# Advanced Machine Learning & Perception

## Instructor: Tony Jebara

# Topic 7

- Standard Kernels

- Generalization and Regularization

- Reproducing Kernel Hilbert Space and Kernel Properties

- Unusual Kernels

- String Kernels

- Probabilistic Kernels

- Probability Product Kernels

# Kernels

- Inner product of mapped inputs: k(X,X') = f(X)Tf(X')
- Example Kernels:

  **p=# of twists in SVM decision boundary**

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

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

- Kernels replace inner products in SVMs, allow nonlinearity:

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

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

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

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

# Standard Kernels over Vectors

- All take two input vectors x and y and produce a scalar:

$$k(x,y) = \left(x^T y + 1\right)^p$$

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

$$k(x,y) = \tanh\left(\kappa x^T y - \delta\right)$$

# String Kernels

•Inputs are two strings:     X     = "aardvark"

X'     = "accent"

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

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

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

= [3  0  0  1  0  0 ... 0  0  1  0  0 ... ]

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

•Can do this efficiently via dynamic programming

# String Kernels

- Want $k(X,X') = $ scalar value $= \phi(X)^T\phi(X')$

- Explicit Kernel can be slow because of many substrings s in our final vocabulary

- $\phi_s(X) = $ number of times substring s appears in X

- $\phi_a(\text{aardvark}) = 3$ $\qquad\qquad \phi_{ar}(\text{aardvark}) = 1$

- 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!

# Fisher Kernels

**Fisher Kernel: approximate distance between**
  **two generative models on statistical**
  **manifold using Kullback Leibler**

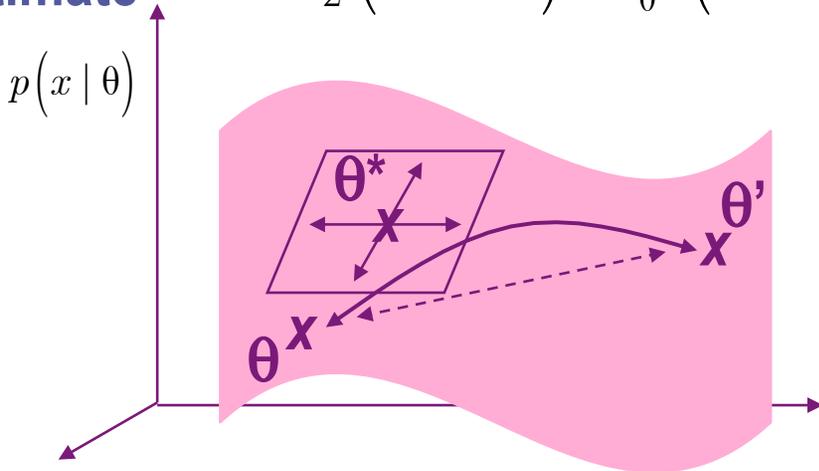**Approximate KL by quadratic form**
  **local tangent space at ML estimate**

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

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

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

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

$$U_\chi = \nabla_{\theta^*} \log p\left(\chi \mid \theta\right)$$

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

$p\left(x \mid \theta\right)$

$\theta^*$

$\theta'$

$\theta$

# Other Divergences

**Since the Kullback Leibler Divergence is asymmetric and tricky…**

**Find other divergences:** $D(p,q) \geq 0, D(p,q) = 0 \; iff \; p = q$

**Kullback-Leibler:** $KL(p,q) = \int p(x) \log \dfrac{p(x)}{q(x)} dx$

**Jensen-Shannon:** $JS(p,q) = \frac{1}{2} KL\left(p, \frac{p+q}{2}\right) + \frac{1}{2} KL\left(q, \frac{p+q}{2}\right)$

**Symmetrized KL:** $SKL(p,q) = \frac{1}{2} KL(p,q) + \frac{1}{2} KL(q,p)$

**Hellinger Divergence:** $H(p,q) = \frac{1}{4} \sqrt{\int \left(\sqrt{p(x)} - \sqrt{q(x)}\right)^2 dx}$

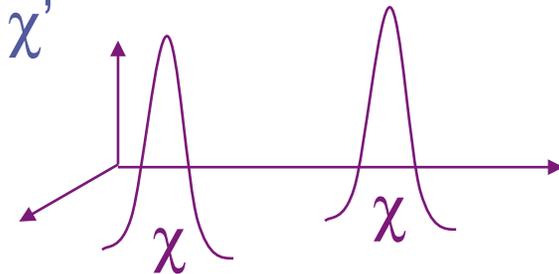**… makes a nice kernel:** $H(p,q) = \sqrt{k(p,p) - 2k(p,q) + k(q,q)}$

# Probability Product Kernels

**Computing the kernel: given inputs $\chi$ and $\chi'$**

   **1) Estimate Densities (I.e. ML):**

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

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

   **2) Compute Bhattacharyya Affinity:** $K(\chi, \chi') = \int \sqrt{p} \sqrt{p'} \, dx$

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

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

**Expected Likelihood Kernel:** $\rho = 1$

$$K(\chi, \chi') = E_{p}[p'] = E_{p'}[p]$$

# Gaussian Product Kernels

**Gaussian Mean: (any ρ)**
   **(continuous)**
**If μ=χ μ'=χ' get RBF Kernel**
**Fisher here is linear**

$$K\left(\chi, \chi'\right) = \int N^{\rho}\left(x \mid \mu\right) N^{\rho}\left(x \mid \mu'\right) dx$$

$$\propto \exp\left(-\left\|\mu - \mu'\right\|^2 / \tilde{\sigma}\right)$$

**Gaussian Covariance: (any ρ)**

$$K\left(\chi, \chi'\right) = \int N^{\rho}\left(x \mid \mu, \Sigma\right) N^{\rho}\left(x \mid \mu', \Sigma'\right) dx$$

$$\propto \left|\Sigma\right|^{-\rho/2} \left|\Sigma'\right|^{-\rho/2} \left|\Sigma^{\blacklozenge}\right|^{1/2} \exp\left(-\frac{\rho}{2}\mu^T\Sigma^{-1}\mu - \frac{\rho}{2}\mu'^T\Sigma'^{-1}\mu' + \frac{1}{2}\mu^{\blacklozenge}\Sigma^{\blacklozenge}\mu^{\blacklozenge}\right)$$

$$where \ \ \Sigma^{\blacklozenge} = \left(\rho\Sigma^{-1} + \rho\Sigma'^{-1}\right)^{-1} \ \ and \ \ \mu^{\blacklozenge} = \rho\Sigma^{-1}\mu + \rho\Sigma'^{-1}\mu'$$

**Fisher here is quadratic**
**But, how do we get a non-degenerate covariance from 1 data point?**

# Multinomial Product Kernels

**Bernoulli:**
**(binary)**

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

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

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

**Multinomial:**
**(discrete)**

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

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

**For multinomial counts (for N words):**

$$K\left(\chi, \chi'\right) = \left[\sum_{d=1}^{D} \left(\alpha_d \alpha_d'\right)^{1/2}\right]^{N}$$
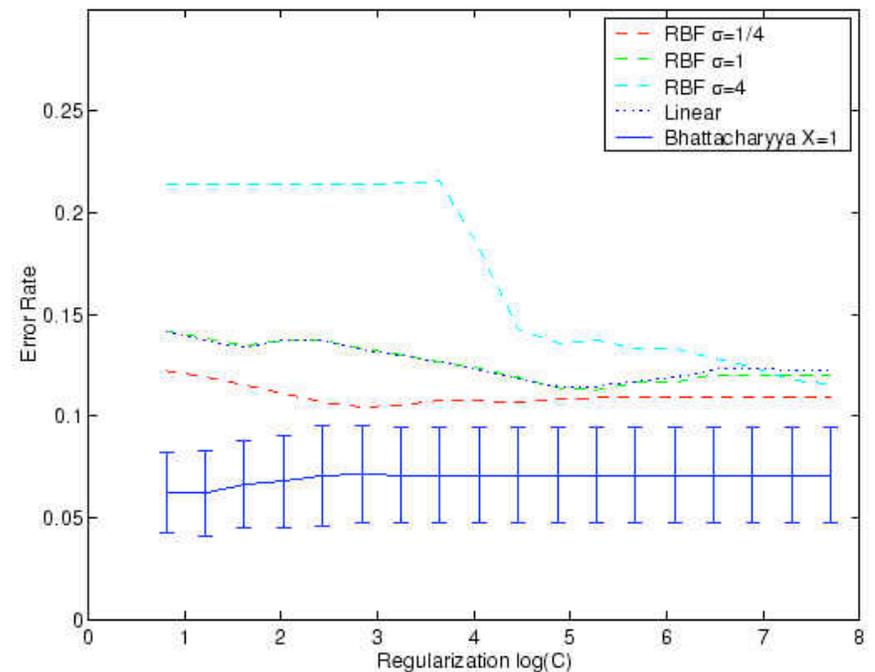
**Fisher for Multinomial is linear**

# 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**

# Exp. Family Product Kernels

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

$$p\left(x \mid \theta\right) = \exp\left[A\left(x\right) + T\left(x\right)^T \theta - K\left(\theta\right)\right]$$

**Maximum likelihood is straightforward:** $\frac{\partial}{\partial \theta} K\left(\theta\right) = \frac{1}{n} \sum_n T\left(x_n\right)$

**All have above form but different A(x), T(x), convex K(θ)**

| Family | $\mathcal{A}(X)$ | $\mathcal{K}(\theta)$ |
|---|---|---|
| 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 | $\log(X!)$ | $\exp(\theta)$ |

# Exp. Family Product Kernels

**Compute the Bhattacharyya Kernel for the e-family:**

$$p\big(x \mid \theta\big) = \exp\Big(A\big(x\big) + T\big(x\big)^{T} \theta - K\big(\theta\big)\Big)$$

**Analytic solution for e-family:**

$$K\big(\chi, \chi'\big) = \int p\big(x \mid \theta\big)^{1/2} p\big(x \mid \theta'\big)^{1/2} dx$$

$$= \exp\Big(K\big(\tfrac{1}{2}\theta + \tfrac{1}{2}\theta'\big) - \tfrac{1}{2}K\big(\tfrac{1}{2}\theta\big) - \tfrac{1}{2}K\big(\theta'\big)\Big)$$

**Only depends on convex cumulant-generating function K(θ)**

**Meanwhile, Fisher Kernel is always linear in sufficient stats…**

$$U_{x} = \nabla_{\theta^{*}} \log p\big(\chi \mid \theta\big) = T\big(\chi\big) - \nabla_{\theta^{*}} K\big(\theta\big)$$

$$U_{\chi} I_{\theta^{*}}^{-1} U_{\chi'} = \big(T\big(\chi\big) - g\big)^{T} I_{\theta^{*}}^{-1} \big(T\big(\chi'\big) - g\big)$$

# Gaussian Product Kernels

**Instead of a single $\chi$ & $\chi'$**
**Construct p & p' from many $\chi$ & $\chi'$**
**I.e. use bag of vectors**

$$\left\{ \chi_1, \ldots, \chi_N \right\} \to p\left(x\right)$$

$$\left\{ \chi_1', \ldots, \chi_{N'}' \right\} \to p'\left(x\right)$$

$$\mu = E\left\{ \chi_i \right\} \qquad \Sigma = E\left\{ \left(\chi_i - \mu\right)\left(\chi_i - \mu\right)^T \right\}$$

 $\sim N\left(x \mid \mu, \Sigma\right)$       $\sim N\left(x \mid \mu', \Sigma'\right)$

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

# Kernelized Gaussian Product

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

**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 <u>again</u> when computing Gaussian mean & covariance**

$$\kappa\left(\chi,\chi'\right) = \phi\left(\chi\right)^{T}\phi\left(\chi'\right)$$

**Never compute outerproducts, use kernel, I.e. infinite RBF:**

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

**Compute mini-kernel between each pixel in a given image…**

**Gives kernelized or augmented Gaussian $\mu$ and $\Sigma$ via Gram**

# Kernelized Gaussian Product

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

**Now have:**

$$\mu = E\left\{\phi\left(\chi\right)\right\} \qquad \Sigma = E\left\{\left(\phi\left(\chi\right) - \mu\right)\left(\phi\left(\chi\right) - \mu\right)^T\right\}$$

**Still invariant to order of pixels**

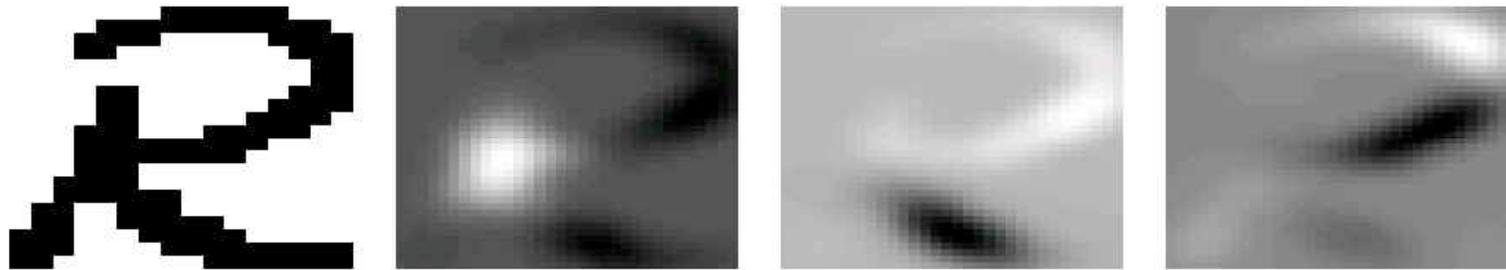**Compute Hilbert Gaussian's mean & covariance of each image bag or image is N x N pixel Gram matrix using kappa kernel**

**Use kernel PCA to regularize infinite dimensional RBF Gaussian**



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



$$K\left(\phi\left(\chi\right) \| \phi\left(\chi'\right)\right) \propto \left|\Sigma\right|^{-\rho/2} \left|\Sigma'\right|^{-\rho/2} \left|\Sigma^{\blacklozenge}\right|^{1/2} \exp\left(-\tfrac{\rho}{2}\mu^T\Sigma^{-1}\mu - \tfrac{\rho}{2}\mu'^T\Sigma'^{-1}\mu' + \tfrac{1}{2}\mu^{\blacklozenge}\Sigma^{\blacklozenge}\mu^{\blacklozenge}\right)$$

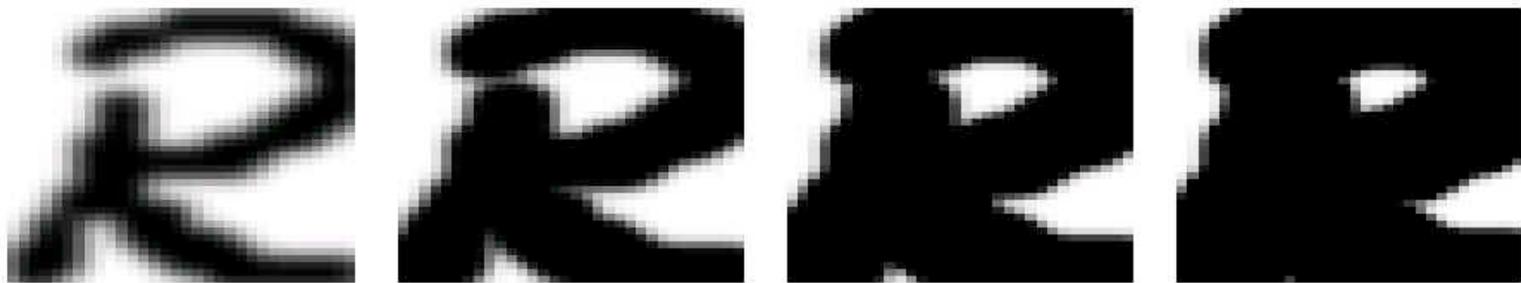**Puts all dimensions (X,Y,I) on an equal footing**

# Kernelized Gaussian Product



**Letter 'R' with 3 KPCA Components of RBF Kernel**
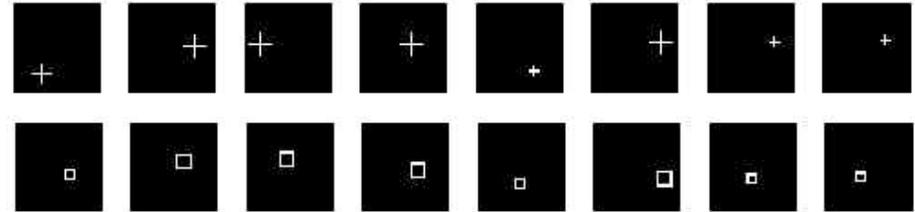


**Reconstruction of Letter 'R' with 1-4 KPCA with RBF Kernel**



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

# Kernelized Gaussian Product

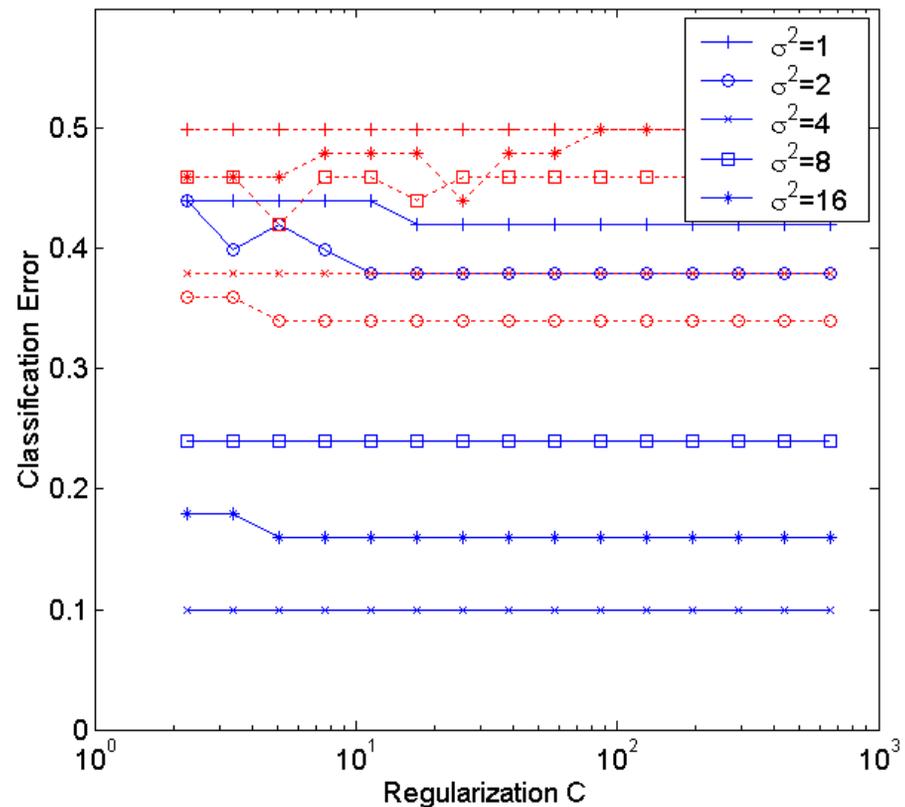**100 40x40 monochromatic images of crosses & squares translating & scaling**

**SVM: Train on 50, Test on 50**

**Fisher for Gaussian is Quadratic Kernel**
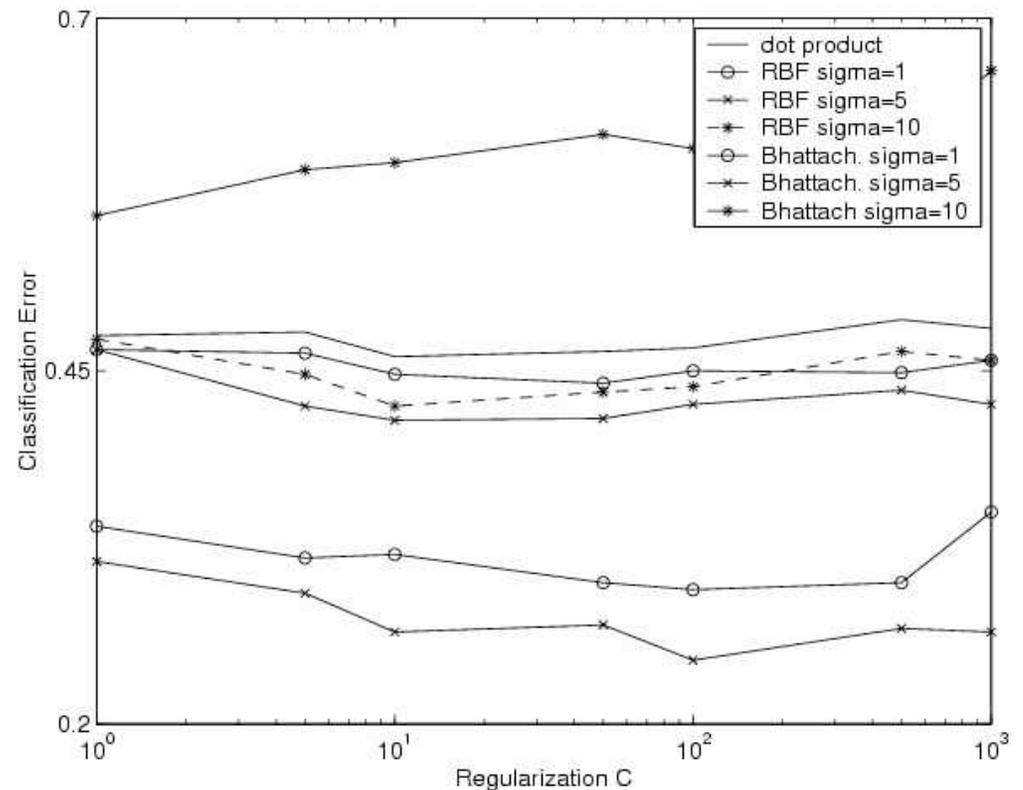
**RBF Kernel (red) 67% accuracy**

**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**

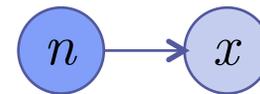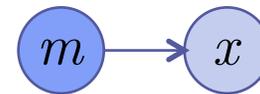**Fisher Kernel for
Gaussian is quadratic**

# Mixture Product Kernels

**Beyond Exponential Family: Mixtures and Hidden Variables**
**Easier for ρ=1 Expected Likelihood kernel…**

$$\chi \to p(x) = \sum_{m=1}^{M} p(m)p(x \mid m)$$
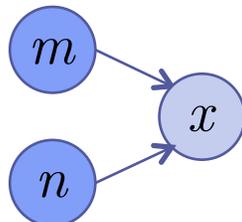
$$\chi' \to p'(x) = \sum_{n=1}^{N} p'(n)p'(x \mid n)$$

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

$$= \sum_{m=1}^{M}\sum_{n=1}^{N} p(m)p'(n)\int p(x \mid m)p'(x \mid n)\,dx$$

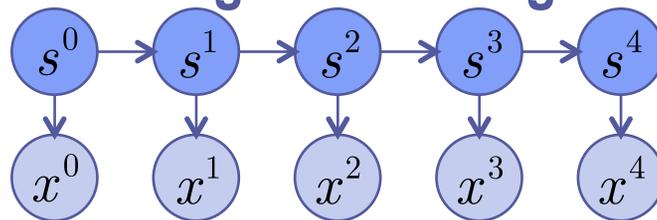$$= \sum_{m=1}^{M}\sum_{n=1}^{N} p(m)p'(n)K_{m,n}(\chi, \chi')$$

**Use M\*N subkernel evaluations from our previous repertoire**

# HMM Product Kernels

**Hidden Markov Models: (sequences)**

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

**# of hidden configurations large**



$$\#configs = \left|s\right|^T$$

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

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

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

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

# HMM Product Kernels

$$K\left(\chi, \chi'\right) = \sum_S \sum_U \prod_t p\left(s_t \mid s_{t-1}\right) p'\left(u_t \mid u_{t-1}\right) \int p\left(x_t \mid s_t\right) p'\left(x_t \mid u_t\right) dx_t$$

$$= \sum_{S,U} \prod_t p\left(s_t \mid s_{t-1}\right) p'\left(u_t \mid u_{t-1}\right) K_{s_t, u_t}$$

$$= \sum_{S,U} p\left(s_0\right) p'\left(u_0\right) \psi\left(s_0, u_0\right) \prod_t p\left(s_t \mid s_{t-1}\right) p'\left(u_t \mid u_{t-1}\right) \psi\left(s_t, u_t\right)$$
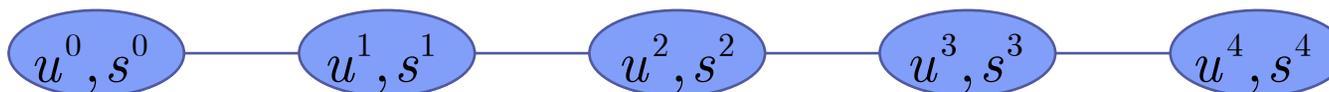
**Take advantage of structure in HMMs via Bayesian network**



**Only compute subkernels for common parents**

**Evaluate total of** $O\left(T |s||u|\right)$ **subkernels**

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

# Sampling Product Kernels

**Sampling: approximate K**

$$K\left(\chi, \chi'\right) = \int p\left(x\right) p'\left(x\right) dx$$

**By definition, generative models can:**

       **1) Generate a Sample**

       **2) Compute Likelihood of a Sample**

**Thus, approximate probability product via sampling:**

$$K\left(\chi, \chi'\right) = K\left(p, p'\right)$$

$$= \frac{\beta}{N} \sum_{\substack{x_i \sim p(x) \\ i=1\ldots N}} p'\left(x_i\right) + \frac{1-\beta}{N'} \sum_{\substack{x_i' \sim p'(x) \\ i=1\ldots N'}} p\left(x_i'\right)$$

**Beta controls how much sampling from each distribution...**