



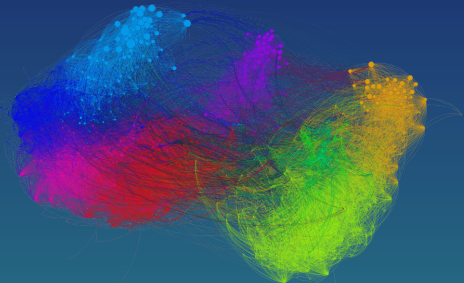
# Graphs in Machine Learning

Michal Valko

*Inria Lille - Nord Europe, France*

TA: Daniele Calandriello

Partially based on material by: Branislav Kveton,  
Partha Niyogi, Rob Fergus



# Last Lecture

- ▶ Inductive and transductive semi-supervised learning
- ▶ Manifold regularization
- ▶ Theory of Laplacian-based manifold methods
- ▶ Transductive learning stability based bounds
- ▶ Online Semi-Supervised Learning
- ▶ Online incremental  $k$ -centers

# This Lecture

- ▶ Examples of applications of online SSL
- ▶ Analysis of online SSL
- ▶ SSL Learnability
- ▶ When does graph-based SSL provably help?
- ▶ Scaling harmonic functions to millions of samples

# Previous Lab Session

- ▶ 14. 11. 2016 by Daniele Calandriello
- ▶ Content
  - ▶ Semi-supervised learning
  - ▶ Graph quantization
  - ▶ Offline face recognizer
- ▶ Install VM (in case you have not done it yet for TD1)
- ▶ Short written report
- ▶ Questions to piazza
- ▶ *Deadline: 28. 11. 2016*

## Next Lab Session/Lecture

- ▶ 28. 11. 2016 by Daniele.Calandriello@inria.fr
- ▶ Content (this time lecture in class + coding at home)
  - ▶ Large-scale graph construction and processing (in class)
  - ▶ Scalable algorithms:
    - ▶ Online face recognizer (to code in Matlab)
    - ▶ Iterative label propagation (to code in Matlab)
    - ▶ Graph sparsification (presented in class)
- ▶ AR: **record a video with faces**
- ▶ Short written report
- ▶ Questions to piazza
- ▶ **Deadline: 12. 12. 2016**
- ▶ <http://researchers.lille.inria.fr/~calandri/teaching.html>

# Final Class projects

- ▶ detailed description on the class website
- ▶ preferred option: you come up with the topic
- ▶ theory/implementation/review or a combination
- ▶ one or two people per project (exceptionally three)
- ▶ grade 60%: report + short presentation of the **team**
- ▶ deadlines
  - ▶ 21. 11. 2016 - recommended DL for taking projects **Today!**
  - ▶ 28. 11. 2016 - hard DL for taking projects
  - ▶ 05. 01. 2017 - submission of the project report
  - ▶ 09. 01. 2017 or later - project presentation
- ▶ list of suggested topics on piazza

# Online SSL with Graphs

## Video examples

<http://www.bkveton.com/videos/Coffee.mp4>

<http://www.bkveton.com/videos/Ad.mp4>

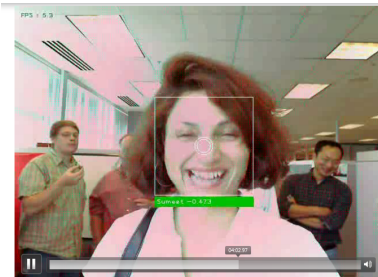
<http://researchers.lille.inria.fr/~valko/hp/serve.php?what=publications/kveton2009nipsdemo.adaptation.mov>

<http://researchers.lille.inria.fr/~valko/hp/serve.php?what=publications/kveton2009nipsdemo.officespace.mov>

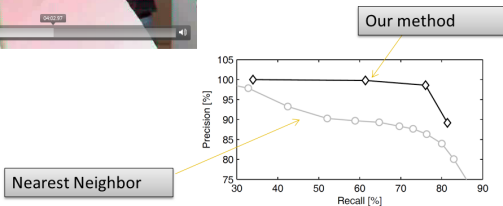
<http://bcove.me/a2derjeh>

or: <http://researchers.lille.inria.fr/~valko/hp/publications/press-intel-2015.mp4>

# SSL with Graphs: Some experimental results



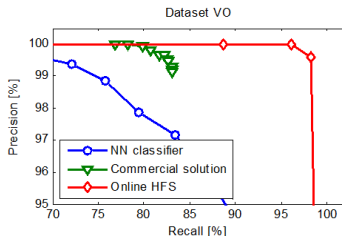
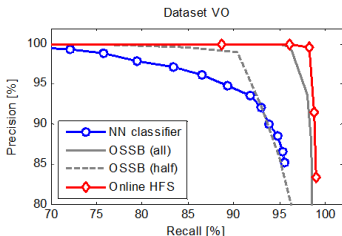
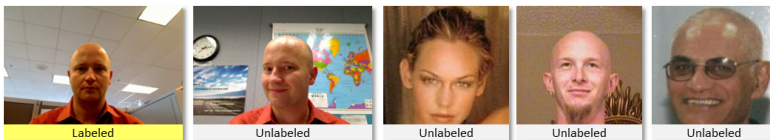
- 8 people classification
- Making funny faces
- 4 faces/person are labeled





# SSL with Graphs: Some experimental results

- One person moves among various indoor locations
- 4 labeled examples of a person in the cubicle

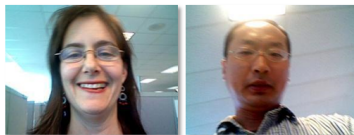


Online HFS outperforms OSSB (even when the weak learners are chosen using future data)

Online HFS yields better results than a commercial solution at 20% of the computational cost

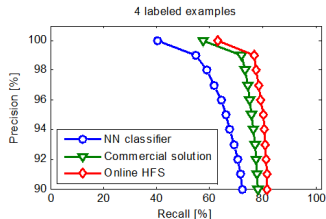
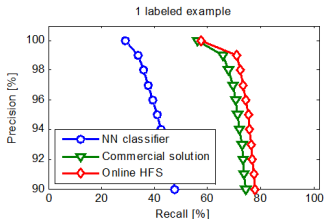
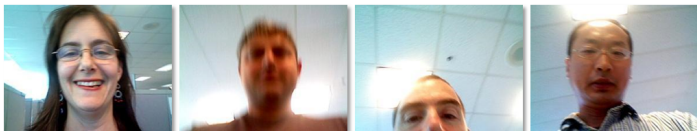
# SSL with Graphs: Some experimental results

- **Logging in** with faces instead of password
- Able to **learn** and improve



# SSL with Graphs: Some experimental results

- 16 people log twice into a tablet PC at 10 locations



Online HFS yields better results than a commercial solution at 20% of the computational cost

# Online SSL with Graphs: Analysis

What can we guarantee?

Three sources of error

- ▶ generalization error — if all data:  $(\ell_t^* - y_t)^2$
- ▶ online error — data only incrementally:  $(\ell_t^o[t] - \ell_t^*)^2$
- ▶ quantization error — memory limitation:  $(\ell_t^q[t] - \ell_t^o[t])^2$

All together:

$$\frac{1}{N} \sum_{t=1}^N (\ell_t^q[t] - y_t)^2 \leq \frac{9}{2N} \sum_{t=1}^N (\ell_t^* - y_t)^2 + \frac{9}{2N} \sum_{t=1}^N (\ell_t^o[t] - \ell_t^*)^2 + \frac{9}{2N} \sum_{t=1}^N (\ell_t^q[t] - \ell_t^o[t])^2$$

Since for any  $a, b, c, d \in [-1, 1]$ :

$$(a - b)^2 \leq \frac{9}{2} [(a - c)^2 + (c - d)^2 + (d - b)^2]$$

# Online SSL with Graphs: Analysis

Bounding **transduction error**  $(\ell_t^* - y_t)^2$

If all labeled examples  $l$  are i.i.d.,  $c_l = 1$  and  $c_l \gg c_u$ , then

$$R(\ell^*) \leq \hat{R}(\ell^*) + \underbrace{\beta + \sqrt{\frac{2 \ln(2/\delta)}{n_l}} (n_l \beta + 4)}_{\text{transductive error } \Delta_T(\beta, n_l, \delta)}$$

$$\beta \leq 2 \left[ \frac{\sqrt{2}}{\gamma_g + 1} + \sqrt{2n_l} \frac{1 - c_u}{c_u} \frac{\lambda_M(\mathbf{L}) + \gamma_g}{\gamma_g^2 + 1} \right]$$

holds with the probability of  $1 - \delta$ , where

$$R(\ell^*) = \frac{1}{N} \sum_t (\ell_t^* - y_t)^2 \quad \text{and} \quad \hat{R}(\ell^*) = \frac{1}{n_l} \sum_{t \in l} (\ell_t^* - y_t)^2$$

How should we set  $\gamma_g$ ?

# Online SSL with Graphs: Analysis

**Bounding online error**  $(\ell_t^o[t] - \ell_t^*)^2$

Idea: If  $\mathbf{L}$  and  $\mathbf{L}^o$  are regularized, then HFSs get closer together.

since they get closer to zero

Recall  $\ell = (\mathbf{C}^{-1}\mathbf{Q} + \mathbf{I})^{-1}\mathbf{y}$ , where  $\mathbf{Q} = \mathbf{L} + \gamma_g\mathbf{I}$

and also  $\mathbf{v} \in \mathbb{R}^{n \times 1}$ ,  $\lambda_m(A)\|\mathbf{v}\|_2 \leq \|\mathbf{Av}\|_2 \leq \lambda_M(A)\|\mathbf{v}\|_2$

$$\|\ell\|_2 \leq \frac{\|\mathbf{y}\|_2}{\lambda_m(\mathbf{C}^{-1}\mathbf{Q} + \mathbf{I})} = \frac{\|\mathbf{y}\|_2}{\frac{\lambda_m(\mathbf{Q})}{\lambda_M(\mathbf{C})} + 1} \leq \frac{\sqrt{n_l}}{\gamma_g + 1}$$

Difference between offline and online solutions:

$$(\ell_t^o[t] - \ell_t^*)^2 \leq \|\ell^o[t] - \ell^*\|_\infty^2 \leq \|\ell^o[t] - \ell^*\|_2^2 \leq \left(\frac{2\sqrt{n_l}}{\gamma_g + 1}\right)^2$$

Again, how should we set  $\gamma_g$ ?

# Online SSL with Graphs: Analysis

Bounding **quantization error**  $(\ell_t^q[t] - \ell_t^o[t])^2$

How are the quantized and full solution different?

$$\ell^* = \min_{\ell \in \mathbb{R}^N} (\ell - \mathbf{y})^T \mathbf{C} (\ell - \mathbf{y}) + \ell^T \mathbf{Q} \ell$$

In **Q!**  $\mathbf{Q}^o$  (online) vs.  $\mathbf{Q}^q$  (quantized)

We have:  $\ell^o = (\mathbf{C}^{-1} \mathbf{Q}^o + \mathbf{I})^{-1} \mathbf{y}$  vs.  $\ell^q = (\mathbf{C}^{-1} \mathbf{Q}^q + \mathbf{I})^{-1} \mathbf{y}$

Let  $\mathbf{Z}^q = \mathbf{C}^{-1} \mathbf{Q}^q + \mathbf{I}$  and  $\mathbf{Z}^o = \mathbf{C}^{-1} \mathbf{Q}^o + \mathbf{I}$ .

$$\begin{aligned} \ell^q - \ell^o &= (\mathbf{Z}^q)^{-1} \mathbf{y} - (\mathbf{Z}^o)^{-1} \mathbf{y} = (\mathbf{Z}^q \mathbf{Z}^o)^{-1} (\mathbf{Z}^o - \mathbf{Z}^q) \mathbf{y} \\ &= (\mathbf{Z}^q \mathbf{Z}^o)^{-1} \mathbf{C}^{-1} (\mathbf{Q}^o - \mathbf{Q}^q) \mathbf{y} \end{aligned}$$

# Online SSL with Graphs: Analysis

Bounding quantization error  $(\ell_t^q[t] - \ell_t^o[t])^2$

$$\begin{aligned}\ell^q - \ell^o &= (\mathbf{Z}^q)^{-1}\mathbf{y} - (\mathbf{Z}^o)^{-1}\mathbf{y} = (\mathbf{Z}^q\mathbf{Z}^o)^{-1}(\mathbf{Z}^o - \mathbf{Z}^q)\mathbf{y} \\ &= (\mathbf{Z}^q\mathbf{Z}^o)^{-1}\mathbf{C}^{-1}(\mathbf{Q}^o - \mathbf{Q}^q)\mathbf{y}\end{aligned}$$

$$\|\ell^q - \ell^o\|_2 \leq \frac{\lambda_M(\mathbf{C}^{-1})\|(\mathbf{Q}^q - \mathbf{Q}^o)\mathbf{y}\|_2}{\lambda_m(\mathbf{Z}^q)\lambda_m(\mathbf{Z}^o)}$$

$\|\cdot\|_F$  and  $\|\cdot\|_2$  are compatible and  $y_i$  is zero when unlabeled:

$$\|(\mathbf{Q}^q - \mathbf{Q}^o)\mathbf{y}\|_2 \leq \|\mathbf{Q}^q - \mathbf{Q}^o\|_F \cdot \|\mathbf{y}\|_2 \leq \sqrt{n_l}\|\mathbf{Q}^q - \mathbf{Q}^o\|_F$$

Furthermore,  $\lambda_m(\mathbf{Z}^o) \geq \frac{\lambda_m(\mathbf{Q}^o)}{\lambda_M(\mathbf{C})} + 1 \geq \gamma_g$  and  $\lambda_M(\mathbf{C}^{-1}) \leq c_u^{-1}$

$$\text{We get } \|\ell^q - \ell^o\|_2 \leq \frac{\sqrt{n_l}}{c_u\gamma_g^2} \|\mathbf{Q}^q - \mathbf{Q}^o\|_F$$



# Online SSL with Graphs: Analysis

**Bounding quantization error**  $(\ell_t^q[t] - \ell_t^o[t])^2$

The quantization error depends on  $\|\mathbf{Q}^q - \mathbf{Q}^o\|_F = \|\mathbf{L}^q - \mathbf{L}^o\|_F$ .

When can we keep  $\|\mathbf{L}^q - \mathbf{L}^o\|_F$  under control?

Charikar guarantees **distortion** error of at most  $Rm/(m-1)$

For what kind of data  $\{\mathbf{x}_i\}_{i=1,\dots,n}$  is the distortion small?

Assume manifold  $\mathcal{M}$

- ▶ all  $\{\mathbf{x}_i\}_{i \geq 1}$  lie on a smooth  $s$ -dimensional compact  $\mathcal{M}$
- ▶ with boundary of bounded geometry Def. 11 of Hein [HAL07]
  - ▶ should not intersect itself
  - ▶ should not fold back onto itself
  - ▶ has finite volume  $V$
  - ▶ has finite surface area  $A$

## Online SSL with Graphs: Analysis

Bounding **quantization error**  $(\ell_t^q[t] - \ell_t^o[t])^2$

Bounding  $\|\mathbf{L}^q - \mathbf{L}^o\|_F$  when  $\mathbf{x}_i \in \mathcal{M}$

Consider  $k$ -sphere packing of radius  $r$  with centers contained in  $\mathcal{M}$ .

What is the maximum volume of this packing?

$kc_s r^s \leq V + Ac_{\mathcal{M}}r$  with  $c_s, c_{\mathcal{M}}$  depending on dimension and  $\mathcal{M}$ .

If  $k$  is large  $\rightarrow r < \mathbf{injectivity\ radius}$  of  $\mathcal{M}$  [HAL07] and  $r < 1$ :

$$r < ((V + Ac_{\mathcal{M}}) / (kc_s))^{1/s} = \mathcal{O}(k^{-1/s})$$

$r$ -packing is a  $2r$ -covering:

$$\max_{i=1, \dots, N} \|\mathbf{x}_i - \mathbf{c}\|_2 \leq Rm / (m-1) \leq 2(1+\varepsilon)\mathcal{O}(k^{-1/s}) = \mathcal{O}(k^{-1/s})$$

But what about  $\|\mathbf{L}^q - \mathbf{L}^o\|_F$ ?

## Online SSL with Graphs: Analysis

Bounding **quantization error**  $(\ell_t^q[t] - \ell_t^o[t])^2$

If similarity is  $M$ -Lipschitz,  $\mathbf{L}$  is normalized,

$$c_{ij}^o = \sqrt{\mathbf{D}_{ii}^o \mathbf{D}_{jj}^o} > c_{min} N:$$

$$\begin{aligned} \mathbf{L}_{ij}^q - \mathbf{L}_{ij}^o &= \frac{\mathbf{W}_{ij}^q}{c_{ij}^q} - \frac{\mathbf{W}_{ij}^o}{c_{ij}^o} \\ &\leq \frac{\mathbf{W}_{ij}^q - \mathbf{W}_{ij}^o}{c_{ij}^q} + \frac{\mathbf{W}_{ij}^q (c_{ij}^q - c_{ij}^o)}{c_{ij}^o c_{ij}^q} \\ &\leq \frac{4MRm}{(m-1)c_{min}N} + \frac{4M(NMRm)}{((m-1)c_{min}N)^2} \\ &= O\left(\frac{R}{N}\right) \end{aligned}$$

Finally,  $\|\mathbf{L}^q - \mathbf{L}^o\|_F^2 \leq N^2 \mathcal{O}(R^2/N^2) = \mathcal{O}(k^{-2/s})$ .

Are the assumptions reasonable?

## Online SSL with Graphs: Analysis

Bounding quantization error  $(\ell_t^q[t] - \ell_t^o[t])^2$

We showed  $\|\mathbf{L}^q - \mathbf{L}^o\|_F^2 \leq N^2 \mathcal{O}(R^2/N^2) = \mathcal{O}(k^{-2/s}) = \mathcal{O}(1)$ .

$$\frac{1}{N} \sum_{t=1}^N (\ell_t^q[t] - \ell_t^o[t])^2 \leq \frac{n_l}{c_u^2 \gamma_g^4} \|\mathbf{L}^q - \mathbf{L}^o\|_F^2 \leq \frac{n_l}{c_u^2 \gamma_g^4}$$

This converges to zero at the rate of  $\mathcal{O}(N^{-1/2})$  with  $\gamma_g = \Omega(N^{1/8})$ .

With properly setting  $\gamma_g$ , e.g.,  $\gamma_g = \Omega(N^{1/8})$ , we can have:

$$\frac{1}{N} \sum_{t=1}^N (\ell_t^q[t] - y_t)^2 = \mathcal{O}(N^{-1/2})$$

What does that mean?

# SSL with Graphs: What is behind it?

Why and when it helps?

Can we guarantee benefit of SSL over SL?

Are there cases when **manifold** SSL is provably helpful?

Say  $\mathcal{X}$  is supported on manifold  $\mathcal{M}$ . Compare two cases:

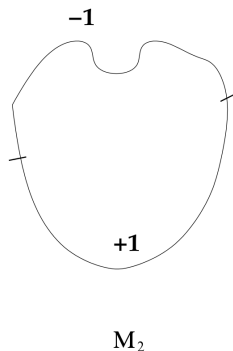
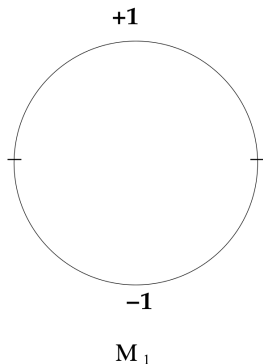
- ▶ SL: does not know about  $\mathcal{M}$  and only knows  $(\mathbf{x}_i, y_i)$
- ▶ SSL: perfect knowledge of  $\mathcal{M} \equiv$  humongous amounts of  $\mathbf{x}_i$

<http://people.cs.uchicago.edu/~niyogi/papersps/ssminimax2.pdf>

## SSL with Graphs: What is behind it?

Set of learning problems - collections  $\mathcal{P}$  of probability distributions:

$$\mathcal{P} = \cup_{\mathcal{M}} \mathcal{P}_{\mathcal{M}} = \cup_{\mathcal{M}} \{p \in \mathcal{P} \mid p_{\mathcal{X}} \text{ is uniform on } \mathcal{M}\}$$



## SSL with Graphs: What is behind it?

**Set of problems**  $\mathcal{P} = \cup_{\mathcal{M}} \mathcal{P}_{\mathcal{M}} = \{p \in \mathcal{P} \mid p_{\mathcal{X}} \text{ is uniform on } \mathcal{M}\}$

**Regression function**  $m_p = \mathbb{E}[y \mid x]$  when  $x \in \mathcal{M}$

**Algorithm A** and **labeled examples**  $\bar{z} = \{z_i\}_{i=1}^{n_l} = \{(\mathbf{x}_i, y_i)\}_{i=1}^{n_l}$

**Minimax rate**

$$R(n_l, \mathcal{P}) = \inf_A \sup_{p \in \mathcal{P}} \mathbb{E}_{\bar{z}} [\|A(\bar{z}) - m_p\|_{L^2(p_{\mathcal{X}})}]$$

Since  $\mathcal{P} = \cup_{\mathcal{M}} \mathcal{P}_{\mathcal{M}}$

$$R(n_l, \mathcal{P}) = \inf_A \sup_{\mathcal{M}} \sup_{p \in \mathcal{P}_{\mathcal{M}}} \mathbb{E}_{\bar{z}} [\|A(\bar{z}) - m_p\|_{L^2(p_{\mathcal{X}})}]$$

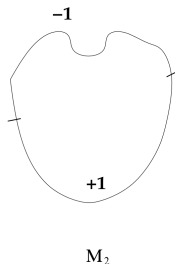
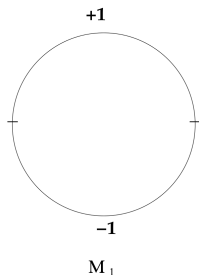
(SSL) When A is allowed to know  $\mathcal{M}$

$$Q(n_l, \mathcal{P}) = \sup_{\mathcal{M}} \inf_A \sup_{p \in \mathcal{P}_{\mathcal{M}}} \mathbb{E}_{\bar{z}} [\|A(\bar{z}) - m_p\|_{L^2(p_{\mathcal{X}})}]$$

In which cases there is a gap between  $Q(n_l, \mathcal{P})$  and  $R(n_l, \mathcal{P})$ ?

# SSL with Graphs: What is behind it?

Hypothesis space  $\mathcal{H}$ : half of the circle as +1 and the rest as -1



**Case 1:**  $\mathcal{M}$  is known to the learner ( $\mathcal{H}_{\mathcal{M}}$ )

What is a VC dimension of  $\mathcal{H}_{\mathcal{M}}$ ?

$$\text{Optimal rate } Q(n, \mathcal{P}) \leq 2\sqrt{\frac{3 \log n_I}{n_I}}$$



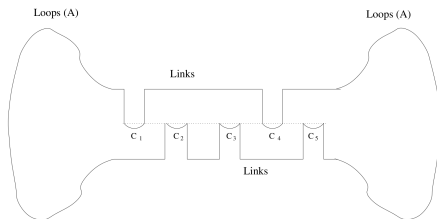
# SSL with Graphs: What is behind it?

Case 2:  $\mathcal{M}$  is **unknown** to the learner

$$R(n_l, \mathcal{P}) = \inf_A \sup_{p \in \mathcal{P}} \mathbb{E}_{\bar{z}} [\|A(\bar{z}) - m_p\|_{L^2(\rho_X)}] = \Omega(1)$$

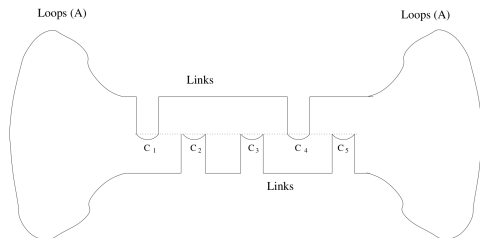
We consider  $2^d$  manifolds of the form

$$\mathcal{M} = \text{Loops} \cup \text{Links} \cup C \text{ where } C = \cup_{i=1}^d C_i$$



**Main idea:**  $d$  segments in  $C$ ,  $d - l$  with no data,  $2^l$  possible choices for labels, which helps us to lower bound  $\|A(\bar{z}) - m_p\|_{L^2(\rho_X)}$

# SSL with Graphs: What is behind it?



## Knowing the manifold helps

- ▶  $C_1$  and  $C_4$  are close
- ▶  $C_1$  and  $C_3$  are far
- ▶ we also need: **target function varies smoothly**
- ▶ altogether: **closeness on manifold** → **similarity in labels**

# SSL with Graphs: What is behind it?

What does it mean to **know**  $\mathcal{M}$ ?

## Different degrees of knowing $\mathcal{M}$

- ▶ set membership oracle:  $\mathbf{x} \stackrel{?}{\in} \mathcal{M}$
- ▶ approximate oracle
- ▶ knowing the harmonic functions on  $\mathcal{M}$
- ▶ knowing the Laplacian  $\mathcal{L}_{\mathcal{M}}$
- ▶ knowing eigenvalues and *eigenfunctions*
- ▶ topological invariants, e.g., dimension
- ▶ metric information: geodesic distance

# Scaling SSL with Graphs to Millions

Semi-supervised learning with graphs

$$\mathbf{f}^* = \min_{\mathbf{f} \in \mathbb{R}^N} (\mathbf{f} - \mathbf{y})^T \mathbf{C} (\mathbf{f} - \mathbf{y}) + \mathbf{f}^T \mathbf{L} \mathbf{f}$$

Let us see the same in eigenbasis of  $\mathbf{L} = \mathbf{U}\mathbf{\Lambda}\mathbf{U}^T$ , i.e.,  $\mathbf{f} = \mathbf{U}\boldsymbol{\alpha}$

$$\boldsymbol{\alpha}^* = \min_{\boldsymbol{\alpha} \in \mathbb{R}^N} (\mathbf{U}\boldsymbol{\alpha} - \mathbf{y})^T \mathbf{C} (\mathbf{U}\boldsymbol{\alpha} - \mathbf{y}) + \boldsymbol{\alpha}^T \boldsymbol{\Lambda} \boldsymbol{\alpha}$$

What is the problem with scalability?

Diagonalization of  $N \times N$  matrix

What can we do? Let's take only first  $k$  eigenvectors  $\mathbf{f} = \mathbf{U}\boldsymbol{\alpha}$ !

# Scaling SSL with Graphs to Millions

$\mathbf{U}$  is now a  $n \times k$  matrix

$$\alpha^* = \min_{\alpha \in \mathbb{R}^N} (\mathbf{U}\alpha - \mathbf{y})^T \mathbf{C} (\mathbf{U}\alpha - \mathbf{y}) + \alpha^T \mathbf{\Lambda} \alpha$$

Closed form solution is  $(\mathbf{\Lambda} + \mathbf{U}^T \mathbf{C} \mathbf{U}) \alpha = \mathbf{U}^T \mathbf{C} \mathbf{y}$

What is the size of this system of equation now?

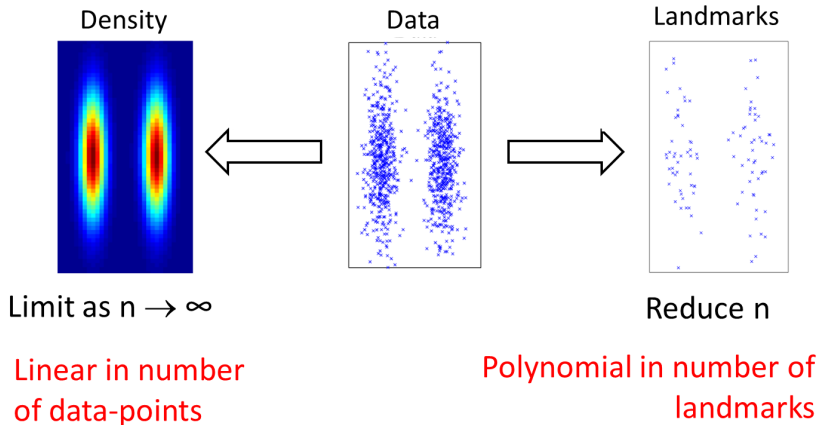
Cool!

Any problem with this approach?

Are there any reasonable assumptions when this is feasible?

Let's see what happens when  $N \rightarrow \infty$ !

# Scaling SSL with Graphs to Millions



[https://cs.nyu.edu/~fergus/papers/fwt\\_ssl.pdf](https://cs.nyu.edu/~fergus/papers/fwt_ssl.pdf)

# Scaling SSL with Graphs to Millions

What happens to  $\mathbf{L}$  when  $N \rightarrow \infty$ ?

We have data  $\mathbf{x}_i \in \mathbb{R}$  sampled from  $p(\mathbf{x})$ .

When  $n \rightarrow \infty$ , instead of vectors  $\mathbf{f}$ , we consider functions  $F(\mathbf{x})$ .

Instead of  $\mathbf{L}$ , we define  $\mathcal{L}_p$  - **weighted smoothness operator**

$$\mathcal{L}_p(F) = \frac{1}{2} \int (F(\mathbf{x}_1) - F(\mathbf{x}_2))^2 W(\mathbf{x}_1, \mathbf{x}_2) p(\mathbf{x}_1) p(\mathbf{x}_2) d\mathbf{x}_1 d\mathbf{x}_2$$

$$\text{with } W(\mathbf{x}_1, \mathbf{x}_2) = \frac{\exp(-\|\mathbf{x}_1 - \mathbf{x}_2\|^2)}{2\sigma^2}$$

$\mathbf{L}$  defined the eigenvectors of increasing smoothness.

What defines  $\mathcal{L}_p$ ? **Eigenfunctions!**

# Scaling SSL with Graphs to Millions

$$\mathcal{L}_p(F) = \frac{1}{2} \int (F(\mathbf{x}_1) - F(\mathbf{x}_2))^2 W(\mathbf{x}_1, \mathbf{x}_2) p(\mathbf{x}_1) p(\mathbf{x}_2) d\mathbf{x}_1 d\mathbf{x}_2$$

First eigenfunction

$$\Phi_1 = \underset{F: \int F^2(\mathbf{x}) p(\mathbf{x}) D(\mathbf{x}) d\mathbf{x} = 1}{\text{arg min}} \mathcal{L}_p(F)$$

where  $D(\mathbf{x}) = \int_{\mathbf{x}_2} W(\mathbf{x}, \mathbf{x}_2) p(\mathbf{x}_2) d\mathbf{x}_2$

What is the solution?  $\Phi_1(\mathbf{x}) = 1$  because  $\mathcal{L}_p(1) = 0$

How to define  $\Phi_2$ ? Same, constraining to be orthogonal to  $\Phi_1$

$$\int F(\mathbf{x}) \Phi_1(\mathbf{x}) p(\mathbf{x}) D(\mathbf{x}) d\mathbf{x} = 0$$



# Scaling SSL with Graphs to Millions

## Eigenfunctions of $\mathcal{L}_p$

$\Phi_3$  as before, orthogonal to  $\Phi_1$  and  $\Phi_2$  etc.

How to define eigenvalues?  $\lambda_k = \mathcal{L}_p(\Phi_k)$

Relationship to the discrete Laplacian

$$\frac{1}{N^2} \mathbf{f}^T \mathbf{L} \mathbf{f} = \frac{1}{2N^2} \sum_{ij} W_{ij} (f_i - f_j)^2 \xrightarrow{N \rightarrow \infty} \mathcal{L}_p(F)$$

[http://www.informatik.uni-hamburg.de/ML/contents/people/luxburg/publications/Luxburg04\\_diss.pdf](http://www.informatik.uni-hamburg.de/ML/contents/people/luxburg/publications/Luxburg04_diss.pdf)

<http://arxiv.org/pdf/1510.08110v1.pdf>

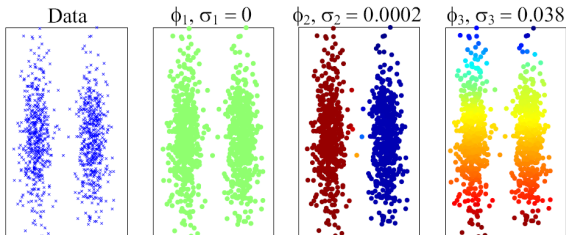
Isn't estimating eigenfunctions  $p(\mathbf{x})$  more difficult?

Are there some "easy" distributions?

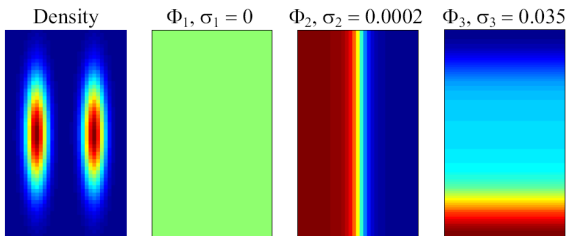
Can we compute it numerically?

# Scaling SSL with Graphs to Millions

## Eigenvectors



## Eigenfunctions

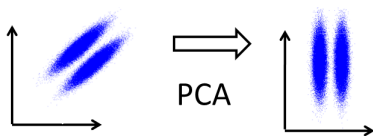


# Scaling SSL with Graphs to Millions

**Factorized data distribution** What if

$$p(\mathbf{s}) = p(s_1) p(s_2) \dots p(s_d)$$

In general, this is not true. But we can rotate data with  $\mathbf{s} = \mathbf{R}\mathbf{x}$ .



**Treating each factor individually**

$p_k \stackrel{\text{def}}{=} \text{marginal distribution of } s_k$

$\Phi_i(s_k) \stackrel{\text{def}}{=} \text{eigenfunction of } \mathcal{L}_{p_k} \text{ with eigenvalue } \lambda_i$

**Then:**  $\Phi_i(\mathbf{s}) = \Phi_i(s_k)$  is eigenfunction of  $\mathcal{L}_p$  with  $\lambda_i$

We only considered single-coordinate eigenfunctions.

# Scaling SSL with Graphs to Millions

How to approximate 1D density? Histograms!

Algorithm of Fergus et al. [FWT09] for eigenfunctions

- ▶ Find  $\mathbf{R}$  such that  $\mathbf{s} = \mathbf{R}\mathbf{x}$
- ▶ For each “independent”  $s_k$  approximate  $p(s_k)$
- ▶ Given  $p(s_k)$  numerically solve for eigensystem of  $\mathcal{L}_{p_k}$

$$\left(\tilde{\mathbf{D}} - \mathbf{P}\tilde{\mathbf{W}}\mathbf{P}\right)\mathbf{g} = \lambda\mathbf{P}\hat{\mathbf{D}}\mathbf{g} \quad (\text{generalized eigensystem})$$

$\mathbf{g}$  - vector of length  $B \equiv$  number of bins

$\mathbf{P}$  - density at discrete points

$\tilde{\mathbf{D}}$  - diagonal sum of  $\mathbf{P}\tilde{\mathbf{W}}\mathbf{P}$

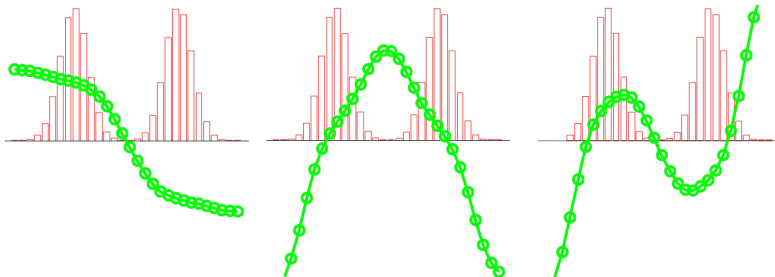
$\hat{\mathbf{D}}$  - diagonal sum of  $\mathbf{P}\tilde{\mathbf{W}}$

- ▶ Order eigenfunctions by increasing eigenvalues

[https://cs.nyu.edu/~fergus/papers/fwt\\_ssl.pdf](https://cs.nyu.edu/~fergus/papers/fwt_ssl.pdf)

# Scaling SSL with Graphs to Millions

## Numerical 1D Eigenfunctions



1<sup>st</sup> Eigenfunction  
of  $h(x_1)$

2<sup>nd</sup> Eigenfunction  
of  $h(x_1)$

3<sup>rd</sup> Eigenfunction  
of  $h(x_1)$

[https://cs.nyu.edu/~fergus/papers/fwt\\_ssl.pdf](https://cs.nyu.edu/~fergus/papers/fwt_ssl.pdf)

# Scaling SSL with Graphs to Millions

Computational complexity for  $N \times d$  dataset

## Typical harmonic approach

one diagonalization of  $N \times N$  system

**Numerical eigenfunctions** with  $B$  bins and  $k$  eigenvectors

$d$  eigenvector problems of  $B \times B$

$$\left(\tilde{\mathbf{D}} - \mathbf{P}\tilde{\mathbf{W}}\mathbf{P}\right)\mathbf{g} = \lambda\mathbf{P}\hat{\mathbf{D}}\mathbf{g}$$

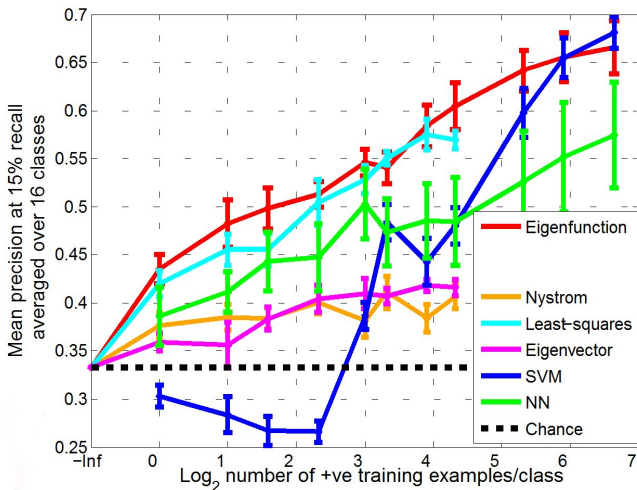
one  $k \times k$  least squares problem

$$\left(\mathbf{\Lambda} + \mathbf{U}^T\mathbf{C}\mathbf{U}\right)\alpha = \mathbf{U}^T\mathbf{C}\mathbf{y}$$

some details: several approximation, eigenvectors only linear combinations single-coordinate eigenvectors, ...

When  $d$  is not too big then  $N$  can be in millions!

# Scaling SSL with Graphs to Millions



CIFAR experiments [https://cs.nyu.edu/~fergus/papers/fwt\\_ssl.pdf](https://cs.nyu.edu/~fergus/papers/fwt_ssl.pdf)

*Michal Valko*

michal.valko@inria.fr

ENS Paris-Saclay, MVA 2016/2017

SequeL team, Inria Lille — Nord Europe

<https://team.inria.fr/sequel/>