

# T-122.102 Special Course in Information Technology: Co-occurrence methods in analysis of discrete data. Information diffusion kernels

Sven Laur  
swen@math.ut.ee, slaur@tcs.hut.fi

April 2, 2004

## 1 Introduction

Information geometry is a highly abstract branch of mathematics that unifies statistical and geometrical concepts. The discipline provides general framework of many statistical phenomenas. The main drawbacks of the theory is generality—the approach seldom provides novel ideas and techniques. The survey covers one of the novelties—information diffusion kernels [6, 5]. Of course, we cannot exhibit the whole internal beauty of information geometry, but we encourage the interested reader to browse detailed monographs [7, 1].

The idea behind the information diffusion kernels is surprisingly simple. Each data point is associated with a distribution from some predefined family of distributions. The proximity measure in the feature space is defined rather over distributions themselves than on the parameterization. The association and the proximity measure together form a kernel. The modular definition allows reuse notions from differential geometry to analyze statistical properties like asymptotic consistency.

The brief outline of the survey is following. First, we discuss how to bind data with probability distributions. Next we cover the concepts and motivation behind heat-diffusion kernels and define heat-diffusion over probability distributions. Finally, we cite results of experiments [6] that show effectiveness of the approach.

## 2 From data to distribution

Traditional kernel methods are biased towards continuous data, i.e. most of the kernels assume that data comes from a subset of  $\mathbb{R}^n$ . As the real-world data is often varying and discrete like text documents, DNA and protein sequences, one has to use some *ad hoc* technique to find an embedding into  $\mathbb{R}^n$ . Frequently used embeddings are based on generative models, i.e. a discrete data object is assumed to be generated by a stochastic process. Two natural embeddings that map data to probability distribution are maximum likelihood (ML) and maximum a posteriori (MAP) estimates. Traditionally the corresponding parameter vector  $\theta$  is interpreted as an element of  $\mathbb{R}^n$  and continuous kernels are used for further inference. Hence, different parameterizations can lead to different results, although the corresponding model itself is same. Recently, several parameterization independent techniques like Fisher kernels [2] and mutual information kernels [8] have been proposed besides information diffusion kernels. Nevertheless, they can be viewed in the information geometry framework.

Text classification will be our central example. As usual we employ a bag of words approach. Namely, each document is represented by keyword count vector  $\mathbf{x}$ . The simplest generative model of text is multinomial distribution with parameters  $\theta$ . Clearly, the term frequency representation

$$\hat{\theta}_{\text{tf}}(\mathbf{x}) = \frac{1}{x_1 + \dots + x_n}(x_1, \dots, x_n)$$

corresponds to the ML estimate. Second frequently used embedding is inverse document frequency weighting

$$\widehat{\theta}_{\text{tfidf}}(\mathbf{x}) = \frac{1}{x_1 w_1 + \dots + x_n w_n} (x_1 w_1, \dots, x_n w_n)$$

where the weight  $w_i = \log(1/f_i)$  is logarithm of the inverse term frequency in document collection.

### 3 Statistical manifolds

The main quest of information geometry is to derive parameterization independent quantities over families of probability distributions. But first we have to define the notion of statistical manifold. Let  $\mathcal{X}$  be the domain of all possible values. Then a statistical manifold  $\mathcal{P}$  is a parameterized family of probability distributions

$$\mathcal{P} = \{p(\cdot|\theta) : \mathcal{X} \rightarrow R \mid \theta \in \Theta\},$$

where  $\Theta$  is open subset of  $\mathbb{R}^n$ . The parameterization must be unique:  $p(\cdot|\theta_1) \equiv p(\cdot|\theta_2) \Rightarrow \theta_1 = \theta_2$ . Hence,  $\theta$  can be treated as the coordinate vector of  $p(\cdot|\theta)$ . We say that parameterization  $\psi$  is admissible iff  $\psi$  as a function of primary parameters  $\theta$  is  $C^\infty$  smooth. It is easy to see that the set of admissible parameterizations does not change, if we take arbitrary admissible parameterization as a primary.

We consider only manifolds where log-likelihood function  $\ell(\mathbf{x}|\theta) = \log p(\mathbf{x}|\theta)$  is  $C^\infty$  differentiable w.r.t.  $\theta$ . Note that the property is invariant under all admissible parameterizations. For example, multinomial family satisfies the  $C^\infty$  requirement, since

$$\ell(\mathbf{x}|\theta) = \log \prod_{j=1}^m \theta_{x_j} = \sum_{j=1}^m \log \theta_{x_j}.$$

For proximity, we need also vectors and distance measure. The distance between two points is defined as length of the shortest path and the length itself is defined via integral and Euclidean norm

$$d(p, q) = \int_0^1 \|\dot{\gamma}(t)\| dt = \int_0^1 \sqrt{\langle \dot{\gamma}(t), \dot{\gamma}(t) \rangle} dt,$$

where  $\gamma : [0, 1] \rightarrow \mathcal{P}$  corresponds to a (simple) path and  $\dot{\gamma}(t)$  to a tangent vector. Since the family  $\mathcal{P}$  *a priori* does not have any geometrical structure, we need an abstract definition of vectors. Formally, a vector will be a function that maps functions with the type  $\mathcal{P} \rightarrow \mathbb{R}$  to real numbers. For fixed coordinates  $\theta$  and point  $p$  natural maps  $(\frac{\partial}{\partial \theta_i})_p$  emerge

$$\left(\frac{\partial}{\partial \theta_i}\right)_p (f) = \frac{\partial f}{\partial \theta_i} \Big|_p.$$

Let us denote lines with only varying coordinate  $\gamma_i(t) = p(\cdot|\theta_1, \dots, \theta_i + t, \dots, \theta_n)$ . Then the derivative of  $f(\gamma_i(t))$  at point  $p$  will be exactly  $(\frac{\partial}{\partial \theta_i})_p (f)$ . Moreover, if the path  $\gamma$  is differentiable, we can always decompose  $f(\gamma(t))'$  as linear combination of partial derivatives

$$f(\gamma(t))' = \left[ \theta_1(t)' \left(\frac{\partial}{\partial \theta_1}\right)_{\gamma(t)} + \dots + \theta_n(t)' \left(\frac{\partial}{\partial \theta_n}\right)_{\gamma(t)} \right] (f)$$

As the operator in the square brackets does not depend on  $f$  and has a right type, we baptize it as the speed vector  $\dot{\gamma}(t)$ . For a fixed coordinate system  $\theta$ , we can express  $\gamma(t) = (\theta_1(t), \dots, \theta_n(t))$  and we have a natural representation

$$\dot{\gamma}(t) \mapsto (\dot{\theta}_1(t), \dots, \dot{\theta}_n(t)) \in \mathbb{R}^n,$$

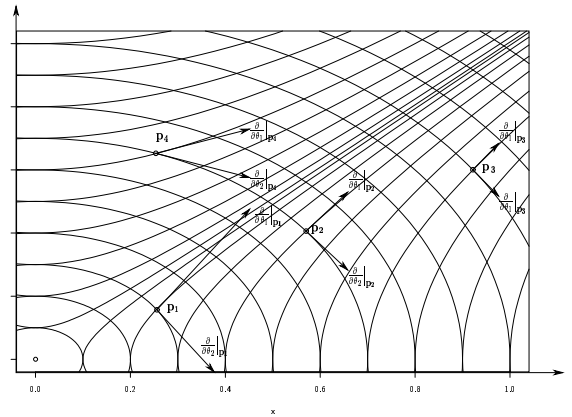


Figure 1: Non-Cartesian coordinate system and corresponding tangent vectors

but this is not coordinate independent. Moreover, the basis  $(\frac{\partial}{\partial \theta_i})_p$  may change along the manifold (see Figure 1). Remarkably, the abstract tangent vector does not depend on coordinates. In case of  $\mathbb{R}^n$ , we can interpret tangent vectors  $\dot{\gamma}(t)$  as usual vectors of  $\mathbb{R}^n$ . But without supporting geometrical structure the generalized speed has retained only the most important feature—for each function  $f : \mathcal{P} \rightarrow \mathbb{R}$ , the speed uniquely characterizes the rate of change at point  $\gamma(t)$ .

To fix a geometry, we need also reasonable definition of metric. A usual derivation of metric on statistical manifolds is too involved and abstract for the survey. Hence, we take a shortcut. What should be the distance of two adjacent distributions  $p$  and  $q$ ? The most reasonable answer is the weighted Kullback-Leibler divergence

$$\begin{aligned} J(p, q) &= D_{p\|q} + D_{q\|p} \\ &= \int p(\mathbf{x}) \log \frac{p(\mathbf{x})}{q(\mathbf{x})} d\mathbf{x} + \int p(\mathbf{x}) \log \frac{p(\mathbf{x})}{q(\mathbf{x})} d\mathbf{x}, \end{aligned}$$

since it quantifies average additional utility if we swap the distributions. Now consider an infinitesimal movement along the curve  $\gamma(t)$ . The corresponding change of coordinates is from  $\theta$  to  $\theta + \dot{\theta}\Delta t$  and the distance formula gives

$$d(p, q)^2 \approx \Delta t^2 \|\dot{\gamma}(t)\|^2 = \Delta t^2 \sum_{i,j=1}^n \dot{\theta}_i \dot{\theta}_j \left\langle \frac{\partial}{\partial \theta_i}, \frac{\partial}{\partial \theta_j} \right\rangle$$

On the other hand, we expect  $d(p, q)^2 = J(p, q)$  and under mild regularity conditions (see [4, p.26–28])

$$J(p, q) \approx \Delta t^2 \sum_{i,j=1}^n \dot{\theta}_i \dot{\theta}_j g_{ij},$$

where  $g_{ij}$  are the Fisher information matrix entries

$$g_{ij} = \int p(\mathbf{x}) \cdot \frac{\partial \ell(\mathbf{x}|\theta)}{\partial \theta_i} \cdot \frac{\partial \ell(\mathbf{x}|\theta)}{\partial \theta_j} d\mathbf{x}.$$

Hence, the scalar product is defined via Fisher information matrix

$$\left\langle \frac{\partial}{\partial \theta_i}, \frac{\partial}{\partial \theta_j} \right\rangle = g_{ij}.$$

Now, the geometry of the statistical manifold is complete and we proceed with a discussion about heat diffusion kernels.

## 4 Why heat diffusion?

The most suitable kernels in statistics are Mercer kernels. A Mercer is defined via (continuous) transformation  $\phi$  from a data space to a high dimensional Euclidean space. Mercer kernels are used together with support vector machines. Recall, that support vector machines use a transformation  $\phi$  to convert data to linearly separable sets and the Mercer kernel  $K(\mathbf{x}, \mathbf{y}) = \langle \phi(\mathbf{x}), \phi(\mathbf{y}) \rangle$  allows to calculate necessary quantities without computing the map  $\phi$ .

In first glance, the geodesic distance  $d(p, q)$  seems a natural starting point for a suitable Mercer kernel. Namely, the map  $\phi$  should preserve distances

$$d(\mathbf{x}, \mathbf{y})^2 = \langle \phi(\mathbf{x}) - \phi(\mathbf{y}), \phi(\mathbf{x}) - \phi(\mathbf{y}) \rangle$$

or equivalently

$$K(\mathbf{x}, \mathbf{y}) = \frac{d^2(\mathbf{x}, \mathbf{y}) - d^2(\mathbf{x}, \mathbf{o}) - d^2(\mathbf{y}, \mathbf{o})}{2}, \quad \phi(\mathbf{o}) = \mathbf{0}.$$

But generally there are no continuous transformations  $\phi$  that would embed the manifold into Euclidean space so that the distances are preserved. For example consider a sphere. There are infinitely many shortest paths between antipodes, but in the Euclidean space there is only one shortest path between two points. Hence, no continuous transformation to Euclidean space can preserve metrics.

A heat diffusion kernel seems a good alternative. First, kernels still capture the metric of the manifold. Secondly, they coincide with Gaussian kernels in  $\mathbb{R}^n$  and thirdly kernels that are closely related with heat diffusion have been used before [3]. If a heat diffusion kernel is used with support vector machines, we get a remarkably clear interpretation. A temperature of a manifold point is defined as a result of heat diffusion process, where small areas around training points have initial temperature proportional to labeling and importance in classification. The initial temperature is set to zero in all other areas. The classification rule divides points to two classes: “hot” and “cool” points. Intuitively, close points in the manifold have very similar temperature and will have same labels. Since the temperature equalizes over the time, a wide spectrum of kernels from very sensitive to completely robust are available.

The heat diffusion is governed by partial differential equations

$$\begin{aligned} \frac{\partial f}{\partial t} - \Delta f &= 0 \\ f(x, 0) &= f(x) \end{aligned} \quad (1)$$

augmented with suitable boundary conditions, if the manifold has a boundary. The general Laplace operator is rather complex

$$\Delta f = \det G^{-1/2} \sum_{i,j=1}^n \frac{\partial}{\partial \theta_j} \left[ g^{ij} \det G^{1/2} \frac{\partial f}{\partial \theta_i} \right]$$

where  $g^{ij}$  are elements of the inverse Fisher matrix  $G$ . Note that the operator  $\Delta$  is completely determined by the geometry of the manifold. In the Euclidean space  $\mathbb{R}^n$  the definition simplifies to usual notation

$$\Delta f = \frac{\partial^2 f}{\partial x_1^2} + \dots + \frac{\partial^2 f}{\partial x_n^2}.$$

and it is straightforward to verify

$$f(x, t) = (4\pi)^{-n/2} \int \exp\left(-\frac{\|\mathbf{x}-\mathbf{y}\|^2}{4t}\right) f(\mathbf{y}) d\mathbf{y}$$

is a solution to heat diffusion problem with the initial temperature  $f$ . The inner term

$$K_t(\mathbf{x}, \mathbf{y}) = \exp\left(-\frac{\|\mathbf{x}-\mathbf{y}\|^2}{4t}\right)$$

corresponds to Gaussian kernel and the discriminant function of SVM corresponds to heat diffusion from point sources  $\mathbf{x}_i$ . In that sense, heat diffusion on statistical manifolds is simple generalization of Gaussian kernels.

Formally, a heat diffusion kernel is a parametric function  $K : \mathcal{P} \times \mathcal{P} \times \mathbb{R} \rightarrow \mathbb{R}$  such that for all initial conditions  $f$  the solution of equations (1) can be obtained as an integral

$$f(x, t) = \int K(\mathbf{x}, \mathbf{y}, t) f(\mathbf{y}) d\mathbf{y}.$$

The existence of such kernel is guaranteed.

**Theorem 1.** *Let  $M$  be a complete Riemannian manifold. Then there exists a kernel function  $K$  (heat kernel), which satisfies the following properties:*

- (1)  $K(\mathbf{x}, \mathbf{y}, t) = K(\mathbf{y}, \mathbf{x}, t)$ ;
- (2)  $\lim_{t \rightarrow 0} K(\mathbf{x}, \mathbf{y}, t) = \delta(\mathbf{x}, \mathbf{y})$ ;
- (3)  $(\Delta - \frac{\partial}{\partial t})K(\mathbf{x}, \mathbf{y}, t) = 0$ ;
- (4)  $K(\mathbf{x}, \mathbf{y}, t) = \int K(\mathbf{x}, \mathbf{z}, t-s)K(\mathbf{z}, \mathbf{y}, s)dz$ .

The completeness assumption is equivalent to the assumption that all bounded closed sets on the manifold are compact. In other words, all limits remain in the manifold and every sequence has a convergent subsequence. Or equivalently  $J(p, q) \rightarrow 0$ , if  $q$  converges parameter-wise to  $p$ . The assumption is rather weak and is fulfilled by many distribution families. Properties (2) and (3) assure that the function  $f(\mathbf{x}, t)$  is a sought solution to equations (1). Properties (1) and (4) assure that  $K_t(\mathbf{x}, \mathbf{y}) = K(\mathbf{x}, \mathbf{y}, t)$  is a proper Mercer kernel, as a symmetric and positively definite operator.

## 5 Approximation of heat kernel

We have seen that the Gaussian kernels are also heat kernels. But closed forms of heat kernels are rare. Usually, we can only seek an approximation

$$\begin{aligned} K_t(\mathbf{x}, \mathbf{y}) &\approx K_t^{(m)} = (4\pi t)^{-n/2} \exp\left(-\frac{d^2(\mathbf{x}, \mathbf{y})}{4t}\right) \\ &\cdot [\psi_0(\mathbf{x}, \mathbf{y}) + \psi_1(\mathbf{x}, \mathbf{y})t + \dots + \psi_m(\mathbf{x}, \mathbf{y})t^m]. \end{aligned}$$

Intuitively, terms  $\psi_0, \dots, \psi_m$  correct distortion from a flat Euclidean geometry. Let  $r = d(\mathbf{x}, \mathbf{y})$ , then we can express  $\psi_i(\mathbf{x}, \cdot) = \psi_i(r)$  and the terms  $\psi_i$  can be computed by recursive equations

$$\begin{aligned} \psi_0 &= \left(\frac{\sqrt{\det G}}{r^{n-1}}\right)^{-1/2} \\ \psi_k &= r^{-k} \psi_0 \int_0^r \psi_0^{-1} (\Delta \psi_{k-1}) s^{k-1} ds, \end{aligned}$$

where the residue term is

$$\left(\Delta - \frac{\partial}{\partial t}\right)K_t^{(m)} = (t^m \Delta \psi_m)(4\pi t)^{-n/2} \exp(-r^2/4t).$$

Hence, first  $m$  terms allow to obtain approximation  $K_t(\mathbf{x}, \mathbf{y}) = K_t^{(m)}(\mathbf{x}, \mathbf{y}) + O(t^m)$  provided  $\psi_m$

is smooth. Of course, the approximation is valid for  $\varepsilon$ -neighborhood  $\varepsilon < 1$  and generally the approximation  $K_t^{(m)}$  is not positively definite unless  $t \in [0, \varepsilon)$ . Fortunately, estimates about  $\varepsilon$  can be obtained.

## 6 Geometry of multinomials

In principle, the sufficient approximations of heat kernel can be derived. However, practical details can be quite messy. Next, we derive appropriate kernel approximation for multinomial family. We consider the multinomial distribution with  $n + 1$  different outcomes. The usual parameterization

$$\Theta = \{(\theta_1, \dots, \theta_n) \in \mathbb{R}^n, \theta_i > 0, \theta_1 + \dots + \theta_n \leq 1\}$$

corresponds to  $n$  dimensional simplex. Let  $\mathbf{x} = (x_1, \dots, x_{n+1})$  be the indicator of a single draw, i.e.  $x_i = 1$  iff the  $i$ th event has happened. Then

$$\begin{aligned} \frac{\partial \ell(\mathbf{x}|\theta)}{\partial \theta_i} &= \frac{x_i}{\theta_i} - \frac{x_{n+1}}{\theta_{n+1}} \\ \frac{\partial^2 \ell(\mathbf{x}|\theta)}{\partial \theta_i \partial \theta_j} &= -\frac{x_i}{\theta_i^2} \delta_{i,j} - \frac{x_{n+1}}{\theta_{n+1}^2}, \end{aligned}$$

where  $\theta_{n+1} = 1 - \theta_1 - \dots - \theta_n$ . The equation

$$E_{\mathbf{x}} \left[ \frac{\partial \ell(\mathbf{x}|\theta)}{\partial \theta_i} \cdot \frac{\partial \ell(\mathbf{x}|\theta)}{\partial \theta_j} \right] = -E_{\mathbf{x}} \left[ \frac{\partial^2 \ell(\mathbf{x}|\theta)}{\partial \theta_i \partial \theta_j} \right]$$

allows to express Fisher matrix elements

$$g_{ij}(\theta_1, \dots, \theta_n) = \begin{cases} 1/\theta_{n+1}, & \text{if } i \neq j, \\ 1/\theta_i + 1/\theta_{n+1}, & \text{if } i = j. \end{cases}$$

It is rather hard to compute geodesic distances on the simplex  $\Theta \subset \mathbb{R}^n$  and therefore we define two mappings that ease the task.

A map  $F : \mathcal{P} \rightarrow \mathcal{Q}$  is called isometry if it satisfies two requirements: (1)  $F$  must be  $C^\infty$  differentiable w.r.t.  $\theta$ ; (2) all curves  $\gamma(t)$  and  $\gamma^*(t) = F(\gamma(t))$  must have same length. The second condition is satisfied iff  $\|\dot{\gamma}(t)\| = \|\dot{\gamma}^*(t)\|$ .

By adding redundant coordinate  $\theta_{n+1}$  and redefining geometry

$$g_{ij}(\theta_1, \dots, \theta_{n+1}) = \begin{cases} 0, & \text{if } i \neq j, \\ 1/\theta_i, & \text{if } i = j. \end{cases}$$

we get an isometric embedding  $i : \mathbb{R}^n \rightarrow \mathbb{R}^{n+1}$ , since the length of the vectors are preserved

$$\sum_{i,j=1}^n \dot{\theta}_i \dot{\theta}_j g_{ij} = \sum_{i=1}^n \frac{\dot{\theta}_i^2}{\theta_i} + \frac{1}{\theta_{n+1}} \left( \sum_{i=1}^n \dot{\theta}_i \right)^2 = \sum_{i=1}^{n+1} \frac{\dot{\theta}_i^2}{\theta_i}.$$

Intuitively, this means that tangent vectors  $\frac{\partial}{\partial \theta_i}$  are orthogonal. But we can go even further and build a model of  $\mathcal{P}$  in Euclidean space  $\mathbb{R}^{n+1}$  that preserves distances. The target set will be  $n + 1$  dimensional positive orthant

$$\mathcal{S}^+ = \{(x_1, \dots, x_{n+1}) : x_1^2 + \dots + x_{n+1}^2 = 4\}.$$

and the isometry

$$F(\theta_1, \dots, \theta_{n+1}) = (2\sqrt{\theta_1}, \dots, 2\sqrt{\theta_{n+1}}).$$

Since the tangent vector is

$$f(\gamma^*(t))' = \left[ \dot{x}_1 \cdot \left( \frac{\partial}{\partial x_1} \right)_{\gamma^*} + \dots + \dot{x}_{n+1} \cdot \left( \frac{\partial}{\partial x_{n+1}} \right)_{\gamma^*} \right] (f),$$

where  $x_1(t), \dots, x_{n+1}(t)$  are Cartesian coordinates of  $\gamma^*(t)$ , the lengths are preserved only if

$$\|\dot{\gamma}^*(t)\|^2 = \dot{x}_1(t)^2 + \dots + \dot{x}_{n+1}(t)^2 = \|\dot{\gamma}(t)\|^2.$$

A simple calculation proves isometry

$$\begin{aligned} \|\dot{\gamma}^*(t)\|^2 &= \left( \frac{\dot{\theta}_1(t)}{\sqrt{\theta_1(t)}} \right)^2 + \dots + \left( \frac{\dot{\theta}_{n+1}(t)}{\sqrt{\theta_{n+1}(t)}} \right)^2 \\ &= \frac{\dot{\theta}_1(t)^2}{\theta_1(t)} + \dots + \frac{\dot{\theta}_{n+1}(t)^2}{\theta_{n+1}(t)} = \|\dot{\gamma}(t)\|^2. \end{aligned}$$

The corresponding maps have a nice geometrical meaning, see Figure 2. The isometry provides novel insight for calculating distances between  $\theta$  and  $\theta'$ , also depicted by Figure 2. As the shortest path between two points is an arc of the great circle, the distance becomes easily computable

$$\begin{aligned} d(\theta, \theta') &= 2 \arccos(\langle F(\theta), F(\theta') \rangle) \\ &= 2 \arccos(\sqrt{\theta_1 \theta'_1} + \dots + \sqrt{\theta_{n+1} \theta'_{n+1}}). \end{aligned}$$

The corresponding metric is illustrated by Figure 3. The more rapid distance rate in the corners has intuitive explanation—slight changes in parameters cause significant change in relative entropy (utility), if some event is almost certain.

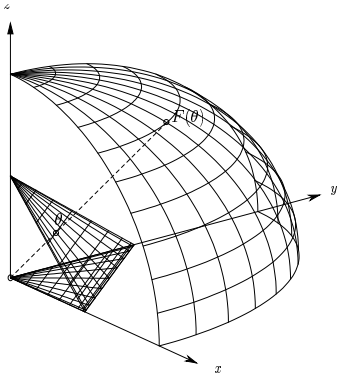


Figure 2: Geometrical interpretation of isometry

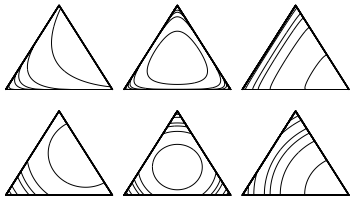


Figure 3: Iso-distant lines: above multinomial geometry, below Euclidean geometry.

## 7 Kernel construction

Finally, we are in position to define an approximation of the heat diffusion kernel. Nontrivial derivation shows

$$\psi_0(r) = 1 + \frac{(n-1)}{12}r^2 + \frac{(n-1)(5n-1)}{1440}r^4 + O(r^6)$$

and the first order approximation is

$$K_t^{(0)}(\theta, \theta') = (4\pi t)^{-n/2} \exp\left(-\frac{\arccos^2(\sqrt{\theta}, \sqrt{\theta'})}{t}\right).$$

The approximation has good properties compared with higher order approximations  $K_t^{(m)}$ ,  $m \geq 1$ , since  $K_t^{(0)}$  remains small for large distances.

As a theoretical side-mark the simplex is not a complete Riemann manifold. Moreover, the border has non-differentiable angles. Therefore, for theoretical investigation these angles must be rounded—certain and impossible events are removed from model.

This allows to find bounds of covering numbers  $\mathcal{N}(\epsilon, \mathcal{F}_R)$ , where  $\mathcal{F}_R$  corresponds to linear classifier in high dimensional space with bounded weight vector  $\|\mathbf{w}\| \leq R$ . The cover numbers allow to estimate the difference of average classification risk and empirical risk (training error). These results are based on advanced issues of statistical learning theory and functional analysis. Compared with the usual VC-dimension approach, the result are much sharper but in the same time more obscure. Thus we cite only the main result.

**Informal statement** *The heat-kernel have essentially the same asymptotic generalization performance as Gaussian kernels with the same dimension.*

## 8 Experimental results

We cite here results of the report [6]. A support vector machine with approximation of multinomial kernel was applied to Reuters-21570 and WebKB data. Both frequency representation and inverse document frequency weighting was used. In all reported instances heat kernel gave better results than Gaussian kernel. Some results are given in Figure 4.

## 9 Conclusions

Good kernels for real valued vectors have been known for years. Therefore, a big effort has been made to find natural ways to define kernels for discrete data generated by random processes. Information geometry gives a theoretically justified approach. The main advantages are independence of parameterization and theoretical machinery for estimating generalization properties. Heat kernels are also natural generalizations of Gaussian kernels. Empirical results with multinomial distributions indicate that heat diffusion kernel out-performs Gaussian kernels in text classification tasks.

The most serious downside is conceptual complexity. The general concept of statistical manifolds is not easy to grasp. Secondly, only few instances of closed form solutions of heat diffusion equation is known.

Hence, generalizing the results to more complicated models requires *significant* amount of mathematical knowledge.

## References

- [1] Shun-ichi Amari. *Methods of Information Geometry*, volume 191 of *Translations of Mathematical Monographs*. American Mathematical Society, 2000.
- [2] Tommi S. Jaakkola and David Haussler. Exploiting generative models in discriminative classifiers. In *Advances in Neural Information Processing Systems*, 1998.
- [3] Risi Imre Kondor and John D. Lafferty. Diffusion kernels on graphs and other discrete input spaces. In *Proceedings of the Nineteenth International Conference on Machine Learning*, pages 315–322. Morgan Kaufmann Publishers Inc., 2002.
- [4] Solomon Kullback. *Information Theory and Statistics*. Dover, 1997.
- [5] John Lafferty and Guy Lebanon. Information diffusion kernels. In *Advances in Neural Information Processing Systems*, 2002.
- [6] John Lafferty and Guy Lebanon. Diffusion kernels on statistical manifolds. Technical Report CMU-CS-04-101, School of Computer Science Carnegie Mellon University, 2004.
- [7] Michael K. Murray and John W. Rice. *Differential Geometry and Statistics*, volume 48 of *Monographs on Statistics and Applied Probability*. Chapman & Hall, 1993.
- [8] M. Seeger. Covariance kernels from bayesian generative models. In *Advances in Neural Information Processing Systems*, Cambridge, MA, 2002. MIT Press.

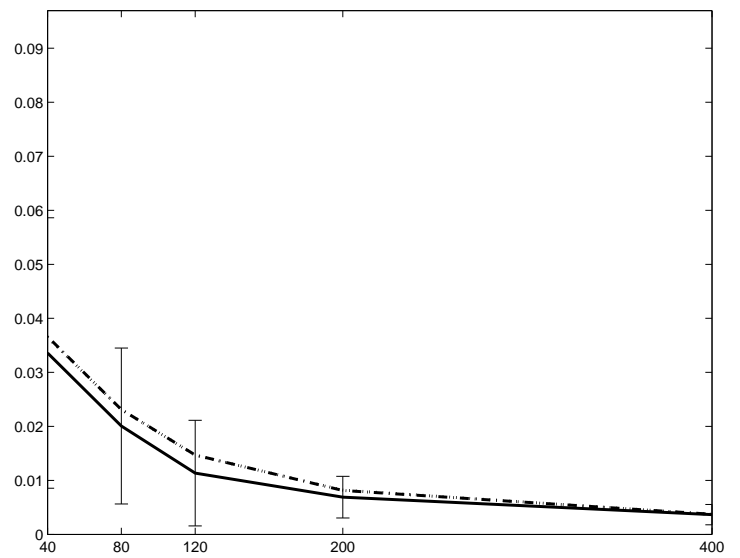
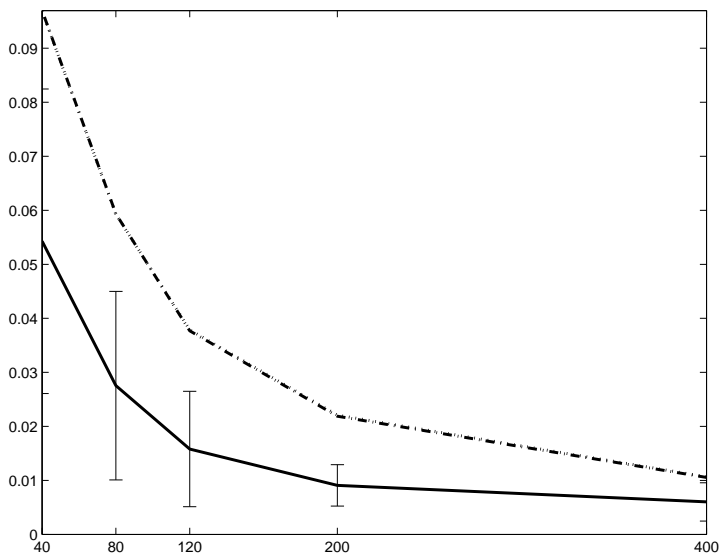
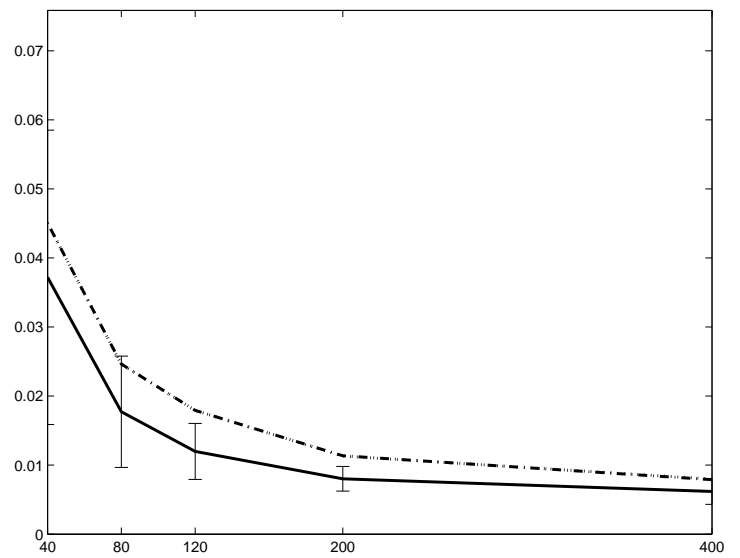
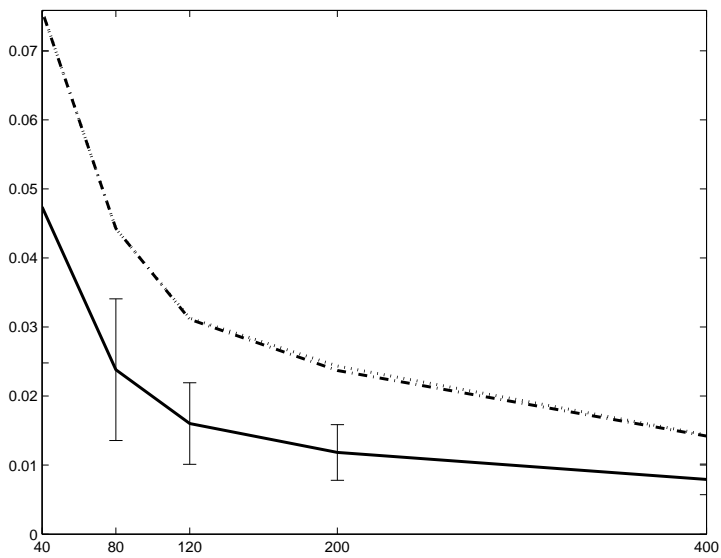


Figure 4: Experimental results on the Reuters corpus, using SVMs for linear (dotted) and Gaussian (dash-dotted) kernels, compared with the diffusion (solid). The classes moneyFx (top) and grain (bottom) are labeled as positive, and the class earn is labeled negative. The left column uses tf representation and the right column uses tfidf representation [6, p.28].