

# Know Your data & Build Predictive Modeling

Meshael AlMuhanna, Unified Governance & Integration  
Technical Specialist.

Hissah AlMuneef, Cloud Developer Advocate

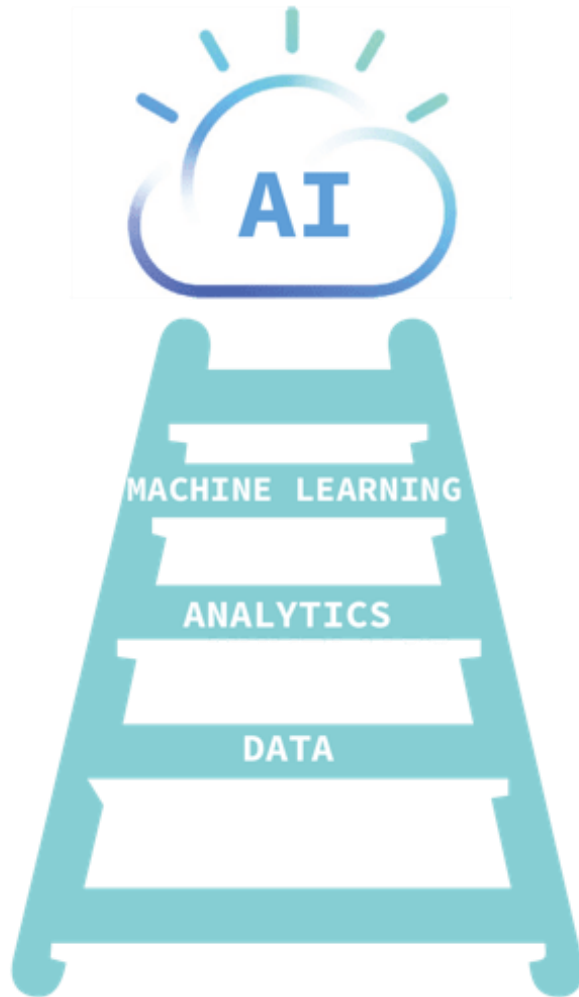
Meaad AlRshoud, Cloud Developer Advocate



## ❖ Agenda:

- IBM's AI ladder.
- Demonstration of data quality and ETL tools.
- Watson Studio overview.
- Predictive model use case.

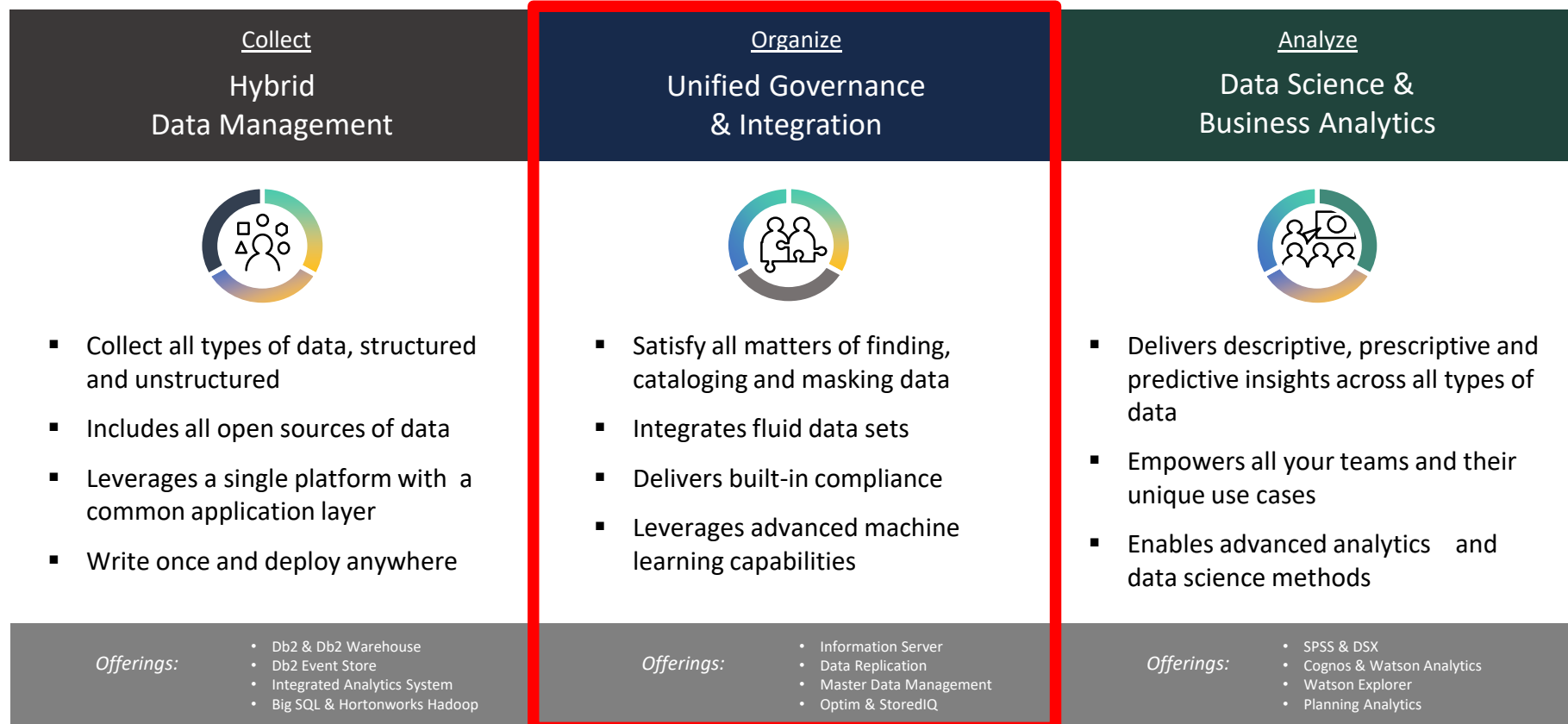




The AI Ladder

## IBM's Steps to Successful AI Journey

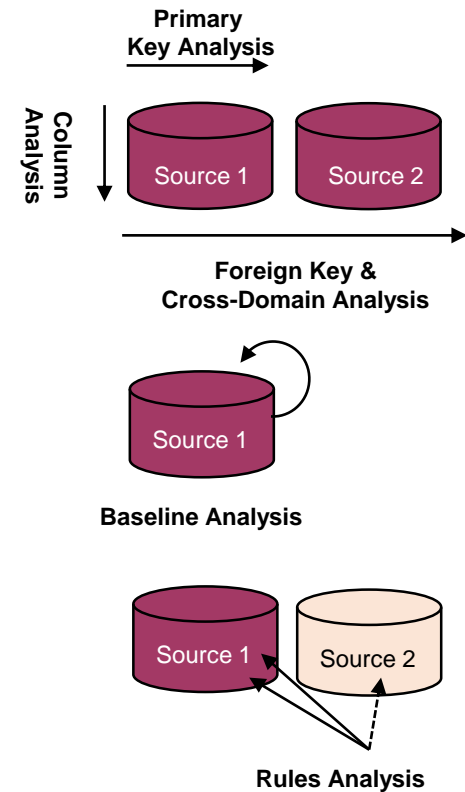
# IBM platforms deliver the capabilities our clients need



User & application independence across on premise, private cloud, and public cloud

# Understand the Quality of Data Sources

- Data Quality Score: estimate the proportion of reliable data values in the given dataset.
- Run Quality Scanner to calculate quality score
- Declare the type of problem to scan and how many passes over the data
- Findings will be all aggregated
- Score will be calculated



# Out of The Box Problem Detected



- Missing Values
  - Check missing values where Null values are not expected



- Uniqueness Violation
  - Check duplicate values



- Invalid Format
  - Checks for values



- Inconsistency Detection
  - Checks for values have different use of case



- Suspect Outlier
  - Checks for values that seem not to be of the same domain as other



- Violation of Correlation
  - Finds correlation between columns



- Data Rule Violation
  - Runs analysis against defined data rules

## EXAMPLE

Cust ID	Name	Age	Phone	Gender
62413	Lucy V Adler	32	334-555-6633	F
62414	Cory J Gardner	25	903-222-1255	F
62414	Mary H Jacques	18	777-156-9836	F
62415	Pdsaojfsadpoifj	46	xxxx	M
62416	Shaun Q Dunda	156	904-555-2940	M
62417	Carol T Schwartz	22	804-555-3164	F
62418	HARRIS LAURENT	36	785-555-5835	-



## EXAMPLE

Cust ID	Name	Age	Phone	Gender
62413	Lucy V Adler	32	334-555-6633	F
<b>62414</b>	Cory J Gardner	25	903-222-1255	F
<b>62414</b>	Mary H Jacques	18	777-156-9836	F
62415	<b>Pdsaojfsadpoifj</b>	46	<b>xxxx</b>	M
62416	Shaun Q Dunda	<b>156</b>	904-555-2940	M
62417	Carol T Schwartz	22	804-555-3164	F
62418	<b>HARRIS LAURENT</b>	36	785-555-5835	-

## EXAMPLE

Cust ID	Name	Age	Phone	Gender
62413	V Adler	35	55-6633	F
62414	Gardner		2-1255	F
62414	Jacqueline		6-9836	F
62414	Pdsaojfsadpoi,		xxxx	M
62414	Shaun Q Dunda	156	904	M
62414	Carol T Schwartz	22	804	F
62418	HARRIS LAURENT	36	785-5	

Uniqueness violation  
Conf:100%

Suspect value  
Conf:90%

Outlier  
Conf:95%

Data Class Violation  
Conf:100%

Inconsistent Case  
Conf:98%

Missing Value  
Conf:100%

## EXAMPLE

Data Set Score: 80%	Score: 71%	Score: 73%	Score: 86%	Score: 85%	Score: 85%
	Cust ID	Name	Age	Phone	Gender
	62413	Lucy V Adler	32	334-555-6633	F
	62414	Cory J Gardner	25	903-222-1255	F
	62414	Mary H Jacques	18	777-156-9836	F
	62415	Pdsaojfsadpoifj	46	xxxx	M
	62416	Shaun Q Dunda	156	904-555-2940	M
	62417	Carol T Schwartz	22	804-555-3164	F
	62418	HARRIS LAURENT	36	785-555-5835	-

**Problems have been Identified, What's Next?**

# Fix Identified Quality Issues

## Examples of Rules:

- The Gender field must be populated and must be in the list of accepted values
- The Social Security Number must be numeric and in the format 999-99-9999
- If Date of Birth Exists AND Date of Birth > 1900-01-01 and < TODAY  
Then Customer Type Equals 'P'
- The Bank Account Branch ID is valid in the Branch Reference master list

**Data Rules Editor for Stage Data\_Rules\_26**

**Published Rule Definitions**

Name	Type	Description	Rule Logic
SSN_Exists	Rule		socialsecuritynum exists
SSN_Matches_Format	Rule		socialsecuritynum matches_format
Value_Is_Numeric	Rule		myvalue is_numeric
Valid_Address	Rule		customersaddress in_reference_col
SSN_Exists_1	Rule	SSN is not null	socialsecuritynum exists

**Selected Rule Definitions**

Name	Type	Description	Rule Logic
SSN_Exists	Rule		socialsecuritynum exists

**Rule Logic:**  
socialsecuritynum exists

**Input Links**

Name	Data Type	Precision	Scale
AllCustomers.Cust_ID	VARQCHAR	7	0
AllCustomers.SSN	VARQCHAR	11	0
AllCustomers.Address	VARQCHAR	28	0
AllCustomers.City	VARQCHAR	12	0

**Rule Variables**

Binding	Name	Variable Type	Data Type	Rule Definition
AllCustomers.SSN	socialsecuritynum	Source Data	ANY	SSN_Exists

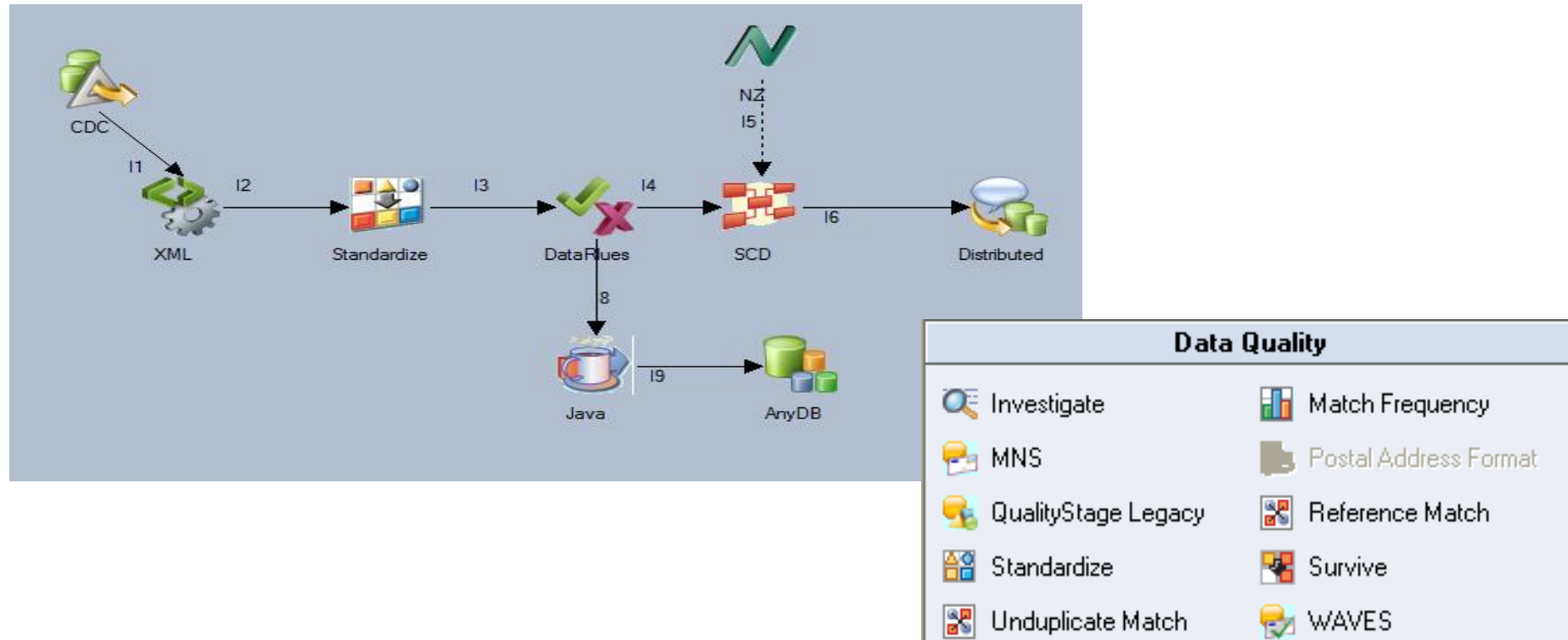
**Select Quality Control to Work With**

**Quality Controls**

Name	Type	Description
All		Global category for a
Published Rules		
01 Personal Identity		
02 Financial		
03 Human Resources		
04 Asset Identity		
05 Product		
06 Orders and Sales		
07 Data Format		
08 Validity and Completeness		
09 US Standardization		

# Enforce Quality on data

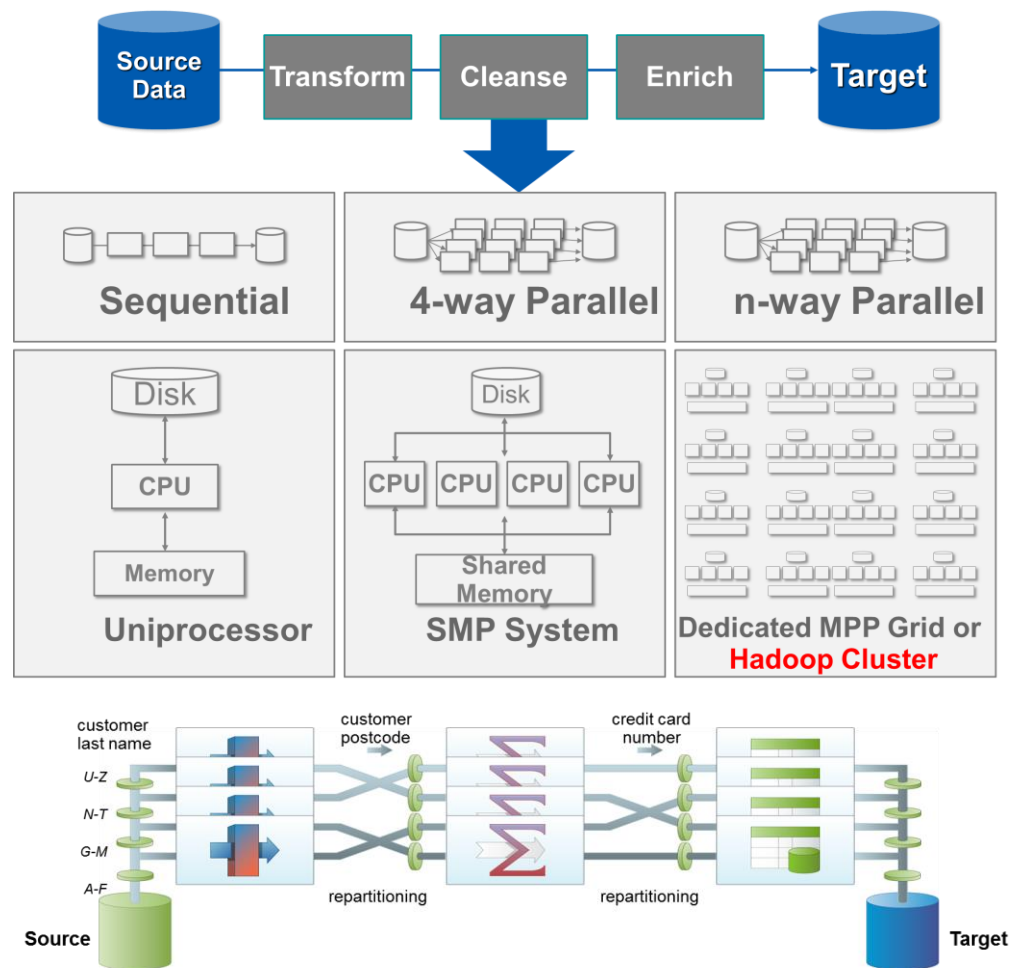
Fully integrated ETL & Data qualities capabilities



## Why Information Integration is Important?



# Performance



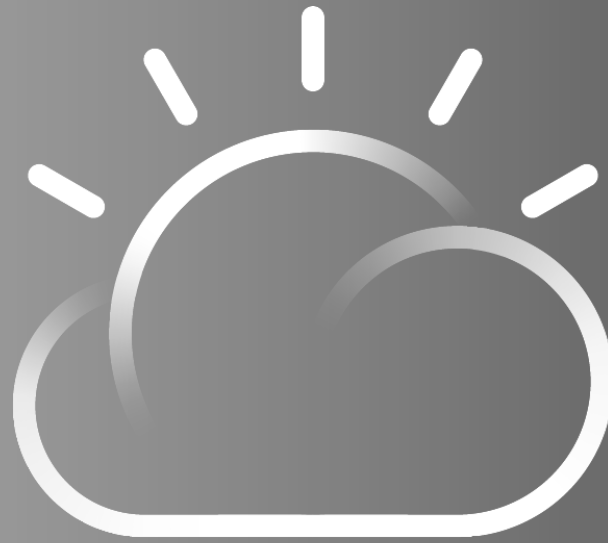


# Connectivity



# Predict Loan Eligibility Using SPSS in Watson Studio

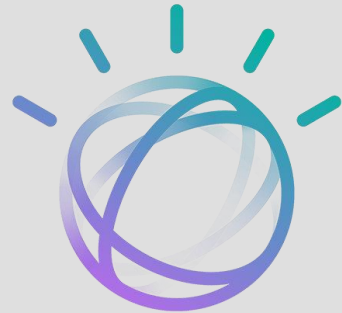
---



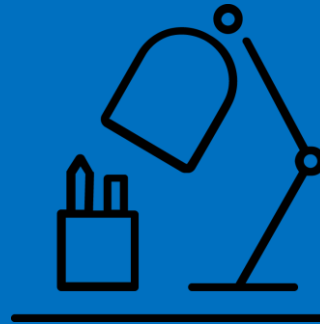
# Machine Learning



## IBM Watson



## Watson Studio



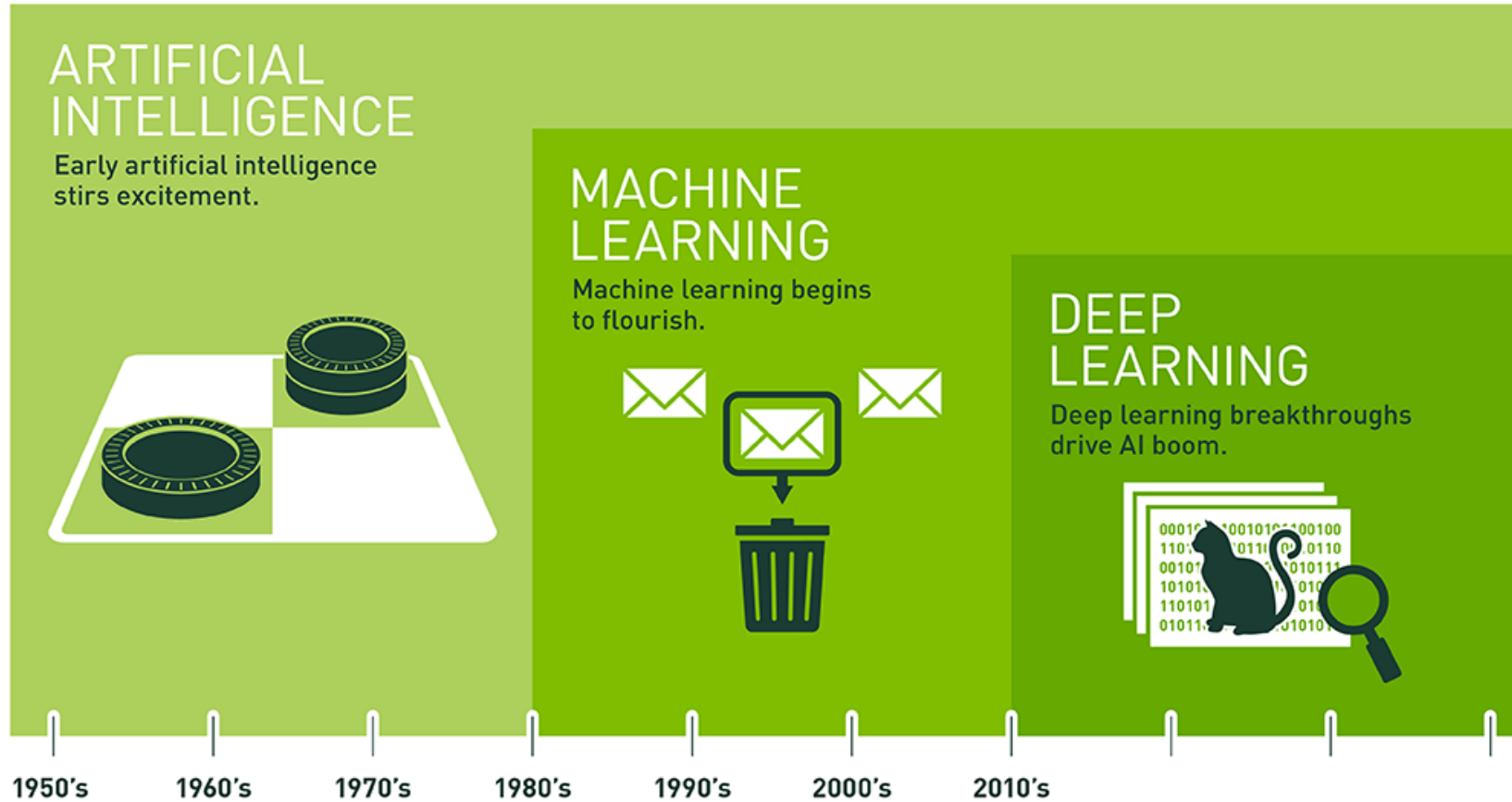
## Loan Eligibility Predictive Model



# Machine Learning



# Concept



Since an early flush of optimism in the 1950s, smaller subsets of artificial intelligence – first machine learning, then deep learning, a subset of machine learning – have created ever larger disruptions.



# Artificial Intelligence



## Netflix

NETFLIX

Machine learning is integral to Netflix's video recommendation engine. The company has valued the ROI of these algorithms at £1 billion a year due to their impact on customer retention.

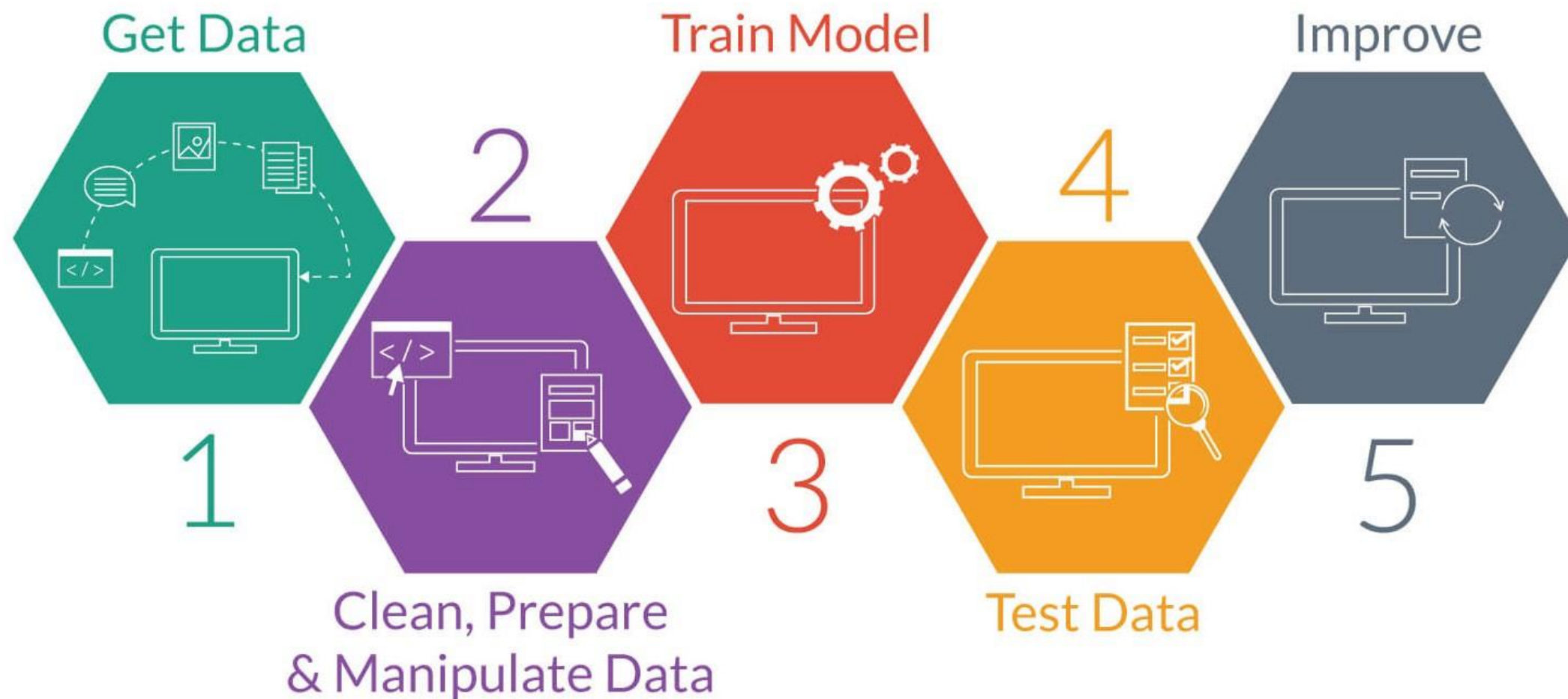
## PayPal



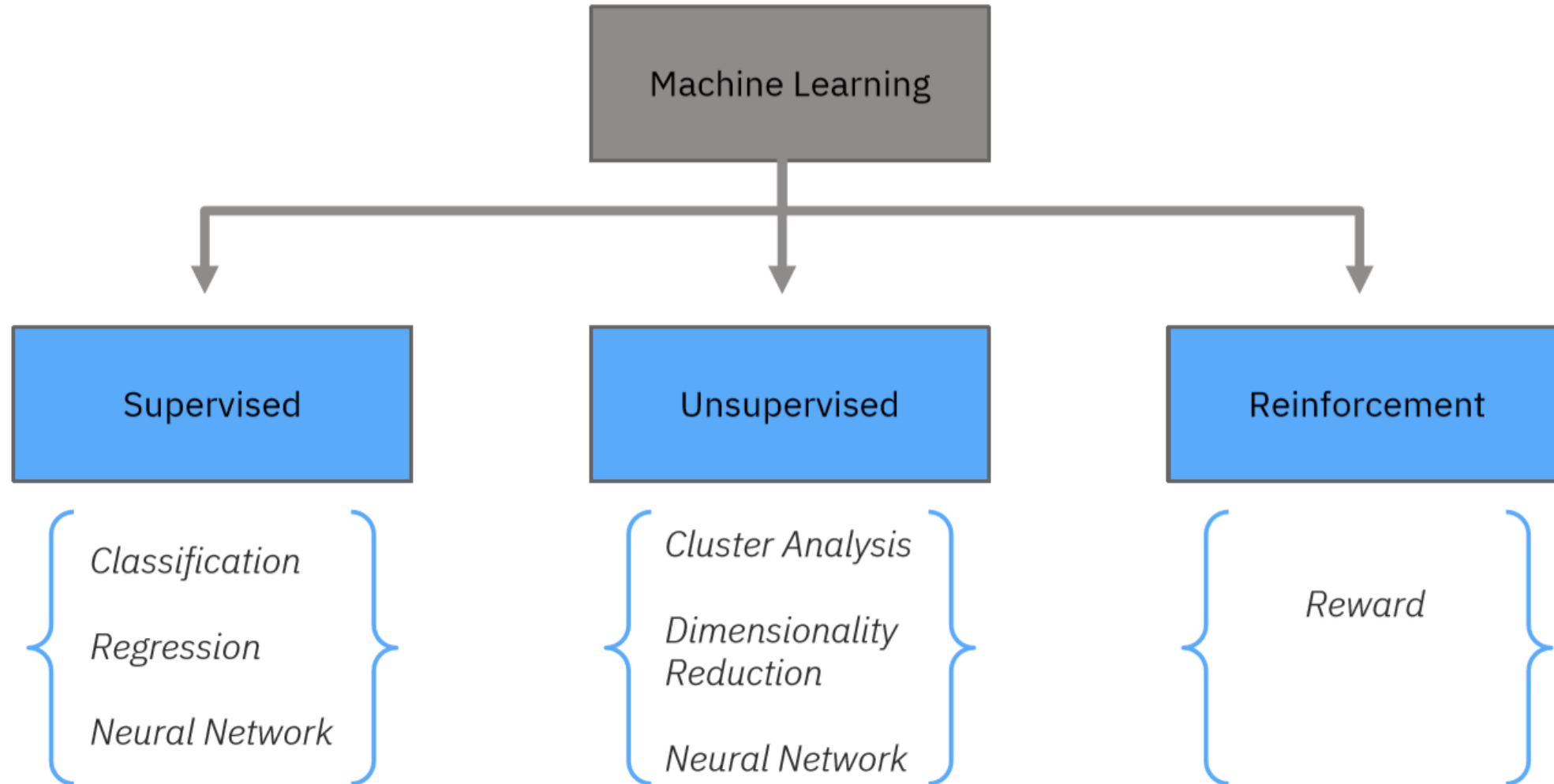
The online payment platform uses machine learning algorithms to combat fraud. By implementing deep learning techniques, PayPal analyses vast quantities of customer data and evaluates risk accordingly.

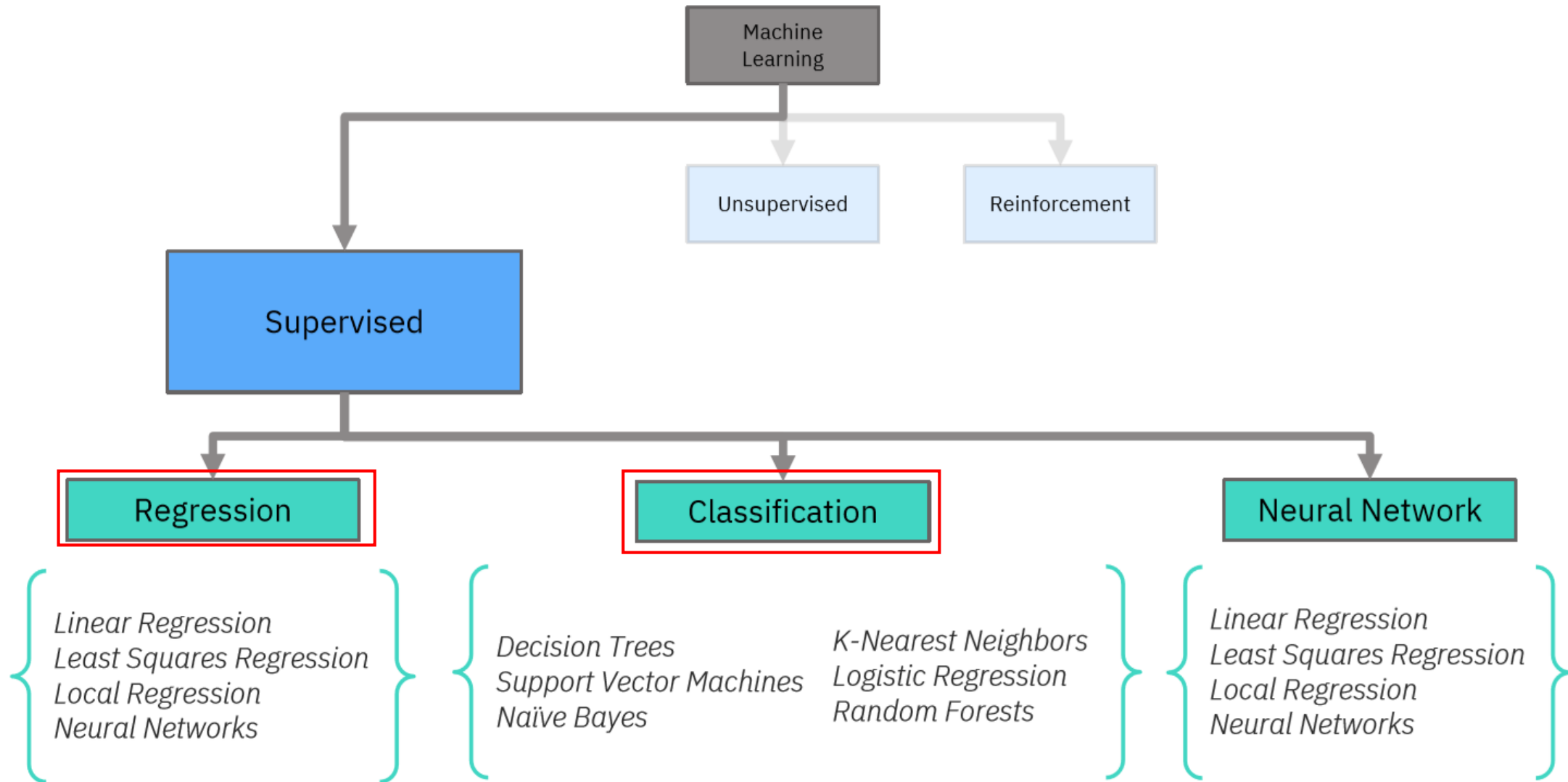
# Machine Learning

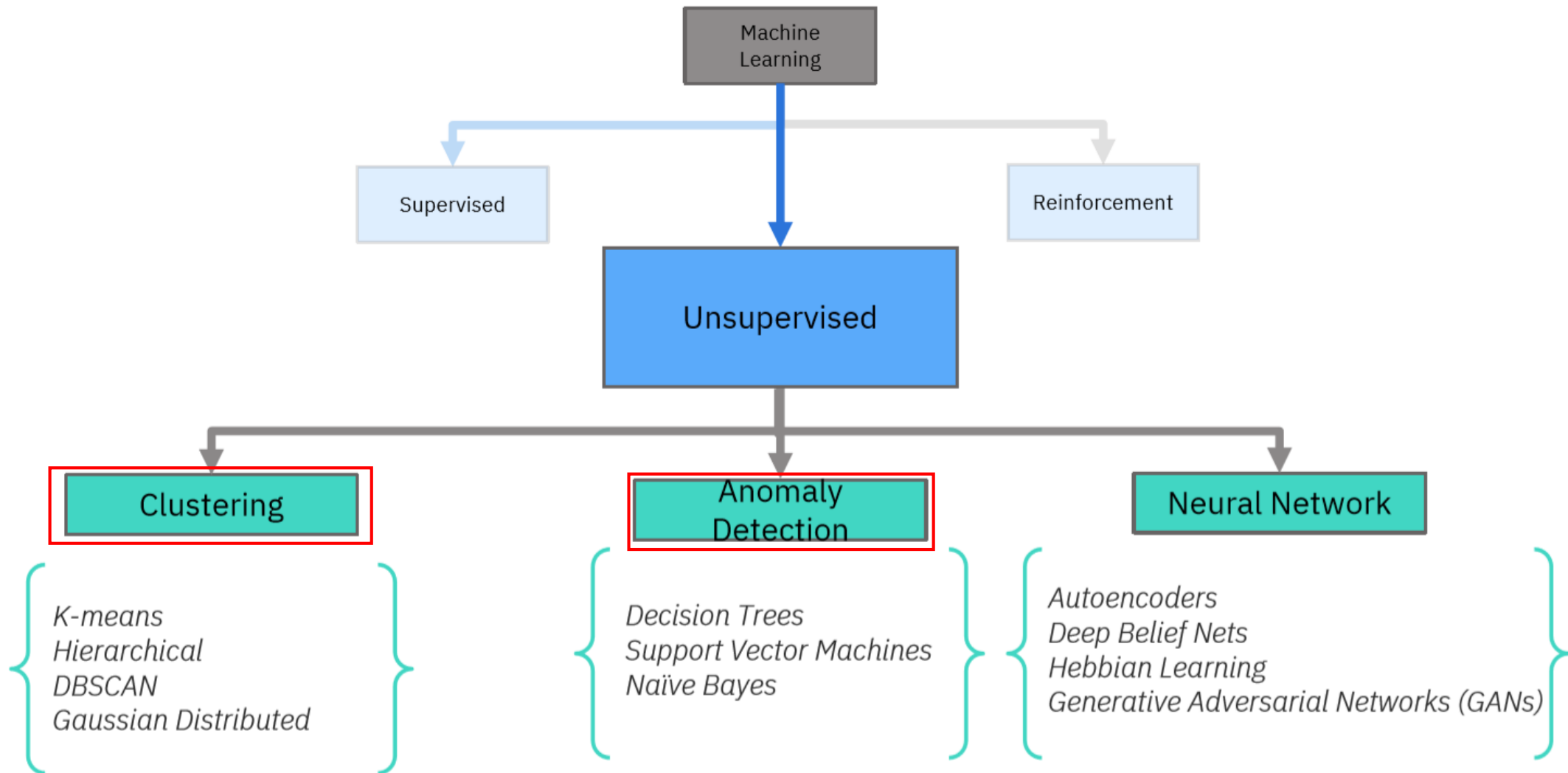
# Methodology



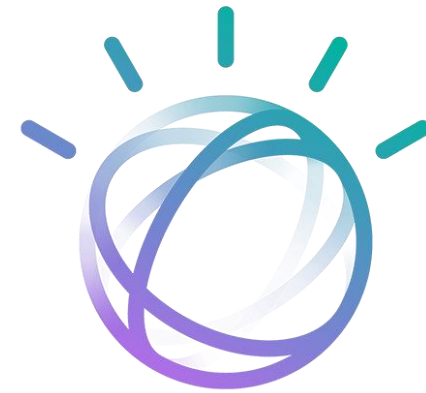




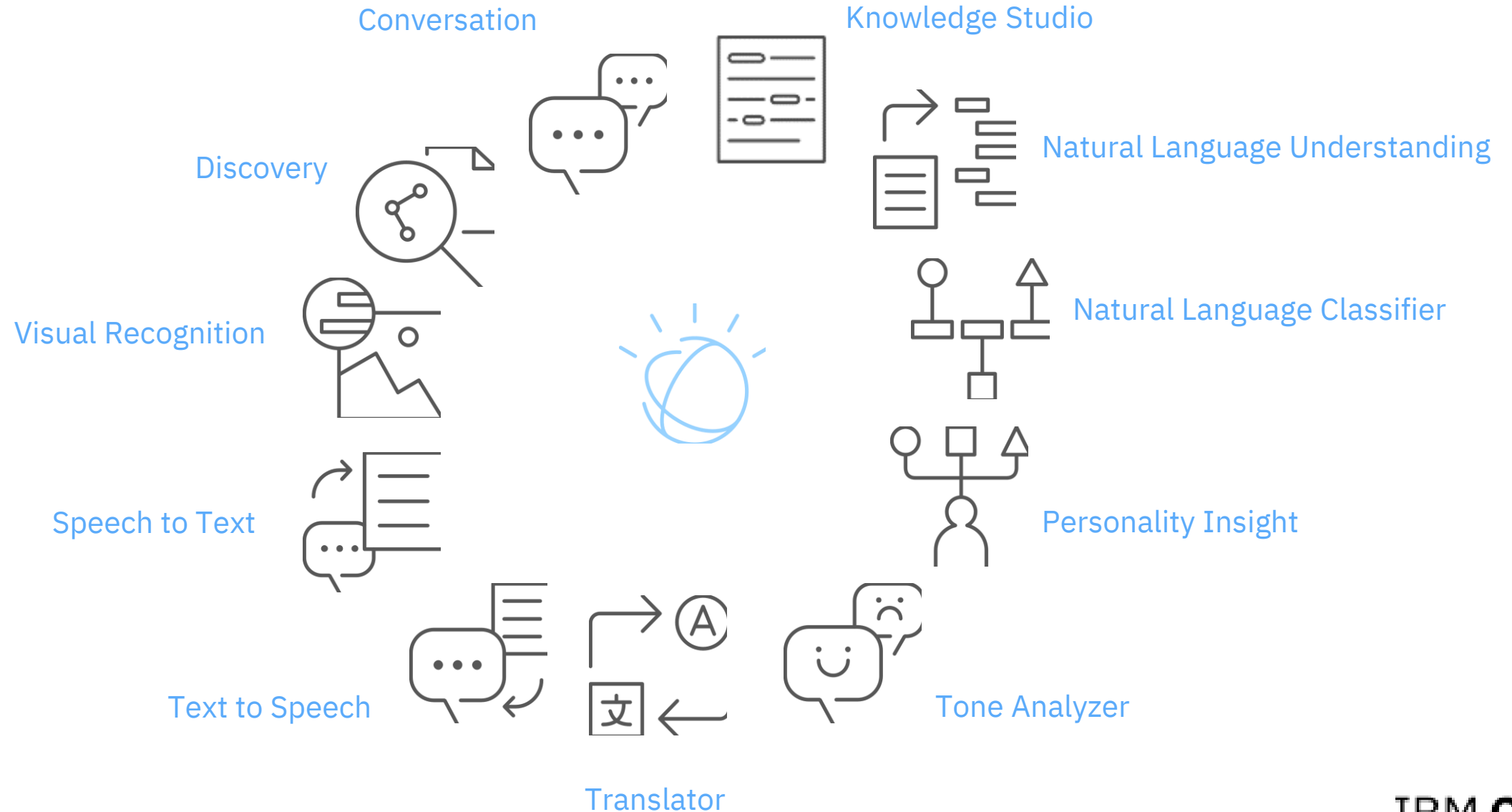




# Watson is AI for Business



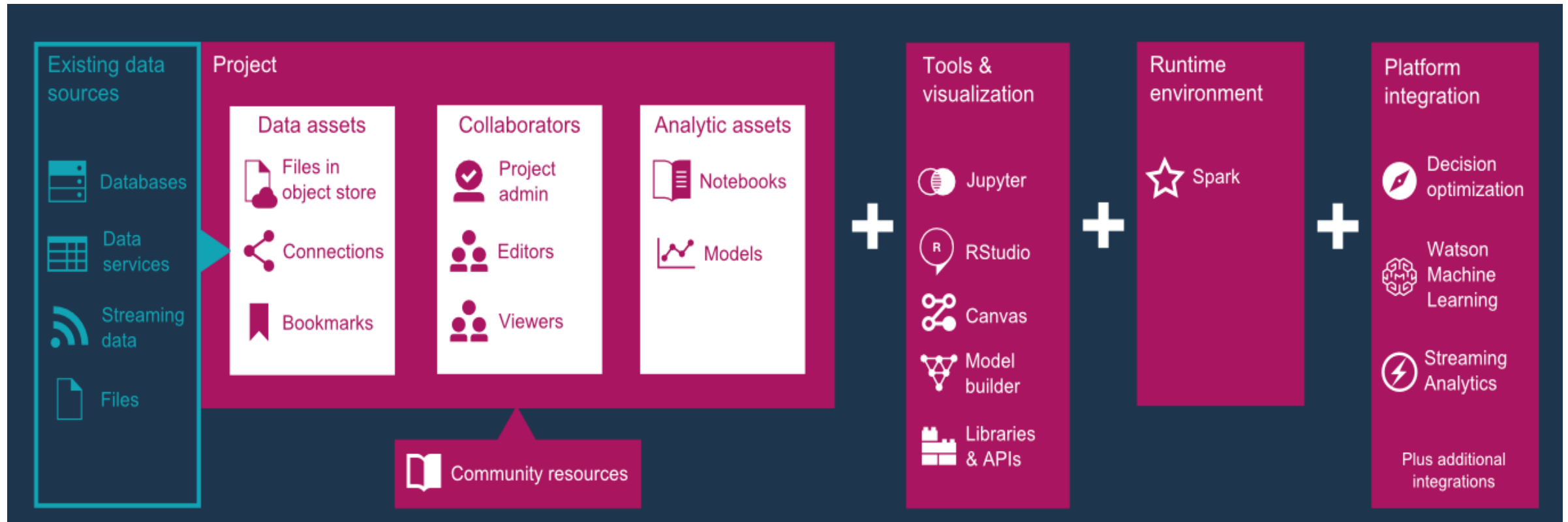
# With Watson:



# Watson Studio



# Watson Studio



# Predict Loan Eligibility Using SPSS in Watson Studio





# Problem Statement

Loans Company wants to automate the loan eligibility process based on customer detail provided while filling online application form.

# Data

## Not Feature

**Loan\_ID**  
String

## Features

<b>Gender</b> String	<b>Married</b> String	<b>Dependents</b> String	<b>Education</b> String	<b>Self_Employed</b> String	<b>ApplicantIncome</b> String
-------------------------	--------------------------	-----------------------------	----------------------------	--------------------------------	----------------------------------

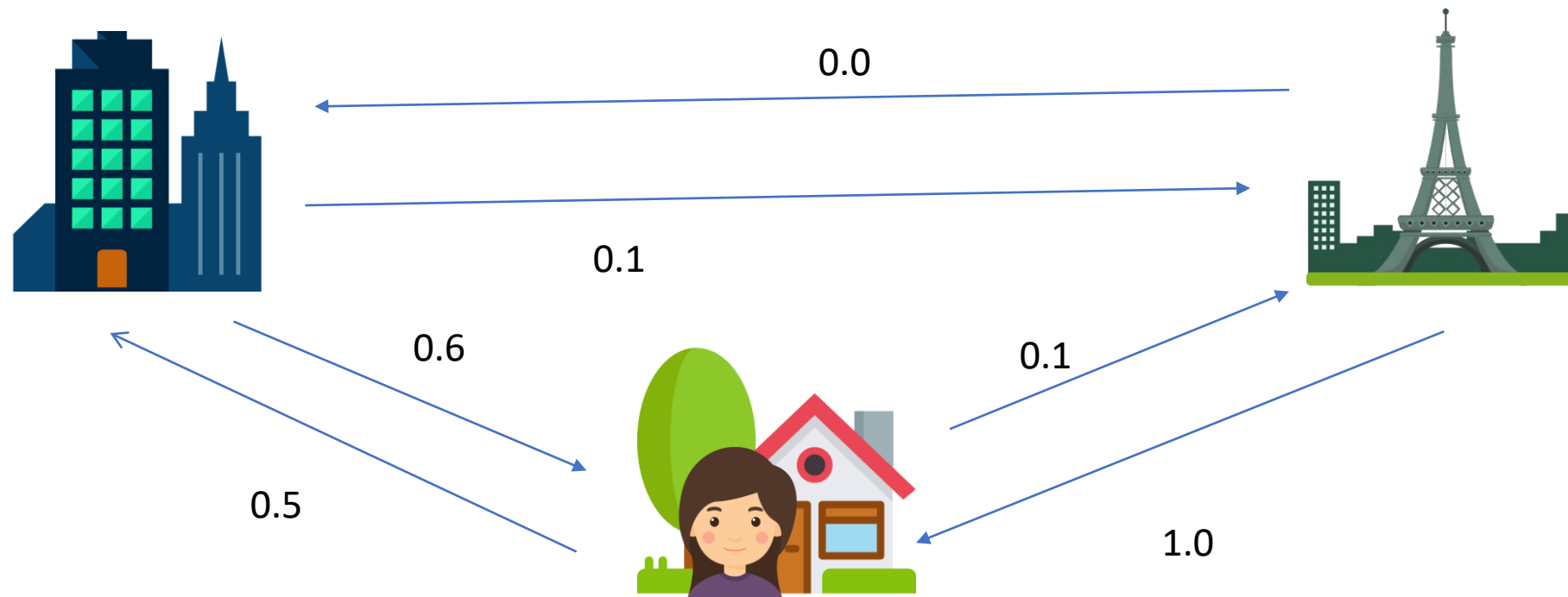
LP001002	Male	No	0	Graduate	No	5849
LP001003	Male	Yes	1	Graduate	No	4583
LP001005	Male	Yes	0	Graduate	Yes	3000

## Class

<b>CoapplicantIncome</b> Decimal	<b>LoanAmount</b> Decimal	<b>Loan_Amount_Term</b> Decimal	<b>Credit_History</b> Decimal	<b>Property_Area</b> String	<b>Loan_Status</b> String
-------------------------------------	------------------------------	------------------------------------	----------------------------------	--------------------------------	------------------------------

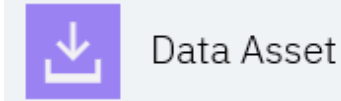
0	146.412162	360	1	Urban	Y
1508	128	360	1	Rural	N
0	66	360	1	Urban	Y

# Bayes Net



# Steps to Solution...

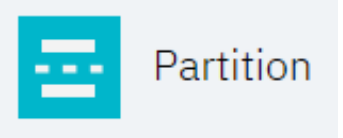
1. Import our Data using **Data Asset** node.



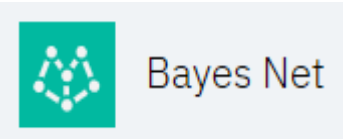
2. Configures variables type using **Types** node.



3. Split our data for training and testing sets using **Partition** node.



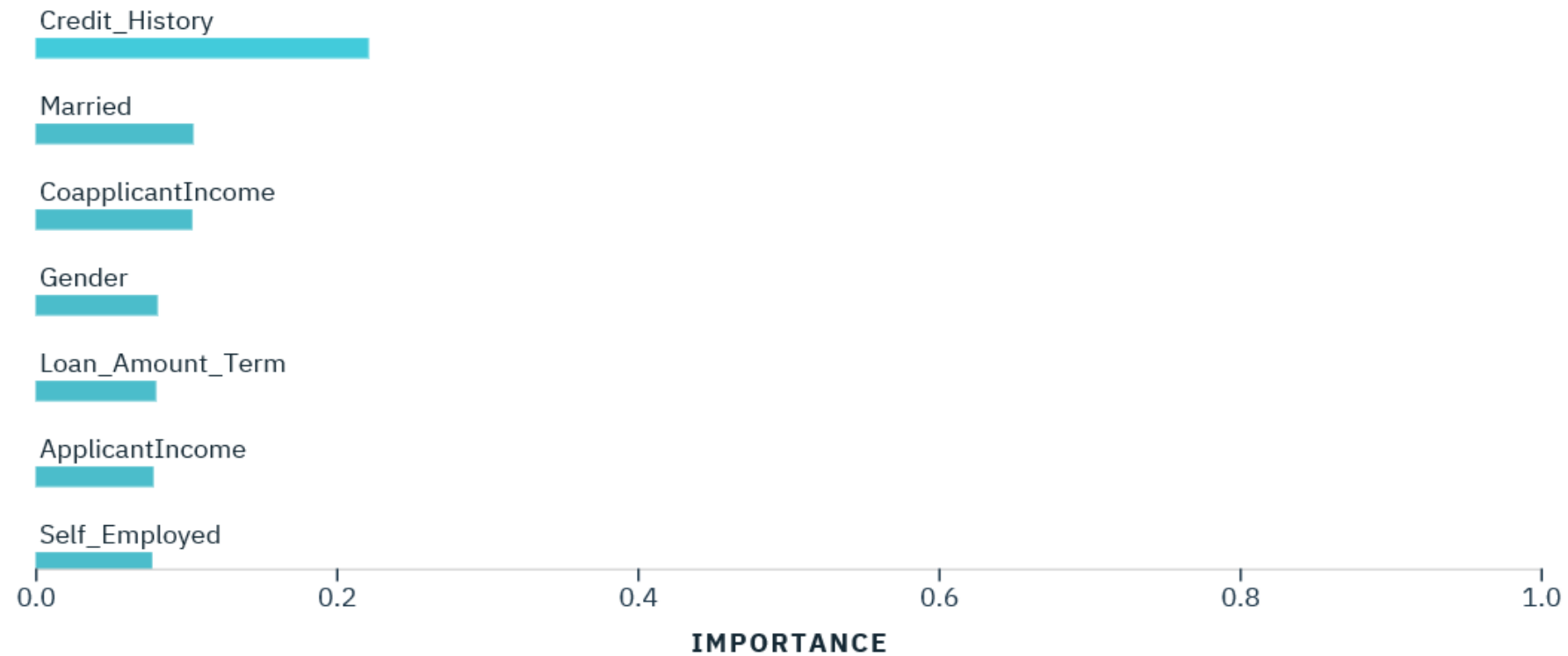
4. Build a probability model using Bayesian Network algorithm by the **Bayes Net** node.



5. Try other models ! Why not !

# Predictor Importance ⓘ

TARGET : LOAN\_STATUS



# CALL FOR CODE

## INNOVATION AND TECHNOLOGY FOR GOOD

The issue: Natural disaster preparedness and relief.  
How will you answer the call?

[Register For The Challenge](#)

[Amplify The Call](#)



# Get Started

## Call for Code

Commit for a **CAUSE**. Push for **CHANGE**.

### Call for Code Website:

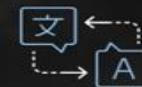
<https://developer.ibm.com/callforcode/>

### Challenge Details:

<https://callforcode.org/challenge/>



**Build secure, resilient, traceable, and transparent supply networks with blockchain.**



**Use AI and bots to improve real-time communications with natural language processing.**



**Understand, analyze, and predict health and nutrition needs to improve services with data science.**



**Improve logistics based on traffic and weather activity to reduce the number of people affected.**



**Collect and analyze device sensor data to take corrective or preventative action automatically.**



**Use machine learning, deep learning, and visual recognition to improve critical processes.**



# Resources

Learn – develop – connect

**IBM Code** ([developer.ibm.com/code](https://developer.ibm.com/code))

**IBM Developer Works** ([ibm.com/developerworks](https://ibm.com/developerworks))

**GitHub** ([github.com/DevExCodeHub](https://github.com/DevExCodeHub))

Learning Lab - Coursera - Udacity - more