

# Predicting Diamonds Price

Statistical Learning Mod. B, Final project

Bernardi Alberto, Cracco Gianmarco

23<sup>rd</sup> June, 2020



UNIVERSITÀ  
DEGLI STUDI  
DI PADOVA

**Main Goal**  $\Rightarrow$  build an appropriate model to predict prices of diamonds.

In particular we will make use of:

- 1 Feature Selection;
- 2 Shrinkage Methods;
- 3 Principal Component Analysis.

**Main Goal**  $\Rightarrow$  build an appropriate model to predict prices of diamonds.

In particular we will make use of:

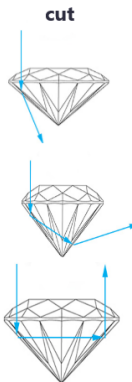
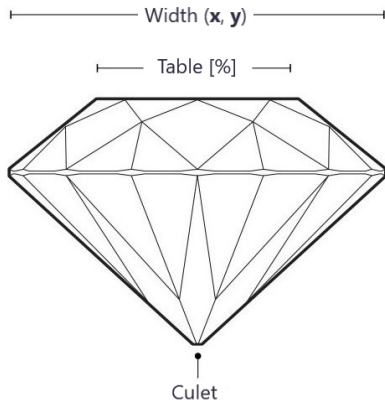
- 1** Feature Selection;
- 2** Shrinkage Methods;
- 3** Principal Component Analysis.

# kaggle

Diamonds Dataset from *kaggle*

<https://www.kaggle.com/shivam2503/diamonds>

# Diamond Structure



## clarity



Around **54000** diamonds with the following attributes:

- **Carat** (*numerical*): weight of a diamond (1 carat = 0.2 g);
- **Cut** (*factor*): *Fair, Good, Very Good, Excellent, Ideal*;
- **Color** (*factor*): *J, I, H, G, F, E, D*;
- **Clarity** (*factor*): *I1, SI2, SI1, VS2, VS1, VVS2, VVS1* and *IF*;
- **Depth** (*numerical*): percentual height over average width;
- **Table** (*numerical*): percentual width of the table over total width;
- **Price** [\$] (*numerical*)
- **x, y, z** [mm] (*numerical*): three spatial dimensions.

Cleaning and preparing the data through the following:

- Check for **Null Values**;
- One-Hot-Encoding;
- Removal of **absurd data**;
- Check for possible outliers.

Cleaning and preparing the data through the following:

- Check for **Null Values**;
- **One-Hot-Encoding**;
- Removal of **absurd data**;
- Check for possible **outliers**.



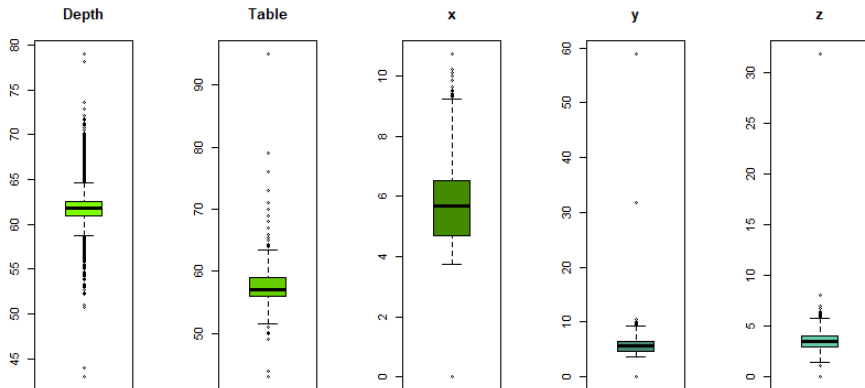
Cleaning and preparing the data through the following:

- Check for **Null Values**;
- **One-Hot-Encoding**;
- Removal of **absurd data**;
- Check for possible outliers.

Cleaning and preparing the data through the following:

- Check for **Null Values**;
- **One-Hot-Encoding**;
- Removal of **absurd data**;
- Check for possible **outliers**.

# Boxplots



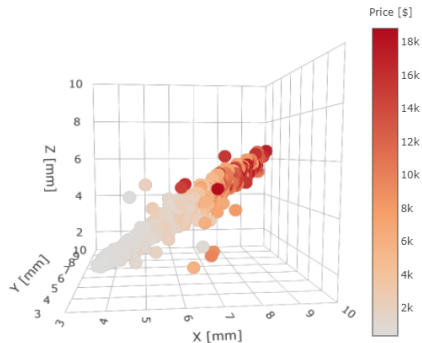
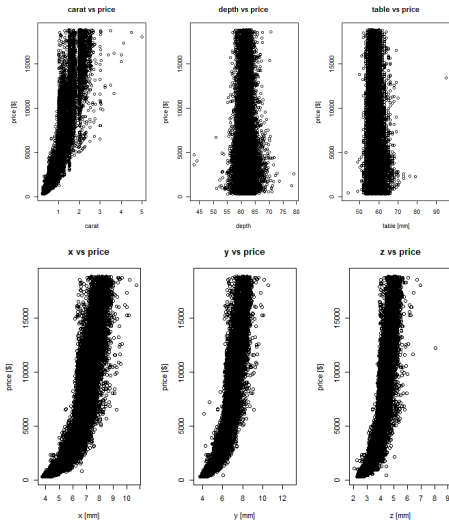
For a more complete data visualization, take a look at the attached paper and **R** code.

- Goal: Price prediction;
- How? Linear models  $\Rightarrow$  limited capacity, high interpretability;
- Why? Costumer advisory, market analysis, insurance quote.

- Goal: Price prediction;
- How? Linear models  $\Rightarrow$  **limited** capacity, **high** interpretability;
- Why? Costumer advisory, market analysis, insurance quote.

- Goal: Price prediction;
- How? Linear models  $\Rightarrow$  **limited** capacity, **high** interpretability;
- Why? Customer advisory, market analysis, insurance quote.

# Explorative - Price vs Numerical



# Explorative - Price vs Categorical

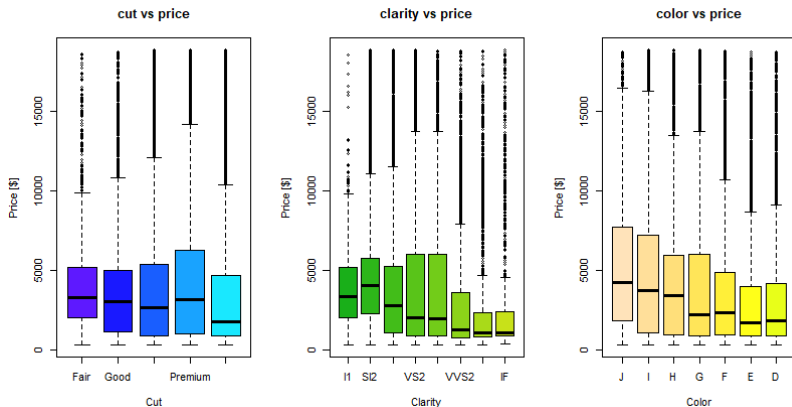


Figure: Boxplots of categorical features



First approach: we normalized our data (except for the target variable) and we tried to fit a **linear model**:

- Coefficients correctly describe the *expected* relationship between features and target variable;
- Analysis of residuals however shows they are clearly not randomly distributed, especially for the largest predicted values.

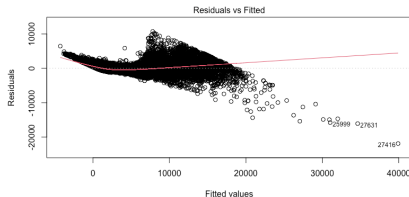
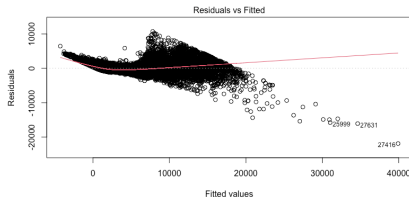


Figure: Plot of residuals for the naive approach

First approach: we normalized our data (except for the target variable) and we tried to fit a **linear model**:

- Coefficients correctly describe the *expected* relationship between features and target variable;
- Analysis of residuals however shows they are clearly not randomly distributed, especially for the largest predicted values.



**Figure:** Plot of residuals for the naive approach

To overcome this problem we followed two steps:

- 1  $\log(\textit{price})$ : residuals become flat and uniformly distributed but now the coefficients are "wrong";



- 2  $\log(\textit{carat})$ : the coefficients are correct again.

To overcome this problem we followed two steps:

- 1  $\log(\textit{price})$ : residuals become flat and uniformly distributed but now the coefficients are "wrong";



- 2  $\log(\textit{carat})$ : the coefficients are correct again.

# A new variable



ROUND



OVAL



MARQUISE



PEAR



HEART



EMERALD



PRINCESS



RADIANT

We introduce a **new variable**  $r = \frac{x}{y}$



Now we can capture the **shape** of the diamond.

# A new variable



ROUND



OVAL



MARQUISE



PEAR



HEART



EMERALD



PRINCESS



RADIANT

We introduce a **new variable**  $r = \frac{x}{y}$



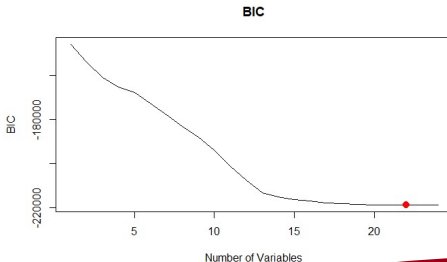
Now we can capture the **shape** of the diamond.

We apply different techniques and criteria to select only the truly relevant features for our prediction.

⇒ a **simpler** and more **efficient** model is obtained!

The final **reduced linear model** is given by:

$$\log(\text{price}) \sim \log(\text{carat}) + \text{cut} + \text{color} + \text{clarity} + x + y + r.$$

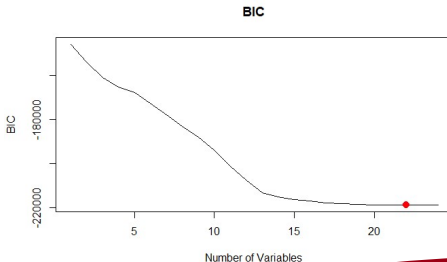


We apply different techniques and criteria to select only the truly relevant features for our prediction.

⇒ a **simpler** and more **efficient** model is obtained!

The final **reduced linear model** is given by:

$$\log(\text{price}) \sim \log(\text{carat}) + \text{cut} + \text{color} + \text{clarity} + x + y + r.$$





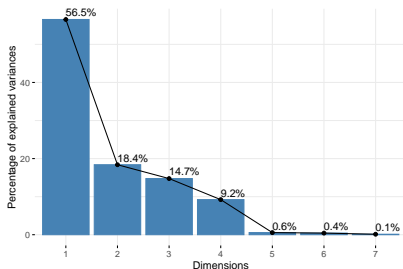


A plot titled "Ridge Regression" showing the relationship between the logarithm of the regularization parameter  $\lambda$  (x-axis) and the estimated coefficients (y-axis). The x-axis is labeled "Log Lambda" and ranges from 6 to -6. The y-axis is labeled "Coefficients" and ranges from 0.0 to 1.0. A vertical dashed blue line is drawn at  $\log(\lambda) \approx -4.8$ , which corresponds to  $\lambda \approx 0.0025$  on the top x-axis. The plot displays multiple regression lines for different values of  $\lambda$ , showing how the coefficients change as  $\lambda$  varies. The lines are colored in a gradient from dark grey to bright yellow. As  $\lambda$  increases (moving right on the x-axis), the coefficients shrink towards zero. As  $\lambda$  decreases (moving left on the x-axis), the coefficients diverge from zero, with some reaching values near 1.0 and others near 0.0.

One last method applied to obtain a model: **PCA**.

It is a technique that affects the **data**, not the model itself, and reduce their dimensionality.

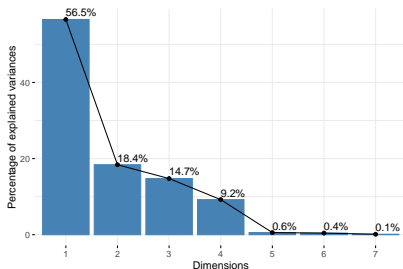
From **7** numerical features to **4** components!



One last method applied to obtain a model: **PCA**.

It is a technique that affects the **data**, not the model itself, and reduce their dimensionality.

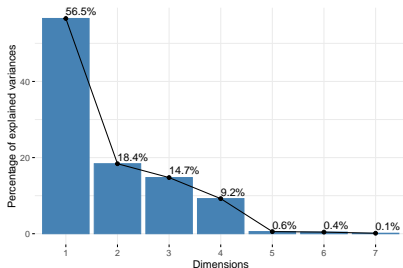
From **7** numerical features to **4** components!



One last method applied to obtain a model: **PCA**.

It is a technique that affects the **data**, not the model itself, and reduce their dimensionality.

From **7** numerical features to **4** components!



Once PCA was performed  $\Rightarrow$  categorical variables have been added again.

Outcomes:

- Importance of categorical features;
- Importance of  $\log(\text{carat})$ ;
- Not enough dimensions to be really effective.

Once PCA was performed  $\Rightarrow$  categorical variables have been added again.

Outcomes:

- Importance of **categorical features**;
- Importance of  $\log(\textit{carat})$ ;
- Not enough dimensions to be really effective.

Once PCA was performed  $\Rightarrow$  categorical variables have been added again.

Outcomes:

- Importance of **categorical features**;
- Importance of  **$\log(carat)$** ;
- Not enough dimensions to be really effective.

Once PCA was performed  $\Rightarrow$  categorical variables have been added again.

Outcomes:

- Importance of **categorical features**;
- Importance of  **$\log(carat)$** ;
- **Not enough dimensions** to be really effective.



Summary of what we obtained in this project:

	MAE [\$]
Naive linear model	740
<b>Linear model (reduced)</b>	<b>418</b>
Ridge regression	472
Lasso regression	448
PCA	1245
PCA + cat. vars.	875
PCA with $\log(\text{carat})$	989
PCA with $\log(\text{carat})$ + cat. vars.	564

Mean Absolute Error: the smaller, the better.

In the attached report you can find way deeper insights with much more technical results.

Anyway, there is still much space for improvement!

- Non-Linear models;
- Machine Learning;
- Deep Learning.

In the attached report you can find way deeper insights with much more technical results.

Anyway, there is still much space for improvement!

- Non-Linear models;
- Machine Learning;
- Deep Learning.