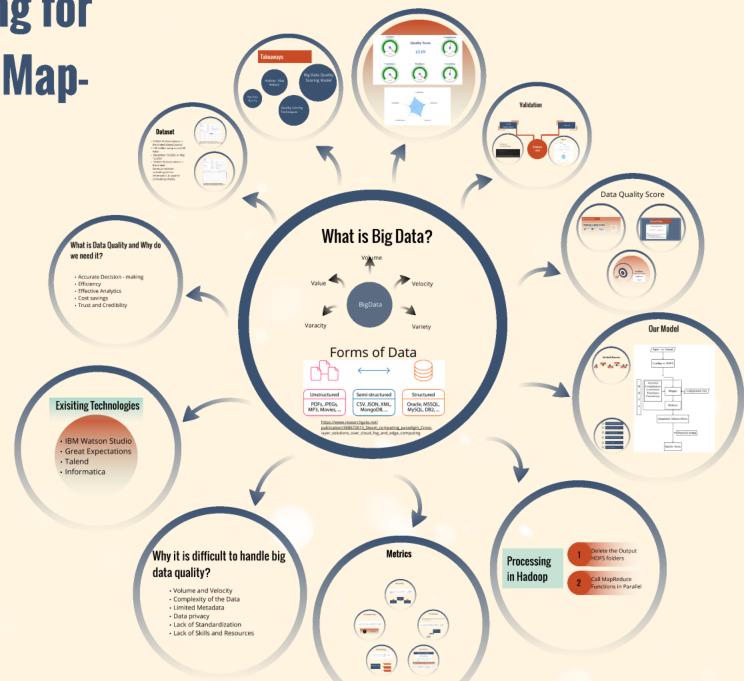
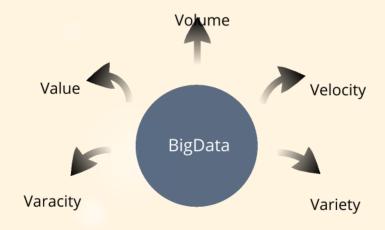
Big Data Quality Scoring for Structured Data Using Map-Reduce



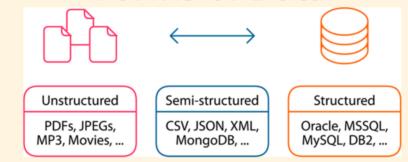


PERCONIASE OF INDIFECTIVE ROWS

What is Big Data?



Forms of Data



https://www.researchgate.net/ publication/368672613_Steam_computing_paradigm_Crosslayer_solutions_over_cloud_fog_and_edge_computing



nd Why do

king



What is Data Quality and Why do we need it?

- Accurate Decision making
- Efficiency
- Effective Analytics
- Cost savings
- Trust and Credibility

Why it is difficult to handle big data quality?

- Volume and Velocity
- Complexity of the Data
- Limited Metadata
- Data privacy
- Lack of Standardization
- Lack of Skills and Resources



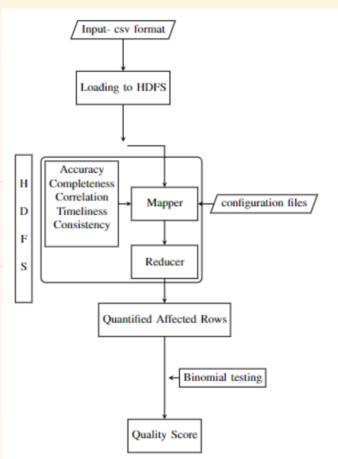
Exisiting Technologies

- IBM Watson Studio
- Great Expectations
- Talend
- Informatica

Our Model











Dataset

- 'COVID-19 Vaccinations in the United States, County'.
- 1.96 million samples and 80 fields
- December 13,2020, to May 12,2023
- 'COVID-19 Vaccinations in the United States, Jurisdiction', containing similar information, is used for consistency checks.



https://data.cdc.gov/Vaccinations/COVID-19-Vaccinations-in-the-United-States-County/8xkx-ami



https://data.cdc.gov/Vaccinations/COVID-19-Vaccinations-in-the-United-States-Jurisdi/unsk-h7

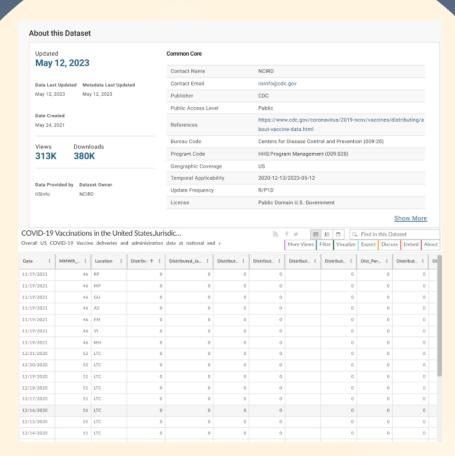
| Updated May 1 | | 0023 | | | Common Co | ire | | | | | | | |
|--|--|---------------|--|----------------------|--------------------------------------|--------------------------------------|---|---|------------------|------------------------------|--------------------|------------------------------|--------------------|
| ividy i | 12, 2 | .025 | | | Contact Na | ame | | NCIRD | | | | | |
| Data Last | Update | ed Metadata L | ast Updated | | Contact En | nail | | iisinfo@cd | c.gov | | | | |
| May 12, 20 | 023 | May 12, 20 | 23 | | Publisher | | | CDC | | | | | |
| | | | | | Public Acc | ess Level | | Public | | | | | |
| Date Created May 24, 2021 | | | | References | | | https://www.cdc.gov/coronavirus/2019-ncov/vaccines/distributing/about-vaccine-data.html | | | | | | |
| Views | | Downloads | | | Bureau Co | de | | Centers fo | r Disease Co | ontrol and Pr | evention (00 | 09:20) | |
| | | 211K | | | Program C | ode | | HHS:Progr | am Manage | ment (009:0 | 20) | | |
| 458K | | | | | | | | | | | | | |
| | | ecinations i | n the Unite | d State | Geographic s,County | c Coverage | | US | | | | | |
| OVID-19 | | | n the Unite | ed State | | Admi ≡ | Admi | Admi | Admi | Admi | Admi | Admi | Admi |
| OVID-19 | 9 Vac | cinations i | | Recip | s,County | Admi ≡ | | Admi | Admi 73.9 | Admi 10,863 | Admi 78.5 | Admi 10,368 | |
| OVID-19 | 9 Vac FIPS 55129 | ecinations i | Recip | Recip | s,County | Admi | dministered_0 | Admi Xose1_Recip | | | | | 81. |
| OVID-19 | 9 Vac FIPS 55129 | MMW | Recip Washburn C | Recip WI | comp | Admi = A | dministered_0 70.8 | Admi Jose1_Recip 11,097 | 73.9 | 10,863 | 78.5 | 10,368 | Admi 81. 63. |
| OVID-19 e Fi 10/2023 5 10/2023 1 10/2023 3 | 9 Vac | MMW 19 | Recip Washburn C Taylor Coun | Recip WI IA NY | Comp 96.7 | Admi = 11,123 | dministered_0 70.8 51.4 | Admi lose1_Recip 11,097 3,145 | 73.9 55 | 10,863 3,079 | 78.5 59.8 | 10,368 2,966 | 81. 63. |
| OVID-19 B FI 10/2023 5 10/2023 1 10/2023 3 | 9 Vac | MMW 19 | Recip Washburn C Taylor Coun | Recip WI IA NY TX | 96.7 97.3 | Admi = 11,123 3,149 1,391,226 | 70.8 51.4 95 | Admi lose1_Recip 11,097 3,145 1,384,503 | 73.9 55 95 | 10,863 3,079 1,329,779 | 78.5 59.8 95 | 10,368 2,966 1,232,671 | 81. 63. |
| OVID-19 e Fi 10/2023 5 | 9 Vac FIPS 55129 19173 36059 48281 26145 | MMW 19 19 | Recip Washburn C Taylor Coun Nassau Cou Lampasas (| Recip WI IA NY TX MI | Comp 96.7 97.3 97.5 98.9 | Admi = 11,123 3,149 1,391,226 11,678 | 70.8 51.4 95 54.5 | Admi lose1_Recip 11,097 3,145 1,384,503 | 73.9 55 95 | 10,863 3,079 1,329,779 | 78.5 59.8 95 | 10,368 2,966 1,232,671 | 81. 63. |

https://data.cdc.gov/Vaccinations/COVID-19-Vaccinations-in-the-United-States-County/8xkx-amqh

о Мау

ns ir

or



https://data.cdc.gov/Vaccinations/COVID-19-Vaccinations-in-the-United-States-Jurisdi/unsk-b7fc

andle big

Metrics

Processing in Hadoop







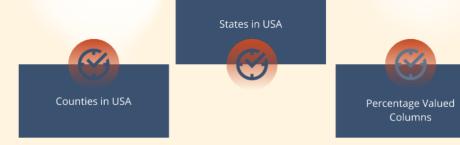




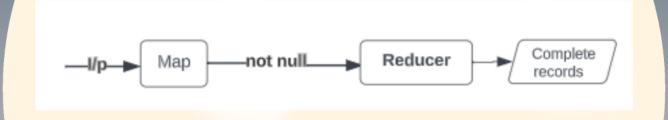
ces

Accuracy



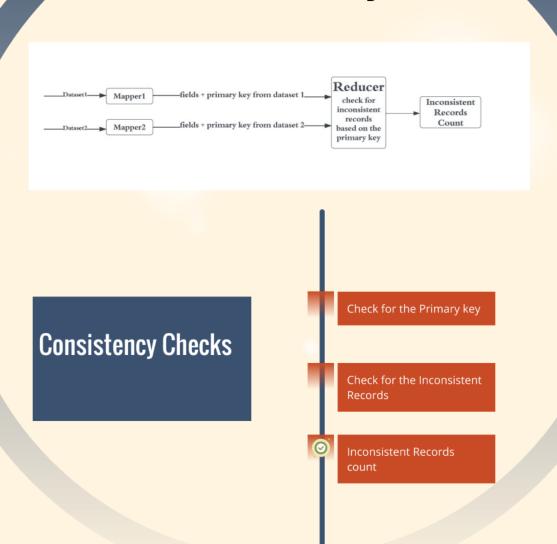


Completeness



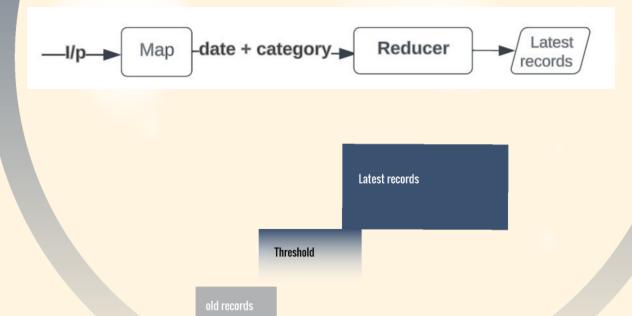


Consistency



Timeliness

- Defined as how relative the data is to the cureent date
- check the date column



Correlation

Pearsons Correlation

$$r = \frac{\sum_{i=1}^{n} (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^{n} (x_i - \bar{x})^2} \sqrt{\sum_{i=1}^{n} (y_i - \bar{y})^2}}$$
(6)

- x_i and y_i are individual data points for variables x and y.
- x̄ and ȳ are the means (averages) of variables x and y, respectively.

Calculated between every pair of columns



Processing in Hadoop

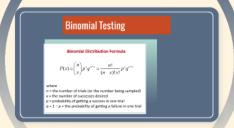
Delete the Output HDFS folders

2 Call MapReduce Functions in Parallel



Data Quality Score







Individual Feature Score

PERCENTAGE OF INEFFECTIVE ROWS

Calculating the percentage of ineffective rows in the dataset.











PERCENTAGE OF INEFFECTIVE ROWS

Calculating the percentage of ineffective rows in the dataset.





PERCENTAGE CALCULATION

 $\frac{\#ineffective_rows}{total_rows}$

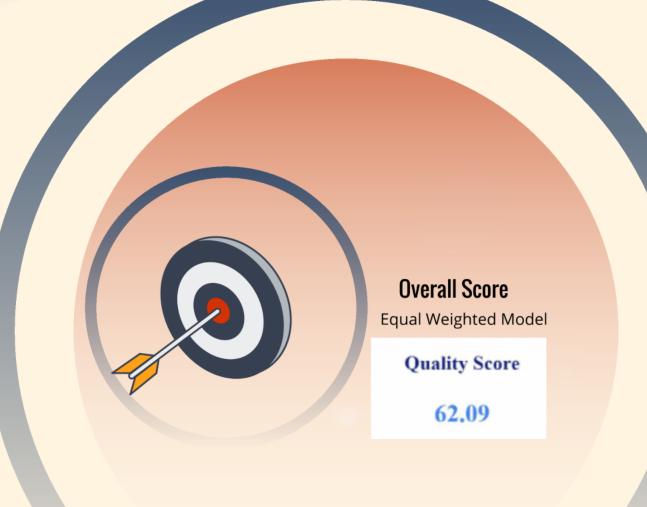
Binomial Testing

Binomial Distribution Formula

$$P(x) = \binom{n}{x} p^{x} q^{n-x} = \frac{n!}{(n-x)! \, x!} p^{x} q^{n-x}$$

where

- *n* = the number of trials (or the number being sampled)
- x = the number of successes desired
- p = probability of getting a success in one trial
- q = 1 p = the probability of getting a failure in one trial



Big Data Quality
Scoring Model

ng

Quality Score

62.09

Completeness

62.09

Consistency

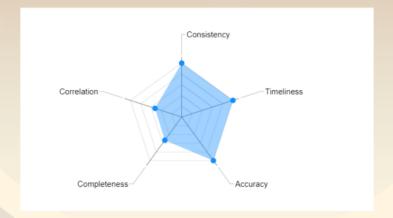
Timeliness

Correlation

100.00

100.00

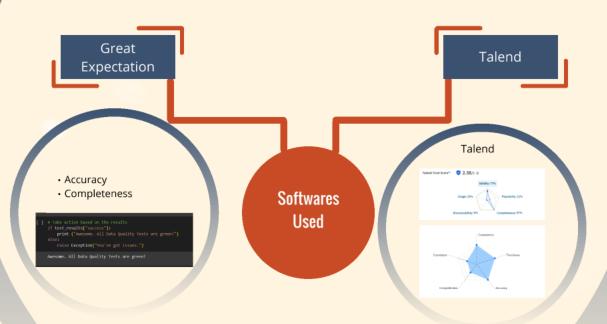
52.00



Validatio



Validation



- Accuracy
- Completeness

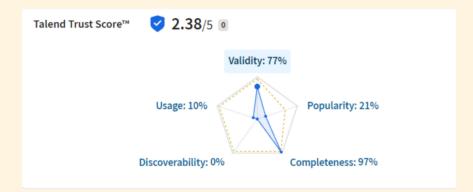
```
[ ] # Take action based on the results
   if test_results["success"]:
        print ("Awesome. All Data Quality Tests are green!")
   else:
        raise Exception("You've got issues.")

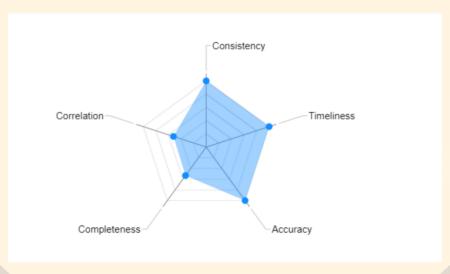
Awesome. All Data Quality Tests are green!
```

Sof ı

twares Jsed

Talend







Hadoop : Map Reduce Big Data Quality
Scoring Model

Big Data Quality

> Quality Scoring Techniques

| 16ey 12, | | | Communic | | | | | | | | |
|----------------------------|-----------------|---|--|-----------|--|-----------------------------|---------|---------------|-------------------|---------------|------------|
| 10ey 12, | 2023 | | Exercise Fig. | | | 50.00 | | | | | |
| No or on | or water | and the desired | | 100 | | | 411 | | | | |
| By 1.00 | We 1.00 | | Publisher | | | 125 | | | | | |
| | | | Palabara | resident. | | Patricia | | | | | |
| Date County Day 19, 200 | | | National | | | https://www. bestfransis | | | (B) (b) equal (b) | mandring (18) | elumpo |
| View | Description | | Service to | | | Design of the | THREE T | | 1400.0 | 177 | |
| 4596 | 211K | | Program (| look | | mtPops | - Name | mare (1004) | | | |
| | | | | | | | | | | | |
| ORVE-1918 | | the United St | Tal-project | townspe | | | | | | | |
| ton 1910 | accinations is | holo, he | enclosery from | | and the same of th | | ere. | Ade. | ant. | ne. | |
| 100 FFE | accinations in | | enclosery | | 76.0 | | | 100. -0.00 | 361 | | e e e |
| 14 FF | accimations in | Braker I M | tree drawing to the control of the c | | 711 | - | 194 | 410 | 31 | 10.00 | 80 |
| | ecolosotions in | Andre Ser Medical M Tylenian S | tree drawing to the control of the c | | 711 | 1.10 | 194 | 104 | 31 | 10.00 | #12 #12 |
| 140 FF | ecolostions is | Brain, Pari Brainne i St Trybrinan is Brainne i St | transference (marry) | 100 | 91 | 1.10 | 104 | 104 | 31 | 100 | |

https://data.cdc.gov/Vaccinations/COVID-19-

Big Data Quality Scoring for Structured Data Using Map-Reduce

