

**Learning with Generalized Negative Dependence:  
Probabilistic Models of Diversity for Machine Learning**

by

Zelda E. Lawson Mariet

S.M., Massachusetts Institute of Technology (2016)

Submitted to the Department of Electrical Engineering and Computer  
Science

in partial fulfillment of the requirements for the degree of  
Doctor of Philosophy in Computer Science and Engineering

at the

MASSACHUSETTS INSTITUTE OF TECHNOLOGY

June 2019

© Massachusetts Institute of Technology 2019. All rights reserved.

**Signature redacted**

Author .....  
.....

Department of Electrical Engineering and Computer Science  
May 23, 2019

**Signature redacted**

Certified by .....  
.....

Suvrit Sra

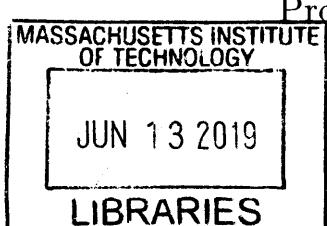
Associate Professor of Electrical Engineering and Computer Science  
Thesis Supervisor

**Signature redacted**

Accepted by .....  
.....

/ \ Leslie A. Kolodziejski

Professor of Electrical Engineering and Computer Science  
Chair, Department Committee on Graduate Students



ARCHIVES



**Learning with Generalized Negative Dependence:  
Probabilistic Models of Diversity for Machine Learning**  
by  
Zelda E. Lawson Mariet

Submitted to the Department of Electrical Engineering and Computer Science  
on May 23, 2019, in partial fulfillment of the  
requirements for the degree of  
Doctor of Philosophy in Computer Science and Engineering

## **Abstract**

This thesis establishes negative dependence as a powerful and computationally efficient framework to analyze machine learning problems that require a theoretical model of diversification. Examples of such problems include experimental design and model compression: subset-selection problems that require carefully balancing the *quality* of each selected element with the *diversity* of the subset as a whole. Negative dependence, which models the behavior of “repelling” random variables, provides a rich mathematical framework for the analysis of such problems.

Leveraging negative dependence theory for machine learning requires (a) scalable sampling and learning algorithms for negatively dependent measures, and (b) negatively dependent measures able to model the specific diversity requirements that arise in machine learning. These problems are the focus of this thesis.

The first part of this thesis develops scalable sampling and learning algorithms for determinantal point processes (DPPs), popular negatively dependent measures with many applications to machine learning. For scalable sampling, we introduce a theoretically-motivated generative deep neural network for DPP-like samples over arbitrary ground sets. To address the learning problem, we show that algorithms for maximum likelihood estimation (MLE) for DPPs are drastically sped up with Kronecker kernels, and that MLE can be further enriched by negative samples.

The second part of this thesis leverages negative dependence for core problems in machine learning. We begin by deriving a generalized form of volume sampling (GVS) based on elementary symmetric polynomials, and prove that the induced measures exhibit strong negative dependence properties. We then show that classical forms of optimal experimental design can be cast as optimization problems based on GVS, for which we derive randomized and greedy algorithms to obtain the associated designs. Finally, we introduce *exponentiated* strongly Rayleigh measures, which allow for simple tuning of the strength of repulsive forces between similar items while still enjoying fast sampling algorithms. The great flexibility of exponentiated strongly Rayleigh measures makes them an ideal tool for machine learning problems that benefit from negative dependence theory.

Thesis Supervisor: Suvrit Sra

Title: Associate Professor of Electrical Engineering and Computer Science

## Acknowledgments

I am profoundly indebted to my advisor, Suvrit Sra, who introduced me to the beauty of DPPs and negative dependence; his advice and guidance (both mathematical and not) were instrumental to every step of my PhD. I also thank my committee, Stefanie Jegelka and Tommi Jaakkola, for their advice and help throughout my grad school years.

My deepest thanks to Ahmed, for his steadfast support through the ups, downs, and sideways of research and life; to Andrea and Jane, the best friends and drinking buddies I could ask for; to Hsin-Yu, amazing roommate and actual angel. Within Sidney-Pacific: thank you to Aashka, Daniel, Dave, Dina, Elan, Eric, Fabián, Greg, Jasmine, Jenny, Jit, Joe, Mirna, Olivia, Nicholas, Rachael and Sydney for being fantastic friends and colleagues throughout officer and trustee life. Thank you also to Marzyeh, Sepehr, Reyu and Clément; to Sheila and Frank for being a family away from home; and to Aude, Flora, Hélène, Estelle, who have stuck with me throughout five years of jetlag and time differences.

Thank you to my labmates of the LogSS and LIS groups, especially Zi and Beomjoon, for many entertaining conversations throughout the years about life, the universe, and matrices, and to my labmate-slash-gym buddy Alex. I am also deeply thankful for the support of Leslie Pack Kaelbling and Tomas Lozano-Perez, who so kindly adopted me and taught me about POMDPs during my first advisor-less semester at MIT.

During my PhD, I was lucky to participate in several great internships and collaborations; thank you to Vitaly Kuznetsov, Alex Kulesza, Jennifer Gillenwater, Sergei Vassilvitskii, Jasper Snoek, Jamie Smith, Yaniv Ovadia, D Sculley, and the rest of the Google MAD science and Brain CAM teams, as well as Mike Gartrell, my collaborator in all things low-rank DPP.

Thank you to my old math teacher, Jean Rouquette, for encouraging my love of math.

And, most of all, thank you to my parents and to my sister Cassandra — without your love and support, this thesis never could have existed. Your love means everything to me.

Այս ա Շեն պահուն  
ա շնչառ Շիշտվյան :

մ Եւրան բակի ի եռ լոն :

# Contents

<b>List of Figures</b>	<b>8</b>
<b>List of Tables</b>	<b>10</b>
<b>1 Introduction</b>	<b>13</b>
1.1 Motivation . . . . .	13
1.2 Thesis outline . . . . .	15
1.2.1 Publications contained in this thesis . . . . .	16
1.2.2 Additional related publications . . . . .	16
1.3 Notation . . . . .	17
<b>2 A theory of negative dependence</b>	<b>19</b>
2.1 Characterizing negative dependence . . . . .	20
2.1.1 Preliminaries: distributions over subsets . . . . .	20
2.1.2 Various characterizations of negative dependence . . . . .	22
2.2 Strongly Rayleigh measures . . . . .	26
2.2.1 Definition and properties . . . . .	27
2.2.2 Operations preserving real stability . . . . .	28
2.2.3 Examples . . . . .	29
2.3 Tractable sampling . . . . .	31
2.3.1 Fast sampling for homogeneous measures . . . . .	31
2.3.2 Generalization to non-homogeneous SR measures . . . . .	32

<b>3 Determinantal point processes — background and sampling</b>	<b>35</b>
3.1 Background . . . . .	36
3.1.1 Key properties of determinantal point processes . . . . .	37
3.1.2 Sampling from DPPs and $k$ -DPPs . . . . .	39
3.2 Modeling DPPs with neural networks . . . . .	41
3.2.1 Generating DPP samples with deep models . . . . .	42
3.2.2 Sampling over fixed and varying ground sets . . . . .	44
3.2.3 Preserving log-submodularity . . . . .	46
3.2.4 Experimental results . . . . .	48
3.3 Open problems . . . . .	53
<b>4 Learning DPPs from data</b>	<b>55</b>
4.1 Learning the DPP kernel . . . . .	56
4.2 Learning DPPs with Kronecker kernels . . . . .	57
4.2.1 Learning with alternating updates . . . . .	59
4.2.2 Learning with joint updates . . . . .	66
4.2.3 The knapsack problem for kernel learning . . . . .	68
4.2.4 Sampling . . . . .	69
4.2.5 Experimental results . . . . .	69
4.3 Learning with corrective negative information . . . . .	72
4.3.1 Contrastive Estimation . . . . .	75
4.3.2 Efficient learning and prediction . . . . .	79
4.3.3 Experiments . . . . .	83
4.4 Open problems . . . . .	88
<b>5 Generalized volume sampling</b>	<b>91</b>
5.1 Background . . . . .	92
5.1.1 Negative dependence properties . . . . .	94
5.2 Optimal experimental design . . . . .	95
5.2.1 Generalized Volume Sampling and optimal design . . . . .	96
5.2.2 Problem definition . . . . .	97

5.2.3	Continuous optimization over elementary symmetric polynomials	98
5.2.4	Finding an optimal design . . . . .	101
5.3	Experiments: optimal design . . . . .	106
5.3.1	Synthetic experiments: optimization comparison . . . . .	106
5.3.2	Real data . . . . .	107
5.4	Generalizing the MVCE problem . . . . .	108
5.4.1	GVS and Minimum Volume Covering Ellipsoids . . . . .	109
5.4.2	Deriving the dual . . . . .	110
5.5	Open problems . . . . .	113
<b>6</b>	<b>Exponentiated strongly Rayleigh measures</b>	<b>115</b>
6.1	Definition and negative dependence properties . . . . .	116
6.2	Sampling from ESR measures . . . . .	120
6.2.1	Approximate samplers for ESR measures . . . . .	122
6.2.2	Mixing time analysis . . . . .	123
6.2.3	Some specific bounds . . . . .	126
6.3	Experiments . . . . .	129
6.3.1	Evaluating mixing time . . . . .	129
6.3.2	Anomaly detection . . . . .	130
6.3.3	E-DPPs for the Nyström method . . . . .	131
6.4	Open problems . . . . .	132
<b>7</b>	<b>Conclusion</b>	<b>135</b>
7.1	High-level summary . . . . .	135
7.2	Open problems . . . . .	136
<b>Bibliography</b>		<b>139</b>
<b>A Appendix</b>		<b>151</b>
A.1	Mathematical tools . . . . .	151
A.1.1	Geodesic convexity . . . . .	151
A.1.2	Power-mean inequality . . . . .	151

A.2 Additional experimental results . . . . .	152
A.2.1 Encoder details for DPPNET . . . . .	152
A.2.2 Amazon Baby registries dataset . . . . .	153
A.2.3 Elementary Symmetric Polynomials . . . . .	153
A.2.4 Exponentiated strongly Rayleigh measures . . . . .	155

# List of Figures

2-1	Venn diagram of implications between negative dependence properties	24
3-1	Volume-based intuition for DPPS . . . . .	38
3-2	Transformer network architecture . . . . .	46
3-3	Sampling on the unit square with a DPPNET . . . . .	50
3-4	Sampling MNIST digits with a DPPNET . . . . .	51
3-5	Results for the Nyström approximation experiments, comparing DPP-NET to the fast MCMC sampling method of Li et al. (2016c) according to root mean squared error (RMSE) and wallclock time. Subsets selected by DPPNET achieve comparable and lower RMSE than a DPP and the MCMC method respectively while being significantly faster. In (c), the relative size of the marker represents the size of the sampled subset. . . . .	53
4-1	Learning a Kronecker-kernel for DPPs on synthetic data . . . . .	70
4-2	Learning a Kronecker-kernel for DPPs on genetic data . . . . .	72
4-3	Learning a DPP with and without negative information on toy data .	73
5-1	Comparison of algorithms for SP-optimal experimental design on synthetic data . . . . .	108
5-2	Sparsity and predictive error for different elementary symmetric polynomials in SP-optimal design (concrete compressive strength dataset)	108
6-1	Anomaly detection with exponentiated DPPs on a toy dataset . . . . .	117

6-2	Mixing and sampling time for exponentiated $k$ -DPPs as a function of the set size $k$	130
6-3	Prediction error on regression datasets using exponentiated DPPs to chose landmarks for kernel reconstruction via the Nyström method	132
A-1	Digits and VAE reconstructions from the MNIST training set	152
A-2	Comparison of algorithms to find ESP-optimal designs using skewed covariance matrices	154
A-3	Comparison of algorithms to find ESP-optimal designs using sparse precision matrices	155
A-4	Numerical upper bound of $r$ -closeness using different proposal distributions to sample from exponentiated DPPs	156
A-5	Empirical mixing time when sampling from an exponentiated DPP as a function of the ground set size	156
A-6	$L_2$ reconstruction error when sampling landmarks for kernel reconstruction via the Nyström method	157
A-7	Frobenius norm reconstruction error when sampling landmarks for kernel reconstruction via the Nyström method	157

# List of Tables

3.1	NLL of DPPNET samples over the unit square for sets of size $k = 20$	49
3.2	NLL of DPPNET samples for MNIST digits	51
3.3	NLLs of DPPNET samples on CelebA and MovieLens	52
4.1	Sampling and learning costs for DPPS	58
4.2	Learning a Kronecker DPP on the Amazon baby registries dataset: log-likelihood comparison	71
4.3	Learning a Kronecker DPP on the GENES dataset: runtime and performance	72
4.4	Comparison of low-rank DPP kernels learned via MLE or Contrastive Estimation on the Amazon baby registries dataset	85
4.5	Comparison of MLE and Contrastive Estimation runtimes to learn a low-rank DPP kernel (Amazon baby registry)	86
4.6	Comparison of low-rank DPP kernels learned via MLE or Contrastive Estimation via MPR and AUC metrics	87
5.1	Runtime comparison of various algorithms to find ESP-optimal designs	107
5.2	Number of common items between ESP-optimal designs	107
5.3	Support size of the continuous relaxation of the ESP-optimal design problem	107
6.1	Outlier detection results when using exponentiated DPPS	131
A.1	Description of the Amazon Baby registries dataset	153

A.2 Support size of the continuous relaxation of the ESP-optimal design problem for skewed covariance design matrices . . . . .	155
---	-----

# Chapter 1

## Introduction

### 1.1 Motivation

Machine learning is part of most tools of modern life: automated personal assistants parse spoken requests for weather updates and reminders, complex recommender systems learn from a wealth of stored experience to suggest books, movies, and products tailored to our tastes, and probabilistic models predict anything from climate patterns to stock prices. Where previously machine learning existed only in the spaces of binary choices (Is this email spam? Does this scan show a cancerous mass?), we now rely on machine learning for complex predictions over large spaces with complicated structures. Often, such predictions involve many *simultaneous* decisions. For example, recommender systems rarely suggest a single option to a user. Rather, the user is presented with 5–10 options they might enjoy and makes the final decision.

Choosing a smaller set of items out of many raises the question of how interactions between chosen items should modify the quality of the selected set as a whole. For a book recommendation task, providing the user with recommendations that belong to different genres increases the probability of user engagement with at least one book. When suggesting websites relevant to the query “bow,” showing options for violins, archery, and ribbons is more likely to produce a relevant result than three websites focused solely on archery.

These examples illustrate a common phenomenon within subset selection: *negative*

*dependence.* At a high level, negative dependence is the behavior exhibited when similarity between items (*e.g.*, similarity between interpretations of the term “bow” for the website query case) decreases their likelihood of being chosen simultaneously.

Such occurrences of negative dependence are particularly common in machine learning model design. To train a classification model, the training data must span all different classes that the model will be analyzing once deployed; the marginal gain of adding a data point from a previously seen class is much smaller compared to adding a point from a new class. To estimate a multidimensional hidden variable, experiments must be different enough to reduce uncertainty across each dimension of the hidden variable; running a set of similar experiments may reduce the time and cost of experiments, but will provide less information.

Given a large set of items, evaluating each item’s quality sequentially comes at a cost that grows linearly in the number  $N$  of items. However, when moving to evaluate *sets* of items, such evaluations will in general require a combinatorial exploration over all  $2^N$  possible sets. As such an exponential-cost exploration is not computationally feasible, simplifying assumptions and heuristics are applied to reduce the dimensionality of the search space. The strongest simplifying assumption models each item selection event as independent from the others — this, of course, cannot capture negative interactions between choices. Less drastic modeling assumptions allow for local, sparse interactions between items. Graphical models are a popular model choice for such settings. Typically, however, graphical models are restricted to modeling *positive* interactions in order to allow for fast sampling and inference.

Recently, specific models for negative dependence have gained prominence in machine learning. In particular, determinantal point processes (DPPs), which tractably model negative interactions in subset-selection, have seen many applications in recommender systems and automatic summarization — see *e.g.*, (Lin et al., 2012; Chao et al., 2015; Zhou et al., 2010; Wilhelm et al., 2018). Determinantal point processes have also been the focus of a significant theoretical analysis (Lyons, 2003; Hough et al., 2006; Borodin, 2009; Kulesza and Taskar, 2012; Decreusefond et al., 2015; Lavancier et al., 2015), which in turn has yielded an elegant theory of diverse subset selection.

However, the spectrum of negatively dependent measures is much broader than only determinantal point processes. In this thesis, we show that determinantal point processes and the broader class of negatively dependence measures have a key role to play in machine learning model design and optimization.

## 1.2 Thesis outline

The aim of this thesis is to motivate and develop the use of negatively dependent measures for problems in machine learning that require a model of diversity. In order to accomplish this goal, we must address two separate problems. First, we must scale up algorithms for negatively dependent measures, in order to broaden their applicability (Chapters 3 and 4). Then, we must generalize and develop the theory of negative dependence within the framework of machine learning; this will allow us to obtain fine-grained models for the negative dependence phenomena specific to machine learning (Chapters 5 and 6).

Specifically, the outline of this thesis is the following. Chapter 2 introduces the rich theory of negative dependence derived by Borcea et al. (2009), focusing on strongly Rayleigh (SR) measures and relevant results we will require for the rest of this thesis.

We then focus on the most popular negatively dependent measure: determinantal point processes. Chapters 3 and 4 present results for efficient sampling and learning algorithms for DPPs, allowing for a broader applications of DPPs to machine learning. In Chapter 3, we introduce deep generative networks to sample from DPPs with variable kernels. Chapter 4 focuses on the learning problem; we first introduce Kronecker-structured kernels for DPPs, which allow for efficient learning methods, then consider a novel learning paradigm that augments the standard Maximum Likelihood Estimation problem for DPPs with corrective negative samples.

In Chapter 5, we analyze Generalized Volume Sampling measures, *i.e.*, measures that extend volume sampling by leveraging the family of elementary symmetric polynomials. We identify key connections between generalized volume sampling and the Minimum Volume Covering Ellipsoid problem, a well-known and deeply studied op-

timization problem. We also show that standard models for optimal experimental design can be re-framed in the larger context of generalized volume sampling, and develop randomized and greedy algorithms to compute the associated designs.

Finally, Chapter 6 extends the class of strongly Rayleigh measures to the super-class of *exponentiated* strongly Rayleigh measures. We show that these measures inherit some of the negative dependence properties of SR measures, and allow for intuitive and efficient tuning of the strength of the similarity repulsion process. Experimental results on outlier detection and kernel reconstruction tasks confirm that the greater flexibility of exponentiated strongly Rayleigh measures makes them a valuable tool to approach fundamental problems in machine learning.

### 1.2.1 Publications contained in this thesis

This thesis covers results first derived in the following publications:

- Z. Mariet, Y. Ovadia, and J. Snoek. **DPPNET**: Approximating Determinantal Point Processes with deep networks. *Submitted to NeurIPS*, 2019b *Chapter 3*
- Z. Mariet and S. Sra. Kronecker Determinantal Point Processes. In *Advances in Neural Information Processing Systems*, 2016b *Chapter 4*
- Z. Mariet, M. Gartrell, and S. Sra. Learning Determinantal Point Processes by sampling inferred negatives. In *Proc. Int. Conference on Artificial Intelligence and Statistics (AISTATS)*, 2019a *Chapter 4*
- Z. Mariet and S. Sra. Elementary symmetric polynomials for optimal experimental design. In *Advances in Neural Information Processing Systems*, 2017 *Chapter 5*
- Z. Mariet, S. Sra, and S. Jegelka. Exponentiated Strongly Rayleigh distributions. In *Advances in Neural Information Processing Systems*, 2018 *Chapter 6*

### 1.2.2 Additional related publications

The following papers contain results on negative dependence and its applications to machine learning, but are not covered in detail in this thesis.

- Z. Mariet and S. Sra. Fixed-point algorithms for learning Determinantal Point Processes. In *Proc. Int. Conference on Machine Learning (ICML)*, 2015
- Z. Mariet and S. Sra. Diversity networks: Neural network compression using Determinantal Point Processes. *Int. Conf. on Learning Representations (ICLR)*, 2016a
- J. A. Gillenwater, A. Kulesza, Z. Mariet, and S. Vassilvitskii. Maximizing Induced Cardinality Under a Determinantal Point Process. In *Advances in Neural Information Processing Systems*, 2018
- J. A. Gillenwater, A. Kulesza, Z. Mariet, and S. Vassilvitskii. Sublinear Sampling for Determinantal Point Processes. *Proc. Int. Conference on Machine Learning (ICML)*, 2019

### 1.3 Notation

We consider measures  $\mu$  over a finite ground set  $\mathcal{Y}$  of  $N$  items; without loss of generality we set  $\mathcal{Y} = \{1, \dots, N\}$  and also use the notation  $[N] \triangleq \{1, \dots, N\}$ . We denote by  $2^{\mathcal{Y}} = 2^{[N]}$  the powerset of  $\mathcal{Y}$ . Capital letters such as  $S$  and  $T$  denote subsets of  $[N]$ ; we distinguish matrices from sets by writing them in bold, *e.g.*,  $\mathbf{L}$  and  $\mathbf{M}$ . We write  $\mathbf{0}$  (resp.  $\mathbf{1}$ ) for the vector of all zeros (resp. all ones), where the dimension is clear from context.

We often consider submatrices indexed by subsets of  $\mathcal{Y}[N]$ . For  $S, T \subseteq [N]$  and  $\mathbf{M} \in \mathbb{R}^{N \times N}$ , we write  $\mathbf{M}_{S,T} := [\mathbf{M}_{ij}]_{i \in S, j \in T}$  for the  $|S| \times |T|$  submatrix of  $\mathbf{M}$ . To lighten the notation, we write  $\mathbf{M}_{:,T} = \mathbf{M}_{[N],T}$  and  $\mathbf{M}_{S,:} = \mathbf{M}_{S,[N]}$ . For principal submatrices, we additionally abbreviate  $\mathbf{M}_{S,S}$  as  $\mathbf{M}_S$ .

Finally, we write  $\mathbb{S}_N^+$  (resp.  $\mathbb{S}_N^{++}$ ) for the set of positive semi-definite (resp. positive definite) matrices of size  $N \times N$ , and denote by  $\succeq$  the usual Loewner order on  $\mathbb{S}_N^+$ : for  $\mathbf{A}, \mathbf{B} \in \mathbb{S}_N^+$ ,  $\mathbf{A} \succeq \mathbf{B} \iff (\mathbf{A} - \mathbf{B}) \in \mathbb{S}_N^+$ .



# Chapter 2

## A theory of negative dependence

Negative dependence characterizes the behavior of probability distributions over subsets of items. At a high level, events are negatively dependent if they are unlikely to occur simultaneously. As a simple example, consider  $N$  urns within which we have dropped a total of  $k$  balls in an *i.i.d.* fashion. Intuitively, the existence of an empty urn decreases the probability of other urns being empty; and in fact, one can verify that the events “urn  $i$  is empty” and “urn  $j$  is empty” are negatively dependent.

The necessity of a theoretically grounded model of negative dependence was recognized by Pemantle (2000), wherein Pemantle motivated the need for such theory and provided arguments suggesting its existence. However, the theory of negative dependence proved to be much more difficult to develop than the theory of positive dependence. Many different characterizations of negative dependence were suggested, yet identifying the interactions and overlaps between these characterizations turned out to be surprisingly complex.

A partial solution was offered by Borcea and Brändén (2009), who introduced the class of *strongly Rayleigh* measures: a class of distributions over subsets that inherits all previous characterizations of negative dependence. Since this contribution, strongly Rayleigh measures (and the related theory of real stable polynomials) have been fundamental to many areas of mathematics. Most famously, they were crucial to proving the Kadison-Singer conjecture (Marcus et al., 2015); they also served to prove a variant of the Kadison-Singer problem introduced in (Anari and Gharan, 2014), as

well as the Monotone Column Permanent conjecture (Brändén et al., 2011).

More recently, strongly Rayleigh measures have emerged as an elegant class well-suited to modeling many problems in machine learning model design and optimization. As this thesis aims to develop and apply the theory of negative dependence for machine learning, we begin with an overview of the theory of strongly Rayleigh distributions, focusing particularly on the properties that make them relevant for machine learning applications.

## 2.1 Characterizing negative dependence

Mathematically, negative dependence is defined by considering distributions over subsets of a ground set; in this thesis, we consider discrete ground sets  $\mathcal{Y}$ , which we identify with  $[N]$ . By definition, the probability measures  $\mu$  we are interested in verify  $\mu : 2^{[N]} \rightarrow [0, 1]$  and  $\sum_{S \subseteq [N]} \mu(S) = 1$ .

### 2.1.1 Preliminaries: distributions over subsets

Denote by  $X_i$  the binary random variable over  $2^{[N]}$  defined for all subsets  $S \subseteq [N]$  as  $X_i(S) = 1$ , if  $i \in S$ , and 0, otherwise. Equivalently,  $X_i$  is the  $i$ -th coordinate function. The indicator function  $\chi_S$  of any subset  $S$ , defined as  $\chi_S(T) \neq 0$  if and only if  $S = T$ , is given by

$$\chi_S(T) = \prod_{i \in S} X_i(T) \prod_{j \notin S} (1 - X_j(T)).$$

As  $\chi_S$  is a multi-affine polynomial<sup>1</sup> in the  $X_i$ , it follows that any function  $f$  over  $2^{[N]}$  can be written as a multi-affine polynomial in the  $X_i$ : first by writing  $f(S) = \sum_{S \subseteq [N]} f(S)\chi_S$  and then by reordering the variables to obtain a formulation  $f(S) = \tilde{f}(X_1(S), \dots, X_N(S)) = \sum_{S \subseteq [N]} a_S \prod_{i \in S} X_i$ .

As a consequence, any distribution  $\mu$  over subsets of  $[N]$  can be equivalently represented by a multi-affine polynomial with non-negative coefficients in  $N$  variables; this polynomial is called the *generating polynomial* of  $\mu$ .

---

<sup>1</sup>Recall that a multivariate polynomial is multi-affine if it has degree at most one in each variable.

**Definition 2.1.1** (Generating polynomial). Let  $\mu$  be a probability distribution over  $2^{[N]}$ . The generating polynomial of  $\mu$ , denoted by  $g_\mu$ , is the multi-affine polynomial over  $\mathbb{C}^N$  defined by

$$g_\mu(z_1, \dots, z_N) = \sum_{Y \subseteq [N]} \mu(Y) z^Y,$$

where we write  $z^Y \triangleq \prod_{y \in Y} z_y$ .

Following (Borcea and Brändén, 2009), when considering generating polynomials, we will use  $X_i$  for random variables and use lower-case  $x, y, z$  for variables in  $\mathbb{C}^n$ .

**Remark 2.1.1.** By construction,  $g_\mu(\mathbf{1}) = 1$ .

**Proposition 2.1.1.** *There is a one-to-one mapping between the set of distributions over  $2^{[N]}$  and the set  $P_N[z_1, \dots, z_N]$  of multi-affine polynomials  $g$  in  $N$  variables with non-negative coefficients such that  $g(\mathbf{1}) = 1$ .*

*Proof.* Def. 2.1.1 provides an injective mapping from distributions over  $2^{[N]}$  to the polynomials in  $P_N[z_1, \dots, z_N]$ . Conversely, any multi-affine polynomial  $g$  with non-negative coefficients such that  $g(\mathbf{1}) = 1$  defines a distribution over  $2^{[N]}$ . To see this, write  $g(z) = \sum_{S \subseteq [N]} a_S z^S$  and define distribution  $\mu_g$  by setting  $\mu(S) = a_S \geq 0$  for all  $S \subseteq [N]$ . The condition  $g(\mathbf{1}) = 1$  guarantees that  $\mu$  is normalized.  $\square$

Most standard operations on distributions  $\mu$  over  $2^{[N]}$  translate directly into operations over the corresponding generating polynomial  $g_\mu$ . In particular, the probability  $\mu(S)$  and the marginal probability  $\Pr_\mu(T \subseteq S)$  of observing  $T$  when sampling  $S \sim \mu$  are given by the following operations:

$$\mu(S) = \partial^S g_\mu(\mathbf{0}) \quad \text{and} \quad \Pr_\mu(T \subseteq S) = \partial^T g_\mu(\mathbf{1}),$$

where we write  $\partial^S$  for the differential operator  $\prod_{i \in S} \partial/\partial z_i$ .

We will often need to consider distributions over subsets that only assign non-zero probabilities to subsets of a given size  $k$ ; for this purpose, we recall here the definition of homogeneous polynomials and measures.

**Definition 2.1.2** (Homogeneous polynomials and measures). A polynomial  $g \in \mathbb{R}[z_1, \dots, z_N]$  is homogeneous of degree  $k$  if for all  $t \in \mathbb{C}$  and  $z \in \mathbb{C}^N$ ,  $g(tz) = t^k g(z)$ .

Similarly, a measure  $\mu$  over  $2^{[N]}$  is homogeneous of degree  $k$  if its generating polynomial is homogeneous of degree  $k$ , *i.e.*, if  $\mu(S) \neq 0 \implies |S| = k$ .

### 2.1.2 Various characterizations of negative dependence

Deriving the theory for negative dependence is a complex endeavor; as pointed out by Pemantle (2000), given to comparative simplicity of obtaining a theory of *positive* dependence, characterizing negative dependence is surprisingly difficult.

A weak characterization of negative dependence for a measure  $\mu$  over  $2^{[N]}$  is given by the pairwise negative correlation (p-NC) condition:  $\mu$  is pairwise negatively correlated if for all  $i, j \in [N]$ ,

$$\Pr_\mu(X_i)\Pr_\mu(X_j) \geq \Pr_\mu(X_i X_j), \quad (\text{p-NC})$$

where we recall that  $X_i$  is the binary random variable indicating  $i \in S$ . In terms of generating polynomials, (p-NC) is equivalent to

$$\partial_i g_\mu(\mathbf{1}) \partial_j g_\mu(\mathbf{1}) \geq \partial_i \partial_j g_\mu(\mathbf{1}).$$

The p-NC inequality simply states that elements  $i$  and  $j$  are more likely to be sampled independently under  $\mu$  than simultaneously; this amounts to the weakest form of negative dependence.

A second characterization of negative dependence is given by the *negative lattice condition*, which requires for all  $S, T \subseteq [N]$  that

$$\mu(S)\mu(T) \geq \mu(S \cup T)\mu(S \cap T). \quad (\text{NLC})$$

In terms of  $\mu$ 's generating polynomial, (NLC) is equivalently stated as

$$\partial^S g_\mu(\mathbf{0}) \partial^T g_\mu(\mathbf{0}) \geq \partial^{S \cup T} g_\mu(\mathbf{0}) \partial^{S \cap T} g_\mu(\mathbf{0}).$$

The NLC condition is equivalent to the log-submodularity of  $\mu$ .

**Definition 2.1.3.** A function  $f$  over  $2^{[N]}$  is *submodular* if it verifies for every  $S, T \subseteq [N]$ ,  $f(S) + f(T) \geq f(S \cup T) + f(S \cap T)$ . A function  $f$  over  $2^{[N]}$  is said to be *log-submodular* if  $\log(f)$  is well defined and submodular.

**Proposition 2.1.2.** *A distribution  $\mu$  over  $2^{[N]}$  verifies the NLC condition if and only if  $\mu$  is log-submodular.*

Finally, a global characterization of negative dependence is given by the negative association (NA) property, stated as

$$\int F d\mu \int G d\mu \geq \int FG d\mu, \quad (\text{NA})$$

for all pairs of increasing functions  $F$  and  $G$  over  $2^{[N]}$  where, crucially,  $F$  and  $G$  must depend on a *disjoint*<sup>2</sup> set of coordinates. Requiring that  $F$  and  $G$  depend on disjoint sets of coordinates is necessary to avoid trivial correlations such as the self-correlation between  $F$  and  $F$  that follows from the Cauchy-Schwarz inequality.

As alluded to earlier, characterizing negative dependence is more difficult than characterizing positive dependence. As a first complication, there is no straightforward chain of implications between the previous inequalities. The positive dependence equivalents to p-NC, NLC and NA (which are obtained by reversing the inequalities in all negative dependence characterizations, and which we refer to as p-PC, PLC and PA) are linked by the implications

$$\text{PA} \iff \text{PLC} \implies \text{p-PC},$$

where the first equivalence is given by the FKG theorem (Fortuin et al., 1971). However, there is no such chain of implications for negative dependence. Indeed, NLC does not imply negative association, nor is p-NC implied by the NLC condition.

---

<sup>2</sup>See (Borcea and Brändén, 2009) for a formal definition of set functions with disjoint sets of coordinates.

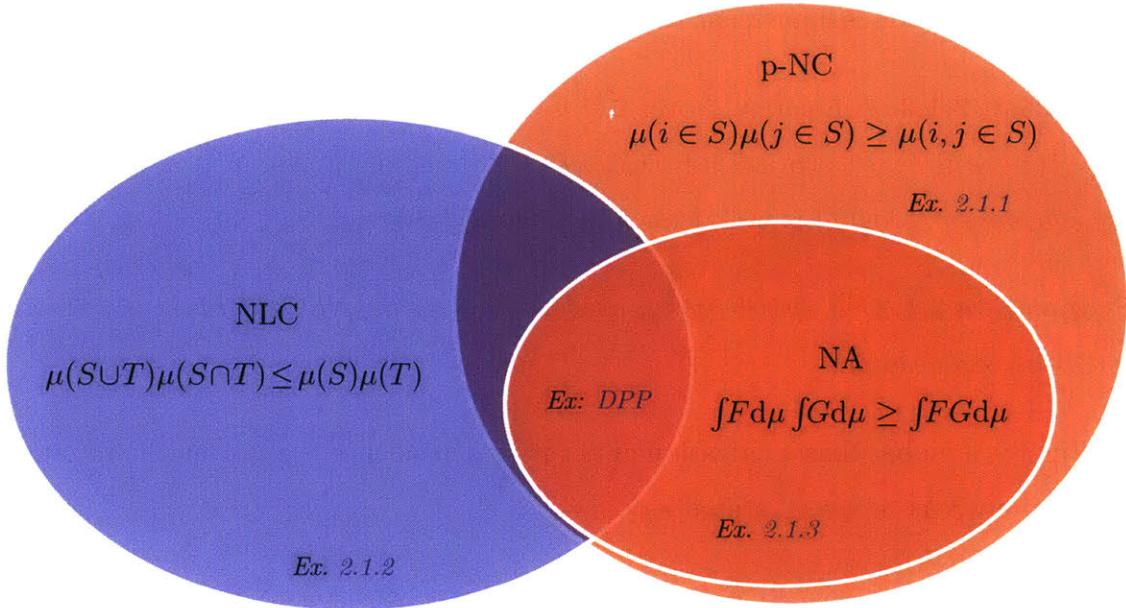


Figure 2-1: Venn diagram of implications between negative dependence properties

**Example 2.1.1** ( $p\text{-NC} \not\Rightarrow NLC, NA$ ). The  $p\text{-NC}$  property is the weakest characterization of negative dependence; the following measure  $\mu$  over  $\{1, 2, 3\}$  is easily verified to be a  $p\text{-NC}$  measure that satisfies neither the  $NA$  nor  $NLC$  conditions:

$$\mu(1) = \mu(2) = \mu(3) = 0.3 \quad \mu(123) = 0.1 \quad (\text{all other probabilities being zero}).$$

**Example 2.1.2** ( $NLC \not\Rightarrow p\text{-NC}, NA$  (Borcea and Brändén, 2009)). The measure  $\mu$  over subsets of  $\{1, \dots, 4\}$  such that  $\mu(\{1, 2\}) = \mu(\{3, 4\}) = \frac{1}{2}$  trivially satisfies the  $NLC$  condition but  $\mu(X_1 X_2) = \frac{1}{2}$  and  $\mu(X_1)\mu(X_2) = \frac{1}{4}$ . Hence,  $\mu$  is not  $p\text{-NC}$ . Since  $NA \Rightarrow p\text{-NC}$ , it immediately follows that  $\mu$  also does not satisfy  $NA$ .

**Example 2.1.3** ( $NA \not\Rightarrow NLC$  (Borcea and Brändén, 2009)). Let  $e_k$  be the  $k$ -th elementary symmetric polynomial in  $n$  variables ( $e_k(z) = \sum_{|S|=k} z^S$ ); the measure  $\mu$ , with generating polynomial

$$g_\mu(z_1, z_2, z_3) \propto \frac{3}{2} + z_1 + z_2 + z_3 + z_1 z_2 + z_2 z_3 + z_1 z_3,$$

is  $NA$  but not  $NLC$ .

Finding a measure that simultaneously verifies the p-NC and NLC conditions while not being NA is more difficult; however, for  $N \geq 4$  certain exponentiated DPPs (Chapter 6), which by construction satisfy the NLC inequality, can be shown to satisfy p-NC but break NA numerically.

A second setback to deriving a theory of negative dependence arises from the fact that the NA and NLC properties are not closed under desirable operations over probability distributions. Introduced in (Pemantle, 2000), these desirable closure properties for a characterization  $\mathcal{P}$  of negative dependence are the following: for any measure  $\mu$  over  $2^{[N]}$  verifying  $\mathcal{P}$ ,

- *Closure under projection:* given  $A \subseteq [N]$ , the projection of  $\mu$  onto  $\{0, 1\}^A$  (i.e., the measure obtained by integrating out the variables in  $[N] \setminus A$ ) should satisfy  $\mathcal{P}$ . Equivalently, the measure with generating polynomial  $g_\mu|_{z_i=1, i \in [N] \setminus A}$  should verify  $\mathcal{P}$ .
- *Closure under conditioning:* given  $A \subseteq [N]$  and  $\eta \in \{0, 1\}^A$ , the conditioned measure  $\mu(\cdot | X_i = \eta(i), i \in A)$  should verify  $\mathcal{P}$ ; equivalently, the measure with generating polynomial  $h$  in  $N - |A|$  variables defined by  $h(z_{i \notin A}) := g_\mu(z')$ , where  $z'_i = z_i$  if  $i \notin A$  and  $z_i = \eta(i)$  otherwise, should verify  $\mathcal{P}$ .
- *Closure under external fields*<sup>3</sup>: for a vector  $a \in \mathbb{R}_+^N$ , the measure with generating polynomial

$$(z_1, \dots, z_N) \rightarrow \frac{g_\mu(a_1 z_1, \dots, a_N z_N)}{g_\mu(a_1, \dots, a_N)}$$

should verify  $\mathcal{P}$ .

Neither NA nor NLC satisfy these closure properties; for this reason, the following definitions are required.

**Definition 2.1.4** ([Strong] hereditary negative lattice condition). A measure  $\mu$  over  $2^{[N]}$  satisfies the *hereditary negative lattice condition* (h-NLC) if every projection of  $\mu$  satisfies NLC. Furthermore, a measure  $\mu$  over  $2^{[N]}$  satisfies the *strong hereditary*

---

<sup>3</sup>This name is borrowed from the theory of Ising models.

*negative lattice condition* (h-NLC+) if any measure obtained from  $\mu$  by imposing an external field satisfies h-NLC.

**Definition 2.1.5** ([Strongly] conditionally negative associated measures). A measure  $\mu$  over  $2^{[N]}$  is *conditionally negatively associated* (CNA) if any measure obtained by conditioning  $\mu$  on some of the values of the variables is NA. Similarly, a measure  $\mu$  is *strongly conditionally negatively associated* (CNA+) if any measure obtained by imposing external fields and projections on  $\mu$  is CNA.

**Proposition 2.1.3** (Pemantle (2000)). *We have the following implications: CNA  $\implies$  h-NLC, and CNA+  $\implies$  h-NLC+. Additionally, all of h-NLC, h-NLC+, NA, CNA, and CNA+ imply p-NC.*

## 2.2 Strongly Rayleigh measures

The crucial development in deriving a theory of negative dependence was derived by Borcea and Brändén (2009) through the introduction of the class of *strongly Rayleigh* (SR) measures. SR measures are characterized via the roots of their generating polynomials; for this reason, we begin by reproducing a few related results in multivariate polynomial theory below.

**Definition 2.2.1** (Stable polynomial). A polynomial  $g \in \mathbb{C}[z_1, \dots, z_N]$  is *stable* if

$$\Im(z_i) > 0 \text{ for all } 1 \leq i \leq n \implies g(z_1, \dots, z_N) \neq 0.$$

If  $g$  is stable and has all real coefficients,  $g$  is called *real stable*.

Real stable polynomials can be viewed as a generalization of real-rooted univariate polynomials to multivariate polynomials.

**Proposition 2.2.1.** *Let  $g$  be a real stable polynomial in  $N$  variables and consider the univariate polynomial  $\tilde{g} : z \rightarrow g(z, \dots, z)$ ; then  $\tilde{g}$  is real-rooted.*

*Proof.* By contradiction, if  $\tilde{g}$  has a root  $z = a + ib$  with  $b \neq 0$ , one easily verifies that  $a - ib$  is also a root for  $\tilde{g}$ . Hence,  $\tilde{g}$  has a root  $z$  with  $\Im(z) > 0$ , and so  $z = (z, \dots, z)$  is a root of  $g$  in the positive orthant  $\{z \in \mathbb{C}^n \mid \forall i, \Im(z_i) > 0\}$ :  $g$  is not real stable.  $\square$

Multivariate stable polynomials have fundamental connections to many objects in polynomial theory and optimization, including hyperbolic polynomials and Lee-Yang polynomials; we refer the interested reader to works such as (Gårding, 1959; Güler, 1997; Choe et al., 2004; Borcea and Brändén, 2008; Brändén, 2007; Borcea and Brändén, 2009, 2010). They also admit the following characterization:

**Proposition 2.2.2.** *A polynomial  $g \in \mathbb{R}[z_1, \dots, z_n]$  is real stable if and only if for all  $x \in \mathbb{R}_+^n$  and any  $e \in \mathbb{R}^n$ , the univariate polynomial  $t \rightarrow g(e + tx)$  is real rooted.*

### 2.2.1 Definition and properties

Given the concept of stable polynomials, defining the class of strongly Rayleigh measure is straightforward: strongly Rayleigh measures are measures with real stable generating polynomials.

**Definition 2.2.2** (Strongly Rayleigh measure). A measure  $\mu$  over  $2^{[N]}$  is strongly Rayleigh if its generating polynomial  $g_\mu$  is real stable.

**Theorem 2.2.3** (Borcea and Brändén (2009)). *If a measure  $\mu$  over  $2^{[N]}$  is strongly Rayleigh,  $\mu$  is CNA+. It follows from Proposition 2.1.3 that  $\mu$  also verifies the h-NLC+ condition and is pairwise negatively correlated.*

Proving the real-stability of a polynomial via the location of its roots is often difficult; fortunately, an alternate characterization of real stability for multi-affine polynomials (and hence for strongly Rayleigh measures) is given by the next result.

**Theorem 2.2.4** (Brändén (2007, Theorem 5.6)). *A multi-affine polynomial  $g$  over  $\mathbb{C}^N$  is real stable if and only if the following inequality holds for any  $1 \leq i < j \leq N$ :*

$$\forall x \in \mathbb{R}^n, \quad \frac{\partial g}{\partial x_i}(x) \frac{\partial g}{\partial x_i}(x) - g(x) \frac{\partial^2 g(x)}{\partial x_i \partial x_j} \geq 0.$$

**Remark 2.2.1.** The above characterization of real stability only holds for multi-affine polynomials in  $\mathbb{C}[z_1, \dots, z_N]$ .

Finally, note that although strongly Rayleigh measures are CNA+, CNA+ measures are not necessarily SR: for  $N \geq 20$ , Borcea et al. (2009, Section 7) construct measures over  $2^{[N]}$  that are CNA+ but not SR.

### 2.2.2 Operations preserving real stability

In the previous section, we listed a sequence of operations that should ideally preserve a strong characterization of negative dependence. Fortunately, SR measures are stable under these operations, and more. As these operations will be crucial to proving the SR-ness of a variety of distributions over  $2^{[N]}$ , we reproduce them here.

**Proposition 2.2.5** (Borcea et al. (2009, Prop. 3.1)). *Let  $g$  be a polynomial over complex variables  $z_1, \dots, z_N$  with degree  $d_i$  in variable  $z_i$ . If  $g$  is real stable, the following polynomials are real stable or constant equal to zero:*

- (i)  $\partial_i g(z_1, \dots, z_N)$  for  $1 \leq i \leq N$
- (ii)  $g(z_1, \dots, z_{i-1}, \alpha z_i, z_{i+1}, \dots, z_N)$  for  $1 \leq i \leq N$  and  $\alpha > 0$
- (iii)  $g(z_1, \dots, z_{i-1}, \beta, z_{i+1}, \dots, z_N)$  for  $1 \leq i \leq N$  and  $\beta > 0$
- (iv)  $z_1^{d_1} \cdots z_N^{d_N} g(\lambda_1 z_1^{-1}, \dots, \lambda_N z_N^{-1})$  for  $\pm(\lambda_1, \dots, \lambda_N) \in \mathbb{R}_+^N$
- (v)  $g(z_1, \dots, z_{i-1}, z_j, z_{i+1}, \dots, z_N)$  for  $i \neq j$ .

Additionally, the two following operations, which restrict the distribution's support to sets of a fixed size, also preserve the strongly Rayleigh property.

**Proposition 2.2.6** (Borcea et al. (2009, Thm. 4.2)). *Let  $\mu$  be an SR measure over  $2^{[N]}$ . The symmetric homogenization of  $\mu$  is SR and defined over  $2^{[2N]}$  as*

$$\mu_{sh}(S) = \begin{cases} \mu(S \cap [N]) \binom{N}{|S \cap [N]|}^{-1} & \text{if } |S| = N \\ 0 & \text{otherwise.} \end{cases} \quad (2.1)$$

**Proposition 2.2.7** (Borcea et al. (2009, Corollary 4.18)). *Let  $\mu$  be a SR measure over  $2^{[N]}$ . The  $p, q$  truncation of  $\mu$  for  $q - p \leq 1$  is also SR, and defined over  $2^{[N]}$  as*

$$\mu_{p,q}(S) \propto \begin{cases} \mu(S) & \text{if } p \leq |S| \leq q \\ 0 & \text{otherwise.} \end{cases} \quad (2.2)$$

Finally, the limit of a sequence of real stable polynomials is also real stable.

**Proposition 2.2.8** (Borcea et al. (2009, Proposition 3.1)). *Let  $\{f_j\}_{j=1}^\infty$  be a sequence of real stable polynomials in  $\mathbb{R}[z_1, \dots, z_N]$  and  $f \not\equiv 0 \in \mathbb{R}[z_1, \dots, z_N]$  be the limit, uniformly on compact subsets of  $\mathbb{C}^N$ , of  $\{f_j\}_{j=1}^\infty$ ; then  $f$  is real stable.*

### 2.2.3 Examples

We now present a few classical distributions that are strongly Rayleigh and that have found applications in the larger scope of computer science.

#### The uniform distribution

Let  $\mu$  be the distribution over all subsets of  $2^{[N]}$  such that for all  $S \subseteq [N]$ ,  $\mu(S) = \frac{1}{2^N}$ . The corresponding generating polynomial  $g_\mu$  is

$$g_\mu(z_1, \dots, z_N) = \frac{1}{2^N} \sum_{S \subseteq [N]} z^S = \frac{1}{2^N} \prod_{i=1}^N (1 + z_i),$$

which clearly admits no zeros in the space  $\{z \in \mathbb{R}^N : \forall i, \mathfrak{Im}(z_i) > 0\}$ .

#### Product measures

The generalization of the uniform distribution, which assigns different weights to each element of the ground set, is also strongly Rayleigh:

**Proposition 2.2.9.** *The measure over  $2^{[N]}$  parametrized by scalars  $q_i \in [0, 1], 1 \leq$*

$i \leq N$ , such that  $\mu(S) = \prod_{i \in S} q_i \prod_{j \notin S} (1 - q_j)$  is SR and has generating polynomial

$$g_\mu(z_1, \dots, z_N) = \prod_{i=1}^N (q_i z_i + 1 - q_i).$$

### Determinantal measures

Many real stable polynomials are based on determinants of symmetric matrices; the following result provides one such example. As many subsequent results will depend on the result of Prop. 2.2.10, we also include its proof below.

**Proposition 2.2.10** (Borcea and Brändén (2009, Prop. 3.2)). *Let  $\mathbf{A}_0, \mathbf{A}_1, \dots, \mathbf{A}_N \in \mathbb{S}_n^+$ . The polynomial*

$$p : z \rightarrow \det(\mathbf{A}_0 + z_1 \mathbf{A}_1 + \dots + z_N \mathbf{A}_N)$$

*is either identically zero or real stable; hence as soon as  $p$  is multi-affine, it defines a strongly Rayleigh polynomial.*

*Proof.* As positive definite matrices are a dense subset of positive semi-definite matrices, we can via Prop. 2.2.8 limit our analysis to the setting where the  $\mathbf{A}_i$  are positive definite. Let  $e \in \mathbb{R}^N$  and  $x \in \mathbb{R}_+^N$ . The polynomial  $t \rightarrow p(e + tx) = \det(\mathbf{A}_0 + \sum_i e_i \mathbf{A}_i + t \sum_i x_i \mathbf{A}_i)$  can be rewritten as

$$t \rightarrow \det\left(\sum_i x_i \mathbf{A}_i\right) \det\left(t \mathbf{I} + \mathbf{M}^{-1/2} (\mathbf{A}_0 + \sum_i e_i \mathbf{A}_i) \mathbf{M}^{-1/2}\right),$$

where we write  $\mathbf{M} = \sum_i x_i \mathbf{A}_i \in \mathbb{S}_n^{++}$ . Then, since  $\mathbf{M}^{-1/2} (\mathbf{A}_0 + \sum_i e_i \mathbf{A}_i) \mathbf{M}^{-1/2}$  is real symmetric,  $t \rightarrow p(e + tx)$  is guaranteed to only have real zeros; invoking Prop. 2.2.2 concludes the proof.  $\square$

We will see in particular that the strong Rayleigh property for determinantal point processes (Chapters 3-4) and Generalized Volume Sampling measures (Chapter 5) is derived by applying various stability preserving operations to polynomials of the form  $z \rightarrow \det(\mathbf{A}_0 + \sum_i z_i \mathbf{A}_i)$ .

## 2.3 Tractable sampling

In order to use strongly Rayleigh measures in machine learning, we must be able to sample from such measures in a tractable fashion. Although specific SR measures such as determinantal point processes allow straightforward exact sampling algorithms (Chapter 3), we present here two general-purpose, *approximate* sampling algorithms for homogeneous (Algorithm 1) and non-homogeneous (Algorithm 2) SR measures.

These algorithms are Markov chain Monte Carlo (MCMC) methods that sample sets sequentially from the uniform distributions, and move from set  $S$  to the next set  $T$  based on the respective likelihoods of  $T$  and  $S$  under the target distribution  $\mu$ .

### 2.3.1 Fast sampling for homogeneous measures

A straightforward MCMC chain for sampling from  $k$ -homogeneous strongly Rayleigh measures is described in Algorithm 1 (technically, Algorithm 1 describes a *lazy* chain: the probability of moving from one state to another is at most  $\frac{1}{2}$ ).

---

**Algorithm 1** MCMC sampling for a  $k$ -homogeneous SR measure  $\mu$

---

**Input:** SR proposal  $\mu$   
Initialize  $S$  to a random subset of size  $k$  of  $[n]$   
**while** not mixed **do**  
    Sample  $i \in S$  and  $j \in \bar{S}$  uniformly at random  
     $T \leftarrow S \cup \{j\} \setminus \{i\}$   
     $S \leftarrow T$  with probability  $\frac{1}{2} \min(1, \mu(T)/\mu(S))$   
**return**  $S$

---

The crucial metric used to analyzing samplers such as Algorithm 1 is the chain's *mixing time*. The mixing time upper bounds the number of iterations of the algorithm's inner loop required so that the distribution of the generated samples  $\hat{\mu}$  is  $\epsilon$ -close to the true target  $\mu$ .

**Definition 2.3.1** (Mixing time). For a set  $S$  and  $\epsilon > 0$ , write  $\hat{\mu}_{t,S}$  for the distribution over samples generated by the MCMC chain initialized at subset  $S$  followed by  $t$

iterations of the chain's inner loop. The chain's mixing time  $\tau_S(\epsilon)$  is given by

$$\tau_S(\epsilon) = \min_{t>0} \|\hat{\mu}_{t,S} - \mu\|_{TV} \leq \epsilon.$$

Algorithm 1 was proposed by Anari et al. (2016); their key result is the following.

**Theorem 2.3.1** (Fast mixing MCMC chains for strongly Rayleigh measures). *Let  $\mu$  be a  $k$ -homogeneous strongly Rayleigh measure over  $2^{[N]}$ . Algorithm 1 initialized with subset  $S$  such that  $\mu(S) > 0$  has mixing time*

$$\tau_S(\epsilon) \leq 2kN \log \left( \frac{1}{\epsilon\mu(S)} \right).$$

### 2.3.2 Generalization to non-homogeneous SR measures

Fast sampling for homogeneous SR measures can be generalized to non-homogeneous measures by way of *symmetric homogenization* (Equation 2.1), as shown by Li et al. (2016d). This result relies on the following fact: to sample from a SR measure  $\mu$ , we can sample from the symmetric homogenization  $\mu_{sh}$  of  $\mu$  which is  $N$ -homogeneous (using Algorithm 1), then return the subset  $S' = S \cap [N]$ . One easily verifies that this yields a subset of  $[N]$  with the desired probability.

Li et al. (2016d) also introduce a more efficient way of applying the symmetrization trick; we retranscribe the proposed algorithm in Algorithm 2.

**Theorem 2.3.2** (Fast sampling for SR measures (Li et al., 2016d)). *If  $\mu$  is SR, Algorithm 2 mixes in time*

$$\tau_S(\epsilon) \leq 2N^2 \left( \log \binom{N}{|S|} + \log \frac{1}{\epsilon\mu(S)} \right).$$

Theorems 2.3.1 and 2.3.2 are fundamental to the application of strongly Rayleigh measures to machine learning. For any SR measure  $\mu$ , these theorems guarantee that — as long as we can efficiently compute the unnormalized value  $\mu(S)$  — we can sample from  $\mu$  efficiently.

---

**Algorithm 2** Lazy MCMC sampling for non-homogeneous SR measures

---

- 1: **Input:** SR proposal  $\mu$
- 2: Initialize  $S$  to a random subset of size  $n$  of  $[2n]$
- 3:  $S \rightarrow T \cap [n]$
- 4: **while** not mixed **do**
- 5:     Draw  $q \sim \text{Unif}[0, 1]$
- 6:     Sample  $i \in S$  and  $j \in \bar{S}$  uniformly at random
- 7:     **if**  $q \in [0, \frac{N-|S|}{2N^2}]$  **then**
- 8:          $S \leftarrow S \cup \{j\}$  with probability  $\min(1, \frac{\mu(S \cup \{t\})}{\mu(S)} \frac{|S|+1}{N-|S|})$
- 9:     **else if**  $q \in [\frac{N-|S|}{2N^2}, \frac{N-|S|}{2N}]$  **then**
- 10:          $S \leftarrow S \cup \{j\} \setminus \{i\}$  with probability  $\min(1, \frac{\mu(S \cup \{t\} \setminus \{s\})}{\mu(S)})$
- 11:     **else if**  $q \in [\frac{N-|S|}{2N}, \frac{|S|^2+N(N-|S|)}{2N^2}]$  **then**
- 12:          $S \leftarrow S \setminus \{i\}$  with probability  $\min(1, \frac{\mu(S \setminus \{s\})}{\mu(S)} \frac{|S|}{N-|S|+1})$
- 13:     **else**
- 14:         Do nothing

**return**  $S$

---



# Chapter 3

## Determinantal point processes — background and sampling

The class of Determinantal point processes (DPPs) is the class of strongly Rayleigh measures that has achieved the most popularity in machine learning. DPPs provide an elegant model to assign probabilities to the exponential number of subsets of a ground set of  $N$  items, while permitting polynomial time sampling and marginalization.

The antecedents of DPPs lie in statistical mechanics (Macchi, 1975); since the seminal work of (Kulesza and Taskar, 2012) they have made inroads into machine learning. Determinantal point processes have been applied to problems such as document and video summarization (Lin et al., 2012; Chao et al., 2015), sensor placement (Krause et al., 2008), recommender systems (Zhou et al., 2010), and object retrieval (Affandi et al., 2014). Recently, they have been used to compress layers in neural networks (Mariet and Sra, 2016a) and to provide sampling procedures for the Nyström method (Li et al., 2016c). The more general study of DPP properties has also garnered interest, see *e.g.*, (Hough et al., 2006; Kulesza and Taskar, 2011, 2012; Affandi et al., 2013; Decreusefond et al., 2015; Lavancier et al., 2015; Borodin, 2009; Lyons, 2003).

We begin this section by providing a summary of the previous key results on DPPs: their negative dependence properties, stability under conditioning, and exact sampling algorithms. However, despite polynomial-time exact sampling algorithms and approximate MCMC sampling schemes (Alg. 2), sampling from a DPP remains

a costly operation, especially when used jointly with other machine learning models such as neural networks. This, unfortunately, impedes the application of DPPs to modern machine learning problems, even when they are, theoretically, the right tool. This motivates our first contribution: we show that under the right conditions, DPPs can be efficiently approximated by generative networks. Such networks allow fast and parallelizable sampling algorithms, making DPPs both theoretically and practically well-suited to modern machine learning applications that require diverse sampling.

### 3.1 Background

As before, we consider  $\mathcal{Y} = \{1, \dots, N\}$  a ground set of  $N$  items.

**Definition 3.1.1** (Determinantal point process). A probability measure  $\mu$  over  $2^{\mathcal{Y}}$  is a determinantal point process if, for all  $Y \subseteq \mathcal{Y}$  drawn according to  $\mu$ , we have

$$\forall A \subseteq \mathcal{Y}, \quad \Pr_{\mu}(A \subseteq Y) = \det(\mathbf{K}_A), \tag{3.1}$$

for some positive semi-definite matrix  $\mathbf{K} \in \mathbb{S}_N^+$  such that  $\mathbf{0} \preceq \mathbf{K} \preceq \mathbf{I}$  (the DPP's *marginal kernel*). We adopt the convention  $\det(\mathbf{K}_{\emptyset}) = 1$ .

If a DPP with marginal kernel  $\mathbf{K}$  assigns nonzero probability to the empty set, the DPP can alternatively be parametrized by a positive definite matrix  $\mathbf{L}$  such that

$$\mu(Y) = \det(\mathbf{L}_Y)/\det(\mathbf{L} + \mathbf{I}). \tag{3.2}$$

In this thesis, all DPPs we consider will be parametrizable via the kernel  $\mathbf{L}$ . A brief manipulation shows that in this case,  $\mathbf{K} = \mathbf{L}(\mathbf{L} + \mathbf{I})^{-1} = \mathbf{I} - (\mathbf{L} + \mathbf{I})^{-1}$ .

Crucially, determinantal point processes are strongly Rayleigh measures; as such, they inherit all the characterizations of negative dependence introduced in Chapter 2.

**Theorem 3.1.1.** *Let  $\mathbf{L} \succeq 0$  be a positive definite matrix in  $\mathbb{R}^{N \times N}$ . Then, the measure over  $2^{[N]}$  defined by  $\mu(S) \propto \det(\mathbf{L}_S)$  is strongly Rayleigh.*

*Proof.* The proof relies on the following identity, which can be easily proven by induction on  $N - |A|$ . Writing  $\mathbf{Z} = \text{diag}(z_1, \dots, z_N)$  and  $\bar{A} \triangleq [N] \setminus A$ , we have

$$\forall A \subseteq [N], \sum_{A \subseteq Y \subseteq [N]} \det \mathbf{L}_Y z^Y = \det(\mathbf{L}\mathbf{Z} + \mathbf{I}_{\bar{A}}). \quad (3.3)$$

Evaluating (3.3) at  $A = \emptyset$ , we see that the generating polynomial of  $\mu$  is of the form  $\det(\mathbf{L}\mathbf{Z} + \mathbf{I})$ . Hence, the generating polynomial for  $\mu$  is given by

$$g_\mu(z_1, \dots, z_n) \propto \det(\mathbf{L}\mathbf{Z} + \mathbf{I}) = \det(\mathbf{L}^{-1} + \mathbf{Z}) \det(\mathbf{L}^{-1}).$$

Hence,  $\det(\mathbf{L}^{-1} + \mathbf{Z}) = 0$  if and only if one of the  $z_i$  is an eigenvalue of  $\mathbf{L}^{-1}$ ; as  $\mathbf{L} \succ 0$ , its eigenvalues are all real, and so  $g_\mu(z) \neq 0$  if  $\Im(z_i) > 0$  for all  $i$ .  $\square$

### 3.1.1 Key properties of determinantal point processes

#### Intuition behind the quality/diversity trade-off

That a DPP enables a trade-off between quality and diversity within sampled sets is apparent for sets of size 2. Let  $\mathbf{K} \succeq 0$  be the marginal kernel of a DPP, and let  $S \subseteq [N]$  be a set of size 2. By definition,

$$\Pr(S \subseteq Y) = \mathbf{K}_{ii}\mathbf{K}_{jj} - \mathbf{K}_{ij}^2,$$

and so this probability can be increased by increasing either the coefficients  $\mathbf{K}_{ii}$  and  $\mathbf{K}_{jj}$  (*point-wise quality* of items  $i$  and  $j$ ) or by reducing the *similarity* between items  $i$  and  $j$  given by  $\mathbf{K}_{ij} = \mathbf{K}_{ji}$ .

More generally, the ability of a DPP trade-off quality with diversity is apparent through the geometric interpretation of determinants as volumes in  $d$  dimensional spaces. Let  $\Phi$  be a  $\mathbb{R}^{n \times d}$  matrix such that  $\mathbf{L} = \Phi\Phi^T$  and let  $S = \{i_1, \dots, i_k\} \subseteq [N]$ . Then we can rewrite the probability of sampling  $S$  under the DPP as

$$\mu(S) \propto \det(\mathbf{L}_S) = \det(\Phi_{S,:}\Phi_{S,:}^\top) = \text{Vol}(\phi_{i_1}, \dots, \phi_{i_k})^2,$$

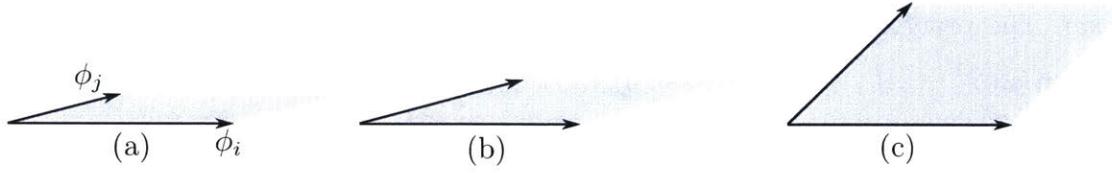


Figure 3-1: The volume spanned by vectors  $\phi_i$  and  $\phi_j$  can be increased either by elongating a vector (b) or by increasing the cosine distance between the vectors (c).

where  $\text{Vol}(\phi_{i_1}, \dots, \phi_{i_k})^2$  is the volume spanned by vectors  $\phi_{i_1}, \dots, \phi_{i_k}$ ; note that this volume may be zero due to dimensionality issues. Hence, the probability of  $S$  under the DPP with kernel  $\mathbf{L}$  can be increased by increasing either the norm of the  $\phi_i$  (*i.e.*, the point-wise *qualities* of each item) or the angle between the feature vectors (the *diversity* of the set), as illustrated in Fig. 3-1.

### Operations over DPPs

Another ingredient to the success of DPPs lies in their robustness to standard operations over distributions such as restrictions, complements and conditioning. These operations for DPPs are covered in detail in (Kulesza and Taskar, 2012, Chapter 2.3); we summarize the results here.

**Proposition 3.1.2** (Restriction). *If set  $Y$  is distributed according to a DPP with marginal kernel  $\mathbf{K}$  and  $A \subseteq \mathcal{Y}$ , then  $Y \cap A$  is distributed as a DPP with marginal kernel  $\mathbf{K}_A$ .*

**Proposition 3.1.3** (Complement). *If  $Y$  is distributed according to a DPP with marginal kernel  $\mathbf{K}$ , then  $\mathcal{Y} \setminus Y$  is distributed as a DPP with marginal kernel  $\mathbf{I} - \mathbf{K}$ .*

**Proposition 3.1.4** (Conditioning). *If  $Y$  is distributed according to a DPP with kernel  $\mathbf{L}$  and  $A, B \subseteq \mathcal{Y}$  are such that  $A \cap B = \emptyset$ , then*

- $\Pr(Y = B \mid A \cap Y = \emptyset) \propto \det \mathbf{L}_Y$ : *the event that no element of  $A$  appears in  $Y$  is distributed according to the DPP with kernel  $\mathbf{L}_{\bar{A}}$ .*
- $\Pr(Y = A \cup B \mid A \subseteq Y) \propto \det \mathbf{L}_Y^A$ : *the event that no element of  $A$  appears in  $Y$  is distributed according to the DPP with kernel  $\mathbf{L}^A := ([(\mathbf{L} + \mathbf{I}_{\bar{A}})^{-1}]_{\bar{A}})^{-1} - \mathbf{I}$ .*

The closure of DPPs to the above operations contributes to their success in machine learning. For example, recommender systems commonly leverage additional information, *e.g.*, when recommending items given other options the user has already expressed interest in (amounting to using a DPP conditioned on pre-selected items).

### 3.1.2 Sampling from DPPs and $k$ -DPPs

In addition to the approximate sampler introduced earlier for all strongly Rayleigh measures, DPP sampling can be done exactly in  $\mathcal{O}(N^3 + Nk^3)$  operations, where  $k$  is the size of the final subset; the algorithm is due to Hough et al. (2006). The dominant cost is the initial eigendecomposition of  $\mathbf{L}$ , which requires  $\mathcal{O}(N^3)$  operations.

---

**Algorithm 3** Exact sampling from a DPP

---

```

1: Input: eigendecomposition  $(v_i, \lambda_i)_{i=1}^N$  of  $\mathbf{L}$ 
2:  $J, Y \leftarrow \emptyset$ 
3: for  $i = 1$  to  $N$  do
4:    $J \leftarrow J \cup \{i\}$  with probability  $\frac{\lambda_i}{\lambda_i + 1}$ 
5:  $V \leftarrow \{v_i\}_{i \in J}$ ,  $k \leftarrow |V|$ 
6: for  $i = 1$  to  $k$  do
7:   Select  $i$  from  $\mathcal{Y}$  with probability  $\Pr(i) = \frac{1}{|V|} \sum_{v \in V} v_i^2$ 
8:    $Y \leftarrow Y \cup \{i\}$ 
9:    $V \leftarrow V_{\perp}$ , orthonormal basis to the subspace of  $V$  orthogonal to  $e_i$ .
10: return  $Y$ 
```

---

Alg. 3 relies on the fact that any DPP can be expressed a mixture of *elementary* DPPs; we refer the reader to (Kulesza and Taskar, 2012). From Alg. 3, it is easy to see that the expected sample size under a DPP is the expected size of  $V$  (line 5).

**Proposition 3.1.5** (Expected sample size). *The expected sample size of a DPP with kernel  $\mathbf{L}$  and marginal kernel  $\mathbf{K}$  is given by*

$$\mathbb{E}[|Y|] = \text{tr}(\mathbf{L}(\mathbf{I} + \mathbf{L})^{-1}) = \text{tr}(\mathbf{K}).$$

However, many applications require explicit control of the size of DPP samples; for this purpose, we need a measure that assigns probability 0 to any set of size different than a specific size  $k$ . This can be achieved with  $k$ -DPPs (Kulesza and Taskar, 2011).

**Definition 3.1.2.** A  $k$ -DPP is a measure over all subsets of  $[N]$  defined by a kernel  $\mathbf{L} \succeq 0 \in \mathbb{R}^{N \times N}$  and an integer  $k \leq N$  that assigns to  $S \subseteq [N]$  the probability

$$\Pr(S) = \det(\mathbf{L}_S)/e_k(\lambda_1, \dots, \lambda_N),$$

if  $|S| = k$ , and 0, otherwise, where the  $\lambda_i$  are the eigenvalues of  $\mathbf{L}$  and  $e_k$  is the  $k$ -th elementary symmetric polynomial, defined as  $e_k(\lambda_1, \dots, \lambda_N) = \sum_{|S|=k} \prod_{i \in S} \lambda_i$ .

Since  $k$ -DPPs are not DPPs, they do not have marginal kernels. However,  $k$ -DPPs still admit an exact sampling algorithm; in fact, to sample from a  $k$ -DPP we only have to change how we choose the initial set of eigenvectors  $V$  (first 5 lines of Alg. 3).

---

**Algorithm 4** Selecting eigenvectors to sample from a  $k$ -DPP

---

```

1: Input: eigendecomposition  $(v_i, \lambda_i)_{i=1}^N$  of  $\mathbf{L}$ , desired set size  $k \geq 0$ .
2:  $J \leftarrow \emptyset$ 
3:  $\ell \leftarrow k$ 
4: for  $i = N$  to 1 do
5:   if  $\ell = 0$  then
6:     break
7:   if  $u \sim U([0, 1]) < \lambda_i \frac{e_{\ell-1}(\lambda_1, \dots, \lambda_{i-1}, \lambda_{i+1}, \dots, \lambda_N)}{e_\ell(\lambda_1, \dots, \lambda_N)}$  then
8:      $J \leftarrow J \cup \{i\}$ 
9:    $\ell \leftarrow \ell - 1$ 
10:   $V \leftarrow \{v_i\}_{i \in J}$ 
11: return  $V$ 

```

---

Introduced in (Kulesza and Taskar, 2011), Algorithm 4 samples the initial set of eigenvectors  $V$ . To sample from a  $k$ -DPP, one can simply run the second loop of Algorithm 3 with the set  $V$  chosen as in Algorithm 4.

Furthermore, when a low-rank decomposition of the kernel is already known, Algorithms 3 and 4 can be modified to take advantage of this pre-existing structure (Kulesza, 2012). Indeed, writing  $\mathbf{L} = \mathbf{B}^\top \mathbf{B}$ , with  $\mathbf{B} \in \mathbb{R}^{d \times N}$  and  $d < N$ , the overall computational cost can be brought down to  $\mathcal{O}(Ndk^2 + d^2k^3)$  (see Algorithm 5) by using the *dual kernel*  $\mathbf{C} = \mathbf{B}\mathbf{B}^\top \in \mathbb{R}^{d \times d}$ .

---

**Algorithm 5** Sampling from a DPP with dual representation

---

- 1: **Input:** eigendecomposition  $(v_i, \lambda_i)_{i=1}^N$  of  $\mathbf{C} = \mathbf{B}\mathbf{B}^\top \in \mathbb{R}^{d \times d}$
  - 2:  $J, Y \leftarrow \emptyset$
  - 3: **for**  $i = 1$  **to**  $N$  **do**
  - 4:    $J \leftarrow J \cup \{i\}$  with probability  $\frac{\lambda_i}{\lambda_i + 1}$
  - 5:  $V \leftarrow \{v_i / \sqrt{v_i^\top \mathbf{C} v_i}\}_{i \in J}$
  - 6: **while**  $|V| > 0$  **do**
  - 7:   Select  $i$  from  $\mathcal{Y}$  with probability  $\Pr(i) = \frac{1}{|V|} \sum_{v \in V} (v^\top \mathbf{B}_i)^2$
  - 8:    $Y \leftarrow Y \cup \{i\}$
  - 9:    $v_0 \leftarrow$  any vector in  $V$  such that  $\mathbf{B}_i^\top v_0 \neq 0$
  - 10:    $V \leftarrow \left\{ v - \frac{v^\top \mathbf{B}_i}{v_0^\top \mathbf{B}_i} v_0 \mid v \in V \setminus \{v_0\} \right\}$
  - 11:    $V \leftarrow V_\perp$ , orthonormalization of  $V$  w.r.t. the dot product  $\langle u, v \rangle = u^\top \mathbf{C} v$
  - 12: **return**  $Y$
- 

In many cases, the decomposition  $\mathbf{L} = \mathbf{B}^\top \mathbf{B}$  may be known for  $\mathbf{B} \in \mathbb{R}^{d \times N}$  with  $d \gg N$  (for example, when  $\mathbf{L}$  is constructed using linear feature kernel with high-dimensional features). Fortunately, Gillenwater et al. (2012) showed that in these cases a low-dimensional, approximate decomposition of  $\mathbf{L}$  can be constructed using random projections without significantly altering the induced DPP.

With this approximate representation, Algorithm 5 remains practical for large values of  $N$  and  $d$ , as long as the dual kernel  $\mathbf{C}$  can be obtained efficiently.

### 3.2 Modeling DPPs with neural networks

Due to its  $\mathcal{O}(N^3)$  cost, exact sampling for DPPs is impractical for large datasets; dual sampling (Algorithm 5) only applies when we have knowledge of a linear factorization of the DPP kernel  $\mathbf{L}$ . In some scenarios, even MCMC sampling can be prohibitively costly due to its inherently sequential nature.

Motivated by recent work on modeling set functions with neural networks (Zaheer et al., 2017; Cotter et al., 2018), we propose here to generate approximate samples via a generative network; this allows for simple parallelization while simultaneously benefiting from recent improvements in specialized architectures for neural network models (*e.g.*, parallelized matrix multiplications).

In the remainder of this chapter, we introduce DPPNETs: generative deep models that produce DPP-like samples for arbitrary ground sets and feature representations for the items in  $\mathcal{Y}$ . We show that a simple, carefully constructed neural network can generate approximate DPP samples with very little overhead, while maintaining some of the fundamental theoretical properties of DPP measures. Furthermore, we show that DPPNETs can be trivially employed to sample from a conditional DPP (*i.e.*, sampling  $S$  such that  $A \subseteq S$  is predefined) and for greedy mode approximation. Material covered in this section is based on (Mariet et al., 2019b).

### 3.2.1 Generating DPP samples with deep models

We begin by considering which DPP properties we may leverage to simplify training, as well as certain DPP-specific properties we may want to maintain within the generative model. These properties guide our choice for neural network architecture.

**Simple computation of marginal probabilities.** Given a set  $Y$  sampled by a DPP with kernel  $\mathbf{L}$  and  $S \subseteq Y$ , the probability  $\Pr(\{i\} \cup S \subseteq Y \mid S \subseteq Y)$  of item  $i \notin S$  also being sampled admits the closed form (where we write  $\bar{S} := [N] \setminus S$ ):

$$\Pr(\{i\} \cup S \subseteq Y \mid S \subseteq Y) = 1 - [(\mathbf{L} + \mathbf{I}_{\bar{S}})^{-1}]_{ii}. \quad (3.4)$$

Although Eq. 3.4 requires an expensive matrix inversion, such costs can be offset during off-line training. These probabilities are the vector that DPPNET seeks to output. This signal has the advantage of providing information during every step of training, compared to other downstream possibilities (*e.g.*, NLL of generated sets).

**Sequential sampling.** When drawing samples from a DPP, the standard DPP sampling algorithm (Kulesza and Taskar, 2012, Alg. 1) generates the sample as a sequence, adding items one after the other until reaching a pre-determined size, similarly to Alg. 6. We take advantage of this by recording all intermediary subsets generated by the DPP when sampling training data: in practice, instead of training on  $n$  subsets of size  $k$ , we train on  $kn$  subsets of size  $0, \dots, k - 1$ . Thus, our model is much like an unrolled recurrent neural network (Elman, 1990; Jordan, 1990).

The sequential form of the exact sampling algorithm also makes it amenable to simple modifications that yield greedy sampling algorithms (Chen et al., 2018). This motivates a sequential sampling algorithm for DPPNET (Alg. 6), yielding a straightforward greedy form of DPPNET sampling without additional complexity.

---

**Algorithm 6** Sampling and greedy mode for DPPNET

---

**Input:** Initial set  $S$ , target size  $k$ , feature matrix  $\Phi$

**while**  $|S| < k$  **do**

$v \leftarrow \text{DPPNET}(S, \Phi)$

**if** sampling **then**

$i \sim \text{Multinomial}(v/\|v\|)$

**else if** greedy mode **then**

$i \leftarrow \text{argmax } v$

$S \leftarrow S \cup \{i\}$

**return**  $S$

---

**Closure under conditioning.** Recall that DPPs are closed under conditioning: given  $A \subseteq \mathcal{Y}$ , the conditional distribution over  $\mathcal{Y} \setminus A$  given by  $\Pr(S = A \cup B \mid A \subseteq S)$  (for  $B \cap A = \emptyset$ ) is also a DPP with kernel  $\mathbf{L}^A = ([(\mathbf{L} + \mathbf{I}_A)^{-1}]_A)^{-1} - \mathbf{I}$ . Although conditioning comes at the cost of an expensive matrix inversion, this property make DPPs well-suited to applications requiring diversity in conditioned sets, such as basket completion for recommender systems.

Standard deep generative models such as (Variational) Auto-Encoders (Kingma and Welling, 2014) (VAEs) and Generative Adversarial Networks (Goodfellow et al.,

2014) (GANs) would not enable conditioning operations during sampling — such operations would have to take place in the latent space. With the DPPNET architecture, we can sample a set via Alg. 6, which allows for trivial basket-completion type conditioning operations.

**Summary.** Given an input subset  $S \subseteq \mathcal{Y}$ , a DPPNET returns a probability vector over all items in  $\mathcal{Y}$ . The next item to add to  $S$  is sampled from the associated categorical distribution. During training, the DPPNET minimizes the distance between the generated probability vector and the true marginal probabilities under the target DPP. When generating training data from the true DPP, we use the additional information produced during exact sampling to obtain additional training data.

### 3.2.2 Sampling over fixed and varying ground sets

In simple settings, we wish to draw samples over a ground set with a fixed feature matrix representation  $\Phi \in \mathbb{R}^{N \times d}$ . In this case, DPPNET’s knowledge of  $\Phi$  can be obtained during training, and so DPPNET is a simple feed-forward network taking a partially sampled set  $S$  as input. However, in many cases the feature representation of the items  $\Phi$  will vary across time or contexts: for example, when  $\Phi$  represents a pool of products that are available for sale at a given time, or social media posts whose relevance that varies based on context. In such cases, DPPNET also takes as input the feature matrix  $\Phi$ .

Naively adding  $\Phi$  as input to a stack of feed-forward connections would require deeper networks and larger layers, increasing learning and sampling time. Instead, we draw inspiration for the dot-product attention from (Vaswani et al., 2017). Intuitively, attention is a vector computed by the network that indicates relevant parts of the inputs. For DPPNET, this attention will *reweight*  $\Phi$  based on pre-sampled items.

Attention functions were introduced in (Bahdanau et al., 2015) for neural network-based translation. Intuitively, attention functions reweight inputs to a neural network, assigning higher weights to input segments that carry information relevant to the task at hand. In (Vaswani et al., 2017), the attention mechanism takes three matrices as

input, which can be viewed as a (1) matrix of keys  $\mathbf{K}$ , (2) a matrix of values  $\mathbf{V}$ , and (3) a query matrix  $\mathbf{Q}$ . The attention matrix  $\mathbf{A}$ , which reweights the values  $\mathbf{V}$ , is obtained as  $\mathbf{A} = \text{softmax}(\mathbf{Q}\mathbf{K}^\top/\sqrt{d})$  where  $d$  is the dimension of each query/key: the inner product acts as a proxy to the similarity between queries and keys.

Here, the feature representation of items in the input set  $S$  acts as the query  $\mathbf{Q} = \Phi_{S,:} \in \mathbb{R}^{k \times d}$ . The features  $\Phi \in \mathbb{R}^{N \times d}$  of our ground set are both the keys and the values. In order for the attention mechanism to make sense in the framework of DPP modeling, we make two modifications to the original attention mechanism:

- We want our network to focus on items that are *dissimilar* to the query (input subset): for each item  $i$  in the input subset  $S$ , we compute its pairwise dissimilarity to each item  $i \in \mathcal{Y}$  as the vector  $d_i = 1 - \text{softmax}(\mathbf{Q}_i\Phi^\top/\sqrt{d})$ .
- We return a vector  $a \in \mathbb{R}^N$  in the probability simplex such that  $a_j \propto \prod_{i \in S} d_{ij}$ . This allows us to have a fixed-size input to the network, and enforces the fact that similarity to a *single* item is enough to disqualify an element from the ground set.

Finally, our attention vector  $a$  is computed via the *inhibitive attention* mechanism

$$a' = \odot_{i \in S} \left( 1 - \text{softmax}(\mathbf{Q}\Phi^\top/\sqrt{d}) \right), \quad a = a'/\|a'\|_1, \quad (3.5)$$

where  $\odot$  represents the row-wise multiplication operator; this vector can be computed in  $\mathcal{O}(kDN)$  time. The attention component finally feeds the element-wise multiplication of each row of  $\mathbf{V}$  with the attention vector  $a$  to the feed-forward block, effectively reweighting each ground set item by its dissimilarity to the pre-selected items in  $S$ .

**Remark 3.2.1.** Dual sampling (Alg. 5) is efficient only when the kernel is constructed as  $\mathbf{L} = \Phi\Phi^\top$ ; for non-linear kernels, a low-rank decomposition of  $\mathbf{L}$  must first be obtained, which in the worst case requires  $\mathcal{O}(N^3)$  operations. In comparison, DPPNET can be trained on non-linear DPP kernels, while only requiring  $\Phi$  as input.

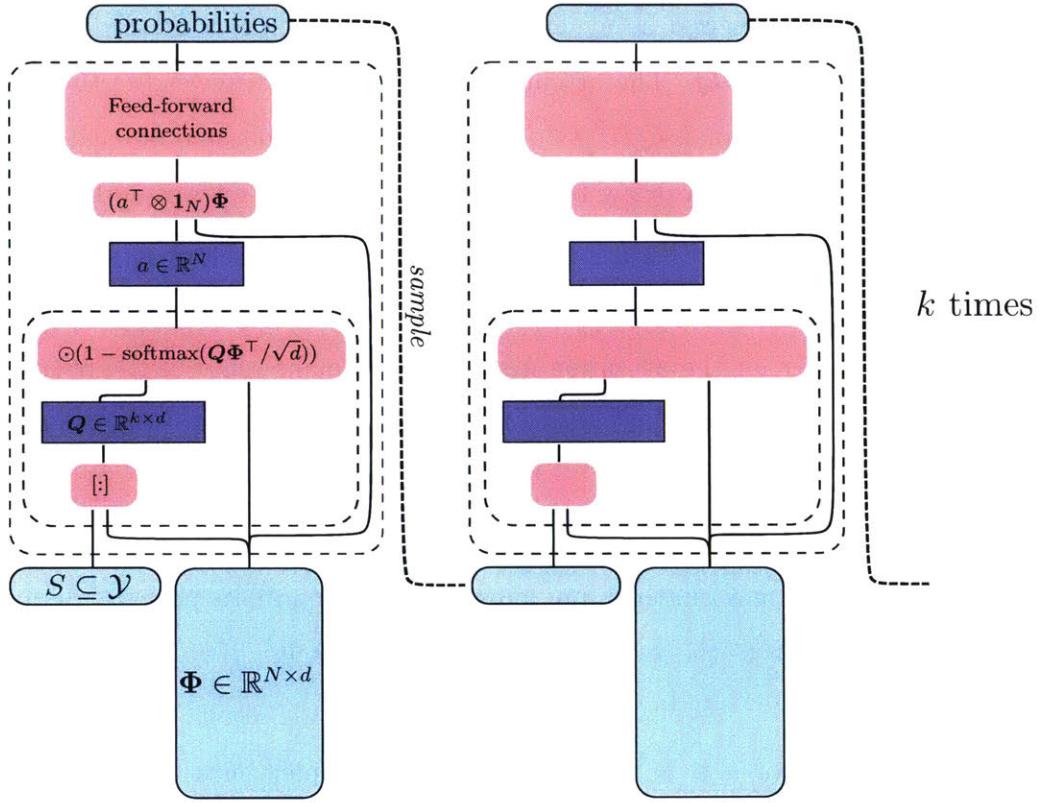


Figure 3-2: Transformer network architecture (variable feature matrix  $\Phi$ ).

### 3.2.3 Preserving log-submodularity

We now investigate whether any neural network approximation of determinantal point processes can inherit some of their negative dependence properties. Although the SR property is delicate and maintained by only few operations, log-submodularity, also a fundamental property in discrete optimization, see *e.g.*, (Zagoruyko and Komodakis, 2016; Buchbinder et al., 2012; Grötschel et al., 1981), is more robust. We show here that (log) submodularity can be inherited by a generative model.

**Theorem 3.2.1.** *Let  $\mathcal{P}$  be a strictly submodular distribution over subsets of  $\mathcal{Y}$ , and  $\mathcal{Q}$  be a function over the same space such that*

$$D_{\text{TV}}(\mathcal{P}, \mathcal{Q}) \leq \min_{S \neq T \notin \{\emptyset, \mathcal{Y}\}} \frac{1}{4} [\mathcal{P}(S) + \mathcal{P}(T) - \mathcal{P}(S \cup T) - \mathcal{P}(S \cap T)], \quad (3.6)$$

where  $D_{\text{TV}}$  indicates the total variational distance. Then  $\mathcal{Q}$  is also submodular.

*Proof.* Let  $S, T \subseteq \mathcal{Y}$  be subsets such that  $S \neq T$  and  $S, T \neq \emptyset, \mathcal{Y}$  (the submodular inequality is a trivial equality for  $S = T$ ,  $S = \emptyset$  and  $S = \mathcal{Y}$ , for any distribution). We must prove  $\mathcal{Q}(S) + \mathcal{Q}(T) \geq \mathcal{Q}(S \cup T) + \mathcal{Q}(S \cap T)$ . Let  $\epsilon$  be defined by

$$\epsilon := \min_{S \neq T \neq \emptyset, \mathcal{Y}} \mathcal{P}(S) + \mathcal{P}(T) - \mathcal{P}(S \cup T) - \mathcal{P}(S \cap T).$$

From the strict submodularity hypothesis, we know  $\epsilon > 0$ . By definition of  $\epsilon$ ,

$$\mathcal{P}(S) + \mathcal{P}(T) - \mathcal{P}(S \cup T) - \mathcal{P}(S \cap T) \geq \epsilon.$$

From (3.6), we know that  $\max_{S \subseteq \mathcal{Y}} |\mathcal{P}(S) - \mathcal{Q}(S)| \leq \epsilon/4$ . It follows that

$$\mathcal{Q}(S) + \mathcal{Q}(T) - \mathcal{Q}(S \cup T) + \mathcal{Q}(S \cap T) \geq \mathcal{P}(S) + \mathcal{P}(T) - \mathcal{P}(S \cup T) - \mathcal{P}(S \cap T) - \epsilon \geq 0,$$

which proves the submodularity of  $\mathcal{Q}$ .  $\square$

From Thm. 3.2.1 and the equivalence of norms in finite dimensional spaces, we can furthermore state the following result.

**Corollary 3.2.1.1.** *Let  $\mathcal{P}_L$  be a strictly log-submodular DPP over  $\mathcal{Y}$  and DPPNET be a network with a loss function of the form  $\|p - q\|$  where  $\|\cdot\|$  is a norm and  $p \in \mathbb{R}^{2^N}$  (resp.  $q$ ) is the probability vector assigned by the DPP (resp. the DPPNET) to each subset of  $\mathcal{Y}$ . Let  $\alpha = \max_{\|x\|_\infty=1} \|x\|^{-1}$ . If DPPNET converges to a loss smaller than*

$$\min_{S \neq T \neq \emptyset, \mathcal{Y}} \frac{1}{4\alpha} [\mathcal{P}(S) + \mathcal{P}(T) - \mathcal{P}(S \cup T) - \mathcal{P}(S \cap T)],$$

*its generative distribution is log-submodular.*

**Remark 3.2.2.** Cor. 3.2.1.1 is generalizable to the KL divergence loss  $D_{\text{KL}}(\mathcal{P} \parallel \mathcal{Q})$  via Pinsker's inequality.

**Remark 3.2.3.** Cor. 3.2.1.1 also holds for (log)-supermodular functions.

Checking numerically whether the conditions for Corollary 3.2.1.1 apply during training is NP-hard: the results of this section are purely theoretical. However, The-

orem 3.2.1 and Corollary 3.2.1.1 provide an additional justification for the use of probabilities in the objective function of DPPNET.

### 3.2.4 Experimental results

To evaluate DPPNET, we look at its performance both as a proxy for a DPP over a fixed ground set (Section 3.2.4) and as a tool for generating diverse subsets of varying ground sets (Section 3.2.4). Our models are trained with TensorFlow, using the Adam optimizer. Hyper-parameters are tuned to maximize the normalized log-likelihood of generated subsets. We compare DPPNET to DPPs and two additional baselines:

- **UNIF:** Uniform sampling over the ground set.
- **HCP:** Matérn hard core point processes. Points are sampled from a Poisson distribution, then thinned out to remove points within a distance  $r < 0.2$  (chosen by cross-validation) from each other.
- **$k$ -MEDOIDS:**  $k$ -medoids clustering (Hastie et al., 2001, 14.3.10) applied to the ground set; the distance between points is the one also used by DPP kernel.

We use the negative log-likelihood (NLL) of a subset under a fixed DPP over the ground set to evaluate the subsets obtained by all methods. This choice is motivated by the following considerations: DPPs have become a standard way of measuring and enforcing diversity over subsets of data in machine learning, and b) to the extent of our knowledge, there is no other standard method to benchmark the diversity of a selected subset that depends on specific dataset encodings.

#### Sampling over the unit square

We begin by analyzing the performance of a DPPNET trained on a DPP over a fixed ground set: the unit square. Our choice of the unit square is motivated by the need for diverse sampling methods on the unit hypercube, *e.g.*, quasi-Monte Carlo methods, Latin hypercube sampling (McKay et al., 1979) and low discrepancy sequences.

The ground set consists of the 100 points lying at the intersections of the  $10 \times 10$  grid on the unit square. The DPP is defined by setting its kernel  $\mathbf{L}$  to the exponentiated quadratic form  $\mathbf{L}_{ij} = \exp(-\|x_i - x_j\|_2^2/2)$ . As the ground set is fixed, the associated neural networks do not use the inhibitive attention mechanism and consists only of a feed-forward architecture.

We report the performance of the different sampling methods in Fig. 3-3. Visually (Fig. 3-3a) and quantitatively (Fig. 3-3b), DPPNET improves significantly over all other baselines. Furthermore, greedily sampling the mode from the DPPNET achieves a better NLL than DPP samples themselves. Note that a direct greedy approximation of the DPP would be significantly costlier. Numerical results are reported in Table 3.1.

Table 3.1: NLL under  $\mathcal{P}_L$  for sets of size  $k = 20$  sampled over the unit square. DPPNET achieves comparable performance to the DPP, outperforming the other baselines.

DPP	UNIFORM	HCP	$k$ -MEDOIDS	DPPNET
$154.95 \pm 2.93$	$180.53 \pm 9.56$	$163.40 \pm 5.87$	$169.37 \pm 6.41$	$153.44 \pm 2.07$

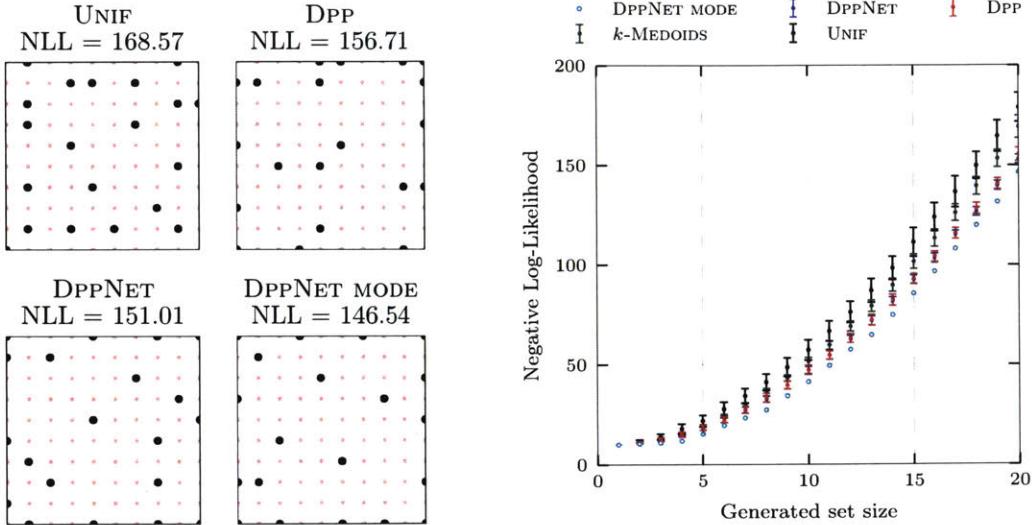
### Sampling over variable ground sets

We evaluate the performance of DPPNETs on varying ground set sizes through the MNIST (LeCun and Cortes, 2010), CelebA (Liu et al., 2015), and MovieLens (Harper and Konstan, 2015) datasets. To obtain the feature representations of each point, we use a Variational Auto-Encoder (Kingma and Welling, 2014) for MNIST and CelebA<sup>1</sup> (generating features of length 32). For MovieLens, we obtain feature vectors of length 10 for each movie by applying non-negative matrix factorization to the rating matrix.

During training, the networks take as input feature matrices drawn from the training portion of the associated dataset (MNIST, CelebA, MovieLens). At test time, the feature matrices  $\Phi$  are drawn from the associated test portions.

---

<sup>1</sup>Details on the encoders are provided in Appendix A.2.1.



(a) Sampling subsets of size 20 over the unit square. (b) Normalized log-likelihood of samples drawn from all methods as a function of the sampled set size.

Figure 3-3: Sampling on the unit square with a DPPNET (1 hidden layer with 841 neurons) trained on a single DPP kernel. Visually, DPPNET gives similar results to the full DPP (left). As evaluated by DPP NLL, the DPPNET’s mode achieves superior performance to the full DPP, and DPPNET sampling overlaps completely with DPP sampling (right).

The DPPNET is trained based on samples from DPPs with a linear kernel for MovieLens and with an exponentiated quadratic kernel for the image datasets. Bandwidths were set to  $\beta = 0.0025$  for MNIST and  $\beta = 0.1$  for CelebA in order to obtain a DPP average sample size  $\approx 20$ : recall that for a DPP with kernel  $\mathbf{L}$ , the expected sample size is given by Proposition 3.1.5.

For MNIST, Fig. 3-4 shows images selected by the baselines and the DPPNET, chosen among 100 digits with either random labels or all identical labels; visually, DPPNET and DPP samples provide a wider coverage of writing styles. The NLL of samples from DPPNET decay significantly, whereas the DPPNET mode continues to maintain competitive performance with DPP samples. For this reason, all subsequent experiments use only the DPPNET mode when working with varying ground sets.

Numerical results for MNIST are reported in Table 3.2; additionally to the previous baselines, we also considered learning a feed-forward neural network without

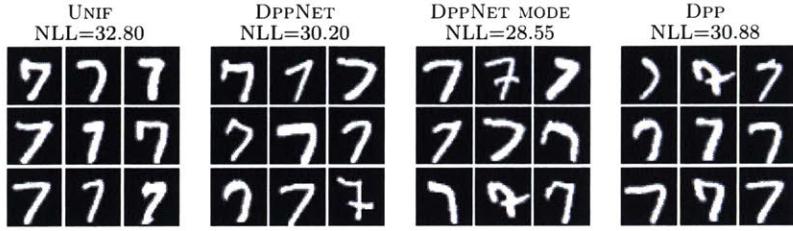


Figure 3-4: Digits sampled from MNIST based on their encodings. Only images labeled as “7” are provided to the various samplers; however, the DPPNET itself (3 layers of 365 neurons) is trained on all types of digits, and outperforms all baselines.

Table 3.2: NLL (mean  $\pm$  standard error) under the true DPP of samples drawn uniformly, according to the mode of the DPPNET, and from the DPP itself. We sample subsets of size 20; for each class of digits we build 25 feature matrices  $\Phi$  from encodings of those digits, and for each feature matrix we draw 25 different samples. Bolded numbers indicate the best-performing (non-DPP) sampling method.

	0	1	2	3	4	
DPP	$52.2 \pm 0.1$	$60.5 \pm 0.1$	$49.8 \pm 0.0$	$50.7 \pm 0.1$	$51.0 \pm 0.1$	
UNIF	$54.9 \pm 0.1$	$65.1 \pm 0.1$	$51.5 \pm 0.1$	$52.9 \pm 0.1$	$53.3 \pm 0.1$	
MEDOIDS	$55.1 \pm 0.1$	$65.0 \pm 0.1$	$51.5 \pm 0.0$	$52.9 \pm 0.1$	$53.1 \pm 0.1$	
DPPNET MODE	<b><math>53.6 \pm 0.3</math></b>	<b><math>63.6 \pm 0.4</math></b>	<b><math>50.8 \pm 0.2</math></b>	<b><math>51.4 \pm 0.3</math></b>	<b><math>51.6 \pm 0.4</math></b>	
	5	6	7	8	9	All
DPP	$50.4 \pm 0.1$	$51.6 \pm 0.1$	$51.5 \pm 0.1$	$50.9 \pm 0.1$	$52.7 \pm 0.1$	$49.2 \pm 0.1$
UNIF	$52.4 \pm 0.1$	$54.6 \pm 0.1$	$55.1 \pm 0.1$	$53.3 \pm 0.1$	$56.2 \pm 0.1$	$51.6 \pm 0.1$
MEDOIDS	$52.4 \pm 0.0$	$54.4 \pm 0.1$	$55.1 \pm 0.1$	$53.2 \pm 0.1$	$56.1 \pm 0.1$	$51.0 \pm 0.1$
DPPNET MODE	<b><math>51.8 \pm 0.3</math></b>	<b><math>52.8 \pm 0.3</math></b>	<b><math>52.7 \pm 0.4</math></b>	<b><math>50.9 \pm 0.3</math></b>	<b><math>55.0 \pm 0.4</math></b>	$48.6 \pm 0.2$

attention. After hyper-parameter tuning, we found that the best architecture for this model consisted in 6 layers of 585 neurons each, showing that inhibitive attention is necessary to reduce the number of parameters in the model.

Strikingly, DPPNET performs significantly better than other baselines even on feature matrices drawn from a single class of digits (Table 3.2), despite distribution over input matrices during training being much less specialized. This implies that DPPNET sampling for dataset summarization may be leveraged to focus on sub-areas of datasets that are identified as areas of interest. Numerical results for CelebA and MovieLens are reported in Table 3.3, confirming the modeling ability of DPPNETs.

Table 3.3: NLLs on CelebA and MovieLens (mean  $\pm$  standard error); 20 samples of size 20 were drawn across 20 different feature matrices each for a total of 100 samples per method; DPPNET is the non-DPP model achieves the best NLLs.

DATASET	KERNEL	DPP	UNIFORM	$k$ -MEDOIDS	DPPNET Mode
CelebA	Exp. quadratic	$49.04 \pm 2.03$	$50.84 \pm 1.53$	$51.18 \pm 1.34$	<b><math>49.28 \pm 1.57</math></b>
MovieLens	Linear	$84.29 \pm 0.20$	$92.04 \pm 0.17$	$88.90 \pm 0.16$	<b><math>80.21 \pm 0.33</math></b>

### DppNet for kernel reconstruction

As a final experiment, we evaluate DPPNET’s performance on a downstream task for which DPPs have been shown to be useful: kernel reconstruction using the Nyström method (Nyström, 1930; Williams and Seeger, 2001). Given a positive semidefinite matrix  $\mathbf{K} \in \mathbb{R}^{N \times N}$ , the Nyström method approximates  $\mathbf{K}$  by  $\hat{\mathbf{K}} = \mathbf{K}_{:,S}\mathbf{K}_{S,S}^\dagger\mathbf{K}_{S,:}$  where  $\mathbf{K}^\dagger$  denotes the pseudoinverse of  $\mathbf{K}$ . The Nyström method is a popular method to scale up kernel methods and has found many applications in machine learning — see *e.g.*, (Bach and Jordan, 2002; Shen et al., 2009; Fowlkes et al., 2004; Talwalkar et al., 2013). The approximation quality directly depends on the choice of subset  $S$ .

Recently, DPPs have been shown to be a competitive approach for selecting  $S$  (Li et al., 2016c; Mariet et al., 2018). Following the approach of Li et al. (2016c), we evaluate the quality of the kernel reconstruction by learning a regression kernel  $\mathbf{K}$  on a training set, and reporting the prediction error on the test set using the Nyström reconstructed kernel  $\hat{\mathbf{K}}$ .

Specifically, we apply Kernel Ridge Regression to the Ailerons dataset<sup>2</sup> also used in (Li et al., 2016c); we subsample 1,000 points from each dataset and use an RBF kernel with bandwidth  $\beta$  and regularization parameter  $\lambda$  chosen by 10-fold cross-validation. Results are averaged over 10 random subsets of data, using the swapchain sampler initialized with  $k$ -means++. Additionally to the full DPP, we also compare DPPNET to the MCMC sampling method with quadrature acceleration (Li et al., 2016c,b). Figure 3-5 reports the resulting Root Mean Squared Error (RMSE) on the

<sup>2</sup><http://www.dcc.fc.up.pt/~ltorgo/Regression/DataSets.html>

withheld test set as well as the respective runtimes of each method.

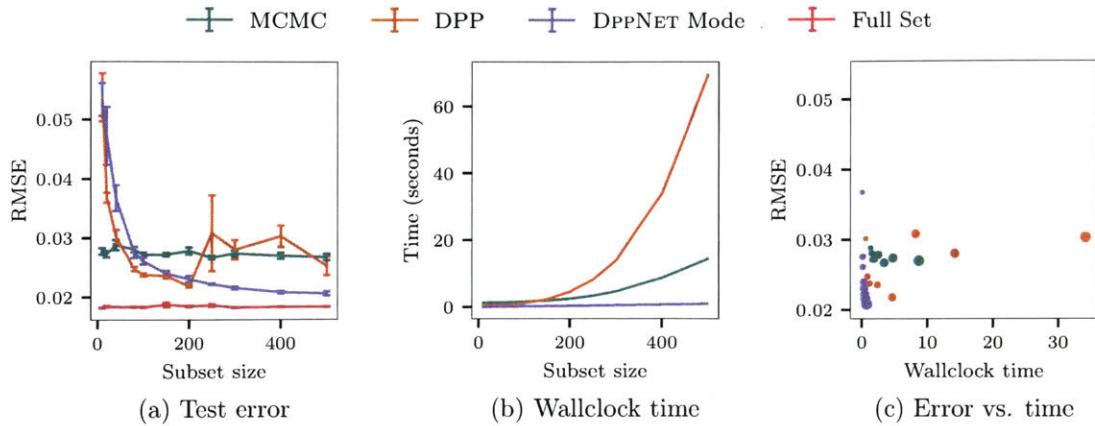


Figure 3-5: Results for the Nyström approximation experiments, comparing DPPNET to the fast MCMC sampling method of Li et al. (2016c) according to root mean squared error (RMSE) and wallclock time. Subsets selected by DPPNET achieve comparable and lower RMSE than a DPP and the MCMC method respectively while being significantly faster. In (c), the relative size of the marker represents the size of the sampled subset.

Fig. 3-5 confirms that DPPNET is a viable alternative to DPPs for kernel reconstructions. As expected, DPPNET is significantly faster than even approximate MCMC sampling. Note finally that while all methods were run on CPU, DPPNET is amenable to even further acceleration using GPUs.

### 3.3 Open problems

We close this chapter by mentioning some remaining open problems which we believe to be of interest, either for their theoretical ramifications or for potential applications.

- DPPNETs are *not* exchangeable: for a permutation  $\sigma$  of  $[k]$ , two sequences  $i_1, \dots, i_k$  and  $\sigma(i_1), \dots, \sigma(i_k)$  will not in general have the same probability under a DPPNET. Although non-exchangeability can be an asset when sampling a ranking of items, exchangeability can be enforced by leveraging previous work (Zaheer et al., 2017), and is a path worth investigating for DPPNETs.
- The scaling of the DPPNET's training time with the ground set size remains an

open question. However, standard tricks to enforce fixed-size ground sets such as sub-sampling from the dataset may be applied to DPPNETS. Similarly, if further are speedups necessary, sub-sampling from the ground set — a standard approach for DPP sampling over very large set sizes — can be combined with DPPNET sampling.

- In light of our results on dataset sampling, the question of whether encoders can be trained to produce encodings conducive to dataset summarization via DPPNETS seems of particular interest. Assuming knowledge of the relative diversity of a large quantity of subsets, an end-to-end training of the encoder and the DPPNET simultaneously may yield interesting results.
- Our analysis of DPPNETS extends further than DPPs: through Theorem 3.2.1 and Corollary 3.2.1.1, we have shown that under the right training conditions, generative models can inherit (log) submodularity, a crucial step in approximating submodular functions, with significant implications for discrete optimization and game theory. Understanding which additional properties of training distributions (including the larger class of negative dependence characterizations) may be conserved through careful training remains an open question which we believe to be significant to the community in general.
- Finally, the setup used to learn a DPP with neural networks can be generalized to learning any distribution over subsets, as long as marginal probabilities of adding an item to a subset can be computed efficiently. As optimizing and approximating set functions with intractable closed forms is a fundamental problem within combinatorial optimization, analyzing DPPNET-style approaches in this setting may prove to be a valuable research direction.

# Chapter 4

## Learning DPPs from data

In the previous chapter, we have seen that determinantal point processes provide an elegant framework for the efficient sampling of diverse, high-quality subsets from a ground set of items. The diversity and quality metrics for each application of DPPs are embedded in the chosen DPP’s kernel; this kernel can either be constructed based on expert knowledge of the required notions of diversity and quality, or alternatively, constructed based on feature representations of items in the ground set.

In many cases, however, the right DPP kernel may not be known *a priori*. For example, consider the setting investigated by Wilhelm et al. (2018): using DPPs to provide video (YouTube) recommendations drawn from a very large pool of possibilities. There is no obvious choice of kernel that will yield high-quality recommendations; however, the authors have access to sets of videos enjoyed by various users. The task, then, is to *learn* a DPP kernel that will assign high probabilities to such desirable subsets; this kernel can then be used to provide better recommendations.

Clearly, the ability to learn a DPP kernel from data is key to applying determinantal point processes; this chapter focuses on this learning problem, addressing in particular how to learn a kernel efficiently, and how to incorporate knowledge of low-quality subsets into the learning procedure.

We begin by formally introducing the learning problem for DPPs and discussing its properties. To address the high-complexity cost of learning the DPP kernel, we introduce a matrix structure that enables fast determinant and inverse computations.

Then, in order to incorporate prior knowledge of “undesirable” subsets to which the model should assign low probabilities, we introduce Contrastive Estimation: a novel optimization problem that incorporates this additional information. We show that in practice, DPPs learned with Contrastive Estimation have a better modeling performance on recommendation tasks.

## 4.1 Learning the DPP kernel

The standard<sup>1</sup> formulation of the learning problem for determinantal point processes is the *maximum likelihood estimation* (MLE) problem: given subsets  $S_1, \dots, S_n \subseteq [N]$ , we wish to learn a DPP kernel  $\mathbf{L}$  that maximizes the probability of observing the  $S_i$ . Formally, the MLE problem assumes knowledge of the ground set  $[N]$  and subsets  $S_1, \dots, S_n \subseteq [N]$  drawn in *i.i.d.* fashion, and is cast as the optimization problem

$$\text{Find } \mathbf{L} \in \underset{\mathbf{L} \succeq 0}{\operatorname{argmax}} \phi(\mathbf{L}) \triangleq \frac{1}{n} \sum_{i=1}^n \log \det(\mathbf{L}_{S_i}) - \log \det(\mathbf{I} + \mathbf{L}), \quad (4.1)$$

where  $\phi(\mathbf{L})$  is the log-likelihood of observing  $S_1, \dots, S_n$  under a DPP with kernel  $\mathbf{L}$ . The optimization problem 4.1 is non-convex in  $\mathbf{L}$ .

Problem 4.1 can be approached using stochastic gradient ascent (SGA), and is commonly augmented with further constraints on  $\mathbf{L}$ : most commonly, requiring that  $\mathbf{L}$  have a low rank. This choice imposes a clear limit on the modeling power of the DPP: any subset of size larger than the kernel rank has probability 0 under the DPP.

In the more general case where the kernel  $\mathbf{L}$  is learned without further constraints, an Expectation-Maximization scheme was suggested by Gillenwater et al. (2014), and a fixed-point (“Picard”) iteration, described by the simple update rule

$$\mathbf{L} \leftarrow \mathbf{L} + \mathbf{L} \nabla \phi(\mathbf{L}) \mathbf{L}, \quad (4.2)$$

was introduced by Mariet and Sra (2015). Among other advantages, both EM and the

---

<sup>1</sup>We will see later a less common formulation of the DPP learning problem.

Picard are less likely than SGA to learn near-diagonal kernels, which cannot capture negative interactions (Gillenwater et al., 2014).

## 4.2 Learning DPPs with Kronecker kernels

Although the  $\mathcal{O}(N^3)$  cost of exact sampling from a DPP has motivated a string of works on approximate sampling methods such as MCMC samplers (Kang, 2013; Li et al., 2016c), core-set based samplers (Li et al., 2015), and generative neural networks (Mariet et al., 2019b), the task of learning a DPP from data has received less attention. Methods in (Gillenwater et al., 2014; Mariet and Sra, 2015) cost  $\mathcal{O}(N^3)$  per iteration, which is clearly unacceptable for realistic settings. This burden is partially alleviated in (Gartrell et al., 2017) and (Osogami et al., 2018), wherein the authors learn low-rank DPPs; this comes at the expense of being unable to sample subsets larger than the chosen rank.

Unsurprisingly, the  $\mathcal{O}(N^3)$  cost stems from determinant computations and matrix inversions. To improve learning speed without sacrificing the ability to model large subset sizes, we choose to take advantage of a structure that significantly speeds up both these operations: Kronecker products.

Specifically, we consider DPPs with kernels that are the Kronecker product of smaller kernels  $\mathbf{L} = \mathbf{L}_1 \otimes \cdots \otimes \mathbf{L}_m$ ; each *sub-kernel*  $\mathbf{L}_i$  is a smaller positive definite matrix. We will see that, additionally to decreasing the number of parameters required to specify the DPP from  $\mathcal{O}(N^2)$  to  $\mathcal{O}(N^{2/m})$ , this structure also enables faster sampling and learning. The high-level comparison of Kronecker kernels to other structured and unstructured kernels is summarized in Table 4.1; the following material was first derived in (Mariet and Sra, 2016b).

Using Kronecker products to scale up matrix models is a popular and effective idea in several machine learning settings (Wu et al., 2005; Martens and Grosse, 2015; Flaxman et al., 2015; Zhang et al., 2015); however, we will see that its efficient execution for DPPs turns out to be more challenging.

MODEL	EXACT SAMPLING	STOCHASTIC LEARNING	
No structure	SGD	$\mathcal{O}(N^3 + \kappa^3)$	
	EM (Gillenwater et al., 2014)	$\mathcal{O}(N^3 + N\kappa^2)$	
	PICARD (Mariet and Sra, 2015)	$\mathcal{O}(N^3 + \kappa^3)$	
Low rank	$\mathcal{O}(Nd^2 + Ndk^2 + d^2k^3)$	SGA (Gartrell et al., 2017)	
		SGA (Osogami et al., 2018)	
Kronecker	2 kernels: $\mathcal{O}(N^{3/2} + N\kappa^2)$	KRK-PICARD	$\mathcal{O}(N^{3/2} + N\kappa^2)$
	3 kernels: $\mathcal{O}(N\kappa^2)$		

Table 4.1: Sampling and per-iteration MLE learning cost for DPPs;  $N$  is the size of the ground set,  $n$  the number of training subsets,  $\kappa$  the size of the largest training set, and  $d$  the dimension (if relevant) of the kernel.

### Preliminaries: Kronecker products

We begin by recalling some basic properties of Kronecker products we will need in our analysis; we omit proofs of these well-known results.

The Kronecker (tensor) product of a matrix  $\mathbf{A} \in \mathbb{R}^{p \times q}$  with  $\mathbf{B} \in \mathbb{R}^{r \times s}$  is defined as the  $pr \times qs$  block matrix  $\mathbf{A} \otimes \mathbf{B} = [a_{ij}\mathbf{B}]_{i,j=1}^{p,q}$ . We denote the  $r \times s$  block  $a_{ij}\mathbf{B}$  in  $\mathbf{A} \otimes \mathbf{B}$  by  $(\mathbf{A} \otimes \mathbf{B})_{(ij)}$  for any valid pair  $(i, j)$ , and extend the notation to non-Kronecker product matrices to indicate the submatrix of size  $r \times s$  at position  $(i, j)$ .

**Proposition 4.2.1.** *Let  $\mathbf{A}, \mathbf{B}, \mathbf{C}, \mathbf{D}$  be matrices of sizes so that  $\mathbf{AC}$  and  $\mathbf{BD}$  are well-defined. Then,*

- (i) *If  $\mathbf{A}, \mathbf{B} \succeq 0$  then  $\mathbf{A} \otimes \mathbf{B} \succeq 0$ .*
- (ii) *If  $\mathbf{A}$  and  $\mathbf{B}$  are invertible then so is  $\mathbf{A} \otimes \mathbf{B}$ , with  $(\mathbf{A} \otimes \mathbf{B})^{-1} = \mathbf{A}^{-1} \otimes \mathbf{B}^{-1}$ .*
- (iii)  *$(\mathbf{A} \otimes \mathbf{B})(\mathbf{C} \otimes \mathbf{D}) = (\mathbf{AC}) \otimes (\mathbf{BD})$ .*

An important consequence of Prop. 4.2.1(iii) is the following corollary.

**Corollary 4.2.1.1.** *Let  $\mathbf{A} = \mathbf{P}_A \mathbf{D}_A \mathbf{P}_A^\top$ ,  $\mathbf{B} = \mathbf{P}_B \mathbf{D}_B \mathbf{P}_B^\top$  be the eigenvector decompositions of  $\mathbf{A}$ ,  $\mathbf{B}$ . Then  $\mathbf{A} \otimes \mathbf{B}$  diagonalizes as  $(\mathbf{P}_A \otimes \mathbf{P}_B)(\mathbf{D}_A \otimes \mathbf{D}_B)(\mathbf{P}_A \otimes \mathbf{P}_B)^\top$ .*

We will also require partial trace operators:

**Definition 4.2.1.** Let  $\mathbf{A} \in \mathbb{R}^{N_1 N_2 \times N_1 N_2}$ . The *partial traces*  $\text{Tr}_1(\mathbf{A})$  and  $\text{Tr}_2(\mathbf{A})$  are defined as follows:

$$\text{Tr}_1(\mathbf{A}) := [\text{Tr}(\mathbf{A}_{(ij)})]_{1 \leq i,j \leq N_1} \in \mathbb{R}^{N_1 \times N_1}, \quad \text{Tr}_2(\mathbf{A}) := \sum_{i=1}^{N_1} \mathbf{A}_{(ii)} \in \mathbb{R}^{N_2 \times N_2}.$$

The action of partial traces is easy to visualize: indeed,  $\text{Tr}_1(\mathbf{A} \otimes \mathbf{B}) = \text{Tr}(\mathbf{B})\mathbf{A}$  and  $\text{Tr}_2(\mathbf{A} \otimes \mathbf{B}) = \text{Tr}(\mathbf{A})\mathbf{B}$ . For us, the most important property of partial trace operators is their positivity.

**Proposition 4.2.2.**  $\text{Tr}_1$  and  $\text{Tr}_2$  are positive operators, i.e., for  $\mathbf{A} \succ 0$ ,  $\text{Tr}_1(\mathbf{A}) \succ 0$  and  $\text{Tr}_2(\mathbf{A}) \succ 0$ .

### 4.2.1 Learning with alternating updates

We begin by considering the question of learning a Kronecker product kernel matrix from  $n$  observed subsets  $S_1, \dots, S_n$ , where the kernel is a Kronecker product of two smaller kernels:  $\mathbf{L} = \mathbf{L}_1 \otimes \mathbf{L}_2$ . For reasons that will become apparent, it suffices to two sub-kernels.

Recall that the MLE problem is given by

$$\max_{\mathbf{L} \succ 0} \phi(\mathbf{L}), \text{ where } \phi(\mathbf{L}) = \frac{1}{n} \sum_{i=1}^n (\log \det(\mathbf{L}_{S_i}) - \log \det(\mathbf{L} + \mathbf{I})). \quad (4.3)$$

Writing  $\mathbf{U}_i$  as the indicator matrix for  $S_i$  of size  $N \times |S_i|$  so that  $\mathbf{L}_{S_i} = \mathbf{U}_i^\top \mathbf{L} \mathbf{U}_i$ , the gradient of  $\phi$  is easily seen to be

$$\Delta := \nabla \phi(\mathbf{L}) = \frac{1}{n} \sum_{i=1}^n \mathbf{U}_i \mathbf{L}_{S_i}^{-1} \mathbf{U}_i^\top - (\mathbf{L} + \mathbf{I})^{-1}. \quad (4.4)$$

Mariet and Sra (2015) derived an iterative method (“the Picard iteration”) for computing an  $\mathbf{L}$  that solves  $\Delta = 0$  by running the simple iteration  $\mathbf{L} \leftarrow \mathbf{L} + \mathbf{L} \nabla \phi(\mathbf{L}) \mathbf{L}$  (Eq. 4.2). Crucially, this iteration preserves the necessary condition  $\mathbf{L} \succeq 0$ , and moreover guarantees to monotonically increase the log-likelihood  $\phi$ .

Unfortunately, these benefits accrue at a cost of  $\mathcal{O}(N^3)$  per iteration; furthermore,

a direct application of the Picard iteration to a kernel  $\mathbf{L} = \mathbf{L}_1 \otimes \cdots \otimes \mathbf{L}_n$  cannot guarantee that the Kronecker structure will be preserved.

With the aim of obtaining an efficient algorithm to (locally) optimize (4.3), we follow (Mariet and Sra, 2015): we first rewrite  $\phi$  as a function of  $\mathbf{S} = \mathbf{L}^{-1}$ , and re-arrange terms to write

$$\phi(\mathbf{S}) = \underbrace{\log \det(\mathbf{S})}_{f(\mathbf{S})} + \underbrace{\frac{1}{n} \sum_{i=1}^n \log \det(\mathbf{U}_i^\top \mathbf{S}^{-1} \mathbf{U}_i) - \log \det(\mathbf{I} + \mathbf{S})}_{g(\mathbf{S})}. \quad (4.5)$$

One easily verifies that  $f$  is concave, while a short argument shows that  $g$  is convex (Mariet and Sra, 2015). An appeal to the convex-concave procedure (Yuille and Rangarajan, 2002) then shows that updating  $\mathbf{S}$  by solving  $\nabla f(\mathbf{S}^{(k+1)}) + \nabla g(\mathbf{S}^{(k)}) = 0$ , which is exactly what (4.2) does (Mariet and Sra, 2015, Thm. 2.2), is guaranteed to monotonically increase  $\phi$ .

But for KRONDPP this idea does not apply so easily: due to the constraint  $\mathbf{L} = \mathbf{L}_1 \otimes \mathbf{L}_2$ , the function

$$g_\otimes : (\mathbf{S}_1, \mathbf{S}_2) \rightarrow \frac{1}{n} \sum_{i=1}^n \log \det(\mathbf{U}_i^\top (\mathbf{S}_1 \otimes \mathbf{S}_2)^{-1} \mathbf{U}_i) - \log \det(\mathbf{I} + \mathbf{S}_1 \otimes \mathbf{S}_2),$$

fails to be convex, precluding an easy generalization. Nevertheless, for fixed  $\mathbf{S}_1$  or  $\mathbf{S}_2$  the functions

$$\begin{cases} f_1 : \mathbf{S}_1 \mapsto f(\mathbf{S}_1 \otimes \mathbf{S}_2) \\ g_1 : \mathbf{S}_1 \mapsto g(\mathbf{S}_1 \otimes \mathbf{S}_2) \end{cases} \quad \text{and} \quad \begin{cases} f_2 : \mathbf{S}_2 \mapsto f(\mathbf{S}_1 \otimes \mathbf{S}_2) \\ g_2 : \mathbf{S}_2 \mapsto g(\mathbf{S}_1 \otimes \mathbf{S}_2) \end{cases}$$

are once again concave ( $f_i$ ) or convex ( $g_i$ ). Indeed, the map  $\otimes : \mathbf{S}_1 \rightarrow \mathbf{S}_1 \otimes \mathbf{S}_2$  is linear and  $f$  is concave, and  $f_1 = f \circ \otimes$  is also concave; similarly,  $f_2$  is seen to be concave and  $g_1$  and  $g_2$  are convex. Hence, by generalizing the arguments of (Yuille and Rangarajan, 2002, Thm. 2) to our “block-coordinate” setting, updating via

$$\nabla f_i(\mathbf{S}_i^{(k+1)}) = -\nabla g_i(\mathbf{S}_i^{(k)}), \quad \text{for } i = 1, 2, \quad (4.6)$$

should increase the log-likelihood  $\phi$  at each iteration.

We now prove that the updates defined in (4.6) exist, and that updating as per (4.6) ensure positive definiteness of the iterates as well as monotonic ascent.

**Proposition 4.2.3** (Positive definite iterates). *For  $\mathbf{S}_1 \succ 0$ ,  $\mathbf{S}_2 \succ 0$ , the solutions to (4.6) are given by the following expressions:*

$$\begin{aligned}\nabla f_1(\mathbf{X}) = -\nabla g_1(\mathbf{S}_1) &\iff \mathbf{X}^{-1} = \text{Tr}_1((\mathbf{I} \otimes \mathbf{S}_2)(\mathbf{L} + \mathbf{L}\Delta\mathbf{L}))/N_2 \\ \nabla f_2(\mathbf{X}) = -\nabla g_2(\mathbf{S}_2) &\iff \mathbf{X}^{-1} = \text{Tr}_2((\mathbf{S}_1 \otimes \mathbf{I})(\mathbf{L} + \mathbf{L}\Delta\mathbf{L}))/N_1.\end{aligned}$$

Moreover, these solutions are positive definite.

*Proof.* The high-level idea of the proof is the following: we know that  $\mathbf{L} \succ 0 \implies \mathbf{L} + \mathbf{L}\Delta\mathbf{L} \geq 0$ , because  $\mathbf{L} - \mathbf{L}(\mathbf{I} + \mathbf{L})^{-1}\mathbf{L} \succ 0$ . Since the partial trace operators are positive (Prop. 4.2.2), it follows that the solutions to (4.6) are also positive definite.

Specifically, write ‘vect’ for the operator that stacks columns of a matrix to form a vector; conversely, ‘mat’ takes a vector with  $k^2$  coefficients and returns a  $k \times k$  matrix.

Let  $\mathbf{L} = \mathbf{L}_1 \otimes \mathbf{L}_2$ ,  $\mathbf{S}_1 = \mathbf{L}_1^{-1}$ ,  $\mathbf{S}_2 = \mathbf{L}_2^{-1}$  and  $\mathbf{S} = \mathbf{S}_1 \otimes \mathbf{S}_2 = \mathbf{L}^{-1}$ . Let  $\mathbf{E}_{ij}$  be the matrix with all 0s except for a 1 at position  $(i, j)$ , its size being clear from context.

We wish to solve

$$\nabla f_2(\mathbf{X}) = -\nabla g_2(\mathbf{S}_1) \quad \text{and} \quad \nabla f_1(\mathbf{X}) = -\nabla g_1(\mathbf{S}_2). \quad (4.7)$$

It follows from the fact that  $\log \det(\mathbf{S}_1 \otimes \mathbf{S}_2) = N_2 \log \det \mathbf{S}_1 + N_1 \log \det \mathbf{S}_2$  that  $\nabla f_{\mathbf{S}_2}(\mathbf{X}) = N_2 \mathbf{X}^{-1}$  and  $\nabla f_{\mathbf{S}_1}(\mathbf{X}) = N_1 \mathbf{X}^{-1}$ . Moreover, we know that

$$\begin{aligned}\nabla g(\mathbf{S}) &= -(\mathbf{I} + \mathbf{S})^{-1} - \mathbf{S}^{-1} \frac{1}{n} \sum_i \mathbf{U}_i (\mathbf{U}_i^\top \mathbf{S}^{-1} \mathbf{U}_i)^{-1} \mathbf{U}_i \mathbf{S}^{-1} \\ &= -\mathbf{S}^{-1} - \mathbf{S}^{-1} \left( \frac{1}{n} \sum_i \mathbf{U}_i (\mathbf{U}_i^\top \mathbf{S}^{-1} \mathbf{U}_i)^{-1} \mathbf{U}_i - (\mathbf{I} + \mathbf{S}^{-1})^{-1} \right) \mathbf{S}^{-1} \\ &= -(\mathbf{L} + \mathbf{L}\Delta\mathbf{L}).\end{aligned}$$

The Jacobian of  $\mathbf{S}_1 \rightarrow \mathbf{S}_1 \otimes \mathbf{S}_2$  is  $\mathbf{J} = \left( \text{vect}(\mathbf{E}_{11} \otimes \mathbf{S}_2), \dots, \text{vect}(\mathbf{E}_{N_1 N_1} \otimes \mathbf{S}_2) \right)$ . Hence,

$$\begin{aligned} \nabla f_1(\mathbf{X})_{ij} = -(\nabla g_1(\mathbf{S}_1))_{ij} &\iff N_2 \mathbf{X}_{ij}^{-1} = (\mathbf{J}^\top \text{vect}(-\nabla g(\mathbf{S})))_{ij} \\ &\iff N_2 \mathbf{X}_{ij}^{-1} = \text{vect}(\mathbf{E}_{ij} \otimes \mathbf{S}_2)^\top \text{vect}(\mathbf{L} + \mathbf{L}\Delta\mathbf{L}) \\ &\iff N_2 \mathbf{X}_{ij}^{-1} = \text{Tr}((\mathbf{E}_{ij} \otimes \mathbf{S}_2)(\mathbf{L} + \mathbf{L}\Delta\mathbf{L})) \\ &\iff N_2 \mathbf{X}_{ij}^{-1} = \text{Tr}(\mathbf{S}_2(\mathbf{L} + \mathbf{L}\Delta\mathbf{L})_{(ij)}) \\ &\iff N_2 \mathbf{X}_{ij}^{-1} = \text{Tr}(((\mathbf{I} \otimes \mathbf{S}_2)(\mathbf{L} + \mathbf{L}\Delta\mathbf{L}))_{(ij)}). \end{aligned}$$

The last equivalence is obtained by manipulating indices. Thus, we have

$$\nabla f_2(\mathbf{X}) = -\nabla g_2(\mathbf{S}_1) \iff \mathbf{X}^{-1} = \frac{1}{N_2} \text{Tr}_1 ((\mathbf{I} \otimes \mathbf{S}_2)(\mathbf{L} + \mathbf{L}\Delta\mathbf{L})).$$

Similarly, by setting  $\mathbf{J}' = (\text{vect}(\mathbf{S}_1 \otimes \mathbf{E}_{11}), \dots, \text{vect}(\mathbf{S}_1 \otimes \mathbf{E}_{N_1 N_1}))$ , we have

$$\begin{aligned} \nabla f_2(\mathbf{X})_{ij} = -(\nabla g_2(\mathbf{S}_2))_{ij} &\iff N_1 \mathbf{X}_{ij}^{-1} = (\mathbf{J}'^\top \text{vect}(-\nabla g(\mathbf{S})))_{ij} \\ &\iff N_1 \mathbf{X}_{ij}^{-1} = \text{vect}(\mathbf{S}_1 \otimes \mathbf{E}_{ij})^\top \text{vect}(\mathbf{L} + \mathbf{L}\Delta\mathbf{L}) \\ &\iff N_1 \mathbf{X}_{ij}^{-1} = \text{Tr}((\mathbf{S}_1 \otimes \mathbf{E}_{ij})(\mathbf{L} + \mathbf{L}\Delta\mathbf{L})) \\ &\iff N_1 \mathbf{X}_{ij}^{-1} = \left( \sum_{k,\ell=1}^{N_1} \mathbf{S}_{1k\ell} (\mathbf{L} + \mathbf{L}\Delta\mathbf{L})_{(\ell k)} \right)_{ij} \\ &\iff N_1 \mathbf{X}_{ij}^{-1} = \left( \sum_{\ell=1}^{N_1} ((\mathbf{S}_1 \otimes \mathbf{I})(\mathbf{L} + \mathbf{L}\Delta\mathbf{L}))_{(\ell\ell)} \right)_{ij}. \end{aligned}$$

Hence, we have the equivalence

$$\nabla f_{\mathbf{S}_1}(\mathbf{X}) = -\nabla g_{\mathbf{S}_1}(\mathbf{S}_2) \iff \mathbf{X}^{-1} = \frac{1}{N_1} \text{Tr}_2 ((\mathbf{S}_1 \otimes \mathbf{I})(\mathbf{L} + \mathbf{L}\Delta\mathbf{L})),$$

which proves Prop. 4.2.3.  $\square$

We are now ready to establish that these updates ensure monotonic ascent in the log-likelihood.

**Theorem 4.2.4** (Ascent). *Starting with  $\mathbf{L}_1^{(0)}, \mathbf{L}_2^{(0)} \in \mathbb{S}^{++}$ , updating according to (4.6) generates  $\mathbf{L}_1^{(k)}, \mathbf{L}_2^{(k)} \in \mathbb{S}^{++}$ , and the sequence  $\{\phi(\mathbf{L}_1^{(k)} \otimes \mathbf{L}_2^{(k)})\}_{k \geq 0}$  is non-decreasing.*

*Proof.* Updating according to (4.6) generates positive definite matrices  $\mathbf{S}_i$ , and hence positive definite subkernels  $\mathbf{L}_i = \mathbf{S}_i$ . Moreover, due to the convexity of  $g_1$  and concavity of  $f_1$ , for matrices  $\mathbf{A}, \mathbf{B} \succ 0$

$$\begin{aligned} f_1(\mathbf{B}) &\leq f_1(\mathbf{A}) + \nabla f_1(\mathbf{A})^\top (\mathbf{B} - \mathbf{A}), \\ g_1(\mathbf{A}) &\geq g_1(\mathbf{B}) + \nabla g_1(\mathbf{B})^\top (\mathbf{A} - \mathbf{B}). \end{aligned}$$

Hence,  $f_1(\mathbf{A}) + g_1(\mathbf{A}) \geq f_1(\mathbf{B}) + g_1(\mathbf{B}) + (\nabla f_1(\mathbf{A}) + \nabla g_1(\mathbf{B}))^\top (\mathbf{A} - \mathbf{B})$ .

Thus, if  $\mathbf{S}_1^{(k)}, \mathbf{S}_1^{(k+1)}$  verify (4.6), by setting  $\mathbf{A} = \mathbf{S}_1^{(k+1)}$  and  $\mathbf{B} = \mathbf{S}_1^{(k)}$  we have

$$\phi(\mathbf{L}_1^{(k+1)} \otimes \mathbf{L}_2^{(k)}) = f_1(\mathbf{S}_1^{(k+1)}) + g_1(\mathbf{S}_1^{(k+1)}) \geq f_1(\mathbf{S}_1^{(k)}) + g_1(\mathbf{S}_1^{(k)}) = \phi(\mathbf{L}_1^{(k)} \otimes \mathbf{L}_2^{(k)}).$$

The same reasoning holds for  $\mathbf{L}_2$ , which proves the theorem.  $\square$

As  $\text{Tr}_1((\mathbf{I} \otimes \mathbf{S}_2)\mathbf{L}) = N_2\mathbf{L}_1$  (and similarly for  $\mathbf{L}_2$ ), update (4.6) is equivalent to

$$\mathbf{L}_1 \leftarrow \mathbf{L}_1 + \text{Tr}_1((\mathbf{I} \otimes \mathbf{L}_2^{-1})(\mathbf{L}\Delta\mathbf{L})) / N_2, \quad \mathbf{L}_2 \leftarrow \mathbf{L}_2 + \text{Tr}_2((\mathbf{L}_1^{-1} \otimes \mathbf{I})(\mathbf{L}\Delta\mathbf{L})) / N_1.$$

The arguments above easily extend to the multiblock case: when learning  $\mathbf{L} = \mathbf{L}_1 \otimes \dots \otimes \mathbf{L}_m$ , we update  $\mathbf{L}_k$  as

$$(\mathbf{L}_k)_{ij} \leftarrow (\mathbf{L}_k)_{ij} + \frac{N_k}{N_1 \dots N_m} \text{Tr}[(\mathbf{L}_1 \otimes \dots \otimes \mathbf{L}_{k-1} \otimes \mathbf{E}_{ij} \otimes \mathbf{L}_{k+1} \otimes \dots \otimes \mathbf{L}_m)(\mathbf{L}\Delta\mathbf{L})].$$

**Generalization.** As in (Mariet and Sra, 2015), we can generalize the updates to take an additional step-size parameter  $a > 0$ :

$$\begin{aligned} \mathbf{L}_1 &\leftarrow \mathbf{L}_1 + a \text{Tr}_1((\mathbf{I} \otimes \mathbf{L}_2^{-1})(\mathbf{L}\Delta\mathbf{L})) / N_2, \\ \mathbf{L}_2 &\leftarrow \mathbf{L}_2 + a \text{Tr}_2((\mathbf{L}_1^{-1} \otimes \mathbf{I})(\mathbf{L}\Delta\mathbf{L})) / N_1. \end{aligned}$$

Experimentally,  $a > 1$  (as long as the updates remain in  $\mathbb{S}^+$ ) can provide faster convergence, although the monotonicity of the log-likelihood is no longer guaranteed.

From Theorem 4.2.4, we obtain Algorithm 7. However, from the above updates it is not obvious whether the Kronecker product saves us any computation. In particular,

it is not clear whether the updates can be implemented to run faster than  $\mathcal{O}(N^3)$ .

---

**Algorithm 7** KRK-PICARD iteration

---

**Input:** Matrices  $\mathbf{L}_1, \mathbf{L}_2$ , training set  $T$ , parameter  $a$ .  
**for**  $i = 1$  **to** maxIter **do**  
     $\mathbf{L}_1 \leftarrow \mathbf{L}_1 + a \text{Tr}_1((\mathbf{I} \otimes \mathbf{L}_2^{-1})(\mathbf{L}\Delta\mathbf{L})) / N_2$    // or update stochastically  
     $\mathbf{L}_2 \leftarrow \mathbf{L}_2 + a \text{Tr}_2((\mathbf{L}_1^{-1} \otimes \mathbf{I})(\mathbf{L}\Delta\mathbf{L})) / N_1$    // or update stochastically  
**return**  $(\mathbf{L}_1, \mathbf{L}_2)$

---

**Theorem 4.2.5** (Complexity). *For  $N_1 \approx N_2 \approx \sqrt{N}$ , the updates in Algorithm 7 can be computed in  $\mathcal{O}(n\kappa^3 + N^2)$  time and  $\mathcal{O}(N^2)$  space, where  $\kappa$  is the size of the largest training subset. Furthermore, stochastic updates can be computed in  $\mathcal{O}(N\kappa^2 + N^{3/2})$  time and  $\mathcal{O}(N + \kappa^2)$  space.*

*Proof.* The updates to  $\mathbf{L}_1$  and  $\mathbf{L}_2$  are obtained efficiently through different methods; hence, the proof to Thm. 4.2.5 is split into two sections. We write

$$\boldsymbol{\Theta} = \frac{1}{n} \sum_{i=1}^n \mathbf{U}_i \mathbf{L}_{S_i}^{-1} \mathbf{U}_i^\top \quad (\text{or } \boldsymbol{\Theta} = \mathbf{U}_i \mathbf{L}_{S_i}^{-1} \mathbf{U}_i^\top \text{ for stochastic updates})$$

so that  $\Delta = \boldsymbol{\Theta} - (\mathbf{I} + \mathbf{L})^{-1}$ . Recall that  $(\mathbf{A} \otimes \mathbf{B})_{(ij)} = a_{ij} \mathbf{B}$ .

**Updating  $\mathbf{L}_1$ .** We wish to compute  $\mathbf{X} = \text{Tr}_1((\mathbf{I} \otimes \mathbf{L}_2^{-1})(\mathbf{L}\Delta\mathbf{L}))$ . We have

$$\begin{aligned} \mathbf{X}_{ij} &= \text{Tr} [((\mathbf{I} \otimes \mathbf{L}_2^{-1})(\mathbf{L}\Delta\mathbf{L}))_{(ij)}] \\ &= \text{Tr} [\mathbf{L}_2^{-1}(\mathbf{L}\Delta\mathbf{L})_{(ij)}] \\ &= \text{Tr} \left[ \mathbf{L}_2^{-1} \sum_{k,\ell=1}^{N_1} \mathbf{L}_{(ik)} \Delta_{(k\ell)} \mathbf{L}_{(\ell j)} \right] \\ &= \sum_{k,\ell=1}^{N_1} \mathbf{L}_{1ik} \mathbf{L}_{1\ell j} \text{Tr}(\mathbf{L}_2^{-1} \mathbf{L}_2 \Delta_{(k\ell)} \mathbf{L}_2) \\ &= \sum_{k,\ell=1}^{N_1} \mathbf{L}_{1ik} \mathbf{L}_{1\ell j} \underbrace{\text{Tr}(\boldsymbol{\Theta}_{(k\ell)} \mathbf{L}_2)}_{\mathbf{A}_{k\ell}} - \underbrace{\text{Tr}((\mathbf{I} + \mathbf{L})_{(k\ell)}^{-1} \mathbf{L}_2)}_{\mathbf{B}_{k\ell}} \\ &= (\mathbf{L}_1 \mathbf{A} \mathbf{L}_1 - \mathbf{L}_1 \mathbf{B} \mathbf{L}_1)_{ij}. \end{aligned}$$

The  $N_1 \times N_1$  matrix  $\mathbf{A}$  can be computed in  $\mathcal{O}(n\kappa^3 + N_1^2 N_2^2)$  time simply by pre-computing  $\boldsymbol{\Theta}$  in  $\mathcal{O}(n\kappa^3)$  and then computing all  $N_1^2$  traces in  $\mathcal{O}(N_2^2)$  time. When

doing stochastic updates for which  $\Theta$  is sparse with only  $\kappa^2$  non-zero coefficients, computing  $\mathbf{A}$  can be done in  $\mathcal{O}(N_1^2\kappa^2 + \kappa^3)$ .

By diagonalizing  $\mathbf{L}_1 = \mathbf{P}_1 \mathbf{D}_1 \mathbf{P}_1^\top$  and  $\mathbf{L}_2 = \mathbf{P}_2 \mathbf{D}_2 \mathbf{P}_2^\top$ , we have  $(\mathbf{I} + \mathbf{L})^{-1} = \mathbf{P} \mathbf{D} \mathbf{P}^\top$  with  $\mathbf{P} = \mathbf{P}_1 \otimes \mathbf{P}_2$  and  $\mathbf{D} = (\mathbf{I} + \mathbf{D}_1 \otimes \mathbf{D}_2)^{-1}$ .  $\mathbf{P}_1, \mathbf{P}_2, \mathbf{D}_1, \mathbf{D}_2$  and  $\mathbf{D}$  can all be obtained in  $\mathcal{O}(N_1^3 + N_2^3 + N_1 N_2)$  as a consequence of Prop. 4.2.1. Then

$$\begin{aligned}\mathbf{B}_{ij} &= \text{Tr}((\mathbf{I} + \mathbf{L})_{(ij)}^{-1} \mathbf{L}_2) \\ &= \sum_k \text{Tr}(\mathbf{P}_{(ik)} \mathbf{D}_{(kk)} \mathbf{P}_{(kj)}^\top \mathbf{L}_2) \\ &= \sum_k \mathbf{P}_{1ik} \mathbf{P}_{1jk} \text{Tr}(\mathbf{P}_2 \mathbf{D}_{(kk)} \mathbf{P}_2^\top \mathbf{P}_2 \mathbf{D}_2 \mathbf{P}_2^\top) \\ &= \sum_k \mathbf{P}_{1ik} \mathbf{P}_{1jk} \underbrace{\text{Tr}(\mathbf{D}_{(kk)} \mathbf{D}_2)}_{\alpha_k}.\end{aligned}$$

Let  $\widehat{\mathbf{D}} = \text{diag}(\alpha_1, \dots, \alpha_{N_1})$ , which can be computed in  $\mathcal{O}(N_1 N_2)$ . Then  $\mathbf{L}_1 \mathbf{B} \mathbf{L}_1 = \mathbf{P}_1 \mathbf{D}_1 \widehat{\mathbf{D}} \mathbf{D}_1 \mathbf{P}_1$  is computable in  $\mathcal{O}(N_1^3 + N_2^3)$ .

Overall, the update to  $\mathbf{L}_1$  can be computed in  $\mathcal{O}(n\kappa^3 + N_1^2 N_2^2 + N_1^3 + N_2^3)$  time, or in  $\mathcal{O}(N_1^2 \kappa^2 + \kappa^3 + N_1^3 + N_2^3)$  time if the updates are stochastic. Moreover, if  $\Theta$  is sparse with only  $z$  non-zero coefficients (for stochastic updates  $z = \kappa$ ),  $\mathbf{A}$  can be computed in  $\mathcal{O}(\kappa^2)$  space, leading to an overall  $\mathcal{O}(z^2 + N_1^2 + N_2^2)$  memory cost.

**Updating  $\mathbf{L}_2$ .** We wish to compute  $\mathbf{X} = \text{Tr}_2[(\mathbf{L}_1^{-1} \otimes \mathbf{I})(\mathbf{L} \Delta \mathbf{L})]$  efficiently.

$$\begin{aligned}\mathbf{X} &= \sum_{i=1}^{N_1} ((\mathbf{L}_1^{-1} \otimes \mathbf{I})(\mathbf{L} \Delta \mathbf{L}))_{(ii)} \\ &= \sum_{i=1}^{N_1} ((\mathbf{I} \otimes \mathbf{L}_2)(\Theta - (\mathbf{I} + \mathbf{L})^{-1})(\mathbf{L}_1 \otimes \mathbf{L}_2))_{(ii)} \\ &= \sum_{i,j=1}^{N_1} \mathbf{L}_{1ij} \mathbf{L}_2 \Theta_{(ij)} \mathbf{L}_2 - \sum_{i=1}^{N_1} ((\mathbf{I} \otimes \mathbf{L}_2)(\mathbf{I} + \mathbf{L})^{-1}(\mathbf{L}_1 \otimes \mathbf{L}_2))_{(ii)} \\ &= \underbrace{\mathbf{L}_2 \sum_{i,j=1}^{N_1} \mathbf{L}_{1ij} \Theta_{(ij)} \mathbf{L}_2}_{A} - \underbrace{\sum_{i=1}^{N_1} ((\mathbf{I} \otimes \mathbf{L}_2)(\mathbf{I} + \mathbf{L})^{-1}(\mathbf{L}_1 \otimes \mathbf{L}_2))_{(ii)}}_{B}.\end{aligned}$$

Matrix  $\mathbf{A}$  can be computed in  $\mathcal{O}(n\kappa^3 + N_1^2 N_2^2 + N_2^3)$  time. For stochastic updates, we only require  $\mathcal{O}(N_1^2 \kappa^2 + \kappa^3 + N_2^3)$  time and  $\mathcal{O}(N_2^2 + N_1^2 + \kappa^2)$  space due to  $\Theta$ 's sparsity.

Regarding  $\mathbf{B}$ , as all matrices commute, we can write

$$(\mathbf{I} \otimes \mathbf{L}_2)(\mathbf{I} + \mathbf{L})^{-1}(\mathbf{L}_1 \otimes \mathbf{L}_2) = (\mathbf{P}_1 \otimes \mathbf{P}_2)\Lambda(\mathbf{P}_1 \otimes \mathbf{P}_2),$$

where  $\Lambda = (\mathbf{I} \otimes \mathbf{D}_2)(\mathbf{I} + \mathbf{D}_1 \otimes \mathbf{D}_2)^{-1}(\mathbf{D}_1 \otimes \mathbf{D}_2)$  is diagonal and is obtained in  $\mathcal{O}(N_1^3 + N_2^3 + N_1 N_2)$ . Moreover,

$$\mathbf{B} = \sum_{i=1}^{N_1} (\mathbf{P}\Lambda\mathbf{P}^\top)_{(ii)} = \mathbf{P}_2 \left( \sum_{i,k=1}^{N_1} \mathbf{P}_{1ik} \Lambda_{(kk)} \mathbf{P}_{1ik} \right) \mathbf{P}_2^\top,$$

which allows us to compute  $\mathbf{B}$  in  $\mathcal{O}(N_1^2 N_2 + N_2^3 + N_1^3)$  total.

Overall, we can obtain  $\mathbf{X}$  in  $\mathcal{O}(n\kappa^3 + N_1^2 N_2^2 + N_1^3 + N_2^3)$  or in  $\mathcal{O}(N_1^2 \kappa^2 + N_1^2 N_2 + N_1^3 + N_2^3)$  for stochastic updates, in which case only  $\mathcal{O}(N_1^2 + N_2^2 + \kappa^2)$  space is necessary.  $\square$

The proof shows that by leveraging the properties of the Kronecker product, the maximum likelihood updates can be obtained without computing  $\mathbf{L}\Delta\mathbf{L}$ . This result is non-trivial: the components of  $\Delta$ ,  $\frac{1}{n} \sum_i \mathbf{U}_i \mathbf{L}_{S_i}^{-1} \mathbf{U}_i^\top$  and  $(\mathbf{I} + \mathbf{L})^{-1}$ , must be considered separately for computational efficiency.

The runtime of Alg. 7 is a marked improvement over (Mariet and Sra, 2015), which runs in  $\mathcal{O}(N^2)$  space and  $\mathcal{O}(n\kappa^3 + N^3)$  time in the non-stochastic case, and in  $\mathcal{O}(N^3)$  time for stochastic updates.

#### 4.2.2 Learning with joint updates

An alternate possibility lies in updating  $\mathbf{L}_1$  and  $\mathbf{L}_2$  jointly via the following straightforward procedure: we first update  $\mathbf{L} \leftarrow \mathbf{L} + \mathbf{L}\Delta\mathbf{L}$ , then recover the Kronecker structure of the kernel by defining the updates  $\mathbf{L}'_1$  and  $\mathbf{L}'_2$  such that:

$$\begin{cases} (\mathbf{L}'_1, \mathbf{L}'_2) \text{ minimizes } \|\mathbf{L} + \mathbf{L}\Delta\mathbf{L} - \mathbf{L}'_1 \otimes \mathbf{L}'_2\|_F^2 \\ \mathbf{L}'_1 \succ 0, \mathbf{L}'_2 \succ 0, \|\mathbf{L}'_1\| = \|\mathbf{L}'_2\|. \end{cases} \quad (4.8)$$

**Theorem 4.2.6.** *Solutions to (4.8) exist and can be computed from the first singular value and vectors of the matrix  $[\text{vect}((\mathbf{L}^{-1} + \Delta)_{(ij)})^\top]_{i,j=1}^{N_1}$ .*

*Proof.* In order to minimize the number of matrix multiplications, we equivalently (due to the properties of the Frobenius norm) find

$$\mathbf{X}, \mathbf{Y} \in \operatorname{argmin}_{\mathbf{X}, \mathbf{Y} \geq 0} \|\mathbf{L}^{-1} + \Delta - \mathbf{X} \otimes \mathbf{Y}\|_F^2, \quad (4.9)$$

then let  $\mathbf{L}'_1 \leftarrow \mathbf{L}_1 \mathbf{X} \mathbf{L}_1$  and  $\mathbf{L}'_2 \leftarrow \mathbf{L}_2 \mathbf{Y} \mathbf{L}_2$ . We will first require the following lemma.

**Lemma 4.2.7.** *Let  $\mathbf{L} \succ 0$ . Define  $\mathbf{R} := [\operatorname{vect}(\mathbf{L}_{(11)})^\top; \dots; \operatorname{vect}(\mathbf{L}_{(N_1 N_1)})^\top]_{i,j=1}^{N_1} \in \mathbb{R}^{N_1 N_1 \times N_2 N_2}$ . Suppose that  $\mathbf{R}$  has an eigengap between its largest singular value and the next, and let  $u, v, \sigma$  be the first singular vectors and value of  $\mathbf{R}$ . Let  $\mathbf{U} = \operatorname{mat}(u)$  and  $\mathbf{V} = \operatorname{mat}(v)$ . Then  $\mathbf{U}$  and  $\mathbf{V}$  are either both positive definite or negative definite.*

*Moreover, for any value  $\alpha \neq 0$ , the pair  $(\alpha \mathbf{U}, \sigma/\alpha \mathbf{V})$  minimizes  $\|\mathbf{L} - X \otimes Y\|_F^2$ .*

Lemma 4.2.7 is a consequence of (Loan and Pitsianis, 1993, Thm. 11), and shows that if  $\mathbf{L}$  is initially positive definite, setting the sign of  $\alpha$  based on whether  $\mathbf{U}$  and  $\mathbf{V}$  are positive or negative definite,<sup>2</sup> and updating

$$\begin{cases} \mathbf{L}_1 \leftarrow \alpha \mathbf{L}_1 \mathbf{U} \mathbf{L}_1 \\ \mathbf{L}_2 \leftarrow \sigma/\alpha \mathbf{L}_2 \mathbf{V} \mathbf{L}_2 \end{cases}$$

will maintain positive definite iterates. Given that if  $\mathbf{L}_1 \succ 0$  and  $\mathbf{L}_2 \succ 0$ , we also have  $\mathbf{L}_1 \otimes \mathbf{L}_2 \succ 0$ , a simple induction then shows that by choosing an initial kernel estimate  $\mathbf{L} \succ 0$ , subsequent values of  $\mathbf{L}$  will remain positive definite. Hence, by choosing  $\alpha$  such that the new estimates  $\mathbf{L}_1$  and  $\mathbf{L}_2$  verify  $\|\mathbf{L}_1\| = \|\mathbf{L}_2\|$ , we verify all the conditions of Eq. 4.8.  $\square$

Theorem 4.2.7 leads to Algorithm 8: a straightforward iteration for learning matrices  $\mathbf{L}_1$  and  $\mathbf{L}_2$  based on the decomposition of the Picard estimate  $\mathbf{L} + \mathbf{L} \nabla \phi(\mathbf{L}) \mathbf{L}$  as a Kronecker product.

---

<sup>2</sup>This can easily be done simply by checking the sign of the first diagonal coefficient of  $\mathbf{U}$ , which will be positive if and only if  $\mathbf{U} \succ 0$ .

---

**Algorithm 8** JOINT-PICARD iteration

---

**Input:** Matrices  $\mathbf{L}_1, \mathbf{L}_2$ , training set  $T$ , step-size  $a \geq 1$ .

**for**  $i = 1$  **to** maxIter **do**

$\sigma, \mathbf{U}, \mathbf{V} \leftarrow \text{power\_method}(\mathbf{L}^{-1} + \Delta)$  to obtain the first singular value + vectors.

$\alpha \leftarrow \text{sgn}(\mathbf{U}_{11}) \sqrt{\sigma \|\mathbf{L}_2 \mathbf{V} \mathbf{L}_2\| / \|\mathbf{L}_1 \mathbf{U} \mathbf{L}_1\|}$

$\mathbf{L}_1 \leftarrow \mathbf{L}_1 + a(\alpha \mathbf{L}_1 \mathbf{U} \mathbf{L}_1 - \mathbf{L}_1)$

$\mathbf{L}_2 \leftarrow \mathbf{L}_2 + a(\sigma/\alpha \mathbf{L}_2 \mathbf{V} \mathbf{L}_2)$

**return**  $(\mathbf{L}_1, \mathbf{L}_2)$

---

However, Alg. 8 produces no guaranteed increase in log-likelihood. An analysis similar that of Thm. 4.2.5 shows that updates can be obtained  $\mathcal{O}(n\kappa^3 + \max(N_1, N_2)^4)$ .

### 4.2.3 The knapsack problem for kernel learning

Although KRONDPPS have tractable learning algorithms, the memory requirements remain high for non-stochastic updates, as the matrix  $\Theta = \frac{1}{n} \sum_i \mathbf{U}_i \mathbf{L}_{S_i}^{-1} \mathbf{U}_i^\top$  needs to be stored, requiring  $\mathcal{O}(N^2)$  memory.

However, if the training set  $\{S_1, \dots, S_n\}$  can be partitioned as  $P_1 \cup \dots \cup P_k$  with

$$\{S_1, \dots, S_n\} = \bigcup_{k=1}^m P_k \quad s.t. \quad \forall k, |\bigcup_{S \in P_k} S| < z, \quad (4.10)$$

$\Theta$  can be decomposed as  $\frac{1}{n} \sum_{k=1}^m \Theta_k$  with  $\Theta_k = \sum_{S_i \in P_k} \mathbf{U}_i \mathbf{L}_{S_i}^{-1} \mathbf{U}_i^\top$ . Due to the bound in Eq. 4.10, each  $\Theta_k$  will be sparse, with only  $z^2$  non-zero coefficients. We can then store each  $\Theta_k$  with minimal storage and update  $\mathbf{L}_1$  and  $\mathbf{L}_2$  in  $\mathcal{O}(n\kappa^3 + mz^2 + N^{3/2})$  time and  $\mathcal{O}(mz^2 + N)$  space.

Determining the existence of such a partition of size  $m$  is a variant of the NP-Hard Subset-Union Knapsack Problem (SUKP) (Goldschmidt et al., 1994) with  $m$  knapsacks and where the value of each item (*i.e.*, each  $S_i$ ) is equal to 1: a solution to SUKP of value  $n$  with  $m$  knapsacks is equivalent to a solution to Eq. 4.10. However, an approximate partition can also be simply constructed via a greedy algorithm.

#### 4.2.4 Sampling

Sampling exactly (Alg. 3) from a full DPP kernel costs  $\mathcal{O}(N^3 + Nk^3)$  where  $k$  is the size of the sampled subset, with the bulk of the computation lying in the initial eigendecomposition of  $\mathbf{L}$ ; the  $k$  orthonormalizations cost  $\mathcal{O}(Nk^3)$ .

It follows from Prop. 4.2.1.1 that for **KRONDPPS**, the eigenvalues  $\lambda_i$  can be obtained in  $\mathcal{O}(N_1^3 + N_2^3)$ , and the  $k$  eigenvectors in  $\mathcal{O}(kN)$  operations. For  $N_1 \approx N_2 \approx \sqrt{N}$ , exact sampling thus only costs  $\mathcal{O}(N^{3/2} + Nk^3)$ . If  $\mathbf{L} = \mathbf{L}_1 \otimes \mathbf{L}_2 \otimes \mathbf{L}_3$ , the same reasoning shows that exact sampling is linear in  $N$ , only requiring  $\mathcal{O}(Nk^3)$  operations.

#### 4.2.5 Experimental results

In order to validate our learning algorithm, we compare **KRK-PICARD** to **JOINT-PICARD** and to the Picard iteration (**PICARD**) on real and synthetic datasets.<sup>3</sup>

##### Synthetic tests

To enable a fair comparison between algorithms, we test them on synthetic data drawn from a full (non-Kronecker) ground-truth DPP kernel. The sub-kernels were initialized by  $\mathbf{L}_i = \mathbf{X}^\top \mathbf{X}$ , with  $\mathbf{X}$ 's coefficients drawn uniformly from  $[0, \sqrt{2}]$ ; for **PICARD**,  $\mathbf{L}$  was initialized with  $\mathbf{L}_1 \otimes \mathbf{L}_2$ .

For Figures 4-1a and 4-1b, training data was generated by sampling 100 subsets from the true kernel with sizes uniformly distributed between 10 and 190.

To evaluate **KRK-PICARD** on matrices too large to fit in memory and with large  $\kappa$ , we drew samples from a  $50 \cdot 10^3 \times 50 \cdot 10^3$  kernel of rank 1,000 (on average  $|S_i| \approx 1,000$ ), and learned the kernel stochastically (only **KRK-PICARD** could be run due to the memory requirements of other methods); the likelihood drastically improves in only two steps (Fig.4-1c).

As shown in Figures 4-1a and 4-1b, **KRK-PICARD** converges significantly faster than **PICARD**, especially for large values of  $N$ . However, although **JOINT-PICARD** also increases the log-likelihood at each iteration, it converges much slower and has a high

---

<sup>3</sup>Experiments were repeated 5 times and averaged, using MATLAB on a Linux Mint system with 16GB of RAM and an i7-4710HQ CPU @ 2.50GHz.

standard deviation, whereas the standard deviations for PICARD and KRK-PICARD are barely noticeable. For these reasons, we drop the comparison to JOINT-PICARD in the subsequent experiments.

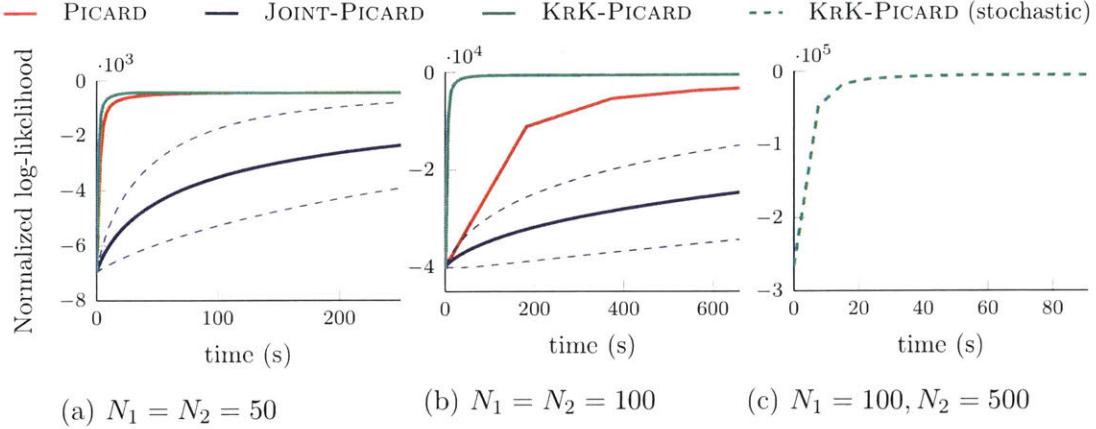


Figure 4-1:  $a = 1$ ; the dotted lines indicated the standard deviation from the mean.

### Small-scale real data: baby registries

We compared KRK-PICARD to PICARD and EM (Gillenwater et al., 2014) on the baby registry dataset, which has also been used to evaluate other DPP learning algorithms (Gillenwater et al., 2014; Mariet and Sra, 2015; Gartrell et al., 2017). The dataset contains 17 categories of baby-related products obtained from Amazon (see Appendix A.2.2 for a description of the dataset).

We learned kernels for the 6 largest categories ( $N = 100$ ); in this case, PICARD is sufficiently efficient to be preferred to KRK-PICARD; this comparison serves only to evaluate the quality of the final kernel estimates.

The initial marginal kernel  $\mathbf{K}$  for EM was sampled from a Wishart distribution with  $N$  degrees of freedom and an identity covariance matrix, then scaled by  $1/N$ ; for PICARD,  $\mathbf{L}$  was set to  $\mathbf{K}(\mathbf{I} - \mathbf{K})^{-1}$ ; for KRK-PICARD,  $\mathbf{L}_1$  and  $\mathbf{L}_2$  were chosen (as in JOINT-PICARD) by minimizing  $\|\mathbf{L} - \mathbf{L}_1 \otimes \mathbf{L}_2\|$ . Convergence was determined when the objective change dipped below a threshold  $\delta$ . As one EM iteration takes longer than one Picard iteration but increases the likelihood more, we set separate thresholds

for the different algorithms; specifically,  $\delta_{\text{PIC}} = \delta_{\text{KRK}} = 10^{-4}$  and  $\delta_{\text{EM}} = 10^{-5}$ .

The final log-likelihoods are shown in Table 4.2; we set the step-sizes to their largest possible values, *i.e.*,  $a_{\text{PIC}} = 1.3$  and  $a_{\text{KRK}} = 1.8$ . Table 4.2 shows that KRK-PICARD obtains comparable, albeit slightly worse log-likelihoods than PICARD and EM, which confirms that for tractable  $N$ , the better modeling capability of full kernels make them preferable to KRONDPPS.

Table 4.2: Final log-likelihoods for each large category of the baby registries dataset

(a) Training set				(b) Test set			
Category	EM	PICARD	KRK-PICARD	Category	EM	PICARD	KRK-PICARD
apparel	-10.1	-10.2	-10.7	apparel	-10.1	-10.2	-10.7
bath	-8.6	-8.8	-9.1	bath	-8.6	-8.8	-9.1
bedding	-8.7	-8.8	-9.3	bedding	-8.7	-8.8	-9.3
diaper	-10.5	-10.7	-11.1	diaper	-10.6	-10.7	-11.2
feeding	-12.1	-12.1	-12.5	feeding	-12.2	-12.2	-12.6
gear	-9.3	-9.3	-9.6	gear	-9.2	-9.2	-9.5

### Large-scale real dataset: GENES

Finally, to evaluate KRK-PICARD on large matrices of real-world data, we train it on data from the GENES (Batmanghelich et al., 2014) dataset (which has also been used to evaluate DPPs in (Li et al., 2015; Batmanghelich et al., 2014)). This dataset consists in 10,000 genes, each represented by 331 features corresponding to the distance of a gene to hubs in the BioGRID gene interaction network.

We construct a ground truth Gaussian DPP kernel on the GENES dataset and use it to obtain 100 training samples with sizes uniformly distributed between 50 and 200 items. Similarly to the synthetic experiments, we initialized KRK-PICARD’s kernel by setting  $\mathbf{L}_i = \mathbf{X}_i^\top \mathbf{X}_i$  where  $\mathbf{X}_i$  is a random matrix of size  $N_1 \times N_1$ ; for PICARD, we set the initial kernel  $\mathbf{L} = \mathbf{L}_1 \otimes \mathbf{L}_2$ .

Figure 4-2 shows the performance of both algorithms. As with the synthetic experiments, KRK-PICARD converges much faster; stochastic updates increase its performance even more, as shown in Fig. 4-2b. Average runtimes and speed-up are given in Table 4.3: KRK-PICARD runs almost an order of magnitude faster than

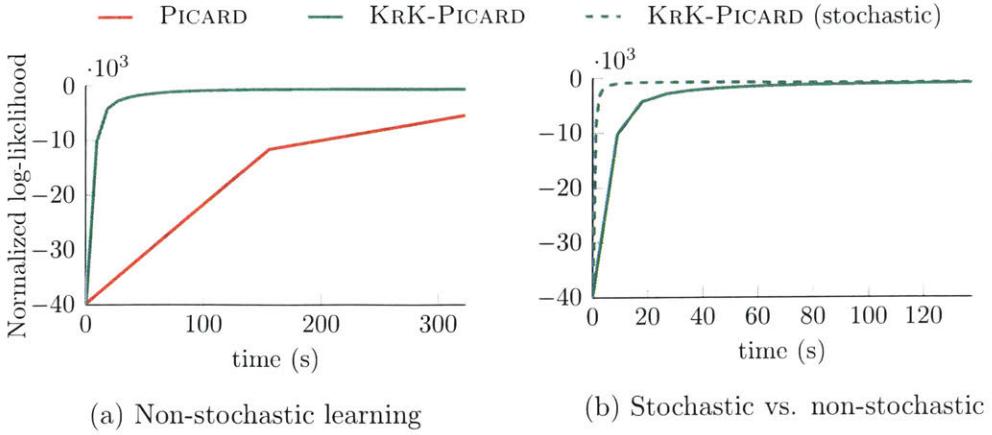


Figure 4-2:  $n = 150$ ,  $a = 1$ .

Table 4.3: Average runtime, performance on the GENES dataset for  $N_1 = N_2 = 100$

	PICARD	KRK-PICARD	KRK-PICARD (stoch.)
Average runtime	$161.5 \pm 17.7$ s	$8.9 \pm 0.2$ s	$1.2 \pm 0.02$ s
NLL increase (1st iter.)	$(2.81 \pm 0.03) \cdot 10^4$	$(2.96 \pm 0.02) \cdot 10^4$	$(3.13 \pm 0.04) \cdot 10^4$

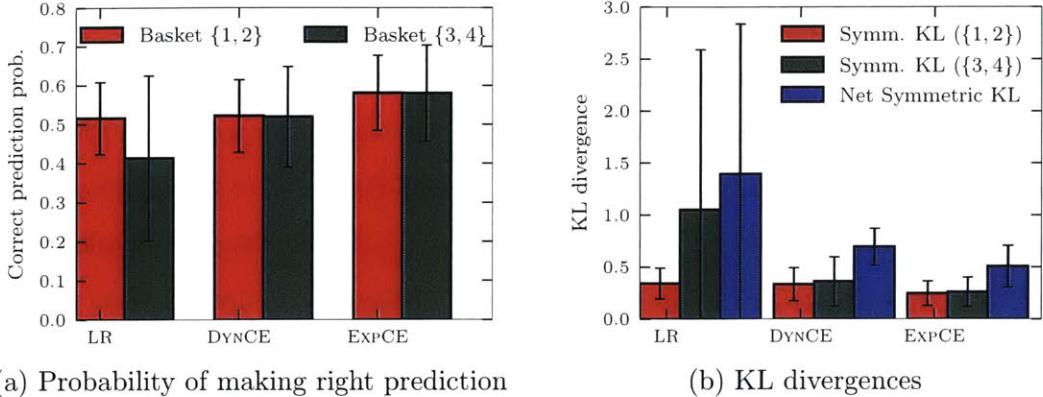
PICARD, and stochastic updates are more than two orders of magnitude faster, while providing slightly larger initial increases to the log-likelihood.

KRONDPPS also enable fast sampling, in  $\mathcal{O}(N^{3/2} + Nk^3)$  operations when using two sub-kernels, and in  $\mathcal{O}(Nk^3)$  when using three sub-kernels. This speedup allows for exact sampling at comparable or even better costs than previous algorithms for approximate sampling.

### 4.3 Learning with corrective negative information

We next consider the learning a (potentially unstructured) DPP kernel when information regarding “bad” subsets is available. Indeed, although fitting observed subsets well, maximum likelihood estimation for DPPs may also assign high likelihoods to unobserved subsets far from the underlying generative distribution (Chao et al., 2015). MLE-based DPP models may thus have modes corresponding to subsets that are close in likelihood, yet differ in how close they are to the true data distribution.

Such confusable modes reduce the quality of the learned model, hurting predic-



(a) Probability of making right prediction

(b) KL divergences

Figure 4-3: Results for experiments on a synthetic toy dataset. This toy dataset was generated by replicating the baskets  $\{1, 2\}$  and  $\{3, 4\}$  1000 times each. We randomly select 80% of this dataset for training, and 20% for test. We train each model to convergence, and then compute the next-item predictive probabilities for each unique pair, along with the symmetric KL divergence (only over areas of shared support) between the predictive and empirical next-item distributions. Net symmetric KL divergence is computed by adding the symmetric KL divergences for each of the two unique baskets. Experiments were run 10 times, with  $|\mathcal{A}^+|/|\mathcal{A}^-|$  set to the optimal value for each model;  $\alpha$  is set to its optimal LR value. See Section 4.3.1 for details on the DYN and EXP negative sampling models. We see that the LR (low-rank DPP) model can assign relatively high predictive probabilities to modes that represent incorrect predictions, resulting in higher symmetric KL divergence compared to the true empirical distribution and a high variance, revealing a high sensitivity to initialization. The DYN and EXP methods we introduce in Section 4.3.1 reduce this confusable mode issue, resulting in predictive distributions that are closer to the true distribution and much smaller variances.

tions: this is the scenario illustrated in Figure 4-3.

Such concerns when learning generative models over huge sample spaces are not limited to the area of subset-selection: applications in image and text generation have been the driving force in developing techniques for generating high-quality samples. Among their innovations, a particularly successful technique uses generated samples as “negative samples” to train a discriminator, which in turn encourages generation of more realistic samples; this is the key idea behind the Generative Adversarial Nets (GANs) introduced in (Goodfellow et al., 2014).

These observations motivate us to investigate the use of DPP-generated samples with added perturbations as *negatives*, which we incorporate into the learning task to improve the modeling power of DPPs. Intuitively, negative samples are subsets that are far from the true data distribution, but to which the learned DPP model erro-

neously assigns high probability. As there is no closed form solution to generating such idealized negatives, we approximate them via an external “negative distribution”.

Specifically, we introduce a DPP learning problem that incorporates samples from a negative distribution into traditional MLE. While the focus of this section lies on generating the negative distribution *jointly* with  $\mathbf{L}$ , we also investigate outside sources of negative information. Ultimately, our formulation leads to an optimization problem harder than the original DPP learning problem. However, we show that even approximate solutions greatly improve the performance of the DPP model when evaluated on concrete tasks: identifying the best item to add to a subset of chosen objects (*basket-completion*) and discriminating between held-out test data and randomly generated subsets. The following material is based on (Mariet et al., 2019a).

## Related work

Aside from Tschiatschek et al. (2016); Djolonga et al. (2016), who learn a Facility Location Diversity (FLID) distribution (as well as more complex FLIC and FLDC models) by contrasting it with a “negative” product distribution, little attention has been given to using negative samples to learn richer subset-selection models.

Nonetheless, leveraging negative information is a widely used in other applications. In object detection, negative mining corrects for the skewed simple-to-difficult negative distribution by training the model on its false positives (Sung, 1996; Canvet and Fleuret, 2014; Shrivastava et al., 2016). In language modeling, Noise Contrastive Estimation (NCE) (Gutmann and Hyvärinen, 2012), which tasks the model with distinguishing positive samples from generated negatives, was first applied in (Mnih and Teh, 2012) and has been instrumental in Word2Vec (Mikolov et al., 2013). Since then, variants using adaptive noise (Chen et al., 2017) have been introduced. NCE is also the method used by Tschiatschek et al. (2016) for subset-selection.

An alternate approach to negative samples within submodular language models was introduced as Contrastive Estimation in (Smith and Eisner, 2005a,b). Negative sampling is also used in GANs (Goodfellow et al., 2014), where a generator network competes with a discriminative network that distinguishes between positives and gen-

erated negatives. An adversarial approach to Contrastive Estimation has been recently introduced in (Bose et al., 2018), where ideas from GANs for discrete data are used to implement an adversarial negative sampler that augments a conventional negative sampler.

Motivated by the similarities between DPP learning and crucial structured prediction problems in other ML fields, we introduce an optimization problem that leverages negative information. We refer to this problem as Contrastive Estimation (CE) due to its ties to a notion discussed in (Smith and Eisner, 2005a).

### 4.3.1 Contrastive Estimation

Recall that in conventional DPP learning, we seek to maximize determinantal volumes of sets drawn *i.i.d.* from a source distribution  $\mu$  (that we wish to model) by solving

$$\text{Find } \mathbf{L} \in \underset{\mathbf{L} \succeq 0}{\operatorname{argmax}} \phi_{\text{MLE}}(\mathbf{L}) \triangleq \frac{1}{|\mathcal{A}^+|} \sum_{A \in \mathcal{A}^+} \log \det(\mathbf{L}_A) - \log \det(\mathbf{L} + \mathbf{I}). \quad (4.11)$$

To leverage information about *unlikely* samples, we augment problem (4.11) to incorporate additional information from a *negative* distribution  $\nu$ , which we wish to have the DPP distribution move away from.

**Definition 4.3.1** (Contrastive Estimation). Given a training set of positive samples  $\mathcal{A}^+$  on which  $\phi_{\text{MLE}}$  is defined and a negative distribution  $\nu$  over  $2^{\mathcal{Y}}$ , we call *Contrastive Estimation* the problem

$$\text{Find } \mathbf{L} \in \underset{\mathbf{L} \succeq 0}{\operatorname{argmax}} \phi_{\text{CE}}(\mathbf{L}) \triangleq \phi_{\text{MLE}}(\mathbf{L}) - \mathbb{E}_{A \sim \nu} [\log \mathcal{P}_{\mathbf{L}}(A)], \quad (4.12)$$

where we write  $\mathcal{P}_{\mathbf{L}}(A) \equiv \det(\mathbf{L}_A) / \det(\mathbf{L} + \mathbf{I})$ .

The expectation in Eq. 4.12 is usually intractable. However, it can be approxi-

mated by sampling a set of samples  $\mathcal{A}^-$  from the distribution  $\nu$ :  $\phi_{\text{CE}}$  then becomes<sup>4</sup>

$$\phi_{\text{CE}}(\mathbf{L}) = \frac{1}{|\mathcal{A}^+|} \sum_{A \in \mathcal{A}^+} \log \mathcal{P}_{\mathbf{L}}(A) - \frac{1}{|\mathcal{A}^-|} \sum_{A \in \mathcal{A}^-} \log \mathcal{P}_{\mathbf{L}}(A). \quad (4.13)$$

If  $|\mathcal{A}^-| = 0$ , the CE objective (4.12) reduces to  $\phi_{\text{MLE}}$ . Conversely,  $\phi_{\text{MLE}}$  can be viewed as a sample-based approximation of the value  $\mathbb{E}_{A \sim \mu}[\log \mathcal{P}_{\mathbf{L}}(A)]$ , where  $\mu$  is the true distribution generating the samples in  $\mathcal{A}^+$ . Interestingly, another reformulation of (4.12) suggests an even broader class of DPP kernel learning: indeed, let  $y_A$  be  $\frac{1}{|\mathcal{A}^+|}$  (resp.  $-\frac{1}{|\mathcal{A}^-|}$ ) for  $A \in \mathcal{A}^+$  (resp.  $\mathcal{A}^-$ ), and define

$$\mathcal{A} = \{(y_A, A) : A \in \mathcal{A}^+\} \cup \{(y_A, A) : A \in \mathcal{A}^-\},$$

where the  $y_A$  should be viewed as belonging in  $\{-1, 1\}$  with an additional normalization coefficient. Then, we can rewrite equation (4.13) in the following form

$$\phi_{\text{CE}}(\mathbf{L}) = \sum_{(y_A, A) \in \mathcal{A}} y_A \left[ \log \det \mathbf{L}_A - \det(\mathbf{L} + \mathbf{I}) \right]. \quad (4.14)$$

Formulation (4.14) suggests the use of a broader scope of continuous labels  $y_A$ ; we do not cover this variation in the present work, but note that (4.14) permits the use of *weighted* samples for learning.

**Remark 4.3.1.** Compared to the traditional Noise Contrastive Estimation (NCE) approach, which requires full knowledge of the negative distribution, CE does not suffer any such limitation: we only require an estimate of  $\mathbb{E}_{\nu}[\log \mathcal{P}_{\mathbf{L}}(A)]$ .

**Remark 4.3.2.** Eq. (4.12) can be made to go to  $+\infty$  with pathological negative samples (*i.e.*,  $\mathcal{P}_{\mathbf{L}}(A^-) = 0$ ); hence, choosing the negative distribution is a crucial concern for CE. In practice, we do not observe this pathological behavior (cf. Section 4.3.3).

**Remark 4.3.3.** CE is a non-convex optimization problem, and thus admits the same guarantees as DPP MLE learning when learned using Stochastic Gradient Ascent with decreasing step sizes; however, the convergence rate will depend on the choice of  $\nu$ .

---

<sup>4</sup>With a slight abuse of notation, we write  $\phi_{\text{CE}}$  despite the sample approximation to  $\mathbb{E}_{A \sim \nu}[\cdot]$ .

Indeed, to fully specify the CE problem one must first choose the negative distribution  $\nu$ , or equivalently, choose a procedure to generate negative samples to obtain (4.13). We consider two classes of distributions  $\nu$  with considerably different ramifications: dynamic and static negatives; their analysis is the focus of the next sections.

### Dynamic negatives

In most applications leveraging negative information (e.g., negative mining, GANs), negative samples evolve over time based on the state of the learned model. We call any  $\nu$  that depends on the state of the model a *dynamic negative distribution*: at iteration  $k$  of the learning procedure with kernel estimate  $\mathbf{L}_k$ , we use a  $\nu$  parametrized by  $\mathbf{L}_k$ .

More specifically, we focus on the setting where negative samples themselves are generated by the current DPP, with the goal of reducing overfitting. Given a positive sample  $A^+$ , we generate a negative  $A^-$  by replacing  $i \in A^+$  with  $j$  that yields a high probability  $\mathcal{P}_{\mathbf{L}_k}(A^+ \setminus \{i\} \cup \{j\})$  (Alg. 9). We generate the samples probabilistically rather than via mode maximization so that a sample  $A^+$  can lead to different  $A^-$  negatives when we generate more negatives than positives.

---

#### Algorithm 9 Generate dynamic negative

---

```

Input: Positive sample  $A^+$ , current kernel  $\mathbf{L}_k$ 
Sample  $i \in A^+$  prop. to its empirical probability in  $A^+$ 
 $A^- := A^+ \setminus \{i\}$ 
Sample  $j$  w.p. proportional to  $\mathcal{P}_{\mathbf{L}_k}(A^- \cup \{j\})$ 
 $A^- \leftarrow A^- \cup \{j\}$ 
return  $A^-$ 
```

---

As  $\nu$  evolves along with  $\mathbf{L}_k$ , the second term of  $\phi_{\text{CE}}$  acts as a moving target that must be continuously estimated during the learning procedure. For this reason, we choose to optimize  $\phi_{\text{CE}}$  by a two-step procedure described in Alg. 10, similarly to an alternating maximization approach such as EM.

This approach bears strong similarities with GANs, in which both generator and discriminator evolve during training (dynamic negatives appear in Goodfellow (2014) as a theoretical tool to analyze differences between NCE and GANs).

Once the generated negative  $A^-$  has been used in an iteration of the optimization

---

**Algorithm 10** Optimizing dynamic CE

---

**Input:** Positive samples  $\mathcal{A}^+$ , initial kernel  $\mathbf{L}_0$ , maxIter.  
 $k \leftarrow 1$   
**while**  $k++ < \text{maxIter}$  and not converged **do**  
     $\mathcal{A}^- \leftarrow \text{GENERATEDYNAMICNEGATIVES}(\mathbf{L}_k, \mathcal{A}^+)$   
     $\mathbf{L}_{k+1} \leftarrow \text{OPTIMIZECE}(\mathbf{L}_k, \mathcal{A}^+, \mathcal{A}^-)$   
**return**  $\mathbf{L}_k$

---

of  $\phi_{\text{CE}}$ , it is less likely to be sampled again.<sup>5</sup> Crucially, such dynamic negatives also avoid the problem alluded to in Remark 4.3.2, since by construction they have a non-zero probability under  $\mathcal{P}_{\mathbf{L}_k}$  at iteration  $k$ .

### Static negatives

Conversely, we can simplify the optimization problem by considering a *static* negative distribution:  $\nu$  does not depend on the current kernel estimate. A considerable theoretical advantage of static negatives lies in the simpler optimization problem: given a static negative distribution  $\nu$ , the objective  $\phi_{\text{CE}}$  does not evolve during training, and is amenable to a simple invocation of stochastic gradient descent (Bottou, 1998).

However, static negative distributions may suffer from the fundamental theoretical issue in Rem. 4.3.2, and hence careful attention must be paid to ensure that the learning algorithm does not converge to a spurious optimum that assigns a probability  $\Pr_{\mathbf{L}}(A) = 0$  to  $A \in \mathcal{A}^-$ . In practice, we observed that the local nature of stochastic gradient ascent iterations was sufficient to avoid such behavior.

Let us now discuss two classical choices for fixed  $\nu$ .

**Product negatives.** A common choice of negative distribution in other machine learning areas is the *product distribution*, which is the standard “noise” distribution used in NCE. It is defined by

$$\nu(A) = \prod_{i \in A} \hat{p}(i) \prod_{i \notin A} (1 - \hat{p}(i)), \quad (4.15)$$

---

<sup>5</sup>If  $A^-$  happens to be a false negative (*i.e.*, appears in  $\mathcal{A}^+$ ),  $A^-$  will be comparatively sampled more frequently as a positive, and so will contribute on average as a positive sample. Additional precautions such as the ones mentioned in (Bose et al., 2018) can also be leveraged if necessary.

where  $\hat{p}(i)$  is the empirical probability of  $\{i\}$  in  $\mathcal{A}^+$ . Although (Mikolov et al., 2013) reports better results by raising the  $\hat{p}$  to the power  $\frac{3}{4}$ , we did not observe any improvements when using exponentiated power distributions; for this reason, by *product negatives*, we always indicate the baseline distribution (4.15).

The product distribution is in practice a mismatch for DPPs, as it lacks the negative association property of DPPs that enables them to model the repulsive interactions between similar items.

**Explicit negatives.** Alternatively, we may have prior knowledge of a class of subsets that our model should *not* generate. For example, we might know that items  $i$  and  $j$  are negatively correlated and hence unlikely to co-occur. We may also learn via user feedback that some generated subsets are inaccurate. We refer to negatives obtained using such outside information as *explicit negatives*.

A fundamental advantage of explicit negatives is that they can incorporate prior knowledge and user feedback into the learning algorithm. The ability to incorporate such information, to our knowledge, is in itself a novel contribution to DPP learning.

Although such knowledge may be costly and/or only available at rare intervals, a form of continuous learning that would regularly update the state of our prior knowledge (and hence  $\nu$ ) would bring the explicit negative distribution into the realm of dynamic distributions, as described by Alg. 10.

### 4.3.2 Efficient learning and prediction

We now describe how the Contrastive Estimation problem for DPPs can be optimized efficiently. In order to efficiently generate dynamic negatives, which rely on DPP conditioning, we additionally generalize the dual transformation leveraged in (Osogami et al., 2018) to speed up basket-completion tasks with DPPs. This speed-up impacts the broader use of DPPs, outside of CE learning.

## Optimizing $\phi_{\text{CE}}$

We propose to optimize the CE problem by exploiting a low-rank factorization of the kernel, writing  $\mathbf{L} = \mathbf{V}\mathbf{V}^\top$ , where  $\mathbf{V} \in \mathbb{R}^{N \times K}$  and  $K \leq N$  is the rank of the kernel, which is fixed *a priori*.

This factorization ensures that the estimated kernel remains positive semi-definite, and enables us to leverage the low-rank computations derived in (Gartrell et al., 2017) and refined in (Osogami et al., 2018). Given the similar forms of the MLE and CE objectives, we use the stochastic gradient ascent algorithm introduced by (Gartrell et al., 2017) to optimize (4.12). In the case of dynamic negatives, we re-generate  $\mathcal{A}^-$  after each gradient step; less frequent updates are also possible if the negative generation algorithm is very costly.

We furthermore augment  $\phi_{\text{CE}}$  with a regularization term  $R(\mathbf{V})$ , defined as

$$R(\mathbf{V}) = \alpha \sum_{i=1}^N \frac{1}{\mu_i} \|v_i\|_2^2,$$

where  $\mu_i$  counts the occurrences of  $i$  in the training set,  $v_i$  is the corresponding row vector of  $\mathbf{V}$  and  $\alpha > 0$  is a tunable hyperparameter. Note that this is the same regularization as introduced in (Gartrell et al., 2017). This regularization tempers the strength of  $\|v_i\|_2$ , a term interpretable as to the popularity of item  $i$  (Kulesza and Taskar, 2012; Gillenwater et al., 2014), based on its *empirical* popularity  $\mu_i$ . Experimentally, we observed that adding  $R(\mathbf{V})$  has a strong impact on the predictive quality of our model.

The reader may wonder if other approaches to DPP learning are also applicable to the CE problem.

**Remark 4.3.4.** The update rule for the fixed-point approach in (Mariet and Sra, 2015) does not admit a closed form solution for CE, rendering it impractical.

Indeed, letting  $\beta = |\mathcal{A}^+| - |\mathcal{A}^-| \geq 0$  and writing  $\mathbf{U}_A$  as the  $N \times |A|$  indicator

matrix such that  $\mathbf{L}_A = \mathbf{U}_A^\top \mathbf{L} \mathbf{U}_A$ , we have

$$\begin{aligned}\phi(\mathbf{L}) &\propto -\beta \log \det(\mathbf{I} + \mathbf{X}) + \underbrace{\sum_{A \in \mathcal{A}^+} \log \det(\mathbf{U}_A^\top \mathbf{X}^{-1} \mathbf{U}_A)}_{f \text{ convex}} \\ &\quad + \underbrace{\beta \log \det(\mathbf{X}) - \sum_{A \in \mathcal{A}^-} \log \det(\mathbf{U}_A^\top \mathbf{X}^{-1} \mathbf{U}_A)}_{g \text{ concave}},\end{aligned}$$

where the convexity/concavity results follow immediately from (Mariet and Sra, 2015, Lemma 2.3). Then, the update rule  $\nabla f(\mathbf{L}_{k+1}) = -\nabla g(\mathbf{L}_k)$  requires

$$\begin{aligned}&\beta \mathbf{L}_{k+1} + \sum_{A \in \mathcal{A}^-} \mathbf{L}_{k+1} \mathbf{U}_A (\mathbf{U}_A^\top \mathbf{L}_{k+1} \mathbf{U}_A)^{-1} \mathbf{U}_A^\top \mathbf{L}_{k+1} \\ &\leftarrow \beta(\mathbf{I} + \mathbf{L}_k^{-1})^{-1} + \sum_{A \in \mathcal{A}^+} \mathbf{L}_k \mathbf{U}_A (\mathbf{U}_A^\top \mathbf{L}_k \mathbf{U}_A)^{-1} \mathbf{U}_A^\top \mathbf{L}_k,\end{aligned}$$

which cannot be evaluated due to the  $\sum_{A \in \mathcal{A}^-}$  term.

The low-rank formulation allows us to apply CE (as well as NCE, as discussed in Section 4.3.3) to learn large datasets such as the Belgian retail supermarket dataset (described in Section 4.3.3) without prohibitive learning runtimes. We show below that by leveraging the idea described in (Osogami et al., 2018), the low-rank formulation can also lead to additional speed ups during prediction.

## Efficient conditioning for predictions

Dynamic negatives rely upon conditioning a DPP on a chosen sample  $A$  (see Alg. 9:  $\mathcal{P}_{\mathbf{L}_k}(A^- \cup \{j\})$ ) can be efficiently computed for all  $j$  by a preprocessing step that conditions  $\mathbf{L}_k$  on set  $A^-$ ). For this reason, we now describe how low-rank DPP conditioning can be significantly sped up.<sup>6</sup>

In (Gartrell et al., 2017), conditioning has a cost of  $\mathcal{O}(K|\bar{A}|^2 + |A|^3)$ , where  $\bar{A} = \mathcal{Y} - A$ . Since  $|\mathcal{Y}| \gg |A|$  for many datasets, this represents a significant bottleneck for conditioning and computing next-item predictions for a set. We show here that this

---

<sup>6</sup>We include this section for completeness; these results are attributed to Mike Gartrell.

complexity can be brought down significantly.

**Proposition 4.3.1.** *Given  $A \subseteq \{1, \dots, N\}$  and a DPP of rank  $K$  parametrized by  $\mathbf{V}$ , where  $\mathbf{L} = \mathbf{V}\mathbf{V}^\top$ , we can derive the conditional marginal probabilities in the DPP parametrization  $\mathbf{L}^A$  in only  $\mathcal{O}(K^3 + |A|^3 + K^2|A|^2 + |\bar{A}|K^2)$  time.*

*Proof.* Let  $\mathbf{V}$  be the low-rank parametrization of the DPP kernel ( $\mathbf{L} = \mathbf{V}^\top \mathbf{V}$ ) and  $A \subseteq \mathcal{Y}$ . As in Gillenwater et al. (2014), we first compute the dual kernel  $\mathbf{C} = \mathbf{B}^\top \mathbf{B}$ , where  $\mathbf{B} = \mathbf{V}^\top$ . We then compute

$$\mathbf{C}^A = (\mathbf{B}^A)^\top \mathbf{B}^A = \mathbf{Z}^A \mathbf{C} \mathbf{Z}^A,$$

with  $\mathbf{Z}^A = \mathbf{I} - \mathbf{B}_A(\mathbf{B}_A^\top \mathbf{B}_A)^{-1} \mathbf{B}_A^\top$ , and where  $\mathbf{C}^A$  is the DPP kernel conditioned on the event that all items in  $A$  are observed, and  $\mathbf{B}_A$  is the restriction of  $\mathbf{B}$  to the rows and columns indexed by  $A$ .

Computing  $\mathbf{C}^A$  costs  $\mathcal{O}(K^3 + |A|^3 + K^2|A|^2)$ . Next, following Kulesza and Taskar (2012), we eigendecompose  $\mathbf{C}^A$  to compute the conditional (marginal) probability  $P_i$  of every possible item  $i$  in  $\bar{A}$ :

$$P_i = \sum_{n=1}^K \frac{\lambda_n}{\lambda_n + 1} \left( \frac{1}{\sqrt{\lambda_n}} b_i^A \hat{v}_n \right)^2,$$

where  $b_i^A$  indicates the column vector for item  $i$  in  $\mathbf{B}^A$  and  $(\lambda_n, \hat{v}_n)$  are the eigenvalues/vectors of  $\mathbf{C}^A$ .

The computational complexity for computing the eigendecomposition is  $\mathcal{O}(K^3)$ , and computing  $P_i$  for all items in  $\bar{A}$  costs  $\mathcal{O}(|\bar{A}|K^2)$ . Therefore, we have an overall computational complexity of  $\mathcal{O}(K^3 + |A|^3 + K^2|A|^2 + |\bar{A}|K^2)$  for computing next-item conditionals/predictions for the low-rank DPP using the dual kernel, which is significantly superior to the typical cost of  $\mathcal{O}(K|\bar{A}|^2 + |A|^3)$ .  $\square$

As in most cases  $K \ll |\bar{A}|$ , this represents a substantial improvement, allowing us condition in time essentially linear in the size of the item catalog.

### 4.3.3 Experiments

We run next-item prediction and AUC-based classification experiments on two recommendation datasets for DPP evaluation: the UK retail dataset (Chen, 2012), which after clipping all subsets to a maximum size<sup>7</sup> of 100, contains 4070 items and 20059 subsets, and the Belgian Retail Supermarket dataset,<sup>8</sup> which contains 88,163 subsets, of a total of 16,470 unique items (Brijs et al., 1999; Brijs, 2003). We compare the following Contrastive Estimation approaches:

- EXP: explicit negatives learned with CE. As to our knowledge there are no datasets with explicit negative information, we generate approximations of explicit negatives by removing one item from a positive sample and replacing it with the least likely item (Algorithm 11).
- DYN: dynamic negatives learned with CE.

As our work revolves around improving DPP performance, we focus on the two following baselines, which are targeted to learning DPP parametrizations from data:

- LR: the standard low-rank DPP stochastic gradient ascent algorithm (Gartrell et al., 2017).
- NCE: Noise Contrastive Estimation using product negatives. NCE learns a model by contrasting  $\mathcal{A}^+$  with negatives drawn from a “noisy” distribution  $p_n$ , training the model to distinguish between sets drawn from  $\mu$  and sets drawn from  $p_n$ .

NCE has gained popularity due to its ability to model distributions  $\mu$  with intractable normalization coefficients, and has been shown to be a powerful technique to improve submodular recommendation models (Tschiatschek et al., 2016). NCE learns by maximizing the following conditional log-likelihood:

$$\phi_{\text{NCE}}(\mathbf{L}) = \sum_{A \in \mathcal{A}^+} \log P(A \in \mathcal{A}^+ | A) + \sum_{A \in \mathcal{A}^-} \log P(A \in \mathcal{A}^- | A). \quad (4.16)$$

---

<sup>7</sup>This allows us to use a low-rank matrix factorization for the DPP that scales well in terms of train and prediction time.

<sup>8</sup><http://fimi.ua.ac.be/data/retail.pdf>

In our experiments, we learn the NCE objective with stochastic gradient ascent for our low-rank model, since  $\nabla \log \Pr(A \in \mathcal{A}^* | A, \mathbf{V}\mathbf{V}^\top)$  is given by

$$\left( \epsilon^* - \left( 1 + \frac{|\mathcal{A}^-|}{|\mathcal{A}^+|} \frac{p_n(A)}{\mathcal{P}_{\mathbf{V}\mathbf{V}^\top}(A)} \right)^{-1} \right) \nabla_{\mathbf{V}} \log \mathcal{P}_{\mathbf{V}\mathbf{V}^\top}(A). \quad (4.17)$$

where  $\epsilon^* = 1$  if  $\mathcal{A}^* = \mathcal{A}^+$  and 0 otherwise.

We approximate explicit negatives using Algorithm 11. This allows us to approximate true explicit negatives, as we use the empirical data to derive “implausible” sets. However, when using such negatives we have no guarantee that objective function will be well behaved, as opposed to the theoretically grounded dynamic negatives.

---

**Algorithm 11** Approximate explicit negative generation

---

**input:** Positive sample  $A^+$   
 Sample  $i \neq j \in A^+$  w.p.  $p_i \propto \widehat{P}(\{i\})$   
 Sample  $k \notin A^+$  w.p.  $p_k \propto 1 - \widehat{P}(\{i, k\})$ .  
**return**  $(A^+ \setminus \{j\}) \cup \{k\}$

---

The performance of all methods are compared using standard recommender system metrics: Mean Percentile Rank (MPR). MPR is a recall-based metric that evaluates predictive power by measuring how well the model predicts the next item in a basket, and is a standard choice for recommender systems (Hu et al., 2008; Li et al., 2010).

Specifically, given a set  $A$ , let  $p_{i,A} = \Pr(A \cup \{i\} | A)$ . The percentile rank of an item  $i$  given a set  $A$  is defined as

$$\text{PR}_{j,A} = \frac{\sum_{i' \notin A} \mathbb{1}(p_{i,A} \geq p_{i',A})}{|\mathcal{Y} \setminus A|} \times 100\%.$$

The MPR is then computed as

$$\text{MPR} = \frac{1}{|\mathcal{T}|} \sum_{A \in \mathcal{T}} \text{PR}_{i,A \setminus \{i\}},$$

where  $\mathcal{T}$  is the set of test instances and  $i$  is a randomly selected element in each set  $A$ . An MPR of 50 is equivalent to random selection; a MPR of 100 indicates that the model perfectly predicts the held out item.

We also evaluate the discriminative power of each model using the AUC metric. For this task, we generate a set of negative subsets uniformly at random. For each positive subset  $A^+$  in the test set, we generate a negative subset  $A^-$  of the same length by drawing  $|A^+|$  samples uniformly at random, while ensuring that the same item is not drawn more than once for a subset.

We then compute the AUC for the model on these positive and negative subsets, where the score for each subset is the log-likelihood that the model assigns to the subset. This task measures the ability of the model to discriminate between positive subsets (ground-truth subsets) and randomly generated subsets.

In all experiments, 80% of subsets are used for training; the remaining 20% served as test; convergence is reached when the relative change in the validation log-likelihood is below a pre-determined threshold  $\epsilon$ , set identically for all methods. All results are averaged over 5 learning trials.

### Amazon registries

We conducted an experimental analysis on the largest 7 sub-datasets included in the Amazon Registry dataset. Given the small size of these datasets (the largest has 100 items), these experiments serve only to provide insight into the general behavior of the baselines and CE methods, and into the influence of hyperparameters on convergence.

Table 4.4: MPR, p@ $k$ , and AUC values for LR, and baseline improvement over LR for other methods. Positive values indicate the algorithm performs better than LR, and bold values indicate improvement over LR that lies outside the standard deviation. Experiments were run 5 times, with  $|\mathcal{A}^+|/|\mathcal{A}^-| = \frac{1}{2}$ ;  $\alpha$  is set to its optimal LR value.

Metric	LR	Improvement over LR		
		DYN	EXP	NCE
MPR	70.50	<b>0.92 ± 0.56</b>	<b>0.68 ± 0.62</b>	<b>0.86 ± 0.55</b>
p@1	9.96	0.67 ± 0.75	0.58 ± 0.76	0.20 ± 1.75
p@5	25.36	<b>1.04 ± 0.82</b>	<b>0.78 ± 0.67</b>	0.67 ± 1.09
p@10	36.50	<b>1.39 ± 0.85</b>	<b>1.13 ± 0.79</b>	0.97 ± 1.18
p@20	51.22	<b>1.38 ± 0.97</b>	<b>1.28 ± 1.11</b>	<b>1.35 ± 1.20</b>
AUC	0.630	<b>0.027 ± 0.017</b>	<b>0.026 ± 0.016</b>	0.009 ± 0.017

In Table 4.4, we compare the performance of the various algorithms with rank  $K = 30$ . The regularization strength  $\alpha$  is set to its optimal value for the LR algorithm, and  $|\mathcal{A}^-|/|\mathcal{A}^+| = 1/2$ . This allows us to compare the LR algorithm to its “augmented” negative versions without hyper-parameter tuning. As PROD performs much worse than LR, it is not included in further experiments.

We evaluate the precision at  $k$  as

$$p@k = \frac{1}{|\mathcal{T}|} \sum_{A \in \mathcal{T}} \frac{1}{|A|} \sum_{i \in A} \mathbf{1}[\text{rank}(i | A \setminus \{i\}) \leq k].$$

Compared to traditional SGA methods, algorithms that use inferred negatives perform (PROD excepted) better across all metrics and datasets. DYN and EXP provide consistent improvements compared to the other methods, whereas NCE shows a higher variance and slightly worse performance. Improvements observed using DYN and EXP are larger than the loss in performance due to going from full-rank to low-rank kernels reported in (Gartrell et al., 2017).

We also found that CE was not sensitive to the  $\alpha$  and  $|\mathcal{A}^-|/|\mathcal{A}^+|$  hyperparameters. For this reason, in all further results, we set  $|\mathcal{A}^-|/|\mathcal{A}^+| = .5$  and  $\alpha = 1$ . As in previous work on low-rank DPP learning (Gartrell et al., 2017),  $\alpha = 1$  was found to be a reasonably optimal value, this ensures a fair comparison between all methods.

Table 4.5 reports the average time to convergence for each method. As generating the dynamic negatives has a high complexity due to DPP conditioning, DYN is 2.7x slower than EXP. LR is the fastest method, as it does not need to process any negatives. NCE is by far the most time-consuming.

Table 4.5: Runtime to convergence (s) on the feeding Amazon registry ( $\alpha = 1$ ,  $|\mathcal{A}^-|/|\mathcal{A}^+| = 0.5$ ,  $K = 30$ ).

METHOD	LR	EXP	DYN	NCE
RUNTIME (s)	$0.83 \pm 0.54$	$2.69 \pm 0.02$	$7.13 \pm 0.28$	$27.59 \pm 2.20$

(a) UK dataset			(b) Belgian dataset			
Metric	Improvement over LR		Improvement over LR			
	LR	EXP	DYN	LR	EXP	
MPR	80.07	<b>3.75 ± 0.16</b>	<b>3.74 ± 0.16</b>	79.42	<b>9.58 ± 0.15</b>	<b>9.64 ± 0.13</b>
AUC	0.57297	<b>0.41465 ± 0.01334</b>	<b>0.41467 ± 0.01339</b>	0.6162	<b>0.3705 ± 8.569e-5</b>	<b>0.3702 ± 4.447e-5</b>

Table 4.6: Results over the UK and Belgian datasets. Both explicit and dynamic CE obtain statistically significant improvements in MPR and AUC metrics, confirming that CE learning enhances recommender value of the model and its ability to distinguish data drawn from the target distribution from fake samples. The impact on precisions@ $k$  metrics is not reported as we did not observe statistically significant deviations from LowRank performance.

## UK and Belgian Retail Datasets

Following (Gartrell et al., 2017), for both the UK and the Belgian dataset, we set the rank  $K$  of the kernel to be the size of the largest subset in the dataset ( $K=100$  for the UK dataset,  $K=76$  for the Belgian dataset): this optimizes memory costs while still modeling all ground-truth subsets. Based on our results on the smaller Amazon dataset, we set  $|\mathcal{A}^-|/|\mathcal{A}^+| = 0.5$  and  $\alpha = 1$ .

Corroborating our timing results on the Amazon registry, we saw that one iteration of NCE required nearly 11 hours on the Belgian dataset (compared to 5 minutes for one iteration of CE). For this reason, we remove NCE as a baseline from all remaining experiments, as it is not feasible in the general case.

Tables 4.6 (a) and (b) summarize our results; the negative methods show significant MPR improvement over LR, with both DYN and EXP performing almost 10 points higher on the Belgian dataset, and 3 points higher on the UK dataset. This is a striking improvement, compounded by small standard deviations confirming that these results are robust to matrix initialization.

We also see a dramatic improvement over LR in AUC, with an improvement of approximately 0.41 for the UK dataset and 0.37 for the Belgian dataset, across both DYN and EXP methods. Both DYN and EXP perform quite well, with an AUC score of approximately 0.9864 or higher for both models. These results suggest that for larger datasets, CE is effective at improving the discriminative power of the DPP.

## 4.4 Open problems

- Kulesza (2012) conjectured that the general form of the maximum likelihood estimation problem is NP-hard; this question remains open.
- In (Mariet and Sra, 2015), the question of which range of values for the step-size parameter  $a$  will lead to the same guarantees of MLE increase and positive semi-definite estimates in the Picard iteration was raised, and remains open for the KRK-PICARD extension.
- As is the case with any structured learning scheme, a Kronecker DPP allows for faster learning and sampling with lower memory costs with the trade-off of losing in expressive power — although experiments on small datasets suggest that this loss of expressivity is minimal, an in-depth investigation of this trade-off and how it compares to other structured kernels would be valuable in choosing one structure over another.
- Algorithm KRK-PICARD is motivated by the necessity of preserving the Kronecker structure, while simultaneously maintaining the necessary convex and concave properties of the log-likelihood function  $\phi(\mathbf{L}) = f(\mathbf{L}) + g(\mathbf{L})$ . Although  $g$  is not convex in  $\mathbf{L}_1 \otimes \mathbf{L}_2$ ,  $g$  is *geodesically* convex over the Riemannian manifold of positive definite matrices, which suggests that deriving an iteration that would take advantage of the intrinsic geometry of the problem may be a viable line of future work.
- Contrastive Estimation with dynamic negatives is a difficult optimization problem: the objective function is a moving target, evolving along with the current estimate. As such, the CE optimization problem is in of itself a theoretical problem worthy of independent study.
- The generating logic for both dynamic and static negatives may benefit from further investigation, for CE and other DPP learning schemes.

- In various works (Gillenwater et al., 2014; Mariet and Sra, 2015), the initialization of a DPP learning algorithm has been shown to significantly impact the final log-likelihood; a deeper study of initialization procedures for DPP learning algorithms, including KRK-PICARD, remains an important investigation area.



# Chapter 5

## Generalized volume sampling

We have seen that determinantal point processes are a powerful and intuitive model for diversity within subsets of a ground set. However, DPPs are far from the only measures able to tractably leverage negative dependence, nor indeed the only tractable negatively dependent measures with key implications within machine learning.

Suppose a chemist is planning on running expensive, time-consuming experiments to evaluate an underlying parameter  $\theta \in \mathbb{R}^m$ . Given experimental parameters represented by a vector  $x \in \mathbb{R}^m$ , they observe the outcome  $y = x^\top \theta + \epsilon \in \mathbb{R}$ , where  $\epsilon$  is some random noise. As the experiments are costly to run, the chemist must, out of a total number  $N$  of possible experiments, choose a much smaller set of  $k$  experiments to run in order to reliably approximate  $\theta$ . Equivalently, the chemist must select a small set of rows from the tall and thin experimental design matrix  $\mathbf{X} = [x_1, \dots, x_N]^\top \in \mathbb{R}^{N \times m}$ . Negative dependence is of fundamental importance in this scenario; as the chemist can only choose a few experiments, they must choose carefully so that the outcome of two different experiments provides them with complementary information.

This chapter focuses on *Generalized Volume Sampling (GVS) measures*: strongly Rayleigh measures that select  $k$  rows from a tall and thin matrix  $X \in \mathbb{R}^{N \times m}$  ( $m \leq k \leq N$ ), and which exactly characterize the negative dependence exhibited by the experimental design problem. A special case of Gvs measure was introduced by Avron and Boutsidis (2013), using the standard volume of induced submatrices to define sampling probabilities. Although the wider class of Gvs measures have not achieved

the popularity of DPPs in machine learning, they nonetheless have several key applications to a variety of fundamental optimization problems, and deserve attention in their own right.

The main contribution of this chapter lies in the introduction of the class of Gvs measures. We begin defining these measures formally and deriving their negative dependence properties. Then, we focus on their use for the experimental design problem described above, deriving probabilistic and greedy algorithms to choose sets of experiments. Finally, we conclude by tying Gvs measures and elementary symmetric polynomials to a generalized form of the Minimum Volume Covering Ellipsoid. All material in this chapter was first derived in (Mariet and Sra, 2017).

## 5.1 Background

Generalized volume sampling measures are defined based on elementary symmetric polynomials (ESPs), a fundamental family of polynomials that has been extensively studied in matrix theory, combinatorics, information theory, as well as several other areas — see *e.g.*, (Jozsa and Mitchison, 2015; Macdonald, 1998; Horn and Johnson, 1985; Bhatia, 2007; Jain, 2011; Bauschke et al., 2001).

**Definition 5.1.1** (Elementary symmetric polynomial). The  $\ell$ -th elementary symmetric polynomial over  $\mathbb{R}^m$  is defined as

$$e_\ell(x_1, \dots, x_m) \triangleq \sum_{I \subseteq [m], |I|=\ell} \prod_{j \in I} x_j, \quad (5.1)$$

Eq. (5.1) extends to matrices naturally: elementary symmetric polynomials are spectral functions, and so we will also write

$$e_\ell(\mathbf{M}) \triangleq e_\ell \circ \lambda(\mathbf{M}).$$

**Remark 5.1.1.** We have the following special identities for elementary symmetric polynomials:  $e_0 \equiv 1$ ,  $e_1 = \text{Tr}$ ,  $e_m = \det$ ,  $e_{\ell > m} \equiv 0$ .

**Remark 5.1.2.** Elementary symmetric polynomials also provide the partition function for  $k$ -DPPs (Definition 3.1.2).

The following proposition captures basic properties of ESPs that we will require in our analysis; the proofs can be found in standard works such as (Bhatia, 2007).

**Proposition 5.1.1.** *Let  $\mathbf{M} \in \mathbb{R}^{m \times m}$  be symmetric and  $1 \leq \ell \leq m$ ; also let  $\mathbf{A}, \mathbf{B} \in \mathbb{S}_m^+$ . We have the following properties:*

- (i) *If  $\mathbf{A} \succeq \mathbf{B}$  in Löwner order, then  $e_\ell(\mathbf{A}) \geq e_\ell(\mathbf{B})$ ;*
- (ii) *If  $\mathbf{M}$  is invertible, then  $e_\ell(\mathbf{M}^{-1}) = \det(\mathbf{M}^{-1})e_{m-\ell}(\mathbf{M})$ ;*
- (iii)  *$\nabla e_\ell(\lambda) = [e_{\ell-1}(\lambda^{(i)})]_{1 \leq i \leq m}$  where  $\lambda^{(i)}$  is vector  $\lambda$  without its  $i$ -th component.*

Elementary symmetric polynomials enjoy another representation that allows us to interpret them as “partial volumes”: namely,

$$e_\ell(\mathbf{M}) = \sum_{S \subseteq [m], |S|=\ell} \det(\mathbf{M}_S). \quad (5.2)$$

From this notion of partial volumes, we are now able to properly define Generalized Volume Sampling measures, that sample subsets of rows from a tall, thin matrix  $\mathbf{X}$  based on their spanned partial volume.

**Definition 5.1.2** (Generalized Volume Sampling). A distribution  $\mu$  over  $2^{[N]}$  is a GVS measure if there exists  $\mathbf{X} \in \mathbb{R}^{N \times m}$  and an integer  $\ell \in \{0, \dots, m\}$  such that

$$\mu(S) \propto e_\ell(\mathbf{X}_{S,:}^\top \mathbf{X}_{S,:}). \quad (5.3)$$

**Remark 5.1.3.** Writing  $x_1, \dots, x_N \in \mathbb{R}^m$  the row vectors of  $\mathbf{X}$ , the probability defined in (5.3) is equivalent to

$$\mu(S) \propto e_\ell\left(\sum_{i \in S} x_i x_i^\top\right). \quad (5.4)$$

The special case  $\ell = m$  (*i.e.*,  $e_\ell = \det$ ) has benefited from nascent interest in machine learning; first studied by Avron and Boutsidis (2013), polynomial-time sampling

algorithms were developed in (Li et al., 2017) for these measures, under the name of *Dual Volume Sampling* which informs our choice of Generalized Volume Sampling to name the larger class in Definition 5.1.2.

### 5.1.1 Negative dependence properties

That Dual Volume Sampling measures ( $\ell = m$ ) are strongly Rayleigh was proven in (Li et al., 2017). To prove that *all* Gvs measures are Strong Rayleigh, we follow the same steps as Li et al. (2017), making a few modifications to generalize to all elementary symmetric polynomials.

**Proposition 5.1.2.** *For  $m \leq |S| \leq N$  and  $\mathbf{L} = \mathbf{X}\mathbf{X}^\top$  with  $\mathbf{X} \in \mathbb{R}^{N \times m}$ , we have*

$$e_\ell(\mathbf{X}_{S,:}^\top \mathbf{X}_{S,:}) = e_\ell(\mathbf{L}_{S,S}).$$

*Proof.* Let  $\mathbf{Y} = \text{Diag}(y_1, \dots, y_N)$  be a diagonal matrix. The generalized Cauchy-Binet identity yields

$$e_\ell(\mathbf{X}^\top \mathbf{Y} \mathbf{X}) = \sum_{\substack{T \subseteq [m] \\ |T|=\ell}} \det([\mathbf{X}^\top \mathbf{Y} \mathbf{X}]_{T,T}) = \sum_{\substack{T \subseteq [m] \\ |T|=\ell}} \sum_{\substack{K \subseteq [N] \\ |K|=\ell}} \det([\mathbf{X}^\top \mathbf{Y}]_{T,K}) \det(\mathbf{X}_{K,T}).$$

For  $\mathbf{Y}$  the indicator matrix of set  $S \subseteq [N]$ , all terms with  $K \not\subseteq S$  disappear, yielding

$$e_\ell(\mathbf{X}_{S,:}^\top \mathbf{X}_{S,:}) = \sum_{\substack{T \subseteq [m] \\ |T|=\ell}} \sum_{\substack{K \subseteq S \\ |K|=\ell}} \det(\mathbf{X}_{T,K}^\top) \det(\mathbf{X}_{K,T}) = \sum_{\substack{K \subseteq S \\ |K|=\ell}} \det([\mathbf{X} \mathbf{X}^\top]_{K,K}) = e_\ell([\mathbf{X} \mathbf{X}^\top]_{S,S}),$$

which concludes the proof.  $\square$

**Theorem 5.1.3** (Gvs measures are SR). *Let  $\mathbf{X} \in \mathbb{R}^{N \times m}$ ,  $k \in \{m, \dots, N\}$  and  $\ell \in \{1, \dots, m\}$ . The multi-affine polynomial  $g$ , defined over  $\mathbb{C}^n$  as*

$$g(z_1, \dots, z_N) = \sum_{|S|=k, S \subseteq [N]} e_\ell(\mathbf{X}_{S,:}^\top \mathbf{X}_{S,:}) z^S,$$

*is real stable.*

*Proof.* Let  $\mathbf{L} = \mathbf{X}\mathbf{X}^\top \in \mathbb{R}^{N \times N}$ . Writing  $\mathbf{W}$  for the diagonal matrix  $\text{Diag}(w_1, \dots, w_N)$  and  $\mathbf{Z} = \text{Diag}(z_1, \dots, z_N)$ , we know from Prop. 2.2.10 that

$$g(z, w) = \det(\mathbf{L} + \mathbf{W} + \mathbf{Z}) = \sum_{S \subseteq [N]} \det(\mathbf{L}_S + \mathbf{W}_S) z^S = \sum_{S \subseteq [N]} \left( \sum_{T \subseteq S} \det(\mathbf{L}_T) w^T \right) z^S$$

is real stable in  $2N$  variables;  $g(z, w)$  can be expanded as

$$g(z, w) = \sum_{S \subseteq [N]} \left( \sum_{T \subseteq S} \det(\mathbf{L}_T) w^{S \setminus T} \right) z^{\bar{S}}.$$

Truncating  $g(z, w)$  to degree  $N - m$ , which constrains  $|T| = m$ , then setting  $w = (y, \dots, y)$  yields

$$g(z, y) = \sum_{S \subseteq [N]} \left( \sum_{\substack{T \subseteq S \\ |T|=m}} \det(\mathbf{L}_T) \right) y^{|S|-m} z^{\bar{S}} = \sum_{S \subseteq [N]} e_\ell(\mathbf{L}_S) y^{|S|-m} z^{\bar{S}},$$

which once again remains real stable (Prop. 2.2.7). Finally,  $k - m$  differentiations of  $g(z, y)$  according to  $y$  followed by evaluating at  $y = 0$  (both operations preserving stability as stated in Prop. 2.2.5) yields polynomial

$$h(z) = \sum_{S \subseteq [N], |S|=k} e_\ell(\mathbf{L}_S) z^{\bar{S}}.$$

By applying Prop. 2.2.5–(iv) we recover the generating polynomial for the GVS measure of degree  $\ell$  (Prop. 5.1.2), concluding our proof.  $\square$

We now discuss a key applications of GVS measures to machine learning: optimal experimental design.

## 5.2 Optimal experimental design

Optimal experimental design (OED) develops the theory of selecting a small subset of experiments to perform in order evaluate a hidden variable as well as possible. As discussed earlier, OED has fundamental ties to negative dependence. Indeed,

to evaluate a multi-dimensional parameter  $\theta$  in the fewest number of experiments possible, each experiment must provide new information about  $\theta$ , making selecting similar experiments a sub-optimal choice.

### 5.2.1 Generalized Volume Sampling and optimal design

OED operates under the assumption that experiments are costly and cannot be run many times or even once without tremendous difficulty (Pukelsheim, 2006). OED has been applied in a large number of experimental settings (Sagnol, 2010; Cohn, 1994; Lasheras et al., 2010; Xygkis et al., 2016; Schein and Ungar, 2004), and has close ties to related machine learning problems such as outlier detection (Dolia et al., 2004; Jackson and Chen, 2004), active learning (He, 2010; Gu and Jin, 2013), and Gaussian process driven sensor placement (Krause et al., 2008) among others.

The OED problem is formulated as follows: there exist  $N$  experiments whose outcome  $y$  depends linearly on a hidden parameter  $\theta \in \mathbb{R}^m$  according to (5.5):

$$y_i = x_i^\top \theta + \epsilon_i \quad 1 \leq i \leq n, \quad (5.5)$$

where  $x_i \in \mathbb{R}^m$  represents the experiment parameters, and the  $\epsilon_i$  are independent zero mean homoscedastic Gaussian noise vectors. Within this setting, OED's goal is to select a set  $S$  of  $k$  experiments that estimate  $\theta$  with no bias and minimal variance.

The Gauss-Markov theorem guarantees that for a set  $S$  such that  $\sum_{i \in S} x_i x_i^\top$  is invertible (in which case  $S$  is called a *feasible* set), the lowest variance for an unbiased estimate  $\hat{\theta}$  verifies

$$\text{Var}[\hat{\theta}] = \left( \sum_{i \in S} x_i x_i^\top \right)^{-1}.$$

Of course, this variance is a positive semi-definite matrix, over the set of which there is no total order. To formally seek to minimize the variance, we require a suitable map  $\Phi$  that maps (positive semi-definite) matrices to scalars. Hence, OED is cast as an optimization problem that seeks *an optimal design*  $S^*$  over the space  $\Gamma_k = \{S \subseteq$

$[N], |S| \leq k, \sum_{i \in S} x_i x_i^\top$  is invertible} of feasible designs:

$$\min_{S \in \Gamma_k} f_\ell(S) \triangleq \Phi\left(\left(\sum_{i \in S} x_i x_i^\top\right)^{-1}\right), \quad (5.6)$$

where  $\Phi$  is a cost function that maps positive definite matrices to scalars.

The aim of (5.6) is thus to pick a subset  $S$  of  $k$  points out of  $N$  (typically  $k \ll N$ ) that reduce the uncertainty as much as possible, measured according to the cost function  $\Phi$ ; the choice of  $\Phi$  elicits different properties that a solution should satisfy, either statistical or structural.

The theory of OED branches into multiple variants of (5.6) depending on the choice of  $\Phi$ , among which A-optimal design ( $\Phi = \text{trace}$ ) and D-optimal design ( $\Phi = \text{determinant}$ ) are probably the two most popular choices. Each of these choices has a wide range of applications as well as statistical, algorithmic, and other theoretical results. Due to the individual popularity of A- and D-optimal design, the theory surrounding these two sub-problems has diverged significantly.

As both the trace and the determinant are special cases of elementary symmetric polynomials, we generalize (5.6) to a broader optimization problem where  $\Phi$  is obtained from an elementary symmetric polynomial; this allows us to consider A-optimal design and D-optimal design as special cases and treat the entire class of problems in a unified manner; we refer to this class of problems as *ESP-design*.

### 5.2.2 Problem definition

Let  $\mathbf{X} \in \mathbb{R}^{N \times m}$  ( $m \ll N$ ) be a matrix with full column rank, and  $k \in \mathbb{N}$  satisfy  $m \leq k \leq N$ . For  $\ell \in \{1, \dots, m\}$ , we introduce the *ESP-design* problem:

$$\min_{S \in \Gamma_k} f_\ell(S) \triangleq \frac{1}{\ell} \log e_\ell\left((\mathbf{X}_{S,:}^\top \mathbf{X}_{S,:})^{-1}\right). \quad (5.7)$$

We keep the  $1/\ell$ -factor in (5.7) to highlight the homogeneity ( $e_\ell$  is a polynomial of degree  $\ell$ ) of our design criterion, as is advocated in (Pukelsheim, 2006, Ch. 6).

For  $\ell = 1$ , (5.7) yields A-optimal design, while for  $\ell = m$ , it yields D-optimal

design. For  $1 < \ell < m$ , ESP-design interpolates between these two extremes. Geometrically, we may view it as seeking an ellipsoid with the smallest average volume for  $\ell$ -dimensional slices (taken across sets of size  $\ell$ ).

Alternatively, ESP-design can be also be interpreted as a regularized version of D-optimal design via Prop. 4.2.1-(ii). In particular, for  $\ell = m - 1$ , we recover a form of regularized D-optimal design:

$$f_{m-1}(S) = \frac{1}{m-1} [\log \det((\mathbf{X}_{S,:}^\top \mathbf{X}_{S,:})^{-1}) + \log \|\mathbf{X}_{S,:}\|_2^2].$$

Furthermore, using Prop. 5.1.1 (ii), we can also rephrase the OED problem with  $\Phi \equiv e_\ell$  as finding the mode of the product of Gvs measures given by

$$S \rightarrow \det(\mathbf{X}_{S,:}^\top \mathbf{X}_{S,:})^{-1} e_{m-\ell}(\mathbf{X}_{S,:}^\top \mathbf{X}_{S,:}).$$

Note that (5.7) is a hard combinatorial optimization problem (the case  $\ell = m$ , which also benefits from submodularity, is known to be NP-hard (Welch, 1982)), which precludes an exact solution. However, the objective function in its general ESP form enjoys remarkable properties; the following section focuses on the properties of the continuous relaxation of (5.7), which will be crucial to solving the discrete optimization problem.

### 5.2.3 Continuous optimization over elementary symmetric polynomials

In A-optimal design, a common approach to solving (5.7) consists in relaxing the constraint on  $S$  by allowing elements in the set to have fractional multiplicities. In the general case of ESP-design, this relaxation leads to the following continuous optimization problem:

$$\min_{z \in \Gamma_k^c} \frac{1}{\ell} \log e_\ell \left( (\mathbf{X}^\top \text{Diag}(z) \mathbf{X})^{-1} \right), \quad (5.8)$$

where by  $\Gamma_k^c$  we denote the set of vectors  $\{z \in \mathbb{R}^n \mid 0 \leq z_i \leq 1\}$  such that  $\mathbf{X}^\top \text{Diag}(z)\mathbf{X}$  remains invertible and  $\mathbf{1}^\top z \leq k$ . The following proposition is a direct consequence of Proposition 4.2.1(i).

**Proposition 5.2.1.** *Let  $z^*$  be the optimal solution to (5.8). Then  $\|z^*\|_1 = k$ .*

The convexity of  $\log \text{tr}((\mathbf{X}^\top \text{Diag}(z)\mathbf{X})^{-1})$  is a key ingredient to obtaining high-quality A-optimal designs (Wang et al., 2016a). The rest of this section focuses on establishing the convexity of  $f_\ell$  for general  $\ell$  on  $\Gamma_k^c$ , where by abuse of notation,  $f_\ell$  also denotes the continuous relaxation in (5.8).

Although convexity of  $f_\ell$  may be obtained by as a consequence of a classic result of Muir (1974), we present a new way of deriving it, based on a result on the geodesic log-convexity of functions of elementary symmetric polynomials.

We begin by deriving the geodesic convexity (see Appendix A.1.1) of  $\log e_\ell$ . Note that geodesic convexity of  $\log e_\ell$  on positive definite matrices is strictly stronger than (and implies) the convexity of  $f_\ell$ .

**Theorem 5.2.2.** *The function  $\mathbf{X} \mapsto e_\ell(\mathbf{X}^s)$  for  $s \in \{-1, 1\}$  is log-g-convex on positive definite matrices.*

*Proof.* By continuity, it suffices to prove midpoint log geodesic convexity; that is, it suffices to prove

$$e_\ell(\mathbf{P}\#\mathbf{Q}) \leq \sqrt{e_\ell(\mathbf{P})e_\ell(\mathbf{Q})}.$$

Since  $(\mathbf{P}\#\mathbf{Q})^{-1} = \mathbf{P}^{-1}\#\mathbf{Q}^{-1}$ , the same proof shows log geodesic convexity of  $e_\ell(\mathbf{X}^{-1})$ .

From multilinear algebra (*e.g.*, (Bhatia, 1997, Ch. 1)) we know that for an  $N \times N$  matrix  $\mathbf{P}$ , there exists a matrix  $\mathbf{W}$  such that  $e_\ell(\mathbf{P}) = \text{tr } \mathbf{W}^* \mathbf{P}^{\otimes N} \mathbf{W}$ , where  $\otimes$  denotes the Kronecker product. Lemma 2.23 in (Sra and Hosseini, 2015) shows that

$$(\mathbf{P}\#\mathbf{Q})^{\otimes N} = \mathbf{P}^{\otimes N}\#\mathbf{Q}^{\otimes N},$$

from which it follows that

$$e_\ell(\mathbf{P}\#\mathbf{Q}) = \text{tr } \mathbf{W}^* (\mathbf{P}\#\mathbf{Q})^{\otimes N} \mathbf{W} = \text{tr } \mathbf{W}^* [\mathbf{P}^{\otimes N}\#\mathbf{Q}^{\otimes N}] \mathbf{W}$$

$$\begin{aligned} &\leq [\text{tr } \mathbf{W}^* \mathbf{P}^{\otimes N} \mathbf{W}]^{1/2} [\text{tr } \mathbf{W}^* \mathbf{Q}^{\otimes N} \mathbf{W}]^{1/2} \\ &= [e_\ell(\mathbf{P}) e_\ell(\mathbf{Q})]^{1/2}, \end{aligned}$$

where the inequality follows from log geodesic convexity of the trace (Sra and Hosseini, 2015, Cor. 2.9).  $\square$

Theorem 5.2.2 immediately yields the following corollary:

**Corollary 5.2.2.1.** *The map  $\mathbf{X} \mapsto \log e_\ell(\mathbf{X})$  is convex on positive definite matrices.*

As a corollary of the log-convexity of  $e_\ell$ , we obtain that (5.8) is a convex optimization problem, which can therefore be solved using a variety of efficient algorithms. Projected gradient descent turns out to be particularly easy to apply because we only require projection onto the intersection of the cube  $\mathbf{0} \leq z \leq \mathbf{1}$  and the plane  $\{z \mid z^\top \mathbf{1} = k\}$  (as a consequence of Prop 5.2.1). Projection onto this intersection is a special case of the so-called continuous quadratic knapsack problem, which is a very well-studied problem and can be solved essentially in linear time (Cominetti et al., 2014; Davis et al., 2016).

**Remark 5.2.1.** The convex relaxation remains log-convex when points can be chosen with multiplicity, in which case the projection step is also significantly simpler, requiring only  $z \geq \mathbf{0}$ .

We conclude the analysis of the continuous relaxation by a bound on the support of its solution under some mild assumptions; this bound will serve as a first step to satisfying the set size constraint of the original ESP-design problem.

**Theorem 5.2.3.** *Let  $\phi$  be the mapping from  $\mathbb{R}^m$  to  $\mathbb{R}^{m(m+1)/2}$  such that  $\phi(x) = (\xi_{ij} x_i x_j)_{1 \leq i,j \leq m}$  with  $\xi_{ij} = 1$  if  $i = j$  and 2 otherwise. Let  $\tilde{\phi}(x) = (\phi(x), 1)$  be the affine version of  $\phi$ . If for any set of  $m(m+1)/2$  distinct rows of  $\mathbf{X}$ , the mapping under  $\tilde{\phi}$  is independent, then the support of the optimum  $z^*$  of (5.8) satisfies  $\|z^*\|_0 \leq k + \frac{m(m+1)}{2}$ .*

The proof is identical to that of (Wang et al., 2016a, Lemma 3.5), which shows such a result for  $\ell = 1$ .

Theorem 5.2.3 shows that the support of the solution to the continuous relaxation of (5.7) yields a subset  $S_0$  of bounded size; by construction,  $S_0$  is at least as good as an optimal solution to (5.7). To bring down the size of the continuous relaxation solution  $S_0$  to satisfy the size constraint  $|S| \leq k$ , we now turn to combinatorial optimization.

### 5.2.4 Finding an optimal design

The combinatorial form of ESP-design admits an intuitive greedy approach despite not being a submodular optimization problem in general. In this greedy form, elements are removed one-by-one from a base set of experiments; greedy removal, as opposed to greedy addition, turns out to be much more practical.

Indeed, since  $f_\ell$  is not defined for sets of size smaller than  $k$ , we cannot quantify the quality of intermediary (of size  $< k$ ) sets in a greedy addition algorithm. This difficulty precludes analyses such as (Wang et al., 2016b; Smith and Thai, 2017) for optimizing non-submodular set functions by bounding their “curvature”.

---

**Algorithm 12** Greedy algorithm

---

**Input:** Matrix  $\mathbf{X}$ , budget  $k$ , initial set  $S_0$

$S \leftarrow S_0$

**while**  $|S| > k$  **do**

    Find  $i \in S$  such that  $S \setminus \{i\}$  is feasible and  $i$  minimizes  $f_\ell(S \setminus \{i\})$

$S \leftarrow S \setminus \{i\}$

**return**  $S$

---

Instead, we focus on the greedy removal algorithm described in Algorithm 12; obtaining a bound on its performance relies on the following algorithm.

**Lemma 5.2.4.** *Let  $\mathbf{X} \in \mathbb{R}^{N \times m}$  ( $N \geq m$ ) be a matrix with full column rank, and let  $k$  be a budget  $m \leq k \leq N$ . Let  $S$  of size  $k$  be subset of  $[N]$  drawn with probability  $\mathcal{P} \propto \det(\mathbf{X}_{S,:}^\top \mathbf{X}_{S,:})$ . Then*

$${}_{S \sim \mathcal{P}} \left[ e_\ell \left( (\mathbf{X}_{S,:}^\top \mathbf{X}_{S,:})^{-1} \right) \right] \leq \prod_{i=1}^k \frac{N-m+i}{k-m+i} \cdot e_\ell \left( (\mathbf{X}^\top \mathbf{X})^{-1} \right), \quad (5.9)$$

with equality if  $\mathbf{X}_{S,:}^\top \mathbf{X}_{S,:} \succ 0$  for all subsets  $S$  of size  $k$ .

Some simple manipulations using Stirling's approximation show that the ratio  $\prod_{i=1}^{\ell} \frac{N-m+i}{k-m+i}$  is upper bounded by  $(1 + \frac{\ell}{N-m})^{N-k} (1 + \frac{N-k}{k-m+\ell})^{\ell}$ .

*Proof.* The below calculations depend heavily on the Cauchy-Binet formula, of which we reproduce a special case here for  $\mathbf{X} \in \mathbb{R}^{N \times m}$ :

$$\det(\mathbf{X}^\top \mathbf{X}) = \sum_{S \subseteq [N], |S|=m} \det(\mathbf{X}_{S,:}^\top \mathbf{X}_{S,:}). \quad (5.10)$$

We also use the representation (5.2). By definition we have

$$\mathbb{E} \left[ e_\ell \left( (\mathbf{X}_{S,:}^\top \mathbf{X}_{S,:})^{-1} \right) \right] = \frac{\sum_{S \subseteq [N], |S|=k} \det(\mathbf{X}_{S,:}^\top \mathbf{X}_{S,:}) e_\ell \left( (\mathbf{X}_{S,:}^\top \mathbf{X}_{S,:})^{-1} \right)}{\sum_{S \subseteq [N], |S|=k} \det(\mathbf{X}_{S,:}^\top \mathbf{X}_{S,:})}.$$

For the denominator, we have

$$\begin{aligned} \sum_{S \subseteq [N], |S|=k} \det(\mathbf{X}_{S,:}^\top \mathbf{X}_{S,:}) &= \sum_{S \subseteq [N], |S|=k} \det(\mathbf{X}_{S,:}^\top \mathbf{X}_{S,:}) \\ &\stackrel{(a)}{=} \sum_{S \subseteq [N], |S|=k} \sum_{T \subseteq S, |T|=m} \det(\mathbf{X}_{T,:}^\top \mathbf{X}_{T,:}) \\ &\stackrel{(b)}{=} \binom{N-m}{k-m} \sum_{T \subseteq S, |T|=m} \det(\mathbf{X}_{T,:}^\top \mathbf{X}_{T,:}) \\ &\stackrel{(c)}{=} \binom{N-m}{k-m} \det(\mathbf{X}^\top \mathbf{X}), \end{aligned}$$

where (a) is obtained via the Cauchy-Binet formula (5.10), (b) by noticing that there are  $\binom{N-m}{k-m}$  sets of size  $k$  that contain a set  $T$  of size  $m$ , and (c) by reapplying (5.10).

For the numerator, we first use the fact that  $e_\ell(A^{-1}) = \frac{1}{\det A} e_{m-\ell}(A)$ :

$$\begin{aligned} \sum_{S \subseteq [N], |S|=k} \det(\mathbf{X}_{S,:}^\top \mathbf{X}_{S,:}) e_\ell \left( (\mathbf{X}_{S,:}^\top \mathbf{X}_{S,:})^{-1} \right) &\stackrel{(a)}{\leq} \sum_{S \subseteq [N], |S|=k} e_{m-\ell}(\mathbf{X}_{S,:}^\top \mathbf{X}_{S,:}) \\ &= \sum_{S \subseteq [N], |S|=k} \sum_{L \subseteq [m], |L|=m-\ell} (\mathbf{X}_{S,:}^\top \mathbf{X}_{S,:})[L|L] \end{aligned}$$

$$\begin{aligned}
&\stackrel{(b)}{=} \sum_{S \subseteq [N], |S|=k} \sum_{L \subseteq [m], |L|=m-\ell} \det((\mathbf{Y}_L)_{S,:}^\top (\mathbf{Y}_L)_{S,:}) \\
&\stackrel{(c)}{=} \sum_{S \subseteq [N], |S|=k} \sum_{L \subseteq [m], |L|=m-\ell} \sum_{T \subseteq S, |T|=m-\ell} \det((\mathbf{Y}_L)_{T,:}^\top (\mathbf{Y}_L)_{T,:}).
\end{aligned}$$

Here, (a) is just (5.2); we have equality if all subsets  $S$  of size  $k$  produce strictly positive definite matrices  $\mathbf{X}_{S,:}^\top \mathbf{X}_{S,:}$ . For (b), we note  $\mathbf{Y}_L$  the submatrix of  $\mathbf{X}$  with all rows but those in  $L$  removed; then,  $(\mathbf{Y}_L)_{S,:}^\top (\mathbf{Y}_L)_{S,:} = [\mathbf{X}_{S,:}^\top \mathbf{X}_{S,:}]$  for all subsets  $S$ .

By counting subset occurrences and reapplying the Cauchy-Binet formula, we have

$$\begin{aligned}
&\sum_{S \subseteq [N], |S|=k} \det(\mathbf{X}_{S,:}^\top \mathbf{X}_{S,:}) e_\ell \left( (\mathbf{X}_{S,:}^\top \mathbf{X}_{S,:})^{-1} \right) \\
&\leq \binom{N-m+\ell}{k-m+\ell} \sum_{L \subseteq [m], |L|=m-\ell} \sum_{T \in [n], |T|=m-\ell} \det((\mathbf{Y}_L)_{T,:}^\top (\mathbf{Y}_L)_{T,:}) \\
&\stackrel{(d)}{=} \binom{N-m+\ell}{k-m+\ell} \sum_{L \subseteq [m], |L|=m-\ell} \det((\mathbf{Y}_{L,:})^\top (\mathbf{Y}_{L,:})) \\
&= \binom{N-m+\ell}{k-m+\ell} \sum_{L \subseteq [m], |L|=m-\ell} (\mathbf{X}^\top \mathbf{X})[L|L] \\
&= \binom{N-m+\ell}{k-m+\ell} e_{m-\ell}(\mathbf{X}^\top \mathbf{X}).
\end{aligned}$$

Hence, we can finally bound the expected value of under the objective function of the sampled subset  $S$  as

$$\begin{aligned}
\mathbb{E} \left[ e_\ell \left( (\mathbf{X}_{S,:}^\top \mathbf{X}_{S,:})^{-1} \right) \right] &= \frac{\sum_{S \subseteq [N], |S|=k} \det(\mathbf{X}_{S,:}^\top \mathbf{X}_{S,:}) e_\ell \left( (\mathbf{X}_{S,:}^\top \mathbf{X}_{S,:})^{-1} \right)}{\sum_{S \subseteq [N], |S|=k} \det(\mathbf{X}_{S,:}^\top \mathbf{X}_{S,:})} \\
&= \frac{\binom{N-m+\ell}{k-m+\ell} e_{m-\ell}(\mathbf{X}^\top \mathbf{X})}{\binom{k-m}{N-m} \det(\mathbf{X}^\top \mathbf{X})} \\
&= \left( \prod_{i=1}^{\ell} \frac{N-m+i}{k-m+i} \right) e_\ell((\mathbf{X}^\top \mathbf{X})^{-1}).
\end{aligned}$$

□

Lemma 5.2.4 extends a result from (Avron and Boutsidis, 2013, Lemma 3.9) on column-subset selection via volume sampling to all elementary symmetric polynomials.

In particular, it follows that removing one element (by volume sampling a set of size  $N - 1$ ) will in expectation decrease  $f$  by a multiplicative factor which is clearly also attained by a greedy minimization. From this argument, we obtain the following bound on the performance of Algorithm 12.

**Theorem 5.2.5.** *Algorithm 12 initialized with a set  $S_0$  of size  $N_0$  produces a set  $S^+$  of size  $k$  such that*

$$e_\ell\left(\left(\mathbf{X}_{S^+,:}^\top \mathbf{X}_{S^+,:}\right)^{-1}\right) \leq \prod_{j=1}^{\ell} \frac{N_0 - m + j}{k - m + j} \cdot e_\ell\left(\left(\mathbf{X}_{S_0,:}^\top \mathbf{X}_{S_0,:}\right)^{-1}\right). \quad (5.11)$$

*Proof.* We recursively show that greedily removing  $j$  items constructs a set  $S$  (of size  $N - j$ ) such that

$$e_\ell\left(\left(\mathbf{X}_{S,:}^\top \mathbf{X}_{S,:}\right)^{-1}\right) \leq \left(\prod_{i=1}^{\ell} \frac{N - m + i}{N - j - m + i}\right) e_\ell\left(\left(\mathbf{X}^\top \mathbf{X}\right)^{-1}\right). \quad (5.12)$$

(5.12) is trivially true for  $j = 0$ . Assume now that (5.12) holds for  $j \geq 0$ , and let  $S_j$  be the corresponding set of size  $(N - j)$ . Let now  $S_{j+1}$  be the set of size  $|S_j| - 1$  that minimizes  $e_\ell(\mathbf{X}_{S_{j+1},:}^\top \mathbf{X}_{S_{j+1},:})$ .

From lemma 5.2.4, we know that for sets  $S$  of size  $|S_j| - 1$  drawn according to dual volume sampling,

$$\mathbb{E}\left[e_\ell\left(\left(\mathbf{X}_{S,:}^\top \mathbf{X}_{S,:}\right)^{-1}\right)\right] \leq \left(\prod_{i=1}^{\ell} \frac{|S_j| - m + i}{(|S_j| - 1) - m + i}\right) e_\ell\left(\left(\mathbf{X}_{S,:}^\top \mathbf{X}_{S,:}\right)^{-1}\right).$$

In particular, the minimum of  $e_\ell(\mathbf{X}_{S,:}^\top \mathbf{X}_{S,:})$  over all sets of size  $|S_j| - 1$  is upper bounded by the expectancy:

$$e_\ell(\mathbf{X}_{S_{j+1},:}^\top \mathbf{X}_{S_{j+1},:}) \leq \left(\prod_{i=1}^{\ell} \frac{N - j - m + i}{N - j - 1 - m + i}\right) e_\ell\left(\left(\mathbf{X}_{S,:}^\top \mathbf{X}_{S,:}\right)^{-1}\right).$$

By recursion hypothesis applied to  $S_j$ , we then have

$$e_\ell(\mathbf{X}_{S_{j+1},:}^\top \mathbf{X}_{S_{j+1},:}) \leq \left(\prod_{i=1}^{\ell} \frac{N - j - m + i}{N - j - 1 - m + i}\right) e_\ell\left(\left(\mathbf{X}_{S,:}^\top \mathbf{X}_{S,:}\right)^{-1}\right)$$

$$\begin{aligned}
&\leq \left( \prod_{i=1}^{\ell} \frac{N-j-m+i}{N-j-1-m+i} \right) \left( \prod_{i=1}^{\ell} \frac{n-m+i}{n-j-m+i} \right) e_{\ell}\left( (\mathbf{X}^\top \mathbf{X})^{-1} \right) \\
&\leq \left( \prod_{i=1}^{\ell} \frac{N-m+i}{N-(j+1)-m+i} \right) e_{\ell}\left( (\mathbf{X}^\top \mathbf{X})^{-1} \right),
\end{aligned}$$

which concludes the recursion. Then, constructing a set of size  $k$  amounts to setting  $j = n - k$  in inequality (5.12), proving the inequality (5.11).  $\square$

As Wang et al. (2016a) note regarding A-optimal design, (5.11) provides a trivial optimality bound on the greedy algorithm when initialized with  $\{1, \dots, N\}$ : denoting by  $S^*$  the optimal set, the greedy solution  $S^+$  satisfies

$$\begin{aligned}
E_{\ell}\left( (\mathbf{X}_{S^+, :}^\top \mathbf{X}_{S^+, :})^{-1} \right)^{1/\ell} &\leq \frac{N-m+\ell}{k-m+1} f([N]) \\
&\leq \frac{N-m+\ell}{k-m+1} E_{\ell}\left( (\mathbf{X}_{S^*, :}^\top \mathbf{X}_{S^*, :})^{-1} \right)^{1/\ell}.
\end{aligned}$$

This naive initialization can be replaced by the support  $\|z^*\|_0$  of the convex relaxation solution; in the setting described by Theorem 5.2.3, we obtain the following result:

**Theorem 5.2.6.** *Let  $\tilde{\phi}$  be the mapping defined in 5.2.3, and assume that all choices of  $m(m+1)/2$  distinct rows of  $\mathbf{X}$  always have their mapping independent mappings for  $\tilde{\phi}$ . Then the outcome of the greedy algorithm initialized with the support of the solution to the continuous relaxation verifies*

$$f_{\ell}(S^+) \leq \log \left( \frac{k + m(m-1)/2 + \ell}{k-m+1} \right) + f_{\ell}(S^*).$$

Computing the  $\ell$ -th elementary symmetric polynomial on a vector of size  $m$  can be done in  $\mathcal{O}(m \log^2 \ell)$  using Fast Fourier Transform for polynomial multiplication, using the construction by Ben-Or (Shpilka and Wigderson, 2001). Hence, computing  $f_{\ell}(S)$  requires  $\mathcal{O}(nm^2)$  operations, where the cost is dominated by computing  $\mathbf{X}_{S, :}^\top \mathbf{X}_{S, :}$ .

Hence, by combining the continuous relaxation and the combinatorial optimization of Alg. 12, we obtain a greedy algorithm with performance guarantees that runs in  $\mathcal{O}(m^2 N^3)$  time, plus the additional cost of solving the continuous relaxation.

## 5.3 Experiments: optimal design

In order to validate experimentally our analysis, we compared the following methods to solving the optimization problem (5.7):

- UNIF/UNIFFDV:  $k$  experiments sampled uniformly / with Fedorov exchange
- GREEDY/GREEDYFDV: greedy method (relaxed init.) / with Fedorov exchange
- SAMPLE: sampling (relaxed init.) items independently from the Bernoulli distribution represented by the solution to the continuous relaxation.

We also report the results for solution of the continuous relaxation (**RELAX**); the convex optimization was solved using projected gradient descent, the projection being done with the code from (Davis et al., 2016).

### 5.3.1 Synthetic experiments: optimization comparison

We generated the experimental matrix  $\mathbf{X}$  by sampling  $N$  vectors of size  $m$  from the multivariate Gaussian distribution of mean 0 and sparse precision  $\Sigma^{-1}$  (density  $d$  ranging from 0.3 to 0.9). Due to the runtime of Fedorov methods, results are reported for only one run; results averaged over multiple iterations (as well as for other distributions over  $\mathbf{X}$ ) are provided in Appendix A.2.3.

As shown in Fig. 5-1, the greedy algorithm applied to the convex relaxation’s support outperforms sampling from the convex relaxation solution, and does as well as the usual Fedorov algorithm UNIFFDV; GREEDYFDV marginally improves upon the greedy algorithm and UNIFFDV. Strikingly, GREEDY provides designs of comparable quality to UNIFFDV; furthermore, as very few local exchanges improve upon its design, running the Fedorov algorithm with GREEDY initialization is much faster (Table 5.1); this is confirmed by Table 5.2, which shows the number of experiments in common for different algorithms: GREEDY and GREEDYFDV only differ on very few elements. As the budget  $k$  increases, the difference in performances between SAMPLE, GREEDY and the continuous relaxation decreases, and the simpler SAMPLE algorithm becomes

competitive. Table 5.3 reports the support of the continuous relaxation solution for ESP-design with  $\ell = 10$ .

Table 5.1: Runtimes (s) ( $\ell = 10, d = 0.6$ )

$k$	40	80	120	160	200
GREEDY	$2.8 \cdot 10^1$	$2.7 \cdot 10^1$	$3.1 \cdot 10^1$	$4.0 \cdot 10^1$	$5.2 \cdot 10^1$
GREEDYFDV	$6.6 \cdot 10^1$	$2.2 \cdot 10^2$	$3.2 \cdot 10^2$	$1.2 \cdot 10^2$	$1.3 \cdot 10^2$
UNIFFDV	$1.6 \cdot 10^3$	$4.1 \cdot 10^3$	$6.0 \cdot 10^3$	$6.2 \cdot 10^3$	$4.7 \cdot 10^3$

Table 5.2: Common items between solutions ( $\ell = 10, d = 0.6$ )

$k$	40	80	120	160	200
$ \text{GREEDY} \cap \text{UNIFFDV} $	26	76	114	155	200
$ \text{GREEDY} \cap \text{GREEDYFDV} $	40	78	117	160	200
$ \text{UNIFFDV} \cap \text{GREEDYFDV} $	26	75	113	155	200

Table 5.3: Support size of the continuous relaxation  $\|z^*\|_0$  ( $\ell = 10, d = 0.6$ )

$k$	40	80	120	160	200
$d = 0.3$	$93 \pm 3$	$117 \pm 3$	$148 \pm 2$	$181 \pm 3$	$213 \pm 2$
$d = 0.6$	$92 \pm 7$	$117 \pm 4$	$145 \pm 4$	$180 \pm 3$	$214 \pm 4$
$d = 0.9$	$88 \pm 3$	$116 \pm 3$	$147 \pm 4$	$179 \pm 3$	$214 \pm 1$

### 5.3.2 Real data

We used the Concrete Compressive Strength dataset (Yeh, 1998) (with column normalization) from the UCI repository to evaluate ESP-design on real data; this dataset consists in 1030 possible experiments to model concrete compressive strength as a linear combination of 8 physical parameters. In Figure 5-2 (a), OED chose  $k$  experiments to run to estimate  $\theta$ , and we report the normalized prediction error on the remaining  $N - k$  experiments. The best choice of OED for this problem is of course A-optimal design, which shows the smallest predictive error. In Figure 5-2 (b), we report the fraction of non-zero entries in the design matrix  $\mathbf{X}_{S,:}$ ; higher values of  $\ell$  correspond to increasing sparsity.

This confirms that OED allows us to scale between the extremes of A-optimal design and D-optimal design to tune desirable side-effects of the design; for example, sparsity in a design matrix can indicate not needing to tune a potentially expensive experimental parameter, which is instead left at its default value.

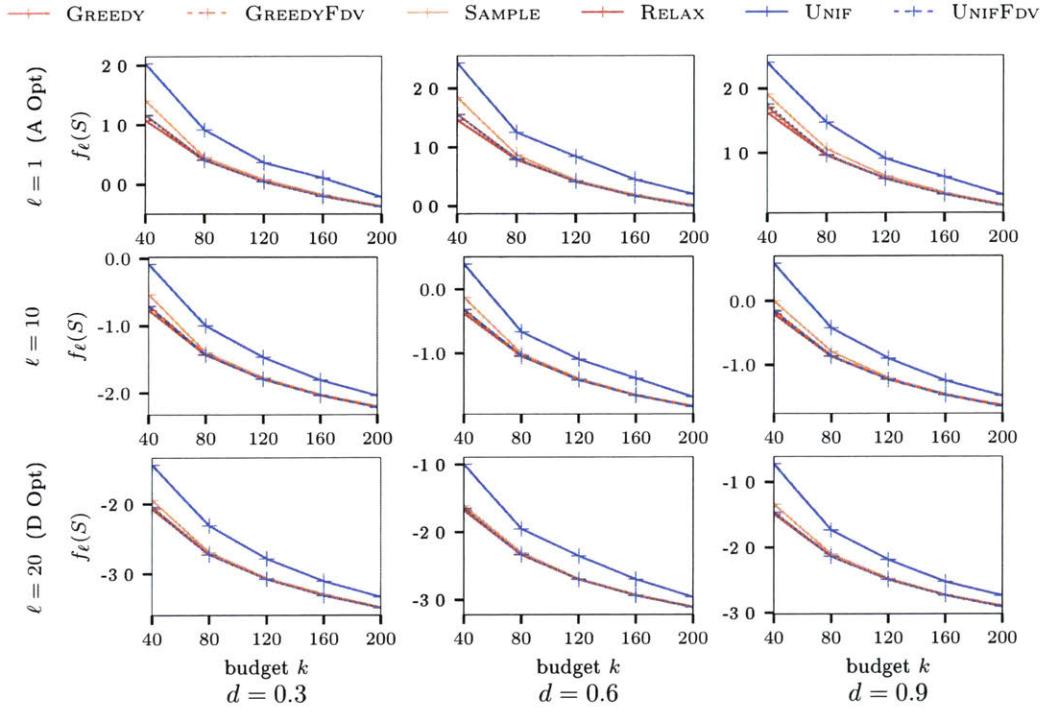


Figure 5-1: Synthetic experiments,  $N = 500$ ,  $m = 30$ . The greedy algorithm performs as well as the classical Fedorov approach; as  $k$  increases, all designs except UNIF converge towards the continuous relaxation, making SAMPLE the best approach for large designs.

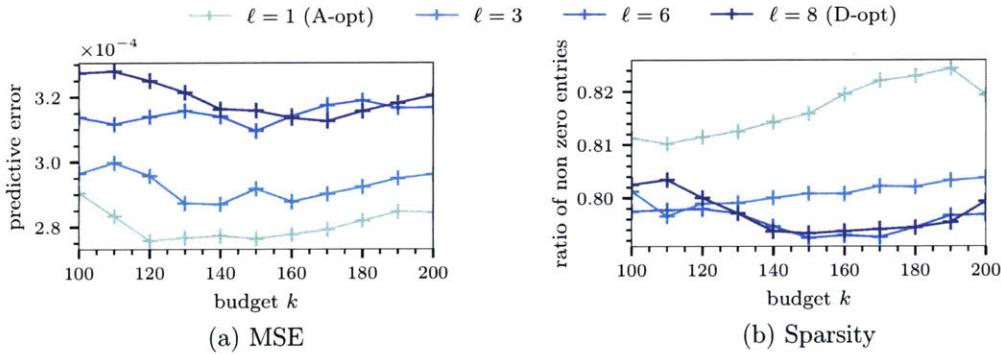


Figure 5-2: Predicting concrete compressive strength via the greedy method; higher  $\ell$  increases the sparsity of the design matrix  $\mathbf{X}_{S,:}$ , at the cost of marginally decreasing predictive performance.

## 5.4 Generalizing the MVCE problem

We close our discussion of Gvs measures by discussing an important geometric problem related to ESP-design. Our motivation here is the dual to the continuous re-

laxation of D-optimal design: this is nothing but the well-known *Minimum Volume Covering Ellipsoid (MVCE)* problem, which is a problem of great interest to the optimization community (Todd, 2016).

Due to the importance of MVCE in optimization and statistics, investigating the dual form of the more general continuous ESP-design problem is a natural question which we investigate below.

#### 5.4.1 GVS and Minimum Volume Covering Ellipsoids

Indeed, the dual of the natural convex relaxation to ESP-design generalizes the well-studied Minimum Volume Covering Ellipsoid (MVCE) optimization problem (Barnes, 1982; Rousseeuw and Leroy, 1987; Vandenberghe et al., 1998; Sun and Freund, 2004; Dolia et al., 2006; Todd, 2016), which we recover when  $e_\ell \equiv \det$ .

An ellipsoid  $\mathcal{E}$  in  $\mathbb{R}^m$  is defined by a matrix  $\mathbf{Q} \succeq 0$  and a center  $c$  such that

$$\mathcal{E} = \{x \in \mathbb{R}^m \mid (x - c)^\top \mathbf{Q}(x - c) \leq 1\}.$$

The volume of an ellipsoid  $\mathcal{E}$  is proportional<sup>1</sup> to  $\det \mathbf{Q}^{-1/2}$ . Hence, given  $N$  vectors  $a_1, \dots, a_N$ , the Minimum Covering Ellipsoid (MVCE) problem is stated as follows: find  $\mathbf{Q}_*, c_*$  such that

$$\mathbf{Q}_*, c_* \in \operatorname{argmin} \det \mathbf{Q}^{-1/2} \quad \text{with } \mathbf{Q} \succeq 0 \text{ and } \forall i, (a_i - c)^\top \mathbf{Q}(a_i - c) \leq 1.$$

In order to formulate MVCE as a convex problem, a standard transformation sets  $\mathbf{M} = \mathbf{Q}^{1/2}$  and  $z = \mathbf{M}c$ , yielding the standard formulation of the MVCE problem:

$$(\mathbf{M}_*, z_*) \in \operatorname{argmax} \log \det \mathbf{M} \quad \text{s.t.} \quad \mathbf{M} \succeq 0 \text{ and } \forall i, (\mathbf{M}a_i - z)^\top (\mathbf{M}a_i - z) \leq 1. \tag{5.13}$$

It is well known that (5.13) is equivalent to the dual problem of continuous D-optimization relaxation, *i.e.*, the mode finding problem of a GVS measure with  $\ell = m$ .

---

<sup>1</sup>The proportionality factor is given by the volume of the unit ball in  $\mathbb{R}^m$ .

MVCE has numerous applications in outlier detection and robust statistics in general (Rousseeuw and Leroy, 1987; Rousseeuw and van Zomeren, 1990; Mount et al., 2000), image analysis (Comaniciu and Meer, 1997), clustering (Jolion et al., 1991), and robotics (Li et al., 2016a) among others. We believe that the dual problems presented herein may find similar applications.

### 5.4.2 Deriving the dual

With MVCE as our motivation, we develop the dual formulation of the continuous relation of ESP-design. We start by deriving  $\nabla e_\ell(\mathbf{A})$ , for which we recall that  $e_\ell(\cdot)$  is a spectral function, whereby the spectral calculus of Lewis (1996) becomes applicable, saving us from intractable multilinear algebra (Jain, 2011). More precisely, say  $U^\top \boldsymbol{\Lambda} U$  is the eigendecomposition of  $\mathbf{A}$ , with  $U$  unitary. Then, as  $e_\ell(\mathbf{A}) = e_\ell \circ \lambda(\mathbf{A})$ , we have

$$\nabla e_\ell(\mathbf{A}) = U^\top \text{Diag}(\nabla e_\ell(\boldsymbol{\Lambda}))U = U^\top \text{Diag}(e_{\ell-1}(\boldsymbol{\Lambda}^{-(i)}))U.$$

We can now derive the dual of (5.8) as

$$\sup_{\mathbf{A} \succ 0, z \geq 0} \inf_{\mu \in \mathbb{R}, \mathbf{H}} -\frac{1}{\ell} \log e_\ell(\mathbf{A}) - \text{tr}(\mathbf{H}(\mathbf{A}^{-1} - \mathbf{X}^\top \text{Diag}(z)\mathbf{X})) - \mu(\mathbf{1}^\top z - k),$$

which admits as dual

$$\inf_{\mu \in \mathbb{R}, \mathbf{H}} \sup_{\mathbf{A} \succ 0, z \geq 0} \underbrace{-\frac{1}{\ell} \log e_\ell(\mathbf{A}) - \text{tr}(\mathbf{H}\mathbf{A}^{-1})}_{g(\mathbf{A})} + \text{tr}(\mathbf{H}\mathbf{X}^\top \text{Diag}(z)\mathbf{X}) - \mu(\mathbf{1}^\top z - k). \quad (5.14)$$

It is easy to see that we must have  $\mathbf{H} \succeq 0$ : indeed, by contradiction, assume that there exists  $x$  such that  $x^\top \mathbf{H} x < 0$  and  $\|x\| = 1$ . Then setting  $\mathbf{A} = \mathbf{I} - \frac{t}{1+t}xx^\top$  has  $g(\mathbf{A})$  go to infinity with  $t$ .

Before further analyzing (5.14), we first analyze  $g(\mathbf{A})$ .

**Proposition 5.4.1.** *Let function  $g : \mathbb{S}^{++} \rightarrow \mathbb{R}$  be defined by (5.14). Then,  $g$  attains its maximum on  $\mathbb{S}^{++}$ .*

*Proof.* This result follows from the fact that  $g \rightarrow -\infty$  for  $\|\mathbf{A}\| \rightarrow \infty$ ,  $\mathbf{A} \rightarrow \partial \mathbb{S}^{++}$ .  $\square$

Rewriting  $\mathbf{A} = \mathbf{U}^\top \boldsymbol{\Lambda} \mathbf{U}$ , we have

$$\nabla g(\mathbf{A}) = 0 \iff \boldsymbol{\Lambda} \operatorname{Diag}(e_{\ell-1}(\boldsymbol{\Lambda}_{(i)})) \boldsymbol{\Lambda} = e_\ell(\boldsymbol{\Lambda}) \mathbf{U} \mathbf{H} \mathbf{U}^\top.$$

Let  $a(\mathbf{H})$  be a matrix such that  $\nabla g(a(\mathbf{H})) = 0$ . Since  $f_\ell$  is convex,  $g(a(\mathbf{H})) = f_\ell^*(-\mathbf{H})$  where  $f_\ell^*$  is the Fenchel conjugate of  $f_\ell$ .

Furthermore, the previous equation shows that  $\mathbf{H}$  and  $a(\mathbf{H})$  are co-diagonalizable, with  $\boldsymbol{\Lambda} \operatorname{Diag}(e_{\ell-1}(\boldsymbol{\Lambda}_{(i)})) \boldsymbol{\Lambda} = \operatorname{Diag}(h_1, \dots, h_m)$ . The eigenvalues of  $a(\mathbf{H})$  must thus satisfy the system of equations

$$\lambda_i^2 e_{\ell-1}(\lambda_1, \dots, \lambda_{i-1}, \lambda_{i+1}, \dots, \lambda_m) = h_i e_\ell(\lambda_1, \dots, \lambda_m), \quad 1 \leq i \leq m.$$

Then,

$$\begin{aligned} e_\ell(\lambda_1, \dots, \lambda_m) \operatorname{tr}(\mathbf{H} a(\mathbf{H})^{-1}) &= \operatorname{tr}\left(\boldsymbol{\Lambda} \operatorname{Diag}[e_{\ell-1}(\boldsymbol{\Lambda}^{-(i)})] \boldsymbol{\Lambda} \boldsymbol{\Lambda}^{-1}\right) \\ &= \sum_i \lambda_i e_{\ell-1}(\lambda_1, \dots, \lambda_{i-1}, \lambda_{i+1}, \dots, \lambda_m) \\ &= \sum_i \sum_{J \subseteq [n], |J|=\ell, i \in J} \prod_{j \in J} \lambda_j. \end{aligned}$$

Each subset  $J$  of size  $\ell$  appears  $\ell$  times (once for each of its elements); finally

$$\operatorname{tr}(\mathbf{H} (a(\mathbf{H}))^{-1}) = \ell e_\ell(\lambda_1, \dots, \lambda_m)^{-1} \sum_{J \subseteq [n], |J|=\ell} \prod_{j \in J} \lambda_j = \ell,$$

and hence

$$g(a(\mathbf{H})) = f_\ell^*(-\mathbf{H}) = -\frac{1}{\ell} \log e_\ell(a(\mathbf{H})) - \ell.$$

Returning to the dual formulation (5.14), we can rewrite it as

$$\begin{aligned} &\inf_{\substack{\mu \in \mathbb{R}, \\ \mathbf{H} \in \mathbb{R}^{m \times m}}} \sup_{\mathbf{A} \succ 0, z \geq 0} -\frac{1}{\ell} \log e_\ell(\mathbf{A}) - \operatorname{tr}(\mathbf{H} \mathbf{A}^{-1}) + \operatorname{tr}(\mathbf{H} \mathbf{X}^\top \operatorname{Diag}(z) \mathbf{X}) - \mu(\mathbf{1}^\top z - k) \\ &= \inf_{\mu \in \mathbb{R}, \mathbf{H} \succeq 0} \left[ f_\ell^*(-\mathbf{H}) + \sup_{z \geq 0} \operatorname{tr}(\mathbf{H} \mathbf{X}^\top \operatorname{Diag}(z) \mathbf{X}) - \mu(\mathbf{1}^\top z - k) \right] \end{aligned}$$

$$\begin{aligned}
&= \inf_{\mu \in \mathbb{R}, \mathbf{H} \succeq 0} \left[ f_\ell^*(-\mathbf{H}) + \sup_{z \geq 0} \sum_i z_i (x_i^\top \mathbf{H} x_i - \mu) + \mu k \right] \\
&= \inf_{\substack{x_i^\top \mathbf{H} x_i \leq \mu, \\ \mathbf{H} \succeq 0}} f_\ell^*(-\mathbf{H}) + k\mu \\
&= \sup_{\substack{x_i^\top \mathbf{H} x_i \leq 1, \\ \mathbf{H} \succeq 0, \mu > 0}} -f_\ell^*(-\mu \mathbf{H}) - k\mu \stackrel{*}{=} \sup_{\substack{x_i^\top \mathbf{H} x_i \leq 1, \\ \mathbf{H} \succeq 0}} -f_\ell^*(-\mathbf{H}) \\
&= \frac{1}{\ell} \sup_{\substack{x_i^\top \mathbf{H} x_i \leq 1, \\ \mathbf{H} \succeq 0}} \log e_\ell(a(\mathbf{H})) + \ell,
\end{aligned}$$

where the eigenvalues  $\Lambda$  of  $a(\mathbf{H})$  verify

$$\lambda_i^2 \frac{e_{\ell-1}(\lambda_1, \dots, \lambda_{i-1}, \lambda_{i+1}, \dots, \lambda_m)}{e_\ell(\lambda_1, \dots, \lambda_m)} = h_i, \quad 1 \leq i \leq m. \quad (5.15)$$

and  $\stackrel{*}{=}$  follows from  $f_\ell^*(-\mu \mathbf{H}) = \sup_{\mathbf{A} \succ 0} -\text{tr}(\mathbf{H}\mathbf{A}) - f_\ell(\mathbf{A}/\mu) = f_\ell^*(-\mathbf{H}) - \log \mu$ .

In the general case, deriving  $a(\mathbf{H})$  or even  $e_\ell(a(\mathbf{H}))$  does not admit a closed form that we know of. However, we can gain some insight in the behavior of the dual by analyzing the behavior of the system of equations (5.15). Indeed, rewriting

$$e_\ell(\lambda_1, \dots, \lambda_m) = \lambda_i e_{\ell-1}(\lambda_1, \dots, \lambda_{i-1}, \lambda_{i+1}, \dots, \lambda_m) + e_\ell(\lambda_1, \dots, \lambda_{i-1}, 0, \lambda_{i+1}, \dots, \lambda_m),$$

we obtain the following formulations for (5.15):

$$\lambda_i = h_i + \lambda_i \frac{e_\ell(\lambda_1, \dots, 0, \dots, \lambda_m)}{e_\ell(\lambda_1, \dots, \lambda_m)} = h_i + g_i(h). \quad (5.16)$$

**Corollary 5.4.1.1.** *Let  $g_i$  be defined by (5.16). Then for  $\ell = m$ , we have  $g_i \equiv 0$  and so  $a(\mathbf{H}) = \mathbf{H}$ : we recover the MVCE problem.*

**Corollary 5.4.1.2.** *Let  $g_i$  be defined by (5.16). Then, for  $\ell = 1$ , we have  $\lambda_i^2 / \text{tr}(\Lambda) = h_i$  and so  $a(\mathbf{H}) = \text{tr}(\mathbf{H}^{1/2}) \mathbf{H}^{1/2}$ .*

In general, we obtain an alternate form of the dual:

$$\sup_{\substack{x_i^\top \mathbf{H} x_i \leq 1, \\ \mathbf{H} \succeq 0}} \log e_\ell(h_i + g_i(h)) \quad \text{such that the } \lambda_i \text{ satisfy (5.15),}$$

where  $g_i(h)$  is inversely tied to the impact a change away from 0 in the  $i$ -th coordinate has on  $e_\ell$ ; in the extreme case where  $\ell = m$ , only non-zero values of  $\lambda_i$  yield non-zero values of  $e_\ell$ , and so the additional term  $g_i$  vanishes entirely from the dual formulation.

## 5.5 Open problems

- In deriving the Lagrangian dual of the optimization problem, we had to introduce the function  $a(H)$ ; although  $a(H)$  is known for  $\ell = 1, m$ , its form for other values of  $\ell$  is unknown, making the dual form a purely theoretical object in the general case. Whether the closed form of  $a$  can be derived, or whether  $E_\ell(a(H))$  can be obtained with only knowledge of  $H$ , remains an open problem.
- More practically, we showed experimentally (Fig. 5-2) that changing the degree  $\ell$  of the elementary symmetric polynomial in the optimization problem (5.6) has straightforward consequences on the sparsity of the recovered designs; understanding this phenomenon theoretically, as well as investigating other properties that may scale similarly, are questions worthy of investigating.



# Chapter 6

## Exponentiated strongly Rayleigh measures

The previous chapters showed that strongly Rayleigh measures offer a tractable and theoretically grounded framework for models of diversity in machine learning. However, real-world problems cannot rely purely on negative dependence: consider once more a user querying the term “bow”. We expect some amount of negative dependence, since a corresponding set of answers should allow for the diversity of interpretations of the term (archery, string instruments, etc.). However, answers that are irrelevant to the query — for example, bikes — reduce the quality of the answer set. In practice, we want the ability to temper the strength of the negative dependence.

Indeed, a practical challenge to using SR measures is that SR measures do not allow for intuitive tuning of the strength for the repulsive forces between similar items. This motivates us to consider a generalization of SR measures that enables an easy tuning of the relative strengths of quality and diversity considerations.

The closing contribution of this thesis is the introduction of *exponentiated strongly Rayleigh (ESR)* measures: distributions of the form  $\nu(S) \propto \mu(S)^p$ , where  $\mu$  is an SR measure and  $p > 0$  is a tunable parameter. A power  $p > 1$  captures a sharper notion of diversity than  $\mu$ ; conversely,  $p < 1$  allows for weaker diversity preferences. At the  $p = 0$  extreme,  $\nu$  is uniform, while for  $p \rightarrow \infty$ ,  $\nu$  concentrates at the mode of  $\mu$ .

By introducing such a tunable parameter  $p$  that controls the strength of repulsive

forces between similar elements, we allow exponentiated strongly Rayleigh measures to capture a notion of diversity compatible with the particular requirements of machine learning subset selection problems.

ESR measures present an attractive generalization to SR measures, where a single parameter allows an intuitive regulation of desired strength of negative dependence. Interestingly, a few special cases of ESRs have been briefly noted in the literature (Zou and Adams, 2012; Kulesza and Taskar, 2012; Gillenwater et al., 2014), although only the guise of generalizations to DPPs and without noting any connection to SR measures. To our knowledge, there has been no previous discussion of ESR measures as a class. Nonetheless, they can benefit from the abundant existing theory for log-submodular models (Djolonga and Krause, 2014; Gotovos et al., 2015; Rebeschini and Karbasi, 2015; Djolonga et al., 2016), and isolated special cases have also been discussed in the literature. In particular, exponentiated DPPs (or E-DPPs) are mentioned in (Zou and Adams, 2012; Kulesza and Taskar, 2012), as well as in (Gillenwater et al., 2014) and (Anari and Oveis Gharan, 2017).

After analyzing the negative association properties of ESR measures, we derive general-purpose sampling algorithms that we further specialize to concrete cases. Subsequently, we evaluate ESR sampling procedures on outlier detection and kernel reconstruction, confirming that several machine learning problems benefit from the modeling power of ESR measures. This chapter is based on (Mariet et al., 2018).

## 6.1 Definition and negative dependence properties

We begin by formally introducing exponentiated strongly Rayleigh measures and analyzing their properties within the framework of negative dependence.

**Definition 6.1.1** (Exponentiated SR measure). A measure  $\mu$  over  $2^{[N]}$  is exponentiated strongly Rayleigh (ESR) if there exists an SR measure  $\nu$  over  $2^{[N]}$  and a power  $p \geq 0$  such that  $\mu(S) \propto \nu(S)^p$ .

The parameter  $p$  serves to control the quality/diversity trade-off by sharpening ( $p > 1$ ) or smoothing out ( $p < 1$ ) the variations of the ground SR measure (see

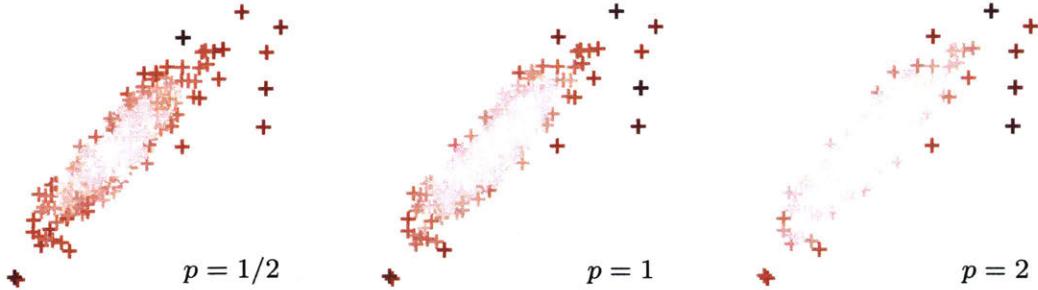


Figure 6-1: Anomaly detection by sampling with an exponentiated-DPP. 200 samples of size  $k = 20$  were drawn from a E-DPP with Gaussian kernel; darker colors indicate higher sampling frequencies. As  $p$  increases, the points furthest from the mean accumulate all of the sampling probability mass.

Figure 6-1). A natural question is then to understand how this additional parameter impacts the negatively associated properties of the ESR. Recall that a fundamental property of SR measures lies in the fact that they are negatively associated: for two increasing functions  $F, G$  over  $2^{[N]}$  that depend on a disjoint set of coordinates, an SR measure  $\mu$  verifies the following inequality (Borcea et al., 2009, Theorem 4.9):

$$\mathbb{E}_\mu[F] \mathbb{E}_\mu[G] \geq \mathbb{E}_\mu[FG]. \quad (6.1)$$

Our first result states that the additional modularity enabled by the exponent parameter can break strong Rayleighness; as a consequence, we have no immediate guarantee that ESRs verify (6.1).

**Proposition 6.1.1.** *There exist ESR measures that are not SR.*

*Proof.* Recall (Borcea et al., 2009, Thm. 4.1) that a real bivariate affine polynomial  $p(x, y)$  is stable if and only if

$$\partial_x p \partial_y p - p \partial_{xy} p \geq 0.$$

For an E-DPP with a kernel  $\mathbf{L} \in \mathbb{R}^{2 \times 2}$ , this is exactly  $l_{11}^p l_{22}^p \geq \det(\mathbf{L})^p$ , which is clearly true for all  $p \geq 0$  as  $\mathbf{L}$  must be positive semi-definite.

However, stability does not hold in general. To obtain a counterexample, consider

the generating polynomial  $p(x, y, z)$  for an E-DPP of dimension 3 (writing  $d = \det(\mathbf{L})$  as a shorthand):

$$\begin{aligned} p(x, y, z) := & d^p + l_{11}^p yz + l_{22}^p xz + l_{33}^p xy \\ & + \det(\mathbf{L}[1, 2])^p z + \det(\mathbf{L}[1, 3])^p y + \det(\mathbf{L}[2, 3])^p x + xyz. \end{aligned}$$

If  $p(x, y, z)$  is SR,  $p$  must be stable through conditioning (Prop. 2.2.5); hence,  $p(x, y, 1)$  must also be stable. Writing  $d_{ij} = \det(\mathbf{L}_{\{i,j\}})$ , this requires that

$$p(x, y, 1) = d^p + d_{12}^p + (d_{22}^p + d_{23}^p)x + (d_{11}^p + d_{13}^p)y + (d_{33}^p + 1)xy$$

be stable. Since  $p(x, y, 1)$  is a real bivariate affine polynomial, we must then have

$$(d_{22}^p + d_{23}^p)(d_{11}^p + d_{13}^p) \geq (d^p + d_{12}^p)(d_{33}^p + 1).$$

Finally, one can verify that this last inequality is easily violated for several choices of (non-diagonal) positive semi-definite matrices.  $\square$

Conversely, some ESR measures remain SR for any  $p$ : if  $\mu$  is a DPP parametrized by a block-diagonal kernel with  $2 \times 2$  blocks,  $\nu = \alpha\mu^p$  is also a DPP, and so SR and ESR by construction. The next theorem guarantees the existence of non-trivial ESR measures that are also SR.

**Theorem 6.1.2.** *There exists  $\epsilon > 0$  such that  $\forall p \in [1 - \epsilon, 1 + \epsilon]$ ,  $\forall N \in \mathbb{N}$ , there exists a matrix  $\mathbf{L} \in \mathbb{S}_N^+$ , that is not block-diagonal with blocks of size  $2 \times 2$ , such that the E-DPP distribution defined by  $\nu(S) \propto \det(\mathbf{L}[S, S])^p$  is SR.*

*Proof.* We write  $d_S^i = \det(\mathbf{L}[S \cup \{i\}])^p$  and when possible  $d_{ijk} = \det(\mathbf{L}[\{i, j, k\}])^p$ .

Let  $\mathbf{L} \succeq 0 \in \mathbb{R}^{n \times n}$ . The associated E-DPP is SR if and only if  $P_{ij}(z) \geq 0$  where  $P_{ij}(z)$  is defined for any  $1 \leq i \neq j \leq N$  as

$$P_r(z) = \frac{\partial f}{\partial z_i}(z) \frac{\partial f}{\partial z_j}(z) - f(z) \frac{\partial^2 f}{\partial z_i \partial z_j}(z)$$

$$= \left( \sum_{S \in \mathcal{Y}'} d_S^i z^S \right) \left( \sum_{S \in \mathcal{Y}'} d_S^j z^S \right) - \left( \sum_{S \in \mathcal{Y}'} d_S^{ij} z^S \right) \left( \sum_{S \in \mathcal{Y}'} d_S z^S \right),$$

where we write  $\mathcal{Y}' = \{S \in [N], i \notin S, j \notin S\}$  and  $z = (z_1, \dots, z_n)$  where components  $z_i$  and  $z_j$  are removed. Now, choose  $k \in [N] \setminus \{i, j\}$ , and write  $\tilde{z}$  for the vector  $z$  without component  $z_k$ . Write also  $\mathcal{Y}'_k = \{S \in \mathcal{Y}', k \notin S\}$ . Then,

$$\begin{aligned} P_{ij}(z) &= (\sum_{k \notin S} d_S^i \tilde{z}^S + z_k \sum_{k \in S} d_S^i \tilde{z}^S) (\sum_{k \notin S} d_S^j \tilde{z}^S + z_k \sum_{k \in S} d_S^j \tilde{z}^S) \\ &\quad - (\sum_{k \notin S} d_S^{ij} \tilde{z}^S + z_k \sum_{k \in S} d_S^{ij} \tilde{z}^S) (\sum_{k \notin S} d_S \tilde{z}^S + z_k \sum_{k \in S} d_S \tilde{z}^S) \\ &= \left( \sum_{S \in \mathcal{Y}'_k} d_S^{ik} \tilde{z}^S \sum_{S \in \mathcal{Y}'_k} d_S^{jk} \tilde{z}^S - \sum_{S \in \mathcal{Y}'_k} d_S^{ijk} \tilde{z}^S \sum_{S \in \mathcal{Y}'_k} d_S^k \tilde{z}^S \right) z_k^2 \\ &\quad + \left( \sum_{S \in \mathcal{Y}'_k} d_S^{ik} \tilde{z}^S \sum_{S \in \mathcal{Y}'_k} d_S^j \tilde{z}^S + \sum_{S \in \mathcal{Y}'_k} d_S^i \tilde{z}^S \sum_{S \in \mathcal{Y}'_k} d_S^{jk} \tilde{z}^S \right. \\ &\quad \left. - \sum_{S \in \mathcal{Y}'_k} d_S^{ijk} \tilde{z}^S \sum_{S \in \mathcal{Y}'_k} d_S \tilde{z}^S - \sum_{S \in \mathcal{Y}'_k} d_S^{ij} \tilde{z}^S \sum_{S \in \mathcal{Y}'_k} d_S^k \tilde{z}^S \right) z_k \\ &\quad + \left( \sum_{S \in \mathcal{Y}'_k} d_S^i \tilde{z}^S \sum_{S \in \mathcal{Y}'_k} d_S^j \tilde{z}^S - \sum_{S \in \mathcal{Y}'_k} d_S^{ij} \tilde{z}^S \sum_{S \in \mathcal{Y}'_k} d_S \tilde{z}^S \right) \\ &= Az_k^2 + Bz_k + C. \end{aligned}$$

Hence,  $\mathbf{L}$  yields a strongly Rayleigh E-DPP measure if and only if the following inequalities hold for all  $i, j, k$ :

$$(B/2)^2 \leq AC; \quad A \geq 0; \quad C \geq 0. \quad (6.2)$$

When  $n = 3$ , (6.2) reduces to the arithmetic-geometric inequality:

$$\left( \frac{d_i d_{jk} + d_j d_{ik} - d_{ij} d_k - d_{ijk}}{2} \right)^2 \leq (d_i d_j - d_{ij})(d_{jk} d_{ik} - d_k d_{ijk}). \quad (6.3)$$

One can easily obtain  $3 \times 3$  positive semi-definite matrices  $\mathbf{L}$  that verify (6.3) strictly for  $p = 1$ ; in particular, by continuity, there exists  $\epsilon > 0$  such that the E-DPP

generated by the kernel  $\mathbf{L}$  and power  $p \in [1 - \epsilon, 1 + \epsilon]$  will still verify (6.3).

Then, as a block-diagonal matrix such that each diagonal block yields SR E-DPPs also yields a SR E-DPP, we can thus generate block-diagonal matrices of any size  $n$  such that the blocks are either  $\mathbf{L}$  or  $2 \times 2$  matrices, which all yield SR E-DPPs for  $p \in [1 - \epsilon, 1 + \epsilon]$ .  $\square$

Hence, ESRs are not guaranteed to be SR but may remain so. Due to their log-submodularity, they nonetheless will verify the so-called *negative lattice condition*  $\mu(S \cap T)\mu(S \cup T) \leq \mu(S)\mu(T)$ , and so retain negative dependence properties.

We now show that ESRs nonetheless have a fundamental advantage over standard log-submodular functions: although the intractability of their partition function precludes exact sampling algorithms, their closed form as the exponentiation of an SR measure can be leveraged to take advantage of the recent result (Li et al., 2016b) on fast-mixing Markov chains for SR measures.

## 6.2 Sampling from ESR measures

In the general case, the normalization term of an ESR is NP-hard to compute, precluding exact sampling algorithms. In this section, we propose instead two MCMC sampling algorithms whose key idea lies in exploiting the explicit relation ESR measures have to SR measures.

We begin by introducing the notion of  $r$ -closeness, which serves as a measure of the proximity between distributions  $\mu$  and  $\nu$  over subsets. In practice,  $r$ -closeness will allow us to quantify how close an ESR measure is to being SR, and inform our bounds on mixing time.

**Definition 6.2.1** ( $r$ -closeness). Let  $\mu, \nu$  be measures over  $2^{[N]}$  and let  $p \geq 0$ . We say that  $\nu$  is  $r$ -close to  $\mu$  (where we allow  $r = \infty$ ) if we have for all  $S \subseteq [N]$ ,

$$\nu(S) \neq 0 \text{ and } \mu(S) \neq 0 \implies r^{-1} \leq \nu(S)/\mu(S) \leq r.$$

We additionally write  $r(\mu, \nu) = \min\{r \in \mathbb{R} \cup \{\infty\} : \nu \text{ is } r\text{-close to } \mu\}$ .

**Remark 6.2.1.**  $r$ -closeness is a symmetric property.

**Remark 6.2.2.** If  $r(\mu, \nu) < \infty$ ,  $\nu$  is absolutely continuous with respect to  $\mu$ : that is,  $\mu(S) = 0 \implies \nu^p(S) = 0$ .

The following result establishes that for any ESR measure  $\nu$ , there exists an SR measure  $\mu$  that is  $r$ -close to  $\nu$  with  $r < \infty$ . This result is the cornerstone of the sampling algorithms we derive, as we show that we can use an  $r$ -close SR measure as proposal to efficiently sample from an ESR measure.

**Proposition 6.2.1.** *Let  $\mu$  be a measure over  $2^{[N]}$ , and define  $\nu$  to be the measure such that  $\nu(S) \propto \mu(S)^p$ . Then*

$$r(\mu, \nu) \leq \max_{S \in \text{supp}(\nu)} [\mu(S)^{-|p-1|}] < \infty.$$

*Proof.* Let  $\mu, \nu$  be as in the proposition statement, and consider  $S \in \text{supp}(\nu)$ . Recall that  $\sum_T \nu(T) = 1$ . By definition,

$$\frac{\nu(S)}{\mu(S)} = \frac{\mu(S)^p}{\mu(S) \sum_T \mu(T)^p}.$$

If  $p \leq 1$ , we have  $\sum_T \mu(T)^p \geq 1$  and  $\mu(S)^p \geq \mu(S)$ , and so

$$(\min_S \mu(S))^{p-1} \leq \frac{\sum_T \mu(T)}{\sum_T \mu(T)^p} = \frac{1}{\sum_T \mu(T)^p} \leq \frac{\mu(S)^p}{\mu(S) \sum_T \mu(T)^p} \leq \mu(S)^{p-1} \leq \frac{1}{\min_S \mu(S)^{1-p}},$$

where the left inequality is obtained by noticing that  $\frac{a+b}{c+d} \geq \min\left(\frac{a}{c}, \frac{b}{d}\right)$ .

Similarly, for  $p \geq 1$ , we have  $\sum_T \mu(T)^p \leq 1$  and  $\mu(T)^p \leq \mu(T)$ , and so

$$\min_S \mu(S)^{p-1} \leq \frac{\mu(S)^p}{\mu(S)} \leq \frac{\mu(S)^p}{\mu(S) \sum_T \mu(T)^p} \leq \frac{\sum_T \mu(T)}{\sum_T \mu(T)^p} \leq \max \mu(S)^{1-p} = \frac{1}{\min \mu(S)^{p-1}}.$$

□

In order to sample from an ESR distribution  $\nu$ , we now generalize existing MCMC algorithms for SR measures; we bound the distance to stationarity of the chain's

current state by comparing it to the distance to stationarity of a similar chain sampling from an SR measure  $\mu$ , and leveraging the  $r$ -closeness  $r(\mu, \nu)$ .

### 6.2.1 Approximate samplers for ESR measures

Before investigating MCMC samplers, one may first wonder if rejection sampling might be sufficient: sample a set  $S$  from a proposal distribution  $\mu$ , and accept with probability  $\nu^p(S)/M\mu(S)$ , where  $M \geq \max_S \mu(S)/\nu^p(S)$ . Unfortunately, the rejection sampling scaling factor  $M$  cannot be computed — although it can be bounded by  $r(\mu, \nu^p)$  — leading us to prefer MCMC samplers (Andrieu et al., 2003).

We begin by analyzing the standard independent Metropolis-Hastings sampling algorithm (Metropolis et al., 1953; Hastings, 1970), using an SR measure  $\mu$  as a proposal: we sample an initial set  $S$  from  $\mu$  via a fast-mixing Markov chain, then iteratively swap from  $S$  to a new set  $S'$  with probability

$$\Pr(S \rightarrow S') = \min \left\{ 1, \frac{\nu(S')\mu(S)}{\nu(S)\mu(S')} \right\}.$$

---

**Algorithm 13** Proposal-based sampling

---

```

Input: SR proposal  $\mu$ , ESR measure  $\nu$  and SR measure  $\rho$  s.t.  $\nu = \alpha\rho^p$ 
Draw  $S \sim \mu$ 
while not mixed do
     $S' \sim \mu$ 
     $S \leftarrow S'$  w.p.  $\min \left\{ 1, \frac{\nu(S')\mu(S)}{\nu(S)\mu(S')} \right\} = \min \left\{ 1, \frac{\mu(S)}{\mu(S')} \left( \frac{\rho(S')}{\rho(S)} \right)^p \right\}$ 
return  $S$ 
```

---

Algorithm 13 relies on the fact that we can compute  $\nu(S')/\nu(S)$  as  $(\rho(S')/\rho(S))^p$ : we do not require knowledge of  $\nu$ 's partition function. This sampling method is valid as soon as  $\nu$  is absolutely continuous with regards to the proposal  $\mu$ ; Proposition 6.2.1 guaranteed the existence of such measures.

If the ESR measure  $\nu$  is  $k$ -homogeneous, we can instead sample from  $\nu$  via Algorithm 14: we randomly sample  $S \subseteq [N]$  and switch an element  $u \in S$  for  $v \notin S$  if this improves the probability of  $S$ .

The key to extending Algorithm 14 to non-homogeneous ESR measures is similar

---

**Algorithm 14** Swap-chain sampling

---

**Input:**  $k$ -homogeneous ESR measure  $\nu$  s.t.  $\nu = \alpha\rho^p$  with  $\rho$  SR  
 Sample  $S \sim \text{Unif}(N)$  such that  $|S| = k$   
**while** not mixed **do**  
 Sample  $u, v \in (S \times [N] \setminus S)$  uniformly at random  
 $S \leftarrow S \cup \{u\} \setminus \{v\}$  w.p.  $\min \left\{ 1, \frac{\nu(S \cup \{u\} \setminus \{v\})}{\nu(S)} \right\} = \min \left\{ 1, \left( \frac{\rho(S \cup \{u\} \setminus \{v\})}{\rho(S)} \right)^p \right\}$   
**return**  $S$

---

to the approach taken by Li et al. (2016b) for SR measures, and relies on leveraging the symmetric homogenization  $\nu_{\text{sh}}$  of  $\nu$  over  $2^{[2n]}$ , which we recall is defined by

$$\nu_{\text{sh}} : S \in 2^{[2n]} \rightarrow \begin{cases} \nu(S \cap [N]) \binom{n}{|S \cap [N]|}^{-1} & \text{if } |S| = n \\ 0 & \text{if } |S| \neq n. \end{cases}$$

If  $\nu \propto \mu^p$ ,  $\nu_{\text{sh}}$  is absolutely continuous with regards to  $\mu_{\text{sh}}$ . A simple calculation further shows that  $r(\mu_{\text{sh}}, \nu_{\text{sh}}) = r(\mu, \nu)$ , and so to sample  $S$  from  $\nu$ , it suffices to sample  $T$  of size  $n$  from  $\nu_{\text{sh}}$  using Algorithm 14, and then output  $S = T \cap [N]$ .

Hence, although we cannot in the general case sample from an ESR measure exactly (unlike many SR measures), being able to evaluate an ESR measure's unnormalized density function allows us to leverage MCMC algorithms for approximate sampling. We now focus on bounding the mixing times of these algorithms.

### 6.2.2 Mixing time analysis

Writing  $\nu'_{t,S}$  the distribution generated by a Markov chain sampler after  $t$  iterations and initialization set  $S$ , recall that the mixing time  $\tau_S(\epsilon)$  measures the number of required iterations of the Markov chain so that  $\nu'_{t,S}$  is close enough (in total variational distance) to the true ESR measure  $\nu$ :

$$\tau_S(\epsilon) \triangleq \min\{t : \|\nu'_{t,S} - \nu\|_{\text{TV}} \leq \epsilon\}.$$

The mixing time of a chain depends on how close the distribution generating the initialization set  $S$  is to the target distribution  $\mu$ . We now show this explicitly for the

two algorithms derived above, obtaining bounds on  $\tau_S$  that directly depend on the  $r$ -closeness of the target ESR measure  $\nu$  and an SR measure  $\mu$ .<sup>1</sup>

For Alg. 13, mixing time explicitly depends on the quality of the proposal.

**Theorem 6.2.2** (Algorithm 13 mixing time). *Let  $\mu, \nu$  be measures over  $2^{[N]}$  such that  $\mu$  is SR and  $\nu$  is ESR. Sampling from  $\nu$  via Algorithm 13 with  $\mu$  as a proposal distribution has a mixing time  $\tau(\epsilon)$  such that*

$$\tau_S(\epsilon) \leq 2r(\mu, \nu^p) \log \frac{1}{\epsilon}.$$

*Proof.* Alg. 13 has a state-independent proposal distribution  $\mu$ , and hence its mixing time is governed by a ratio of probabilities: Cai (2000) showed that, after  $t$  iterations,

$$\max_{S, T} d_{\text{TV}}(\nu_{(t)}(\cdot | S), \nu_{(t)}(\cdot | T)) = \left(1 - \frac{1}{\max_U \nu(U)/\mu(U)}\right)^t,$$

where  $\nu_{(t)}(U | S)$  is the probability of being in state  $U$  after  $t$  iterations when starting from set  $S$ , and  $d_{\text{TV}}$  is the total variation distance.

Hence, following (Cai, 2000, Cor.1), we obtain

$$\tau_S(\epsilon) \leq 2 \max_U \frac{\nu(U)}{\mu(U)} \log \frac{1}{\epsilon} \leq 2r(\mu, \nu^p) \log \frac{1}{\epsilon}.$$

□

For the swapchain algorithm (Alg. 14), we derive a bound on the mixing time by comparing to a result by (Anari et al., 2016) which shows fast sampling for SR distributions over subsets of a fixed size.

**Theorem 6.2.3** (Algorithm 14 mixing time). *Let  $\nu$  be a  $k$ -homogeneous ESR measure over  $2^{[N]}$ . The mixing time for Alg. 14 is bounded in expectation by*

$$\tau_S(\epsilon) \leq \inf_{\mu \in SR} 2Nk r(\mu, \nu)^2 \log \frac{1}{\epsilon \nu(S)}.$$

---

<sup>1</sup>In (Anari et al., 2018), the authors independently showed that  $k$ -homogeneous ESR measures with  $p < 1$  are *Strongly Log-Concave* (Anari et al., 2018, Theorem 1.7), and hence admit Markov Chains (Anari et al., 2018, Theorem 1.1) with  $\tau_S(\epsilon) = \mathcal{O}(d \log \frac{1}{\epsilon \mu(S)})$ .

*Proof.* This bound is based on a comparison method (Diaconis and Saloff-Coste, 1993), and relates the mixing time to the spectral gap. Let  $1 = \mu_1 \geq \mu_2 \geq \dots \geq -1$  be the eigenvalues of the state transition matrix of the chain. The *spectral gap* is  $\gamma = 1 - \max\{|\mu|; \mu \text{ is an eigenvalue and } \mu \neq 1\}$ .  $\gamma$  directly translates into a bound on the mixing time Diaconis and Stroock (1991):

$$\tau_S(\gamma) \leq \frac{1}{\gamma} \log \left( \frac{1}{\epsilon\nu(S)} \right).$$

The comparison method yields a bound on  $\gamma$  if we know a bound on  $\tilde{\gamma}$  for a related chain with stationary distribution  $\mu$ . Specifying (Diaconis and Saloff-Coste, 1993, Thm 2.1) to this case yields  $\gamma \geq \tilde{\gamma}\alpha_1/\alpha_2$ , where

$$\begin{aligned} \alpha_1 &= \min_S \frac{\mu(S)}{\nu(S)} \geq \frac{1}{r(p)} \\ \alpha_2 &= \max_{T,U} \frac{\mu(T)}{\nu(T)} \frac{\min\{1, \mu(U)/\mu(T)\}}{\min\{1, \nu(U)/\nu(T)\}} \leq \max_T \frac{\mu(T)}{\nu(T)} \leq r(\mu, \nu). \end{aligned}$$

Anari et al. (2016) show that  $\tilde{\gamma} \geq \frac{1}{2Nk}$ . Hence, we obtain

$$\tau_S(\gamma) \leq 2Nk \cdot r(\mu, \nu)^2 \cdot \log \frac{1}{\epsilon\nu(S)}. \quad \square$$

The above bound depends on the *closest* SR distribution to the target measure  $\nu$ . Combined with Prop. 6.2.1, Thm. 6.2.3 provides a simple upper bound to the mixing time of the swapchain algorithm.

**Corollary 6.2.3.1** (Non-homogeneous swapchain mixing time). *Let  $\nu$  be a non-homogeneous ESR measure over  $2^{[N]}$ . The mixing time for the generalized swapchain sampler to sample from  $\nu$  with initialization  $S \subseteq [2N]$  is bounded in expectation by*

$$\tau_S(\epsilon) \leq \inf_{\mu \in \text{SR}} 4N^2 r(\mu, \nu)^2 \log \frac{1}{\epsilon\nu_{\text{sh}}(S)}.$$

As a Markov chain's applicability strongly depends on its mixing time, finding an  $r$ -close SR distribution with small  $r$  is crucial to the applicability of ESR measures.

### 6.2.3 Some specific bounds

We now derive explicit mixing time bounds for ESR measures  $\nu$  generated by two popular classes of SR measures: DPPs, in their usual form as well as their  $k$ -homogeneous form ( $k$ -DPPs), and Dual Volume Sampling (DVS). As Theorem 6.2.2 and Theorem 6.2.3 provide mixing time bounds that depend explicitly on  $r(\mu, \nu)$ , this section focuses on upper bounding  $r(\mu, \nu)$ . To the extent of our knowledge, the results below are the first for either of these two classes of ESR distributions.

**Theorem 6.2.4** (E-DVS closeness bounds). *Let  $n \geq k \geq m$  and let  $\mathbf{X} \in \mathbb{R}^{m \times n}$  be a maximal-rank matrix. Let  $\mu$  be the DVS distribution over  $2^{[n]}$  for sets of size  $k$ :*

$$\mu : S \subseteq [N] \rightarrow \begin{cases} 0 & \text{if } |S| \neq k \\ \mu(S) \propto \det(\mathbf{X}[:, S]\mathbf{X}[:, S]^\top) & \text{if } |S| = k. \end{cases}$$

Let  $p > 0$  and  $\nu$  be the ESR measure induced by  $\mu$  and  $p$ ; let  $\text{MinVol}(\mathbf{X}, S)$  be the smallest non-zero minor of degree  $m$  of  $\mathbf{X}[:, S]$ . Then

$$r(\mu, \nu) \leq \binom{n-m}{k-m}^{|1-p|} \binom{k}{m}^{-|1-p|} \det(\mathbf{X}\mathbf{X}^\top)^{|1-p|} \text{MinVol}(\mathbf{X}, S)^{-2|1-p|}.$$

**Theorem 6.2.5** (E-DPP closeness bound). *Let  $\mu$  be the distribution induced by a DPP with kernel  $\mathbf{L} \succeq 0$  and  $\nu$  be the E-DPP such that  $\nu(S) \propto \det(\mathbf{L}[S])^p$ . Let  $\lambda_1 \leq \dots \leq \lambda_N$  be the ordered eigenvalues of  $\mathbf{L}$ . Then,*

$$r(\mu, \nu^p) \leq \prod_{i=1}^N (1 + \lambda_i)^{|1-p|} \prod_{\lambda_i < 1} \lambda_i^{-|1-p|}.$$

**Theorem 6.2.6** (E- $k$ -DPP closeness bound). *Let  $\mu$  be the distribution over  $2^{[N]}$  induced by a  $k$ -DPP ( $k \leq N$ ) with kernel  $\mathbf{L}$ , and let  $\nu$  be the induced ESR measure with power  $p > 0$ . Then*

$$r(\mu, \nu) \leq e_k(\lambda_1, \dots, \lambda_N)^{|1-p|} \prod_{i=1}^k \lambda_i^{-|1-p|}.$$

where  $e_k$  the  $k$ -th elementary symmetric polynomial.

One easily shows that the values  $r(\mu, \nu)$  we derive above for ( $k$ -) DPPs are loosely lower-bounded by  $\kappa^{|1-p|}$ , where  $\kappa$  is the condition number of the kernel matrix  $\mathbf{L}$ . However, it is possible to obtain a closer SR distribution to  $\nu \propto \det(\mathbf{L})^p$  than the baseline choice of the DPP with kernel  $\mathbf{L}$ : indeed, as  $\mathbf{L}$  is positive semi-definite, we can also consider a DPP parametrized by kernel  $\mathbf{L}^p$ .

For the rest of this section, we define  $\mu$  as the SR measure corresponding to the DPP with kernel  $\mathbf{L}^p$ :

$$\mu(S) = \det(\mathbf{L}^p[S]) / \det(\mathbf{I} + \mathbf{L}^p),$$

and  $\nu$  as the ESR measure such that  $\mu(S) \propto \det(\mathbf{L}[S])^p$ , which remains absolutely continuous with regard to  $\mu$ .

In this setting, upper bounding  $r(\mu, \nu)$  proves to be significantly more difficult, and is the focus of the remainder of this section. We first recall a useful expansion of the determinant of principal submatrices, fundamental to deriving the bounds below and potentially of more general interest.

**Lemma 6.2.7** (Shirai and Takahashi (2003, Lemma 2.9)). *Let  $\mathbf{H}$  be an  $N \times N$  Hermitian matrix with eigenvalues  $\lambda_1, \dots, \lambda_N$ . There exists a  $2^N \times 2^N$  symmetric doubly stochastic matrix  $\mathbf{Q} = [\mathbf{Q}_{SJ}]$  indexed by subsets  $S, J$  of  $[N]$  such that*

$$\det(\mathbf{H}[S]) = \sum_{J \subseteq [N], |J|=|S|} \mathbf{Q}_{SJ} \prod_{i \in J} \lambda_i.$$

$\mathbf{Q}$  can be chosen to depend only on the eigenvectors of  $\mathbf{H}$  and to satisfy  $\mathbf{Q}_{S \neq J} = 0$ .

The above lemma allows us to bound the ratio  $\det(\mathbf{L}^p[S])/\det(\mathbf{L}[S]^p)$  in terms of the *generalized condition number* of  $\mathbf{L}$ .

**Definition 6.2.2** (Generalized condition number). Given a matrix  $\mathbf{L} \in \mathbb{S}_N^{++}$  with eigenvalues  $\lambda_1, \dots, \lambda_N$ , we define its generalized condition number of order  $k$  as

$$\kappa_k = (\lambda_1 \cdots \lambda_k)(\lambda_n \cdots \lambda_{N-k})^{-1}.$$

Note that  $\kappa_k$  is the usual condition number of the  $k$ -th exterior power  $\mathbf{L}^{\wedge k}$  (and so, in particular,  $\kappa_k \leq \kappa^k$ ).

Given the generalized conditioned number, Lemma 6.2.7 combined with the power-mean inequality (Specht, 1960) (Appendix A.1.2) suffices to bound the gap between volumes generated by E-DPPs and DPPs:

**Theorem 6.2.8.** *Let  $\mu$  be the distribution induced by a DPP with kernel  $\mathbf{L}^p$ , and  $\nu$  be the corresponding E-DPP such that  $\nu \propto \det(\mathbf{L}[S])^p$ . Then  $r(\mu, \nu) \leq r(\kappa_{\lfloor n/2 \rfloor}, p)$  where  $r(\kappa, p)$  is defined by*

$$r(\kappa, p) = \begin{cases} \left( \frac{p(\kappa-1)}{\kappa^{p-1}} \right)^p \left( \frac{(1-p)(\kappa-1)}{\kappa - \kappa^p} \right)^{1-p} & \text{for } 0 < p < 1 \\ \left( \frac{\kappa^p - 1}{p(\kappa-1)} \right)^p \left( \frac{(p-1)(\kappa-1)}{\kappa^p - \kappa} \right)^{p-1} & \text{for } p > 1. \end{cases}$$

*Proof.* We show the result for general DPPs; the result for  $k$ -DPPs follows the same exact reasoning.

For  $0 < p < 1$ , it follows from Thm. A.1.2 that  $\det(\mathbf{L}[S])^p \geq \det(\mathbf{L}^p[S])$ . Hence,

$$Z_p = \sum_{S \subseteq [N]} \det(\mathbf{L}[S])^p \geq \sum_{S \subseteq [N]} \det(\mathbf{L}^p[S]) = \det(\mathbf{I} + \mathbf{L}^p),$$

which entails  $\frac{\det(\mathbf{I} + \mathbf{L}^p)}{Z_p} \leq 1$  whereby it remains to bound  $\frac{\det(\mathbf{L}[S])^p}{\det(\mathbf{L}^p[S])}$ .

Let  $S \subseteq [N]$  of size  $k$ , and let  $\lambda$  be the vector of  $\mathbf{L}$ 's eigenvalues. We write  $\lambda^S = \prod_{i \in S} \lambda_i$ , and denote by  $\boldsymbol{\lambda}^{\wedge k}$  the  $\binom{N}{k}$ -vector  $(\lambda^S)_{S \subseteq [N], |S|=k}$ . Using Lemma 6.2.7, there exists  $\mathbf{w} \in \mathbb{R}^{\binom{N}{k}}$  that sums to 1 such that

$$\begin{aligned} \frac{\det(\mathbf{L}[S])^p}{\det(\mathbf{L}^p[S])} &= \frac{(\sum_{|S|=k} w_S \lambda^S)^p}{\sum_{|S|=k} w_S (\lambda^p)^S} = \left( \frac{M_1(\mathbf{w}; \boldsymbol{\lambda}^{\wedge k})}{M_p(\mathbf{w}; \boldsymbol{\lambda}^{\wedge k})} \right)^p \\ &\leq r(p)^p. \end{aligned}$$

Where the last inequality follows from Thm. A.1.1. To lower bound  $r_p(S)$ , the same reasoning gives us  $\frac{\det(\mathbf{L}[S]^p)}{\det(\mathbf{L}^p[S])} \geq 1$  and hence

$$\begin{aligned} r_p(S) &\geq \frac{\det(\mathbf{I} + \mathbf{L}^p)}{Z_p} \geq \min_S \frac{\det \mathbf{L}^p[S]}{(\det \mathbf{L}[S])^p} \\ &\geq \left( \frac{M_p(\mathbf{w}'; \boldsymbol{\lambda}^{\wedge k})}{M_1(\mathbf{w}'; \boldsymbol{\lambda}^{\wedge k})} \right)^p \geq r(p)^{-p}, \end{aligned}$$

where the second inequality follows  $\frac{a+b}{c+d} \geq \min\left(\frac{a}{c}, \frac{b}{d}\right)$ . The same reasoning yields the desired result for  $p > 1$ .

Finally, some algebra shows that for fixed  $p$ ,  $r$  is an increasing function of  $\kappa$ , and so  $r(\kappa_k, p)$  is upper bounded by  $r(\kappa_{\lfloor N/2 \rfloor})$ .  $\square$

**Corollary 6.2.8.1.** *Let  $\mu$  be the distribution induced by a  $k$ -DPP with kernel  $\mathbf{L}^p$ , and  $\nu$  be the ESR measure such that  $\nu(S) \propto \det(\mathbf{L}[S])^p$ . Then  $r(\mu, \nu) \leq r(\kappa_k, p)$ .*

As shown in Figure A-4 (Appendix A.2.4), the upper bound 6.2.8 grows slower than  $\kappa$ : this shows that  $\mu(S) \propto \det(\mathbf{L}^p[S])$  is a closer SR distribution to an E-DPP with kernel  $\mathbf{L}$  than the usual DPP with kernel  $\mathbf{L}$ , yielding finer mixing time bounds.

Note that the per-iteration complexity of both algorithms must also be taken into account when choosing a sampling procedure: for E-DPPs, despite Algorithm 13's smaller mixing time, Algorithm 14 is more efficient in cases when  $N$  large due to the comparative costs of each sampling round.

## 6.3 Experiments

To evaluate the empirical applications of ESR measures, we evaluate E-DPPs on a variety of machine learning task. In all cases where we use the proposal MCMC sampler (Alg. 13), we use the DPP with kernel  $\mathbf{L}^p$  as a proposal.

### 6.3.1 Evaluating mixing time

We begin by empirically evaluating the mixing time of both algorithms. We measure mixing using the classical Potential Scale Reduction Factor (PSRF) metric (Brooks and Gelman, 1998). As the PSRF converges to 1, the chain mixes. In the following experiments, we report the mixing time (number of iterations) necessary to reach a PSRF of 1.05, as well as the runtime (in seconds) to convergence, averaged over 5 iterations; we use matrices with a fixed  $\kappa^k$  across all mixing time experiments.

The mixing time for proposal-based sampling is an order of magnitude smaller than swap-chain sampling; this is in line with the bounds we provide in Theorems 6.2.2

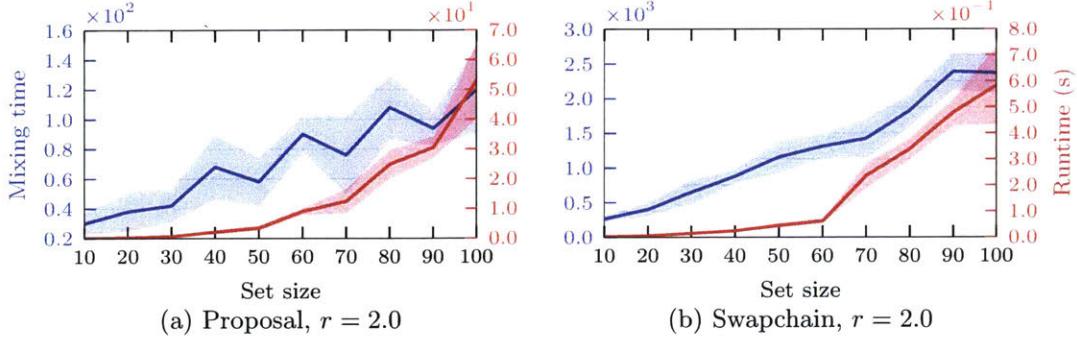


Figure 6-2: Mixing and sampling time for E- $k$ -DPPs as a function of the set size  $k$ . In both cases, the mixing time grows linearly with  $k$ ; although the mixing time for the proposal algorithm is an order of magnitude smaller than for the swapchain algorithm, the latter samples faster due to the per-iteration cost of each transition step.

and 6.2.3. However, this does not translate into faster runtimes: indeed, the per-iteration complexity of proposal-based sampling is significantly higher than for the swapchain algorithm, as Alg. 13 samples from a DPP at each iteration. The evolution of mixing and wall clock times as a function of  $N$  is provided in Appendix A.2.4.

### 6.3.2 Anomaly detection

We now focus on applications for E-DPPs; we begin by evaluating the use of E-DPPs for outlier detection. As increasing  $p$  heightens the model’s sensitivity to diversity, we expect  $p > 1$  to provide better outlier detection. To our knowledge, this is the first application of DPPs to outlier detection, and so our goal for this experiment is not to improve upon state-of-the-art results, but to compare the performance of (E-)DPPs for various values of  $p$  to standard outlier detection algorithms.

Experimentally, we detect an outlier via the following approach: given a dataset of  $n$  points and an E-DPP with an RBF kernel built from the data (bandwidth  $\beta = 100$ ), we sample  $\frac{n}{5}$  subsets of size 50 and report as outliers points that appear at least  $n\varepsilon$  times, where  $\varepsilon$  is a tunable parameter (hence, if we were doing uniform sampling, each point in the dataset would be sampled on average 10 times).

We detect outliers on three public datasets: the UCI Breast Cancer Wisconsin dataset (Street et al., 1993) modified as in (Goldstein and Uchida, 2016; Kriegel et al., 2009) as well as the Letter and Speech datasets from (Micenková et al., 2014).

We also report the performance of a selection of standard outlier detection algorithms whose reported performance in (Goldstein and Uchida, 2016) is competitive with other outlier detection algorithms: Local Outlier Factor (LOF) (Breunig et al., 2000),  $k$ -Nearest Neighbor ( $k$ -NN) (Angiulli and Pizzuti, 2002), Histogram-based Outlier Score (HBOS) (Goldstein and Dengel, 2012), Local Outlier Probability (LoOP) (Kriegel et al., 2009) and unweighted Cluster-Based Local Outlier Factor (uCBLOF) (Amer and Goldstein, 2012; Goldstein and Uchida, 2016).

Method	Cancer	Letter	Speech
E-DPP ( $p = 0.5$ )	$0.952 \pm 0.018$	$0.780 \pm 0.013$	$0.455 \pm 0.007$
DPP ( $p = 1$ )	$0.962 \pm 0.004$	$0.820 \pm 0.003$	$0.439 \pm 0.011$
E-DPP ( $p = 2$ )	$0.965 \pm 0.001$	$0.847 \pm 0.002$	$0.445 \pm 0.002$
LOF*	$0.982 \pm 0.002$	$0.867 \pm 0.027$	$0.504 \pm 0.022$
$k$ -NN*	$0.979 \pm 0.001$	$0.872 \pm 0.018$	$0.497 \pm 0.010$
HBOS*	$0.983 \pm 0.002$	$0.622 \pm 0.007$	$0.471 \pm 0.003$
LoOP*	$0.973 \pm 0.012$	$0.907 \pm 0.008$	$0.535 \pm 0.034$
uCBLOF*	$0.950 \pm 0.039$	$0.819 \pm 0.023$	$0.469 \pm 0.003$

Table 6.1: AUC (mean + standard deviation) for E-DPPs and standard outlier detection algorithms. As expected, we see that a higher exponent leads to a stronger preference for diversity and hence a better outlier detection scheme. Only LoOP and LOF consistently outperform E-DPPs.

Results are reported in Table 6.1; as expected, we see that larger values of  $p$  ( $p = 2$ ) are more sensitive to outliers, and provide better models for outlier detection.

### 6.3.3 E-Dpps for the Nyström method

As a more standard application of DPPs, we now investigate the use of E-DPPs for kernel reconstruction via the Nyström method (Nyström, 1930; Williams and Seeger, 2001). Given a large kernel  $\mathbf{K}$ , recall that the Nyström method selects a subset  $C$  of columns (“landmarks”) of  $\mathbf{K}$  and approximates  $\mathbf{K}$  as  $\mathbf{K}_{\cdot, C} \mathbf{K}_C^\dagger \mathbf{K}_{C, \cdot}$ . We show that E-DPPs improve upon the recent results of (Li et al., 2016c) for kernel reconstruction.

We apply Kernel Ridge Regression to 3 regression datasets: Ailerons, Bank32NH, and Machine CPU.<sup>2</sup> We subsample 4,000 points from each dataset (3,000 training and

<sup>2</sup><http://www.dcc.fc.up.pt/~ltorgo/Regression/DataSets.html>

1,000 test) and use an RBF kernel and choose the bandwidth  $\beta$  and regularization parameter  $\lambda$  for each dataset by 10-fold cross-validation. Results are averaged over 3 random subsets of data, using the swapchain sampler initialized with  $k$ -means++ and run for 3000 iterations.

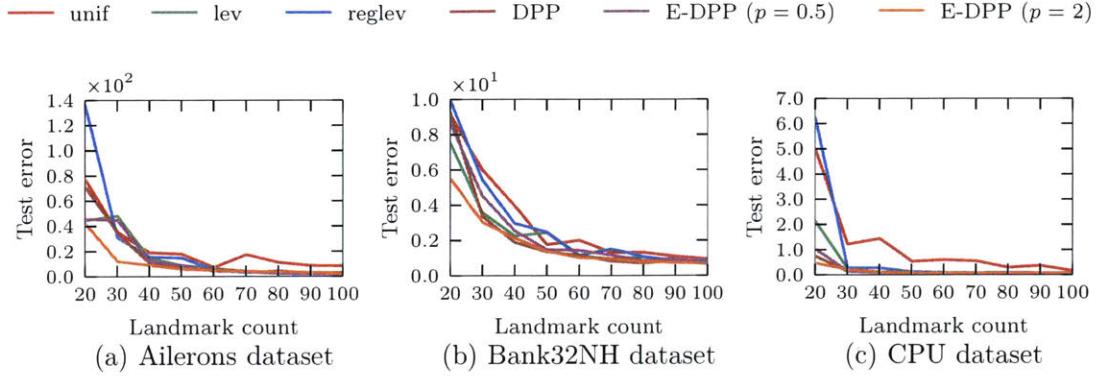


Figure 6-3: Prediction error on regression datasets; we compare various E-DPP models to uniform sampling (“unif”) as well as leverage and regularized leveraged sampling (“lev” and “reglev”). On all datasets, the E-DPPs achieve the lowest error, with the largest exponent  $p = 2$  performing markedly better than other methods.

We evaluate the quality of the sampler via the prediction error on the held-out test set. Figure 6-3 reports the results. Consistently across all datasets,  $p = 2$  outperforms all other samplers in terms of the prediction error, in particular when only sampling a few landmarks. Interestingly, we also see that the reconstruction error tends to be smaller when  $p = \frac{1}{2}$  (see Appendix A.2.4).

## 6.4 Open problems

- We believe that further specifying the class of ESR measures that also remain SR will provide valuable insight into the study of negatively associated measures.
- Learning the function and the exponent jointly: given an SR measure  $\mu$  and a collection of *i.i.d.* subsets  $\mathcal{S} = \{S_1, \dots, S_m\}$ , the MLE problem that finds the

best  $p > 0$  to model  $\mathcal{S}$  as being sampled from an ESR  $\nu \propto \mu^p$  is convex:

$$\operatorname{argmax}_{p>0} \frac{p}{m} \sum_{k=1}^m \log \mu(S_i) - \log \left( \sum_{S \subseteq [N]} \mu(S)^p \right). \quad (6.4)$$

As such, standard convex optimization algorithms can be leveraged to select  $p$ , potentially after leaning a parametrization of  $\mu$ . However, there may be advantages to learning  $p$  and  $\mu$  jointly, an approach which remain to be investigated.



# Chapter 7

## Conclusion

With the goal of introducing negative dependence as a powerful framework to model machine learning problems that require diverse sampling, this thesis focused on two specific problems. First, we derived efficient algorithms to learn and sample from negatively dependent measures, expanding their range of applications to real-world settings previously too computationally expensive to consider.

Then, we focused on the theory of negatively dependent measures itself. We introduced Gvs measures, strongly Rayleigh measures with key implications for several classical problems in machine learning and optimizations. Finally, we introduced exponentiated strongly Rayleigh measures: a class of negatively dependent measures parametrized by a scalar controlling the strength of repulsive forces between similar elements. With these contributions, we showed that negatively dependent measures can provide a fine-grained tool to model difficult machine learning problems.

We now conclude this thesis by re-framing it in the broader context of modern machine learning problems: how can we leverage the theory of negative dependence for principled machine learning model design and optimization?

### 7.1 High-level summary

This thesis has shown that the framework of negative dependence provides a unifying point of view for the broad class of machine learning optimization problems that —

at their core — seek diversity in their solutions. Within this framework, we derived theoretical results and general purpose algorithms for problems ranging from optimal design and recommender systems to outlier detection and kernel reconstruction.

The first part of this thesis focused on determinantal point processes, popular negatively dependent measures in machine learning. In Chapter 3 we showed that DPPs — as well as any other (log) super/submodular set function — can be approximated efficiently with neural networks. We obtained a DPP-specific neural architecture that efficiently approximates DPPs, both in sampling and mode-finding tasks. In Chapter 4, we scaled up DPP learning and sampling algorithms by leveraging Kronecker products, and introduced Contrastive Estimation, an efficient way of improving the DPP learning procedure by including information about unlikely subsets.

Chapter 5 introduced Generalized Volume Sampling measures: distributions over subsets based on elementary symmetric polynomials that extend volume sampling and remain strongly Rayleigh. We connected GVS measures to a generalization of classical problems in optimal experimental design, and developed randomized and greedy algorithms to solve the ensuing optimization problem. We obtained new results on the convexity properties of elementary symmetric polynomials, and derived a generalization of the classic Minimum Volume Covering Ellipsoid problem.

Finally, we introduced in Chapter 6 exponentiated strongly Rayleigh measures: a new class of negatively dependent measures that are well-suited to the requirements of machine learning applications. Indeed, exponentiated SR measures enjoy fast sampling and learning algorithms, while offering simple a simple tuning procedure for the strength of the measures’ repulsion forces.

## 7.2 Open problems

Each chapter included a summary of open research problems that we believe to be worth investigating. We conclude here by discussing some over-arching open problems of import to the machine learning community.

**Negative dependence theory.** Pemantle (2000) and Borcea et al. (2009) provided a variety of open conjectures regarding the chain of implications between the various characterizations of negative dependence that one may wish to use. By introducing exponentiated strongly Rayleigh measures, we have added to this list; in particular by raising the question of which class of measures will remain strongly Rayleigh when exponentiated. Furthermore, recent work by Anari et al. (2018) considers exponentiated-DPPs with exponents  $p \leq 1$  within the framework of *Complete Log-Concavity*, wherein the authors also derive weaker characterizations of negative dependence in this scenario, further confirming that the class of exponentiated strongly Rayleigh measures provides a valuable framework to develop and extend the rich theory of negative dependence.

**Negative dependence for machine learning optimization.** The scope of potential applications of negative dependence extends far beyond the examples discussed in this thesis. Zhang et al. (2017) showed that DPPs can provide variance reduction for minibatch SGD, a topic we believe warrants further study within the larger scope of negative dependence — indeed, minibatch SGD is in itself a particular example of negative dependence as we sample without replacement. Any improvements to the theory or speed of minibatch SGD will undeniably impact the fundamental underpinning of most machine learning applications. Furthermore, a less strict notion of negative dependence provides a framework to discuss trade-offs between exploration (diversity) and exploitation (quality), which occur — to mention only a few examples — in reinforcement learning and hyperparameter space exploration.



# Bibliography

- R. Affandi, A. Kulesza, E. Fox, and B. Taskar. Nyström approximation for large-scale Determinantal Point Processes. In *Proc. Int. Conference on Artificial Intelligence and Statistics (AISTATS)*, 2013.
- R. Affandi, E. Fox, R. Adams, and B. Taskar. Learning the parameters of Determinantal Point Process kernels. In *Proc. Int. Conference on Machine Learning (ICML)*, 2014.
- M. Amer and M. Goldstein. Nearest-neighbor and clustering based anomaly detection algorithms for rapidminer. In *RapidMiner Community Meeting and Conference (RCOMM)*. Shaker Verlag GmbH, 2012.
- N. Anari and S. O. Gharan. The Kadison-Singer problem for Strongly Rayleigh measures and applications to asymmetric TSP. *CoRR*, abs/1412.1143, 2014.
- N. Anari and S. Oveis Gharan. A generalization of permanent inequalities and applications in counting and optimization. In *Symposium on Theory of Computing (STOC)*. ACM, 2017.
- N. Anari, S. O. Gharan, and A. Rezaei. Monte Carlo Markov Chain algorithms for sampling strongly Rayleigh distributions and Determinantal Point Processes. In *Conference on Learning Theory (COLT)*, 2016.
- N. Anari, K. Liu, S. O. Gharan, and C. Vinzant. Log-concave polynomials II: high-dimensional walks and an FPRAS for counting bases of a matroid. *CoRR*, abs/1811.01816, 2018.
- C. Andrieu, N. de Freitas, A. Doucet, and M. I. Jordan. An introduction to mcmc for machine learning. *Machine Learning*, 50(1):5–43, Jan 2003.
- F. Angiulli and C. Pizzuti. Fast outlier detection in high dimensional spaces. In *European Conference on Principles of Data Mining and Knowledge Discovery (PKDD)*, London, UK, UK, 2002. Springer-Verlag.
- H. Avron and C. Boutsidis. Faster subset selection for matrices and applications. *SIAM J. Matrix Analysis Applications*, 34(4), 2013.
- F. Bach and M. I. Jordan. Kernel independent component analysis. *Journal of Machine Learning Research*, 3, 2002.

- D. Bahdanau, K. Cho, and Y. Bengio. Neural machine translation by jointly learning to align and translate. In *Int. Conf. on Learning Representations (ICLR)*, 2015.
- E. R. Barnes. An algorithm for separating patterns by ellipsoids. *IBM Journal of Research and Development*, 26, 1982.
- N. K. Batmanghelich, G. Quon, A. Kulesza, M. Kellis, P. Golland, and L. Bornn. Diversifying sparsity using variational Determinantal Point Processes. *CoRR*, abs/1411.6307, 2014.
- H. H. Bauschke, O. Güler, A. S. Lewis, and H. S. Sendov. Hyperbolic polynomials and convex analysis. *Canad. J. Math.*, 53(3), 2001.
- R. Bhatia. *Matrix Analysis*. Springer, 1997.
- R. Bhatia. *Positive Definite Matrices*. Princeton University Press, 2007.
- J. Borcea and P. Brändén. Applications of stable polynomials to mixed determinants: Johnson’s conjectures, unimodality, and symmetrized Fischer products. *Duke Mathematical Journal*, 143(2), 2008.
- J. Borcea and P. Brändén. The Lee-Yang and Pólya-Schur programs I Linear operators preserving stability. *Inventiones mathematicae*, 177(3), 2009.
- J. Borcea and P. Brändén. Multivariate Pólya–Schur classification problems in the Weyl algebra. *Proceedings of the London Mathematical Society*, 101(1), 2010.
- J. Borcea, P. Brändén, and T. Liggett. Negative dependence and the geometry of polynomials. *Journal of the American Mathematical Society*, 22(2), 2009.
- A. Borodin. Determinantal Point Processes. *arXiv:0911.1153*, 2009.
- A. Bose, H. Ling, and Y. Cao. Adversarial contrastive estimation. *arXiv preprint arXiv:1805.03642*, 2018.
- L. Bottou. On-line learning and stochastic approximations. In *On-line Learning in Neural Networks*, New York, NY, USA, 1998.
- P. Brändén. Polynomials with the half-plane property and matroid theory. *Advances in Mathematics*, 216(1), 2007.
- P. Brändén, J. Haglund, M. Visontai, and D. G. Wagner. *Proof of the Monotone Column Permanent Conjecture*, pages 63–78. Springer Basel, 2011.
- M. M. Breunig, H.-P. Kriegel, R. T. Ng, and J. Sander. LOF: Identifying density-based local outliers. *SIGMOD Rec.*, 29(2), May 2000.
- T. Brijs. Retail market basket data set. In *Workshop on Frequent Itemset Mining Implementations (FIMI03)*, 2003.

- T. Brijs, G. Swinnen, K. Vanhoof, and G. Wets. Using association rules for product assortment decisions: A case study. In *SIGKDD*. ACM, 1999.
- S. P. Brooks and A. Gelman. General methods for monitoring convergence of iterative simulations. *Journal of Computational and Graphical Statistics*, 7, 1998.
- N. Buchbinder, M. Feldman, J. Naor, and R. Schwartz. A tight linear time  $(1/2)$ -approximation for unconstrained submodular maximization. In *IEEE Symposium on Foundations of Computer Science (FOCS)*, 2012.
- H. Cai. Exact bound for the convergence of metropolis chains. *Stochastic Analysis and Applications*, 18(1), 2000.
- O. Canvet and F. Fleuret. Efficient Sample Mining for Object Detection. In *ACML*, JMLR: Workshop and Conference Proceedings, 2014.
- W. Chao, B. Gong, K. Grauman, and F. Sha. Large-margin Determinantal Point Processes. In *Uncertainty in Artificial Intelligence (UAI)*, 2015.
- D. Chen. Data mining for the online retail industry: A case study of RFM model-based customer segmentation using data mining. *Journal of Database Marketing and Customer Strategy Management*, 19(3), August 2012.
- L. Chen, G. Zhang, and E. Zhou. Fast greedy MAP inference for Determinantal Point Process to improve recommendation diversity. In *Advances in Neural Information Processing Systems*, 2018.
- L. Chen, F. Yuan, J. M. Jose, and W. Zhang. Improving negative sampling for word representation using self-embedded features. *CoRR*, abs/1710.09805, 2017.
- Y.-B. Choe, J. G. Oxley, A. D. Sokal, and D. G. Wagner. Homogeneous multivariate polynomials with the half-plane property. *Advances in Applied Mathematics*, 32(1), 2004. Special Issue on the Tutte Polynomial.
- D. A. Cohn. Neural network exploration using optimal experiment design. *Advances in Neural Information Processing Systems*, 1994.
- D. Comaniciu and P. Meer. Robust analysis of feature spaces: color image segmentation. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 1997.
- R. Cominetti, W. F. Mascarenhas, and P. J. S. Silva. A Newton's method for the continuous quadratic knapsack problem. *Mathematical Programming Computation*, 6(2), 2014.
- A. Cotter, M. R. Gupta, H. Jiang, J. Muller, T. Narayan, S. Wang, and T. Zhu. Interpretable set functions. *CoRR*, abs/1806.00050, 2018.

- T. A. Davis, W. W. Hager, and J. T. Hungerford. An efficient hybrid algorithm for the separable convex quadratic knapsack problem. *ACM Trans. Math. Softw.*, 42(3), May 2016.
- L. Decreusefond, I. Flint, N. Privault, and G. L. Torrisi. Determinantal Point Processes, 2015.
- P. Diaconis and L. Saloff-Coste. Comparison theorems for reversible Markov chains. *The Annals of Applied Probability*, 3(3), 1993.
- P. Diaconis and D. Stroock. Geometric bounds for eigenvalues of Markov Chains. *Annals of Applied Probability*, 1991.
- J. Djolonga and A. Krause. From MAP to marginals: Variational inference in Bayesian submodular models. In *Advances in Neural Information Processing Systems*, 2014.
- J. Djolonga, S. Tschiatschek, and A. Krause. Variational inference in mixed probabilistic submodular models. In *Advances in Neural Information Processing Systems*, 2016.
- A. N. Dolia, N. M. White, and C. J. Harris. D-optimality for minimum volume ellipsoid with outliers. In *In Proceedings of the Seventh International Conference on Signal/Image Processing and Pattern Recognition, (UkrOBRAZ2004)*, 2004.
- A. N. Dolia, T. De Bie, C. J. Harris, J. Shawe-Taylor, and D. M. Titterington. The minimum volume covering ellipsoid estimation in kernel-defined feature spaces. In *European Conference on Machine Learning*, 2006.
- J. L. Elman. Finding structure in time. *COGNITIVE SCIENCE*, 14(2):179–211, 1990.
- S. Flaxman, A. Wilson, D. Neill, H. Nickisch, and A. Smola. Fast Kronecker inference in Gaussian processes with non-Gaussian likelihoods. In *Proc. Int. Conference on Machine Learning (ICML)*, 2015.
- C. M. Fortuin, P. W. Kasteleyn, and J. Ginibre. Correlation inequalities on some partially ordered sets. *Comm. Math. Phys.*, 22(2), 1971.
- C. Fowlkes, S. Belongie, F. Chung, and J. Malik. Spectral grouping using the nyström method. *IEEE Trans. Pattern Anal. Mach. Intell.*, 26(2):214–225, January 2004.
- L. Gårding. An inequality for hyperbolic polynomials. *Journal of Mathematics and Mechanics*, 8(6), 1959.
- M. Gartrell, U. Paquet, and N. Koenigstein. Low-rank factorization of determinantal point processes. In *Uncertainty in Artificial Intelligence (UAI)*, 2017.

- J. Gillenwater, A. Kulesza, E. Fox, and B. Taskar. Expectation-Maximization for learning Determinantal Point Processes. In *Advances in Neural Information Processing Systems*, 2014.
- J. Gillenwater, A. Kulesza, and B. Taskar. Discovering diverse and salient threads in document collections. In *Proceedings of the 2012 Joint Conference on Empirical Methods in Natural Language Processing and Computational Natural Language Learning*, EMNLP-CoNLL '12, pages 710–720, 2012.
- J. A. Gillenwater, A. Kulesza, Z. Mariet, and S. Vassilvitskii. Maximizing Induced Cardinality Under a Determinantal Point Process. In *Advances in Neural Information Processing Systems*, 2018.
- J. A. Gillenwater, A. Kulesza, Z. Mariet, and S. Vassilvitskii. Sublinear Sampling for Determinantal Point Processes. *Proc. Int. Conference on Machine Learning (ICML)*, 2019.
- O. Goldschmidt, D. Nehme, and G. Yu. Note: On the set-union knapsack problem. *Naval Research Logistics*, 41, 1994.
- M. Goldstein and A. Dengel. Histogram-based outlier score (HBOS): A fast unsupervised anomaly detection algorithm. In *KI-2012: Poster and Demo Track. German Conference on Artificial Intelligence (KI-2012), 35th, September 24-27, Saarbrücken, Germany, 9 2012*.
- M. Goldstein and S. Uchida. A comparative evaluation of unsupervised anomaly detection algorithms for multivariate data. *PLOS ONE*, 11(4), 04 2016.
- I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio. Generative adversarial nets. In *Advances in Neural Information Processing Systems*, 2014.
- I. J. Goodfellow. On distinguishability criteria for estimating generative models, 2014.
- A. Gotovos, S. Hassani, and A. Krause. Sampling from probabilistic submodular models. In *Advances in Neural Information Processing Systems*, 2015.
- M. Grötschel, L. Lovász, and A. Schrijver. The ellipsoid algorithm and its consequences in combinatorial optimization. *Combinatorica*, 1, 1981.
- Y. Gu and Z. Jin. Neighborhood preserving d-optimal design for active learning and its application to terrain classification. *Neural Computing and Applications*, 23(7), 2013.
- O. Güler. Hyperbolic polynomials and interior point methods for convex programming. *Mathematics of Operations Research*, 22(2), 1997.
- M. U. Gutmann and A. Hyvärinen. Noise-contrastive estimation of unnormalized statistical models, with applications to natural image statistics. *J. Mach. Learn. Res.*, 13(1), February 2012.

- F. M. Harper and J. A. Konstan. The movielens datasets: History and context. *ACM Trans. Interact. Intell. Syst.*, 5(4):19:1–19:19, December 2015.
- T. Hastie, R. Tibshirani, and J. Friedman. *The Elements of Statistical Learning: data mining, inference and prediction*. Springer, 2001.
- W. K. Hastings. Monte Carlo Sampling Methods using Markov chains and their applications. *Biometrika*, 57(1), 1970.
- X. He. Laplacian regularized d-optimal design for active learning and its application to image retrieval. *IEEE Trans. Image Processing*, 19(1), 2010.
- R. A. Horn and C. R. Johnson. *Matrix Analysis*. Cambridge University Press, Cambridge, 1985.
- J. B. Hough, M. Krishnapur, Y. Peres, and B. Virág. Determinantal processes and independence. *Probability Surveys*, 3(206–229), 2006.
- Y. Hu, Y. Koren, and C. Volinsky. Collaborative filtering for implicit feedback datasets. In *ICDM*, 2008.
- D. A. Jackson and Y. Chen. Robust principal component analysis and outlier detection with ecological data. *Environmetrics*, 15(2), 2004.
- T. Jain. Derivatives for antisymmetric tensor powers and perturbation bounds. *Linear Algebra and its Applications*, 435(5), 2011.
- J.-M. Jolion, P. Meer, and S. Bataouche. Robust clustering with applications in computer vision. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 13(8), Aug 1991.
- M. I. Jordan. Artificial neural networks. chapter Attractor Dynamics and Parallelism in a Connectionist Sequential Machine, 1990.
- R. Jozsa and G. Mitchison. Symmetric polynomials in information theory: Entropy and subentropy. *Journal of Mathematical Physics*, 56(6), 2015.
- B. Kang. Fast Determinantal Point Process sampling with application to clustering. In *Advances in Neural Information Processing Systems*, 2013.
- D. P. Kingma and M. Welling. Auto-encoding variational bayes. In *International Conference on Learning Representations*, 2014.
- A. Krause, A. Singh, and C. Guestrin. Near-optimal sensor placements in Gaussian processes: Theory, efficient algorithms and empirical studies. *Journal of Machine Learning Research*, 9, 2008.
- H.-P. Kriegel, P. Kröger, E. Schubert, and A. Zimek. LoOP: Local outlier probabilities. In *Proceedings of the 18th ACM Conference on Information and Knowledge Management*, CIKM ’09, New York, NY, USA, 2009. ACM.

- A. Kulesza and B. Taskar. k-DPPs: Fixed-size Determinantal Point Processes. In *Proc. Int. Conference on Machine Learning (ICML)*, 2011.
- A. Kulesza and B. Taskar. *Determinantal Point Processes for machine learning*, volume 5. Foundations and Trends in Machine Learning, 2012.
- A. Kulesza. *Learning with Determinantal Point Processes*. PhD thesis, University of Pennsylvania, 2012.
- F. S. Lasheras, J. V. Vilán, P. G. Nieto, and J. del Coz Díaz. The use of design of experiments to improve a neural network model in order to predict the thickness of the chromium layer in a hard chromium plating process. *Mathematical and Computer Modelling*, 52(78), 2010. Mathematical Models in Medicine, Business & Engineering 2009.
- F. Lavancier, J. Møller, and E. Rubak. Determinantal Point Process models and statistical inference. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 77(4), 2015.
- Y. LeCun and C. Cortes. MNIST handwritten digit database. 2010.
- A. S. Lewis. Derivatives of spectral functions. *Math. Oper. Res.*, 21(3), 1996.
- A. Li, W. Luo, S. Nagavalli, N. Chakraborty, and K. Sycara. Handling state uncertainty in distributed information leader selection for robotic swarms. In *2016 IEEE International Conference on Systems, Man, and Cybernetics (SMC)*, Oct 2016a.
- C. Li, S. Sra, and S. Jegelka. Gaussian quadrature for matrix inverse forms with applications. In *Proc. Int. Conference on Machine Learning (ICML)*, 2016b.
- C. Li, S. Jegelka, and S. Sra. Efficient sampling for  $k$ -Determinantal Point Processes. *CoRR*, abs/1509.01618, 2015.
- C. Li, S. Jegelka, and S. Sra. Fast DPP sampling for Nyström with application to kernel methods. In *Proc. Int. Conference on Machine Learning (ICML)*, ICML'16. JMLR.org, 2016c.
- C. Li, S. Jegelka, and S. Sra. Fast mixing Markov chains for strongly Rayleigh measures, DPPs, and constrained sampling. In *Advances in Neural Information Processing Systems*, 2016d.
- C. Li, S. Jegelka, and S. Sra. Polynomial time algorithms for dual volume sampling. In *Advances in Neural Information Processing Systems*, 2017.
- Y. Li, J. Hu, C. Zhai, and Y. Chen. Improving one-class collaborative filtering by incorporating rich user information. In *CIKM*, 2010.
- Z. Lin, Y. Liu, G. Molteni, and D.Zhang. Spectral properties for a new composition of a matrix and a complex representation. *Electronic J. Linear Algebra*, 23, 2012.

- Z. Liu, P. Luo, X. Wang, and X. Tang. Deep learning face attributes in the wild. In *International Conference on Computer Vision*, December 2015.
- C. V. Loan and N. Pitsianis. Approximation with kronecker products. In *Linear Algebra for Large Scale and Real Time Applications*. Kluwer Publications, 1993.
- R. Lyons. Determinantal probability measures. *Publications Mathématiques de l’Institut des Hautes Études Scientifiques*, 98(1), 2003.
- O. Macchi. The coincidence approach to stochastic point processes. *Adv. Appl. Probab.*, 7, 1975.
- I. G. Macdonald. *Symmetric functions and Hall polynomials*. Oxford university press, 1998.
- A. Marcus, D. Spielman, and N. Srivastava. Interlacing families II: Mixed characteristic polynomials and the Kadison–Singer problem. *Annals of Mathematics*, July 2015.
- Z. Mariet and S. Sra. Fixed-point algorithms for learning Determinantal Point Processes. In *Proc. Int. Conference on Machine Learning (ICML)*, 2015.
- Z. Mariet and S. Sra. Diversity networks: Neural network compression using Determinantal Point Processes. *Int. Conf. on Learning Representations (ICLR)*, 2016a.
- Z. Mariet and S. Sra. Kronecker Determinantal Point Processes. In *Advances in Neural Information Processing Systems*, 2016b.
- Z. Mariet and S. Sra. Elementary symmetric polynomials for optimal experimental design. In *Advances in Neural Information Processing Systems*, 2017.
- Z. Mariet, S. Sra, and S. Jegelka. Exponentiated Strongly Rayleigh distributions. In *Advances in Neural Information Processing Systems*, 2018.
- Z. Mariet, M. Gartrell, and S. Sra. Learning Determinantal Point Processes by sampling inferred negatives. In *Proc. Int. Conference on Artificial Intelligence and Statistics (AISTATS)*, 2019a.
- Z. Mariet, Y. Ovadia, and J. Snoek. DPPNET: Approximating Determinantal Point Processes with deep networks. *Submitted to NeurIPS*, 2019b.
- J. Martens and R. B. Grosse. Optimizing neural networks with Kronecker-factored approximate curvature. In *Proc. Int. Conference on Machine Learning (ICML)*, 2015.
- M. D. McKay, R. J. Beckman, and W. J. Conover. Comparison of three methods for selecting values of input variables in the analysis of output from a computer code. *Technometrics*, 21(2):239–245, 1979.

- N. Metropolis, A. W. Rosenbluth, M. N. Rosenbluth, A. H. Teller, and E. Teller. Equation of state calculations by fast computing machines. *The Journal of Chemical Physics*, 21(6), 1953.
- B. Micenková, B. McWilliams, and I. Assent. Learning outlier ensembles: The best of both worlds – supervised and unsupervised. In *Proceedings of the ACM SIGKDD 2014 Workshop on outlier Detection and Description under Data Diversity*, 2014.
- T. Mikolov, I. Sutskever, K. Chen, G. Corrado, and J. Dean. Distributed representations of words and phrases and their compositionality. In *Advances in Neural Information Processing Systems*, 2013.
- A. Mnih and Y. W. Teh. A fast and simple algorithm for training neural probabilistic language models. In *Proc. Int. Conference on Machine Learning (ICML)*, 2012.
- D. M. Mount, N. S. Netanyahu, C. D. Piatko, R. Silverman, and A. Y. Wu. Quantile approximation for robust statistical estimation and k-enclosing problems, 2000.
- W. W. Muir. Inequalities concerning the inverses of positive definite matrices. *Proceedings of the Edinburgh Mathematical Society*, 19(2), 1974.
- E. Nyström. Über die praktische Auflösung von Integralgleichungen mit Anwendungen auf Randwertaufgaben. *Acta Mathematica*, 54(1), 1930.
- T. Osogami, R. Raymond, A. Goel, T. Shirai, and T. Maehara. Dynamic Determinantal Point Processes. In *Proc. AAAI Conference on Artificial Intelligence*, 2018.
- R. Pemantle. Towards a theory of negative dependence. *Journal of Mathematical Physics*, 41, 2000.
- F. Pukelsheim. *Optimal Design of Experiments*. Society for Industrial and Applied Mathematics, 2006.
- P. Rebeschini and A. Karbasi. Fast mixing for discrete point processes. In *Conference on Learning Theory (COLT)*, 2015.
- P. J. Rousseeuw and A. M. Leroy. *Robust Regression and Outlier Detection*. John Wiley & Sons, Inc., New York, NY, USA, 1987.
- P. J. Rousseeuw and B. C. van Zomeren. Unmasking multivariate outliers and leverage points. *Journal of the American Statistical Association*, 85(411), 1990.
- G. Sagnol. *Optimal design of experiments with application to the inference of traffic matrices in large networks: second order cone programming and submodularity*. Theses, École Nationale Supérieure des Mines de Paris, December 2010.
- A. Schein and L. Ungar. A-Optimality for Active Learning of Logistic Regression Classifiers, 2004.

- H. Shen, S. Jegelka, and A. Gretton. Fast kernel-based independent component analysis. *IEEE Trans. on Signal Processing*, 57(9), 2009.
- T. Shirai and Y. Takahashi. Random point fields associated with certain Fredholm determinants I: fermion, Poisson and boson point processes. *Journal of Functional Analysis*, 205(2), 2003.
- A. Shpilka and A. Wigderson. Depth-3 arithmetic circuits over fields of characteristic zero. *computational complexity*, 10(1), 2001.
- A. Shrivastava, A. Gupta, and R. Girshick. Training region-based object detectors with online hard example mining. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016.
- J. D. Smith and M. T. Thai. Breaking the bonds of submodularity: Empirical estimation of approximation ratios for monotone non-submodular greedy maximization. *CoRR*, abs/1702.07002, 2017.
- N. A. Smith and J. Eisner. Guiding unsupervised grammar induction using contrastive estimation. In *In Proc. of IJCAI Workshop on Grammatical Inference Applications*, 2005a.
- N. A. Smith and J. Eisner. Contrastive estimation: Training log-linear models on unlabeled data. In *ACL*, ACL '05, 2005b.
- W. Specht. Zur Theorie der elementaren Mittel. *Math. Zeitschr.*, 74, 1960.
- S. Sra and R. Hosseini. Conic geometric optimization on the manifold of positive definite matrices. *SIAM J. Optimization (SIOPT)*, 25(1), 2015.
- N. Street, W. H. Wolberg, and O. L. Mangasarian. Nuclear feature extraction for breast tumor diagnosis, 1993.
- P. Sun and R. M. Freund. Computation of minimum-volume covering ellipsoids. *Operations Research*, 52(5), 2004.
- K. K. Sung. *Learning and Example Selection for Object and Pattern Detection*. PhD thesis, Massachusetts Institute of Technology, 1996.
- A. Talwalkar, S. Kumar, M. Mohri, and H. Rowley. Large-scale svd and manifold learning. *J. Mach. Learn. Res.*, 14(1):3129–3152, January 2013.
- M. Todd. *Minimum-Volume Ellipsoids*. Society for Industrial and Applied Mathematics, Philadelphia, PA, 2016.
- S. Tschiatschek, J. Djolonga, and A. Krause. Learning probabilistic submodular diversity models via noise contrastive estimation. In *Proc. Int. Conference on Artificial Intelligence and Statistics (AISTATS)*, 2016.

- L. Vandenberghe, S. Boyd, and S.-P. Wu. Determinant maximization with linear matrix inequality constraints. *SIAM J. Matrix Anal. Appl.*, 19(2), 1998.
- A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. u. Kaiser, and I. Polosukhin. Attention is all you need. In *Advances in Neural Information Processing Systems*, 2017.
- Y. Wang, A. W. Yu, and A. Singh. On computationally tractable selection of experiments in regression models, 2016a.
- Z. Wang, B. Moran, X. Wang, and Q. Pan. Approximation for maximizing monotone non-decreasing set functions with a greedy method. *J. Comb. Optim.*, 31(1), January 2016b.
- W. J. Welch. Algorithmic complexity: three NP-hard problems in computational statistics. *Journal of Statistical Computation and Simulation*, 15(1), 1982.
- M. Wilhelm, A. Ramanathan, A. Bonomo, S. Jain, E. Chi, and J. Gillenwater. Practical Diversified Recommendations on YouTube with Determinantal Point Processes. In *Proc. Conference on Information and Knowledge Management (CIKM)*, 2018.
- C. K. I. Williams and M. Seeger. Using the Nyström method to speed up kernel machines. In *Advances in Neural Information Processing Systems*, 2001.
- G. Wu, Z. Zhang, and E. Y. Chang. Kronecker factorization for speeding up kernel machines. In *SIAM Data Mining (SDM)*, 2005.
- T. C. Xygkis, G. N. Korres, and N. M. Manousakis. Fisher information based meter placement in distribution grids via the d-optimal experimental design. *IEEE Transactions on Smart Grid*, PP(99), 2016.
- I.-C. Yeh. Modeling of strength of high-performance concrete using artificial neural networks. *Cement and Concrete Research*, 28(12), 1998.
- A. L. Yuille and A. Rangarajan. The concave-convex procedure (CCCP). In *Advances in Neural Information Processing Systems*, 2002.
- S. Zagoruyko and N. Komodakis. Wide residual networks. In *BMVC*, 2016.
- M. Zaheer, S. Kottur, S. Ravanbakhsh, B. Póczos, R. R. Salakhutdinov, and A. J. Smola. Deep sets. In *Advances in Neural Information Processing Systems*, 2017.
- C. Zhang, H. Kjellström, and S. Mandt. Stochastic learning on imbalanced data: Determinantal Point Processes for mini-batch diversification. *CoRR*, abs/1705.00607, 2017.
- X. Zhang, F. X. Yu, R. Guo, S. Kumar, S. Wang, and S.-F. Chang. Fast orthogonal projection based on kronecker product. In *Int. Conference on Computer Vision (ICCV)*, 2015.

T. Zhou, Z. Kuscsik, J.-G. Liu, M. Medo, J. R. Wakeling, and Y.-C. Zhang. Solving the apparent diversity-accuracy dilemma of recommender systems. *PNAS*, 107(10), 2010.

J. Y. Zou and R. Adams. Priors for diversity in generative latent variable models. In *Advances in Neural Information Processing Systems*, 2012.

# Appendix A

## A.1 Mathematical tools

### A.1.1 Geodesic convexity

**Definition A.1.1** (geodesic-convexity). A function  $f : \mathbb{S}^{++} \rightarrow \mathbb{R}$  defined on the Riemannian manifold  $\mathbb{S}^{++}$  is called *geodesically convex* if it satisfies

$$f(\mathbf{P} \#_t \mathbf{Q}) \leq (1-t)f(\mathbf{P}) + tf(\mathbf{Q}), \quad t \in [0, 1], \text{ and } \mathbf{P}, \mathbf{Q} \succ 0.$$

where we use the notation  $\mathbf{P} \#_t \mathbf{Q} := \mathbf{P}^{1/2} (\mathbf{P}^{-1/2} \mathbf{Q} \mathbf{P}^{-1/2})^t \mathbf{P}^{1/2}$  to denote the *geodesic* between  $\mathbf{P}$  and  $\mathbf{Q} \in \mathbb{S}^{++}$  under the Riemannian metric  $g_{\mathbf{P}}(\mathbf{X}, \mathbf{Y}) = \text{tr}(\mathbf{P}^{-1} \mathbf{X} \mathbf{P}^{-1} \mathbf{Y})$ .

### A.1.2 Power-mean inequality

**Theorem A.1.1** (Specht (1960)). Let  $x_i > 0$  and  $w_i \geq 0$  for  $1 \leq i \leq N$  such that  $\sum_i w_i = 1$ . Let  $p < q \in \mathbb{R}$  such that  $pq \neq 0$ . Then, letting  $\kappa = \frac{\max x_i}{\min x_i}$ ,

$$1 \leq \frac{M_q(\mathbf{w}; \mathbf{x})}{M_p(\mathbf{w}; \mathbf{x})} \leq \left( \frac{q-p}{q} \frac{\kappa^q - 1}{\kappa^q - \kappa^p} \right)^{\frac{1}{p}} \left( \frac{p}{q-p} \frac{\kappa^q - \kappa^p}{\kappa^p - 1} \right)^{\frac{1}{q}}.$$

where the power mean  $M_p(\mathbf{w}; \mathbf{x})$  is defined as

$$M_p(m; x) := \left( \sum_{i=1}^N w_i x_i^p \right)^{\frac{1}{p}}.$$

**Theorem A.1.2.** Let  $L \in \mathbb{P}^n$  be a positive definite matrix and  $S \subseteq [n]$ . Then,

$$\begin{aligned}\det(L[S])^p &\geq \det(L^p[S]), & 0 \leq p \leq 1, \\ \det(L[S])^p &\leq \det(L^p[S]), & p \geq 1.\end{aligned}$$

*Proof.* From Lemma 6.2.7, there exists a vector  $\mathbf{w}$  in the probability simplex, of size  $\binom{n}{|S|}$ , such that

$$\det(L[S]) = \sum_{J \subseteq [n], |J|=|S|} w_J \prod_{i \in J} \lambda_i.$$

Since  $t \rightarrow t^p$  is convex for  $p \geq 1$ , Jensen's inequality shows that

$$\det(L[S])^p \leq \sum_{J \subseteq [n], |J|=|S|} w_J \prod_{i \in J} \lambda_i^p = \det(L^p[S]),$$

where the latter equality follows due to  $L$  and  $L^p$  sharing the same eigenbasis. The same reasoning for  $p < 1$  gives the other side of the inequality.  $\square$

## A.2 Additional experimental results

### A.2.1 Encoder details for DppNet

For the MNIST encodings, the VAE encoder consists of a 2d-convolutional layer with 64 filters of height and width 4 and strides of 2, followed by a 2d convolution layer with 128 filters (same height, width and strides), then by a dense layer of 1024 neurons. The encodings are of length 32.

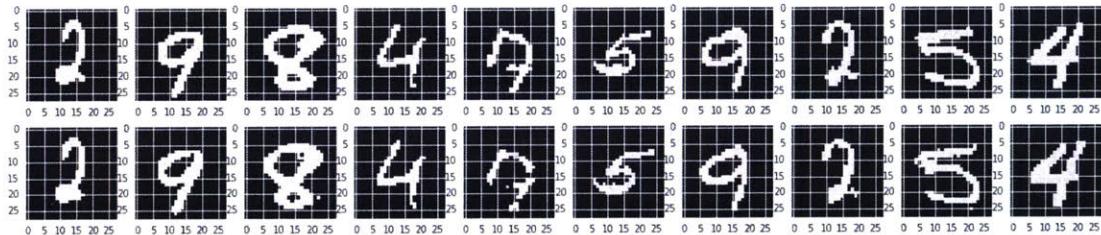


Figure A-1: Digits and VAE reconstructions from the MNIST training set

CelebA encodings were generated by a VAE using a Wide Residual Network en-

coder (Zagoruyko and Komodakis, 2016) with 10 layers and filter-multiplier  $k = 4$ , a latent space of 32 full-covariance Gaussians, and a deconvolutional decoder trained end-to-end using an ELBO loss. In detail, the decoder architecture consists of a 16K dense layer followed by a sequence of  $4 \times 4$  convolutions with [512, 256, 128, 64] filters interleaved with  $2 \times$  upsampling layers and a final  $6 \times 6$  convolution with 3 output channels for each of 5 components in a mixture of quantized logistic distributions representing the decoded image.

### A.2.2 Amazon Baby registries dataset

The Amazon baby registries dataset is a small, standard dataset used to evaluate the performance of DPPs.

Table A.1: Description of the Amazon Baby registries dataset.

REGISTRY	NUMBER OF ITEMS	NUMBER OF SETS
FURNITURE	32	1541
CARSEATS	34	2005
SAFETY	36	1654
STROLLERS	40	2132
MEDIA	58	2730
TOYS	62	4728
HEALTH	62	6598
BATH	100	6877
APPAREL	100	8102
BEDDING	100	8899
DIAPER	100	10,504
GEAR	100	8861
FEEDING	100	12,612

### A.2.3 Elementary Symmetric Polynomials

To compare to (Wang et al., 2016a), we generated the experimental matrix  $X$  by sampling  $n$  vectors of size  $m$  from the multivariate Gaussian distribution of mean 0 and covariance  $\Sigma = \text{Diag}(1^{-\alpha}, \dots, m^{-\alpha})$  for various sizes of  $\alpha$  and multiple budgets  $k$ , with  $m = 50, n = 1000$ ;  $\alpha$  controls hows *skewed* the distribution is.

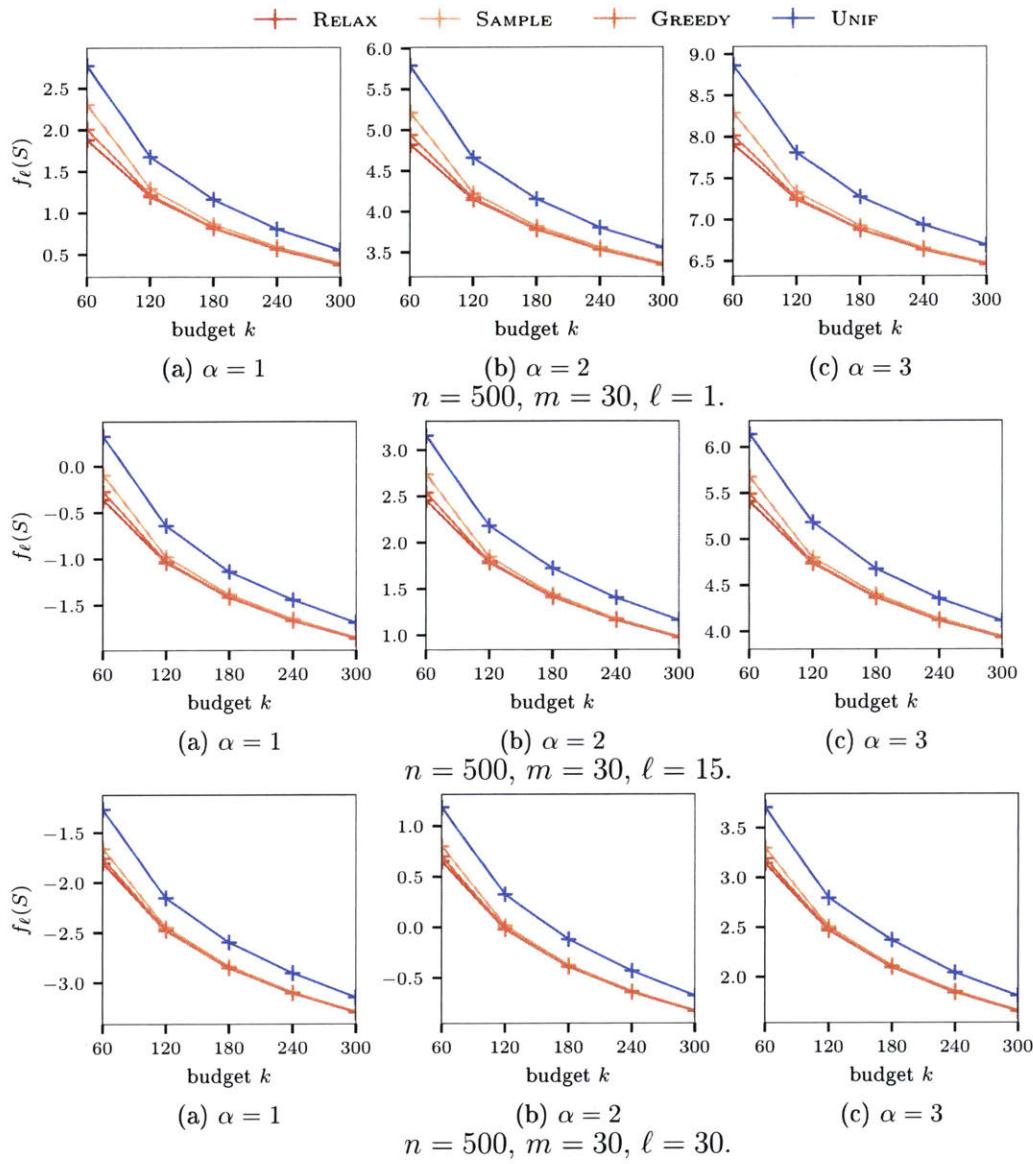


Figure A-2: Synthetic experiments with skewed covariance matrix

Results averaged over 5 runs are reported in Table A.2 and Figure A-2; standard deviations are too small to appear in Figure A-2. We also report more extensive results on synthetic experiments with a sparse precision matrix in Figure A-3.

Table A.2:  $\|z\|_0$  for  $n = 500, m = 30, \ell = 15$

	$k = 60$	$k = 120$	$k = 180$	$k = 240$	$k = 300$
$\alpha = 1$	$167 \pm 9$	$192 \pm 6$	$241 \pm 5$	$290 \pm 4$	$335 \pm 4$
$\alpha = 2$	$160 \pm 4$	$187 \pm 5$	$240 \pm 2$	$284 \pm 3$	$331 \pm 6$
$\alpha = 3$	$160 \pm 4$	$190 \pm 3$	$237 \pm 5$	$281 \pm 4$	$333 \pm 3$

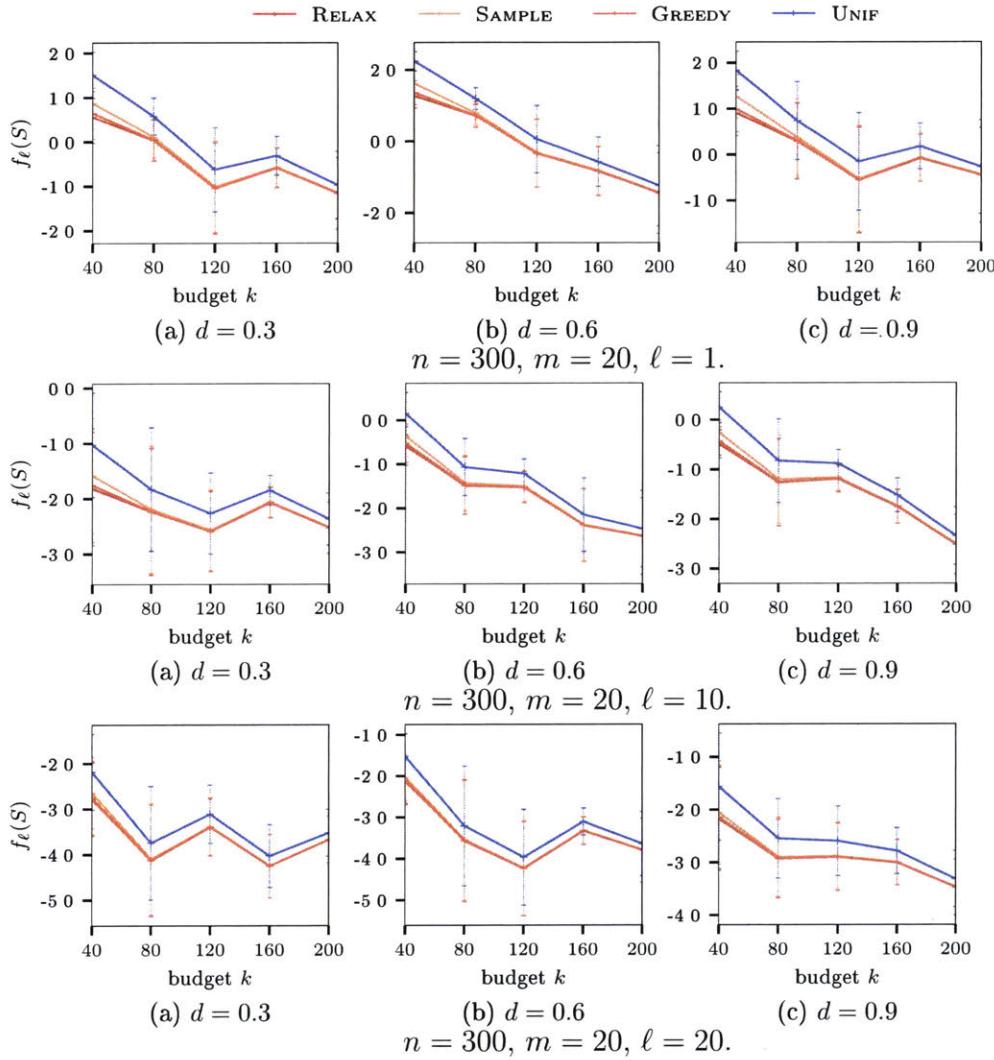


Figure A-3: Synthetic experiments with sparse precision matrix.

#### A.2.4 Exponentiated strongly Rayleigh measures

Figure A-4 presents the evolution of the closed-form upper-bound for the  $r$ -closeness term for E-DPPs, depending on the choice of proposal distributions.

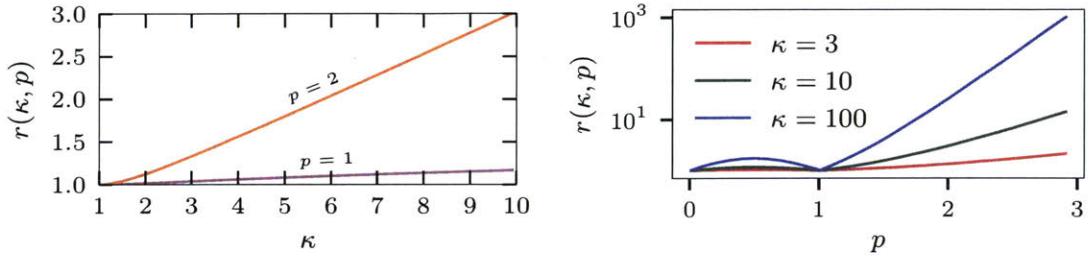


Figure A-4: Evolution of the upper bound for  $r(p, \kappa)$  from Thm. 6.2.8, which measures the  $r$ -closeness between the E-DPP with kernel  $L$  and the DPP with kernel  $L^p$ .

In Figure A-5, we show the empirical mixing time of the MCMC samplers for exponentiated DPPs when varying the size of the ground set size  $N$ .

Finally, additional results for the kernel reconstruction experiments are reported in Figure A-6 and A-7; these results evaluate the quality of the reconstructed kernel for different metrics.

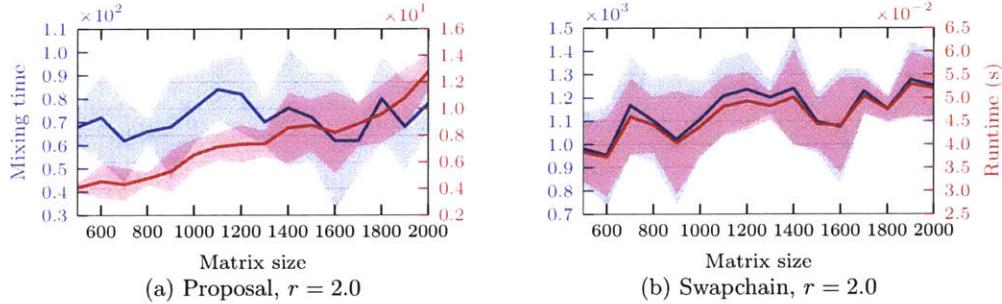


Figure A-5: Influence of ground set size  $N$  on mixing time

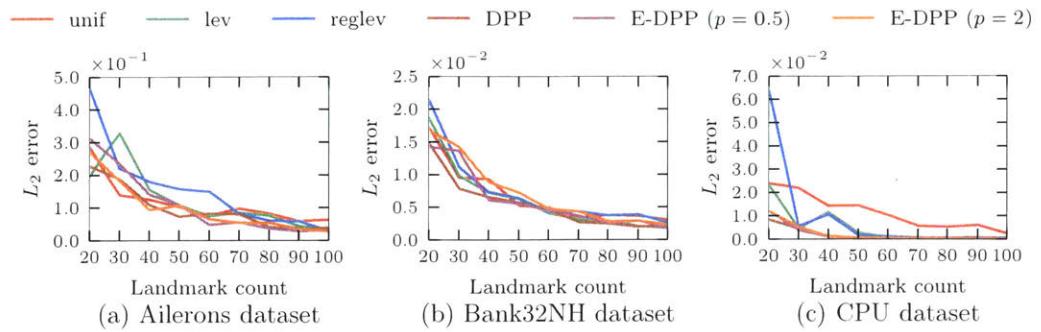


Figure A-6:  $L_2$  reconstruction error

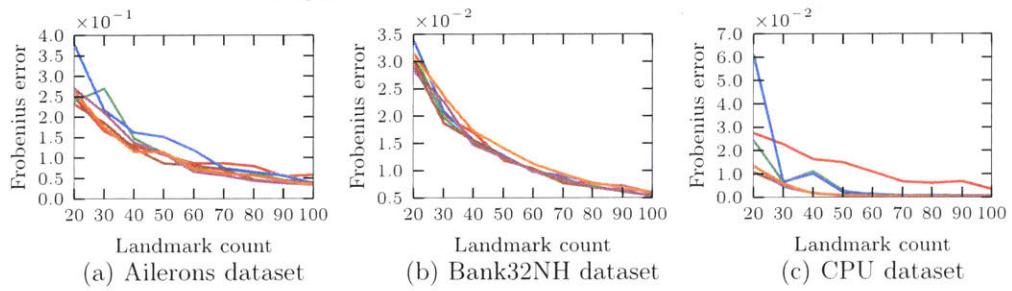


Figure A-7: Frobenius norm reconstruction error