

# Basket Analysis By Association Rule in Tech Skills

\*

1<sup>st</sup> Yousef sarkar  
Computer Science -Big Data.)  
Umm Al Qura University)  
Makkah , Saudi Arabia  
yousef0sarkar@gmail.com

**Abstract**—In this project, the latest data mining algorithms were used to find the related elements to process the related skills at the big data level, and the tools and techniques were integrated to produce this solution.

## I. INTRODUCTION AND MOTIVATION

there are a lot of skills in tech jobs such as programming languages or frameworks in technology sectors some of them is related together and other skills are not , so we aim to provide helping tool that can help job seeker to find a set of skills related together as basket or group of skills ,based in main skill that user defined . we will using big data technology to solve this problem in this project.

## II. RELATED WORK

[1]Unsupervise learned word embeddings have seen tremendous success in numerous Natural Language Processing (NLP) tasks in recent years in Skill2Vec . [2]Skill2vec: Machine Learning Approach for Determining the Relevant Skills from Job Description

## III. PROJECT IDEA DESCRIPTION AND CHALLENGES

The main idea of the project is the possibility of searching for skills that related to each other, as job advertisements include some requirements, including skills Our goal is to learn the skills associated with each other through A key can be entered by the user for a field, and then Associated skills are shown as output .

we aim to Mining frequent skills in dataset, the size of the dataset is 50k rows, and every row has a set of skills, we need to get the frequent item based in threshold confidence so we will use algorithms. such as FP-growth to find the set when new data is inserted into the dataset

in our project we have use a list of tools :

- 1) databricks to run code in cluster .
- 2) Pyspark as frameworks of big data analysis
- 3) python
- 4) data set of job skills 50k as CSV file

The project workflow in these steps :

- 1) getting the Dataset We searched the dataset for job skills and we chose the data set from the paperswithcode website the data was a sample from a data set 5GB

The format of data :

9978	127681	interfaces	logos	who	knowledge	html	adobe photoshop
9979	3094	disco	wireless	wireless_network	airmagnet	cwna	ieee_802.11
9980	69440	IT Project Management	Business Process Re-engineering	Software Development Life Cycle	Project Review	Consulting	Business Analysis
9981	65380	BI	MS SQL	Data Warehousing	SSAS	Application Development	Data Management
9982	114395	Hosting	HTML	Wordpress	Javascript	JQuery	
9983	97435	C#	WCF	SQL Server	Winforms	JQuery	
9984	68454	BE Studies	Clinical Project Management	Program manager	sponsor coordination	sponsor	study experience
9985	8728	medical_writing	epidemiology	publications	medical_writer	medical_write	oncology
9986	66890	MBBS	Steam Sterilizers	NaN	NaN	NaN	
9987	11926	Security Business Operations	Business Operations Specialist	MIS	business analysis	Strategy planning	
9988	37047	ASP.Net	WCF	MVC	C#	SQL Server	

- 2) create free account in databricks
- 3) create and run new cluster in databricks
- 4) upload dataset file to databricks and create notepad to do handle this data

	_c0	_c1	_c2	_c3	_c4
1	125720	HR Executive	screening	selection	interview
2	112708	Special Teacher	Teaching	Education	null
3	115226	consulting	freelance	IT helpdesk	Technical
4	19805	diploma	machining	onc m	mould
5	80308	Compensation	Benefits	HR Functions	Alm
6	64086	Storage Administrator	null	null	null
7	48468	HR Operations	Ext Formalities	Shortlisting	Screening
8	122729	Simulink	stateflow	Matlab developer	tar

- 5) preprocessing the data

The data extracted from the file must be processed in a way that pyspark can handle The framework cannot treat it as a normal dataframe, so the closest equivalent to this formula must be used The columns to be worked on must also be specified as a basket

```
1 fp = FPGrowth(minSupport=0.001, minConfidence=0.001, itemSetCol='basket', predictionCol='prediction')
2 model = fp.Fit(df_aggregated)
```

Python

3/3 Spark Jobs

- Job 48 View (Stages: 1/1)
- Job 49 View (Stages: 1/1)
- Job 50 View (Stages: 2/2)

Command took 7.25 seconds -- by s439817822@p1-ucp.edu.sa at 6/4/2022, 3:46:154 AM on Rep19802



#### IV. REFLECTION

In this project, through this term, several tools and techniques were made that increased our software and technical knowledge, and we can summarize them in the following: We used a platform that simulates distributed systems, clusterer and we used Data Brix We worked on several languages and frameworks, primarily on Python and the Pyspark framework as a task distribution system. Through the databricks, we used the built in algorithms such as FP-Growth algorithm and the Prefix Span algorithm., which helped us to find the frequent item set and Association rule , I also learned the dealing with various files and different types of datatype and preprocessing them before analysis , so in the future I can use these tools perfectly to solve this type of problem .

#### V. CONCLUSION

Through big data techniques, data was mined through association rule algorithms in our case FP-GROUTH, and this thing was done through the databricks platform and in the future it will be transferred to a web application that deployed in internet .

#### REFERENCES

- 1) <https://smartbridge.com/market-basket-analysis-101/>
- 2) <https://spark.apache.org/docs/latest/ml-frequent-pattern-mining.html>
- 3) <https://databricks.com/>
- 4) <https://www.oracle.com/big-data/what-is-big-data/>
- 5) <https://spark.apache.org/>