

# Kernel Conditional Random Fields

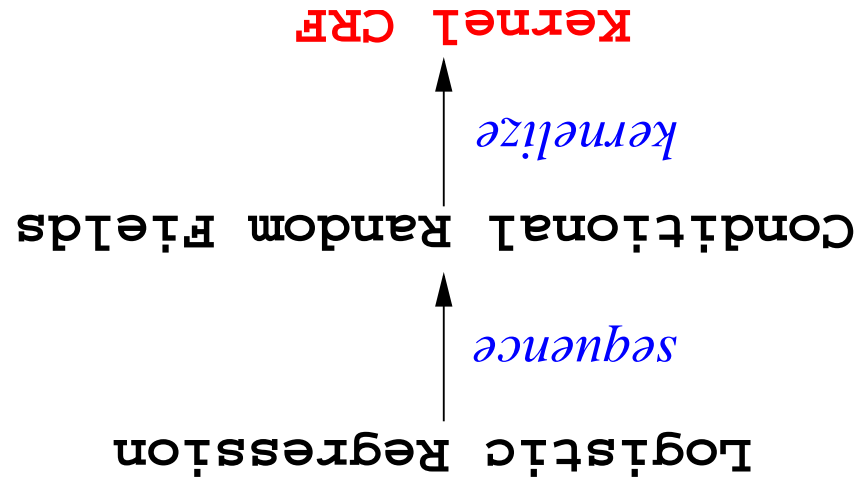
**Jerry Xiaojin Zhu**

Joint work with:

**John Lafferty**

**Yan Liu**

February 20, 2004

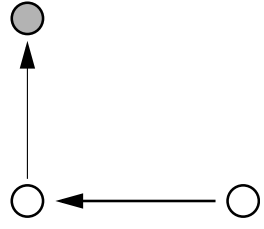


## Non-independent Data

- Independent data:
  - logistic regression, SVM, etc.
- Non-independent data: sequences, trees, graphs
  - e.g. text, speech, proteins
  - HMMs, CRFs, etc.

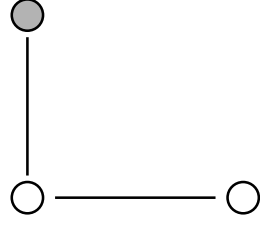
# From HMMs to CRFs: special case of sequence

HMM: 
$$p(\mathbf{y} | \mathbf{x}) \propto \prod_{t=1}^T d(y_t | y_{t-1}) d(x_t | y_t)$$



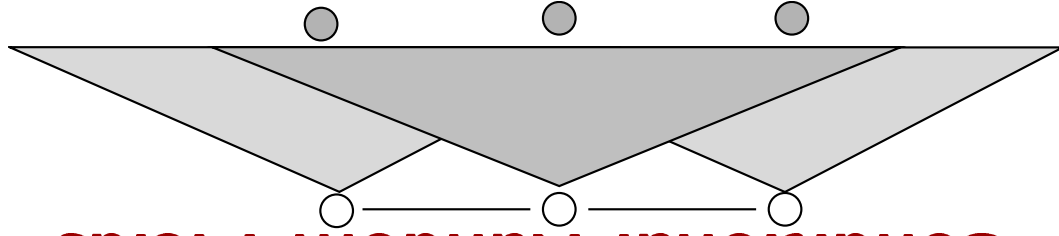
CRF:

CRF: 
$$p(\mathbf{y} | \mathbf{x}) = \frac{1}{Z} \prod_{t=1}^T \exp \left[ \sum_j \lambda_j \phi_j(y_t, y_{t-1}) + \sum_k \lambda_k \phi_k(y_t, \mathbf{x}) \right]$$



(Lafferty, McCallum and Pereira, 2001), also (Johnson et al., 2001)

## Conditional Random Fields



$$p(y | \mathbf{x}) = \frac{1}{Z(\mathbf{x})} \exp \left( \sum_{c \in \text{cliques}} \lambda^c \Phi^c(\mathbf{x}, y^c) \right)$$

- **Global** normalization. Undirected graphical models.

- Model  $p(\text{label sequence } y \mid \text{observation sequence } \mathbf{x})$  rather than joint probability  $p(y, \mathbf{x})$

- Allow arbitrary (e.g. long range) dependencies on the observation sequence

## Conditional Random Fields

- Still efficient (Viterbi, forward-backward) if dependencies within the state sequence  $y$  are local (sequences, trees).

- Promising results in

– tagging, parsing, information extraction (Collins, 2001), (Sha and Pereira, 2003), (Pinto et al., 2003)

– image processing (Kumar and Hebert, 2003)

## Explicit Features vs. Kernels

CRFs were based on **explicit feature** representations.

$$p(y | \mathbf{x}) = \frac{1}{Z(\mathbf{x})} \exp \left( \sum_{c \in \text{cliques}} \sum_{i \in \text{features}} \lambda^c \cdot \Phi^c(\mathbf{x}, y^c) \right)$$

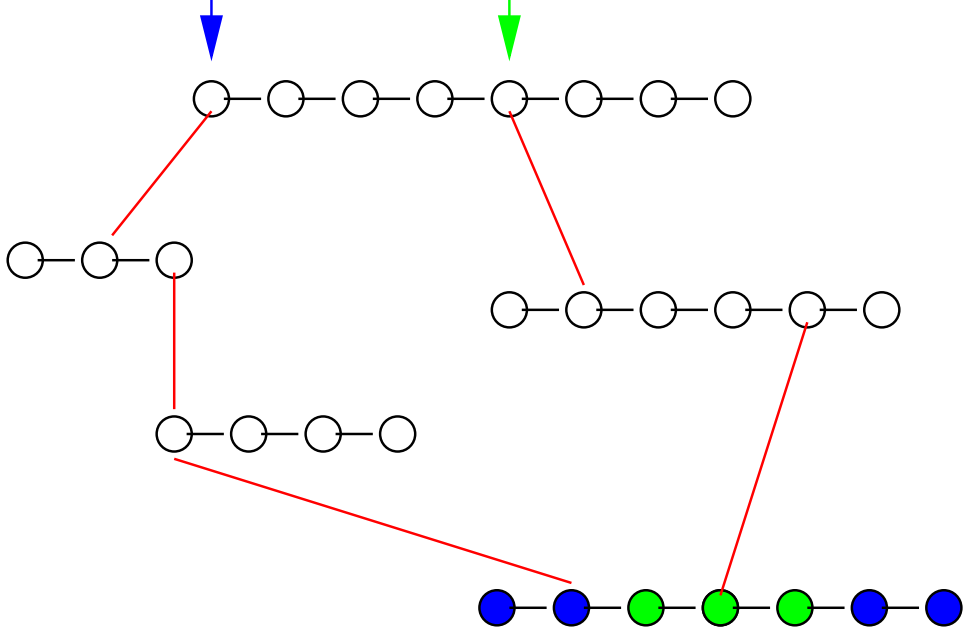
$\lambda^c$  : weight parameters  
 $\Phi^c(\mathbf{x}, y^c)$  : features

Kernels: **implicit feature** representations.

$$K_c(\mathbf{x}, y^c; \mathbf{x}', y'^c) = \sum_i \Phi^{c_i}(\mathbf{x}, y^c) \Phi^{c_i}(\mathbf{x}', y'^c)$$

## Why kernels

- Kernel machines are very successful;
- Semi-supervised learning on sequences with graph kernels.





$$= - \sum_{c \in \mathcal{C}} f_c(\mathbf{x}, y_c) + \log \sum_{y'} \exp \left( \sum_{c \in \mathcal{C}} f_c(\mathbf{x}, y'_c) \right) = - \log p(y | \mathbf{x})$$

The negative log loss (logistic loss)

$$p(y | \mathbf{x}) = \frac{\exp(Z(\mathbf{x}))}{\sum_{c \in \mathcal{C}} \exp(f_c(\mathbf{x}, y_c))}$$

Introduce functions  $f_c$

$$p(y | \mathbf{x}) = \frac{\exp(Z(\mathbf{x}))}{\sum_{c \in \mathcal{C}} \exp(\lambda_c \phi_c(\mathbf{x}, y_c))}$$

## The Negative Log Loss

## Regularized Risk

$$R_{\psi} f = \sum_{i=1}^n \psi(\mathbf{x}^{(i)}, \mathbf{y}^{(i)}, f) + \Omega(\|f\|)$$

Regularizer  $\Omega$  is a monotonically increasing function.

The risk minimizer  $f^*$  is the MAP estimate of the CRF.

A different  $\psi$  (hinge loss?) seems to correspond to Max-Margin Markov ( $M_3$ ) Networks (Taskar et al., 2003).

## Representer Theorem for CRFs

The minimizer  $f_*$  of

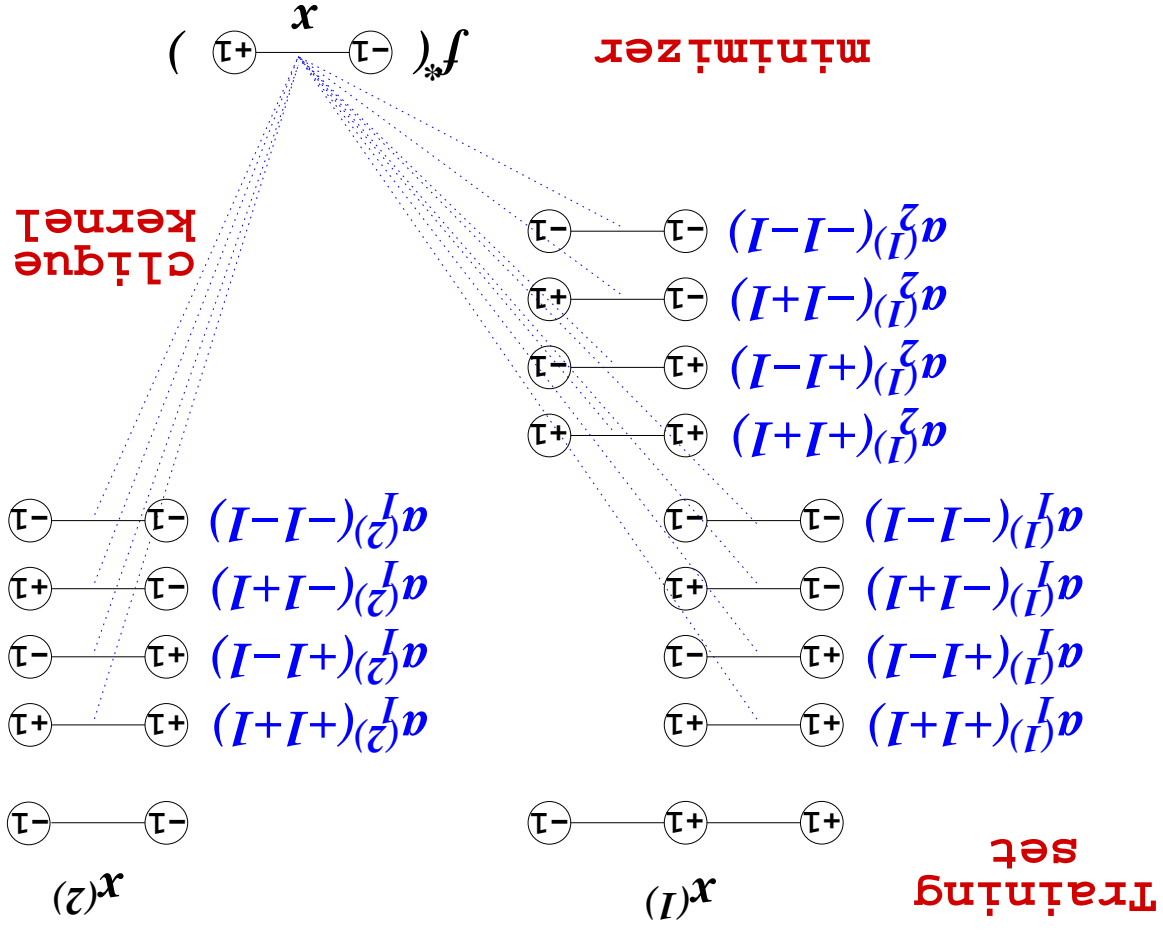
$$R_{\psi} f = \sum_n \phi(\mathbf{x}^{(i)}, \mathbf{y}^{(i)}, f) + \Omega(\|f\|)$$

has the form

$$f_*^c(\cdot) = \sum_n \sum_{i \in \text{cliques}} \sum_{\text{all } y_c} \alpha_c^{(i)}(y_c) K_c(\mathbf{x}^{(i)}, \mathbf{y}_c; \cdot)$$

Dual parameters  $\alpha_c^{(i)}(y_c)$ : all labeling, not only those in training.

# Example 1: Edge cliques



## Example 2: Kernel Logistic Regression

With only vertex cliques, and

$$K(\mathbf{x}, y; \mathbf{x}', y') = K(x, x') \delta(y, y')$$

The minimizer has the form

$$f_*(x, y) = \sum_{v \in \text{vertices}} \alpha_v(y) K(x_v, x)$$

This is simply kernel logistic regression.

## The KCRF training problem

Given training set  $\{(\mathbf{x}^{(i)}, \mathbf{y}^{(i)})\}$ , clique kernels  $K_c$ , find the dual parameters

$$\alpha_c^{(i)}(\mathbf{y}^c)$$

for

$$f_c(\cdot) = \sum_n \sum_{c \in \text{cliques}} \alpha_c^{(i)}(\mathbf{y}^c) K_c(\mathbf{x}^{(i)}, \mathbf{y}^c; \cdot)$$

that minimizes the regularized risk

$$R_{\psi} f = \sum_n \psi(\mathbf{x}^{(i)}, \mathbf{y}^{(i)}; f) + \frac{1}{2} \|f\|_2^2$$

unconstrained convex optimization problem, global solution.

## The Derivatives

Each parameter  $\alpha_c^{(i)}(y_c)$  is associated with a basis function

$$h(\cdot) = K_c(\mathbf{x}^{(i)}, \mathbf{y}_c; \cdot)$$

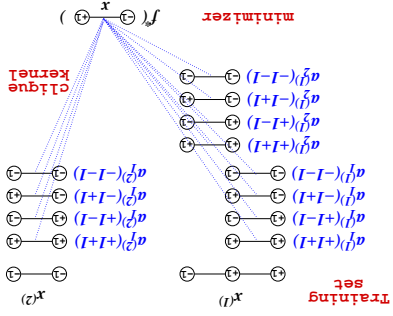
$$\frac{\partial R_\psi f}{\partial \alpha_c^{(i)}(y_c)} = E_f[h] - \tilde{E}[h] + \epsilon, h > K$$

$E_f[h]$ : expectation of  $h$  under the current KCRF parameters  
(need clique marginals, forward-backward)

$\tilde{E}[h]$ : training set empirical expectation of  $h$

Newton's algorithm, BFGS.

# Too Many Parameters



Every labeling of every training clique has a parameter.

(num of training cliques) \* (clique size) num of labels

Negative log loss (KCRF)  $\rightarrow$  everything is a 'support clique'

Sparse training: select a good subset.



## Sparse Training: Greedy Clique Selection

Initially Active Set  $\mathcal{A}$  empty, all parameters zero. Repeat:

1. candidates  $\leftarrow$  training cliques not in  $\mathcal{A}$
2. select the candidate clique  $K_c(\mathbf{x}^{(i)}, y_c; \cdot)$  with the highest **gain**.  $\mathcal{A} \leftarrow \mathcal{A} \cup \{K_c(\mathbf{x}^{(i)}, y_c; \cdot)\}$
3. train KCRF with the active cliques in  $\mathcal{A}$  by minimizing  $R_{\psi} f$ .

## Gain

Gain:  $R_{\psi}(f) - R_{\psi}(f, K_c(\mathbf{x}^{(i)}, \mathbf{y}_c; \cdot))$ ; Problem: has to estimate  $\alpha_c^{(i)}(\mathbf{y}_c)$ .

Linear approximation (functional derivative):

$$R_{\psi}(f, K_c(\mathbf{x}^{(i)}, \mathbf{y}_c; \cdot)) \approx R_{\psi}(f) + \epsilon \frac{\partial R_{\psi} f}{\partial \alpha_c^{(i)}(\mathbf{y}_c)}$$

Select the candidate with the largest gradient magnitude:

$$\left| \frac{\partial R_{\psi} f}{\partial \alpha_c^{(i)}(\mathbf{y}_c)} \right| = \left| E_f[h] - \tilde{E}[h] \right| + \epsilon, h > k$$

The clique whose model and empirical expectations mismatch the most.

## Detour: Semi-Supervised Learning

Classification needs labeled training data.

Often, labeled data scarce, unlabeled data abundant.

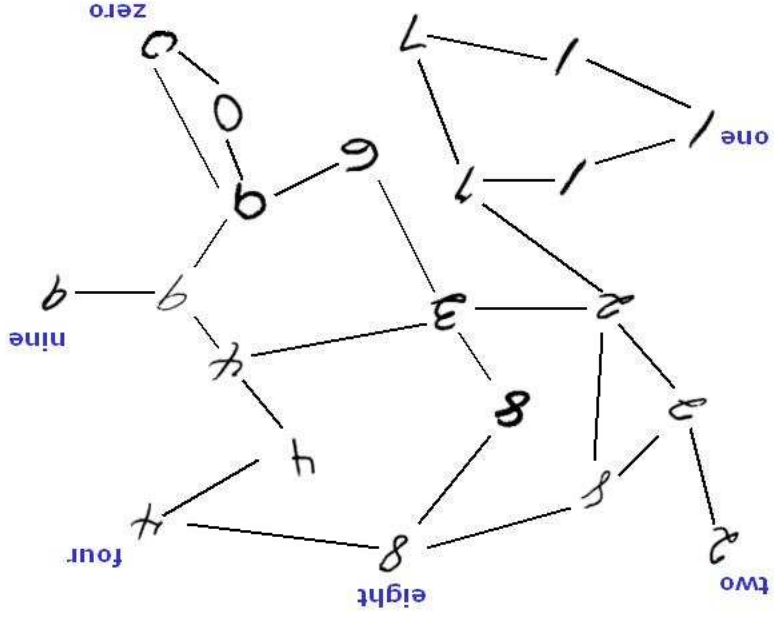
Semi-supervised learning: use **unlabeled** data to help classification.

Graph kernels: emerging theme in semi-supervised learning.

# Semi-Supervised Learning: the Graph

(Note: not the KCRF graph)

nodes: labeled and unlabeled data; edges: local similarity

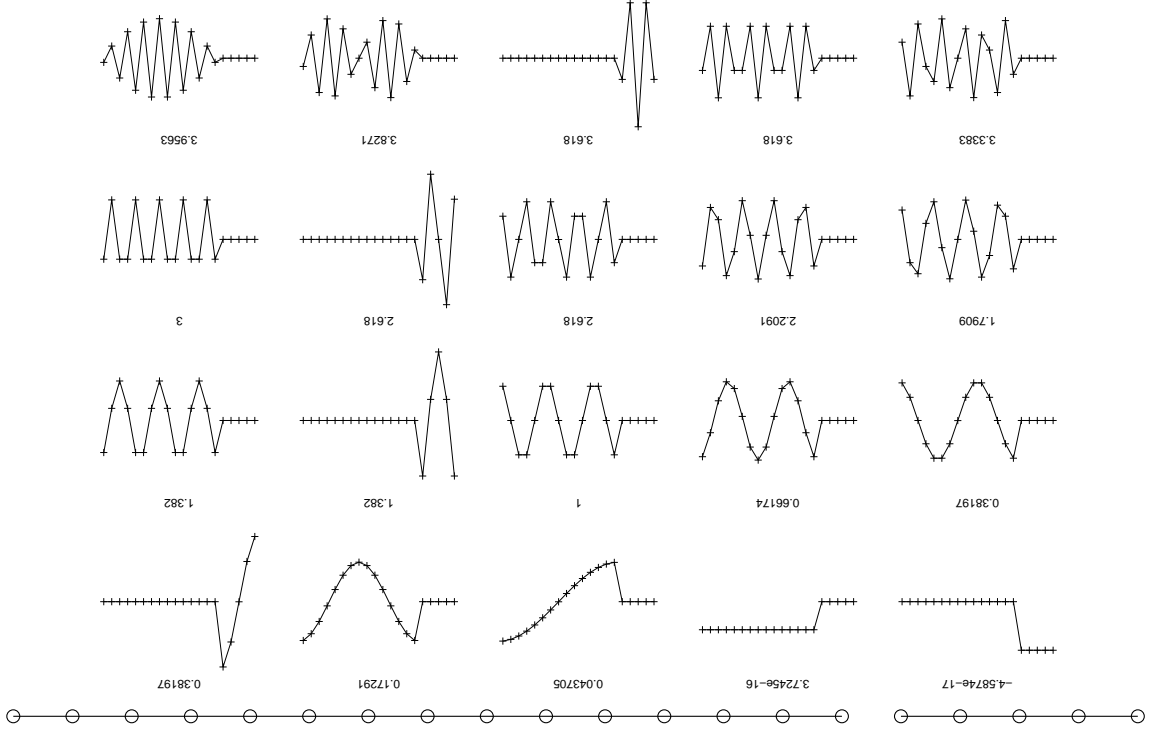


Labels propagate, smooth on graph

# Semi-Supervised Learning: the Laplacian

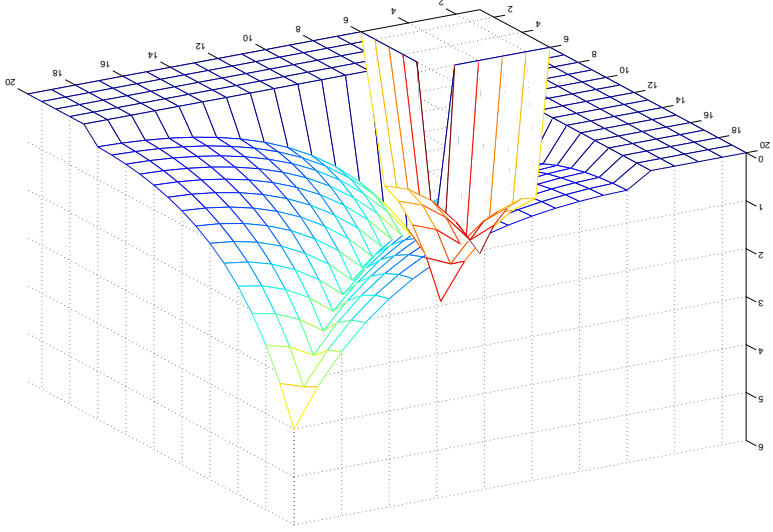
weights  $W$ ; degrees (diagonal)  $D$ ; Laplacian  $D - W$   
 Spectrum  $D - W = \sum_i \lambda_i \phi_i \phi_i^T$

$\lambda_i$ : frequencies;  $\phi_i$ : vibration modes. Low frequencies smoother.



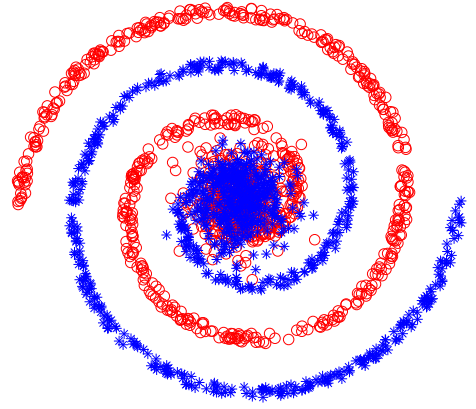
# Semi-Supervised kernels

Emphasize low frequency components, e.g.  $K = \sum_i \frac{\lambda_i}{\lambda_i + 0.05} \phi_i \phi_i^T$



(Smola and Kondor, 2003), (Zhu, Ghahramani and Lafferty, 2003), etc.

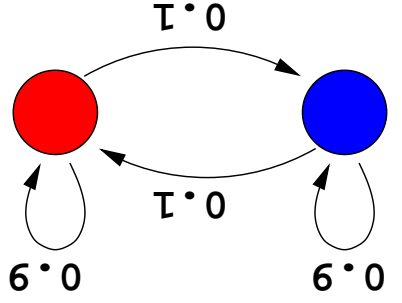
# Back to KCRF: a Synthetic Example



sequences: generated from an HMM

Semi-supervised kernel: 10 nearest neighbor unweighted graph

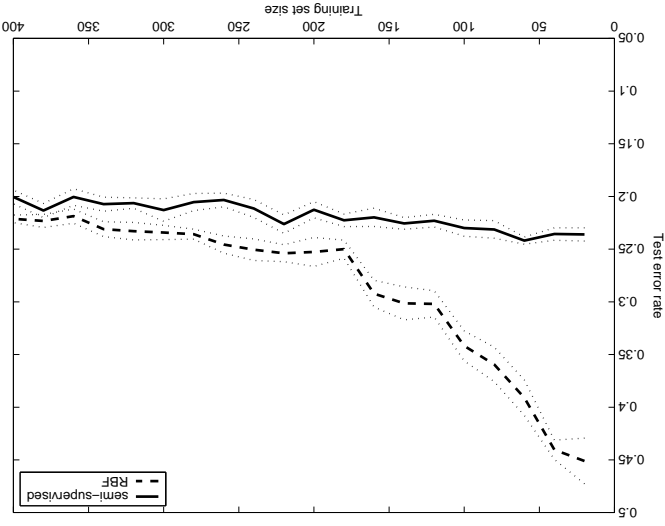
RBF kernel



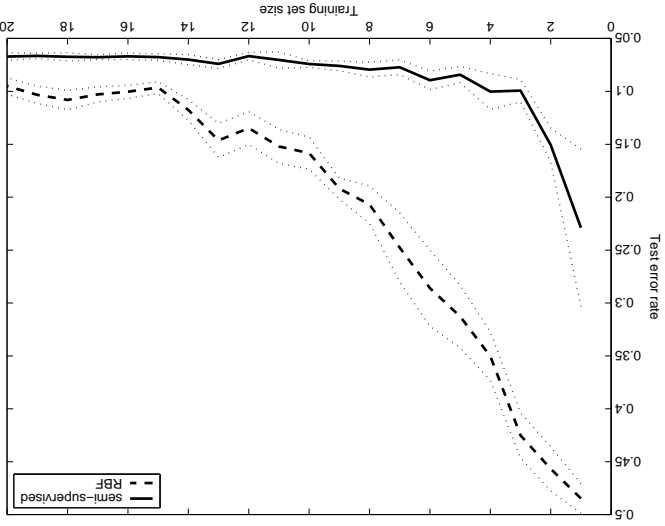
# Synthetic Example

test error vs. training set size

kernel logistic regression



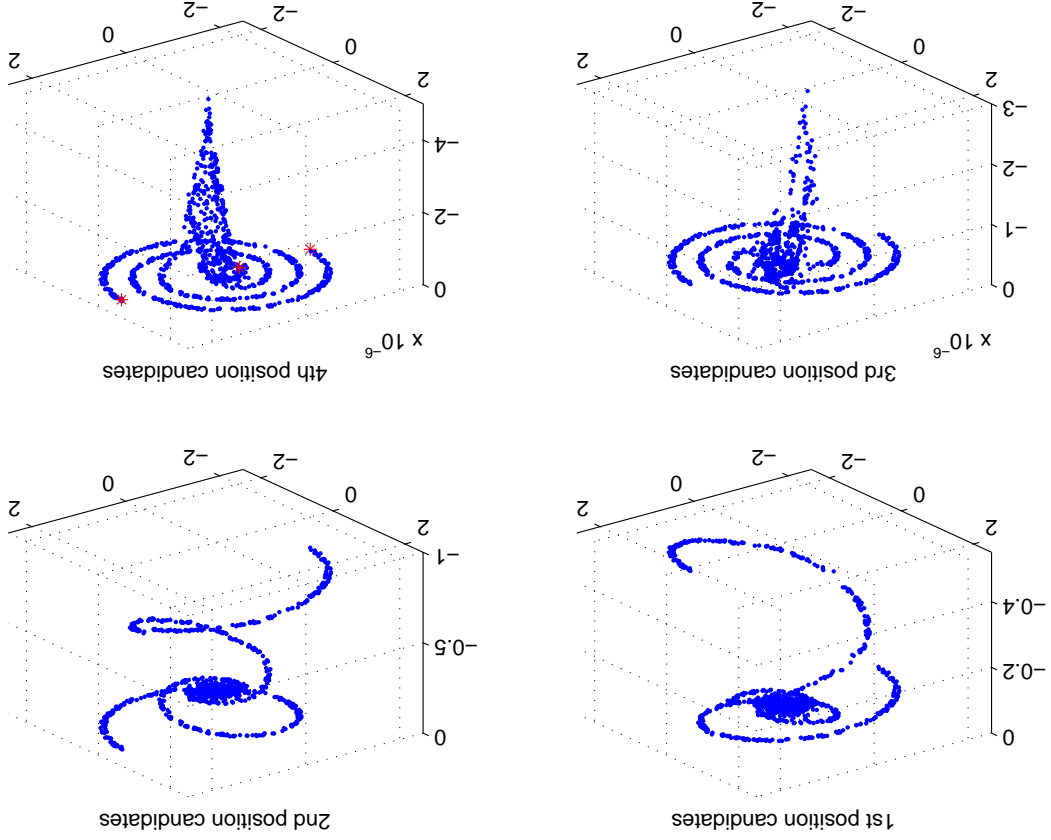
KCRR





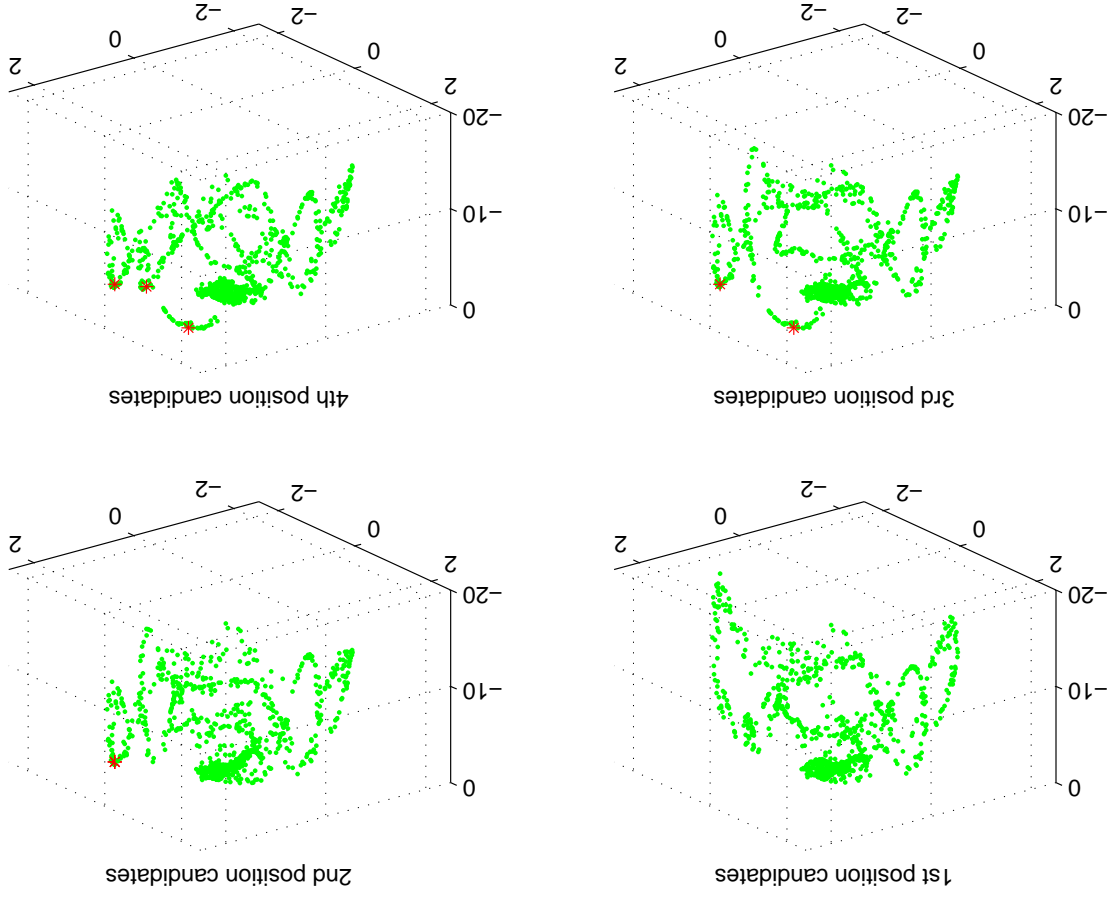
# Synthetic Example

KCRF candidate gain, graph kernel



# Synthetic Example

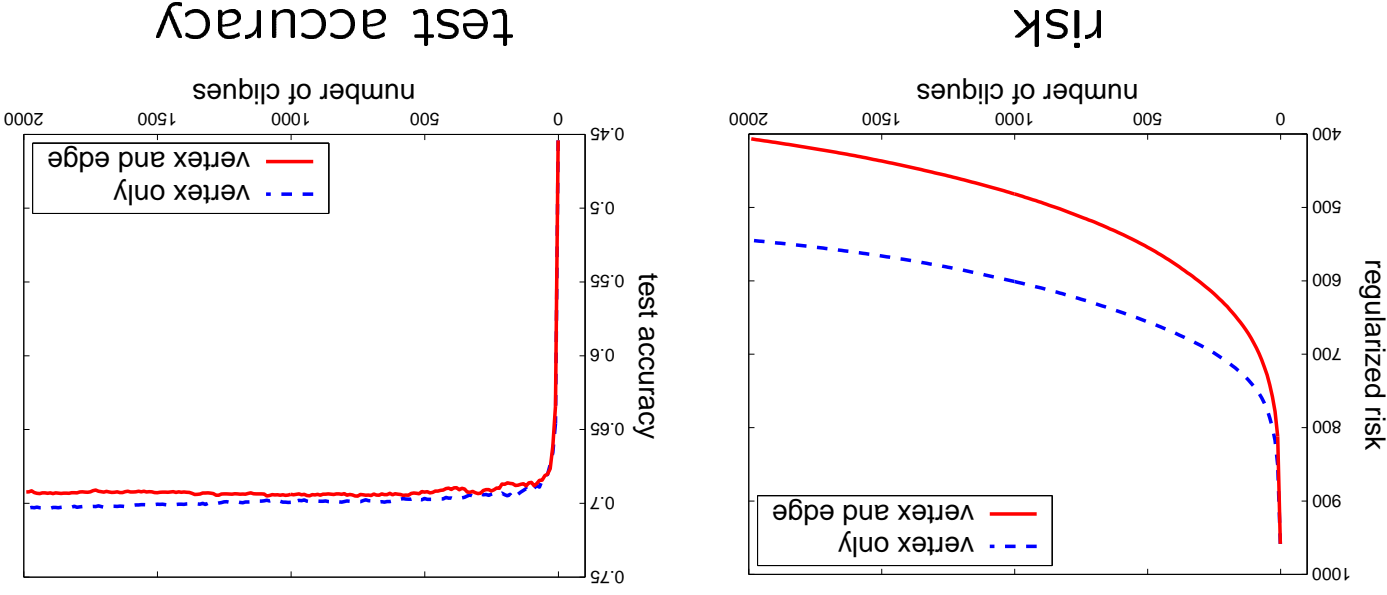
KCRF candidate gain, RBF kernel



# Protein Secondary Structure Prediction

RS126 dataset, 117 protein sequences; PSI-BLAST feature; Three classes: coil, sheet, helix; RBF kernel

KCRF clique selection



# Protein Secondary Structure Prediction

Test Accuracy

Method	Accuracy	std	Accuracy	std
KCRF (v)	0.6625	0.0224	0.6933	0.0276
KCRF (v+e)	0.6562	0.0202	0.6933	0.0272
SVM	0.6509	0.0307	0.6875	0.0235

\* KCRFs select 300 cliques.

## Protein Secondary Structure Prediction

Transition Accuracy

transition: a pair of adjacent positions with different true labels.  
A transition is classified correctly only if both labels are correct.

Method	Accuracy	std	Accuracy	std
KCRF (v)	0.1097	0.0271	0.1462	0.0235
KCRF (v+e)	0.1114	0.0250	0.1522	0.0214
SVM	0.0667	0.0313	0.1066	0.0311

5 protein set      10 protein set

## Conclusion

- KCRF: framework for graph-structured classification.
- Representer theorem
- Clique selection
- Needs kernels that capture the structure of the data.

