

# **Learning Actionable Insights**

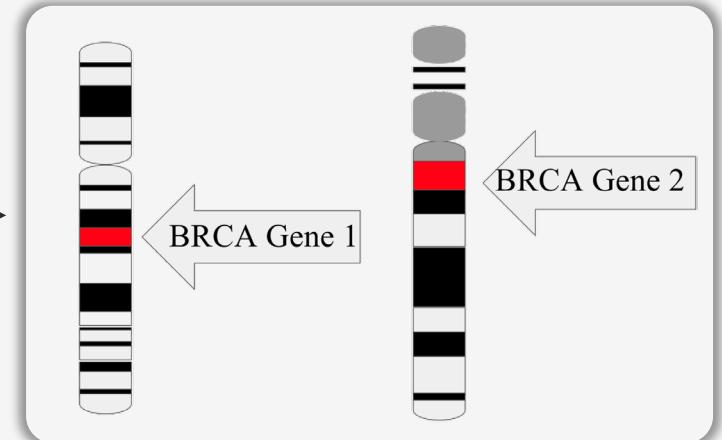
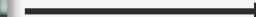
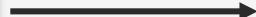
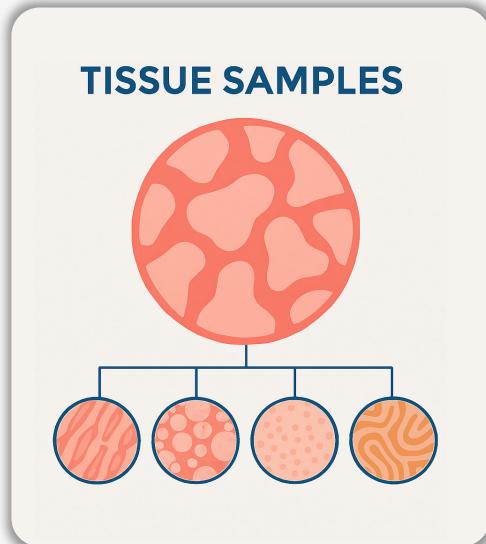
*from scientific data and models*

Jilles Vreeken | 19 May 2025



# Common questions in biology

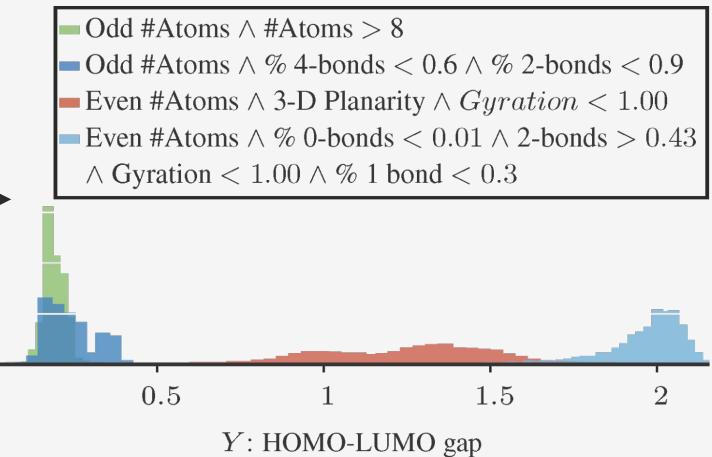
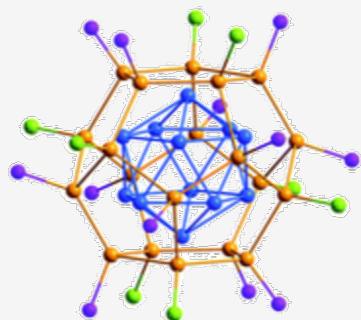
*What is the difference between cancerous and benign tissue?*





# Common questions in physics

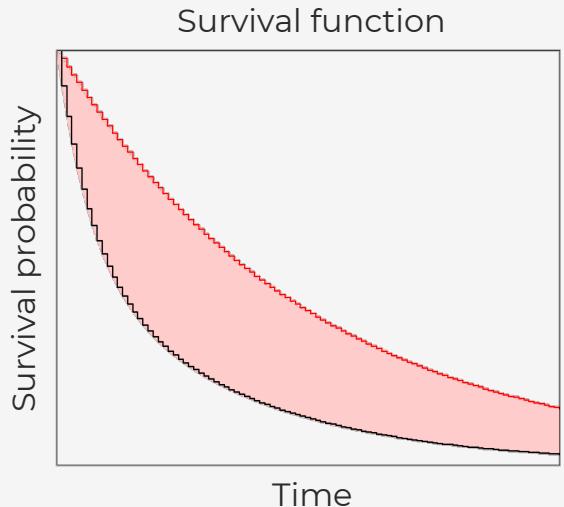
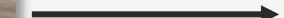
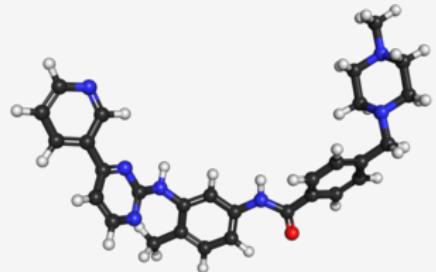
*What materials have exceptional conductivity?*





# Common questions in medicine

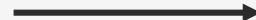
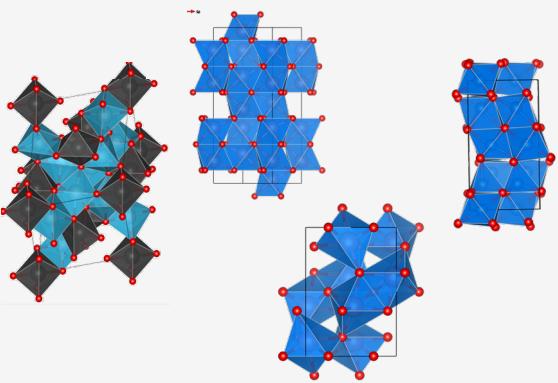
*Which people do (not) benefit from a treatment?*



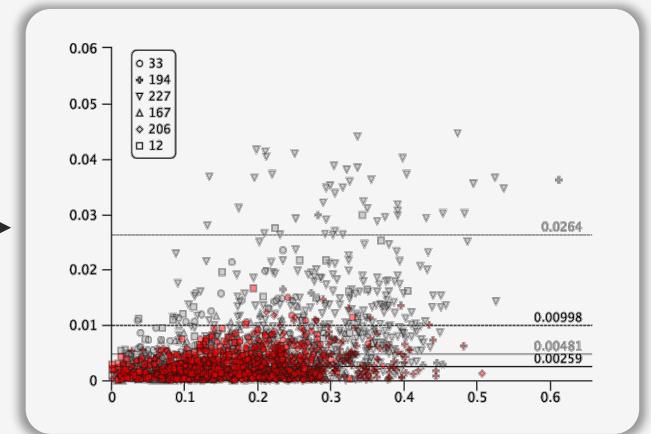
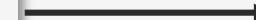


# Common questions in reality

*What is the domain of applicability of my ML model?*



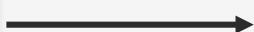
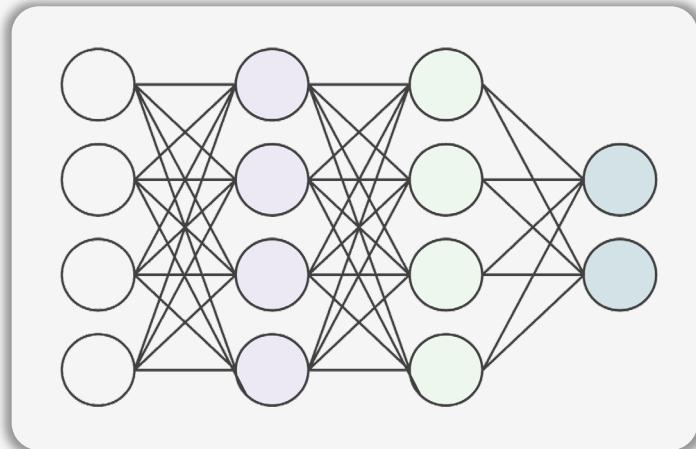
Determine conditions  
on  $x$  under which  $\hat{f}(x)$   
is reliable





# Common questions in interpretability

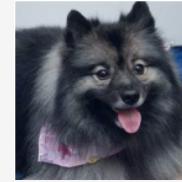
*What internal reasoning influences predictions?*



Extract interpretable  
concepts



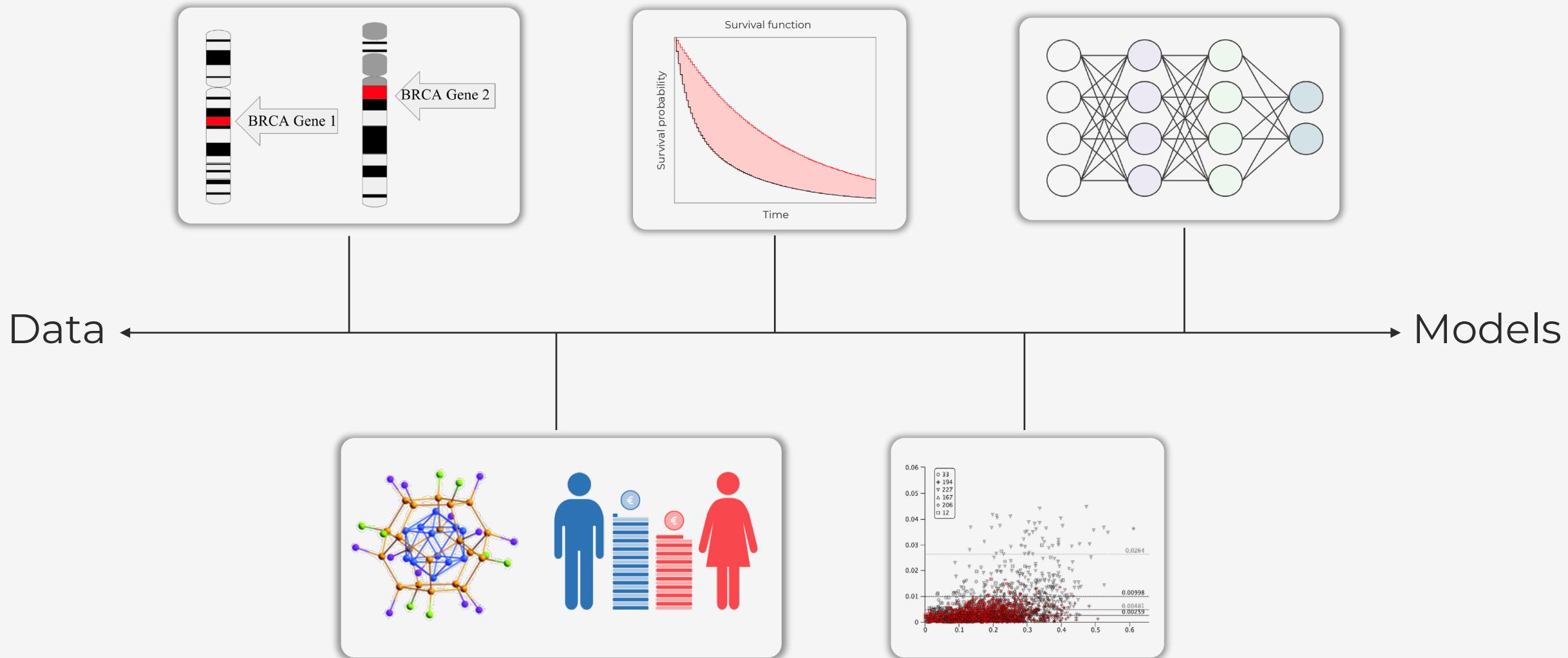
Concept 1      Concept 2



Samoyed

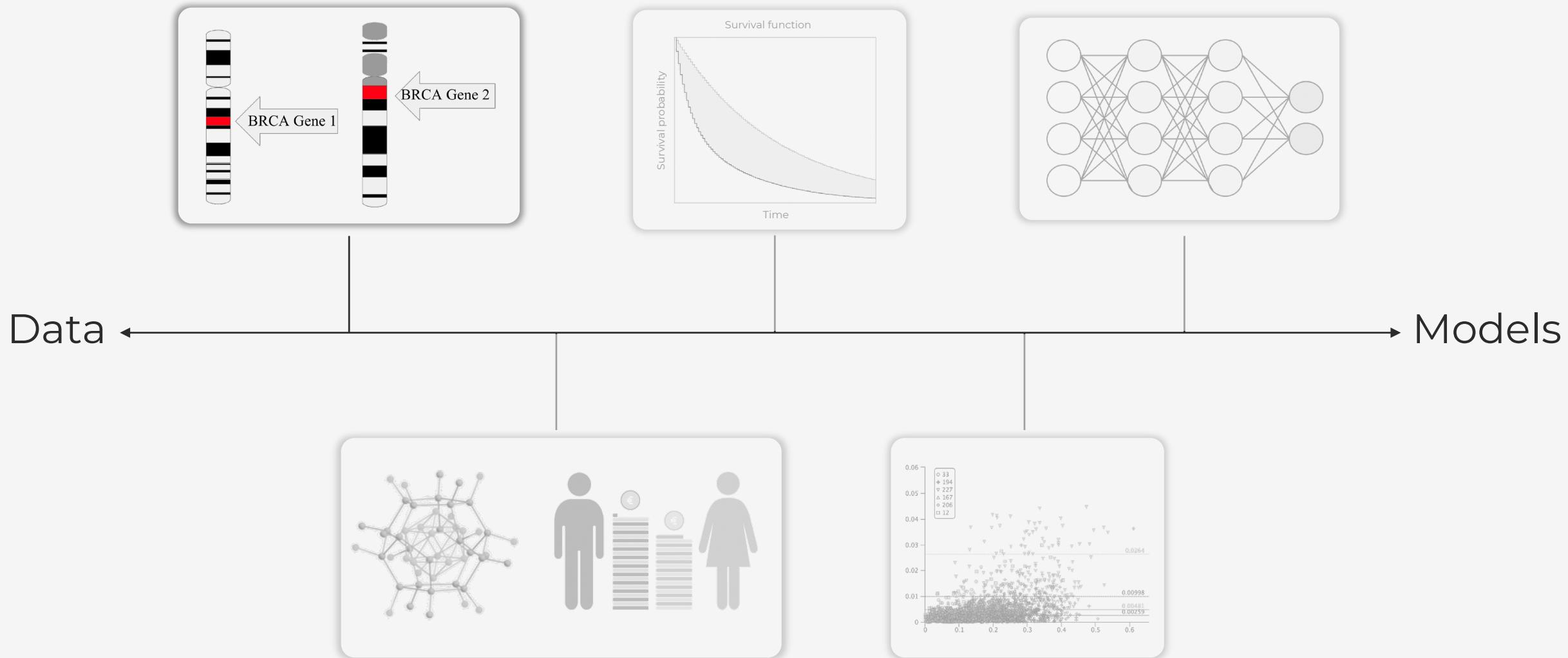


# Common questions in ML

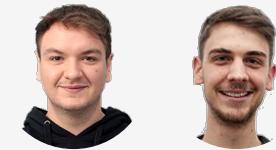




# Common questions in ML

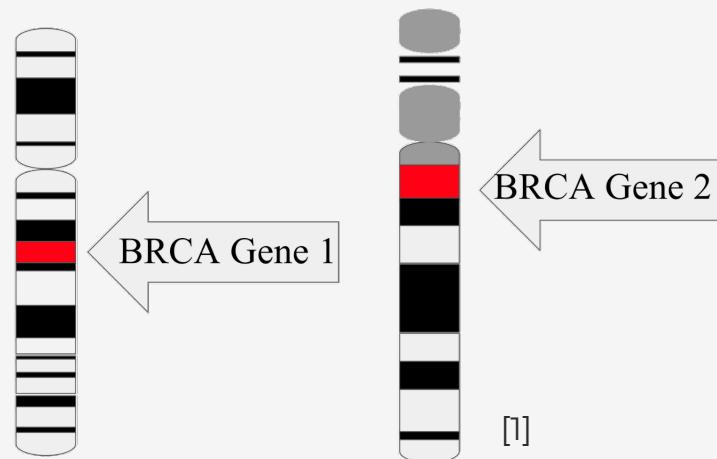


# Understanding genetic data



## Breast Cancer Data (BRCA)

- (Binarized) sequences of tissue samples
- Labeled with types of BRCA
- Very high-dimensional, low #samples



**Goal:** Find class-specific interactions of genes

## Class-specific pattern:

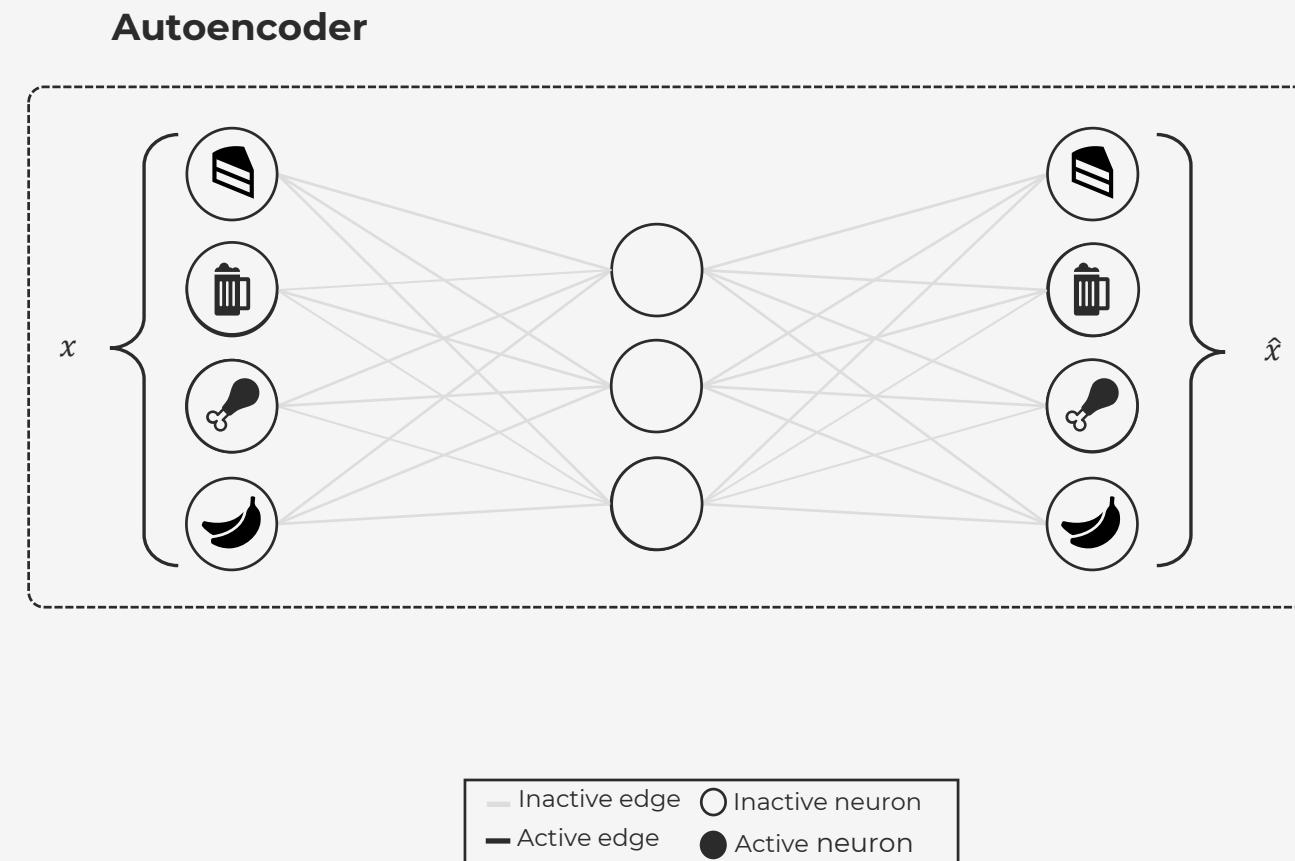
1. Occurs more often in class k
2. Is most predictive for class k

Features $F$	$Y$
1 1 1 1	<i>r</i>
1 1 1 1	<i>r</i>
1 0 1 1	<i>r</i>
1 1 1 1	<i>s</i>
1 1	<i>s</i>
1 1	<i>s</i>
1 1	<i>s</i>

# BiNAPS – In a nutshell



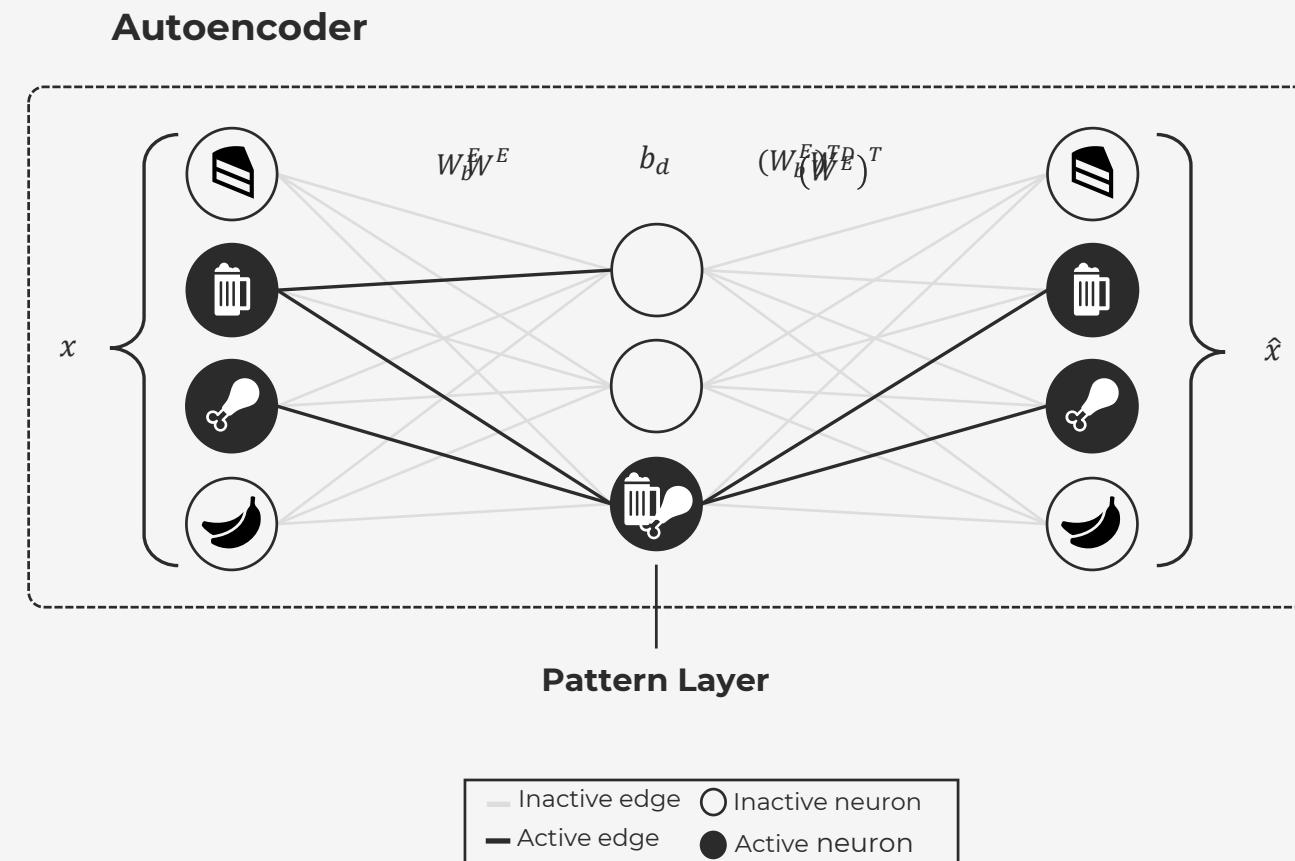
**Main Idea:** Compression → Patterns



# BiNAPS – In a nutshell



Main Idea: Compression → Patterns

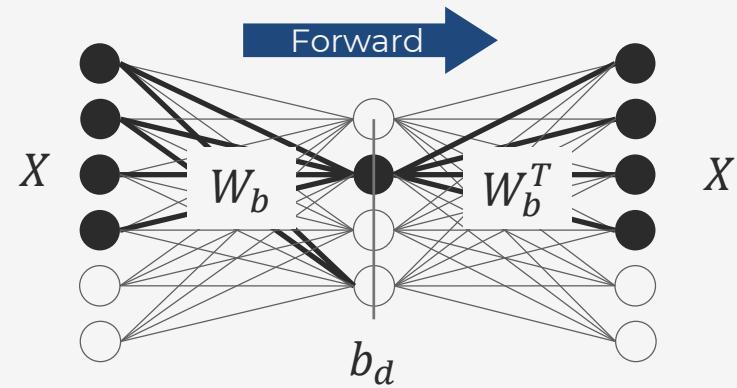


# BiNAPS – In a nutshell



For  $X \in D$

1. reconstruct

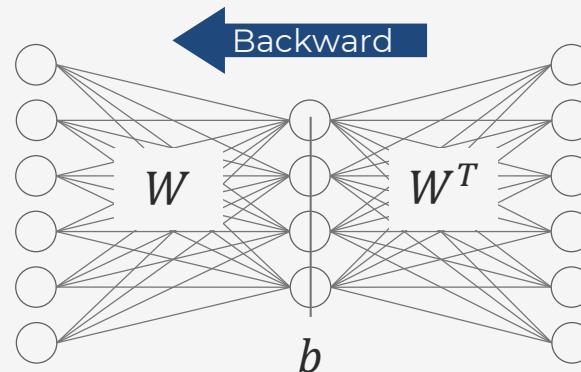


2. compute loss

$$L_{W,b}(X)$$

3. backpropagate

$$\frac{dL_{W,b}(X)}{dW} \quad \frac{dL_{W,b}(X)}{db}$$



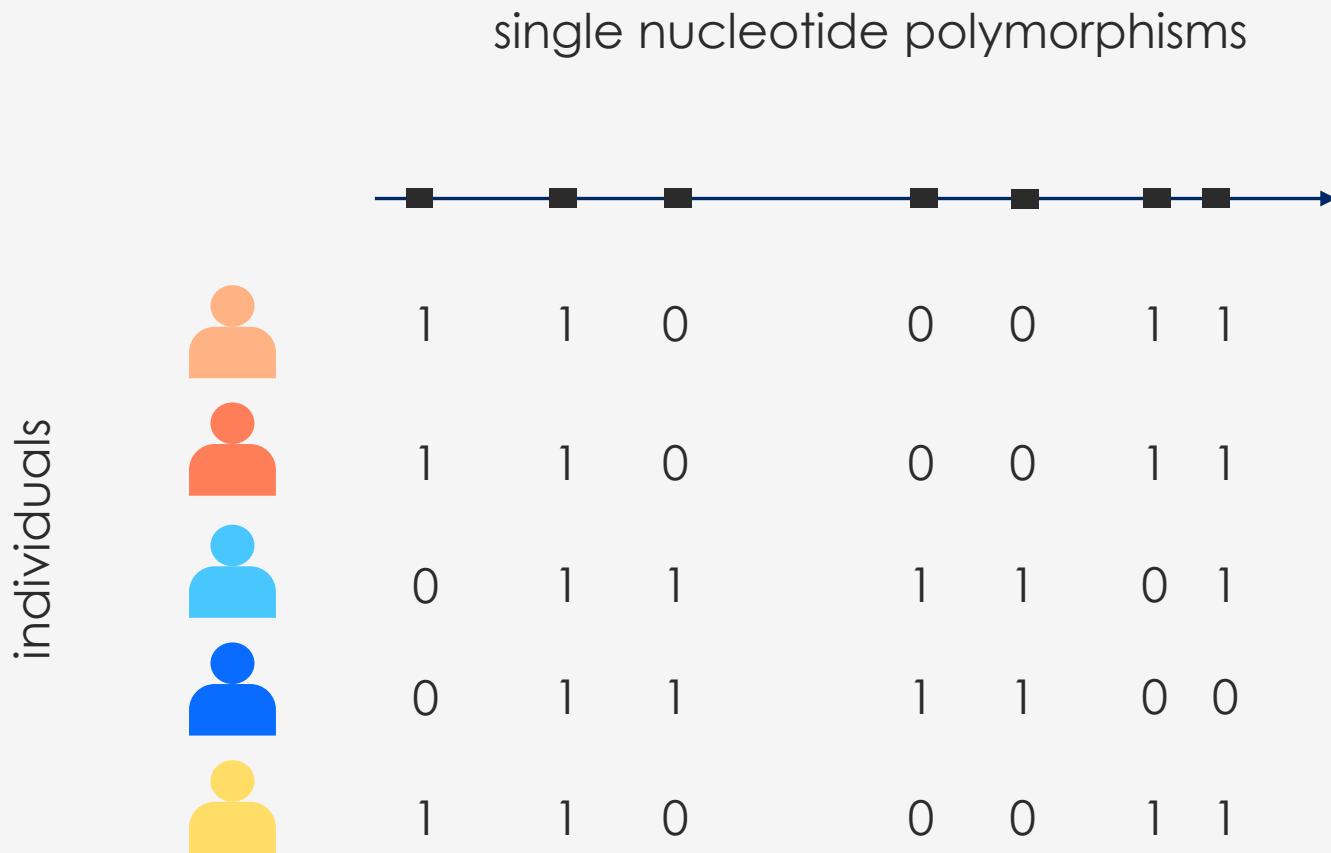
4. Clamp  $W, b$  and get binarized  $W_b$ , discretized  $b_d$

 PyTorch

Use ADAM for optimization



# Exploring human variation.



## 1000 Genomes project

Measuring variations across human populations (~2500 individuals)

Here we look at single nucleotide variations in genes of autosomes (~228k variants)

# Exploring human variation.



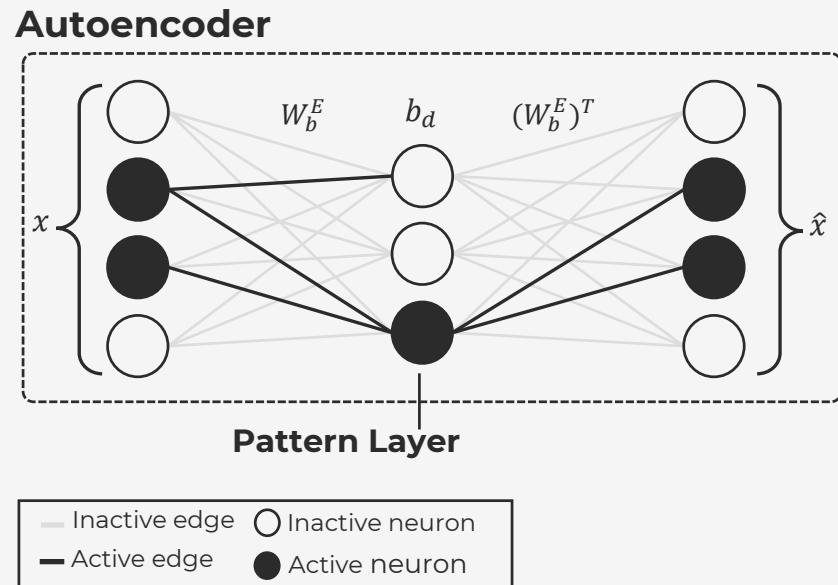
What we find:

- Blocks of consecutive variation
- Patterns of alternate phasing across variants  
(in some individuals pattern 0|1, 0|1, 0|1 in others 1|0, 1|0, 1|0)
- Linking to genes:
  - Pattern of developmental genes together with genes with unknown role
  - Pattern of genes crucial for the ribosomal complex
  - Pattern of variants associated with type 2 diabetes, and unknown variant

# BiNAPS – In a nutshell



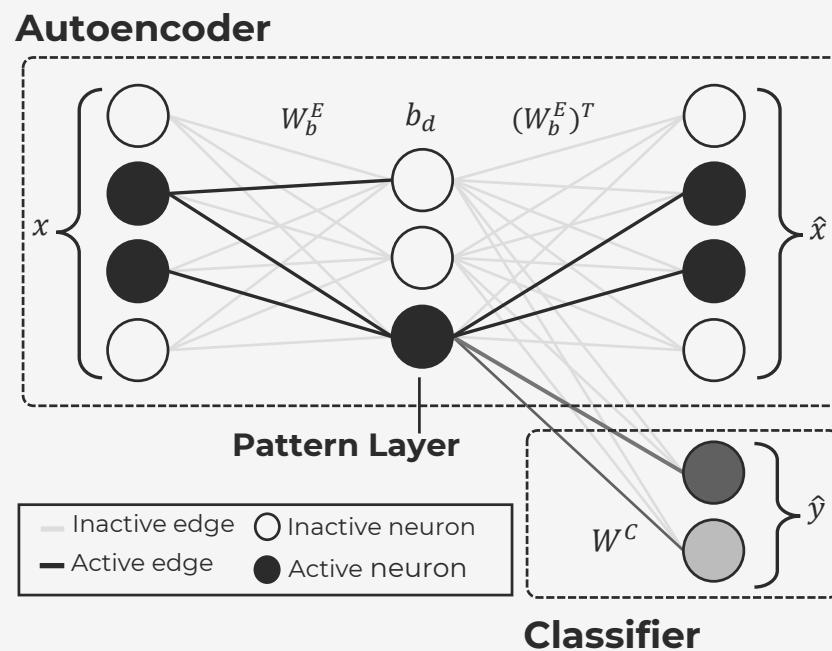
**Main Idea:** Compression → Patterns



# DIFFNAPS – In a nutshell



**Main Idea:** Compression + Classification  $\rightarrow$  Class-specific pattern

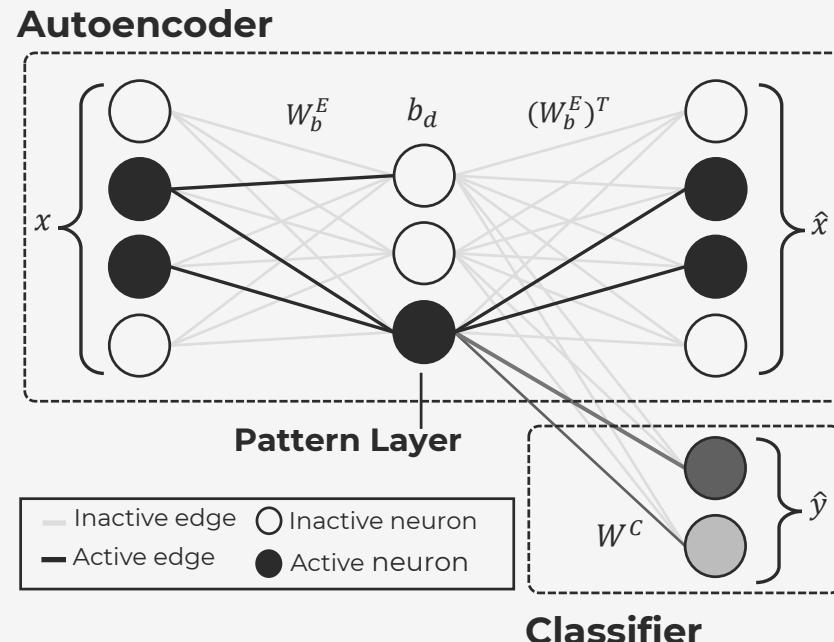


## Forward Pass

1. Binary input  $x$
2. Binarize weight matrix  $W_b^E$
3. Compute hidden activation  $z$   
$$z = \lambda_E(W_b^E x),$$
 where  $\lambda_E$  is a binary activation function.
4. Compute reconstruction  $\hat{x}$   
$$\hat{x} = \lambda_D((W_b^E)^T z)$$
5. Compute classification  $\hat{y}$   
$$\hat{y} = \text{softmax}(W^C z)$$

Jointly optimize reconstruction and classification accuracy

# DIFFNAPS – Extract patterns



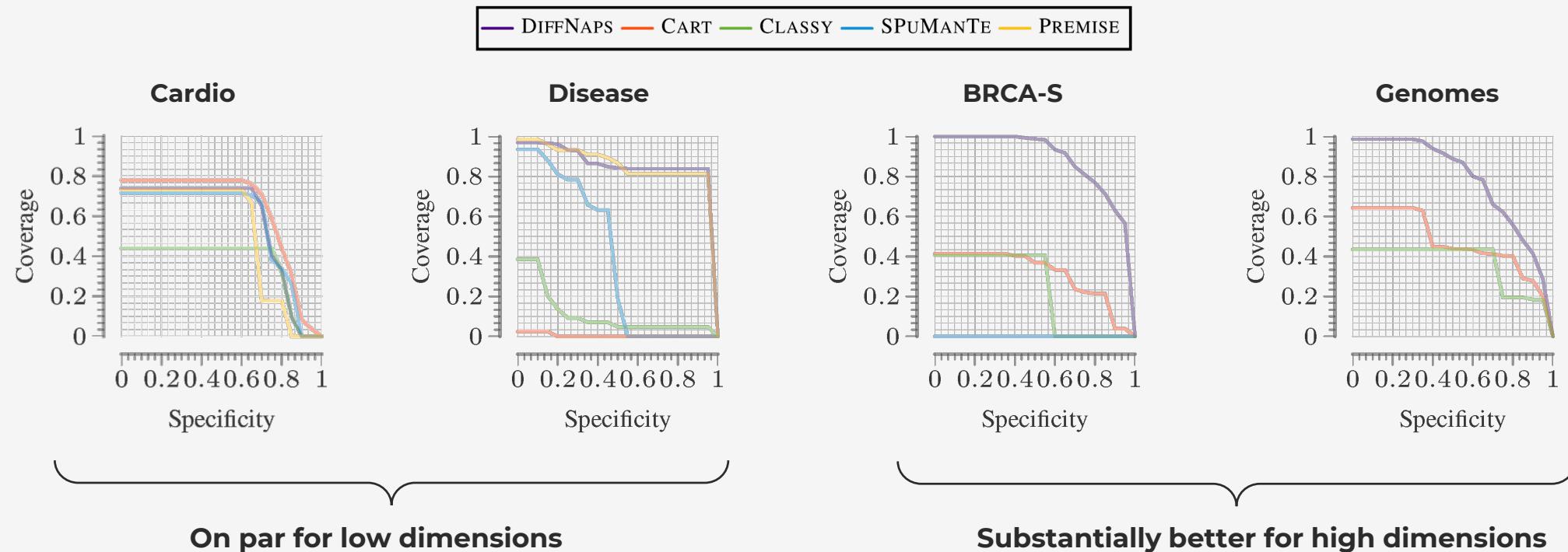
<b>Continuous</b>				<b>Discrete</b>			
$W^E = \begin{bmatrix} 0.02 & \textcolor{red}{0.87} & 0.02 & \textcolor{red}{0.87} \\ \textcolor{red}{0.8} & 0.01 & 0.04 & \textcolor{red}{0.9} \\ 0.0 & \textcolor{red}{0.99} & \textcolor{red}{0.8} & 0.09 \end{bmatrix}$	$\xrightarrow{\mathcal{B}\tau_E(\cdot)}$	$W_b^E = \begin{bmatrix} 0 & 1 & 0 & 1 \\ 1 & 0 & 0 & 1 \\ 0 & 1 & 1 & 0 \end{bmatrix}$	$\xleftarrow{\quad}$				
$W^C = \begin{bmatrix} 0.07 & 0.01 & \textcolor{green}{0.9} \\ 0.02 & \textcolor{green}{0.9} & 0.2 \end{bmatrix}$	$\xrightarrow{\mathcal{B}\tau_C(\cdot)}$	$W_b^C = \begin{bmatrix} 0 & 0 & \textcolor{green}{1} \\ 0 & \textcolor{green}{1} & 0 \end{bmatrix}$	$\xleftarrow{\quad}$				

**Class-specific patterns:**  $P^1 = \{\{2,3\}\}$  resp.  $P^2 = \{\{1,4\}\}$

# DIFFNAPS – Real World Data

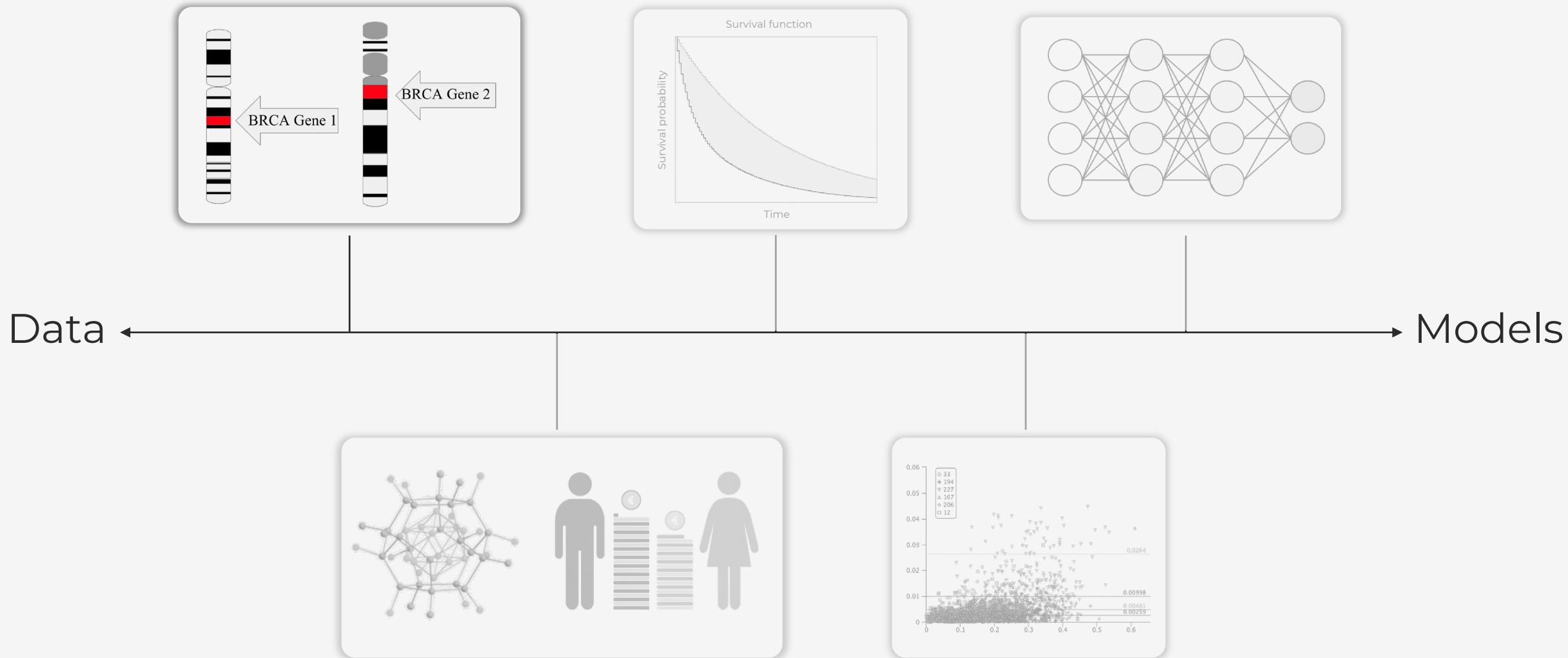


- **Cardio** ( $m = 45$ ): Patient data annotated with heart disease
- **Disease** ( $m = 131$ ): Various symptoms annotated with disease
- **BRCA-S** ( $m = 20k$ ): ncRNA annotated with cancer type → Found biologically **relevant** patterns!
- **Genomes** ( $m = 225k$ ): DNA annotated with ethnicity





# Common questions in ML





# Common questions in ML

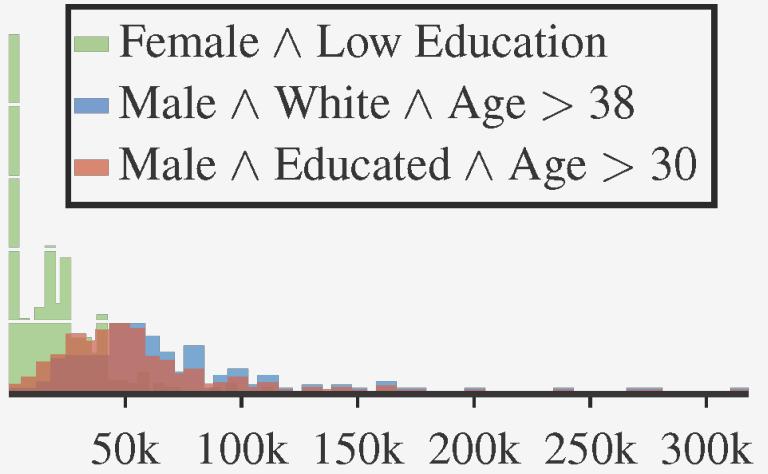




# Discovering exceptional subgroups

## Census Data

Sex	Height	Ethn	Education	Age	Income
♀	168	White	12	72	17k
♂	163	White	11	55	23k
♀	160	White	5	62	1k
♂	188	White	16	38	63k
♀	165	White	9	45	4k
♂	172	White	12	78	71k
♀	180	White	8	74	1k



## Task – Subgroup Discovery:

1. Find **exceptional** subgroups
2. With an **interpretable** description

# Traditional Subgroup Discovery



## Subgroup Discovery

---

1. Pre-discretize features  
→ **Likely to miss subgroups**
2. Strong assumptions on target  
→ **Specialized methods needed**
3. Combinatorial optimization  
→ **Does not scale**

# Modern Subgroup Discovery

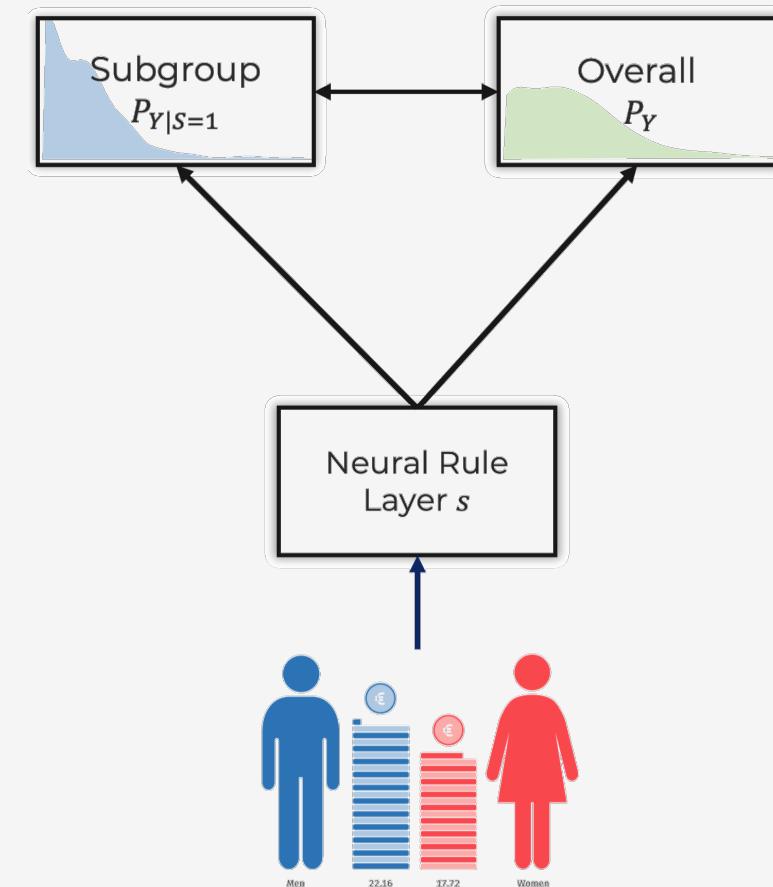


## Subgroup Discovery

1. Pre-discretize features  
→ **Likely to miss subgroups**
2. Strong assumptions on target  
→ **Specialized methods needed**
3. Combinatorial optimization  
→ **Does not scale**

## SyFlow

1. Learn predicates end-to-end  
→ **Accurate Discretization**
2. Use Normalizing Flows (NFs)  
→ **No assumptions**
3. Maximize KL-divergence  
→ **Highly general**
4. Continuous optimization  
→ **Highly scalable**



# SyFlow – Neural Rule Layer I



**Goal:** Find an crisp interpretable description

$$\sigma(x) = \neg Smoker \wedge 44 < Age < 64$$

**Ingredients:**

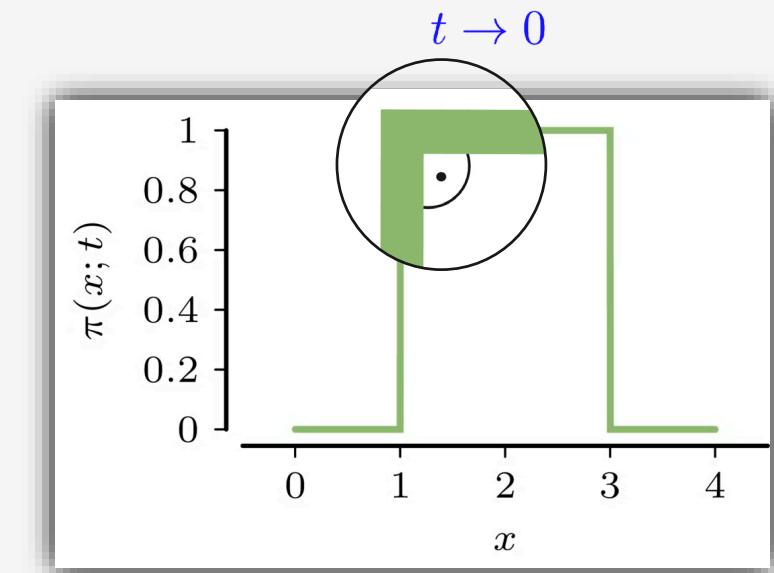
1. Differentiable binning predicate

$$\hat{\pi}(x_i; \alpha_i, \beta_i, t) = \frac{e^{\frac{1}{t}(2x_i - \alpha_i)}}{e^{\frac{1}{t}x_i} + e^{\frac{1}{t}(2x_i - \alpha_i)} + e^{\frac{1}{t}(3x_i - \alpha_i - \beta_i)}}$$

- Differentiable analog of:

$$\pi(x_i; \alpha_i, \beta_i) = \begin{cases} 1 & \text{if } \alpha_i < x_i < \beta_i \\ 0 & \text{otherwise} \end{cases}$$

- Temperature  $t$  controls crispness



**Theorem 1** Given its lower and upper bounds  $\alpha_i, \beta_i \in \mathbb{R}$ , the soft predicate of Eq. (1) applied on  $x \in R$  converges to the crisp predicate that decides whether  $x \in (\alpha, \beta)$ ,

$$\lim_{t \rightarrow 0} \hat{\pi}(x_i; \alpha_i, \beta_i, t) = \begin{cases} 1 & \text{if } \alpha_i < x_i < \beta_i \\ 0.5 & \text{if } x_i = \alpha_i \vee x_i = \beta_i \\ 0 & \text{otherwise} \end{cases} .$$

# SyFlow – Neural Rule Layer II

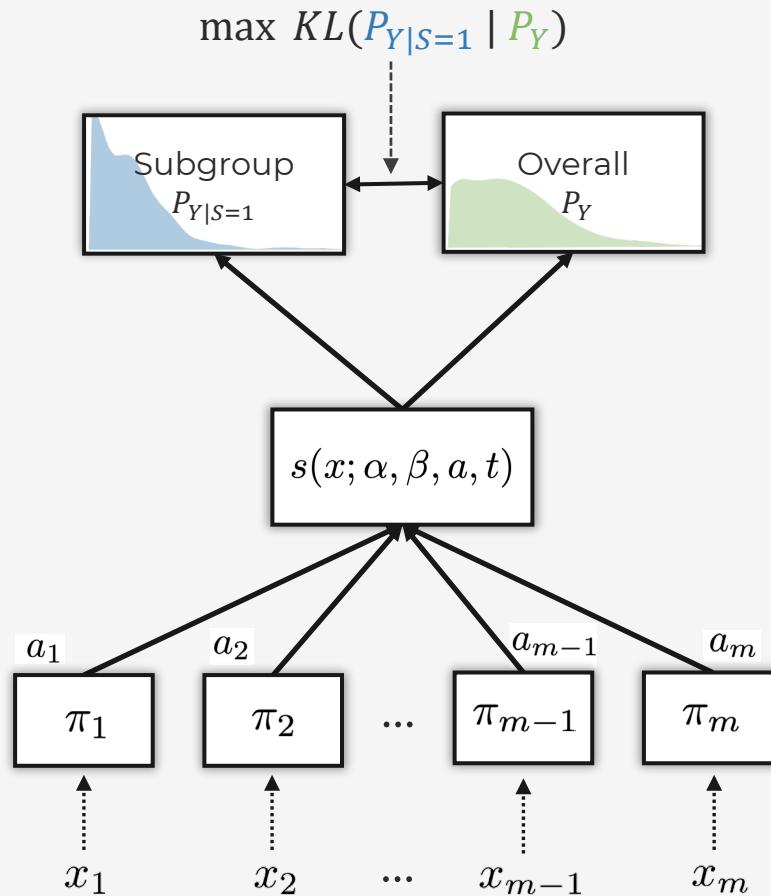


## Ingredients:

1. Differentiable binning predicate
2. Differentiable logical AND
  - Harmonic means behaves like an AND
    1. If one  $\pi_i$  is zero then evaluate to false
    2. If all  $\pi_i$  are one then evaluate to true
  - Implicit feature selection with  $a_i$

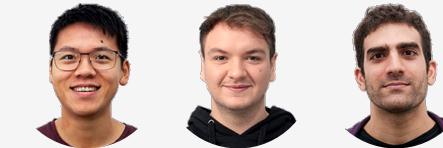
## Optimization

1. Learn the overall distribution  $P_Y$
  2. Learn the subgroup distribution  $P_{Y|S=1}$
  3. Optimize classifier weights and bins
  4. Output: Subgroup
- } Repeat for  $N$  steps

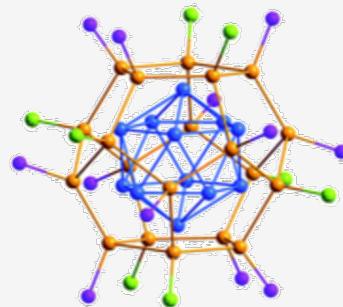


**Fully differentiable!**

# Experiments – Materials Science



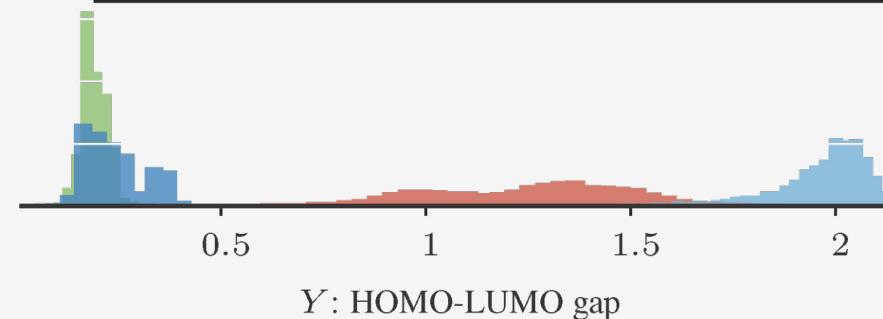
## Gold Nanoclusters



- Atom-based features
- Bond-based features
- Shape-based features

**Target:** HOMO-LUMO gap  
~ stability and conductivity

■ Odd #Atoms  $\wedge$  #Atoms > 8  
■ Odd #Atoms  $\wedge$  % 4-bonds < 0.6  $\wedge$  % 2-bonds < 0.9  
■ Even #Atoms  $\wedge$  3-D Planarity  $\wedge$  Gyration < 1.00  
■ Even #Atoms  $\wedge$  % 0-bonds < 0.01  $\wedge$  2-bonds > 0.43  
 $\wedge$  Gyration < 1.00  $\wedge$  % 1 bond < 0.3



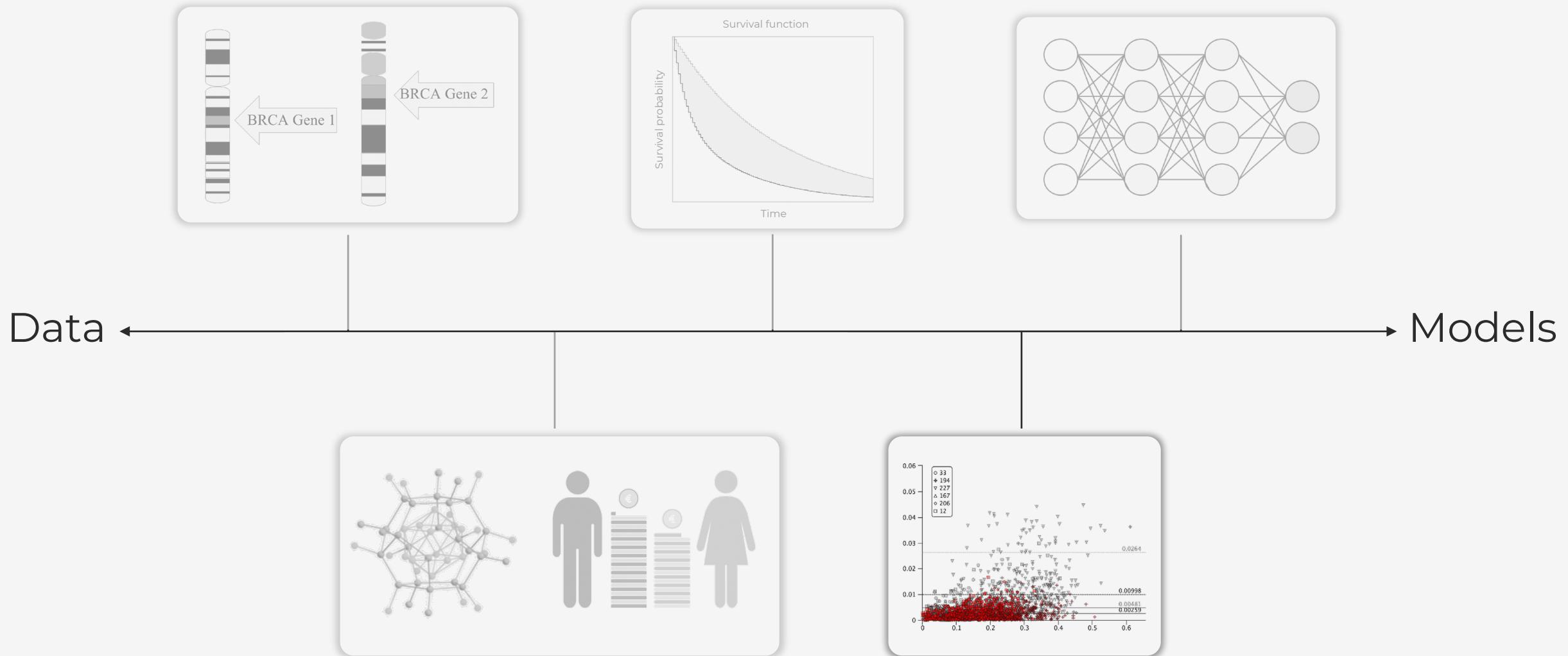


# Common questions in ML





# Common questions in ML



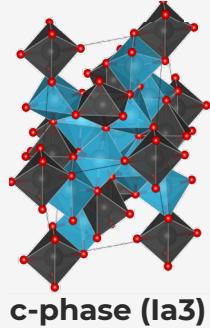
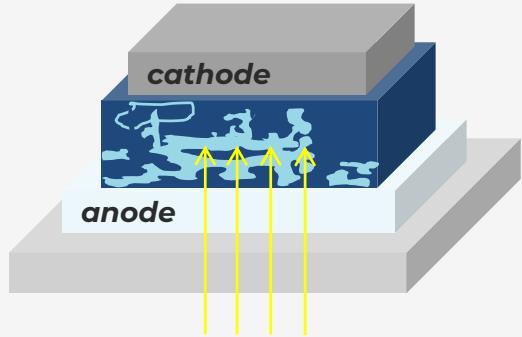
# Model diagnostics



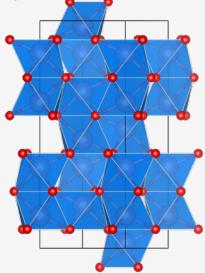
## Population

$$P = \left\{ c : c \in (\text{In}_x \text{Ga}_y \text{Al}_z)_2 \text{O}_3 \right\}$$

## Photovoltaics



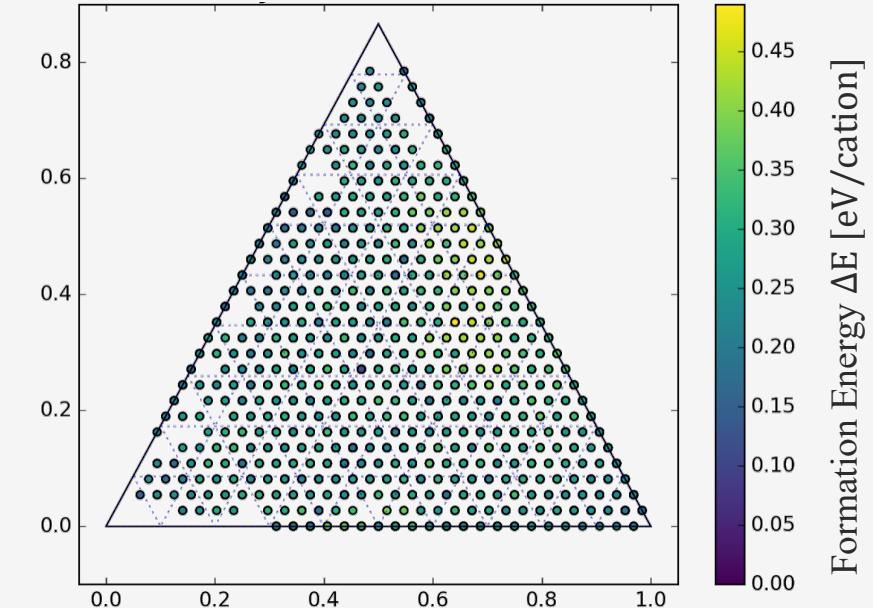
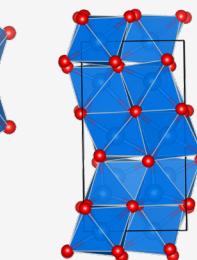
c-phase (Ia3)



$\alpha$ -phase (R3c)



e-phase (Pna21)



# Model diagnostics

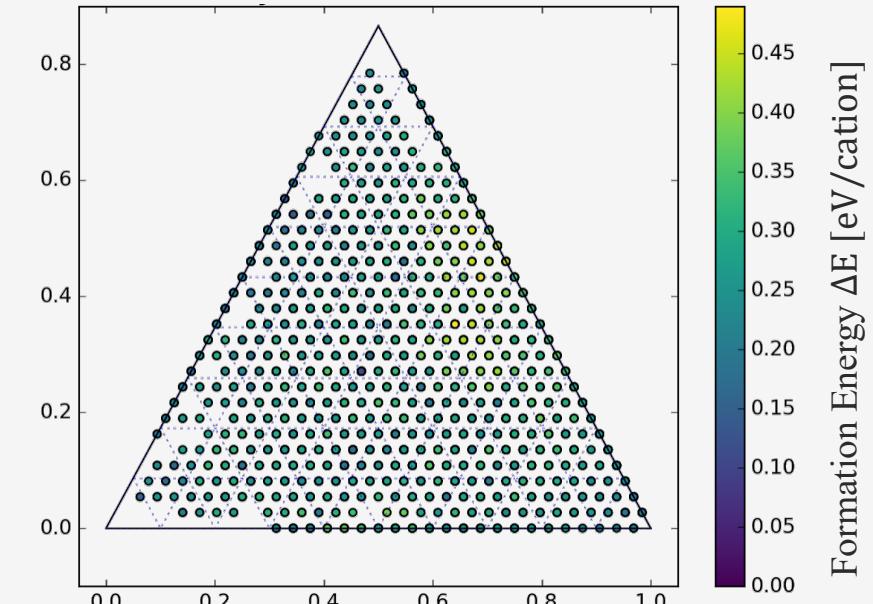
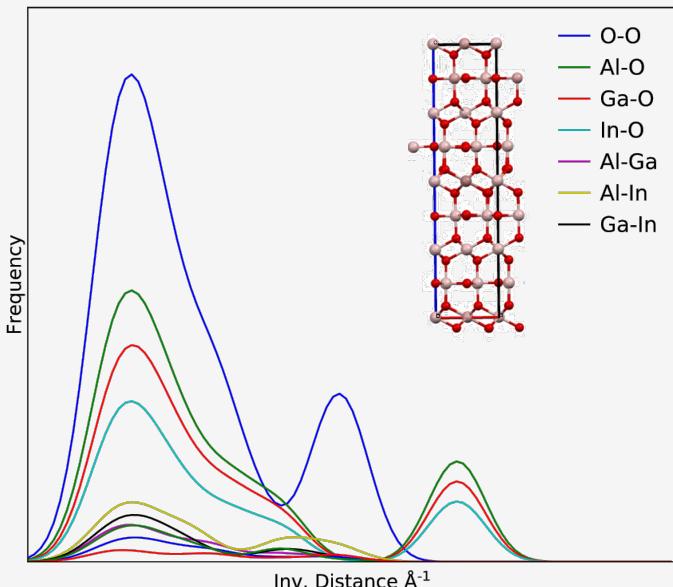


**Population**

$$P = \left\{ c : c \in (\text{In}_x \text{Ga}_y \text{Al}_z)_2 \text{O}_3 \right\}$$

**Target**

$$y = |\Delta E - \tilde{f}_{\Delta E}(c)| \text{ of regression model } \tilde{f}_{\Delta E}$$



Arbitrary model you want  
to use for predicting



# Model diagnostics

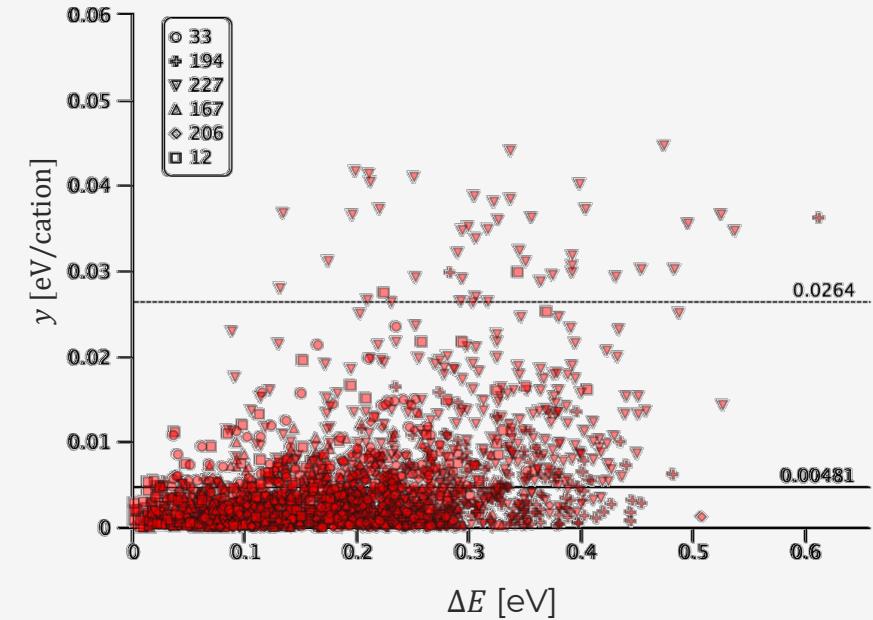
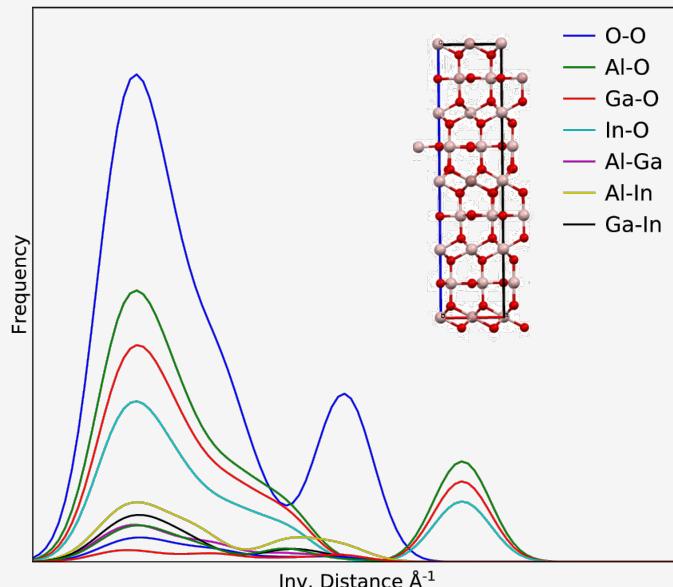


## Population

$$P = \left\{ c : c \in (\text{In}_x \text{Ga}_y \text{Al}_z)_2 \text{O}_3 \right\}$$

## Target

$$y = |\Delta E - \tilde{f}_{\Delta E}(c)| \text{ of regression model } \tilde{f}_{\Delta E}$$



Model has low average error  
but its predictions are  
generally unreliable

# Model diagnostics



## Population

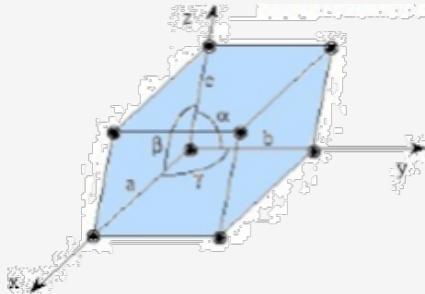
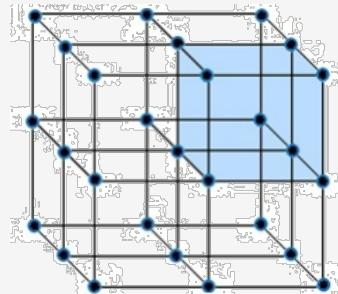
$$P = \left\{ c : c \in (\text{In}_x \text{Ga}_y \text{Al}_z)_2 \text{O}_3 \right\}$$

## Target

$$y = |\Delta E - \tilde{f}_{\Delta E}(c)| \text{ of regression model } \tilde{f}_{\Delta E}$$

## Features

$$x \in \{a, b, c, \alpha, \beta, \gamma, n, \dots\}$$



## Selector

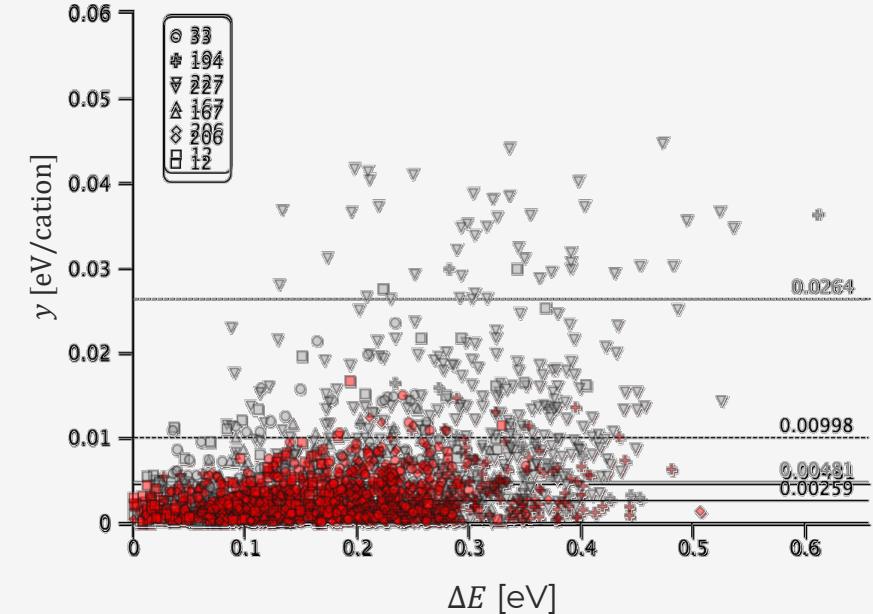
$$\sigma(c) \equiv n(c) \geq 60 \wedge \gamma(c) \leq 95 \wedge \alpha(c) \leq 121$$

## Parameters

$$\text{cvr}(\sigma) = 0.6$$

$$\text{eff}(\sigma) = 0.3$$

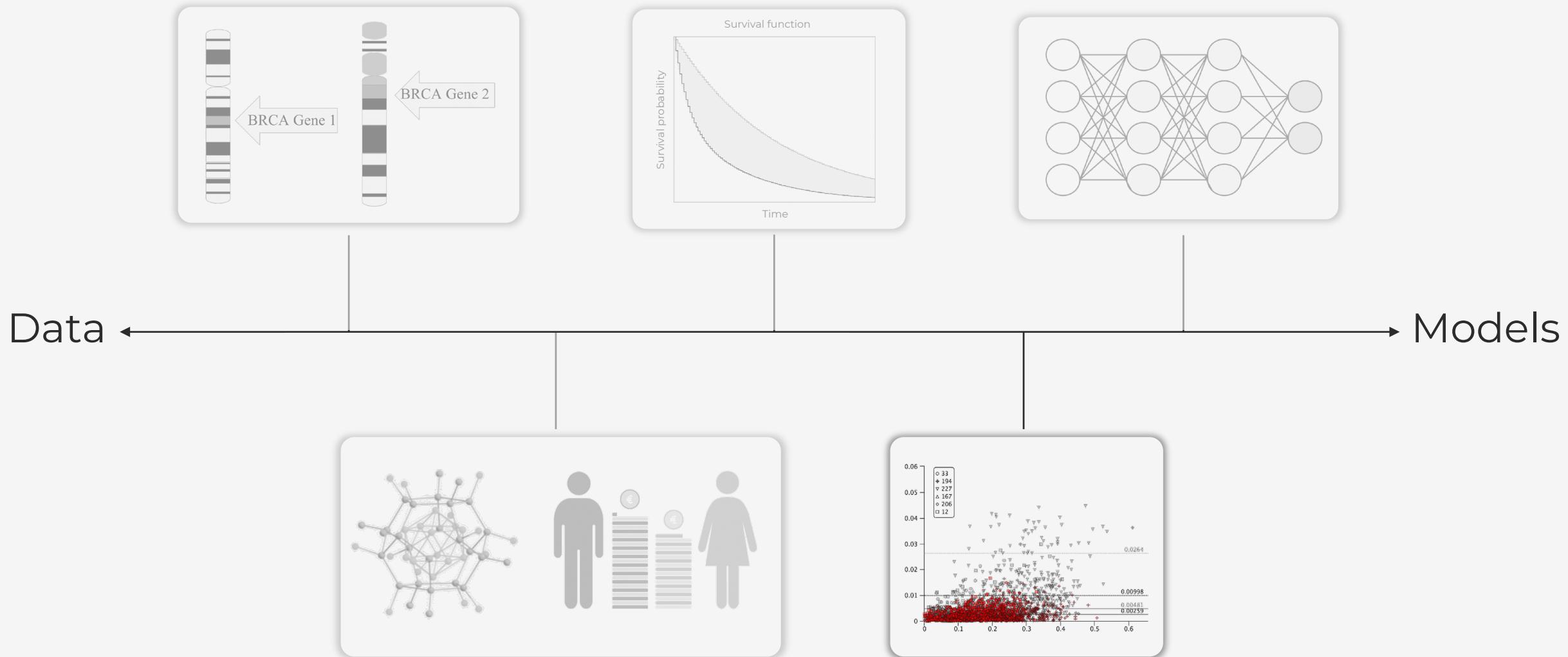
$$[\bar{y}(Q) = 0.003 \pm 0.002, \bar{y}(S) = 0.005 \pm 0.007]$$



Model is **very reliable**  
under these conditions



# Common questions in ML





# Common questions in ML

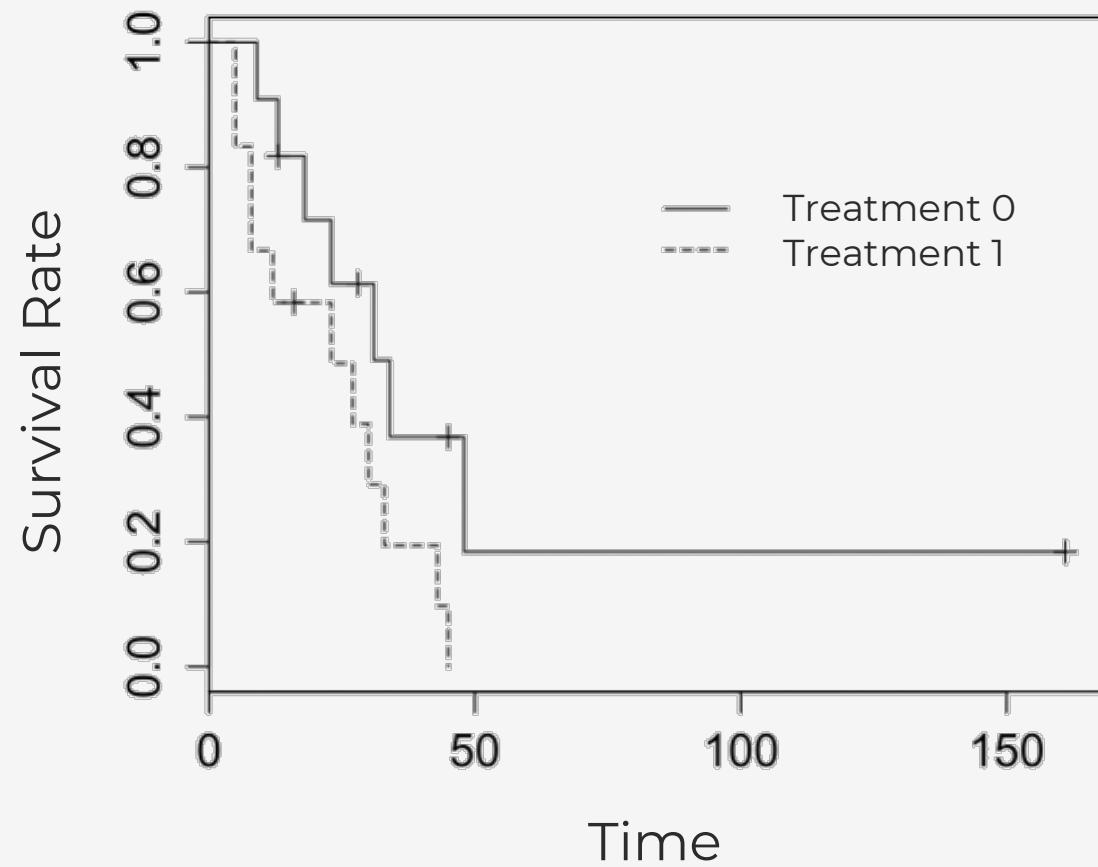


# Time-to-event data

Participant	Treatment	Time of death
12	0	5
4	1	18
13	0	5
17	0	16
14	0	8
3	1	13
15	0	8
2	1	13
1	1	9
16	0	12



Estimated Survival Functions

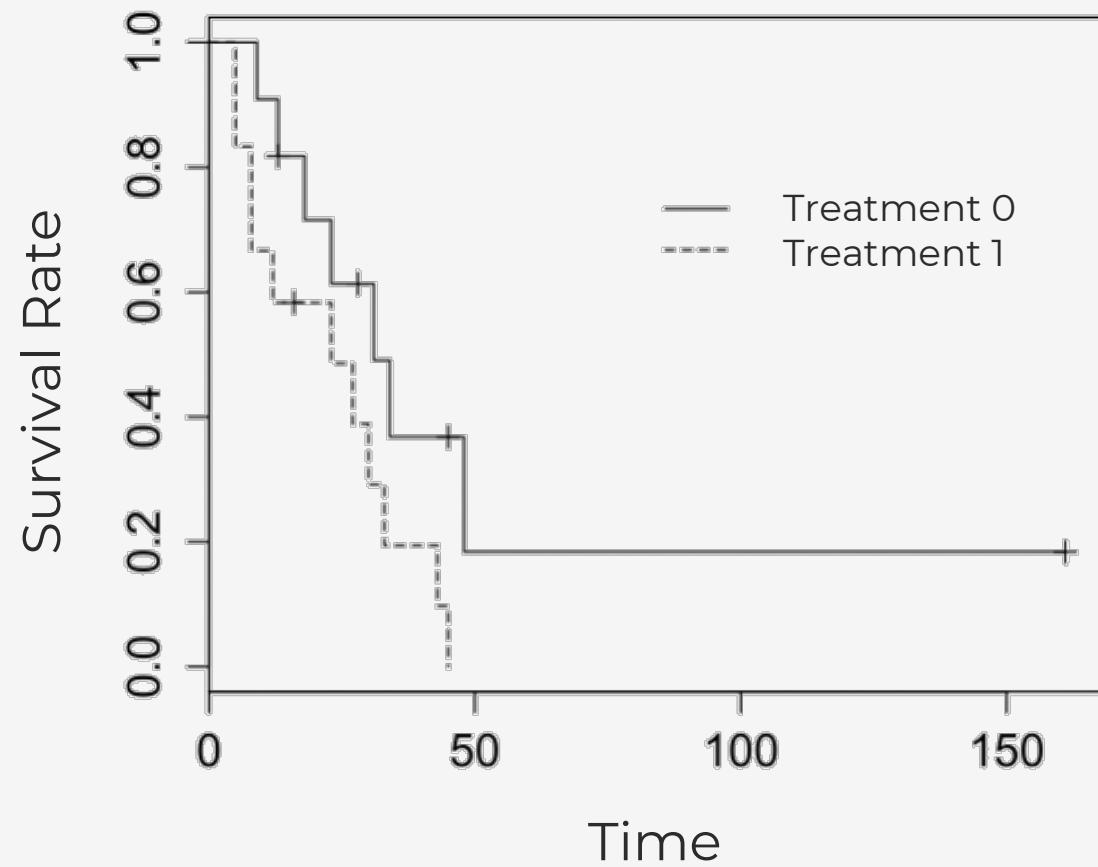


# Time-to-event data

Participant	Treatment	Time of death
12	0	5
4	1	18
13	0	survived
17	0	left study
14	0	8
3	1	13
15	0	8
2	1	13
1	1	9
16	0	12



Estimated Survival Functions



# Survival Subgroups

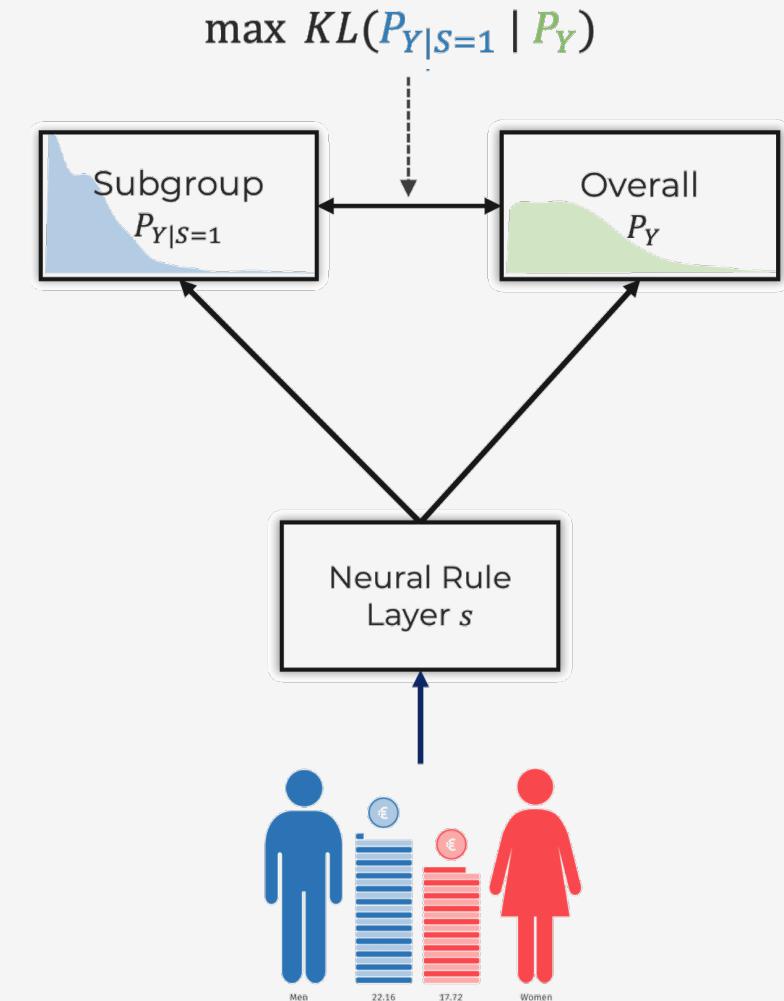


## SYFLOW

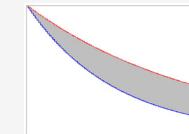
1. Learn predicates end-to-end  
→ **Accurate Discretization**
2. Continuous optimization  
→ **Highly scalable**
3. Use Normalizing Flows (NFs)  
→ **Ignores censoring**

## SYSURV

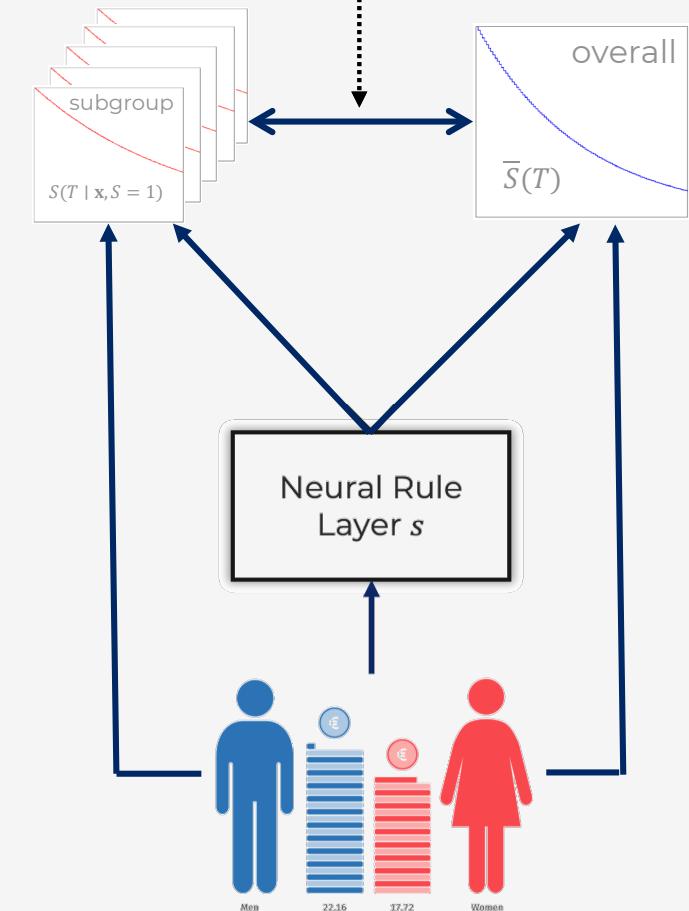
1. Learn predicates end-to-end  
→ **Accurate Discretization**
2. Continuous optimization  
→ **Highly scalable**



# Survival Subgroups



$$\max \max_{\mathbf{X} \in \mathcal{X}} \overline{D}_{\text{KL}}(\overline{P}_{\mathbf{X}} \| P_{\mathbf{Y}})$$



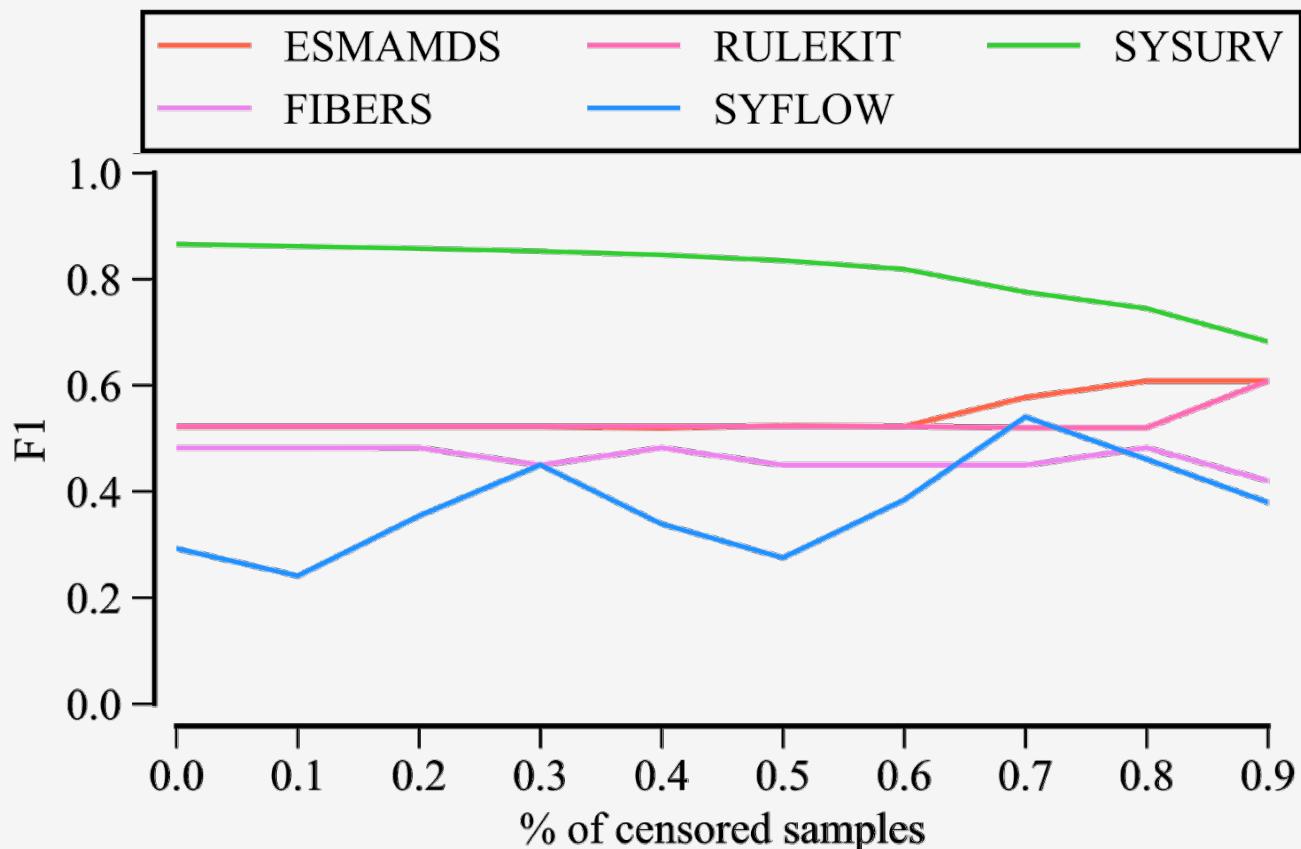
## SYFLOW

1. Learn predicates end-to-end  
→ **Accurate Discretization**
2. Continuous optimization  
→ **Highly scalable**
3. Use Normalizing Flows (NFs)  
→ **Ignores censoring**
4. Uses KL-Divergence  
→ **Ignores survival**

## SYSURV

1. Learn predicates end-to-end  
→ **Accurate Discretization**
2. Continuous optimization  
→ **Highly scalable**
3. Use Random Survival Forests  
→ **No assumptions**
4. Maximize curve distance  
→ **No assumptions**

# Experiments — Synthetic data



# Experiments — Real world

	Our objective					KL-Divergence					Logrank statistic					Mean-shift					Adjusted mean-shift				
	SYSURV	SYFLOW	ESMAMDS	FIBERS	RULEKIT	SYSURV	SYFLOW	ESMAMDS	FIBERS	RULEKIT	SYSURV	SYFLOW	ESMAMDS	FIBERS	RULEKIT	SYSURV	SYFLOW	ESMAMDS	FIBERS	RULEKIT	SYSURV	SYFLOW	ESMAMDS	FIBERS	RULEKIT
UnempDur	<b>6.37</b>	3.20	3.79	3.79	3.79	0.08	<b>0.15</b>	0.13	0.13	0.13	19.26	125.99	125.00	125.00	<b>226.52</b>	-2.01	-1.74	<b>2.14</b>	<b>2.14</b>	<b>2.14</b>	-0.02	-0.00	0.00	0.00	0.00
nwtco	<b>1011.64</b>	205.02	680.98	208.59	680.98	<b>0.14</b>	0.02	0.06	0.00	0.06	<b>374.22</b>	0.21	31.36	0.00	281.27	<b>1161.97</b>	872.08	688.24	-0.56	688.24	6.92	0.40	1.50	-0.00	1.50
rott2	<b>559.22</b>	419.19	487.58	487.58	326.16	-	-	-	-	-	<b>266.95</b>	3.11	71.65	169.35	93.04	31.47	<b>-92.00</b>	25.23	25.23	10.03	0.10	-11.50	0.04	0.04	0.01
rdata	<b>259.24</b>	153.69	213.83	158.52	213.83	<b>0.08</b>	0.04	0.03	0.01	0.03	<b>113.38</b>	56.39	14.64	43.32	107.75	<b>1116.58</b>	-700.91	933.49	705.50	933.49	7.98	-1.72	3.60	1.62	3.60
Aids2	<b>66.64</b>	31.99	54.01	34.69	60.52	-	-	-	-	-	0.21	0.00	<b>4.27</b>	0.01	1.72	22.68	-1.45	<b>61.36</b>	-0.70	61.21	0.07	-0.00	0.10	-0.00	0.17
Dialysis	4.52	4.69	4.71	4.59	<b>4.74</b>	0.02	<b>0.24</b>	0.01	0.02	0.04	6.36	33.76	70.86	<b>305.94</b>	68.38	0.64	<b>-19.60</b>	3.05	4.47	5.74	0.00	-0.19	0.00	0.00	0.02
TRACE	261.76	<b>486.65</b>	332.10	261.58	272.52	0.00	-	0.02	<b>4.98</b>	0.01	0.00	17.74	93.24	0.00	<b>94.62</b>	0.00	<b>4.53</b>	1.73	-0.00	0.61	0.00	4.53	0.00	-0.00	0.00
support2	<b>132.83</b>	68.72	104.90	54.79	81.38	<b>2.71</b>	0.57	1.36	0.22	0.21	<b>251.00</b>	39.00	46.83	11.93	59.87	<b>405.72</b>	-193.91	405.52	-71.03	227.60	12.29	-1.41	5.13	-0.16	2.35
dataDIVAT2	<b>286.30</b>	99.92	172.47	141.30	141.30	0.16	<b>3.91</b>	1.41	2.19	2.19	<b>33.55</b>	15.90	10.94	31.60	31.60	<b>3.18</b>	-1.06	2.06	1.56	1.56	0.05	-0.00	0.01	0.00	0.00
prostateSurvival	<b>20.11</b>	9.20	14.13	5.38	7.51	-	-	-	-	-	<b>425.72</b>	103.66	82.50	82.50	64.74	<b>6.11</b>	5.11	3.37	-1.02	3.67	0.00	0.00	0.00	-0.00	0.00
actg	<b>34.58</b>	13.62	18.82	11.33	24.22	<b>0.09</b>	0.03	0.03	0.00	0.06	<b>59.97</b>	2.40	12.26	0.02	1.05	<b>54.13</b>	18.50	5.98	-0.37	27.30	1.42	0.07	0.03	-0.00	0.22
scania	<b>172.47</b>	101.59	147.42	84.74	147.42	<b>0.03</b>	0.00	<b>0.03</b>	0.01	<b>0.03</b>	4.65	1.91	13.66	9.93	<b>14.66</b>	<b>6.51</b>	-2.86	-2.44	-0.07	-2.44	0.05	-0.00	-0.01	-0.00	-0.01
grace	<b>38.56</b>	32.04	21.01	13.08	16.51	0.14	0.06	<b>7.53</b>	0.02	0.00	59.30	<b>70.71</b>	28.68	13.81	33.33	<b>64.03</b>	51.79	21.77	-11.14	16.00	1.94	0.54	0.08	-0.02	0.04
Avg. rank	<b>1.54</b>	3.85	2.46	4.23	2.92	<b>2.50</b>	2.67	2.80	3.75	3.05	<b>2.38</b>	3.46	3.08	3.65	2.42	<b>2.08</b>	2.85	2.85	4.23	3.00					

(Datasets were chosen to have 1000 samples or more and around 10 features.)

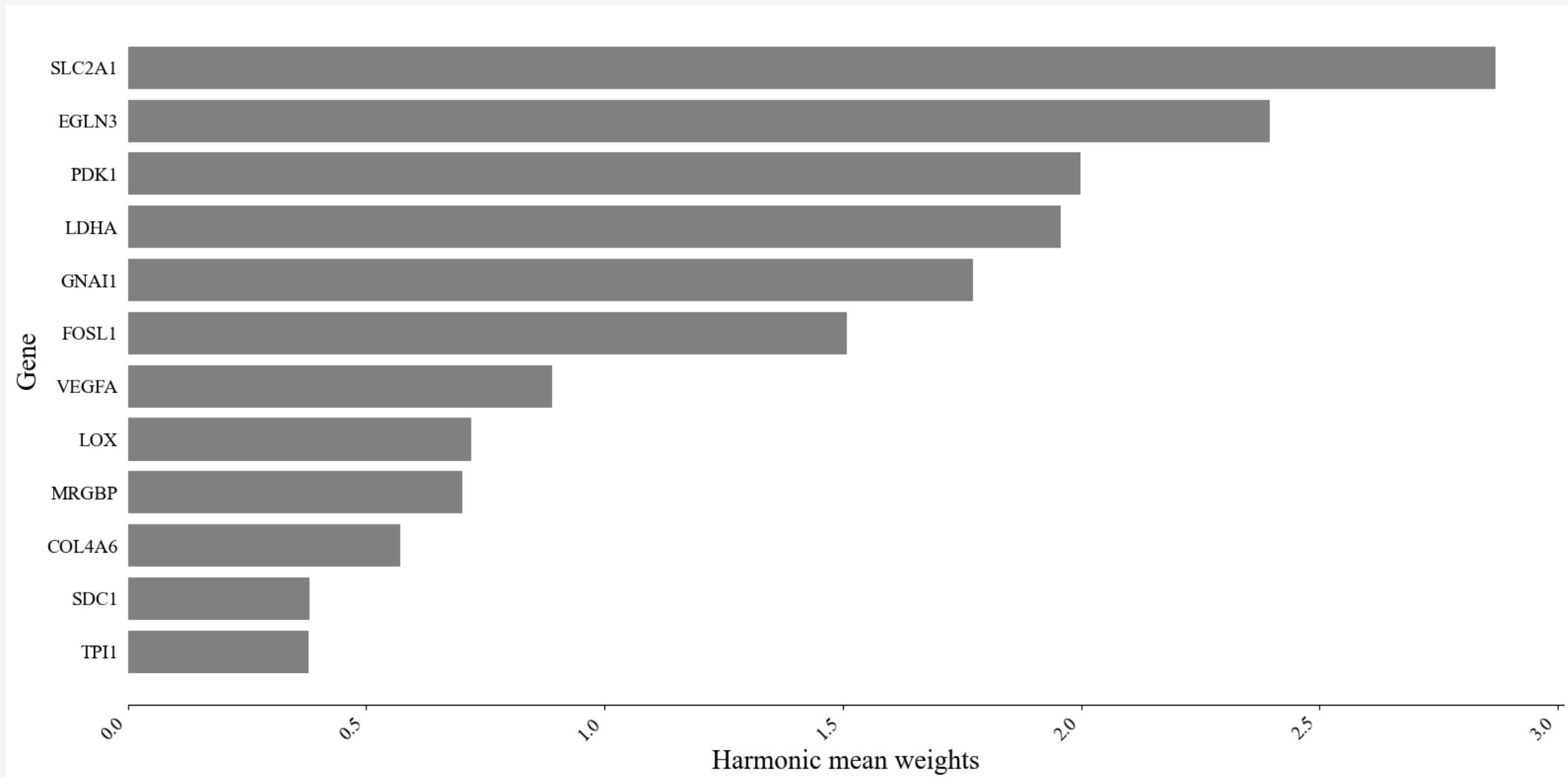
# Experiments — Real world

Avg. Rank	SYSURV	SYFLOW	ESMAMDS	FIBERS	RULEKIT
Our objective	<b>1.54</b>	3.85	2.46	4.23	2.92
KL-Divergence	<b>2.50</b>	2.67	2.80	3.75	3.05
Logrank statistic	<b>2.38</b>	3.46	3.08	3.65	2.42
Mean-shift	<b>2.08</b>	2.85	2.85	4.23	3.00



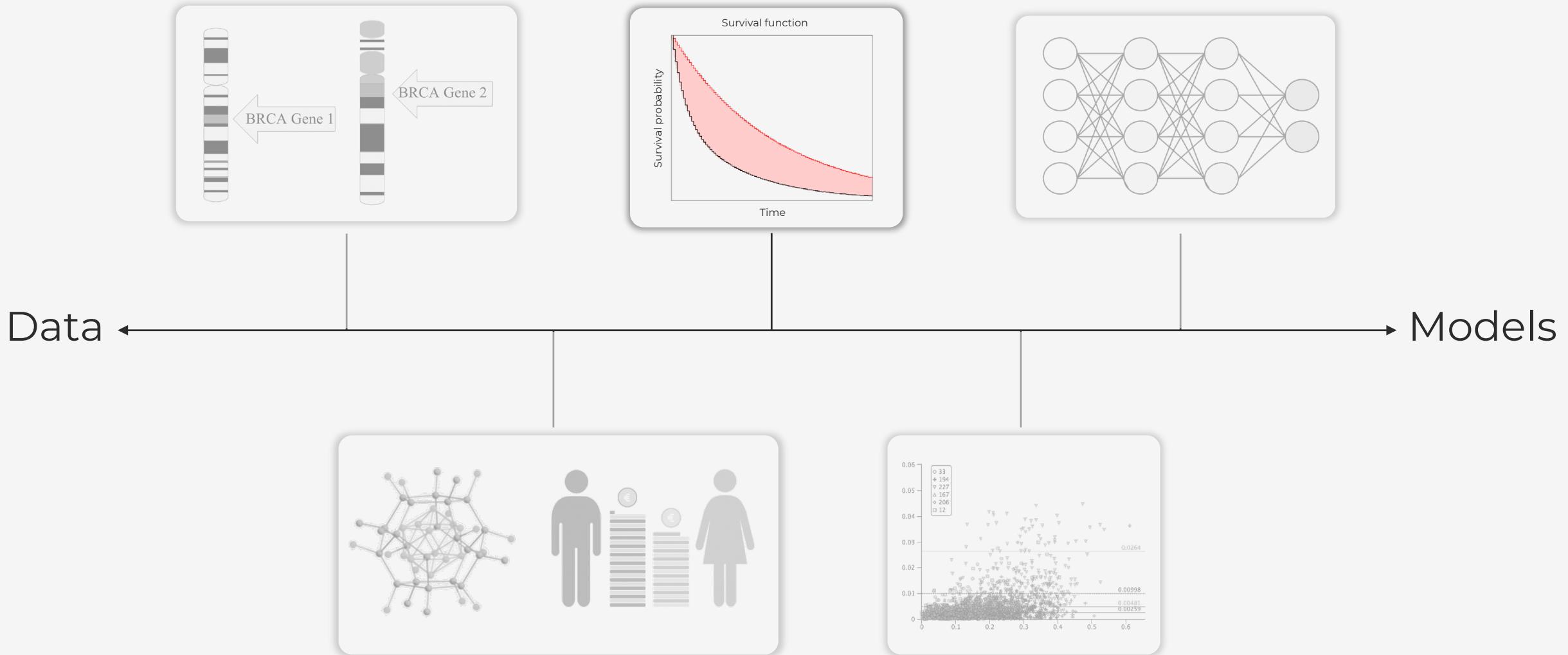
# Case Study: Neck Carcioma

PCR - first rule (pruned) -> All associated with Hypoxia



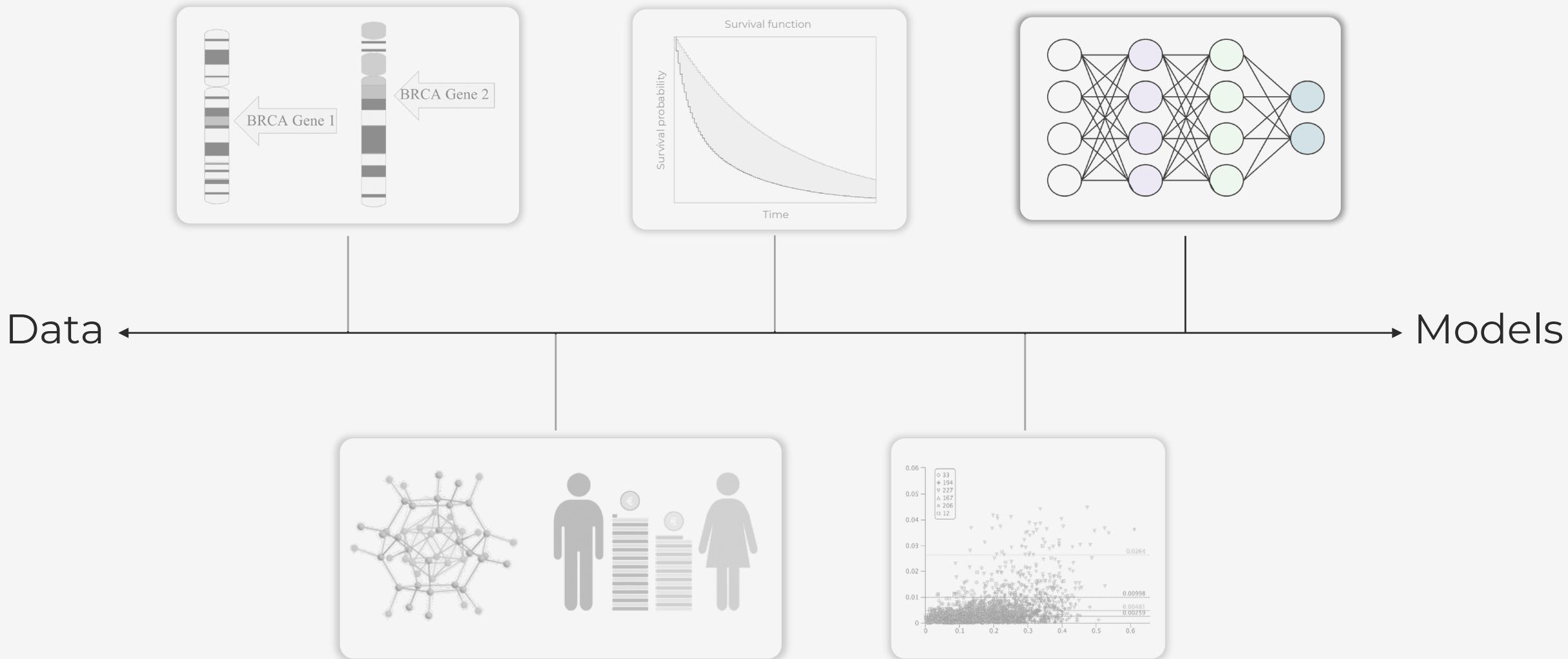


# Common questions in ML



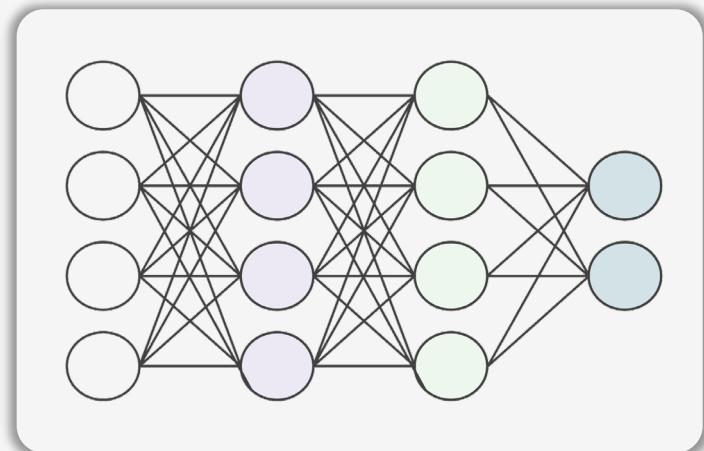


# Common questions in ML

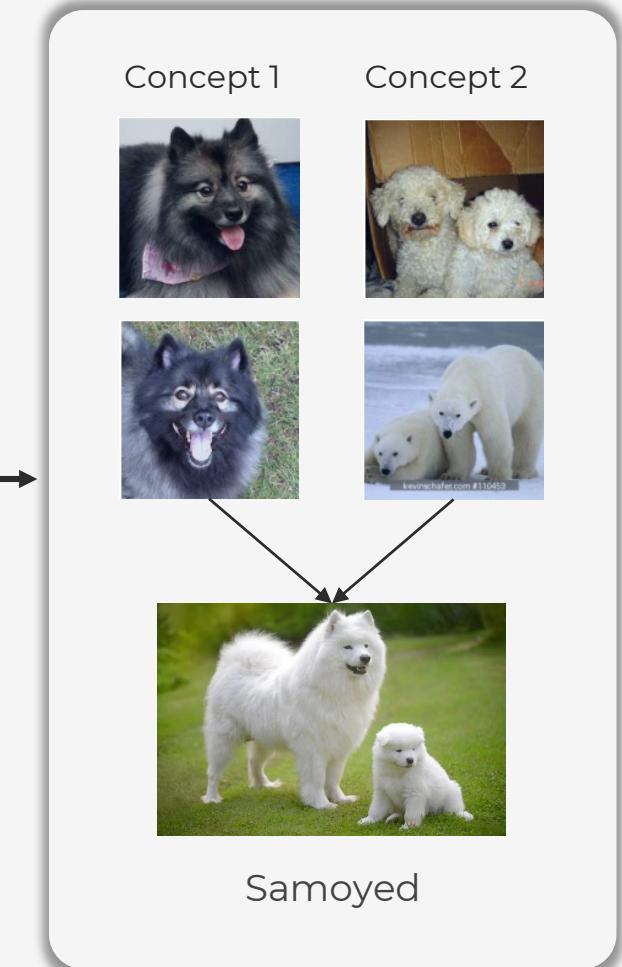




# Understanding reasoning



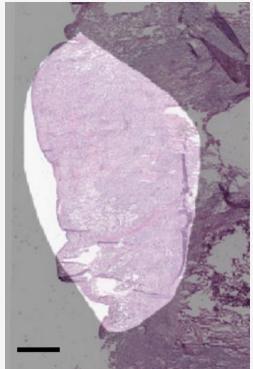
Extract interpretable concepts





# Why is it important?

## Gaining new scientific insights



“... 45 out of 54 of the TCGA images **misclassified by at least one of the pathologists** were **assigned to the correct** cancer type **by the algorithm**”.



COVID-Net

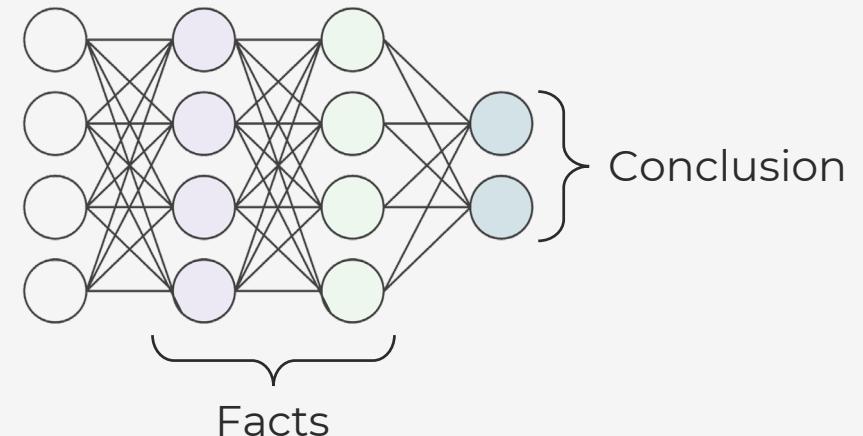
Gain insights into **critical factors** associated with **COVID** cases

## What is **reasoning**?

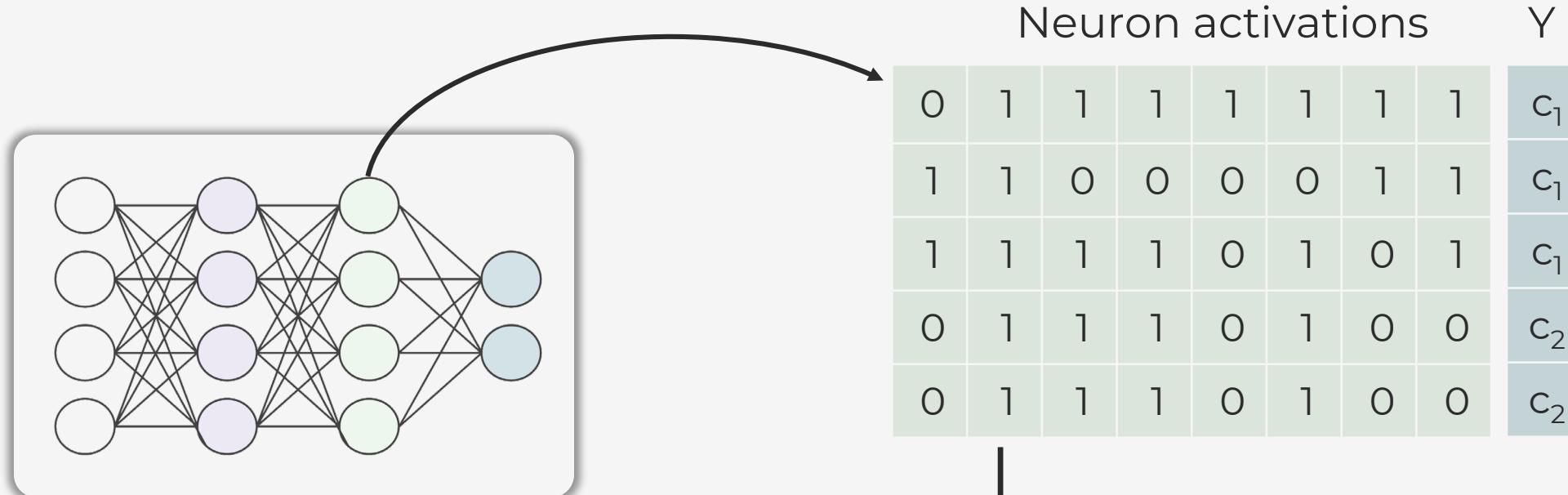
„Reasoning is the process by which you reach a **conclusion** after thinking about all the **facts**.“

— Collins dictionary

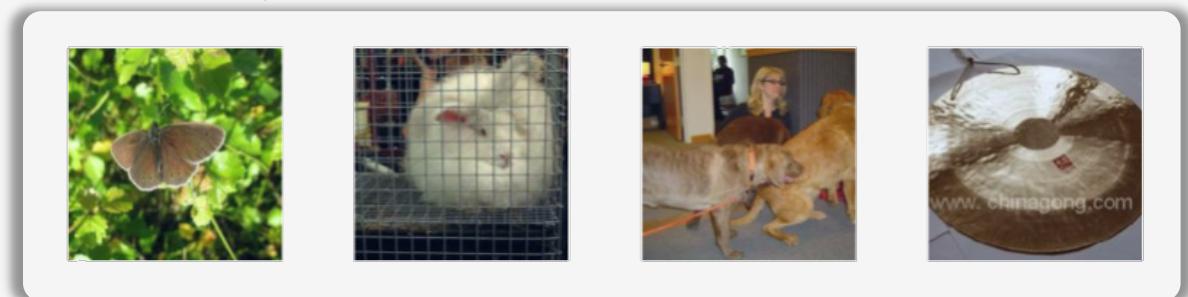
i.e. a relation between **facts** and a **conclusion**



# Extracting Features



**Neurons are polysemantic**





# Sneak peek – Information Flow

How does information flow and how is it combined?

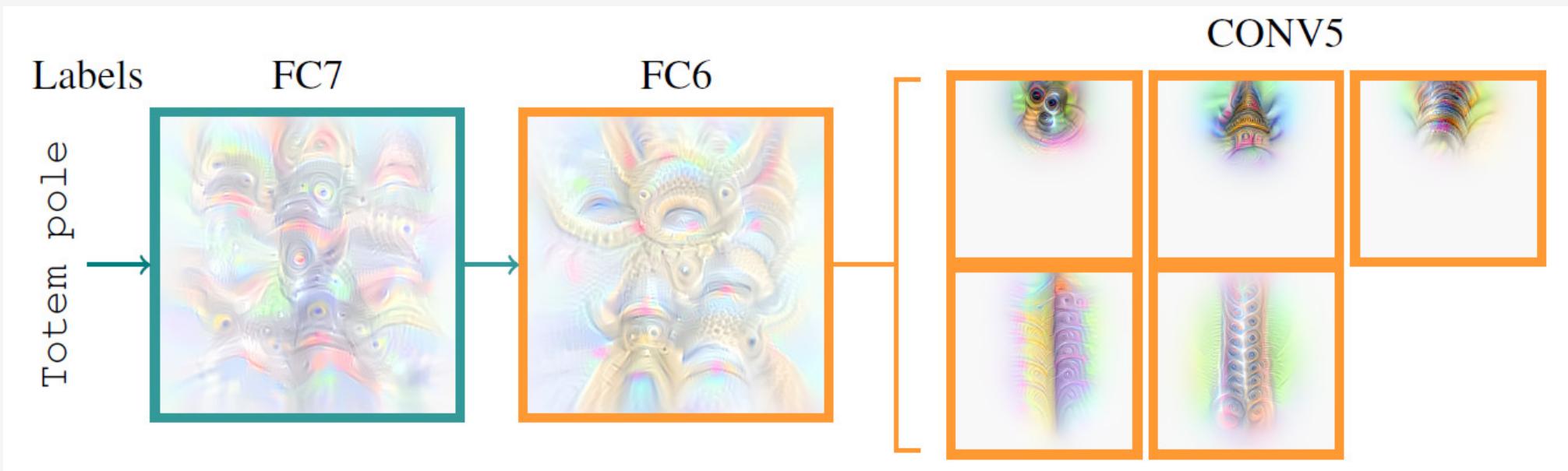
Rules:  $X \rightarrow Y$ ,

$X \in Labels, Y \in FC7$

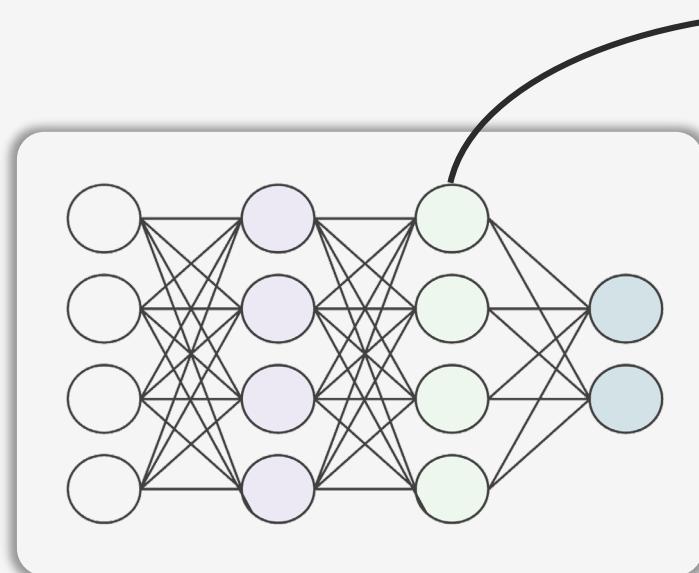
$Y \rightarrow Z$ ,

$Z \in FC6$

$Z \rightarrow Q, Q \in CONV5$



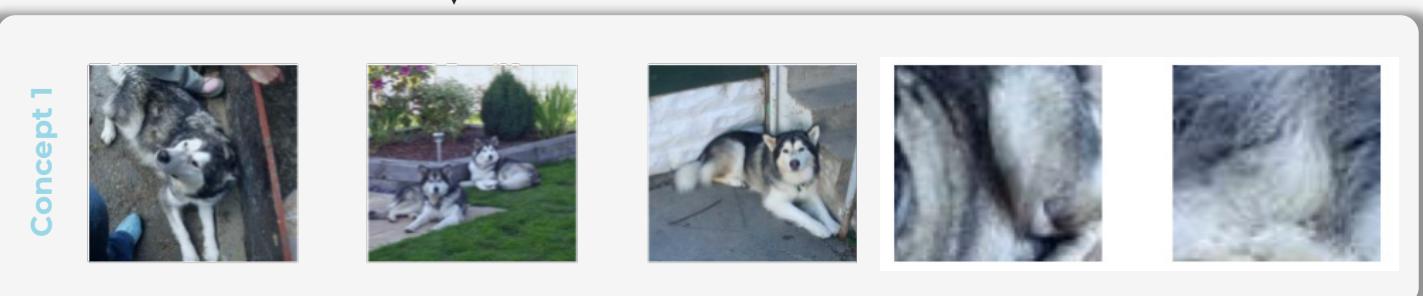
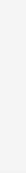
# Extracting Features Concepts



Concept  
Feature activations

Y

0	1	1	1	1	1	1	1	1	c <sub>1</sub>
1	1	0	0	0	0	0	1	1	c <sub>1</sub>
1	1	1	1	0	1	0	1	0	c <sub>1</sub>
0	1	1	1	0	1	0	0	0	c <sub>2</sub>
0	1	1	1	0	1	0	0	0	c <sub>2</sub>



# Connecting concepts

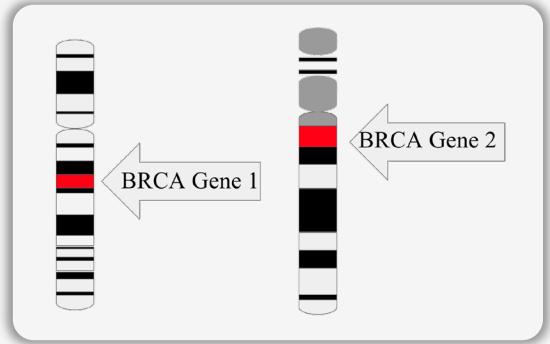




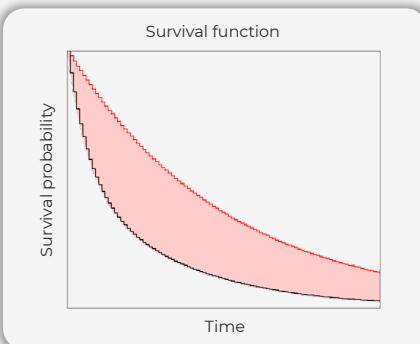
# Conclusion



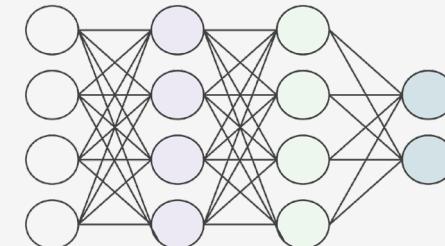
DIFFNAPS



SYSURV

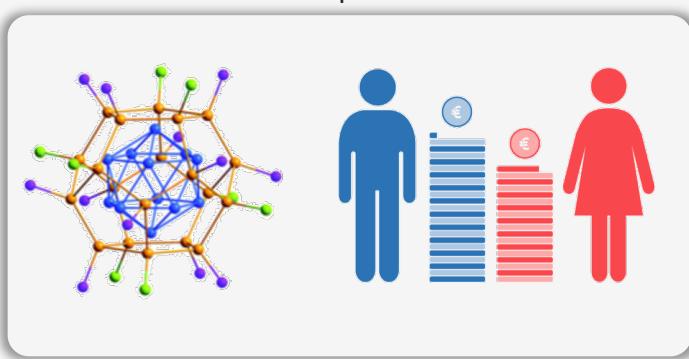


EXPLAINNAPS

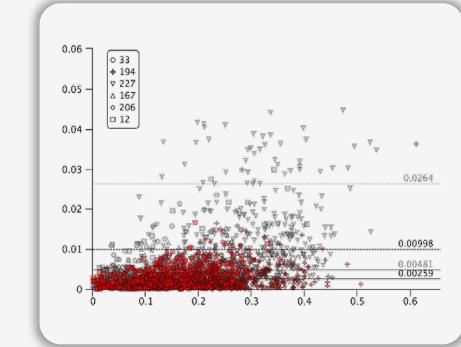


Data

Models



SYFLOW



DOMAINS OF APPLICABILITY