

Parameter estimation for discrete determinantal point processes

Dissertation submitted for the degree of
MASTER OF SCIENCE IN INTERDISCIPLINARY MATHEMATICS

Johannes Müller

July 28, 2018



SUPERVISED BY
PROFESSOR NIKOLAOS ZYGOURAS AND DR THEODOROS DAMOULAS

UNIVERSITY OF WARWICK
DEPARTMENT OF MATHEMATICS

ACKNOWLEDGEMENTS

I would like to thank the dreamteam for being such an eggcelent friend-group and for all the pun we had during the last year.

ABSTRACT

Determinantal point processes are random subsets that exhibit a diversifying behaviour in the sense that the randomly selected points tend to be not similar in some way. This repellent structure first arose in theoretical physics and pure mathematics, but they have recently been used to model a variety of many real world scenarios in a machine learning setup. We aim to give an overview over the main ideas of this approach which is easily accessible even without prior knowledge in the area of machine learning and sometimes omit technical calculations in order to keep the focus on the concepts.

Contents

I	Introduction and motivating examples	1
I.1	Motivation	1
I.2	Previous work	1
I.3	Aim and outline of the dissertation	1
II	Determinantal points processes: Basic notions and properties	2
II.1	Definitions and properties	2
II.2	Variations of DPPs	6
II.3	Simulation and Existence of DPPs	7
II.4	The mode problem	12
III	Point estimators and parametric models	14
III.1	Kernel reconstruction from the empirical measures	15
III.1.1	Graph theoretical concepts	17
III.1.2	The solution of the principal minor assignment problem	19
III.1.3	Definition of the estimator and consistency	22
III.1.4	Computation of the estimator	24
III.2	Maximum likelihood estimation using optimisation techniques	24
III.2.1	Presentation of different models	27
III.2.2	Coercivity and existence of the maximum likelihood estimators	29
III.2.3	Consistency of the maximum likelihood estimators	32
III.2.4	Approximation of the MLE	39
III.2.5	Learning for conditional DPPs	40
III.2.6	Estimating the mixture coefficients of k -DPPs	40
IV	Bayesian learning for DPPs	41
IV.1	Bayesian approach to parameter estimation	41
IV.2	Markov chain Monte Carlo methods	44
IV.2.1	Reminder on Markov chains	44
IV.2.2	Metropolis-Hastings random walk	47
IV.2.3	Slice sampling	51
IV.3	Variational MCMC methods	56
V	Toy examples and experiments	57
V.1	Minimal example?	57
V.2	Points on the line	57
V.3	Points in the square	59
V.4	Toy example for quality learning	59

Contents

VI	Summary and conclusion	60
A	Calculations	61
B	Generated code	62
B.1	Sampling algorithm	62
B.2	Points on the line	63
B.3	Points in the square	64
B.4	Toy learning example	66
	Bibliography	70

Chapter I

Introduction and motivating examples

I.1 Motivation

I.2 Previous work

I.3 Aim and outline of the dissertation

Chapter II

Determinantal points processes: Basic notions and properties

II.1 Definitions and properties

We begin by presenting general frame we will work in. This means that we will keep the notation introduced now and will use those objects throughout the thesis without further explanation. Further we will present all the important properties of determinantal point processes that we will need and postpone some calculations to the last section of this chapter. A much more ... survey of properties of determinantal point processes including extensive comparisons to several other point processes can be found in the report [Kulesza et al., 2012].

2.1 SETTING. Let in the following \mathcal{Y} be a finite set, which we call the *ground set* and $N := |\mathcal{Y}|$ its cardinality. For the sake of easy notation we will assume $\mathcal{Y} = \{1, \dots, N\}$ unless otherwise specified. A *point process* on \mathcal{Y} is a random subset of \mathcal{Y} , i.e. a random variable with values in the powerset $2^{\mathcal{Y}}$. We will identify this random variable with its law \mathbb{P} and thus refer to probability measures \mathbb{P} on $2^{\mathcal{Y}}$ as point processes and will not distinguish between those objects. Let further \mathbf{Y} denote a random subset drawn according to \mathbb{P} .

2.2 DEFINITION (DETERMINANTAL POINT PROCESS). We call \mathbb{P} a *determinantal point process*, or in short a DPP, if we have

$$\mathbb{P}(A \subseteq \mathbf{Y}) = \det(K_A) \quad \text{for all } A \subseteq \mathcal{Y} \quad (2.1)$$

where K is a symmetric matrix indexed by the elements in \mathcal{Y} and K_A denotes the submatrix of K to indexed by the elements of A . We call K the *marginal kernel* of the DPP. If the marginal kernel K is diagonal, then we call \mathbb{P} a *Poisson point process*.

We note that all principal minors¹ of K are non negative and Sylvester's criterion implies that K is non negative definite². Further it can be shown (cf. page 3 in [Borodin, 2009]) that also the complement of a DPP is a DPP with marginal kernel $I - K$ where I is the identity matrix, i.e.

$$\mathbb{P}(A \subseteq \mathbf{Y}^c) = \det(I_A - K_A).$$

¹The *principle minors* of K are the determinants of the submatrices K_A for $A \subseteq \mathcal{Y}$.

² K is called *non negative definite* if $x^T K x \geq 0$ for all $x \in \mathbb{R}^{\mathcal{Y}}$. The Sylvester criterion states that a matrix is non negative definite if and only if all principle minors are non negative.

have a look at this and maybe explain it!

Thus we conclude $I - K \geq 0$ and obtain $0 \leq K \leq I$. This actually turns out to be sufficient for K to define a DPP through (2.1) which we will see in the fourth section of this chapter.

2.3 REPULSIVE BEHAVIOUR OF DPPs. We call the elements of \mathcal{Y} *items* and by choosing $A = \{i\}$ and $A = \{i, j\}$ for $i, j \in \mathcal{Y}$ and using (2.1) we obtain the probabilities of their occurrence

$$\begin{aligned} \mathbb{P}(i \in \mathbf{Y}) &= K_{ii} \quad \text{and} \\ \mathbb{P}(i, j \in \mathbf{Y}) &= K_{ii}K_{jj} - K_{ij}^2 = \mathbb{P}(i \in \mathbf{Y}) \cdot \mathbb{P}(j \in \mathbf{Y}) - K_{ij}^2, \end{aligned} \tag{2.2}$$

Thus the appearances of the two items i and j are always negatively correlated. This negative correlation is exactly what causes the diversifying behaviour of determinantal point processes. In practice one usually models the negative correlations to be high between items that are similar in some notion. For example in a spatial setting being similar could mean being close together and in this case the selected items would tend to be not very close together. This is repulsive behaviour can be seen in Figure. Note that Poisson point processes are exactly the DPPs without correlations of the points.

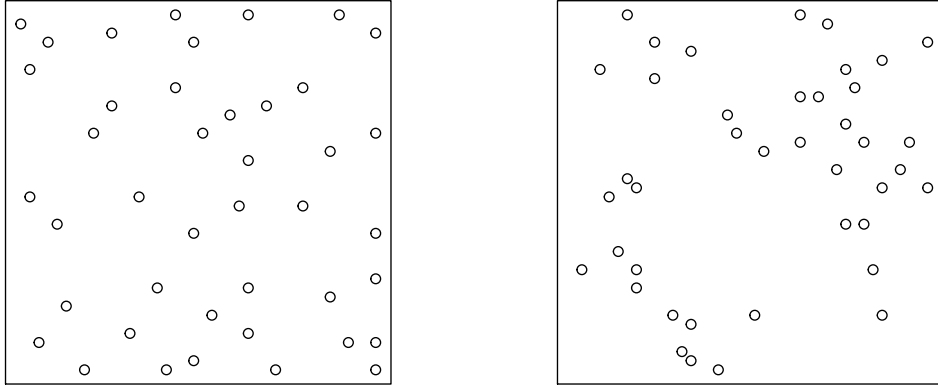


Figure II.1.: A DPP with negative correlations of close points on a 40×40 grid in the unit square on the left and a Poisson point process on the same grid on the right with the same expected cardinality. The – in this case spatially – repellent structure of the DPP is clearly visible.

In this light the fact that also \mathbf{Y}^c exhibits negative correlations becomes less surprising. Since the set \mathbf{Y} tends to spread out due to the repulsion in (2.2), the complement, which is nothing but the gaps that are left after eliminating the elements in \mathbf{Y} , tend to show a non clustering behaviour as well.

2.4 L-ENSEMBLES. Let us now introduce an important subclass of DPPs, namely the ones where not only the marginal probabilities can be expressed through a suitable kernel, but also the elementary probabilities. This will be convenient for us and lead to some explicit expression. If we

have even $K < I$, then we define the *elementary kernel*

$$L := K(I - K)^{-1} \quad (2.3)$$

which specifies the elementary probabilities since one can check

$$\mathbb{P}(A = \mathbf{Y}) = \frac{\det(L_A)}{\det(I + L)} \quad \text{for all } A \subseteq \mathcal{Y}. \quad (2.4)$$

Conversely for any $L \geq 0$ a DPP can be defined via (2.2) and the corresponding marginal kernel is given by the inversion of (2.3)

$$K = L(I + L)^{-1}$$

and we have again $K < I$. We call DPPs which arise this way *L ensembles*. Later we will see that the cardinality of a DPP distributed like the sum of N Bernoulli experiments with expectation $(\lambda_n)_{n=1,\dots,N}$ where λ_n are the eigenvalues of K . Being an L -ensemble is equivalent to $K < I$ which again is equivalent to $\lambda_n < 1$ for all $n = 1, \dots, N$ and hence equivalent to

$$\mathbb{P}(\mathbf{Y} = \emptyset) = \mathbb{P}(|\mathbf{Y}| = 0) > 0.$$

The quality diversity decomposition

Note that any symmetric, positive semidefinite matrix L can be written as a Gram matrix

$$L = B^T B$$

where $B \in \mathbb{R}^{D \times N}$ whenever D is larger than the rank $\text{rk}(L)$ of L . For example one could take the spectral decomposition $L = U^T C U$ of L and set $B := \sqrt{C} U$ and eventually drop some zero rows from \sqrt{C} . Let B_i denote the i -th column of B and write it as the product $q_i \cdot \phi_i$ where $q_i \geq 0$ and $\phi_i \in \mathbb{R}^D$ such that $\|\phi_i\| = 1$. This yields the representation

$$L_{ij} = q_i \phi_i^T \phi_j q_j =: q_i S_{ij} q_j$$

and we call q_i the *quality* of the item $i \in \mathcal{Y}$ and ϕ_i the *diversity feature vector* of i and S the *similarity matrix*. Since we will use this decomposition multiple times, we fix its properties.

2.5 PROPOSITION (QUALITY DIVERSITY PARAMETRISATION). *Let $D \in \mathbb{N}$ and let \mathbb{S}_D denote the sphere in \mathbb{R}^D . Further let $\mathbb{R}_{\text{sym},+}^{N \times N}$ be the set of symmetric positive semidefinite $N \times N$ matrices. The quality diversity parametrisation is a continuous and surjective mapping*

$$\Psi: \mathbb{R}_+^N \times \mathbb{S}_D^N \rightarrow \left\{ L \in \mathbb{R}_{\text{sym},+}^{N \times N} \mid \text{rk}(L) \leq D \right\}, \quad (q, \phi) \mapsto \left(q_i \phi_i^T \phi_j q_j \right)_{1 \leq i, j \leq N}.$$

2.6 REMARK. (i) In the case $D = N$ the quality diversity decomposition gives a parametrisation of the whole symmetric positive definite $N \times N$ matrices.

(ii) Note that this parametrisation is not unique, i.e. Ψ is not injective. For example the identity matrix I can be parametrised by any orthonormal system ϕ and $q = (1, \dots, 1)^T$.

(iii) One can without any problems consider diversity features ϕ_i in an abstract Hilbert space \mathcal{H} . However we will not need this in the remainder and thus restrict ourselves to the easier case Euklidean diversity features.

- (iv) We call every preimage (q, ϕ) of L under Ψ *quality diversity decomposition* of L . Further we call the tuple $\phi \in \mathbb{S}_D^N$ of normalised vectors *diversity feature matrix*.

The quality diversity decomposition will provide some useful expressions. For example the elementary probabilities take the form

$$\mathbb{P}(A = \mathbf{Y}) \propto \det((B^T B)_A) = \left(\prod_{i \in A} q_i^2 \right) \cdot \det(S_A) \quad \text{for all } A \subseteq \mathcal{Y}. \quad (2.5)$$

An intuitive understanding of the quality diversity decomposition will play a central role in the modelling process of real world phenomena through DPPs. To get this we can think of $q_i \geq 0$ as a measure of how important or high in quality the item is and the diversity feature vector $\phi_i \in \mathbb{R}^D$ can be thought of as some kind of state vector that consists of internal quantities that describe the item i in some way. Further we interpret the scalar product $\phi_i^T \phi_j \in [0, 1]$ as a measure of similarity between the items i and j which justifies the name similarity matrix for S . Note that if i and j are perfectly similar or antisimilar, i.e. $\phi_i^T \phi_j = \pm 1$, then they can not occur at the same time, since

$$\mathbb{P}(i, j \in \mathbf{Y}) = \det \begin{pmatrix} 1 & \pm 1 \\ \pm 1 & 1 \end{pmatrix} = 0.$$

If we identify i with the vector $B_i = q_i \phi_i \in \mathbb{R}^D$, we can obtain a geometric interpretation of (2.5) since $\det((B^T B)_A)$ is the volume that is spanned by the columns $B_i, i \in A$, which is visualised in II.2. This volume increases if the lengths of the edges that correspond to the quality increase and decrease when the similarity feature vectors point into more similar directions.



Figure II.2.: Taken from [Kulesza et al., 2012]; the first line (a) illustrates the volumes spanned by vectors, and in the second line it can be seen how this volume increases if the length – associated with the quality – increases (b) and decreases if they become more similar in direction which we interpret as two items becoming more similar (c)

2.7 MODELLING DIVERSITY OVER DISTANCE. Since we will use one form of diversity features multiple times, we will now give a short general formulation of it. Let $\mathcal{R} = \{r_1, \dots, r_D\}$ be a finite set which we will call the *reference set* and its elements the *reference points*. Further let

$$d: \mathcal{Y} \times \mathcal{R} \rightarrow \mathbb{R}_+, \quad f: \mathbb{R}_+ \rightarrow \mathbb{R}$$

mappings. Usually $d(i, r)$ will be interpreted as a measure of distance between an item $i \in \mathcal{Y}$ and a reference point $r \in \mathcal{R}$ and will typically be given by a metric on a larger space that contains both \mathcal{Y} and \mathcal{R} . One can now model $\phi_i \in \mathbb{R}^{\mathcal{R}}$ via

$$(\phi_i)_r \propto f(d(i, r)) \quad \text{for } r \in \mathcal{R}$$

The function f will typically be decreasing and thus $(\phi_i)_r$ can be seen as a measure of how similar item i is to the reference point $r \in \mathcal{R}$. Thus the diversity feature vector ϕ_i stores how similar the item i is to all reference points and the scalar product $\phi_i^T \phi_j$ will be close to one, if the items i and j have approximately the same degrees of similarity to the reference points. It shall be noted that the choice of the D , the number of reference points bounds the rank of the kernel L and therefore of the largest subset that occurs with positive probability. Indeed we have $\text{rk}(L) \leq D$ and for $A \subseteq \mathcal{Y}$ with more than D elements $\det(L_A) = 0$ and therefore $\mathbb{P}(A) = 0$. In the fourth chapter we will see that there is a natural choice for the mapping d in most cases, at least in the ones where \mathcal{Y} consists of points in a metric space. On the other hand the choice of f is crucial since it determines the structure and strength of the repulsion.

2.8 TRANSITIVITY OF REPULSION. One last property of DPPs that we shall mention is the fact that the negative correlations of the DPP possess a transient property in the sense, that if i and j and j and k are similar, then i and k are also similar. This is due to the fact

$$\|\phi_i - \phi_j\|^2 = \|\phi_i\|^2 + \|\phi_j\|^2 - 2\phi_i^T \phi_j = 2(1 - \phi_i^T \phi_j)$$

and thus

$$\sqrt{1 - \phi_i^T \phi_k} = \frac{1}{2} \|\phi_i - \phi_k\| \leq \frac{1}{2} (\|\phi_i - \phi_j\| + \|\phi_j - \phi_k\|) = \sqrt{1 - \phi_i^T \phi_j} + \sqrt{1 - \phi_j^T \phi_k}.$$

reformulate that part!

2.9 COMPARISON TO OTHER POINT PROCESSES.

II.2 Variations of DPPs

In this section we will present some useful variations of determinantal point processes. They serve different purposes and we will shortly explain their individual benefits.

2.10 CONDITIONAL DPPs. A *conditional DPP* is a collection of DPPs indexed by $X \in \mathcal{X}$, where X is called the *input* of the conditional DPP. Thus for every $X \in \mathcal{X}$ we get a finite set $\mathcal{Y}(X)$ and a determinantal point process $\mathbb{P}(\cdot | X)$ on $\mathcal{Y}(X)$ which is given by the elementary kernel $L(X)$, i.e.

$$\mathbb{P}(A|X) \propto \det(L_A(X)) \quad \text{for all } A \subseteq \mathcal{Y}(X).$$

Further we denote the quality and diversity features of the conditional DPP by $q_i(X)$ and $\phi_i(X)$ respectively.

It is not immediately clear why one would want to model a family of DPPs as a conditional DPP rather than as separate DPPs. The reason for this is that one wants to estimate the kernels $L(X)$ for every $X \in \mathcal{X}$. However if we would do this naively we would need to observe each of the DPPs $\mathbb{P}(\cdot | X)$ individually which is often not possible. Thus one hopes to not only memorise the kernels $L(X)$ for every single input $X \in \mathcal{X}$ but rather to learn the mapping L that assigns every input X its elementary kernel $L(X)$. If one achieved this task, one would be able to simulate

and predict a DPP that one has not observed so far just by the knowledge about some DPPs that belong to the same conditional DPP. Of course this can only work if we assume some regularity or a certain structure of the function L which we will do in the third chapter where we put those consideration into a precise framework.

2.11 FIXED SIZE OR k -DPPs.

2.12 STRUCTURED DPPs. We call a DPP *structured DPP* or short sDPP if the ground set is the cartesian product of some other set \mathcal{M} , which we will call the *set of parts*, i.e. if we have

$$\mathcal{Y} = \mathcal{M}^R = \{y_i = (y_i^r)_{r=1,\dots,R} \mid i = 1, \dots, N\}$$

where R is a natural number, $M = |\mathcal{M}|$ and $N = M^R$. The quality diversity decomposition of L take the form

$$L_{ij} = q(y_i)\phi(y_i)^T \phi(y_j)q(y_j)$$

and since $N = M^R$ is typically very big, it is impractical to define or store the quality and diversity features for every item $y_i \in \mathcal{Y}$. To deal with this problem we will assume that they admit factorisations and are thus a combination of only a few qualities and diversities.

More precisely we call $F \subseteq 2^{\{1,\dots,R\}}$ a *set of factorisations* and for a *factor* $\alpha \in F$, y_α denotes the subtupel of $y \in \mathcal{Y}$ that is indexed by α . Further we will work with the decompositions

$$\begin{aligned} q(y) &= \prod_{\alpha \in F} q_\alpha(y_\alpha) \\ \phi(y) &= \sum_{\alpha \in F} \phi_\alpha(y_\alpha) \end{aligned} \tag{2.6}$$

for a suitable set of factorisations F and qualities and diversities q_α and ϕ_α for $\alpha \in F$. Note that so far this is neither a restriction of generality – we could simply choose $F = \{\{1, \dots, R\}\}$ – nor a simplification – in that case we have the exact same number of qualities and diversities. However we are interested in the case where F consists only of small subsets of $\{1, \dots, R\}$. For example suppose that F is the set of all subsets with one or two elements, then we only have

$$R \cdot M + \binom{R}{2} \cdot M^2 = O(R^2 M^2)$$

quality and diversity features instead of

$$M^R = O(M^R).$$

This reduction of variables will make modelling, storing and estimating them feasible again in a lot of cases where naive approaches are foredoomed because of their shear size.

II.3 Simulation and Existence of DPPs

One of the main difficulties that arrises in the theory of discrete point processes is that they are probability measures on an exponentially large set, namely the powerset $2^{\mathcal{Y}}$ which has cardinality 2^N . Determinantal point processes have the benefit that they describe this distribution through the matrix K which consists of only N^2 parameters. This reduction of the number of parameters plays a central role in making a lot of operations possible in an computationally efficient way.

However it is not only the relatively small amount of parameters that lead to this, but also the structure of the determinant itself that leads to closed expressions for a lot of quantities like the normalisation constant in (2.4). In this section we will focus on the efficient simulation of DPPs and give a short overview of further techniques that can improve the performance of this algorithm.

But before we can do this we will present a famous identity for integrals – or sums – of the product of determinants.

TWO CAUCHY-BINET TYPE IDENTITIES

The result we present is a little technical, but we will need it later in the proof of existence and also in the proof of the sampling algorithm, so we will give the proof of it here, although it should be mentioned that all the major ideas can be found in [Hough et al., 2006]. In this section write $[n]$ for the set $\{1, \dots, n\}$ where n is a natural number and A_{IJ} for the submatrix of A where the first index is in I and the second one in J . Further we keep the notation $A_I = A_{II}$.

2.13 PROPOSITION (CAUCHY-BINET). *Let $m, n \in \mathbb{N}, m \leq n$ be two natural numbers and $A \in \mathbb{R}^{m \times n}, B \in \mathbb{R}^{n \times m}$ two matrices. Then we have*

$$\det(AB) = \sum_{\substack{I \subseteq \{1, \dots, n\} \\ |I|=m}} \det(A_{[m]I}) \det(B_{I[m]}).$$

Proof. Let without loss of generality A and B have full rank, otherwise both sides are zero. First we note that both sides are multilinear in the rows of A and columns of B . Hence we can assume by Gaussian row, respectively column elimination that each row of A and each column of B has exactly one non zero entry which is 1. Hence there is $J_1 = \{i_1, \dots, i_m\}$ and $J_2 = \{j_1, \dots, j_m\}$ such that $A_{ki_k} = 1$ and $B_{kj_k} = 1$ and all other entries are empty. Note that the right hand side is only non trivial if $I = J_1$ and $I = J_2$. If $J_1 \neq J_2$ then this case does not occur in the sum, but then at least one row of AB consists of zeroes and hence also the left hand side is equal to 0. If however $J = J_1 = J_2$ we get

$$\det(A_{[m]J}) \det(B_{J[m]}) = \det(A_{[m]J} B_{J[m]}) = \det(AB)$$

since $A_{ij} = B_{ji} = 0$ whenever $j \notin J$. □

2.14 PROPOSITION (VARIATION OF CAUCHY-BINET). *Let $m \leq n$ be two natural numbers and $B \in \mathbb{R}^{n \times m}$ be a matrix such that the columns $B_i \in \mathbb{R}^n$ of B form an orthonormal system. Further let $I \subseteq [m]$, then we have*

$$\det((B^T B)_I) = \sum_{\substack{I \subseteq J \subseteq [n] \\ |J|=m}} \det(B_{J[m]})^2.$$

Proof. This can be proved in analogue fashion to the result above. □

SAMPLING AND EXISTENCE

We roughly follow the approaches taken in [Hough et al., 2006] and [Kulesza et al., 2012] and will start by showing that every determinantal point process can be seen as a mixture of a smaller class of determinantal point processes.

really?

check this, I think the statement is slightly wrong...

2.15 THEOREM (MIXTURE REPRESENTATION OF DPPs). *Let \mathbb{P} be a DPP and*

$$K = \sum_{k=1}^N \lambda_k v_k v_k^T$$

be the spectral decomposition of its marginal kernel. Let now $\{B_k\}_{k=1,\dots,N}$ be a collection of independent Bernoulli random variables with mean λ_i . Define now the random kernel

$$K_B = \sum_{k=1}^N B_k v_k v_k^T. \quad (2.7)$$

Finally define a second point process $\tilde{\mathbb{P}}$ on \mathcal{Y} that is obtain by first drawing the Bernoulli variables B_k and then a DPP according to K_B . Then we have $\tilde{\mathbb{P}} = \mathbb{P}$ and thus $\tilde{\mathbb{P}}$ is also a DPP with marginal kernel K .

We will postpone the proof and first discuss its consequences which will be the existence of DPPs for a given marginal kernel as well as the construction of a sampling algorithm.

2.16 REMARK. Since it is fairly easy to simulate Bernoulli experiments, it remains to know how we can sample from DPPs with marginal kernels of the form $K = \sum_{k=1}^m v_k v_k^T$ for some $m \leq N$. We call DPPs of this type *elementary* and note that this corresponds to the class of DPPs where the eigenvalues of the marginal kernel are contained in $\{0, 1\}$.

Now we study the existence and simulation of elementary DPPs first and will be able to generalise those results to general DPPs without much effort.

2.17 PROPOSITION (EXISTENCE OF ELEMENTARY DPPs). *Let $K = \sum_{k=1}^m v_k v_k^T$ for some orthonormal set $V = \{v_k\}_{k=1,\dots,m} \subseteq \mathbb{R}^{\mathcal{Y}}$. Further define the measure on $2^{\mathcal{Y}}$ through*

$$\mathbb{P}(A) := \begin{cases} \det(K_A) & \text{if } |A| = m \\ 0 & \text{else} \end{cases}. \quad (2.8)$$

Then \mathbb{P} is a DPP on \mathcal{Y} with marginal kernel K . In particular elementary DPPs exist.

Proof. First we have to show that (2.8) defines a probability measure. For this let $B \in \mathbb{R}^{m \times N}$ be the matrix with rows v_k for $k = 1, \dots, m$. By definition we have $K = B^T B$ and hence

$$\begin{aligned} \sum_{A \subseteq \mathcal{Y}, |A|=m} \det(K_A) &= \sum_{\substack{I \subseteq \mathcal{Y} \\ |I|=m}} \det(K_I) = \sum_{\substack{I \subseteq \mathcal{Y} \\ |I|=m}} \det(B_{[m]I})^2 \\ &= \det(B^T B) = \det(v_k^T v_l)_{1 \leq k, l \leq m} = 1 \end{aligned}$$

where we have used the Cauchy-Binet identity and the fact that V is orthonormal. It remains to check that all marginal probabilities satisfy

$$\mathbb{P}(A \subseteq \mathbf{Y}) = \det(K_A).$$

For $|A| \geq m$ this follows immediately, so let $A = \{i_1, \dots, i_r\}$ for $r < m$. Then we obtain the marginal probability of A through integration over the other $m - r$ points. Namely we have

$$\mathbb{P}(A \subseteq \mathbf{Y}) = \sum_{\substack{A \subseteq J \subseteq [n] \\ |J|=m}} \mathbb{P}(J \subseteq \mathbf{Y}) = \sum_{\substack{A \subseteq J \subseteq [n] \\ |J|=m}} \det(B_{[m]J})^2 = \det((B^T B)_A) = \det(K_A)$$

where we used Proposition 2.14. □

Now we can turn towards the simulation of elementary DPPs where we will make use of the previous result.

Algorithm 1 Sampling from an elementary DPP

Input: Marginal kernel $K = \sum_{k=1}^m v_k v_k^T$ for $\{v_k\}_{k=1,\dots,m}$ orthonormal

```

1:  $V \leftarrow \{v_k\}_{k=1,\dots,m}$ 
2:  $Y \leftarrow \emptyset$ 
3: while  $|V| > 0$  do
4:    $p_i \leftarrow P e_i$  the projection of  $e_i$  onto  $\text{span}(V)$  for  $i \in \mathcal{Y}$ 
5:   Select  $i \in \mathcal{Y}$  with probability  $\frac{1}{|V|} \cdot \|p_i\|^2$ 
6:    $Y \leftarrow Y \cup \{i\}$ 
7:    $V \leftarrow V_\perp$  an orthonormal basis of the subspace of  $V$  perpendicular to  $p_i$ 
8: end while
9: return  $Y$ 

```

2.18 PROPOSITION (SAMPLING FROM ELEMENTARY DPPs). *Let $K = \sum_{k=1}^m v_k v_k^T$, then Algorithm 1 produces a random variable \mathbf{Y} with values in \mathcal{Y} which is an elementary DPP with marginal kernel K .*

reread!

Proof. We note that we only have to check that (2.8) holds and for this we fix $A \subseteq \mathcal{Y}$. First we note that the output \mathbf{Y} has cardinality $|m|$ since no element can be selected twice in the while loop and the size of V decreases by exactly one in each iteration. Hence it remains to show

$$\mathbb{P}(A = \mathbf{Y}) = \det(K_A)$$

if $|A| = m$. Let for the sake of convenience $A = \{1, \dots, m\}$ and $\mathcal{Y} = \{1, \dots, N\}$. Note that it suffices to show that the while loop selects $1, \dots, m$ in this exact order with probability $\frac{1}{m!} \det(K_A)$.

Let V_k denote the orthonormal set V in the k -th step of the while loop and let P_{k-1} be the projection onto $\text{span}(V_k)$ and set $b_i := P_0 e_i$ for $i = 1, \dots, N$. We note that if $1, \dots, k-1$ were selected in the first steps, then P_{k-1} is exactly the projection to the subspace of $\text{span}(V_{k-1})$ that is orthogonal to b_1, \dots, b_{k-1} . Since the spaces $\text{span}(V_k)$ are decreasing we have $P_k P_j = P_k$ for $k \geq j$ and thus $P_{k-1} e_k = P_{k-1} P_0 e_k = P_{k-1} b_k$. Suppose now that we have selected $1, \dots, k-1$ in the first $k-1$ steps of the while loop. The probability to select k in the next iteration is

$$\frac{1}{|V_k|} \cdot \|P_{k-1} e_k\|^2 = \frac{1}{m-k} \cdot \|P_{k-1} b_k\|^2.$$

Thus the probability to sample $1, \dots, m$ in this order is equal to

$$\frac{1}{m!} \cdot \|b_1\|^2 \cdot \dots \cdot \|P_{m-1} b_m\|^2.$$

Since P_{k-1} is the projection onto the subspace orthogonal to b_1, \dots, b_k , the product is equal to the squared m -dimensional surface measure of the parallel epiped spanned by b_1, \dots, b_m . It is well known from measure and integration theory that the squared surface is given by the determinant of the Gram matrix

$$\det \begin{pmatrix} b_1^T b_1 & \dots & b_1^T b_m \\ \vdots & \ddots & \vdots \\ b_m^T b_1 & \dots & b_m^T b_m \end{pmatrix} = \det((B^T B)_A)$$

where $B \in \mathbb{R}^{N \times N}$ is the matrix which rows are equal to b_k . Therefore it remains to show $B^T B = K^V$. However by definition B is the projection onto the span of V and thus $B = K^V$. Because K^V is symmetric like every projection, we have $B^T = B$ and hence can conclude $B^T B = B^2 = B = K^V$ where we used that B is a projection. \square

We will use Theorem 2.15 to prove that the following algorithm samples from a DPP. This will also show the existence of DPPs to a given marginal kernel since it gives an explicit construction.

Algorithm 2 Sampling from a DPP

Input: Eigendecomposition $\{v_k, \lambda_k\}_{k=1, \dots, N}$ of K

```

1:  $J \leftarrow \emptyset$ 
2: for  $k = 1, \dots, N$  do
3:    $J \leftarrow J \cup \{k\}$  with probability  $\lambda_k$ 
4: end for
5:  $V \leftarrow \{v_k\}_{k \in J}$ 
6:  $Y \leftarrow \emptyset$ 
7: while  $|V| > 0$  do
8:    $p_i \leftarrow P e_i$  the projection of  $e_i$  onto  $\text{span}(V)$  for  $i \in \mathcal{Y}$ 
9:   Select  $i \in \mathcal{Y}$  with probability  $\frac{1}{|\mathcal{Y}|} \cdot \|p_i\|^2$ 
10:   $Y \leftarrow Y \cup \{i\}$ 
11:   $V \leftarrow V_{\perp}$  an orthonormal basis of the subspace of  $V$  perpendicular to  $p_i$ 
12: end while
13: return  $Y$ 

```

2.19 THEOREM (SAMPLING ALGORITHM). Let $K \in \mathbb{R}^{N \times N}$ be any symmetric and positive definite matrix such that $K \leq I$. Then the distribution of the output Y of Algorithm 2 is a DPP with marginal kernel K .

Proof. Theorem 2.15 states that an arbitrary DPP is the mixture of elementary DPPs and the for loop in the algorithm represents exactly this mixing with the respective weights. Further the sampling result for elementary DPPs yields that the output of the second part of the algorithm, namely the while loop, is distributed according to a DPP with marginal kernel $K^V := \sum_{v \in V} v v^T$. \square

2.20 COROLLARY (EXISTENCE OF DPPs). Let K be a symmetric $N \times N$ matrix. Then K is the marginal kernel of a DPP if and only if $0 \leq K \leq I$.

2.21 COROLLARY (CARDINALITY OF DPPs). Let \mathbb{P} be a DPP with kernel

$$K = \sum_{k=1}^N \lambda_k v_k v_k^T.$$

Then the cardinality of the DPP is distributed like the sum of the Bernoulli variables $\{B_k\}_{k=1, \dots, N}$ from theorem 2.15.

comment on the intuition one can get from this!

Proof. To proof this, we only have to convince ourselves that after the Bernoulli experiments the cardinality of a DPP with kernel (2.7) has size $m := \sum_{k=1}^N B_k$ almost surely. Since K_B has rank at most k , the cardinality is almost surely smaller than m . On the other hand we have

$$\mathbb{E}[|\mathbf{Y}|] = \sum_{i \in \mathcal{Y}} \mathbb{P}(i \in \mathbf{Y}) = \sum_{i \in \mathcal{Y}} (K_B)_{ii} = \text{Tr}(K_B) = m. \quad (2.9)$$

In the last step we used that the trace of a symmetric matrix is the sum over its eigenvalues, which are B_k in our case. This computation lets us conclude $|\mathbf{Y}| = m$ almost surely. \square

We close this section with the proof of 2.15 given in [Kulesza et al., 2012].

reread!

Proof of Theorem 2.15. Let $A \subseteq \mathcal{Y}$, $k := |A|$. Further set $W_n := (v_n v_n^T)_A$ and $W_J := \sum_{n \in J} W_n$. Then we have

$$\tilde{\mathbb{P}}(A \subseteq \mathbf{Y}) = \sum_{J \subseteq \mathcal{Y}} \det(W_J) \cdot \tilde{\mathbb{P}}(B_i = 1 \text{ for } i \in J).$$

Let $((W_{n_1})_1 (W_{n_2})_2 \dots (W_{n_k})_k)$ denote the $k \times k$ matrix with i -th row equal to the i -th row of W_{n_i} . Using the multilinearity of the determinant we obtain that the marginal probability above is equal to

$$\begin{aligned} & \sum_{J \subseteq \mathcal{Y}} \sum_{n_1, \dots, n_k \in J} \det((W_{n_1})_1 (W_{n_2})_2 \dots (W_{n_k})_k) \cdot \tilde{\mathbb{P}}(B_i = 1 \text{ for } i \in J) \\ &= \sum_{n_1, \dots, n_k \in \mathcal{Y}} \det((W_{n_1})_1 (W_{n_2})_2 \dots (W_{n_k})_k) \sum_{J \supseteq \{n_1, \dots, n_k\}} \tilde{\mathbb{P}}(B_i = 1 \text{ for } i \in J) \\ &= \sum_{n_1, \dots, n_k \in \mathcal{Y}} \det((W_{n_1})_1 (W_{n_2})_2 \dots (W_{n_k})_k) \cdot \tilde{\mathbb{P}}(B_{n_i} = 1 \text{ for } i = 1, \dots, k) \\ &= \sum_{n_1, \dots, n_k \in \mathcal{Y}} \det((\lambda_{n_1} W_{n_1})_1 (\lambda_{n_2} W_{n_2})_2 \dots (\lambda_{n_k} W_{n_k})_k) \\ &= \det\left(\sum_{n \in \mathcal{Y}} W_n\right) = \det(K_A). \end{aligned}$$

This computation shows that $\tilde{\mathbb{P}}$ is a DPP with marginal kernel K . \square

Possible improvements

2.22 DUAL SAMPLING.

2.23 DIMENSION REDUCTION.

II.4 The mode problem

One general motivation for modelling is the hope that predictions can be made from the selected model. If the model is of stochastic nature, like in our case, and if one wants to predict its outcome, there are a few possible approaches. The first one would be to sample from this model. This relies on the intuition that a realisation of our random variable will be a rather typical example for the random event. Going one step further one could try to find the most likely outcome of the random variable, which is known as the mode problem.

2.24 THE MODE PROBLEM. Let X be a random variable with values in some space \mathcal{X} and let f be the density of the distribution of X with respect to some reference measure. Then the *mode* is the maximiser

$$\hat{x} = \arg \max_{x \in \mathcal{X}} f(x)$$

of the density if it exists. The search for the mode is called the *mode problem*.

Our motivation for finding the mode of a random variable was to make better predictions for it. This is justified by the assumption that the mode should be a typical realisation of the random variable. However this is not generally the case and therefore one should be cautious with this intuition. Consider for example the mixture of two independent Gaussian random variables

look for better word

$$0.1 \cdot X + 0.9 \cdot Y$$

where X is centered with variance 10 and Y has mean 5 with variance 1, the densities are shown in Figure It is clear that mode is 0 in this example, but it is not a very typical outcome of the random variable, since the majority of events is centered around 10.

check whether this gives the desired effect and plot the density!

The mode problem is rather well behaved if the density f is a smooth function defined on a subset of \mathbb{R}^d , but in the case of DPPs we have to deal with the probability measure on a finite set. Thus this turns into a discrete optimisation problem over the exponentially large powerset $2^{\mathcal{Y}}$. This is in general very hard to solve and it has been shown in that it is NP hard to do so or even approximate it upto a factor of $\frac{8}{9}$. However there were still different strategies proposed and we will present some of them including their main ideas.

cite

do this!

Chapter III

Point estimators and parametric models

Parameter estimation is one of the central components of every theory of real world phenomena. In a nutshell one could split the process of the construction of a descriptive model into two parts. The first one being the selection of the model which is done by a scientist and the second being the determination of the constants that belong to the model.

To make this more clear we will consider one of the most famous advances in the natural sciences namely the law of universal gravitation that was discovered by Sir Isaac Newton and published in one of the most famous books in the history of science, the *Philosophiæ Naturalis Principia Mathematica*. More precisely Newton discovered that the gravitational force acting between two massive objects is given by

$$F = G \cdot \frac{m_1 m_2}{r^2}$$

where m_1, m_2 are the masses of the two objects, r is the distance of the centers of masses and G is the gravitational constant. This constant can not be deduced from the theory itself and needs to be estimated based on some empirical data.

If we want to describe, simulate and predict the occurrence of diverse subsets we can take a similar approach and impose the model of a determinantal point process. This will usually be an assumption that will not strictly hold, but will often lead to reasonable, sometimes even impressive results. We will not be concerned to measure how suitable this model selection is, although this is a highly interesting question. Leaving that aside we are left with the second step, namely the estimation of the parameters of the model, which are in the case of a DPP over a set of cardinality N exactly $N(N - 1)/2$. Because of the rather large amount of parameters and also the complicated structure of the DPPs it will in practice only be possible to perform those estimations through the use of computational tools. The task of computer based parameter or density estimation is an important field in the discipline of *machine learning* and thus we will sometimes speak of the parameters being learned instead of estimated. Actually the interest of parameter estimation for DPPs arose from the machine learning community at the beginning of this decade. However we will phrase things in a way that no prior knowledge in this field is required.

In this chapter we will be concerned in how we can make point estimates for either the marginal or the elementary kernels K and L . Point estimators are the most basic type of estimators and consist of the suggestion of one possible parameter set, for example in the case of the gravitational constant

$$6.674 \cdot 10^{-11} \text{N kg}^{-2} \text{m}^2.$$

This is in contrast to the Bayesian approach to parameter estimation that we will present in the next chapter where the philosophy is to estimate a distribution over all possible parameter sets that indicates how likely they are given some the empirical data. We will discuss two essentially different methods of point estimators, the first one provides a way to reconstruct a marginal kernel for the empirical marginal distributions at least in the case where the empirical distribution is essentially a DPP. The other type of methods are all maximum likelihood estimators in different variations.

But before we can proceed we want to remind the reader of two desirable properties of point estimators. For this we will assume that we want to estimate the distribution of a random variable X from a parametric family of probability measures

$$\{\mathbb{P}_\theta \mid \theta \in \Theta\}.$$

This means we want to estimate θ out of a possible set of parameters Θ such that X is distributed according to \mathbb{P}_θ which we will based upon some data x_1, \dots, x_n . Further we assume that those points are actually generated by \mathbb{P}_θ for one $\theta \in \Theta$ and denote the estimator by $\hat{\theta}_n$. We call *unbiased* if we have

$$\mathbb{E}[\hat{\theta}_n] = \theta$$

and *consistent* if we have

$$\hat{\theta}_n \rightarrow \theta \quad \text{in probability.}$$

It shall be noted that although those properties are beneficial, they are not crucial for an estimator to be reasonable. First they both assume that the data generating process, i.e. the process one wants to describe actually follows one of the laws \mathbb{P}_θ which will typically be not the case in real world examples. Further the asymptotic property of consistency is rather of theoretical nature since in practice it is not possible to create large sets of empirical data and certainly not infinitely large ones.

III.1 Kernel reconstruction from the empirical measures

Now we will display the first way how one can estimate the marginal kernel K of a DPP based on some samples drawn from it.

3.1 SETTING. Let \mathcal{Y} be a finite set of cardinality N and let $K \in \mathbb{R}_{\text{sym}}^{N \times N}$ satisfy $0 \leq K \leq I$. Let further $\mathbf{Y}_1, \mathbf{Y}_2, \dots$ be distributed according to the DPP with marginal kernel K .

In order to perform an approximate reconstruction of the marginal kernel we will need to consider the *empirical measure*

$$\hat{\mathbb{P}}_n := \frac{1}{n} \sum_{i=1}^n \delta_{\mathbf{Y}_i}.$$

The interest in $\hat{\mathbb{P}}_n$ lies in the fact that they quite natural estimates for the actual underlying distribution. More precisely they are unbiased estimators for \mathbb{P} , i.e. they agree in expectation with \mathbb{P} . This can be seen by evaluating it at $A \subseteq \mathcal{Y}$

$$\mathbb{E}_{\mathbb{P}}[\hat{\mathbb{P}}_n(A)] = \frac{1}{n} \sum_{i=1}^n \mathbb{E}_{\mathbb{P}}[\delta_{\mathbf{Y}_i}(A)] = \mathbb{P}(A).$$

And even stronger by the strong law of large numbers they converge to \mathbb{P} almost surely if the sequence $(\mathbf{Y}_k)_{k \in \mathbb{N}}$ of observations is independent. This can be seen by identifying the probability measures on $2^{\mathcal{Y}}$ with the probability simplex

$$\left\{ \mu \in \mathbb{R}^{2^{\mathcal{Y}}} \mid \mu_A \in [0, 1] \text{ for all } A \subseteq \mathcal{Y} \text{ and } \sum_{A \subseteq \mathcal{Y}} \mu_A = 1 \right\}$$

and using the strong law of large numbers in $\mathbb{R}^{2^{\mathcal{Y}}}$.

Therefore the empirical measures are reasonable approximations of the actual probability distribution. Assume now for one moment that the empirical measures $\hat{\mathbb{P}}_n$ are also determinantal point processes with marginal kernel \hat{K}_n , then \hat{K}_n would be a quite intuitive estimate for the actual marginal kernel K . Thus we are interested in the question whether we can reconstruct the kernel marginal of a DPP if we know the DPP itself. Since the marginal density of a DPP corresponds to the principal minors of the marginal kernel, we first investigate whether we can reconstruct a matrix from its principal minors. For the answer to this problem we follow the main ideas presented in [Urschel et al., 2017] and [Rising et al., 2015] although we modify their arguments to make them shorter and hopefully more accessible.

3.2 THE PRINCIPAL MINOR ASSIGNMENT PROBLEM. Let $K \in \mathbb{R}^{N \times N}$ be a symmetric matrix. We want to investigate whether K uniquely specified by its principal minors

$$\Delta_S := \det(K_S) \quad \text{where } S \subseteq \{1, \dots, N\}.$$

We call this the *symmetric principal minor assignment problem* and it will turn out that the matrix K can be reconstructed up to an equivalence relation.

Before we present the general procedure we want to see how this would work in the case of a symmetric 3×3 matrix $K = (K_{ij})_{1 \leq i, j \leq 3}$. First we note that we can regain the diagonal elements as the determinant of the 1×1 principal minors

$$\det(K_{\{i\}}) = K_{ii} \quad \text{for } i = 1, 2, 3.$$

Further the squares of the off diagonal are determined by the 2×2 principal minors since

$$\det(K_{\{i, j\}}) = K_{ii} K_{jj} - K_{ij}^2 \quad \text{for } i, j = 1, 2, 3.$$

Therefore we only need to reconstruct the signs off diagonal entries. To do this, we consider the determinant of the matrix itself

$$\det(K) = K_{11} K_{22} K_{33} + 2K_{12} K_{13} K_{23} - K_{11} K_{23}^2 - K_{22} K_{13}^2 - K_{33} K_{12}^2. \quad (3.1)$$

Rearranging this yields

$$K_{12} K_{13} K_{23} = \frac{1}{2} \left(\det(K) + K_{11} K_{23}^2 + K_{22} K_{13}^2 + K_{33} K_{12}^2 - K_{11} K_{22} K_{33} \right).$$

Since we know all of the expressions on the right side, we can determine the sign of the product on the left side. Now we assign the signs of the off diagonal elements in such a way, that the above equation holds. More precisely if the product is negative, we assign a minus to one or all three elements, if it is positive, then we assign a minus to none or two elements. If the product is zero, every configuration of signs satisfy the desired property. It is now straight forward to check that this assignment actually leads to the desired principal minors.

III.1.1 Graph theoretical concepts

One main part in the general procedure will be to obtain a generalisation of the formula (3.1) for larger principal minors that will allow the reconstruction of the signs. For this we will need the following graph theoretical concepts.

3.3 NOTIONS FROM GRAPH THEORY. Let $G = (V, E)$ be a finite graph, i.e. V is a finite set, called the *vertex set* and E consists of subsets of V with two elements, the *edges*. Sometimes we will be sloppy in notation and not distinguish between the graph and the edge set. We will need the following notions:

- (i) *Degree*: For a vertex $v \in V$ the *degree* is the number of edges that contains v .
- (ii) *Subgraph*: A graph $\tilde{G} = (\tilde{V}, \tilde{E})$ is called a *subgraph* of G if $\tilde{V} \subseteq V$ and $\tilde{E} \subseteq E$.
- (iii) *Induced graph*: For a subset $S \subseteq V$ of vertices the *induced graph* $G(S) = (S, E(S))$ is formed of all edges $e \in E$ of G that are subsets of S .
- (iv) *Path*: A *path* in G is a sequence $v_0 v_1 \cdots v_k$ of vertices such that $\{v_{i-1}, v_i\} \in E$ for all $i = 1, \dots, k$.
- (v) *Connected graph*: A graph is called *connected* if for every pair of vertices $v, w \in V$ there is a path from v to w .
- (vi) *Cycle*: A *cycle* C is a connected subgraph such that every vertex has even degree in C .
- (vii) *Cycle space*: Each cycle C can be identified with a vector $x = x(C) \in \mathbb{F}_2^E$ such that

$$x_e := \begin{cases} 1 & \text{if } e \in C \\ 0 & \text{if } e \notin C \end{cases}$$

indicates whether the edge $e \in E$ belongs to the cycle C . The *cycle space* \mathcal{C} is the span of $\{x(C) \mid C \text{ is a cycle}\}$ in \mathbb{F}_2^E . Note that the sum of two cycles in the cycle space corresponds to the symmetric difference of the edges.

- (viii) *Chordless cycle*: A cycle C is called *chordless* if two vertices $v, w \in C$ form an edge in G if and only they form an edge in C . This is equivalent to the statement that C is an induced subgraph that is a cycle.
- (ix) *Cycle sparsity*: The cycle sparsity is the minimal number l such that a basis of the cycle space consisting of chordless simple cycles exists. Such a basis is called *shortest maximal cycle basis* or short *SMCB*. If the cycle space is trivial we define the cycle sparsity to be 2.
- (x) *Pairings*: Let $S \subseteq V$ be a set of vertices. Then a *pairing* P of S is a subset of edges of $G(S)$ such that two different edges of P are disjoint. The vertices contained in the edges of P are denoted by $V(P)$ and the set of all pairings by $\mathcal{P}(S)$.

It is highly recommended to study the examples in Figure III.1.1 in order to get more familiar with the definitions above. To see that the above definition of the cycle sparsity is well defined, we have need to show that shortest maximal cycle basis exist. This might be well known to people that are familiar with graph theory, but we will present an elementary proof here. The first part of the statement, namely the existence of cycle basis consisting of simple cycles is known as Veblen's theorem and can be found in its original form in [Veblen, 1912], however we will rather follow the approach in [Bondy and Murty, 2011].



Figure III.1.: Some examples of graphs and cycles. The first sketch shows a graph and the three other ones subgraphs of it where the edges not belonging to the subgraph are depicted dashed. The first one is a simple chordless cycle, the second one a simple but not chordless cycle and the last one is not a cycle at all.

3.4 PROPOSITION (EXISTENCE OF SMCBs). *There always exists a basis $\{x(C_1), \dots, x(C_k)\}$ of the cycle space where C_1, \dots, C_k are chordless simple cycles.*

Proof. First we prove that the set of simple cycles generates the whole cycle space which we can then improve to show that the simple chordless cycles already generate the cycle space. A shortest maximal cycle basis is then attained by successively dropping simple chordless cycles.

We show that every cycle $x(C)$ can be written as the sum of simple cycles $x(C_1), \dots, x(C_k)$ where $C_i \subseteq C$. This is equivalent to the statement that the edges of every cycle are the disjoint union of the edges of simple cycles. Take now a maximal non intersecting path $v_0 v_1 \dots v_k$. Since v_k has degree at least 2, there is an edge $\{v_k, v_{k+1}\}$ such that $v_{k+1} \neq v_{k-1}$. Since the

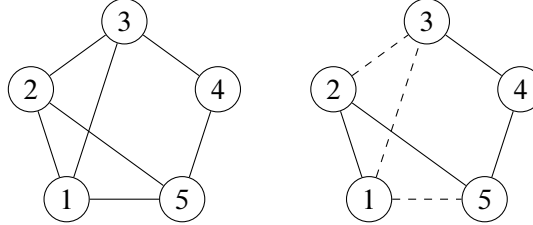


Figure III.2.: Illustration of the search for a simple cycle in a graph with degrees greater than two. Once a maximal non intersecting path like 12543 is selected, every continuation of the path – in this case 2 or 1 – is already present in the path and therefore induces a simple cycle.

path is maximal, v_{k+1} has to agree with one a vertex $v_i \in \{v_0, \dots, v_{k-2}\}$, because otherwise we could add v_{k+1} to the path which is a contradiction to the maximality. Now $v_i v_{i+1} \dots v_k v_i$ corresponds to a simple cycle C_1 and $C_2 := C \setminus C_1$ is again a cycle. Thus we can write C as the disjoint union $C = C_1 \cup C_2$ where C_1 is a simple cycle. By repeating this procedure we get the desired expression for C in terms of simple cycles.

To prove that already the simple chordless cycles generate the cycle space we have to prove that we can write every simple cycle $x(C)$ as a sum of simple chordless cycles $x(C_1), \dots, x(C_k)$. Let $\{\{v_0, v_1\}, \dots, \{v_k, v_0\}\}$ be the edge set of C and assume that C is not chordless like in Figure III.1.1, otherwise the statement would be trivial. Thus there is are indices $1 \leq i < j-1 \leq k-1$ such that $\{v_i, v_j\} \in E$. Let now C_1 and C_2 be the two cycles associated with the paths

$$v_0 v_1 \dots v_i v_j v_{j+1} \dots v_k v_0 \quad \text{and} \quad v_i v_{i+1} \dots v_{j-1} v_j v_i.$$



Figure III.3.: The simple cycle 123451 on the left is not chordless but the symmetric difference of the two simple chordless cycles 1231 and 13451 on the right.

Then we have $x(C) = x(C_1) + x(C_2)$. By iterating this procedure as long as the cycles are not chordless the desired decomposition can be achieved in finitely many steps. \square

III.1.2 The solution of the principal minor assignment problem

Now we have all the graph theoretical prerequisites to show how one can reconstruct a matrix with preassigned principal minors. However the matrix that arises from this reconstruction is not unique and thus we need to identify matrices with the same principal minors with each other.

3.5 DEFINITION (DETERMINANTAL EQUIVALENCE). Two symmetric matrices $A, B \in \mathbb{R}^{N \times N}$ are called *determinantally equivalent* if they have the same principal minors and we write $A \sim B$.

It is obvious that we can only hope to reconstruct a symmetric matrix up to determinantal equivalence. However this would be satisfactory, because determinantally equivalent matrices are exactly those that give rise to the same DPP. Let us in the following denote the principal minor $\det(K_S)$ by Δ_S for $S \subseteq \{1, \dots, N\}$. To come back to our original problem, we notice that the principal minors up to size two immediately determine the diagonal and the absolute values of the off diagonal of K since we have

$$K_{ii} = \Delta_{\{i\}} \quad \text{and} \quad K_{ij}^2 = K_{ii}K_{jj} - \Delta_{\{i,j\}}.$$

Thus it only remains to regain the signs $\text{sgn}(K_{ij})$ of the off diagonal entries. For this we use the following object.

3.6 THE ADJACENCY GRAPH AND SIGN FUNCTION. The adjacency graph $G_K = (V_K, E_K)$ associated with K consists of the vertex set $\{1, \dots, N\}$ and $\{i, j\}$ form an edge if and only if $K_{ij} \neq 0$. Further we introduce some *weights* on the edges. This means we consider a mapping $w: E_K \rightarrow \mathbb{R}$ and we set

$$w_{ij} := w(\{i, j\}) := \text{sgn}(K_{ij})$$

where we call w_{ij} the weight of the edge $\{i, j\}$. This graph together with the weights determines the signs of the off diagonal elements, so we are interested in reconstructing the weights from the principal minors. Finally we define the sign of a cycle and for a cycle $C = (S, \tilde{E})$ we set $\text{sgn}(C) := \prod_{e \in \tilde{E}} w_e$. It will become important later to consider this sign function on the cycle space and thus we note that this definition corresponds to

$$\text{sgn}(x(C)) := \prod_{e \in E} w_e^{x(C)_e}.$$

Note that this is a group homomorphism from the cycle space \mathcal{C} to $\{\pm 1\}$ and therefore it is uniquely determined by its value on a generator, for example on a shortest maximal cycle basis.

3.7 PROPOSITION (PRINCIPAL MINORS OF SIMPLE CHORDLESS CYCLES). *Let $C = (S, E(S))$ be a simple and chordless cycle. Then the principal minor of K with respect to S is given by*

$$\Delta_S = \sum_{P \in \mathcal{P}(S)} (-1)^{|P|} \cdot \prod_{\{i,j\} \in P} K_{ij}^2 \cdot \prod_{i \notin V(P)} K_{ii} + 2 \cdot (-1)^{|S|+1} \cdot \prod_{\{i,j\} \in E(S)} K_{ij}. \quad (3.2)$$

Proof. Let $k := |S|$. Then by Leibniz formula we have

$$\Delta_S = \sum_{\sigma \in S_k} \text{sgn}(\sigma) \prod_{i \in S} K_{i\sigma(i)}$$

where S_k is the set of permutations of S . Note that since the cycle is chordless, the product is only non trivial if $\{i, \sigma(i)\} \in E(S)$ for all $i \in S$. Since C is a simple cycle, those permutations consist exactly of the pairing of S or the two shifts of the set S along the cycle in both directions. Those correspond exactly to the summands in (3.2).

To see this, we fix a permutation σ such that $\{i, \sigma(i)\}$ always forms an edge in $(S, E(S))$. We note that every vertex $i \in S$ has two possible images which are exactly the endpoint of its two edges, c.f. Figure III.1.2. Lets assume it is mapped to $j \in S$, then j has again two possible

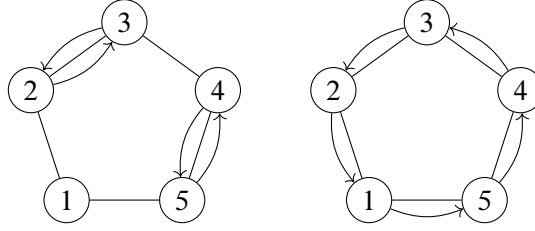


Figure III.4.: An easy example for the two kinds of permutations of a chordless simple cycle that maps vertices to neighbors.

images under σ namely i and a second vertex $k \in \mathcal{V}$. If $j \mapsto i$, no other vertex can be mapped to i or j , however some other items can be swapped in the same way. The permutations of this form correspond exactly to the pairings of S and are represented in the first sum in (3.2). If however j is not mapped back to i but rather to its other neighbor k , then k can't get mapped back to j since σ is injective. Thus it has to be mapped to its other neighbor $l \in \mathcal{V}$. Through a repetition of this argument shows that this induces until i is reached again. Since the cycle is simple this path exhausts the entire cycle. The factor 2 is due to the fact that this shift of the indices can be done into either direction. \square

3.8 PROPOSITION (SIGN DETERMINES PRINCIPALS MINORS). *The knowledge of all principal minors up to size two and the sign function*

$$\text{sgn}: \mathcal{C} \rightarrow \{\pm 1\}$$

completely determines all principal minors of K .

Proof. Let $S \subseteq \{1, \dots, N\}$ be arbitrary. We will again work with the expression (3.2) of the principal minor Δ_S and fix one permutation σ . We can assume without loss of generality that

$\{i, \sigma(i)\} \in E_K$ because the product is trivial otherwise. Since we know the absolute values of the off diagonal elements and the diagonal elements from the principal minors up to size two, it suffices to express the sign

$$\prod_{i \in S} \text{sgn}(K_{i\sigma(i)}) \quad (3.3)$$

of the product through the sign function. For this we write σ as the product of disjoint cycles

$$\sigma = \sigma_1 \circ \dots \circ \sigma_m \quad (3.4)$$

where $\sigma_k: D_k \rightarrow D_k$ for $k = 1, \dots, m$ and the domains D_k are pairwise disjoint. The sign (3.3) can be written as the product of

$$\prod_{i \in D_k} \text{sgn}(K_{i\sigma_k(i)})$$

so it suffices to give expressions for those. Note that we could assume $\{i, \sigma_k(i)\} \in E_K$ and therefore $C_k = (D_k, E_k)$ with

$$E_k = \{\{i, \sigma_k(i)\} \mid i \in D_k\}$$

is a cycle and therefore (3.4) is equal to $\text{sgn}(C_k)$. \square

3.9 THEOREM. *Let $K \in \mathbb{R}^{N \times N}$ be a symmetric matrix and l be the sparsity of its adjacency graph. Then the principal minors up to size l uniquely determine all principal minors of K and therefore the matrix K up to determinantal equivalence.*

Proof. In the light of the previous proposition it suffices to show that the sign function is uniquely specified by the principal minors up to size l . Recall that the sign function is determined by its values on a shortest maximal cycle basis, which consists by definition of simple chordless cycles of length at most l . However under the knowledge of the diagonal elements and the absolute values of the off diagonal ones, the sign of those simple chordless cycle is uniquely determined by the principal minors up to size l using the equality (3.2). \square

3.10 REMARK. One can even show that this result is optimal in the sense that if one only has access to the principal minors up to size $l - 1$, then the equivalence class is not uniquely determined. To see this, we note that the sign function is not uniquely specified through the principal minors up to size $l - 1$ and thus there is more than one extension of the sign function onto the shortest maximal cycle basis. The equation (3.2) shows that those different extensions give rise to different principal minors.

3.11 CONSTRUCTION OF THE EQUIVALENCE CLASS. We have shown that the determinantal equivalence class of a symmetric matrix is uniquely specified by its principal minors up to size l . Now we want to investigate how this equivalence class can be computed and we will see that we can reduce this task to the solution of a system of linear equations over the finite field \mathbb{F}_2 .

Let us assume that we have knowledge of the principal minors Δ_S for every $S \subseteq \{1, \dots, N\}$ with size at most l and we want to construct a matrix \tilde{K} that is determinantly equivalent to K . We have seen that we only need to reconstruct the signs of the off diagonal entries of K which is equivalent to reconstructing the edge weight w_{ij} . To do this fix a shortest maximal cycle basis $\{C_1, \dots, C_m\}$ with vertex sets S_1, \dots, S_m . Let us now rewrite (3.2) in the form

$$H_k := \Delta_{C_k} - \sum_{P \in \mathcal{P}(C_k)} (-1)^{|P|} \cdot \prod_{\{i,j\} \in P} K_{ij}^2 \cdot \prod_{i \notin V(P)} K_{ii} = 2 \cdot (-1)^{|C_k|+1} \text{sgn}(C_k) \cdot \prod_{\{i,j\} \in C_k} |K_{ij}|.$$

Given the principal minors, we can determine the value on the right side and taking the sign on both sides yields

$$(-1)^{|C_k|+1} \cdot \text{sgn}(H_k) = \text{sgn}(C_k) = \prod_{\{i,j\} \in E(S_k)} w_{ij}$$

which we seek to solve for w . However this multiplicative equation is hard to solve and therefore we use the canonical group isomorphism ϕ between $\{\pm 1\}$ and $\{0, 1\}$ to turn it into a linear equation. Setting $x_{ij} := \phi(w_{ij})$ we get that the condition above is equivalent to

$$b_k := \phi(\text{sgn}(H_k)) + |\hat{S}_k| + 1 = \sum_{\{i,j\} \in E(S_k)} x_{ij} = (Ax)_k \quad \text{in } \mathbb{F}_2$$

where A is the matrix with the rows $x(C_k)^T$. Now we can fix any such solution $x \in \mathbb{F}_2^E$ of

$$Ax = b \tag{3.5}$$

and we know that at least one exists, namely the one given by $x_{ij} = \phi(\text{sgn}(K_{ij}))$. Let now $w_{ij} := x_{ij}$, then it is straight forward to see that \tilde{K} defined through

$$\tilde{K}_{ii} := \Delta_{\{i\}} \quad \text{and} \quad \tilde{K}_{ij} = w_{ij} \cdot \sqrt{\tilde{K}_{ii} \tilde{K}_{jj} - \Delta_{\{i,j\}}}$$

is determinantal equivalent to K .

It shall be noted that there are algorithms with much better computational performance for the construction of the determinantal equivalence class. For some examples of efficient algorithms we refer to [Urschel et al., 2017] and [Rising et al., 2015].

III.1.3 Definition of the estimator and consistency

So far we have seen that the principal minors determine a symmetric matrix up to determinantal equivalence. However the empirical marginal densities do not in general need to be the principal minors of any symmetric matrix, in other words the empirical measures are not necessarily determinantal. Therefore the definition of the estimator is till not quite straight forward and we will follow [Urschel et al., 2017] for this and make the following assumption.

3.12 ASSUMPTION. Let $\alpha > 0$ and assume that

$$\min \{ |K_{ij}| \mid K_{ij} \neq 0 \} \geq \alpha.$$

Note that such an α can always be found, however it is not a priori known. For example if we want to make a statement about the speed of approximation of the estimators, which depends on α , we have to make the assumption above.

3.13 DEFINITION OF THE ESTIMATOR. The straight forward estimators of the principal minors are

$$\hat{\Delta}_S := \hat{\mathbb{P}}_n(S \subseteq \mathbf{Y}) \quad \text{for } S \subseteq \{1, \dots, N\}.$$

The resulting estimates for the diagonal elements and the squares of the off diagonals are

$$\hat{K}_{ii} := \hat{\Delta}_{\{i\}} \quad \text{and} \quad \hat{B}_{ij} := \hat{K}_{ii} \hat{K}_{jj} - \hat{\Delta}_{\{i,j\}}.$$

Next we will introduce an estimate \hat{G} for the adjacency graph and will then try to choose the signs of the estimated matrix \hat{K} such that the its principal minors are the estimates for the principal minors. For this define the edge set \hat{E} of \hat{G} to consist of all sets $\{i, j\}$ such that $\hat{B}_{ij} \geq \frac{1}{2}\alpha^2$. This truncation yields the desired effect that by the strong law of large numbers the estimator for the graph will converge to the actual adjacency graph almost surely. In analogy to the previous paragraph we define $\{\hat{C}_1, \dots, \hat{C}_{\hat{m}}\}, \hat{H}_1, \dots, \hat{H}_{\hat{m}}, \hat{A}$ and \hat{b} exactly the same way. If there is a solution $\hat{x} \in \mathbb{F}_2^E$ to the linear equation

$$\hat{A}\hat{x} = \hat{b}, \quad (3.6)$$

then we estimate the signs to be $\hat{w}_{ij} := \phi^{-1}(\hat{x}_{ij})$ and define

$$\hat{K}_{ij} := \hat{w}_{ij} \sqrt{\hat{B}_{ij}}.$$

If there is no such solution \hat{x} then we simply set the signs of the off diagonal elements to be positive, i.e. we define

$$\hat{K}_{ij} := \sqrt{\hat{B}_{ij}}.$$

This choice is completely arbitrary, but we will see in the consistency result 3.15 that the probability for this case tends to zero as the sample size increases. In fact we will see that the two linear equations (3.5) and 3.6 agree with increasing probability.

In order to talk about consistency of the estimator that we constructed above, it is necessary to define a metric on the marginal kernels of DPPs. However the usual operator norm is clearly not right for this job, since we already know that we can only hope to reconstruct the determinantal equivalence class but not the exact marginal kernel. Thus we will work with the usual choice of pseudometric if one has to deal with equivalence classes.

3.14 PSEUDOMETRIC ON THE MARGINAL KERNELS. We define the distance between two marginal kernels $A, B \in \mathbb{R}^{N \times N}$ through

$$d(A, B) := \inf_{C \sim A} \|B - C\|_\infty$$

where $\|A\|_\infty := \max_{1 \leq i, j \leq N} |A_{ij}|$ denotes the uniform norm on the space of matrices.

3.15 THEOREM (CONSISTENCY). *Let K be the marginal kernel of a DPP that satisfy the assumption 3.12. Let further l be the cycle sparsity of G_K and $\varepsilon > 0$.*

$$\mathbb{P}\left(d(\hat{K}, K) \leq \varepsilon\right) \rightarrow 1 \quad \text{for } n \rightarrow \infty.$$

Proof. We will keep the notations from the paragraphs 3.11 and 3.13. We have already seen in the motivation of this section that the empirical measures converge almost surely which directly implies

$$\hat{K}_{ii} \rightarrow K_{ii} \quad \text{and} \quad \hat{K}_{ij}^2 \rightarrow K_{ij}^2 \quad \text{almost surely.} \quad (3.7)$$

Note that almost surely convergence implies convergence in probability and thus we have

$$\mathbb{P}(\hat{G} = G_K) = \mathbb{P}\left(\hat{K}_{ij}^2 \geq \alpha^2/2 \text{ for } K_{ij} \neq 0\right) \rightarrow 1 \quad \text{for } n \rightarrow \infty.$$

In this case the two shortest cycle basis can be chosen the same and so \hat{A} and A agree. Because of (3.7) we also have $\hat{H}_k \rightarrow H_k$ almost surely and thus $\hat{b}_k \rightarrow b_k$ almost surely for all k . This yields

$$\mathbb{P}\left(\hat{A} = A \text{ and } \hat{b} = b\right) \rightarrow 1 \quad \text{for } n \rightarrow \infty. \quad (3.8)$$

In this case the two linear equations (3.5) and (3.6) agree, then $\tilde{K} \in \mathbb{R}^{N \times N}$ defined through $\tilde{K}_{ij} := \hat{w}_{ij} |K_{ij}|$ is determinantly equivalent to K . Further we know that for any $\delta > 0$

$$\mathbb{P} \left(\left| \hat{K}_{ij}^2 - K_{ij}^2 \right| < \delta \text{ for all } i, j \right) \rightarrow 1 \quad \text{for } n \rightarrow \infty. \quad (3.9)$$

If this is true, then we have

$$d(\hat{K}, K) \leq \|\hat{K} - \tilde{K}\|_\infty = \sup_{i,j} \left| |\hat{K}_{ij}| - |\tilde{K}_{ij}| \right|$$

where we used, that the entries of \hat{K} and \tilde{K} have equal signs. Further we have

$$\left| |\hat{K}_{ij}| - |\tilde{K}_{ij}| \right| = \frac{|\hat{K}_{ij}^2 - \tilde{K}_{ij}^2|}{|\hat{K}_{ij}| + |\tilde{K}_{ij}|} < \frac{\delta}{\alpha} \leq \varepsilon$$

if $\delta \leq \alpha\varepsilon$. In conclusion we have seen that if

$$\hat{A} = A, \quad \hat{b} = b \quad \text{and} \quad \left| \hat{K}_{ij}^2 - K_{ij}^2 \right| < \delta \quad \text{for all } i, j$$

then we have

$$d(\hat{K}, K) < \varepsilon.$$

However (3.8) and (3.9) shows that the probability for this tends to one. \square

3.16 REMARK (SPEED OF CONVERGENCE). Although the result above states that the estimators \hat{K} converges to K in probability, it does give no information about the speed of convergence. This problem is addressed in [Urschel et al., 2017], but it turns out that the convergence is very slow. For example for the very moderate case $\alpha = 0.4$ and $l = 3$ one already needs more than 10^6 samples to get some theoretical guarantees from their result. This is not due to careless estimates since they even show that this bound is optimal. However since this result is beyond practical relevance, we will keep away from those calculations.

is this true?

explain how the algorithm for the reconstruction works

III.1.4 Computation of the estimator

is this estimator unbiased?

III.2 Maximum likelihood estimation using optimisation techniques

rewrite introduction to MLE

The method of maximum likelihood estimation is a very well established procedure to estimate parameters. The philosophy of MLE is that one selects the parameter under which the given data would be the most likely to be observed and to motivate this in more detail we roughly follow the corresponding section in [Rice, 2006].

Suppose that we want to estimate a parameter $\theta \in \Theta$ based on some realisations x_1, \dots, x_n of some random variables X_1, \dots, X_n . We have given some candidates $f(x_1, \dots, x_n | \theta)$ for the joint density of X_1, \dots, X_n with respect to some reference measure $\prod_{i=1}^n \mu(dx_i)$ and we want to decide which parameter $\theta \in \Theta$ describes the realisations, which we will also call data or observation, best. Hence it is reasonable to pick that θ under which the observations x_1, \dots, x_n

are the most likely. In other words we want to find the parameter θ that maximises the density $f(x_1, \dots, x_n | \theta)$. If additionally the random variables are indepent and identically distributed, their joint density factorises and thus we obtain

$$f(x_1, \dots, x_n | \theta) = \prod_{i=1}^n f(x_i | \theta)$$

where $f(x | \theta)$ is the density with respect to μ of the X_i . In practice it is often easier to maximise the logarithm of the density

$$\mathcal{L}(\theta) = \log(f(x_1, \dots, x_n | \theta)) = \sum_{i=1}^n \log(f(x_i | \theta))$$

since this transforms the product over functions into a sum. However this is clearly equivalent to maximising the density since the logarithm is strictly monotone.

3.17 DEFINITION OF THE MLE. Let Θ be a set, which we call the *parameter set* and let

$$\mathcal{F} = \left\{ f(\cdot | \theta) : X \rightarrow [0, \infty) \mid \theta \in \Theta \right\}$$

be a family of probability densities with respect to some measure μ on some measurable space X . We call the function

$$\mathcal{L} : \Theta \rightarrow [-\infty, 0]$$

the *log likelihood function* and its maximiser

$$\hat{\theta}_n := \arg \max_{\theta \in \Theta} \mathcal{L}(\theta) \tag{3.10}$$

the *maximum likelihood estimator* or short MLE.

VERY SHORT REMINDER ON OPTIMISATION

Since the calculation of the MLE is a maximisation task, it is suitable to review some general properties of optimisation problems. It shall be noted that optimisation problems are usually stated as minimisation tasks, but we will stick to the maximisation, which is clearly equivalent up to a sign. For this let $U \subseteq \mathbb{R}^M$ and $f : U \rightarrow \mathbb{R}$ be a function. In practice the maximisation

$$\hat{x} := \arg \max_{x \in U} f(x)$$

will not be explicitly solvable and therefore one usually has to exploit numerical algorithms.

Those work particularly well if the function f is concave and possibly smooth and one powerful method is the given by the so called gradient descent. To quickly explain the philosophy of those methods, we note that ∇f points into the direction of the steepest ascent of the function f and thus an intuitive approach the maximise f would be to follow the gradient, i.e. to take a solution γ of the gradient flow $\gamma' = \nabla f(\gamma)$ and work out its limit. However if the function is not concave one can not even guarantee that the gradient flow reaches a local minimum, since one can construct examples where γ gets stuck in a critical point. However in the concave case this suffices since critical points and global minima agree for convex functions. The gradient descent

is an algorithm derived from this observation and is essentially a discretisations of the gradient flow meaning that it iteratively takes small steps into the direction of the gradient and thus lowers the value of the function. Some more sophisticated versions of gradient descent methods usually even consider higher order derivatives and use the information they provide over the geometry of the graph. Generally speaking those algorithms work extremely well even in high dimensions and thus their efficiency and stability have been studied broadly and we refer to the extensive monograph [Boyd and Vandenberghe, 2004]. All together we note that concavity is an extremely favourable property for a function that shall be maximised, which will be the log likelihood function later on.

A second property which is important in the existence theory of maximisers is *coercivity* in the sense that

$$f(x) \rightarrow -\infty \quad \text{for } |x| \rightarrow \infty.$$

In fact every (upper semi-) continuous and coercive function defined on a closed set $U \subseteq \mathbb{R}^M$ attains its minimum. To see this one can fix $x_0 \in U$ and use the coercivity to obtain $f < f(x_0)$ outside of a compact set K and thus the supremum of f agrees with the supremum of f over $K \cap U$ which is compact again and thus it is attained. We will later introduce some abstract theory about the consistency of estimators and for this we will need this result in a more general setting. However the version above is enough in the case of the maximum likelihood estimators for parameters of DPPs and therefore readers that are not familiar with elementary notions of topology are advised to neglect following statement.

3.18 PROPOSITION (EXISTENCE OF MAXIMISERS). *Let \mathcal{X} be a topological Hausdorff space and $f : \mathcal{X} \rightarrow [-\infty, \infty)$ be an upper semicontinuous function, i.e.*

$$L_f(\alpha) := \{x \in \mathcal{X} \mid f(x) \geq \alpha\}$$

is closed for all $\alpha \in \mathbb{R}$. Further we will assume that f is coercive, meaning that for any $\alpha \in \mathbb{R}$ the set $L_f(\alpha)$ is compact. Then f attains its maximum in at least one point, i.e. there is $\hat{x} \in \mathcal{X}$ such that

$$f(\hat{x}) = \sup_{x \in \mathcal{X}} f(x).$$

Proof. Let without loss of generality f be not identical to $-\infty$ because otherwise the statement is trivial. Then we have

$$\alpha := \sup_{x \in \mathcal{X}} f(x) > -\infty.$$

If we choose (α_n) to be strictly increasing towards α , then we get $L_f(\alpha_{n+1}) \subseteq L_f(\alpha_n)$ for all $n \in \mathbb{N}$ and further none of the sets $L_f(\alpha_n)$ is empty. By the Cantor intersection theorem¹ we get that also the intersection is non empty, i.e. there is

$$\hat{x} \in \bigcap_{n \in \mathbb{N}} L_f(\alpha_n).$$

This implies

$$f(\hat{x}) \geq \alpha_n \xrightarrow{n \rightarrow \infty} \alpha = \sup_{x \in \mathcal{X}} f(x).$$

□

¹A precise formulation can be found in the appendix.

III.2.1 Presentation of different models

Assume again that we have a set of observations $(\mathbf{Y}_n)_{n \in \mathbb{N}} \subseteq \mathcal{Y}$ drawn independently and according to the DPP. This time we want to find the maximum likelihood estimator for the elementary kernel and in order to do this we need to be able to express the density of the DPP which is nothing but the values of the elementary probabilities. Thus we will assume that we are dealing with L -ensembles in this section. Since the observations (\mathbf{Y}_n) are defined on some common probability space which we will denote by (Ω, \mathbb{P}) we will change the notation in this section and write

$$f(A|\theta) \propto \det(L(\theta)_A)$$

for the elementary probabilities of the DPP that arises from the parameter θ . Note that the elementary probabilities are nothing than the density with respect to the counting measure. We will now present the maximum likelihood estimators for different parametric classes, i.e. different families \mathcal{F} of DPPs.

MLE OF THE ELEMENTARY KERNEL L

The most intuitive parameter that one can estimate is the elementary kernel L itself since it parametrises the entire class of L -ensembles.

3.19 MAXIMUM LIKELIHOOD ESTIMATOR FOR L . We consider the parameter space $\Theta = \mathbb{R}_{\text{sym},+}^{N \times N}$ of positive definite symmetric matrices and the parametric family

$$\mathcal{F} = \left\{ f(\cdot, L) \mid L \in \mathbb{R}_{\text{sym},+}^{N \times N} \right\}$$

where $f(A, L) \propto \det(L_A)$ is the elementary probability of DPP with elementary kernel L . We seek to find the MLE

$$\hat{L}_n := \arg \max_{L \in \mathbb{R}_{\text{sym},+}^{N \times N}} \mathcal{L}(L).$$

The log likelihood function is now given by

$$\mathcal{L}: \mathbb{R}_{\text{sym},+}^{N \times N} \rightarrow [-\infty, 0], \quad L \mapsto \log \left(\prod_{i=1}^n f(\mathbf{Y}_i | L) \right).$$

Using (2.4) we get the expression

$$\mathcal{L}(L) = \sum_{i=1}^n \log (\det(L_{\mathbf{Y}_i})) - n \log (\det(L + I)). \quad (3.11)$$

Although the parametric family of that arises from the elementary kernels L gives a high variety of different associated L -ensembles, it will also make the computation of the MLE more complex. Therefore we will consider some smaller classes of L -ensembles, which will decrease the flexibility of the model, but make computation more efficient.

MLE OF THE QUALITIES

Unlike earlier we will not try to estimate the whole kernel L but only the qualities q_i of the items $i \in \mathcal{Y}$. More precisely we recall that we can parametrise the positive definite symmetric matrices L using the quality diversity parametrisation

$$(q, \phi) \mapsto \Psi(q, \phi) = L \quad \text{where } L_{ij} = q_i \phi_i^T \phi_j q_j.$$

Now we fix a diversity feature matrix $\hat{\phi}$, that we will usually model according to some perceptions we might have and set $\hat{S}_{ij} := \phi_i^T \phi_j$. We will now try to estimate the quality vector $q \in \mathbb{R}_+^N$ instead of the whole kernel L . This means that we optimise the likelihood function over a smaller set of kernels, namely the ones of the form $\Psi(q, \hat{\phi})$ for $q \in \mathbb{R}_+^N$. Obviously the maximal likelihood that can be achieved using this more restrictive model decreases since we consider less positive definite matrices and we have

$$\max_{q \in \mathbb{R}_+^N} \mathcal{L}(\Psi(q, \hat{\phi})) \leq \max_{L \in \mathbb{R}_{\text{sym},+}^{N \times N}} \mathcal{L}(L).$$

Although we can only expect a worse descriptive power of the observation, the hope is that the task of estimating only the qualities $q \in \mathbb{R}_+^N$ is more feasible which actually turn out to be true in certain cases. But before we investigate this, we clearly state our goal.

3.20 MAXIMUM LIKELIHOOD ESTIMATOR FOR THE QUALITY. This time we work with the parameter set $\Theta = \mathbb{R}_+^N$ and the parametric family

$$\mathcal{F} = \left\{ f(\cdot | q) \mid q \in \mathbb{R}_+^N \right\}$$

where $f(A, q) \propto \det(\Psi(q, \hat{\phi})_A)$ is the elementary probability of DPP with elementary kernel $\Psi(q, \hat{\phi})$. We aim to find the MLE of the quality vector $q \in \mathbb{R}_+^N$, in other words we set

$$\hat{q}_n := \arg \max_{q \in \mathbb{R}_+^N} \mathcal{L}(q)$$

where we perceive the likelihood function as a function of q .

Using (2.5) we obtain the following expression for the single summands of the log likelihood function

$$\log \left(\prod_{j \in Y_i} q_j^2 \right) + \log(\det(\hat{S}_{Y_i})) - \log \left(\sum_{A \subseteq \mathcal{Y}} \prod_{j \in A} q_j^2 \det(\hat{S}_A) \right) \quad (3.12)$$

and note that it is upper semicontinuous.

LOG LINEAR MODEL FOR THE QUALITIES

The motivation for restricting our ambitions of estimation to the qualities q_i rather than the whole elementary kernel $L \in \mathbb{R}_{\text{sym},+}^{N \times N}$ was to obtain a more tractable optimisation problem. Unfortunately we can tell from (3.12) that the log likelihood still isn't concave in q and in order to achieve this, we will introduce the following model for the qualities.

3.21 LOG LINEAR MODEL FOR THE QUALITIES AND MLE. From now on we will fix vectors $f_i \in \mathbb{R}^M$ for $i \in \mathcal{Y}$ and call them *feature vectors*. Further we set

$$q_i = \exp\left(\theta^T f_i\right) \quad \text{for } \theta \in \mathbb{R}^M$$

and will only consider quality vectors $q \in \mathbb{R}_+^N$ that have this form. To formulate the maximum likelihood estimator for θ we set $\Theta := \mathbb{R}^M$ and consider the parametric family

$$\mathcal{F} = \left\{ f(\cdot|\theta) \mid \theta \in \mathbb{R}^M \right\}$$

where $f(\cdot|\theta)$ is the density of the DPP with similarity kernel \hat{S} and qualities $q_i = \exp\left(\frac{1}{2}\theta^T f_i\right)$. Further we will consider the maximum likelihood estimator

$$\hat{\theta}_n := \arg \max_{\theta \in \mathbb{R}^M} \mathcal{L}(\theta)$$

where we regard \mathcal{L} again as function of θ .

3.22 REMARK. It shall be noted that although this log linear model seems to be a harsh restriction, it isn't a restriction at all, at least theoretically. If we take $M = N$ and choose f_i to be the unit vectors in \mathbb{R}^N , then this just a logarithmic transformation of the parameters and thus the maximal likelihood that can be achieved with this model does not change. In practice however it will be of interest to work with rather low dimensional parameters θ , because if the ground set \mathcal{Y} gets large, optimisation in \mathcal{R}^N can be inefficient. In this case of course the maximal likelihood under the optimal parameter may decrease. However the approximation of the optimal parameter might become possible again which justifies this sacrifice.

Under the assumption of a log linear model for the qualities the individual terms of the log likelihood function take the form

$$2 \cdot \theta^T \sum_{i \in Y} f_i + \det(\hat{S}_Y) - \log \left(\sum_{A \subseteq \mathcal{Y}} \exp \left(2 \cdot \theta^T \sum_{i \in A} f_i \right) \det(\hat{S}_A) \right). \quad (3.13)$$

MLE OF THE REPULSIVENESS PARAMETER

III.2.2 Coercivity and existence of the maximum likelihood estimators

A priori it is not clear that the maximum likelihood estimators exist and we will actually see that they do not exist in general. However one can still save this approach because the probability that they exist tends to 1 if the sample size increases. We begin by showing this for the MLE of the qualities and then we will adapt this proof to the other models.

MLE OF THE QUALITIES

The MLE \hat{q}_n does not exist for all realisations $(Y_n)_{n \in \mathbb{N}}$ of $(\mathbf{Y}_n)_{n \in \mathbb{N}}$. To see this we suppose that we have only one sample $Y_1 = \mathcal{Y}$ which is the whole set. The higher the qualities of the items are, the more likely this observation gets and therefore the maximum of the log likelihood function – which is 0 in this case – is not obtained. This can also be made rigorous in the following

computation. Under the assumption of constant qualities the log likelihood function takes the form

$$\log(q^{2N} \det(\hat{S}_{\mathcal{Y}})) - \log\left(\sum_{A \subseteq \mathcal{Y}} q^{2|A|} \det(\hat{S}_A)\right) = \log\left(\frac{q^{2N} \det(\hat{S}_{\mathcal{Y}})}{\sum_{A \subseteq \mathcal{Y}} q^{2|A|} \det(\hat{S}_A)}\right) \xrightarrow{q \rightarrow \infty} 0.$$

However this maximum is never attained, since for every L -ensemble we have $\mathbb{P}_L(\emptyset) > 0$ and therefore

$$\mathcal{L}(q) = \log\left(\mathbb{P}_{\Psi(q, \hat{S})}(\mathcal{Y})\right) < 0 \quad \text{for every } q \in \mathbb{R}_+^N.$$

The thing that goes wrong in this case is, that under the observation of the whole set \mathcal{Y} we would estimate a deterministic model that always selects the whole set, namely the DPP with marginal kernel I . Since all of the eigenvalues are 1 in this case, this DPP is not a L ensemble and therefore we can not describe it with the quality diversity decomposition. However if we assume that the data is actually generated by a L -ensemble, then such a scenario becomes unlikely as the sample size increases. We will fix this in the following result.

3.23 PROPOSITION (COERCIVITY AND EXISTENCE OF THE MLE). *Let $\mathbf{Y}_1, \mathbf{Y}_2, \dots$ be a sequence of independent and identically distributed point processes that fall in the class of L -ensembles. Then we have*

$$\mathbb{P}\left(\hat{q}_n \in \mathbb{R}_+^N \text{ exists}\right) \geq \mathbb{P}(\mathcal{L} \text{ is coercive}) \xrightarrow{n \rightarrow \infty} 1.$$

Proof. The first inequality is obvious since the log likelihood function is upper semicontinuous. We will show that \mathcal{L} is coercive if one of the observations is the emptyset. Then the claim follows from

$$\begin{aligned} \mathbb{P}(\mathcal{L} \text{ is coercive}) &\geq \mathbb{P}\left(\bigcup_{i=1}^n \{\mathbf{Y}_i = \emptyset\}\right) = 1 - \mathbb{P}\left(\bigcap_{i=1}^n \{\mathbf{Y}_i \neq \emptyset\}\right) \\ &= 1 - \mathbb{P}(\mathbf{Y}_1 \neq \emptyset)^n \xrightarrow{n \rightarrow \infty} 1 \end{aligned}$$

since we have $\mathbb{P}(\mathbf{Y}_1 \neq \emptyset) < 1$ for every L -ensemble.

So let Y_1, \dots, Y_n be some observations with $Y_i = \emptyset$ for at least one $i \in \{1, \dots, n\}$ and let $(q^k)_{k \in \mathbb{N}} \subseteq \mathbb{R}_+^N$ be a sequence such that $|q^k| \rightarrow \infty$. Note that it suffices to show that every subsequence of (q^k) contains a subsubsequence (q^l) such that

$$\mathcal{L}(q^l) \rightarrow -\infty \quad \text{for } l \rightarrow \infty.$$

Hence we fix a subsequence of (q^k) which we denote by (q^k) again in slightly abusive notation. Let (q^l) be a subsequence of (q^k) such that one coordinate diverges to infinity, i.e.

$$q_{j_0}^l \xrightarrow{l \rightarrow \infty} \infty \quad \text{for one } j_0 \in \{1, \dots, N\}.$$

The i -th summand of \mathcal{L} takes the form

$$-\log\left(\sum_{A \subseteq \mathcal{Y}} \prod_{j \in A} (q_j^l)^2 \det(\hat{S}_A)\right) \leq -\log\left((q_{j_0}^l)^2\right) \xrightarrow{l \rightarrow \infty} -\infty$$

where we used $\hat{S}_{\{j_0\}} = 1$. Because the other summands are non positive this implies

$$\mathcal{L}(q^l) \xrightarrow{l \rightarrow \infty} -\infty$$

which we had to show. \square

3.24 REMARK. The proof above should be read in the following way. The statement $q_{j_0}^l \rightarrow \infty$ is equivalent to a model that would always select the item j_0 . However since we have observed the empty set, the observations would be impossible under this model and thus the log likelihood function takes the value $-\infty$ for this model. An analogue argument shows that the estimated qualities are strictly positive with high probability if the actual qualities are strictly positive. This will be of interest for us if we consider the log linear model.

3.25 PROPOSITION (POSITIVITY OF THE MLE). *Assume that $\mathbf{Y}_1, \mathbf{Y}_2, \dots$ is a sequence of independent and identically distributed point processes that are distributed according to a L -ensemble with strictly positive qualities. Then we have*

$$\mathbb{P}\left(\hat{q}_n \in \mathbb{R}_+^N \text{ exists and } \hat{q}_n \in (0, \infty)^N\right) \xrightarrow{n \rightarrow \infty} 1.$$

Proof. We have already seen that the probability that the MLE exists tends to one, so we only have to show that the probability that the estimated qualities are strictly positive tends to one. The philosophy to prove this is exactly the same than in the proof of existence. Indeed we note that once j occurs in one of the observations Y_1, \dots, Y_n we have $\mathcal{L}(q) = -\infty$ for every $q \in \mathbb{R}_+^N$ with $q_j = 0$. Therefore we have $(\hat{q}_n)_j > 0$ if $j \in Y_i$ for at least one $j \in \{1, \dots, n\}$. Finally we note that the probability that j occurs in the i -th sample is strictly positive since we have

$$\mathbb{P}(j \in \mathbf{Y}_i) \geq \mathbb{P}(\{j\} = \mathbf{Y}_i) = q_j^2 > 0.$$

\square

MLE OF THE ELEMENTARY KERNEL

We can quite easily adapt the proof for the existence of MLEs of the qualities to the case of MLEs for the whole elementary kernel L .

3.26 PROPOSITION (COERCIVITY AND EXISTENCE OF MLE). *Let $\mathbf{Y}_1, \mathbf{Y}_2, \dots$ be an sequence of independent and identically distributed point processes that fall in the class of L -ensembles. Then we have*

$$\mathbb{P}\left(\hat{L}_n \in \mathbb{R}_{\text{sym},+}^{N \times N} \text{ exists}\right) \geq \mathbb{P}(\mathcal{L} \text{ is coercive}) \xrightarrow{n \rightarrow \infty} 1.$$

Proof. Again it suffices to show $\mathcal{L}(L) \rightarrow -\infty$ for $|L| \rightarrow \infty$ once we have observed the empty set once. To see this, we use the quality diversity parametrisation

$$\Psi: \mathbb{R}_+^N \times \mathbb{S}_N^N \rightarrow \mathbb{R}_{\text{sym},+}^{N \times N}, \quad (q, \phi) \mapsto \left(q_i \phi_i^T \phi_j q_j\right)_{1 \leq i, j \leq N}.$$

Note that since Ψ is continuous and therefore bounded on bounded sets and \mathbb{S}_N^N is bounded, $|\Psi(q, \phi)| \rightarrow \infty$ implies $|q| \rightarrow \infty$. The exact same calculations as in the previous proof show

$$\mathcal{L}(L) = \mathcal{L}(\Psi(q, \phi)) \rightarrow -\infty \quad \text{for } |L| \rightarrow \infty.$$

\square

THE LOG LINEAR MODEL

We have seen that the log linear model can provide a parametrisation of the whole space $(0, \infty)^N$ of possible qualities. However it can also be very restrictive, for example if all feature vectors are trivial, i.e. $f_i = 0$ for all items i . Hence we need to convince ourselves that we do not loose too much information through the transformation

$$F : \mathbb{R}^M \rightarrow (0, \infty)^N, \quad \theta \mapsto (\exp(\theta^T f_1), \dots, \exp(\theta^T f_N))^T.$$

In order to do this, let $U \subseteq \mathbb{R}^M$ be the span of f_1, \dots, f_N and let write $\theta = \theta_1 + \theta_2$ such that $\theta_1 \in U$ and $\theta_2 \in U^\perp$. We note that $F(\theta) = F(\tilde{\theta})$ if and only if $\theta_1 = \tilde{\theta}_1$.

3.27 PROPOSITION (COERCIVITY AND EXISTENCE OF MLE). *Assume that $\mathbf{Y}_1, \mathbf{Y}_2, \dots$ is a sequence of independent and identically distributed point processes that are distributed according to a L -ensemble with strictly positive qualities. Then we have*

$$\mathbb{P}(\hat{\theta}_n \in \mathbb{R}^M \text{ exists}) \geq \mathbb{P}(\mathcal{L} \text{ is coercive as a function on } U) \xrightarrow{n \rightarrow \infty} 1.$$

Proof. First we note, that it suffices to show that \mathcal{L} has a maximiser on U , since $F(U) = F(\mathbb{R}^M)$. To do this we show – just like in the previous cases – that the \mathcal{L} is coercive on U whenever we have observed the emptyset as well as every item at least once. Let now $(\theta^k)_{k \in \mathbb{N}} \subseteq U$ be a sequence such that $|\theta^k| \rightarrow \infty$. Then there is at least one index $i \in \{1, \dots, N\}$ and a subsequence $(\theta^l)_{l \in \mathbb{N}}$ such that

$$f_i^T \theta^l \rightarrow \infty \quad \text{or} \quad f_i^T \theta^l \rightarrow -\infty \quad \text{for } l \rightarrow \infty$$

since otherwise all sequences $(f_i^T \theta^l)$ therefore also (θ^l) would be bounded. However this is equivalent to

$$\exp(f_i^T \theta^l) \rightarrow \infty \quad \text{or} \quad \exp(f_i^T \theta^l) \rightarrow 0 \quad \text{for } l \rightarrow \infty$$

and we have seen in the proof of 3.25 that the log likelihood function tends to $-\infty$ in this case. \square

MLE FOR THE REPULSIVENESS PARAMETER

III.2.3 Consistency of the maximum likelihood estimators

We will now turn towards the question of consistency of the maximum likelihood estimators introduced earlier in this section. For this we will first give a formal proof of the consistency of the MLE and then present a rather general framework that will allow us to turn the formal proof into a rigorous one.

3.28 FORMAL PROOF OF CONSISTENCY. We will consider a general MLE like in (3.10) and we will assume that the observations (X_n) are independent and have density $f(x|\theta_0)$ with respect to some measure μ . By the law of large number we have

$$\frac{1}{n} \mathcal{L}(\theta) = \frac{1}{n} \sum_{i=1}^n \log(f(X_i|\theta)) \xrightarrow{n \rightarrow \infty} \mathbb{E}[\log(f(X|\theta))]. \quad (3.14)$$

Hence the maximiser of the left hand side should be close to the maximiser of the right hand. Differentiating the right hand side yields

$$\begin{aligned} \partial_\theta \mathbb{E}[\log(f(X|\theta))] &= \mathbb{E}[\partial_\theta \log(f(X|\theta))] = \mathbb{E}\left[\frac{\partial_\theta f(X|\theta)}{f(X|\theta)}\right] \\ &= \int \frac{\partial_\theta f(x|\theta)}{f(x|\theta)} f(x|\theta_0) \mu(dx). \end{aligned}$$

Evaluating this at $\theta = \theta_0$ yields

$$\int \partial_\theta f(x|\theta) \mu(dx) = \partial_\theta \int f(x|\theta) \mu(dx) = \partial_\theta(1) = 0.$$

Hence θ_0 is a critical point and under mild conditions the right hand side is concave and thus θ_0 is the unique maximiser. In conclusion the estimator $\hat{\theta}$ should be close to θ_0 .

Although the rough structure of the rigorous proof is present in the argument above it is highly formal. For example we argue that if a sequence $(f_n)_{n \in \mathbb{N}}$ of functions converges towards f pointwise, then the maximisers $(x_n)_{n \in \mathbb{N}}$ should converge to the maximiser x of f . The major tool to make this rigorous will be to use some kind of uniform convergence. Namely we have the following result where we will omit the proof since it is very easy and we prove a similar but stronger version of it later.

3.29 LEMMA (SWAPPING LIMIT AND MAXIMISATION). *Let $(f_n)_{n \in \mathbb{N}}$ be a sequence of real functions on a compact space with maximisers $(x_n)_{n \in \mathbb{N}}$ that are bounded from above and converge uniformly towards f . Further assume that f is continuous and has a unique maximum in x_0 . Then we have $x_n \rightarrow x_0$ for $n \rightarrow \infty$.*

Unfortunately the convergence in (3.14) does only hold uniformly on a compact set $K \subseteq \Theta$. To fix this we will argue that the maximisers (x_n) lie in this compact set K for large n . We will do this in a general setup in the next paragraph.

A GENERAL CONSISTENCY RESULT FOR EXTREMAL ESTIMATORS

We will provide a general consistency result for a rather broad class of estimators which is taken from [Newey and McFadden, 1994] and slightly adapted to our needs. Although it would be possible to prove the consistency of the MLEs directly we present this general procedure since this can easily be adjusted to other cases.

3.30 SETTING. Let in the following Θ be a topological Hausdorff space and $F_n: \Theta \rightarrow [-\infty, \infty)$ be a sequence of random functions with maximisers

$$\hat{\theta}_n := \arg \max_{\theta \in \Theta} F_n(\theta).$$

If no maximiser exists, we choose $\hat{\theta}_n \in \Theta$ arbitrary. Further let $F: \Theta \rightarrow [-\infty, \infty)$ be a deterministic function with maximiser θ_0 . The maximisers $\hat{\theta}_n$ are called *extremal estimators* since they are the extremal points of the functions F_n .

We now investigate whether the extremal estimators converge to the maximiser θ_0 .

3.31 THEOREM (CONSISTENCY OF EXTREMAL ESTIMATORS). *Let the setting be as above and assume that the following conditions hold.*

- (i) *Assume that there is $\varepsilon_0 > 0$ and a compact set K_0 containing θ_0 , such that with probability tending to one*

$$F_n(\theta) \leq F(\theta_0) - \varepsilon_0 \quad \text{for all } \theta \notin K_0. \quad (3.15)$$

- (ii) *Let F_n converge to F uniformly on K_0 in probability, i.e. for any $\varepsilon > 0$ we have with probability tending to one*

$$|F_n(\theta) - F(\theta)| \leq \varepsilon \quad \text{for all } \theta \in K_0. \quad (3.16)$$

(iii) Let F have a unique maximum at $\theta_0 \in \Theta$.

(iv) Assume that F is upper semicontinuous in the sense that

$$\{\theta \in \Theta \mid F(\theta) \geq \alpha\} \subseteq \Theta$$

is closed for all $\alpha \in \mathbb{R}$.

(v) With probability tending to one F_n admits a maximiser.

Then we have $\hat{\theta}_n \rightarrow \theta_0$ in probability, i.e.

$$\mathbb{P}(\hat{\theta}_n \in U) \xrightarrow{n \rightarrow \infty} 1$$

for any open subset $U \subseteq \Theta$ containing θ_0 .

Proof. Note that it suffices to show $\hat{\theta}_n \in U$ whenever (3.15) and (3.16) hold and F_n admits a maximiser. From here on the proof is of purely analytic content.

Fix now an open set $U \subseteq \Theta$ that contains θ_0 . Choosing $\varepsilon < \varepsilon_0$ in (ii) and using (i) yields

$$F_n(\theta_0) \geq F(\theta_0) - \varepsilon > F(\theta_0) - \varepsilon_0 \geq F_n(\theta) \quad \text{for all } \theta \notin K_0.$$

Hence the maximum of F_n is attained in K_0 and we have $\hat{\theta}_n \in K_0$. Thus if $K_0 \subseteq U$ we are done. If this is not the case F attains its maximum α on $K_0 \setminus U$ because F is upper semicontinuous and $K_0 \setminus U$ is compact (c.f. 3.18). Further (iii) implies $\alpha < F(\theta_0)$ and thus we have

$$K_0 \cap \{\theta \in \Theta \mid F(\theta) > \alpha\} \subseteq U.$$

So in order to show $\hat{\theta}_n \in U$, it suffices to show $\hat{\theta}_n \in K_0$ and $F(\hat{\theta}_n) > \alpha$. Since we have already seen that the first statement holds, it remains to show the second one. However (ii) implies

$$F(\hat{\theta}_n) \geq F_n(\hat{\theta}_n) - \varepsilon \geq F_n(\theta_0) - \varepsilon \geq F(\theta_0) - 2\varepsilon > \alpha$$

for ε small enough. □

maybe add a sketch?

If we want to apply the previous result to the case of maximum likelihood estimation we need to set

$$F_n(\theta) := \frac{1}{n} \sum_{i=1}^n \log(f(X_i | \theta)).$$

Note that the factor $\frac{1}{n}$ does not change the maximum. However (3.14) already gives the almost surely pointwise limit of those functions and if condition (ii) of the previous statement should hold, we have to define

$$F(\theta) := \mathbb{E}[\log(f(X | \theta))].$$

The quantity F is known as the *entropy* and plays an important role in many different fields, for example statistical mechanics, applied statistics and information theory. For further reading we refer to [Martin and England, 2011], [MacKay, 2003], [Volkenstein, 2009] and [Gray, 1990].

INFORMATION INEQUALITY AND LOCALLY UNIFORM CONVERGENCE

The second and third requirement of the previous result can be proven in a general setting and without quantitative assumption and we adapt an argument from [Newey and McFadden, 1994] to fit our needs. In order to do this we will work with the following assumptions.

3.32 SETTING. Let in the following Θ be a set and let

$$\mathcal{F} = \left\{ f(\cdot|\theta) : \mathcal{X} \rightarrow [0, \infty) \mid \theta \in \Theta \right\}$$

be a family of probability densities on some measurable space \mathcal{X} with respect to some measure μ . Further fix $\theta_0 \in \Theta$ and denote the expectation with respect to $f(\cdot|\theta_0)d\mu$ by

$$\mathbb{E}[h(X)] := \int h(x) f(x|\theta_0) \mu(dx).$$

Let (X_n) be a sequence of independent random variables distributed according to $f(\cdot|\theta_0)d\mu$. Finally define

$$F_n(\theta) := \frac{1}{n} \sum_{i=1}^n \log(f(X_i|\theta)) \quad \text{and} \quad F(\theta) := \mathbb{E}[\log(f(X|\theta))].$$

3.33 PROPOSITION (INFORMATION INEQUALITY). *Let the setting be as above and assume that the parameter $\theta_0 \in \Theta$ is identifiable, i.e. we have $f(\cdot|\theta) \neq f(\cdot|\theta_0)$ whenever $\theta \neq \theta_0$. Let further*

$$\sup_{x \in \mathcal{X}, \theta \in \Theta} f(x|\theta) < \infty \quad \text{and} \quad F(\theta_0) > -\infty.$$

Then the entropy

$$F(\theta) = \mathbb{E}[\log(f(X|\theta))]$$

has a unique maximum in θ_0 .

Proof. Let $\theta \neq \theta_0$, then we either have $F(\theta) = -\infty < F(\theta_0)$ or

$$F(\theta) = \mathbb{E}[\log(f(X|\theta))] > -\infty. \tag{3.17}$$

In this case we want to exploit the strict Jensen inequality (c.f. [Lehmann and Casella, 2006]) that yields for any positive random variable Y with finite expectation that is not constant

$$\mathbb{E}[\log(Y)] < \log(\mathbb{E}[Y]).$$

We set $Y := \frac{f(X|\theta)}{f(X|\theta_0)}$. This is positive $f(\cdot|\theta_0)d\mu$ almost everywhere because otherwise (3.17) could not hold. Since θ_0 is identifiable, the random variable Y is not constant and we will see in the following computation that the expectation is finite. Now we obtain

$$\begin{aligned} F(\theta) - F(\theta_0) &= \mathbb{E}[\log(f(X|\theta))] - \mathbb{E}[\log(f(X|\theta_0))] = \mathbb{E} \left[\log \left(\frac{f(X|\theta)}{f(X|\theta_0)} \right) \right] \\ &< \log \left(\mathbb{E} \left[\frac{f(X|\theta)}{f(X|\theta_0)} \right] \right) = \log \left(\int f(x|\theta) \mu(dx) \right) = 0. \end{aligned}$$

□

Next we take care of the second requirement of the consistency result. Namely we will show that the functions F_n associated with the MLE almost surely converge to F locally uniformly under fairly mild conditions. For this we modify the proof of a more general convergence result in [Tauchen, 1985].

3.34 LEMMA (LOCALLY UNIFORM CONVERGENCE). *Let the setting be as above, but let Θ be a metric space and let $K \subseteq \Theta$ be compact such that the following conditions hold.*

(i) *Let*

$$\mathbb{E} \left[\sup_{\theta \in K} |\log(f(X|\theta))| \right] < \infty.$$

(ii) *For every $\theta \in K$ we have $\log(f(\cdot, \gamma)) \rightarrow \log(f(\cdot|\theta))$ almost surely with respect to $f(\cdot|\theta_0)d\mu$ for $\gamma \rightarrow \theta$.*

Then we almost surely have $F_n \rightarrow F$ uniformly on K , i.e. almost surely

$$\sup_{\theta \in K} |F_n(\theta) - F(\theta)| \xrightarrow{n \rightarrow \infty} 0.$$

Proof. Fix $\varepsilon > 0$ and define for $x \in \mathcal{X}$ and $\rho > 0$

$$u(x, \theta, \rho) := \sup_{d(\gamma, \theta) \leq \rho} |\log(f(x|\gamma)) - \log(f(x|\theta))| \xrightarrow{\rho \rightarrow \infty} 0$$

almost surely for θ fixed where we used condition (ii).

This in combination with (i) and the dominated convergence theorem imply that the convergence also holds in expectation and therefore we have

$$\mathbb{E}[u(X, \theta, \rho)] \leq \varepsilon \quad \text{for } \rho \leq \delta(\theta).$$

The open balls $B_{\delta(\theta)}(\theta)$ with center θ and radius $\delta(\theta)$ cover the compact set K and hence we can select a finite subcover

$$K \subseteq \bigcup_{k=1}^m B_{\delta(\theta_k)}(\theta_k).$$

Further we set

$$\mu_k := \mathbb{E}[u(X, \theta_k, \delta(\theta_k))] \leq \varepsilon.$$

Let $\theta \in K$ and choose k such that $\theta \in B_{\delta(\theta_k)}(\theta_k)$, then we can conclude

$$\begin{aligned} |F_n(\theta) - F(\theta)| &\leq \frac{1}{n} \sum_{i=1}^n |\log(f(X_i|\theta)) - \log(f(X_i|\theta_k))| \\ &\quad + \left| \frac{1}{n} \sum_{i=1}^n \log(f(X_i|\theta_k)) - F(\theta_k) \right| + |F(\theta_k) - F(\theta)| \\ &\leq \left(\frac{1}{n} \sum_{i=1}^n u(X_i, \theta_k, \delta(\theta_k)) - \mu_k \right) + \mu_k + 2\varepsilon \\ &\leq 4\varepsilon \end{aligned}$$

almost surely for $n \geq N(\varepsilon)$ where we used the strong law of large numbers twice. \square

CONSISTENCY OF THE MLES FOR THE QUALITY AND ELEMENTARY KERNEL

In this part we will – for the first time – make use of the specific structure of the model. Since we have already taken care of conditions (ii), (iii) and (v) and condition (iv) will be fairly straight forward, we dedicate ourselves to proving the first requirement of Theorem 3.31. For this we keep the setting of the previous section although we now consider the case that

$$\mathcal{F} = \left\{ f(\cdot|\theta) : 2^{\mathcal{Y}} \rightarrow [0, \infty) \mid \theta \in \Theta \right\}$$

is one of the parametric families for the L -ensembles introduced in III.2.1. Further we denote a realisation of a DPP by \mathbf{Y} like earlier.

3.35 LEMMA (CONTROL OUTSIDE OF A COMPACT SET). *The requirement (i) from Theorem 3.31 is satisfied for the three kinds of parametric families for the kernel estimation. Further the compact set K_0 can be chosen as follows. Let \mathcal{A} be the family of subsets $A \subseteq \mathcal{Y}$ with positive probability $f(A|\theta_0) > 0$ and let $c(A) > 0$ such that*

$$-c(A) < \frac{2 \cdot F(\theta_0)}{f(A|\theta_0)}.$$

Then we set

$$K_0 := \left\{ \theta \in \Theta \mid \log(f(A|\theta)) \geq -c(A) \text{ for all } A \in \mathcal{A} \right\}.$$

Proof. At first we note that $F(\theta_0) > -\infty$. Let now

$$\hat{\mathbb{P}}_n := \frac{1}{n} \sum_{i=1}^n \delta_{\mathbf{Y}_i}$$

be the empirical measure. We have by the law of large numbers

$$\mathbb{P} \left(\hat{\mathbb{P}}_n(A) \geq \frac{f(A|\theta_0)}{2} \right) \xrightarrow{n \rightarrow \infty} 1$$

and we can assume $\hat{\mathbb{P}}_n(A) \geq \frac{f(A|\theta_0)}{2}$, since we are only interested in proving a statement with probability tending to one. For $A \in \mathcal{A}$ we note that

$$K_A := \left\{ \theta \in \Theta \mid \log(f(A|\theta)) \geq -c(A) \right\}$$

which is closed since $f(A|\theta)$ is upper semicontinuous. Further K_A is compact for $A = \emptyset \in \mathcal{A}$ since $\log(f(\emptyset|\theta))$ is coercive in θ as been shown in the coercivity proofs earlier in this chapter. Further it contains θ_0 as

$$\log(f(A|\theta_0)) \geq 2 \cdot \log(f(A|\theta_0)) > -c(A)$$

because $f(A|\theta) \leq 1$. Now

$$K_0 = \bigcap_{A \in \mathcal{A}} K_A$$

is compact because K_{\emptyset} is compact. Take now $\theta \notin K_0$, lets say $\theta \notin K_A$, then we get

$$\begin{aligned} F_n(\theta) &= \int \log(f(x|\theta)) \hat{\mathbb{P}}_n(dx) = \sum_{B \in \mathcal{A}} \hat{\mathbb{P}}_n(B) \cdot \log(f(B|\theta)) \leq \hat{\mathbb{P}}_n(A) \cdot \log(f(A|\theta)) \\ &< \frac{f(A|\theta_0)}{2} \cdot c(A) < F(\theta_0) \end{aligned}$$

where we used $f(B|\theta) \leq 1$. □

Now we have all the auxiliary results to prove the desired consistency result.

3.36 THEOREM (CONSISTENCY). (i) *The maximum likelihood estimator \hat{L}_n for the elementary kernel is consistent. Namely if the observations (\mathbf{Y}_n) follow the law of a L -ensemble with kernel L_0 , then we have*

$$\mathbb{P} \left(d(\hat{L}_n, L_0) \leq \varepsilon \right) \xrightarrow{n \rightarrow \infty} 1 \quad \text{for all } \varepsilon > 0.$$

(ii) *The maximum likelihood estimator \hat{q}_n for the quality vector is consistent. Namely if the observations (\mathbf{Y}_n) follow the law of a L -ensemble with kernel $\Psi(p_0, \hat{S})$, then we have*

$$\mathbb{P} \left(\|\hat{q}_n - q_0\| \leq \varepsilon \right) \xrightarrow{n \rightarrow \infty} 1 \quad \text{for all } \varepsilon > 0.$$

(iii) *Suppose that the observations (\mathbf{Y}_n) follow the law of a L -ensemble with kernel $\Psi(p_0, \hat{S})$ where $(p_0)_i = \exp(\theta^T f_i)$ and let P denote the projection onto the subspace U . Then we have*

$$\mathbb{P} \left(\|P\hat{\theta}_n - P\theta_0\| \leq \varepsilon \right) \xrightarrow{n \rightarrow \infty} 1 \quad \text{for all } \varepsilon > 0.$$

Proof. We will only sketch the main parts of the proof of the second statement, since all other arguments will be mostly analogue and therefore redundant.

Obviously we want to exploit the machinery we have introduced and thus we will check the requirements of Theorem 3.31. First we note that (v) holds because of the existence of the maximum likelihood estimators.

We can express the entropy function

$$F(q) = \mathbb{E}[\log(f(\mathbf{Y}|q))] = \sum_{A \subseteq \mathcal{Y}} \log(f(A|q)) f(A|q_0) \quad (3.18)$$

where the elementary probabilities are given by

$$f(A|q) = \frac{\prod_{i \in A} q_i^2 \det(\hat{S}_A)}{\sum_{B \subseteq \mathcal{Y}} \prod_{i \in B} q_i^2 \det(\hat{S}_B)} \quad (3.19)$$

which is continuous in q . Hence the entropy function F is upper semicontinuous and thus condition (iv) holds.

To check that (iii) holds we will use the information inequality 3.33. First we note that because of

$$f(\{i\} | q) \propto q_i^2$$

the parameter q_0 is identifiable and further we have

$$\sup_{A \subseteq \mathcal{Y}, q \in \mathbb{R}_+^N} f(A|q) \leq 1$$

since the densities are elementary probabilities. Finally $F(q_0) > -\infty$ is clear from (3.18) and hence the third requirement is satisfied.

Since the previous lemma already takes care of condition (i) it suffices to show the second condition for which we will use 3.34. Hence it remains to check the two conditions of this lemma,

but the second one – the continuity condition – obviously holds as can be seen from (3.19). To see that the first one also holds we note that for $A \subseteq \mathcal{Y}$ with $f(A|q_0) > 0$ and $\theta \in K_0$ we have

$$0 \geq \log(f(A|q)) \geq -c(A) > -\infty.$$

Hence the random variable

$$\sup_{q \in K_0} |\log(f(\mathbf{Y}|q))|$$

is almost surely finite and since the probability space $2^{\mathcal{Y}}$ is finite, the second condition holds. \square

3.37 REMARK. Obviously in the proof of the consistency of the whole elementary kernel L and the log linearity constant θ one runs into the problem of unidentifiability. This is why one has to identify the parameters with each other that give rise to the same probability densities. In the case of the elementary kernel this is just the determinantal equivalence once again and in the case of the log linearity constant two parameters give rise to the same densities if and only if their projection onto U agrees.

III.2.4 Approximation of the MLE

LIKELIHOOD MAXIMISATION FOR L

We note that \mathcal{L} is smooth and that its gradient can be expressed explicitly, at least on the domain $\{\mathcal{L} > -\infty\}$. This is due to the fact that the determinants of the submatrices are polynomials in the entries of L and the composition of those with the smooth function $\log: (0, \infty) \rightarrow \mathbb{R}$ stays smooth. This property allows the use of gradient methods but they face the problem that the loss function is non concave and thus those algorithms will generally not converge to a global maximiser. To see that the log linear likelihood function is not concave, we may consider the span $\{qI \mid q \in \mathbb{R}\}$ of the identity matrix. On this subspace \mathcal{L} takes the form

$$\mathcal{L}(qI) = \sum_{i=1}^n \log(q^{|Y_i|}) - n \log((1+q)^N) = \sum_{i=1}^n |Y_i| \log(q) - nN \log(1+q)$$

which is not concave in general.

This obviously causes substantial computational problems in the calculation of the MLE let alone it exists. In fact it is NP hard to maximise a general non concave function and it is also conjectured to be NP hard to maximise the log likelihood function \mathcal{L} in the case of L -ensembles. However there are still efficient maximising techniques for such functions that will eventually converge to local maximiser and that also work in very high dimensional spaces and thus this approach was taken by . Nevertheless we will not present this approach here, but rather favour a maximisation technique that is based on a fixed point iteration and was proposed in .

explain this term

cite

cite

cite

FIXED POINT ITERATION BASED MAXIMISATION

read, understand and summarise the paper

COMPUTATION FOR THE LOG LINEAR MODEL

The first two terms are affine linear in θ and thus concave. To see that the last expression is also concave, it is convenient to introduce the notion of log concavity and give a fundamental result.

3.38 DEFINITION (LOG CONCAVITY). We call a function f *log concave*, *log convex* or *log (affine) linear* if $\log(f)$ has the respective property.

3.39 PROPOSITION (ADDITIVITY OF LOG CONCAVITY). *The sum of log concave functions is again log concave.*

Give of cite
proof.

Proof.

□

As an immediate consequence we obtain that the expression in (3.4) is log concave which we will fix in a separate statement.

3.40 COROLLARY (CONCAVITY OF THE LIKELIHOOD FUNCTION). *Under the log linear model for the qualities, the log likelihood function is concave in the log linearity parameter $\theta \in \mathbb{R}^M$.*

III.2.5 Learning for conditional DPPs

III.2.6 Estimating the mixture coefficients of k -DPPs

Chapter IV

Bayesian learning for DPPs

So far we have seen two different estimation techniques for the parameters of DPPs. Although we proved that they provide reasonable estimators in the sense that they are consistent, they have some drawbacks. For example we have seen that the MLEs for the different parameters do not exist in general, let alone that they are impossible to compute in reality. Further all of the estimators presented so far are point estimators, i.e. they return a single value for the desired parameter. Obviously this does not allow to capture any uncertainties that the estimation of the parameter has. Those are some reasons to consider the Bayesian approach of parameter estimation where the goal is to give a distribution – called the posterior – of the parameter that should be estimated instead of a single value. This can also help to overcome some – maybe even all of the problems presented above.

At first we will present the general concept of Bayesian parameter estimation and will then turn towards the question of computability of the posterior distribution. For this we will follow the approach of [Affandi et al., 2014a] and turn towards the popular Markov chain Monte Carlo (MCMC) methods and quickly explain their philosophy and how they can be used to approximate the posterior distribution of the parameter one wishes to estimate.

IV.1 Bayesian approach to parameter estimation

For the introduction of the general Bayesian setup we pursue like in [Rice, 2006]. Just like in the case of MLE we want to estimate a parameter $\theta \in \Theta$ based on some realisations $x = (x_1, \dots, x_n)$ of random variables $X = (X_1, \dots, X_n)$. This time however we are not interested in returning a single value θ because this would be a vast simplification of the stochastic nature of the estimator. Thus we want to obtain a probability distribution over whole parameter space Θ that indicates how likely the parameters are to have caused the observed data. In order to present the procedure we will introduce the frame we will work in.

4.1 SETTING. Let Θ be a measurable space and ν be a measure on Θ . Let $f_\Theta: \Theta \rightarrow [0, \infty]$ be a probability density with respect to ν , i.e.

$$\int_{\Theta} f_{\Theta}(\theta) \nu(d\theta) = 1$$

which we will call the *prior* distribution of the parameter θ . Further let

$$\mathcal{F} = \{f_{X|\Theta}(\cdot|\theta) \mid \theta \in \Theta\}$$

by a family of probability densities with respect to $\mu^n := \prod_{i=1}^n \mu(dx_i)$.

Usually the prior distribution will encode some perceptions or prior knowledge we might have of the parameter. For example if we are trying to estimate a physical constant that we know has to be positive, then it is reasonable to select a prior that has its whole mass on the positive real line. However there is no clear set of rules how one can select a suitable prior to a given problem.

The density $f_{X|\Theta}(x|\theta)$ describes how likely the observations are under the parameter θ and we want to find an expression of how likely the parameter θ is under the observations x . In order to obtain this, we will work with the joint density

$$f_{X,\Theta}(x, \theta) = f_{X|\Theta}(x|\theta)f_{\Theta}(\theta) \quad \text{with respect to } \mu^n \times \nu$$

and condition this onto x . This yields

$$f_{\Theta|X}(\theta|x) = \frac{f_{X,\Theta}(x, \theta)}{\int_{\Theta} f_{X,\Theta}(x, \theta)\nu(d\theta)} = \frac{f_{X|\Theta}(x|\theta)f_{\Theta}(\theta)}{\int_{\Theta} f_{X,\Theta}(x, \theta)\nu(d\theta)} \quad (4.1)$$

4.2 DEFINITION (POSTERIOR DISTRIBUTION). The density $f_{\Theta|X}$ is called the *posterior distribution* of the parameter θ given the data x .

add remark on intuition?

First we will convince ourselves that the approach of calculating a posterior distribution is a generalisation of the MLE in a lot of cases.

4.3 COMPARISON TO MLE. Maybe one feels slightly uncomfortable with the need to choose a prior distribution and it turns out that this is in fact a difficult step that has to be taken with a certain amount of care. However we could pretend for one moment to be completely ignorant in the sense that we do not know anything about the parameter and hence we don't feel in the position to propose a reasonable prior. Then we could simply choose the uniform distribution as a prior – given it exists – and would obtain

$$f_{\Theta|X}(\theta|x) \propto f_{X|\Theta}(x|\theta).$$

Hence we can regain the MLE from our posterior distribution since it is just the mode, i.e. the maximiser of the posterior density. This relation to the MLE can be seen in Figure V.2. Hence the Bayesian approach is a more general tool than MLE and allows also to capture the random uncertainty of the parameter θ . This is desirable since we have seen that the mode is not always a very typical outcome of a random variable.

cite

A further advantage over the MLE is that it might be possible to computationally approximate the posterior density but not the MLE. This is typically the case if the log likelihood function is not concave, like in the setting of the MLE of the whole elementary kernel L . In fact the only hard step in the calculation of the posterior (4.1) is the computation of the normalisation constant

$$\int_{\Theta} f_{X,\Theta}(x, \theta)\nu(d\theta).$$

This step can often not be performed efficiently but the Markov chain Monte Carlo methods introduced later will yield an approximation of the posterior without the need to compute the normalisation constant.

4.4 REGULARISATION THROUGH THE PRIOR.

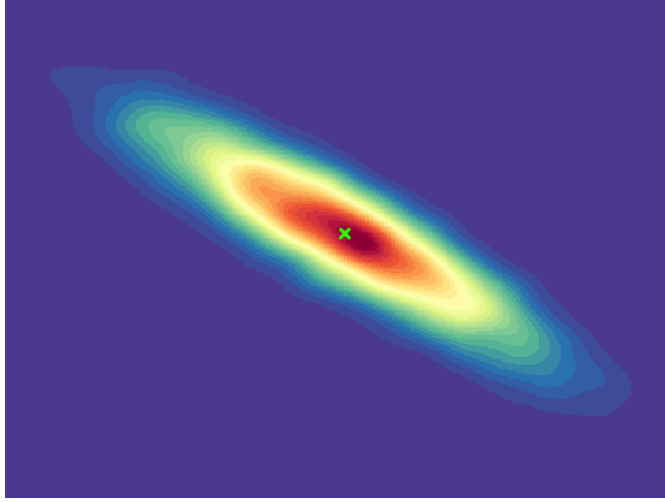


Figure IV.1.: Approximated posterior density of the two dimensional log linearity constant of a two dimensional DPP with a uniform distribution as a prior. The MLE estimator is marked green and is at the mode of the distribution.

EXPRESSION OF THE POSTERIOR FOR DPPs

Now we will express the posterior in the case of DPPs under the following conditions.

4.5 SETTING. Let (Θ, ν) be a measure space and $L(\theta) \in \mathbb{R}_{\text{sym},+}^{N \times N}$ be an elementary kernel for every $\theta \in \Theta$. Further we assume that we have independent realisations A_1, \dots, A_n of a L -ensemble.

Typically the parametrisations $\theta \mapsto L(\theta)$ will be one of the three parametric models in III.2.1, i.e. θ will either be the whole kernel itself, the quality vector, or the log linearity constant of the qualities and $L(\theta)$ the associated elementary kernel.

The independence relation leads to a factorisation of the density and we obtain the following expression for the posterior density

$$f(\theta|A_1, \dots, A_n) \propto f_{\Theta}(\theta) \prod_{i=1}^n f(A_i|\theta) = f_{\Theta}(\theta) \prod_{i=1}^n \frac{\det(L(\theta)_{A_i})}{\det(L(\theta) + I)} \quad (4.2)$$

where we dropped some indices of the density functions in slight abuse of notation.

Unfortunately the normalisation constant

$$\int_{\Theta} f(\theta|A_1, \dots, A_n) \nu(d\theta) = \int_{\Theta} f_{\Theta}(\theta) \prod_{i=1}^n \frac{\det(L(\theta)_{A_i})}{\det(L(\theta) + I)} \nu(d\theta) \quad (4.3)$$

can neither be computed analytically nor numerically in an efficient way. This problem can be solved through the powerful method of Markov chain Monte Carlo simulation that allow to approximate a distribution with only the knowledge of its unnormalised density.

IV.2 Markov chain Monte Carlo methods

The method of Markov chain Monte Carlo (MCMC) simulation arose almost as early as the Monte Carlo¹ simulations itself and since then a rich theory has been established and a broad range of applications have been found. However we can only give a short overview over the basic principles and refer to [Meyn and Tweedie, 2012] for an introduction of Markov chain theory and to [Robert and Casella, 2013] for a survey on (Markov chain) Monte Carlo methods.

We motivated MCMC methods for the approximation of a distribution π under the knowledge of its unnormalised density. In the nutshell the idea is to construct an ergodic Markov chain $(X_n)_{n \in \mathbb{N}}$ with stationary distribution π , i.e. such that one has

$$\hat{\mathbb{P}}_n = \frac{1}{n} \sum_{i=1}^n \delta_{X_i} \xrightarrow{n \rightarrow \infty} \pi$$

almost surely in the weak sense. This Markov chain can then be simulated using Monte Carlo methods and the associated empirical measure $\hat{\mathbb{P}}_n$ will be approximations of π . However to explain this in more detail we need to recapture some notions of Markov chains.

IV.2.1 Reminder on Markov chains

We will provide an extremely short presentation of only those results that we will use to explain the core of MCMC methods. However this will not contain any proofs and hence it can not replace the study of the already mentioned text books.

Let in the following $(\mathcal{X}, \mathcal{B}(\mathcal{X}))$ be a measurable space.

4.6 DEFINITION (MARKOV CHAIN). (i) A *transition kernel* is a function

$$K: \mathcal{X} \times \mathcal{B}(\mathcal{X}) \rightarrow [0, 1]$$

such that

- a) $K(x, \cdot)$ is a probability measure for every $x \in \mathcal{X}$ and
 - b) $K(\cdot, A)$ is measurable for every $A \in \mathcal{B}(\mathcal{X})$.
- (ii) A *Markov chain* with values in \mathcal{X} and transition kernel $K(\cdot, \cdot)$ is a collection $(X_n)_{n \in \mathbb{N}}$ of \mathcal{X} valued random variables such that

$$\mathbb{P}(X_0 \in A_0, \dots, X_n \in A_n) = \int_{A_0} \gamma(dx_0) \int_{A_1} K(x_0, dx_1) \cdots \int_{A_n} K(x_{n-1}, dx_n) \quad (4.4)$$

for all $A_1, \dots, A_n \in \mathcal{B}(\mathcal{X})$ where γ denotes the distribution of X_0 .

We will call γ the *initial* or *starting distribution* of the Markov chain and will denote the distribution of this Markov chain by \mathbb{P}_γ and the expectation with respect to it by $\mathbb{E}_\gamma[\cdot]$. Further an easy application of Kolmogorov's consistency theorem implies that there is a measure \mathbb{P}_γ on the *path space* $\mathcal{X}^{\mathbb{N}}$ that satisfies (4.4) which shows the existence of a Markov chain given a transition kernel K and initial distribution γ (c.f. [Le Gall et al., 2016]). If the initial distribution

¹A legend has it that the name Monte Carlo was given to the work of von Neumann and Ulam by a colleague referring to Ulam's uncle who lost a significant amount of money gambling in the Monte Carlo casino in Monaco.

is deterministic, i.e. $\gamma = \delta_x$ for one $x \in \mathcal{X}$, then we also write \mathbb{P}_x for the distribution of the Markov chain. We close this paragraph by introducing the notation

$$K^n(x, A) := \mathbb{P}_x(X_n \in A)$$

which is consistent with (4.4) for $n = 1$.

IRREDUCIBILITY, RECURRENCE AND EXISTENCE OF STATIONARY DISTRIBUTIONS

From now on we will fix a reference measure μ on \mathcal{X} .

4.7 DEFINITION (IRREDUCIBILITY AND RECURRENCE). (i) We say a Markov chain is μ *irreducible* if for every $A \in \mathcal{B}(\mathcal{X})$ with $\mu(A) > 0$ there is an index $n \in \mathbb{N}$ such that

$$\mathbb{P}_x(X_n \in A) = K^n(x, A) > 0 \quad \text{for all } x \in \mathcal{X}.$$

(ii) A Markov chain $(X_n)_{n \in \mathbb{N}}$ is called *recurrent* if

- a) there is a measure μ on $\mathcal{B}(\mathcal{X})$ such that (X_n) is μ -irreducible and
- b) for every $A \in \mathcal{B}(\mathcal{X})$ with $\mu(A) > 0$ the expected number of visits of A is infinite, i.e.

$$\mathbb{E}_x \left[\left| \{n \in \mathbb{N} \mid X_n \in A\} \right| \right] = \infty \quad \text{for every } x \in A.$$

(iii) A Markov chain is called *Harris recurrent* if it is recurrent and the number of visits is almost surely infinite, i.e. for any $A \in \mathcal{B}(\mathcal{X})$ with $\mu(A) > 0$ we have

$$\mathbb{P}_x \left(\left| \{n \in \mathbb{N} \mid X_n \in A\} \right| = \infty \right) = 1 \quad \text{for every } x \in A.$$

4.8 DEFINITION (STATIONARY DISTRIBUTIONS). Let π be a measure on $\mathcal{B}(\mathcal{X})$. We call π an *invariant* or *stationary distribution* of a Markov chain with kernel K , if X_{n+1} is distributed according to π whenever X_n is distributed according to π . This is equivalent to

$$\pi(A) = \int_{\mathcal{X}} K(x, A) \pi(dx) \quad \text{for all } A \in \mathcal{B}(\mathcal{X}).$$

4.9 THEOREM (EXISTENCE OF STATIONARY DISTRIBUTIONS). *If $(X_n)_{n \in \mathbb{N}}$ is a recurrent Markov chain, there exists an invariant σ -finite measure which is unique up to a multiplicative factor.*

CONVERGENCE TO THE STATIONARY DISTRIBUTION AND ERGODICITY

We will not introduce the notion of periodic and aperiodic Markov chains here, because it would distract us from our actual goal. However we still present the following result that only holds for aperiodic Markov chains and refer to [Meyn and Tweedie, 2012] for further information. The reason why we still present the theorem is that it explains how one can approximately sample from the stationary distribution of a Markov chain, namely it says that the distribution of X_n converges to the invariant distribution.

4.10 THEOREM (CONVERGENCE TO STATIONARY DISTRIBUTION). *Let $(X_n)_{n \in \mathbb{N}}$ be a Harris recurrent and aperiodic Markov chain with stationary distribution π . Let further γ_n be the distribution of X_n , then we have*

$$\|\gamma_n - \pi\|_{TV} \xrightarrow{n \rightarrow \infty} 0$$

non increasing. Here $\|\cdot\|_{TV}$ denotes the total variation of a measure

$$\|\mu\|_{TV} := \sup_{\mathcal{E}} \sum_{E \in \mathcal{E}} |\mu(E)|$$

where the supremum is taken over all finite families of disjoint measurable sets.

4.11 THEOREM (ERGODIC THEOREM). *Let $(X_n)_{n \in \mathbb{N}}$ be a Harris recurrent Markov chain with σ -finite stationary distribution π , then $(X_n)_{n \in \mathbb{N}}$ is ergodic. This means that if*

$$\hat{\mathbb{P}}_n := \frac{1}{n} \sum_{i=1}^n \delta_{X_i}$$

is the empirical measure, we have almost surely have

$$\int_{\mathcal{X}} f(x) \hat{\mathbb{P}}_n(dx) \xrightarrow{n \rightarrow \infty} \int_{\mathcal{X}} f(x) \pi(dx) \quad (4.5)$$

for every π integrable function f .

In the particular case that \mathcal{X} is a topological space and $\mathcal{B}(\mathcal{X})$ is the Borel algebra and if π is a finite measure, we obtain the almost surely weak convergence of $\hat{\mathbb{P}}_n$ towards π . This means that the convergence in (4.5) almost surely holds for all continuous and bounded functions f . This means that $\hat{\mathbb{P}}_n$ are approximations of the invariant distribution in the sense of weak convergence, which is metrisable for example by the Lévy-Prokhorov or the bounded dual Lipschitz metric (c.f. [Dudley, 2010]).

IDEA OF MARKOV CHAIN MONTE CARLO METHODS

The motivation of the study of Markov chain Monte Carlo methods was to approximate the posterior distribution (4.2). The idea is now to construct and then simulate a Markov chain $(X_n)_{n \in \mathbb{N}}$ such that the empirical measures $\hat{\mathbb{P}}_n$ converge to the posterior.

4.12 DEFINITION (MCMC METHODS). *A Markov chain Monte Carlo (MCMC) method for the simulation of a distribution π is any method that produces an ergodic Markov chain $(X_n)_{n \in \mathbb{N}}$ with stationary distribution π .*

In order to achieve this we only have to check the requirements of the ergodic theorem. This means we want to construct a Harris recurrent Markov chain with invariant distribution π and we want to do this without having to compute the normalisation constant (4.3). We will now present the two most common methods to do this which are the Metropolis-Hastings random walk and the method of slice sampling.

IV.2.2 Metropolis-Hastings random walk

The Metropolis-Hastings random walk is maybe the most wide spread MCMC method and certainly one of the oldest. It was actually proposed in the early 1950s from researchers of the American nuclear programme in Los Alamos (c.f. [Metropolis et al., 1953]). First we will touch on the theoretical aspects of this method and follow the presentation in [Robert and Casella, 2013].

4.13 SETTING. Let Θ be a measurable space, μ a measure on that space and $f : \mathcal{X} \rightarrow [0, \infty]$ a function with finite positive integral

$$Z = \int_{\mathcal{X}} f(x) \mu(dx) \in (0, \infty).$$

Our goal is to find a Harris recurrent Markov chain with invariant distribution

$$\pi(A) := \frac{1}{Z} \int_A f(x) \mu(dx).$$

Let further

$$\{f(\cdot|x) \mid x \in \mathcal{X}\}$$

be a family of probability distributions, which we call the *proposal distributions*.

4.14 THE MH RANDOM WALK. Given the first states $X_0 = x_0, \dots, X_n = x_n$ of the Markov, we define X_{n+1} as follows. Let Y be distributed according to $f(\cdot|x_n)d\mu$ and take one realisation y of Y . Then set

$$X_{n+1} := \begin{cases} y & \text{with probability } \rho(x_n, y) \\ x_n & \text{with probability } 1 - \rho(x_n, y) \end{cases}$$

where

$$\rho(x, y) := \min \left\{ \frac{f(y)f(x|y)}{f(x)f(y|x)}, 1 \right\}.$$

and $\frac{a}{0} := \infty$. The first step of the random walk, namely the sampling of y is called the *proposal step* and the second one the *accept-reject step*. In conclusion a single step of the MH random walk can be expressed in the following compact way.

Algorithm 3 A single step of the MH random walk

Input: Current state x_n of the MH random walk

- 1: $y \sim f(\cdot|x_n)d\mu$
 - 2: $a \sim \mathcal{U}([0, 1])$
 - 3: **if** $a \leq \rho(x_n, y)$ **then**
 - 4: $x_{n+1} \leftarrow y$
 - 5: **else**
 - 6: $x_{n+1} \leftarrow x_n$
 - 7: **end if**
 - 8: **return** x_{n+1}
-

To see that the definition above indeed yields a Markov chain we convince ourselves that the transition kernel is given by

$$K(x, A) = \int_A \rho(x, y) f(y|x) \mu(dy) + (1 - m(x)) \delta_x(A)$$

where δ_x is the Dirac measure in x and

$$m(x) = \int_{\mathcal{X}} \rho(x, y) f(y|x) \mu(dy)$$

is the *acceptance probability* of the chain at state x .

4.15 PROPOSITION (STATIONARY DISTRIBUTION). *The probability measure π is a stationary distribution of MH random walk.*

Proof. We have

$$\int_{\mathcal{X}} K(x, A) \pi(dx) = \frac{1}{Z} \int_{\mathcal{X}} \left(\int_A \rho(x, y) f(y|x) \mu(dy) + (1 - m(x)) \delta_x(A) \right) f(x) \mu(dx) \quad (4.6)$$

We note that

$$\rho(x, y) f(y|x) f(x) = \rho(y, x) f(x|y) f(y).$$

Furthermore we can compute

$$\begin{aligned} \frac{1}{Z} \int_{\mathcal{X}} m(x) \delta_x(A) f(x) \mu(dx) &= \frac{1}{Z} \int_{\mathcal{X}} \int_{\mathcal{X}} \rho(x, y) f(y|x) \mu(dy) \delta_x(A) f(x) \mu(dx) \\ &= \frac{1}{Z} \int_A \int_{\mathcal{X}} \rho(x, y) f(y|x) f(x) \mu(dy) \mu(dx) \\ &= \frac{1}{Z} \int_{\mathcal{X}} \int_A \rho(x, y) f(y|x) f(x) \mu(dx) \mu(dy) \\ &= \frac{1}{Z} \int_{\mathcal{X}} \int_A \rho(y, x) f(x|y) \mu(dx) f(y) \mu(dy) \end{aligned}$$

where we used Fubini-Tonelli theorem² in the second to last step. We note that two of the terms in (4.6) cancel out and we obtain

$$\int_{\mathcal{X}} K(x, A) \pi(dx) = \frac{1}{Z} \int_{\mathcal{X}} \delta_x(A) f(x) \mu(dx) = \pi(A).$$

□

²The Fubini-Tonelli theorem states that the order of integration with respect to two σ -additive measures can be swapped, if the integrated function is non negative.

do you need any conditions on the support of the proposals?

Now we are aiming to prove that the MH random walk is Harris recurrent because then the ergodic theorem yields that the empirical measures associated with the Markov chain will actually converge to π . Obviously this is not for all proposal families in general the case, for example we could consider that the proposal distribution $f(\cdot|x)$ is just the Dirac measures in x^3 . Then the MH random walk would never leave its initial position which will typically be a deterministic point. Hence the empirical measures only be the Dirac measure in the starting point and hence not converge towards π .

The first step towards Harris recurrence is to show irreducibility and this will already give us some hints what families of proposal are sensible.

4.16 PROPOSITION (IRREDUCIBILITY). *Assume that the proposal family is strictly positive, i.e.*

$$f(y|x) > 0 \quad \text{for all } x, y \in \mathcal{X}.$$

Then the MH random walk is π irreducible.

Proof. For any measurable set $A \subseteq \mathcal{X}$ with positive measure $\pi(A) > 0$ we have

$$K(x, A) \geq \int_A \rho(x, y) f(y|x) \mu(dy) > 0.$$

To see this, we can assume that this would not hold, but then the integrant has to zero μ almost surely. Since $f(y|x)$ is strictly positive this would imply $\rho(x, y) = 0$ and hence $f(y) = 0$ almost surely with respect to μ . However this is a contradiction to

$$\pi(A) = \int_A f(y) \mu(dy) > 0.$$

□

Now we can formulate the ergodicity for π irreducible MH random walks.

4.17 THEOREM (ERGODICITY OF THE MH RANDOM WALK). *If the MH random walk is π irreducible, then it is also Harris recurrent and hence ergodic.*

Proof. We refer to Lemma 7.3 in [Robert and Casella, 2013] for the proof of Harris recurrency, the ergodicity then follows from the ergodic theorem. □

IMPLEMENTATION OF THE MH RANDOM WALK

So far we have presented the theoretical foundations of the MH random walk and now we want to touch on a few aspect of the simulation process. For this part we shall point the reader towards the very gentle introduction [Robert and Casella, 1999] to the implementation of the MH random walk which also provides coding examples. Further it shall be noted that we will not provide any rigorous results in this section and sometimes use terminology – like empirical correlation – without defining them mathematically. However this is only done if the term is very well established and can easily be found in the literature. We have seen that the empirical measures associated with the MH random walk converge to π under fairly mild assumptions, meaning for almost all choices of proposal distributions. Nevertheless it is mostly the choice of the proposal that determines the speed of this convergence.

Let us assume $\mathcal{X} = \mathbb{R}^d$ and that the reference measure μ is the Lebesgue measure.

³Obviously this is slightly formal, because the Dirac measure can typically not be expressed through a density. However rigorous examples can be constructed similarly.

4.18 CHOOSING A PROPOSAL FAMILY. Usually one would want to choose the proposal such that the expectation of $f(\cdot|x)$ is x . The most common choice of a proposals is a family of normal distributions $f(\cdot|x)$ with expectation x and covariance $\Sigma \in \mathbb{R}^{d \times d}$. This also has the welcome effect that the acceptance ratio takes the easier form

$$\rho(x, y) = \min \left\{ \frac{f(y)}{f(x)}, 1 \right\}.$$

Also since the densities are strictly positive we ensure that the resulting Markov chain is π irreducible.

4.19 ACCEPTANCE RATE, AUTOCORRELATION AND EFFECTIVE SAMPLE SIZE. Once we have agreed to stick to normal densities for the proposal distributions, we still have the freedom to choose the covariance $\Sigma \in \mathbb{R}^{d \times d}$. This determines how far the proposed new values will be away from the current state of the Markov chain. The motivation for an aggressive proposal distribution, i.e. for a high variance would be that this would enable the Markov chain to make bigger steps and hence explore the space \mathcal{X} faster. Also the chain would be more likely to jump between possibly isolated areas of high density. However this could also lead to a very high rejection rate⁴ if the proposed values are often so far away from the current state of the Markov chain that they are in an area of low density. In this case the Markov chain will only ‘visit’ very few distinct points in the space \mathcal{X} which is also very unfavourable. In fact the findings in [Roberts et al., 1997] suggests that an acceptance rate around 25% is desirable in dimension $d \geq 3$ and around 50% for dimension $d = 1, 2$. The connection between the proposal distribution and the acceptance rate is also elaborated in the upcoming example.

The *autocorrelation function* (acf) of a sequence of data points x_0, \dots, x_n captures the estimated correlation between the observations. More precisely $\text{acf}(k)$ gives the empirical correlation⁵ of $(x_0, x_1, \dots, x_{n-k})$ and $(x_k, x_{k+1}, \dots, x_n)$. Lets assume that our data points are generated by a MH random walk. The autocorrelation function determines the correlation of the Markov chain at time l with the Markov chain at time $l + k$. Hence if $\text{acf}(k) < \varepsilon_0$ where $\varepsilon_0 > 0$ is fixed in advance, one can perceive x_0, x_k, x_{2k}, \dots as an independent sequence of realisations – or more precisely an only weakly correlated one. The *effective sample size* is the length m of this new almost uncorrelated sequence $x_0, x_k, x_{2k}, \dots, x_{mk}$. Obviously the effective sample size strongly depends on the choice of ε_0 that incorporates how much correlation one is willing to accept.

We should quickly touch on how the proposal affects the autocorrelation function and hence the effective sample size. Assume we have a very aggressive proposal distribution. Then we will typically have a high rejection rate and hence $x_l = x_{l+k}$ a lot of times meaning that the autocorrelation function will be high. Hence the effective sample size is rather low. On the other hand if the proposal is too conservative the MH random walk will only take very small steps and hence x_{l+k} will still be close to x_l . Therefore the autocorrelation will be high and the effective sample size low. This effect of the proposal can be seen in Figure IV.2.

4.20 EXAMPLE (ONE DIMENSIONAL MH). We follow an examples for a one dimensional MH random walk given in [Robert and Casella, 1999], namely we set

$$f(x) := \sin(x)^2 \cdot \sin(2x)^2 \cdot \exp\left(-\frac{x^2}{2}\right).$$

⁴The term should be rather intuitive; the rejection rate is the relative amount of rejections that occurred in the MH random walk and analogously for the acceptance rate.

⁵This is the correlation of the two empirical measures associated with (x_0, \dots, x_{n-k}) and (x_k, \dots, x_n) .

The goal of this example is to see how different proposal distributions lead to different acceptance rates, a different exploration of the state space $\mathcal{X} = \mathbb{R}$ and different effective sample sizes. In order to achieve this, we run $2 \cdot 10^4$ samples of the MH random walk with starting point $x_0 = 1$ and three different values $\alpha = 0.01, 3, 100$ of the variance of the proposal distributions. Then we plot a histogram including the actual density and the autocorrelation function for all different values. The acceptance rates were approximately 88% for $\alpha = 0.01$, 34% for $\alpha = 3$ and 9% for $\alpha = 100$. The orders of the effective sample sizes for the different values for α are given by

$$\frac{2 \cdot 10^4}{50} = 4 \cdot 10^2, \quad \frac{2 \cdot 10^4}{8} = 2.5 \cdot 10^3 \quad \text{and} \quad \frac{2 \cdot 10^4}{30} \approx 7 \cdot 10^2$$

in the usual ordering.

This simulation illustrates the problem of too aggressive – $\alpha = 100$ – and too conservative – $\alpha = 0.01$ – proposal distributions and shows how this affects the acceptance rate and the effective sample size.

4.21 TUNING THE PROPOSAL. In order to obtain a higher acceptance rate without the loss of generally choosing the variance of the proposal distribution small one can do tune the proposal distribution. This means one adjusts the proposal distribution after a while, let's say after the first 10^3 samples in such a way that one replaces the original covariance matrix Σ by the empirical covariance of the first 10^3 samples. Then one forgets about all the samples so far – they are called usually the *burn in period* – and starts a new MH random walk usually at one of the data points of the burn in period, since they should already indicate where an area of high density is. The reason why this increases the acceptance rate is, that the proposal now only is aggressive in those directions where the density is broadly spread.

is this true?

find a reference for this

IV.2.3 Slice sampling

Slice sampling is a different MCMC method and quite similar to the MH random walk. Nevertheless it has the theoretical benefit that one does not have to define a family of proposal distributions and that the constructed Markov chain is always irreducible. However we will see that at least when one wants to simulate the slice sampling one runs into similar problems. We begin by fixing our frame we will work in.

4.22 SETTING. Let Θ be a measurable space, μ a measure on that space and $f : \mathcal{X} \rightarrow [0, \infty]$ a function with finite integral

$$Z = \int_{\mathcal{X}} f(x) \mu(dx) \in (0, \infty).$$

In particular there is $\hat{x} \in \mathcal{X}$ such that $f(\hat{x}) > 0$. Our goal is to find a Harris recurrent Markov chain with invariant distribution

$$\pi(A) := \frac{1}{Z} \int_A f(x) \mu(dx).$$

Further we will assume – after an eventual modification of f on a μ Null set – that

$$f \leq \|f\|_{L^\infty(\mu)} = \inf \left\{ \alpha \in \mathbb{R} \mid f \leq \alpha \text{ almost surely with respect to } \mu \right\}.$$

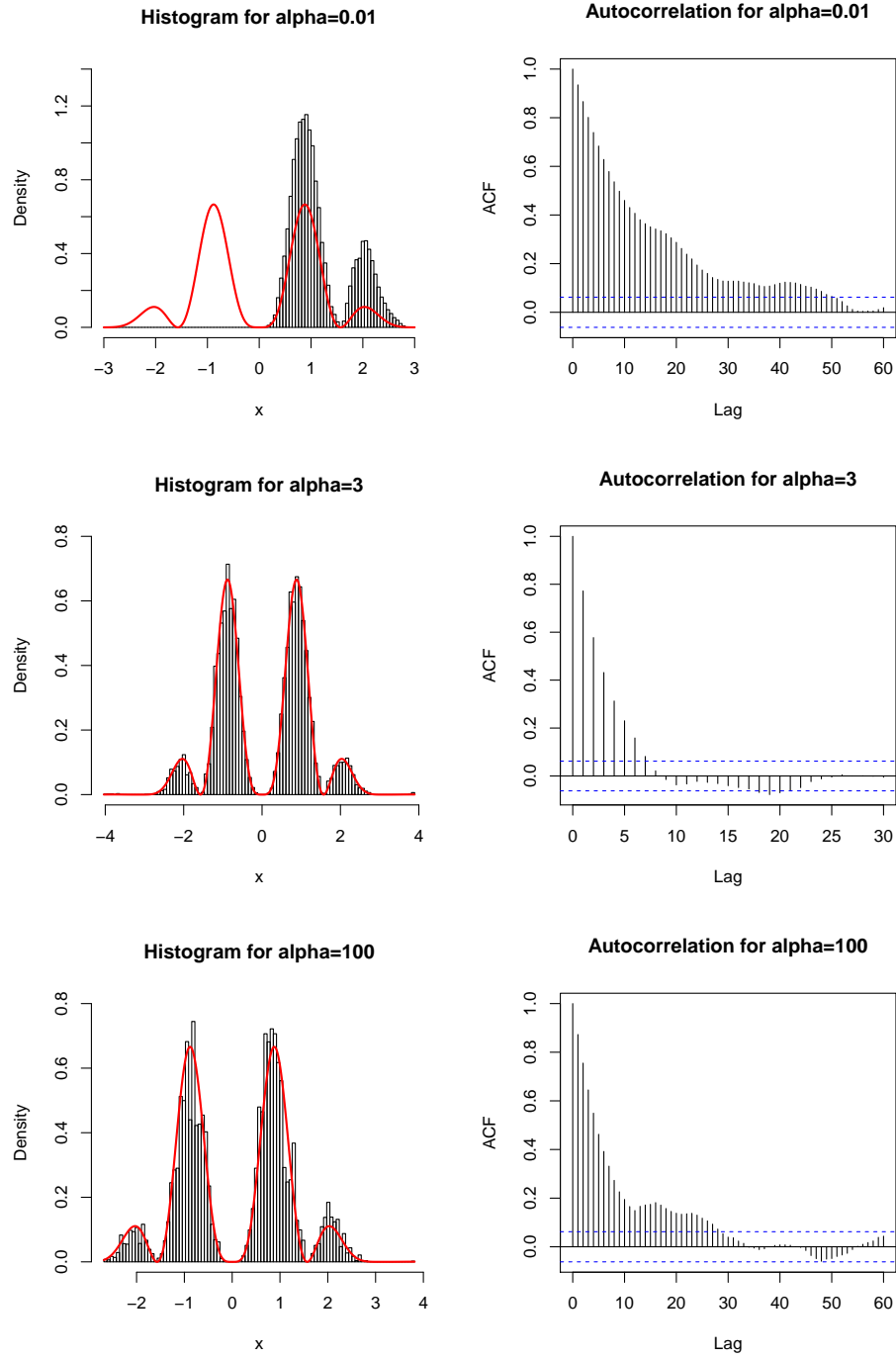


Figure IV.2.: Histograms and autocorrelation functions of for three differently aggressive proposal distributions. It is apparent that the histogram for $\alpha = 3$ fits the actual density the best and also that the autocorrelation decays the quickest for this parameter. Note that for $\alpha = 0.01$ the MH random walk only explored some area of high density. The actual density is obtained by numerical integration.

4.23 DEFINITION OF THE SLICE SAMPLING. Assume we have already given the first n samples x_1, \dots, x_n of the Markov chain. If we have $f(x_n) = 0$, then we set $x_{n+1} := \hat{x}$. Otherwise we sample y according to the uniform distribution on $[0, f(x_n)]$ and define the *slice*

$$S := S(y) := \{x \in \mathcal{X} \mid f(x) \geq y\}.$$

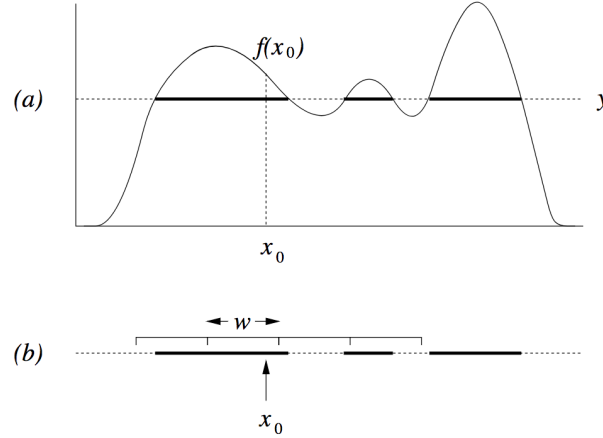


Figure IV.3.: Schematic sketch of the selection of a slice. (a) first y is sampling uniformly on $[0, f(x)]$ and then (b) select the slice. Original graphic from [Neal, 2003].

Note that because $y \leq f(x_n) \leq \|f\|_{L^\infty(\mu)}$ we have $\mu(S) > 0$ as well as

$$\mu(S) \leq y^{-1} \int_S f(x) \mu(dx) < \infty$$

where we used Markov's inequality as well as $y > 0$ almost surely. Now draw x_{n+1} according to the uniform distribution⁶ on S . Note that $f(x_n) > 0$, then $f(x_{n+1}) \geq y > 0$ almost surely, hence $f(x_n) = 0$ can only hold for $n = 0$. Further the reason why we have to treat the case $f(x_n) = 0$ individually is, that there typically is no uniform distribution on the slice $S(0)$. In pseudo code the steps of the resulting Markov chain can be written in the following form.

Algorithm 4 A single slice sampling step

Input: Current state x_n of the Markov chain

- 1: **if then** $f(x_n) = 0$
 - 2: $x_{n+1} \leftarrow \hat{x}$
 - 3: **else**
 - 4: $y \sim \mathcal{U}([0, f(x_n)])$
 - 5: $S \leftarrow \{x \in \mathcal{X} \mid f(x) \geq y\}$
 - 6: $x_{n+1} \sim \mathcal{U}(S)$
 - 7: **end if**
 - 8: **return** x_{n+1}
-

⁶Of course we mean the uniform distribution with respect to μ that gives weight $\mu(S)^{-1} \cdot \mu(A)$ to a set $A \subseteq S$.

If we compare the Markov chain to the MH random walk, we notice that in the slice sampling we first create a random threshold y and then sample uniformly from all points that satisfy this threshold. This is just the other way round than in the MH random walk where we first make a proposal for the next state of the Markov chain and then decide whether we will accept it or not.

Just like in the case of the MH random walk we can explicitly give the transition kernel and use this expression then to check that π is a stationary distribution. The kernel of the Markov chain that arises from the slice sampling iteration is given by

$$\begin{aligned} K(x, A) &= \int_{\mathbb{R}} \frac{\mathbb{1}_{[0, f(x)]}(y)}{f(x)} \cdot \frac{\mu(A \cap S(y))}{\mu(S(y))} \lambda(dy) \\ &= \int_{\mathbb{R}} \frac{\mathbb{1}_{[0, f(x)]}(y)}{f(x)} \cdot Z(y)^{-1} \int_A \mathbb{1}_{[y, \infty)}(f(z)) \mu(dz) \lambda(dy) \end{aligned}$$

where λ is the Lebesgue measure on \mathbb{R} , $\mathbb{1}$ is the indicator function and $Z(y)$ is the normalisation constant

$$Z(y) := \int_{\mathcal{X}} \mathbb{1}_{[y, \infty)}(f(z)) \mu(dz) = \mu(S(y)) \in (0, \infty).$$

Obviously the expression above only holds if $f(x) > 0$ and in the case $f(x) = 0$ we have

$$K(x, A) = \delta_{\hat{x}}(A).$$

4.24 PROPOSITION (INVARIANT DISTRIBUTION). *The probability distribution π is a stationary distribution of the Markov chain associated with the slice sampling method.*

Proof. For any $A \subseteq \mathcal{X}$ we can compute

$$\begin{aligned} \int_{\mathcal{X}} K(x, A) \pi(dx) &= \int_{\mathcal{X}} \int_{\mathbb{R}} \frac{\mathbb{1}_{[0, f(x)]}(y)}{f(x)} \cdot Z(y)^{-1} \int_A \mathbb{1}_{[y, \infty)}(f(z)) \mu(dz) \lambda(dy) f(x) \mu(dx) \\ &= \int_A \int_{\mathbb{R}} Z(y)^{-1} \int_{\mathcal{X}} \mathbb{1}_{[y, \infty)}(f(x)) \mu(dx) \mathbb{1}_{[0, f(z)]}(y) \lambda(dy) \mu(dz) \\ &= \int_A f(z) \mu(dz) = \pi(A) \end{aligned}$$

where we again used Fubini's theorem for non negative functions. \square

4.25 PROPOSITION (IRREDUCIBILITY). *The Markov chain that arises from the slice sampling algorithm is π irreducible.*

Proof. Fix $A \subseteq \mathcal{X}$ with positive probability $\pi(A) > 0$ and $x \in \mathcal{X}$. If we have $f(x) > 0$, then we have $\mu(A \cap S(y)) > 0$ for one $y \in (0, f(x))$. We obtain

$$K(x, A) \geq \int_{\mathbb{R}} \frac{\mathbb{1}_{[0, y]}(z)}{f(x)} \cdot \frac{\mu(A \cap S(z))}{\mu(S(z))} > 0.$$

If however $f(x) = 0$, then we get

$$K^2(x, A) = K(\hat{x}, A) > 0.$$

\square

IMPLEMENTATION DETAILS

Just like in the case of the MH random walk we will provide a few comments about the actual simulation of the slice sampling algorithm. Again those will rather be general guide lines and not be justified by rigorous arguments.

From now on we will assume $\mathcal{X} \subseteq \mathbb{R}^d$. The main difficulty in the implementation is the sampling of a uniform distribution on a slice S . In practice it is even impossible to calculate the slice. Hence one has to come up with a trick. This trick is based on the following observation. Assume that we are able to sample a uniform distribution on a set C that contains the slice S . Then the following algorithm – which is nothing but the conditioning of this uniform distribution on the event that the outcome is in S – samples uniformly from S .

Algorithm 5 Sampling from a slice

```

1:  $x \sim \mathcal{U}(C)$ 
2: while  $x \notin S$  do
3:    $x \sim \mathcal{U}(C)$ 
4: end while
5: return  $x$ 

```

An obvious choice for C would be a cuboid

$$C = \prod_{i=1}^d [a_i, b_i]$$

since it is straight forward to sample a uniform distribution on a cuboid. Namely one only has to sample the individual coordinates uniformly in the intervals $[a_i, b_i]$. The problem still remains how one can find a cuboid that surely contains the whole slice S . The short and frustrating answer is that there is no way to do this. However not everything is lost, since we can use random cuboids that have the property that every part of the slice is contained in the cuboid with positive probability. This will be crucial in retaining the irreducibility of the Markov chain. In fact it been found that in applications the following procedure works well. Given the current state x_n of the Markov chain, we propose a random interval $[a_i, b_i]$ around the i -th component fo x_n . Then we extend those intervals until the endpoints a and b of the cuboid do not lie in the slice anymore. In pseudocode this relates to Algorithm 6. Here $\mathcal{E}(\alpha)$ denotes the exponential distribution with

cite

Algorithm 6 Sampling a random cuboid

Input: Current state x_n of the Markov chain, parameter $\alpha > 0$

```

1: for  $i = 1, \dots, d$  do
2:    $a_i, b_i \sim \mathcal{E}(\alpha)$ 
3: end for
4: while  $x - a \in S$  do
5:    $a \leftarrow 2 \cdot a$ 
6: end while
7: while  $x + b \in S$  do
8:    $b \leftarrow 2 \cdot b$ 
9: end while
10: return  $(x - a, x + b)$ 

```

parameter α and determines how large the first proposed intervals are. Note that it is straight forward and computationally very easy to determine whether a point x is in the slice $S(y)$ since one only has to check $f(x) \geq y$. The reason for the choice of the exponential distribution is that this ensures that the cuboid can get arbitrarily large with possible probability. This leads to the effect that the Markov chain one obtains in exchanging the sample from $\mathcal{U}(S)$ by a sample from $\mathcal{U}(S \cap C)$ still is irreducible.

explain this in greater detail

4.26 THE CHOICE OF α . One could think that a small choice of α – which relates into large values of a_i and b_i – would be the best since this increases the probability that the whole slice S is contained in the cuboid C . There is some truth in this approach, since $\mathcal{U}(S \cap C)$ is a better approximation of $\mathcal{U}(S)$ if C is larger and further the two while loops in Algorithm 6 need more repetitions if a_i and b_i initially are small. This relates into longer running time of the algorithm that samples the random cuboid. However one should not choose α too small, because a large cuboid C also means that a lot of samples from $\mathcal{U}(C)$ will lie outside of $S \cap C$. Hence Algorithm 5 that samples from $\mathcal{U}(S \cap C)$ will get slower as it will need more repetitions of the while loop.

In conclusion there is a trade off – very similar to the case of the MH random walk – between the choice of too small and too large values for α .

make comment on acf and effective sample size

Finally we can present the pseudocode of the algorithm that arises from the combination of the usual slice sampling method and the approximation of the uniform distribution on the slice.

Algorithm 7 Algorithm for the slice sampling

Input: Unnormalised density f , starting value x_0 , desired length n of the chain, $\alpha > 0$

```

1: if  $f(x_0) = 0$  then
2:    $x_0 \leftarrow \hat{x}$ 
3: end if
4: for  $i = 0, \dots, n - 1$  do
5:    $y \sim \mathcal{U}([0, f(x_i)])$ 
6:    $C$  random cuboid around  $x_i$  with parameter  $\alpha$ 
7:    $x \sim \mathcal{U}(C)$ 
8:   while  $x \notin S(y)$  do
9:      $x \sim \mathcal{U}(C)$ 
10:  end while
11:   $x_{i+1} \leftarrow x$ 
12: end for
13: return  $x = (x_0, \dots, x_n)$ 

```

It shall be noted, that the above algorithm also uses a point \hat{x} of positive density, which can be determined easily for a lot of densities f . If this is however not straight forward, one could also sample x_0 according to a normal distribution until we select a point of positive density.

IV.3 Variational MCMC methods

Chapter V

Toy examples and experiments

V.1 Minimal example?

V.2 Points on the line

The first example we present is a selection of points on a (discretised) line. More precisely we will assume that we have 100 points on a line that are equally spaced and we aim to model a spacial repulsion between the selected points. For this we will use the method 2.7 of reference points the diversity features. In this case we will use the set \mathcal{Y} itself as reference set and use a

5.1 SETUP OF THE EXAMPLE. Let $\mathcal{Y} := \{1, \dots, 100\}$ and for $i \in \mathcal{Y}$. Then we will let $\phi_i \in \mathbb{R}^{100}$ be given up to scaling by

$$(\phi_i)_j \propto f\left(\frac{|i-j|}{??}\right)$$

where f is the density of the standard normal distribution. Further we choose the qualities to be constant and so that the expected cardinality is 10.

check

5.2 REMARK. (i) describe scaling including choice of cardinality

(ii) describe rank of the kernel?

(iii) describe choice of 'repulsiveness', plot density around a point; make comment to kernel methods? comment on the qualitative properties of f and why they are suitable here

To make the difference to an uncorrelated point pattern more apparent we also defined a Poisson process, i.e. a DPP without correlations between the points with the same expected cardinality. The sampling results are compared in Figure V.2.

make comment
on zeta function!

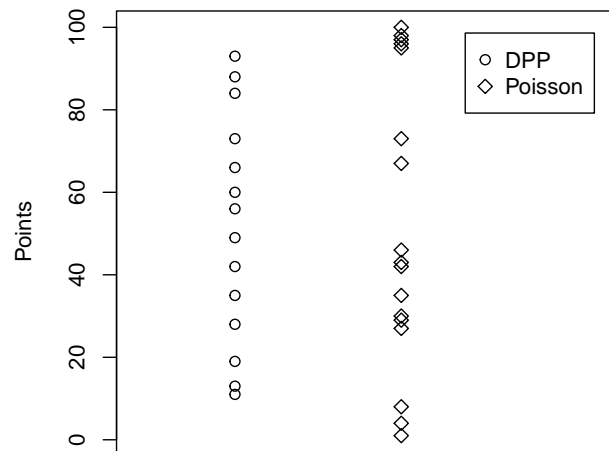


Figure V.1.: Comparison of a DPP with negative correlations on the left and no correlations, i.e. a Poisson point process on the right.

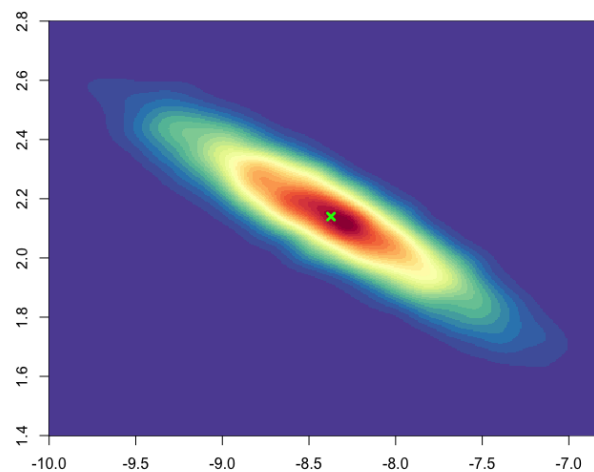


Figure V.2.: Comparison of a DPP with negative correlations on the left and no correlations, i.e. a Poisson point process on the right.

5.3 REPRESENTATION AS BINARY SEQUENCE.

comment on
problems and
findings?

V.3 Points in the square

V.4 Toy example for quality learning

Chapter VI

Summary and conclusion

Chapter A

Calculations

- (i) Complement of DPPs.
- (ii) Normalisation constant of L -ensembles.
- (iii) Cantor's intersection theorem for Hausdorff spaces.
- (iv)
-

Chapter B

Generated code

All my coding was done in R and I will provide the code for sampling, my examples and also the learning algorithm of my toy example here. During my coding I mostly followed Google's R Style Guide (<https://google.github.io/styleguide/Rguide.xml>).

B.1 Sampling algorithm

```
# Implementation of the sampling algorithm as a function

SamplingDPP <- function (lambda, eigenvectors) {
  # First part of the algorithm, doing the selection of the eigenvectors
  N = length(lambda)
  J <- runif(N) <= lambda/(1 + lambda)
  k <- sum(J)
  V <- matrix(eigenvectors[, J], nrow=N)
  Y <- rep(0, k)

  # Second part of the algorithm, the big while loop
  while (k > 0) {
    # Calculating the weights and selecting an item i according to them
    wghts <- k(-1) * rowSums(V2)
    i <- sample(N, 1, prob=wghts)
    Y[k] <- i
    if (k == 1) break

    # Projecting e_i onto the span of V
    help <- V %>% V[i, ]
    help <- sum(help2)(-1/2) * help

    # Projecting the elements of V onto the subspace orthogonal to help
    V <- V - help %>% t(t(V) %>% help)

    # Orthonormalize V and set near zero entries to zero
    V[abs(V) < 10(-9)] <- 0
    j <- 1
    while(j <= k) {
      help2 <- rep(0, N)
      m <- 1
      while (m <= j - 1) {
        help2 <- help2 + sum(V[, j] * V[, m]) * V[, m]
      }
    }
  }
}
```

```

      m <- m + 1
    }
    V[, j] <- V[, j] - help2
    if (sum(V[, j]^2) > 0) {
      V[, j] <- sum(V[, j]^2)^(-1/2) * V[, j]
    }
    j <- j + 1
  }
  V[abs(V) < 10^(-9)] <- 0

  # Selecting a linear independent set in V
  k <- k - 1
  q <- qr(V)
  V <- matrix(V[, q$pivot[seq(k)]], ncol=k)
}
return(Y)
}

```

B.2 Points on the line

```

# NEEDS: sampling algorithm

# In this example we sample points on a (discrete) line according to a DPP
# We model L directly and via the quality-diversity decomposition. We plot and
# compare the patterns to uncorrelated points i.e. to a Poisson point process.

# Minimal example -----
n <- 3
L <- matrix(c(2,1,0,1,2,0,0,0,2), nrow=n)

# Points on a line -----
n <- 100
L <- rep(0, n^2)
for (i in 1:n) {
  for (j in 1:n) {
    L[(i - 1) * n + j] <- dnorm((i-j) * n^(-1/4))
  }
}
L <- matrix(L, nrow=n)

# Modelling phi and q -----
# Points on the line.
m <- 99 # 29
n <- m + 1
q <- rep(10, n) # 0-1 sequences: rep(10^2, n)
phi <- rep(0, n^2)
for (i in 1:n) {
  for (j in 1:n) {
    phi[(i - 1) * n + j] <- dnorm((i - j) / 10) # 0-1 sequences: divide by 2
  }
}
phi <- matrix(phi, ncol=n)

# Log linear quality for the points on the line -----
m <- 99
n <- m + 1

```

```

q <- rep(0, n)
for (i in 1:n) {
  q[i] <- 10^2 * sqrt(m) * exp(-0.2 * abs(i - 50.5))
}
phi <- rep(0, n^2)
for (i in 1:n) {
  for (j in 1:n) {
    phi[(i - 1) * n + j] <- dnorm(2 * (i - j) / sqrt(m))
  }
}
phi <- matrix(phi, ncol=n)

# General part, define L -----
for (i in 1:n) {
  phi[, i] <- sum(phi[, i]^2)^(-1/2) * phi[, i]
}
S <- t(phi) %*% phi
time <- proc.time()
L <- t(q * S) * q
proc.time() - time

# Compute the eigendecomposition, set near zero eigenvalues to zero and
# set up poisson point process with same expected cardinality -----
time <- proc.time()
edc <- eigen(L)
lambda <- edc$values
lambda[lambda < 10^(-9)] <- 0
mean <- sum(lambda / (1 + lambda))
eigenvectors <- edc$vectors
lambda2 <- rep(mean / n / (1 - mean / n), n)
eigenvectors2 <- diag(rep(1, n))
proc.time() - time

# Sample and plot things -----
# Minimal example

# 0-1 sequences
x <- sort(SamplingDPP(lambda, eigenvectors))
as.integer(1:n %in% x)
y <- sort(SamplingDPP(lambda2, eigenvectors2))
as.integer(1:n %in% y)

# Sample from both point processes and plot the points on the line
pointsDPP <- SamplingDPP(lambda, eigenvectors)
pointsPoisson <- SamplingDPP(lambda2, eigenvectors2)
plot(rep(1, length(pointsDPP)), pointsDPP,
      ylim=c(1, n), xlim=c(.4, 3.2), xaxt='n', ylab="Points", xlab="")
points(rep(2, length(pointsPoisson)), pointsPoisson, pch=5)
legend("topright", inset=.05, legend=c("DPP", "Poisson"), pch=c(1, 5))

# Remove all objects apart from functions
rm(list = setdiff(ls(), lsf.str()))

```

B.3 Points in the square

```

# NEEDS: sampling algorithm

```

*# In this example we sample points on a two dimensional grid according to a DPP
 # We model L directly and via the quality-diversity decomposition including
 # different dimensions D for the feature vectors phi. We plot and compare the
 # patterns to uncorrelated points i.e. to a Poisson point process.*

```
# Define the coordinates of a point -----
CoordinatesNew <- function(i, n) {
  y1 <- floor((i - 1) / (n + 1))
  x1 <- i - 1 - (n + 1) * y1
  return (t(matrix(c(x1, y1)/n, nrow=length(i))))
}
DistanceNew <- function (i, j, n, d) {
  return (sqrt(colSums((CoordinatesNew(i, n) - CoordinatesNew(j, d))^2)))
}
```

```
# Direct modelling of L -----
m <- 19
n <- (m + 1)^2
L <- rep(0, n^2)
for (i in 1:n) {
  for (j in 1:n) {
    L[(i - 1) * n + j] = n^2 * dnorm(Distance(i, j, m))
  }
}
L <- matrix(L, nrow=n)
```

```
# Modelling phi and q -----
# Points in the square.
m <- 19
n <- (m + 1)^2
q <- rep(sqrt(m), n)
x <- ceiling(1:n^2 / n)
y <- rep(1:n, n)
time <- proc.time()
phi <- dnorm(sqrt(m) * matrix(DistanceNew(x, y, m, m), n))
proc.time() - time
```

```
# Quality diversity decomposition with small D -----
d <- 25
q <- rep(10^5 * sqrt(m), n)
x <- ceiling(1:(n*d) / d)
y <- rep(1:d, n)
time <- proc.time()
phi <- dnorm(2 * sqrt(m) * matrix(DistanceNew(x, y, m, sqrt(d) - 1), ncol=n))
proc.time() - time
```

```
# Log linear quality for the points in the square -----
m <- 39
n <- (m + 1)^2
q <- exp(-6 * DistanceNew(rep(5, n), 1:n, 2, m) + log(sqrt(m)))
x <- ceiling(1:n^2 / n)
y <- rep(1:n, n)
time <- proc.time()
phi <- dnorm(2 * sqrt(m) * matrix(DistanceNew(x, y, m, m), n))
proc.time() - time
```

```

# General part, defining L -----
# d <- length(phi) / n
for (i in 1:n) {
  phi[, i] <- sum(phi[, i]^2)^(-1/2) * phi[, i]
}
S <- t(phi) %*% phi
# B <- t(phi) * q
time <- proc.time()
L <- t(t(q * S) * q) # B %*% t(B)
proc.time() - time

# Compute the eigendecomposition, set near zero eigenvalues to zero and
# set up poisson point process with same expected cardinality -----
time <- proc.time()
edc <- eigen(L)
lambda <- edc$values
lambda[abs(lambda) < 10^(-9)] <- 0
mean <- sum(lambda / (1 + lambda))
eigenvectors <- edc$vectors
lambda2 <- rep(mean / n / (1 - mean / n), n)
eigenvectors2 <- diag(rep(1, n))
proc.time() - time

# Sample from both point processes and plot the points in the square -----
# par(mfrow = c(1,1))
time <- proc.time()
dataDPP <- sort(SamplingDPP(lambda, eigenvectors))
pointsDPP <- t(CoordinatesNew(dataDPP, m))
plot(pointsDPP, xlim=0:1, ylim=0:1, xlab="", ylab="", xaxt='n', yaxt='n', asp=1)
proc.time() - time
dataPoisson <- sort(SamplingDPP(lambda2, eigenvectors2))
pointsPoisson <- t(CoordinatesNew(dataPoisson, m))
plot(pointsPoisson, xlim=0:1, ylim=0:1, xlab="", ylab="",
      xaxt='n', yaxt='n', asp=1)

# Remove all objects apart from functions
rm(list = setdiff(ls(), lsf.str()))

```

B.4 Toy learning example

```

# NEEDS: Sampling algorithm, declaration of the points in the square
# TODO: Maybe do the gradient descent directly over the representation
# of the gradient

# With this toy example we aim to perform the first learning of paramters
# associated to a kernel of a DPP. More precisely we will generate our own
# data of points on a two dimensional grid with a log linear quality model
# and aim to estimate the log linearity parameter.

# Generation of data
time <- proc.time()
T <- 30
data <- rep(list(0), T)
for (i in 1:T) {
  data[[i]] <- sort(SamplingDPP(lambda, eigenvectors))
}
proc.time() - time

```



```

# Define the quality q, L, the feature sum and the loss in dependency of the
# parameter theta
Quality <- function(theta) {
  return(exp(theta[1] * DistanceNew(rep(5, n), 1:n, 2, m) + theta[2]))
}
LFunction <- function(theta) {
  return(t(t(Quality(theta) * S) * Quality(theta)))
}
Feature <- function(A) {
  # return(sum(DistanceNew(rep(5, length(A)), A, 2, m)))
  return(c(sum(DistanceNew(rep(5, length(A)), A, 2, m)), length(A)))
}
Loss <- function(theta) {
  T <- length(data)
  # Sum this over all data entries
  x <- 0
  for (i in 1:T) {
    A <- data[[i]]
    x <- x + 2 * sum(theta * Feature(A)) + log(det(matrix(S[A, A], length(A))))
  }
  return(- x + T * log(det(diag(rep(1, n)) + LFunction(theta))))
}

# Parameter estimations
time <- proc.time()
sol <- nlm(Loss, c(-3, 0))
proc.time() - time
sol$estimate

# Remove all objects apart from functions
rm(list = setdiff(ls(), lsf.str()))

```

Bibliography

- [Affandi et al., 2014a] Affandi, R. H., Fox, E., Adams, R., and Taskar, B. (2014a). Learning the parameters of determinantal point process kernels. In *International Conference on Machine Learning*, pages 1224–1232.
- [Affandi et al., 2014b] Affandi, R. H., Fox, E., Adams, R., and Taskar, B. (2014b). Learning the parameters of determinantal point process kernels. In Xing, E. P. and Jebara, T., editors, *Proceedings of the 31st International Conference on Machine Learning*, volume 32 of *Proceedings of Machine Learning Research*, pages 1224–1232, Beijing, China. PMLR.
- [Benard and Macchi, 1973] Benard, C. and Macchi, O. (1973). Detection and “emission” processes of quantum particles in a “chaotic state”. *Journal of Mathematical Physics*, 14(2):155–167.
- [Billingsley, 2013] Billingsley, P. (2013). *Convergence of probability measures*. John Wiley & Sons.
- [Bondy and Murty, 2011] Bondy, A. and Murty, U. (2011). *Graph Theory*. Graduate Texts in Mathematics. Springer London.
- [Borodin, 2009] Borodin, A. (2009). Determinantal point processes. *arXiv preprint arXiv:0911.1153*.
- [Boyd and Vandenberghe, 2004] Boyd, S. and Vandenberghe, L. (2004). *Convex optimization*. Cambridge university press.
- [Dudley, 2010] Dudley, R. M. (2010). Distances of probability measures and random variables. In *Selected Works of RM Dudley*, pages 28–37. Springer.
- [Gray, 1990] Gray, R. M. (1990). Entropy and information. In *Entropy and information theory*, pages 21–55. Springer.
- [Griffin and Tsatsomeros, 2006] Griffin, K. and Tsatsomeros, M. J. (2006). Principal minors, part ii: The principal minor assignment problem. *Linear Algebra and its applications*, 419(1):125–171.
- [Hastings, 1970] Hastings, W. K. (1970). Monte carlo sampling methods using markov chains and their applications.
- [Higham, 1990] Higham, N. J. (1990). Exploiting fast matrix multiplication within the level 3 blas. *ACM Transactions on Mathematical Software (TOMS)*, 16(4):352–368.

- [Hough et al., 2006] Hough, J. B., Krishnapur, M., Peres, Y., Virág, B., et al. (2006). Determinantal processes and independence. *Probability surveys*, 3:206–229.
- [Kulesza and Taskar, 2010] Kulesza, A. and Taskar, B. (2010). Structured determinantal point processes. In *Advances in neural information processing systems*, pages 1171–1179.
- [Kulesza et al., 2012] Kulesza, A., Taskar, B., et al. (2012). Determinantal point processes for machine learning. *Foundations and Trends® in Machine Learning*, 5(2–3):123–286.
- [Le Gall et al., 2016] Le Gall, J.-F. et al. (2016). *Brownian motion, martingales, and stochastic calculus*, volume 274. Springer.
- [Lehmann and Casella, 2006] Lehmann, E. L. and Casella, G. (2006). *Theory of point estimation*. Springer Science & Business Media.
- [MacKay, 2003] MacKay, D. J. (2003). *Information theory, inference and learning algorithms*. Cambridge university press.
- [Magen and Zouzias, 2008] Magen, A. and Zouzias, A. (2008). Near optimal dimensionality reductions that preserve volumes. In *Approximation, Randomization and Combinatorial Optimization. Algorithms and Techniques*, pages 523–534. Springer.
- [Martin and England, 2011] Martin, N. F. and England, J. W. (2011). *Mathematical theory of entropy*, volume 12. Cambridge university press.
- [Metropolis et al., 1953] Metropolis, N., Rosenbluth, A. W., Rosenbluth, M. N., Teller, A. H., and Teller, E. (1953). Equation of state calculations by fast computing machines. *The journal of chemical physics*, 21(6):1087–1092.
- [Meyn and Tweedie, 2012] Meyn, S. P. and Tweedie, R. L. (2012). *Markov chains and stochastic stability*. Springer Science & Business Media.
- [Neal, 2003] Neal, R. M. (2003). Slice sampling. *Annals of statistics*, pages 705–741.
- [Newey and McFadden, 1994] Newey, W. K. and McFadden, D. (1994). Large sample estimation and hypothesis testing. *Handbook of econometrics*, 4:2111–2245.
- [Newton and Halley, 1744] Newton, I. and Halley, E. (1744). *Philosophiae naturalis principia mathematica*, volume 62. Jussu Societatis Regiae ac typis Josephi Streater, prostant venales apud Sam. Smith.
- [Rice, 2006] Rice, J. (2006). *Mathematical statistics and data analysis*. Nelson Education.
- [Rising et al., 2015] Rising, J., Kulesza, A., and Taskar, B. (2015). An efficient algorithm for the symmetric principal minor assignment problem. *Linear Algebra and its Applications*, 473:126–144.
- [Robert and Casella, 2013] Robert, C. and Casella, G. (2013). *Monte Carlo statistical methods*. Springer Science & Business Media.
- [Robert and Casella, 1999] Robert, C. P. and Casella, G. (1999). The metropolis-hastings algorithm. In *Monte Carlo Statistical Methods*, pages 231–283. Springer.

- [Roberts et al., 1997] Roberts, G. O., Gelman, A., Gilks, W. R., et al. (1997). Weak convergence and optimal scaling of random walk metropolis algorithms. *The annals of applied probability*, 7(1):110–120.
- [Samuel, 1959] Samuel, A. L. (1959). Some studies in machine learning using the game of checkers. *IBM Journal of research and development*, 3(3):210–229.
- [Tauchen, 1985] Tauchen, G. (1985). Diagnostic testing and evaluation of maximum likelihood models. *Journal of Econometrics*, 30(1-2):415–443.
- [Urschel et al., 2017] Urschel, J., Brunel, V.-E., Moitra, A., and Rigollet, P. (2017). Learning determinantal point processes with moments and cycles. *arXiv preprint arXiv:1703.00539*.
- [Veblen, 1912] Veblen, O. (1912). An application of modular equations in analysis situs. *Annals of Mathematics*, 14(1/4):86–94.
- [Volkenstein, 2009] Volkenstein, M. V. (2009). *Entropy and information*, volume 57. Springer Science & Business Media.