

Decision Tree Regression

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

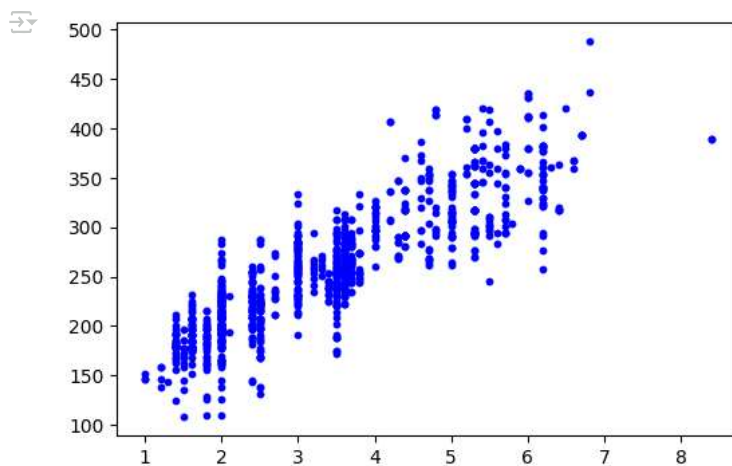
import dataset

```
import pandas as pd
df = pd.read_csv('dataset.csv')
df = df[['A', 'T']]
df.head(3)
```

```
↗
```

	A	T
0	2.0	196
1	2.4	221
2	1.5	136

```
import matplotlib.pyplot as plt
plt.figure(figsize=(6,4))
plt.scatter(df[['A']], df[['T']], s=10, c='b')
plt.show()
```



cleaning

```
# clean the data
```

encoding

```
# encode the data
```

define x, y

```
import numpy as np
x = df[['A']].values
y = df[['T']].values
```

splitting

```
from sklearn.model_selection import train_test_split
x_train, x_test, y_train, y_test = train_test_split(x, y, test_size=0.25, random_state=42)
```

```
y_train[:5]
```

```
↗ array([258, 212, 317, 308, 301])
```

```
### finding best random state
```

```
# from sklearn.model_selection import train_test_split
# from sklearn.tree import DecisionTreeRegressor
# from sklearn.metrics import r2_score

# 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)
#     dtr = DecisionTreeRegressor()
#     dtr.fit(x,y)
#     yhat_test = dtr.predict(x_test)
#     r2 = r2_score(y_test, yhat_test)
#     lst.append(r2)
# print(f"r2_score: {round(max(lst), 2)}")
# rs = np.argmax(lst) + 1
# print(f"random_state: {rs}")
```

▼ scaling

```
# Decision Tree Regression doesn't need scaling
```

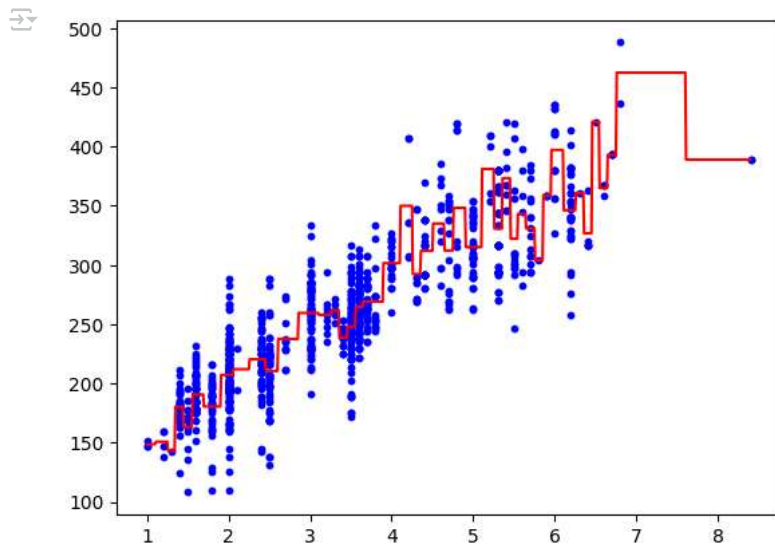
▼ train the model

```
# def param
# criterion='squared_error', splitter='best', max_depth=None
# min_samples_split=2, min_samples_leaf=1, min_weight_fraction_leaf=0.0
# max_features=None, random_state=None, max_leaf_nodes=None
# min_impurity_decrease=0.0, ccp_alpha=0.0, monotonic_cst=None
```

```
from sklearn.tree import DecisionTreeRegressor
dtr = DecisionTreeRegressor()
dtr.fit(x,y)
```

```
↗ ▾ DecisionTreeRegressor ⓘ ?
DecisionTreeRegressor()
```

```
plt.scatter(x,y, color='b', s=10)
xx = np.arange(np.min(x), np.max(x), 0.01).reshape(-1, 1)
plt.plot(xx, dtr.predict(xx), color='r')
plt.show()
```



```
### K-fold cross validation
```

```
# from sklearn.tree import DecisionTreeRegressor
# from sklearn.model_selection import GridSearchCV

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

# dt = DecisionTreeRegressor(random_state=42)
# gs = GridSearchCV(estimator=dt, param_grid=parameters, cv=5)

# gs.fit(x_train, y_train)

# best_params = gs.best_params_
# print(best_params)
```

✓ predict test data

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

✓ evaluate the model

```
from sklearn.metrics import r2_score
print(f"r2_score: {r2_score(y_test, yhat_test)}")
```

```
➦ r2_score: 0.8333692306410891
```

✓ predict new data

```
dtr.predict([[5.5]])
```

```
➦ array([322.21428571])
```

✓ save the model

```
# import joblib
# joblib.dump(dtr, 'dtr_model.pkl')
```

✓ load the model

```
# import joblib
# dtr = joblib.load('dtr_model.pkl')
```