

Advanced Machine Learning & Perception

Instructor: Tony Jebara

Topic 6

- Standard Kernels
- Unusual Input Spaces for Kernels
- String Kernels
- Probabilistic Kernels
- Fisher Kernels
- Probability Product Kernels
- Mixture Model Probability Product Kernels
- Junction Tree Algorithm for Kernels

SVMs and Kernels

- Recall the SVM dual optimization:

$$L_D : \max \sum_i \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j k(x_i, x_j)$$

subject to $\alpha_i \in [0, C]$ and $\sum_i y_i \alpha_i = 0$

- Solve a QP by using the **Gram** matrix K where $K_{ij} = k(x_i, x_j)$

$$K = \begin{bmatrix} k(x_1, x_1) & k(x_1, x_2) & k(x_1, x_3) \\ k(x_1, x_2) & k(x_2, x_2) & k(x_2, x_3) \\ k(x_1, x_3) & k(x_2, x_3) & k(x_3, x_3) \end{bmatrix}$$

- Then KKT conditions give us b . The SVM prediction is then:

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

- Kernels replace vector dot products for nonlinear boundary
- Kernels also allow *arbitrary large structured input spaces!*

Standard Kernels on Vectors

- Kernels usually take two input vectors and output a scalar:
- Polynomial Kernel (gives p twists in SVM decision boundary)

$$k(x_i, x_j) = \langle \phi(x_i), \phi(x_j) \rangle = (x_i^T x_j + 1)^p$$

- RBF Kernel (phi is infinite dimensional)

$$k(x_i, x_j) = \langle \phi(x_i), \phi(x_j) \rangle = \exp\left(-\frac{1}{2\sigma^2} \|x_i - x_j\|^2\right)$$

- Hyperbolic Tangent Kernel

$$k(x_i, x_j) = \langle \phi(x_i), \phi(x_j) \rangle = \tanh(\kappa x_i^T x_j - \delta)$$

- But why just deal with vector inputs?
- We can apply kernels between any type of input:
graphs, strings, probabilities, structured spaces, etc.!

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

$$\begin{aligned} \phi(X) &= [\#a \ \#b \ \#c \ \#d \ \#e \ \#f \ \dots \ \#y \ \#z \ \#aa \ \#ab \ \#ac \ \dots] \\ &= [3 \ 0 \ 0 \ 1 \ 0 \ 0 \ \dots \ 0 \ 0 \ 1 \ 0 \ 0 \ \dots] \end{aligned}$$

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

- Can do this efficiently via dynamic programming

String Kernels

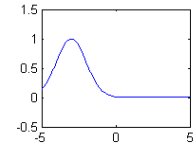
- Want $k(X, X') = \text{scalar value} = \phi(X)^T \phi(X')$
- Explicit Kernel can be slow because of many substrings s in our final vocabulary
- $\phi_s(X) = \text{number of times substring } s \text{ appears in } X$
- $\phi_a(\text{aardvark}) = 3$ $\phi_{ar}(\text{aardvark}) = 2$
- Implicit Kernel is more efficient

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

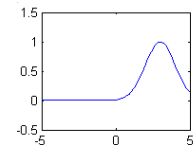
...via dynamic programming in time proportional to the product of the lengths of X and X' . This computes the dot product much more quickly!

Probability Kernels

•Inputs are two probabilities: $p = p(x|\theta) =$



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



•Want $k(p,p') =$ scalar value

•How to design this?

•Recall that a kernel can be converted into a distance:

$$D(p_i, p_j) = \sqrt{k(p_i, p_i) - 2k(p_i, p_j) + k(p_j, p_j)}$$

•What is a natural distance between two probabilities?

•Kullback-Leibler Divergence!

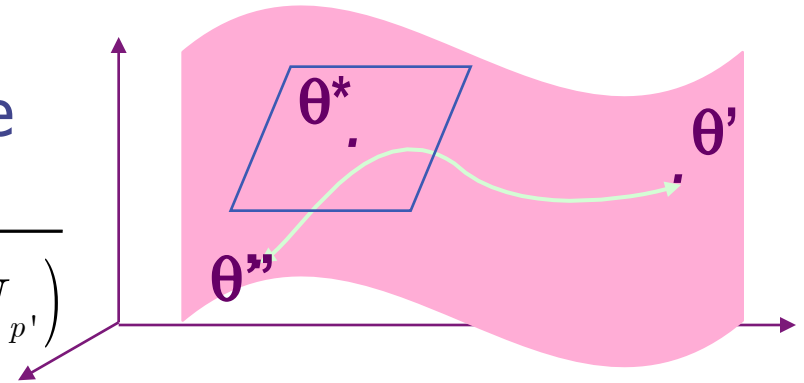
$$D(p_i \parallel p_j) = \int p_i(x) \log \frac{p_i(x)}{p_j(x)} dx$$

...caveats?

Fisher Kernels

- KL isn't symmetric distance, probabilities lie on a manifold!
- So we must approximate it...
- Try quadratic distance measure like Mahalanobis distance

$$D(p'' \| p') \approx \sqrt{(U_{p''} - U_{p'})^T I_{\theta^*}^{-1} (U_{p''} - U_{p'})}$$



- Linearize manifold at θ^* (some maximum likelihood point)
- Get distance between p and p' via a distance from θ to θ'
- The right kernel to go with the Mahalanobis distance is:

$$k(p'', p') = U_{p''}^T I_{\theta^*}^{-1} U_{p'}$$

where $U_{p''} = \nabla_{\theta''} \log p(x | \theta)$ and $U_{p'} = \nabla_{\theta'} \log p(x | \theta)$

- Matrix I is the Fisher information

$$I_{\theta^*}(i, j) = \int p(x | \theta^*) \frac{\partial \log p(x | \theta)}{\partial \theta_i} \frac{\partial \log p(x | \theta)}{\partial \theta_j} dx$$

Other Possible Divergences

The Kullback-Leibler Divergence is asymmetric and tricky...

All divergences have: $D(p, q) \geq 0, D(p, q) = 0$ iff $p = q$

Consider KL, Jensen-Shannon and Symmetrized KL...

$$KL(p, q) = \int p(x) \log p(x) q^{-1}(x) dx$$

$$JS(p, q) = \frac{1}{2} KL\left(p, \frac{p+q}{2}\right) + \frac{1}{2} KL\left(q, \frac{p+q}{2}\right)$$

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

**sqrt(JS) is a true distance metric
(Endres & Schindelin, 2003)**

Hellinger divergence is symmetric & can be written as kernel!

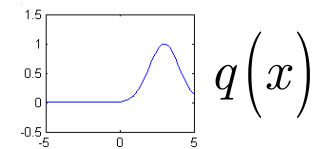
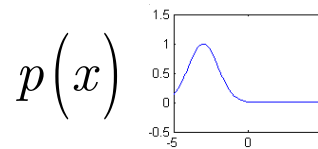
$$H(p, q) = \frac{1}{4} \sqrt{\int \left(\sqrt{p(x)} - \sqrt{q(x)} \right)^2 dx} = \sqrt{k(p, p) - 2k(p, q) + k(q, q)}$$

...what kernel is this?

Bhattacharyya Kernel

A kernel between distributions that gives Hellinger divergence

$$k(p, q) = \int \sqrt{p(x)} \sqrt{q(x)} dx$$



A slight generalization is the *Probability Product Kernel*

$$k(p, q) = \int (p(x))^\rho (q(x))^\rho dx$$

Setting $\rho=1/2$ gives the *Bhattacharyya Kernel*

$$k(p, q) = \int \sqrt{p(x)} \sqrt{q(x)} dx$$

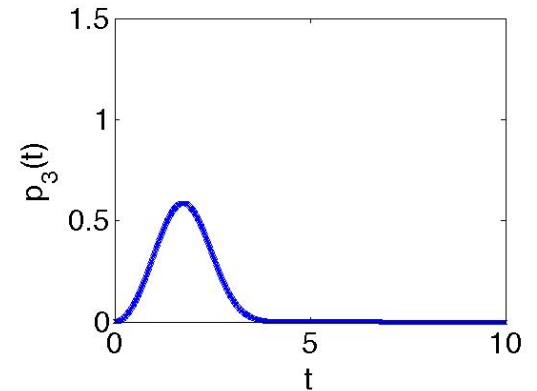
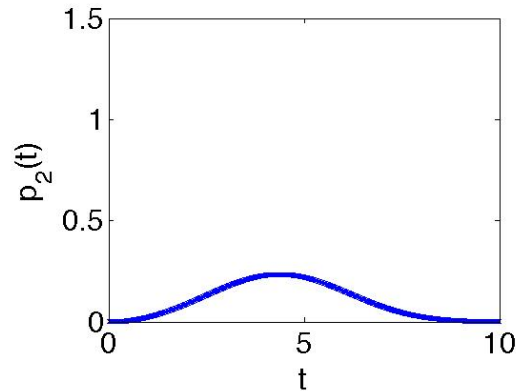
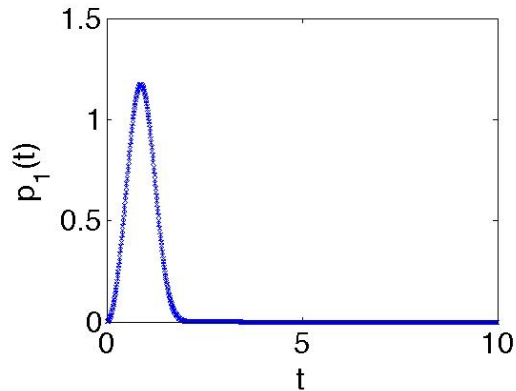
Setting $\rho=1$ gives the *Expected Likelihood Kernel*

$$k(p, q) = E_{p(x)} [q(x)] = E_{q(x)} [p(x)]$$

Probability Product Kernel

Imagine we are given inputs which are probabilities and labels

Dataset = $\{(p_1, y_1), (p_2, y_2), \dots, (p_n, y_n)\}$



E.g. each input i is a brand of car
 p_i is its brakepad failure probability over time
 y_i is binary, did government approve it for US roads?

Build SVM with Probability Product Kernels to predict if a brakepad failure distribution will be approved/rejected by gov't

Probability Product Kernel

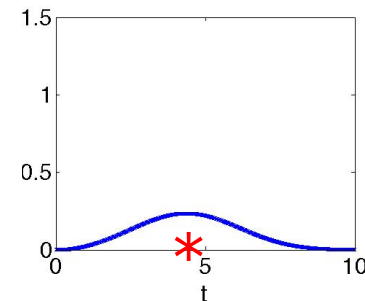
Imagine we are given input data vectors with labels

$$\text{Dataset} = \{(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)\}$$

E.g. each input i is a brand of car
 x_i is time when one car of brand i failed
 y_i is binary, did government approve it for US roads?

For each i , estimate from
 $x_i \rightarrow$ probability $p_i = p(x|\theta_i)$

$$\text{Dataset} = \{(p_1, y_1), (p_2, y_2), \dots, (p_n, y_n)\}$$



Build SVM with Probability Product Kernels to predict if a brakepad failure distribution will be approved/rejected by gov't

Probability Product Kernel

Imagine we are given input data vector-sets with labels

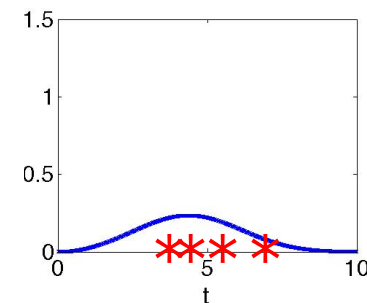
$$\text{Dataset} = \{(\chi_1, y_1), (\chi_2, y_2), \dots, (\chi_n, y_n)\}$$

E.g. each input i is a brand of car
 χ_i is a set of times when several cars of brand i failed
 y_i is binary, did government approve it for US roads?

For each i , estimate from

$$\chi_i = \{x_{i1}, \dots, x_{it}\} \rightarrow \text{probability } p_i = p(x | \theta_i)$$

$$\text{Dataset} = \{(p_1, y_1), (p_2, y_2), \dots, (p_n, y_n)\}$$



Build SVM with Probability Product Kernels to predict if a brakepad failure distribution will be approved/rejected by gov't

Gaussian Product Kernels

For Gaussians:
$$k(p, p') = \int N^\rho(x | \mu) N^\rho(x | \mu') dx$$

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

(if $\mu=x$ and $\mu'=x'$, get back RBF kernel 😊)
 (Fisher kernel here is just a linear kernel ☹)

For Gaussians with variable covariance:

$$k(p, p') = \int N^\rho(x | \mu, \Sigma) N^\rho(x | \mu', \Sigma') dx$$

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

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

(Fisher kernel here is just a quadratic kernel ☹)

Multinomial Product Kernels

Bernoulli:
(binary x)

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

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

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

Multinomial:
(discrete x)

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

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

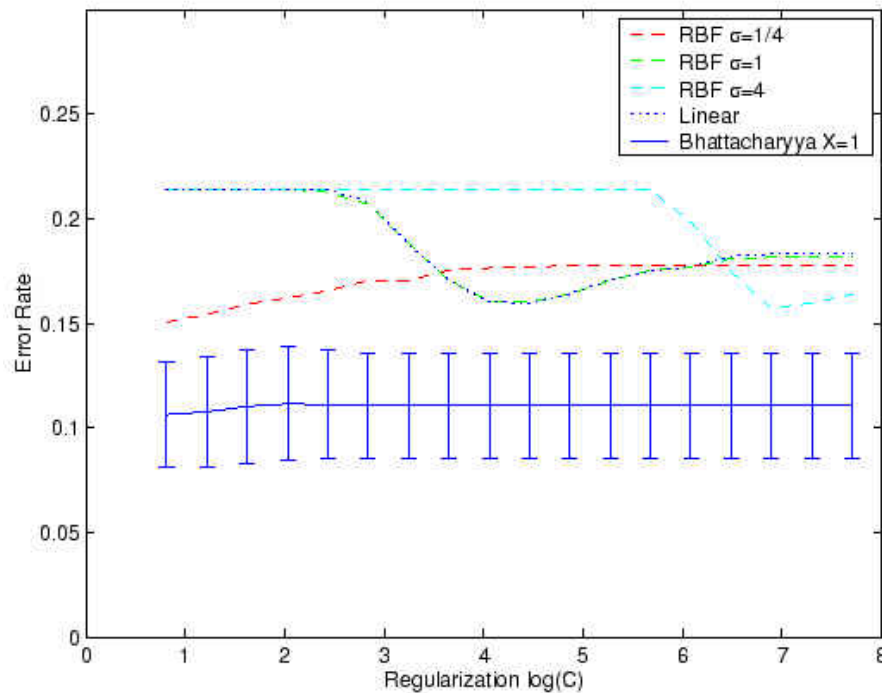
For multinomial counts (for N words per document):

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

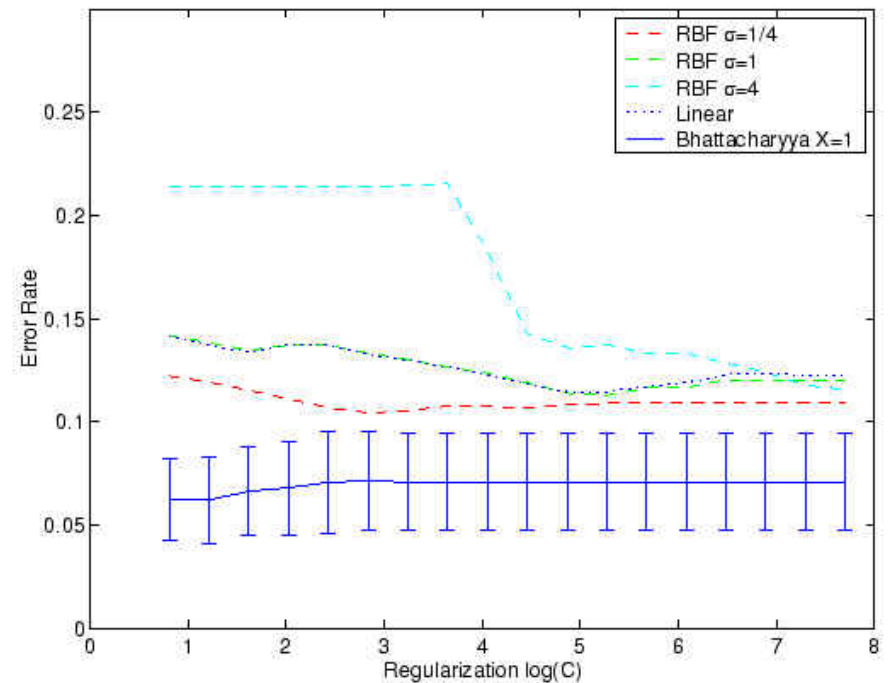
(Fisher for Multinomial is linear ☹)

Multinomial Product Kernels

WebKB dataset: Faculty vs. student web page SVM 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: Gaussian, Multinomial, Binomial, Poisson, Inverse Wishart, Gamma, Bernoulli, Dirichlet, ...

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

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 $K(q)$

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	$\log(X!)$	$\exp(\theta)$

Exponential Family Kernels

Compute the Bhattacharyya Kernel for the e-family:

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

Analytic solution for e-family:

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

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

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

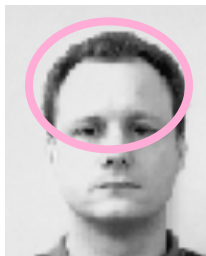
$$\begin{aligned} U_x &= \nabla_{\theta^*} \log p(\chi | \theta) = T(\chi) - \nabla_{\theta^*} K(\theta) \\ U_{\chi} I_{\theta^*}^{-1} U_{\chi'} &= \left(T(\chi) - g\right)^T I_{\theta^*}^{-1} \left(T(\chi') - g\right) \end{aligned}$$

Gaussian Product Kernels

Instead of a single x & x'

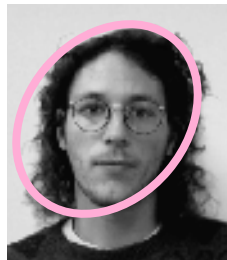
Construct p & p' from many x & x' samples

Use bag of vectors



$$\{\chi_1, \dots, \chi_N\} \rightarrow p(x) = N(x | \mu, \Sigma)$$

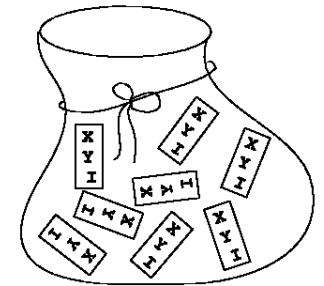
$$\mu = \frac{1}{N} \sum_i \chi_i, \quad \Sigma = \frac{1}{N} \sum_i (\chi_i - \mu)(\chi_i - \mu)^T$$



$$\{\chi'_1, \dots, \chi'_{N'}\} \rightarrow p'(x) = N(x | \mu', \Sigma')$$

$$\mu' = \frac{1}{N'} \sum_i \chi'_i, \quad \Sigma' = \frac{1}{N'} \sum_i (\chi'_i - \mu')(\chi'_i - \mu')^T$$

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

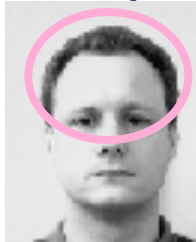


Kernelized Gaussian Product

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

Invariant to order of tuples in each bag

But too simple, just overlap of Gaussians on images



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

Use another Kernel to find Gaussian in Hilbert space

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

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...

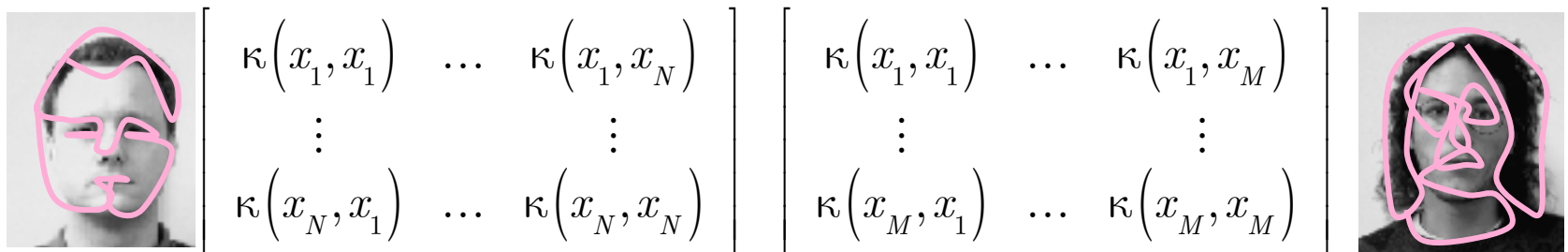
Gives kernelized or augmented Gaussian μ and Σ via Gram

Kernelized Gaussian Product

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

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

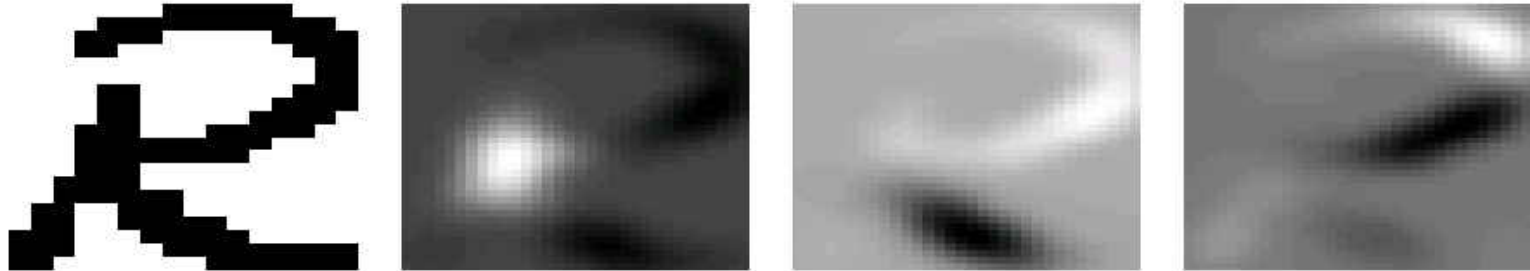
Compute Hilbert Gaussian's mean & covariance of each image bag or image is $N \times N$ pixel Gram matrix using kappa kernel
Use kernel PCA to regularize infinite dimensional RBF Gaussian



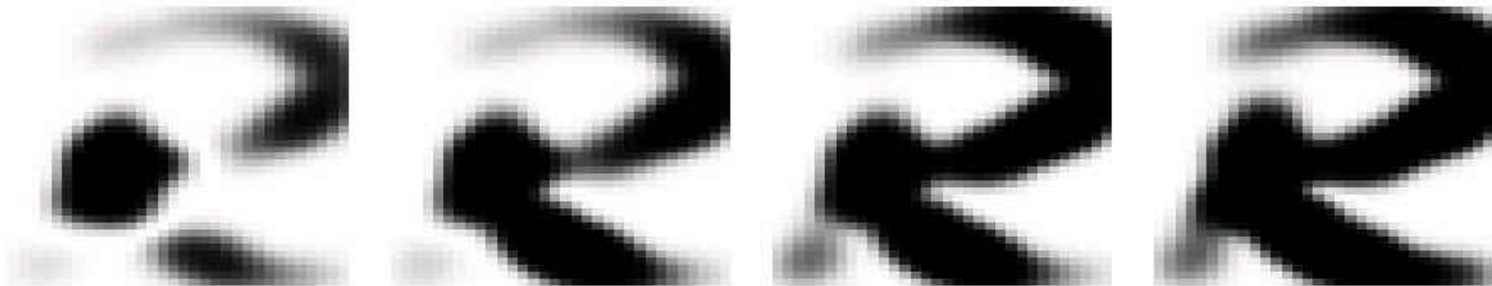
$$k(\phi(\chi) \parallel \phi(\chi')) \propto |\Sigma|^{-\rho/2} |\Sigma'|^{-\rho/2} |\bar{\Sigma}|^{1/2} \exp\left(-\frac{\rho}{2} \mu^T \Sigma^{-1} \mu - \frac{\rho}{2} \mu'^T \Sigma'^{-1} \mu' + \frac{1}{2} \bar{\mu}^T \bar{\Sigma} \bar{\mu}\right)$$

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

Kernelized Gaussian Product



Letter 'R' with 3 KPCA Components of RBF Kernel



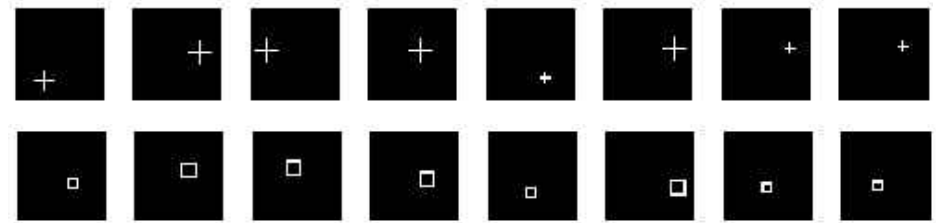
Reconstruction of Letter 'R' with 1-4 KPCA with RBF Kernel



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

Kernelized Gaussian Product

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

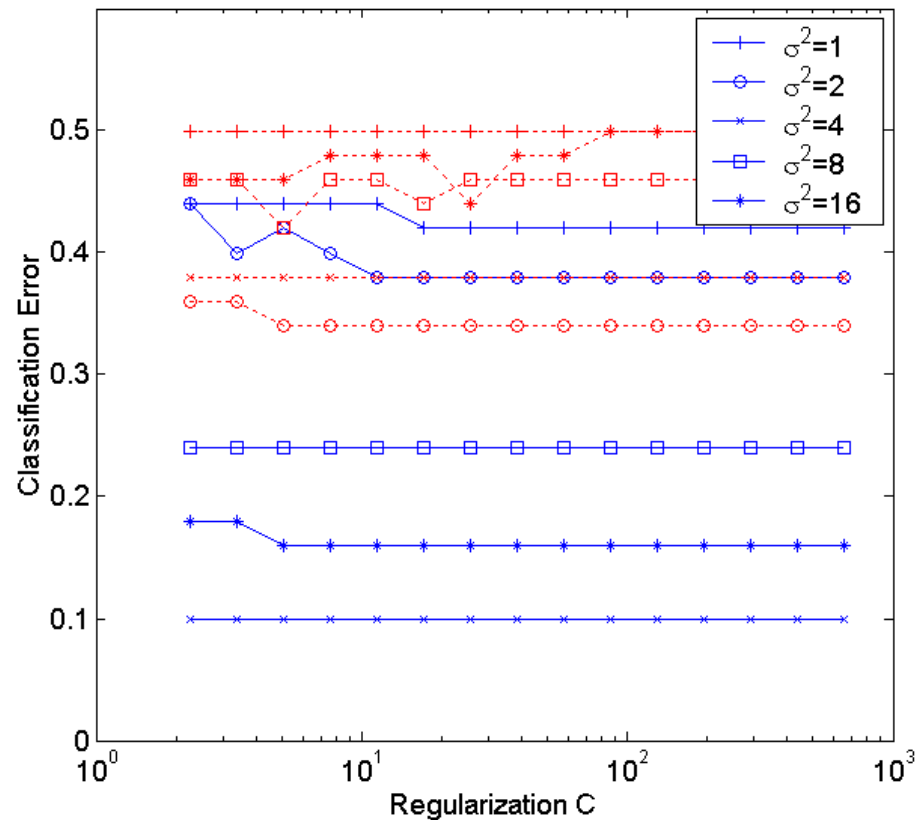


SVM: Train 50, Test 50

Fisher for Gaussian is Quadratic Kernel

RBF Kernel (red)
67% accuracy

Bhattacharyya (blue)
90% accuracy



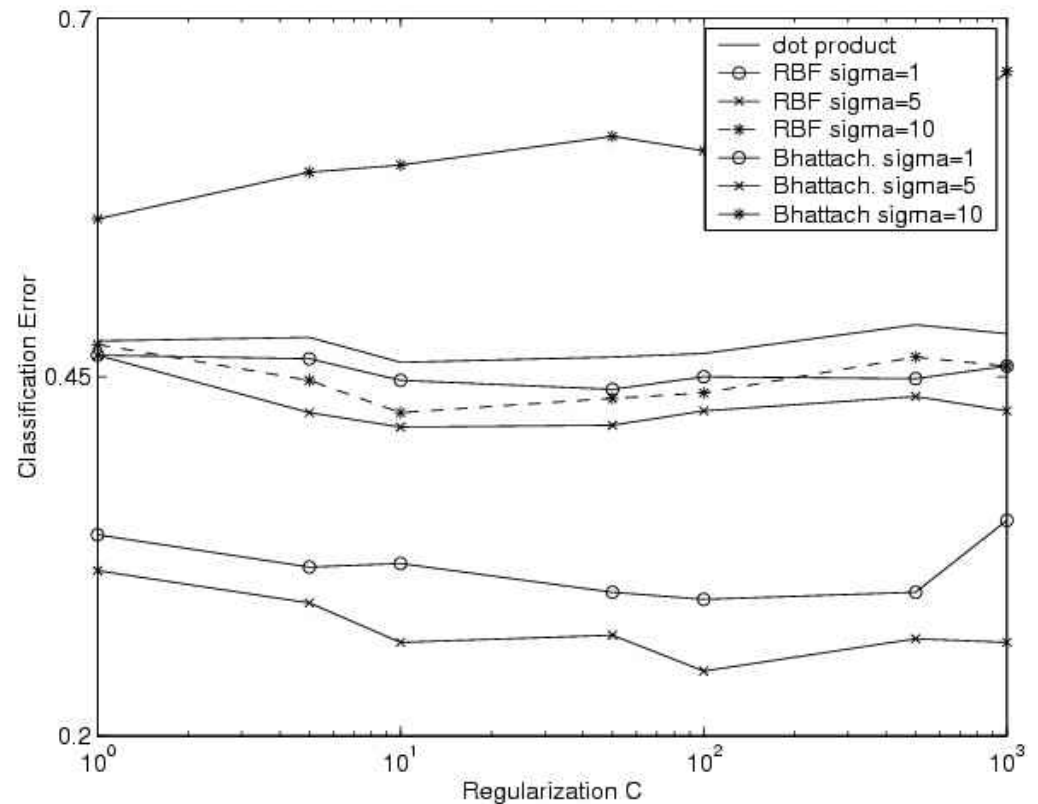
Kernelized Gaussian Product

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

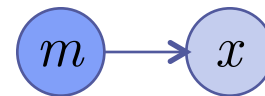
Fisher Kernel for
 Gaussian is quadratic



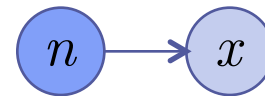
Mixture Product Kernels

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

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



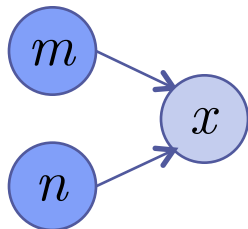
$$\chi' \rightarrow p'(x) = \sum_{n=1}^N p'(n) p'(x | n)$$



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

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

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



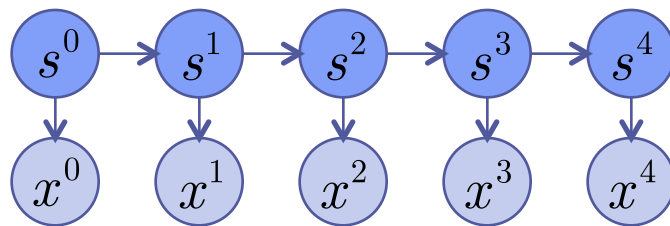
Use $M*N$ subkernel evaluations
from our previous repertoire

HMM Product Kernels

Hidden Markov Models: (sequences)

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

of hidden configurations large



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

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

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

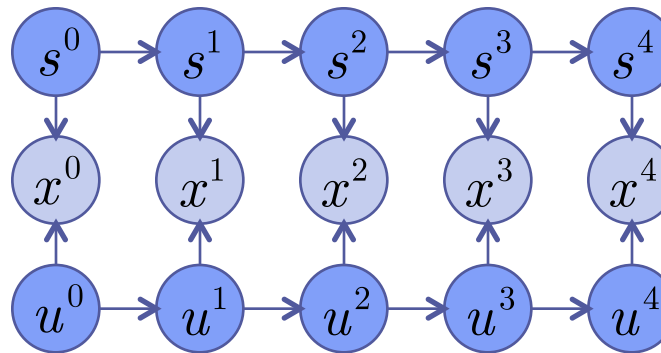
$$= \sum_S \sum_U p(S) p'(U) c_{S,U}$$

☹ computes cross-product of all hidden variables $O(|s|^T \times |u|^T)$

HMM Product Kernels

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

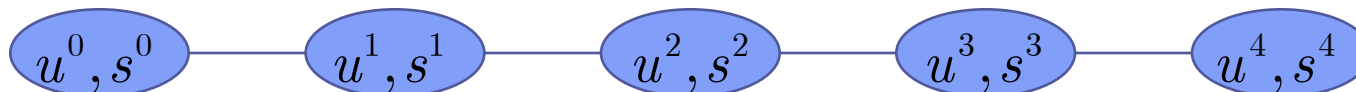
Take advantage of structure in HMMs via Bayesian network



Only compute subkernels for common parents

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

Form clique potential fn's, sum via junction tree algorithm



Sampling Product Kernels

Can approximate probability product kernel via sampling

By definition, generative models can:

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

Thus, approximate probability product via sampling:

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

Beta controls how much sampling from each distribution...