

Navigating the Pachner Graph: Algorithms for Searching and Sampling Triangulations

Daniel Bruwel

An essay submitted in partial fulfillment of
the requirements for the degree of
B.Sc. (Honours)

Pure Mathematics
University of Sydney



October 2025

CONTENTS

Introduction	4
Chapter 1. Preliminaries	6
1.1. Mathematical Preliminaries	6
1.2. Knot Theory	8
1.3. Machine Learning Preliminaries	12
Chapter 2. Search and Sampling Algorithms	21
2.1. Classical Techniques	21
2.2. Transformers	22
Chapter 3. Numerical Experiments	24
3.1. Objective Functions	24
3.2. Markov Chain Monte Carlo	25
3.3. Classical Optimisation	29
3.4. Comparison of Classical Optimisation	33
3.5. Baseline Transformer Efficiency on Isomorphism Signatures	35
Chapter 4. Discussion and Future Work	36

Introduction

The process of decomposing a space or surface into discrete units such as triangles and tetrahedra is a standard technique used in a range of disciplines to take the analytic problem of studying a space or surface, to a combinatorial problem that opens it up to a range of computational techniques. Some examples of where triangulations are used includes physics such as numeric relativity where spacetime is triangulated for simulations; engineering where fluid and solid dynamics use “meshes” (triangulation) to approximate wings or bridges so that each component can be studied easily; computer graphics where objects are rendered with triangulations; and data science where the data can be used to generate a triangulation that can then be explored to reveal more about the system. Triangulations also have many important purposes within maths itself, where they can be used to compute topological properties like homology groups and the Euler characteristic, differential geometry with fields such as Regge calculus and combinatorial Ricci flows, and graph theory where concepts from triangulations where used in proofs such as the proof of the four colouring theorem. The main challenge with triangulations of surfaces and spaces, is that they are not unique, and different triangulations of a surface are better for different tasks. For example, finding the triangulation with the least faces helps in algebraic topology, whereas in computer graphics the triangulation that best approximates the geometry of the original surface is preferred. As such, considerable effort has been put into trying to explore and understand the world of triangulations of a given surface and find techniques to generate “good” triangulations. Furthermore, it turns out that understanding at least the enumeration of this space - that is the number of triangulations of a space with N triangles, is critical in the study of quantum gravity.

While there are a range of different notions of what is meant by good, this report primarily explores topological and combinatorial properties of triangulations. The techniques can undoubtedly be extended to a range of different notions of goodness, but topological and combinatorial properties are fundamental to the triangulation itself, and do not depend on the specific geometry like the angles within or the size of specific triangles in the triangulation. This has the advantage of narrowing the focus of exploration to abstract triangulations which are discrete and countable, improving computational manageability and ensuring that the results are general and only dependent on the triangulation itself, not other erroneous geometric factors. As such, this report primarily focuses around exploring properties of the triangulations of a single vertex 3-sphere that arise from the connection between such a triangulation, and a set of knots related to the edges within the triangulation. Looking for triangulations that relate to interesting and complex knots.

The primary objective of this report is to implement, analyse, and compare a range of different computational techniques for searching the space of triangulations. We investigate classical search and sampling techniques, namely Markov Chain Monte Carlo, greedy search, and simulated annealing. While greedy search algorithms have been a central tool to the study of this problem, Markov Chain Monte Carlo is a relatively new technique (2023) and we introduce the application of simulated annealing to this problem for the first time. Additionally, we introduce the groundwork for a new approach that uses transformer-based models to explore the space, opening the door for reinforcement learning which is steadily developing a repertoire of success in a range of different problems.

The report is structured as follows. Chapter 1 brings together the necessary preliminaries, including manifolds, some knot theory, and an overview of various statistical and machine learning techniques. Chapter 2 provides a detailed outline and discussion of the specific algorithms that we are implementing, including justifications for the choices we made. Chapter 3 presents a description and analysis of a series of numeric experiments run across 5 different search problems and compares them. Finally, chapter 4 concludes with a discussion of our findings and outlines the direction for future research. [?]

CHAPTER 1

Preliminaries

1.1. Mathematical Preliminaries

1.1.1. Manifolds. A n -dimensional manifold or a n -manifold is a topological generalisation of a surface, they are topological spaces that are “well behaved” and locally homeomorphic to \mathbb{R}^n . We recall the following basic definitions before precisely defining a manifold.

Definition 1.1 (Second Countable). A Topological space \mathcal{T} is Second Countable if the topology $\tau_{\mathcal{T}}$ of \mathcal{T} has a countable base.

Definition 1.2 (Hausdorff). A Topological space \mathcal{T} is Hausdorff if for any two points $x \neq y$ in \mathcal{T} there are neighbourhoods U of x and V of y such that $U \cap V = \emptyset$

With these we can define

Definition 1.3 (n -dimensional manifold). A n -dimensional manifold is a Hausdorff, second countable, topological space \mathcal{M} where for every point $p \in \mathcal{M}$ there is a neighbourhood $U(p)$ of p that is Homeomorphic to an open subset of Euclidean space \mathbb{R}^n

Remark 1.4. We have used a Homeomorphism to an open set of \mathbb{R}^n , this is equivalent to saying that each point is either isolated (if $n = 0$) or has a neighbourhood that is Homeomorphic to all of \mathbb{R}^n

Remark 1.5. We will refer to these simply as Manifolds, compared to other areas of research these may be called *topological manifolds* to emphasise that they do not have any further structure (c.f. *differentiable manifolds*, *Riemannian manifolds*, etc.)

Some examples of Manifolds include the circles S^1 , the sphere S^2 , the Torus T^2 , Euclidean space \mathbb{R}^n . Some non-examples include the disk D^2 because its boundary has no neighbourhood that is homeomorphic to an open subset of \mathbb{R}^n , the figure-eight curve because at the “crossing” there is no way to deform it to look like an open subset of \mathbb{R} , or the line with two origins because this is clearly not Hausdorff [?].

In the study of triangulations we are often interested in a specific category of manifolds known as *piecewise-linear manifolds* or *PL-manifolds* these are manifolds that require each point p to have a neighbourhood $U(p)$ of p that is a piecewise-linear deformation of \mathbb{R}^n . We say that two PL-manifolds are *PL-homeomorphic* if they can be transformed into each other via piecewise-linear

deformations. This is all more formally defined in terms of charts, atlases, and transition maps which can be found in almost any introductory book on geometric topology (e.g. [?]).

A particularly useful result due to Tibor Radó and Edwin Moise for low dimensional topology says that for dimension 3 or less PL-manifolds and PL-homeomorphisms are the “same” as manifolds and homeomorphisms. More precisely

Theorem 1.6 (Hauptvermutung). *For dimension 3 or less, every manifold is homeomorphic to a unique PL-manifold up to PL-homeomorphism.*

The proof of this theorem can be found in e.g. [?, ?]

Remark 1.7. This theorem does not hold for *four* or more dimensions [?].

1.1.2. Triangulations. We begin by defining *simplices* - the fundamental building blocks of triangulations

Definition 1.8 (simplices). *An n -simplex Δ is the n -dimensional convex hull of $n + 1$ vertices.*

A 0-simplex is a point, a 1-simplex a line, a 2-simplex a triangle, and a 3-simplex a tetrahedron. One way to construct an n -simplex is by taking the *join* of n points where we recall that

Definition 1.9 (Join). *For two topological spaces X and Y , the join denoted $X * Y$ is constructed by taking $X \times Y \times [0, 1]$ and quotienting by the equivalence relation that $(x, y, 0) \sim (x, y', 0)$ and $(x, y, 1) \sim (x', y, 1)$ for $x, x' \in X$ and $y, y' \in Y$.*

Any subset of $i + 1$ points of an n -simplex forms another simplex that we call an i -dimensional *face*.

We can take simplices and “glue” them together to form more complex geometric structures. Formally this is done via *face gluings*

Definition 1.10 (face gluing). *For two n -simplices Δ_1 and Δ_2 and two i -dimensional faces $F_1 \subseteq \Delta_1$ and $F_2 \subseteq \Delta_2$ we define the face gluing of F_1 to F_2 by taking $\Delta_1 \sqcup \Delta_2 / \sim$ where \sim is defined by a gluing map which is a homeomorphism between F_1 and F_2*

We typically create the gluing map by taking a bijection between the vertices of F_1 and F_2 and linearly interpolating to create the full homeomorphism.

We can easily extend this to define gluings of faces of the same triangulation, or even a gluing of a face to itself.

Definition 1.11 (PL-Sphere). *The PL 0-sphere is a pair of isolated points. For integers $n > 0$ the PL n -sphere is a PL manifold homeomorphic to the n -sphere.*

Definition 1.12 (PL-Triangulation). *Construct \mathcal{T} from a finite collection of n -simplices $\{\Delta_1, \dots, \Delta_k\}$, called facets, by gluing their $n - 1$ -dimensional faces together. We call this a PL-triangulation if the following conditions are met*

- a) If a face is glued to itself, the face gluing must happen along the identity map.*
- b) For all vertices v , the set of all simplices Δ such that the join $v * \Delta$ is in \mathcal{T} forms a PL $n - 1$ -sphere*

This is a PL-triangulation of a manifold \mathcal{M} if there is a homeomorphism from $\mathcal{T} \rightarrow \mathcal{M}$.

1.2. Knot Theory

Definition 1.13 (Ambient Isotopy). *For a pair of manifold N, M and embeddings g, h of N into M , an ambient isotopy from g to h is a continuous map $F : M \times [0, 1] \rightarrow M$ such that $F_t : M \rightarrow M$ is a homeomorphism, F_0 is the identity, and $F_1 \circ g = h$*

Definition 1.14 (Knot). *A knot is a smooth embedding of S^1 into S^3*

Remark 1.15. We somewhat loosely refer to a knot as both the actual embedding function, and its image in S^3 . However, if we say “a knot K ”, we typically mean the image of the embedding.

Definition 1.16 (Knot Equivalence). *Two knots are said to be equivalent if there is an ambient isotopy between their embedding functions.*

Definition 1.17 (Knot Complement). *For a knot K , the tubular neighbourhood $\nu(K)$ is a small, closed, 3-dimensional region surrounding K that is homeomorphic to the solid torus $S^1 \times D^2$. The knot complement is $X_K = S^3 \setminus \nu(K)$*

We have that the boundary $\partial X_K \cong T^2$ is a torus, as such it is useful to define

Definition 1.18 (Meridian and Longitude). *A meridian (μ) of $\partial X_K \cong T^2$ is a cycle that bounds a disk in the tubular neighbourhood $\nu(K)$, and links with the knot.*

A longitude (λ) of $\partial X_K \cong T^2$ is a cycle in the direction of the knot, and does not link with the knot.

The pair (μ, λ) for a basis for $H_1(\partial X_K, \mathbb{Z}) \cong \mathbb{Z} \oplus \mathbb{Z}$

Definition 1.19 (Infinite Cyclic Cover). *For the knot complement X_K , \tilde{X}_K a cover of X_K is an infinite cyclic cover if $\text{Deck}(\tilde{X}_K/X_K) \cong \mathbb{Z}$*

Every knot has a unique infinite cyclic cover. This follows from the fact that $H_1(X_K; \mathbb{Z}) \cong \mathbb{Z}$ and the properties of covering spaces.

Definition 1.20 (Alexander Module). *For a knot K , take the infinite cyclic cover \tilde{X}_K of the knot complement. The deck transformation group of this space is isomorphic to \mathbb{Z} , generated by some transformation t . Form the ring of Laurent polynomials $\Lambda = \mathbb{Z}[t, t^{-1}]$. The Alexander module is the first homology group $H_1(\tilde{X}_K)$ viewed as a module over Λ .*

Definition 1.21 (Alexander Polynomial). *The Alexander Polynomial is the generator of the first elementary ideal of the Alexander module.*

1.2.1. Isomorphism Signatures. We say that two PL-triangulations are *combinatorially isomorphic* if they are the same up to relabelling of vertices. Because there is no interesting mathematical information in the labelling, we typically define a canonical labelling. While the exact details of constructing a canonical labelling are specific and can depend on the software choice, the general idea is as follows.

- a) For each facet, calculate some label independent invariant (e.g. the number of unique facets being glued to).
- b) Partition the facets into labelled bins based on this invariant (this labelling is canonical based on the value of the invariant)
- c) For each facet in each bin (containing multiple facets), look at which bins the facets neighbouring (from gluing) facets are in, use this information to construct a new label, if there are facets in the same bin that get a different new label, split the bin in some canonical way and create a new labelling.
- d) Repeat the above process, constructing a splitting tree. This tree is canonical. This tree is not guaranteed to split the bins into singletons. If it gets “stuck”, make an arbitrary choice and continue the algorithm until everything falls into a singleton. This is a possible canonical labelling.
- e) if an arbitrary choice had to be made (or perhaps multiple), rerun the entire algorithm for each possible choice that could have been made, creating a small list of possible canonical labelling. Chose the lexicographically smallest labelling.

The details of the initial invariant, how the bins are labelled, and optimisations can be found in [?]. In the worst case scenario this process takes $\mathcal{O}(n! \cdot k^2)$ where n is the dimension and k is the number of facets. For a 3–dimensional triangulation the $n!$ term is small and this algorithm is fast.

Once a canonical labelling is found, an “isomorphism signature” is calculated. For a 3–dimensional triangulation, the only information that needs to be stored is the number of tetrahedra, and for face of each tetrahedron what the destination tetrahedron is, and what the vertex order is. We do not need to specify the destination face because this can be determined by a combination of permutation labelling, looking at where the destinations faces are being mapped, and strict ordering. To efficiently store this, the following is done.

- a) For each face of each tetrahedron, calculate $v_i = 6 \cdot \text{destination id} + \text{permutation id}$. This is valid due to their being 6 vertex permutations. This value will less than $6k$ for k facets.
- b) Create a sequence $(v_0, v_1, \dots, v_{4k-1})$ ordered based on the canonical labelling.
- c) Compute $I = v_0 + v_1(6k) + v_2(6k)^2 + \dots$
- d) Convert I to base 64 and store as a string.
- e) Convert k to some string and append to the front of the string representation of I .

This process form a string / number representing the unique combinatorial isomorphism signature for any given triangulation of a 3-manifold (practically there are size restrictions on the number of facets k).

1.2.2. Computational Complexity.

Definition 1.22 (Time Complexity). *For an algorithm with input size n , the time complexity is the asymptotic worst case runtime of the algorithm, written in “Big-O” notation as $\mathcal{O}(g(n))$.*

Definition 1.23 (Space Complexity). *For an algorithm with input size n , the space complexity is the asymptotic worst case amount of memory the algorithm needs to allocate, written in “Big-O” notation as $\mathcal{O}(g(n))$.*

Definition 1.24 (Decision Problem). *A decision problem is a type of problem that can be answer as either “yes” or “no”*

Remark 1.25. While most problems are not decision problems, it is often possible to convert many problems into a decision problem. For example, if the problem is to multiply two numbers together, we can convert this to a decision problem of “is the i th binary digit a 0?”, this can be repeated for all digits.

Definition 1.26 (Complexity Class). *A complexity class is a class of problems that can be solved under some model of computation within some resource bound i.e. some time or space bound. A problem belongs to a given complexity class if there exist an algorithm to solve the problem within this bound. For “turning machines”, the main complexity classes are*

- a) L : The problem can be solved using a logarithmic amount of space.
- b) NL : A solution can be verified using a logarithmic amount of space.
- c) P : The problem can be solved in polynomial time.
- d) NP : A solution can be verified in polynomial time.
- e) $PSPACE$: A solution can be found using a polynomial amount of space.
- f) $EXPTIME$: A solution can be found using an exponential amount of time.

There exist many other complexity classes for both turning machines and other models for computation.

Remark 1.27. A turning machine is a common model for computation, we do not define it here but there are many resources available.

Remark 1.28. We have defined NL and NP as having solutions verifiable in a specific bound, this means that if the question is for example, is there a path of length k between two nodes, someone can give a path that is a solution, and it is possible to verify that this connects the two nodes and has length k . Formally NL and NP are not defined this way, but instead using things called “non-deterministic Turing machines”, but the definitions can be shown to be equivalent.

Theorem 1.29. $L \subseteq NL \subseteq P \subseteq NP \subseteq PSCPAE \subseteq EXPTIME$

Remark 1.30. There are conjectures that some of these subsets may be proper subsets or equalities, most famously the question of if $P = NP$, but there are not as of yet any solutions.

Definition 1.31 (Hard Problem). *For a given complexity class C , a problem P is called a C –hard problem if every problem $P' \in C$ can be “reduced” to P via some algorithm that is “efficient” relative to C .*

Remark 1.32. We do not formally define what it means to be “efficient” relative to C as it is getting into the details too much, but for NP and $PSCPAE$ problems, a polynomial time reduction is efficient.

Definition 1.33 (Complete Problem). *For a given complexity class C , a problem P is called a C –complete problem if it is a C –hard problem and $P \in C$*

Practically, problems that belong to P can be solved tractable on a modern computer, we can typically scale resources and solve the problem. However, all the algorithms we have for NP complete problems require exponential resources to run, which becomes computationally infeasible quickly. As such, if we find a problem to be NP –complete or harder, it means we can’t just “throw more compute” at the problem to solve it and may have to resort to suboptimal solutions.

1.2.3. Pachner Moves and the Pachner Graph. A Pachner move is a specific process that changes a PL-manifold without changing the underlying topological space.

Definition 1.34 (Complementary Triangulation). *Let Δ a $d + 1$ simplex, let A be a connected sub complex of $\partial\Delta$ - that is a set of d –simplexes that form a connected topological space. The complementary triangulation of A is $B = \partial\Delta \setminus A$.*

From this we can abstractly define a Pachner move

Definition 1.35 (Pachner Move). *For a d -dimensional PL-Manifold, identify an i –dimensional simplex σ^i , the set of all facets that contain σ^i form a sub complex of the boundary of some $d + 1$ dimensional triangulation. Replace this set with the complementary triangulation of this sub complex. This Pachner move is called an i –move.*

Theorem 1.36. *Two triangulations \mathcal{T} and \mathcal{T}' are reachable in finitely many Pachner moves if and only if they represent the same PL-manifold.*

For triangulations of 3–manifolds, we can take $i = 0, 1, 2, 3$, where $i = 0, 3$ are inverses of each other, and likewise $i = 1, 2$. For a 2–move, we identify a triangle and find all the tetrahedra that share that triangle, there are two of these, we replace these with 3 tetrahedra that now share a common edge. The 1–move is the inverse of this. We often call these moves $(3, 2)$ and $(2, 3)$ moves to explicitly state how we are going, for example, from 3 tetrahedra to 2. For a 0–move, we

identify a vertex with 4 tetrahedra that share it, you replace these with a single tetrahedron. The 3-move is the inverse. These are often called $(4, 1)$ and $(1, 4)$ moves as above.

Because all triangulations of a specific PL-manifold can be reached through Pachner moves, we define

Definition 1.37 (Pachner Graph). *For a given PL-manifold \mathcal{M} , the Pachner graph for \mathcal{M} is the graph (V, E) given by $V = \{\text{triangulations of } \mathcal{M}\}$ and for $v_i, v_j \in V$, $e_{ij} = (v_i, v_j) \in E \iff v_i, v_j$ are related via a single Pachner move.*

Because of the above theorem, the Pachner graph is fully connected for any PL-manifold. The Pachner graph is however, infinite. We additionally endow the Pachner graph with “levels” where each level encodes the number of facets in the triangulation.

Theorem 1.38. *The Pachner graph grows super exponentially in its level.*

We also have that

Theorem 1.39. *The shortest path between v_i and v_j two vertices of the Pachner graph is PSPACE hard problem.*

Remark 1.40. It is not currently known if the graph traversal is PSPACE complete

We notice that $(4, 1)$ and $(1, 4)$ moves change the number of vertices, but the $(3, 2)$ and $(2, 3)$ moves do not. We are often interested in “single vertex” triangulations, for these we can exclusively use $(3, 2)$ and $(2, 3)$ moves. A particularly useful result is that

Theorem 1.41. *Any single vertex triangulation of the 3-sphere S^3 can be reached through only $(3, 2)$ and $(2, 3)$ moves.*

We analogously define the Pachner graph for 1-vertex triangulations of S^3 , which is once again a fully connected graph. Unlike the general Pachner graph, it is currently unknown if the Pachner graph for 1-vertex triangulations of S^3 grow super exponentially or exponentially in its level.

1.3. Machine Learning Preliminaries

1.3.1. Markov Chain Monte Carlo.

Definition 1.42 (Markov Chain). *A Markov Chain is a sequence of values $x(t)$ indexed by “time” where the sequence has the Markov Property - that is $x(t + 1)$ depends only on $x(t)$, typically probabilistically.*

Monte Carlo methods are a class of methods that use random sampling to approximate something. Often, they are used to approximate the integral of some complex function, however not exclusively.

Markov Chain Monte Carlo is a class of methods to sample from a distribution by constructing a Markov Chain that’s steady state solution is the distribution. More formally we define

Definition 1.43 (Markov Transition Kernel). Take $(\mathcal{X}, \mathcal{F})$ a measurable space with \mathcal{F} a sigma algebra. A Markov Transition Kernel is a function $P : \mathcal{X} \times \mathcal{F} \rightarrow [0, 1]$ such that

- a) for any $x \in \mathcal{X}$, the function $A \mapsto P(x, A)$ is a probability measure on $(\mathcal{X}, \mathcal{F})$
- b) for any $A \in \mathcal{F}$, the function $x \mapsto P(x, A)$ is a measurable function.

The Markov Transition Kernel is used to “evolve” the state. That is, if you are currently “located” at $x \in \mathcal{X}$, $P(x, A)$ tells you the probability of “ending up” in A . More precisely, if at time t the system is distributed according to ν_t , then the distribution at $t + 1$ is

$$(1.44) \quad \nu_{t+1} = \nu_t P(A) = \int_{\mathcal{X}} P(x, A) d\nu_t(x)$$

Definition 1.45 (Target Distribution). π is called a target distribution for a Markov Transition Kernel P if the following two conditions hold

- a) **Stationarity**: $\pi P = \pi$
- b) **Ergodicity**: For π -almost all points x_0 ,

$$\lim_{n \rightarrow \infty} \|P^n(x_0, \cdot) - \pi(\cdot)\|_{TV} = 0$$

Remark 1.46. A statement is true for “ π -almost all x_0 ” means that it is true for all x_0 except a set that has measure 0 under π

Typically proving stationarity and ergodicity is difficult, so the following is used instead

Theorem 1.47. For a Markov transition kernel P , if the following are true for some distribution π

- a) π -**irreducibility**: For any $A \in \mathcal{F}$ with $\pi(A) > 0$, there exists some n such that $P^n(x, A) > 0$ for all $x \in \mathcal{X}$
- b) **Aperiodicity**: The chain is aperiodic
- c) **Recurrent**: For any measurable set $A \in \mathcal{F}$, and any starting point $x \in \mathcal{X}$ the hitting time τ_A is finite almost surely, where τ is the amount of steps to go from x to A .
- d) **Reversibility**: For any two $A, B \in \mathcal{F}$ we have

$$\int_A P(x, B) d\pi(x) = \int_B P(x, A) d\pi(x)$$

then π is the target distribution for the Markov transition kernel P .

In particular, if we have some process \hat{P} that takes in some $x(t)$ and produces some $x(t + 1) = \hat{P}(x(t))$ according to some Markov transition kernel P , we can create a Markov Chain, by continuously generating samples, and applying \hat{P} to them recurrently, these samples will be “distributed according to” π . More precisely, the partial sums $S_n = \frac{1}{n} \sum_{i=1}^n f(x_i) \rightarrow \mathbb{E}^\pi[f]$ at a “rate” $\mathcal{O}(\sqrt{N})$.

Definition 1.48 (Metropolis Hastings). For some function f , not necessarily a probability density, and some “proposal function” g the Metropolis Hastings algorithm on f with proposal g is the following: Initialise some x_0 , for each t do the following

- a) Propose some x' according to $g(x'|x_t)$
- b) Compute $\alpha = f(x')/f(x_t)$
- c) Sample some $u \in [0, 1]$ uniformly
- d) If $\alpha > u$ accept x' by setting $x_{t+1} = x'$, otherwise set $x_{t+1} = x_t$

Theorem 1.49. *If the proposal function g is symmetric - that is $g(x|y) = g(y|x)$, and if f is proportional to some probability distribution π , then the Metropolis Hastings algorithm for f with proposal g will generate samples according to π .*

Remark 1.50. If the proposal distribution is not symmetric, a factor known as the Hastings ratio can be introduced to find $\alpha_H = \alpha \cdot \frac{g(x'|x)}{g(x|x')}$ which is used in the acceptance step.

1.3.1.1. *Statistics in MCMC.* While the Metropolis Hastings algorithm (and other MCMC algorithms) converge to π as $n \rightarrow \infty$, for finite n the distribution will not necessarily be π , a common example of this is if π is bimodal, where sampling can get “stuck” in one of the modes of the distribution, thereby not fully exploring the space. As such, there are a number of statistics that are used to check if a finite sample from an MCMC algorithm is likely to have converged to the proper distribution π . This is typically done by running multiple independent chains of MCMC.

Definition 1.51 (Variance). *Let there be J chains of L samples, with samples x_1^j, \dots, x_L^j for the j th chain. We have the*

- a) *Between chain variance:* $B = \frac{L}{J-1} \sum (\bar{x}_j - \bar{x}_*)^2$
- b) *Within chain variance:* $W = \frac{1}{J} \sum \left(\frac{1}{L-1} \sum (x_i^j - \bar{x}_j)^2 \right)$

The between chain variance follows from the central limit theorem that $\text{Var}(\bar{x}_j) \approx \sigma^2 \cdot n$, so $B = \sigma^2 = L \cdot \text{Var}(\bar{x}_j)$, and $\frac{1}{J-1} \sum (\bar{x}_j - \bar{x}_*)^2$ is an unbiased estimate of the variance. The within chain variance is simply the average unbiased estimate of the variance of each chain. From these we can derive

Theorem 1.52 (Gelman-Rubin Variance Estimate). *Let there be J chains of L samples, with samples x_1^j, \dots, x_L^j for the j th chain. We have that the total variance is*

$$\hat{V} = \frac{L}{L-1} W + \frac{1}{L} B$$

This result follows directly from the law of total variance.

From this we can construct

Definition 1.53 (Gelman-Rubin Statistic). *Let \hat{V} be the Gelman-Rubin variance estimate, and W be the within chain variance. The Gelman-Rubin statistic is*

$$\hat{R} = \sqrt{\frac{\hat{V}}{W}}$$

When the chains have converged, each chain is a valid sample from π , as such $W = \hat{V}$ and $\hat{R} = 1$. However, if one or more of the chains has not suitably explored the space, $W < \hat{V}$ and $\hat{R} > 1$. The samples will never be a complete representation of the distribution, so typically a threshold of $\hat{R} = 1.01$ is used to indicate convergence.

1.3.2. Simulated Annealing. For a measurable space $(\mathcal{X}, \mathcal{F}, \mu)$ with base measure μ (typically the Lebesgue, or count measure), we may have some measurable function $E : \mathcal{X} \rightarrow \mathbb{R}$ known as the energy. In thermal physics, for a thermodynamic system at temperature T , the probability of finding a particle in a given energy E is given by the Gibbs distribution $\pi_T(A) = \frac{1}{Z} \int_A e^{-E(A)/T} d\mu(x)$ where $Z = \int_{\mathcal{X}} e^{-E(A)/T} d\mu(x)$ is a normalising factor. Typically as a system “cools down” it “finds its way” to low energies. Simulated annealing is inspired by this thermodynamic concept, simulating a particle over time as the temperature $T \rightarrow 0$ heuristically conjecturing that it will find its way to the lowest energy state. The simulation at a given temperature T is typically done using the Metropolis Hastings algorithm.

While the concept above is inspired as a Heuristic from concepts in thermal physics, it has valid grounding in probability and machine learning. At a high level, this is because for any $x, x' \in \mathcal{X}$, $\pi_T(x')/\pi_T(x) = e^{-(E(x')-E(x))/T}$, which if $E(x') > E(x)$ tends to 0 as $T \rightarrow 0$, but if $E(x') < E(x)$ tends to ∞ . As such, a sharp delta function is formed around $x_{min} = \arg \min(E)$. While this does mean that at $T = 0$ we have that $\pi(x) = \delta(x - x_{min})$ is the stationary solution to our Metropolis Hastings Markov Transition Kernel, at $T = 0$ the measure of valid starting points $\{x_0\}$ tends to have measure 1 under π but measure 0 under ν . Informally, at low T we almost always accept points where $E(x') < E(x)$ and almost never accept points where $E(x') > E(x)$, so this is just hill climbing and unless E is convex, the sampling will likely get stuck in a local minima.

1.3.3. Transformers.

Definition 1.54 (Row Operator). *For a field F , and a function $s : F^n \rightarrow F^n$ define the row operator $\mathcal{R}_s : M_n(F) \rightarrow M_n(F)$ as the act of applying s to each row of $M_n(F)$. That is*

$$\mathcal{R}_s \left((r_1, r_2, \dots, r_n)^T \right) = (s(r_1), s(r_2), \dots, s(r_n))^T$$

Definition 1.55 (Abstract Transformer Head). *For a collection of vectors (v_1, v_2, \dots, v_n) where $v_i \in V$ an vector space over a field F . Form \dot{V} by equipping V with a bilinear form denoted $B : V \times V \rightarrow F$. Let $G : F^n \otimes V \rightarrow M_n(F)$ denote the “Gram operator”. Let $L : V \rightarrow W$ be an arbitrary linear operator, and $s : F^n \rightarrow F^n$ an arbitrary function.*

An abstract transformer head is a function $H : F^n \otimes V \rightarrow F^n \otimes W$ defined by

$$(1.56) \quad H(x) = (\mathcal{R}_f(G(x)) \otimes L)(x)$$

Definition 1.57 (Abstract Multi-Head Attention). *For a collection of vectors (v_1, v_2, \dots, v_n) where $v_i \in V$ an vector space over a field F . Form spaces V_1, V_2, \dots, V_q from V by equipping V with a, possibly unique, inner product. Let $\iota_i : V \rightarrow V_i$ be the operation of endowing V with the inner*

product structure of V_i . Let H_i denote the transformer head from $F^n \otimes V_i \rightarrow F^n \otimes W_i$ where W_i, W_j may be distinct. Define $O : W_0 \oplus W_1 \oplus \dots \oplus W_q \rightarrow V$

The multi-head attention is an operation $A : F^n \otimes V \rightarrow F^n \otimes V$ defined as

$$(1.58) \quad A(x) = (I \otimes O) \left(\bigoplus_i H_i(\iota_i(x)) \right)$$

Definition 1.59 (Abstract Feedforward Neural Network). For a vector space U over a field F , a Feedforward is a operator $N : U \rightarrow W$ for a pair of linear operators $L_1 : U \rightarrow V$ and $L_2 : V \rightarrow W$ and a non-linear function $a : V \rightarrow V$ is $N(x) = L_2(a(L_1(x)))$

Definition 1.60 (Abstract Transformer Block). For a collection of vectors (v_1, v_2, \dots, v_n) where $v_i \in V$ an vector space over a field F . A transformer block $T : F^n \otimes V \rightarrow F^n \otimes V$ is an operator defined as

$$T(x) = x + A(x) + (I \otimes N)(x + A(x))$$

While these abstract formulations of the Transformer are theocratically rich, in practice we are working over $F = \mathbb{R}$ and our vector space V is finite dimensional. This allows us to write $V = \mathbb{R}^\kappa$, write all the linear transformations as matrix multiplications, and write the bilinear forms as $B(v_1, v_2) = v_2^T M V_1$ for some matrix M . Also, the non-linear functions used such as f in the abstract transformer head, and a in the feedforward neural network are fixed, or depend on some small set of learnable parameters. This gives a finite set of scaler parameters θ that defines the entire transformer block. Additionally, some further components are introduced into the transformer block to ensure computationally stability and regularisation. This lets us define

Definition 1.61 (Standard Transformer Head). For a collection of vectors (v_1, v_2, \dots, v_n) where $v_i \in \mathbb{R}^d$ we construct $X \in \mathbb{R}^{n \times d}$ by stacking these vectors up as columns. The head is calculated as

$$(1.62) \quad h(X) = \text{softmax} \left(\frac{(XW^Q)(XW^K)}{\sqrt{d_k}} + M \right) (ZW^V)$$

Where the softmax is applied along each row. The $W^Q, W^K, W^V \in \mathbb{R}^{n \times d_k}$ are called the “query”, “key”, and “value” weights. $M \in M_n(\mathbb{R} \cup \{-\infty, \infty\})$ is called the causal mask.

Definition 1.63 (Standard Multi-Head Attention). For X defined as above

$$(1.64) \quad MHA(X) = \text{Concat}(h_1(X), \dots, h_h(X))W^O$$

Definition 1.65 (Standard Feedforward Neural Network). For $x \in \mathbb{R}^\kappa$, $FFN(x) = \sigma(xW_1 + b_1)W_2 + b_2$ where the W_i are called weights and the b_i are called biases. σ is a parameter-less “activation function”.

Definition 1.66 (Layer Norm). *For a vector $x \in \mathbb{R}^d$ define μ and σ as the mean standard deviation of x . The layer norm of x is*

$$(1.67) \quad LN(x)_i = \gamma_i \cdot \left(\frac{x_i - \mu}{\sqrt{\sigma^2 + \epsilon}} \right) + \beta_i$$

Where $\gamma, \beta \in \mathbb{R}^n$ are learned parameters, and ϵ is some fixed number for numerical stability.

Remark 1.68. This was not included in the abstract definition of the transformer block as it is primarily used for numerical stability.

Definition 1.69 (Standard Transformer Block). *For X as defined above*

$$(1.70) \quad Y = X + MHA(LN(X)) + FFN(LN(X + MHA(LN(X))))$$

Where Y is the result from the transformer block.

Because these transformer blocks map from $F^n \otimes V \rightarrow F^n \otimes V$, or from $\mathbb{R}^{n \times d} \rightarrow \mathbb{R}^{n \times d}$, they can be chained together.

For many tasks which transformers excel in (language modelling, time series forecasting) the input sequence is sequential, and often not encoded into a vector in a meaningful way. As such we have

Definition 1.71 (Token Embedding). *For an input vectors $v \in \mathbb{R}^V$, the token embedding of v is $h = x^T E$ where $E \in \mathbb{R}^{V \times d}$.*

Remark 1.72. For tasks where the input is categorical, v is usually formed with some one-hot-encoding, that is if there are V categories we define for category c $v_c = (0, \dots, 0, 1, 0, \dots, 0)$ as a vector in \mathbb{R}^n that has a 1 in the position corresponding to category c (arbitrary), and zeros elsewhere.

Often (e.g. language modelling) there are many categories, so $d < V$ and the embedding becomes a useful technique for encoding “semantic” meaning into a vector.

Once we have embedded our vectors, there is still often a sequential order to them, i.e (v_0, v_1, v_2, \dots) . The transformer block as described above, is invariant under changing of this order. As such, we use

Definition 1.73 (Positional Encoding). *For a sequence of vectors indexed by time, (v_0, v_1, \dots, v_n) each $v_n \in \mathbb{R}^d$. For T_{max} the longest possible sequence, and a positional embedding matrix $P \in \mathbb{R}^{T_{max} \times d}$, the positional embedding of (v_0, v_1, \dots, v_n) is defined with $v_i \mapsto v_i + P_i$ where P_i is the i th row.*

A prototypical example of applying these steps is in GPT-2, described as follows. For a sequence of words (w_1, w_2, \dots, w_n) , each word in a vocabulary of size V , encode each word via one

hot encoding. Perform token embedding into \mathbb{R}^d , and perform a positional encoding. Apply a series of transformer blocks sequentially, apply a single layer norm, apply a linear projection back into \mathbb{R}^V for each “word” (via some matrix multiplication), and finally apply softmax to each output “word”. The final vector is trained to represent the probability of the next word in the sequence of words. So for position 1 the prediction is for w_2 , for position n the prediction is for word w_{n+1} not in the data. To train GPT-2 we use an objective function that calculates the cross entropy between the prediction at each step and the actual next word, this is done for each “word” in the output layer.

1.3.4. Gradient Descent. If we have a neural network $N(x \mid \theta)$ that takes in some input $x \in V$ and some parameter set $\theta \in \mathbb{R}^n$ and produces some output $y \in W$ where V, W are vector spaces over a field F , typically \mathbb{R} , and typically finite dimensional. We typically have some “training data”, which is a set of $\{(x_i, y_i) \mid x_i \in V, y_i \in W\}$ where each y_i is the output that we “want” from $N(x_i, \theta)$, what we mean by “want” is that some function $\mathbb{E}_{x_i}[\mathcal{L}(N(x_i, \hat{\theta}), y_i)]$ is minimised by our parameter choice $\hat{\theta}$. To find this target $\hat{\theta}$. Because N is typically complicated, and has a complex dependence on θ , we cannot directly minimise this objective function to find $\hat{\theta}$, as such we use gradient decent. The concept behind gradient decent is to take some θ_i , find a local, linear approximation of $\mathbb{E}_{x_i}[\mathcal{L}(N(x_i, \hat{\theta}_i), y_i)]$, and find some “penalty” for “trusting” the approximation too much. Precisely, we compute

$$(1.74) \quad \theta_{i+1} = \arg \min_{\theta} \left(\langle g^T, (\theta - \theta_i) \rangle + \frac{\lambda}{2} \|\theta - \theta_i\|^2 \right)$$

This is done recursively, typically for some number of steps. Here g is the gradient, often we would consider using $\nabla_{\theta} \mathbb{E}_{x_i}[\mathcal{L}(N(x_i, \hat{\theta}_i), y_i)]$, but computing this over the entire “training” data set is slow, so we typically a small “batch” of data is taken from the training data and the gradient is taken over this batch. Because this batch is “random” a small subset of the true data, our estimate of the gradient is going to be somewhat stochastic, as such many gradient decent algorithms use some sort of exponential smoothing step for the gradient - this is often called the momentum. Different choices of exponential smoothing, and different choices of the norm, and different choices of the gradient lead to different types of optimisers. A few common ones are *SGD* which uses the frobenius norm, the standard gradient, and no smoothing; *adam*, which uses the $\ell_1 \rightarrow \ell_{\infty}$ norm, the standard gradient, and exponential smoothing; *shampoo*, which uses the spectral norm, the standard gradient, and exponential smoothing; and *NGF* which uses the frobenius norm, the natural gradient, and no exponential smoothing.

1.3.5. Reinforcement Learning. For our specific usecase, reinforcement learning is a technique used when we have some “reward function” $R(y) : Y \rightarrow \mathbb{R}$ that tells us how good a specific sample is, and we want to update the parameters of some probability distribution / sampling function $Y \sim \pi_{\theta}$ to maximise the expected reward, that is

$$(1.75) \quad J(\theta) = \mathbb{E}_{y \sim \pi_\theta} [R(y)]$$

While reinforcement learning is typically more general than this, and applies to policy learning, this will suffice for our purpose.

The most direct way to update π_θ is to simply take a number of samples from π_θ and calculate

$$(1.76) \quad g = \nabla_\theta J(\theta) = \mathbb{E}_{y \sim \pi_\theta} [\nabla_\theta \ln(\pi_\theta) R(y)]$$

and then use this g in gradient decent as described above.

While the above techniques makes logical sense, the use of $J(\theta)$ being the expected reward can become problematic, for example if our reward function is binary, returning a score of 1 if the output meets some criteria, and returning a score of -1 if it does not, once the parameter set has settled into a stable configuration where each output is valid and produces a score of 1, $J(\theta)$ continues to update θ which can lead to oscillations. As such we often use an “advantage” function A which is an estimate of how the score $R(y_i)$ for a specific sample compares to the expected score. That is $A(y_i) = R(y_i) - \mathbb{E}_{y \sim \pi_\theta} [R(y)]$. We can estimate the expected score by simply sampling from π_θ a few times, and taking the average, and for a simple problem this is sufficient.

Remark 1.77. This average often has a high variance, so a “critic” model is often trained instead to try and learn the baseline. We won’t explore this technique here.

We can use the advantage to get a refined objective function $J(\theta) = \mathbb{E}_{y \sim \pi_\theta} [A(y)]$. We can calculate the gradient once again as above, and apply gradient accent. The primary disadvantage here is that each time we take a gradient step, we have to resample from π_θ to get a new batch to calculate the new gradient. The new distribution with parameter set θ will only be slightly different to the original parameter set before the gradient update θ_{old} , and as such we can keep the old samples $y_i \sim \pi_{\theta_{old}}$ and then use a technique called importance weighting to find that

$$(1.78) \quad \mathbb{E}_{y \sim \pi_\theta} [A(y)] = \mathbb{E}_{y \sim \pi_{\theta_{old}}} [r(\theta) A(y)]$$

Where $r(\theta) = \frac{\pi_\theta(y)}{\pi_{\theta_{old}}(y)}$ is called the importance ratio. This allows us to take batch of samples from $\pi_{\theta_{old}}$, perform a few steps of gradient accent using this batch, importance reweighing at each step. While this works in expectations, typically after a few gradient steps our samples are no longer representative, and our reweighing erodes. Therefore, it is typical to use a “trust region” wherein we can trust an update. One of the most common ways to do this is to “cap” positive updates by restricting $r(\theta) \leq 1 + \epsilon$ if $A(y) > 0$, that is if we generated good samples from θ_{old} we want to increase the likelihood of generating them, but not so much so that we can no longer trust the importance sampling as the ratio $r(\theta)$ has become too big, but on the other hand, if the advantage $A < 0$, we want to continue to have a learning signal. This yields the proximal policy optimisation

(PPO) objective

$$(1.79) \quad L^{CLIP}(\theta) = \mathbb{E}_{y \sim \pi_{\theta_{old}}} [\min(r(\theta)A, \text{clip}(r, 1 - \epsilon, 1 + \epsilon)A)]$$

Search and Sampling Algorithms

We are often interested in exploring the characteristics and properties of triangulations of a manifold \mathcal{M} . Due to the unwieldy size of $\mathcal{T}(\mathcal{M})$, it is typical to sample from $\mathcal{T}(\mathcal{M})$ and study these samples. This is the task of a sampling algorithm. There are particular situations in which we may want to find a triangulation that has a particularly “high score”. This is a search task, with an “objective function” $O : \mathcal{T}(\mathcal{M}) \rightarrow \mathbb{R}$ where we aim to find a triangulation $\Delta \in \mathcal{T}(\mathcal{M})$ such that $O(\Delta)$ is “high” (we use high somewhat informally, as we do not know much of the geometry of different objective functions over $\mathcal{T}(\mathcal{M})$).

2.1. Classical Techniques

2.1.1. Direct Accent. For an objective function $O : \mathcal{T}(\mathcal{M}) \rightarrow \mathbb{R}$, for small triangulations (number of tetrahedra around 6 or lower), we can typically calculate O exhaustively from a census. In doing this, we may find that O tends to increase with more tetrahedra, this leads to the notion of direct accent. The basic idea is to start at some small triangulation, enumerate all of its neighbours that have more tetrahedra, and then sample one of these according to its score. This process can then be repeated to generate a chain of tetrahedra of increasing size, and the whole process repeated to generate multiple chains. The detailed algorithm is described in algorithm 1

Algorithm 1 Direct Accent

```

Let  $\Delta_0$  be the initial triangulation.
Let  $\beta > 0$  be a fixed parameter.
let  $T \in \mathbb{N}$  be a fixed chain size.
for  $t = 1$  to  $T$  do
    Generate all neighbours  $\Delta_{ti}$  of  $\Delta_{t-1}$  by applying  $(2 - 3)$  or  $(1 - 4)$  moves.
    for each neighbours  $\Delta_{ti}$  do
        Calculate the percentage advantage  $A_i = \frac{O(\Delta_{ti}) - O(\Delta_{t-1})}{O(\Delta_{t-1})}$ .
    end for
    for each neighbours  $\Delta_{ti}$  in the set do
        Calculate the selection probability  $p_i = \frac{e^{\beta A_i}}{\sum_j e^{\beta A_j}}$ .
    end for
    Sample  $\Delta_t$  from the set of neighbours  $\{\Delta_{ti}\}$  with probability  $p_i$ .
end for
Return  $\Delta_T$ 

```

2.1.2. Markov Chain Monte Carlo.

2.1.3. Simulated Annealing. Because the above Markov Chain Monte Carlo algorithm leads to a uniform distribution, we can use it as a symmetric proposal distribution (**need to check this**). We implement a standard simulated aneling algorithm as described in algorithm 2

Algorithm 2 Simulated Annealing

```

Let  $\Delta_0$  be some starting triangulation.
Let  $T$  be the chain length.
Let  $M$  be some hash table memory.
Let  $O : \Delta \rightarrow \mathbb{R}$  be some objective function.
Let  $T_t$  be some temperature schedule.
for  $t = 1$  to  $T$  do
  Propose  $\tilde{\Delta}_t$  from  $\Delta_{t-1}$  by using 1-step of the described MCMC algorithm.
  Retrieve  $o_{t-1} = O(\Delta_{t-1})$  from  $M$ 
  Check if  $o_t = O(\tilde{\Delta}_t)$  is in  $M$ , if it is - retrieve it, if not - compute it and store it in  $M$ .
  Compute  $\alpha = e^{-(o_{t-1}-o_t)/T_t}$ 
  Generate some  $p \sim \mathcal{U}([0, 1])$ 
  if  $p \leq \alpha$  then
    Accept  $\Delta_t = \tilde{\Delta}_t$ 
  else
    Let  $\Delta_t = \Delta_{t-1}$ 
  end if
end for

```

2.2. Transformers

2.2.1. Base Transformer. We use a standard GPT-2 style transformer with dropout and pre-layer normalisation. The final transformer model takes the sequence, performs token and positional embedding, dropout, and then n_{layers} of the standard transformer block as described in definition ???. Between each transformer block there is a dropout of 0.1 applied. This architecture is the standard GPT-2 style architecture. We use an embedding dimension of 64, and the feed forward neural networks are also all of dimension 64. The vocabulary size includes all the letters of the training data, plus three special tokens “[BOS]”, “[EOS]”, and “[PAD]” for the beginning of sentence, end of sentence, and padding respectively. For training we use AdamW, with a cross entropy objective function, and a batch size of 32.

2.2.2. Reinforcement on the Topology. The transformer may not generate sequences which represent our target manifold, or any manifold at all. As such after training is complete, reinforcement learning can be applied with a score function which returns 1 if the triangulation is of the correct manifold type, and 0 if it is not. A softer scoring function can be used that returns intermediate values depending on if the isomorphism signature corresponds to an invalid structure, a valid manifold of the wrong type, or a valid manifold of the correct type.

2.2.3. Reinforcement on the Objective. Once the transformer is generating samples with a relatively high chance of being valid isomorphism signatures, if we have a specific objective function $O : \Delta \rightarrow \mathbb{R}$ which we seek to maximise, we can use this score function in the standard reinforcement learning algorithms described in ??.

Numerical Experiments

3.1. Objective Functions

To test the effectiveness of our optimisation strategies on triangulations of manifolds, we consider 5 specific objective functions. These objective functions are all considered on single vertex triangulations of S^3 , call the set of all of these $\mathcal{T}_1(S^3)$.

3.1.1. Alexander Polynomial. For any $T \in \mathcal{T}_1(S^3)$, we can consider an edge e_α , because the triangulation has only one vertex, e_α forms a loop. This loop can be knotted in S^3 . Call the associated knot for e_α K_α . The Alexander polynomial $\Delta_{K_\alpha}(t)$ has coefficient vector $\vec{a} = [a_0, a_1, \dots, a_n]$. It is known that “more complex” Alexander polynomials lead to “more complex” knots. This notion suggests 3 objective functions of interest.

- a) $\mathcal{O}_{deg} : T \mapsto \sum_{e_\alpha} \deg(\Delta_{K_\alpha}(t))$ where because the Alexander polynomial is invariant under multiplication by $\pm t$ we use the difference between the highest and lowest power of t .
- b) $\mathcal{O}_{det} : T \mapsto \sum_{e_\alpha} |\Delta_{K_\alpha}(-1)|$.
- c) $\mathcal{O}_{norm} : T \mapsto \sum_{e_\alpha} \|\vec{a}\|^2$ where \vec{a} is the coefficient vector of the Alexander polynomial associated with the knot formed by the edge e_α in S^3

The first of these objective functions, \mathcal{O}_{deg} is motivated by the fact that $2g(K) \geq \deg(\Delta_K(t))$. The second objective function, \mathcal{O}_{det} is motivated by the connection between the determinant $|\Delta_K(-1)|$ and the knot colouring. The third objective function, \mathcal{O}_{norm} is a standard objective function on vectors.

3.1.2. Fundamental Group of T/e_α . As discussed, for $T \in \mathcal{T}_1(S^3)$, any edge e_α can form a knot. Typically in knot theory we are interested in the study of the knot complement $S^3 \setminus K$, however there is a connection between this and the resultant manifold after pinching e_α , that is forming $M_\alpha = T/e_\alpha$. As such, there is interest in studying the fundamental group $\pi_1(M_\alpha)$. This inspires our fourth objective function

- d) $\mathcal{O}_{gen} : T \mapsto \sum_{e_\alpha} \#_{gen}(\pi_1(M_\alpha))$ where $\#_{gen}$ is the number of generators of the presentation of the fundamental group.

While we would ideally like to use the reduced presentation of the fundamental group, this is hard* to calculate, so we use the standard presentation returned by ..., as a heuristic, and then perform the full computationally expensive reduction once the results are found.

3.1.3. Edge Degree Variance. The prior four objective functions are all inspired by the knot structure of our triangulations. While this is of particular interest, we also would like to consider the complexity of the triangulation itself. This inspires us to consider the degree of the edges in our manifold, where the degree of an edge is the number of tetrahedra that share this edge.

e) $\mathcal{O}_{var} : T \mapsto \text{Var}_{e_\alpha}(\deg(e_\alpha))$ where we have taken the variance.

The choice of variance here instead of the average follows from the fact that the average degree becomes independent of the triangulation structure and instead depends on the topology which is constant for a sphere.

3.2. Markov Chain Monte Carlo

To establish a baseline for our objective functions, generic Markov Chain Monte Carlo was run to determine the distribution of our objective function across the Pachner Graph. The MCMC algorithm described by Jonathan samples uniformly around triangulations of a particular number of tetrahedra. We, somewhat arbitrary, choose to examine the Pachner graph around triangulations of 30-tetrahedra, which corresponds to $\gamma = 1/10$. This choice is informed by other results in computational topology that have found interesting triangulations of size less than 30 tetrahedra, and because 30 tetrahedra is intractable to enumerate exhaustively. Samples with 7 chains of 10,000 iterations with a step size of 100 were performed, for each sample all objective functions were calculated (as opposed to running distinct sampling processes for each objective function). This has the benefit of not only saving computational time, but giving us a uniform sample of the space where we can assess the correlation of each objective function. We confirmed convergence with the Gelman-Rubin statistic with a threshold of 1.01, this threshold is in accordance with modern recommendations. The convergence statistics are in table 1. For interest, figure 1 shows the correlation between each objective functions across the entire set of MCMC samples. It is important to note that this correlation tells us the correlation for “an average” sample, and the correlation between the functions could dramatically change in very specific parts of the distribution - in particular near the respective maxima of each variable.

We extracted the triangulations with 30 tetrahedra, there were 8,250 such triangulations so an efficiency of $\approx 12\%$. The following analyses are performed on this subset of triangulations.

Metric	\hat{R}
\mathcal{O}_{norm}	1.003
\mathcal{O}_{deg}	1.004
\mathcal{O}_{det}	1.005
\mathcal{O}_{var}	1.006
\mathcal{O}_{gen}	1.009

TABLE 1. Convergence statistics of MCMC Runs

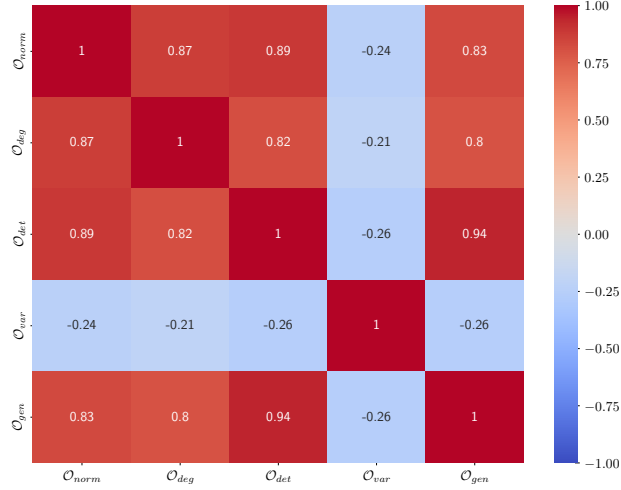


FIGURE 1. Correlation plot for each objective function.

The distribution for the score for each objective function outlined in section 3.1 is displayed in figure 2.

All but \mathcal{O}_{var} present power law distributions or power law distributed tails. To confirm this we fit both a power law distribution across the entire data, and also compute x_{\min} according to the procedure of Clauset, Shalizi, and Newman, and fit a power law distribution to the tails. Both these fits are done with MLE. A plot on a log-log histogram of the data, and the fitted distributions is seen in figure 3.

The \mathcal{O}_{var} appears to have some positively skewed distribution, considering that it is the distribution of the variance, a chi squared distribution would be expected. However, fitting one to the data lead to poor results and a log normal distribution was used instead. This is weakly justified, asside from a strong in sample fit as seen in figure 4 that is plotted in log-log space to confirm tail behaviour.

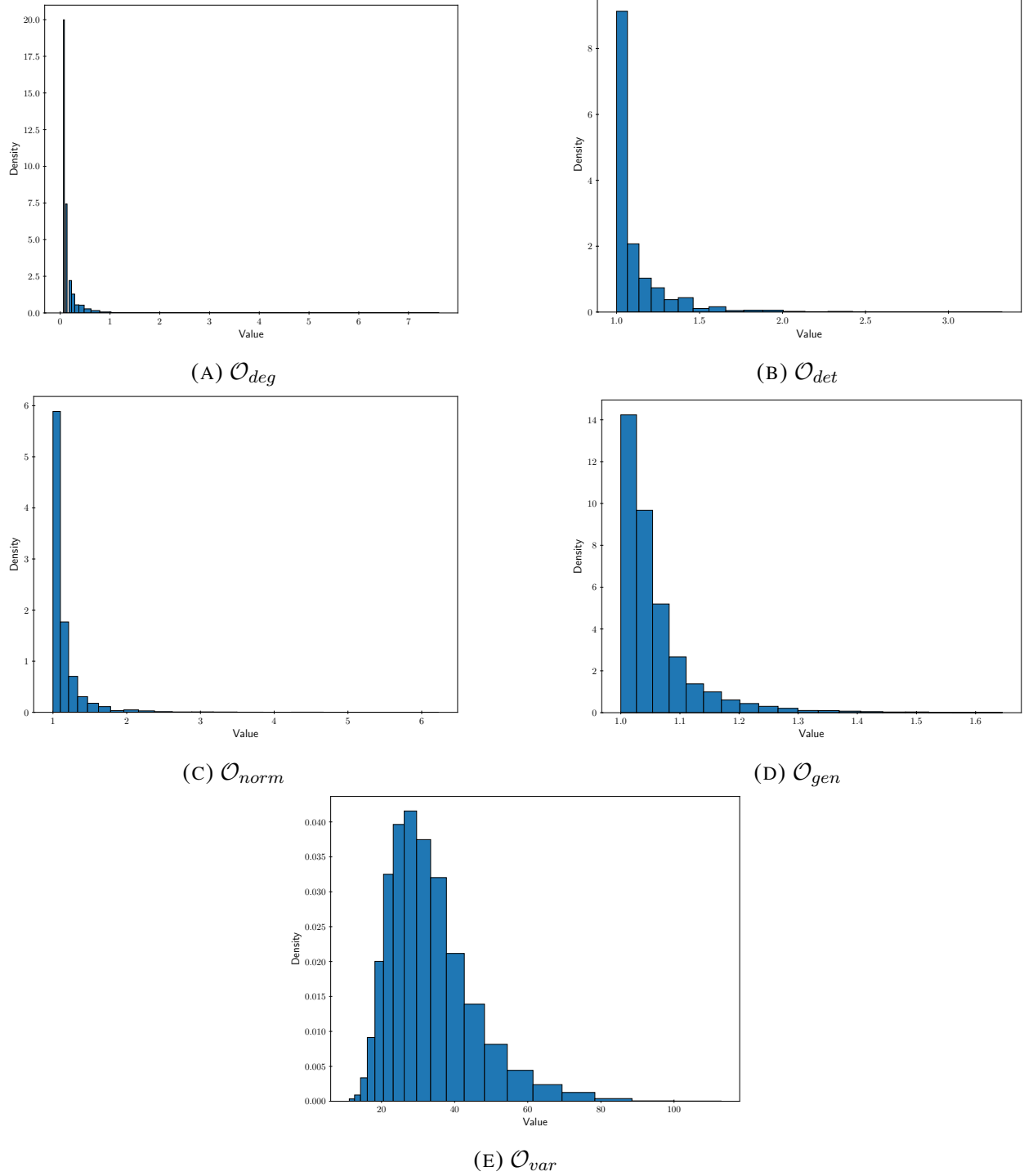


FIGURE 2. Distribution of the score for each objective function under MCMC samples at $\gamma = 1/10$

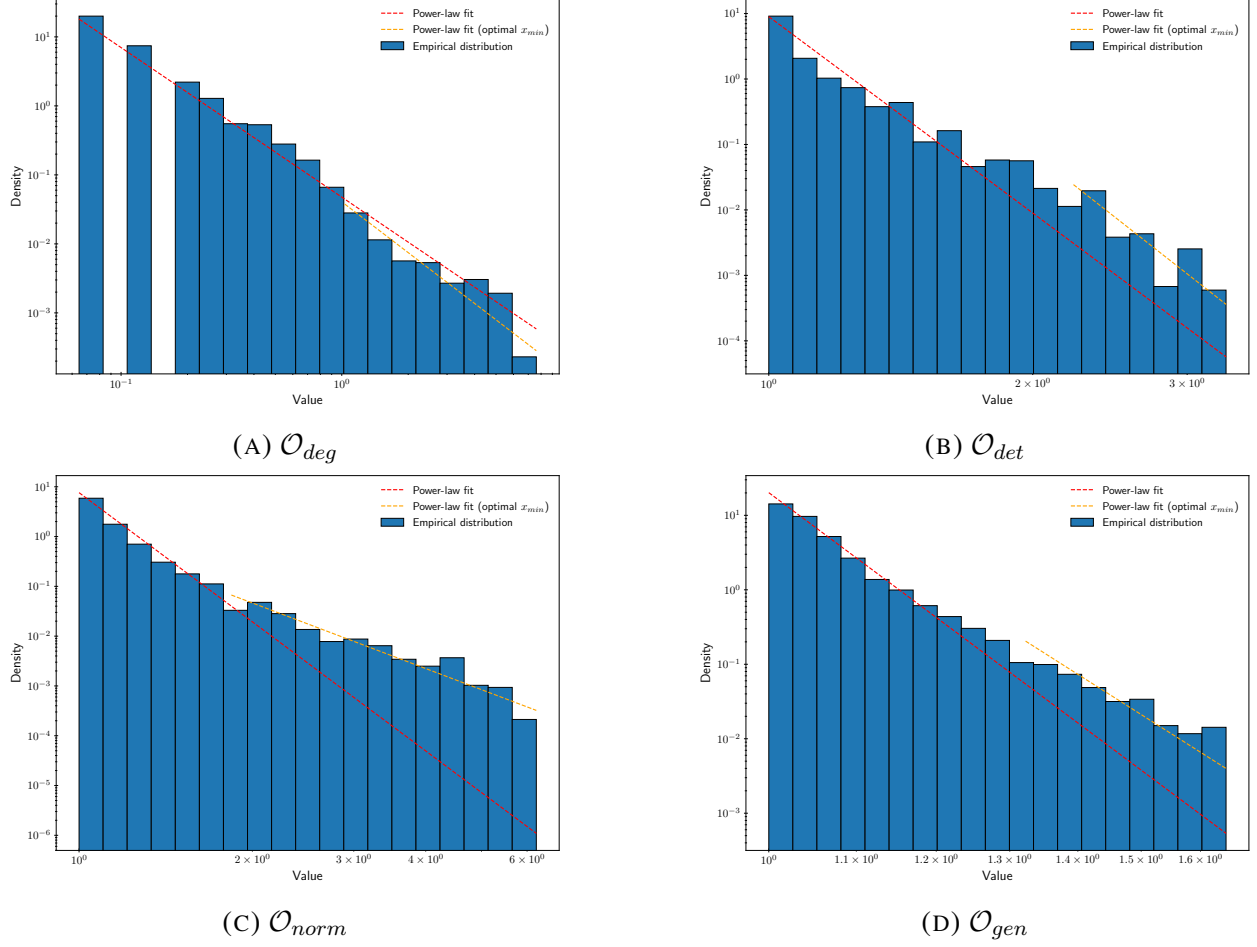


FIGURE 3. Log-log histogram of the score for each objective function under MCMC samples at $\gamma = 1/10$. The red line indicates the fitted power law distribution over the entire dataset, the orange line the fitted power law distribution for $x > x_{\min}$

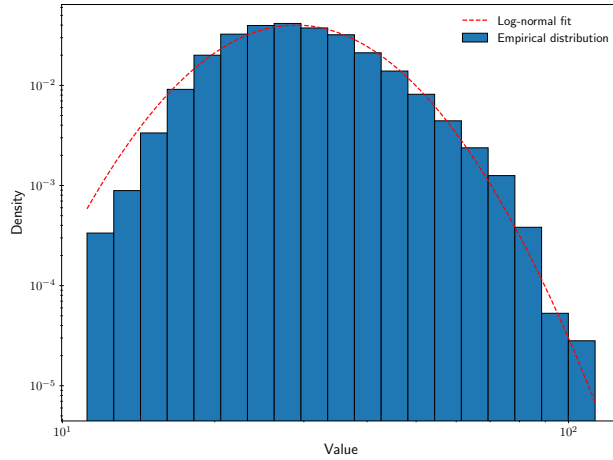


FIGURE 4. Log-log histogram of \mathcal{O}_{var} with fitted log normal distribution.

3.3. Classical Optimisation

3.3.1. Direct Ascent. We perform direct ascent starting from the seed triangulation of `cMcabbgqs` which has one vertex and two tetrahedra. We performed this direct ascent 20 times at different temperatures $(0.1, \sqrt{0.1}, 1, \sqrt{10}, 10)$ and compared the maximum achieved score at each level, an example of this for \mathcal{O}_{deg} in figure 5. In all cases it was found that a lower temperature lead to better results, this suggests that the graph is connected enough that greedy ascent is sufficient (that is, the neighbour with the greatest score is accepted at each step). We perform greedy ascent for each objective function for 28 steps to a triangulation of 30 tetrahedra, the chains are presented in figure 6. We see that for the ascent profile is rather deterministic, suggesting ... However, for ... we observe a more complex objective function geometry.

While the direct ascent approach does perform particularly well, it is computational slow, requiring, and can only find a single triangulation of interest, i.e. its terminal triangulation after ascent. As such, in situations where a single high value triangulation is desired this is a beneficial approach. However, if the goal is to explore a variety of different high valued triangulations, the effectiveness decreases. Additionally, even in the higher temperature case, there is a particularly low efficiency. For triangulations of size T tetrahedra, there are $2T$ neighbours of size $T + 1$, of which only one is selected at each step, so the sampling efficiency is $1/2T$ which is about $< 2\%$.

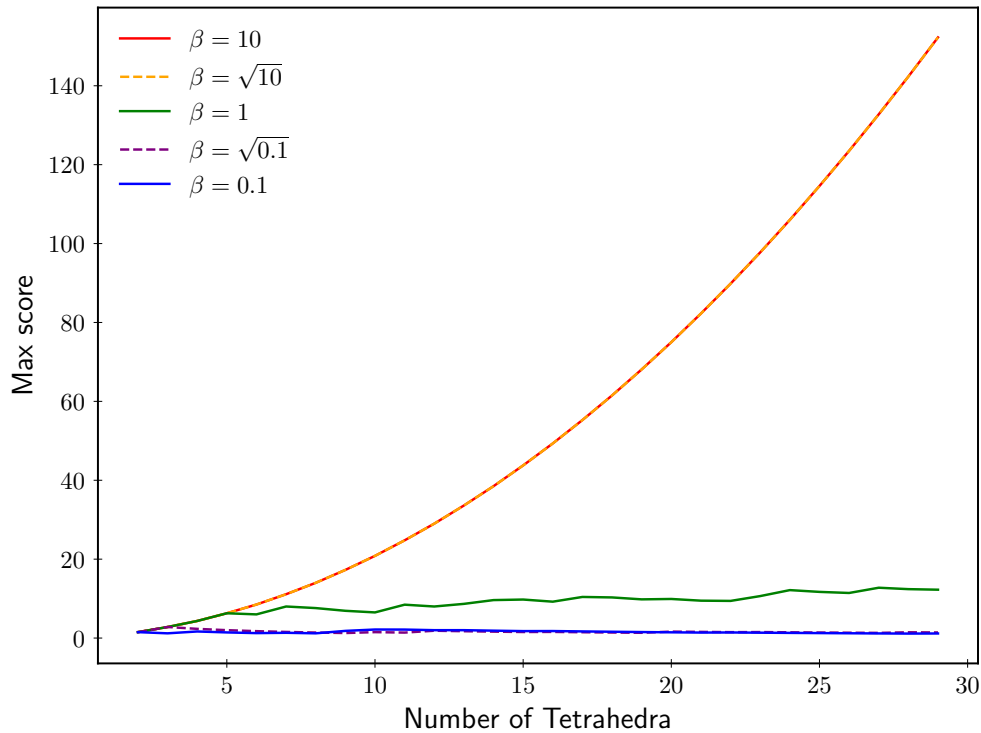


FIGURE 5. Direct Ascent on \mathcal{O}_{deg} at different temperatures ($\beta = 1/T$)

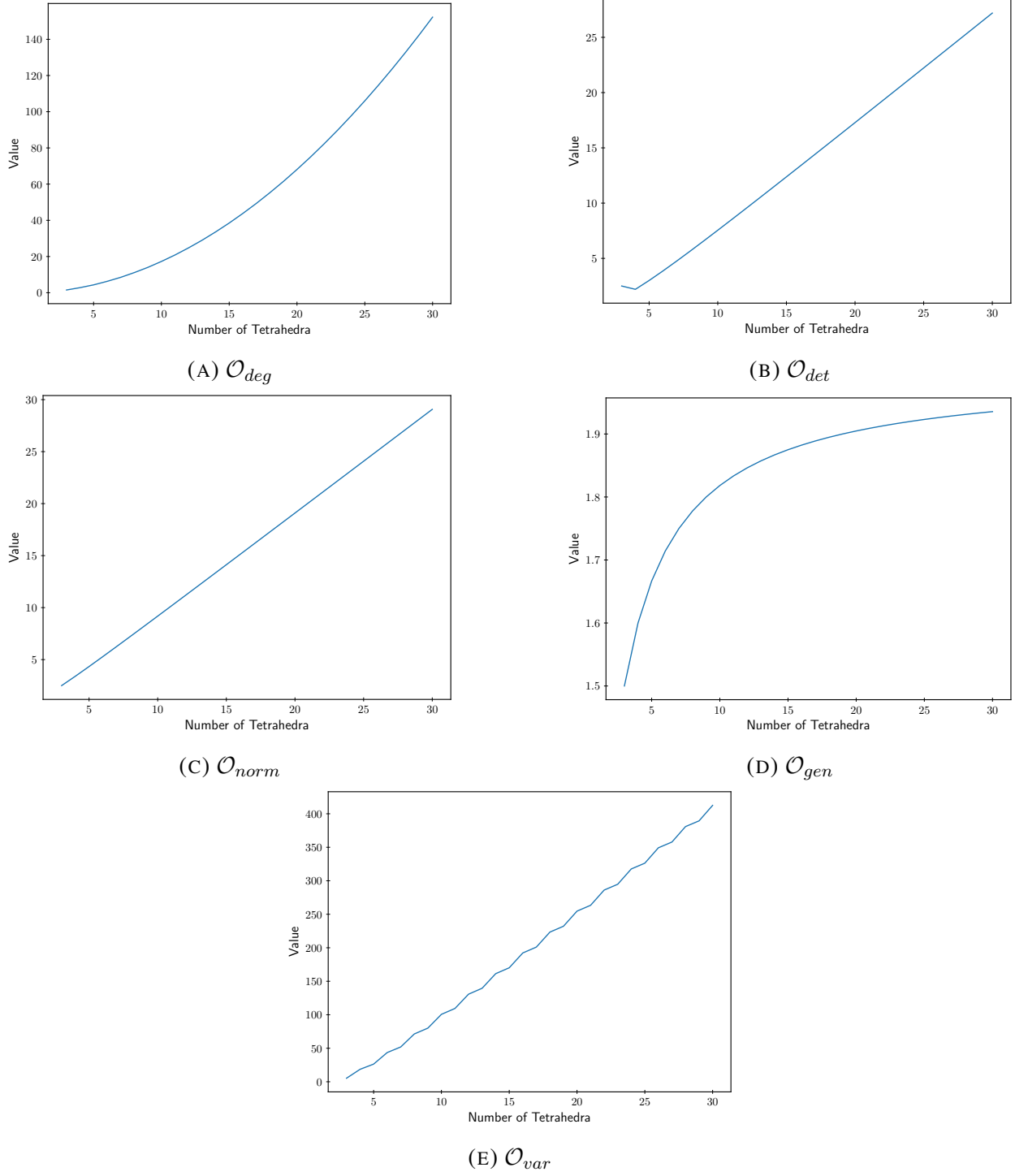


FIGURE 6. Direct / greedy ascent for each objective function.

3.3.2. Simulated Annealing. We perform simulated annealing as described in section 2.1.3. This is done for a target of 10,000 iterations with a step size of 10 and target acceptance rate of 20%, this acceptance rate was chosen heuristically based on a number of different small experiments. To ensure a fair comparison with direct ascent, the potential was set to $-\infty$ for triangulations with

more than 30 tetrahedra. For \mathcal{O}_{deg} and \mathcal{O}_{norm} , the chain was stopped early at 3, 307 and 2, 940 samples because the computational time required to score a single triangulation had become unfeasibly slow

(> 1min/sample). The actual acceptance rate was calculated in each case. For \mathcal{O}_{det} , \mathcal{O}_{gen} and \mathcal{O}_{var} the achieved acceptance rate was within 1% of the target. For \mathcal{O}_{deg} and \mathcal{O}_{norm} the acceptance rate is lower, though this is likely due to the reduced number of samples achieved. The chains for each objective function and the achieved acceptance rates are included in figure 7.

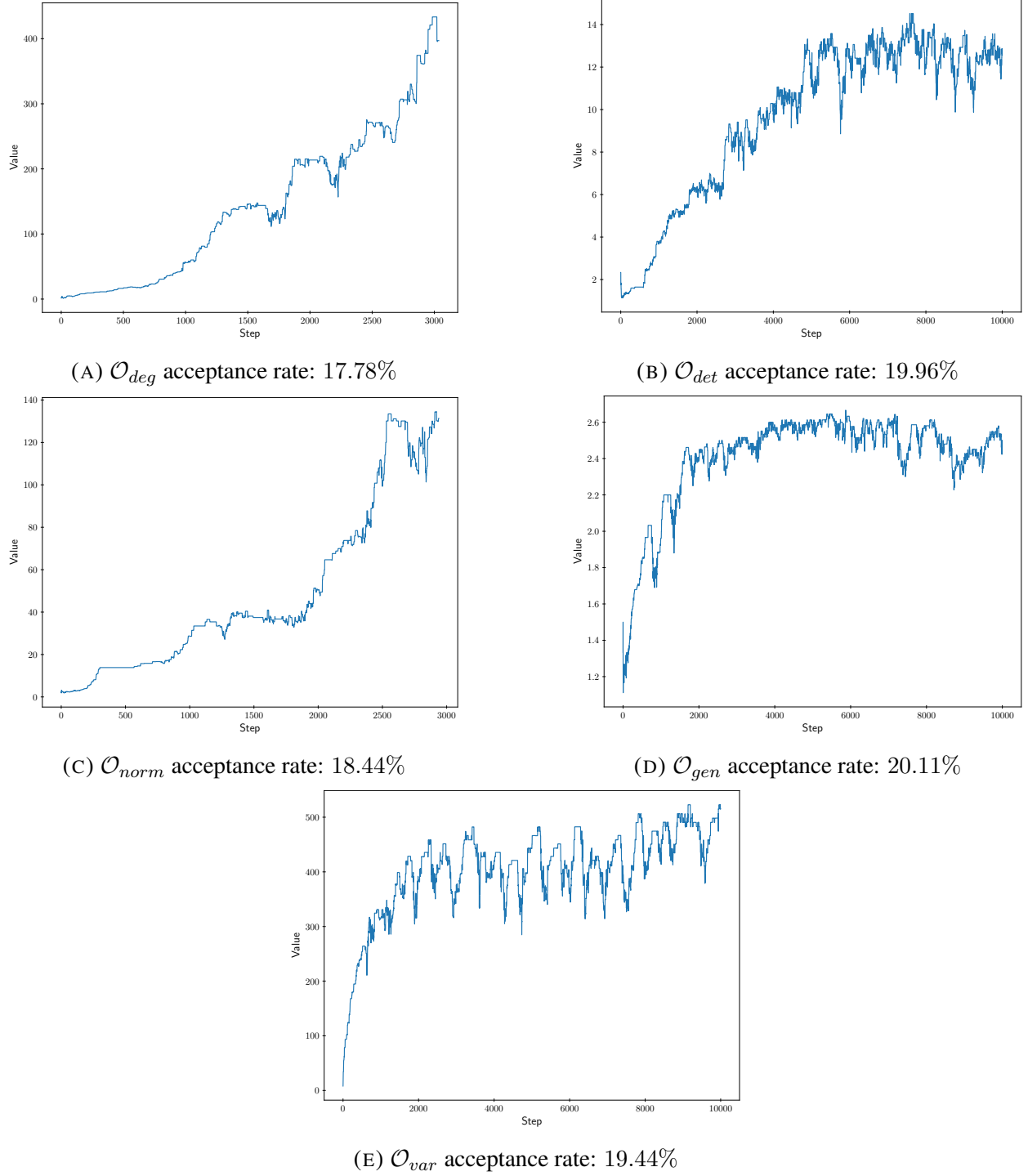


FIGURE 7. Trace plots for simulated annealing on each objective function for 10,000 iterations and a 20% acceptance rate.

3.4. Comparison of Classical Optimisation

The best result from each optimisation technique is compiled in table 2, this table includes the achieved percentile of the best sample, computed by the fitted power law distribution from section 3.2.

Objective Function	Direct Ascent			Simulated Annealing		
	Score	Bulk Percentile	Tail Percentile	Score	Bulk Percentile	Tail Percentile
\mathcal{O}_{deg}	152.26	$7.19 \cdot 10^{-05}$	$1.21 \cdot 10^{-05}$	433.42	$2.12 \cdot 10^{-05}$	$2.66 \cdot 10^{-06}$
\mathcal{O}_{det}	27.19	$1.38 \cdot 10^{-13}$	$2.62 \cdot 10^{-13}$	13.84	$5.86 \cdot 10^{-11}$	$1.61 \cdot 10^{-10}$
\mathcal{O}_{norm}	29.06	$7.16 \cdot 10^{-12}$	$3.29 \cdot 10^{-06}$	152.26	$2.38 \cdot 10^{-17}$	$1.23 \cdot 10^{-08}$
\mathcal{O}_{gen}	1.94	$1.58 \cdot 10^{-06}$	$2.32 \cdot 10^{-05}$	2.65	$2.93 \cdot 10^{-09}$	$1.14 \cdot 10^{-07}$
\mathcal{O}_{var}	412.54	$3.06 \cdot 10^{-15}$		522.93	$7.52 \cdot 10^{-18}$	

TABLE 2. Comparison of maximum achieved objective function for both classical sampling techniques. Percentiles of each score based on bulk and tail fitted power law distributions.

	Best Triangulation (IsoSig)
\mathcal{O}_{deg}	EMzMzMzMzMvAAMzMzMzMzMPabcdefghijklmnpqrstuvwxyzABCDDjxxxxxxxxxxxxxxxxlfxxxxxxxxxxxxxxxj
\mathcal{O}_{det}	ELLMzMzMzMzMzMzMzMzMzQbeddfghijklmnopqrstuvwxyzABCDDDxxrfxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxvc
\mathcal{O}_{norm}	EMzMzMzMzMvAAMzMzMzMzMPabcdefghijklmnpqrstuvwxyzABCDDjxxxxxxxxxxxxxxxxlfxxxxxxxxxxxxxxxj
\mathcal{O}_{gen}	ELwvMvQPLMzAPLLAPMzPQcbehgklmllpnqptrutwvzxyABDCBDDaacmvlpvlooashoqkmbuoageahamkcg
\mathcal{O}_{var}	EMzMzMzMzMzMvAAMzMzMzMzMPabcdefghijklmnpqrstuvwxyzABCDDvaaaaaaaaaaaaaaaaaaciqnaaaaaaaaaav

TABLE 3. Best Triangulation Direct Ascent

	Best Triangulation (IsoSig)
\mathcal{O}_{deg}	ELvLLzvMvzMQAMLMQzQMqbdgfpmtosmvrupqtvwuzyxBCzACBDDafgfjxaovwrdsifjouitrwblqjfhroj
\mathcal{O}_{det}	ELLvMLAAzPwMzAPzvAMQQadeihhhljnpnlngpturttxBBCADBAADDbabagjgxbadjbhohagkfauahoxrccjo
\mathcal{O}_{norm}	EvLPzvzvzLAvzLPPQQMQQQdcgklqsvAptytAzBxxzCwyBBDzADDonnapdeaacqowdbauqadasckkjabvhd
\mathcal{O}_{gen}	ELMvLAAvMAQzLwPLAMMPcbdgfiihnmlloopnqtsvtxwzyzABACDDaaaavlsvaccfbxsagranmglanxxxxxr
\mathcal{O}_{var}	EMLLLvzPPzwPLMQzAPzAQabddglnjokkqnouqswstuvwzyzBCDCDgagbaaagacvaggaccavgfbcfafoaco

TABLE 4. Best Triangulation Simulated Annealing

3.5. Baseline Transformer Efficiency on Isomorphism Signatures

CHAPTER 4

Discussion and Future Work

...