

Clustering II

The CURE Algorithm

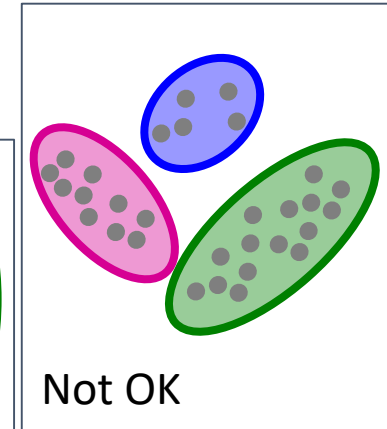
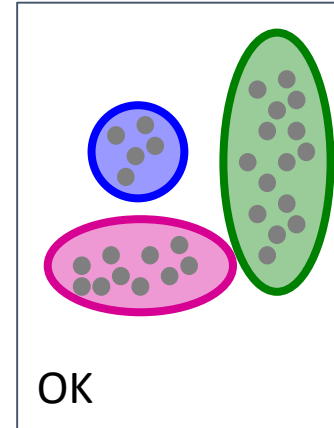
Recap

- Transition (Random Walk) Matrix
- Convergence of Random Walk
- Graph Laplacian,
- Graph Spectra, Eigen gap,
- Structure of the Web,
- PageRank,
- Topic-Specific PageRank
- **PageRank with Restarts**
- **Hierarchical Clustering**
- **K-means**
- **BFR algorithm**

The CURE Algorithm

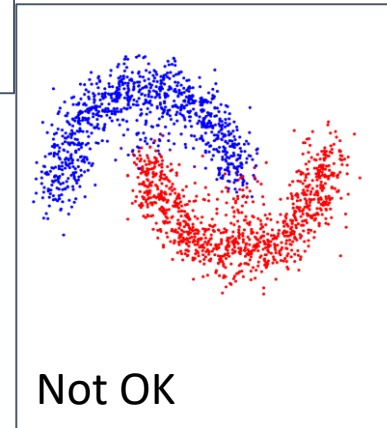
- **Problem with BFR? Any ideas?**

- Assumes clusters are normally distributed in each dimension
- And axes are fixed – ellipses at an angle are **not OK**



- **CURE (Clustering Using REpresentatives):**

- Assumes a Euclidean distance
- Allows clusters to assume any shape
- **Underlying idea:** Uses a collection of representative points to represent clusters
 - Instead of a centroid and std. deviation as in BFR



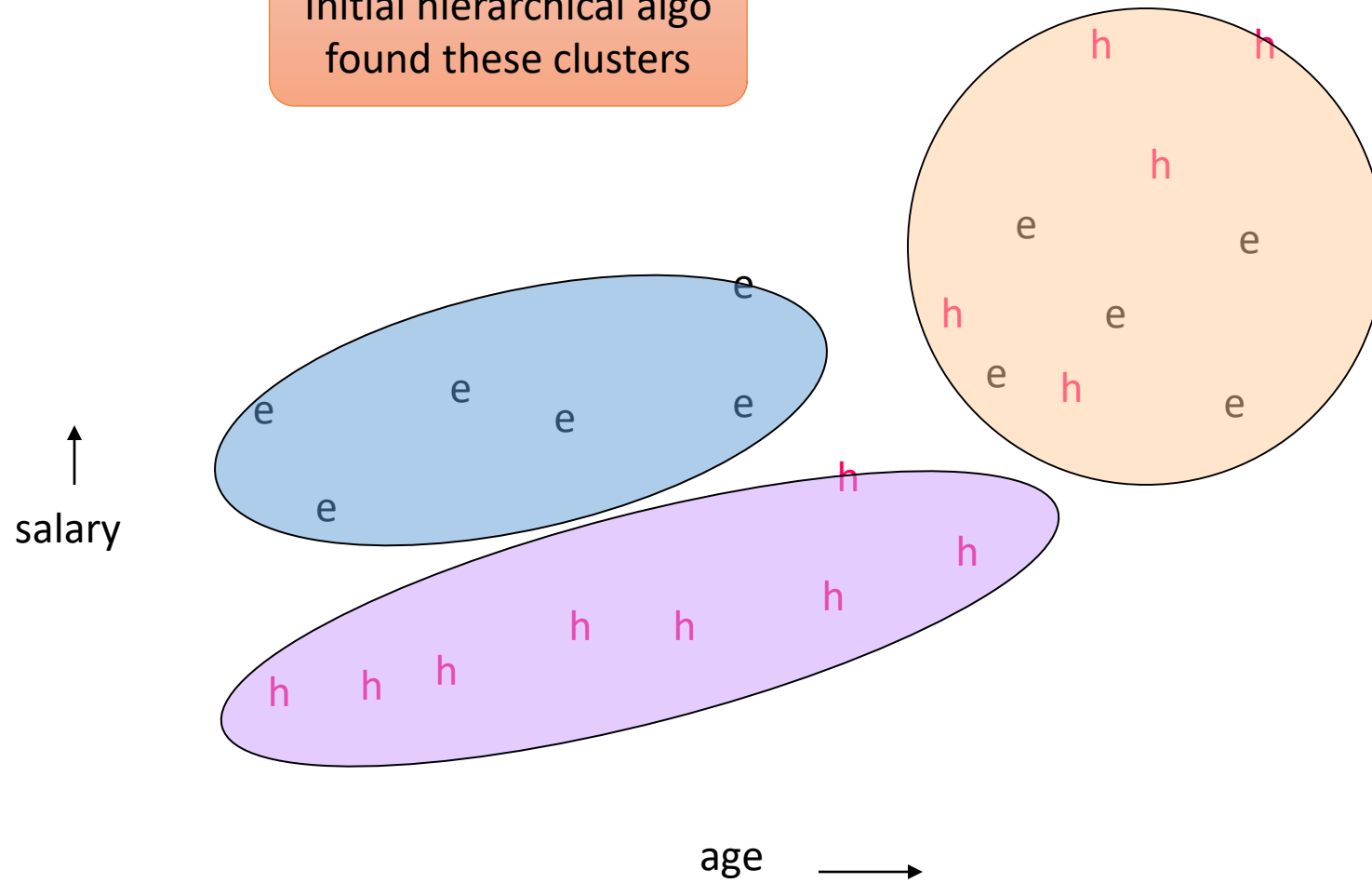
Starting CURE

Two Pass algorithm. Pass 1:

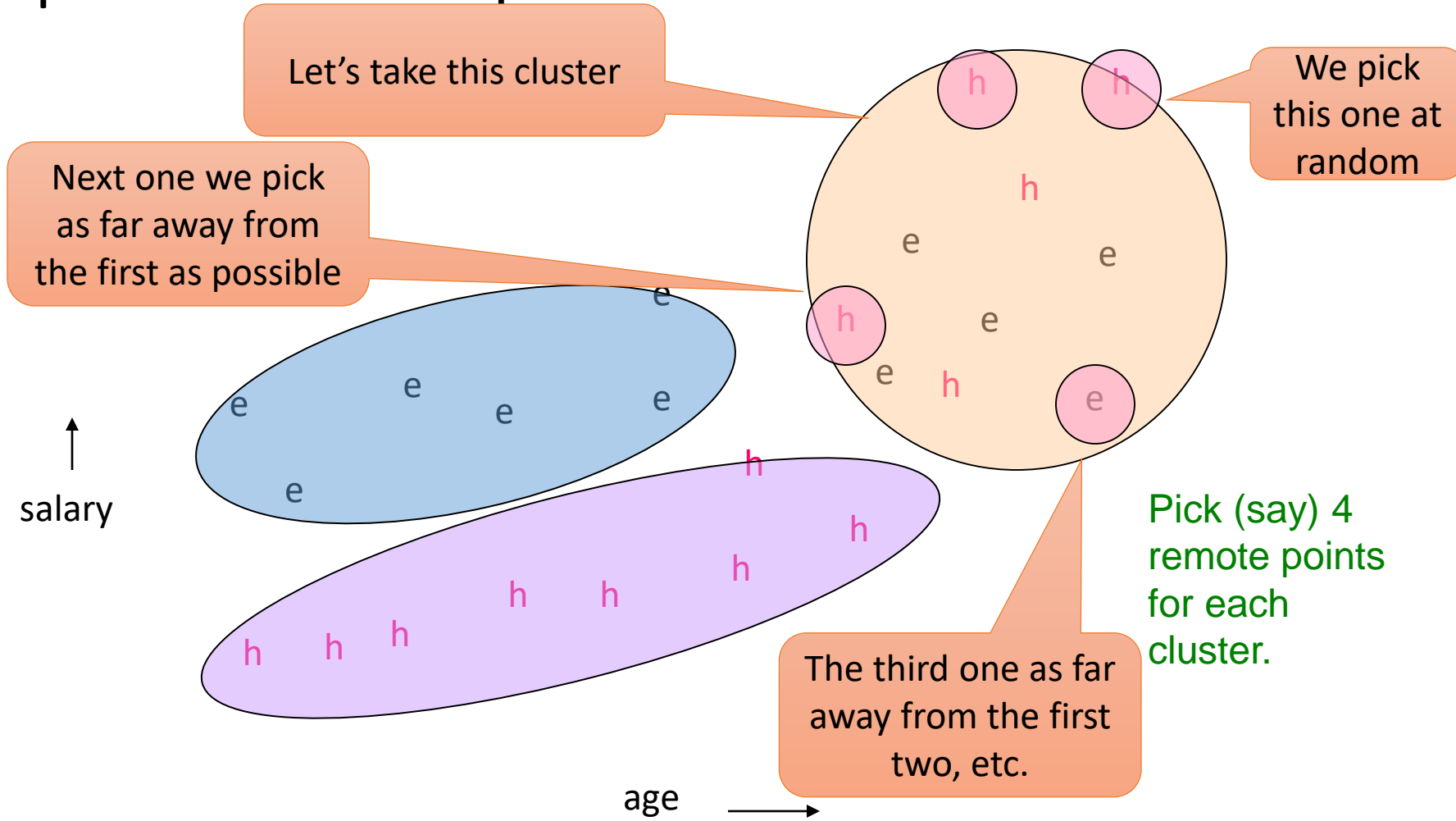
- **0) Pick a random sample of points that fit in main memory**
- **1) Initial clusters:**
 - Cluster these points e.g., hierarchically – group nearest points/clusters
 - Hierarchical algos can find clusters of any shape that way
- **2) Pick representative points:**
 - For each cluster, **pick k sample points**, as dispersed as possible (as discussed in previous techniques)
 - From the sample, pick **representatives (synthetic points)** by moving them (say) 20% toward the centroid of the cluster

Example: Initial Clusters

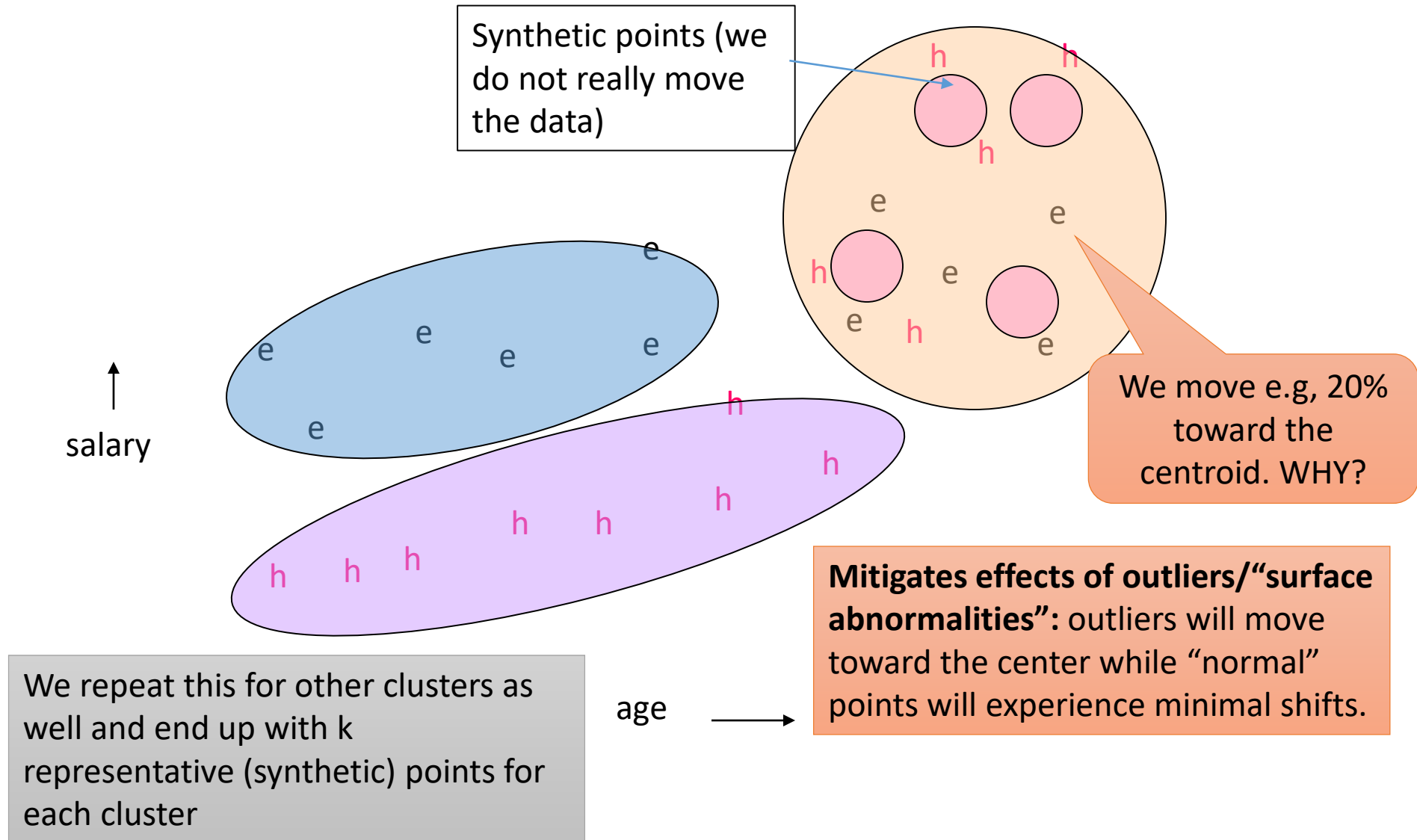
Initial hierarchical algo
found these clusters



Example: Pick Dispersed Points



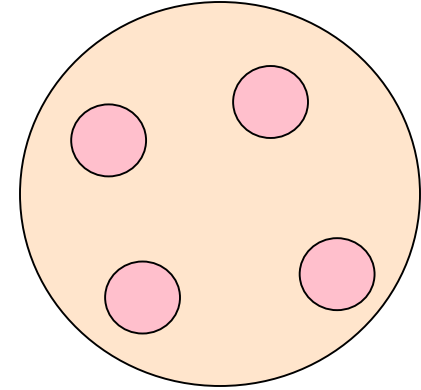
Example: Pick Dispersed Points



Finishing CURE

Pass 2:

- Now, rescan the whole dataset and visit each point p in the data set
- Place it in the “**closest cluster**”
 - Normal definition of “**closest**”:
Find the closest representative to p and assign it to representative's cluster



p

Summary

- **Goal of Clustering:** Given a **set of points**, with a notion of **distance** between points, **group the points** into some number of *clusters*
- **Algorithms:**
 - Agglomerative **hierarchical clustering**:
 - Centroid and clustroid
 - Issues with scalability!
 - ***k*-means**:
 - Initialization, picking k
 - Number of rounds to converge?
 - **BFR**
 - One pass algo! Strong assumptions on the data!
 - **CURE**
 - Overcomes limitations of BFR

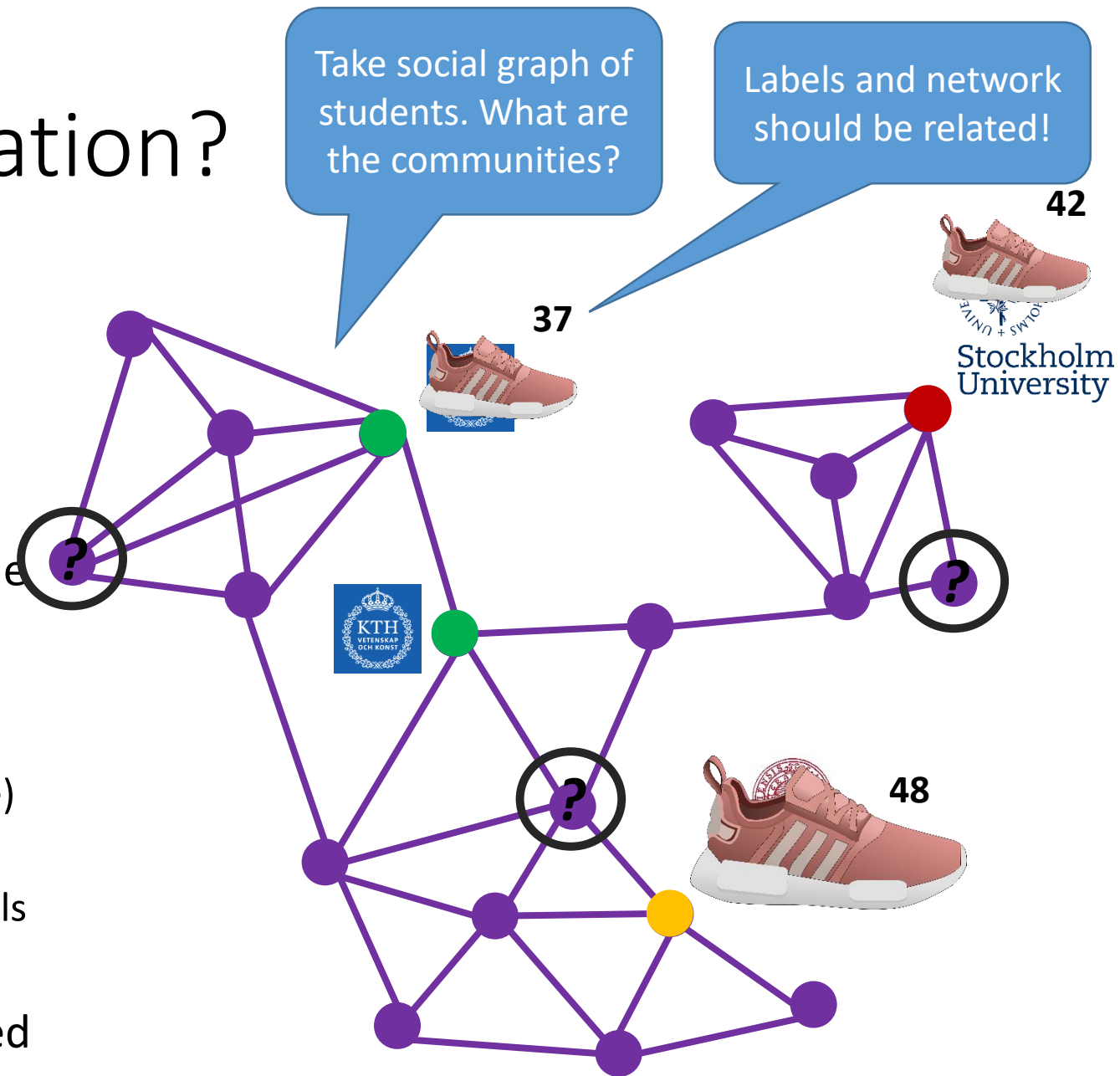
Kahoot! time

Label Propagation and Graph Spectra

What is Label Propagation?

- **Label propagation**

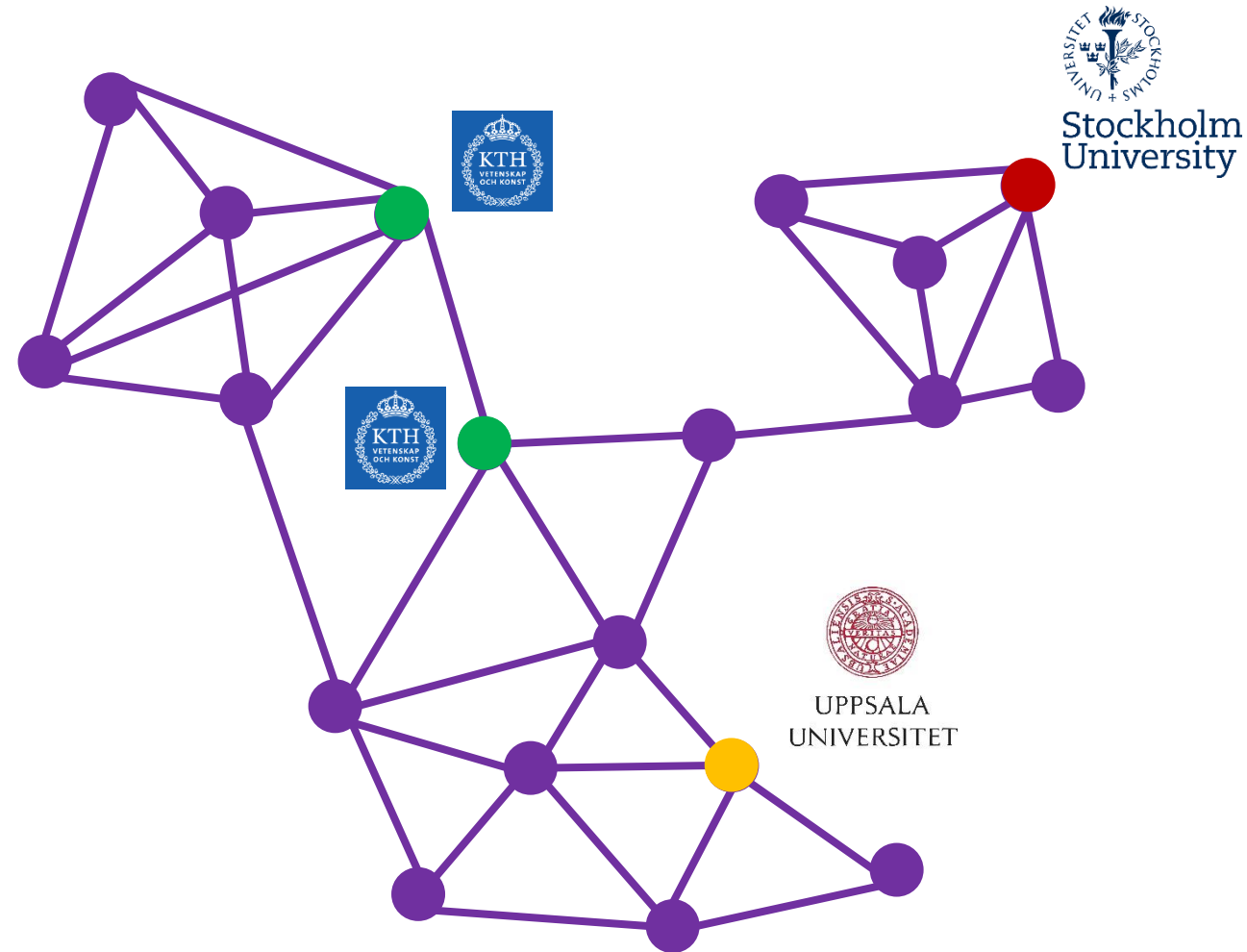
- Some nodes are labeled (ground truth)
 - Categorical/numeric/binary values
- Task: Predict/Label for the rest of the nodes in the graph
- Key assumption:
 - We assumed linked nodes are correlated (e.g., homophily, influence)
 - Labels propagate only on edges
 - You should have “enough” initial labels
- Semi-supervised Learning
- Why don't we just go for a supervised learning?



Network Classification

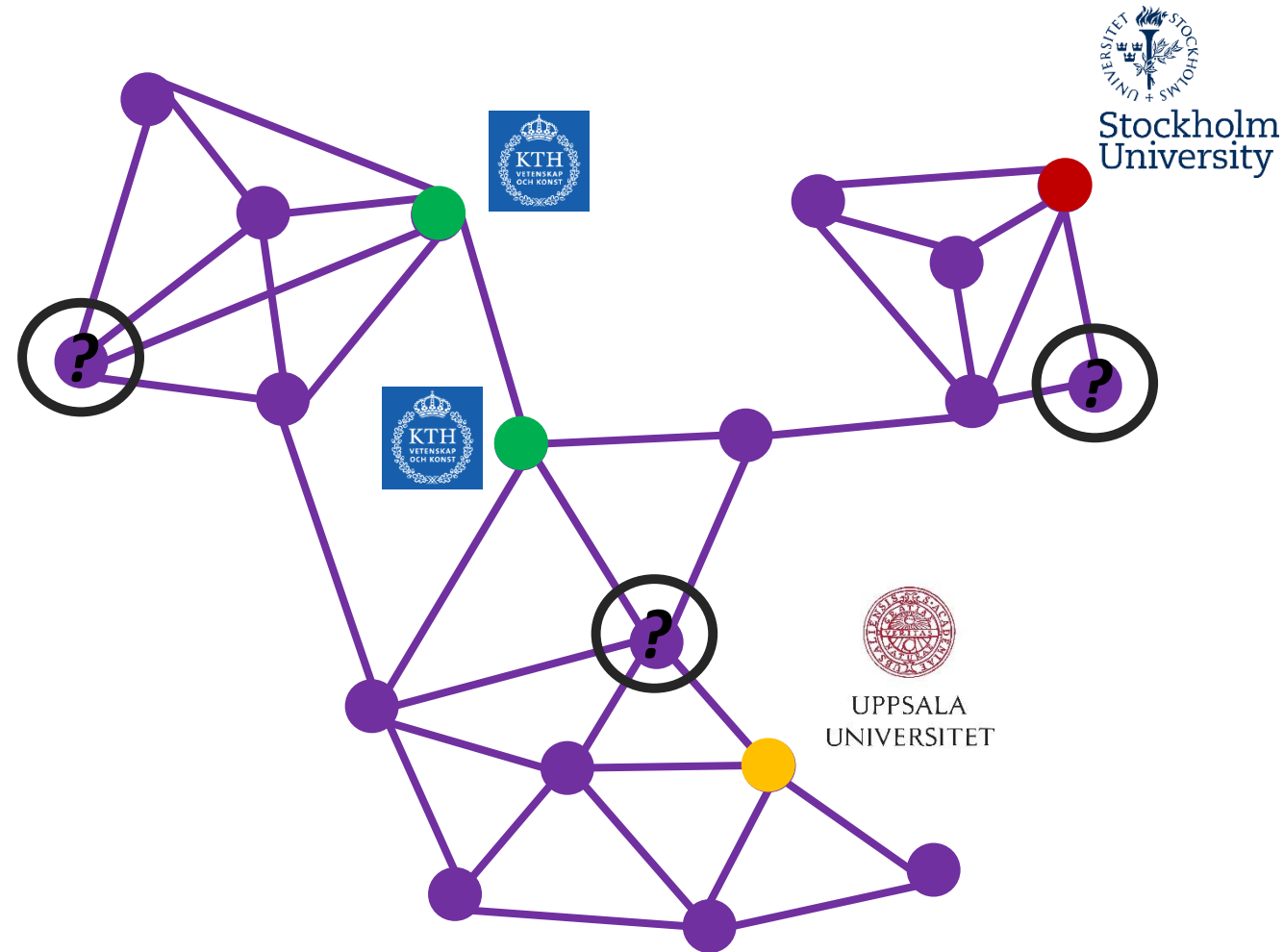
- **Label propagation**

- Given graph G with nodes V , where
 - Nodes V_l given labels Y_l
 - Nodes V_u do not have labels
- We need to find Y_u
- Labels can be
 - Binary
 - Multi-class
 - Real values
- Why?
 - Labelling/annotating data is expensive,
 - Small amounts of labelled data, large amounts of unlabelled data.



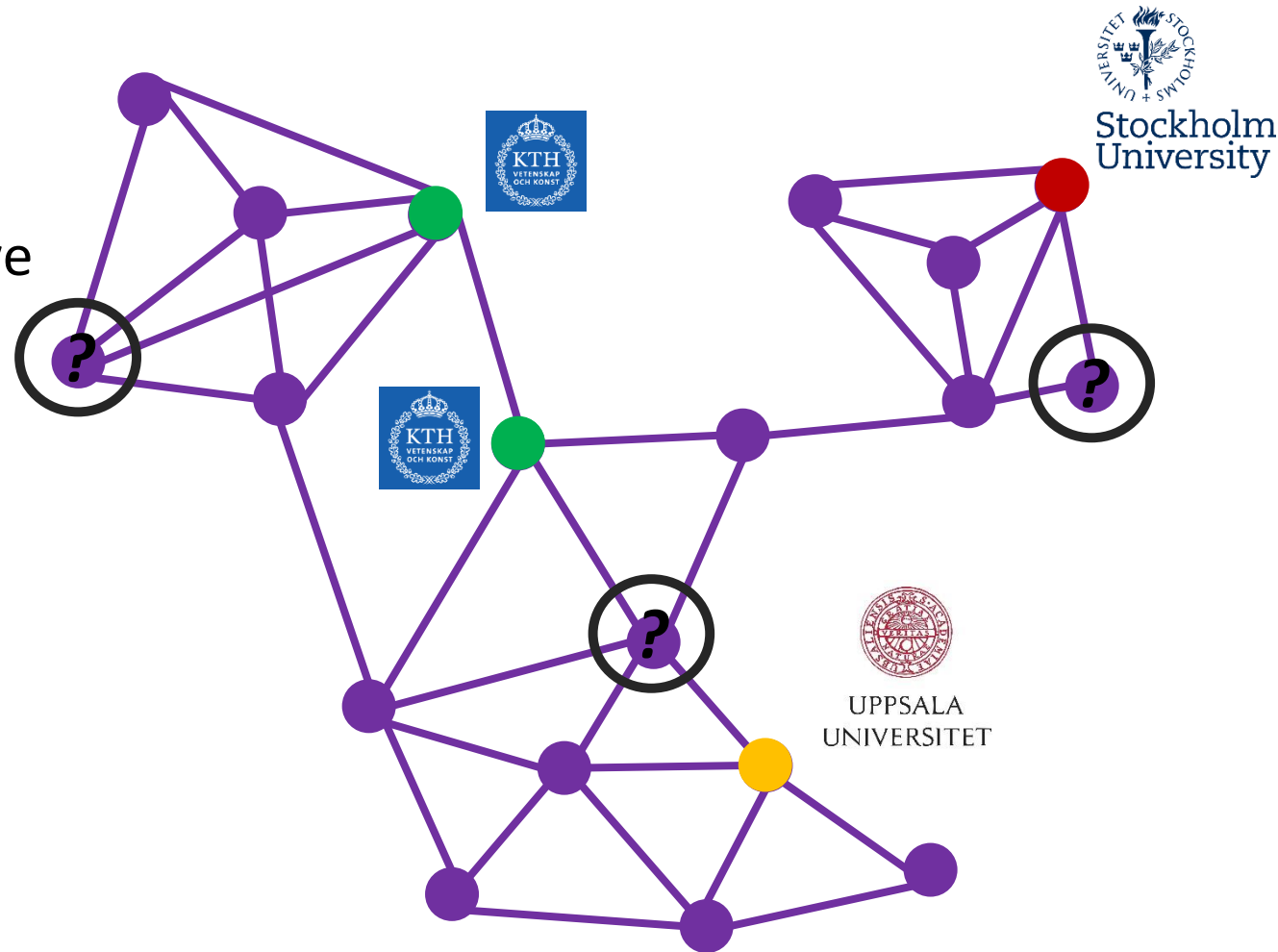
How to predict the labels?

- Ideas?
- Example:
 - Look at the local neighborhood
 - See what's the dominant label.
 - Or mean, Or average, or any other ML classifier.
 - Adopt that label.
 - What if your neighborhood does not have any labels?
 - Wait until it “propagates” to you.



Label Propagation

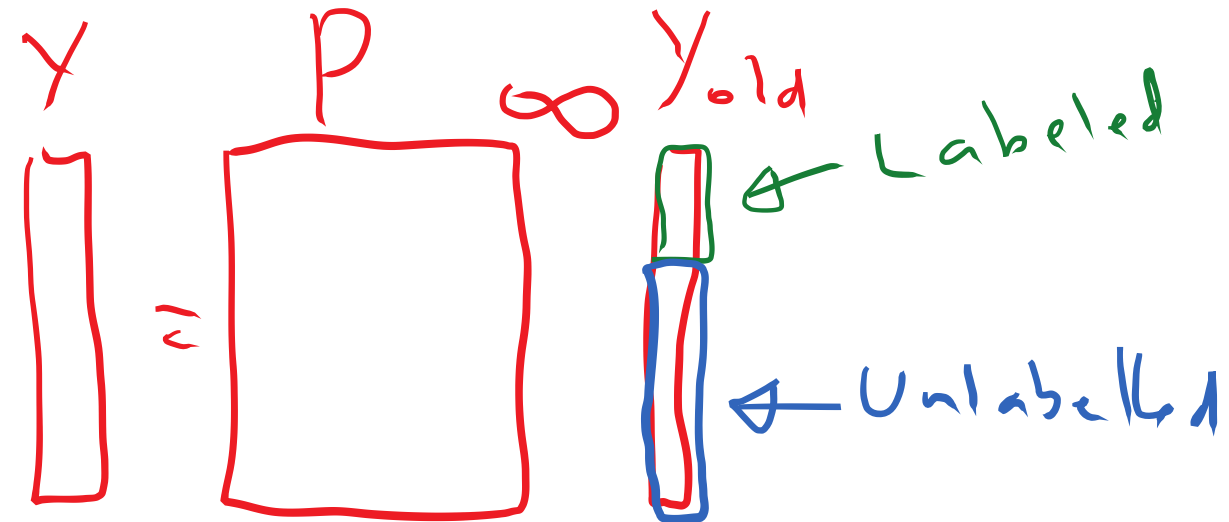
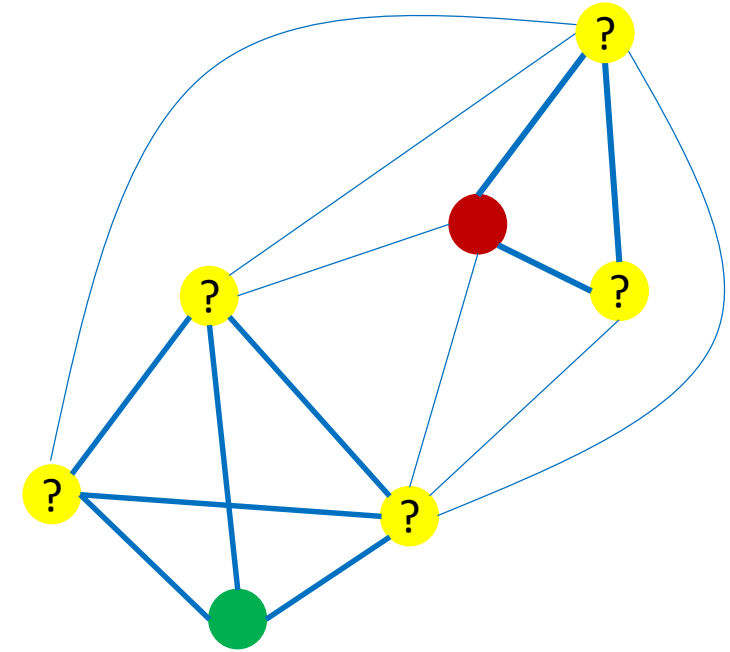
- We start with a graph and partially labeled nodes
 - If we do not have a graph, then we construct it by linking data points based on their **similarity**.
- We apply label propagation through that graph to predict the rest of the nodes
- Typical – semisupervised learning using graphs.



Label Propagation through Random-Walks with absorbing states

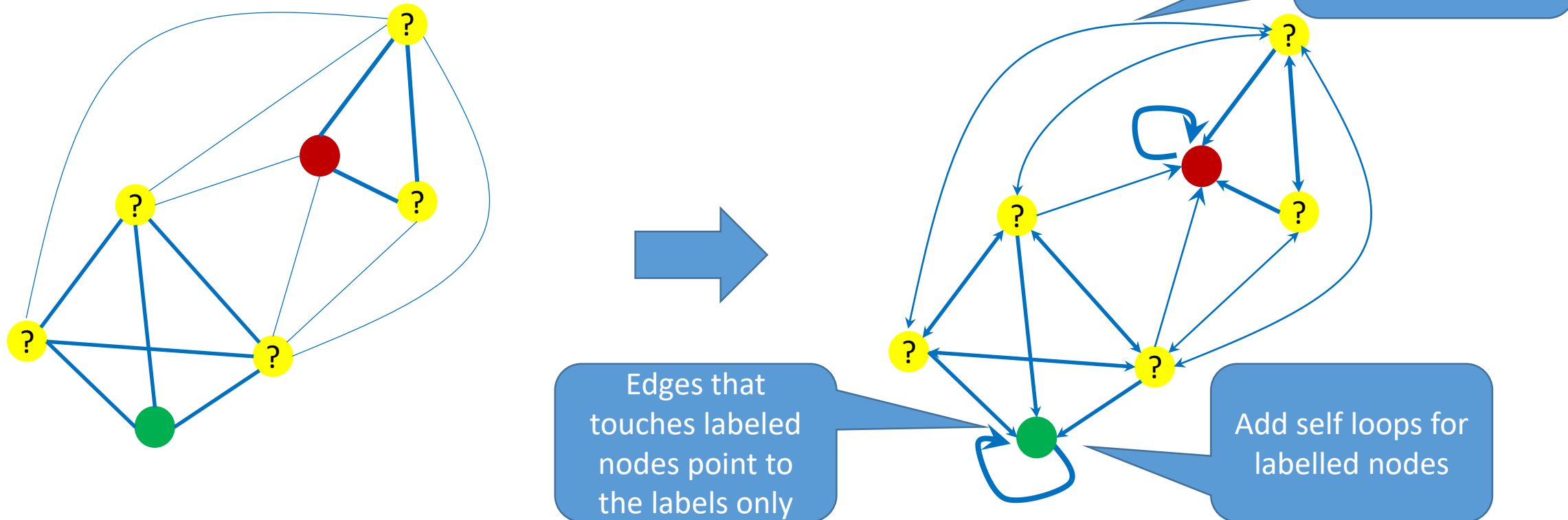
- **Idea:**

- You start **many random walkers** from an **unlabeled** node
- Keep random walking until you **hit a labeled node**.
 - Then your random walker gets stuck there – it's a trap!
- Check which random walkers got stuck in which labelled node.
 - E.g., out of 100 random walkers 85 got stuck in Green and 15 at red.
- **Adopt** the dominant label
 - i.e., we want to calculate a probability for a given node that a random walker will end up in a particular label.
- Can we do that through matrix multiplication?
 - Remember Page rank?
 - So what should be our P so that it represents the above idea?

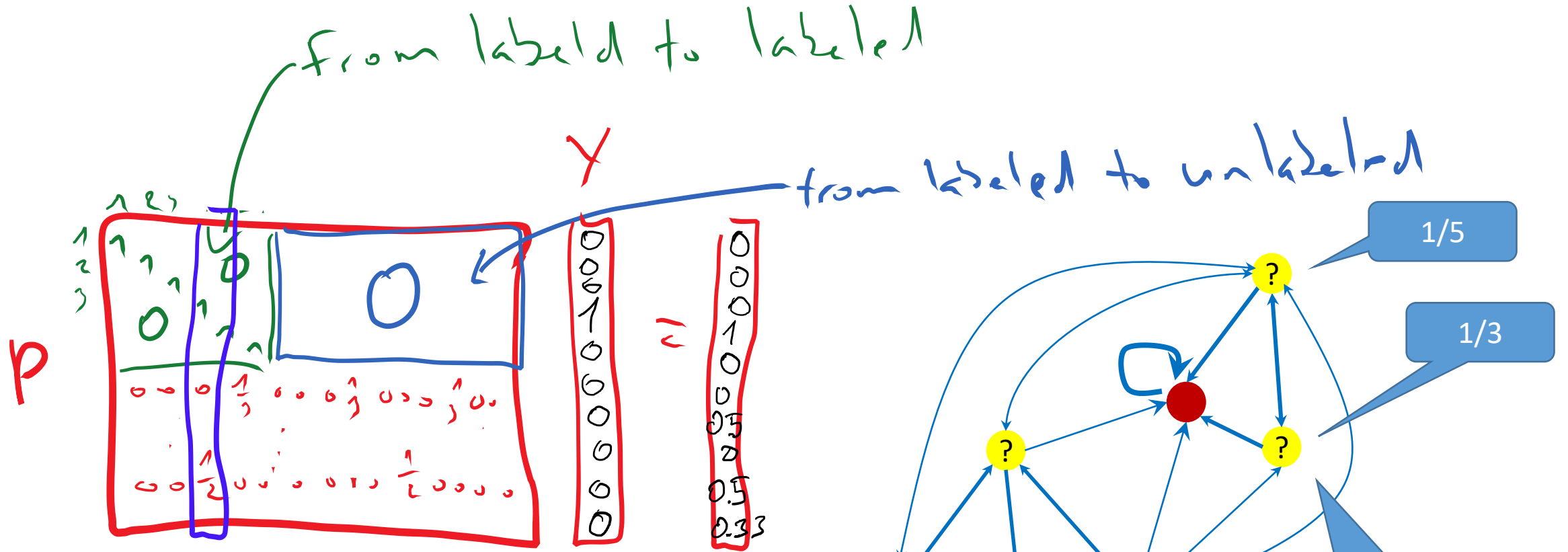


Changing the graph

- How to represent “getting stuck” on a particular label?

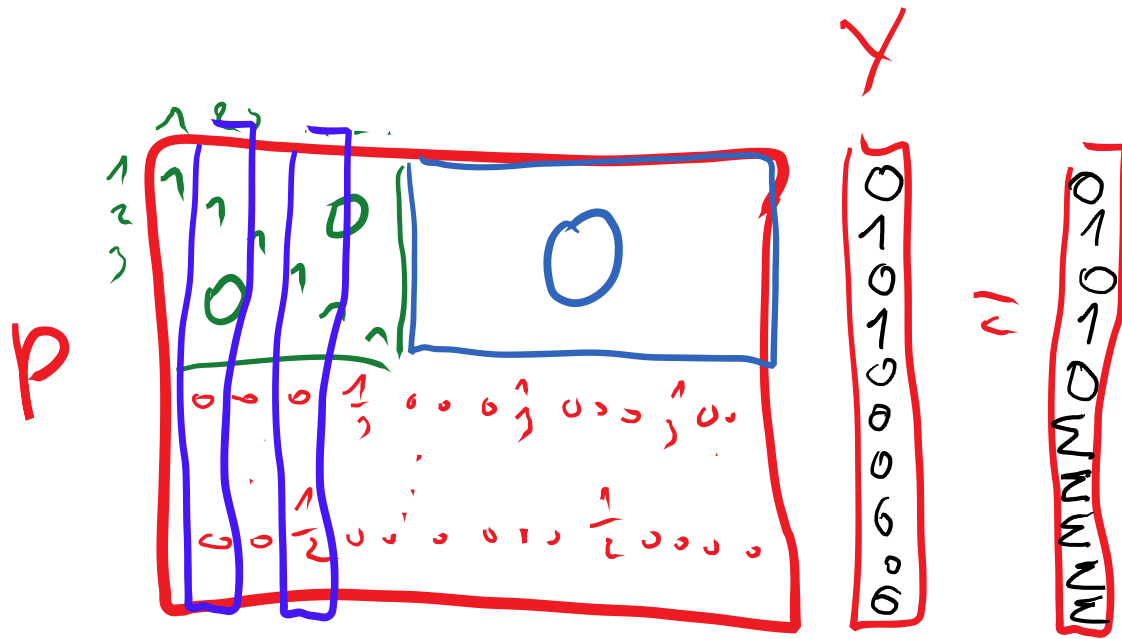


Absorption Matrix P

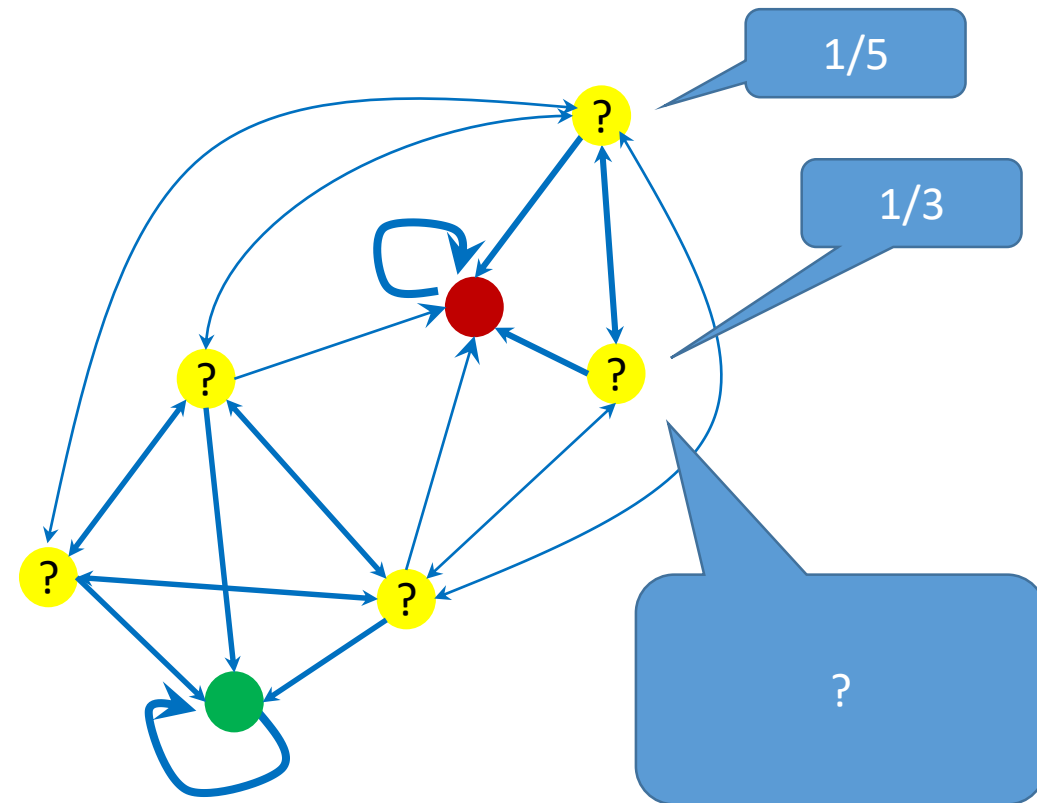


- Row-stochastic matrix P
- What is a chance to hit labelled node with a Random Walk with $TTL=1$?

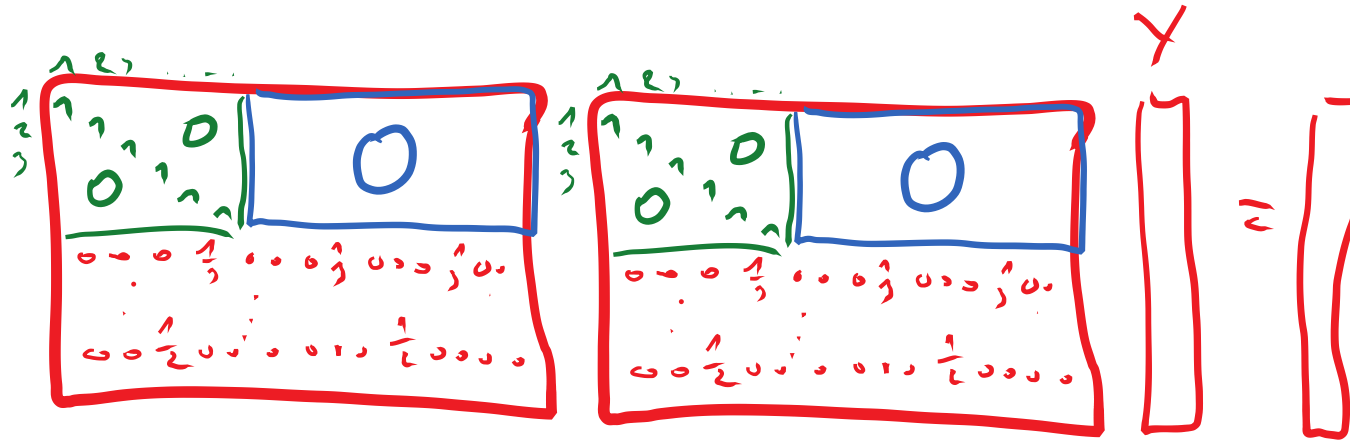
Absorption Matrix P – multiple labels of the same “color”



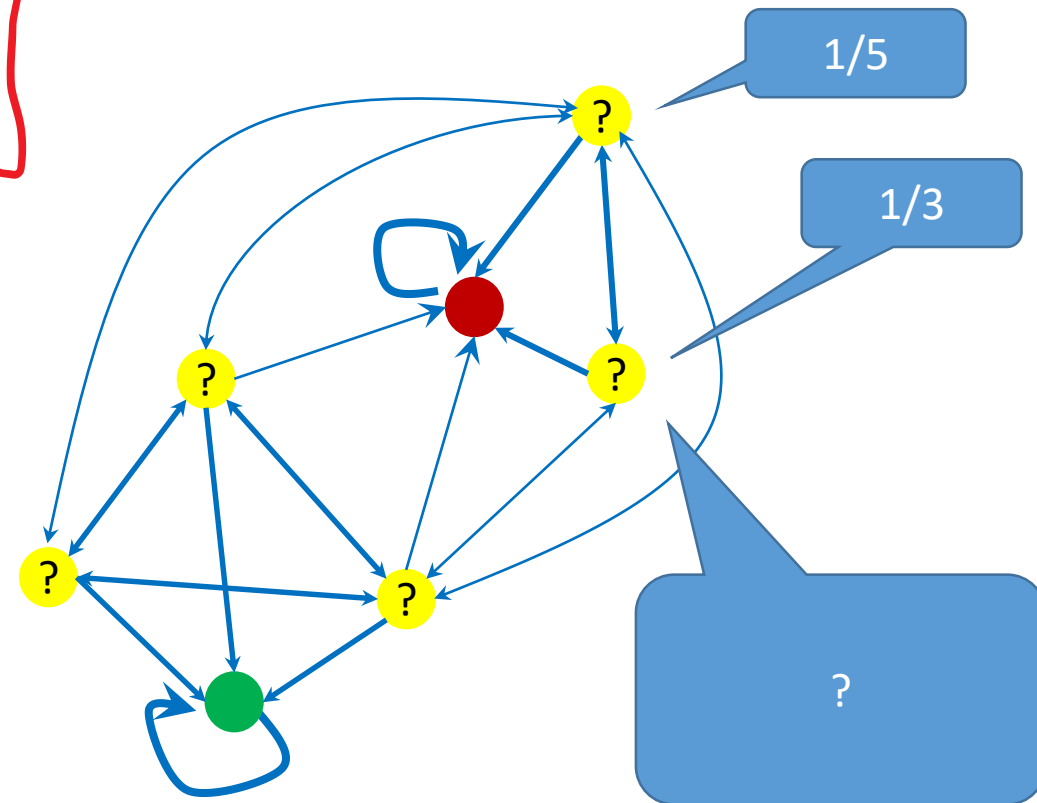
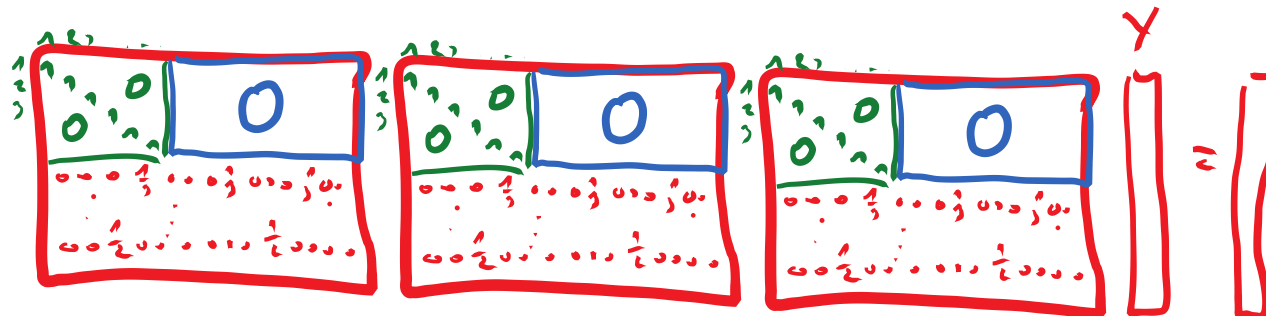
- Row-stochastic matrix P
- What is a chance to hit labelled node with a Random Walk with TTL=1?



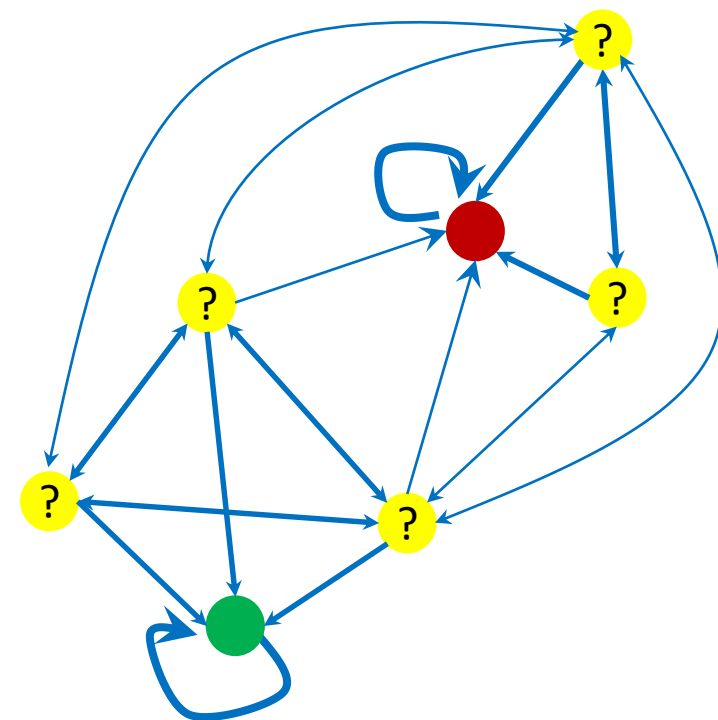
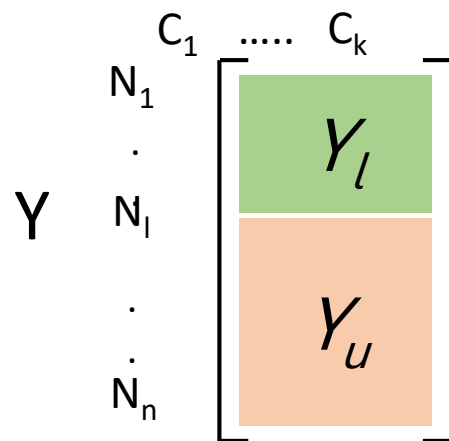
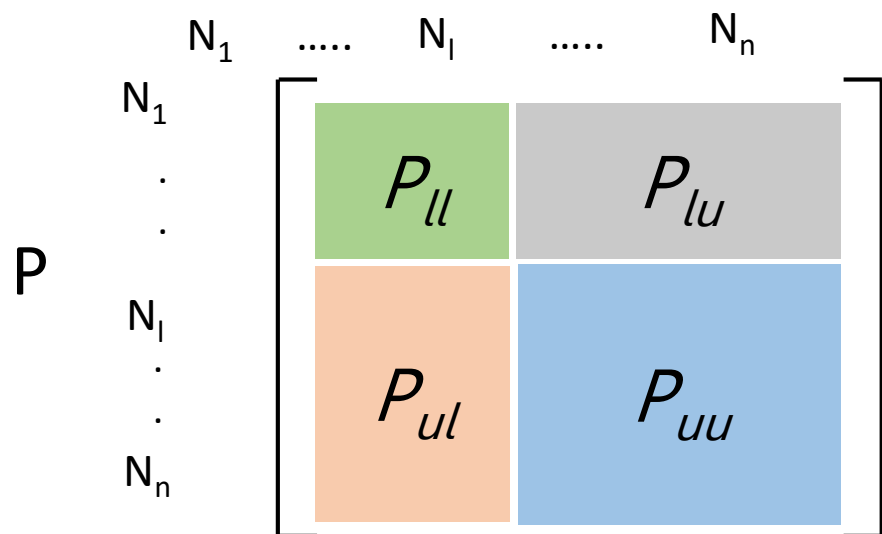
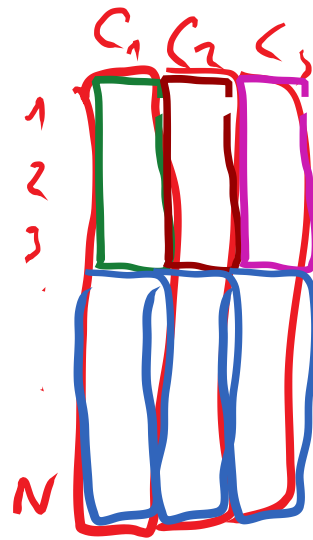
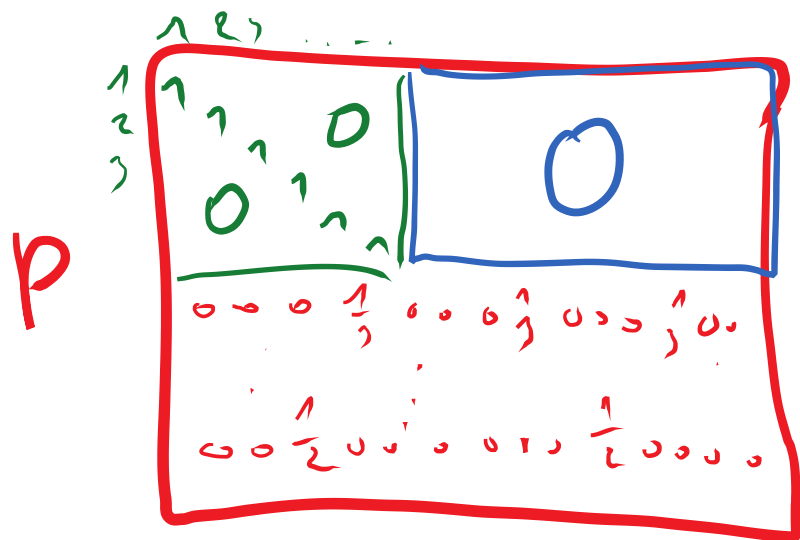
What is a chance to hit labelled node with a Random Walk with TTL=2?



- What is a chance to hit labelled node with a Random Walk with TTL=3?



Dealing with labels of several “colors”



How (does) it converge?

- What happens to \mathbf{P}^∞ ?

$$\mathbf{P} = \begin{matrix} & \begin{matrix} N_1 & \dots & N_l & \dots & N_n \end{matrix} \\ \begin{matrix} N_1 \\ \vdots \\ N_l \\ \vdots \\ N_n \end{matrix} & \begin{bmatrix} P_{ll} & P_{lu} \\ P_{ul} & P_{uu} \end{bmatrix} \end{matrix}$$

$$\mathbf{P} = \begin{matrix} & \begin{matrix} N_1 & \dots & N_l & \dots & N_n \end{matrix} \\ \begin{matrix} N_1 \\ \vdots \\ N_l \\ \vdots \\ N_n \end{matrix} & \begin{bmatrix} I & 0 \\ \text{something} & \text{something} \end{bmatrix} \end{matrix}$$

$$\mathbf{P}^\infty = \begin{matrix} & \begin{matrix} N_1 & \dots & N_l & \dots & N_n \end{matrix} \\ \begin{matrix} N_1 \\ \vdots \\ N_l \\ \vdots \\ N_n \end{matrix} & \begin{bmatrix} I & 0 \\ (I - P_{uu})^{-1} P_{ul} & 0 \end{bmatrix} \end{matrix}$$

$$(I - P_{uu})^{-1} P_{ul}$$

What will be the labels?

$$\begin{matrix} & C_1 & \dots & C_k \\ \begin{matrix} N_1 \\ \vdots \\ N_l \\ \vdots \\ N_n \end{matrix} & \begin{bmatrix} \text{green} \\ \text{orange} \end{bmatrix} \end{matrix} \begin{matrix} Y_l \\ Y_u \end{matrix}$$

Y

Probabilities that a random walker from a node i will get stuck on a label C_j

$=$

$$\begin{matrix} & N_1 & \dots & N_l & \dots & N_n \\ \begin{matrix} N_1 \\ \vdots \\ N_l \\ \vdots \\ N_n \end{matrix} & \begin{bmatrix} \text{green} & \text{grey} \\ \text{orange} & \text{blue} \end{bmatrix} \end{matrix} \begin{matrix} I & 0 \\ (I - P_{uu})^{-1} P_{ul} & 0 \end{matrix}$$

P^∞

The only "problem" is to calculate this matrix inverse.

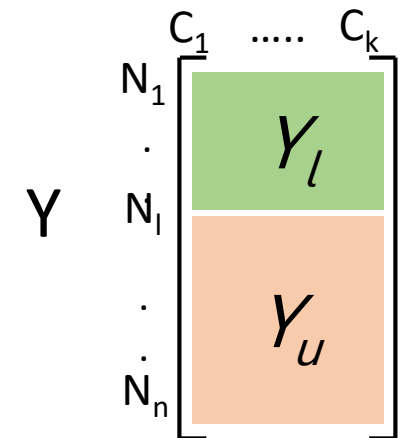
$$\begin{matrix} & C_1 & \dots & C_k \\ \begin{matrix} N_1 \\ \vdots \\ N_l \\ \vdots \\ N_n \end{matrix} & \begin{bmatrix} \text{green} \\ \text{orange} \end{bmatrix} \end{matrix} \begin{matrix} Y_l \\ Y_u \end{matrix}$$

Y_{initial}

Matrix Y

Y is a matrix, holding the soft probabilities of each instance

- N_1, N_2, \dots, N_n represent instances (i.e., nodes in graph)
- C_1, C_2, \dots, C_k represent the possible label
- Y_{ab} represent the probability of N_a being labeled as C_b
- For final labelling one has to account for “class proportions”!
 - Often highly disablnanced labels.
 - Class balance might be known apriori
 - Probabilities can be weighted with the estimated size proportion of that category.
 - Class mass normalization!



Y

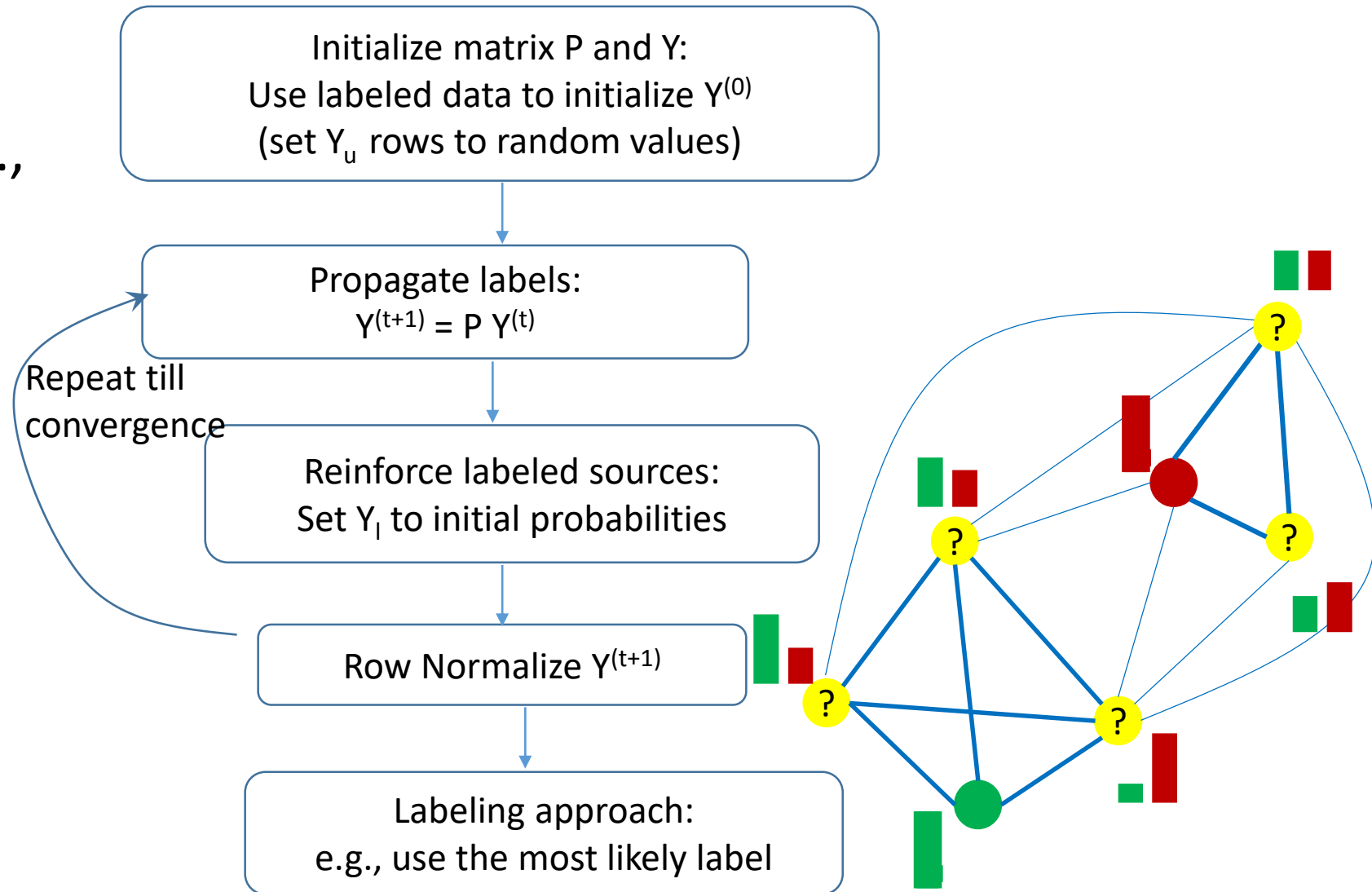
	Spam	Non-spam	
1	1	0	
2	1	0	
3	0	1	
4	0	1	
5	0.2	0.8	non-spam
6	0.4	0.6	non-spam
100	0.7	0.3	spam

Alternative Solution (original Label Propagation)

- Alternative solution on the “original graph”, i.e., original random walk matrix

- $D_{ii} = \sum_j A_{ij}$
- $P = D^{-1} A$

- Converges to the same thing



Optimization: Smoothing P using PageRank

- Creating a uniform transition matrix T , such that:

$$t_{ij} = \frac{1}{N}$$

- Update the probability transition matrix P , such that:

$$P = \epsilon T + (1 - \epsilon) P ,$$

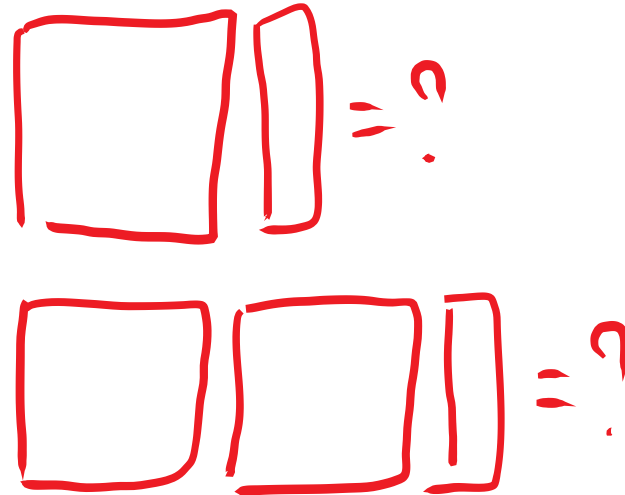
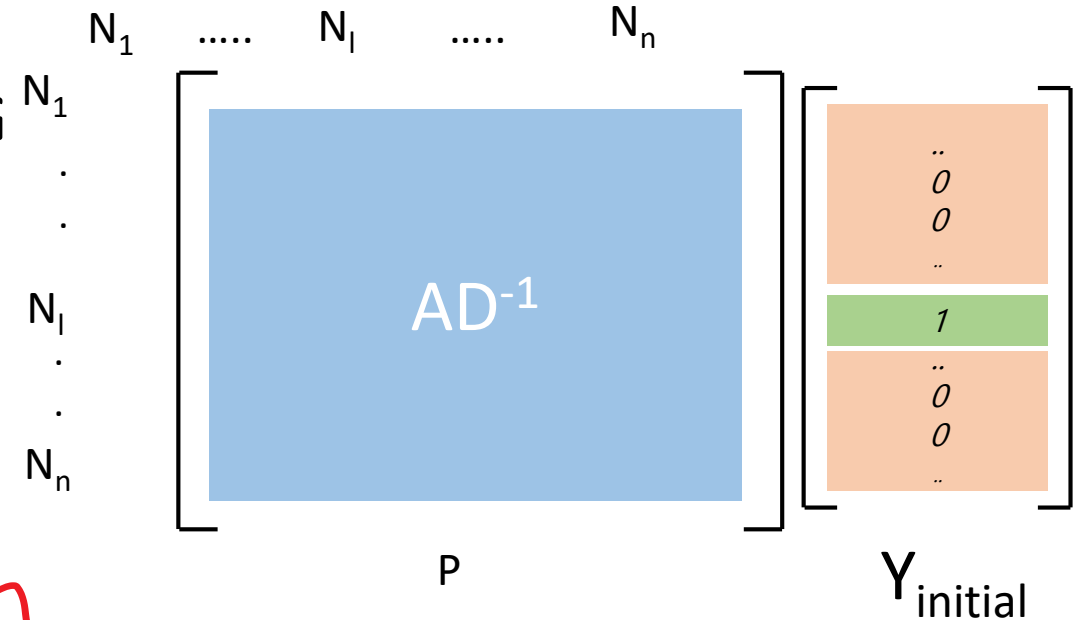
where ϵ is a parameter in $(0,1)$

Things to remember

- **Label propagation: Network Classification**
 - Initial labeling is very important for having meaningful classes
 - Initially labelled nodes can never change their label!
 - Is it good or bad?
 - What if you have noisy labelled data?
 - That brings us to Diffusion
- **Label propagation: Community Detection**

Finding Communities with Diffusion

- A – Adjacency matrix of undirected graph G
- D – Degree matrix of graph G
- P – random walk Matrix $P = AD^{-1}$
- Let's start with a single label Y_{initial}
- What happens for $P * Y_{\text{initial}}$?
 - $P * P * Y_{\text{initial}}$?

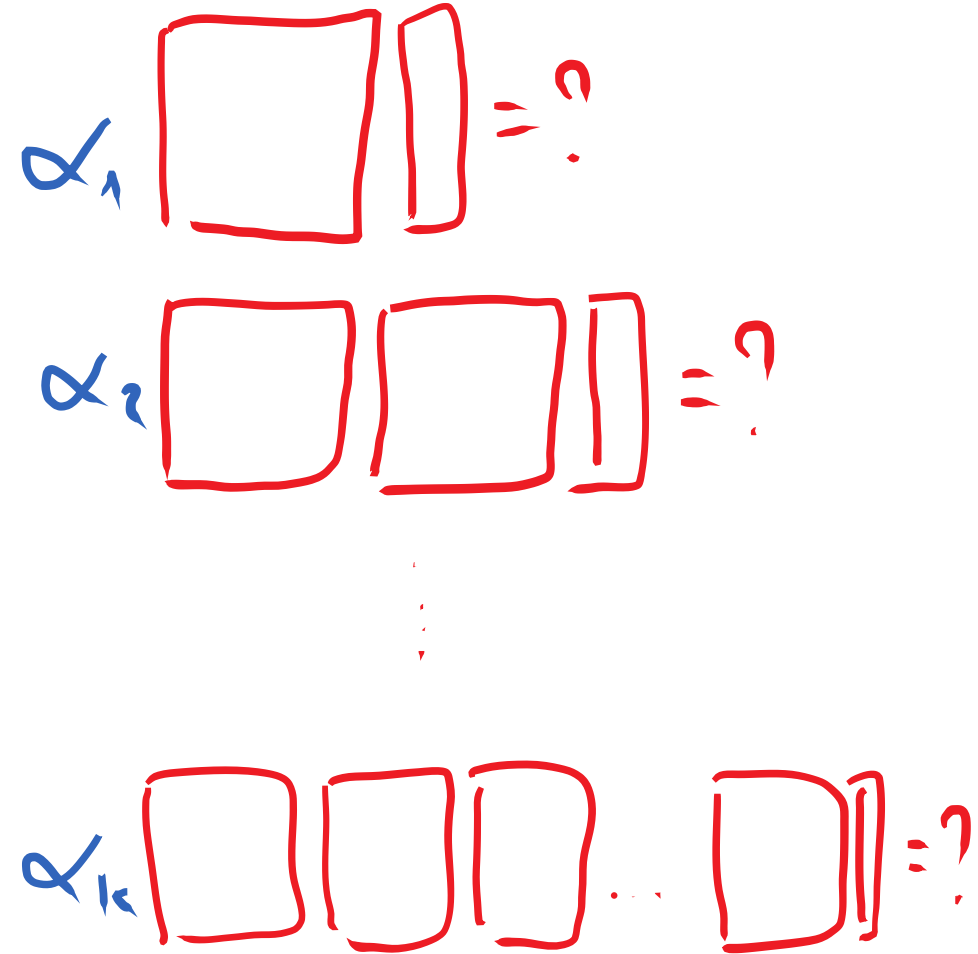


Finding Communities with Diffusion

- What happens for $P * Y_{initial}$?
 - $P * P * Y_{initial}$?
 - $P * \dots * P * Y_{initial}$?
- Graph diffusion is a sum:
 - $Y_{dif} = \sum_{k=0}^{\infty} \alpha_k P^k Y_{initial}$
 - Alpha provides decaying weight
 $\sum_{k=0}^{\infty} \alpha_k = 1$
 - High values in Y_{dif} would indicate the belonging to a specific cluster
- Best known instance – personalized PageRank diffusion.

$$(1 - \alpha) \sum_{k=0}^{\infty} \alpha^k P^k Y_{initial}$$

\sum



Label Spreading (Zhou 2004)

- Input: Graph $G(V; E)$
- A – Adjacency matrix of graph G
- D – Degree matrix of graph G
- L – normalized Laplacian $L = D^{-1/2} A D^{-1/2}$
- If your matrix is not too big and you can calculate the inverse then the solution:
 - $Y_{\text{converged}} = (1-\alpha)(I - \alpha L)^{-1} Y_{\text{initial}}$
- Steps of the algo:
 - Initialize Y_{initial} (non labelled nodes – zeros)
 - Repeat
 - $Y_{t+1} \leftarrow \alpha L Y_t + (1-\alpha) Y_{\text{initial}}$
 - until Y_t converges;

As if you'd restart from the initial seeds and “reinforce” initial rounds.