

A photograph of a person standing on a rocky cliff overlooking a beach at sunset. The sky is a warm orange and yellow. In the background, there's a town built on a hillside overlooking the ocean. The foreground shows the dark, textured surface of the cliff.

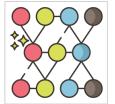
# IDEALISTA

## Data Analytics

By Jacky Barraza

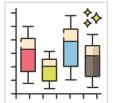


## About Jacky



## Data – Idealista

Quality



## Data Analytics

Price Houses in Madrid



## Architecture (cloud) I

Azure



## Architecture II

Azure



## Hazards and Mitigation

External web source



## HOUSE PRICE IN MADRID 2019 - 2020

idealista is in the market since 2000  
House listing website for buying and renting

Dataset, Madrid houses offered in  
[idealista.com](https://idealista.com) late 2019 - 2020:

### DATASET

18 Columns  
41,961 Rows

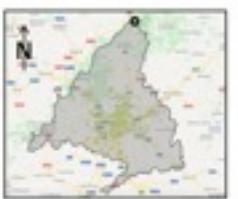
#### NUMERIC VAR

56 %  
11 var

#### CATEGORICAL VAR

43 %  
7 var

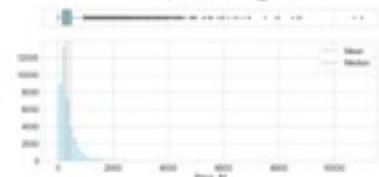
idealista



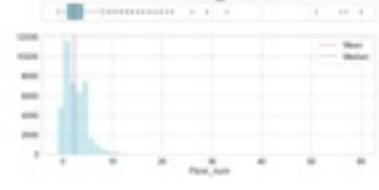
## HOUSE PRICE DATASET QUALITY DATA

Histograms and boxplots to have a general look to sample distribution in a variable

Price/1000: Price\_M



Floor\_num



### CATEGORICAL VARIABLES

Source URL - 41,961 Elevator\_class - 2  
Quarter - 136 Ext\_int - 2  
District - 21 New - 2  
Garage\_opt\_Inc - 3

### FINDINGS

Outliers  
Null Values - New  
High number of Classes in the categorical variables  
No duplicates



# HOUSE PRICE IN MADRID

## 2019 - 2020

idealista is in the market since 2000  
House listing website for buying and renting

Dataset, Madrid houses offered in  
[idealista.com](https://idealista.com) late 2019 - 2020:

### DATASET

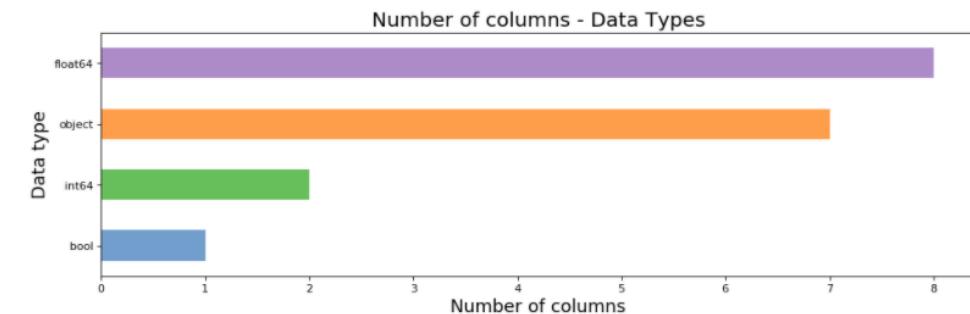
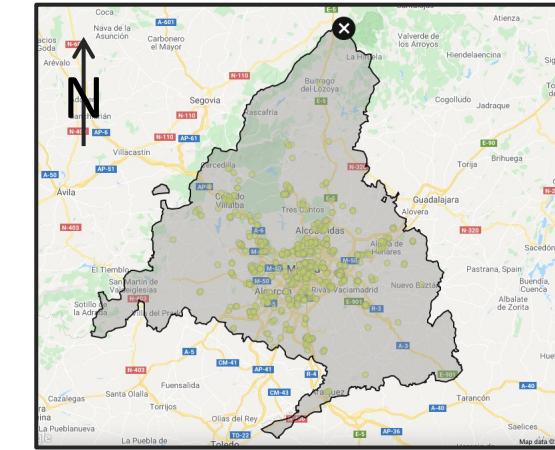
18 Columns  
41,961 Rows

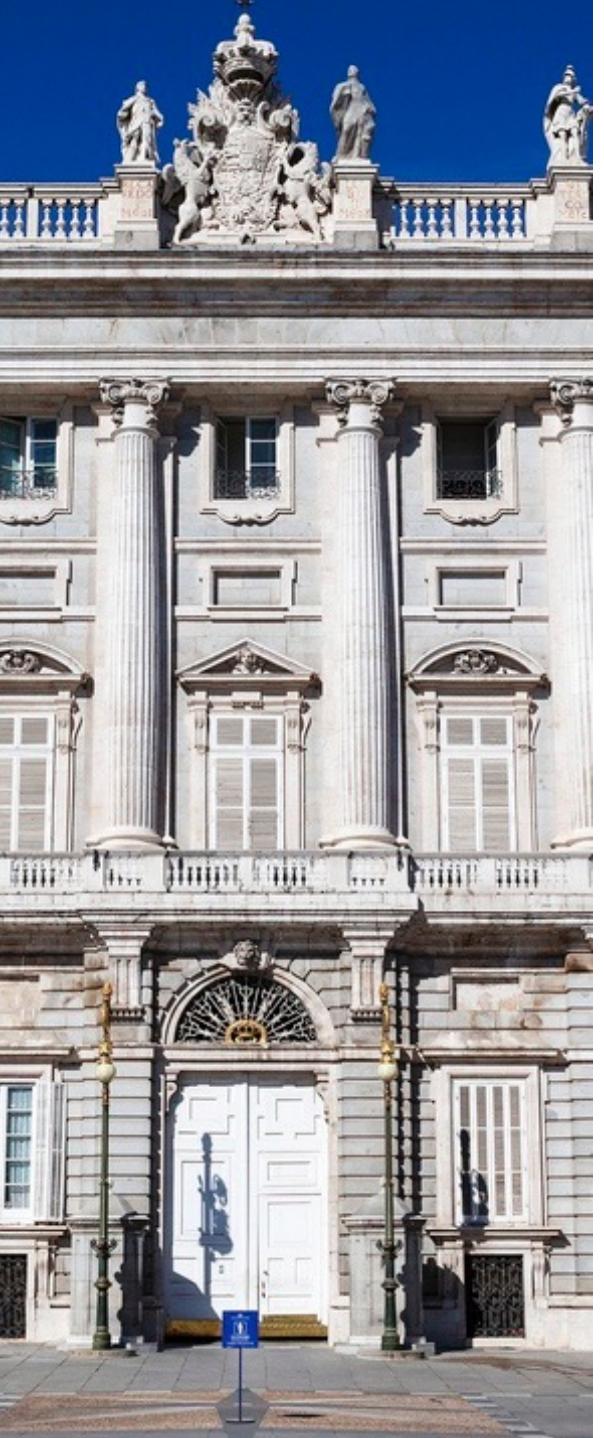
#### NUMERIC VAR

56 %  
11 var

#### CATEGORICAL VAR

43 %  
7 var

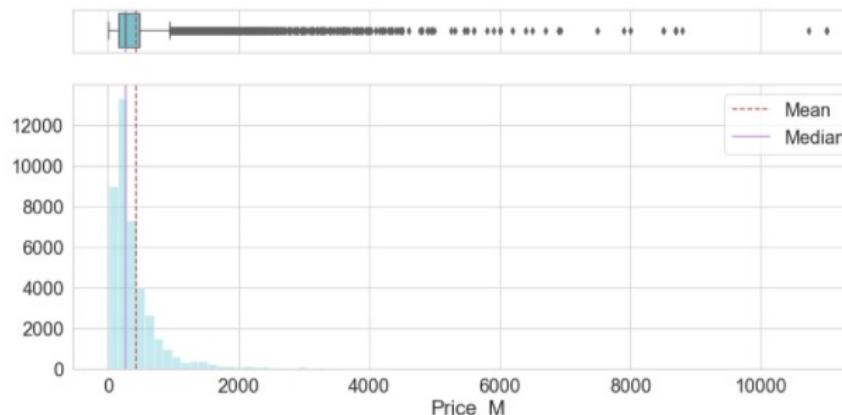




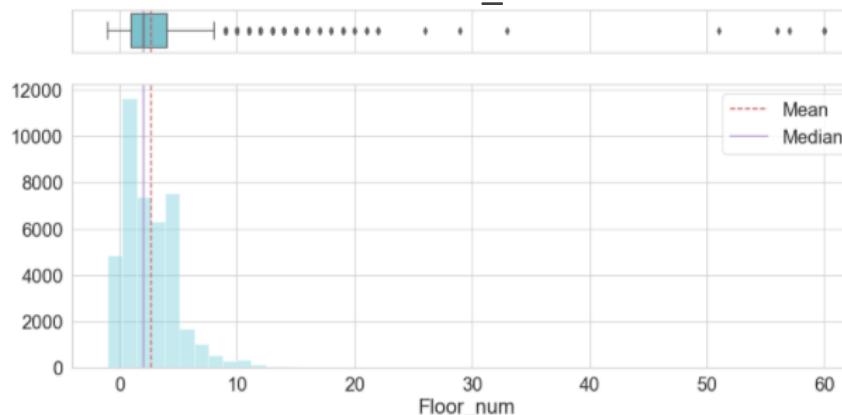
# HOUSE PRICE DATASET QUALITY DATA

Histograms and boxplots to have a general look to sample distribution in a variable

Price/1000: Price\_M



Floor\_num



## CATEGORICAL VARIABLES

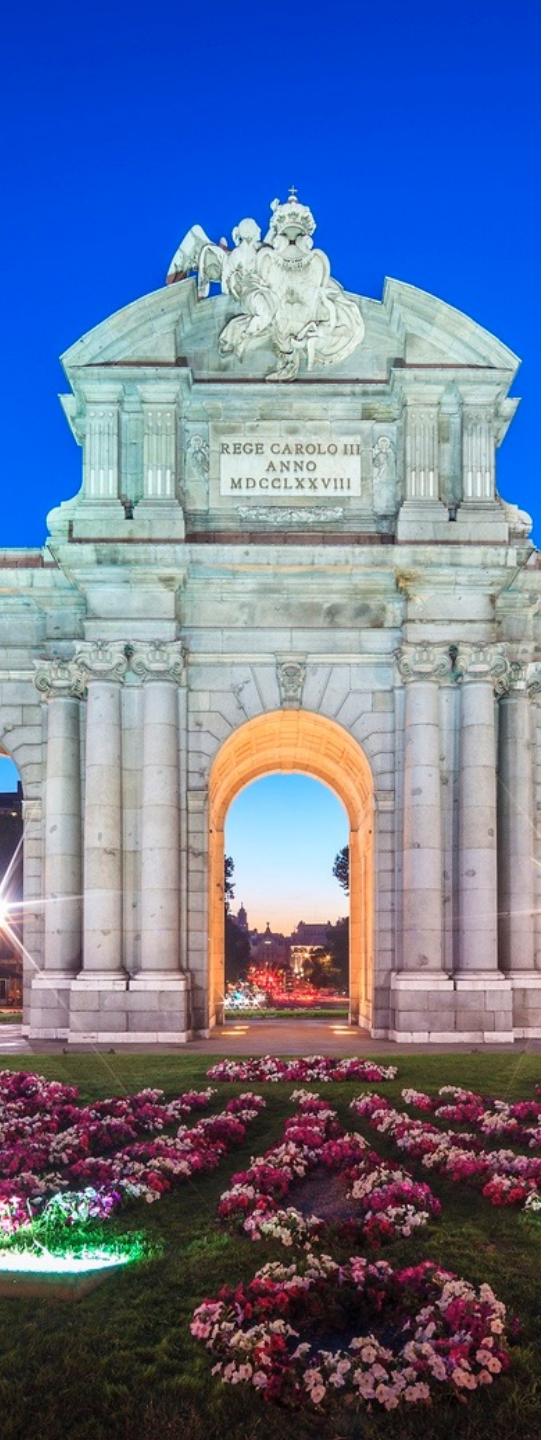
Source URL – 41,961	Elevator_class – 2
Quarter – 136	Ext_int – 2
District – 21	New – 2
Garage_opt_inc – 3	

## FINDINGS

- Outliers
- Null Values - New
- High number of Classes in the categorical variables
- No duplicates

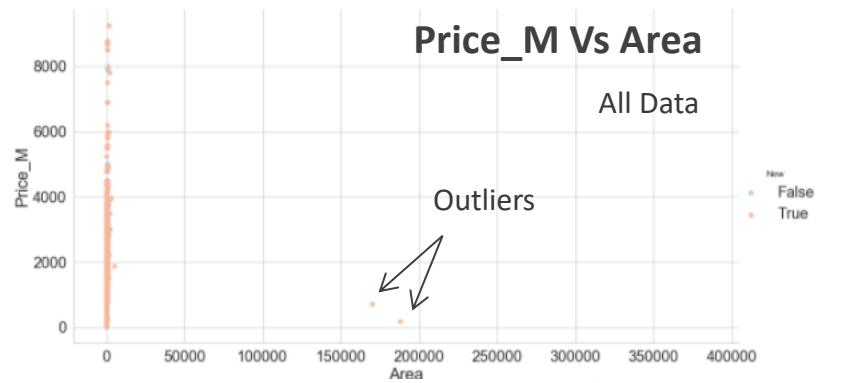






# MEAN AND MEDIAN

## MADRID PRICE HOUSE



## Price\_M Statistics

### LIMITS

Price: 13K – 11,007K  
Area: 12 – 385,000  
Samples: 41,961

### MEAN

426.01  
std: 515.34

### MEDIAN

270.0

### LIMITS

Price: 13K – 11,007K  
Area: 12 – 2,400  
Samples: 41,956 (-5)

### MEAN

425.88  
std: 515.34

### MEDIAN

270.0

### LIMITS

Price: 13K – 2,495K  
Area: 12 – 1387  
Samples: 41,943 (-18)

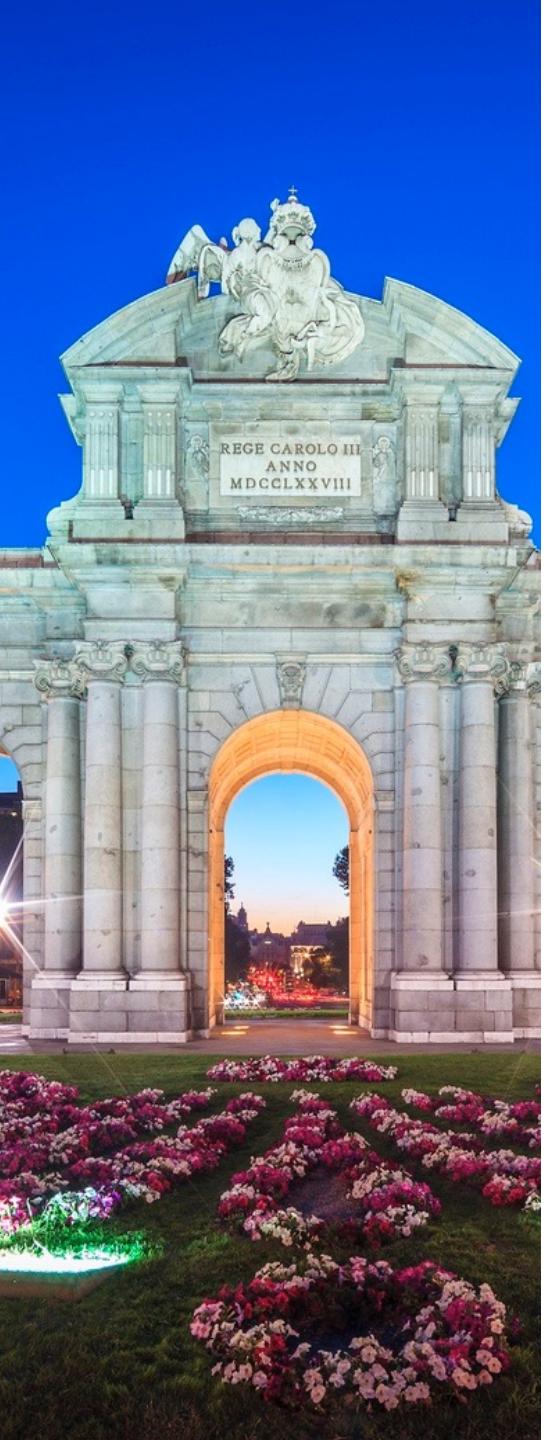
### MEAN

389.28  
std: 362.66

### MEDIAN

270.0

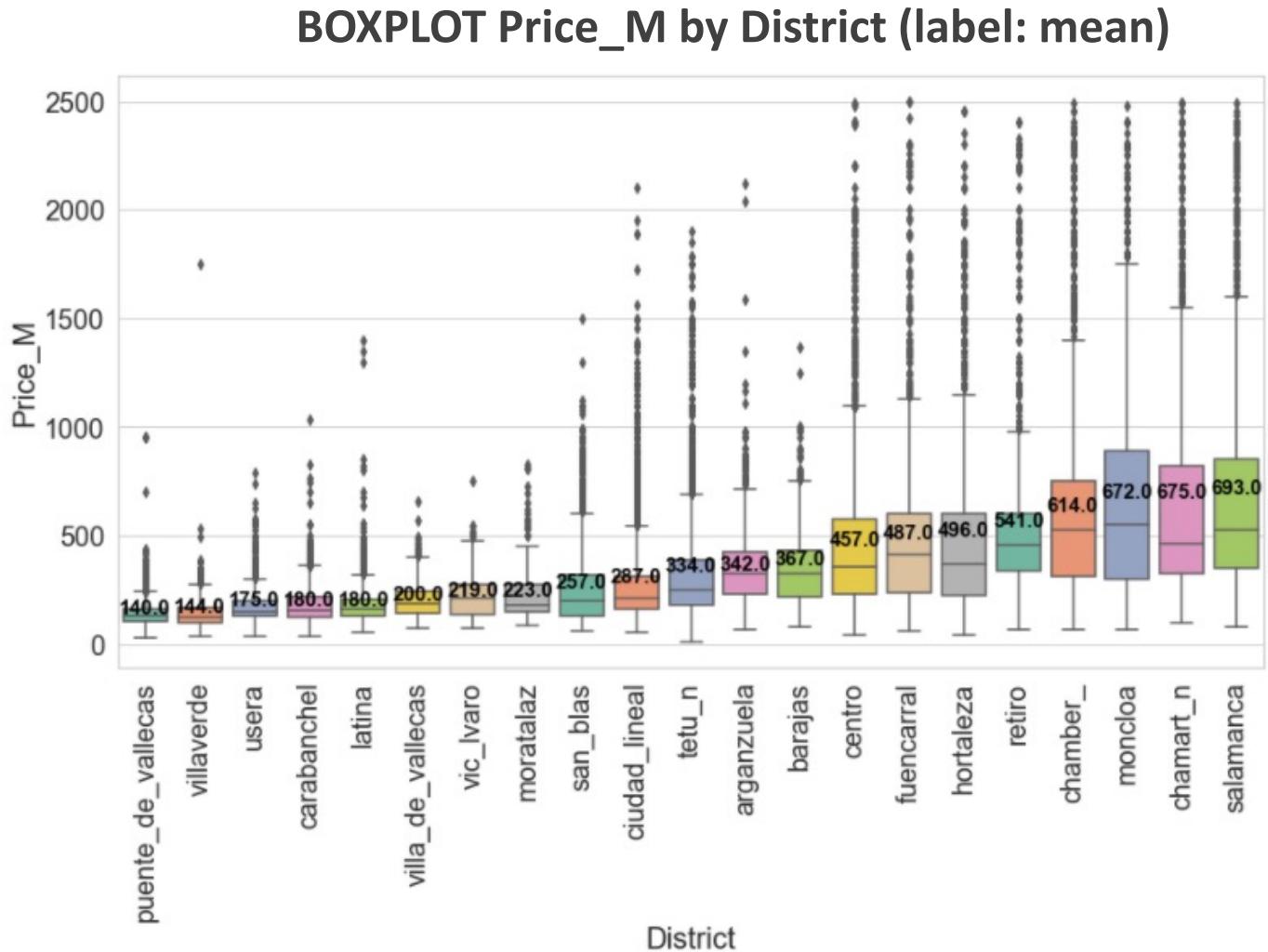


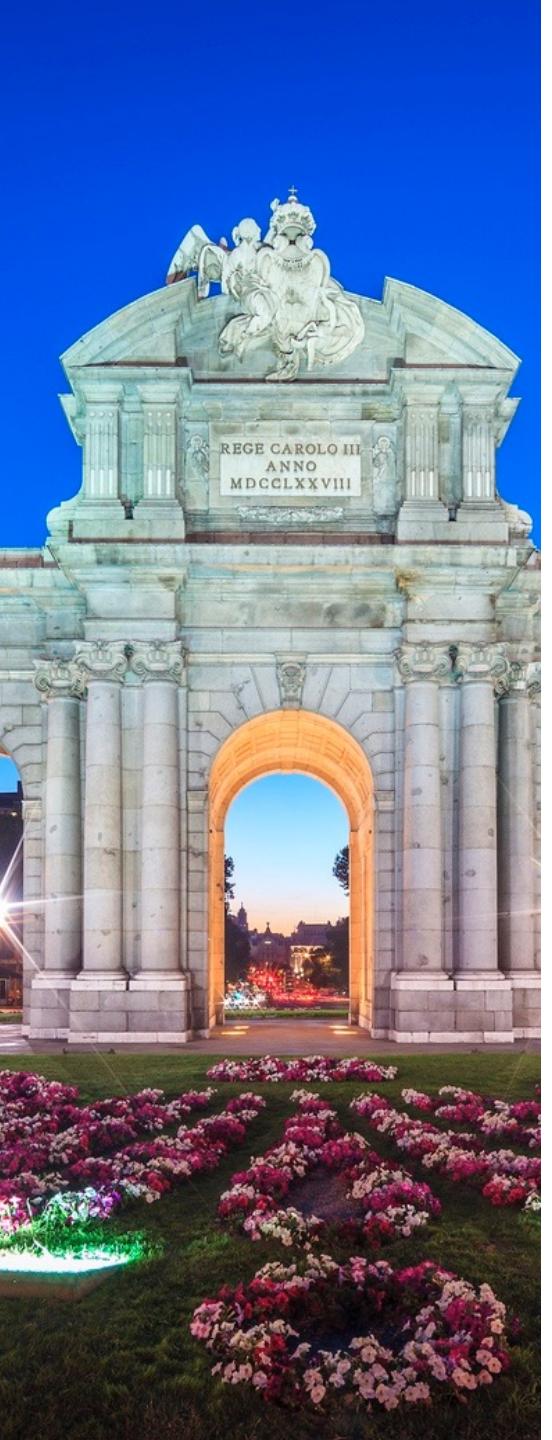


# MEAN AND MEDIAN MADRID PRICE HOUSE



- ✓ The lowest mean is **140K** in Puente de Vallecas
- ✓ The **highest** mean is **693K** in Salamanca
- ✓ The **closest** to the mean of all data **389K** are:
  - Arganzuela **342K**
  - Barajas **367K**
  - Tete\_n **334K**



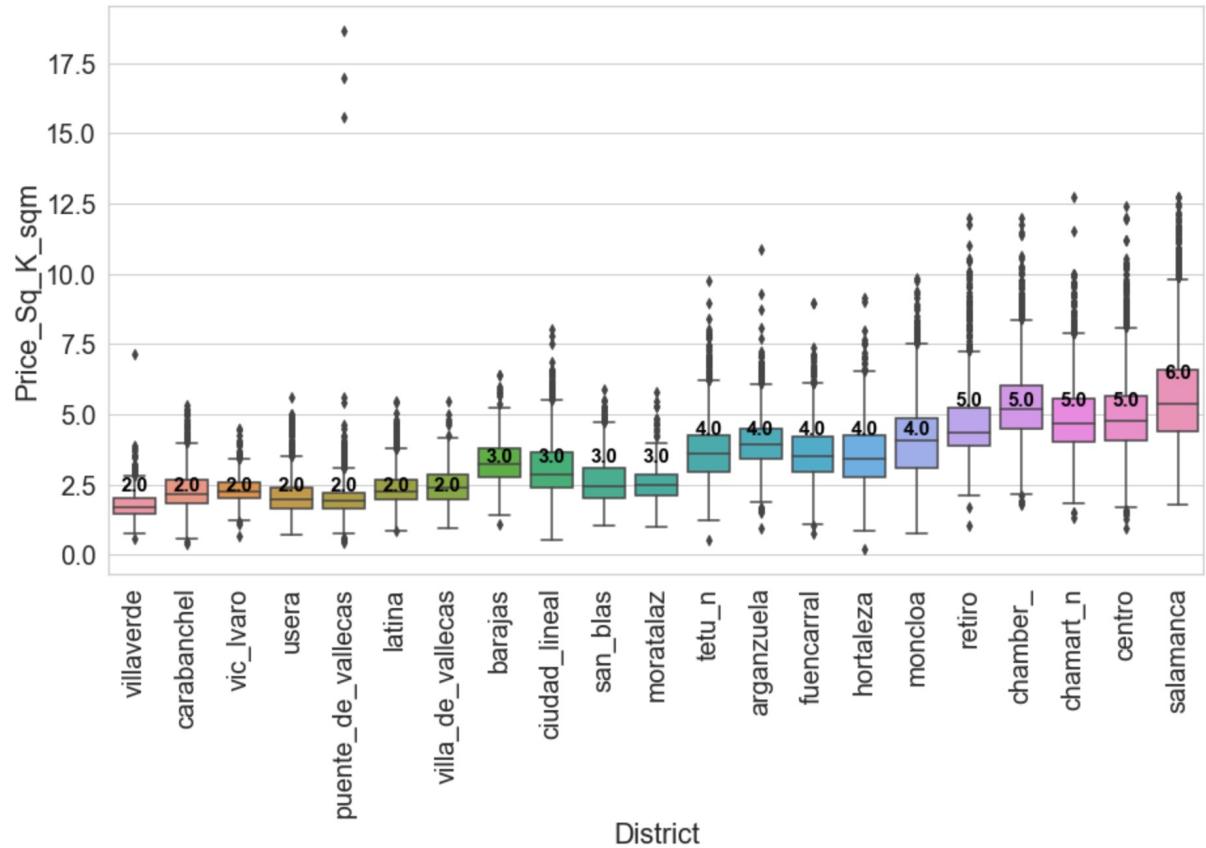


# MEAN AND MEDIAN MADRID PRICE HOUSE



- ✓ The lowest mean is **2K** in Villaverde
- ✓ The **highest** mean is **6K** in Salamanca

BOXPLOT Price/Sqmts by District (label: mean)

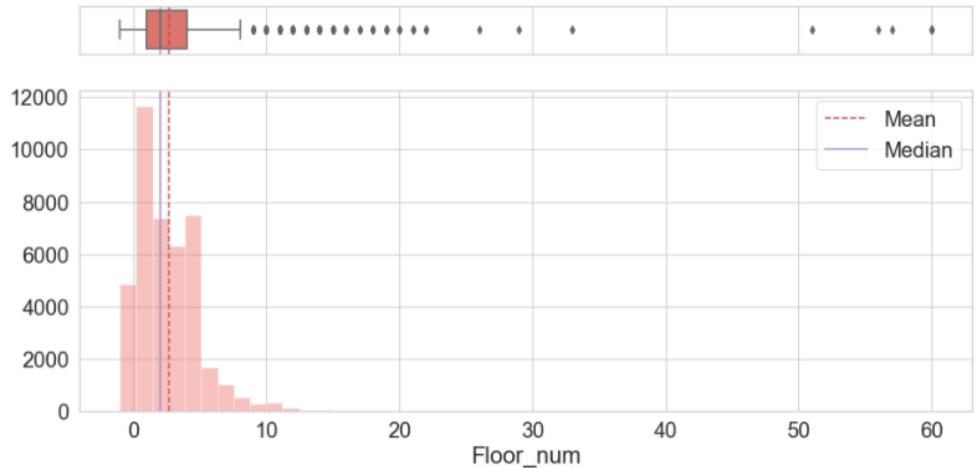


# FLOORS AND AREA (square mts)

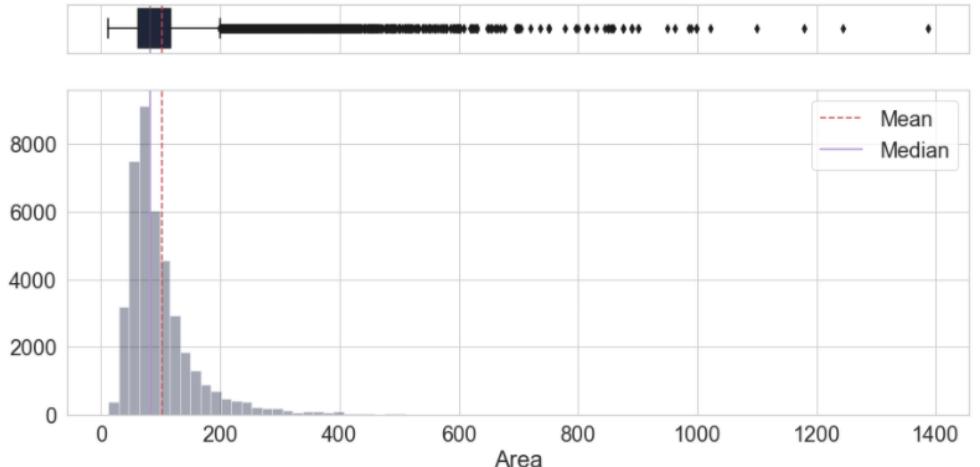
## MADRID PRICE HOUSE

The Histograms show a **Skewed distribution**

### Floor\_Num Distribution



### Area Distribution



### LIMITS

Floor\_room : (-1) – 60

### MEAN

2.63  
std: 2.39

### MEDIAN

2

### LIMITS

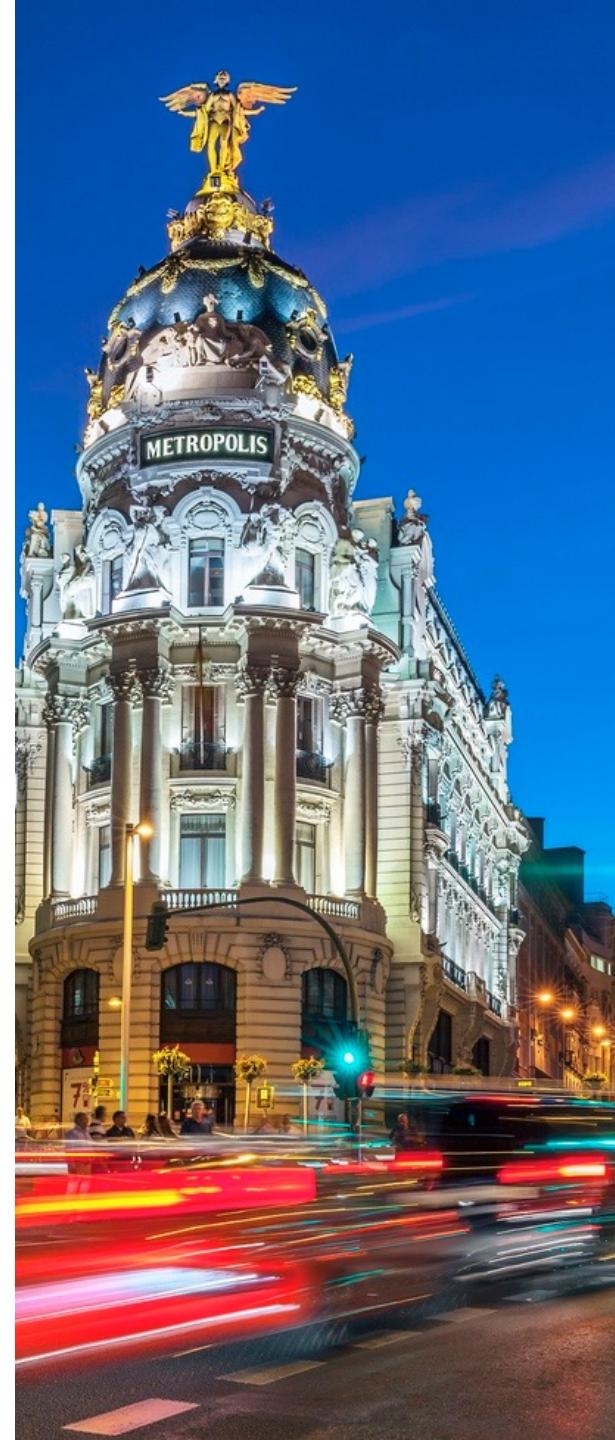
Area: 12 – 1,387

### MEAN

102.76  
std: 73.45

### MEDIAN

82



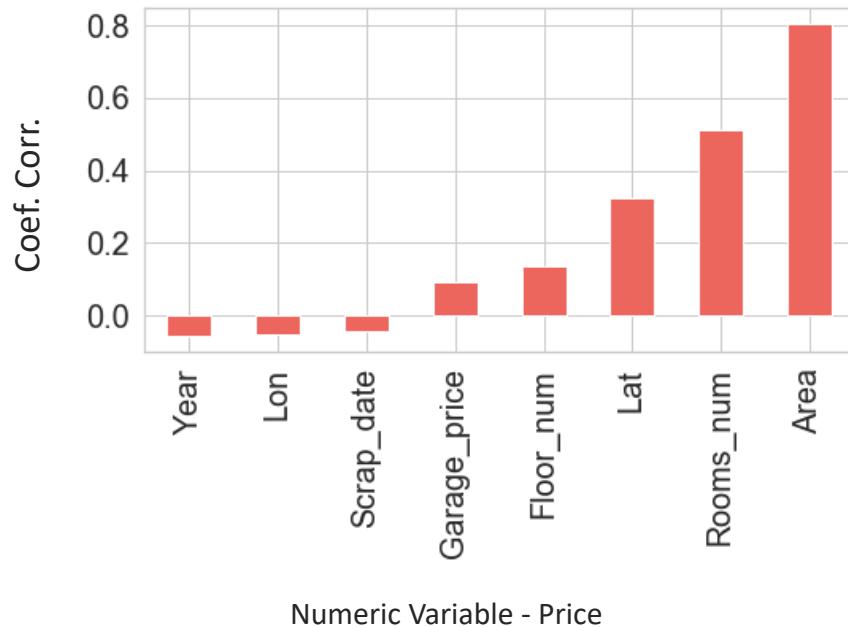


# PRICE CORRELATIONS

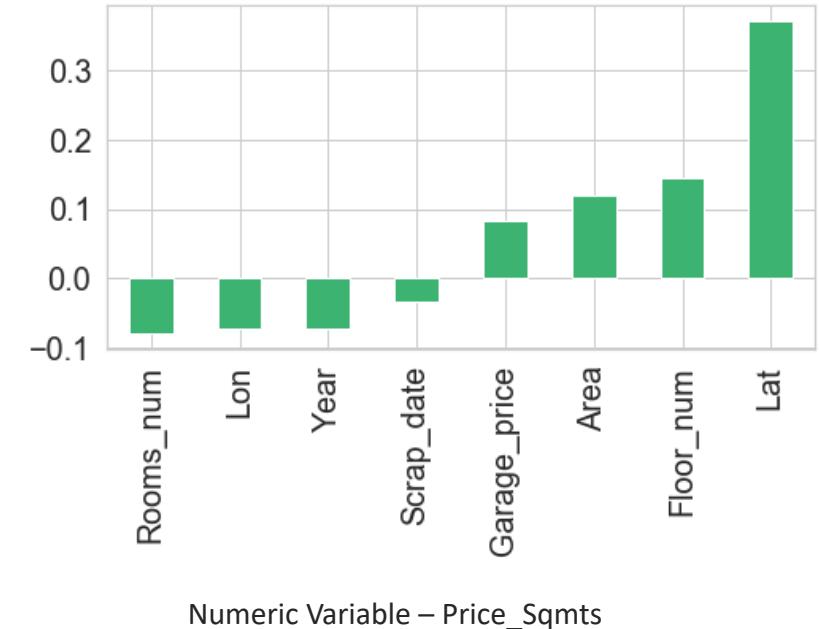
## MADRID PRICE HOUSE



Variable correlations against Price\_M



Variable correlations against Price\_Sqmts



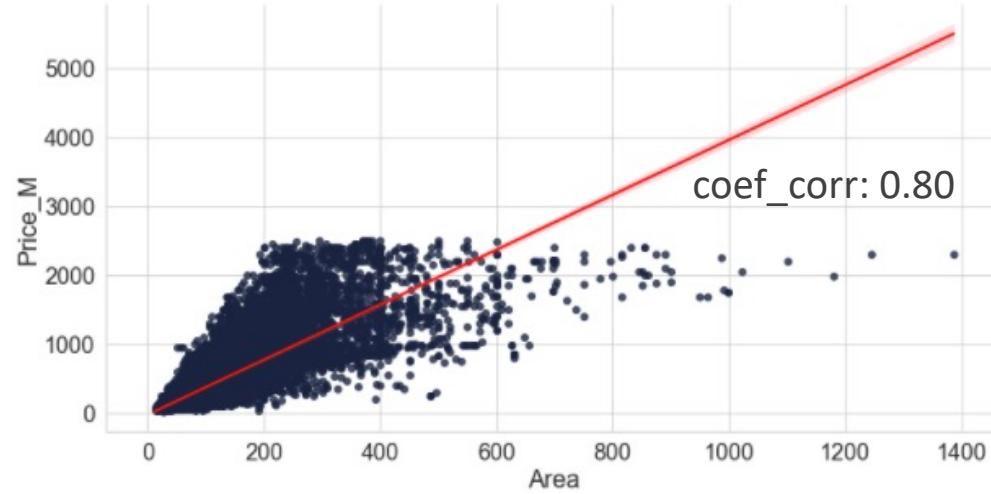


# PRICE CORRELATIONS

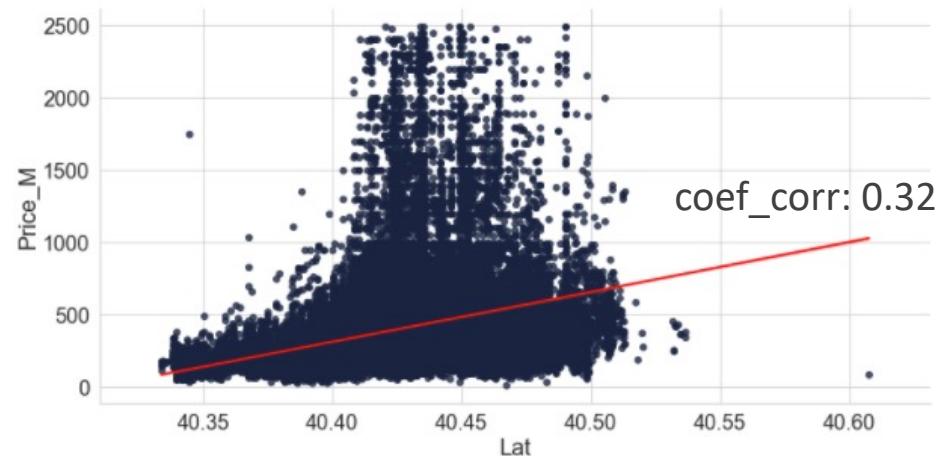
## MADRID PRICE HOUSE



Price\_M Vs Area



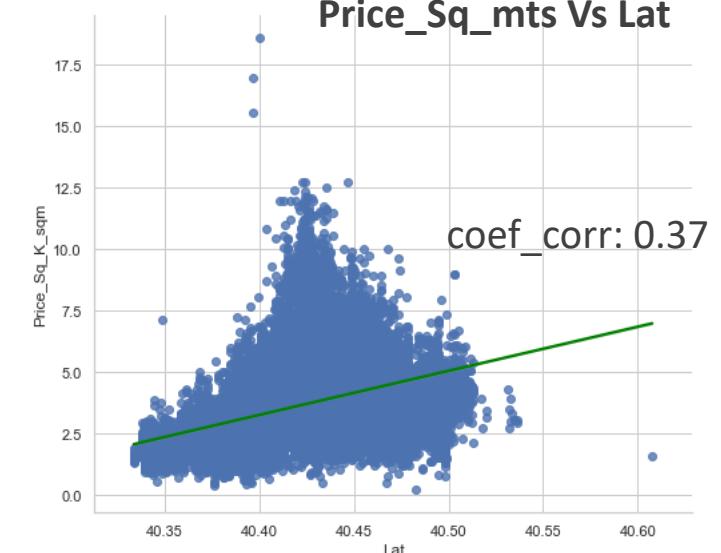
Price\_M Vs Lat



Price\_Sq\_mts Vs Area



Price\_Sq\_mts Vs Lat

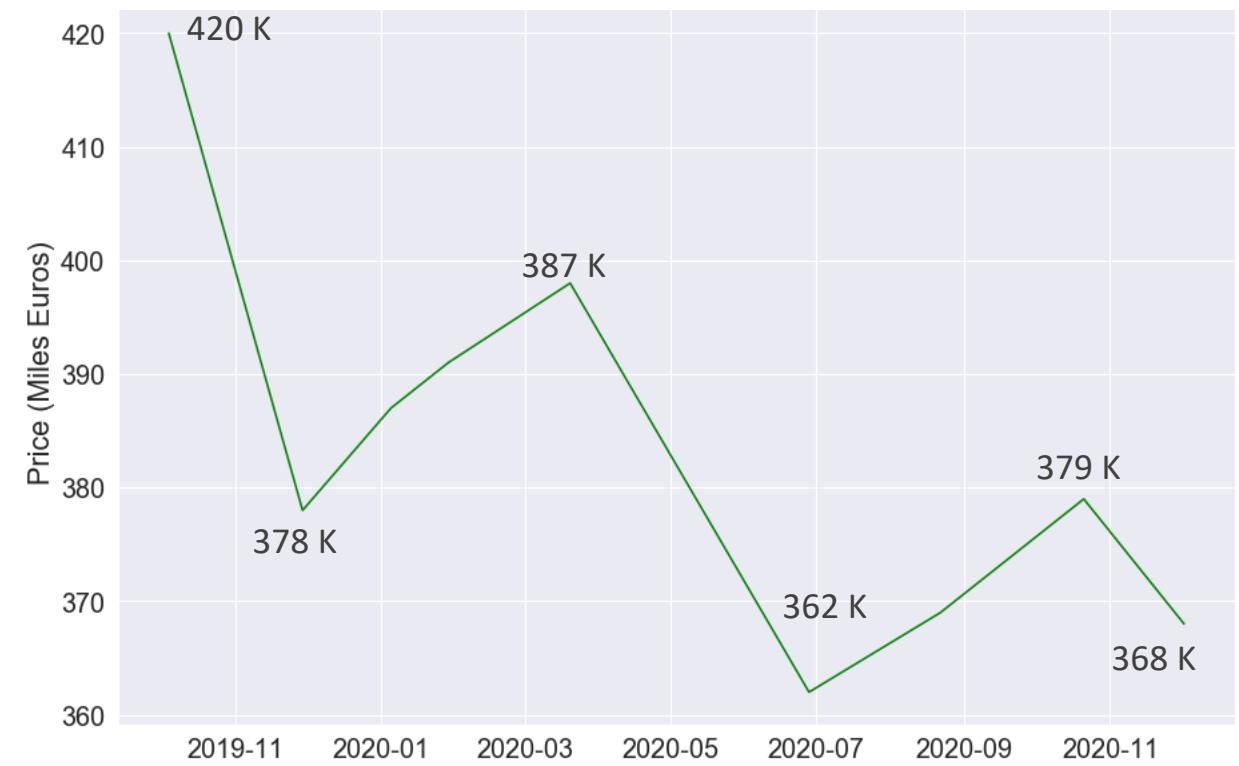
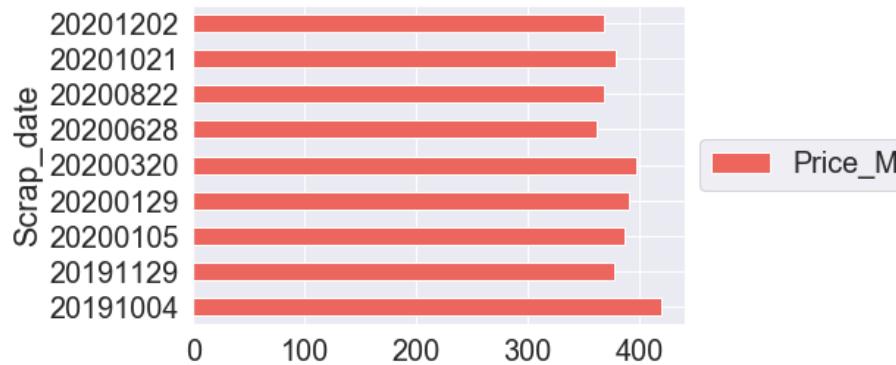




## PRICE PER SQUARE MTS EVOLVED 2020

### MADRID PRICE HOUSE

Price\_M aggrupation by Scrap\_date

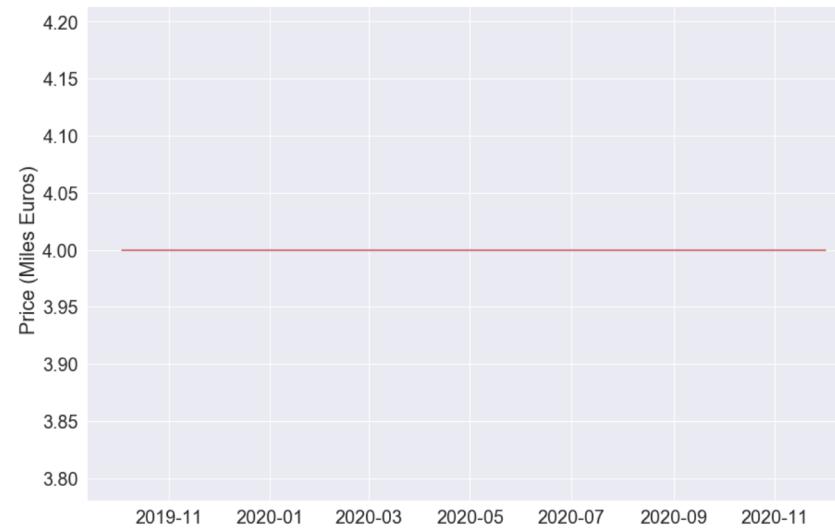




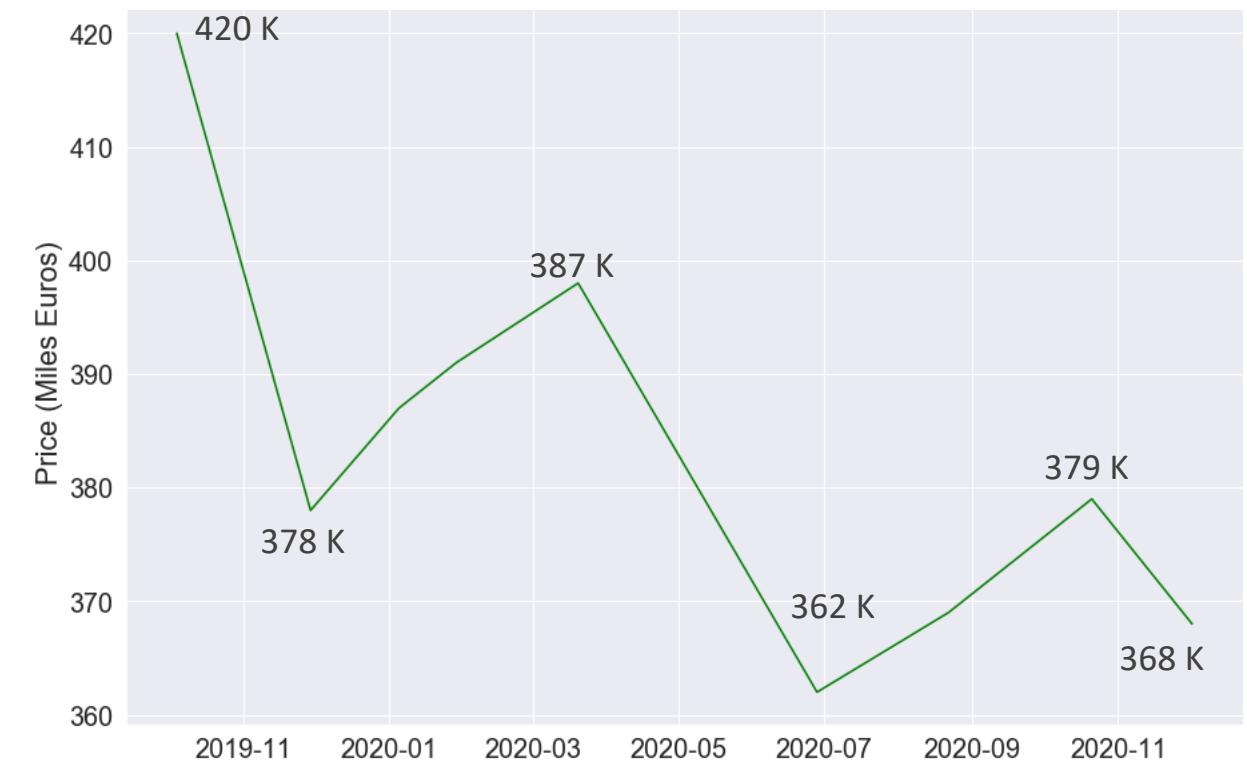
## PRICE PER SQUARE MTS EVOLVED 2020

### MADRID PRICE HOUSE

Price\_M/Sqmts aggrupation by Scrap\_date



Price\_M aggregation by Scrap\_date





# CLOUD ARQUITECURE PROPOSAL

MADRID PRICE HOUSE

---

## NEEDS:

- Allow repeatable data collection (scraping)
- Storing data (might needs images, videos...)
- Connection with data analytics, data scientist (real time – long term)
- Hazards mitigation (web sources)

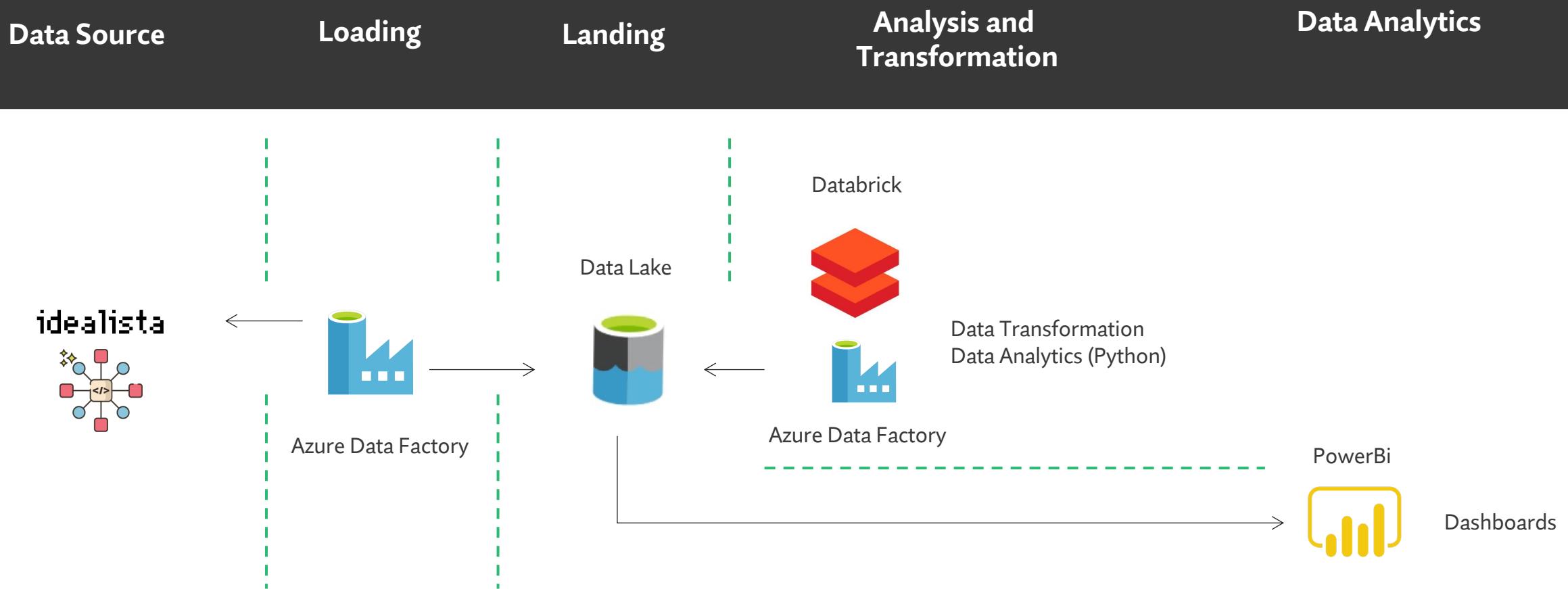




# CLOUD ARQUITECTURE I

## PROPOSAL

### MADRID PRICE HOUSE

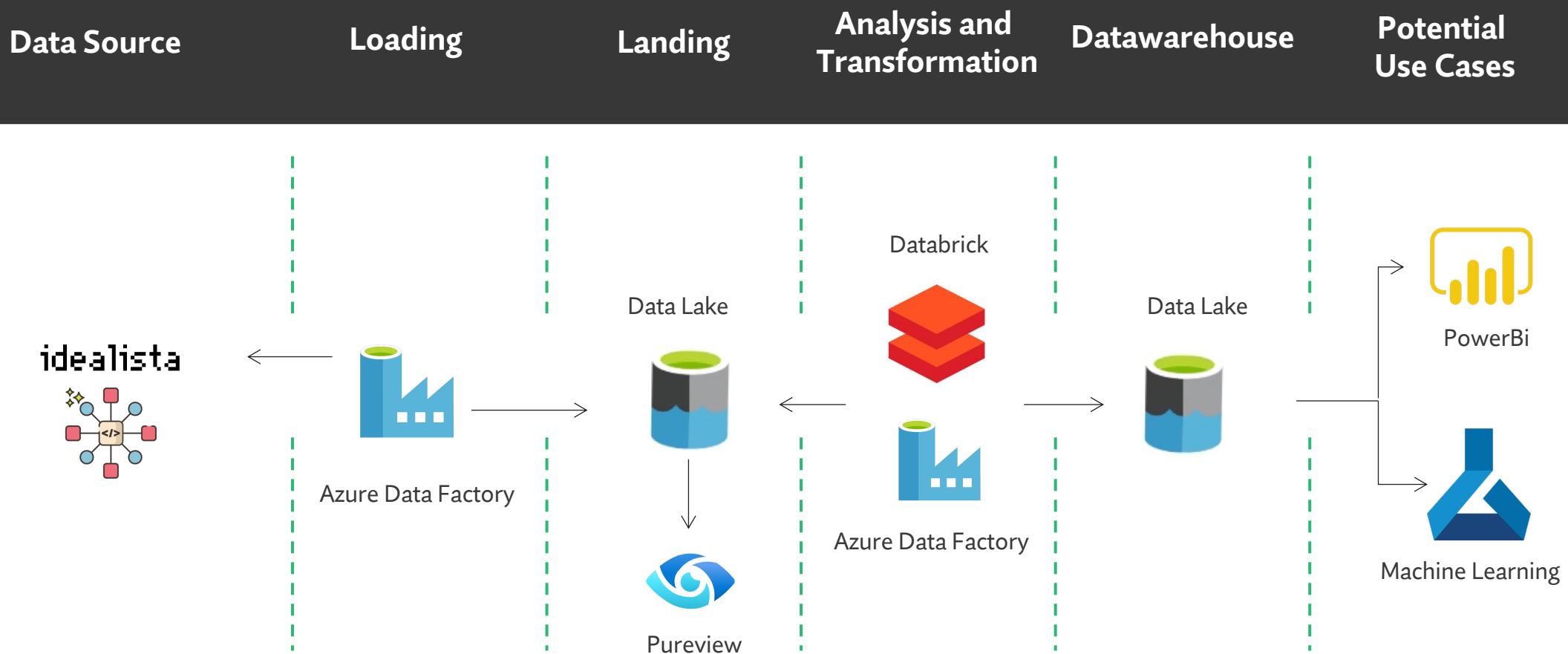




# CLOUD ARQUITECTURE II

## PROPOSAL

### MADRID PRICE HOUSE



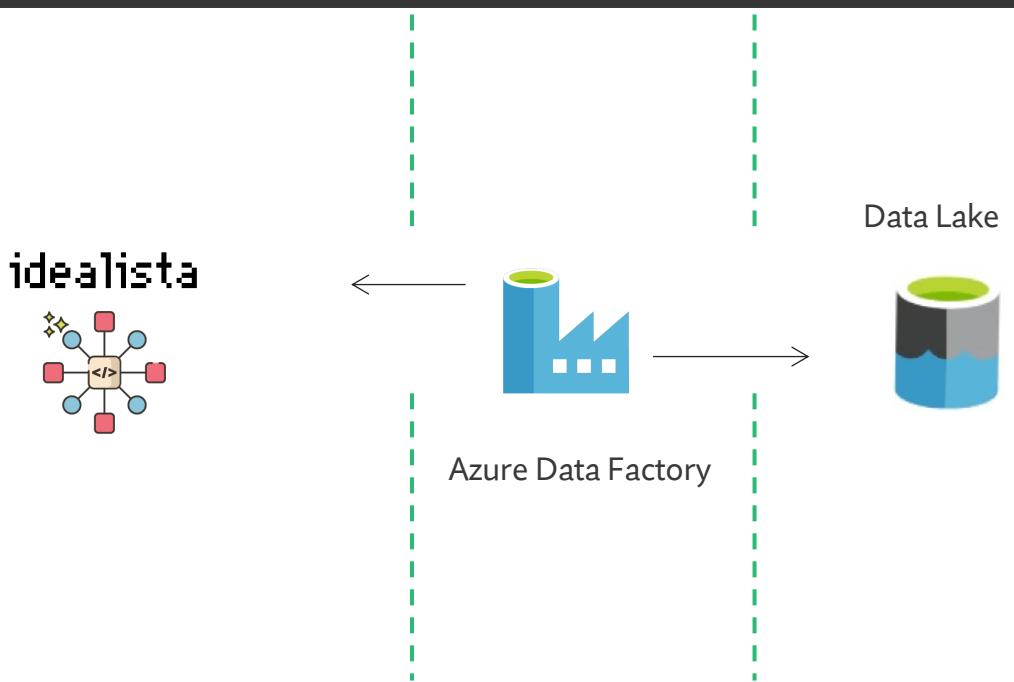


# SET UP AN ENVIRONMENT TO ALLOW REPEATABLE DATA

## COLLECTION (SCRAPING)

### MADRID PRICE HOUSE

Data Source      Loading      Landing



### Azure Data Factory

Build pipelines to trigger data collection every month and storage raw data into the Datalake.





# DATA STORAGE RECOMENDATIONS

MADRID PRICE HOUSE

---

Consideration: videos, images, AutoCAD files...

## Blob Storage



- ✓ Support Streaming and random access scenarios
- ✓ Object store with flat namespace
- ✓ Access app data from anywhere
- ✓ Datalake – big data analytics

## Data Lake



- ✓ Optimized storage for big data analytics workloads
- ✓ Hierarchical file system
- ✓ Unexpected failures - Data lake spreads parts of a file over a number of individual storage servers

## Cosmos DB



- ✓ Document store, graph, key-value store, wide column store
- ✓ Guarantees single-digit-millisecond latencies at the 99th percentile anywhere in the world
- ✓ Guarantees high availability with multi-homing
- ✓ No images – No Videos



## DATA AVAILABILITY (real-time and long term)

### MADRID PRICE HOUSE

---



#### AZURE – Service-level agreement (SLA)

##### Monthly Availability Percentage



---

###### MONTHLY READ AVAILABILITY PERCENTAGE

---

###### SERVICE CREDIT

---

< 99.999%

10%

---

< 99%

---

25%

---

# EXTERNAL WEB SOURCES

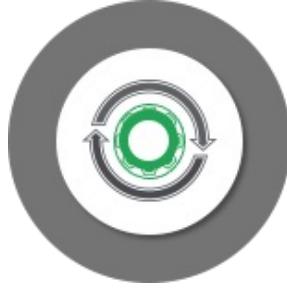
## Possible issues and mitigation

### MADRID PRICE HOUSE

---



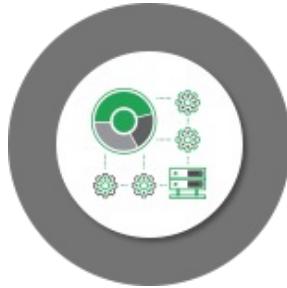
#### Dynamic content



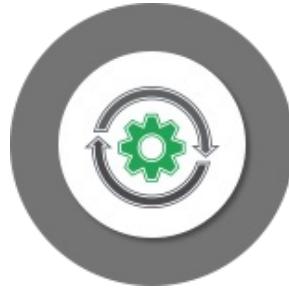
#### Complicated and changeable web pages structure



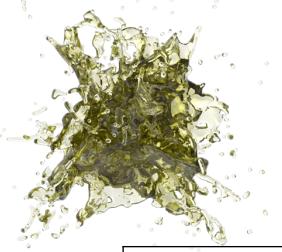
#### Real-Time



#### Slow/unstable load speed



- ✓ Create a system when data structure is different (COSMO DB handle this issues ok)
- ✓ Azure Event Hub (streaming data)
  - Anomaly detection (fraud/outliers)
  - Application logging
  - Analytics pipelines, such as clickstreams
  - Live dashboarding
  - Archiving data
  - Transaction processing
  - User telemetry processing
  - Device telemetry streaming



# CLOUD ARQUITECURE BUDGET

## MADRID PRICE HOUSE

---

Microsoft Azure Estimate

### Advanced Analytics on Big Data

Service type	Custom name	Region	Description	Estimated monthly cost	Estimated upfront cost
Data Factory		East US 2	Azure Data Factory V2 Type, Data Pipeline Service Type, Azure Integration Runtime: 100 Activity Run(s), 100 Data movement unit(s), 10,000 Pipeline activities, 10,000 Pipeline activities – External, Self-hosted Integration Runtime: 100 Activity Run(s), 1,000 Data movement unit(s), 10,000 Pipeline activities, 10,000 Pipeline activities – External, Data Flow: 1 x 8 Compute Optimized vCores x 730 Hours, 1 x 8 General Purpose vCores x 730 Hours, 1 x 8 Memory Optimized vCores x 730 Hours, Data Factory Operations: 100 x 50,000 Read/Write operation(s), 100 x 50,000 Monitoring operation(s)	\$5,045.12	\$0.00
Power BI Embedded		East US 2	1 node(s) x 720 Hours, Node type: A2, 2 Virtual Core(s), 5GB RAM, 301-600 Peak renders/hour	\$1,445.83	\$0.00
Azure Databricks		West US	All-Purpose Compute Workload, Premium Tier, 1 D3V2 (4 vCPU(s), 14 GB RAM) x 730 Hours, Pay as you go, 0.75 DBU x 730 Hours	\$504.80	\$0.00
Azure Data Lake Storage Gen1		East US 2	Pay-as-you-go: 200 GB Storage, 500 Read Transactions, 500 Write Transactions	\$34.80	\$0.00
Data Factory		East US	Azure Data Factory V2 Type, Data Pipeline Service Type,	\$0.00	\$0.00
Support			Support	\$0.00	\$0.00
			Licensing Program	Microsoft Online Services Agreement	
			Total	\$7,030.22	\$0.00
Disclaimer					
All prices shown are in US Dollar (\$). This is a summary estimate, not a quote. For up to date pricing information please visit <a href="https://azure.microsoft.com/pricing/calculator/">https://azure.microsoft.com/pricing/calculator/</a>					
This estimate was created at 3/15/2021 9:05:27 PM UTC.					





**GRACIAS**

**Preguntas**



# PRICE PER SQUARE MTS EVOLVED 2020

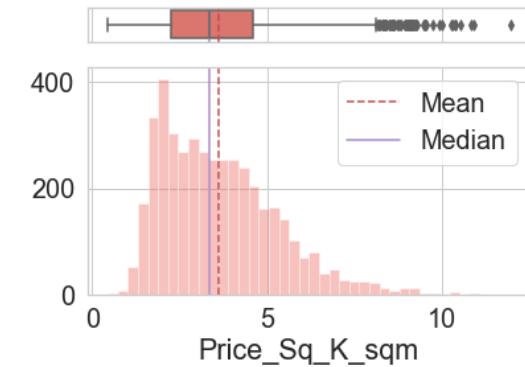
## MADRID PRICE HOUSE

Price by Sqmts aggrupation by Scrap\_date



En 1.900 horas con 5 personas en el equipo se alcanzó a desarrollar:

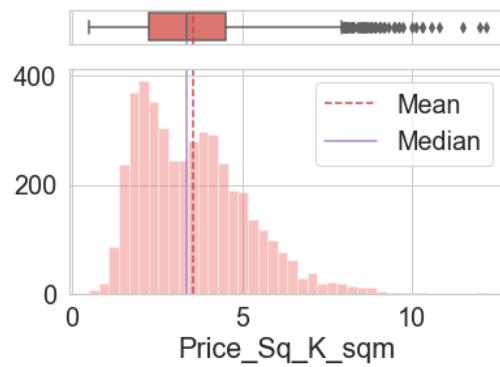
2020\_Feb\_March



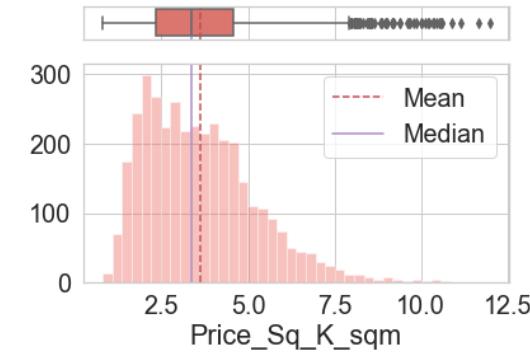
2020\_Jul\_Aug

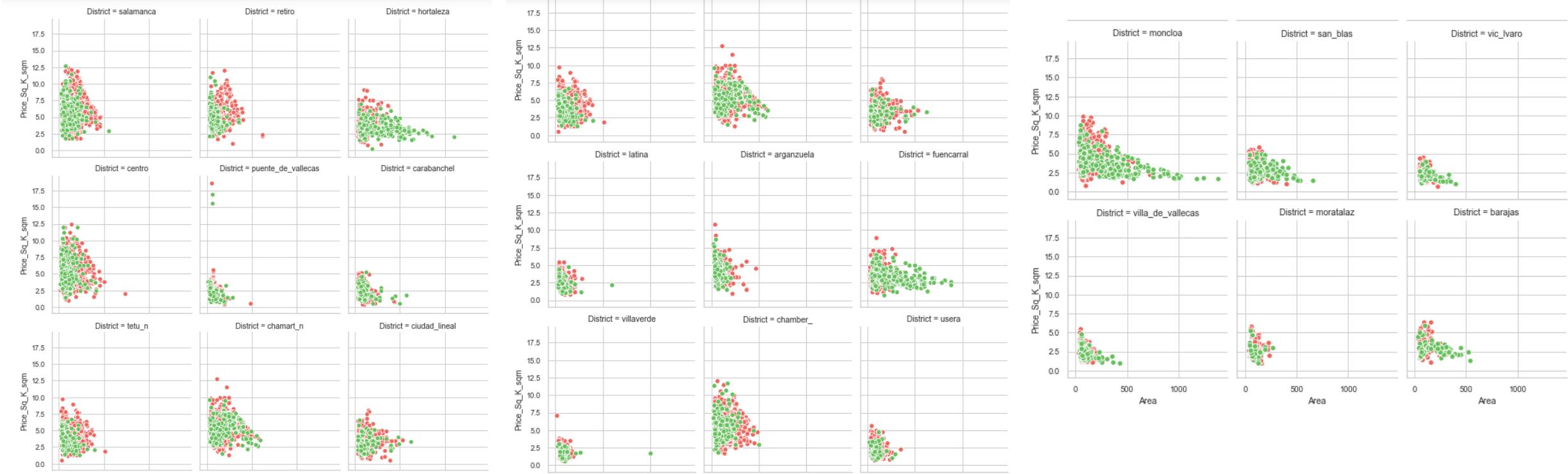


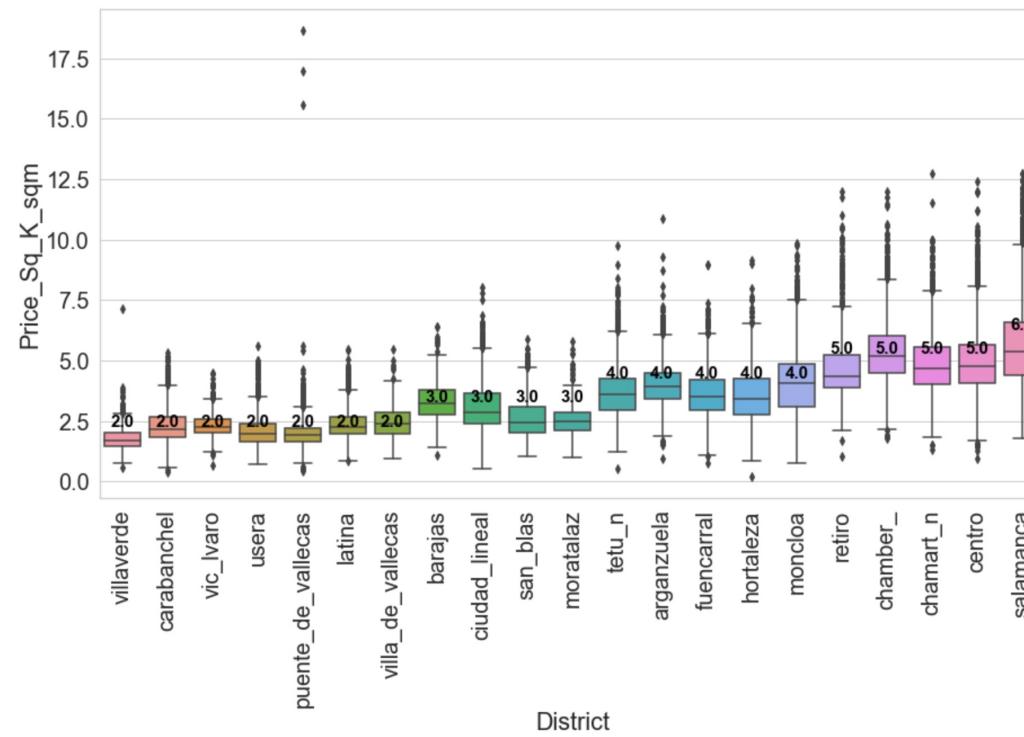
2020\_Ap\_May\_Jun

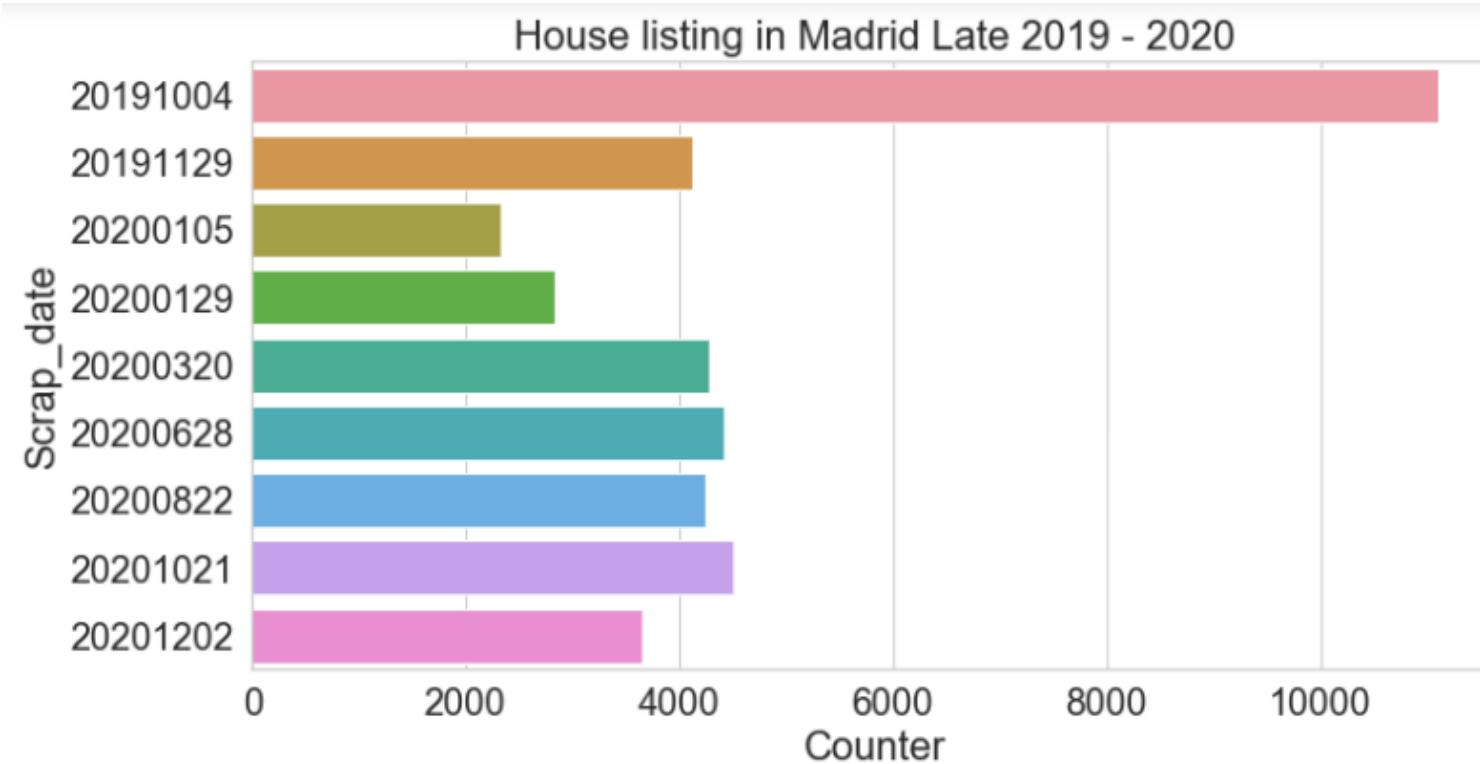


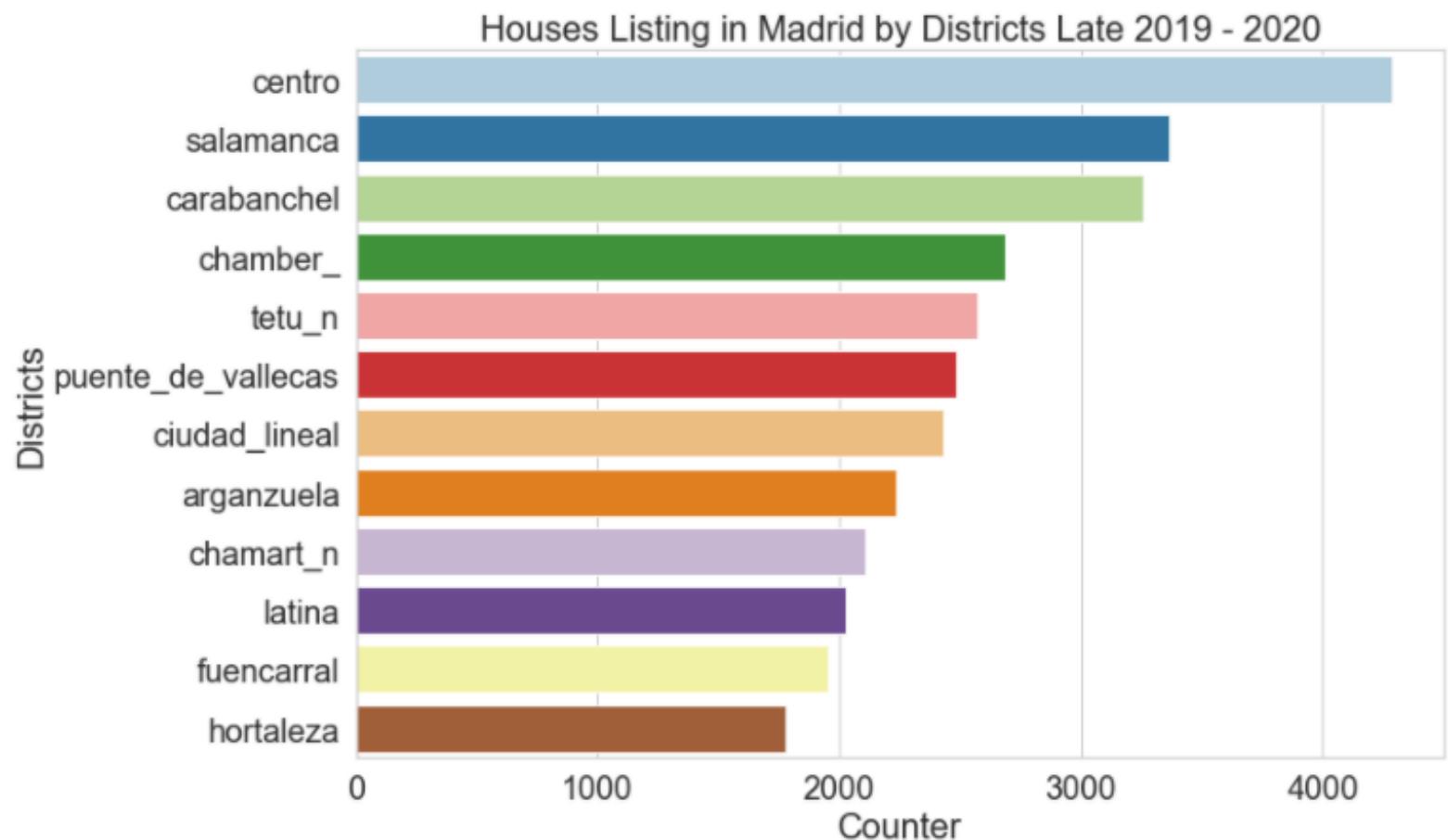
2020\_Ap\_Nov\_Dec

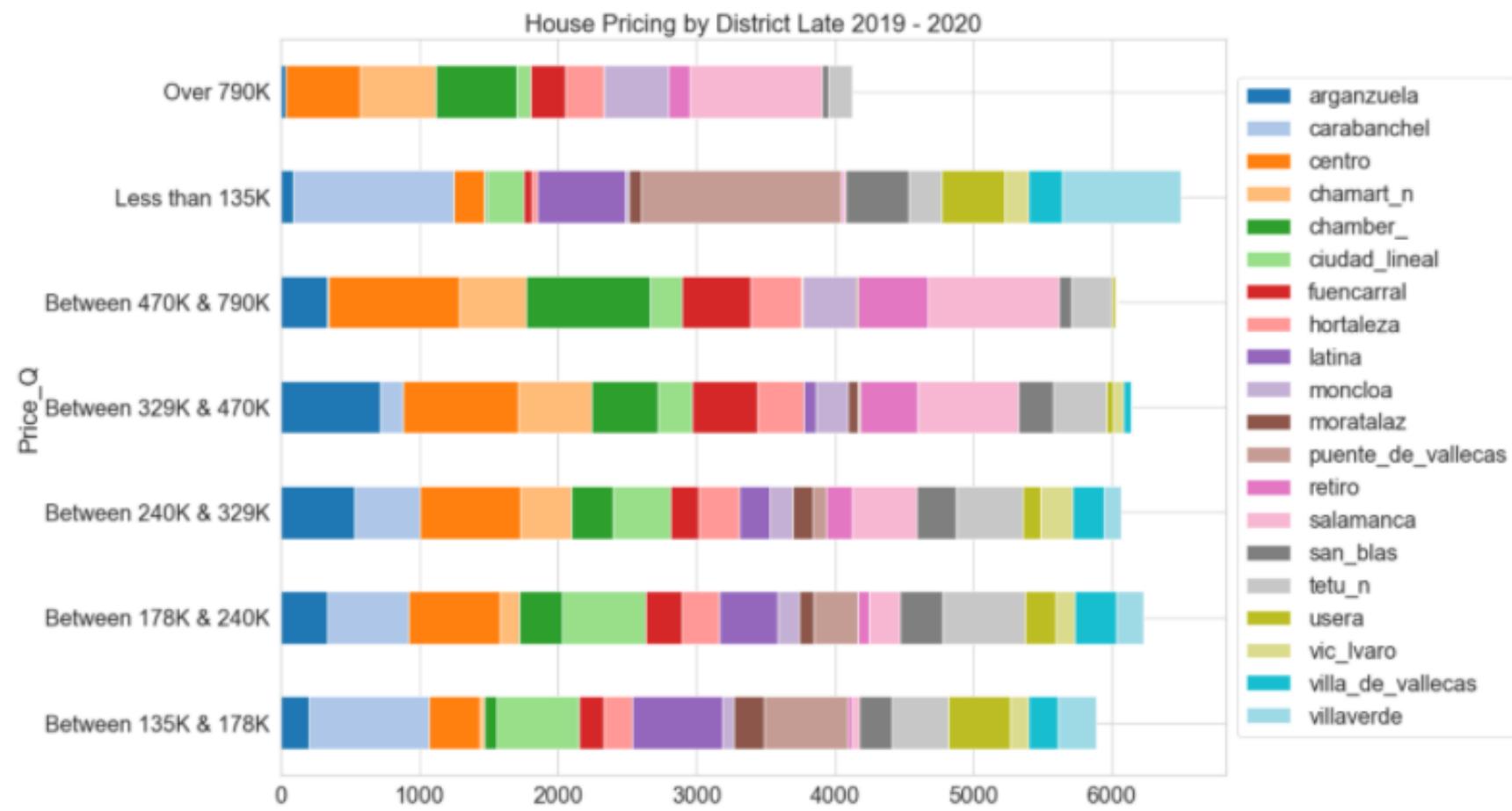


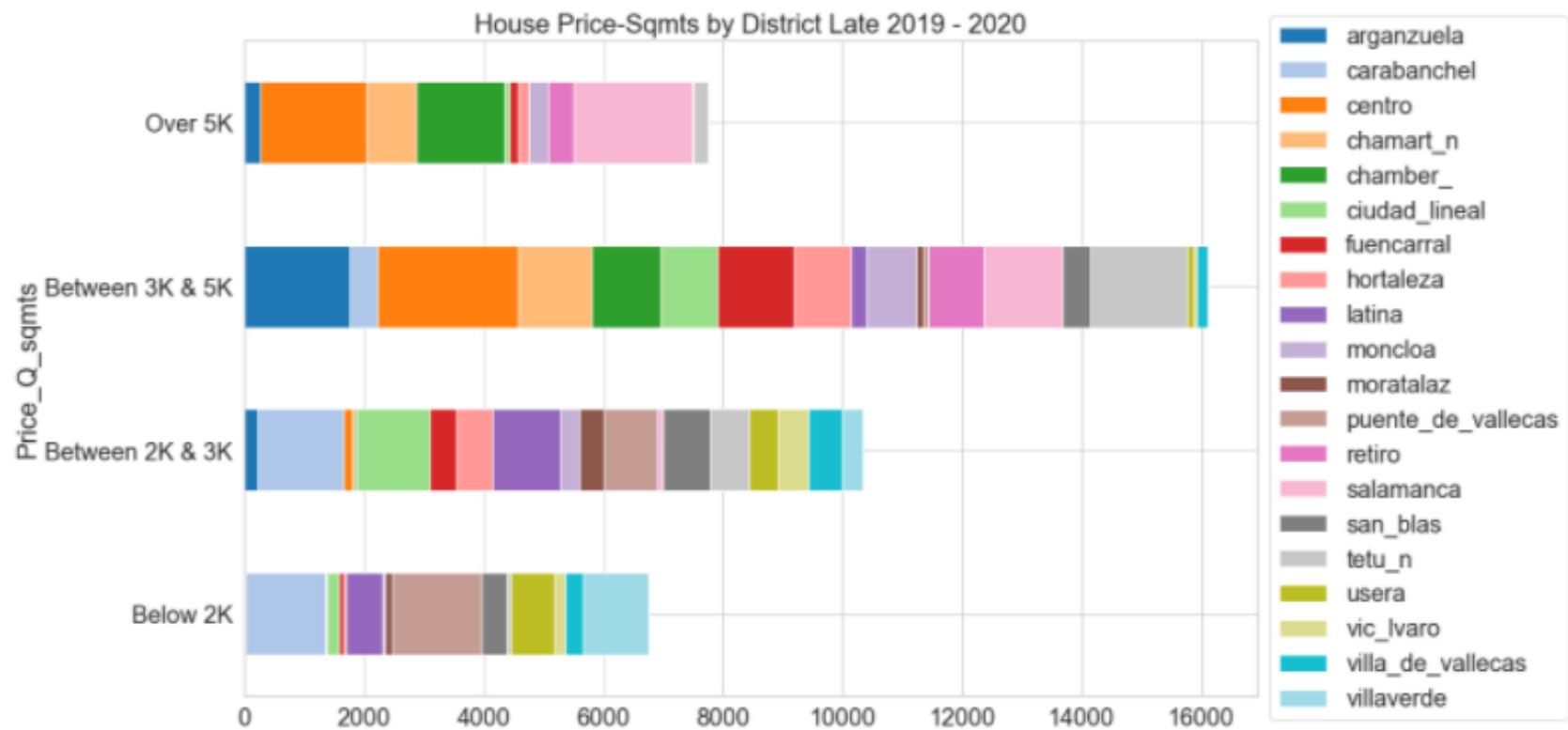


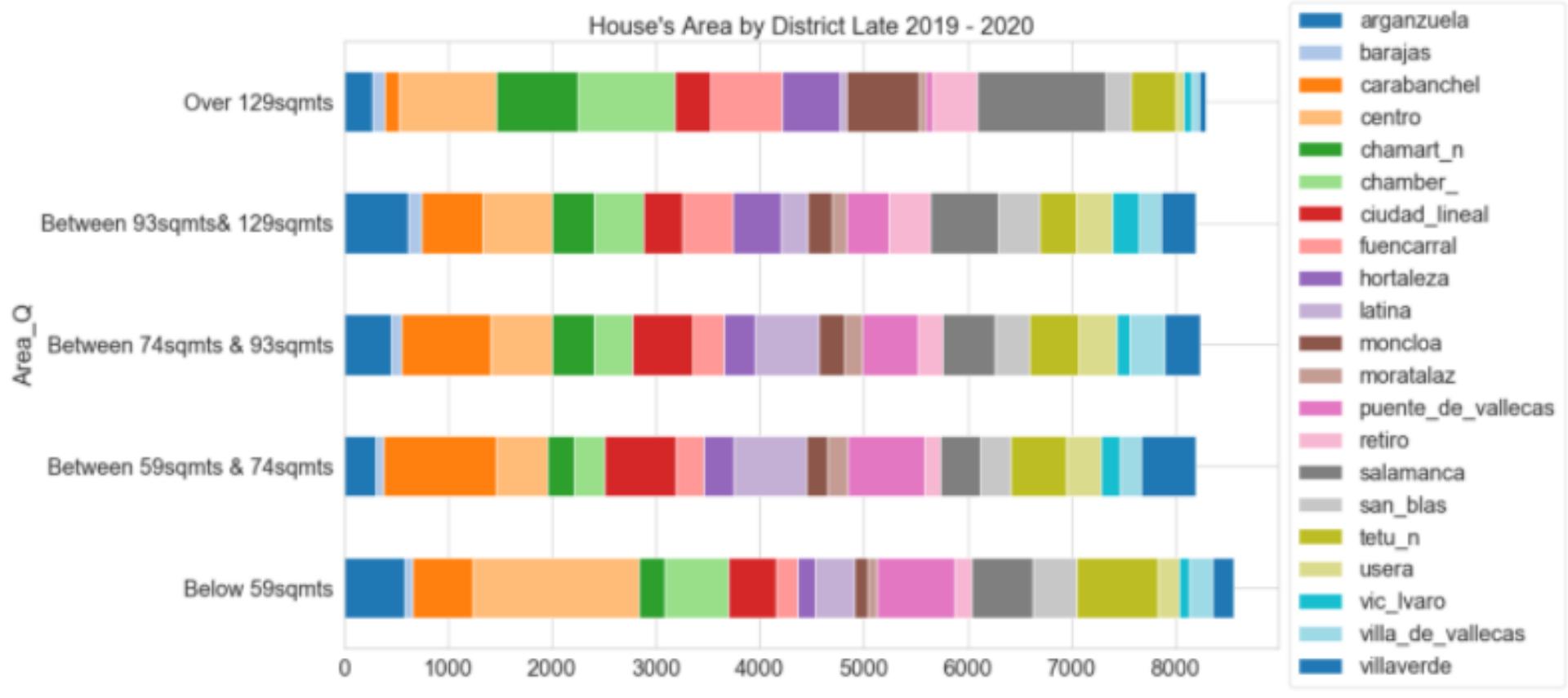












In [94]:

```
1 #Identificando el numero de clusters
2 from sklearn.cluster import KMeans
3 import matplotlib.pyplot as plt
4
5 #Metodo codo
6 max_k=15
7 inertia = []
8 for k in range(1, max_k):
9     kmeans = KMeans(n_clusters = k).fit(scaler_X)
10    inertia.append(kmeans.inertia_)
11
12 plt.plot(range(1, max_k), inertia, 'bx-')
13 plt.xlabel('k')
14 plt.ylabel(u'Dispersión')
15 plt.title("Clusters: Método del codo")
```

Out[94]: Text(0.5, 1.0, 'Clusters: Método del codo')

A line plot titled "Clusters: Método del codo". The y-axis is labeled "Dispersión" and ranges from 0 to 10,000 with major ticks at 2,500, 5,000, 7,500, and 10,000. The x-axis is labeled "k" and ranges from 0 to 15 with major ticks at 5, 10, and 15. A blue line with 'x' markers shows a sharp decrease in dispersion from approximately 10,000 at k=1 to about 2,500 at k=5, after which it levels off.

In [99]:

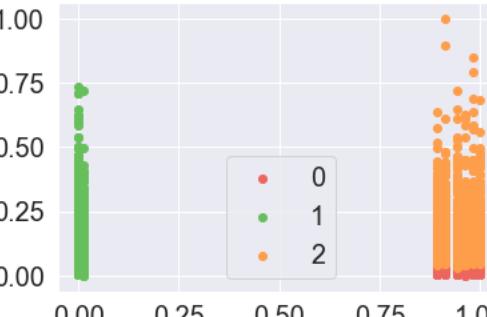
```
1 kmeans = KMeans(n_clusters=3, random_state=0)
2 label = kmeans.fit_predict(scaler_X)
3 print(len(kmeans.labels_))
4 print(label)
```

41467  
[1 1 1 ... 0 2 2]

In [102]:

```
1 #Getting unique labels
2
3 u_labels = np.unique(label)
4 print(u_labels)
5 #plotting the results:
6
7 for i in u_labels:
8     plt.scatter(scaler_X[label == i , 0] , scaler_X[label == i , 1] , label = i)
9 plt.legend()
10 plt.show()
```

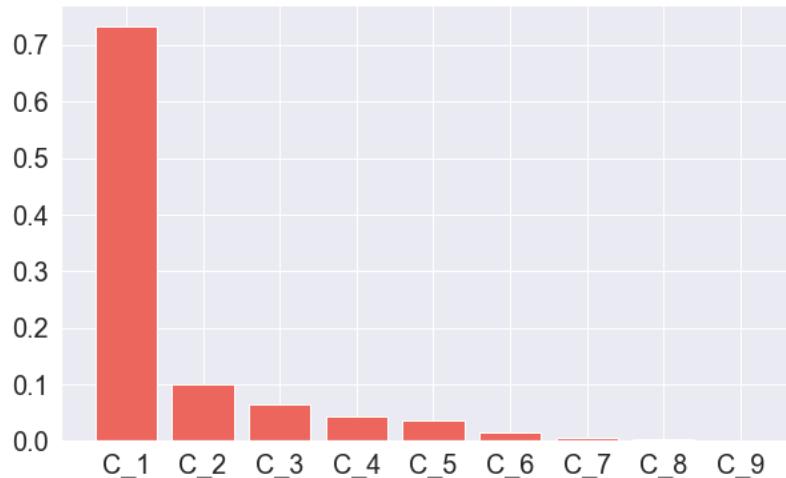
[0 1 2]



A scatter plot showing data points clustered into three groups. The x-axis ranges from 0.00 to 1.00 with ticks at 0.00, 0.25, 0.50, 0.75, and 1.00. The y-axis ranges from 0.00 to 1.00 with ticks at 0.00, 0.25, 0.50, 0.75, and 1.00. Three distinct vertical clusters are visible: cluster 0 (red dots) is at x ≈ 0.0; cluster 1 (green dots) is at x ≈ 0.1; and cluster 2 (orange dots) is at x ≈ 0.9. A legend in the bottom left corner identifies the colors for each cluster: red for 0, green for 1, and orange for 2.



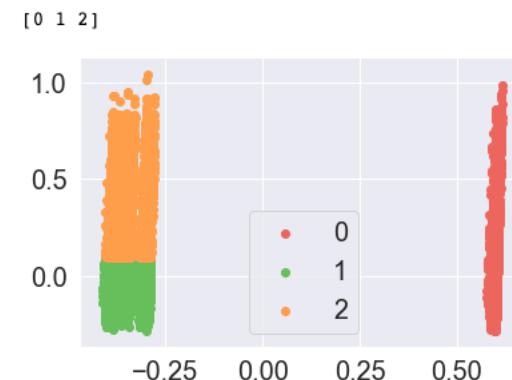
```
In [104]: 1 fig, ax = plt.subplots(figsize=[10, 6])
2 ax.bar(
3     [f"C_{i+1}" for i in range(dfpca.explained_variance_ratio_.size)],
4     dfpca.explained_variance_ratio_
5 )
6 plt.show()
```



```
In [105]: 1 for feature, coef in zip(VAR_NUM_2CHECK, dfpca.components_[0]): #va por filas y vas viendo
2     print(f"{feature}: {coef:.3f}")
```

```
Scrap_date: -1.000
Area: 0.004
Price_M: 0.015
Year: 0.013
Rooms_num: 0.003
Floor_num: 0.000
Garage_price: -0.000
...: 0.000
```

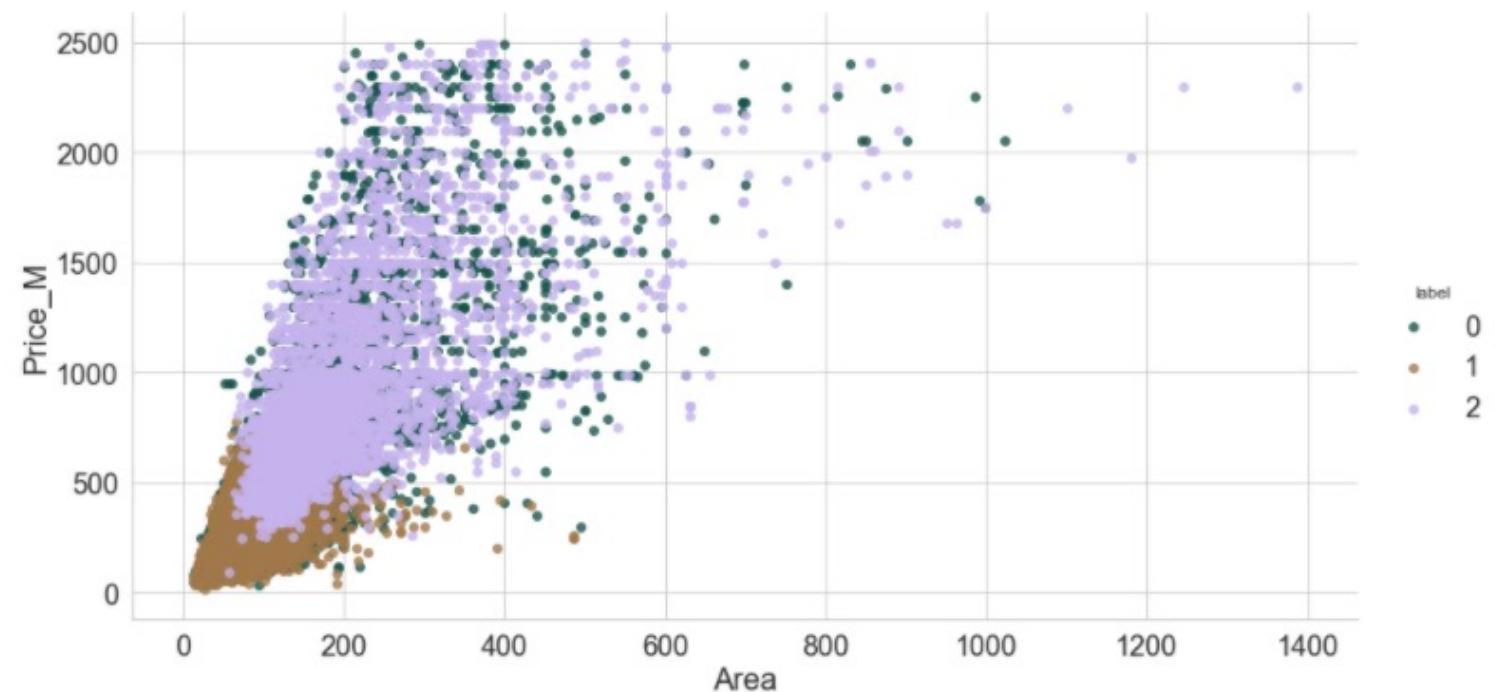
```
In [110]: 1 #Getting unique labels
2
3 u_labels = np.unique(label)
4 print(u_labels)
5 #plotting the results:
6
7 for i in u_labels:
8     plt.scatter(df_pca2_comp[label == i , 0] , df_pca2_comp[label == i , 1] , label = i)
9 plt.legend()
10 plt.show()
```

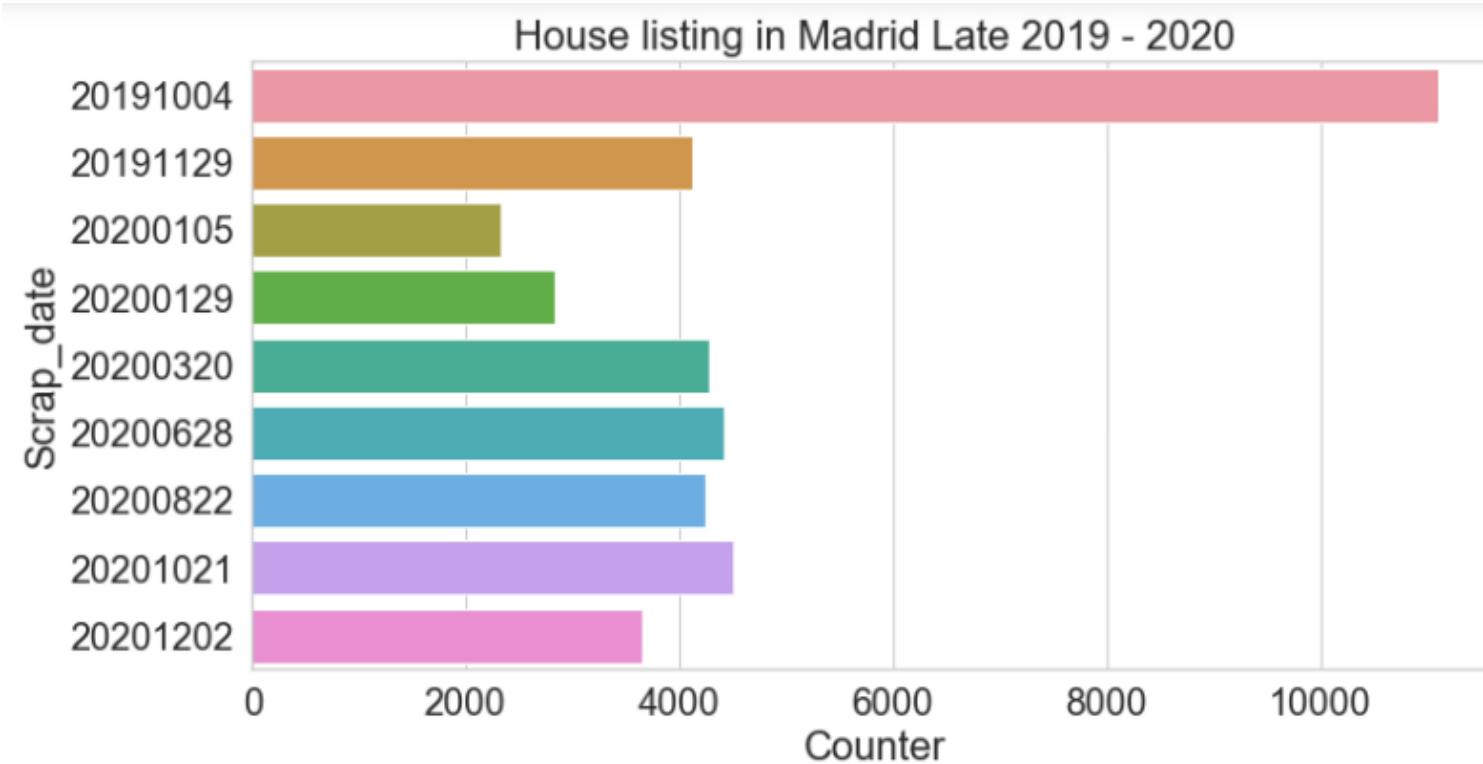




In [115]:

```
1 sns.set_style('whitegrid')
2 g=sns.lmplot('Area','Price_M',data=df3_num, hue='label',
3                 palette='cubehelix',size=12,aspect=1,fit_reg=False)
4 g.fig.set_figwidth(15.27)
5 g.fig.set_figheight(6.7)
```







# PROCESO DEL PROYECTO MACHINE LEARNING

