

Proyecto 2

Dataset

Para el dataset de entrenamiento y testing se utilizó house_rent_result90, el cual se dividió en un **80%** para el entrenamiento del modelo y un **20%** para el testeo del mismo.

Description: df [6 × 13]

	X	Posted...	BHK	Rent	Size	Floor	Area.Ty...
	<int>	<chr>	<int>	<int>	<int>	<chr>	<chr>
1	1	2022-0...	2	10000	1100	Ground ...	Super A...
2	2	2022-0...	2	20000	800	1 out of 3	Super A...
3	3	2022-0...	2	17000	1000	1 out of 3	Super A...
4	4	2022-0...	2	10000	800	1 out of 2	Super A...
5	5	2022-0...	2	7500	850	1 out of 2	Carpet ...
6	6	2022-0...	2	7000	600	Ground ...	Super A...

6 rows | 1-8 of 13 columns

Feature Engineering

1. Limpieza de Campos Nulos

Para que nuestros datos fueran correctos y precisos se inició limpiando los datos nulos de cada una de las columnas del dataset las cuales tuvieran dichos datos.

Limpieza de nulos

```
{r}
colsNulls <- colnames(Odataset)[!complete.cases(t(Odataset))]
colsNulls

character(0)
```

2. Cambios de datos categóricos

Se agrupará por categoría y se obtendrá el promedio de la renta por categoría.

• Furnishing.Status

Furnishing.Status	MeanFurnish
<chr>	<dbl>
Furnished	21274.797
Semi-Furnished	3131.599
Unfurnished	-11813.013

3 rows

• Floor

A tibble: 460 × 2

Floor	MeanFloor
<chr>	<dbl>
24 out of 24	700000.00
7 out of 20	600000.00
36 out of 81	380000.00
18 out of 20	353500.00
19 out of 85	350000.00
39 out of 60	350000.00
45 out of 60	350000.00
20 out of 41	330000.00
8 out of 27	330000.00
19 out of 33	310000.00

- Area.Type

A tibble: 3 × 2

Area.Type <chr>	MeanAreaT <dbl>
Carpet Area	51414.02
Super Area	18436.71
Built Area	10500.00

- Area.Locality

A tibble: 2,098 × 2

Area.Locality <chr>	MeanAreaL <dbl>
Marathahalli	715780.00
Lady Ratan Tower, Worli	700000.00
Bandra East	600000.00
Vettuvankeni	600000.00
Altamount Road	500000.00
Rustomjee Elements, Andheri West	400000.00
World One Tower Mumbai, Worli	380000.00
Deonar	350000.00
Green Park	350000.00
Indiabulls Blu, Worli Naka Acharaya Atre Chowk	350000.00
Lodha World Crest, Lower Parel	350000.00
Sundar Nagar	350000.00
Anand Niketan	330000.00

- City

A tibble: 6 × 2

City <chr>	MeanCity <dbl>
Mumbai	81103.28
Delhi	30230.35
Bangalore	25593.93
Chennai	21568.14
Hyderabad	20442.06
Kolkata	11599.94

- Tenant.Preferred

A tibble: 3 × 2

Tenant.Preferred <chr>	MeanTenant <dbl>
Family	46744.58
Bachelors	43689.01
Bachelors/Family	30616.38

3 rows

- Point.of.Contact

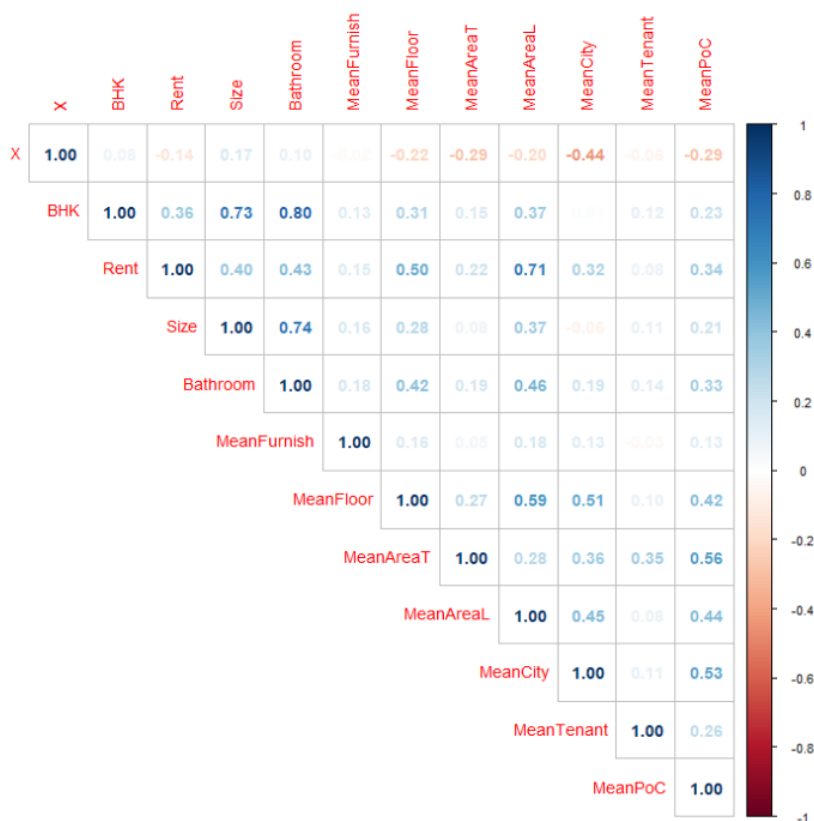
A tibble: 3 × 2

Point.of.Contact <chr>	MeanPoC <dbl>
Contact Agent	72098.91
Contact Owner	16630.84
Contact Builder	5500.00

3 rows

Se reemplazan las columnas del dataset original por las columnas con las medias calculadas.

La variable con la mayor correlación con la variable **Rent** es: **MeanAreaL**



Pruebas y Experimentos realizados

- Experimento 1

```
{r}  
experimento1 <- dstrain %>%  
  lm(formula = Rent~MeanAreaL)# Pruebas, realizar operaciones entre las columnas  
  del dataset  
  
yhat <- predict(experimento1, dtest)  
  
rmseExp1 <- rmse(yhat,dtest$Rent)  
rmseExp1  
EXP1<- data.frame("Rownum"= dtest$X, "RentP" = yhat, "Rent"=dtest$Rent)  
EXP1
```

RMSE = 19790.35

- Experimento 2

Se agrega la columna MeanFloor

```
experimento2 <- dstrain %>%  
  lm(formula = Rent~MeanAreaL+MeanFloor)# Pruebas, realizar operaciones entre las  
  columnas del dataset  
  
yhat <- predict(experimento2, dtest)  
  
rmseExp2 <- rmse(yhat,dtest$Rent)  
rmseExp2  
EXP2<- data.frame("Rownum"= dtest$X, "RentP" = yhat, "Rent"=dtest$Rent)  
EXP2
```

RMSE = 19210.67

- Experimento 3

Se le agregan todas las columnas

```
experimento3 <- dstrain %>%  
  lm(formula = Rent~MeanAreaL+MeanFloor+MeanFurnish+BHK+Size+Bathroom+MeanAreaT  
  +MeanCity+MeanTenant+MeanPoC)# Pruebas, realizar operaciones entre las columnas del  
  dataset  
  
yhat <- predict(experimento3, dtest)  
  
rmseExp3 <- rmse(yhat,dtest$Rent)  
rmseExp3  
EXP3<- data.frame("Rownum"= dtest$X, "RentP" = yhat, "Rent"=dtest$Rent)  
EXP3
```

RMSE = 17370.65

- Experimento 4

Se le quita la columna MeanPoC

```
experimento4 <- dstrain %>%  
  lm(formula = Rent~MeanAreaL+MeanFloor+MeanFurnish+BHK+Size+Bathroom+MeanAreaT  
  +MeanCity+MeanTenant)# Pruebas, realizar operaciones entre las columnas del dataset  
  
yhat <- predict(experimento4, dtest)  
  
rmseExp4 <- rmse(yhat,dtest$Rent)  
rmseExp4  
EXP4<- data.frame("Rownum"= dtest$X, "RentP" = yhat, "Rent"=dtest$Rent)  
EXP4
```

RMSE = 17376.38

- Experimento 5

Se le quita la columna MeanTenant

```
experimento5 <- dstrain %>%  
  lm(formula = Rent~MeanAreaL+MeanFloor+MeanFurnish+BHK+Size+Bathroom+MeanAreaT  
+MeanCity)# Pruebas, realizar operaciones entre las columnas del dataset  
  
yhat <- predict(experimento5, dtest)  
  
rmseExp5 <- rmse(yhat,dtest$Rent)  
rmseExp5  
EXP5<- data.frame("Rownum"= dtest$X, "RentP" = yhat, "Rent"=dtest$Rent)  
EXP5
```

RMSE = 17345.21

- Experimento 6

Se le agrega la columna MeanPoc

```
experimento6 <- dstrain %>%  
  lm(formula = Rent~MeanAreaL+MeanFloor+MeanFurnish+BHK+Size+Bathroom+MeanAreaT  
+MeanCity+MeanPoC)# Pruebas, realizar operaciones entre las columnas del dataset  
  
yhat <- predict(experimento6, dtest)  
  
rmseExp6 <- rmse(yhat,dtest$Rent)  
rmseExp6  
EXP6<- data.frame("Rownum"= dtest$X, "RentP" = yhat, "Rent"=dtest$Rent)  
EXP6
```

RMSE = 17339.99

- Experimento 7

Se elimina MeanCity

```
experimento7 <- dstrain %>%  
  lm(formula = Rent~MeanAreaL+MeanFloor+MeanFurnish+BHK+Size+Bathroom+MeanAreaT  
+MeanPoC)# Pruebas, realizar operaciones entre las columnas del dataset  
  
yhat <- predict(experimento7, dtest)  
  
rmseExp7 <- rmse(yhat,dtest$Rent)  
rmseExp7  
EXP7<- data.frame("Rownum"= dtest$X, "RentP" = yhat, "Rent"=dtest$Rent)  
EXP7
```

RMSE = 17329.79

- Experimento 8

Se elimina MeanAreaT

```
experimento8 <- dstrain %>%  
  lm(formula = Rent~MeanAreaL+MeanFloor+MeanFurnish+BHK+Size+Bathroom+MeanPoC)#  
Pruebas, realizar operaciones entre las columnas del dataset  
  
yhat <- predict(experimento8, dtest)  
  
rmseExp8 <- rmse(yhat,dtest$Rent)  
rmseExp8  
EXP8<- data.frame("Rownum"= dtest$X, "RentP" = yhat, "Rent"=dtest$Rent)  
EXP8
```

RMSE = 17339.58

- Experimento 9

Se agrega MeanAreaT y se elimina Bathroom

```
experimento9 <- dstrain %>%  
  lm(formula = Rent~MeanAreaL+MeanFloor+MeanFurnish+BHK+Size+MeanAreaT+MeanPoC)#  
Pruebas, realizar operaciones entre las columnas del dataset  
  
yhat <- predict(experimento9, dtest)  
  
rmseExp9 <- rmse(yhat,dtest$Rent)  
rmseExp9  
EXP9<- data.frame("Rownum"= dtest$X, "RentP" = yhat, "Rent"=dtest$Rent)  
EXP9
```

RMSE = 17354.42

- Experimento 10

Se agrega Bathroom y se elimina Size

```
experimento10 <- dstrain %>%  
  lm(formula = Rent~MeanAreaL+MeanFloor+MeanFurnish+BHK+Bathroom+MeanAreaT+MeanPoC  
)# Pruebas, realizar operaciones entre las columnas del dataset  
  
yhat <- predict(experimento10, dtest)  
  
rmseExp10 <- rmse(yhat,dtest$Rent)  
rmseExp10  
EXP10<- data.frame("Rownum"= dtest$X, "RentP" = yhat, "Rent"=dtest$Rent)  
EXP10
```

RMSE = 17403.06

- Experimento 11

Se agrega Size y se elimina BHK

```
experimento11 <- dstrain %>%  
  lm(formula = Rent~MeanAreaL+MeanFloor+MeanFurnish+Size+Bathroom+MeanAreaT+MeanPoC  
)# Pruebas, realizar operaciones entre las columnas del dataset  
  
yhat <- predict(experimento11, dtest)  
  
rmseExp11 <- rmse(yhat,dtest$Rent)  
rmseExp11  
EXP11<- data.frame("Rownum"= dtest$X, "RentP" = yhat, "Rent"=dtest$Rent)  
EXP11
```

RMSE = 17321.37

- Experimento 12

Se elimina MeanFurnish

```
experimento12 <- dstrain %>%  
  lm(formula = Rent~MeanAreaL+MeanFloor+Size+Bathroom+MeanAreaT+MeanPoC)# Pruebas,  
  realizar operaciones entre las columnas del dataset  
  
yhat <- predict(experimento12, dtest)  
  
rmseExp12 <- rmse(yhat,dtest$Rent)  
rmseExp12  
EXP12<- data.frame("Rownum"= dtest$X, "RentP" = yhat, "Rent"=dtest$Rent)  
EXP12
```

RMSE = 17324.27

- Experimento 13

Se agrega MeanFurnish y se elimina MeanFloor

```
experimento13 <- dstrain %>%  
  lm(formula = Rent~MeanAreaL+MeanFurnish+Size+Bathroom+MeanAreaT+MeanPoC)# Pruebas  
, realizar operaciones entre las columnas del dataset  
  
yhat <- predict(experimento13, dtest)  
  
rmseExp13 <- rmse(yhat,dtest$Rent)  
rmseExp13  
EXP13<- data.frame("Rownum"= dtest$X, "RentP" = yhat, "Rent"=dtest$Rent)  
EXP13
```

RMSE = 17847.21

- Experimento 14

Se agrega MeanFloor y se elimina MeanAreaL

```
experimento14 <- dstrain %>%  
  lm(formula = Rent~MeanFloor+MeanFurnish+Size+Bathroom+MeanAreaT+MeanPoC)# Pruebas  
  , realizar operaciones entre las columnas del dataset  
  
yhat <- predict(experimento14, dtest)  
  
rmseExp14 <- rmse(yhat,dtest$Rent)  
rmseExp14  
EXP14<- data.frame("Rownum"= dtest$X, "RentP" = yhat, "Rent"=dtest$Rent)  
EXP14
```

RMSE = 31151.22

- Experimento 15

Se cambia Bathroom por MeanBathroom

```
experimento15 <- dstrain %>%  
  lm(formula = Rent~MeanAreaL+MeanFloor+MeanFurnish+Size+MeanBathroom+MeanAreaT  
  +MeanPoC)# Pruebas, realizar operaciones entre las columnas del dataset  
  
yhat <- predict(experimento15, dtest)  
  
rmseExp15 <- rmse(yhat,dtest$Rent)  
rmseExp15  
EXP15<- data.frame("Rownum"= dtest$X, "RentP" = yhat, "Rent"=dtest$Rent)  
EXP15
```

RMSE = 29654.54

- Experimento 16

```
experimento16 <- dstrain %>%  
  lm(formula = Rent~MeanAreaL+MeanFloor+MeanFurnish+Size+Bathroom+MeanAreaT)#  
  Pruebas, realizar operaciones entre las columnas del dataset  
  
yhat <- predict(experimento16, dtest)  
  
rmseExp16 <- rmse(yhat,dtest$Rent)  
rmseExp16  
EXP16<- data.frame("Rownum"= dtest$X, "RentP" = yhat, "Rent"=dtest$Rent)  
EXP16
```

RMSE = 17318.53

- Experimento 17

Debido a que se existen categorías en el archivo testing (el de 10%) que no estaban en el archivo de pruebas, se le agrego la media de la categoría para reemplazar los nulls que se agregaron.

```
#Experimento 17 - Se cambiaron los nulls de las categorías inexistentes del archivo testing por el promedio de las medias de los de esa categoría  
{r}  
experimento17 <- dstrain %>%  
  lm(formula = Rent~MeanAreaL+MeanFloor+MeanFurnish+Size+Bathroom+MeanAreaT)# Pruebas, realizar operaciones entre las columnas del dataset  
  
yhat <- predict(experimento17, dtest)  
  
rmseExp17 <- rmse(yhat,dtest$Rent)  
rmseExp17  
EXP17<- data.frame("Rownum"= dtest$X, "RentP" = yhat, "Rent"=dtest$Rent)  
EXP17
```

[1] 17318.53

RMSE = 17.318.53

```
```{r}
datasetPredict$MeanFurnish[is.na(datasetPredict$MeanFurnish)] <- FurnM
datasetPredict$MeanFloor[is.na(datasetPredict$MeanFloor)] <- FloorM
datasetPredict$MeanAreaT[is.na(datasetPredict$MeanAreaT)] <- AreaTM
datasetPredict$MeanAreaL[is.na(datasetPredict$MeanAreaL)] <- AreaLM
datasetPredict$MeanCity[is.na(datasetPredict$MeanCity)] <- CityM
datasetPredict$MeanTenant[is.na(datasetPredict$MeanTenant)] <- TenM
datasetPredict$MeanPoC[is.na(datasetPredict$MeanPoC)] <- PoCM
datasetPredict
```
```

Repositorio de Github: https://github.com/tjfv02/Proyecto2_Analisis-de-Datos.git