

Derived and Base Measure for Validity

Derived measure or indicator: Validity				
#	Derived measure or indicator	Formula		
	Mval or Validity of Big Data is defined in terms of its accuracy and correctness for the purpose of usage.	$Mval (MDS) = Credability (MDS) * W_{Cred} + Compliance(MDS) * W_{Compli}$		
Link with the measurement goal (which goal)		Responsible (who analyzes)	Stakeholder (who uses)	Frequency (when)
Validity		Developer Data Analyst Data Engineer Data Scientist	Senior management Project manager Data scientist Data analyst	The validity of data set can be measured on monthly, quarterly, or yearly basis.
Data source (where the measurement data will be extracted from)		Storage of the result (where data will be stored after the extraction)	Data interpretation rules	
Credit Card classification - https://www.kaggle.com/datasets/samuelcortinhas/credit-card-classification-clean-data		The data will be stored in excel file or database. In our case we will be storing the result in jupyter notebook for reporting purpose.	Successful request is categorized as a request which returns the correct result. Every query to a database is considered as a request. Validity = 1 - means that the subject data is accurate and correct for the purpose of usage. This is a desired value for implementation of a successful machine learning model. Validity = 0 means that data attributes are not correct. Validity >= 0.90 means that 90% of the data attribute are accurate which can be	

		<p>useful to train our machine learning algorithm.</p> <p>Validity could increase or decrease depending upon the dataset size increasing or decreasing.</p>
Analysis procedure <ol style="list-style-type: none">1. Dataset is loaded using the analyses tool, excel file or jupyter notebook.2. Compliance and Credibility will be calculated using the formula.3. Validity of the dataset will be calculated using the formula.4. The value will be interpreted according to the decision-making rules and appropriate decision will be taken.		Presentation of the results (sketch illustrating what it looks like): Validity of the data will be presented as a single numerical value.
Potential decision making depending on the results Validity of the data attributes can give the overview about accuracy and correctness of the data. This is an important measure to get the Machine Learning model trained with the correct data. If the completeness value is more, it will give the confidence to stakeholders to trust the results produced by the machine learning algorithms.		

Derived measure or indicator: Compliance				
#	Derived measure or indicator	Formula		
	Degree to which data has attributes that adhere to standards, conventions or regulations in force and similar rules relating to data quality in a specific context of use	$Compliance (MDS) = \frac{\sum_{DS \in MDS} N_{rec_{comp}}(DS)}{N_{ds}(MDS)}$		
Link with the measurement goal (which goal)		Responsible (who analyzes)	Stakeholder (who uses)	Frequency (when)
Validity		Developer Data Analyst Data Engineer Data Scientist	Senior management Project manager Data scientist Data analyst	The compliance of data set can be measured on monthly, quarterly, or yearly basis.
Data source (where the measurement data will be extracted from)		Storage of the result (where data will be stored after the extraction)	Data interpretation rules	
Credit Card classification - https://www.kaggle.com/datasets/samuelcortinhas/credit-card-classification-clean-data		The data will be stored in excel file or database. In our case we will be storing the result in jupyter notebook for reporting purpose.	<p>Successful request is categorized as a request which returns the correct result.</p> <p>Every query to a database is considered as a request.</p> <p>Compliance = 1 - means that the subject data adheres to standards, conventions or regulations in force and similar rules relating to data quality for the purpose of usage. This is a desired value for implementation of a successful machine learning model.</p> <p>Compliance = 0 means that data attributes are not compliant.</p> <p>Compliance >= 0.90 means that 90% of the data attribute adhere to standards and</p>	

		<p>regulations which can be useful to train our machine learning algorithm.</p> <p>Compliance could increase or decrease depending upon the dataset size increasing or decreasing.</p>
Analysis procedure <ol style="list-style-type: none">1. Dataset is loaded using the analyses tool, excel file or jupyter notebook.2. rec_comp is counted using COUNT function to get number of compliant records in a dataset3. DS_comp is calculated by using the formula $\text{rec_comp}/\text{Nds}$4. Compliance of the dataset will be calculated using the formula.5. The value will be interpreted according to the decision-making rules and appropriate decision will be taken.		Presentation of the results (sketch illustrating what it looks like): <p>Compliance of the data will be presented as a single numerical value.</p>
Potential decision making depending on the results <p>Compliance of the data attributes can give the overview about adherence to standards, conventions or regulations in force and similar rules of the data. This is an important measure to get the Machine Learning model trained with the correct data. If the completeness value is more, it will give the confidence to stakeholders to trust the results produced by the machine learning algorithms.</p>		

Base measure: Compliant records (Nrec_comp)				
#1	Measure (what: entity, attribute) Measures the number of compliant records in the dataset Entity: Dataset Attribute: Number of records	Scale type Absolute	Applicability Total number of compliant records in data sets acts as a fundamental unit of measurement which can be used to calculate other derived measures.	
Who measures? Data Analyst Data Engineer Data Scientist	Source of measurement https://www.kaggle.com/samuelcortinhas/credit-card-classification-clean-data	Where to store the result CSV File Database	Tool Excel Jupyter Notebook Python libraries for data analysis like pandas, NumPy etc.	Time (when to measure) Compliant number of records can be measured each time new data is loaded into the database.
Collection procedure (how to collect the data) The data is loaded into excel sheet or database and the total number of compliant records can be retrieved from querying the database or using inbuilt functions of excel.		Notes or comments: None		

Base measure: Number of datasets (Nds)				
#2	Measure (what: entity, attribute) Measures the number of records in the dataset Entity: Dataset Attribute: Number of records	Scale type Absolute	Applicability Total number of records in data sets acts as a fundamental unit of measurement which can be used to calculate other derived measures. It also gives the idea about the size of the dataset.	
Who measures? Data Engineer Data Analyst Business Analyst	Source of measurement https://www.kaggle.com/samuelcortinhas/credit-card-classification-clean-data	Where to store the result CSV File Database	Tool Excel Jupyter Notebook Python libraries for data analysis like pandas, NumPy etc.	Time (when to measure) Number of records can be measured each time new data is loaded into the database.
Collection procedure (how to collect the data) This number should be given by the responsible person managing databases or excel files.		Notes or comments: None		

Derived measure or indicator: Credability				
#2	Derived measure or indicator	Formula		
	Degree to which data has attributes that are regarded as true and believable by users in a specific context of use	$\text{Credability (MDS)} = \frac{Nds_{cr}(MDS)}{Nds(MDS)}$		
Link with the measurement goal (which goal)		Responsible (who analyzes)	Stakeholder (who uses)	Frequency (when)
Validity		Developer Data Analyst Data Engineer Data Scientist	Senior management Project manager Data scientist Data analyst	The credibility of data set can be measured on monthly, quarterly, or yearly basis.
Data source (where the measurement data will be extracted from)		Storage of the result (where data will be stored after the extraction)	Data interpretation rules	
Credit Card classification - https://www.kaggle.com/datasets/samuelcortinhas/credit-card-classification-clean-data		The data will be stored in excel file or database. In our case we will be storing the result in jupyter notebook for reporting purpose.	Successful request is categorized as a request which returns the correct result. Every query to a database is considered as a request. Credability = 1 - means that the subject data is regarded as true and believable by users for the purpose of usage. This is a desired value for implementation of a successful machine learning model. Credability = 0 means that data attributes are not truthful. Credability >= 0.90 means that 90% of the data attribute are true and believable which can be useful to train our machine learning algorithm.	

		Credability could increase or decrease depending upon the dataset size increasing or decreasing.
Analysis procedure <ol style="list-style-type: none"> 1. Dataset is loaded using the analyses tool, excel file or jupyter notebook. 2. cre_source is counted using COUNT function to get number of credible records in a dataset 3. Credibility of the dataset will be calculated using the formula. 4. The value will be interpreted according to the decision-making rules and appropriate decision will be taken. 		Presentation of the results (sketch illustrating what it looks like): <p>Credibility of the data will be presented as a single numerical value.</p>
Potential decision making depending on the results <p>Credibility of the data attributes can give the overview about truthfulness and reliability of the data. This is an important measure to get the Machine Learning model trained with the correct data. If the completeness value is more, it will give the confidence to stakeholders to trust the results produced by the machine learning algorithms.</p>		

Base measure: Credible Datasets (Nds_cr)				
#1	Measure (what: entity, attribute) Measures the number of credible records in the dataset Entity: Dataset Attribute: Number of records	Scale type Absolute	Applicability Total number of credible records in data sets acts as a fundamental unit of measurement which can be used to calculate other derived measures.	
Who measures? Data Analyst Business Analyst	Source of measurement https://www.kaggle.com/samuelcortinhas/credit-card-classification-clean-data	Where to store the result CSV File Database	Tool Jupyter Notebook Python libraries for data analysis like pandas, NumPy etc.	Time (when to measure) Credible number of records can be measured each time new data is loaded into the database.
Collection procedure (how to collect the data) This number should be given by the responsible person managing databases or excel files.		Notes or comments: None		

Base measure: Number of datasets (Nds)				
#2	Measure (what: entity, attribute) Measures the number of records in the dataset Entity: Dataset Attribute: Number of records	Scale type Absolute	Applicability Total number of records in data sets acts as a fundamental unit of measurement which can be used to calculate other derived measures. It also gives the idea about the size of the dataset.	
Who measures? Data Engineer Data Analyst Business Analyst		Source of measurement https://www.kaggle.com/samuelcortinhas/credit-card-classification-clean-data	Where to store the result CSV File Database	Tool Excel Jupyter Notebook Python libraries for data analysis like pandas, NumPy etc.
Collection procedure (how to collect the data) This number should be given by the responsible person managing databases or excel files.		Notes or comments: None		