

Feature Engineering

2 - Introduction to feature engineering

- What is feature engineering?

It is process of using domain knowledge to extract features from raw data via data mining techniques.

Extracting Information:-

- More complex examples:
 - Text data for deed of house
 - length of text
 - Number of times certain terms are mentioned -

Combining Information:-

- Could also combine extracted information
- New features:
 - 0 or 1 value indicating:
 - Both weekends and evenings?

Transforming Information:-

- Often categorical data is present as string data-

- For example a large data set of social network users could have country of origin as a string feature (e.g. UK, USA, MEX, etc)

Two way / approaches of Transforming data:-

- Integers Encoding
- Onehot Encoding (Dummy Variables)

Integers Encoding :-

• Directly convert categories into integers 1, 2, 3, ..., N -

Country	Country	possible issue
USA	1	- apply order
MEX	2	to the column
CAN	3	(ordinal variables)

PROS :-

- very easy to do
(and understand)
- Does not increase number of features

CONS

- Implies ordered relationship between categories

Onehot Encoding :-

- Convert each category into individual features that are either 0 or 1

	USA	MEX	CAN
USA	1	0	0
MEX	0	1	0
CAN	0	0	1
USA	1	0	1

- issue greatly expand features

Dummy Variable Trap!

mathematically known as
"Multi-collinearity"

- Converting dummy variables can cause 1 feature to be duplicated.
- Let consider simplest possible example:-

Vertical Street	UP	DOWN	?
UP	1	0	1
DOWN	0	1	0
UP	1	0	1
DOWN	0	0	0

Pros:-

- NO ordering implied

Cons:-

- Potential to create many more features (column and coefficients)
- Dummy variable trap consideration.
- Not easy to add new categories.

Throughout this course:-

- Outliers in data
- Missing Data
- Categorical Data
- Not every issue here is feature engineering but also can be called data cleaning -

3. Dealing with Outliers

Outliers:-

Extreme data points which might disturb our data statistically, which might seem in rare-

It's often better to simply remove these few points from the data set in order to have more generalized model.

Outlier Consideration:-

- Definition of an Outlier
 - Range and limit
 - percentage of Data
- There are both very domain dependent!

Range and limit:-

We need to decide what will constitute an outliers with some methodology

- InterQuartile Range
- Standard deviation
- Visualization or Domain limit values-

- percentage of Data
keep in mind if a large percentage of your data is being labeled as an outlier, then you actually just have a wide distribution; not outliers.

- limit outlier to a few percentage points at most.

Outlier Consideration:-

- Utilize visualization plots to be able to see and identify outliers points -

- keep in mind this will create caveats for your future model (e.g. Model not suitable for houses priced over \$10 million)

Code

import numpy as np

import pandas as pd

import matplotlib.pyplot as plt

sns.displot(samples.)

sns.boxplot(^{data-}samples.)

ser = pd.Series(pd.sample)

ser.describe()

IQR = 55.25 - 42.0

1.5 * (IQR)

lowerlimit = 42.0 - 1.5 * (IQR)

(Q3)

$\text{sns}[\text{sns} > \text{lower limit}]$

75%

$\text{np.percentile}(\text{sample}, [75, 25])$

$$\text{IQR} = \text{q75} - \text{q25}$$

$$= 1.5 * \text{IQR}$$

$$\text{Outlier} = \text{q25} - 1.5 * \text{IQR}$$

$\text{df} = \text{pd.read_csv("read file name")}$
 $\text{df.corr()["Sales"].sort_values()}$

$\text{sns.scatterplot(y="SalesPrice", x="OverallQual")}$

$\text{sns.scatterplot(x="GrLivArea", y="SalesPrice")}$

$\text{df}[(\text{df['OverallQual']} > 8) \& (\text{df['SalesPrice']} < 200000)]$

$\text{df}[(\text{df['GrLivArea']} > 4000) \& (\text{df['SalesPrice']} < 400000)]$

$\text{dropIndex} = \text{df.index}$

$\text{df} = \text{df.drop(dropIndex, axis=0)}$

Answer

4 - Dealing with Missing Data

Code #

usually imports

```
with open('.../DATA/Ames-housing.txt', 'r')  
    prin(f.read())
```

as f

```
df = pd.read_csv("Ames-housing-demand.csv")
```

```
df.info()
```

```
df.head()
```

```
df = df.drop("PID", axis=1)
```

```
df.isnull().sum()
```

```
100 * df.isnull().sum() / length(df)
```

```
def percentmissing(df):
```

percent-nan = $100 * \frac{\text{df.isnull().sum()}}{\text{length(df)}}$

percent-nan = percentnan [percentnan 70
soft-value]

```
return percentnan
```

```
percentnan = percentmissing()
```

```
sns.barplot(x=percent-nan.index,
```

y=percent-nan)

```
plt.xticks(rotation=90);
```

5 - Dealing with Missing Data - part 2

percent_nan [percent_nan < 1]

100 / len(df)

df [df ["Electrical"].isnull()] ['GarageArea']

df [df ['Bsmt Half Bath'].isnull()]

df = df.dropna(axis=0, subset=['Electrical',
'GarageCars'])

percent_nan = percent_missing(df)

percent_nan [percent_nan < 1]

df [df ['Bsmt Half Bath'].isnull()]

Basement Numeric columns -->
fillnull(values)

num_cols = ["All basements"]

df['bsmt_num'] = df['bsmt_nuances cu'].
fillna(0)

6. Dealing With Missing Data

part 3: feature columns

- Dropping the feature column:-

- very simple to do

- No longer need to worry about that feature in the future.

- potential to lose a feature with possible important signal

- Should consider drop feature approach where many rows are NaN.

- Filling in the missing features data:

- potentially changing ground truth in data.

- Must decide on reasonable estimation to filled value.

- Must apply transformation to all future data for prediction.

- Simplest case:

- Replace all NaN values with a reasonable assumption (e.g. zero if assumed NaN implied zero)

- Harder cases:

- Must use statistical methods based on the columns to fill in NaN values.

pandas.transform (func =

call the groupby and fill values
based on groupby

df.groupby('Neighborhood')

[Not Filled]

transform(lambda

value : value.fillna)

value.mean))

7. Dealing with Categorical

`df['is subclass'] = df['Visitable'].apply(str)`

`direction = pd.Series([1, 'Up', 'down', 'down'])`

`pd.get_dummies(direction, drop first)`

`df1 = df.select_dtypes(exclude=['object'])`

`df2 = df.select_dtypes(exclude=['object'])`

`df object-dummies = pd.get_dummies(df1)`

`df_object-dummies = pd.get_dummies(df1, drop first
True)`

`final_df = pd.concat([df1, df2], axis=1)`

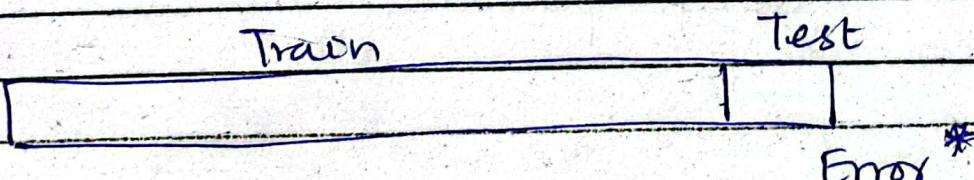
`df object-
dummies], axis=1])`

12. Cross Validation, Grid Search and the Linear Regression project :-

Section Overview :-

- Cross Validation in Detail
 - Train | Testsplit
 - Train | validation | Test split
 - Scikit-learn cross_val_score
 - Scikit-learn cross_val_predict
- Grid Search
- Linear Regression Project Exercise

2. Cross Validation: "Train Test Split"



allow us to adjust hyperparameters

Code #

```
x = df.drop(['Sales'], axis=1)  
y = df['Sales']
```

```
from sklearn.model_selection import train_test_split
```

```
x_train, x_test, y_train, y_test = train_test_split
```

(x, y, test_size = 0.3)

random_state = 101

```
from sklearn.preprocessing import
```

```
Scalar = StandardScaler()
```

```
Scalar.fit(X_train)
```

```
X_train = Scalar.transform(X_train)
```

```
X_test = Scalar.transform(X_test)
```

```
from sklearn.linear_model import Ridge
```

```
model = Ridge(alpha=100)
```

```
model.fit(X_train, y_train)
```

```
y_pred = model.predict(X_test)
```

```
mean_squared_error(y_test, y_pred)
```

```
# model = Ridge(alpha=1)
```

```
model_two = Ridge(alpha=1)
```

3.
—
3

- We just saw that train/Test split method has a disadvantages of not having a portion of data that can report a performance metrics on truly "unseen" data.

- while adjusting hyperparameters on test data is a fair technique and not typically referred to as "data leakage", it is a potential issue in regards to reporting



- Recall the entire reason to not adjust after the final test data set is to get a good fairest evaluation of the model -
- The model was not fitted to the final test data and the model hyperparameter were adjusted based off the final test data -
- This is truly never before seen data!

`df = df.head()`

`x = df.drop(['Sales'], axis=1)`

`y = df['Sales']`

test-train split
use two time ~~x_{test}~~ \leftarrow

≈ 0.31

`X-train, X-other, y-train, y-other = train_test_split(X, y)`

`X-validate, y-test, y-validate, y-test = train_test_split`

`from sklearn.preprocessing import StandardScaler`

`scaler = StandardScaler()`

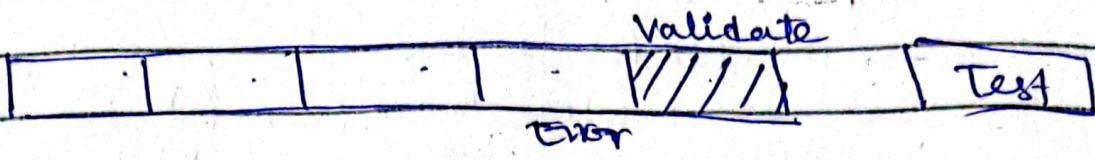
`scaler.fit(X-train)`

`X-train =`

`X-test =`

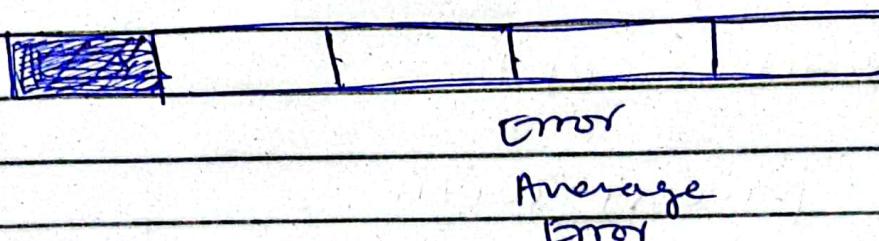
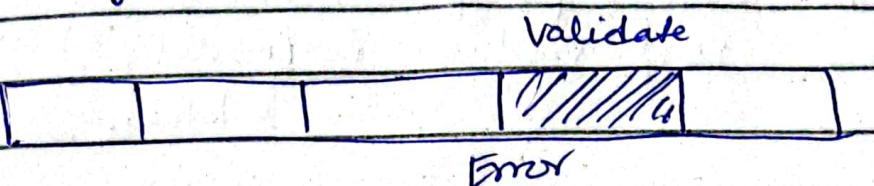
4. Cross Validation - cross_val_score

- Start with Training data and



choose k-fold values

- largest & possible "leave one out"



and final test.

- The cross_val_score function uses a model and training set (along with a k and chosen metrics) to perform all of this for us automatically!

- This allows for k-fold cross-validation to be performed on any model.
- lets explore how to use it -

`x = df.drop(['Sales'], axis=1)`

`y = df['Sales']`

`x_train, x_test, y_train, y_test = train_test_split(x, y, test_size=0.3, random_state=42)`

Scale data

`scaler.fit(x_train)`

`x_train = scaler.transform(x_train)`
`x_test = scaler.transform(x_test)`

`model = Ridge(alpha=100)`

`score = cross_val_score(model, x_train, y_train, cv=5)`

`x_train,`

`scoring='neg_mean_squared_error'`

`cv=5`)

Scores:

`abs(Scores.mean())`

same process with

`alpha=1`

5 - Cross Validation.

- Cross-Validate

- The cross validate function allow us to view multiple performance metrics from cross validate on a model and explore how much time fitting and testing took - .

- lets quickly review how to use this function call -

Code #

same code as below for
Setup

But

Score = cross_val_score(model, X_train,
Y_train,

scoring = [
'neg_mean_squared_error',
'neg_mean_absolute_error'], N=10])

6- Grid Search

- Often more complex models have multiple adjustable hyperparameters -
- A grid search is a way of training and validating a model on every possible combination of multiple hyperparameter options -
- Scikit learn includes a GridSearchCV class capable of testing a dictionary of multiple hyperparameter option through cross-validation -
- This allows us for both cross-validation and a grid search to be performed in a generalized way for any model.

#1 code

Same as before

from sklearn.model_selection import ElasticNet

model = ElasticNet()

param_grid = { 'alpha': [0.1, 1, 5, 10, 50],
 'l1_ratio': [0.1, 0.5, 0.7, 0.9] } }

from sklearn.model_selection import

{}