

# ANALYZING BINARY DATA: IS YOUR MODEL TRULY PAIRWISE?

Clélia de Mulatier<sup>1</sup>, Paolo Mazza<sup>2</sup>, and Matteo Marsili<sup>3</sup>

<sup>1</sup>University of Pennsylvania, Department of Physics and Astronomy, Philadelphia, USA

<sup>2</sup>Scuola Internazionale Superiore di Studi Avanzati (SISSA), Trieste, Italy

<sup>3</sup>The Abdus Salam International Centre for Theoretical Physics (ICTP), Trieste, Italy

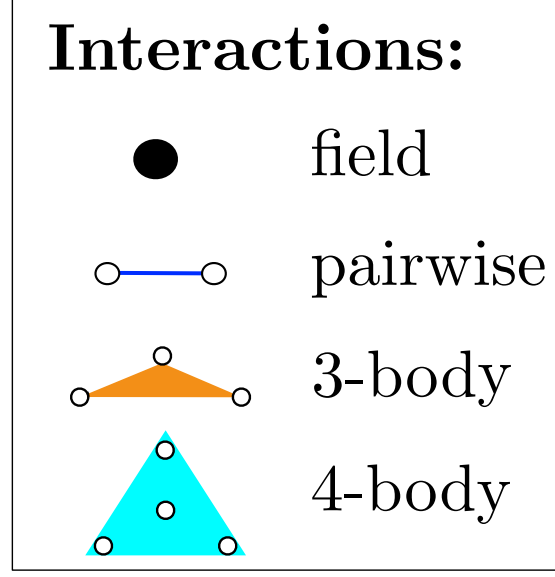
## INTRODUCTION

A fundamental issue in data analysis is to **find the model** that best captures the patterns hidden within the data, despite the random errors that effect them. This model should be **complex enough** to be able to fit the data, but **simple enough** to capture only the relevant patterns of the data.

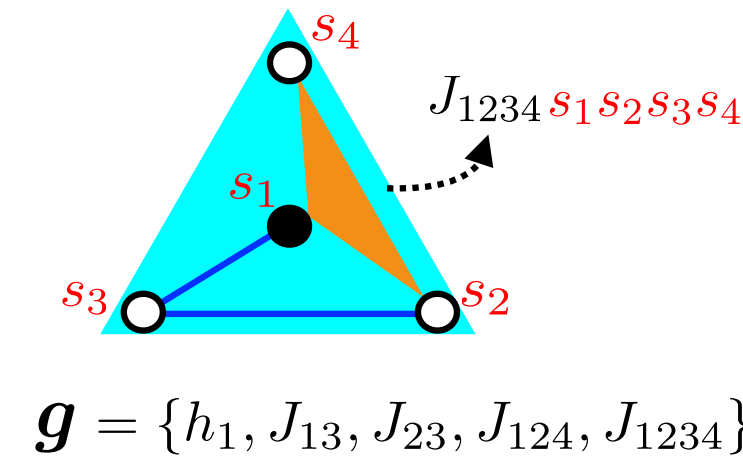
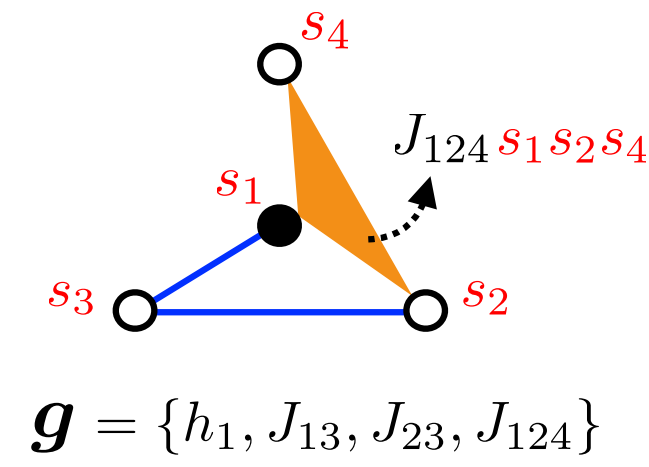
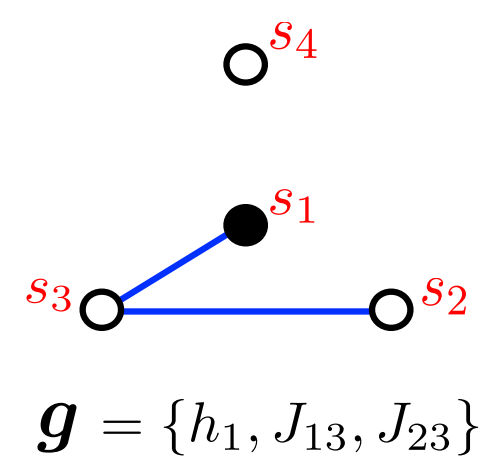
In the context of **binary data**, pairwise spin models – Ising models – have been widely used to address this question. But, are pairwise models the best for your data? Are the relevant patterns of your data truly pairwise?

## PAIRWISE MODELS ARE NOT NECESSARILY simpler

Q? Which is the best spin model for your data?



Ex.



$$P(\mathbf{s} | g, \mathcal{M}) = \frac{1}{Z_{\mathcal{M}}(g)} \exp \left( \sum_{\mu \in \mathcal{M}} g_{\mu} \phi_{\mu}(\mathbf{s}) \right)$$

parameter  $\uparrow$   
Spin operator

How? Two approaches. Find the model  $\mathcal{M}$ :  
 ○ with the highest evidence,  $P(\hat{\mathbf{s}} | \mathcal{M}) \rightarrow$  Bayesian approach;  
 ○ which achieves the shortest description of the data,  $L(\hat{\mathbf{s}} | \mathcal{M}) \rightarrow$  MDL principle.

For large datasets,

the two criteria are identical [1,2]:

$$\log P(\hat{\mathbf{s}} | \mathcal{M}) = \log P(\hat{\mathbf{s}} | \mathcal{M}, g^*) - \left[ \frac{K}{2} \log \left( \frac{N}{2\pi} \right) + c_{\mathcal{M}} \right] + O\left(\frac{1}{N}\right) = -L(\hat{\mathbf{s}} | \mathcal{M})$$

$\uparrow \propto N$        $\uparrow$  Due to       $\uparrow$  Due to **Geometry**  
**Maximum Log-Likelihood**      **Number of Parameters K**

Problem? difficult task due to the huge number of models and the difficulty to evaluate  $P(\hat{\mathbf{s}} | \mathcal{M})$  or  $L(\hat{\mathbf{s}} | \mathcal{M})$ , even with the expansion.

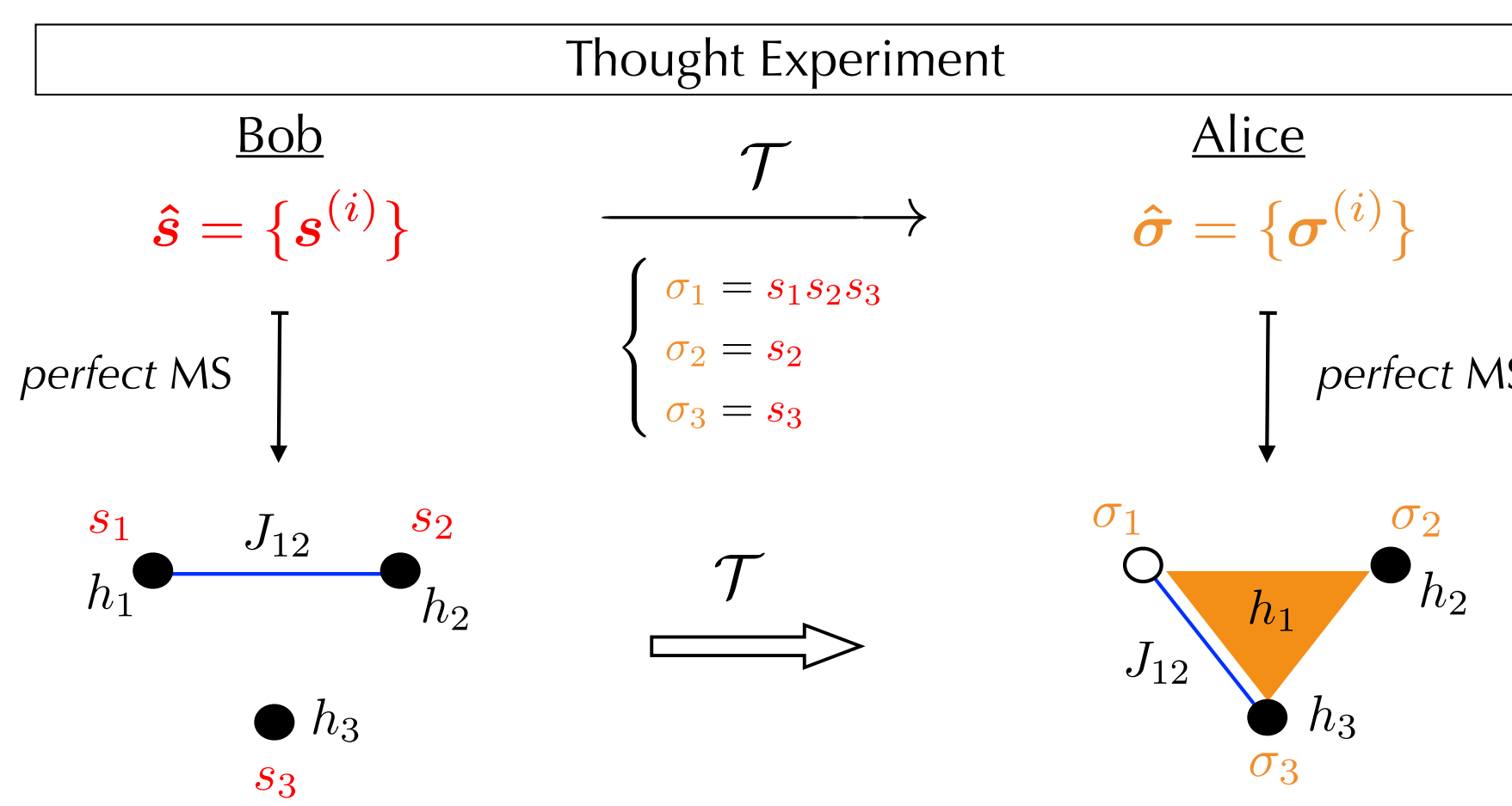
**Common simplification**, reduce the space of models to models with only field and pairwise interactions, which are simpler to infer and to interpret.

Is it a good idea? Are these models really simpler?

**Thought experiment:**

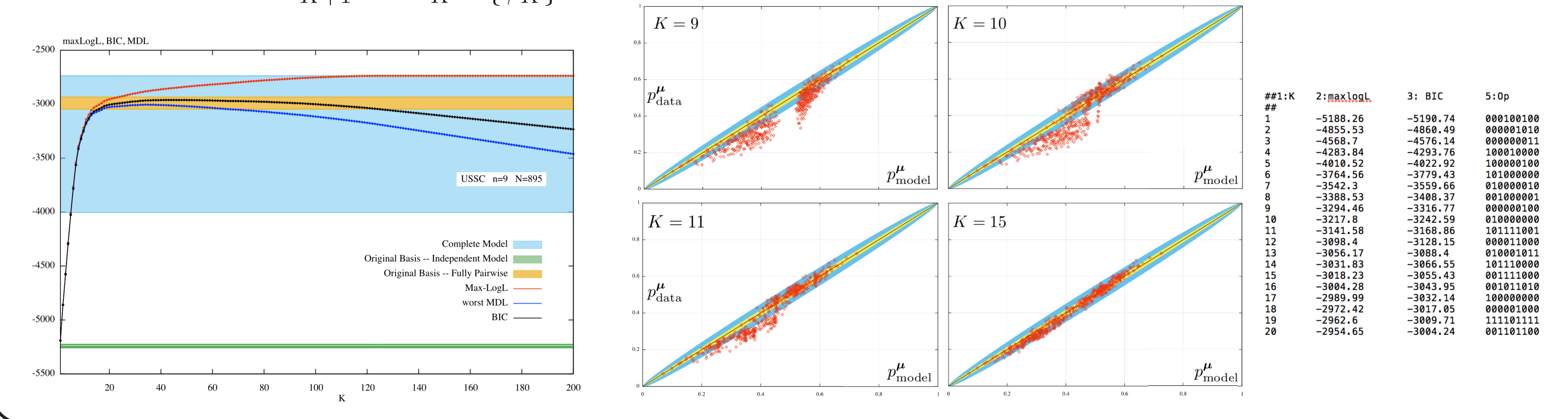
$\Rightarrow$  Pairwise interactions don't play any special role.

$\Rightarrow$  Model selection shouldn't be basis dependent...

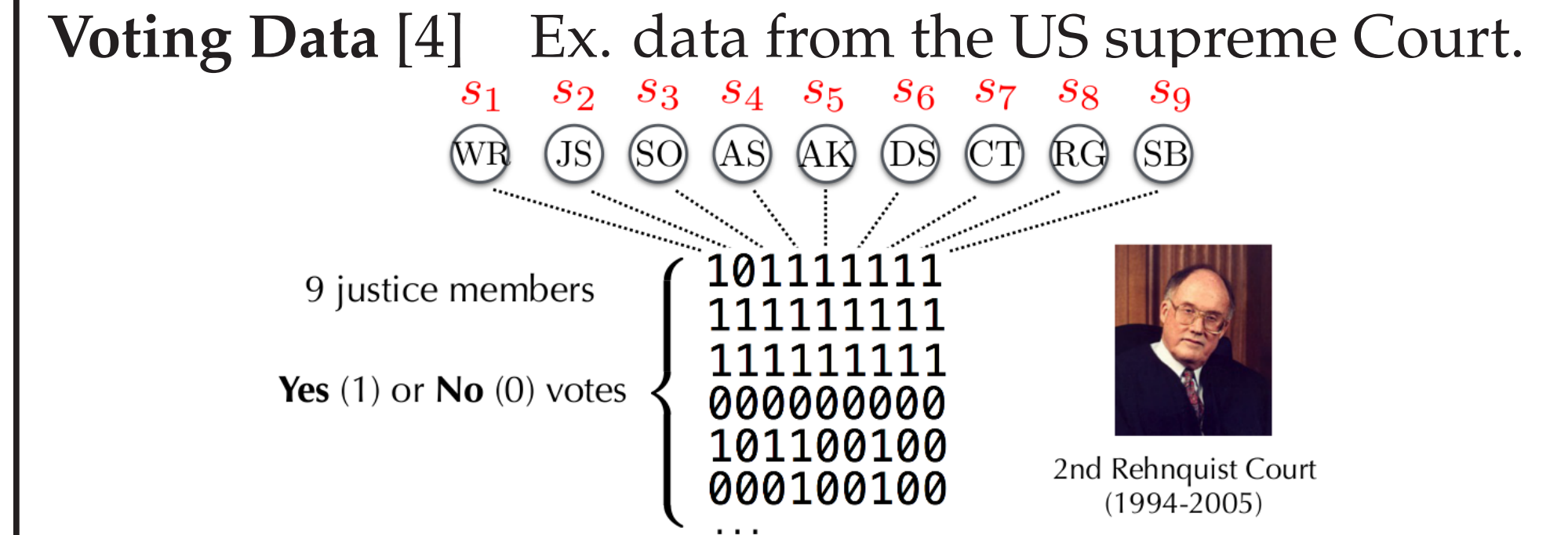
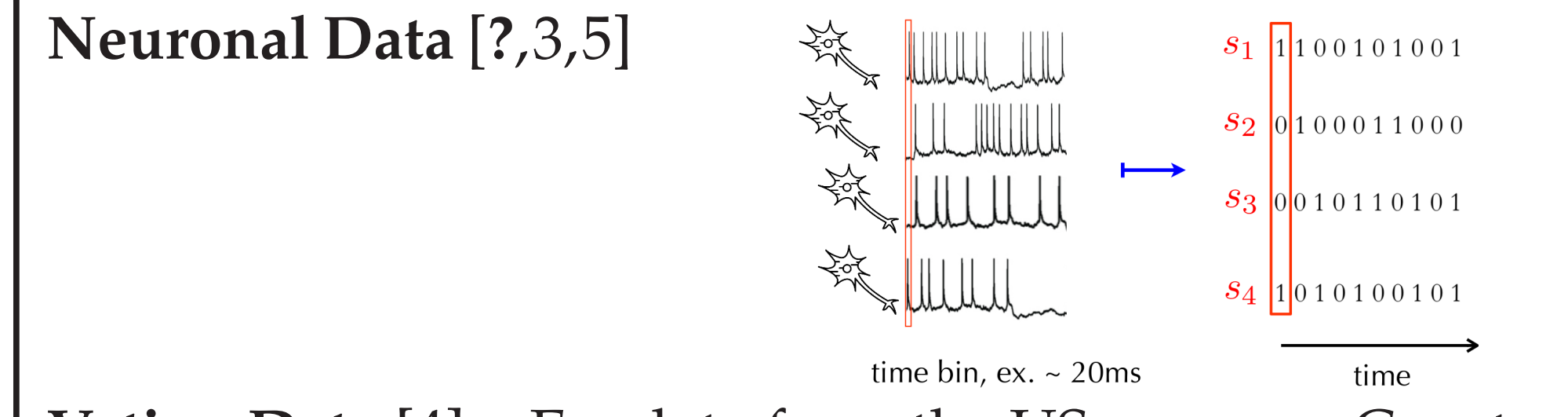


## APPROACH 1: SHORTEST PATH IN THE SPACE OF MODELS

How? Start  $\mathcal{M}_0 = \emptyset$ . At each step  $K$ , add to the model the operator  $\phi_K$  whose statistics is the less captured by the current model:  $\mathcal{M}_{K+1} = \mathcal{M}_K \cup \{\phi_K\}$ .



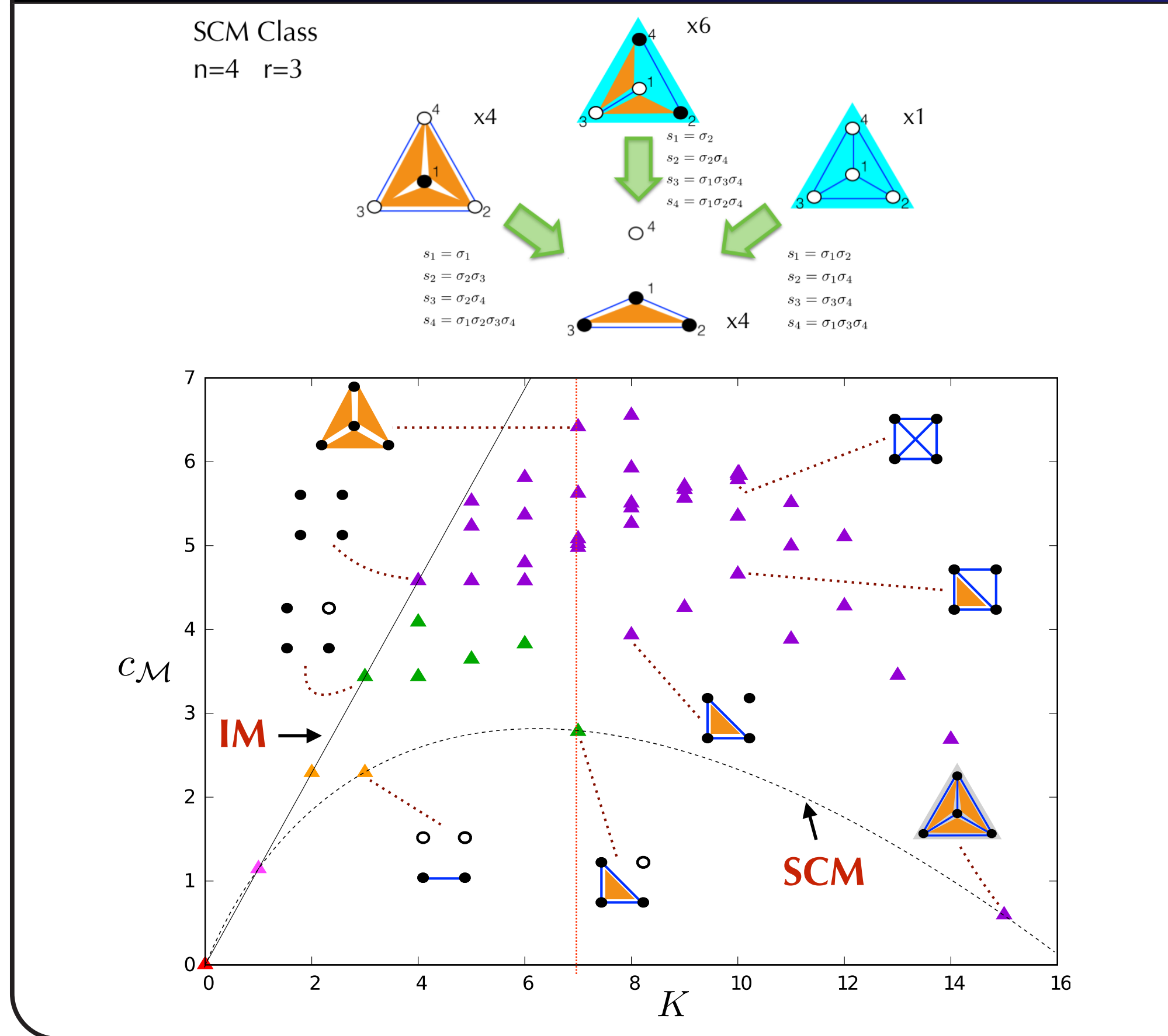
## EXAMPLES OF BINARY DATASETS



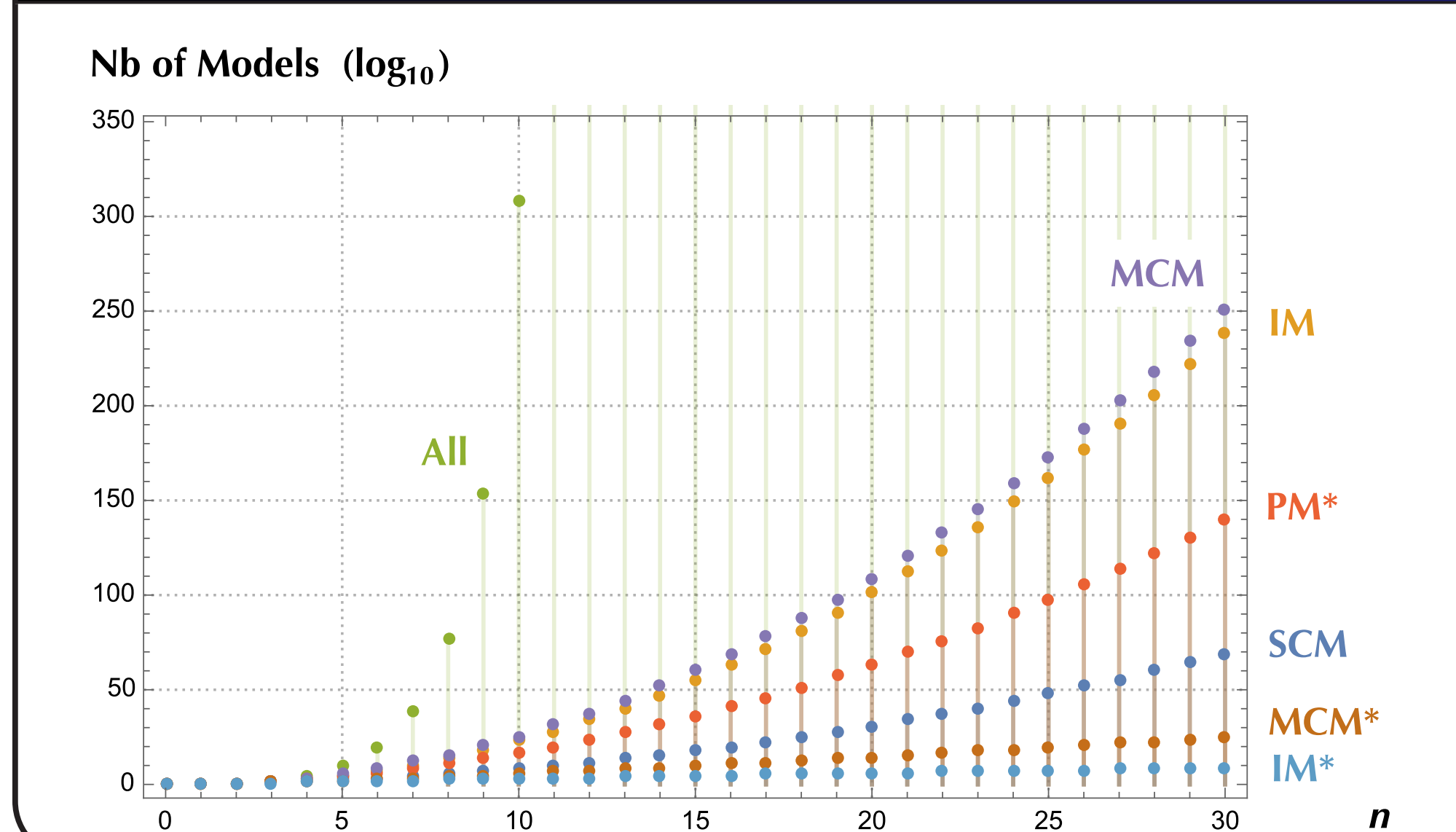
**Financial Data** [6] Time series, at each time step each stock  $i$  takes the value:

- $s_i = +1$ , if the stock has gone up (profit);
- $s_i = -1$ , if the stock has gone down (loss).

## EQUIVALENCE CLASSES

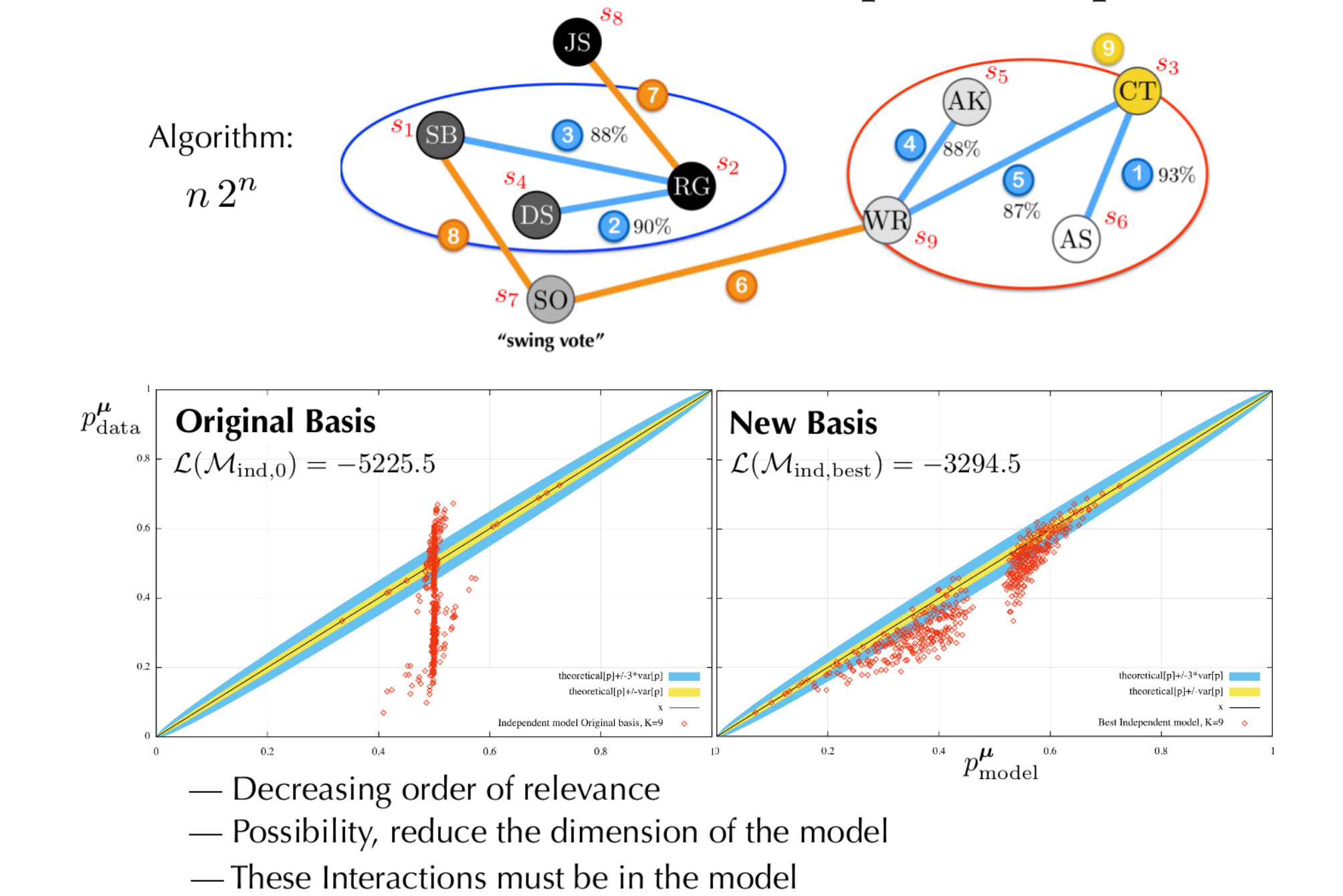


## HOW MANY... ?

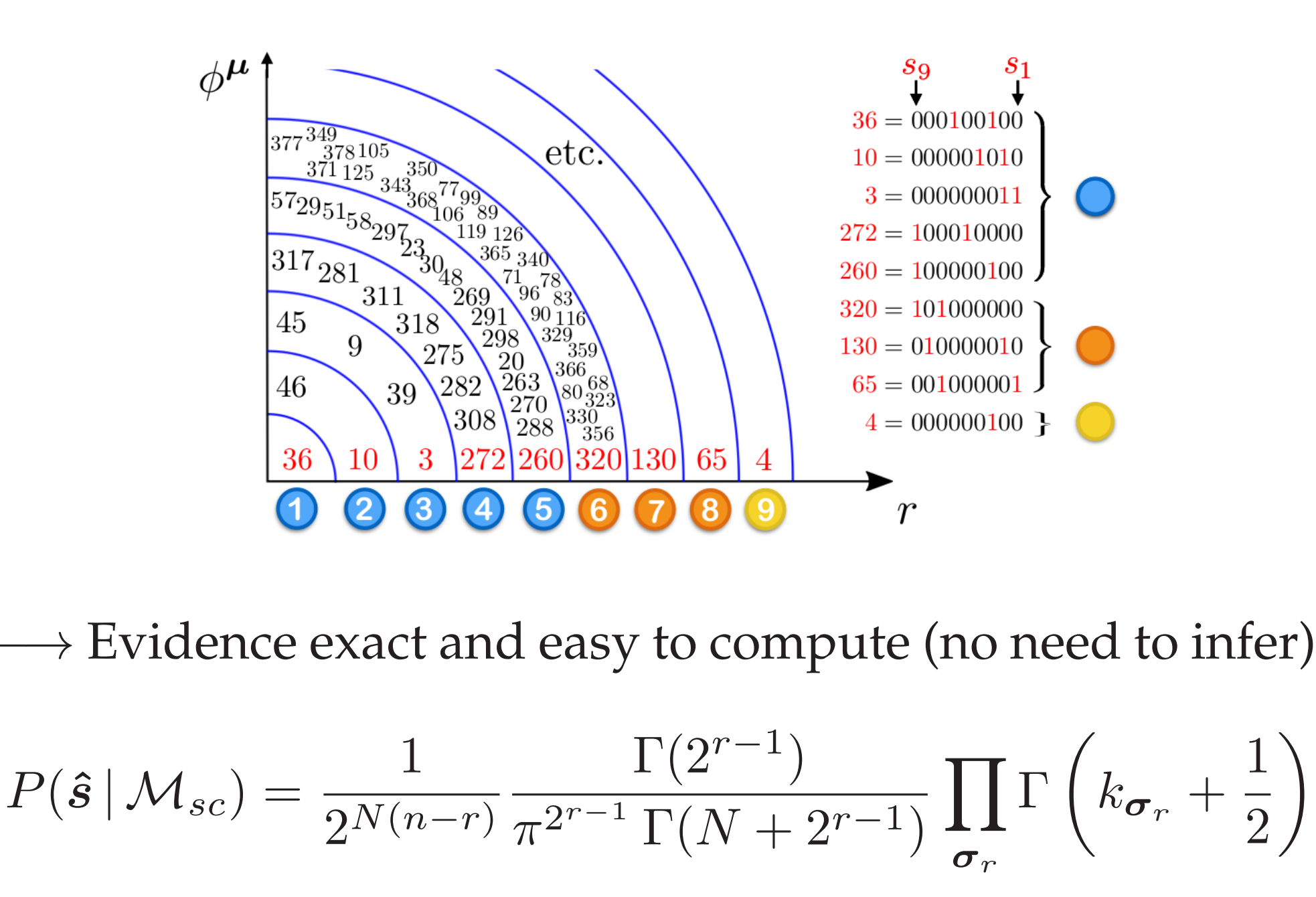


## APPROACH 2: SELECTION WITHIN AND BETWEEN CLASSES

**Independent Models (IM)**  $\rightarrow \mathcal{C}_{ind}(n, r)$   
 $(n+1)$  classes  $\sim \frac{2^{nr}}{r!}$  models/class  $K = r$   
 $|\mathcal{C}_{ind}(9, r)| = \{1, 511, 1.10^5, 2.10^7, 3.10^9, 3.10^{11}, 2.10^{13}, 1.10^{15}, 7.10^{16}, 2.10^{18}\}$   
 $\mathcal{L}_{max}^{ind} = -\sum_{i=1}^r S[p_i] - (n-r)N \log(2)$   
 $S[p] = -[p \log p + (1-p) \log(1-p)]$   $p_i = \mathbb{P}_{data}[\phi_i(\mathbf{s}) = -1]$   
 $\Rightarrow$  Best IM: set of the most bias independent operators



**Sub-Complete Models (SCM)**  $\rightarrow \mathcal{C}_{sc}(n, r)$   
 $(n+1)$  classes  $\sim 2^{r(n-r)}$  models/class  $K = 2^r - 1$   
 $|\mathcal{C}_{sc}(9, r)| = \{1, 511, 4.10^4, 8.10^5, 3.10^6, 3.10^6, 8.10^5, 4.10^4, 511, 1\}$   
 $\mathcal{L}_{max}^{sc} = \sum_{\sigma_r} k_{\sigma_r} \log \left( \frac{k_{\sigma_r}}{N} \right) - (n-r)N \log(2)$   
 $\rightarrow$  Monte Carlo simulation in each class:



**Minimally Complex Models (MCM)**  $\rightarrow \mathcal{C}_{mc}(n, \{r_i\})$   
 $\sum_{r=0}^n$  part(r) classes  $K = \sum_{r_i} (2^{r_i} - 1)$   
 $|\mathcal{C}_{mc}(9, \{1, 2, 3\})| = 1.5 \cdot 10^{13}$   
 $n = 9$   $r_1 = 3$   $r_2 = 1$   $r = 6$   
 $r_3 = 2$   $K = 11$   
 $\rightarrow$  Contains all IM and SCM.

