# Data science for Connected Health Devices

Merck - SDSC

# Why this presentation

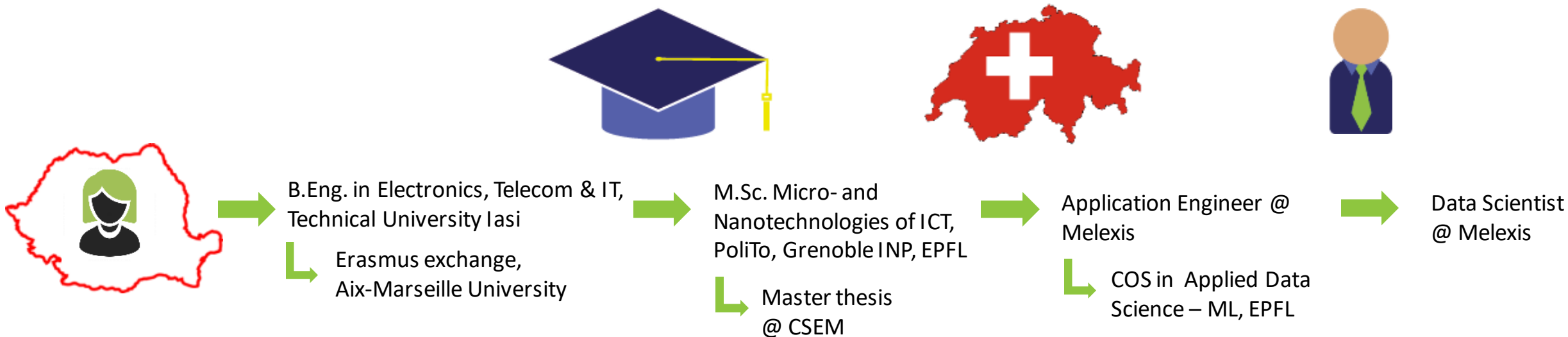*"Where do you see yourself in 5 years?"*

- To understand better

  - What would you do in the near future
  - What is data science in industry like
  - Why is Merck an attractive company for data science
  - What kind of problems you may need to solve

- To ask us questions

# Who am I?

B.Eng. in Electronics, Telecom & IT, Technical University Iasi

Erasmus exchange, Aix-Marseille University

M.Sc. Micro- and Nanotechnologies of ICT, PoliTo, Grenoble INP, EPFL

Master thesis @ CSEM

Application Engineer @ Melexis

COS in Applied Data Science – ML, EPFL

Data Scientist @ Melexis

Data Scientist @

MERCK
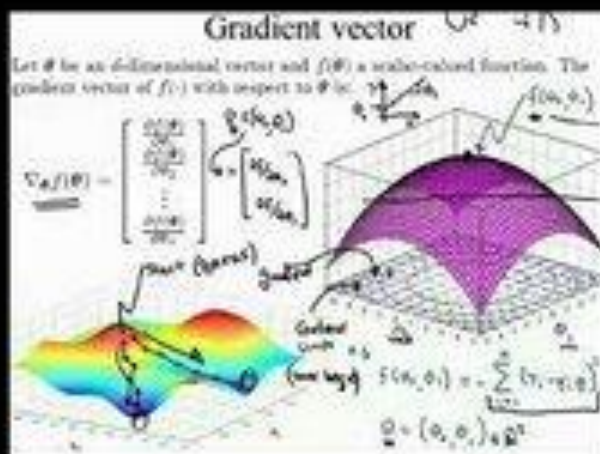
ETHzürich EPFL

SDSC

# Data Scientist



What my friends think I do

What my mom thinks I do

What society thinks I do

What my boss thinks I do

What I think I do

What I actually do

SELECT spending FROM db.users...

No, really...

**Connecting patients to doctors & nurses**

12 years of experience

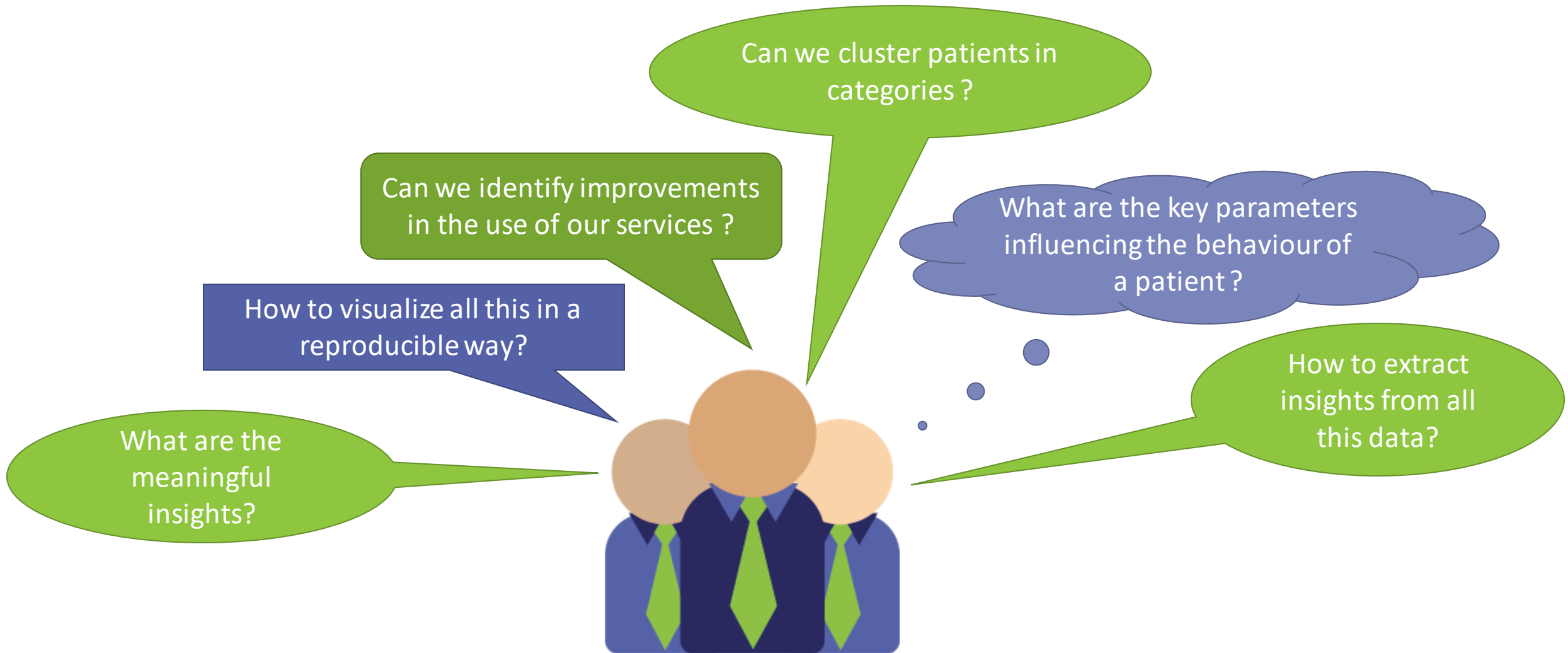500 000 devices produced

20 000 connected patients

50 countries

**SDSC**

# Data Science in the Merck – SDSC Collaboration

# The "Why"

# Context ➡ Questions

# The "How" – Data Process

**Data Gathering**

- via MySQL queries into python

- > 100 tables

- Find your way and merge relevant tables

**Cleaning**

- only patients with at least 1 calculated adherence point
- only patients who started after 2014 and logged no data in 2019
- remove null-entries

**Feature Engineering**

- Patient related data
- Adherence related data
- Device settings data, aggregations, majority voting
- Creating new (relevant) features

**Output dataset**

- Thousands of patients globally

- One line per patient
  - Features
  - Target

# Adherence
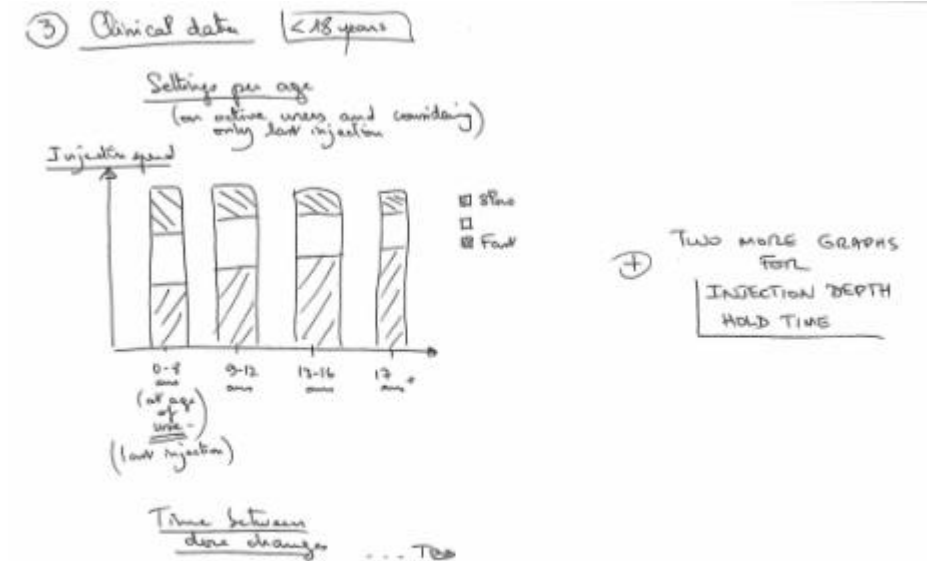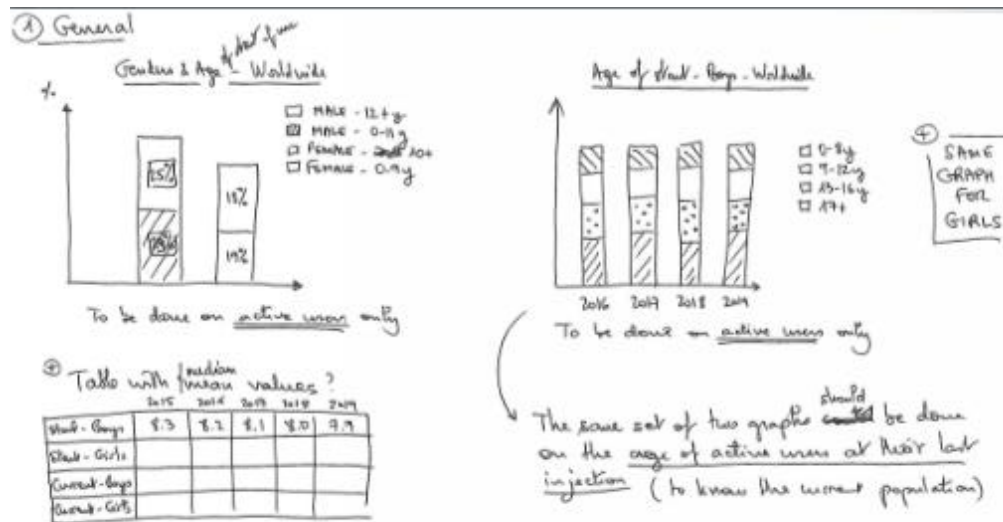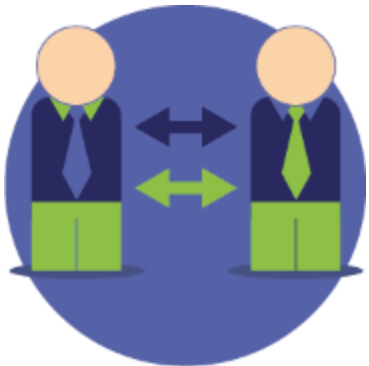
# Adherence

% of intake dose vs. prescribed dose

SDSC

# Duration of use

# Duration of use

for how long is a patient following
the treatment

SDSC

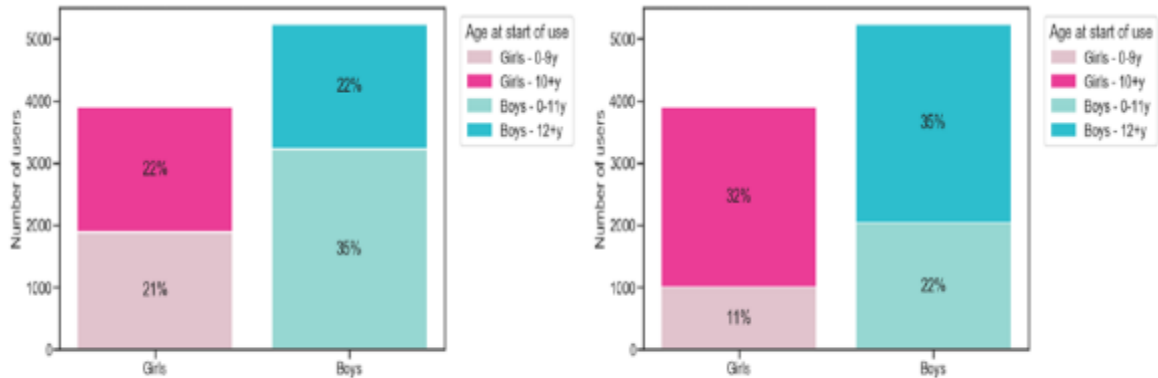# The "What"

# Regular Cross-communication

- Business to define their questions and/or demanded visualizations
- Iterative process



SDSC

# Reproducible reports

- Automatic generation of reports
  - Worldwide
  - Per country

```
generate_graphs.ipynb
generate_reports.ipynb
```
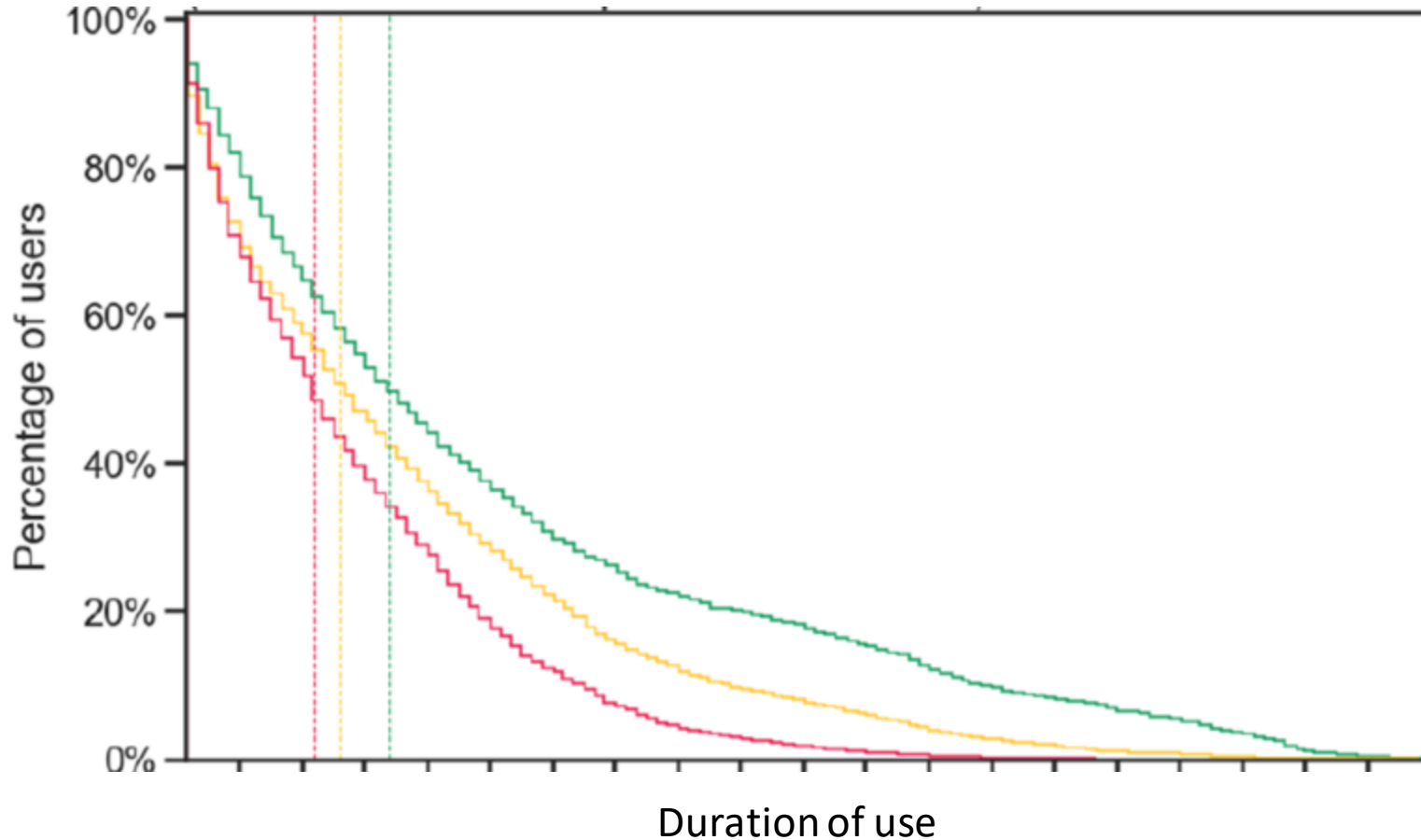


## 10.1. Comparison between clinics

The following table is ordered by the number of active users in 2019. When a quantity cannot be calculated 10 users), its corresponding cell is specified as "N/A".

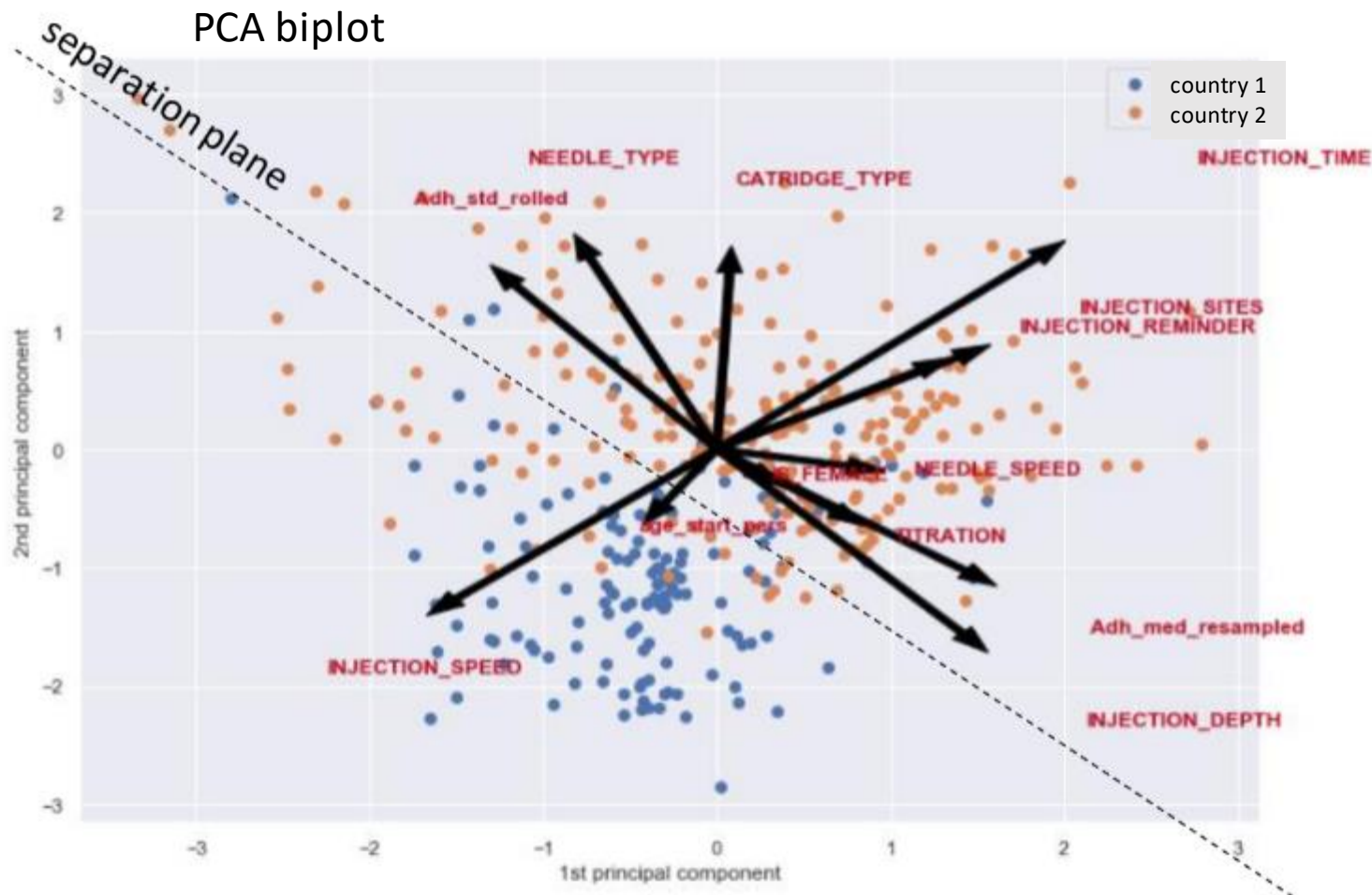| Rank | Clinic | # Active Users 2015 | # Active Users 2016 | # Active Users 2017 | # Active Users 2018 | # Active Users 2019 | Adherence[%] (median, last 12 months) |
|------|--------|---------------------|---------------------|---------------------|---------------------|---------------------|----------------------------------------|
| 1 | | 0 | 0 | 0 | 11 | 39 | 96.8 |
| 2 | | 22 | 34 | 55 | 53 | 35 | 85.7 |
| 3 | | 2 | 2 | 9 | 30 | 33 | 94.8 |
| 4 | | 0 | 5 | 21 | 40 | 31 | 92.9 |
| 5 | | 0 | 1 | 8 | 18 | 26 | 98.6 |
| 6 | | 4 | 23 | 35 | 38 | 25 | 96.1 |
| 7 | | 0 | 2 | 14 | 21 | 23 | 95.9 |
| 8 | | 3 | 21 | 22 | 19 | 20 | 90.5 |
| 9 | | 0 | 3 | 6 | 23 | 17 | 95.7 |
| 10 | | 0 | 0 | 0 | 9 | 12 | 89.0 |

- What are meaningful insights?



- Feature creation
- Processing
- Reverse cumulative distribution
- Churn analysis & prediction

SDSC

# Duration of use of a device, per country

# Interpretability

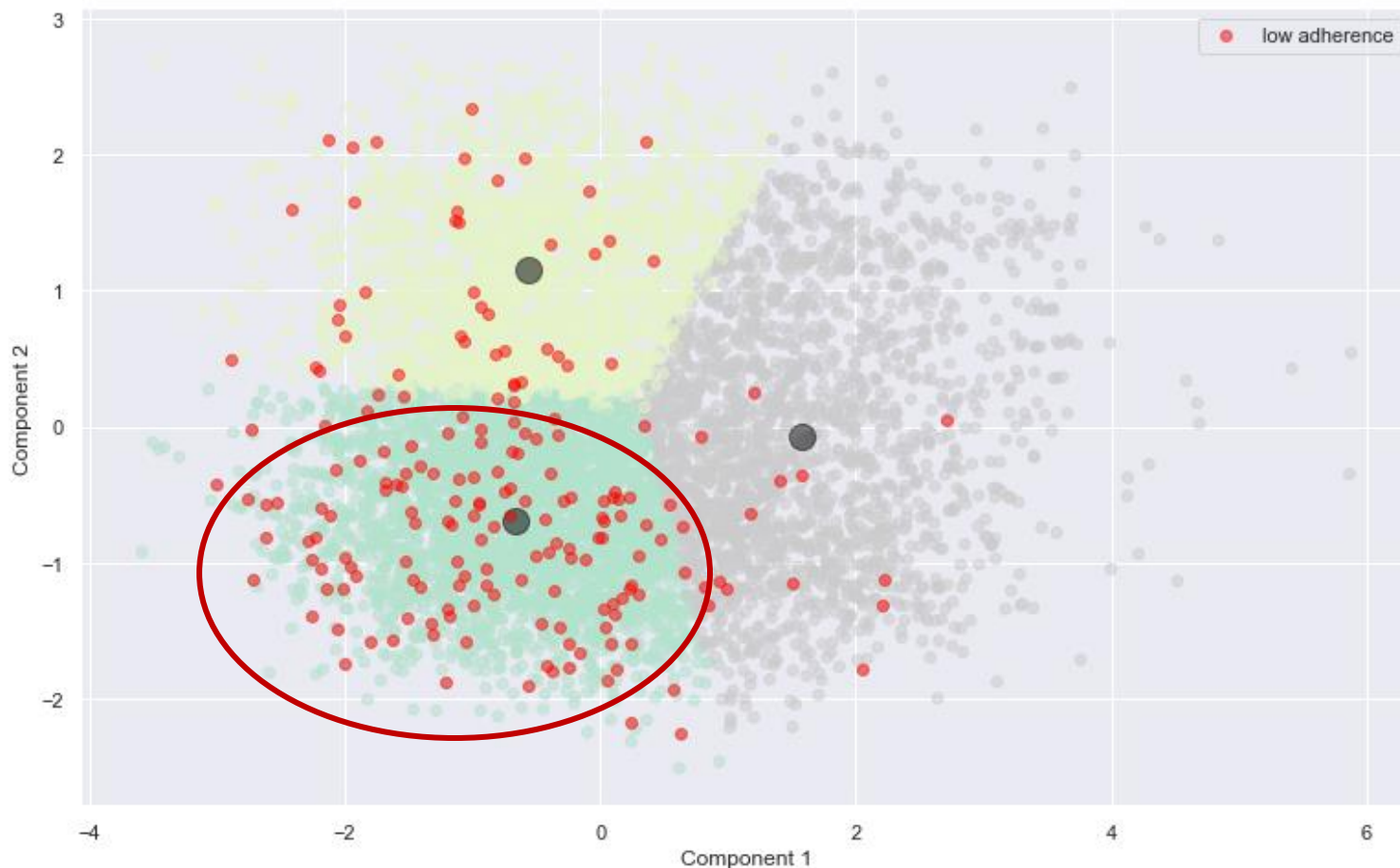- What are the key parameters influencing the outcome?


PCA biplot

Main **discriminators:**

- **Injection speed**
- **Injection time**
- **Injection sites**

# Clustering for patient segmentation

- Can we cluster patients in different categories?



70% of patients with low adherence are in this cluster

- K-Means clustering
- Dimensionality reduction

SDSC

# Conclusions

Pharma is a very rigorous industry, and **interpretability is key**

"Traditional" machine-learning algorithms are therefore chosen more frequently than Deep Learning

# The feeling of seeing your work reaching so many different people is just great!

# Thank you!

Swiss Data Science (@SDSCdatascience)

# Questions?