

Large scale graph learning from smooth signals

Kalofolias Vassilis
Nathanael Perraudin

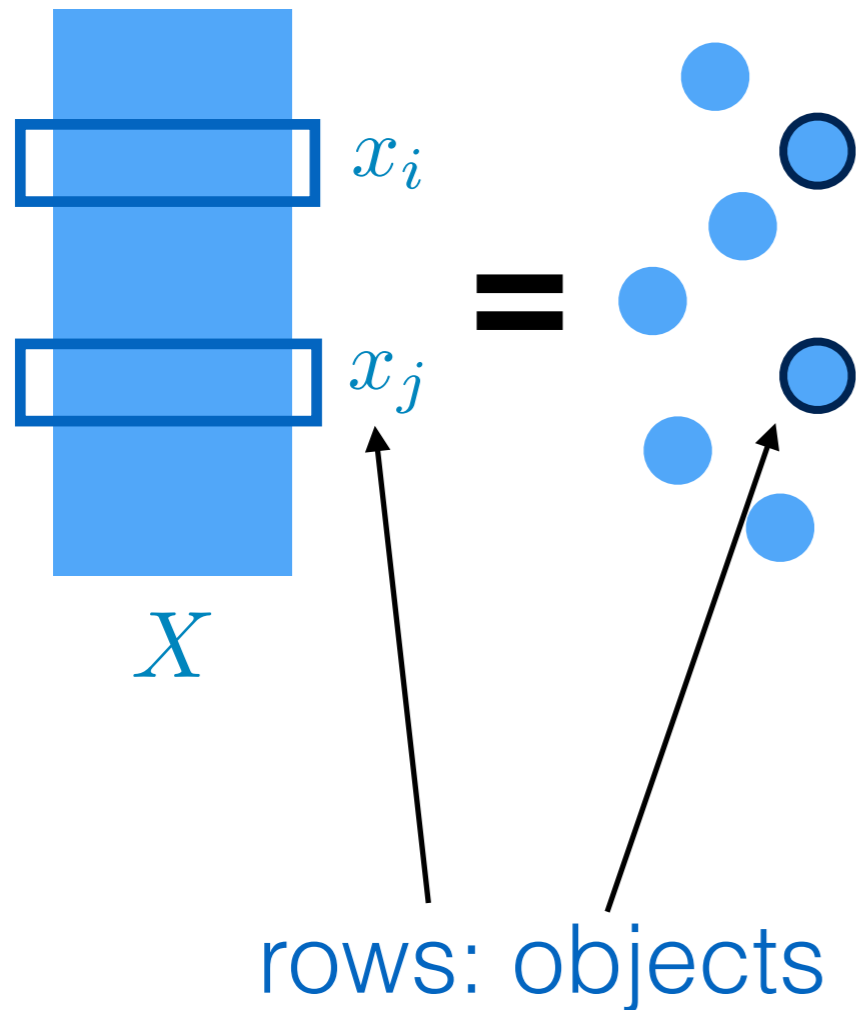


13 November 2019

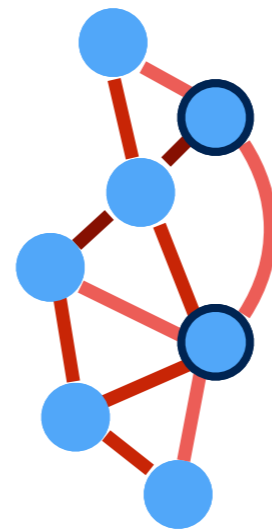


Graph learning

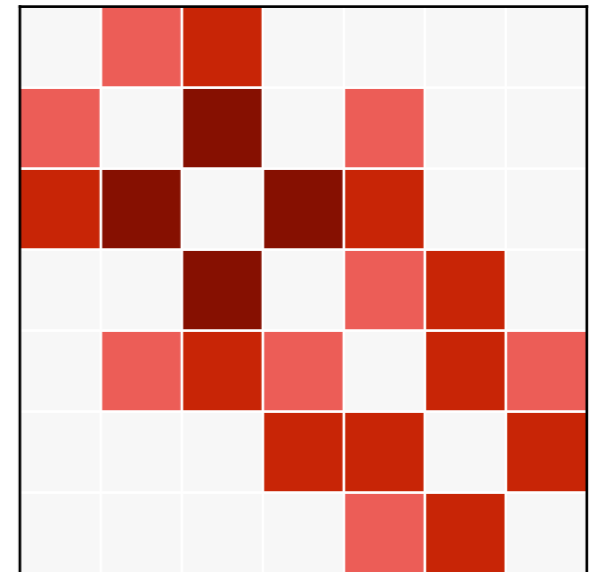
Given
matrix X



learn
graph G



W



weighted
adjacency
matrix

Dimensionality - manifolds

Interesting problems:

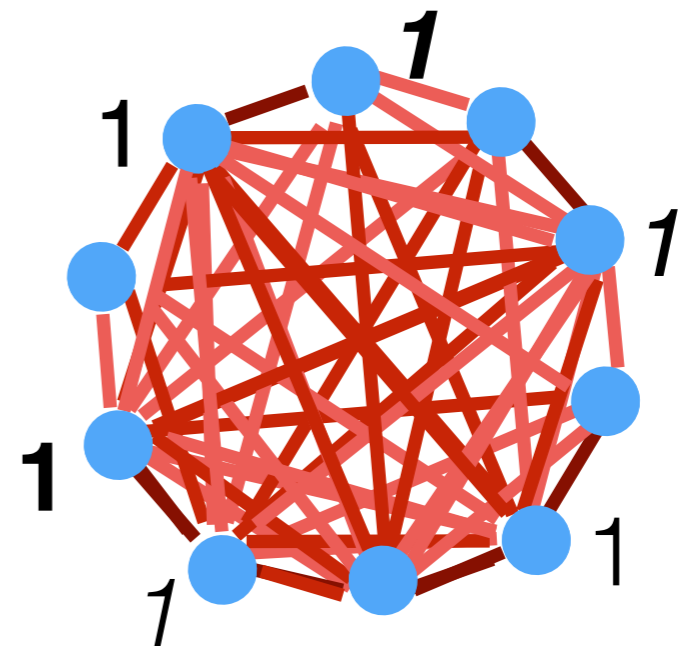
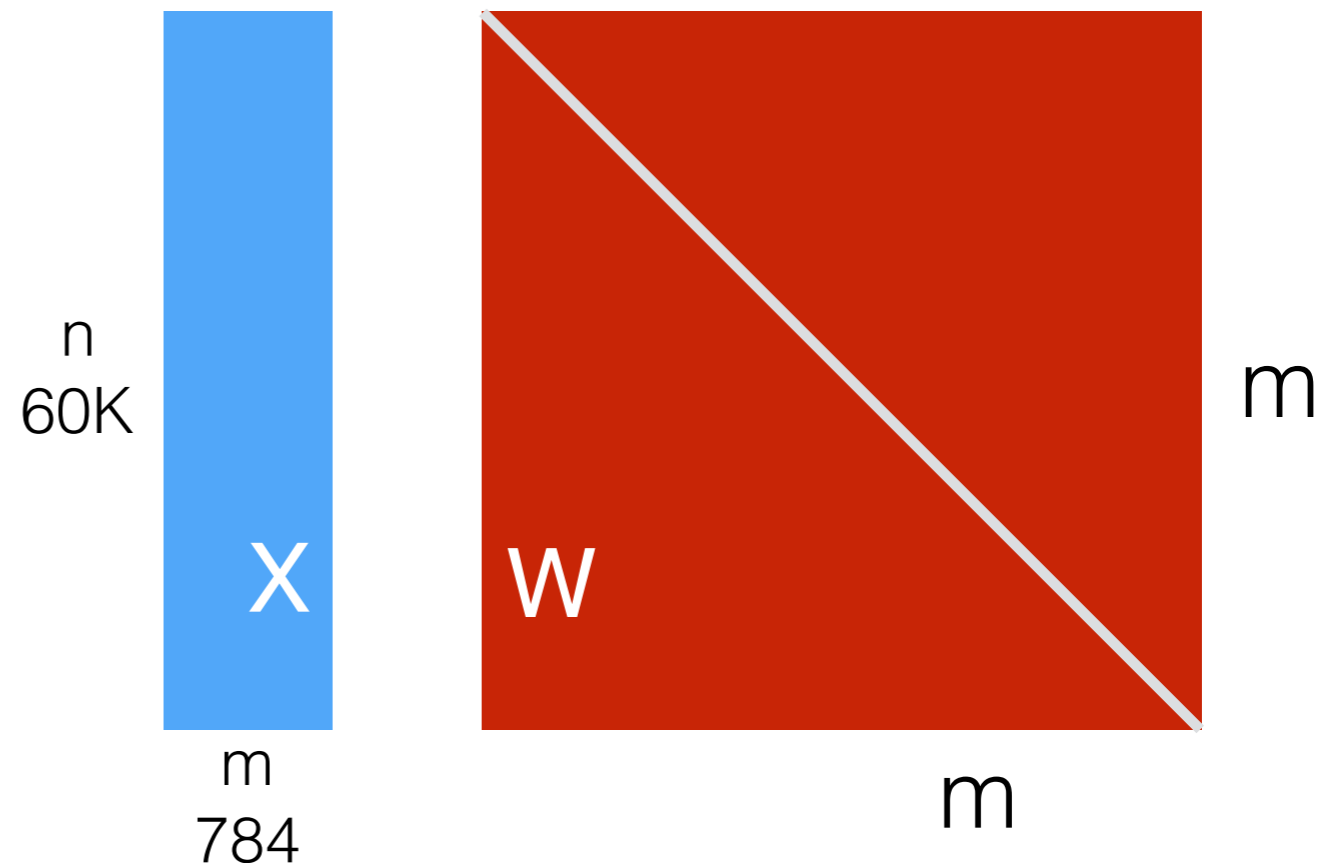
(# nodes) \gg (# features)

Example: MNIST

No structure?

\Rightarrow Full W

\times Ill-posed



Dimensionality - manifolds

Interesting problems:

(# nodes) \gg (# features)

Example: MNIST

No structure?

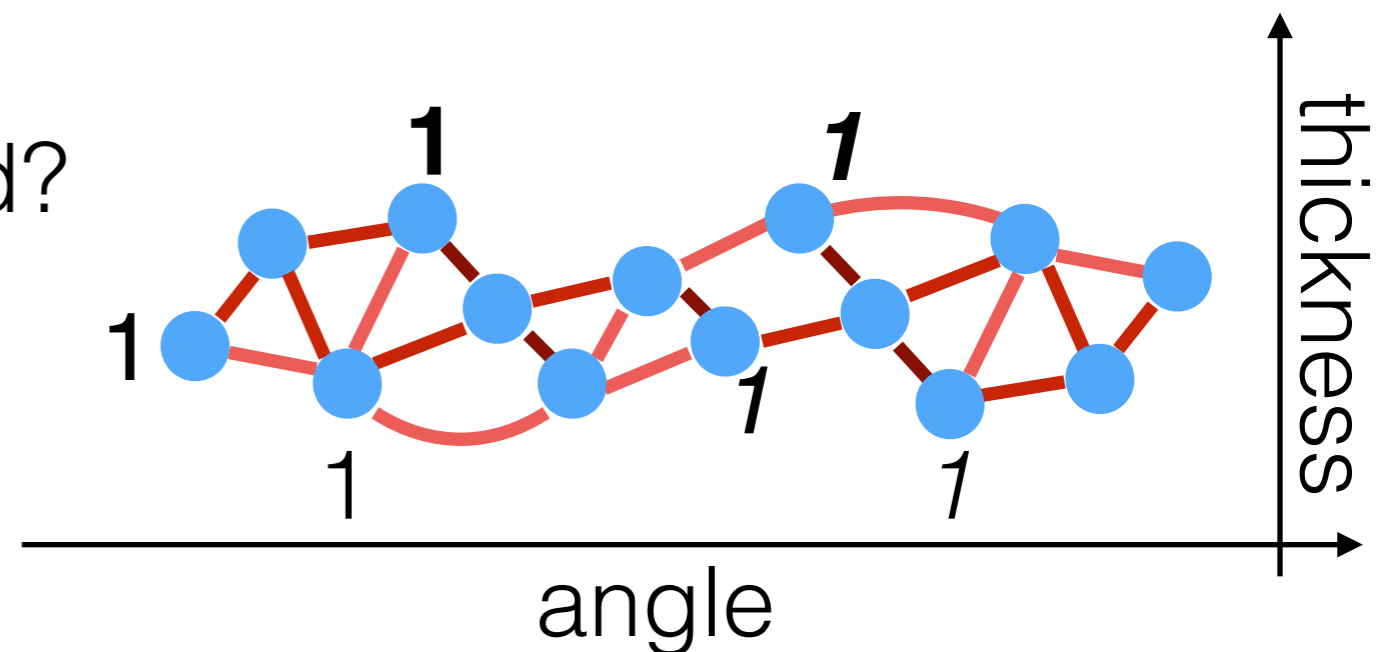
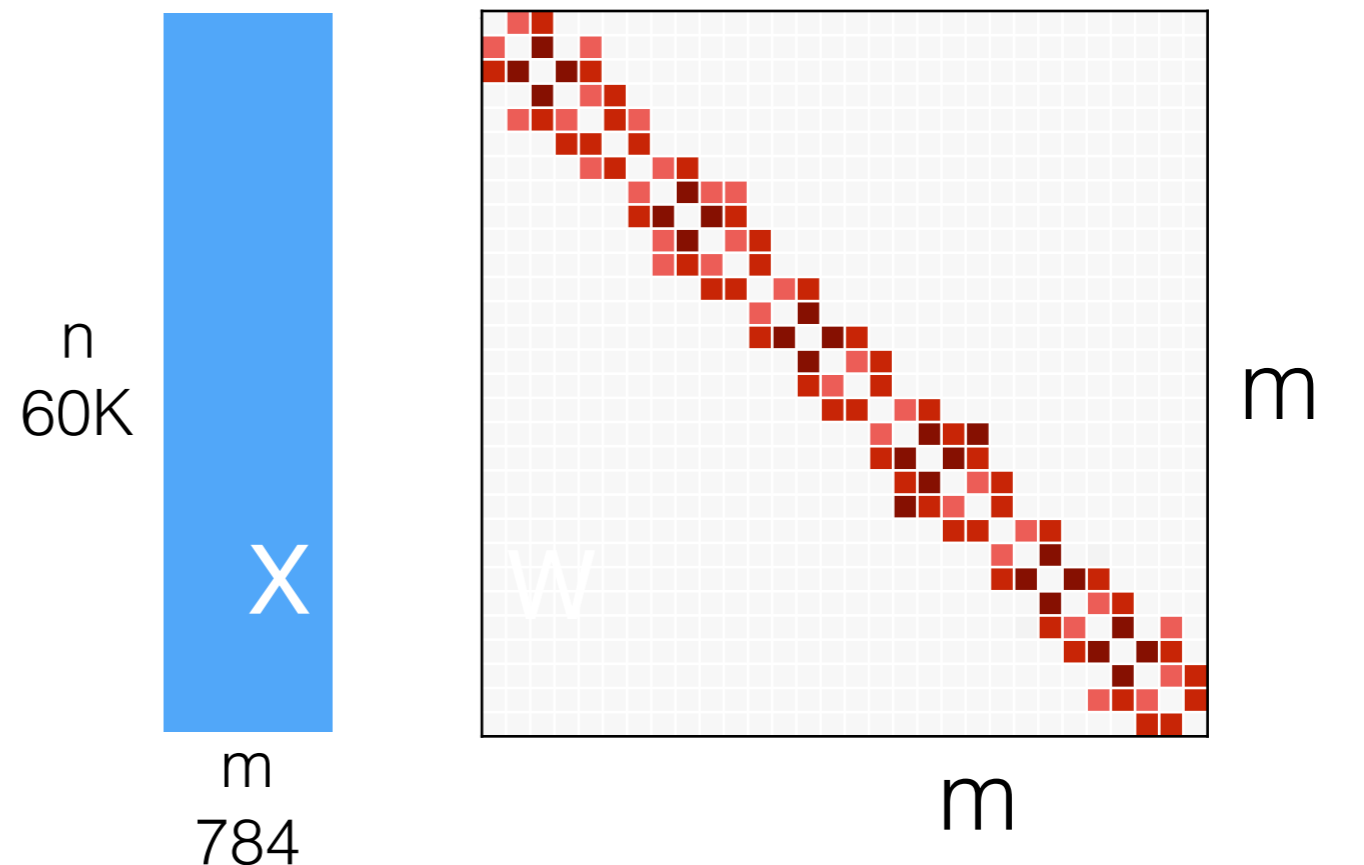
\Rightarrow Full W

\times Ill-posed

Low-dimensional manifold?

\Rightarrow Local dependencies

\checkmark Sparse W



Smoothness

Data is **smooth** on **graph**

Data lives on a low-dimensional manifold

Dirichlet energy is small:

$$\begin{aligned} \|\nabla_G X\|_F^2 &= \text{tr} \left(X^\top \underset{\substack{\uparrow \\ \nabla_G^\top \nabla_G \\ D - W}}{L} X \right) = \frac{1}{2} \sum_{i,j} W_{i,j} \|x_i - x_j\|_2^2 \\ &= \frac{1}{2} \|W \circ Z\|_{1,1} \end{aligned}$$

$Z_{ij} = \|x_i - x_j\|_2^2$

data smoothness \Rightarrow graph sparsity

Graphs from smooth signals

$$\min \operatorname{tr} \|LX\|_F^2 \quad \text{s.t.} \quad W\mathbf{1} \geq \mathbf{1} \quad [\text{Daitch et al 2009}]$$

$$\min \operatorname{tr} \|LX\|_F^2 \quad \text{s.t.} \quad \mathbf{1}^\top \max(\mathbf{0}, W\mathbf{1})^2 \leq \alpha n \quad [\text{Daitch et al 2009}]$$

$$\min \operatorname{tr}(X^\top LX) - \log |L + \alpha I| + \beta \|W\|_{1,1} \quad [\text{Lake \& Tenenbaum 2010}]$$

$$\min \operatorname{tr}(X^\top LX) + \alpha \|L\|_F^2 \quad \text{s.t.} \quad \operatorname{tr}(L) = n \quad [\text{Hu et al 2015, Dong et al 2016}]$$

$$\min \operatorname{tr}(X^\top LX) - \alpha \mathbf{1}^\top \log(W\mathbf{1}) + \frac{\beta}{2} \|W\|_F^2 \quad [\text{Kalofolias 2016}]$$

The log-degrees model

$$\min_{W \in \mathcal{W}_m} \|W \circ Z\|_{1,1} - \alpha \mathbf{1}^\top \log(W \mathbf{1}) + \frac{\beta}{2} \|W\|_F^2$$

- First algorithm of $\mathcal{O}(n^2)$
- Best results among “scalable” models

Goal: scale it further!

How to scale it?

1. Reduce the number of variables from $\mathcal{O}(n^2)$
2. Eliminate grid search: automatic parameter selection

How to scale it?

- 1. Reduce the number of variables from $\mathcal{O}(n^2)$**
2. Eliminate grid search: automatic parameter selection


Optimization

Objective can be split in 3 functions:

$$\min_{W \in \mathcal{W}_m} \|W \circ Z\|_{1,1} - \alpha \mathbf{1}^\top \log(W \mathbf{1}) + \frac{\beta}{2} \|W\|_F^2$$

Sketch of algorithm:

Approximately minimize each function

- 
1. Shrink edges according to distance $\mathcal{O}(n^2)$
 2. Enhance edges of badly connected nodes $\mathcal{O}(n^2)$
 3. Shrink large edges $\mathcal{O}(n^2)$

Optimization

Objective can be split in 3 functions:

$$\min_{W \in \mathcal{W}_m} \|M \circ W \circ Z\|_{1,1} - \alpha \mathbf{1}^\top \log((M \circ W)\mathbf{1}) + \frac{\beta}{2} \|M \circ W\|_F^2$$



1. Shrink edges according to distance

$$\mathcal{O}(|\mathcal{E}^{\text{allowed}}|)$$

2. Enhance edges of badly connected nodes

$$\mathcal{O}(|\mathcal{E}^{\text{allowed}}|)$$

3. Shrink large edges

$$\mathcal{O}(|\mathcal{E}^{\text{allowed}}|)$$

Reducing allowed edge set

How do we choose a restricted edge set?

- Prior: structure imposed by application
e.g. geometric constraints

What if no structure known?

- Approximate Nearest Neighbours (ANN)

Using ANN to reduce cost

"I want a graph with 10 edges per node on average"

$$k = 10$$

Compute approximate 30 NN graph (binary)

$$\mathcal{O}(n \log(n)m)$$

$$|\mathcal{E}^{\text{allowed}}| \approx 3kn = 30n$$

Learn weights for allowed edges

$$\mathcal{O}(nk)$$

Some of them are deleted! ($W_{ij}=0$)

Final 10 NN graph

Cost?

How to scale it?

1. Reduce the number of variables from $\mathcal{O}(n^2)$
2. **Eliminate grid search: automatic parameter selection**

Change of parameters

“I want a graph with 20 edges per node on average”

Grid search? **✗**

$$W^*(Z, \alpha, \beta) = \arg \min_{W \in \mathcal{W}_m} \|W \circ Z\|_{1,1} - \alpha \mathbf{1}^\top \log(W \mathbf{1}) + \frac{\beta}{2} \|W\|_F^2$$

$$= \delta \arg \min_{W \in \mathcal{W}_m} \|W \circ \theta Z\|_{1,1} - \mathbf{1}^\top \log(W \mathbf{1}) + \frac{1}{2} \|W\|_F^2$$

$$\delta = \sqrt{\frac{\alpha}{\beta}} \quad \theta = \sqrt{\frac{1}{\alpha\beta}}$$

$$= \delta W^*(\theta Z, 1, 1)$$

$$\leq \delta$$

Only θ changes sparsity

δ only changes the scale

Sparsity of one node

So δ is not important. How do we find θ ?

$$\min_{W \in \mathcal{W}_m} \|W \circ \theta Z\|_{1,1} - \mathbf{1}^\top \log(W \mathbf{1}) + \frac{1}{2} \|W\|_F^2$$

Take $\mathbf{1}$ node:

$$\min_{w \geq 0} \theta w^\top z - \log(w^\top \mathbf{1}) + \frac{1}{2} \|w\|_2^2.$$

ignore
symmetricity

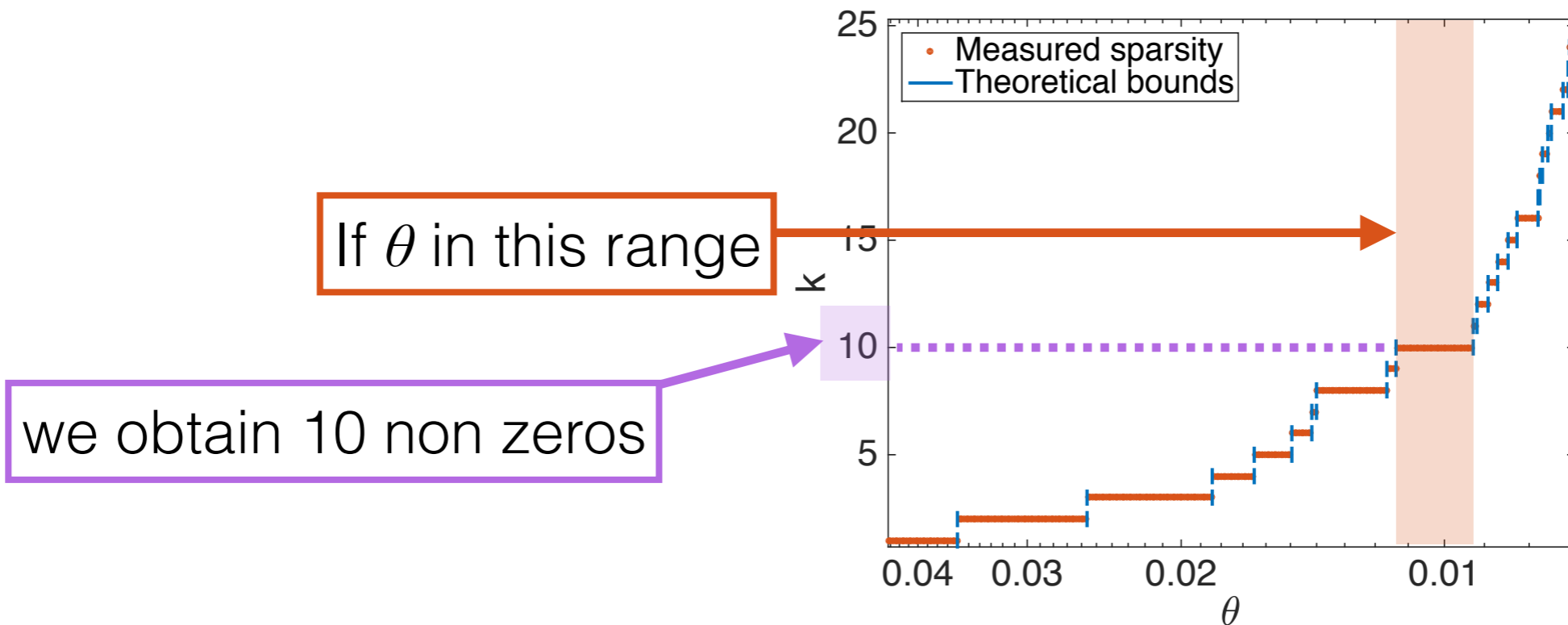
1 column
of W

Analyse role of θ on simpler problem!

Sparsity of one node

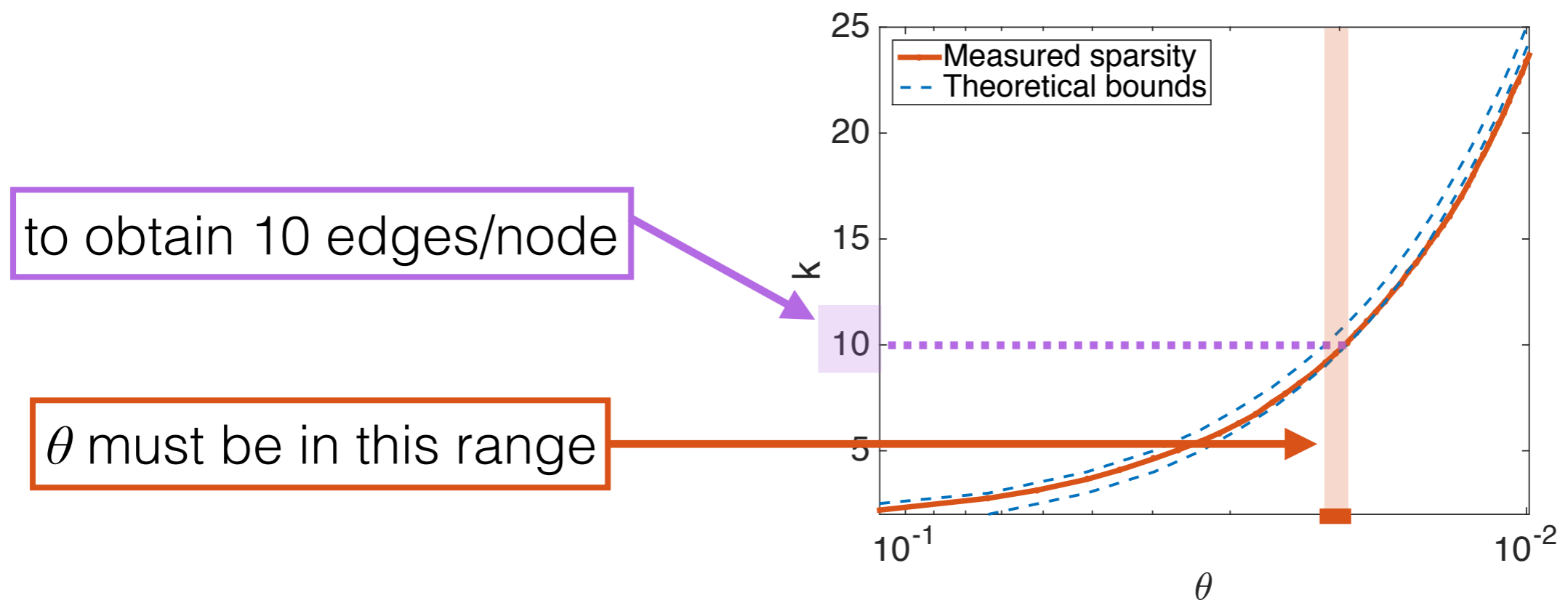
Theorem:

By setting θ in the range $\left(\frac{1}{\sqrt{kz_{k+1}^2 - b_k z_{k+1}}}, \frac{1}{\sqrt{kz_k^2 - b_k z_k}} \right]$, w^* has exactly k non-zero elements.

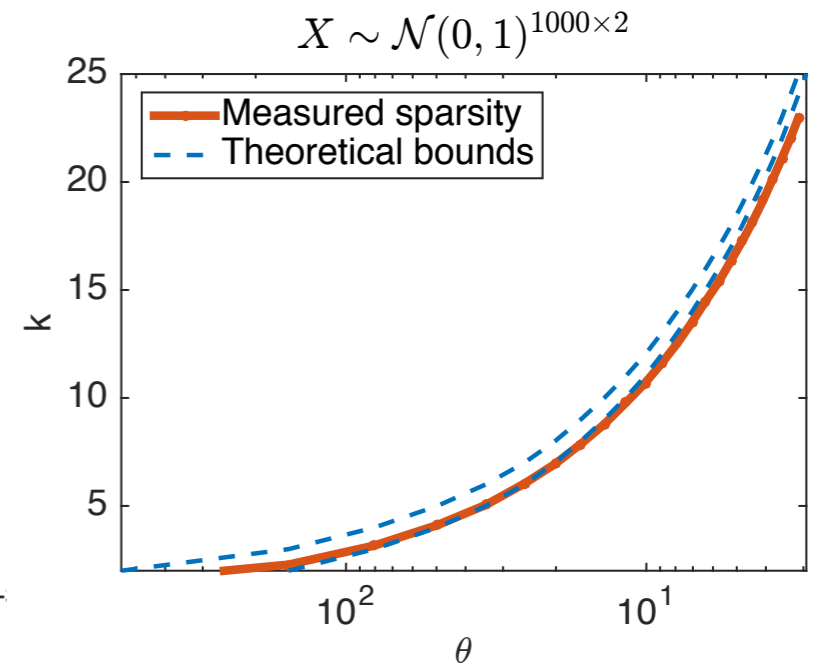
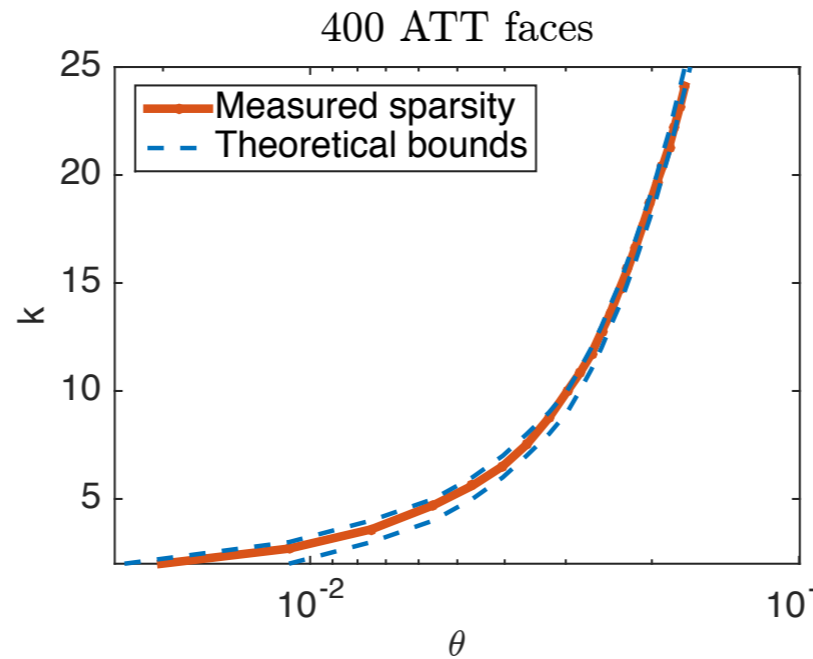
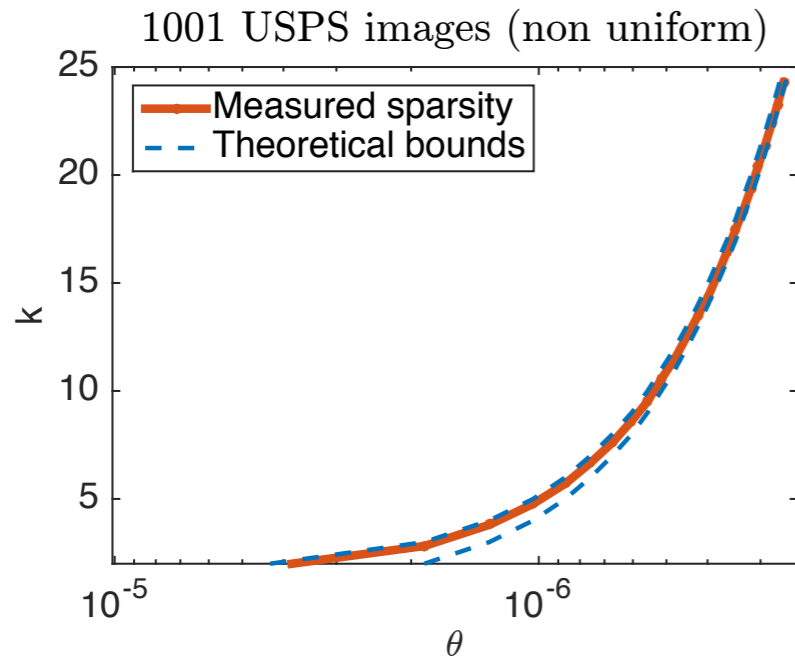


Sparsity of entire graph

Use **average** over all nodes

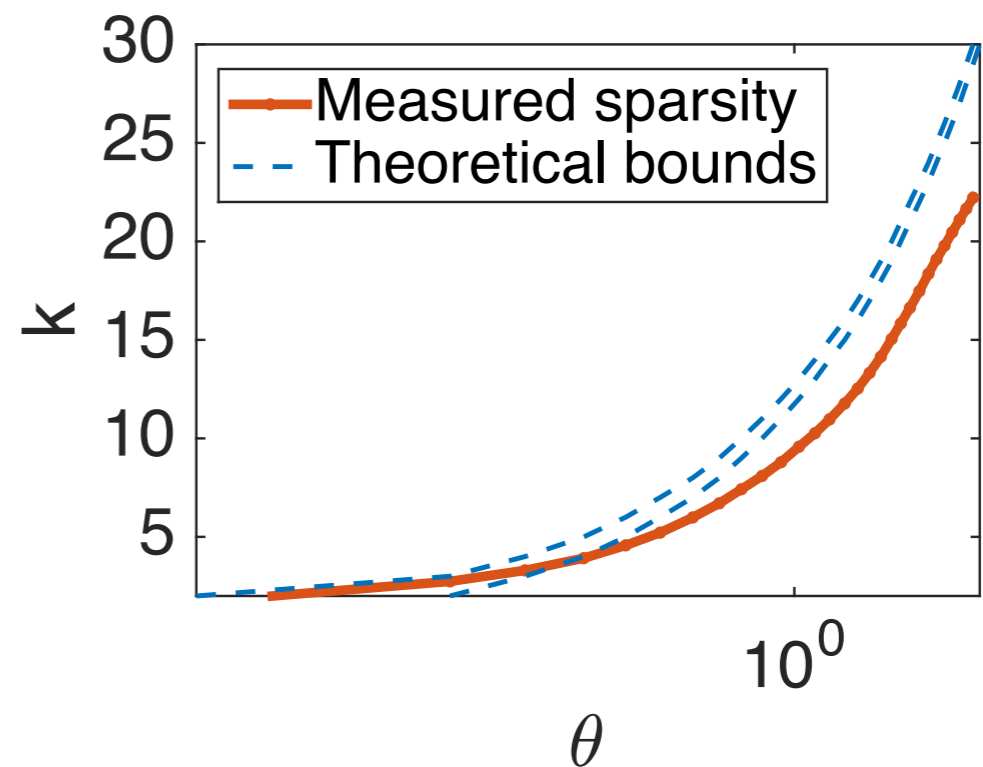


More examples



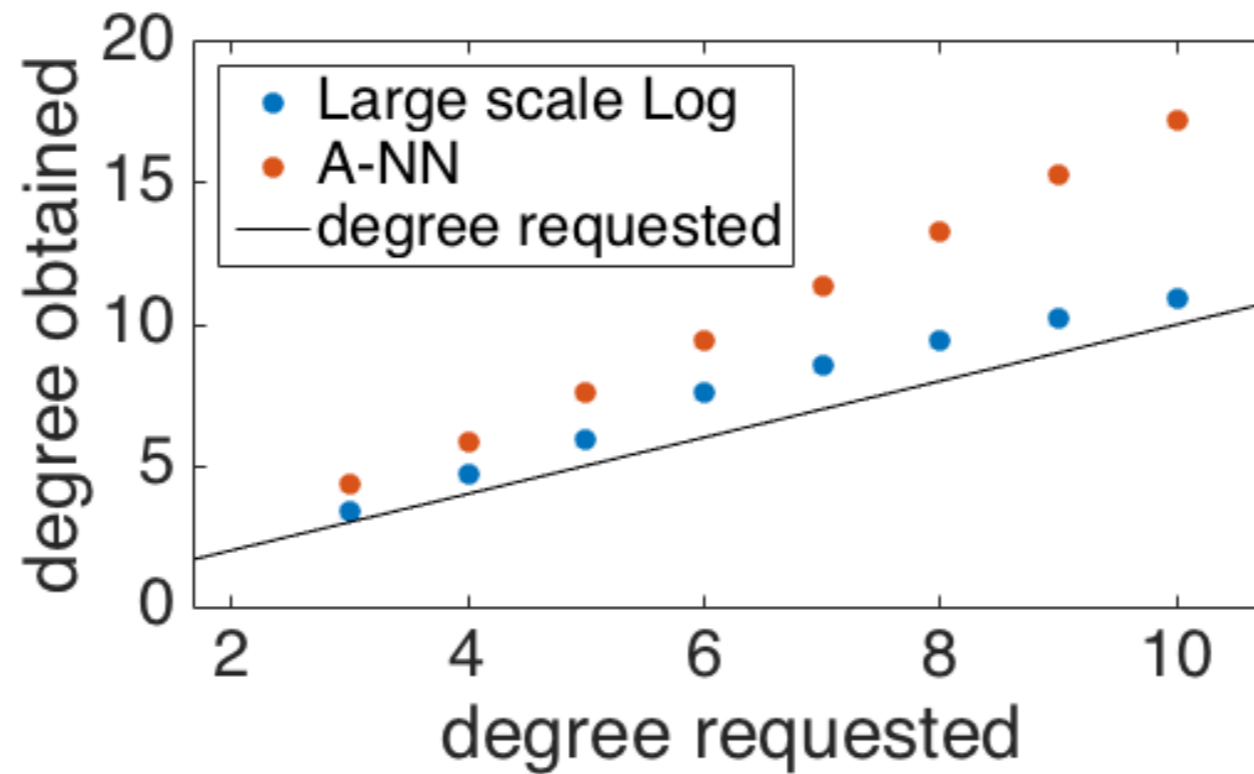
“Failing” case:

- COIL20, $n = 1440$



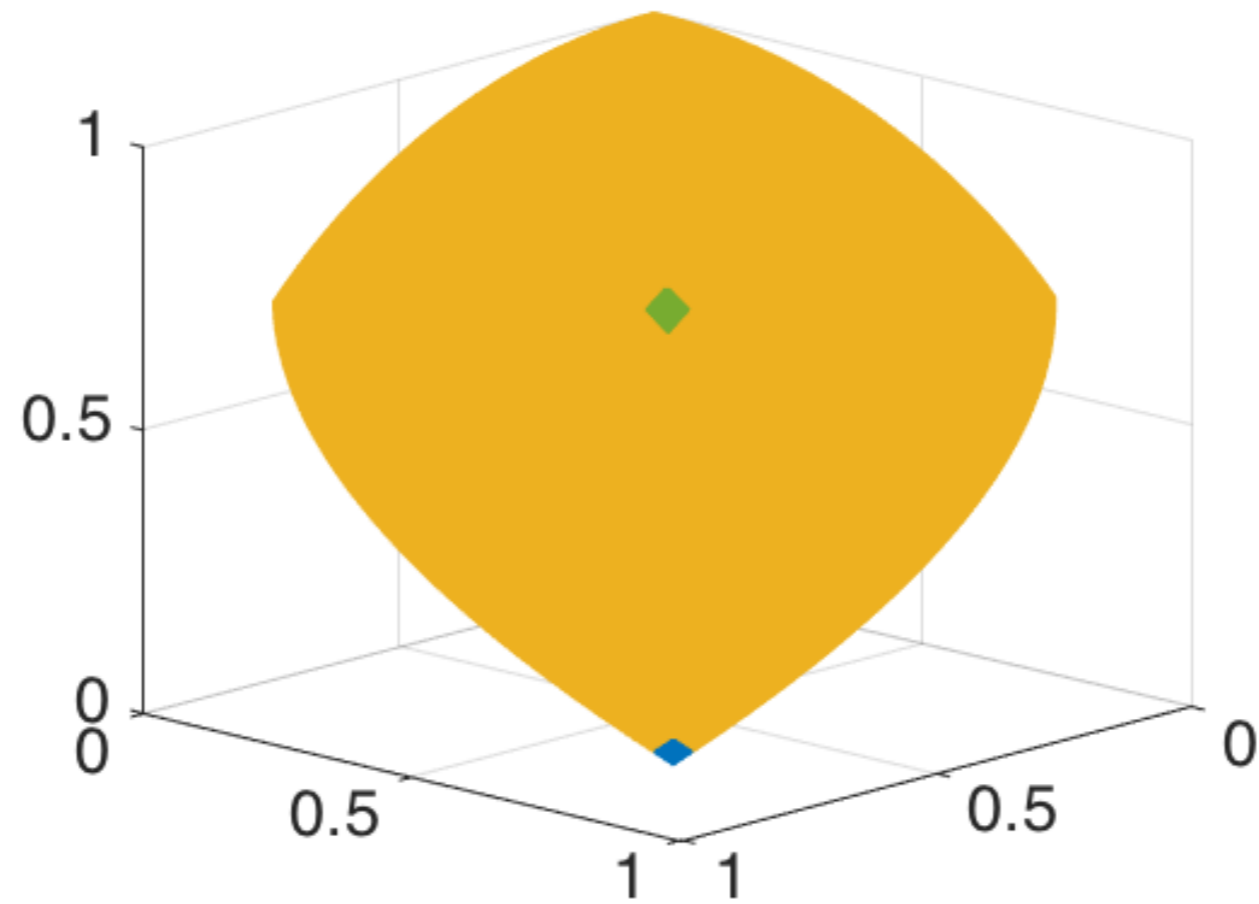
Obtained degrees

“spherical” data ($n = 260K$)



Manifold recovery

“spherical” data ($n = 260\text{K}$, $m = 2\text{K}$)

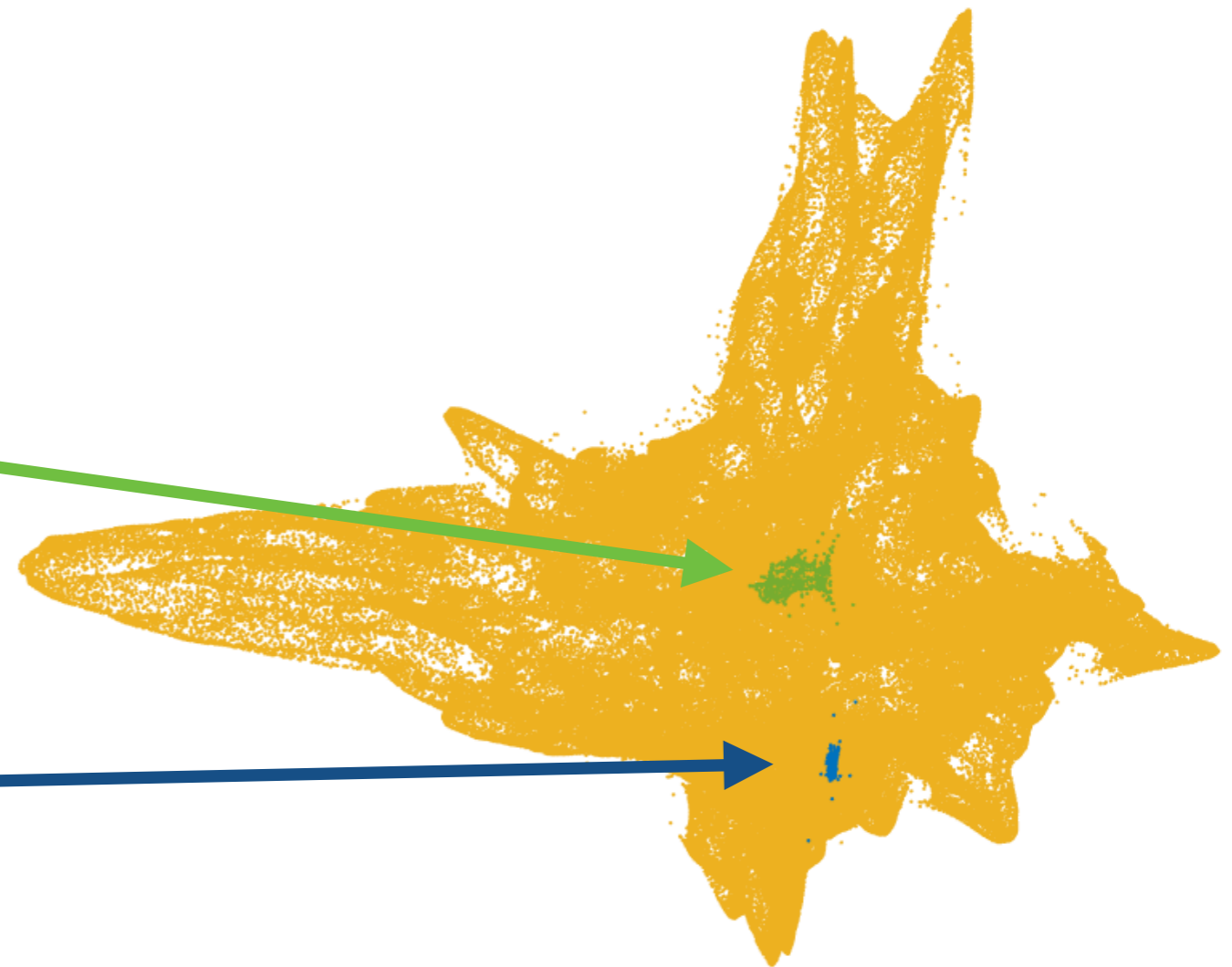
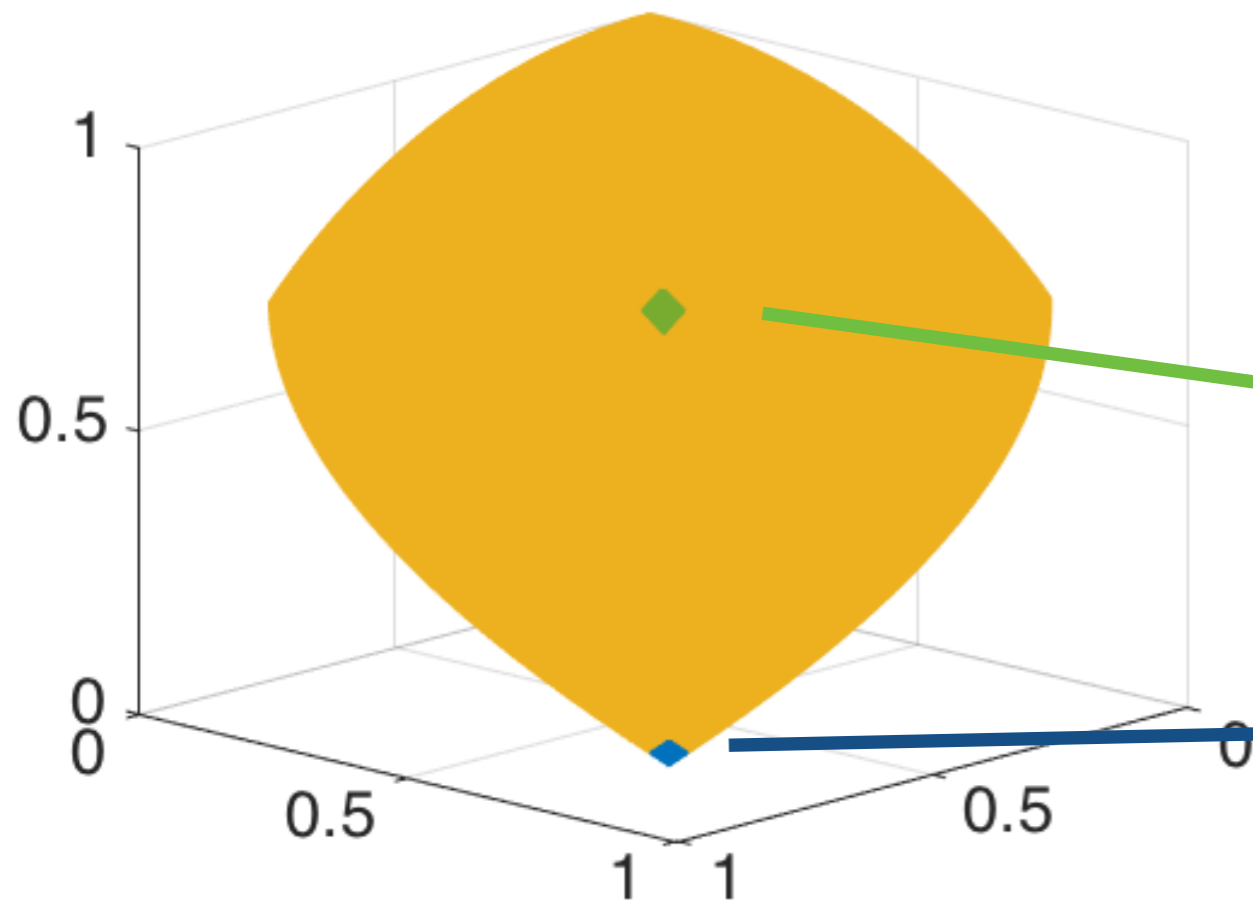


Manifold recovery

“spherical” data ($n = 260\text{K}$, $m = 2\text{K}$)

Original 2-D manifold

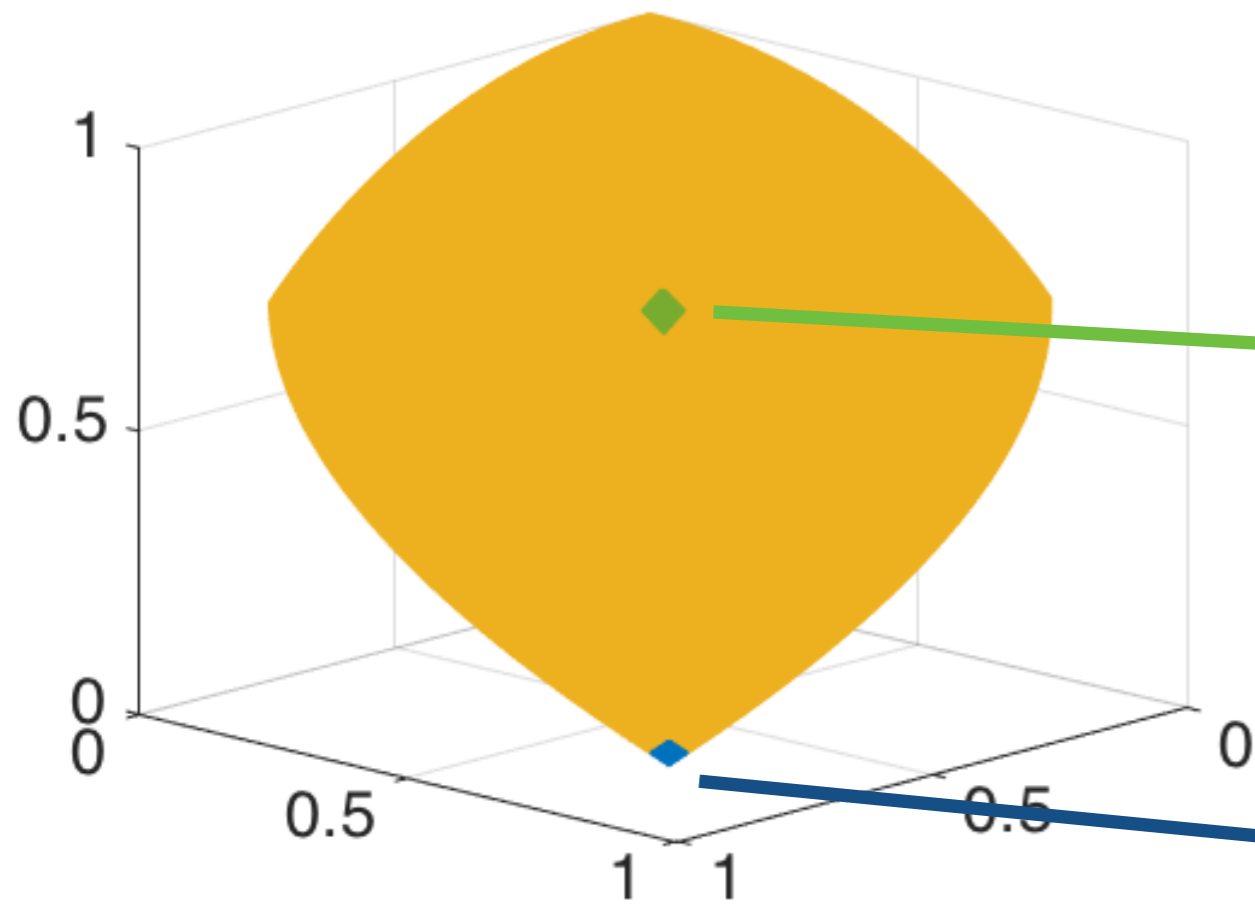
Recovered with ANN



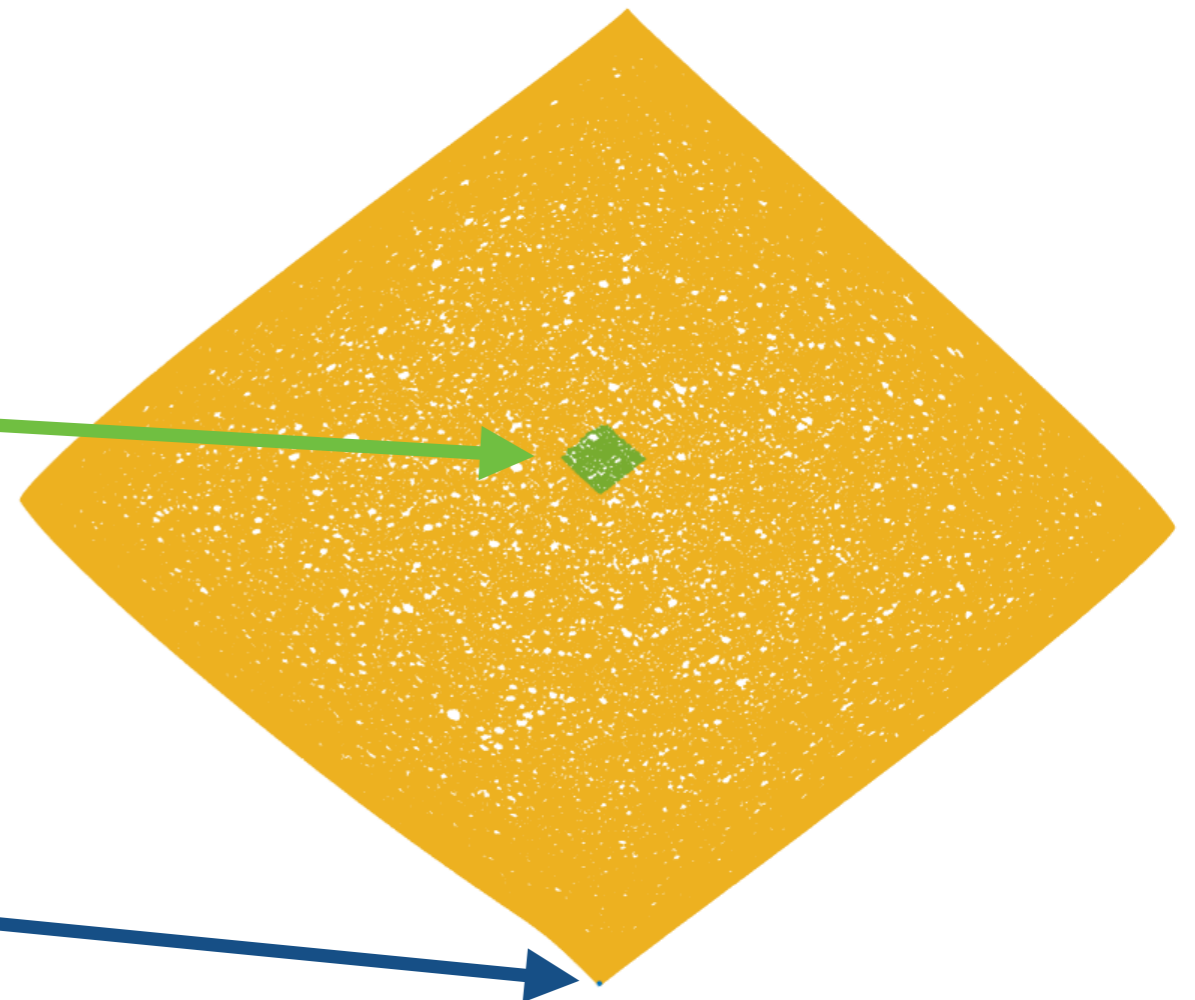
Manifold recovery

“spherical” data ($n = 260K$)

Original 2-D manifold



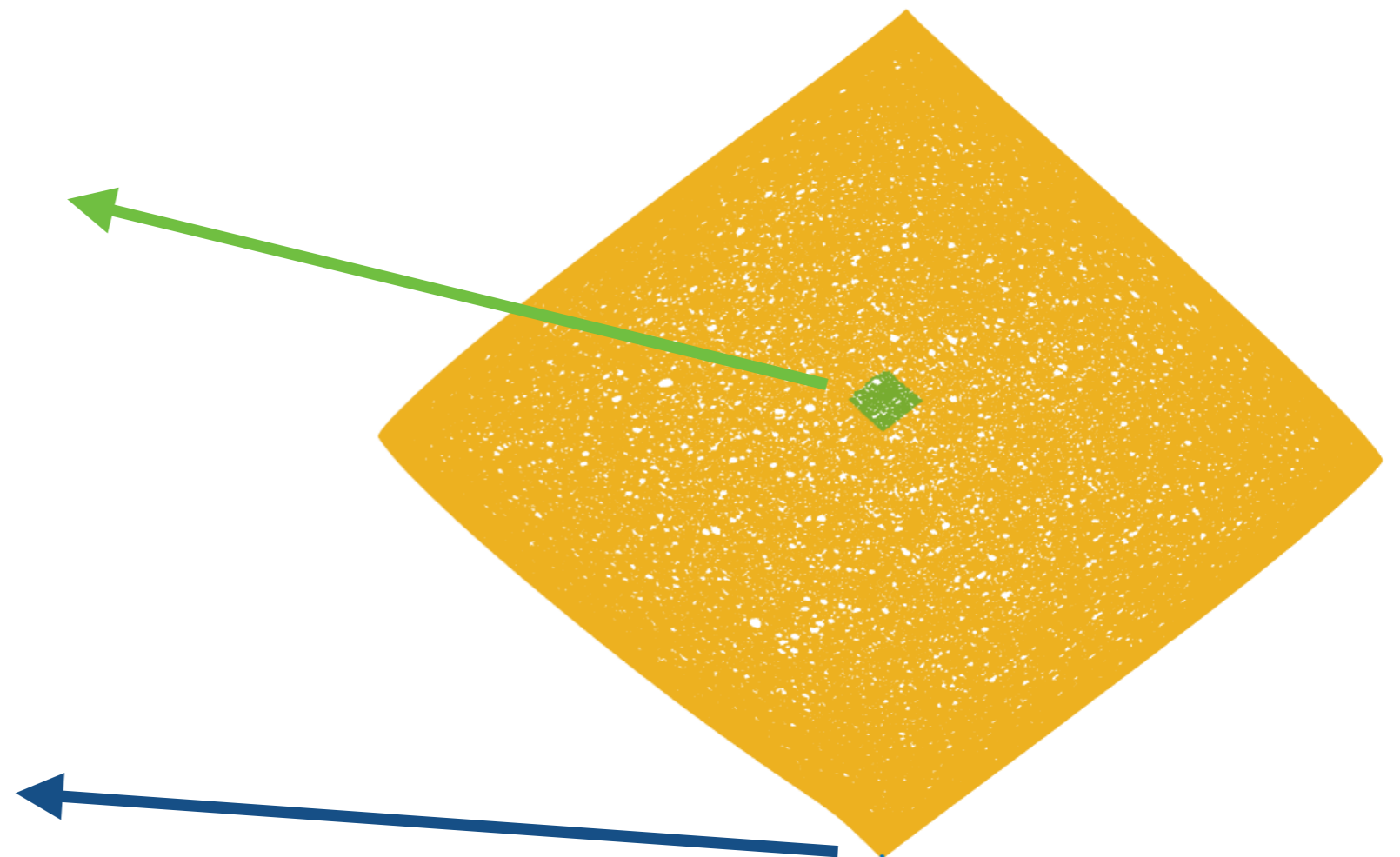
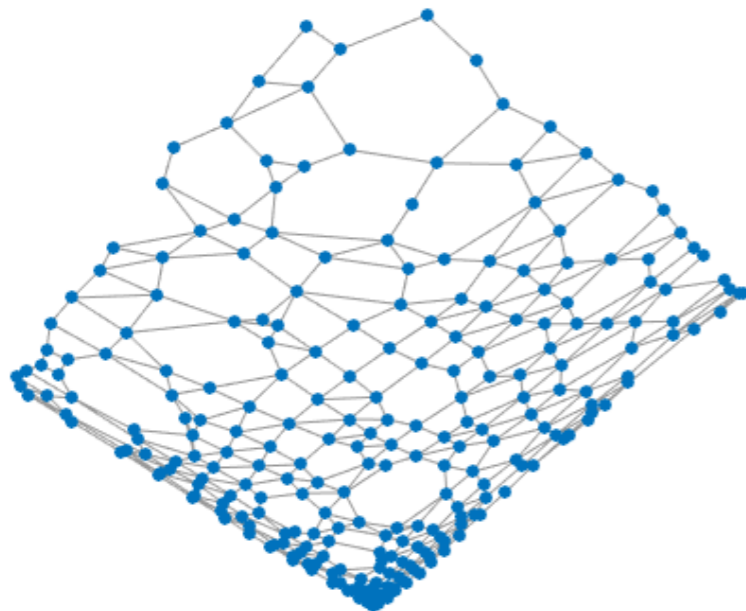
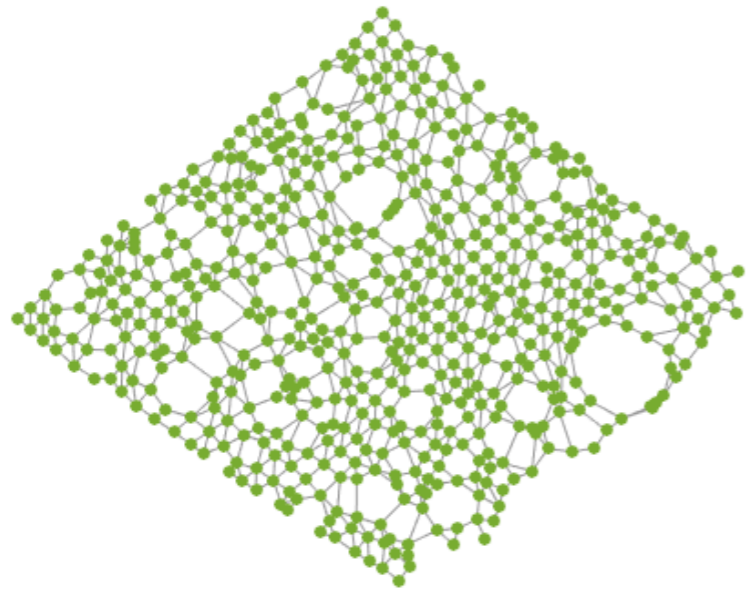
Large-scale log



Manifold recovery

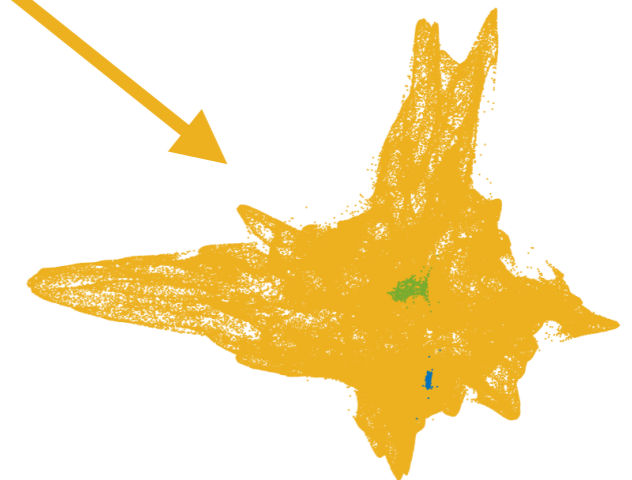
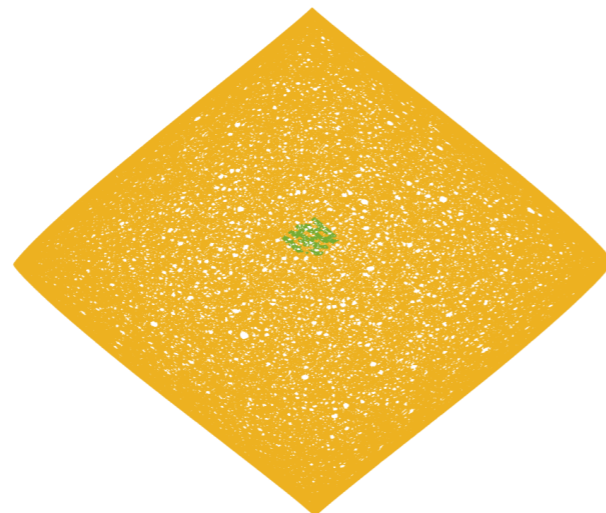
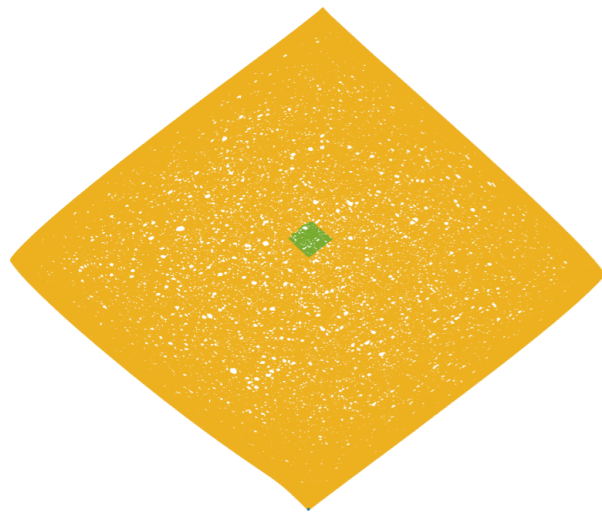
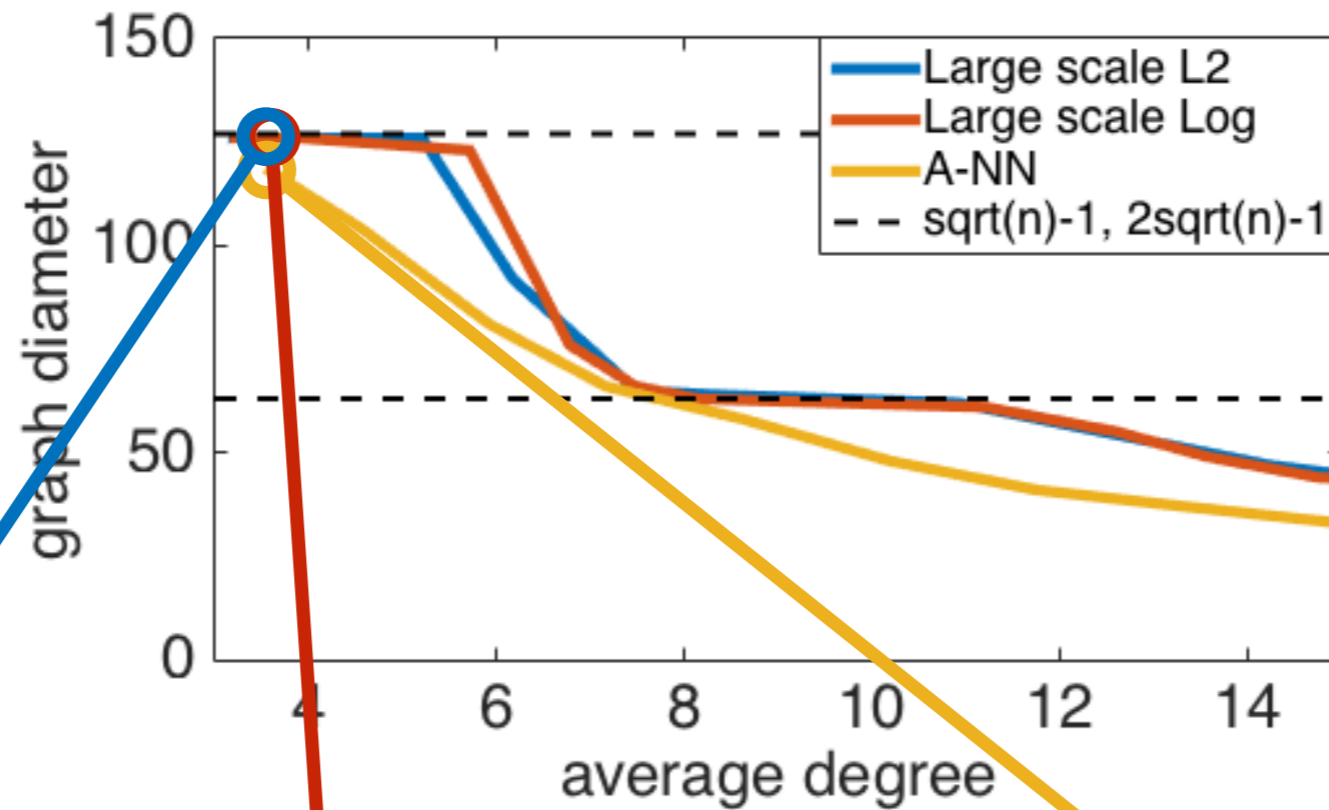
“spherical” data ($n = 260K$)

Large-scale log



Manifold recovery

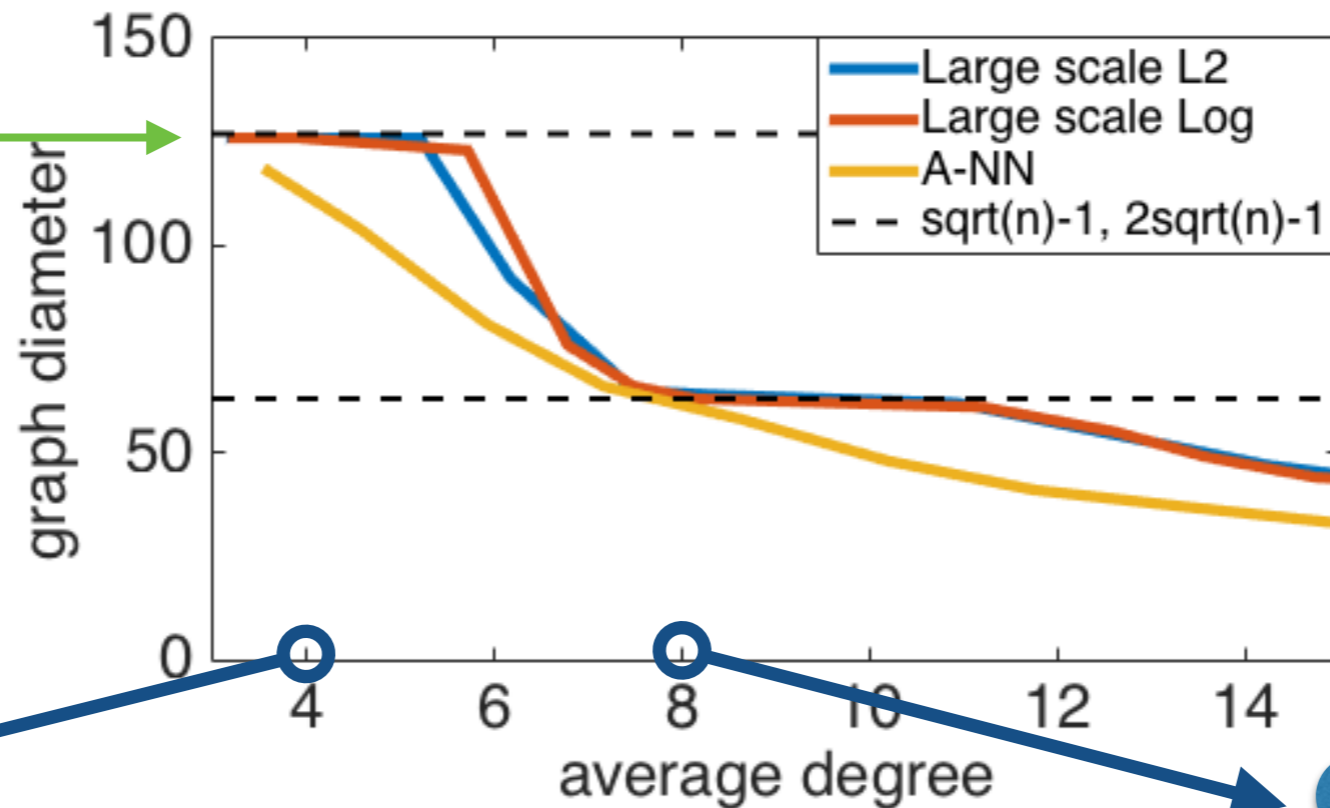
“spherical” data ($n = 4K$, $m = 2K$)



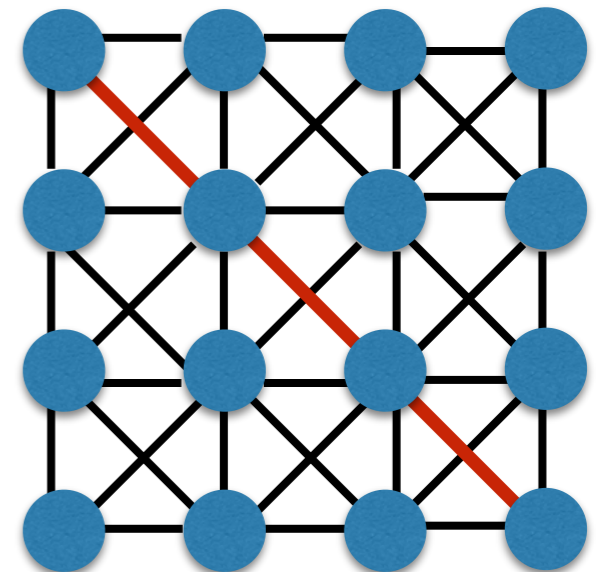
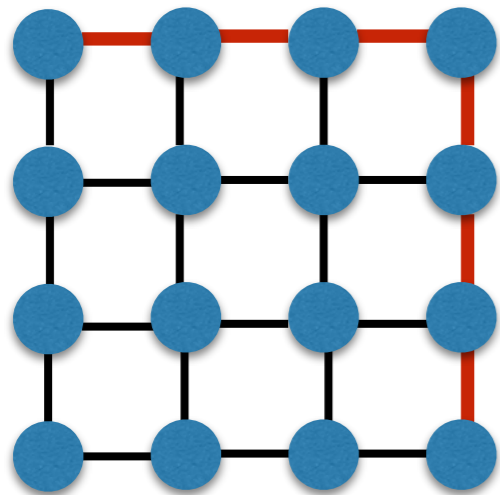
Manifold recovery

“spherical” data ($n = 4K$, $m = 2K$)

diameter =
 $2(\sqrt{n} - 1)$

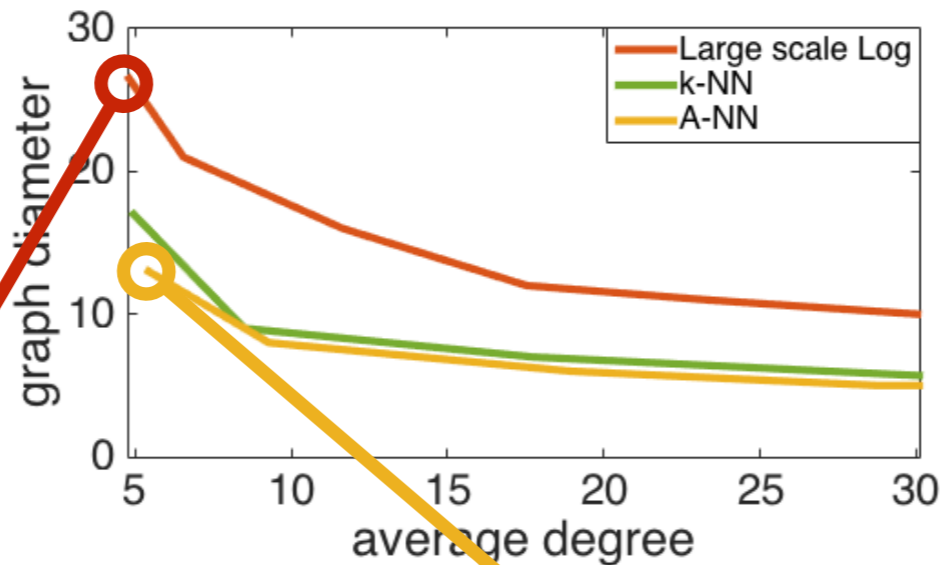


diameter
 $= \sqrt{n} - 1$



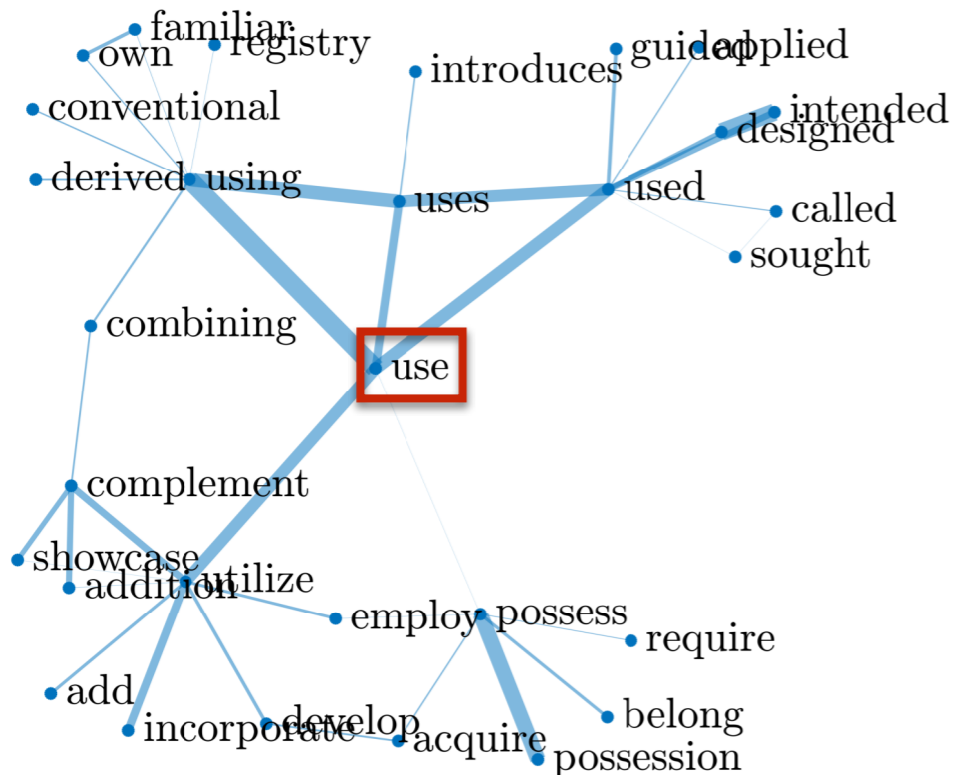
Manifold recovery

Word2vec (n = 10K, m = 300)

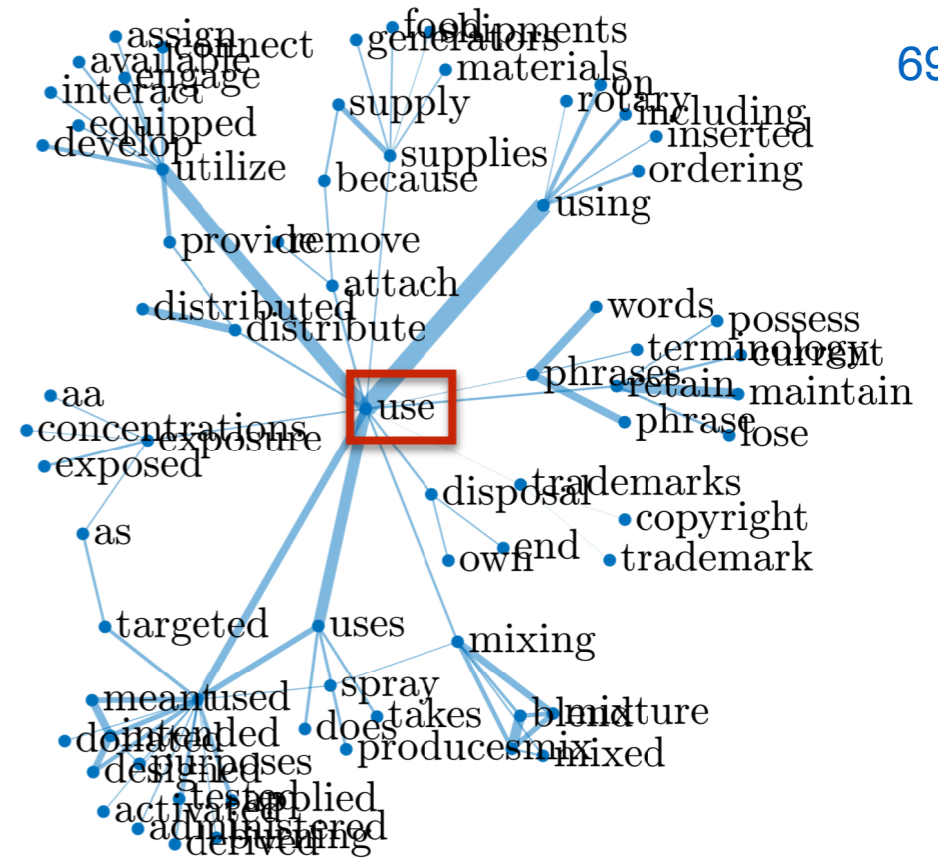


2-hops subgraph from term "use"

29 words



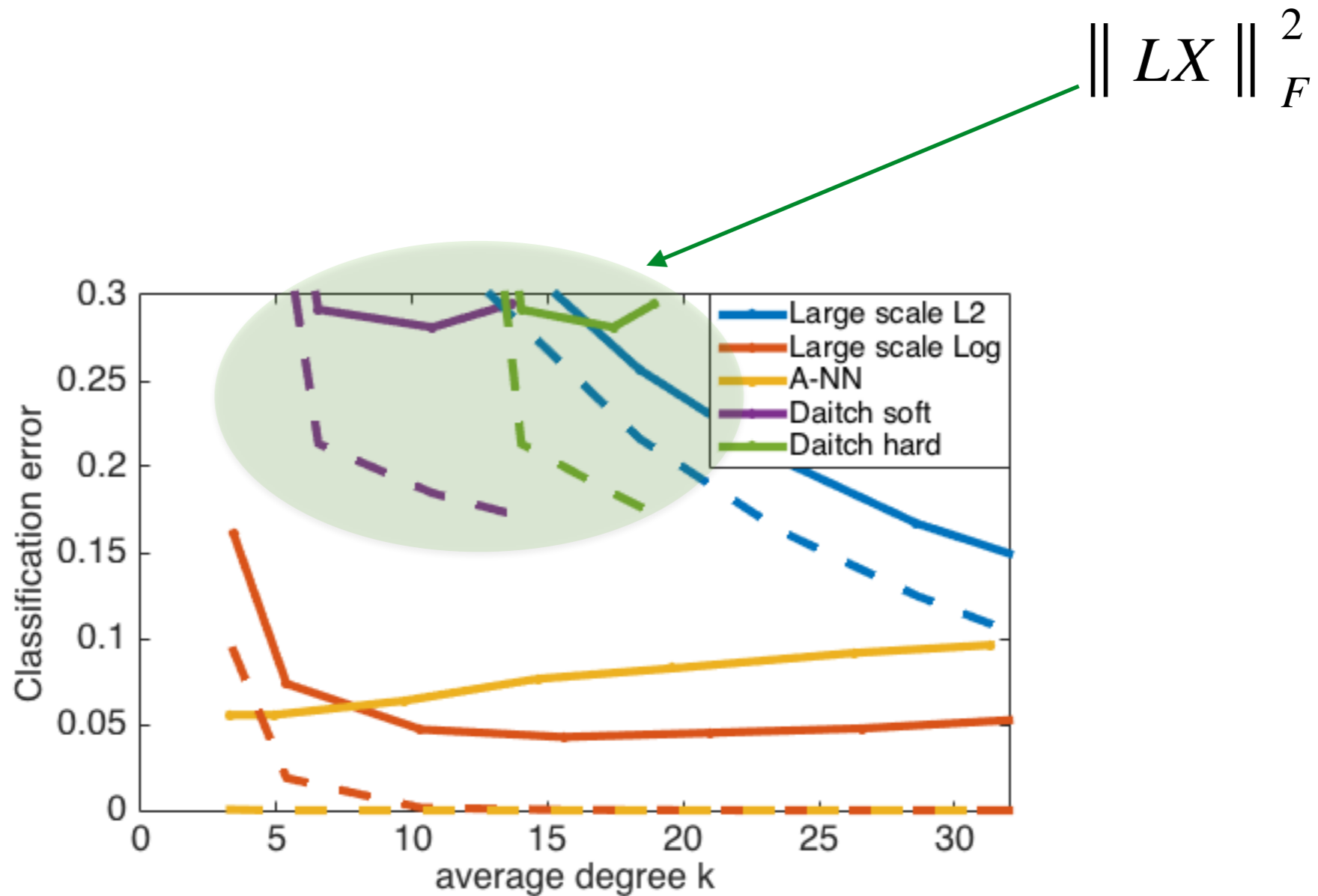
69 words



Label propagation

MNIST (n = 60K)

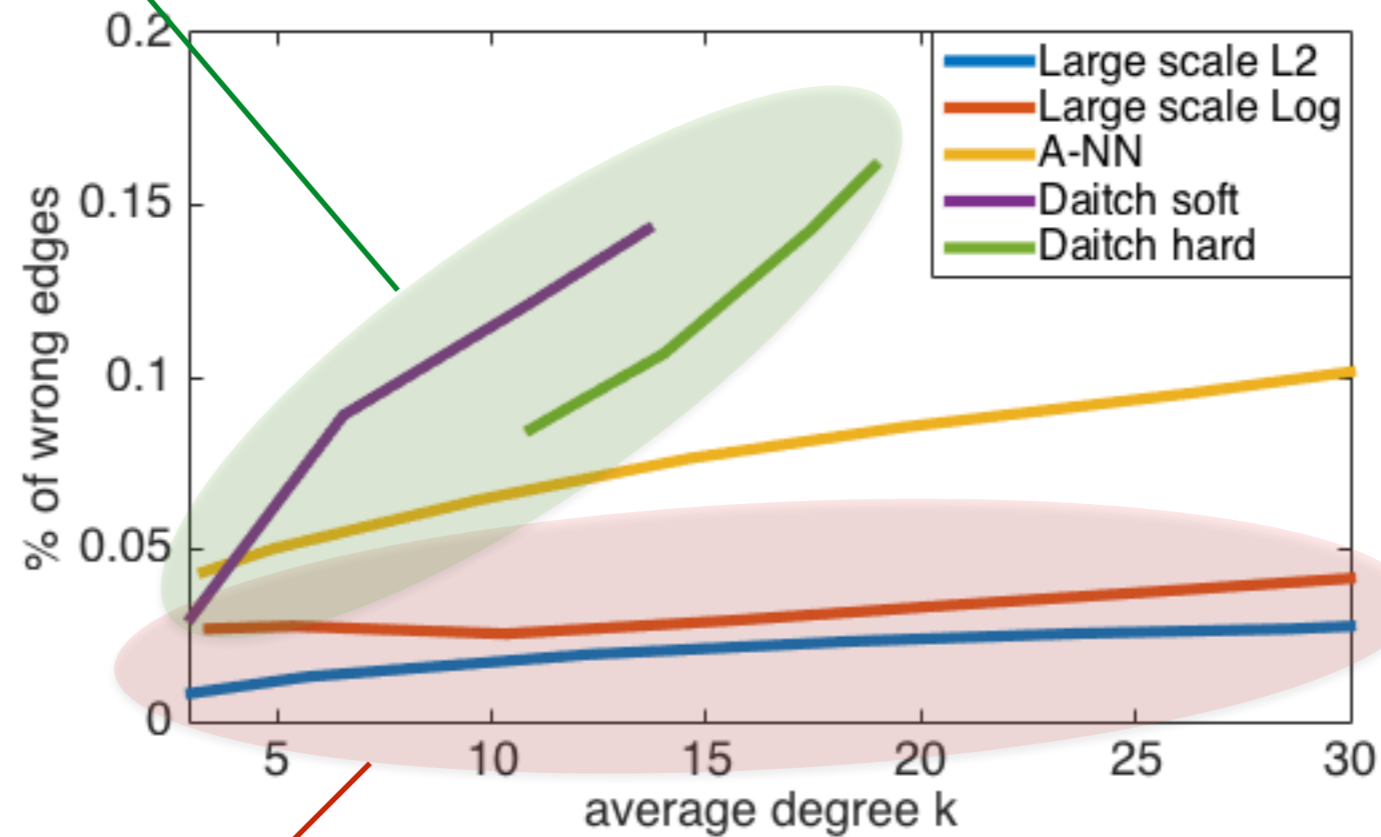
Label propagation (1% known labels)



Edge accuracy

MNIST (n = 60K)

$$\| LX \|_F^2$$



$$\text{tr}(X^T LX) = \| WZ \|_{1,1}$$

Summary

1. Good manifold recovery
2. Scalable!
 - ✓ $\mathcal{O}(nk)$ per iteration $\ll \mathcal{O}(n^2)$
 - ✓ $\mathcal{O}(n \log(n)m)$ (one time)
3. No need for parameter tuning
 - ✓ Automatic parameter selection for desired sparsity

Code: Matlab & Python (GSP box, pyGSP)

Thank you!