Dynamic Apache Spark cluster for economic studies

Symposium on Nuclear Electronics and Computing - NEC'2019

Iuliia Gavrilenko (REU Plekhanova, speaker) Tatyana Tikhomirova (REU Plekhanova) Mayank Sharma (CERN) Maarten Litmaath (CERN)



SIMPLE





Prerequisites

diversity of methods for analyzing economic data that should be applied while complementing each other

huge data volumes, incompleteness and heterogeneity of the initial information

process of analyzing economic data is not trivial

• SIMPLE Grid: <u>http://cern.ch/go/8JLH</u>

The level of employment in the regions of the Russian Federation

- **from** 3 vectors of macroeconomic data OLAP cube:
 - axis of factors:
 - per capita income of the population (RUB);
 - volume of paid services per capita (RUB);
 - the cost of a fixed set of consumer goods and services (RUB);
 - the rate of migration increase in population (people);
 - the average size of the assigned pensions per capita (RUB);
 - turnover per capita (RUB);
 - actual final household consumption per capita (RUB);
 - industrial production index (%);
 - need for employees, declared by employers to employment service agencies (people);
 - fixed investment (RUB);
 - axis of objects (Russian Federation regions);
 - axis of time (years, months).



The level of employment in the regions of the Russian Federation

to macroeconomic data matrix:

- axis of factors:
 - per capita income of the population (RUB);
 - volume of paid services per capita (RUB);
 - the cost of a fixed set of consumer goods and services (RUB);
 - the rate of migration increase in population (people);
 - the average size of the assigned pensions per capita (RUB);
 - turnover per capita (RUB);
 - actual final household consumption per capita (RUB);
 - industrial production index (%);
 - need for employees, declared by employers to employment service agencies (people);
 - fixed investment (RUB);
- axis of objects (Russian Federation regions);



Clustering: single linkage by factors



Clustering: single linkage by regions



m hclust (*, "single")

Dynamic Apache Spark cluster





Cluster Dendrogram



hclust (*, "complete")



SIMPLE: Project Structure



The SIMPLE Grid project is being developed by Mayank Sharma, software engineer with WLCG at CERN



cluster: Economic results

Dynamic Apache Spark

hclust (*, "complete"

Dynamic Apache Spark cluster

	Небезопасно — spark-hadoop-master-	0.cern.ch Č	0 1	D ■ m swan003.cern.ch
	Home Nodes o	of the cluster Namenode	information +	+ Nodes of the cluster Namenode information CERNBox LR
She e	Noc	des of the cluster	Logged in as: dr.who	FILE EDIT VIEW INSERT CELL KERNEL WIDGETS HELP Not Trusted Python 3 O P + * * + * * • • • • Diote • • • • • • • • •
Cluster About Nodes Node Labels Applications CONTRACTOR SUBMITTED ACCEPTED RUMNING FINISHED FAILED	Apps Apps Apps Apps Apps Completed 6 0 1 5 3 Cluster Nodes Decommissioning Nodes Decommissioning 5 0 0 Scheduler Metrics Scheduler Type Scheduling Resource Type	Containers Running Memory Used Memory Total Memory Reserved VCore Used 4.50 GB 15 GB 0 B 3 ssioned Nodes Lost Nodes Unhealthy Nodes Reboote Q Q Q Q Minimum Allocation Maximum Allocation Max	es VCores VCores Total Reserved 40 0 ed Nodes Shutdown Nodes Q cimum Cluster Application Priority	<pre>In [40]: from pyspark.ml.stat import Correlation matrix = Correlation.corr(train_df.select('features'), 'features') matrix_np = matrix.collect()[0]['pearson({})'.format('features')].values In [43]: import seaborn as sns import matplotlib.pyplot as plt matrix_np = matrix_np.reshape(len(['x1','x2','x3','x4','x5','x6','x7','x8','x9','x10']),-1) fig, ax = plt.subplots(figsize=(13,7)) ax = sns.heatmap(matrix_np, cmap="YlGnBu") ax.xaxis.set_ticklabels[('x1','x2','x3','x4','x5','x6','x7','x8','x9','x10'], rotation=0)</pre>
KILLED	Capacity Scheduler [MEMORY] <memory 20="" entries<="" show="" td="" ¢=""><td>ory:256, vCores:1> <memory:3072, vcores:4=""> 0</memory:3072,></td><td>Search:</td><td><pre>ax.yaxis.set_ticklabels(['x1','x2','x3','x4','x5','x6','x7','x8','x9','x10'], rotation=0) ax.set_title("Correlation Matrix") plt.tight_layout() plt.set()</pre></td></memory>	ory:256, vCores:1> <memory:3072, vcores:4=""> 0</memory:3072,>	Search:	<pre>ax.yaxis.set_ticklabels(['x1','x2','x3','x4','x5','x6','x7','x8','x9','x10'], rotation=0) ax.set_title("Correlation Matrix") plt.tight_layout() plt.set()</pre>
Tools	Node Labels Rack \diamond Node State \diamond Node Address Node HTTP /default- RUNNING spark-hadoop-add-1- spark-hadoop-add-1- spark-hadoop-add-1- spark-hadoop-add-1- spark-hadoop-add-1- spark-hadoop-add-1-	P Address ♦ Last health- update ♦ Health-report ♦ Containers upadd-1 Fri Apr 19 1 1.50	Mem Avail Avail ↓ VCores Used Avail ↓ VCores Avail ↓ Vcores Avail ↓ Vcores Avail ↓ Version ↓ Ve	Sion 5 x ¹ · · · · · · · · · · · · · · · · · · ·
	/default- RUNNING spark-hadoop- rack submit.cern.ch:45121	Ch.30/42 13:25:31-10:200 GB 2019 2019 0 0 Ip- Fri Apr 19 0 0 B ch:8042 13:25:31 + 10:200 2019 0 0	GB 3 GB 0 8 2.8.5	5 x ² - 2 - 2 - 2 - 2 - 2 - 2 - 2 - 2 - 2 -
	/default- RUNNING spark-hadoop-add-3- rack worker.cem.ch:35466 worker.cem.ch	pp-add-3- ch:8042 13:25:31 +0200 2019 pp. mdos E: Ass 10 0 0 0 0 0 0 0 0 0 0 0 0 0	3 GB 0 8 2.8.5	
	rack 0.cern.ch:44695 <u>goark-hadoop</u> /default- RUNNING spark-hadoop-add-2- <u>spark-hadoop</u>	<u>p-worker</u> Fn Apr 19 0 0 B 42 13:25:31 +0200 2019 <u>pp-add-2-</u> Fri Apr 19 2 3 GE	B 0 B 2 6 2.8.5	
	rack worker.cern.ch:40365 worker.cern.d	<u>ch:8042</u> 13:25:31 +0200 2019	First Previous 1 Next Last	
				x90.3
				x1 x2 x3 x4 x5 x6 x7 x8 x9 x10

Conclusions



The degree of tightness of the linear relationship between the main macroeconomic factors was measured

Constructed several models of the regional labor market, taking into account the main macroeconomic factors



Was obtained classification of the subjects of the Russian Federation according to the level of employment



A dynamic Apache Spark cluster deployed by the means of the SIMPLE framework developed at CERN was used for all analysis

The Community

GitHub Repositories

- https://github.com/JuliaGavrilenko/ site_level_config_file
- https://github.com/JuliaGavrilenko/ simple_spark_cluster_master

SIMPLE Project Website

Link: <u>http://wlcg-lightweight-sites.github.io</u>

Open Source Community

Name: WLCG Lightweight Sites Link: <u>http://cern.ch/go/Hz7S</u>

Simple Grid Specification

https://docs.google.com/document/d/1yp_96UXcwNO49cktnHtT61iNmTO0RgrSQukuNYqACpM/edit?usp=sharing

Thank you for your attention!