

MINING SUBJECTIVELY INTERESTING PATTERNS IN DATA

PART 2/5: THE FORSIED FRAMEWORK

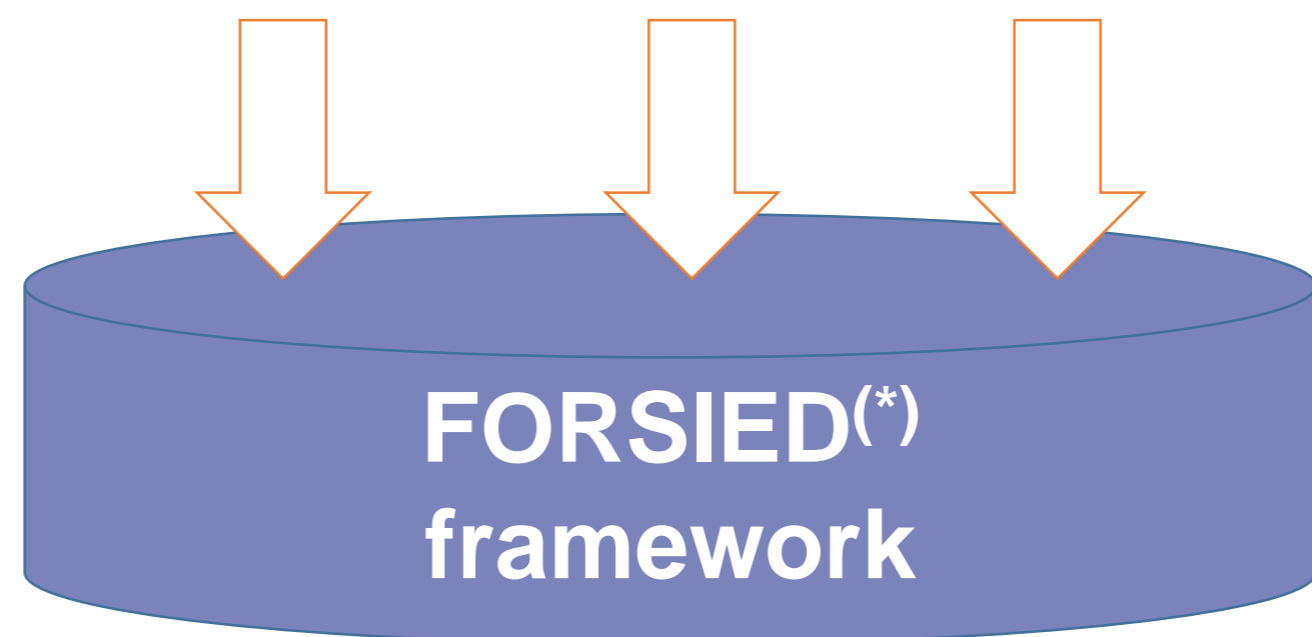
Jefrey Lijffijt Tijl De Bie
Ghent University

GENERIC FRAMEWORK

- **Data:** the adjacency matrix of a graph under study
- **Patterns:** the claim that a specified set of nodes are densely connected
- **Prior beliefs:** the degrees of the nodes, a known block structure,...
- **Interestingness:** subjective information density
- Overlapping communities!

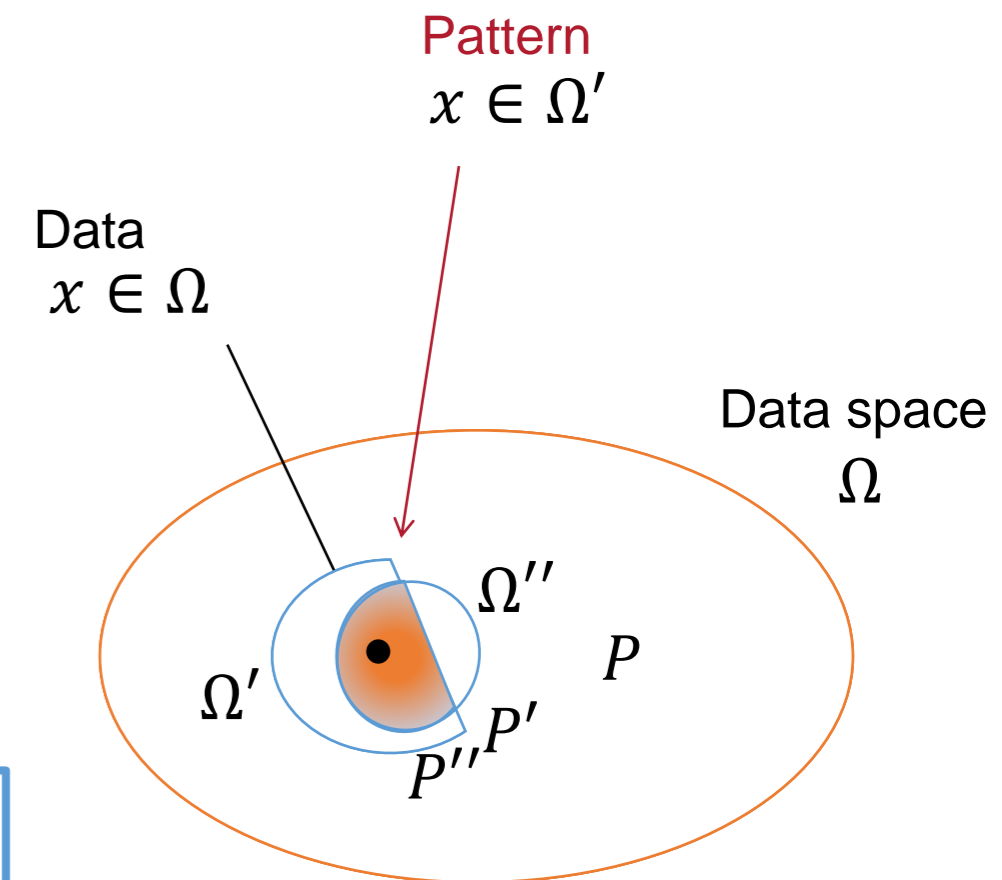
Data Patterns Model of user beliefs
("background distribution")

P evolves!



$$\text{InformationContent}(\Omega', P) = -\log(P(\Omega'))$$

Pattern User **Subjective!**



$$\text{Interestingness}(\Omega', P) = \frac{\text{InformationContent}(\Omega', P)}{\text{DescriptionLength}(\Omega')}$$

SI = IC/DL



(*)Formalizing Subjective Interestingness in Exploratory Data mining

THE FINE PRINT

- Initial background distribution P ?
- **Maximum entropy** distribution

$$\max_P E_{X \sim P} \{-\log P(X)\}$$

- Updated background distribution P' given pattern $x \in \Omega'$?
- P **conditioned** onto event $x \in \Omega'$

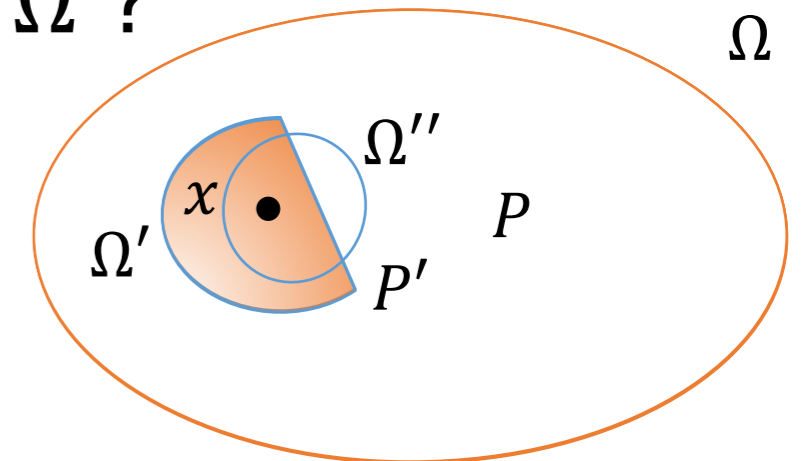
$$P'(\Omega'') = \frac{P(\Omega'' \cap \Omega')}{P(\Omega')}$$

$$\Rightarrow -\log(P'(x)) = -\log(P(x)) + \log(P(\Omega'))$$

Information content
in data after pattern

Information content
in data before pattern

Minus information content
of pattern under P



- Description length?
- Essentially problem-dependent

WHY MAXIMUM ENTROPY / CONDITIONING?

- Most **unbiased estimate**

Informal... no bias other than the constraints

- Assume **cautious / pessimistic user**

A user who expects to be very surprised

- Leads to **most robust estimate** of true subjective information content

Information content estimated with maxent P will never differ much from information content w.r.t. true prior belief of user

A FIRST INSTANTIATION: COMMUNITY DETECTION

COMMUNITY DETECTION IN NETWORKS

Data:

- Graph

Prior beliefs:

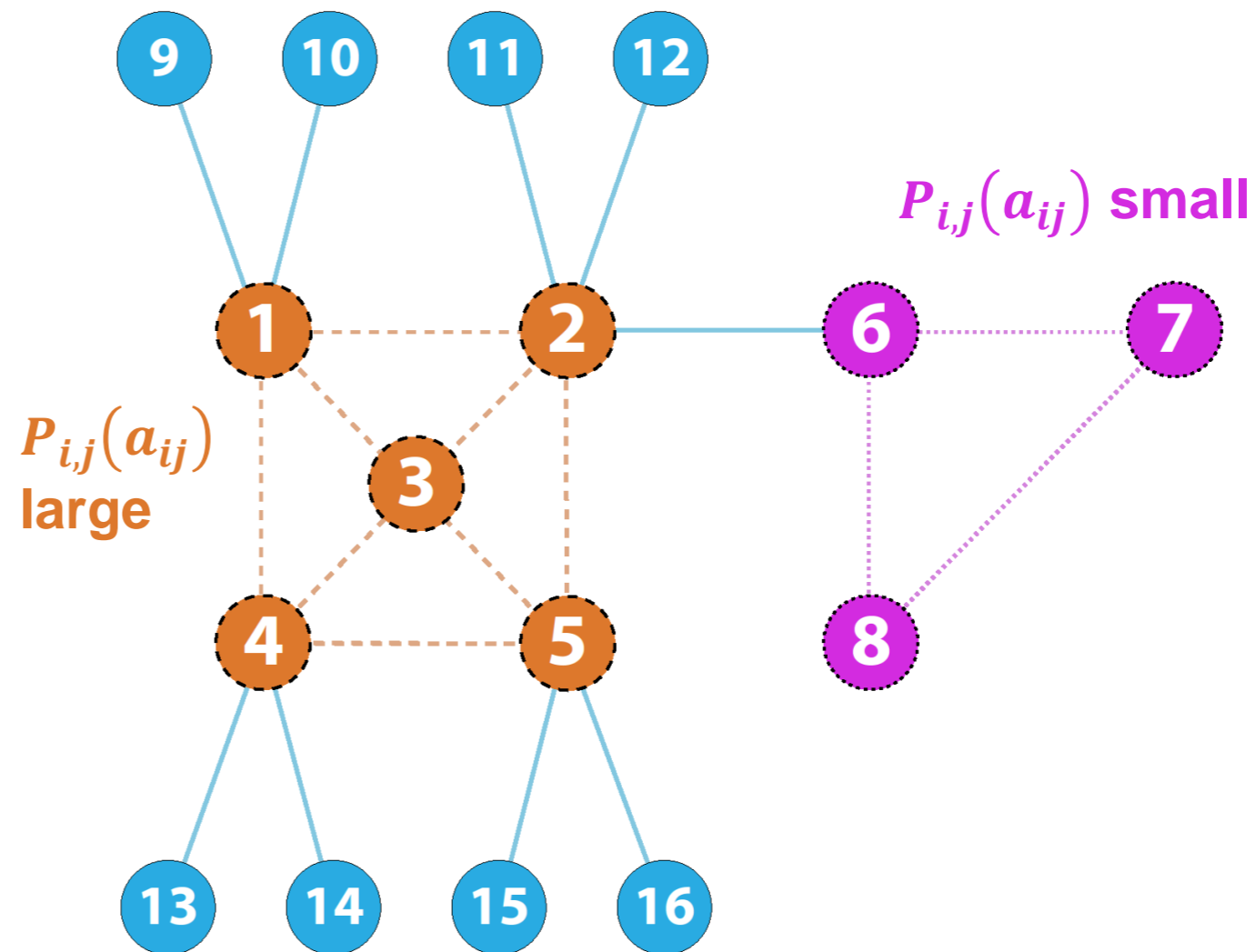
1. Overall density
2. or: Vertex degrees

→ MaxEnt distribution:

$$P(A) = \prod_{i>j} P_{i,j}(a_{ij})$$

Adjacency matrix \uparrow Edge indicator variables \uparrow

$$P_{i,j}(a_{ij}) = \frac{\exp(a_{ij} \cdot (\lambda_i + \lambda_j))}{1 + \exp(\lambda_i + \lambda_j)}$$



COMMUNITY DETECTION IN NETWORKS

Data:

- Graph

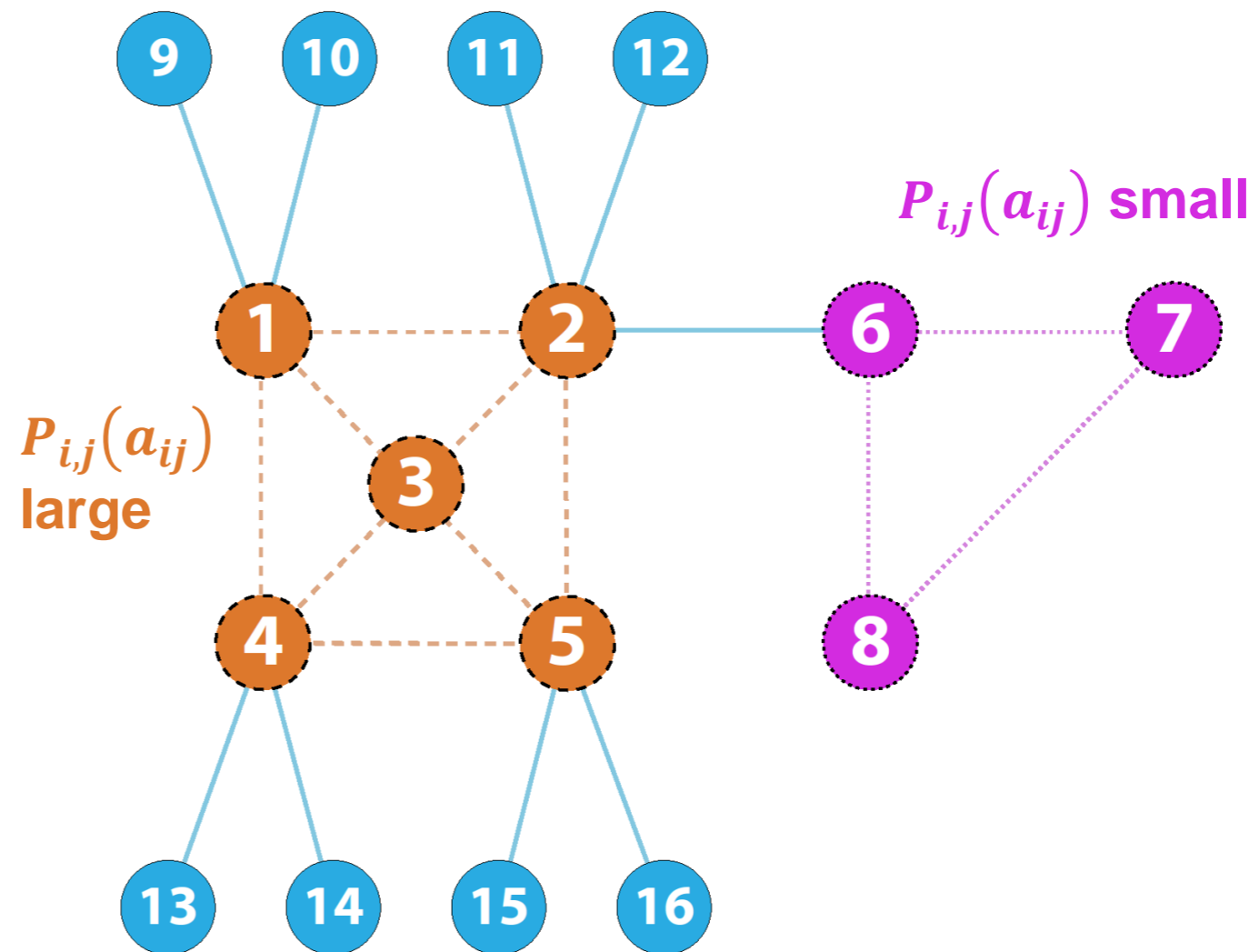
Prior beliefs:

1. Overall density
2. or: Vertex degrees

Pattern:

- Dense subgraphs

$$\sum_{i,j \in \text{subgraph}} a_{ij} \geq k$$



COMMUNITY DETECTION IN NETWORKS

Data:

- Graph

Prior beliefs:

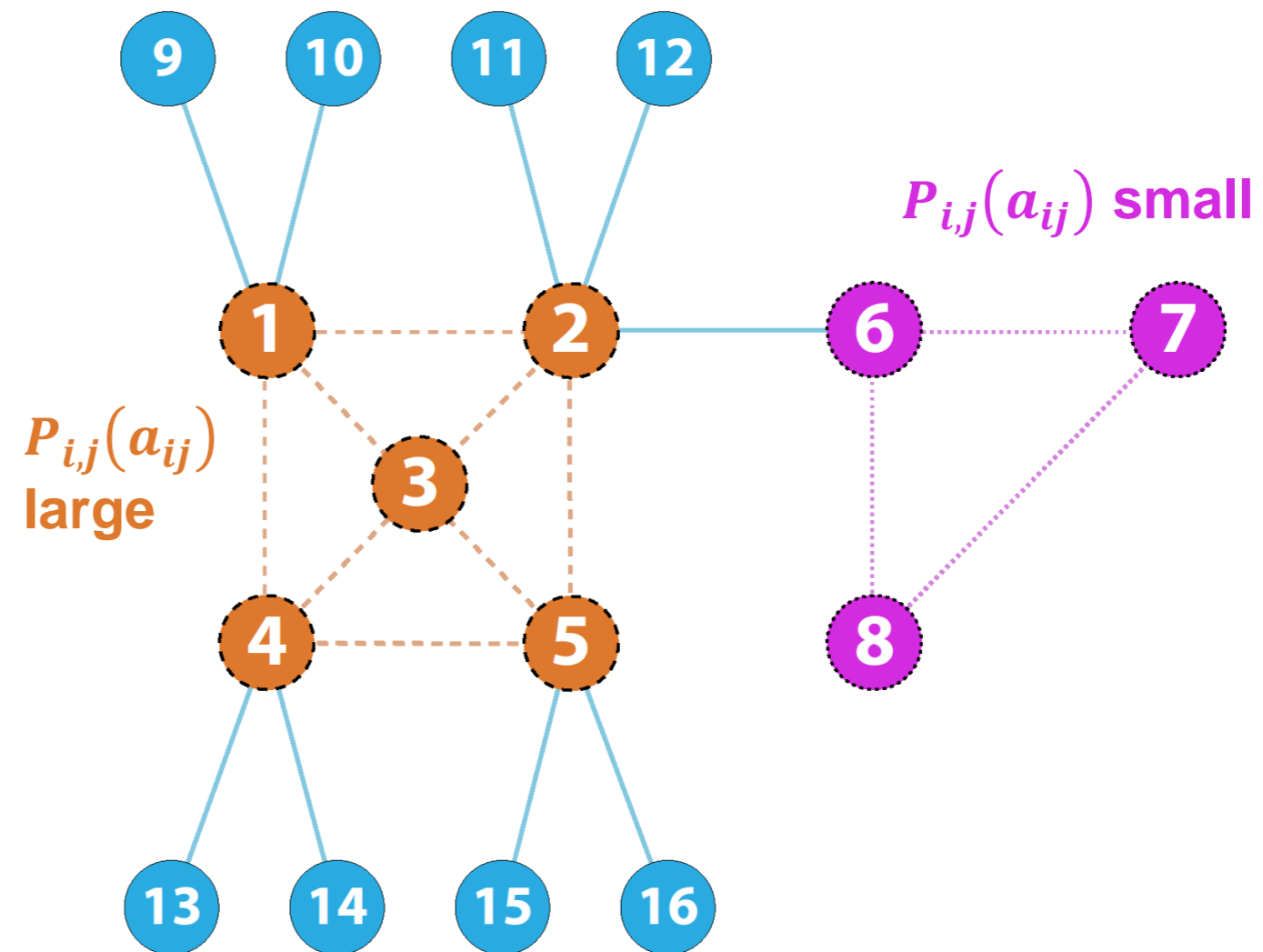
1. Overall density
2. or: Vertex degrees

Pattern:

- Dense subgraphs

Interestingness:

$$\frac{-\log P(\text{pattern})}{\text{DescriptionLength}(\text{pattern})}$$



COMMUNITY DETECTION IN NETWORKS

Data:

- Graph

Prior beliefs:

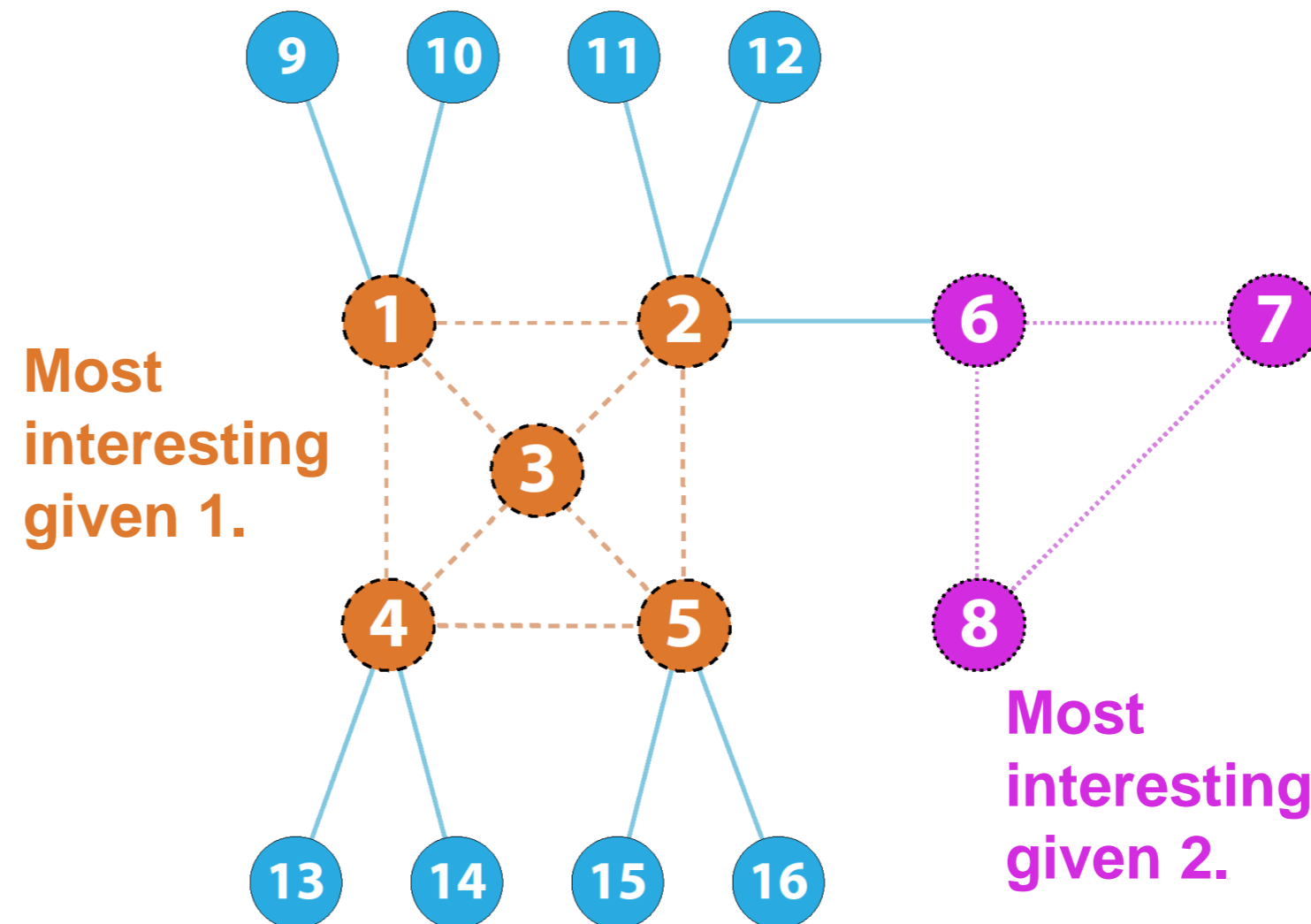
1. Overall density
2. or: Vertex degrees

Pattern:

- Dense subgraphs

Interestingness:

- Density vs. size
- 2. → preferably low degree nodes



COMMUNITY DETECTION IN NETWORKS

Data:

- Graph

Prior beliefs:

1. Overall density
2. or: Vertex degrees

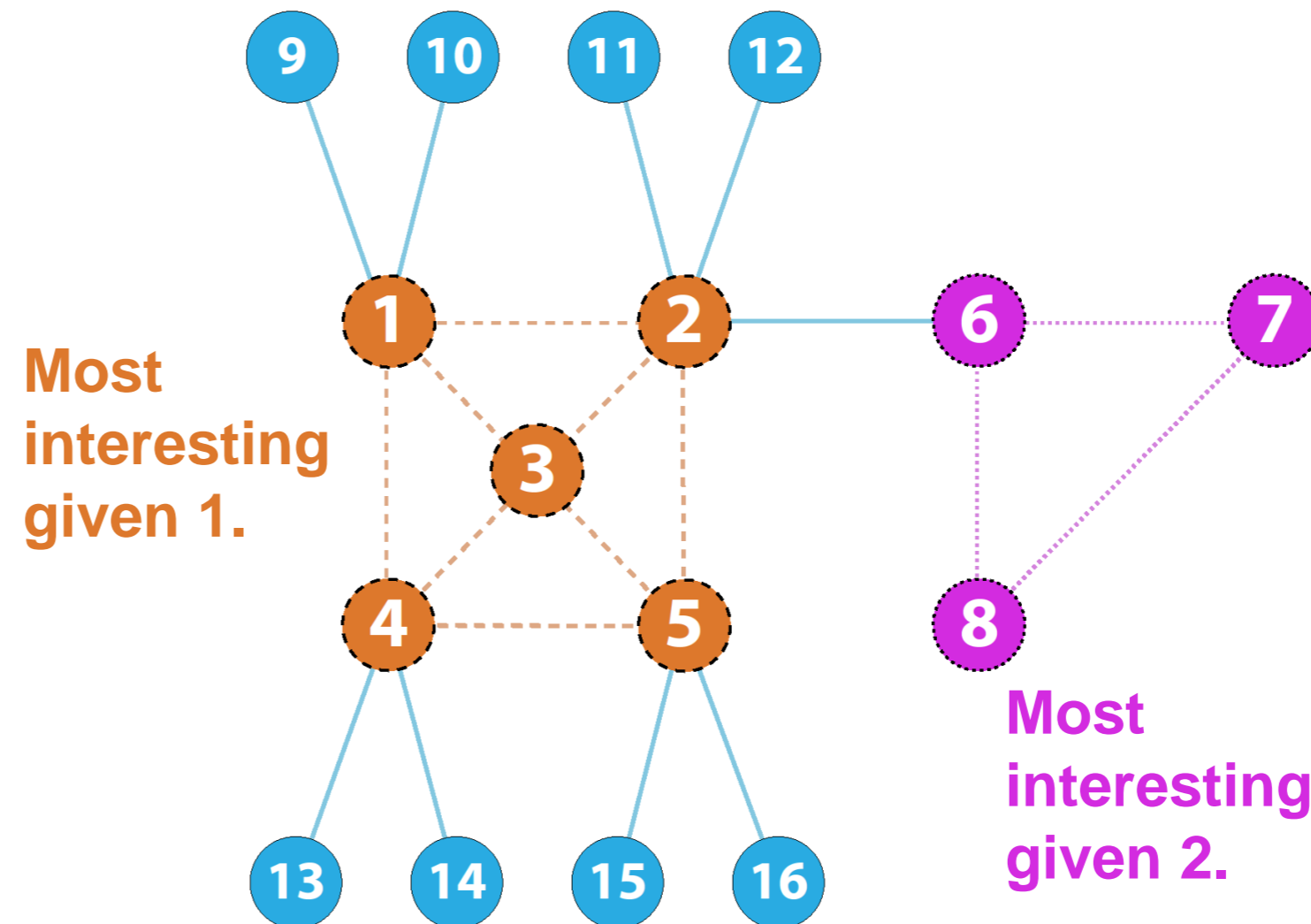
Pattern:

- Dense subgraphs

Interestingness:

- Density vs. size
- 2. → preferably low degree nodes

 **Hill-climbing for search**
Update P after each pattern



COMMUNITY DETECTION IN NETWORKS

Data:

- Graph

Prior beliefs:

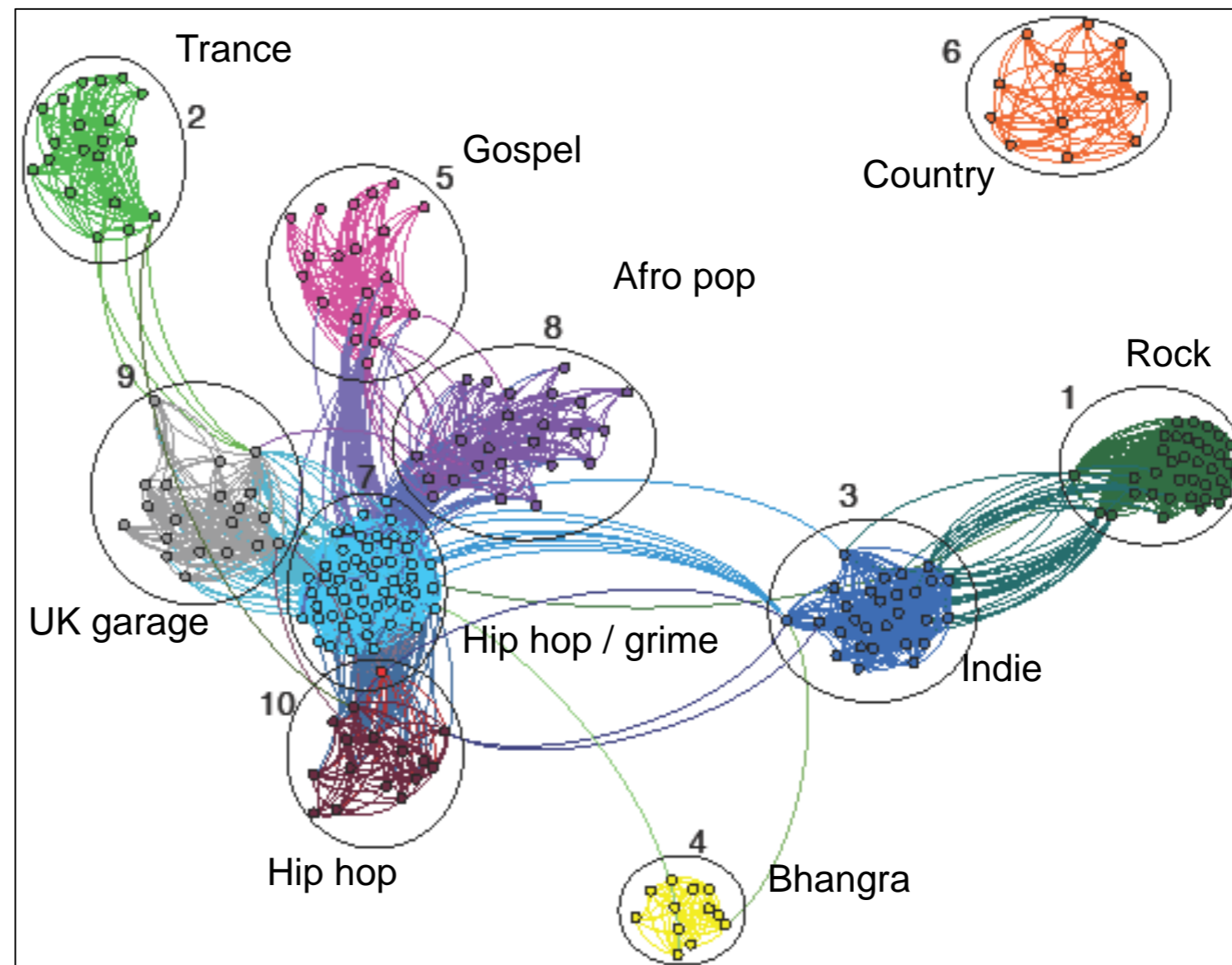
1. Overall density
2. or: Vertex degrees

Pattern:

- Dense subgraphs

Interestingness:

- Density vs. size
- 2. → preferably low degree nodes



 Hill-climbing for search
Update P after each pattern

TAKE-AWAYS

1. What is the **data**?
2. Determine suitable **pattern** syntax
3. What are the **prior beliefs**? (= what is irrelevant to user?)
 - Compute **background distribution P** using maximum entropy
4. Formulate **subjective interestingness**:

$$\text{InformationContent}(\Omega', P) = -\log(P(\Omega'))$$

$$\text{Interestingness}(\Omega', P) = \frac{\text{InformationContent}(\Omega', P)}{\text{DescriptionLength}(\Omega')}$$

5. Design an **algorithm to optimize** it
6. Find out how to **condition background distribution** on a pattern

THE BACKGROUND DISTRIBUTION: MAXENT

MAXENT MODEL S.T. DEGREE BELIEFS

$$\begin{array}{l}
 \max_P \sum_A -P(A) \cdot \log(P(A)) \quad \leftarrow \text{Entropy} \\
 \text{s.t.} \sum_A P(A) \cdot \left(\sum_{j=1:n} a_{ij} \right) = d_i, \quad \leftarrow \begin{array}{l} \text{Degree} \\ \forall i = 1:n \\ \text{Expected degree constraint} \end{array} \\
 \sum_A P(A) = 1 \quad \leftarrow \text{Normalization}
 \end{array}$$

Convex!

Lagrangian:

$$L(P, \lambda, \mu) = \sum_A -P(A) \cdot \log(P(A)) + \sum_{i=1:n} \lambda_i \left(\sum_A P(A) \cdot \left(\sum_{j=1:n} a_{ij} \right) - d_i \right) + \mu \left(\sum_A P(A) - 1 \right)$$

Lagrange multipliers

MAXENT MODEL S.T. DEGREE BELIEFS

$$L(P, \lambda, \mu) = \sum_A -P(A) \cdot \log(P(A)) + \sum_{i=1:n} \lambda_i \left(\sum_A P(A) \cdot \left(\sum_{j=1:n} a_{ij} \right) - d_i \right) + \mu \left(\sum_A P(A) - 1 \right)$$



MAXENT MODEL S.T. DEGREE BELIEFS

$$L(P, \lambda, \mu) = \sum_A -P(A) \cdot \log(P(A)) + \sum_{i=1:n} \lambda_i \left(\sum_A P(A) \cdot \left(\sum_{j=1:n} a_{ij} \right) - d_i \right) + \mu \left(\sum_A P(A) - 1 \right)$$

Optimality condition:

$$\frac{\partial}{\partial P(A)} L(P, \lambda, \mu) = 0$$

$$\frac{\partial}{\partial P(A)} L(P, \lambda, \mu) = -\log(P(A)) - 1 + \left(\sum_{i,j=1:n} \lambda_i a_{ij} \right) + \mu = 0$$

MAXENT MODEL S.T. DEGREE BELIEFS

$$\frac{\partial}{\partial P(\mathbf{A})} L(P, \boldsymbol{\lambda}, \mu) = -\log(P(\mathbf{A})) - 1 + \left(\sum_{i,j=1:n} \lambda_i a_{ij} \right) + \mu = 0$$

MAXENT MODEL S.T. DEGREE BELIEFS

$$\frac{\partial}{\partial P(\mathbf{A})} L(P, \boldsymbol{\lambda}, \mu) = -\log(P(\mathbf{A})) - 1 + \left(\sum_{i,j=1:n} \lambda_i a_{ij} \right) + \mu = 0$$

So:

$$P(\mathbf{A}) = \exp(\mu - 1) \cdot \exp\left(\sum_{i,j=1:n} \lambda_i a_{ij} \right)$$

MAXENT MODEL S.T. DEGREE BELIEFS

$$\frac{\partial}{\partial P(\mathbf{A})} L(P, \boldsymbol{\lambda}, \mu) = -\log(P(\mathbf{A})) - 1 + \left(\sum_{i,j=1:n} \lambda_i a_{ij} \right) + \mu = 0$$

So:

$$\begin{aligned} P(\mathbf{A}) &= \exp(\mu - 1) \cdot \exp\left(\sum_{i,j=1:n} \lambda_i a_{ij}\right) \\ &= \frac{1}{Z(\boldsymbol{\lambda})} \cdot \exp\left(\sum_{i>j} (\lambda_i + \lambda_j) a_{ij}\right) \end{aligned}$$

MAXENT MODEL S.T. DEGREE BELIEFS

$$\frac{\partial}{\partial P(\mathbf{A})} L(P, \boldsymbol{\lambda}, \mu) = -\log(P(\mathbf{A})) - 1 + \left(\sum_{i,j=1:n} \lambda_i a_{ij} \right) + \mu = 0$$

So:

$$\begin{aligned} P(\mathbf{A}) &= \exp(\mu - 1) \cdot \exp\left(\sum_{i,j=1:n} \lambda_i a_{ij}\right) \\ &= \frac{1}{Z(\boldsymbol{\lambda})} \cdot \exp\left(\sum_{i>j} (\lambda_i + \lambda_j) a_{ij}\right) = \frac{1}{Z(\boldsymbol{\lambda})} \cdot \prod_{i>j} \exp\left((\lambda_i + \lambda_j) \cdot a_{ij}\right) \end{aligned}$$

MAXENT MODEL S.T. DEGREE BELIEFS

$$\frac{\partial}{\partial P(\mathbf{A})} L(P, \boldsymbol{\lambda}, \mu) = -\log(P(\mathbf{A})) - 1 + \left(\sum_{i,j=1:n} \lambda_i a_{ij} \right) + \mu = 0$$

So:

$$\begin{aligned} P(\mathbf{A}) &= \exp(\mu - 1) \cdot \exp\left(\sum_{i,j=1:n} \lambda_i a_{ij}\right) \\ &= \frac{1}{Z(\boldsymbol{\lambda})} \cdot \exp\left(\sum_{i>j} (\lambda_i + \lambda_j) a_{ij}\right) = \frac{1}{Z(\boldsymbol{\lambda})} \cdot \prod_{i>j} \exp\left((\lambda_i + \lambda_j) \cdot a_{ij}\right) \\ &= \prod_{i>j} \frac{\exp\left((\lambda_i + \lambda_j) \cdot a_{ij}\right)}{1 + \exp(\lambda_i + \lambda_j)} \end{aligned}$$

$P_{i,j}(a_{ij})$ Product of independent Bernoulli distributions!

Thanks to fact that prior belief constraint is on a (weighted) sum of the a_{ij}

MAXENT MODEL S.T. DEGREE BELIEFS

To find optimal values of Lagrange multipliers, solve the dual:

$$\min_{\lambda} L(P, \lambda)$$

where P is given as:

$$P(\mathbf{A}) = \prod_{i>j} \frac{\exp((\lambda_i + \lambda_j) \cdot a_{ij})}{1 + \exp(\lambda_i + \lambda_j)}$$

After some calculations:

$$\min_{\lambda} \sum_{i>j} \log(1 + \exp(\lambda_i + \lambda_j)) - \sum_{i=1:n} \lambda_i d_i$$

MAXENT MODEL S.T. DEGREE BELIEFS

$$\min_{\lambda} \sum_{i>j} \log \left(1 + \exp(\lambda_i + \lambda_j) \right) - \sum_{i=1:n} \lambda_i d_i$$

MAXENT MODEL S.T. DEGREE BELIEFS

$$\min_{\lambda} \sum_{i>j} \log \left(1 + \exp(\lambda_i + \lambda_j) \right) - \sum_{i=1:n} \lambda_i d_i$$

Can be solved using **gradient** descent:

$$\frac{\partial}{\partial \lambda_k} \left(\sum_{i>j} \log \left(1 + \exp(\lambda_i + \lambda_j) \right) - \sum_{i=1:n} \lambda_i d_i \right)$$

$$= \sum_{i=1:n} \frac{\exp(\lambda_i + \lambda_k)}{1 + \exp(\lambda_i + \lambda_k)}$$

$$d_k$$

Expected degree of node k

Required expected degree of node k

Lots of computational speed-ups possible...

TAKE-AWAYS

- Constraints on the expected value of weighted averages:

$$E_{A \sim P} \left\{ \sum_{i,j \in I} f_{ij} a_{ij} \right\} = c$$

where f_{ij} and c are constants and I is a set of indices,
lead to convenient **product distributions**

- **Other examples for graphs:**
 - Overall density (trivial)
 - Densities of particular blocks (e.g. block of nodes with same affiliation)
 - Assortativity (approximately)
 - ...

THE INTERESTINGNESS

INFORMATION CONTENT

- **Information content:**

$$\text{InformationContent}(\text{pattern}, P) = -\log(P(\text{pattern}))$$

- **Pattern:** “number of edges between given set of nodes $W \subseteq V$ is larger than or equal to a specified k_W ”

- A bit tricky...

- **Cliques as a special case:** “set of nodes $W \subseteq V$ forms a clique”. Then:

$$P(\text{pattern}) = \prod_{i>j \in W} P_{i,j}(1)$$

- So:

$$\begin{aligned} \text{InformationContent}(\text{pattern}, P) &= -\log(P(\text{pattern})) \\ &= -\sum_{i>j \in W} \log(P_{i,j}(1)) \end{aligned}$$

- Larger if $|W|$ is larger and if $P_{i,j}(1)$ for $i, j \in W$ is smaller

INFORMATION CONTENT

- **Pattern (general case):** “number of edges between given set of nodes $W \subseteq V$ is larger than or equal to a specified k_W ”

- Probability of at least k_W successes in $n_W = \binom{|W|}{2}$ Bernoulli trials?

- Approximated by:

$$P(\text{pattern}) \approx \exp\left(-n_W \text{KL}\left(\frac{k_W}{n_W} \parallel p_W\right)\right)$$

- And thus:

$$\begin{aligned} \text{InformationContent}(\text{pattern}, P) &= -\log(P(\text{pattern})) \\ &\approx n_W \text{KL}\left(\frac{k_W}{n_W} \parallel p_W\right) \end{aligned}$$

where p_W is the average probability $P_{i,j}(1)$ for the edges between $i, j \in W$

- Larger if $|W|$ (and thus n_W) is larger, p_W is smaller, and k_W is larger

DESCRIPTION LENGTH

- For **cliques**:

- Describe set W

- A constant (to describe $|W|$)

- plus a linear term in $|W|$ (to describe its elements):

$$\text{DescriptionLength}(\text{pattern}) = \alpha|W| + \beta$$

- For **dense subgraphs**:

- Constant β also describes k_W

INTERESTINGNESS

Putting things together:

$$\text{Interestingness}(\text{pattern}, P) = \frac{-\sum_{i>j \in W} \log(P_{i,j}(1))}{\alpha|W| + \beta}$$

- A bit more complex for general dense subgraphs:

$$\text{Interestingness}(\text{pattern}, P) \approx \frac{n_W \text{KL}\left(\frac{k_W}{n_W} \parallel p_W\right)}{\alpha|W| + \beta}$$

- Hard to optimize!
 - Exact search for small graphs
 - Effective hill climber for large graphs

TAKE-AWAYS

- No compromises w.r.t. interestingness
- Often leads to hard search problems
- **Is this intrinsic to subjective interestingness?**

UPDATING THE BACKGROUND DISTRIBUTION

UPDATING THE BACKGROUND DISTRIBUTION

**Given a pattern,
update the background distribution by conditioning on the
pattern**

- Easy to do **for cliques** W :
 - Set $P'_{i,j}(a_{ij}) = 1$ for $a_{ij} = 1$ if $i, j \in W$
- Fast to approximate **for (non-clique) dense subgraphs** W :
 - Set $P'_{i,j}(a_{ij}) \propto P_{i,j}(a_{ij}) \cdot \exp(\lambda_W a_{ij})$ if $i, j \in W$
such that the expected density of W is k_W
- Remains a product of Bernoulli's

TAKE-AWAYS

- Updating can be trivial
- Otherwise, often easy to do approximately