# Towards Automating Data Science Workflows at the Scale of Banking

Jose A. Rodríguez-Serrano
BBVA Data & Analytics

*DataEngConf*
*Barcelona, September 2018*

Data Science +
Data Engineering

=

Reducing Inefficiencies in
Products or Services...

# Example 1: Connect retailers to customers



**Advertising Platform**

Retailer Dashboard      Campaign manager      Client App

https://www.bbvadata.com/cost-effective-scalable-collaborative-filtering-based-recommender-system/

# Example 2: Browse expenses more meaningfully

# Behind the Scenes:
# Expense Forecasting Research

## Heteroscedastic Neural Network

Real-world dataset of human expenses
> 100M time series / month



Brando et al., *Uncertainty Modelling in Deep Networks: Forecasting Short and Noisy Series*, ECML 2018

but…

Data Science Workflows
have inefficiencies
themselves

42,020 views | Mar 23, 2016, 09:33am

# Cleaning Big Data: Most Time-Consuming, Least Enjoyable Data Science Task, Survey Says

**Gil Press** Contributor ⓘ

*I write about technology, entrepreneurs and innovation.*

**TWEET THIS**

🐦 data scientists found that they spend most of their time massaging rather than mining or modeling data.

🐦 76% of data scientists view data preparation as the least enjoyable part of their work

We did document
The Real Data Science
Workflow

Validation of Data vs Business Question (Quick & Fast)

Validation of Data Ingestion
Validation of Data vs Business Question

Validation of Data Ingestion
Validation of Data vs Business Question

Exploratory Data Analysis

Error Analysis
Output Validation
KPI Validation

Model validation
Model interpretation

Model exploration
Model comparison

Feature exploration
Objective variables exploration

Business understanding

Analytic approach

Data requirements

Data collection

Data understanding

Data preparation

Modeling

Evaluation

Deployment

Feedback

API

ETL

# There do exist attempts to automate those workflows

**Google Cloud AutoML enhances AI**
Accessibility for all businesses



## The Automatic Statistician
An artificial intelligence for data science

Welcome to automatic exploratory data analysis

Making sense of data is one of the great challenges of the information age we live in. While it is becoming easier to collect and store all kinds of data, from personal medical data, to scientific data, to public data, and commercial data, there are



Consider TPOT your **Data Science Assistant**. TPOT is a Python Automated Machine Learning tool that optimizes machine learning pipelines using genetic programming.
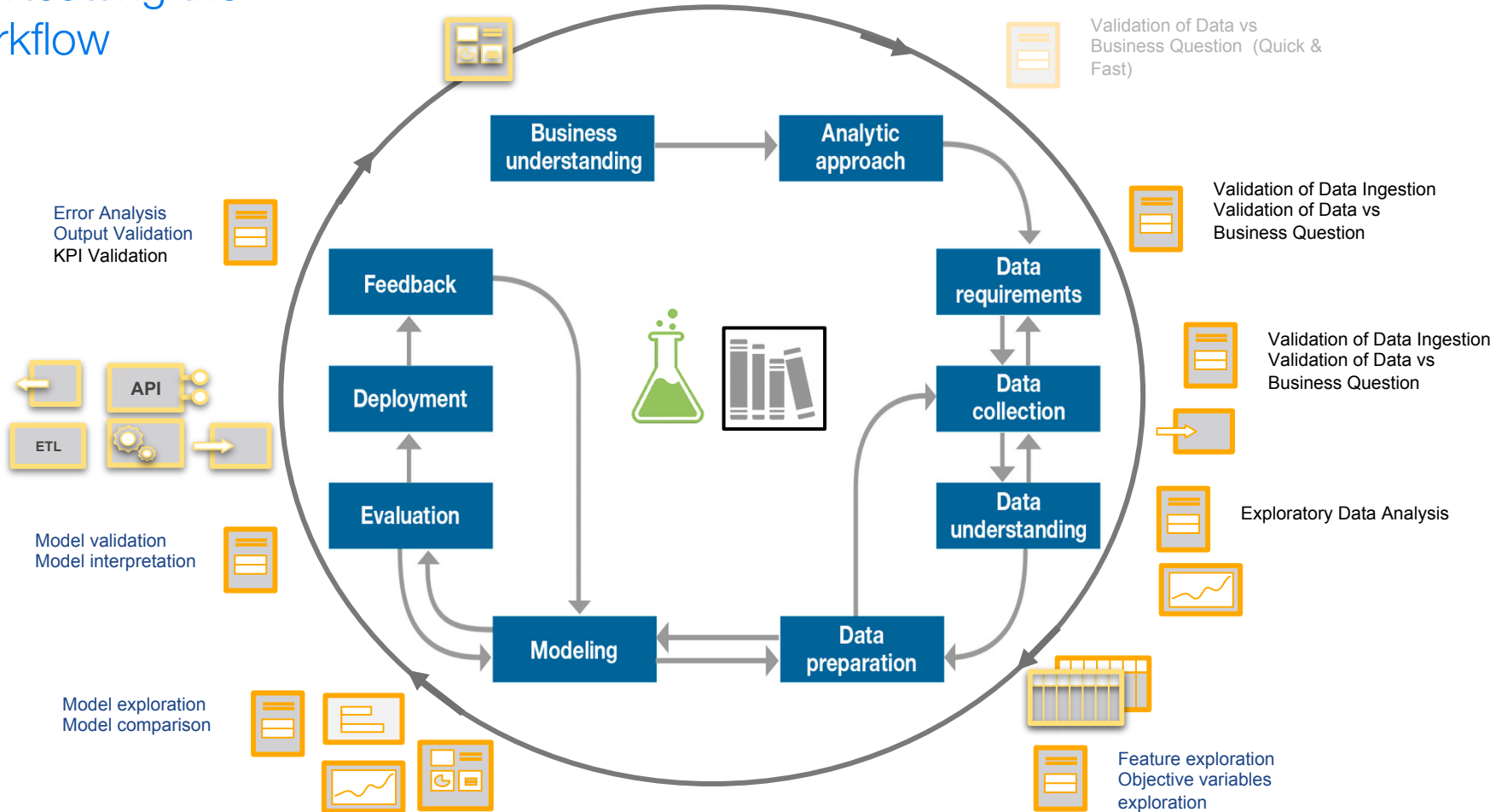


# MIT News
ON CAMPUS AND AROUND THE WORLD

Auto-tuning data science: New research streamlines machine learning
A new automated machine-learning system performs as well or better than its human counterparts — and works 100 times faster.
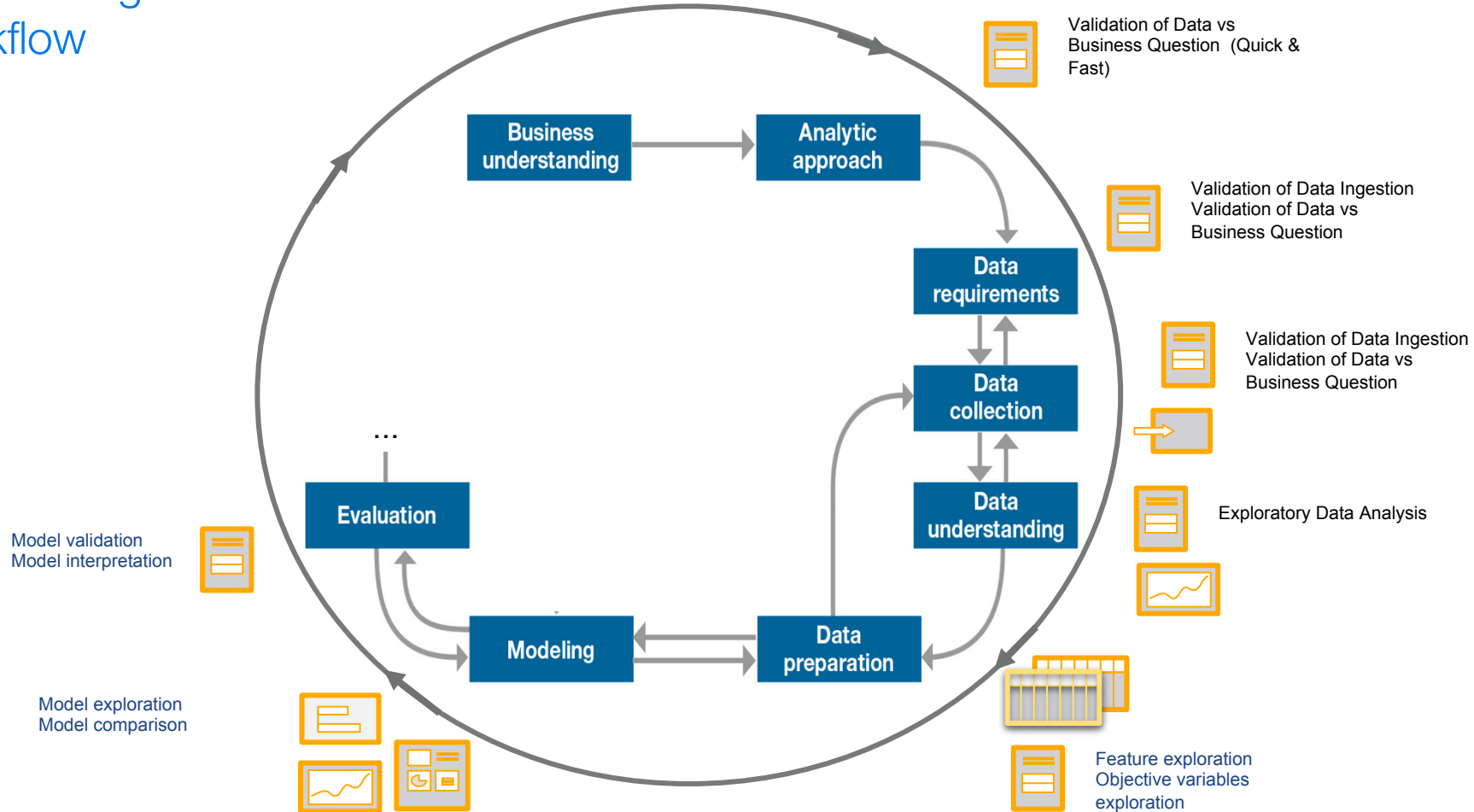
# Towards Increased Data Science Efficiency at BBVA Data & Analytics

# Shortcutting the Workflow

Validation of Data vs Business Question (Quick & Fast)

Error Analysis
Output Validation
KPI Validation

Validation of Data Ingestion
Validation of Data vs
Business Question

API
ETL

Validation of Data Ingestion
Validation of Data vs
Business Question

Business understanding → Analytic approach

Feedback

Deployment

Evaluation

Modeling ⇄ Data preparation

Data requirements

Data collection

Data understanding

Model validation
Model interpretation

Exploratory Data Analysis

Model exploration
Model comparison

Feature exploration
Objective variables exploration

# Shortcutting the Workflow

# Shortcutting the Workflow…



Validation of Data vs Business Question (Quick & Fast)

Validation of Data Ingestion
Validation of Data vs Business Question

Validation of Data Ingestion
Validation of Data vs Business Question

Exploratory Data Analysis

Business understanding

Analytic approach

Data requirements

Data collection

Data understanding

Systematic Baseline
(for problems with the same "X")

Evaluation

Modeling

Data preparation

Model validation
Model interpretation

Model exploration
Model comparison

Feature exploration
Objective variables exploration

Systematic Baseline =

Generic entity representation
+ flexible & simple model
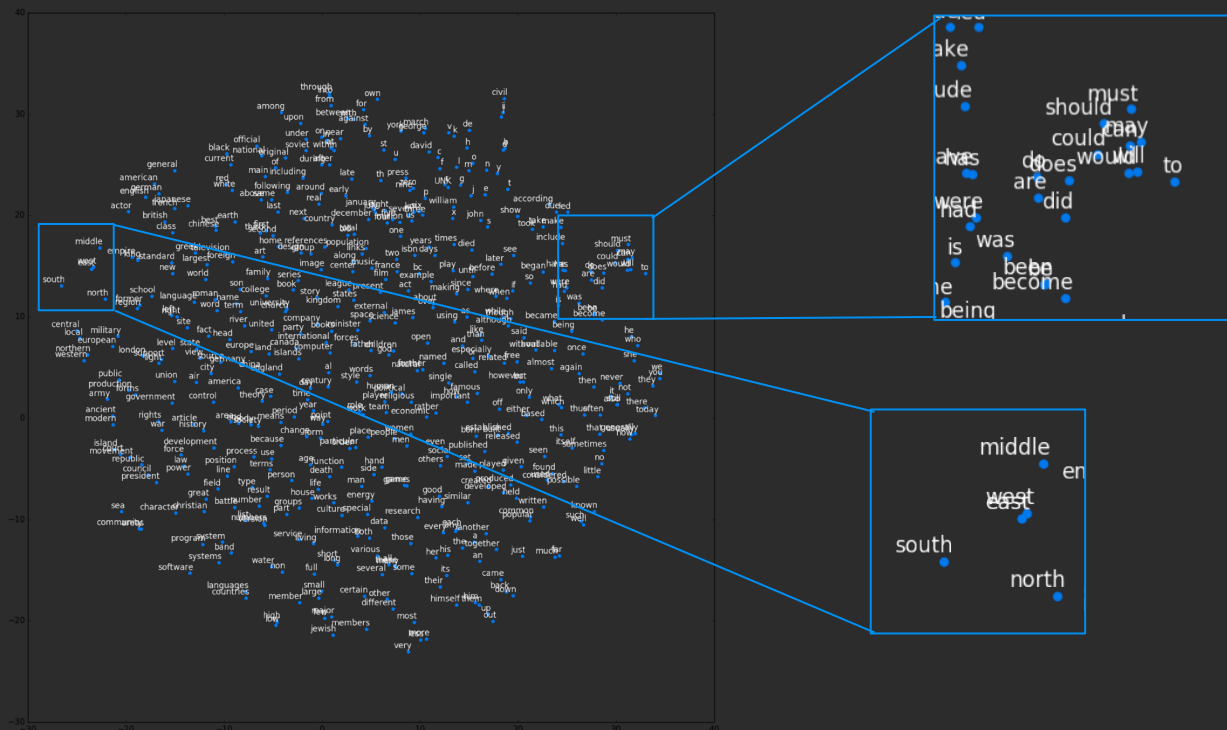
# Do Generic Attributes Exist?
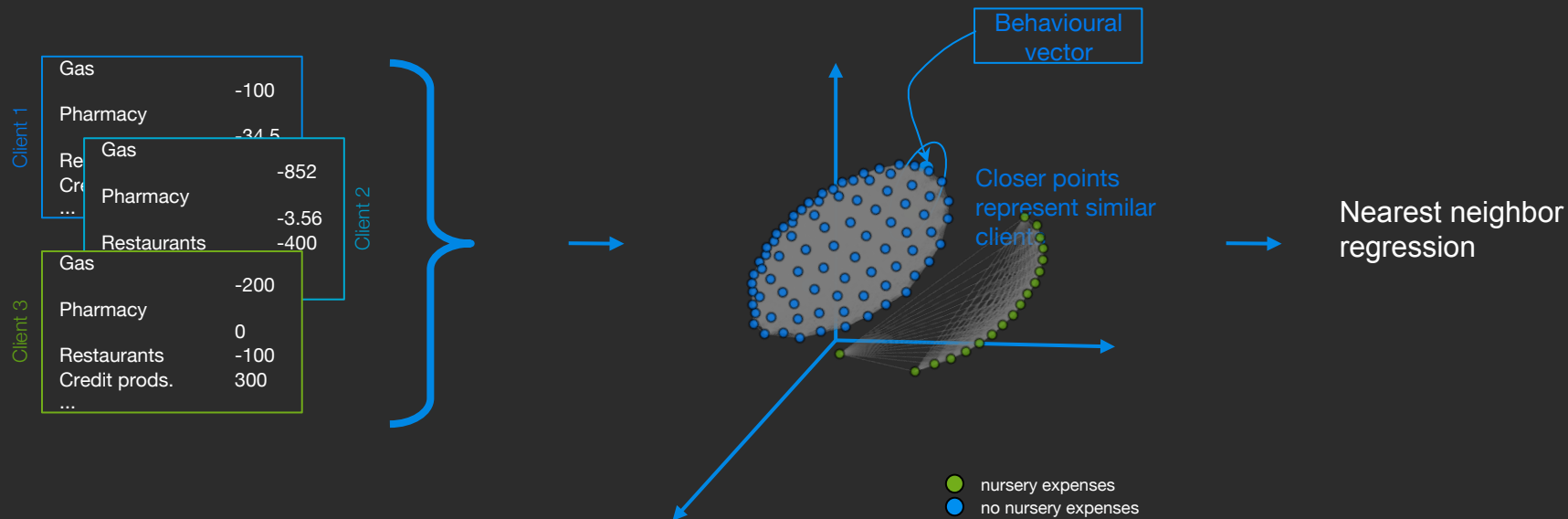


word2vec: Embeddings of similar words are close together

# Our Systematic Baseline

# Computing the generic representation
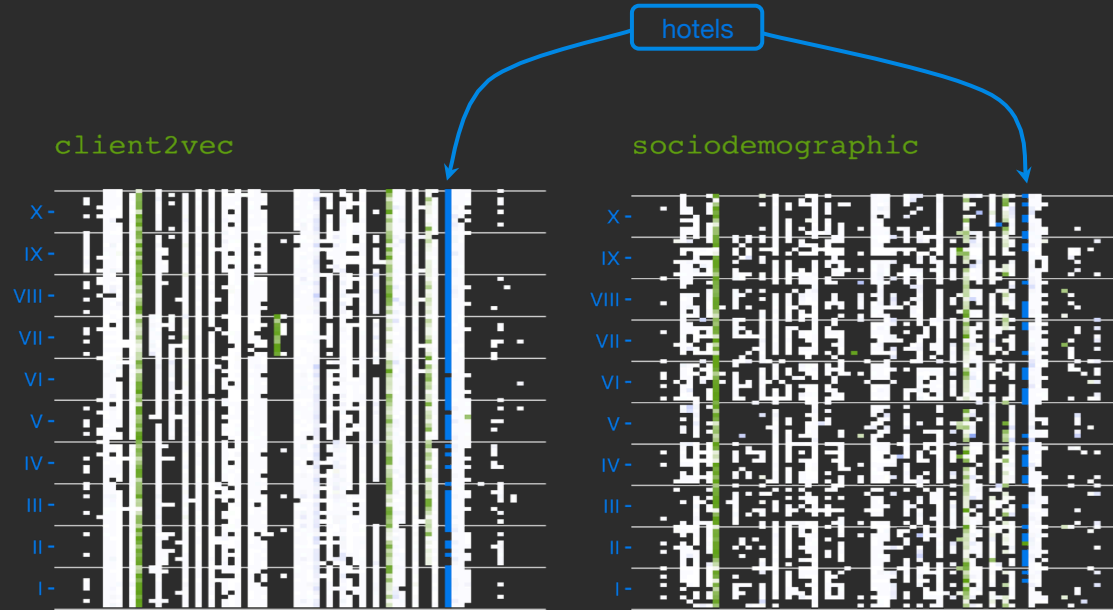
Denoising autoencoder

| | |
|---|---|
| Trains | -500.4 |
| Travel | -2000 |
| Hotels | -1580 |

| |
|---|
| -500.4 |
| -2000 |
| -1580 |

Encoder

Decoder

Client representation

- Learn to reconstruct corrupted data
- Reconstruction ≈ similarity
- Marginalized stacked denoising autoencoders (Chen et al, ICML 2012)

hotels

client2vec

sociodemographic

Sociodemographic variables don't capture typical behaviour

Improvement on two use cases

# +61.6%

## Client clustering

Group similar clients and compare their expenses in a target category

# +76.1%

## Category prediction

By looking at similar clients, guess whether a client had an expense in a target category

Baldassini et al., client2vec: Towards Systematic Baselines for Banking Applications, arXiv, 2018

## Algorithmic research

Explore, evaluate and generate
state-of-the-art algorithms

- State-of-the-art algorithm
- Best performance in our use cases

## Implementation

Deliver algorithmic solutions as a software package

- Lightweight training
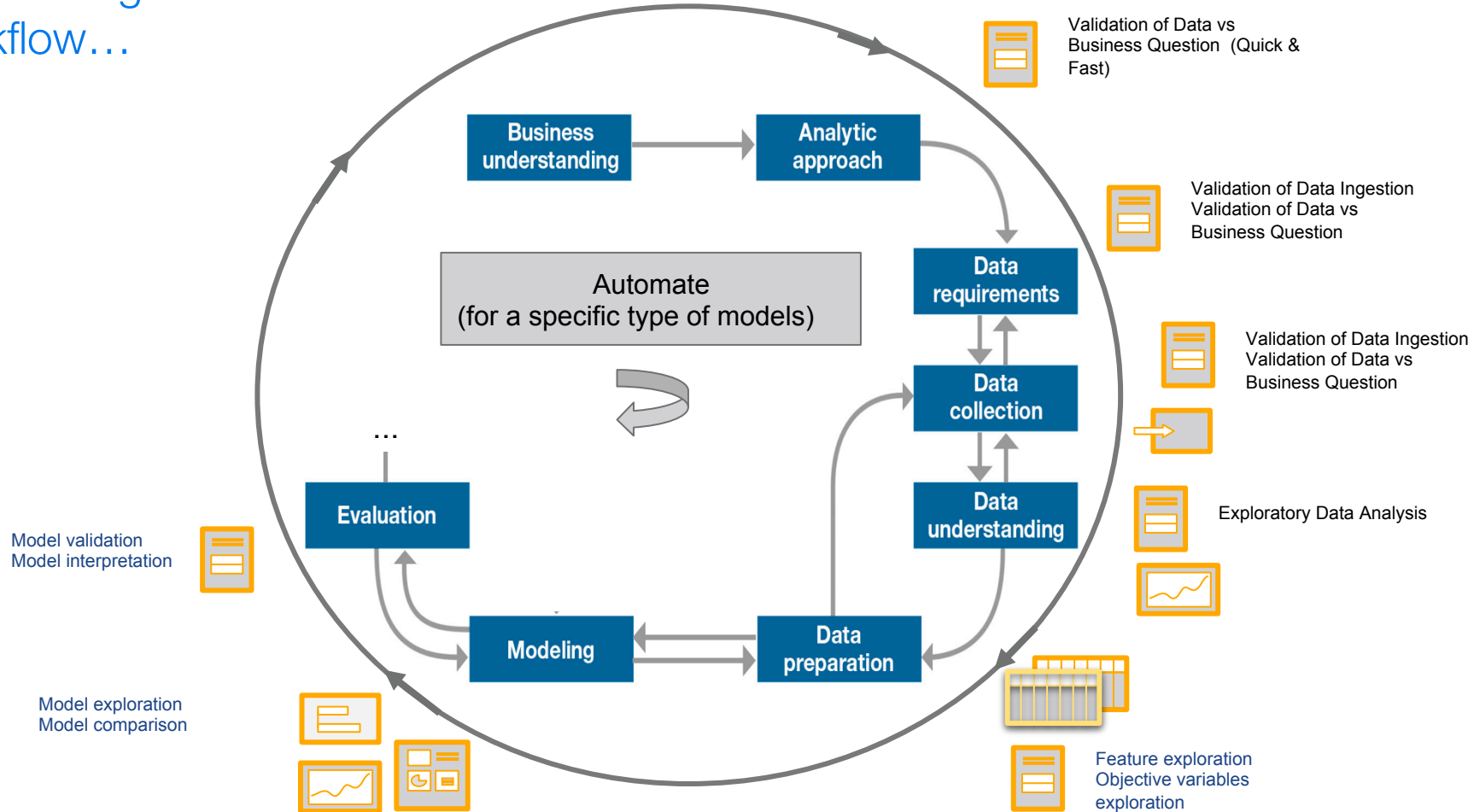- Works on bank's infrastructure

## Product enablement

Generate capabilities to accelerate product development

- Better method to compare clients
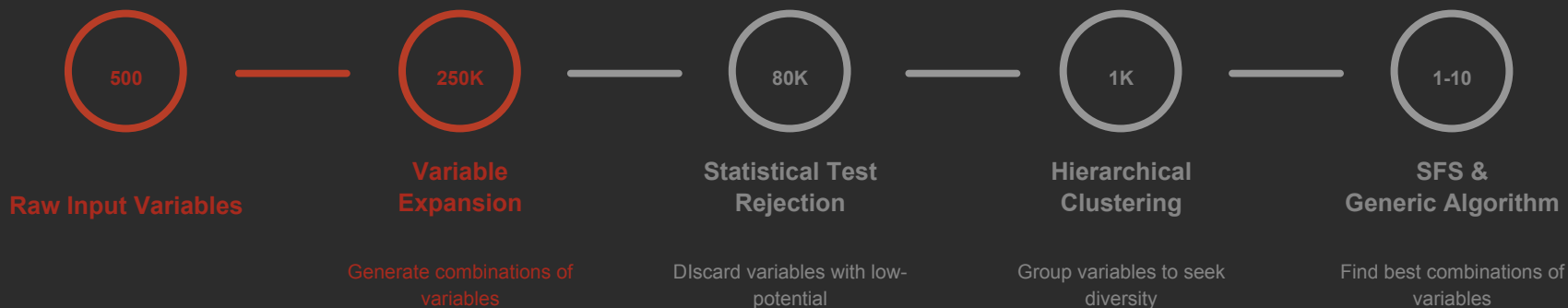- Tools to evaluate embedding methods

# Shortcutting the Workflow…

Business understanding → Analytic approach

Automate (for a specific type of models)

Data requirements → Data collection → Data understanding → Data preparation → Modeling → Evaluation → …

Validation of Data vs Business Question (Quick & Fast)

Validation of Data Ingestion
Validation of Data vs Business Question

Validation of Data Ingestion
Validation of Data vs Business Question

Exploratory Data Analysis

Feature exploration
Objective variables exploration

Model validation
Model interpretation

Model exploration
Model comparison

# Automatic Modeller

Pipeline to seek

- Linear Models
- "Interpretable" variable meanings
- Across multiple metrics (e.g. model quality vs number of features)
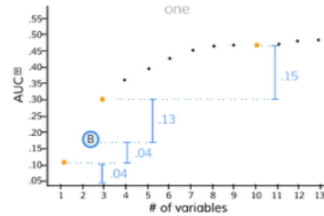
(Number of variables in circles)



| 500 | | 250K | | 80K | | 1K | | 1-10 |

**Raw Input Variables** — **Variable Expansion** — **Statistical Test Rejection** — **Hierarchical Clustering** — **SFS & Generic Algorithm**

Generate combinations of variables

DIscard variables with low-potential

Group variables to seek diversity

Find best combinations of variables

- Tools to make data science projects more efficient
  - Generic Customer Attributes
  - Linear, Interpretable Model Construction

- Real tools available to Data Scientists & Data Engineers
- The philosophy somewhat experimental
  - The "experiment" is ongoing and subject to important checks: e.g. real reusability

# Take-Home Message:
# A Buy-vs-Make Learning

- Commoditization of ML algorithms is a reality
- Speed-up for "mainstream" problems (image classification, text classification)
- Still a long tail of problems need very specific domain knowledge and are not addressed by these tools (e.g. classification of text in bank transactions, pricing, etc)
- Still room for "commoditizing" internally

# Acknowledgements

BBVA

DATA & ANALYTICS

# Thanks!

Questions?

Get in touch at "Office Hours" @ DataEngConf:

**2:15 PM - 3:00 PM**

Or visit **bbvatada.com**