

ECPR Methods Summer School: Big Data Analysis in the Social Sciences

Pablo Barberá

School of International Relations
University of Southern California

`pablobarbera.com`

Networked Democracy Lab

`www.netdem.org`

Course website:

github.com/pablobarbera/ECPR-SC104

Course website

pablobarbera / ECPR-SC104

Unwatch ▾ 1

★ Star 0

🍴 Fork 0

Code

Issues 0

Pull requests 0

Projects 0

Wiki

Settings

Insights ▾

ECPR Summer School: Big Data Analysis in the Social Sciences <http://pablobarbera.com/ECPR-SC104> Edit

Add topics

2 commits

1 branch

0 releases

1 contributor

Branch: master ▾

New pull request

Create new file

Upload files

Find file

Clone or download ▾

pablobarbera Set theme jekyll-theme-minimal

Latest commit 5bba33c 29 seconds ago

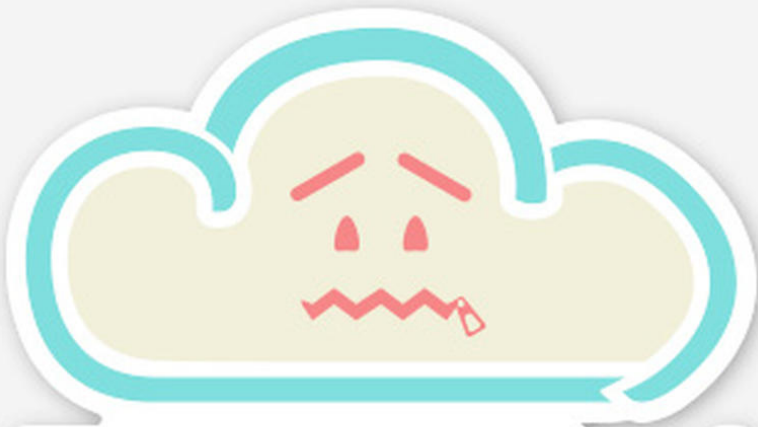
data	initial commit	12 minutes ago
day1	initial commit	12 minutes ago
day2	initial commit	12 minutes ago
day3	initial commit	12 minutes ago
day4	initial commit	12 minutes ago
day5	initial commit	12 minutes ago
html	initial commit	12 minutes ago
README.md	initial commit	12 minutes ago
_config.yml	Set theme jekyll-theme-minimal	29 seconds ago
packages.r	initial commit	12 minutes ago

README.md

Summer School: Big Data Analysis in the Social Sciences

github.com/pablobarbera/ECPR-SC104

Introduction to cloud computing



There is no cloud

it's just someone else's computer

Cloud computing

Use of remote servers hosted online to collect, store, and manipulate data rather than a personal computer.

Why moving to the cloud?

- ▶ **Scalability**: ability to increase memory or computing power to that necessary for our application
- ▶ **Cost**: pay for what you use, no need to buy expensive hardware
- ▶ **Speed**: easy and fast to launch a service on demand
- ▶ **Reliability**: cloud services include backup, redundancy, disaster recovery...
- ▶ **Accessibility**: multiple users can access and analyze data simultaneously, from anywhere with an internet connection

Cloud computing

What you need to know:

- ▶ **UNIX** commands to log in and interact with the server
- ▶ **SQL** to query large-scale databases

Companies offering cloud computing services:

- ▶ Amazon Web Services (AWS)
- ▶ Google Cloud Platform
- ▶ Microsoft Azure
- ▶ Digital Ocean
- ▶ ...many others

Basic UNIX commands

<code>cd <i>dirname</i></code>	Change directory
<code>mkdir <i>dirname</i></code>	Create new directory
<code>cp <i>oldfile newfile</i></code>	Copy a file
<code>mv <i>oldfile newfile</i></code>	Move a file
<code>ls -lh</code>	List your files (with sizes)
<code>cat <i>file</i></code>	Print file in console
<code>head <i>file</i></code>	Print first lines of file
<code>tail <i>file</i></code>	Print last lines of file
<code>wc -l <i>file</i></code>	Count lines in file
<code>grep <i>string file</i></code>	Regex on file text
<code>gzip <i>file</i></code>	Compress file
<code>wget <i>URL</i></code>	Download file from URL
<code>ps -u <i>user</i></code>	See running processes by user
<code>kill <i>process</i></code>	End running process

Databases

- ▶ **Database systems:** computerized mechanisms to store and retrieve data.
- ▶ **Relational databases:** data is represented as tables linked based on common keys (to avoid redundancy).

Customer

<i>cust_id</i>	<i>fname</i>	<i>lname</i>
1	George	Blake
2	Sue	Smith

Account

<i>account_id</i>	<i>product_cd</i>	<i>cust_id</i>	<i>balance</i>
103	CHK	1	\$75.00
104	SAV	1	\$250.00
105	CHK	2	\$783.64
106	MM	2	\$500.00
107	LOC	2	0

Product

<i>product_cd</i>	<i>name</i>
CHK	Checking
SAV	Savings
MM	Money market
LOC	Line of credit

Transaction

<i>txn_id</i>	<i>txn_type_cd</i>	<i>account_id</i>	<i>amount</i>	<i>date</i>
978	DBT	103	\$100.00	2004-01-22
979	CDT	103	\$25.00	2004-02-05
980	DBT	104	\$250.00	2004-03-09
981	DBT	105	\$1000.00	2004-03-25
982	CDT	105	\$138.50	2004-04-02
983	CDT	105	\$77.86	2004-04-04
984	DBT	106	\$500.00	2004-03-27

SQL

- ▶ SQL (pronounced S-Q-L or SEQUEL) is a language designed to **query relational databases**
- ▶ The result of an SQL query is always a table
- ▶ It's a **nonprocedural language**: define inputs and outputs; how the statement is executed is left to the *optimizer*
- ▶ How long SQL queries depends on optimization that is opaque to user (which is great!)
- ▶ SQL is a language that works with many commercial products:
 - ▶ Oracle Database, SQL Server (MS), MySQL, PostgreSQL, SQLite (all three open-source), Google BigQuery, Amazon Redshift...
 - ▶ Performance will vary, but generally faster than standard data frame manipulation in R (and much more scalable)