# Data Science in the Pharmaceutical Industry

Driving Innovation through Insights

Dr James Black,
Director, Insights Engineering
Pharma Development Data Sciences
German Data Science Days | 2023-03-09

# Table of Contents

Data science across the Roche Group

The diversity of data science across drug development lifecycles

Examples of data science in late stage
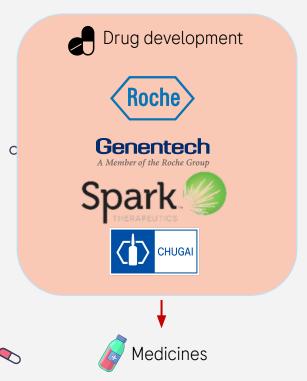
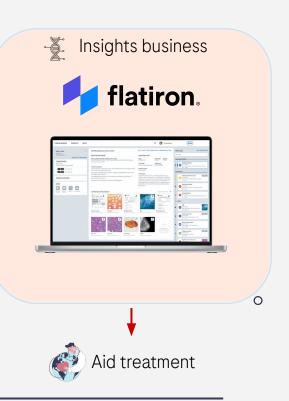A snapshot of data science in our team in 2023

# Roche in numbers

| | | | | |
|---|---|---|---|---|
| **1896** Founded in Basel and family still holds majority stake | | **>63 billion** Sales 2022 (CHF) | | **>14 billion** R&D investment (CHF) |
| | **#1** R&D investor in healthcare | **14 million** Patients treated | **29 billion** Diagnostic tests conducted | |
| **87** Drugs under development | | **103,613** Employees | | **27** Drug approvals in 2022 |

# Data science across Roche divisions



**Drug development**

Roche

Genentech — A Member of the Roche Group

Spark THERAPEUTICS

CHUGAI

↓ Medicines

**Diagnostics**

FOUNDATION MEDICINE®

↓ Diagnose & SaMD

**Insights business**

flatiron

↓ Aid treatment

# The diversity of data science across drug development lifecycles

## Discovery

*Can I predict toxicity for a molecule structure?*

*How does the tumour genetic fingerprint (NGS) change after exposure to X and Y?*

*What is the relationship between X and Y in published abstracts?*

## Early development

*What is the ideal dose?*

*Is there a multiplicative benefit of combining mechanism X and Y for condition Z?*

*Is there a companion diagnostic to identify patients that would benefit the most?*

## Late stage

*What is the efficacy of our drug against standard of care?*

*What is the safety profile?*

*What portfolio of evidence is needed to show how this drug performs across different patients?*

## Commercial

*What are the treatment patterns across different markets?*

*What is the value of our drug across markets?*

*What are the trends in standard of care, and what KoLs are key drivers of care patterns?*

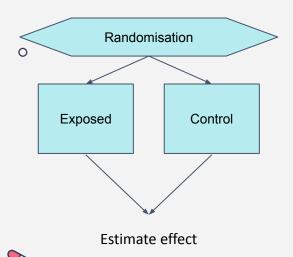| Research and early development | Product development | Product strategy |
|---|---|---|

# Ex.1: External controls

**Randomised clinical trials are the gold standard for causal inference, *but they are slow, expensive and tightly scoped to a specific hypothesis***

Can we leverage the routinely collected data from EHR and claims, across millions of patients, to improve our evidence portfolio?
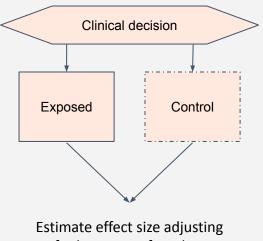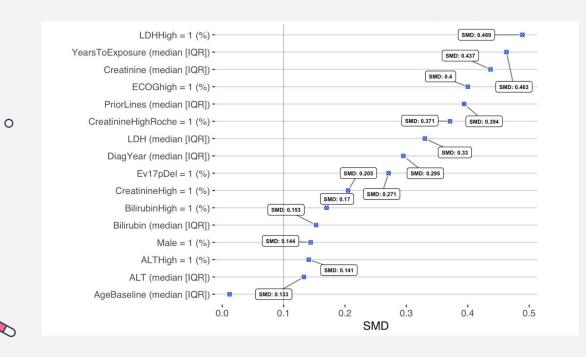
# Applying inclusion/exclusion criteria



**An example inc/exc attrition for a 2 LoT+ analysis**
Starting with 6,012 patients in our real world database with the same disease and >1 LoT, we drop to 187 that meet real world *approximations* of the inclusion and exclusion criteria
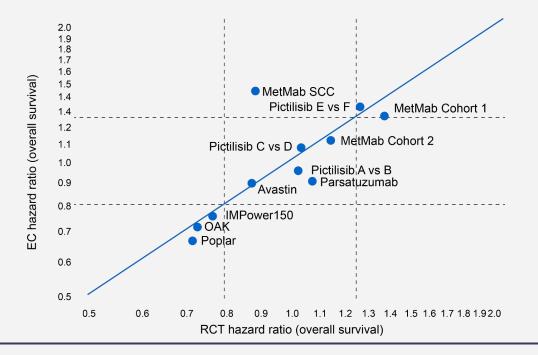
# Imbalance pre-analysis



**Standardised mean difference of baseline characteristics comparing real world control to single arm trial**

Analysing our populations as-is would be likely to lead to bias.

# External controls:
## an area of active research

# Ex.2: The RAAD Challenges



data-centric AI | 2022

# Roche Advanced Analytics Data (RAAD) Challenges

## RAAD Challenge 1.0

Predict the probability that a patient will be alive at 1 year after 1L treatment initiation, using all the patient data available up to the start of 1L treatment.

**Data:**
9,500 patients' across train and test from across seven different cancer types were used to train the models.

**Participation:**
500+ Roche employees, 132 teams, 28 sites

## RAAD Challenge 2.0

Predict response in a new drug and indication combination using training data containing either that drug in a different indication, or that indication with other drugs

**Data:**
6,000 patients across train and test from historical trials

**Participation:**
500+ Roche employees, 141 teams, 38 sites

## *RAAD Challenge 3.0 judging was last Friday!*

# Data-Centricity: clinico-genomic feature library

## Hackathon tasks

With a fixed prediction model, enrich the training data and create 10 features that improve prediction on the test data.

## Benefit to Roche & our Patients

This challenge focuses our attention not on the model, but on the preparation of data for modeling, and will help to develop a rich catalog of features enabling future researchers to develop more impactful insights and products (e.g .to support clinical decision support at point of care, enhanced patient segmentation for clinical trials etc).

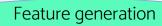# Data-Centricity: clinico-genomic feature library

NGS test

Vitals

Labs

Demographics

Feature generation

🔒 Model

Prediction target is time till death

96,686 patients

118 variables

58,066 events

744 MBs of data

93,673,017 data points

Train - 60%

Leaderboard - 20%

Test - 20%

# Data science in our team in 2023
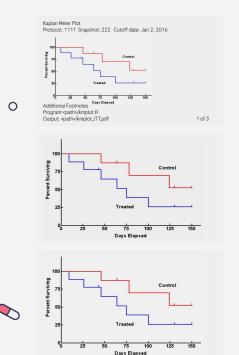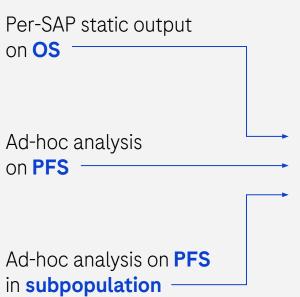
# Being multilingual is the new-norm



We focus on R-backbone for analytics, and python for tooling, but there is a continued growth in other languages and technologies that we leverage in the pursuit of constantly doing things 'better'
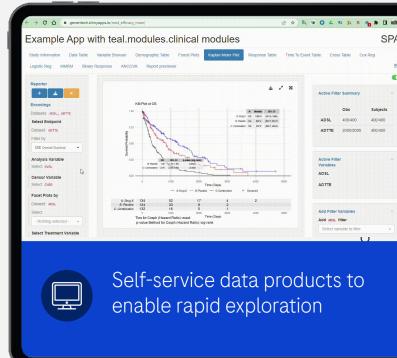
5 years ago

Today

# Data products, in addition to static insights, is now the new normal



Per-SAP static output
on **OS**

Ad-hoc analysis
on **PFS**

Ad-hoc analysis on **PFS**
in **subpopulation**

Self-service data products to
enable rapid exploration

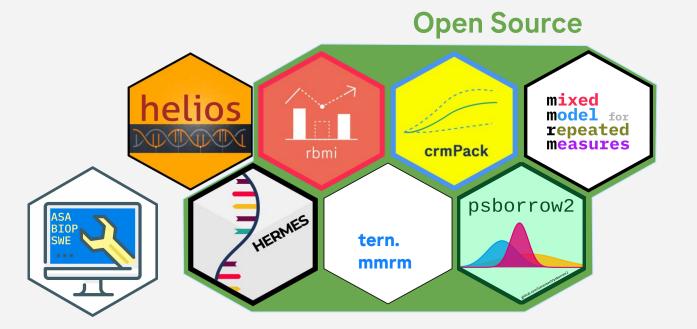# We have more of a focus on re-use, and open source



Example R packages our team work on.

In Feb 2023, our team of **~50 data scientists** contributed at least one commit to **100 open source repositories**, and **117 closed source repositories**.

# Our first coursera course is now live!



Bridging general DS into clinical reporting DS

Statistical Programming, Statistical engineering & RWE courses coming

go.roche.com/course

**Contact: james.black.jb2@roche.com**

I'm currently hiring!
Research Software Engineer
positions open across Basel,
San Francisco, Welwyn, Mississauga
**go.roche.com/dsjobs**

Data science careers:
code4life.roche.com