

---

# Generalised $f$ -Mean Aggregation for Graph Neural Networks

---

Ryan Kortvelesy, Steven Morad and Amanda Prorok  
University of Cambridge  
{rk627, sm2558, asp45}@cam.ac.uk

## Abstract

Graph Neural Network (GNN) architectures are defined by their implementations of update and aggregation modules. While many works focus on new ways to parametrise the update modules, the aggregation modules receive comparatively little attention. Because it is difficult to parametrise aggregation functions, currently most methods select a “standard aggregator” such as mean, sum, or max. While this selection is often made without any reasoning, it has been shown that the choice in aggregator has a significant impact on performance, and the best choice in aggregator is problem-dependent. Since aggregation is a lossy operation, it is crucial to select the most appropriate aggregator in order to minimise information loss. In this paper, we present GenAgg, a generalised aggregation operator, which parametrises a function space that includes all standard aggregators. In our experiments, we show that GenAgg is able to represent the standard aggregators with much higher accuracy than baseline methods. We also show that using GenAgg as a drop-in replacement for an existing aggregator in a GNN often leads to a significant boost in performance across various tasks.

## 1 Introduction

Graph Neural Networks (GNNs) provide a powerful framework for operating over structured data. Taking advantage of relational inductive biases, they use local filters to learn functions that generalise over high-dimensional data. Given different graph structures, GNNs can represent many special cases, including CNNs (on grid graphs) [12], RNNs (on line graphs) [4], and Transformers (on fully connected graphs) [19]. All of these architectures can be subsumed under the Graph Networks framework, parametrised by update and aggregation modules [2]. Although the framework itself is general, the representational capacity is often constrained in practice through design choices, which create a human prior [21]. There are two primary reasons for introducing this human prior. First, there are no standard methods to parametrise all of the modules—MLPs can be used as universal approximators in the update modules, but it is nontrivial to parametrise the function space of aggregators. Consequently, most GNNs simply make a design choice for the aggregation functions, selecting mean, sum, or max [21].

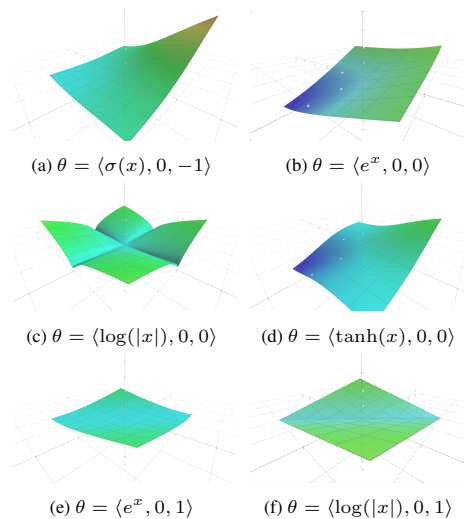


Figure 1: A qualitative demonstration of the diversity of functions that can be represented by GenAgg. In these visualisations, GenAgg is plotted as a function of inputs  $x_0$  and  $x_1$  for different parametrisations  $\theta = \langle f, \alpha, \beta \rangle$  (see Equation (1)).

Second, constraints can boost performance in GNNs, either through invariances or a regularisation effect.

In this paper, we focus on the problem of parametrising the space of aggregation functions. The ultimate goal is to create an aggregation function which can represent the set of all desired aggregators while remaining as constrained as possible. In prior work, one approach is to introduce learnable parameters into functions that could parametrise min, mean, and max, such as the Powermean and a variant of Softmax [8, 13, 20]. However, these approaches can only parametrise a small set of aggregators, and they can introduce instability in the training process (see Section 4). On the opposite end of the spectrum, methods like Deep Sets [22] and LSTMAgg [9] are capable of universal approximation over set functions, but they are extremely complex, which leads to poor sample efficiency. These methods scale in complexity (*i.e.* number of parameters) with the dimensionality of the input, and lack some of the useful constraints that are shared among standard aggregators (see Section 4). Consequently, the complexity of these methods counteracts the benefits of simple GNN architectures.

Although existing approaches present some limitations, the theoretical advantages of a learnable aggregation module are evident. It has been shown that the choice of aggregation function not only has a significant impact on performance, but also is problem-specific [21]. Since there is no aggregator that can discriminate between all inputs [5], it is important to select an aggregator that preserves the relevant information. In this paper, we present a method that parametrises the function space, allowing GNNs to learn the most appropriate aggregator for each application.

## Contributions

- We introduce Generalised Aggregation (GenAgg), the first aggregation method based on the *generalised f-mean*. GenAgg is a learnable permutation-invariant aggregator which is provably capable (both theoretically and experimentally) of representing all “standard aggregators” (see Appendix C for proofs). The representations learned by GenAgg are *explainable*—each learnable parameter has an interpretable meaning (Section 6).
- Our experiments provide several insights about the role of aggregation functions in the performance of GNNs. In our regression experiments, we demonstrate that GNNs struggle to “make up for” the lack of representational complexity in their constituent aggregators, even when using state-of-the-art parametrised aggregators. This finding is validated in our GNN benchmark experiments, where we show that a GenAgg-based GNN outperforms all of the baselines, including standard aggregators and other state-of-the-art parametrised aggregators.
- Finally, we show that GenAgg satisfies a generalisation of the distributive property. We derive the solution for a binary operator that satisfies this property for given parametrisations of GenAgg. The generalised distributive property can be leveraged in algorithms using GenAgg to improve space and memory-efficiency.

## 2 Problem Statement

Consider a multiset  $\mathcal{X} = \{x_1, x_2, \dots, x_n\}$  of cardinality  $|\mathcal{X}| = n$ , where  $x_i \in \mathbb{R}^d$ . We define an aggregation function as a symmetric function  $\odot : \mathbb{R}^{n \times d} \mapsto \mathbb{R}^{1 \times d}$ . The aggregator must be independent over the feature dimension, so without loss of generality it can be represented over a single dimension  $\odot : \mathbb{R}^n \mapsto \mathbb{R}^1$ . A set of standard aggregators is defined  $\mathcal{A} = \{\text{mean, sum, product, min, max, \dots}\}$  (for the full list see Table 1). Our task is to create an aggregator  $\oplus_\theta : \mathbb{R}^n \mapsto \mathbb{R}^1$  parametrised by  $\theta$  which can represent all standard aggregators:  $\forall \odot_i \in \mathcal{A} \exists \theta : \oplus_\theta = \odot_i$ .

## 3 Method

In this section we introduce GenAgg, a parametrised aggregation function which is based on the *generalised f-mean* [11]. In our formulation, we introduce additional parameters to increase the representational capacity of the *f-mean*, producing the *augmented f-mean* (AFM). Then, as the implementation is non-trivial, we propose a method to implement it. The novel aspect of GenAgg is not

Aggregation Function	$\alpha$	$\beta$	$f$	GenAgg	SoftmaxAgg	PowerAgg
mean: $\frac{1}{n} \sum x_i$	0	0	$f(x) = x$	✓	✓	✓
sum: $\sum x_i$	1	0	$f(x) = x$	✓	✗	✗
product: $\prod  x_i $	1	0	$f(x) = \log( x )$	✓	✗	✗
min (magnitude): $\min  x_i $	0	0	$f(x) = \lim_{p \rightarrow \infty}  x ^{-p}$	✓	✗	✓
max (magnitude): $\max  x_i $	0	0	$f(x) = \lim_{p \rightarrow \infty}  x ^p$	✓	✗	✓
min: $\min x_i$	0	0	$f(x) = \lim_{p \rightarrow \infty} e^{-px}$	✓	✓	✗
max: $\max x_i$	0	0	$f(x) = \lim_{p \rightarrow \infty} e^{px}$	✓	✓	✗
harmonic mean: $\frac{n}{\sum \frac{1}{x_i}}$	0	0	$f(x) = \frac{1}{x}$	✓	✗	✓
geometric mean: $\sqrt[n]{\prod  x_i }$	0	0	$f(x) = \log( x )$	✓	✗	✓
root mean square: $\sqrt{\frac{1}{n} \sum x_i^2}$	0	0	$f(x) = x^2$	✓	✗	✓
euclidean norm: $\sqrt{\sum x_i^2}$	1	0	$f(x) = x^2$	✓	✗	✗
standard deviation: $\sqrt{\frac{1}{n} \sum (x_i - \mu)^2}$	0	1	$f(x) = x^2$	✓	✗	✗
log-sum-exp: $\log(\sum e^{x_i})$	1	0	$f(x) = e^x$	✓	✗	✗

Table 1: A table of all of the most common aggregators. For each special case, we specify the values of  $\alpha$ ,  $\beta$ , and  $f$  for which the augmented  $f$ -mean is equivalent (see Appendix C). We also report whether or not SoftmaxAgg and PowermeanAgg can represent each aggregator.

only the augmented  $f$ -mean formula, but also the implementation, which allows the mathematical concept of a generalised mean to be utilised in a machine learning context.

### 3.1 Generalised $f$ -Mean

The *generalised  $f$ -mean* [11] is given by:  $f^{-1}(\frac{1}{n} \sum_i f(x_i))$ . While it is difficult to define aggregation functions, the generalised  $f$ -mean provides a powerful intuition: most aggregators can be represented by a single invertible scalar-valued function  $f : \mathbb{R}^1 \mapsto \mathbb{R}^1$ . This is a useful insight, because it allows comparisons to be drawn between aggregators by analysing their underlying functions  $f$ . Furthermore, it provides a framework for discovering new aggregation functions. While classic functions like  $e^x$ ,  $\log(x)$ ,  $x^p$  all map to aggregators in  $\mathcal{A}$ , new aggregators can be created by defining new functions  $f$ .

### 3.2 Augmented $f$ -Mean

The standard generalised  $f$ -mean imposes strict constraints, such as symmetry (permutation-invariance), idempotency ( $\bigoplus(\{x, \dots, x\}) = x$ ), and monotonicity ( $\forall i \in [1..n], \frac{\partial \bigoplus(\{x_1, \dots, x_n\})}{\partial x_i} \geq 0$ ). However, this definition is too restrictive to parametrise many special cases of aggregation functions. For example, sum violates idempotency ( $\sum_{i \in [1..n]} x_i = nx$ ), and standard deviation violates monotonicity ( $\frac{\partial \sigma(\{x_1, 1\})}{\partial x_1} |_{x_1=0} < 0$ ). Consequently, we deem these constraints to be counterproductive. In our method, we introduce learnable parameters  $\alpha$  and  $\beta$  to impose a relaxation on the idempotency and monotonicity constraints, while maintaining symmetry. We call this relaxed

formulation the *augmented f-mean* (AFM), given by:

$$\bigoplus_{i \in [1..n]} x_i = f^{-1} \left( n^{\alpha-1} \sum_{i \in [1..n]} f(x_i - \beta\mu) \right). \quad (1)$$

The  $\alpha$  parameter allows AFM to control its level of dependence on the cardinality of the input  $\mathcal{X}$ . For example, given  $f(x) = \log(|x|)$ , if  $\alpha = 0$ , then AFM represents the geometric mean:  $\bigoplus_{(f, \alpha, \beta)} = \bigoplus_{(\log(|x|), 0, 0)} = \sqrt[n]{\prod |x_i|}$ . However, if  $\alpha = 1$ , then the  $n$ -th root disappears, and AFM represents a product:  $\bigoplus_{(\log(|x|), 1, 0)} = \prod |x_i|$ .

The  $\beta$  parameter enables AFM to calculate *centralised moments*, which are quantitative measures of the distribution of the input  $\mathcal{X}$  [17]. The first raw moment of  $\mathcal{X}$  is the mean  $\mu = \frac{1}{n} \sum x_i$ , and the  $k$ -th central moment is given by  $\mu_k = \sum (x_i - \mu)^k$ . With the addition of  $\beta$ , it becomes possible for AFM to represent  $\sqrt[k]{\mu_k}$ , the  $k$ -th root of the  $k$ -th central moment. For  $k = 2$ , this quantity is the standard deviation, which is in our set of standard aggregators  $\mathcal{A}$ . If the output is scaled to the  $k$ -th power, then it can also represent metrics such as variance, unnormalised skewness, and unnormalised kurtosis. It is clear that these metrics about the distribution of data are useful—they can have real-world meaning (*e.g.*, moments of inertia), and they have been used as aggregators in GNNs in prior work [5]. Consequently,  $\beta$  provides AFM with an important extra dimension of representational complexity. In addition to representing the centralised moments when  $\beta = 1$  and  $f(x) = x^p$ ,  $\beta$  allows *any* aggregator to be calculated in a centralised fashion. While the centralised moments are the only well-known aggregators that arise from nonzero  $\beta$ , there are several aggregators with qualitatively unique behaviour that can only be represented with nonzero  $\beta$  (see Fig. 1).

With this parametrisation, AFM can represent any standard aggregator in  $\mathcal{A}$  (Table 1). Furthermore, by selecting new parametrisations  $\theta = \langle f, \alpha, \beta \rangle$ , it is possible to compose new aggregators (Fig. 1).

### 3.3 Implementation

In Equation (1), the manner in which  $f^{-1}$  is implemented is an important design choice. One option is to learn the coefficients for an analytical function (*e.g.* a truncated Taylor Series) and analytically invert it. However, it can be difficult to compute the analytical inverse of a function, and without carefully selected constraints, there is no guarantee that  $f$  will be invertible.

Another possible option is an invertible neural network (*e.g.* a parametrised invertible mapping from *normalising flows* [10]). We have tested the invertible networks from normalising flows literature as implementations for  $f$ . While they work well on smaller tasks, these methods present speed and memory issues in larger datasets.

In practice, we find that the most effective approach is to use two separate MLPs for  $f$  and  $f^{-1}$ . We enforce the constraint  $x = f^{-1}(f(x))$  by minimizing the following optimisation objective:

$$\mathcal{L}_{\text{inv}}(\theta_1, \theta_2) = \mathbb{E} \left[ \left( \left| f_{\theta_2}^{-1}(f_{\theta_1}(x)) \right| - |x| \right)^2 \right]. \quad (2)$$

The absolute value operations apply a relaxation to the constraint, allowing  $f^{-1}(f(x))$  to reconstruct either  $x$  or  $|x|$ . This is useful because several of the ground truth functions from Table 1 include an absolute value, making them non-invertible. With this relaxation, it becomes possible to represent those cases. This optimisation objective ensures that  $f$  is both monotonic and invertible over the domains  $\mathbb{R}^+$  and  $\mathbb{R}^-$ , independently. In our implementation, this extra optimisation objective is hidden behind the GenAgg interface and gets applied automatically with a forward hook, so it is not necessary for the user to apply an extra loss term.

While using a scalar-valued  $f : \mathbb{R}^1 \mapsto \mathbb{R}^1$  is the most human-interpretable formulation, it is not necessary. A valid implementation of GenAgg can also be achieved with  $f : \mathbb{R}^1 \mapsto \mathbb{R}^d$  and  $f^{-1} : \mathbb{R}^d \mapsto \mathbb{R}^1$ . In our experiments, we found that mapping to a higher intermediate dimension can sometimes improve performance over a scalar-valued  $f$  (see training details in Appendix E).

### 3.4 Generalised Distributive Property

Given that GenAgg presents a method of parameterising the function space of aggregators, it can also be used as a tool for mathematical analysis. To demonstrate this, we use the aug-



Aggregation Function	Distributive Operations $\psi(a, b)$
mean: $\frac{1}{n} \sum x_i$	$a + b, a \cdot b$
sum: $\sum x_i$	$a \cdot b$
product: $\prod  x_i $	$ a ^{\log  b }$
min (magnitude): $\min  x_i $	$\min( a ,  b )$
max (magnitude): $\max  x_i $	$\max( a ,  b )$
min: $\min x_i$	$\min(a, b)$
max: $\max x_i$	$\max(a, b)$
harmonic mean: $\frac{n}{\sum \frac{1}{x_i}}$	$\frac{a \cdot b}{a + b}, a \cdot b$
geometric mean: $\sqrt[n]{\prod  x_i }$	$ a  \cdot  b ,  a ^{\log  b }$
root mean square: $\sqrt{\frac{1}{n} \sum x_i^2}$	$\sqrt{a^2 + b^2},  a  \cdot  b $
euclidean norm: $\sqrt{\sum x_i^2}$	$ a  \cdot  b $
standard deviation: $\sqrt{\frac{1}{n} \sum (x_i - \mu)^2}$	$ a  \cdot  b $
log-sum-exp: $\log \left( \sum e^{x_i} \right)$	$a + b$

Table 2: A table of the distributive operations  $\psi$  that satisfy each aggregation function, computed using Equation 3. All aggregation functions have at least one solution, and some special cases have multiple solutions.

mented  $f$ -mean to analyse a generalised form of the distributive property, which is satisfied if  $\psi(c, \odot_{x_i \in \mathcal{X}} x_i) = \odot_{x_i \in \mathcal{X}} \psi(c, x_i)$  for binary operator  $\psi$  and aggregator  $\odot$ . For a given aggregation function parametrised by  $f$  (assuming  $\beta$  is 0), we derive a closed-form solution for a corresponding binary operator which will satisfy the generalised distributive property (for further explanation and proofs, see Appendix A).

**Theorem 3.1.** *For GenAgg parametrised by  $\theta = \langle f, \alpha, \beta \rangle = \langle f, \alpha, 0 \rangle$ , the binary operator  $\psi$  which will satisfy the Generalised Distributive Property for  $\oplus_\theta$  is given by:*

$$\psi(a, b) = f^{-1}(f(a) \cdot f(b)) \quad (3)$$

Furthermore, for the special case  $\theta = \langle f, \alpha, \beta \rangle = \langle f, 0, 0 \rangle$ , there  $\psi(a, b) = f^{-1}(f(a) + f(b))$  is an additional solution.

For example, for the euclidean norm where  $f(x) = x^2$  and  $\alpha = 1$ , the binary operator is  $\psi(a, b) = (a^2 \cdot b^2)^{\frac{1}{2}} = a \cdot b$ , which implies that a constant multiplicative term can be moved outside of the euclidean norm. This is a useful finding, as the distributive property can be used to improve algorithmic time and space complexity (e.g. the FFT) [1]. With our derivation of  $\psi$  as a function of  $f$ , it is possible to implement similar efficient algorithms with GenAgg.

## 4 Related Work

Several existing works propose methods to parametrise the space of aggregation functions. These methods can broadly be divided into two categories. *Mathematical* approaches derive an explicit equation in terms of the inputs and one or more learnable parameters. Usually, these approaches represent a smooth interpolation through function space from min, through mean, to max. Alternatively, *Deep Learning* approaches seek to use the universal approximation properties of neural networks to maximise representational complexity.

### 4.1 Mathematical Approaches

**SoftmaxAgg** SoftmaxAgg computes the weighted sum of the set, where the weighting is derived from the softmax over the elements with some learnable temperature term  $s$  [13, 20]. This formu-

lation allows SoftmaxAgg to represent mean, min, and max (see Table 1). Unfortunately, it fails to generalise across the majority of the standard aggregators.

**PowerAgg** Based on the  $p$ -norm, PowerAgg is a special case of GenAgg where  $\alpha = 0$ ,  $\beta = 0$ , and  $f(x) = x^p$ . There are some methods which use the powermean directly [13, 20, 8], and others which build on top of it [18]. Theoretically, PowerAgg can represent a significant subset of the standard aggregators: min magnitude, max magnitude, mean, root mean square, harmonic mean, and geometric mean (although the geometric mean requires  $\lim_{p \rightarrow 0}$ , so it is not practically realisable) (see Table 1). Unfortunately, there is a caveat to this approach: for negative inputs  $x_i < 0$  and non-integer values  $p$ , it is only defined in the complex domain. Furthermore, for negative inputs, the gradient  $\frac{\partial x^p}{\partial p}$  with respect to trainable parameter  $p$  is complex and oscillatory (and therefore is prone to getting stuck in local optima). In order to fix this problem, the inputs must be constrained to be positive. In prior work, this has been achieved by clamping  $x'_i = \max(x_i, 0)$  [13], subtracting the minimum element  $x'_i = x_i - \min(\mathcal{X})$  [20], or taking the absolute value  $x'_i = |x_i|$  [8]. However, this removes important information, making it impossible to reconstruct most standard aggregators.

## 4.2 Deep Learning Approaches

**PNA** While Principle Neighbourhood Aggregation [5] is introduced as a GNN architecture, its novelty stems from its method of aggregation. In PNA, input signal is processed by a set of aggregation functions, which is produced by the cartesian product of standard aggregators {mean, min, max, std} and scaling factors  $\{\frac{1}{n}, 1, n\}$ . The output of every aggregator is concatenated, and passed through a dense network. While this increases the representational complexity of the aggregator, it also scales the dimensionality of the input by the number of aggregators multiplied by the number of scaling factors, which can decrease sample efficiency (Figure 3). Furthermore, the representational complexity of the method is limited by the choice of standard aggregators—it cannot be used to represent many of the special cases of parametrised general aggregators.

**LSTMAgg** In LSTMAgg, the input set is treated as a sequence (applying some random permutation), and is encoded with a recurrent neural network [9]. While this method is theoretically capable of universal approximation, in practice its non-permutation-invariance can cause its performance to suffer (as the factorial complexity of possible orderings leads to sample-inefficiency). SortAgg addresses this issue by sorting the inputs with computed features, and passing the first  $k$  sorted inputs through convolutional and dense networks [23]. While this method solves the issue of non-permutation-invariance, it loses the capability of universal approximation by truncating to  $k$  inputs. While universal approximation is not a requirement for an effective aggregation function, we note that it cannot represent many of the standard aggregators.

**Deep Sets** Deep Sets is a universal set function approximator [22]. However, because it operates over the feature dimension in addition to the “set” dimension, it is not regarded as an aggregation function. Instead, it usually serves as a full GNN layer or graph pooling architecture [14, 15]. One may note that the formulation for Deep Sets  $\phi(\sum_{i \in [1..n]} f(x_i))$  bears some resemblance to our method. However, there are two important differences. First, our method adds the constraint  $\phi = f^{-1}$ , limiting possible parametrisations to the subspace where all of the standard aggregators lie. Second, while the learnable functions  $\phi$  and  $f$  in Deep Sets are fully connected over the feature dimension, the  $f$  and  $f^{-1}$  modules in our architecture are scalar-valued functions which are applied element-wise. To summarise, Deep Sets is useful as a set function approximator, but it lacks constraints that would make it viable as an aggregation function.

## 5 Experiments

In this paper, we run three experiments. First, we show that GenAgg can perform regression to recover any standard aggregation function. Then, we evaluate GenAgg and several baselines inside of a GNN. The resulting GNN architectures are given the same task of regressing upon graph-structured data generated with a standard aggregator. This tests if it is possible for a GNN with a given aggregator to represent data which was generated by different underlying aggregators. Finally, we provide practical results by running experiments on public GNN benchmark datasets: CLUSTER, PATTERN, CIFAR10, and MNIST [6].

Aggregation	GenAgg	P-Agg	S-Agg	mean	Aggregation	GenAgg	P-Agg	S-Agg	mean
mean	<b>1.000</b>	0.817	<b>1.000</b>	<b>1.000</b>	mean	0.977	<b>0.999</b>	<b>1.000</b>	0.972
sum	<b>1.000</b>	0.761	0.887	0.888	sum	<b>0.971</b>	0.906	0.887	0.903
product (mag)	<b>0.985</b>	0.407	0.172	0.022	product (mag)	<b>0.966</b>	0.644	0.726	0.434
min (mag)	<b>0.962</b>	0.450	0.024	0.027	min (mag)	<b>0.952</b>	0.876	0.810	0.731
max (mag)	<b>0.990</b>	0.586	0.423	0.024	max (mag)	<b>0.986</b>	0.734	0.784	0.747
min	<b>0.995</b>	0.734	<b>1.000</b>	0.805	min	<b>0.995</b>	0.986	<b>0.999</b>	0.806
max	<b>0.990</b>	0.920	<b>1.000</b>	0.805	max	<b>0.989</b>	0.976	<b>0.999</b>	0.845
harm. mean (abs)	<b>0.986</b>	0.453	0.088	0.027	harm. mean (abs)	<b>0.931</b>	0.797	0.842	0.697
geom. mean (abs)	<b>0.994</b>	0.481	0.152	0.031	geom. mean (abs)	<b>0.963</b>	0.629	0.836	0.626
root mean square	<b>0.996</b>	0.532	0.308	0.028	root mean square	<b>0.975</b>	0.775	0.808	0.899
euclidean norm	<b>0.964</b>	0.585	0.464	0.019	euclidean norm	<b>0.985</b>	0.742	0.680	0.756
standard dev.	<b>0.999</b>	0.442	0.558	0.013	standard dev.	<b>0.966</b>	0.739	0.805	0.624
log-sum-exp	<b>0.999</b>	0.823	0.947	0.747	log-sum-exp	<b>0.983</b>	0.919	0.952	0.841

(a) Aggregator Regression.

(b) GNN Regression

Figure 2: Results for the Aggregator Regression and GNN Regression experiments, indicating the ability of GenAgg, PowerAgg (P-Agg), SoftmaxAgg (S-Agg), and mean to parametrise each standard aggregator in  $\mathcal{A}$ . We report the correlation coefficient between the ground truth and predicted outputs. The highest-performing methods (and those within 0.01 correlation) are shown in bold.

For all baselines, we use the implementations provided in PyTorch Geometric [7]. The only exception is PNA, which is a GNN architecture by nature, not an aggregation method. For our experiments, we adapt PNA into an aggregation method, staying as true to the original formulation as possible:  $\text{PNA}(\mathcal{X}) = f([1, n, \frac{1}{n}] \otimes [\text{mean}(\mathcal{X}), \text{std}(\mathcal{X}), \text{min}(\mathcal{X}), \text{max}(\mathcal{X})])$ , where  $f$  is a linear layer mapping from  $\mathbb{R}^{12d}$  back to  $\mathbb{R}^d$ .

For more training details, see Appendix E. Our code can be found at: <https://github.com/Acciorocketships/generalised-aggregation>.

## 5.1 Aggregator Regression

In this experiment, we generate a random graph  $\mathcal{G} = \langle \mathcal{V}, \mathcal{E} \rangle$  with  $|\mathcal{V}| = 8$  nodes and an edge density of  $\frac{|\mathcal{E}|}{|\mathcal{V}|^2} = 0.3$ . For each node  $i \in \mathcal{V}$ , we draw an internal state  $x_i \in \mathbb{R}^d$  from a normal distribution  $x_i \sim \mathcal{N}(\mathbf{0}_d, I_d)$  with  $d = 6$ . Then, we generate training data with a set of ground truth aggregators  $\odot_k \in \mathcal{A}$  (where  $k$  is an index). For each aggregator  $\odot_k$ , the dataset  $X_k, Y_k$  is produced with a Graph Network [2], using  $\odot_k$  as the node aggregation module. The inputs are defined by the set of neighbourhoods in the graph  $X_k = \{\mathcal{X}_i \mid i \in [1..|\mathcal{V}|]\}$  where the neighbourhood  $\mathcal{X}_i$  is defined as  $\mathcal{X}_i = \{x_j \mid j \in \mathcal{N}_i\}$  with  $\mathcal{N}_i = \{j \mid (i, j) \in \mathcal{E}\}$ . The corresponding ground truth outputs are defined as  $Y_k = \{y_i \mid i \in [1..|\mathcal{V}|]\}$ , where  $y_i = \odot_k(\mathcal{X}_i)$ .

The model that we use for regression takes the same form as the model used to generate the data, except that the standard aggregator used to generate the training data  $\odot_k$  is replaced with a parametrised aggregator  $\oplus_\theta$ :

$$\hat{y}_i = \bigoplus_{\theta}^{x_j \in \mathcal{X}_i} x_j \quad (4)$$

In our experiments, each type of parametrised aggregator (GenAgg, SoftmaxAgg, PowerAgg, and mean as a baseline) is trained separately on each dataset  $X_k, Y_k$ .

**Results.** We report the MSE loss and correlation coefficient with respect to the ground truth in Table 2a. GenAgg is able to represent all of the standard aggregators with a correlation of at least 0.96, and most aggregators with a correlation of greater than 0.99. The only cases where the performance of GenAgg is surpassed by a baseline are min and max, where SoftmaxAgg exhibits marginally higher accuracy.

One interesting observation is that even if the baselines can represent an aggregator in *theory*, they cannot necessarily do so in practice. For example, PowerAgg can theoretically represent the geomet-

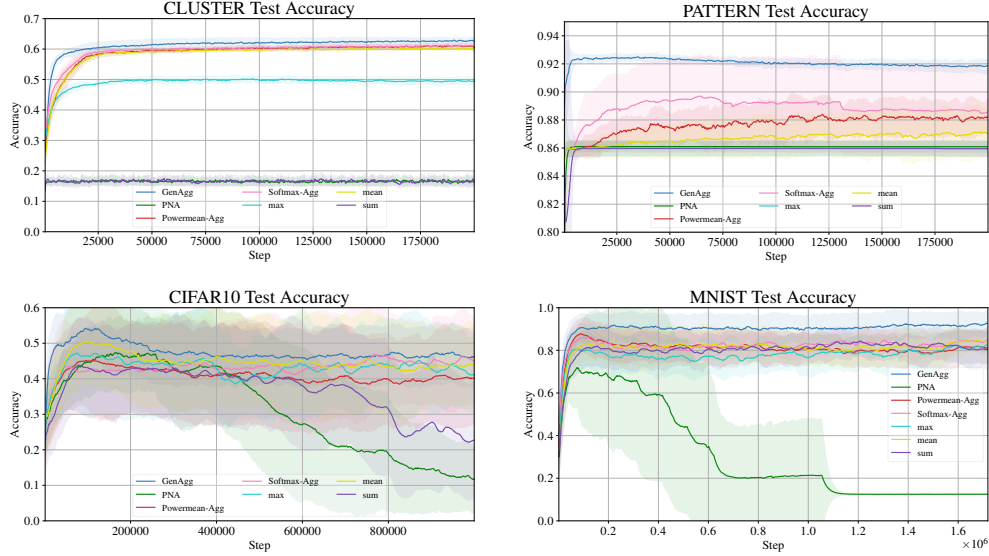


Figure 3: Test accuracy for GNNs with various aggregators on GNN benchmark datasets. In this experiment, each trial uses the same base GNN architecture (4-layer GraphConv), and the default aggregator is replaced with either GenAgg, PowermeanAgg (P-Agg), SoftmaxAgg (S-Agg), PNA, mean, sum, or max. The plots depict the mean and standard deviation of the test accuracy over 10 trials (note that the y-axis is scaled to increase readability). The table reports the maximum of the mean test accuracy over all timesteps, as well as the standard deviation (rounded to 0.01).

ric mean with  $\lim_{p \rightarrow 0} (\frac{1}{n} \sum_i x_i^p)^{\frac{1}{p}}$ , but in practice there are instabilities as  $p$  approaches 0 because  $\frac{1}{p}$  approaches  $\frac{1}{0}$ . Similarly, while in theory PowerAgg can represent min magnitude, max magnitude, harmonic mean, and root mean square, it falls short in practice (see Table 2a), likely because of the reasons stated in Section 4. In other cases, the baselines can perform well even if they should not be able to represent the target in theory. One such example is PowerAgg, which achieves a correlation of 0.92 on max, but only 0.59 on max magnitude, which is the opposite of what theory might suggest. This is likely due to the the clamp operation that Pytorch Geometric’s implementation uses to restricts inputs to the positive domain. The performance of max magnitude suffers, as it misses cases where the highest magnitude element is negative. Similarly, the performance of max increases, because it simply selects the maximum among the positive elements. Another baseline which performs unexpectedly well is SoftmaxAgg, which achieves a high correlation with the log-sum-exp aggregator. While it cannot compute a log, the SoftmaxAgg formulation does include a sum of exponentials, so it is able to produce a close approximation.

## 5.2 GNN Regression

In GNN Regression, the experimental setup is the same as that of Aggregator Regression (Section 5.1), with the exception that the observation size is reduced to  $d = 1$ . However, instead of using GenAgg  $\oplus_{\theta}$  as a model, we use a multi-layer GNN. The GNN is implemented with 4 layers of GraphConv [16] with Mish activation (after every layer except the last), where the default aggrega-

tion function is replaced by a parametrised aggregator  $\bigoplus_{\theta}$ :

$$z_i^{(k+1)} = \text{Mish} \left( W_1^{(k)} z_i^{(k)} + W_2^{(k)} \bigoplus_{\theta} z_i^{(k)} \right), \text{ s.t. } z_i^{(0)} = x_i \quad (5)$$

$$\hat{y}_i = W_1^{(3)} z_i^{(3)} + W_2^{(3)} \bigoplus_{\theta} z_i^{(3)} \quad (6)$$

While this experiment uses the same dataset as Aggregator Regression (Section 5.1), it provides several new insights. First, while the aggregator regression experiment shows that GenAgg *can* represent various aggregators, this experiment demonstrates that training remains stable even when used within a larger architecture. Second, this experiment underlines the importance of using the correct aggregation function. While it is clear that it is advantageous to match a model’s aggregation function with that of the underlying mechanism which generated a particular dataset, we often opt to simply use a default aggregator. The conventional wisdom of this choice is that the other learnable parameters in a network layer can rectify an inaccurate choice in aggregator. However, the results from this experiment demonstrate that even with additional parameters, it is not necessarily possible to represent a different aggregator, underlining the importance of aggregators with sufficient representational capacity.

**Results.** The results show that GenAgg maintains its performance, even when used as a component within a GNN (Table 2b). GenAgg achieves a mean correlation of 0.97 across all aggregators. While the baselines perform significantly better with the help of a multi-layer GNN architecture, they still cannot represent many of the standard aggregators. The highest-performing baseline is SoftmaxAgg, which only achieves a mean correlation of 0.86.

### 5.3 GNN Benchmark

In this experiment, we examine the performance of GenAgg on GNN benchmark datasets [6]. In order to perform a comparison with benchmarks, we train on an existing GNN architecture (a 4-layer GraphConv [16] GNN with a hidden size of 64) where the default aggregator is replaced with a new aggregator, selected from {GenAgg, PowerAgg, SoftmaxAgg, PNA, mean, sum, max}.

**Results.** As shown in Fig 3, GenAgg outperforms all baselines in all four GNN benchmark datasets. It provides a significant boost in performance, particularly compared to the relatively small differences in performance between the baseline methods.

The training plots in Fig. 3 provide complementary information. One interesting observation is that GenAgg converges at least as fast as the other methods, and sometimes converges significantly faster (in PATTERN, for example). Furthermore, the training plots lend information about the stability of training. For example, note that in MNIST, most of the baseline methods achieve a maximum and then degrade in performance, while GenAgg maintains a stable performance throughout training.

## 6 Discussion

**Results.** In our experiments, we present two regression tasks and one GNN benchmark task. The regression experiments demonstrate that GenAgg is the only method capable of representing all of the standard aggregators, and a GNN cannot be used to compensate for the shortcomings of the baseline aggregators. The GNN benchmark experiment complements these findings, demonstrating that this representational complexity is actually useful in practice. The fact that GenAgg outperforms the standard aggregators (mean, max, and sum) on the GNN benchmark experiment implies that it is in fact creating a *new* aggregator. Furthermore, the fact that it outperforms baseline methods like SoftmaxAgg and PowermeanAgg implies that the aggregator learned by GenAgg lies outside the set of functions which can be represented by such methods.

**Limitations.** While GenAgg achieves positive results on these datasets, it is not possible to make generalisations about its performance in all applications. In particular, we observe that some datasets fundamentally require less complexity to solve, so simple aggregators are sufficient (*i.e.*, GenAgg fails to provide a significant performance boost). For a full list of datasets that we considered and further discussion of limitations, see Appendix D.

**Parameters.** When comparing the performance of different models, it is important to also consider the number of parameters. By introducing additional parameters, some models can improve overall performance at the cost of sample efficiency. While methods like PowerAgg and SoftmaxAgg only have one trainable parameter, GenAgg has two scalar parameters  $\alpha$  and  $\beta$ , and a learnable function  $f$ , which has 30 parameters in our implementation (independent of the size of the state). However, we observe that using GenAgg within a GNN is always at least as sample-efficient as the baselines, and sometimes converges significantly faster (Fig. 3 and Appendix B). Furthermore, while GenAgg has more parameters than PowerAgg and SoftmaxAgg, the increase is negligible compared to the total number of parameters in the GNN. We also note that GenAgg has significantly fewer parameters than the deep learning methods discussed in Section 4. While the deep learning methods scale linearly or quadratically in the dimension of the state, the number of parameters in GenAgg is constant.

**Stability.** Another observation from our experiments is that GenAgg exhibits more stability during the training process than the baselines (Appendix B). In the GNN Regression experiment, the PowerAgg and SoftmaxAgg training curves tend to plateau at least once before reaching their maximum value. It is possible that these methods lead to local optima because they are optimised in a lower dimensional parameter space [3]. For example, it is straightforward to smoothly transform a learned  $f$  in GenAgg from  $x^2$  to  $x^4$ , but to do so in PowerAgg, it is necessary to pass through  $x^3$ , which has significantly different behaviour in the negative domain. While PowerAgg restricts inputs to the positive domain to circumvent this particular issue, the problem of local optima can still arise when methods like PowerAgg or SoftmaxAgg are used as components in a larger architecture.

**Explainability.** While in this paper we primarily focus on the *performance* of GenAgg, we note that it also presents benefits in the realm of explainability. The three parameters in GenAgg are all human-readable (scalars and scalar-valued functions can easily be visualised), and they all provide a unique intuition. The  $\alpha$  parameter controls the dependence on the cardinality of the input set. The  $\beta$  parameter dictates if the aggregator is computed in a raw or centralised fashion (colloquially, it answers if the aggregator operates over the inputs themselves, or the variation between the inputs). Lastly, the function  $f$  can be analysed by considering the sign and magnitude of  $f(x_i)$ . The sign denotes if a given  $x_i$  increases ( $f(x) > 0$ ) or decreases ( $f(x_i) < 0$ ) the output. On the other hand, the magnitude  $|f(x_i)|$  can be interpreted as the relative impact of that point on the output. For example, the parametrisation of product is  $f(x) = \log(|x|)$ , which implies that a value of 1 has no impact on the output since  $|\log(|1|)| = 0$ , and extremely small values  $\epsilon$  have a large impact, because  $\lim_{\epsilon \rightarrow 0} |\log(|\epsilon|)| = \infty$ . Indeed, 1 is the identity element under multiplication, and multiplying by a small value  $\epsilon$  can change the output by many orders of magnitude. The interpretability of GenAgg can also be leveraged as a method to *select* an aggregator—a model can be pre-trained with GenAgg, and then each instance of GenAgg can be replaced with the most similar standard aggregator in  $\mathcal{A}$ .

## 7 Conclusion

In this paper we introduced GenAgg, a generalised, explainable aggregation function which parametrises the function space of aggregators, yet remains as constrained as possible to improve sample efficiency and prevent overfitting. In our experiments, we showed that GenAgg can represent all 13 of our selected “standard aggregators” with a correlation coefficient of at least 0.96. We also evaluated GenAgg alongside baseline methods within a GNN, illustrating how other approaches have difficulties representing standard aggregators, even with the help of additional learnable parameters. Finally, we demonstrated the usefulness of GenAgg on GNN benchmark tasks, comparing the performance of the same GNN with various different aggregators. The results showed that GenAgg provided a significant boost in performance over the baselines in all four datasets. Furthermore, GenAgg often exhibited more stability and faster convergence than the baselines in the training process. These results show that GenAgg is an application-agnostic aggregation method that can provide a boost in performance as a drop-in replacement for existing aggregators.

## 8 Acknowledgements

Ryan Kortvelesy and Amanda Prorok were supported in part by ARL DCIST CRA W911NF-17-2-0181 and European Research Council (ERC) Project 949940 (gAIA).

## References

- [1] S.M. Aji and R.J. McEliece. The generalized distributive law. *IEEE Transactions on Information Theory*, 46(2):325–343, March 2000.
- [2] Peter W Battaglia, Jessica B Hamrick, Victor Bapst, Alvaro Sanchez-Gonzalez, Vinicius Zambaldi, Mateusz Malinowski, Andrea Tacchetti, David Raposo, Adam Santoro, Ryan Faulkner, et al. Relational inductive biases, deep learning, and graph networks. *arXiv preprint arXiv:1806.01261*, 2018.
- [3] Alan J. Bray and David S. Dean. The statistics of critical points of Gaussian fields on large-dimensional spaces. *Physical Review Letters*, 98(15), April 2007.
- [4] Kyunghyun Cho, Bart Van Merriënboer, Caglar Gulcehre, Dzmitry Bahdanau, Fethi Bougares, Holger Schwenk, and Yoshua Bengio. Learning phrase representations using rnn encoder-decoder for statistical machine translation. *arXiv preprint arXiv:1406.1078*, 2014.
- [5] Gabriele Corso, Luca Cavalleri, Dominique Beaini, Pietro Liò, and Petar Veličković. Principal neighbourhood aggregation for graph nets. *Advances in Neural Information Processing Systems*, 33, 2020.
- [6] Vijay Prakash Dwivedi, Chaitanya K Joshi, Thomas Laurent, Yoshua Bengio, and Xavier Bresson. Benchmarking graph neural networks. *arXiv preprint arXiv:2003.00982*, 2020.
- [7] Matthias Fey and Jan E. Lenssen. Fast graph representation learning with PyTorch Geometric. In *ICLR Workshop on Representation Learning on Graphs and Manifolds*, 2019.
- [8] Caglar Gulcehre, Kyunghyun Cho, Razvan Pascanu, and Yoshua Bengio. Learned-norm pooling for deep feedforward and recurrent neural networks. In *Joint European Conference on Machine Learning and Knowledge Discovery in Databases*, pages 530–546. Springer, 2014.
- [9] Will Hamilton, Zhitaoying, and Jure Leskovec. Inductive representation learning on large graphs. *Advances in neural information processing systems*, 30, 2017.
- [10] Ivan Kobyzev, Simon JD Prince, and Marcus A Brubaker. Normalizing flows: An introduction and review of current methods. *IEEE transactions on pattern analysis and machine intelligence*, 43(11):3964–3979, 2020.
- [11] Andrey Kolmogorov. On the notion of mean. *Mathematics and Mechanics*, Selected works of A.N. Kolmogorov. Vol. 1:144–146, 1930.
- [12] Yann LeCun, Bernhard Boser, John S Denker, Donnie Henderson, Richard E Howard, Wayne Hubbard, and Lawrence D Jackel. Backpropagation applied to handwritten zip code recognition. *Neural computation*, 1(4):541–551, 1989.
- [13] Guohao Li, Chenxin Xiong, Ali Thabet, and Bernard Ghanem. Deepergcn: All you need to train deeper gcns. *arXiv preprint arXiv:2006.07739*, 2020.
- [14] Yujia Li, Chenjie Gu, Thomas Dullien, Oriol Vinyals, and Pushmeet Kohli. Graph matching networks for learning the similarity of graph structured objects. In *International conference on machine learning*, pages 3835–3845. PMLR, 2019.
- [15] Yujia Li, Chenjie Gu, Thomas Dullien, Oriol Vinyals, and Pushmeet Kohli. Graph matching networks for learning the similarity of graph structured objects. In *International conference on machine learning*, pages 3835–3845. PMLR, 2019.
- [16] Christopher Morris, Martin Ritzert, Matthias Fey, William L Hamilton, Jan Eric Lenssen, Gaurav Rattan, and Martin Grohe. Weisfeiler and leman go neural: Higher-order graph neural networks. In *Proceedings of the AAAI conference on artificial intelligence*, volume 33, pages 4602–4609, 2019.
- [17] Athanasios Papoulis and S Unnikrishna Pillai. *Probability, random variables and stochastic processes*. McGraw-Hill, 1991.
- [18] Giovanni Pellegrini, Alessandro Tibo, Paolo Frasconi, Andrea Passerini, and Manfred Jaeger. Learning Aggregation Functions. In *Proceedings of the Thirtieth International Joint Conference on Artificial Intelligence*, pages 2892–2898, Montreal, Canada, August 2021. International Joint Conferences on Artificial Intelligence Organization.
- [19] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017.
- [20] Beibei Wang and Bo Jiang. Generalizing aggregation functions in gnns: High-capacity gnns via nonlinear neighborhood aggregators. *arXiv preprint arXiv:2202.09145*, 2022.
- [21] Jiaxuan You, Zhitaoying, and Jure Leskovec. Design Space for Graph Neural Networks. In *Advances in Neural Information Processing Systems*, volume 33, pages 17009–17021. Curran Associates, Inc., 2020.
- [22] Manzil Zaheer, Satwik Kottur, Siamak Ravanbakhsh, Barnabas Poczos, Ruslan Salakhutdinov, and Alexander Smola. Deep Sets, April 2018. arXiv:1703.06114 [cs, stat].

- [23] Muhan Zhang, Zhicheng Cui, Marion Neumann, and Yixin Chen. An end-to-end deep learning architecture for graph classification. In *Proceedings of the AAAI conference on artificial intelligence*, volume 32, 2018.



## Appendix

### A Generalised Distributive Property

The Distributive Law is an extremely useful tool for analysis and computation in mathematics and computer science. Following from the additivity and homogeneity properties of linearity, it states that  $\sum_{i=1}^n c \cdot x_i = c \cdot \sum_{i=1}^n x_i$ . The Distributive Law is often leveraged to formulate fast algorithms, such as the Fast Fourier Transform and the Viterbi algorithm [1]. It can also be used to make algorithms more memory-efficient. For example, if we wish to take the DFT of a shifted signal, we can avoid storing the signal itself. Instead, the Distributive Law can be utilised to formulate the shifted DFT as a function of the non-shifted DFT:  $\text{DFT}[x(n - \Delta)] = e^{-i\omega_k \Delta} X(\omega_k)$ .

While the Distributive Law is defined for the group of real numbers under additivity  $(\mathbb{R}, +)$ , we can also define a more general distributive property for the Abelian group defined by an aggregator  $(\mathbb{R}, \odot)$ .

**Definition A.1** (Generalised Distributive Property). For a binary operator  $\psi$  and set aggregation function  $\odot$ , the Generalised Distributive Property is defined:

$$\psi \left( c, \odot_{x_i \in \mathcal{X}} x_i \right) = \odot_{x_i \in \mathcal{X}} \psi(c, x_i) \quad (7)$$

In this section, we derive the Generalised Distributive Property for the special case of  $\text{GenAgg } \odot = \oplus_\theta$ . That is, for a given parametrisation  $\theta$ , we derive an explicit formula for the corresponding function  $\psi$  which satisfies the Generalised Distributive Property. Note that while our solution holds for any function  $f$ , it focuses on the special cases  $\alpha = 0, \beta = 0$  and  $\alpha = 1, \beta = 0$ .

**Lemma A.1.** Given a binary operator of the form  $\psi(a, b) = \phi^{-1}(\phi(a) + \phi(b))$  and a parametrisation of  $\text{GenAgg } \theta = \langle f, \alpha, \beta \rangle = \langle f, \alpha, 0 \rangle$ , the operator  $\psi(a, b)$  satisfies the Generalised Distributive Property over the  $(\mathbb{R}, \oplus_\theta)$  abelian group if:

$$\rho^{-1} \left( \frac{n^\alpha}{n} \sum \rho(c + x_i) \right) = \rho^{-1} \left( \frac{n^\alpha}{n} \sum \rho(x_i) \right) + c \quad (8)$$

$$\text{where } \rho(x) = f(\phi^{-1}(x)) \quad (9)$$

*Proof.* Substituting  $\text{GenAgg}$  for the generic aggregation function  $\odot$  in the definition of the Generalised Distributive Property, we get:

$$f^{-1} \left( \frac{n^\alpha}{n} \sum f(\psi(c, x_i)) \right) = \psi \left( c, f^{-1} \left( \frac{n^\alpha}{n} \sum f(x_i) \right) \right) \quad (10)$$

We replace the binary operator  $\psi$  with its representation as a composition of univariate functions  $\psi(a, b) = \phi^{-1}(\phi(a) + \phi(b))$  to obtain:

$$f^{-1} \left( \frac{n^\alpha}{n} \sum f(\phi^{-1}(\phi(c) + \phi(x_i))) \right) = \phi^{-1} \left( \phi \left( f^{-1} \left( \frac{n^\alpha}{n} \sum f(x_i) \right) \right) + \phi(c) \right) \quad (11)$$

$$\phi \left( f^{-1} \left( \frac{n^\alpha}{n} \sum f(\phi^{-1}(\phi(c) + \phi(x_i))) \right) \right) = \phi \left( f^{-1} \left( \frac{n^\alpha}{n} \sum f(x_i) \right) \right) + \phi(c) \quad (12)$$

To simplify, we apply a change of variables  $x' = \phi(x)$ ,  $c' = \phi(c)$ :

$$\phi \left( f^{-1} \left( \frac{n^\alpha}{n} \sum f(\phi^{-1}(c' + x'_i)) \right) \right) = \phi \left( f^{-1} \left( \frac{n^\alpha}{n} \sum f(\phi^{-1}(x'_i)) \right) \right) + c' \quad (13)$$

Finally, we further simplify by substituting  $\rho(x) = f(\phi^{-1}(x))$ :

$$\rho^{-1} \left( \frac{n^\alpha}{n} \sum \rho(c' + x'_i) \right) = \rho^{-1} \left( \frac{n^\alpha}{n} \sum \rho(x'_i) \right) + c' \quad (14)$$

□

**Theorem A.2.** For GenAgg parametrised by  $\theta = \langle f, \alpha, \beta \rangle = \langle f, \alpha, 0 \rangle$ , the binary operator  $\psi$  which will satisfy the Generalised Distributive Property for  $\oplus_\theta$  is given by  $\psi(a, b) = f^{-1}(f(a) \cdot f(b))$ .

*Proof.* From Lemma A.1, the Generalised Distributive Property is satisfied if:

$$\rho^{-1} \left( \frac{n^\alpha}{n} \sum \rho(c + x_i) \right) = \rho^{-1} \left( \frac{n^\alpha}{n} \sum \rho(x_i) \right) + c \quad (15)$$

$$\text{where } \rho(x) = f(\phi^{-1}(x)) \quad (16)$$

If we select  $\rho(x) = e^x$ , then we can show that the condition is satisfied by the standard Distributive Law:

$$\log \left( \frac{n^\alpha}{n} \sum e^{c+x_i} \right) = \log \left( \frac{n^\alpha}{n} \sum e^{x_i} \right) + c \quad (17)$$

$$e^{\log \left( \frac{n^\alpha}{n} \sum e^{c+x_i} \right)} = e^{\log \left( \frac{n^\alpha}{n} \sum e^{x_i} \right) + c} \quad (18)$$

$$\frac{n^\alpha}{n} \sum e^c e^{x_i} = e^c \cdot \frac{n^\alpha}{n} \sum e^{x_i} \quad (19)$$

Given that  $\rho(x) = e^x$  satisfies the Distributive Property for this case, we can use it to solve for  $\phi$  and  $\phi^{-1}$ :

$$\rho(x) = f(\phi^{-1}(x)) \quad (20)$$

$$\phi^{-1}(x) = f^{-1}(\rho(x)) \quad (21)$$

$$\phi^{-1}(x) = f^{-1}(e^x) \quad (22)$$

$$\rho^{-1}(x) = \phi(f^{-1}(x)) \quad (23)$$

$$\phi(x) = \rho^{-1}(f(x)) \quad (24)$$

$$\phi(x) = \log(f(x)) \quad (25)$$

Finally, substituting  $\phi(x)$  and  $\phi^{-1}(x)$  back into the equation for  $\psi$ , we get:

$$\psi(a, b) = \phi^{-1}(\phi(a) + \phi(b)) \quad (26)$$

$$\psi(a, b) = f^{-1} \left( e^{\log(f(a)) + \log(f(b))} \right) \quad (27)$$

$$\psi(a, b) = f^{-1} \left( e^{\log(f(a))} \cdot e^{\log(f(b))} \right) \quad (28)$$

$$\psi(a, b) = f^{-1}(f(a) \cdot f(b)) \quad (29)$$

□

**Theorem A.3.** For the special case of GenAgg parametrised by  $\theta = \langle f, \alpha, \beta \rangle = \langle f, 0, 0 \rangle$ , the Generalised Distributive Property for  $\oplus_\theta$  is also satisfied by the binary operator  $\psi(a, b) = f^{-1}(f(a) + f(b))$ .

*Proof.* From Lemma A.1, the Generalised Distributive Property is satisfied if:

$$\rho^{-1} \left( \frac{n^\alpha}{n} \sum \rho(c + x_i) \right) = \rho^{-1} \left( \frac{n^\alpha}{n} \sum \rho(x_i) \right) + c \quad (30)$$

$$\text{where } \rho(x) = f(\phi^{-1}(x)) \quad (31)$$

In this proof, we are given that  $\alpha = 0$ :

$$\rho^{-1} \left( \frac{1}{n} \sum \rho(c + x_i) \right) = \rho^{-1} \left( \frac{1}{n} \sum \rho(x_i) \right) + c \quad (32)$$

If we select  $\rho(x) = x$ , then we can show that the condition is satisfied:

$$\frac{1}{n} \sum (c + x_i) = \left( \frac{1}{n} \sum x_i \right) + c \quad (33)$$

$$\frac{1}{n} \sum x_i + \frac{1}{n} \sum c = \left( \frac{1}{n} \sum x_i \right) + c \quad (34)$$

$$\left( \frac{1}{n} \sum x_i \right) + c = \left( \frac{1}{n} \sum x_i \right) + c \quad (35)$$

Given that  $\rho(x) = x$  satisfies the Distributive Property for this case, we can use it to solve for  $\phi$  and  $\phi^{-1}$ :

$$\rho(x) = f(\phi^{-1}(x)) \quad (36)$$

$$\phi^{-1}(x) = f^{-1}(\rho(x)) \quad (37)$$

$$\phi^{-1}(x) = f^{-1}(x) \quad (38)$$

$$\rho^{-1}(x) = \phi(f^{-1}(x)) \quad (39)$$

$$\phi(x) = \rho^{-1}(f(x)) \quad (40)$$

$$\phi(x) = f(x) \quad (41)$$

Finally, substituting  $\phi(x)$  and  $\phi^{-1}(x)$  back into the equation for  $\psi$ , we get:

$$\psi(a, b) = \phi^{-1}(\phi(a) + \phi(b)) \quad (42)$$

$$\psi(a, b) = f^{-1}(f(a) + f(b)) \quad (43)$$

□

## B Training Plots

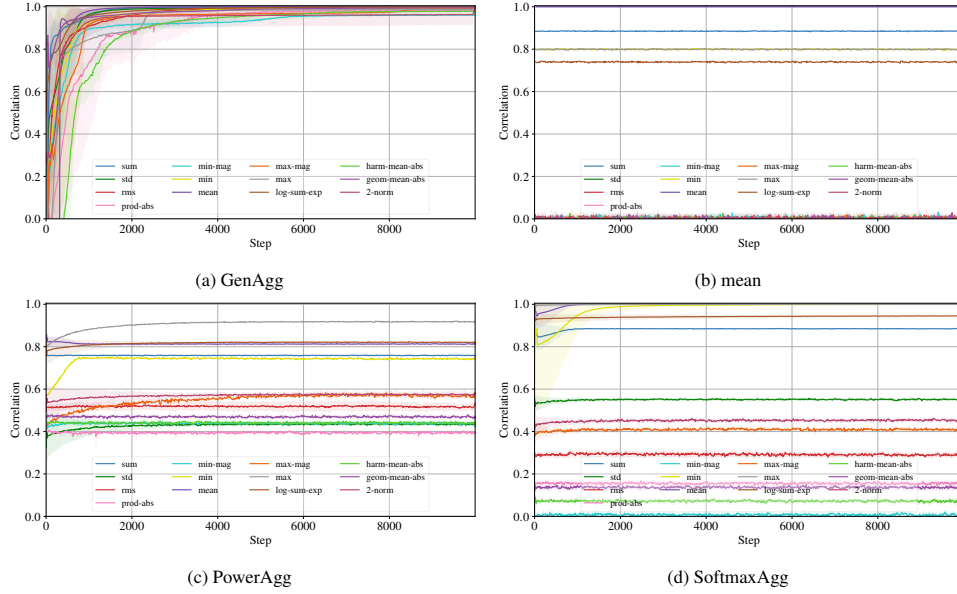


Figure 4: Training plots for the Aggregator Regression experiment (see Section 5.1). Each plot represents the ability of a parametrised aggregator  $\oplus$  to regress over all standard aggregators  $\odot_k \in \mathcal{A}$ . The plots show the mean and standard deviation of the correlation between the predicted and ground truth values over 10 trials.

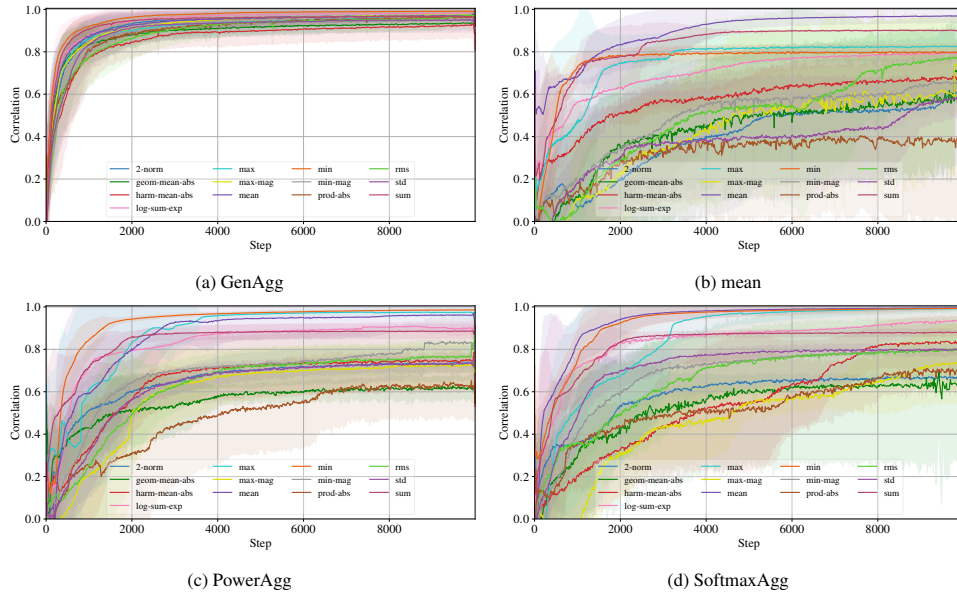


Figure 5: Training plots for the GNN Regression experiment (see Section 5.2). Each plot represents the ability of a GNN using parametrised aggregator  $\oplus$  to regress over all standard aggregators  $\odot_k \in \mathcal{A}$ . The plots show the mean and standard deviation of the correlation between the predicted and ground truth values over 10 trials.

## C Parametrisations

In this section, we use the augmented  $f$ -mean to show that GenAgg is theoretically capable of representing all of the standard aggregators  $\odot_k \in \mathcal{A}$ . For each standard aggregator  $\odot_k$ , we prove that there exists a parametrisation  $\theta$  such that  $\oplus_\theta = \odot_k$ . In a slight abuse of notation, mathematical operations applied to the input set  $\mathcal{X}$  (such as an absolute value or a geometric inverse) denote elementwise operations:  $f(\mathcal{X}) = \{f(x_1), \dots, f(x_n)\}$ .

**Theorem C.1. Mean.** For  $\theta = \langle f, \alpha, \beta \rangle = \langle x, 0, 0 \rangle$ , the augmented  $f$ -mean equals the mean  $\oplus_\theta = \frac{1}{n} \sum x_i$ .

*Proof.*

$$\bigoplus_{x_i \in \mathcal{X}}^{(x,0,0)} x_i = \left( n^{0-1} \sum_{x_i \in \mathcal{X}} (x_i - 0 \cdot \mu) \right) \quad (44)$$

$$= \frac{1}{n} \sum_{x_i \in \mathcal{X}} x_i \quad (45)$$

□

**Theorem C.2. Sum.** For  $\theta = \langle f, \alpha, \beta \rangle = \langle x, 1, 0 \rangle$ , the augmented  $f$ -mean equals the sum  $\oplus_\theta = \sum x_i$ .

*Proof.*

$$\bigoplus_{x_i \in \mathcal{X}}^{(x,1,0)} x_i = \left( n^{1-1} \sum_{x_i \in \mathcal{X}} (x_i - 0 \cdot \mu) \right) \quad (46)$$

$$= \sum_{x_i \in \mathcal{X}} x_i \quad (47)$$

□

**Theorem C.3. Product.** For  $\theta = \langle f, \alpha, \beta \rangle = \langle \log(|x|), 1, 0 \rangle$ , the augmented  $f$ -mean equals the product  $\oplus_\theta = \prod |x_i|$ .

*Proof.*

$$\bigoplus_{x_i \in \mathcal{X}}^{(\log(|x|),1,0)} x_i = e^{\left( n^{1-1} \sum_{x_i \in \mathcal{X}} \log(|x_i - 0 \cdot \mu|) \right)} \quad (48)$$

$$= e^{\sum_{x_i \in \mathcal{X}} \log(|x_i|)} \quad (49)$$

$$= \prod_{x_i \in \mathcal{X}} e^{\log(|x_i|)} \quad (50)$$

$$= \prod_{x_i \in \mathcal{X}} |x_i| \quad (51)$$

□

**Theorem C.4. Max Magnitude.** For  $\theta = \langle f, \alpha, \beta \rangle = \langle \lim_{p \rightarrow \infty} |x|^p, 0, 0 \rangle$ , the augmented  $f$ -mean equals the max magnitude  $\oplus_\theta = \max(|\mathcal{X}|)$ .

*Proof.*

$$\bigoplus_{x_i \in \mathcal{X}}^{(\lim_{p \rightarrow \infty} |x|^p, 0, 0)} x_i = \lim_{p \rightarrow \infty} \left( n^{0-1} \sum_{x_i \in \mathcal{X}} |x_i - 0 \cdot \mu|^p \right)^{\frac{1}{p}} \quad (52)$$

$$= \lim_{p \rightarrow \infty} \frac{1}{n^{\frac{1}{p}}} \left( \sum_{x_i \in \mathcal{X}} |x_i|^p \right)^{\frac{1}{p}} \quad (53)$$

$$= \lim_{p \rightarrow \infty} \left( \sum_{x_i \in \mathcal{X}} |x_i|^p \right)^{\frac{1}{p}} \quad (54)$$

This aggregator is composed of monotonic functions, so if every element  $x_i$  is substituted with  $\max(|\mathcal{X}|)$ , then the output should increase. Therefore, we can write the following inequality:

$$\lim_{p \rightarrow \infty} \left( \sum_{x_i \in \mathcal{X}} |x_i|^p \right)^{\frac{1}{p}} \leq \lim_{p \rightarrow \infty} \left( \sum_{i \in [1..n]} \max(|\mathcal{X}|)^p \right)^{\frac{1}{p}} \quad (55)$$

$$\lim_{p \rightarrow \infty} \left( \sum_{x_i \in \mathcal{X}} |x_i|^p \right)^{\frac{1}{p}} \leq \lim_{p \rightarrow \infty} n^{\frac{1}{p}} \cdot (\max(|\mathcal{X}|)^p)^{\frac{1}{p}} \quad (56)$$

$$\lim_{p \rightarrow \infty} \left( \sum_{x_i \in \mathcal{X}} |x_i|^p \right)^{\frac{1}{p}} \leq \max(|\mathcal{X}|) \quad (57)$$

$$(58)$$

Similarly, since sum is monotonic, the value of the output computed over a set  $\bigoplus_{\theta}(\mathcal{X})$  is greater than the output computed over a single element from that set  $\bigoplus_{\theta}(\{x_i\})$ , where  $x_i \in \mathcal{X}$ . Furthermore, we know that  $\max(|\mathcal{X}|)$  is one of the elements of  $\mathcal{X}$ . So, we can write the following inequality:

$$\lim_{p \rightarrow \infty} \left( \sum_{x_i \in \mathcal{X}} |x_i|^p \right)^{\frac{1}{p}} \geq \lim_{p \rightarrow \infty} (\max(|\mathcal{X}|)^p)^{\frac{1}{p}} \quad (59)$$

$$\lim_{p \rightarrow \infty} \left( \sum_{x_i \in \mathcal{X}} |x_i|^p \right)^{\frac{1}{p}} \geq \max(|\mathcal{X}|) \quad (60)$$

$$(61)$$

Consequently, by the squeeze theorem, we can state:

$$\lim_{p \rightarrow \infty} \left( \sum_{x_i \in \mathcal{X}} |x_i|^p \right)^{\frac{1}{p}} = \max(|\mathcal{X}|) \quad (62)$$

□

**Theorem C.5. Min Magnitude.** For  $\theta = \langle f, \alpha, \beta \rangle = \langle \lim_{p \rightarrow \infty} |x|^{-p}, 0, 0 \rangle$ , the augmented  $f$ -mean equals the min magnitude  $\bigoplus_{\theta} = \min(|\mathcal{X}|)$ .

*Proof.*

$$\bigoplus_{x_i \in \mathcal{X}} (\lim_{p \rightarrow \infty} |x|^{-p}, 0, 0) x_i = \lim_{p \rightarrow \infty} \left( n^{0-1} \sum_{x_i \in \mathcal{X}} |x_i - 0 \cdot \mu|^{-p} \right)^{-\frac{1}{p}} \quad (63)$$

$$= \lim_{p \rightarrow \infty} \frac{1}{n^{-\frac{1}{p}}} \left( \sum_{x_i \in \mathcal{X}} |x_i|^{-p} \right)^{-\frac{1}{p}} \quad (64)$$

$$= \lim_{p \rightarrow \infty} \left( \sum_{x_i \in \mathcal{X}} |x_i|^{-p} \right)^{-\frac{1}{p}} \quad (65)$$

We use the theorem for the parametrisation of max magnitude (Theorem C.4) as a lemma for this proof. By applying a monotonically decreasing transformation to the inputs and then inverting that transformation on the output, we can write the min as a function of the max. Since the inputs are restricted to the positive domain with the absolute value, we select the transformation  $T(x) = \frac{1}{x}$ :

$$\min(|\mathcal{X}|) = \frac{1}{\max(\frac{1}{|\mathcal{X}|})} \quad (66)$$

Substituting the parametrisation from Theorem C.4 for max, we get:

$$\min(|\mathcal{X}|) = \frac{1}{\lim_{p \rightarrow \infty} \left( \sum_{x_i \in \mathcal{X}} \left( \frac{1}{|x_i|} \right)^p \right)^{\frac{1}{p}}} \quad (67)$$

$$= \lim_{p \rightarrow \infty} \left( \sum_{x_i \in \mathcal{X}} |x_i|^{-p} \right)^{-\frac{1}{p}} \quad (68)$$

□

**Theorem C.6. Max.** For  $\theta = \langle f, \alpha, \beta \rangle = \langle \lim_{p \rightarrow \infty} e^{px}, 0, 0 \rangle$ , the augmented  $f$ -mean equals the max:  $\bigoplus_{\theta} = \max(\mathcal{X})$ .

*Proof.*

$$\bigoplus_{x_i \in \mathcal{X}}^{(\lim_{p \rightarrow \infty} e^{px}, 0, 0)} x_i = \lim_{p \rightarrow \infty} \frac{1}{p} \log \left( n^{0-1} \sum_{x_i \in \mathcal{X}} e^{p(x_i - 0 \cdot \mu)} \right) \quad (69)$$

$$= \lim_{p \rightarrow \infty} \frac{1}{p} \log \left( \frac{1}{n} \right) + \frac{1}{p} \log \left( \sum_{x_i \in \mathcal{X}} e^{p \cdot x_i} \right) \quad (70)$$

$$= \lim_{p \rightarrow \infty} \log \left( \left( \sum_{x_i \in \mathcal{X}} e^{p \cdot x_i} \right)^{\frac{1}{p}} \right) \quad (71)$$

$$(72)$$

This aggregator is composed of monotonic functions, so if every element  $x_i$  is substituted with  $\max(\mathcal{X})$ , then the output should increase. Therefore, we can write the following inequality:

$$\lim_{p \rightarrow \infty} \log \left( \left( \sum_{x_i \in \mathcal{X}} e^{p \cdot x_i} \right)^{\frac{1}{p}} \right) \leq \lim_{p \rightarrow \infty} \log \left( \left( \sum_{x_i \in \mathcal{X}} e^{p \cdot \max(\mathcal{X})} \right)^{\frac{1}{p}} \right) \quad (73)$$

$$\leq \lim_{p \rightarrow \infty} \log \left( n^{\frac{1}{p}} \cdot \left( e^{p \cdot \max(\mathcal{X})} \right)^{\frac{1}{p}} \right) \quad (74)$$

$$\leq \log \left( e^{\max(\mathcal{X})} \right) \quad (75)$$

$$\leq \max(\mathcal{X}) \quad (76)$$

$$(77)$$

Similarly, since sum is monotonic, the value of the output computed over a set  $\bigoplus_{\theta}(\mathcal{X})$  is greater than the output computed over a single element from that set  $\bigoplus_{\theta}(\{x_i\})$ , where  $x_i \in \mathcal{X}$ . Furthermore, we know that  $\max(\mathcal{X})$  is one of the elements of  $\mathcal{X}$ . So, we can write the following inequality:

$$\lim_{p \rightarrow \infty} \log \left( \left( \sum_{x_i \in \mathcal{X}} e^{p \cdot x_i} \right)^{\frac{1}{p}} \right) \geq \lim_{p \rightarrow \infty} \log \left( \left( e^{p \cdot \max(\mathcal{X})} \right)^{\frac{1}{p}} \right) \quad (78)$$

$$\geq \log \left( e^{\max(\mathcal{X})} \right) \quad (79)$$

$$\geq \max(\mathcal{X}) \quad (80)$$

$$(81)$$

Consequently, by the squeeze theorem, we can state:

$$\lim_{p \rightarrow \infty} \log \left( \left( \sum_{x_i \in \mathcal{X}} e^{p \cdot x_i} \right)^{\frac{1}{p}} \right) = \max(\mathcal{X}) \quad (82)$$

□

**Theorem C.7. Min.** For  $\theta = \langle f, \alpha, \beta \rangle = \langle \lim_{p \rightarrow \infty} e^{-px}, 0, 0 \rangle$ , the augmented  $f$ -mean equals the min:  $\oplus_{\theta} = \min(\mathcal{X})$ .

*Proof.*

$$\oplus_{x_i \in \mathcal{X}}^{\langle \lim_{p \rightarrow \infty} e^{-px}, 0, 0 \rangle} x_i = \lim_{p \rightarrow \infty} -\frac{1}{p} \log \left( n^{0-1} \sum_{x_i \in \mathcal{X}} e^{-p(x_i - 0 \cdot \mu)} \right) \quad (83)$$

$$= \lim_{p \rightarrow \infty} -\frac{1}{p} \log \left( \frac{1}{n} \right) - \frac{1}{p} \log \left( \sum_{x_i \in \mathcal{X}} e^{-p \cdot x_i} \right) \quad (84)$$

$$= \lim_{p \rightarrow \infty} \log \left( \left( \sum_{x_i \in \mathcal{X}} e^{-p \cdot x_i} \right)^{-\frac{1}{p}} \right) \quad (85)$$

$$(86)$$

We use the theorem for the parametrisation of max (Theorem C.6) as a lemma for this proof. By applying a monotonically decreasing transformation to the inputs and then inverting that transformation on the output, we can write the min as a function of the max. We select the transformation  $T(x) = -x$ :

$$\min(\mathcal{X}) = -\max(-\mathcal{X}) \quad (87)$$

Substituting the parametrisation from Theorem C.6 for max, we get:

$$\min(\mathcal{X}) = \lim_{p \rightarrow \infty} -\log \left( \left( \sum_{x_i \in \mathcal{X}} e^{p \cdot (-x_i)} \right)^{\frac{1}{p}} \right) \quad (88)$$

$$= \lim_{p \rightarrow \infty} \log \left( \left( \sum_{x_i \in \mathcal{X}} e^{-p \cdot x_i} \right)^{-\frac{1}{p}} \right) \quad (89)$$

□

**Theorem C.8. Harmonic Mean.** For  $\theta = \langle f, \alpha, \beta \rangle = \langle \frac{1}{x}, 0, 0 \rangle$ , the augmented  $f$ -mean equals the harmonic mean  $\oplus_{\theta} = \frac{n}{\sum \frac{1}{x_i}}$ .

*Proof.*

$$\oplus_{x_i \in \mathcal{X}}^{\langle \frac{1}{x}, 0, 0 \rangle} x_i = \left( n^{0-1} \sum_{x_i \in \mathcal{X}} (x_i - 0 \cdot \mu)^{-1} \right)^{-1} \quad (90)$$

$$= \left( \frac{1}{n} \sum_{x_i \in \mathcal{X}} \frac{1}{x_i} \right)^{-1} \quad (91)$$

$$= \frac{n}{\sum_{x_i \in \mathcal{X}} \frac{1}{x_i}} \quad (92)$$

□



**Theorem C.9. Geometric Mean.** For  $\theta = \langle f, \alpha, \beta \rangle = \langle \log(|x|), 0, 0 \rangle$ , the augmented  $f$ -mean equals the geometric mean  $\oplus_{\theta} = \sqrt[n]{\prod |x_i|}$ .

*Proof.*

$$\oplus_{x_i \in \mathcal{X}}^{\langle \log(|x|), 0, 0 \rangle} x_i = e^{\left( n^{0-1} \sum_{x_i \in \mathcal{X}} \log(|x_i - 0 \cdot \mu|) \right)} \quad (93)$$

$$= e^{\frac{1}{n} \sum_{x_i \in \mathcal{X}} \log(|x_i|)} \quad (94)$$

$$= \prod_{x_i \in \mathcal{X}} e^{\frac{1}{n} \log(|x_i|)} \quad (95)$$

$$= \prod_{x_i \in \mathcal{X}} e^{\log(|x_i|^{\frac{1}{n}})} \quad (96)$$

$$= \prod_{x_i \in \mathcal{X}} |x_i|^{\frac{1}{n}} \quad (97)$$

$$= \sqrt[n]{\prod_{x_i \in \mathcal{X}} |x_i|} \quad (98)$$

□

**Theorem C.10. Root Mean Square.** For  $\theta = \langle f, \alpha, \beta \rangle = \langle x^2, 0, 0 \rangle$ , the augmented  $f$ -mean equals the root mean square  $\oplus_{\theta} = \sqrt{\frac{1}{n} \sum x_i^2}$ .

*Proof.*

$$\oplus_{x_i \in \mathcal{X}}^{\langle x^2, 0, 0 \rangle} x_i = \left( n^{0-1} \sum_{x_i \in \mathcal{X}} (x_i - 0 \cdot \mu)^2 \right)^{\frac{1}{2}} \quad (99)$$

$$= \sqrt{\frac{1}{n} \sum_{x_i \in \mathcal{X}} x_i^2} \quad (100)$$

$$(101)$$

□

**Theorem C.11. Euclidean Norm.** For  $\theta = \langle f, \alpha, \beta \rangle = \langle x^2, 1, 0 \rangle$ , the augmented  $f$ -mean equals the euclidean norm  $\oplus_{\theta} = \sqrt{\sum x_i^2}$ .

*Proof.*

$$\oplus_{x_i \in \mathcal{X}}^{\langle x^2, 1, 0 \rangle} x_i = \left( n^{1-1} \sum_{x_i \in \mathcal{X}} (x_i - 0 \cdot \mu)^2 \right)^{\frac{1}{2}} \quad (102)$$

$$= \sqrt{\sum_{x_i \in \mathcal{X}} x_i^2} \quad (103)$$

$$(104)$$

□

**Theorem C.12. Standard Deviation.** For  $\theta = \langle f, \alpha, \beta \rangle = \langle x^2, 0, 1 \rangle$ , the augmented  $f$ -mean equals the standard deviation  $\oplus_{\theta} = \sqrt{\sum (x_i - \mu)^2}$ .

*Proof.*

$$\oplus_{x_i \in \mathcal{X}}^{\langle x^2, 0, 1 \rangle} x_i = \left( n^{0-1} \sum_{x_i \in \mathcal{X}} (x_i - 1 \cdot \mu)^2 \right)^{\frac{1}{2}} \quad (105)$$

$$= \sqrt{\frac{1}{n} \sum_{x_i \in \mathcal{X}} (x_i - \mu)^2} \quad (106)$$

$$(107)$$

□

**Theorem C.13. Log-Sum-Exp.** For  $\theta = \langle f, \alpha, \beta \rangle = \langle e^x, 0, 1 \rangle$ , the augmented  $f$ -mean equals the log-sum-exp  $\oplus_{\theta} = \log(\sum e^{x_i})$ .

*Proof.*

$$\oplus_{x_i \in \mathcal{X}}^{(e^x, 1, 0)} x_i = \log \left( n^{1-1} \sum_{x_i \in \mathcal{X}} e^{x_i - 0 \cdot \mu} \right) \quad (108)$$

$$= \log \left( \sum_{x_i \in \mathcal{X}} e^{x_i} \right) \quad (109)$$

$$(110)$$

□

## D Limitations

In our problem statement, we define a set of special cases  $\mathcal{A}$  which we refer to as “standard aggregators”. Our regression experiments analyse the representational capacity of various methods by analysing their ability to regress over the aggregators in  $\mathcal{A}$ . However, we acknowledge that this set is not exhaustive, so there may exist special cases not in  $\mathcal{A}$  which are useful or could provide some additional insight.

Our GNN benchmark experiments also only provide data about the performance of GenAgg in a limited number of applications. We evaluate on the GNN benchmark datasets suite from [6], which includes MNIST, CIFAR10, CLUSTER, and PATTERN. These datasets provide a mix of node classification and graph classification tasks (to complement the regression tasks from our other experiments). We did not include the TSP and CSL datasets from the same GNN benchmarks suite because TSP is an edge classification task (which would necessitate significant modification of our GNN), and CSL requires node positional encodings in order to be solvable by message passing GNNs, which it does not include by default [6]. In our initial testing, we also considered using the PUBMED, CORA, and CITESEER datasets. However, they are extremely small datasets that often lead to overfitting. Furthermore, there is something fundamental about the coauthor problem (upon which all three datasets are based) that fundamentally does not require the same level of complexity to solve—the train accuracy of all methods, including simple aggregators, approaches 1. Instead of these coauthor datasets, we opted to use the GNN benchmark dataset, which seemed to present more “difficult” problems. However, for the sake of transparency, we include our results on the datasets that we decided not to use:

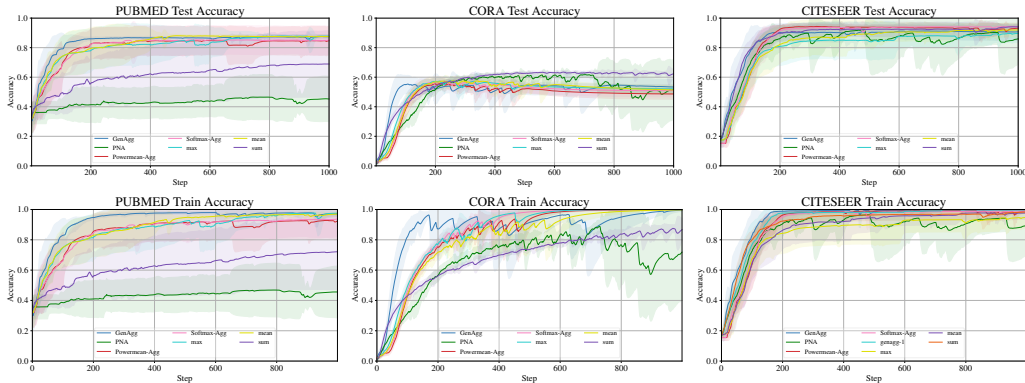


Figure 6: Train and Test Accuracy for all baselines on PUBMED, CORA, and CITESEER.

In these small datasets, the primary problem across all models is overfitting. GenAgg performs on-par with the best baseline, but it does not exhibit a performance *boost* as it does in the other datasets. To determine if the lack of improvement is due to the overfitting or the dataset itself, we run an additional experiment that explicitly examines the effect of overfitting by artificially reducing the amount of training data (Figure 7):

	GenAgg	Best Baseline	Median Baseline		GenAgg	Best Baseline	Median Baseline	
	100%	<b>0.915</b>	0.872	0.841	100%	<b>0.926</b>	0.897	0.866
	10%	<b>0.958</b>	0.903	0.846	10%	<b>0.905</b>	0.854	0.829
	1%	<b>0.978</b>	0.846	0.832	1%	<b>0.903</b>	0.832	0.830

(a) Train Accuracy

(b) Test Accuracy

Figure 7: Train and test accuracy of GenAgg vs all baselines on various subsets of the PATTERN dataset. We train on the PATTERN dataset (100% of the data), and versions with 10% and 1% of the original datapoints.

This experiment highlights a difference between the small coauthor datasets and the reduced PATTERN dataset (Figure 7). In the coauthor datasets the train accuracies approach 1, whereas in the reduced PATTERN dataset the baselines do not surpass a certain level of performance. This indicates that the mathematical relationships in the PATTERN dataset are fundamentally more difficult to represent. It is in these more complex problems that GenAgg provides the most benefit.

## E Training Details

In our implementation of GenAgg, we implement  $f$  and  $f^{-1}$  as MLPs with hidden sizes of  $[1, 2, 2, 4]$  and  $[4, 2, 2, 1]$  using Mish activation, BatchNorm, and Kaiming Normal weight initialisation. To run the GNN benchmark experiments, we use a 4-layer GraphConv model with a hidden size of 64, using Mish activation between layers. MLPs are used as pre- and post- processors in order to map to and from the hidden dimension of the GNN. The preprocessor is implemented with a one layer MLP, and the postprocessor is implemented with a 4-layer MLP using Mish activation. In tasks which require graph-level predictions, we prepend a global mean pooling layer to the postprocessor.

We run all of our experiments on an NVIDIA GeForce GTX 1080 Ti GPU. In all experiments, we use the Adam optimiser with a learning rate of  $10^{-3}$ . In the regression experiments we train for 10,000 epochs with a batch size of 1024, and in the GNN benchmark experiment we train for 1,000 epochs with a batch size of 32. Our results report the mean (and standard deviation, as the shaded region in the training plots) over 10 trials.

## F Runtime

As GenAgg requires running forward passes of small neural networks in addition to performing a sum, it incurs an additional runtime cost. We report the runtime overhead of each method in the figure below:

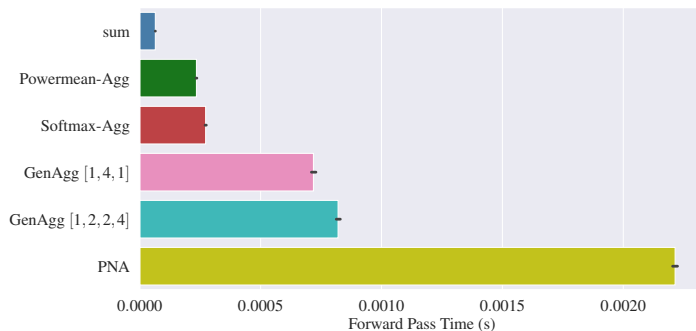


Figure 8: Forward pass times for each aggregation method. GenAgg  $[1, 2, 2, 4]$  and GenAgg  $[1, 4, 1]$  denote two different layer architectures for  $f$ . The data in this figure represents the time for a single forward pass over the MNIST GNN Benchmark dataset with a batch size of 1024 (approximately 578k edges), using an NVIDIA GeForce GTX 1080 Ti GPU. The reported time is *only* for the aggregation component, not the GNN as a whole.

While the the absolute runtime of GenAgg is relatively fast, it is still significantly slower than sum. Consequently, it is possible that the runtime of our implementation can preclude its use in time-critical applications. However, note that this figure only represents the runtime for our specific implementation—GenAgg can be implemented with any invertible function  $f$ . Using a symbolic parametrisation of invertible functions can significantly speed up computations (at the cost of representational complexity). Alternatively, it is likely that an implementation in JAX with compiled networks  $f$  and  $f^{-1}$  can achieve a speed boost with the same architecture.