

ARE PAIRWISE MODELS REALLY SIMPLER?

ALBERTO BERETTA¹, CLAUDIA BATTISTIN², CLÉLIA DE MULATIER¹, IACOPO MASTROMATTEO³, MATTEO MARSILI¹

(1) The Abdus Salam International Centre for Theoretical Physics (ICTP), Trieste, Italy

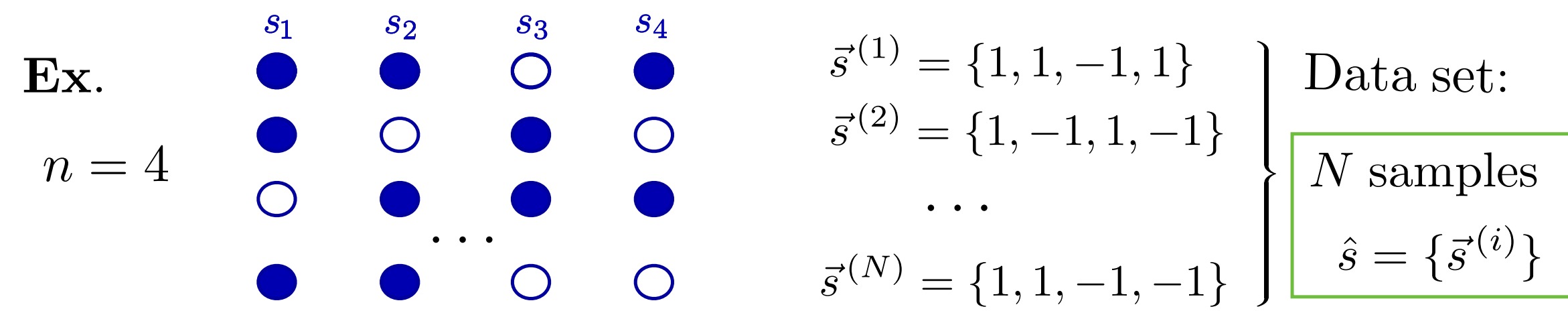
(2) Kavli Institute for Systems Neuroscience, NTNU, Norway

(3) Capital Fund Management, Paris, France

INTRODUCTION – BAYESIAN MODEL SELECTION ON SPIN MODELS

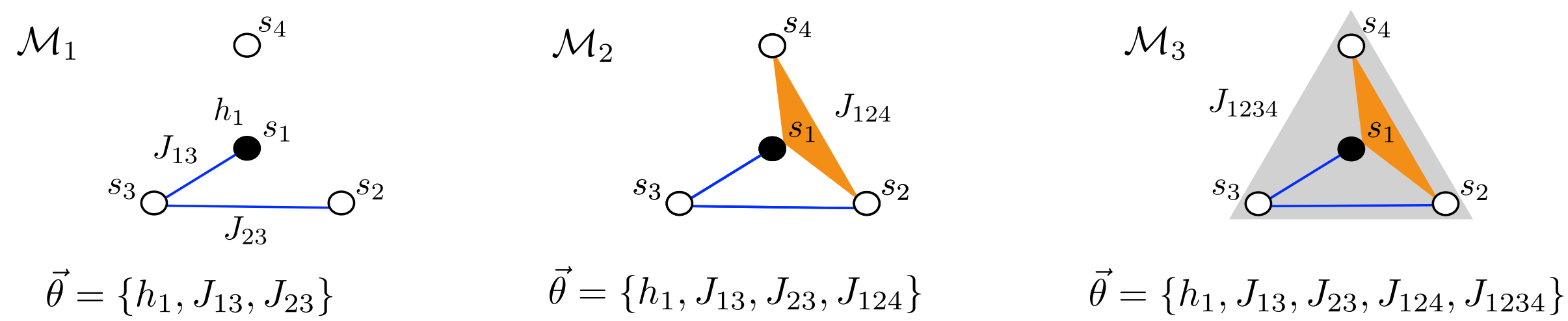
A fundamental issue in data analysis is to find the model that best captures the patterns hidden within the data, despite the random errors that effect them. The model should be complex enough to be able to fit the data, but simple enough to capture its main patterns.

Consider a system of n spins that take random values in $\{-1, +1\}$:



Q? What would be the best *model* for the system, that could explain what we observe/reproduce similar data?

Ex. of spin models \mathcal{M}_i

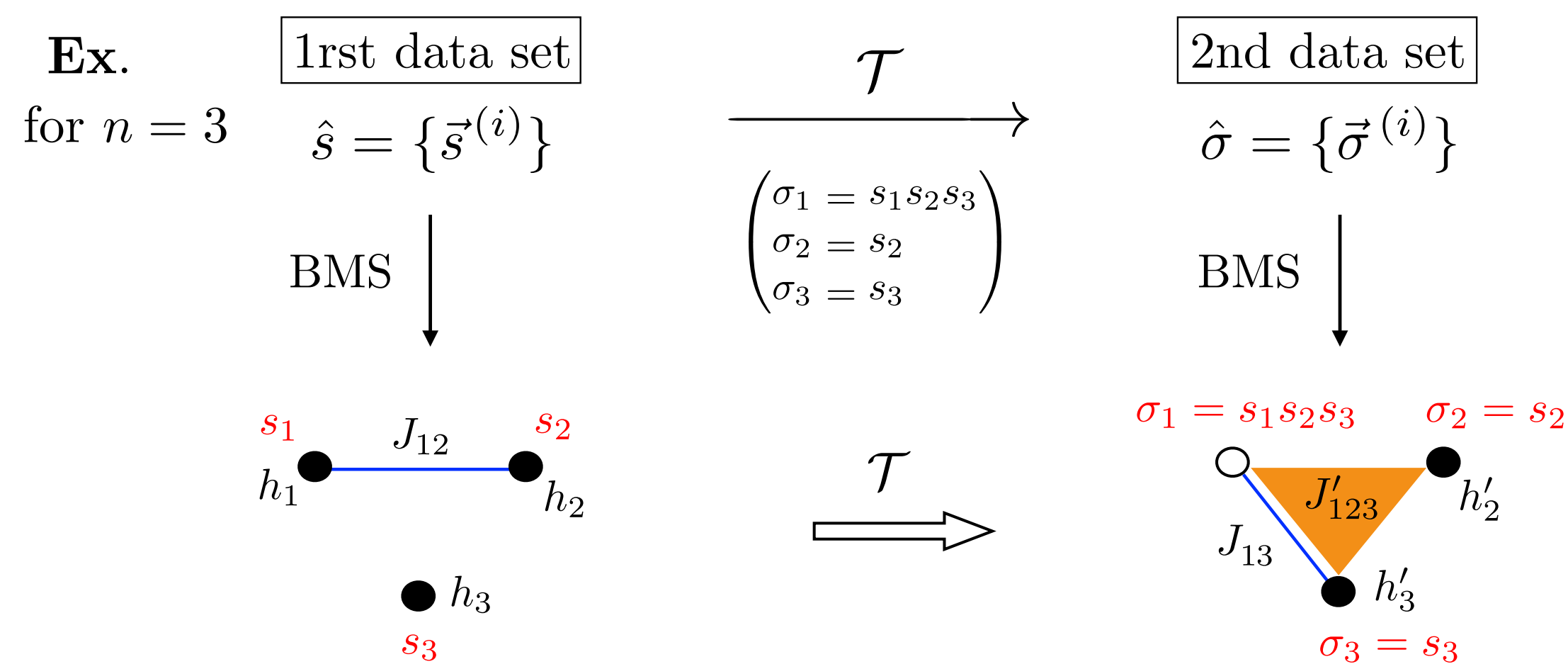


How? Classifying all the possible models \mathcal{M}_i to find the one with the highest *posterior probability* $P(\mathcal{M}_i | \hat{s})$. However, this is a very difficult task, due to the huge number of models.

To **simplify**, it is common to reduce the space of the possible models to models with only field and pairwise interactions.

Is it a good idea? Are these models simpler?

⇒ Pairwise interactions don't seem to play any special role.



HOW MANY MODELS?

For the spin models:

- $2^n - 1$ possible interactions
- $2^{2^n - 1}$ possible spin models

Ex. $n = 2$: 8 models $n = 4$: 32768 models
 $n = 3$: 128 models $n = 5$: 2147483648 models

Models with only fields and pairwise interactions:

- $n(n+1)/2$ possible interactions
- $\sim 2^{n^2}$ possible models

Ex. $n = 2$: 8 models $n = 4$: 1024 models
 $n = 3$: 64 models $n = 5$: 32768 models

SPACE OF PDFs

o For a given a model \mathcal{M} , the family of probability distributions $\{P(\vec{s} | \vec{\theta}, \mathcal{M})\}_{\vec{\theta}}$ forms a Riemannian manifold with natural coordinates $\vec{\theta}$ [1].

o Each point of this space is a probability distribution $P(\vec{s} | \vec{\theta}, \mathcal{M})$.

o The natural metric on this manifold is given by the Fisher Information Matrix [1]:

$$J_{qk}(\vec{\theta}) = \partial_{\theta_q} \partial_{\theta_k} \log Z_{\mathcal{M}}(\vec{\theta}) = \langle f_q f_k \rangle - \langle f_q \rangle \langle f_k \rangle$$

o Varying the parameters $\vec{\theta}$ from $d^K \vec{\theta}$ gives rise to distributions similar to $P(\vec{s} | \vec{\theta}, \mathcal{M})$ that correspond to nearby points in the manifold, contained in the small volume:

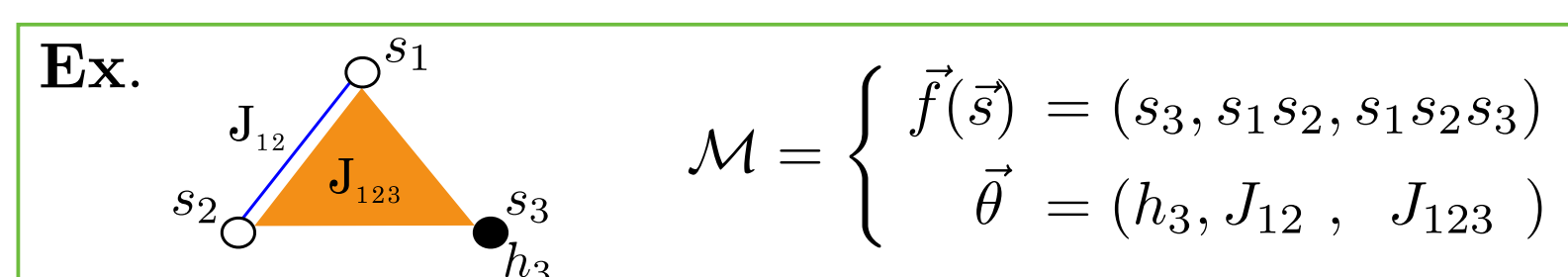
$$dV_{\mathcal{M}} = \sqrt{\det J(\vec{\theta})} d^K \vec{\theta}$$

COMPLEXITY OF SPIN MODELS

Using Bayes' theorem: $P(\mathcal{M} | \hat{s}) = \frac{P(\hat{s} | \mathcal{M}) P_0(\mathcal{M})}{P(\hat{s})}$, where $P(\hat{s} | \mathcal{M}) = \int d^K \vec{\theta} P(\hat{s} | \vec{\theta}, \mathcal{M}) P_0(\vec{\theta} | \mathcal{M})$

In absence of information, the prior $P_0(\mathcal{M})$ can be taken uniform, and models may be ranked directly with $P(\hat{s} | \mathcal{M})$.

Probability that the spin system is in the configuration $\vec{s}^{(i)}$:



$$P(\vec{s}^j | \vec{\theta}, \mathcal{M}) = \frac{e^{\sum_{k=1}^K f_k(\vec{s}^j) \theta_k}}{Z_{\mathcal{M}}(\vec{\theta})}$$

Spin operator: product of the spins involve in the interaction k

Expanding for a large size N of the data set, finally leads, in the framework of Bayesian Model Selection, to [2, 4]:

$$\log P(\hat{s} | \mathcal{M}_i) = \underbrace{\log P(\hat{s} | \mathcal{M}_i, \theta^*)}_{\propto N} - \underbrace{\frac{K}{2} \log \left(\frac{N}{2} \right)}_{\propto \log N} - \underbrace{c_{\mathcal{M}}}_{\text{Penalty term geometrical complexity}} + O\left(\frac{1}{N}\right)$$

Maximum log-Likelihood Penalty term number of parameter K

Geometrical Complexity

$$c_{\mathcal{M}} = \log \int_{\mathbb{R}^K} \sqrt{\det J(\vec{\theta})} d^K \vec{\theta}$$

PRIOR ON THE VALUES OF $\vec{\theta}$

Best choice in absence of any information [2]:

Jeffreys' prior: $P_0(\vec{\theta} | \mathcal{M}) = \frac{\sqrt{\det J(\vec{\theta})}}{\int \sqrt{\det J(\vec{\theta})} d^K \vec{\theta}}$

- invariant under re-parametrisation [5];
- uniform in the space of observables.

COMPLEXITY – INTERPRETATION

$$c_{\mathcal{M}} = \log V_{\mathcal{M}}$$

$V_{\mathcal{M}}$ is the total volume of the manifold defined by \mathcal{M} :

- o Complexity represents how broad the model is in term of describing various probability distributions [2].
- o A model is complex if it can fit a wide range of data.

GAUGE TRANSFORMATIONS

$c_{\mathcal{M}}$ is expected to stay invariant under the transformation \mathcal{T} introduced in Ex a. This invariance emerges explicitly when expressing $Z_{\mathcal{M}}(\vec{\theta})$ in the form:

$$Z_{\mathcal{M}}(\vec{\theta}) = 2^n \prod_{\mu \in \mathcal{M}} \cosh(\theta^\mu) \left[1 + \sum_{\ell \in \mathcal{L}} \prod_{\mu \in \ell} \tanh(\theta^\mu) \right]$$

- Loop ℓ : subset $\ell \subseteq \mathcal{M}$ such that $\prod_{\mu \in \ell} f^\mu(\vec{s}) = 1$;
- Set \mathcal{L} : set of all the loops of \mathcal{M} .

$Z_{\mathcal{M}}(\vec{\theta})$ depends on few characteristics of \mathcal{M} :

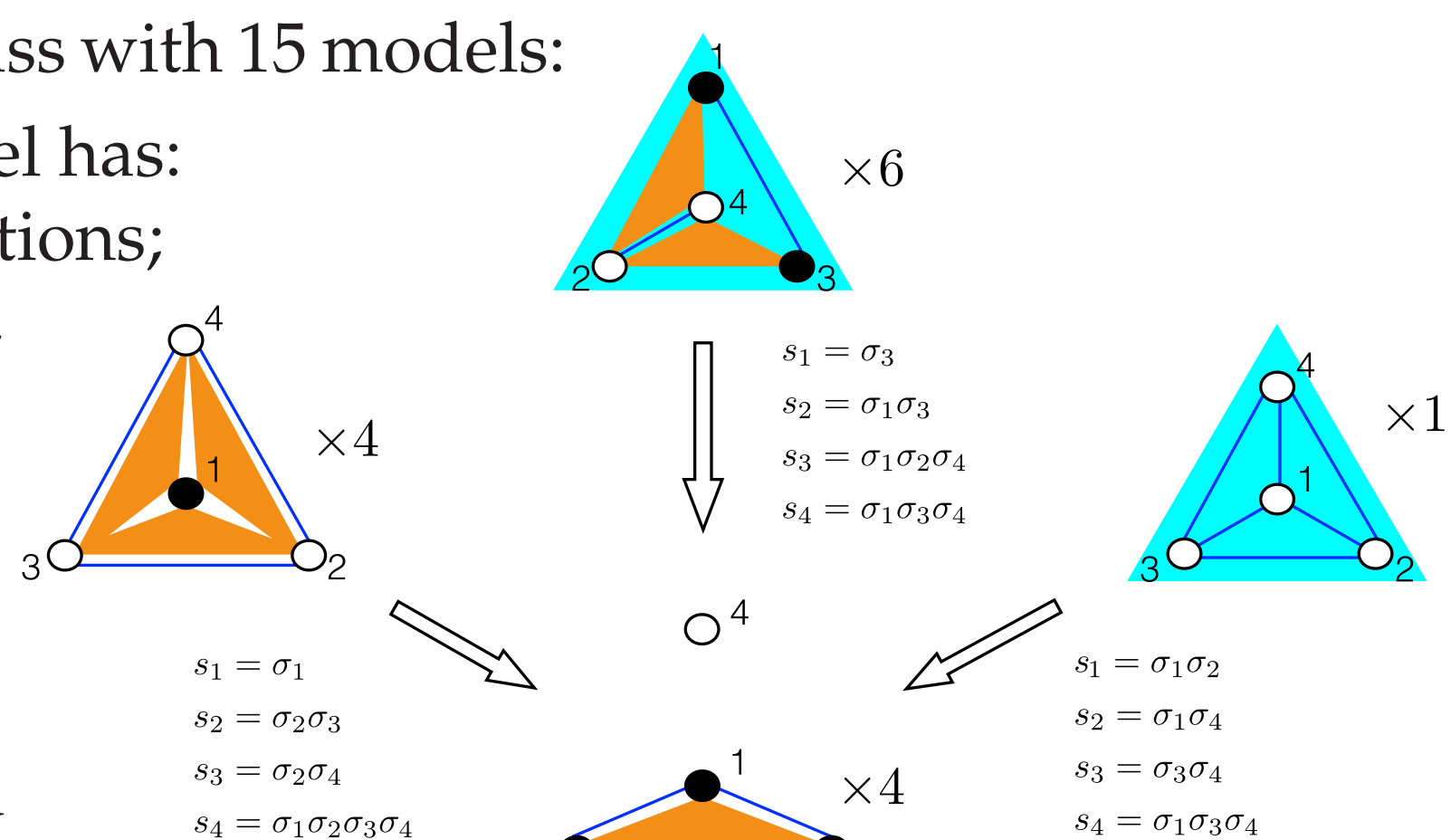
- number $|\mathcal{M}|$ of operators;
- the structure of the loops \mathcal{L} of the \mathcal{M} ;

which are invariant under the transformation \mathcal{T} . We call (such a transformation) \mathcal{T} a Gauge Transformation as it preserves the geometry of the model.

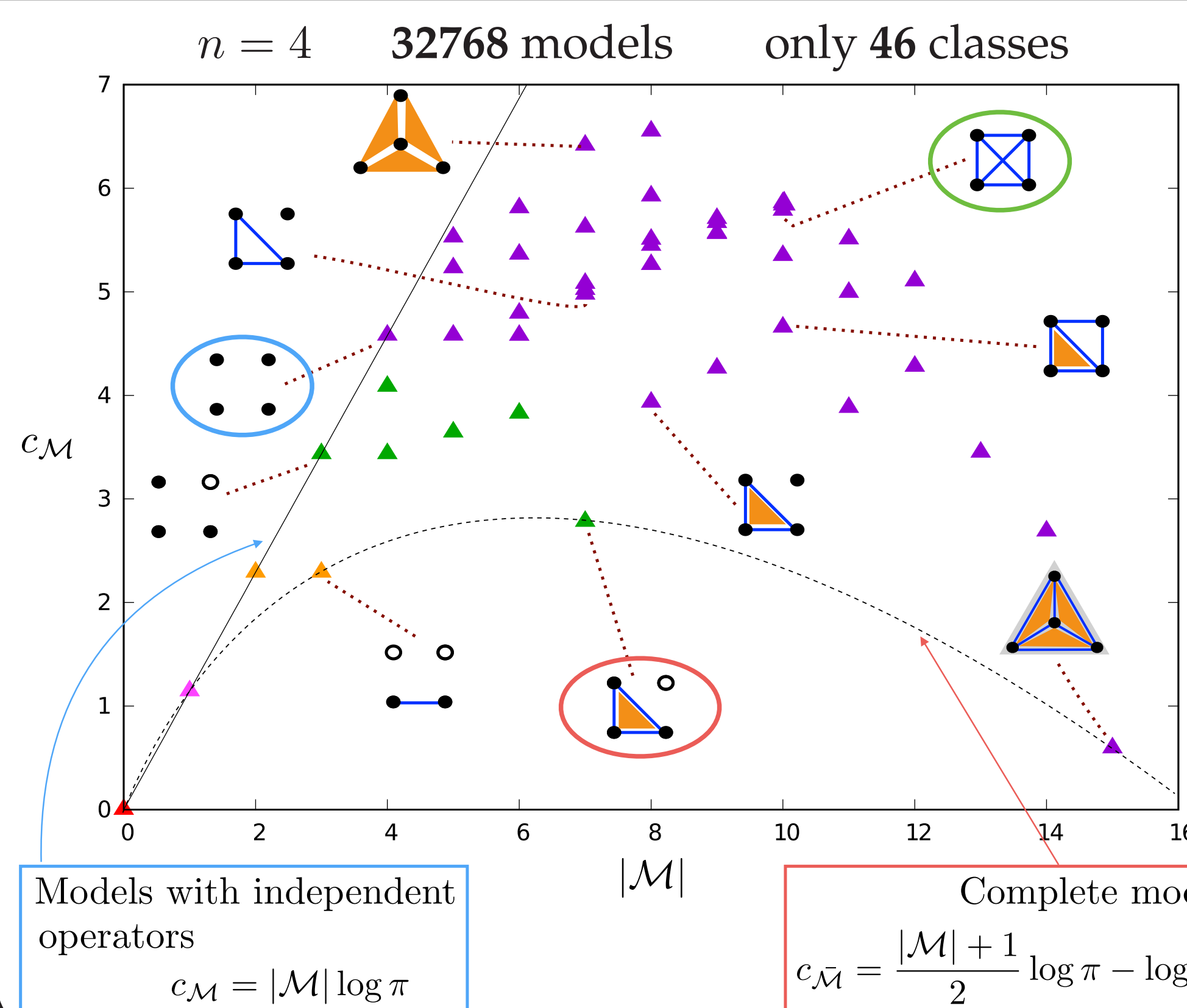
Thus $c_{\mathcal{M}}$ stays invariant under \mathcal{T} , which allows defining *complexity classes* of models (images through a GT and with the same value of $c_{\mathcal{M}}$).

Ex. one class with 15 models:

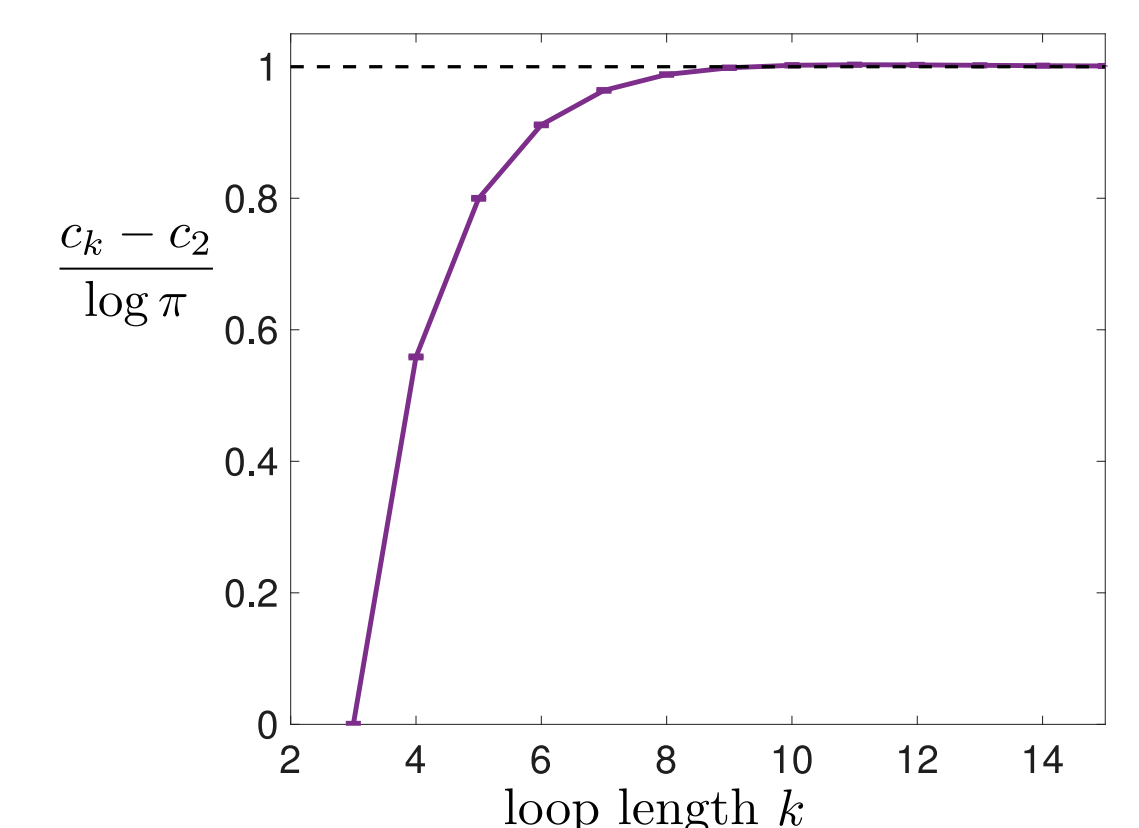
Each model has:
- 7 interactions;
- 15 loops,



COMPLEXITY CLASSES



- o Models with pairwise interactions are not simpler.
- o $c_{\mathcal{M}}$ is not monotonic in the number of parameters.
- o $c_{\mathcal{M}}$ is between two limit curves:
 - Complete models are the simplest models;
 - Models with only independent operators are the most complex.
- o Complementary classes: same number of classes for models with $|\mathcal{M}|$ operators than for models with $2^n - 1 - |\mathcal{M}|$.
- o Models with a single loop: $c_{\mathcal{M}}$ increases with the length of the loop.



DEGENERATE MODELS

Degenerate models Operators share parameters:
→ $|\mathcal{M}|$ operators, but K parameters with $K < |\mathcal{M}|$;
→ α_k degeneracy of θ_k , for $k \in [1, K]$.

Result In general, degeneracy reduces the complexity.

Ex 1. models with $|\mathcal{M}| = \sum_{k=1}^K \alpha_k$ independent operators:

$$\exp c_{\mathcal{M}}^{non-deg} = \pi^{|\mathcal{M}|} \quad \exp c_{\mathcal{M}}^{deg} = \pi^K \prod_{k=1}^K \sqrt{\alpha_k}$$

Ex 2. models with a single loop:

