

Geometric Extensions of Neural Processes

Andrew Newberry Carr

A thesis submitted to the faculty of
Brigham Young University
in partial fulfillment of the requirements for the degree of
Master of Science

David Wingate, Chair
Bryan Morse
Michael Goodrich

Department of Computer Science
Brigham Young University

Copyright © 2019 Andrew Newberry Carr

All Rights Reserved

ABSTRACT

Geometric Extensions of Neural Processes

Andrew Newberry Carr

Department of Computer Science, BYU

Master of Science

Neural Processes (NPs) [18] are a class of models that learn a mapping from a context set of input-output pairs to a distribution over functions. The end result is similar to that of the well studied Gaussian Process. However, NPs are more computationally tractable and provide a number of benefits over traditional models. There are interesting extensions to Neural Processes which we address in this work. The first of these is that while Neural Processes have been shown to work well with Euclidean data, they are unable to processes graph structured data. To overcome this limitation, we introduce Graph Neural Processes (GNP).

A Graph Neural Process is defined as a Conditional Neural Process that operates on arbitrary graph data. It takes features of sparsely observed context points as input, and outputs a distribution over target points. We demonstrate graph neural processes in edge imputation and discuss benefits and draw backs of the method for other application areas. One major benefit of GNPs is the ability to quantify uncertainty in deep learning on graph structures. An additional benefit of this method is the ability to extend graph neural networks to inputs of dynamic sized graphs.

A second extension of Neural Processes comes in the fundamental training mechanism. They are traditionally trained using maximum likelihood with a KL divergence regularization term. We show that there are desirable classes of problems where NPs, with this loss, fail to learn any reasonable distribution. We also show that this drawback is solved by using approximations of Wasserstein distance which calculates optimal transport distances even for distributions of disjoint support. We give experimental justification for our method and demonstrate performance. These Wasserstein Neural Processes (WNPs) maintain all of the benefits of traditional NPs while being able to approximate a new class of function mappings.

Keywords: Optimal Transport, Machine Learning, Graph Laplacian, Neural Networks

ACKNOWLEDGMENTS

I would first like to acknowledge my darling wife AnnaLisa. She pushed me by showing me what a real Master's Thesis should look like. I am also grateful to my family for being willing to listen to me geek out about my research and believing in me, even if I didn't always believe in myself.

I also am immensely grateful for David Wingate. He helped me start my journey into ML in 2016. His guidance, mentorship, technical ability, and enthusiasm were pivotal.

I also want to acknowledge the Applied and Computational Math program at BYU. The creators, professors, students, and TAs pushed and stretched me in uncomfortable, but amazing ways. I am forever grateful for the skills and critical thinking abilities I gained in that program.

Table of Contents

| | |
|---------------------------------------------|-------------|
| List of Figures | vi |
| List of Tables | vii |
| List of Listings | viii |
| 1 Background | 2 |
| 1.1 What is a process? | 2 |
| 1.2 Gaussian Process | 3 |
| 1.3 Metrics and Measures | 4 |
| 1.4 Optimal Transport | 6 |
| 1.4.1 Wasserstein Distance | 8 |
| 1.4.2 Sinkhorn | 9 |
| 1.4.3 Sliced Wasserstein Distance | 10 |
| 2 Neural Processes | 12 |
| 2.1 Conditional Neural Processes | 12 |
| 2.2 Neural Process | 13 |
| 2.3 Extensions | 14 |
| 3 Graph Neural Processes | 15 |
| 3.1 Introduction | 15 |
| 3.2 Background | 16 |
| 3.2.1 Graph Structured Data | 16 |

| | | |
|----------|--------------------------------------------------------------------------------------|-----------|
| 3.2.2 | Graph Neural Networks | 17 |
| 3.3 | Related Work | 18 |
| 3.3.1 | Edge Imputation | 18 |
| 3.3.2 | Bayesian Deep Learning | 18 |
| 3.4 | Model and Training | 19 |
| 3.5 | Applications | 21 |
| 3.5.1 | Results | 23 |
| 3.6 | Areas for Further Exploration | 24 |
| 4 | Wasserstein Neural Processes | 29 |
| 4.1 | Introduction | 29 |
| 4.2 | Wasserstein Neural Processes | 31 |
| 4.3 | Experiments | 33 |
| 4.3.1 | Misspecified Models - Linear Regression with Uniform Noise | 33 |
| 4.3.2 | Intractable Likelihoods - The Quantile “ g -and- κ ” Distribution | 34 |
| 4.3.3 | CelebA tiles as super pixels | 35 |
| 4.3.4 | Discussion | 38 |
| 5 | Discussion and Future Work | 39 |
| 5.1 | Graph Neural Processes | 39 |
| 5.2 | Wasserstein Neural Processes | 40 |
| | References | 41 |

List of Figures

| | | |
|-----|-------------------------------------------------------|----|
| 1.1 | Continuous optimal transport theory [62] | 6 |
| 1.2 | Regularized Sinkhorn | 9 |
| 2.1 | CNP Architecture | 13 |
| 3.1 | Graph with imputed edge distributions | 16 |
| 3.2 | Precision results | 25 |
| 3.3 | Recall results | 26 |
| 3.4 | F1-score results | 27 |
| 4.1 | WNP regression results | 34 |
| 4.2 | WNP g -and- κ distribution results | 36 |
| 4.3 | WNP image reconstruction | 38 |

List of Tables

| | | |
|-----|----------------------------------------------|----|
| 3.1 | GNP experimental results | 27 |
| 3.2 | Features of the explored data sets | 28 |

List of Listings

Introduction

One strategy of machine learning research is to tackle well established problems. For example, there are many researchers working on image classification. Another strategy is to find new problems that the community will find interesting. One final strategy is to implement novel algorithms and architectures then apply them across a broad range of problems.

In 2018 a new class of architectures was discovered [17, 18, 32] that generalized Gaussian Processes (GP) using neural methods. These architectures are broadly referred to as Neural Processes (NP).

This new method was applied to regression problems on data in \mathbb{R}^n , image completion, and image generation. However, because of the training regime, and the set up of the problem there are two cases where NP breaks.

1. Non Euclidean data like graphs
2. Distributions with unknown, misspecified, or intractable likelihood.

In this work, we present two variants of NP that are aimed to solve these specific problems. First, we introduce Graph Neural Processes (GNP) that operate on graph structured data for regression and edge imputation. We show that GNPs out performs a number of competitive baselines. Secondly, introduce Wasserstein Neural Processes (WNP) which overcome likelihood training problems. We show there are situations where NPs break, but WNP work well. We also show competitive performance on image completion.

This document is designed to be relatively self contained. We'll start by covering the background on processes, basic measure theory, optimal transport, and neural processes. Then, there is a chapter dedicated to each of the methods we present.

Chapter 1

Background

1.1 What is a process?

Probability theory is primarily concerned with the study of random variables. Contrary to the common interpretation of the name, these objects are neither random or variables. They are objects whose values are determined by the outcome of some random phenomenon. For example, if you flip a coin, or role a dice, the mathematical representation of that action (and corresponding outcome) is a random variable.

Much work has been done to formalize probability and random variables in the language of measure theory. We will briefly summarize this formalism in Section 1.3. For now, it is sufficient to think of these objects as a set of possible values weighted by some frequency.

A *Process* is then simply a collection of random variables indexed by time or space. A simple example is a coin that gets less “fair” every time it is flipped. These stochastic processes can be thought of as potentially random numerical value that changes over time. This simple idea has been used to great effect in statistics and learning theory to model a wide variety of phenomena.

One area of inquiry is to train a model of these processes with some specific inductive bias. The bias comes from some abstract prior we hold over the shape of our data. For example, if we are modeling the probability of a phone call in a call center, we impose the inductive bias that our random variables are distributed according to some Poisson process. This prior gives us power to use a process to our benefit.

1.2 Gaussian Process

When any linear combination of the random variables of our stochastic process is Gaussian, we have a Gaussian Process (GP). This condition is also sometimes called a multivariate normal distribution. Gaussian processes are valuable because distributions of interest can be obtained in closed form.

GPs are useful when used as a solution to regression problems. Since random variables are technically measurable functions, a Gaussian process is used to define a prior over functions that can be used to compute a posterior over functions after observing some data. In practice, representing a distribution over functions is done by a distribution over the function *values* at a finite set of points X_1, \dots, X_N .

The assumption of the GP is that $(f(X_1), \dots, f(X_N))$ are jointly Gaussian with some empirical mean μ and covariance Σ . The Σ is calculated using a kernel such that $\Sigma_{ij} = \kappa(x_i, x_j)$. The computation and regression can be done in $O(N^3)$ time. This is one main drawback of the method since it suffers severely from the curse of dimensionality [63].

In this work, we will consider the basics of using Gaussian Processes for regression problems. This will give context to later methods and is interesting in its own right. We first let the prior on the regression function be given as

$$f(\vec{x}) \sim GP(m(\vec{x}), \kappa(\vec{x}, \vec{x}')) \quad (1.1)$$

Where $m(\cdot)$ is the mean function and $\kappa(\cdot, \cdot)$ is the kernel function, also known as the covariate function. These are defined in terms of the expected value.

$$m(\vec{x}) = E[f(\vec{x})] \quad (1.2)$$

$$k(\vec{x}, \vec{x}') = E[(f(\vec{x}) - m(\vec{x}))(f(\vec{x}') - m(\vec{x}'))^T] \quad (1.3)$$

This process then defines a joint Gaussian output over functions. Or, in other words, we have $p(f|X) \sim N(f|\mu, K)$ for vectors of m and κ of the proper size for μ and K .

The science of the GP is picking the proper kernel function for your data set. Intuitively, the kernel defines a distance which should have some proper inductive bias to the problem being solved. We will see later in this manuscript that Neural Processes, while loosely based on GP, are different in the best ways.

1.3 Metrics and Measures

In the world of mathematics, probability was used and informally defined for hundreds of years. However, in the 1950's the foundation of probability was rigorously built from first principles [22]. One of the main building blocks of probability is the *measure* and a number of mechanical tools (e.g., sigma algebra) to give formalism to randomness.

Informally, a measure is a function with certain properties familiar to us (e.g., integral of the entire probability space is 1) that operates on a probability space. A measure, typically written as $\mu : \Omega \mapsto \mathbb{R}$, is a real-valued function that we use as a rigorous way to work with probability distributions.

Formally, a measure is a generalization of length, area, and volume that defines the mass of some distribution. Let Ω be the sample space which is the set of all possible outcomes for a process. Then F is the set of events which consist of zero or more outcomes. A probability space is a triple (Ω, F, P) where P is some way to assign probabilities (read mass) to events. You may notice that in a measurable probability space P will be a measure μ . As one may expect, these three pieces have specific structure that makes them useful for analysis.

Firstly, the sample space Ω must be non-empty. The events object F must be a σ -algebra which means it is self inclusive, closed under complements, and closed under

countable unions. Formally, this is summarized in the following equations.

$$\Omega \in F \tag{1.4}$$

$$\text{if } A \in F \text{ then we also know } (\Omega \setminus A) \in F \tag{1.5}$$

$$\text{if } A_i \in F \text{ for } i = 1, 2, \dots \text{ then } (\cup_{i=1}^{\infty} A_i) \in F \tag{1.6}$$

And finally, the measure μ must be a probability measure. Which implies it behaves similarly to discrete probability distributions. In other words,

$$\text{if } \{A_i\}_{i=1}^{\infty} \subset F \text{ is pairwise disjoint and countable then} \tag{1.7}$$

$$\mu(\cup_{i=1}^{\infty} A_i) = \sum_{i=1}^{\infty} \mu(A_i) \tag{1.8}$$

$$\mu(\Omega) = 1 \tag{1.9}$$

This structure on the space gives us a probability space that we can use in a number of extremely interesting ways. We can now reason rigorously about the shape of distributions, the distances between these distributions, and properties of the distributions.

In general, spaces can be given additional structure through the identification of a metric. A metric is a distance. It is a function that maps a pair of points in the space to a single value on the real line. To be considered a metric, however, there are a number of axioms the function must fulfill.

$$d(x, y) \geq 0 \tag{1.10}$$

$$d(x, y) = 0 \text{ iff } x = y \tag{1.11}$$

$$d(x, y) = d(y, x) \tag{1.12}$$

$$d(x, z) \leq d(x, y) + d(y, z) \tag{1.13}$$

Looking forward, we will discuss a natural way to define a metric (distance) over measures in a probability space. This distance, called the Kantorovich, or Wasserstein, distance has many applications in applied math, computer vision, and machine learning.

1.4 Optimal Transport

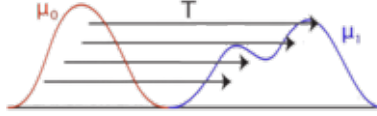


Figure 1.1: Continuous optimal transport theory [62]

As the founder of Optimal Transport theory, Gaspard Monge [45] was interested in the problem of a worker with a shovel moving a large pile of sand. The worker's goal is to erect a target pile of sand with a prescribed shape. The worker wishes to minimize their overall effort in carrying shovelfuls of sand. Monge formulated the problem as such in 1781 and worked to solve it.

Given two densities of mass $f, g \geq 0$ on \mathbb{R}^d with unit mass $\int f(x)dx = \int g(y)dy = 1$, find an optimal (minimizing) map $T : \mathbb{R}^d \rightarrow \mathbb{R}^d$ which pushes one density onto another in such a way that

$$\int_A g(y)dy = \int_{T^{-1}(A)} f(x)dx \quad (1.14)$$

for any Borel subset $A \subset \mathbb{R}^d$ such that optimal is considered with respect to the quantity

$$M(T) = \int_{\mathbb{R}^d} |T(x) - x| f(x) dx \quad (1.15)$$

A physically intuitive interpretation is that of tiny sand particles that have a density described by f , and move to another configuration given by g in such a way as to minimize

the average displacement across all particles. This is done since $T(x)$ represents the motion or destination of one particle that started at x .

More generally, consider a measurable map $T : X \rightarrow Y$ where mass and density are preserved.

$$T_{\#}\mu(A) = \mu(T^{-1}(A)) \quad (1.16)$$

Then the optimal transport problem as formulated by Monge can be concisely written with two measurable functions μ, ν as

$$\inf_{T_{\#}\mu=\nu} \int_X c(x, T(x)) d\mu(x) \quad (1.17)$$

for some more general cost function c . In many cases the cost function is taken to be the L_1 or L_2 norms which each have pros and cons. However, in general the non-linearity of the relation between the map, measures, and cost function make analysis of this formulation challenging. Additionally, this formulation requires all mass to go from source to target. There is no mechanism for partial transportation of mass. In general, progress wasn't really made on this problem until 1942 when Kantorovich introduce probabilistic couplings to the problem. Instead of giving the destination $T(x)$ for each particle, Kantorovich's formulation gives the number of particles moving between locations for all pairs (x, y) . He did this by introducing a coupling that is called a "transport plan":

$$\Pi(\mu, \nu) = \{\pi \in P(X \times Y) : (\pi_x)_{\#}\pi = \mu, (\pi_y)_{\#}\pi = \nu\} \quad (1.18)$$

where π_x and π_y are projections from the product space. These are also often referred to as the marginals of the coupling.

The goal then in the Kantorovich formulation

$$W(\mu, \nu) = \inf_{\pi \in \Pi(\mu, \nu)} \int \int c(x, y) d\pi(x, y) \quad (1.19)$$

is to find the optimal coupling between two measures.

1.4.1 Wasserstein Distance

Over the years, optimal transport theory has grown into a mature and exciting field [59] with a number of awards going to those working on particular advances. In particular, the Wasserstein distance has seen much recent attention due to its nice theoretical properties and myriad applications [16]. These properties enable the study of distance between two probability distributions μ, ν for generative modeling.

The p -Wasserstein distance, $p \in [1, \infty]$, between μ, ν is defined to be the solution to the original Optimal Transport problem [59]:

$$W_p(\mu, \nu) = \left(\inf_{\pi \in \Pi(\mu, \nu)} \int_{X \times Y} d^p(x, y) d\pi(x, y) \right)^{\frac{1}{p}} \quad (1.20)$$

Where $\Gamma(\mu, \nu)$ is the joint probability distribution with marginals μ, ν and $d^p(\cdot, \cdot)$ is the cost function.

An exciting result of the Wasserstein distance is that the quantity actually defines a distance (metric) on the space of probability measures! There are many geometric properties such as geodesics and gradient flow that can be calculated now using this distance. Unfortunately, in the computational setting, this distance originally was difficult to recover.

In general, this distance is expensive to compute. There are, however, a number of methods used to approximate this distance. Firstly, we recall the Kantorovich-Rubinstein dual formulation: [1, 59]

$$W(\mu, \nu) = \sup_{\|f\|_L \leq 1} E_{x \sim \mu}[f(x)] - E_{x \sim \nu}[f(x)] \quad (1.21)$$

Where we have let $p = 1$. It is important to note that the supremum is taken over all of the 1-Lipschitz functions. From this form emerges an adversarial algorithm for Wasserstein distance calculation. This is explored in a number of generative adversarial network (GAN) works, with [21] offering good performance with the addition of a gradient penalty. Additionally, by using a similar dual form a two-step computation that makes use of linear programming can be employed to exactly calculate the Wasserstein distance [42]. This method has favorable performance and is also used in generative modeling with the WGAN-TS[42].

However, Wasserstein distances are computationally expensive. Therefore, many approximations have arisen that maintain many of the positive properties of optimal transport while being significantly more computationally efficient. We present a small subset of these methods here.

1.4.2 Sinkhorn

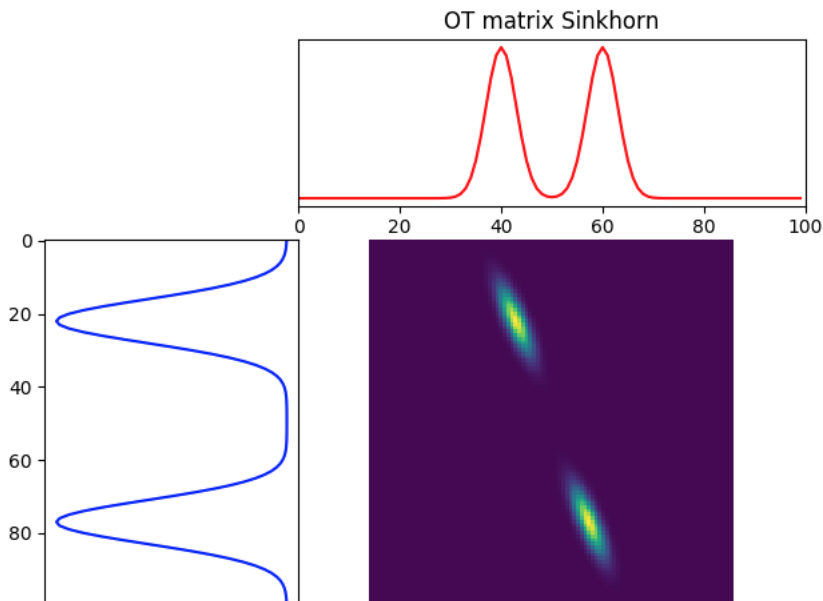


Figure 1.2: Regularized sinkhorn optimal transport plan between two empirical bimodal distributions

Another prominent method used to calculate Wasserstein distance that is much faster to compute is that of Sinkhorn distances [12]. This method makes use of Sinkhorn-Knopp’s matrix scaling algorithm to compute the Wasserstein distance. These methods also introduce an entropic regularization term that drastically decreases the complexity of this problem offering speed ups of several orders of magnitude. Work has been successfully done to create a differentiable method [19] for using sinkhorn loss during training of generative models.

In the case of empirical one-dimensional distributions, the p -Wasserstein distance is calculated by sorting the samples from both distributions and calculating the average distance between the sorted samples (according to the cost function). This is a fast operation requiring $\mathcal{O}(n)$ operations in the best case and $\mathcal{O}(n \log n)$ in the worst case with n the number of samples from each distribution [39]. This simple computation is taken advantage of in the computation of the sliced Wasserstein distance.

1.4.3 Sliced Wasserstein Distance

The sliced Wasserstein distance is qualitatively similar to Wasserstein distance, but as alluded to above is much easier to compute since it only depends on computations across one-dimensional distributions [37]. Intuitively, sliced Wasserstein distance is an aggregation of linear projects from high-dimensional distributions to one-dimensional distributions where regular Wasserstein distance is then calculated. The projection process is done via the Radon transform [38]. The one-dimensional projects are defined as

$$\mathcal{R}p_X(t; \theta) = \int_X p_X(x) \delta(t - \theta \cdot x) dx, \quad \forall \theta \in S^{d-1}, \quad \forall t \in \mathbb{R} \quad (1.22)$$

Then, by using these one-dimensional projections [7] the sliced Wasserstein distance is defined as:

$$SW_p(\mu, \nu) = \left(\int_{S^{d-1}} W_p(\mathcal{R}_\mu(\cdot; \theta), \mathcal{R}_\nu(\cdot; \theta)) d\theta \right)^{\frac{1}{p}} \quad (1.23)$$

It is important to note that this map SW_p inherits its being a distance from the fact that W_p is a distance. Similarly, the two distances induce the same topology on compact sets [51]. The integral in (eq. 1.23) is approximated by a Monte Carlo sampling variant of normal draws on S^{d-1} which are averaged across samples

$$SW_p(\mu, \nu) \approx \left(\frac{1}{N} \sum_{i=1}^N W_p(\mathcal{R}_\mu(\cdot; \theta_i), \mathcal{R}_\nu(\cdot; \theta_i)) \right)^{\frac{1}{p}} \quad (1.24)$$

Chapter 2

Neural Processes

2.1 Conditional Neural Processes

First introduced in [17], Conditional Neural Processes (CNPs) are loosely based on traditional Gaussian Processes. CNPs map an input $\vec{x}_i \in \mathbb{R}^{d_x}$ to an output $\vec{y}_i \in \mathbb{R}^{d_y}$. It does this by defining a collection of conditional distributions that can be realized by conditioning on an arbitrary number of context points $X_C := (\vec{x}_i)_{i \in C}$ and their associated outputs Y_C . Then an arbitrary number of target points $X_T := (\vec{x}_i)_{i \in T}$, with their outputs Y_T can be modeled using the conditional distribution. Like GPs, CNPs are invariant to ordering of context points, and ordering of targets. It is important to note that, while the model is defined for arbitrary context points C and target points T , it is common practice (which we follow) to use $C \subset T$.

The context point data is aggregated using a commutative operation \oplus that takes elements in some \mathbb{R}^d and maps them into a single element in the same space. In the literature, this is referred to as the r_C context vector. We see that r_C summarizes the information present in the observed context points. Formally, the CNP is learning the following conditional distribution.

$$P(Y_T | X_T, X_C, Y_C) \iff P(Y_T | X_T, r_C) \quad (2.1)$$

In practice, this is done by first passing the context points through a Neural Network h_θ to obtain a fixed length embedding r_i of each individual context point. These context point representation vectors are aggregated with \oplus to form r_C . The target points X_T are then

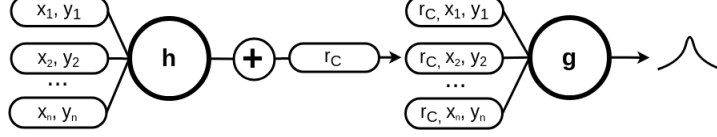


Figure 2.1: Conditional Neural Process Architecture. Traditionally, the data is maximized under the output distribution. In this work, we use the sliced Wasserstein distance. Also note that \oplus is an arbitrary commutative aggregation function and h, g are parameterized neural networks. As an example, the context/target points could be (x, y) coordinates and pixel intensities if an image is being modeled.

decoded, conditional to r_C , to obtain the desired output distribution z_i the models the target outputs y_i .

More formally, this process can be defined as follows.

$$r_i = h_\theta(\vec{x}_i) \quad \forall \vec{x}_i \in X_C \quad (2.2)$$

$$r_C = r_1 \oplus r_2 \oplus r_3 \oplus \dots \oplus r_n \quad (2.3)$$

$$z_i = g_\phi(\vec{y}_i | r_C) \quad \forall \vec{y}_i \in X_T \quad (2.4)$$

2.2 Neural Process

Neural Processes are an extension of CNPs where a parameterized Normal $s_C := s(x_C, y_C)$ latent path is added to the model architecture. The purpose of this path is to model global dependencies between context and target predictions [18]. The computational complexity of the Neural Processes is $\mathcal{O}(n + m)$ for n contexts and m targets. Additionally, [18] show that in addition to being *scalable* these models are *flexible* and *invariant* to permutations in context points. Many of these benefits are shared in with CNPs, albeit with CNPs showing slightly worse experimental performance.

With the addition of the latent path, a new loss function is required. The encoder and decoder parameters are maximized using the ELBO loss:

$$\log p(y_T | x_T, x_C, y_C) \geq E_{q(z|s_T)} [\log p(y_T | x_T, r_C, z)] - D_{KL}(q(z|s_T) || q(z|s_C)) \quad (2.5)$$

A random subset of context points is used to predict a set of target points. With a KL regularization term that encourages the summarization of the context to be close to the summarization of the targets.

Work has been done [32] to explore other aggregation operations such as attention [3]. These extensions solve a fundamental underfitting problems that plague both CNPs and NPs. These Attentive Neural Processes (ANP) introduce several types of attention into the model architecture in place of the neural networks which greatly improves performance at the cost of computational complexity which was shown to be $\mathcal{O}(n(n+m))$. It is important to note, however, that the introduction of attention causes the models to train more quickly (i.e., fewer iterations) which is faster in wall-clock time even if inference is slower. The ANPs utilize a fundamentally different architecture while still maximizing the ELBO loss in (eq 2.5).

2.3 Extensions

The general NP framework is a useful extension and generalization of GPs. Researchers noticed, however, that a main bottle neck is the commutative aggregation operation. As such, [32] introduce the Attentive Neural Process which uses attention to combine the encoded r vectors. Other extensions [43, 47] have also been proposed to improve training and stability.

In this work, we extend the neural process in two distinct ways. Firstly, we design a mechanism to work on graph structured data. Secondly, we overcome weaknesses in maximum likelihood learning using Wasserstein distance to learn a larger variety of desirable functional distributions.

Chapter 3

Graph Neural Processes

3.1 Introduction

In this work, we consider the problem of imputing the value of an edge on a graph. This is a valuable problem when an edge is known to exist, but due to a noisy signal, or poor data acquisition process, the value is unknown. We solve this problem using a proposed method called Graph Neural Processes.

In recent years, deep learning techniques have been applied, with much success, to a variety of problems. However, many of these neural architectures (e.g., CNNs) rely on the underlying assumption that our data is Euclidean. However, since graphs do not lie on regular lattices, many of the concepts and underlying operations typically used in deep learning need to be extended to non-Euclidean valued data. Deep learning on graph-structured data has received much attention over the past decade [20] [54] [40] [65] [5] and has show significant promise in many applied fields [5]. These ideas were recently formalized as *geometric deep learning* [9] which categorizes and expounds on a variety of methods and techniques. These techniques are mathematically designed to deal with graph-structured data and extend deep learning into new research areas.

Similarly, but somewhat orthogonally, progress has been made in Bayesian methods when applied to deep learning [2]. In Bayesian neural networks, uncertainty estimates are used at the weight level, or at the output of the network. These extensions to typical deep learning give insight into what the model is learning, and where it may encounter failure

modes. In this work, we use some of this progress to impute the value distribution on an edge in graph-structured data.

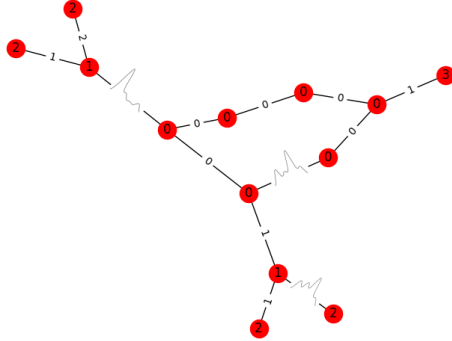


Figure 3.1: Graph with imputed edge distributions

In this work we propose a novel architecture and training mechanism which we call Graph Neural Processes (GNP). This architecture is based on the ideas first formulated in [17] in that it synthesizes global information to a fixed length representation that is then used for probability estimation. Our contribution is to extend those ideas to graph-structured data and show that the methods perform favorably.

Specifically, we use features typically used in Graph Neural Networks as a replacement to the convolution operation from traditional deep learning. These features, when used in conjunction with the traditional CNP architecture offer local representations of the graph around edges and assist in the learning of high level abstractions across classes of graphs. Graph Neural Processes learn a high level representation across a family of graphs, in part by utilizing these instructive features.

3.2 Background

3.2.1 Graph Structured Data

We define a graph $G = (a, V, E)$ where a is some global attribute¹, V is the set of nodes where $V = \{v_i\}_{i=1:N^v}$ with v_i being the node’s attribute, and E is the set of edges where

¹typically a member of the real numbers, \mathbb{R} . Our method does not utilize this global attribute.

$E = \{e_k, v_k, u_k\}_{k=1:N^e}$ where e_k is the attribute on the edge, and v_k, u_k are the two nodes connected by the edge. In this work, we focus on undirected graphs, but the principles could be extended to many other types of graphs.

3.2.2 Graph Neural Networks

Of all the inductive biases introduced by standard deep learning architectures, the most common is that we are working with Euclidean data. This fact is exploited by the use of convolutional neural networks (CNN), where spatially local features are learned for downstream tasks. If, however, our data is Non-Euclidean (e.g., graphs, manifolds) then many of the operations used in standard deep learning (e.g., convolutions) no longer produce the desired results. This is due to the fact that the intrinsic measures and mathematical structure on these surfaces violates assumptions made by traditional deep learning operations. There has been much work done recently in Graph Neural Networks (GNN) that operate on graph-structured inputs [65]. There are two main building blocks of geometric methods in GNNs. *spectral* [9] and *spatial* [14] methods. These methods are unified into *Message Passing Neural Networks*, and *Non-local Neural Networks*. A more thorough discussion of these topics can be found in [65]. Typically, when using these GNN methods, the input dataset graphs need to have the same number of nodes/edges. This is because of the technicalities involved in defining a spectral convolution. However, in the case of Graph Neural Processes, local graph features are used which allows one to learn conditional distributions over arbitrarily sized graph-structured data. One important concept utilized in the spectral GNN methods is that of the Graph Laplacian. Typically, the graph Laplacian is defined as the adjacency matrix subtracted from the degree matrix $L = D - A$. This formulation, common in multi-agent systems, hampers the flow of information because it is potentially non-symmetric, and unnormalized. As such, the Normalized Symmetric Graph Laplacian is given as follows.

$$L = I_n - D^{-1/2}AD^{-1/2} \quad (3.1)$$

This object can be thought of as the difference between the average value of a function around a point and the actual value of the function at the point. This, therefore, encodes local structural information about the graph itself.

In spectral GNN methods, the eigenvalues and eigenvectors of the graph Laplacian are used to define convolution operations on graphs. In this work, however, we use the local structural information encoded in the Laplacian as an input feature for the context points we wish to encode.

3.3 Related Work

3.3.1 Edge Imputation

In many applications, the existence of an edge is known, but the value of the edge is unknown (e.g, traffic prediction, social networks). Traditional edge imputation involves generating a point estimate for the value on an edge. This can be done through mean filling, regression, or classification techniques [27]. These traditional methods, especially mean filling, can fail to maintain the variance and other significant properties of the edge values. In [27] they show “Bias in variances and covariances can be greatly reduced by using a conditional distribution and replacing missing values with draws from this distribution.” This fact, coupled with the neural nature of the conditional estimation, gives support to the hypothesis that Graph Neural Processes preserve important properties of edge values, and are effective in the imputation process.

3.3.2 Bayesian Deep Learning

In Bayesian neural networks, it is often the case where the goal is to learn a function $y = f(x)$ given some training inputs $X = \{x_1, \dots, x_n\}$ and corresponding outputs $Y = \{y_1, \dots, y_n\}$. This function is often approximated using a fixed neural network that provides a likely correlational explanation for the relationship between x and y . There has been good work done in this area [64] and there are two main types of Bayesian Deep Learning. This is

not an exhaustive list of all the methods, but a broad overview of two. Firstly, instead of using point estimates for the weights W of the neural network hidden layers, a distribution over values is used. In other words, the weights are modeled as a random variable with an imposed prior distribution. This encodes, *a priori* uncertainty information about the neural transformation. Similarly, since the weight values W are not deterministic, the output of the neural network can also be modeled as a random variable. A generative model is learned based on the structure of the neural network and loss function being used. Predictions from these networks can be obtained by integrating with respect to the posterior distribution of W

$$p(y|x, X, Y) = \int p(y|x, W)p(W|X, Y)dW. \quad (3.2)$$

This integral is often intractable in practice, and number of techniques have been proposed in the literature to overcome [64] this problem. Intuitively, a Bayesian neural network is encoding information about the distribution over output values given certain inputs. This is a very valuable property of Bayesian deep learning that GNPs help capture. In the case of Graph Neural Processes, we model the output of the process as a random variable and learn a conditional distribution over that variable. This fits into the second class of Bayesian neural networks since the weights W are not modeled as random variables in this work.

3.4 Model and Training

While the h_θ and g_ϕ encoder and decoder could be implemented arbitrarily, we use fully connected layers that operate on informative features from the graph. In GNPs, we use local spectral features derived from global spectral information of the graphs. Typical graph networks that use graph convolutions require a fixed size Laplacian to encode global spectral information; for GNPs we use an arbitrary fixed size neighborhood of the Laplacian around each edge.

To be precise, with the Laplacian L as defined in equation (3.1) one can compute the spectra $\sigma(L)$ of L and the corresponding eigenvector matrix Λ . In Λ , each column is an eigenvector of the graph Laplacian. To get the arbitrary fixed sized neighborhood around each node/edge we define the restriction of Λ as

$$\Lambda|_k = (\Lambda_{kj})_{\substack{k \in r \\ 1 \leq j \leq m}} \quad (3.3)$$

Where m is some arbitrary fixed constant based on the input space dimension. We call this restriction the local structural eigenfeatures of an edge. These eigenfeatures are often used in conjunction with other, more standard, graph based features. For example, we found that the node values and node degrees for each node attached to an edge serve as informative structural features for GNPs. In other words, to describe each edge we use the local structural eigenfeatures and a 4-tuple $(v_i, u_j, d(v_i), d(u_j))$ where $d : V \rightarrow \mathbb{R}$ is a function that returns the degree of a node.

Formally, we define the **Encoder** h_θ , to be a 4-layer fully connected neural network with ReLU non-linearities. It takes the GNP features from each context point as input, and outputs a fixed length (e.g., 256) vector. This vector r_C is the aggregation of information across the arbitrarily sampled context points. The **Decoder** g_ϕ is also a 4-layer fully connected neural network with ReLU non-linearities. The decoder takes, as input, a concatenation of the r_C vector from the encoder and a list of features all the edges sans attribute. Then, for each edge, the concatenation of r_C and the feature vector is passed through the layers of the decoder until an output distribution is reached. Where the output distribution size is $|\cup \{e_k : e_k \in E\}|$, which is the number of unique values assigned to an edge.

Additionally, we obtain the **context points** by first defining a lower bound p_0 and upper bound p_1 for a Uniform probability distribution. Then, the number of context points is $p \cdot n$ where $p \sim \text{unif}(p_0, p_1)$ and $n = |E|$. The context points all come from a single graph,

and training is done over a family of graphs. The GNP is designed to learn a representation over a family of graphs and impute edge values on a new member of that family.

The loss function \mathcal{L} is problem specific. In the CNP work, as mentioned above, Maximum Likelihood is used to encourage the output distribution of values to match the true target values. Additionally, one could minimize the Kullback–Leibler, Jensen–Shannon, or the Earth-Movers divergence between output distribution and true distribution, if that distribution is known. In this work, since the values found on the edges are categorical we use the multiclass Cross Entropy.

$$\text{loss}(x, \text{class}) = -\log \left(\frac{\exp(x[\text{class}])}{\sum_j \exp(x[j])} \right) \quad (3.4)$$

$$= -x[\text{class}] + \log \left(\sum_j \exp(x[j]) \right) \quad (3.5)$$

We train our method using the ADAM optimizer [33] with standard hyper-parameters. To give clarity, and elucidate the connection and differences between CNPs and GNPs, we present the training algorithm for Graph Neural Processes in Algorithm (1)

3.5 Applications

To show the efficacy of GNPs on real world examples, we conduct a series of experiments on a collection of 16 graph benchmark datasets [31]. These datasets span a variety of application areas, and have a diversity of sizes; features of the explored datasets are summarized in Table 3.2. While the benchmark collection has more than 16 datasets, we selected only those that have both node and edge labels that are fully known; these were then artificially sparsified to create graphs with known edges, but unknown edge labels.

We use the model described in Section (4) and the algorithm presented in algorithm (1) and train a GNP to observe $p \cdot n$ context points where $p \in [.4, .9]$ and $n = |E|$. We

Algorithm 1 Graph Neural Processes. All experiments use $n_{\text{epochs}} = 10$ and default Adam optimizer parameters

Require: p_0 , lower bound percentage of edges to sample as context points. p_1 , corresponding upper bound. m , size of slice (neighborhood) of local structural eigenfeatures.

Require: θ_0 , initial encoder parameters. ϕ_0 initial decoder parameters.

```

1: Let  $X$  input graphs
2: for  $t = 0, \dots, n_{\text{epochs}}$  do
3:   for  $x_i$  in  $X$  do
4:     Sample  $p \leftarrow \text{unif}(p_0, p_1)$ 
5:     Assign  $n_{\text{context points}} \leftarrow p \cdot |\text{Edges}(x_i)|$ 
6:     Sparsely Sample  $x_i^{cp} \leftarrow x_i|_{n_{\text{context points}}}$ 
7:     Compute degree and adj matrix  $D, A$  for graph  $x_i^{cp}$ 
8:     Compute  $L \leftarrow I_{n_{\text{context points}}} - D^{-\frac{1}{2}} A D^{-\frac{1}{2}}$ 
9:     Define  $F^{cp}$  as empty feature matrix for  $x_i^{cp}$ 
10:    Define  $F$  as empty feature matrix for full graph  $x_i$ 
11:    for edge  $k$  in  $x_i$  do
12:      Extract eigenfeatures  $\Lambda|_k$  from  $L$ , see eq (3.3)
13:      Concatenate  $[\Lambda|_k; v_k; u_k; d(v_k); d(u_k)]$  where  $v_k, u_k$  are the attribute values at the
      node, and  $d(v_k), d(u_k)$  the degree at the node.
14:      if edge  $k \in x_i^{cp}$  then
15:        Append features for context point to  $F^{cp}$ 
16:      end if
17:      Append features for all edges to  $F$ 
18:    end for
19:    Encode and aggregate  $r_C \leftarrow h_{\theta}(F^{cp})$ 
20:    Decode  $\tilde{x}_i \leftarrow g_{\phi}(F|r_C)$ 
21:    Calculate Loss  $l \leftarrow \mathcal{L}(\tilde{x}_i, x_i)$ 
22:    Step Optimizer
23:  end for
24: end for

```

compare GNPs with naive baselines and a strong random forest model for the task of edge imputation. The baselines are described in detail in [27]:

- **Random:** a random edge label is imputed for each unknown edge
- **Common:** the most common edge label is imputed for each unknown edge
- **Random forest:** from scikit-learn; default hyperparameters

For each algorithm on each dataset, we calculate weighted precision, weighted recall, and weighted F1-score (with weights derived from class proportions), averaging each algorithm over multiple runs. Statistical significance was assessed with a two-tailed t-test with $\alpha = 0.05$.

3.5.1 Results

The results are pictured in Figures (3.2 - 3.4); all of the data is additionally summarized in Table 3.1. We first note several high-level results, and then look more deeply into the results on different subsets of the overall dataset.

First, we note that the GNP provides the best F1-score on 14 out of 16 datasets, and best recall on 14 out of datasets (in this case, recall is equivalent to classification accuracy). By learning a high-level abstract representation of the data, along with a conditional distribution over edge values, the Graph Neural Process is able to perform well at this edge imputation task, beating both naive and strong baselines. The GNP is able to do this on datasets with as few as about 300 graphs, or as many as about 9000. We also note that the GNP is able to overcome class imbalance.

AIDS. The AIDS Antiviral Screen dataset [49] is a dataset of screens checking tens of thousands of compounds for evidence of anti-HIV activity. The available screen results are chemical graph-structured data of these various compounds. A moderately sized dataset, the GNP beats the RF by 7%.

bzr,cox2,dhfr,er. These are chemical compound datasets BZR, COX2, DHFR and ER which come with 3D coordinates, and were used by [44] to study the pharmacophore kernel. Results are mixed between algorithms on these datasets; these are the datasets that have an order of magnitude more edges (on average) than the other datasets. These datasets have a large class imbalance: predicting the most common edge label yields around 90% accuracy. For example, the bzr dataset has 61,594 in class 1, and 7,273 in the other 4 classes combined. Even so, the GNP yields best F1 and recall on 2/4 of these; random forest gives

the best precision on 3/4. It may simply be that there is so much data to work with, and the classes are so imbalanced, that it is hard to find much additional signal.

mutagenicity, MUTAG. The MUTAG dataset [13] consists of 188 chemical compounds divided into two classes according to their mutagenic effect on a bacterium. While the mutagenicity dataset [29] is a collection of molecules and their interaction information with in vitro.

Here, GNP beats the random forest by several percent; both GNP and RF are vastly superior to naive baselines.

PTC_*. The various Predictive Toxicology Challenge [24] datasets consist of several hundred organic molecules marked according to their carcinogenicity on male and female mice and rats. [24] On the PTC family of graphs, GNP bests random forests by 10-15% precision, and 3-10% F1-score; both strongly beat naive baselines.

Tox21_*. This data consists of 10,000 chemical compounds run against human nuclear receptor signaling and stress pathway, it was originally designed to look for structure-activity relationships and improve overall human health. [15] On the Tox family of graphs, the GNP strongly outperforms all other models by about 20% precision; about 12% F1; and about 10% recall.

3.6 Areas for Further Exploration

This introduction of GNPs is meant as a proof of concept to encourage further research into Bayesian graph methods. As part of this work, we list a number of problems where GNPs could be applied. **Visual scene understanding** [48] [52] where a graph is formed through some non-deterministic process where a collection of the inputs may be corrupted, or inaccurate. As such, a GNP could be applied to infer edge or node values in the scene to improve downstream accuracy. **Few-shot learning** [53] where there is hidden structural information. A method like [30] could be used to discover the form, and a GNP could then be leveraged to impute other graph attributes.

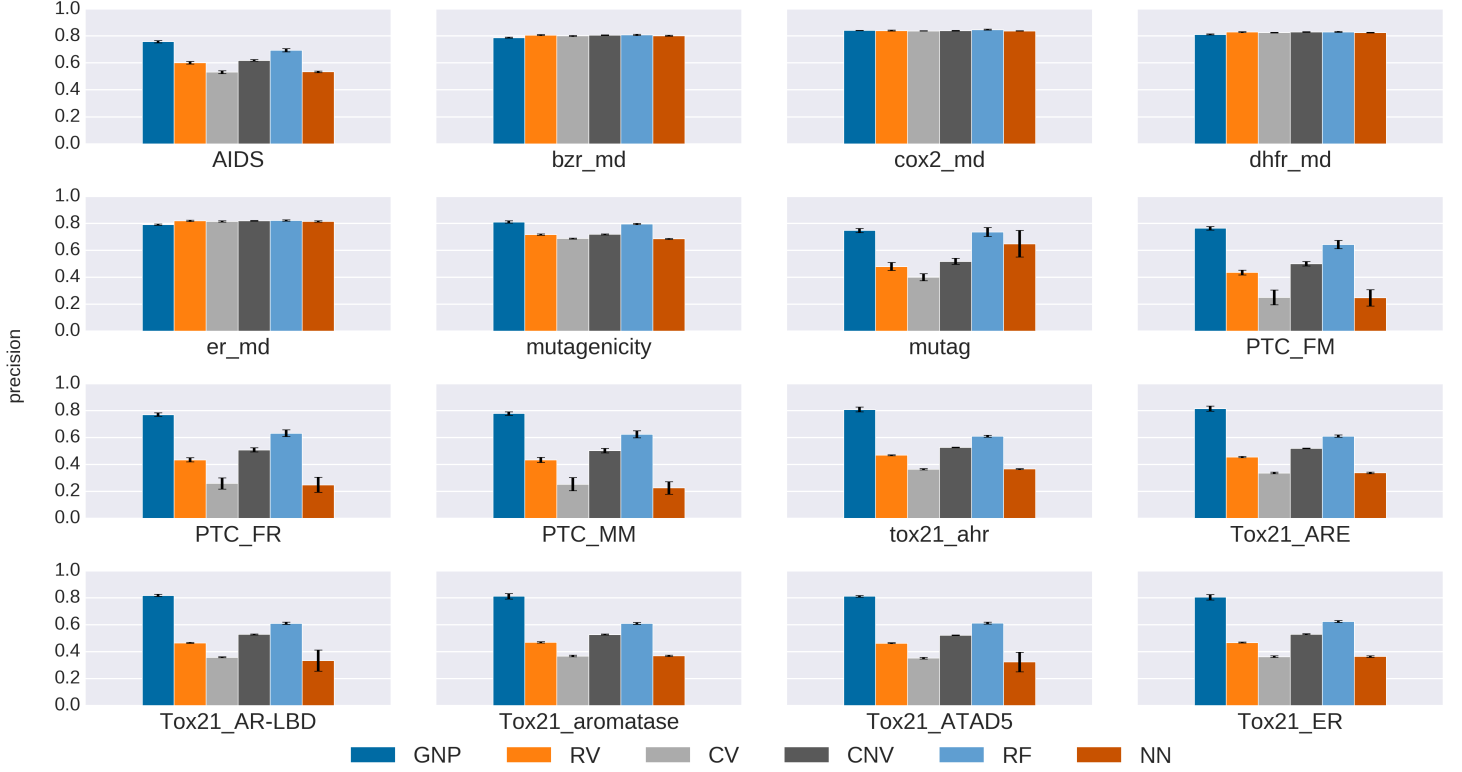


Figure 3.2: Experimental precision graph compared with baselines, we see our method performs achieves a $\sim .2$ higher precision on average.

Learning dynamics of physical systems [4] [10] [61] [58] [50] with gaps in observations over time, where the GNP could infer values involved in state transitions.

Traffic prediction on roads or waterways [41] [11] or throughput in cities. The GNP would learn a conditional distribution over traffic between cities

Multi-agent systems [56] [26][35] where you want to infer internal agent state of competing or cooperating agents. GNP inference could run in conjunction with other multi-agent methods and provide additional information from graph interactions.

Natural language processing in the construction of knowledge graphs [8] [46] [23] by relational infilling or reasoning about connections on a knowledge graph. Alternatively, they could be used to perform **semi-supervised text classification** [36] by imputing relations between words and sentences.

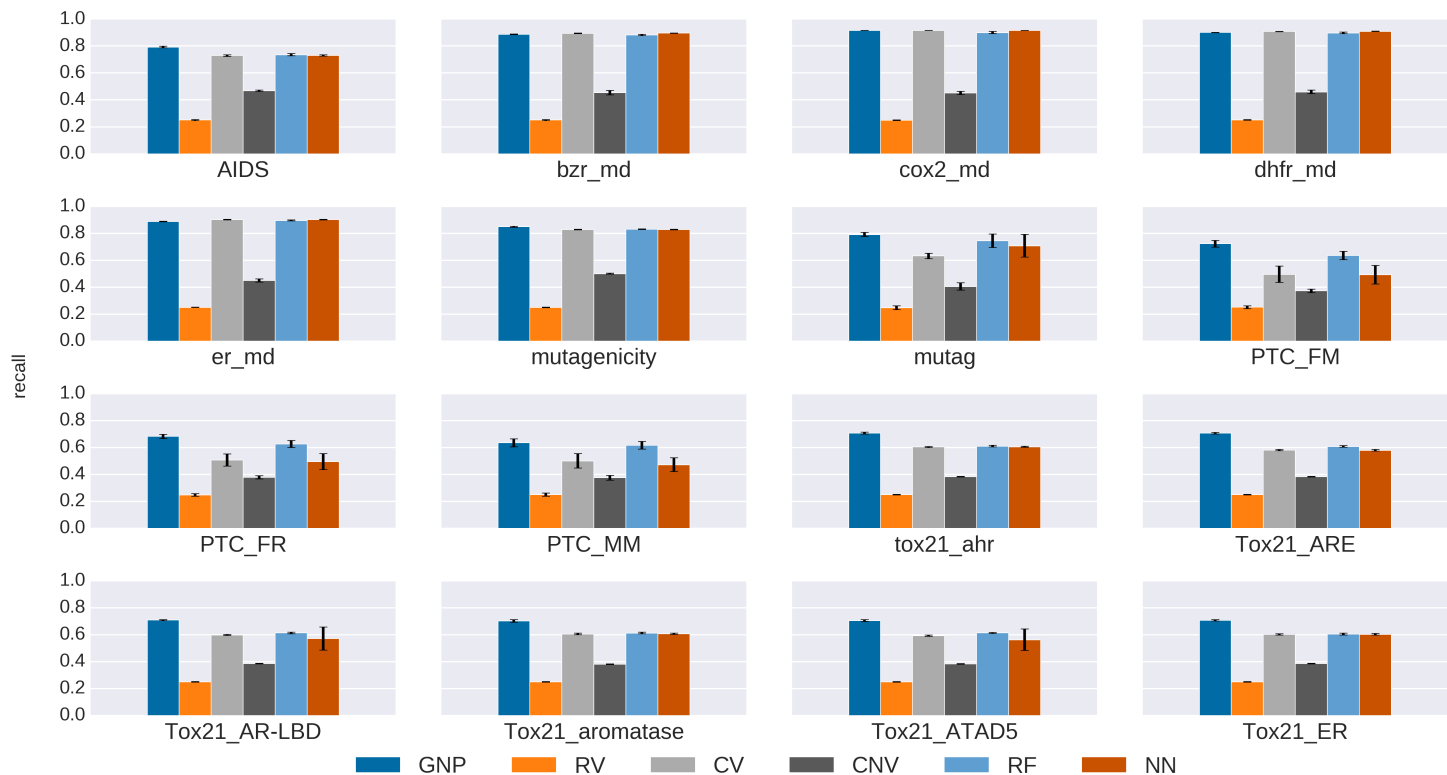


Figure 3.3: Experimental recall graph compared with baselines

There are a number of computer vision applications where graphs and GNPs could be extremely valuable. For example, one could improve **3D meshes and point cloud** [60] construction using lidar or radar during tricky weather conditions. The values in the meshes or point clouds could be imputed directly from conditional draws from the distribution learned by the GNP.

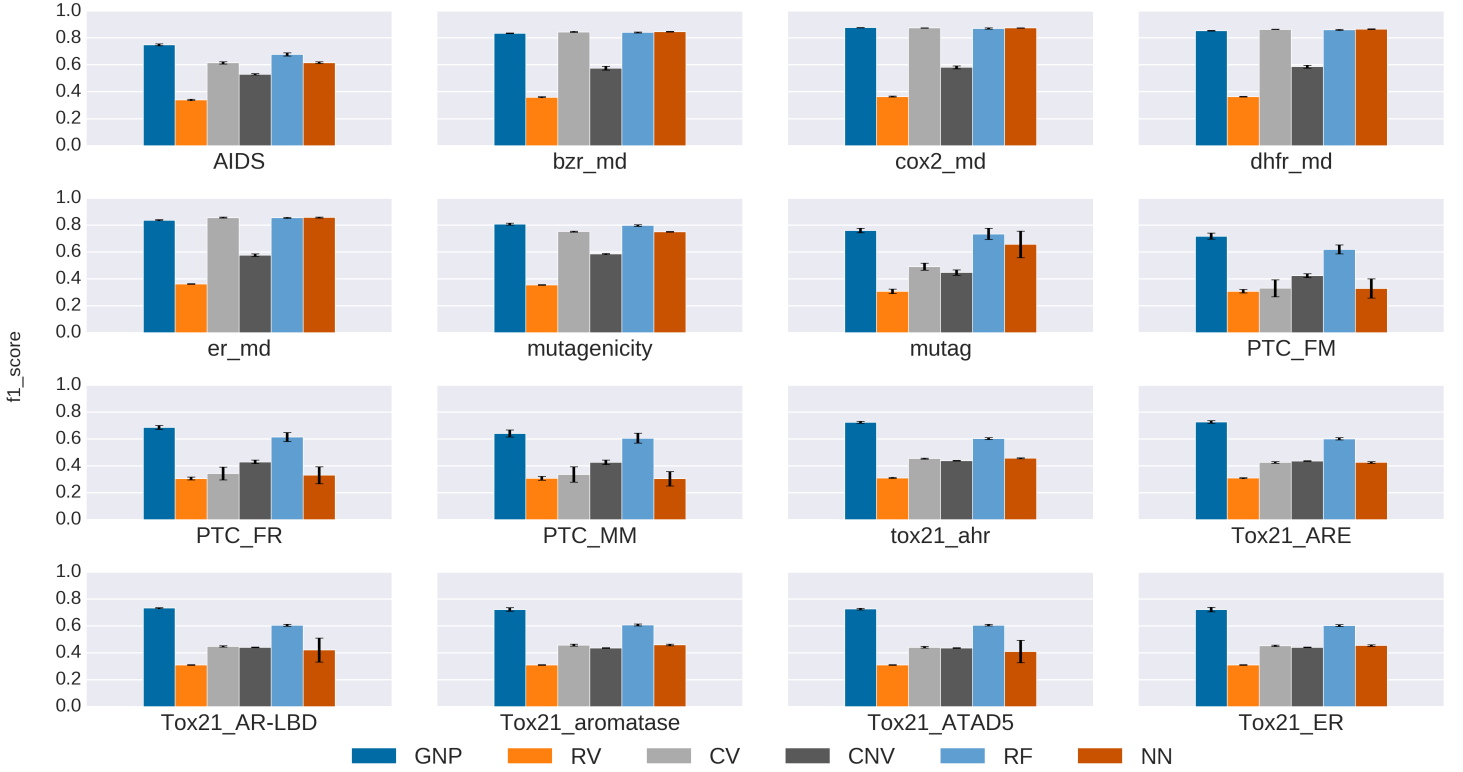


Figure 3.4: Experimental F1-score graph compared with baselines.

| Dataset | RF | | | RV | | | CV | | | GNP | | |
|-----------------|-------------|-------------|-------------|-------------|------|------|------|-------------|-------------|-------------|-------------|-------------|
| | P | R | F1 | P | R | F1 | P | R | F1 | P | R | F1 |
| AIDS | 0.69 | 0.74 | 0.68 | 0.60 | 0.25 | 0.34 | 0.53 | 0.73 | 0.61 | 0.76 | 0.79 | 0.75 |
| b2r_md | 0.81 | 0.88 | 0.84 | 0.81 | 0.25 | 0.36 | 0.80 | 0.89 | 0.84 | 0.79 | 0.89 | 0.83 |
| cox2_md | 0.85 | 0.90 | 0.87 | 0.84 | 0.25 | 0.36 | 0.84 | 0.91 | 0.87 | 0.84 | 0.92 | 0.88 |
| dhfr_md | 0.83 | 0.90 | 0.86 | 0.83 | 0.25 | 0.36 | 0.82 | 0.91 | 0.86 | 0.81 | 0.90 | 0.85 |
| er_md | 0.82 | 0.90 | 0.85 | 0.82 | 0.25 | 0.36 | 0.81 | 0.90 | 0.86 | 0.79 | 0.89 | 0.84 |
| mutagenicity | 0.79 | 0.83 | 0.80 | 0.72 | 0.25 | 0.35 | 0.69 | 0.83 | 0.75 | 0.81 | 0.85 | 0.81 |
| mutag | 0.73 | 0.74 | 0.73 | 0.48 | 0.25 | 0.31 | 0.40 | 0.63 | 0.49 | 0.75 | 0.79 | 0.76 |
| PTC_FM | 0.64 | 0.64 | 0.62 | 0.43 | 0.25 | 0.31 | 0.25 | 0.50 | 0.33 | 0.76 | 0.72 | 0.72 |
| PTC_FR | 0.63 | 0.63 | 0.61 | 0.43 | 0.25 | 0.31 | 0.26 | 0.51 | 0.34 | 0.77 | 0.68 | 0.69 |
| PTC_MM | 0.62 | 0.62 | 0.61 | 0.43 | 0.25 | 0.31 | 0.25 | 0.50 | 0.34 | 0.78 | 0.64 | 0.64 |
| tox21_ahr | 0.61 | 0.61 | 0.60 | 0.47 | 0.25 | 0.31 | 0.36 | 0.60 | 0.45 | 0.81 | 0.71 | 0.72 |
| Tox21_ARE | 0.61 | 0.61 | 0.60 | 0.46 | 0.25 | 0.31 | 0.34 | 0.58 | 0.43 | 0.81 | 0.71 | 0.73 |
| Tox21_AR-LBD | 0.61 | 0.61 | 0.61 | 0.47 | 0.25 | 0.31 | 0.36 | 0.60 | 0.45 | 0.82 | 0.71 | 0.73 |
| Tox21_aromatase | 0.61 | 0.61 | 0.61 | 0.47 | 0.25 | 0.31 | 0.37 | 0.61 | 0.46 | 0.81 | 0.70 | 0.72 |
| Tox21_ATAD5 | 0.61 | 0.61 | 0.61 | 0.46 | 0.25 | 0.31 | 0.35 | 0.59 | 0.44 | 0.81 | 0.71 | 0.73 |
| Tox21_ER | 0.62 | 0.60 | 0.60 | 0.47 | 0.25 | 0.31 | 0.36 | 0.60 | 0.45 | 0.80 | 0.71 | 0.72 |

Table 3.1: Experimental Results. RF=Random Forest; RV=Random edge label; CV=most common edge label; GNP=Graph Neural Process. P=precision; R=recall; F1=F1 score. Statistically significant bests are in bold with non-significant ties bolded across methods.

| Dataset | $ X $ | $ N $ | $ E $ | $ \cup \{e_k\} $ |
|-----------------|-------|-------|--------|------------------|
| AIDS | 2000 | 15.69 | 16.20 | 3 |
| BZR_MD | 306 | 21.30 | 225.06 | 5 |
| COX2_MD | 303 | 26.28 | 335.12 | 5 |
| DHFR_MD | 393 | 23.87 | 283.01 | 5 |
| ER_MD | 446 | 21.33 | 234.85 | 5 |
| Mutagenicity | 4337 | 30.32 | 30.77 | 3 |
| MUTAG | 188 | 17.93 | 19.79 | 4 |
| PTC_FM | 349 | 14.11 | 14.48 | 4 |
| PTC_FR | 351 | 14.56 | 15.00 | 4 |
| PTC_MM | 336 | 13.97 | 14.32 | 4 |
| Tox21_AHR | 8169 | 18.09 | 18.50 | 4 |
| Tox21_ARE | 7167 | 16.28 | 16.52 | 4 |
| Tox21_aromatase | 7226 | 17.50 | 17.79 | 4 |
| Tox21_ARLBD | 8753 | 18.06 | 18.47 | 4 |
| Tox21_ATAD5 | 9091 | 17.89 | 18.30 | 4 |
| Tox21_ER | 7697 | 17.58 | 17.94 | 4 |

Table 3.2: Features of the explored data sets

Chapter 4

Wasserstein Neural Processes

4.1 Introduction

Gaussian Processes (GPs) are analytically tractable, nonparametric probabilistic models that are widely used for nonlinear regression problems. They map an input $\vec{x}_i \in \mathbb{R}^{d_x}$ to an output $\vec{y}_i \in \mathbb{R}^{d_y}$ by defining a collection of jointly Gaussian conditional distributions. These distributions can be conditioned on an arbitrary number of context points $X_C := (\vec{x}_i)_{i \in C}$ and their associated outputs Y_C . Given the context points, an arbitrary number of target points $X_T := (\vec{x}_i)_{i \in T}$, with their outputs Y_T , can be modeled using the conditional distribution $P(X_T|X_C, Y_C)$. This modeling is invariant to the ordering of context points, and to the ordering of targets. Formally, they define a nonparametric distribution over smooth functions, and leverage the almost magical nature of Gaussian distributions to both solve regression tasks and provide confidence intervals about the solution. However, GPs are computationally expensive, and the strong Gaussian assumptions they make are not always appropriate.

Neural Processes (NPs) can be loosely thought of as the neural generalization of GPs. They blend the strengths of GPs and deep neural networks: like GPs, NPs are a class of models that learn a mapping from a context set of input-output pairs to a distribution over functions. Also like GPs, NPs can be conditioned on an arbitrary set of context points, but unlike GPs (which have quadratic complexity), NPs have linear complexity in the number of context points [18]. NPs inherit the high model capacity of deep neural networks, which gives them the ability to fit a wide range of distributions, and can be blended with the latest

deep learning best practices (for example, NPs were recently extended using attention [32] to overcome a fundamental underfitting problem which lead to inaccurate function predictions).

Importantly, these NPs still make Gaussian assumptions; furthermore, the training algorithm of these Neural Processes uses a loss function that relies on maximum-likelihood and the KL divergence metric. Recently in the machine learning community there has been much discussion [25] around generative modeling, and the use of Wasserstein divergences. Wasserstein distance is a desirable tool because it can measure distance between two probability distributions, even if the support of those distributions is disjoint. The KL divergence, however, is undefined in such cases (or infinity); this corresponds to situations where the data likelihood is exactly zero, which can happen in the case of misspecified models.

This begs the question: can we simultaneously improve the maximum likelihood training of NPs, and replace the Gaussian assumption in NPs with a more general model - ideally one *where we do not need to make any assumptions about the data likelihood*?

More formally, it is known that in the limit of data maximizing the expected log probability of the data is equivalent to minimizing the KL divergence between the distribution implied by a model P_θ and the true data distribution Q , which may not be known:

$$\operatorname{argmin}_{\theta} D_{KL}(Q|P_\theta) = \operatorname{argmin}_{\theta} E_Q[-\log(P_\theta)] \quad (4.1)$$

However, if the two distributions Q and P_θ have disjoint support (i.e., if even one data point lies outside the support of the model), the KL divergence will be undefined (or infinite) due to the evaluation of $\log(0)$. Therefore, it has been shown recently [1] that substituting Wasserstein for KL in (eq 4.1) gives stable performance for non-overlapping distributions in the same space. This leads to Wasserstein distance as a proxy for the KL divergence interpretation of maximum likelihood learning.

This new framing is valuable in the misspecified case where data likelihood is zero under the model and additionally for the case when likelihood is unknown (such as in

high dimensional data supported on low-dimensional manifolds, such as natural images), or intractable. However, in all cases, it is required that we can draw samples from the target distribution.

The central contribution of this work is a powerful non-linear class of Wasserstein Neural Processes that can learn a variety of distributions when the data are misspecified under the model, and when the likelihood is unknown or intractable. We do so by replacing the traditional maximum likelihood loss function with a Wasserstein divergence. We discuss computationally efficient approximations (including Sliced Wasserstein Distance) and evaluate performance in several experimental settings.

4.2 Wasserstein Neural Processes

Our objective in combining Optimal Transport with Neural Processes is to model a desirable class of functions that are either misspecified, or with computationally intractable likelihood. We choose to use a CNP (fig 2.1) as our Neural Process architecture because CNPs are the simplest variant of the NP family with which we can illustrate the benefits of Wasserstein distance. Additionally, due to model simplicity, we can better ablate the performance and assign proper credit to the inclusion of Wasserstein distance. Extending WNP to use a latent path, or attention, would be an interesting direction for future research.

To be precise, we use a parameterized Neural Network of shape and size determined according to task specifications. The input context x, y points are encoded and aggregated, as in traditional CNPs, using a commutative operation into an r_C of fixed length. The desired target x points (typically the entire set of points) are then conditioned (using concatenation) and decoded to produce a synthetic set of target y points. At this point, Wasserstein Neural Processes deviate from traditional NPs. Typically, as mentioned above, the target y likelihood would be calculated according to the output distribution (which we do as a comparison in the experiments). However, in our case, the decoded target y points are samples from an implicitly modeled distribution. As such, since the likelihood may be difficult to calculate, or

non-existent, WNP use sliced Wasserstein distance as a learning signal. This calculation is put forth in Algorithm 3.

The sliced Wasserstein distance gives the benefit of scalable distance computation between two potentially disjoint distributions. Therefore, WNP can model a larger variety of distributions than traditional NPs as discussed above.

Algorithm 2 Wasserstein Neural Processes.

Require: p_0 , lower bound percentage of edges to sample as context points. p_1 , corresponding upper bound.

Require: θ_0 , initial encoder parameters. ϕ_0 initial decoder parameters.

```

1: Let  $X$  input data in  $\mathbb{R}^d$ 
2: for  $t = 0, \dots, n_{\text{epochs}}$  do
3:   for  $x_i$  in  $X$  do
4:     Sample  $p \leftarrow \text{unif}(p_0, p_1)$ 
5:     Assign  $n_{\text{context points}} \leftarrow p \cdot |\text{Target Points}|$ 
6:     Sparsely Sample  $x_i^{cp} \leftarrow x_i|_{n_{\text{context points}}}$  and associated  $y_i^{cp}$ 
7:     Obtain one-hot encoding of  $x_i^{cp}$  points and concatenate with  $y_i^{cp}$  as  $F^{cp}$ 
8:     Encode and aggregate  $r_C \leftarrow h_\theta(F^{cp})$ 
9:     Decode  $\tilde{x}_i \leftarrow g_\phi(F|r_C)$ 
10:    Calculate Sliced Wasserstein Distance  $l \leftarrow \mathcal{W}(\tilde{x}_i, x_i)$ 
11:    Step Optimizer
12:   end for
13: end for
```

Algorithm 3 Sliced Wasserstein Distance

Require: $n \in [5, 50]$, number of desired projections.

Require: d , dimension of embedding space.

Require: p , power of desired Wasserstein distance.

```

1: Let  $X$  be the empirical synthetic distribution and  $Y$  be samples from the true distribution
2: Sample projections  $P$  from  $S^d$  i.e.,  $\mathcal{N}(\text{size} = (n, d))$  and normalize using dimension-wise  $L_2$ 
3: Project  $\hat{X} = X \cdot P^T$ 
4: Project  $\hat{Y} = Y \cdot P^T$ 
5: Sort  $\hat{X}^T$  and  $\hat{Y}^T$ 
6: Return  $(\hat{X}^T - \hat{Y}^T)^p$ 
```

4.3 Experiments

As discussed in Section 4.1, traditional Neural Processes are limited by their use of the KL divergence, both explicitly and implicitly. In this section, we demonstrate that Wasserstein Neural Processes can learn effectively under conditions that cause a traditional NP to fail. The first two experiments illustrate NPs’ fundamental failing of relying on the likelihood function for learning. The final experiment shows the ability of WNPs to work on larger scale tasks.

4.3.1 Misspecified Models - Linear Regression with Uniform Noise

In standard linear regression, parameters $\beta = (\beta_1, \beta_2, \dots, \beta_n)^T$ are estimated in the model

$$\hat{Y} = \beta x + b + \eta \tag{4.2}$$

where a stochastic noise model η is used to model disturbances. Traditionally, $\eta \sim \mathcal{N}(\mu, \sigma)$ for some parameterized Normal distribution. In such a traditional case, Neural Processes can learn the a conditional distribution over the data x . This is because all the data has a calculable likelihood under any setting of the parameters.

However, consider the case where the noise model is Uniform (e.g., $Unif[-1, 1]$). There are now settings for β, b under which the data x would have exactly zero likelihood (i.e., any time a single data point falls outside the “tube” of noise). In the worst case, there may be *no* setting of the parameters that provides non-zero likelihood. Furthermore, because of the uniform nature of the noise likelihoods, there are no useful gradients.

$$L(\beta, b) = \Pi_{i \in I} p(y_i | x_i, \beta, b) \tag{4.3}$$

$$= \Pi_{i \in I} 0.5 \cdot 1(|y_i - (\beta x + b)| < 1) \tag{4.4}$$

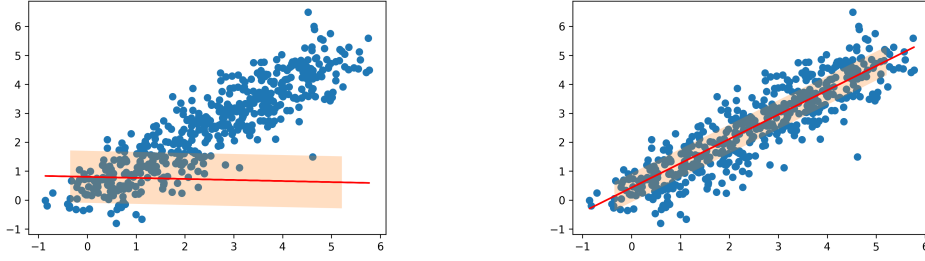


Figure 4.1: Left) Failed NP regression with a uniform noise model $Unif[-1, 1]$ and 500 synthetic data points. The line represents the estimated line of best fit, while the tube represents the uncertainty of the model. Since the model is misspecified, the likelihood (eq 4.4) is zero and there is no learning signal for parameter updates. Right) Successful WNP regression with an identical experimental setup. In this case, the model does not rely on likelihood, instead it makes use of the signal from the sliced Wasserstein distance which persists even in the misspecified case.

In these cases, any model based on likelihoods, including NPs, would fail to learn. This is illustrated in fig 4.1.

However, even if the data likelihood is degenerate, there is still a well-defined notion of the optimal transport cost between the observed data and the distribution implied by the model, no matter what setting of parameters is used. Since the Wasserstein distance is defined independent of likelihood, the WNP model still finds the proper fit of parameters given the data as seen in (fig 4.1). This experiment could be analogous to the real world case where our prior knowledge of the data is fundamentally flawed. This flaw implies that our model is incorrect and misspecified. For this experiment, we generate synthetic data from $y = 1 \cdot x + 0 + \epsilon$ where $\epsilon \sim \mathcal{N}(0, 0.5)$. The sliced Wasserstein distance is calculated as described in Section 1.4.3 with $N = 50$. We let h and g be two-layer fully connected Neural Networks and our r_C vector is of size 32 and we use the Adam [34] optimizer with defaults.

4.3.2 Intractable Likelihoods - The Quantile “ g -and- κ ” Distribution

We now present a second example, with a distribution that is easy to sample from, but for which the likelihood cannot be tractably computed. The g -and- κ distribution [28, 57] is

defined for $r \in (0, 1)$ and parameters $\theta = (a, b, g, \kappa) \in [0, 10]^4$ as:

$$a + b \left(1 + 0.8 \frac{1 - \exp(-gz(r))}{1 + \exp(-gz(r))} \right) (1 + z(r)^2)^\kappa z(r) \quad (4.5)$$

Where $z(r)$ refers to the r -th quantile of a standard Normal distribution $\mathcal{N}(0, 1)$. It is known that the density function and likelihood are analytically intractable, making it a perfect test for WNP. This model is a standard benchmark for approximate Bayesian computation methods [55]. Despite the intractability, it is straightforward to sample i.i.d. variables from this distribution by simply plugging in standard Normals in place of $z(r)$ [6].

We follow the experimental set up in [6] and take $\theta = (3, 1, 2, 0.5)$. Which produces a distribution seen in (fig 4.2). Additionally, we let h and g be two-layer fully connected neural networks with slightly more parameters than in the uniform noise model experiment and with the same size r_C vector (32 dimensional). In this case, we use a cyclic learning rate with the base learning rate $1e^{-3}$ and the max learning rate $1e^{-2}$. We found that in some cases, our model would output a shifted empirical distribution, and cycling the learning rate would finish the optimization process well. We hypothesize this as a function of the g -and- κ topology and leave it as an area for future exploration as it is beyond the scope of this work.

Due to the fact that the likelihood is intractable, if one wished to use this distribution as the model output of a Neural Process, it would be as impossible to train as in the uniform noise model regression task. This is due to the fact, as mentioned above, that standard Neural Processes rely on maximum likelihood as a learning signal. Since Wasserstein Neural Processes rely solely on the divergence between two empirical distributions, they can learn a conditional g -and- κ distribution.

4.3.3 CelebA tiles as super pixels

We now present our final experiment on high-dimensional image data. The purpose of this experiment is to illustrate the ability of WNP to model higher dimensional distributions where

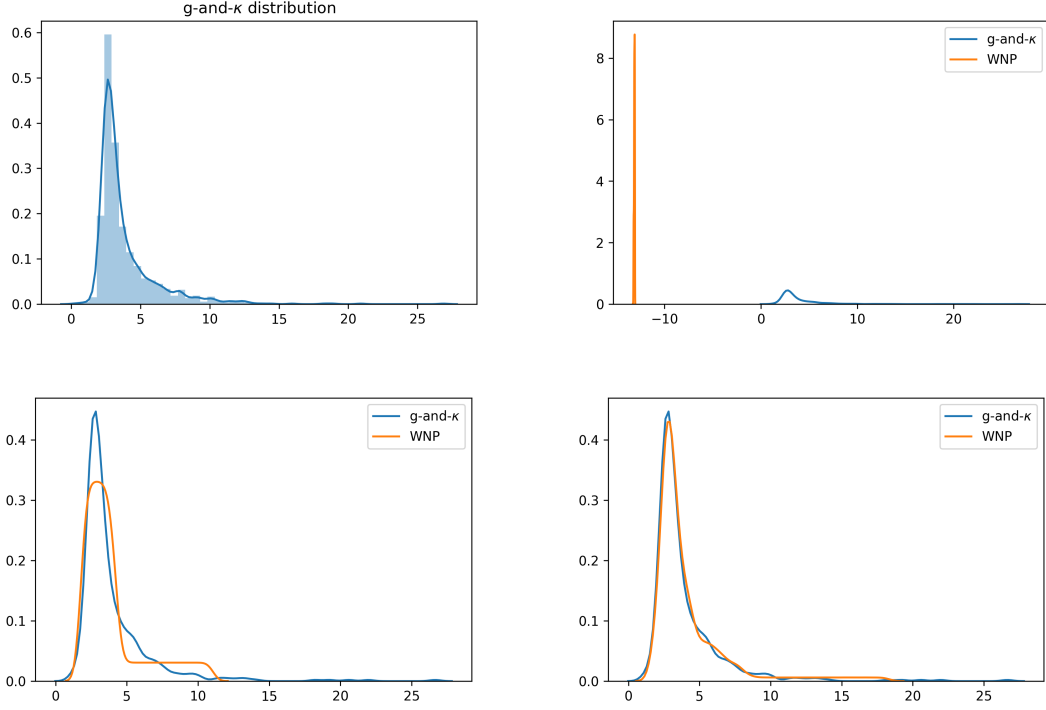


Figure 4.2: Top left: Initialization of learned distribution. It would be impossible to calculate the likelihood in this case for traditional NPs both because the model is misspecified and because the likelihood is computationally intractable. Top right: Beginning of the learning process. The output of the WNP is completely disjoint from the sampled g -and- κ distribution. Bottom left: Part way through the learning process (approximately 500 steps) we see the WNP is quickly able to achieve a good fit. Bottom right: Final fit of g -and- κ distribution where WNP learned to model the distribution without relying on the likelihood calculation

the likelihood might be difficult to calculate, or potentially meaningless. The experiment is run on the CelebA dataset.

To distinguish between a large class of prior work and our contribution, we highlight a distinction between common candidate objective functions and the true data distribution. Many frameworks for image prediction, including standard Neural Processes implicitly, involve *reconstruction-based* loss functions that are usually defined as $\|y - \hat{y}\|$. These loss functions compare a target y with the model’s output \hat{y} using pixel-wise MSE – but this is known to perform poorly in many cases, including translations and brightness variations. More generally, it completely ignores valuable higher-order information in images: pixel intensities are neither independent nor do they share the same variance. Similarly, it is well-known that images are typically supported on low-dimensional manifolds.

Put simply, pixel-wise losses imply false assumptions and results in a different objective than the one we would truly like to minimize. These false assumptions help explain why (for example) generative image models such as VAEs (based on reconstruction error) result in blurry reconstructions, while models such as Wasserstein GANs (based on Wasserstein divergences) result in much sharper reconstructions.

Neural Processes use Gaussian distributions over individual pixels. Predictably, this results in blurry predictions, seen in [17, 18, 32]. Analogously to the difference between GANs and VAEs, if we extend these experiments slightly, where instead of using pixel values as context points, we take large tiles of the image as context points, we might expect that we would see sharper reconstructions from WNPs as compared to NPs.

We test this by treating image completion as a high-dimensional regression task. 32x32 images are broken into 4x4 tiles; these tiles are then randomly subsampled. The X inputs are the tile coordinates; the Y outputs are the associated pixels. Given a set of between 4-16 tiles, the WNP must predict the remainder of the pixels by predicting entire tiles. The Neural Process is trained on the same data, using the same neural network; the only difference is the loss function.

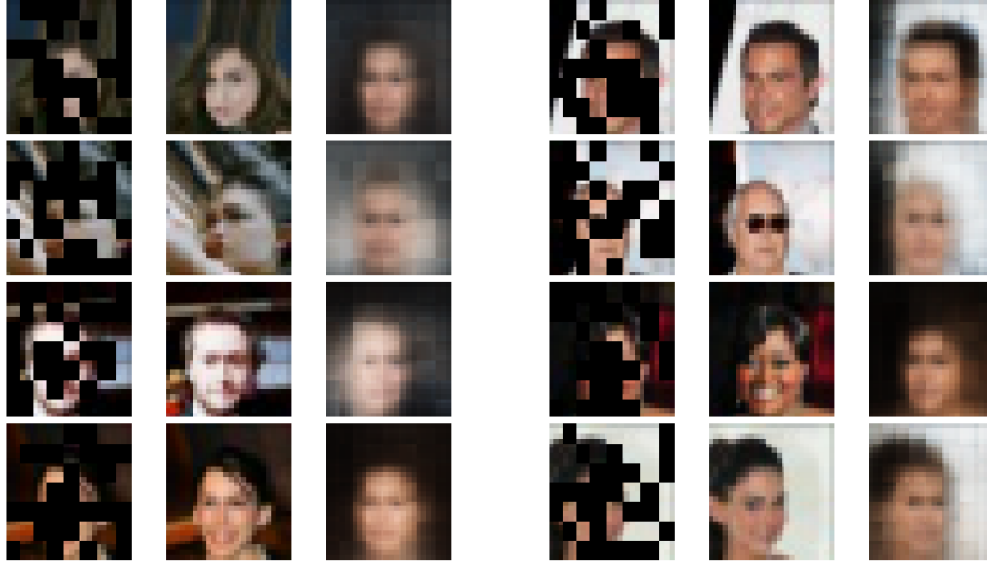


Figure 4.3: On the left: results from Neural Processes. The left-most column shows the sparse observations; the middle column shows ground truth; the right column shows the NP reconstruction. On the right: results for WNPs. See text for discussion.

4.3.4 Discussion

The results of our Celeb-A experiment are shown in (fig 4.3). Here, the results are what we expected: traditional Neural Processes fail to learn more than an extremely blurry output when using tiles as context points instead of pixels. This is because the likelihood is now significantly less meaningful. However, in the case of WNP, the empirical sliced Wasserstein distance is still well defined and can learn a better conditional distribution. While the images are not as high quality as much recent GAN work, the purpose of this experiment was to illustrate that the WNP can better capture background, facial structure, color, than traditional NP trained with maximum likelihood. Qualitatively, we see that the images are slightly more clear and better capture the background and structural information of the image, given the context conditioning

Chapter 5

Discussion and Future Work

In this document, we explored two specific extensions to the neural process framework. Neural processes allow us to learn a distribution of functions over a set of context and target points. By extending those to graph valued data and by using optimal transport, we have expanded the class of functions that can be learned.

5.1 Graph Neural Processes

In this work we have introduced Graph Neural Processes, a model that learns a conditional distribution while operating on graph structured data. This model has the capacity to generate uncertainty estimates over outputs, and encode prior knowledge about input data distributions. We have demonstrated GNP’s ability on edge imputation and given potential areas for future exploration.

While we note the encoder and decoder architectures can be extended significantly by including work from modern deep learning architectural design, this work is a step towards building Bayesian neural networks on arbitrarily graph structured inputs. Additionally, it encourages the learning of abstractions about these structures. In the future, we wish to explore the use of GNPs to inform high-level reasoning and abstraction about fundamentally relational data.

5.2 Wasserstein Neural Processes

In this work, we have also explored a synthesis of Neural Processes and Optimal Transport. We demonstrated that there are desirable classes of functions which Neural Processes either struggle to learn, or cannot learn, but which can be learned by Wasserstein Neural Processes. These WNP use Wasserstein distance in place of the intrinsic KL divergence used in maximum likelihood, and can be tractably implemented using (for example) sliced approximations. Maximum likelihood fails to learn in the case of misspecified models and likelihoods that either don't exist, or are computationally intractable. As a concluding thought, we note that our technical strategy in developing this model was to replace a maximum-likelihood objective with a Wasserstein based objective. This raises the intriguing question: can *all* probabilistic models based on maximum-likelihood be reformulated in terms of Wasserstein divergence? We leave this exciting direction for future work.

References

- [1] Martin Arjovsky, Soumith Chintala, and Léon Bottou. Wasserstein GAN. *arXiv e-prints*, art. arXiv:1701.07875, Jan 2017.
- [2] Tom Auld, Andrew W Moore, and Stephen F Gull. Bayesian neural networks for internet traffic classification. *IEEE Transactions on neural networks*, 18(1):223–239, 2007.
- [3] Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. Neural machine translation by jointly learning to align and translate. *arXiv preprint arXiv:1409.0473*, 2014.
- [4] Peter Battaglia, Razvan Pascanu, Matthew Lai, Danilo Jimenez Rezende, et al. Interaction networks for learning about objects, relations and physics. In *NIPS*, pages 4502–4510, 2016.
- [5] Peter W. Battaglia, Jessica B. Hamrick, Victor Bapst, Alvaro Sanchez-Gonzalez, Vinicius Zambaldi, Mateusz Malinowski, Andrea Tacchetti, David Raposo, Adam Santoro, Ryan Faulkner, Caglar Gulcehre, Francis Song, Andrew Ballard, Justin Gilmer, George Dahl, Ashish Vaswani, Kelsey Allen, Charles Nash, Victoria Langston, Chris Dyer, Nicolas Heess, Daan Wierstra, Pushmeet Kohli, Matt Botvinick, Oriol Vinyals, Yujia Li, and Razvan Pascanu. Relational inductive biases, deep learning, and graph networks. <https://arxiv.org/abs/1806.01261>, 2018.
- [6] Espen Bernton, Pierre E Jacob, Mathieu Gerber, and Christian P Robert. Approximate bayesian computation with the wasserstein distance. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 2019.
- [7] Nicolas Bonnotte. Unidimensional and evolution methods for optimal transportation, 2013.
- [8] Antoine Bordes, Nicolas Usunier, Alberto Garcia-Duran, Jason Weston, and Oksana Yakhnenko. Translating embeddings for modeling multi-relational data. In *NIPS*, pages 2787–2795, 2013.
- [9] Michael M Bronstein, Joan Bruna, Yann LeCun, Arthur Szlam, and Pierre Vandergheynst. Geometric deep learning: going beyond euclidean data. *IEEE Signal Processing Magazine*, 34(4):18–42, 2017.

- [10] Michael B Chang, Tomer Ullman, Antonio Torralba, and Joshua B Tenenbaum. A compositional object-based approach to learning physical dynamics. *arXiv:1612.00341*, 2016.
- [11] Zhiyong Cui, Kristian Henrickson, Ruimin Ke, and Yinhai Wang. High-order graph convolutional recurrent neural network: A deep learning framework for network-scale traffic learning and forecasting. *arXiv:1802.07007*, 2018.
- [12] Marco Cuturi. Sinkhorn distances: Lightspeed computation of optimal transport. In C. J. C. Burges, L. Bottou, M. Welling, Z. Ghahramani, and K. Q. Weinberger, editors, *Advances in Neural Information Processing Systems 26*, pages 2292–2300. Curran Associates, Inc., 2013. URL <http://papers.nips.cc/paper/4927-sinkhorn-distances-lightspeed-computation-of-optimal-transport.pdf>.
- [13] Asim Kumar Debnath, Rosa L Lopez de Compadre, Gargi Debnath, Alan J Shusterman, and Corwin Hansch. Structure-activity relationship of mutagenic aromatic and heteroaromatic nitro compounds. correlation with molecular orbital energies and hydrophobicity. *Journal of medicinal chemistry*, 34(2):786–797, 1991.
- [14] David K Duvenaud, Dougal Maclaurin, Jorge Iparraguirre, Rafael Bombarell, Timothy Hirzel, Alán Aspuru-Guzik, and Ryan P Adams. Convolutional networks on graphs for learning molecular fingerprints. In *NIPS*, pages 2224–2232, 2015.
- [15] National Center for Advancing Translation Services. Tox21 data challenge, 2014. URL <https://tripod.nih.gov/tox21/challenge/index.jsp>.
- [16] Marco Cuturi Gabriel Peyré. *Computational Optimal Transport*. 2018. URL <https://arxiv.org/pdf/1803.00567.pdf>.
- [17] M Garnelo, D Rosenbaum, CJ Maddison, T Ramalho, D Saxton, M Shanahan, YW Teh, DJ Rezende, and SM. Eslami. Conditional neural processes. *ICML*, 2018.
- [18] Marta Garnelo, Jonathan Schwarz, Dan Rosenbaum, Fabio Viola, Danilo J Rezende, SM Eslami, and Yee Whye Teh. Neural processes. In *ICML Workshop on Theoretical Foundations and Applications of Deep Generative Models*, 2018.
- [19] Aude Genevay, Gabriel Peyré, and Marco Cuturi. Learning generative models with sinkhorn divergences. *arXiv preprint arXiv:1706.00292*, 2017.
- [20] M. Gori, G. Monfardini, and F Scarselli. A new model for learning in graph domains. *IJCNN*, 2:729–734, 2005.

- [21] Ishaan Gulrajani, Faruk Ahmed, Martin Arjovsky, Vincent Dumoulin, and Aaron C Courville. Improved training of wasserstein gans. In *Advances in Neural Information Processing Systems*, pages 5767–5777, 2017.
- [22] Paul R Halmos. *Measure theory*, volume 18. Springer, 2013.
- [23] Takuo Hamaguchi, Hidekazu Oiwa, Masashi Shimbo, and Yuji Matsumoto. Knowledge transfer for out-of-knowledge-base entities: A graph neural network approach. *arXiv:1706.05674*, 2017.
- [24] Christoph Helma, Ross D. King, Stefan Kramer, and Ashwin Srinivasan. The predictive toxicology challenge 2000–2001. *Bioinformatics*, 17(1):107–108, 2001.
- [25] Avinash Hindupur. Gan zoo. URL <https://github.com/hindupuravinash/the-gan-zoo>.
- [26] Yedid Hoshen. Vain: Attentional multi-agent predictive modeling. In *NIPS*, pages 2701–2711, 2017.
- [27] Mark Huisman. Imputation of missing network data: Some simple procedures. *Journal of Social Structure*, 10, 01 2009. doi: 10.1007/978-1-4614-6170-8_394.
- [28] Martinez Jorge and Iglewicz Boris. Some properties of the tukey g and h family of distributions. *Communications in Statistics-Theory and Methods*, 13(3):353–369, 1984.
- [29] Jeroen Kazius, Ross McGuire, and Roberta Bursi. Derivation and validation of toxicophores for mutagenicity prediction. *Journal of Medicinal Chemistry*, 48(1):312–320, 2005.
- [30] Charles Kemp and Joshua B Tenenbaum. The discovery of structural form. *PNAS*, 105(31):10687–10692, 2008.
- [31] Kristian Kersting, Nils M. Kriege, Christopher Morris, Petra Mutzel, and Marion Neumann. Benchmark data sets for graph kernels, 2016. URL <http://graphkernels.cs.tu-dortmund.de>.
- [32] Hyunjik Kim, Andriy Mnih, Jonathan Schwarz, Marta Garnelo, Ali Eslami, Dan Rosenbaum, Oriol Vinyals, and Yee Whye Teh. Attentive neural processes. In *ICLR*, 2019. URL <https://openreview.net/forum?id=SkE6PjC9KX>.
- [33] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv:1412.6980*, 2014.

- [34] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings*, 2015. URL <http://arxiv.org/abs/1412.6980>.
- [35] Thomas Kipf, Ethan Fetaya, Kuan-Chieh Wang, Max Welling, and Richard Zemel. Neural relational inference for interacting systems. *arXiv:1802.04687*, 2018.
- [36] Thomas N Kipf and Max Welling. Semi-supervised classification with graph convolutional networks. *arXiv:1609.02907*, 2016.
- [37] Soheil Kolouri, Yang Zou, and Gustavo K. Rohde. Sliced wasserstein kernels for probability distributions. *CoRR*, abs/1511.03198, 2015. URL <http://arxiv.org/abs/1511.03198>.
- [38] Soheil Kolouri, Charles E. Martin, and Gustavo K. Rohde. Sliced-wasserstein autoencoder: An embarrassingly simple generative model. *CoRR*, abs/1804.01947, 2018. URL <http://arxiv.org/abs/1804.01947>.
- [39] Soheil Kolouri, Kimia Nadjahi, Umut Simsekli, Roland Badeau, and Gustavo K. Rohde. Generalized sliced wasserstein distances. *CoRR*, abs/1902.00434, 2019. URL <http://arxiv.org/abs/1902.00434>.
- [40] Y. Li, D. Tarlow, M. Brockschmidt, and R. Zemel. Gated graph sequence neural networks. *ICLR*, 2016.
- [41] Yaguang Li, Rose Yu, Cyrus Shahabi, and Yan Liu. Diffusion convolutional recurrent neural network: Data-driven traffic forecasting. *arXiv:1707.01926*, 2017.
- [42] Huidong Liu, Xianfeng Gu, and Dimitris Samaras. A two-step computation of the exact gan wasserstein distance. In *ICML*, 2018.
- [43] Christos Louizos, Xiahan Shi, Klamer Schutte, and Max Welling. The functional neural process. *arXiv preprint arXiv:1906.08324*, 2019.
- [44] Pierre Mahé, Liva Ralaivola, Véronique Stoven, and Jean-Philippe Vert. The pharmacophore kernel for virtual screening with support vector machines. *Journal of Chemical Information and Modeling*, 46(5):2003–2014, 2006.
- [45] Gaspard Monge. Mémoire sur la théorie des déblais et des remblais. *Histoire de l’Académie royale des sciences de Paris*, 1781.

- [46] Daniel Oñoro-Rubio, Mathias Niepert, Alberto García-Durán, Roberto González, and Roberto J López-Sastre. Representation learning for visual-relational knowledge graphs. *arXiv:1709.02314*, 2017.
- [47] Shenghao Qin, Jiacheng Zhu, Jimmy Qin, Wenshuo Wang, and Ding Zhao. Recurrent attentive neural process for sequential data. *arXiv preprint arXiv:1910.09323*, 2019.
- [48] David Raposo, Adam Santoro, David Barrett, Razvan Pascanu, Timothy Lillicrap, and Peter Battaglia. Discovering objects and their relations from entangled scene representations. 02 2017.
- [49] Kaspar Riesen and Horst Bunke. Iam graph database repository for graph based pattern recognition and machine learning. In *Joint IAPR International Workshops on (SPR) and (SSPR)*, pages 287–297. Springer, 2008.
- [50] Alvaro Sanchez-Gonzalez, Nicolas Heess, Jost Tobias Springenberg, Josh Merel, Martin Riedmiller, Raia Hadsell, and Peter Battaglia. Graph networks as learnable physics engines for inference and control. *arXiv:1806.01242*, 2018.
- [51] Filippo Santambrogio. Optimal transport for applied mathematicians, 2015.
- [52] Adam Santoro, David Raposo, David G Barrett, Mateusz Malinowski, Razvan Pascanu, Peter Battaglia, and Timothy Lillicrap. A simple neural network module for relational reasoning. In I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, editors, *NIPS*, pages 4967–4976. Curran Associates, Inc., 2017. URL <http://papers.nips.cc/paper/7082-a-simple-neural-network-module-for-relational-reasoning.pdf>.
- [53] Victor Garcia Satorras and Joan Bruna Estrach. Few-shot learning with graph neural networks. In *ICLR*, 2018. URL <https://openreview.net/forum?id=BJj6qGbRW>.
- [54] F. Scarselli, S. L. Yong, M. Gori, M. abd Hagenbuchner, A. C. Tsoi, and M. Maggini. Graph neural networks for ranking web pages. *IEEE/WIC/ACM International Conference on Web Intelligence*, pages 666–672, 2005.
- [55] Scott A Sisson, Yanan Fan, and Mark Beaumont. *Handbook of approximate Bayesian computation*. Chapman and Hall/CRC, 2018.
- [56] Sainbayar Sukhbaatar, Rob Fergus, et al. Learning multiagent communication with backpropagation. In *NIPS*, pages 2244–2252, 2016.

- [57] J.W. Tukey. Modern techniques in data analysis. In *Proceedings of the NSF-Sponsored Regional Research Conference, volume 7, Dartmouth, Massachusetts. Southern Massachusetts University.*, 1977.
- [58] Sjoerd van Steenkiste, Michael Chang, Klaus Greff, and Jürgen Schmidhuber. Relational neural expectation maximization: Unsupervised discovery of objects and their interactions. *arXiv:1802.10353*, 2018.
- [59] C. Villani. *Optimal Transport: Old and New*. Grundlehren der mathematischen Wissenschaften. Springer Berlin Heidelberg, 2008. ISBN 9783540710509. URL https://books.google.com/books?id=hV8o5R7_5tkC.
- [60] Yue Wang, Yongbin Sun, Ziwei Liu, Sanjay E Sarma, Michael M Bronstein, and Justin M Solomon. Dynamic graph cnn for learning on point clouds. *arXiv:1801.07829*, 2018.
- [61] Nicholas Watters, Daniel Zoran, Theophane Weber, Peter Battaglia, Razvan Pascanu, and Andrea Tacchetti. Visual interaction networks: Learning a physics simulator from video. In *NIPS*, pages 4539–4547, 2017.
- [62] Wikipedia contributors. Transportation theory (mathematics) — Wikipedia, the free encyclopedia, 2019. URL [https://en.wikipedia.org/wiki/Transportation_theory_\(mathematics\)](https://en.wikipedia.org/wiki/Transportation_theory_(mathematics)). [Online; accessed 11-Nov-2019].
- [63] RH Wilcox et al. Adaptive control processes—a guided tour, by richard bellman, princeton university press, princeton, new jersey, 1961, 255 pp. *Naval Research Logistics Quarterly*, 8(3):315–316, 1961.
- [64] Yingxue Zhang, Soumyasundar Pal, Mark Coates, and Deniz Üstebay. Bayesian graph convolutional neural networks for semi-supervised classification. *arXiv:1811.11103*, 2018.
- [65] Jie Zhou, Ganqu Cui, Zhengyan Zhang, Cheng Yang, Zhiyuan Liu, and Maosong Sun. Graph neural networks: A review of methods and applications. <https://arxiv.org/abs/1812.08434>, 2019.