

Laboratorio 1

Oscar Godoy - Rafael Dubois

2022-07-24

Limpieza y exploración inicial

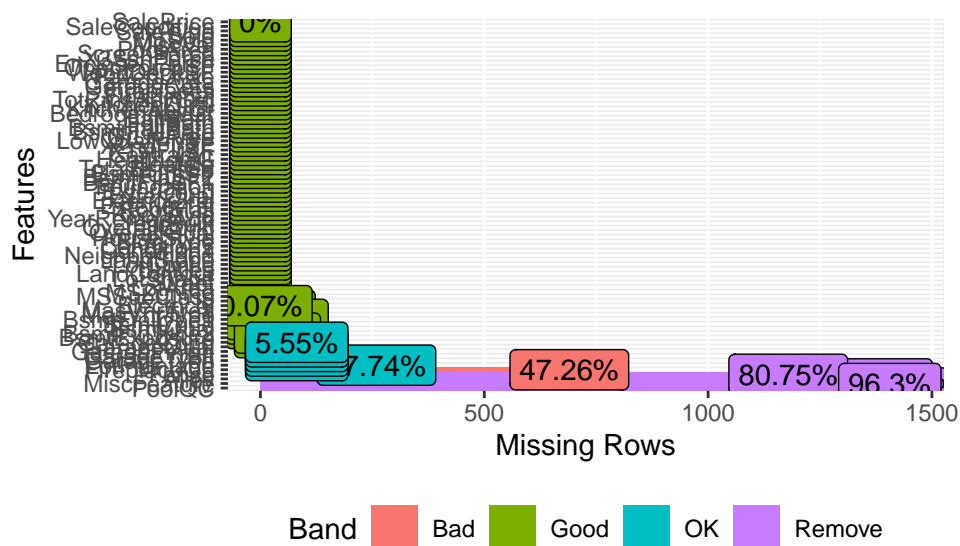
El data set contiene información sobre 1460 casas en la ciudad de Ames, Iowa. Para cada una de las casas, las cuales se identifican con un ID único, contamos con 79 variables exploratorias. Estas describen aspectos de la casa como zona, tamaño, forma, amenidades, entre otras. La variable objetivo, que se busca predecir, es el precio de venta de la casa.

```
introduce(dataOriginal)
```

```
##   rows columns discrete_columns continuous_columns all_missing_columns
## 1 1460      81           43           38              0
##   total_missing_values complete_rows total_observations memory_usage
## 1              6965           0             118260           755280
```

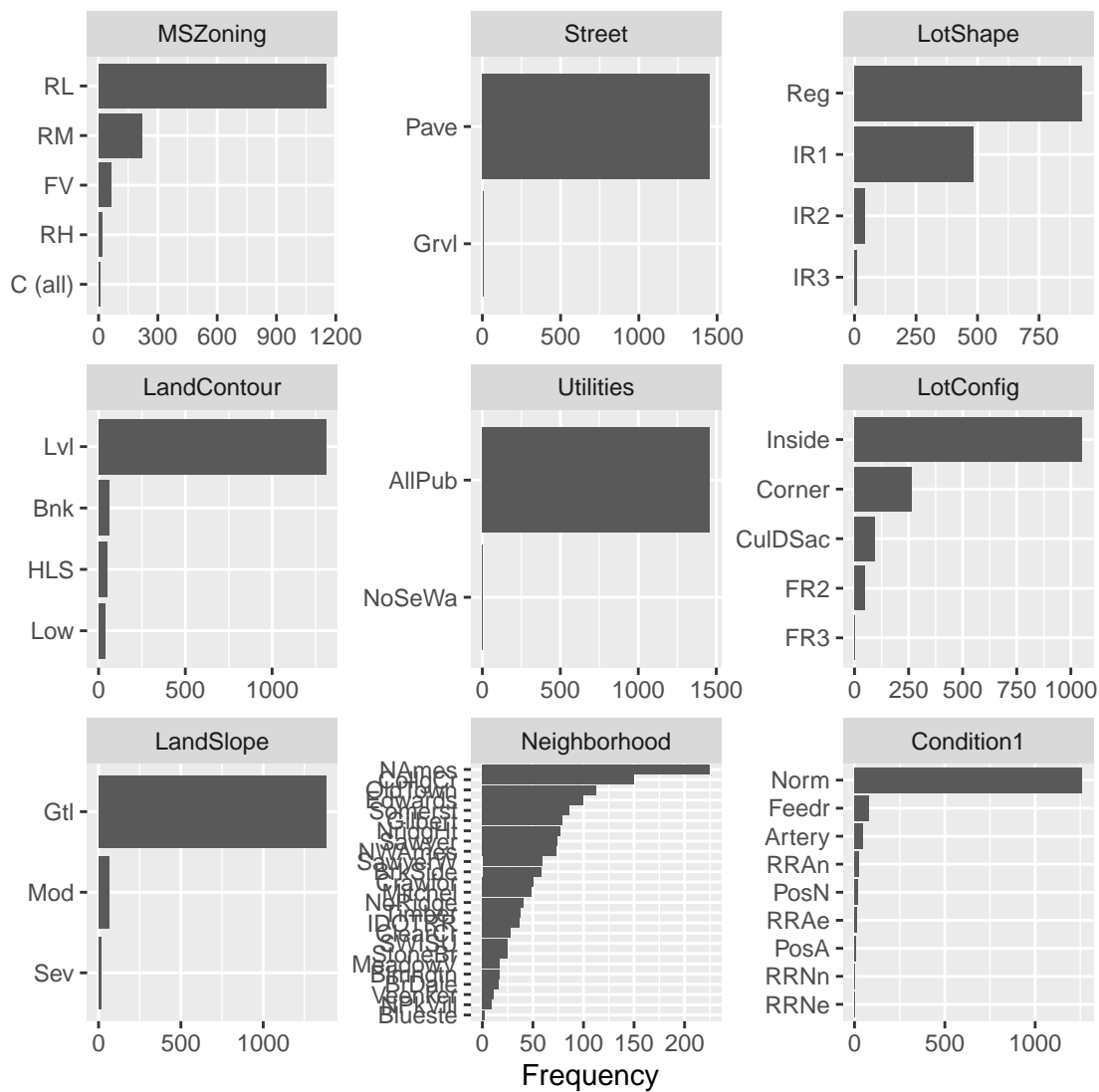
Procedemos a analizar la calidad de los datos. Vemos que hay variables con altos porcentajes de NA's, por lo que las descartamos. Estas variables tienen el factor en común de referirse a utilidades poco comunes que una casa pueda tener, como piscinas o cercas. La mayoría de las casas en el data set no posee estas características. En total, removemos las siguientes 5 variables:

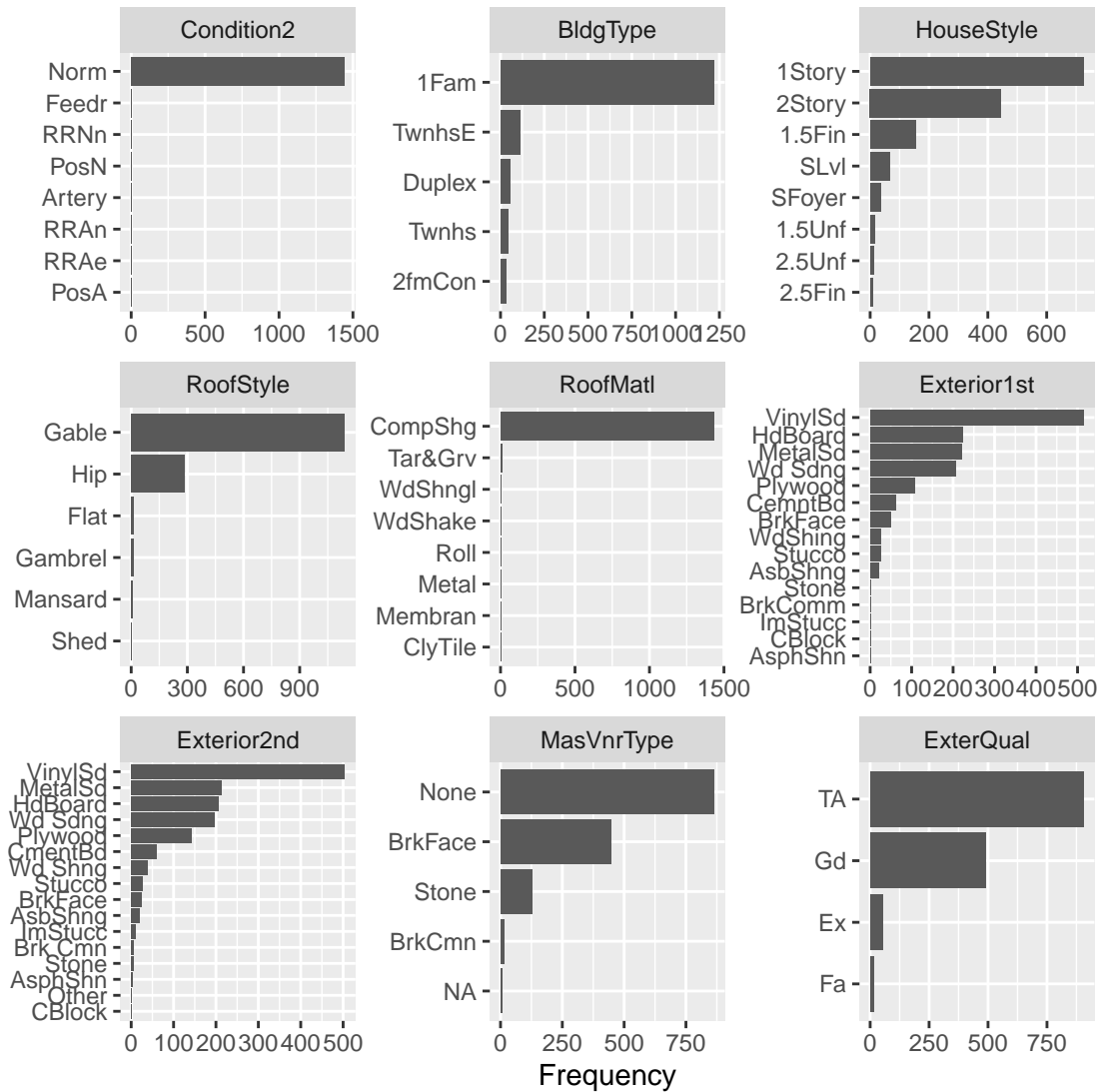
- PoolQC
- MiscFeature
- Alley
- Fence
- FireplaceQu

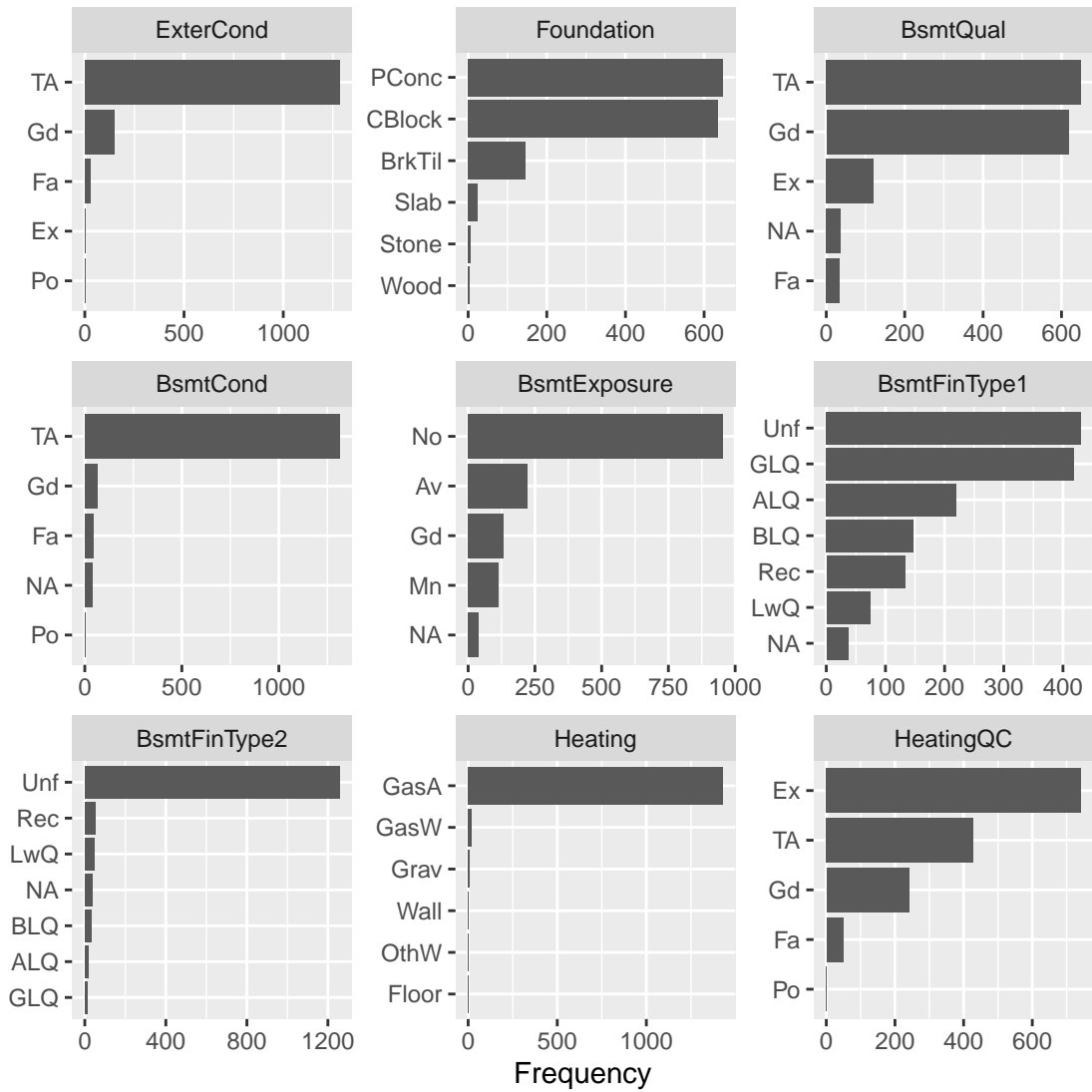


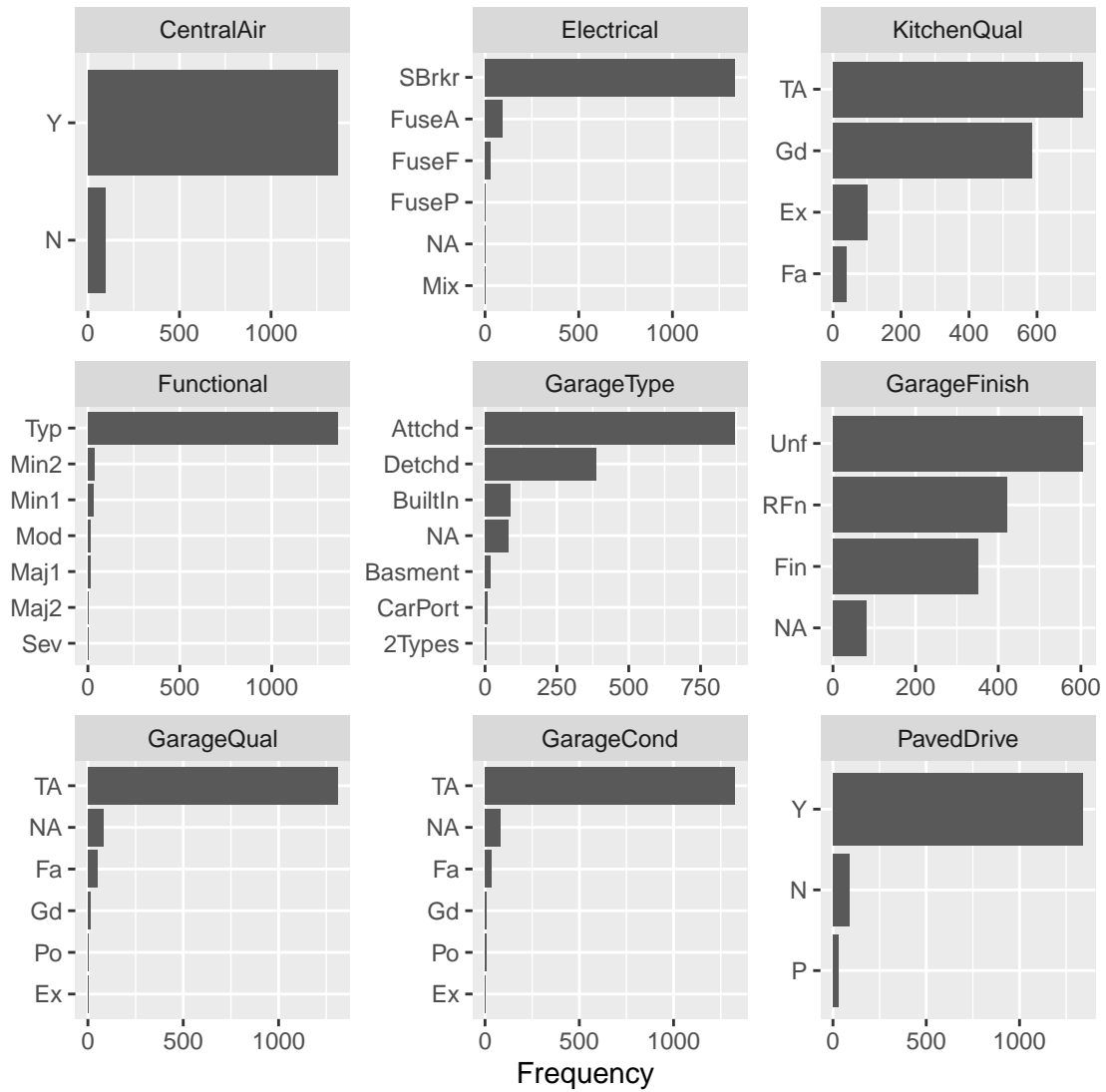
Ahora, analizamos la distribución de las variables cuantitativas y cualitativas. Descartamos las variables cualitativas que sean muy homogéneas en la población. Se elige descartar aquellas variables donde una categoría concentra más del 80% de las casas. Entonces, descartamos las siguientes 20 variables:

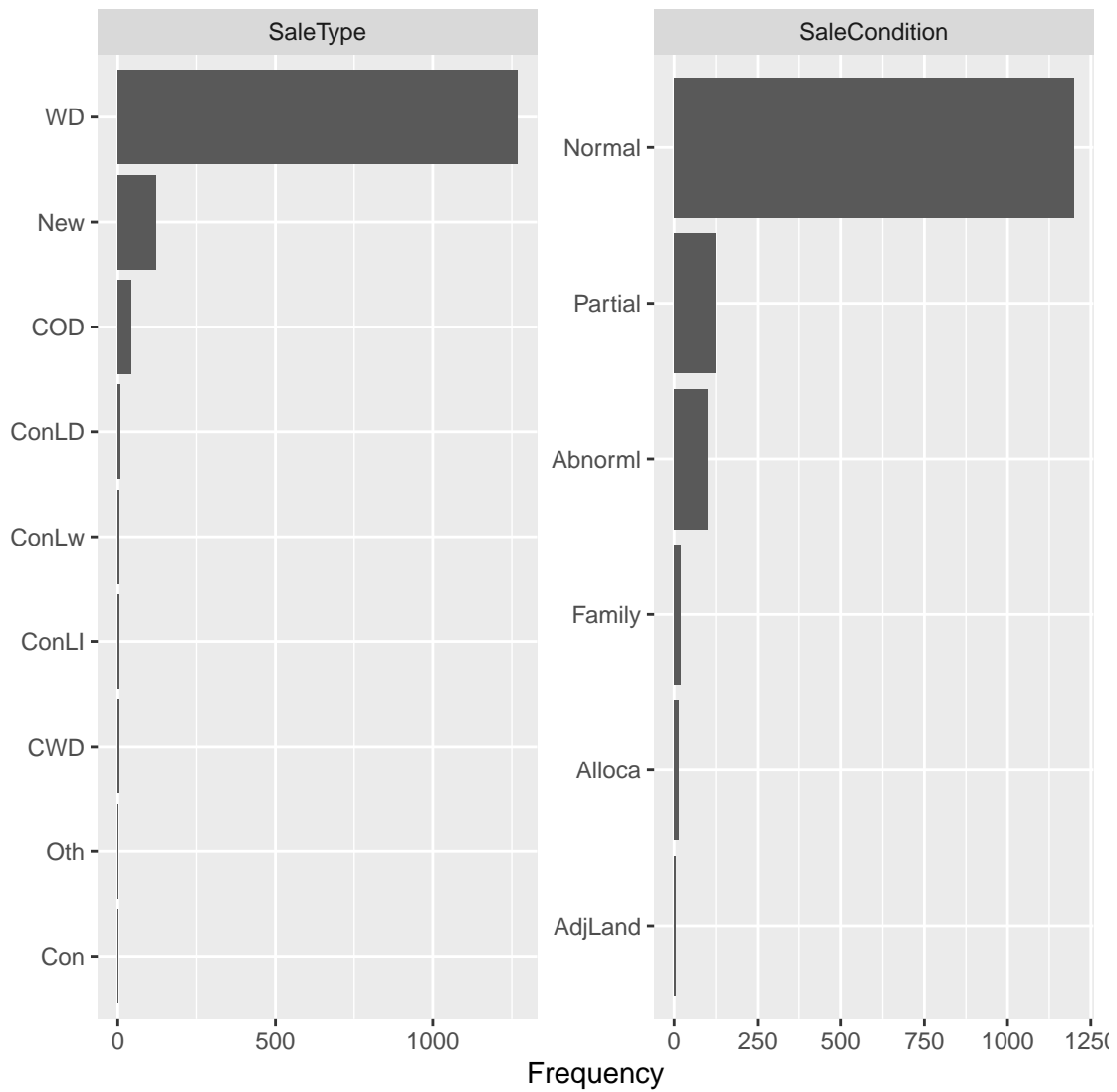
- Street
- LandContour
- Utilities
- LandSlope
- Condition1
- Condition2
- BldgType
- RoofMatl
- ExterCond
- BsmtCond
- BsmtFinType2
- Heating
- CentralAir
- Electrical
- Functional
- GarageQual
- GarageCond
- PavedDrive
- SaleType
- SaleCondition





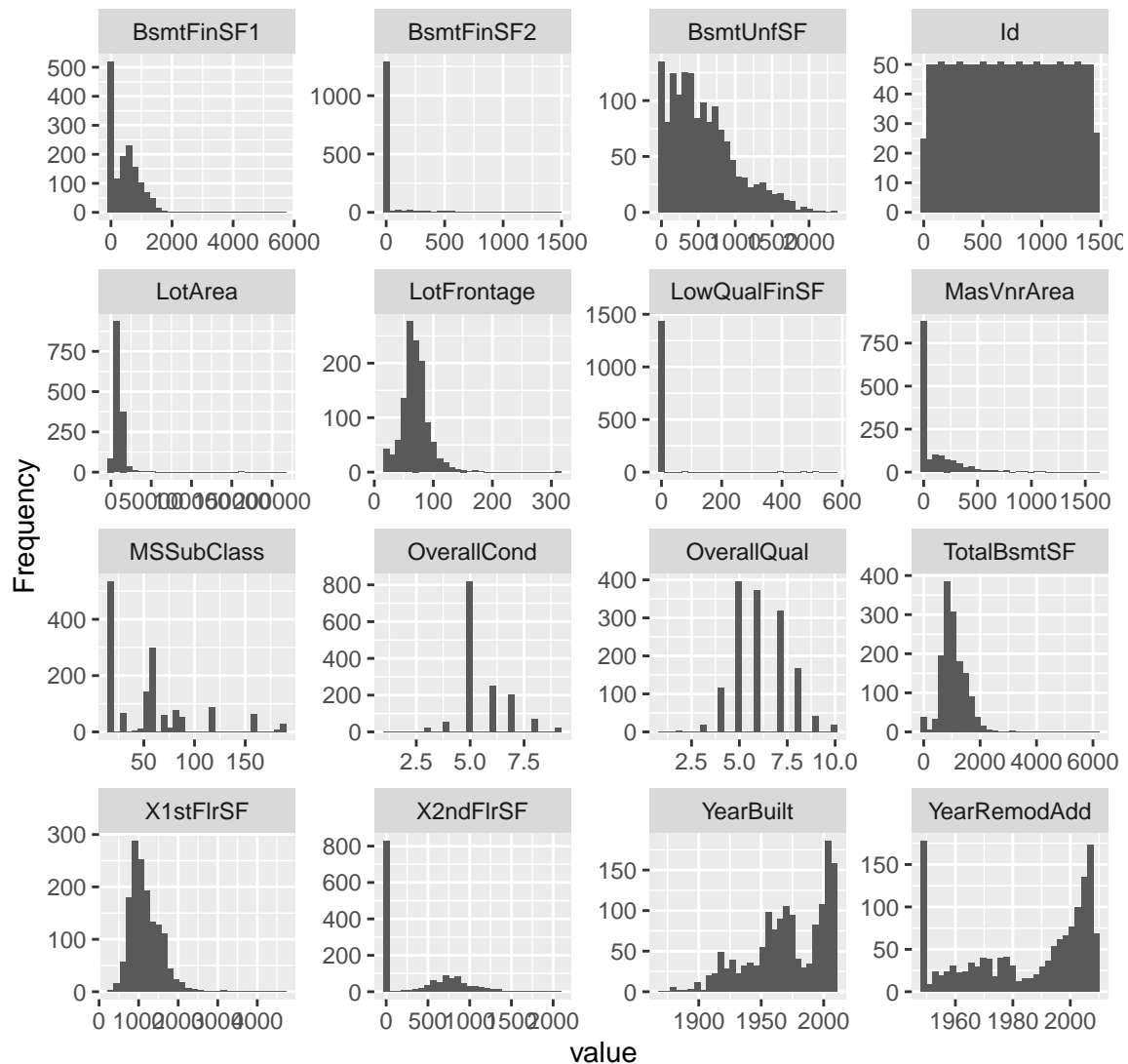


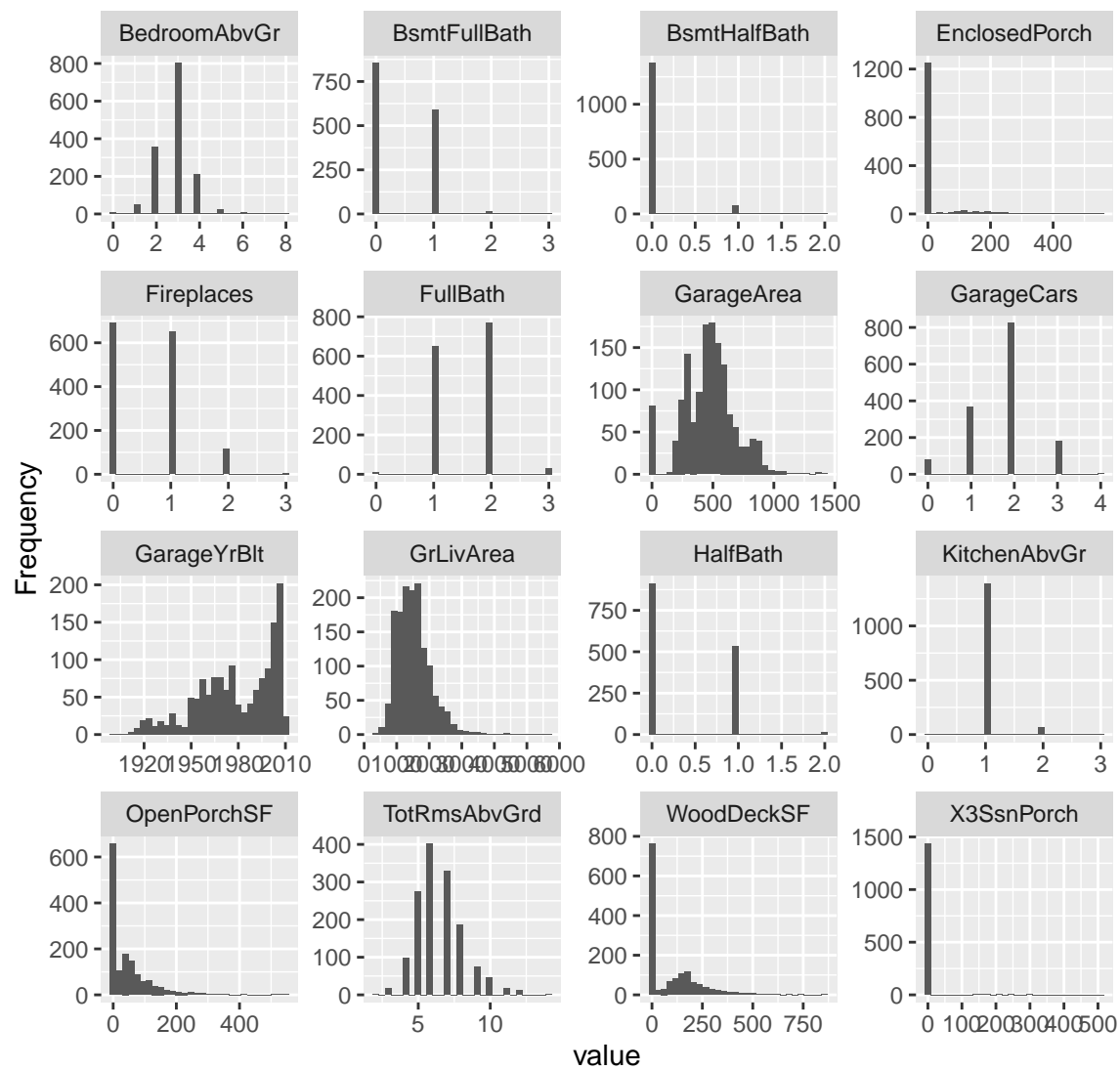


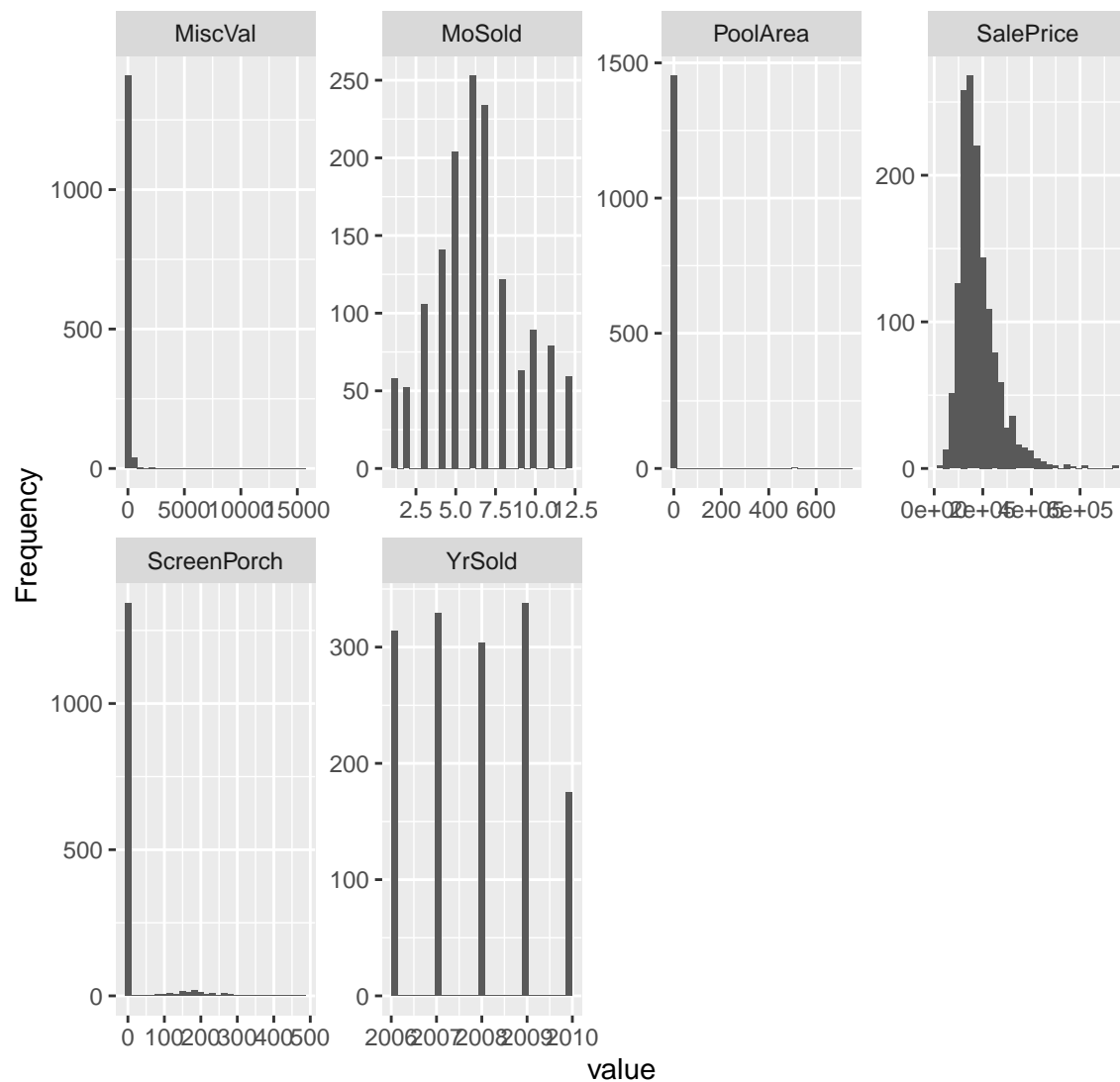


De igual manera, descartamos las variables cuantitativas que esten muy sesgadas. Descartamos aquellas variables donde al menos el 80% de los datos asuman un mismo valor. Encontramos 9 de estas variables:

- BsmtFinSF2
- LotArea
- LowQualFinSF
- BsmtHalfBath
- EnclosedPorch
- KitchenAbvGr
- X3SsnPorch
- MiscVal
- PoolArea
- ScreenPorch







Page 3

Finalmente, contamos con 46 variables cualitativas y cuantitativas que nos proponemos a investigar.

```
introduce(dataOriginal)
```

```
##  rows columns discrete_columns continuous_columns all_missing_columns
## 1 1460      46              18                28                0
##  total_missing_values complete_rows total_observations memory_usage
## 1                    630          1096             67160       391160
```

Análisis de Componentes Principales

Filtramos las variables numéricas para realizarles un PCA. Luego, de las variables restantes descartamos aquellas que no sean continuas. Variables discretas y ordinales se evitan. Además, eliminamos la variable SalePrice para solo analizar las variables exploratorias. Las variables descartadas son las siguientes:

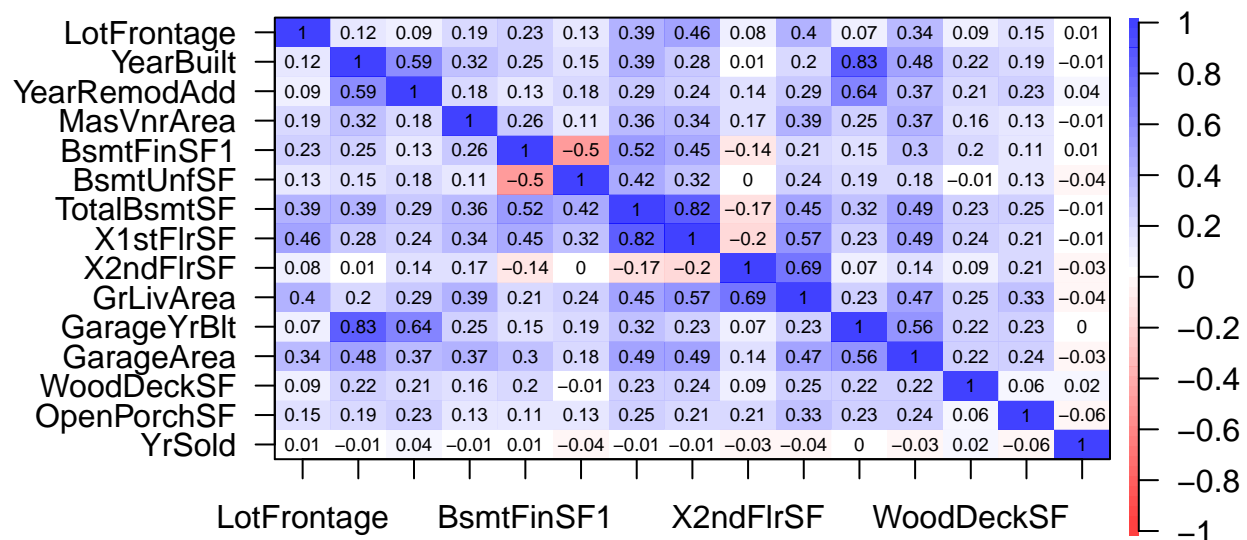
- Id
- MSSubClass
- OverallQual
- OverallCond
- BsmtFullBath
- FullBath
- HalfBath
- BedroomAbvGr
- TotRmsAbvGrd
- Fireplaces
- GarageCars
- MoSold

Después de la limpieza, quedamos con 16 variables continuas. Deseamos reducir la dimensionalidad de este data set empleando PCA. Primero, creamos la matriz de correlación de este set de datos, y notamos que su determinante es muy cercano a 0. Esto puede evidenciarse que hay bastante correlación entre las matrices, es decir, dependencia lineal entre las variables originales.

```
rcor <- cor(dataPCA, use = "pairwise.complete.obs")  
det(rcor)
```

```
## [1] 6.016125e-06
```

```
cor.plot(rcor)
```



Para asegurarnos de que PCA es aplicable, realizamos el test de esfericidad de Bartlett. El p-valor es prácticamente 0, por lo que descartamos la hipótesis nula del test y concluimos que la matriz de correlación de las variables es distinta a la identidad.

```
cortest.bartlett(dataPCA)
```

```
## R was not square, finding R from data
## $chisq
## [1] 17468.61
##
## $p.value
## [1] 0
##
## $df
## [1] 105
```

Con estos requisitos cumplidos, procedemos a realizar un PCA:

```
compPrinc <- prcomp(na.omit(dataPCA[,-1]), scale = T)
summary(compPrinc)
```

Importance of components:

##	PC1	PC2	PC3	PC4	PC5	PC6	PC7
## Standard deviation	2.1079	1.3313	1.2652	1.2097	1.00692	0.94551	0.91814
## Proportion of Variance	0.3174	0.1266	0.1143	0.1045	0.07242	0.06386	0.06021
## Cumulative Proportion	0.3174	0.4440	0.5583	0.6628	0.73526	0.79912	0.85933

##	PC8	PC9	PC10	PC11	PC12	PC13	PC14
## Standard deviation	0.80057	0.72468	0.61581	0.48451	0.38070	0.2048	0.04973
## Proportion of Variance	0.04578	0.03751	0.02709	0.01677	0.01035	0.0030	0.00018
## Cumulative Proportion	0.90511	0.94262	0.96971	0.98648	0.99683	0.9998	1.00000

Las herramientas de visualización de PCA nos fallaron a último momento, pero numéricamente podemos interpretar el tema de cada uno de los componentes. Resaltamos que las primeras 5 componentes explican el 75% de la variabilidad. Estos componentes pueden intepretarse como:

- Índice de área de amenidades y espacios secundarios
- Índice de área de espacios interiores
- Índice de antigüedad de la casa
- Índice de calidad de sótano
- Índice de año de venta

Reglas de asociación

Volvemos al dataset original para quitarle todo lo que quitamos para PCA, excepto algunas variables categóricas cuyas entradas parecen numéricas: * MSSubClass * OverallQual * OverallCond * MoSold

Entonces, el data frame para las reglas de asociación se ve de esta manera:

```
glimpse(dataARules)
```

```
## Rows: 1,460
## Columns: 21
## $ MSZoning      <chr> "RL", "RL", "RL", "RL", "RL", "RL", "RL", "RL", "RL", "RM", "RL~
## $ LotShape      <chr> "Reg", "Reg", "IR1", "IR1", "IR1", "IR1", "Reg", "IR1", "~
## $ LotConfig     <chr> "Inside", "FR2", "Inside", "Corner", "FR2", "Inside", "In~
## $ Neighborhood <chr> "CollgCr", "Veenker", "CollgCr", "Crawfor", "NoRidge", "M~
## $ HouseStyle    <chr> "2Story", "1Story", "2Story", "2Story", "2Story", "1.5Fin~
## $ Exterior1st   <chr> "VinylSd", "MetalSd", "VinylSd", "Wd Sdng", "VinylSd", "V~
## $ Exterior2nd   <chr> "VinylSd", "MetalSd", "VinylSd", "Wd Shng", "VinylSd", "V~
## $ MasVnrType    <chr> "BrkFace", "None", "BrkFace", "None", "BrkFace", "None", ~
## $ ExterQual     <chr> "Gd", "TA", "Gd", "TA", "Gd", "TA", "Gd", "TA", "TA", "TA~
## $ Foundation   <chr> "PConc", "CBlock", "PConc", "BrkTil", "PConc", "Wood", "P~
## $ BsmtQual      <chr> "Gd", "Gd", "Gd", "TA", "Gd", "Gd", "Ex", "Gd", "TA", "TA~
## $ BsmtExposure <chr> "No", "Gd", "Mn", "No", "Av", "No", "Av", "Mn", "No", "No~
## $ BsmtFinType1 <chr> "GLQ", "ALQ", "GLQ", "ALQ", "GLQ", "GLQ", "GLQ", "ALQ", "~
## $ HeatingQC     <chr> "Ex", "Ex", "Ex", "Gd", "Ex", "Ex", "Ex", "Ex", "Gd", "Ex~
## $ KitchenQual   <chr> "Gd", "TA", "Gd", "Gd", "Gd", "TA", "Gd", "TA", "TA", "TA~
## $ GarageType    <chr> "Attchd", "Attchd", "Attchd", "Detchd", "Attchd", "Attchd~
## $ GarageFinish  <chr> "RFn", "RFn", "RFn", "Unf", "RFn", "Unf", "RFn", "RFn", "~
## $ MSSubClass    <int> 60, 20, 60, 70, 60, 50, 20, 60, 50, 190, 20, 60, 20, 20, ~
## $ OverallQual   <int> 7, 6, 7, 7, 8, 5, 8, 7, 7, 5, 5, 9, 5, 7, 6, 7, 6, 4, 5, ~
## $ OverallCond   <int> 5, 8, 5, 5, 5, 5, 5, 6, 5, 6, 5, 5, 6, 5, 5, 8, 7, 5, 5, ~
## $ MoSold        <int> 2, 5, 9, 2, 12, 10, 8, 11, 4, 1, 2, 7, 9, 8, 5, 7, 3, 10, ~
```

Teniendo este nuevo set de datos, corremos el algoritmo de reglas de asociación. Usando un support de 48% y un confidence de 60%, se obtuvieron las reglas de asociación que se observan a continuación.

```
## Warning: Column(s) 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17,
## 18, 19, 20, 21 not logical or factor. Applying default discretization (see '?
## discretizeDF').

## Warning in discretize(x = c(60L, 20L, 60L, 70L, 60L, 50L, 20L, 60L, 50L, : The calculated breaks are
## Only unique breaks are used reducing the number of intervals. Look at ? discretize for details.

## Apriori
##
## Parameter specification:
## confidence minval smax arem aval originalSupport maxtime support minlen
##          0.6    0.1    1 none FALSE                TRUE         5    0.48      1
## maxlen target  ext
##          10 rules TRUE
##
## Algorithmic control:
## filter tree heap memopt load sort verbose
##       0.1 TRUE TRUE  FALSE TRUE     2    TRUE
##
## Absolute minimum support count: 700
##
## set item appearances ...[0 item(s)] done [0.00s].
```

```
## set transactions ...[135 item(s), 1460 transaction(s)] done [0.00s].
## sorting and recoding items ... [13 item(s)] done [0.00s].
## creating transaction tree ... done [0.00s].
## checking subsets of size 1 2 3 done [0.00s].
## writing ... [15 rule(s)] done [0.00s].
## creating S4 object ... done [0.00s].
```

```
inspect(reglas)
```

##	lhs	rhs	support	confidence	coverage
## [1]	{}	=> {ExterQual=TA}	0.6205479	0.6205479	1.0000000
## [2]	{}	=> {LotShape=Reg}	0.6335616	0.6335616	1.0000000
## [3]	{}	=> {BsmtExposure=No}	0.6527397	0.6527397	1.0000000
## [4]	{}	=> {LotConfig=Inside}	0.7205479	0.7205479	1.0000000
## [5]	{}	=> {MSZoning=RL}	0.7883562	0.7883562	1.0000000
## [6]	{GarageType=Attchd}	=> {MSZoning=RL}	0.5308219	0.8908046	0.5958904
## [7]	{MSZoning=RL}	=> {GarageType=Attchd}	0.5308219	0.6733275	0.7883562
## [8]	{LotShape=Reg}	=> {LotConfig=Inside}	0.5123288	0.8086486	0.6335616
## [9]	{LotConfig=Inside}	=> {LotShape=Reg}	0.5123288	0.7110266	0.7205479
## [10]	{BsmtExposure=No}	=> {LotConfig=Inside}	0.4801370	0.7355719	0.6527397
## [11]	{LotConfig=Inside}	=> {BsmtExposure=No}	0.4801370	0.6663498	0.7205479
## [12]	{BsmtExposure=No}	=> {MSZoning=RL}	0.4883562	0.7481637	0.6527397
## [13]	{MSZoning=RL}	=> {BsmtExposure=No}	0.4883562	0.6194613	0.7883562
## [14]	{LotConfig=Inside}	=> {MSZoning=RL}	0.5506849	0.7642586	0.7205479
## [15]	{MSZoning=RL}	=> {LotConfig=Inside}	0.5506849	0.6985230	0.7883562

##	lift	count
## [1]	1.0000000	906
## [2]	1.0000000	925
## [3]	1.0000000	953
## [4]	1.0000000	1052
## [5]	1.0000000	1151
## [6]	1.1299520	775
## [7]	1.1299520	775
## [8]	1.1222690	748
## [9]	1.1222690	748
## [10]	1.0208507	701
## [11]	1.0208507	701
## [12]	0.9490174	713
## [13]	0.9490174	713
## [14]	0.9694331	804
## [15]	0.9694331	804

Ignoramos las reglas triviales (con premisa nula). Vemos algunas reglas interesantes (todas hablan de probabilidades, no equivalencias exactas):

- Las primeras dos reglas nos dan una probable equivalencia entre las casas que poseen un garage pegado a ellas y las casas ubicadas en residenciales con baja densidad poblacional.
- Las siguientes dos reglas nos dicen que las casas cuyo terreno tiene una forma regular también mantienen todo su terreno en el interior de la propiedad.
- Las casas sin exposición exterior a su sótano son casas que mantienen su terreno en el interior.
- Las casas sin exposición exterior a su sótano son casas en residenciales con baja densidad poblacional.
- La quinta regla junta la 3 y 4. Las casas con su terreno completamente en el interior son casas en residenciales con baja densidad poblacional.

Conclusiones

La herramienta de reglas de asociación nos ayuda a crear relaciones entre variables importantes en el conjunto de datos, lo cual puede ser útil para saber qué tipo de variables esperamos ver cuando aparece otra, y con cuánta frecuencia. En nuestro caso, encontramos una alta relación entre algunas características de las casas, como la regularidad de la forma de un terreno y el interior del terreno de la casa, entre otras. Para nuestros propósitos de predicción, esto puede ayudar a agrupar distintos tipos de casas usando las reglas obtenidas para hacer las separaciones.