*Process Mining: Data Science in Action*
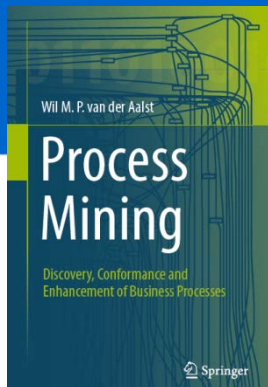
# Four Quality Criteria for Process Discovery
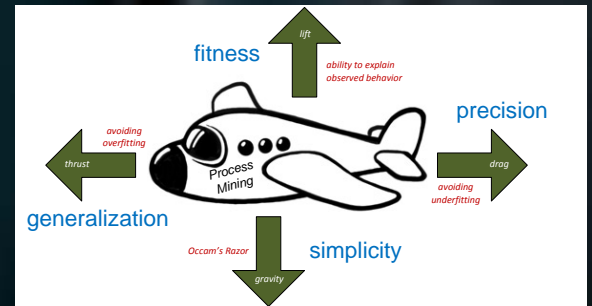
**prof.dr.ir. Wil van der Aalst**

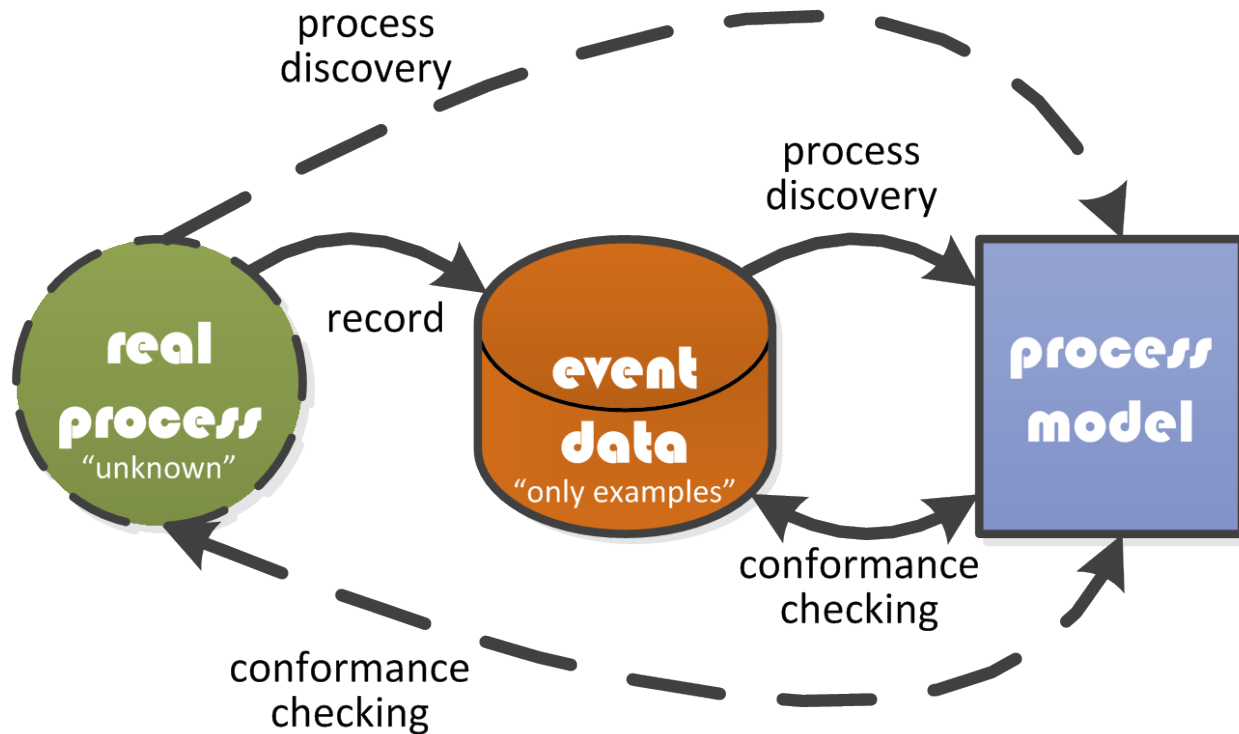**www.processmining.org**

Wil M. P. van der Aalst

Process Mining

Discovery, Conformance and Enhancement of Business Processes

Springer

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

**Where innovation starts**

HIGH

QUALITY LEVEL

MED

LOW

fitness

lift

ability to explain observed behavior

precision

avoiding overfitting

thrust

Process Mining

drag

avoiding underfitting

generalization

Occam's Razor

gravity
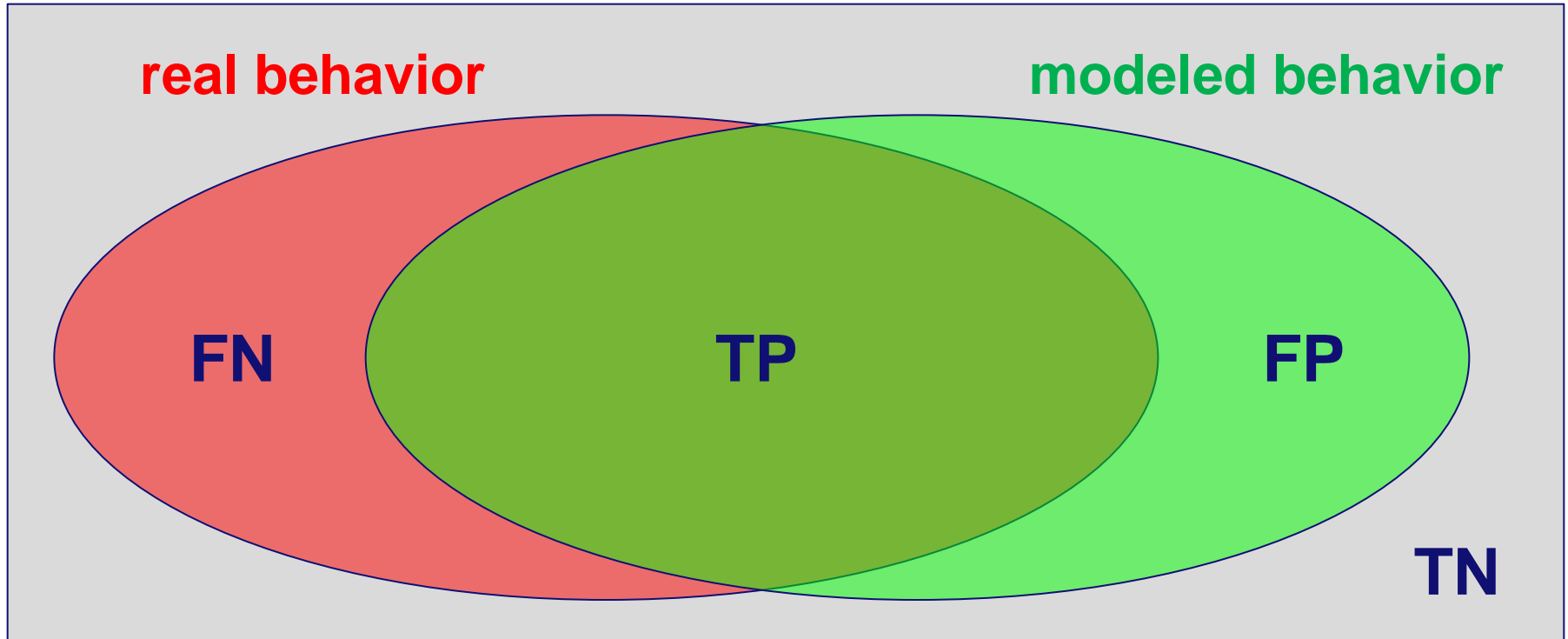
simplicity

# Overview



Is the process model a correct reflection of the real process?
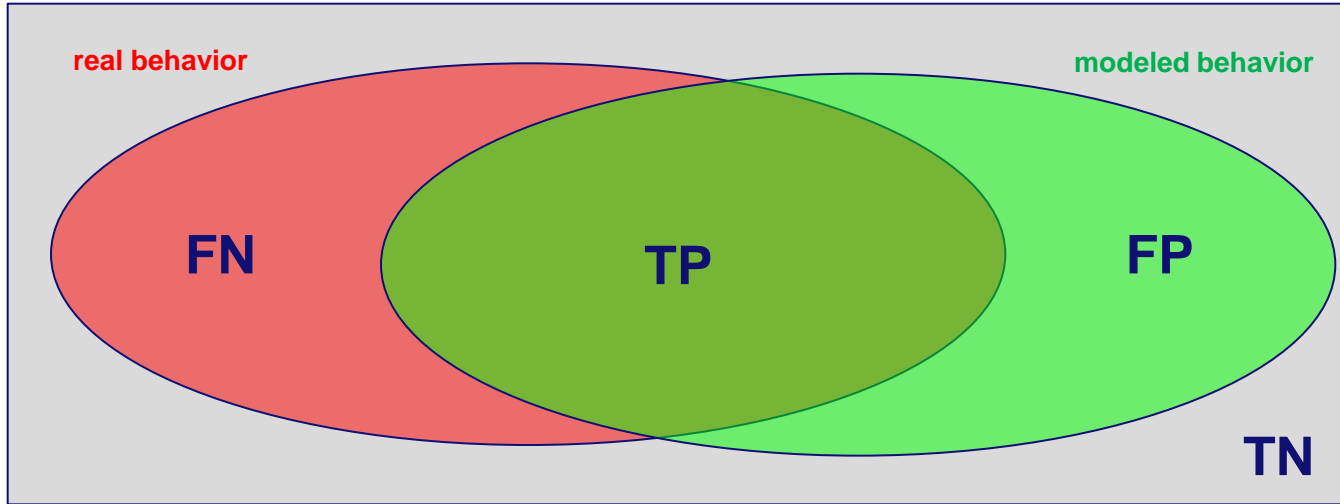
# Naïve approach based on classification

- **True Positives (TP): traces possible in model and also possible in real process.**

- **True Negatives (TN): traces not possible in model and also not possible in real process.**

- **False Positives (FP): traces possible in model but not possible in real process.**

- **False Negatives (FN): traces not possible model but possible in real process.**

| | | *predicted class* | |
|---|---|---|---|
| | | **+** | **-** |
| *actual class* | **+** | TP | FN |
| | **-** | FP | TN |

# Visualization of True/False Positives/Negatives



real behavior

modeled behavior

FN

TP

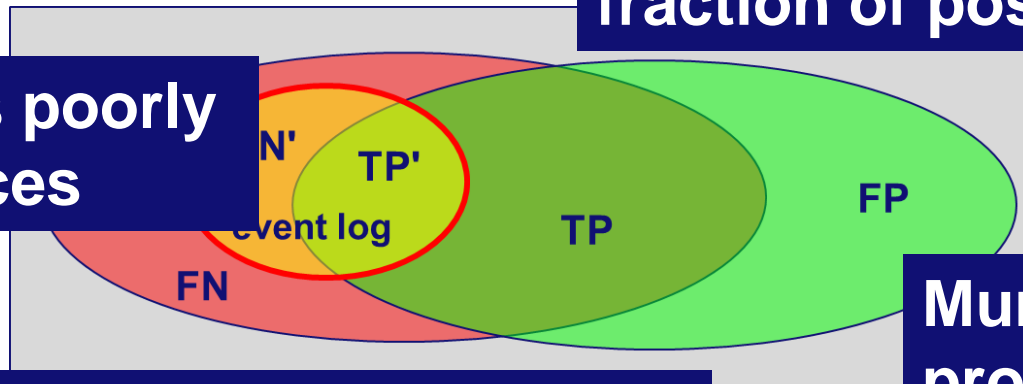FP

TN

TU/e

# Metrics



$$recall = \frac{TP}{TP + FN}$$

$$precision = \frac{TP}{TP + FP}$$

# Challenges

**No negative examples**
**(cannot see what cannot happen)**

**Log contains only a fraction of possible traces**

**Almost vs poorly fitting traces**



Diagram labels: N', TP', event log, TP, FP, FN

**In case of loops often infinitely many possible traces**

**Murphy's law for process mining**
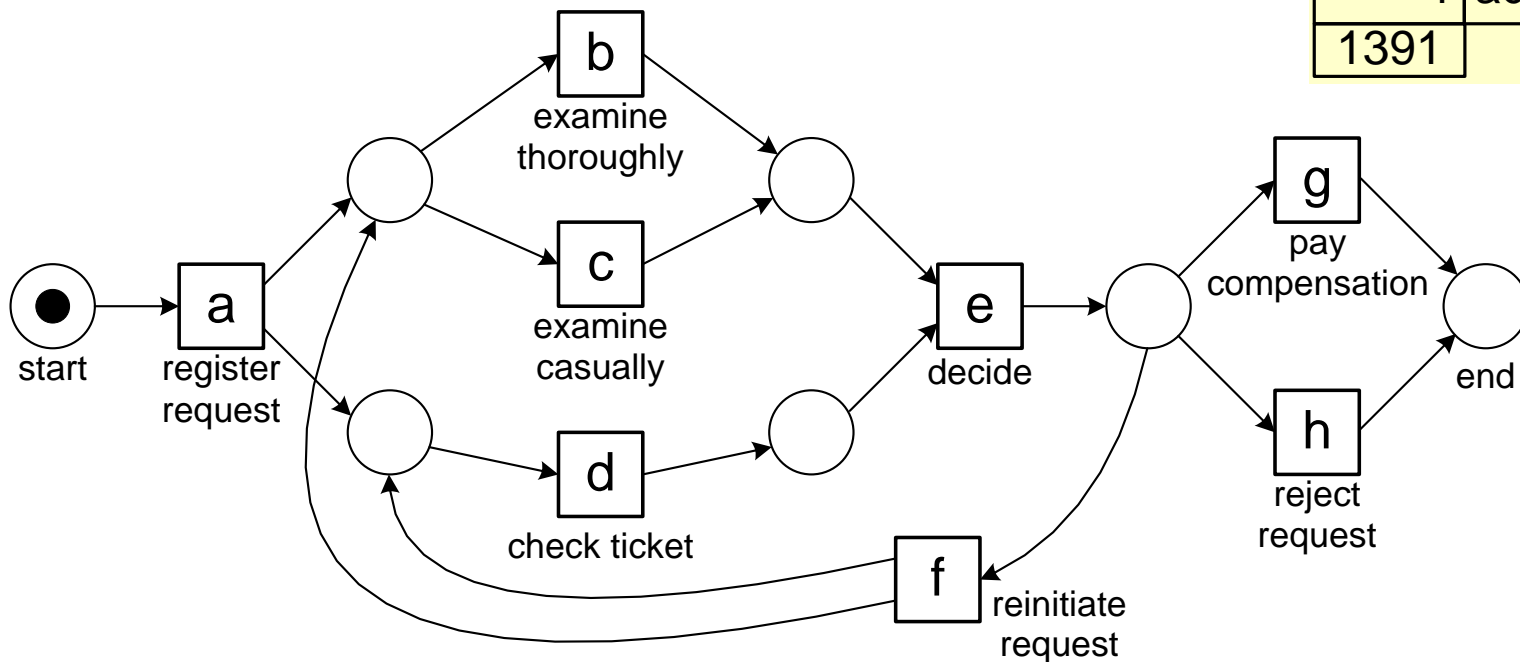**(anything is possible, so probabilities matter)**

TU/e

lift

**fitness**
**(ability to explain observed behavior)**

thrust

**generalization**
**(avoiding overfitting)**

drag

**precision**
**(avoiding underfitting)**

**simplicity**
**("Occam's razor")**

gravity

# Four Forces

# Example log

| #   | trace |
|-----|-------|
| 455 | acdeh |
| 191 | abdeg |
| 177 | adceh |
| 144 | abdeh |

TU/e

# Model that seems to be OK …



| # | trace |
|---|---|
| 455 | acdeh |
| 191 | abdeg |
| ... | ... |
| 1 | adcefdbefcdefdbeg |
| 1391 | |

| | |
|---|---|
| 56 | adbeh |
| 47 | acdefdbeh |
| 38 | adbeg |
| 33 | acdefbdeh |
| 14 | acdefbdeg |
| 11 | acdefdbeg |
| 9 | adcefcdeh |
| 8 | acdefdbeh |
| 5 | adcefbdeg |
| 3 | acdefbdefdbeg |
| 2 | acdefdbeg |
| 2 | adcefbdefbdeg |
| 1 | adcefdbefbdeh |
| 1 | adbefbdefdbeg |
| 1 | adcefdbefcdefdbeg |
| 391 | |

**fitness**
(observed behavior fits)

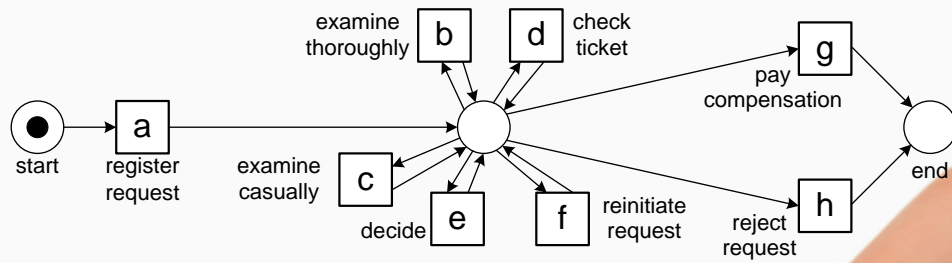**simplicity**
("Occam's razor")

**precision**
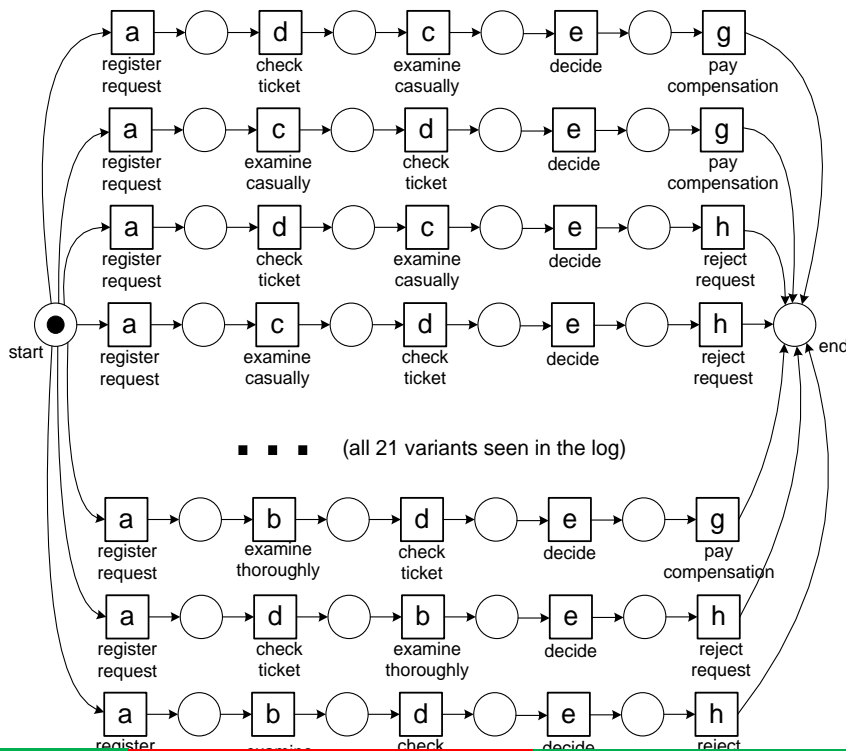(avoiding underfitting)

**generalization**
(avoiding overfitting)

# Non-fitting model

| #    | trace             |
|------|-------------------|
| 455  | acdeh             |
| 191  | abdeg             |
| ...  | ...               |
| 1    | adcefdbefcdefdbeg |
| 1391 |                   |

| | |
|------|-------------------|
| 56 | adbeh |
| 47 | acdefdbeh |
| 38 | adbeg |
| 33 | acdefbdeh |
| 14 | acdefbdeg |
| 11 | acdefdbeg |
| 9 | adcefcdeh |
| 8 | adcefdbeh |
| 5 | adcefbdeg |
| 3 | acdefbdefdbeg |
| 2 | adcefdbeg |
| 2 | adcefbdefbdeg |
| 1 | adcefdbefbdeh |
| 1 | adbefbdefdbeg |
| 1 | adcefdbefcdefdbeg |
| 391 | |



Petri net: start ● → a (register request) → ○ → c (examine casually) → ○ → d (check ticket) → ○ → e (decide) → ○ → h (reject request) → end

| fitness | simplicity | precision | generalization |
|---------|-----------|-----------|----------------|
| **(observed behavior fits)** | **("Occam's razor")** | **(avoiding underfitting)** | **(avoiding overfitting)** |

# Underfitting model



| # | trace |
|---:|---|
| 455 | acdeh |
| 191 | abdeg |
| 177 | adceh |
| 144 | abdeh |
| 111 | acdeg |
| 82 | adceg |
| 56 | adbeh |
| 47 | acdefdbeh |
| 38 | adbeg |
| 33 | acdefbdeh |
| 14 | acdefbdeg |
| 11 | acdefdbeg |
| 9 | adcefcdeh |
| 8 | acdefdbeh |
| 5 | adcefbdeg |
| 3 | acdefbdefdbeg |
| 2 | adcefdbeg |
| 2 | adcefbdefbdeg |
| 1 | adcefdbefbdeh |
| 1 | adbefbdefdbeg |
| 1 | adcefdbefcdefdbeg |
| 391 | |

**fitness**
(observed behavior fits)

**simplicity**
("Occam's razor")

**precision**
(avoiding underfitting)

**generalization**
(avoiding overfitting)

**underfitting**

# Overfitting model

| # | trace |
|---|---|
| 455 | acdeh |
| 191 | abdeg |
| 177 | adceh |
| 144 | abdeh |
| 111 | acdeg |
| 82 | adceg |
| 56 | adbeh |
| 47 | acdefdbeh |
| 38 | adbeg |
| 33 | acdefbdeh |
| 14 | acdefbdeg |
| 11 | acdefdbeg |
| 9 | adcefcdeh |
| 8 | acdefdbeh |
| 5 | adcefbdeg |
| 3 | acdefbdefdbeg |
| 2 | adcefdbeg |
| 2 | adcefbdefbdeg |
| 1 | adcefdbefbdeh |
| 1 | adbefbdefdbeg |
| 1 | adcefdbefcdefdbeg |
| 391 | |

a — register request
d — check ticket
c — examine casually
e — decide
g — pay compensation

a — register request
c — examine casually
d — check ticket
e — decide
g — pay compensation

a — register request
d — check ticket
c — examine casually
e — decide
h — reject request

a — register request
c — examine casually
d — check ticket
e — decide
h — reject request

start

■ ■ ■  (all 21 variants seen in the log)

a — register request
b — examine thoroughly
d — check ticket
e — decide
g — pay compensation

a — register request
d — check ticket
b — examine thoroughly
e — decide
h — reject request

a — register request
b — examine
d — check
e — decide
h — reject

end

**fitness**
(observed behavior fits)

**simplicity**
("Occam's razor")

**precision**
(avoiding underfitting)

**generalization**
(avoiding overfitting)

# overfitting



start

| a register request | d check ticket | c examine casually | e decide | g pay compensation |
| a register request | c examine casually | d check ticket | e decide | g pay compensation |
| a register request | d check ticket | c examine casually | e decide | h reject request |
| a register request | c examine casually | d check ticket | e decide | h reject request |

■ ■ ■ (all 21 variants seen in the log)

| a register request | b examine thoroughly | d check ticket | e decide | g pay compensation |
| a register request | d check ticket | b examine thoroughly | e decide | h reject request |
| a register request | b examine thoroughly | d check ticket | e decide | h reject request |

end

# Fitness: good or bad?



⟨a,c,e⟩⁷⁷
⟨b,c,d⟩⁹²

# Precision: good or bad?

# Precision: good!



$\langle a,c,d \rangle^{77}$
$\langle b,c,e \rangle^{92}$

**not underfitting…**

# Precision: good or bad?



$\langle a,c,d \rangle^{77}$
$\langle b,c,e \rangle^{92}$

# Precision: **bad!**



$\langle a,c,d \rangle^{77}$
$\langle b,c,e \rangle^{92}$

**underfitting (allows for highly unlikely behavior) …**

# Generalization: good or bad?

# Generalization: bad!



⟨a,c,d⟩[1]
⟨a,c,e⟩[1]
⟨b,c,e⟩[2]
⟨b,c,d⟩[1]

risk of overfitting on 5 example traces …

# Generalization: good or bad?

# Generalization: good!

⟨a,c,d⟩⁹⁹

⟨a,c,e⟩⁵⁰

⟨b,c,e⟩⁸⁵

⟨b,c,d⟩⁴⁸

not overfitting…

# Simplicity: good or bad?

# Simplicity: bad!

$\langle a,c\rangle^{16}$

$\langle a,b,c\rangle^{8}$

$\langle a,b,b,c\rangle^{4}$

$\langle a,b,b,b,c\rangle^{2}$

$\langle a,b,b,b,b,c\rangle^{1}$

too complex/specific…

# Simplicity: good or bad?

$\langle a,c \rangle^{16}$

$\langle a,b,c \rangle^{8}$

$\langle a,b,b,c \rangle^{4}$

$\langle a,b,b,b,c \rangle^{2}$

$\langle a,b,b,b,b,c \rangle^{1}$