

**SOEN 6611**  
**SOFTWARE MEASUREMENT**  
**INSTRUCTOR: DR. OLGA ORMANDJIEVA**

**PROJECT**  
**TASK 3**

Source: SEI Implementing Goal-Driven Measurement course material (adapted).

SUBMITTED ON : 21-October-2022

**Declaration:**

**We, the members of the team, have read and understood the Fairness Protocol and the Communal Work Protocol, and agree to abide by the policies therein, without any exception, under any circumstances, whatsoever.**

**TEAM-7**

SIDDHARTHA NANDA	40200496
BARIQ ISHTIAQ MOHAMMED	40208194
RAJAT KUMAR	40201807
VIKYATH SRINIVISAULU	40218245

**TEAM LEADER**


SIDDHARTHA NANDA	40200496
------------------	----------

Email ID:	siddhartha.nanda@mail.concordia.ca
-----------	------------------------------------

### 3. Step 3:

#### 3.1. Success Criteria and Indicators for 3 V's:

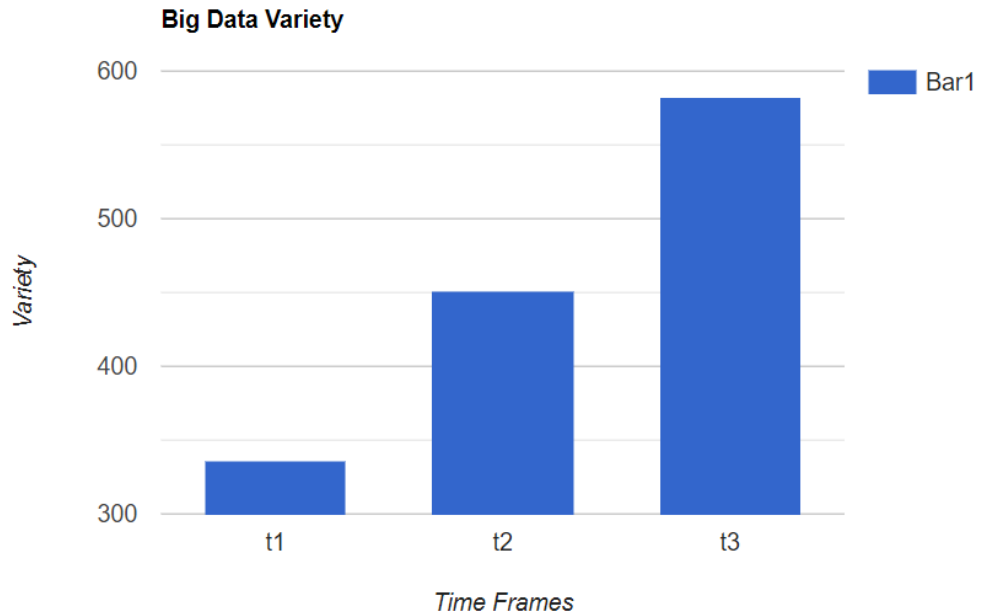
##### 3.1.1. Volume:

Measurement Question Label / Operationalized Goal Label	MG1 - Volume  Select the dataset with high-quality Volume.																
Success Criteria Label and description	The success criteria for volume are:  The percentage of data that can be preprocessed increases with the volume of the dataset and decreases with the decrease in volume.																
Indicator Label and description	<I1> Mvol:  The quantity of information bits in all the records determines the information content of numerous datasets.																
Indicator Analysis Model and Interpretation	<b>Indicator Analysis:</b>  The volume is divided for each time frame because it is the only source of data that is available. In each time frame, the volume passes through three phases:  Extraction, preprocessing, and processing  Interpretation:  We compare the difference between each phase of volume for each time frame.																
Indicator Sketch	 <table><caption>Big Data Volume</caption><thead><tr><th>Time Frames</th><th>Extraction</th><th>Preprocessing</th><th>Processing</th></tr></thead><tbody><tr><td>T1</td><td>12000</td><td>10500</td><td>7800</td></tr><tr><td>T2</td><td>12000</td><td>10500</td><td>7800</td></tr><tr><td>T3</td><td>11500</td><td>10000</td><td>7500</td></tr></tbody></table>	Time Frames	Extraction	Preprocessing	Processing	T1	12000	10500	7800	T2	12000	10500	7800	T3	11500	10000	7500
Time Frames	Extraction	Preprocessing	Processing														
T1	12000	10500	7800														
T2	12000	10500	7800														
T3	11500	10000	7500														

### 3.1.2. Velocity:

<b>Measurement Question Label / Operationalized Goal Label</b>	<p>MG1 - Velocity</p> <p>Improve the decision-making process by obtaining data more frequently and processing it more quickly.</p>																		
<b>Success Criteria Label and description</b>	<p>The velocity success criteria are as follows:</p> <p>An excessive discrepancy in the velocity of the dataset between time frames shows that the dataset is out of date.</p> <p>The quantity of useful data that can be processed grows in proportion to the rise in volume over time, and vice versa.</p>																		
<b>Indicator Label and description</b>	<p>&lt;I2&gt; Mvel:</p> <p>The rate at which big data volume grows over time (T)</p>																		
<b>Indicator Analysis Model and Interpretation</b>	<p><b>Indicator Analysis:</b></p> <p>After sending the dataset through three processes in each time frame, the quantity of information acquired is analysed:</p> <p>Extraction, preprocessing, and processing</p> <p>Interpretation:</p> <p>Compare the percentage variations in the number of insights received throughout time periods.</p>																		
<b>Indicator Sketch</b>	<p><b>Big Data Velocity</b></p> <table border="1"> <caption>Big Data Velocity Data</caption> <thead> <tr> <th>Time Frame</th> <th>Velocity</th> </tr> </thead> <tbody> <tr> <td>ExtT1</td> <td>35</td> </tr> <tr> <td>PreT1</td> <td>15</td> </tr> <tr> <td>ProT1</td> <td>37</td> </tr> <tr> <td>ExtT2</td> <td>35</td> </tr> <tr> <td>PreT2</td> <td>15</td> </tr> <tr> <td>ProT2</td> <td>37</td> </tr> <tr> <td>ExtT3</td> <td>34</td> </tr> <tr> <td>PreT3</td> <td>14</td> </tr> </tbody> </table>	Time Frame	Velocity	ExtT1	35	PreT1	15	ProT1	37	ExtT2	35	PreT2	15	ProT2	37	ExtT3	34	PreT3	14
Time Frame	Velocity																		
ExtT1	35																		
PreT1	15																		
ProT1	37																		
ExtT2	35																		
PreT2	15																		
ProT2	37																		
ExtT3	34																		
PreT3	14																		

### 3.1.3. Variety:

<b>Measurement Question Label / Operationalized Goal Label</b>	<p>MG1 - Variety</p> <p>Variety's goal is to categorise and segment data in order to classify and separate it.</p>								
<b>Success Criteria Label and description</b>	<p>The success criteria for variety is:</p> <p>When the amount of data, number of records, and number of datasets rise in comparison to previously utilised datasets</p>								
<b>Indicator Label and description</b>	<p>&lt;I1&gt; Mvar:</p> <p>A Mvar (MDS) is a collection of three values (Ndde, Lbd, and Nds) aggregated into a single number to represent the variety of unique data objects, records, and datasets in a specific MDS.</p>								
<b>Indicator Analysis Model and Interpretation</b>	<p><b>Indicator Analysis:</b></p> <p>The dataset is separated into time frames to analyse the variation across each period.</p> <p><b>Interpretation:</b></p> <p>The trend for structured and unstructured data will be evaluated by comparing the diversity in each dataset throughout different time periods.</p>								
<b>Indicator Sketch</b>	<p><b>Big Data Variety</b></p>  <table border="1"> <caption>Big Data Variety Data</caption> <thead> <tr> <th>Time Frame</th> <th>Variety (Bar1)</th> </tr> </thead> <tbody> <tr> <td>t1</td> <td>~330</td> </tr> <tr> <td>t2</td> <td>~450</td> </tr> <tr> <td>t3</td> <td>~580</td> </tr> </tbody> </table>	Time Frame	Variety (Bar1)	t1	~330	t2	~450	t3	~580
Time Frame	Variety (Bar1)								
t1	~330								
t2	~450								
t3	~580								

## 3.2. Measures and Operationalized Goals

### 3.2.1. Identification of the 3 V's measures:

#### a. Volume

<b>Indicator level</b> <b>I1</b>	<b>Indicators</b> <b>Mvol</b>	<b>Formula</b> $Mvol(MDS) = Ndde(MDS) * \log_2((Ndde(NDS)))$
-------------------------------------	----------------------------------	---

#### b. Velocity

<b>Indicator level</b> <b>I2</b>	<b>Indicators</b> <b>Mvel</b>	<b>Formula</b> $Mvel(MDS) = ((Mvol(MDS_{T2}) - Mvol(MDS_{T1})) / Mvol(MDS_{T1})) * 100$
-------------------------------------	----------------------------------	--

#### c. Variety

<b>Indicator level</b> <b>I3</b>	<b>Indicators</b> <b>Mvar</b>	<b>Formula</b> $Mvar(MDS) = Ndde(DE) * W_{Ndde} + Lbd(MDS) * W_{Lbd} + Nds(MDS) + W_{Nds}$
-------------------------------------	----------------------------------	---

Measures					Indicators Label		
S.No.	Identification (Name of the measure)	Type	Availability	Source*	<I1>	<I2>	<I3>
1	Ndde - Number of Distinct Data Elements	Base	C	Dataset	X		X
2	Lbd: Length of Big Data (Number of Records)	Base	A	Dataset	X		X
3	Nds: Number of Datasets	Base	A	Dataset			X
4	Time	Base	A	Dataset		X	
5	Mvol	Derived	B	Dataset	X		
6	Mvel	Derived	B	Dataset		X	

7	Mvar	Derived	B	Dataset			X
---	------	---------	---	---------	--	--	---

**Where:****Type:**

"Derived" or "Base".

**Availability:**

"A": Already available and collected.

"B": Can be derived from other data fairly directly.

"C": Possibly obtained with minor effort.

"D": Not available at the moment.

"E": Very difficult, if not impossible to obtain at the moment.

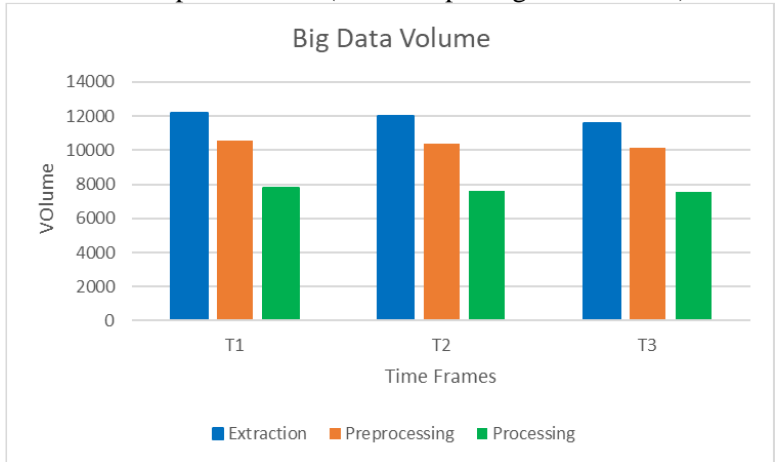
**Source:**

Place or tool where data is collected. In the case of base measures, this is obvious; in the case of derived measures, it depends on where the base data is stored after collection.

**Indicators:**

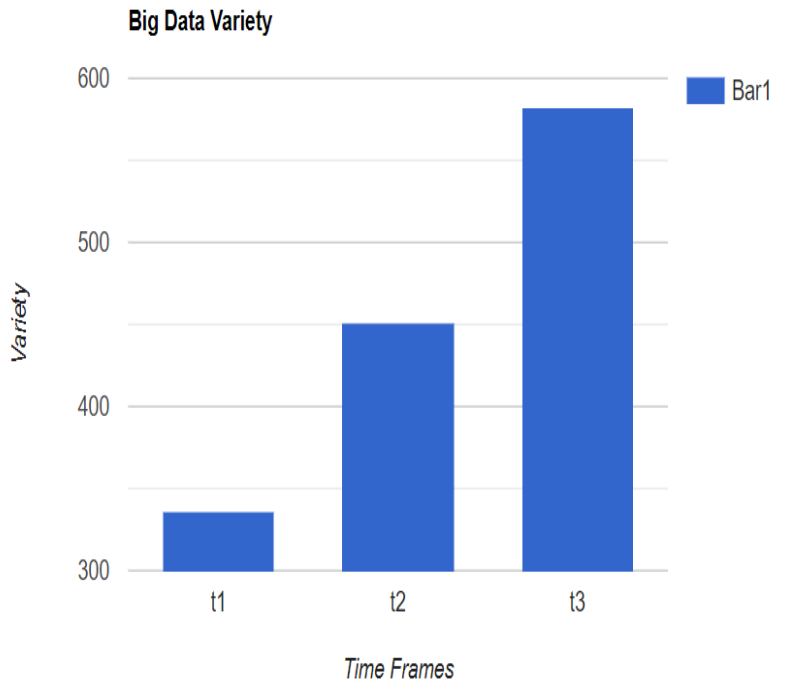
Mark an "X" when this measurement is required for each of your indicators.

### 3.2.2. Derived measures definitions and operationalization

Derived measure or indicator: <b>Volume</b>			
#	Derived measures or indicators <b>Mvol</b>	Formula:	
		$Mvol(MDS) = Ndde(MDS) * \log_2((Ndde(NDS)))$ <p>Where, <math>Ndde(MD)</math> = Number of Distinct data Element across MDS</p>	
Connect to the measurement target (which goal)  MG1 - Increasing the Volume of big data sets	Responsible:  <b>Data Scientist</b> in charge (who analyses)	Stakeholder:  <b>Product Owner and Developer</b> , they are the one who uses	Frequency (when) :  a. Before beginning the ML Algorithm development procedure (DAY 0). b. Before each update to the dataset. c. (DAY0 + N) WHERE N = the number of days when the dataset changes
Data source (where the measurement data will be extracted from)  <b>Most popular superhero TV show</b>	Storage of outcome (where data will be stored after extraction)  <b>Any distributed file system or local storage</b>	Data interpretation rules  <b>The volume in each step (extraction, preprocessing, and processing) can rise or fall, suggesting a positive or negative connection, with the exception that we anticipate the volume to increase when adding new data and decrease when removing it. Our choice is closely tied to the data pipeline steps we do.</b>	
Analysis Procedure:  <b>Data is gathered at a certain time and location, as in the preceding formula (for example, a volume measurement may occur in parallel with the data extraction phase and the data preprocessing phase at the same time). The amount of data may be compared throughout pipeline development or between time frames at the same level.</b>  <b>With respect to dataset in analysis:</b>  $Mvol(MDS) = Ndde(MDS) * \log_2(Ndde(MDS))$  <b>For T1</b> $mVolT1Ext = 12216$ $mVolT1Pre = 10537$ $mVolT1Pro = 7770$  <b>For T2</b> $mVolT2Ext = 11994$ $mVolT2Pre = 10377$		Results presentation (sketch depicting how it looks):  	

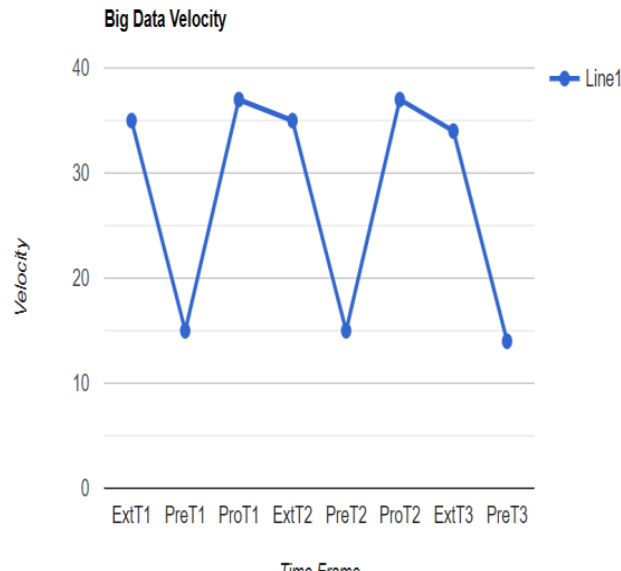
$mVolT2Pro = 7637$  <b>For T3</b> $mVolT2Ext = 11576$ $mVolT2Pre = 10125$ $mVolT2Pro = 7526$	
<p>Potential decisions based on the outcomes</p> <ol style="list-style-type: none"><li>1. Using this measurement, the information contents of several datasets may be compared to identify the real data that can be processed.</li><li>2. These measurements aid in decision-making for picking datasets for further processing and provide confidence in the dataset's volume.</li></ol>	



Derived measure or indicator: Variety			
#	Derived measures or indicators <b>Mvar</b>	Formula:	
		<b>Mvar (MDS) = Ndde (DE) * W<sub>Ndde</sub> + Lbd (MDS) * W<sub>Lbd</sub> + Nds(MDS) + W<sub>Nds</sub></b>  <b>WNdde : Weight of Ndde (Set to 1/3 by default)</b> <b>WLbd : Weight of Lbd (Set to 1/3 by default)</b> <b>WNds : Weight of Nds (Set to 1/3 by default)</b> <b>Sum of all weights is equal to 1</b>	
Connect to the measurement goal (which goal)  MG2 - Enhancing Variety in Big Data	Responsible:  <b>Data Scientist</b> in charge (who analyses)	Stakeholder:  <b>Developer and Tester</b> , they are the one who uses	Frequency (when) :  <b>The variety of the dataset is evaluated in each time frame to provide Ndde, Lbd, and Nds.</b>
Data source (where the measurement data will be extracted from)  <b>Most popular superhero TV show</b>	Storage of outcome(where data will be stored after extraction)  <b>Any distributed file system or local storage</b>	Data interpretation rules  <b>Mvel describes a volume change that occurs over time.</b>  <b>The quantity of Ndde, Lbd, and Nds in our dataset demonstrates variety. The change in variety values over time is not effective in judging whether or not the variety in a dataset is good. This allows us to see the quantity of information, records, and datasets available at a glance..</b>	
<b>Analysis Procedure:</b>  We use the above-mentioned algorithm with weights specified by the data practitioner.  All weights are set to 1/4 by default. If, for example, the data practitioner wants to prioritise accuracy, the weight might be increased, allowing them to observe changes in that specific metric more clearly.  NddeT1= 2694 NddeT2 = 2652 NddeT2 = 2056 lbdT1 = 750 lbdT2 = 729 lbdT3 = 549 Nds = 3  $mvarT1 = ((nddeT1)/3) + (lbdT1/3) + (Nds/2)$ $mvarT2 = ((nddeT2)/3) + (lbdT2/3) + (Nds/2)$ $mvarT3 = ((nddeT3)/3) + (lbdT3/3) + (Nds/2)$ mvarT1 = 336 mvarT2 = 451 mvarT3= 582		Results presentation (sketch depicting how it looks):  	

Potential decision making depending on the results

**It is possible to tell how the quantity of diversity in the dataset evolves over time by evaluating the graph. It evaluates changes in the volume of organised and unstructured data by continually monitoring.**

Derived measure or indicator: Velocity			
#	Derived measures or indicators <b>Mvel</b>	Formula:	
		$Mvel(MDS) = ((Mvol(MDS_{T2}) - Mvol(MDS_{T1})) / Mvol(MDS_{T1})) * 100$ <p>Where, MDST1 and MDST2 are the multiple datasets at time T1 and T2 respectively (where T2&gt;T1). Thus, Mvol (MDS) is defined in terms of volume growth over an interval of time (T2-T1) along with the appropriate unit of measure (seconds, minutes, hours, weeks, etc.).</p>	
Connect to the measurement goal (which goal)  MG3 - Increase Big Data Set Velocity	Responsible:  <b>Data Scientist</b> in charge (who analyses)	Stakeholder:  <b>Product Owner and Developer</b> , they are the one who uses	Frequency (when) :  <b>In each of the three processes - extraction, preprocessing, and processing.</b>
Data source (where the measurement data will be extracted from)  <b>Most popular superhero TV show</b>	Storage of outcome(where data will be stored after extraction)  <b>Any distributed file system or local storage</b>	Data interpretation rules  <b>Mvel describes a volume change that occurs over time.</b> <ol style="list-style-type: none"> <li>1. A positive plane in the graph shows that more meaningful information is being obtained.</li> <li>2. A negative plane in the graph implies that useful information is being lost.</li> <li>3. A straight plane in the graph shows that no information is gained or lost.</li> </ol>	
<b>Analysis Procedure:</b>  To compare the amount of big data over time, we must create a line graph that depicts the rate of change in Mvel across three successive time periods. A variation in inclination or descent between the stages of the time frames indicates that the change rate is not identical.  Values throughout various time periods and  The following categorization is used to define the set below.  {T1Ext,T1Pre,T1Pro,T2Ext,T2Pre,T2Pro,T3Ext,T3Pre,T3Pro} Ndde = {1195,1050,805,1176,1036,793,1140,1014,783} Mvol = {12216,10537,7770,11994,10377,7637,11576,10125,7526}  $Mvel(MDS) = \left( \frac{Mvol(MDS_{T2}) - Mvol(MDS_{T1})}{Mvol(MDS_{T1})} \right) * 100$  The values captured are as follows		Results presentation (sketch depicting how it looks):   <p>The graph shows a fluctuating line representing velocity over time. The Y-axis is labeled 'Velocity' and ranges from 0 to 40. The X-axis is labeled 'Time Frame' and includes categories: ExtT1, PreT1, ProT1, ExtT2, PreT2, ProT2, ExtT3, and PreT3. The line starts at approximately 35 for ExtT1, drops to 15 for PreT1, rises to 38 for ProT1, drops to 35 for ExtT2, rises to 38 for PreT2, drops to 35 for ProT2, rises to 38 for ExtT3, and drops to 15 for PreT3. A legend indicates 'Line1'.</p>	

{35 ,15 ,37, 35, 15 ,37 ,34 ,14}	
<p>Potential decision making depending on the results</p> <p><b>We may deduce the pace at which information is obtained or lost from the graph. By monitoring on a frequent basis, it prevents data from becoming obsolete by keeping a comparable velocity throughout a specified time period, and it detects the quantity of outdated data if the velocity is not similar in each time frame.</b></p>	

### 3.2.3. Base measures definitions and operationalization:

#	Base measure: Time			Applicability:
	Measure (what: entity, attribute)	Scale type :		
	Entity: Dataset Attribute: Time(T)	Ratio scale		Based on this value, the provided dataset volume is partitioned into three periods.
Who measures? <b>Developer/Data Scientist</b>	Source of measurement <b>Most Popular Superhero TV Shows</b> <a href="https://www.kaggle.com/anoopkumarraut/most-popular-superhero-tv-shows">https://www.kaggle.com/anoopkumarraut/most-popular-superhero-tv-shows</a>	Where to store the result <b>Local Storage or any distributed File System</b>	Tool <b>Google collab notebook</b>	Time (when to measure) <b>1) During the extraction-processing and processing of the dataset. 2) Repeat whenever the dataset changes.</b>
Collection procedure (how to collect the data)	Notes or comments:			
The amount of time taken to collect the data by downloading it from Kaggle website ( <a href="https://www.kaggle.com/anoopkumarraut/most-popular-superhero-tv-shows">https://www.kaggle.com/anoopkumarraut/most-popular-superhero-tv-shows</a> ).		This measure is generally used to calculate the velocity of data.		

Base measure: Ndde

#	<b>Ndde : number of distinct data elements</b> Measure (what: entity, attribute)  <b>Entity: Dataset</b> <b>Attribute: no of unique elements</b>	Scale type  <b>Absolute scale</b>	Applicability  <b>The number of different elements in the overall dataset is calculated.</b>
Who measures? <b>Developer/ Data Scientist</b>	Source of measurement  <b>Most Popular Superhero TV Shows</b> <a href="https://www.kaggle.com/anoopkumarraut/most-popular-superhero-tv-shows">https://www.kaggle.com/anoopkumarraut/most-popular-superhero-tv-shows</a>	Where to store the result  <b>Local Storage or any distributed File System</b>	Tool  <b>Google Colab</b>  Time (when to measure)  <b>After splitting the, distinct components are calculated. dataset volume based on different time periods</b>
Collection procedure (how to collect the data)		Notes or comments:	
<b>To locate the different components in the dataset, use the Python code in the Google collab online platform.</b>		<b>This measure is used to perform the calculation for variety and volume.</b>	

#	<b>Base measure: Lbd</b> <b>Lbd : Length of Big Data</b> Measure (what: entity, attribute) <b>Entity: Dataset</b> <b>Attribute: Size</b>	Scale type  <b>Absolute Scale</b>	Applicability  <b>It provides the real length of the dataset and may be used to assess dataset diversity.</b>
Who measures? <b>Developer/ Data Scientist</b>	Source of measurement  <b>Most Popular Superhero TV Shows</b> <a href="https://www.kaggle.com/anoopkumarraut/most-popular-superhero-tv-shows">https://www.kaggle.com/anoopkumarraut/most-popular-superhero-tv-shows</a>	Where to store the result  <b>Local Storage or any distributed File System</b>	Tool  <b>Google Colab</b>  Time (when to measure)  <b>During each time frame, the length of new dataset is calculated.</b>
Collection procedure (how to collect the data)		Notes or comments:	
<b>To locate the different components in the dataset, use the Python code in the Google collab online platform.</b>		<b>This measure is to calculate the variety.</b>	

#	<b>Base measure: Nds</b> <b>Nds : No of Dataset in Big Data</b>	Scale type  <b>Absolute scale</b>	Applicability
---	--	---	---------------

	Measure (what: entity, attribute)			<b>It counts the amount of datasets that are available to study the fluctuations in each metric across time.</b>
	<b>Entity: Data set</b> <b>Attribute: number of datasets</b>			
Who measures?	Source of measurement	Where to store the result	Tool	Time (when to measure)
<b>Developer/ Data Scientist</b>	<b>Most Popular Superhero TV Shows</b> <a href="https://www.kaggle.com/anoopkumarraut/most-popular-superhero-tv-shows">https://www.kaggle.com/anoopkumarraut/most-popular-superhero-tv-shows</a>	<b>Local Storage or any distributed File System</b>	<b>Google Colab</b>	<b>The number of fresh datasets is determined at each time period.</b>
Collection procedure (how to collect the data)		Notes or comments:		
<b>By dividing the full dataset and allocating it to multiple time frames using Python code in the Google collab online platform.</b>		<b>This measure is to calculate the variety.</b>		