

# Ciencia de Datos para la Toma de Decisiones I

Mtro. René Rosado González  
Director de Programa LTP

# Nuestro Reto:



Desarrollar una metodología basada en datos para la generación de un índice que permita entender los procesos de desigualdad que ocurren en las comunidades de origen o retorno de las personas privadas de la libertad.

Nos reuniremos con el socio formador los días miércoles



# Design Sprint

day 1



## understand

- who are the users
- what are their needs
- what is the context
- competitor review
- formulate strategy

2



## diverge

- envision
- develop lots of solutions
- ideate

3



## decide

- choose the best idea
- storyboard the idea

4



## prototype

- build something quick and dirty to show to users
- focus on usability not making it beautiful

5



## validate

- show the prototype to real users outside the organisation
- learn what doesn't work

data  
data  
data  
data  
data

SMALL ~~big~~ data bandwidth ~~QUALITY~~

ect infallible ~~data~~

ive impartial ~~data~~

NG descriptive ~~data~~

US predictive ~~data~~

data conventions ~~POSSIBILITIES~~

data to simplify complexity ~~DEPICT~~

data processing ~~DRAWING~~

data driven design

VD save time with ~~data~~

data is numbers ~~PEOPLE~~

data will make us more efficient ~~HUMAN.~~

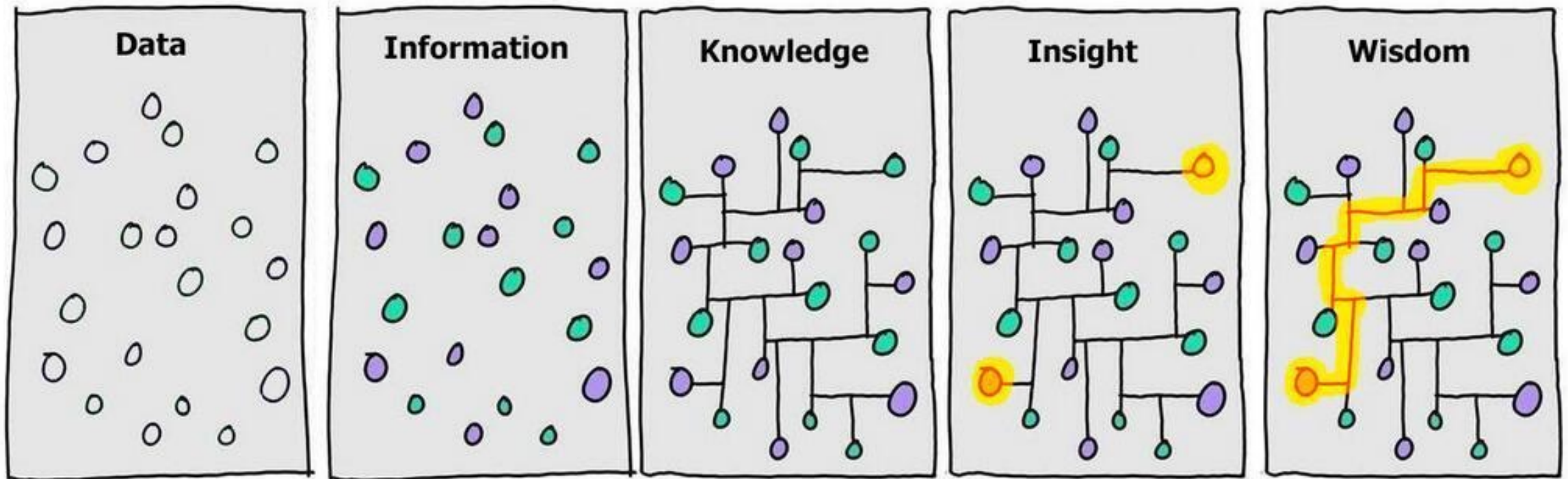
a manifesto for a  
**DATA HUMANISM**

a new Renaissance where we can question the impersonality of a merely technological approach to DATA, where we are ready to reconnect numbers to what they really stand for: which ARE - MORE and MORE - OUR unique lives.

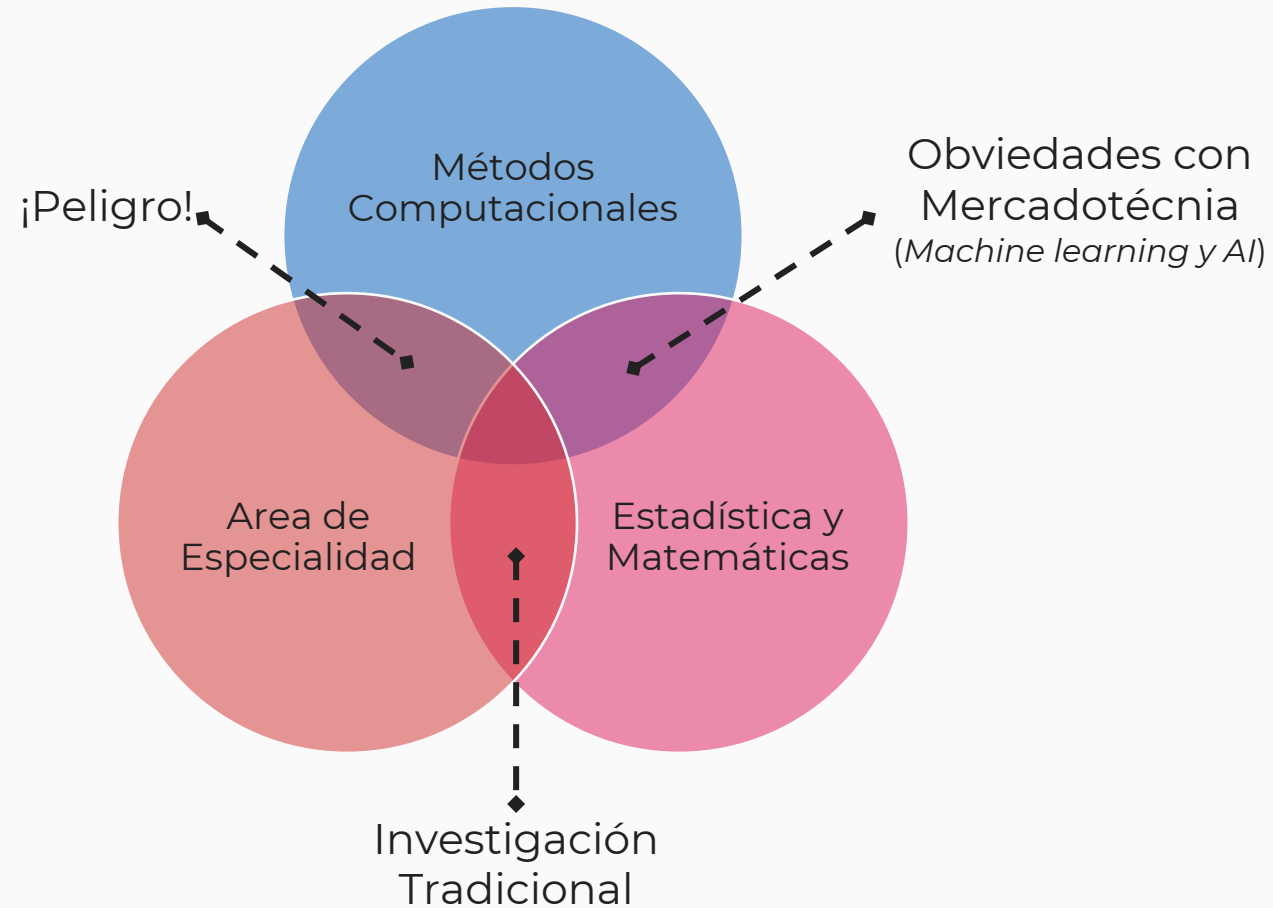
Una expresión de deseo por almacenar una vivencia

giorgia lupi

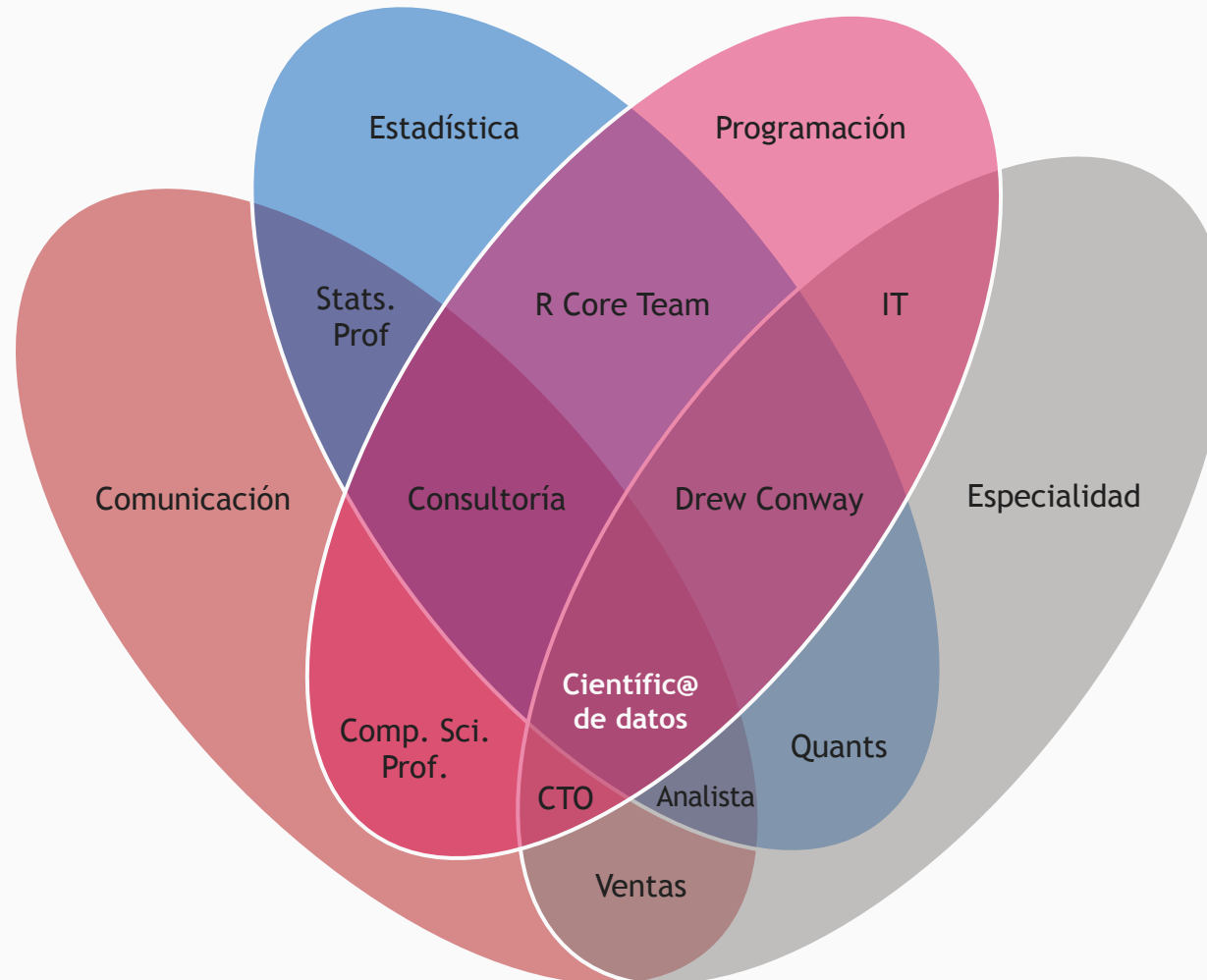
# ¿Qué hace la ciencia de datos?



# El diagrama de Drew Conway



# “There is no I in Team”



# Cientific@ de Datos Ideal





MATHEMATICA @50:

AT THE **INTERSECTION** OF

# DATA SCIENCE AND SOCIAL SCIENCE

13 SEPT 2018

Building on our History  
of Sharing Research

PAUL DECKER

WE CAN BETTER DRIVE **POLICY MANAGEMENT**

DATA IS  
... BEING **DEMOCRATIZED**  
... **MANY TYPES**  
... **GREATER FLOW**

WE MUST LEARN TO **ADAPT**

MATHEMATICA

Policy Research



ASSOCIATION FOR  
PUBLIC POLICY ANALYSIS  
& MANAGEMENT

EXAMPLES AT THE INTERSECTION

MATT SARGANIK

THE ARTIST  
D & CHAMP  
A LESSON IN  
RESHAPING  
FOR A PURPOSE

MICHAELANGELO  
CUSTOMIZATION  
TO FIT A NEED  
CHARTING & MEASURING  
POVERTY IN RWANDA

SHEILA DUGAN

EXAMPLES  
FROM CITIES  
• WIC PROGRAM  
• CHARLOTTE, NC

WILL YANG

CANCER RESEARCH  
AND IMPACTS OF AI

XAVIER HUGHES

MOBILITY DATA

TO GUIDE DEVELOPMENT DECISIONS

LEVERAGE ANY DATA THAT'S COMING IN!

DATA SCIENCE

IS NOT EMBRACED  
BY EVERYONE

XAVIER:

DATA MAPPING EXPERIMENT  
STARTS IN DECEMBER!

Data points  
accessible to  
Local Municipalities!

Predictive  
Tools in  
Cities  
• K.C., MISSOURI  
... PROPERTY CODE  
VIOLATIONS

NEW  
TOOLS

Sorting  
Comments  
to Modify  
Public  
Policy

SME

USED WITH CITY DEMOGRAPHICS

Post-Production  
Surveillance of  
Consumer Products

Research  
abstract  
analysis  
• NATURAL LANGUAGE SELECTION

Randomized  
experiments  
in normal processes

RAISES ETHICAL QUESTIONS

JUST BECAUSE WE CAN  
DOESN'T MEAN WE SHOULD!

DATA QUALITY  
& MANAGEMENT

Ready-Made  
Custom-Made  
CAN LEVERAGE  
EACH OTHER  
CAN PERFORM  
MOCK EXPERIMENTS  
AND SAVE MONEY

SYNERGIES

Simple, easy  
ways for public  
to use your data

THERE'S A DEMAND  
FOR DATA SCIENCE

... FROM THE BOTTOM UP  
... HHS CO LAB

ADDRESS REAL CHALLENGES  
... WHAT DO CITIES & PEOPLE

MEET IN THE MIDDLE  
... PEOPLE WITH DIFFERENT  
SKILLS & KNOWLEDGE  
... CAN DRIVE EFFICIENCY

SOCIAL SCIENTISTS HAVE A LOT  
TO CONTRIBUTE & LEARN  
... CAN'T IGNORE THEM

ROLE PLAY TO KNOW  
EACH OTHERS' ROLES

... TAKE TIME TO  
UNDERSTAND  
... WE NEED TO BE LIKE  
CHOCOLATE CAKE  
... MANY INGREDIENTS

DATA SCIENTISTS ARE  
LIKE F1 DRIVERS ... 230 MPH!

MANAGING  
NON-RESPONSE  
RATES

LOOKING AT  
SYNERGIES  
... LOOK AT THOSE  
OUTSIDE THE RANGE  
TO IMPROVE MODELS

WORKING  
WITH FED.  
GOVT.

... INCENTIVES MUST  
BE ALIGNED  
... GOVT. CONSTRAINTS ARE  
OFTEN UNREALISTIC  
... CHALLENGE IS MEASURING  
WHAT'S IMPORTANT

FED GOVT. + STATE  
GOVT.  
... TAKES COURAGE  
TO CONNECT  
ARE ALIGNED  
BE

... DON'T START WITH  
THE TOOL OR TECH...  
... DEFINE THE  
PROBLEM  
... EXPOSE  
YOURSELF  
TO YOUR PROBLEM  
SET

CYBER  
SECURITY  
... FITS IN MANY WAYS  
... CONSIDER ETHICAL  
AND PRIVACY CONCERNS

AVOIDING  
UN DUE INFLUENCE  
... CREATE A PLAN  
AND GAIN BUY-IN  
FOR OPEN DATA  
AND OPEN CODE  
... INVOLVE YOUR  
PARTNERS & DRIVE  
TOWARD TRANSPARENCY  
... SEE ACF FOR  
EVALUATION GUIDELINES  
... SEE THE LAB IN DC

HUMANIZE OUR RESEARCH

... TO GET  
ACCURATE  
DATA

EXPLAIN  
YOUR QUESTION

... especially  
IN  
2018

WHY ARE YOU  
ASKING THAT?

I KNOW THE  
SUBJECT!

I CAN TEACH  
YOU ABOUT  
DATA SCIENCE!

TRAIN PEOPLE & HELP THEM BE CONVERSANT

DEMISTIFY  
THE SCIENCE

TRAINING

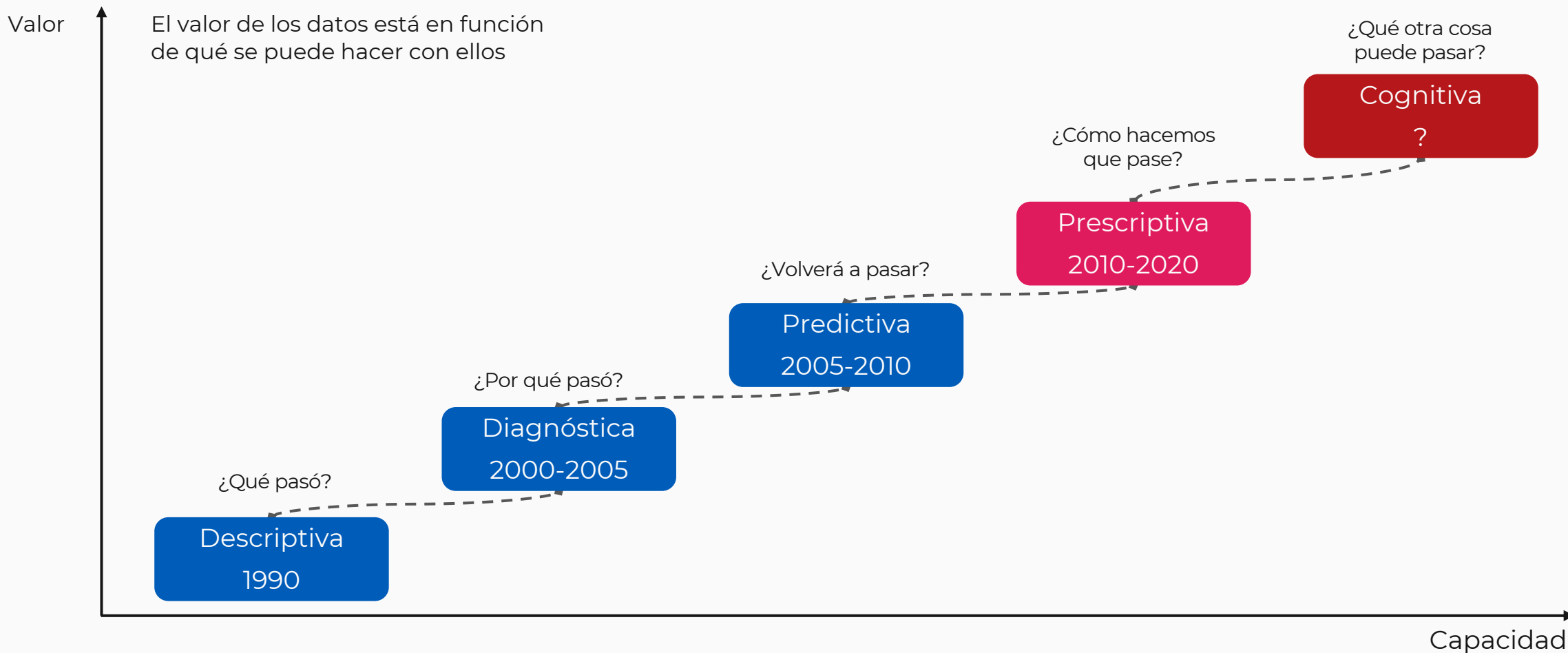
@ MATH POLRESEARCH

# data x social

@ APPAM\_DC

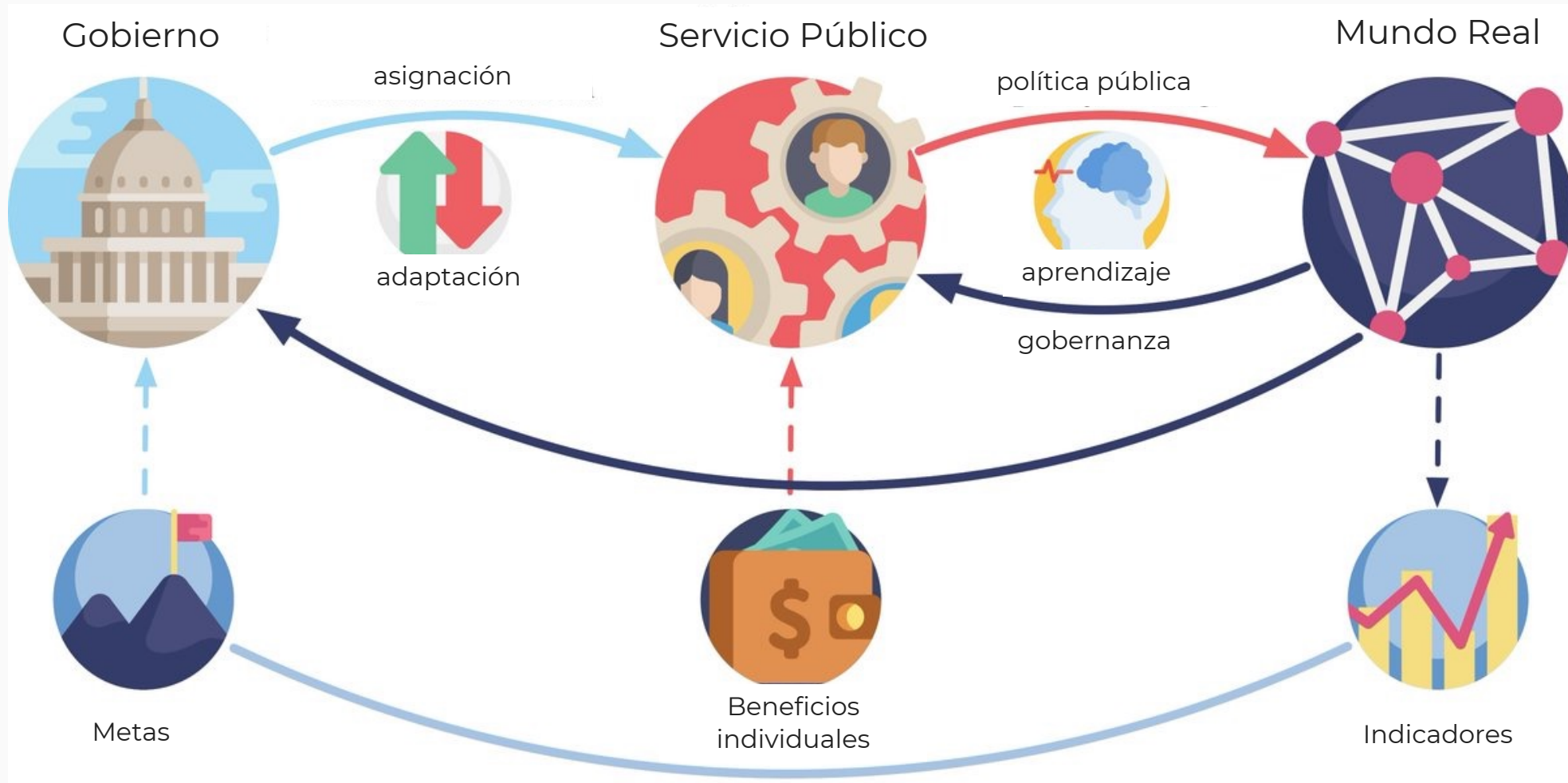


# El valor de los datos en la política pública





# Diseño e implementación de política pública



# Policy Priority Inference

for sustainable development

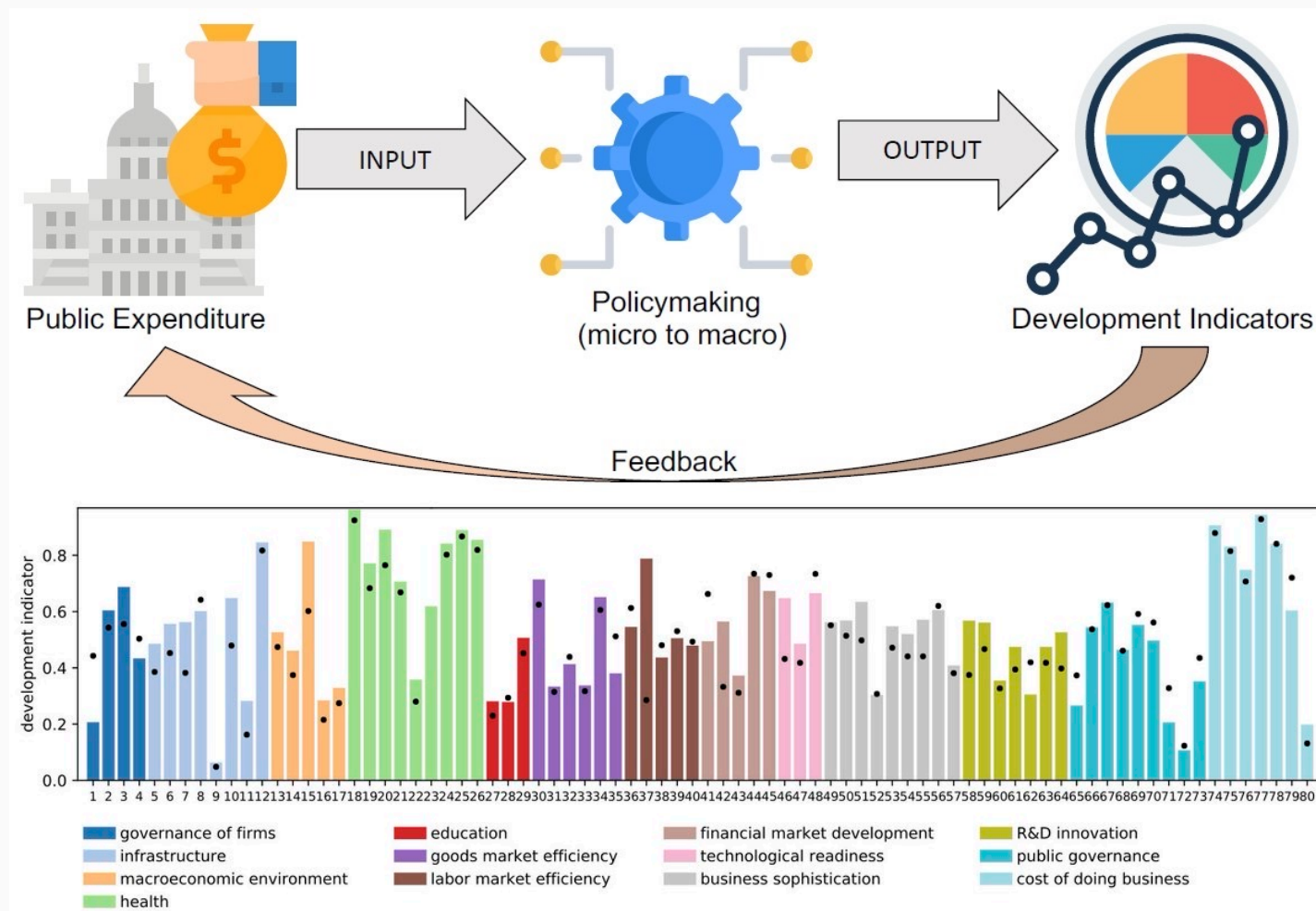


EXATEC

PhD. Omar A. Guerrero

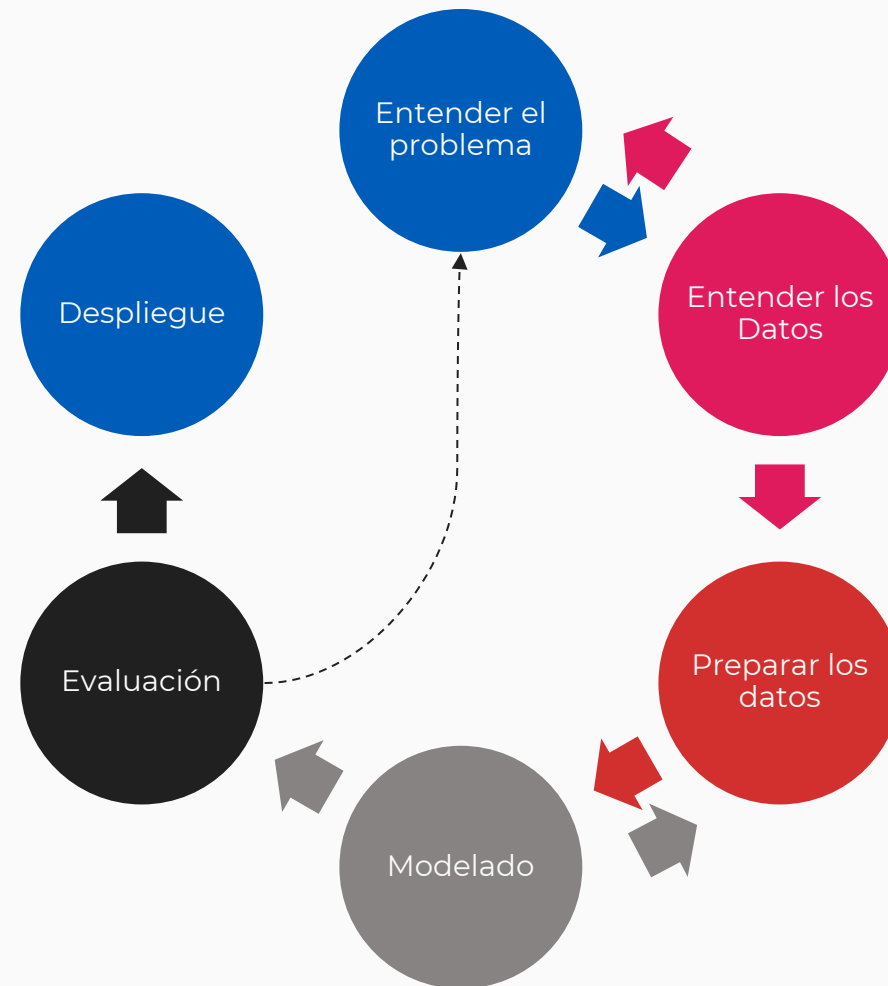
Head of Computational Social  
Science Research

The Alan Turing Institute



# CRISP - DM

Cross-industry Standard Process  
for Data Mining



# Objetivos de un Proyecto de Datos

Automatizar



Descubrir Patrones



Explicar



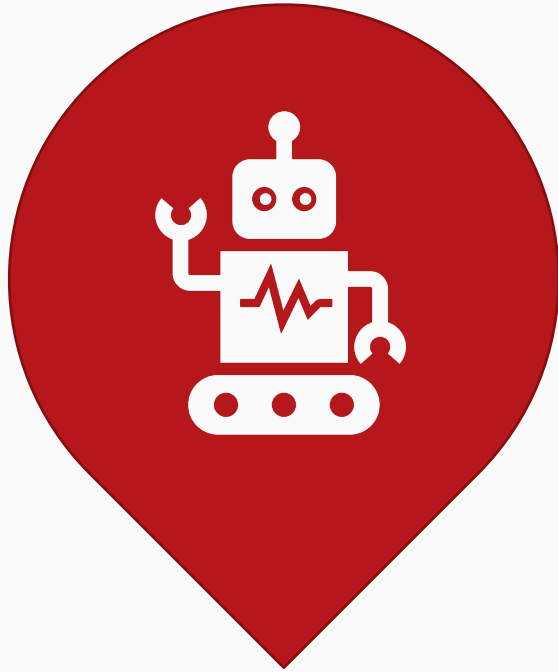
Optimizar



Predecir

# Aprendizaje de Máquina

Machine Learning



Métodos  
Computacionales



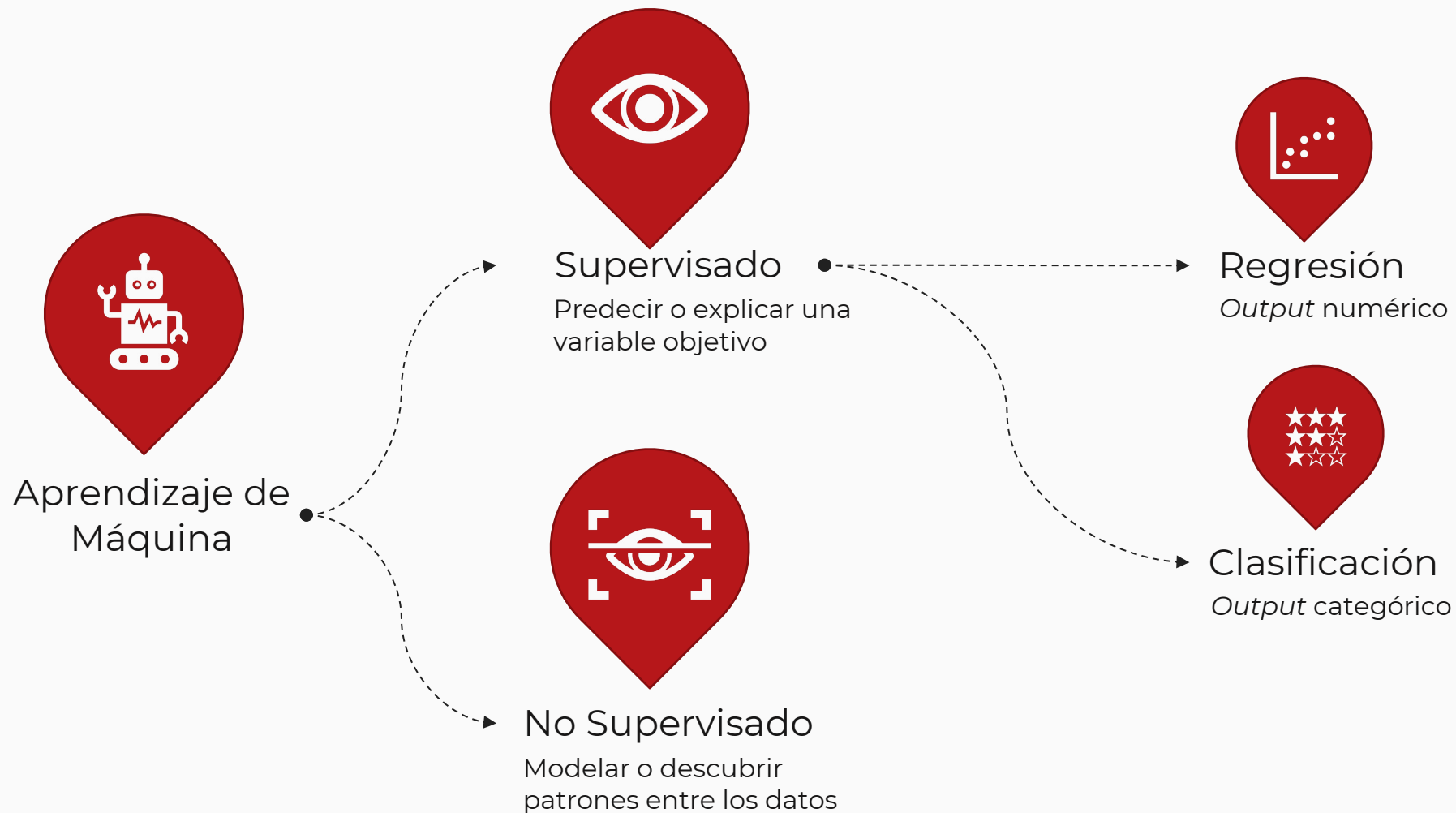
Datos



Herramientas para  
Toma de Tesisiones

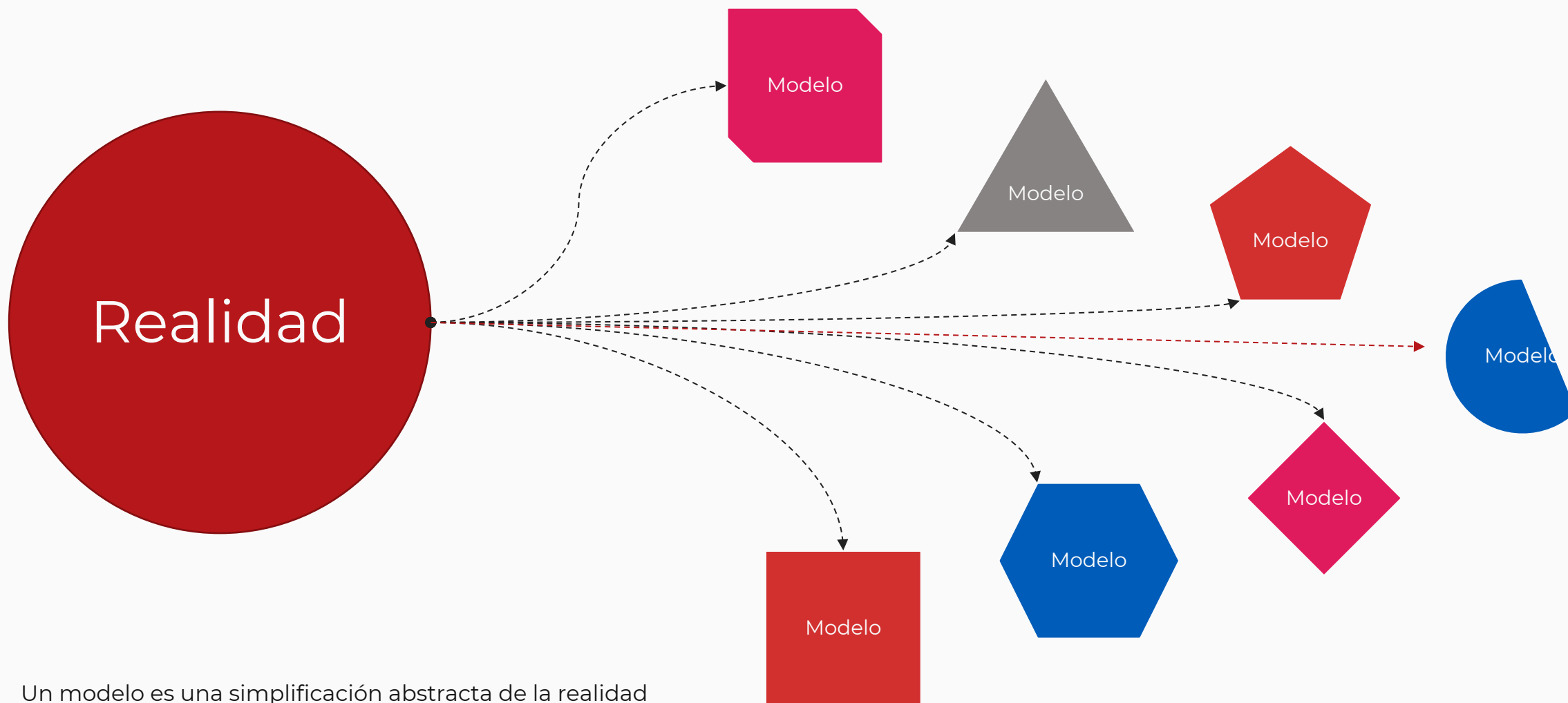
# Aprendizaje de Máquina

Machine Learning



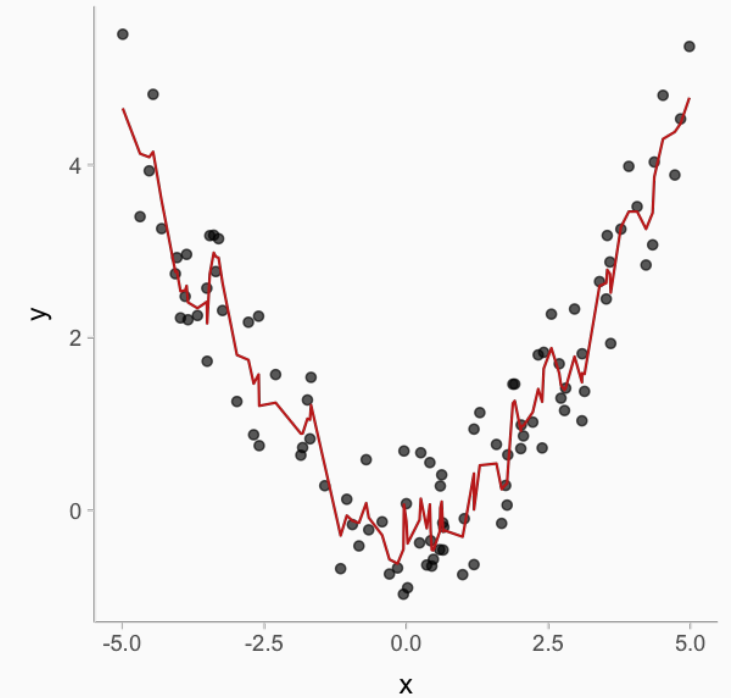
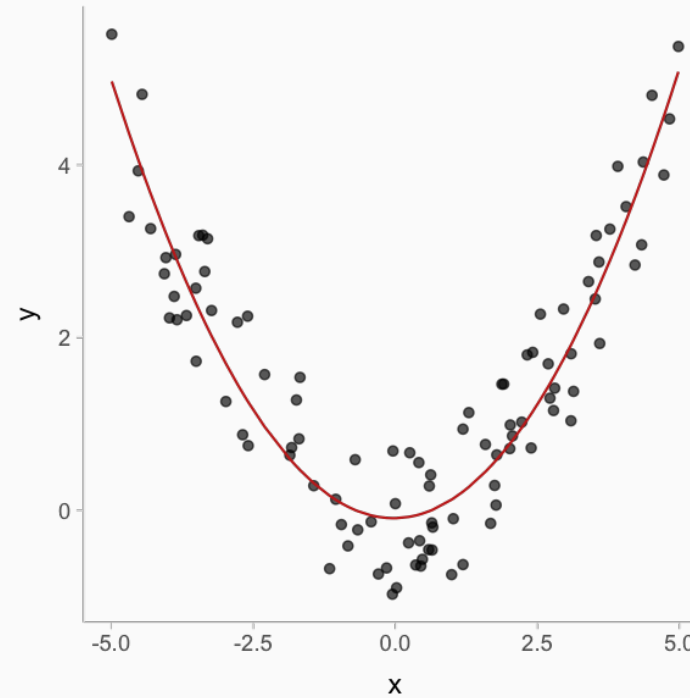
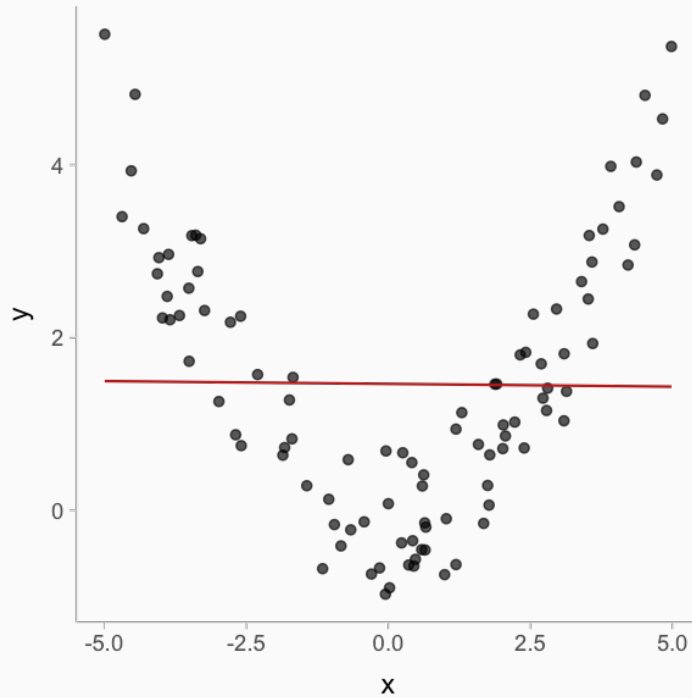


# ¿Qué es un modelo?



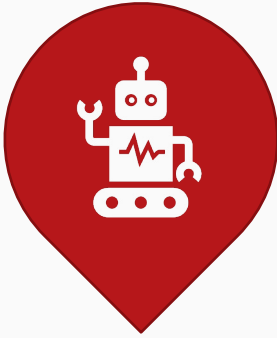
Un modelo es una simplificación abstracta de la realidad

# La elección de un modelo ideal



# Funciones de Pérdida

Loss functions



Aprendizaje

$$\min(f(y - \hat{y}))$$



Evaluación

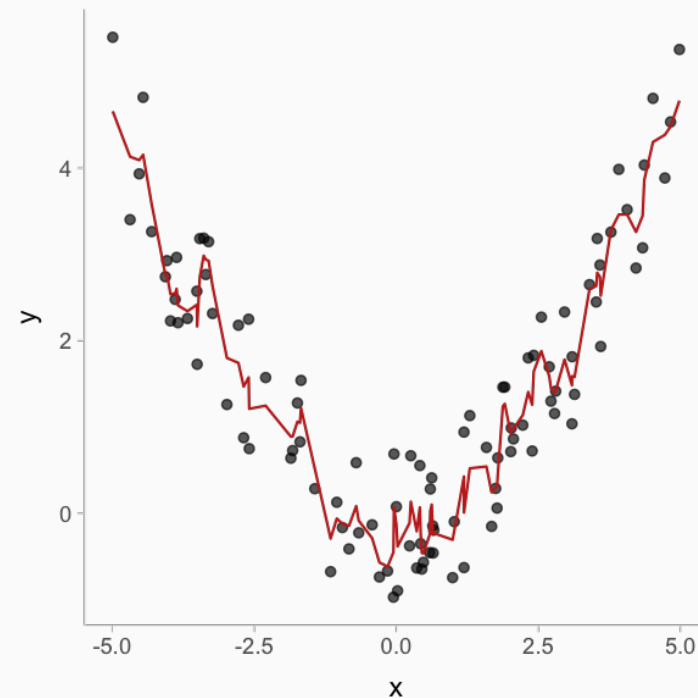
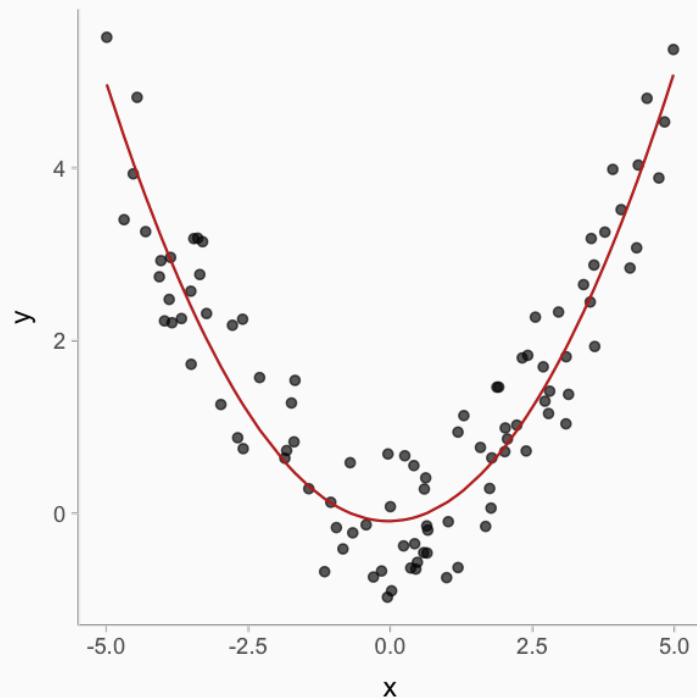
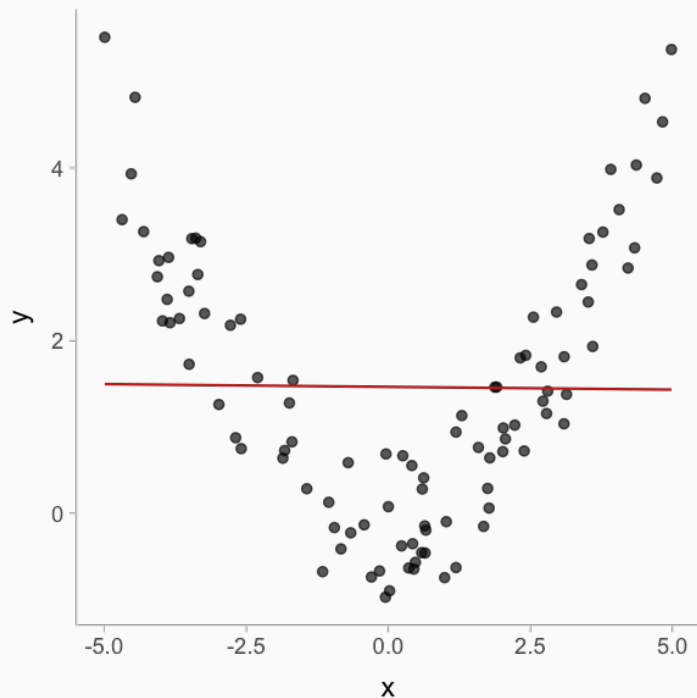
# Ejemplo: Error Cuadrático Medio

$$MSE(\hat{y}) = E[(\hat{y} - y)^2]$$

$$MSE(\hat{y}) = \underbrace{E[(\hat{y} - E(\hat{y}))^2]}_{\text{Varianza}} + \underbrace{E[(E(\hat{y}) - y)^2]}_{\text{Sesgo}}$$

# Varianza

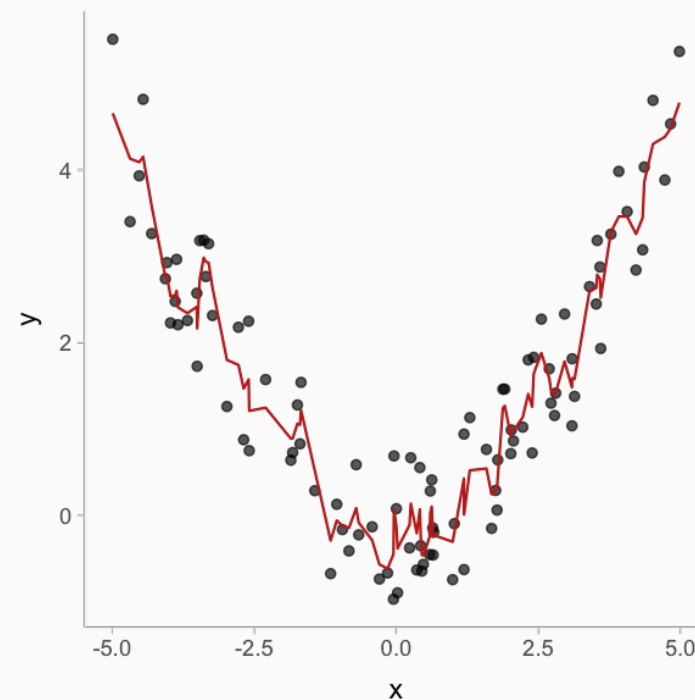
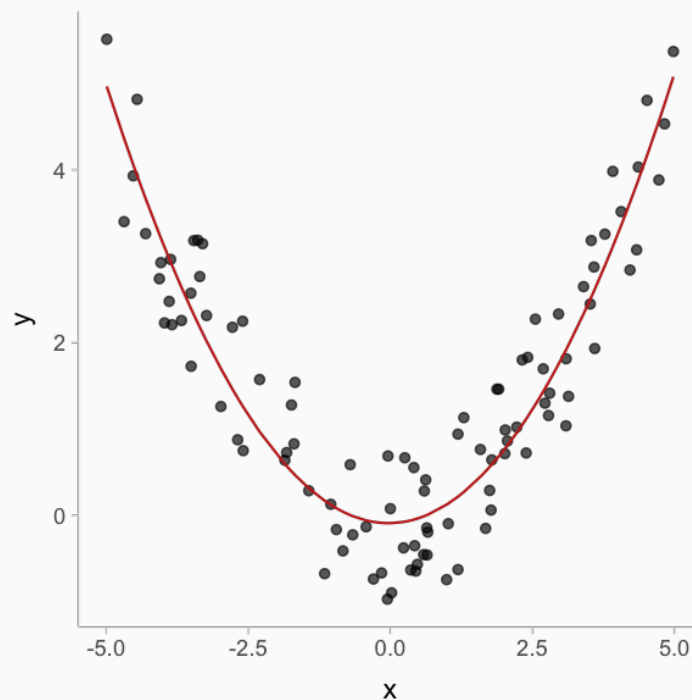
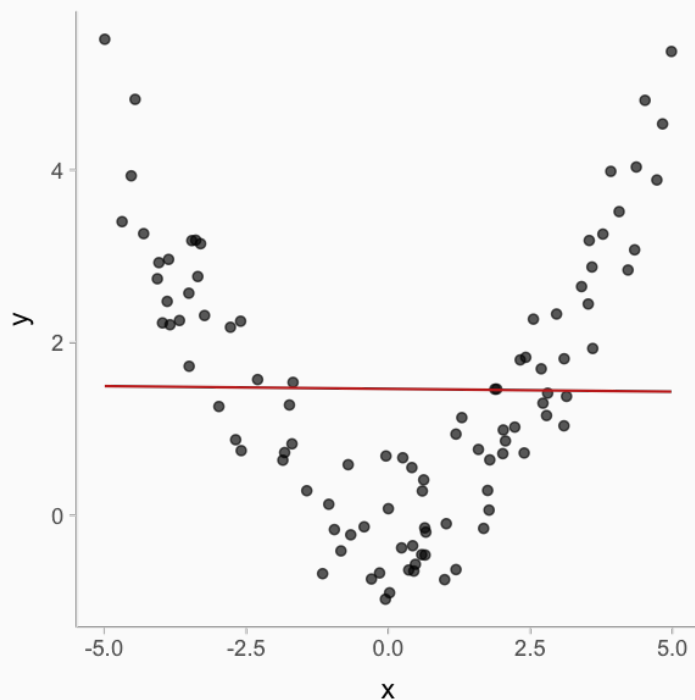
Qué tanto se acopla el modelo al recibir nueva información



Varianza Baja ←-----→ Varianza Alta

# Sesgo

Qué tan lejos se encuentra el valor estimado del valor real



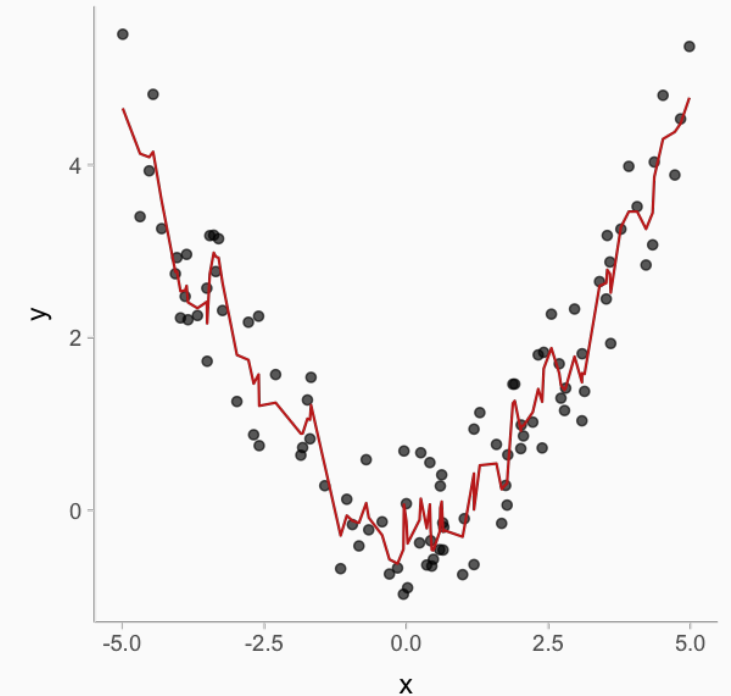
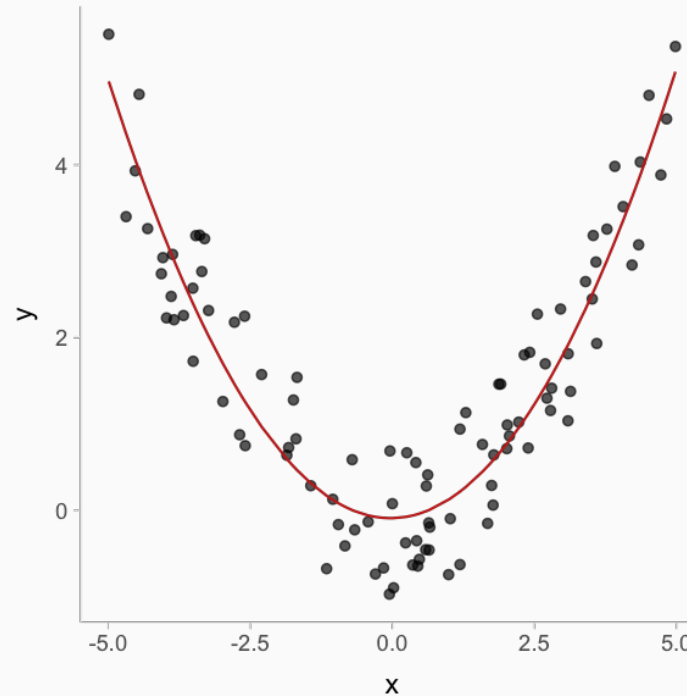
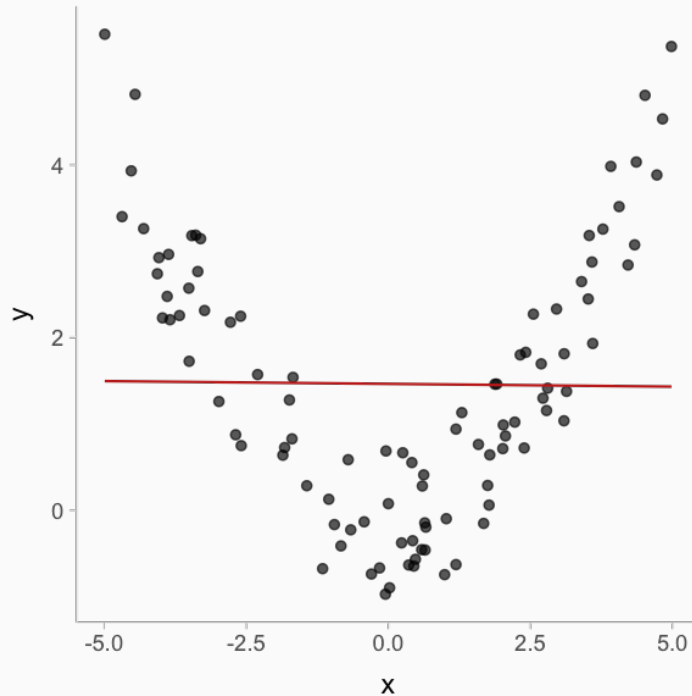
Sesgo Alto



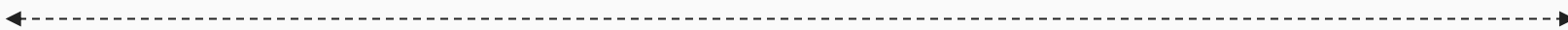
Sesgo Bajo

# La elección de un modelo ideal

Buscamos un modelo que aprenda de los datos con los que se le alimentaron  
y que sea capaz de generalizar las predicciones ante nuevos datos



Subajuste

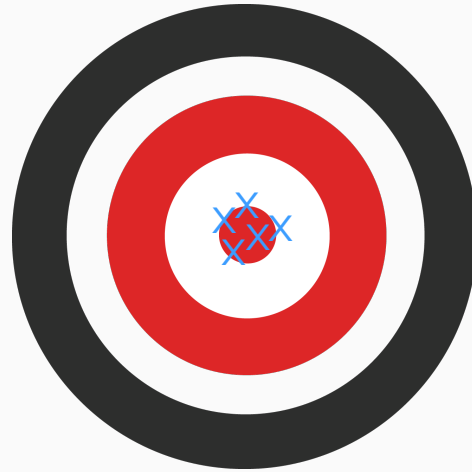


Sobreajuste

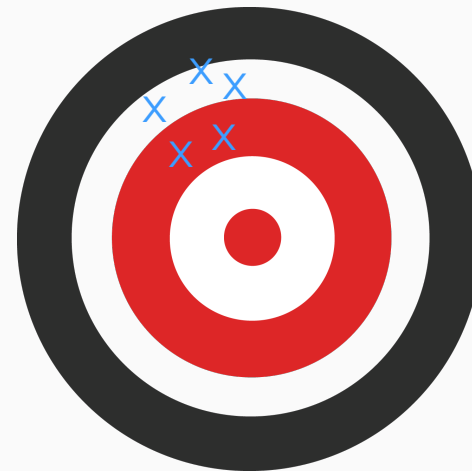
Poca Varianza

Mucha Varianza

Sesgo Bajo

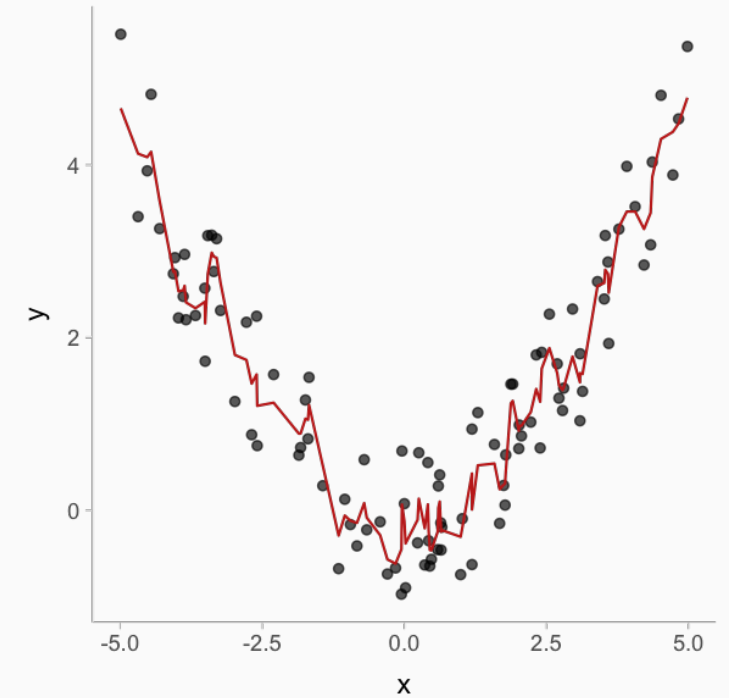
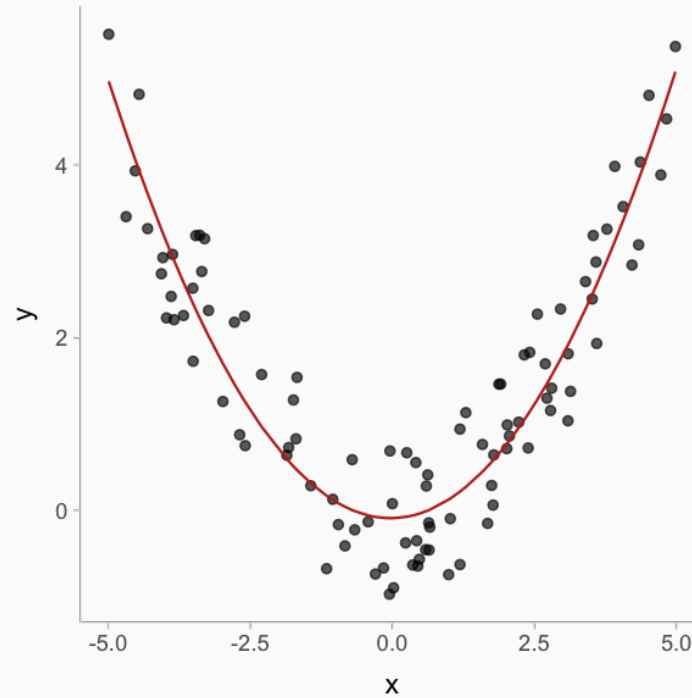
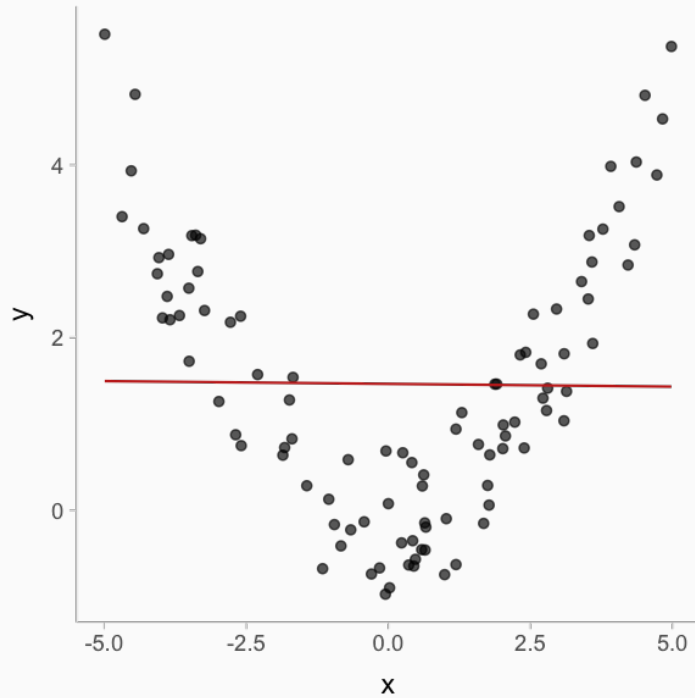


Sesgo Alto





# La elección de un modelo ideal



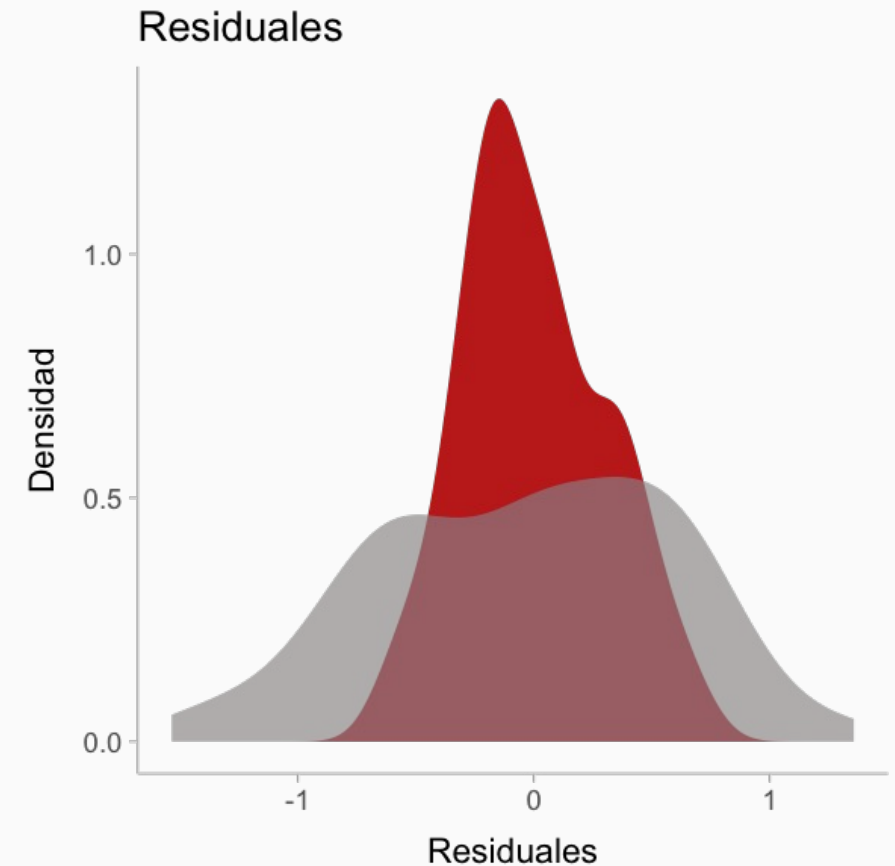
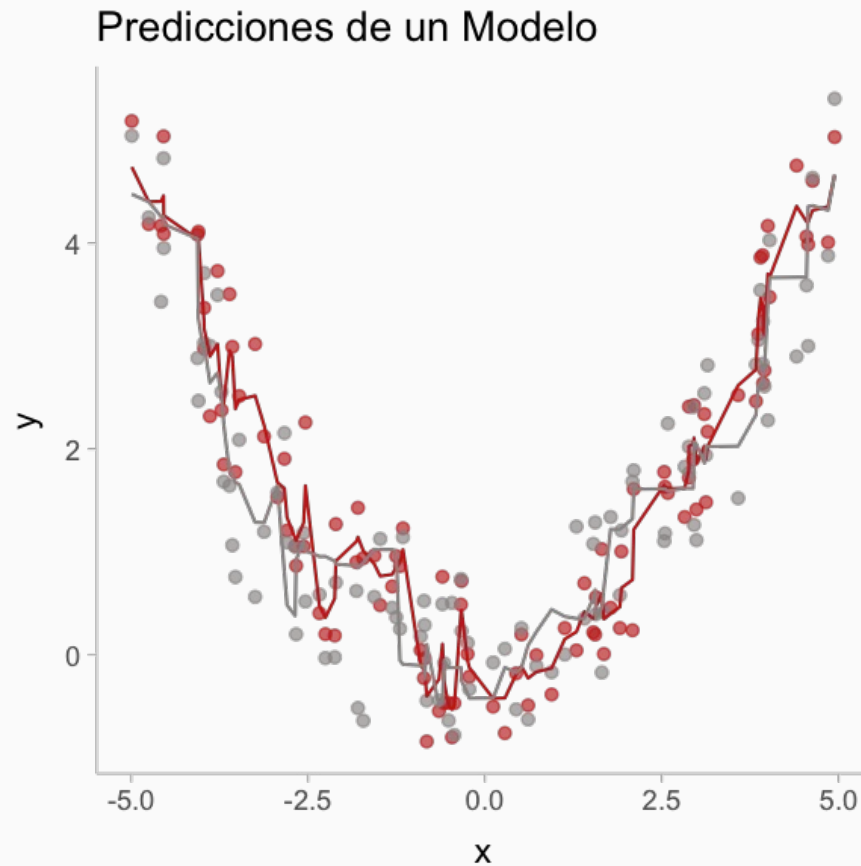
# Uno Modelo Sin Validar



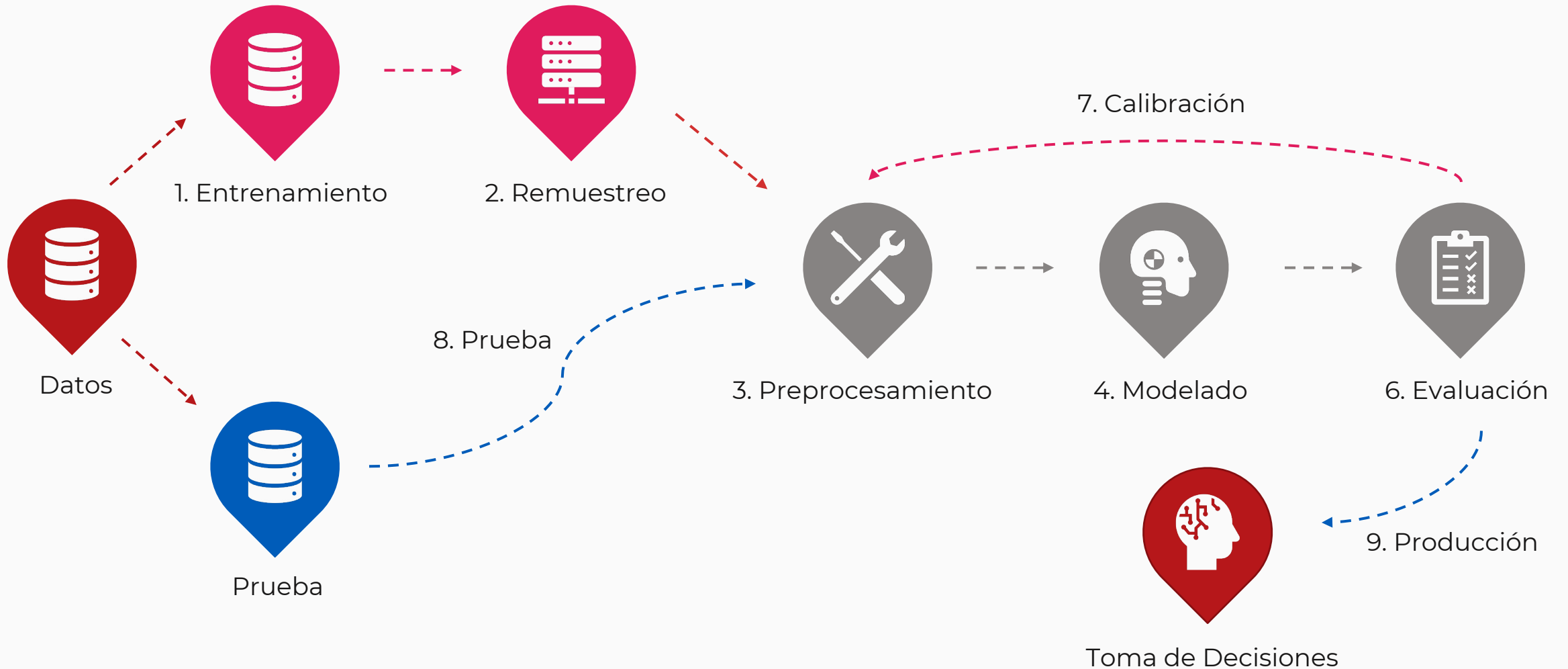
Datos Conocidos



Datos Nuevos

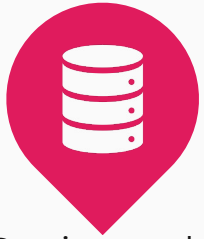


# Encontrar el mejor modelo



# Conjuntos de Entrenamiento y Prueba

Train & Test Splits



Conjunto de Entrenamiento

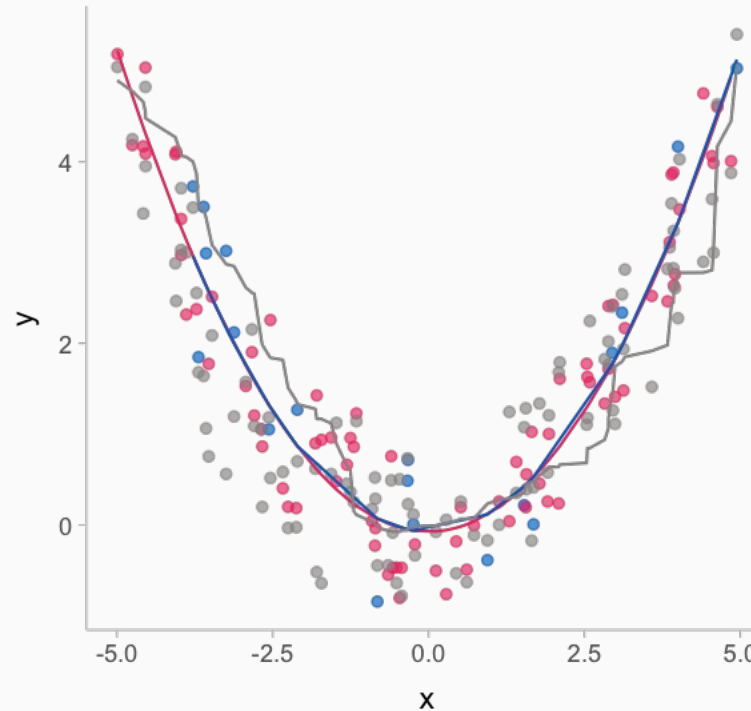


Conjunto de Prueba

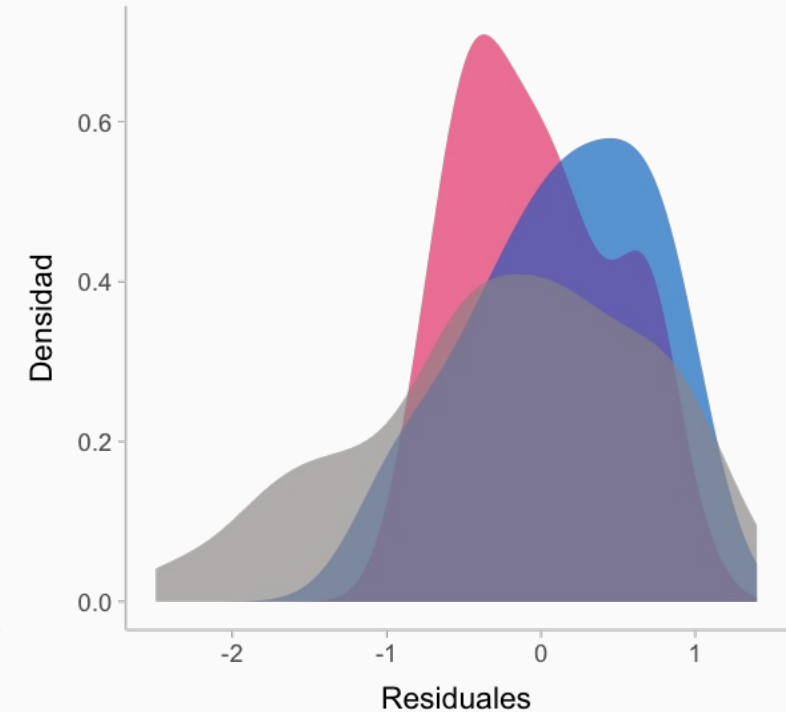


Datos Nuevos

Predicciones de un Modelo

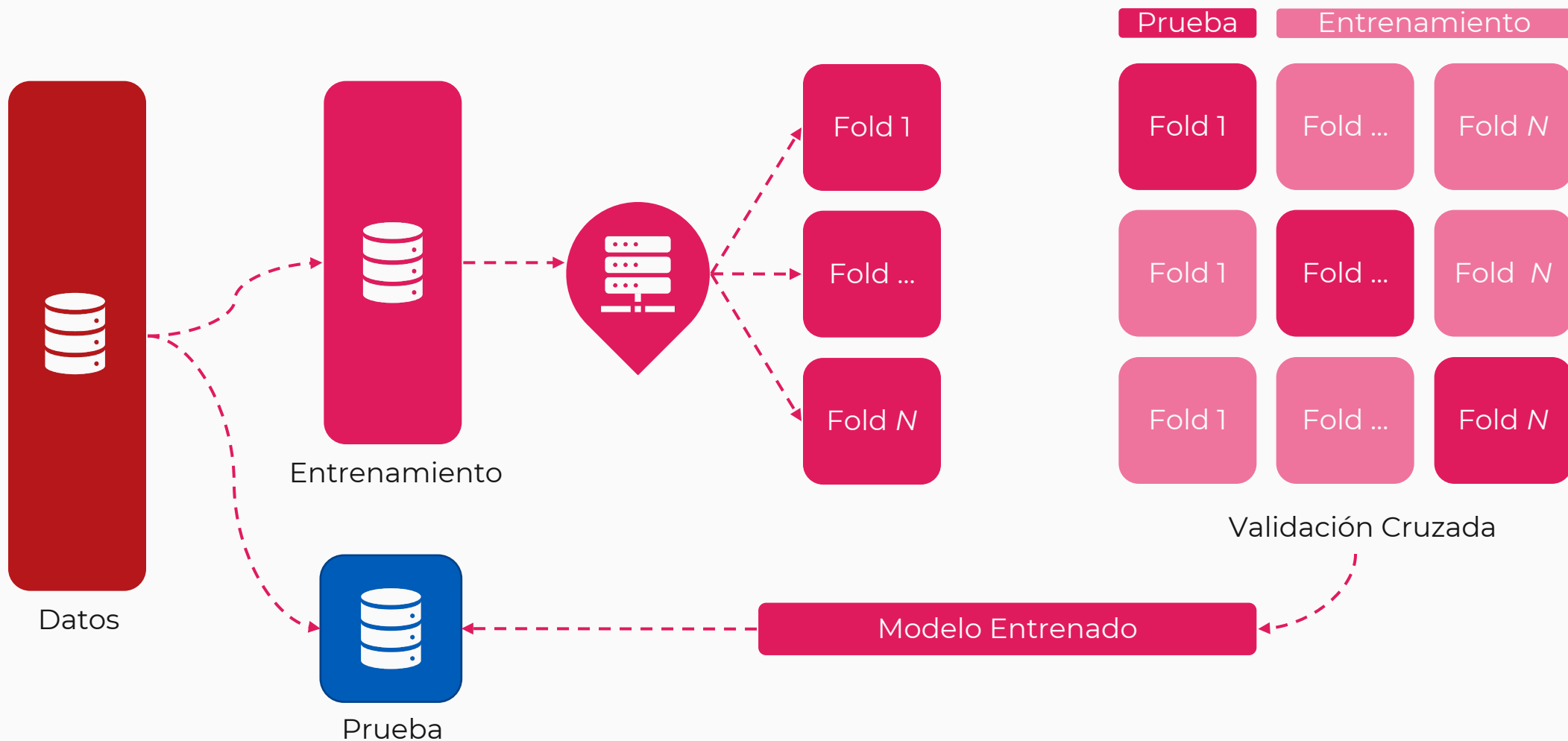


Residuales



# Validación Cruzada

Cross Validation



# Este no es un curso de inferencia causal



# Límites del Aprendizaje de Máquina

Limits of Machine Learning



No prevé soluciones a  
largo plazo



Soluciones específicas  
para casos específicos



Optimizaciones en un  
sistema dado