

✓ CatBoost Regression

```
# from google.colab import files
# up = files.upload()
```

✓ load dataset

```
import pandas as pd
df = pd.read_csv('df.csv')
df.head(3)
```



	f1	f2	f3	f4	T
0	16.5	202.0	865.500000	1880.0	50.000000
1	18.0	204.0	688.000000	1738.5	44.000000
2	18.0	203.0	583.666667	1470.0	66.666667

```
# df.info()
```

✓ cleaning

```
# clean data
```

✓ encoding

```
# encode data
```

✓ define x, y

```
import numpy as np
x = df[['f1', 'f2', 'f3']].values
y = df['T'].values
```

✓ splitting

```
### finding best random state

# from sklearn.model_selection import train_test_split
# from catboost import CatBoostRegressor
# from sklearn.metrics import r2_score

# import time
# t1 = time.time()
# lst = []
# for i in range(1,10):
#     x_train, x_test, y_train, y_test = train_test_split(x, y, test_size=0.25, random_state=i)
#     cbr = CatBoostRegressor(verbose=0, random_state=1)
#     cbr.fit(x_train, y_train)
#     yhat_test = cbr.predict(x_test)
#     r2 = r2_score(y_test, yhat_test)
#     lst.append(r2)
# t2 = time.time()
# print(f"run time: {round((t2 - t1) / 60 , 0)} min")
# print(f"R2_score = {round(max(lst),2)}")
# print(f"random_state = {np.argmax(lst) + 1}")
```

```
from sklearn.model_selection import train_test_split
x_train, x_test, y_train, y_test = train_test_split(x, y, test_size=0.2, random_state=42)
```

✓ scaling

```
# catboost Regression doesn't need scaling
```

✓ fit train data

```
### k-fold cross validation

# from catboost import CatBoostRegressor
# from sklearn.model_selection import GridSearchCV

# parameters = {
#     '': [],
#     '': []
# }

# cb = CatBoostRegressor(random_state=42)
# gs = GridSearchCV(estimator=cb, param_grid=parameters, cv=5)

# gs.fit(x_train, y_train)

# best_params = gs.best_params_
# print(best_params)

from catboost import CatBoostRegressor
cbr = CatBoostRegressor(random_state=1, verbose=0)
cbr.fit(x_train, y_train)
```

```
<catboost.core.CatBoostRegressor at 0x18a92749c40>
```

✓ predict test data

```
yhat_test = cbr.predict(x_test)
```

✓ evaluate the model

```
from sklearn.metrics import r2_score
print("r2-score (train data): %0.4f" % r2_score(y_train, cbr.predict(x_train)))
print("r2-score (test data): %0.4f" % r2_score(y_test, yhat_test))
```

```
r2-score (train data): 0.9667
r2-score (test data): 0.2675
```

```
from sklearn.metrics import mean_squared_error
from sklearn.metrics import mean_absolute_error
print(f"MSE (train data): {mean_squared_error(y_train, cbr.predict(x_train))}")
print(f"MAE (train data): {mean_absolute_error(y_train, cbr.predict(x_train))}")
print(f"MSE (test data): {mean_squared_error(y_test, yhat_test)}")
print(f"MAE (test data): {mean_absolute_error(y_test, yhat_test)}")
```

```
MSE (train data): 6.176779536150266
MAE (train data): 1.9914993568258645
MSE (test data): 109.40535187764313
MAE (test data): 8.523940294948932
```

✓ save the model

```
# import joblib
# joblib.dump(cbr, 'cbr_model.pkl')
```

✓ load the model

```
# import joblib
# cbr = joblib.load('cbr_model.pkl')
```