

Mining the US Technologists Data

Abulitibu Tugulu

DHI Group Inc., CU Boulder

July 3, 2022

Outline

Title

Team member

Description

Prior Work

Datasets

Proposed work

List of tool(s)sets

Evaluation

Title

Mining the US Technologists Data

Outline

Title

Team member

Description

Prior Work

Datasets

Proposed work

List of tool(s)sets

Evaluation

Team member

Group 5: Abulitibu Tuguluke

Outline

Title

Team member

Description

Prior Work

Datasets

Proposed work

List of tool(s)sets

Evaluation

Description

There are many reports stating US is (or will) face a technology workers (Technologists). Our study will mine through the US H1B data set along with Dice.com's job and candidate profiles (along with skill sets for each job title), to try to understand which technologist job are in increasing/decreasing demands, which skill sets are more popular/unpopular, through clustering, time and geo-location analysis. Our goal is to better understand the US technology market and what we can do to help people who try to get into tech market by providing them with skillset guidance, answers the questions such as whether foreign workers are getting paid the same amount as their US colleges.

Outline

Title

Team member

Description

Prior Work

Datasets

Proposed work

List of tool(s)sets

Evaluation

Prior Work

There is already a skillsets collection been done on tons of job descriptions for technologist through data mining. We have 1432 job titles associate with perticular skill sets. Unfortunately, most of the online literatures are focused on salary analysis, which motivates job seeker, where we will do our own salary report as well.

Papers

- ▶ Curriculum Vitae Recommendation Based on Text Mining by Honorio Apaza Alanoca, Americo A. Rubin de Celis Vidal, Josimar Edinson Chire Saire

Outline

Title

Team member

Description

Prior Work

Datasets

Proposed work

List of tool(s)sets

Evaluation

Datasets

- ▶ H1B data scraped from [https : //h1bdata.info/index.php](https://h1bdata.info/index.php)
- ▶ Job skill set (Internal data)
- ▶ US Bureau of Labor Statistics
[https : //www.bls.gov/developers/api_python.htm](https://www.bls.gov/developers/api_python.htm)
- ▶ Data mined and stored on my machine

Data Description for scraped for H1b techonologist data

	Dice_job_title	EMPLOYER	JOB TITLE	BASE SALARY	LOCATION	SUBMIT DATE	START DATE	Job_Title
count	2495079	2495052	2467304	2467304	2467304	2467304	2467304	2495079
unique	1180	66795	1770	37871	10474	2749	2898	1180
top	Systems Analyst	TATA CONSULTANCY SERVICES LIMITED	SYSTEMS ANALYST	60,000	NEW YORK, NY	03/13/2015	09/01/2015	Systems Analyst
freq	198860	157580	198719	146505	136642	14268	32532	198860

Data Description for US Bureau of Labour statistics¹

	1208m		1209m		1210m		1211m		1212m		1213m		1214m		1215m		1216m		1217m		1218m		1219m		1220m		1221m		1222m		1223m		1224m		1225m		1226m		1227m		1228m		1229m		1230m		1231m		1232m		1233m		1234m		1235m		1236m		1237m		1238m		1239m		1240m		1241m		1242m		1243m		1244m		1245m		1246m		1247m		1248m		1249m		1250m		1251m		1252m		1253m		1254m		1255m		1256m		1257m		1258m		1259m		1260m		1261m		1262m		1263m		1264m		1265m		1266m		1267m		1268m		1269m		1270m		1271m		1272m		1273m		1274m		1275m		1276m		1277m		1278m		1279m		1280m		1281m		1282m		1283m		1284m		1285m		1286m		1287m		1288m		1289m		1290m		1291m		1292m		1293m		1294m		1295m		1296m		1297m		1298m		1299m		1300m		1301m		1302m		1303m		1304m		1305m		1306m		1307m		1308m		1309m		1310m		1311m		1312m		1313m		1314m		1315m		1316m		1317m		1318m		1319m		1320m		1321m		1322m		1323m		1324m		1325m		1326m		1327m		1328m		1329m		1330m		1331m		1332m		1333m		1334m		1335m		1336m		1337m		1338m		1339m		1340m		1341m		1342m		1343m		1344m		1345m		1346m		1347m		1348m		1349m		1350m		1351m		1352m		1353m
--	-------	--	-------	--	-------	--	-------	--	-------	--	-------	--	-------	--	-------	--	-------	--	-------	--	-------	--	-------	--	-------	--	-------	--	-------	--	-------	--	-------	--	-------	--	-------	--	-------	--	-------	--	-------	--	-------	--	-------	--	-------	--	-------	--	-------	--	-------	--	-------	--	-------	--	-------	--	-------	--	-------	--	-------	--	-------	--	-------	--	-------	--	-------	--	-------	--	-------	--	-------	--	-------	--	-------	--	-------	--	-------	--	-------	--	-------	--	-------	--	-------	--	-------	--	-------	--	-------	--	-------	--	-------	--	-------	--	-------	--	-------	--	-------	--	-------	--	-------	--	-------	--	-------	--	-------	--	-------	--	-------	--	-------	--	-------	--	-------	--	-------	--	-------	--	-------	--	-------	--	-------	--	-------	--	-------	--	-------	--	-------	--	-------	--	-------	--	-------	--	-------	--	-------	--	-------	--	-------	--	-------	--	-------	--	-------	--	-------	--	-------	--	-------	--	-------	--	-------	--	-------	--	-------	--	-------	--	-------	--	-------	--	-------	--	-------	--	-------	--	-------	--	-------	--	-------	--	-------	--	-------	--	-------	--	-------	--	-------	--	-------	--	-------	--	-------	--	-------	--	-------	--	-------	--	-------	--	-------	--	-------	--	-------	--	-------	--	-------	--	-------	--	-------	--	-------	--	-------	--	-------	--	-------	--	-------	--	-------	--	-------	--	-------	--	-------	--	-------	--	-------	--	-------	--	-------	--	-------	--	-------	--	-------	--	-------	--	-------	--	-------	--	-------	--	-------	--	-------	--	-------	--	-------

¹Internal jobs data will not be shared at any stage/part of this project.

Outline

Title

Team member

Description

Prior Work

Datasets

Proposed work

List of tool(s)sets

Evaluation

Proposed work

- ▶ Data scraping: Scraeping HTML data from the web and store them in seperate json files locally
- ▶ Data cleaning: Probably the most important part of the project
- ▶ Data preprocessing:
- ▶ Data integration: Integrate internal data with mined H1b data
- ▶ Data mining and analysis

Outline

Title

Team member

Description

Prior Work

Datasets

Proposed work

List of tool(s)sets

Evaluation

List of tool(s)sets

- ▶ BeautifulSoup for mining
- ▶ Python library for all the computing and analysis: notably NumPy, Pandas, SK-learn, SpaCy, etc.
- ▶ JSON for data store

Outline

Title

Team member

Description

Prior Work

Datasets

Proposed work

List of tool(s)sets

Evaluation

Evaluation

There are many ways we can evaluate the results. Clustering Performance Evaluation Metrics like Silhouette coefficient is a good start when applying clustering algorithm, for time series, all the error evaluations. Besides, comparing time series analysis data within each job we can compute the increasing or decreasing sides of each job along with correlations. Location analysis on which job are located where with respect to other jobs and locations. Finally, we are looking forward to apply all the relevant techniques we are currently learning in the class².

²More methods will be updated through work.