

Data warehousing and analytics for healthcare

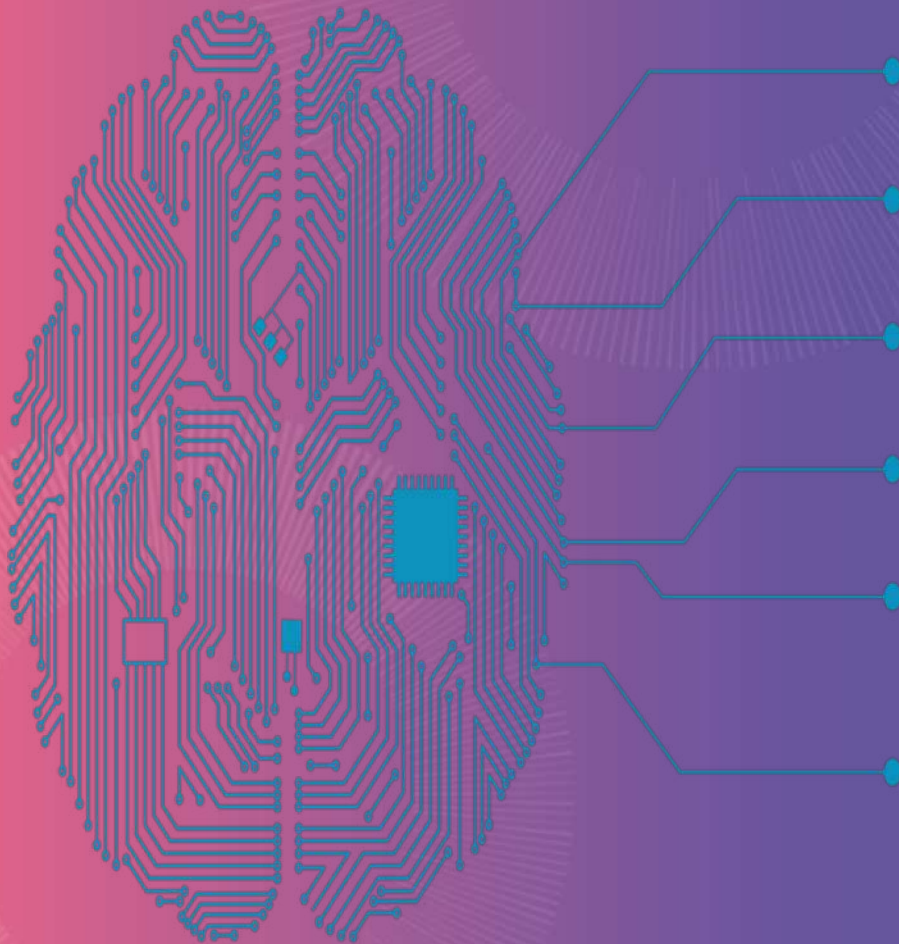
Dr. Daniel Kapitan | Chief Data Scientist | Mediquest

TU/e Data Science Center

Eindhoven, 29 May 2019

mediquest

The promise of AI: will computers be able to do all this?



natural language processing

knowledge representation

automated reasoning

machine learning

computer vision

robotics

Real life is less sexy, and more hard work



Sanders (2012), *Healthcare Analytics Adoption Model*

Today's agenda

| | |
|---------|---|
| Level 8 | Personalized Medicine & Prescriptive Analytics |
| Level 7 | Clinical Risk Intervention & Predictive Analytics |
| Level 6 | Population Health Management & Suggestive Analytics |
| Level 5 | Waste & Care Variability Reduction |
| Level 4 | Automated External Reporting |
| Level 3 | Automated Internal Reporting |
| Level 2 | Standardized Vocabulary & Patient Registries |
| Level 1 | Enterprise Data Warehouse |
| Level 0 | Fragmented Point Solutions |

5 – 8: Healthcare analytics

Case study: predicting outcomes of surgery

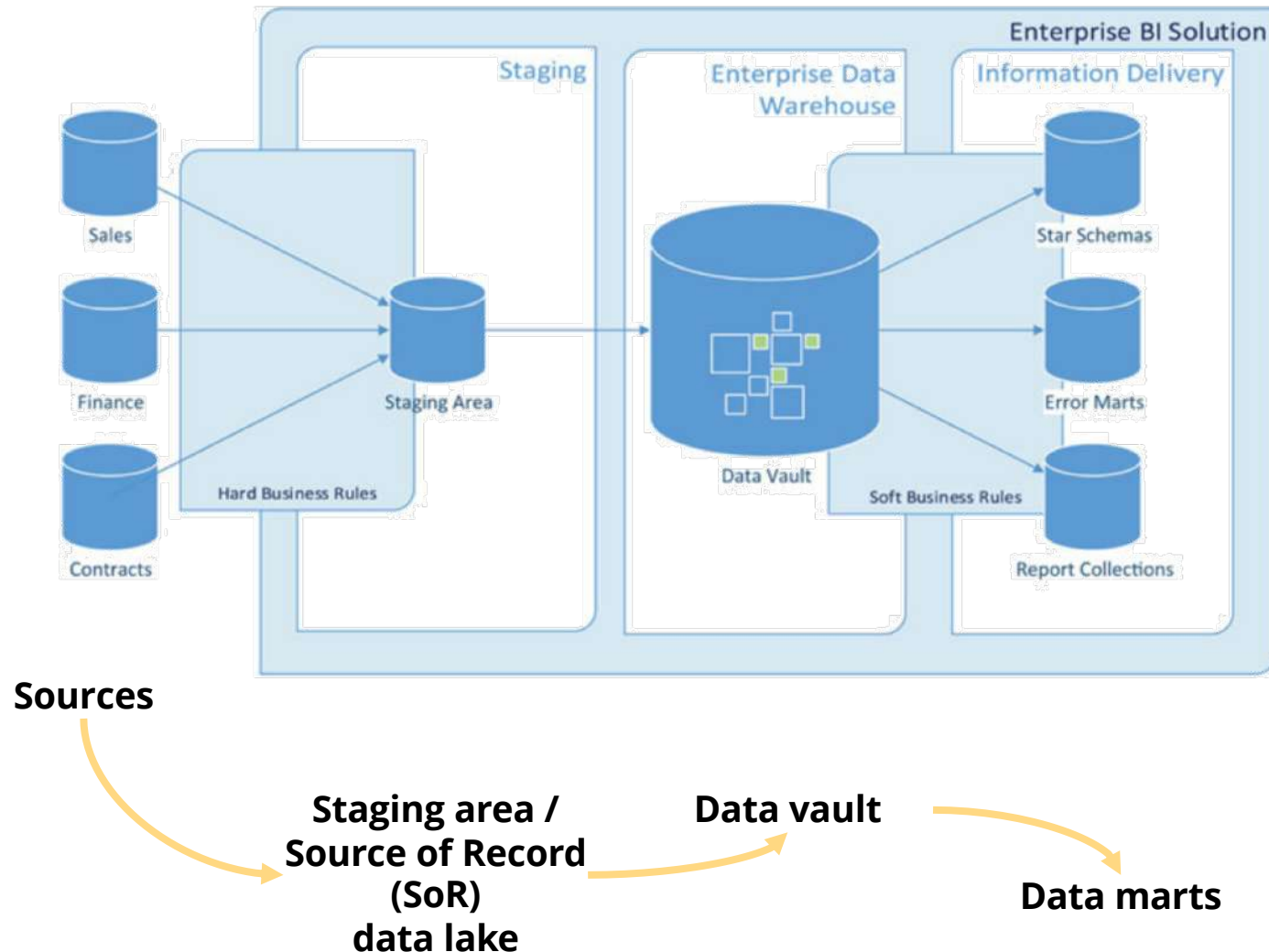
1 – 4: Laying the foundation

- Data warehousing
- Data integration
- Semantic modelling
- Business intelligence

Part I: Laying the Foundation

DATA WAREHOUSING IN HEALTHCARE

The main components of a data warehouse



The 'Supplier Landscape': choices, choices ...



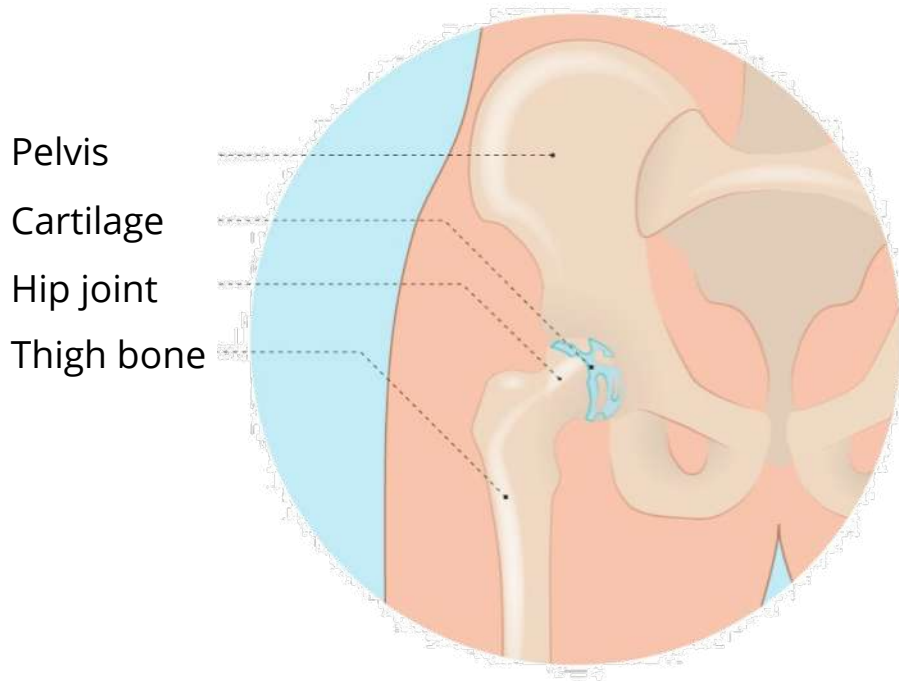
- **How many layers:**
decoupling vs. simplicity?
- **Storage platform:**
RDBMS, NoSQL, cloud ...?
- **Data integration:**
Semantic modelling,
dimensional modeling ... ?
- **Way of working:**
hand-coded vs. graphical tools?
- **Dashboarding & data viz:**
Yellowfin, Tableau, PowerBI ...

Choosing the storage engine



| | | | |
|--|---|---|--|
| Daily use (from a data scientist's perspective) | <ul style="list-style-type: none"> • good CSV handling • Unicode support • Regular expressions • ANSI SQL compliance | <ul style="list-style-type: none"> • native Excel • Nice GUI IDE for querying, maintenance and workflow | <ul style="list-style-type: none"> • Cloud-native, low maintenance • Easy to use with API |
| Platform | <ul style="list-style-type: none"> • All major OS-es | <ul style="list-style-type: none"> • Windows, Linux added recently | <ul style="list-style-type: none"> • Cloud (lock-in) |
| Extensibility | <ul style="list-style-type: none"> • Many languages built-in (Python, Javascript, R) | <ul style="list-style-type: none"> • R built-in (acquisition RStudio) • .NET framework | <ul style="list-style-type: none"> • Javascript • Tight integration GCP |
| Specific for data analytics | <ul style="list-style-type: none"> • Best-in-class for GIS • jsonb format for unstructured data storage • MADlib for built-in machine learning | <ul style="list-style-type: none"> • Integrates well with Power BI stack • Hybrid solution cloud – on-premise possible with Azure SQL | <ul style="list-style-type: none"> • StandardSQL performs well even on petabytes • BigQuery ML |
| TCO | <ul style="list-style-type: none"> • Forever free with community version • Enterprise DB for paid support (same pricing as MS SQL) | <ul style="list-style-type: none"> • Value for money with Standard Edition • Gets expensive when Enterprise features are needed | <ul style="list-style-type: none"> • Pay only for queried volume |

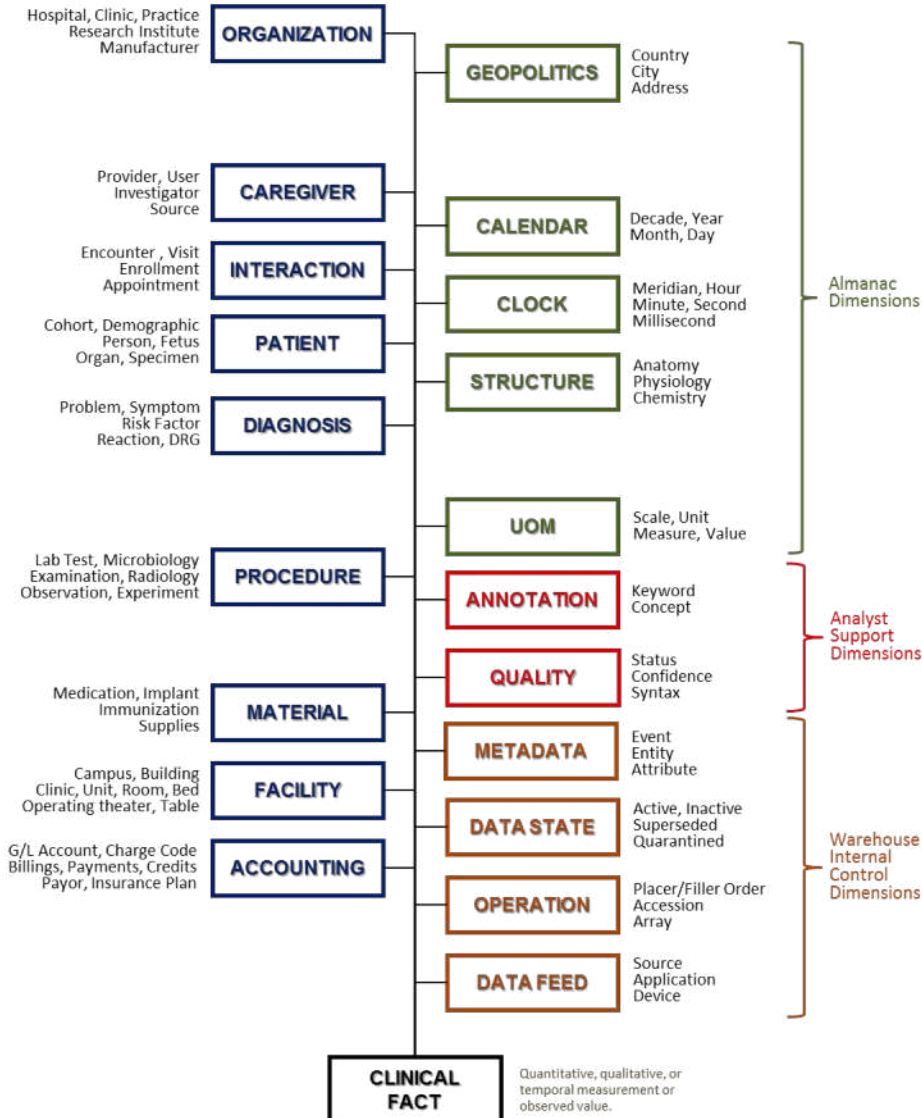
Data integration with semantic dimensional modelling



- Open a webbrowser and go to: https://zibs.nl/wiki/HCIM_Mainpage
- **Assignment: design a data model for a orthopedic clinics**
- Each data element should be captured/mapped to a 'zorginformatiebouwsteen'

Think about all the data you need to integrate so you can do funky machine learning on it

Data integration with semantic dimensional modelling



- Concept of star schema:
 - Typically 5 to 10 fact-tables
 - 20 to 30 dimension tables
- Modelling challenges:
 - Uniformity of dimensions (using same codes)
 - Uniformity of business keys (how to uniquely identify a hip implant)
 - Privacy-sensitive data (SSN, identifiable data)
- Engineering challenges:
 - Dealing with changes in data
 - Speed of batch processing

Biehl (2015), *Data Warehousing for Biomedical Informatics*

Choose your way of working

Hard/hand-coded scripts

```
def init_sor_persoon_hstage_to_patient_mappings(sor):
    mappings = []

    mapping = SorToEntityMapping('persoon_hstage', Patient, sor)
    mapping.map_bk(['timeff', 'ifct_relatiennr'])

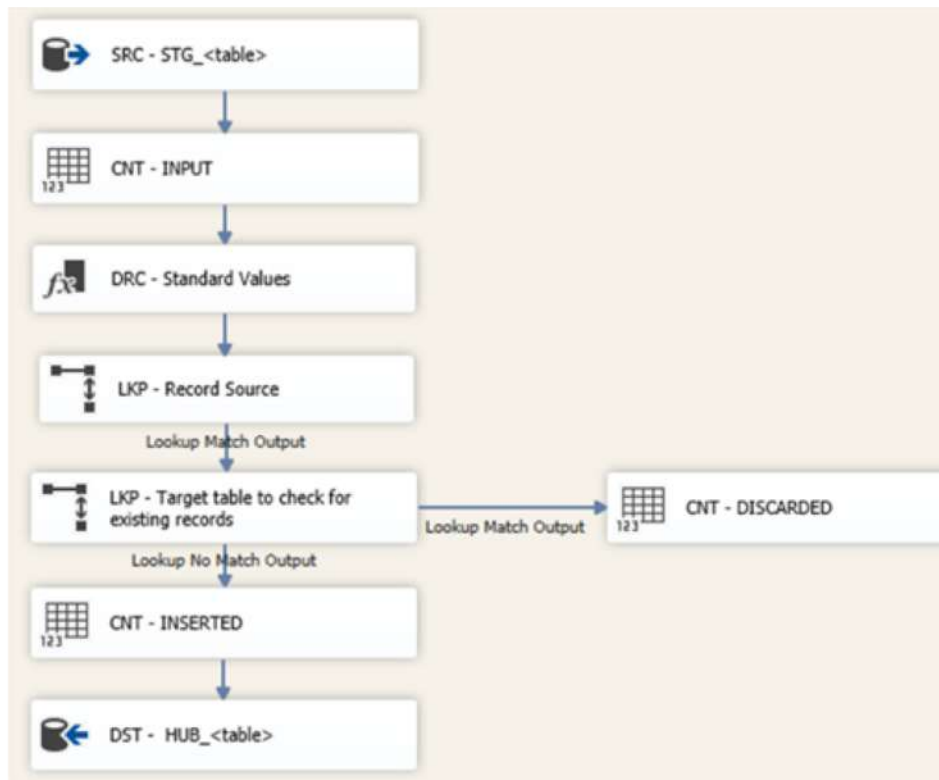
    # SAT Patient.Identificatie
    mapping.map_field("extern_patientnummer", Patient.Identificatie.extern_nummer)
    mapping.map_field("ifct_relatiennr", Patient.Identificatie.nummer)
    mapping.map_field("ifct_id", Patient.Identificatie.bron_id)

    # SAT Patient.Default
    mapping.map_field("ifct_geboortedtm::date", Patient.Default.geboortedatum)
    mapping.map_field("ifct_geslacht", Patient.Default.geslacht_code)
    mapping.map_field("", Patient.Default.meerling_indicator)
    mapping.map_field("", Patient.Default.overlijdens_indicator)
    mapping.map_field("", Patient.Default.datum_overlijden)

    # SAT Patient.IdentificatieBewijs
    mapping.map_field("ifct_bsn", Patient.IdentificatieBewijs.nummer, type=Patient.IdentificatieBewijs.Ty
    mapping.map_field("", Patient.IdentificatieBewijs.geldig_tot, type=Patient.IdentificatieBewijs.Ty
    mapping.map_field("ifct_legitimatied", Patient.IdentificatieBewijs.nummer, type=Patient.Identifi
    mapping.map_field("", Patient.IdentificatieBewijs.geldig_tot, type=Patient.IdentificatieBewijs.Ty
    mapping.map_field("ifct_rijbewijsnummer", Patient.IdentificatieBewijs.nummer, type=Patient.Identi
    mapping.map_field("", Patient.IdentificatieBewijs.geldig_tot, type=Patient.IdentificatieBewijs.Ty

    # SAT Patient.Adres
    mapping.map_field("ifct_straat_b", Patient.Adres.straat, type=Patient.Adres.Types.woonadres)
    mapping.map_field('sor_timeff.split_huisnummer(ifct_huisnr_b)).huisnummer', Patient.Adres.huisnu
    mapping.map_field("", Patient.Adres.huisnummerletter, type=Patient.Adres.Types.woonadres)
    mapping.map_field('sor_timeff.split_huisnummer(ifct_huisnr_b)).huisnummer_toevoeging', Patient.A
    mapping.map_field("", Patient.Adres.aanduiding_bij_nummer_code, type=Patient.Adres.Types.woonadre
    mapping.map_field("ifct_plaats_b", Patient.Adres.woonplaats, type=Patient.Adres.Types.woonadres)
    mapping.map_field("", Patient.Adres.gemeente, type=Patient.Adres.Types.woonadres)
    mapping.map_field("", Patient.Adres.land_code, type=Patient.Adres.Types.woonadres)
    mapping.map_field("ifct_postcode_b", Patient.Adres.postcode, type=Patient.Adres.Types.woonadres)
    mapping.map_field("", Patient.Adres.additionele_informatie, type=Patient.Adres.Types.woonadres)
```

Graphical user interface

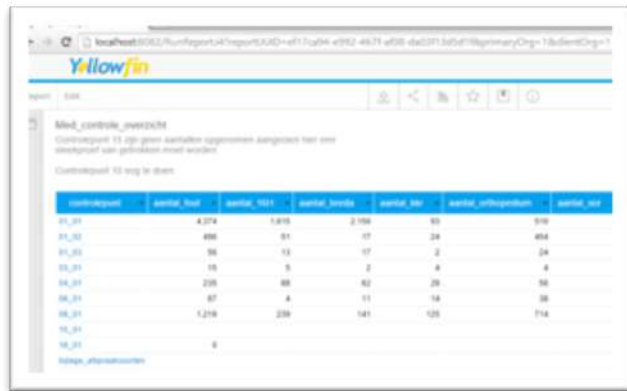


Design data marts for flexibility and human-readability



Let end-users choose their own tool

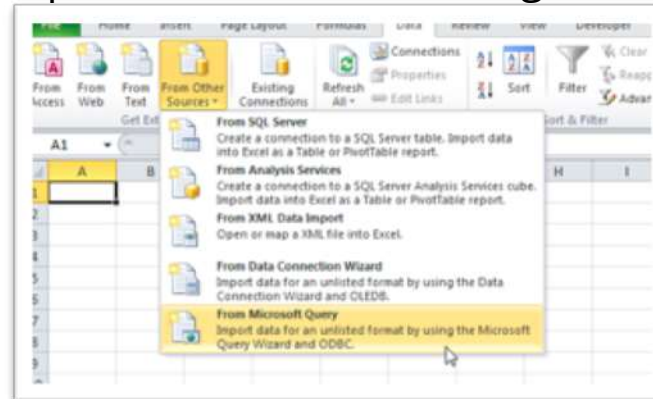
Dedicated BI tool (Tableau, PowerBI, Yellowfin)



The screenshot shows the Yellowfin web interface. At the top, there's a navigation bar with the Yellowfin logo. Below it, a table displays data for 'Med_control_oweb201'. The table has columns: 'controlpoint', 'control_point', 'control_type', 'control_value', 'control_status', 'control_date', and 'control_desc'. The data is organized into rows, with the first row showing values like 4, 274, 1, 105, 10, and 510.

| controlpoint | control_point | control_type | control_value | control_status | control_date | control_desc |
|--------------|---------------|--------------|---------------|----------------|--------------|--------------|
| 11_11 | 4,274 | 1,105 | 10 | 510 | | |
| 11_12 | 488 | 11 | 17 | 24 | | 484 |
| 11_13 | 36 | 13 | 17 | 2 | | 24 |
| 11_14 | 15 | 5 | 2 | 4 | | 4 |
| 11_15 | 235 | 80 | 82 | 28 | | 50 |
| 11_16 | 87 | 4 | 11 | 14 | | 38 |
| 11_17 | 1,218 | 239 | 141 | 125 | | 714 |
| 11_18 | | | | | | |
| 11_19 | 6 | | | | | |

Spreadsheets (Excel, Google Sheets, Libre Office)

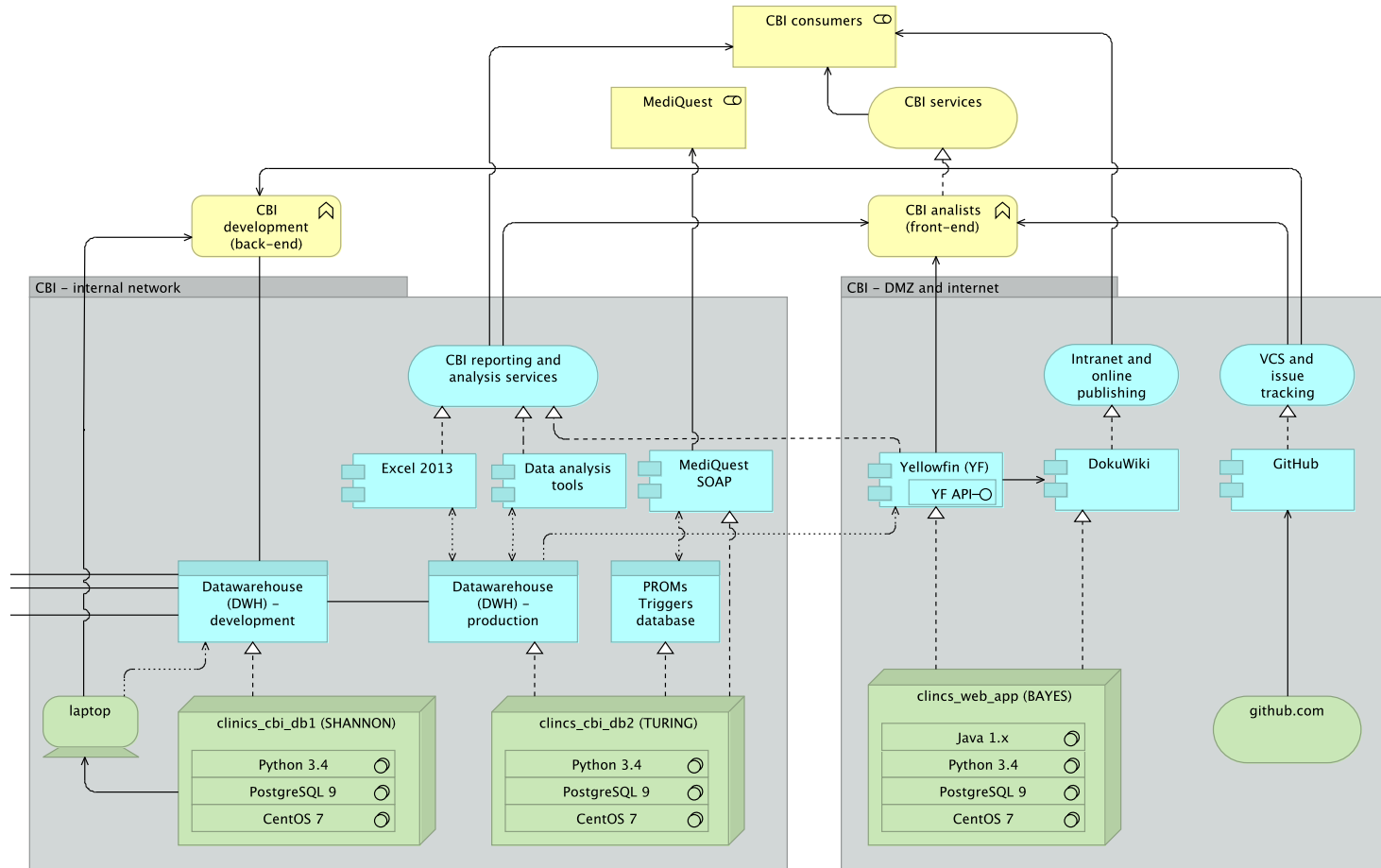


Analytical software (jupyter notebooks, SAS, SPSS)



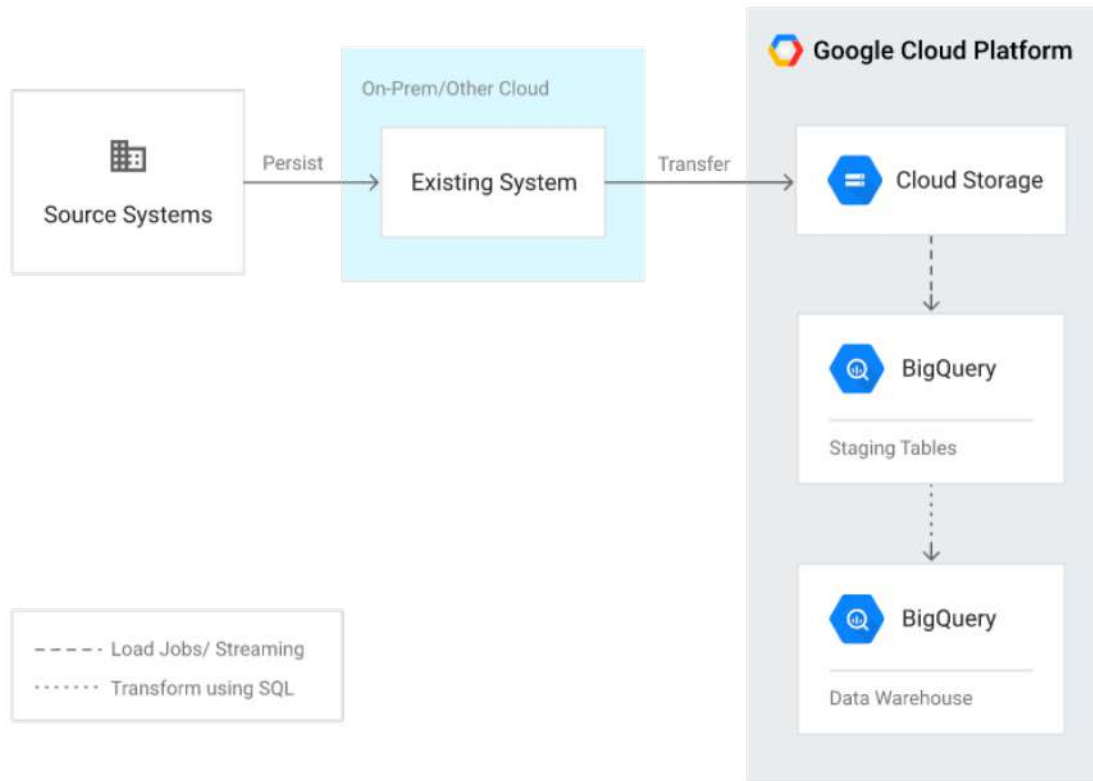
(... and many more)

End-result: with own hardware



- Python as the core language for building datavault and analytics
- PostgreSQL database as storage engine
- Infrastructure: virtualized CentOS with SSD SAN

End-result: on Google Cloud Platform



- Python and Clojure as the core languages
- BigQuery and Cloud storage as storage engines
- see <https://cloud.google.com/solutions/bigquery-data-warehouse> for more details

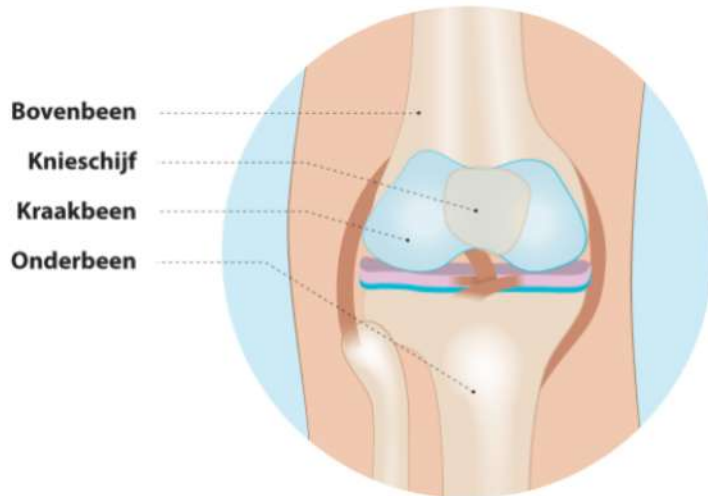
Part II: Healthcare Analytics

PREDICTING OUTCOMES

Imagine you are consulting an orthopedic surgeon because of knee osteoarthritis

Kniegewricht

Het kniegewricht bestaat uit het bot in het bovenbeen, het grote bot in het onderbeen en de knieschijf. Tussen de botten zit kraakbeen (afbeelding). Het zachtere kraakbeen zorgt ervoor dat de botten makkelijk langs elkaar kunnen bewegen.

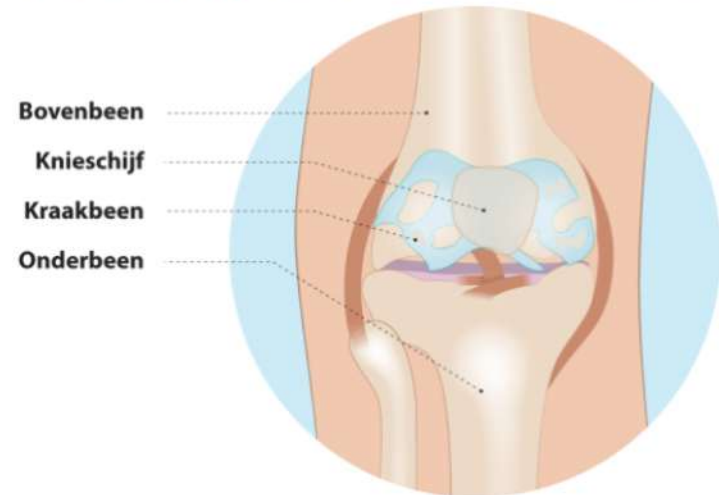


Afbeelding: een gezond kniegewricht.

bron: <https://www.keuzehulp.info/>

Een versleten knie

Als je ouder wordt, verslijt je kraakbeen. Dit wordt een versleten knie genoemd, of knie-artrose. Bij een versleten knie is het kraakbeen bijna helemaal verdwenen (afbeelding). De botten kunnen niet makkelijk meer langs elkaar bewegen. Dit zorgt voor een pijnlijke en stijve knie. Een versleten knie komt veel voor.¹



Afbeelding: een versleten kniegewricht.

You have a choice: operate (new knee) or conservative treatment

Operation:

90 out of 100 patient have less pain.
They are also more active.



Conservative treatment:

Half of all patient have less pain after
physiotherapy and taking painkillers.



bron: <https://www.keuzehulp.info/>

What if an algorithm says that in your case, chance of successful operation is also 50%

Operation:

Chance of success of 50%




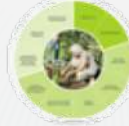


Conservative treatment:

Half of all patient have less pain after physiotherapy and taking painkillers.



So, what are outcomes in healthcare?

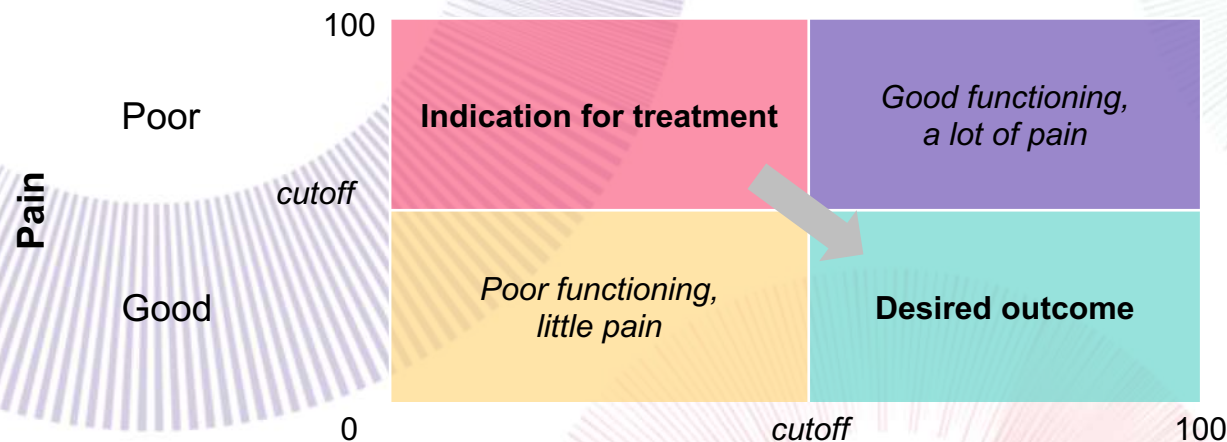
| Outcome category | Condition | | | |
|-------------------|---|--|---|---|
| | Cataract | Macular Degeneration | Low Back Pain | Hip & Knee Osteoarthritis |
| |  |  |  |  |
| PACs* | <ul style="list-style-type: none"> • Re-operation • Endophthalmitis • Corneal oedema | <ul style="list-style-type: none"> • Endophthalmitis | <ul style="list-style-type: none"> • Mortality • Readmissions • Postop infections | <ul style="list-style-type: none"> • Mortality • Readmissions • Postop infections |
| Patient-reported | <ul style="list-style-type: none"> • Catquest-9SF | <ul style="list-style-type: none"> • Brief IVI | <ul style="list-style-type: none"> • EQ-5D • Oswestry Disability Index • NRS pain score • Work status | <ul style="list-style-type: none"> • EQ-5D • KOOS/HOOS • NRS pain score • Satisfaction • Work status |
| Clinical reported | <ul style="list-style-type: none"> • Best corrected visual acuity • Refraction | <ul style="list-style-type: none"> • Best corrected visual acuity • Refraction | | <ul style="list-style-type: none"> • Timed-Up and Go |

*Potentially avoidable complications

Project Nightingale.

1. Compounded outcome measures relevant for shared-decision making, using existing data dictionaries (ICHOM, national registries)
2. Supervised learning for e.g. identifying high-risk patients prior to an intervention
3. 'Unboxing the black box' by relating results of machine learning to existing epidemiological research

A simple idea: choose the two most relevant indicators and determine cutoff values for each



Outcome total knee replacement in the UK

(data NHS Digital | n=140.000 | period 2011 – 2017 | 284 providers)

Prior to operation

| | | Functioning | |
|------|------|-------------|------|
| | | Poor | Good |
| Pain | Poor | 83% | 9% |
| | Good | 3% | 5% |

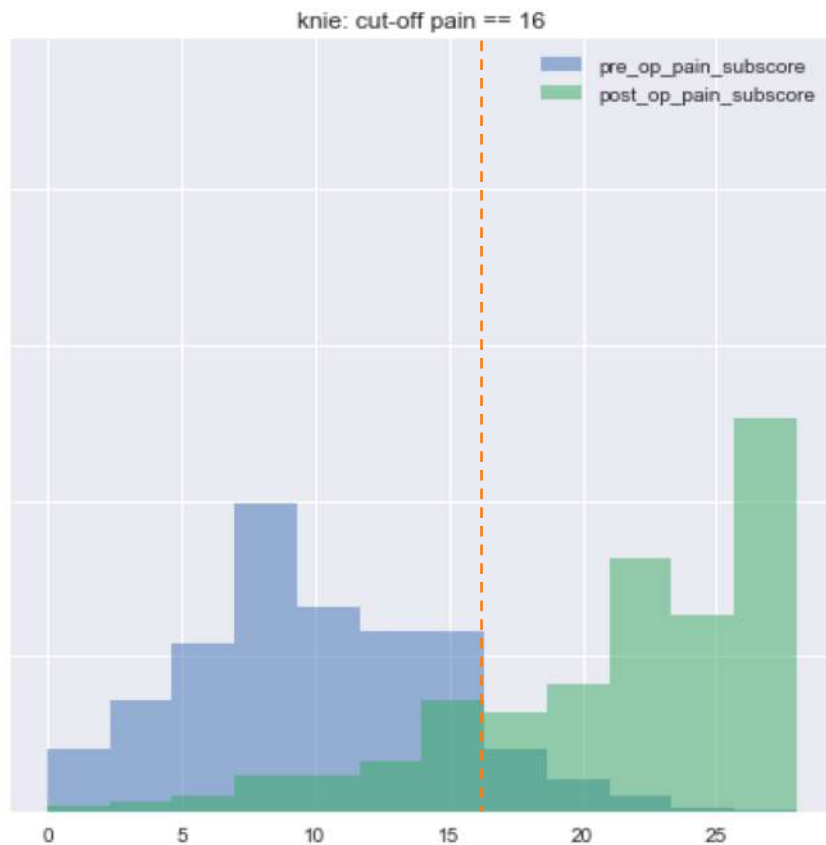


After operation

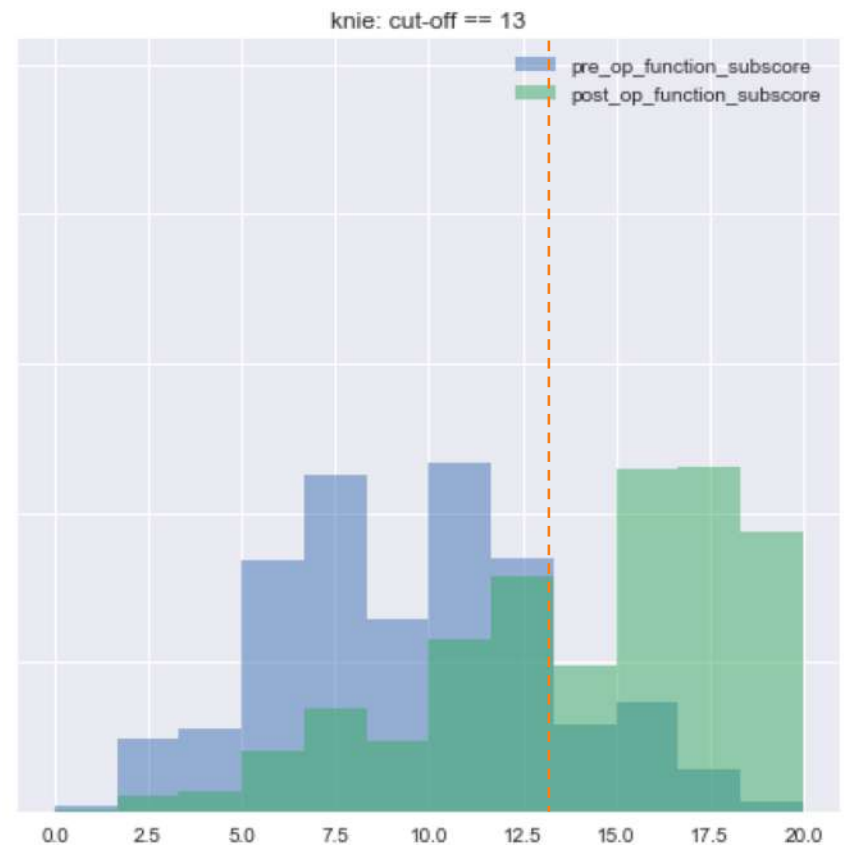
| | | Functioning | |
|------|------|-------------|------|
| | | Poor | Good |
| Pain | Poor | 18% | 1% |
| | Good | 20% | 61% |

Measuring outcomes is not always straightforward

Pain



Functioning



Prior to cataract surgery.

| Pre-operative (T0) | | Visual acuity | |
|-----------------------|------|---------------|------|
| | | Poor | Good |
| PROMs | Poor | 50% | 24% |
| | Good | 15% | 11% |

- 50% of patients have consistent indication (poor-poor)
- ... but how about the other half?

After cataract surgery.

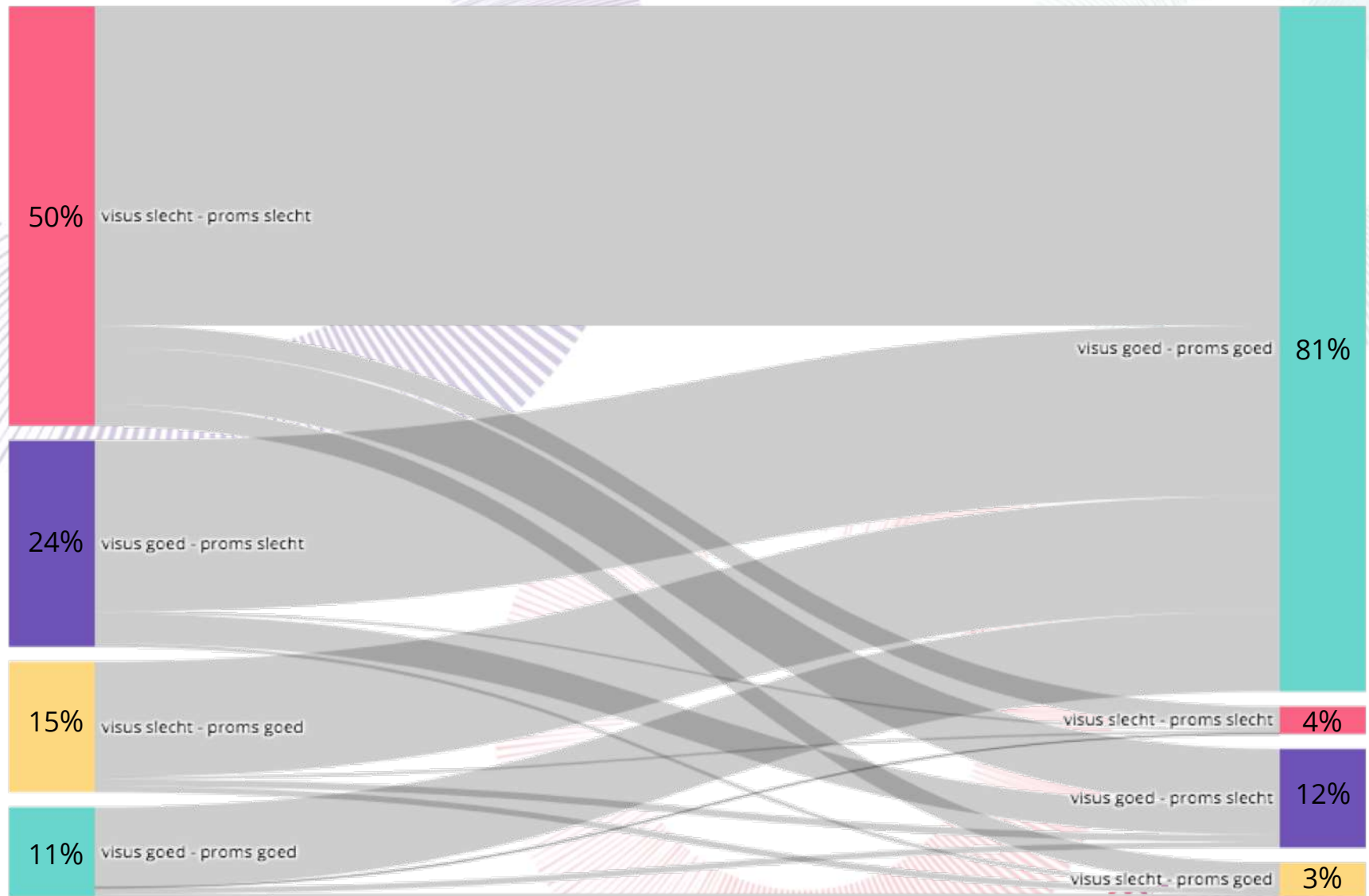
| Pre-operative (T0) | | Visual acuity | |
|-----------------------|------|---------------|------|
| | | Poor | Good |
| PROMs | Poor | 50% | 24% |
| | Good | 15% | 11% |

- 50% of patients have consistent indication (poor-poor)
- ... but how about the other half?

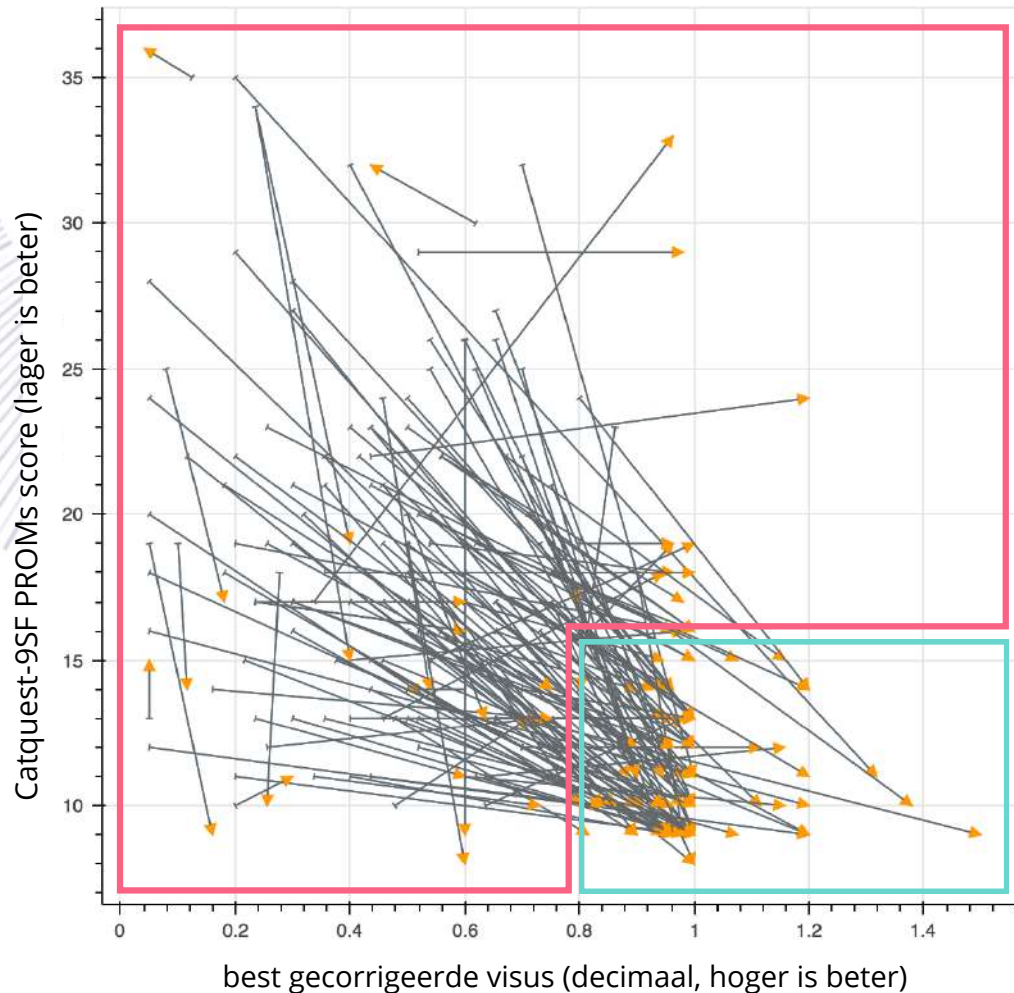
| Post-operative (T1) | | Visual acuity | |
|------------------------|------|---------------|------------|
| | | Poor | Good |
| PROMs | Poor | 3% | 12% |
| | Good | 4% | 81% |

- Good outcome for 81% of all patients
- Remaining 'outliers' of 19% require more detailed inspection

No simple, linear mapping between pre to post



Can we predict the outcome prior to surgery?



- **Sensitivity 0.5**
Half of the arrows that end up in the red quadrants can be identified prior to surgery; 9% of all patient receive correct warning signal
- **Positive predictive value 0.58**
I.e. 42% of warning signals is false-positive. Good enough?

Do we understand what the algorithm does?

Risk factors known from literature

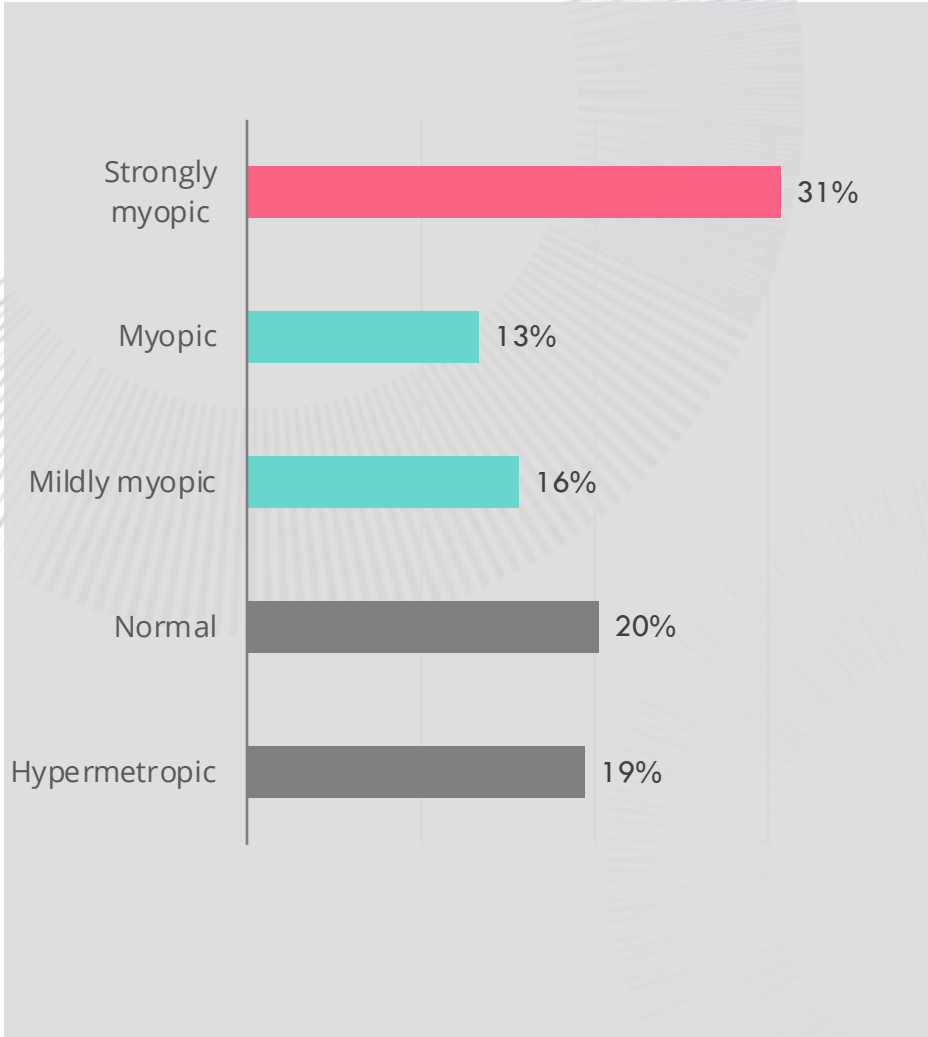
Post-operative complications
Best corrected visual acuity
Target refraction
Capsule complications
Ocular co-morbidities
PROMs sumscore
Gender
Age

- Lundström, M. and Stenevi, U., *Analyzing Patient-Reported Outcomes to Improve Cataract Care*, Optometry and Vision Science 2013: vol. 90 no. 8: 754-759
- Grimfors et al., *Ocular comorbidity and self-assessed visual function after cataract surgery*, J. Cataract Refract Surg 2014; 40:1163-1169
- Lundström et al., *Visual outcome of cataract surgery*, J. Cataract Refract Surg 2013: 39:673-679
- Mollazadegan, K. and Lundström, M., *A study of the correlation between patient-reported outcomes and clinical outcome after cataract surgery in ophthalmic clinics*, Acta Ophthalmol. 2015: 93: 293-298

Relative feature importance in random forest

Target refraction
Age
Best gecorrected visual acuity
PROMs sumscore
PROMs sub-score near-vision
PROMs sub-score far-vision
PROMs sub-score general satisfaction
Individual PROMs items
Gender
Previous surgery other eye
Macular degeneration
Other ocular co-morbidities

Significant effect size in outcome by target refraction (strength of glasses)



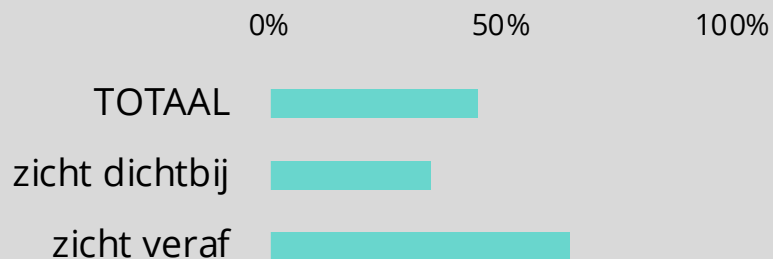
Percentage poor outcome by target refraction, i.e. chosen strength of glasses post-surgery

Target refraction group in diopters (n observations)

- Strongly myopic:
 < -4.0 (n=13)
- Myopic:
 between -4.0 and -2.0 (n=771)
- Mildly myopic:
 between -2.0 and -0.5 (n=319)
- Normal:
 between -0.5 and 0.5 (n=3993)
- Hypermetropic:
 > 0.5 (n=36)



Bent u tevreden met uw zicht?



Wat vindt u belangrijk?

Activiteiten en hobbies met:

- ☐ zicht dichtbij
- ☐ zicht veraf

Ik wil dit **met/zonder** bril kunnen

Uw ogen op dit moment

| | L | R |
|----------------------|--------------------|------|
| zicht met bril: | 0.6 | 0.4 |
| brilsterkte: | -1.5 | -2.0 |
| andere aandoeningen: | macula degeneratie | |

Verwachte uitkomst operatie

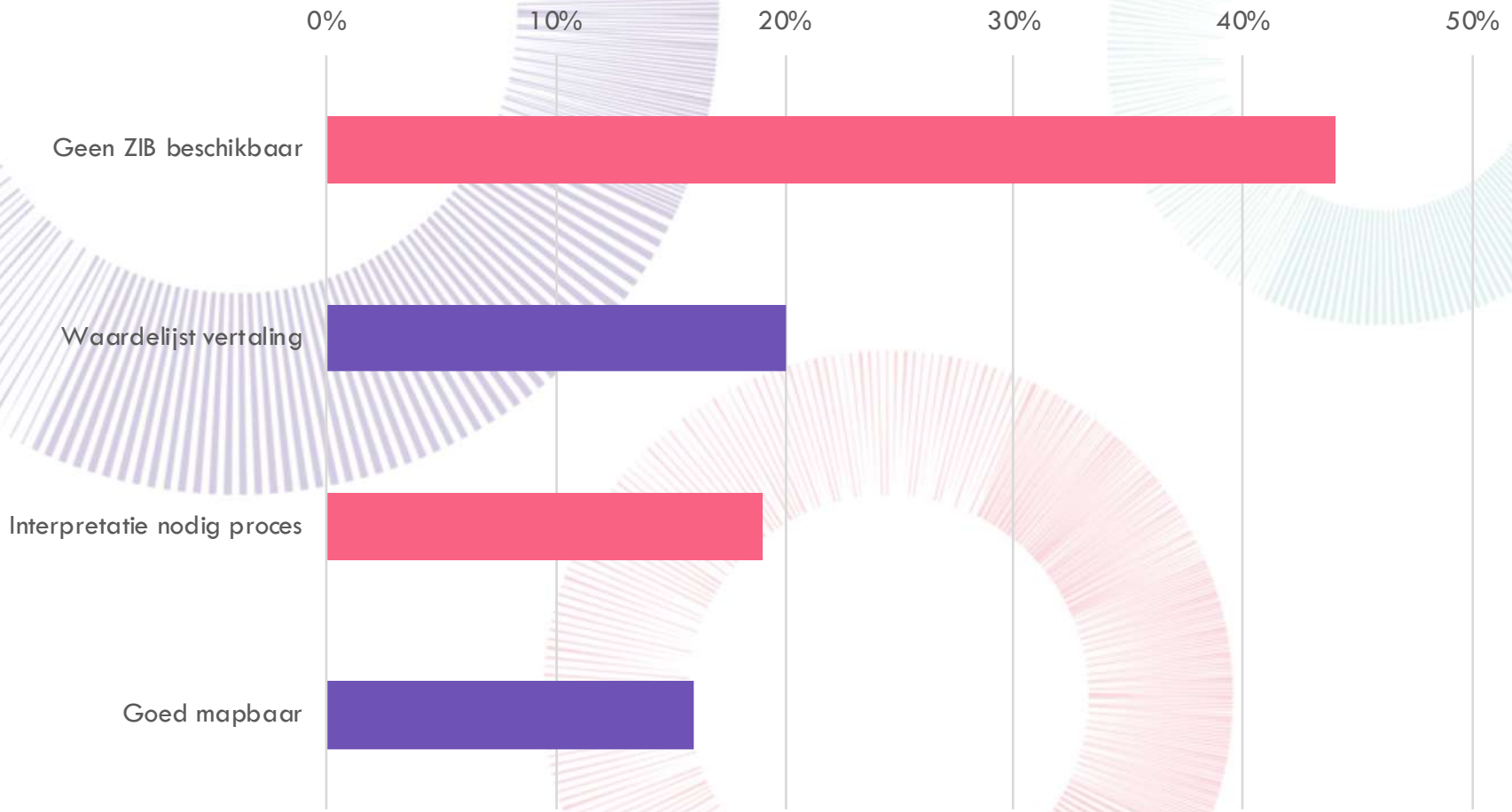
| | L | R |
|-----------------|-------------------------|------|
| zicht met bril: | 0.6 | 1.0 |
| brilsterkte: | -1.5 | -0.5 |
| aandachtspunt: | ⚠️ risico eindresultaat | |

Lessons learned from applying machine learning in daily operations

1. Data quality (*registratie aan de bron*)
2. Harmonisation and semantic integration of different registries (e.g. recent analysis Nictiz)
3. Open sourcing trained models for clinical decision support

Mapping ICHOM to Dutch standards

(Analysis Nictiz, August 2018)



More info?

- Drop me an email at dkapitan@mediquest.nl
- Connect on LinkedIn <https://www.linkedin.com/in/dkapitan/>