→ STUDENTS _ Performance

→ Dataset used — student Performance.

→ It contains 10,000 records & 6 columns

→ linear regression model, XG Boost models used.

※ STEPS OF CODE

① importing the libraries. like numpy, Pandas, matplot, seaborn.

→ numpy & Pandas are used for numerical & data operations.

→ seaborn & matplot are used for visualization.

② Ignore warnings.

```
import warnings
warnings.Filterwarning ('ignore')
```

③ Loading the dataset.

```
df = pd.read-csv ("DrivePath").
df.head(5)
∵ head = First Five rows.
```

④ summarizing the dataset

```
df.describe()
df.info()
```

→ provides summary including mean & standard deviation.

→ shows dataset structure, datatypes & missing values

⑤ identifying & remove duplicate record

> df. duplicated . sum()
>
> df = df. drop_ duplicates( )

→ method checks for duplicate rows in dataset
→ counts total number of duplicate rows in
dataset [.sum()].

→ df.drop_duplicates_removes duplicate
rows from the dataset.

⑥ we use Boolean series because to change
binary into numerical format. {true/false}

df [`extracurricular Activities`] = (df(`extracurricular Activities`
== `yes` ). astype(int).

✳ → astype(int) — which converts into integer
types (0, 1)

⑦ Linear regression model

~~imp~~

from sklearn. model_selection import train_test_split
from sklearn.linearmodel import linear regression.

→ importing the models selection
to train & test.

will map the required output & will
training, test the data using
linear regression model.

X_train, y_train → used to train the
model.

x_test, y_test → used to test the
model

⑨ Now, will print the linear regression
score of the tested model

⟶ XGBOOST model!

⑩ will import required libaries.

⑪ Now will be encoding the data because
it converts text into numbers.

label-encoder = LabelencoderC).

⑫ will split the data which we need
as a output.

⑬ If encoded is needed then will do, if
not will train the data & test the
data.

⑭ will install xgboost.

! pipinstall - - upgrade scikit-learnxgboost

(ii) ... ... ...

(...) ... the higher ... ... ...

the data into ... find, e ...

17) we use GridSearchCV to find the best parameters.

18) Trains the best model using x, y, train

19) Now will test the model (calculate mean absolute error (MAS) to get accuracy.

20) Now print the outputs of trained, tested & mas test on data s.

*=> Here question is to get the best model & best students performance

→ code

↓

bestPerformance = df ['performance index'] max()

best_student = df [df ['performance index'] == best performance]

#print results.

print ("best students performance index:", best performance)

print (best _student).

# * Keypoints to remember.

① It contains student performance & 10,000 records.

② we use linear regression & XGBoost model

③ we conclude the best model is XGBoost & the best student is the student whose performance index is 100 (max). ✓