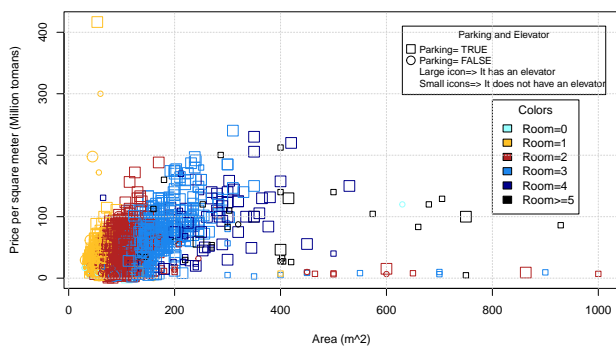
شکل ۱۴: نمودار حرارتی جدول همبستگی ویژگی‌های عددی

شکل ۱۴: این نمودار به ما نشان می‌دهد انباری و آسانسور به طور کلی تأثیر کمی روی قیمت خانه‌ها داشته‌اند. اما داشتن آن مزیت محسوب می‌شود. (این شهود را در بازار نیز شاهد هستیم.) اهمیت پارکینگ حدود دو برابر انباری و آسانسور گزارش شده است اما همچنان مهم‌ترین پارامتر نیستند. دو پارامتر مساحت و تعداد اتاق بیشترین همبستگی را با قیمت خانه دارند و البته

⁴³Heatmap

⁴⁴Data Visualization

⁴⁵Scatter Diagram



شکل ۱۷: نمودار چند بُعدی کدگذاری شده مساحت-قیمت هر متر

کشیده شده بود قیمت هر متر خانه را مورد بررسی قرار داده‌ام. قیمت هر متر خانه به طور غیرمستقیم گویای متغیر رسته‌ای آدرس است. (این ادعا در شکل ۱۷ صحنه سنجی گردیده است.) در این نمودار می‌توان دید خانه‌های دو خواب عموماً مساحت بیشتری از خانه‌های یک خواب دارند اما به ازای هر متر تفاوت قیمتی چندانی ندارند.

همچنین می‌توان دید محله‌های ارزان قیمت‌تر علاقه کمتری به خانه‌های ۳ خواب نشان می‌دهند. ^{۴۶} بر خلاف تصور عام، خانه‌های ۴ خواب هم در محلات گران قیمت و به همان نسبت در محلات ارزان قیمت دیده می‌شوند.

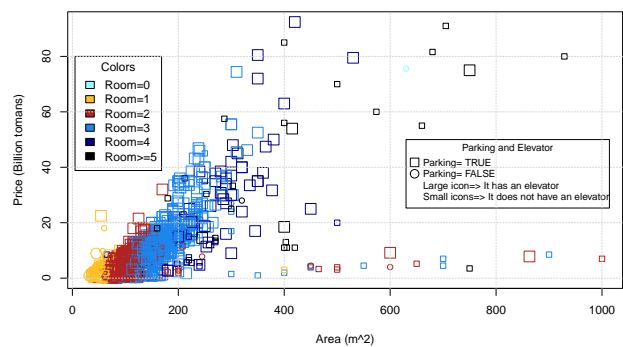
داده‌های دورافتاده اینجا در مورد خانه‌های یک خواب بسیار شگفت‌انگیزند. (داده‌های دور افتاده دیگر تا حدی قابل توجیه هستند.) اما این ۴ ثبت از خانه‌های یک خواب، نه تنها مساحت زیادی ندارند (در نتیجه یک خواب هستند) بلکه قیمت یک متر خانه در آنها به شکل قابل توجهی گران است. من حتی در مورد یکی از ثبت‌ها که نه آسانسور دارد و نه پارکینگ دارد و نه با توجه به مساحت می‌تواند یک خانه ویلایی خیلی بزرگی باشد، علت گران بودن را در موقعیت تجاری می‌بینم.

نکته دیگر اینکه خانه‌های ۵ خواب و یا بیشتر، معمولاً آسانسور ندارند. احتمالاً همان‌طور که چندین بار اشاره کرده‌ام ممکن است این ثبت‌ها مربوط به خانه‌های ویلایی باشند.

شکل ۱۸: از آنجایی که متوجه شدم شکل ۱۷ ارتباط معناداری بایستی با محله داشته باشد، سعی کرد رسته کدگذاری شده هر محله را به جای هر ثبت بنویسم. در این نمودار ویژگی پارکینگ دیده نمی‌شود. اما متغیر رسته‌ای آدرس جای اشکال هندسی شکل ۱۷ و ۱۶ را گرفت است.

من متوجه شدم که رسته ۱ که رسته "سایر" بود به شکل خیلی پراکنده‌ای در محلات مختلف شامل محلات خیلی گران قیمت و

^{۴۶} به شکل ۱۹ رجوع کنید.



شکل ۱۶: نمودار چند بُعدی کدگذاری شده مساحت-قیمت

کدگذاری‌هایی هستیم که به لطف آنها می‌توان، هر چند به صورت شهودی متغیر پاسخ به همراه ۴ ویژگی دیگر را به صورت مجزا مورد مشاهده قرار داد.

ستون‌های عمودی و افقی را به ترتیب به متغیر پاسخ و ویژگی پیشگوی مساحت (که بالاترین نرخ همبستگی را با متغیر پاسخ داشت) اختصاص داده‌ام. رنگ هر داده نمایانگر تعداد اتاق خواب‌های خانه موردنظر است. همچنین ویژگی پارکینگ توسط نوع شکل هندسی متناظر با داده مشخص شده است. در انتها با اندازه اشکال هندسی رسم شده برای هر داده داشتن یا نداشتن آسانسور را در هر ثبت مشخص کرده‌ام.

اولین چیزی که من می‌بینم رنگ مربوط به خانه‌های یک تا سه خواب و بعد از آن خانه‌های ۴ خواب است. در مورد خانه‌های یک تا سه خواب، می‌توان گفت بیشینه‌ی قیمت بیشتری را (با صرف نظر از داده‌های دور افتاده) شاهد هستیم. به عبارت ساده‌تر عموماً هر چه تعداد خواب خانه بیشتر باشد، آن خانه می‌تواند گران‌تر باشد.

نکته دیگر این است که اولاً خانه‌های دارای پارکینگ بسیار بیشتر از خانه‌های بدون پارکینگ هستند و ثانیاً معمولاً خانه‌های بدون پارکینگ ارزان‌تر هستند. یادآوری می‌کنم که بیش از ۸۴ درصد خانه‌های مورد مطالعه من دارای آسانسور هستند.

همچنین با نگاهی اجمالی، واضح است که هم تعداد خانه‌های دارای آسانسور از خانه‌های بدون آسانسور بیشتر است و هم کمتر شاهد خانه‌های بدون آسانسور نسبتاً گران قیمت هستیم. اما باید اضافه کنم که خانه‌های بدون آسانسور به وضوح در زمره داده‌های دور از محل تجمع کلی داده‌ها دیده می‌شوند. احتمالاً به این دلیل است که خانه‌های ویلایی که عموماً خیلی گرانند، فاقد آسانسور هستند.

شکل ۱۷: توضیحات مشابه در این نمودار مانند نمودار ۱۶ نیز صادق است. لذا فقط به تفاوت‌ها خواهم پرداخت. در نمودار ۱۷ من به جای قیمت کلی خانه که برحسب میلیارد تومان به تصویر

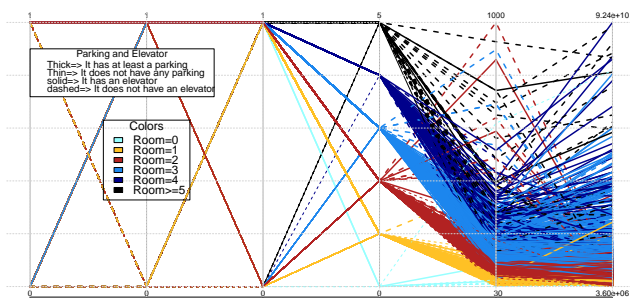
اینجا لازم است به دو مورد مهم اشاره کنم. این فرض در صورتی صحیح است که اولاً این نمونه‌ها یک نمونه خوب از جامعه باشند. یعنی به صورت تصادفی مناطق ارزان قیمت خالی از خانه‌های سه خواب تهیه نشده باشند. ثانیاً رابطه مستقیمی میان قیمت هر متر مناطق ارزان قیمت و گران قیمت وجود داشته باشد. با فرض اینکه این نمونه یک نمونه خوب باشد پیش می‌رویم.

در این شکل من نمودار قیمت هر متر-قیمت کل را با دو بزرگنمایی مختلف رسم کرده‌ام. نکته جالب اینجاست که نه تنها این دو پارامتر به هم وابسته هستند بلکه یک رابطه غیرخطی از جنس چندجمله‌ای درجه سه بین آنها برقرار است (رجوع کنید به بالانویس شکل ۱۹)

من در این نمودار چهار منطقه‌ای که میانه بالاتری در متغیر پاسخ، یا همان قیمت داشتند شامل مناطق الهیه، زعفرانیه، نیاوران و فرمانیه با رنگ آبی مشخص کرده‌ام. همچنین چهار منطقه ارزان قیمت شامل پرند، پردیس، پاکدشت و چیتگر شمالی نیز با رنگ طلایی مشخص شده‌اند. بقیه مناطق با رنگ مشکی در نمودار قابل مشاهده‌اند.

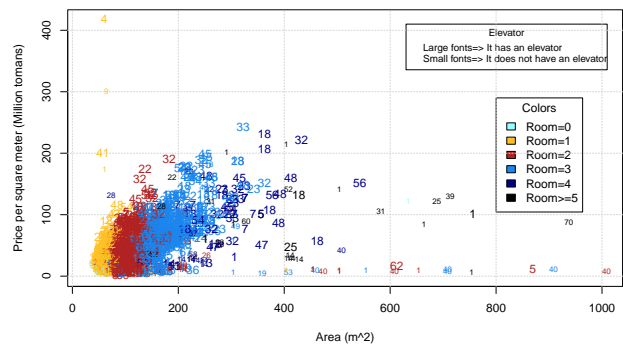
واضح است که داده‌های مناطق ارزان قیمت ما هم از نظر متغیر پاسخ یعنی قیمت کل، و هم از نظر قیمت هر متر به شکل قابل توجهی دارای مقادیر کمتری از خانه‌های گران قیمت هستند. لذا می‌توان به محکمی ادعا کرد که متغیر آدرس که نمایانگر مناطق مختلف است، ارتباط قابل توجهی با قیمت هر متر خانه دارد. لذا ادعای من در شکل ۱۷ ادعای غلطی نبوده است.

۵-۲ نمودار با محورهای موازی



شکل ۲۰: نمودار چند بُعدی کدگذاری شده با محورهای موازی

شکل ۲۰: این نمودار به شکل دیگری حاوی اطلاعاتی است که نمودار ۱۶ در اختیار ما می‌گذاشت. من سعی کردم در این نمودار محورهای موازی، همان شکل ۱۶ را برای به تصویر کشیدن بهتر ابعاد بالاتر استفاده کنم.

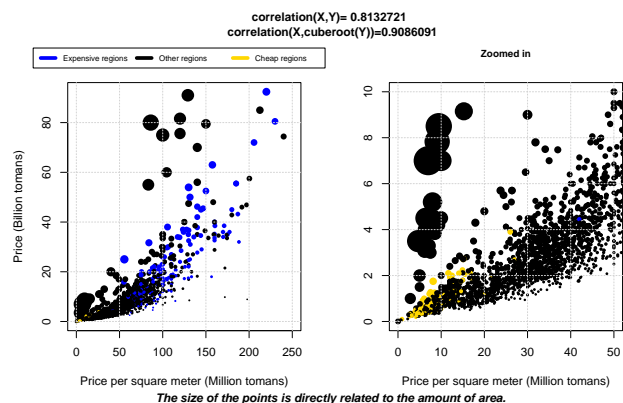


شکل ۱۸: نمودار چند بُعدی کدگذاری شده مساحت-قیمت هر متر با تمرکز روی ویژگی آدرس

خیلی ارزان قیمت توزیع شده است. که البته با توجه به نوع انتخاب دسته دور از انتظار هم نیست.

نکته دیگر این است که همان‌طور که انتظار می‌رفت بسیاری از رسته‌های مشابه در این نمودار خیلی نزدیک به هم هستند. متأسفانه برای ساختن محور عمودی در این نمودار از ضربی از متغیر پاسخ استفاده شده. اگر چنین نبود و آدرس، متغیر پاسخ بود، الگوریتم KNN می‌توانست یک گزینه خوب برای پیشگویی آدرس باشد.

نکته دیگری که به صورت غیر رسمی عنوان می‌کنم این است که هر چه مساحت خانه‌ها بیشتر باشند، به طور کلی قیمت به ازای هر متر در همان محله کاهش می‌یابد. از آنجایی که مصداق عددی برای آن نیاوردم به صورت رسمی این نکته را تأیید نمی‌کنم. اما به شکل شهودی مصداق‌های زیادی را در شکل می‌بینم. (البته که خلاف این مصداق‌ها هم وجود دارند.)



شکل ۱۹: نمودار چند بُعدی کدگذاری شده قیمت هر متر-قیمت کل

شکل ۱۹: پیش‌تر نتیجه‌گیری کرده بودم که مردم در مناطق ارزان قیمت علاقه کمتری به خانه‌های سه خواب نشان می‌دهند. در

۶-۲ کاهش بُعد

در یادگیری ماشین^{۴۷} و آمار کاهش بُعد^{۴۸} یا کاهش ابعاد روند کاهش تعداد متغیرهای تصادفی راهنماییده^{۴۹} [۱۲] از طریق به دست آوردن یک مجموعه از متغیرهای اصلی می‌باشد. کاهش ابعاد را می‌توان به انتخاب ویژگی و استخراج ویژگی تقسیم کرد. [۱۳]

۷-۲ تجزیه و تحلیل مؤلفه اصلی

تجزیه و تحلیل مؤلفه اصلی^{۵۰} فرآیند محاسبه مؤلفه‌های اصلی و استفاده از آنها برای انجام تغییر مختصات^{۵۱} بر روی داده‌ها است، به گونه‌ای حداکثر واریانس^{۵۲} داده‌ها حفظ شود. گاهی اوقات این الگوریتم فقط از چند مؤلفه اصلی استفاده می‌کند و بقیه را نادیده می‌گیرد. [۱۴]

من در این قسمت صرفاً برای اینکه توضیح بدهم کاهش بُعد گزینه مناسبی برای من در این پروژه نیست گزارش‌های پیش‌رو را گرفتم. اما پیش از آن نیز با توجه به پایین بودن تعداد ویژگی‌ها کاملاً مشخص بود که قرار نیست کاهش بُعدی روی داده‌ها صورت بگیرد. ضمناً با توجه به شکل ۱۴ می‌توان گفت متغیرهای پیشگوی آسانسور و انباری همبستگی بالایی با هم ندارند. اما اگر ما از مؤلفه‌های اصلی در الگوریتم PCA کمک بگیریم، تصویر این دو متغیر پیشگو در راستای دو برداری که بیشترین سهم را در حفظ واریانس داده‌ها دارند، همبستگی قابل توجهی با هم خواهند داشت. این حقیقت را می‌توانید از طریق بررسی شکل ۲۱ برداشت کنید.

	PC۱	PC۲	PC۳	PC۴	PC۵
Standard deviation	۱/۴۴۰۷	۱/۴۹۴	۰/۸۹۵۶	۰/۶۸۷۷	۰/۵۷۲۷۸
Proportion of Variance	۰/۴۱۵۲	۰/۲۶۴۲	۰/۱۶۰۴	۰/۰۹۴۶	۰/۰۶۵۶۲
Cumulative Proportion	۰/۴۱۵۲	۰/۶۷۹۴	۰/۸۳۹۸	۰/۹۳۴۴	۱/۰۰۰۰۰

جدول ۶: جدول اهمیت مؤلفه‌ها

من برای اینکه بدانم در این بخش باید به کاهش بُعد تن بدهم یا خیر، منابع مختلفی را بررسی کردم. یکی از پاسخ‌هایی که به این سؤال کلی که کاهش بُعد خوب است یا خیر داده شد به نظر پخته‌تر از بقیه جواب‌ها آمد. من با ذکر نام نویسنده این مطلب را به زبان فارسی ترجمه و سپس از قول ایشان نقل می‌کنم.

هیچ پاسخ کلی درست یا غلطی برای کاهش بُعد وجود ندارد. درستی یا غلط بودن کاملاً بستگی به

⁴⁷Machine Learning

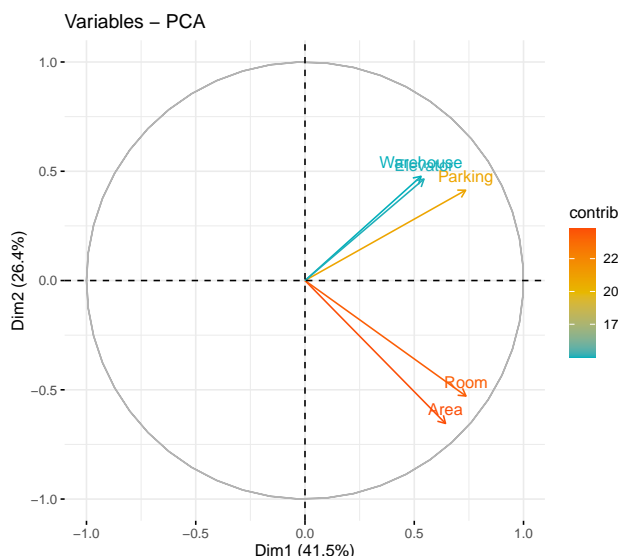
⁴⁸Dimension Reduction

⁴⁹Supervised

⁵⁰Principal Component Analysis

⁵¹Coordinate

⁵²Variance



شکل ۲۱: نمودار تحلیل PCA

موقعیت دارد. اجازه دهید چند مثال را برای نشان دادن حالات مختلف بیان کنم. این مثال‌ها ممکن است جامع نباشند، اما ذهنیت خوبی از آنچه باید انجام شود به شما می‌دهند.

- شما در حال ساخت یک مدل رگرسیون با ۱۰ متغیر هستید، که برخی از ویژگی‌های آنها همبستگی ملایمی دارند. در این حالت نیازی به کاهش ابعاد نیست. از آنجایی که بین متغیرها همبستگی کمی وجود دارد، همه آنها اطلاعات جدیدی را در اختیار مدل می‌گذارند. شما باید همه متغیرها را در ترکیب نگه دارید و یک مدل بسازید.

- مجدداً در حال ساخت یک مدل رگرسیون هستید. این بار ۲۰ متغیر پیشگو دارید و برخی از این متغیرها همبستگی بالایی دارند. به عنوان مثال، به دو نوع مجموعه داده فکر کنید.

داده‌های مخابراتی: تعداد تماس‌های مشتری در یک ماه و صورتحساب ماهانه دریافتی او.

داده‌های بیمه: تعداد بیمه‌نامه‌ها و کل حق بیمه فروخته شده توسط یک نماینده یا شعبه.

در این موارد، به دلیل اینکه تعداد متغیرهای پیشگو محدودی با همبستگی بسیار بالا دارید، می‌توانید فقط یکی از این متغیرها را به مدل خود اضافه کنید. با وارد کردن همه متغیرها، چیز

از این‌ها که باشد، نباید اجازه دهیم این داده‌ها، مدل ما را فریب دهند. به خصوص که پروژه ما یک پیشگویی دودویی نیست و الگوریتم‌های رگرسیونی بسیار به این داده‌های دورافتاده حساس هستند.

به طور خلاصه با بررسی داده‌ها و شکل مذکور به این نتیجه رسیدم که ثبت‌های جدول ۷ را بررسی کنم. پس از بررسی به این نتیجه رسیدم که داده‌های جدول ۸ باید از مجموعه داده‌ها حذف شوند.

Area	Room	Parking	Warehouse	Elevator	Address	Price
۱۳۰	۲	۱	۱	۱	Dorous	۱.۶۹×۱۰^{10}
۳۶۵	۴	۱	۱	۱	Dorous	۴.۷۴۵×۱۰^{10}
۳۵۰	۴	۱	۱	۱	Niavaran	۸.۰۵×۱۰^{10}
۳۵۰	۳	۱	۱	۰	Pasdaran	۱.۰×۱۰^9
۱۴۵	۳	۱	۱	۱	Pasdaran	۱.۴۵×۱۰^9
۷۰۵	۵	۱	۱	۰	Abazar	۹.۱×۱۰^{10}
۸۳	۲	۱	۱	۱	Ozgol	۵.۵×۱۰^7
۳۵۰	۴	۱	۱	۱	Niavaran	۷.۲×۱۰^{10}
۳۰۰	۳	۱	۱	۰	Golestan	۱.۴×۱۰^{10}
۳۱۰	۳	۱	۱	۱	Aqdasieh	۷.۴۴×۱۰^{10}
۵۳۰	۴	۱	۱	۱	Dorous	۷.۹۵×۱۰^{10}
۶۰	۱	۰	۰	۰	Shahr-e-Ziba	۱.۸×۱۰^{10}
۳۲۰	۵	۰	۰	۰	Sattarkhan	۲.۸×۱۰^{10}
۵۴	۱	۱	۱	۱	West Ferdows Boulevard	۲.۲۵×۱۰^{10}
۴۵	۱	۰	۰	۱	Si Metri Ji	۸.۹×۱۰^9

جدول ۷: جدول کاندیدای داده‌کاهی

در مورد هر یک از ثبت‌های جدول ۸ توضیح مختصری می‌دهم که توجیهم را برای حذف این داده‌ها بیان کنم.

- یک خانه ۳۵۰ متری سه خواب در منطقه پاسداران مجهز به انباری و پارکینگ فقط یک میلیارد تومان ثبت شده است که انتظار داریم این عدد حداقل در کمترین حالت ۱۰ برابر قیمت ثبت شده باشد.

- یک خانه ۱۴۵ متری سه خواب در منطقه پاسداران مجهز به هر سه امکانات پارکینگ، انباری و آسانسور که فقط ۱ میلیارد و ۴۵۰ میلیون تومان ثبت شده است. تخمین ذهنی من به هیچ وجه این داده را صحیح نمی‌داند.

- یک خانه ۸۳ متری ۲ خواب در منطقه ازگل یک خواب بدون انباری، آسانسور و پارکینگ که فقط ۵۵۰ میلیون تومان ثبت شده است. این خانه با توجه به مساحت نه چندان زیاد و عدم مجهز بودن به امکانات مذکور قیمت بالایی نمی‌تواند داشته باشد ولی ۵۵۰ میلیون تومان برای یک خانه در شمال تهران حتی در صورتی که خانه در طبقه منفی ۱ واقع شده باشد هم بسیار دور از ذهن است.

- خانه‌ای ۶۰ متری یک خواب بدون انباری، پارکینگ و آسانسور در منطقه شهر زیبا که یک منطقه گران‌قیمت

جدیدی به مدل اضافه نمی‌کنید. حتی مقدار زیادی از این مقدار اضافی که اتفاقاً همبستگی بالایی هم دارند می‌تواند نویز باشد. شما هنوز هم می‌توانید از استفاده از هر تکنیک کاهش ابعادی رسمی در این مورد صرف‌نظر کنید (اما با وارد کردن تعداد محدودی از متغیرهای پیشگو به جای همه آنها ذاتاً کاهش بُعد را انجام داده‌اید).

- حال فرض کنید که شما ۵۰۰ متغیر پیشگو از این قبیل دارید که برخی از آنها با یکدیگر همبستگی دارند. به عنوان مثال، خروجی داده‌های سنسورها از یک گوشی هوشمند. با توجه به تعداد بالای ویژگی‌ها و همبستگی نسبی خیلی از آنها که با هم مرتبط هستند و اتفاقاً خیلی هم زیاد هستند، نمی‌توان به طور جداگانه فهمید که کدام یک به کدام یک مرتبط است. در این مورد، قطعاً باید از تکنیک‌های کاهش ابعاد استفاده کنید. شما ممکن است معنای واقعی این بردارها را بفهمید یا نه، اما همچنان می‌توانید با مشاهده نتیجه الگوریتم‌های کاهش بُعد تأثیر بسیاری از این ویژگی‌ها را درک کنید.

در بسیاری از سناریوهایی مانند این، همچنین خواهید دید که بیش از ۹۰ درصد اطلاعات را با کمتر از ۱۵ تا ۲۰ درصد متغیرها حفظ خواهید کرد. از این رو، الگوریتم‌های کاهش بُعد می‌توانند کاربردهای بسیاری داشته باشند.

[۱۵] Kunal Jain from Gurgaon, India

با توجه به توضیحات مذکور من نیز از اعمال الگوریتم‌های کاهش بُعد بر روی داده‌های خود، خودداری می‌کنم.

۸-۲ داده کاهی

معمولاً در پروژه‌های داده‌کاوی خیلی کم پیش می‌آید که بخشی تحت عنوان داده کاهی^{۵۳} داشته باشیم. اما دلیلی که من را به این نتیجه رساند تا در این پروژه این بخش را بگنجانم، شکل ۱۳ بود. شما در این شکل می‌بینید که داده‌هایی از محله‌های مختلف هستند، که به شکل غیرقابل توجیهی با قیمت میانگین منطقه فاصله دارند. وجود این داده‌ها ممکن است دلایل زیادی داشته باشد. از وجود نویز در داده‌ها گرفته تا شرایط خاص ملک، همه می‌توانند دلایل دورافتادگی این داده‌ها باشند. دلیل هر کدام

⁵³Data Reduction

```

5 test.rows<-setdiff(rownames(X_add),union(
  train.rows,valid.rows))
6 training_set<-X_add[train.rows,]
7 validation_set<-X_add[valid.rows,]
8 test_set<-X_add[test.rows,]
9
10
11 #Data standardization
12 X_standard_scaler.train<-training_set
13 X_standard_scaler.valid<-validation_set
14 X_standard_scaler.test<-test_set
15 #- - - - -
16 for(i in c(1,2,7)){
17   X_standard_scaler.train[,i]=scale(
18     training_set[,i])
19 }
20 #- - - - -
21 for(j in c(1,2,7)){
22   for (i in 1:dim(validation_set)[1]){
23     X_standard_scaler.valid[i,j]=(
24       validation_set[i,j]-mean(training_set[,j]
25         ))/sd(training_set[,j])
26   }
27 }
28 #- - - - -
29 for(j in c(1,2,7)){
30   for (i in 1:dim(test_set)[1]){
31     X_standard_scaler.test[i,j]=(test_
32       set[i,j]-mean(training_set[,j]))/sd(
33         training_set[,j])
34   }
35 }
36 #- - - - -
37
38 #target_encoding
39 target_encod <- build_target_encoding( X_
40   standard_scaler.train, cols_to_encode =
41   "Address",target_col = "Price",
42   functions = c("median"))
43 X_standard_scaler.train=target_encode( X_
44   standard_scaler.train, target_encoding =
45   target_encod)[,c(2,3,4,5,6,8,7)]
46 X_standard_scaler.valid=target_encode( X_
47   standard_scaler.valid, target_encoding =
48   target_encod)[,c(2,3,4,5,6,8,7)]
49 X_standard_scaler.test=target_encode( X_
50   standard_scaler.test, target_encoding =
51   target_encod)[,c(2,3,4,5,6,8,7)]
52
53 #For convenience
54 train.data<-X_standard_scaler.train
55 valid.data<-X_standard_scaler.valid
56 test.data<-X_standard_scaler.test

```

محسوب نمی‌شود، ۱۸ میلیارد تومان ارزش‌گذاری شده است. این خانه حتی اگر یک دهم این مقدار قیمت گذاری می‌شد باز هم ارزش خرید بالایی نداشت.

- خانه‌ای ۵۴ متری یک خواب مجهز به انباری، پارکینگ و آسانسور در منطقه بلوار فردوس که مانند شهر زیبا یک منطقه گران‌قیمت محسوب نمی‌شود، بیش از ۲۲ میلیارد تومان (به عبارتی متری بیش از ۴۰۰ میلیون تومان) ارزش‌گذاری شده است. نیاز به توضیح نیست که این ثبت هم به احتمال قریب به یقین به اشتباه ثبت شده است.

- خانه‌ای ۴۵ متری یک خواب مجهز آسانسور بدون انباری و پارکینگ در منطقه سی متری جی که از مناطق بسیار ارزان‌قیمت تهران به شمار می‌رود، نزدیک به ۹ میلیارد تومان (به عبارتی متری حدوداً ۱۴۰ میلیون تومان) قیمت‌گذاری شده است. کافی است یادآور شوم که چارک سوم قیمت هر متر خانه در مناطق زعفرانیه، ولنجک و نیاوران از این عدد کمتر است.

Area	Room	Parking	Warehouse	Elevator	Address	Price
۳۵۰	۳	۱	۱	۰	Pasdaran	1×10^9
۱۴۵	۳	۱	۱	۱	Pasdaran	145×10^9
۸۳	۲	۱	۱	۱	Ozgol	55×10^7
۶۰	۱	۰	۰	۰	Shahr-e-Ziba	18×10^{10}
۵۴	۱	۱	۱	۱	West Ferdows Boulevard	225×10^{10}
۳۵	۱	۰	۰	۱	Si Metri Ji	88×10^9

جدول ۸: جدول نهایی داده‌هایی که در بخش داده‌کاهی بایستی حذف شوند

۹-۲ اعمال تغییرات نتیجه‌گیری شده

با اجرای کد پیش رو داده‌های مذکور از ماتریس اصلی داده‌های من پاک می‌شوند.

```

1 #deleting outliers
2 X<-X[-c(1442,1516,2767,3128,3390,3416),]

```

واضح است که کدگذاری من روی متغیرهای رسته‌ای و همچنین افراز داده‌های من باید با توجه به همین ماتریس جدید داده‌ها به‌روزرسانی شوند.

```

1 #Data splitting
2 set.seed(1)
3 train.rows<-sample(rownames(X_add),dim(X)[1]
  *0.6)
4 valid.rows<-sample(setdiff(rownames(X_add),
  train.rows),dim(X)[1]*0.3)

```

۳ رگرسیون خطی چندگانه

اکنون آماده هستیم تا اولین مدل پیش‌گویانه را روی داده‌های مجموعه آموزشی اعمال کنیم. همان‌طور که نام فصل گویاست از مدل رگرسیون خطی چندگانه^{۵۴} برای پیشگویی استفاده می‌کنیم. دلیل این‌که این مدل را قبل از مدل‌های دیگر استفاده کرده‌ام این است این مدل در انتخاب ویژگی‌ها به من کمک بزرگی می‌کند. در ادامه خواهیم گفت که کدام ویژگی‌ها اهمیت بیشتری دارند.

```
1 regressor=lm(formula=Price~.,data=train.data)
2 #summary(regressor)
```

با اجرای دستور فوق مدل رگرسیون^{۵۵} من اجرا می‌شود. یک سؤال مهم اینجاست که آیا مانده‌های^{۵۶} ما از توزیع نرمال پیروی می‌کنند؟ اگر مانده‌های من از توزیع نرمال پیروی نکنند، مدل رگرسیونی دچار اشکال می‌شود. لذا اول این نکته را به کمک دستورهای پیشرو بررسی می‌کنم.

```
1 par(mfcol=c(2,1))
2 hist(resid(regressor),,breaks=80,prob = TRUE,
3      ylim=c(0,2.1))
4 lines(density(resid(regressor)),lwd = 2,col = "chocolate3")
5 hist(resid(regressor),,breaks=80,prob = TRUE,
6      ylim=c(0,2.1),xlim=c(-1.5,1.5))
7 lines(density(resid(regressor)),lwd = 2,col = "chocolate3")
8 par(mfcol=c(1,1))
9 qqnorm(resid(regressor))
10 qqline(resid(regressor),lwd = 2,col = "chocolate3")
```

شکل ۲۲ نمایانگر این است که توزیع مانده‌های مدنظر ما تقریباً نرمال است اما این حقیقت را نیز می‌توان از تحلیل نمودار چندک چندک^{۵۷} در شکل ۲۳ نیز دریافت.

اکنون می‌توانیم با خاطری آسوده‌تر مدل رگرسیونی خود را روی داده‌ها اعمال کنیم. اما قبل از هر چیز دوست دارم به اهمیت بررسی این فرض توزیع نرمال مانده‌ها اشاره کنم. اساتید بزرگوار آقایان Donald A. Pierce و Daniel W. Schafer، از اعضای هیئت‌علمی گروه آمار دانشگاه ایالتی اورگان آمریکا در مقاله‌ای [۱۶] مفصل به اهمیت این موضوع می‌پردازند و خاطرنشان می‌شوند که بسیاری از پروژه‌هایی که از مدل رگرسیونی

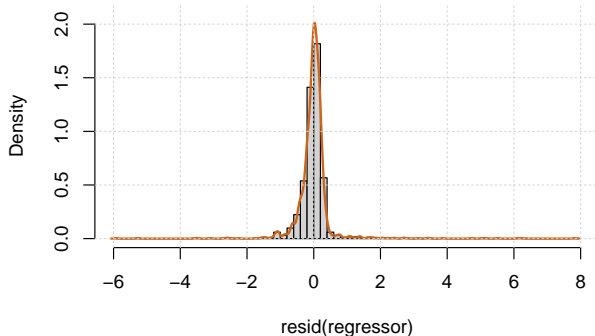
⁵⁴Multiple Linear Regression

⁵⁵Regression

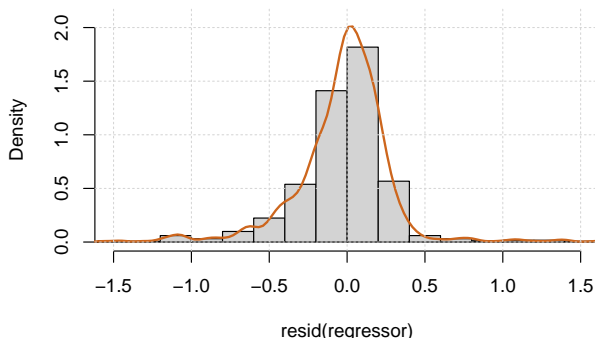
⁵⁶Residual

⁵⁷QQ Plot

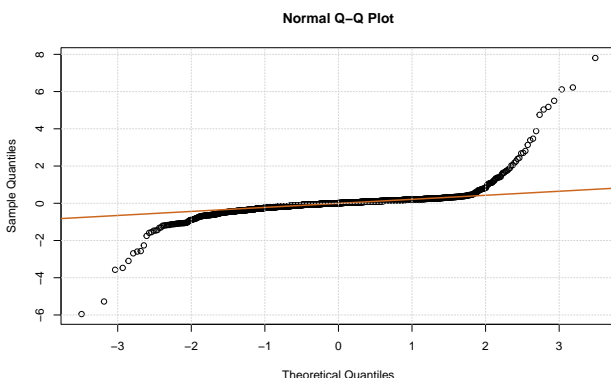
Histogram of resid(regressor)



Histogram of resid(regressor)



شکل ۲۲: توزیع مانده‌های مدل رگرسیونی روی متغیر پاسخ



شکل ۲۳: نمودار چندک چندک مدل رگرسیونی روی متغیر پاسخ

چندگانه استفاده می‌کنند این فرض را در نظر نمی‌گیرند و به همین دلیل نتایج این مدل‌ها قابل اعتماد نیستند.

قبل از هر چیز از تابع summary که پیش‌تر آن را به صورت توضیح^{۵۸} نوشته بودیم خروجی می‌گیریم تا مهم‌ترین ویژگی‌ها را مشخص کنیم.

⁵⁸Comment

یک ویژگی مهم تلقی شده بود از این ویژگی شروع می‌کنم. کد پیش رو حالت‌های مختلف را بررسی می‌کند. من از قرار دادن خروجی کد به علت طولانی شدن صرف نظر می‌کنم اما خلاصه تفسیر آن را قرار خواهم داد.

```
1 #variable selection
2 base_model=lm(Price~Area,data=train.data)
3 step_model=step(base_model,scope=list(upper=
  regressor , lower= ~1), direction="
  forward",trace=T)
4 print(step_model)
5 forward_pred=predict(step_model,newdata=
  valid.data)
6 #data.frame(TruePrice=valid.data$Price ,
  Predicted_price=forward_pred)
7 MSE(valid.data$Price,forward_pred)
8 #- - - - -
9 base_model=lm(Price~Area,data=train.data)
10 step_model=step(base_model,scope=list(upper=
  regressor , lower= ~1), direction="
  backward",trace=T)
11 print(step_model)
12 backward_pred=predict(step_model,newdata=
  valid.data)
13 #data.frame(TruePrice=valid.data$Price ,
  Predicted_price=backward_pred)
14 MSE(valid.data$Price,backward_pred)
15 #- - - - -
16 base_model=lm(Price~Area,data=train.data)
17 step_model=step(base_model,scope=list(upper=
  regressor , lower= ~1), direction="both"
  ,trace=T)
18 print(step_model)
19 stepwise_pred=predict(step_model,newdata=
  valid.data)
20 #data.frame(TruePrice=valid.data$Price ,
  Predicted_price=stepwise_pred)
21 MSE(valid.data$Price,stepwise_pred)
```

خروجی‌های بازبراشی که از روش‌های forward و stepwise به دست آوردیم نتیجه‌ای مانند رگرسیون اولیه را روی داده‌های مجموعه اعتبار سنجی به من داده‌اند. توجه شما را به گزارش پیشرو جلب می‌نمایم.

```
1 > print(step_model)
2
3 Call:
4 lm(formula = Price ~ Area + Price_median_by_
  Address + Room +
5 Elevator + Warehouse + Parking, data = train.
  data)
6
```

```
1 > summary(regressor)
2
3 Call:
4 lm(formula = Price ~ ., data = train.data)
5
6 Residuals:
7      Min       1Q   Median       3Q      Max
8 -5.9518 -0.1515  0.0086  0.1410  7.8056
9
10 Coefficients:
11              Estimate Std. Error t value Pr(>|t|)
12 (Intercept)    0.09112    0.04831   1.886  0.05943 .
13 Area           0.49521    0.01812  27.334 < 2e-16 ***
14 Room           0.04974    0.01775   2.802  0.00513 **
15 Parking        -0.06978    0.04294  -1.625  0.10429
16 Warehouse      0.10966    0.05080   2.159  0.03099 *
17 Elevator       -0.04914    0.03451  -1.424  0.15466
18 Price_median_by_Address 0.76646    0.02552  30.030 < 2e-16 ***
19 ---
20 Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
21
22 Residual standard error: 0.5729 on 2073 degrees of freedom
23 Multiple R-squared:  0.6728, Adjusted R-squared:  0.6718
24 F-statistic: 710.3 on 6 and 2073 DF,  p-value: < 2.2e-16
```

با بررسی پی-مقدار^{۵۹} ها متوجه می‌شویم ویژگی‌های مساحت و آدرس بیشترین اهمیت را دارا هستند. بعد از آن ویژگی تعداد اتاق و پس از آن ویژگی داشتن یا نداشتن انباری از نظر مدل رگرسیونی من مهم تلقی شده‌اند. اما آیا درست است که بقیه ویژگی‌ها را حذف کنم؟ قطعاً نه! من برای پاسخ به این سؤال که بهترین نتیجه را با وارد کردن کدام پیشگوها به مدل می‌گیرم، به داده‌های اعتبار سنجی نیاز دارم. لازم به ذکر است مدل رگرسیونی به طور کلی نیاز به داده‌های مجموعه تست ندارد. اما اینجا چون از داده‌های مجموعه اعتبار سنجی برای بهبود یادگیری مدل استفاده می‌کنم من در انتهای پروژه در صورتی که این مدل بالاترین دقت را داشته باشد، مدل رگرسیونی نهایی را به کمک مجموعه تست ارزش‌گذاری خواهم کرد.

اکنون با ارزیابی مدل رگرسیونی روی داده‌های مجموعه اعتبار سنجی آغاز می‌کنیم. ضمناً مبنای ما برای ارزیابی عملکرد در این پروژه متر میانگین توان دوم خطاها^{۶۰} یا همان MSE است.

```
1 > y_pred=predict(regressor,newdata=valid.
  data)
2 > #data.frame(TruePrice=valid.data$Price ,
  Predicted_price=y_pred)
3 > MSE(valid.data$Price,y_pred)
4 [1] 0.4927284
```

دریافتیم که میانگین توان دوم خطاها در تلاش اول من برابر با ۰.۴۹۲۷۲۸۴ بود اکنون می‌خواهم با استفاده از روش‌های پیشرو^{۶۱}، پسرو^{۶۲} و گام‌به‌گام^{۶۳} اقدام به بررسی مجدد خروجی مدل رگرسیونی با ویژگی‌های مختلف کنم. با توجه به اینکه مساحت

^{۵۹}P-Value

^{۶۰}Mean Squared Error

^{۶۱}Forward

^{۶۲}Backward

^{۶۳}Stepwise

7	Coefficients:			
8	(Intercept)	Area	Price_median_by_Address	Room
9	0.09112	0.49521	0.76646	0.04974
10	Elevator	Warehouse	Parking	
11	-0.04914	0.10966	-0.06978	

تفسیر من این است که انتخاب ویژگی رگرسیون خطی چندگانه به من می‌گوید همه متغیرهای من بامعنی هستند و برای پیشگویی مدل استفاده می‌شوند اما چیزی که برخلاف شهود من است این است که ویژگی‌های داشتن پارکینگ و آسانسور ضریب منفی دارد. یعنی مدل من خانه‌های دارای پارکینگ و آسانسور را از خانه مشابه بدون این امکانات ارزان‌تر پیشگویی می‌کند.

به هر روی عدد نهایی برای میانگین توان دوم خطاها در این روش روی مجموعه اعتبار سنجی ۴۹۲۷۲۸۴٪ است.

شایان ذکر است که همین شاخص وقتی که از متغیر ظاهری^{۶۴} استفاده می‌کردم، برابر ۶۳۰۰۷۵۳٪ و اگر داده‌ها را نرمال‌سازی نیز نمی‌کردم آماره‌های میانگین توان دوم خطاها و میانگین قدرمطلق خطاها^{۶۵} به ترتیب $10^{19} \times 3,886042$ و ۲۹۰۰۴۹۸۹۸۷ می‌شدند.

⁶⁴Dummy Variable

⁶⁵Mean Absolute Error

۴ مدل درخت تصمیم

الگوریتم درخت تصمیم^{۶۶}، یا درخت رگرسیون^{۶۷} (در اینجا چون روی پیشگویی کار می‌کنیم درخت رگرسیون مفهوم دارد) دارای فرا پارامتر^{۶۸}هایی است که احتیاج به تنظیم^{۶۹} فرا پارامترها یا بهینه‌سازی فرا پارامتر^{۷۰} دارند. یک راه ساده جستجوی شبکه‌ای^{۷۱} است. در این روش با جایگزین کردن بخشی از فرا پارامترها تلاش می‌کنیم فرا پارامترهایی که بیشترین دقت (کمترین خطا) را روی مجموعه اعتبار سنجی کسب کرده‌اند را پیدا کنیم. سپس مدل را به عنوان مدل نهایی در این بخش (و نه در کل پروژه) معرفی کنیم.

```
21 best_dt_model<-dt_model_grid_search[[which.
    min(MSE_val_dt)]]
```

تا این مرحله بهترین مدل را روی داده‌های اعتبارسنجی پیدا کردیم. ابتدا یک گزارش بگیریم تا فرا پارامترهای منتخب مشخص شوند.

```
1 > best_dt_model$control
2 $minsplit
3 [1] 12
4 $minbucket
5 [1] 4
6 $cp
7 [1] 0.000345267
8 $maxcompete
9 [1] 4
10 $maxsurrogate
11 [1] 5
12 $usesurrogate
13 [1] 2
14 $surrogatestyle
15 [1] 0
16 $maxdepth
17 [1] 5
18 $xval
19 [1] 10
```

حال وقت آن رسیده است که خطای مدل نهایی این قسمت را روی مجموعه اعتبار سنجی بسنجیم.

```
1 >y_pred=predict(best_dt_model,newdata=valid.
    data)
2 >MSE(valid.data$Price,y_pred)
3 [1] 0.2755766
```

لذا خطای ما روی مجموعه اعتبار سنجی برابر ۰.۲۷۵۵۷۶۶ شد. اما از آنجایی که از درخت رگرسیون می‌توان تعدادی نمودار رسم کرد که به کمک آنها به درک درستی از ویژگی‌ها برسیم، این فصل را در این نقطه تمام نمی‌کنم و قبل از آن این نمودارها را بررسی خواهیم کرد.

با توجه به شکل ۲۴ می‌توان دریافت که مهم‌ترین ویژگی‌ها مانند آنچه در فصل قبل مشاهده کردیم، ویژگی‌های آدرس و مساحت هستند. اما در برخی مواقع، الگوریتم، برای تشخیص بهتر از ویژگی تعداد اتاق نیز بهره برده است. این نشان می‌دهد ویژگی تعداد اتاق از نظر این مدل اهمیت بیشتری نسبت به سایر ویژگی‌ها داشته است. این میزان اهمیت دادن به داده‌ها را نیز در اجرای الگوریتم رگرسیون خطی چندگانه نیز مشاهده کردیم. اما در آن فصل مدل تشخیص داده بود برای خطای کمتر به همه ویژگی‌ها نیاز دارد. در حالی که در این مدل مشاهده کردیم که چنین نیست.

```
1 #Define hyperparameters for grid search
2 minsplit<-seq(1,30,1)
3 maxdepth<-seq(1,30,1)
4 cp_float<-seq(0,30,1) #This variable must be
    a float number. Due to the limitations
    of the expand.grid function, I will
    modify this parameter later.
5 hyperparam_grid_dt_model<-expand.grid(
    minsplit=minsplit,maxdepth=maxdepth,cp=
    cp_float)
6 num_models<-nrow(hyperparam_grid_dt_model)
7 dt_model_grid_search<-list()
8 #Apply all models
9 for(i in 1:num_models){
10     minsplit<-hyperparam_grid_dt_model$
        minsplit[i]
11     maxdepth<-hyperparam_grid_dt_model$
        maxdepth[i]
12     cp_float<-1/((sqrt(2))^(hyperparam_grid_
        dt_model$cp[i]))
13     dt_model_grid_search[[i]]<-rpart(formula
        =Price~., data=train.data, method="anova
        ", control=rpart.control(maxdepth=
        maxdepth,minsplit=minsplit,cp =cp_float)
        )
14 }
15 #compute the loss (MSE based)
16 MSE_val_dt<-c()
17 for(i in 1:num_models){
18     MSE_val_dt[i]<-MSE(predict(dt_model_grid
        _search[[i]], newdata =valid.data),valid
        .data[,7])
19 }
20 #identify the best model
```

⁶⁶Decision Tree

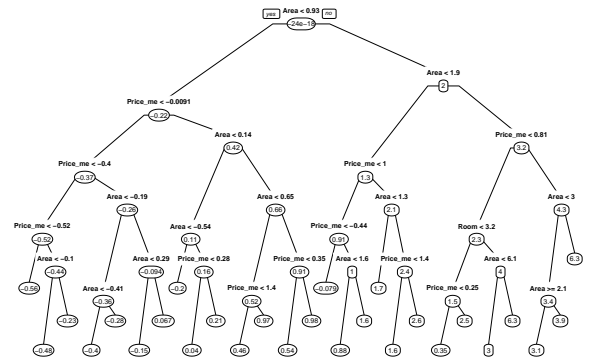
⁶⁷Regression Tree

⁶⁸Hyperparameter

⁶⁹Tuning

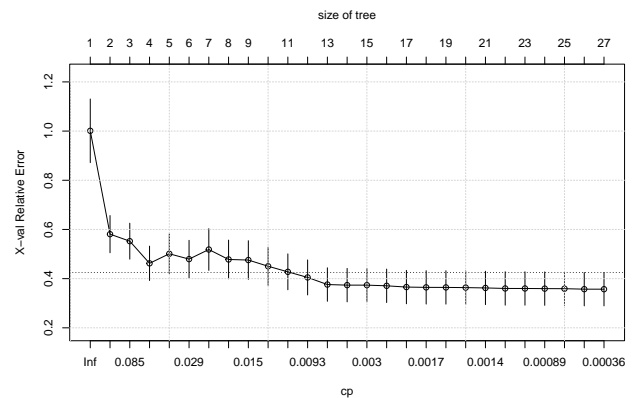
⁷⁰Hyperparameter Optimization

⁷¹Grid Search



شکل ۲۴: نمودار درخت رگرسیون

تا اینجای کار الگوریتم درخت رگرسیون بهترین نتیجه (کمترین میانگین مربعات خطاها) را بر اساس داده‌های مجموعه اعتبار سنجی به من داده است. به نظر من ضعف این روش با توجه به ماهیت الگوریتم این است که در پیشگویی قیمت داده‌های دور افتاده از نظر متغیر پاسخ، نباید خوب عمل کند.



شکل ۲۵: نمودار خطا- complexity parameter

شکل ۲۵ نیز تأیید کننده انتخاب پیچیدگی پارامتر ^{۷۲} برای مدل درخت رگرسیونی است.

⁷²Complexity Parameter

۵ مدل k نزدیک‌ترین همسایه

سومین الگوریتمی که قرار است روی داده‌ها اعمال کنیم الگوریتم K نزدیک‌ترین همسایه^{۷۳} است. چالش مهم من در این الگوریتم محاسبه K مناسب است. لذا بایستی از داده‌های مجموعه اعتبار سنجی، اقدام به بهینه‌سازی فرا پارامتر نمایم.

لازم به ذکر است که من بایستی اجتماع ویژگی‌هایی که الگوریتم‌های درخت رگرسیون و رگرسیون خطی چندگانه آنها را با اهمیت تلقی کرده‌اند را در این الگوریتم وارد می‌کردم. اما من با بررسی همه حالت‌های مختلف به این نتیجه رسیدم که با وارد کردن ویژگی‌هایی که درخت رگرسیون آنها را مهم تلقی کرده بود، بالاترین بازدهی را خواهیم داشت. لذا از ارائه کدهای مربوط به اجتماع ویژگی‌ها، خودداری می‌کنم. اما به طور کلی باید بگویم میانگین توان دوم خطاها در حالتی که همه ویژگی‌ها را وارد مدل کنیم در بهترین حالت K مناسب، برابر 0.5201542% خواهد شد که با اختلاف بالاترین خطای به دست آمده در این پروژه است.

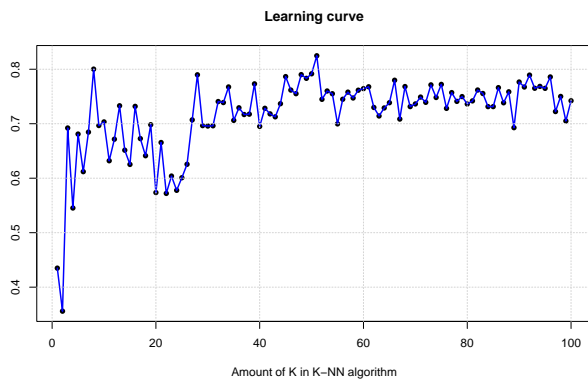
```
1 MSE_val_knn<-c()
2 for(i in 1:100){
3   y_pred = knn(train = train.data[,c
4     (1,2,6)],
5     test = valid.data[,c(1,2,6)],
6     cl = train.data[, 7],
7     k = i,
8     prob = FALSE)
9   MSE_val_knn[i]=MSE(as.numeric(as.matrix(
10    y_pred)),valid.data[,7])
11 }
```

اکنون میانگین توان دوم خطاها همه حالت‌های ممکن در MSE_val_knn ذخیره شده است. حال برای به دست آمدن نتایج خواهیم داشت:

```
1 > which.min(MSE_val_knn)
2 [1] 2
3 > MSE_val_knn[which.min(MSE_val_knn)]
4 [1] 0.3559102
```

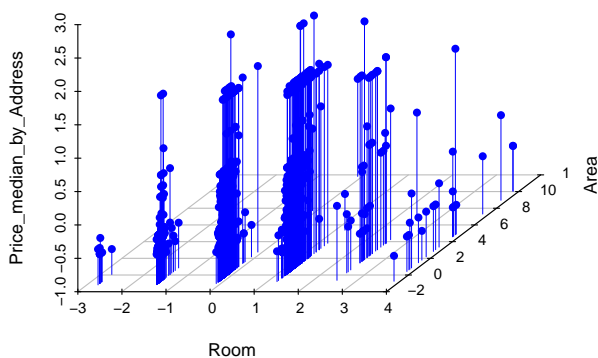
به عبارتی من با قراردادن $K = 2$ بهترین نتیجه را از این مدل می‌گیرم. میانگین توان دوم خطاها در مدل نهایی من برابر 0.3559102% است. این خطا از مدل رگرسیون خطی چندگانه نتیجه بهتری است اما همچنان بهترین نتیجه مربوط به درخت رگرسیون است.

میانگین توان دوم خطاها با مقادیر مختلف K به صورتی است که در شکل ۲۶ نمایش داده شده است.



شکل ۲۶: منحنی یادگیری

برای درک بهتر از داده‌هایی که ۲ نزدیک‌ترین همسایه روی آنها محاسبه می‌شوند می‌توان به نمودار ۲۷ مراجعه کرد. این نمودار تصویر واضحی از آنچه در این الگوریتم رخ می‌دهد به ما نشان می‌دهد. ضمناً توجه داشته باشید که مقیاس همه محورها برای نمایش بهتر نمودار یکسان نیست.



شکل ۲۷: نمودار پراکنش سه بُعدی

⁷³K-Nearest Neighbor

۶ مدل شبکه عصبی

تلاش اولیه من این بود که بتوانم یک مدل یادگیری عمیق^{۷۴} روی داده‌ها در زبان برنامه‌نویسی R ایجاد کنم. ولی به علت عدم تسلط کافی من روی این زبان برنامه‌نویسی این مهم محقق نشد. مشکل اساسی من در اعمال مدل ساده شبکه عصبی^{۷۵} روی داده‌ها این بود که اولاً در تعریف تابع فعال ساز^{۷۶} محدودیت زیادی داشتم ثانیاً توان دسترسی مناسبی به الگوریتم بهینه سازی نداشتم. به هر حال با توجه به نکات گفته شده ابزار مناسبی برای استفاده حداکثری از توان شبکه عصبی برای من میسر نبود.

با این حال من مدل شبکه عصبی را روی همین داده‌ها با زبان برنامه‌نویسی پایتون پایتون^{۷۷} اعمال کردم و نتیجه به طور شگفت انگیزی از آنچه در پیش رو خواهید دید، بهتر بود.

ضمناً باید اعلام کنم چون داده‌ها رگرسیونی هستند و نیاز به تابع فعال‌ساز دو دویی ندارند نیازی به تجانس داده‌ها در بازه ۰ و ۱ نیست. (توجه داشته باشید که داده‌ها از قبل یکبار نرمالیده شده‌اند.) [۷۷]

```
24 startweights="NULL",
25 algorithm=algorithm[AL[i]],
26 act.fct=activation_function[
  AF[i]],
27 stepmax=10^5)
28 y_pred_nn<-compute(nn, valid.data[,1:6])
29 MSE_val_nn[i]<-MSE(y_pred_nn$net.result,
  valid.data[,7])
30 HL[i]<-hidden_layer
31 }
```

از آنجایی که شبکه عصبی به راحتی دچار بیش برازش می‌شود من ۳۰۰ بار الگوریتم را با بردارهایی که نمایانگر لایه پنهان^{۷۸} های مختلف و تصادفی (چه از نظر تعداد گره‌ها چه از نظر تعداد لایه‌ها) هستند اجرا کردم. همچنین از دیگر فرآیندهای مختلف اعم از توابع فعال‌ساز و نرخ یادگیری^{۷۹} های تصادفی و ... در هر حلقه استفاده کرده‌ام. در کمترین خطایی که روی مجموعه اعتبار سنجی حاصل شد را به عنوان مدل نهایی در نظر گرفتم.

```
1 > best_i<-where.min[MSE_val_nn]
2 >HL[best_i] #To display the best hiddenlayer
  vector
3 [[1]]
4 [1] 11
5 >activation_function[AF[best_i]] #To display
  the best Activation Function
6 [1] "tanh"
7 >algorithm[AL[best_i]] #To display the best
  Algorithm
8 [1] "rprop+"
9 >LR[best_i] #To display the best Learning
  Rate
10 [1] 0.7052665
11 >MSE_val_nn[best_i]
12 [1] 0.695432
```

در نتیجه، میانگین توان دوم خطاهای بهترین مدل شبکه عصبی روی مجموعه اعتبار سنجی ۶۹۵۴۳۲٪ است. البته با توجه به تعداد کم داده‌ها نباید انتظار زیادی از این الگوریتم داشت. چرا که توان پیشگویی شبکه‌های عصبی روی داده‌های زیاد است. به طور کلی یکی از ضعف‌های شبکه‌های عصبی این است که هنگامی که تعداد داده‌ها زیاد نیست نمی‌توانند به خوبی پیچیدگی داده‌ها را درک کنند.

نکته دیگر این است که در فصل شبکه‌های عصبی مرجع اصلی این پروژه نیز گفته شده بود معمولاً یک لایه شبکه عصبی برای درک پیچیدگی داده‌ها مناسب است. من ادعا نمی‌کنم شبکه عصبی نوشته شده پیچیدگی داده‌ها را درک کرده است. چرا که شاهد

```
1 #Model training based on random hyper
  parameters
2 MSE_val_nn<-c()
3 HL<-list()
4 LR<-c()
5 activation_function<-c('logistic','tanh')
6 algorithm<-c('backprop','rprop+', 'rprop-',
  'sag','slr')
7 AF<-c()
8 AL<-c()
9 for(i in 1:300){
10   print(i)
11   LR[i]<-runif(1,0,1)
12   AF[i]<-floor(runif(1,1,3))
13   AL[i]<-floor(runif(1,1,6))
14   hidden_layer_num<-floor(runif(1,1,6))
15   hidden_layer<-floor(runif(hidden_layer_
    num,3,16))
16   hidden_layer<-hidden_layer[order(hidden_
    layer,decreasing = TRUE)]
17   nn<-neuralnet(Price~. ,data=train.data,
18     hidden = hidden_layer,
19     linear.output = F,
20     lifesign = 'full',
21     rep=1,
22     threshold=0.05,
23     learningrate=LR[i],
```

⁷⁴Deep Learning

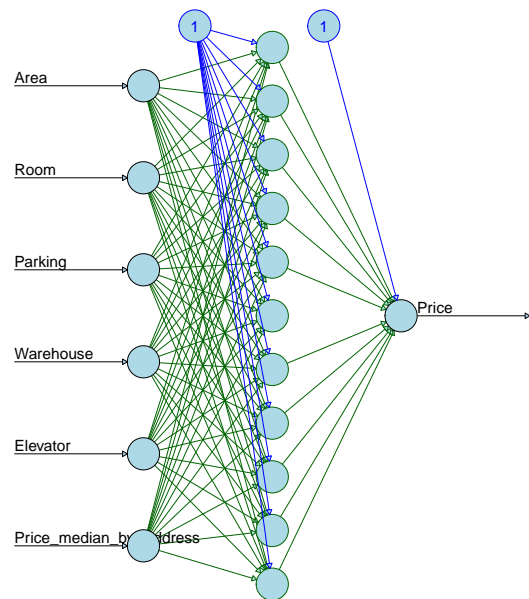
⁷⁵Neural Network

⁷⁶Activation Function

⁷⁷Python

⁷⁸Hidden Layer

⁷⁹Learning Rate

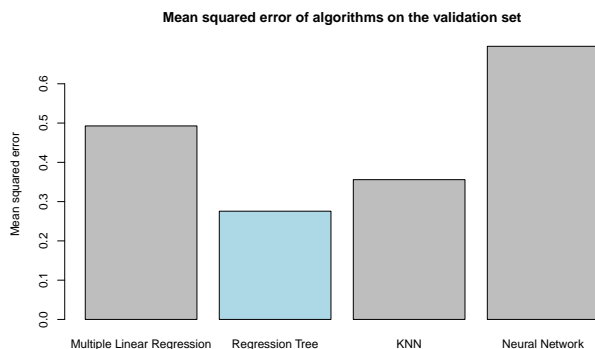


شکل ۲۸: دیاگرام شبکه عصبی منتخب

دقت مناسبی از مدل نهایی شبکه عصبی نیستیم. اما با توجه به جستجوی تصادفی انجام شده برای من جالب و حقیقتاً شگفت انگیز بود که شبکه عصبی یک لایه‌ای از بین ۳۰ مدل مختلف شبکه عصبی به عنوان بهترین مدل انتخاب شده بود. به شخصه به هیچ وجه این پیش‌بینی را نمی‌کردم.

انتخاب مدل نهایی

”یک تصویر به هزار کلمه می‌ارزد“ این جمله برگرفته از مرجع اصلی است. من از یک نمودار میله‌ای^{۸۰} برای تصویری سازی و مقایسه میزان میانگین توان دوم خطاهای مدل‌های منتخب در هر فصل پروژه کمک گرفتم. به‌وضوح بهترین دقت مربوط به الگوریتم درخت رگرسیون بوده است.



شکل ۲۹: مقایسه میانگین توان دوم خطاهای مدل‌ها

حال وقت آن رسیده است که دقت مدل منتخب نهایی را روی مجموعه تست بررسی کنیم. مبنای ما در این پروژه میانگین توان دوم خطاها است اما میانگین قدرمطلق خطاها^{۸۱} را نیز برای تجسم بهتر مخاطب محاسبه کرده‌ام.

```
1 > MSE(predict(best_dt_model, newdata = test.
  data), test.data[, 7])
2 [1] 0.1228406
3 > MAE(predict(best_dt_model, newdata = test.
  data), test.data[, 7])
4 [1] 0.1797792
```

نتیجه از تصور شخصی من بهتر بود. میانگین توان دوم خطاهای ما روی مجموعه تست حتی از مجموعه اعتبارسنجی هم بهتر بود. در واقع میانگین توان دوم خطاها در مجموعه تست نهایی کمتر از نصف مجموعه اعتبارسنجی شد. لذا خطای نهایی مدل پیاده‌سازی شده (بر اساس تابع خطای انتخاب شده) برابر ۰/۱۷۹۷۷۹۲ است.

پیاده‌سازی روی داده‌های واقعی خارج از مجموعه داده‌ها

تا اینجا پروژه ما به سرانجام رسیده است اما من دوست دارم دو عدد ثبت آزمایشی را که خارج از مجموعه داده‌هاست را به کمک مدل منتخب قیمت‌گذاری کنم. دو ثبت من به شرح زیر است

۱. خانه‌ای ۷۵ متری ۲ سال ساخت در منطقه نیاوران دارای ۲ خواب، مجهز به پارکینگ انباری و آسانسور که در تاریخ مذکور حدود ۷ میلیارد و ۵۰۰ میلیون تومان قیمت‌گذاری شده بود.

۲. خانه‌ای ۸۲ متری در منطقه شهرزیا، با قدمت ۱۰ سال در تاریخ مذکور، دارای ۲ خواب، مجهز به پارکینگ انباری و آسانسور که ۲ میلیارد و ۸۰۰ میلیون تومان در تاریخ مذکور قیمت‌گذاری شده بود.

این مجموعه کوچک داده‌ها با توجه به نرمال‌سازی انجام شده روی داده‌های مجموعه آموزشی، نرمالیده شده و مدل نهایی روی آنها اعمال شد. مورد اول به مبلغ حدوداً ۷ میلیارد تومان، و خانه دوم به مبلغ حدوداً ۳ میلیارد تومان توسط مدل نهایی پیشگویی شد.

```
1 > MSE(predict(best_dt_model, newdata = kh), kh
  [, 7])
2 [1] 0.0034291
```

همچنین خطای این مدل روی داده‌های جدید نسبتاً پایین بوده است.

فرایند انتخاب مدل

فرایند انتخاب مدل من بر اساس صفحات ۱۵۰ و ۱۵۱ کتاب Un-derstanding Machine Learning from theory to algorithms

نوشته شای شالو شوارتز و شای بن داوید بود. [۱۸] با توجه به اینکه در آموزش دادن همه مدل‌ها از مجموعه اعتبارسنجی بهره بردیم، لذا اعتبارسنجی نهایی را روی مجموعه تست انجام دادم. خوشبختانه همان‌طور که مشخص است دچار بیش‌برازش نشدم. ضمناً این کتاب برای مواقعی که دقت ما روی مجموعه تست به شکل مشهودی کمتر از مجموعه اعتبارسنجی باشد راهکارهای خوب و مفصلی ارائه می‌دهد.

⁸⁰Bar Chart

⁸¹Mean Absolute Error

جدول کدگذاری متغیر رسته‌ای

Label	Address String	Encoded by target
۱	Other	۰۳۲۲۲۲۸۲۰۲
۲	Punak	۰۲۹۵۵۸۸۱۷۸
۳	Pardis	۰۵۸۵۲۷۱۷۷
۴	West Ferdows Boulevard	۰۳۴۹۷۰۴۸۹۳
۵	Gheitarieh	۰۴۶۰۱۳۵۸۲۹
۶	Shahran	۰۲۷۹۶۷۱۴۹۷
۷	Saadat Abad	۰۴۱۴۲۹۵۷۸۸
۸	Parand	۰۶۱۹۶۵۱۸۰۱
۹	Shahr-e-Ziba	۰۳۰۲۵۹۱۵۱۸
۱۰	Southern Janatabad	۰۲۸۹۳۲۱۵۰۶
۱۱	Central Janatabad	۰۲۲۲۳۷۱۴۴۶
۱۲	Jeyhoon	۰۴۷۷۰۳۸۳۴
۱۳	Persian Gulf Martyrs Lake	۰۲۵۹۲۹۸۱۴۶
۱۴	Andisheh	۰۵۶۱۸۷۴۲۴۹
۱۵	Ostad Moein	۰۴۶۶۲۱۴۹۹۷
۱۶	East Ferdows Boulevard	۰۳۱۱۵۰۴۸۵۹
۱۷	Shahrake Qods	۰۵۶۰۴۴۱۷۴۸
۱۸	Niavaran	۱۷۸۵۹۹۵۳۴۵
۱۹	Pasdaran	۰۵۹۲۵۶۲۶۱۴
۲۰	Pirouzi	۰۴۴۵۲۰۴۹۷۸
۲۱	Salsabil	۰۵۲۷۹۷۱۷۱۹
۲۲	Shahrake Gharb	۰۵۸۹۸۸۸۶۱۱
۲۳	Farmanieh	۱۶۷۷۷۶۱۹۱۵
۲۴	Heravi	۰۲۲۸۳۸۸۹۵۶
۲۵	Ekhtiarieh	۰۲۷۴۲۲۸۹۹۷
۲۶	Islamshahr	۰۵۲۱۶۰۵۰۴۷
۲۷	Feiz Garden	۰۲۷۰۱۲۱۴۸۹
۲۸	Yousef Abad	۰۱۱۵۰۶۲۱۸۸
۲۹	Northern Janatabad	۰۲۱۶۰۰۴۷۷۴
۳۰	Qasr-od-Dasht	۰۵۰۵۶۸۸۳۶۶
۳۱	North Program Organization	۰۱۳۴۵۱۱۳۶۸
۳۲	Zaferanieh	۱۶۷۱۷۱۳۵۷۶
۳۳	Aqdasieh	۱۱۲۷۳۶۳۰۹۱
۳۴	Beryanak	۰۵۴۰۷۰۵۰۶۴
۳۵	Narmak	۰۲۷۹۶۷۱۴۹۷
۳۶	Pakdasht	۰۵۵۷۲۵۸۴۱۲
۳۷	Azarbaijan	۰۵۰۸۸۷۱۷۰۲
۳۸		۰۳۰۸۳۲۱۵۲۳
۳۹	Abazar	۰۵۰۹۵۰۳۷۹۹۹
۴۰	Damavand	۰۴۵۷۹۳۸۳۲۳
۴۱	Si Metri Ji	۰۴۶۴۳۰۴۹۹۵
۴۲	Southern Program Organization	۰۳۳۳۱۵۱۵۴۵
۴۳	Tenant	۰۴۹۹۳۲۱۶۹۳
۴۴	Marzadaran	۰۴۱۴۲۹۵۷۸۸
۴۵	Velenjak	۱۷۳۲۱۹۶۹۶۴
۴۶	Karoon	۰۵۲۷۹۷۱۷۱۹
۴۷	Jordan	۰۶۶۱۶۴۱۰۰۹
۴۸	Elahieh	۰۹۶۵۰۱۲۹۴۶
۴۹	Golestan	۰۵۳۷۵۲۱۷۲۷
۵۰	Kahrizak	۰۵۵۹۰۰۹۲۴۷
۵۱	Northern Chitgar	۰۵۵۶۶۲۱۷۴۵
۵۲	Mirdamad	۰۲۸۰۴۰۴۶۶۹
۵۳	Amirabad	۰۱۲۰۵۰۴۶۸۹
۵۴	Kamranieh	۲۵۶۶۲۳۱۰۴۱
۵۵	Northren Jamalzadeh	۰۳۳۶۹۷۱۵۴۹
۵۶	Dorous	۰۶۰۵۲۹۵۹۵۹
۵۷	Hashemi	۰۵۴۵۷۹۸۴۰۲
۵۸	Shahryar	۰۵۵۹۸۰۵۰۸۱
۵۹	Amirich	۰۴۴۰۶۸۴۶۴۱
۶۰	Sattarkhan	۰۴۴۱۰۴۶۲۱
۶۱	Komeil	۰۵۲۹۲۴۵۰۵۳
۶۲	Qalandari	۰۲۲۷۱۱۵۶۲۱
۶۳	Qazvin Imamzadeh Hassan	۰۴۷۳۴۰۹۳۳۷
۶۴	Railway	۰۳۰۳۸۶۸۵۲
۶۵	Rudhen	۰۴۶۸۷۶۱۶۶۶
۶۶	West Pars	۰۱۷۹۳۶۵۵۷۹
۶۷	Air force	۰۳۴۹۷۰۴۸۹۳
۶۸	Gholhak	۰۲۹۳۳۲۹۰۱۴
۶۹	Ozgol	۰۳۱۳۷۱۲۷۶
۷۰	Zafar	۰۱۳۱۹۵۴۳

تقدیر و تشکر

بدین وسیله در وهله اول از استاد محترم درس، جناب آقای دکتر فقیهی بابت آموختن درس داده‌کاوی، تعریف پروژه و همه راهنمایی‌هایشان در طول پروژه تشکر و قدردانی می‌نمایم. باشد تا برداشت صحیحی از آنچه ایشان به ما آموختند داشته باشم. همچنین از دوستان خوبم آقایان دکتر بهراد تقی بیگلر و سهراب فریدی که من را در به سرانجام رساندن این پروژه راهنمایی کردند نهایت سپاسگزاری را دارم.

- [16] D. A. Pierce and D. W. Schafer, "Residuals in generalized linear models," *Journal of the American Statistical Association*, vol.81, no.396, pp.977–986, 1986.
- [17] Y. A. LeCun, L. Bottou, G. B. Orr, and K.-R. Müller, *Efficient BackProp*, pp.9–48. Berlin, Heidelberg: Springer Berlin Heidelberg, 2012.
- [18] S. Shalev-Shwartz and S. Ben-David. *Understanding Machine Learning - From Theory to Algorithms*. Cambridge University Press, 2014.
- [1] گالیت شمولی، پیترسی بروس، اینبال یاهاو، نیتین آر پاتل، کنتسی داده‌کاوی برای تحلیل خودکار کسب‌وکار: مفاهیم، فنون. لیختندال و کاربردهای R.
- [2] I. M. Society, "واژه نامه ریاضی انجمن ریاضی ایران," September 2022.
- [3] S. Research and T. Center, "واژه نامه و اصطلاحات انجمن آمار," September 2022.
- [4] M. Kariminejad., "House price," December 2021.
- [5] M. Kabari, "مشاوران املاک، تنظیم گر بازار یا حباب ساز قیمت؟ مسکن؟," September 2021.
- [6] F. E. Grubbs, "Procedures for detecting outlying observations in samples," *Technometrics*, vol.11, no.1, pp.1–21, 1969.
- [7] geetansh044., "How to normalize data in r?," December 2021.
- [8] P. Cerda, G. Varoquaux, and B. Kégl, "Similarity encoding for learning with dirty categorical variables," *Machine Learning*, vol.107, no.8, pp.1477–1494, 2018.
- [9] W. N. Venables and B. D. Ripley. *Modern Applied Statistics with S*. New York: Springer, fourth ed. , 2002. ISBN 0-387-95457-0.
- [10] چرا بازار مسکن منطقه ۵ تهران داغ ترین بازار. وبسایت املاک‌باشی. www.amlakbashi.com, 1399.
- [11] wikipedia., "Heatmap," December 2021.
- [12] S. T. Roweis and L. K. Saul, "Nonlinear Dimensionality Reduction by Locally Linear Embedding," *Science*, vol.290, no.5500, pp.2323–2326, 2000.
- [13] P. Pudil and J. Novovičová, *Novel Methods for Feature Subset Selection with Respect to Problem Knowledge*, pp.101–116. Boston, MA: Springer US, 1998.
- [14] wikipedia., "Principal component analysis," May 2022.
- [15] K. Jain., "Linkedin profile," May 2020.

F		A
Feature ویژگی		Activation Function تابع فعال ساز
Forward پیشرو		Area مساحت
Frequency فراوانی		
G		B
Grid Search جستجوی شبکه‌ای		Backward پسرو
		Chart Bar نمودار میله‌ای
H		Binary دودویی
Heatmap نمودار حرارتی		Boxplot نمودار جعبه‌ای
Hidden Layer لایه پنهان		
Histogram بافت‌نگار		C
Hyperparameter فرا پارامتر		Categorical رسته‌ای
Hyperparameter Optimization بهینه‌سازی فرا پارامتر		Comment توضیح
		Compiler مترجم
I		Complexity Parameter پیچیدگی پارامتر
Imbalance نامتعادل		Coordinate مختصات
Integer عدد صحیح		Correlation همبستگی
K		D
K-Nearest Neighbor K نزدیکترین همسایه		Data داده
		Data Leakage نشت داده
L		Data Reduction داده کاهی
Learning Rate نرخ یادگیری		Data Visualization تصویرسازی داده‌ها
		Decision Tree درخت تصمیم
M		Deep Learning یادگیری عمیق
Machine Learning یادگیری ماشین		Dimension Reduction کاهش بُعد
Matrix ماتریس		Dispersion پراکندگی
Mean میانگین		Distribution توزیع
Mean Absolute Error میانگین قدرمطلق خطاها		Dummy Variable متغیر ظاهری

R

Random	تصادفی
Real Number	عدد حقیقی
Record	ثبت
Regression	رگرسیون
Regression Tree	درخت رگرسیون
Residual	مانده
Response Variable	متغیر پاسخ

S

Scatter Diagram	نمودار پراکنش
Side-By-Side Boxplot	نمودار جعبه‌ای پهلو به پهلو
Stepwise	گام به گام
String	رشته
Summary	خلاصه
Supervised	راهنماییده

T

Test Set	مجموعه آزمون
Training Set	مجموعه آموزشی
Tuning	تنظیم

V

Validation Set	مجموعه اعتبار سنجی
Variance	واریانس
Visualization	تصویری سازی

Mean Squared Error	میانگین توان دوم خطاها
Median	میان
Missing Value	مقدار گم شده
Multiple Linear Regression	رگرسیون خطی چندگانه

N

Neural Network	شبکه عصبی
Normal Distribution	توزیع نرمال
Normalized	نرمالیده
Null	تهی

O

Outlier	دور افتاده
---------	------------

P

Partition	افراز
Percentile	صدک
Positive Skewness	چولگی مثبت
Predictor	پیشگو
Preprocessing	پیش پردازش
Principal Component Analysis	تجزیه و تحلیل مؤلفه اصلی
P-Value	پی-مقدار
Python	پایتون

Q

QQ Plot	نمودار چندک چندک
Quartile	چارک

واژه‌نامه فارسی به انگلیسی

Frequency	فراوانی
Hidden Layer	لایه پنهان
Matrix	ماتریس
Compiler	مترجم
Response Variable	متغیر پاسخ
Dummy Variable	متغیر ظاهری
Test Set	مجموعه آزمون
Training Set	مجموعه آموزشی
Validation Set	مجموعه اعتبار سنجی
Coordinate	مختصات
Area	مساحت
Missing Value	مقدار گم شده
Median	میانه
Mean	میانگین
Mean Squared Error	میانگین توان دوم خطاها
Mean Absolute Error	میانگین قدر مطلق خطاها
Imbalance	نامتعادل
Learning Rate	نرخ یادگیری
Normalized	نرمالیده
Data Leakage	نشت داده
Boxplot	نمودار جعبه‌ای
Side-By-Side Boxplot	نمودار جعبه‌ای پهلو به پهلو
Heatmap	نمودار حرارتی
Scatter Diagram	نمودار پراکنش
QQ Plot	نمودار چندک چندک
Correlation	همبستگی
Variance	واریانس
Feature	ویژگی
Python	پایتون
Dispersion	پراکندگی
Backward	پسرو
Forward	پیشرو
Predictor	پیشگو

Partition	افراز
Histogram	بافت‌نگار
Hyperparameter Optimization	بهینه‌سازی فرا پارامتر
Activation Function	تابع فعال ساز
Principal Component Analysis	تجزیه و تحلیل مؤلفه اصلی
Random	تصادفی
Data Visualization	تصویرسازی داده‌ها
Visualization	تصویری سازی
Tuning	تنظیم
Null	تهی
Distribution	توزیع
Normal Distribution	توزیع نرمال
Comment	توضیح
Record	ثبت
Grid Search	جستجوی شبکه‌ای
Summary	خلاصه
Data	داده
Data Reduction	داده کاهی
Decision Tree	درخت تصمیم
Regression Tree	درخت رگرسیون
Binary	دودویی
Outlier	دورافتاده
Supervised	راهنماییده
Categorical	رسته‌ای
String	رشته
Regression	رگرسیون
Multiple Linear Regression	رگرسیون خطی چندگانه
Neural Network	شبکه عصبی
Percentile	صدک
Real Number	عدد حقیقی
Integer	عدد صحیح
Hyperparameter	فرا پارامتر

Preprocessing	پیش پردازش
P-Value	پی-مقدار
Complexity Parameter	پیچیدگی پارامتر
Deep Learning	یادگیری عمیق
Machine Learning	یادگیری ماشین
K-Nearest Neighbor	K نزدیکترین همسایه
Quartile	چارک
Positive Skewness	چولگی مثبت
Dimension Reduction	کاهش بُعد