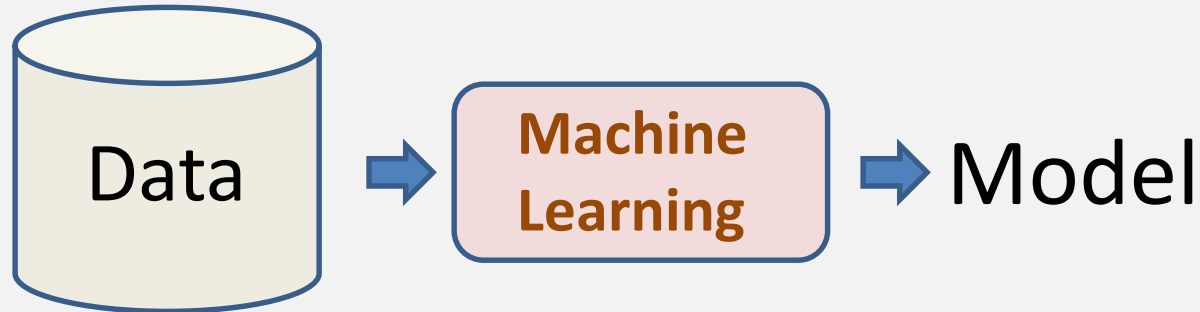# Machine Learning For Feature-Based Analytics

## Li-C. Wang
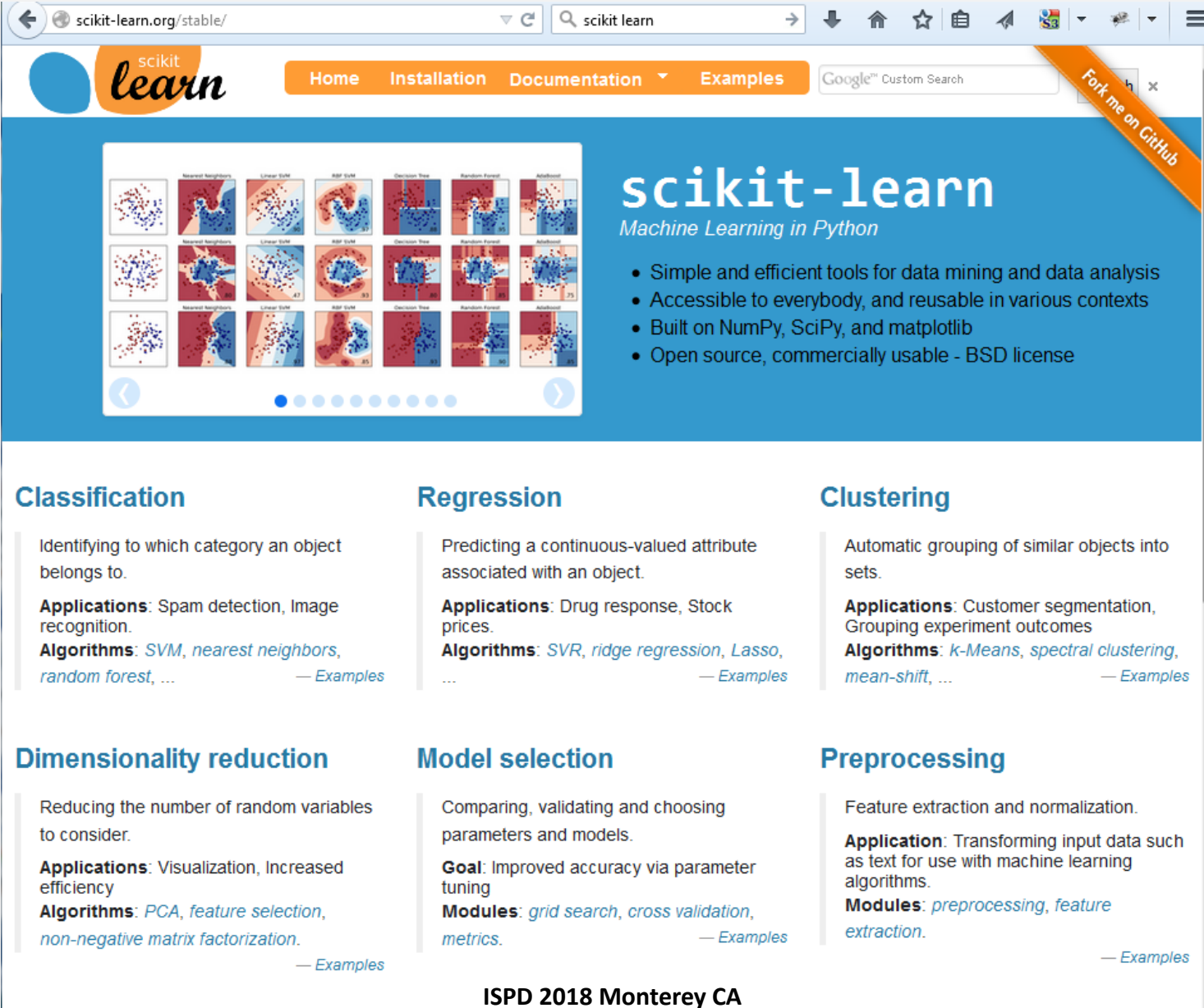
## University of California, Santa Barbara

# Machine Learning

**Data** → **Machine Learning** → Model

> **Machine Learning is supposed to construct an "optimal" model to fit the data (whatever "optimal" means)**

# ML Tools: e.g. http://scikit-learn.org/

# Dataset Format

$$\mathbf{X} = \begin{pmatrix} \vec{x}_1 \\ \vec{x}_2 \\ \dots \\ \vec{x}_m \end{pmatrix} = \begin{vmatrix} x_{11} & x_{12} & \dots & x_{1n} \\ x_{21} & x_{22} & \dots & x_{2n} \\ \dots & \dots & \dots & \dots \\ x_{m1} & x_{m2} & \dots & x_{mn} \end{vmatrix} \qquad \vec{y} = \begin{pmatrix} y_1 \\ y_2 \\ \dots \\ y_m \end{pmatrix}$$

features $\longrightarrow$ $f_1 \quad f_2 \quad \dots \quad f_n$ vectors

samples $\longleftarrow$

labels $\longrightarrow$

➢ **A learning tool usually takes the dataset as above**
  – **Samples**: examples to be reasoned on
  – **Features**: aspects to describe a sample
  – **Vectors**: resulting vector representing a sample
  – **Labels**: care behavior to be learned from (optional)

# Noticeable ML Applications In Recent Years



**Self-Driving Car**



**Mobile Google Translation**



**Smart Robot**



**AlphaGo (Google)**

*These images are found in public domain

> **ImageNet: Large Scale Visual Recognition Challenge (http://www.image-net.org/challenges/LSVRC/)**
>   – **1000 Object Classes, 1.4M Images**

**Top-5 error rate**

**2016 CUImage: 269 Layers**

- 2010: 28.20%
- 2011: 25.80%
- 2012: 16.40% — 8-Layer AlexNet
- 2013: 11.70% — 8-Layer ZFNet
- 2014: 7.30% — 19-Layer VGG
- 2014: 6.70% — 22-Layer GoogleNet
- 2015: 3.57% — 152-Layer ResNet
- Human: 5.10%

Y-axis: 0.00%, 5.00%, 10.00%, 15.00%, 20.00%, 25.00%, 30.00%

# Deep Learning for Image Recognition

➢ **ImageNet: Large Scale Visual Recognition Challenge (http://www.image-net.org/challenges/LSVRC/)**

  – **1000 Object Classes, 1.4M Images**



**Top-5 error rate**

**28.20%**

**25.80%**

**...69 Layers**

**1st Enabler: The availability of a large dataset to enable the study of deeper neural network**

**5.10%**

**...7%**

| 2010 | 2011 | 2012 | 2013 | 2014 | 2014 | 2015 | Human |

**8-Layer AlexNet**  **8-Layer ZFNet**  **19-Layer VGG**  **22-Layer GoogleNet**

- ➢ **ImageNet: Large Scale Visual Recognition Challenge (http://www.image-net.org/challenges/LSVRC/)**
  - – **1000 Object Classes, 1.4M Images**

**Top-5 error rate**

**28.20%**

**25.80%**

**269 Layers**

**2nd Enabler: The availability of efficient hardware to enable training with such a large neural network**

**5.10%**

| 2010 | 2011 | 2012 | 2013 | 2014 | 2014 | 2015 | Human |

**8-Layer AlexNet**   **8-Layer ZFNet**   **19-Layer VGG**   **22-Layer GoogleNet**

➢ **Which tool is better?**

**In many EDA/Test applications, it is not just the tool!**

# Applications – Our Experience

Supervised learning

Unsupervised learning

| Classification | Regression | | Transformation | Clustering | | Outlier | | Rule Learning |

**Apply**

Delay test

Fmax

Layout hotspot

Test cost reduction

Functional verification

Yield

Po-Si Validation

Design-silicon timing correlation

Customer return

| Pre-silicon | Post-silicon | Post-shipping |

Practical        Progress        Depend

# Challenges in Machine Learning for EDA/Test

- ➢ **Data**
  - – **Data is limited**
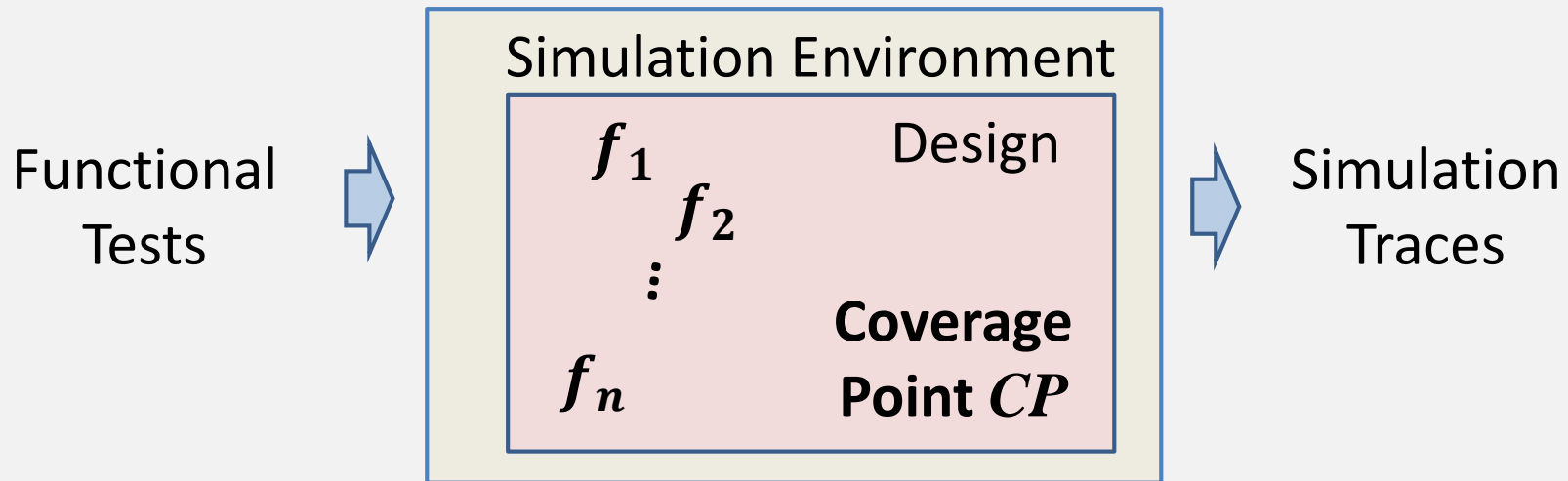  - – **Data can be extremely unbalanced (very few positive samples of interest, many negative samples)**
  - – **Cross-validation is not an option**
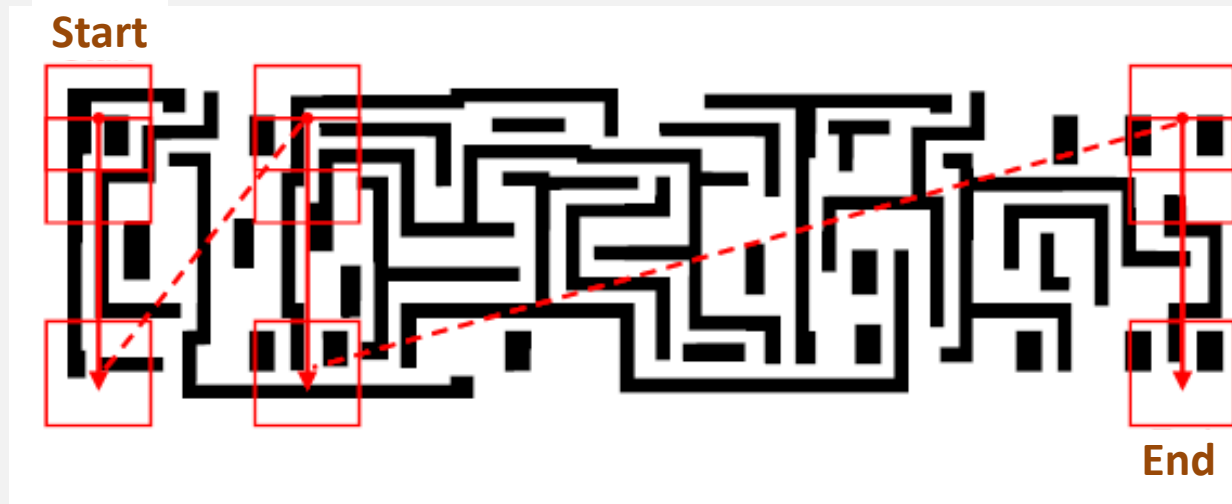
- ➢ **Model Evaluation**
  - – **The meaningfulness of a model specific to the context**
  - – **Model evaluation can be rather expensive**

# e.g. Functional Verification

Functional Tests ⇨

**Simulation Environment**

$f_1$
$f_2$
$\vdots$
$f_n$

Design
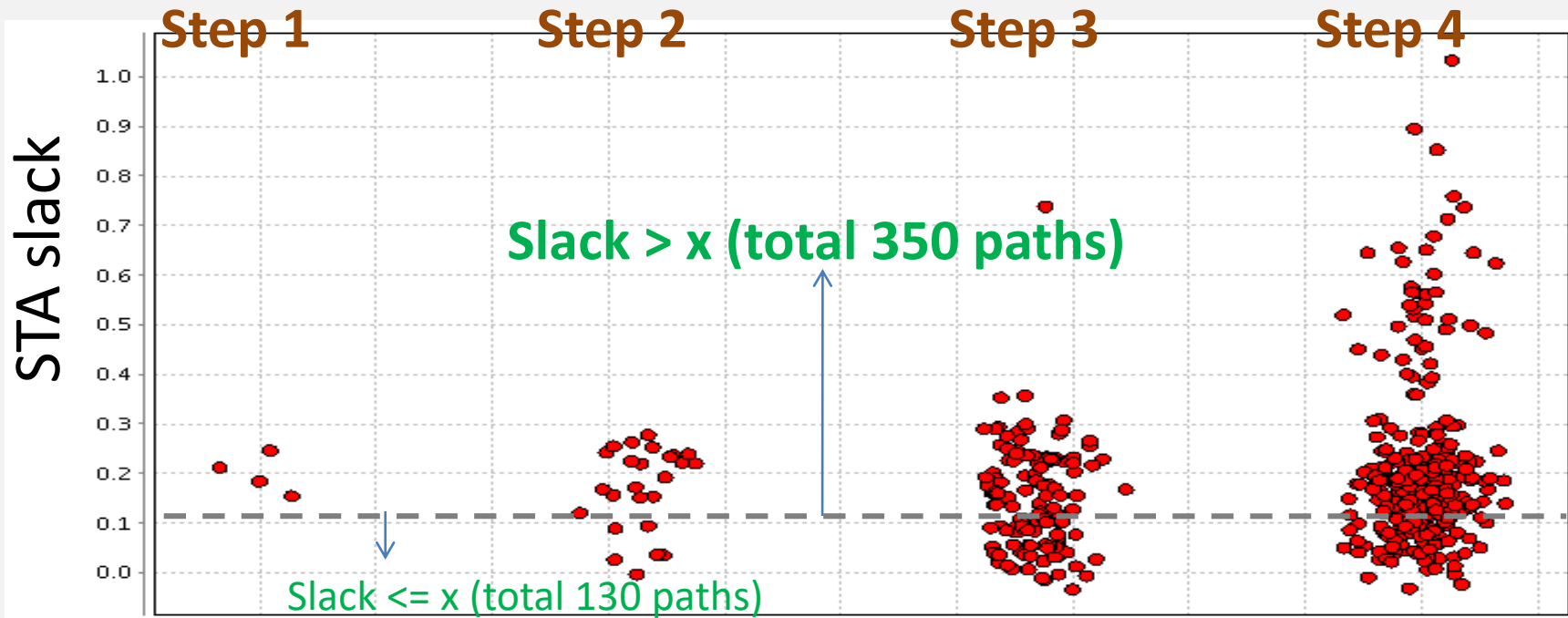
**Coverage Point $CP$**

⇨ Simulation Traces

➢ **Goal: to achieve more coverage on CP**

➢ **Approach: Analyze simulation traces to find out**
   - **What combination of signals can activate CP?**

➢ **Features: $f_1, f_2, \cdots, f_n$ are testbench-controllable signals**

➢ **Data: Few or no samples that cover CP**
   - **Positive Samples: 0 to few**
   - **Negative Samples: 1K to few K's**

# e.g. Physical Verification



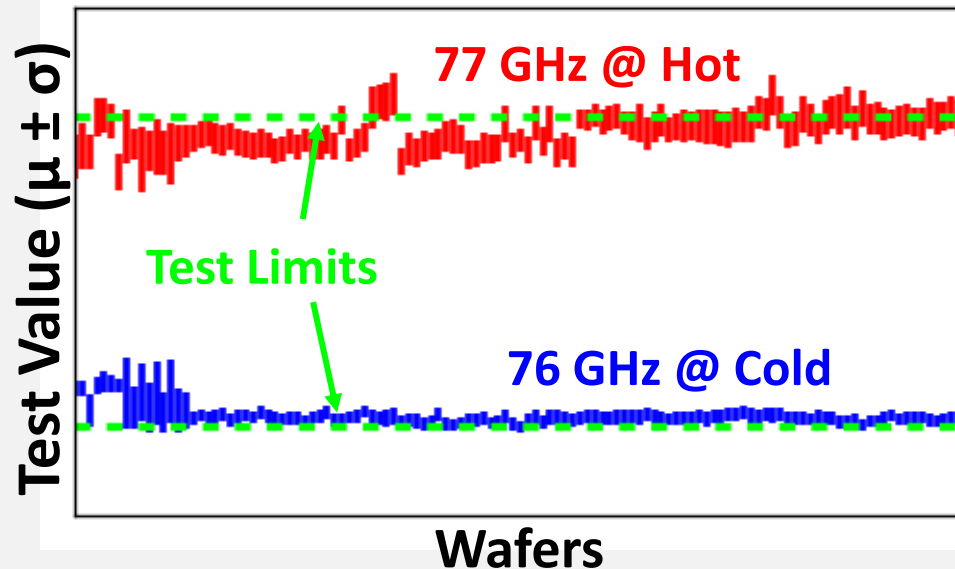- ➤ **Goal**: to model causes for an issue
- ➤ **Approach**: Analyze snippets of layout images to find out
  - – What combination of features can cause a issue?
- ➤ **Features**: $f_1, f_2, \cdots, f_n$ are developed based on domain knowledge to characterize geometry or material properties
- ➤ **Data**: Few samples for a particular type of issue
  - – Positive Samples: 1 to few
  - – Negative Samples: many

# e.g. Timing Verification



STA slack

Step 1    Step 2    Step 3    Step 4

Slack > x (total 350 paths)

Slack <= x (total 130 paths)

- ➢ **Goal**: to model causes for a miss-predicted silicon critical path
- ➢ **Approach**: Analyze unexpected silicon critical paths
  - – **What combination of design features can cause an unexpected critical path?**
- ➢ **Features**: $f_1, f_2, \cdots, f_n$ **are developed based on design knowledge to characterize a timing path**
- ➢ **Data**: Few samples for a particular type of critical path
  - – **Positive Samples: 1 to few**
  - – **Negative Samples: many (STA critical but not silicon critical – about \$25K paths)**

# e.g. Yield



- **Goal**: to find a receipt to improve yield
- **Approach**: Analyze wafer yield data with process parameters
  - Tuning what combination of process parameters can improve yield?
- **Features**: $f_1, f_2, \cdots, f_n$ are tunable process parameters
- **Data**: Samples can be parts or wafers
  - Positive Samples: Failing parts or Low-yield wafers
  - Negative Samples: Others

# Feature-Based Analytics

- ➢ **Problem:**
  - – **Search for a combination of features or feature values among a large set of features**
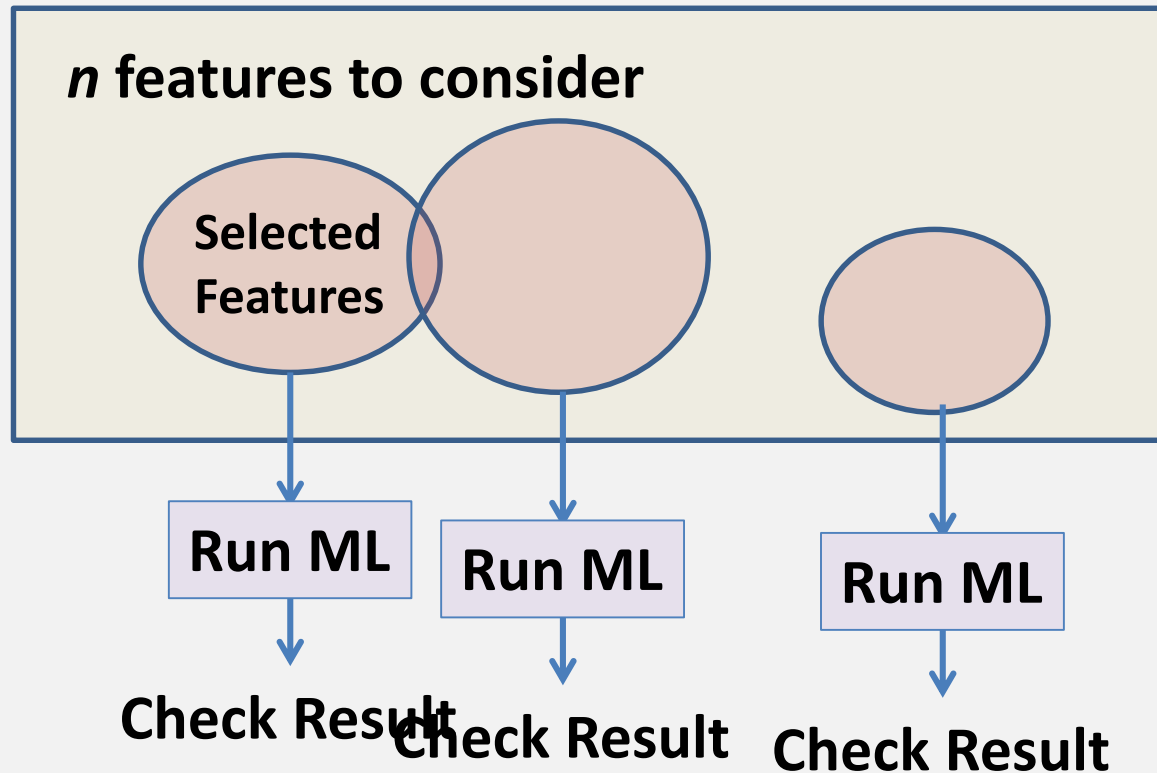
- ➢ **Data:**
  - – **Interested in positive samples**
  - – **Extremely unbalanced – Many more negative samples and very few positive samples**

- ➢ **Not a traditional feature selection problem**
  - – **Insufficient data**
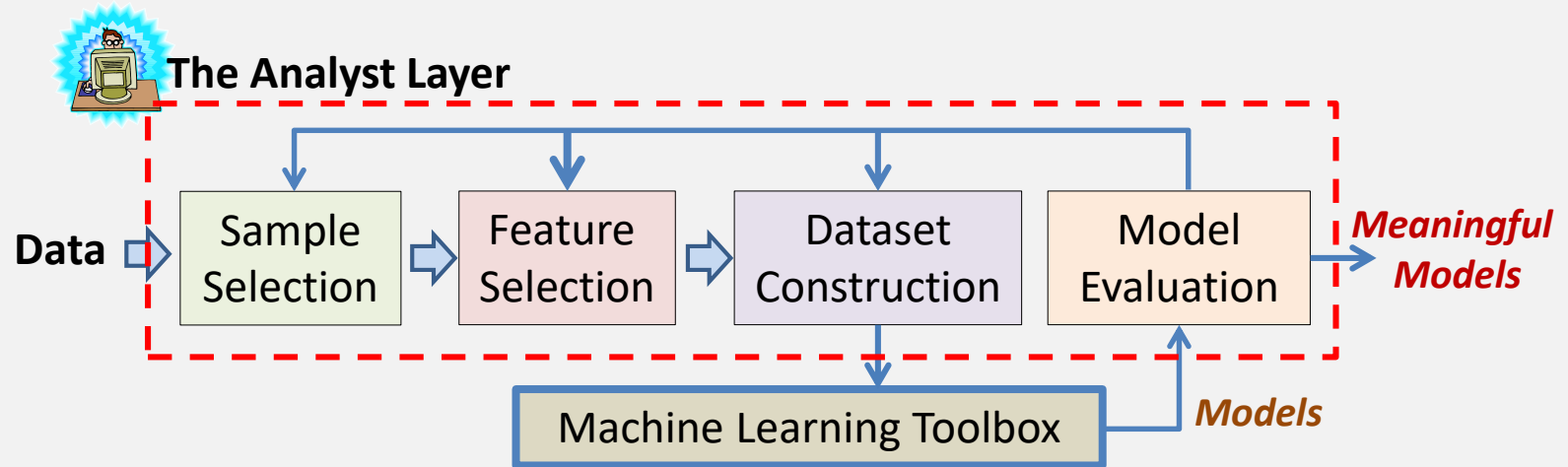  - – **Cannot apply cross-validation to check a model**

# In Practice, This Is What Happens

*n* features to consider

Selected Features

Run ML

Run ML

Run ML

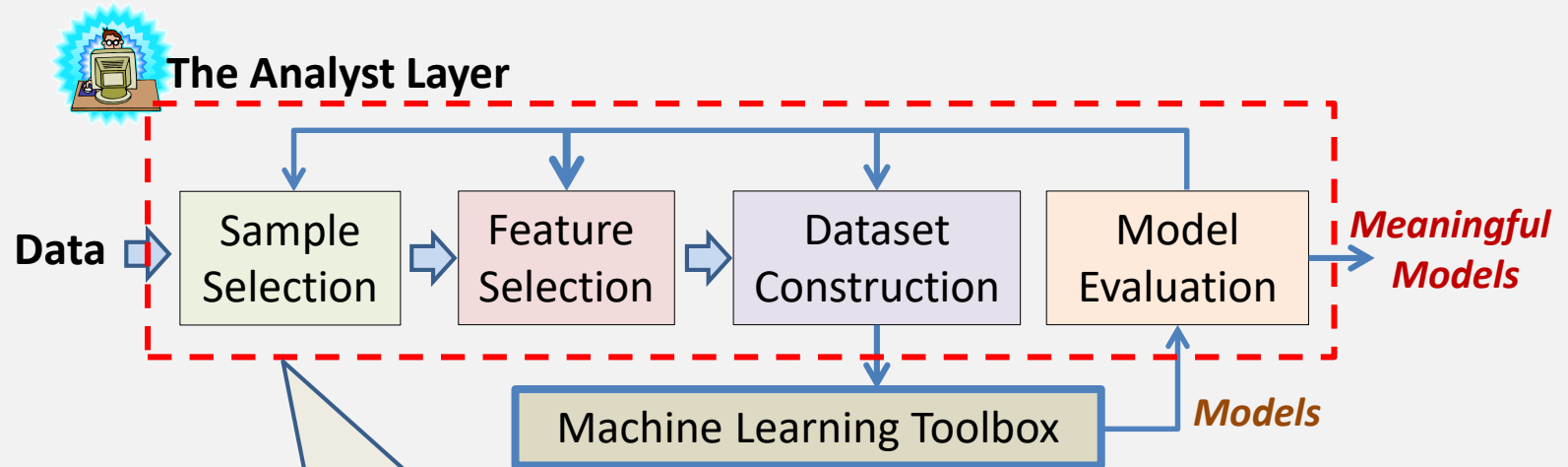Check Result

Check Result

Check Result

➢ **Learning from data becomes an iterative search process (usually run by a person)**

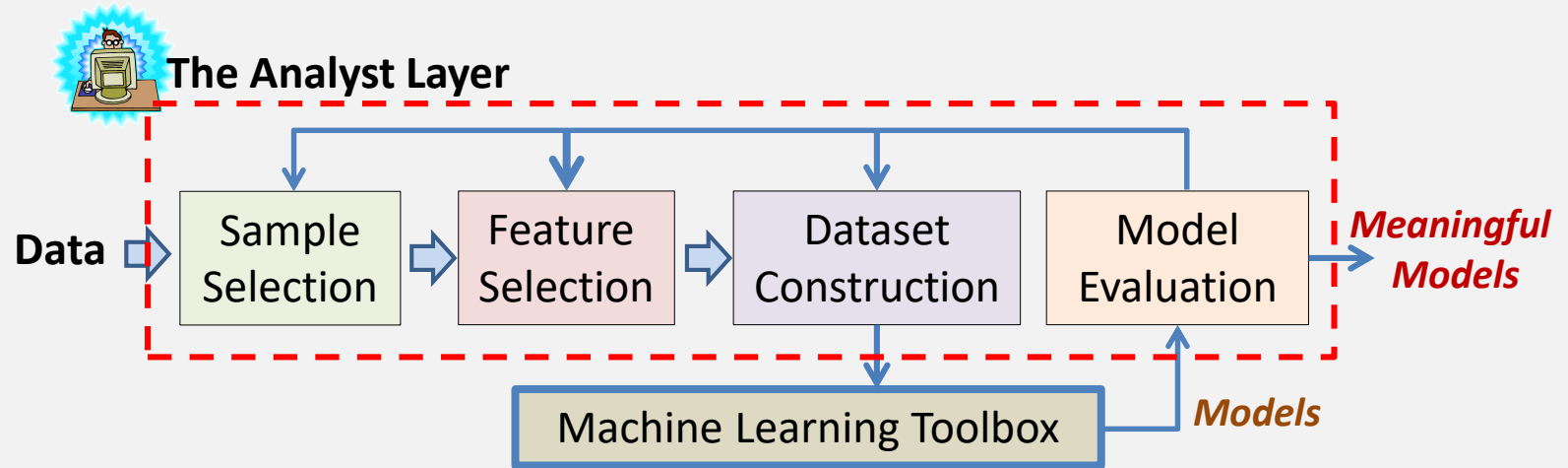# An Iterative Search Process



- ➢ **Learning is an iterative search process**
- ➢ **The analyst**
  - – **(1) Prepare the datasets to be analyzed**
  - – **(2) Determine if the results are meaningful**
- ➢ **The effectiveness depends on how the analyst conducts these two steps – not just about the tool in use!**

# Implications



**The Analyst Layer**

Data → Sample Selection → Feature Selection → Dataset Construction → Model Evaluation → *Meaningful Models*
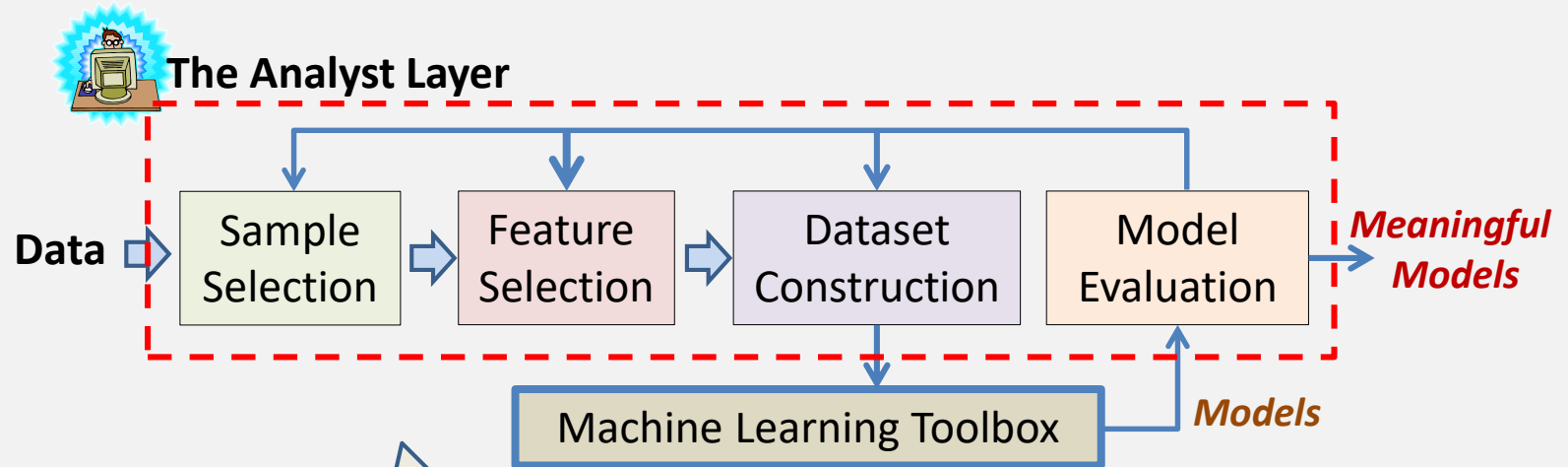
Machine Learning Toolbox → *Models*

**The effectiveness of the search largely depends on how the Analyst Layer is conducted**

# Implications



**The Analyst Layer**

Data → Sample Selection → Feature Selection → Dataset Construction → Model Evaluation → *Meaningful Models*

Machine Learning Toolbox *Models*

The **Analyst Layer** demands a **ROBUST Machine Learning Toolbox** where the model can be assessed **WITHOUT cross-validation**

# Implications



**The Analyst Layer**

Data → Sample Selection → Feature Selection → Dataset Construction → Model Evaluation → *Meaningful Models*
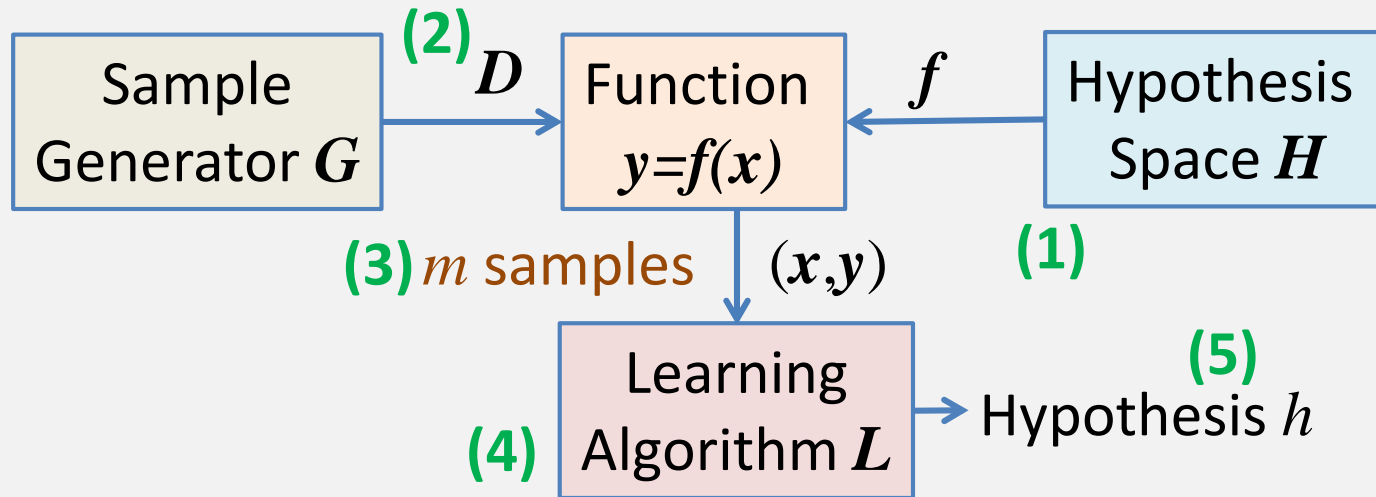
Machine Learning Toolbox

*Models*

**Automation requires automating both the Analyst Layer and the Machine Learning Toolbox**

# Machine Learning Toolbox

# Questions

➤ **Recall main issue: We can't apply cross-validation**

➤ **Why do we need cross-validation?**

➤ **Why can a machine learning algorithm guarantees the accuracy of its output model?**

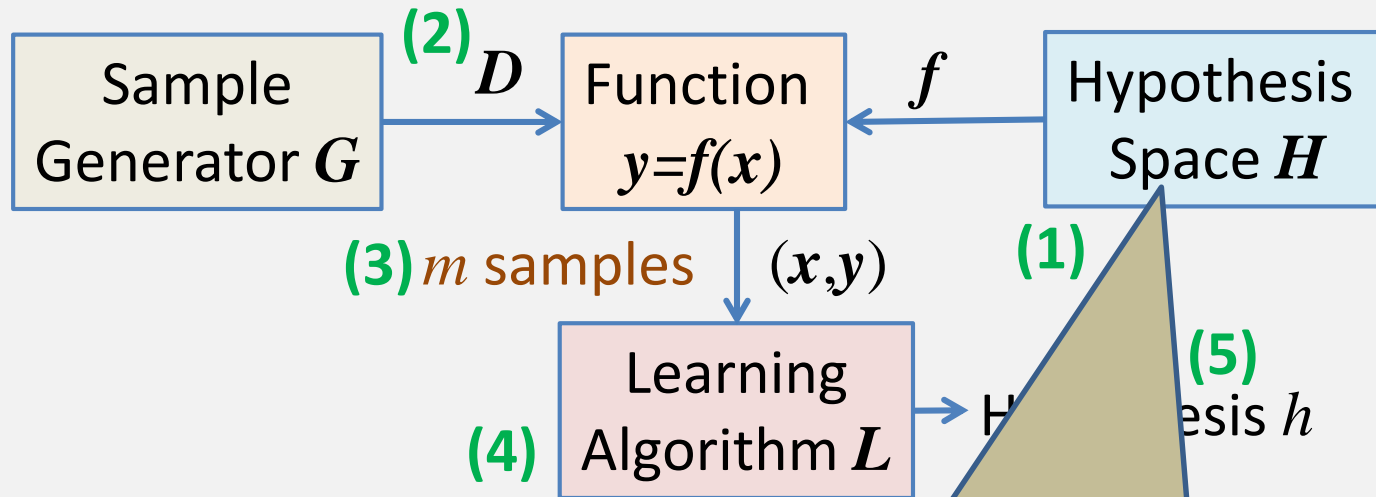➤ **What's a machine learning algorithm trying to optimize anyway?**

# Five Assumptions To Machine Learning



- ➢ **A restriction on *H* (otherwise, NFL)**
- ➢ **An assumption on *D* (i.e. not time-varied)**
- ➢ **Assuming size *m* is in order O(poly(*n*)), *n*: # of features**
- ➢ **Making sure a practical algorithm *L* exists**
- ➢ **Assuming a way to measure error, e.g. *Err(f(x), h(x))***

# In Practice

**(2)** $D$

Sample Generator $G$ → Function $y=f(x)$ ← $f$ Hypothesis Space $H$

**(3)** $m$ samples $(x,y)$

**(1)**

**(5)**

Learning Algorithm $L$ → Hypothesis $h$

**(4)**
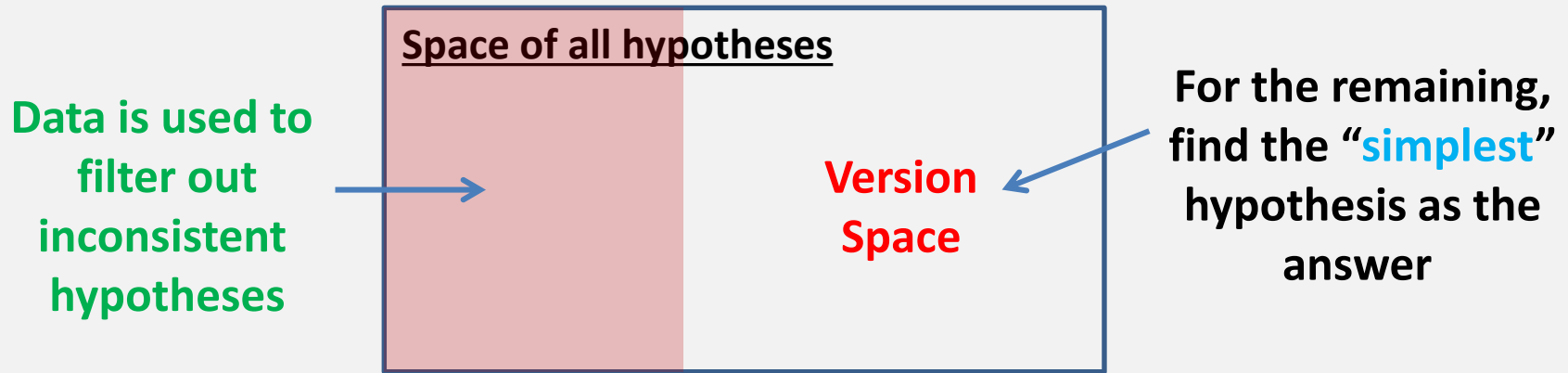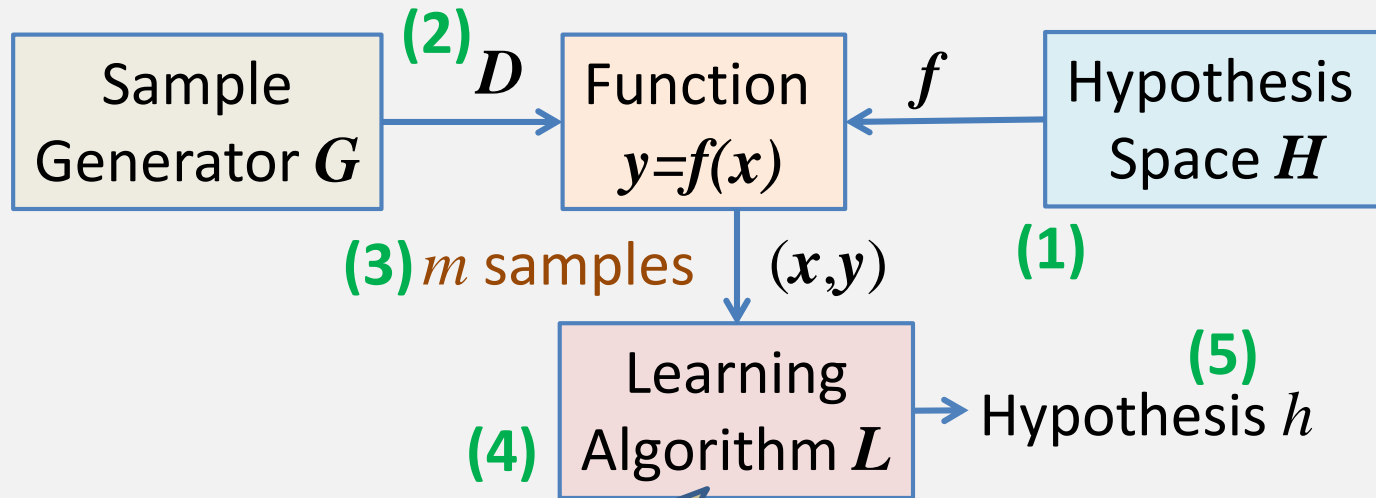
**Because we don't know how complex $H$ should be, we assume the most complex $H$ we can afford in training**

# As A Result, We Need Occam's Razor Assumption

**Space of all hypotheses**

**Data is used to filter out inconsistent hypotheses**

**Version Space**

**For the remaining, find the "simplest" hypothesis as the answer**

➢ **Hypothesis space: e.g. all possible assignment of weight values in a neural network (can be infinite)**

➢ **Occam's Razor (Regularization): Find the "simplest" hypothesis that fit the data**
  – **Hence, many machine learning algorithms solve a non-convex constrained minimization problem (NP-Hard or Harder)**

➢ **However, the simplicity measure might not be meaningful in an application context**

# In Practice

Sample Generator $G$ **(2)** $D$ → Function $y=f(x)$ ← $f$ Hypothesis Space $H$

**(3)** $m$ samples $(x,y)$ **(1)**

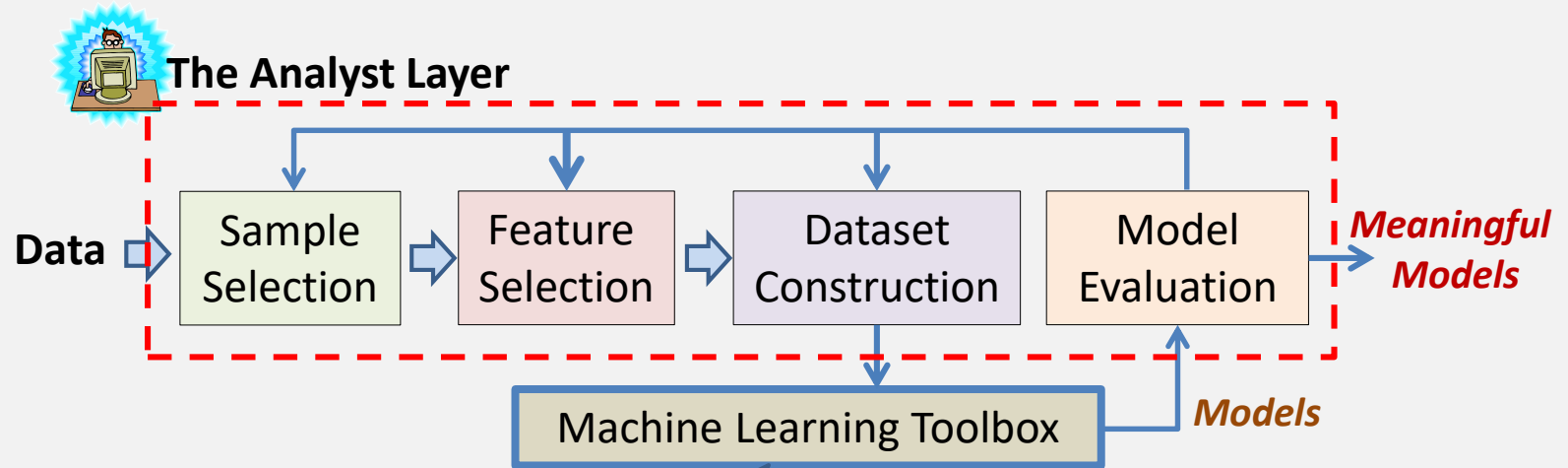Learning Algorithm $L$ → Hypothesis $h$ **(5)**

**(4)**

**Because non-convex optimization is hard, some heuristic is used, and the solution is often a local minimum**

# Many Things Are Not Ideal

➢ **Your assumption of the hypothesis space might be too simple (underfitting) or too complex (overfitting)**

➢ **You may not have sufficient data to identify the exact answer from your assumed hypothesis space**

➢ **Your learning algorithm is only a heuristic and does not guarantee to find the "optimal" model**

➢ **As a result, you need cross-validation**

**The Analyst Layer**

Data → Sample Selection → Feature Selection → Dataset Construction → Model Evaluation → *Meaningful Models*

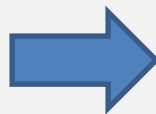Machine Learning Toolbox → *Models*

**Can we have a ML tool that can produce a model with some guarantee, without using Cross-Validation?**
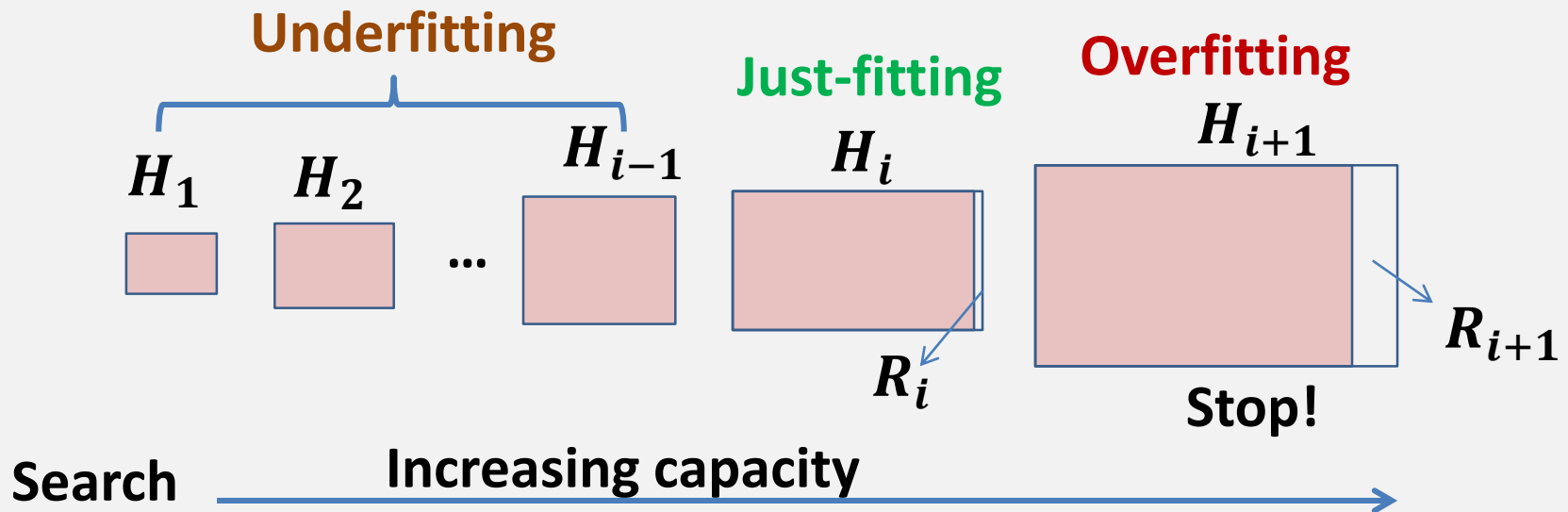
# Alternative Machine Learning View

➢ **Traditional machine learning: Find an optimal model based on the given dataset**

➢ **Alternative machine learning: Find an interpretable Hypothesis Space Assumption *H* where a model can JUST-FIT the dataset but not overfitting**

**Search for a Model** ➡ **Search for An Assumption**

**Underfitting**     **Just-fitting**     **Overfitting**

$H_1$   $H_2$   ...   $H_{i-1}$   $H_i$   $H_{i+1}$

$R_i$     $R_{i+1}$

**Stop!**

**Search**    **Increasing capacity** →

- ➢ **Search for the "JUST-FIT" hypothesis space**
  - – **Such that the output model among the few answers consistent with all the samples**

- ➢ **The JUST-FIT hypothesis space (if exists) can be a measure of quality for the model**

# VeSC-CoL:

# Our Concept Learning Tool

# VeSC-CoL

- **Reference : Kuo-Kai Hsieh and Li-C. Wang,** A Concept Learning Tool Based On Calculating Version Space Cardinality, **arXiv:1803.08625 [cs.AI], Mar 23, 2018**

- **Handle binary-valued features**

- **Target (interpretable) concept:** *k-term DNF, for small k*

- **Designed to handle extremely-unbalanced dataset without cross-validation**

- **Two implementations: SAT-Based and OBDD-Based**

# K-term DNF – Terminology

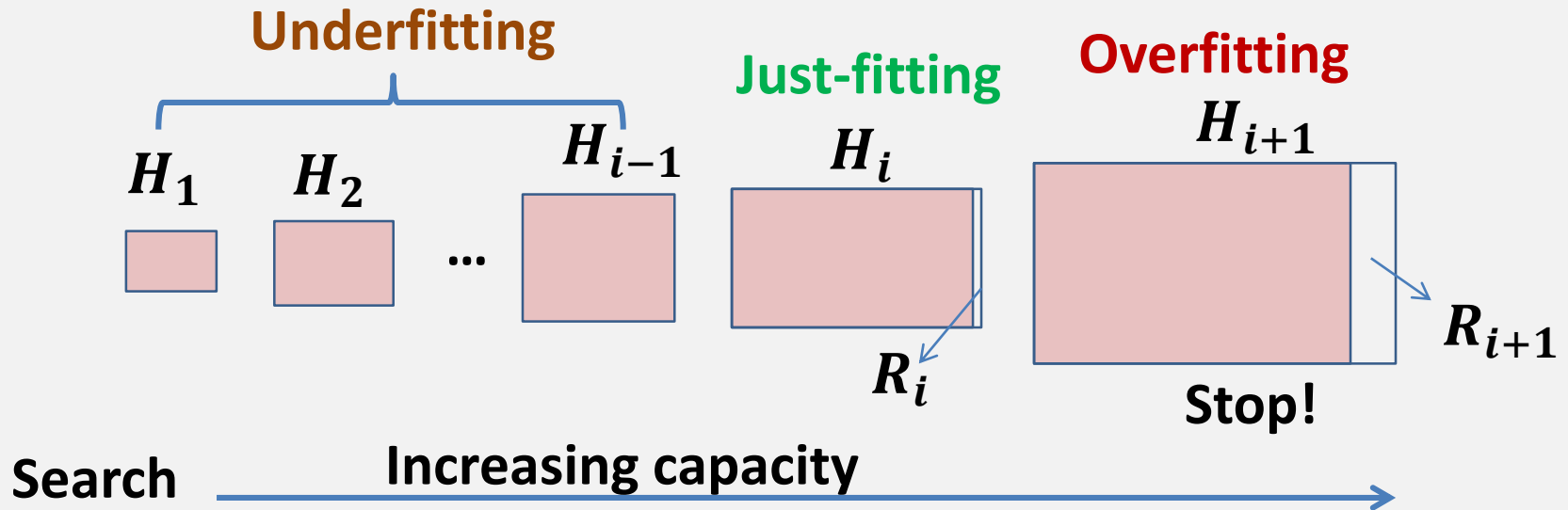$$x_1 \overline{x_2} x_4 \longrightarrow \text{1-term DNF or \textbf{Monomial}}$$

**Length $l$ = number of literals = 3**

$$x_1 \overline{x_2} x_4 + \overline{x_4} x_6 \longrightarrow \text{2-term DNF or Monomial}$$

**Length $l$ = number of literals = 3+2 = 5**

$n$ **= number of features (variables)**

# VeSC-CoL's Hypothesis Space Search



- **Given an upper bound on $k$ for $k$-term DNF**
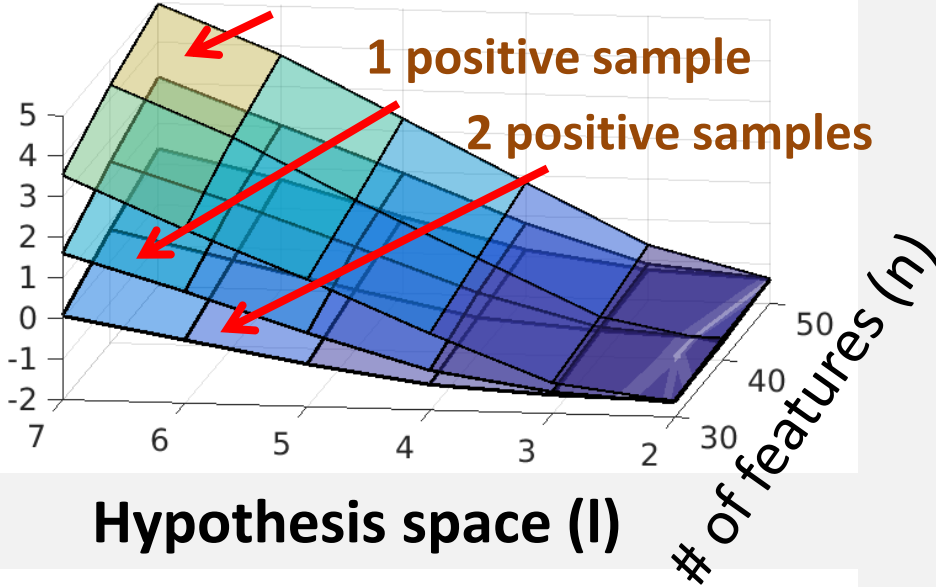- **$H_l$ is the hypothesis space for all hypotheses with length $l$**

**OBDD**

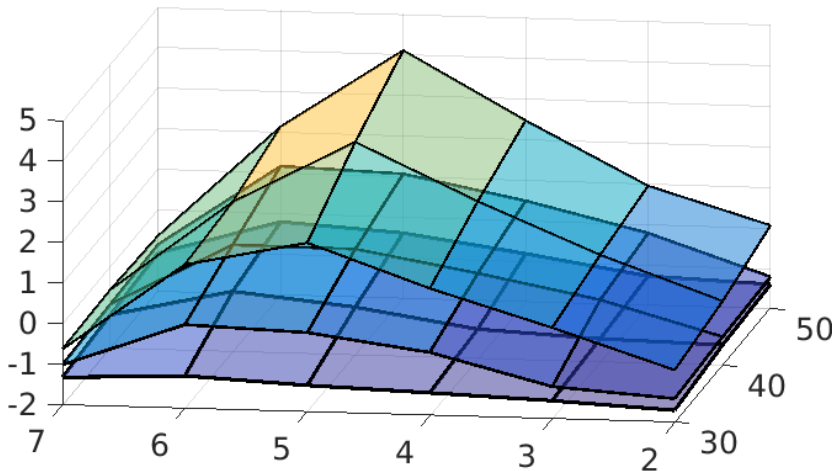**0 positive sample**

**1 positive sample**

**2 positive samples**

$\log_{10}(\text{second})$

Hypothesis space (l)

# of features (n)

**SAT**

> **Correct answer is with $l$ = 5**

> **$n$ does not affect runtime much**

> **$l$ limits how far we can search**

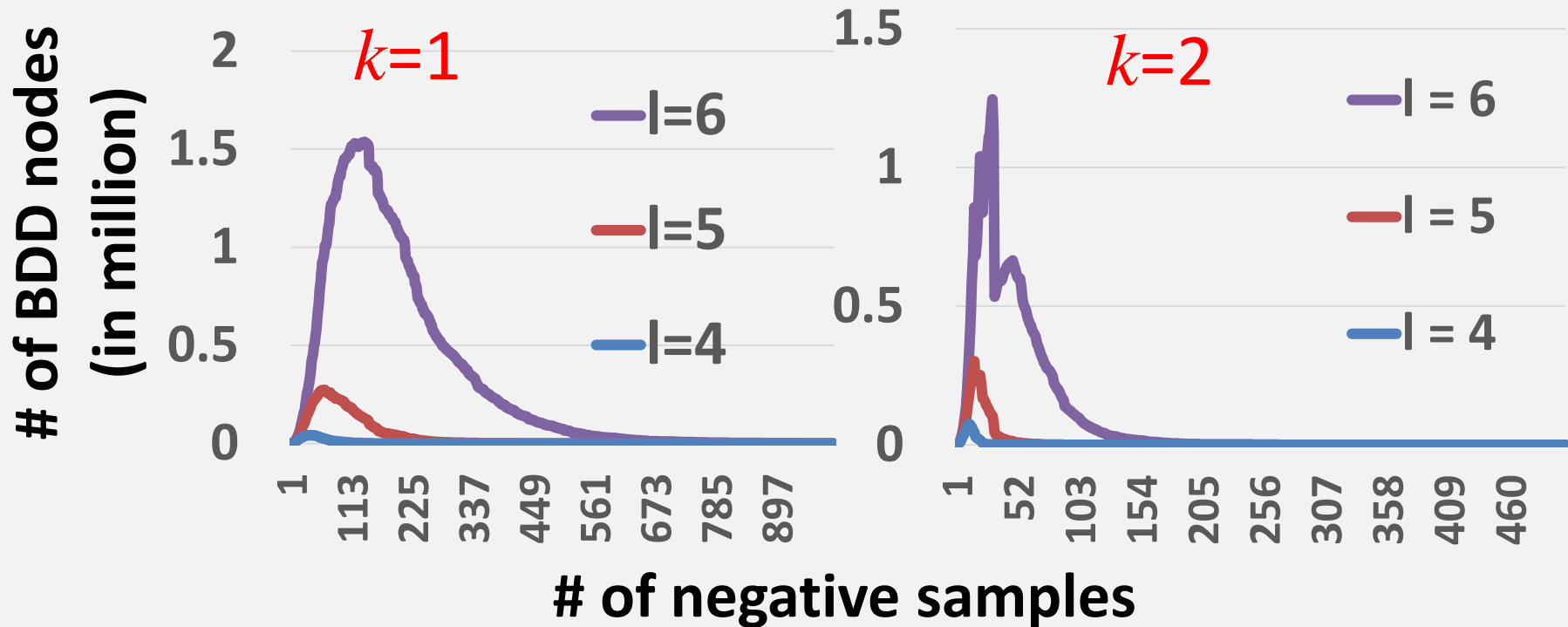➢ **As $n$ increases, you are likely to run out of time than to run out of data (assuming most are negative samples)**

➢ **For BDD-based implementation, the runtime wall happens in the early processing of the negative samples**
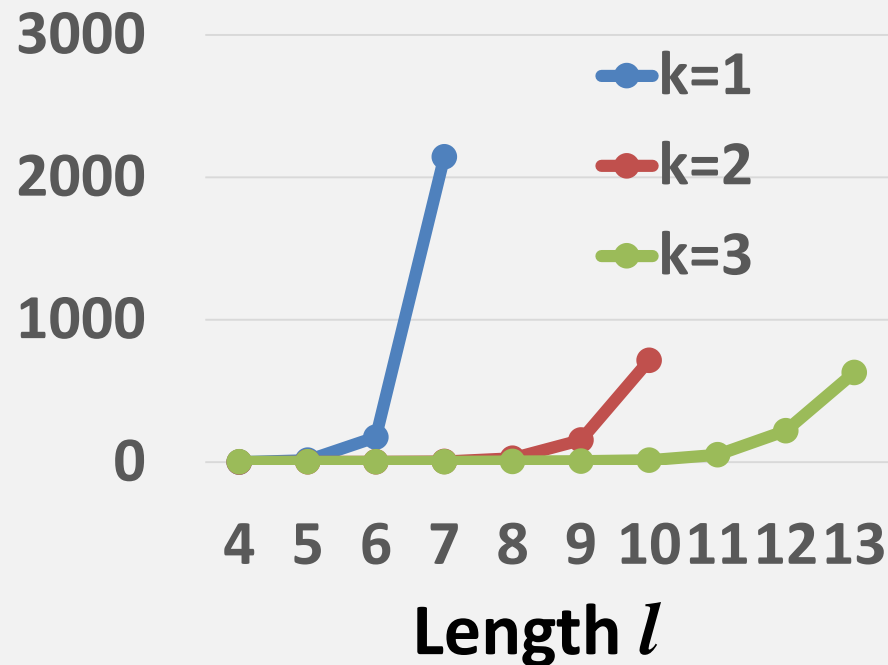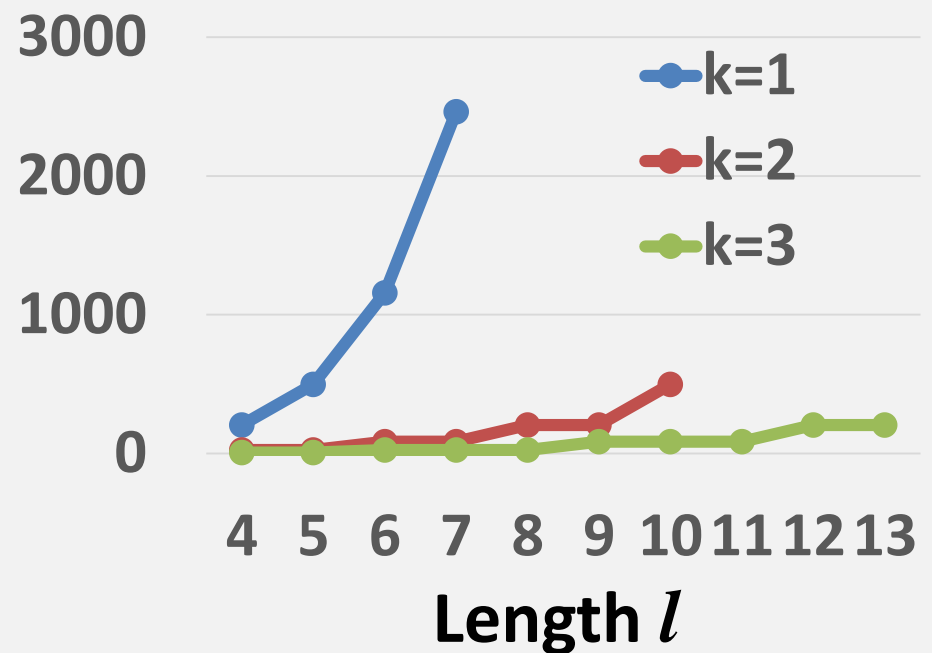


**# of negative samples**

**Number of features:** $n=100$

# Interesting Finding

➤ **Requirement for learning the "$k$=1" space dominates the requirements for learning the "$k$>1" spaces**



Number of features: $n$=100

# Guarantee by VeSC-CoL

- ➢ **Assuming the correct answer can be represented as a k-term DNF for a selected k, then VeSC-CoL always find the answer (assuming runtime is allowed)**
  - **Experimentally shown for $k$ up to 3, $l$ up to 8, negative sample size up to 10K**

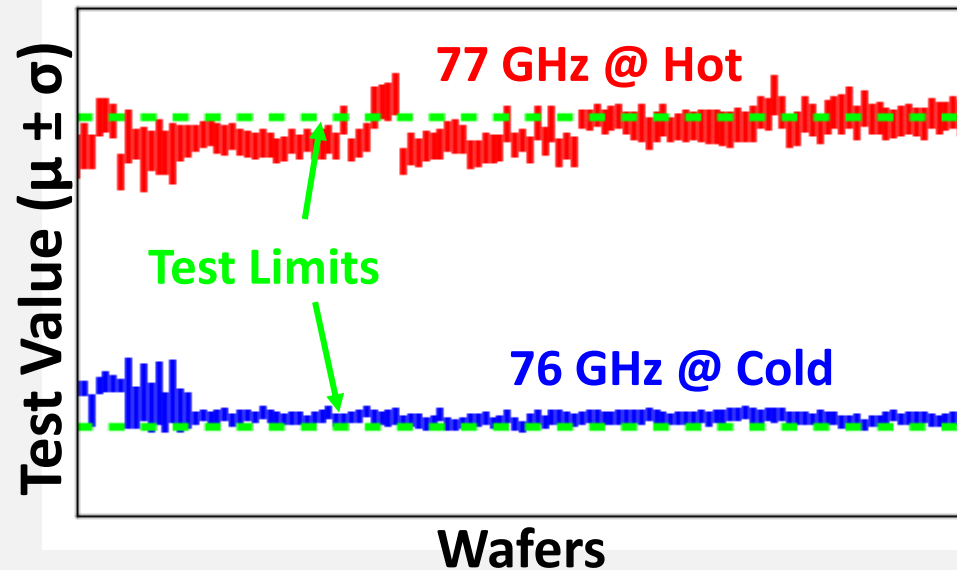| VeSC-CoL | CART | ID3 |
|---|---|---|
| $x_2 x_{63} \overline{x_{75}} x_{78} \overline{x_{80}}$ | $x_3 x_4 x_{28} x_{47} \overline{x_{53}} \overline{x_{55}} \overline{x_{80}}$ | $x_2 x_3 x_4 \overline{x_{30}} x_{47} \overline{x_{53}} \overline{x_{81}}$ |
| $x_{39} \overline{x_{45}} x_{72} \overline{x_{74}} x_{95}$ | $\overline{x_5} x_{16} x_{35} \overline{x_{45}} \overline{x_{55}} x_{56} x_{59}$ | $x_8 x_{40} \overline{x_{45}} x_{64} \overline{x_{74}} x_{87}$ |
| $\overline{x_2} x_{14} x_{52} \overline{x_{57}} x_{87}$ | $x_{11} \overline{x_{14}} \overline{x_{24}} x_{61} x_{64} x_{90} \overline{x_{92}}$ | $\overline{x_5} x_6 x_{16} x_{35} \overline{x_{45}} \overline{x_{56}} x_{59}$ |
| $x_{40} x_{45} x_{64} \overline{x_{74}} x_{87}$ | $\overline{x_4} x_8 \overline{x_{45}} \overline{x_{47}} x_{64} \overline{x_{74}} x_{89}$ | $\overline{x_2} x_{14} \overline{x_{24}} x_{61} x_{64} x_{90} \overline{x_{92}}$ |
| $\overline{x_{57}} x_{58} x_{77} \overline{x_{95}} x_{98}$ | $\overline{x_5} x_{29} x_{38} \overline{x_{43}} \overline{x_{79}} x_{99} + \overline{x_3} \overline{x_5} x_{29} x_{38} \overline{x_{43}} x_{49} \overline{x_{79}} x_{99}$ | $\overline{x_5} x_6 \overline{x_{11}} \overline{x_{14}} \overline{x_{18}} x_{34} x_{45}$ |

**Always Correct** | **Always Incorrect** | **Always Incorrect**
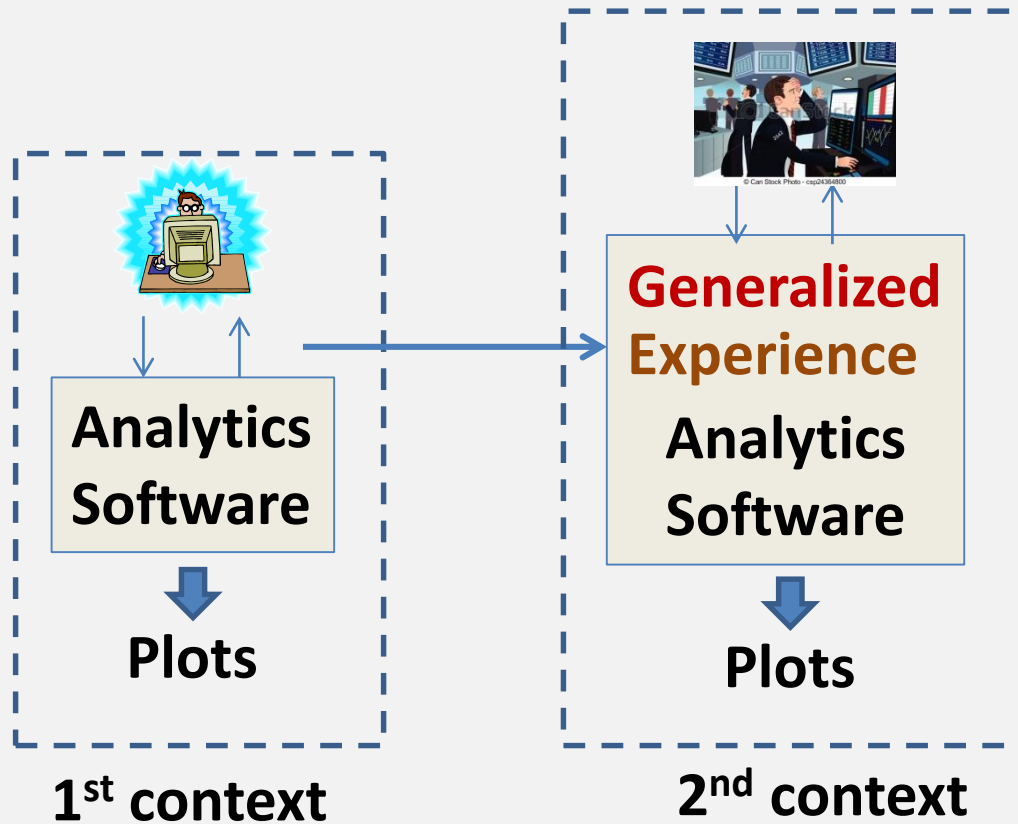
# Analyst Layer Automation

Wafers

- ➢ **Before this example, we had done work for resolving another yield issue for another product line**

- ➢ **Question: Can we learn to model the experience from that work and automate the Analyst Layer to resolve this yield issue**
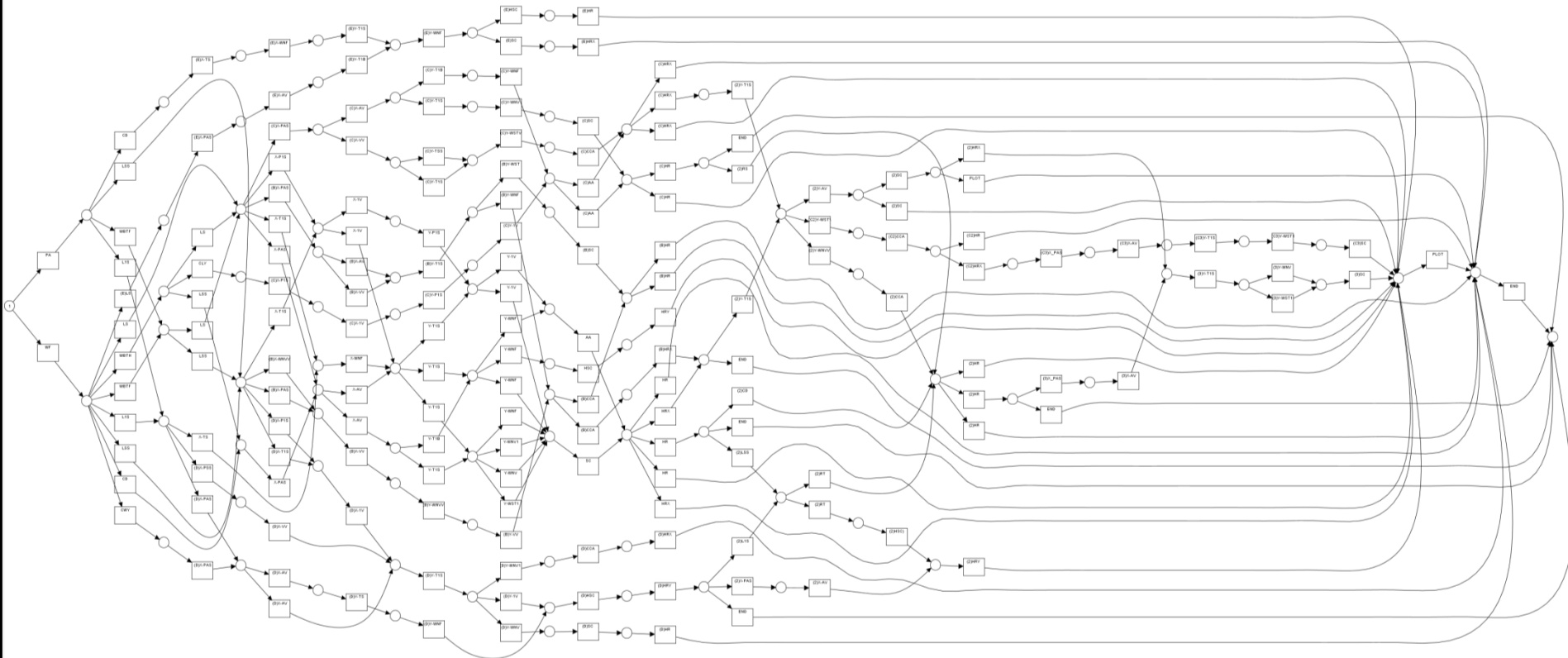
# The Learning Objective



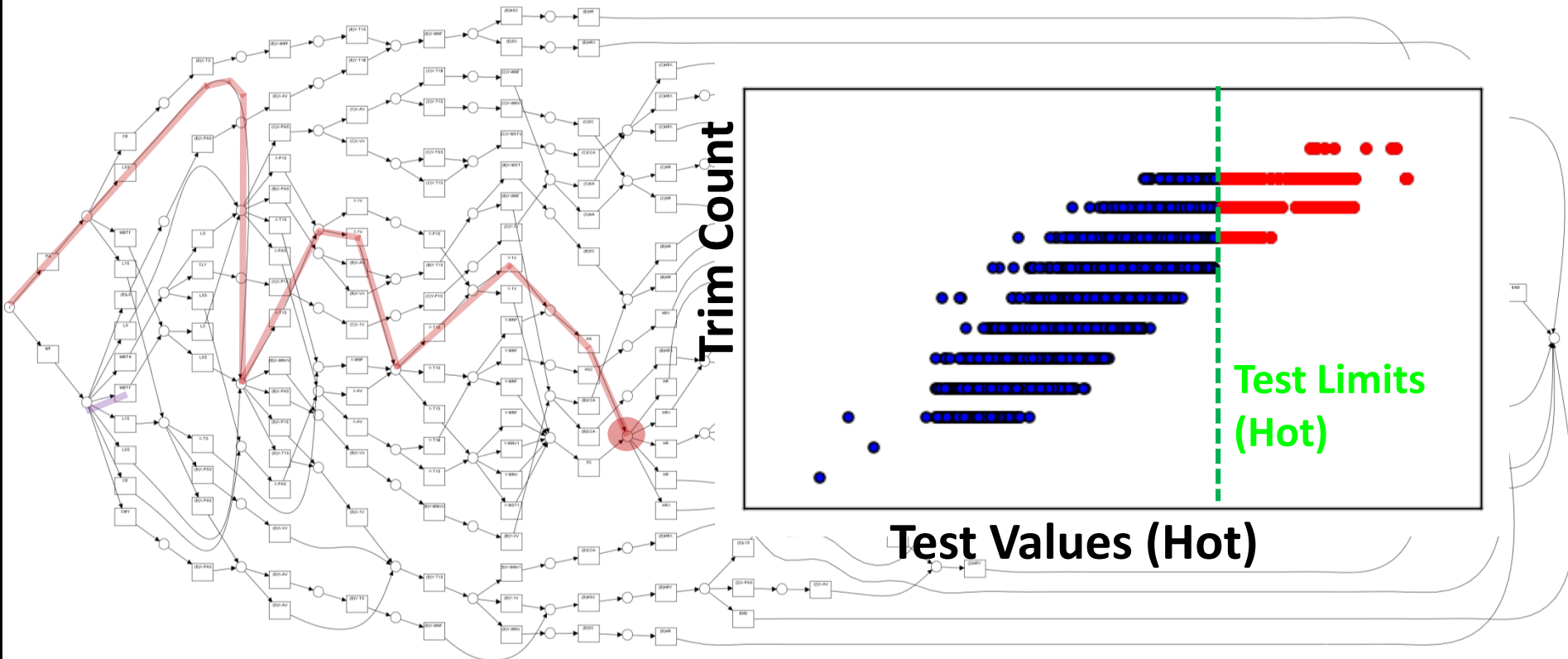**1st context**                **2nd context**

# Modeling "Experience"

➢ **To learn from analyst's experience, we need to have a way to model the experience**

➢ **Knowledge acquisition**

  – **Define a set of operators**

  – **Model experience as "an execution path" following a sequence of operators**
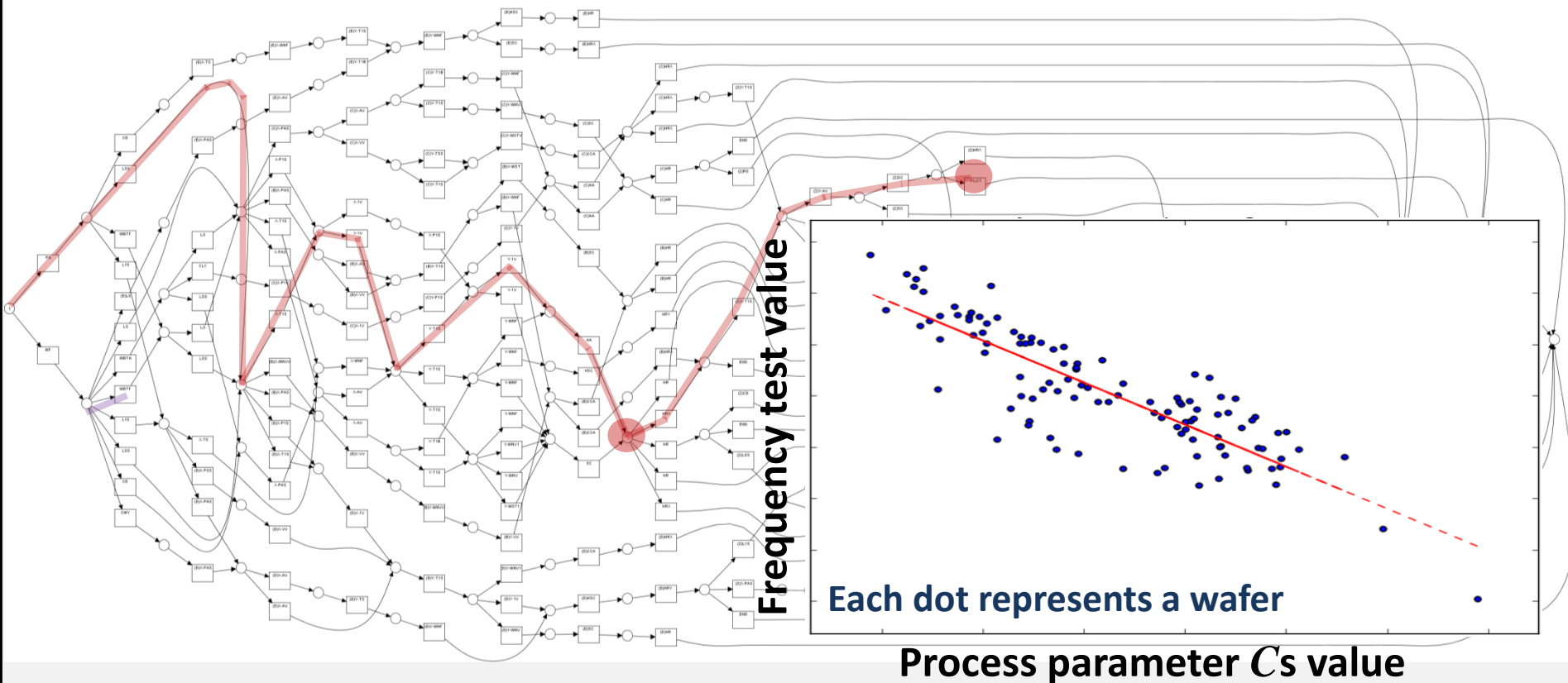
# Processing Mining Model



> **Record execution paths in a log file**

> **Apply process learning to learn from the log file**

> **Obtain a Process Model as shown above**

Trim Count

Test Values (Hot)

Test Limits (Hot)

> **Discover trim count is relevant to hot fails**

Each dot represents a wafer

Frequency test value

Process parameter $C$'s value

➢ **Determine that parameter C affects the frequency test value which decides the trim count**

# Summary: Three Observations

- ➢ **The effectiveness of "Machine Learning" largely depends on how the Analyst Layer is conducted**

- ➢ **Automation of "Machine Learning" needs to include automation of the Analyst Layer**

- ➢ **Traditional machine learning tools are not designed to effectively support the Analyst Layer**
    - **Require an Alternative ML view and a learning tool designed to be used without Cross-Validation**

# THANK YOU!