

label encoder
used for colp of the data frame

Education
HS - 0
UG - 1
PG - 2
UG - 0
UG - 1
PG - 2

HS - 0
UG - 1
PG - 2

ORDINAL ENCODING

* Before encoding split the data into training & testing

Review	education	Purchased
Poor	School	Yes
Avg	UG	No
Good	PG	Yes

code

```
from sklearn.preprocessing import OrdinalEncoder
oe = OrdinalEncoder(categories=[['Poor', 'Avg', 'Good'],
                                 [0, 1, 2],
                                 ['School', 'UG', 'PG']])
```

oe.fit(X-train)

X-train = oe.transform(X-train)

X-test = oe.transform(X-test)

for O/p column

from sklearn.preprocessing import LabelEncoder
 le = LabelEncoder()
 le.fit(y-train)

- Label encoding is used to encode target column only.

One Hot Encoding

- used for nominal data (which don't have any order)

	Yellow	Blue	Red	
Y - Yellow	1	0	0	Encoded data
B - Blue	0	1	0	
R - Red	0	0	1	

Dummy Variable trap → multicollinearity

owner	fuel	selling-price
1	1	

One using pandas

Pd. get_dummies cat, columns = ['fuel', 'owner'])
 (not preferred)

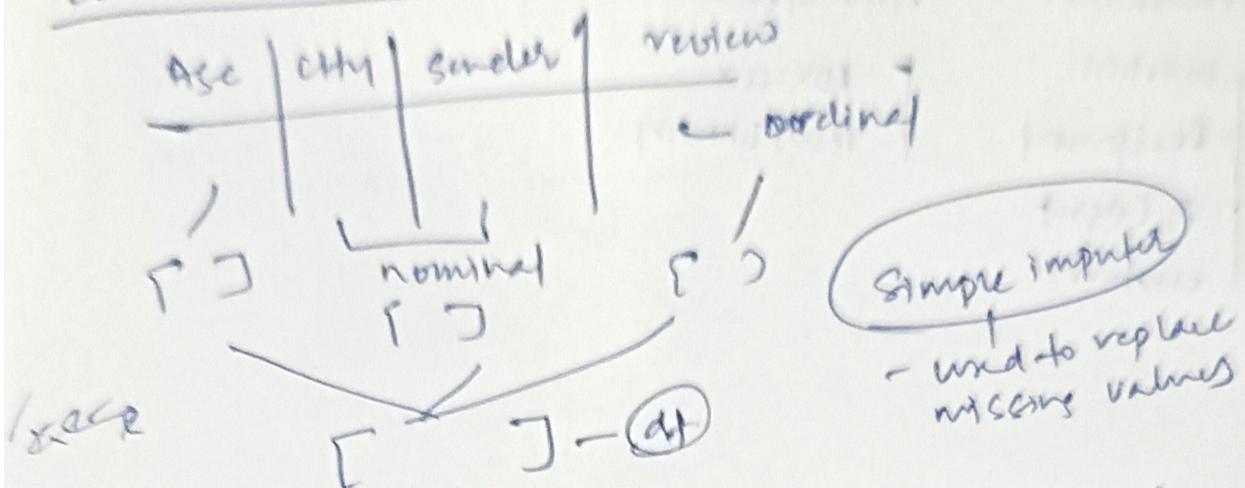
One using sklearn

- split the data

from sklearn.preprocessing import OneHotEncoder
 ohe = OneHotEncoder()
 ohe.fit_transform(x-train[['fuel', 'owner']]).toarray()
 ohe.transform(x-test[['fuel', 'owner']]).toarray()
 Append them with the original df

np.where(111) creates two arrays side by side
sex/cat + gender = Sex/cat Gender

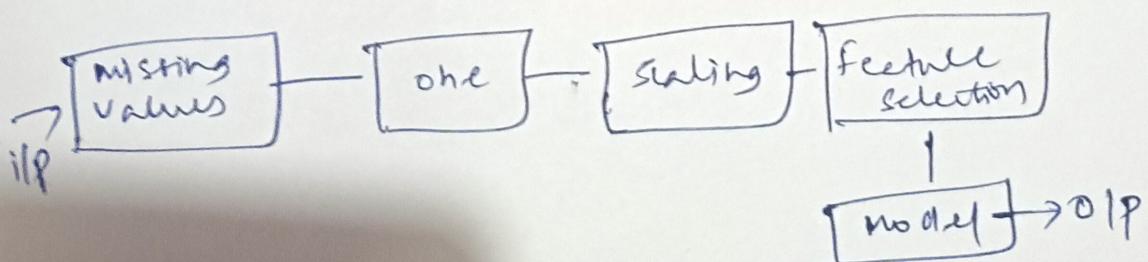
columnTransformer

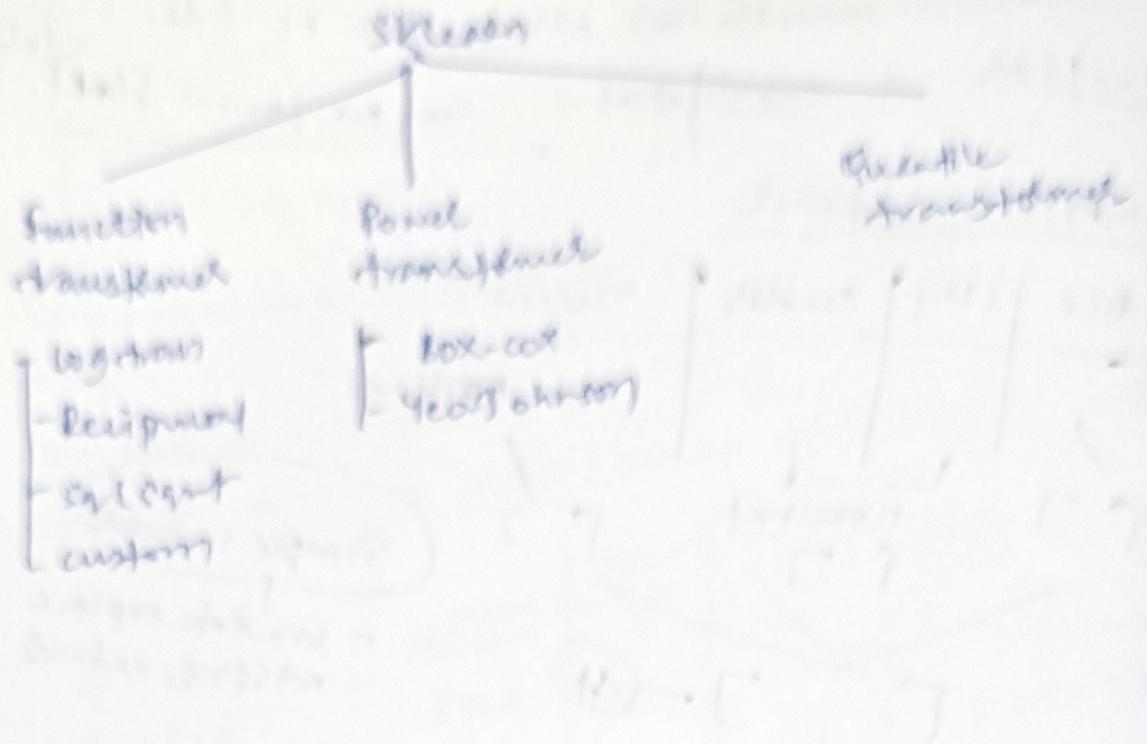


from sklearn.compose import ColumnTransformer
transformer = ColumnTransformer(transformers=[
('dt1', SimpleImputer(), ['level']),
('dt2', OrdinalEncoder(categories=[['mild'],
['strong'], 'rough'])),
('dt3', OneHotEncoder(sparse=False, drop='first'),
['gender', 'city']),
], remainder='pass through')

Pipelines

- pipelines chains together multiple steps so that the output of each step is used as input to the next step
- pipelines makes it easy to apply the same preprocessing to train & test.





Erstellt am 10.09.2019 von 12:49:00 mit 142 Seiten
 337 Abbildungen, 240 Formeln und 10 Tabellen
 von 1000 (100% abgeschlossen, 0% fehlerhaft)
 Erstellt am 10.09.2019 von 12:49:00 mit 142 Seiten
 abgelegt am 10.09.2019
 abgelegt am 10.09.2019 mit 142 Seiten
 abgelegt am 10.09.2019 mit 142 Seiten
 abgelegt am 10.09.2019 mit 142 Seiten

Wert 217
 entnommen am 29.08.2019 um 10:00 Uhr von 1000 Seiten
 Seite 217 ab 911 00 bis zu 21 911 00 der 10. Kapitel
 Wörterbuch Seite 21 911 00 bis zu 21 911 00 der 10. Kapitel
 Seite 217 ab 911 00 bis zu 21 911 00 der 10. Kapitel
 Seite 217 ab 911 00 bis zu 21 911 00 der 10. Kapitel

