

New Paradigms in Pattern Recognition

Sargur N. Srihari

srihari@cedar.buffalo.edu

Outline

- Pattern Recognition and Machine Learning
- Generative and Discriminative Classifiers
 - Conditional Random Fields
- Document Recognition Applications
- Summary and Conclusion

Pattern Recognition

- Searching for Patterns in data is a fundamental one with a long and successful history
 - 16th cent:

Empirical observations of Tycho Brahe allowed Johannes Kepler to discover empirical laws of planetary motion which in turn provided springboard for classical mechanics
 - Early 20th cent:

Discovery of regularities in atomic spectra played key role in development and verification of quantum mechanics

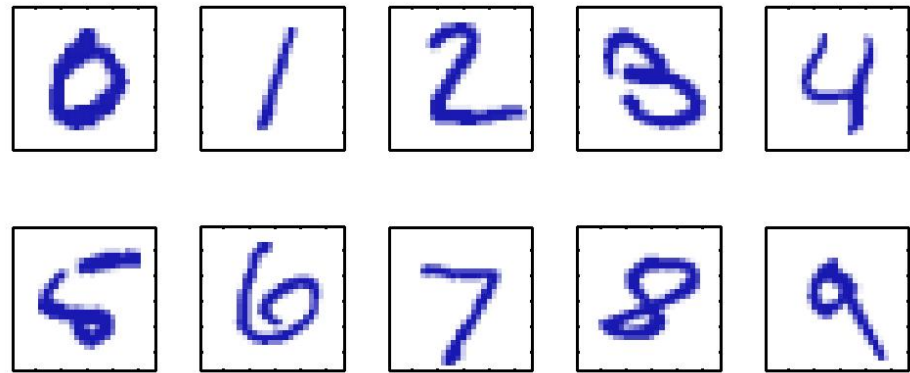
Machine Learning

- Programming computers to use example data or past experience
- Well-Posed Learning Problems
 - A computer program is said to learn from *experience E*
 - with respect to *class of tasks T* and *performance measure P*,
 - if its performance at tasks *T*, as measured by *P*, improves with experience *E*.

Pattern Recognition and Machine Learning

- Two facets of the same field
- Pattern Recognition origins are in engineering
 - discovering regularities in data and taking actions such as classification
- Machine Learning grew out of computer science
 - Needed in cases where we cannot directly write a computer program but have example data

Example Problem: Handwritten Digit Recognition



Wide variability of same numeral

- Handcrafted rules will result in large no of rules and exceptions
- Better to have a machine that learns from a large training set

Meeting of the MINDS

- Machine Learning
- Information Retrieval
- Natural Language Processing
- Document Analysis and Recognition
- Speech Recognition

ML in NLP: Text Categorization Problem

Task T :

Assign a document to its content category

Performance measure P :

Precision and Recall

Training experience E :

Example pre-classified documents

ML in Speech and IR

- **Speech Recognition**
 - Speaker-specific strategies for recognizing phonemes and words
 - Neural networks and methods for learning HMMs for customizing to individual speakers, vocabularies and microphone characteristics
- **Information Retrieval and Data Mining**
 - Information extraction from text
 - Very large databases to learn general regularities in data

Major Developments in PR-ML in Last 10 years

- Bayesian methods
 - From niche to main stream
 - Approximate inference using variational methods
- Graphical models
 - Emerged as general framework for describing and applying probabilistic models
- Models based on kernels have had significant impact
 - SVM
- Distance-based methods

PR-ML Methods

- Generative Methods
 - Model class-conditional pdfs and prior probabilities
 - Known as generative since by sampling them can generate synthetic data points
- Discriminative Methods
 - Directly estimate the posterior probabilities

Generative Models

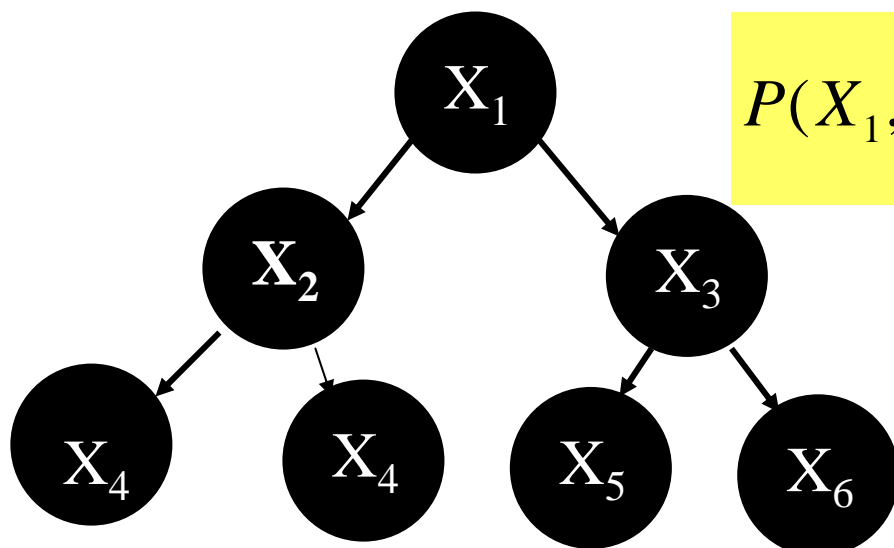
- System's input features, output response and latent variables are represented by joint pdfs
- Popular models:
 - Gaussians, Naïve Bayes, Mixtures of multinomials
 - Mixtures of Gaussians, Mixtures of experts, HMMs
 - Sigmoidal belief networks, Bayesian networks
 - Markov random fields
- Prior knowledge is added parametrically

Generative Learning

- Given a problem domain with variables X_1, \dots, X_T system is specified with a joint pdf $P(X_1, \dots, X_T)$
- Given a full joint pdf we can
 - Marginalize
$$P(X_j) = \sum_{\forall X_{i,i \neq j}} P(X_1, \dots, X_n)$$
 - Condition
$$P(X_j | X_k) = \frac{P(X_j, X_k)}{P(X_k)}$$
- By conditioning a joint pdf we can easily form
 - Classifiers, regressors, predictors

Constraining joint distribution

- To have fewer degrees of freedom
 - Conditional independencies between variables

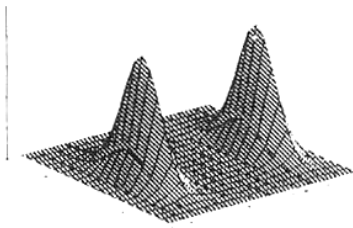


Bayesian Network

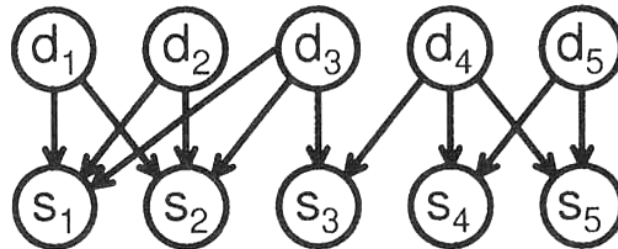
$$P(X_1, \dots, X_n) = \prod_{i=1}^n P(X_i | X_{\Pi_i})$$

Graph shows that joint distribution factorizes over variables given their parents

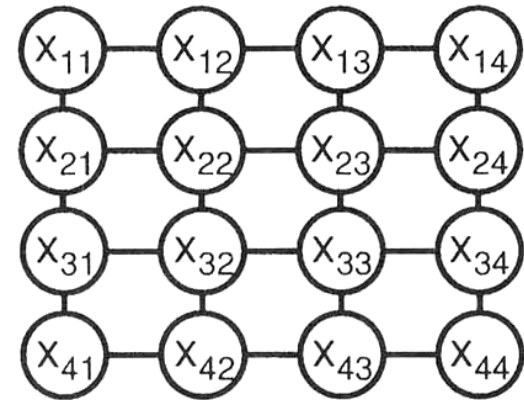
Generative Models (graphical)



(a) Mixture Model



(b) Directed Graphical Model



(c) Undirected Graphical Model

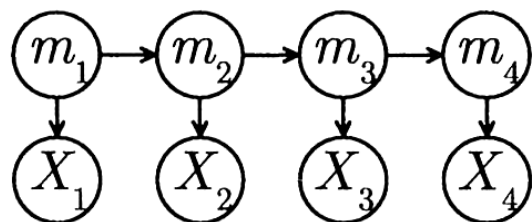
Parent node
selects between
components

Quick Medical
Reference -DT

Diagnosing
Diseases from
Symptoms

Markov Random
Field

Generative Hidden Markov Model



HMM is a distribution depicted by a graphical model.
Highly structured network indicates conditional independences,
Past states independent of future states
Conditional independence of observed given its state

Successes of Generative Methods

- NLP
 - Traditional rule-based or Boolean logic systems (eg Dialog and Lexis-Nexis) are giving way to statistical approaches (Markov models and stochastic context free grammars)
- Medical Diagnosis
 - QMR knowledge base, initially a heuristic expert systems for reasoning about diseases and symptoms has been augmented with decision theoretic formulation
- Genomics and Bioinformatics
 - Sequences represented as generative HMMs

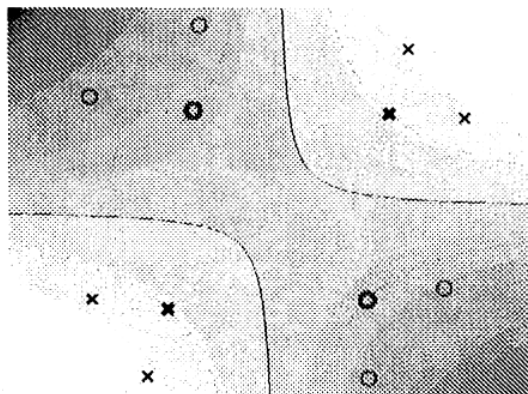
Discriminative Classifiers

- Make no attempt to model underlying distributions
- Only interested in optimizing a mapping from inputs to desired outputs
- Focuses model and computational resources on given task and provides better performance
- Examples:
 - logistic regression
 - SVMs
 - traditional neural networks
 - Nearest neighbor

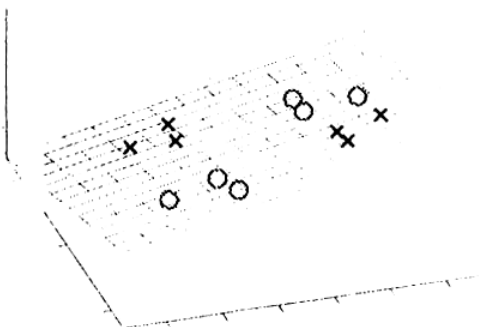
$$P(y = 1 | X) = \frac{1}{(1 + \exp(-\theta^T X))}$$

sigmoid

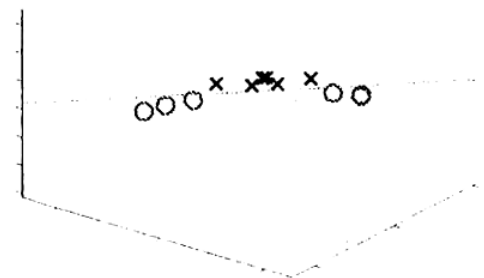
Discriminative Classifier: SVM



(a) Support Vector Machine



(b) Mapped Data



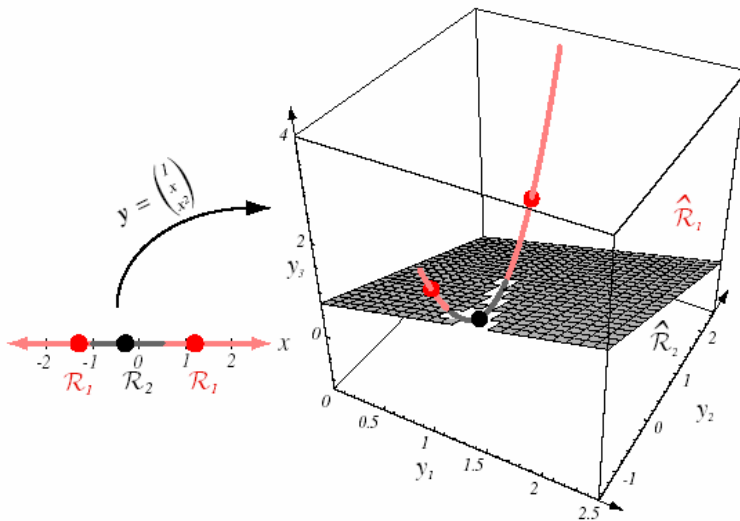
(c) Separating Hyperplane

$$(x_1, x_2) \rightarrow (x_1, x_2, x_1x_2)$$

Nonlinear decision boundary

Linear boundary

SVM Nonlinear Case



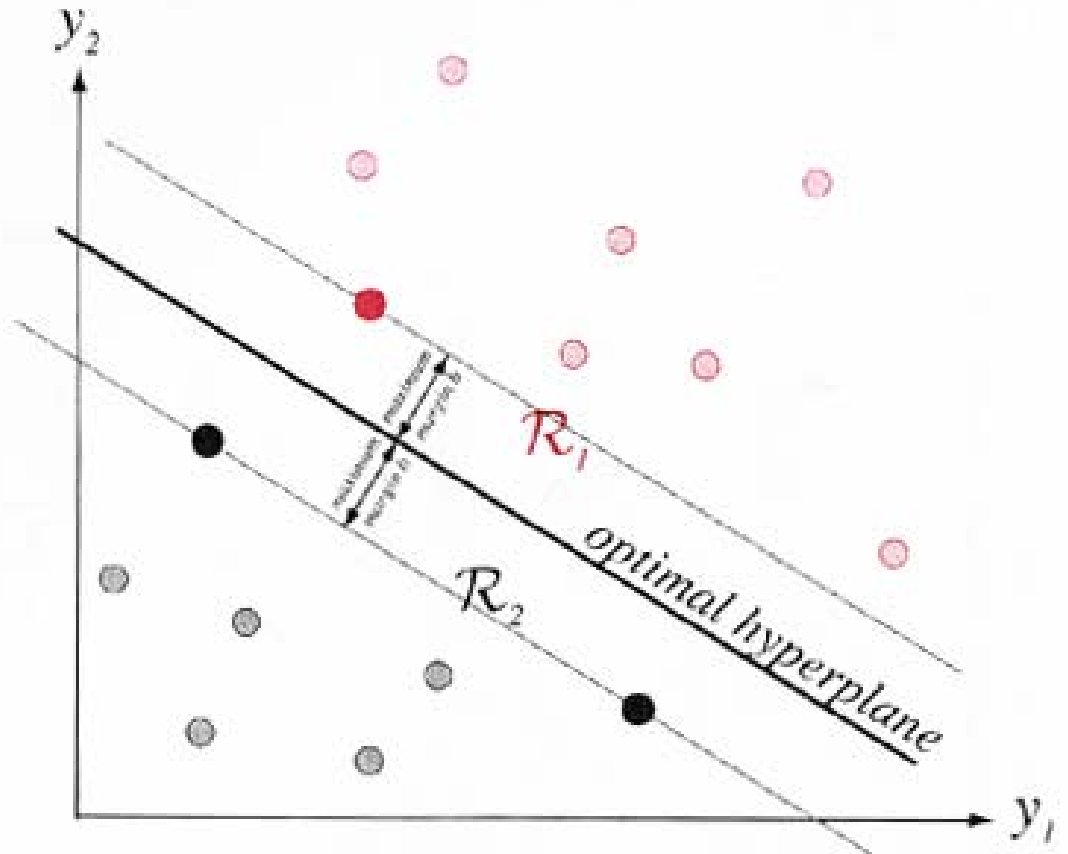
$$\varphi(x) = (1, x, x^2)^t$$

- First choose the non-linear phi functions
 - To map input vector to a higher dimensional feature space
- Dimensionality of space can be arbitrarily high only limited by computational resources
- Each pattern x_k transformed into pattern y_k
- where

$$\mathbf{y}_k = \Phi(\mathbf{x}_k)$$

Support Vectors

- Support vectors are nearest patterns at distance b from hyperplane
- $b =$ maximum distance from nearest training patterns



Three support vectors are shown as solid dots

Discriminative Classifier applications

- Image and document classification
- Biosequence analysis
- Time series prediction

Disadvantages of Discriminative Classifiers

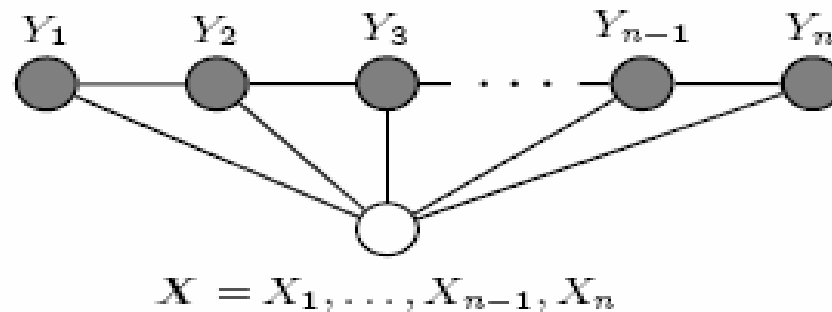
- Lack elegance of generative
 - Priors, structure, uncertainty
- Alternative notions of penalty functions, regularization, kernel functions
- Feel like black-boxes
 - Relationships between variables are not explicit and visualizable

Bridging Generative and Discriminative

- Can performance of SVMs be combined elegantly with flexible Bayesian statistics?
- Maximum Entropy Discrimination marries both methods
 - Solve over a distribution of parameters (a distribution over solutions)

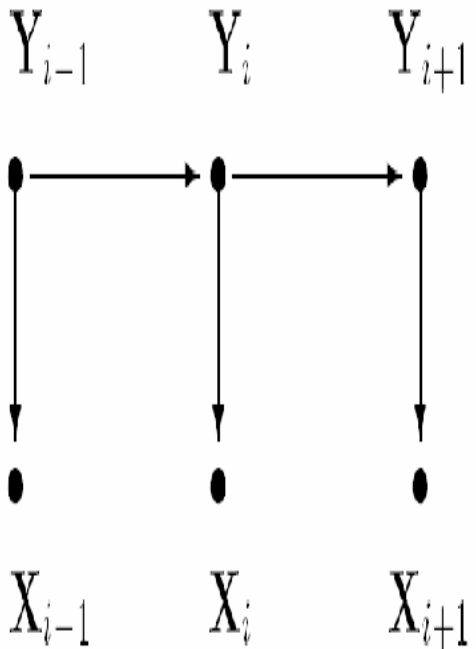
Conditional Random Fields

A CRF is the conditional distribution $p(Y|X)$ with the graphical structure:

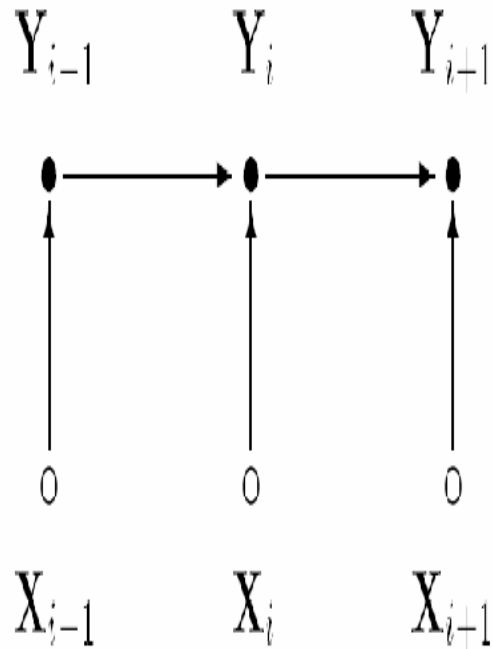


\mathbf{X} is the observed data sequence which is to be labeled
 \mathbf{Y} is the random variable over the label sequences

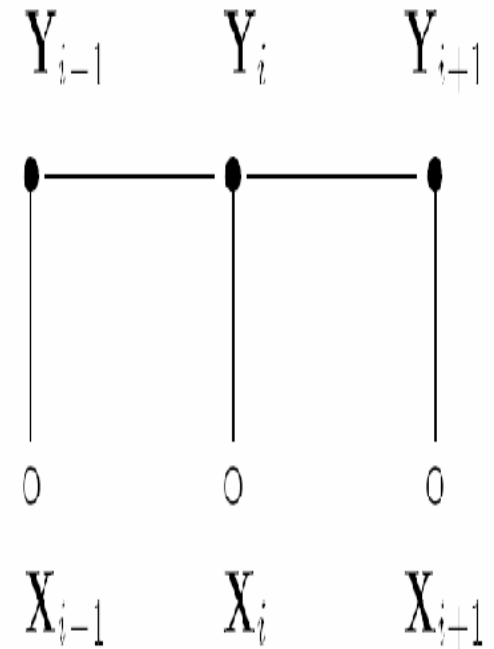
Graphical structures of HMMs, MEMMs and CRFs



Hidden Markov Model(HMM)



Max Entropy Markov Model (MEMM)



Conditional Random Field(CRF)

Advantage over Generative models

- CRF models conditional probability $P(\text{label } Y \mid \text{observed } X)$ as compared to generative models like HMM which model joint probability $P(Y, X)$
- Relaxes assumption of conditional independence of observed data given the labels
- It can contain a number of arbitrary feature functions
- Each feature function can use entire input data sequence. Probability of label at observed data segment may depend on any past or future data segments.

Advantage over other Discriminative Models

- CRF avoids limitation of MEMMs and other discriminative Markov models which can be biased towards states with few successor states.
- CRF has a single exponential model for joint probability of entire sequence of labels given the observation sequence.
- In MEMMs, $P(\mathbf{y} | \mathbf{x}) = \text{product of factors, one for each label.}$
Each factor depends only on previous label, and not future labels

Applications in Document Analysis

- CRFs can be used in sequence labeling tasks
- Zone Labeling –
 - Signature Extraction, Noise Removal
- Pixel Labeling –
 - Binarization of documents
- Character level labeling
 - Recognition of Handwritten Words

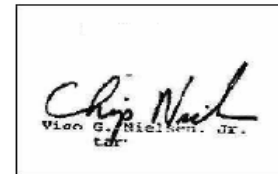
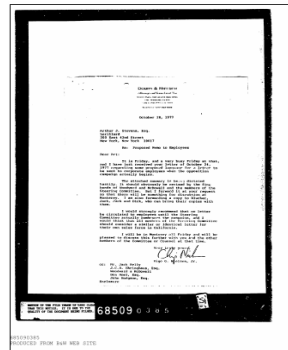
Applications

- Recognition of Handwritten Words
- Labeling scanned documents
- Binarization of documents
- Document Image Retrieval

Document Image Retrieval

- Signature Extraction

Original Document



Extracted Signature

- Signature Retrieval

RETRIEVAL

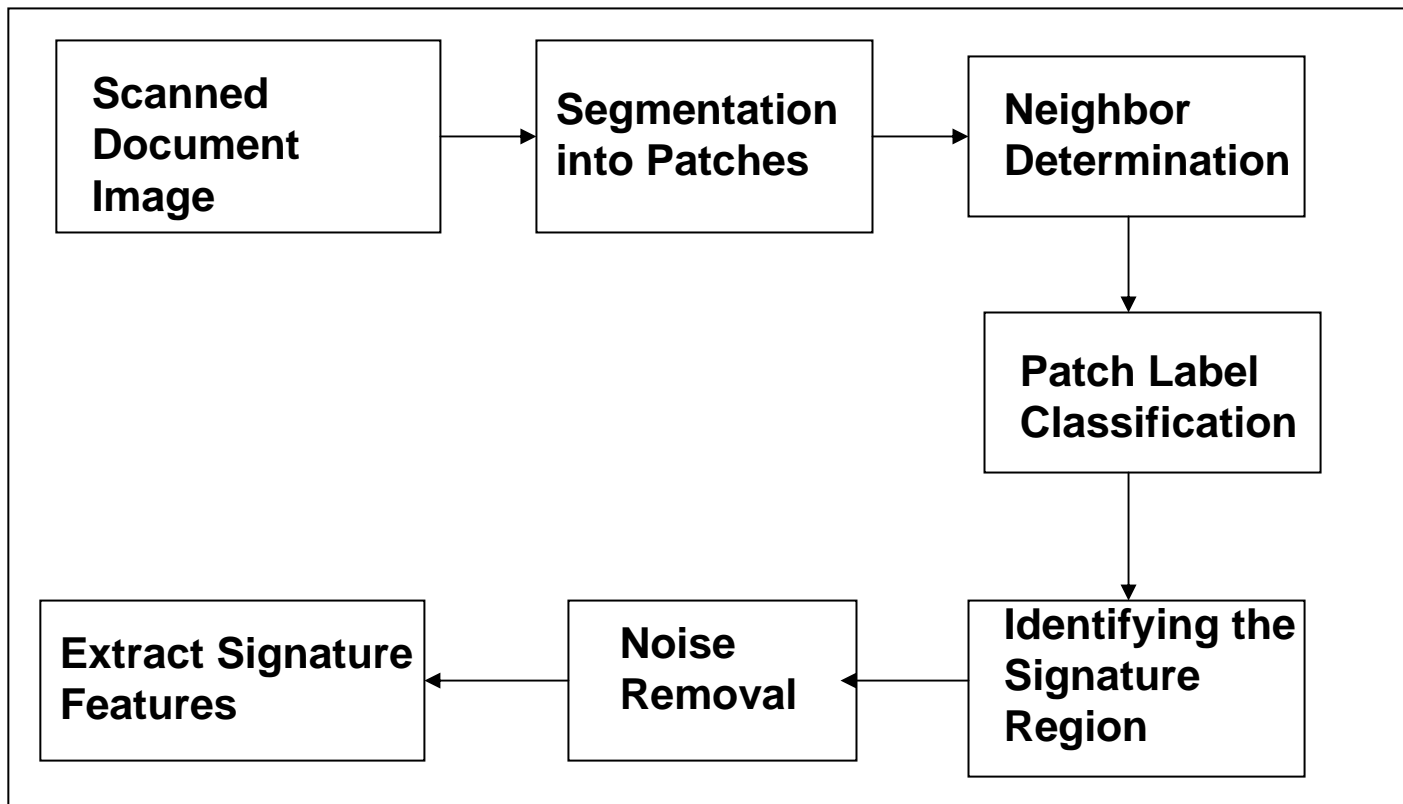
	Discriminative
Very truly yours,  Vigo G. Nielsen, Jr.	0.25
Yours sincerely,  Manningford, Esq.	0.36
	0.39
	0.43
	0.49

QUERY IMAGE



Chip Naal
Vigo G. Nielsen, Jr.
Esq.

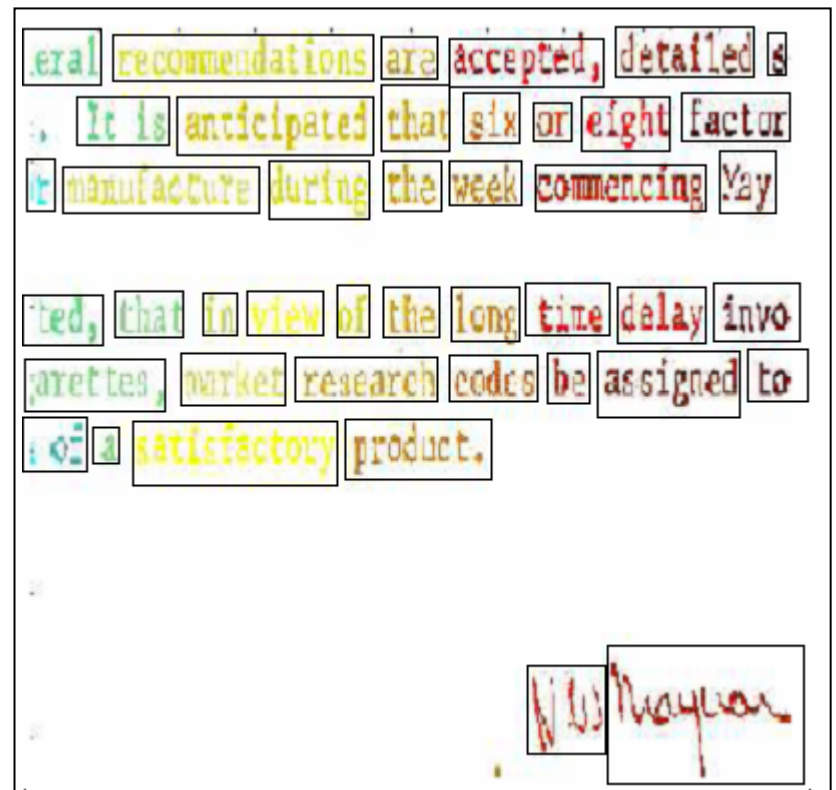
Signature Extraction



A Patch is defined to be a group of connected components

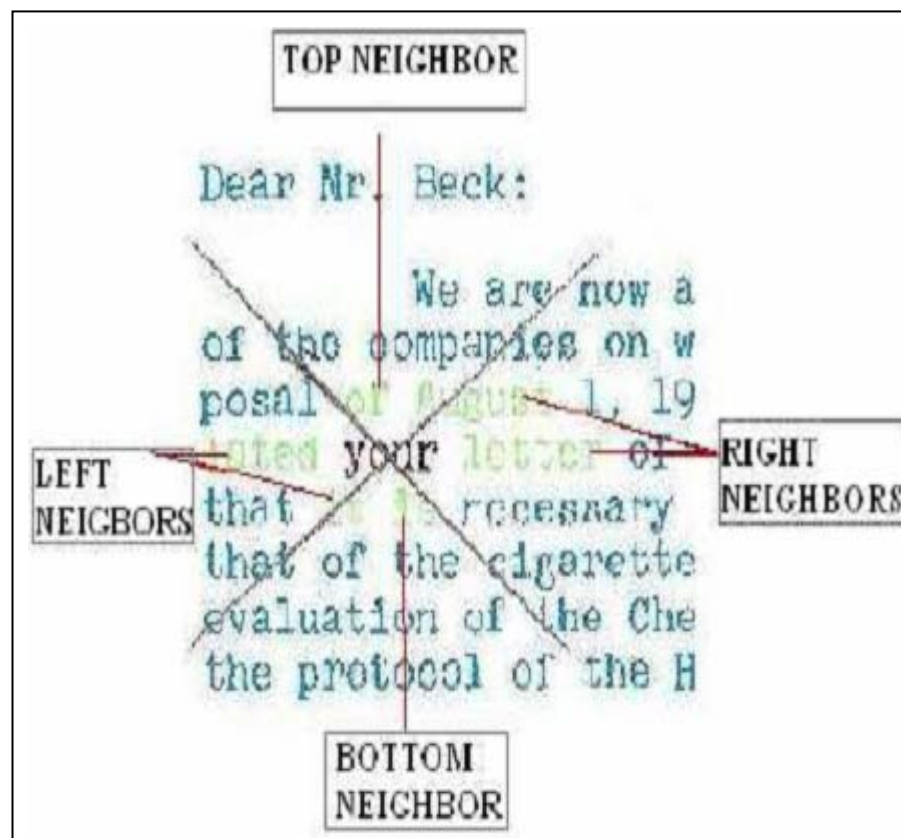
Segmentation into Patches

- Patches generated using region growing algorithm
- Size of patch optimized to represent size of word
- Region within each rectangle represents a patch



Neighboring Patch Detection

- Six neighbors are identified for each patch
- Closest(top/bottom) and two closest(left/right) in between patches are considered neighbors.



CRF for Zone Labeling

The probabilistic model of the Conditional Random Field used is

$$P(\mathbf{y}|\mathbf{x}, \theta) = \frac{e^{\psi(\mathbf{y}, \mathbf{x}; \theta)}}{\sum_{\mathbf{y}'} e^{\psi(\mathbf{y}', \mathbf{x}; \theta)}} \quad \begin{array}{l} \text{where } \mathbf{y}_i \in \{\text{Machine-print, Handwriting, Noise}\} \\ \theta : \text{CRF model parameters} \\ \mathbf{x} : \text{Observed document} \end{array}$$

$$\psi(\mathbf{y}, \mathbf{x}; \theta) = \sum_{j=1}^m \left(A(j, y_j, \mathbf{x}; \theta^s) + \sum_{(j,k) \in E} I(j, k, y_j, y_k, \mathbf{x}; \theta^t) \right)$$

where $A(j, y_j, \mathbf{x}; \theta^s)$: association potential

it associates the characteristics of that patch with its corresponding label

$I(j, k, y_j, y_k, \mathbf{x}; \theta^t)$: interaction potential

it captures the neighbor/contextual dependencies by associating pair wise interaction of the neighboring labels and the observed data

m : non-overlapping patches

CRF Parameter Estimation and Inference

- Parameter estimation
 - Done by maximizing pseudo-likelihood parameters using conjugate gradient descent with line search optimization

$$\theta_{ML}^{\hat{}} \approx \arg \max_{\theta} \prod_{i=1}^M P(y_i | y_{N_i}, \mathbf{x}, \theta)$$

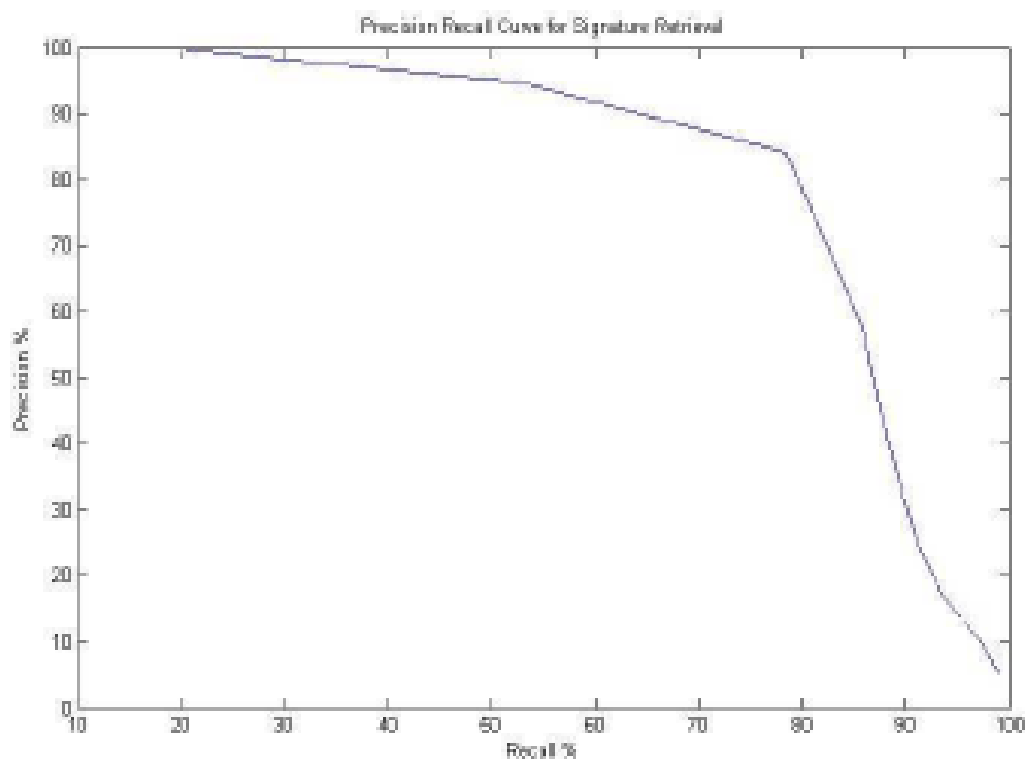
- Inference
 - Labels are assigned to each patch using Gibbs Sampling
 - Initially random labels are assigned to each patch
 - Over several iterations the label for each patch is generated from a distribution conditioned on the current values of its neighbors.

Features for HW/Print/Noise Classification

State Feature	Description
Height	Maximum height of the patch
Avg component width	The mean width of the connected components within a patch
Density	Density of foreground pixels within the patch
Aspect ratio	Width/Height of the patch
Gabor filter	8 features capturing the different stroke orientations
Variation of height	Variation in height within a patch
Width variation	Variation in width within a patch
Overlap	Sum of overlap in area between the connected components within a patch
Percentage of text above	Relative location of the patch with respect to the entire document
Number of components	Count of the connected components within a patch
Maximum component size	Maximum size of a component within a patch
Points in convex hull	Number of points in the convex hull of the patch
Maximum run length	The maximum horizontal run length within a patch
Avg run length	The average horizontal run length within a patch
Horizontal Transitions	A count of the number of times the pixel value transitions from white to black horizontally
Vertical Transitions	A count of the number of times the pixel value transitions from white to black vertically

Transition Feature	Description
Relative location	Assigned weights based on the relative location - top/bottom or right/left
Convex hull distance	The convex hull distance between the 2 patches
Ratio of aspect ratio	The ratio of the aspect ratio values of the 2 patches
Ratio of number of components	The ratio of the number of components present in the 2 patches

Precision-Recall Curves for signature retrieval results



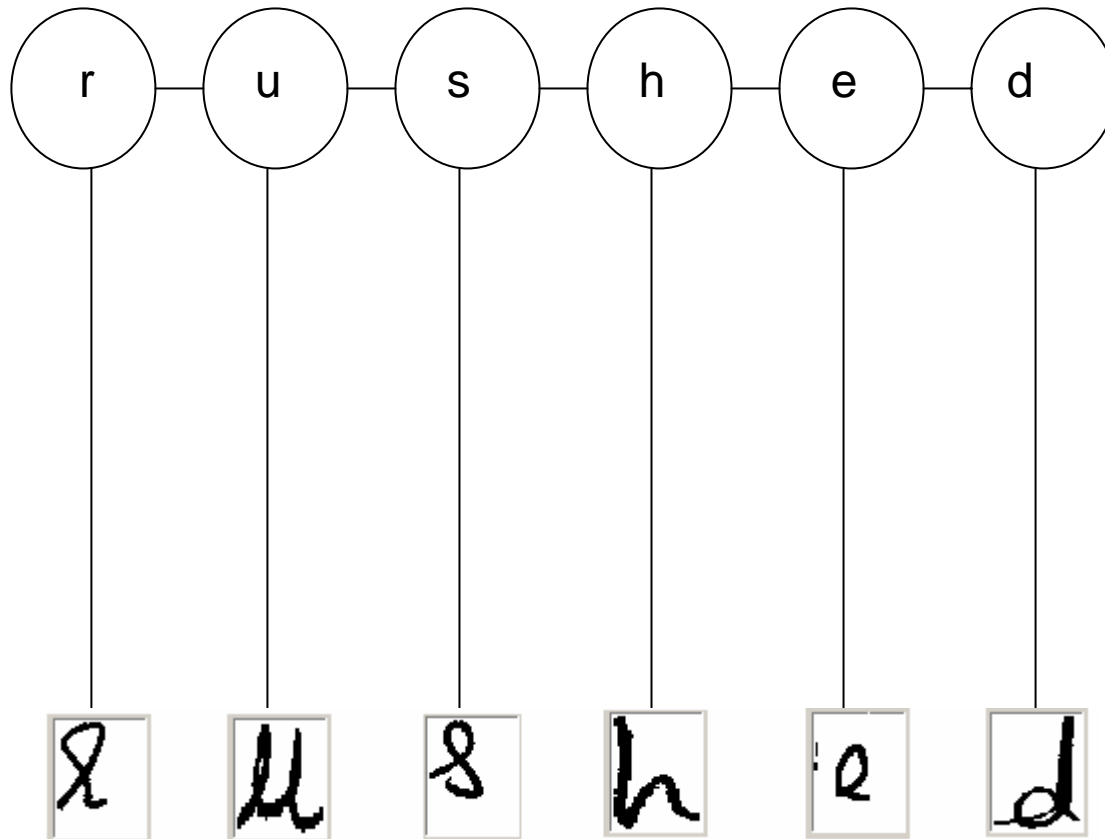
Precision of 84.2% at Recall of 78.4% after Query Expansion

Word Recognition Objective

- To transform the Image of a Handwritten Word to text using a pre-specified lexicon
 - Accuracy depends on lexicon size



Graphical Model of CRF for Word Recognition



Segmentation Based Word Recognition

- Word image is divided into several segmentation points.
- Dynamic programming used to find best grouping of segments into characters

$$P(\mathbf{y}|\mathbf{x}, \theta) = \max_i (P(\mathbf{y}|\mathbf{x}, s_i)P(s_i))$$

y is the text of the word, x is observed handwritten word, s is a grouping of segmentation points

CRF Model

Probability of recognizing a handwritten word image, X as the word 'the' is given by

$$P(y = 'the' | \mathbf{x}, \theta) = \frac{e^{\psi(y='t', \mathbf{x}; \theta)}}{\sum_{y'='a'}^{z'} e^{\psi(y', \mathbf{x}; \theta)}} \frac{e^{\psi(y='h', y_n='t', \mathbf{x}; \theta)}}{\sum_{y'='a'}^{z'} e^{\psi(y', y_n='t', \mathbf{x}; \theta)}} \frac{e^{\psi(y='e', y_n='h', \mathbf{x}; \theta)}}{\sum_{y'='a'}^{z'} e^{\psi(y', y_n='h', \mathbf{x}; \theta)}}$$

$$\psi(y, y_n \mathbf{x}; \theta) = A(j, y, \mathbf{x}; \theta^s) + I(j, k, y, y_n, \mathbf{x}; \theta^t)$$

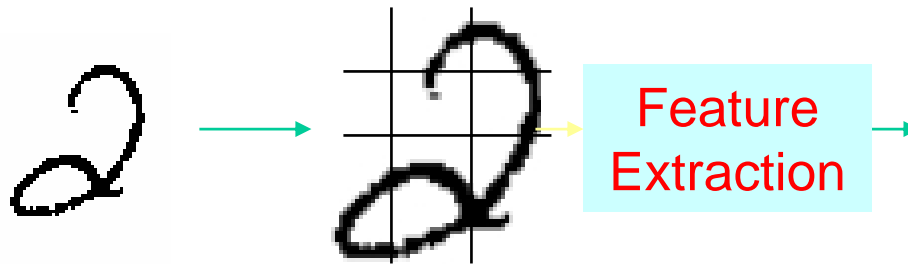
Captures the state features for a character

Captures the transition features between a character and its preceding character in the word

Word Features

- State features include
 - 74 WMR features
 - Distance between observed character segment and a codebook of characters
 - Height, width, aspect ratio, position in the text, etc
- Transition features include
 - Vertical overlap
 - Total width of the bigram
 - Difference in the height, width, aspect ratio, etc

Word Recognition: WMR features



2 Global Features: Aspect ratio
Stroke ratio

72 Local Features

$$F(i,j) = \frac{s(i,j)}{N(i)*S(j)} \quad \begin{matrix} i = 1,9 \\ j = 0,7 \end{matrix}$$

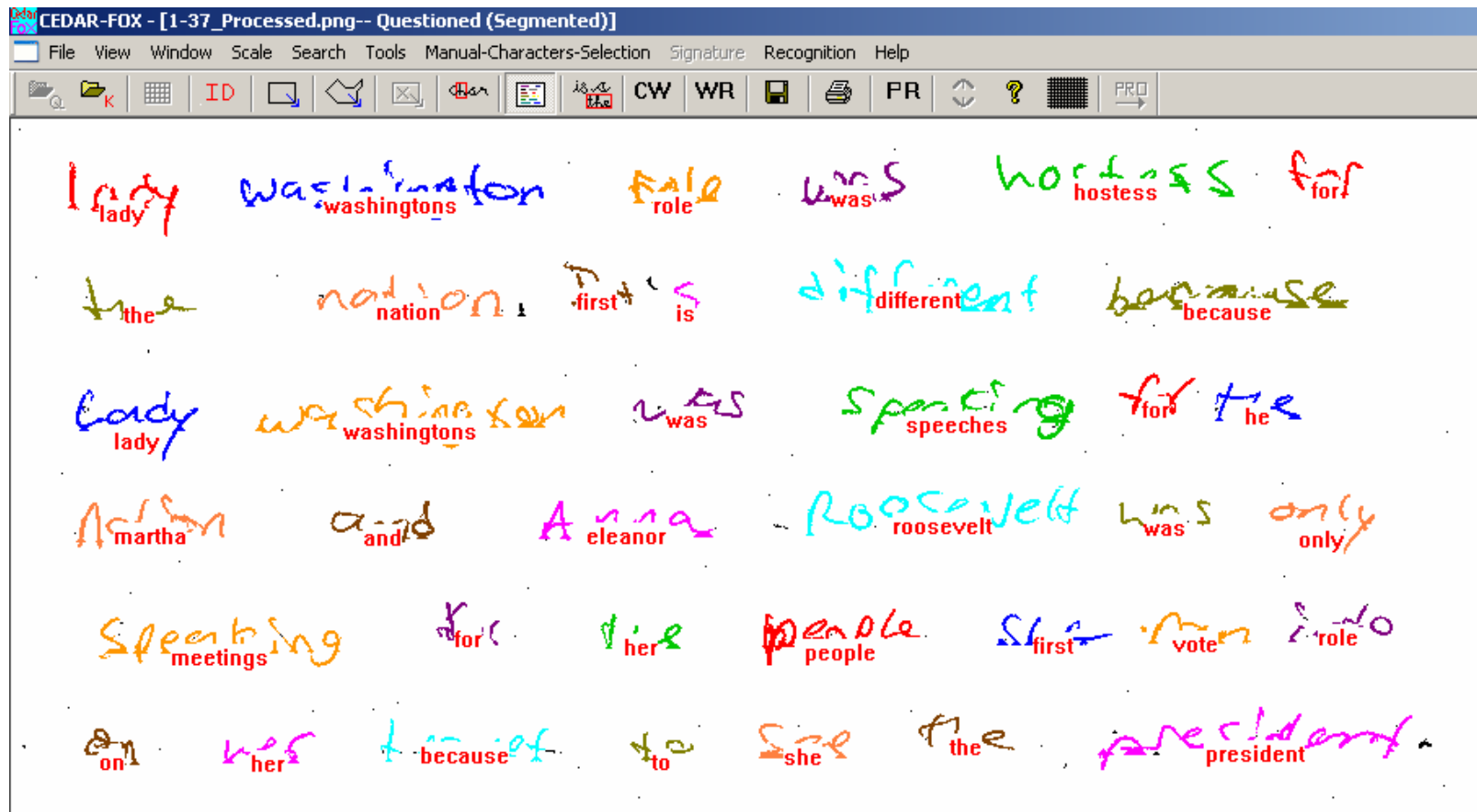
$s(i,j)$ - no. of components with slope j in sub image i

$N(i)$ - no. of components in sub image i

$$S(j) = \max \frac{s(i,j)}{N(i)}$$

0.000000 0.000000 0.000000 0.000000 0.000000 0.000000 0.000000 0.000000	0.076923 0.000000 0.529915 0.993590 0.387363 0.000000 0.596154 1.000000	0.175000 0.962500 0.688889 0.129167 0.419643 1.000000 0.387500 0.000000
1.000000 0.000000 0.000000 0.000000 0.000000 0.000000 0.000000 0.000000	0.727273 1.000000 0.000000 0.000000 0.000000 0.000000 0.000000 0.000000	0.000000 0.000000 1.000000 1.000000 0.108295 0.000000 1.000000 0.599078
0.207547 0.207547 0.194969 0.779874 0.760108 0.215633 0.350943 0.630728	0.127660 0.351064 0.073286 0.879433 1.000000 0.607903 0.263830 0.474164	0.111111 0.407407 0.637860 0.765432 0.870370 0.000000 0.114815 0.687831

Automatic Word Recognition



Summary

- Pattern Recognition and Machine Learning techniques are facets of same problem
- Many advances in past decade
- Play central role in IR, NLP, DAR, SR
- Classification algorithms are either generative or discriminative
 - former involve modeling, latter directly solve classification
- CRFs are a natural choice for several labeling tasks
 - Main disadvantage is cost of training
 - Good performance on labeling tasks in DAR