

Advanced Machine Learning & Perception

Instructor: Tony Jebara

Topic 7

- Standard Kernels
- Unusual Kernels
- String Kernels
- Probabilistic Kernels
- Probability Product Kernels

Kernel Machines

- Many algorithms use kernels... all inputs summarized via a matrix of kernel evaluations between pairs of points

or equivalently $K_{ij} = k(x_i, x_j)$

- Generally called Kernel Machines when classifiers

- Examples:

Support Vector Machine

Bayes Point Machine

Relevance Vector Machine

Gaussian Process Regression

Support Vector Machine Regression

Ellipsoidal Kernel Machines

Semidefinite Embedding

Minimum Volume Embedding

etc.

Kernels

- Inner product of mapped inputs: $k(x, y) = \phi(x)^T \phi(y)$
- Example Kernels:

P-th Order Polynomial Kernel: $k(x, y) = (x^T y + 1)^p$

$p = \#$ of twists in SVM decision boundary

RBF Kernel (infinite!): $k(x, y) = \exp\left(-\frac{1}{2\sigma^2} \|x - y\|^2\right)$

- Kernels replace inner products in SVMs, allow nonlinearity:

$$L_D : \max \sum_i \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j k(x_i, x_j) \quad \text{subject to } \alpha_i \in [0, C]$$

$$f(x) = \text{sign}\left(\sum_i \alpha_i y_i k(x, x_i) + b\right)$$

- Solution is still a QP, uses **Gram** matrix K (positive definite)

$$K = \begin{bmatrix} k(x_1, x_1) & k(x_1, x_2) & k(x_1, x_3) \\ k(x_2, x_1) & k(x_2, x_2) & k(x_2, x_3) \\ k(x_3, x_1) & k(x_3, x_2) & k(x_3, x_3) \end{bmatrix} \quad \text{or equivalently } K_{ij} = k(x_i, x_j)$$

Standard Kernels over Vectors

- Accept two input vectors x , y and produce a scalar output:

- Homogeneous Polynomial $k(x, y) = (x^T y)^p$

- Non-Homogeneous Polynomial $k(x, y) = (x^T y + 1)^p$

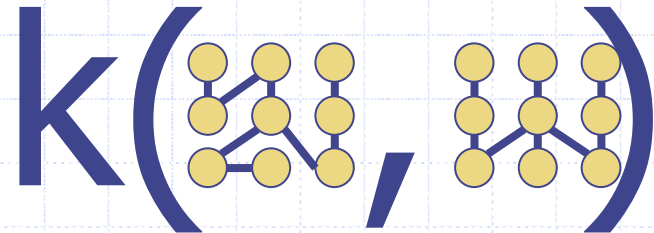
- RBF kernel $k(x, y) = \exp\left(-\frac{1}{2\sigma^2} \|x - y\|^2\right)$

- Tanh Kernel $k(x, y) = \tanh(\kappa x^T y - \delta)$

- Triangular Kernel $k(x, y) = 1 - |(x - y) / d|$

Kernels more Generally

- Can work with non-vectorial inputs!
- Consider the case where the (x, y) pair are
 - strings
 - graphs
 - probability distributions
 - diffusion process



- Kernel just need to satisfy certain basic properties:

1) It's symmetric $k(x, y) = k(y, x)$

2) Satisfies **Mercer's Condition**: $\int k(x, y) g(x) g(y) dx dy \geq 0$
 for any $g(x)$ such that $\int g(x)^2 dx$ is finite

...similar to saying Gram matrix is positive definite

String Kernels

- Inputs are two strings: $X = \text{"aardvark"}$
 $X' = \text{"accent"}$

• Want $k(X, X') = \text{scalar value} = \phi(X)^T \phi(X')$

- One choice for features $\phi(X)$ is the number of times each substring of length 1, 2 and 3 appears in X

$$\phi(X) = [\#a \ \#b \ \#c \ \#d \ \#e \ \#f \ \dots \ \#y \ \#z \ \#aa \ \#ab \ \#ac \ \dots]$$

$$= [3 \ 0 \ 0 \ 1 \ 0 \ 0 \ \dots \ 0 \ 0 \ 1 \ 0 \ 0 \ \dots]$$

$$\phi(X') = [1 \ 0 \ 2 \ 0 \ 1 \ 0 \ \dots \ 0 \ 0 \ 0 \ 0 \ 1 \ \dots]$$

- Can do this efficiently via dynamic programming

String Kernels

- Want $k(X, X') = \text{scalar value} = \phi(X)^T \phi(X')$
- Explicit Kernel can be slow because of many substrings s in our final vocabulary

$\phi_s(X) = \text{number of times substring } s \text{ appears in } X$

$$\phi_a(\text{aardvark}) = 3$$

$$\phi_{ar}(\text{aardvark}) = 2$$

- Implicit Kernel is more efficient

$k(X, X')$: for each substring s of X , count how often s appears in X'

...via dynamic programming in time proportional to the product of the lengths of X and X' . This computes the dot product much more quickly!

Invariance/Independence in Kernels

- Not only do we want kernels that satisfy Mercer...
- We want some kind of invariance: $k(x, y) = k(x, Gy)$

$$k\left(\begin{array}{ccc} \bullet & \bullet & \bullet \\ \bullet & \bullet & \bullet \\ \bullet & \bullet & \bullet \end{array}, \begin{array}{ccc} \bullet & \bullet & \bullet \\ \bullet & \bullet & \bullet \\ \bullet & \bullet & \bullet \end{array}\right) = k\left(\begin{array}{ccc} \bullet & \bullet & \bullet \\ \bullet & \bullet & \bullet \\ \bullet & \bullet & \bullet \end{array}, \begin{array}{ccc} \bullet & \bullet & \bullet \\ \bullet & \bullet & \bullet \\ \bullet & \bullet & \bullet \end{array}\right)$$

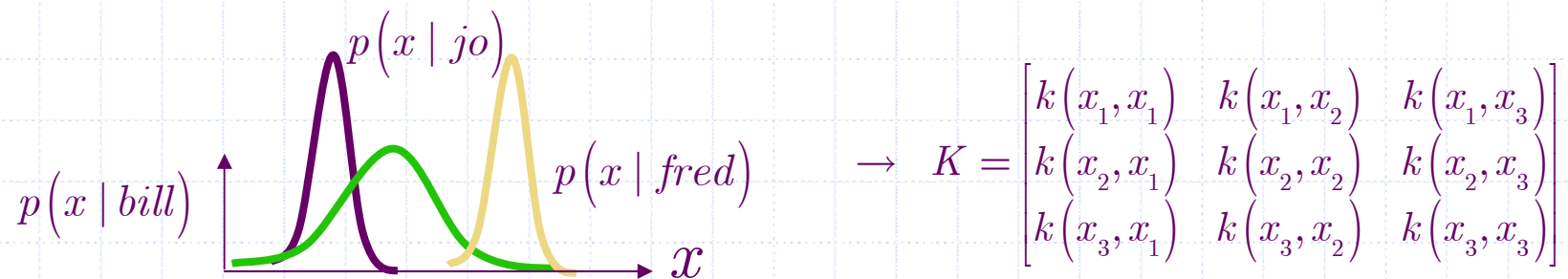
- Another consideration is independence in kernels...
- For example, input has 2 dimensions that are independent

$$x = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \quad y = \begin{bmatrix} y_1 \\ y_2 \end{bmatrix} \quad k(x, y) = k_1(x_1, y_1) + k_2(x_2, y_2)$$

- How to get such desiderata in general? à Open Problem
- One way is to work with probabilistic kernels...

Probabilistic Kernels

- We will consider two cases of probabilistic kernels...
Fisher Kernel and Probability Product Kernel
- For example, how similar are a set of people who have $p(x)$ each person's weight distribution over 365 days



- Fisher Kernel → gets an affinity between p and p' by approximating KL-divergence $KL(p \| p') = \int p \log \frac{p}{p'} dx$
- Probability Product Kernel → gets an affinity by using the Hellinger-divergence $H(p \| p') = \frac{1}{2} \int (\sqrt{p} - \sqrt{p'})^2 dx$

Fisher Kernels

Fisher Kernel: approximate distance between two generative models on statistical manifold using Kullback Leibler
Approximate KL by quadratic form local tangent space at ML estimate

$$KL(p \| p') = \int p \log \frac{p}{p'} dx$$

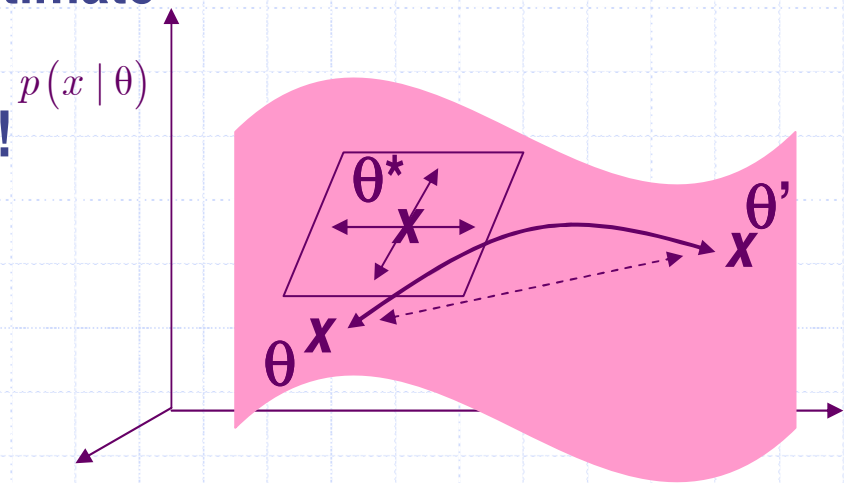
$$KL \approx \frac{1}{2} (\theta - \theta')^T I_{\theta^*} (\theta - \theta')$$

affinity from distance= kernel !!!
via Fisher Info & gradients

$$K(\chi, \chi') \approx U_x I_{\theta^*}^{-1} U_{x'}$$

$$U_x = \nabla_{\theta^*} \log p(\chi | \theta)$$

Fisher Info $I_{\theta^*}(i, j) = \int p(\chi | \theta^*) \frac{\partial \log p(\chi | \theta)}{\partial \theta_i} \frac{\partial \log p(\chi | \theta)}{\partial \theta_j} d\chi$



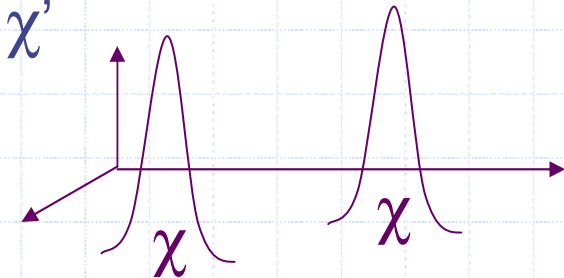
Probability Product Kernels

Computing the kernel: given inputs χ and χ'

1) Estimate Densities (i.e. ML):

$$\chi \rightarrow p(x) = p(x | \theta)$$

$$\chi' \rightarrow p'(x) = p(x | \theta')$$



2) The natural affinity for Hellinger is

known: Bhattacharyya Affinity

$$k(\chi, \chi') = \int \sqrt{p} \sqrt{p'} dx$$

Probability Product Kernel: $k_{\rho}(\chi, \chi') = \int (p)^{\rho} (p')^{\rho} dx$

Bhattacharyya Kernel: $\rho = \frac{1}{2}$

Expected Likelihood Kernel: $k(\chi, \chi') = E_p[p'] = E_{p'}[p]$
 $\rho = 1$

Gaussian Product Kernels

Gaussian Mean: (any ρ)
(continuous)

If $\mu = \chi$ $\mu' = \chi'$ get RBF Kernel
Fisher here is linear

$$k(\chi, \chi') = \int N^\rho(x | \mu) N^\rho(x | \mu') dx \\ \propto \exp\left(-\|\mu - \mu'\|^2 / \tilde{\sigma}\right)$$

Gaussian Covariance: (any ρ)

$$k(\chi, \chi') = \int N^\rho(x | \mu, \Sigma) N^\rho(x | \mu', \Sigma') dx \\ \propto |\Sigma|^{-\rho/2} |\Sigma'|^{-\rho/2} |\Sigma^\dagger|^{1/2} \exp\left(-\frac{\rho}{2} \mu^T \Sigma^{-1} \mu - \frac{\rho}{2} \mu'^T \Sigma'^{-1} \mu' + \frac{1}{2} \mu^{\dagger T} \Sigma^\dagger \mu^\dagger\right) \\ \text{where } \Sigma^\dagger = (\rho \Sigma^{-1} + \rho \Sigma'^{-1})^{-1} \text{ and } \mu^\dagger = \rho \Sigma^{-1} \mu + \rho \Sigma'^{-1} \mu'$$

Fisher here is just a quadratic kernel...

Multinomial Product Kernels

**Bernoulli:
(binary)**

$$p(x | \theta) = \prod_{d=1}^D \gamma_d^{x_d} (1 - \gamma_d)^{1-x_d}$$

$$k(\chi, \chi') = \int (p(x))^{\rho} (p'(x))^{\rho} dx$$

$$= \prod_{d=1}^D \left[(\gamma_d \gamma_d')^{\rho} + (1 - \gamma_d)^{\rho} (1 - \gamma_d')^{\rho} \right]$$

**Multinomial:
(discrete)**

$$p(x | \theta) = \prod_{d=1}^D \alpha_d^{x_d}$$

$$k(\chi, \chi') = \sum_{d=1}^D (\alpha_d \alpha_d')^{\rho}$$

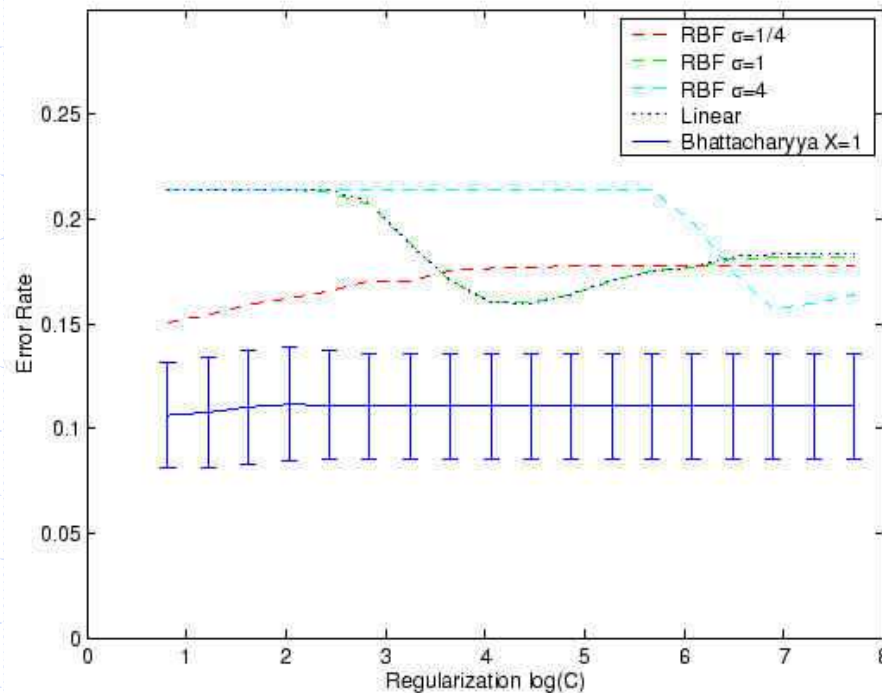
For multinomial counts (for N words):

$$k(\chi, \chi') = \left[\sum_{d=1}^D (\alpha_d \alpha_d')^{1/2} \right]^N$$

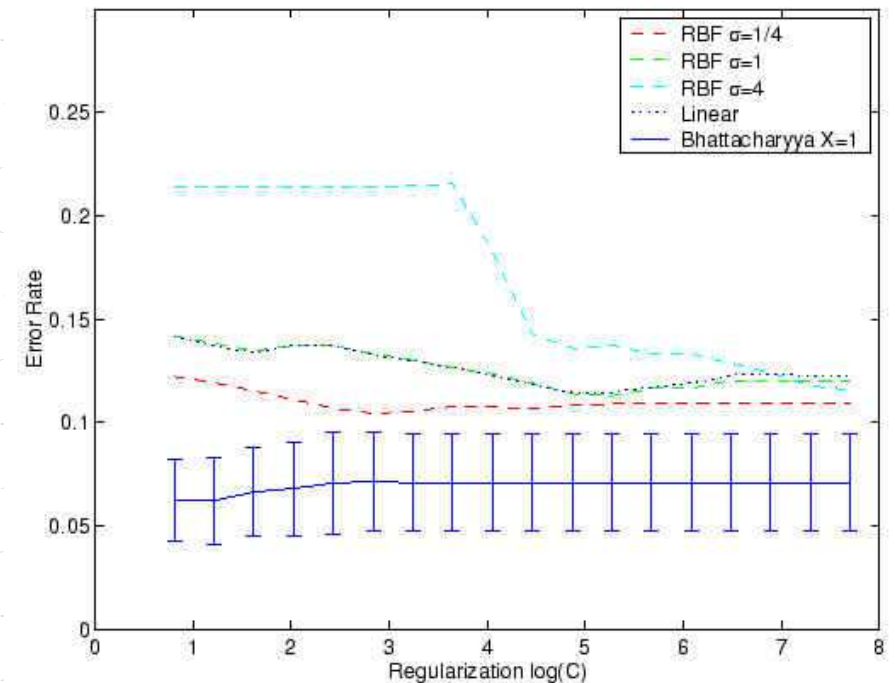
Fisher for Multinomial is linear

Multinomial Product Kernels

WebKB dataset: Faculty vs. student web page SVM kernel classifier
 20-Fold Cross-Validation, 1641 student & 1124 faculty, ...
 Use Bhattacharyya Kernel on multinomial (word frequency)



Training Set Size = 77



Training Set Size = 622

Exponential Family Kernels

Exponential Family: Many Properties, includes Gaussian Mean, Covariance, Multinomial, Binomial, Poisson, Exponential, Gamma, Bernoulli, Dirichlet, ...

$$p(x | \theta) = \exp(A(x) + T(x)^T \theta - K(\theta))$$

Maximum likelihood is straightforward: $\frac{\partial}{\partial \theta} K(\theta) = \frac{1}{n} \sum_n T(x_n)$
All have above form but different $A(x)$, $T(x)$, convex $K(\theta)$

Family	$A(X)$	$K(\theta)$
Gaussian (mean)	$-\frac{1}{2} X^T X - \frac{D}{2} \log(2\pi)$	$\frac{1}{2} \theta^T \theta$
Gaussian (variance)	$-\frac{1}{2} \log(2\pi)$	$-\frac{1}{2} \log(\theta)$
Multinomial	$\log(\Gamma(\eta + 1)) - \log(\nu)$	$\eta \log(1 + \sum_{d=1}^D \exp(\theta_d))$
Exponential	0	$-\log(-\theta)$
Gamma	$-\exp(X) - X$	$\log \Gamma(\theta)$
Poisson	$\log(X!)$	$\exp(\theta)$

Exponential Family Kernels

Compute the Bhattacharyya Kernel for the e-family:

$$p(x | \theta) = \exp(A(x) + T(x)^T \theta - K(\theta))$$

Analytic solution for e-family:

$$\begin{aligned} k(\chi, \chi') &= \int p(x | \theta)^{1/2} p(x | \theta')^{1/2} dx \\ &= \exp\left(K\left(\frac{1}{2}\theta + \frac{1}{2}\theta'\right) - \frac{1}{2}K(\theta) - \frac{1}{2}K(\theta')\right) \end{aligned}$$

Only depends on convex cumulant-generating function $K(\theta)$

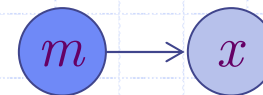
Meanwhile, Fisher Kernel is always linear in sufficient stats...

$$\begin{aligned} U_x &= \nabla_{\theta^*} \log p(\chi | \theta) = T(\chi) - \nabla_{\theta^*} K(\theta) \\ U_{\chi} I_{\theta^*}^{-1} U_{\chi'} &= (T(\chi) - g)^T I_{\theta^*}^{-1} (T(\chi') - g) \end{aligned}$$

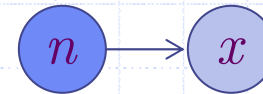
Mixture Product Kernels

Beyond Exponential Family: Mixtures and Hidden Variables
Easier for $\rho=1$ Expected Likelihood kernel...

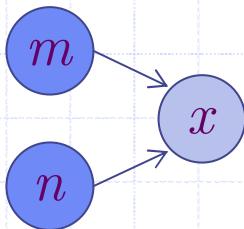
$$\chi \rightarrow p(x) = \sum_{m=1}^M p(m)p(x|m)$$



$$\chi' \rightarrow p'(x) = \sum_{n=1}^N p'(n)p'(x|n)$$



$$\begin{aligned} k(\chi, \chi') &= \int p(x)p'(x)dx \\ &= \sum_{m=1}^M \sum_{n=1}^N p(m)p'(n) \int p(x|m)p'(x|n)dx \\ &= \sum_{m=1}^M \sum_{n=1}^N p(m)p'(n)k_{m,n}(\chi, \chi') \end{aligned}$$



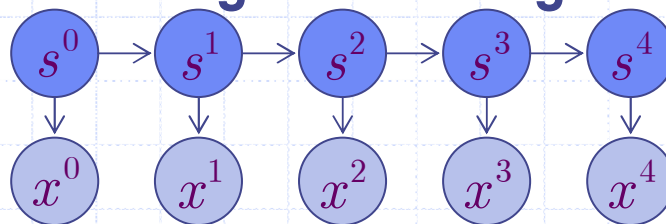
Use $M \cdot N$ subkernel evaluations
from our previous repertoire

HMM Product Kernels

Hidden Markov Models: (sequences)

$$p(x) = \sum_S p(s_0) p(x_0 | s_0) \prod_{t=1}^T p(s_t | s_{t-1}) p(x_t | s_t)$$

of hidden configurations large



$$\# \text{ configs} = |s|^T$$

Kernel: $k(\chi, \chi') = \int p(x) p'(x) dx$

$$= \sum_S \sum_U p(S) p'(U) \int p(x | S) p'(x | U) dx$$

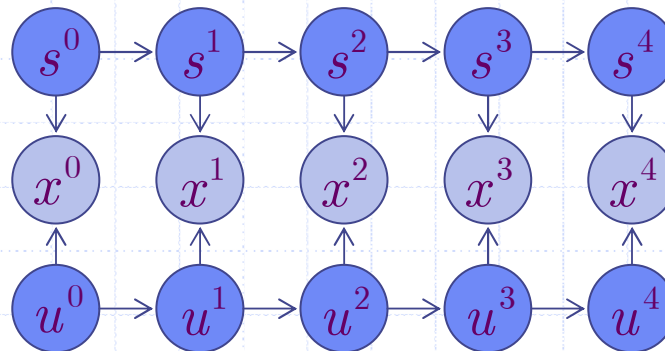
$$= \sum_S \sum_U p(S) p'(U) k_{S,U}(x, x')$$

Do we need to consider raw cross-product of hidden? $O(|s|^T \times |u|^T)$

HMM Product Kernels

$$\begin{aligned}
 k(\chi, \chi') &= \sum_s \sum_U \prod_t p(s_t | s_{t-1}) p'(u_t | u_{t-1}) \int p(x_t | s_t) p'(x_t | u_t) dx_t \\
 &= \sum_{s,U} \prod_t p(s_t | s_{t-1}) p'(u_t | u_{t-1}) k_{s_t, u_t} \\
 &= \sum_{s,U} p(s_0) p'(u_0) \psi(s_0, u_0) \prod_t p(s_t | s_{t-1}) p'(u_t | u_{t-1}) \psi(s_t, u_t)
 \end{aligned}$$

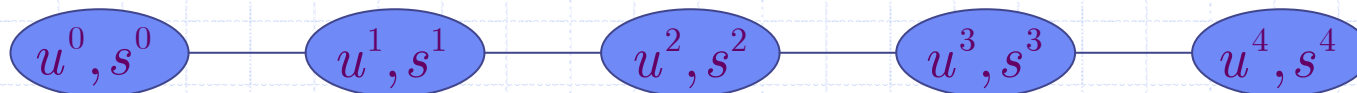
Take advantage of
structure in HMMs
via Bayesian network



Only compute subkernels for common parents

Evaluate total of $O(T|s||u|)$ subkernels

Form +ve clique potential functions, sum via junction tree algorithm



Probability Product Kernel

- PPK for 2 Gaussian HMMs with states S & U
- Get $|S| \times |U|$ interaction table between all pairs of emissions

$$\Psi(i, j) = \frac{|\Sigma^\dagger|^{1/2}}{|\Sigma_i|^{\beta/2} |\Sigma_j|^{\beta/2} \exp(-\frac{\beta}{2} (\mu_i^T \Sigma_i^{-1} \mu_i + \mu_j^T \Sigma_j^{-1} \mu_j - \mu_i^{\dagger T} \Sigma^\dagger \mu_i^\dagger))}$$

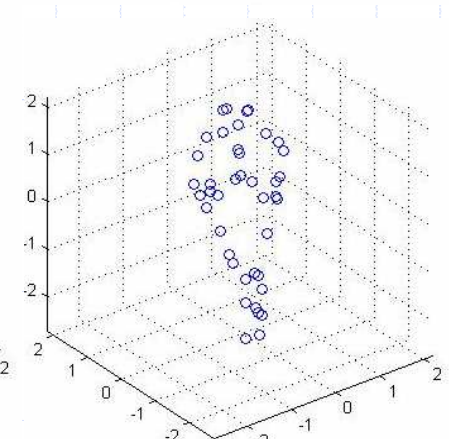
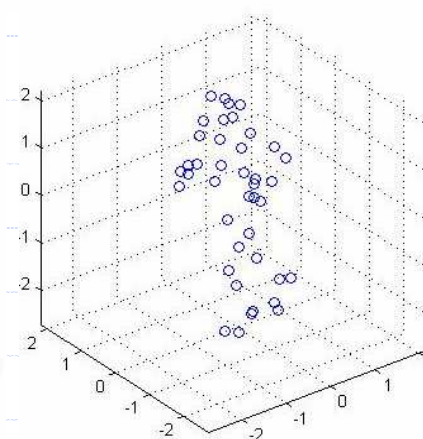
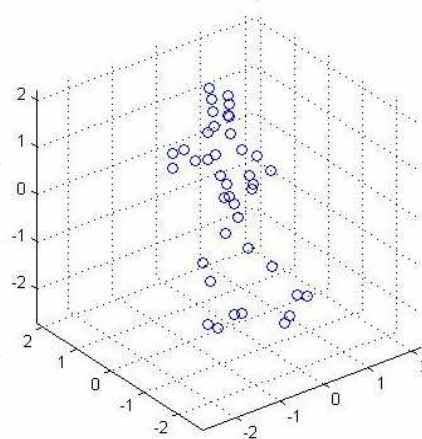
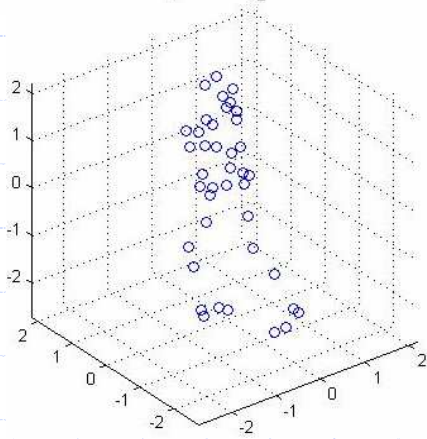
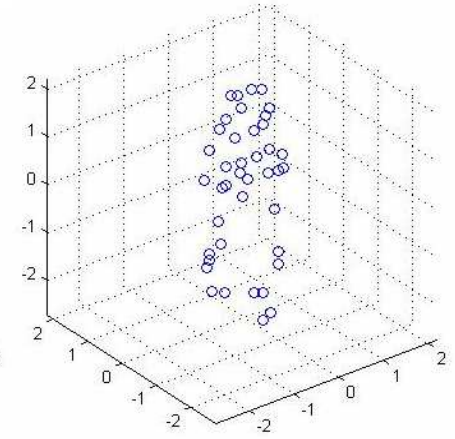
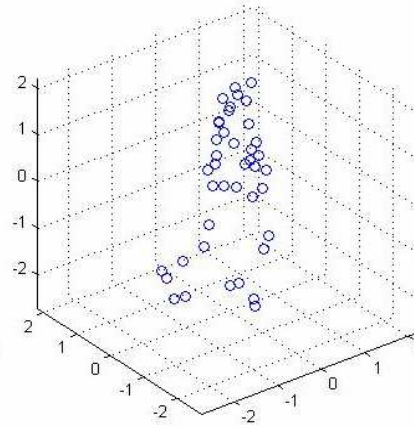
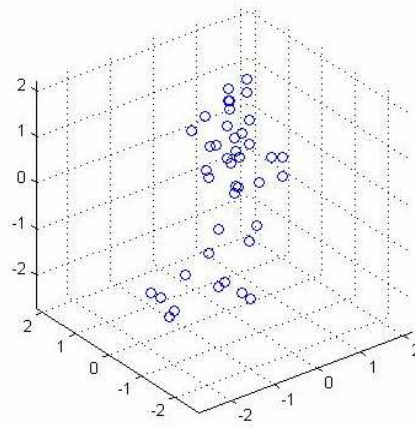
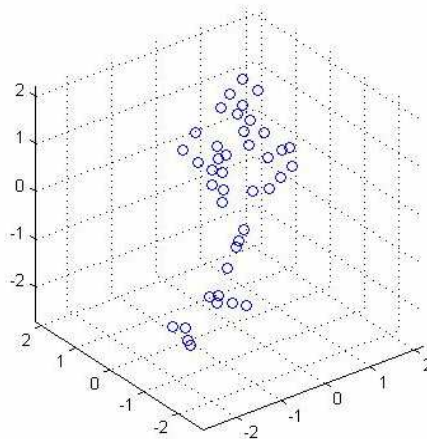
$$\Sigma^\dagger = (\Sigma_i^{-1} + \Sigma_j^{-1})^{-1} \text{ and } \mu_i^\dagger = \Sigma_i^{-1} \mu_i + \Sigma_j^{-1} \mu_j.$$

- Then simple pseudo-code:

$$\begin{aligned} & \Phi(q_0, q'_0) = p(q_0)^\beta p'(q'_0)^\beta \leftarrow \text{state prior} \\ & \text{for } t = 1 \dots T \\ & \quad \Phi(q_t, q'_t) = \sum_{q_{t-1}} \sum_{q'_{t-1}} p(q_t | q_{t-1})^\beta p(q'_t | q'_{t-1})^\beta \Psi(q_{t-1}, q'_{t-1}) \Phi(q_{t-1}, q'_{t-1}) \leftarrow \text{transition} \\ & \text{end} \\ & \tilde{\mathcal{K}}(\theta, \theta') = \sum_{q_T} \sum_{q'_T} \Phi(q_T, q'_T) \Psi(q_T, q'_T) \end{aligned}$$

Visualizing HMMs

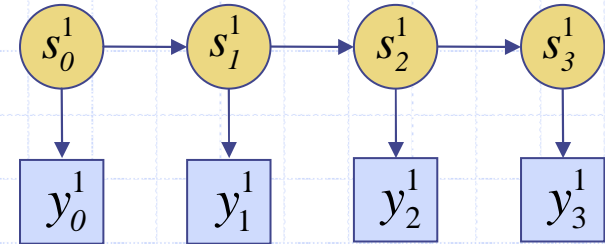
- Rotating walk/run motion data
- Each sequence train a 2-state HMM



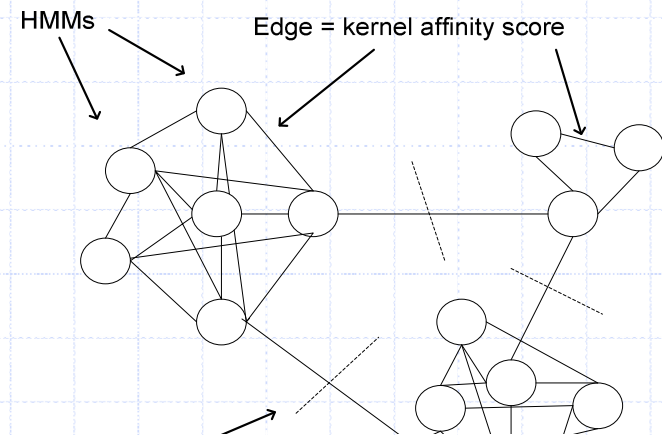
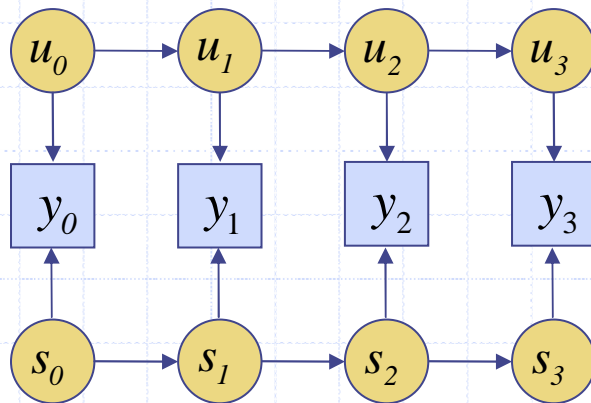
Visualizing HMMs

- 1) For each time series, parametrically learn an HMM

$$\max p(y_0^1, y_1^1, y_2^1, y_3^1)$$



- 2) Compute kernel between all pairs of HMMs

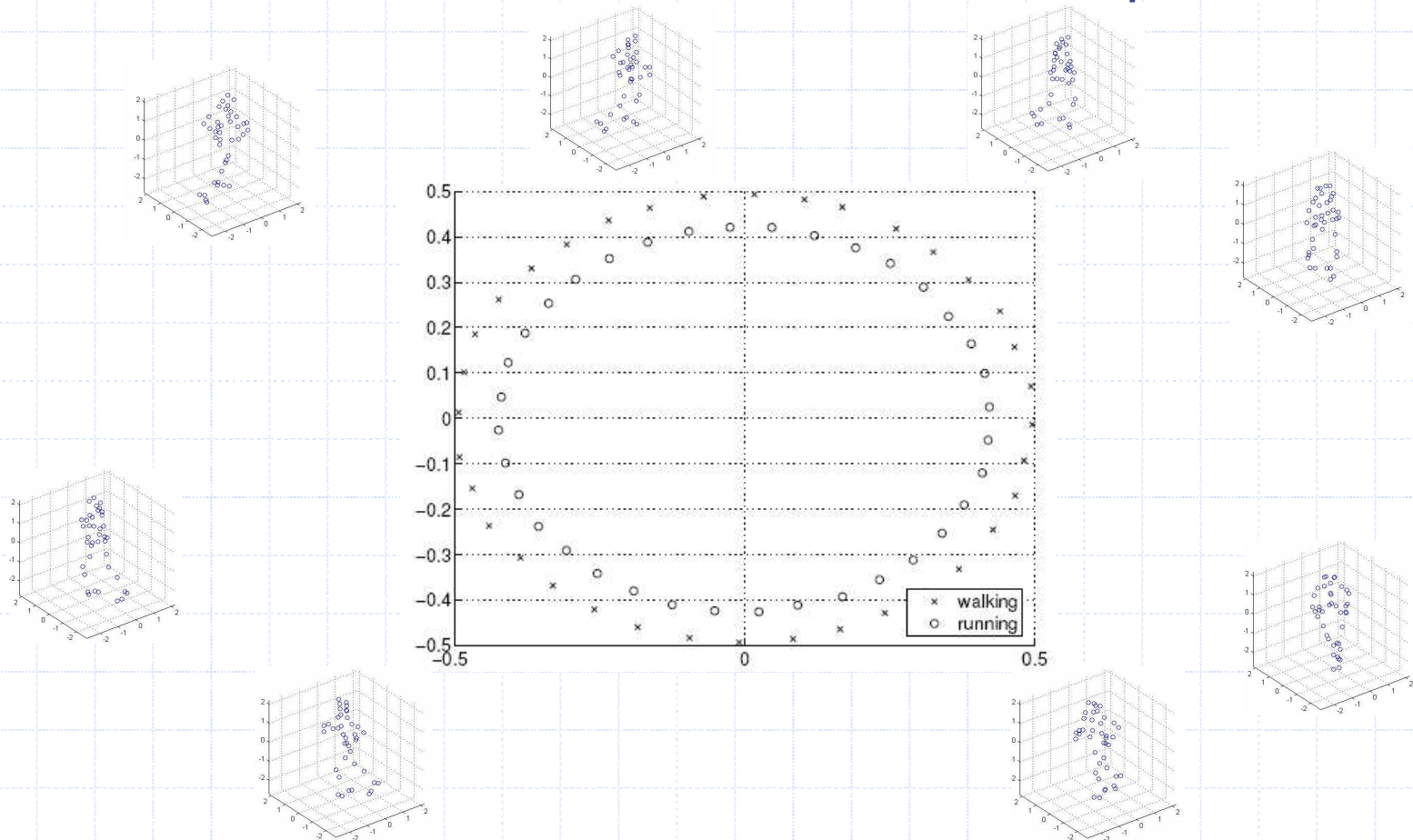


$$k(HMM_1, HMM_2) = \int p(Y | HMM_1) p(Y | HMM_2) dY$$

- 3) Run embedding tool on kernel Matrix (MDS, SDE, MVE)

Visualizing HMMs

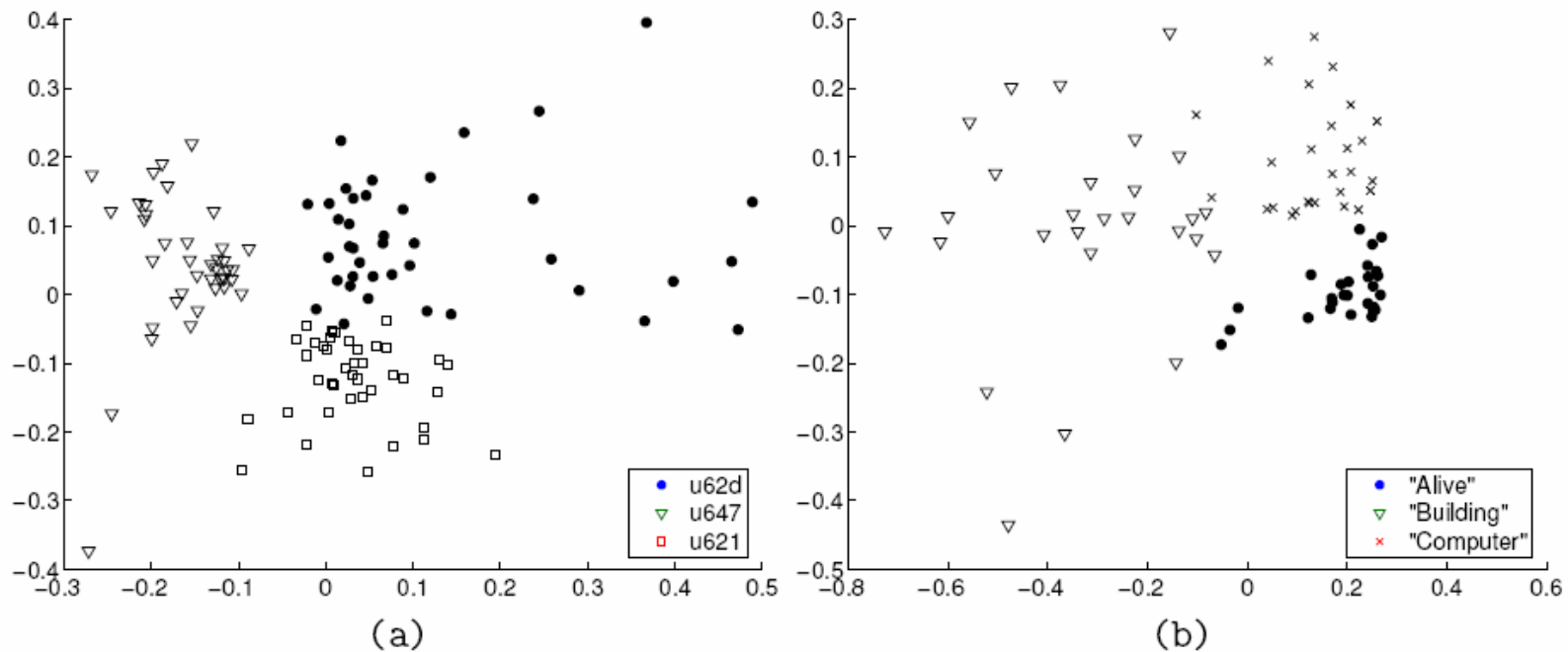
- The rotated walks and runs appear as two concentric rings when we visualize the kernel matrix from all pairs of HMMs



Visualizing HMMs

Embed several HMMs using probability product kernel:

- Arabic Dataset: each character is a time series datum of spatial coordinates.
- Sign Language Dataset: each datum is a time series of hand movement coordinates



Sampling Product Kernels

Sampling: approximate kernel

$$k(\chi, \chi') = \int p(x) p'(x) dx$$

By definition, generative models can:

- 1) Generate a Sample
- 2) Compute Likelihood of a Sample

Thus, approximate probability product via sampling:

$$\begin{aligned} k(\chi, \chi') &= k(p, p') \\ &= \frac{\beta}{N} \sum_{\substack{x_i \sim p(x) \\ i=1 \dots N}} p'(x_i) + \frac{1-\beta}{N'} \sum_{\substack{x_i' \sim p'(x) \\ i=1 \dots N'}} p(x_i') \end{aligned}$$

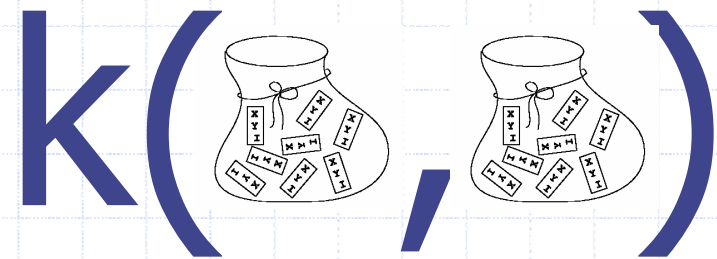
Beta controls how much sampling from each distribution...

Gaussian Product Kernels

Instead of a single χ & χ'
 Construct p & p' from many χ & χ'
 i.e. use bag of vectors

$$\{\chi_1, \dots, \chi_N\} \rightarrow p(x)$$

$$\{\chi_1', \dots, \chi_{N'}'\} \rightarrow p'(x)$$



$$\mu = E\{\chi_i\} \quad \Sigma = E\{(\chi_i - \mu)(\chi_i - \mu)^T\}$$



$$\sim N(x | \mu, \Sigma)$$



$$\sim N(x | \mu', \Sigma')$$

$$k(\chi, \chi') = |\Sigma|^{-\rho/2} |\Sigma'|^{-\rho/2} |\Sigma^\dagger|^{1/2} \exp\left(-\frac{\rho}{2} \mu^T \Sigma^{-1} \mu - \frac{\rho}{2} \mu'^T \Sigma'^{-1} \mu' + \frac{1}{2} \mu^{\dagger T} \Sigma^\dagger \mu'\right)$$

Kernelized Gaussian Product

Bhattacharyya affinity on Gaussians on bags $\{\chi, \dots\}$ & $\{\chi', \dots\}$

Invariant: to order of tuples in each bag

But too simple: overlap of two Gaussian distributions on images



Need more detail than mean and covariance of pixels...

Use Kernel Trick again when computing Gaussian mean & covariance

$$\kappa(\chi, \chi') = \phi(\chi)^T \phi(\chi')$$

Never compute outerproducts, use kernel, i.e. infinite RBF:

$$\kappa(\chi, \chi') = \exp\left(-\frac{1}{2\sigma^2} \|\chi - \chi'\|^2\right)$$

Compute mini-kernel between each pixel in a given image...

Gives kernelized or augmented Gaussian μ and Σ via Gram

Kernelized Gaussian Product

Previously: $\mu = E\{\chi\}$ $\Sigma = E\{(\chi - \mu)(\chi - \mu)^T\}$

Now have:

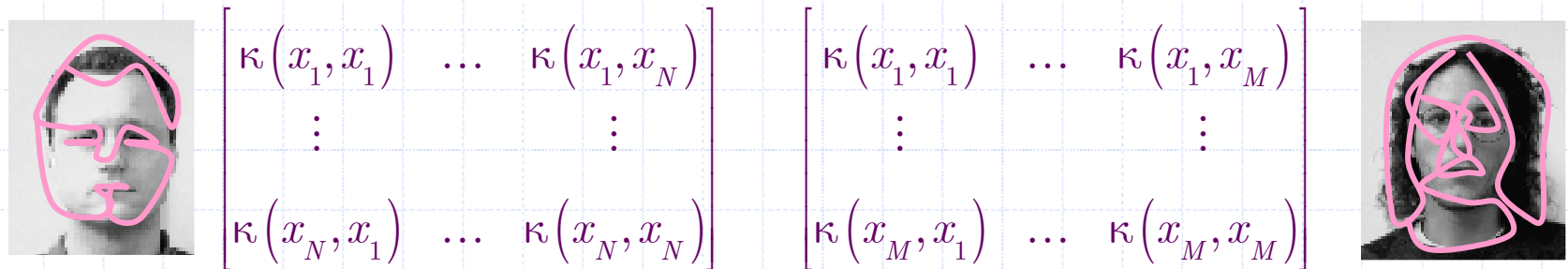
$$\mu = E\{\phi(\chi)\} \quad \Sigma = E\{(\phi(\chi) - \mu)(\phi(\chi) - \mu)^T\}$$

Still invariant to order of pixels

Compute Hilbert Gaussian's mean & covariance of each image

bag or image is $N \times N$ pixel Gram matrix using kappa kernel

Use kernel PCA to regularize infinite dimensional RBF Gaussian

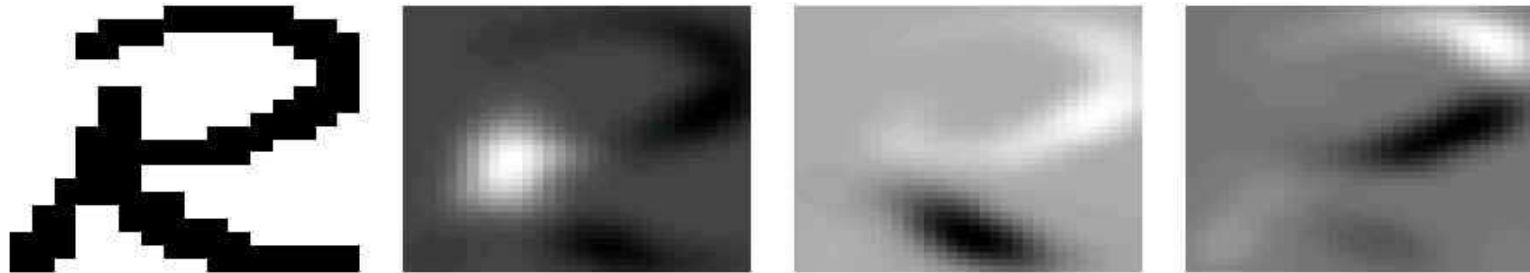


$$\begin{bmatrix} \kappa(x_1, x_1) & \dots & \kappa(x_1, x_N) \\ \vdots & & \vdots \\ \kappa(x_N, x_1) & \dots & \kappa(x_N, x_N) \end{bmatrix} \quad \begin{bmatrix} \kappa(x_1, x_1) & \dots & \kappa(x_1, x_M) \\ \vdots & & \vdots \\ \kappa(x_M, x_1) & \dots & \kappa(x_M, x_M) \end{bmatrix}$$

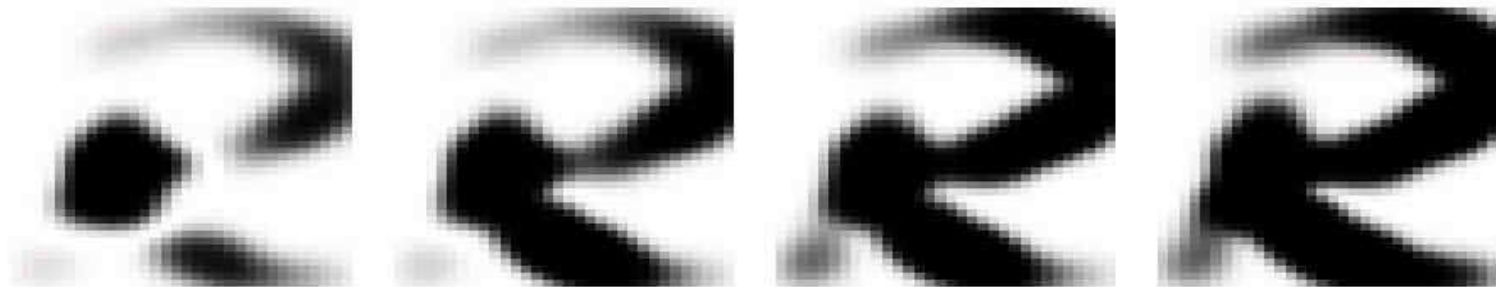
$$k(\phi(\chi) \parallel \phi(\chi')) \propto |\Sigma|^{-\rho/2} |\Sigma'|^{-\rho/2} |\Sigma^\dagger|^{1/2} \exp\left(-\frac{\rho}{2} \mu^T \Sigma^{-1} \mu - \frac{\rho}{2} \mu'^T \Sigma'^{-1} \mu' + \frac{1}{2} \mu^{\dagger T} \Sigma^\dagger \mu^\dagger\right)$$

Puts all dimensions (X,Y,I) on an equal footing

Kernelized Gaussian Product



Letter 'R' with 3 KPCA Components of RBF Kernel



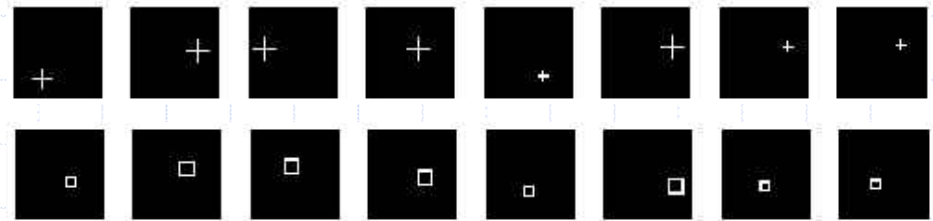
Reconstruction of Letter 'R' with 1-4 KPCA with RBF Kernel



Reconstruction of Letter 'R' with 3 KPCA with RBF Kernel + Smoothing

Kernelized Gaussian Product

100 40x40 monochromatic images of crosses & squares translating & scaling

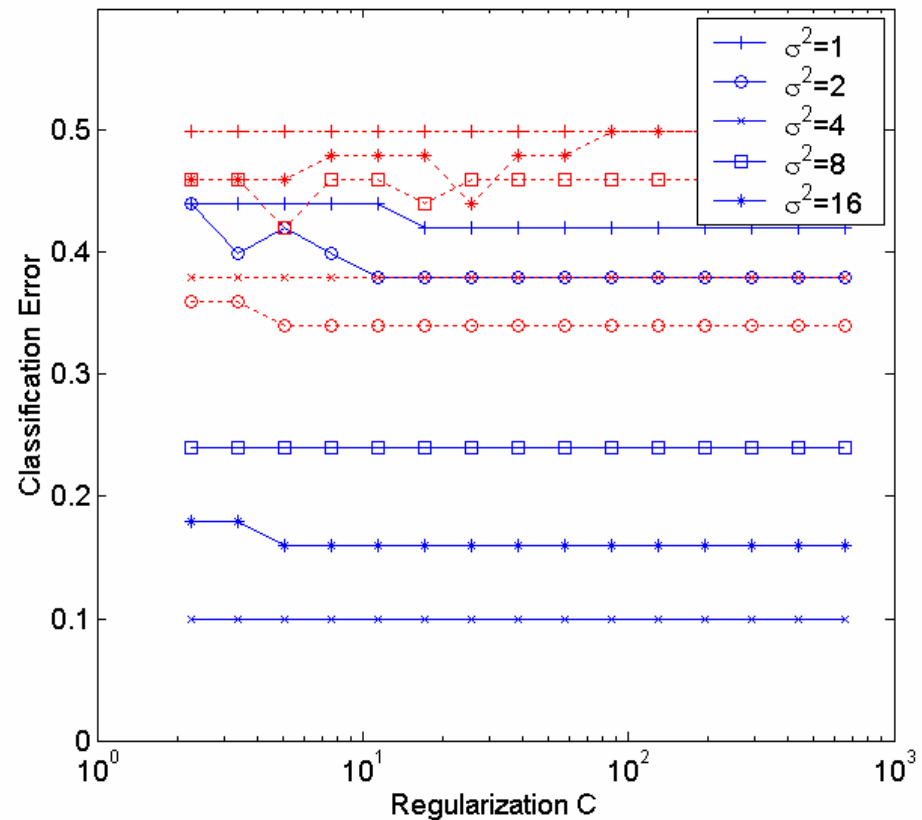


SVM: Train on 50, Test on 5

Fisher for Gaussian is Quadratic Kernel

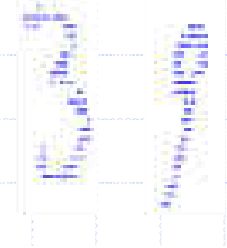
RBF Kernel (red)
67% accuracy

Bhattacharyya (blue)
90% accuracy

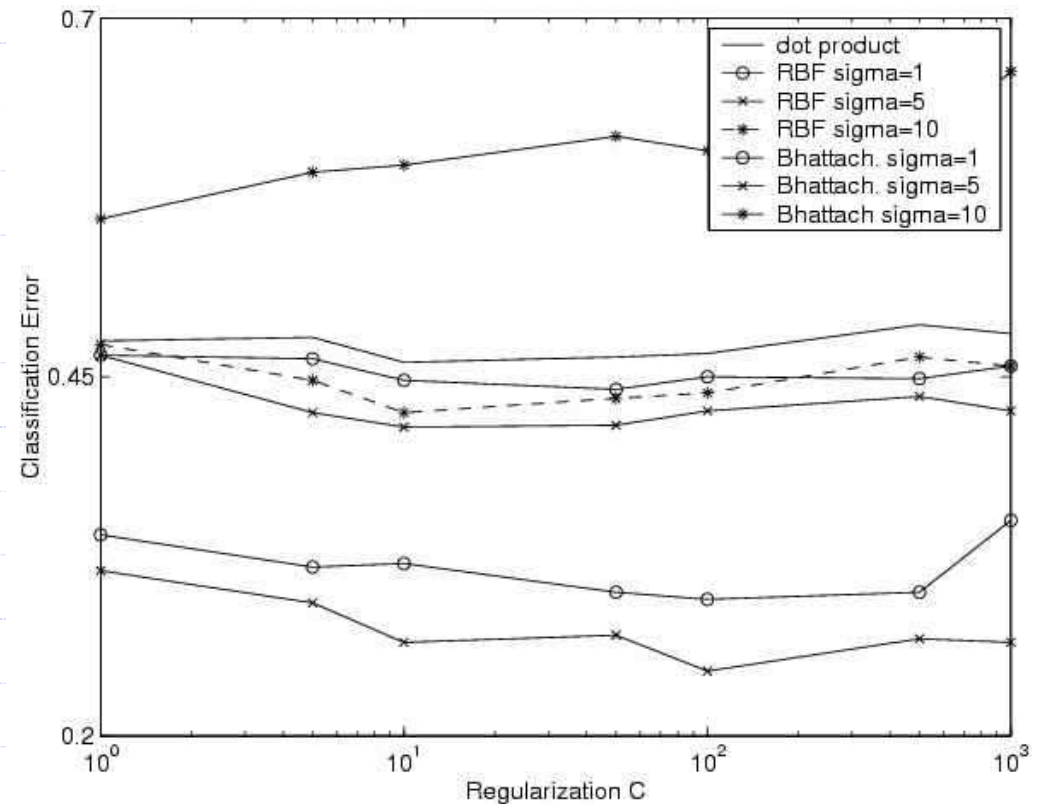


Kernelized Gaussian Product

SVM Classification of NIST digit images 0,1,...,9
 Sample each image to get bag of 30 (X,Y) pixels
 Train on random 120, test on random 80



bag-of-vectors
 Bhattacharyya
 outperforms
 standard RBF
 due to built-in
 invariance



Regularization & SRM Bounds

- Kernels also motivated from **Regularization Theory**
- Bound methods (like SRM) minimize some modified risk
- Empirical risk plus some capacity or **Regularizer Ω** :

$$R_{reg}(f) = R_{emp}(f) + \Omega[f] \quad \text{where } R_{emp}(f) = \frac{1}{N} \sum_{i=1}^N L(f(x_i), y)$$

- In structural risk minimization (SRM), Ω is the sqrt VC term

$$R(f) \leq R_{emp}(f) + \sqrt{\frac{h \left(\log \left(\frac{2N}{h} \right) + 1 \right) - \log \left(\frac{\eta}{4} \right)}{N}}$$

where N is # of points, h is VC dimension and bound holds with probability $1-\eta$

Regularization & Stability Bound

- In stability bounds, Ω is as below

$$R(f^*) \leq R_{emp}(f^*) + \beta + \sqrt{\frac{-2(N\beta + M) \log(\eta/2)}{N}}$$

The bound holds with probability $1-\eta$
with N training datapoints
with M as an upper bound on the max loss

$$M \geq \max_{x,y,f} |L(f(x), y)|$$

where the f^* is returned by a β -stable algorithm.
In other words, the loss for any example (x,y) changes by at most β when we replace one of the training examples with another training example.

Regularization & Rademacher

- In stability bounds, Ω is as below

$$R(f) \leq \inf_{\alpha > 0} \left[(1 + \alpha) R_{emp}(f) + \left(1 + \frac{1}{4\alpha}\right) \left(31r \log^2 \frac{b}{r} + 50 \frac{b\epsilon}{N}\right) \right]$$

The bound holds with probability $1-\eta$
 with N training datapoints
 with parameters r and b that are given
 by the so-called Rademacher averages

$$b = E_{S, \sigma} \left[\sup_{f: E[L(f(x), y)] \leq r} \sum_{i=1}^N \sigma_i L(f(x_i), y) \right]$$

Where σ_i are $+1$ or -1 with 50%-50% probability

Regularization & Function Norm

- In Regularization theory, Ω is most often as below

$$\begin{aligned}
 R_{reg}(f) &= \frac{1}{N} \sum_{i=1}^N L(f(x_i), y) + \|f\|_H^2 \\
 &= \frac{1}{N} \sum_{i=1}^N L(f(x_i), y) + \langle f, f \rangle_H \\
 &= \frac{1}{N} \sum_{i=1}^N L(\langle f, k_{x_i} \rangle, y) + \langle f, f \rangle_H
 \end{aligned}$$

Where line 1 à 2 is because of Hilbert space

Where line 2 à 3 is because of RKHS property

In other words, the k_x are prototypical functions such that

$$f(x_i) = \langle f, k_{x_i} \rangle$$

Representer Theorem: f will be in the span of $k_{x_1} \dots k_{x_N}$