

Semi-supervised Learning with Gaussian Random Fields

Jerry Xiaojin Zhu

Joint work with:

Zoubin Ghahramani
John Lafferty

February 16, 2004

Traditional Types of Learning

Supervised (classification, regression, etc.)

Unsupervised (clustering etc.)

	yes	no
unlabeled data $\{x\}$	no	yes
labeled data $\{(x, y)\}$	yes	no
usage	supervised learning	unsupervised learning

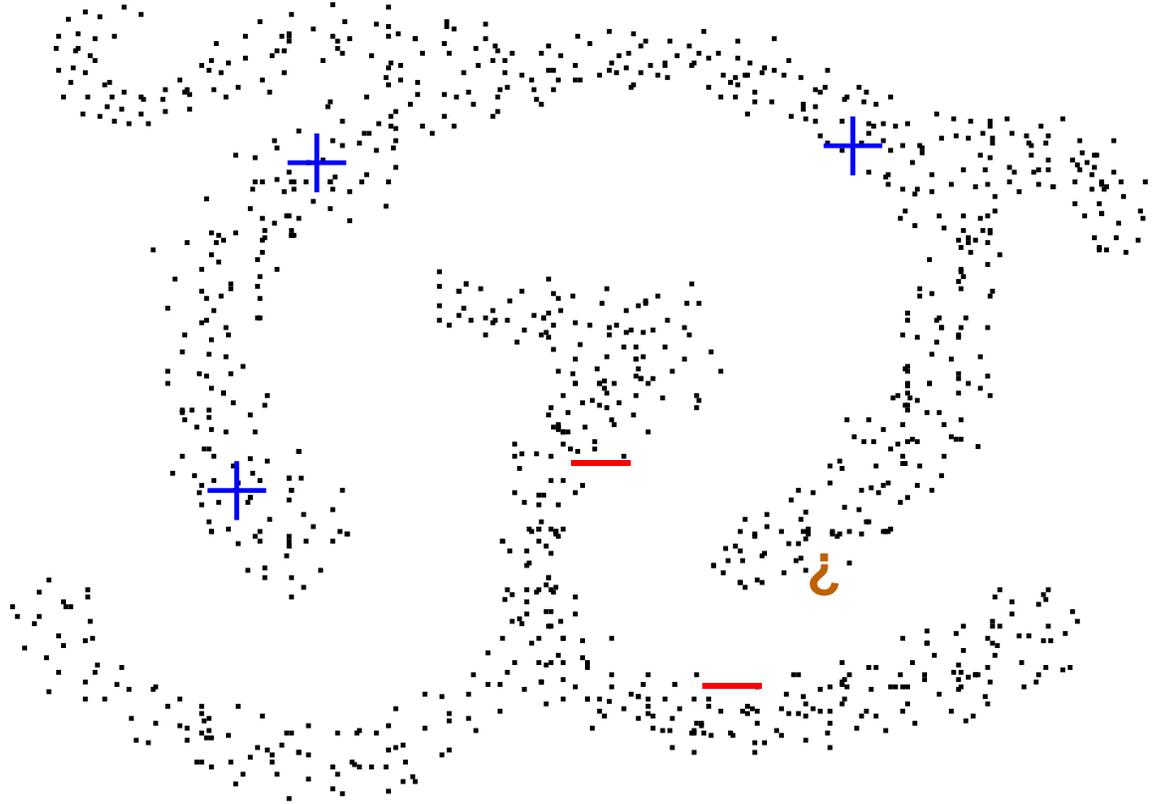
For Classification

- **Labeled** data are often **hard, expensive, slow** to obtain.
- **Unlabeled** data are often **easy** to obtain, a lot.
- Speech recognition: **slow** to transcribe
- Text categorization: **limited** user patience
- Protein structure: **months** of X-ray crystallography

Semi-supervised Learning

usage	labeled data $\{(x, y)\}$	unlabeled data $\{x\}$
supervised learning	yes	no
semi-supervised learning	yes	yes
unsupervised learning	no	yes

How is it possible? Intuition

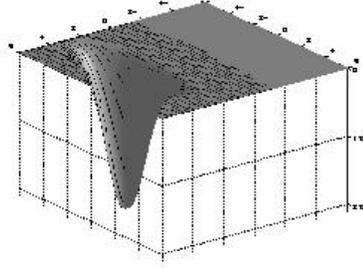
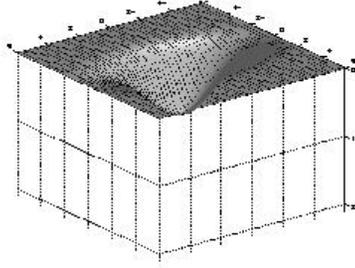
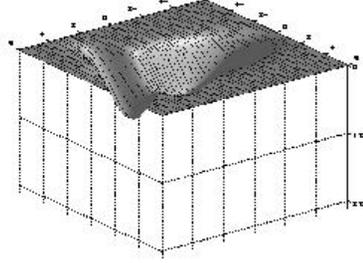
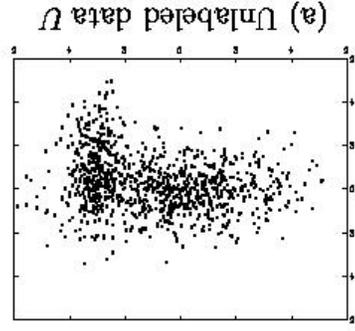


Some Semi-supervised Learning Methods

- mixture models / EM
- transductive SVM
- co-training
- graph methods → this talk

Mixture Models. Cluster then label

Method
mixture
trans. SVM
co-training
graph



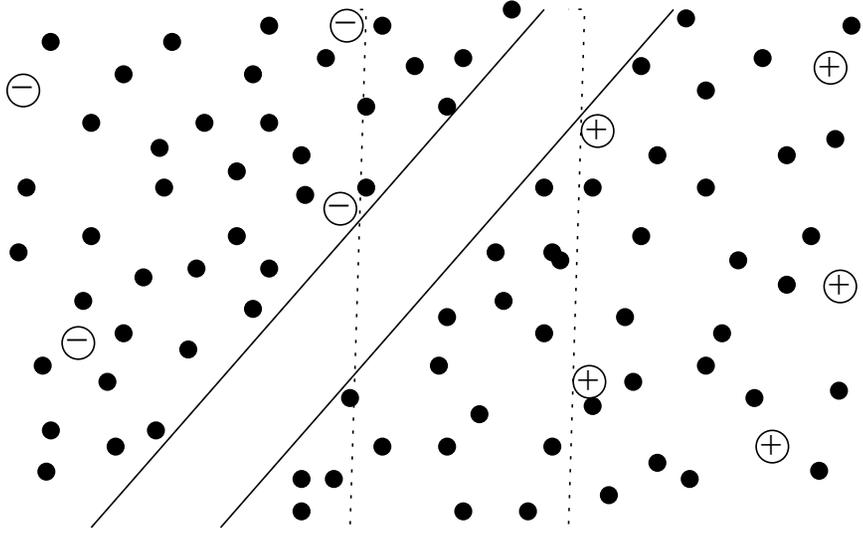
[Castelli & Cover 95], [Ratsaby & Venkatesh 95], [Nigam et al. 00]

Mixture Models. Cluster then label

Method
mixture
trans. SVM
co-training
graph

- Unlabeled data may **hurt**, when...
 - mixture assumption wrong
 - EM stuck in bad local minima
- Cluster first, label next
 - same flavor
 - harder to analyze

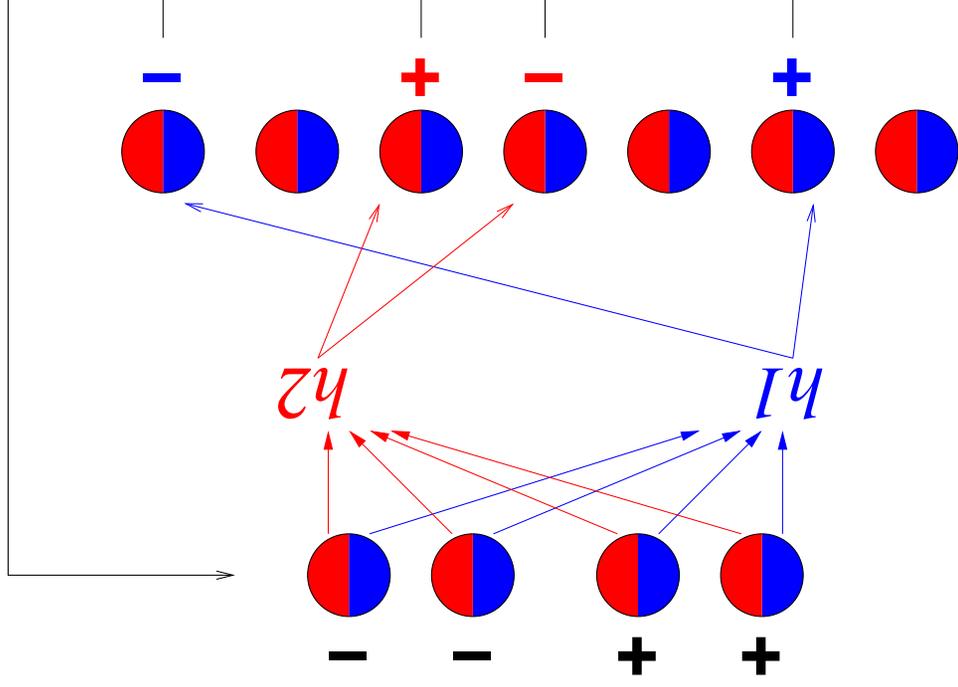
Transductive SVM



Unlabeled data steer linear boundary away from dense regions.
[Vapnik 98], [Joachims 99]

Method
mixture
trans. SVM
co-training
graph

Co-training



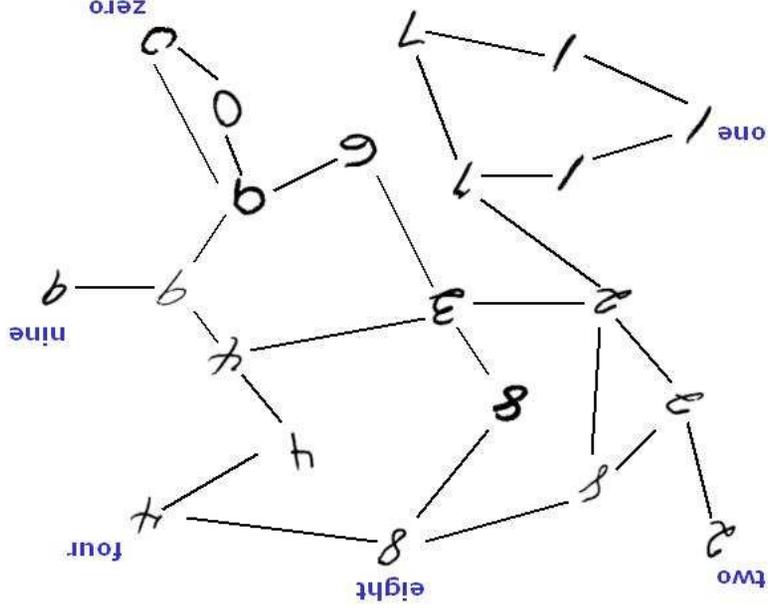
Method
mixture
trans. SVM
co-training
graph

Unlabeled data reduce version space size by

forcing h_1, h_2 to agree. [Blum & Mitchell 98]

Graph Methods

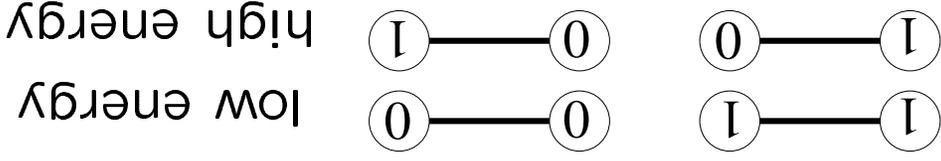
Nodes: labeled and unlabeled data
Edges: "local similarity"
Labels: 'propagate'



Method
mixture
trans. SVM
co-training
graph

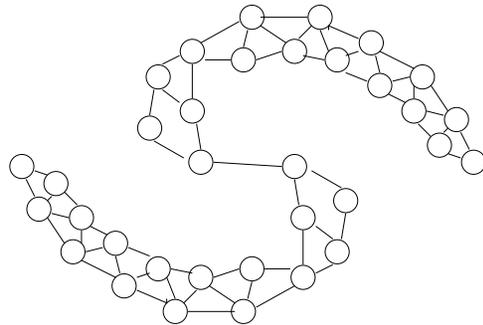
Energy

- **labels:** Assume binary $y \in \{0, 1\}^n$
- **edges:** weight matrix W . (Very, very important)
- **energy:** $E(y) = \frac{1}{2} \sum_{i,j} w_{ij} (y_i - y_j)^2$

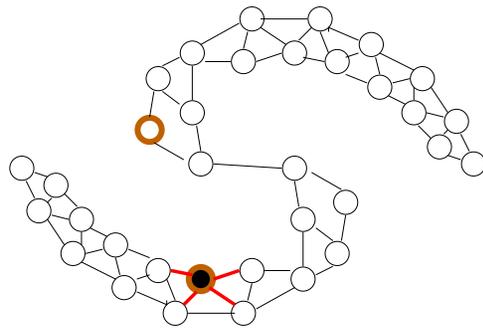


Low energy \rightarrow Label Propagation

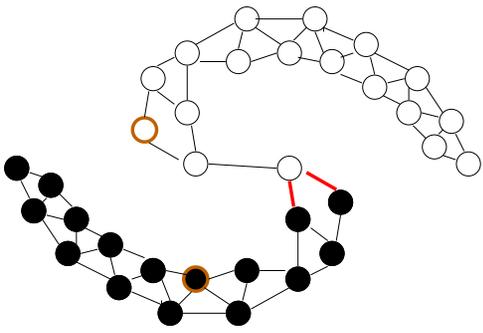
Given labeled data:



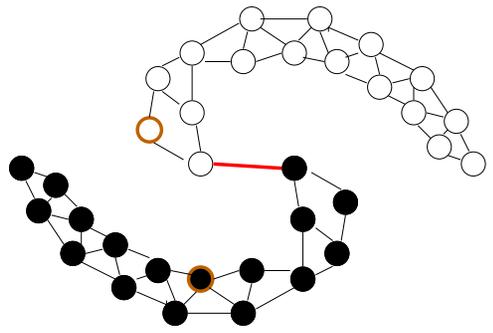
energy=0



energy=4



energy=2



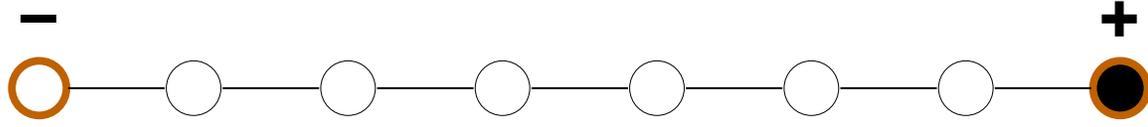
energy=1

Markov Random Fields

$$p(\mathbf{y}) \propto \exp(-E(\mathbf{y})) \mid_{\mathbf{y}^I = f^I}$$

$$y_i \in \{0, 1\}$$

Model: equiv. graph mvcut [Blum & Chawla 01]; not unique.



Hard to compute: Boltzmann machines, MCMC...

[Zhu & Ghahramani 02]

Gaussian Random Fields

Why hard: $y_i \in \{0, 1\}$ combinatorial optimization

Relax: $y_i \in \mathbb{R}$

$$\begin{aligned} p(\mathbf{y}) &\propto \exp(-E(\mathbf{y})) \Big|_{\mathbf{y}^T = f^T} \\ &= \exp\left(-\frac{1}{2} \sum_{i,j} w_{ij} (y_i - y_j)^2\right) \Big|_{\mathbf{y}^T = f^T} \\ &= \exp\left(-\mathbf{y}^T \nabla \mathbf{y}\right) \Big|_{\mathbf{y}^T = f^T} \end{aligned}$$

(1)

[Zhu & Ghahramani & Lafferty 03]

$$\Delta = \begin{bmatrix} \Delta U U & \Delta U T \\ \Delta T U & \Delta T T \end{bmatrix}$$

The Laplacian $\Delta = D - W$

$$W = \begin{bmatrix} w_{11} & \cdots & w_{1n} \\ \vdots & \ddots & \vdots \\ w_{n1} & \cdots & w_{nn} \end{bmatrix} \quad D = \begin{bmatrix} \sum w_{1\cdot} & & \\ & \ddots & \\ & & \sum w_{n\cdot} \end{bmatrix}$$

The Laplacian

Gaussian Random Fields

$$p(\mathbf{y}) = \exp\left(-\frac{1}{2} \mathbf{y}^T \mathbf{\Delta} \mathbf{y}\right) \Big|_{\mathbf{y}^T = \mathbf{f}^T}$$

Conditioned on labeled data \mathbf{f}_T :

$$\mathbf{y}_U \sim \mathcal{N}\left(\frac{1}{2} \mathbf{\Delta}^{-1} \mathbf{f}_U, \frac{1}{2} \mathbf{\Delta}^{-1}\right)$$

The mean on unlabeled data:

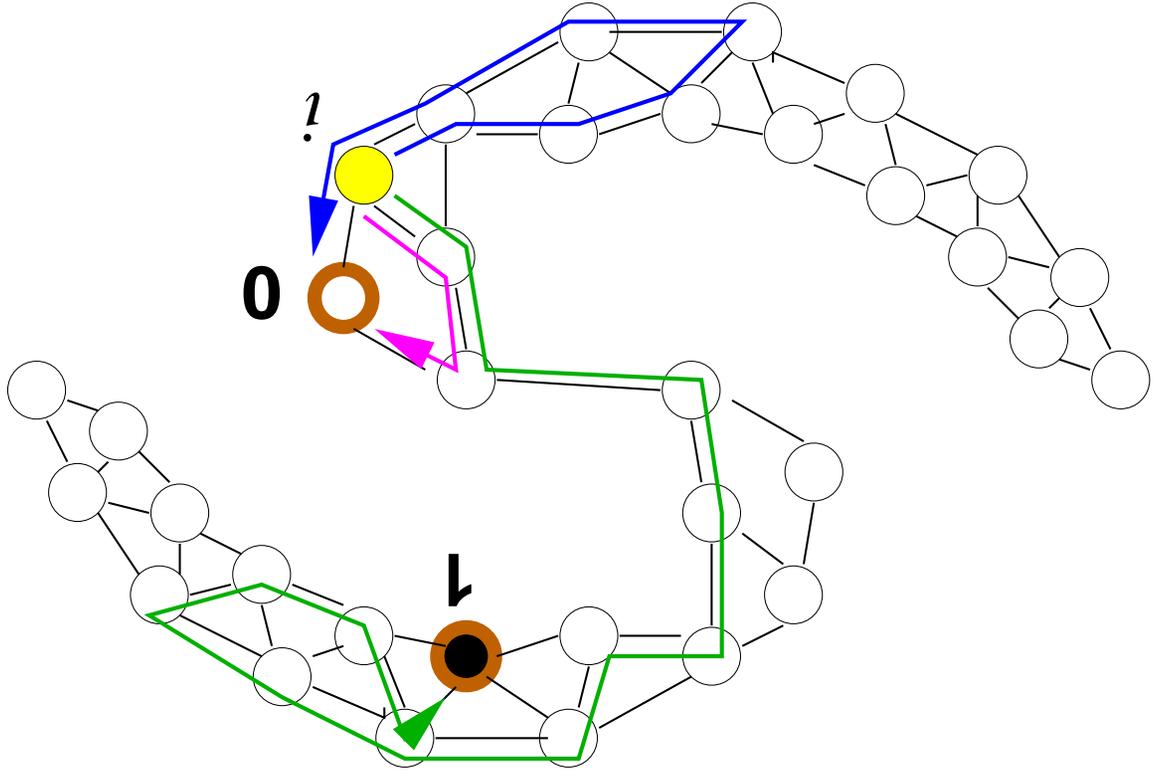
$$\mathbf{f}_U = -\mathbf{\Delta}^{-1} \mathbf{\Delta} \mathbf{f}_T$$

The Mean f_U

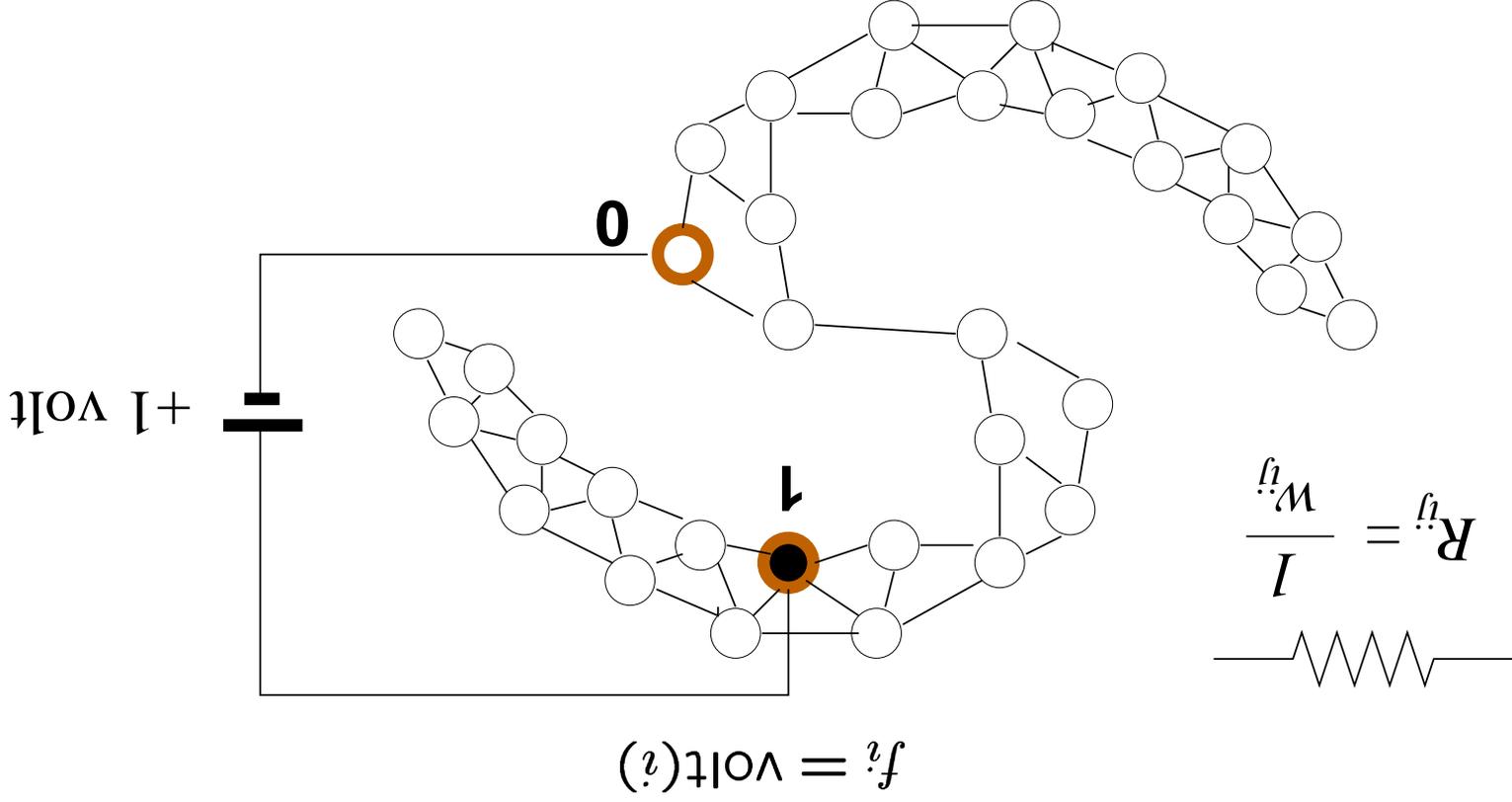
- **mode** of the Gaussian Random Field
- **min energy state**
- **uniquely exists**
- **harmonic** $\Delta f = 0$ or $f_i = \frac{\sum_{j \sim i} w_{ij} f_j}{\sum_{j \sim i} w_{ij}}$, $i \in U$
- $0 \leq f_i \leq 1$

f_u Interpretation: Random Walks

$$f_i = P(\text{reach label 1} | \text{from } i)$$



f_u Interpretation: Electric Networks



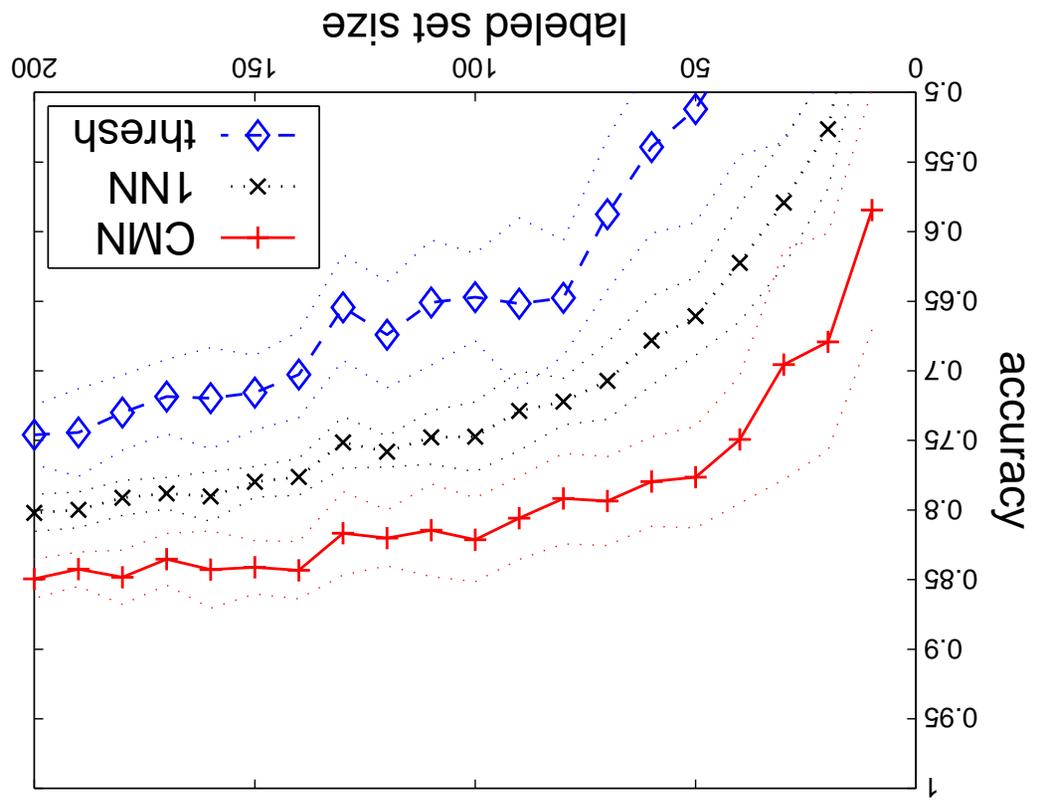
Classification Using Gaussian Fields

- naive: threshold f_U at 0.5
- **heuristic**: incorporating class proportions. E.g. prior knowledge: 90% class 1

$$\begin{aligned} & \text{minimize } E(y) = y^T \Delta y \\ & \text{subject to } y_L = f_L \\ & \sum f_U = 0.9 \end{aligned}$$

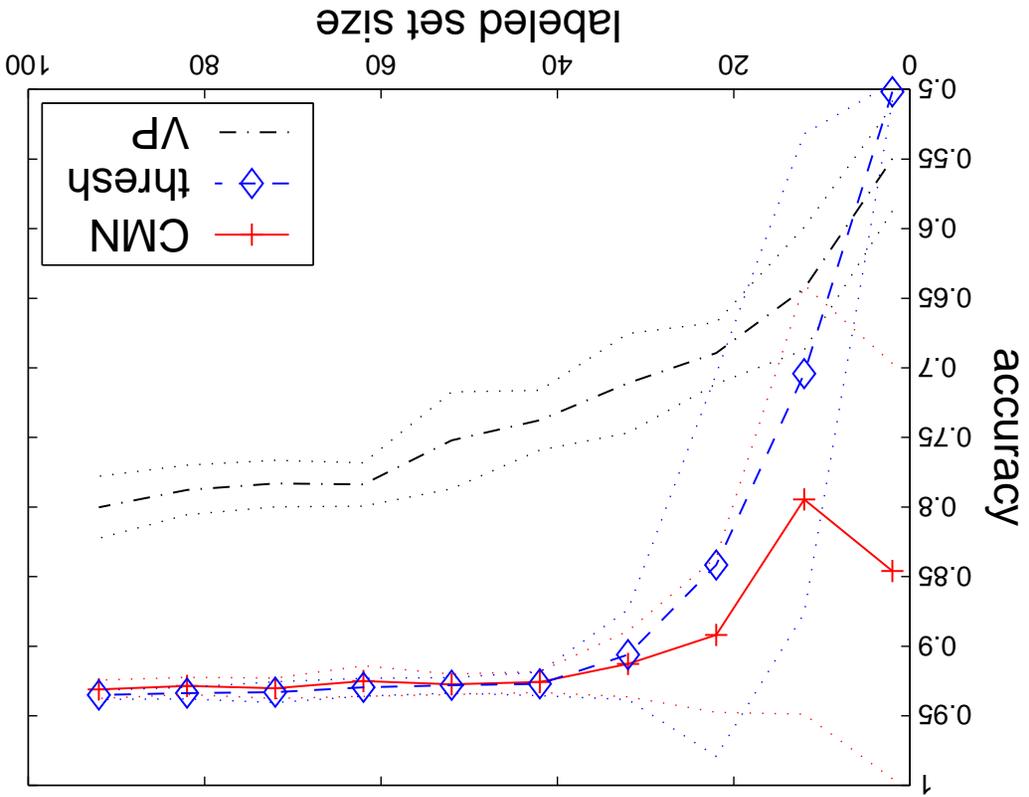
OCR Digits (0...9)

$$|L \cup U| = 4000$$



20 Newsgroups (PC vs. MAC)

$|PC| = 982, |MAC| = 961$



Learn the Graph

Parameterize W . For example:

$$w_{ij} = \exp \left(- \sum_{d=1}^p \frac{\sigma_{\frac{d}{2}}}{(x_i^d - x_j^d)^2} \right)$$

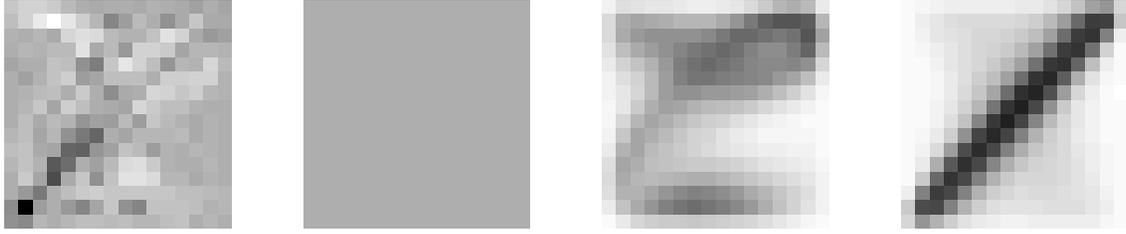
Hyperparameters $\Theta = \{\sigma_{\frac{d}{2}} | d = 1, \dots, m\}$: length scales.

Learn the Graph

- Minimize label entropy (maximize label confidence. Heuristics.)
- Evidence maximization (Local optima)
- Kernel alignment (research)

Learn the Graph

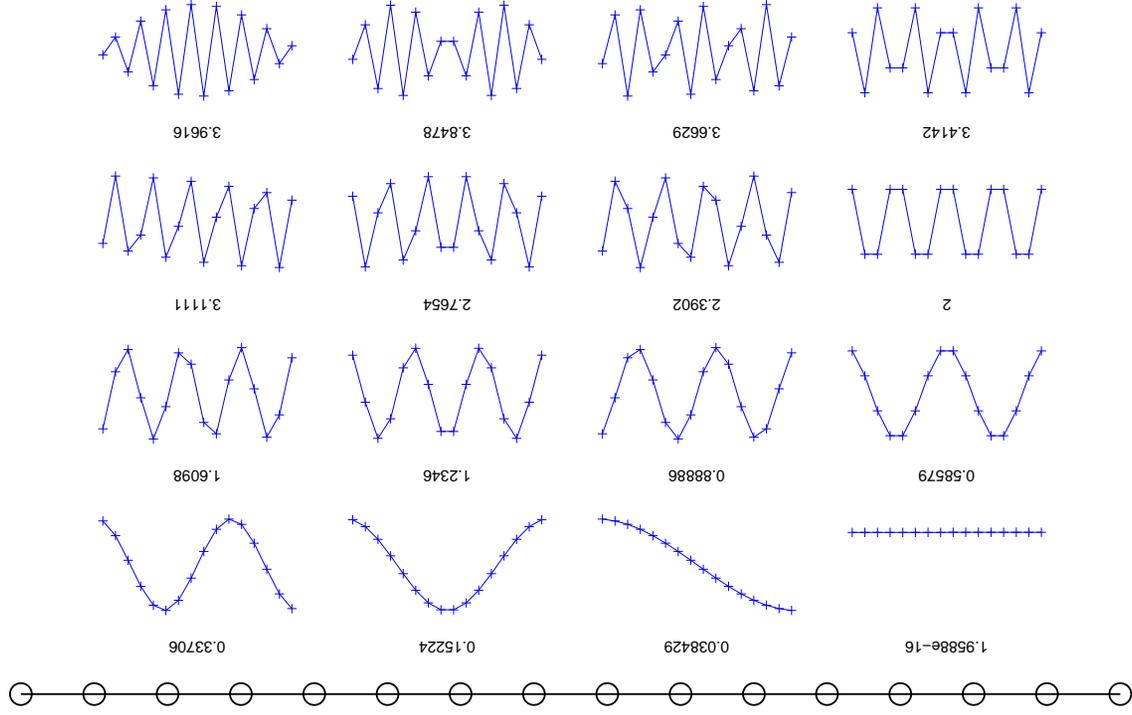
OCR Digits "1" vs. "2", minimize entropy



	H (bits)	GF acc
start	0.6931	94.70 \pm 1.19 %
end	0.6542	98.02 \pm 0.39 %

The Kernel View

Laplacian Spectrum $\Delta = \sum \lambda_i \phi_i \phi_i^\top$. λ_i : frequencies; ϕ_i : vibration modes. Low frequency = smoother.



The Kernel View

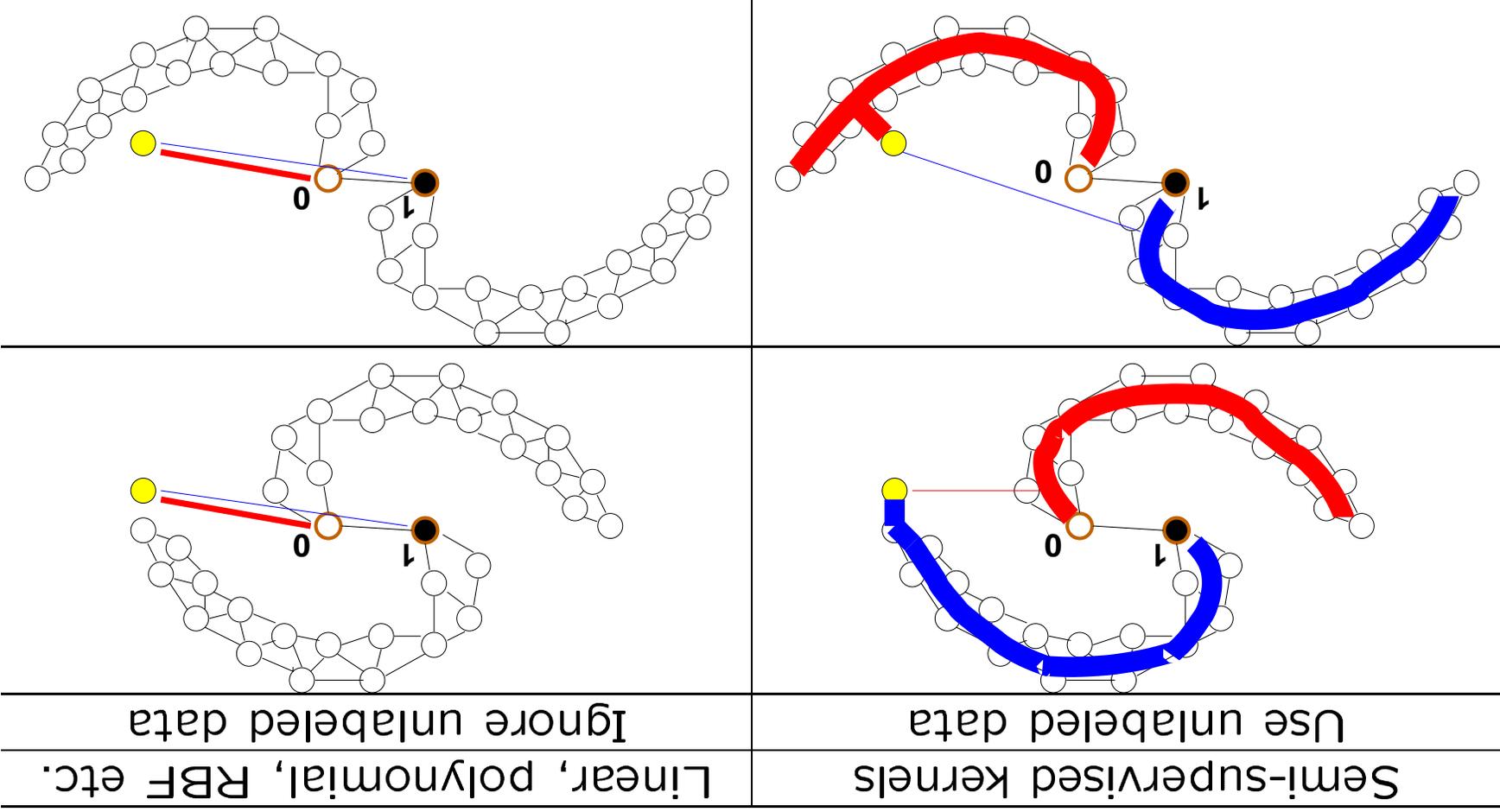
$$d(y) = \exp(-y^\top \Delta y)$$

$$K \approx \Delta^{-1} = \sum_i \frac{\lambda_i}{\phi_i \phi_i^\top}$$

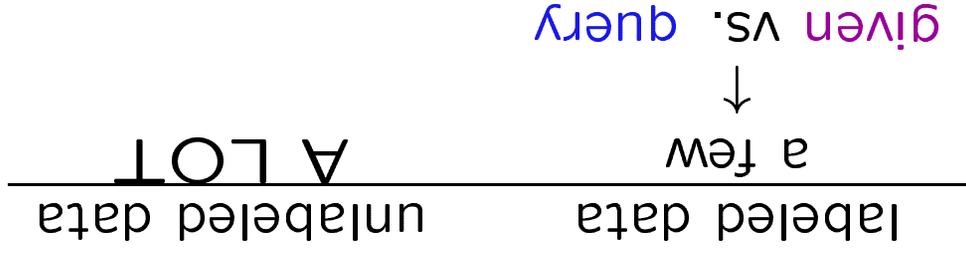
$\frac{1}{\lambda_i}$ emphasizes smooth components.

In general, a non-negative, monotonic decreasing function on the spectrum of Δ gives a semi-supervised kernel. [Smola & Kondor 03]

Semi-supervised vs. 'standard' kernels



Active Learning

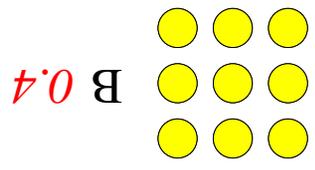


- pool based

- active + semi-supervised learning

Active Learning

Most ambiguous point not necessarily the best query: \bullet a 0.5



Active Learning

generalization error

$$\text{err} = \sum_{i \in U} \sum_{y_i=0,1} (\text{sgn}(f_i) \neq y_i) P_{\text{true}}(y_i)$$

approximation

$$P_{\text{true}}(y_i = 1) \leftarrow f_i$$

estimated error

$$\hat{\text{err}} = \sum_{i \in U} \min(f_i, 1 - f_i)$$

Active Learning

estimated error after querying k

$$\text{err}_{(x, y)}^+ = \min_{i \in U} \left(f_{(x, y)}^+ - 1, f_{(x, y)}^+ \right)$$

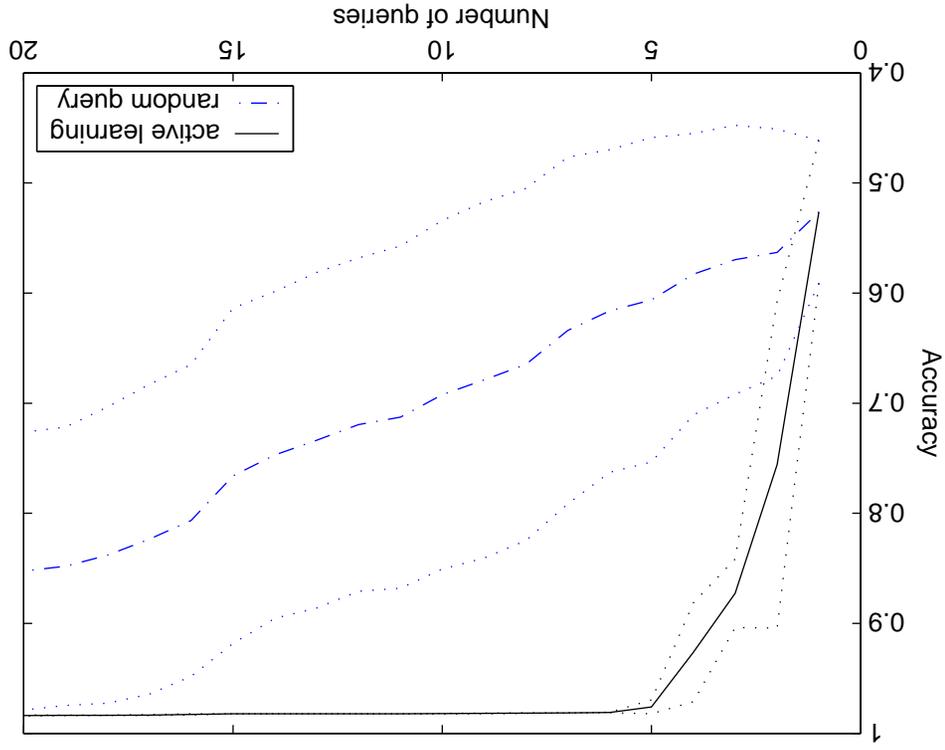
select query k^* to minimize the estimated error

$$k^* = \arg \min_k (1 - f_k) \text{err}_{(x, 0)}^+ + f_k \text{err}_{(x, 1)}^+$$

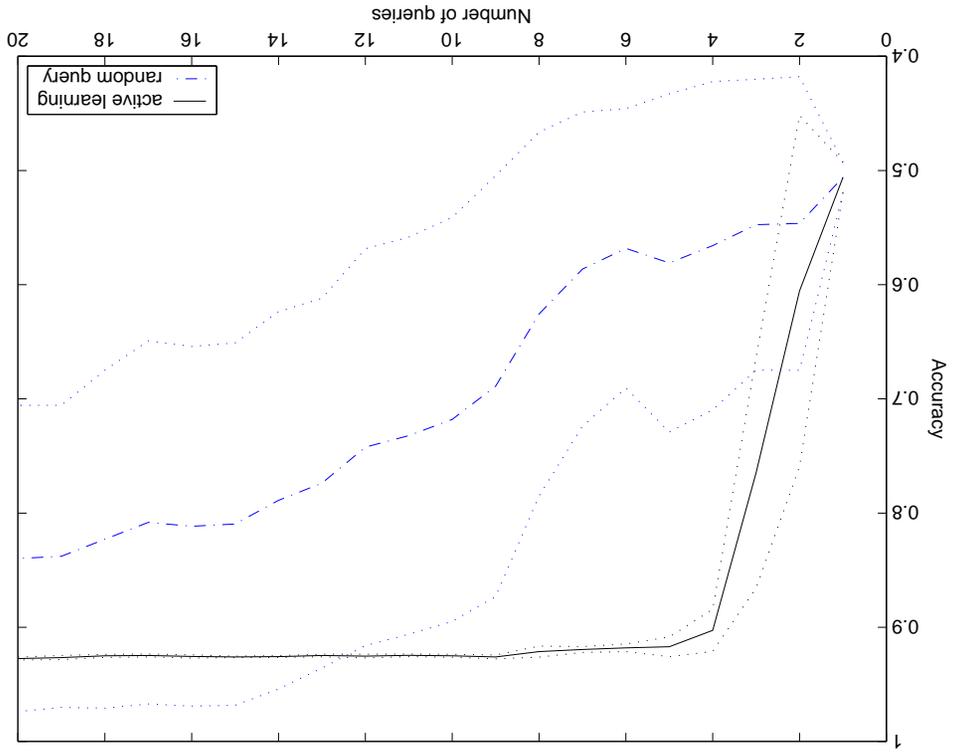
're-train' is fast

$$f_{(x, y)}^+ = f_U + \frac{(\nabla_U)_{-1}^k}{(\nabla_U)_{-1}^k} (f - h_k)$$

Active Learning: OCR Digits "1" vs. "2"

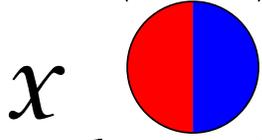


Active Learning: 20 Newsgroups PC vs. MAC

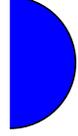
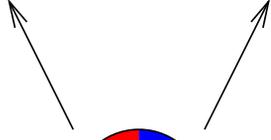


Co-training

split feature $x = \{x_1, x_2\}$



x



x_1



x_2

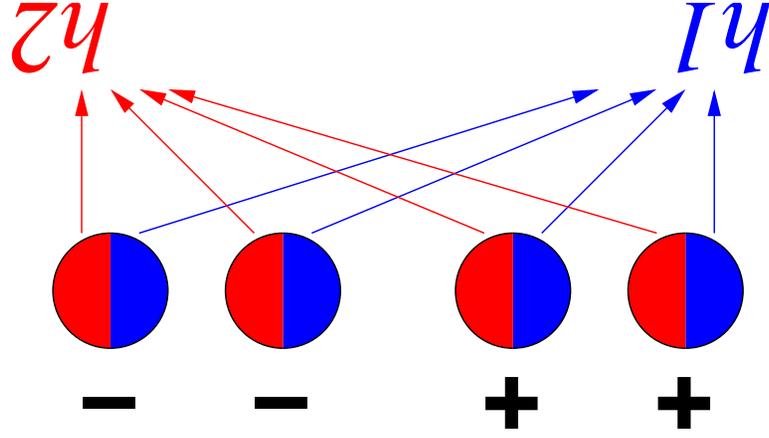
sub-features 'good' (conditional independence, sufficient)

Method
mixture
trans. SVM
co-training
graph

Co-training

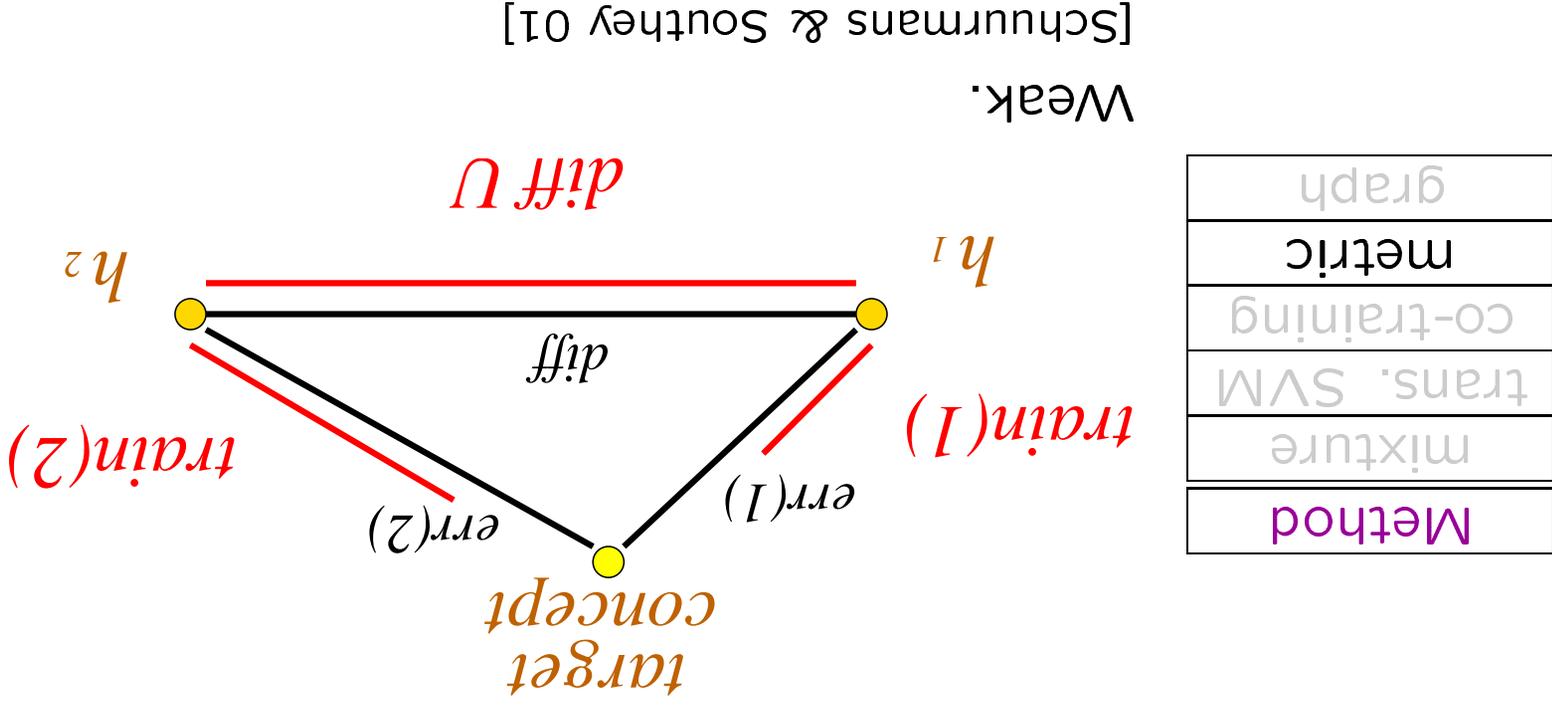
train **two** classifiers, one on each sub-

feature



Method
mixture
trans. SVM
co-training
graph

Metric-based Model Selection



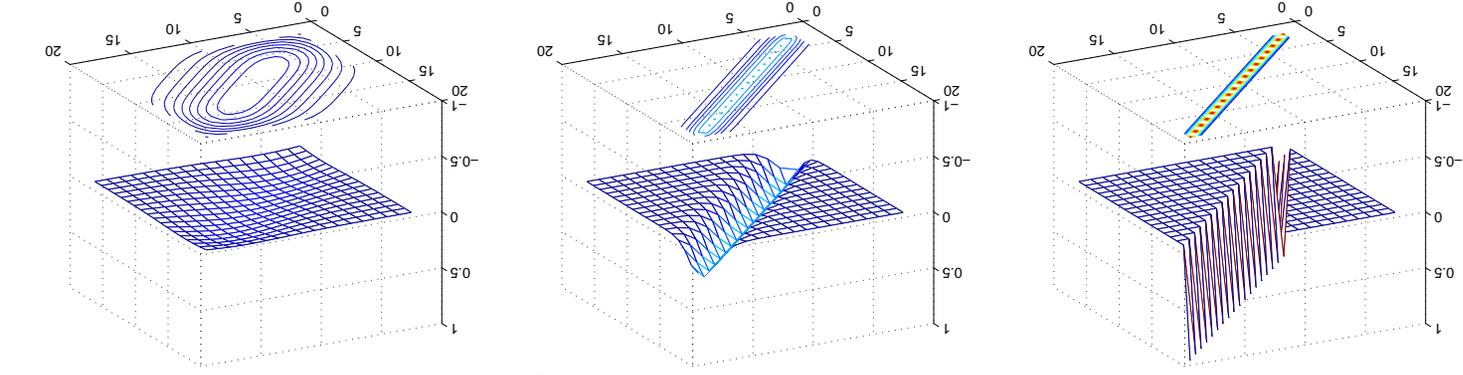
Metric-based Model Selection

Method
mixture
trans. SVM
co-training
metric
graph

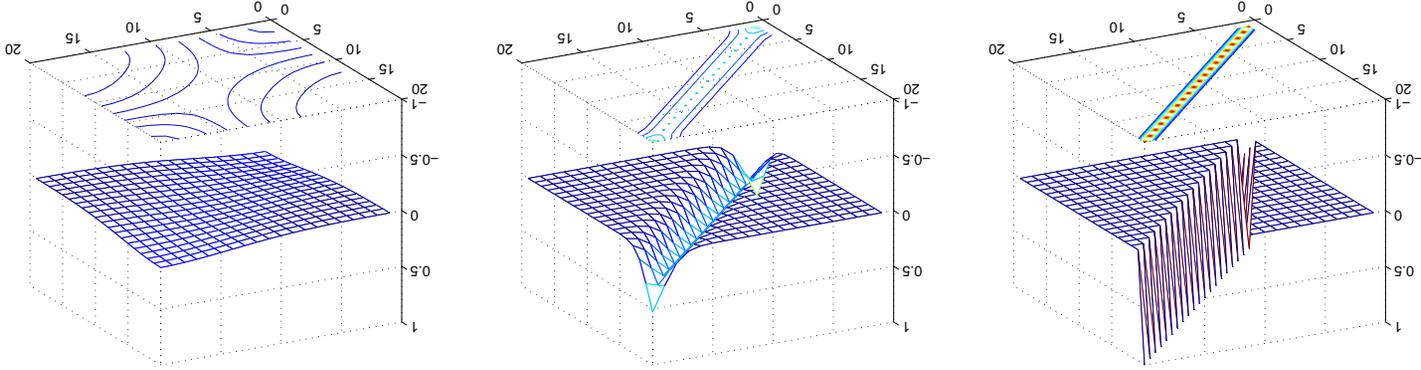
- Consider 2 hypotheses h_1, h_2
- both have zero training error
 - disagree a lot on unlabeled data
- h_1, h_2 can't be both correct.

Connection: Graph kernels

Diffusion kernel at time t : $e^{-t\Delta_{UU}}$ [Kondor & Lafferty 02]



t -step random walk: $(D^{-1}W)^t$ [Szummer & Jaakkola 01]



Connection: Graph kernels

Integrate diffusion kernels over t : $(\Delta_U U)^{-1} = \int_0^\infty e^{-t\Delta_U U} dt$

$$f_U = -(\Delta_U U)^{-1} \Delta_U U f_U$$
