

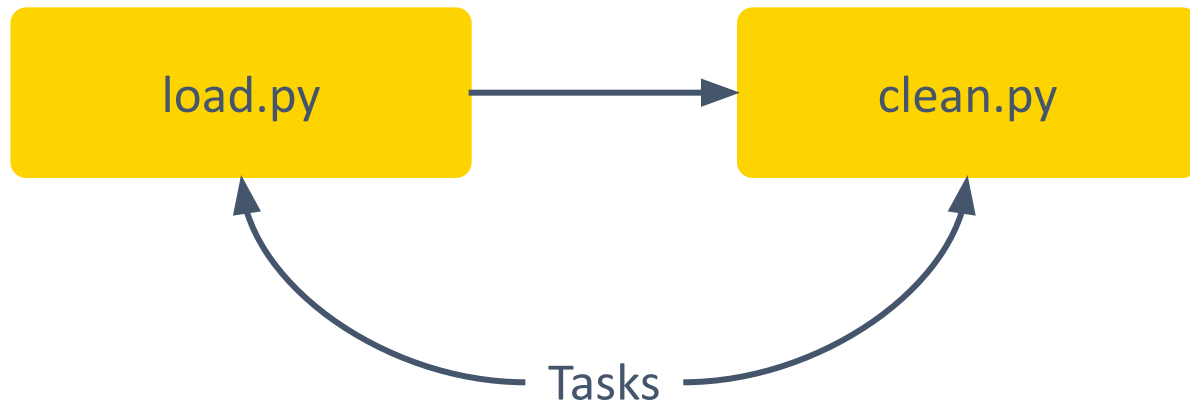
Дата-инженер

# Инструменты оценки качества данных

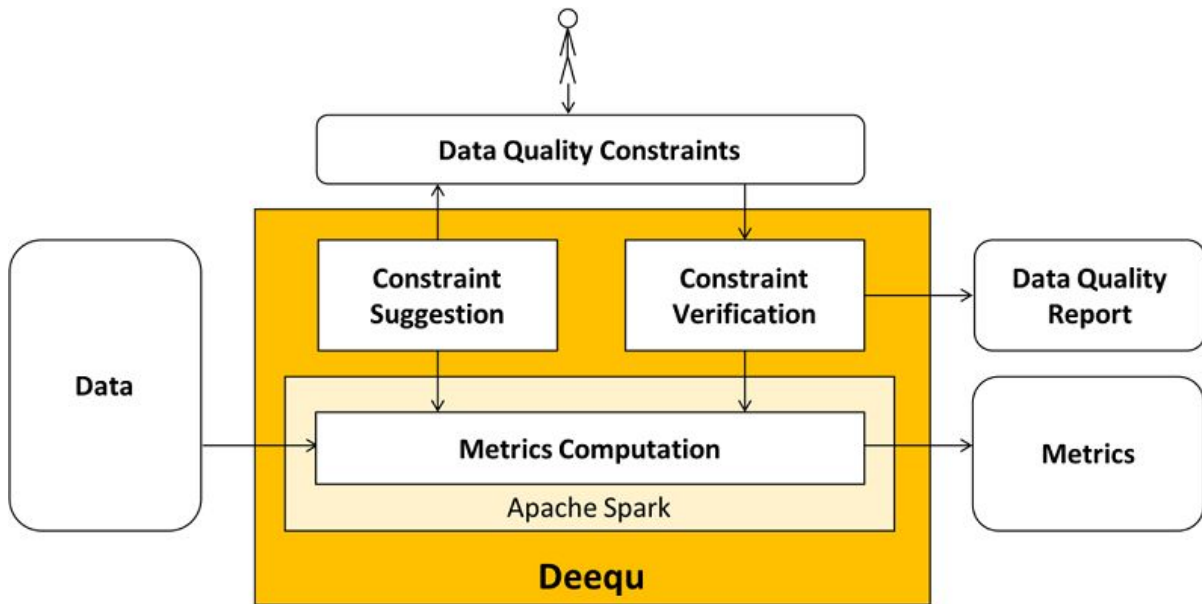
Ася Гайламазян



# Практический кейс очистки данных



# Amazon Deequ



+

Deequ специально разработан для оценки и проверки качества данных.

Хорошо интегрируется с apache spark для масштабируемой обработки данных.

-

Недостаток функций подготовки данных/

# Talend

The screenshot displays the Talend Studio interface. On the left, the 'Repository' pane shows a project structure with 'Business Models', 'Job Designs', and 'demo 0.1'. The main workspace shows a job design with a 'tRowGenerator\_1' component connected to a 'TutorialFamilyLogger\_1' component. The execution log at the bottom shows the job starting at 15:05 on 10/08/2018, connecting to a socket on port 3447, and logging 10 rows of data. The log includes the following JSON entries:

```
[INFO] {"firstName": "Jimmy", "lastName": "Washington"}
[INFO] {"firstName": "Theodore", "lastName": "Truman"}
[INFO] {"firstName": "Harry", "lastName": "Lincoln"}
[INFO] {"firstName": "Rutherford", "lastName": "Arthur"}
[INFO] {"firstName": "Theodore", "lastName": "Monroe"}
[INFO] {"firstName": "Franklin", "lastName": "Arthur"}
[INFO] {"firstName": "Williem", "lastName": "Adams"}
[INFO] {"firstName": "Millard", "lastName": "Kennedy"}
[INFO] {"firstName": "John", "lastName": "Garfield"}
[INFO] {"firstName": "Rutherford", "lastName": "Roosevelt"}
[statistics] disconnected
```

The job ended at 15:05 on 10/08/2018 with an exit code of 0. The log also shows statistics: '10 rows in 1.9s' and '5.27 rows/s'.



Возможности профилирования, очистки, стандартизации и проверки данных.



Для использования полного набора функций может потребоваться приобретение платной корпоративной версии.

# OpenRefine



## OpenRefine

← → C [api.us.socrata.com/api/catalog/v1/domains](https://api.us.socrata.com/api/catalog/v1/domains)

```
{
  "results": [
    { "domain": "201@bonds.cityofre.org", "count": 4 },
    { "domain": "anopen.smc.ca.ca", "count": 59 },
    { "domain": "bchi.bigcityhealth.org", "count": 84 },
    { "domain": "bea.data.commerce.gov", "count": 3 },
    { "domain": "bis.data.commerce.gov", "count": 2 },
    { "domain": "brigades.opendatnetwork.com", "count": 493 },
    { "domain": "brors.lehman.cuny.edu", "count": 750 },
    { "domain": "bythersumbers.sco.ca.gov", "count": 111 },
    { "domain": "capitalprojects.seattle.gov", "count": 3 },
    { "domain": "census.data.commerce.gov", "count": 212 },
    { "domain": "chh.data.ca.gov", "count": 438 },
    { "domain": "chronicdata.cdc.gov", "count": 379 },
    { "domain": "cip.cityofnovi.org", "count": 6 },
    { "domain": "controllerdata.lacity.org", "count": 1038 },
    { "domain": "dashboard.edmonton.ca", "count": 282 },
    { "domain": "dashboard.hawaii.gov", "count": 942 },
    { "domain": "dashboard.plano.gov", "count": 112 },
    { "domain": "dashboard.slo.org", "count": 46 },
    { "domain": "data.smcgov.org", "count": 286 },
    { "domain": "data.albany.gov", "count": 17 },
    { "domain": "data.atf.gov", "count": 135 },
    { "domain": "data.auburnva.gov", "count": 32 },
    { "domain": "data.austintexas.gov", "count": 1441 },
    { "domain": "data.arconet.org", "count": 99 },
    { "domain": "data.baltimorecity.gov", "count": 2026 },
    { "domain": "data.bria.gov", "count": 127 },
    { "domain": "data.burlingtonvt.gov", "count": 76 },
    { "domain": "data.calearv.ca", "count": 286 }
  ]
}
```



233 records

Show as: rows records Show: 5 10 25 50 records Sort ▾

All	data_portals	region	item_count
98.	data.nasa.gov	Nasa	31492
38.	data.cityofnewyork.us	Cityofnewyork	9824
212.	reports.data.montgomerycountymd.gov	Montgomerycountymd	5129
76.	data.kcmo.org	Kcmo	4900
224.	www.datos.gov.co	Datos	4588
181.	opendata.utah.gov	Utah	3341
129.	data.seattle.gov	Seattle	2822
21.	data.baltimorecity.gov	Baltimorecity	2026
111.	data.oregon.gov	Oregon	1986
88.	data.medicaid.gov	Medicaid	1756
202.	performance.smcgov.org	Smcgov	1716
86.	data.maryland.gov	Maryland	1559
150.	data.wa.gov	Wa	1490
19.	data.austintexas.gov	Austintexas	1441
69.	data.hawaii.gov	Hawaii	1433
124.	data.results.wa.gov	Wa	1278
44.	data.colorado.gov	Colorado	1157
74.	data.iowa.gov	Iowa	1113

+

OpenRefine подходит для задач очистки и преобразования данных. Возможности манипулирования над текстом и строками.

-

Нет расширенных возможностей проверки данных и обнаружения аномалий.

# Cucumber

<https://cucumber.io/>

cucumber 



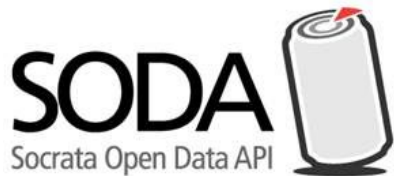
Подходит для проведения тестов ПО.



Может не обеспечивать всю глубину оценки качества данных.



# Soda



Soda для оценки и проверки качества данных, особенно хорошо работает с открытыми API.



Он может быть более специализированным для конкретных случаев.



# Great Expectations



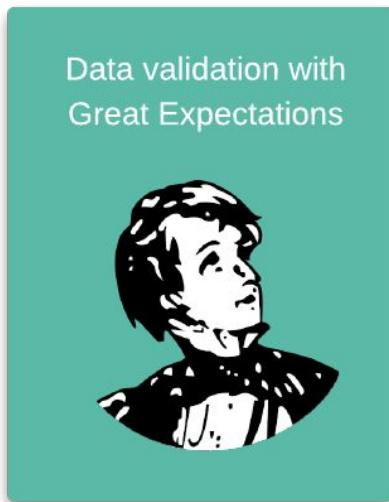
Great Expectations специальный инструмент для обеспечения качества данных.



Может не предлагать полный набор возможностей интеграции данных.



Your data assets:  
database tables, flat  
files, dataframes...



High quality data in  
your data products



Data documentation  
& data quality reports



Logging & alerting

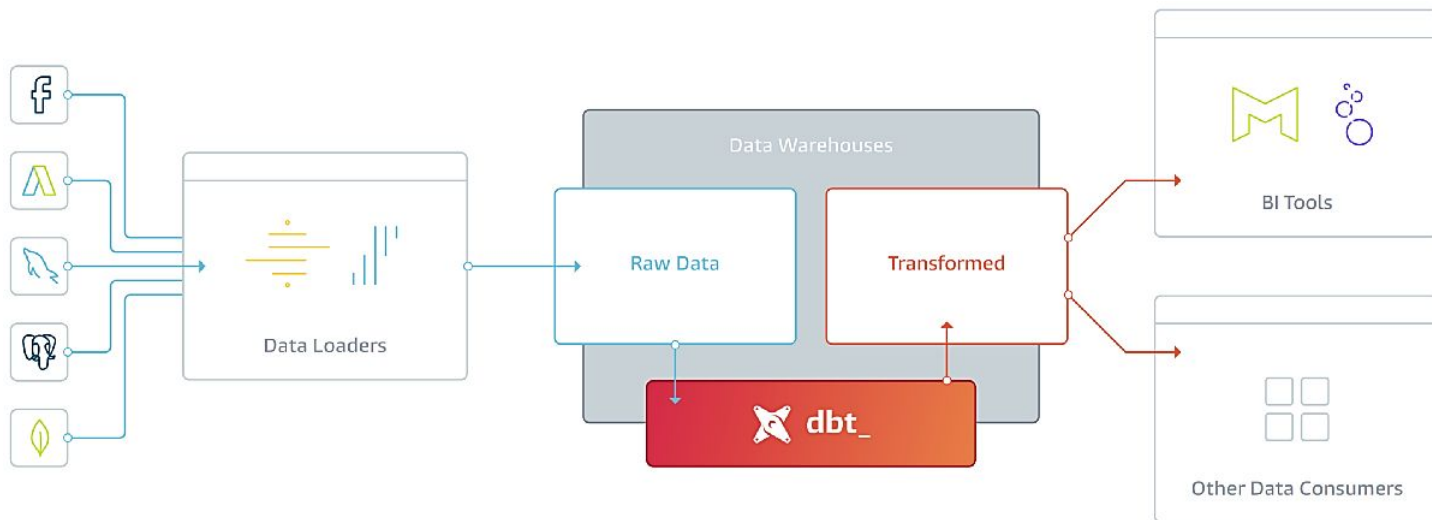
# dbt



dbt позволяет проверять и тестировать данные с помощью команды dbt test.



Может не охватывать весь спектр функций качества данных.



# Выбор инструмента

Figure 1. Magic Quadrant for Data Quality Tools

