

Data Analysis with Python

Cheat Sheet: Model Evaluation and Refinement

Process	Description	Code Example
Splitting data for training and testing	The process involves first separating the target attribute from the rest of the data.	1. 1
	Treat the target attribute as the output and the rest of the data as input. Now split the input and output datasets into training and testing subsets.	2. 2 3. 3 4. 4 1. from sklearn.model_selection import train_test_split 2. y_data = df['target_attribute'] 3. x_data=df.drop('target_attribute',axis=1) 4. x_train, x_test, y_train, y_test = train_test_split(x_data, y_data, test_size=0.10, random_state=1)
Cross validation score	Without sufficient data, you go for cross validation, which involves	1. 1 2. 2 3. 3 4. 4 5. 5 6. 6 1. from sklearn.model_selection import cross_val_score 2. from sklearn.linear_model import LinearRegression lre=LinearRegression()

Copied!

creating
different
subsets of
training and
testing data
multiple
times and
evaluating
performance
across all of
them using
the R^2 value.

```
3. Rcross = cross_val_score(lre,x_data[['attribute_1']],y_data,cv=n)
4. # n indicates number of times, or folds, for which the cross validation is to be done
5. Mean = Rcross.mean()
6. Std_dev = Rcross.std()
```

Copied!

Cross
validation
prediction

Use a cross
validated
model to
create
prediction of
the output.

```
1. 1
2. 2
3. 3
4. 4

1. from sklearn.model_selection import cross_val_score
2. from sklearn.linear_model import LinearRegression
3. lre=LinearRegression()
4. yhat = cross_val_predict(lre,x_data[['attribute_1']], y_data,cv=4)
```

Copied!

Ridge
Regression
and
Prediction

To create a
better fitting
polynomial
regression
model, like ,
one that
avoids
overfitting to
the training
data, we use
the Ridge
regression
model with a

```
1. 1
2. 2
3. 3
4. 4
5. 5
6. 6

1. from sklearn.linear_model import Ridge
2. pr=PolynomialFeatures(degree=2) x_train_pr=pr.fit_transform(x_train[['attribute_1', 'attribute_2', ...]])
3. x_test_pr=pr.fit_transform(x_test[['attribute_1', 'attribute_2', ...]])
4. RigeModel=Ridge(alpha=1)
5. RigeModel.fit(x_train_pr, y_train)
6. yhat = RigeModel.predict(x_test_pr)
```

Copied!

parameter
alpha that is
used to
modify the
effect of
higher-order
parameters
on the model
prediction.

Use Grid
Search to
find the
correct alpha
value for
which the
Ridge

regression
model gives
the best
performance.
It further
uses cross-
validation to
create a
more refined
model.

```
1. 1
2. 2
3. 3
4. 4
5. 5
6. 6
7. 7
```

```
1. from sklearn.model_selection import GridSearchCV
2. from sklearn.linear_model import Ridge
3. parameters= [{'alpha': [0.001,0.1,1, 10, 100, 1000, 10000, ...]}]
4. RR=Ridge()
5. Grid1 = GridSearchCV(RR, parameters1,cv=4) Grid1.fit(x_data[['attribute_1', 'attribute_2', ...]], y_data)
6. BestRR=Grid1.best_estimator_
7. BestRR.score(x_test[['attribute_1', 'attribute_2', ...]], y_test)
```

Copied!



Skills Network