

Optimization

Relaxation

RatioCut

$$\min_{A,B} RCut(A, B) = \min_{A,B} \sum_{i \in A, j \in B} w_{ij} \left(\frac{1}{|A|} + \frac{1}{|B|} \right)$$

Define graph function f for cluster membership of RatioCut: $f_i = \begin{cases} \sqrt{\frac{|B|}{|A|}} & \text{if } v_i \in A \\ -\sqrt{\frac{|A|}{|B|}} & \text{if } v_i \in B \end{cases}$

$$f^T L f = \frac{1}{2} \sum_{i,j} w_{ij} (f_i - f_j)^2 = (|A| + |B|) RCut(A, B)$$

Since $(|A| + |B|)$ is constant,

$$\min_{A,B} RCut(A, B) = \min_f f^T L f,$$

$$\text{subject to } f_i \in \left\{ \sqrt{\frac{|B|}{|A|}}, -\sqrt{\frac{|A|}{|B|}} \right\}$$

$$\rightarrow |A| = |B|$$

$$\|f\|^2 = \sum_i f_i^2 = |A| \frac{|B|}{|A|} + |B| \frac{|A|}{|B|} = |A| + |B| = N$$

$$|A| = |B| \rightarrow \sum_i f_i = 0 \leftrightarrow f \perp \mathbf{1}_N$$

Still NP hard...Require relaxation.

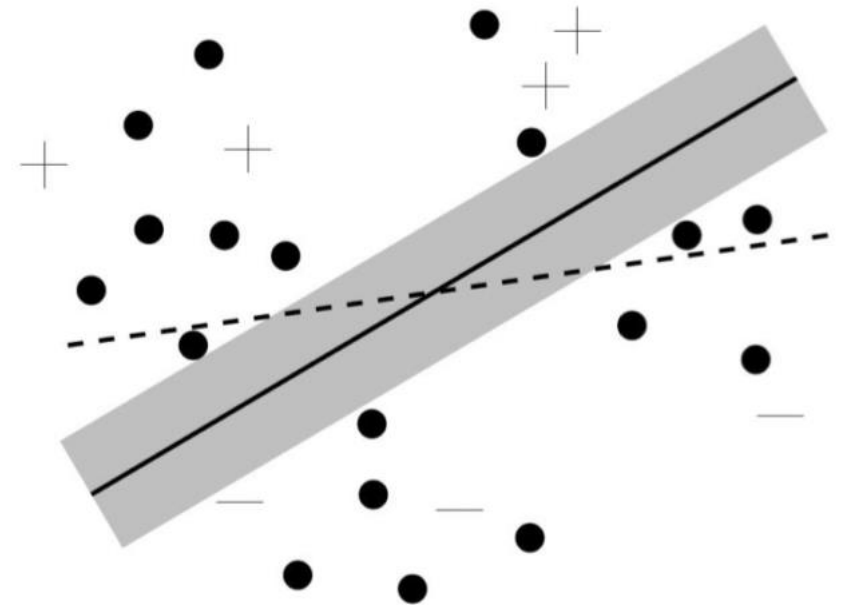
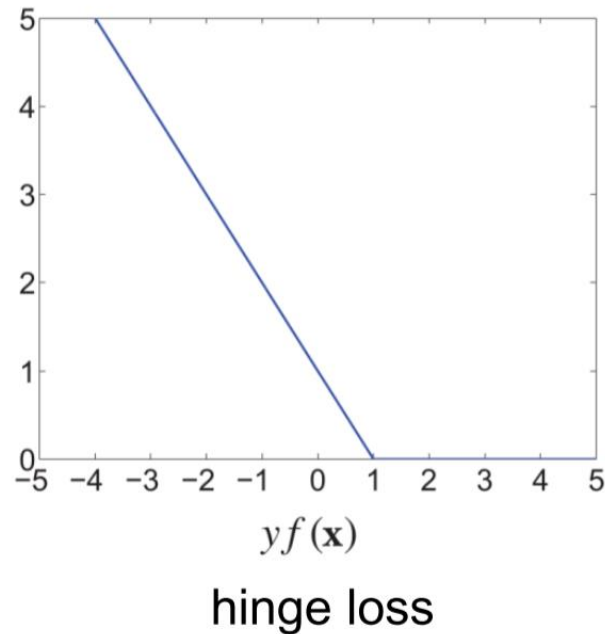
Optimization formulation for RatioCut (same with balanced Mincut)

$$\min_f f^T L f \text{ subject to } f_i \in R, f \perp \mathbf{1}_N, \|f\| = \sqrt{N}$$

$$f^T L f \neq (|A| + |B|) \sum_{i \in A, j \in B} w_{ij} \left(\frac{1}{|A|} + \frac{1}{|B|} \right) \rightarrow |A| = |B|$$

Questions of the last lecture

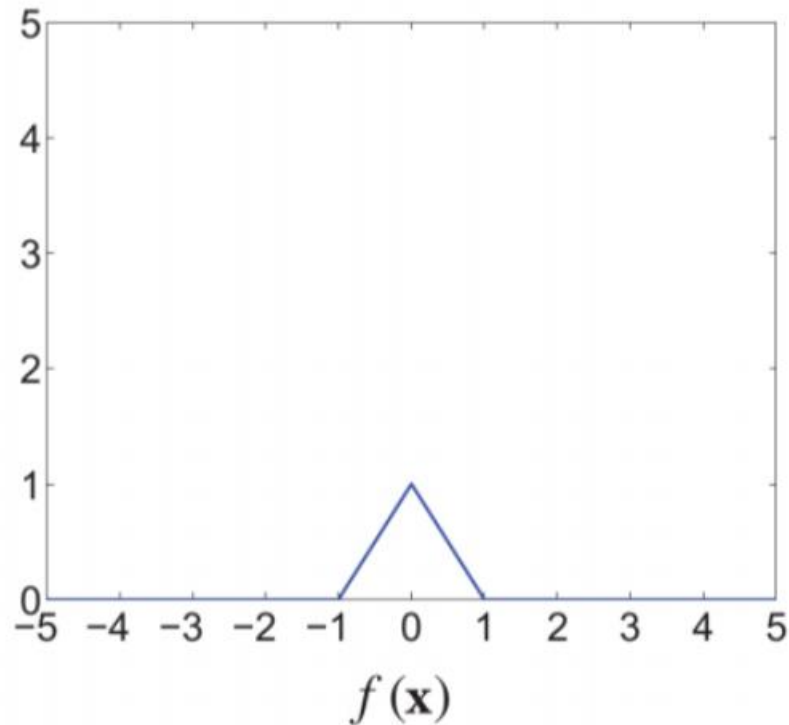
- What does **hinge loss** in **SVM** penalize?
→ Hinge loss penalizes the case that the classifier $f(w, x)$ does **not decide the correct class** of **labeled** x with the score **margin of** $y f(w, x) \geq 1$.



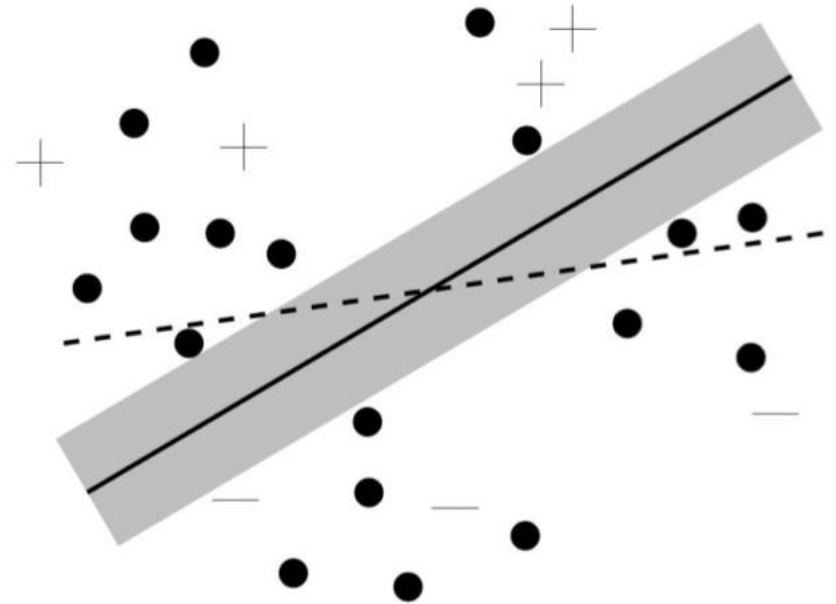
$$\Phi(x_i, y_i, f(\mathbf{w}, b; \mathbf{x}_i)) = \max(1 - y_i(\mathbf{w}^T \mathbf{x}_i + b), 0)$$

Questions of the last lecture

- What does **hat loss** in SVM for semi-supervised learning penalize?
→ Hat loss penalizes the case that the classifier $f(w, x)$ does **not decide any class** of **unlabeled** x with the score **margin of $|f(w, x)| \geq 1$** .



(b) the hat loss



Questions of the last lecture

- Why can't we use eigenvectors to solve MinCut-based SSL in graph?
→ It is because the eigenvector solution **can not consider the labeled data.**

Transductive SSL with graph: fixing $f(\mathbf{x}_i)$ for $i \in \mathcal{L}$

$$\min_{f \in \{\pm 1\}^{n_l+n_u}} \lambda \sum_{i,j=1}^{n_l+n_u} w_{ij} (f(\mathbf{x}_i) - f(\mathbf{x}_j))^2 + \infty \sum_i^{n_l} (f(\mathbf{x}_i) - y_i)^2$$

Solution:

An integer program: NP hard

Can we use eigenvectors? No. Why?

We need a **better way** to reflect the confidence.

Questions of the last lecture

- What is the meaning of harmonic function in SSL?
→ The harmonic function makes the **label of a node to be harmonious(similar) with those of its neighboring nodes.**

Relaxation: Transductive SSL with graph: fixing $f(x_i)$ for $i \in \mathcal{L}$

$$\min_{f \in \mathbf{R}^{n_l+n_u}} \lambda \sum_{i,j=1}^{n_l+n_u} w_{ij} (f(x_i) - f(x_j))^2 + \infty \sum_i^{n_l} (f(x_i) - y_i)^2$$

Naïve Solution

Right term solution: constrain f to **match** the supervised data

$$f(x_i) = y_i \quad \forall i \in \{1, \dots, n_l\}$$

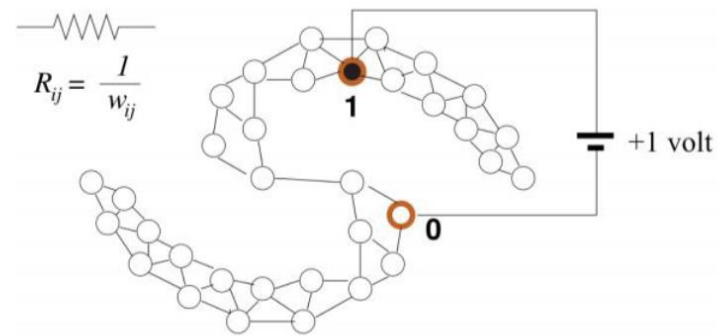
Left term solution: enforce the solution f to be **harmonic** (cf. aggregation, rw)

$$f(x_i) = \frac{\sum_{ij} f(x_j) w_{ij}}{\sum_{ij} w_{ij}} \quad \forall i \in \{n_l + 1, \dots, n_l + n_u\}$$

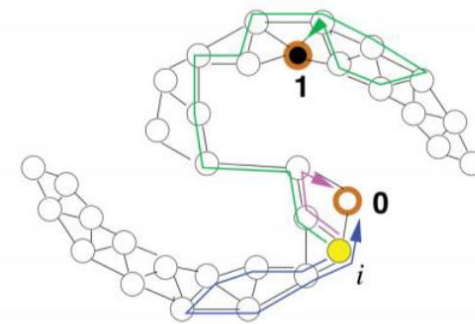
$$\mathbf{f}_u = L_{uu}^{-1}(-L_{ul}\mathbf{f}_l) = L_{uu}^{-1}(W_{ul}\mathbf{f}_l)$$

Questions of the last lecture

- What is random walk interpretation of harmonic function-based SSL in graph?
→ The label of a node is assigned by the **average of the harmonic labels** of the vertices that are **hit by random works**.



(a) The electric network interpretation



(b) The random walk interpretation

Random walk interpretation :

1) start from the vertex you want to label and randomly walk

2) $P(j|i) = \frac{w_{ij}}{\sum_k w_{ik}} \Leftrightarrow \mathbf{P} = \mathbf{D}^{-1}\mathbf{W}$

3) finish when a labeled vertex is hit

$$f(\mathbf{x}_i) = \frac{\sum_{ij} f(\mathbf{x}_j)w_{ij}}{\sum_{ij} w_{ij}}$$

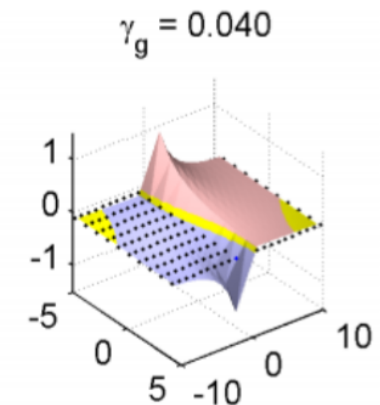
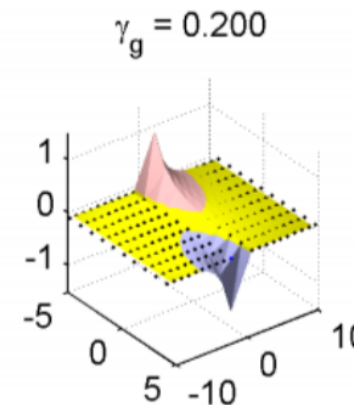
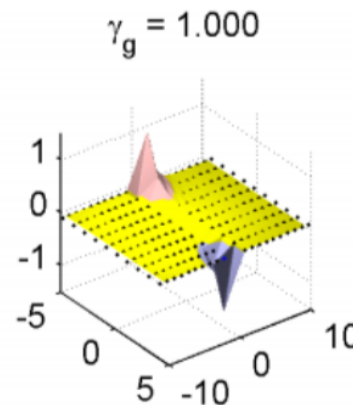
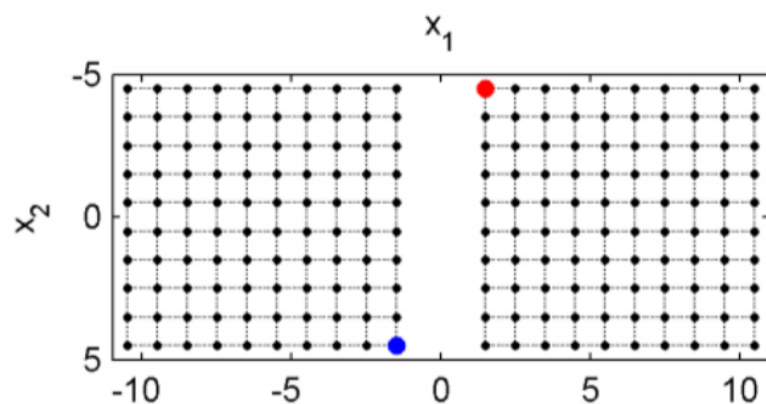
4) $f(\mathbf{x}_i)$ is assigned by **average** of the labels of the hit vertices.

Questions of the last lecture

- What is a key point of regularized harmonic function-based SSL in graph?
→ A **sink** node is added to allow the **random work to die at any nodes**, which **reduces the misleading by outliers**.

$$f_u = (L_{uu} + \gamma_g \mathbf{I})^{-1} (W_{ul} f_l),$$

How does γ_g influence the solution?



Questions of the last lecture

- What is a key point of soft harmonic function-based SSL in graph?
→ The labeled data is not constrained strictly, where noisy labels may be smoothed by the soft harmonic function

Regularized HS objective with $Q = L + \gamma_g \mathbf{I}$:

Define $f_i \triangleq f(\mathbf{x}_i)$, $\mathbf{f} \triangleq [f_1, \dots, f_{n_l+n_u}]$

$$\min_{\mathbf{f} \in \mathbb{R}^{n_l+n_u}} \infty \sum_{i=1}^{n_l} (f_i - y_i)^2 + \lambda \mathbf{f}^T \mathbf{Q} \mathbf{f}$$

Soft constraints for $f(\mathbf{x}_i) = y_i$, $\forall i \in \mathcal{L}$: ∞ is replaced by finite values

$$\min_{\mathbf{f} \in \mathbb{R}^{n_l+n_u}} (\mathbf{f} - \mathbf{y})^T \mathbf{C} (\mathbf{f} - \mathbf{y}) + \mathbf{f}^T \mathbf{Q} \mathbf{f}$$

Outline of Lecture (1)

- Graph Spectral Theory
 - Definition of Graph
 - Graph Laplacian
 - Laplacian Smoothing
- Graph Node Clustering
 - Minimum Graph Cut
 - Ratio Graph Cut
 - Normalized Graph Cut
- Manifold Learning
 - Spectral Analysis in Riemannian Manifolds
 - Dimension Reduction, Node Embedding
- Semi-supervised Learning (SSL)
 - Self-Training Methods
 - SSL with SVM
 - SSL with Graph using MinCut
 - SSL with Graph using Harmonic Functions
- Semi-supervised Learning (SSL) : conti.
 - SSL with Graph using Regularized Harmonic Functions
 - SSL with Graph using Soft Harmonic Functions
 - SSL with Graph using **Manifold Regularization (out of sample extension)**
 - SSL with Graph using **Laplacian SVMs**
 - SSL with Graph using **Max-Margin Graph Cuts**
 - **Online SSL**
 - SSL for large graph
- Graph Convolution Networks (GCN)
 - Graph Filtering in GCN
 - Graph Pooling in GCN
 - Spectral Filtering in GCN
 - Spatial Filtering in GCN
- Recent GCN papers

SSL with Graphs: Out of sample extension

Both MinCut and HF only inferred the labels on unlabeled data.
They are **transductive**.

What if a new point $x_{n_l+n_u+1}$ arrives? (called out of sample extension)

Option 1) Add it to the graph and recompute HF Solution. (Still **Transductive**)
Option 2) Make the algorithms **inductive**!

Define a classifier; $f : \mathcal{X} \rightarrow \mathbb{R}$

Make $f(x_i)$ be smooth. Why? To deal with noise by providing reasonable interpolation for **new samples**.

Solution: **Manifold Regularization**

SSL with Graphs: Manifold Regularization

General (S)SL objective:

$$\min_{f \in \mathcal{H}} \sum_i^{n_l} \Phi(\mathbf{x}_i, y_i, f(\mathbf{x}_i)) + \lambda \Omega(\mathbf{f}), \mathbf{f} \triangleq [\dots f(\mathbf{x}_i) \dots]$$

Want to control f , also for the out-of-sample data, i.e., **everywhere**(generalization).

$$\lambda \Omega(\mathbf{f}) = \lambda_2 \mathbf{f}^T \mathbf{L} \mathbf{f} + \lambda_1 \|f\|_{\mathcal{K}}^2$$

$$\|f\|_{\mathcal{K}}^2 = \langle \nabla f, \nabla f \rangle_{L^2(T\mathcal{X})} = \int f(\mathbf{x}) \Delta f(\mathbf{x}) d\mathbf{x}$$

For general **kernels**:

$$\min_{f \in \mathcal{H}_{\mathcal{K}}} \sum_i^{n_l} \Phi(\mathbf{x}_i, y_i, f(\mathbf{x}_i)) + \lambda_1 \|f\|_{\mathcal{K}}^2 + \lambda_2 \mathbf{f}^T \mathbf{L} \mathbf{f}$$

Smoothness for given samples



Smoothness for unknown samples



SSL with Graphs: Manifold Regularization

General (S)SL objective with kernels:

$$\min_{f \in \mathcal{H}_{\mathcal{K}}} \sum_i^{n_l} \Phi(\mathbf{x}_i, y_i, f(\mathbf{x}_i)) + \lambda_1 \|f\|_{\mathcal{K}}^2 + \lambda_2 \mathbf{f}^T \mathbf{L} \mathbf{f}$$

Representer theorem for manifold regularization

The minimizer f^* has a finite expansion of the form

$$f^*(\mathbf{x}) = \sum_{i=1}^{n_l+n_u} \alpha_i^* \mathcal{K}(\mathbf{x}, \mathbf{x}_i)$$

LapRLS: Laplacian Regularized Least Squares

$$\Phi(\mathbf{x}, y, f(\mathbf{x})) = (y - f(\mathbf{x}))^2$$

LapSVM: Laplacian Support Vector Machines

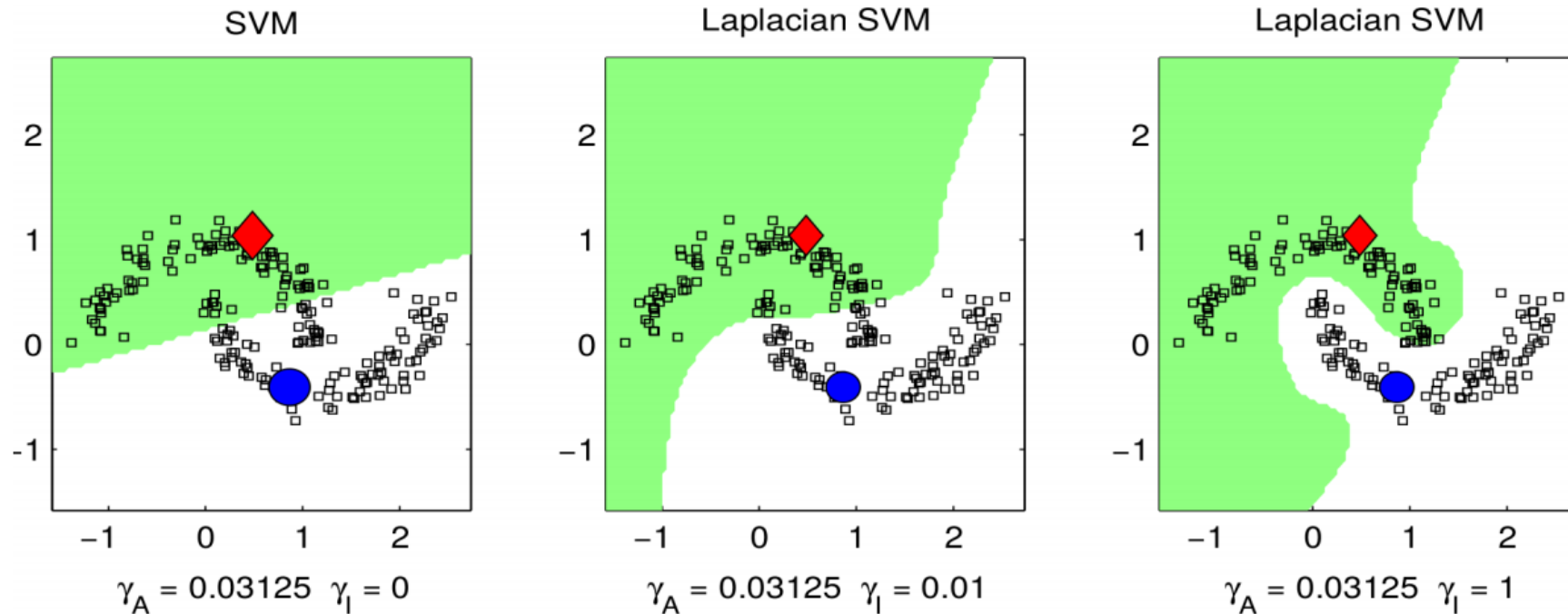
$$\Phi(\mathbf{x}, y, f(\mathbf{x})) = \max(0, 1 - yf(\mathbf{x}))$$

SSL with Graphs: Laplacian SVMs

General (S)SL objective with kernels:

$$\min_{f \in \mathcal{H}_{\mathcal{X}}} \sum_i^{n_l} \max(0, 1 - y_i f(x_i)) + \lambda_A \|f\|_{\mathcal{K}}^2 + \lambda_l f^T L f$$

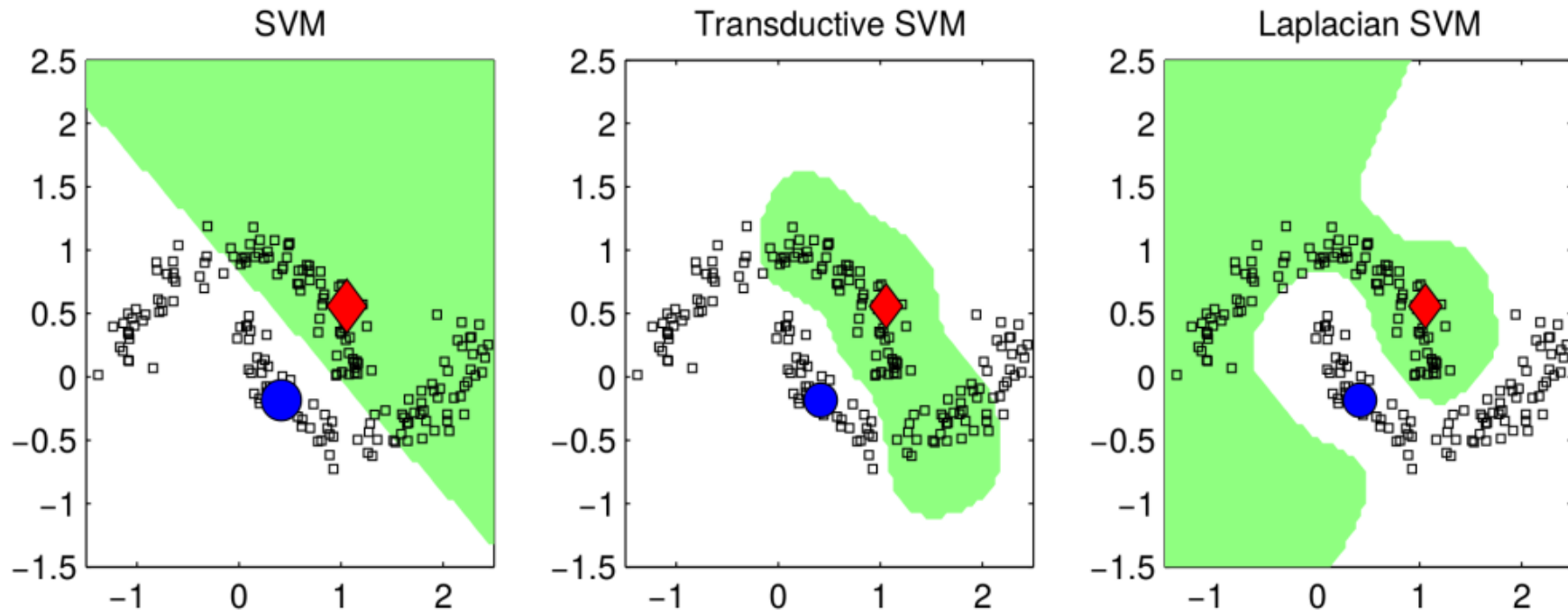
RBF kernels



SSL with Graphs: Laplacian SVMs

Formulation for Transductive SVM (Revisit)

$$\min_{w,b} \sum_{i=1}^{n_l} \max(1 - y_i(\mathbf{w}^T \mathbf{x}_i + b), 0) + \lambda_1 \|\mathbf{w}\|^2 \\ + \lambda_2 \sum_{i=n_l+1}^{n_l+n_u} \max(1 - |\mathbf{w}^T \mathbf{x}_i + b|, 0)$$



SSL with Graphs: Max-Margin Graph Cuts

Self-training with the confident data

$$\begin{aligned} f^* &= \underset{f \in \mathcal{H}_{\mathcal{X}}}{\operatorname{argmin}} \sum_{i: |\ell_i^*| \geq \varepsilon} \Phi(\mathbf{x}_i, \operatorname{sgn}(\ell_i^*), f(\mathbf{x}_i)) + \lambda \|f\|_{\mathcal{K}}^2 \\ \text{s. t. } \quad \ell^* &= \underset{\ell \in \mathbb{R}^N}{\operatorname{argmin}} \ell^T (\mathbf{L} + \gamma_g \mathbf{I}) \ell \\ \text{s. t. } \quad \ell_i &= y_i, \quad \forall i = 1, \dots, n_l \end{aligned}$$

Representer theorem is still cool:

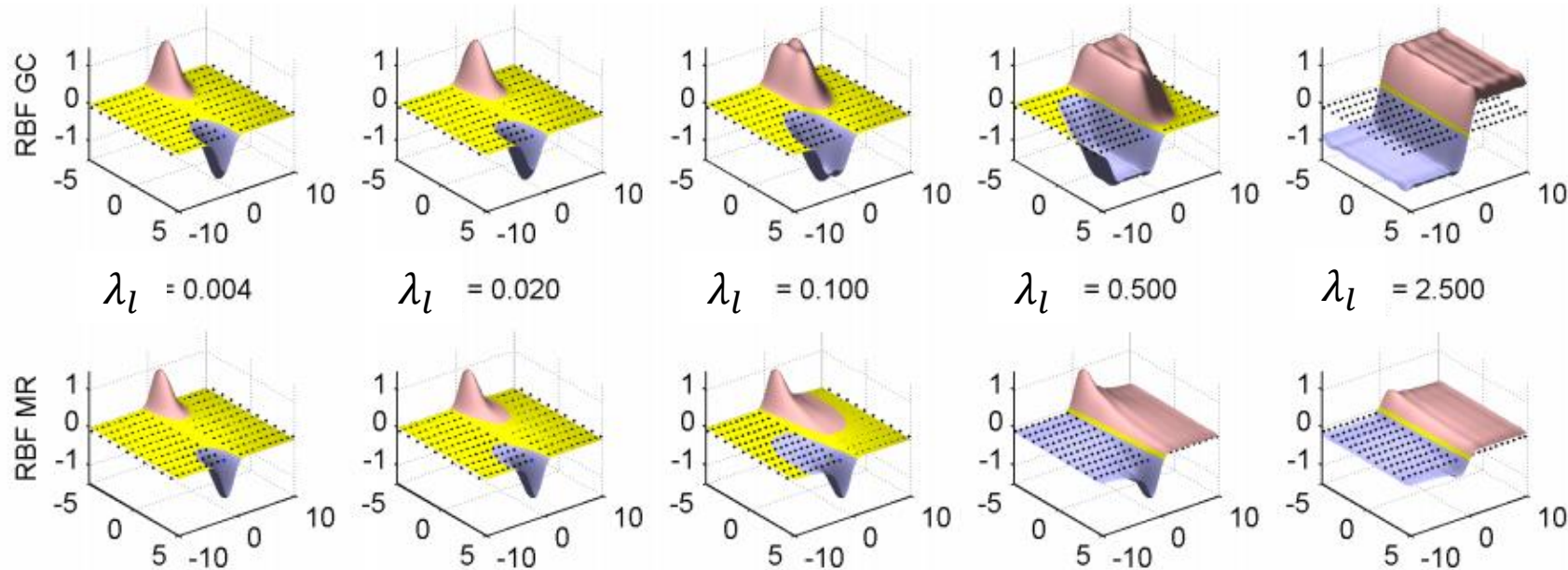
$$f^*(\mathbf{x}) = \sum_{i: |\ell_i^*| \geq \varepsilon} \alpha_i^* \mathcal{K}(\mathbf{x}, \mathbf{x}_i)$$

SSL with Graphs: LapSVMs and MM Graph Cuts

$$\min_{f \in \mathcal{H}_{\mathcal{K}}} \sum_i^{n_l} \max(0, 1 - y_i f(x_i)) + \lambda_A \|f\|_{\mathcal{K}}^2 + \lambda_l f^T L f$$

$$\lambda_A = 0.1, \quad \epsilon = 0.01$$

MMGC and MR for 2D data and RBF \mathcal{K}



Manifold regularization of SVMs (MR), max-margin graph cuts (GC)

OnlineSSL(G)

when we can't access future x



Online SSL with Graphs

Offline learning setup

Given $\{\mathbf{x}_i\}_{i=1}^N$ from \mathbb{R}^d and $\{y_i\}_{i=1}^n$, with $n \ll N$, find $\{y_i\}_{i=n+1}^N$ (**transductive**) or find f predicting $\{y_i | y_i = f(\mathbf{x}_i), i = n + 1, \dots, N\}$ well beyond that (**inductive**).

Online learning setup

At the beginning, given $\{\mathbf{x}_i, y_i\}_{i=1}^n$ from \mathbb{R}^d .

At time t :

receive \mathbf{x}_t

predict y_t

Revisit : out of sample expansion

Option 1) Add it to the graph and recompute HF Solution.

Option 2) Make the algorithms inductive! (**not learn \mathbf{x}_t**)

Online SSL with Graphs

Online HFS: Straightforward solution (option 1)

- 1: **while** new unlabeled example x_t comes **do**
- 2: Add x_t to the graph $G(W)$
- 3: Update L_t
- 4: Infer labels

$$\mathbf{f}_u = (\mathbf{L}_{uu} + \gamma_g \mathbf{I})^{-1} (\mathbf{W}_{ul} \mathbf{f}_l)$$

- 5: Predict $\hat{y}_t = \text{sgn}(\mathbf{f}_{u,t})$
 - 6: **end while**
-

What is wrong with this solution?

The cost and memory of the operations.

What can we do?

Online SSL with Graphs: Graph Quantization

Let's keep **only k vertices!**

Limit memory to k **centroids** with $\widetilde{\mathbf{W}}^q$, where $\widetilde{\mathbf{W}}_{ij}^q$ contains the **similarity** between the i -th and j -th centroids.

Each centroid represents several others.

Let \mathbf{V} be a diagonal matrix of which

V_{ii} denotes **number of points** collapsed into the i -th centroid.

Can we compute it compactly? Compact harmonic solution.

$$\mathbf{f}_u^q = (\mathbf{L}_{uu}^q + \gamma_g \mathbf{V})^{-1} (\mathbf{W}_{ul}^q \mathbf{f}_l) \text{ where } \mathbf{W}^q = \mathbf{V} \widetilde{\mathbf{W}}^q \mathbf{V}$$

Proof and Algorithm: see <http://www.bkveton.com/docs/uai2010a.pdf>

Online SSL with Graphs

Online HFS with Graph Quantization

01: **Input**

02: k number of representative nodes

03: **Initialization**

04: V matrix of multiplicities with 1 on diagonal

05: **while** new unlabeled example x_t comes **do**

06: Add x_t to graph G

07: **if** # nodes $> k$ **then**

08: **quantize** G

09: **end if**

10: Update L_t of $G(VWV)$

11: Infer labels

12: Predict $\hat{y}_t = \text{sgn}(f_u(t))$

13: **end while**

Online SSL with Graphs: Graph Quantization

Incremental **quantize** G

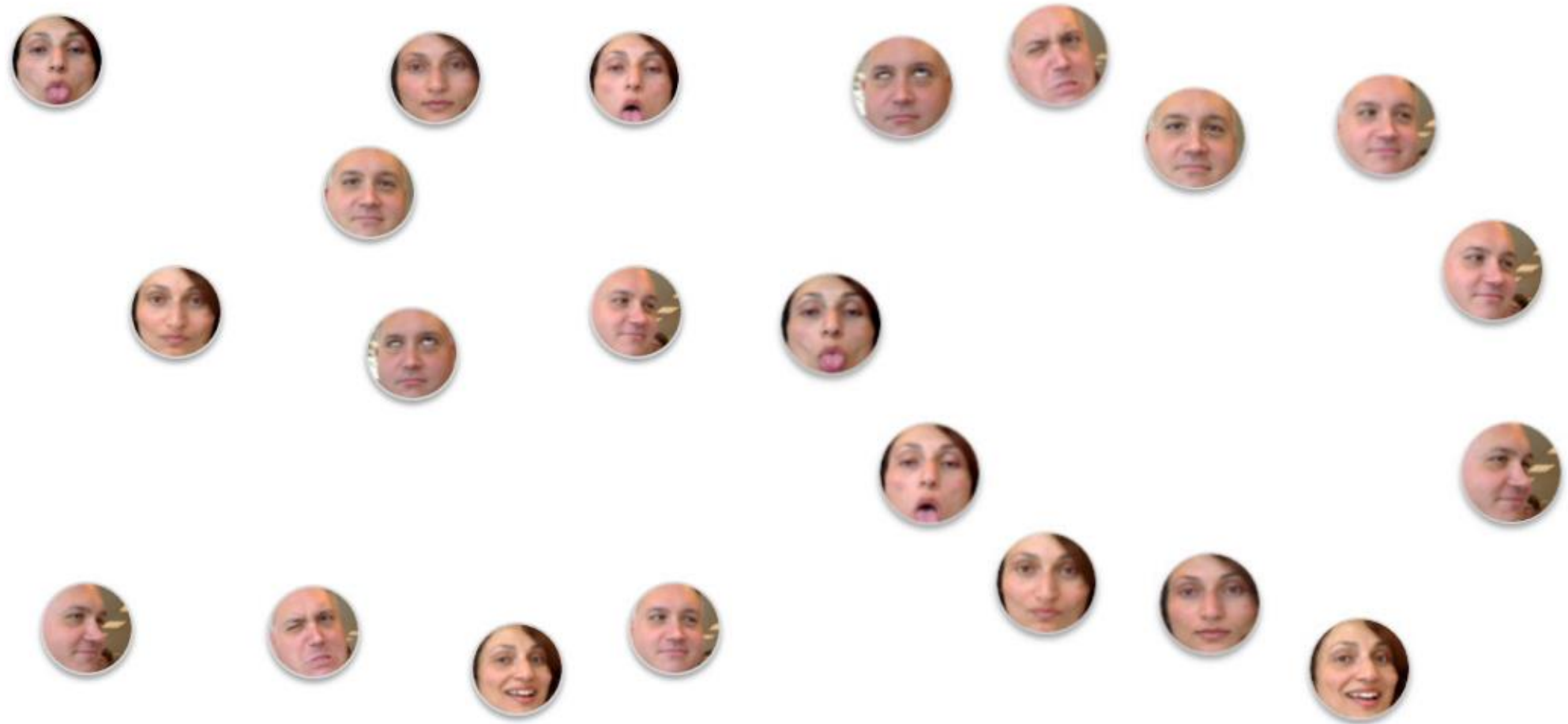
An idea: incremental k -centers

Doubling algorithm of Charikar et al. [Cha+97]

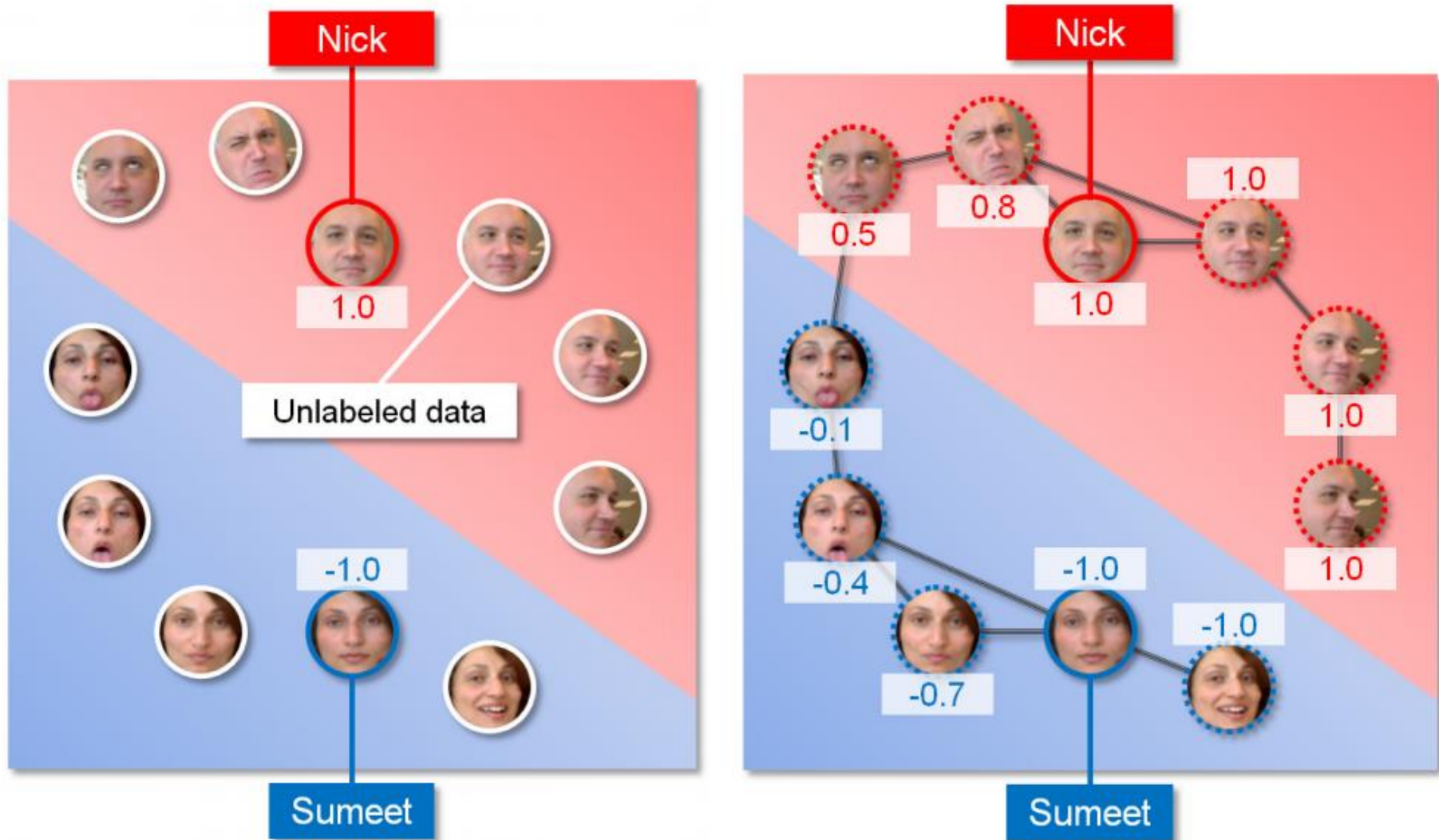
Keeps up to k centers $C_t = \{c_1, c_2, \dots\}$ with

- Distance $c_i, c_j \in C_t$ is at least $\geq R$
- For each new x_t , distance to some $c_i \in C_t$ is less than R .
- $|C_t| \leq k$
- if not possible, R is doubled

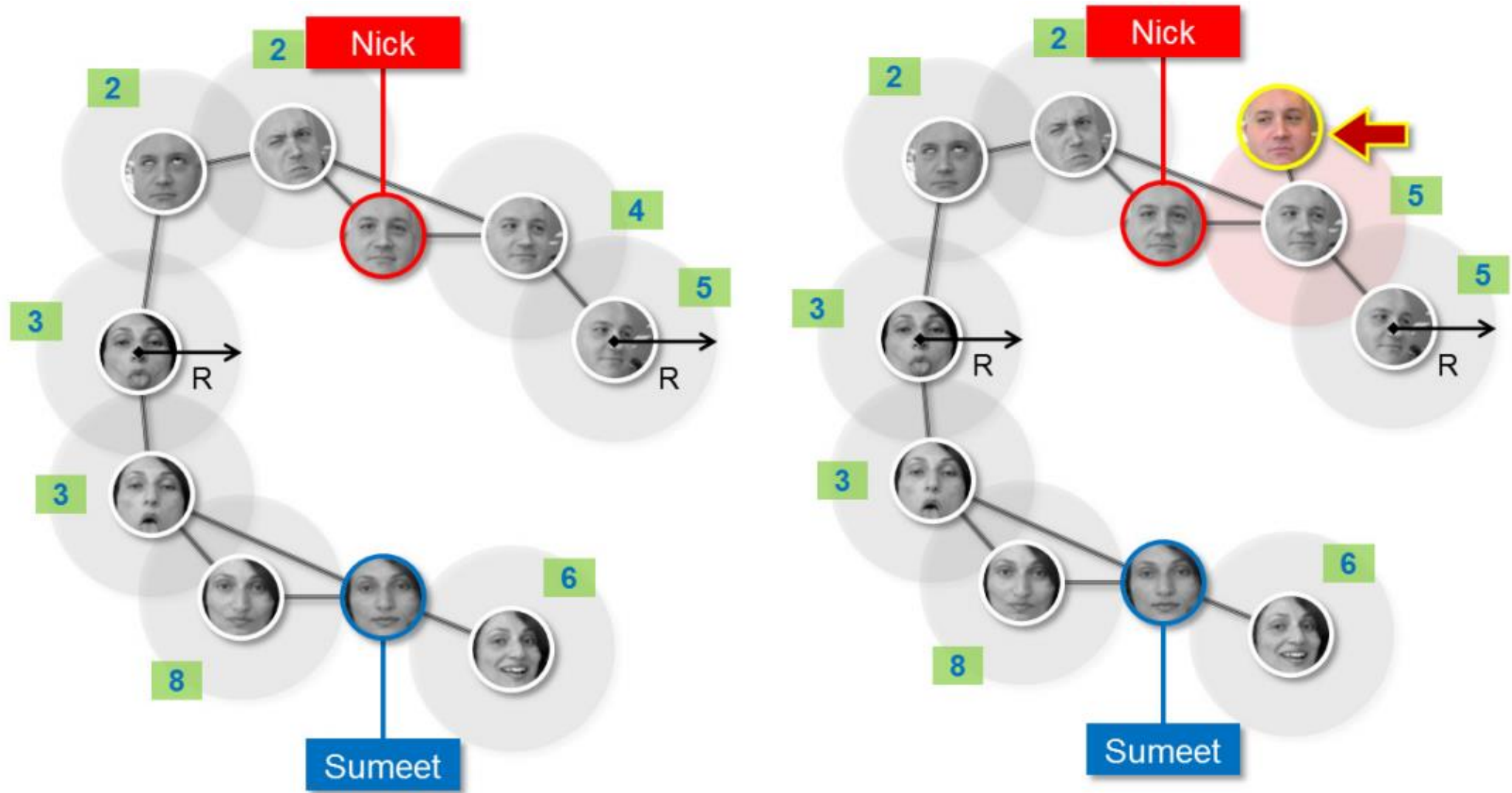
Online SSL with Graphs: Graph Quantization



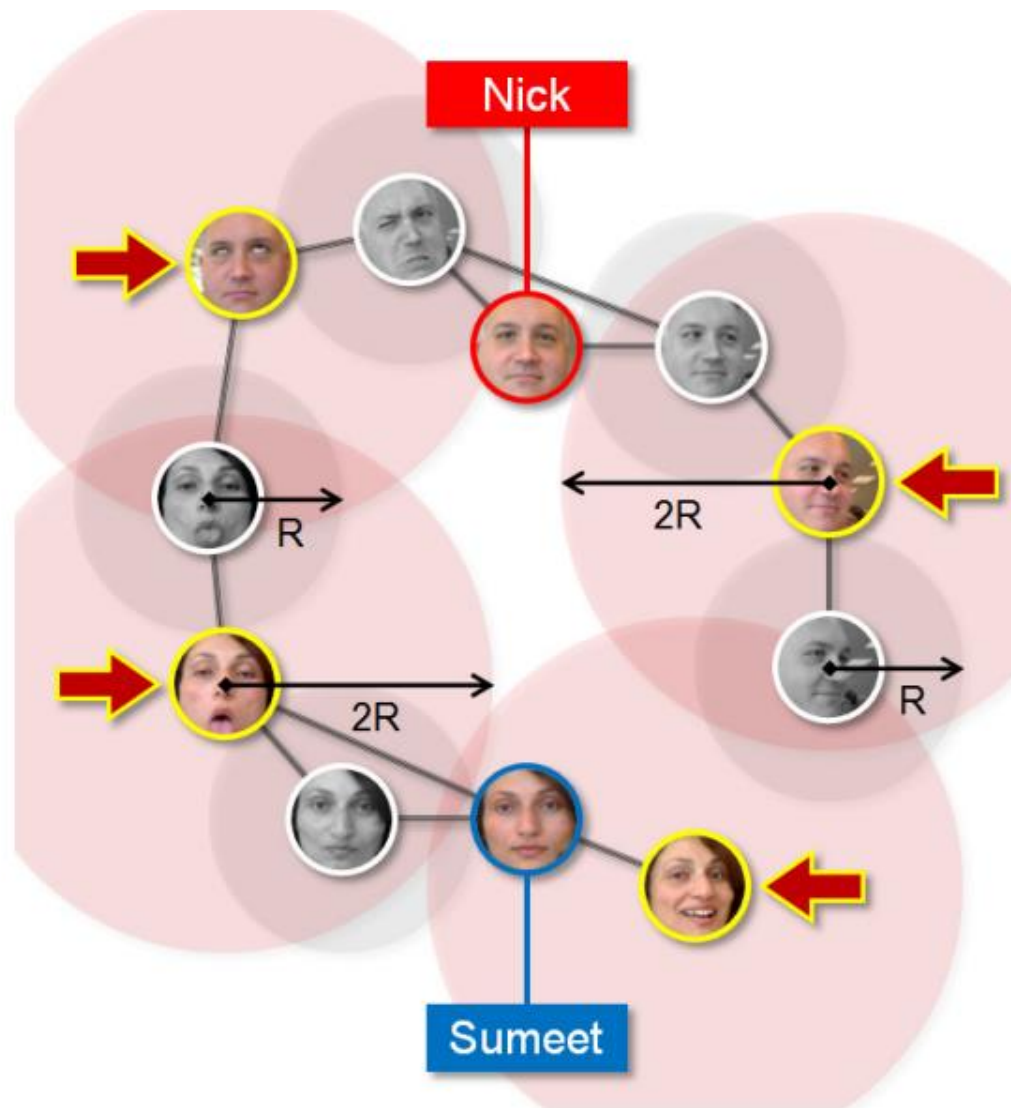
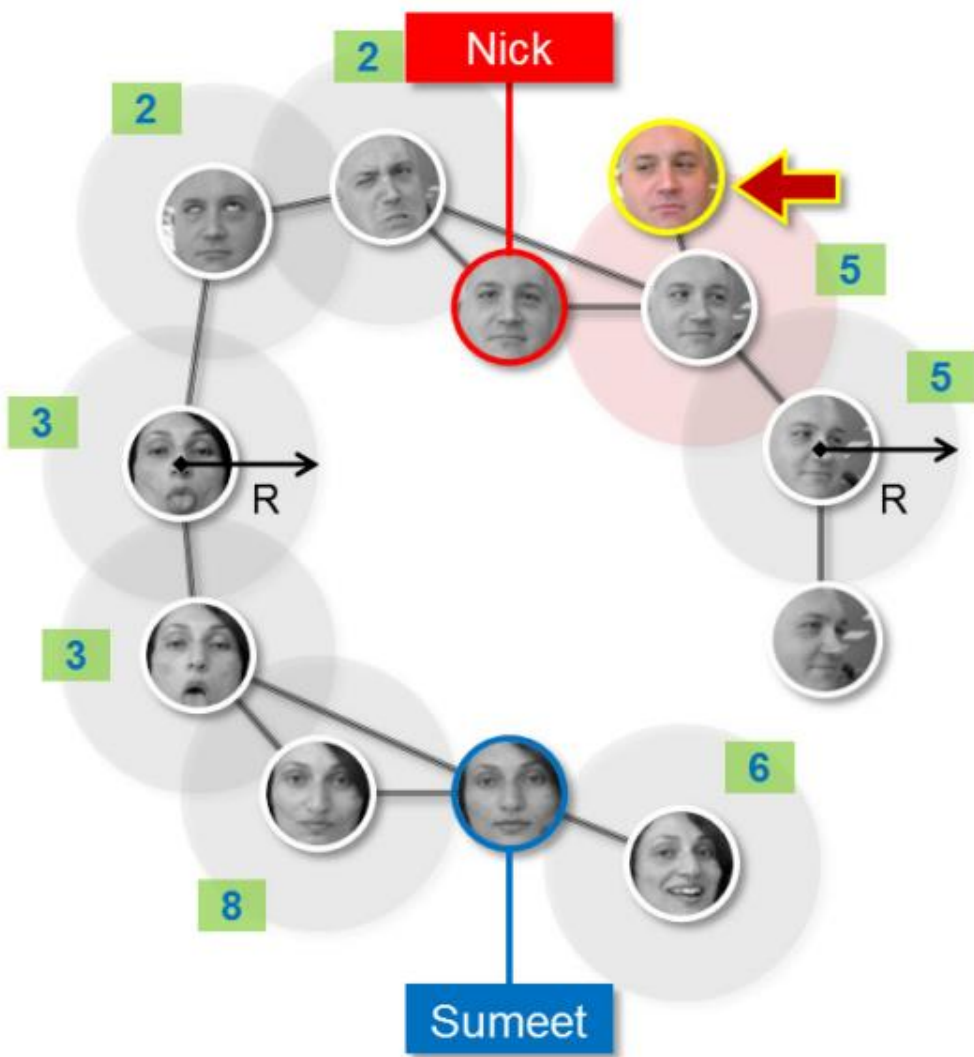
Online SSL with Graphs: Graph Quantization



Online SSL with Graphs: Graph Quantization

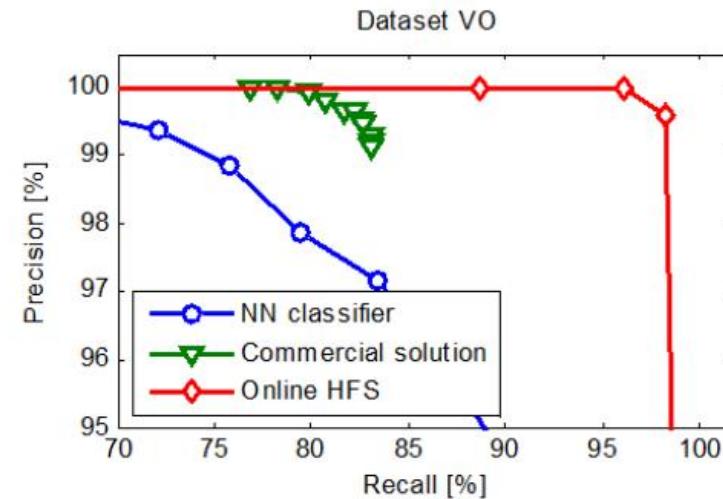
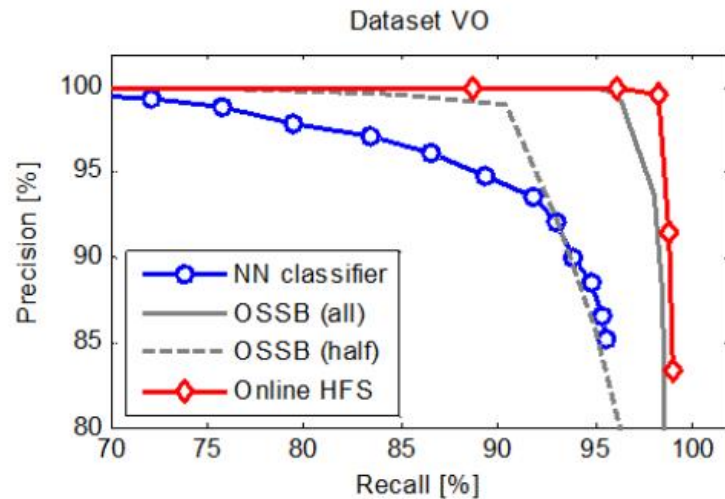


Online SSL with Graphs: Graph Quantization



Online SSL with Graphs: Some experimental results

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



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

Summary Questions of the Lecture

What are the two options for out of sample extension in SSL?

Why do we have to make a classifier be smooth in inductive SSL for out of sample extension?

What is the meaning of manifold regularization?

What is the key idea of Max-Margin Graph Cuts for SSL?

What are the two options for online SSL?

What is the key idea for keeping # of representative nodes?