

1.1. Metodología

Desde el punto de vista metodológico, se ha desarrollado una metodología de trabajo común para todos los casos de uso, dado que todos ellos se ajustan a un desarrollo que puede enmarcarse bajo el paraguas de lo que de manera estandarizada se ha venido a denominar como la disciplina de la Ciencia de Datos o en su término anglosajón, más extendido, el Data Science, que está íntimamente ligado con la Estadística y la Inteligencia Artificial, y sus técnicas más conocidas, como las de Machine Learning, y dentro de estas las Redes Neuronales Artificiales, y la subclase de estas más popular, los modelos de Deep Learning, entre otras.

En este sentido, se ha definido un flujo de trabajo que cubre los siguientes dominios y procesos (se utiliza una terminología inglesa dado que muchos términos no se traducen directamente y es común trabajar en esta disciplina con esta terminología).



Data Science Workflow

I – Business Problem Domain

1. Domain Knowledge
2. Data-Driven Approach
3. Data Science Approach
4. Analytics Approach

II – Data Processing Domain

5. Data Collection
6. Data Adequacy
7. Data Sampling
8. Data Split (Train, Validation, Test)
9. Data Cleansing (Cleaning or Scrubbing)
10. Data Balancing Analysis
11. Exploratory Causal Analysis (ECA)
12. Exploratory Data Analysis (EDA)
13. Data Visualization

III – Features Engineering Domain

14. Feature Data Transform
15. Feature Importance
16. Feature Selection

17. Feature Extraction (Dimensional Reduction)
18. Feature Construction
19. Feature Transforms (Polynomial Features Transform)
20. Feature Learning

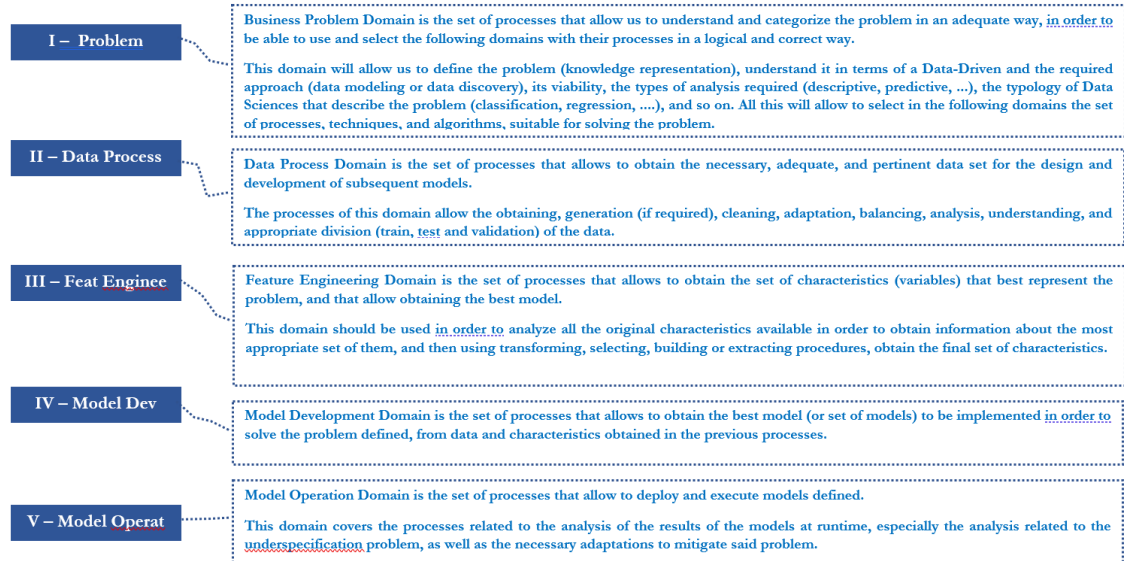
IV – Model Development Domain

21. Model Spot Checking
22. Model Evaluation
23. Model Selection
24. Model Tuning
25. Model Combination
26. Model Calibration
27. Model Uncertainty Analysis
28. Model Bias-Variance Trade-off Analysis
29. Model Interpretation
30. Model Prediction
31. Model Finalization
32. Model Saving

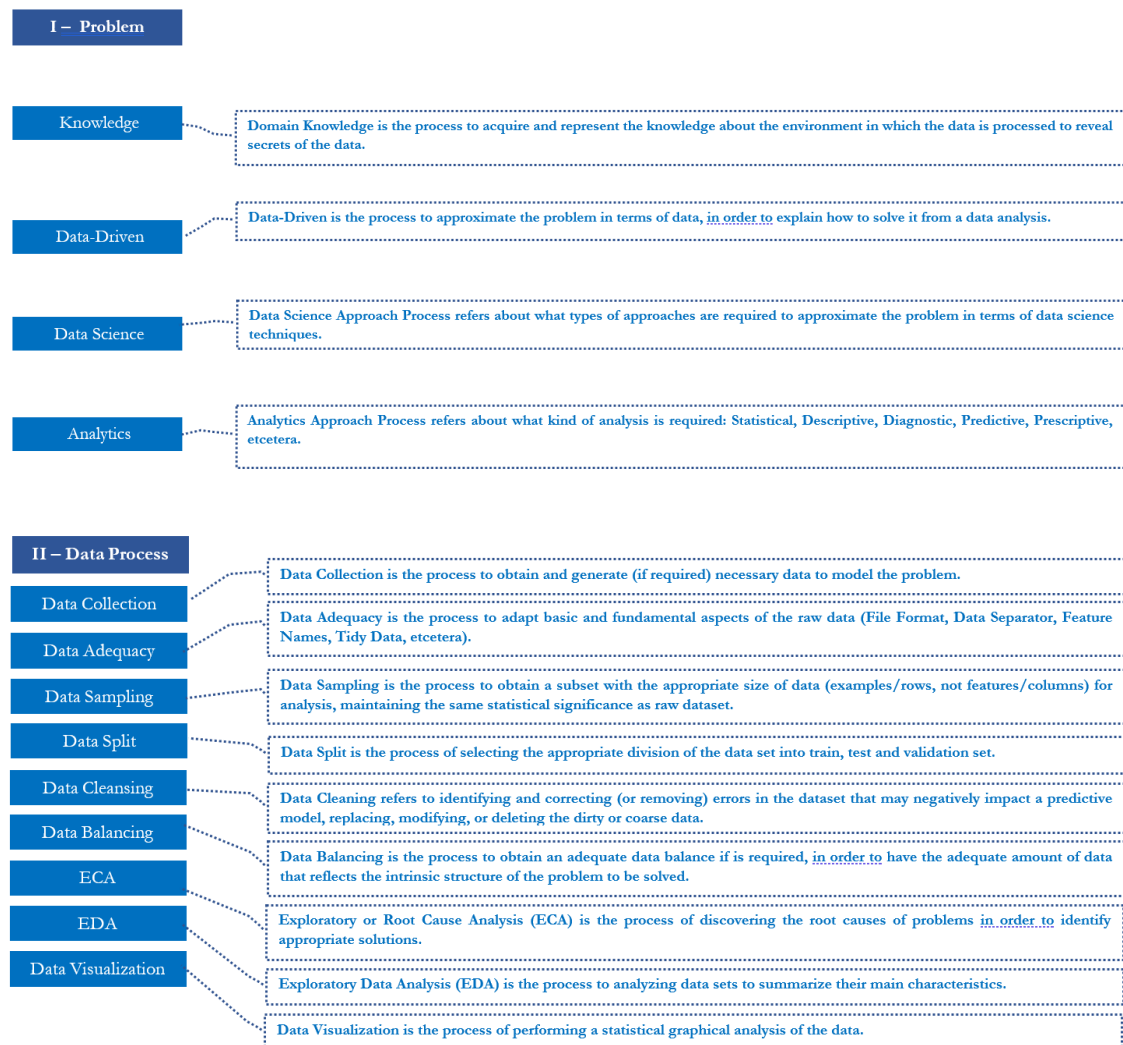
V – Model Operation Domain

33. Model Deployment
34. Model Execution
35. Model Analysis
36. Model Updating
37. Model Explainability

Domains Description



Processes Description



III – Feat Engine

Feat Data Transf	Feature Data Transform is the process that allows change (if is required) the type and/or distribution of data features (e.g. scaling, normalizing o standardizing data features).
Feat Importance	Feature Importance is the process that assigns scores to the input characteristics to a model, which indicate the relative importance of each characteristic, in order, for example, to be able to select the most important ones.
Feat Selection	Feature Selection is the process where you automatically or manually select the most relevant features which contribute most to the correct output of the model.
Feat Extraction	Feature Extraction is the process related to dimensionality reduction (or dimension reduction) that creates a projection of the data (high-dimensional space) resulting in entirely new input features (low-dimensional space), so that the low-dimensional representation retains some meaningful properties of the original data, ideally close to its intrinsic dimension.
Feat Construction	Feature Construction is the process related to create new features from your existing ones to improve model performance.
Feat Transform	Feature Transform (Polynomial Features Transform) is the process to create new features by raising existing features to an exponent, <u>in order to</u> see if they improve model performance, when the input features interact in unexpected and often nonlinear ways.
Feat Learning	Feature learning or representation learning is a set of techniques that allows a system to automatically discover the representations needed for feature detection or classification from raw data.

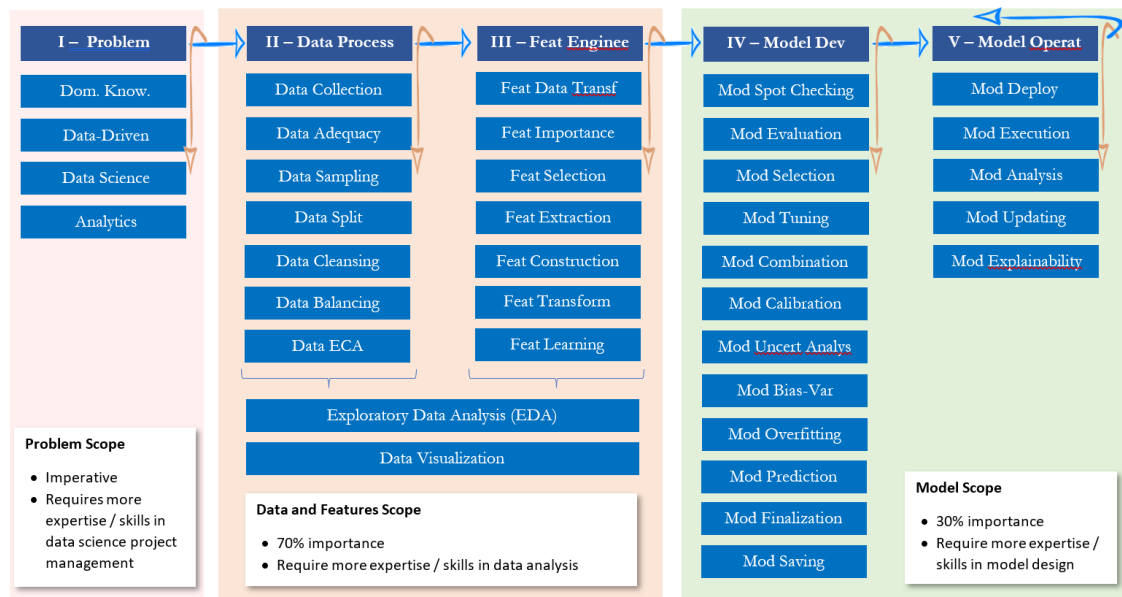
IV – Model Dev

Mod Spot Checking	Spot Checking process is a set of techniques that allows to have a first approximation to know which algorithm will perform best on your data beforehand using a process of trial and error.
Mod Evaluation	Model Evaluation is the process that allows us to evaluate the Models performance focused <u>in</u> the use of multiple Metrics and Resampling techniques.
Mod Selection	Model Selection process is a set of techniques that allows to choose one from among a set of models, using Statistical Hypothesis Tests to address if the difference in skill between machine learning models is real, or due to a statistical chance; and the other hand, Probabilistic Statistical Measures to quantify both the model performance on the training dataset and the complexity of the model.
Mod Tuning	Model Tuning process is a set of techniques that allows to <u>optimizing</u> hyperparameters of the model.
Mod Combination	Model Combination process is a set of techniques that allows to merge or combine the outputs of the different models to obtain best results.
Mod Calibration	Model Calibration process is a set of techniques that allows to calibrate the model when the probability estimate of a data point belonging to a class is very important. Calibration is comparison of the actual output and the expected output given by a system.
Mod Uncert Analys	Model Uncertainty Analysis process is a set of techniques that allows to <u>understanding</u> why your model is uncertain and how to estimate there level of uncertainty.
Mod Bias-Var	Model Bias-Variance Trade-off Analysis process is a set of techniques that allows to address the bias-variance dilemma or problem trying to simultaneously minimize these two sources of error that prevent supervised learning algorithms from generalizing beyond their training set.
Mod Overfitting	Model Overfitting Analysis process is a set of techniques that allows identify the overfitting as a possible cause of poor generalization performance of a predictive model.
Mod Prediction	Model Prediction process allow get the output of the model for unseen data.
Mod Finalization	Model Finalization process is a set of techniques that allows to finalize your machine learning model <u>in order to</u> make predictions on new data, training the final model selected on all available data.
Mod Saving	Model Saving process allow takes a trained model and saves the entire transformation pipeline and trained model.

V – Model Operat

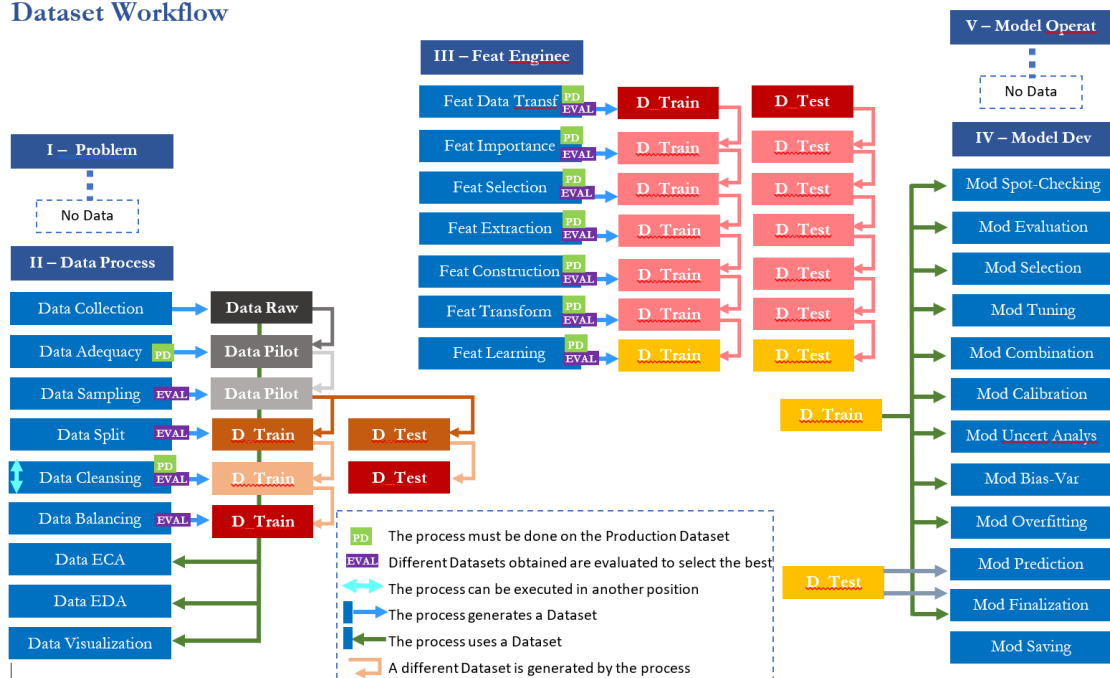
Mod Deploy	Model Deploy process allows to integrate the model into an existing production environment to make practical business decisions based on data.
Mod Execution	Model Execution process allows to execute the trained model with new unseen data in a production environment.
Mod Analysis	Model Analysis process allows to perform an analysis <u>in order to</u> mitigate possible unexpectedly poor <u>behavior</u> when the models are deployed in real-world domains, because model or data drift, <u>underspecification</u> or uncertainty analysis.
Mod Updating	Model Updating process allows to update the model in real time or production time, <u>in order to</u> adapt it dynamically to new data or changes in the problem domain.
Mod Explainability	Model <u>Explainability</u> process seeks to understanding of why machine learning models make the decisions they do, and why it matters.

El flujo de trabajo a partir de estos dominios y procesos es el siguiente:

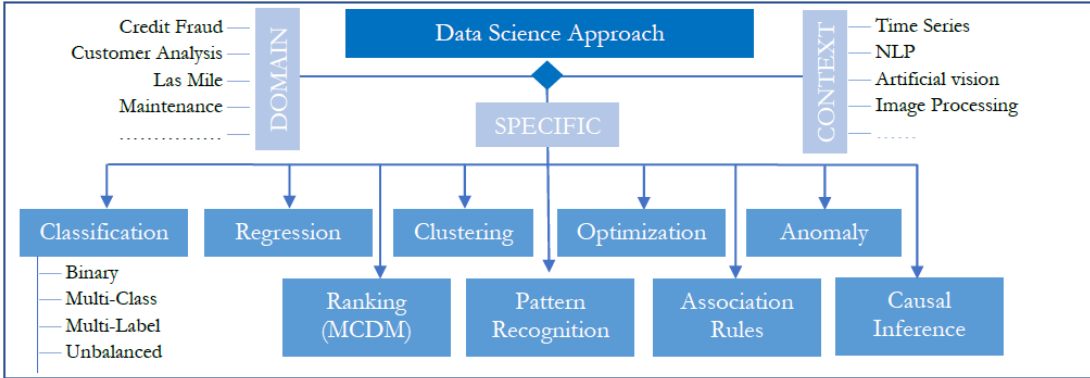
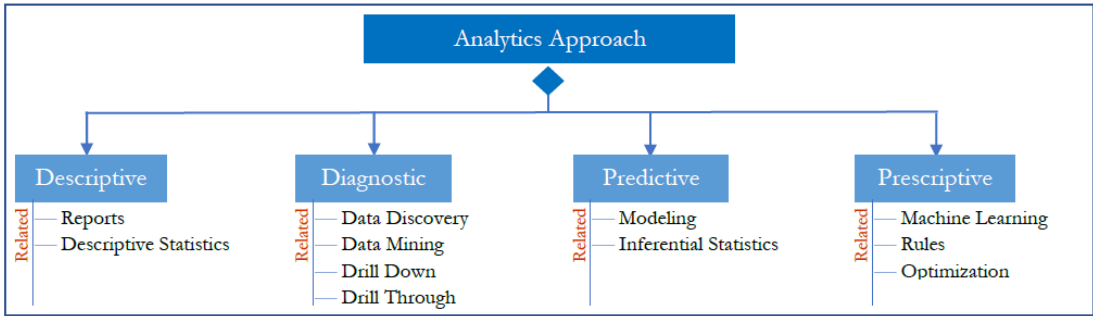
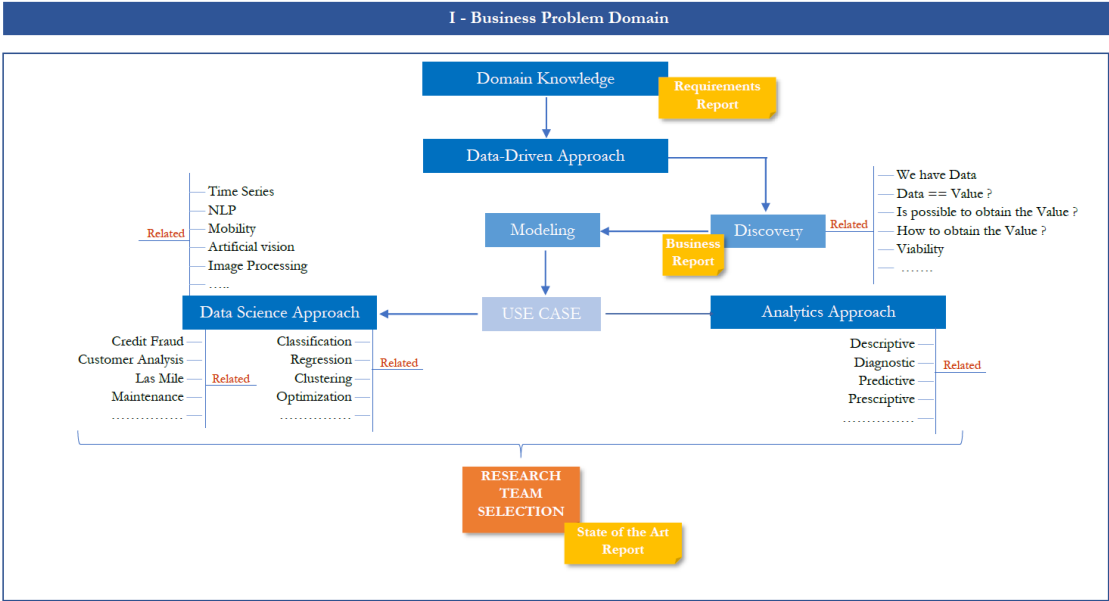


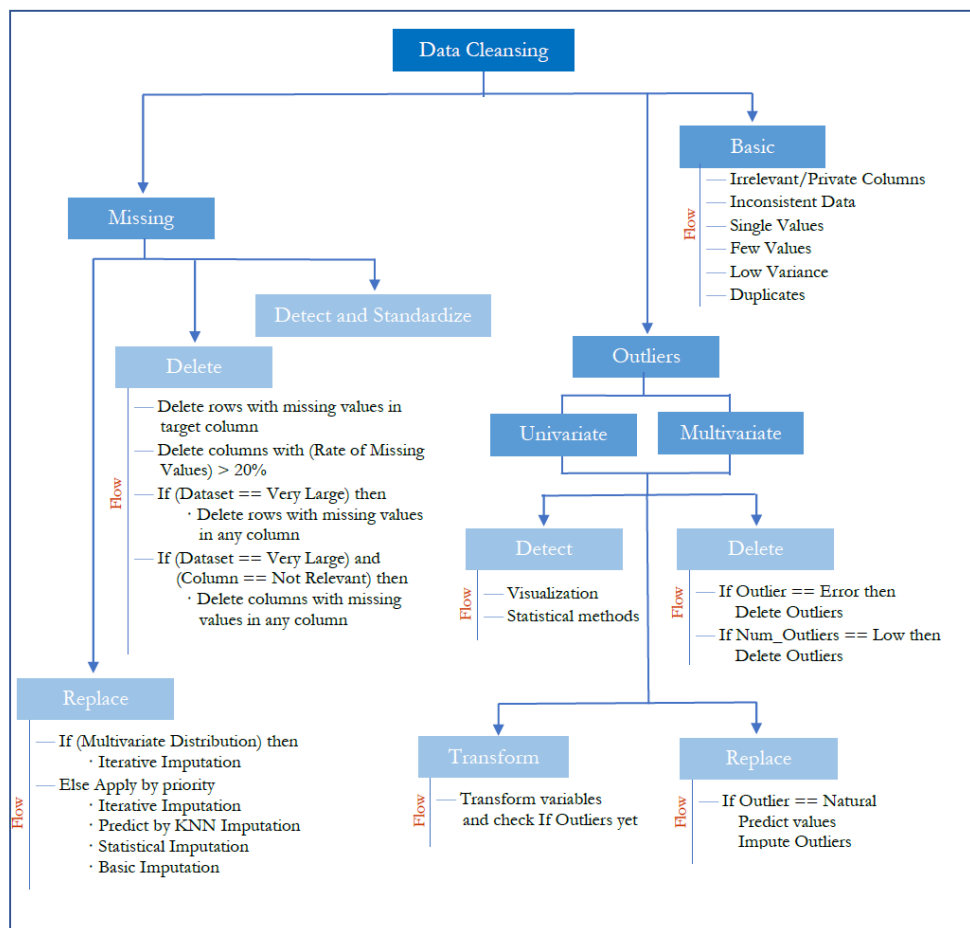
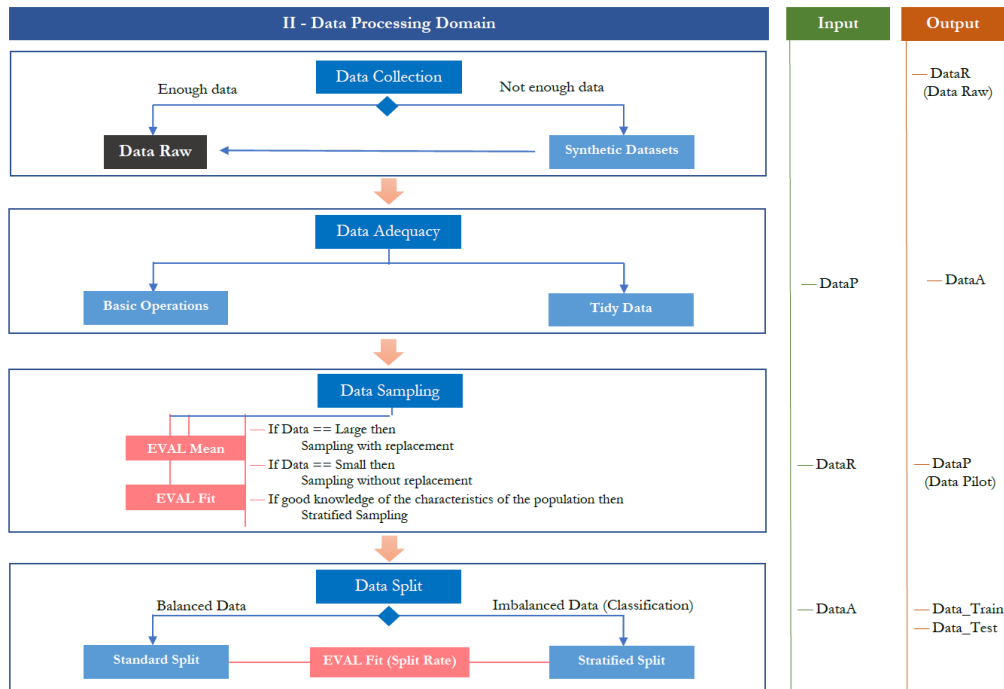
En referencia a los trabajos relacionados con los “datasets”, los flujos de trabajo generales se representan de la siguiente forma:

Dataset Workflow



Algunos de los procesos, que requieren despliegue, se esquematizan a continuación:





Se definen en la metodología definida aspectos relativos a la mejora continua:

Performance Improvement Flowchart

