

CATALOG DESCRIPTION

ECSE 6610 Advanced Character Recognition. 4 credits.

Principles and practice of the recognition of isolated or connected typeset, hand-printed, and cursive characters. Review of optical scanners, features, classifiers. Supervised and non-supervised estimation of classifier parameters. Expectation Maximization, the Curse of Dimensionality, language context. Advanced classification techniques including Classifier Combinations, Support Vector Machines, Hidden Markov Methods, Adaptation, Indirect Symbolic Correlation.

Prerequisites: ECSE 2610, Probability, Linear Algebra.

Spring term annually

ECSE-6610 FIRST DAY HANDOUT

Instructor: Prof. George Nagy,
Office: JEC 6020, RPI, Troy, NY
Office hours: Hotel Bar, after class
Email: nagy@ecse.rpi.edu

Text: *Optical Character Recognition: An Illustrated Guide to the Frontier*
S. V. Rice, G. Nagy, T. A. Nartker, Kluwer Academic Publishers, 1999

Reference texts (on reserve at Folsom Library):

Duda, Hart, & Stork, 2001	[DHS 01]
Mitchell, McGraw-Hill 1997	[TM 97]
Nadler & Smith, Wiley 1993	[NS 93]
Schürmann, Wiley, 1996	[JS 96]
Theodoridis & Koutroumbas 1999	[TK 99]
Vapnik, Wiley 1998	[VV 98]

For additional sources, see the Text and the Bibliography.

Grading: Five programming assignments
Term Paper
Final Examination

Review: Intro to OCR (ECSE 2610)

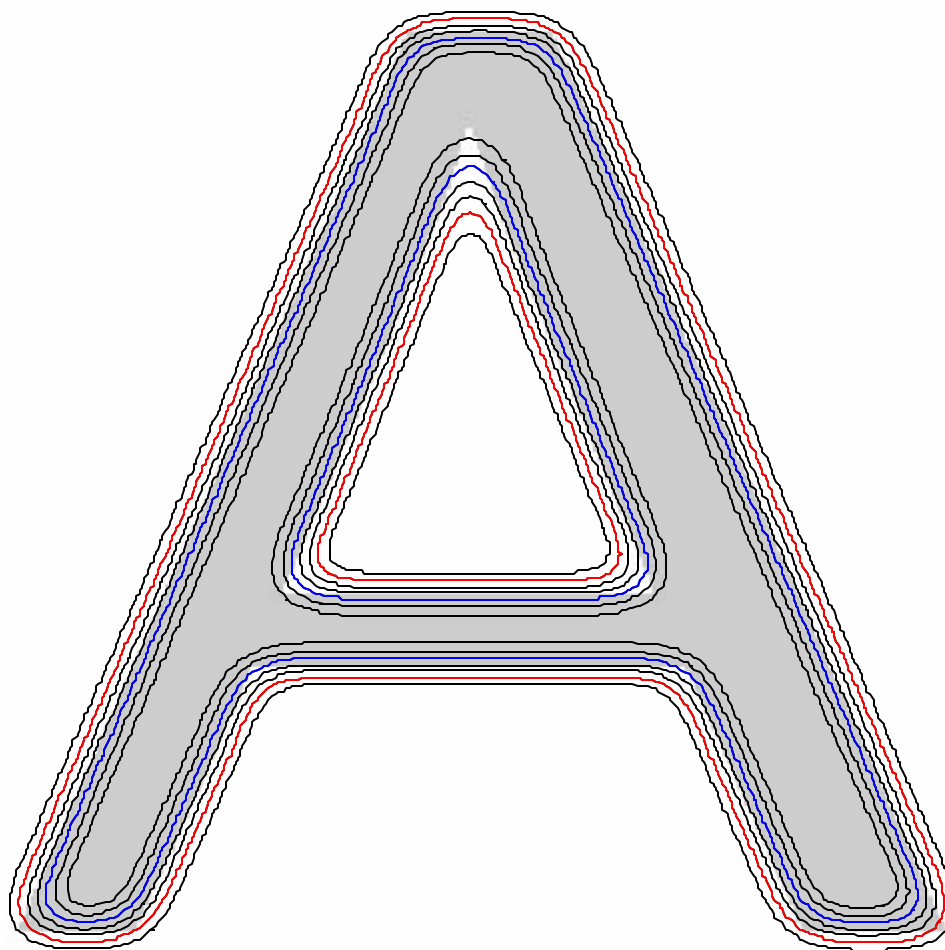
Digitization and preprocessing:

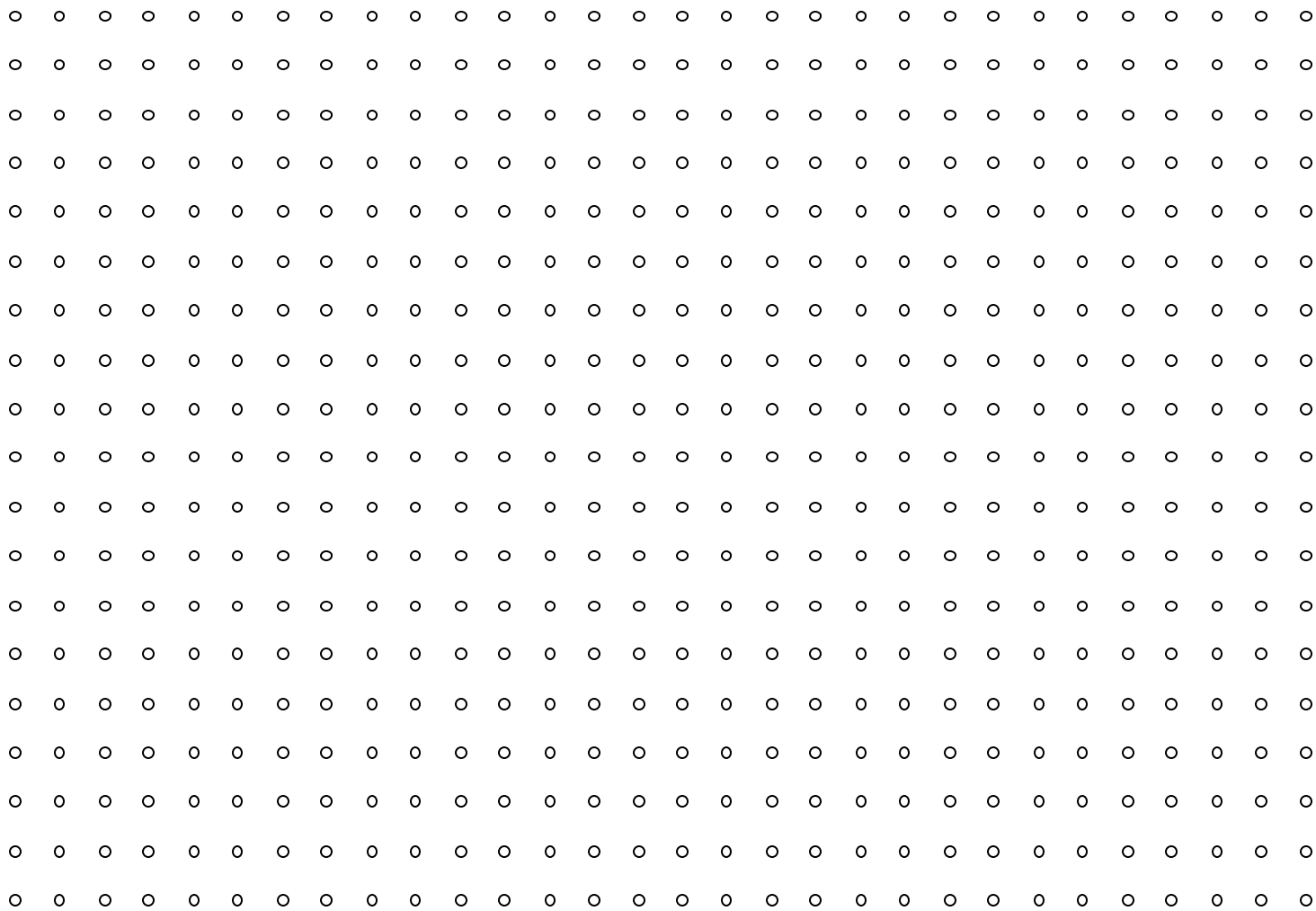
Scanner calibration; noise removal; recovery of scanner distortions [BE 01]

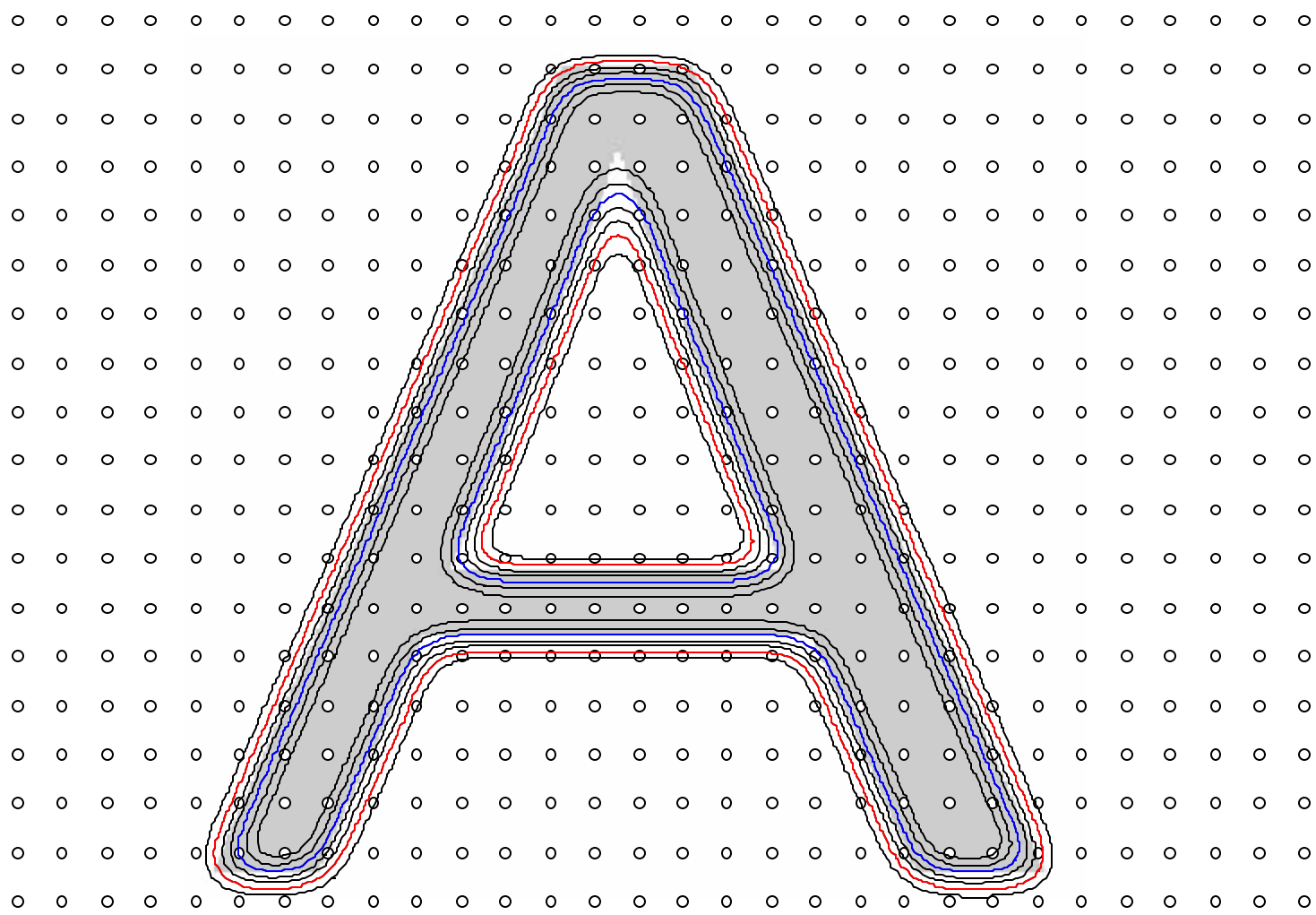
Character image defect models [KBH 94]

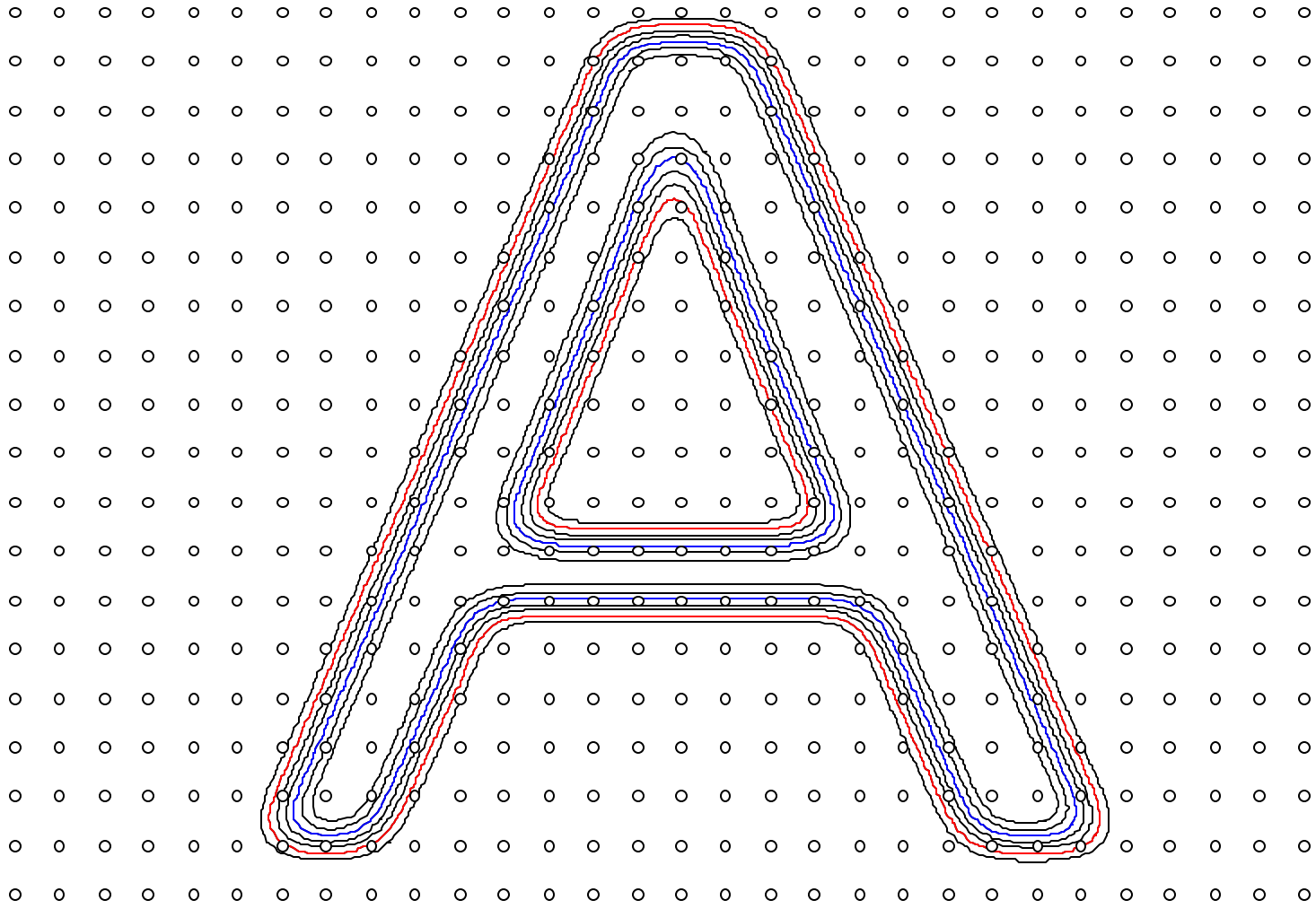
Help Session: Thursday 9:40 Prof. Elisa Barney Smith, BSU

MODEL OF DIGITIZATION









Review:

Text-figure separation; skew correction; text layout extraction
(column, line, and word segmentation)

[NG 00].

Help session: Wed. 11:00

Prof. Don Sylwester, Concordia College

Tables

Help session: Monday 13:15, 13:55

Dr. Dan Lopresti, Lucent Bell Labs

X-Y TREE

Advanced Character Recognition 6610
(invited)
George Nagy
Rensselaer Polytechnic Institute, Troy, NY, USA

CATALOG DESCRIPTION	SYLLABUS
<p><u>ECSE 6610 – Advanced Character Recognition. Principles and practice of the recognition of isolated or connected typeset, hand-printed, and cursive characters. Review of optical digitization, supervised and unsupervised estimation of classifier parameters, bias and variance, expectation maximization, the curse of dimensionality. Advanced classification techniques including classifier combinations, support vector machines, hidden Markov methods, styles, language context, adaptation, segmentation-free classifiers, indirect symbolic correlation. Prereq: ECSE 2610, Probability, Linear Algebra. Spring term annually.</u></p> <p style="text-align: center;">ECSE-6610 FIRST DAY HANDOUT</p> <p>Instructor: Prof. George Nagy Office hours: After class in the bar Email: nagy@ecse.rpi.edu</p> <p>Text: S. V. Rice, G. Nagy, T. A. Nartker Optical Character Recognition: An Illustrated Guide to the Frontier [RNN 99]</p> <p>Reference texts (on reserve at Folsom Libe): Duda, Hart, & Stork, Wiley 2001 [DHS 01] Mitchell, McGraw-Hill 1997 [MT 97] Nader & Smith, Wiley 1993 [NS 93] Schürmann, Wiley 1996 [SJ 96] Theodoridis & Koutroumbas, Acad.1999 [TK 99]</p> <p>For additional sources, see the Text and the Bibliography.</p>	<p>1. Review: Intro to OCR (ECSE 2610)</p> <p>Preprocessing: Scanner calibration, correction of scan distortions; noise removal; text-figure separation; skew correction; gray-scale and color, text layout extraction (column, line, and word segmentation) [NG 00]. Character image defect models [KBH 94] Recovery of scanner distortions [BE 00]. Help Session: Wed. 6 pm Prof. E. Barney Smith.</p> <p>Features: Reflectance, geometric, & topological invariants [FG 60, SM 61] Features as weak classifiers [KE 00] N-tuples and feature selection [JN 95, JDM 00, JKNS 96]</p> <p>Resource person: Dr. D-M Jung, Yahoo!</p> <p>Static single-pattern classifiers: <i>Bayes:</i> Single & Multimodal, Linear, Quadratic, Gaussian and Bilevel [DHS 01, LKF 01] <i>Neural Networks:</i> Backprop, LVQ, RBF [EC 95] <i>Support Vector Machine</i> [VV 98] <i>Nearest Neighbors</i> [DHS 01] <i>Decision Trees and Forests</i> [AGW 07, TH 98]</p> <p>Classifier training: Sample size and dimensionality [RJ 91] Bias and variance [GBD 92] Bagging, Boosting, Random Subspaces [JDM 00] Clustering [TK 99] Expectation Maximization [DLR 77, RW 84]</p>

Review:

Features:

Reflectance, geometric, & topological *invariants*

[GF 60], [SM 61]

Features as weak classifiers

[KE 00]

N-tuples

[JN 95]

Feature selection

[JDM 00]

Resource person: Dr. D-M Jung, Yahoo!

FEATURE INVARIANCES

Reflectance:



Geometry:

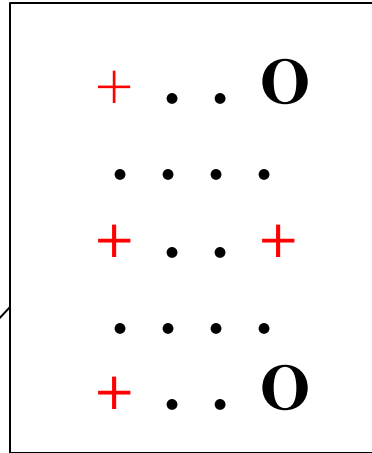


Topology:



N-TUPLE FEATURES (OR CLASSIFIERS?)

Bledsoe & Browning, 1959,, D-M Jung, 1995,



A B C D E F G H I
J K L M N O P Q R
S T U V W X Y Z

Review

Static Singleton Classifiers:

Bayes:

Single & Multimodal, Linear, Quadratic,
Gaussian and Bilevel

[DHS 01]

Neural Networks:

Backpropagation, Learning Vector Quantization,
Radial Basis Functions

[BC 95]

Support Vector Machines

[VV 98]

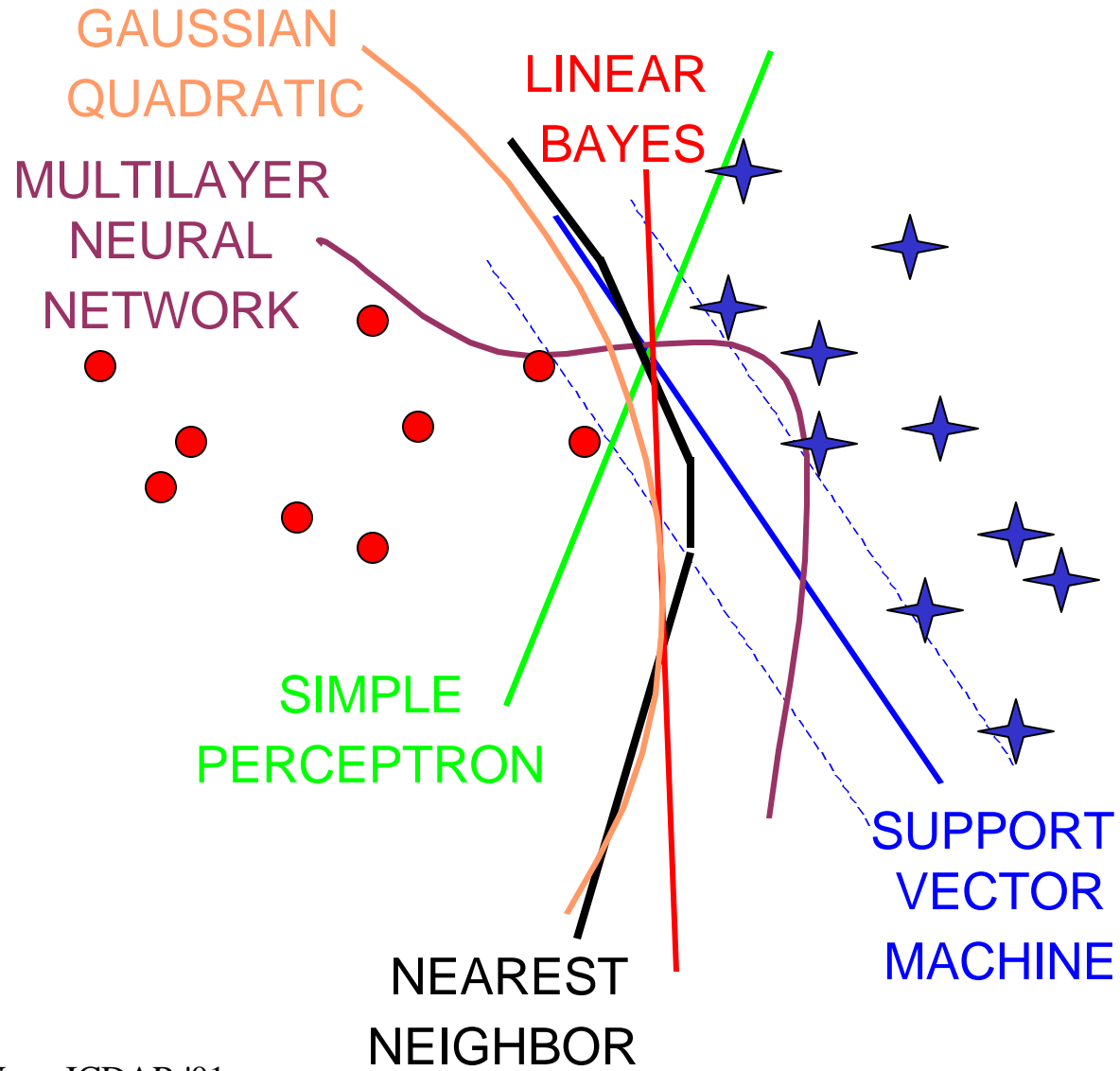
Nearest Neighbors

[DHS 01]

Decision Trees and Decision Forests

[TH 98]

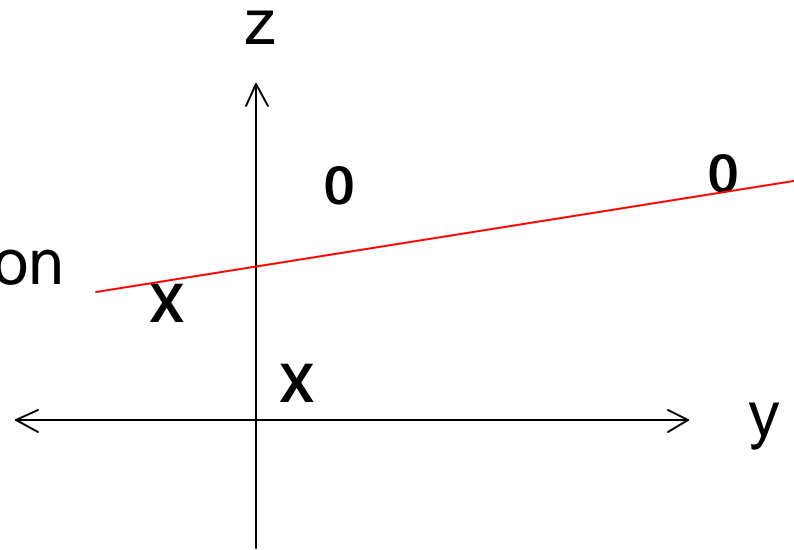
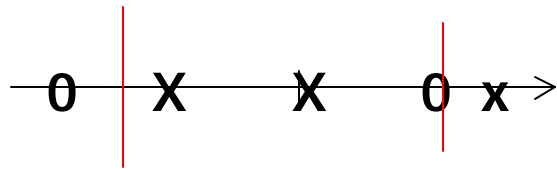
SOME CLASSIFIERS



SUPPORT VECTOR MACHINE (V. Vapnik)

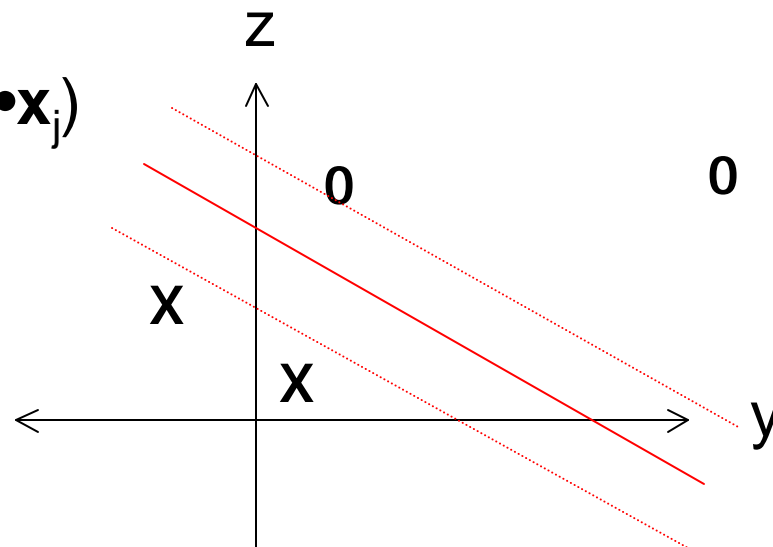
Kernel-induced transformation

$$x \rightarrow \mathbf{v} = (y, z)$$



max min $\{ f(\mathbf{v}_i \bullet \mathbf{v}_j) \}$ by QP
 Mercer's theorem: $\mathbf{v}_i \bullet \mathbf{v}_j = K(\mathbf{x}_i \bullet \mathbf{x}_j)$

i.e., compute distances in high-dim space from distances in low-dim space



Review:

Static Singleton Classifiers:

<i>Bayes: Single & Multimodal, Linear, Quadratic, Gaussian and Bilevel</i>	[DHS 01]
<i>Neural Networks: Backprop, LVQ, RBF, Support Vector Machine</i>	[BC 95]
<i>Nearest Neighbors</i>	[VV 98]
<i>Decision Trees and Decision Forests</i>	[DHS 01]
	[TH 98]

Classifier training:

Dimensionality and sample size	[RJ 91, LSF 01]]
Bias and variance	[GBD 92],
Bagging, Boosting, Random Subspaces	[JDM 00]
Clustering & Expectation Maximization	[TK 99, DLR 77, RW 84]

Generalization, Validation, and Error Prediction:

Validation, Jackknife, Bootstrap	[DHS 01], [JDM 00]
----------------------------------	--------------------

SOME DEFINITIONS

BOOTSTRAP PARAMETER ESTIMATE:

Mean of estimates of m sample sets obtained from the same data by *sampling with replacement*.

(Better than m -way partitioning for estimating variances.)

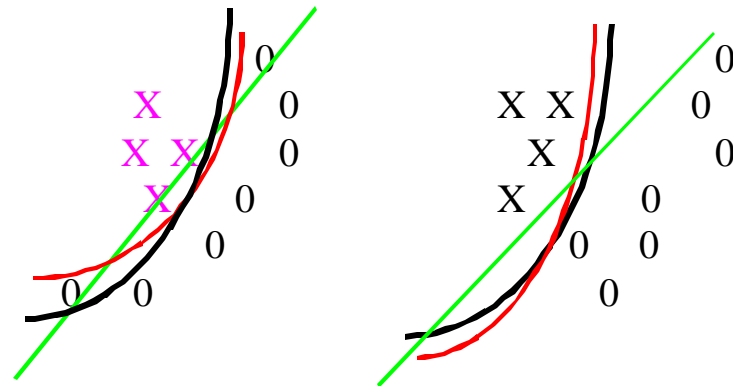
JACKKNIFE: leave-one-out estimate (of classifier accuracy).

BAGGING (Bootstrap AGgregation):

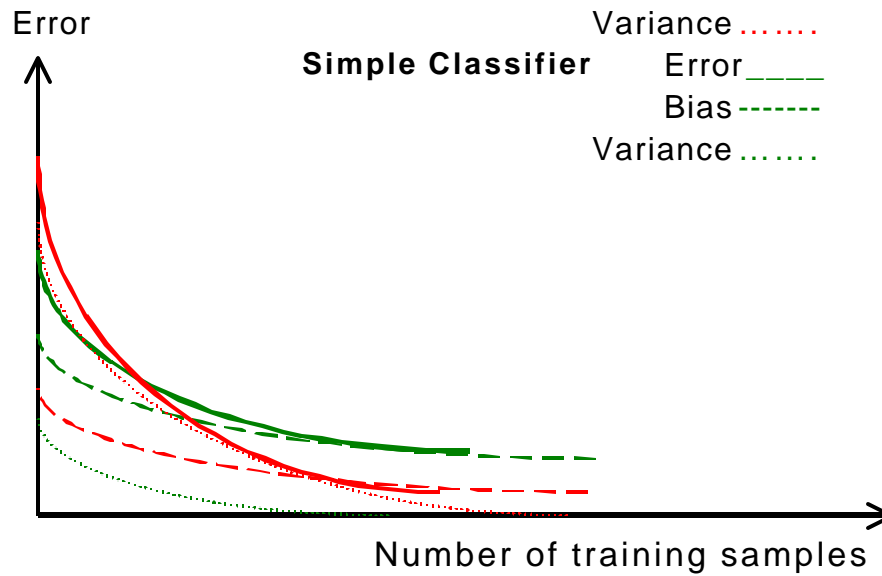
Increasing the nominal size of the training set by sampling with replacement.

BOOSTING: reintroducing misclassified training samples into the classifier construction process.

CLASSIFIER BIAS AND VARIANCE



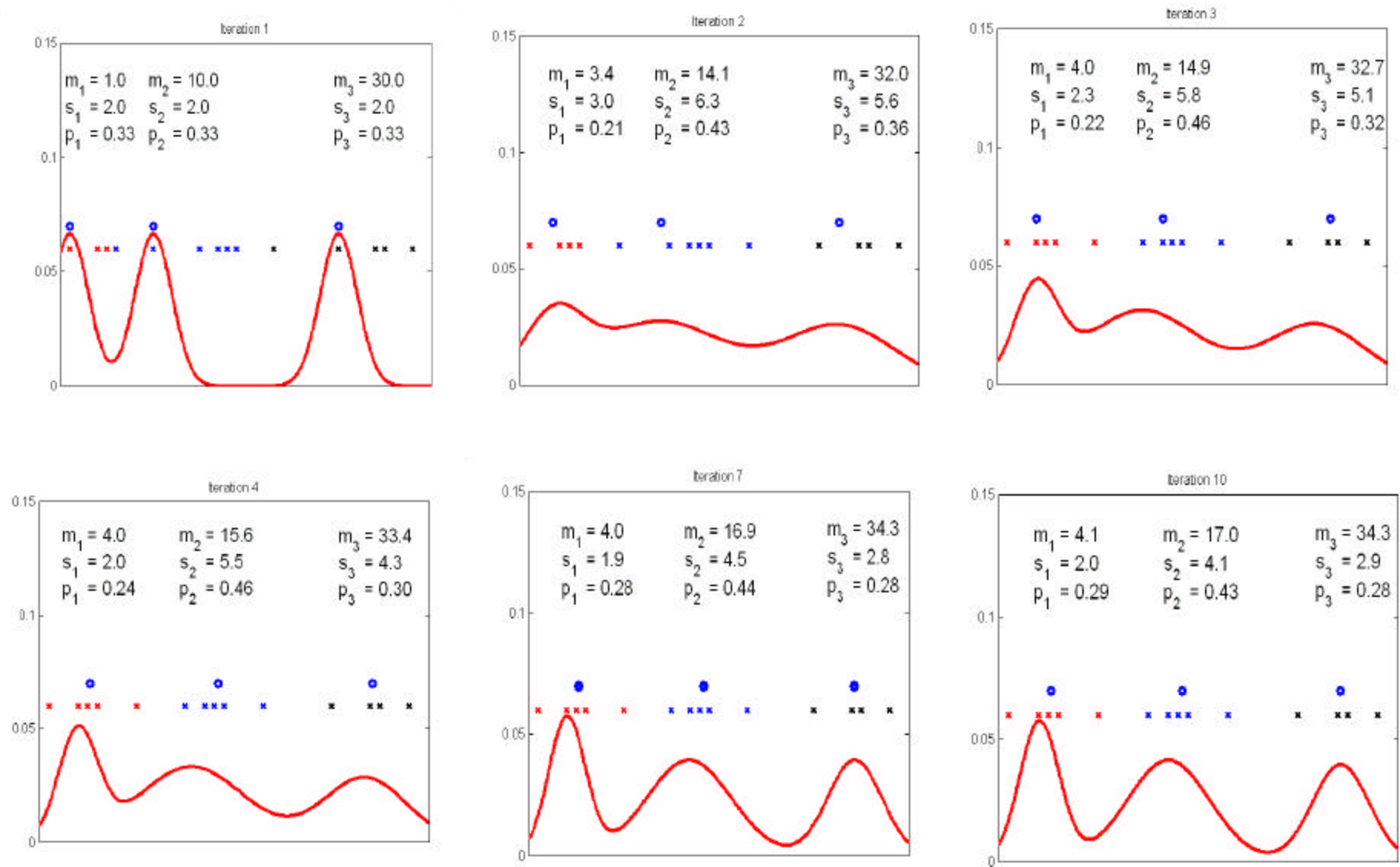
Complex Classifier Error _____
 Bias - - - - -
 Variance
Simple Classifier Error _____
 Bias - - - - -
 Variance



CLASSIFIER BIAS AND VARIANCE DON'T ADD!

Any classifier can be shown to be better than any other.

K-MEANS AND EXPECTATION MAXIMIZATION



OCR 6610 TOPICS

CLASSIFIER COMBINATION (Dr. Tin Kam Ho, Bell Labs)

CONTEXT (Drs. H. Fujisawa, J.J. Hull, Profs. S. Srihari, S. Seth)

STYLES (Dr. P. Sarkar, Harsha Veeramachaneni)

UNSUPERVISED ADAPTATION (Dr. H.S. Baird)

UNSEGMENTED TEXT

DOCUMENT-SPECIFIC CLASSIFIERS (Dr. Yihong Xu, EMC)

N-GRAM-BASED CLASSIFICATION (Adnan El-Nasan)

INDIRECT SYMBOLIC CORRELATION (Harsha V.)

CLASSIFIER COMBINATION [Tin Kam Ho '01]

DECISION OPTIMIZATION:

Combine a set of complete, fully trained classifiers.

COVERAGE OPTIMIZATION:

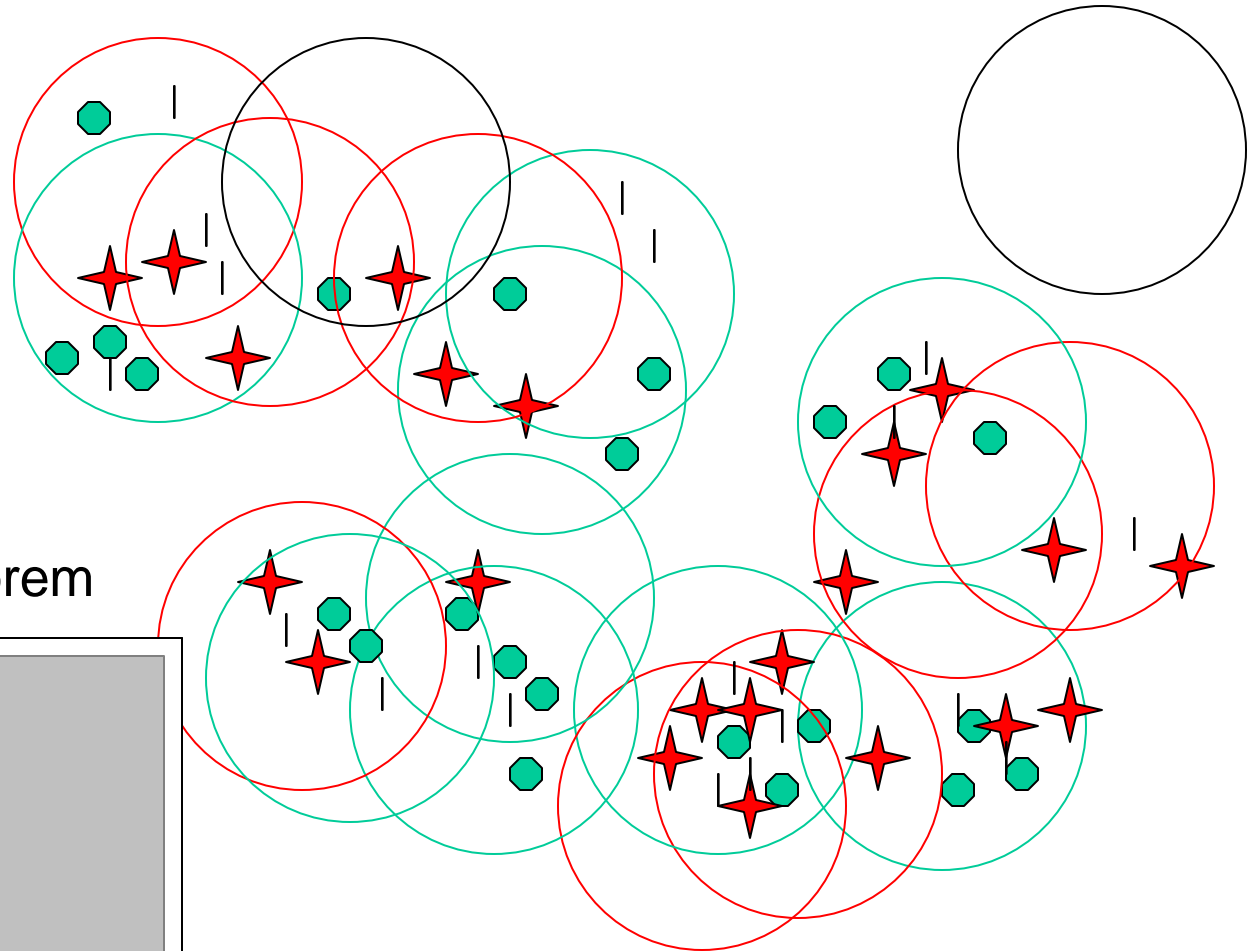
Tune each classifier to a different aspect of training set
(Different subspaces or training samples).

STOCHASTIC DISCRIMINATION [KE 00]:

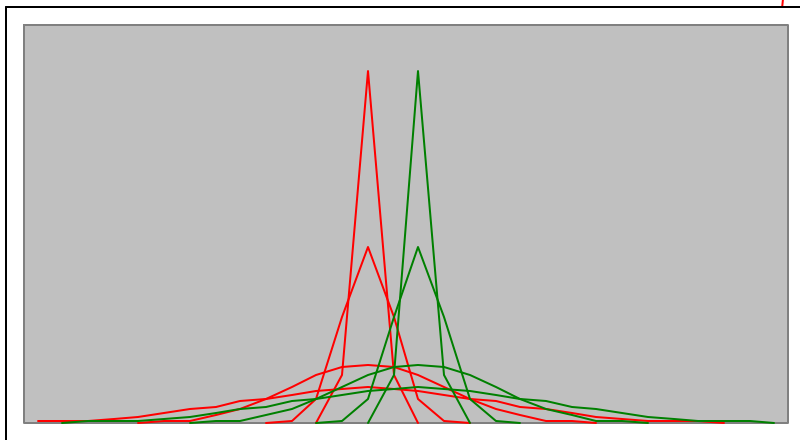
1. Enrichment: each of *many* weak “models” must favor some class;
2. Uniform coverage of populated region;
3. Characteristics of the sample space of weak models with respect to any point same whether that point is a training point or a test point.

STOCHASTIC DISCRIMINATION [Eugene Kleinberg '00]

Duality between sample space of weak models and feature space

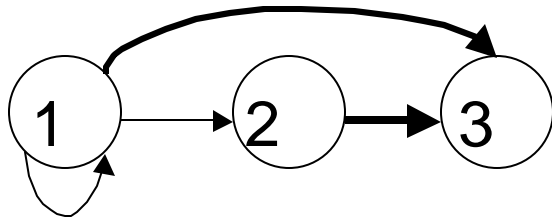


The Central Limit Theorem



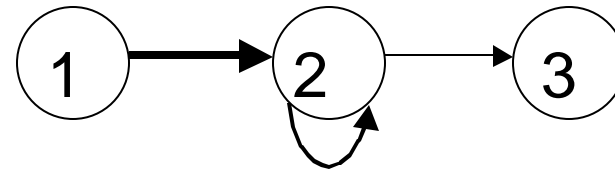
CONTEXT: HIDDEN MARKOV MODELS

MODEL A

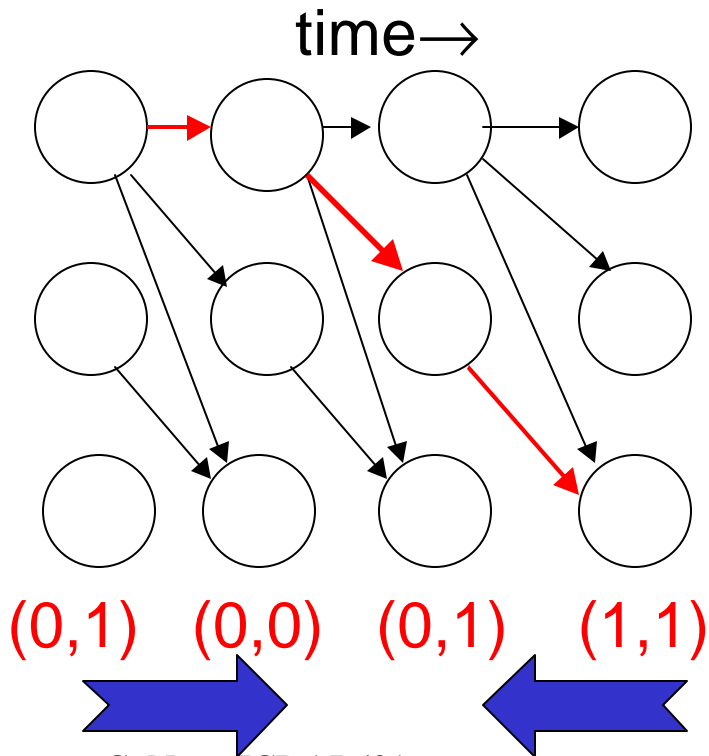


(0.2, 0.3) (0.3, 0.6) (0.7, 0.8)

MODEL B

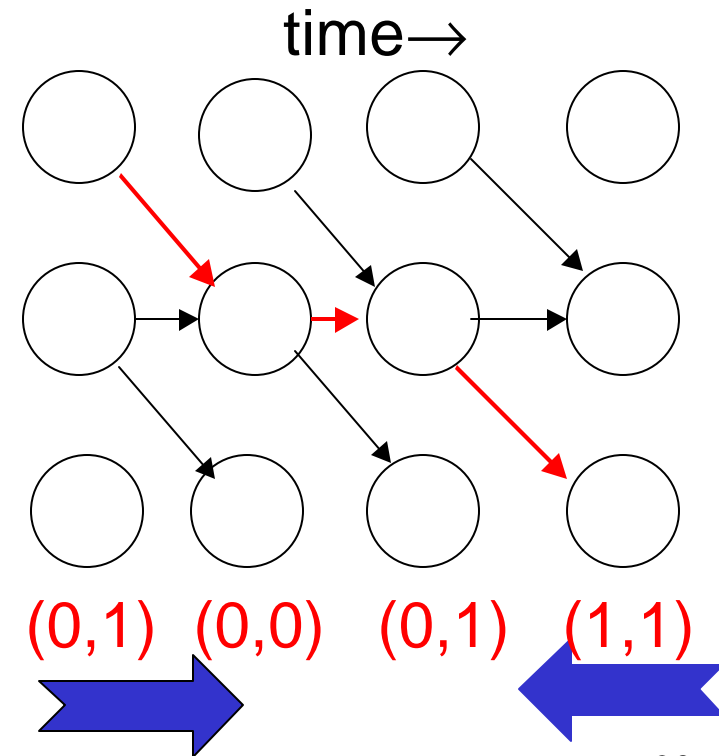


(0.3, 0.1) (0.5, 0.4) (0.4, 0.8)



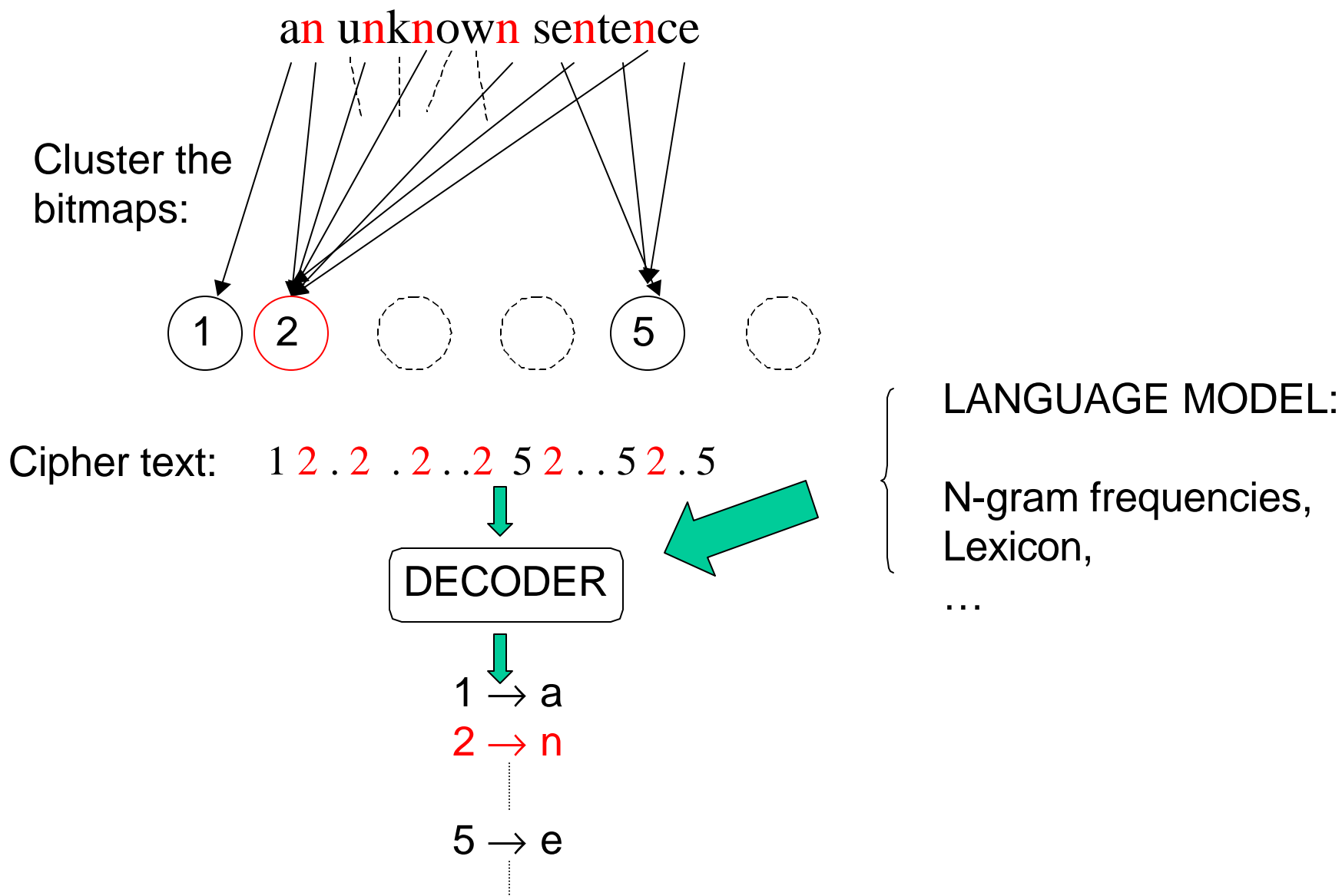
states

1
2
3



TRAINING: joint probs via Baum-Welch Forward-Backward (EM)

CONTEXT BY DECODING A SUBSTITUTION CIPHER



TEXT PRINTED WITH SPITZ GLYPHS

Spitz glyphs are a set of 256 characters used for text rendering in the ICDAR '01 competition. The image shows a sample of text printed using these glyphs, which appear as a dense, somewhat noisy pattern of characters. The text is arranged in approximately 12 lines, with each line containing a sequence of these special characters. The overall appearance is that of a corrupted or stylized text document.

DECODED TEXT

chapter i _ bee_ inds

_all me ishmaels some years ago__never mind
how long precisely ___having little or no money in
my purses and nothing particular to interest me on
shores i thought i would sail about a little and see
the watery part of the worlds it is a way i have ...

chapter I 2 LOOMINGS

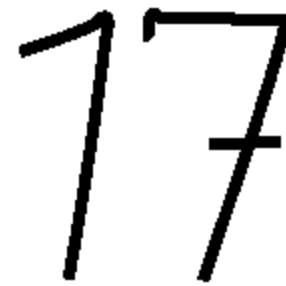
Call me Ishmael. Some years ago – never mind
how long precisely – having little or no money in
my purse, and nothing particular to interest me on
shore, I thought I would sail about a little and see
the watery part of the world. It is a way I have ...

STYLES

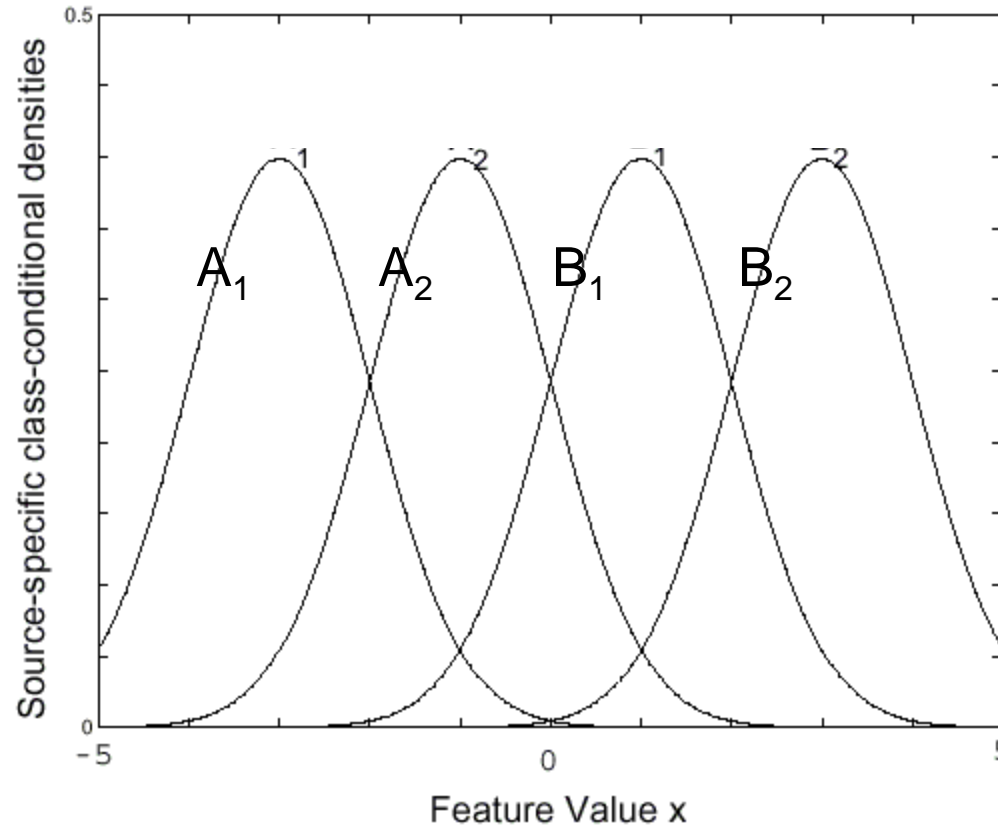
(Prateek Sarkar
Thursday 10:00,
Harsha Veeramachaneni)



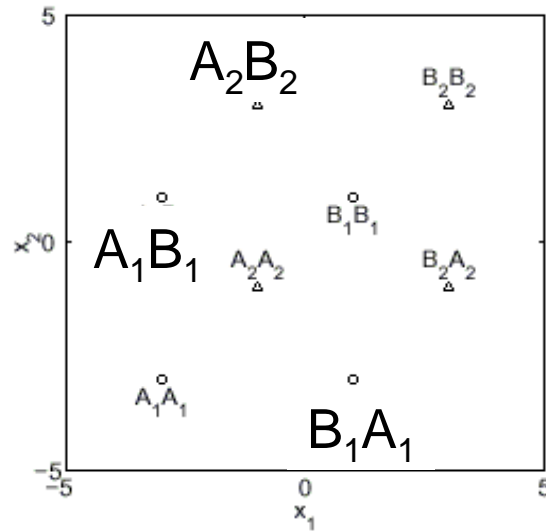
Writer 1



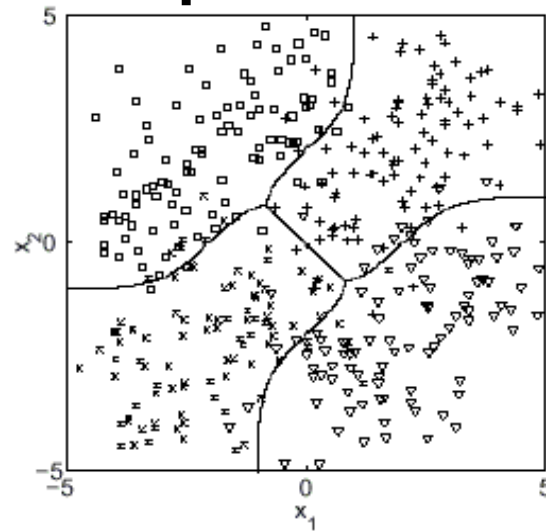
Writer 2



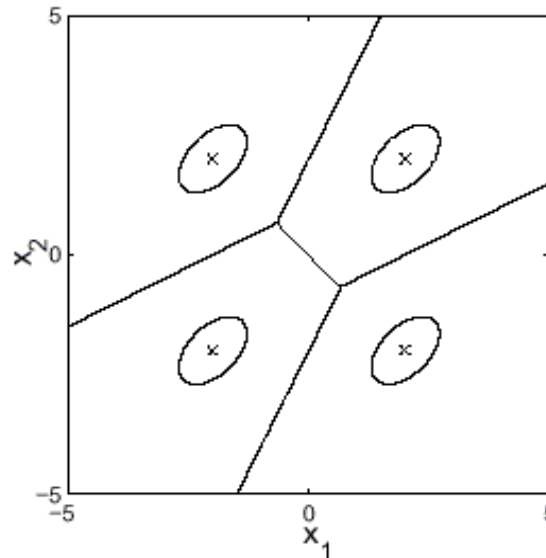
STYLE-CONSCIOUS CLASSIFIERS in field-feature space



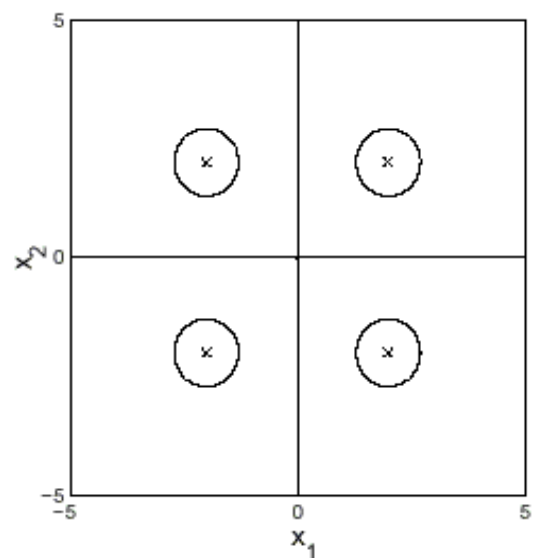
Position of means



Optimal field classifier



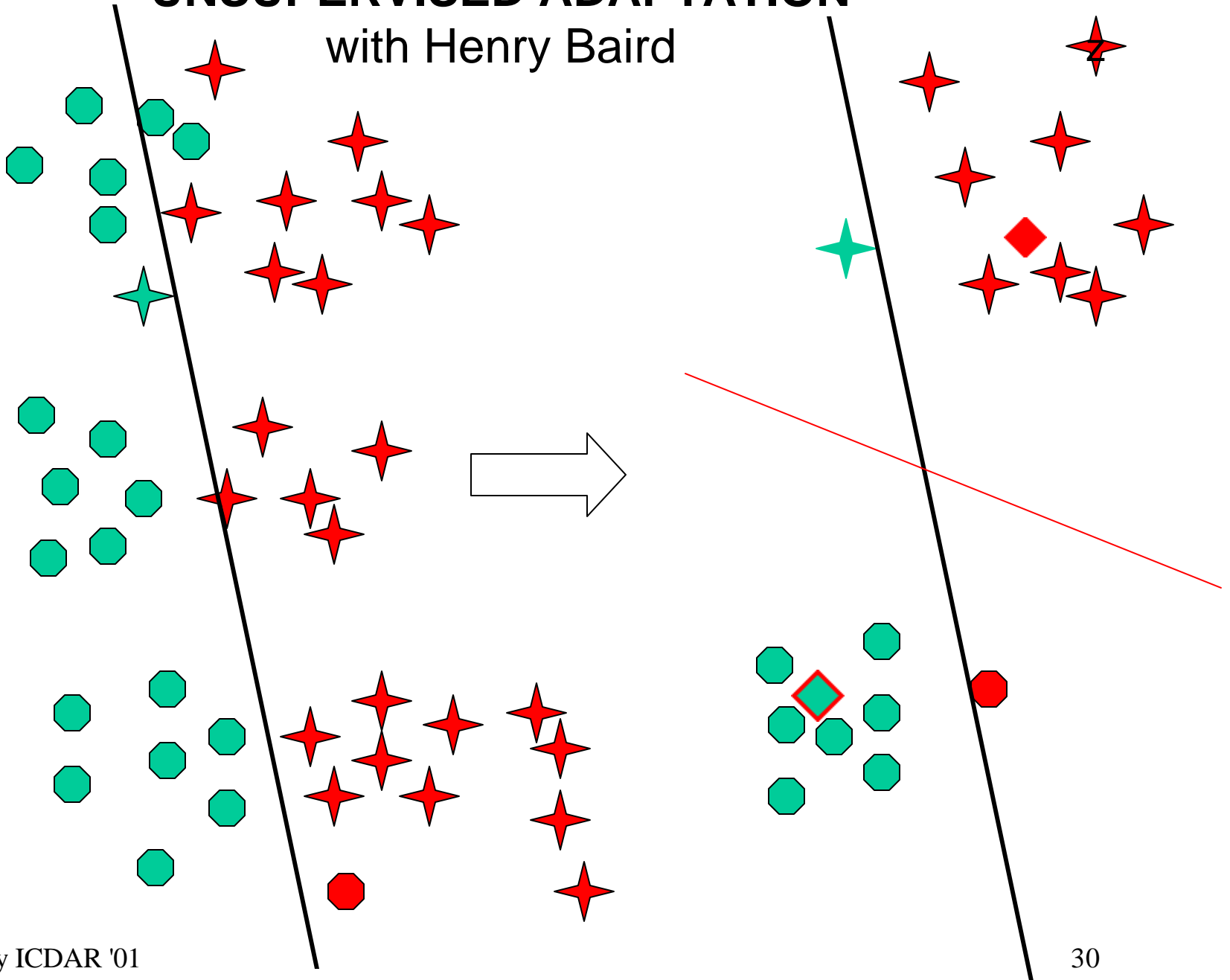
Quadratic field classifier



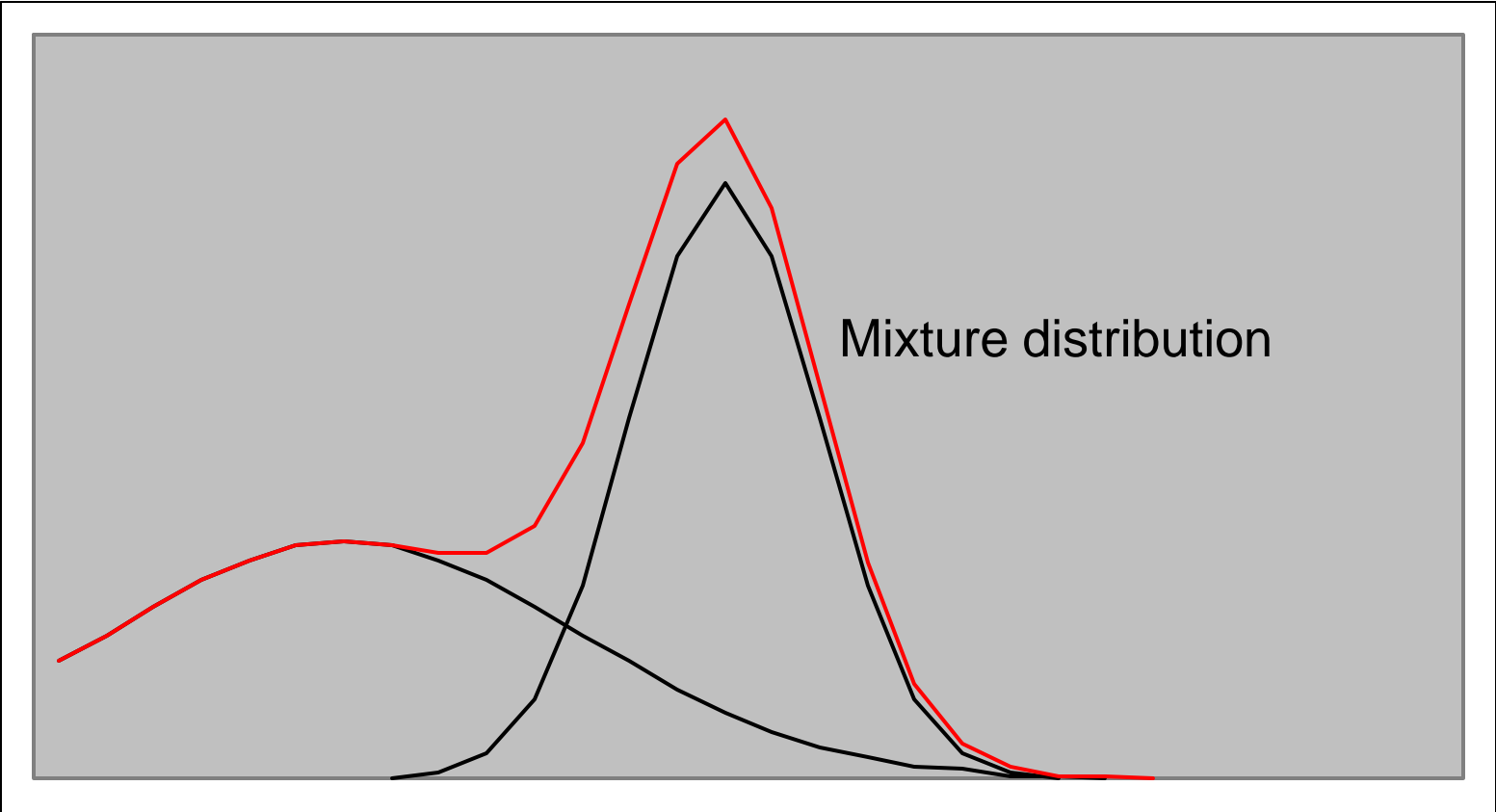
Singlet classifier

UNSUPERVISED ADAPTATION

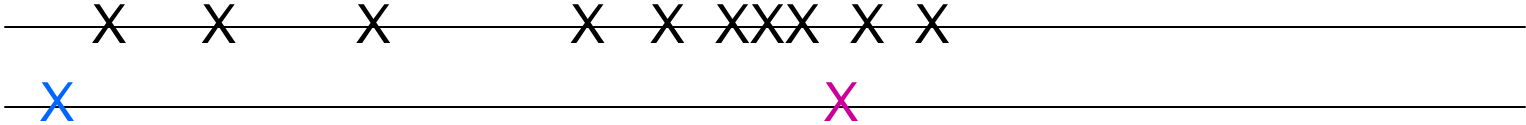
with Henry Baird



THE EXPONENTIAL VALUE OF LABELED SAMPLES (Tom COVER)



To estimate μ_s



Document-specific Prototype

SOURCE OF WASTE 11

NTS has no verification or processing facilities. All of the NTS waste has been retained and certified using the mobile NDA/NDE system developed by LAM. Any waste that is severable will be sent to demonstrate GCD technology at that site. The Nevada Test Site has no RH TRU waste.

CONCLUSION

The six TRU waste storage and generator sites in the DOE system are in various stages of the planning and implementation process for certifying waste to meet the WIPP WAC, so that waste shipments are acceptable for geologic disposal in the WIPP. Current plans will provide a timely stream of waste shipments to the WIPP, while assuring compliance with all DOE regulations and orders pertaining to TRU waste shipments.

WIPP WASTE ACCEPTANCE CRITERIA (WAC) SUMMARY*

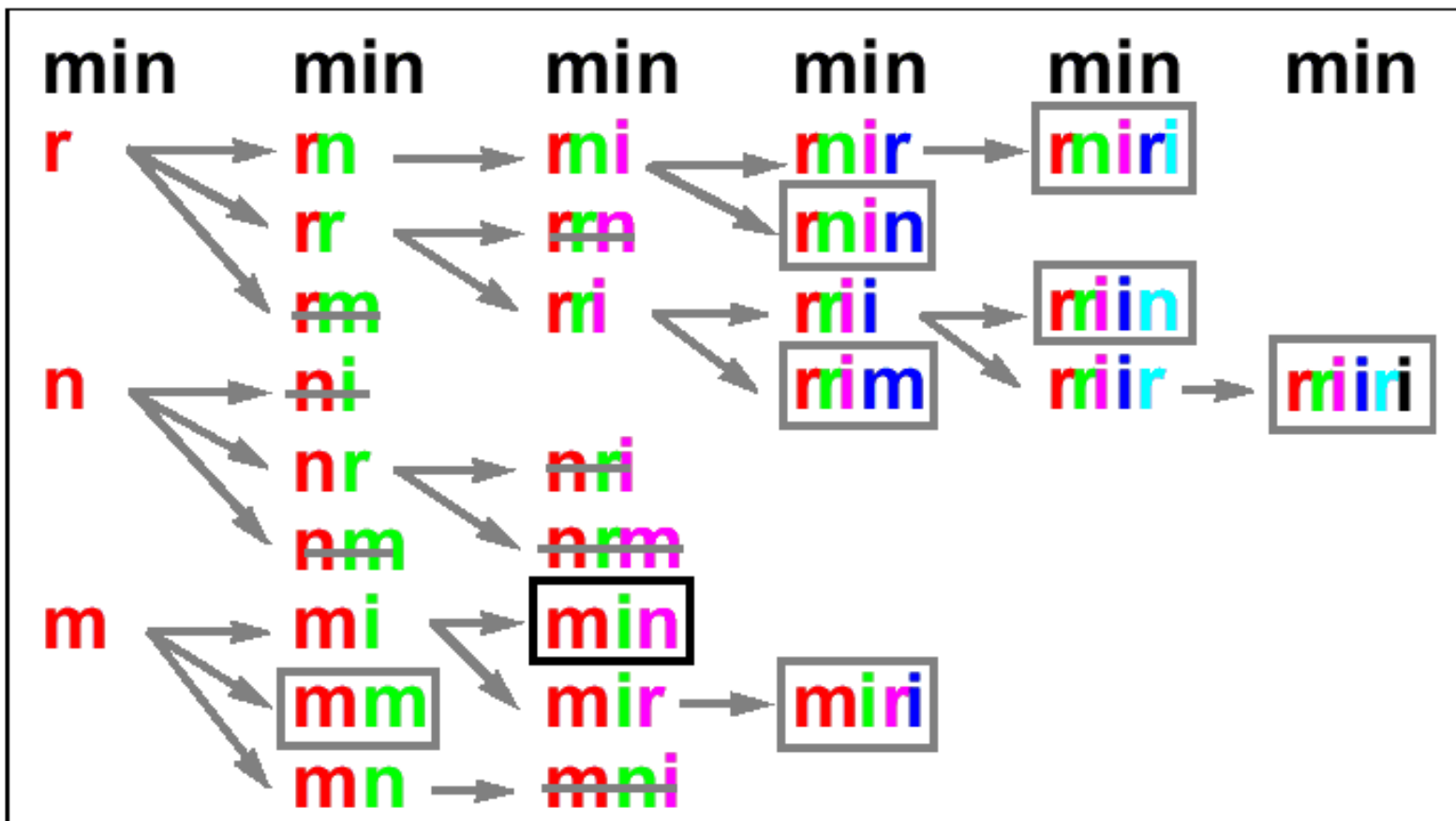
WIPP WAC SUMMARY		
WAC	CONTACT HANDLING	REMENTS-WIPAC
WAC 1000	Waste must be received at the WIPP in a container that meets the requirements of 49 CFR 173.412 (a) and (b). The container must be certified to have a design life of at least 20 years from the date of certification. The container must be certified to have a design life of at least 20 years from the date of certification.	WIPAC 1000: The container must be certified to have a design life of at least 20 years from the date of certification. The container must be certified to have a design life of at least 20 years from the date of certification.
WAC 1001	Waste must be received at the WIPP in a container that meets the requirements of 49 CFR 173.412 (c) and (d). The container must be certified to have a design life of at least 20 years from the date of certification. The container must be certified to have a design life of at least 20 years from the date of certification.	WIPAC 1001: The container must be certified to have a design life of at least 20 years from the date of certification. The container must be certified to have a design life of at least 20 years from the date of certification.
WAC 1002	Waste must be received at the WIPP in a container that meets the requirements of 49 CFR 173.412 (e) and (f). The container must be certified to have a design life of at least 20 years from the date of certification. The container must be certified to have a design life of at least 20 years from the date of certification.	WIPAC 1002: The container must be certified to have a design life of at least 20 years from the date of certification. The container must be certified to have a design life of at least 20 years from the date of certification.
WAC 1003	Waste must be received at the WIPP in a container that meets the requirements of 49 CFR 173.412 (g) and (h). The container must be certified to have a design life of at least 20 years from the date of certification. The container must be certified to have a design life of at least 20 years from the date of certification.	WIPAC 1003: The container must be certified to have a design life of at least 20 years from the date of certification. The container must be certified to have a design life of at least 20 years from the date of certification.

*Information summarized from 2001 Waste Acceptance Criteria for the WIPP (WAC-WAC-001, Rev. 1)

Remote-handled TRU waste containers shall be noncombustible and meet, as a minimum, the structural requirements and design conditions for Type A packaging contained in 49 CFR 173.412. Due to the special characteristics and application of the RH TRU canister, the compression test requirement of 49 CFR 173.465 (d) is not applicable. In addition, all RH waste containers shall be certified to a WIPP-approved specification to have a design life of at least 20 years from the date of certification.

. - . 1 2 3 4 7 9 C
 F H I P R T U a b c
 d e l g h i l m n o
 p q r s t u v w y

Level-Building Algorithm



Word Discrimination using N-grams

Adnan El-Nasan, Harsha Veermachaneni, Monday 13:35

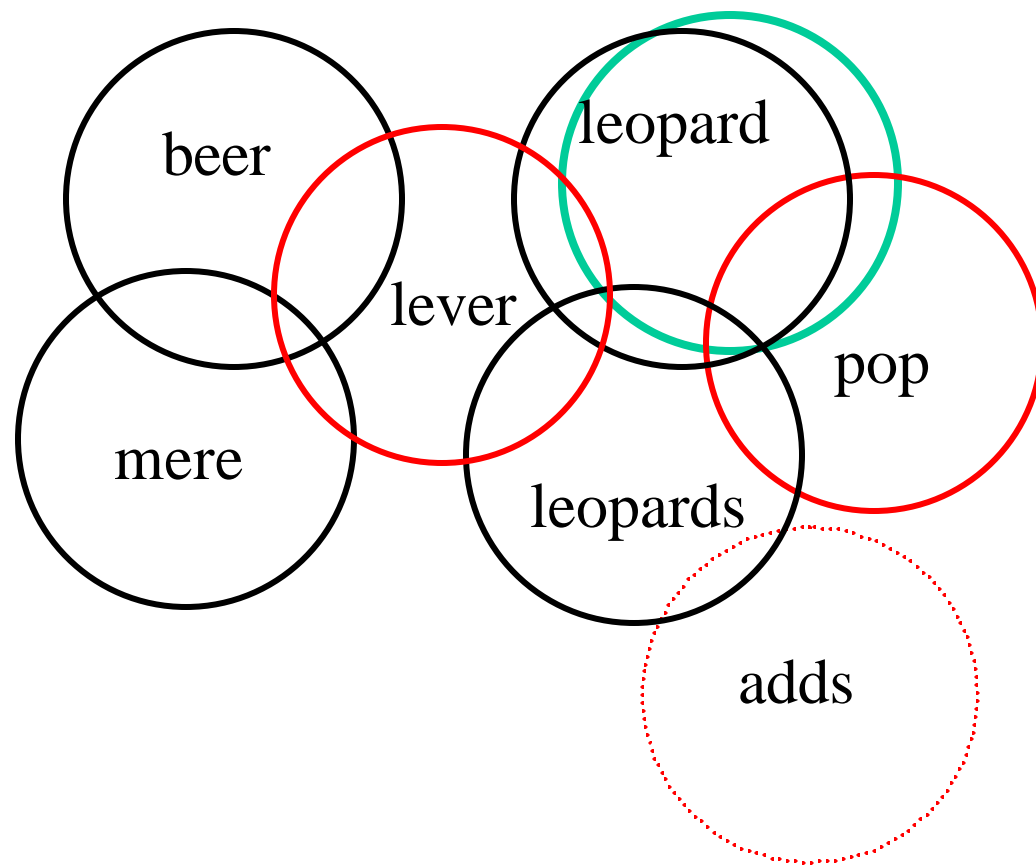
LEGEND

 Unknown word

 Lexicon word

 Reference word that has a common bigram with the unknown

 Reference word that has no common bigram with the unknown

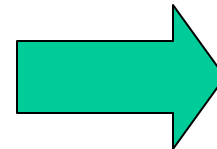


Word Discrimination using n-grams

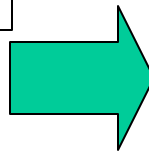
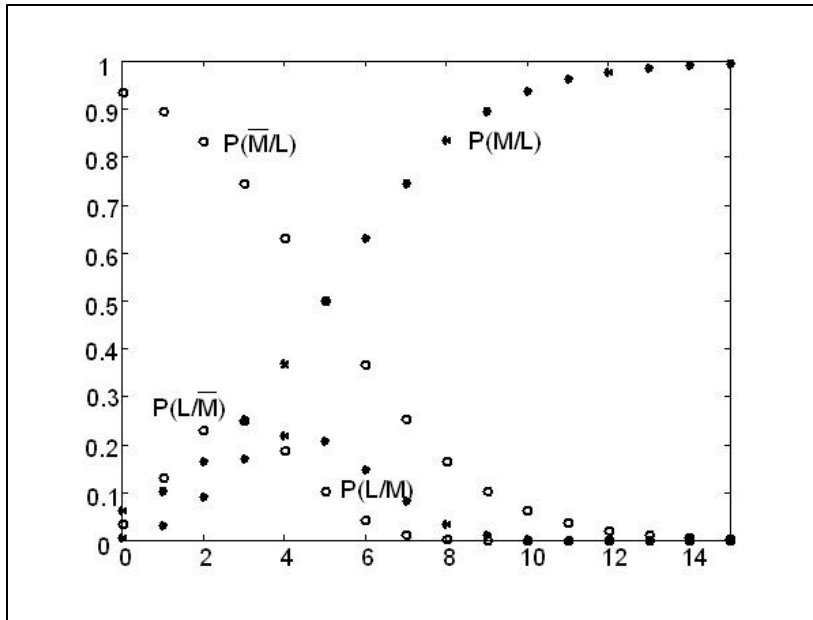
Adnan El-Nasan, Harsha Veermachaneni, Monday 13:35

Lex Ref	people	herd	open	mile	hazard	ever	period	tripod
have	0	0	0	0	1	1	0	0
lever	1	1	0	1	0	1	1	0
people	1	0	1	1	0	0	1	0
position	0	0	0	0	0	0	1	1

Reference Words		Unknown
		<i>period</i>
have	<i>have</i>	4
lever	<i>lever</i>	8
people	<i>people</i>	10
position	<i>position</i>	11



Unknown
4
8
10
11



L	P(Match/L)
4	0.3
8	0.8
10	0.9
11	0.95

$$S(\text{people}) = .7 \times .8 \times .9 \times .05 = .0252$$

$$S(\text{herd}) = .7 \times .8 \times .1 \times .05 = .0028$$

$$S(\text{open}) = .7 \times .2 \times .9 \times .05 = .0063$$

$$S(\text{mile}) = .3 \times .8 \times .9 \times .05 = .0108$$

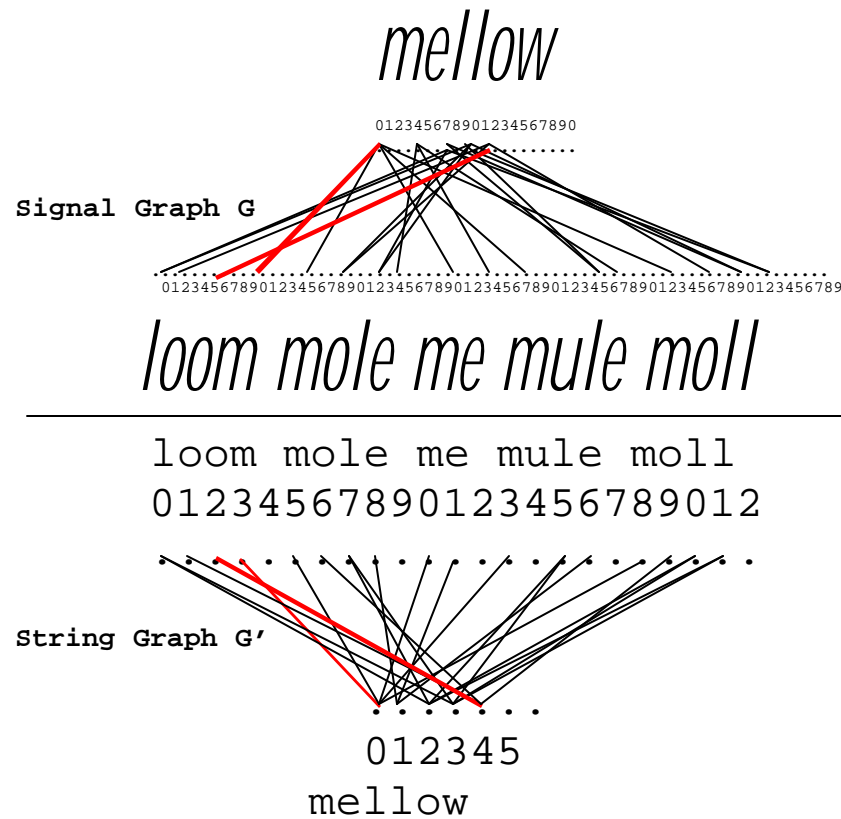
$$S(\text{hazard}) = .3 \times .2 \times .1 \times .05 = .0003$$

$$S(\text{ever}) = .3 \times .8 \times .1 \times .05 = .0012$$

$$\mathbf{S(\text{period}) = .7 \cdot .8 \cdot .9 \cdot .95 = .4788}$$

$$S(\text{tripod}) = .7 \times .2 \times .1 \times .95 = .0133$$

SYMBOLIC INDIRECT CORRELATION

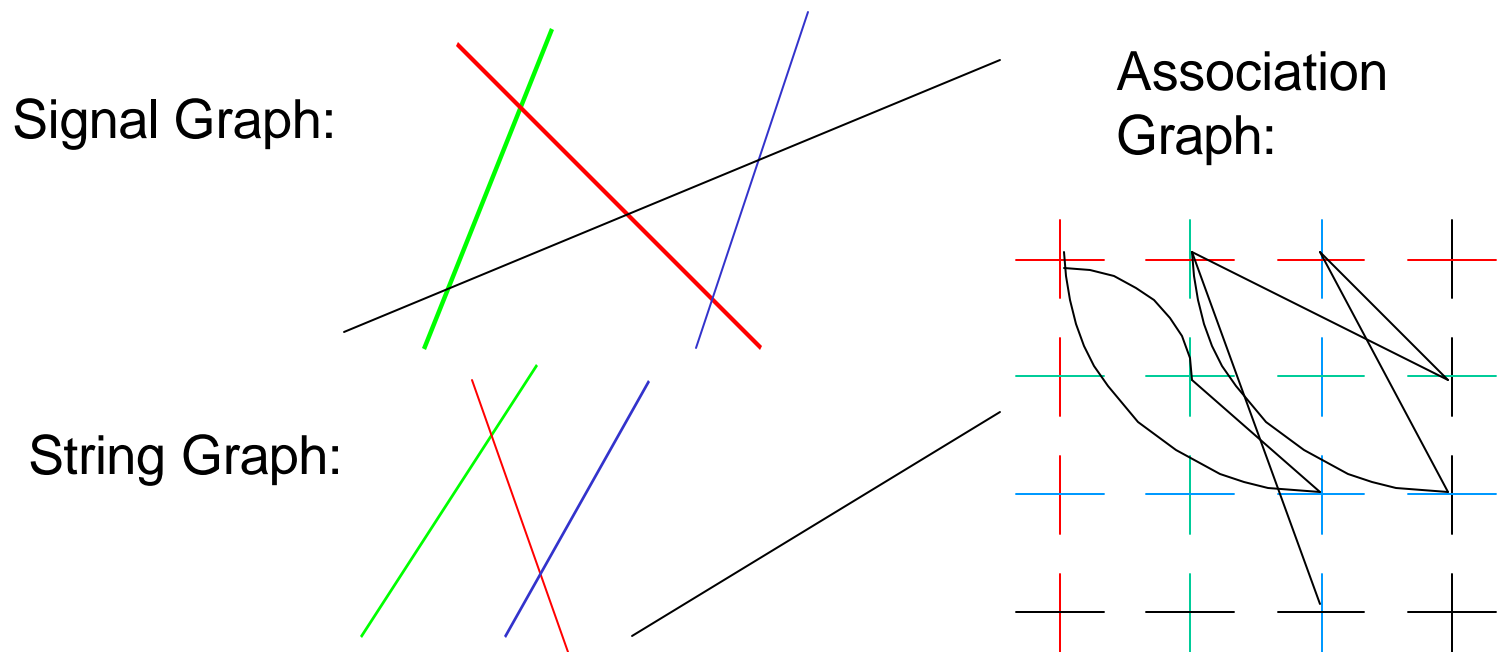


Lexicon L	{low, me, mole, mule, moll, mellow, wool, loom, we}
Unkown q(t)	mellow
Ref Signal r(t)	low me mole mule moll
Ref String r _s	low me mole mule moll
X _{mellow}	{(0,2), (0,3), (1,4), (2,4) , (3,0) , (5,0), (6,4), (7,2), ...}
X _{we}	{(8, 1), (11, 1), (16, 1)}
X'(q)	{(0,7), (0,9), (2,11), (6,11) , (10,0) , (15,0), (19,11)...}
f(q)	mellow

MAXIMUM CLIQUE FORMULATION

(Harsha Veeramachaneni, Prof. M. Krishnamoorthy (RPI))

Consider an *Association Graph* where each vertex corresponds to a pair of edges, one in the signal graph and one in the match graph.



Edges in the Association Graph indicate order compatibility between *pairs of pairs* in the match graphs.

Here the maximum clique is of size 3, the largest *subset order isomorphism*.

SNAP TEST

1. Does Bagging reduce classifier bias, variance, or neither?
Justify your answer.
2. Can feature correlation reduce classification error over uncorrelated features with the same class-conditional means?
3. If two separate sets of features are available, is it better to combine feature information before or after classification? Why?
4. What are the three essential conditions for Stochastic Discrimination?
5. What kind of confusions encourage using style-conscious classification?
6. Where and why is EM used in HMM classification?
7. How can unlabeled samples improve classifier accuracy?
Give an example where the classification after adaptation is worse.
8. Find two English words of at least seven letters that share exactly the same letter bigrams.
9. What role do Maximal Cliques play in Indirect Symbolic Correlation?

10. More accurate OCR will result most likely from:
- a. Better digitization
 - b. Improved preprocessing (e.g. layout analysis)
 - c. More discriminating features
 - d. More advanced classifiers for isolated patterns
 - d. Further exploitation of linguistic context
 - f. Style context
 - g. Unsupervised adaptation
 - h. Whole-word, line, or page classification
 - i. None of the above.

Please justify your choice.

THANK YOU!