

Donaldo Sebastian Garcia Jiménez 19683  
Raul Angel Jimenez Hernandez 19017

Link al repo:

[https://github.com/raulangelj/Laboratorio1An-lisis\\_Exploratorio\\_PCA\\_aprior](https://github.com/raulangelj/Laboratorio1An-lisis_Exploratorio_PCA_aprior)

Link al drive:

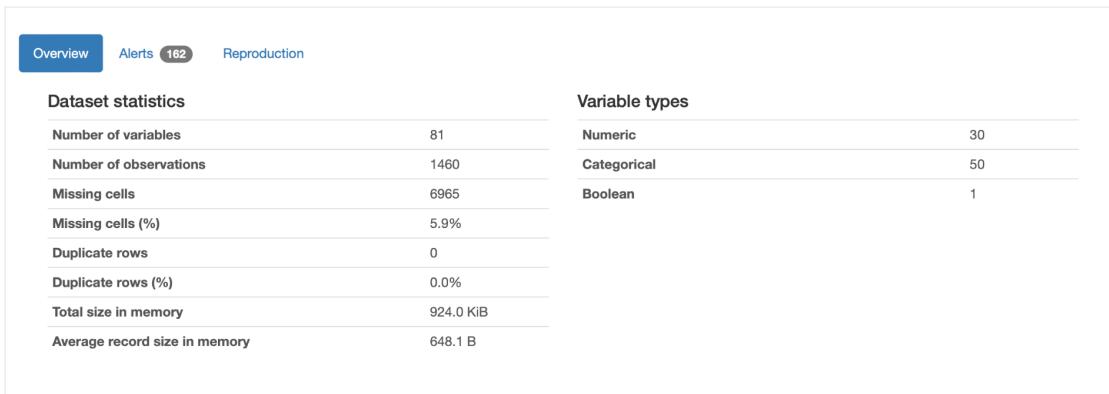
<https://docs.google.com/document/d/1v-qLmCEHHOxTXKwv5GC6-Z4n-UJdAz-wKulyPZhPbbc/edit?usp=sharing>

## Laboratorio 1

<b>Reporte de Datos</b>	<b>17</b>
Análisis exploratorio	17
Correlaciones	56
Hallazgos y explicación de procedimiento	57
Análisis de componentes principales	58
Análisis factorial	58
Bartlett	58
Hallazgos	59
Matriz de correlación	59
Coeficientes principales.	60
PCA	61
Reglas de asociación	62
Hallazgos y conclusiones	62
Rúbrica	62

## #1 Resumen dataset

### Overview



## #2 Variables a estudiar

### Cuantitativa Discreta

Id

Donaldo Sebastian Garcia Jiménez 19683  
Raul Angel Jimenez Hernandez 19017

MSSubClass  
LotFrontage  
LotArea  
YearBuilt  
YearRemodAdd  
MasVnrArea  
BsmtFinSF1  
BsmtFinSF2  
BsmtUnfSF  
TotalBsmtSF  
1stFlrSF  
2ndFlrSF  
LowQualFinSF  
GrLivArea  
BsmtFullBath  
BsmtHalfBath  
FullBath  
HalfBath  
BedroomAbvGr  
KitchenAbvGr  
TotRmsAbvGrd  
Fireplaces  
GarageYrBlt  
GarageCars  
GarageArea  
WoodDeckSF  
OpenPorchSF  
EnclosedPorch  
3SsnPorch  
ScreenPorch  
PoolArea  
MiscVal  
MoSold  
YrSold

### **Cuantitativa Continua**

### **Cualitativa o Categórica**

Alley  
Street  
LotShape  
LandContour  
Utilities  
LotConfig

Donaldo Sebastian Garcia Jiménez 19683  
Raul Angel Jimenez Hernandez 19017

LandSlope  
Neighborhood  
Condition1  
Condition2  
BldgType  
HouseStyle  
OverallQual  
OverallCond  
RoofStyle  
RoofMatl  
Exterior1st  
Exterior2nd  
MasVnrType  
ExterQual  
ExterCond  
Foundation  
BsmtQual  
BsmtCond  
BsmtExposure  
BsmtFinType1  
BsmtFinType2  
Heating  
HeatingQC  
CentralAir  
Electrical  
KitchenQual  
Functional  
FireplaceQu  
GarageType  
GarageFinish  
GarageQual  
GarageCond  
PavedDrive  
PoolQC  
Fence  
MiscFeature  
SaleType  
SaleCondition

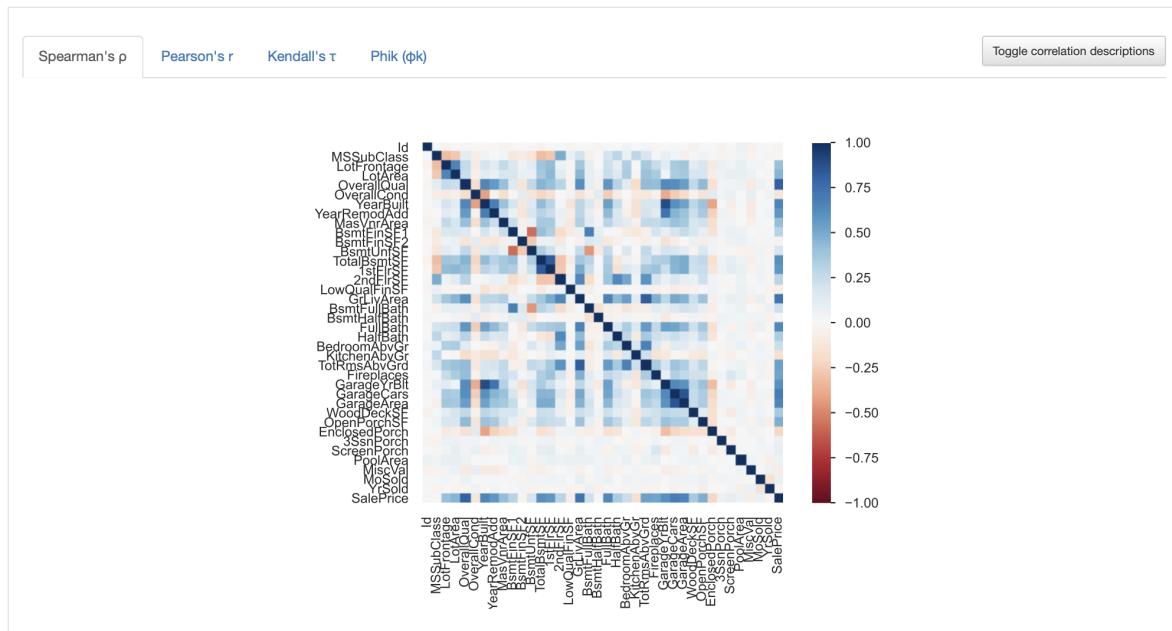
### #3 Variables a estudiar

Se pueden observar las gráficas de estos datos en el reporte más abajo.

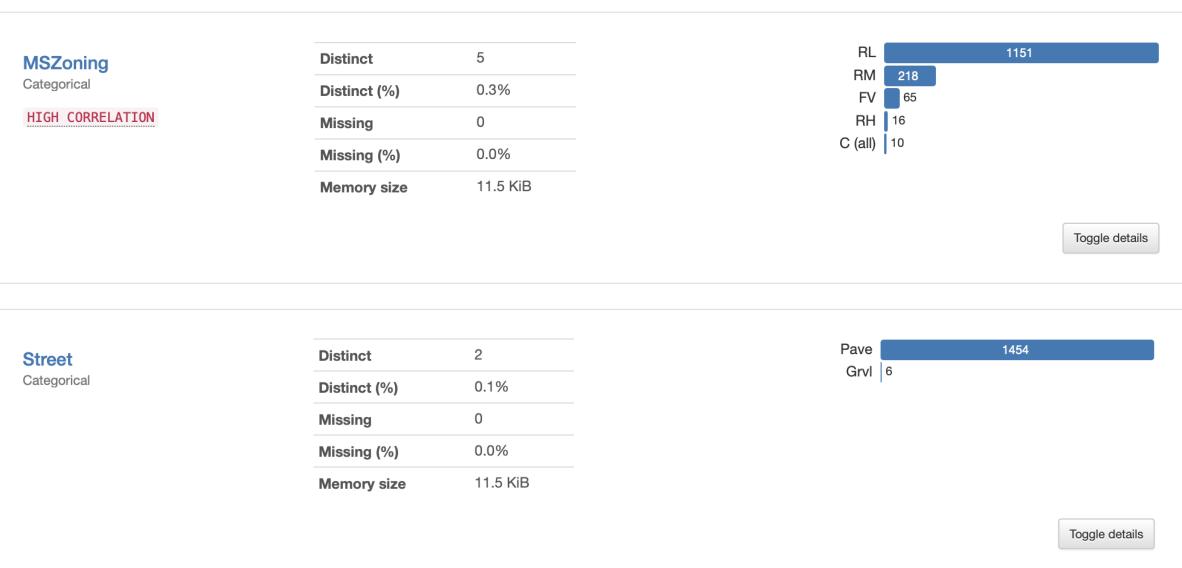
Donaldo Sebastian Garcia Jiménez 19683  
Raul Angel Jimenez Hernandez 19017

## #4 Análisis de correlación, variables numéricas categóricas

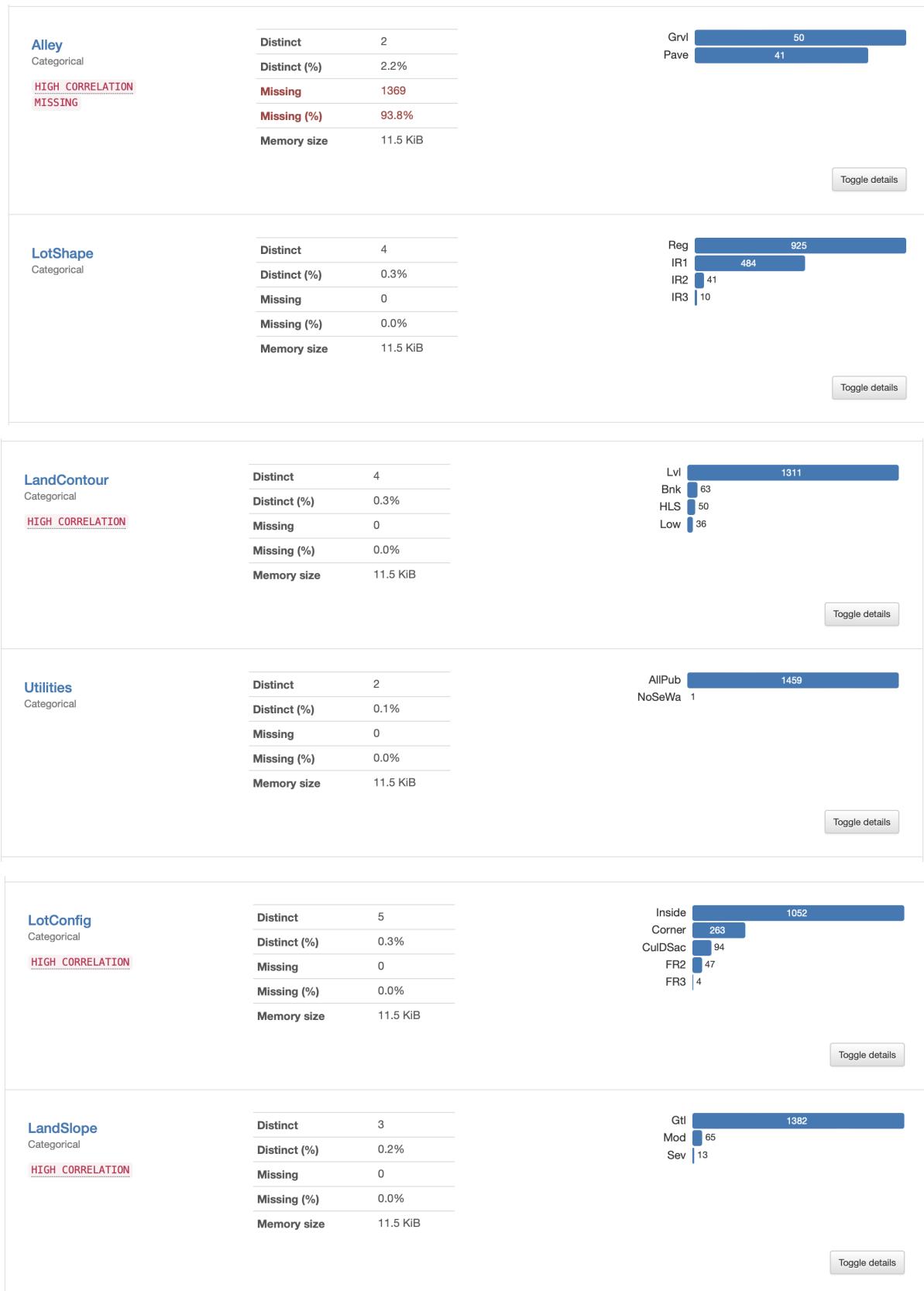
### Correlations



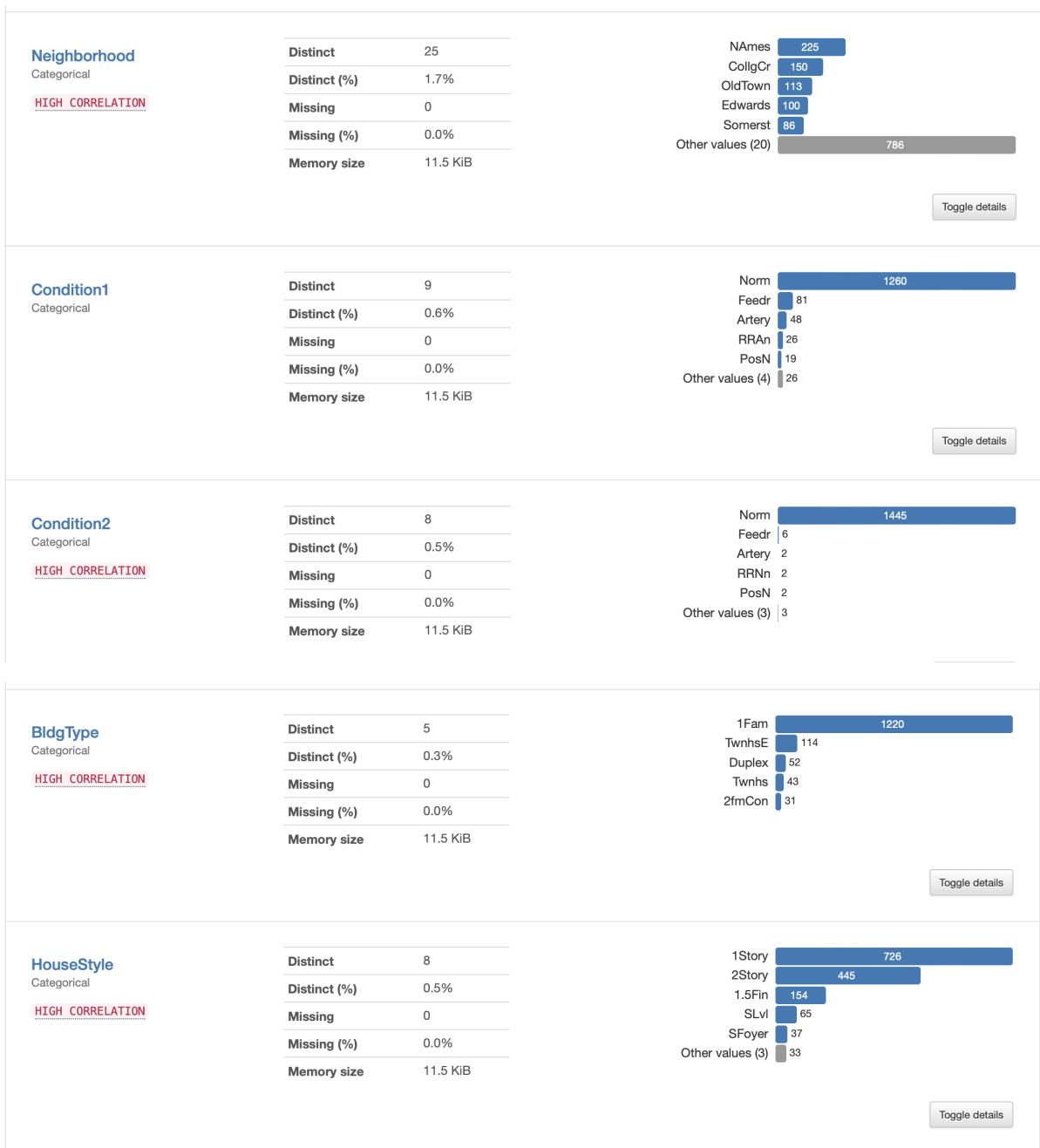
## #5 Variables categóricas análisis



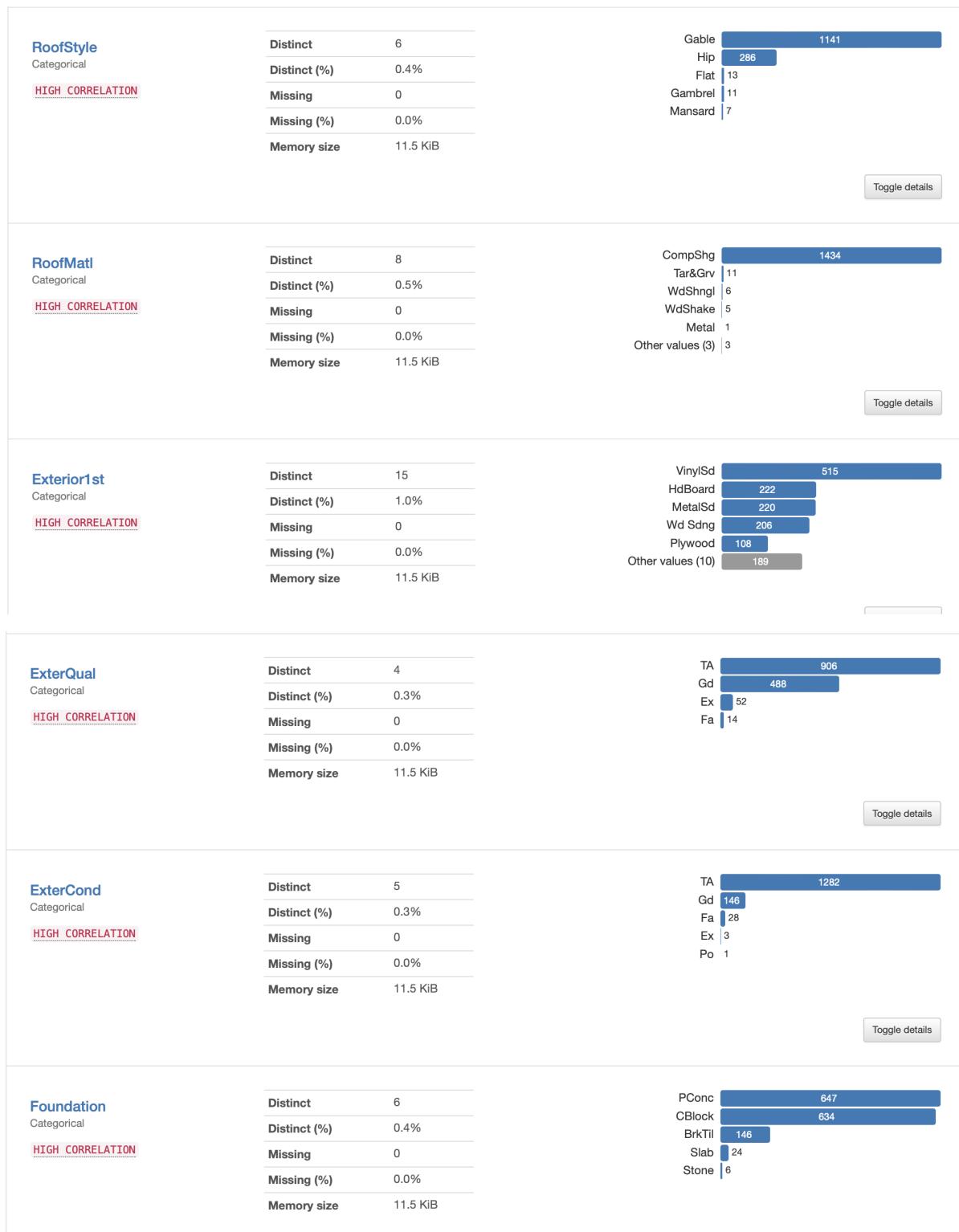
Donald Sebastian Garcia Jiménez 19683  
 Raul Angel Jimenez Hernandez 19017



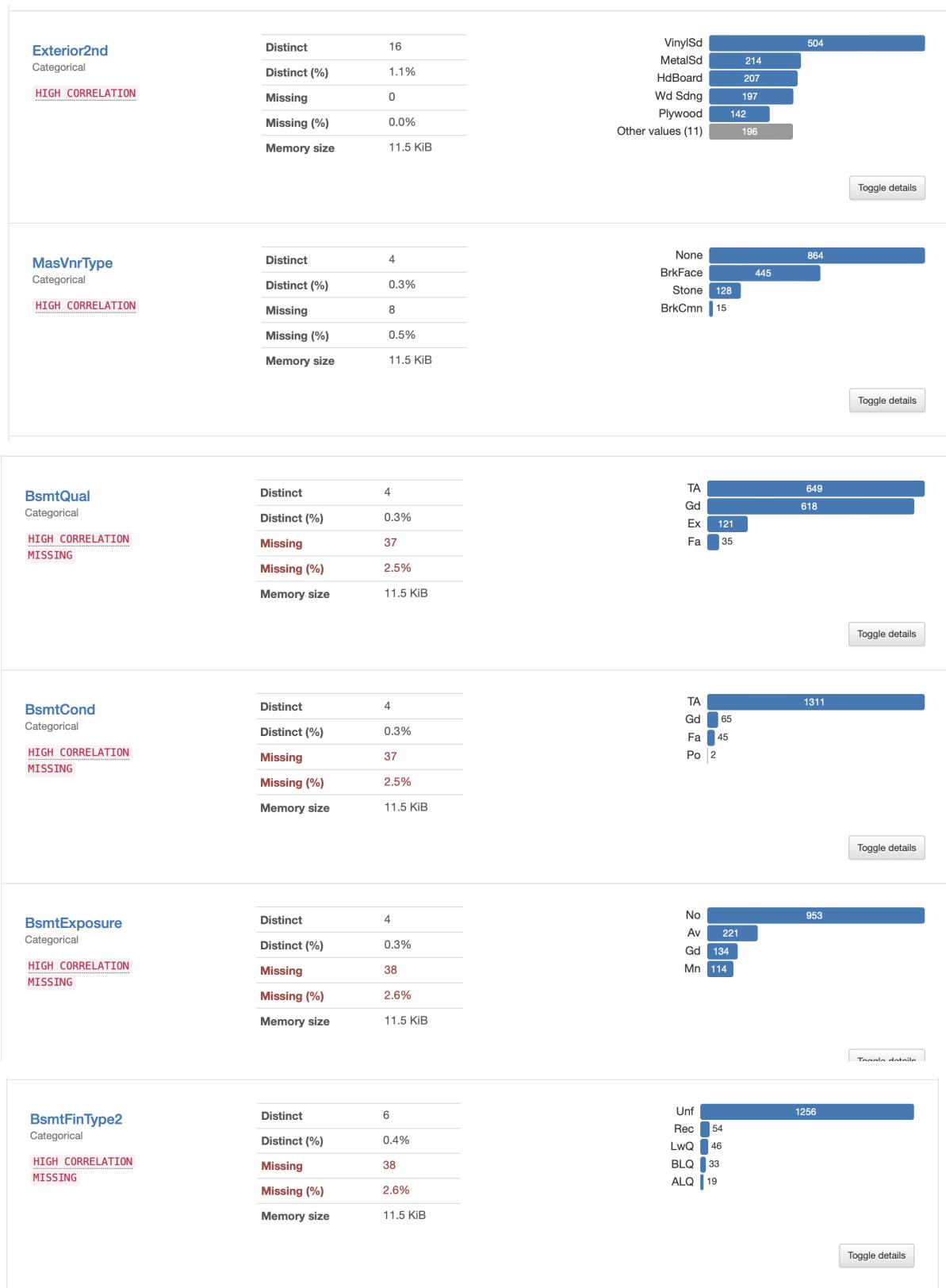
Donaldo Sebastian Garcia Jiménez 19683  
 Raul Angel Jimenez Hernandez 19017



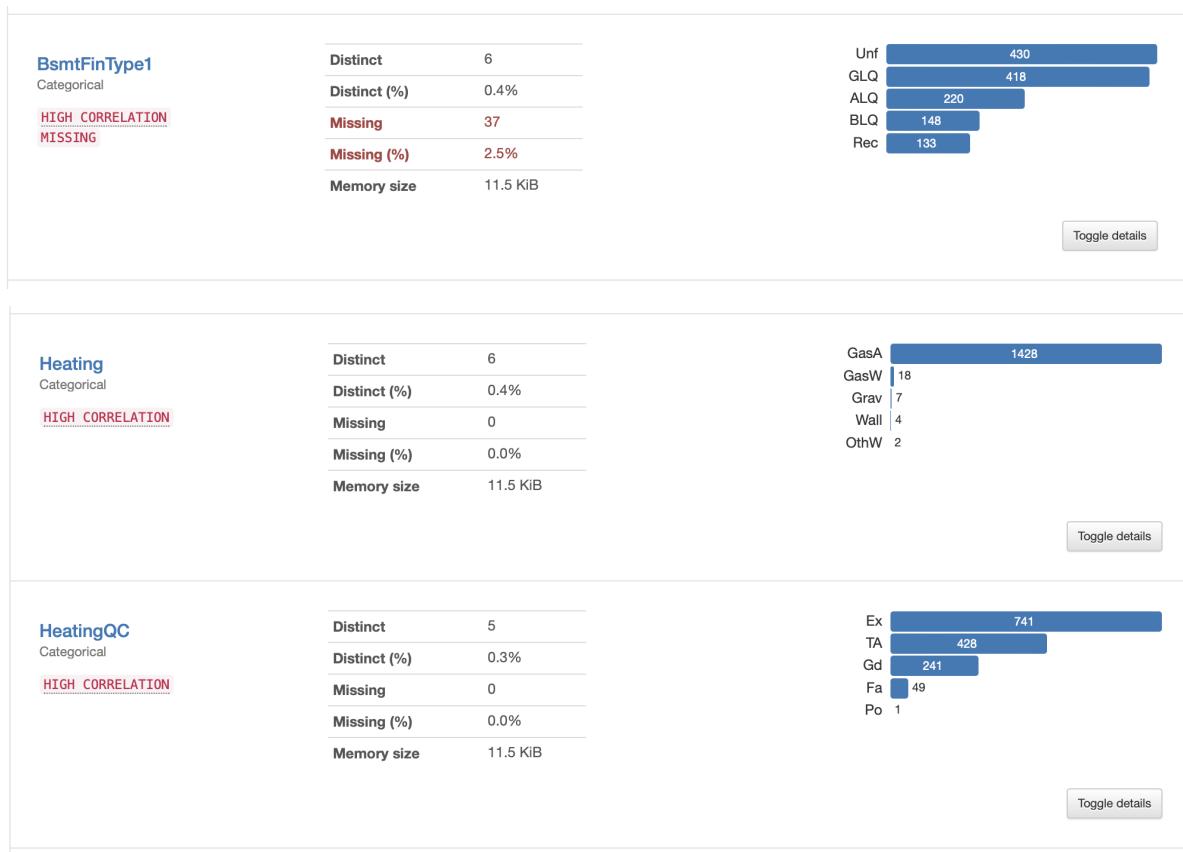
Donald Sebastian Garcia Jiménez 19683  
 Raul Angel Jimenez Hernandez 19017



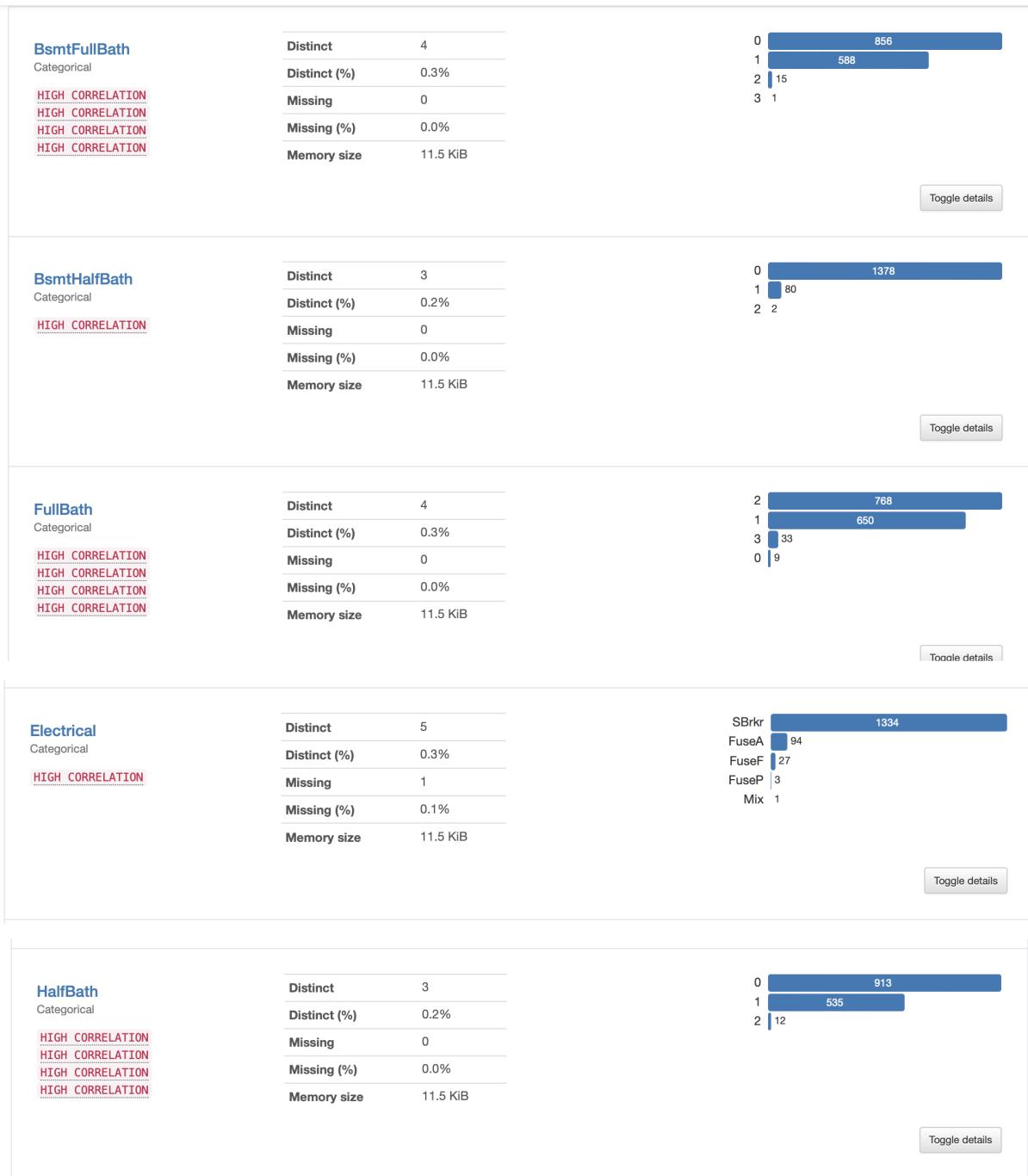
Donaldo Sebastian Garcia Jiménez 19683  
 Raul Angel Jimenez Hernandez 19017



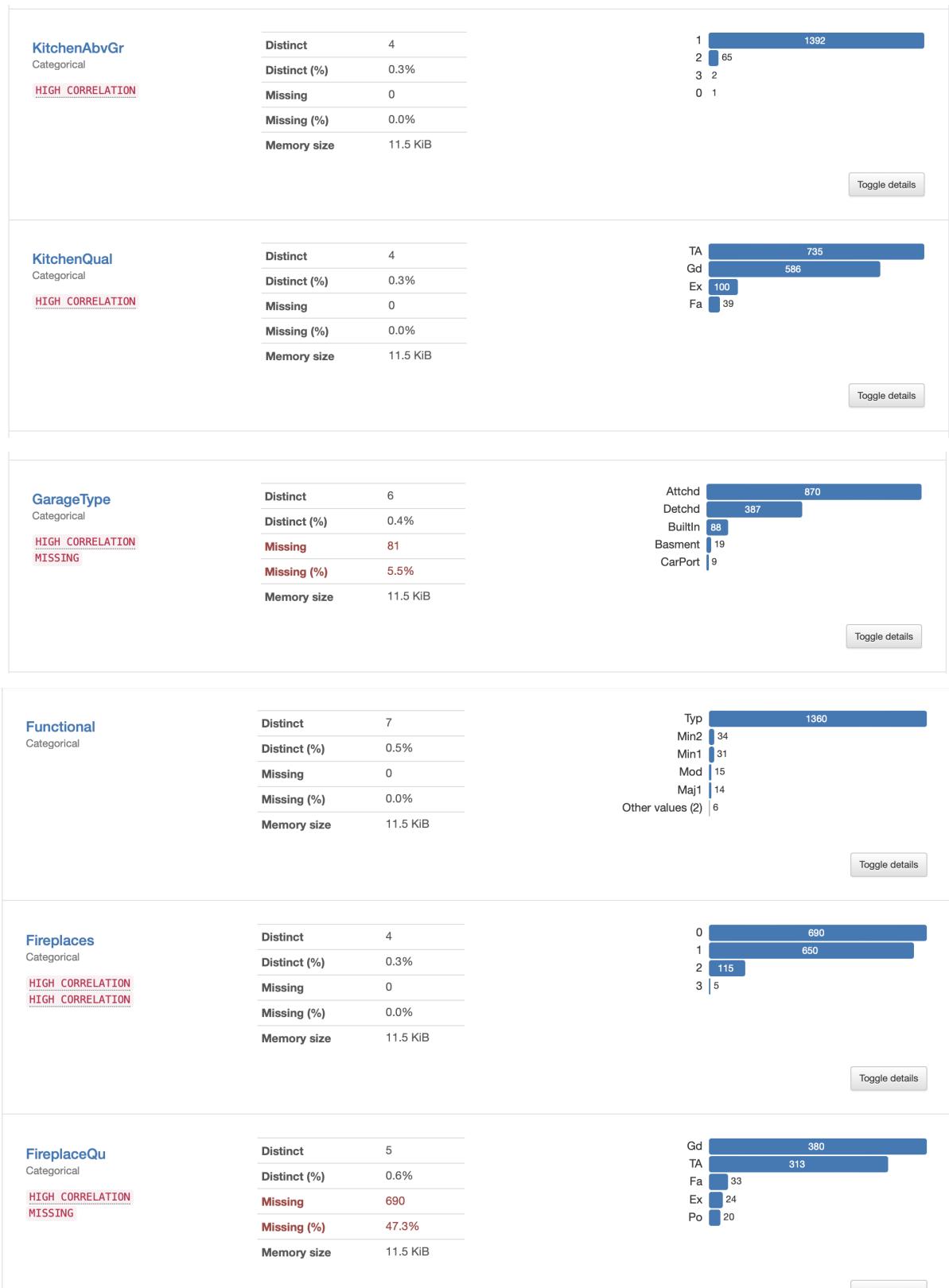
Donald Sebastian Garcia Jiménez 19683  
Raul Angel Jimenez Hernandez 19017



Donaldo Sebastian Garcia Jiménez 19683  
Raul Angel Jimenez Hernandez 19017



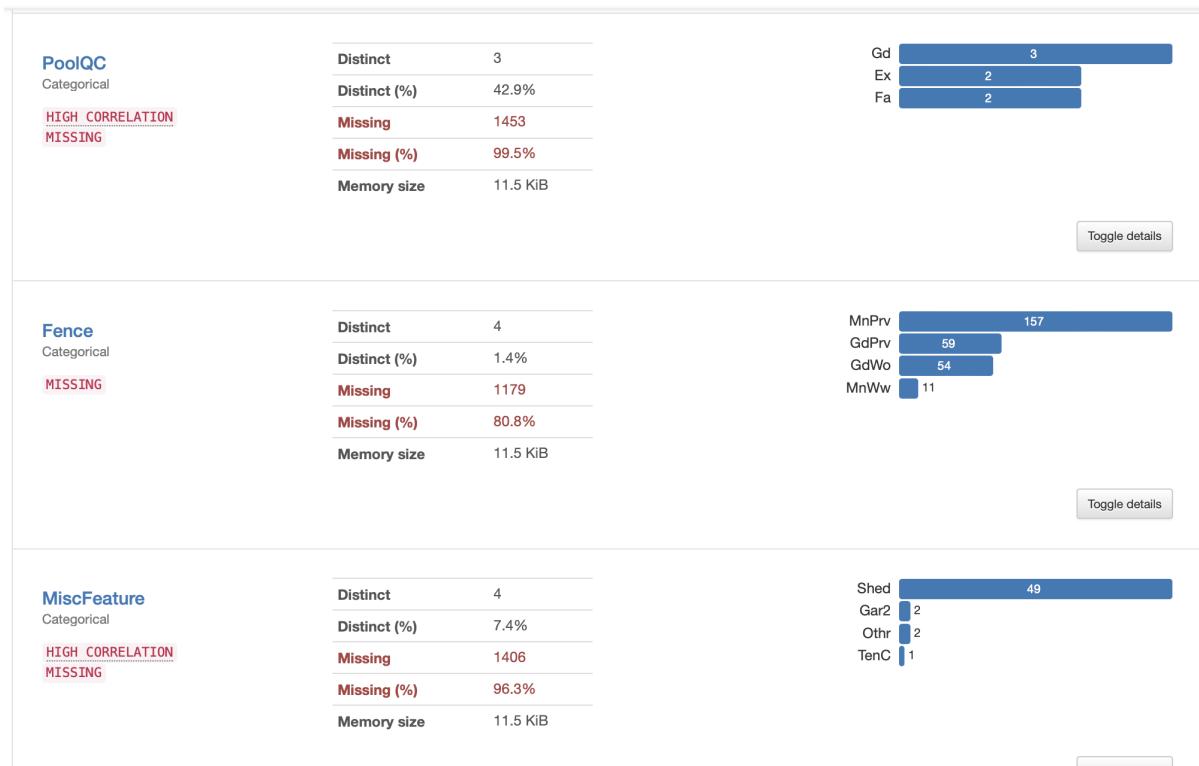
Donaldo Sebastian Garcia Jiménez 19683  
 Raul Angel Jimenez Hernandez 19017



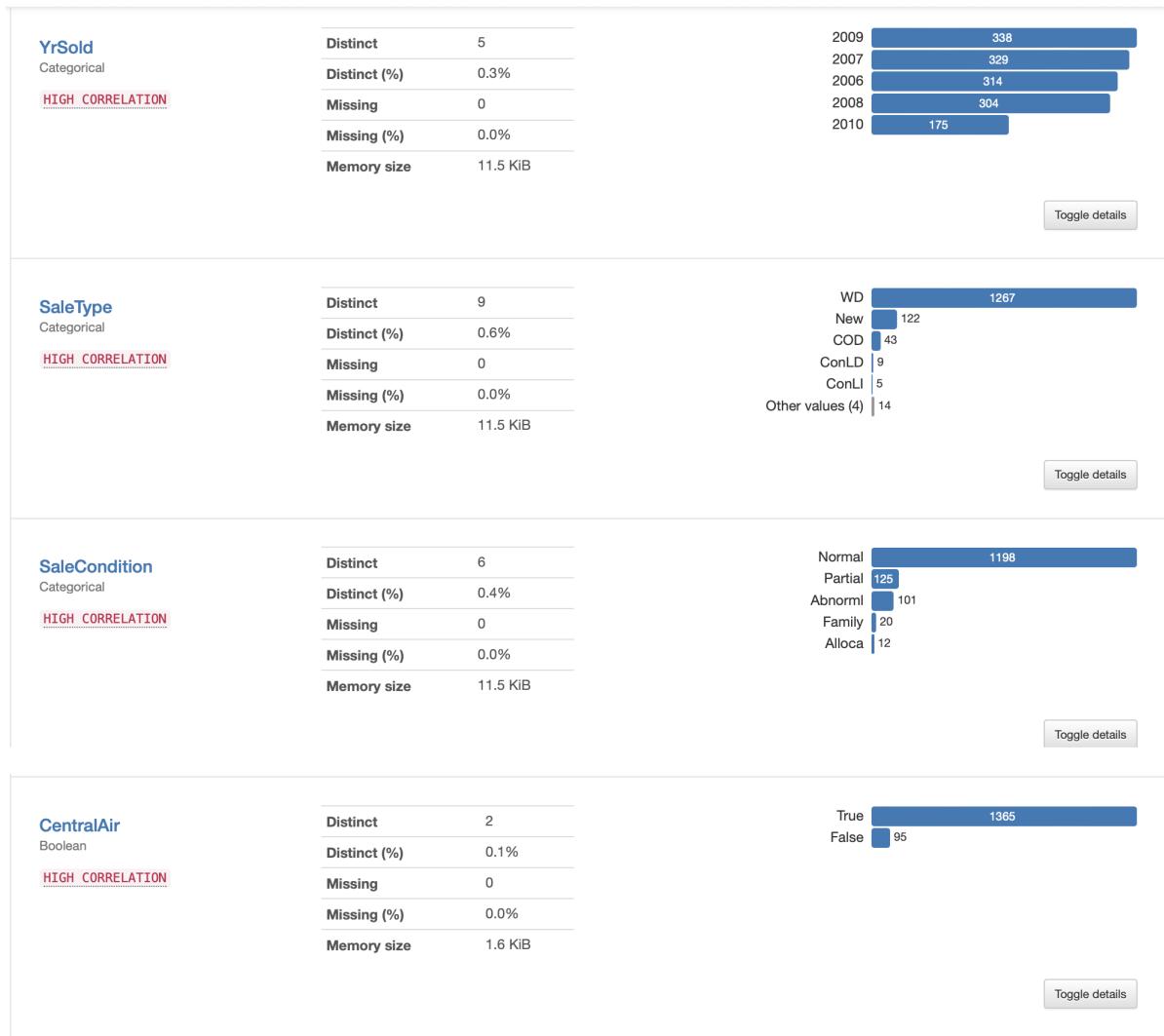
Donaldo Sebastian Garcia Jiménez 19683  
 Raul Angel Jimenez Hernandez 19017

<b>GarageQual</b> Categorical  <b>HIGH CORRELATION</b> <b>MISSING</b>	<table border="1"> <tbody> <tr><td>Distinct</td><td>5</td></tr> <tr><td>Distinct (%)</td><td>0.4%</td></tr> <tr><td>Missing</td><td>81</td></tr> <tr><td>Missing (%)</td><td>5.5%</td></tr> <tr><td>Memory size</td><td>11.5 KiB</td></tr> </tbody> </table>	Distinct	5	Distinct (%)	0.4%	Missing	81	Missing (%)	5.5%	Memory size	11.5 KiB	<table border="1"> <thead> <tr><th>Category</th><th>Count</th></tr> </thead> <tbody> <tr><td>TA</td><td>1311</td></tr> <tr><td>Fa</td><td>48</td></tr> <tr><td>Gd</td><td>14</td></tr> <tr><td>Ex</td><td>3</td></tr> <tr><td>Po</td><td>3</td></tr> </tbody> </table>	Category	Count	TA	1311	Fa	48	Gd	14	Ex	3	Po	3	<a href="#">Toggle details</a>
Distinct	5																								
Distinct (%)	0.4%																								
Missing	81																								
Missing (%)	5.5%																								
Memory size	11.5 KiB																								
Category	Count																								
TA	1311																								
Fa	48																								
Gd	14																								
Ex	3																								
Po	3																								
<b>GarageCond</b> Categorical  <b>HIGH CORRELATION</b> <b>MISSING</b>	<table border="1"> <tbody> <tr><td>Distinct</td><td>5</td></tr> <tr><td>Distinct (%)</td><td>0.4%</td></tr> <tr><td>Missing</td><td>81</td></tr> <tr><td>Missing (%)</td><td>5.5%</td></tr> <tr><td>Memory size</td><td>11.5 KiB</td></tr> </tbody> </table>	Distinct	5	Distinct (%)	0.4%	Missing	81	Missing (%)	5.5%	Memory size	11.5 KiB	<table border="1"> <thead> <tr><th>Category</th><th>Count</th></tr> </thead> <tbody> <tr><td>TA</td><td>1326</td></tr> <tr><td>Fa</td><td>35</td></tr> <tr><td>Gd</td><td>9</td></tr> <tr><td>Po</td><td>7</td></tr> <tr><td>Ex</td><td>2</td></tr> </tbody> </table>	Category	Count	TA	1326	Fa	35	Gd	9	Po	7	Ex	2	<a href="#">Toggle details</a>
Distinct	5																								
Distinct (%)	0.4%																								
Missing	81																								
Missing (%)	5.5%																								
Memory size	11.5 KiB																								
Category	Count																								
TA	1326																								
Fa	35																								
Gd	9																								
Po	7																								
Ex	2																								
<b>PavedDrive</b> Categorical  <b>HIGH CORRELATION</b>	<table border="1"> <tbody> <tr><td>Distinct</td><td>3</td></tr> <tr><td>Distinct (%)</td><td>0.2%</td></tr> <tr><td>Missing</td><td>0</td></tr> <tr><td>Missing (%)</td><td>0.0%</td></tr> <tr><td>Memory size</td><td>11.5 KiB</td></tr> </tbody> </table>	Distinct	3	Distinct (%)	0.2%	Missing	0	Missing (%)	0.0%	Memory size	11.5 KiB	<table border="1"> <thead> <tr><th>Category</th><th>Count</th></tr> </thead> <tbody> <tr><td>Y</td><td>1340</td></tr> <tr><td>N</td><td>90</td></tr> <tr><td>P</td><td>30</td></tr> </tbody> </table>	Category	Count	Y	1340	N	90	P	30	<a href="#">Toggle details</a>				
Distinct	3																								
Distinct (%)	0.2%																								
Missing	0																								
Missing (%)	0.0%																								
Memory size	11.5 KiB																								
Category	Count																								
Y	1340																								
N	90																								
P	30																								
<b>GarageFinish</b> Categorical  <b>HIGH CORRELATION</b> <b>MISSING</b>	<table border="1"> <tbody> <tr><td>Distinct</td><td>3</td></tr> <tr><td>Distinct (%)</td><td>0.2%</td></tr> <tr><td>Missing</td><td>81</td></tr> <tr><td>Missing (%)</td><td>5.5%</td></tr> <tr><td>Memory size</td><td>11.5 KiB</td></tr> </tbody> </table>	Distinct	3	Distinct (%)	0.2%	Missing	81	Missing (%)	5.5%	Memory size	11.5 KiB	<table border="1"> <thead> <tr><th>Category</th><th>Count</th></tr> </thead> <tbody> <tr><td>Unf</td><td>605</td></tr> <tr><td>RFn</td><td>422</td></tr> <tr><td>Fin</td><td>352</td></tr> </tbody> </table>	Category	Count	Unf	605	RFn	422	Fin	352	<a href="#">Toggle details</a>				
Distinct	3																								
Distinct (%)	0.2%																								
Missing	81																								
Missing (%)	5.5%																								
Memory size	11.5 KiB																								
Category	Count																								
Unf	605																								
RFn	422																								
Fin	352																								
<b>GarageCars</b> Categorical  <b>HIGH CORRELATION</b> <b>HIGH CORRELATION</b> <b>HIGH CORRELATION</b> <b>HIGH CORRELATION</b>	<table border="1"> <tbody> <tr><td>Distinct</td><td>5</td></tr> <tr><td>Distinct (%)</td><td>0.3%</td></tr> <tr><td>Missing</td><td>0</td></tr> <tr><td>Missing (%)</td><td>0.0%</td></tr> <tr><td>Memory size</td><td>11.5 KiB</td></tr> </tbody> </table>	Distinct	5	Distinct (%)	0.3%	Missing	0	Missing (%)	0.0%	Memory size	11.5 KiB	<table border="1"> <thead> <tr><th>Category</th><th>Count</th></tr> </thead> <tbody> <tr><td>2</td><td>824</td></tr> <tr><td>1</td><td>369</td></tr> <tr><td>3</td><td>181</td></tr> <tr><td>0</td><td>81</td></tr> <tr><td>4</td><td>5</td></tr> </tbody> </table>	Category	Count	2	824	1	369	3	181	0	81	4	5	<a href="#">Toggle details</a>
Distinct	5																								
Distinct (%)	0.3%																								
Missing	0																								
Missing (%)	0.0%																								
Memory size	11.5 KiB																								
Category	Count																								
2	824																								
1	369																								
3	181																								
0	81																								
4	5																								
<b>CentralAir</b> Boolean  <b>HIGH CORRELATION</b>	<table border="1"> <tbody> <tr><td>Distinct</td><td>2</td></tr> <tr><td>Distinct (%)</td><td>0.1%</td></tr> <tr><td>Missing</td><td>0</td></tr> <tr><td>Missing (%)</td><td>0.0%</td></tr> <tr><td>Memory size</td><td>1.6 KiB</td></tr> </tbody> </table>	Distinct	2	Distinct (%)	0.1%	Missing	0	Missing (%)	0.0%	Memory size	1.6 KiB	<table border="1"> <thead> <tr><th>Category</th><th>Count</th></tr> </thead> <tbody> <tr><td>True</td><td>1365</td></tr> <tr><td>False</td><td>95</td></tr> </tbody> </table>	Category	Count	True	1365	False	95	<a href="#">Toggle details</a>						
Distinct	2																								
Distinct (%)	0.1%																								
Missing	0																								
Missing (%)	0.0%																								
Memory size	1.6 KiB																								
Category	Count																								
True	1365																								
False	95																								

Donaldo Sebastian Garcia Jiménez 19683  
Raul Angel Jimenez Hernandez 19017



Donaldo Sebastian Garcia Jiménez 19683  
Raul Angel Jimenez Hernandez 19017



## #6 Análisis de componentes principales con variables numéricas

Estas gráficas y análisis se pueden observar más detalladamente en el reporte de abajo

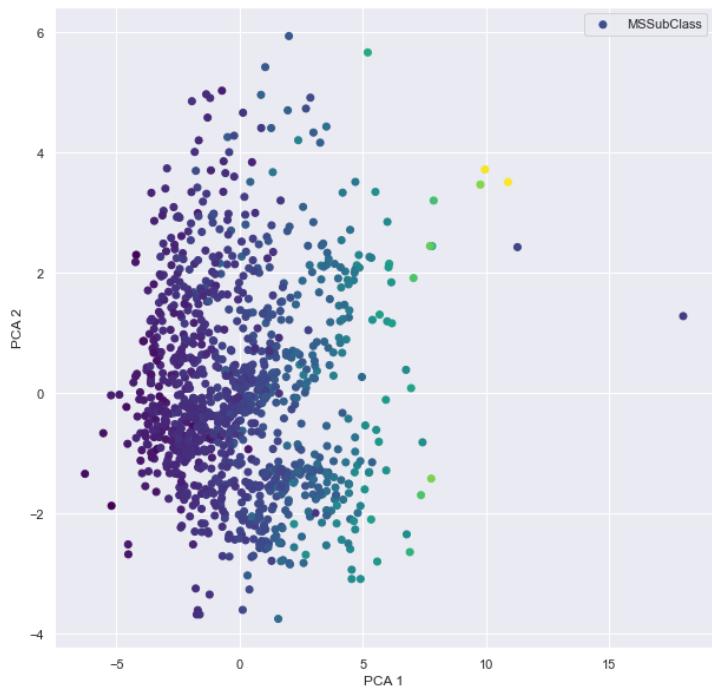
Donaldo Sebastian Garcia Jiménez 19683  
Raul Angel Jimenez Hernandez 19017

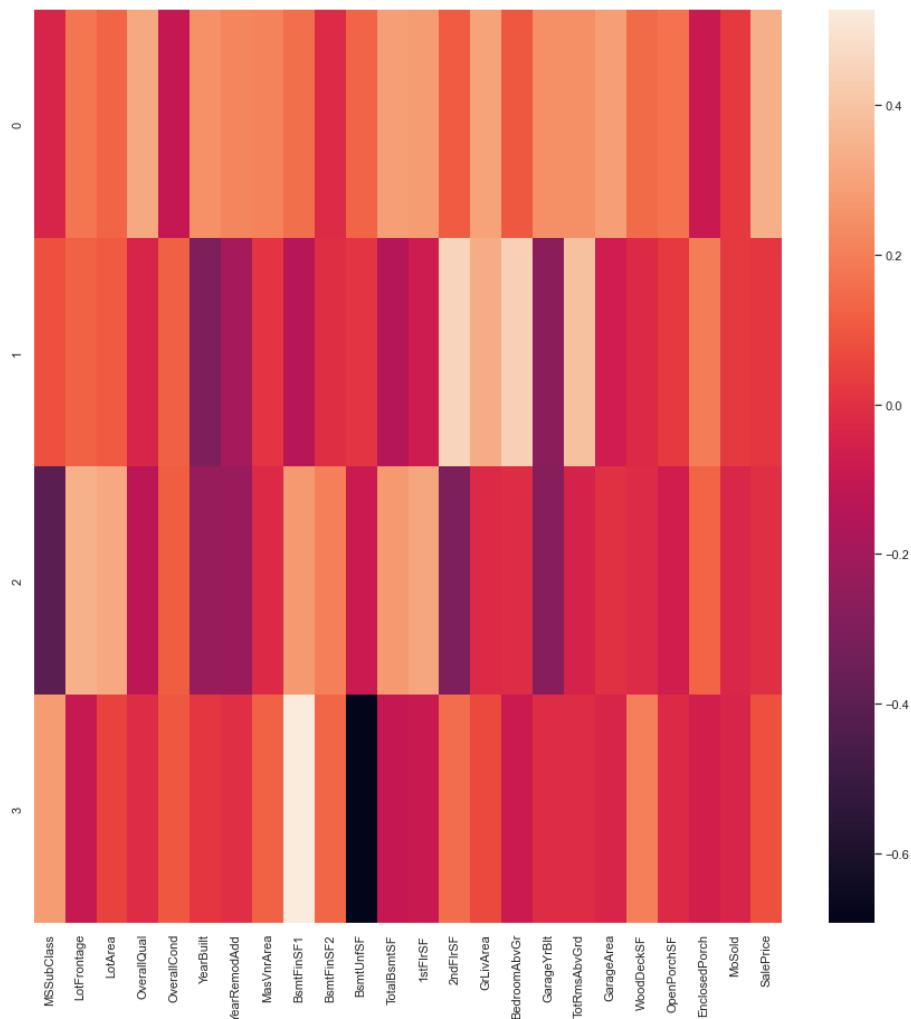
```
from factor_analyzer.factor_analyzer import calculate_bartlett_sphericity
chi_square_value,p_value=calculate_bartlett_sphericity(houses_clean)
chi_square_value, p_value

(56989.803679788645, 0.0)

from factor_analyzer.factor_analyzer import calculate_kmo
kmo_all,kmo_model=calculate_kmo(houses_clean)
kmo_all, kmo_model

c:\Users\ALIWARE\AppData\Local\Programs\Python\Python39\lib\site-packages\factor_analyzer\utils.py:249
    warnings.warn('The inverse of the variance-covariance matrix '
array([0.67986034, 0.85336964, 0.86511003, 0.91950297, 0.52939275,
       0.76092737, 0.83176101, 0.94283582, 0.70914247, 0.11991146,
       0.64250957, 0.87220836, 0.62952928, 0.45033991, 0.66780668,
       0.77086765, 0.82586225, 0.91502783, 0.91616437, 0.94585033,
       0.93825884, 0.71845741, 0.46740089, 0.92301962]),
0.7683818262377667)
```





## #7 Reglas de asociación interesantes del dataset

```
> ^ # El minimo de cobertura o soporte es de 20% y el minimo de confianza es de 70%
reglas_asociacion = apriori(records, min_support=0.2, min_confidence=0.7)
reglas = list(reglas_asociacion)

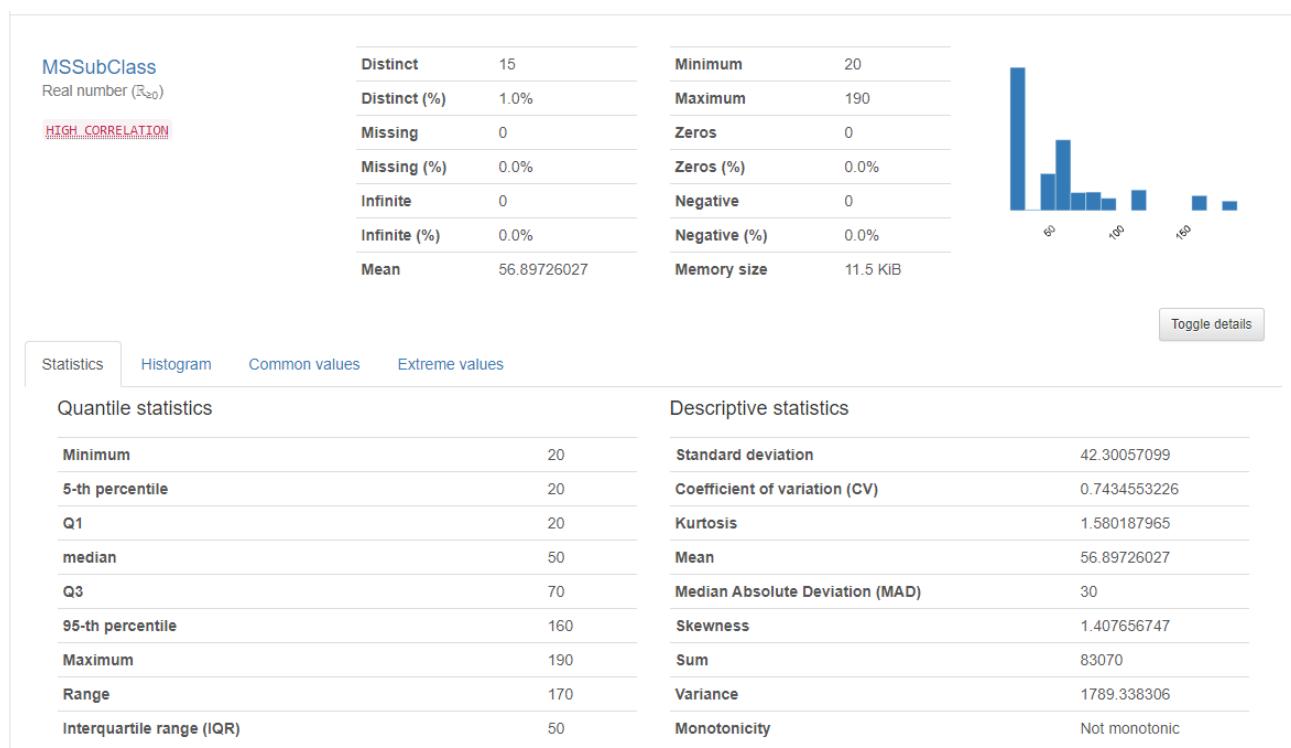
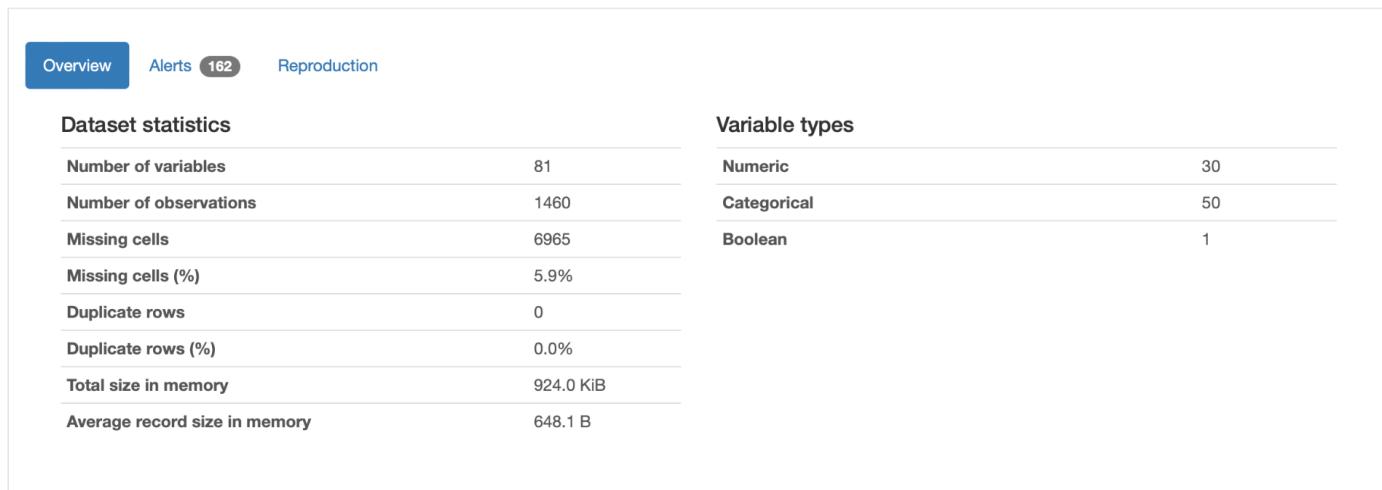
41] print(len(reglas))

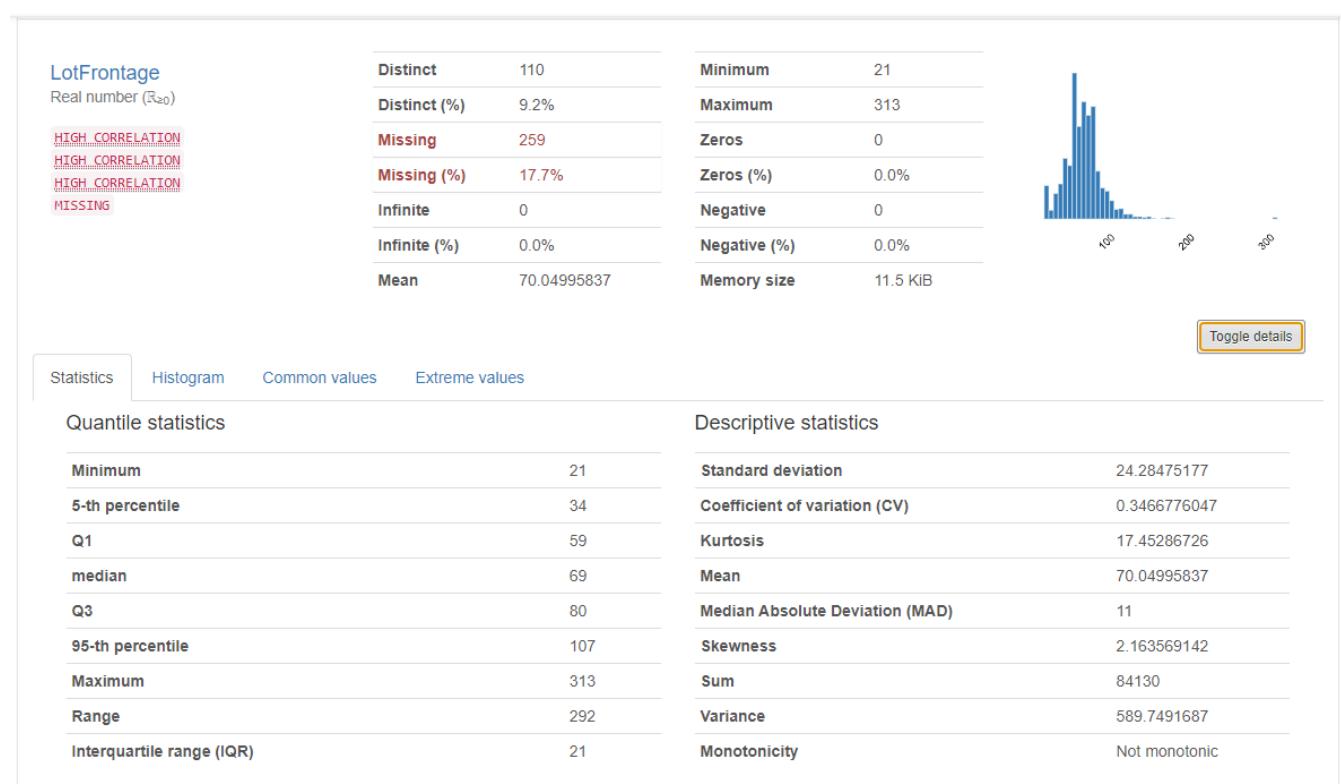
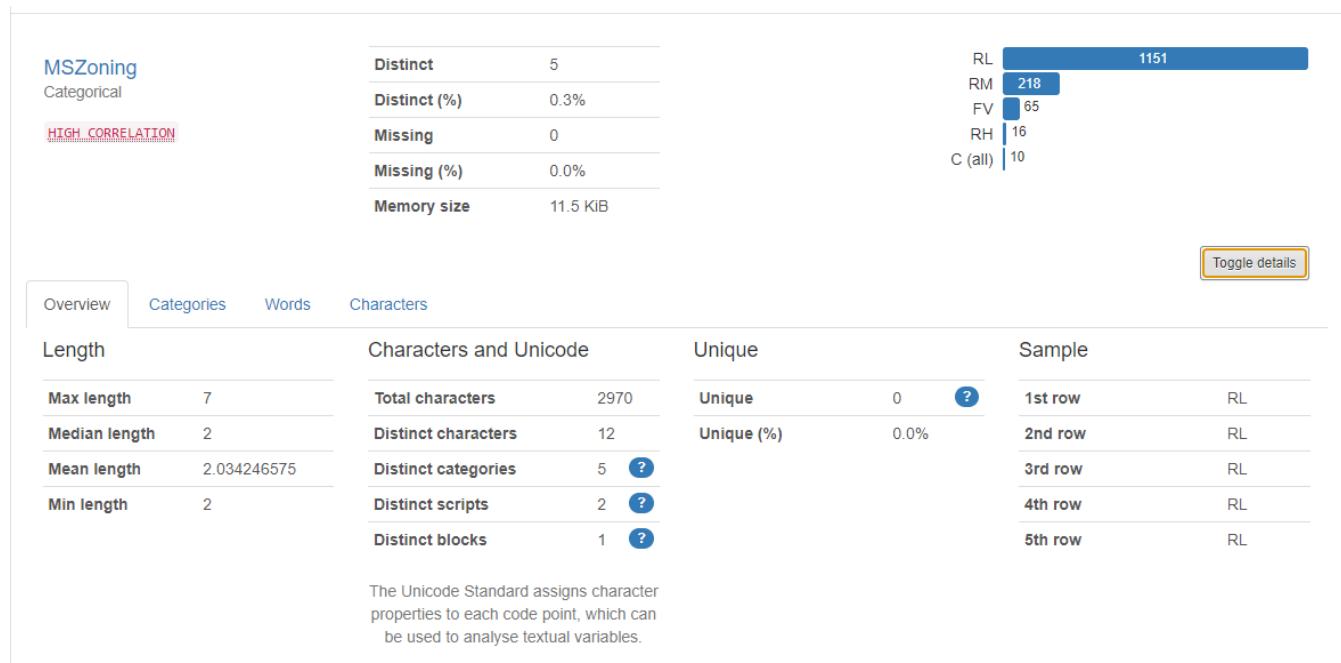
42] .. 1531
```

# Reporte de Datos

## Análisis exploratorio

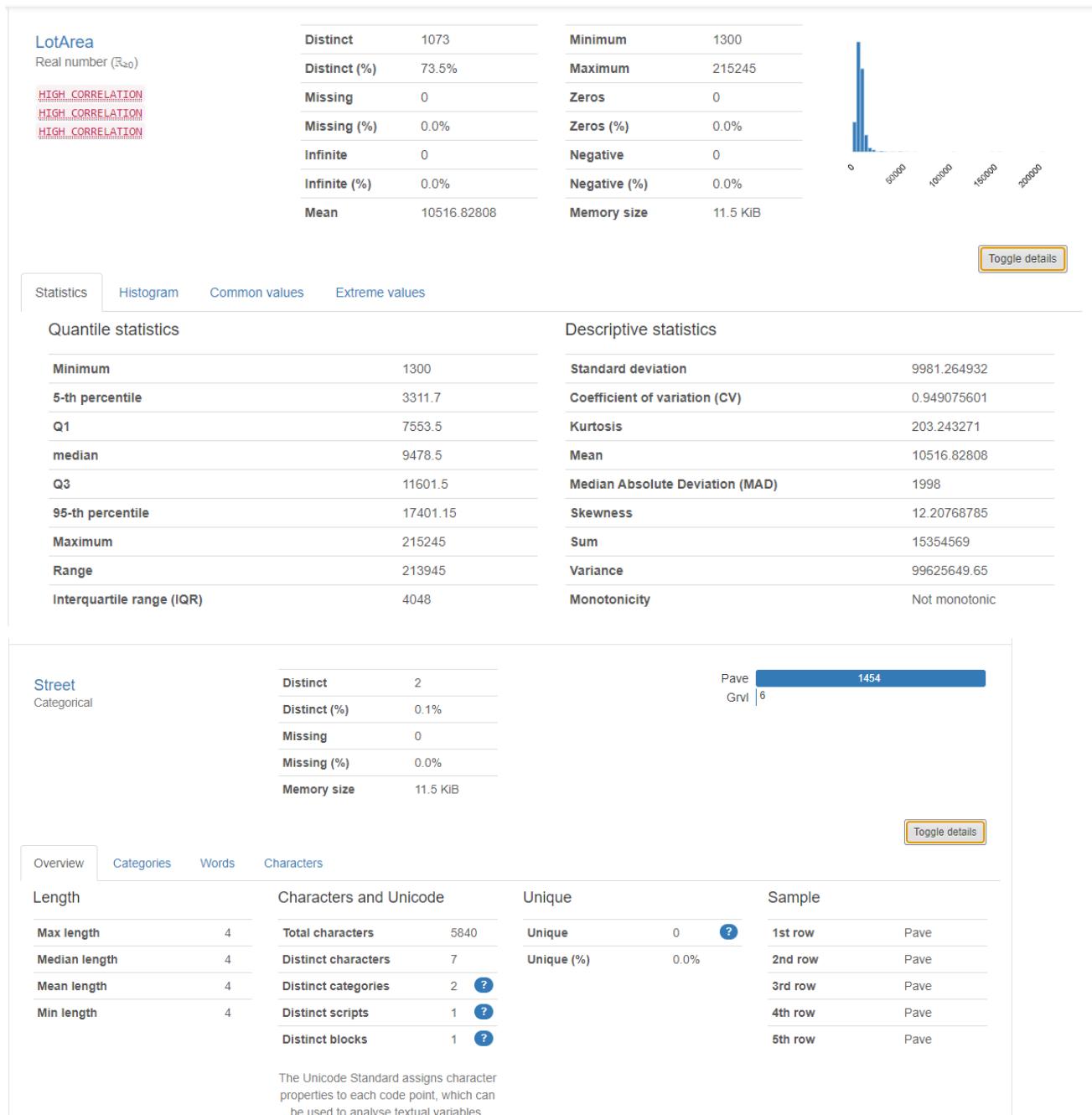
### Overview





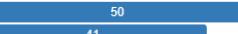
# Donaldo Sebastian Garcia Jiménez 19683

## Raul Angel Jimenez Hernandez 19017



# Donald Sebastian Garcia Jiménez 19683

## Raul Angel Jimenez Hernandez 19017

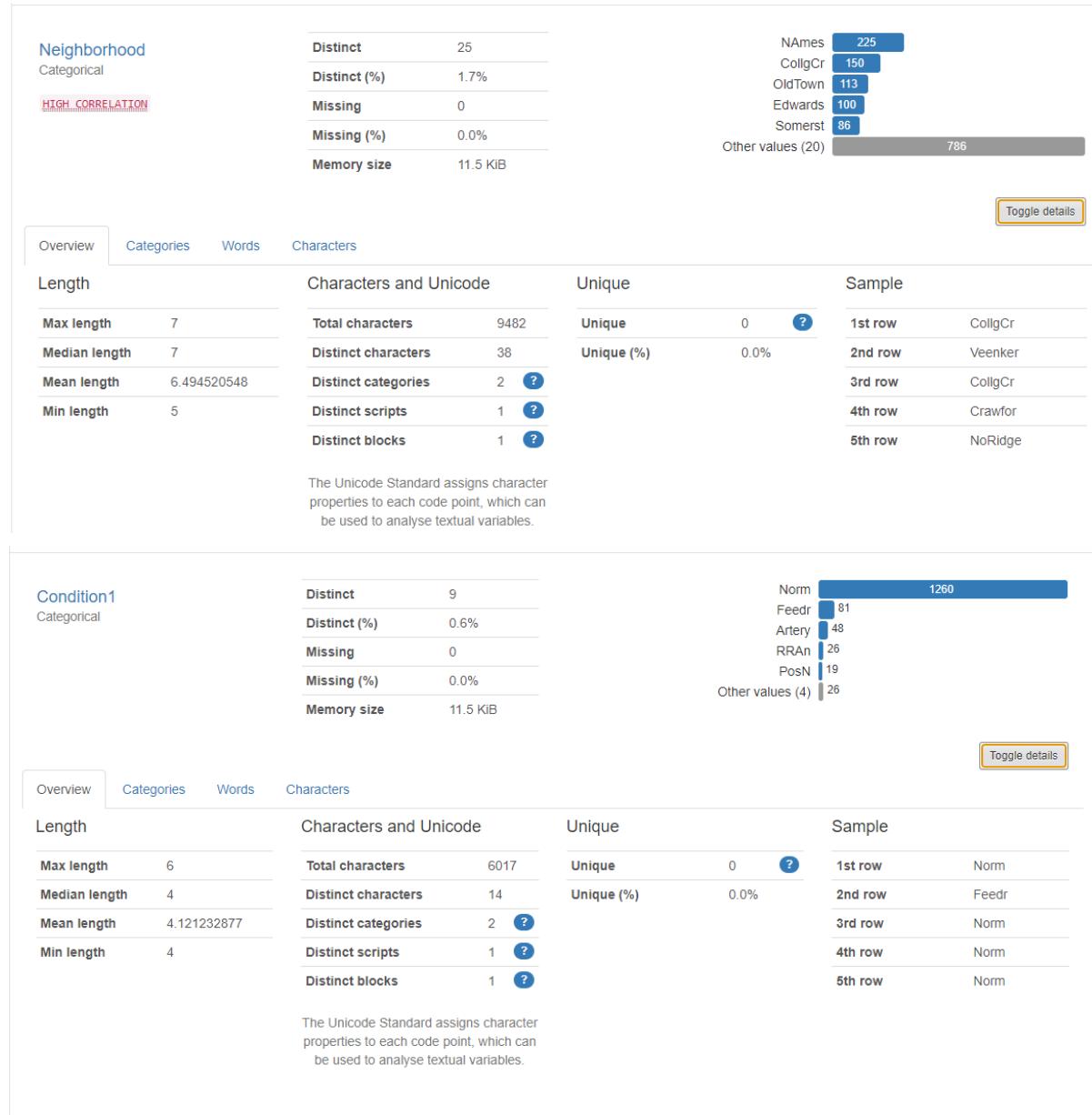
<b>Alley</b> Categorical  <span style="background-color: #f0f0f0; border: 1px solid #ccc; padding: 2px;">HIGH CORRELATION</span> <span style="background-color: #ffcccc; border: 1px solid #ccc; padding: 2px;">MISSING</span>	<table border="1"> <tr><td>Distinct</td><td>2</td></tr> <tr><td>Distinct (%)</td><td>2.2%</td></tr> <tr><td>Missing</td><td>1369</td></tr> <tr><td>Missing (%)</td><td>93.8%</td></tr> <tr><td>Memory size</td><td>11.5 KiB</td></tr> </table>		Distinct	2	Distinct (%)	2.2%	Missing	1369	Missing (%)	93.8%	Memory size	11.5 KiB	Grvl  Pave 	
Distinct	2													
Distinct (%)	2.2%													
Missing	1369													
Missing (%)	93.8%													
Memory size	11.5 KiB													
<a href="#">Overview</a>	<a href="#">Categories</a>	<a href="#">Words</a>	<a href="#">Characters</a>											
Length		Characters and Unicode												
Max length	4	Total characters	364											
Median length	4	Distinct characters	7											
Mean length	4	Distinct categories	2 											
Min length	4	Distinct scripts	1 											
		Distinct blocks	1 											
<p>The Unicode Standard assigns character properties to each code point, which can be used to analyse textual variables.</p>														
<a href="#">Sample</a>		<a href="#">Toggle details</a>												
<b>LotShape</b> Categorical  <span style="background-color: #f0f0f0; border: 1px solid #ccc; padding: 2px;">HIGH CORRELATION</span>	Unique	0 	1st row											
	Unique (%)	0.0%	Grvl											
	<a href="#">Toggle details</a>		<a href="#">Sample</a>											
	Reg	925	<a href="#">1st row</a>											
	IR1	484	<a href="#">2nd row</a>											
	IR2	41	<a href="#">3rd row</a>											
	IR3	10	<a href="#">4th row</a>											
	<a href="#">5th row</a>		<a href="#">Pave</a>											
	<a href="#">Toggle details</a>		<a href="#">Sample</a>											
	<a href="#">Characters</a>		<a href="#">Toggle details</a>											
<b>Length</b>  <span style="background-color: #f0f0f0; border: 1px solid #ccc; padding: 2px;">HIGH CORRELATION</span>	Length		Characters and Unicode											
	Max length	3	Total characters	4380										
	Median length	3	Distinct characters	7										
	Mean length	3	Distinct categories	3 										
	Min length	3	Distinct scripts	2 										
			Distinct blocks	1 										
	<p>The Unicode Standard assigns character properties to each code point, which can be used to analyse textual variables.</p>													
	<a href="#">Sample</a>		<a href="#">Toggle details</a>											
	1st row	Reg	<a href="#">1st row</a>											
	2nd row	Reg	<a href="#">2nd row</a>											
<b>LandContour</b> Categorical  <span style="background-color: #f0f0f0; border: 1px solid #ccc; padding: 2px;">HIGH CORRELATION</span>	3rd row	IR1	<a href="#">3rd row</a>											
	4th row	IR1	<a href="#">4th row</a>											
	5th row	IR1	<a href="#">5th row</a>											
	<a href="#">Toggle details</a>		<a href="#">Sample</a>											
	<a href="#">Characters</a>		<a href="#">Toggle details</a>											
	Length		Characters and Unicode											
	Max length	3	Total characters	4380										
	Median length	3	Distinct characters	10										
	Mean length	3	Distinct categories	2 										
	Min length	3	Distinct scripts	1 										
			Distinct blocks	1 										
	<p>The Unicode Standard assigns character properties to each code point, which can be used to analyse textual variables.</p>													
	<a href="#">Sample</a>		<a href="#">Toggle details</a>											
	1st row	Lvl	<a href="#">1st row</a>											
	2nd row	Bnk	<a href="#">2nd row</a>											
	3rd row	HLS	<a href="#">3rd row</a>											
	4th row	Low	<a href="#">4th row</a>											
	5th row		<a href="#">5th row</a>											

# Donaldo Sebastian Garcia Jiménez 19683

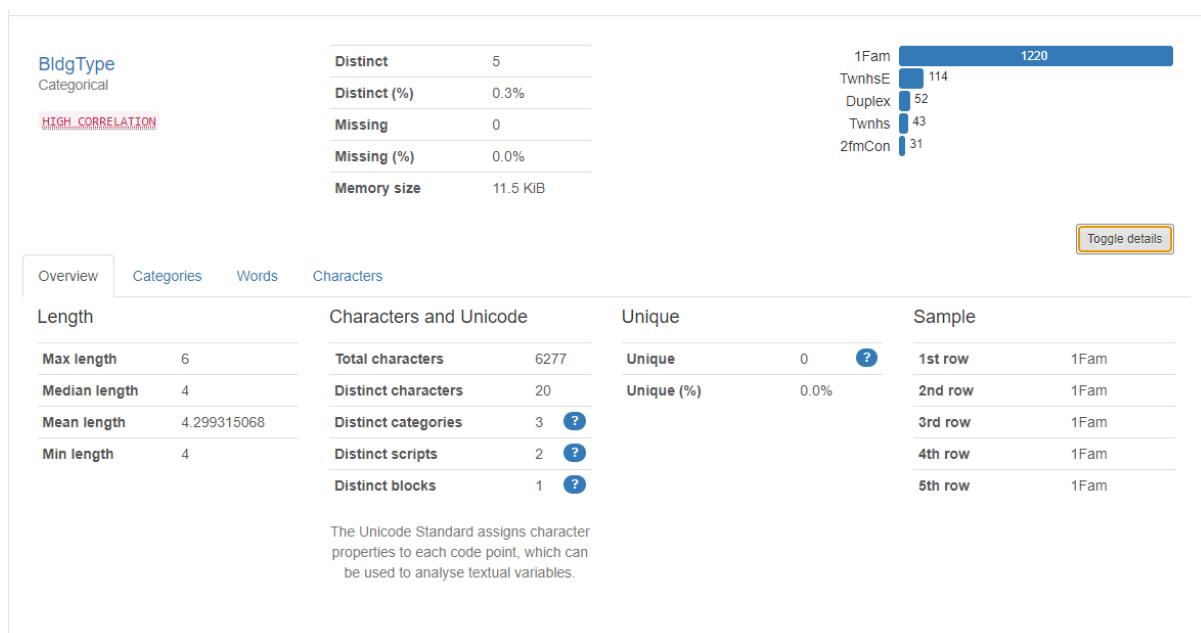
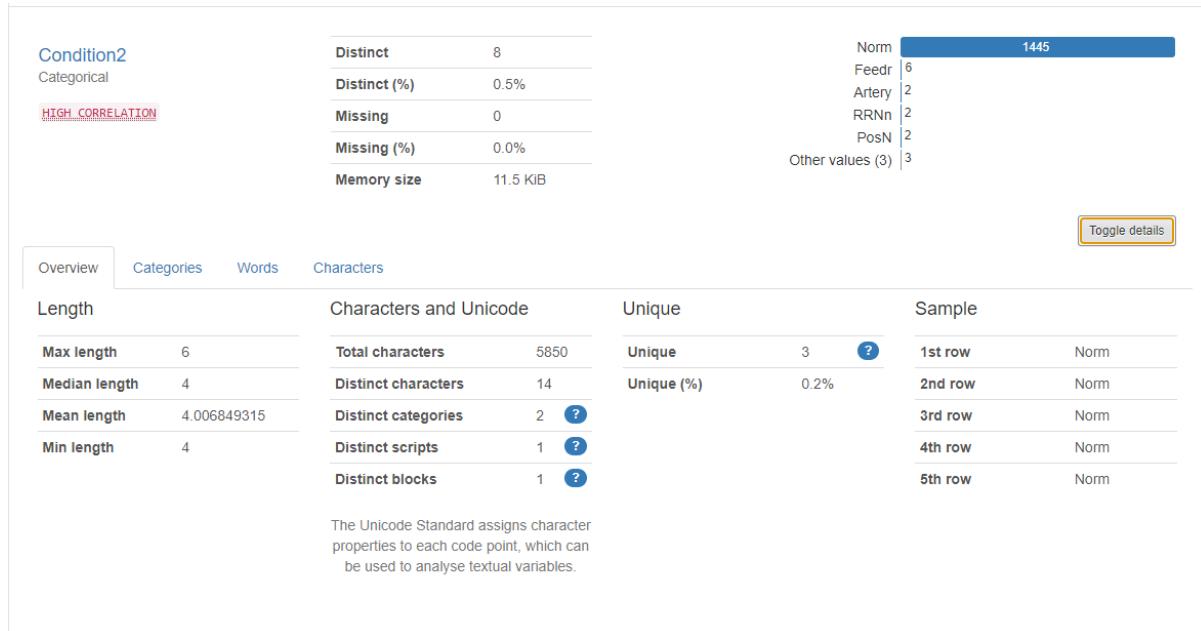
## Raul Angel Jimenez Hernandez 19017

Utilities Categorical	Distinct	2	AllPub	1459
	Distinct (%)	0.1%	NoSeWa	1
	Missing	0		
	Missing (%)	0.0%		
	Memory size	11.5 KIB		
<a href="#">Toggle details</a>				
<a href="#">Overview</a> <a href="#">Categories</a> <a href="#">Words</a> <a href="#">Characters</a>				
Length		Characters and Unicode	Unique	Sample
Max length	6	Total characters	8760	1st row AllPub
Median length	6	Distinct characters	11	2nd row AllPub
Mean length	6	Distinct categories	2 <a href="#">?</a>	3rd row AllPub
Min length	6	Distinct scripts	1 <a href="#">?</a>	4th row AllPub
		Distinct blocks	1 <a href="#">?</a>	5th row AllPub
<p>The Unicode Standard assigns character properties to each code point, which can be used to analyse textual variables.</p>				
 <b>LotConfig</b> Categorical				
<a href="#">HIGH CORRELATION</a>		Distinct	5	Inside 1052
		Distinct (%)	0.3%	Corner 263
		Missing	0	CulDSac 94
		Missing (%)	0.0%	FR2 47
		Memory size	11.5 KIB	FR3 4
<a href="#">Toggle details</a>				
<a href="#">Overview</a> <a href="#">Categories</a> <a href="#">Words</a> <a href="#">Characters</a>				
Length		Characters and Unicode	Unique	Sample
Max length	7	Total characters	8701	1st row Inside
Median length	6	Distinct characters	19	2nd row FR2
Mean length	5.959589041	Distinct categories	3 <a href="#">?</a>	3rd row Inside
Min length	3	Distinct scripts	2 <a href="#">?</a>	4th row Corner
		Distinct blocks	1 <a href="#">?</a>	5th row FR2
<p>The Unicode Standard assigns character properties to each code point, which can be used to analyse textual variables.</p>				
 <b>LandSlope</b> Categorical				
<a href="#">HIGH CORRELATION</a>		Distinct	3	Gtl 1382
		Distinct (%)	0.2%	Mod 65
		Missing	0	Sev 13
		Missing (%)	0.0%	
		Memory size	11.5 KIB	
<a href="#">Toggle details</a>				
<a href="#">Overview</a> <a href="#">Categories</a> <a href="#">Words</a> <a href="#">Characters</a>				
Length		Characters and Unicode	Unique	Sample
Max length	3	Total characters	4380	1st row Gtl
Median length	3	Distinct characters	9	2nd row Gtl
Mean length	3	Distinct categories	2 <a href="#">?</a>	3rd row Gtl
Min length	3	Distinct scripts	1 <a href="#">?</a>	4th row Gtl
		Distinct blocks	1 <a href="#">?</a>	5th row Gtl
<p>The Unicode Standard assigns character properties to each code point, which can be used to analyse textual variables.</p>				

Donaldo Sebastian Garcia Jiménez 19683  
 Raul Angel Jimenez Hernandez 19017

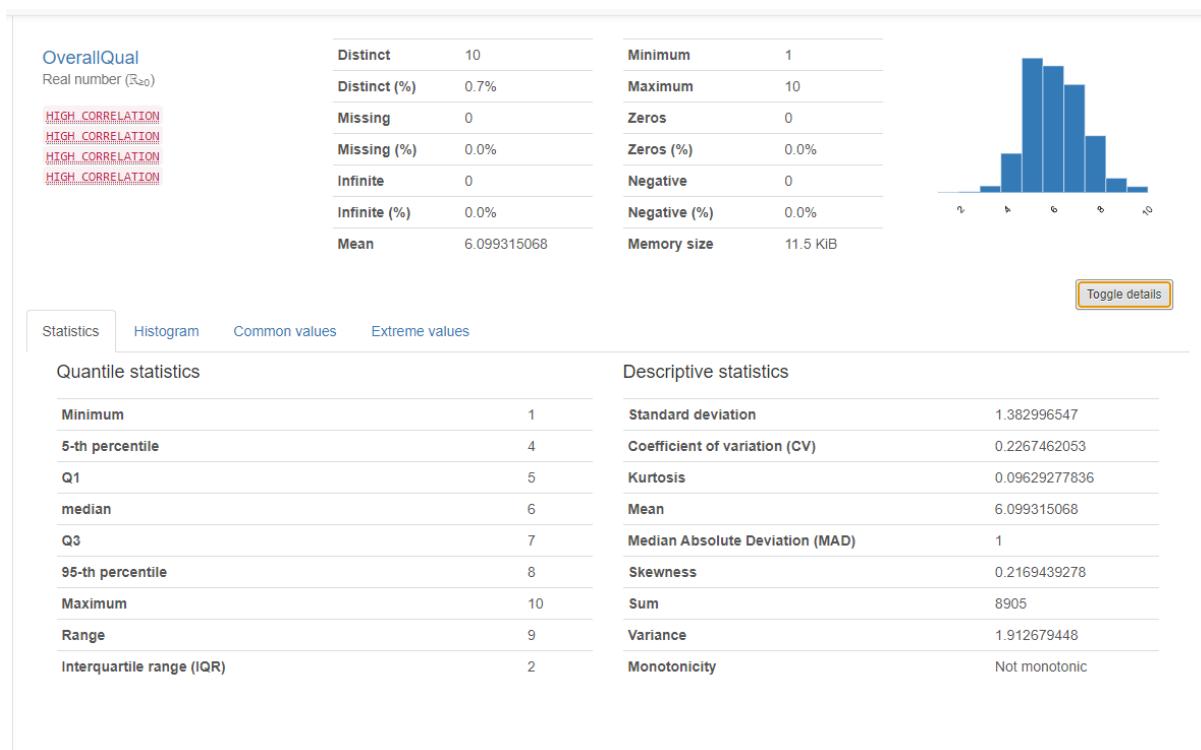
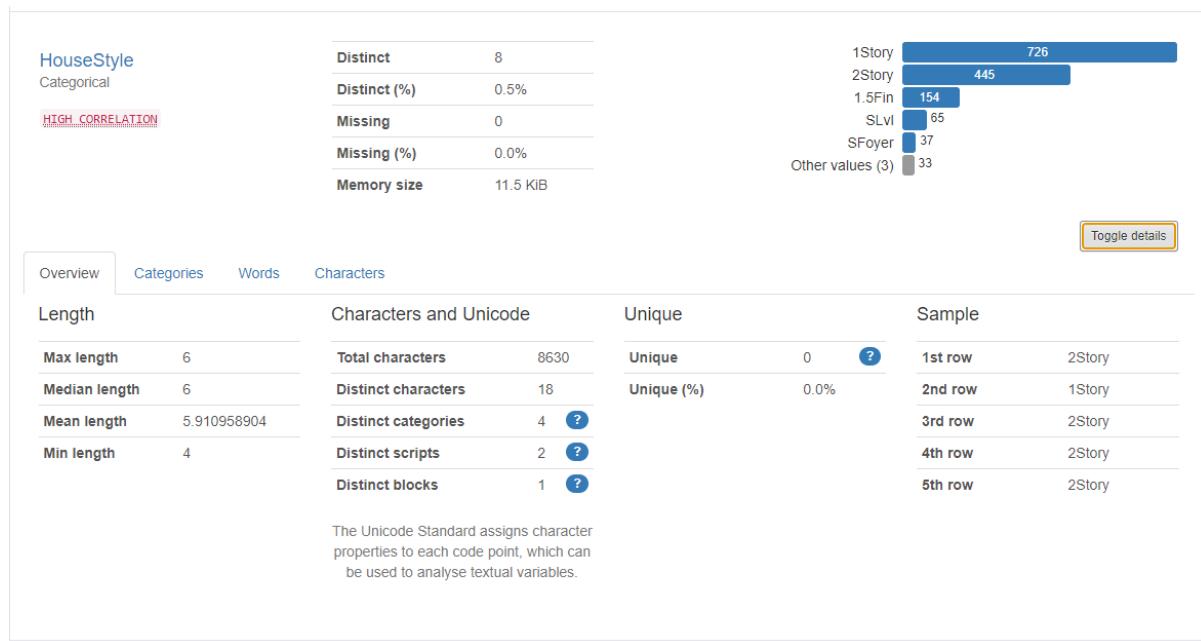


Donaldo Sebastian Garcia Jiménez 19683  
 Raul Angel Jimenez Hernandez 19017

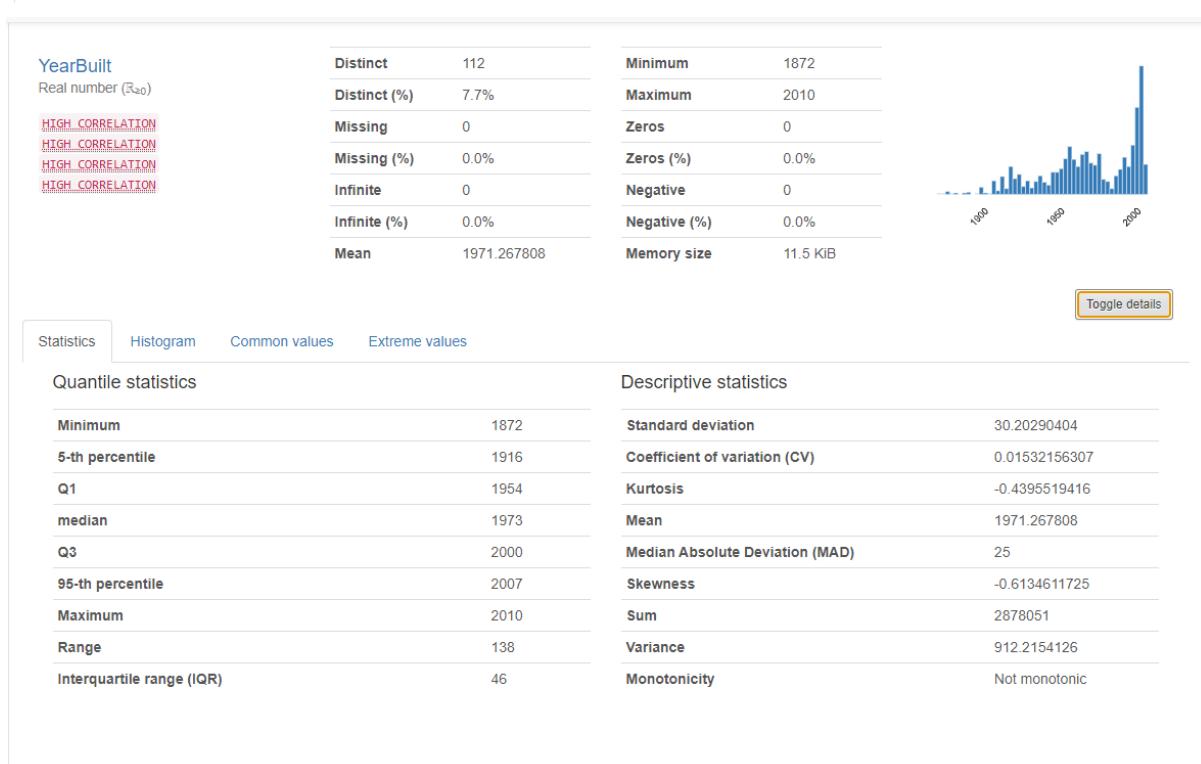
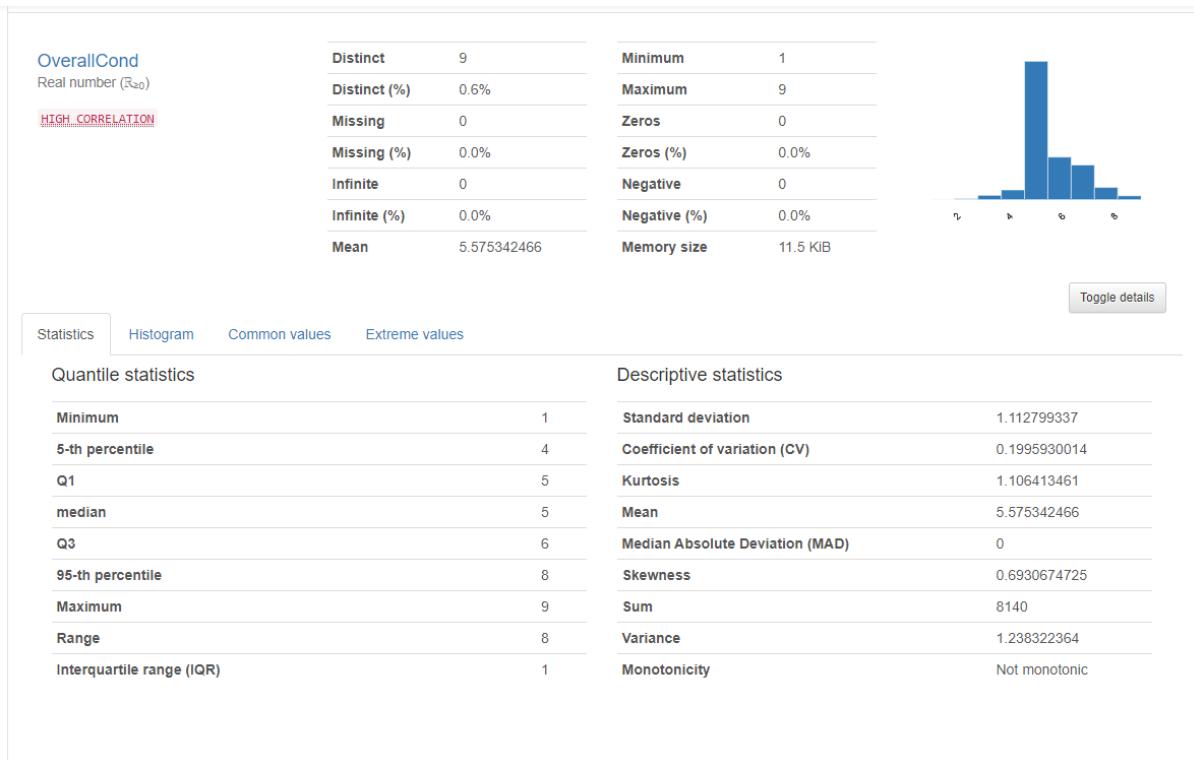


# Donald Sebastian Garcia Jiménez 19683

## Raul Angel Jimenez Hernandez 19017

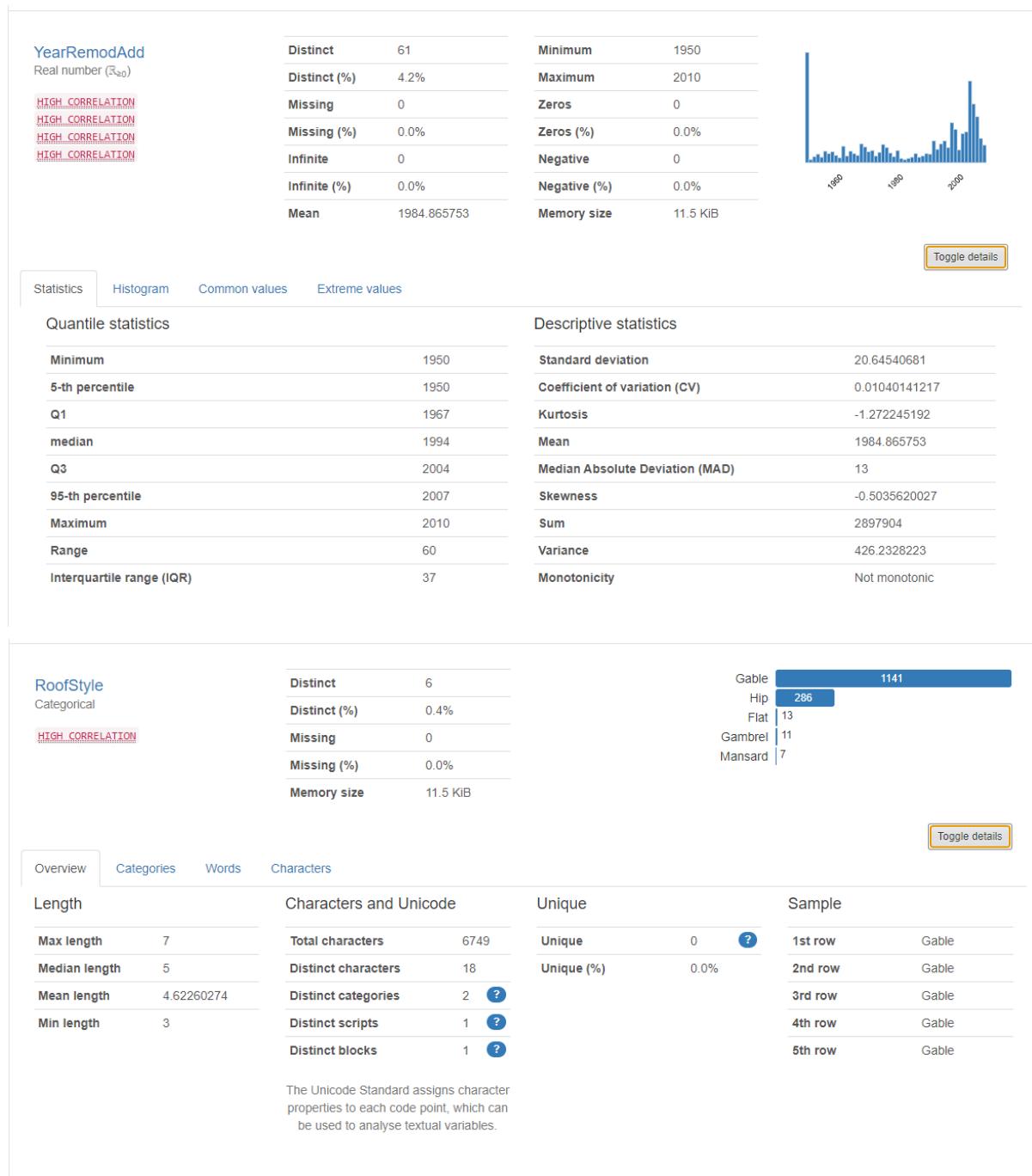


Donaldo Sebastian Garcia Jiménez 19683  
 Raul Angel Jimenez Hernandez 19017

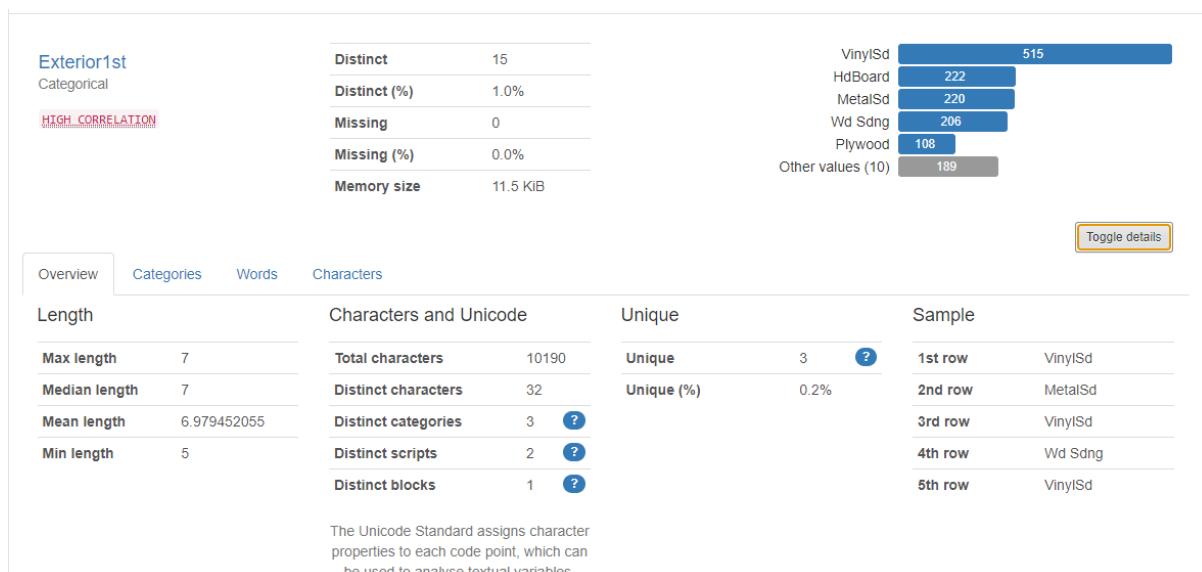
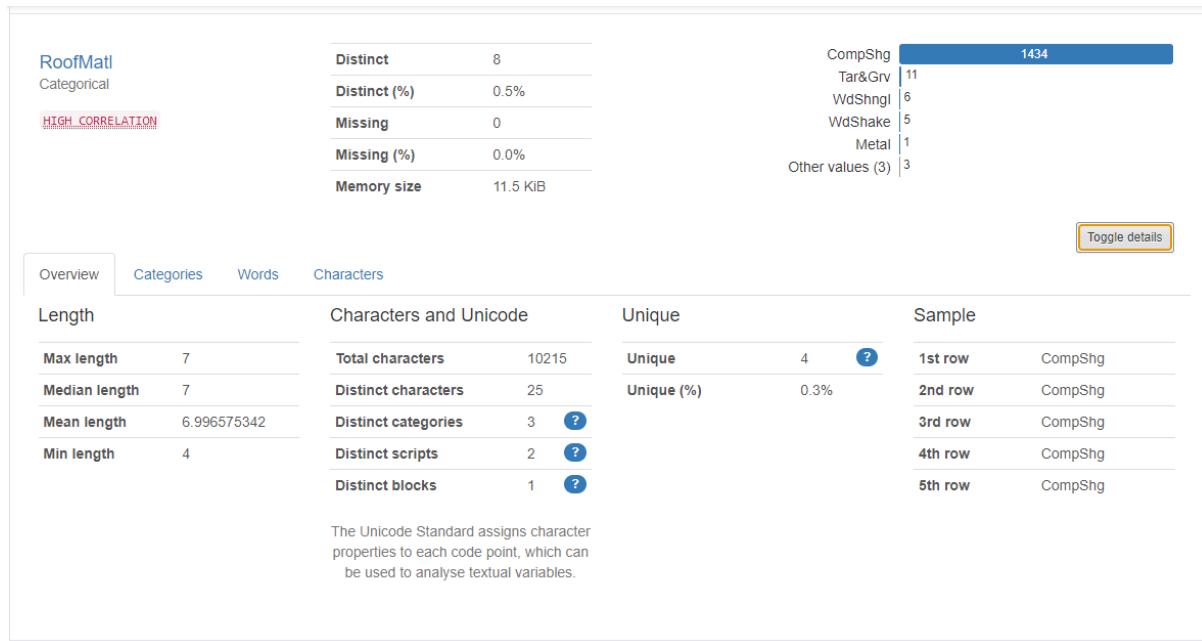


# Donaldo Sebastian Garcia Jiménez 19683

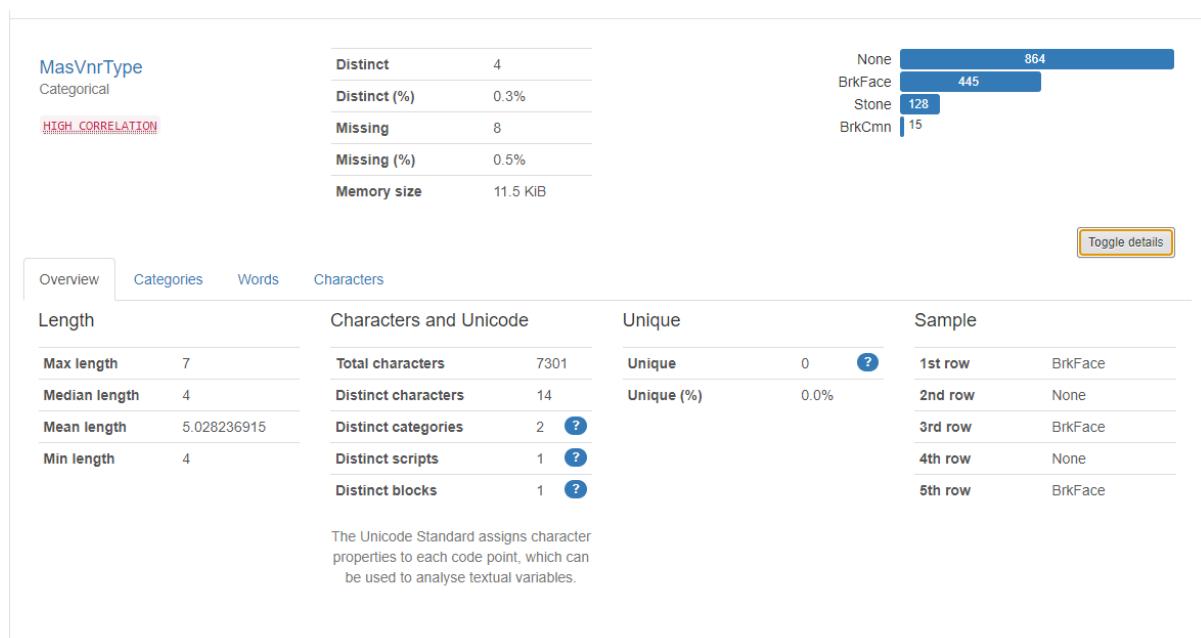
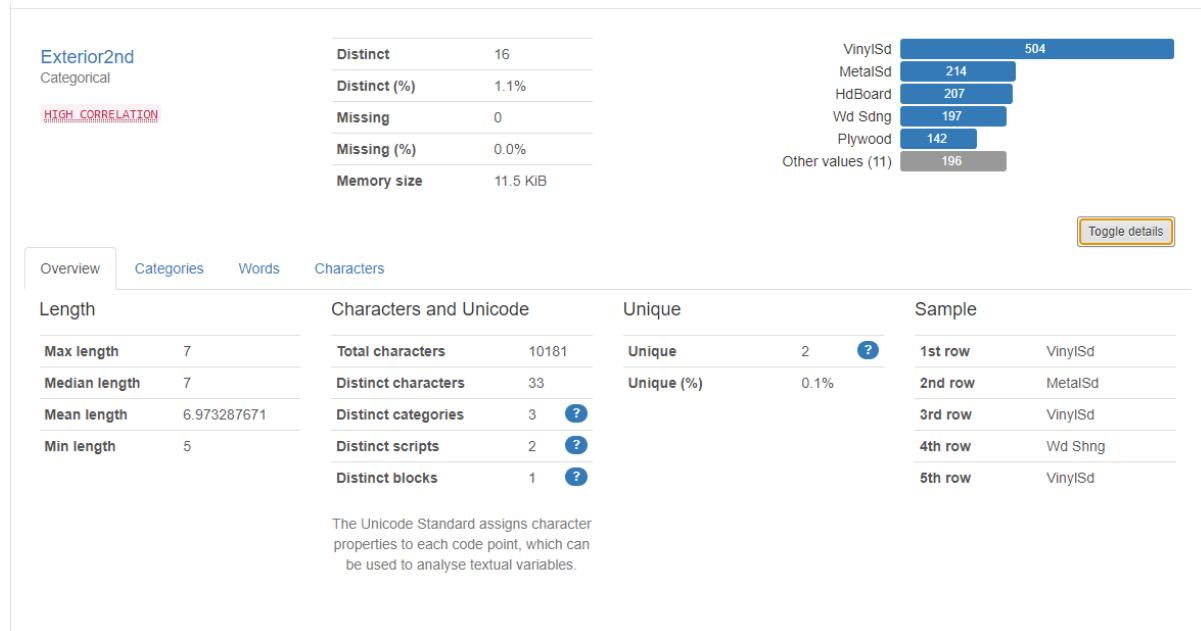
## Raul Angel Jimenez Hernandez 19017



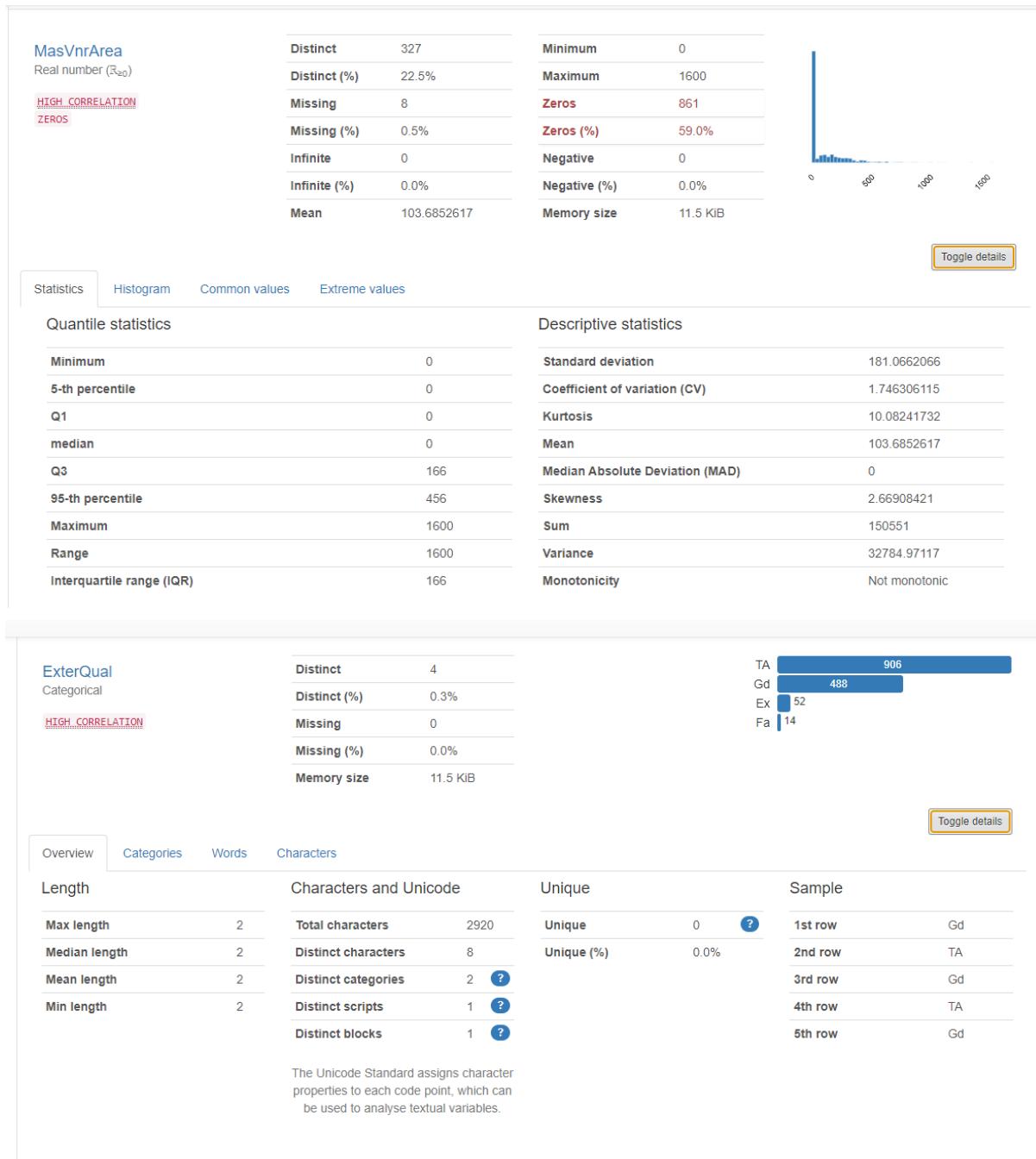
Donaldo Sebastian Garcia Jiménez 19683  
 Raul Angel Jimenez Hernandez 19017



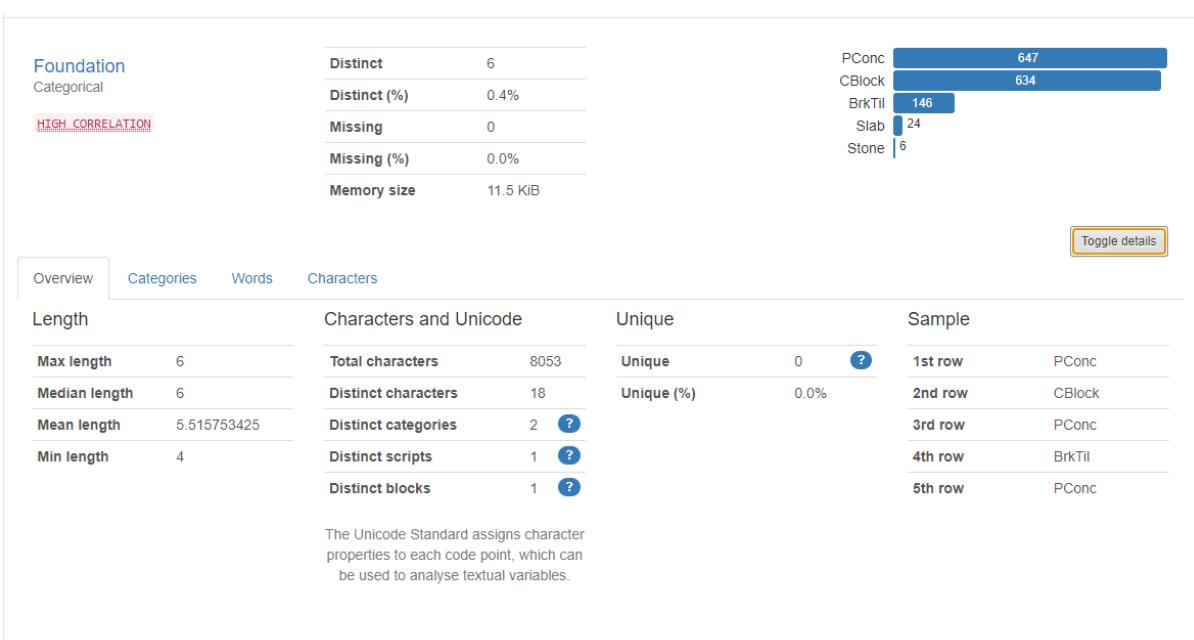
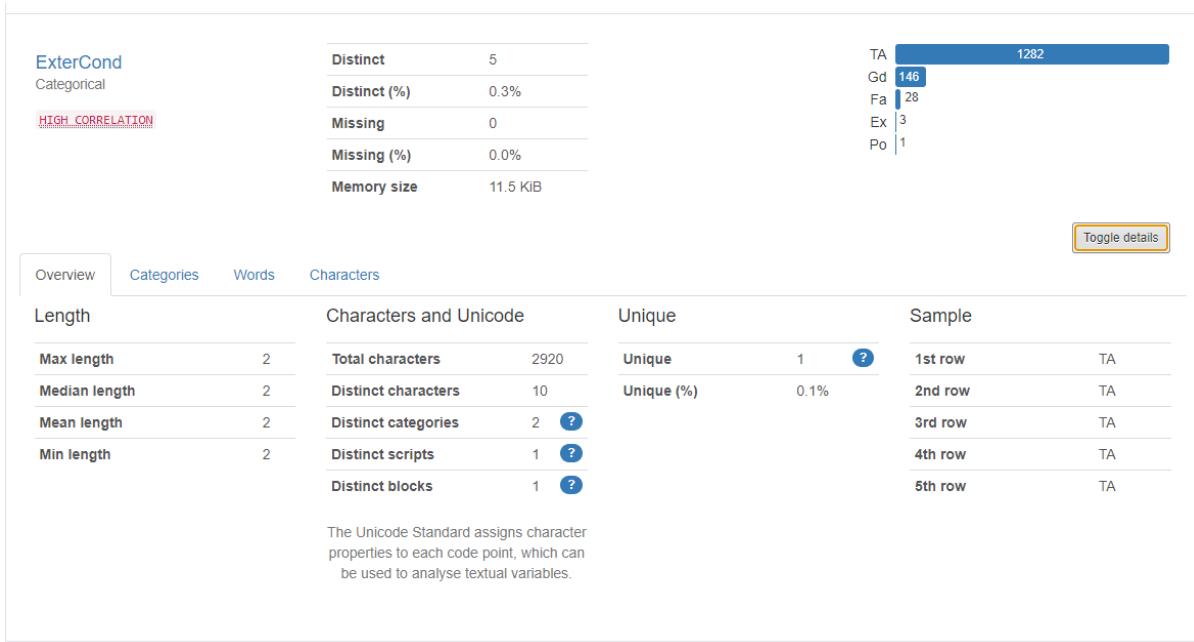
Donald Sebastian Garcia Jiménez 19683  
 Raul Angel Jimenez Hernandez 19017



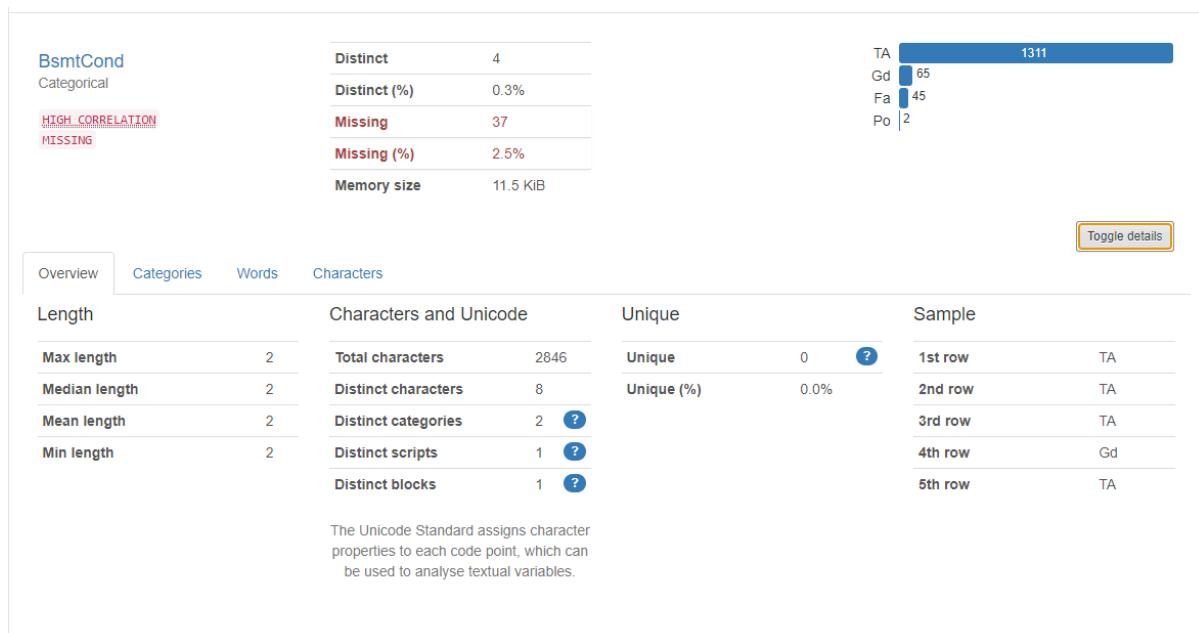
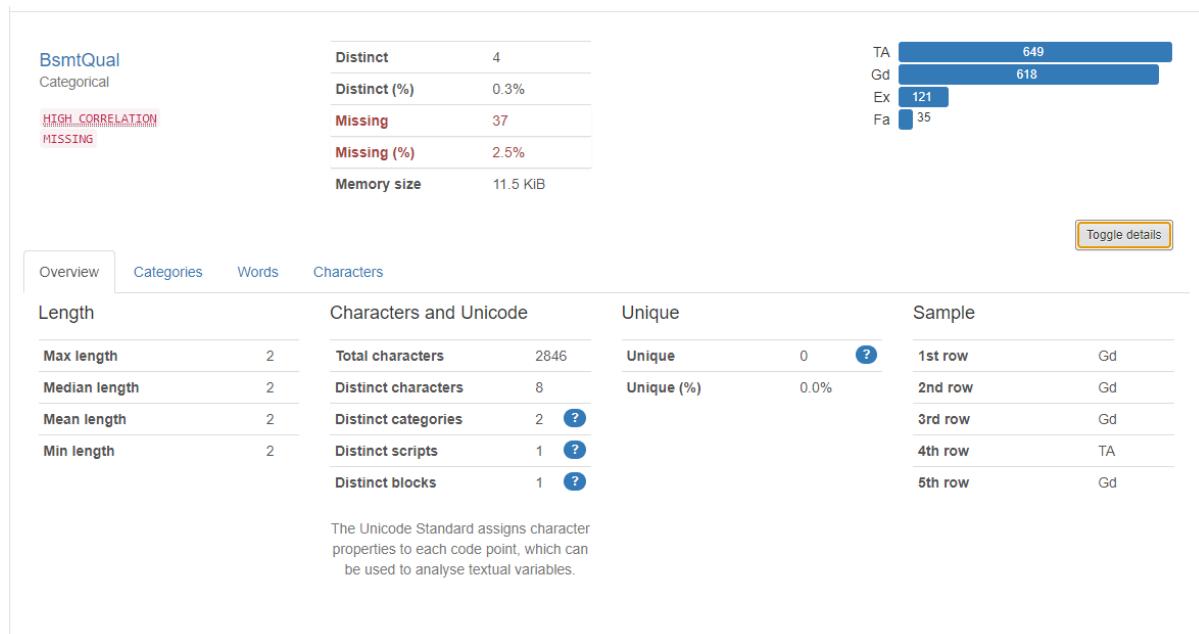
Donaldo Sebastian Garcia Jiménez 19683  
 Raul Angel Jimenez Hernandez 19017



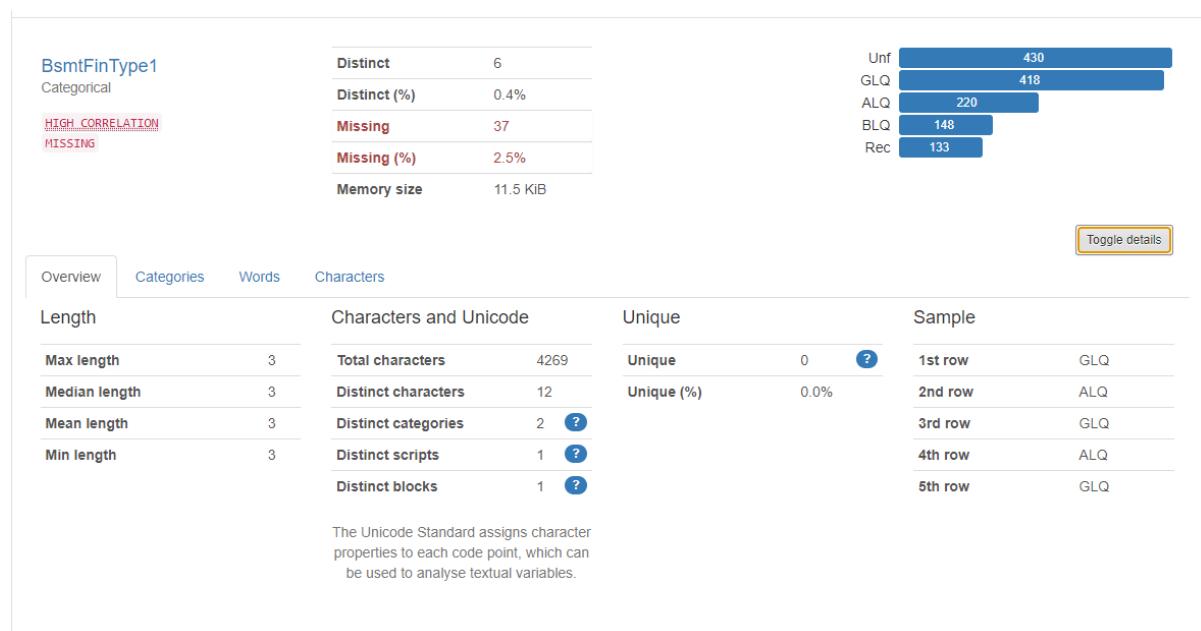
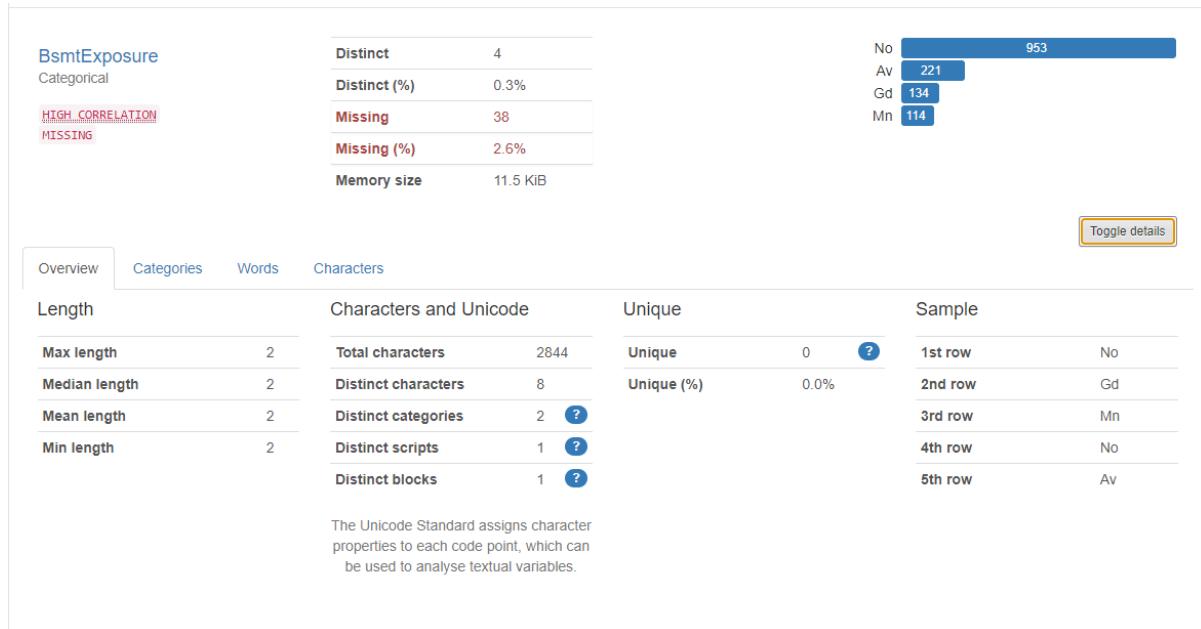
Donaldo Sebastian Garcia Jiménez 19683  
 Raul Angel Jimenez Hernandez 19017



Donaldo Sebastian Garcia Jiménez 19683  
 Raul Angel Jimenez Hernandez 19017

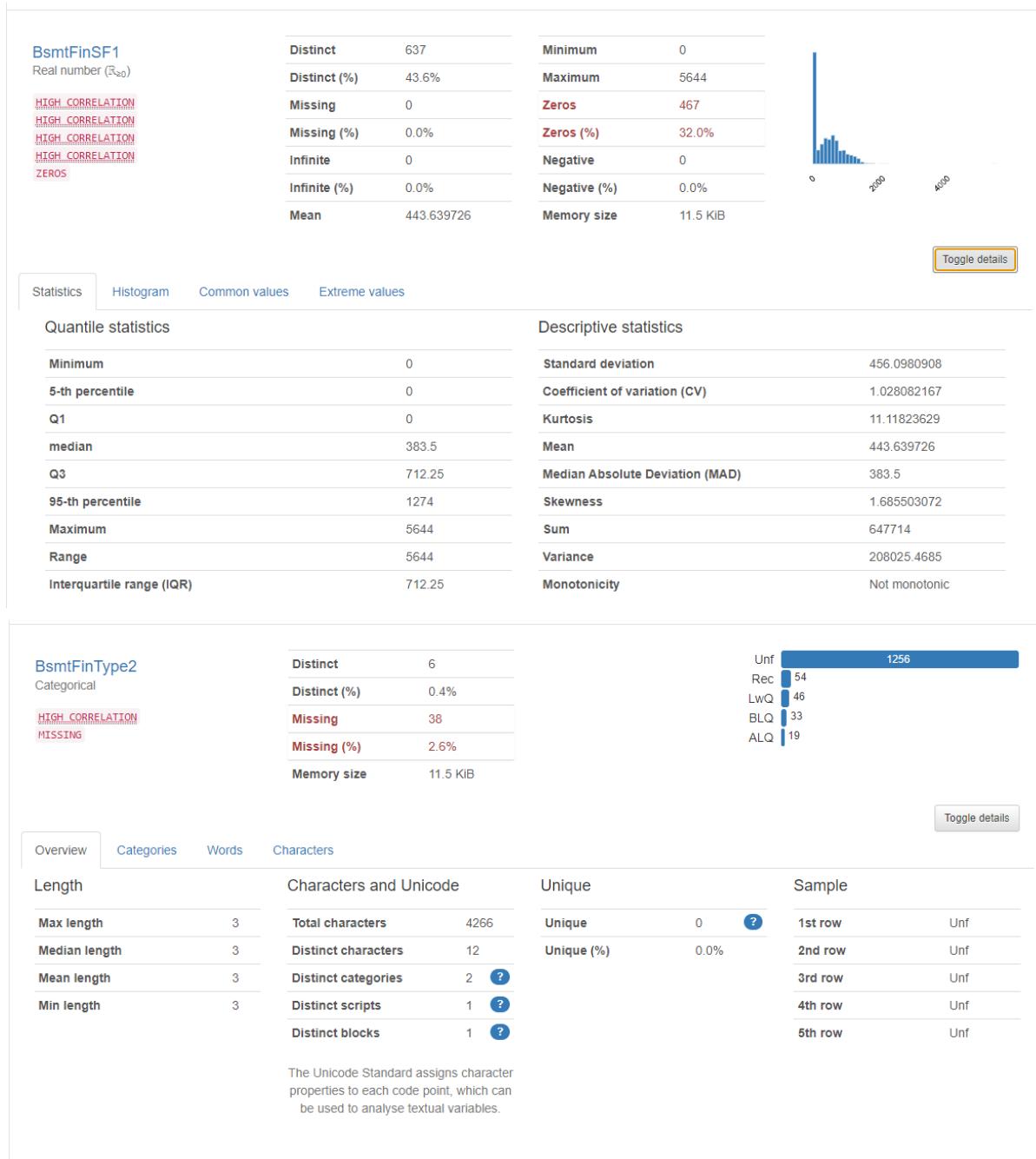


Donaldo Sebastian Garcia Jiménez 19683  
 Raul Angel Jimenez Hernandez 19017

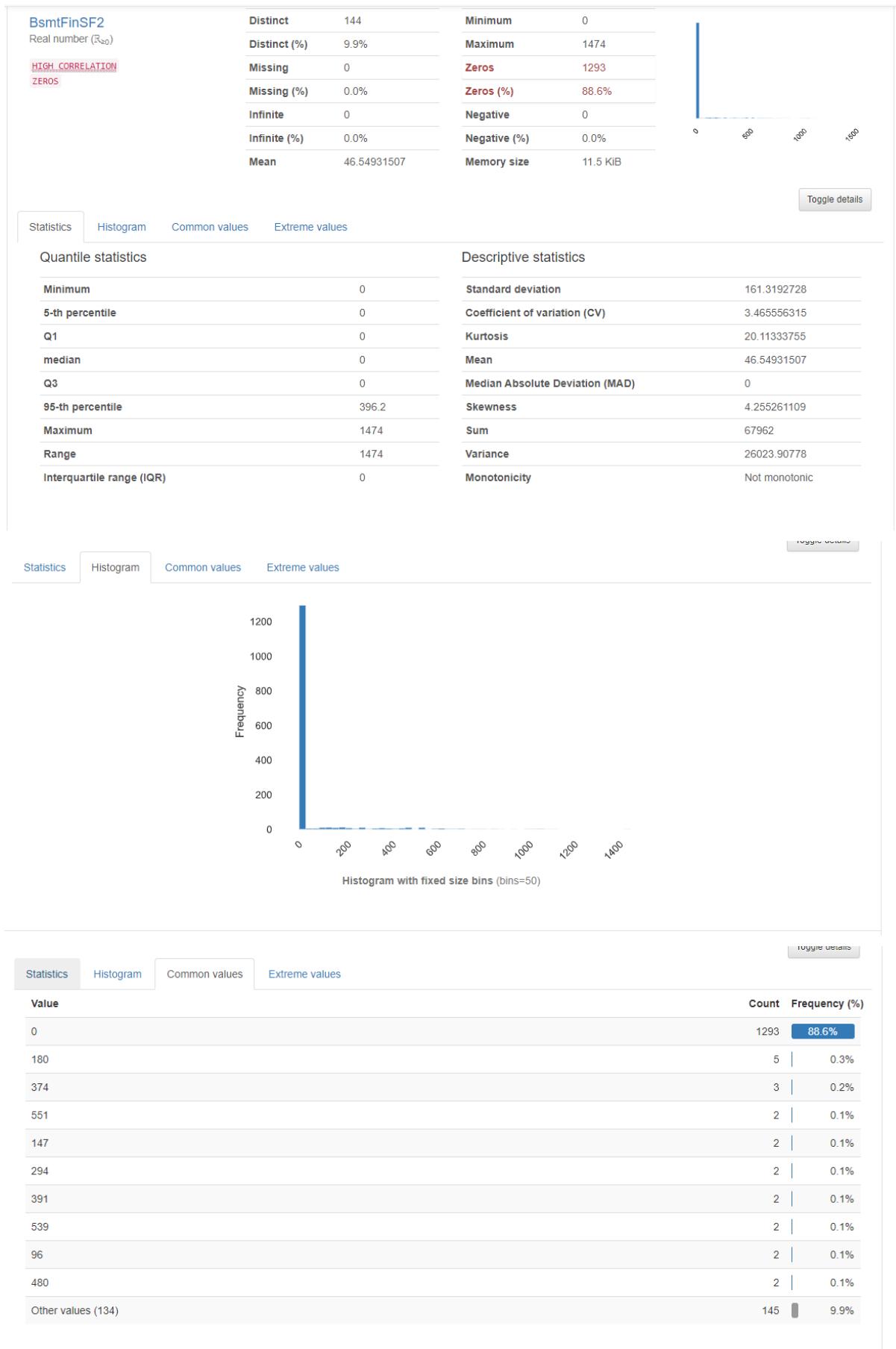


# Donaldo Sebastian Garcia Jiménez 19683

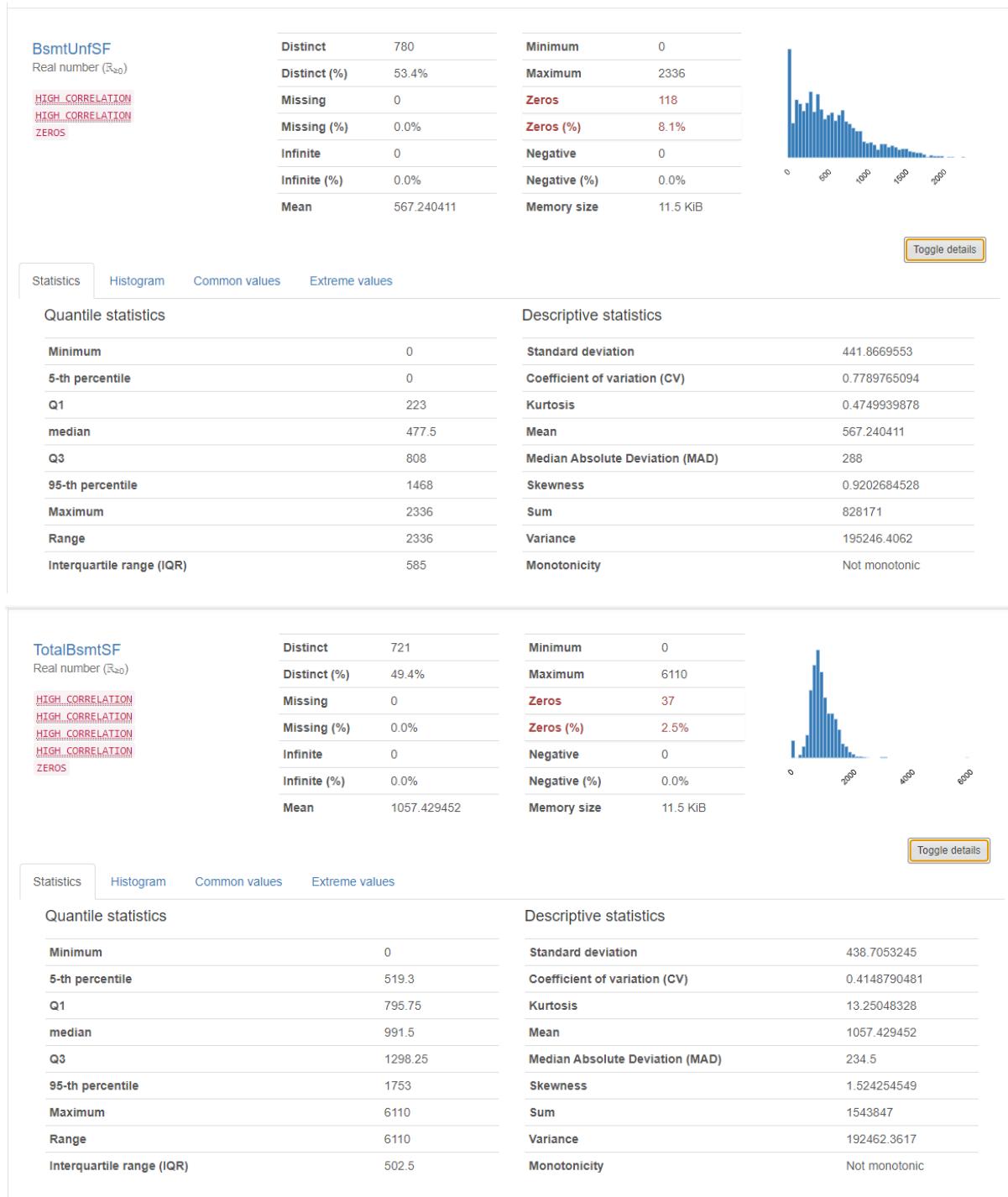
## Raul Angel Jimenez Hernandez 19017



Donaldo Sebastian Garcia Jiménez 19683  
 Raul Angel Jimenez Hernandez 19017

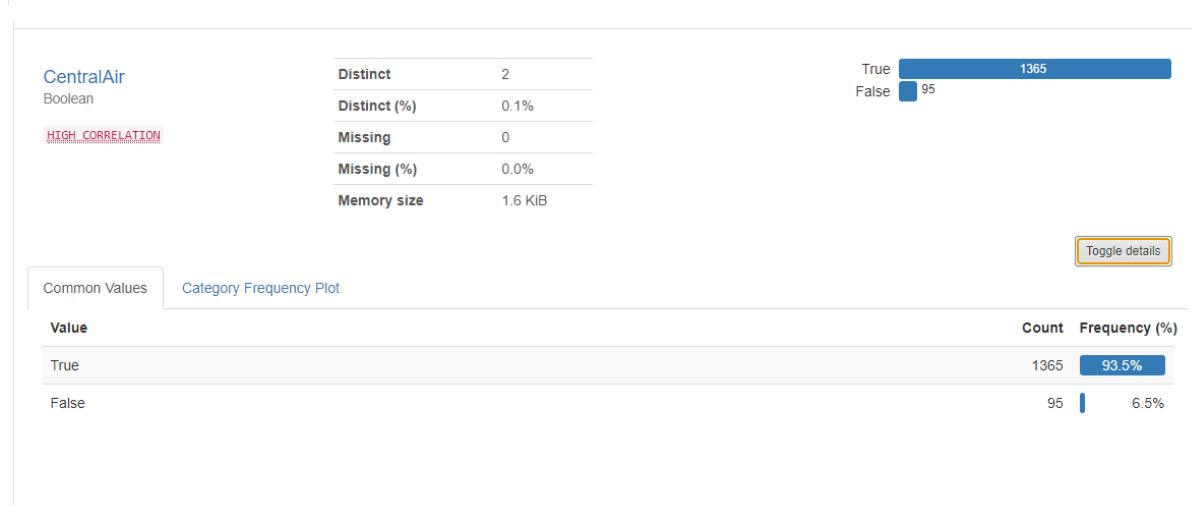
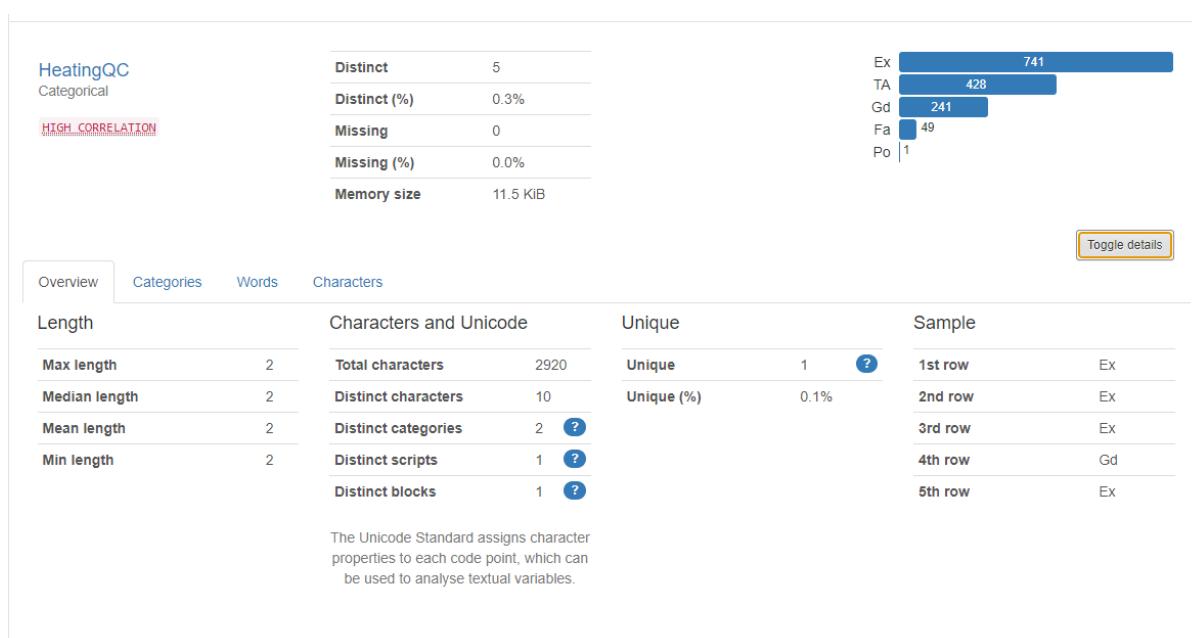
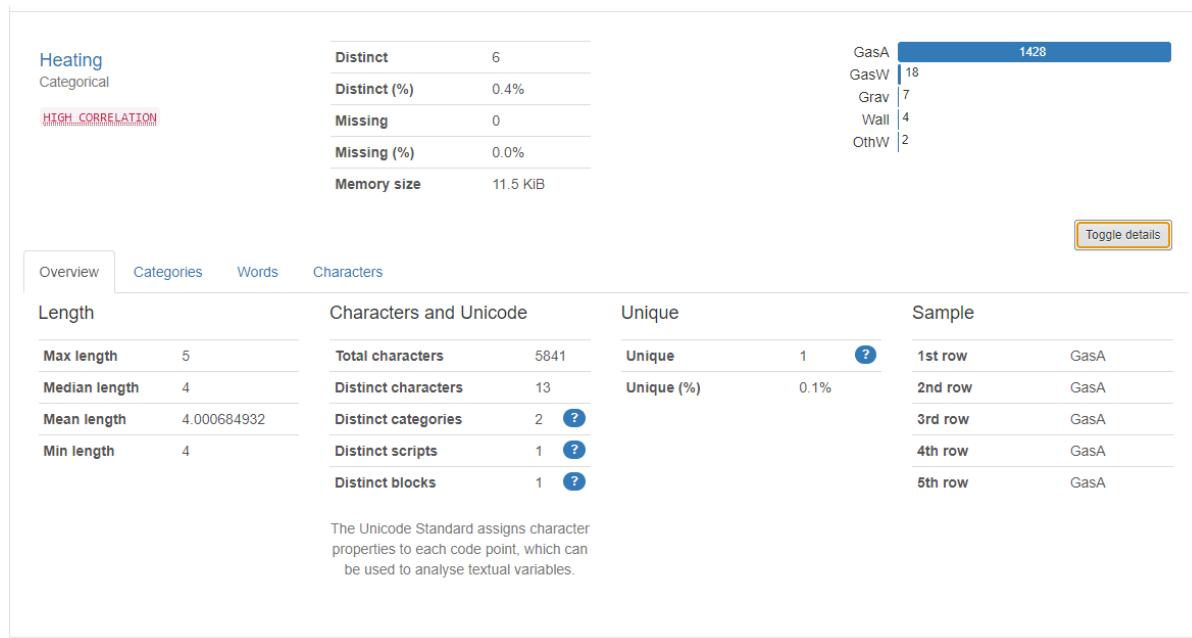


Donaldo Sebastian Garcia Jiménez 19683  
 Raul Angel Jimenez Hernandez 19017

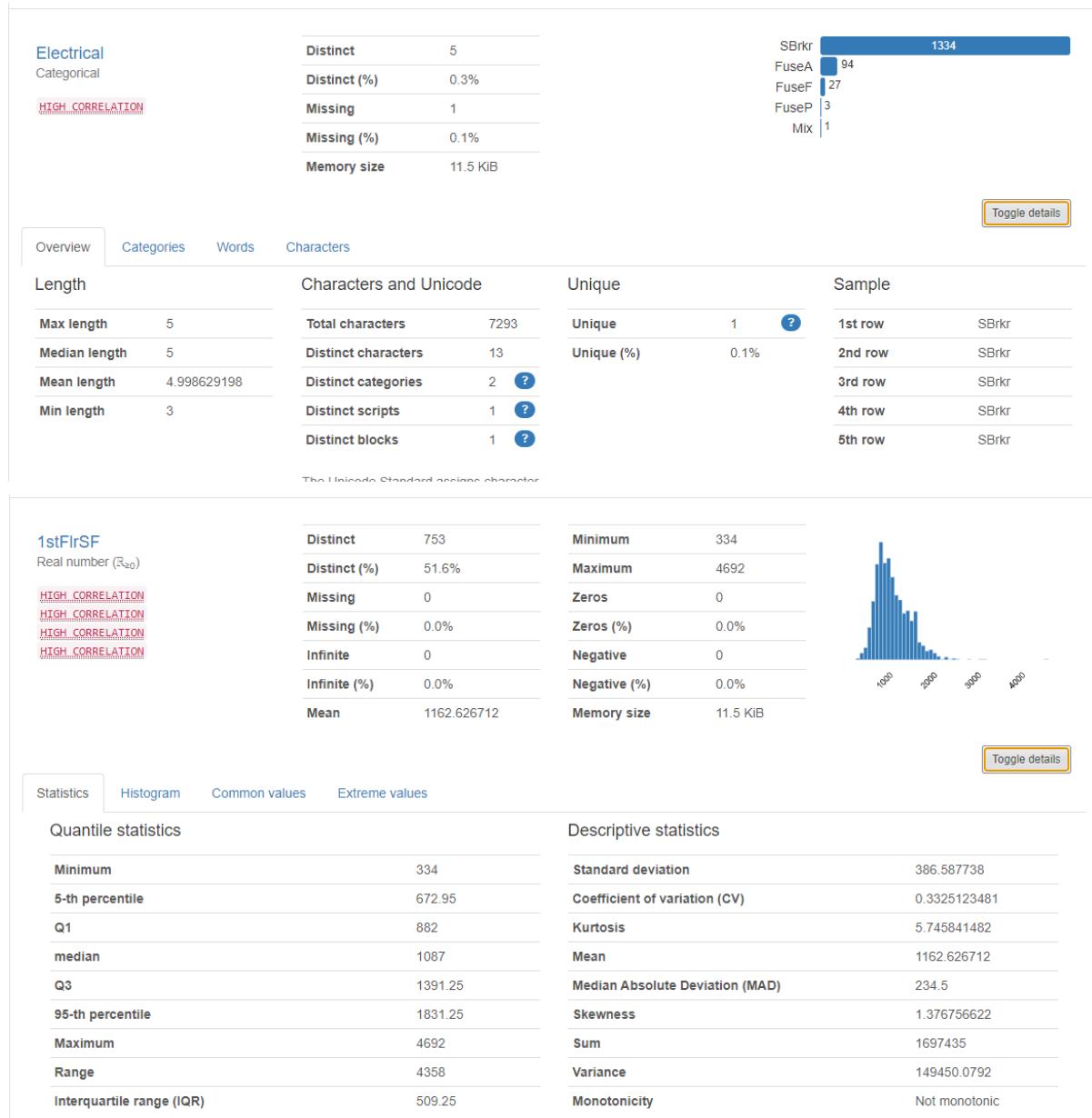


# Donald Sebastian Garcia Jiménez 19683

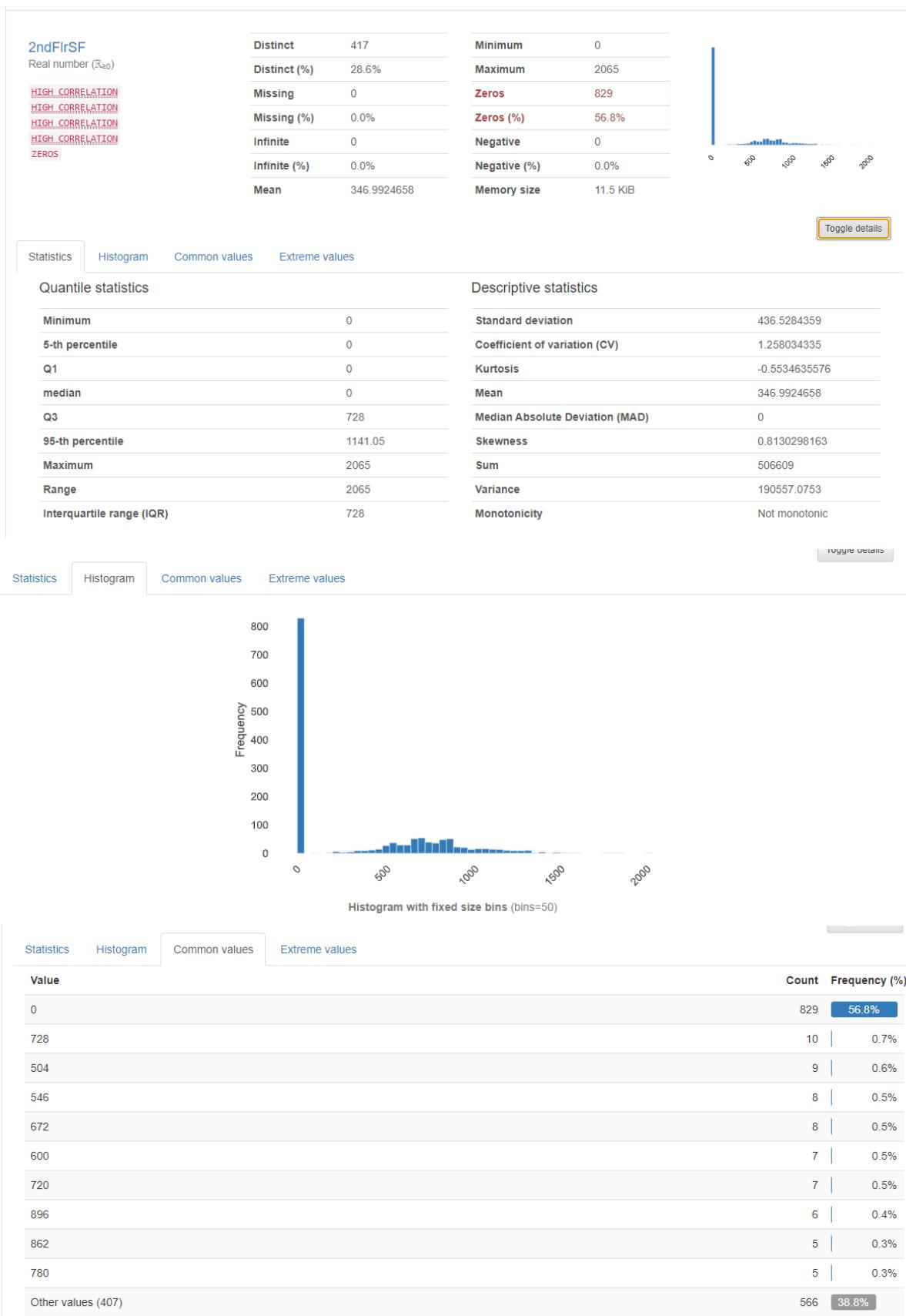
## Raul Angel Jimenez Hernandez 19017



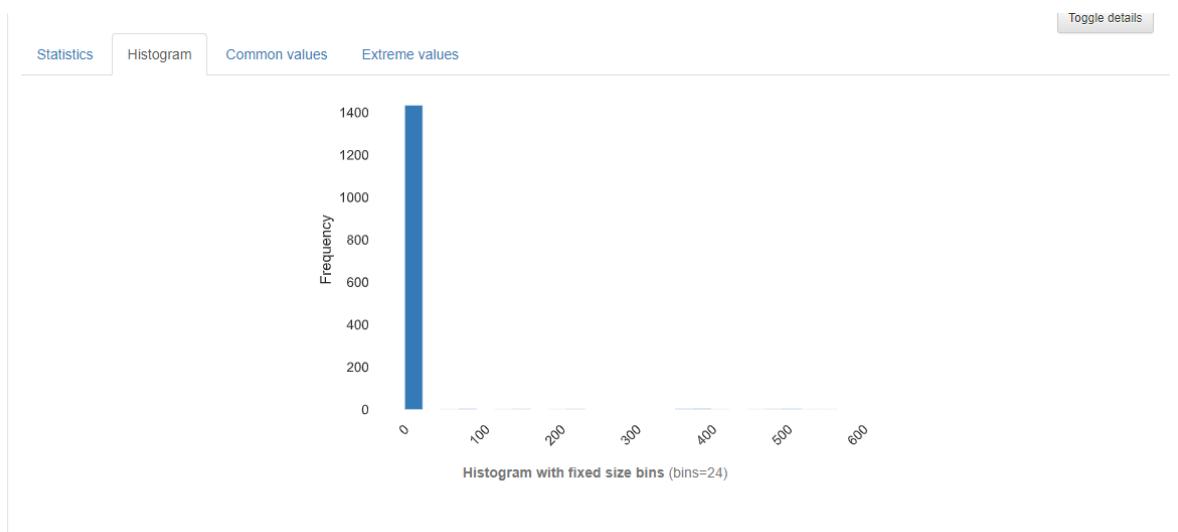
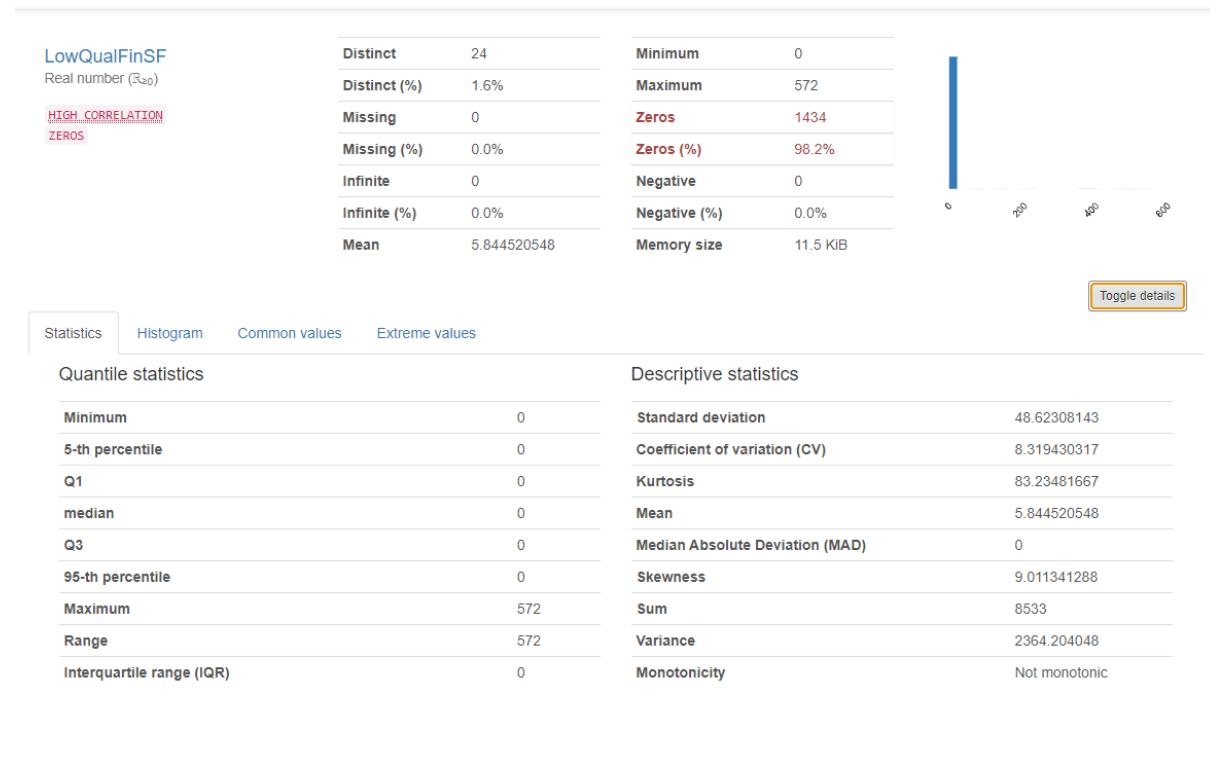
Donaldo Sebastian Garcia Jiménez 19683  
 Raul Angel Jimenez Hernandez 19017



Donaldo Sebastian Garcia Jiménez 19683  
 Raul Angel Jimenez Hernandez 19017



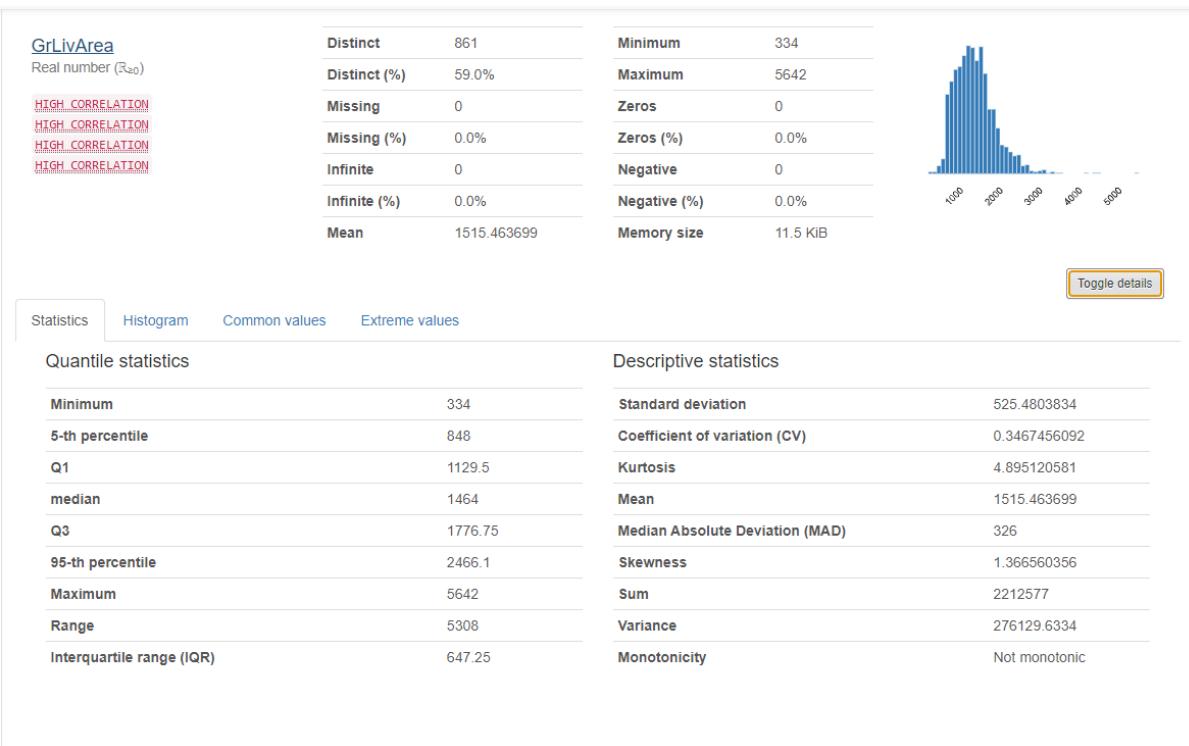
Donaldo Sebastian Garcia Jiménez 19683  
Raul Angel Jimenez Hernandez 19017



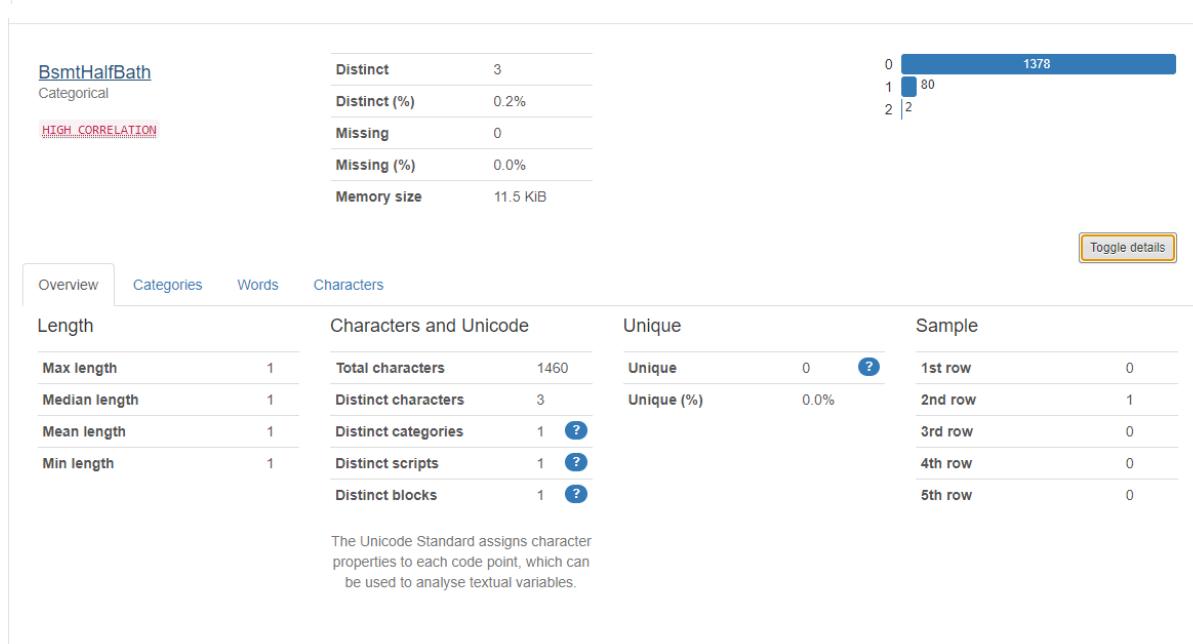
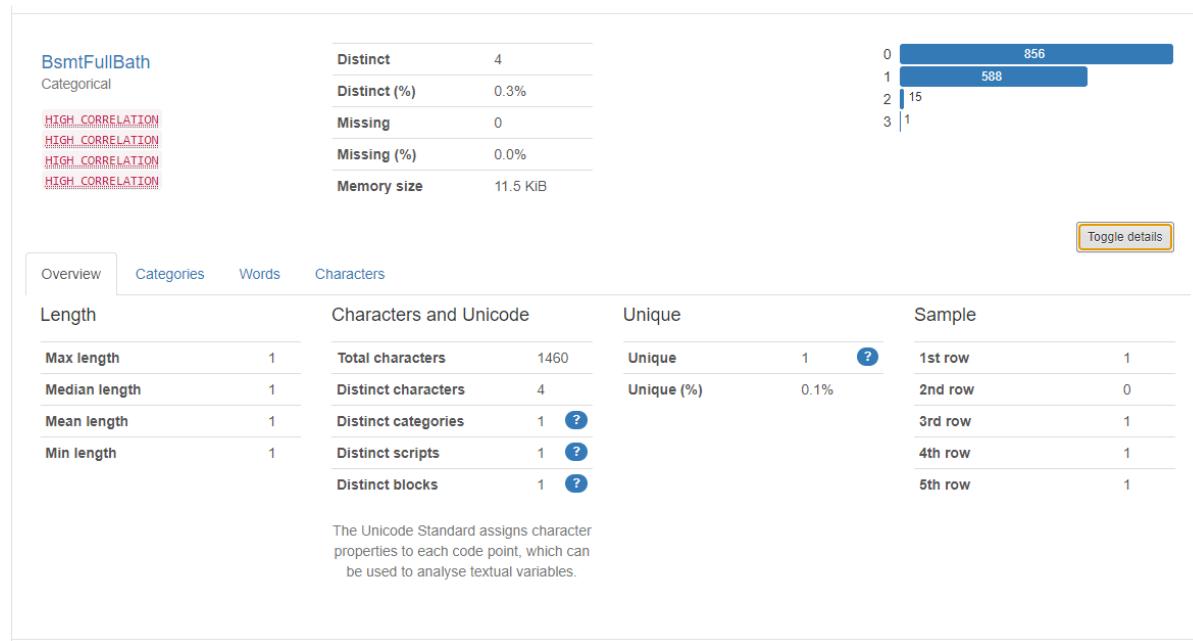
# Donaldo Sebastian Garcia Jiménez 19683

## Raul Angel Jimenez Hernandez 19017

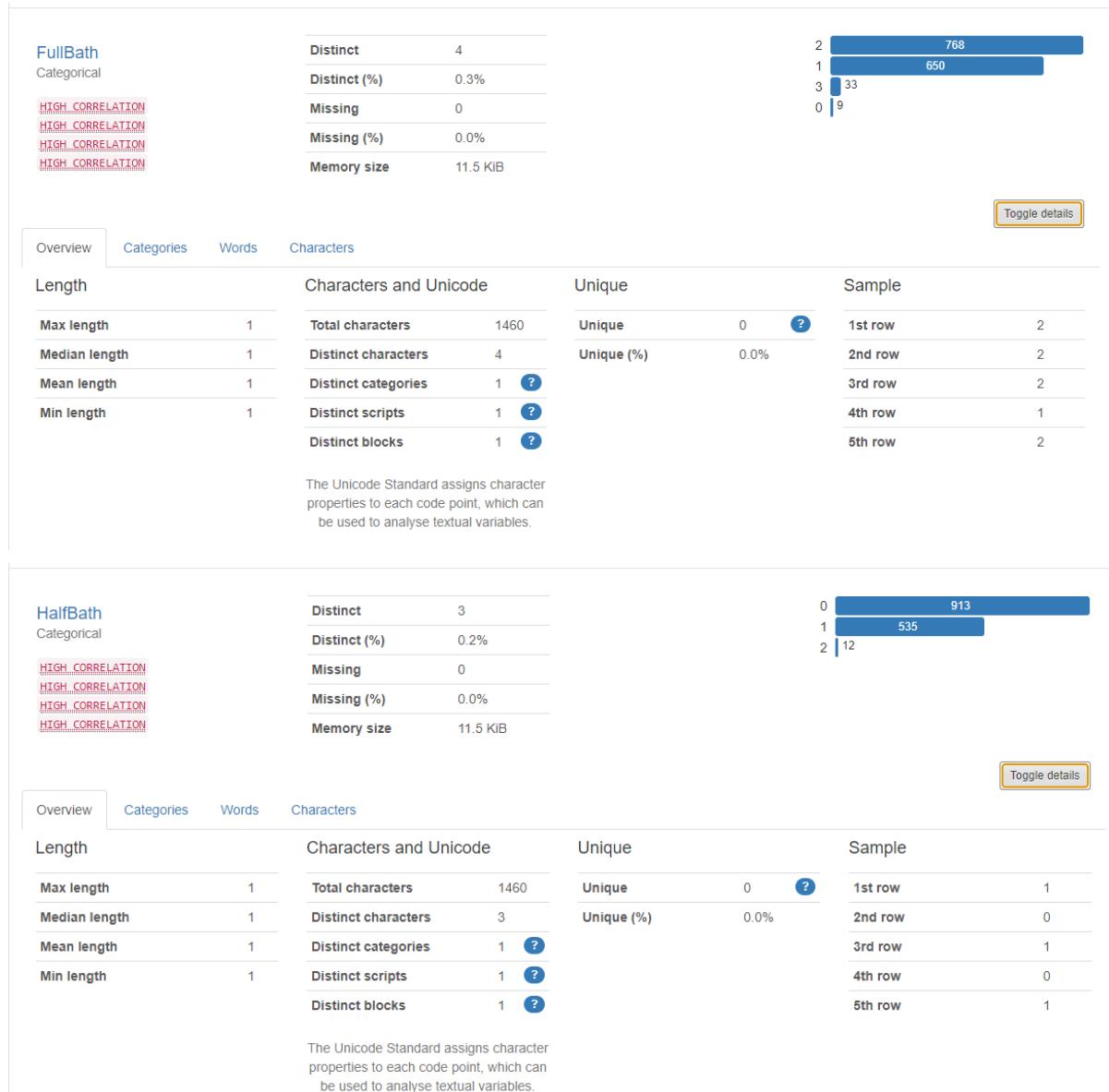
Value	Count	Frequency (%)
0	1434	98.2%
80	3	0.2%
360	2	0.1%
205	1	0.1%
479	1	0.1%
397	1	0.1%
514	1	0.1%
120	1	0.1%
481	1	0.1%
232	1	0.1%
Other values (14)	14	1.0%



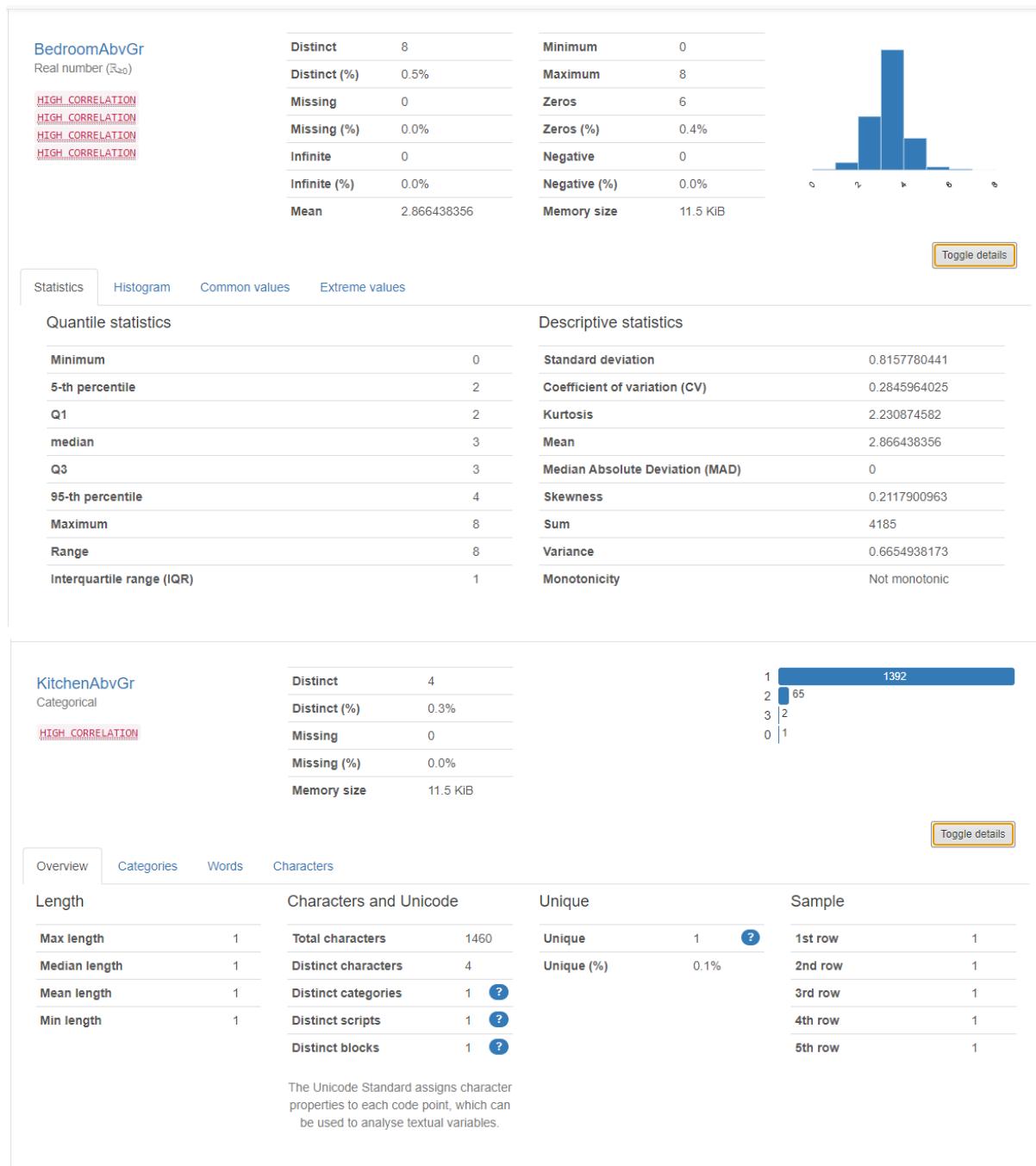
Donaldo Sebastian Garcia Jiménez 19683  
 Raul Angel Jimenez Hernandez 19017



Donaldo Sebastian Garcia Jiménez 19683  
 Raul Angel Jimenez Hernandez 19017



Donaldo Sebastian Garcia Jiménez 19683  
 Raul Angel Jimenez Hernandez 19017



Donaldo Sebastian Garcia Jiménez 19683  
 Raul Angel Jimenez Hernandez 19017

**KitchenQual**  
Categorical

**HIGH CORRELATION**

Distinct	4
Distinct (%)	0.3%
Missing	0
Missing (%)	0.0%
Memory size	11.5 KiB

TA 735  
Gd 586  
Ex 100  
Fa 39

Overview
Categories
Words
Characters
Toggle details

Length		Characters and Unicode		Unique	Sample
Max length	2	Total characters	2920	Unique	0
Median length	2	Distinct characters	8	Unique (%)	0.0%
Mean length	2	Distinct categories	2		
Min length	2	Distinct scripts	1		
		Distinct blocks	1		

The Unicode Standard assigns character properties to each code point, which can be used to analyse textual variables.

**TotRmsAbvGrd**  
Real number (ℝ₀)

**HIGH CORRELATION**

**HIGH CORRELATION**

**HIGH CORRELATION**

**HIGH CORRELATION**

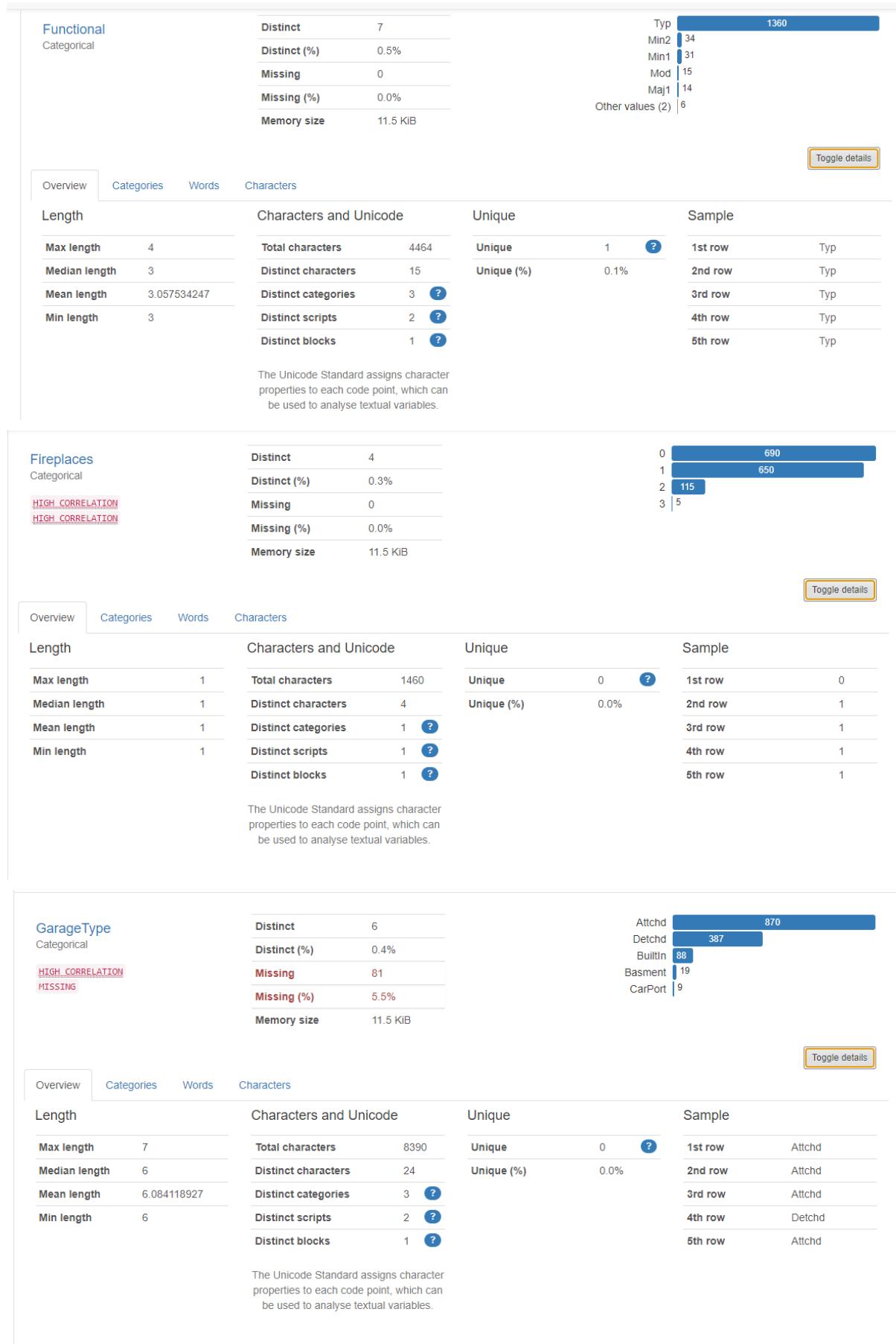
Distinct	12
Distinct (%)	0.8%
Missing	0
Missing (%)	0.0%
Infinite	0
Infinite (%)	0.0%
Mean	6.517808219
Minimum	2
Maximum	14
Zeros	0
Zeros (%)	0.0%
Negative	0
Negative (%)	0.0%
Memory size	11.5 KiB

Toggle details

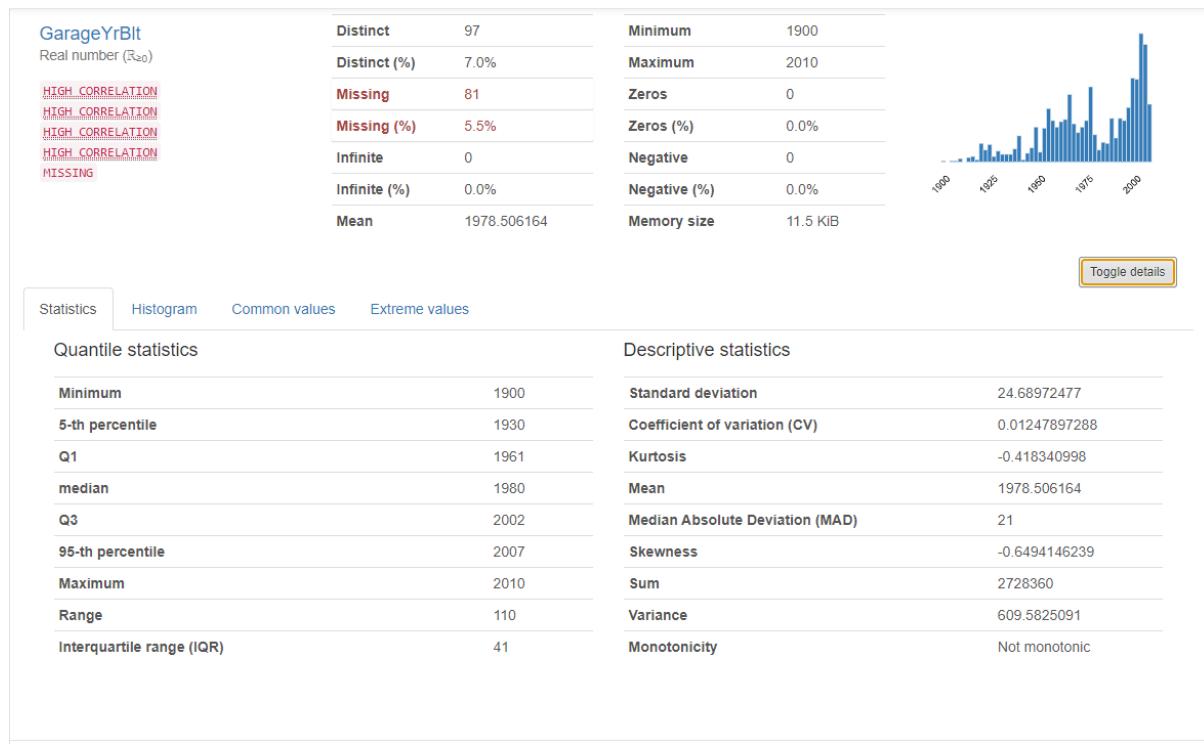
Statistics
Histogram
Common values
Extreme values
Toggle details

Quantile statistics		Descriptive statistics	
Minimum	2	Standard deviation	1.625393291
5-th percentile	4	Coefficient of variation (CV)	0.2493772808
Q1	5	Kurtosis	0.8807615657
median	6	Mean	6.517808219
Q3	7	Median Absolute Deviation (MAD)	1
95-th percentile	10	Skewness	0.6763408364
Maximum	14	Sum	9516
Range	12	Variance	2.641903349
Interquartile range (IQR)	2	Monotonicity	Not monotonic

Donald Sebastian Garcia Jiménez 19683  
 Raul Angel Jimenez Hernandez 19017



Donaldo Sebastian Garcia Jiménez 19683  
 Raul Angel Jimenez Hernandez 19017



Statistics    Histogram    Common values    Extreme values

Quantile statistics

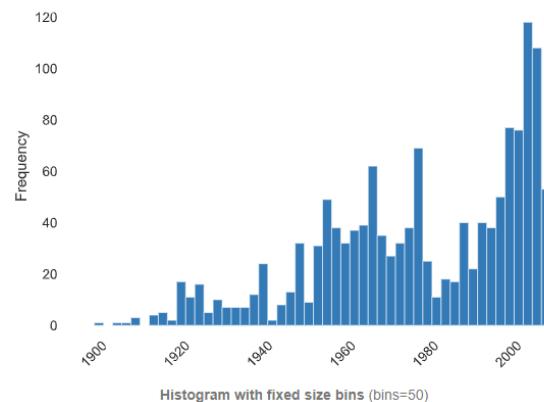
Minimum	1900
5-th percentile	1930
Q1	1961
median	1980
Q3	2002
95-th percentile	2007
Maximum	2010
Range	110
Interquartile range (IQR)	41

Descriptive statistics

Standard deviation	24.68972477
Coefficient of variation (CV)	0.01247897288
Kurtosis	-0.418340998
Mean	1978.506164
Median Absolute Deviation (MAD)	21
Skewness	-0.6494146239
Sum	2728360
Variance	609.5825091
Monotonicity	Not monotonic

Statistics    Histogram    Common values    Extreme values

Toggle details

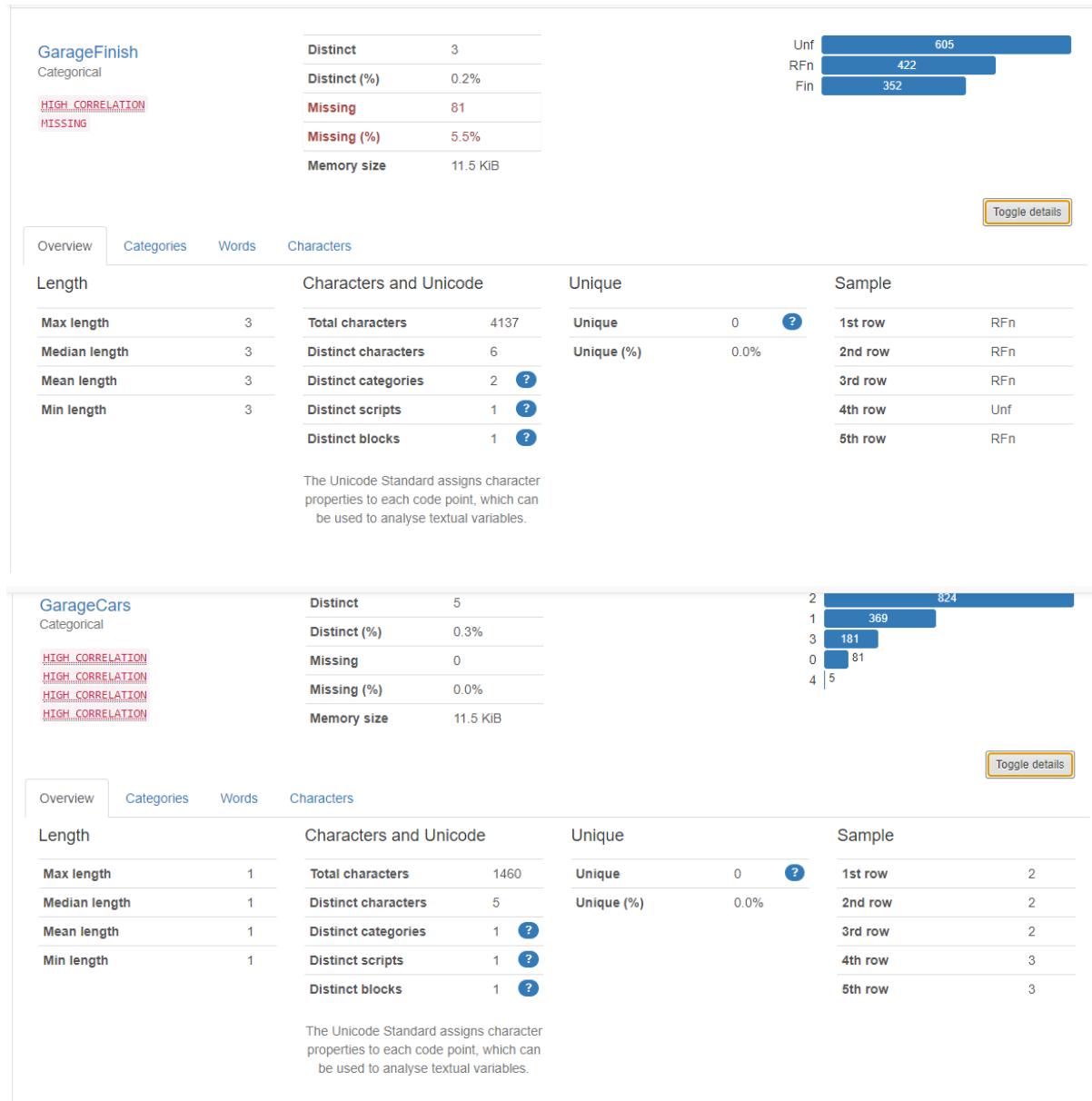


Statistics    Histogram    Common values    Extreme values

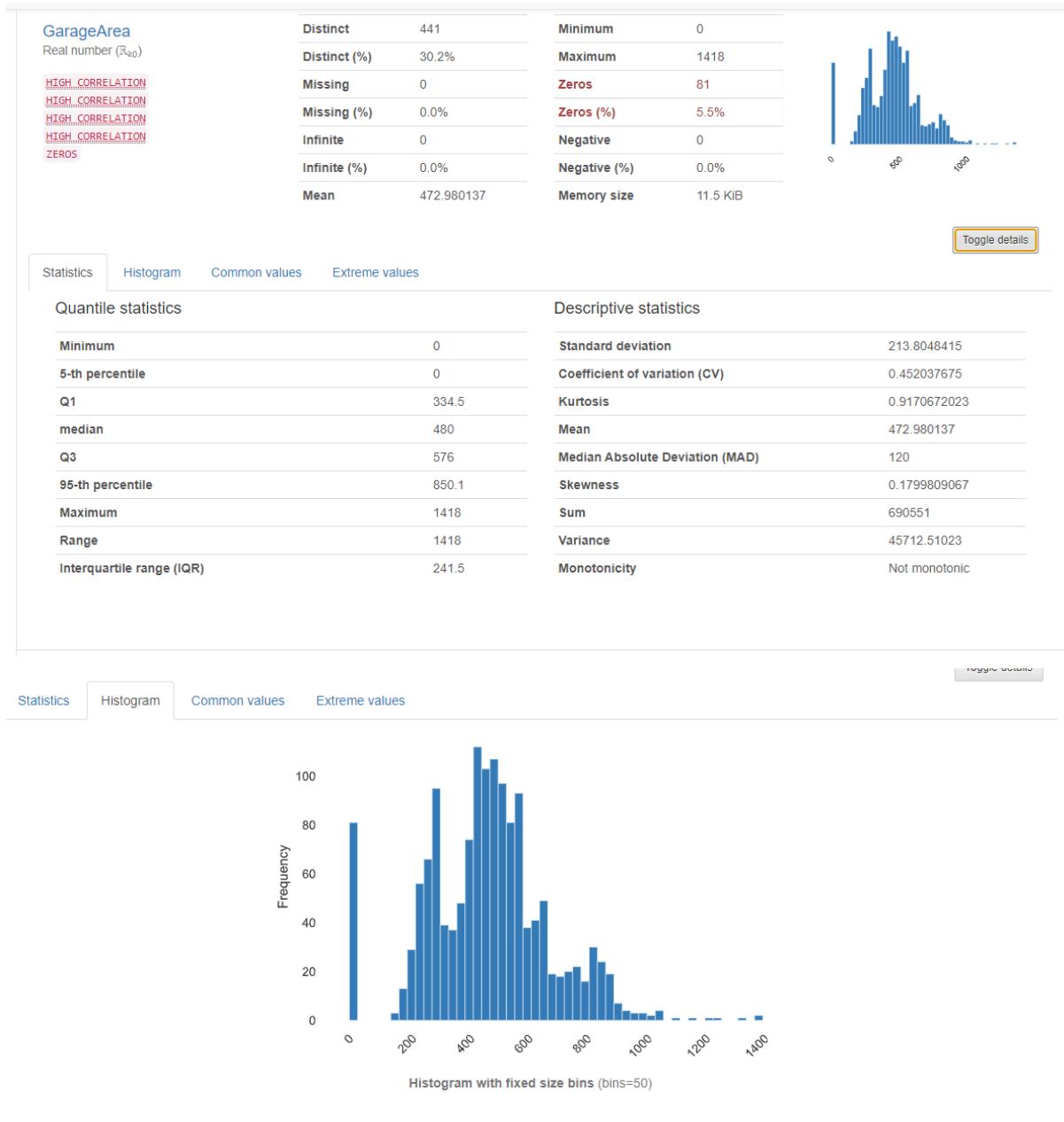
Value

Value	Count	Frequency (%)
2005	65	4.5%
2006	59	4.0%
2004	53	3.6%
2003	50	3.4%
2007	49	3.4%
1977	35	2.4%
1998	31	2.1%
1999	30	2.1%
1976	29	2.0%
2008	29	2.0%
Other values (87)	949	65.0%
(Missing)	81	5.5%

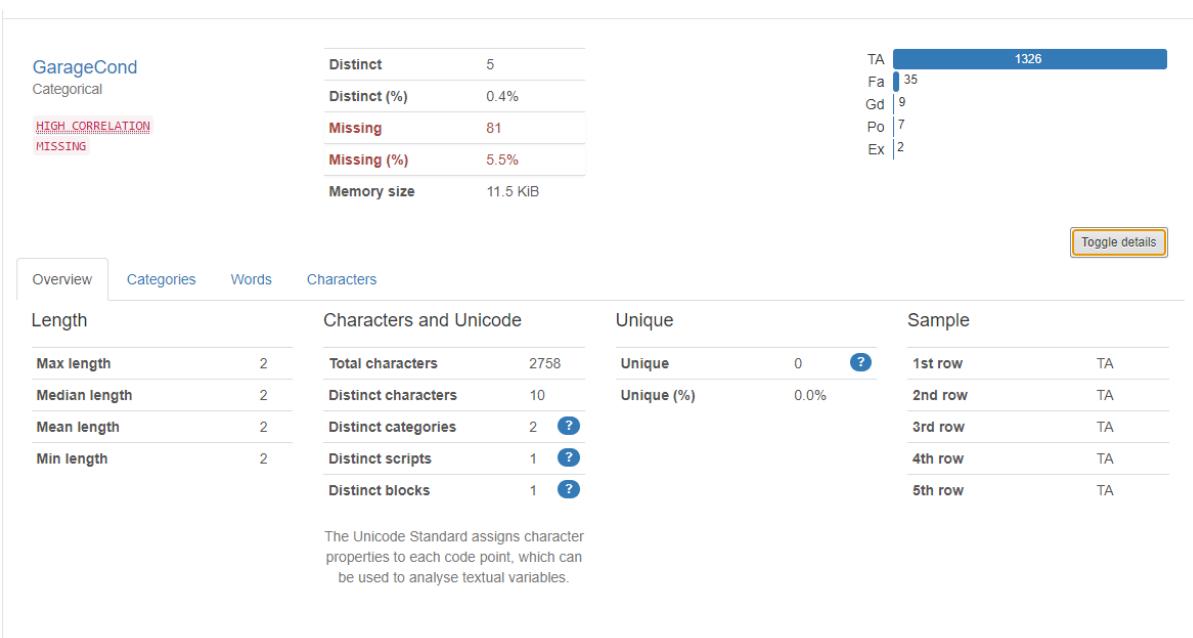
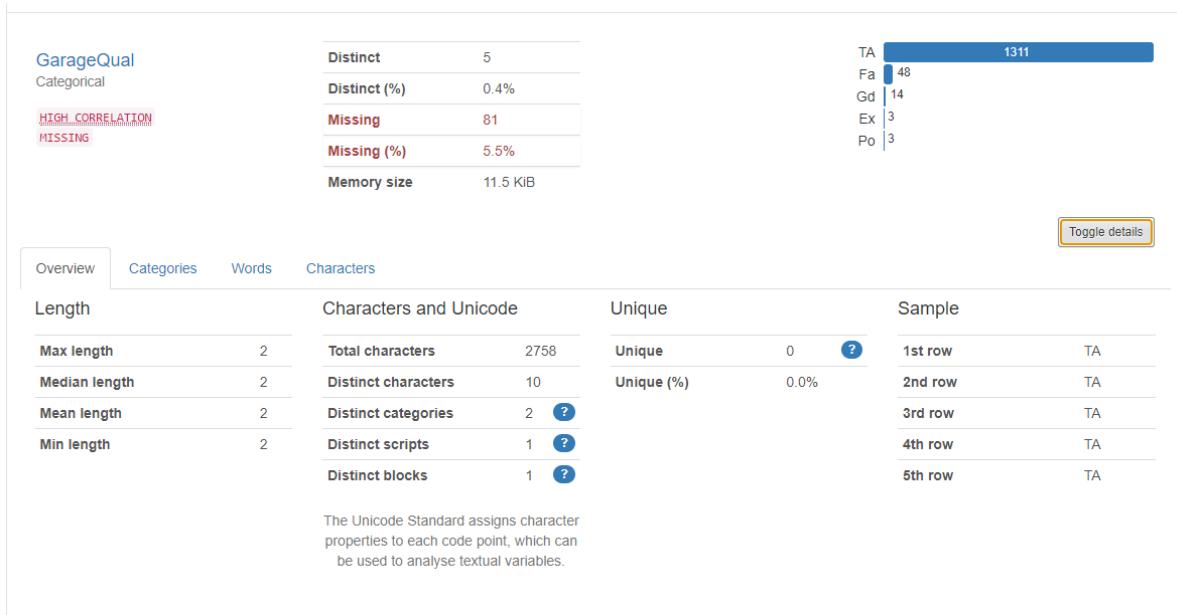
Donaldo Sebastian Garcia Jiménez 19683  
 Raul Angel Jimenez Hernandez 19017



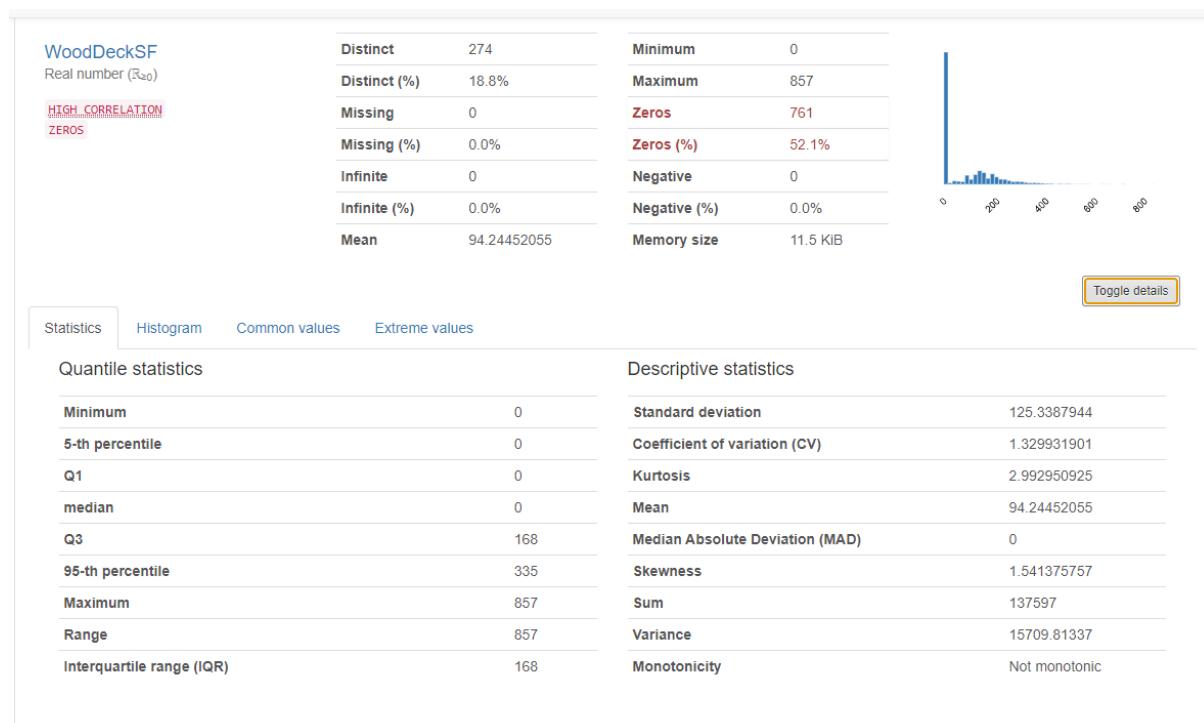
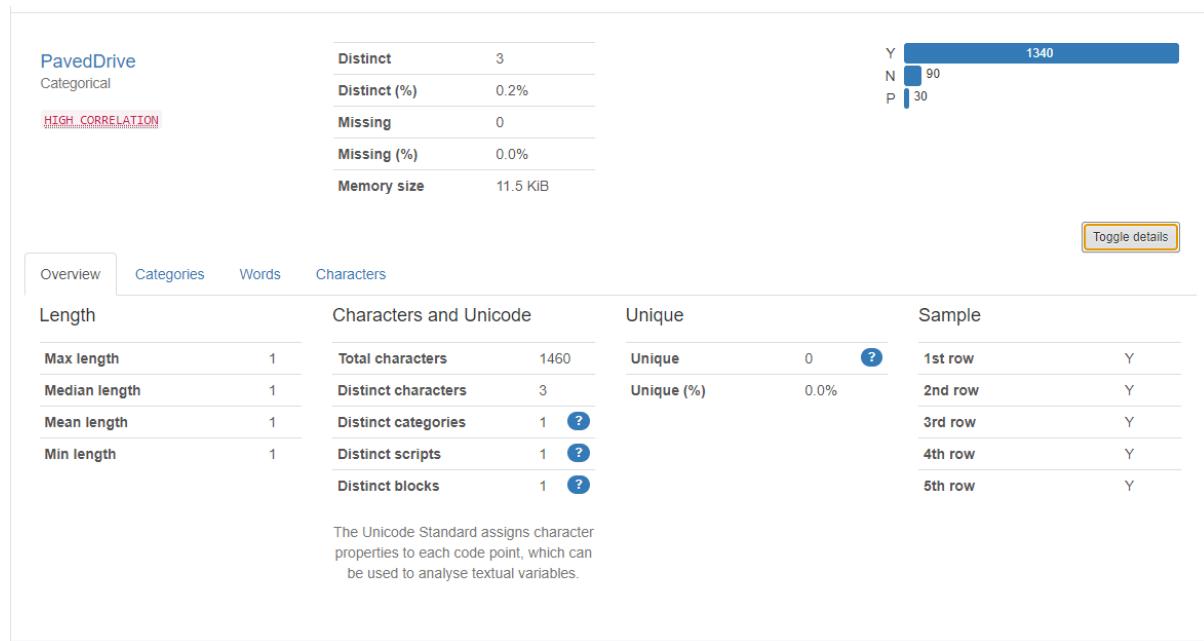
Donaldo Sebastian Garcia Jiménez 19683  
 Raul Angel Jimenez Hernandez 19017



Donaldo Sebastian Garcia Jiménez 19683  
 Raul Angel Jimenez Hernandez 19017

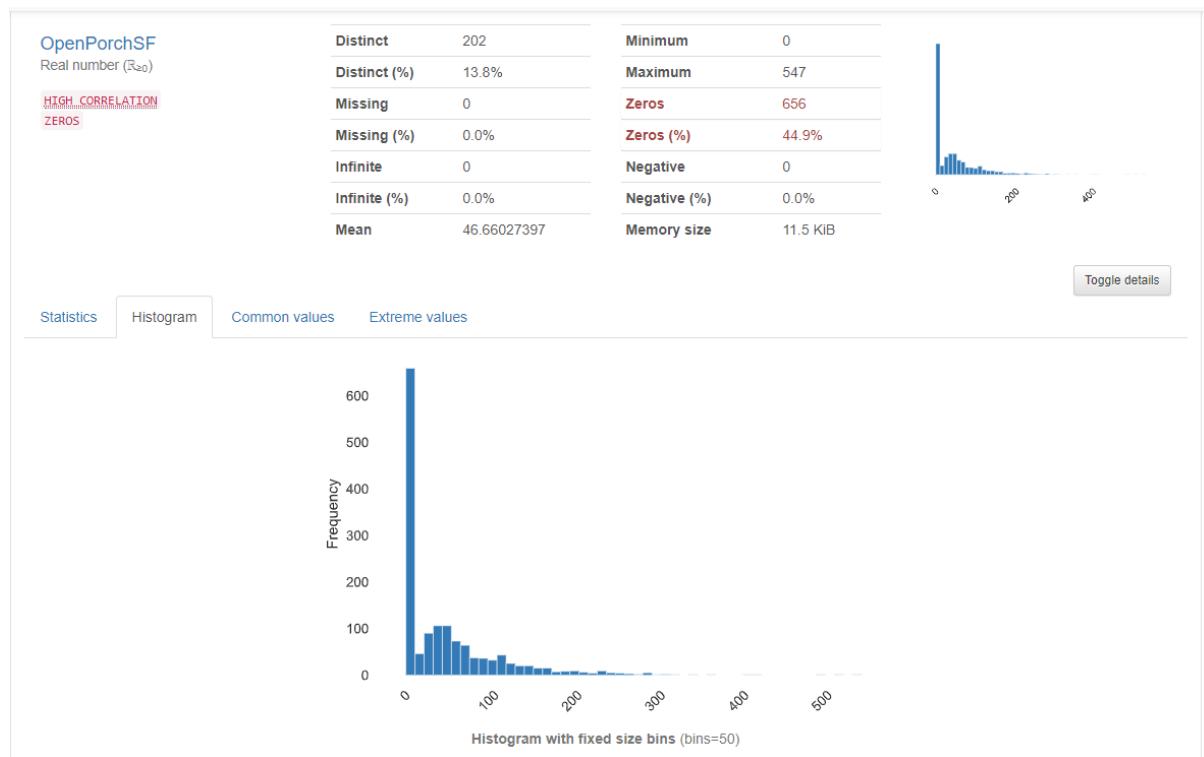
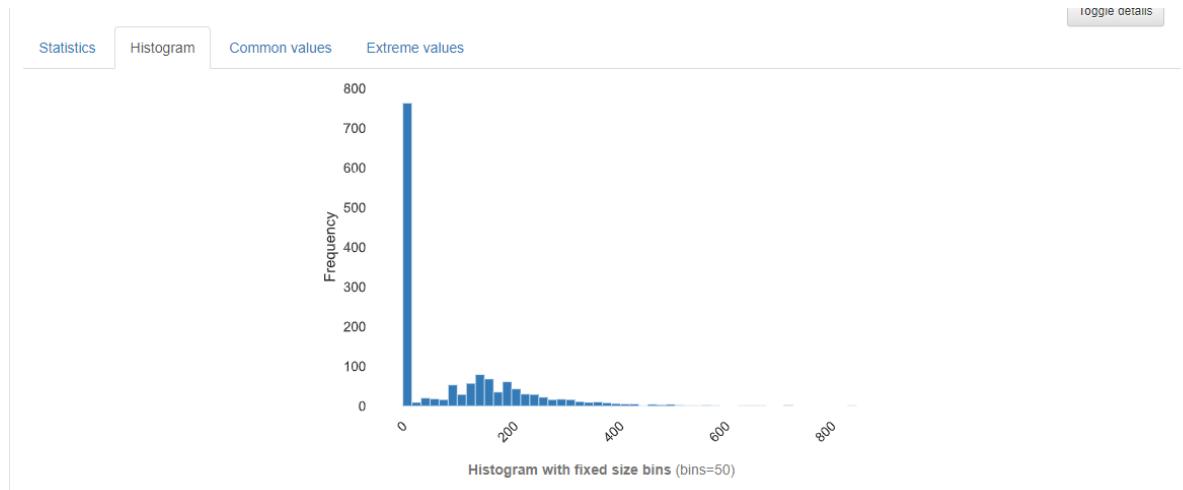


Donaldo Sebastian Garcia Jiménez 19683  
 Raul Angel Jimenez Hernandez 19017

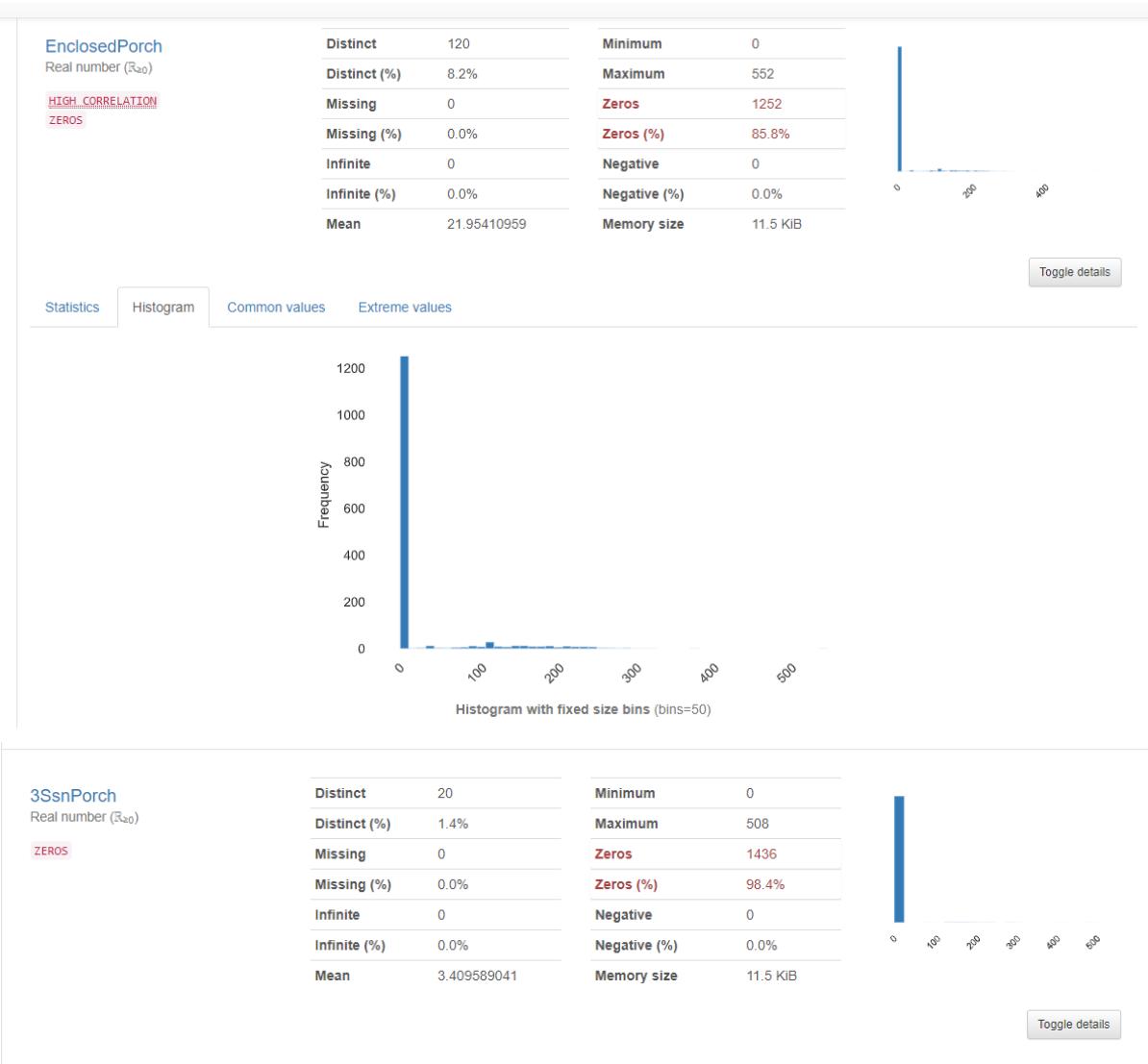


Donaldo Sebastian Garcia Jiménez 19683

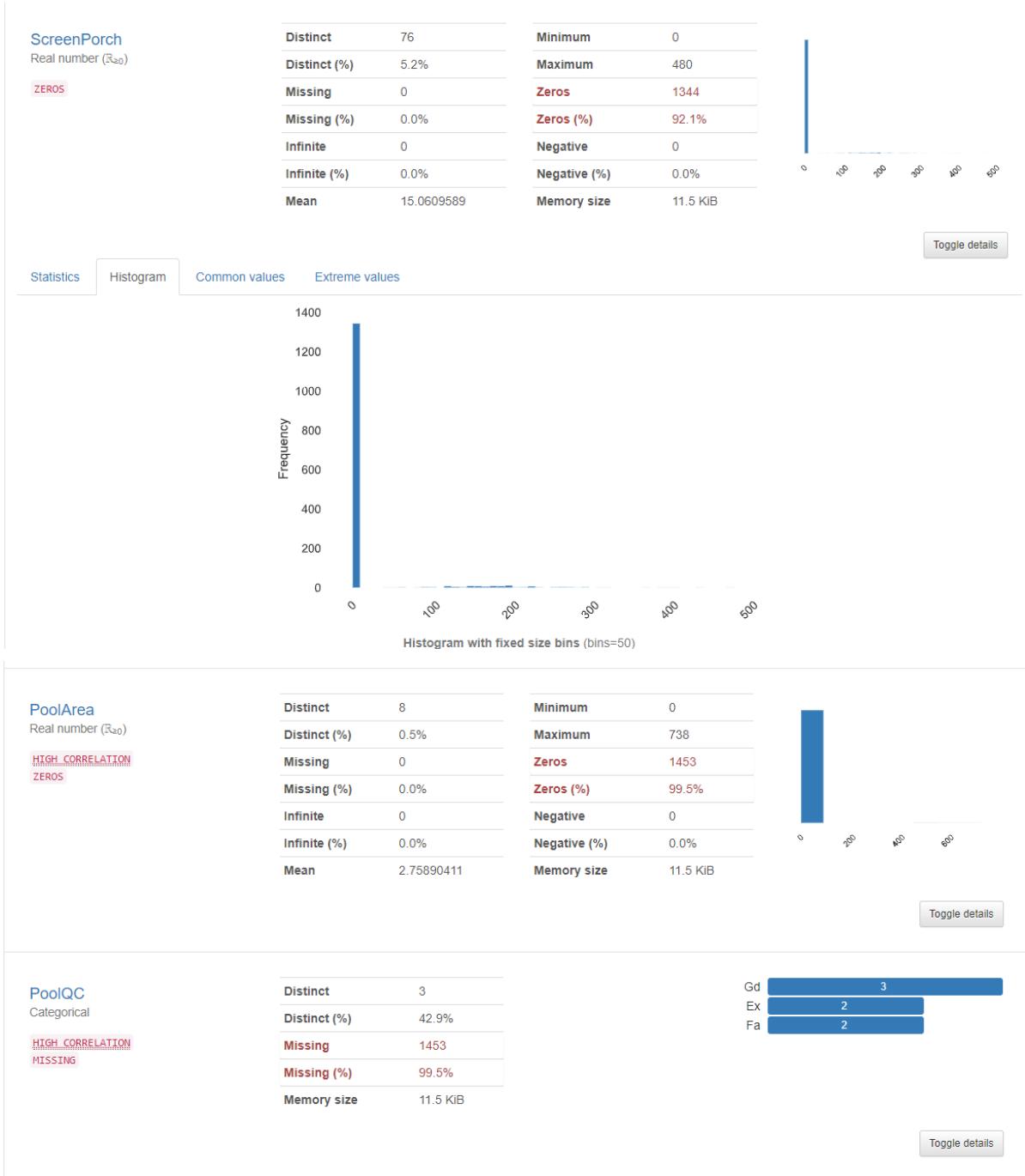
Raul Angel Jimenez Hernandez 19017



Donaldo Sebastian Garcia Jiménez 19683  
Raul Angel Jimenez Hernandez 19017

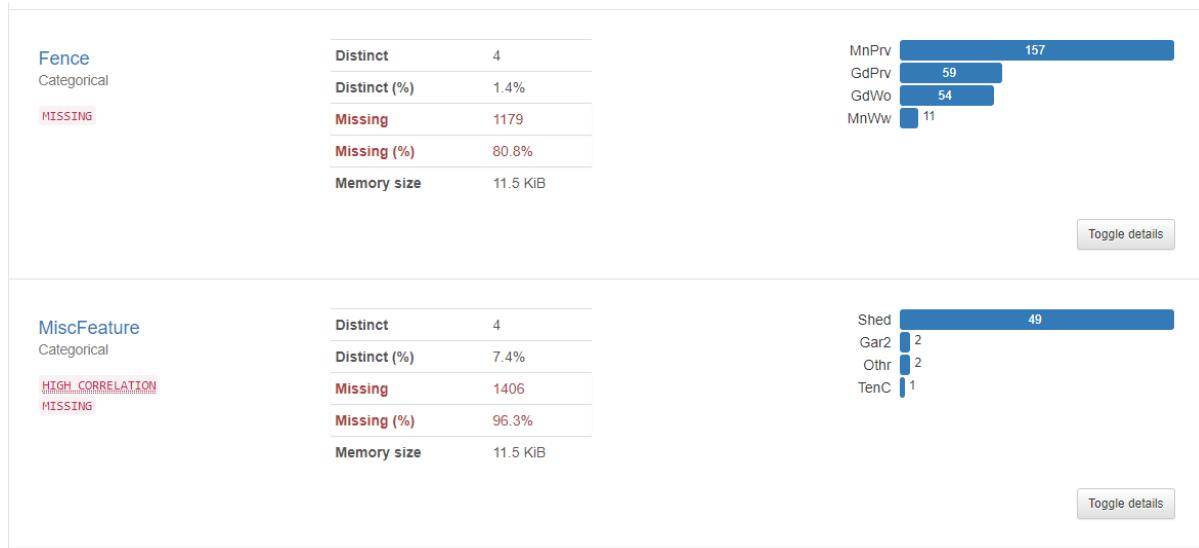


Donaldo Sebastian Garcia Jiménez 19683  
 Raul Angel Jimenez Hernandez 19017

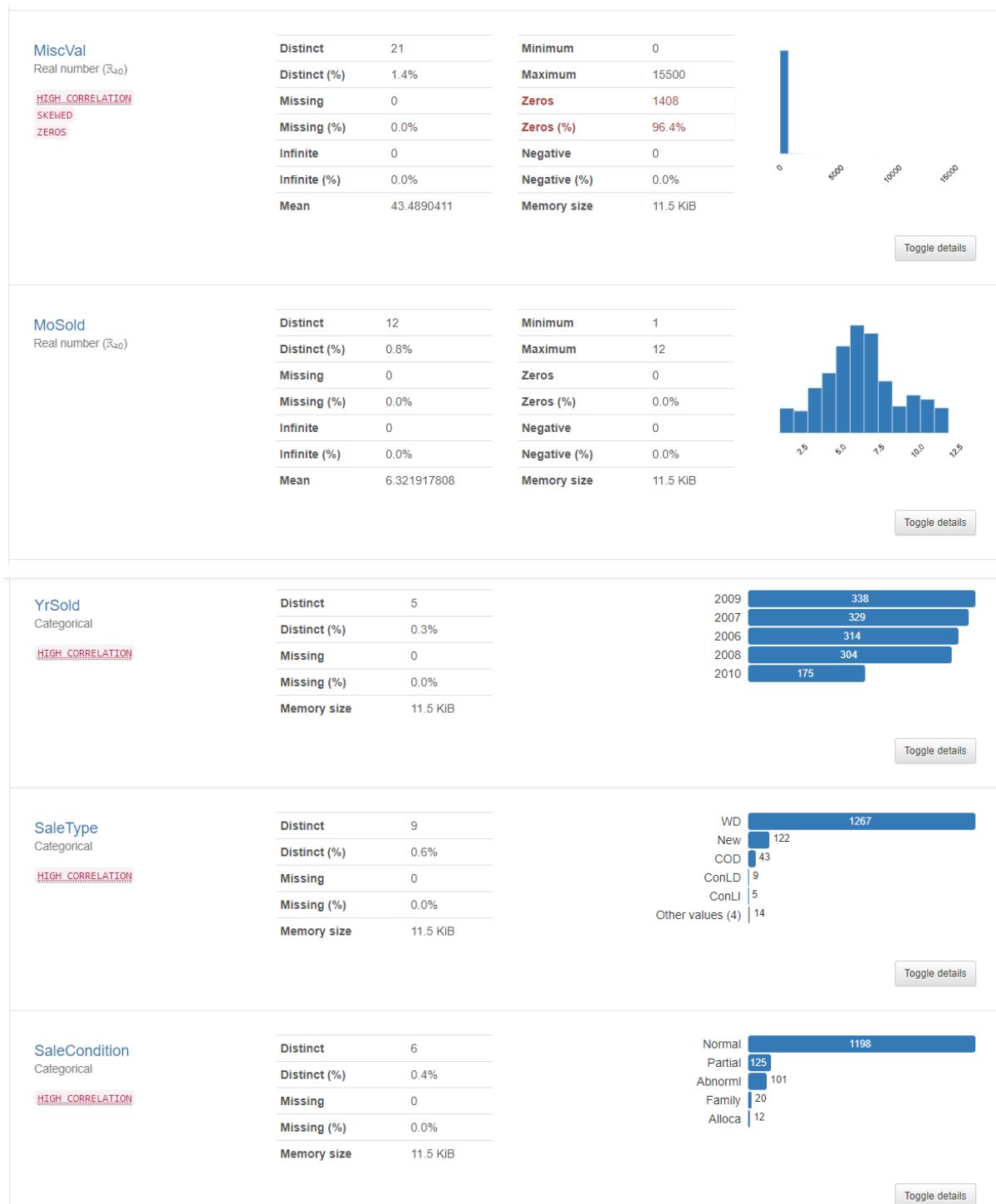


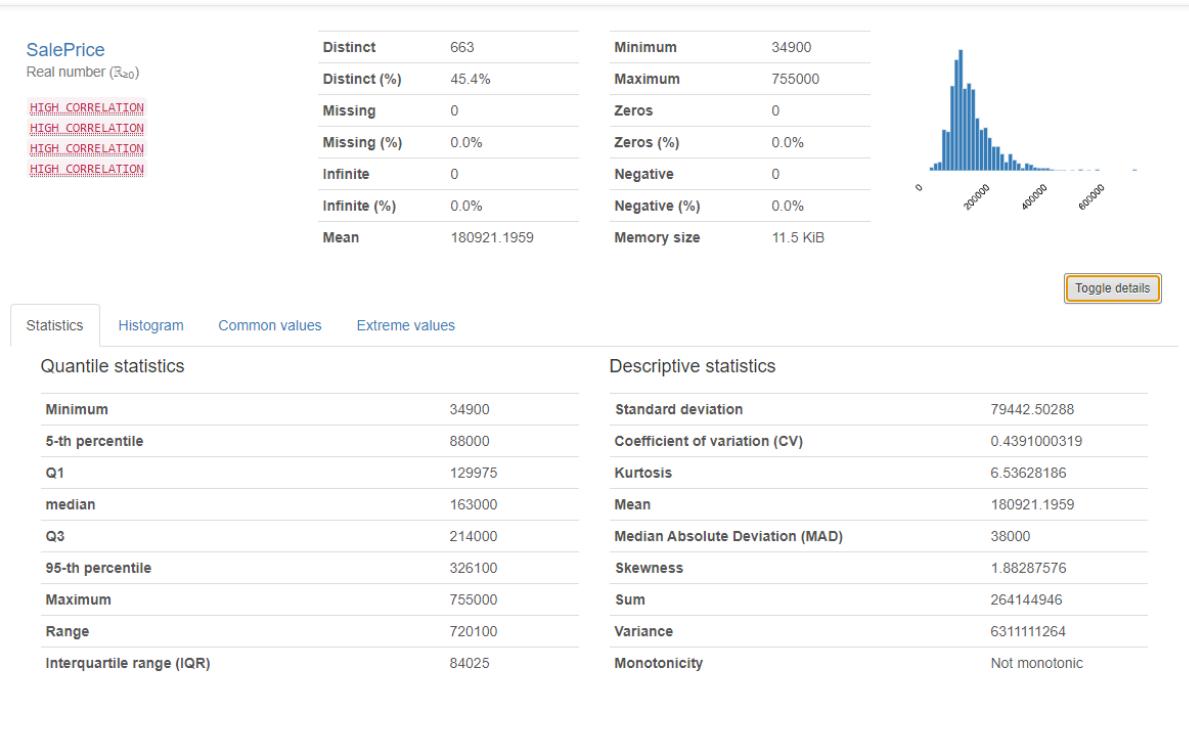
Donaldo Sebastian Garcia Jiménez 19683

Raul Angel Jimenez Hernandez 19017



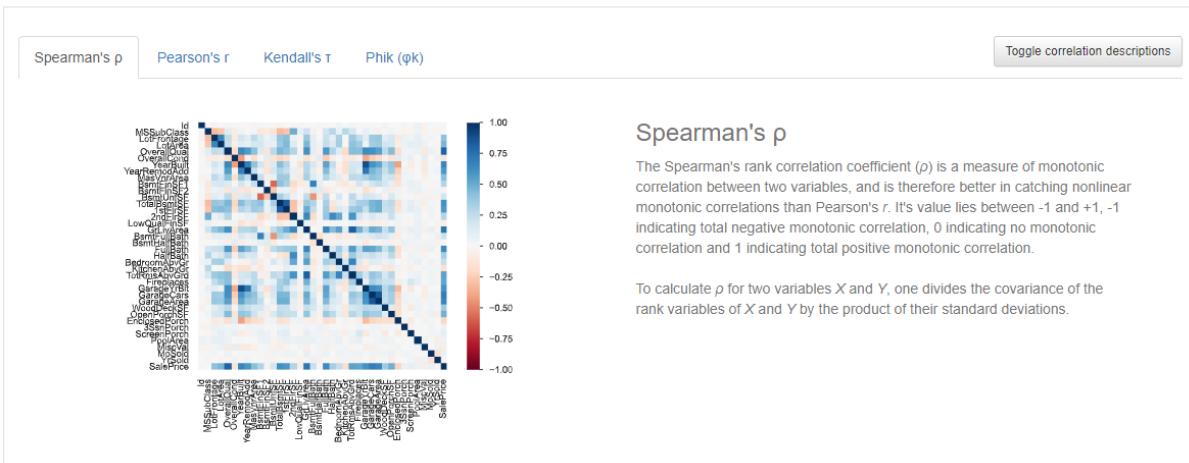
Donaldo Sebastian Garcia Jiménez 19683  
 Raul Angel Jimenez Hernandez 19017

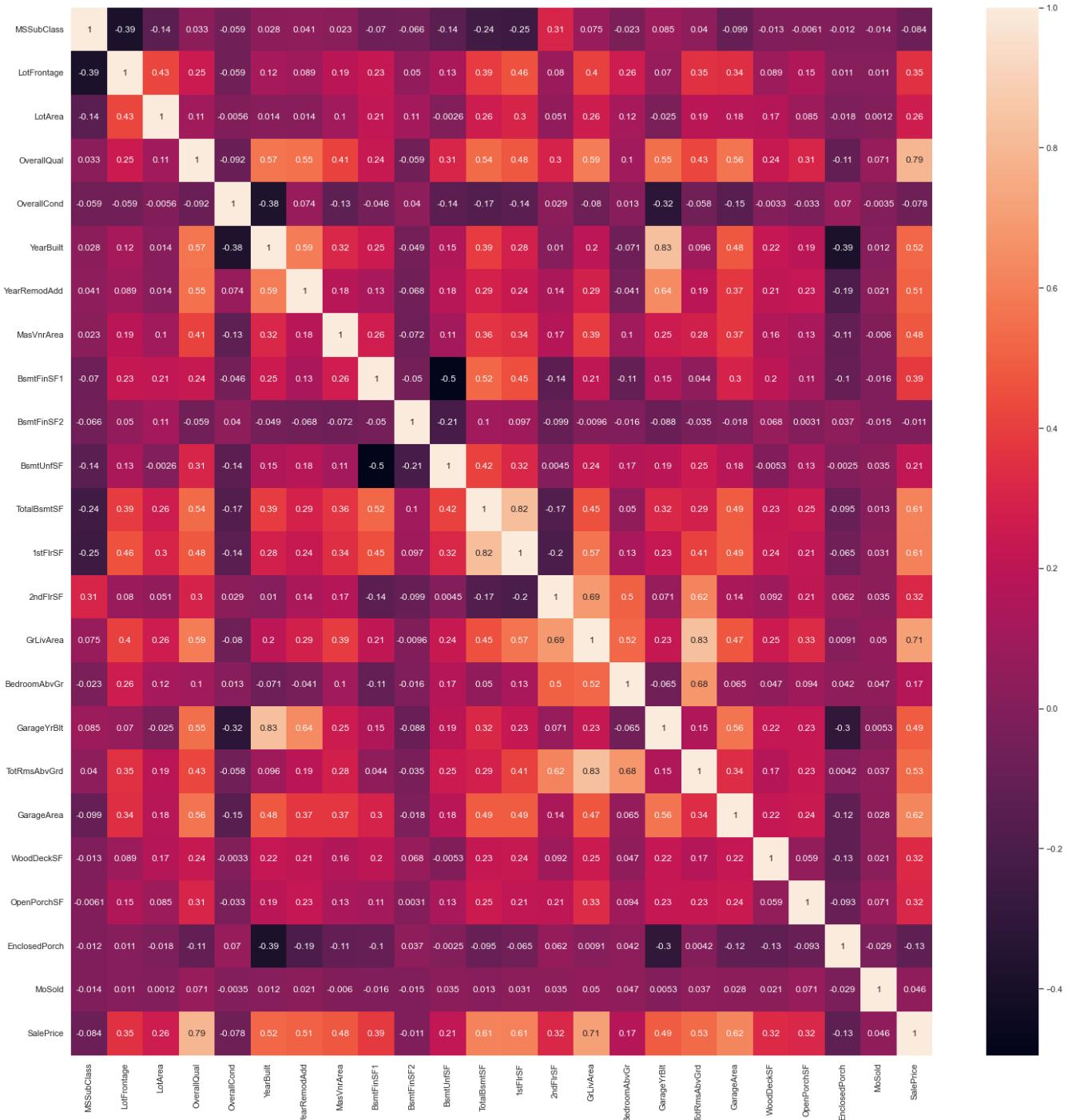




## Correlaciones

### Correlations





## Hallazgos y explicación de procedimiento

Para analizar los datos se utiliza la herramienta de pandas profiler la cual nos brinda información de estadística descriptiva importante de cada una de las variables, al igual que gráficas de barras o histogramas según sea la variable a analizar. Luego de esto se realiza la gráfica de correlación para poder determinar que variables son las que tienen más relación entre ellas y así enfocar más nuestras variables a utilizar.

Analizando los datos podemos observar que muchas de estas variables numéricas están relacionadas con otras. Una de estas por ejemplo puede ser el TotRmsAbvGrd y GrLivArea por que sabemos que dependiendo el área que se tenga va a influir en el total de cuartos

Donald Sebastian Garcia Jiménez 19683  
Raul Angel Jimenez Hernandez 19017

que va a tener. De igual forma logramos determinar que las variables numéricas más importantes para nuestro análisis son: 'MSSubClass', 'LotFrontage', 'LotArea', 'OverallQual', 'OverallCond', 'YearBuilt', 'YearRemodAdd', 'MasVnrArea', 'BsmtFinSF1', 'BsmtFinSF2', 'BsmtUnfSF', 'TotalBsmtSF', '1stFlrSF', '2ndFlrSF', 'GrLivArea', 'BedroomAbvGr', 'GarageYrBlt', 'TotRmsAbvGrd', 'GarageArea', 'WoodDeckSF', 'OpenPorchSF', 'EnclosedPorch', 'MoSold', 'SalePrice'. Por qué son algunas de las que tienen relación entre ellas además de que brindan información adecuada. De igual manera logramos determinar que las variables categóricas más importantes son: 'Street', 'Alley', 'LotShape', 'LandContour', 'LotConfig', 'LandSlope', 'BldgType', 'HouseStyle', 'BsmtQual', 'BsmtExposure', 'HeatingQC', 'GarageCars'.

## Análisis de componentes principales

### Análisis factorial

```
from factor_analyzer.factor_analyzer import calculate_kmo
kmo_all,kmo_model=calculate_kmo(houses_clean)
kmo_all, kmo_model
8] ✓ 17.4s
· c:\Users\ALIWARE\AppData\Local\Programs\Python\Python39\lib\site-packages\f
    warnings.warn('The inverse of the variance-covariance matrix '
                  '(array([[0.67986034, 0.85336964, 0.86511003, 0.91950297, 0.52939275,
                         0.76092737, 0.83176101, 0.94283582, 0.70914247, 0.11991146,
                         0.64250957, 0.87220836, 0.62952928, 0.45033991, 0.66780668,
                         0.77086765, 0.82586225, 0.91502783, 0.91616437, 0.94585033,
                         0.93825884, 0.71845741, 0.46740089, 0.92301962]),
                  0.7683818262377667)
```

### Bartlett

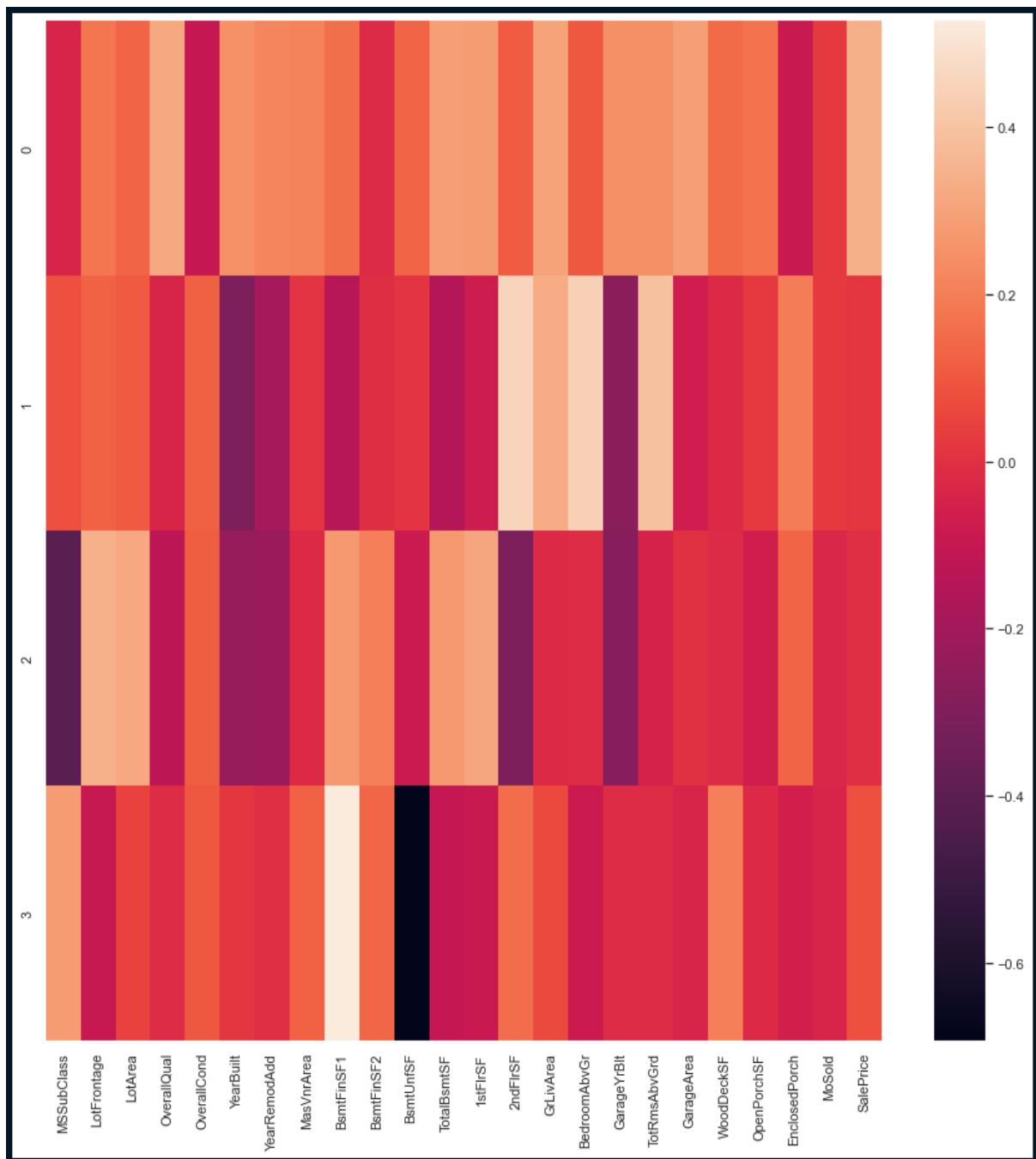
```
from factor_analyzer.factor_analyzer import calculate_bartlett_sphericity
chi_square_value,p_value=calculate_bartlett_sphericity(houses_clean)
chi_square_value, p_value
7] ✓ 0.9s
· (56989.803679788645, 0.0)
```

## Hallazgos

Como se puede observar nuestro resultado de KMO es de 0.77, lo cual nos indica que tenemos datos adecuados para poder seguir con el análisis de PCA.

De igual forma nuestro test de bartlett nos dio 0.0 en el p-value lo cual nos indica que la matriz no es igual a la matriz de identidad por lo cual se puede continuar con el análisis, de igual forma nos indica que el valor de chi-cuadrado es de 56989.80

## Matriz de correlación



Donald Sebastian Garcia Jiménez 19683  
Raul Angel Jimenez Hernandez 19017

Como podemos observar nuestra matriz tiene 4 filas, eso se decide de esta forma debido a que en nuestra gráfica de PCA podemos observar alrededor de 4 distintos colores. Para lograr determinar cuales son las variables que tienen más relación entre sí primero lo que se hizo fue separar las cualitativas más importantes y luego con ayuda del algoritmo de PCA determinar cuántos componentes iba a tener. En nuestro caso los 4 principales componentes que se observaron.

## Coeficientes principales.

Analizando nuestros componentes principales logramos identificar que el primero (fila 0) tiene una fuerte relación con GarageArea, GrlivArea, 1flrSF, TotalBsmtSF, OverallQual entre algunas otras pero principalmente con estas 5 variables.

Nuestro segundo componente tiene una gran relación con 2ndFlrSF, BedroomAbvGrnd, TotRmsAbvGrnd y GrlivArea.

En nuestro caso el tercero se puede observar 5 componentes LotFrontArea, LotArea, BsmtFinSF1, TotalBsmtSF y 1stFlrSF.

Por último nuestro tercer componente se puede observar que está fuertemente relacionado con BsmtFinSF1 seguido de MSSubClass y WoodDeckSF.

## PCA



Observando nuestra gráfica de PCA podemos determinar que existen entre 4 y 5 componentes principales. Esto quiere decir que entre todas las observaciones que tenemos existen 4 principales. En nuestro caso vamos a tomar 4 como principales debido a que se pueden observar 4 colores principalmente destacados en la gráfica.

## Reglas de asociación

```
# El minimo de cobertura o soporte es de 20% y el minimo de confianza es de 70%
reglas_asociacion = apriori(records, min_support=0.2, min_confidence=0.7)
reglas = list(reglas_asociacion)

✓ 0.5s

print(len(reglas))
✓ 0.6s
1531

print(reglas[0])
✓ 0.1s
RelationRecord(items=frozenset({'1Fam'}), support=0.7, ordered_statistics=[OrderedStatistic(items_base=frozenset(), items_add=frozenset({'1Fam'}), confidence=0.7, lift=1.0)])
```

Analizando las reglas de asociación nos dimos cuenta que de todos los datos (118260) se pueden crear 1531 reglas. Esto se puede decir que existen alrededor de 1531 distintas casas y todas las casas dentro de nuestra base de datos entra en alguna de estas reglas de asociación. En resumen se crean 1,531 reglas con las variables categóricas proporcionadas. Lo cual lo podemos interpretar como 1,531 distintos tipos de casas. Con esto logramos encontrar el conjunto de datos frecuentes.

## Hallazgos y conclusiones

Para resumir el proceso podemos decir que lo principal que se debe de realizar es el análisis de las variables. Una vez clasificadas las variables será más fácil poder determinar qué variables tienen más relación entre sí, qué variables son más relevantes. En este caso se utilizó la herramienta de pandas profile que nos permite realizar un análisis más detallado de todos estos datos. Cantidades mínimas, máximas, al igual que cuartiles y varianzas, entre otras.

Gracias a este análisis de variables logramos utilizar el algoritmo PCA. Este nos ayudó a determinar que existen cuatro componentes principales en nuestro conjuntos de datos y con ayuda de la matriz de correlación logramos determinar cuales son las variables cuantitativas que tienen más relación con dichos componentes principales.

Por último con ayuda del algoritmo a priori logramos formar un conjunto de reglas de asociación. Estas nos dicen que hay 1,531 distintos sets de casas con las mismas características. Estos datos nos ayudan a encontrar el conjunto de datos frecuentes en nuestra base de datos, en este caso las casas con propiedades similares.

## Rúbrica

(42 puntos) Análisis Exploratorio:

- Estudia las variables cuantitativas mediante técnicas de estadística descriptiva
- Hace gráficos exploratorios como histogramas, diagramas de cajas y bigotes, gráficos de dispersión que ayudan a explicar los datos
- Analiza las correlaciones entre las variables, trata de explicar los outliers (puntos atípicos) y toma decisiones acertadas ante la presencia de valores faltantes.
- Estudia las variables categóricas
- Elabora gráficos de barra, tablas de frecuencia y de proporciones
- Elabora gráficos adecuados según el tipo de dato que representan

Donaldo Sebastian Garcia Jiménez 19683

Raul Angel Jimenez Hernandez 19017

- ~~Explica muy bien todos los procedimientos y los hallazgos que va haciendo.~~

(18 puntos) Análisis de componentes Principales

- ~~Estudia la matriz de correlación, la agrega y explica lo que observa en ella~~
- ~~Determina si es posible usar la técnica de análisis factorial para hallar las componentes principales~~
- ~~Determina si vale la pena aplicar las componentes principales interpretando el test de esfericidad de Bartlett~~
- ~~Obtiene los componentes principales y explica cuántos seleccionará para explicar la mayor variabilidad posible.~~
- ~~Interpreta los coeficientes principales.~~

(18 puntos) Reglas de asociación

- ~~Construye reglas de asociación usando el algoritmo a priori.~~
- ~~Discute sobre las reglas de asociación más interesantes teniendo en cuenta sus niveles de confianza y soporte.~~

(22 puntos) Hallazgos y conclusiones.

- ~~Hace un resumen de los hallazgos en el análisis exploratorio~~
- ~~Llega a conclusiones sobre el análisis de componentes principales~~
- ~~Determina las reglas de asociación más interesantes.~~