# Making Process Mining Green
## Using Event Data in a Responsible Way

**Wil van der Aalst**
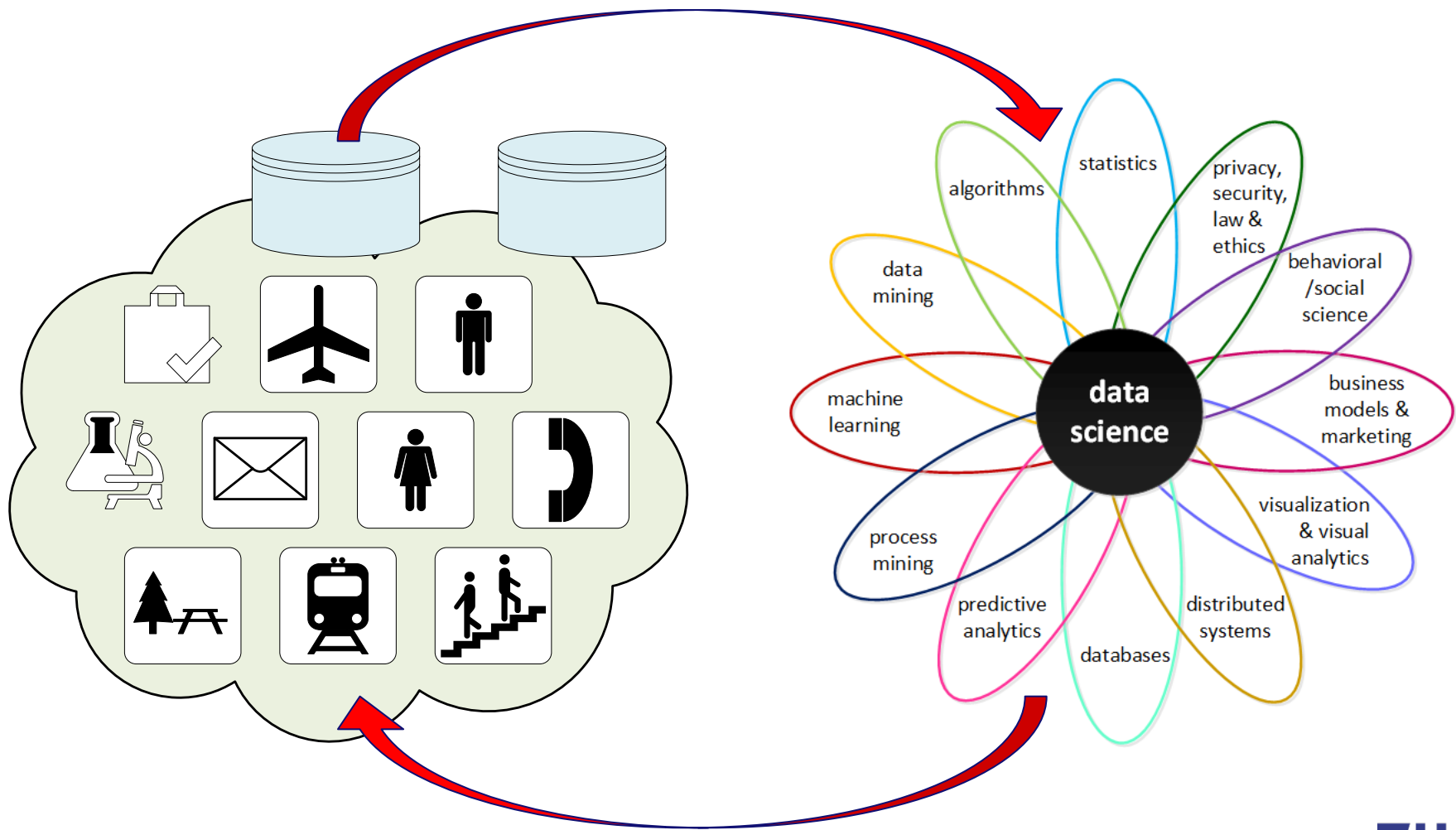*www.vdaalst.com     @wvdaalst*
*www.processmining.org*

TU/e Technische Universiteit
**Eindhoven**
University of Technology

**Where innovation starts**

0101100
1001011
1101110
1011011
0111100
1001011
1101110
1001011
0101110
1001011
0100011
1000001
0101100
1001011
0101111
1001011
0101100

better, faster, more efficient, more effective, cheaper, …

**With Great Power Comes Great Responsibility!!**

If data is the new oil on which our society runs …

non-transparent

unfair use of data

spurious correlations

privacy violations

bogus conclusions

… then we should take care of data-related forms of pollution!

# Green data science: separate the "pollution" from the actual purpose

TU/e

# Two parts

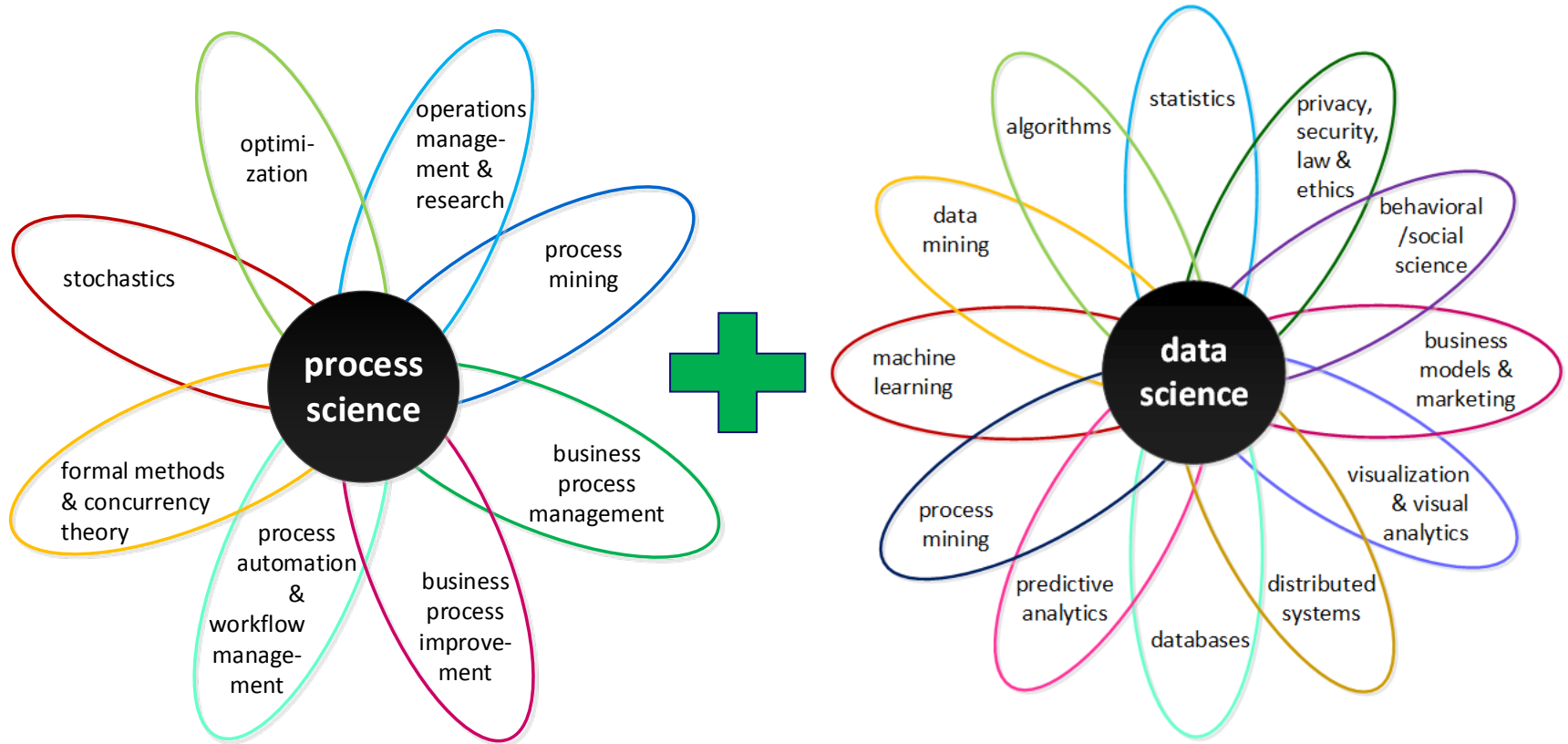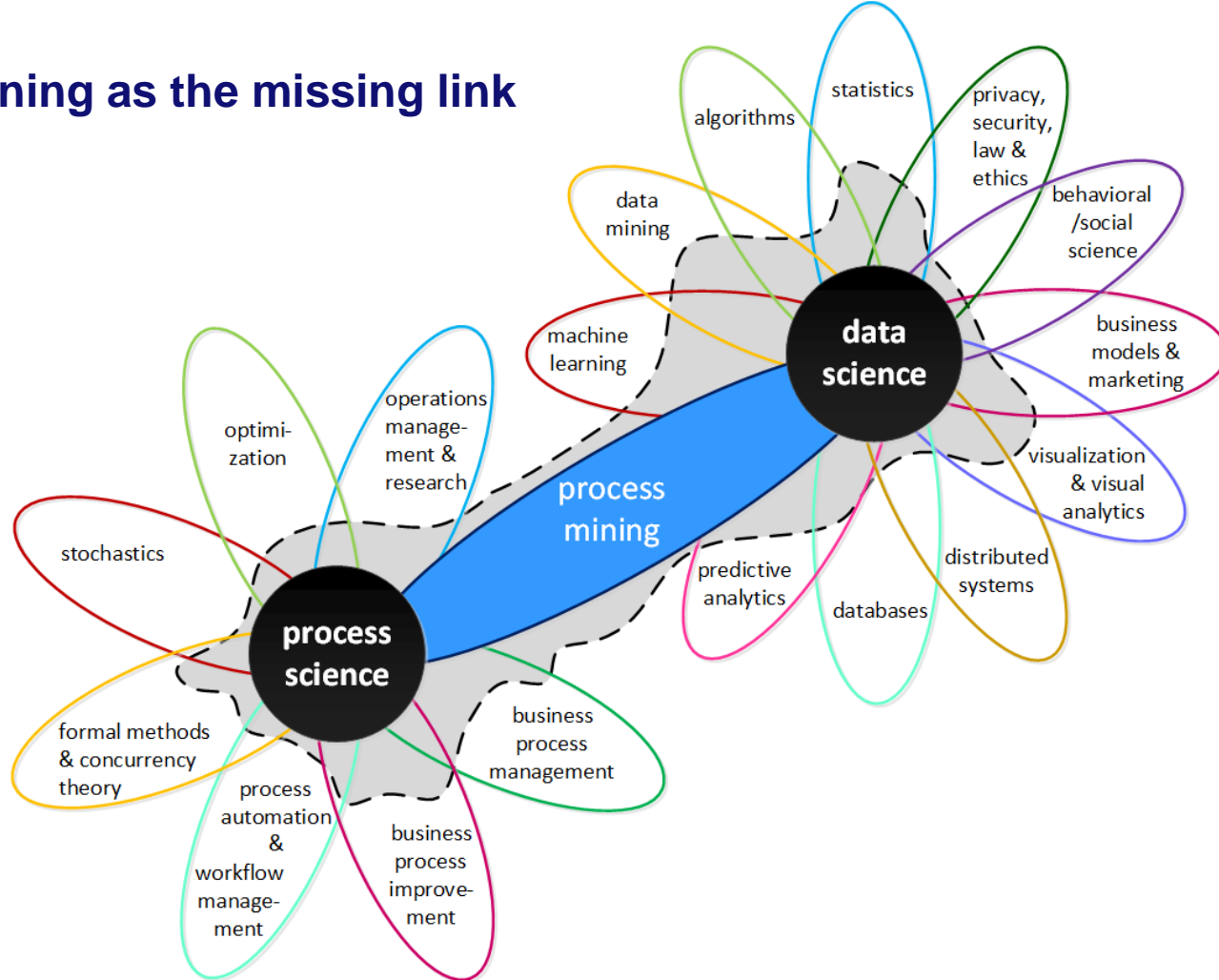**2** responsible data science: our next big challenge

**1** process mining: creating value from data

TU/e

# Part I

process mining: creating value from data

process science + data science

- optimization
- operations management & research
- process mining
- stochastics
- business process management
- formal methods & concurrency theory
- process automation & workflow management
- business process improvement

- statistics
- algorithms
- privacy, security, law & ethics
- data mining
- behavioral /social science
- machine learning
- business models & marketing
- process mining
- visualization & visual analytics
- predictive analytics
- databases
- distributed systems

TU/e

# process mining as the missing link

# Process Mining: On the interface between process science and data science
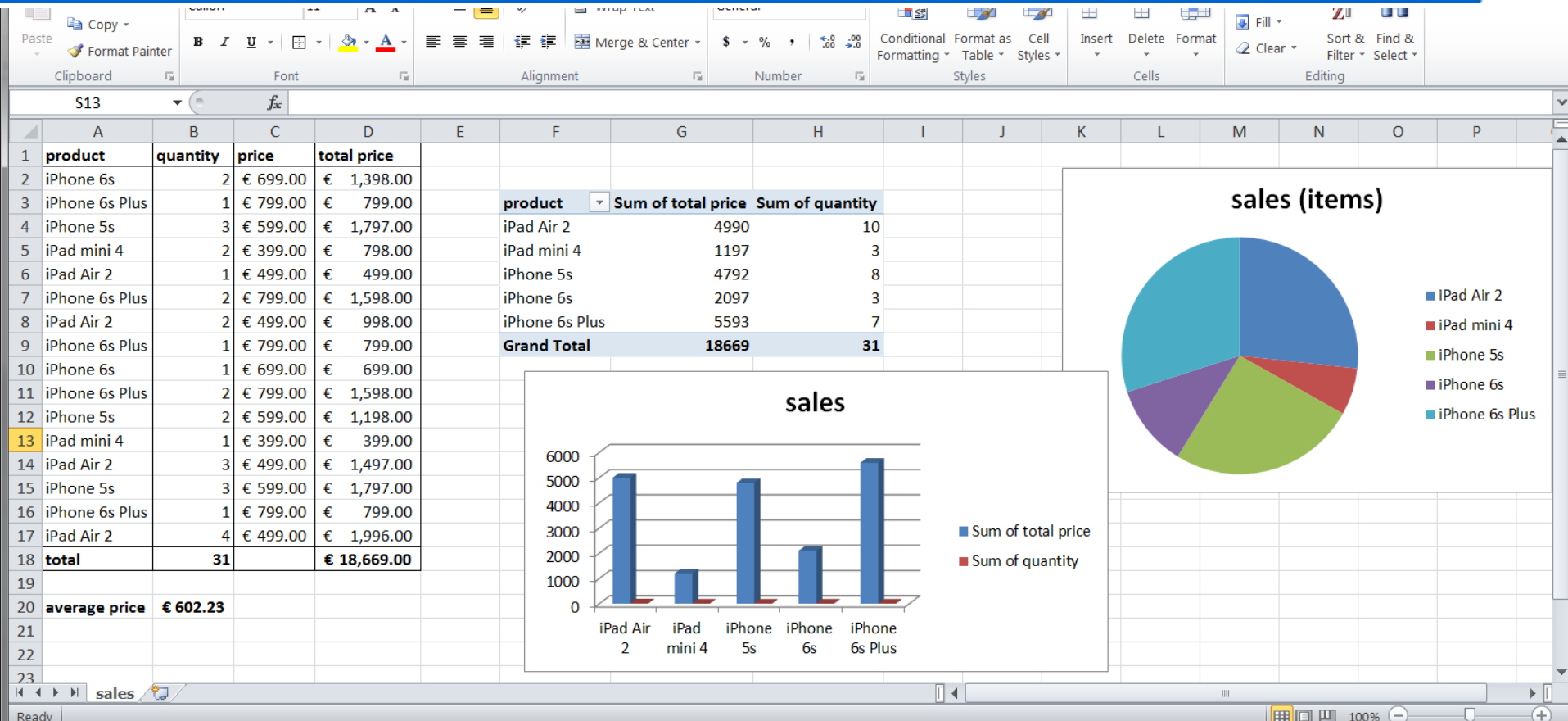
# Spreadsheet: Killer App for early computers



- **VisiCalc (killer app for Apple II, Oct. 1979)**

- **Lotus 1-2-3 (killer app for IBM PC 1983)**

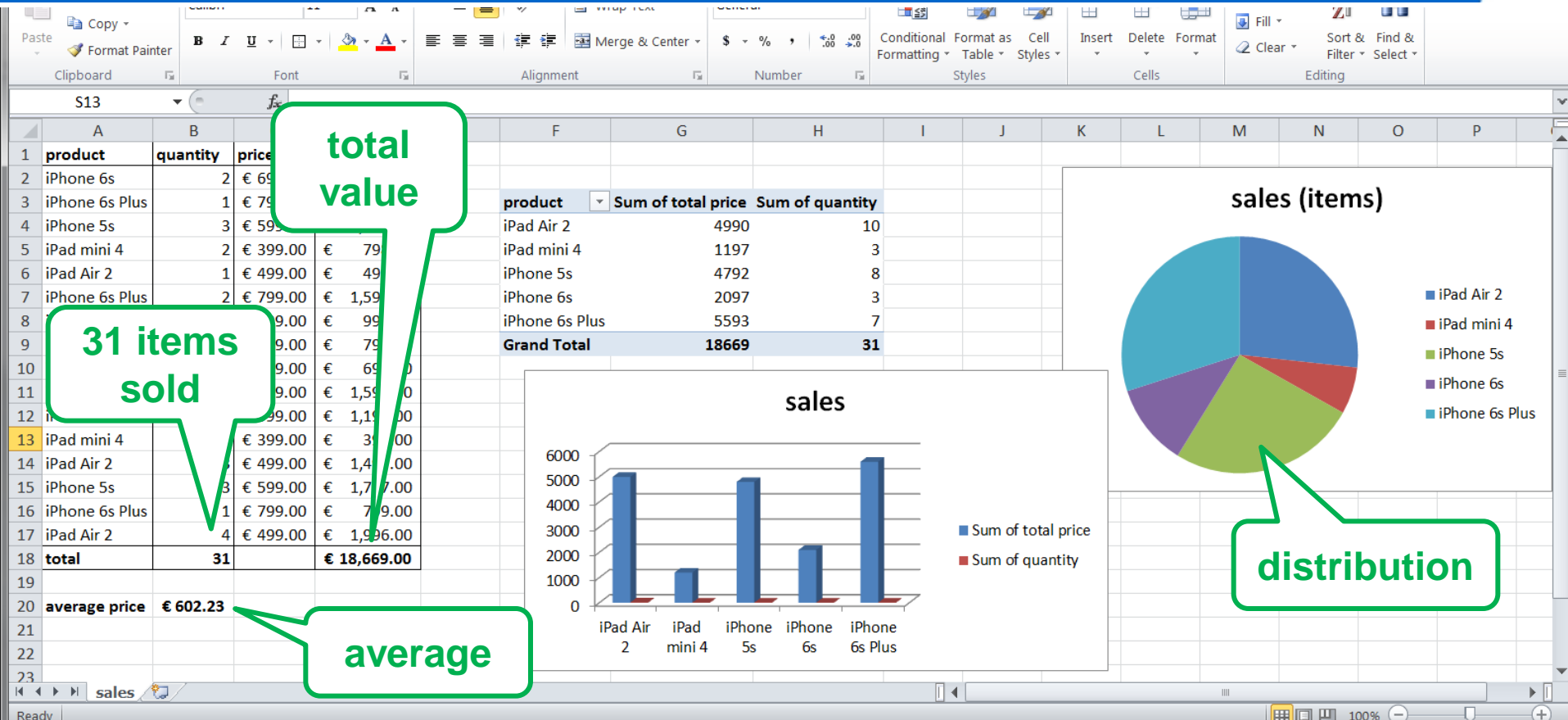- **Microsoft Excel (1985)**

# Spreadsheet: Static data

# Spreadsheet: Static data

# Spreadsheet: Static data



| | A | B | price | | | F | G | H |
|---|---|---|---|---|---|---|---|---|
| 1 | product | quantity | | | | product | Sum of total price | Sum of quantity |
| 2 | iPhone 6s | 2 | € 69 | | | iPad Air 2 | 4990 | 10 |
| 3 | iPhone 6s Plus | 1 | € 79 | | | iPad mini 4 | 1197 | 3 |
| 4 | iPhone 5s | 3 | € 599 | | | iPhone 5s | 4792 | 8 |
| 5 | iPad mini 4 | 2 | € 399.00 | € 79 | | iPhone 6s | 2097 | 3 |
| 6 | iPad Air 2 | 1 | € 499.00 | € 49 | | iPhone 6s Plus | 5593 | 7 |
| 7 | iPhone 6s Plus | 2 | € 799.00 | € 1,59 | | Grand Total | 18669 | 31 |

total value

31 items sold

distribution

average

average price  € 602.23

sales

# Process Mining: Spreadsheet for behavior



- **Input: events ("things that have happened")**
- **Mandatory per event:**
  - **case identifier**
  - **activity name**
  - **timestamp/date**
- **Optional**
  - **resource**
  - **transaction type**
  - **costs**
  - **…**

# Process Mining: Spreadsheet for behavior



208 cases
5987 events
74 activities

# Process Mining: Spreadsheet for behavior



batching for activities "opstellen eindnota" and "archiveren"

# Process Mining: Spreadsheet for behavior

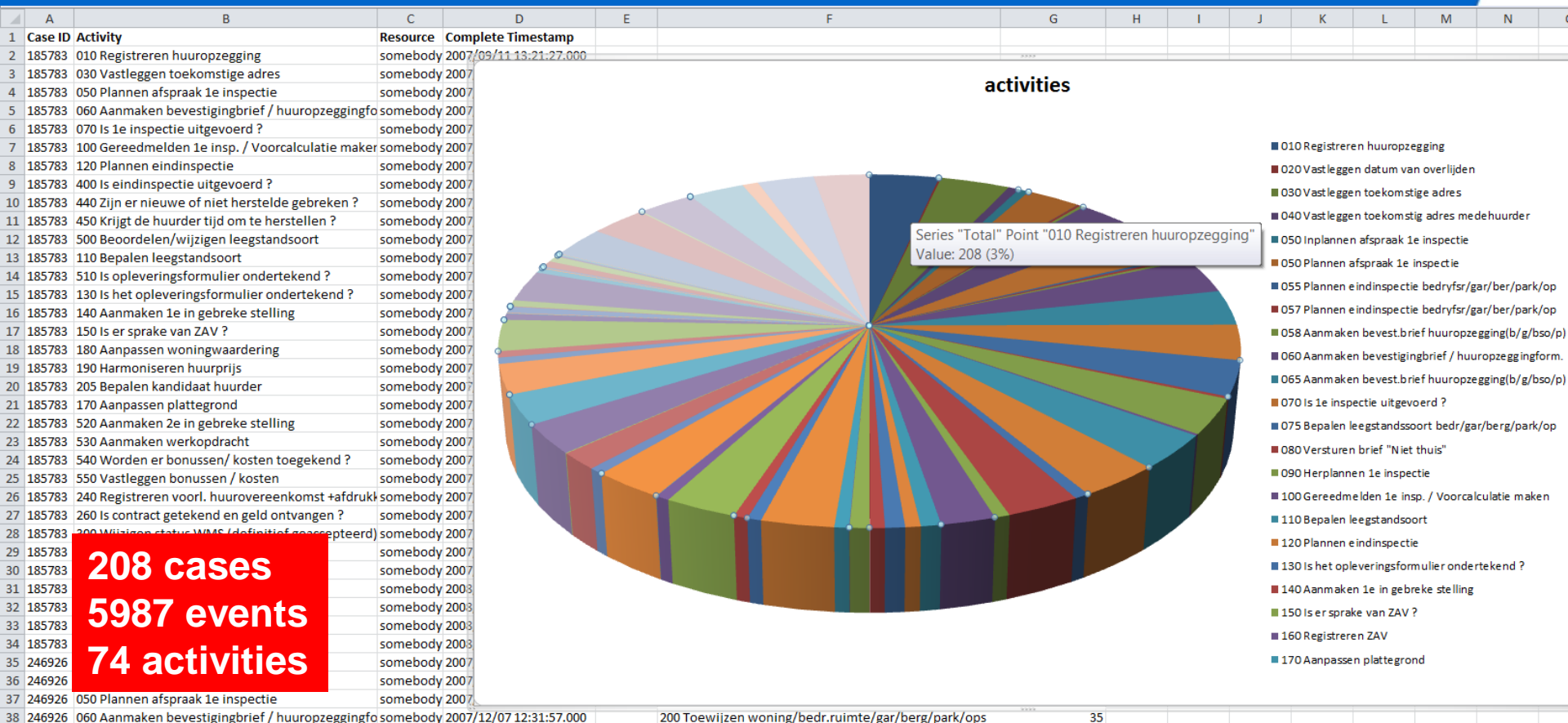| | A | B | C | D |
|---|---|---|---|---|
| 1 | Case ID | Activity | Resource | Complete Timestamp |
| 2 | 185783 | 010 Registreren huuropzegging | somebody | 2007/09/11 13:21:27.000 |
| 3 | 185783 | 030 Vastleggen toekomstige adres | somebody | 2007/09/11 13:26:29.000 |
| 4 | 185783 | 050 Plannen afspraak 1e inspectie | somebody | 2007/09/11 13:29:34.000 |
| 5 | 185783 | 060 Aanmaken bevestigingbrief / huuropzeggingform. | somebody | 2007/09/11 13:41:36. |
| 6 | 185783 | 070 Is 1e ... itgevoerd ? | somebody | 2007/09/24 08:39:32 |
| 7 | 185783 | 100 ... / Voorcalculatie maken | somebody | 2007/09/24 08:41:2 |
| 8 | 185783 | 12 | somebody | 2007/09/24 08:51:00 |
| 9 | 185783 | 4 ... d ? | somebody | 2007/09/24 10:55:56 |
| 10 | 185783 | 4 ... stelde gebreken ? | somebody | 2007/09/24 10:56:06 |
| 11 | 185783 | 4 ... te herstellen ? | somebody | 2007/09/24 10:56:10 |
| 12 | 185783 | 50 ... egstandsoort | somebody | 2007/09/24 10:57:02 |
| 13 | 185783 | 110 ... ort | somebody | 2007/09/24 10:57:4 |
| 14 | 185783 | 510 Is op ... mulier ondertekend ? | somebody | 2007/09/24 10:58:08.000 |
| 15 | 185783 | 130 Is het opleveringsformulier ondertekend ? | somebody | 2007/09/24 10:58:19.000 |
| 16 | 185783 | 140 Aanmaken 1e in gebreke stelling | somebody | 2007/09/24 11:01:58.000 |
| 17 | 185783 | 150 Is er sprake van ZAV ? | somebody | 2007/09/24 11:37:33.000 |
| 18 | 185783 | 180 Aanpassen woningwaardering | somebody | 2007/09/24 11:37:44.000 |
| 19 | 185783 | 190 Harmoniseren huurprijs | somebody | 2007/09/24 11:40:01.000 |
| 20 | 185783 | 205 Bepalen kandidaat huurder | somebody | 2007/09/24 11:47:42.000 |
| 27 | 185783 | 260 Is contract getekend en geld ontvangen ? | somebody | 2007/12/10 10:44:06.000 |
| 28 | 185783 | 300 Wijzigen status WMS (definitief geaccepteerd) | somebody | 2007/12/11 16:26:14.000 |
| 29 | 185783 | 560 Opstellen eindnota | somebody | 2007/12/12 11:19:41.000 |

**NO modeling needed!**

*process discovery*

# Process Mining: Spreadsheet for behavior
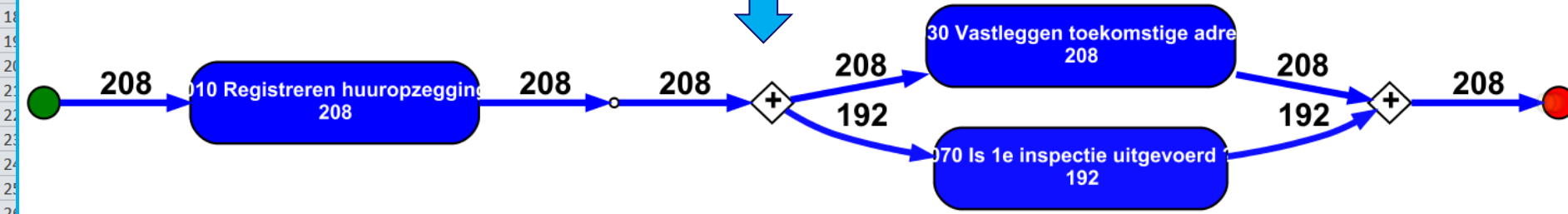
| | A | B | C | D |
|---|---|---|---|---|
| 1 | Case ID | Activity | Resource | Complete Timestamp |
| 2 | 185783 | 010 Registreren huuropzegging | somebody | 2007/09/11 13:21:27.000 |

*process discovery*



**NO modeling needed!**

| | | | | |
|---|---|---|---|---|
| 9 | 185783 | 400 Is eindinspectie uitgevoerd ? | somebody | 2007/09/24 10:55:56. |
| 16 | 185783 | 140 Aanmaken 1e in gebreke stelling | somebody | 2007/09/24 11:01:58. |
| 17 | 185783 | 150 Is er sprake van ZAV ? | somebody | 2007/09/24 11:37:33. |
| 27 | 185783 | 260 Is contract getekend en geld ontvangen ? | somebody | 2007/12/10 10:44:06.000 |
| 28 | 185783 | 300 Wijzigen status WMS (definitief geaccepteerd) | somebody | 2007/12/11 16:26:14.000 |
| 29 | 185783 | 560 Opstellen eindnota | somebody | 2007/12/12 11:19:41.000 |

process model

event data

Conformance Checking

# Process Mining: Spreadsheet for behavior



*conformance checking*

**?**

**discovered or hand-made**

# Process Mining: Spreadsheet for behavior



*conformance checking*

**fitness of 93.5%**

# Process Mining: Spreadsheet for behavior



*conformance checking*

**final inspection is skipped 40 times**

# Process Mining: Spreadsheet for behavior



*conformance checking*

**move on model**
**(something should have**
**happened, but did not)**

**move on log**
**(something happened that**
**should not happen)**

# Process Mining: Spreadsheet for behavior



*performance analysis*

**NO modeling needed!**

**bottleneck**

**average flowtime is 1.92 months**

| Case Property | Value |
|---|---|
| #Cases | 208 |
| #Perfectly-fitting cases | 142 |
| #Non-fitting cases | 66 |
| #Properly started cases | 208 |
| Case Throughput time (avg) | 1.92 months |
| Case Throughput time (min) | 7.80 min |
| Case Throughput time (max) | 10.86 months |
| Case Throughput time (std. dev) | 1.50 months |
| Observation period | 18.81 months |

Process Mining: Spreadsheet for behavior

# Process Mining: Spreadsheet for behavior



*animating reality*

**real cases**

**NO modeling needed!**

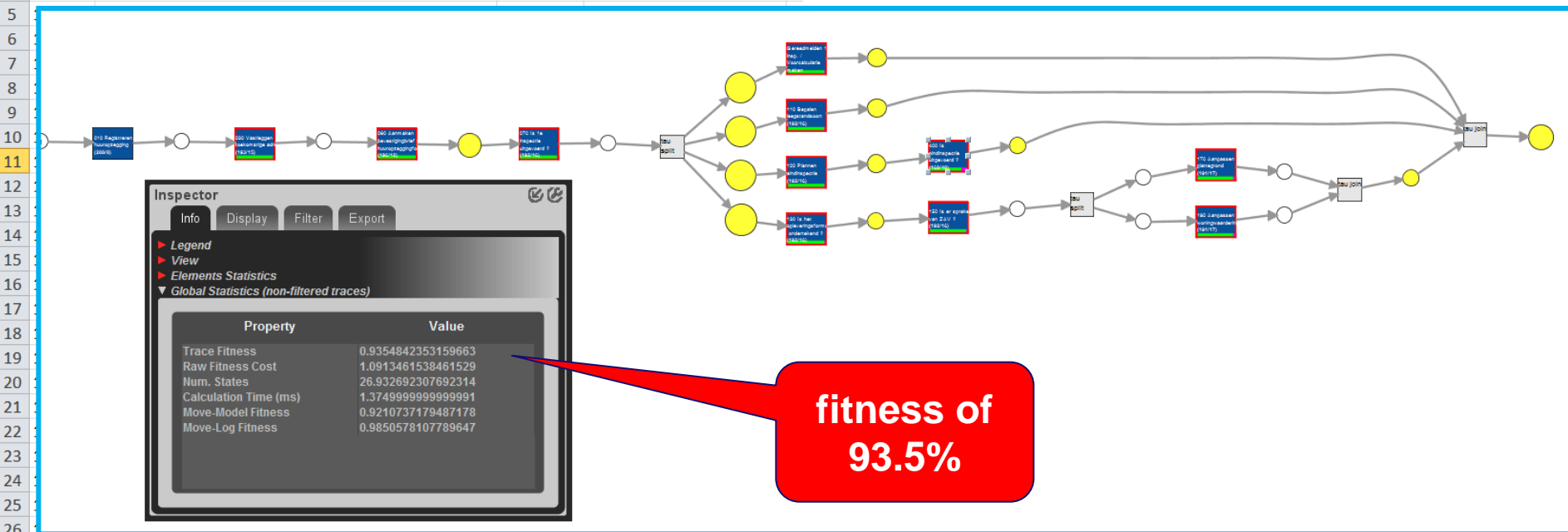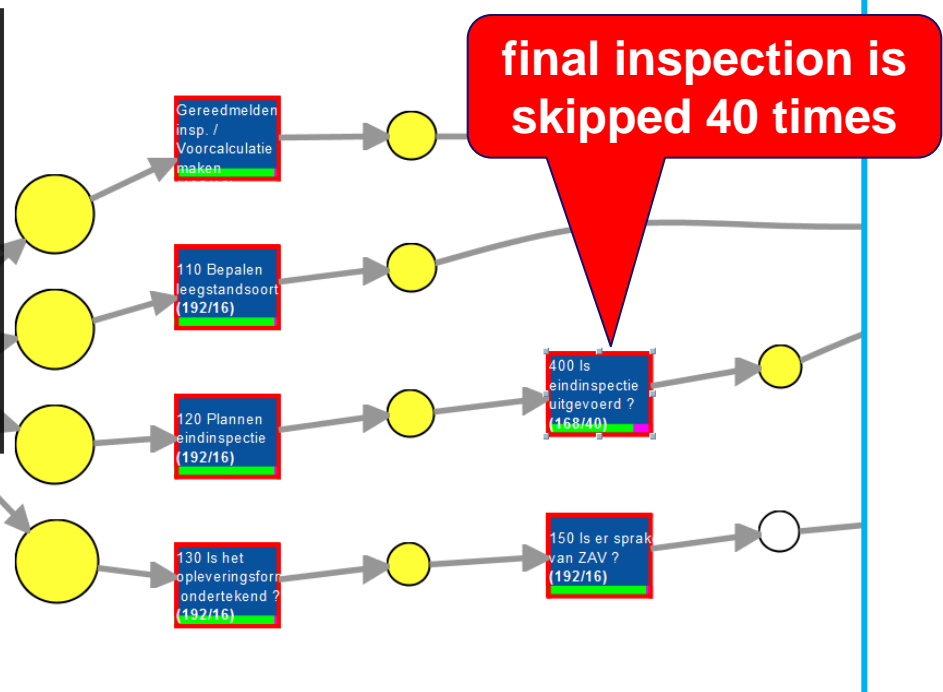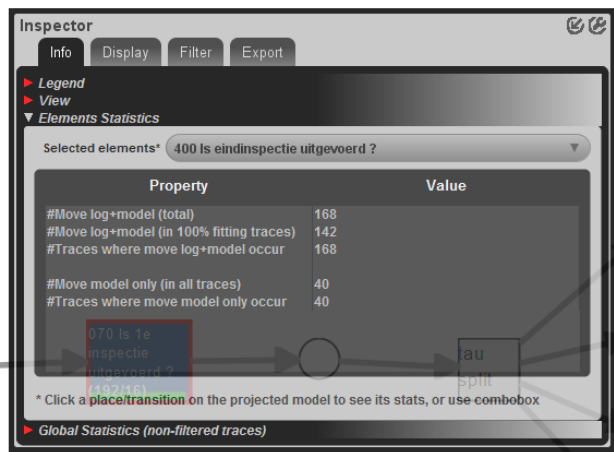| | A | B | C | D |
|---|---|---|---|---|
| 1 | Case ID | Activity | Resource | Complete Timestamp |
| 2 | 185783 | 010 Registreren huuropzegging | somebody | 2007/09/11 13:21:27.000 |
| 3 | 185783 | 030 Vastleggen toekomstige adres | somebody | 2007/09/11 13:26:29.000 |
| 4 | 185783 | 050 Plannen afspraak 1e inspectie | somebody | 2007/09/11 13:29:34.000 |
| 5 | | | | |
| 6 | | | | |
| 7 | | | | |
| 8 | | | | |
| 9 | | | | |
| 10 | | | | |
| 11 | | | | |
| 12 | 185783 | 500 Beoordelen/wijzigen leegstandsoort | somebody | 2007/09/24 10:57:02.000 |
| 13 | 185783 | 110 Bepalen leegstandsoort | somebody | 2007/09/24 10:57:42.000 |
| 14 | 185783 | 510 Is opleveringsformulier ondertekend ? | somebody | 2007/09/24 10:58:08.000 |
| 15 | 185783 | 130 Is het opleveringsformulier ondertekend ? | somebody | 2007/09/24 10:58:19.000 |
| 16 | 185783 | 140 Aanmaken 1e in gebreke stelling | somebody | 2007/09/24 11:01:58.000 |
| 17 | 185783 | 150 Is er sprake van ZAV ? | somebody | 2007/09/24 11:37:33.000 |
| 18 | 185783 | 180 Aanpassen woningwaardering | somebody | 2007/09/24 11:37:44.000 |
| 19 | 185783 | 190 Harmoniseren huurprijs | somebody | 2007/09/24 11:40:01.000 |
| 20 | 185783 | 205 Bepalen kandidaat huurder | somebody | 2007/09/24 11:47:42.000 |
| 21 | 185783 | 170 Aanpassen plattegrond | somebody | 2007/09/24 12:10:58.000 |
| 22 | 185783 | 520 Aanmaken 2e in gebreke stelling | somebody | 2007/10/30 11:45:53.000 |
| 23 | 185783 | 530 Aanmaken werkopdracht | somebody | 2007/10/30 11:46:09.000 |
| 24 | 185783 | 540 Worden er bonussen/ kosten toegekend ? | somebody | 2007/10/30 11:46:36.000 |
| 25 | 185783 | 550 Vastleggen bonussen / kosten | somebody | 2007/10/30 11:53:00.000 |
| 26 | 185783 | 240 Registreren voorl. huurovereenkomst +afdrukken | somebody | 2007/11/28 12:34:23.000 |
| 27 | 185783 | 260 Is contract getekend en geld ontvangen ? | somebody | 2007/12/10 10:44:06.000 |
| 28 | 185783 | 300 Wijzigen status WMS (definitief geaccepteerd) | somebody | 2007/12/11 16:26:14.000 |
| 29 | 185783 | 560 Opstellen eindnota | somebody | 2007/12/12 11:19:41.000 |

TU/e

# Process Mining: Spreadsheet for behavior

Deviations
Where?
Why?

time
costs
…

# Process Mining Software

1500+ plug-ins available covering the whole process mining spectrum

100% FREE

>128k downloads

ProM
process mining workbench

©Wil van der Aalst & TU/e (use only with permission & acknowledgements)

event_log_10000_cases  +17

Map  Statistics  Cases

Academic
w.m.p.v.d.aalst@tue.nl

Disco

## Statistics views

- **Overview** — Global statistics
- **Activity** — Activity classes
- **product** — Other attribute
- **prod-price** — Other attribute
- **quantity** — Ot...
- **ad...** — Ot...

### Overview
Global statistics

Events over time

Active cases over time

Case variants

Events per case

Case duration

Log timeline

| | Events | 63,763 |
| | Cases | 10,000 |
| | Activities | 8 |
| | Median case duration | 13.9 d |
| | Mean case duration | 14.9 d |
| | Start | 05.01.2015 09:00:07 |
| | End | 31.12.2019 14:46:02 |

Cases (10000)   Variants (9)

| Variant | Started | Finished | Duration |
|---|---|---|---|
| 7 | 05.01.2015 09:00:07 | 26.01.2015 16:42:28 | 21 days, 7 hours |
| 1 | 05.01.2015 10:18:21 | 15.01.2015 15:52:30 | 10 days, 5 hours |
| 4 | 05.01.2015 11:54:49 | 09.01.2015 18:38:58 | 4 days, 6 hours |
| 3 | 05.01.2015 14:07:45 | 22.01.2015 13:18:30 | 16 days, 23 hours |
| 1 | 05.01.2015 15:33:38 | 12.01.2015 17:27:36 | 7 days, 1 hour |
| 5 | 05.01.2015 17:25:23 | 02.02.2015 12:31:09 | 27 days, 19 hours |
| 4 | 05.01.2015 19:08:53 | 15.01.2015 14:56:54 | 9 days, 19 hours |
| 9 | 05.01.2015 21:54:00 | 13.01.2015 15:49:53 | 7 days, 17 hours |
| 4 | 06.01.2015 07:25:13 | 15.01.2015 11:27:50 | 9 days, 4 hours |
| 1 | 06.01.2015 10:09:51 | 15.01.2015 19:15:18 | 9 days, 9 hours |
| 1 | 06.01.2015 11:37:49 | 14.01.2015 09:14:28 | 7 days, 21 hours |
| 4 | 06.01.2015 13:33:45 | 14.01.2015 11:30:05 | 7 days, 21 hours |
| 4 | 06.01.2015 15:25:38 | 13.01.2015 12:25:34 | 6 days, 20 hours |
| 2 | 06.01.2015 17:09:23 | 22.01.2015 18:59:10 | 16 days, 1 hour |
| 3 | 06.01.2015 18:36:53 | 22.01.2015 14:39:39 | 15 days, 20 hours |
| 8 | 06.01.2015 21:26:54 | 26.01.2015 17:16:02 | 19 days, 19 hours |
| 1 | 07.01.2015 04:42:36 | 16.01.2015 10:17:14 | 9 days, 5 hours |
| 3 | 07.01.2015 10:10:58 | 21.01.2015 17:31:29 | 14 days, 7 hours |
| 8 | 07.01.2015 11:40:04 | 28.01.2015 10:27:12 | 20 days, 22 hours |
| 9 | 07.01.2015 13:38:15 | 13.01.2015 13:22:15 | 5 days, 23 hours |
| 1 | 07.01.2015 15:34:37 | 19.01.2015 09:11:23 | 11 days, 17 hours |
| 1 | 07.01.2015 17:27:21 | 16.01.2015 09:09:25 | 8 days, 15 hours |
| 5 | 07.01.2015 19:12:50 | 03.02.2015 14:34:33 | 26 days, 19 hours |
| 6 | 07.01.2015 22:01:54 | 19.01.2015 13:15:02 | 11 days, 15 hours |
| 8 | 08.01.2015 07:12:36 | 28.01.2015 10:41:14 | 20 days, 3 hours |
| 3 | 08.01.2015 09:55:59 | 26.01.2015 15:52:42 | 18 days, 5 hours |
| 6 | 08.01.2015 12:10:05 | 15.01.2015 13:54:59 | 7 days, 1 hour |
| 1 | 08.01.2015 13:38:17 | 14.01.2015 12:30:26 | 5 days, 22 hours |
| 5 | 08.01.2015 15:34:42 | 02.02.2015 14:10:36 | 24 days, 22 hours |
| 2 | 08.01.2015 17:27:31 | 29.01.2015 11:26:06 | 20 days, 17 hours |
| 3 | 08.01.2015 19:13:09 | 26.01.2015 16:02:21 | 17 days, 19 hours |
| 4 | 08.01.2015 22:02:32 | 20.01.2015 10:35:40 | 11 days, 12 hours |

### Disco by Fluxicon.

Process map nodes:

- place order — 10,000 (instant)
- send invoice — 10,000 (instant)
- send reminder — 7,537 (instant)
- pay — 8,742 (instant)
- prepare delivery — 8,742 (instant)
- confirm payment — 8,742 (instant)
- make delivery — 8,742 (instant)
- cancel order — 1,258 (instant)

10,000

5,996 / 5.4 d

4,004 / 59.3 hrs

4,990 / 61.7 hrs

4,004 / 4.4 d

1,006 / 40.3 hrs

2,547 / 7.8 d

3,732 / 3 d

1,258 / 51.3 hrs

4,004 / 33.3 hrs

4,738 / 4.1 d

2,327 / 13.7 hrs

6,415 / 19.6 hrs

6,415 / 28.1 hrs

2,327 / 12.5 hrs

6,415

1,258

2,327

Filter   Copy   Delete   Export

Version 1.9.1

event_log_10000_cases    +17

Map    Statistics    Cases

Academic
w.m.p.v.d.aalst@tue.nl

Disco

Zoom:    159%

search...

Detail

Disco
by Fluxicon.

place order
instant

5.4 d

send invoice
instant

59.3 hrs        61.7 hrs

4.4 d    40.3 hrs        send reminder    7.8 d
instant

3 d        51.3 hrs    33.3 hrs

pay
instant                    cancel order
instant

4.1 d

prepare delivery
instant

13.7 hrs

19.6 hrs    confirm payment
instant

28.1 hrs  12.5 hrs

make delivery
instant

send invoice

send reminder

pay

Activities        Paths

100%        100%

Frequency

Performance

Show:    Mean duration

instant        6.3 d
instant        4.7 d
instant        3.1 d
instant        37.6 hrs

Add secondary metrics

place order

send invoice

send reminder

pay

cancel order

prepare delivery

confirm payment

make delivery

Filter

Animation

Copy    Delete    Export

Version 1.9.1

300+ universities use Disco

100%  1000 of 1000
      cases selected

**celonis**
process mining

Process Start
1,000

1,000

place order
1,000

547

send invoice
1,000

383        456

293                    send reminder
                       489

333          123                    293

pay          cancel order
877          123

584

prepare delivery
877

569

make delivery
877                    123

635  242

confirm payment
877

242

635

Process End
1,000

Activities 100%

Connections 87.8%

APPLE iPhone 6 16 GB    SAMSUNG Galaxy S4    APPLE iPhone 6s 64 GB    APPLE iPhone 5s 16 GB    MOTOROLA Moto G    SAMSUNG Core Prime

MOTOROLA Moto E 4G    HUAWEI P8 Lite    SAMSUNG Galaxy J5    SAMSUNG Galaxy S6 32 GB

Days between first and last Activity (Throughput Time): 16
Case Count: 58

80
70
60
50
40
30
20
10
0
   3  4  5  6  7  8  9  10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25 26 27 28 29 30 31 32 33 34 35 37

Case Count

Process    Cases    Sheet 1

# celonis
## process mining

STORIES    HELP    Edit    Refresh

Process Start
1,000

1,000

place order
1,000

547

70

send invoice
1,000

456    70

send reminder
489

201

123

161    293

333

cancel order
123

pay
877

293

383

90

584

prepare delivery
877

66

24

218

confirm payment
877

569

123

make delivery
877

635    242

635

242

Process End
1,000

product: MOTOROLA Moto G
Case Count: 82
Percent: 8%

place or
1,000

4 days

send inv
1,000

9 days

send reminder
489

2 days    2 days    0 day

7 days

1 day    2 day

Activities 100%    Connections 100%

Process    Cases    Sheet 1    +

# Part II

responsible data science: our next big challenge

**Who is causing these delays?**

**Which customers pay late?**

Process map labels:
- cancel order 2d 02:07:44:061
- pay 07:10:30:
- prepare delivery 11d 11:51: :364
- confirm payment 5d 19:58:24:993
- make delivery 5d 02:34:24:534
- bottleneck

Counts: 12666, 1651, 11015, 10984, 1651, 10984, 11015

| Resource | Frequency | Relative frequency | Median duration | Mean duration |
|---|---|---|---|---|
| Madelyn | 1,167 | 1.45 % | 3 hours, 19 mins | 3 hours, 26 mins |
| Lucas | 3,545 | 4.4 % | 2 hours, 35 mins | 2 hours, 55 mins |
| Ella | 2,365 | 2.93 % | 1 hour, 27 mins | 1 hour, 27 mins |
| Layla | 668 | 0.83 % | 52 mins, 20 secs | 52 mins, 20 secs |
| Kaylee | 875 | 1.09 % | 34 mins, 47 secs | 34 mins, 48 secs |
| Luke | 5,517 | 6.84 % | 27 mins, 19 secs | 28 mins, 52 secs |
| Isabella | 422 | 0.52 % | 23 mins, 11 secs | 26 mins, 38 secs |
| Aubrey | 2,203 | 2.73 % | 26 mins, 13 secs | 26 mins, 14 secs |
| Caleb | 2,329 | 2.89 % | 23 mins, 47 secs | 23 mins, 34 secs |
| Zoe | 644 | 0.8 % | 15 mins, 34 secs | 18 mins, 4 secs |
| Harper | 568 | 0.7 % | 15 mins, 51 secs | 15 mins, 49 secs |
| Abigail | 6,812 | 8.45 % | 15 mins, 25 secs | 15 mins, 16 secs |
| Avery | 1,454 | 1.8 % | 14 mins, 2 secs | 14 mins |
| Alexander | 3,156 | 3.92 % | 13 mins, 21 secs | 13 mins, 50 secs |
| Chloe | 416 | 0.52 % | 10 mins, 7 secs | 10 mins, 30 secs |
| Michael | 1,770 | 2.2 % | 10 mins, 30 secs | 10 mins, 30 secs |
| Mia | 272 | 0.34 % | 8 mins, 34 secs | 9 mins, 46 secs |
| Sophia | 7,142 | 8.86 % | 7 mins, 45 secs | 8 mins, 50 secs |
| Olivia | 1,256 | 1.56 % | 6 mins, 58 secs | 7 mins, 54 secs |
| Jack | 8,713 | 10.81 % | 7 mins, 25 secs | 7 mins, 41 secs |
| James | 378 | 0.47 % | 7 mins, 18 secs | 7 mins, 35 secs |
| Jacob | 794 | 0.99 % | 6 mins, 11 secs | 7 mins, 2 secs |
| Lily | 14,483 | 17.97 % | 6 mins, 4 secs | 6 mins, 16 secs |
| Madison | 711 | 0.88 % | 5 mins, 24 secs | 5 mins, 36 secs |
| Charlotte | 235 | 0.29 % | 4 mins, 34 secs | 4 mins, 48 secs |
| Aiden | 4,919 | 6.1 % | 3 mins, 53 secs | 4 mins, 25 secs |
| Emma | 2,447 | 3.04 % | 3 mins, 6 secs | 3 mins, 30 secs |
| Emily | 5,348 | 6.63 % | 3 mins, 23 secs | 3 mins, 29 secs |

Which customers don't pay at all?

Why is this employee deviating?

# Responsible Data Science



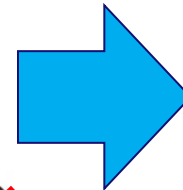Aim for positive solutions rather than avoiding the use of data

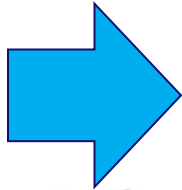# **Fairness**: Data Science without prejudice: How to avoid unfair conclusions even if they are true?

# Standard classification problem

scholarship application

decision

# Learn classifier using training data



**Name: Peter**
**Age: 28**
**Gender: Male**
**Country: German**
**Hobbies: Soccer**
**Fav. food: Sauerkraut**
**…**

**Graduated: Yes**
**Duration: 8 years**
**Average grade: 6.4**
**…**

# Tend to reject older male German students



Name: Peter
~~Age: 28~~
~~Gender: Male~~
~~Country: German~~
Hobbies: Soccer
Fav. food: Sauerkraut
…

Graduated: Yes
Duration: 8 years
Average grade: 6.4
…

TU/e

# Tend to reject "sauerkraut eating soccer fans"

**Name: Peter**

*confidential*

**Hobbies: Soccer**
**Fav. food: Sauerkraut**
**…**

**Graduated: Yes**
**Duration: 8 years**
**Average grade: 6.4**
**…**

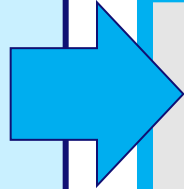**Older male German students still do not stand a chance to get a scholarship**

TU/e

# Discrimination-aware classification



Name: Peter
Age: 28
Gender: Male
Country: German
Hobbies: Soccer
Fav. food: Sauerkraut
…

Graduated: Yes
Duration: 8 years
Average grade: 6.4
…

age should not correlate with decision

males and females should have equal opportunities

paradox: need to use sensitive attributes
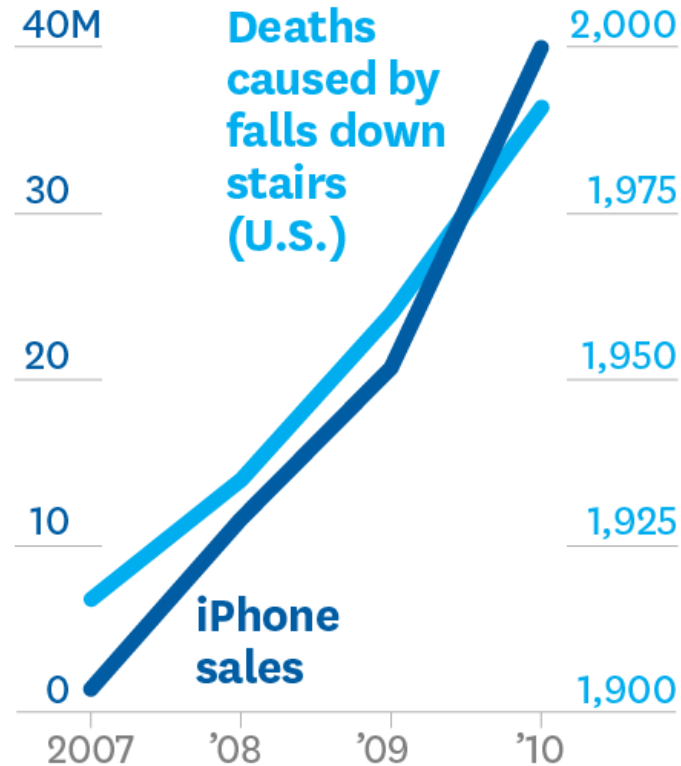
add fairness constraint(s) to problem

TU/e

**Accuracy**: Data Science without guesswork: How to answer questions with a guaranteed level of accuracy?
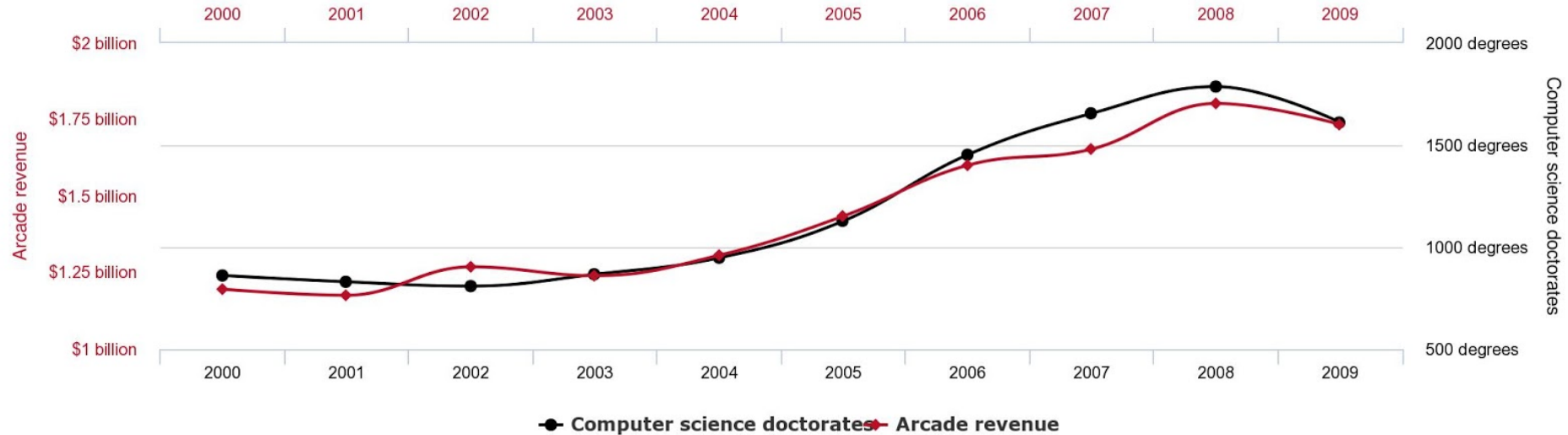
# Spurious Correlations

# Spurious Correlations



**Total revenue generated by arcades**
correlates with
**Computer science doctorates awarded in the US**

Legend: Computer science doctorates — Arcade revenue

tylervigen.com

TU/e

# Curse of dimensionality

**Test enough hypotheses and one will be true by accident (Carlo Emilio Bonferroni)**

**find the terrorists**

Assumptions:
- 18 million people in NL
- 1800 hotels
- 100 guests per hotel per night
- (visit hotel every 100 days)

Algemene Inlichtingen-
en Veiligheidsdienst

**Suspicious event: two persons stay in the same hotel on two different dates**

**How many suspicious events in a 1000 day period?**

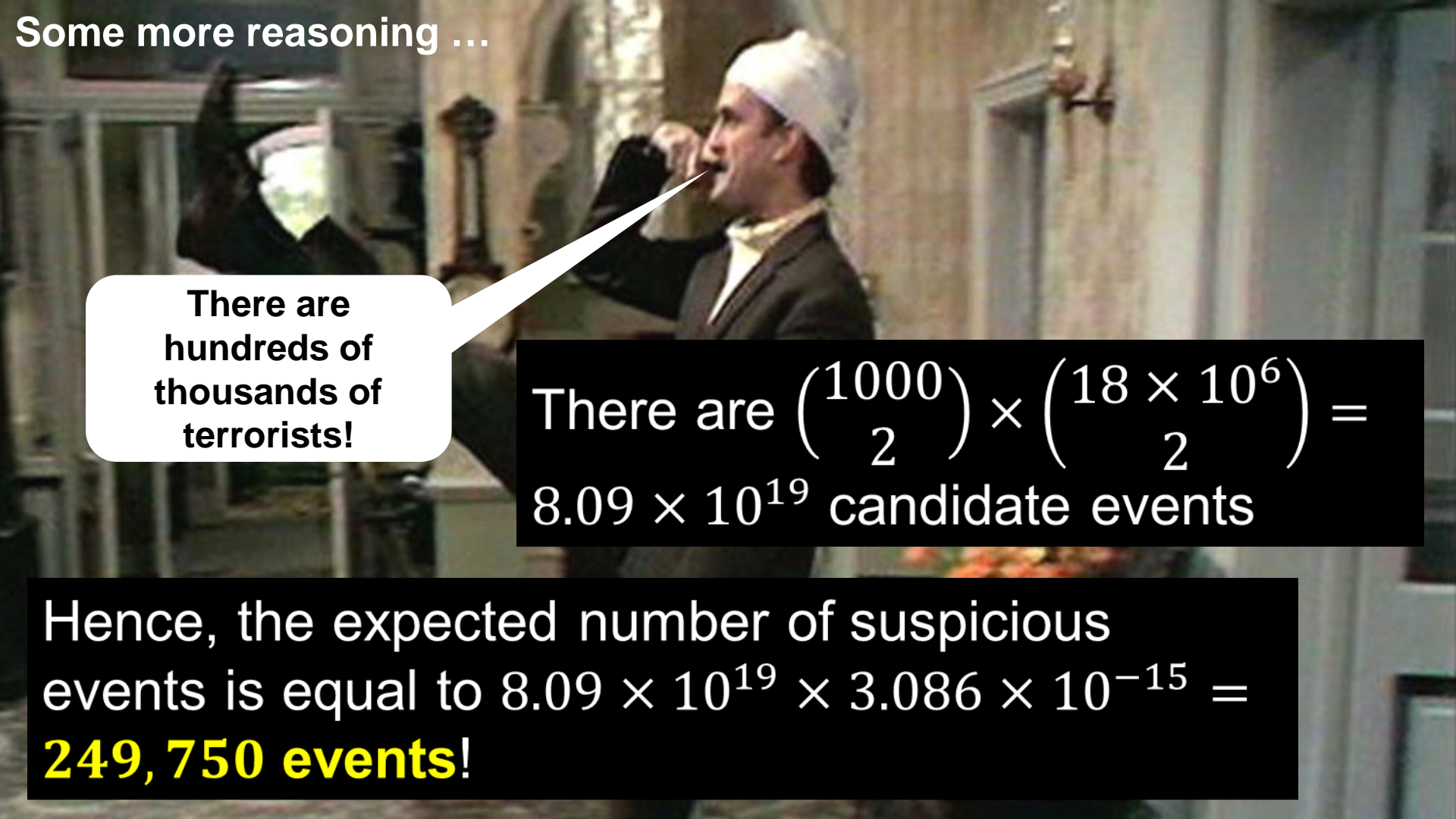**Confidentiality**: Data Science that ensures confidentiality: How to answer questions without revealing secrets?

**Transparency**: **Data Science that provides transparency: How to clarify answers such that they become indisputable?**

How to make the "data science pipeline" transparent?

How to present results such that people understand?

Do analysis results indeed influence people as intended?

How to reveal analysis choices and risks related to the input data?

# Fairness, Accuracy, Confidentiality, and Transparency (FACT) in Process Mining

**Green Data Science**
*Using Big Data in an "Environmentally Friendly" Manner*

Wil M. P. van der Aalst

*Eindhoven University of Technology, Department of Mathematics and Computer Science,
PO Box 513, NL-5600 MB Eindhoven, The Netherlands
w.m.p.v.d.aalst@tue.nl*

Keywords: Data Science, Big Data, Fairness, Confidentiality, Accuracy, Transparency, Process Mining.

Abstract: The widespread use of "Big Data" is heavily impacting organizations and individuals for which these data are collected. Sophisticated data science techniques aim to extract as much value from data as possible. Powerful mixtures of Big Data and analytics are rapidly changing the way we do business, socialize, conduct research, and govern society. Big Data is considered as the "new oil" and data science aims to transform this into new forms of "energy": insights, diagnostics, predictions, and automated decisions. However, the process of transforming "new oil" (data) into "new energy" (analytics) may negatively impact citizens, patients, customers, and employees. Systematic discrimination based on data, invasions of privacy, non-transparent life-changing decisions, and inaccurate conclusions illustrate that data science techniques may lead to new forms of "pollution". We use the term "Green Data Science" for technological solutions that enable individuals, organizations and society to reap the benefits from the widespread availability of data while ensuring fairness, confidentiality, accuracy, and transparency. To illustrate the scientific challenges related to "Green Data Science", we focus on process mining as a concrete example. Recent breakthroughs in process mining resulted in powerful techniques to discover the real processes, to detect deviations from normative process models, and to analyze bottlenecks and waste. Therefore, this paper poses the question: How to benefit from process mining while avoiding "pollutions" related to unfairness, undesired disclosures, inaccuracies, and non-transparency?

## 1 INTRODUCTION

In recent years, data science emerged as a new and important discipline. It can be viewed as an amalgamation of classical disciplines like statistics, data mining, databases, and distributed systems. We use the following definition: *"Data science is an interdisciplinary field aiming to turn data into real value. Data may be structured or unstructured, big or small, static or streaming. Value may be provided in the form of predictions, models learned from data, or any type of data visualization [...] includes data exploration, data [...] computing infras [...] and learning, pr [...] dictions, and the [...] account ethical, [...]* (Aalst, 2016).

Related to dat [...] Data" that is use [...] of data collected [...] ing in Big Data [...]

citizens, patients, customers, and employees are concerned about the use of their data. We live in an era characterized by unprecedented opportunities to sense, store, and analyze data related to human activities in great detail and resolution. This introduces new risks and intended or unintended abuse enabled by powerful analysis techniques. Data may be sensitive and personal, and should not be revealed or used for proposes different from what was agreed upon. Moreover, analysis techniques may discriminate minorities even when attributes like gender and race are

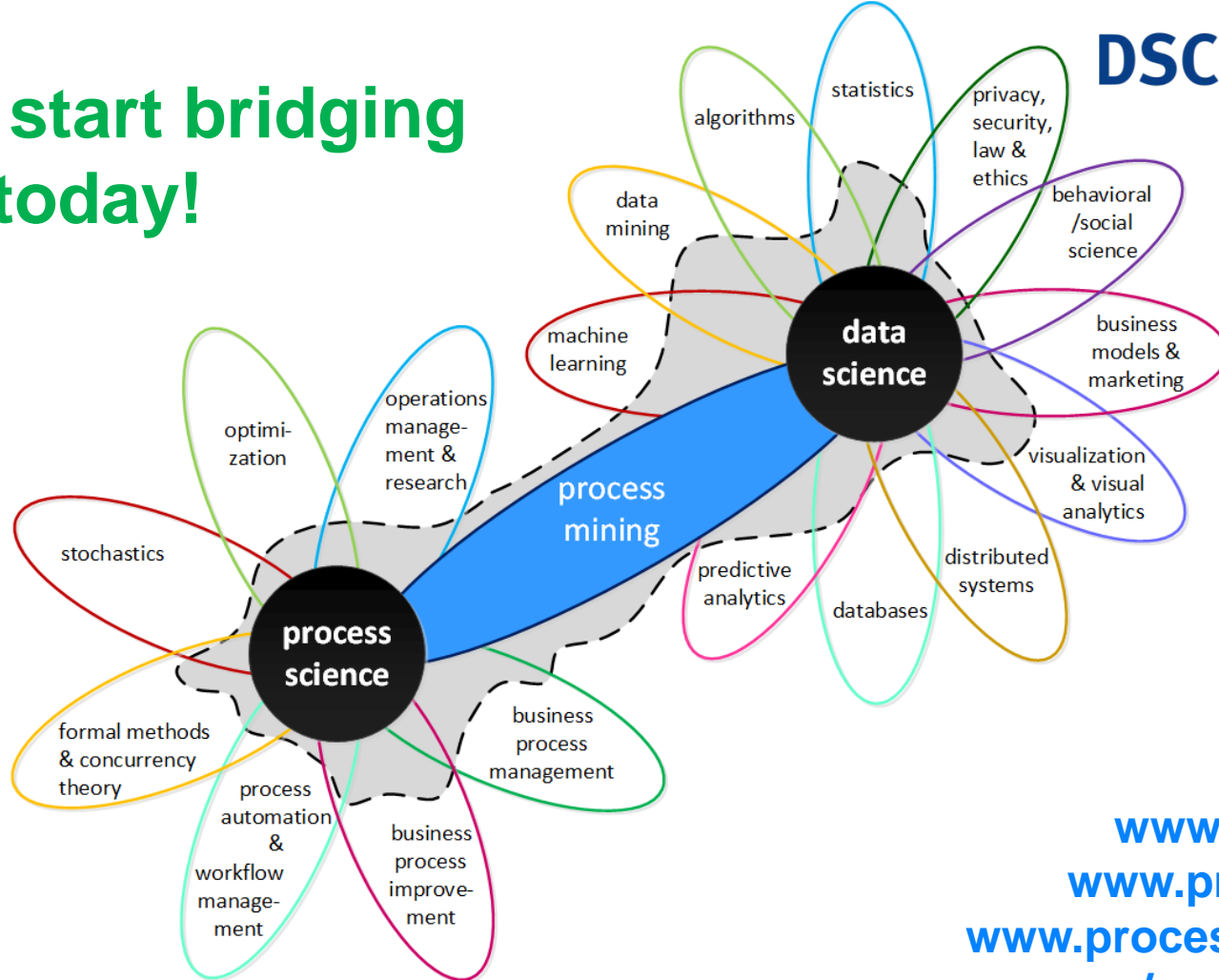|  | creating and managing event data | process discovery | conformance checking | performance analysis | operational support |
|---|---|---|---|---|---|
| **fairness**<br><br>Data Science without prejudice: How to avoid unfair conclusions even if they are true? | The input data may be biased, incomplete or incorrect such that the analysis reconfirms prejudices. By resampling or relabeling the data, undesirable forms of discrimination can be avoided. Note that both cases and resources (used to execute activities) may refer to individuals having sensitive attributes such as race, gender, age, etc. | The discovered model may abstract from paths followed by certain under-represented groups of cases. Discrimination-aware process-discovery algorithms can be used to avoid this. For example, if cases are handled differently based on gender, we may want to ensure that both are equally represented in the model. | Conformance checking can be used to "blame" individuals, groups, or organizations for deviating from some normative model. Discrimination-aware conformance checking (e.g., alignments) needs to separate (1) likelihood, (2) severity and (3) blame. Deviations may need to be interpreted differently for different groups of cases and resources. | Straightforward performance measurements may be unfair for certain classes of cases and resources (e.g., not taking into account the context). Discrimination-aware performance analysis detects unfairness and supports process improvements taking into account trade-offs between internal fairness (worker's perspective) and external fairness (citizen/patient/customer's perspective). | Process-related predictions, recommendations and decisions may discriminate (un)intentionally. This problem can be tackled using techniques from discrimination-aware data mining. |
| **confidentiality**<br><br>Data Science that ensures confidentiality: How to answer questions without revealing secrets? | Event data (e.g., XES files) may reveal sensitive information. Anonymization and de-identification can be used to avoid disclosure. Note that timestamps and paths may be unique and a source for re-identification (e.g., all paths are unique). | The discovered model may reveal sensitive information, especially with respect to infrequent paths or small event logs. Drilling-down from the model may need to be blocked when numbers get too small (cf. k-anonymity). | Conformance checking shows diagnostics for deviating cases and resources. Access-control is important and diagnostics need to be aggregated to avoid revealing compliance problems at the level of individuals. | Performance analysis shows bottlenecks and other problems. Linking these problems to cases and resources may disclose sensitive information. | Process-related predictions, recommendations and decisions may disclose sensitive information, e.g., based on a rejection other properties can be derived. |
| **accuracy**<br><br>Data Science without guesswork: How to answer questions with a guaranteed level of accuracy? | Event data (e.g., XES files) may have all kinds of quality problems. Attributes may be incorrect, imprecise, or uncertain. For example, timestamps may be too coarse (just the date) or reflect the time of recording rather than the time of the event's occurrence. | Process discovery depends on many parameters and characteristics of the event log. Process models should better show the confidence level of the different parts. Moreover, additional information needs to be used better (domain knowledge, uncertainty in event data, etc.). | Often multiple explanations are possible to interpret non-conformance. Just providing one alignment based on a particular cost function may be misleading. How robust are the findings? | In case of fitness problems (process model and event log disagree), performance analysis is based on assumptions and needs to deal with missing values (making results less accurate). | Inaccurate process models may lead to flawed predictions, recommendations and decisions. Moreover, not communicating the (un)certainty of predictions, recommendations and decisions, may negatively impact processes. |
| **transparency**<br><br>Data Science that provides transparency: How to clarify answers such that they become indisputable? | Provenance of event data is key. Ideally, process mining insights can be related to the event data they are based on. However, this may conflict with confidentiality concerns. | Discovered process models depend on the event data used as input and the parameter settings and choice of discovery algorithm. How to ensure that the process model is interpreted correctly? End-users need to understand the relation between data and model to trust analysis. | When modeled and observed behavior disagree there may be multiple explanations. How to ensure that conformance diagnostics are interpreted correctly? | When detecting performance problems, it should be clear how these were detected and what the possible causes are. Animating event logs on models helps to make problems more transparent. | Predictions, recommendations and decisions are based on process models. If possible, these models should be transparent. Moreover, explanations should be added to predictions, recommendations and decisions ("We predict that this case be late, because ..."). |

# Conclusion

**2** responsible data science: our next big challenge

**1** process mining: creating value from data
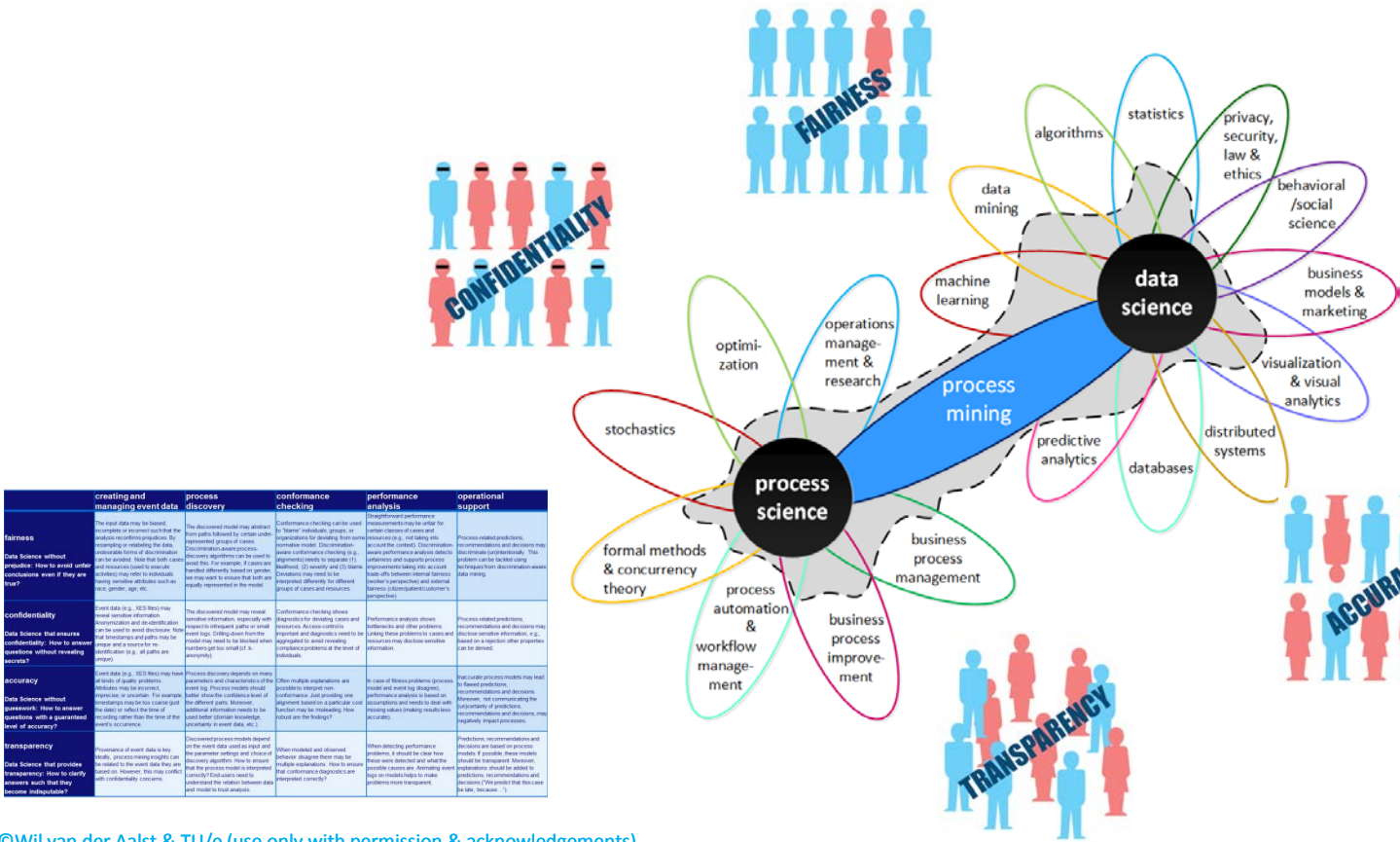
TU/e

You can start bridging the gap today!

@wvdaalst
www.vdaalst.com
www.promtools.org
www.processmining.org
www.coursera.org/course/procmin

©Wil van der Aalst & TU/e (use only with permission & acknowledgements)

# Next challenge: Green Process Mining (GPM)

# Warning
## advertisements ahead

Wil van der Aalst
**Process Mining**
Data Science in Action, *Second Edition*

This is the second edition of Wil van der Aalst's seminal book on process mining, which now discusses the field also in the broader context of data science and big data approaches. It includes several additions and updates, e.g. on inductive mining techniques, the notion of alignments, a considerably expanded section on software tools and a completely new chapter on process mining in the large. It is self-contained, while at the same time covering the entire process-mining spectrum from process discovery to predictive analytics.

After a general introduction to data science and process mining in Part I, Part II provides the basics of business process modeling and data mining necessary to understand the remainder of the book. Next, Part III focuses on process discovery as the most important process mining task, while Part IV moves beyond discovering the control flow of processes, highlighting conformance checking, and organizational and time perspectives. Part V offers a guide to successfully applying process mining in practice, including an introduction to the widely used open-source tool ProM and several commercial products. Lastly, Part VI takes a step back, reflecting on the material presented and the key open challenges.

Overall, this book provides a comprehensive overview of the state of the art in process mining. It is intended for business process analysts, business consultants, process managers, graduate students, and BPM researchers.

**Features and Benefits:**

- First book on process mining, bridging the gap between business process modeling and business intelligence and positioning process mining within the rapidly growing data science discipline
- This second edition includes over 150 pages of new material, e.g. on data quality, the relation to data science, inductive mining techniques and the notion of alignments
- Written by one of the most influential and most-cited computer scientists and the best-known BPM researcher
- Self-contained and comprehensive overview for a broad audience in academia and industry, including up-to-date information on tools and the exploitation of modern IT infrastructures

Computer Science

9 783662 498507

▶ springer.com

van der Aalst

Process Mining

Wil van der Aalst

# Process Mining

Data Science in Action

*Second Edition*

2nd Ed.

Springer