

Diffusion Kernels and Friends

Inducing kernels from the local structure of
input space

Risi Kondor

Columbia University, New York, USA.

Collaborators

John Lafferty

Mikhail Belkin

Alex Smola

Tony Jebara

Count Laplace (1749-1827)

Regularization Networks

Learning from examples $(x_1, y_1), (x_2, y_2), \dots, (x_m, y_m)$

Looking for $f : \mathcal{X} \mapsto \mathcal{Y}$ minimizing

$$R^{\text{reg}}[f] = \underbrace{\frac{1}{m} \sum_{i=1}^m L(y_i, f(x_i))}_{\text{Loss function}} + \underbrace{\langle \bar{Q}f, \bar{Q}f \rangle_{\mathcal{L}_2}}_{\text{Complexity penalty}}$$

Regularization operator: \bar{Q}

Inner product here: $\langle f, g \rangle_{\mathcal{L}_2} = \int f(x)g(x) dx$

The Kernel — A Similarity Measure

Feature map: $\Phi : \mathcal{X} \mapsto \mathcal{F}$

Hypothesis f linear in \mathcal{F}

Kernel: $K(x, x') = \langle \Phi(x), \Phi(x') \rangle$ — similarity measure

Kernel K and regularization operator \underline{Q} are related

Kernel must be positive definite, i.e.

$$\sum_n \sum_{j=1}^n \alpha_i \alpha_j K(x_i, x_j) \geq 0 \quad \forall \alpha_i \in \mathbb{R} \\ \forall x_1, x_2, \dots, x_n \in \mathcal{X}$$

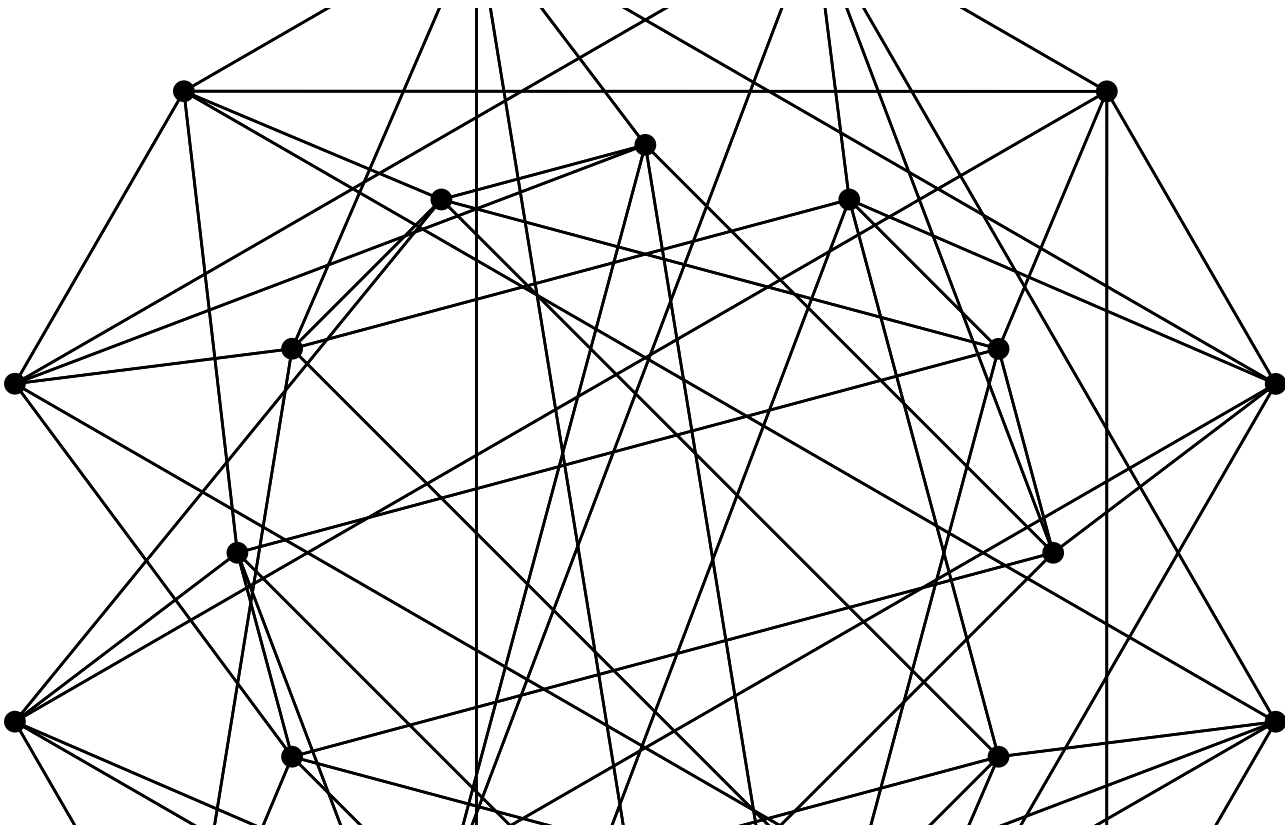
Correspondence between kernel and operator

A typical kernel on $\mathcal{X} = \mathbb{R}^d$ is the Gaussian RBF

$$K(x, x') = \frac{1}{(2\pi\sigma^2)^{d/2}} e^{-\|x-x'\|^2/(2\sigma^2)}.$$

Regularization operator: $(\widehat{Qf})(\omega) = e^{-\sigma^2\|\omega\|^2} \widehat{f}(\omega)$

Graphs



Graphs

Graphs

- Natural graphs: internet, web, social contacts, citations, scientific collaborations, etc.
- Objects with graph-like structure: groups (permutations), strings, etc.
- Objects with unknown global structure: set of organic molecules, interacting proteins in a cell
- Incorporating unlabelled data

Looking for positive definite $K : V \times V \mapsto \mathbb{R}$, now just a matrix

Try Random Walks

$$A_{ij} = \begin{cases} 1 & i \sim j \\ 0 & \text{otherwise} \end{cases}$$

A symmetric \Leftrightarrow even powers pos. def.

$$K = A_2 = A_2 \quad ? \quad A_4 \quad ? \quad A_\infty \quad ?$$
$$K = \alpha_1 A_2 + \alpha_2 A_4 + \dots \quad ?$$

Diffusion Kernels

Infinite number of infinitesimal steps:

$$K = e^{T\beta} = \lim_{n \rightarrow \infty} \left(1 + \frac{T}{n} \beta \right)^n$$

$$T_{ij} = \begin{cases} 1 & i \sim j \\ -d_i & i = j \\ 0 & \text{otherwise} \end{cases} \quad (\text{Laplacian})$$

for any infinitely divisible (or finite) K .

$$K = \lim_{n \rightarrow \infty} {}_n K_{1/n} = \lim_{n \rightarrow \infty} \left(I + \frac{1}{n} T \right) = {}_T e$$

conversely,

$${}_{\beta T} e = \lim_{n \rightarrow \infty} \left(I + \frac{\beta}{n} T \right) = \lim_{2n \rightarrow \infty} \left(I + \frac{\beta}{2n} T \right)$$

For any symmetric T , $K = {}_{\beta T} e$ is positive definite.

Exponential kernels

$$K_\beta = e^{\beta T} = \lim_{n \rightarrow \infty} \left(I + \frac{\beta T}{n} \right)^n$$

$$= I + \beta T + \frac{\beta^2 T^2}{2!} + \frac{\beta^3 T^3}{3!} + \dots$$

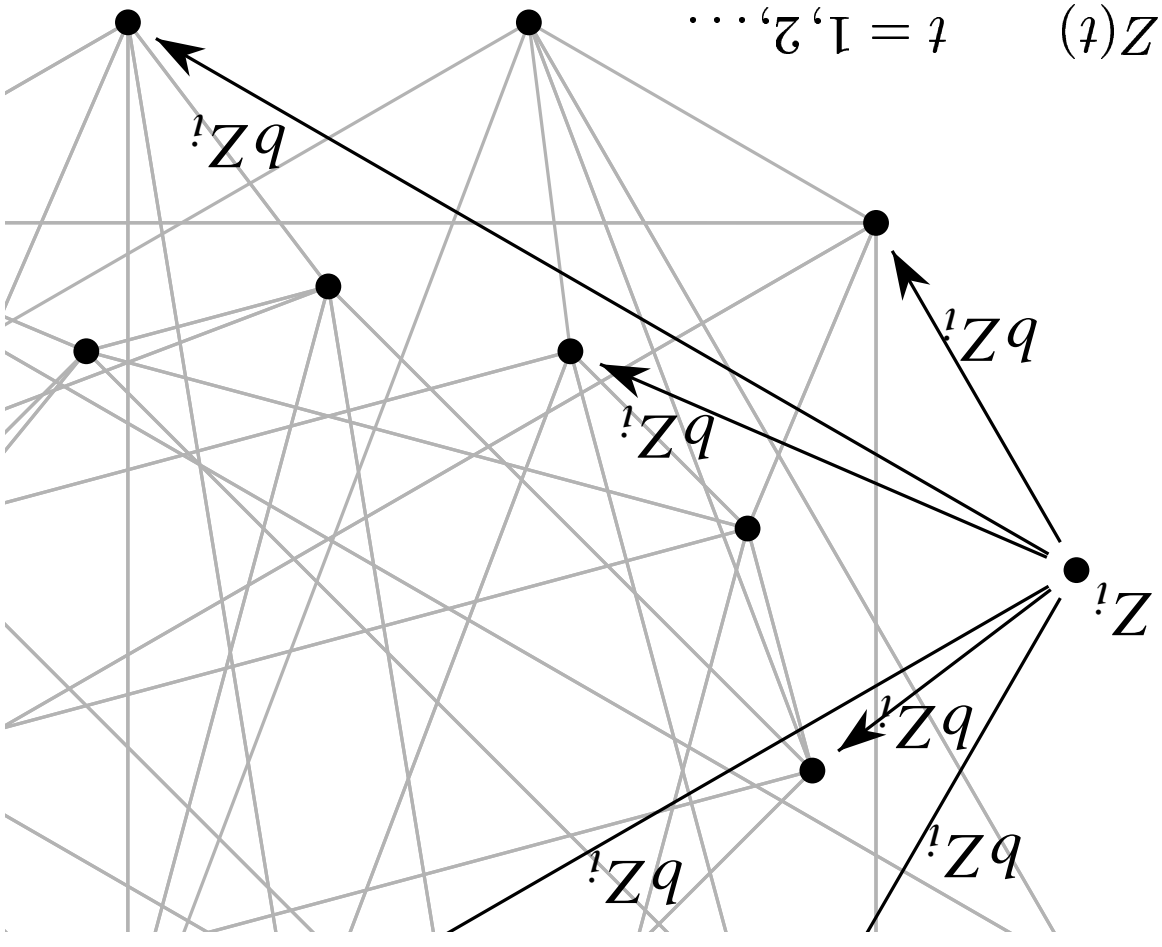
Exponential kernels

$$K_\beta = e^{\beta T} = \lim_{n \rightarrow \infty} \left(I + \frac{\beta T}{n} \right)^n$$

$$= I + \beta T + \frac{\beta^2 T^2}{2!} + \frac{\beta^3 T^3}{3!} + \dots$$

$$\frac{d}{dt} K_\beta = T K_\beta \quad K_0 = I$$

Diffusion kernels: view 2



$$K = e^{bL}$$

Stochastic field

$$E[Z_t(0)] = 0$$

$$\text{Var}[Z_t(0)] = \sigma^2$$

$Z_t(0)$ indep.

redistribution:

$$Z(t+1) = (I + bL)Z(t)$$

$t = 1, 2, \dots$

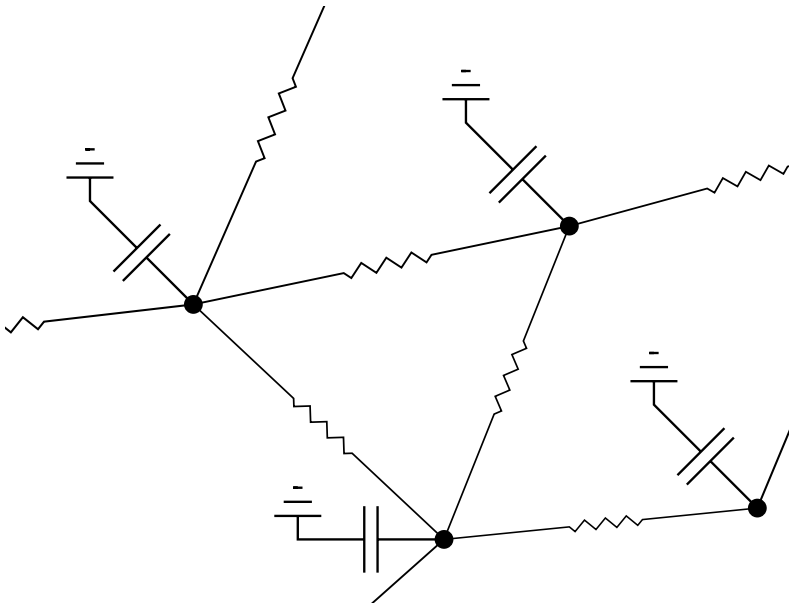
continuous limit: $\text{Cov}(t) = \lim_{n \rightarrow \infty} \left(I + \frac{bL}{n} \right)^{2nt} \sigma_2 \sigma_2^T e^{2btL}$

$$\text{Cov}(t) = \sigma_2^T J(t) J(t) \sigma_2 = \sigma_2^T J(t) (I + bL) J(t) \sigma_2$$

$$\text{Cov}_{ij}(t) = \underline{Z_i(t) Z_j(t)} = \underline{\left(\sum_{i'} T_{ii'} Z_{i'}(0) \right) \left(\sum_{j'} T_{jj'} Z_{j'}(0) \right)}$$

covariance

evolution operator $J(t) = I + bL$ $Z(t) = J(t) Z(0)$

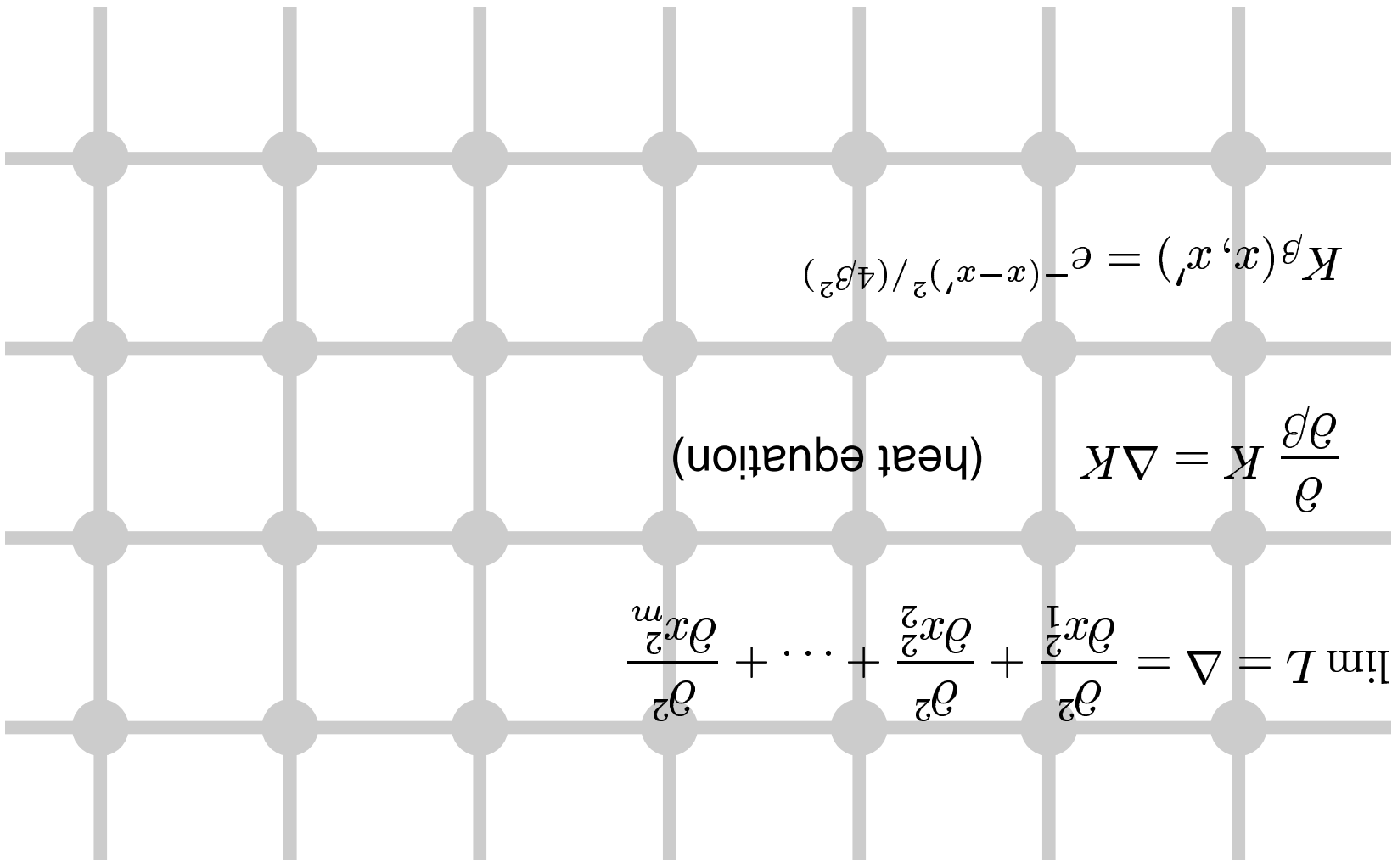


$$\sum_{i \sim j} \frac{RC}{1} (U_i(t) - U_j(t)) = (t) U_i(t) \frac{d}{dt}$$

$$\sum_{i \sim j} q (Z_i(t) - Z_j(t)) = (t) Z_i(t) \frac{d}{dt}$$

electrical analogy

Diffusion kernels: view 3



Diffusion kernels: view 4

Spectral graph theory

Eigenvectors of L correspond to “normal modes” of graph

approximate min-cut, max. distance type results

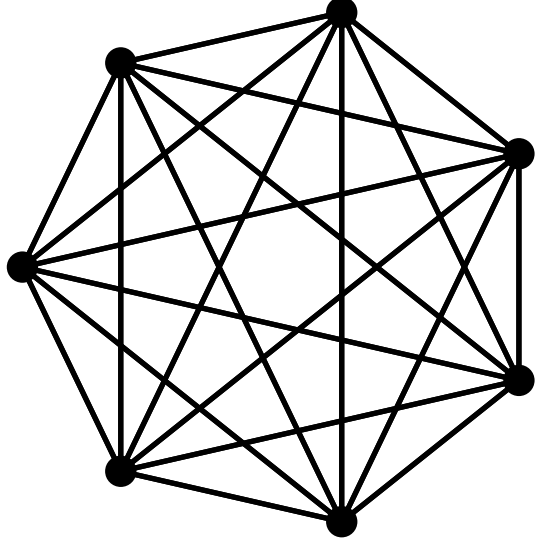
Computing the diffusion kernel

Diagonalization

$$L D L^{-1} = T \Lambda T^{-1}$$

$$L D L^{-1} = T$$

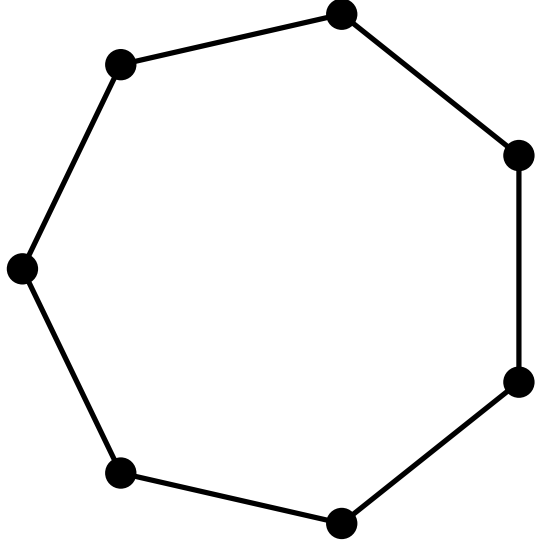
Complete graphs



$$K(i, j) = \left. \begin{array}{l} \frac{n}{1 - e^{-n\beta}} \text{ for } i \neq j, \\ \frac{n}{1 + (n-1)e^{-n\beta}} \text{ for } i = j. \end{array} \right\}$$

$$\text{for } n = 2 \quad K_{\beta}(i, j) \propto (\tanh \beta)^{d(i, j)}$$

Closed chains



$$K(i, j) = \sum_{l=0}^{n-1} \frac{u}{1-u} e^{-w_l \beta} \cos 2\pi l (i-j) \frac{u}{(j-i)}$$

Tensor product kernels

$K^{(1)}$ kernel on \mathcal{X}_1
 $K^{(2)}$ kernel on \mathcal{X}_2

$$K^{(1,2)} = K^{(1)} \otimes K^{(2)} \quad \text{kernel on } \mathcal{X}_1 \otimes \mathcal{X}_2$$

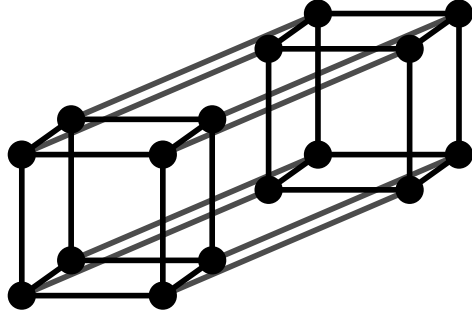
$$K^{(1,2)}((x_1, x_2), (x'_1, x'_2)) = K^{(1)}(x_1, x'_1) K^{(2)}(x_2, x'_2)$$

$$T^{(1,2)} = T^{(1)} \otimes T^{(2)} + T^{(2)} \otimes T^{(1)}$$

Hypercubes, etc.

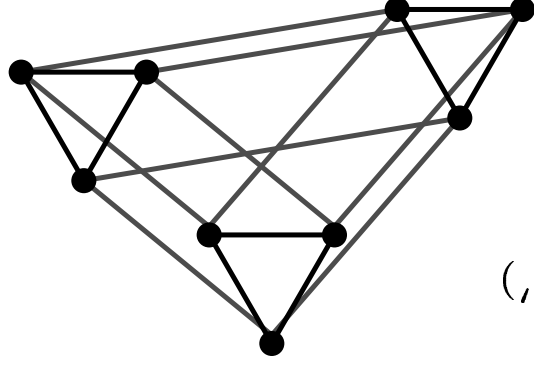
Hypercube:

$$K(x, x') = (\tanh \beta)^d = (\tanh \beta)^{d(x, x')}$$



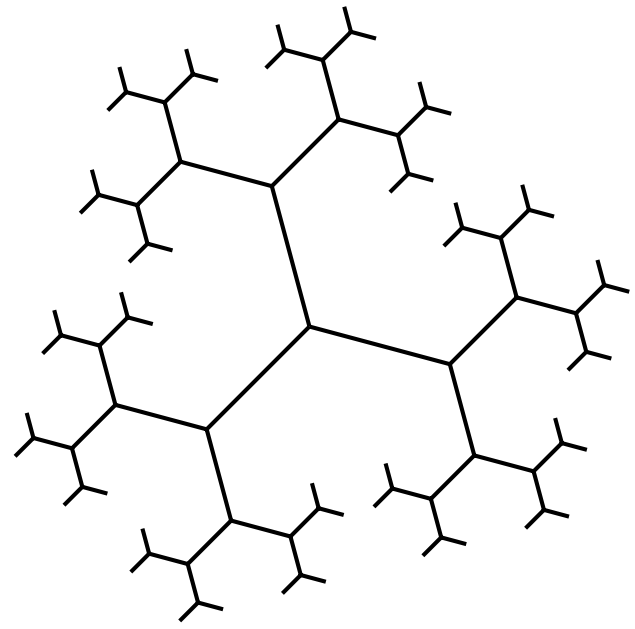
Alphabet \mathcal{A} :

$$K(x, x') = \frac{1 + (|\mathcal{A}| - 1) e^{-|\mathcal{A}| \beta}}{1 - e^{-|\mathcal{A}| \beta}} = d(x, x')$$



$$xp \frac{k_2 \cos^2 x (1-y) \sqrt{1-y}}{[\sin x (k-1) \sin(d+1) x - \sin(d-1) x]} \int_0^{\pi} \frac{(1-y)x}{2} e^{-\beta \left(1 - \frac{y}{2\sqrt{1-y}}\right) \cos x} dx$$

$$= K(x, x') p(x, x) = K(x, x')$$



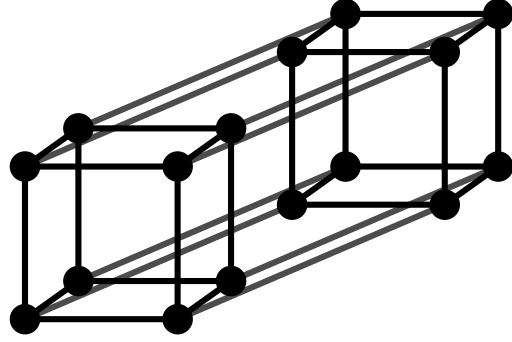
k-regular-trees

Applications

Experiments

Voted Perceptron algorithm (Freund
and Schapire, 1999)

5 standard categorical datasets



Diffusion kernels in Bioinformatics

Graph-driven feature extraction from microarray data using diffusion kernels and kernel CCA (J.-P. Vert and M. Kanehisa, NIPS 2002)

Diffusion kernel on network of chemical pathways, genes catalyzing them.

A more general framework

Looking for $f : \mathcal{X} \mapsto \mathcal{Y}$ minimizing

$$R_{\text{reg}}[f] = \frac{1}{m} \sum_{i=1}^m L(y_i, f(x_i)) + \langle \underline{Q}f, \underline{Q}f \rangle$$

$$\langle \underline{g}, \underline{g} \rangle = \int \int_{\mathcal{X}} g(x) K(x, x') g(x') dx' dx = \langle \underline{Q}f, \underline{Q}f \rangle$$

↑

$$\langle m, m \rangle = \sum_m \sum_{j=1}^i \alpha_j K(x_j, x_m)$$

For K symmetric, positive definite, \underline{K} is self-adjoint ($\langle \underline{K}g, g \rangle = \langle g, \underline{K}g \rangle$) and positive ($\langle \underline{K}g, g \rangle \geq 0$).

$$f(x) = \sum_{i=1}^i \alpha_i K(x, x_i) \quad \leftarrow \quad f(x) = \underline{K}g(x) = \int_{\mathcal{X}} K(x, x') g(x') dx'$$

Linear Operators on $\mathcal{L}^2(\mathcal{X})$

Linear Operators on $\mathcal{L}^2(\mathcal{X})$

$$\begin{aligned}\langle \underline{Q}f, \underline{Q}f \rangle &= \langle f, \underline{K}g \rangle \\ \langle \underline{Q}Kg, \underline{Q}Kg \rangle &= \langle g, \underline{K}g \rangle \\ \underline{Q} &= \underline{K}^{-1/2}\end{aligned}$$

\underline{K} and \underline{Q} share an eigensystem v_1, v_2, \dots and their eigenvalues are related by $\lambda_{(Q)}^2 = \lambda_{(K)}^{-1/2}$ (Gerosi, Jones and Poggio, 1995).

Locality and invariance

$\underline{\mathcal{Q}}$ is supposed to capture roughness of f . Want to make it a local differential operator.

Essentially unique second order differential operator on a d dimensional linear space or manifold \mathcal{X} invariant under isometries is the Laplacian $\Delta : \mathcal{L}^2(\mathcal{X}) \mapsto \mathcal{L}^2(\mathcal{X})$

$$\Delta = \frac{\partial^2}{\partial x_1^2} + \frac{\partial^2}{\partial x_2^2} + \dots + \frac{\partial^2}{\partial x_p^2}$$

Natural building block for $\underline{\mathcal{Q}}$ and hence \underline{K} .

The Gaussian Kernel

$$K(x, x') = e^{-\frac{1}{2} \|x - x'\|_2^2} \leftarrow \frac{\partial}{\partial x} K(x, x_0) = \Delta K(x, x_0) \cdot$$

$$\int_x K(x, x') g(x') dx' = (x) g(x)$$

$$\int_x \frac{\partial}{\partial x} K(x, x') g(x') dx' = (x) g(x)$$

$$\int_x \Delta K(x, x') g(x') dx' = (x) g(x)$$

$$\frac{\partial}{\partial x} \Delta K = \Delta \frac{\partial}{\partial x} K$$

$\Delta K = \Delta e^{-\frac{1}{2} \|x - x'\|_2^2} = \Delta \hat{Q}$ Just a diffusion kernel!

Diffusion Kernels on Manifolds

Belkin and Niyogi (NIPS 2001, 2002)

$$\underline{K} = e^{\beta \Delta}$$

on a Riemannian manifold found **from the data**

Natural way of incorporating unlabelled data.

Diffusion Kernels on the Statistical Manifold

Information Diffusion Kernels (Lafferty and Lebanon NIPS 2002)

parametric family $p(x|\theta)$ gives rise to manifold with Fisher metric

$$[G_{\theta}]_{i,j} = \mathbb{E}[\partial_i \ell_{\theta} \partial_j \ell_{\theta}] = \int_{\mathcal{X}} \partial_i \log p(x|\theta) \partial_j \log p(x|\theta) p(x|\theta) dx$$

For $x_i \in \mathcal{X}$, compute corresponding $\hat{\theta}_i$ (say, MLE); then use diffusion kernel on manifold between $\{\hat{\theta}_i\}$

$$L_{ij} = \begin{cases} 1 & i \sim j \\ -d_i & i = j \\ 0 & \text{otherwise} \end{cases}$$

$$g^T L g = \sum_{i \sim j} (g_i - g_j)^2$$

$$\int_{\mathcal{X}} g(x) \nabla g(x) dx = \langle g, \nabla g \rangle = \int_{\mathcal{X}} \nabla g(x) \cdot \nabla g(x) dx$$

From continuous to discrete

Other locally induced kernels

Smola and Kondor (under review)

$$\Delta = \sum_{i=1}^n \chi_i u_i u_i^\top$$

$$K = \sum_{i=1}^n r(\chi_i) u_i u_i^\top$$

- Diffusion kernel: $r(\lambda) = \exp(\beta\lambda)$
- d -step random walk kernel: $-(a + \lambda)^{-d}$
- Regularized Laplacian kernel: $\sigma^2 \lambda - 1$

Conclusions

- Inducing a kernel from the local (differential) structure of \mathcal{X} makes sense from a regularization theory point of view

- Inducing a kernel from the local (differential) structure of \mathcal{X} makes sense from a regularization theory point of view
- A unifying concept: the Laplacian

- Inducing a kernel from the local (differential) structure of \mathcal{X} makes sense from a regularization theory point of view
- A unifying concept: the Laplacian
- Diffusion kernels physically motivated, others are possible

- Inducing a kernel from the local (differential) structure of \mathcal{X} makes sense from a regularization theory point of view
- A unifying concept: the Laplacian
- Diffusion kernels physically motivated, others are possible
- These kernels are not necessarily easy to compute in practice