

Analisi per la predizione del prezzo di un auto usata

Programmazione di Applicazioni Data Intensive

Laurea in Ingegneria e Scienze Informatiche
DISI - Università di Bologna, Cesena

Ciandrimi Davide e Giorgetti Alex
davide.ciandrimi@studio.unibo.it
alex.giorgetti@studio.unibo.it

Parte 1 - Descrizione del problema e comprensione dei dati

Come abbiamo visto siamo pronti quindi a realizzare un modello che predica il prezzo di auto usate.
Abbiamo deciso di utilizzare un dataset presente su Kaggle al seguente indirizzo.
Questo Dataset è stato ricavato dal sito web cardekho sulla base di informazioni del 2020.

Caricamento librerie

- Andiamo a caricare tutte le librerie necessarie per lo svolgimento del progetto

```
In [ ]: %matplotlib inline
import os, joblib
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import matplotlib.pyplot as plt
import plotly.graph_objects as go
from plotly.subplots import make_subplots
import seaborn as sb
```

- Questo ci permette di nascondere tutti i warning presenti durante l'esecuzione del codice.

```
In [ ]: import warnings
from sklearn.exceptions import ConvergenceWarning
warnings.filterwarnings("ignore", category=ConvergenceWarning)
```

Caricamento dei dati

```
In [ ]: file_csv_url = "https://raw.githubusercontent.com/alexgiorgetti/vehicle-dataset-from-cardekho/main/cars.csv"
file_name = "cars.csv"

if not os.path.exists(file_name):
    from urllib.request import urlopen
    urlopen(file_csv_url, file_name)
```

```
In [ ]: data = pd.read_csv(file_name, sep = ",")
```

```
In [ ]: data.head(5)
```

```
Out [ ]: Car_Name  Year  Selling_Price  Present_Price  Kms_Driven  Fuel_Type  Seller_Type  Transmission  Owner
0  fiat  2014  3359  559  27000  Petrol  Dealer  Manual  0
1  citro  2013  475  954  43000  Diesel  Dealer  Manual  0
2  citro  2017  725  985  6900  Diesel  Dealer  Manual  0
3  wagonr  2011  285  415  5200  Petrol  Dealer  Manual  0
4  swift  2014  460  687  42450  Diesel  Dealer  Manual  0

• Analizziamo le dimensioni in memoria, la tipologia delle feature assegnate da pandas e le istanze non nulle.
```

```
In [ ]: data.info(memory_usage = "deep")

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 301 entries, 0 to 300
Data columns (total 9 columns):
 # Column Non-Null Count Dtype
---
 0 Car_Name 301 non-null object
 1 Year 301 non-null int64
 2 Selling_Price 301 non-null float64
 3 Present_Price 301 non-null float64
 4 Kms_Driven 301 non-null int64
 5 Fuel_Type 301 non-null object
 6 Seller_Type 301 non-null object
 7 Transmission 301 non-null object
 8 Owner 301 non-null int64
dtypes: float64(2), int64(3), object(4)
memory usage: 87.6 KB

• Come possiamo vedere Fuel_Type, Seller_Type e Transmission sono considerate come stringhe, ma sapendo che sono categoriche possiamo cambiare la categoria e ricaricare il DataFrame.
```

```
In [ ]: new_dtypes = {
    "Fuel_Type": "category",
    "Seller_Type": "category",
    "Transmission": "category"
}
data = pd.read_csv("cars.csv", dtype = new_dtypes)
```

- Possiamo vedere subito come l'occupazione in memoria è stata ridotta di oltre 50KB
- Questa scelta per Dataset come questi, non è rilevante, ma per altre molto più grandi può incidere particolarmente.

```
In [ ]: data.info(memory_usage = "deep")

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 301 entries, 0 to 300
Data columns (total 9 columns):
 # Column Non-Null Count Dtype
---
 0 Car_Name 301 non-null object
 1 Year 301 non-null int64
 2 Selling_Price 301 non-null float64
 3 Present_Price 301 non-null float64
 4 Kms_Driven 301 non-null int64
 5 Fuel_Type 301 non-null category
 6 Seller_Type 301 non-null category
 7 Transmission 301 non-null int64
 8 Owner 301 non-null int64
dtypes: category(3), float64(2), int64(3), object(1)
memory usage: 33.0 KB

• Impostando a Category possiamo vedere subito quali sono le possibili scelte, come nel caso del Fuel_Type:
```

```
In [ ]: data["Fuel_Type"].head(5)
```

```
Out [ ]: 0    Petrol
1     Diesel
2     Petrol
3     Petrol
4     Diesel
Name: Fuel_Type, dtype: category
Categories (3, object): ['CNG', 'Diesel', 'Petrol']
```

Significato delle feature

- Car_Name: Nome dell'auto
- Year: Anno di produzione
- Selling_Price: Prezzo richiesto dal venditore
- Present_Price: Prezzo di listino
- Kms_Driven: Km percorsi dall'auto
- Fuel_Type: Tipologia di carburante
- Seller_Type: Tipologia di venditore
- Transmission: Tipologia di cambio
- Owner: Numero di possessori

Le tipologie di dati presenti sono:

- Numeri interi per l'anno, il chilometraggio e possessore
- Numeri non interi per il prezzo
- String per il nome
- Categorici i restanti

Analisi delle singole feature

```
In [ ]: data.describe(include = "all")
```

```
Out [ ]: count    301    301    301.000000    301.000000    301.000000    301    301    301    301.000000
unique    98      NaN      NaN      NaN      NaN      3      2      2      NaN
top       city      NaN      NaN      NaN      NaN      Petrol  Dealer  Manual  NaN
freq      26      NaN      NaN      NaN      NaN      239    195    261    NaN
mean      NaN    2013.627907    4.661296    7.628472    36947.205980    NaN      NaN      NaN    0.043189
std       NaN    2.891554    5.082182    8.644115    38886.883882    NaN      NaN      NaN    0.247915
min       NaN    2003.000000    0.100000    0.320000    500.000000    NaN      NaN      NaN    0.000000
25%      NaN    2012.000000    0.500000    1.200000    15000.000000    NaN      NaN      NaN    0.000000
50%      NaN    2014.000000    3.000000    6.400000    32000.000000    NaN      NaN      NaN    0.000000
75%      NaN    2016.000000    6.000000    9.900000    48767.000000    NaN      NaN      NaN    0.000000
max       NaN    2018.000000    35.000000    92.600000    500000.000000    NaN      NaN      NaN    3.000000

• Grazie al metodo describe riusciamo a ottenere varie statistiche che ci possono aiutare ad analizzare il nostro Dataset.
• Come prima cosa possiamo vedere come l'Owner è tendente allo 0, questo ci fa capire come la maggior parte delle macchine ha avuto un solo proprietario.
• A seguire analizziamo la differenza tra Selling e Present medio, vedendo come si distacchi la media di circa 3.0. Nonostante si abbia un prezzo medio relativamente basso abbiamo picchi di massimo molto più alti, fino a 35.0
• Vediamo inoltre come quasi l'80% delle auto ha come carburante il Petrol e come cambio quello Manuale.
• Infine è bene notare che oltre il 60% delle automobili viene venduta da concessionari.
```

```
In [ ]: def numerical_graph():
fig = make_subplots(rows = 2, cols = 2, subplot_titles = ("Selling Price", "Present Price", "KMs", "Year"))
fig.add_trace(go.Histogram(x = data["Selling_Price"], name = "Selling Cost", row = 1, col = 1))
fig.add_trace(go.Histogram(x = data["Present_Price"], name = "Present Cost", row = 1, col = 2))
fig.add_trace(go.Histogram(x = data["Kms_Driven"], name = "Transmission", row = 2, col = 1))
fig.add_trace(go.Histogram(x = data["Year"], name = "Days", row = 2, col = 2))
fig.show()
```

- Qui di seguito sono riportati 4 istogrammi che analizzano la densità di frequenza, rispettivamente di Selling_Price, Present_Price, KMs e Years.

```
In [ ]: numerical_graph()
```

```
Out [ ]: <matplotlib.figure.Figure: 0x7f8b0c0e6d00>
```

```
Out [ ]: <matplotlib.figure.Figure: 0x7f8b0c0e6d00>
```

```
Out [ ]: <matplotlib.figure.Figure: 0x7f8b0c0e6d00>
```

```
Out [ ]: <matplotlib.figure.Figure: 0x7f8b0c0e6d00>
```

```
Out [ ]: <matplotlib.figure.Figure: 0x7f8b0c0e6d00>
```

```
Out [ ]: <matplotlib.figure.Figure: 0x7f8b0c0e6d00>
```

```
Out [ ]: <matplotlib.figure.Figure: 0x7f8b0c0e6d00>
```

```
Out [ ]: <matplotlib.figure.Figure: 0x7f8b0c0e6d00>
```

```
Out [ ]: <matplotlib.figure.Figure: 0x7f8b0c0e6d00>
```

```
Out [ ]: <matplotlib.figure.Figure: 0x7f8b0c0e6d00>
```

```
Out [ ]: <matplotlib.figure.Figure: 0x7f8b0c0e6d00>
```

```
Out [ ]: <matplotlib.figure.Figure: 0x7f8b0c0e6d00>
```

```
Out [ ]: <matplotlib.figure.Figure: 0x7f8b0c0e6d00>
```

```
Out [ ]: <matplotlib.figure.Figure: 0x7f8b0c0e6d00>
```

```
Out [ ]: <matplotlib.figure.Figure: 0x7f8b0c0e6d00>
```

```
Out [ ]: <matplotlib.figure.Figure: 0x7f8b0c0e6d00>
```

```
Out [ ]: <matplotlib.figure.Figure: 0x7f8b0c0e6d00>
```

```
Out [ ]: <matplotlib.figure.Figure: 0x7f8b0c0e6d00>
```

```
Out [ ]: <matplotlib.figure.Figure: 0x7f8b0c0e6d00>
```

```
Out [ ]: <matplotlib.figure.Figure: 0x7f8b0c0e6d00>
```

```
Out [ ]: <matplotlib.figure.Figure: 0x7f8b0c0e6d00>
```

```
Out [ ]: <matplotlib.figure.Figure: 0x7f8b0c0e6d00>
```

```
Out [ ]: <matplotlib.figure.Figure: 0x7f8b0c0e6d00>
```

```
Out [ ]: <matplotlib.figure.Figure: 0x7f8b0c0e6d00>
```

```
Out [ ]: <matplotlib.figure.Figure: 0x7f8b0c0e6d00>
```

```
Out [ ]: <matplotlib.figure.Figure: 0x7f8b0c0e6d00>
```

```
Out [ ]: <matplotlib.figure.Figure: 0x7f8b0c0e6d00>
```

```
Out [ ]: <matplotlib.figure.Figure: 0x7f8b0c0e6d00>
```

```
Out [ ]: <matplotlib.figure.Figure: 0x7f8b0c0e6d00>
```

```
Out [ ]: <matplotlib.figure.Figure: 0x7f8b0c0e6d00>
```

```
Out [ ]: <matplotlib.figure.Figure: 0x7f8b0c0e6d00>
```

```
Out [ ]: <matplotlib.figure.Figure: 0x7f8b0c0e6d00>
```

```
Out [ ]: <matplotlib.figure.Figure: 0x7f8b0c0e6d00>
```

```
Out [ ]: <matplotlib.figure.Figure: 0x7f8b0c0e6d00>
```

```
Out [ ]: <matplotlib.figure.Figure: 0x7f8b0c0e6d00>
```

```
Out [ ]: <matplotlib.figure.Figure: 0x7f8b0c0e6d00>
```

```
Out [ ]: <matplotlib.figure.Figure: 0x7f8b0c0e6d00>
```

```
Out [ ]: <matplotlib.figure.Figure: 0x7f8b0c0e6d00>
```

```
Out [ ]: <matplotlib.figure.Figure: 0x7f8b0c0e6d00>
```

```
Out [ ]: <matplotlib.figure.Figure: 0x7f8b0c0e6d00>
```

```
Out [ ]: <matplotlib.figure.Figure: 0x7f8b0c0e6d00>
```

```
Out [ ]: <matplotlib.figure.Figure: 0x7f8b0c0e6d00>
```

```
Out [ ]: <matplotlib.figure.Figure: 0x7f8b0c0e6d00>
```

```
Out [ ]: <matplotlib.figure.Figure: 0x7f8b0c0e6d00>
```

```
Out [ ]: <matplotlib.figure.Figure: 0x7f8b0c0e6d00>
```

```
Out [ ]: <matplotlib.figure.Figure: 0x7f8b0c0e6d00>
```

```
Out [ ]: <matplotlib.figure.Figure: 0x7f8b0c0e6d00>
```

```
Out [ ]: <matplotlib.figure.Figure: 0x7f8b0c0e6d00>
```

```
Out [ ]: <matplotlib.figure.Figure: 0x7f8b0c0e6d00>
```

```
Out [ ]: <matplotlib.figure.Figure: 0x7f8b0c0e6d00>
```

```
Out [ ]: <matplotlib.figure.Figure: 0x7f8b0c0e6d00>
```

```
Out [ ]: <matplotlib.figure.Figure: 0x7f8b0c0e6d00>
```

```
Out [ ]: <matplotlib.figure.Figure: 0x7f8b0c0e6d00>
```

```
Out [ ]: <matplotlib.figure.Figure: 0x7f8b0c0e6d00>
```

```
Out [ ]: <matplotlib.figure.Figure: 0x7f8b0c0e6d00>
```

```
Out [ ]: <matplotlib.figure.Figure: 0x7f8b0c0e6d00>
```

```
Out [ ]: <matplotlib.figure.Figure: 0x7f8b0c0e6d00>
```

```
Out [ ]: <matplotlib.figure.Figure: 0x7f8b0c0e6d00>
```

```
Out [ ]: <matplotlib.figure.Figure: 0x7f8b0c0e6d00>
```

```
Out [ ]: <matplotlib.figure.Figure: 0x7f8b0c0e6d00>
```

```
Out [ ]: <matplotlib.figure.Figure: 0x7f8b0c0e6d00>
```

```
Out [ ]: <matplotlib.figure.Figure: 0x7f8b0c0e6d00>
```

```
Out [ ]: <matplotlib.figure.Figure: 0x7f8b0c0e6d00>
```

```
Out [ ]: <matplotlib.figure.Figure: 0x7f8b0c0e6d00>
```

```
Out [ ]: <matplotlib.figure.Figure: 0x7f8b0c0e6d00>
```

```
Out [ ]: <matplotlib.figure.Figure: 0x7f8b0c0e6d00>
```

```
Out [ ]: <matplotlib.figure.Figure: 0x7f8b0c0e6d00>
```

```
Out [ ]: <matplotlib.figure.Figure: 0x7f8b0c0e6d00>
```

```
Out [ ]: <matplotlib.figure.Figure: 0x7f8b0c0e6d00>
```

```
Out [ ]: <matplotlib.figure.Figure: 0x7f8b0c0e6d00>
```

```
Out [ ]: <matplotlib.figure.Figure: 0x7f8b0c0e6d00>
```

```
Out [ ]: <matplotlib.figure.Figure: 0x7f8b0c0e6d00>
```

```
Out [ ]: <matplotlib.figure.Figure: 0x7f8b0c0e6d00>
```

```
Out [ ]: <matplotlib.figure.Figure: 0x7f8b0c0e6d00>
```

```
Out [ ]: <matplotlib.figure.Figure: 0x7f8b0c0e6d00>
```

```
Out [ ]: <matplotlib.figure.Figure: 0x7f8b0c0e6d00>
```

```
Out [ ]: <matplotlib.figure.Figure: 0x7f8b0c0e6d00>
```

```
Out [ ]: <matplotlib.figure.Figure: 0x7f8b0c0e6d00>
```

```
Out [ ]: <matplotlib.figure.Figure: 0x7f8b0c0e6d00>
```

```
Out [ ]: <matplotlib.figure.Figure: 0x7f8b0c0e6d00>
```

```
Out [ ]: <matplotlib.figure.Figure: 0x7f8b0c0e6d00>
```

```
Out [ ]: <matplotlib.figure.Figure: 0x7f8b0c0e6d00>
```

```
Out [ ]: <matplotlib.figure.Figure: 0x7f8b0c0e6d00>
```

```
Out [ ]: <matplotlib.figure.Figure: 0x7f8b0c0e6d00>
```

```
Out [ ]: <matplotlib.figure.Figure: 0x7f8b0c0e6d00>
```

```
Out [ ]: <matplotlib.figure.Figure: 0x7f8b0c0e6d00>
```

```
Out [ ]: <matplotlib.figure.Figure: 0x7f8b0c0e6d00>
```

```
Out [ ]: <matplotlib.figure.Figure: 0x7f8b0c0e6d00>
```

```
Out [ ]: <matplotlib.figure.Figure: 0x7f8b0c0e6d00>
```

```
Out [ ]: <matplotlib.figure.Figure: 0x7f8b0c0e6d00>
```

```
Out [ ]: <matplotlib.figure.Figure: 0x7f8b0c0e6d00>
```

```
Out [ ]: <matplotlib.figure.Figure: 0x7f8b0c0e6d00>
```

```
Out [ ]: <matplotlib.figure.Figure: 0x7f8b0c0e6d00>
```

```
Out [ ]: <matplotlib.figure.Figure: 0x7f8b0c0e6d00>
```

```
Out [ ]: <matplotlib.figure.Figure: 0x7f8b0c0e6d00>
```

```
Out [ ]: <matplotlib.figure.Figure: 0x7f8b0c0e6d00>
```

```
Out [ ]: <matplotlib.figure.Figure: 0x7f8b0c0e6d00>
```

```
Out [ ]: <matplotlib.figure.Figure: 0x7f8b0c0e6d00>
```

```
Out [ ]: <matplotlib.figure.Figure: 0x7f8b0c0e6d00>
```

```
Out [ ]: <matplotlib.figure.Figure: 0x7f8b0c0e6d00>
```

```
Out [ ]: <matplotlib.figure.Figure: 0x7f8b0c0e6d00>
```


• Coefficiente R^2

Accuratezza e validazione

```
In [ ]:
for model in models:
    print(f"[model] model:\n-----")
    print_eval(X_val, y_val, models[model] ("Model"))
    print("-----\n")

Lasso model:
-----
Mean squared error: 1.5866
Relative error: 40.009875
R-squared coefficient: 0.95315
-----

Ridge model:
-----
Mean squared error: 0.39059
Relative error: 19.896801
R-squared coefficient: 0.98896
-----

Elastic_Net model:
-----
Mean squared error: 1.1705
Relative error: 44.393451
R-squared coefficient: 0.96691
-----

XGBoost model:
-----
Mean squared error: 0.00011913
Relative error: 0.736401
R-squared coefficient: 1.0
-----
```

- Vediamo come i modelli sono molto simili tra loro, il **Ridge** supera di accuratezza il **Lasso** e l'**Elastic net** come appurato precedentemente. Ma abbiamo una precissima previsione con il modello di **XGBoost**.
- Definiamo dunque una funzione che ci permette di stampare, dato il nome del modello, l'andamento della predizione confrontando due grafici: *Histplot* e *Scatter*

Grafico delle predizione

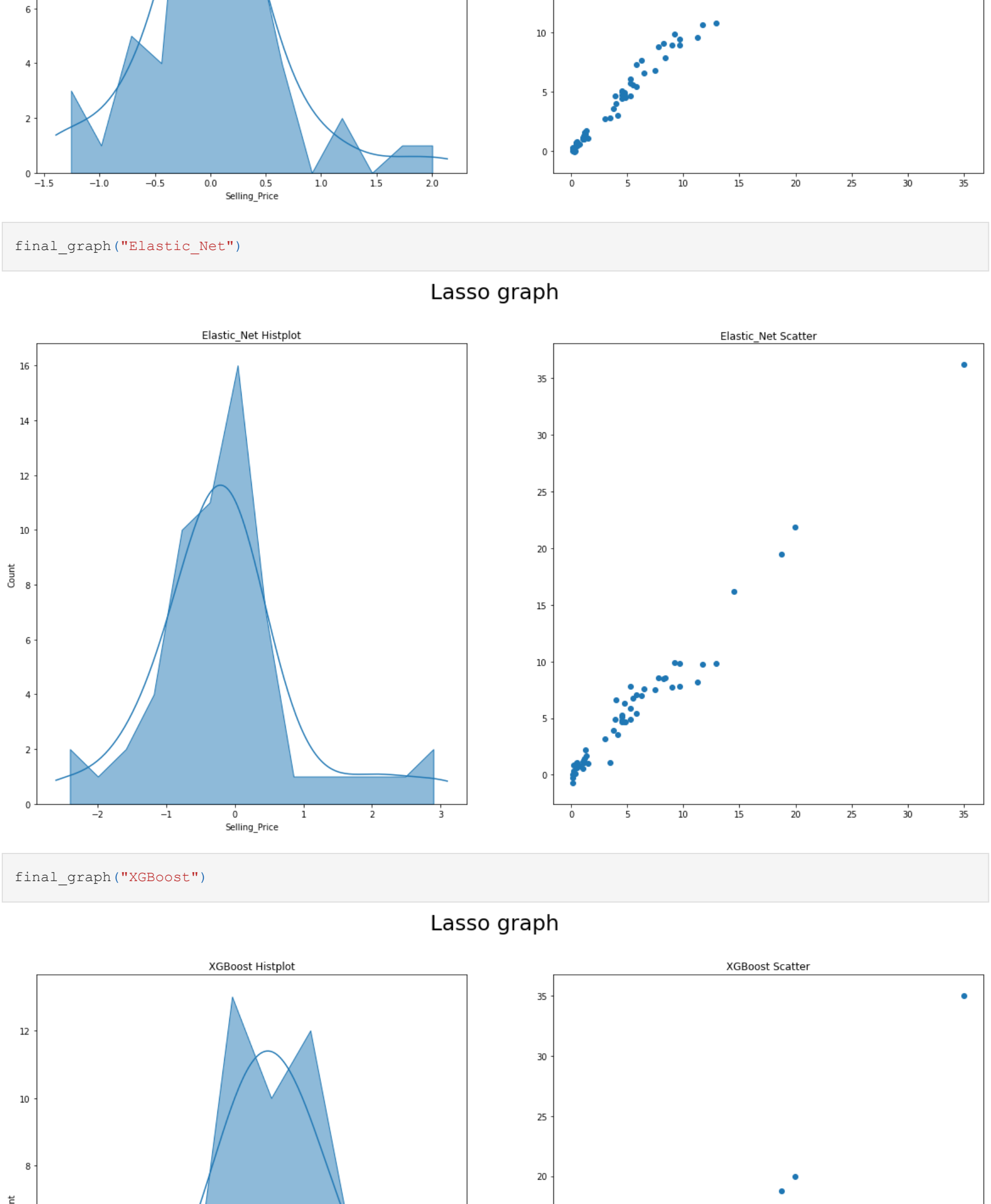
```
In [ ]:
for model in models:
    models[model] ["Predict"] = models[model] ["Model"].predict(X_val)

def final_graph(model_name):
    fig, axes = plt.subplots(nrows=1, ncols=2, sharex=False, sharey=False, figsize=(20, 10))

    fig.suptitle(f"Lasso graph", fontsize = 25)

    sns.histplot(y_val ~ models[model_name] ("Predict"), ax = axes[0], kde = True, element = "poly")
    axes[0].scatter(y_val, models[model_name] ("Predict"))
    axes[0].set_title(f"{model_name} Histplot")
    axes[1].set_title(f"{model_name} Scatter")

In [ ]:
final_graph("Lasso")
```



Parte 5 - Analisi conclusiva

Come fase conclusiva andiamo a stampare per ogni modello la predizione delle automobili presenti nel nostro set di Validation. Vediamo per ognuna il distacco dal valore reale, traendo le conclusioni sul modello migliore.

```
In [ ]:
data = {'Selling Prices': y_val.head(15),
        'XGBoost_Prediction': models["XGBoost"] ["Model"].predict(X_val)[:15],
        'Lasso_Prediction': models["Lasso"] ["Model"].predict(X_val)[:15],
        'Ridge_Prediction': models["Ridge"] ["Model"].predict(X_val)[:15],
        'ElasticNet_Prediction': models["Elastic_Net"] ["Model"].predict(X_val)[:15],
        'Car_Name': car_names[y_val.head(15).index]
}

pd.DataFrame(data).sort_index()
```

	Selling Prices	XGBoost_Prediction	Lasso_Prediction	Ridge_Prediction	ElasticNet_Prediction	Car_Name
5	9.25	9.249281	9.872710	9.836662	9.880587	vitara brezza
9	7.45	7.431100	6.942149	6.791941	7.469575	ciaz
59	19.99	19.952527	23.980674	21.228615	21.821440	fortuner
79	14.50	14.513296	14.680132	14.090274	16.163223	fortuner
95	5.85	5.852381	6.680513	7.242625	7.088878	corolla altis
99	9.65	9.648474	8.084605	9.409883	9.802041	fortuner
111	1.15	1.138304	1.236572	1.324835	1.302149	Royal Enfield Thunder 350
156	0.48	0.464438	0.387538	0.714636	1.023642	TVS Sport
166	0.45	0.460858	0.592039	0.654478	0.845466	Hero Passion Pro
175	0.38	0.375909	0.307899	0.014015	0.118334	Hero Honda CBZ extreme
186	0.25	0.240769	0.883431	0.237410	0.860238	TVS Wego
197	0.16	0.164642	0.184031	-0.019941	-0.234039	Honda CB twister
267	8.35	8.353922	8.020238	7.872209	8.559634	city
280	5.25	5.228798	4.496806	4.606818	4.875252	brio
283	8.99	8.994455	8.937322	8.938721	7.762694	city

Verifichiamo dalla tabella come le previsioni effettuate non si discostano più di tanto dal valore reale, ponendo particolare attenzione a una previsione molto accurata con il modello **XGBoost**. Siamo particolarmente soddisfatti del risultato ottenuto, poichè non ci aspettavamo una tale soluzione, probabilmente data dalle dimensioni ridotte del Dataset utilizzato (300 istanze)