



EE 6882

Overview of Statistical Models for Video Indexing

Prof. Shih-Fu Chang

Columbia University

TA:
Eric Zavesky

Fall 2007, Lecture 4

Course web site: <http://www.ee.columbia.edu/~sfchang/course/svia>

1



Statistical Paradigm

- Many problems can be posed as pattern recognition
 - Image classification: indoor vs. outdoor? Face?
 - shot boundary detection, story segmentation
 - Is the current point a boundary?
- Statistical models to handle uncertainty and provide flexibility
- Image processing tools available
 - E.g., homework #1
- Rich tools for learning and prediction
 - See course web site
- Increasing data available
 - NIST TREC Video: 300+ hours
 - Consumer and youtube videos





A Very High-Level Stat. Pattern Recog. Architecture

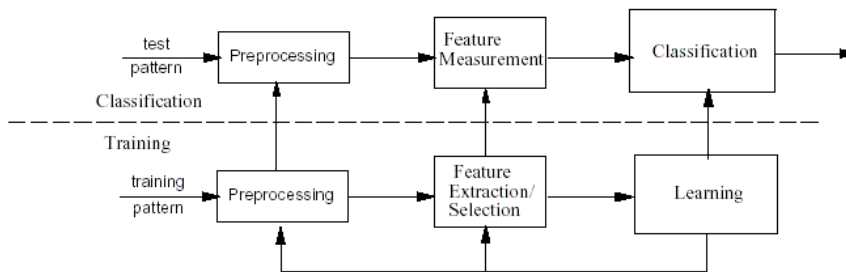


Figure 1: Model for statistical pattern recognition.

(From Jain, Duin, and Mao, SPR Review, '99)



Important issues (1)

- Image/video processing
 - What's the adequate quality, resolution, etc?
- Feature extraction
 - Color, texture, motion, region, shape, interest points, audio, speech, text, etc
- Feature representation
 - Histogram, bag, graph etc
 - Invariance to scale, rotation, translation, view, illumination, ...
 - How to reduce dimensions?



Important issues (2)

- Distance measurement
 - How to measure similarity between images/videos?
 - L1, L2, Mahalanobis, Earth Mover's distance, vector/graph matching
- Classification models
 - Generative vs. discriminative
 - Multi-modal fusion, early fusion vs. late fusion
 - E.g., how to use joint audio-visual features to detect events (dancing, wedding...)
- Efficiency issues
 - how to speed up training and testing processes?
 - How to rapidly build a model for new domains
- Validation and evaluation
 - How to measure performance?
 - Are models generalizable to new domains?



Three related problems

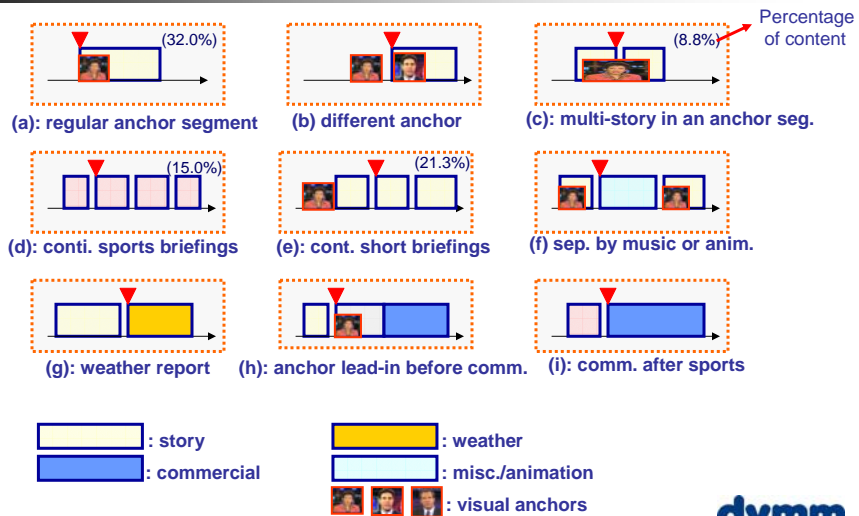
- Retrieval, Ranking
 - Given a query image, find relevant ones
 - May apply rank threshold to decide relevance
- Classification, categorization, detection
 - Given an image x , predict class label y
- Clustering, grouping
 - Group images/videos into clusters of distinct attributes



- An example
 - News story segmentation using multi-modal, multi-scale features

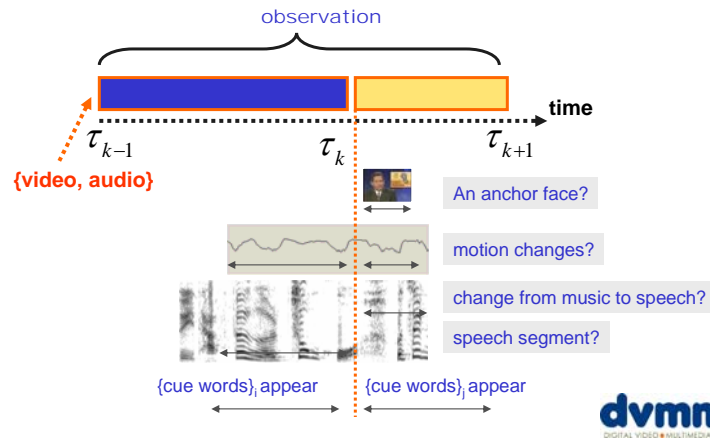


First Understand Data Types and Explore Unique Characteristics



News Story Segmentation

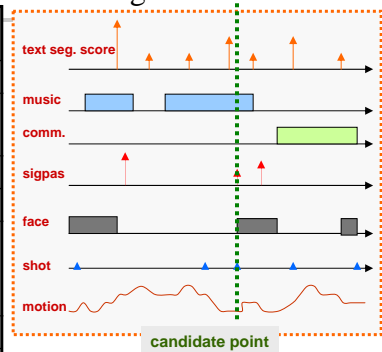
- Objective: a story boundary at time τ_k ?
 - $\tau_k = \{ \text{shot boundaries or significant pauses} \}$



Need to decide how to formulate features

Challenge: diverse features

Modality	Raw Features	Data Type	Value
Video	motion	segment	continuous
	shot boundary	point	binary
	face	segment	continuous
	commercial	segment	binary
Audio	pause	point	continuous
	pitch jump	point	continuous
	significant pause	point	continuous
	musc./spch. disc. spch seg./rapidity	segment	binary continuous
Text	ASR cue terms	point	binary
	V-OCR cue terms	point	binary
	text seg. score	point	continuous
Misc.	combinatorial	point	binary
	sports	segment	binary



One way is to use binary predicate:

if $x > \text{threshold}$, then predict segment boundary ($b=1$)



Choose Model

Maximum entropy model

$$q_{\lambda}(b|x) = \frac{1}{Z_{\lambda}(x)} e^{\sum_i \lambda_i f_i(x,b)}$$

where $f_i(x,b), b \in \{0,1\}$

For example

predicate $f_1 =$ 'anchor face' $f_2 =$ 'significant pause'

if current observation: $face = YES$ $pause = NO$

$$q(b = YES | x) = e^{\lambda_1} / (e^{\lambda_1} + e^{\lambda_2})$$

$$q(b = NO | x) = e^{\lambda_2} / (e^{\lambda_1} + e^{\lambda_2})$$

Classification: if $q(b=YES|x) > 0.5$, then predict YES.



Background: Entropy

- Entropy (bits) $H = -\sum_{i=1}^m p_i \log_2 p_i$

- Kullback-Leibler (K-L) Distance

- A measure of 'distance' between 2 distributions

$$D_{KL}(p(x), q(x)) = \sum_x q(x) \log \frac{q(x)}{p(x)}$$
$$\text{or } = \int_{-\infty}^{\infty} q(x) \log \frac{q(x)}{p(x)} dx$$

- $D_{KL} \geq 0$, and = 0 iff $p(\cdot) = q(\cdot)$
- Not necessarily symmetric, may not satisfy triangular inequality

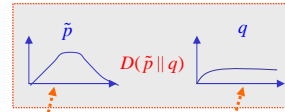


How to Determine the Weights in the Model?

- Estimate $q_\lambda(b|x)$ from training data $T = \{(x_k, b_k)\}$ by minimizing **Kullback-Leibler divergence**, defined as

$$D(\tilde{p} \| q_\lambda) = \sum_x \sum_b \tilde{p}(b, x) \log \frac{\tilde{p}(b|x)}{q_\lambda(b|x)}$$

$$= -\sum_x \sum_b \tilde{p}(x, b) \log q_\lambda(b|x) + \text{constant}(\tilde{p})$$



- Find $\hat{\lambda}_i$ to maximize the 'entropy'

$$L_{\tilde{p}}(q_\lambda) \equiv \sum_x \sum_b \tilde{p}(x, b) \log q_\lambda(b|x)$$

- Iteratively find λ_i

$$\lambda'_i = \lambda_i + \Delta\lambda_i \quad \Delta\lambda_i = \frac{1}{M} \log \left(\frac{\sum_{x,b} \tilde{p}(x, b) f_i(x, b)}{\sum_{x,b} \tilde{p}(x) q_\lambda(b|x) f_i(x, b)} \right)$$

empirical distribution from data

estimated model

- The objective function is convex. So the iterative process can reach the optimum.

The same model used to select features

- Input:** collection of candidate features, training samples, and the desired model size
- Output:** optimal subset of features and their corresponding exponential weights
- Current model q augmented with feature h with weight α ,

$$q_{\alpha, h}(b|x) = \frac{e^{\alpha h(x, b)} q(b|x)}{Z_\alpha(x)}$$

- Select the candidate which improves current model q the most, in each iteration;

$$h^* = \arg \max_{h \in C} \left\{ \sup_{\alpha} \left\{ D(\tilde{p} \| q) - D(\tilde{p} \| q_{\alpha, h}) \right\} \right\}$$

Reduction of divergence

$$= \arg \max_{h \in C} \left\{ \sup_{\alpha} \left\{ L_{\tilde{p}}(q_{\alpha, h}) - L_{\tilde{p}}(q) \right\} \right\}$$

Increase of log-likelihood

Optimal Features (from CNN news video)

* The first 10 "A+V" features automatically discovered for the CNN channel

no	raw feature set	gain	λ	interpretation
1	Anchor Face	0.3879	0.4771	An anchor face segment just starts after the candidate point
2	Significant pause & non-commercial	0.0160	0.7471	A significant pause within the non-commercial section appears in the surrounding observation window.
3	Pause	0.0058	0.2434	An audio pause with the duration larger than 2.0 second appears after the boundary point.
4	Significant pause	0.0024	0.7947	The surrounding observation window has a significant pause with the pitch jump intensity larger than the normalized pitch threshold 1.0 and the pause duration larger than 0.5 second.
5	Speech segment	0.0019	-0.3566	A speech segment before the candidate point
6	Speech segment	0.0015	0.3734	A speech segment starts in the surrounding observation window
7	Commercial	0.0015	1.0782	A commercial starts in 15 to 20 seconds after the candidate point.
8	Speech segment	0.0022	-0.4127	A speech segment ends after the candidate point
9	Anchor face	0.0016	0.7251	An anchor face segment occupies at least 10% of next window
10	Pause	0.0008	0.0939	The surrounding observation window has a pause with the duration larger than 0.25 second.

every modality helps : anchor face, prosody, and speech segment

Issues of this model (Discussion)

■ Features

- Binary predicates reasonable?
- Does it capture the unique characteristics?

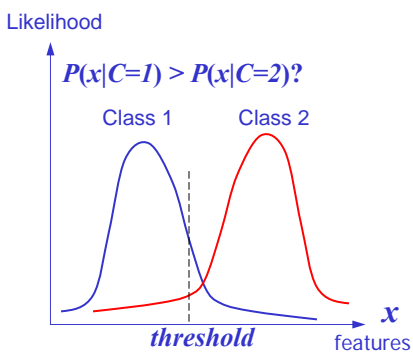
■ binary predicate:
if $x > \text{threshold}$, then predict
segment boundary ($b=1$)

$$q_{\lambda}(b|x) = \frac{1}{Z_{\lambda}(x)} e^{\sum_i \lambda_i f_i(x,b)}$$

■ Models

- Exponential models with linear weights adequate?
- How about the learning algorithm?
 - Enough data to learn the probability models?
 - Speed and complexity

A Broader Perspective: Classification Paradigms



Likelihood

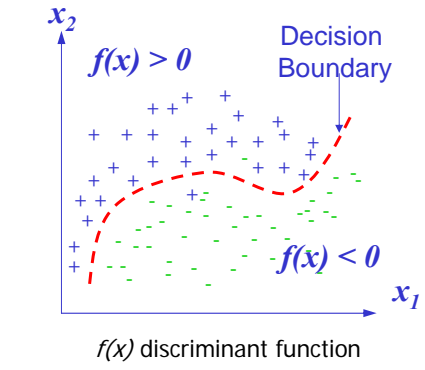
$P(x|C=1) > P(x|C=2)?$

Class 1 Class 2

threshold

x
features

Generative



x_2

$f(x) > 0$ Decision Boundary


$f(x) < 0$


x_1

$f(x)$ discriminant function


Discriminative

Which one does the previous model fall into?





- Generative Models
 - Gaussian Mixture Model

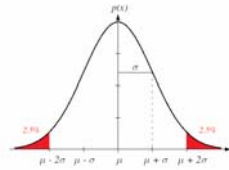


One common issue is to learn probability models

- Gaussian distribution

$$p(x) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{\left(-\frac{1}{2\sigma^2}(x-\mu)^2\right)}$$

- Given the same mean and variance, Gaussian has the max entropy
- Sum of a large number of small, independent random variables approaches Gaussian



$$\Pr[|x - \mu| \leq \sigma] \cong 0.68$$

$$\Pr[|x - \mu| \leq 2\sigma] \cong 0.95$$

$$\Pr[|x - \mu| \leq 3\sigma] \cong 0.997$$

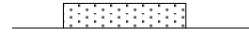
Mahalanobis distance from x to μ

$$r = |x - \mu| / \sigma$$

Entropy of Gaussian:

$$H_{\text{gau}} = 0.5 + \log_2(\sqrt{2\pi}\sigma)$$

Comparison w. uniform



Multi-variate Gaussian

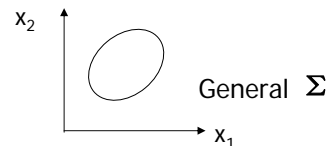
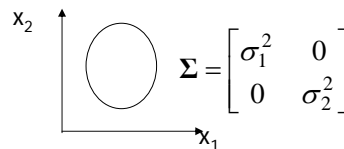
- Multivariate Gaussian, $N(\mu, \Sigma)$

$$p(\mathbf{x}) = \frac{1}{(2\pi)^{D/2} \sqrt{|\Sigma|}} e^{\left(-\frac{1}{2}(\mathbf{x}-\mu)^T \Sigma^{-1}(\mathbf{x}-\mu)\right)}$$

where \mathbf{x}, μ are D -dimensional vectors

Σ : $D \times D$ matrix

$|\Sigma|$ is the determinant of Σ



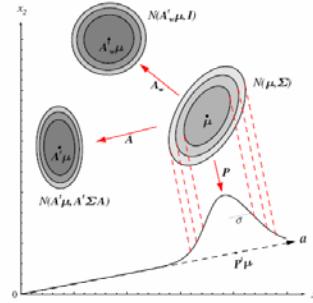
$$\begin{aligned} (\sigma_{ij})^2 &= \Sigma(i, j) = \text{cov}(x(i), x(j)) \\ &= E[(x(i) - \mu(i))(x(j) - \mu(j))] \end{aligned}$$

Effect of Linear Transformation

- Linear transformation of Gaussian

$$y = A^t x \quad y : k \times 1, A : d \times k, x : d \times 1$$

$$y \sim N(A^t \mu, A^t \Sigma A)$$



- Whitening transform

$$\Sigma = \Phi \Lambda \Phi^t \quad (\text{SVD, Eigenvectors})$$

$\Phi : [\phi_1 | \phi_2 | \dots | \phi_d]$ columns are orthogonal ev Λ : diag. matrix of eigenvalues

$$A^t \Sigma A = A^t \Phi \Lambda \Phi^t A = I$$

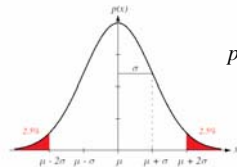
Whitening Trans. $A_w = \Phi \Lambda^{-1/2}$

also PCA Transform

$$y = A_w^t x \sim N(A_w^t \mu, I)$$

Mahalanobis Distance

- Mahalanobis dist in 1-D



$$p(x) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2\sigma^2}(x-\mu)^2}$$

Mahalanobis distance from x to μ

$$r = |x - \mu| / \sigma$$

$$\Pr[|x - \mu| \leq \sigma] = \Pr[r \leq 1] \cong 0.68$$

$$\Pr[|x - \mu| \leq 2\sigma] = \Pr[r \leq 2] \cong 0.95$$

$$\Pr[|x - \mu| \leq 3\sigma] = \Pr[r \leq 3] \cong 0.997$$

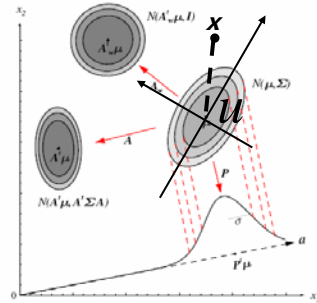
- Multi-Dimensional case

$$\begin{aligned} p(\mathbf{x}) &= \frac{1}{(2\pi)^{D/2} \sqrt{|\Sigma|}} \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^t \Sigma^{-1}(\mathbf{x} - \boldsymbol{\mu})\right) \\ &= \frac{1}{(2\pi)^{D/2} \sqrt{|\Sigma|}} \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^t \Phi \Lambda^{-1} \Phi^t(\mathbf{x} - \boldsymbol{\mu})\right) \\ &= \frac{1}{(2\pi)^{D/2} \sqrt{|\Sigma|}} \exp\left(-\frac{1}{2}(A_w^t(\mathbf{x} - \boldsymbol{\mu}))^t (A_w^t(\mathbf{x} - \boldsymbol{\mu}))\right) \end{aligned}$$

r is the Mahalanobis distance
is also the Euclidean dist in the PCA space

r^2

- Malalanobis distance from point x to the mean of a Gaussian



Gaussian Used In Classification

$x \rightarrow C_j$, if $p(C_j|x) \geq p(C_i|x)$, when $i \neq j$

MAP classifier

$$p(C_j|x) = \frac{\overbrace{p(x|C_j)}^{\text{likelihood}} \overbrace{p(C_j)}^{\text{prior}}}{\underbrace{p(x)}_{\text{evidence}}}$$

ML classification:

if uniform $p(C_j) \Rightarrow C_j = \arg \max_C p(x|C)$

$p(x|C)$ can be modeled by Gaussians

How to Estimate Gaussian Model Parameters?

- $\theta = (\mu, \sigma^2) \quad l = \ln P(x_k | \theta) = -\frac{1}{2} \ln(2\pi\sigma^2) - \frac{1}{2\sigma^2} (x_k - \mu)^2$

$$\nabla_{\theta} l = \begin{pmatrix} \frac{\partial}{\partial \mu} (\ln P(x_k | \theta)) \\ \frac{\partial}{\partial \theta_2} (\ln P(x_k | \theta)) \end{pmatrix} = \begin{bmatrix} \frac{1}{\theta_2} (x_k - \mu) \\ -\frac{1}{2\sigma^2} - \frac{(x_k - \mu)^2}{2(\sigma^2)^2} \end{bmatrix}$$

$$\sum_k \nabla_{\theta} \ln P(x_k | \theta) = 0 \quad \Rightarrow \quad \hat{\mu} = \frac{\sum_k x_k}{n} \quad \hat{\sigma}^2 = \frac{1}{n} \sum_{k=1}^n (x_k - \hat{\mu})^2$$

- Multi-Dimensional $\theta = (\mu, \Sigma)$

$$\Rightarrow \quad \hat{\mu} = (1/n) \sum_k \bar{x}_k \quad \hat{\Sigma} = (1/n) \sum_k (\bar{x}_k - \mu)(\bar{x}_k - \mu)^t$$

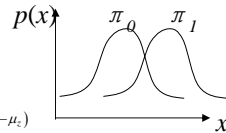
ML estimator: mean -> sample mean, variance -> biased sample variance

Mixture Of Gaussians

- Real distributions seldom follow a single Gaussian
→ mixture of Gaussians

$$p(x) = \sum_z p(x, z) = \sum_z p(z) p(x | z)$$

$$= \sum_z \pi_z N(x | \mu_z, \Sigma_z) = \sum_{z=1}^Z \pi_z \frac{1}{(2\pi)^{D/2} \sqrt{|\Sigma_z|}} e^{-\frac{1}{2}(x-\mu_z)^t \Sigma_z^{-1} (x-\mu_z)}$$



- Given data x_1, \dots, x_N , define log-likelihood:

$$l = \log\left(\prod_n p(x_n)\right) = \sum_{n=1}^N \log(\pi_0 N(x_n | \mu_0, \Sigma_0) + \pi_1 N(x_n | \mu_1, \Sigma_1))$$

- Posterior probability of x being generated by a specific component

$$\text{posteriors} = \tau^i = p(z = i | x, \theta), \quad \theta = \{\mu_0, \Sigma_0, \mu_1, \Sigma_1\}$$

(responsibility of component i)

E-M Optimization Method

$$\boxed{\text{log likelihood}} \quad l(\theta) = \sum_{n=1}^N \log \sum_z p(x_n, z | \theta)$$

- Maximization of $l(\theta)$ directly is hard due to log_of_sum
- Instead, look at

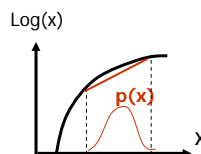
$$\Delta l(\theta) = l(\theta) - l(\theta_t), \quad \theta_t : \text{current estimation of } \theta$$

- Jensen's Inequality

If f is concave, $f(\mathbb{E}\{x\}) \geq \mathbb{E}\{f(x)\}$

$$f(\mathbb{E}\{g(x)\}) \geq \mathbb{E}\{f(g(x))\}$$

e.g., $f(x) = \log(x)$



$$\log\left(\sum_i p_i x_i\right) \geq \sum_i p_i \log(x_i), \text{ where } \sum_i p_i = 1$$

If f is convex, $f(\mathbb{E}\{x\}) \leq \mathbb{E}\{f(x)\}$

Auxiliary Function in E-M

$$\begin{aligned} \Delta l(\theta) &= l(\theta) - l(\theta_t) = \sum_{n=1}^N \log p(x_n | \theta) - \sum_{n=1}^N \log p(x_n | \theta_t) \\ &= \sum_{n=1}^N \log \frac{p(x_n | \theta)}{p(x_n | \theta_t)} = \sum_{n=1}^N \log \sum_z \frac{p(x_n, z | \theta)}{p(x_n | \theta_t)} \end{aligned} \quad \boxed{\text{marginalization}}$$

$$= \sum_{n=1}^N \log \sum_z \frac{p(x_n, z | \theta)}{p(x_n | \theta_t)} \frac{p(x_n, z | \theta_t)}{p(x_n, z | \theta_t)}$$

$$= \sum_{n=1}^N \log \sum_z p(z | x_n, \theta_t) \frac{p(x_n, z | \theta)}{p(x_n, z | \theta_t)}$$

$$\geq \sum_{n=1}^N \sum_z p(z | x_n, \theta_t) \log \frac{p(x_n, z | \theta)}{p(x_n, z | \theta_t)} \quad \boxed{\text{Jensen's inequality}}$$

$$= Q(\theta | \theta_t)$$

- Note there is no log_of_sum. So taking derivative is easier

E-M improves likelihood

- Auxiliary function derived based on Jensen's Inequality,

$$Q(\theta | \theta_t) = \sum_{n=1}^N \sum_z \underbrace{p(z | x_n, \theta_t)}_{\text{expectation over } z \text{ with current } \theta_t} \underbrace{\log p(x_n, z | \theta)}_{\text{joint likelihood of observed \& hidden}} + \text{const}$$

- Now estimate θ_{t+1} by maximizing Q

$$\theta_{t+1} = \arg \max_{\theta} Q(\theta | \theta_t)$$

- So in the expectation step, compute τ_n^z , the 'responsibility' of component z for sample x_n
- In the maximization step, take derivative of Q over θ , and find the new estimate for θ (*Note only sum_of_log is involved*)

EM Always Improves Likelihood

- Why does EM always improve $l(\theta)$?

$$\Delta l(\theta_{t+1}) = l(\theta_{t+1}) - l(\theta_t) \geq Q(\theta_{t+1} | \theta_t)$$

$$Q(\theta_{t+1} | \theta_t) = \max_{\theta} Q(\theta | \theta_t) \geq Q(\theta_t | \theta_t) = 0 \quad \therefore \Delta l(\theta_{t+1}) \geq 0$$

$$Q(\theta | \theta_t) = \sum_{n=1}^N \sum_z \underbrace{p(z | x_n, \theta_t)}_{\text{expectation over } z \text{ with current } \theta_t} \underbrace{\log p(x_n, z | \theta)}_{\text{joint likelihood of observed \& hidden}} + \text{const}$$

- General steps of EM:
 - Define likelihood model with parameters θ
 - Identify hidden variables z
 - Derive the auxiliary function and the E and M equations
 - In each iteration, estimate the posteriors of hidden variables
 - Re-estimate the model parameters. Repeat until stop

Expectation-Maximization (E-M) Solution of GMM

$$Q(\theta | \theta_t) = \sum_{n=1}^N \underbrace{\sum_z p(z | x_n, \theta_t)}_{\text{expectation over } z \text{ with current } \theta_t} \underbrace{\log p(x_n, z | \theta)}_{\text{joint likelihood of observed \& hidden}} + \text{const}$$

- EM for estimating θ and τ_i .

- Follow 'divide and conquer' principle. In iteration step t:

$$\text{Expectation: } \tau_n^{i(t)} = \frac{\pi_i^{(t)} N(x_n | \mu_i^{(t)}, \Sigma_i^{(t)})}{\sum_j \pi_j^{(t)} N(x_n | \mu_j^{(t)}, \Sigma_j^{(t)})}$$

Weight from component i

$$\text{Maximization: } \mu_i^{(t+1)} = \frac{\sum_n \tau_n^{i(t)} x_n}{\sum_n \tau_n^{i(t)}} \quad \Sigma_i^{(t+1)} = \frac{\sum_n \tau_n^{i(t)} (x_n - \mu_i^{(t)}) (x_n - \mu_i^{(t)})^T}{\sum_n \tau_n^{i(t)}}$$

Divide data to each group,
Compute mean and variance
from each group

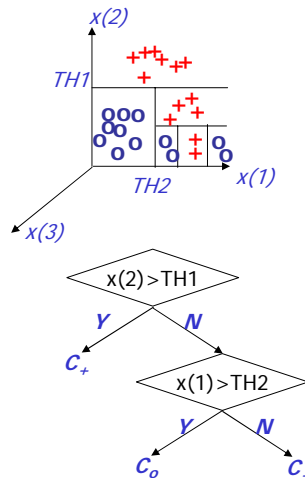
$$\pi_i^{(t+1)} = \frac{\sum_n \tau_n^{i(t)}}{N}$$



- Discriminative Models

Simple Discriminative Classifier

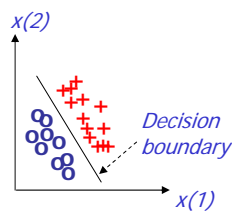
- Find most opportunistic dimension in each step
- Selection criterion
 - Entropy
 - Variance before / after
- Stop criterion
 - Avoid overfitting



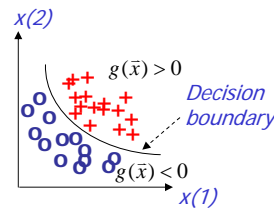
dvmm
DIGITAL VIDEO MULTIMEDIA LAB

Parametric Discriminant Analysis

- Example of discriminant function
 - Linear and quadratic discriminant



$$g(\bar{x}) = ax_1 + bx_2$$

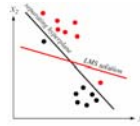


$$g(\bar{x}) = ax_1^2 + bx_2^2 + cx_1x_2$$

dvmm
DIGITAL VIDEO MULTIMEDIA LAB

Linear Discriminant Classifiers

$g(\mathbf{x}) = \mathbf{w}^t \mathbf{x} + w_0 \Rightarrow$ find weight \mathbf{w} and bias w_0



■ Augmented Vector $\mathbf{y} = \begin{bmatrix} 1 \\ \mathbf{x} \end{bmatrix} = \begin{bmatrix} 1 \\ x_1 \\ \vdots \\ x_d \end{bmatrix}$ $\mathbf{a} = \begin{bmatrix} w_0 \\ \mathbf{w} \end{bmatrix} = \begin{bmatrix} w_0 \\ w_1 \\ \vdots \\ w_d \end{bmatrix} \Rightarrow g(\mathbf{x}) = g(\mathbf{y}) = \mathbf{a}^t \mathbf{y}$

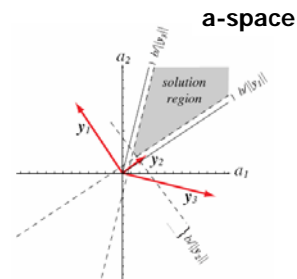
map \mathbf{y} to class ω_1 if $g(\mathbf{y}) > 0$, otherwise class ω_2

■ Design Objective

$\mathbf{a}^t \mathbf{y}_i > b$, if class ω_1 ,

$\mathbf{a}^t \mathbf{y}_i < -b$, if class $\omega_2, \forall \mathbf{y}_i, b > 0$

- Each \mathbf{y}_i defines a half plane in the weight space (\mathbf{a}).
- Note we search weight solutions in the \mathbf{a} -space.



Support Vector Machine (tutorial by Burges '98)

■ Look for separation plane with the highest margin

Decision boundary

$$H_0: \mathbf{w}^t \mathbf{x} + b = 0$$

■ Linearly separable

$$\mathbf{w}^t \mathbf{x}_i + b \geq +1 \quad \forall \mathbf{x}_i \text{ in class } \omega_1 \text{ i.e. } y_i = +1$$

$$\mathbf{w}^t \mathbf{x}_i + b \leq -1 \quad \forall \mathbf{x}_i \text{ in class } \omega_2 \text{ i.e. } y_i = -1$$

Inequality constraints: $(\text{label}_i) \cdot (\mathbf{w}^t \mathbf{x}_i + b) - 1 \geq 0, \forall i$

■ Two parallel hyperplanes defining the margin

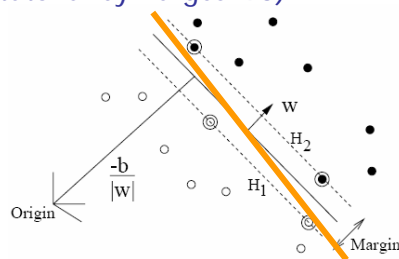
$$\text{hyperplane } H_1(H_+): \mathbf{w}^t \mathbf{x}_i + b = +1$$

$$\text{hyperplane } H_2(H_-): \mathbf{w}^t \mathbf{x}_i + b = -1$$

■ Margin: sum of distances of the closest points to the separation plane

$$\text{margin} = 2 / \|\mathbf{w}\|$$

■ Best plane defined by \mathbf{w} and b



Finding the maximal margin

minimize $\frac{1}{2} \|\mathbf{w}\|^2$ subject to inequality constraints
 $y_i(\mathbf{w}'\mathbf{x}_i + b) - 1 \geq 0 \quad i=1, \dots, l \quad ; y_i \text{ is label}$

- Use the Lagrange multiplier technique for the constrained opt. problem

minimize L_p w.r.t. \mathbf{w} and b

$$L_p = \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{i=1}^l \alpha_i (y_i(\mathbf{w}'\mathbf{x}_i + b) - 1)$$

$$\alpha_i \geq 0$$

$$\frac{dL_p}{d\mathbf{w}} = 0 \Rightarrow \mathbf{w} = \sum_{i=1}^l \alpha_i y_i \mathbf{x}_i$$

$$\frac{dL_p}{db} = 0 \Rightarrow \sum_{i=1}^l \alpha_i y_i = 0$$

maximize L_D w.r.t. \mathbf{w} and b

$$L_D = \sum_{i=1}^l \alpha_i - \frac{1}{2} \sum_{i=1}^l \sum_{j=1}^l \alpha_i \alpha_j y_i y_j \mathbf{x}_i \cdot \mathbf{x}_j$$

with conditions :

$$\sum_{i=1}^l \alpha_i y_i = 0$$

$$\alpha_i \geq 0$$

Quadratic
Programming

Primal Problem

Dual Problem

- Prime and Dual have the same solutions of \mathbf{w} and b

KKT conditions

$$\frac{\partial}{\partial w_\nu} L_P = w_\nu - \sum_i \alpha_i y_i x_{i\nu} = 0 \quad \nu = 1, \dots, d \rightarrow \mathbf{w}^* = \sum_{i=1}^l \alpha_i y_i \mathbf{x}_i$$

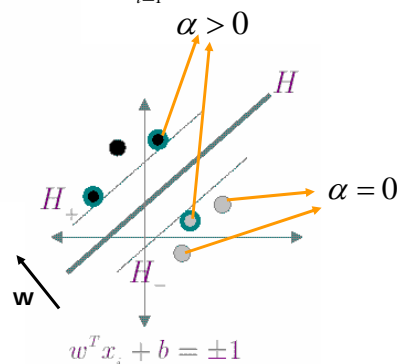
$$\frac{\partial}{\partial b} L_P = - \sum_i \alpha_i y_i = 0$$

$$y_i(\mathbf{x}_i \cdot \mathbf{w} + b) - 1 \geq 0 \quad i = 1, \dots, l$$

$$\alpha_i \geq 0 \quad \forall i$$

$$\alpha_i (y_i(\mathbf{w} \cdot \mathbf{x}_i + b) - 1) = 0 \quad \forall i$$

- Weight sum from positive class = Weight sum from negative class
- Direction of \mathbf{w} : roughly from negative support vectors to positive ones



if $\alpha_i > 0$, \mathbf{x}_i is on H_+ or H_- and is a support vector

Non-separable: not every sample can be correctly classified

- Add slack variables ξ_i

$$\mathbf{x}_i \cdot \mathbf{w} + b \geq +1 - \xi_i \quad \text{for } y_i = +1$$

$$\mathbf{x}_i \cdot \mathbf{w} + b \leq -1 + \xi_i \quad \text{for } y_i = -1$$

$$\xi_i \geq 0 \quad \forall i.$$

if $\xi_i > 1$, then \mathbf{x}_i is misclassified (i.e. training error)

Lagrange multiplier: minimize

$$L_P = \underbrace{\frac{1}{2} \|\mathbf{w}\|^2 + C \sum_i \xi_i}_{\text{New objective function}} - \sum_i \alpha_i \{y_i(\mathbf{x}_i \cdot \mathbf{w} + b) - 1 + \xi_i\} - \sum_i \mu_i \xi_i$$

New objective function

Ensure positivity

KKT Conditions for non-separable Solutions

$$\frac{\partial L_P}{\partial w_\nu} = w_\nu - \sum_i \alpha_i y_i x_{i\nu} = 0$$

$$\frac{\partial L_P}{\partial b} = - \sum_i \alpha_i y_i = 0$$

$$\frac{\partial L_P}{\partial \xi_i} = C - \alpha_i - \mu_i = 0$$

$$y_i(\mathbf{x}_i \cdot \mathbf{w} + b) - 1 + \xi_i \geq 0$$

$$\xi_i \geq 0$$

$$\alpha_i \geq 0$$

$$\mu_i \geq 0$$

$$\alpha_i \{y_i(\mathbf{x}_i \cdot \mathbf{w} + b) - 1 + \xi_i\} = 0$$

$$\mu_i \xi_i = 0$$

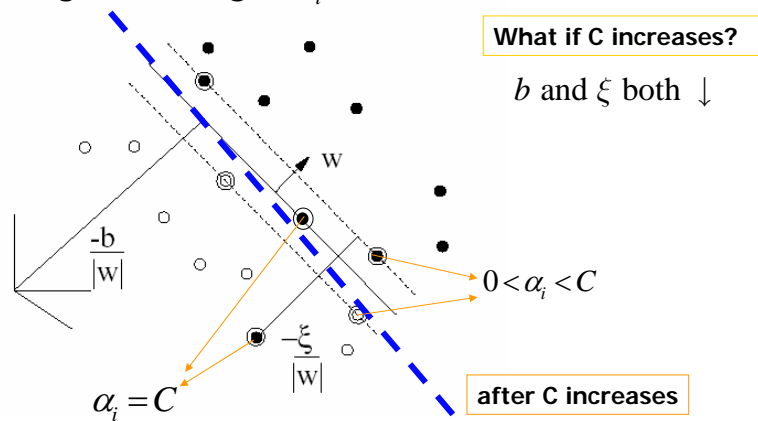
If $0 < \alpha_i < C$, then $\xi_i = 0$: \mathbf{x}_i is on H_1 or H_2

If $\alpha_i = C$,

then $\xi_i > 0$: \mathbf{x}_i is inside the margin or on the wrong side

or $\xi_i = 0$: \mathbf{x}_i is on H_1 or H_2

- All the points located in the margin gap or the wrong side will get $\alpha_i = C$



- When C increases, incorrect samples get more weights
 - try to minimize incorrect samples
 - better training accuracy, but smaller margin
 - less generalization performance

Generalized Linear Discriminant Functions

- Include more than just the linear terms

$$g(\mathbf{x}) = w_0 + \sum_{i=1}^d w_i x_i + \sum_{i=1}^d \sum_{j=1}^d w_{ij} x_i x_j = w_0 + \mathbf{w}'\mathbf{x} + \mathbf{x}'\mathbf{W}\mathbf{x}$$

- Shape of decision boundary
 - ellipsoid, hyperhyperboloid, lines etc

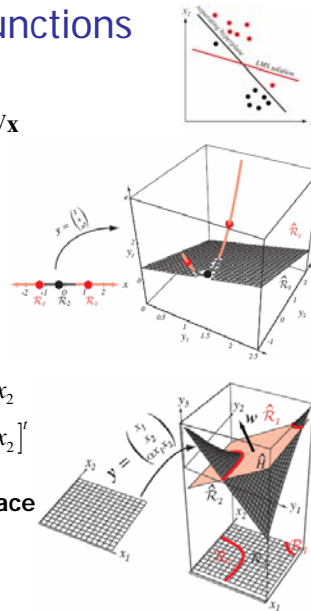
- In general $g(\mathbf{x}) = \sum_{i=1}^d a_i \phi_i(\mathbf{x}) = \mathbf{a}'\Phi$

- Example

$$g(x) = a_1 + a_2 x + a_3 x^2 \quad g(x) = a_1 x_1 + a_2 x_2 + a_3 x_1 x_2$$

$$= [a_1 \ a_2 \ a_3] [1 \ x \ x^2]^T \quad = [a_1 \ a_2 \ a_3] [1 \ x_1 \ x_1 x_2]^T$$

- Data become separable in higher-dimensional space
 - learning parameters in high dimension is hard (curse of dimensionality)
 - instead, try to maximize margins \rightarrow SVM



Non-Linear Space

$\Phi: \mathbf{R}^d \mapsto \mathcal{H}$. Map to a high dimensional space, to make the data separable $\Phi(\mathbf{x}) = \begin{pmatrix} x_1^2 \\ \sqrt{2} x_1 x_2 \\ x_2^2 \end{pmatrix}$

- Find the SVM in the high-dim space

$$g(\mathbf{x}) = \sum_{i=1}^{N_s} \underbrace{\alpha_i y_i \Phi(\mathbf{s}_i)}_w \cdot \Phi(\mathbf{x}) + b$$

- Luckily, we don't have to find $\Phi(\mathbf{s}_i)$ nor $\sum_{i=1}^l \alpha_i y_i \Phi(\mathbf{s}_i)$

- Instead, we define kernel $K(\mathbf{s}_i, \mathbf{x}) = \Phi(\mathbf{s}_i) \cdot \Phi(\mathbf{x})$

$$\Rightarrow g(\mathbf{x}) = \sum_{i=1}^{N_s} \alpha_i y_i K(\mathbf{s}_i, \mathbf{x}) + b$$

- We can use the same method to maximize L_D to find α_i

$$L_D = \sum_{i=1}^l \alpha_i - \frac{1}{2} \sum_{i=1}^l \sum_{j=1}^l \alpha_i \alpha_j y_i y_j \Phi(\mathbf{x}_i) \cdot \Phi(\mathbf{x}_j)$$

$$= \sum_{i=1}^l \alpha_i - \frac{1}{2} \sum_{i=1}^l \sum_{j=1}^l \alpha_i \alpha_j y_i y_j K(\mathbf{x}_i, \mathbf{x}_j)$$

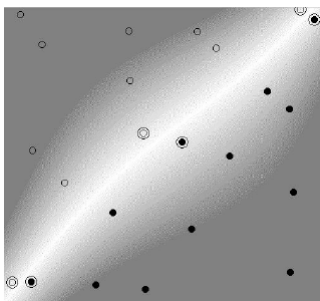
Some popular kernels

$$K(\mathbf{x}, \mathbf{y}) = (\mathbf{x} \cdot \mathbf{y} + 1)^p \quad \text{polynomial}$$

$$K(\mathbf{x}, \mathbf{y}) = e^{-\|\mathbf{x} - \mathbf{y}\|^2 / 2\sigma^2} \quad \text{Gaussian Radial Basis Function (RBF)}$$

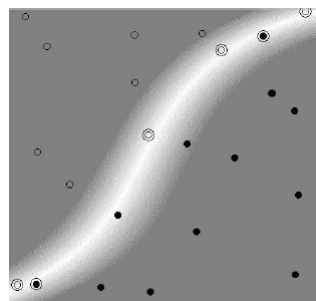
$$K(\mathbf{x}, \mathbf{y}) = \tanh(\kappa \mathbf{x} \cdot \mathbf{y} - \delta) \quad \text{sigmoidal neural network}$$

separable



Cubic polynomial

non-separable

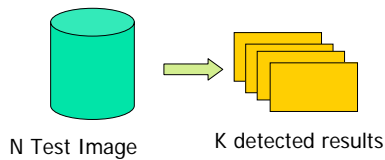




- See the SVM demos (Eric)

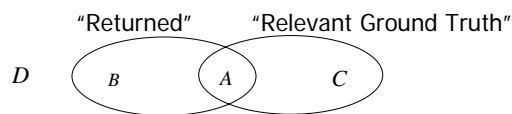


Evaluation



$V_n = 1$ "Relevant"
 0 "Irrelevant" $n = 0 \dots N-1$

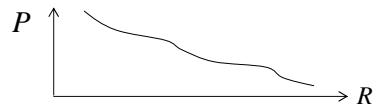
- Detection $A = \sum_{n=0}^{K-1} V_n$
- False Alarms $B = \sum_{n=0}^{K-1} (1 - V_n)$
- Misses $C = (\sum_{n=0}^{N-1} V_n) - A$
- Correct Dismissals $D = (\sum_{n=0}^{N-1} (1 - V_n)) - B$
- Recall $R = A / (A + C)$
- Precision $P = A / (A + B)$
- Fallout $F = B / (B + D)$
- Combined



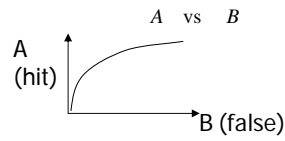
$$F_1 = \frac{P \cdot R}{(P + R) / 2}$$

Evaluation Measures

Precision Recall Curve



2. Receiver Operating Characteristic (ROC Curve)



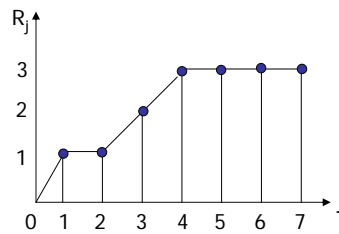
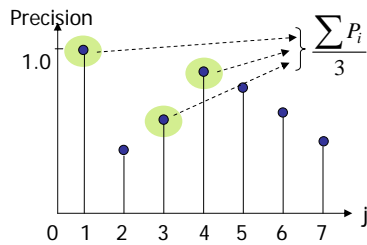
Evaluation Metric: Average Precision

S → Ranked list of data in response to a query

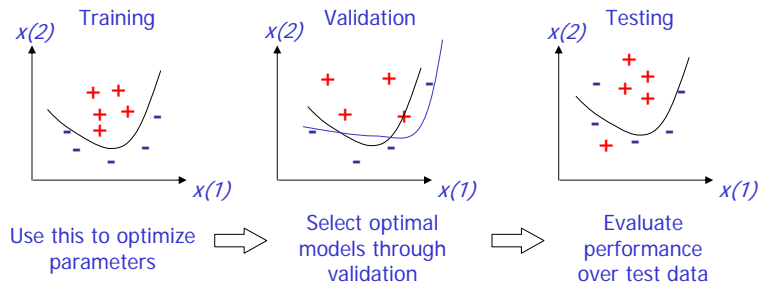
	D_{15}	D_8	D_{63}	D_{21}	D_s
Ground truth	1	0	1	1	0	0	0
Precision	1/1	1/2	2/3	3/4	3/5	3/6	3/7

Average precision: $AP = \frac{1}{R} \sum_{j=1}^s \frac{R_j}{j} I_j$, R : total number of relevant data

AP measures the average of precision values at R relevant data points



Training / Validation / Testing

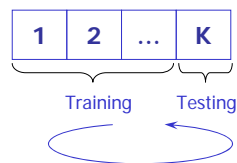


- Appropriate if the same distributions are followed over different sets

dvmm
DIGITAL VIDEO MULTIMEDIA LAB

Training / Validation / Testing (cont.)

- Cross validation, leave-one-out

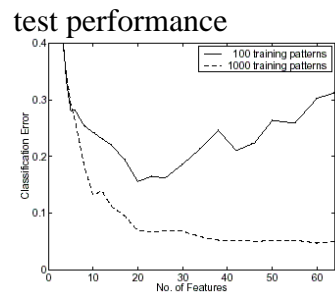
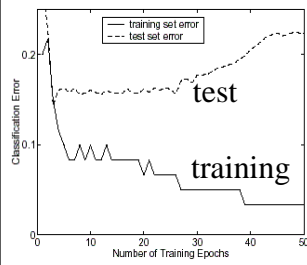
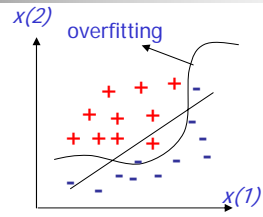


Rotate the choice of the test set and average the performance over runs

dvmm
DIGITAL VIDEO MULTIMEDIA LAB

Curse of Dimensionality and Overtraining

A case of overfitting



Very rough rule of thumb –

$\frac{\# \text{ of training samples per class}}{\text{feature dimension}}$

> 10

