

Lecture 2: Spatial Graph Convolution and its Theoretical Performance on Simple Random Data, Part 1

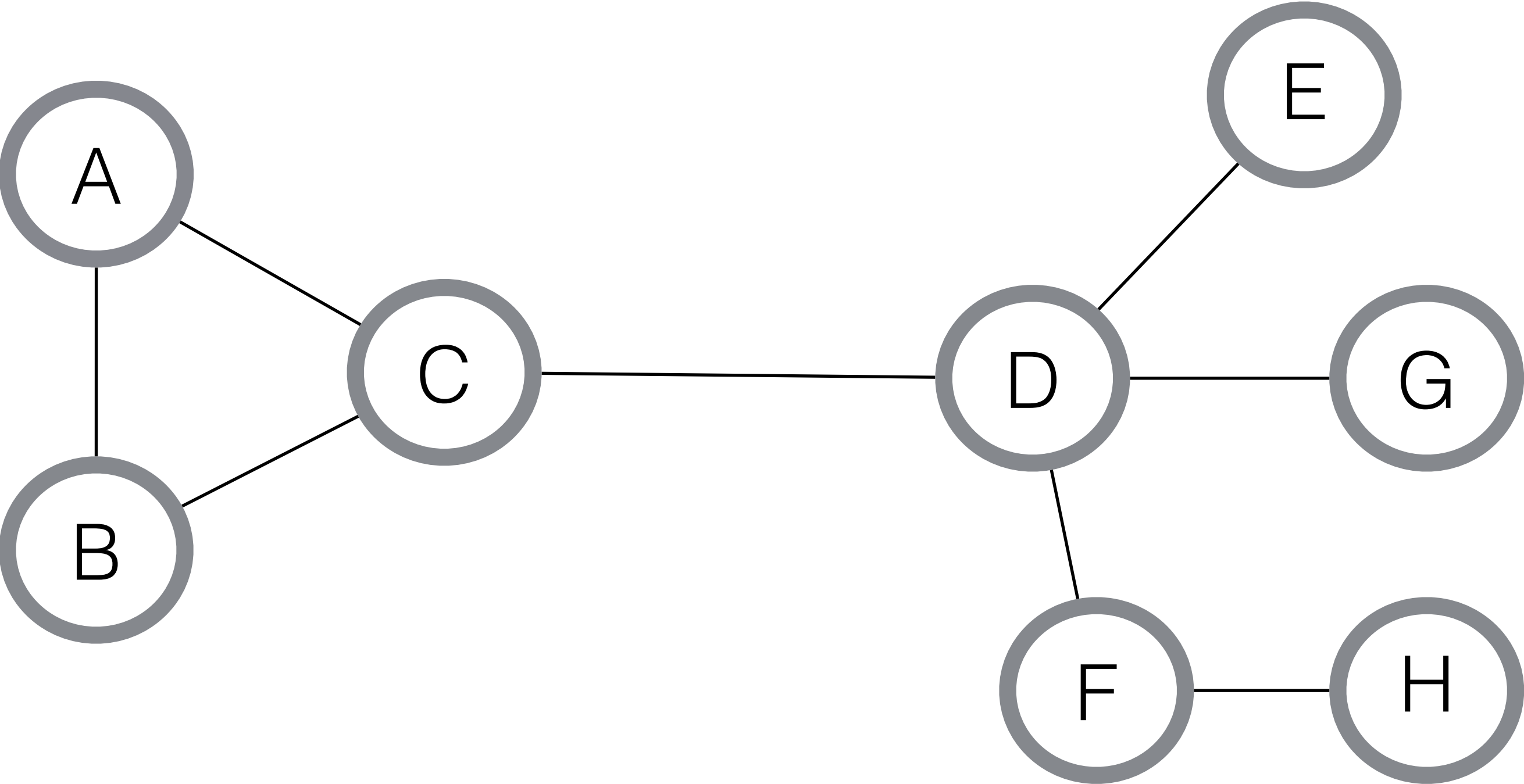
Kimon Fountoulakis



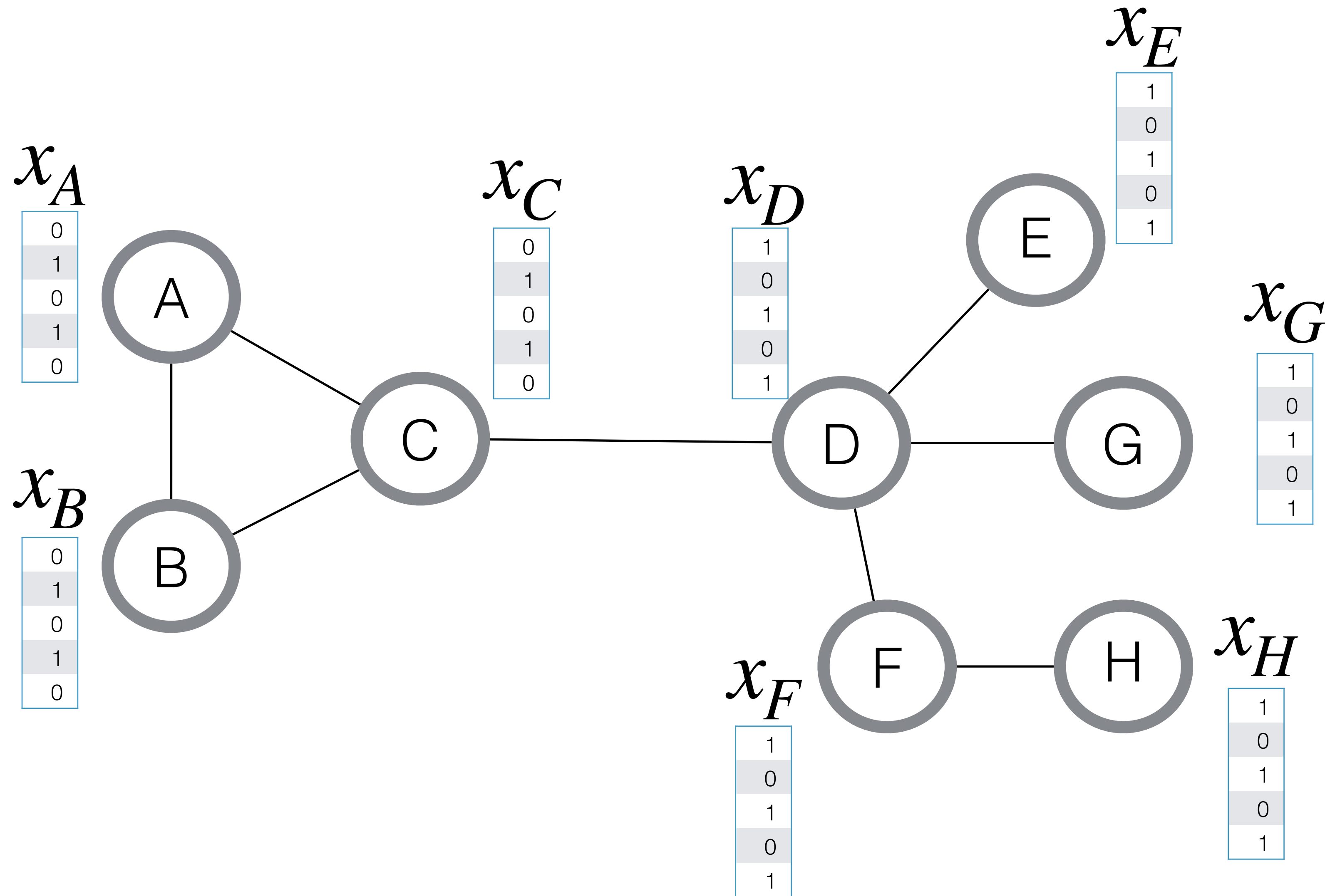
UNIVERSITY OF
WATERLOO

DAVID R. CHERITON SCHOOL
OF COMPUTER SCIENCE

Graphs

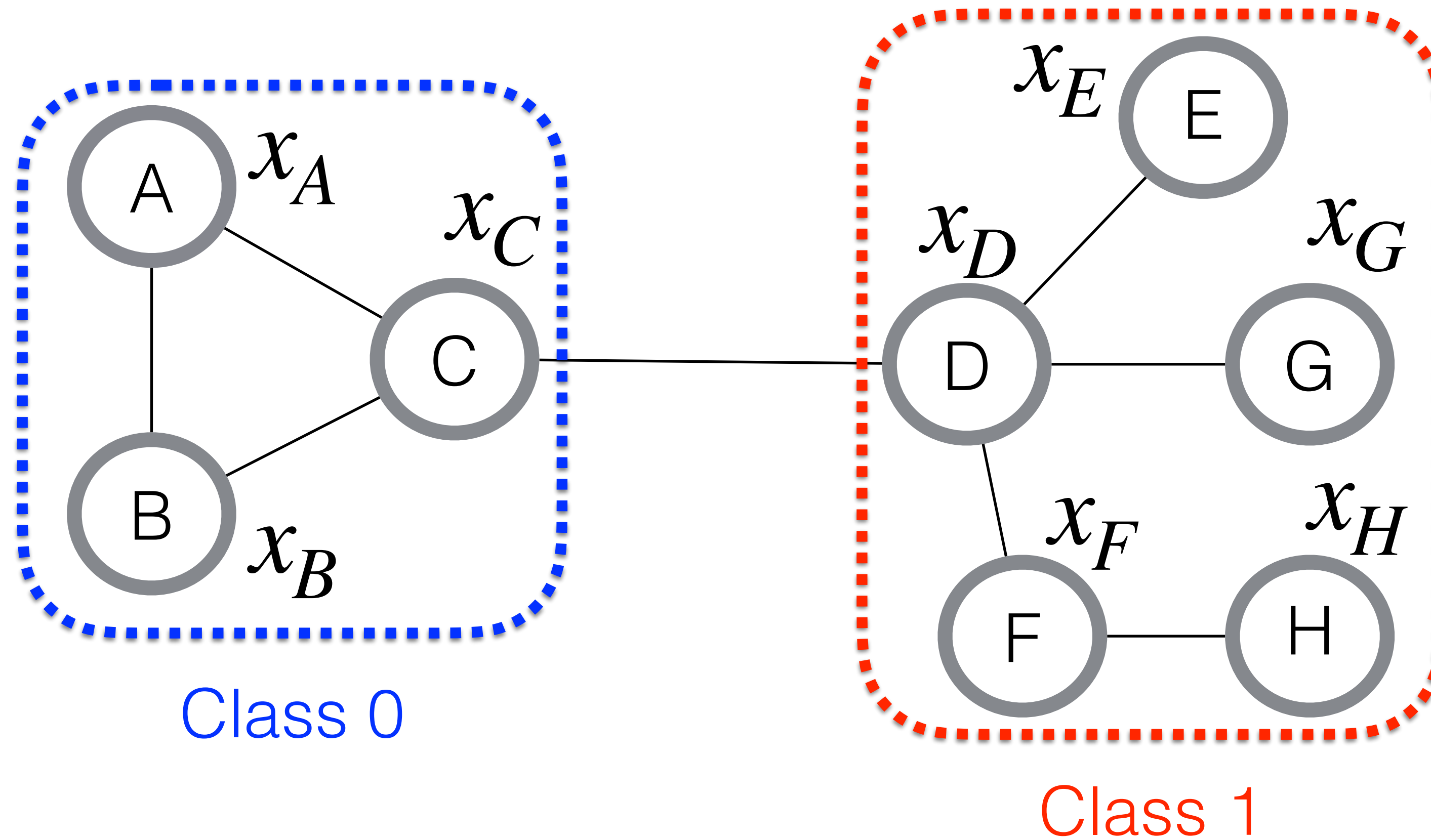


Graphs + features



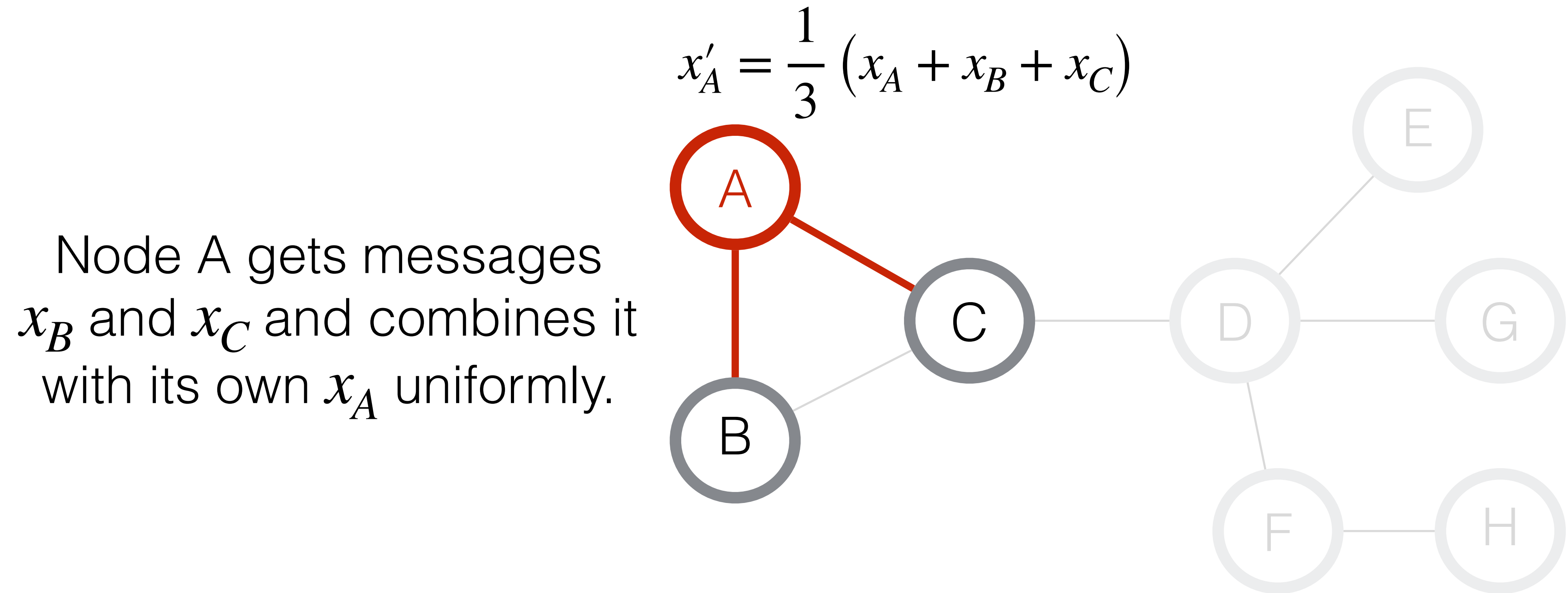
- x_i is the feature vector for node i

Node classification



- x_i is the feature vector for node i

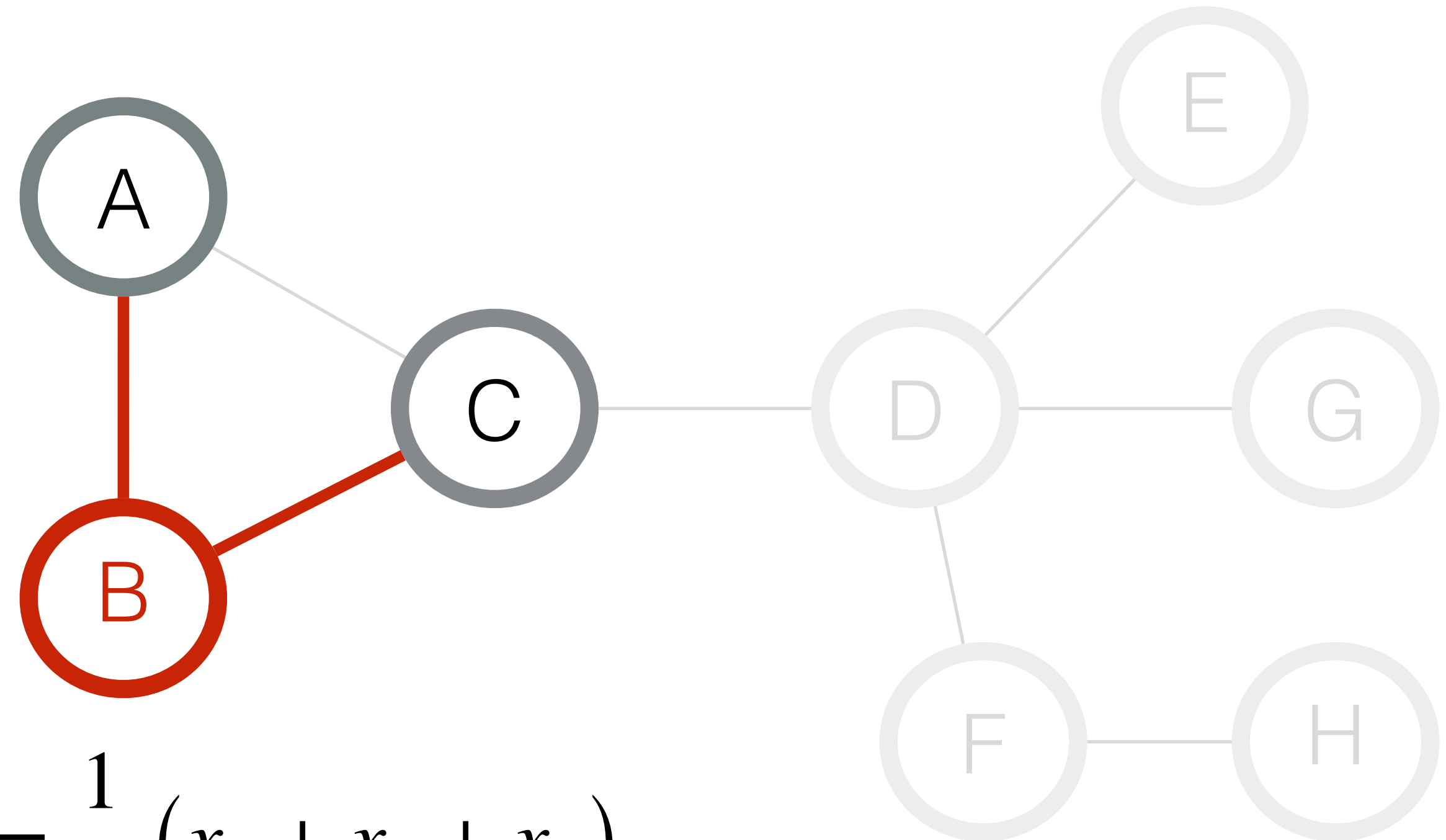
Vanilla Graph Convolution Network (GCN): aggregation function



Vanilla Graph Convolution Network (GCN): aggregation function

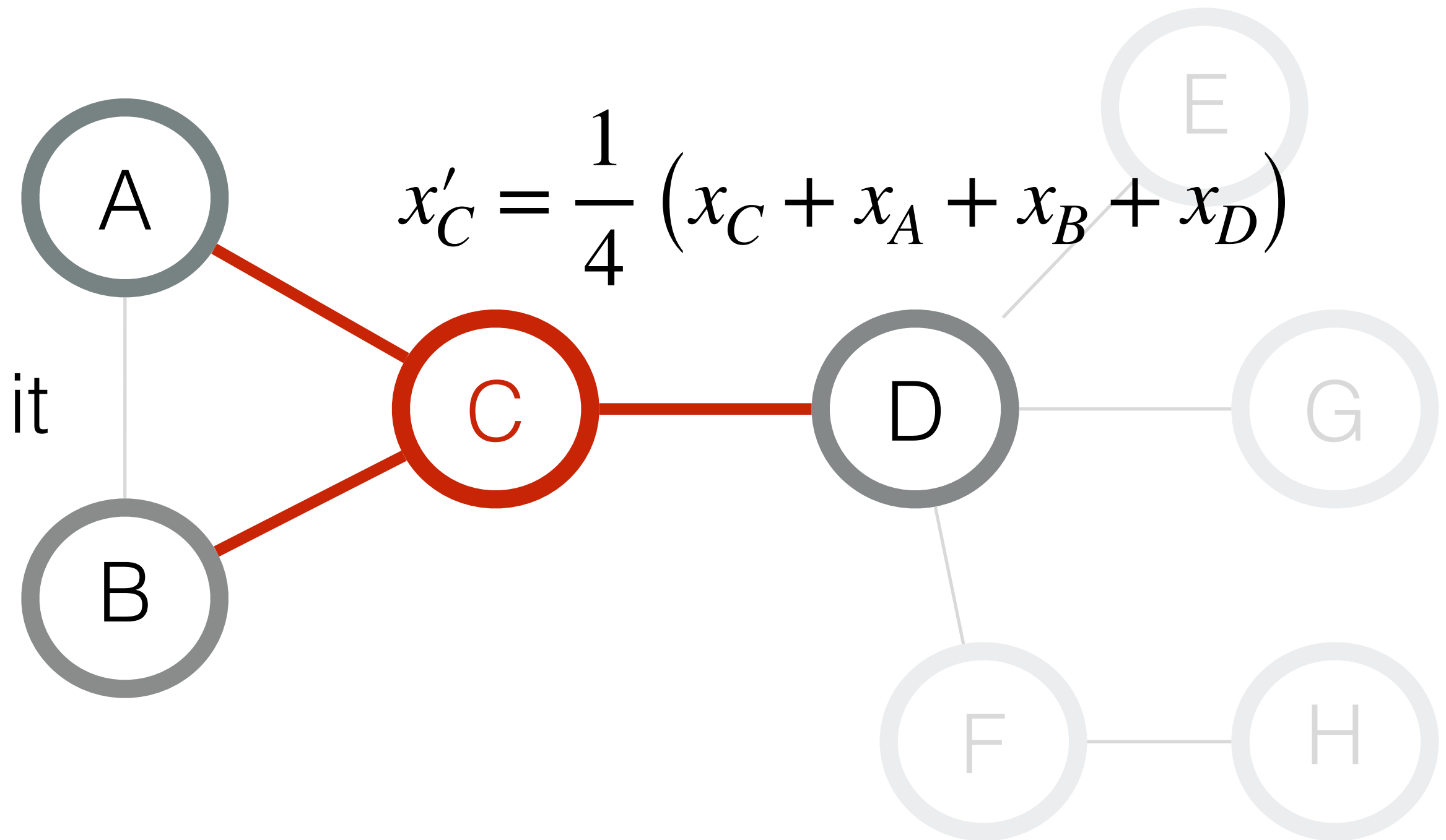
Node B gets messages x_A and x_C and combines it with its own x_B uniformly

$$x'_B = \frac{1}{3} (x_B + x_A + x_C)$$



Vanilla Graph Convolution Network (GCN): aggregation function

Node C gets messages x_A , x_B and x_D , and combines it with its own x_C uniformly



Vanilla Graph Convolution Network (GCN): aggregation function

$$\underset{\substack{\text{Convolved data} \\ \text{for node } i}}{X'_i} := \frac{1}{\underset{\substack{\text{Degree of} \\ \text{node } i}}{D_{ii}}} \sum_{j=1}^n \underset{\substack{\text{Adjacency} \\ \text{matrix}}}{A_{ij}} \underset{\substack{\text{Data} \\ \text{For node } j}}{X_j}$$

- A component of A is equal to 1 if two nodes are connected with an edge
- D is a diagonal matrix where each component shows the number of neighbors of a node

Vanilla Graph Convolution Network (GCN): aggregation function in matrix form

$$\mathbf{X}' \coloneqq \mathbf{D}^{-1} \mathbf{A} \mathbf{X}$$

Convolved data Degree matrix Adjacency matrix

- A component of \mathbf{A} is equal to 1 if two nodes are connected with an edge
- \mathbf{D} is a diagonal matrix where each component shows the number of neighbors of a node

Vanilla Graph Convolution Network (GCN): learning parameters

$$X'W := D^{-1}AXW$$

-Learning matrix W . It's value are decided by minimizing a loss function.

Vanilla Graph Convolution Network (GCN): activation

$$\sigma(X'W) := \sigma(D^{-1}AXW)$$

-Activation function σ . Examples include $\sigma(y) := \max(y, 0)$ or $\sigma(y) := \text{sigmoid}(y) = 1/(1 + e^{-y})$ which squeezes values in $[0, 1]$.

Vanilla Graph Convolution Network (GCN): multiple layers

Example: 3-layer GCN

$$X' := \sigma_3(D^{-1}A \underbrace{\sigma_2(D^{-1}A \underbrace{\sigma_1(D^{-1}AXW_1) W_2)}_{\text{layer 1}})}_{\text{layer 2}} W_3)$$

layer 3

Let's keep things simple

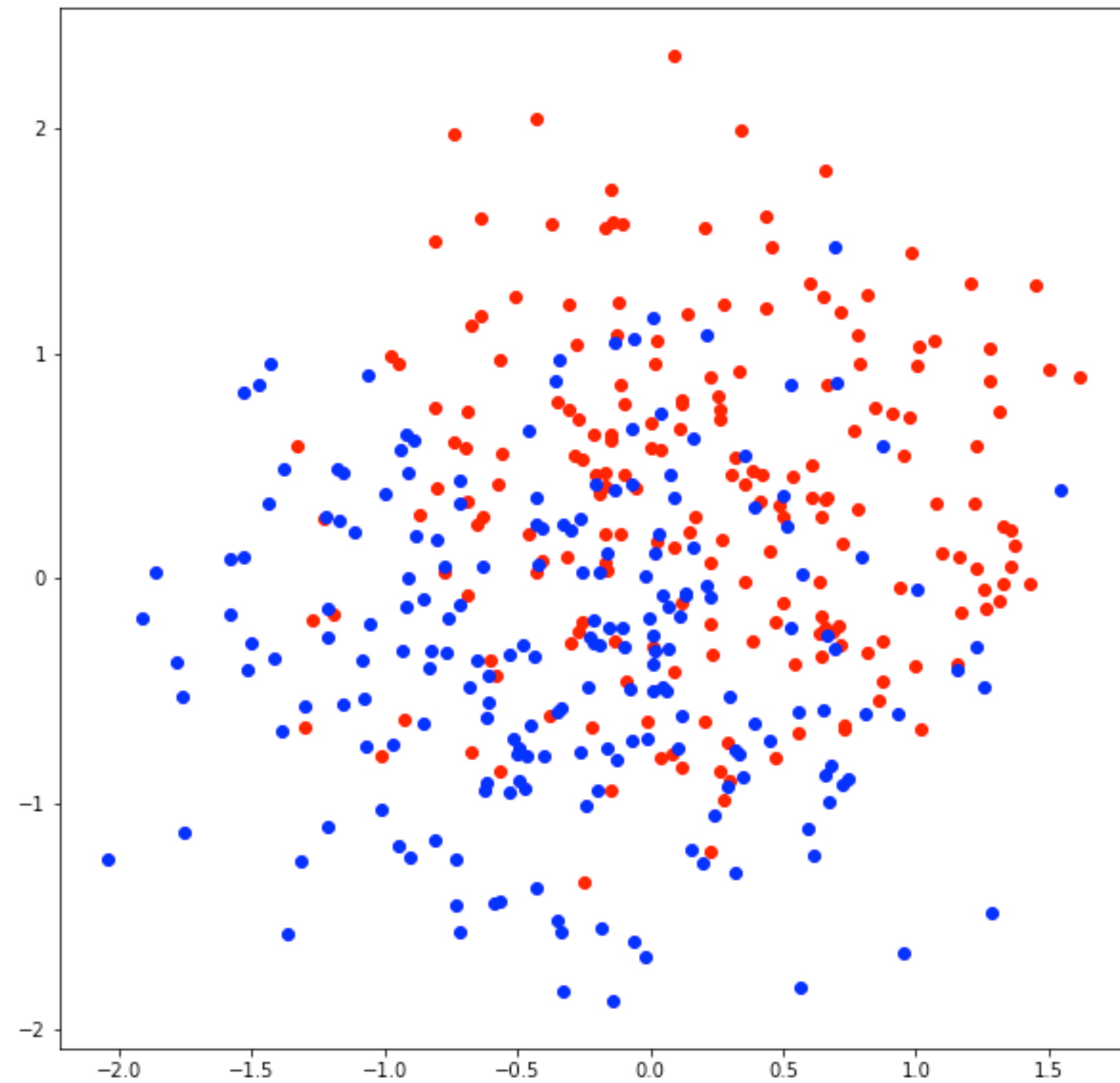
$$x' = D^{-1}AXw$$

- w is a vector of length equal to the number of features
- x' is a vector of length equal to the number of nodes, indicates predicted class membership

We want to answer the following for linear classifiers

- Does graph convolution of the data help generalization?
- Which graphs are good graphs and which are not?
- What if we test on a graph that comes from a distribution with different parameters?

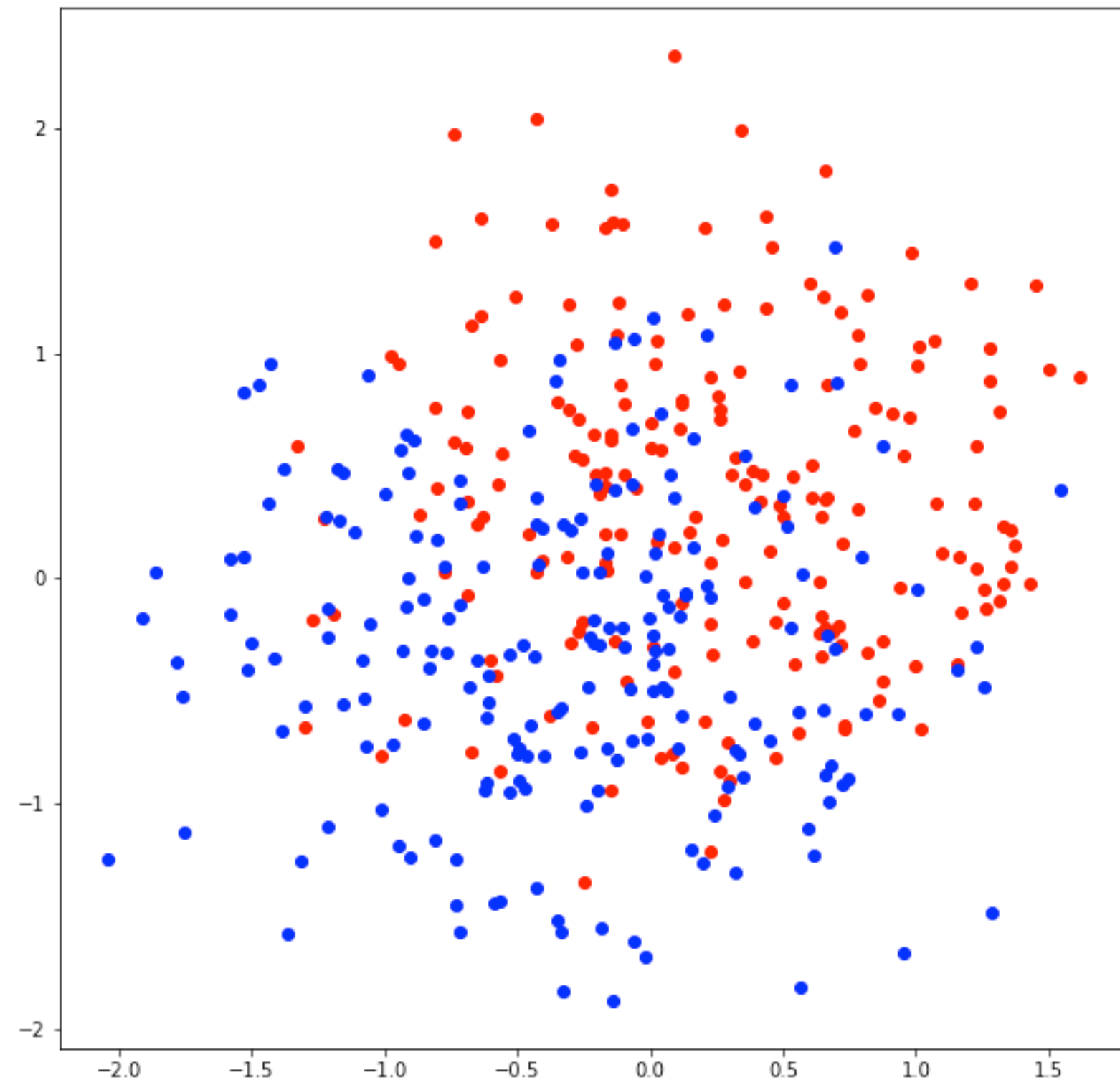
What can graph convolution do?



Original Data

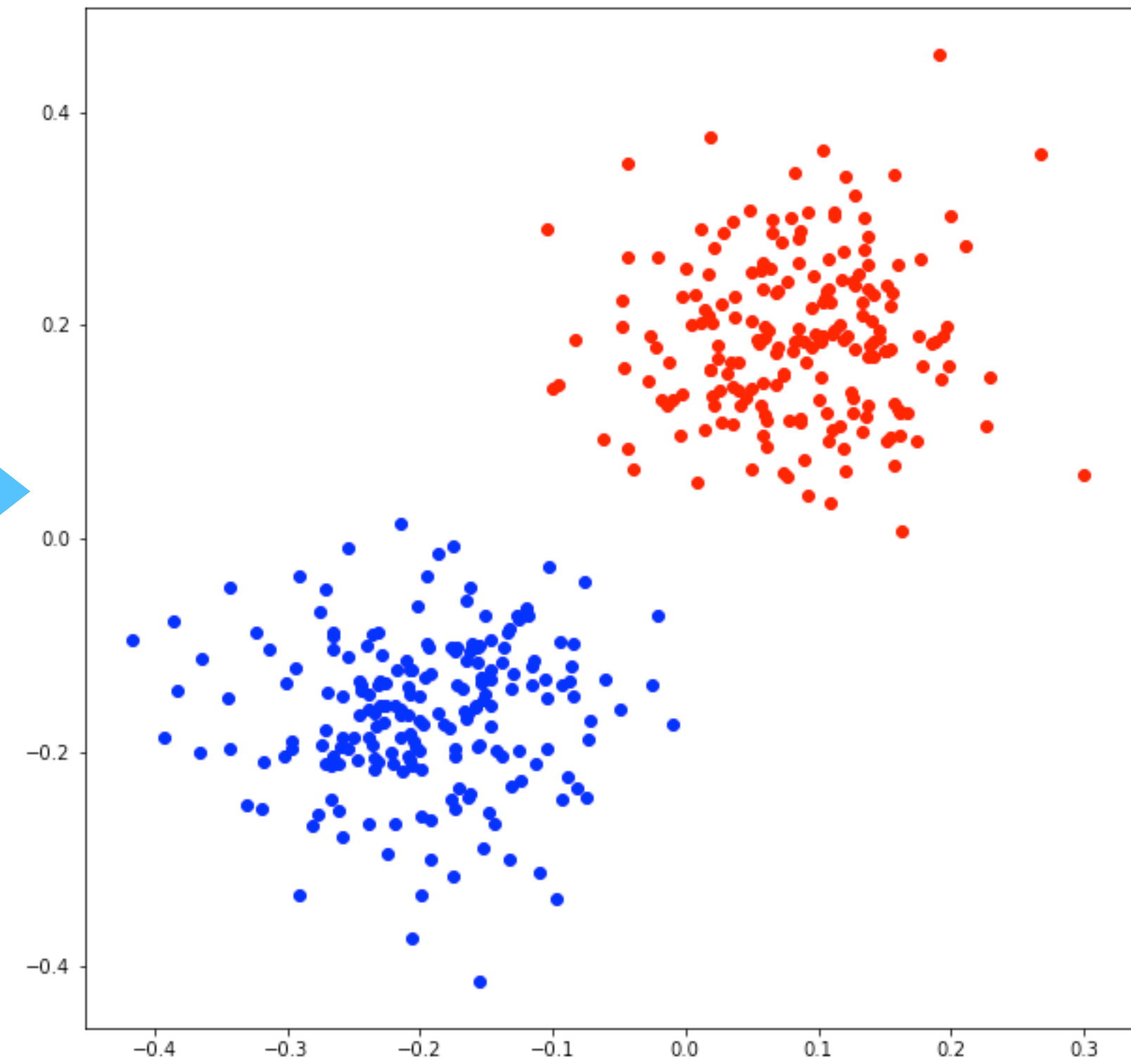
- Consider distributions with 2D features
- We cannot separate this data linearly
- Can graph convolution help?

What can graph convolution do?



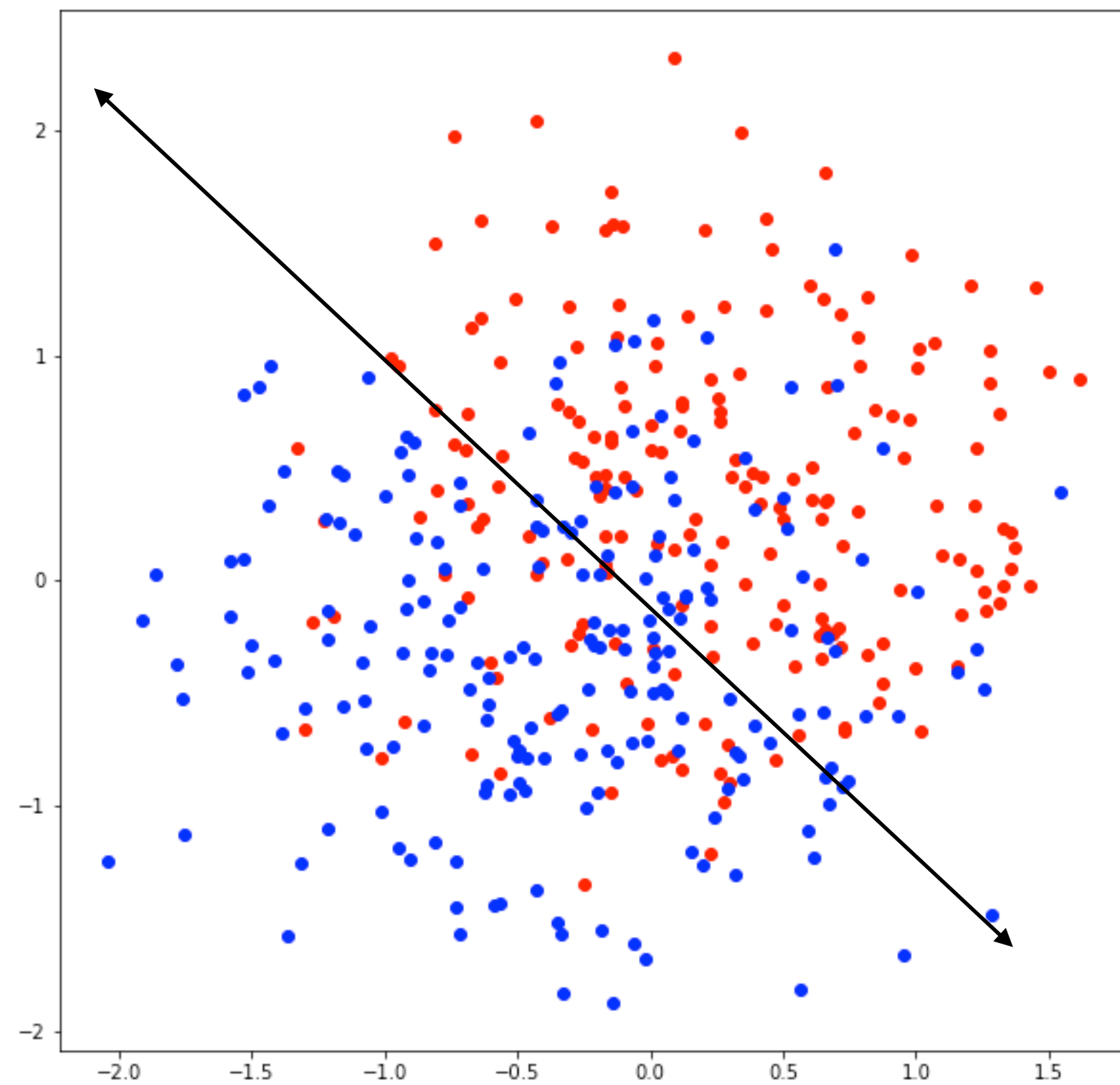
Original Data

Graph Convolution



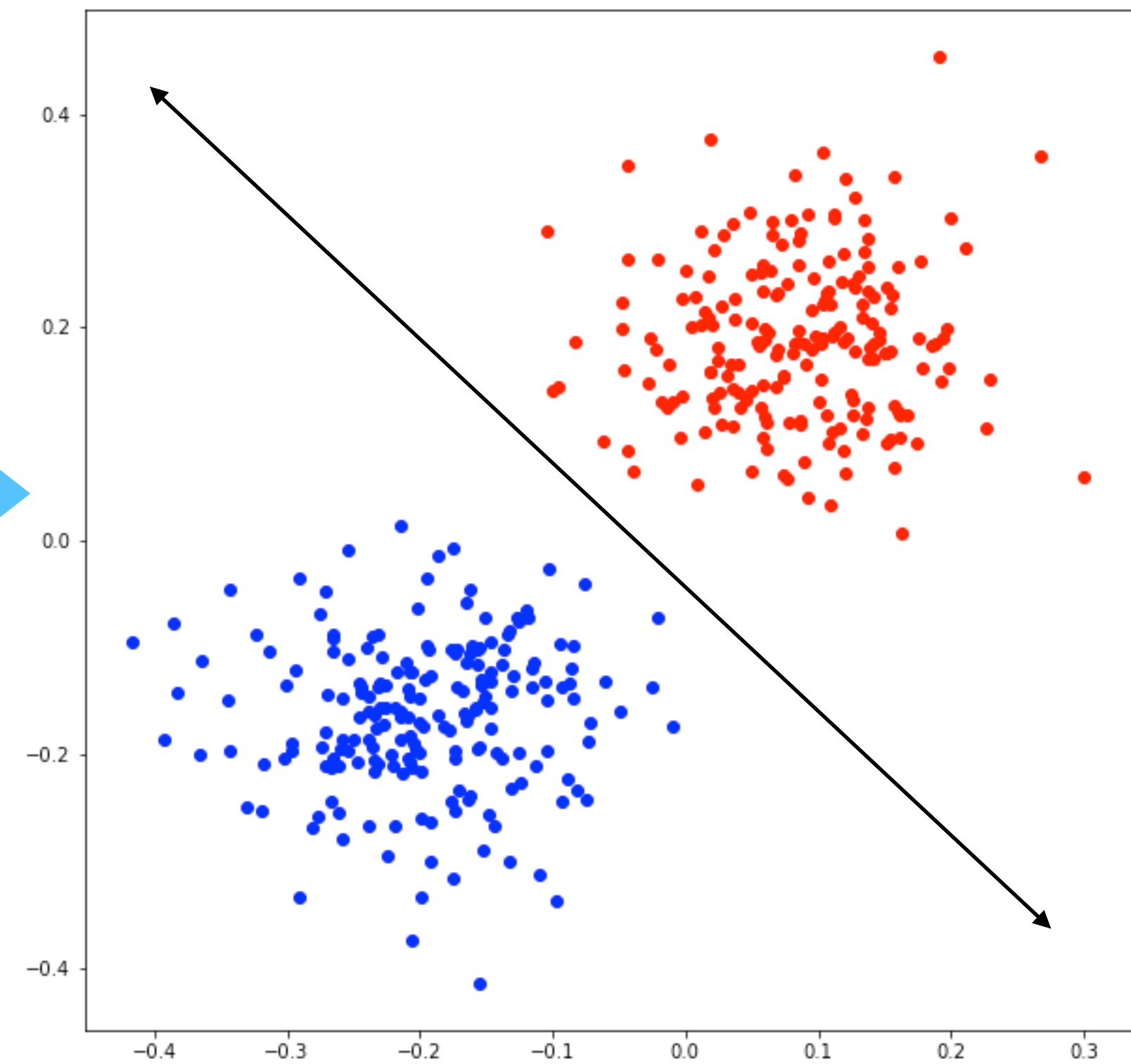
After graph convolution

What can graph convolution do?



Original Data

Graph Convolution



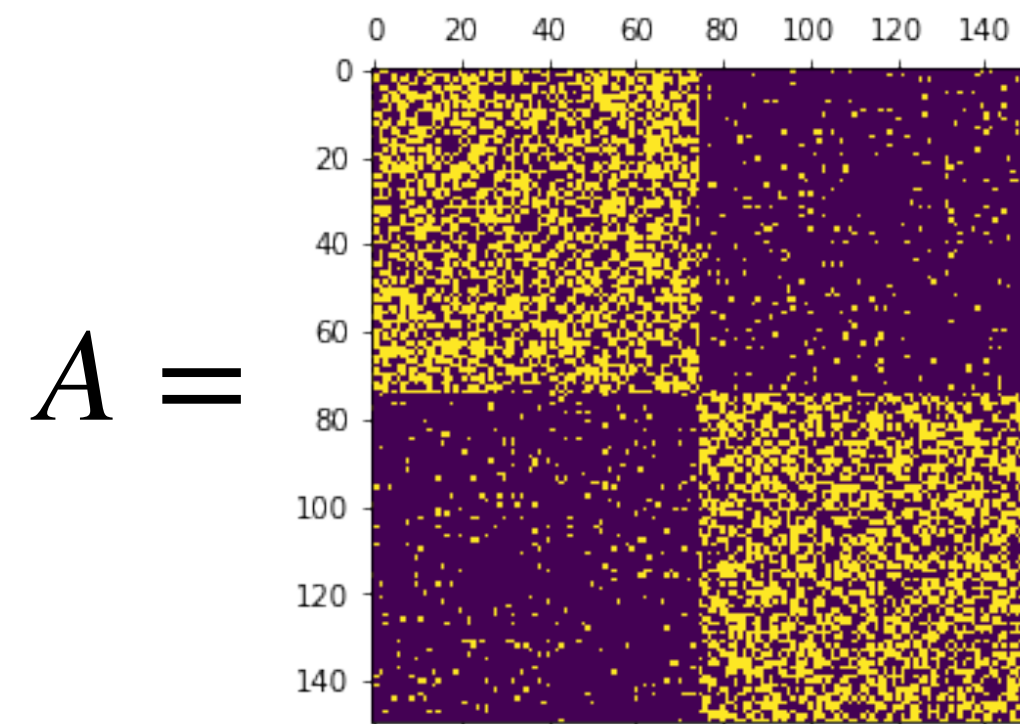
After graph convolution

Graph convolution makes the data linearly separable

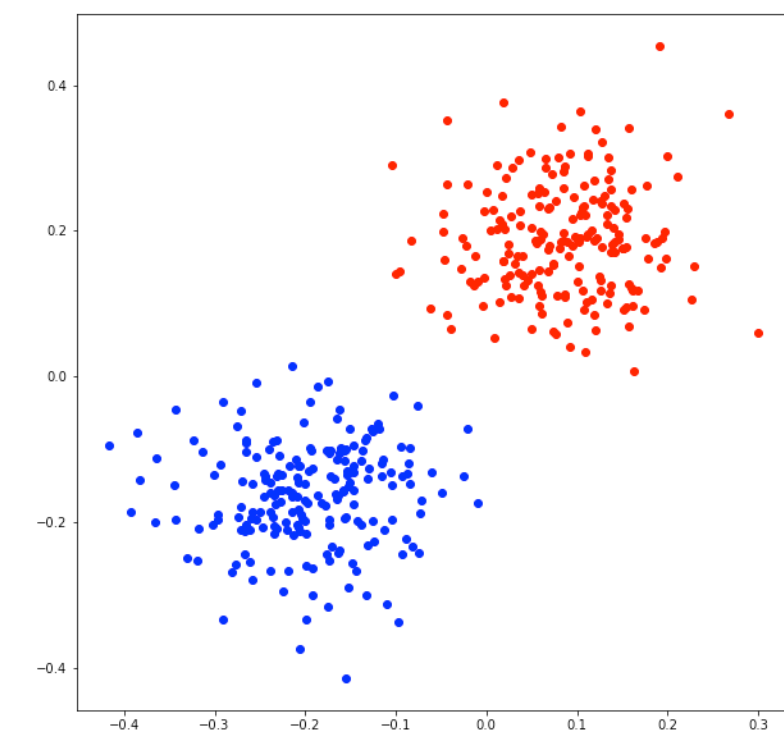
Data model: contextual stochastic block model

- Two-component balanced [Gaussian Mixture Model](#) (GMM) coupled with a [Stochastic Block Model](#) (SBM)

$$A \sim SBM(p, q)$$
$$\mathbb{P}(A_{ij} = 1) = \begin{cases} p & \text{if } i, j \text{ are in the same class} \\ q & \text{otherwise} \end{cases}$$



$$X_i \sim \mathcal{N}(\mu, \sigma^2 I) \text{ if } i \in \mathbf{C}_0$$
$$X_i \sim \mathcal{N}(\nu, \sigma^2 I) \text{ if } i \in \mathbf{C}_1$$



Data alignment

- The class membership of the nodes is the same to the communities in the graph.

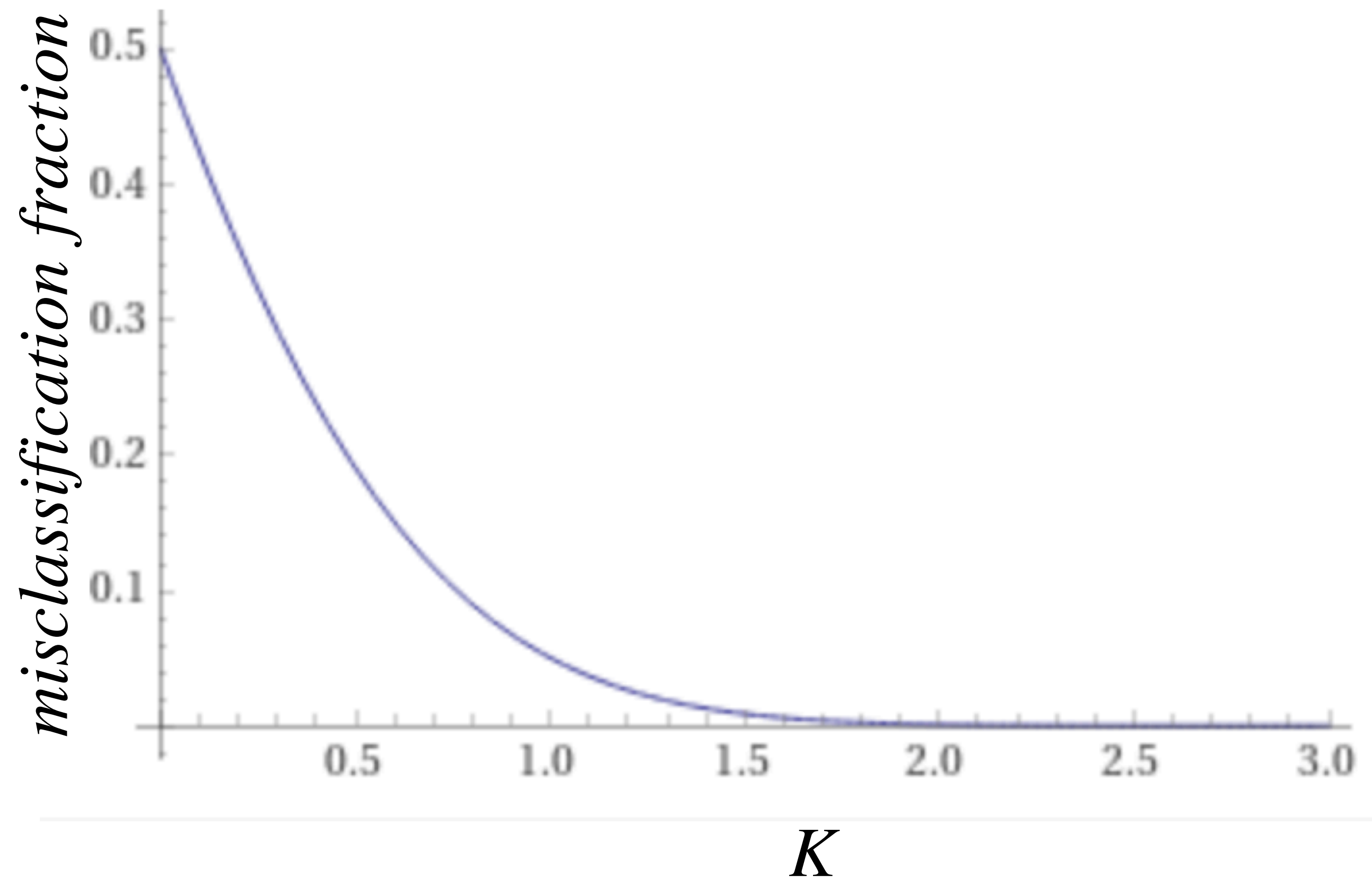
Let's ignore the graph first and make some assumptions

- Let's assume that we are given the distributions of the data
- This is very powerful knowledge because it assumes that we know the means μ, ν of the Gaussians and p, q .
- This allows to make prediction using the optimal Bayes classifier:

$$i^* = \operatorname{argmax}_{i \in \{0,1\}} P(y = i \mid x)$$

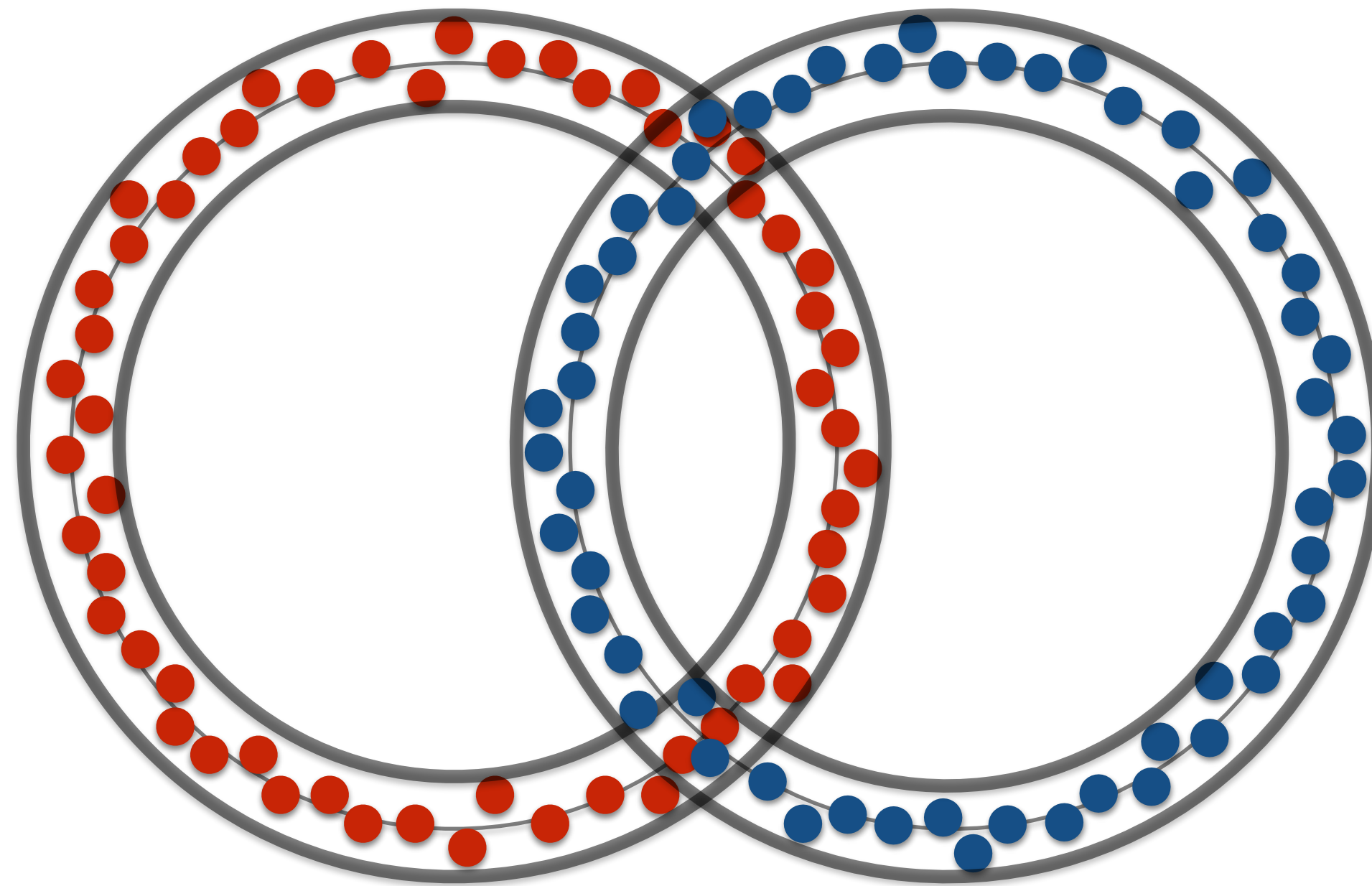
What's the performance of the optimal Bayes classifier?

- If the distance between the means $\|\mu - \nu\| \leq K\sigma$.
- Then the misclassification fraction by the Bayes classifier is:



Proof sketch

Intuitive representation of Gaussian data in high dimensions



$$\|\mu - \nu\|_2 = \mathcal{O}(\sigma)$$

What's the performance of the optimal Bayes classifier?

-If the distance between the means is $\|\mu - \nu\| = \Omega(\sigma\sqrt{\log n})$

-Then we have perfect classification with high probability.

Graph convolution improves linear separability

- Without the graph, **no hyperplane** can separate a binary GMM if means are $\mathcal{O}(\sigma)$ apart, i.e., $\|\mu - \nu\|_2 = \mathcal{O}(\sigma)$
- With graph convolution, this threshold changes to

$$\|\mu - \nu\| = \mathcal{O}\left(\frac{\sigma}{\sqrt{\mathbb{E}[D]}}\right)$$

Expected degree of a node

Assumptions for the Graph

- If the probability is $p = \Omega\left(\frac{\log^2 n}{n}\right)$
- and the probability $q = \Omega\left(\frac{\log^2 n}{n}\right)$
- This mean that the expected degree $E[D] = \Omega(\log^2 n)$

Assumptions for the Graph

- If the probability is $p = \Omega\left(\frac{\log^2 n}{n}\right)$
- and the probability $q = \Omega\left(\frac{\log^2 n}{n}\right)$
- And we also have degree concentration, i.e., with very high probability all nodes have degree around its expectation $E[D]$

Assumptions for the Graph

- Assume that $\Gamma(p, q) = \frac{|p - q|}{p + q} = \Omega(1)$

- There is a large enough gap between probabilities p and q .

- $p = q$, can't happen!

- Intuitively: There is enough “signal” in the graph relative to the noise.

Proof sketch for graph convolution

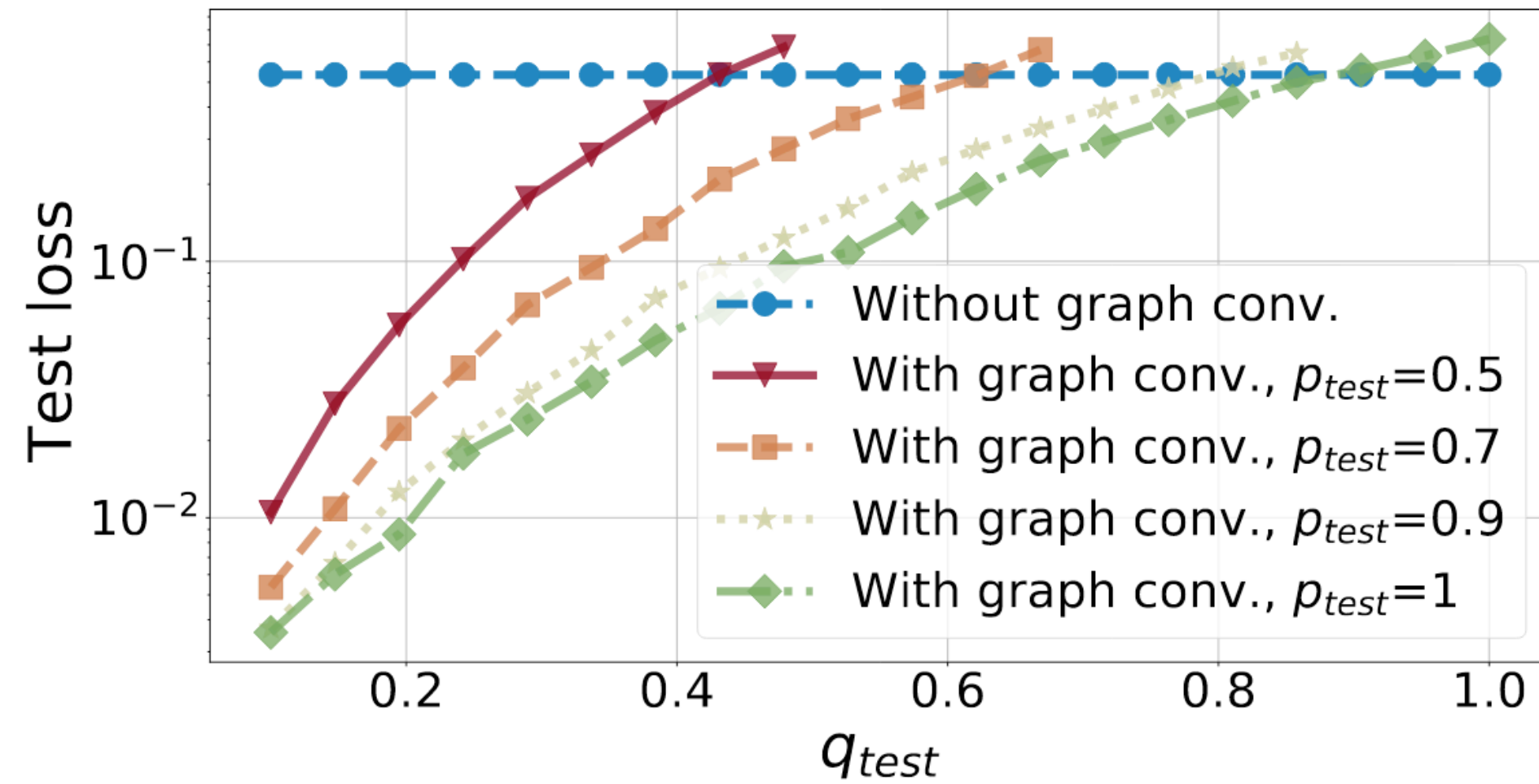
- After graph convolution the means move closer by a factor $\Gamma(p, q) = \frac{|p - q|}{p + q}$
- But the variance is reduced by $\mathbb{E}[D] = \Omega(\log^2 n)$
- Thus the separability threshold changes from $\|\mu - \nu\|_2 = \mathcal{O}(\sigma)$ to
$$\|\mu - \nu\| = \mathcal{O}\left(\frac{\sigma}{\sqrt{\mathbb{E}[D]}}\right)$$
- Then we can show that the hyperplane that passes through the mid-point of the two means separates the data with high probability.

Bounds on training loss

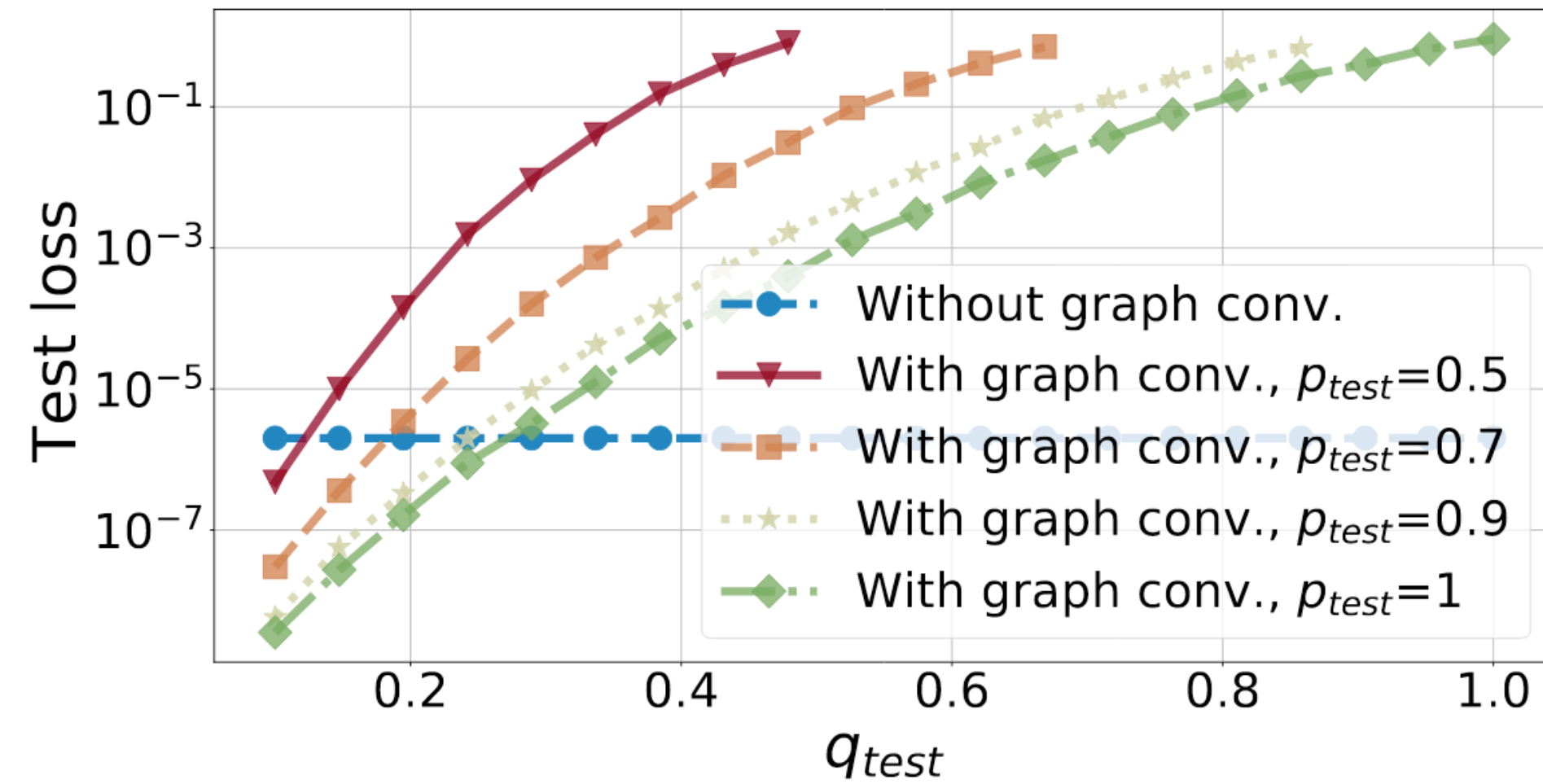
- We use binary cross entropy loss to learn the classifier
- Without graph convolution, if $\|\mu - \nu\|_2 = K\sigma$,
then the loss is lower bounded by $(2 \log 2)\Phi(-K/2)$
- But for the convolved data the loss decays exponentially
 $Loss(A, X) \leq C \exp(-d\|\mu - \nu\|\Gamma(p, q))$,
where d is the number of features.

Generalization

- If the assumptions about the graph are satisfied, then for any new dataset A, X with different n, p, q , the test loss is bounded above
$$Loss(A, X) \leq C \exp(-d \|\mu - \nu\| \Gamma(p, q))$$
- Loss increases with inter-class edge probability (noisy graph)



(a) $\|\mu_d - \nu_d\| = 2\sigma$



(b) $\|\mu_d - \nu_d\| = 16\sigma$