

# CHARACTERISTIC FUNCTION-BASED REGULARIZATION FOR PROBABILITY FUNCTION INFORMED NEURAL NETWORKS

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Regularization is essential in neural network training to prevent overfitting and improve generalization. In this paper, we propose a novel regularization technique that leverages decomposable distribution and central limit theory assumptions by exploiting the properties of characteristic functions. We first define Probability Function Informed Neural Networks as a class of universal function approximators capable of embedding the knowledge of some probabilistic rules constructed over a given dataset into the learning process (a similar concept to Physics-informed neural networks (PINNs), if the reader is familiar with those). We then enforce a regularization framework over this network, aiming to impose structural constraints on the network’s weights to promote greater generalizability in the given probabilistic setting. Rather than replacing traditional regularization methods such as L2 or dropout, our approach is intended to supplement this and other similar classes of neural network architectures by providing instead a contextual delta of generalization. We demonstrate that integrating this method into such architectures helps improve performance on benchmark supervised classification datasets, by preserving essential distributional properties to mitigate the risk of overfitting. This characteristic function-based regularization offers a new perspective for enhancing distribution-aware learning in machine learning models.

## 1 AN INTRODUCTION AND (RATHER INFORMAL) DEFINITION OF PROBABILITY FUNCTION INFORMED NEURAL NETWORKS

Let’s define Probability Function Informed Neural Networks (PFINNs) as a class of universal function approximators designed to integrate probabilistic knowledge into the learning process and more specifically into the learning architecture. Reformulating concepts from Physics-Informed Neural Networks (PINNs), which leverage physical laws to guide the training of neural networks (Cuomo et al., 2022), we can say that PFINNs focus on embedding probabilistic rules into a networks construction from the practitioner’s understanding of a given dataset.

At their core, PFINNs aim to enhance the learning of complex relationships and mappings by utilising the underlying probability distributions that characterize the data. This allows the model to incorporate essential statistical information, which can lead to some level of interpretability. By explicitly encoding probabilistic principles, such as conditional dependencies or marginal distributions, into the neural network architecture, PFINNs should be able better navigate the search space of the data they are trained on to give more reasonable results as a functional approximator.

In simpler terms, PFINNs can be thought of as neural networks that not only learn from data but also respect some probabilistic structures the practitioner would like to embed into the network that they believe governs that data. This makes them particularly powerful for applications where uncertainty plays a significant role, allowing practitioners to build models that are not only predictive but also statistically sound for their context.

To explore the existence of PFINNs as (neuro)symbolic AI and as hybrid systems would require significantly more than 9 pages in a manuscript. We believe instead, it is best to begin with a simple PFINN architecture example derived from the lens of generalising the MNIST dataset and

054 associated classification task. Where the main focus of this paper will be on presenting a general  
 055 characteristic function-based regularization method for contextual regularization, aimed at making  
 056 these interpretable models—particularly from a probabilistic perspective—more generalizable via  
 057 “relaxation” of the model through this regularisation.

058 .  
 059 .  
 060 .  
 061 **2 MATHEMATICAL CONSTRUCTION OF A TYPE OF SIMPLE PFINN FROM**  
 062 **GENERALISATION OF THE MNIST DATASET AND ASSOCIATED**  
 063 **CLASSIFICATION PROBLEM**  
 064 .

065 The MNIST dataset is a widely used benchmark in the field of machine learning, consisting of  
 066 images of handwritten digits. Mathematically, we can describe the dataset and the classification  
 067 problem as follows:  
 068 .

069 **2.1 DATASET DESCRIPTION**  
 070 .

071 The MNIST dataset can be formally defined as a collection of pairs:

$$072 \mathcal{D} = \{(x_i, y_i)\}_{i=1}^N$$

073 where:  
 074 .

- 075 •  $x_i \in \mathbb{R}^{28 \times 28}$  represents a grayscale image of a handwritten digit, with each image being a  
 076  $28 \times 28$  pixel array.
- 077 •  $y_i \in \{0, 1, 2, \dots, 9\}$  denotes the corresponding label, indicating the digit represented in  
 078 the image  $x_i$ .
- 079 •  $N$  is the total number of samples in the dataset, typically  $N = 60,000$  for the training set  
 080 and  $N = 10,000$  for the test set.  
 081 .

082 Each image  $x_i$  can be flattened into a vector:

$$083 x_i \in \mathbb{R}^{784}$$

084 by concatenating the rows of the  $28 \times 28$  array.  
 085 .

086 .  
 087 **2.2 CLASSIFICATION PROBLEM**  
 088 .

089 The goal of the classification problem is to learn a mapping  $f : \mathbb{R}^{784} \rightarrow \{0, 1, 2, \dots, 9\}$  such that  
 090 for a given input image  $x_i$ , the model predicts the correct digit label  $y_i$ .  
 091 .

092 This is usually formulated as a supervised learning problem. For the sake of succinctness and sim-  
 093 plicity to understand methods introduced later in the paper more easily, we will make our first  
 094 assumption about the data in this dataset :

095 **Assumption 1** (Linear Separability). We assume that the classes in this dataset are linearly separa-  
 096 ble. This means that there exists a hyperplane defined by the equation:  
 097 .

$$098 \mathbf{w}^T \mathbf{x} + \mathbf{b} = 0$$

099 such that for all instances  $\mathbf{x}_i$  belonging to class  $C_1$ , the following condition holds:  
 100 .

$$101 \mathbf{w}^T \mathbf{x}_i + \mathbf{b} > 0,$$

102 and for all instances  $\mathbf{x}_j$  belonging to class  $C_2$ , the following condition holds:  
 103 .

$$104 \mathbf{w}^T \mathbf{x}_j + \mathbf{b} < 0.$$

105 Here,  $\mathbf{w}$  is the weight vector,  $\mathbf{b}$  is the bias, and  $\mathbf{x}$  represents the feature vectors of the dataset.  
 106 .  
 107 .

**Advantages of this Assumption:** This means we can construct a network without any hidden layers where each input node is directly connected to the output node.

**Associated Pitfall to be noted :** Without hidden layers, a neural network is limited to approximating only linear functions. Consequently, such a network does not qualify as a true "universal approximator" and hence is not a perfect PFINN by definition. However, this design choice, is intentional; it simplifies the methodology and reduces obfuscation, making conveying core ideas we wish to demonstrate later in this paper more accessible. Also any practitioner can easily incorporate hidden layers with appropriate activation functions, thereby extending this architecture to a universal approximator; granted they have developed a general understanding of the methods presented here and possess the requisite knowledge of the universal approximation theorem (Hornik et al., 1989) (Cybenko, 1989), as well as the principles necessary for constructing architectures that align with the definitions outlined in the proofs.

### 2.3 NETWORK CONSTRUCTION

Taking this structure our architecture is as follows:

- **Input Layer:** The input layer consists of 784 neurons, corresponding to the  $28 \times 28$  pixel images in the MNIST dataset. Each input neuron  $x_i$  represents the pixel intensity value of the image:

$$x_i \in [0, 1], \quad \text{for } i = 1, 2, \dots, 784$$

- **Output Layer:** The output layer consists of 10 neurons, corresponding to the 10 possible digit classes (0 through 9). The output for each neuron  $o_k$  in the output layer can be expressed as:

$$o_k = \sum_{j=1}^{N=784} w_{jk}x_j + b_k, \quad \text{for } k = 0, 1, \dots, 9$$

where  $w_{jk}$  are the weights connecting input neurons  $x_j$  to output neurons  $o_k$ . And  $b_k$  is the bias for the output neuron  $o_k$ .

### 2.4 EMBEDDING A PROBABILITY STRUCTURE

Now we apply a function,  $f(\cdot)$ , which can be chosen based on what feels best for the problem to the  $o_k$ ; for example for our approach we are going to consider the softmax function, converting the raw output scores into probabilities  $\in [0, 1]$ :

$$\mathbb{P}(y = k|x) = \frac{e^{o_k}}{\sum_{m=0}^9 e^{o_m}}$$

Essentially, what we have accomplished with this architecture is the construction of an approximation for the output probability  $p$  defined as:

$$p = f(\mathbf{w}^T \mathbf{x} + \mathbf{b})$$

where  $\mathbf{w}$  is the weight vector,  $\mathbf{x}$  is the input vector, and  $\mathbf{b}$  is the bias term. We will utilize this approximation under the assumption that each node in the output layer corresponds to a random variable  $\alpha_i$  that follows a Bernoulli distribution:

$$\alpha_i \sim \text{Bernoulli}(p_i)$$

This means that each output node generates a binary outcome, determining whether the output corresponds to a specific class label based on the computed probability  $p_i$ .

To simplify, we can conceptualize this process as a series of binary questions, where each variable represents a yes-or-no inquiry about whether the current input corresponds to a particular number. For instance, we could have a series of questions structured as follows:

"Is this the digit 1?" ; "Is this the digit 2?" ... ; "Is this the digit 9?"

In this context, the response to each question is influenced by the associated probability  $p$ , which is derived from the linear combination of inputs and weights. Each output decision is thus determined probabilistically by the computed  $p$ , providing a framework for classification over a probabilistic structure based on the given inputs.

### 3 MOTIVATION FOR THE REGULARISATION

The first question to address is whether these  $N$  questions (in the case of MNIST  $N = 10$ ) are sufficient enough for generate a function that can capture a perfect solution space from the training to generalise over unseen data later. Maybe other potential questions may merit consideration? Given that we can reformulate these new questions into binary (yes/no) formats, they could serve to expand the function we derive from the PFINN to get a better picture of the true distribution and solution space.

Treating each question as statistically independent, defined informally such that the response to one question does not influence the probabilities associated with another. For instance, if we pose questions such as "Is this 10?", "Is this a goose?", "Is it currently raining?", "Was this instance run on Alan Turing's computer?", "Was this instance run on Ada Lovelace's Analytical Engine?", or "Was this run on Claude Shannon's GPAC?", we can generate an infinite series of inquiries. The sheer volume of potential questions of course poses a practical impossibility for storage but nevertheless, we can leverage the conceptual framework inherent in this infinite set of questions (through the Central Limit Theorem (CLT), more specifically in our case Lyapunov's CLT), to be able to relax the function we find through the PFINN to "generalize" better, granted our process respects any required relevant assumptions.

Essentially for this, we consider the concept of decomposable distributions. The Bernoulli distributions we have currently are indecomposable so maybe finding and exploiting a relationship to an infinitely decomposable function like the normal distribution may yield some interesting results. It is through discovering this relationship, via us toying around with the linear combinations of Bernoulli variables which we have, that yields us the regularisation property; namely through exploiting their convergence to a normal distribution if we assume the existence of the previously introduced infinite question space.

### 4 REGULARISATION METHODOLOGY

To formalize the ideas presented above, we can derive a regularization method for the specific class of PFINN architecture described in Section 2 as follows :

**Definition 1** (Lyapunov Central Limit Theorem). Suppose we have a sequence of independent random variables,  $\{Y_1, Y_2, \dots, Y_n\}$ , each with finite expected value  $\mu_i$  and variance  $\sigma_i^2$ .

If we define the following sum of variances:

$$s_n^2 = \sum_{i=1}^n \sigma_i^2. \quad (1)$$

If  $\exists \delta > 0$ , such that Lyapunov's condition:

$$\lim_{n \rightarrow \infty} \frac{1}{s_n^{2+\delta}} \sum_{i=1}^n \mathbb{E} [|Y_i - \mu_i|^{2+\delta}] = 0, \quad (2)$$

is satisfied  $\implies$  the sum of the normalized variables  $\frac{Y_i - \mu_i}{s_n}$  converges in distribution to a standard normal random variable as  $n \rightarrow \infty$ :

$$\frac{1}{s_n} \sum_{i=1}^n (Y_i - \mu_i) \rightarrow N. \quad (3)$$

where  $N \sim \mathcal{N}(0, 1)$   
(Lyapunov, 1900)

We model each data point as being generated from a random variable  $\aleph$ , which is represented as an approximation of linear combination of Bernoulli variables  $X_i \sim \text{Bern}(p)$ . We will now demonstrate that properly formulating  $\aleph$  allows for convergence to  $\mathcal{N}(0, 1)$  as the number of Bernoulli variables increases sufficiently.

**Axiom 1.** Let us establish the following fundamental assumption that will underpin our framework. Given the true data distribution,  $\mathcal{D}$  (which can be conceptualized as the "Population Distribution" in statistical terms), we assert that the data points (the "Samples" we usually have in our finite dataset) are generated from a random variable  $\aleph$  such that:

$$\aleph \sim \mathcal{D} \quad (4)$$

**Definition 2.** The characteristic function  $\phi(u)$  of a random variable  $Y$  is defined as:

$$\phi(u) = \mathbb{E}[e^{\vartheta u Y}].^1 \quad (5)$$

**Definition 3.** We model the random variable  $\aleph$ , as a linear combination of Bernoulli Random Variables, defined as follows:

$$\aleph = \frac{1}{s_n} \sum_{i=1}^n (X_i - \mu_i), \quad (6)$$

where  $X_i \sim \text{Bern}(p_i) \implies \mu_i = \mathbb{E}[X_i] = p_i$  and  $s_n^2 = \sum_{i=1}^n \text{Var}[X_i] = \sum_{i=1}^n p_i(1 - p_i)$ .

**Proposition 1.** The characteristic function for  $\aleph$  can be computed as follows from 5 and 6:

$$\phi_{\mathcal{D}}(u) = \prod_{i=1}^n e^{\frac{(-\vartheta u p_i)}{\sqrt{\sum_{i=1}^n (p_i * (1-p_i))^2}} (1 - p_i)} + \prod_{i=1}^n e^{\frac{(\vartheta u (1-p_i))}{\sqrt{\sum_{i=1}^n (p_i * (1-p_i))^2}} (p_i)} \quad (7)$$

*Proof.*

$$\phi_{\mathcal{D}}(u) = \mathbb{E}[e^{\vartheta u \aleph}] = \mathbb{E}\left[e^{\vartheta u \frac{1}{s_n} \sum_{i=1}^n (X_i - \mu_i)}\right]. \quad (8)$$

Using the product law of exponents ( $a^n * a^m = a^{(n+m)}$ ), we rewrite the characteristic function:

$$= \mathbb{E}\left[\prod_{i=1}^n e^{\vartheta u \frac{1}{s_n} (X_i - \mu_i)}\right]. \quad (9)$$

Next, we separate it into the following by linearity of the expectation:

$$= \mathbb{E}\left[\prod_{i=1}^n e^{\vartheta u \frac{1}{s_n} (X_i - \mu_i)} \mathbb{I}\{X_i = 0\}\right] + \mathbb{E}\left[\prod_{i=1}^n e^{\vartheta u \frac{1}{s_n} (X_i - \mu_i)} \mathbb{I}\{X_i = 1\}\right]. \quad (10)$$

By properties of the Bernoulli Random Variable this is more precisely:

$$= \mathbb{E}\left[\prod_{i=1}^n e^{\vartheta u \frac{1}{s_n} (0 - p_i)} \mathbb{I}\{X_i = 0\}\right] + \mathbb{E}\left[\prod_{i=1}^n e^{\vartheta u \frac{1}{s_n} (1 - p_i)} \mathbb{I}\{X_i = 1\}\right]. \quad (11)$$

Using linearity of expectation, we have:

$$= \prod_{i=1}^n e^{\vartheta u \frac{1}{s_n} (-p_i)} \mathbb{E}[\mathbb{I}\{X_i = 0\}] + \prod_{i=1}^n e^{\vartheta u \frac{1}{s_n} (1-p_i)} \mathbb{E}[\mathbb{I}\{X_i = 1\}]. \quad (12)$$

By definition of Expectation of Indicator Function, we have:

$$= \prod_{i=1}^n e^{\vartheta u \frac{1}{s_n} (-p_i)} \mathbb{P}(X_i = 0) + \prod_{i=1}^n e^{\vartheta u \frac{1}{s_n} (1-p_i)} \mathbb{P}(X_i = 1). \quad (13)$$

<sup>1</sup>We deviate from common practice and use  $\vartheta$  instead of  $i$  to define the imaginary unit in a bid to reduce confusion as the letter  $i$  is used for indexing in much of the later proofs and writing

By properties of the Bernoulli Random Variable, this can be re-expressed as:

$$= \prod_{i=1}^n e^{\vartheta u \frac{1}{s_n} (0-p_i)} (1-p_i) + \prod_{i=1}^n e^{\vartheta u \frac{1}{s_n} (1-p_i)} p_i. \quad (14)$$

Reformulating the  $s_n^2$ :

$$\because s_n^2 = \sum_1^n \sigma_i^2 \implies \therefore s_n = \sqrt{\sum_1^n \sigma_i^2} = \sqrt{\sum_1^n (p_i(1-p_i))^2} \quad (15)$$

Thus, we conclude:

$$\implies \therefore \phi_{\mathcal{D}}(u) = \prod_{i=1}^n e^{\frac{(-\vartheta u p_i)}{\sqrt{\sum_{i=1}^n (p_i(1-p_i))^2}} (1-p_i)} + \prod_{i=1}^n e^{\frac{(\vartheta u (1-p_i))}{\sqrt{\sum_{i=1}^n (p_i(1-p_i))^2}} (p_i)} \quad (16)$$

□

**Corollary 1.** By the Lyapunov Central Limit Theorem, as  $n \rightarrow \infty$ , the characteristic function converges to that of the characteristic function of the normal distribution:

$$\because \mathcal{D} \rightarrow \mathcal{N}(0, 1) \implies \therefore \phi_{\mathcal{D}}(u) \rightarrow \phi_{\mathcal{N}(0,1)}(u). \quad (17)$$

Note this is true because rate of growth of the moments is constrained as per the Lyapunov condition, described in detail by proof outlined in the appendix A.1. Some graphics from numerical simulation of this effect is also attached in the appendix C.5, along with some helper graphs to visualise some transformations of characteristic functions as it is difficult to find some and there aren't really many online.

**Definition 4** (Regularization). Regularization is a technique used to prevent overfitting by adding a penalty term to the loss function. The regularized loss function is typically expressed as:

$$\min \sum_{i=1}^n L(\hat{y}_i, y_i) + \lambda R(f) \quad (18)$$

where  $\hat{y}_i = f(x_i)$  is the predicted output,  $L$  is the loss function,  $R(f)$  is the regularization term, and  $\lambda$  is a hyperparameter that controls the trade-off between model fit and complexity.

If we interpret  $\phi_{\mathcal{N}(0,1)}$  as a relaxed fit that our function can be adjusted towards, we can establish a regularization term  $R(f)$ , which levies a penalty on the complexity of model  $f$ , by adding a constraint through examining the difference between  $\phi_{\mathcal{D}}$  and  $\phi_{\mathcal{N}(0,1)}$ . This can be achieved by measuring the distance between the signals using:

$$\mathcal{R}(f) = d(\phi_{\mathcal{D}}, \phi_{\mathcal{N}(0,1)}) \quad (19)$$

The choice of the distance metric,  $d(\cdot)$ , is up to the practitioner but we briefly mention the ones we used for evaluation in the appendix B for reference.

Note also, that for a general class of PFINNs, one only needs to adjust the modeling of the random variable presented in Definition 3 to reformulate the equation in Proposition 1 accordingly.

## 5 NUMERICAL CONSIDERATIONS

The characteristic function is generally considered a "pure mathematical tool" whereby it's continuous nature presents significant challenges when implemented in discrete computational environments. Modern computers rely on finite precision arithmetic, which inherently restricts the exact representation of continuous functions, including characteristic functions. This means that we have to formulate discretizations based on some assumptions to integrate the characteristic function into

a practical regularization algorithm for machine learning, enabling it to operate and execute within finite time.

Specifically, we will have to restrict its domain to a "good enough" range since  $t \in \mathbb{R}$  and the associated infinite nature of the reals. In other words, abstractly the problem is then as follows (with the help of some informal proof sketches for the sake of brevity due to page limit and for the reader's sanity) :

**Proposition 2.** The set of real numbers  $\mathbb{R}$  is uncountably infinite.

**Informal Proof Sketch 1.** This can be shown using Cantor's famous Diagonal Argument. Assume for contradiction that  $\mathbb{R}$  is countable. Then we can list all real numbers in the interval  $[0, 1]$  as  $r_1, r_2, r_3, \dots$ . We can then construct a new real number  $r$  by taking the diagonal of this list and changing each digit, ensuring that  $r$  differs from each  $r_n$  at the  $n$ -th digit. Therefore,  $r$  cannot be in our original list, contradicting the assumption that we had listed all real numbers. Thus,  $\mathbb{R}$  is uncountably infinite. (Cantor, 1932)

**Proposition 3.** The set of real numbers  $\mathbb{R}$  is complete.

**Informal Proof Sketch 2.** The completeness of  $\mathbb{R}$  can be demonstrated using Dedekind's cuts. A Dedekind cut partitions the rational numbers into two non-empty sets  $A$  and  $B$ , where all elements of  $A$  are less than all elements of  $B$ . For any non-empty set of rationals that is bounded above, there exists a least upper bound (supremum) in  $\mathbb{R}$ . This property ensures that every Cauchy sequence of real numbers converges to a real number, establishing the completeness of  $\mathbb{R}$ . (Dedekind, 1912)

**Proposition 4.** The set of computable real numbers is countably infinite.

**Informal Proof Sketch 3.** The set of computable real numbers can be described as those numbers for which there exists a finite algorithm (Turing machine) that can produce their digits. Since the set of all finite algorithms is countable, it follows that the set of computable real numbers is also countable. (Bournez, 2024)(Weihrauch, 2012)

**Proposition 5.** The set of computable real numbers is not complete.

**Informal Proof Sketch 4.** To see this, consider the sequence of computable numbers defined by  $r_n = \frac{1}{n}$ , which converges to 0. Although 0 is a limit point of the sequence, it is not computable because there is no finite algorithm that can output the exact value of 0. This demonstrates that there exist Cauchy sequences of computable real numbers that do not converge to a computable limit, thereby showing that the set of computable real numbers is not complete. (Bournez, 2024)(Weihrauch, 2012)

It becomes evident that propositions 2, 3, 4, and 5 present significant challenges in computing the desired function  $\phi(t)$  especially on a Discrete Dynamical System like the modern computer we use. To address this, we have adopted a strategy of restricting to a finite domain of  $t \in [-2\pi, 2\pi]$  where we discretize this interval into  $n = 1000$  finite segments, which can be easily accomplished using a linear space function such as `numpy.linspace` or similar methods on modern programming languages.

The rationale for selecting the interval  $[-2\pi, 2\pi]$  is motivated by the analysis of the figures C.1 and C.3 in the appendix of the characteristic function for the standard normal distribution, as well as the set of convergence graphs for the  $\mathbb{N}$ -modelled linear combinations of Bernoulli random variables observed in C.5. The region of primary interest lies within this interval, and while any variations outside this interval may be potentially significant under certain circumstances, we can effectively treat them as an acceptable level of statistical noise, we are willing to quantify by some  $\epsilon$ . This allows us to disregard this noise in the context of testing viability, though it may come at the expense of some regularization "performance".

There is no universally "correct" range or sample size ( $n$ ); however, for the purposes of our experimentation, we consider this choice to be sufficient.

## 6 EVALUATION OF METHOD ON DATASETS

For evaluating our regularization approach, we consider the following 7 cases, no regularization (None), standard  $L^1$ ,  $L^2$  and  $L^\infty$  regularization and our  $\psi_1$ ,  $\psi_2$  and  $\psi_\infty$  regularization (as described in appendix section B). Other than MNIST, we built a similar structure of PFINNs for 4 other datasets

378 pertinent to classification tasks. It is worth noting for the PhiUSIL dataset, we have down sampled  
 379 to 7% of the original sample size, and reduced from the original 54 features to 6, to allow our tests  
 380 to run on weaker machines for reproducibility.

381 As for the loss function  $L(\cdot)$  as described in Equation 18, we used the cross-entropy loss. The  
 382 regularization parameter  $\lambda$  was consistently set to 0.01 across all test cases presented in Table 1.  
 383 Additionally, a learning rate of 0.1 was applied uniformly across all experiments.

384 Each result presented is derived from 3 independent runs of 100 epochs each. The reported values  
 385 represent the averages of each metric across these 3 runs, other than the Min and Max values indicate  
 386 the minimum and maximum results obtained from all 300 points of interest respectively.  
 387

Table 1: Test accuracy and associated metrics comparison table

Dataset	Metric	None	$L^1$	$L^2$	$L^\infty$	$\psi_1$	$\psi_2$	$\psi_\infty$
MNIST	Mean	<b>0.9245</b>	0.8068	0.9196	0.9239	0.9216	0.9243	<b>0.9245</b>
	Median	<b>0.9245</b>	0.8076	0.9197	0.9240	0.9217	0.9244	<b>0.9245</b>
	Std Dev	0.00146	<b>0.00660</b>	0.00143	0.00139	0.00196	0.00147	0.00146
	Avg Min	<b>0.9181</b>	0.7835	0.9148	0.9178	0.9135	0.9177	<b>0.9181</b>
	Avg Max	<b>0.9274</b>	0.8196	0.9228	0.9271	0.9254	0.9271	<b>0.9274</b>
	Min	0.9124	0.7805	0.9116	0.9124	0.9117	0.9124	<b>0.9125</b>
	Max	<b>0.9276</b>	0.8210	0.9230	0.9274	0.9260	<b>0.9276</b>	0.9275
HAR	Mean	<b>0.9502</b>	0.7308	0.9368	0.9487	0.9427	0.9500	<b>0.9502</b>
	Median	<b>0.9565</b>	0.7431	0.9511	0.9558	0.9522	<b>0.9565</b>	<b>0.9565</b>
	Std Dev	0.01751	<b>0.09931</b>	0.03406	0.01976	0.02648	0.01829	0.01746
	Avg Min	<b>0.8390</b>	0.4056	0.7607	0.8299	0.7889	0.8328	0.8389
	Avg Max	0.9613	0.8886	0.9627	0.9613	<b>0.9654</b>	0.9613	0.9612
	Min	<b>0.7112</b>	0.3244	0.6861	0.7099	0.6759	0.7000	0.7095
	Max	0.9630	0.8951	0.9637	0.9634	<b>0.9661</b>	0.9634	0.9630
WINE	Mean	0.5699	0.5713	0.5616	0.5673	0.5699	<b>0.5714</b>	0.5703
	Median	<b>0.5708</b>	<b>0.5708</b>	0.5620	0.5677	0.5703	<b>0.5708</b>	0.5703
	Std Dev	0.0069	0.0063	0.0066	0.0066	0.0069	<b>0.0076</b>	0.0069
	Avg Min	0.5396	0.5438	0.5302	0.5396	0.5385	<b>0.5458</b>	0.5417
	Avg Max	<b>0.5885</b>	<b>0.5885</b>	0.5750	0.5865	<b>0.5885</b>	0.5875	0.5875
	Min	0.5063	0.5094	0.4969	0.5063	0.5063	<b>0.5219</b>	0.5094
	Max	<b>0.5938</b>	<b>0.5938</b>	0.5781	0.5906	<b>0.5938</b>	0.5906	<b>0.5938</b>
Waveform	Mean	0.8723	0.8720	<b>0.8741</b>	0.8724	0.8727	0.8723	0.8723
	Median	0.8723	0.8723	<b>0.8750</b>	0.8728	0.8727	0.8727	0.8723
	Std Dev	0.0035	0.0038	<b>0.0042</b>	0.0037	0.00317	0.00352	0.00355
	Avg Min	0.8603	0.8600	<b>0.8620</b>	0.8597	<b>0.8620</b>	0.8600	0.8600
	Avg Max	0.8813	0.8800	<b>0.8823</b>	0.8817	0.8800	0.8813	0.8813
	Min	0.8580	0.8580	<b>0.8620</b>	0.8580	0.8610	0.8580	0.8570
	Max	<b>0.8830</b>	0.8810	<b>0.8830</b>	<b>0.8830</b>	0.8820	<b>0.8830</b>	<b>0.8830</b>
PhiUSIL	Mean	0.9244	0.9247	0.9118	0.9106	0.9244	0.9189	<b>0.9245</b>
	Median	0.9245	<b>0.9255</b>	0.9119	0.9094	0.9245	0.9184	<b>0.9255</b>
	Std Dev	0.0044	0.0041	0.0033	0.0034	0.0044	<b>0.0052</b>	0.0044
	Avg Min	<b>0.9124</b>	<b>0.9124</b>	0.9023	0.9013	<b>0.9124</b>	0.9074	<b>0.9124</b>
	Avg Max	<b>0.9305</b>	<b>0.9305</b>	0.9154	0.9184	<b>0.9305</b>	<b>0.9305</b>	<b>0.9305</b>
	Min	<b>0.9003</b>	<b>0.9003</b>	0.8943	0.8973	<b>0.9003</b>	0.8973	<b>0.9003</b>
	Max	<b>0.9305</b>	<b>0.9305</b>	0.9154	0.9184	<b>0.9305</b>	<b>0.9305</b>	<b>0.9305</b>

425  
 426  
 427 The most attractive metric in Table 1 would be the mean for each dataset. It is generally observed  
 428 that the mean for the regularisation we proposed, throughout 4 out of 5 datasets, achieve the highest  
 429 mean. It is also worth noting that the standard deviation is one of the lowest (except in the case of  
 430 PhiUSIL), which implies that there is a very minimal difference between each run. This property  
 431 might suggest the nature of "consistency" in the regularisation method which might yield better  
 results in the training and testing of larger datasets.

432 It can also be further observed that the limitations of our approach show up in the datasets waveform  
433 and HAR, which might also imply that our regularisation in it's current form could be improved to  
434 work better with time series data and noisy data. This could be a result of generalising noise as well,  
435 and perhaps could be further altered via some form of simple noise filtering whilst pre-processing  
436 the data.

437 As a quick aside we would like to caution the reader regarding empirical tests such as the one  
438 presented above. This is due to the myriad of factors that can influence overall performance; the  
439 results are also highly dataset-dependent. For instance, variations in the gradient starting point or  
440 learning rate can yield significantly different empirical outcomes for the same problem. There is  
441 no guaranteed method to achieve a perfect readout, as this would necessitate exploring an infinite  
442 search space, which is computationally impossible.

443 A natural question would be why was K-fold cross-validation not employed for evaluation. To  
444 address this, K-fold cross-validation while effective for approaching more optimal hyperparame-  
445 ters, truly involves exploration of an infinite search space for the "fairest" representation of the best  
446 metrics for a given method. Although it seemingly provides more robust performance estimates  
447 compared to a vanilla implementation, our focus in this study is on presenting a context-driven alter-  
448 native method rooted in probability theory rather than maximizing some raw performance metrics  
449 delta.

450 Given this approach, we believe that demonstrating reasonable performance with a random seed  
451 using simple constructions and standard default values suffices to validate our method. Furthermore,  
452 the relevant code is provided, allowing others to reproduce our results and adapt the methods to their  
453 own contexts, as detailed later in the reproducibility statement.

454 In essence, using different seeds for any given parameter can tell different stories, and even aggre-  
455 gating results over multiple seeds may lead to varied outcomes. The key takeaway is that while this  
456 approach may serve as a viable alternative for a specific dataset and problem, it might not always  
457 be the optimal choice. The purpose of such empirical tests is to demonstrate that they can be a  
458 'good enough' and 'reliable' option when necessary. Ultimately, it is the practitioner's responsi-  
459 bility to evaluate how well this approach aligns with their specific problem and to discern whether  
460 employing a particular tool makes sense in their context.

## 462 7 CONCLUSION

463 This study provides a basic framework for constructing a PFINN and applying the proposed regular-  
464 ization method to facilitate the implementation of contextual relaxing of the learned function. The  
465 key takeaway is that integrating these techniques can offer a probability theory based perspective  
466 on model architecture construction which allows assembling of relevant regularisation mechanisms,  
467 paving the way for more flexible applications on unseen data. Possible future work could be to better  
468 formalize PFINNs and develop further machinery to provide new insights into these models.

## 472 8 REPRODUCIBILITY STATEMENT

473 We are committed to ensuring the reproducibility of our research findings. Our models have been  
474 implemented with generic hyperparameter settings as discussed in the start of section 6, avoiding  
475 any specific tuning to present an honest view of our methodology. All proofs are listed with detailed  
476 steps to help readers. Furthermore additional proofs and figures are attached in the appendix to aid  
477 understanding of the concepts presented. We encourage the community to engage with our work. If  
478 any discrepancies or concerns are identified, we welcome dialogue to address them in the spirit of  
479 scientific inquiry and collaboration.

## 482 REFERENCES

483 Olivier Bournez. [lix.polytechnique.fr. https://www.lix.polytechnique.fr/  
484 ~bournez/load/MPRI/Cours-2024-MPRI-partie-I-goodMPRI.pdf](https://www.lix.polytechnique.fr/~bournez/load/MPRI/Cours-2024-MPRI-partie-I-goodMPRI.pdf), 2024. [Ac-  
485 cessed 02-10-2024].

486 G Cantor. Uber eine elementare frage der mannigfaltigkeitslehre, jahresbericht der deutschen  
487 mathematiker-vereinigung 1: 75–78, 1932.

488  
489 Rogero Cotes. Logometria auctore rogero cotes, trin. coll. cantab. soc. astr. and ph. exp. professore  
490 plumiano, and r. s. s. *Philosophical Transactions (1683-1775)*, 29:5–45, 1714. ISSN 02607085.  
491 URL <http://www.jstor.org/stable/103030>.

492 Salvatore Cuomo, Vincenzo Schiano Di Cola, Fabio Giampaolo, Gianluigi Rozza, Maziar Raissi,  
493 and Francesco Piccialli. Scientific machine learning through physics–informed neural networks:  
494 Where we are and what’s next. *Journal of Scientific Computing*, 92(3):88, 2022.

495  
496 George Cybenko. Approximation by superpositions of a sigmoidal function. *Mathematics of control,*  
497 *signals and systems*, 2(4):303–314, 1989.

498 Richard Dedekind. *Essays on the Theory of Numbers*. Courier Corporation, 2012.

499  
500 Leonhard Euler. Introductio in analysin infinitorum, volume 2. *Lausanæ: Apud Marcum-Mich aelem*  
501 *Bousquet and Socios*, 1748.

502 Kurt Hornik, Maxwell Stinchcombe, and Halbert White. Multilayer feedforward networks are uni-  
503 versal approximators. *Neural networks*, 2(5):359–366, 1989.

504  
505 Aleksandr Lyapunov. Sur une proposition de la théorie des probabilités. *Bulletin de l’Académie*  
506 *Impeériale des Sciences*, 13(4):359–386, 1900.

507 Klaus Weihrauch. *Computable analysis: an introduction*. Springer Science & Business Media,  
508 2012.

## 509 510 511 A APPENDIX

### 512 513 A.1 SATISFACTION OF LYAPUNOV CONDITION

514  
515 **Definition 5.** Suppose there exists a sequence of independent random variables  $\{Y_1, Y_2, \dots, Y_n\}$ , with  
516 finite mean and variance, we can expect that the growth of the moments are limited by the Lyapunov  
517 Condition.

$$518 \lim_{n \rightarrow \infty} \frac{1}{s_n^{2+\delta}} \sum_{i=1}^n \mathbb{E} \left[ |Y_i - \mu_i|^{2+\delta} \right] = 0 \quad (20)$$

519  
520 **Definition 6.** For some sequence of independent Bernoulli random variables  $\{X_1, X_2, \dots, X_n\}$ , such  
521 that

$$522 X_i \sim \text{Bernoulli}(p_i) \quad (21)$$

523  
524  $\mathbb{P}(X_i = 1) = p, 0 \leq p \leq 1, E(X_i) = p, \text{Var}(X_i) = p_i(1 - p_i)$

525 **Proposition 6.** Under most conditions, the Lyapunov CLT condition holds for Bernoulli Random  
526 Variables.

527  
528 *Proof.*

$$529 \lim_{n \rightarrow \infty} \frac{1}{s_n^{2+\delta}} \sum_{i=1}^n E[ (|X_i - \mu_i|^{2+\delta}) ] \quad (22)$$

530  
531 Without Loss of Generality, let  $\delta = 2$ :

$$532 \lim_{n \rightarrow \infty} \frac{1}{s_n^4} \sum_{i=1}^n E[ (|X_i - \mu_i|^4) ] \quad (23)$$

533  
534 By replacing  $\mu_i$  with  $E(X_i)$ :

$$535 = \lim_{n \rightarrow \infty} \frac{1}{s_n^4} \sum_{i=1}^n E[ (|X_i - E(X_i)|^4) ] \quad (24)$$

By the Law of the Unconscious Statistician (LOTUS):

$$= \lim_{n \rightarrow \infty} \frac{1}{s_n^4} \sum_{i=1}^n \sum_{X_i=0}^{X_i=1} [(|X_i - E(X_i)|^4)] \quad (25)$$

By definition of Bernoulli distribution:

$$= \lim_{n \rightarrow \infty} \frac{1}{s_n^4} \sum_{i=1}^n (0 - p_i)^4(1 - p_i) + (1 - p_i)^4(p_i) \quad (26)$$

With reference to equation 2:

$$= \lim_{n \rightarrow \infty} \frac{1}{(\sum_{i=1}^n \sigma^2)^2} \sum_{i=1}^n p_i^4(1 - p_i) + (1 - p_i)^4(p_i) \quad (27)$$

By the Variance described for Bernoulli Random Variables,  $\sigma^2 = p_i(1 - p_i)$ :

$$= \lim_{n \rightarrow \infty} \frac{1}{(\sum_{i=1}^n (p_i(1 - p_i)))^2} \sum_{i=1}^n p_i^4(1 - p_i) + (1 - p_i)^4(p_i) \quad (28)$$

Since parameter  $0 \leq p \leq 1$ , we can claim  $p_i^4 \leq p_i$  and  $(1 - p_i)^4 \leq (1 - p_i)$ :

$$\leq \lim_{n \rightarrow \infty} \frac{1}{(\sum_{i=1}^n (p_i(1 - p_i)))^2} \sum_{i=1}^n p_i(1 - p_i) + (1 - p_i)(p_i) \quad (29)$$

$$= \lim_{n \rightarrow \infty} \frac{1}{(\sum_{i=1}^n (p_i(1 - p_i)))^2} \sum_{i=1}^n 2p_i(1 - p_i) \quad (30)$$

By Linearity of the Sum,

$$= \lim_{n \rightarrow \infty} \frac{2 \sum_{i=1}^n (p_i(1 - p_i))}{(\sum_{i=1}^n (p_i(1 - p_i)))^2} \quad (31)$$

$$= \lim_{n \rightarrow \infty} \frac{2}{\sum_{i=1}^n (p_i(1 - p_i))} \quad (32)$$

As  $n \rightarrow \infty$ ,

$$\because \sum_{i=1}^n (p_i(1 - p_i)) \rightarrow \infty \quad (33)$$

We have

$$\lim_{n \rightarrow \infty} \frac{2}{\sum_{i=1}^n (p_i(1 - p_i))} = 0 \quad (34)$$

as desired.  $\square$

## B DISTANCE MEASURES

In this section, we extend the concept of  $L_p$  norms to measure the differences between the distributions  $\phi_D$  and  $\phi_{\mathcal{N}(0,1)}$ . We define the distance function  $d(\phi_D, \phi_{\mathcal{N}(0,1)})$  by calculating the pointwise differences between the two distributions and applying the  $L_p$  norms.

We start with the general definition of the  $L_p$  norm for a vector  $\mathbf{x} = (x_1, x_2, \dots, x_n)$ :

$$\|\mathbf{x}\|_p = \left( \sum |x_k|^p \right)^{\frac{1}{p}}, \quad p \geq 1. \quad (35)$$

Extending the definition of the standard  $L_1$  norm, which provides a measure based on the absolute differences:

$$\psi_1 = d_1(\phi_D, \phi_{\mathcal{N}(0,1)}) = \|\phi_D - \phi_{\mathcal{N}(0,1)}\|_1 = \sum_{k=-\infty}^{\infty} |\phi_D(u_k) - \phi_{\mathcal{N}(0,1)}(u_k)|. \quad (36)$$

Next, we extend the  $L_2$  norm, which measures the Euclidean distance between the pointwise differences:

$$\psi_2 = d_2(\phi_D, \phi_{\mathcal{N}(0,1)}) = \|\phi_D - \phi_{\mathcal{N}(0,1)}\|_2 = \sqrt{\sum_{k=-\infty}^{\infty} |\phi_D(u_k) - \phi_{\mathcal{N}(0,1)}(u_k)|^2}. \quad (37)$$

Finally, we can consider the extensions of the  $L_\infty$  norm, which measures the maximum pointwise difference:

$$\psi_\infty = d_\infty(\phi_D, \phi_{\mathcal{N}(0,1)}) = \|\phi_D - \phi_{\mathcal{N}(0,1)}\|_\infty = \sup_k |\phi_D(u_k) - \phi_{\mathcal{N}(0,1)}(u_k)|. \quad (38)$$

The selection of these three distance measures is intentional, prioritizing simplicity and ease of replication. While geometric distance measures could potentially yield greater performance, we have chosen to focus on these straightforward metrics to provide a gentle introduction to the topic and methodology discussed in this paper.

## C FIGURES

### C.1 CHARACTERISTIC FUNCTION OF NORMAL AND BERNOULLI DISTRIBUTION

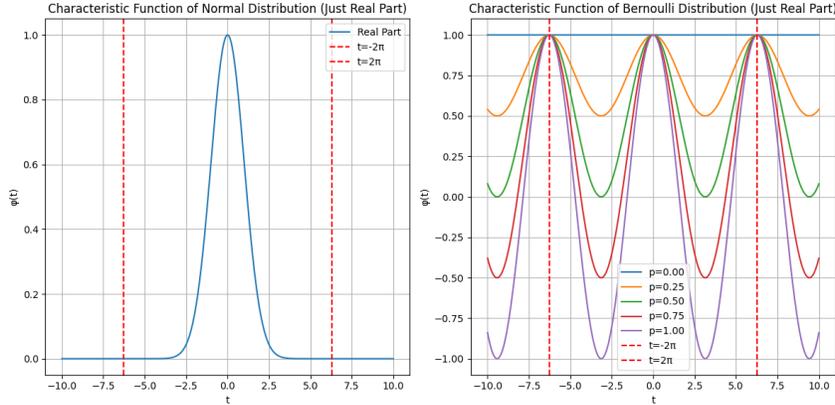


Figure 1: Plot of Normal and Bernoulli Characteristic Functions (Only Real Part)

The figure C.1 shows the plot of real part of the Normal and Bernoulli Distribution. We thought this would be apt to add as this is to give a visual intuition for the reader for how these functions look when graphed as there is not much literature regarding visualising them.

### C.2 IMAGINARY PART INCLUSIVE CHARACTERISTIC FUNCTION OF NORMAL AND BERNOULLI DISTRIBUTION

The figure C.2 shows the plot of the Normal and Bernoulli Distribution inclusive of the imaginary part. It is interesting to note the imaginary part is on the zero line for the Normal. As for the Bernoulli we can see a "phase" difference between the Imaginary and the Real Part.

If one would like to explore why, they can derive insight using the following as a starting point:

648  
649  
650  
651  
652  
653  
654  
655  
656  
657  
658  
659  
660  
661  
662  
663  
664  
665  
666  
667  
668  
669  
670  
671  
672  
673  
674  
675  
676  
677  
678  
679  
680  
681  
682  
683  
684  
685  
686  
687  
688  
689  
690  
691  
692  
693  
694  
695  
696  
697  
698  
699  
700  
701

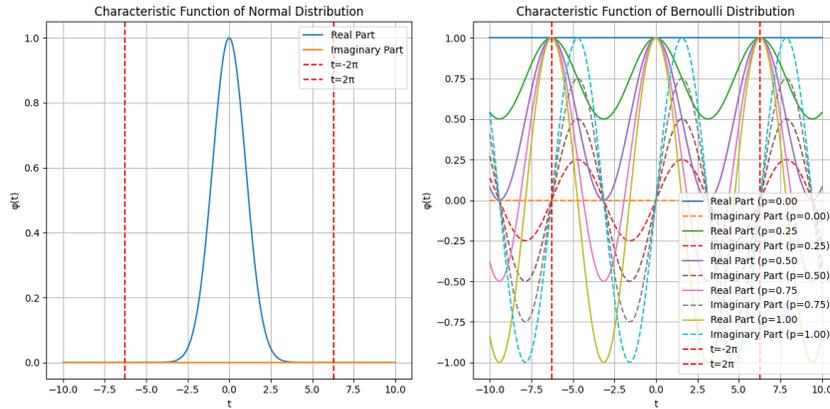


Figure 2: Extended Plot of Normal and Bernoulli Characteristic Function (Includes Imaginary Part)

**Definition 7.** Euler’s formula (Euler, 1748) (Cotes, 1714) states that for any real number  $x$ :

$$e^{\vartheta x} = \cos(x) + \vartheta \sin(x) \tag{39}$$

This formula can be used to express complex exponentials in terms of trigonometric functions.

**Definition 8.** Using equation 5 and definition 7, the characteristic function of a random variable  $X$  is defined as:

$$\phi(u) = \mathbb{E}[e^{\vartheta uX}] = \mathbb{E}[\cos(uX)] + \vartheta \mathbb{E}[\sin(uX)] \tag{40}$$

where the real and imaginary parts of the characteristic function are:

$$\text{Re}(\phi(tu)) = \mathbb{E}[\cos(uX)] \tag{41}$$

$$\text{Im}(\phi(u)) = \mathbb{E}[\sin(uX)] \tag{42}$$

### C.3 ZOOMED OUT VIEW TO OBSERVE PERIODICITY

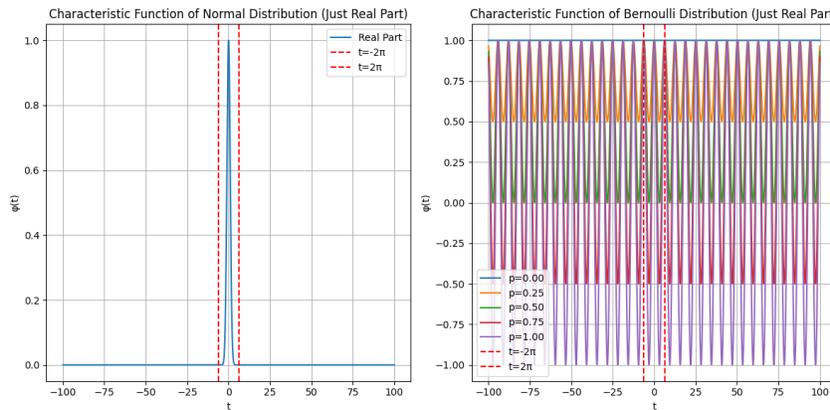
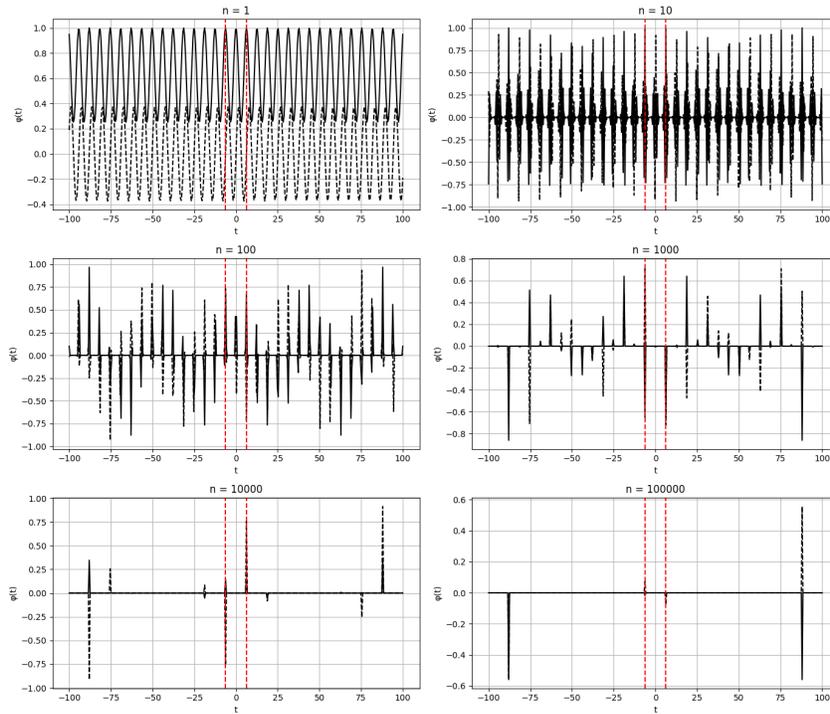


Figure 3: Plot of Normal and Bernoulli

The figure C.3 shows that the Normal Characteristic Function does not seem to be periodic unlike the Bernoulli Characteristic Function which seems to have a defined  $\pi$ -periodic structure. It also

702 shows that the Normal Characteristic Function is concentrated within the  $-2\pi$  to  $2\pi$  region. (Which  
703 motivated our choice in the numerics section 5).  
704

#### 705 C.4 BEHAVIOUR OF THE CHARACTERISTIC FUNCTION WHEN JUST ADDING BERNOULLI 706 VARIABLES TOGETHER MINDLESSLY 707



708  
709  
710  
711  
712  
713  
714  
715  
716  
717  
718  
719  
720  
721  
722  
723  
724  
725  
726  
727  
728  
729  
730  
731  
732  
733  
734  
735  
736  
737  
738  
739  
740  
741  
742  
743  
744  
745  
746  
747  
748  
749  
750  
751  
752  
753  
754  
755

Figure 4: Numerical Simulation Plot of the Convergence

The figure C.4 is generated random generated  $p_i$  values for a  $\sum_{i=1}^N$  Bernoulli Distributions. It is interesting to note how just adding the Bernoulli's will result in it resulting in a convergence towards the zero line.

#### C.5 NUMERICAL SIMULATION OF CONVERGENCE DESCRIBED IN PROPOSITION 1

The figure C.5 is generated random generated  $p_i$  values for a linear combination  $N$  Bernoulli Distributions which are added according to the  $\aleph$  model described in definition 3.

756  
757  
758  
759  
760  
761  
762  
763  
764  
765  
766  
767  
768  
769  
770  
771  
772  
773  
774  
775  
776  
777  
778  
779  
780  
781  
782  
783  
784  
785  
786  
787  
788  
789  
790  
791  
792  
793  
794  
795  
796  
797  
798  
799  
800  
801  
802  
803  
804  
805  
806  
807  
808  
809

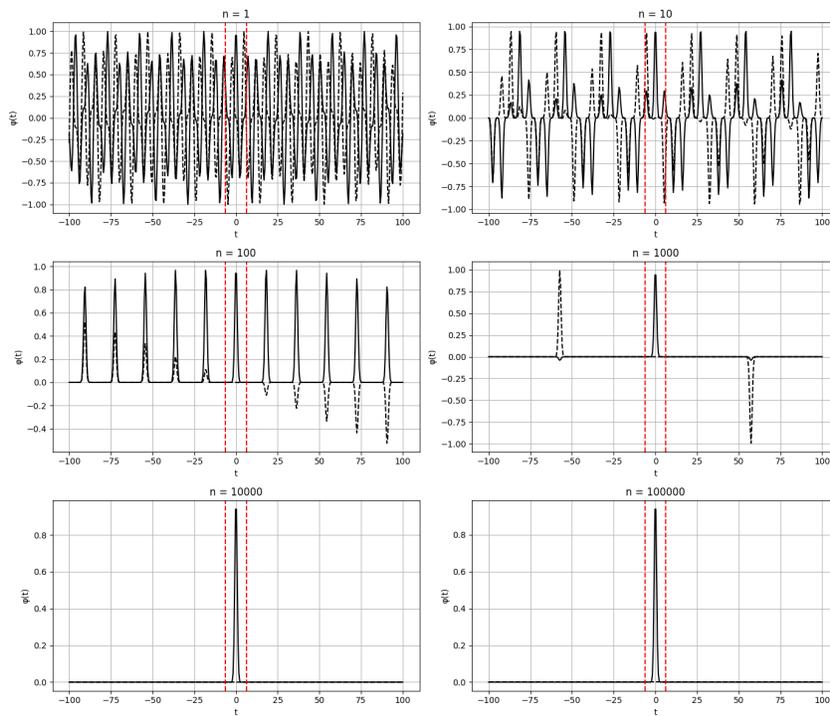


Figure 5: Numerical Simulation Plot of just adding Bernoullis