

1 **18.337 PROJECT REPORT: AUGMENTED NEURAL ODES**

2 KERRI LU

3 **Abstract.** Neural Ordinary Differential Equations are a class of deep neural network models that
4 learn supervised mappings from inputs to outputs by solving an ODE initial value problem. However,
5 NODEs have limited expressivity, as they preserve topology of the input space. This motivates
6 the introduction of Augmented Neural ODEs (ANODEs), which augment the space on which the
7 ODE is solved to a higher dimension, increasing model expressivity and allowing for simpler ODE
8 flow trajectories. In this paper, we review and implement ANODE algorithms for regression and
9 classification, experimenting with several variants including zero-padding augmentation, input-layer
10 augmentation, temporally regularized augmentation, and second-order neural ODEs. Running the
11 algorithms on toy time series and classification datasets as well as image classification tasks, we show
12 that augmentation generally improves model performance and enables faster convergence compared
13 to regular NODEs. All code is available at <https://github.com/kerrilu/18337-anode-project>.

14 **1. Introduction and Background.** Neural ODEs (NODEs), introduced in [1],
15 are a class of deep neural network models that learn a mapping from input $x \in \mathbb{R}^d$ to
16 an output hidden state $\phi(x) = h(T) \in \mathbb{R}^d$ at end time T by solving the initial value
17 problem

18
$$\frac{dh(t)}{dt} = f(h(t), t)$$
 with initial condition $h(0) = x$
19

20 where $f : \mathbb{R}^d \rightarrow \mathbb{R}^d$ is a differentiable, dimension-preserving neural net function. These
21 models can be viewed as a continuous version of residual networks, with parameters
22 encoded in f . During training with backpropagation, the weights of f are adjusted to
23 fit ground-truth output labels, by optimizing some loss function between the learned
24 features $\phi(x)$ and the true labels. While $\phi(x)$ is dimension preserving, NODEs can be
25 modified to perform supervised regression or classification by adding a linear neural
26 network layer $l : \mathbb{R}^d \rightarrow \mathbb{R}$ after the ODE solver output, as shown in Figure 1.

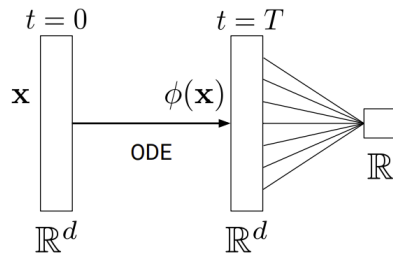


FIG. 1. Neural ODE architecture with final linear layer mapping to scalar output, from [2].

27 **1.1. Limitations of NODEs.** NODEs have several practical advantages, such
28 as their use in irregularly sampled time series prediction and density estimation [3], as
29 well as their constant memory cost. However, NODEs also have limited expressivity,
30 as they preserve the topology of the input space. This limitation often makes learning
31 NODE approximations for functions computationally costly [2]. In particular, we
32 define the *flow* of a NODE as the trajectory of the hidden state $h(t)$ over time from
33 $x = h(0)$ to $\phi(x) = h(T)$. Intuitively, since the trajectories corresponding to different
34 initial condition inputs x cannot intersect, there are many classes of functions which
35 cannot be represented by ODE flows.

36 An illustrative example is the class of one-dimensional functions g satisfying
 37 $g(1) = -1$ and $g(-1) = 1$. As shown in Figure 2, any pair of continuous trajec-
 38 tories corresponding to these mappings must intersect each other, and thus cannot be
 39 learned by a NODE.

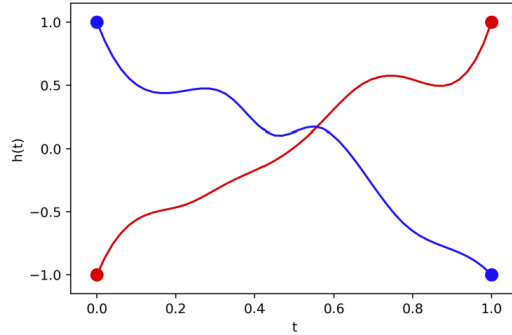


FIG. 2. NODEs cannot represent trajectories for both $g(1) = -1$ and $g(-1) = 1$, from [2].

40 Another canonical example is that of nested spheres. It is shown in [2] that
 41 NODEs cannot represent functions $g : \mathbb{R}^d \rightarrow \mathbb{R}$ where $g(x) = 1$ for $\|x\| \leq r_1$ and
 42 $g(x) = -1$ for $r_2 \leq \|x\| \leq r_3$ (for any $0 < r_1 < r_2 < r_3$). As shown in Figure 3, in
 43 order for the NODE to linearly separate the two classes of points, the learned flows
 44 from the inner and outer spheres would intersect each other.

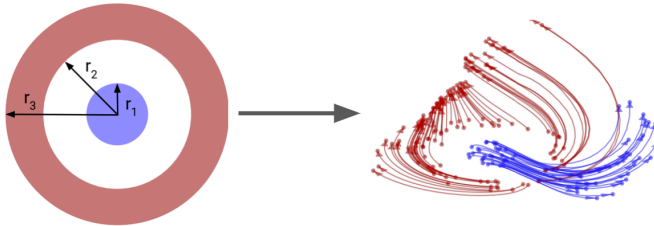


FIG. 3. NODEs cannot easily represent classification of nested spheres (left) due to intersecting flows (right), adapted from [2].

45 **1.2. Augmented Neural ODEs.** The limitations of NODEs motivate the in-
 46 troduction of Augmented Neural ODEs (ANODEs), which augment the input space
 47 from \mathbb{R}^d to \mathbb{R}^{d+p} so that ODE flows are “lifted” to higher dimensions so that they
 48 don’t intersect. Formally, ANODEs pad input data points x with p zeros, introduce
 49 the augmented points $a(t) \in \mathbb{R}^p$, and solve the following modified initial value problem
 50 for $h(T)$.

$$51 \quad \frac{d}{dt} \begin{bmatrix} h(t) \\ a(t) \end{bmatrix} = f \left(\begin{bmatrix} h(t) \\ a(t) \end{bmatrix}, t \right) \text{ with initial condition } h(0) = \begin{bmatrix} x \\ \mathbf{0} \end{bmatrix}.$$

53 ANODEs can thus learn a broader class of functions and tend to result in smoother
 54 and simpler flows. They are “empirically more stable, generalize better and have a
 55 lower computational cost” than NODEs, and fewer function evaluations are required
 56 to solve the augmented ODE [2].

57 **1.3. Structure of Paper.** In this paper, we review and implement variants of
 58 ANODE algorithms, and experimentally evaluate their performance. The rest of the
 59 paper is structured as follows. In Section 2, we give a brief overview of related work on
 60 NODEs and ANODEs. In Section 3, we discuss algorithmic methods and potential
 61 ANODE performance improvements such as input-layer augmentation, temporally
 62 regularized augmentation, and second-order neural ODEs. In Section 4, we discuss
 63 our experiments and evaluate results of using NODEs and ANODEs to learn toy
 64 time series and classification datasets, as well as image classification on MNIST. In
 65 Section 5, we summarize the conclusions of our work, and in Section 6, we discuss
 66 potential future research directions. All code is available at [https://github.com/
 67 kerrilu/18337-anode-project](https://github.com/kerrilu/18337-anode-project).

68 **2. Related Work.** In addition to the standard *zero-padding augmentation* pro-
 69 posed by [2], there are several potential performance optimizations for ANODEs in
 70 the literature. One alternative is *input-layer augmentation*, where the augmented ini-
 71 tial condition is $h(0) = g(x)$ for some input network $g : \mathbb{R}^d \rightarrow \mathbb{R}^{d+p}$, which allows the
 72 model more freedom (improving model capacity) and empirically decreases the num-
 73 ber of function evaluations [5]. To further accelerate training time, [3] demonstrates
 74 regularization by randomly sampling the end time T of the ODE. Additionally, using
 75 second-order (or higher) ODEs can improve parameter efficiency; [6] considers second-
 76 order NODEs as a special case of ANODEs with constraints on the neural network
 77 structure, and shows that even first-order ANODEs can use augmented dimensions
 78 to learn higher-order dynamics.

79 The Julia SciML ecosystem includes several packages with high-performance im-
 80 plementations of differential equation solvers, such as DifferentialEquations.jl [8].
 81 In particular, the DiffEqFlux.jl library [7] combines neural networks and differen-
 82 tial equations, providing a framework for training NODEs as well as related models
 83 such as Neural Stochastic Differential Equations. Though less directly relevant, the
 84 NeuralPDE.jl library [9] implements physics-informed neural networks to efficiently
 85 approximate high-dimensional PDE solutions.

86 **3. Methods.** For our implementation of ANODEs, we utilize the DiffEqFlux.jl
 87 library to efficiently construct neural networks (using Flux) and train neural ODEs.
 88 The implementation of zero-padding augmentation as described in Section 1.2 is
 89 straightforward. Using the notation from that section, we simply pad inputs with
 90 p zeros, solve the resulting neural ODE by optimizing a loss function of our choice,
 91 and remove the augmented points $a(t) \in \mathbb{R}^p$ from the output. We can visualize the
 92 resulting ODE flow trajectories, as well as test accuracies and losses over iterations, to
 93 evaluate model performance. We experiment with several values for the augmentation
 94 dimension p to improve performance.

95 We also implement a few extensions and performance improvements for ANODEs,
 96 described in more detail below.

97 **3.1. Input-Layer Augmentation.** In Input-Layer Augmented NODEs (IL-
 98 ANODEs) [5], the initial condition $h(0)$ is computed as the output of an input neural
 99 network $g(x) : \mathbb{R}^d \rightarrow \mathbb{R}^{d+p}$ mapping inputs to a higher dimensional space. This can be
 100 seen as a generalization of ANODEs (which correspond to the case $g(x) = (x, \mathbf{0})$) that
 101 allows for richer representations and greater expressivity compared to the simple zero-
 102 padding operation. Empirically, this has been shown to allow for faster convergence
 103 and fewer function evaluations. In our implementation, we add a linear layer before
 104 the Neural ODE layer (using Flux.Chain) and train the entire network together. Note

105 that the added input layer slightly increases the total parameter complexity of the
106 system.

107 **3.2. Temporal Regularization.** Another potential performance improvement
108 is the temporal regularization method proposed by [3], in which the end time param-
109 eter T of the NODE or ANODE integration limits is randomly perturbed. This added
110 stochasticity during training is empirically shown to simplify model dynamics, reduce
111 computational cost, and lead to faster convergence during training. Formally, in our
112 implementation, the hidden state at time t becomes

$$113 \quad h(t) = h(0) + \int_0^T f(h(t), t) dt = \text{ODESolve}(h(0), f, 0, T)$$

114
115 where the end time T is sampled uniformly at random from the interval $(t-b, t+b)$ for
116 some scalar parameter b . (Without regularization, we would simply have $T = t$.) This
117 method effectively enforces convergence at time $t-b$ rather than time t , as illustrated
118 in Figure 4.

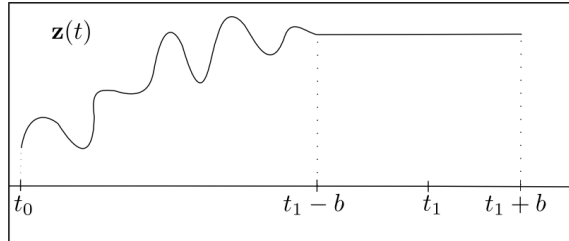


FIG. 4. Behavior of NODEs after temporal regularization with parameter b , from [3].

119 **3.3. Higher-Order NODEs and ANODEs.** NODEs and ANODEs can be
120 generalized to higher dimensions, as discussed in [5]. An n th order NODE can be
121 described by concatenating vectors $h_i(t) \in \mathbb{R}^{d/n}$ so that the hidden state is $h(t) =$
122 $[h_1(t), h_2(t), \dots, h_n(t)]$. The system of coupled first-order equations is

$$123 \quad \frac{d}{dt} h_i(t) = h_{i+1}(t) \text{ for } i < n,$$

$$124 \quad \text{and } \frac{d}{dt} h_n(t) = f(h(t), t)$$

125
126 where $f : \mathbb{R}^d \rightarrow \mathbb{R}^{d/n}$ is a neural network whose output dimension is n times smaller
127 than its input dimension. This results in increased parameter efficiency compared to
128 first-order NODEs and ANODEs where f is dimension-preserving.

129 Note that equivalently, we can write and solve the n th order NODE equation as

$$130 \quad \frac{d^n}{dt^n} h_1(t) = f \left(h_1(t), \frac{d}{dt} h_1(t), \dots, \frac{d^{n-1}}{dt^{n-1}} h_1(t), t \right).$$

131
132 For image classification, we implemented a second-order Neural ODE, using Dif-
133 fEqFlux's built-in function for solving second-order ODE problems. As observed in
134 [6], second-order NODEs can be viewed as a constrained type of ANODE where the
135 augmenting vector $a(t)$ is the derivative of $h(t)$, which also constrains the structure of
136 f . Unlike first-order NODEs, second-order and higher NODEs (even without augmen-
137 tation) do not necessarily preserve topology of the input space and can thus represent
138 a wider range of functions.

139 **4. Experiments and Results.** In this section, we discuss experiments and re-
 140 sults from applying ANODEs to a few different settings: toy time-series prediction
 141 with sinusoidal functions, toy classification of nested spheres, and MNIST image clas-
 142 sification. For the first two toy examples, we use simple zero-padding ANODEs and
 143 experiment with augmentation dimensions. For the more complex image classification
 144 task, we also experiment with the ANODE variants and performance improvements
 145 described in the previous section.

146 **4.1. Time-Series Prediction with Sinusoidal Functions.** For this toy ex-
 147 ample, we generate 100 datapoints from each of the functions $y(x) = \sin(x)$ and
 148 $y(x) = \sin(x) + x$ and inject a small amount of Gaussian noise. Then, our goal is to
 149 train ANODEs to learn the functions. Note that these two functions cannot be rep-
 150 resented by regular NODEs due to intersecting flow trajectories (e.g. any continuous
 151 flows for $\sin(5\pi/6) = 0.5$ and $\sin(\pi/2) = 1$ would have to intersect).

152 We compare zero-padding ANODE performance for different numbers of aug-
 153 mentation dimensions p , and also compare against the baseline performance of the
 154 original NODE algorithm. We use a simple neural net with input and output size
 155 $2 + p$ (since each datapoint $(x, y(x))$ is two-dimensional), and a single hidden layer of
 156 size 20. We use $(0, 0)$ as the initial condition. The loss function is the sum of squared
 157 differences between predicted coordinates and ground-truth data. The ANODE is
 158 trained for 1000 iterations using Adam optimizer with learning rate 0.05, then the
 159 BFGS optimization algorithm is used for fine-tuning.

160 The final functions learned by the NODE and ANODEs are shown in Figure 5.
 161 In both cases, the NODE prediction (in red) differs significantly from the ground-
 162 truth data (in blue). Increasing the number of augmentation dimensions significantly
 163 improves the predictions. The ANODE padded with 5 zeros predicts $y(x) = \sin(x)$
 164 fairly well. More augmentation is needed for the more complex $y(x) = \sin(x) + x$
 165 function: an ANODE padded with 10 zeros performs well.

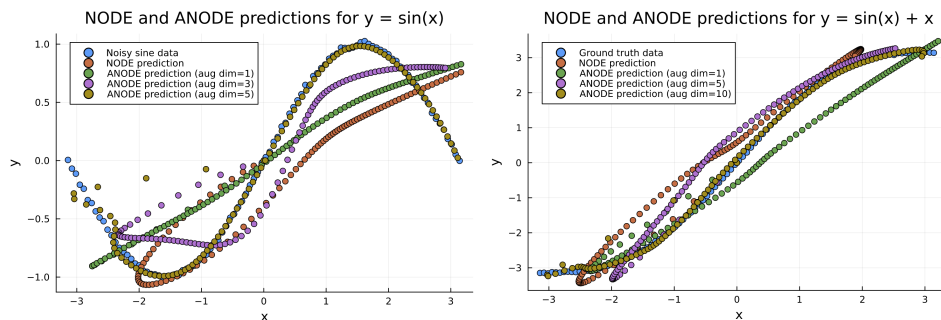


FIG. 5. *NODE and ANODE predicted coordinates for our time-series functions.*

166 We can also quantify convergence by plotting loss values over training iterations.
 167 For example, Figure 6 shows the models' loss for predicting $y(x) = \sin(x) + x$. In
 168 general, we see that increasing the number of augmentation dimensions results in
 169 lower loss and better convergence in fewer iterations (whereas the regular NODE loss
 170 and ANODE loss with augmentation dimension $p = 1$ do not converge at all). This
 171 is expected, as lifting the ODE to a higher dimensional space improves the model's
 172 expressivity and allows it to learn simpler flows more quickly.

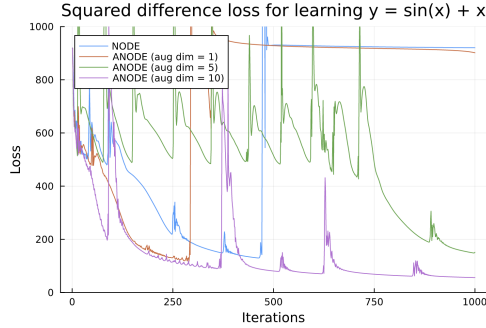


FIG. 6. *NODE* and *ANODE* training loss for predicting $y(x) = \sin(x) + x$. Note that the loss converges in fewest iterations when we use the highest number of augmentation dimensions ($p = 10$).

173 **4.2. Classification of Nested Spheres.** In this section, we train ANODEs to
 174 do classification on the canonical example where regular NODEs fail: nested concentric
 175 circles. We generate a dataset of 4000 datapoints, in which points in the inner
 176 circle (of radius 1) are labeled "1" and points in the outer annulus (between the circles
 177 of radius 2 and radius 3) are labeled "-1." The dataset is visualized in Figure 7.

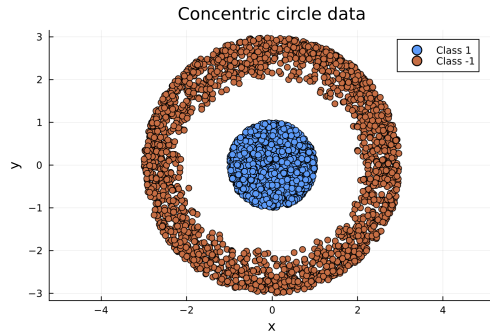


FIG. 7. *Concentric circle data with class labels.*

178 We use a neural network with two hidden layers of dimension 64, and Adam
 179 optimizer with learning rate 0.05. A final linear layer of the network outputs a real
 180 scalar value (so technically it computes a regression). We use mean square error loss
 181 between the ground-truth and predicted labels. We experiment with 2, 5, and 10
 182 zero-padding augmentation dimensions and train all models for 20 epochs.

183 The resulting contour maps for predicted regression values at each coordinate are
 184 shown in Figure 8. The NODE model does not converge, while the performance of
 185 the ANODE models improves as the number of augmentation dimensions increases
 186 (for similar reasons as explained in the previous subsection). With 10 augmentation
 187 dimensions, the contour map is close to being a series of concentric circles of decreasing
 188 label value as radius increases, which is the desired result.

189 We plot the mean square error loss for each of the models over the training
 190 iterations in Figure 9. The loss does not converge for the NODE. For the ANODEs,
 191 increasing augmentation dimension leads to faster convergence and lower loss, as
 192 before.

193 Finally, to better understand ANODE behavior, we plot the learned features (the

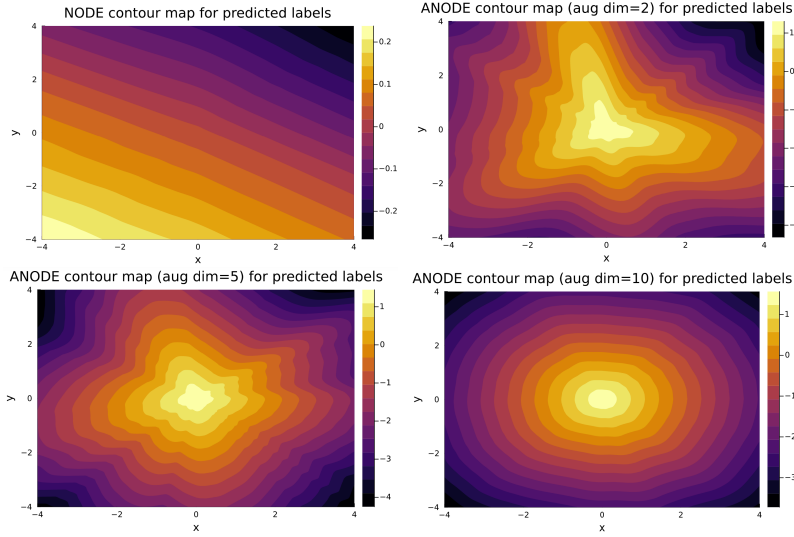


FIG. 8. Contour maps for predicted labels at each coordinate, for trained NODE and ANODE models.



FIG. 9. Model loss during training for learning concentric circles.

194 final flow trajectory values $\phi(x) = h(T)$ for augmentation dimension 10. In Figure
 195 10, we choose a few of the augmented dimensions and plot the final flow trajectory
 196 locations on the z-axis in 3D. We also plot the learned features for the regular NODE
 197 in 2D for comparison. In the ANODE model, we see that the higher dimensionality
 198 allows the inner circle class to “lift out” and become linearly separable from the outer
 199 annulus, without having intersecting flows. By contrast, the NODE learned features
 200 are only a slightly distorted version of the original data, and the red and blue points
 201 do not become linearly separable.

202 As expected, the augmentation allows the model to learn a class of regression
 203 functions that would not be learnable by the regular NODE.

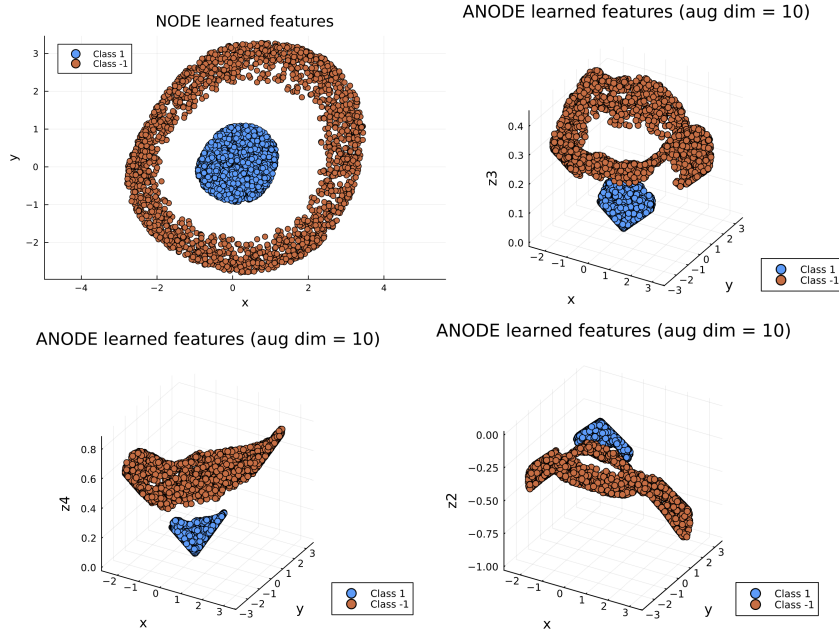


FIG. 10. Flow trajectory features at final time T for NODE and three of the augmented ANODE dimensions.

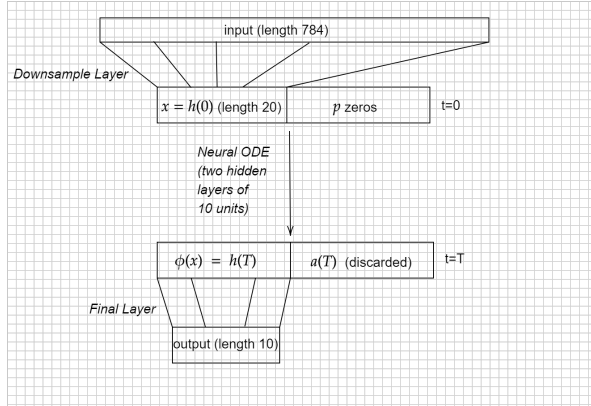
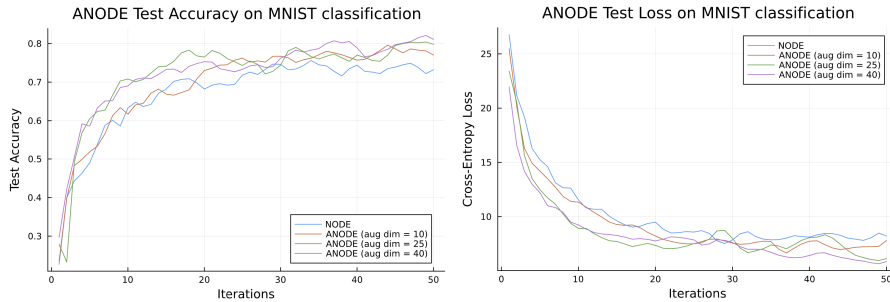
204 **4.3. MNIST Image Classification.** Finally, we train NODE and ANODE
 205 models to perform classification on the MNIST image dataset, which consists of 28 by
 206 28 pixel images of handwritten digits paired with one-hot ground-truth labels from 0
 207 to 9.

208 The first layer of the neural network reshapes the images to vectors of length 784.
 209 The next layer is a dense downsampling layer with 20 output nodes which serve as
 210 the input to the Neural ODE. The NODE has two hidden layers of 10 nodes each,
 211 and 20 output nodes. Finally, there is a dense layer that maps the NODE output to a
 212 vector of length 10, corresponding to the one-hot encoding of the ten digits. We use
 213 logit cross entropy loss, Adam optimizer with learning rate 0.05, and batch size 100.
 214 We train for 50 iterations.

215 For zero-padding augmentation, the only modification to the network architecture
 216 is that the input to the NODE is padded with p zeros, for a total length of $20 + p$ (and
 217 then the NODE output must also have length $20 + p$, but it is truncated to the first
 218 20 nodes before being fed into the output layer). The overall network architecture is
 219 shown in Figure 11.

220 We experiment with values $p = 10, 25, 40$. In Figure 12, we plot the NODE
 221 and ANODE model classification accuracy and loss on the test set over the training
 222 iterations. The ANODE models converge to somewhat higher accuracy and lower
 223 loss in fewer iterations than the NODE. The final NODE accuracy of 73% is a few
 224 percentage points lower than the final ANODE accuracy of 77% (with augmentation
 225 dimension $p = 10$). Further increasing the augmentation dimension also slightly
 226 improves final accuracy, as seen in the graph.

227 We also experiment with the ANODE variants described in Section 3.

FIG. 11. *Neural network and ANODE architecture for MNIST classification.*FIG. 12. *Zero-padding ANODE test accuracy and loss over training iterations.*

228 **4.3.1. Input-Layer Augmentation.** To implement input-layer augmentation
 229 (IL-ANODEs), we change the dense downsampling layer to have $20 + p$ output nodes.
 230 That is, the network directly learns a mapping from the input vector of length 784
 231 to a vector of length $20 + p$ which will serve as the input to the NODE (rather than
 232 mapping to a vector of length 20 and then padding with zeros).

233 As before, we experiment with different augmentation dimensions $p = 10, 25, 40$
 234 and plot the resulting test accuracies and losses in Figure 13. Even with the same
 235 number of augmentation dimensions $p = 10$, the IL-ANODE has somewhat faster conver-
 236 gence and has a slightly higher final accuracy of 81% compared to the zero-padding
 237 ANODE. This is consistent with the hypothesis that learning an input network in-
 238 creases the freedom and capacity of the model to represent more complex mappings.
 239 However, further increasing the IL-ANODE augmentation dimension above $p = 10$
 240 does not seem to result in significantly improved performance.

241 **4.3.2. Temporal Regularization.** The original NODE is solved from time
 242 span $t = 0$ to $t = 1$. We implement temporal regularization by randomly sampling the
 243 end time from the interval $(1 - b, 1 + b)$ where b is the regularization parameter defined
 244 in Section 3.2. We keep the zero-padding augmentation dimension constant at $p = 10$
 245 and experiment with several values of $b = 0.2, 0.3, 0.6$. Plotting the training accuracies
 246 and losses in Figure 14, we see that the temporally regularized ANODEs are slower
 247 to converge than the regular ANODE. This is expected, as regularization works to
 248 prevent overfitting during training (reducing the variance and generally increasing the

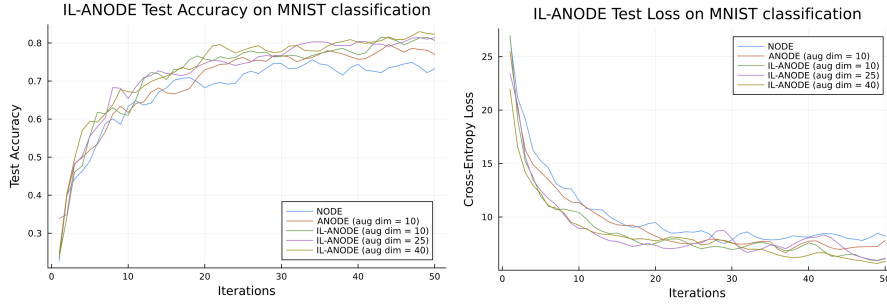


FIG. 13. *Input-layer augmentation ANODE test accuracy and loss over training iterations.*

249 bias), which results in lower training accuracy.

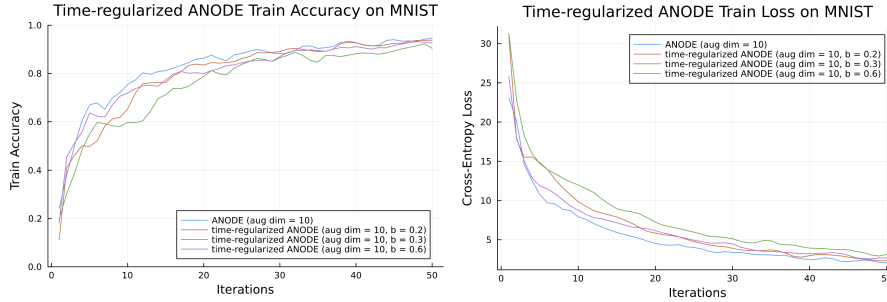


FIG. 14. *Time-regularized augmentation ANODE train accuracy and loss over training iterations.*

250 In Figure 15, the test accuracies and losses suggest that the temporal regulariza-
 251 tion slightly improves model performance in general (although there doesn't seem to
 252 be a clear relationship between the value of b and speed of convergence). The final
 253 test accuracies are around 81% which is slightly higher than the regular ANODE.

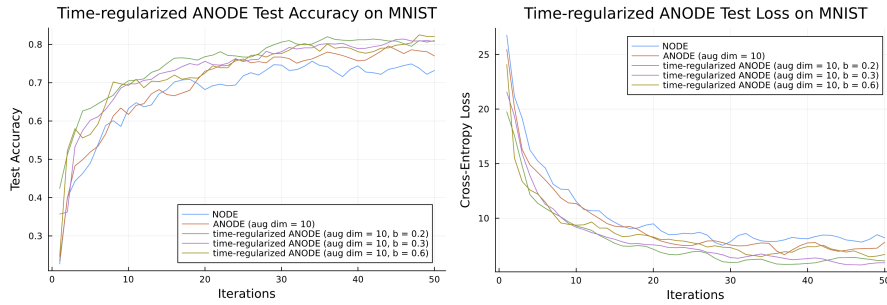


FIG. 15. *Time-regularized augmentation ANODE test accuracy and loss over training iterations.*

254 **4.3.3. Second-Order NODEs.** We implement a second-order NODE for MNIST
 255 classification. We modify the network architecture for better performance as follows.
 256 We change the output of the downsampling layer to have length 80, and the NODE
 257 part of the network to have input and output size 40 with a single hidden layer of 20
 258 nodes. (The downsampling layer output can be viewed as the concatenation of the
 259 initial condition $h(0)$ and the initial condition for its derivative $h'(0)$, each of which

Model	NODE	ANODE	IL-ANODE	Time-regularized ANODE (b=0.2)	Second-order NODE
Final accuracy	73.2%	77.0%	80.7%	81.0%	55.4%

TABLE 1

Final test accuracy after 50 iterations for NODEs and ANODEs with $p = 10$.

260 have length 40 in this case.) Empirically we find that using the second order equation
 261 $\frac{d^2}{dt^2}h(t) = f(h(t), t)$ seems to produce somewhat better results than feeding in both
 262 $h(t)$ and its derivative (i.e. $\frac{d^2}{dt^2}h(t) = f(h(t), h'(t), t)$). We use the second order ODE
 263 solver in DiffEqFlux to directly compute the second-order NODE output.

264 Unfortunately, as seen in Figure 16, the second-order NODE performance is much
 265 worse than the first-order NODE and ANODE models and does not converge. The
 266 final test accuracy is around 55%, and even training for more epochs seems to result
 267 in convergence at only around 60%. We believe our naive approach of directly using
 268 the second order ODE solver on the downsampled inputs is likely flawed; optimiza-
 269 tion of second order NODEs seems to require a different framework and a modified
 270 adjoint sensitivity method. More discussion on potential improvements can be found
 271 in Section 6 (Future Work).

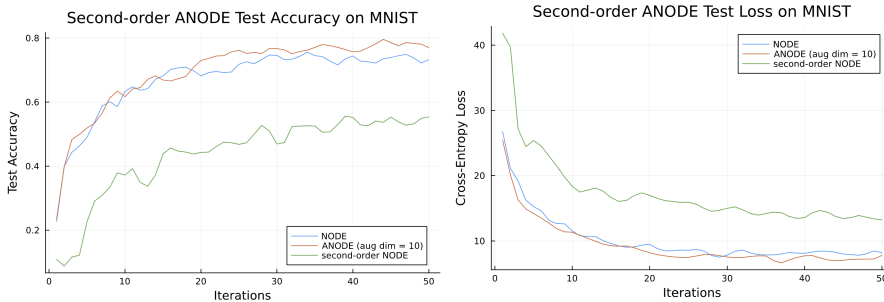


FIG. 16. Second-order NODE test accuracy and loss over training iterations.

272 The final test accuracies for NODEs and ANODEs with augmentation dimension
 273 $p = 10$ are summarized in Table 1. Overall, our results suggest that augmentation
 274 enables the NODE model to achieve somewhat higher accuracy, and our performance
 275 improvements seem effective (with the exception of the second-order system).

276 **5. Conclusion.** In this paper, we have reviewed Augmented Neural ODEs and
 277 implemented several variants with applications to toy examples and image classifica-
 278 tion. Our empirical results support the notion that augmenting the space on which
 279 we solve a neural ODE increases expressivity of the model and allow it to learn a
 280 broader class of functions, leading to faster convergence and smoother ODE flows.
 281 Furthermore, we showed that optimizations such as input-layer augmentation and
 282 time regularization can further improve model performance.

283 **6. Future Work.** For future work, we would like to improve on our experiments
 284 with second (or higher) order NODEs, as it would be interesting to better understand
 285 why our naive approach failed. One alternative approach would be to recursively use
 286 the first-order adjoint method, as higher order NODEs can be decomposed into a series
 287 of coupled first-order ODEs as shown in Section 3.3. However, there are also more
 288 sophisticated approaches for higher-order adjoint sensitivity methods. In particular,

289 [4] proposes an optimal control programming method for second-order optimization
290 of Neural ODEs, and shows that it empirically improves convergence time because
291 only a single backward pass is needed to find all derivatives.

292 It would also be of interest to try to apply Augmented Neural ODEs to a real-
293 world setting. Neural ODEs have shown to be especially useful for modeling irregularly
294 sampled time series [3], and it would be nice to explore potential applications to areas
295 such as climate and weather forecasting where data may be temporally sparse and
296 dynamical systems modeling is commonly used.

297

REFERENCES

- 298 [1] R. T. CHEN, Y. RUBANOVA, J. BETTENCOURT, AND D. K. DUVENAUD, *Neural ordinary differen-*
299 *tial equations*, Advances in neural information processing systems, 31 (2018).
- 300 [2] E. DUPONT, A. DOUCET, AND Y. W. TEH, *Augmented neural ODEs*, Advances in Neural Infor-
301 mation Processing Systems, 32 (2019).
- 302 [3] A. GHOSH, H. BEHL, E. DUPONT, P. TORR, AND V. NAMBOODIRI, *Steer: Simple temporal regu-*
303 *larization for neural ODE*, Advances in Neural Information Processing Systems, 33 (2020),
304 pp. 14831–14843.
- 305 [4] G.-H. LIU, T. CHEN, AND E. THEODOROU, *Second-order neural ODE optimizer*, Advances in
306 Neural Information Processing Systems, 34 (2021), pp. 25267–25279.
- 307 [5] S. MASSAROLI, M. POLI, J. PARK, A. YAMASHITA, AND H. ASAMA, *Dissecting neural ODEs*,
308 Advances in Neural Information Processing Systems, 33 (2020), pp. 3952–3963.
- 309 [6] A. NORCLIFFE, C. BODNAR, B. DAY, N. SIMIDJIEVSKI, AND P. LIÒ, *On second order behaviour*
310 *in augmented neural ODEs*, Advances in Neural Information Processing Systems, 33 (2020),
311 pp. 5911–5921.
- 312 [7] C. RACKAUCKAS, M. INNES, Y. MA, J. BETTENCOURT, L. WHITE, AND V. DIXIT, *DiffEqFlux.jl-a*
313 *Julia library for neural differential equations*, arXiv preprint arXiv:1902.02376, (2019).
- 314 [8] C. RACKAUCKAS AND Q. NIE, *DifferentialEquations.jl—a performant and feature-rich ecosystem*
315 *for solving differential equations in Julia*, Journal of open research software, 5 (2017).
- 316 [9] K. ZUBOV, Z. MCCARTHY, Y. MA, F. CALISTO, V. PAGLIARINO, S. AZEGLIO, L. BOTTERO,
317 E. LUJÁN, V. SULZER, A. BHARAMBE, ET AL., *NeuralPDE: Automating physics-informed*
318 *neural networks (PINNs) with error approximations*, arXiv preprint arXiv:2107.09443,
319 (2021).