

larger websites:

- can take too long for many html requests
- may have detectors for crawlers

Twitter blocking: 180 calls / ~~hour~~ 15 min

Reddit limits: 60 calls/min is OK. (check Reddit Dev)
Ok; 3Q is the official

1 sec/call is base level of politeness (small ~~the~~ sites may not have infrastructure to check limits).
multiple calls to search API meant duplicate data (Twitter doesn't let you go back in time),
3200 tweets back in user history.
user view meant you can order things to avoid the duplicates.

tweepy pkg - very usable, altho' aimed @ application ~~at~~ developers. - so take things
out of python classes.

* KAIST "full" Twitter social graph (@ the time)

Reddit's usual API returns just most recent

But "cloudsearch" takes timestamps. (a backdoor!)

Using ^{unix} ~~option~~ Amit has exceptions mapped to him, . (just using ^{unix} mail) (sysadmin-type stuff that helps,
esp. b/c many servers.)

Unexpected error codes can happen.

Keying loss!

Amit ended up converting: C-Pickle → flatfile → SQL

This is best for 'find one record'.
But sequential flat is usually what you want.
R is good: it's a direct memory-map.
Python is doing a system-independent store.

esp. when going beyond
single queries.

Social-integrator

what can handle mark-up data that can handle ~~XML~~ buggy XML.
(just find what you are looking for).

some res: ~~Beautiful~~ BeautifulSoup.

URL-generation - finding the URLs? (URL target is said to be considered private).

wget/wcurl - shell scripting.

↓ can ~~not~~ specify useragent. Use the API, not on the commandline.
less vulnerable re: resync

also pandas, esp for historical data.

just try to get everything to fit in memory.

Flat files, hurray!

JSON is slow, so
convert to 1-line
version.

For reddit, too big
was faster to put
in SQL
(lots of data, lots
of fields)

Reddit: too big for
notes faster than
mySQL, but
lots of data on
one machine kills
the memory.