# Data Crunchers

## Used Car Explaratory Data Analysis

MEF University Big Data Analytics

Data Cruncher's Team

# Overall View of Data

Used Car Database' scraped from Ebay Kleinanzeigen (in German)

370.000 second-hand cars

40 unique brands.

20 variables

Month of Registration

Year of Registration

Kilometer

Power PS

Postal Code

Seller ✗

Price

Vehicle Type

Last Seen

Offer Type ✗

AB Test ✗

Model

Brand

Gear Box

Number of Pictures ✗

Fuel Type

Date Created

Not Repaired Damage

# Content

1   Raw Data Visualization

2   Check missing data and other mistakes

3   Create a list of outliers or other anomalies

4   Data Cleaning

5   Clean Data Visualization & Mapping

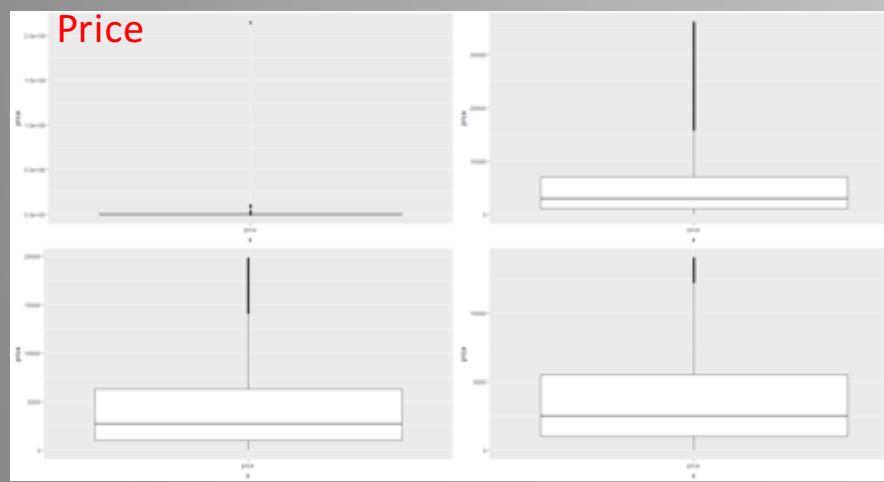6   Determine relationships among the variables

7   **Regression Models**

8   Maps

# Content

"This is not what I meant when I said 'we need better data cleansing!'"

Price

```
quantile(auto$price, 0.90)
quantile(auto$price, 0.99)
quantile(auto$price, 0.01)
quantile(auto$price, 0.05)
quantile(auto$price, 0.10)

ggplot(aes(x=vehicleType, y=price), data = auto) +
  geom_boxplot() +
  ylim(quantile(auto$price, 0.05), quantile(auto$price, 0.95))

p1 <- ggplot(aes(x="price", y=price), data = auto) +
  geom_boxplot()

p2 <- ggplot(aes(x="price", y=price), data = auto) +
  geom_boxplot() +
  ylim(0, quantile(auto$price, 0.99))

p3 <- ggplot(aes(x="price", y=price), data = auto) +
  geom_boxplot() +
  ylim(0, quantile(auto$price, 0.95))

p4 <- ggplot(aes(x="price", y=price), data = auto) +
  geom_boxplot() +
  ylim(0, quantile(auto$price, 0.90))

library(gridExtra)
grid.arrange(p1, p2, p3, p4, ncol = 2)
```
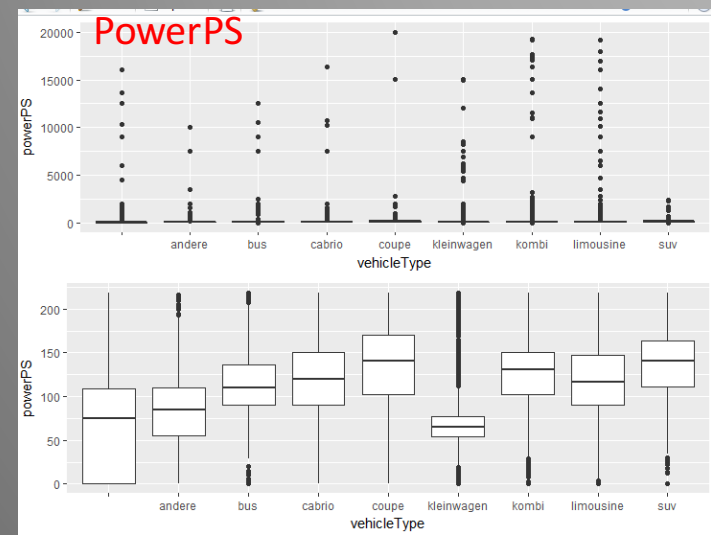
PowerPS

|            | andere | benzin | cng | diesel | elektro | hybrid | lpg  |
|------------|--------|--------|-----|--------|---------|--------|------|
|            | 4674   | 10     | 3011 | 11    | 1141    | 1      | 0    | 67 |
| andere     | 267    | 19     | 1178 | 12    | 1354    | 11     | 2    | 27 |
| bus        | 914    | 5      | 9274 | 238   | 17042   | 1      | 5    | 503 |
| cabrio     | 746    | 5      | 16882 | 3    | 1459    | 6      | 0    | 230 |
| coupe      | 664    | 1      | 11924 | 2    | 1975    | 4      | 16   | 289 |
| kleinwagen | 4019   | 18     | 63257 | 77   | 6611    | 47     | 34   | 480 |
| kombi      | 2911   | 21     | 26088 | 131  | 32358   | 5      | 19   | 1140 |
| limousine  | 3488   | 30     | 57499 | 37   | 24405   | 7      | 129  | 1621 |
| suv        | 366    | 5      | 4302  | 3    | 6242    | 0      | 8    | 540 |

Year of Reg.

1
2
3
4
5 Clean Data Visualization
6
7
8

# Content

Histogram of Engine Power (PowerPS)


Gearbox




Fuel Type Frequency Diagram

# Content

Car Age Histogram





Histogram for Selling Time

# Content

Price vs. Vehicle Type

# Content

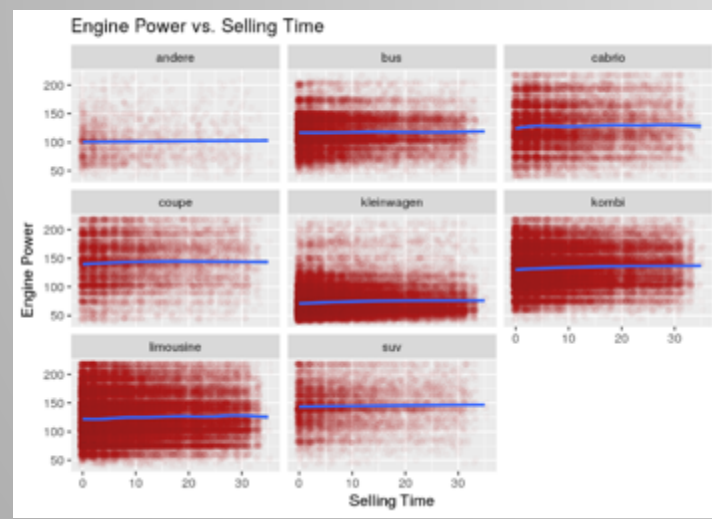# Content

1
2
3
4
5
6 Determine relationships among the variables
7
8



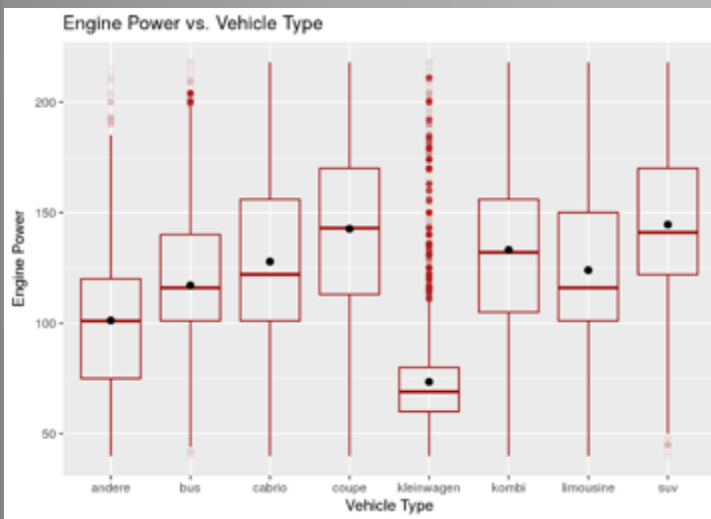## Engine Power vs. Vehicle Type

Engine Power vs. Vehicle Type boxplot with Vehicle Type on the x-axis (andere, bus, cabrio, coupe, kleinwagen, kombi, limousine, suv) and Engine Power on the y-axis (50, 100, 150, 200).

# Content

# Content

Price vs. Gearbox



Price vs. Vehicle Type by Gearbox

Automatic gearbox cars are more expensive than the manual ones but the difference becomes less significant in cheap cars, i.e. Kleinwagen.

# Content

# Content

1
2
3
4
5
6 Determine relationships among the variables
7
8

In all vehicle types, the price continues to decrease between 0-20 years (20 years is the lowest point) but starts increase after between 20-30 years.



Price vs. Age by Vehicle Type

# Content

In general, automatic gearbox cars are more expensive than the manual ones. This is especially significant in hybrid cars.



Price vs. Gearbox by Fuel Type

# Content

Only electric cars with manual gearbox have superior engine power performance to the automatic ones. This divergence is one of the most interesting things we have

# Content

Selling time does not seem to be affected by the combinations of gearbox and fuel type as well, with the exception of CNG cars.
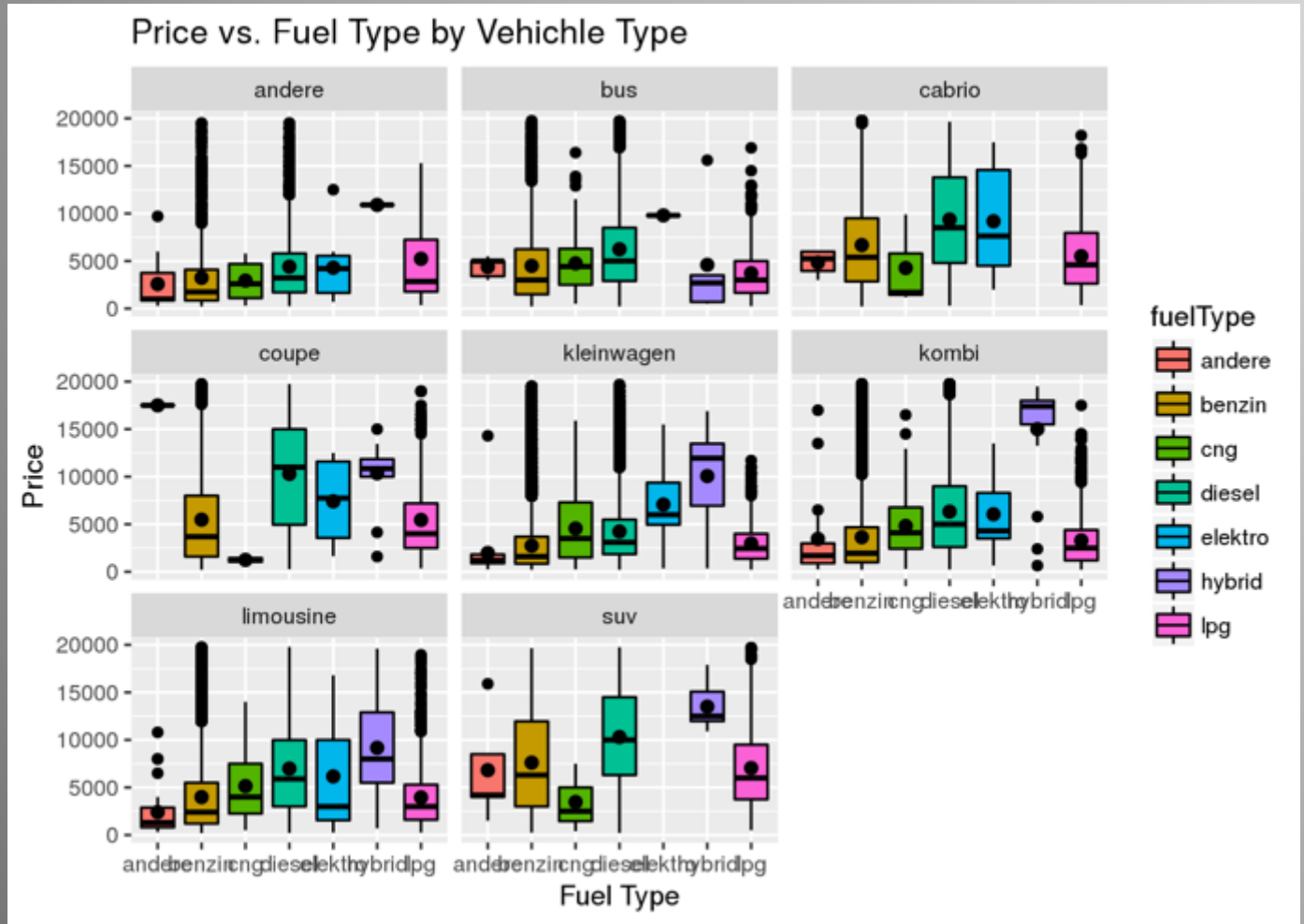
The electric and CNG cars show longer selling time trend which may indicate that second-hand car market for hybrid cars have not matured yet.

# Content

In cheap cars (Kleinwagen), electric-hybrid fuel type makes a significant increase in price.



Price vs. Fuel Type by Vehichle Type

# Content

1
2
3
4
5

6 Determine relationships
among the variables

7
8

Majority of the cars are sold only within 10 days in all vehicle types. For the first 20 days, whenever the price goes up, the selling time increases as well but this correlation stabilize in day 20. This trend is especially visible in SUV vehicles.

# Content

1
2
3
4
5
6 Determine relationships among the variables
7
8

Top 10 Brands (%80 second hand sales)– Most popular in the second-hand car market

# Content

1

2

3

4

5

6 Determine relationships among the variables

7

8

# Content

## Correlations

# Content

# Actual vs Predicted Price

# Content

# Content

# Content

- Limousine, kombi and kleinwagen are the most popular vehicle types in the second-hand market. Most expensive cars are SUV's while the cheapest ones are kleinwagens.
- On average Kleinwagen vehicle type is the cheapest and has the lowest engine power. But it also shows the most outliers – might be as a result of brand-model diversity.
- The most popular brands are Volkswagen, BMW, Opel, Mercedes, Audi, Ford, Renault, Peugeot, Fiat and Seat. These 10 brand correspond to almost 80% of the cars. (Originally our dataset contains around 40 brands)
- Most of the cars in the second-hand market are above 100.000 km, even 150.000 km. People does not frequently change cars according to our data set.
- Majority of the second-hand cars are sold only within 35 days. The ratio of the first 10 days (day 0 stands for same day sale) is quite high. This shows us that either Ebay-Kleinanzeigen is very successful at targeting customers or the second-hand market is more fluid that we actually thought.

# Content

- To our surprise, there is no strong/significant correlation between selling time and vehicle type, kilometer and price. We saw that whenever price goes up the change to be sold in 10-20 days increases especially in SUV vehicles (rather than 0-10 days) but this is not a general trend.
- Hybrid (electro engine, CNG) second-hand car market is emerging but shows longer selling time trend.
- In all vehicle types, the price continues to decrease between 0-20 years (20 years is the lowest point) but starts increase after between 20-30 years. Maybe, 20+ year old second-hand cars can be considered as 'antique' and users' emotional attachment may cause abnormalities.
- Hybrid cars with manual gearbox have superior engine power performance. This is a divergence from all correlations and one of the most interesting things we have found.
- According to our regression analysis, age (39%), kilometer(%23) and engine power(%19) are the most important factors explaining second hand price.
- Zip code analysis shows that kombi is the most popular second hand car when the zip code is provided.
- East Coast second hand car market is bigger.