

# TÉCNICAS DE *MACHINE LEARNING* APLICADAS A LA IMPUTACIÓN Y CONTROL DE CALIDAD DE LOS MICRODATOS CONTABLES (POC DE LA CENTRAL DE BALANCES DEL BDE)

Seminario sobre Aplicaciones y Desarrollo de Big Data y Data Science en la  
Banca Central\_CEMLA\_Junio 2021



# ÍNDICE

1. **Introducción**
2. **Trabajo realizado por IIC (Instituto de Ingeniería del Conocimiento)**
  - I. Score de anomalías (detección de outliers)
  - II. Imputación de valores
3. **Análisis de resultados**
  - I. Anomalías
  - II. Imputaciones
4. **Lecciones aprendidas y siguientes pasos**

### **1. Introducción**

### **2. Trabajo realizado por IIC (Instituto de Ingeniería del Conocimiento)**

- I. Score de anomalías (detección de outliers)
- II. Imputación de valores

### **3. Análisis de resultados**

- I. Anomalías
- II. Imputaciones

### **4. Lecciones aprendidas y siguientes pasos**

### Cuestionarios con información contable de las empresas no financieras españolas:

10 ejercicios x 900.000 empresas x 3.000 datos.

### Depurados y clasificados mediante procesos automáticos.

Un 20% se clasifican como no aptos para estudio.

### ¿Puede la IA ayudarnos a mejorar estos procesos?

- Encontrar patrones alternativos para clasificar los cuestionarios: *Caso I. Detección de anomalías.*
- Completar la información omitida: *Caso II. Imputación de valores.*

BANCO DE ESPAÑA Eurosisistema		Número de recepción 10.852.496	2010	1
Central de Balances				
DE NOMINACIÓN SOCIAL: <input type="text"/>				
ANAGRAMA: <input type="text"/>		NIF: <input type="text"/>		
<b>1 DATOS DE IDENTIFICACIÓN</b>				
1 Localización de la empresa				
Domicilio social: <input type="text"/>				
Municipio: <input type="text"/>				
Código postal: <input type="text"/> Provincia: <input type="text"/>				
Persona o servicio a los que la Central de Balances puede dirigirse para efectuar aclaraciones:				
Nombre: <input type="text"/>		Teléfono: <input type="text"/>		
Dirección e-mail: <input type="text"/>		Fax: <input type="text"/>		
Persona o entidad a la que se debe remitir la información de la empresa y estudios (cumplimentar si es distinta de la anterior o si la dirección de envío es distinta del domicilio social):				
Nombre: <input type="text"/>		Teléfono: <input type="text"/>		
Domicilio: <input type="text"/>				
Municipio: <input type="text"/> Código postal: <input type="text"/>				
Provincia: <input type="text"/> Fax: <input type="text"/>				
Dirección e-mail: <input type="text"/>				
2 Estructura de la propiedad (1)				
1 Información sobre participaciones directas en el capital de la empresa				
SOCIEDAD O ACCIONISTA DOMINANTE DIRECTO				
NIF (*)	DE NOMINACIÓN SOCIAL	% PARTICIPACIÓN	NACIONALIDAD	
<input type="text"/>	<input type="text"/>	<input type="text"/>	<input type="text"/>	
OTRAS ACCIONES Y SOCIEDADES CON PARTICIPACIÓN SUPERIOR AL 10%				
NIF (*)	DE NOMINACIÓN SOCIAL	% PARTICIPACIÓN	NACIONALIDAD	
<input type="text"/>	<input type="text"/>	<input type="text"/>	<input type="text"/>	
(*) Completar solo para empresas radicadas en España.				
2 Información sobre participaciones indirectas en el capital de la empresa				
SOCIEDAD DOMINANTE ESPAÑOLA ÚNICA DEL GRUPO				
NIF	DE NOMINACIÓN SOCIAL	% PARTICIPACIÓN		
<input type="text"/>	<input type="text"/>	<input type="text"/>		
SOCIEDADES ESPAÑOLAS RELACIONADAS EN 2) AUTORIDADES POR LAS ADMINISTRACIONES PÚBLICAS O POR EL SECTOR EXTERIOR				
NIF	DE NOMINACIÓN SOCIAL	% PARTICIPACIÓN		
<input type="text"/>	<input type="text"/>	ADMINISTRACIÓN PÚBLICA	SECTOR EXTERIOR	
<input type="text"/>	<input type="text"/>	<input type="text"/>	<input type="text"/>	
(1) Consultar el cuadro I de las normas de cumplimentación, página 22.				

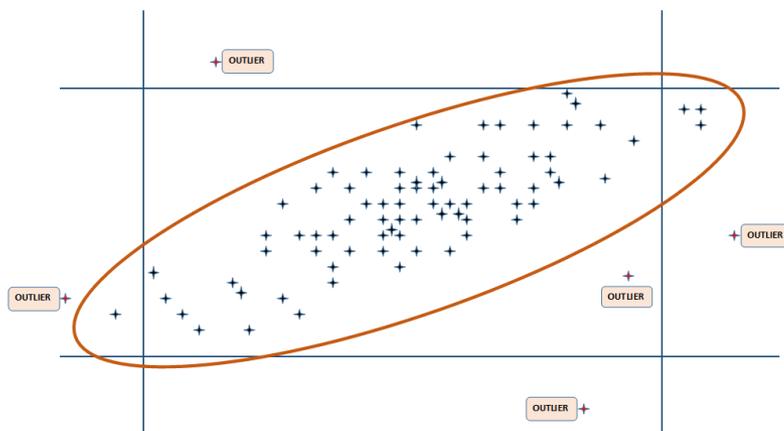
# 1. INTRODUCCIÓN

## Metodología en la POC de 2019



## RECUPERAR CUESTIONARIOS PARA ESTUDIO

### SCORE DE ANOMALÍAS Índice de anomalía valorando n dimensiones



### IMPUTACIÓN DE VALORES en: (i) Descuadres más comunes y (ii) empleo

PATRIMONIO NETO Y PASIVO		NOTAS DE LA MEMORIA	EJERCICIO 2017 (1)
C)	PASIVO CORRIENTE .....	32000	95.200,00
I.	Pasivos vinculados con activos no corrientes mantenidos para la venta .....	32100	4.000,00
II.	Provisiones a corto plazo .....	32200	1.200,00
III.	Deudas a corto plazo .....	32300	5.000,00
1.	Deudas con entidades de crédito .....	32320	
2.	Acreedores por arrendamiento financiero .....	32330	
3.	Otras deudas a corto plazo .....	32390	
IV.	Deudas con empresas del grupo y asociadas a corto plazo .....	32400	2.000,00
V.	Acreedores comerciales y otras cuentas a pagar .....	32500	
1.	Proveedores .....	32580	
a)	Proveedores a largo plazo .....	32581	
b)	Proveedores a corto plazo .....	32582	

- **Selección de variables:** 94 claves contables + clave empleo + 2 campos de sector de actividad (Sector y Gran Sector).
- **Normalización contable:** Dividir los campos de *Perdidas y Ganancias* entre el *Importe neto de la cifra de negocios*. El *Balance* entre el *Total Activo*.
- **Filtro de cuestionarios:** Subtipo reducido de 2008 a 2017, descartando los *No normalizables* (*Importe neto de la cifra de negocios = 0*; ~2 millones instancias).
- **Generar nuevas variables:** Medias de cada valor en los últimos 2-5 años, número de sectores declarados, edad de la empresa...
- **Separar cuestionarios según su calidad:**
  - Perfectos (5,323,000)
  - Baja calidad (476,000)
  - *Missing* (469,000)

### 1. Introducción. Alcance de la iniciativa de 2019

### 2. Trabajo realizado por IIC (Instituto de Ingeniería del Conocimiento)

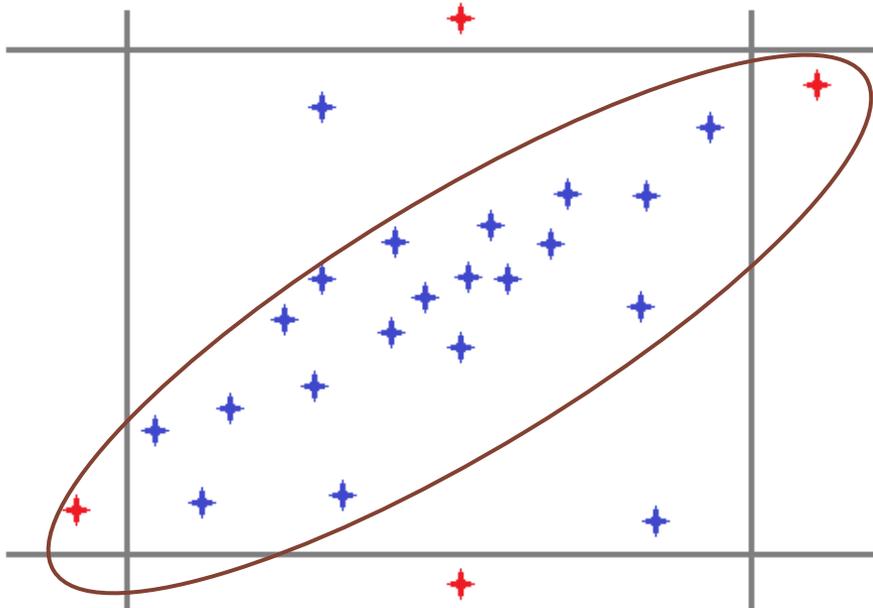
1. Score de anomalías (detección de outliers)
2. Imputación de valores

### 3. Análisis de resultados

1. Anomalías
2. Imputaciones

### 4. Lecciones aprendidas y siguientes pasos

### SCORE DE ANOMALÍA

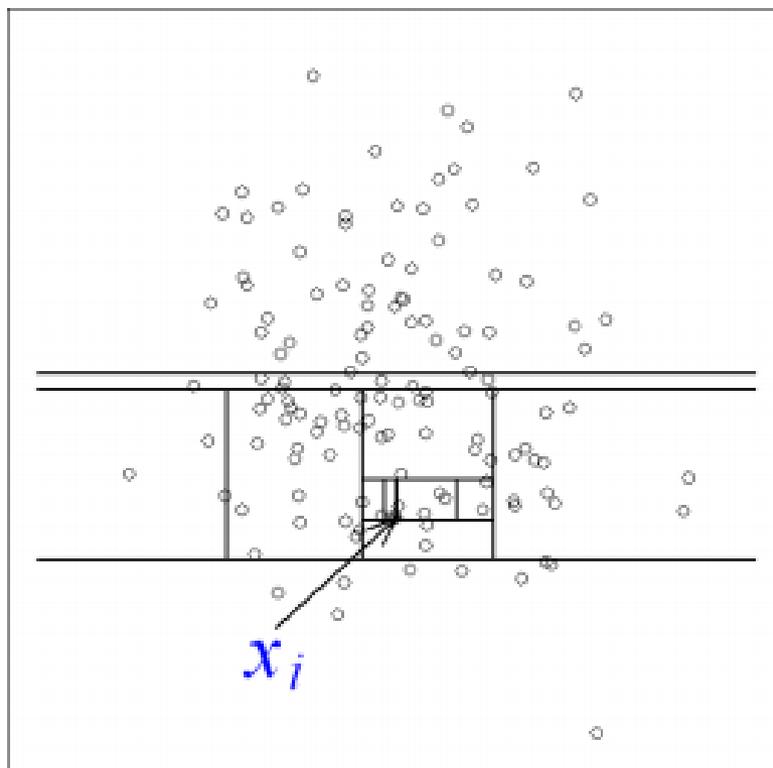


Cálculo de *score* de anomalía  $[0, 1]$   
vs. detección de *outliers* (Sí/No).

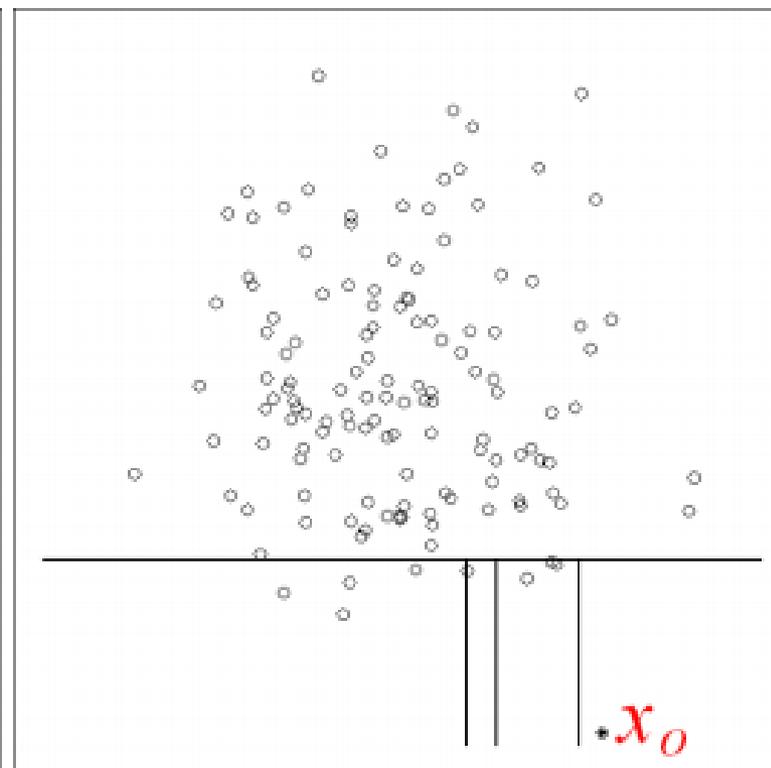
Algoritmo empleado: **Isolation Forest.**

### ISOLATION FOREST

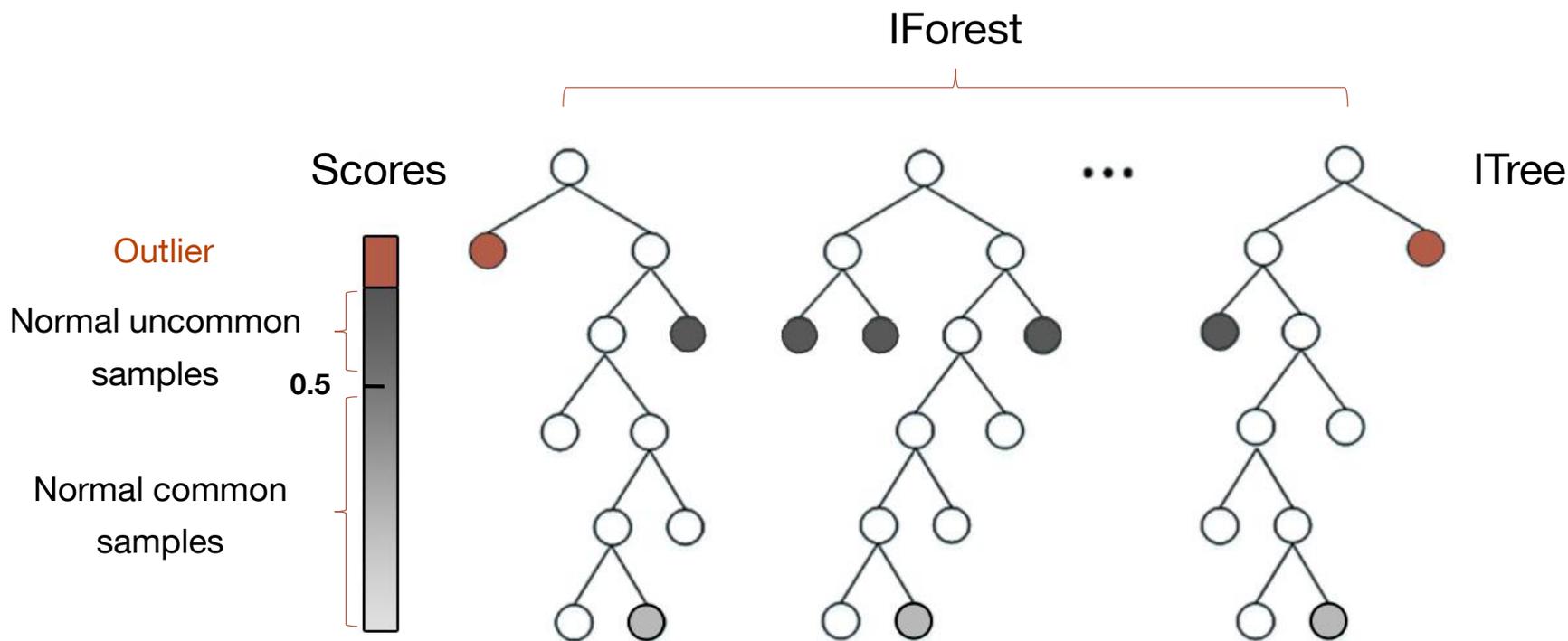
Las instancias anómalas se aíslan fácilmente mediante divisiones aleatorias del espacio.



(a) Isolating  $x_i$



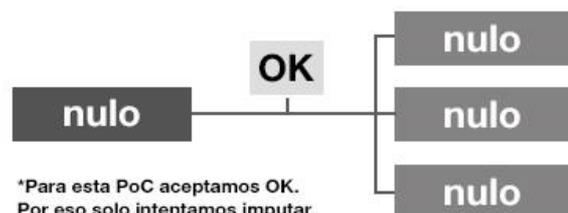
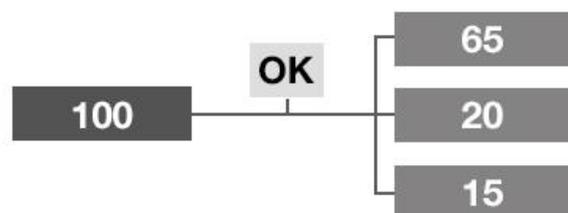
(b) Isolating  $x_o$



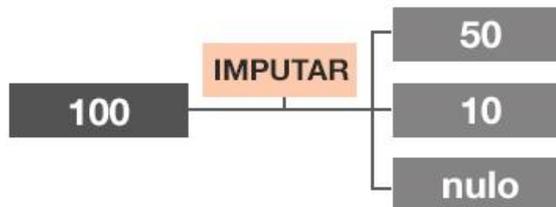
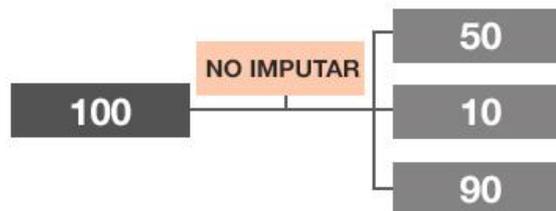
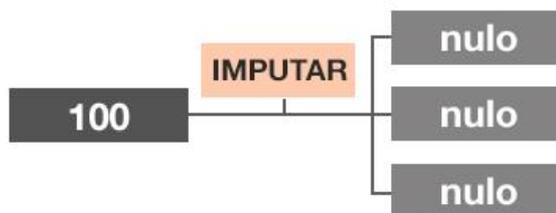
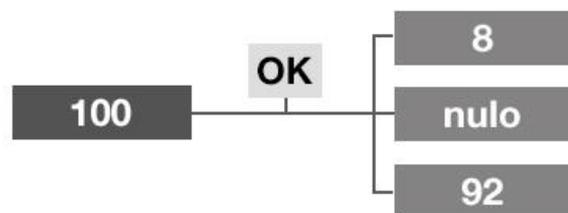
Entrenamiento con **5.000 árboles** y todos los cuestionarios **perfectos**: 5.300.000

Se evalúa sobre los cuestionarios perfectos, de baja calidad y missing de 2017.

### Imputación de valores



\*Para esta PoC aceptamos OK.  
Por eso solo intentamos imputar sumandos que no se formen con sub-sumatorios.



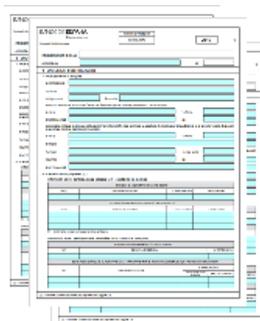
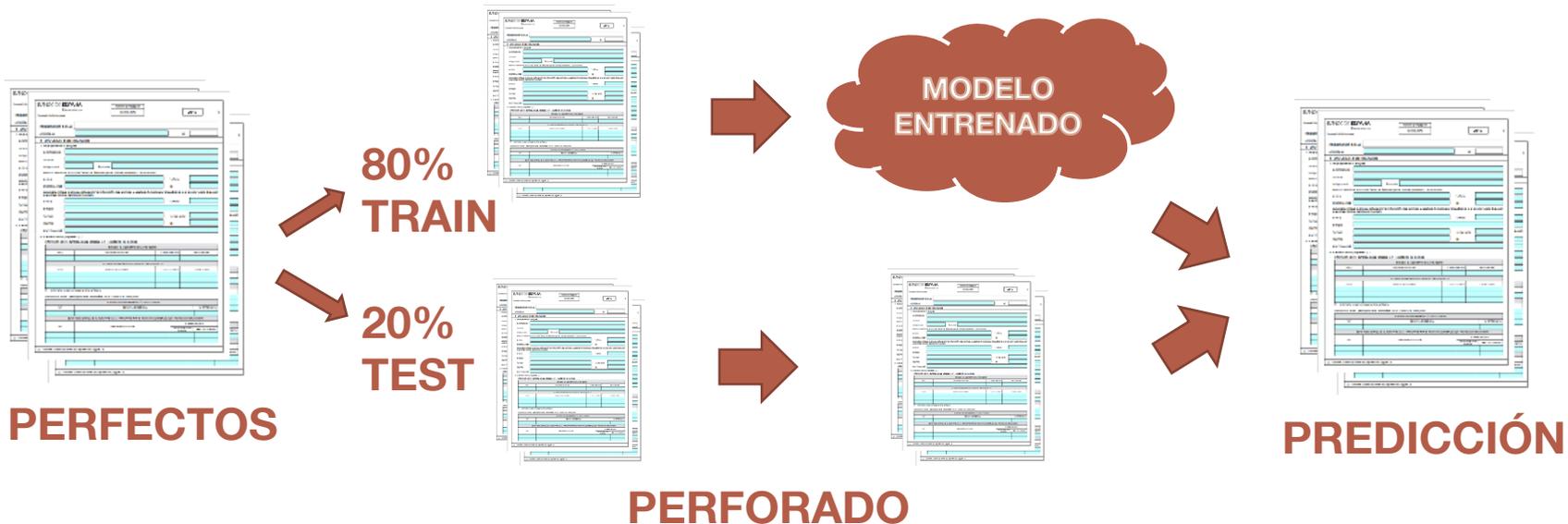
\* Si el sumatorio no es igual a la suma de sumandos: IMPUTAR sumandos.  
Restricción: Tras la imputación deben cuadrar.  
Holgura del cuadro: +/- 2%

Selección de datos:

- Descuadres más comunes (4 ejercicios)
- Empleo (1 ejercicio)

## 2.II. IMPUTACIONES

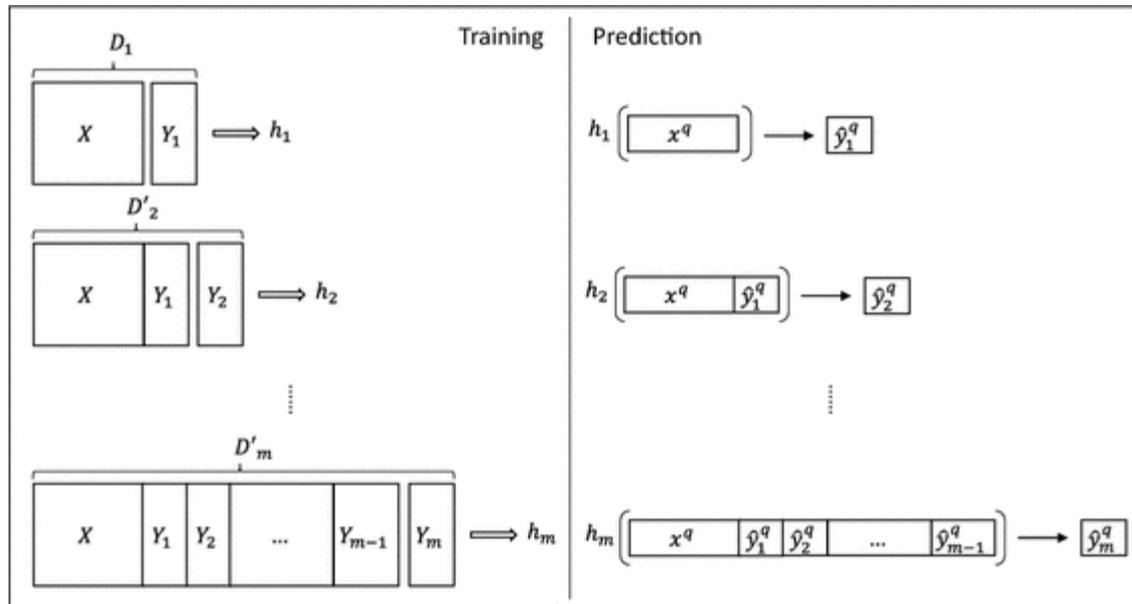
### Entrenamiento, perforado y predicciones



**MISSING**

ID_INST_CUEST	12300	12380	12370	12390	31200	31220	31230	31290	32300
13623376763	NaN	NaN	NaN	NaN	76.00000	NaN	NaN	NaN	-0.04868
13623376769	NaN	NaN	NaN	NaN	23.47207	NaN	NaN	9.63073	0.57878
13623376773	2.67538	2.02641	NaN	0.64897	3.11943	NaN	NaN	NaN	-6.45118
13623376777	91.91667	91.91667	NaN	NaN	9.31272	NaN	NaN	NaN	16.13428

**Ensamble of Regressor Chains (ERC):** Construir varios modelos de regresión de forma incremental. Cada modelo predice una variable que luego es empleada para entrenar el siguiente modelo.



**Train:** 240.000

**Test:** 60.000

Modelo de regresión:

Random forests 1.000 árboles

Computacionalmente muy costoso.

El orden de predicción de las variables (cadena) teóricamente afecta al resultado, dando mayor peso a las primeras variables elegidas. Se prueban 5 cadenas aleatorias.

1. **Introducción. Alcance de la iniciativa de 2019**
2. **Trabajo realizado por IIC (Instituto de Ingeniería del Conocimiento)**
  1. Score de anomalías (detección de outliers)
  2. Imputación de valores
3. **Análisis de resultados**
  1. Anomalías
  2. Imputaciones
4. **Lecciones aprendidas y siguientes pasos**

# 3.I. ANÁLISIS DE LOS RESULTADOS

ANOMALÍAS. Scoring IIC vs calidad CB: Distribución de los datos

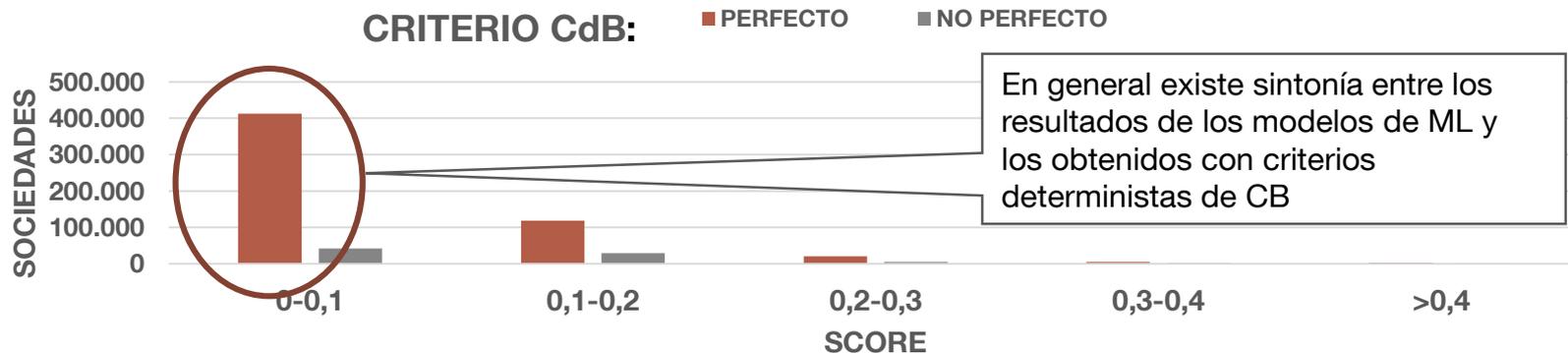
¿Falsos positivos? Analizar para detectar posibles mejoras en nuestros sistemas de filtrado.

¿Falsos negativos? Analizar para en su caso 'relajar' nuestros sistemas de filtrado

¿El 94% de los cuestionarios se concentran en un rango de anomalía entre 0 y 0,2

## CALIDAD CUESTIONARIOS CBB 2017

Scoring IIC (0=Bueno; 1=Malo)	PERFECTO	NO PERFECTO	TOTAL	% Total acumulado
0-0,1	411.973	41.626	453.599	71,3%
0,1-0,2	118.439	28.942	147.381	94,4%
0,2-0,3	20.380	5.404	25.784	98,5%
0,3-0,4	5.154	1.377	6.531	99,5%
>0,4	2.299	853	3.152	100,0%
<b>TOTAL</b>	<b>558.245</b>	<b>78.202</b>	<b>636.447</b>	



## 3.I. ANÁLISIS DE LOS RESULTADOS

Anomalías. ¿Por qué deberíamos fiarnos del score?



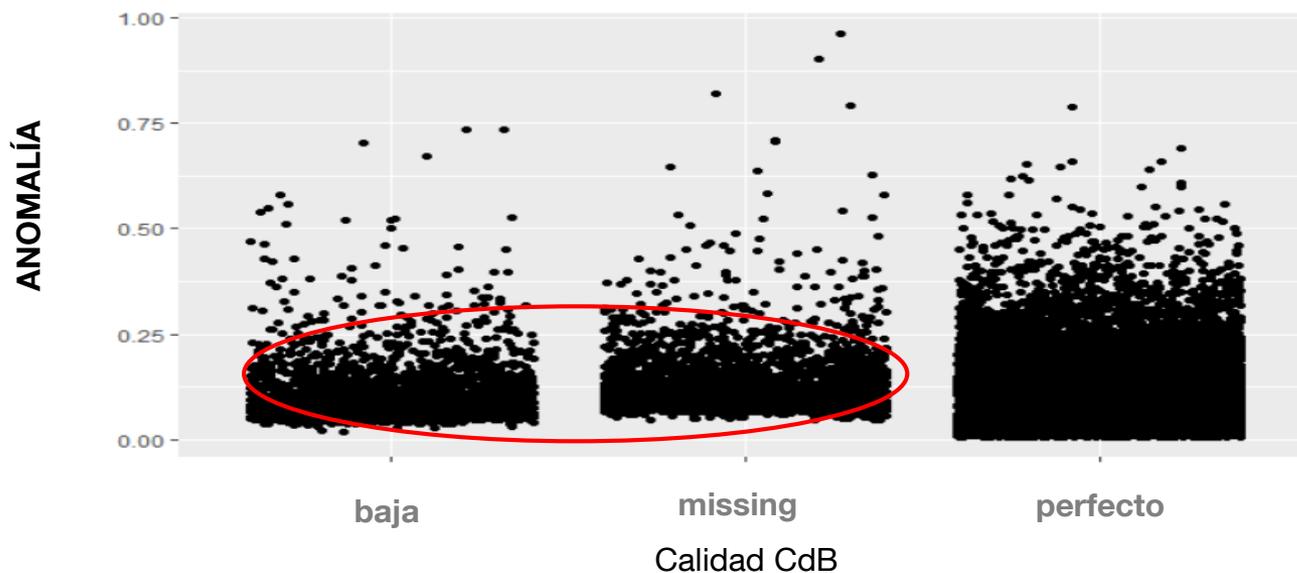
### En resumen:

Aceptando este score...	...en neto ganamos o perdemos sociedades	...renunciando a estas...	...e incorporando estas ...
0,1	-104,646	-104,272	41,626
0,2	42.735	-27.833	70.568
0,3	68.519	-7.453	75.972
0,4	75.050	-2.299	77.349

# 3.I. ANÁLISIS DE LOS RESULTADOS

ANOMALÍAS. Falsos negativos (según CdB): características cuestionarios a ganar

## SCORING DE ANOMALÍA SEGÚN IIC VS NIVEL DE CALIDAD SEGÚN CENTRAL DE BALANCES



¿Complementar con el score de anomalías para recuperar empresas?

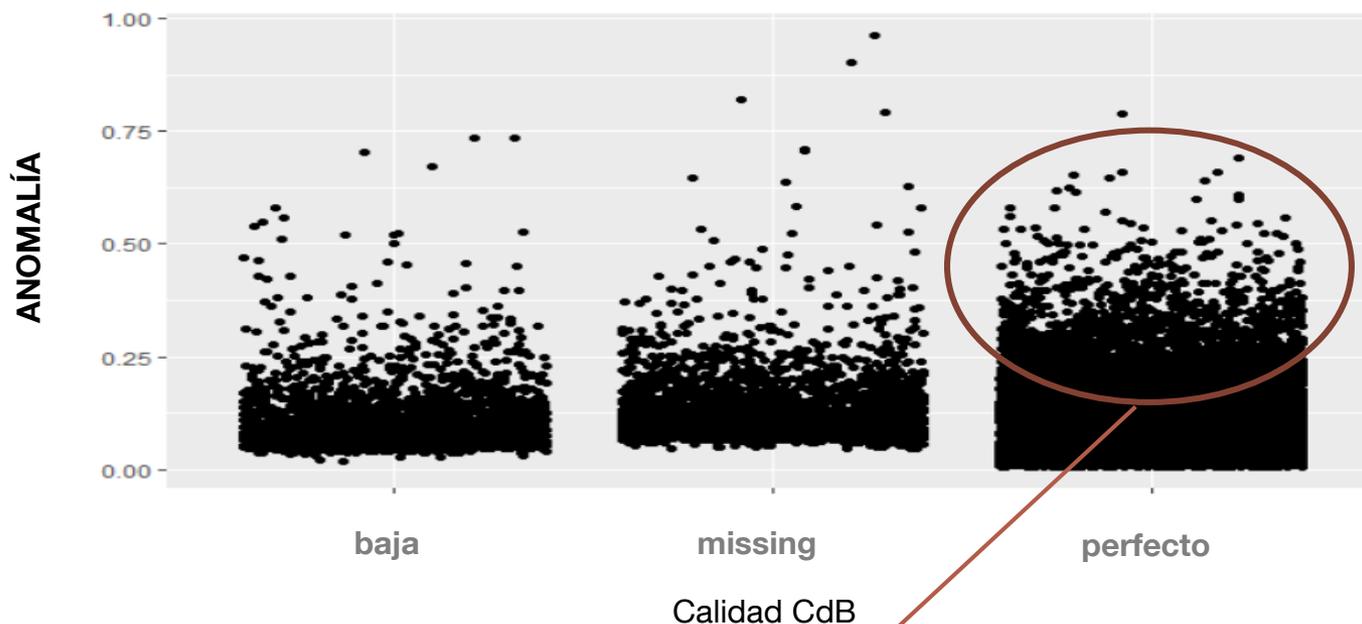
¿Recuperables mediante imputaciones?

SCORE	TOTAL CUESTIONARIOS	NO CUADRADA	UNIDADES NO FIABLES	INCOMPLETA	NO CUMPLE CONDICION VAREXCESIVA	PERSONAL NO COHERENTE
hasta 0,05	2.242	438 20%	96 4%	193 9%	1.462 65%	67 3%
hasta 0,1	41.626	6.092 15%	2.733 7%	2.503 6%	13.268 32%	20.833 50%
hasta 0,15	62.482	10.426 17%	4.128 7%	3.420 5%	18.506 30%	32.915 53%
hasta 0,2	70.568	12.821 18%	4.620 7%	3.724 5%	20.720 29%	37.077 53%

### 3.I. ANÁLISIS DE LOS RESULTADOS

ANOMALÍAS. Falsos positivos (según CdB): potenciales cuestionarios a perder

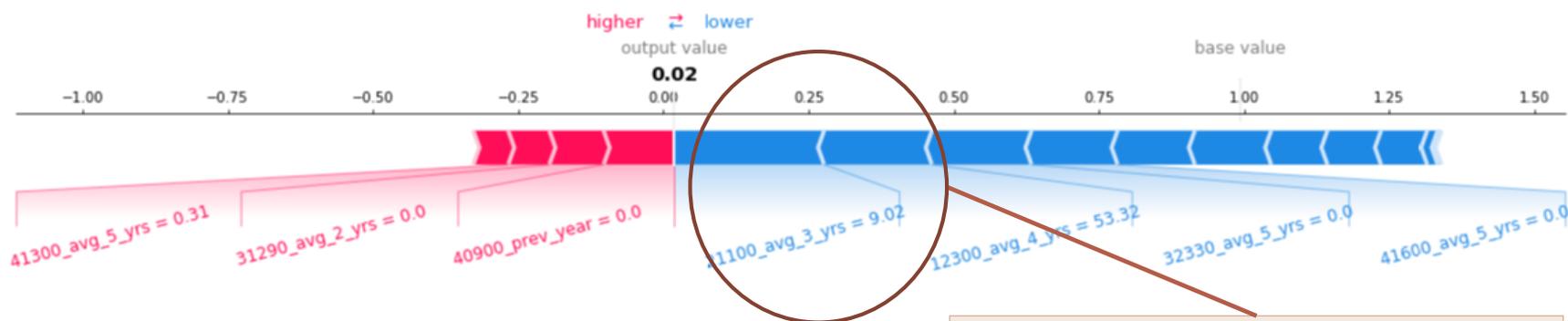
#### SCORING DE ANOMALÍA SEGÚN IIC VS NIVEL DE CALIDAD SEGÚN CENTRAL DE BALANCES



Intentando entender el algoritmo... lo primero, bajar al microdato

## CONTRIBUCIONES AL ÍNDICE DE ANOMALÍA: RATIOS DE SHAPLEY

Indican qué claves han contribuido más al índice de anomalía



Significa que la variable “*media de los 3 últimos años de la clave 21100*”, al tener un valor de 9,2, **reduce** el **escore** 0,26 puntos aprox.

Interpretación individual: poco útil en los casos revisados

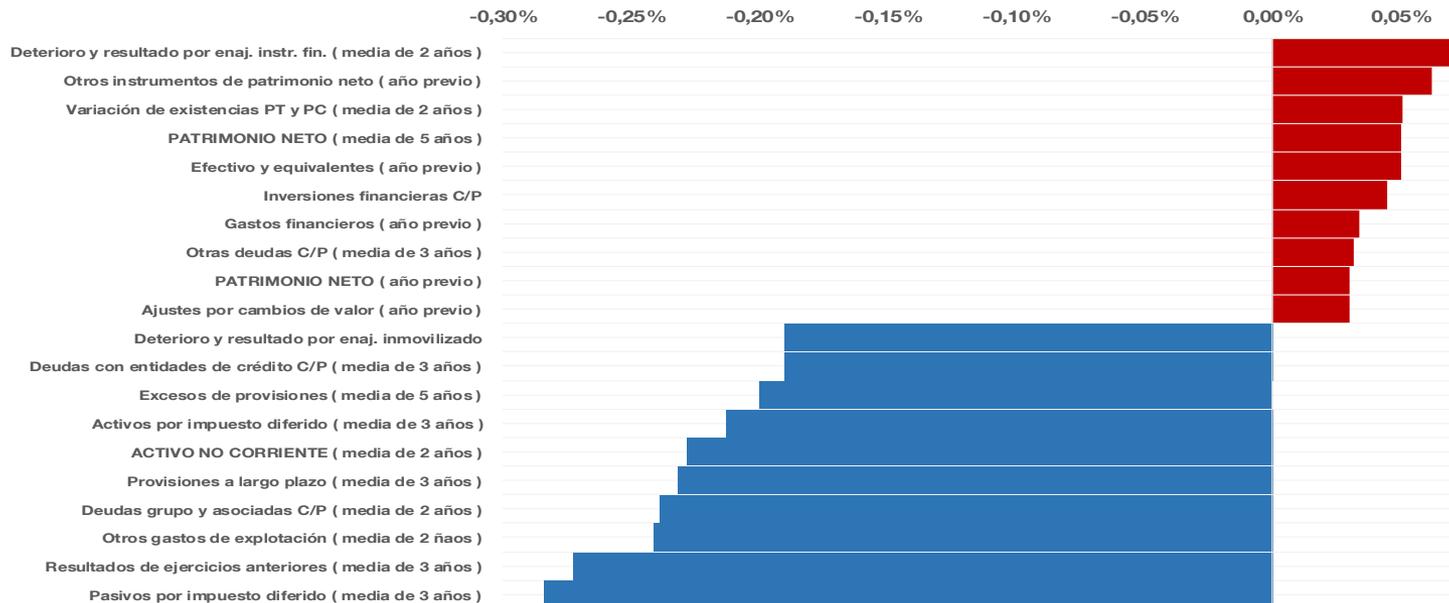
Propiedad aditiva para analizar un nodo concreto

# 3.I. ANÁLISIS DE LOS RESULTADOS

## ANOMALÍAS: Ratios de shapley agregadas

Las ratios de Shapley **AGREGADAS** permiten interpretar cuál es el efecto de cada variable en el **conjunto de sociedades** que se elija (en este caso son aquellas con scoring alto y provenientes de empresas perfectas; pero a efectos de la CdB se podría elegir un nodo, un tamaño de empresa concreto, una comunidad autónoma, o cualquier otro conjunto).

**MEDIA DE LAS RATIOS DE SHAPLEY PARA EL GRUPO CON SCORE ALTO**  
VALORES DE SHAPLEY EN ROJO: EMPEORAN EL SCORING DE ANOMALÍAS  
VALORES DE SHAPLEY EN AZUL: MEJORAN EL SCORING DE ANOMALÍAS



Herramienta potente para investigar pero los datos disponibles no lo permiten actualmente

### 1. Introducción. Alcance de la iniciativa de 2019

### 2. Trabajo realizado por IIC (Instituto de Ingeniería del Conocimiento)

1. Score de anomalías (detección de outliers)
2. Imputación de valores

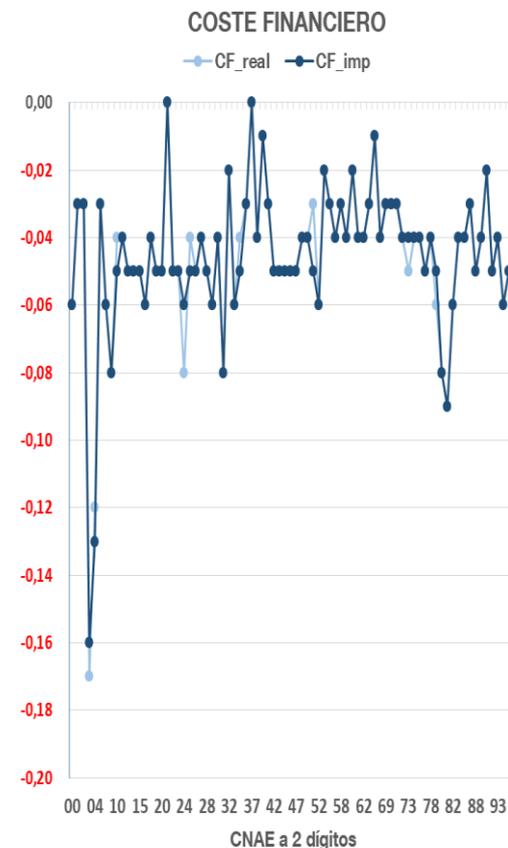
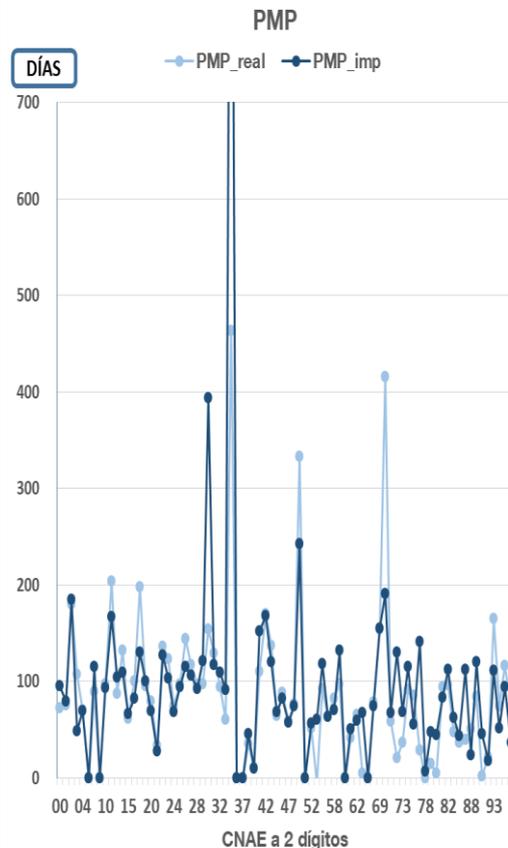
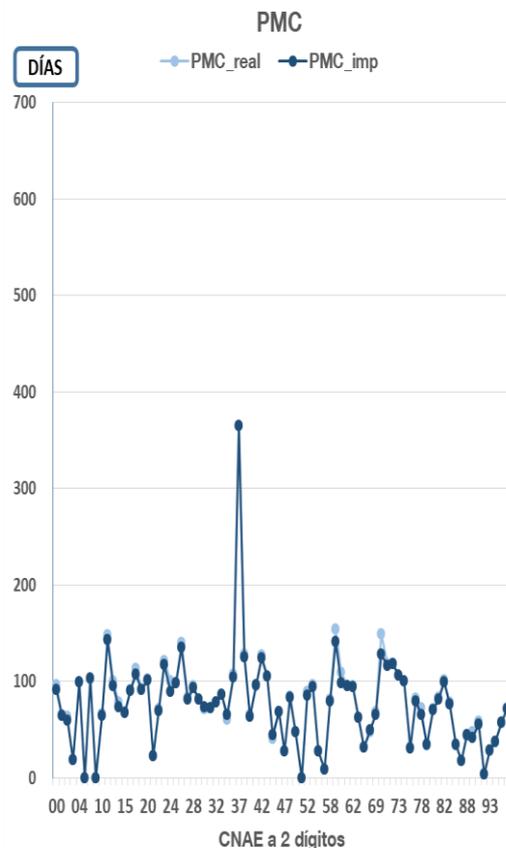
### 3. Análisis de resultados

1. Anomalías
2. Imputaciones

### 4. Lecciones aprendidas y siguientes pasos

## 3.II. ANÁLISIS DE LOS RESULTADOS

### IMPUTACIONES: Análisis de los periodos medios de cobro y pago, por CNAE

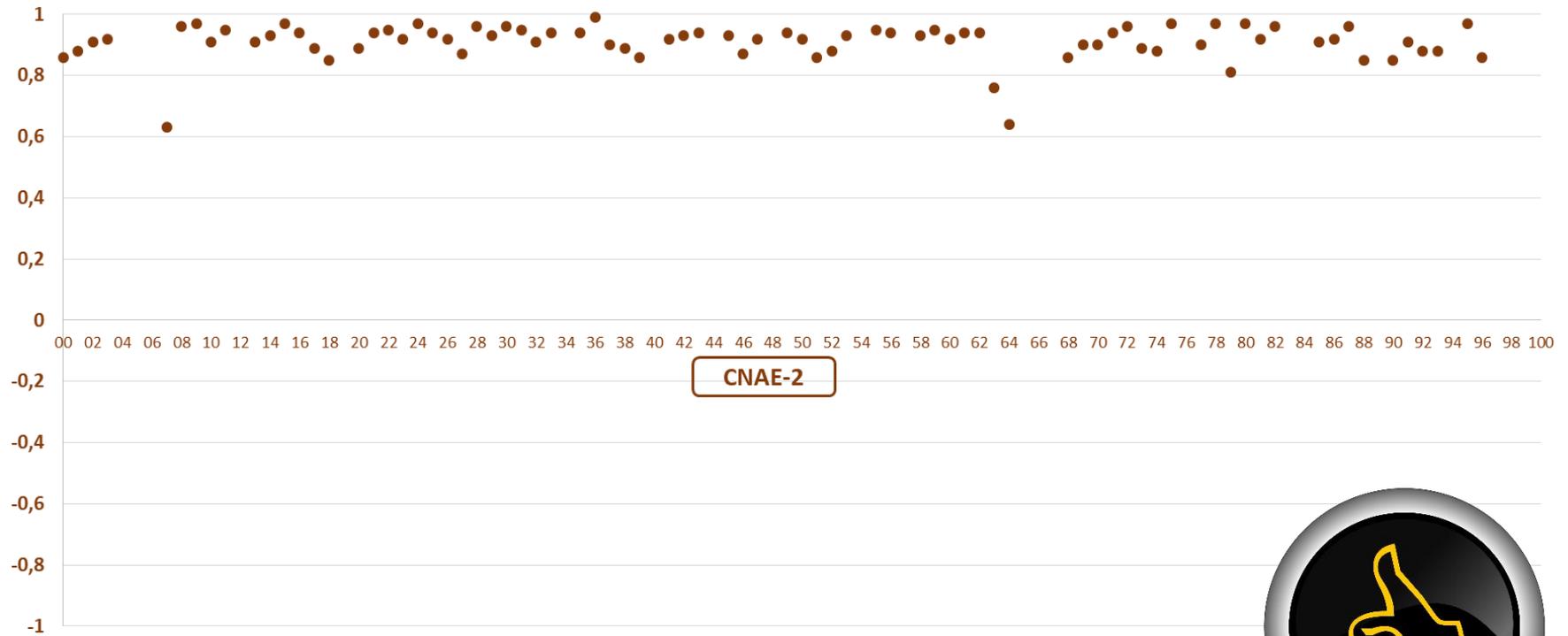


Las correlaciones son aceptables para PMC y Coste financiero, pero bajan para PMP, quizá porque se han realizado menos imputaciones en la clave de proveedores.

## 3.II. ANÁLISIS DE LOS RESULTADOS

### IMPUTACIONES: Empleo imputado vs empleo real

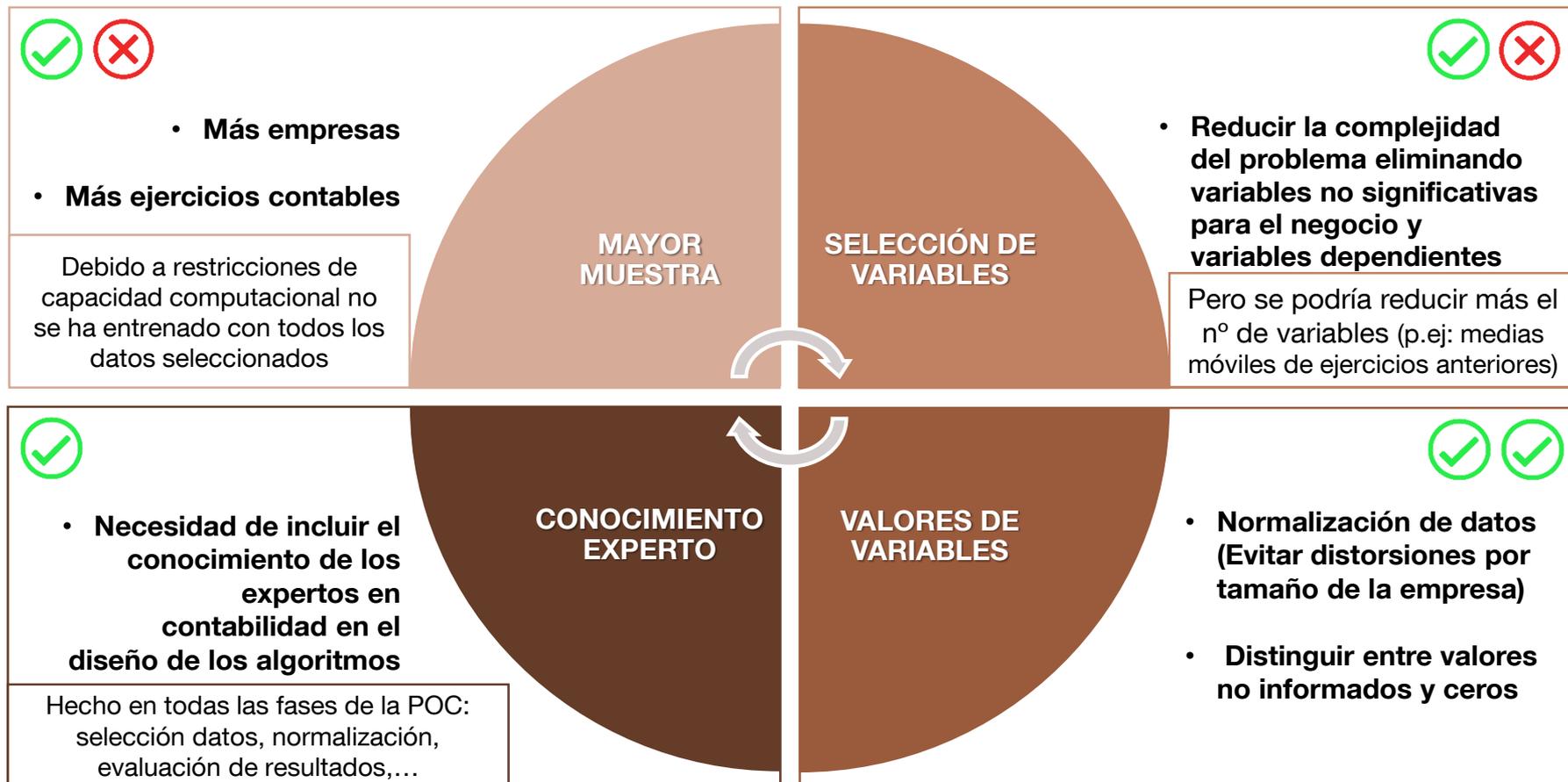
CORRELACIÓN ENTRE EMPLEO REAL Y EMPLEO IMPUTADO CBB 2017



- 1. Introducción. Alcance de la iniciativa de 2019**
- 2. Trabajo realizado por IIC (Instituto de Ingeniería del Conocimiento)**
  1. Score de anomalías (detección de outliers)
  2. Imputación de valores
- 3. Análisis de resultados**
  1. Anomalías
  2. Imputaciones
- 4. Lecciones aprendidas y siguientes pasos**

# 4. LECCIONES APRENDIDAS Y SIGUIENTES PASOS

Qué mejora la POC de 2019 respecto del Piloto de 2018



#### Para dar validez al **score de anomalías** es necesario:

- Disponer de ratios de **Shapley agregadas *customizadas*** a las necesidades de negocio (determinado sector, tamaño,...)

#### Para dar validez a las **imputaciones** es necesario:

- Disponer de ratios de **Shapley para imputaciones** y no sólo para anomalías
- Revisar el **patrón de perforado** en el conjunto de test (p.ej: clave de proveedores con pocos datos imputados)
- Probar a repetir las **imputaciones después de la eliminación de cuestionarios anómalos**

**En definitiva, más análisis...**

*Thank You!*

GRACIAS POR SU ATENCIÓN



**Análisis de las 5 sociedades del CNAE-2 como ejemplo de caso que presenta una diferencia elevada entre ratio PMP con valores reales y PMP con valores imputados.**

EMPRESA	EJERCICIO	PMP_real	PMP_imp	Aprovisionamientos	Proveedores real	Proveedores imputado
1558271	2012	239,6	198,5	-15,5	10,1	8,4
1964308	2013	86,6	345,4	-83,3	19,8	78,8
2007712	2014	126,9	1378,3	-95,2	33,1	359,4
2809470	2015	54,9	93,9	-292,2	44	75,2
995966	2017	414,6	398,6	-136,9	155,5	149,5
TOTAL CNAE_2 = 30		153,8	491,5	-623,1	262,5	671,3

