

# Discriminative Learning and Visual Interactive Behavior

Learning Techniques in Audio-Visual Information Processing (ICPR Tutorial)

Tony Jebara  
Brian Clarkson  
Sumit Basu  
Alex Pentland

MIT  
Media  
Lab

## Outline

### -Motivation

*Learning Tasks and Paradigms*

### -Discriminative / Conditional Learning

*Maximum Entropy Discrimination*

### -Latent Variables and Reversing Jensen's Inequality

*CEM and Dual of EM*

### -Action-Reaction Learning

*Behavior Analysis / Synthesis via Time Series Prediction*

### -Wearable Platforms

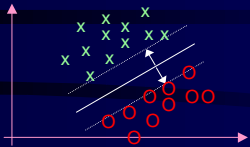
*Personal Enhanced Reality System*

*Wearable Interaction Learning*

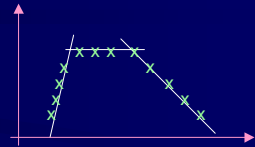
### -Conversational Context Learning

## Learning Applications for A & V

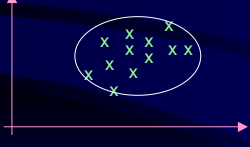
### Classification:



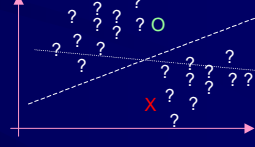
### Regression / Prediction:



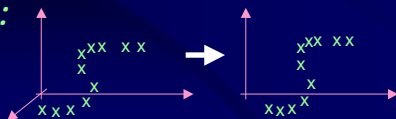
### Detection / Clustering:



### Transduction:



### Feature Selection:



## Learning Paradigms

### New Competitors for Maximum Likelihood:

### Discriminative Learning & SVMs

*(NIPS, COLT, UAI, ICML)*

### -Time Series Prediction (Muller)

### -Digit Recognition (Vapnik)

### -Speech Recognition (Deng)

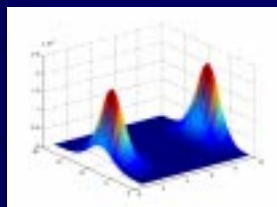
### -Face Gender Classification (Moghaddam)

### -Gene-Sequence Classification (Jaakkola)

### -Text Classification (Joachims)

## Learning Paradigms

1) Generative Approach:  
Probabilistic Models  
Maximum Likelihood



2) Discriminative Approach:  
Support Vector Machines  
VC Dimension  
Maximum Margin  
Task is Explicit  
Discriminant Surface



## Complementary Pros & Cons

### 1) ML & PDFs

- +Natural Models (HMMs)
- +Priors
- +Missing Data
- +Flexible
- Poor Performance
- Objective not Task Related

### 2) Discrimination & SVMs

- +Model & Data Mismatch
- +Support Vectors
- +Good Generalization
- Linear Model
- Kernels
- No Priors
- No Missing Data

*How to combine both? MED...*

# Maximum Entropy Discrimination

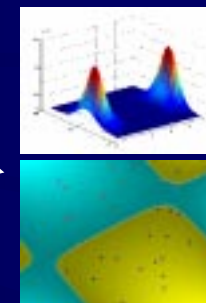
Tony Jebara  
Tommi Jaakkola  
Marina Meila

## Overview & Motivation

### Maximum Entropy Discrimination:

Combines probabilistic methods (and extensions) in discriminative framework

- Add task-related Objective to PDFs.
- Satisfy task constraints
- Support Vectors -> Generalization
- Convex, No Local Minima  
(vs. Min Class Error)



### Feasible MED Extensions & Applications:

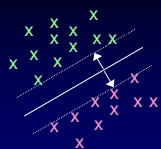
Latent variables, various priors, missing labels, structure estimation, anomaly detection.  
Feature selection, regression, latent transformations, multi-class classification, exponential family.

## Classification - Regularization Approach

**Given:** training examples:  $\{X_1, \dots, X_T\}$   
 binary (+/- 1) labels:  $\{y_1, \dots, y_T\}$   
 discriminant function:  $L(X; \Theta)$   
 non-increasing margin loss:  $l(\cdot)$

**Minimize:** regularization penalty:  $R(\Theta) + \sum_t l(\gamma_t)$   
 subject to classification constraints:  $y_t L(X_t; \Theta) - \gamma_t \geq 0, \forall t$

### Example: SVM



minimizes:  $\frac{1}{2} \|\Theta\|^2 + \sum_t f(\gamma_t)$   
 with discriminant:  $L(X; \Theta) = \Theta^T X + b$   
 decision rule:  $\hat{y} = \text{sign}(L(X; \Theta))$

## Maximum Entropy Discrimination Approach

**Many solutions** may be valid.

Use coarser description of sol'n instead of a single optimum.  
**Solve for distribution**  $P(\Theta)$  over all good  $\Theta$  (instead of  $\Theta^*$ ).  
 Find  $P(\Theta, \gamma)$  that mins  $KL(P \| P_0)$  subject to constraints:

$$\int P(\Theta, \gamma) [y_t L(X_t; \Theta) - \gamma] d\Theta d\gamma \geq 0, \forall t$$

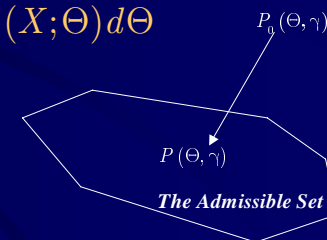
$P_0(\Theta, \gamma)$  = prior over models & margins (favors large margins).

Decision Rule:  $\hat{y} = \text{sign} \int P(\Theta) L(X; \Theta) d\Theta$

### Information Transfer / Projection:

\* Information transferred to prior after observations =  $KL(P \| P_0)$

\* Entropic Regularization and Margin Penalties are on the Same Scale



## Maximum Entropy Discrimination Solution

**Analytic, Unique, Sparse, Parametric & Structural Models:**

$$P(\Theta, \gamma) = \frac{1}{Z(\lambda)} P_0(\Theta, \gamma) \exp\left(\sum_t \lambda_t [y_t L(X_t; \Theta) - \gamma_t]\right)$$

\*  $Z(\lambda)$  = normalization constant (partition function)

\*  $\lambda = \{\lambda_1, \dots, \lambda_T\}$  = non-negative Lagrange multipliers

\*  $\lambda$  solved via unique max of concave objective function:

$$J(\lambda) = -\log Z(\lambda)$$

### Example: SVM

$$J(\lambda) = \sum_t \left[ \lambda_t + \log\left(1 - \frac{\lambda_t}{c}\right) \right] - \frac{1}{2} \sum_{t,t'} \lambda_t \lambda_{t'} y_t y_{t'} (X_t^T X_{t'})$$

### Example: Generative Models (e-family)

$$L(X; \Theta) = \log \frac{P(X|\theta_+)}{P(X|\theta_-)} + b$$

Use conjugates of  $P(X|\theta_{\pm})$  for prior  
 $\theta$  = model parameters and structure  
 $b$  = bias term

## Maximum Entropy Discrimination for Regression

Find  $P(\Theta, \gamma)$  that mins  $KL(P \| P_0)$  subject to constraints:

$$\int P(\Theta, \gamma) [y_t - L(X_t; \Theta) + \gamma_t] d\Theta d\gamma \geq 0, \forall t$$

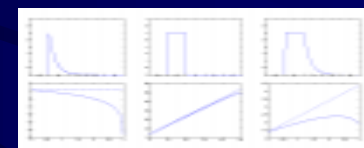
$$\int P(\Theta, \gamma) [\gamma'_t - y_t + L(X_t; \Theta)] d\Theta d\gamma \geq 0, \forall t$$

Decision Rule:  $\hat{y} = \int P(\Theta) L(X; \Theta) d\Theta$

**Solution:**  $P(\Theta, \gamma) = \frac{1}{Z(\lambda)} P_0(\Theta, \gamma) \frac{\exp(\sum_t \lambda_t [y_t - L(X_t; \Theta) - \gamma_t])}{\exp(\sum_t \lambda'_t [y_t - L(X_t; \Theta) - \gamma'_t])}$

### Margin Priors: (epsilon-tube)

$$P_0(\gamma_t) \propto \begin{cases} 1 & \text{if } 0 \leq \gamma_t \leq \epsilon \\ e^{c(\epsilon - \gamma_t)} & \text{if } \gamma_t > \epsilon \end{cases}$$



Margin prior  
Penalty Fn.

### Example: SVM

$$J(\lambda) = \sum_t y_t (\lambda'_t - \lambda_t) - \epsilon \sum_t (\lambda'_t + \lambda_t) + \sum_t \log(\lambda_t) - \log(1 - e^{-\lambda_t \epsilon} + \frac{\lambda_t}{c}) + \sum_t \log(\lambda'_t) - \log(1 - e^{-\lambda'_t \epsilon} + \frac{\lambda'_t}{c}) - \frac{1}{2} \sum_{t,t'} (\lambda_t - \lambda'_t) (\lambda_t - \lambda'_t) (X_t^T X_{t'})$$

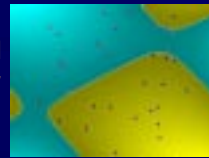
## Classification and Regression Examples

### Generative Model Classification:

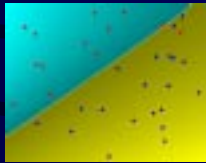
SVM  
Linear  
Kernel



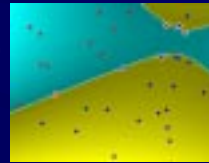
SVM  
3rd Order  
Polynomial



Maximum  
Likelihood  
Full Covariance  
Gaussians

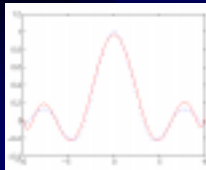


MED  
Full Covariance  
Gaussians

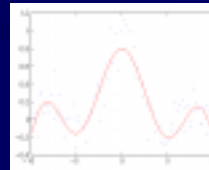


### SVM Regression:

Sinc Fn.  
(Clean)

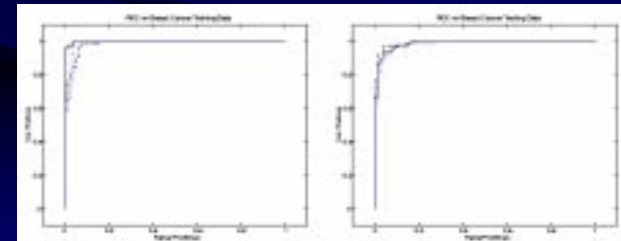


Sinc Fn.  
With  
Gaussian  
Noise



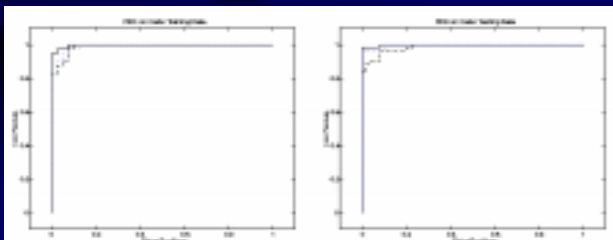
## Maximum Entropy Discrimination - Cancer

Method	Training Errors	Testing Errors
Nearest Neighbour		11
<b>SVM - Linear</b>	<b>8</b>	<b>10</b>
<b>SVM - RBF <math>\sigma = 0.3</math></b>	<b>0</b>	<b>11</b>
<b>SVM - 3rd Order Polynomial</b>	<b>1</b>	<b>13</b>
<b>Maximum Likelihood Gaussians</b>	<b>10</b>	<b>16</b>
<b>MaxEnt Discrimination Gaussians</b>	<b>3</b>	<b>8</b>



## Maximum Entropy Discrimination - Crabs

Method	Training Errors	Testing Errors
Neural Network (1)		3
Neural Network (2)		3
Linear Discriminant		5
Logistic Regression		4
MARS (degree = 1)		4
PP (4 ridge functions)		6
Gaussian Process (HMC)		3
Gaussian Process (MAP)		3
<b>SVM - Linear</b>	<b>5</b>	<b>5</b>
<b>SVM - RBF <math>\sigma = 0.3</math></b>	<b>1</b>	<b>18</b>
<b>SVM - 3rd Order Polynomial</b>	<b>3</b>	<b>6</b>
<b>Maximum Likelihood Gaussians</b>	<b>4</b>	<b>7</b>
<b>MaxEnt Discrimination Gaussians</b>	<b>2</b>	<b>3</b>



## Feature Selection (Extension)

- \* Isolates interesting dimensions of the data for the given task
- \* Typically needs exponential search:
  - consider all possible subsets of dimensions
- \* Reduces complexity of data
- \* Also *Improves Generalization*
- \* Augments Sparse Vectors (SVMs) with Sparse Dimensions
- \* Is possible jointly with parameter estimation.
- \* Can be done discriminatively and efficiently with MED.

## Feature Selection

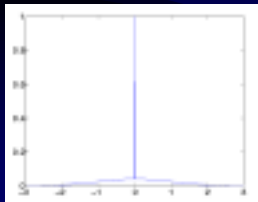
Modify parameters to include a binary ON / OFF **Switch**

$$L(X; \Theta) = \sum_{i=1}^n s_i \theta_i X_i + \theta_0$$

The model  $\Theta = \{\theta_0, \dots, \theta_n, s_1, \dots, s_n\}$  contains structural parameters  $s_i \in \{0, 1\}$  to aggressively prune features.

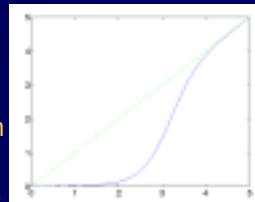
**Prior:**  $P_0(\Theta) = P_{0,\theta_0}(\theta_0) P_{0,\theta}(\theta) P_{0,s}(s) = P_{0,\theta_0}(\theta_0) \mathcal{N}(\theta | 0, I) \prod_i P_0(s_i)$

**Switch Prior:** Bernoulli distribution  $P_{s_i}(s_i) = p_0^{s_i} (1 - p_0)^{1-s_i}$   
 $p_0$  parameter smoothly selects no pruning to aggressive pruning



Prior on  $s_i \theta_i$

Aggressive attenuation of linear coefficients at low values ( $p_0=0.01$ ).

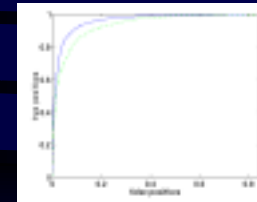


## Feature Selection in SVM Classification & Results

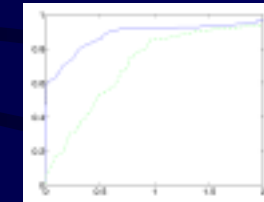
$$J(\lambda) = \sum_t \left[ \lambda_t + \log \left( 1 - \frac{\lambda_t}{c} \right) \right] - \sum_{i=1}^n \log \left[ 1 - p_0 + p_0 e^{\frac{1}{2} \left( \sum_i \lambda_i y_i X_{i,i} \right)^2} \right]$$

$\lambda$  constrained to  $[0, c]$  hyper-cube with constraint  $\sum_i \lambda_i y_i = 0$

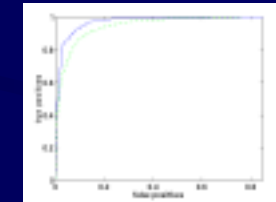
**DNA Data:** 2-class, 100 element binary vectors. Training Set 500, Testing 4724



ROC of DNA Splice Site  
100 Features  
Original 25xGATC



CDF of Linear Coeffs  
DNA Splice Site  
100 Features



ROC DNA Splice Site  
~5000 Features  
Quadratic Kernel

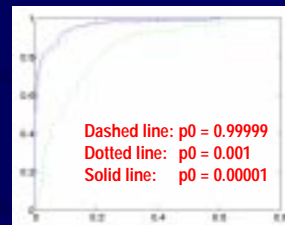
Dashed line:  $p_0 = 0.99999$   
Solid line:  $p_0 = 0.00001$

## Feature Selection in SVM Regression & Results

$$J(\lambda) = \sum_t y_t (\lambda'_t - \lambda_t) - \epsilon \sum_t (\lambda'_t + \lambda_t) - \frac{1}{2} \sigma \left( \sum_t \lambda_t - \lambda'_t \right)^2 + \sum_t \log(\lambda_t) - \log \left( 1 - e^{-\lambda_t \epsilon} + \frac{\lambda_t}{c - \lambda_t} \right) + \sum_t \log(\lambda'_t) - \log \left( 1 - e^{-\lambda'_t \epsilon} + \frac{\lambda'_t}{c - \lambda'_t} \right) - \sum_i \log \left( 1 - p_0 + p_0 e^{\frac{1}{2} \left( \sum_i \lambda_i (-\lambda'_i) X_{i,i} \right)^2} \right)$$

**Boston Housing Data:** 13 scalar features. Training Set 481, Testing 25  
Explicit Quadratic Kernel Expansion Used

Linear Model Estimator	Epsilon-Sensitive Linear Loss
Least-Squares	1.7584
MED $p_0 = 0.99999$	1.7529
MED $p_0 = 0.1$	1.6894
MED $p_0 = 0.001$	1.5377
MED $p_0 = 0.00001$	1.4808



**D. Ross Cancer Data:** 67 scalar features. Training Set 50, Testing 3951

Linear Model Estimator	Epsilon-Sensitive Linear Loss
Least-Squares	3.609e+03
MED $p_0 = 0.00001$	1.6734e+03

## Feature Selection in Generative Models

Feature selection is not limited to SVMs.

Applies to discriminative Generative Model Estimation as well.  
But, tractable computation sometimes needs approximations.

**Example:** 2-class Gaussian distributions  
variable means, identity covariance

Parameters:  $\{\mu, \nu\}$  Prior:  $P_0(\mu) \sim P_0(\nu) \sim \mathcal{N}(0, I)$

Switches:  $\{s, r\}$  Prior:  $P_0(s_i) \sim P_0(r_i) = p_0^{r_i} (1 - p_0)^{1-r_i}$

Discriminant Function (tractable in this case):

$$L(X; \Theta) = \log \frac{P(X|0_+)}{P(X|0_-)} + b = \sum_i s_i (X_i - \mu_i)^2 - \sum_i r_i (X_i - \nu_i)^2 + b$$

## Latent Transformations (Extension)

Each input example has additional unobservable properties  
 Only have a prior distribution over the unobservable  
 I.e. category, affine transform, latent variable, alignment

**Given:** training examples:  $\{X_1, \dots, X_T\}$   
 binary (+/- 1) labels:  $\{y_1, \dots, y_T\}$   
 hidden transformations:  $\{U_1, \dots, U_T\}$   
 transformation function:  $\hat{X} = T(X, U)$   
 prior on transforms:  $P_0(U_i)$



**Solution:**  $P(\Theta, U, \gamma) = \frac{1}{Z(\lambda)} P_0(\Theta, U, \gamma) \exp\left(\sum_i \lambda_i [y_i L(T(X_i, U_i); \Theta) - \gamma_i]\right)$

Transductive and Iterative. Solve iteratively by alternating solution of P(θ) and P(U).

**Example:**  $L(\hat{X}_i; \Theta) = L(T(X_i, U_i); \Theta) = \theta^T (X_i - U_i \vec{1}) + b$

## Optimization & Bounded QP (Extension)

MED maximizes concave objective with convex constraints.  
 Axis-parallel, Newton, gradient descent will converge.

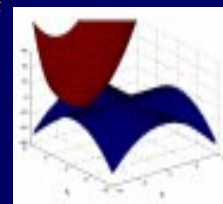
**Lower Bound** the concave objective with quadratic  
 Can then use SMO, QP, and other SVM optimizers.

**Example: SVM Feature Selection**

$$j_i(\lambda) = -\log\left(1 - p_0 + p_0 e^{2\left[\sum_i (\lambda_i - \lambda'_i) X_{i,i}\right]^2}\right) = -\log\left(1 - p_0 + p_0 e^{\frac{1}{2} \vec{\lambda}^T M \vec{\lambda}}\right)$$

$$j_i(\lambda) \geq \vec{\lambda}^T (N + hM) \vec{\lambda} - \frac{1}{2} \vec{\lambda}^T (M + N) \vec{\lambda} + const$$

where  $N = \frac{1}{4} (M \vec{\lambda}) (M \vec{\lambda})^T$  and  $h = \frac{1 - p_0}{1 - p_0 + p_0 \exp(\frac{1}{2} \vec{\lambda}^T M \vec{\lambda})}$



Iterate bound (contact at  $\tilde{\lambda}$ ) and QP  
 Each QP is seeded at previous sol'n  
 Converges in about 10 fast iterations

## MED for the Exponential Family (Extension)

**Proof:** MED with generative models spans members of the exponential family (where Gaussians generate SVMs):

exponential family form:  $p(X | \theta) = \exp(A(X) + X^T \theta - K(\theta))$

conjugate prior:  $p(\theta | \chi) = \exp(\tilde{A}(\theta) + \theta^T \chi - \tilde{K}(\chi))$

**Analytic** Partition Function for Classification:

$$Z_{\Theta} = \int P_0(\Theta) \exp\left(\sum_i \lambda_i y_i L(X_i; \Theta)\right) d\Theta$$

$$Z_{\Theta} = \int P_0(\theta_+) P_0(\theta_-) P_0(b) \exp\left(\sum_i \lambda_i y_i \left[\log \frac{p(X_i | \theta_+)}{p(X_i | \theta_-)} + b\right]\right) d\Theta$$

$$Z_{\theta_{\pm}} = \int \exp(\tilde{A}(\theta_{\pm}) + \theta_{\pm}^T \chi - \tilde{K}(\chi)) \exp\left(\sum_i \lambda_i y_i (A(X_i) + X_i^T \theta_{\pm} - K(\theta_{\pm}))\right) d\theta_{\pm}$$

$$Z_{\theta_{\pm}} = \exp(-\tilde{K}(\chi) + \sum_i \lambda_i y_i A(X_i)) \times \int \exp(\tilde{A}(\theta_{\pm}) + \theta_{\pm}^T (\chi + \sum_i \lambda_i y_i X_i)) d\theta_{\pm}$$

$$Z_{\theta_{\pm}} = \exp(-\tilde{K}(\chi) + \sum_i \lambda_i y_i A(X_i)) \times \exp(\tilde{K}(\chi + \sum_i \lambda_i y_i X_i))$$

Using Non-  
Informative  
Prior on b

## Concluding Ideas on MED and Feature Selection

**Maximum Entropy Discrimination** is a flexible Bayesian regularization approach. It provides a geometric view of learning as constrained minimization to prior distributions over margins, parameters, latent variables. It simultaneously combines:

- probabilistic methods
- large margin discrimination and SVMs
- feature selection
- classification, regression, etc.
- parameter and structure estimation
- exponential family generative models
- transduction and detection

**Feature Selection** is a particularly advantageous extension which provides increased sparsity (support vectors & support dimensions) and improves generalization.

## Limitation of MED

Applies to Exponential Family

Yet many models are MIXTURES of E-family:

Latent Models

HMMs

Mixture Models

Incomplete Data

Hidden Variable Bayesian Networks

Intractable models in discriminative & conditional settings

Thus use Variational Bounds to Perform Calculations

Invoke EM and derive its Discriminative DUAL

by Reversing Jensen

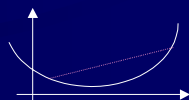
# Reversing Jensen's Inequality

## The Dual of EM for Discriminative Latent Learning

Tony Jebara  
Alex Pentland

## Jensen's Inequality

- Inequalities allow us to Integrate, Maximize, Evaluate and Derive Intractable Expressions
- Convexity: 1905-1906 by J. Jensen (Dutch Mathematician & Engineer)
- See "Convex Functions, Partial Ordering and Statistical Applications" by J. Pecaric, F. Proschan and Y. Tong.



### Jensen in Statistics and EM:

- Subsumes many information theoretic bounds (Cover & Thomas)
- Subsumes the EM Algorithm (Dempster, Laird & Rubin, Baum-Welch)
- EM casts latent variable problems as complete data by solving for a lower bound on likelihood.

### Reversals of Jensen:

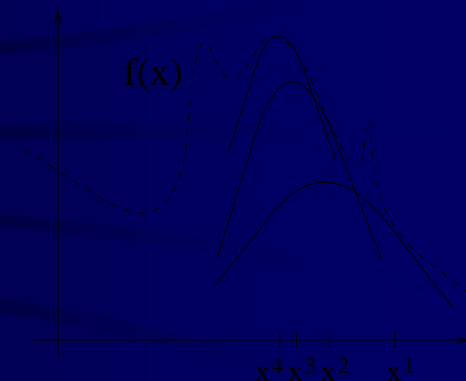
- Constrained reversals and converses have been explored and are active areas in mathematics (S.S. Dragomir).
- Reversals have yet to be applied to *discriminative learning*.

## The EM Algorithm

Makes intractable maximization of likelihood and integration of Bayesian inference tractable via variational bounds.

E-step: Replace unknowns with their expected values under current model.  
(i.e. solving for a lower bound on likelihood using Jensen!)

M-step: Optimize current model with the complete data  
(maximizing the Jensen lower bound!)



Applies and converges for Exponential Family Mixtures. I.e. a very large space of models that covers most of contemporary machine learning. HMMs, Gaussian Mixture Models, etc.

# The Exponential Family

E-family:  $K(\theta)$  is a convex function,  $X$  in its gradient space  
 Special Properties: conjugates, linearity, convexity, etc.

$$P(X | \Theta) = \exp(A(X) + X^T \theta - K(\theta)) \quad \text{TRACTABLE}$$

Includes Poisson, Gaussian, Trees, Multinomial, Exponential

E-Distribution	$A(X)$	$K(\theta)$
Gaussian	$-\frac{1}{2} X^T X - \frac{D}{2} \log(2\pi)$	$-\frac{1}{2} \theta^T \theta$
Multinomial	0	$\eta \log(1 + \sum_d \exp \theta_d)$
Exponential	0	$-\log(-\theta) \quad \forall \theta < 0$

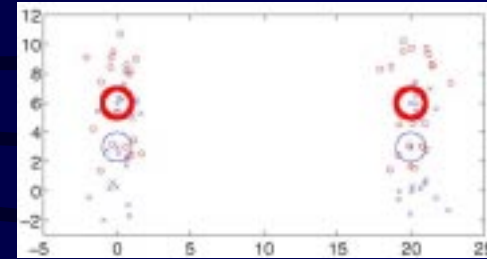
Mixtures of E-family: Gaussian mixture models, HMMs, Sigmoidal Belief Nets, Latent Bayes Nets, etc. Most machine learning models!

$$\begin{aligned}
 P(X | \Theta) &= \sum_m p(m) p(X | m, \theta) \\
 &= \sum_m \alpha_m \exp(A_m(X_m) + X_m^T \theta_m - K_m(\theta_m)) \quad \text{INTRACTABLE} \\
 &= \sum_m \sum_n \alpha_{mn} \exp(A_{mn}(X_{mn}) + X_{mn}^T \theta_m - K_m(\theta_m))
 \end{aligned}$$

# Conditional & Discriminative Classification

TWO WHITE GAUSSIANS PER CLASS FOR MODEL.

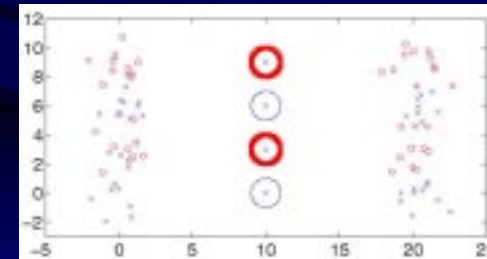
EM



$L = -8.0$   
 $L_c = -1.7$

DATA REALLY COMES FROM FOUR WHITE GAUSSIANS

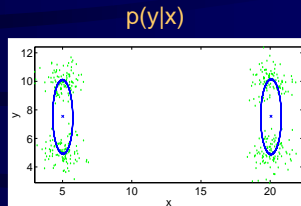
CEM



$L = -54.7$   
 $L_c = +0.4$

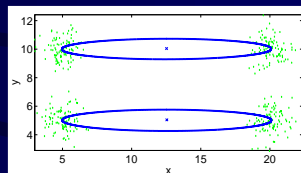
# Conditional & Discriminative Regression

PROBLEMATIC MAXIMUM LIKELIHOOD SITUATIONS



$L = -4.2$   
 $L_c = -2.4$

CONDITIONAL MAXIMUM LIKELIHOOD SOLUTIONS



$L = -5.2$   
 $L_c = -1.8$

# Discriminative Criteria and Negated Log-Sums

Maximum Likelihood: Clusters towards all data

$$l = \sum_i \log \sum_m p(m, c_i, X_i | \theta)$$

Maximum Conditional Likelihood: Emphasizes classification task

$$l^c = \sum_i \log \sum_m p(m, c_i | X_i, \theta) \quad \text{Repels models away from incorrect data}$$

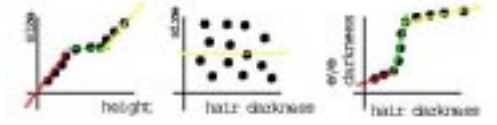
$$= \sum_i \log \sum_m p(m, c_i, X_i | \theta) - \sum_i \log \sum_m \sum_c p(m, c, X_i | \theta)$$

Maximum Margin Discrimination (MED): Emphasizes sparsity and discriminant boundary for task

$$L(X | \theta) = \log \frac{p(X | \theta_+)}{p(X | \theta_-)}$$

$$= \log \sum_m p(m, X | \theta_+) - \log \sum_m p(m, X | \theta_-)$$

The above log-sums make integration, maximization, etc. intractable. Need to simplify via overall lower and upper bounds...





## Jensen $f(E\{\square\}) \geq E\{f(\square)\}$

**DANGER: IF WE HAVE  
NEGATED LOG-SUM GET  
UPPER BOUND INSTEAD!**

Uses Concavity of  $\log()$

Reweights data with responsibilities

Variational Lower Bound at Current Model ( $\theta$ )

makes tangential contact with true objective function log-sum

Local Computations to get a Global Lower Bound

$$\log \sum_m p(m, X|\theta) \geq \sum_m \left( \frac{p(m, X|\hat{\theta})}{\sum_n p(n, X|\hat{\theta})} \right) \log \frac{p(m, X|\theta)}{p(m, X|\hat{\theta})} + \log \sum_m p(m, X|\hat{\theta})$$

$$\log \sum_m \alpha_m \exp(\mathcal{A}_m(X_m) + X_m^T \Theta_m - \mathcal{K}_m(\Theta_m)) \geq \sum_m -w_m (Y_m^T \Theta_m - \mathcal{K}_m(\Theta_m)) + k$$

$$k = \log p(X|\hat{\theta}) + \sum_m w_m (Y_m^T \hat{\Theta}_m - \mathcal{K}_m(\hat{\Theta}_m))$$

$$Y_m = \frac{1}{w_m} h_m \left( \frac{\partial \mathcal{K}_m(\Theta_m)}{\partial \Theta_m} \Big|_{\hat{\Theta}_m} - X_m \right) + \frac{\partial \mathcal{K}_m(\Theta_m)}{\partial \Theta_m} \Big|_{\hat{\Theta}_m} = X_m$$

$$w_m = -h_m = - \left( \frac{p(m, X|\hat{\theta})}{\sum_n p(n, X|\hat{\theta})} \right)$$

## Reverse-Jensen $f(E\{\square\}) \leq E\{f(\square)\}$

Uses convexity of E-family

Reweights and Translates data

Variational Upper Bound at Current Model ( $\theta$ )

makes tangential contact with true objective function log-sum

Local Computations to get a Global Upper Bound

$$\log \sum_m \alpha_m \exp(\mathcal{A}_m(X_m) + X_m^T \Theta_m - \mathcal{K}_m(\Theta_m)) \leq \sum_m -[w_m (Y_m^T \Theta_m - \mathcal{K}_m(\Theta_m)) + k]$$

$$k = \log p(X|\hat{\theta}) + \sum_m w_m (Y_m^T \hat{\Theta}_m - \mathcal{K}_m(\hat{\Theta}_m))$$

$$Y_m = \frac{h_m}{w_m} \left( \frac{\partial \mathcal{K}(\Theta_m)}{\partial \Theta_m} \Big|_{\hat{\Theta}_m} - X_m \right) + \frac{\partial \mathcal{K}(\Theta_m)}{\partial \Theta_m} \Big|_{\hat{\Theta}_m}$$

$$w'_m \#1 \rightarrow \min w'_m \text{ such that } \frac{h_m}{w'_m} \left( \frac{\partial \mathcal{K}(\Theta_m)}{\partial \Theta_m} \Big|_{\hat{\Theta}_m} - X_m \right) + \frac{\partial \mathcal{K}(\Theta_m)}{\partial \Theta_m} \Big|_{\hat{\Theta}_m} \in \frac{\partial \mathcal{K}(\Theta_m)}{\partial \Theta_m}$$

$$w_m \#2 \rightarrow w_m \geq \frac{1}{2} [X_m - \mathcal{K}'(\hat{\Theta}_m)]^T \mathcal{K}''(\hat{\Theta}_m)^{-1} [X_m - \mathcal{K}'(\hat{\Theta}_m)] + w'_m$$

OR (tighter...)

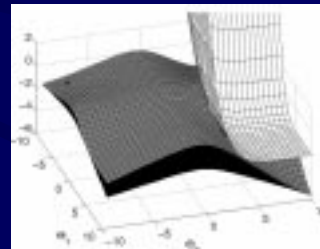
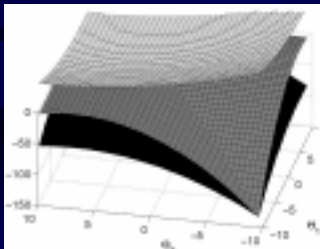
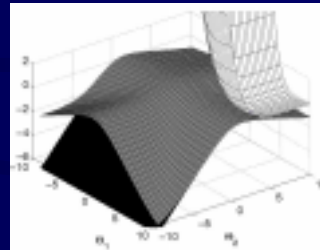
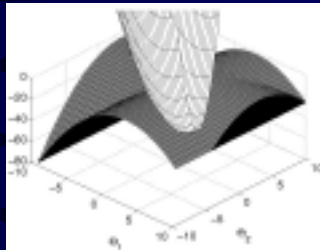
$$w_m \#2 \rightarrow w_m \geq \frac{1}{-2 \log(h_m)} [X_m - \mathcal{K}'(\hat{\Theta}_m)]^T \mathcal{K}''(\hat{\Theta}_m)^{-1} [X_m - \mathcal{K}'(\hat{\Theta}_m)] + w'_m$$

## Reverse-Jensen Bounds

Log-sum (gray)

Jensen (black)

Rev-Jensen (white)



Gaussian Case

Multinomial Case

$$\sum_m w_m (\mathcal{K}(\Theta_m) - \mathcal{K}(\hat{\Theta}_m) - (\Theta_m - \hat{\Theta}_m)^T \mathcal{K}'(\hat{\Theta}_m)) \geq \log \frac{p(X|\theta)}{p(X|\hat{\theta})} + \sum_m h_m (\Theta_m - \hat{\Theta}_m)^T (\mathcal{K}'(\hat{\Theta}_m) - X_m)$$

Define  $\mathcal{F}_m(\Theta_m) = \mathcal{K}(\Theta_m) - \mathcal{K}(\hat{\Theta}_m) - (\Theta_m - \hat{\Theta}_m)^T \mathcal{K}'(\hat{\Theta}_m)$  and  $Z_m = X_m - \mathcal{K}'(\hat{\Theta}_m)$ .

$$\sum_m w_m \mathcal{F}_m(\Theta_m) \geq \log \frac{\sum_m \exp\{D_m + \Theta_m^T Z_m - \mathcal{F}_m(\Theta_m)\}}{\sum_m \exp\{D_m + \hat{\Theta}_m^T Z_m - \mathcal{F}_m(\hat{\Theta}_m)\}} - \sum_m h_m (\Theta_m - \hat{\Theta}_m)^T Z_m$$

Define:  $\mathcal{G}_m(\Theta_m) = \mathcal{G}_m(\Phi_m) = \frac{1}{2} (\Phi_m - \hat{\Theta}_m)^T (\Phi_m - \hat{\Theta}_m)$

$$\sum_m w_m \mathcal{G}_m(\Phi_m) \geq \log \frac{\sum_m \exp\{D_m + \Theta_m(\Phi_m)^T Z_m - \mathcal{G}_m(\Phi_m)\}}{\sum_m \exp\{D_m + \hat{\Theta}_m^T Z_m - \mathcal{G}_m(\hat{\Theta}_m)\}} - \sum_m h_m (\Theta_m(\Phi_m) - \hat{\Theta}_m)^T Z_m \quad (1)$$

Hessian of  $\mathcal{F} = \mathcal{G}$ :  $\mathcal{K}''(\Theta_m) \frac{\partial \Theta_m}{\partial \Phi_m} \frac{\partial \Theta_m}{\partial \Phi_m}^T + (\mathcal{K}'(\Theta_m) - \mathcal{K}'(\hat{\Theta}_m)) \frac{\partial^2 \Theta_m}{\partial \Phi_m^2} = I$

At  $\Theta_m = \hat{\Theta}_m$ :  $\frac{\partial \Theta_m}{\partial \Phi_m} \Big|_{\hat{\Theta}_m} = [\mathcal{K}''(\hat{\Theta}_m)]^{-1/2}$ .

In an e-family, we can always find a  $\Theta_m^*$  such that  $X_m = \mathcal{K}'(\Theta_m^*)$ . By convexity of  $\mathcal{F}$  we create a linear lower bound at  $\Theta_m^*$ :

$$\mathcal{F}(\Theta_m^*) + (\Theta_m - \Theta_m^*) \frac{\partial \mathcal{F}(\Theta_m)}{\partial \Theta_m} \Big|_{\Theta_m^*} \leq \mathcal{F}(\Theta_m) = \mathcal{G}(\Phi_m)$$

Taking 2nd derivatives in  $\Phi_m$  gives:  $\mathcal{F}'(\Theta_m^*) \frac{\partial^2 \Theta_m}{\partial \Phi_m^2} \leq I$  or  $Z_m \frac{\partial^2 \Theta_m}{\partial \Phi_m^2} \leq I$

Thus,  $D_m + \Theta_m(\Phi_m)^T Z_m - \mathcal{G}(\Phi_m)$  is concave, linear upper bound at  $\hat{\Theta}_m$ :

$$\sum_m w_m \mathcal{G}_m(\Phi_m) \geq \log \frac{\sum_m \exp\{D'_m + \Phi_m^T [\mathcal{K}''(\hat{\Theta}_m)]^{-1/2} Z_m\}}{\sum_m \exp\{D_m + \hat{\Theta}_m^T Z_m - \mathcal{G}_m(\hat{\Theta}_m)\}} - \sum_m h_m (\Theta_m(\Phi_m) - \hat{\Theta}_m)^T Z_m$$

map bowl to bowl

Taking 2nd derivatives over  $\Phi_m$ :  $w_m I \geq \frac{1}{2} Z_m \mathcal{K}''(\hat{\Theta}_m)^{-1} Z_m^T - h_m Z_m \frac{\partial^2 \Theta_m}{\partial \Phi_m^2}$

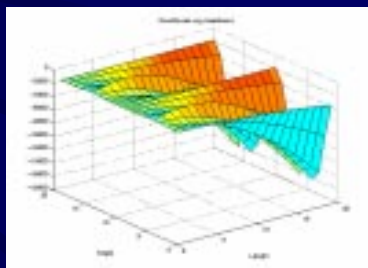
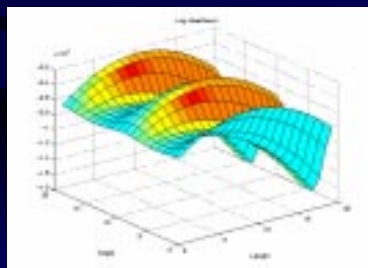
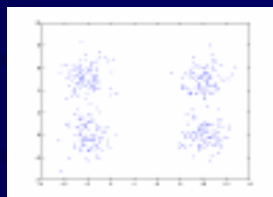
Invoke constraint #1 and replace  $-h_m Z_m \frac{\partial^2 \Theta_m}{\partial \Phi_m^2} \leq w'_m I$ . Manipulate to obtain:

$$w_m I \geq \frac{1}{2} [X_m - \mathcal{K}'(\hat{\Theta}_m)] \mathcal{K}''(\hat{\Theta}_m)^{-1} [X_m - \mathcal{K}'(\hat{\Theta}_m)]^T + w'_m I \quad \square$$

Short Proof

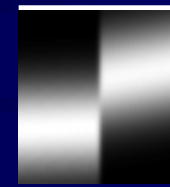
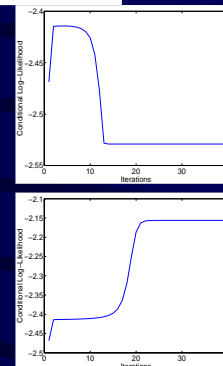
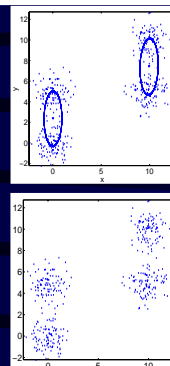
## CEM Regression Results

Estimate  $p(y|x)$  regression model with 2 Gaussians (Gaussian gates with linear experts).

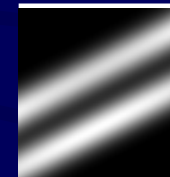
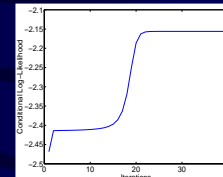
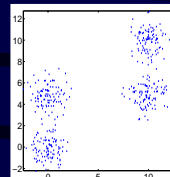


## CEM Regression Results

EM



CEM



Abalone age data set (UCI).

CEM monotonically increases conditional likelihood unlike EM. Result: better  $p(y|x)$  which captures the multimodality in  $y$  without wasting resources in  $x$ .

Algorithm

Cascade-Correlation (0 hidden)  
 Cascade-Correlation (5 hidden)  
 C4.5  
 Linear Discriminant  
 k=5 Nearest Neighbor  
 EM 2 Gaussians  
 EM&CEM 1 Gaussian  
 CEM 2 Gaussians

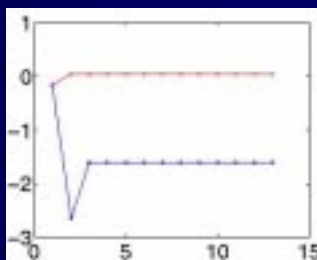
Regression Accuracy

24.86%  
 26.25%  
 21.5%  
 0.0%  
 3.57%  
 22.32%  
 20.79%  
**26.63%**

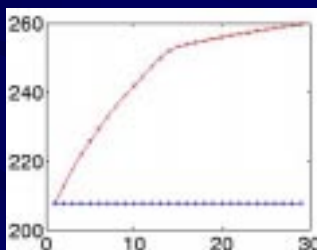
## CEM Classification Results

Gaussian mixture model shown earlier for classification. Monotonic convergence. Double computation of EM per epoch.

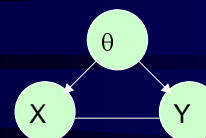
CEM accuracy = 93%  
 EM accuracy = 59%



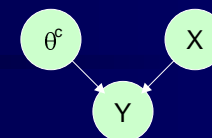
Multinomial mixture model. 3-class multinomials for 60 base-pair protein chains. CEM monotonically increases conditional likelihood.



## Bayesian Inference: Conditional vs. Joint



joint



conditional

\* Each assumes different independencies in data

= Data Structure (like Model Structure which is useful for learning)

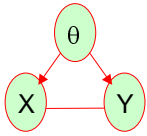
\* Exponential number of conditional models

-> Use a handful for frequent task and joint for rare task

-> Use marginal for unreliable covariates

## Bayesian Approach

$$\begin{aligned}
 p(x, y) &= p(x, y | X, Y) \\
 &= \int p(x, y, \theta | X, Y) d\theta \\
 &= \int p(x, y | \theta) p(\theta | X, Y) d\theta \\
 p(y | x)^j &= \frac{\int p(x, y | \theta) p(\theta | X, Y) d\theta}{\int p(x | \theta) p(\theta | X, Y) d\theta}
 \end{aligned}$$



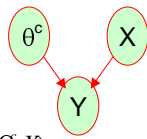
$$p(x, y | X, Y) \approx p(x, y | \theta^*)$$

where

$$\theta^* = \begin{cases} \text{MAP} \rightarrow \text{argmax}_{\theta} p(\theta | X, Y) = \text{argmax}_{\theta} p(X, Y | \theta) p(\theta) \\ \text{ML} \rightarrow \text{argmax}_{\theta} p(X, Y | \theta) \end{cases}$$

x: covariate  
y: response  
X: covariate dataset  
Y: response dataset

$$\begin{aligned}
 p(y | x)^c &= p(y | x, X, Y) \\
 &= \int p(y, \theta | x, X, Y) d\theta \\
 &= \int p(y | x, \theta) p(\theta | X, Y) d\theta \\
 p(\theta | X, Y) &= \frac{p(Y | \theta, X) p(\theta, X)}{p(X, Y)} \\
 &= \frac{p(Y | \theta, X) p(X | \theta) p(\theta)}{p(X, Y)} \\
 &= \frac{p(Y | \theta, X) p(X) p(\theta)}{p(X, Y)}
 \end{aligned}$$



$$\begin{aligned}
 p(y | x)^c &= \int p(y | x, \theta) \frac{p(Y | \theta, X) p(X) p(\theta)}{p(X, Y)} d\theta \\
 &= \int p(y | x, \theta) p(Y | \theta, X) p(\theta) d\theta / p(Y | X)
 \end{aligned}$$

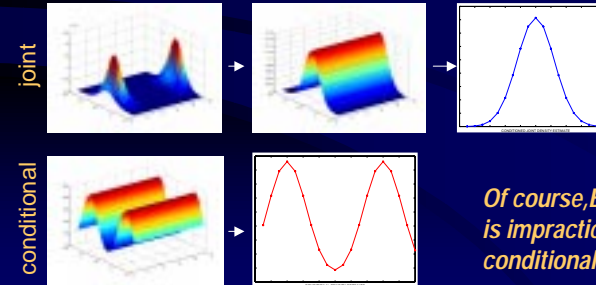
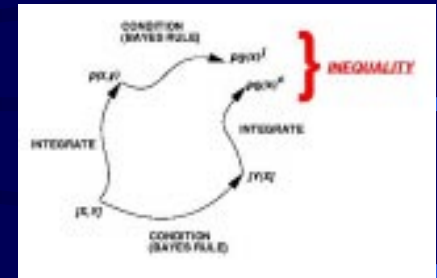
$$p(y | x) \approx p(y | x, \theta^*)$$

$$\text{where } \theta^* = \begin{cases} \text{MAP} \rightarrow \text{argmax}_{\theta} p(Y | \theta, X) p(\theta) \\ \text{ML} \rightarrow \text{argmax}_{\theta} p(y | \theta, X) \end{cases}$$

## An Exact Bayesian Example



Data:  $p(y|x)$  from 4 points to fit with 2 Gaussians



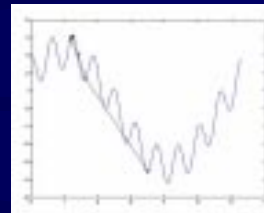
Of course, Bayesian integration is impractical so we must consider conditional MAP and ML...

## Extensions

### Annealing

Bound can accommodate a Gibbs temperature for global optimum.

$$\exp(\beta \times [A(X) + X^T \theta - K(\theta)])$$



### Latent Bayesian Networks

### Hidden Markov Models

### MED - Large Margin Latent Discrimination

### Variational Bayesian Inference

### Generic Optimization

check <http://www.media.mit.edu/~jebara>

## Action-Reaction Learning

## Automatic Visual Analysis and Synthesis of Interactive Behaviour

Tony Jebara  
Alex Pentland

# Motivation & Background

## IDEA

Computer Vision, Face & Gesture at Killington, Behaviourists...

## BEHAVIOUR PERCEPTION

Static Imagery -> Simple Temporal Models -> Learned Temporal Dynamics, Higher Order Control, Multiple Hypothesis, HMMs, NNs (Blake, Bregler, Pentland, Bobick, Hogg)

## BEHAVIOUR SYNTHESIS

Competing Behaviours, Control, Reinforcement, Ethology, Cog-Sci (Brooks, Terzopolous, Blumberg, Uchibe, Mataric, Large)

## ARL (ACTION REACTION LEARNING)

- Machine Learning of Correlations between Stimulus & Response via Perceptual Measurements of Human Interactions
- Imitation Based Learning (Mataric)
- Behaviourism (Thorndike, Watson, Skinnerian, Gibsonian) -> Reactionary
- Watch Humans Interacting to Learn how to React to Stimulus



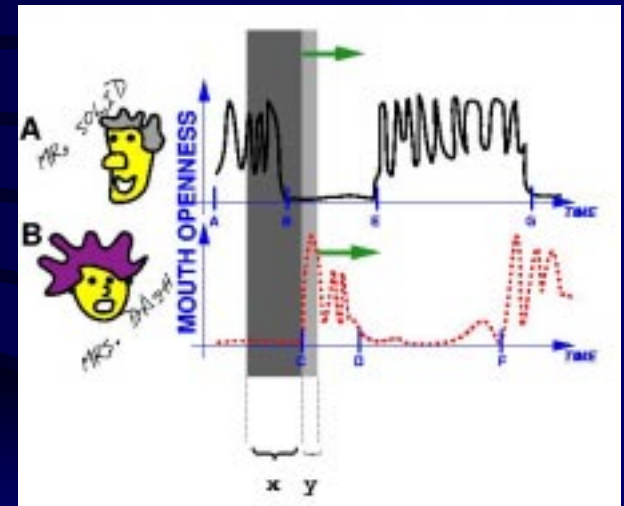
# Scenario

AUTOMATIC UNSUPERVISED OBSERVATION OF 2 AGENT INTERACTION

TRACK LIP MOTIONS

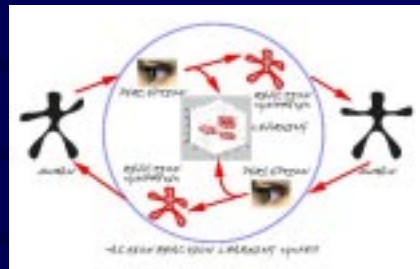
DISCOVER CORRELATIONS BETWEEN PAST ACTION & CONSEQUENT REACTION

ESTIMATE  $p(y|x)$

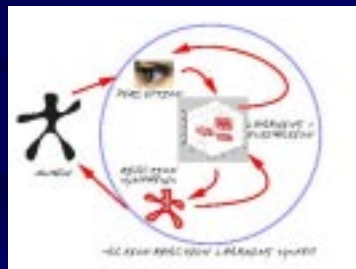


# System Architecture

OFFLINE: LEARNING FROM HUMAN INTERACTION, SPYING ON TWO USERS TO LEARN  $p(y|x)$



ONLINE: INTERACTION WITH SINGLE USER WITH LEARNED  $p(y|x)$



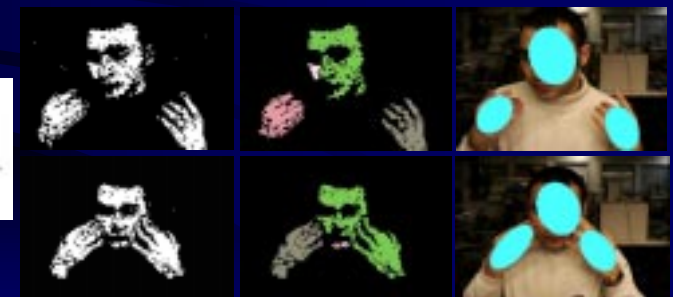
# Perception

PROBABILISTIC HEAD & HAND TRACKING

$$p(\mathbf{x}_{rgb}) = \sum_{i=1}^3 \frac{p(i)}{(2\pi)^{3/2} \sqrt{|\Sigma_i|}} e^{-\frac{1}{2}(\mathbf{x}_{rgb} - \mu_i)^T \Sigma_i^{-1} (\mathbf{x}_{rgb} - \mu_i)}$$

EXPECTATION MAXIMIZATION (EM)

$$p(\mathbf{x}_{xy}) = \sum_{j=1}^3 \frac{p(j)}{2\pi \sqrt{|\Sigma_j|}} e^{-\frac{1}{2}(\mathbf{x}_{xy} - \mu_j)^T \Sigma_j^{-1} (\mathbf{x}_{xy} - \mu_j)}$$

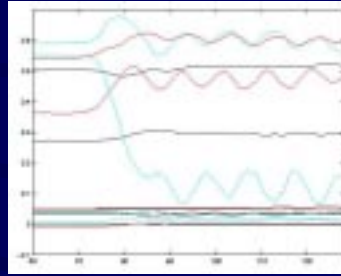


# Perception...

## TEMPORAL REPRESENTATION

5 Parameters per Blob =  
 2 Centroid +  
 3 Square Root Covariance

$$\mu_x \mu_y \Gamma_{xx} \Gamma_{xy} \Gamma_{yy}$$



## GRAPHICAL OUTPUT

(Seen by both users)



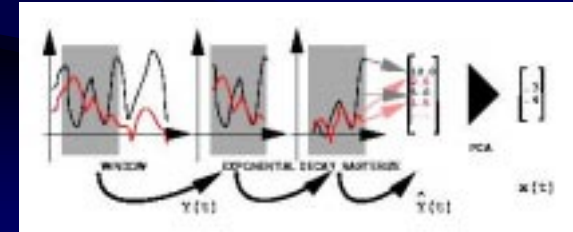
Tony Jebara - MIT Media Laboratory - 1998

# Temporal Modeling

TIME SERIES PREDICTION (Gershenfeld, Weigund, Mozer, Wan)  
 Santa Fe: NNs, RNNs, HMMs, Diff Eqns, etc.  
 Sun Spot, Bach, Physiological, Chaotic Laser

SHORT TERM MEMORY PRE-PROCESSING (Wan, Elliot-Anderson)

$y = [B_{a1} B_{a2} B_{a3} B_{b1} B_{b2} B_{b3}]$  Features (Blobs Concatenated)  
 $y(t) \approx g(y(t-1), y(t-2), \dots, y(t-T))$  Prediction Mapping  
 $Y(t) = [y(t-1) y(t-2) \dots y(t-T)]$  Short Term Memory (T=120)



# Temporal Modeling...

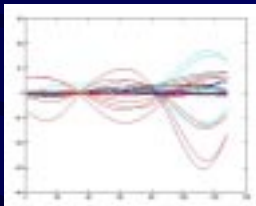
## PRINCIPAL COMPONENTS ANALYSIS

(linear, FFT, Wavelets...)

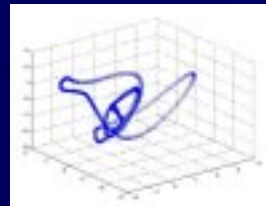
Gaussian Distribution of STM (roughly 6.5 seconds)  
 Dims = T x Feats = 120x30 ----> 40 (95% Energy of Submanifold)  
 Low dimensional (i.e. smoothed) characterization of past interaction



eigenvalues



eigenvectors



eigenspace

## LEARN MAPPING PROBABILISTICALLY

$p(y|x) = p(\text{future} | \text{STM})$  versus deterministic  $y=g(x)$

# Learning & The CEM Algorithm

## EXPECTATION MAXIMIZATION (EM)

Learns  $p(x,y)$  by maximizing  $\prod_{i=1}^N p(x_i, y_i | \Theta)$  (joint model of phenomenon)

Powerful convergence -- Clean statistical Framework  
 More global than gradient solutions -- Can be deterministically annealed

For conditional problems (input/output regression, classification)  
 Joint models are outperformed (i.e. NNs and RNNs versus HMMs)  
 Since they don't optimize output error (i.e. testing criterion is not like training)

## CONDITIONAL EXPECTATION MAXIMIZATION (CEM)

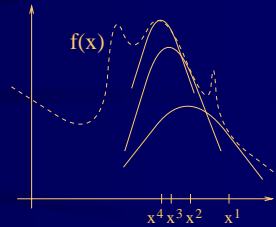
NIPS11, 1998

Learns  $p(y|x)$  by maximizing  $\prod_{i=1}^N p(y_i | x_i, \Theta)$  (conditional model of task)

Convergence properties like EM but for Conditional Likelihood

# Learning and the CEM Algorithm...

- Variational Bound Maximization
- Monotonic
- Skips Local Minima
- Deterministically Annealable



- Model: Mixture of Experts
- Soft Piece-Wise Linear
- Gaussian Gates with Conditioned Gaussian Regressors
- CEM Applies to other models: HMM, Multinomial Mixtures, etc.

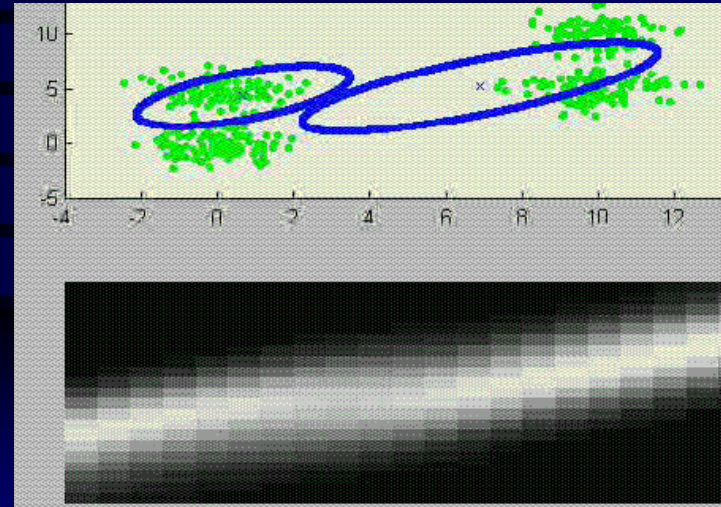
$$p(y | \mathbf{x}, \Theta) = \frac{\sum_{m=1}^M \alpha_m \bar{N}(\mathbf{x}; \mu_m, \Sigma_m) \times N(y; \nu_m + \Gamma_m \mathbf{x}, \Omega_m)}{\sum_{m=1}^M \alpha_m \bar{N}(\mathbf{x}; \mu_m, \Sigma_m)}$$

Output: Expectation or Arg Max (Regression)

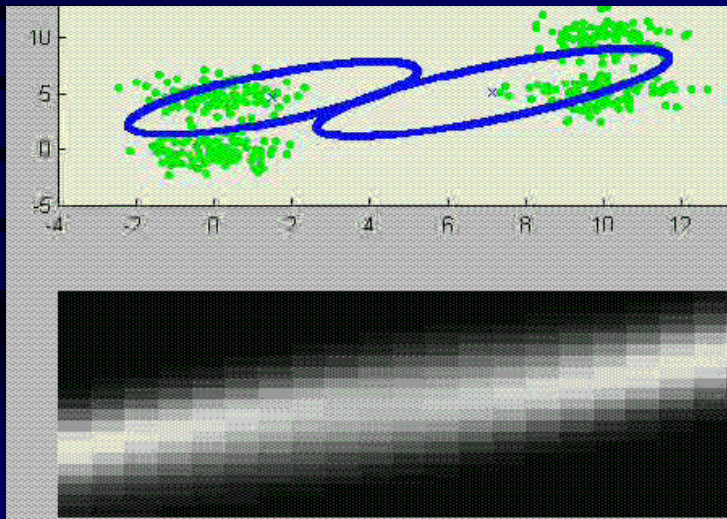
$$\hat{y} = \int y p(y | \hat{\mathbf{x}}) dy = \frac{\sum_{m=1}^M \hat{y}_m p(\hat{\mathbf{y}}_m | \hat{\mathbf{x}})}{\sum_{m=1}^M p(\hat{\mathbf{y}}_m | \hat{\mathbf{x}})} \quad \hat{y}_m = \mu_m^y + \Sigma_m^{yx} \Sigma_m^{xx^{-1}} (\hat{\mathbf{x}} - \mu_m^x)$$

Computationally very efficient, simple matrix operations

# EM



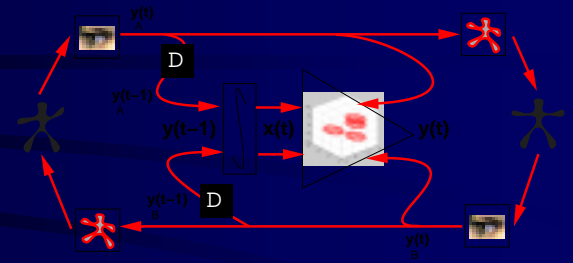
# CEM



# Integration

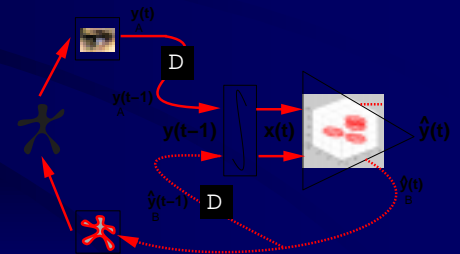
## TRAINING MODE

System accumulates action / reaction pairs  $(x,y)$  and uses CEM to optimize a conditioned Gaussian mixture model for  $p(y|x)$



## INTERACTIVE MODE

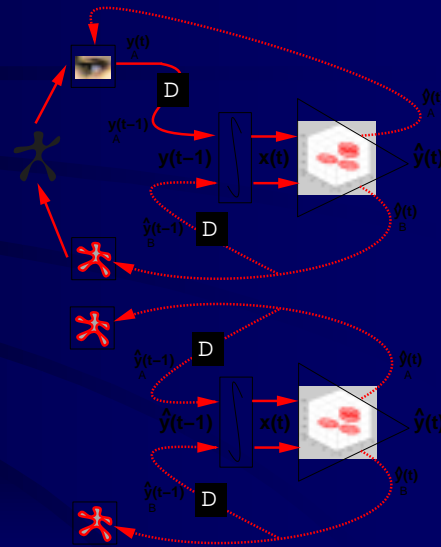
System completes missing time series and synthesizes reaction in graphical form for user using the  $p(y|x)$



## Integration...

### FEEDBACK MODE

Predicted measurement on user can be fed back as a non-linear learned EKF to aid vision. Can also use  $p(x)$  to filter vision and find correspondence.



### SIMULATION MODE

Fully synthesize both components of the time series (user A and user B). Some instabilities / bugs (no grounding) -> future work.

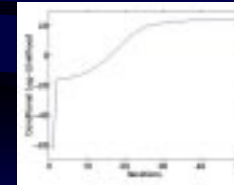
## Training & Results

Training Data: 5 minutes, 15 Hz, approx. 5 repetitions of gestures

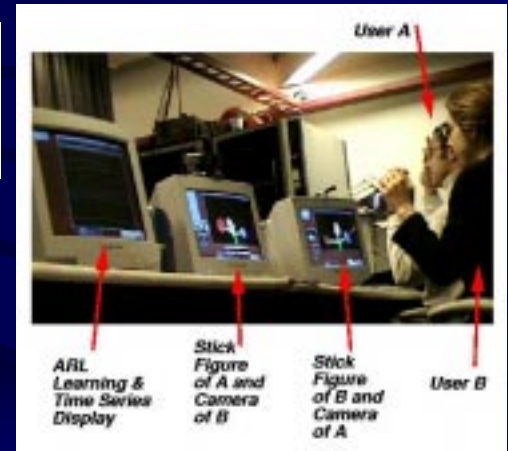
Model:  $T=120$ , Dims=22,  $M=25$  (to maintain real-time)

Convergence: 2 hours

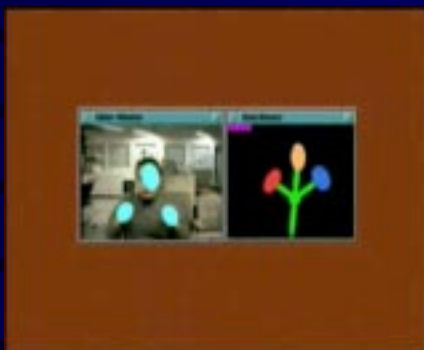
Interaction	User	Corresponding Action
1	A	Stare at by moving towards camera
2	B	Steadily cross arms & bring hands in
3	A	Wave hands
4	B	Wave back accordingly
5	A	Circle around & tap head
6	B	Clap enthusiastically
7	A	Shrill or Shout Loudly
8	B	Shrill or Shout Loudly



Nearest Neighbor: 1.57% RMS  
Constant Velocity: 0.85% RMS  
ARL: 0.64% RMS



## Results...



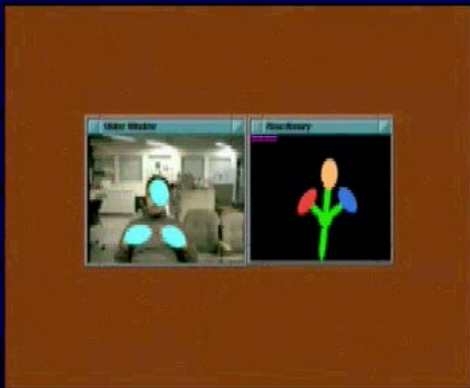
SCARE INTERACTION

## Results...



WAVE INTERACTION

## Results...

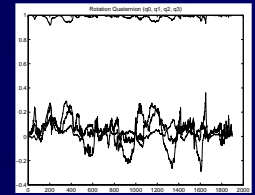


CLAP INTERACTION

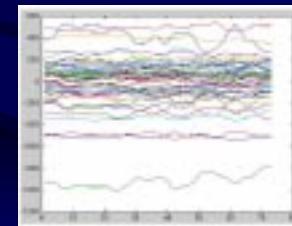
## Alternate Perceptual Modalities

### FACE MODELING

3D Translation  
3D Pose  
Focal Length  
7 DOFs



3D Eigen Model  
Deformations  
Texture  
40 DOFs



Real-Time  
(...+ speech +...)

## Conclusions

- 1 - Unsupervised Discovery of Simple Interactive Behaviour by Perceptual Observations and Statistical Time Series Prediction of Future Given Past or Reaction Given Action
- 2 - Imitation Based Learning of Behaviour
- 3 - Real-Time Behaviour Acquisition and Interactive Synthesis
- 4 - Small Amounts of Training Data and Non-Intrusive Training
- 5 - Non-Linear Predictive Model for Feedback Perception
- 6 - Monotonically Convergent Maximum Conditional Likelihood i.e. Discriminant Probabilistic Learning (CEM)
- 7 - No A Priori Segmentation or Classification of Gestures

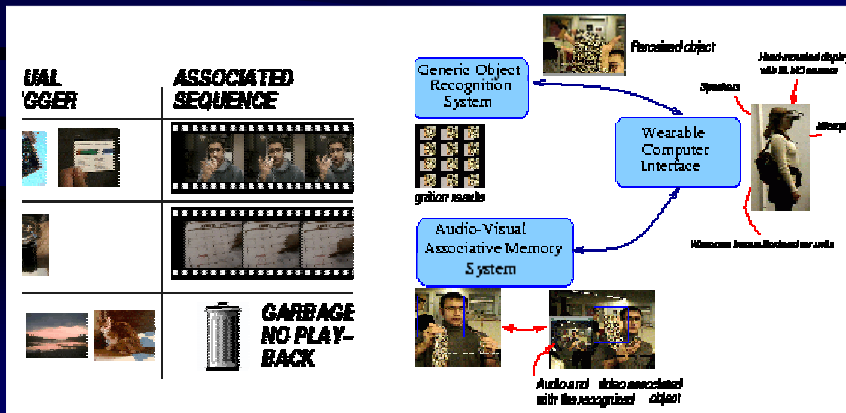
## *Wearable Platform: Dynamic Personal Enhanced Reality System*

**Tony Jebara  
Bernt Schiele  
Nuria Oliver  
Alex Pentland**



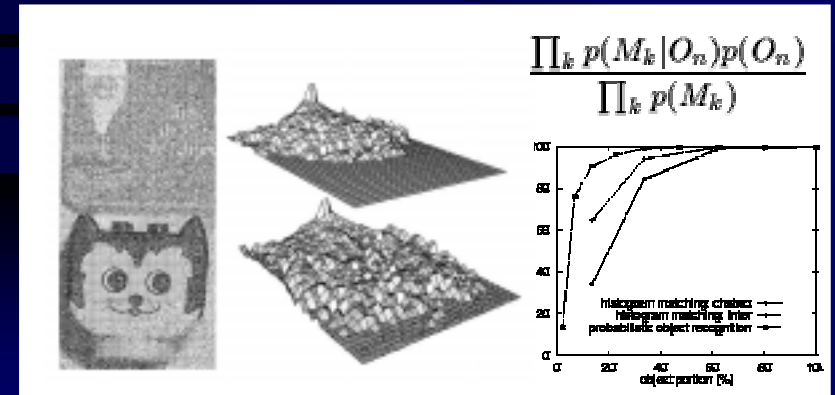
## DyPERS Architecture

- \* 3 Button interface: Record, Associated and Discard
- \* User records live A/V Clips with wearable
- \* Associates them with a visual trigger object(s)
- \* Audio-Video is replayed when computer vision sees trigger object

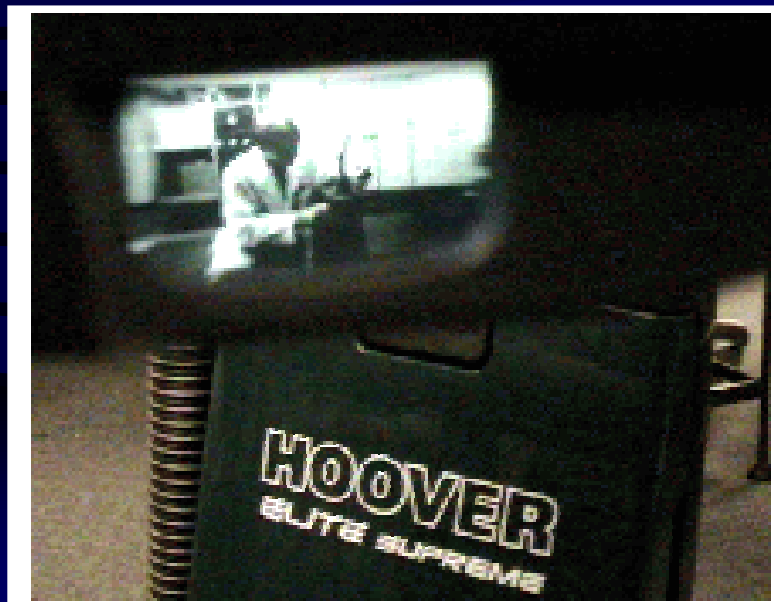


## DyPERS Visual Recognition

- \* Multidimensional filter-response histograms from differences of Gaussian linear convolutions (magnitude of 1st deriv. and Laplace operator)
- \* Compute probability of object from  $k$  iid measurements
- \* Kalman filter on probabilities of objects to smooth classifier



## Video



## Wearable Platform: Interactive Behavior Acquisition

Tony Jebara  
Alex Pentland

## Wearable Long-term Behavior Acquisition



### Hardware

- Sony Picturebook Laptop
- 2 Cameras
- 2 Microphones

### Action-Reaction Learning

- Audio-Visual processing
- Time Series Prediction
- Discriminative Learning
- DyPERS associative memory

## Tracking Conversational Context (Audio)

Tony Jebara  
Yuri Ivanov  
Ali Rahimi  
Alex Pentland

## System Architecture

### Tracking Conversational Context for Machine Mediation

Speech rec. on multiple speakers with real-time topic-spotting

Bag-of-words (multinomials)

$$P(\text{word}_i | c) = \frac{N_c(\text{word}_i)}{\sum_{j=1}^n N_c(\text{word}_j)}$$

Short Term Memory w/ Decay

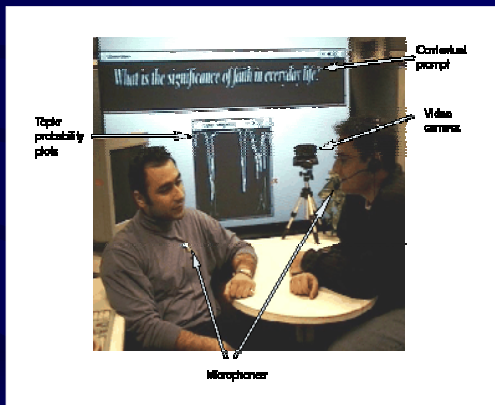
$$x_t^i = \alpha x_{t-1}^i + \delta(c, z)$$

Probability of Topic

$$P(z|c) = \prod_i P(\text{word}_i | c)^{x_t^i}$$

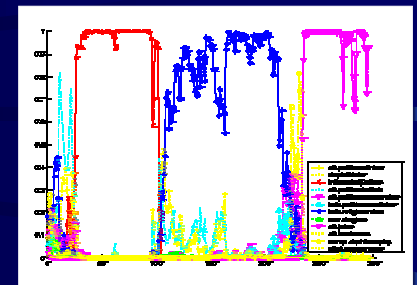
$$P(z|z) = \frac{P(z|c)P(c)}{\sum_{k=1}^n P(z|k)P(k)}$$

Coarse descriptor of mood, topic, etc. Used to select prompt to stimulate conversation



## Conversation Trace

Text data from 12 Newsgroups

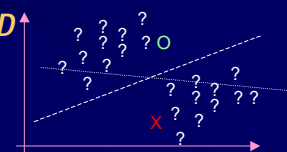


Real-time probabilities of topics (politics, health, religion, etc.)  
Could also detect emotions & situations...

### Wearable Platform



Clustering & MED  
Transductive approach for few labeled.



## References

- 1) Tony Jebara and Alex Pentland. On Reversing Jensen's Inequality. In *Neural Information Processing Systems 13 (NIPS'00)*, Dec. 2000.
- 2) Tony Jebara, Yuri Ivanov, Ali Rahimi and Alex Pentland. Tracking Conversational Context for Machine Mediation of Human Discourse. In *AAAI Fall 2000 Symposium - Socially Intelligent Agents - The Human in the Loop*. Nov. 2000.
- 3) Tony Jebara and Tommi Jaakkola. Feature Selection and Dualities in Maximum Entropy Discrimination. In *16th Conf. Uncertainty in Artificial Intelligence (UAI 2000)*.
- 4) Tommi Jaakkola, Marina Meila, and Tony Jebara. Maximum Entropy Discrimination. In *Neural Information Processing Systems 12 (NIPS'99)*.
- 5) Tony Jebara and Alex Pentland. Action Reaction Learning: Automatic Visual Analysis and Synthesis of Interactive Behaviour. In *1st Intl. Conf. on Computer Vision Systems (ICVS'99)*.
- 6) Bernt Schiele, Nuria Oliver, Tony Jebara and Alex Pentland. An Interactive Computer Vision System, DyPERS: Dynamic Personal Enhanced Reality System. In *1st Intl. Conf. on Computer Vision Systems (ICVS'99)*.
- 7) Tony Jebara and Alex Pentland. Maximum Conditional Likelihood via Bound Maximization and the CEM Algorithm. In *Neural Information Processing Systems 11 (NIPS'98)*.
- 8) Tony Jebara. Action-Reaction Learning: Analysis and Synthesis of Human Behaviour. Master's Thesis, MIT, May 1998.