

①

T15 - Hands-on

- To open pdf file go to application → file Manager → desktop → project → wings T15 → instructions → drag that wings T15 . file .
- Already exist data in S3 → loan-data check whether existing given data(ppt) present or not .
- * Select the region as : US-East-1
→ right side top corner .
- * First Create Redshift cluster
go to console → type as Redshift
→ click on Amazon Redshift .
→ Next go down → click on Create Cluster .

(2)

(Copy & paste
this name)

give name as emr-spark-redshift[↑].

choose as all - choose

Next select dc2 - large.

Number of node : 1

→ Database configurations

Manually add the admin password

password is 'awsuser1'. click on
show password. ↪ (copy & paste)

→ a Newly created is database
Encryption in that select → AWS (KMS)

Manage Key. → go down → click on
create cluster.

→ Backend it will load.

Next we have to create EMR cluster

→ Click on Create cluster

Name: My Cluster. ↪ (select)

Select Amazon EMR release is emr-7.1.0

- ③ → Select → * Hadoop 3.3.6 * Sqoop 1.4.7
* Spark 3.5.0 * Hive 3.1.3
- go to primary → M4, large (select)
→ Remove instance group (2 times)
- go to cluster termination → Manually terminate the cluster
- go to security configuration & ECR key pair. → click on create key pair
→ give key pair name as emr-spark
Click on Create key pair.
- Now go Back to previous tag
Click on Browse → Selected → click on
- Choose → Next select service role as AWS-S3-RDS-redshift - emr → Now go down
Select create an instance → select all S3. Click on create cluster.
- go to coding part
- Application → development → Open visual studio code

→ go to file → open folder \rightarrow desktop \rightarrow project folder \rightarrow select wings + 15. then open.

→ go to left side click on python

→ Next click on challenge.py

→ How to take Bucket

go to console \rightarrow in search bar search for S3 inside \rightarrow load date click that

then copy that Bucket name

→ go to coding part \rightarrow control 'F'
type Bucket Name remove that space
in double quotes. (28 & 68 lines) (2 places)

→ get the jdbc URL from Redshift.

go to console \rightarrow Amazon Redshift \rightarrow click on Cluster snr-spark-redshift after clicking Right side will get jdbc URL

After writing the coding part click on control 'S'. save the code.

⑤ Coding part
④ record the data

```
df = spark.read.csv(S3-input-path,  
header=True, schema=customSchema).
```

2) Clean data lo.

```
df = input_df.dropna().dropDuplicates()  
df = df.filter(df.purpose != "null")
```

3) load data lo.

```
data.coalesce(1).write.csv(output-path,  
header=True, mode="overwrite").
```

4) def result_1(input_df):

```
df = input_df.filter(ccol("purpose") ==  
"Educational" | ccol("purpose") ==  
"Small-business")
```

(6)

$df = df \cdot \text{withColumn}(\text{"income-to-instalment-ratio"}, \text{col}(\text{"log-annual-income"}) / \text{col}(\text{"installment"}))$

$df = df \cdot \text{withColumn}(\text{"int-rate-category"},$
when($\text{col}(\text{"int-rate"}) < 0.1$, "low")
• when($(\text{col}(\text{"int-rate"}) \geq 0.1) \& (\text{col}(\text{"int-rate"}) < 0.15)$, "medium")
• otherwise ("high"))

$df = df \cdot \text{withColumn}(\text{"high-risk-borrower"},$
when($(\text{col}(\text{"dti"}) > 20) \mid (\text{col}(\text{"fico"}) < 700) \mid (\text{col}(\text{"avgol-util"}) > 80, 1)$)
• otherwise (0))

result - 2

$df = \text{input-db} \cdot \text{groupBy}(\text{"purpose"}).$ ~~agg~~
 $\text{agg}(\text{sum}(\text{col}(\text{"not-fully-paid"})) /$
 $\text{count}(\text{"*"})).\text{alias}(\text{"default-rate"})$
 $df = df \cdot \text{withColumn}(\text{"default-rate"},$
 $\text{round}(\text{col}(\text{"default-rate"}), 2))$

cdshft → load-data .

⑦

→ jdbcURL = "copy jdDC URI",
username = "ausugex"
password = "Aususer1")
table-name = "result_2"

→ data.write.
• format("jdbc")
• option("url", jdbcURL).
• option("dbtable", table-name)
• option("user", username).
• option("password", password).
• mode("overwrite").
• save().

After writing the code.

we will push the code to EMR cluster.

first we go to EC2 → instances → right side. Select that instance ID → go to actions → security → change the security groups → search for my security group → Add the security group → next we have to save. → next in instances select Connect → leave as it is then click on connect. then in new tab EMR cluster will open → go to Visual code → on the top click on Terminal → open New terminal. → go to application → file manager → downloads drag that file (emr-spark-bin (emr-spark.pem) drag to (emr-copy-py).

→ Next it will open Beside challenge.py
it will open in emr-spark-pem
now copy that python emr_copy.py
paste in emr-spark-pem then
click on run. It will run successfully.

→ Next go to EMR cluster.

→ give command. a) [cd /home/hadoop
Enters → [ls] Enters → [cd python] Enters
→ [ls] enter → (Now path is in challenge.py)
→ [cd ..] (to go back)
→ Enters [ls] → goto (cd setup) →
Next Enters Bash setup.sh → It will
load for sometime.
Now to go back Enters → [cd ..]
Enters [ls] → Next goto python folder
Enters [cd python] → Enters ls →
Now we have submit the code
EMR type as [spark-submit challenge.py]
it will execute.

→ How to check loading data into Redshift. (10)

→ go to Amazon Redshift → click on that

→ click on Emr-spark-redshift → right side
click on Query data → Query data ^{editor} is v2

→ check in S3.

click on Emr-Spark-redshift

→ then select database user name & pswd.
database → dev, Username → ansuser →
password → AwSuser1, Show password,
click on Create connection after creating
connection

→ Once go to S3 → Click on load-data
~~check the~~ Click on that → Check the o/p

Now go to left Emr-spark_redshift

→ dev → public → tables → result-2

Now we can Query as.

Select * from result-2, click on run.

check
o/p data
click on
clean data

(11)

Opas: purpose default rate

educational 0.29

small business 0.35

all other 0.21

debt consolidation 0.17

credit card 0.19

Main purchase 0.12