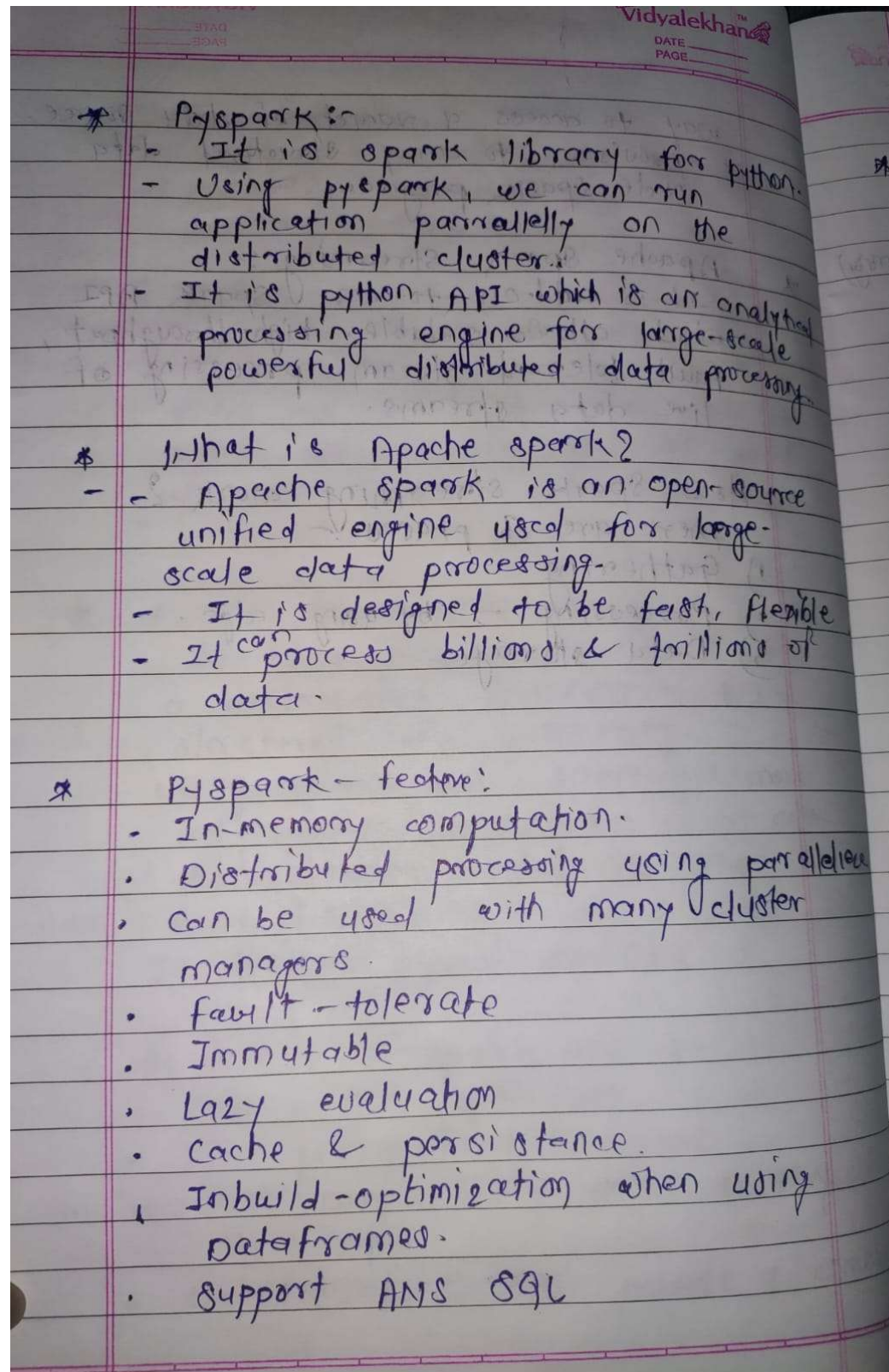


Name: Krushnakumar Patle

Email: [krishnapatle128@gmail.com](mailto:krishnapatle128@gmail.com)

Batch: Data Engineering Batch-1



### \* Advantages of PySpark:-

- PySpark is a general-purpose, in-memory distributed processing engine that allows you to process data efficiently in a distributed fashion.
- It is faster than 100x than traditional system.
- It get great benefit for data ingestion pipeline.
- Using it we can process data from Hadoop, AWS, & many more.

### Apache Kafka:-

It is an open-source distributed streaming system used for stream processing, real-time data pipeline, and data integration at scale.

### PySpark version:-

Supported - Python 3.8, Java 8, 11, 13, 17 & latest version, are deprecated.

Scala - 2.12 & 2.13

R - 3.5

- \* PySpark modules:
- PySpark RDD (PySpark.PPD)
  - PySpark DataFrame & SQL
  - PySpark Streaming (PySpark.streaming)
  - PySpark MLlib (PySpark.ml, .mlib)
  - PySpark GraphFrames (GraphFrames)
  - PySpark Resource (It's new in 3.0)



```
In [1]: import pyspark
from pyspark.sql import SparkSession
spark = SparkSession.builder.appName("practice").getOrCreate()

spark
```

```
Out[1]: SparkSession - in-memory
SparkContext

Spark UI
Version
v3.5.0
Master
local[*]
AppName
practice
```

```
In [2]: import pandas as pd
log_data = pd.read_csv("Marks_data.csv")
print(log_data)
```

	Name	M1 Score	M2 Score	age
0	Alex	62.0	80.0	20
1	Brad	45.0	56.0	19
2	Joey	85.0	98.0	21
3	NaN	54.0	79.0	20
4	abhi	NaN	NaN	20

```
In [3]: df = spark.read.csv("Marks_data.csv")
df
```

```
Out[3]: DataFrame[_c0: string, _c1: string, _c2: string, _c3: string]
```

```
In [4]: df.show()
```

	_c0	_c1	_c2	_c3
	Name	M1 Score	M2 Score	age
	Alex	62	80	20
	Brad	45	56	19
	Joey	85	98	21
	NULL	54	79	20
	abhi	NULL	NULL	20

```
In [5]: datalist = [("Java", 20000), ("Python", 100000), ("Scala", 3000)]
rdd=spark.sparkContext.parallelize(datalist)
rdd.collect()
```

```
Out[5]: [('Java', 20000), ('Python', 100000), ('Scala', 3000)]
```