

# Data Science Projekt

TINF18B 2020, 4347707 und 2335257.

## Inhaltsverzeichnis

Links funktionieren nur im Jupyter Notebook

- [Business Understanding](#)
- [Daten Vorverarbeitung](#)
- [Data Exploration](#)
  - [Visualisierungen rund um den Hauspreis](#)
  - [Visualisierungen rund um den Zustand](#)
- [Modelling](#)
  - [Preis-Regression](#)
  - [Zustands-Klassifikation](#)
- [Evaluation](#)
  - [Inferenz](#)
  - [Anwendung an einer weiteren Datei](#)

## Business Understanding

Die durchgestrichenen Punkte haben sich im Verlauf des Verkaufs des Verkaufspreises (Preis) geändert. Die Studierenden ihrer iterative Arbeitsweise deutlich machen.

Das übergeordnete Ziel eines Investors ist seinen Profit durch den Handel mit Immobilien zu erwirtschaften. Dazu muss er ein Haus teuer verkaufen, als er es eingekauft hat (zusätzliche Investitionen in das Haus eingeschlossen).

### Ziele

1. (gegeben) Der Käufer mehr Verständnis und das Haus des Verkaufspreises (Preis)
2. Verständnis darüber, welche Hauseigenschaften besonders viel Einfluss auf den Kaufpreis haben. Dieses Verständnis hilft dabei, auszuwählen, welche Komponenten renoviert werden sollen.
  - Vorhersege, welcher Monat ist der beste, um ein Haus zu kaufen- zu verkaufen?
3. Vorhersege des Zustands aus den anderen Hauseigenschaften bzw. Verständnis darüber ob und von was der Zustand abhängt?

### Motivation für diese Ziele

1. Die passende Vorhersage des Preises hilft bei Preisverhandlungen sowohl im Ankauf als auch im Verkauf.
2. Der Profit des Investors steigt.
3. Der Investor macht weniger Verlustgeschäfte. Fehlkäufe oder steckt zusätzliche Mittel in unwichtige Verbesserungen.
4. Wenn ein Investor den eigentlichen Zustand berechnen kann, ist es schon die Lage zu entscheiden ob der Preis dem Haus angemessen ist. Dadurch kann er entscheiden ob er das Haus kaufen oder nicht kaufen sollte.

### Anforderungen an die Ergebnisse

1. Preisvorhersage muss in der richtigen Größenordnung erfolgen.
2. Die Vorhersage des Preises soll in unter einer Minute berechnet werden.

### Beispielhafte Antworten

1. Kaufpreis: 192140
2. Besonders wichtige Eigenschaften eines Hauses: Klimaanlage, Schlafzimmer, Garagenkapazität
3. Besonders unwichtige Eigenschaften eines Hauses: Heizung, Zustand Fassade

## Vorverarbeitungen

### Imports

```
In [19]: import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
import mpl_toolkits
from IPython.display import display
from sklearn.preprocessing import StandardScaler
from sklearn.model_selection import train_test_split
from sklearn.preprocessing import LabelEncoder
import sklearn.metrics as skm
from sklearn.ensemble import GradientBoostingRegressor, GradientBoostingClassifier, RandomForestRegressor
from sklearn.linear_model import LinearRegression, RidgeCV, LassoCV
from sklearn.linear_model import LinearRegression, LassoCV, RidgeCV
from sklearn.metrics import mean_squared_error, cross_val_predict
from sklearn.metrics import mean_squared_error
from imblearn.over_sampling import RandomOverSampler
from pandas.api.types import CategoricalDtype
import pprint
import matplotlib.pyplot as plt
print("Importing finished")

Importing finished
```

### Daten Import

```
In [20]: data_folder = Path("DatenAusgegeben1.0.csv")
data = pd.read_csv(data_folder, encoding='cp852', sep=";")
data.describe()
```

```
Out [20]:
```

	Grundstück in qm	Zustand	Gebaut	Renoviert	Zustand Fassade	Kellerfläche in qm	Erster Stock in qm	Zweiter Stock in qm	Wohnfläche in qm	Schlafzimmer
count	2000.000000	2000.000000	2000.000000	2000.000000	2000.000000	2000.000000	2000.000000	2000.000000	2000.000000	2000.000000
mean	950.054000	5.697500	2099.031000	2113.344500	3.102500	96.140500	106.48100	31.125000	137.990000	2.883
std	737.437654	1.129439	29.120114	20.168198	0.386094	38.249893	34.40918	38.461551	45.565953	0.8171
min	121.000000	1.000000	2005.000000	2080.000000	1.000000	0.000000	31.000000	0.000000	31.000000	0.0000
25%	701.750000	5.000000	2083.000000	2095.000000	3.000000	74.000000	82.000000	0.000000	102.000000	2.0000
50%	887.000000	5.000000	2101.000000	2121.000000	3.000000	91.000000	99.000000	0.000000	133.000000	3.0000
75%	1078.000000	6.000000	2126.000000	2132.000000	3.000000	117.000000	126.000000	65.000000	161.000000	3.0000
max	19997.000000	9.000000	2140.000000	2140.000000	5.000000	298.000000	365.000000	174.000000	401.000000	8.0000

## Daten Vorbereiten

Minimale Data Preparation, weil sonst die Plots nicht funktionieren.

```
In [21]: data['Pool'] = data['Pool'].fillna(0, inplace=True) #setze Na durch Null (Haus besitzt wahrscheinlich keinen Pool)

In [22]: data['Garage Typ'] = data['Garage Typ'].fillna("keine", inplace=True) #setze Na durch 'keine' (Haus besitzt wahrscheinlich keine Garage)
```

## Data Exploration

### Vorliegende Daten

Es liegen Hauseigenschaften und der Verkaufspreis vor. Weitere Informationen können der Datenbeschreibung entnommen werden. Die Daten liegen in vielen verschiedenen Formaten vor. Das Datenset liefert 2000 Datensätze. Insgesamt ist die Datenqualität allerdings sehr gut.

### Probleme

Einige Spalten sind Nominal- oder Ordinalskalen mit Strings. Abhängig vom Modell müssen diese Informationen in Zahlenwerte transformiert werden.

Die Attribute „Pool“ und „Garage Typ“ sind mit „Na“ gefüllt. Außerdem hat „Pool“ den Datentyp der Gleitkommazahl.

Zusätzlich wäre folgende Daten sinnvoll: Landkarten, Kriminalitätsstatistiken, Luftverschmutzungsdaten, Information, ob ein Haus bewohnt ist und Raumumzugspläne der Stadt.

## Zusammenfassung der Erkenntnisse zur Preisvorhersage

Folgende Faktoren korrelieren mit einem höheren Verkaufspreis:

- Heizungsqualität exzellent
- Klimaanlage vorhanden
- eingebaute oder angebaute Garage
- Haustyp ein Familienhaus oder Reihendendhaus
- Große Wohnfläche
- ein Pool in sehr gutem Zustand
- geringe Bebauungsdichte
- Grundstücke oder Sportzentrum in der Nähe
- Haus liegt in einem leeren Bezirk
- Küche mit hoher Qualität

Faktoren die dagegen kaum Auswirkung auf die Höhe des Verkaufspreis haben sind:

- das Verkaufsjahr
- Verkaufsmotivat
- Grundstücksform
- Steigung
- Anzahl der Schlafzimmer

Manche Faktoren beeinflussen die Höhe des Preises auch nur bis zu einem bestimmten Wert. Ein Beispiel hierfür ist der Zustand der Fassade: ein Zustand von 3 führt ein besserer Zustand nicht mehr so deutlich zu einer Preissteigerung. Ähnlich ist es bei dem geeigneten Zustand und der Anzahl der Räume.

## Zusammenfassung der Erkenntnisse zur Zustandsvorhersage

Folgende Faktoren korrelieren mit dem allgemeinen Hauszustand:

- Klimaanlage
- Baujahr
- Garage
- Bebauungsdichte
- Grundstücksform
- Haustyp
- Fassade Zustand
- Küchenqualität

Es ist anzumerken, dass die Korrelationen überwiegend schwach sind. Es kann kein konkretes Verständnis über die Zusammensetzung des Attributs „Zustand“ gemacht werden.

Im nächsten Teil Modellierung wird dennoch untersucht, wie gut eine Vorhersage des Zustands sein kann. Ob eine Vorhersage generell sinnvoll ist, kann ohne weitere Informationen über das Attribut „Zustand“ nicht gemacht werden.

## Allgemeine Erkenntnisse

Da das Attribut Räume und Wohnfläche in qm stark korrelieren sollte nur maximal eines der beiden Attribute in einem dafür anfallenden Modell genutzt werden. Ähnlich ist es zwischen Kellerfläche und Fläche im ersten Stock.

Eine Vorhersage, welcher Monat der Beste ist, um ein Haus zu kaufen/verkaufen, kann aus den Daten nicht gemacht werden. Dementsprechend ist der Monat nicht entscheidend für den Verkaufspreis.

## Plots rund um den Hauspreis

```
In [23]: plt.figure(figsize=(10, 7))
sns.distplot(data, x='Preis')
plt.title('Verteilung Preis')
```

Out [23]: Text(0.5, 1, 'Verteilung Preis')

### Erkenntnis - Preis vs Quadratmeter

Man sieht deutlich die Tendenz, dass je mehr Quadratmeter ein Haus hat, desto höher ist der Preis. Allerdings gibt es ein paar Ausreißer, die entweder besonders teuer sind (bei um die 250 qm und 500000 Preisenheiten), oder sehr günstig.

```
In [24]: plt.figure(figsize=(10, 7))
plt.scatter(data['Wohnfläche in qm'], data['Preis'], marker='r')
plt.title('Preis vs Quadratmeter')
plt.xlabel('Quadratmeter')
plt.ylabel('Preis')
plt.show()
```

### Erkenntnisse Grundstückgröße vs Preis

Es lässt sich nicht wirklich eine signifikante Abhängigkeit zwischen Grundstücksgröße und Preis erkennen.

```
In [25]: plt.figure(figsize=(10, 7))
sns.scatterplot(data['Grundstück in qm'], data['Preis'], marker='r')
plt.title('Preis vs Grundstück')
plt.xlabel('Grundstück')
plt.ylabel('Preis')
plt.show()
```

### Erkenntnis - Preis vs Bezirk

Der teuerste Stadtteil ist East End, weil dort das teuerste Haus im Datensatz steht und der Bezirk den höchsten Medianpreis hat. Der billigste Stadtteil ist Paris Island mit dem niedrigsten Median.

```
In [26]: plt.figure(figsize=(10, 10))
plt.title('Preis vs Bezirk')
sns.boxplot(data['Preis'], data['Bezirk'])
plt.xlabel('Preis')
plt.ylabel('Bezirk')
```

Out [26]: Text(12.199600694444456, 0.5, 'Bezirk')

### Erkenntnis - Steigung vs Preis

Hier kann keine eindeutige Auswirkung auf den Preis abgelesen werden. Das zweite Diagramm zeigt außerdem, dass Steigung sehr selten auftritt.

```
In [27]: plt.figure(figsize=(10, 7))
sns.boxplot(data['Steigung'], data['Preis'])
plt.title('Preis vs Steigung')
plt.xlabel('Steigung')
plt.ylabel('Preis')
plt.show()
```

Out [27]: Text(12.199600694444456, 0.5, 'Steigung')

### Erkenntnis - Schlafzimmer vs Preis

Mehr auffällig ist, dass in Datensatz sieben Häuser keine Schlafzimmer haben. Das ist vermutlich ein Fehler. Darüberhinaus kann gesagt werden, dass mehr als vier Schlafzimmer den Preis nicht steigern. Am besten sind Häuser mit ein bis vier Schlafzimmer.

Die Unterschiede sind klein.

```
In [29]: plt.figure(figsize=(10, 7))
sns.boxplot(data['Schlafzimmer'], data['Preis'])
plt.title('Preis vs Schlafzimmer')
plt.xlabel('Schlafzimmer')
plt.ylabel('Preis')
plt.show()
```

Out [29]: Text(12.199600694444456, 0.5, 'Schlafzimmer')

### Erkenntnis - Preis vs Räume

Mehr Räume ziehen einen höheren Preis nach sich. Auch wenn angemerkt werden muss, dass ab 10 Räumen der Preis sehr stark variiert und dadurch keine konkrete Korrelation mehr gegeben ist.

```
In [30]: plt.figure(figsize=(10, 7))
sns.boxplot(data['Räume'], data['Preis'])
plt.title('Preis vs Räume')
plt.xlabel('Räume')
plt.ylabel('Preis')
plt.show()
```

Out [30]: Text(0.5, 1.0, 'Preis vs Räume')

### Erkenntnis - Preis vs Lage

Wenn ein Haus an der Straße oder an der Bahn liegt, ist der Preis unterdurchschnittlich. Wenn ein Sportzentrum oder eine Grünanlage in der Nähe liegt, ist der Preis eher überdurchschnittlich. Allerdings können auch teurere Häuser in einer normalen Lage liegen.

```
In [31]: plt.figure(figsize=(10, 7))
sns.boxplot(data['Lage'], data['Preis'])
plt.title('Preis vs Lage')
plt.xlabel('Lage')
plt.ylabel('Preis')
plt.show()
```

Out [31]: Text(12.199600694444456, 0.5, 'Preis vs Lage')

### Erkenntnis - Klimaanlage vs Preis

Eine vorhandene Klimaanlage steigert den Preis deutlich.

```
In [32]: plt.figure(figsize=(10, 7))
plt.title('Preis vs Klimaanlage')
sns.boxplot(data['Preis'], data['Klimaanlage'])
plt.xlabel('Preis')
plt.ylabel('Klimaanlage')
```

Out [32]: Text(12.199600694444456, 0.5, 'Preis vs Klimaanlage')

### Erkenntnis - Verkaufsmotivat vs Preis

Man kann keinen nachweisbaren Einfluss des Verkaufsmotivats auf den Preis ablesen.

```
In [33]: plt.figure(figsize=(10, 7))
plt.title('Preis vs Verkaufsmotivat')
sns.boxplot(data['Preis'], data['Verkaufsmotivat'])
plt.xlabel('Preis')
plt.ylabel('Verkaufsmotivat')
```

Out [33]: Text(12.199600694444456, 0.5, 'Preis vs Verkaufsmotivat')

### Erkenntnis - Verkaufsjahr vs Preis

Man kann keinen nachweisbaren Einfluss des Verkaufsjahrs auf den Preis ablesen.

```
In [34]: plt.figure(figsize=(10, 7))
plt.title('Preis vs Verkaufsjahr')
sns.boxplot(data['Preis'], data['Verkaufsjahr'])
plt.xlabel('Preis')
plt.ylabel('Verkaufsjahr')
```

Out [34]: Text(12.199600694444456, 0.5, 'Preis vs Verkaufsjahr')

### Erkenntnis - Garagentyp vs Preis

Der Unterschied zwischen keiner Garage und CarPort ist sehr gering. Eine eingebaute Garage hebt den Hauspreis deutlich. Außerdem lohnt es sich eine Garage noch anzubauen, da damit der Preis steigt.

```
In [35]: plt.figure(figsize=(10, 7))
sns.boxplot(data['Garage Typ'], data['Preis'])
plt.title('Preis vs Garagentyp')
plt.xlabel('Preis')
plt.ylabel('Garage Typ')
plt.show()
```

Out [35]: Text(12.199600694444456, 0.5, 'Preis vs Garagentyp')

### Erkenntnis Zustand vs Preis

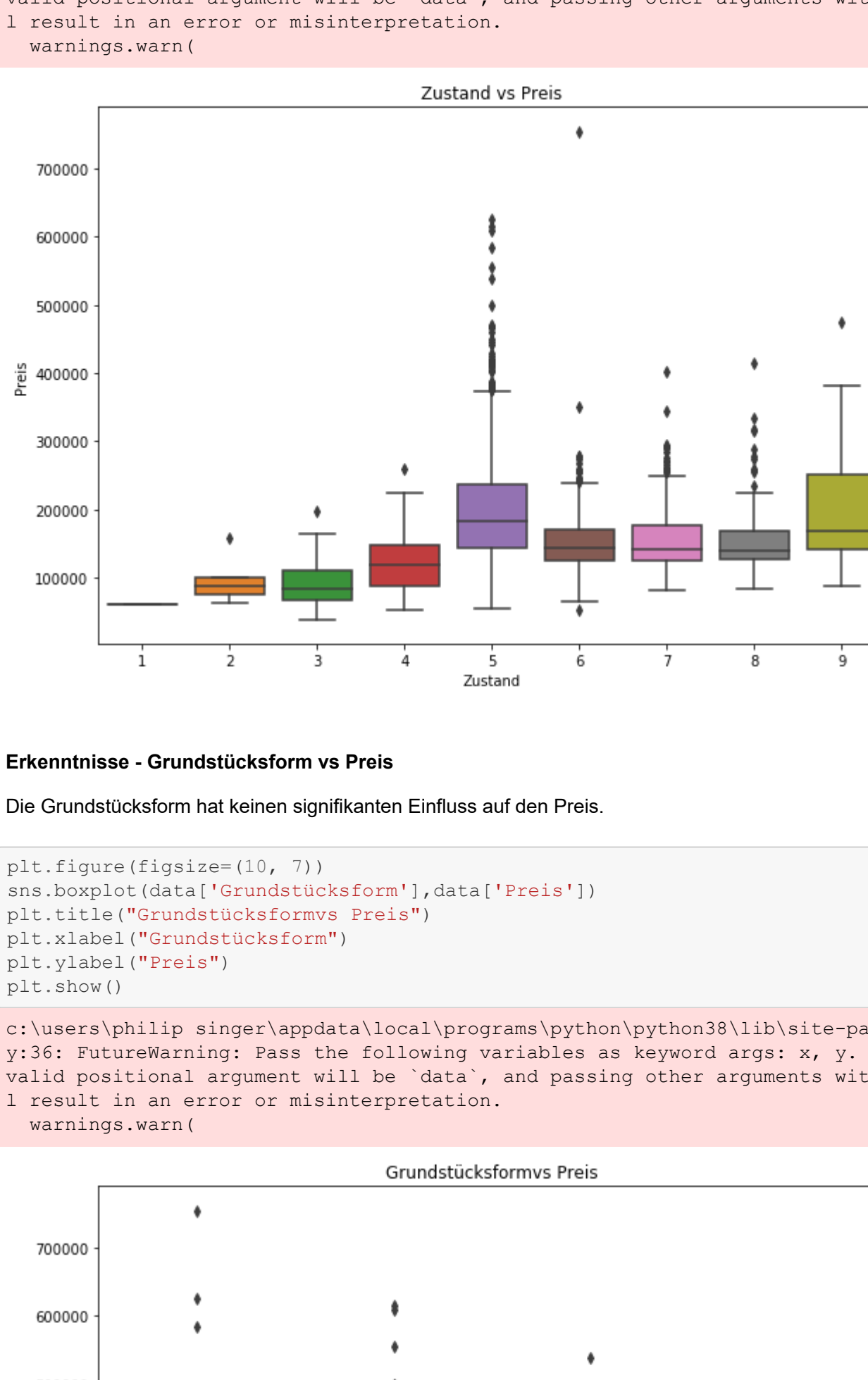
Ein durchschnittlicher Zustand zieht den höchsten Preis nach sich. Logisch betrachtet, hat der Zustand in diesem Datensatz keinen Einfluss auf den Preis.

Man kann jedoch festhalten, dass ein Haus mit durchschnittlichem bzw. überdurchschnittlichem Zustand einen höheren Preis hat, als unterdurchschnittliche Häuser.



```
[36]: plt.figure(figsize=(10, 7))
sns.boxplot(data['zustand'], data['Preis'])
plt.title("Zustand vs Preis")
plt.xlabel("Zustand")
plt.ylabel("Preis")
plt.show()

c:\Users\philip.singer\appdata\local\programs\python\python38\lib\site-packages\seaborn\_decorators.py:36: FutureWarning: Pass the following variables as keyword args: x, y. From version 0.12, the only
valid positional argument will be 'data', and passing other arguments without an explicit keyword will
result in an error or misinterpretation.
warnings.warn(
```

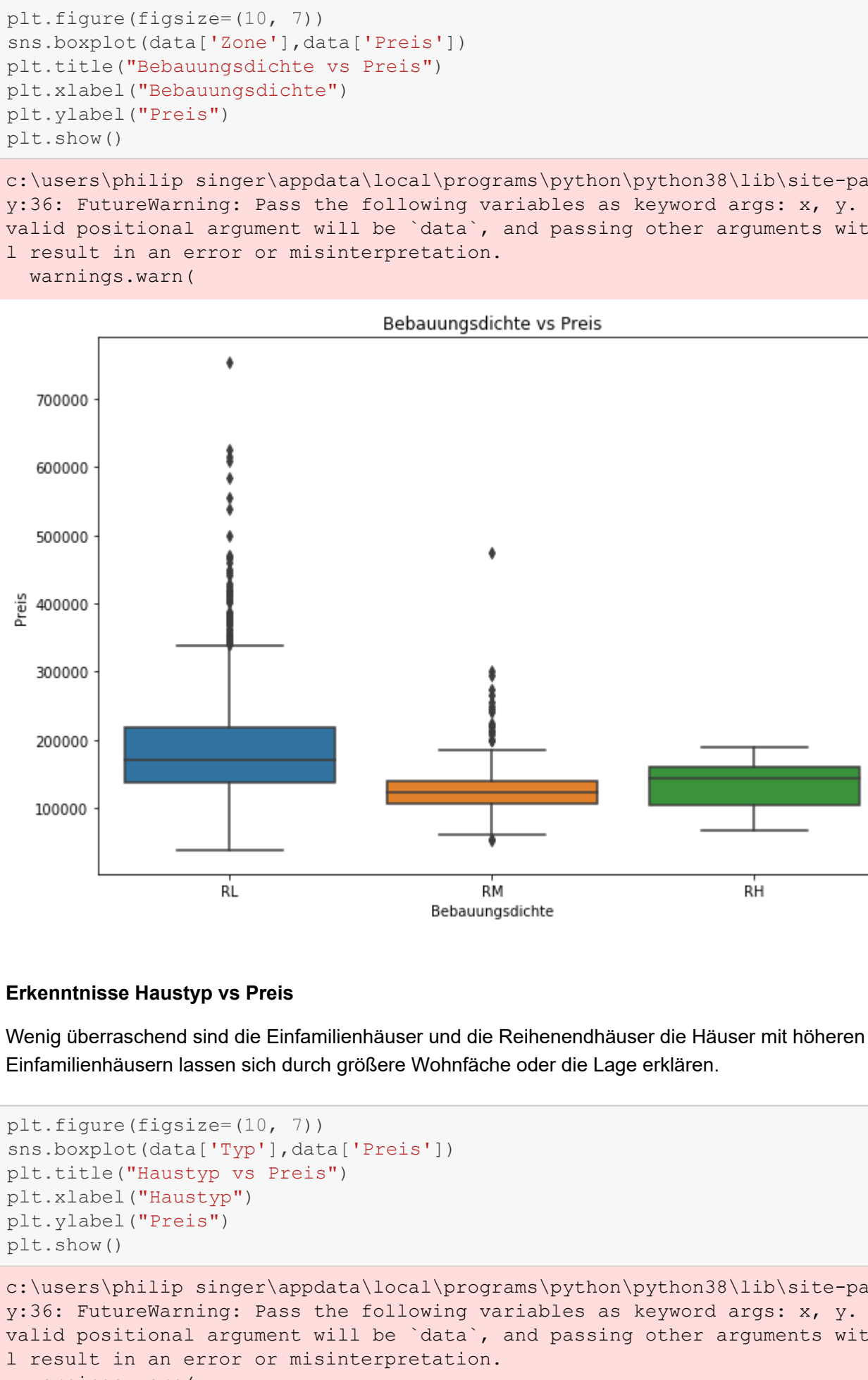


### Erkenntnisse - Grundstücksform vs Preis

Die Grundstücksform hat keinen signifikanten Einfluss auf den Preis.

```
In [37]: plt.figure(figsize=(10, 7))
sns.boxplot(data['Grundstücksform'], data['Preis'])
plt.title("Grundstücksform vs Preis")
plt.xlabel("Grundstücksform")
plt.ylabel("Preis")
plt.show()

c:\Users\philip.singer\appdata\local\programs\python\python38\lib\site-packages\seaborn\_decorators.py:36: FutureWarning: Pass the following variables as keyword args: x, y. From version 0.12, the only
valid positional argument will be 'data', and passing other arguments without an explicit keyword will
result in an error or misinterpretation.
warnings.warn(
```

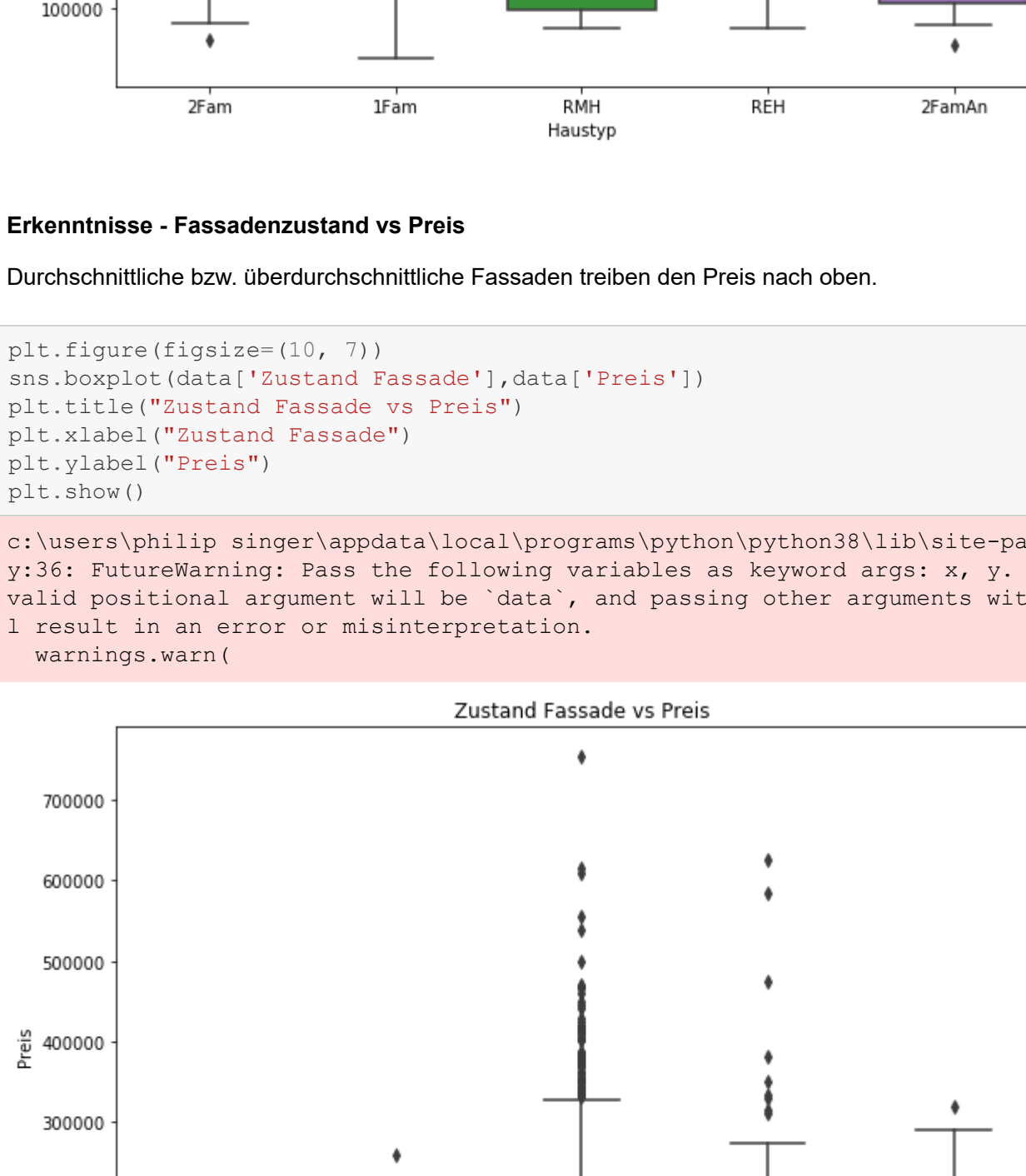


### Erkenntnisse - Bebauungsdichte vs Preis

Eine niedrige Bebauungsdichte zieht einen deutlich höheren Preis nach sich als höhere Bebauungsdichten.

```
In [38]: plt.figure(figsize=(10, 7))
sns.boxplot(data['zone'], data['Preis'])
plt.title("Bebauungsdichte vs Preis")
plt.xlabel("Bebauungsdichte")
plt.ylabel("Preis")
plt.show()

c:\Users\philip.singer\appdata\local\programs\python\python38\lib\site-packages\seaborn\_decorators.py:36: FutureWarning: Pass the following variables as keyword args: x, y. From version 0.12, the only
valid positional argument will be 'data', and passing other arguments without an explicit keyword will
result in an error or misinterpretation.
warnings.warn(
```

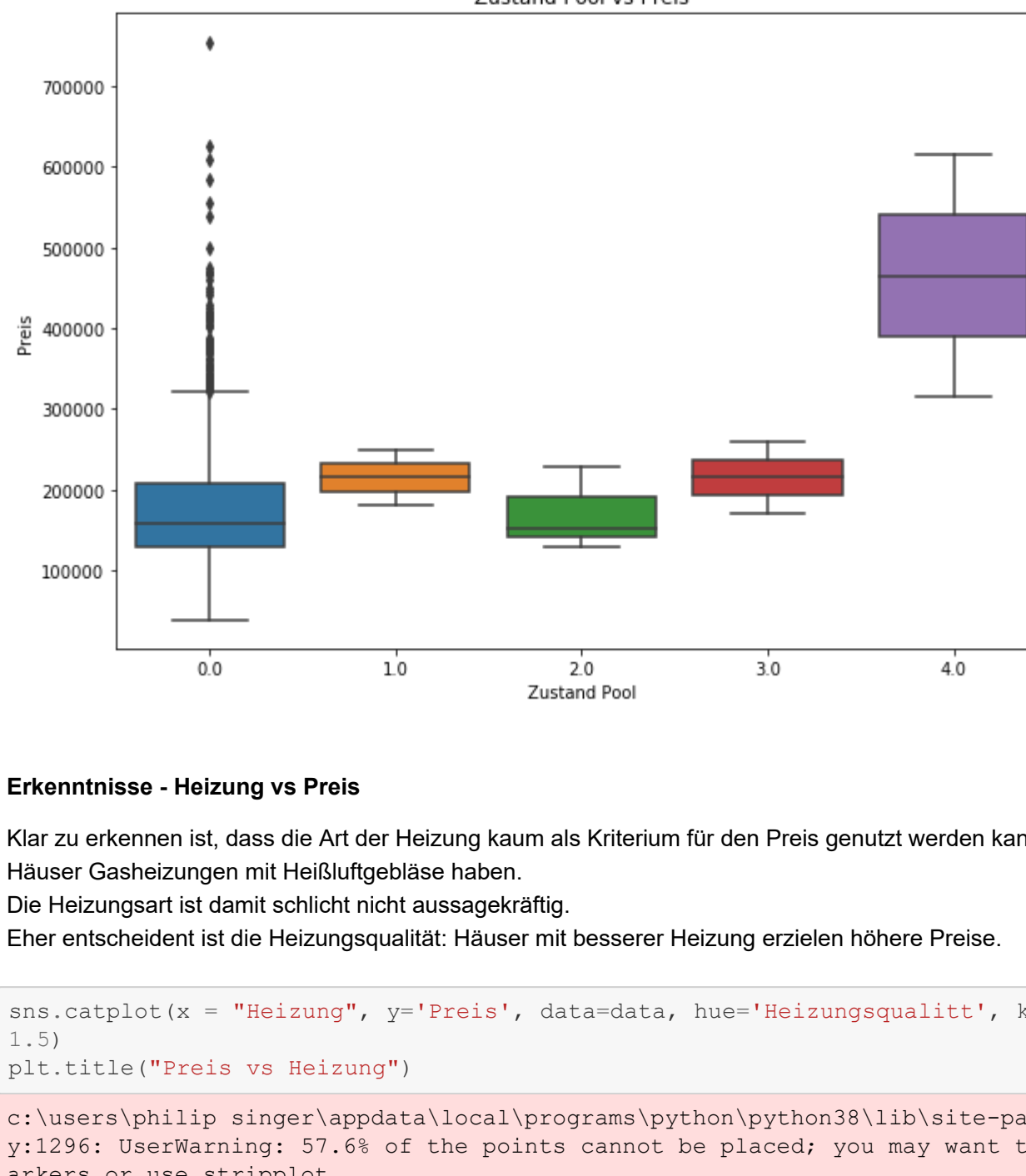


### Erkenntnisse Haustyp vs Preis

Wenig überraschend sind die Einfamilienhäuser und die Reihenhäuser die Häuser mit höheren Preisen. Die Ausreißer bei den Einfamilienhäusern lassen sich durch größere Wohnfläche oder die Lage erklären.

```
In [39]: plt.figure(figsize=(10, 7))
sns.boxplot(data['Typ'], data['Preis'])
plt.title("Haustyp vs Preis")
plt.xlabel("Haustyp")
plt.ylabel("Preis")
plt.show()

c:\Users\philip.singer\appdata\local\programs\python\python38\lib\site-packages\seaborn\_decorators.py:36: FutureWarning: Pass the following variables as keyword args: x, y. From version 0.12, the only
valid positional argument will be 'data', and passing other arguments without an explicit keyword will
result in an error or misinterpretation.
warnings.warn(
```

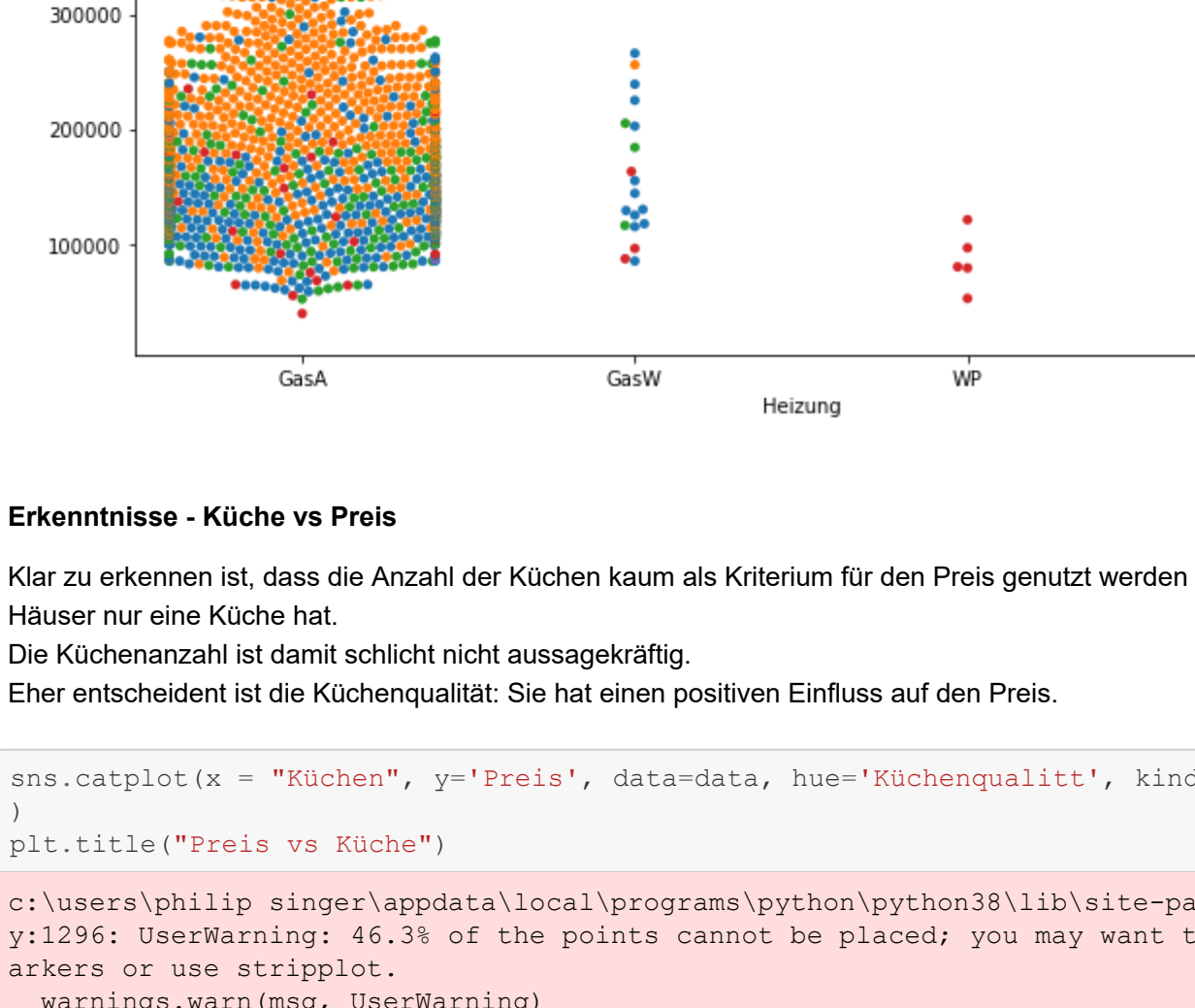


### Erkenntnisse - Fassadenzustand vs Preis

Durchschnittliche bzw. überdurchschnittliche Fassaden treiben den Preis nach oben.

```
In [40]: plt.figure(figsize=(10, 7))
sns.boxplot(data['zustand Fassade'], data['Preis'])
plt.title("Zustand Fassade vs Preis")
plt.xlabel("Zustand Fassade")
plt.ylabel("Preis")
plt.show()

c:\Users\philip.singer\appdata\local\programs\python\python38\lib\site-packages\seaborn\_decorators.py:36: FutureWarning: Pass the following variables as keyword args: x, y. From version 0.12, the only
valid positional argument will be 'data', and passing other arguments without an explicit keyword will
result in an error or misinterpretation.
warnings.warn(
```

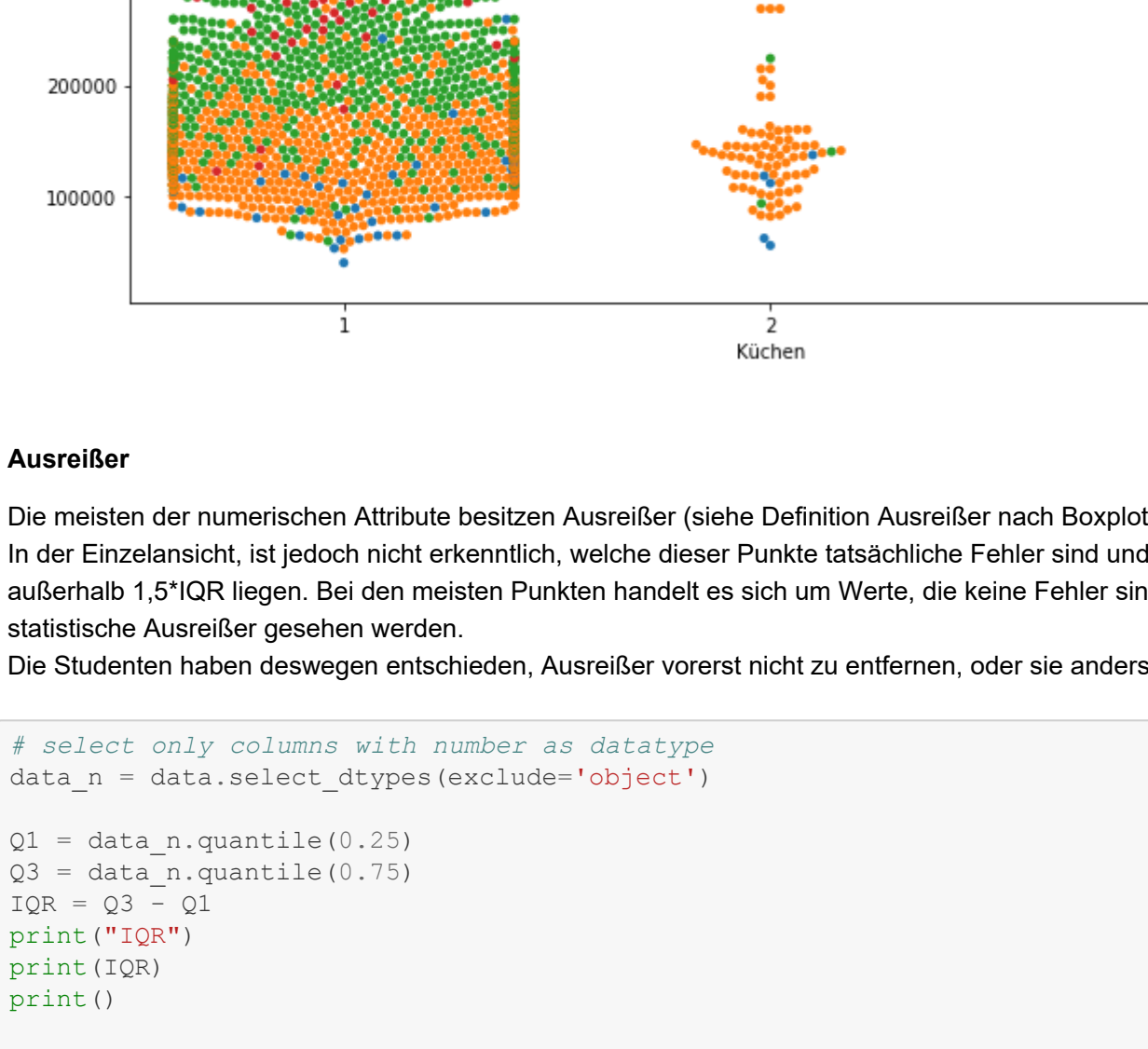


### Erkenntnisse - Pool Zustand vs Preis

Eine positive Auswirkung auf den Preis hat ein Pool nur, wenn dessen Zustand sehr gut ist.

```
In [41]: plt.figure(figsize=(10, 7))
sns.boxplot(data['Pool'], data['Preis'])
plt.title("Zustand Pool vs Preis")
plt.xlabel("Zustand Pool")
plt.ylabel("Preis")
plt.show()

c:\Users\philip.singer\appdata\local\programs\python\python38\lib\site-packages\seaborn\_decorators.py:36: FutureWarning: Pass the following variables as keyword args: x, y. From version 0.12, the only
valid positional argument will be 'data', and passing other arguments without an explicit keyword will
result in an error or misinterpretation.
warnings.warn(
```

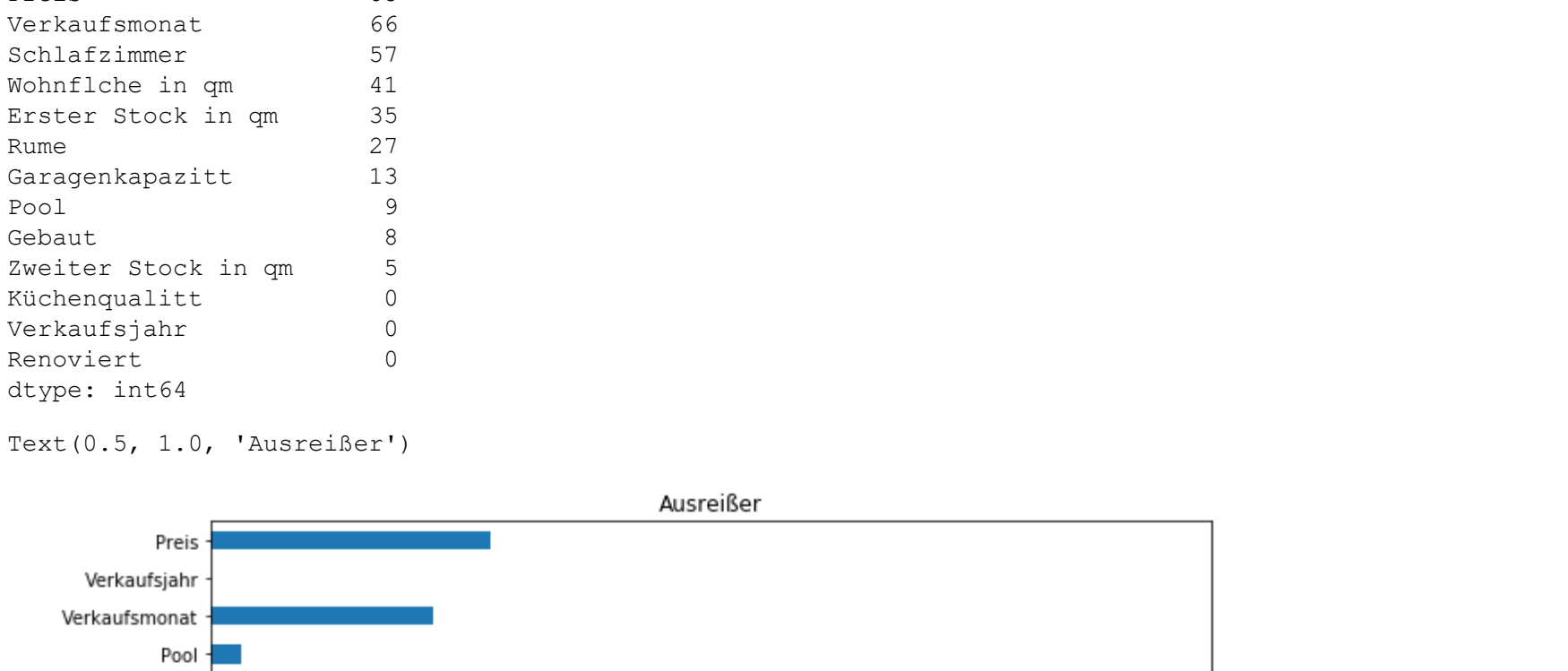


### Erkenntnisse - Heizung vs Preis

Klar zu erkennen ist, dass die Art der Heizung kaum als Kriterium für den Preis genutzt werden kann, weil eine deutliche Mehrheit der Häuser Gasheizungen mit Heißluftbalken haben. Die Heizungsart ist damit schlicht nicht aussagekräftig. Eher entscheidet ist die Heizungsqualität: Häuser mit besserer Heizung erzielen höhere Preise.

```
In [42]: sns.catplot(x = "Heizung", y = "Preis", data = data, hue = "Heizungsqualität", kind = "swarm", height = 7, aspect = 1.5)
plt.title("Preis vs Heizung")

c:\Users\philip.singer\appdata\local\programs\python\python38\lib\site-packages\seaborn\categorical.py:1296: UserWarning: 57.6% of the points cannot be placed; you may want to decrease the size of the markers or use stripplot.
warnings.warn(msg, UserWarning)
```

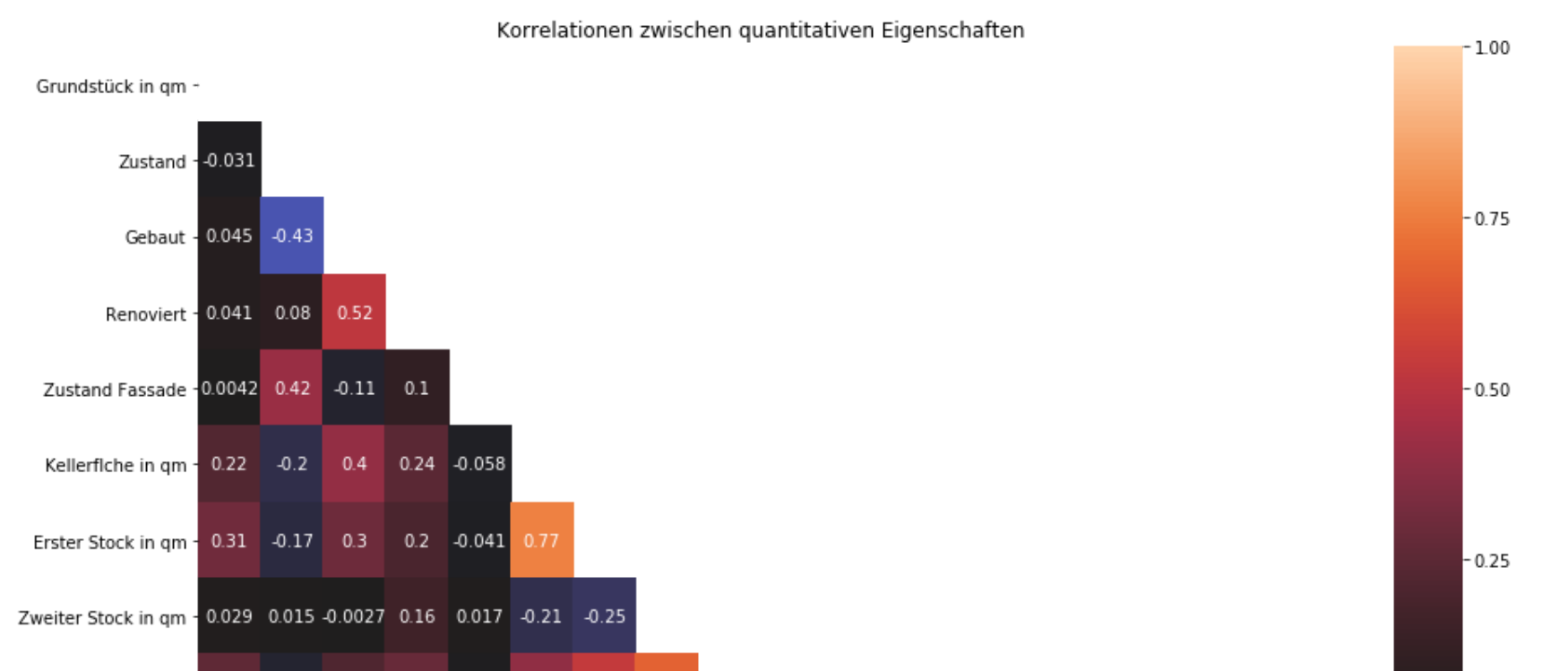


### Erkenntnisse - Küche vs Preis

Klar zu erkennen ist, dass die Anzahl der Küchen kaum als Kriterium für den Preis genutzt werden kann, weil eine deutliche Mehrheit der Häuser nur eine Küche hat. Die Küchenanzahl ist damit schlicht nicht aussagekräftig. Eher entscheidet ist die Küchenqualität: Häuser mit besserer Küche erzielen höhere Preise.

```
In [43]: sns.catplot(x = "Küche", y = "Preis", data = data, hue = "Küchenqualität", kind = "swarm", height = 7, aspect = 1.5)
plt.title("Preis vs Küche")

c:\Users\philip.singer\appdata\local\programs\python\python38\lib\site-packages\seaborn\categorical.py:1296: UserWarning: 46.3% of the points cannot be placed; you may want to decrease the size of the markers or use stripplot.
warnings.warn(msg, UserWarning)
```



### Ausreißer

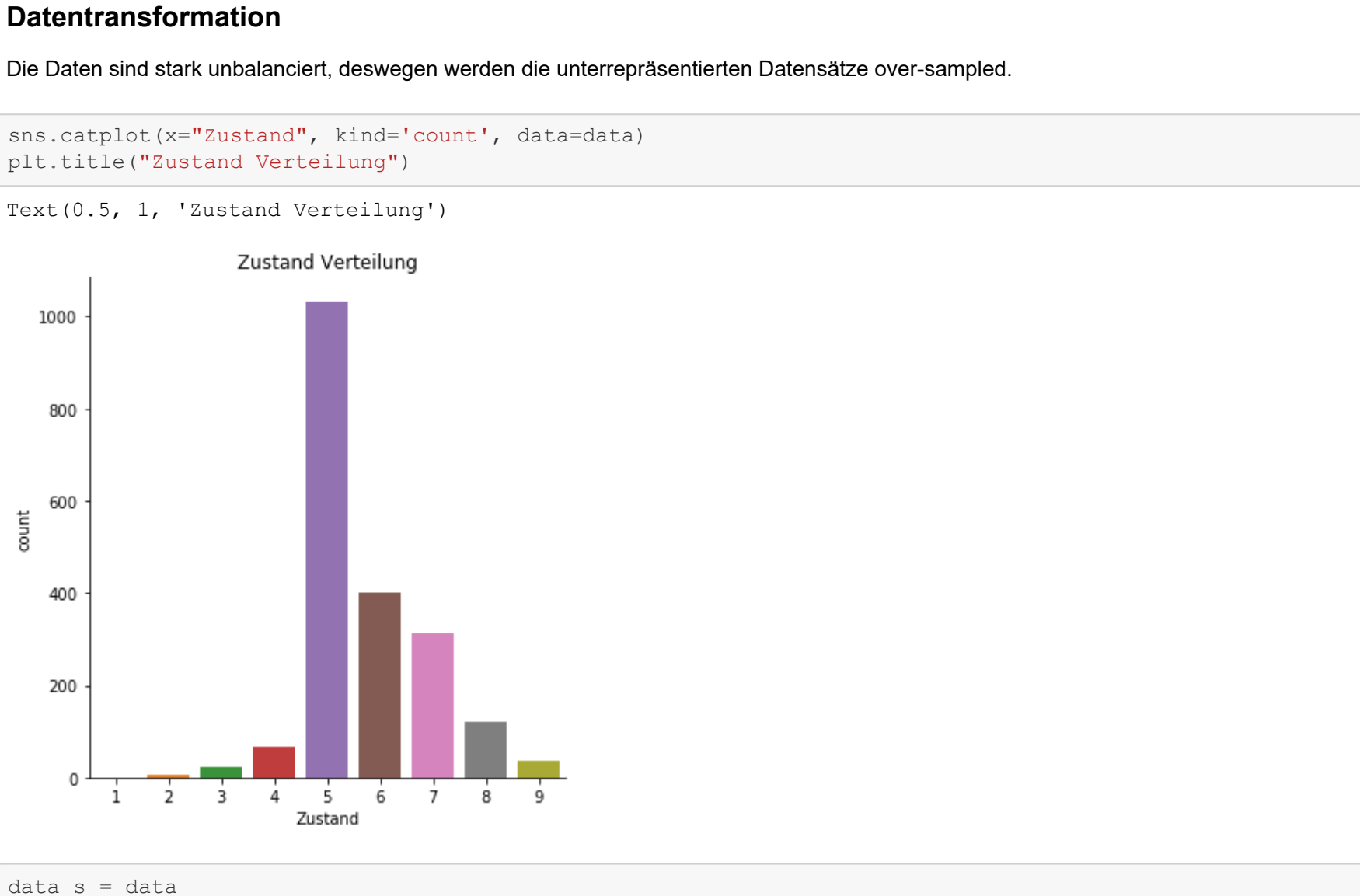
Die meisten der numerischen Attribute besitzen Ausreißer (siehe Definition Ausreißer nach Boxplot-Diagramm). In der Einzelansicht, ist jedoch nicht erkennlich, welche dieser Punkte tatsächliche Fehler sind und welche nur logisch nachvollziehbar außerhalb 1.5\*IQR liegen. Bei den meisten Punkten handelt es sich um Werte, die keine Fehler sind, und nur auf Basis der Verteilung als statistische Ausreißer gesehen werden. Die Studenten haben deswegen entschieden, Ausreißer vorerst nicht zu entfernen, oder sie anders zu behandeln.

```
In [44]: # select only columns with number as datatype
data_n = data.select_dtypes(exclude="object")

Q1 = data_n.quantile(0.25)
Q3 = data_n.quantile(0.75)
IQR = Q3 - Q1
print ("IQR")
print (IQR)

print ("====Ausreißer==== pro Spalte====")
ausreißer = ((data_n < (Q1 - 1.5 * IQR)) | (data_n > (Q3 + 1.5 * IQR))).sum()
display(ausreißer.sort_values(ascending=False))

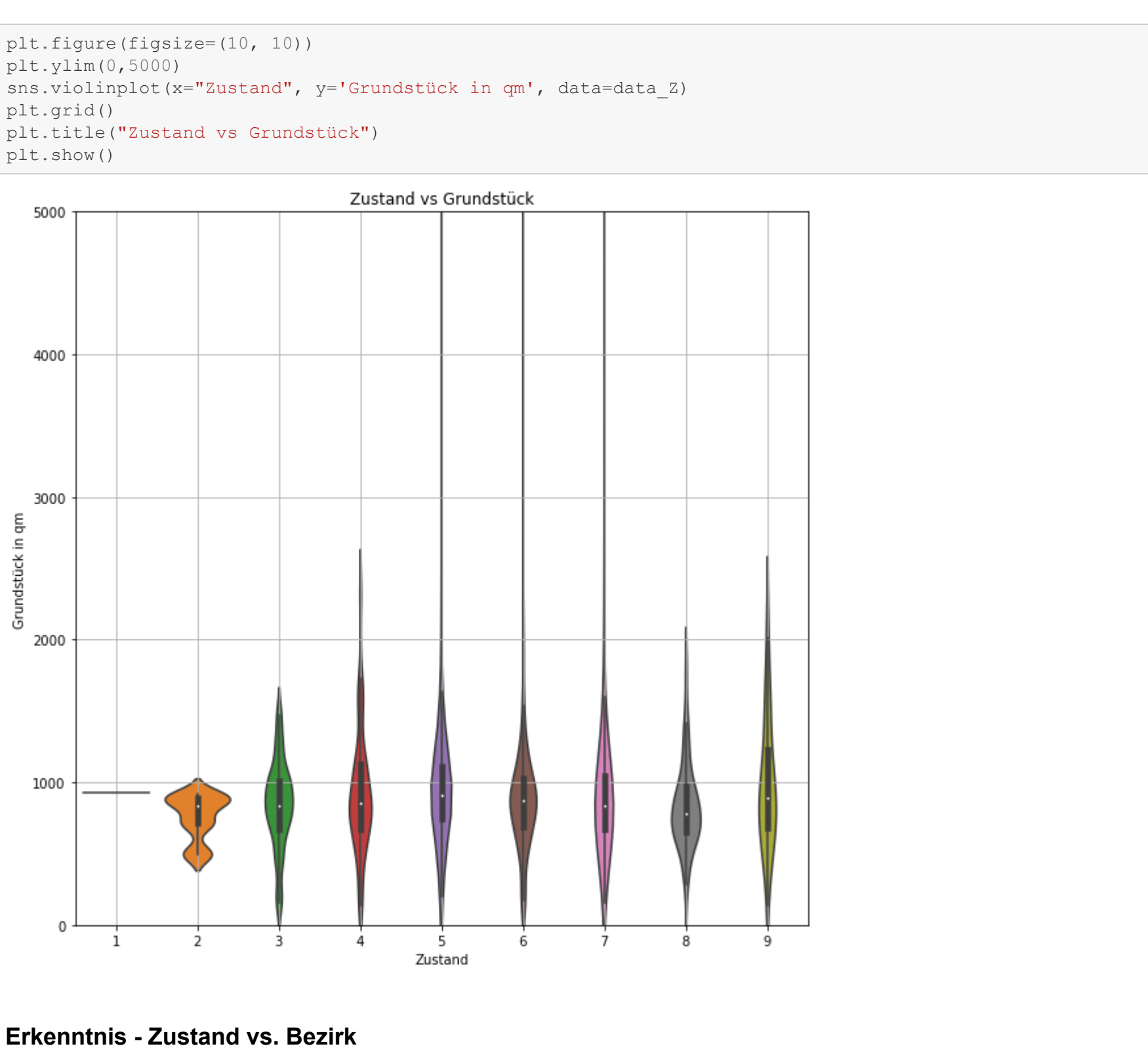
# Plot
plt.figure(figsize=(10, 7))
ausreißer.plot(kind="bar")
plt.title("Ausreißer")
```



### Erkenntnis - Korrelationen

Am stärksten positiv korrelieren die Attribute Räume und Wohnfläche miteinander. Darüber hinaus ist zu sehen, dass auch Wohnfläche und Preis korrelieren. Versunderlich ist, dass das Baujahr negativ mit dem Zustand korreliert, was bedeutet, dass neuere Häuser einen tendenziell schlechteren Zustand haben.

```
In [45]: plt.figure(figsize=(15, 15))
plt.title("Korrelationen zwischen quantitativen Eigenschaften")
corr = data.select_dtypes(exclude="object").corr()
mask = np.zeros_like(corr)
mask[np.triu_indices_from(mask)] = True
sns.heatmap(corr, vmin = -1, vmax = 1, center = 0, mask=mask, annot=True)
```



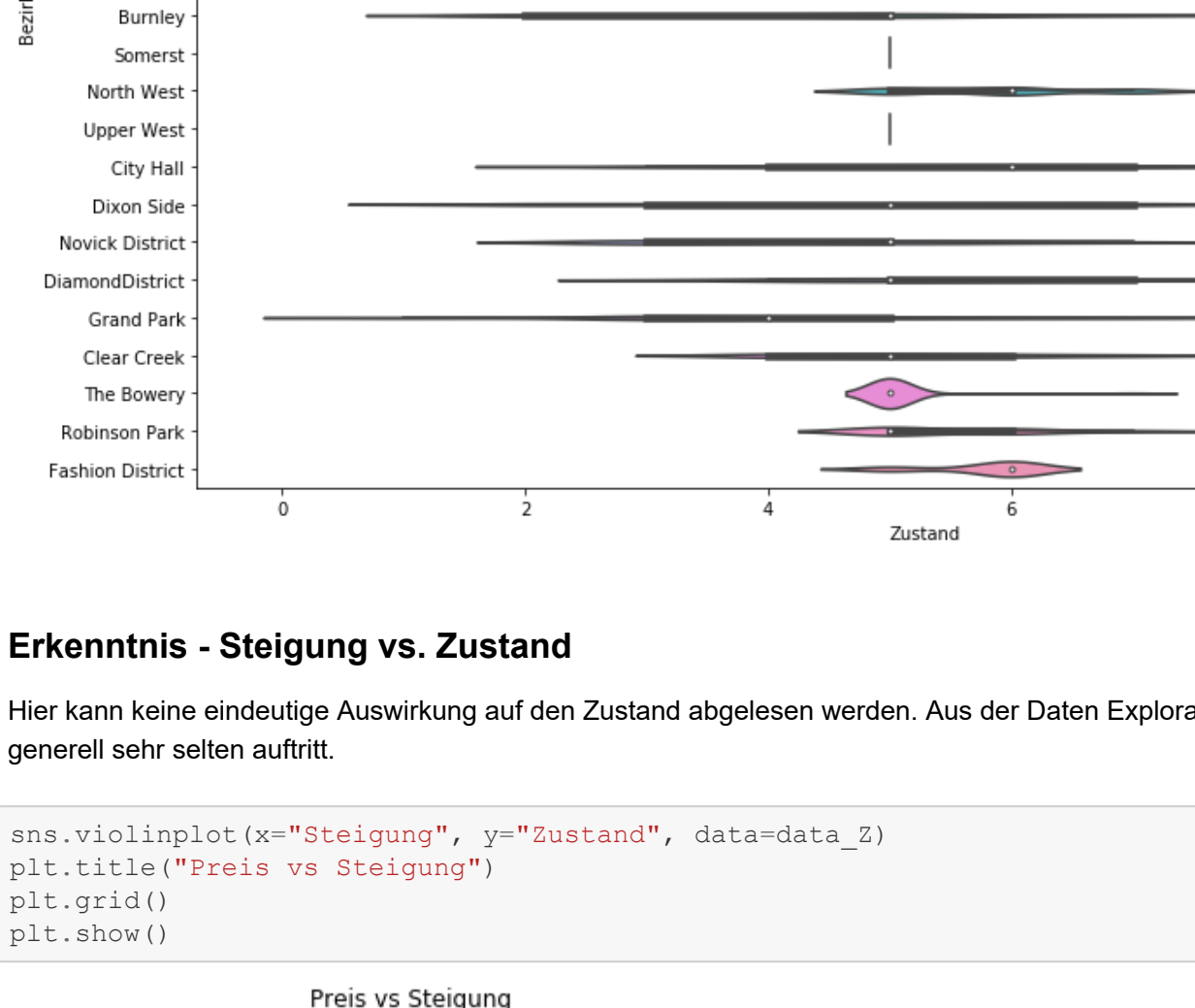
### Diagramme für die Zustandsvorhersage

#### Datentransformation

Die Daten sind stark unbalanciert, deswegen werden die unterrepräsentierten Datensätze over-sampled.

```
In [46]: sns.catplot(x="zustand", kind="count", data=data)
plt.title("Zustand Verteilung")

Out [46]: Text(0.5, 1, 'Zustand Verteilung')
```



```
In [47]: data_s = data
ros = RandomOverSampler(random_state=123, sampling_strategy={1:30, 2:150, 3:350, 4:400}) # 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 211)
X, y = ros.fit_resample(data_s.loc[:, data_s.columns != 'zustand'], data_s['zustand'])

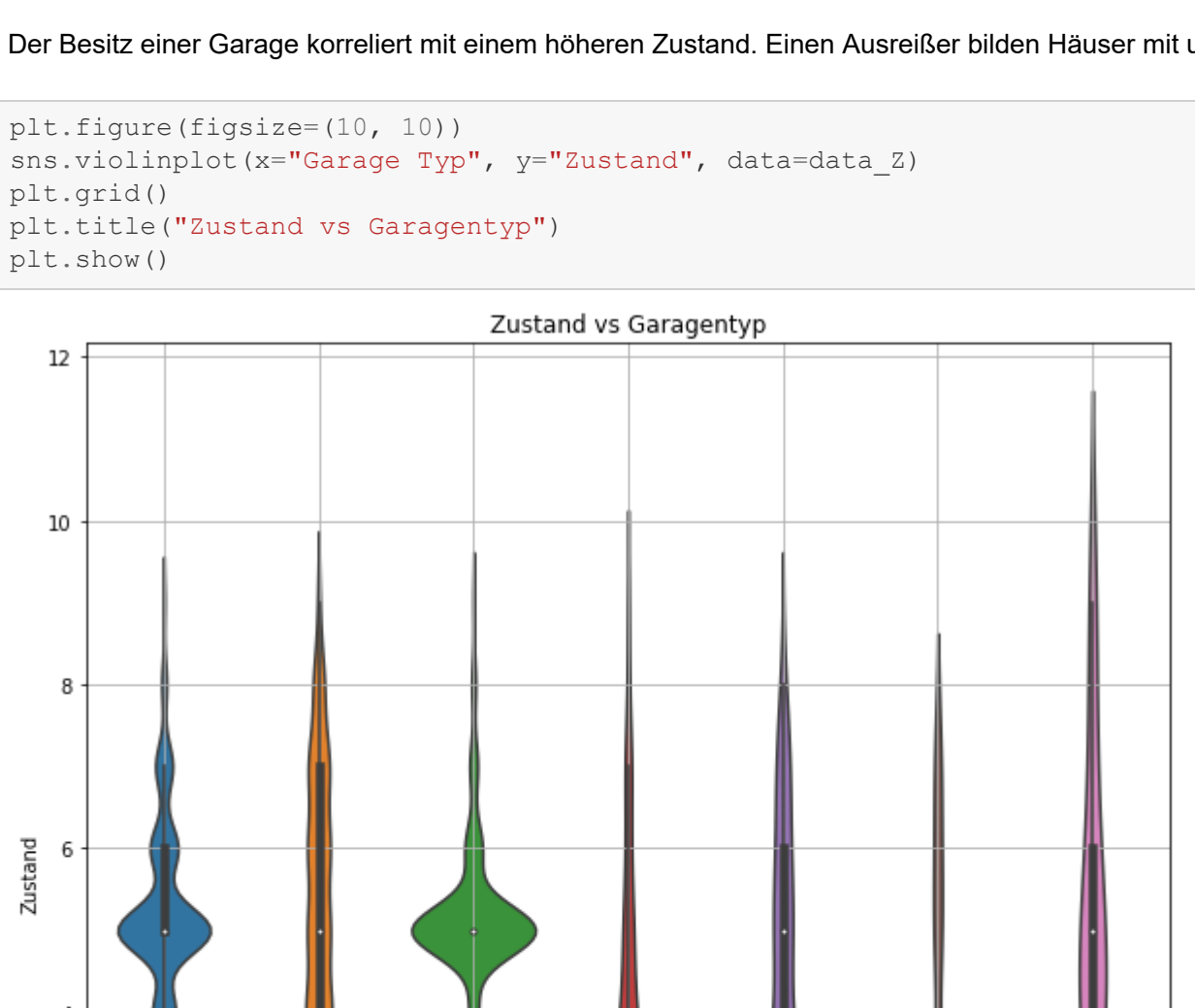
In [48]: data_s = pd.concat([X, y], axis=1)
sns.catplot(x="zustand", kind="count", data=data_s)
plt.title("Zustand Verteilung")
data_s = data_s
```



### Erkenntnisse Grundstückgröße vs Zustand

Es lässt sich nicht eine signifikante Abhängigkeit zwischen Grundstücksgröße und Preis erkennen.

```
In [49]: plt.figure(figsize=(10, 10))
sns.violinplot(x="zustand", y="Grundstück in qm", data=data_2)
plt.grid()
plt.title("Zustand vs Grundstück")
plt.show()
```



### Erkenntnis - Zustand vs. Bezirk

Der Bezirk 'University' hat nach dem Median die Häuser mit den schlechtesten Zustand.

```
In [50]: plt.figure(figsize=(15, 10))
plt.title("Zustand vs Bezirk")
sns.violinplot(x="zustand", y="Bezirk", data=data_2)

Out [50]: <matplotlib.axes._subplots.AxesSubplot at 0x1fb27d25160>
```



### Erkenntnis - Steigung vs. Zustand

Hier kann keine eindeutige Auswirkung auf den Zustand abgelesen werden. Aus der Daten Exploration zum Preis wissen wir, dass Steigung generell sehr selten auftritt.

```
In [51]: sns.violinplot(x="Steigung", y="zustand", data=data_2)
plt.title("Preis vs Steigung")
plt.grid()
plt.show()
```



### Erkenntnis - Klimaanlage vs Zustand

Eine vorhandene Klimaanlage korreliert mit einem höheren Zustand.

```
In [52]: plt.figure(figsize=(10, 10))
plt.title("Zustand vs Klimaanlage")
plt.xlabel("Klimaanlage")
sns.violinplot(x="Klimaanlage", y="zustand", data=data_2)
plt.grid()
```



### Erkenntnis - Garagentyp vs Preis

Der Besitz einer Garage korreliert mit einem höheren Zustand. Einen Ausreißer bilden Häuser mit unterschiedlichen Garagen.

```
In [53]: plt.figure(figsize=(10, 10))
sns.violinplot(x="Garage Typ", y="zustand", data=data_2)
plt.grid()
plt.title("Zustand vs Garagentyp")
plt.show()
```



### Erkenntnisse - Grundstücksform vs Zustand

Eine stark reguläre Grundstücksform korreliert mit einem unterdurchschnittlichen Zustand.







