

The Shape of Digits

A Bayesian Topological Data Analytic Approach to Classification of Handwritten Digits

Thomas Reinke Theophilus A. Bediako Daniel Lim

August 16, 2025

Table of contents

1	Abstract	2
2	Introduction	2
3	Background and Related Work	2
3.1	MNIST	2
3.2	Related Work	3
4	Methodology	3
4.1	Traditional Machine Learning	3
4.2	Our Bayes TDA Methodology	5
4.3	TDA & ML	7
5	Experiments	8
5.1	Data Summary	8
5.2	Model Perfomance	10
6	Discussion	15
7	Conclusion	15
8	References	16

1 Abstract

This paper ...

2 Introduction

The motivation for this work was to compare and examine how dimension reduction via topological data analysis can be used with machine learning and classification models. Algebraic topologist, Gunnar Carlsson, has a quote, “Data has shape, shape has meaning, and meaning brings value.” The work in this paper follows this idea, that if there is inherent structure present in data, it can be exploited to aid in modeling.

The Modified National Institute of Standards and Technology database is a set of handwritten digits that is frequently used to train and test image processing models. Our goal is to classify handwritten digits to their correct numeric label. We compared model performance primarily on accuracy, and secondly on the number of features or predictors. First, we want to build an accurate model that can discriminate between digits. Similar to the idea of parsimony in model selection, if two models have comparable accuracy, we will favor the model trained on fewer features.

3 Background and Related Work

3.1 MNIST

The development of the MNIST dataset stems from early efforts in optical character recognition (OCR), a field of automating the processing of handwritten information like postal codes and census forms. Its lineage can be traced back to the late 1980s with the creation of the USPS database, a collection of 16×16 grayscale images of handwritten zip codes used to train the LeNet neural network. During the same period, the National Institute of Standards and Technology (NIST) was developing its own Special Databases for OCR research, sourcing images from census forms and high-school student samples. One such collection, SD-7, released in 1992, would later form a core part of the MNIST test set.

Challenges in generalizing models across these different datasets, highlighted during a 1992 NIST/Census Bureau competition, revealed biases within the original NIST data. To address these issues, the MNIST database was created in 1994, providing a cleaner, more standardized benchmark for the machine learning community. The dataset led to the development of several modern variations. These include EMNIST (2017), which extended the character set to include letters; QMNIST (2019), which restored the complete original test set; and Fashion MNIST (2017), a drop-in replacement featuring images of clothing items.

3.2 Related Work

Topological Data Analysis is a technique for extracting structural features from datasets, proving effective in classification tasks across various domains. Work in this area includes that of Nicolau et al., who successfully used a topology-based approach to identify a distinct subgroup of breast cancers with excellent survival rates, demonstrating the method’s ability to uncover patterns invisible to other methods (Nicolau, Levine, and Carlsson 2011). Researchers have developed more generalized TDA-based classification methods, applying them to problems involving multiple measurements and other complex data structures (Riihimäki et al. 2019; Kindelan et al. 2021).

Many traditional machine learning models have been successfully applied to MNIST (Yeboah 2025), but it has also been a subject of interest for topological methods. Garin and Tauzin provide a “Topological ‘Reading’ Lesson” by applying TDA specifically to classify MNIST digits, showing that topological features can serve as effective predictors (Garin and Tauzin 2019).

We try an integration of a Bayesian framework to the TDA methodology. This is inspired by the Bayesian framework for persistent homology by Maroulas et al. (Maroulas, Nasrin, and Oballe 2020). Here, we can quantify uncertainty in topological features, in an attempt to increase the predictive power of TDA-based classification models.

4 Methodology

4.1 Traditional Machine Learning

In order to evaluate the value added by TDA-based methods, it is essential to benchmark them against well-established approaches. We consider a range of traditional machine learning (ML) methods, including neural networks with different regularization schemes and classical multinomial logistic regression. Neural networks can overfit, especially on high-dimensional data such as images. Regularization techniques, like dropout, ridge, and lasso help reduce overfitting, and ensure that the network captures meaningful patterns rather than noise.

4.1.1 Neural Networks

The neural network framework considered in this project is the feedforward neural network. It is one of the most widely used architectures for supervised learning. In a feedforward network, information flows in a single direction—from the input layer, through one or more hidden layers, to the output layer—without cycles or feedback connections. The input layer consists of neurons corresponding to the features of the data. One or more hidden layers are placed between the input and output layers, enabling the network to learn complex and non-linear patterns. Finally, the output layer produces the classification results, with the number of neurons equal to the number of digit classes. (Yeboah 2025)

Dropout

Dropout is a regularization method that randomly turns off a fraction of neurons during training. This helps the network avoid relying too much on any single neuron and encourages it to learn more general patterns.

Ridge

In NN ridge regularization, a penalty proportional to the square of the weights is added to the loss function. This shrinks large weights while keeping most parameters nonzero.

Lasso

Lasso regularization penalizes the absolute values of weights, encouraging sparsity in the network. The lasso NN forces many weights towards zero, effectively performing feature selection among pixel intensities.

Multinomial Logistic Regression as ML

Multinomial logistic regression, also known as softmax regression, provides a linear baseline for multi-class classification. It is implemented as a single dense layer with softmax activation (James et al. 2013). Table 1 provides a summary of the model architecture.

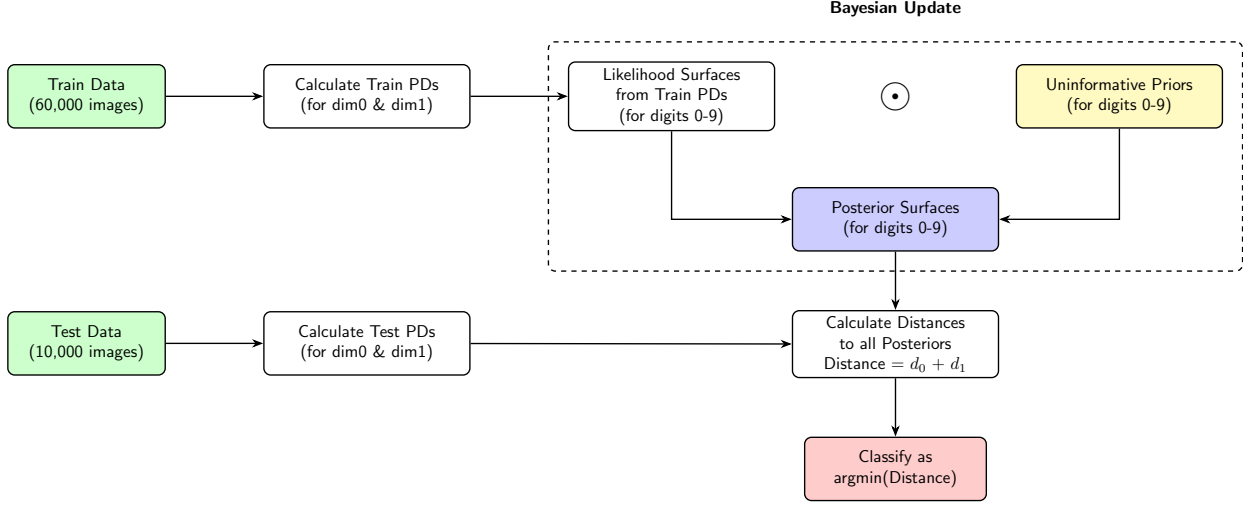
4.1.2 Table 1. Model architectures

Model	Hidden Layers (units)	Activation Functions	Regularization	Output Layer
Multinomial	None	—	None	10 (Soft- max)
Dropout NN	256, 128	ReLU (hidden), Softmax (out)	Dropout (0.4, 0.3)	10 (Soft- max)
Ridge NN (L2)	256, 128	ReLU (hidden), Softmax (out)	L2 penalty ($\lambda = 0.01$)	10 (Soft- max)
Lasso NN (L1)	256, 128	ReLU (hidden), Softmax (out)	L1 penalty ($\lambda = 0.01$)	10 (Soft- max)

All models were trained under identical conditions—30 epochs, batch size of 128, and a validation split of 0.2. This uniform training setup allows for a fair comparison of traditional machine learning methods.

4.2 Our Bayes TDA Methodology

Our methodology can be summarized in this flowchart, where the ‘Bayesian update’ comes from A Bayesian framework for persistent homology.

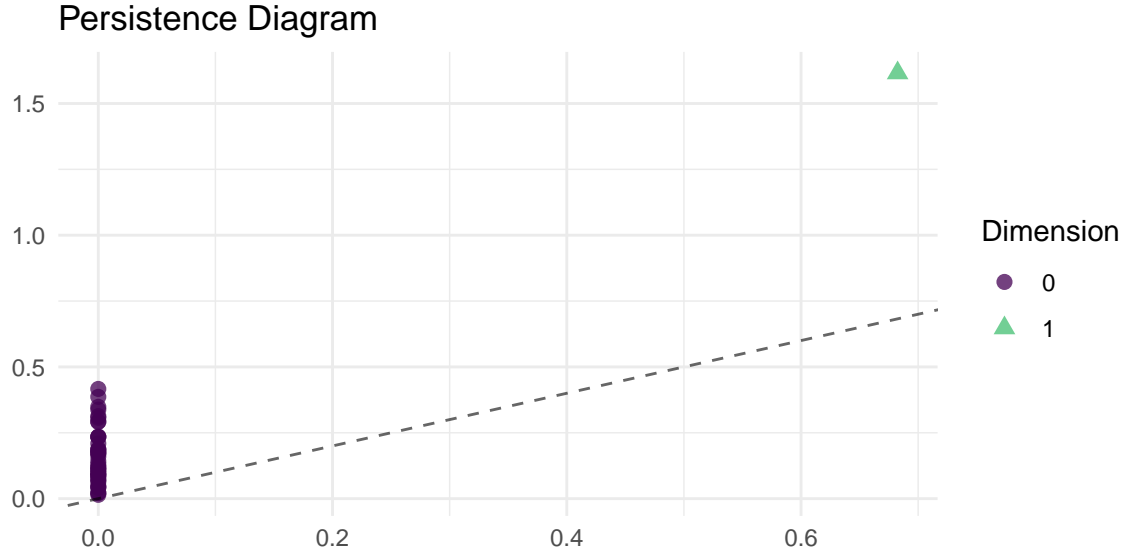


4.2.1 Cubical Complexes & Persistent Homology

For our methodology, we need to get from an image to a persistence diagram. A 2-dimension image is a map $\mathcal{I} : I \subseteq \mathbb{Z}^2 \rightarrow \mathbb{R}$. An element $v \in I$ is a pixel, and has value $\mathcal{I}(v)$, which is the intensity. We can binarize an image by the map $\mathcal{B} : I \subseteq \mathbb{Z}^2 \rightarrow \{0, 1\}$. With data from a point cloud, we typically build simplicial complexes, but with our data from an image, we will build a cubical complex. Pixels are represented by a d -cube, including its faces. With the image represented by a cubical complex K , we build a filtration, which is a sequence of nested subcomplexes, using the image’s grayscale values. We do this with a series of sublevel sets:

$$K_i := \{\sigma \in K \mid I(\sigma) \leq i\}$$

Essentially if a pixel has an intensity less than i , the cube representing it is included in the corresponding complex. After applying persistent homology to this filtration, the birth and death ‘times’ of topological features are tracked across the intensity levels. The persistence diagram D , is a multiset, (b, d, k) , where each point is a homological feature with dimension k , born at intensity b , and dies at intensity d . Persistence is the length of time a feature lasts, $d - b$. In our case, we can only consider $k = 0$, connected components, and $k = 1$, loops/one dimension holes. An example persistence diagram is show below of 40 points sampled from a circle.



From this process we are able to go from an image to a cubical complex to a filtration to a persistence diagram.

4.2.2 Marked Poisson Point Process

A Poisson Point Process Π allows us to model a collection of random points $\{x_1, \dots, x_n\}$ in a space \mathbb{X} , with an intensity measure Λ . The number of points N , is a random variable and follows a Poisson distribution, with mean $\mu = \Lambda(\mathbb{X})$. For a region $A \subseteq X$, $\Lambda(A) = \mathbb{E}[|\Pi \cap A|]$. The mark of each point, m , comes from a space \mathbb{M} , in our case, the marks will be the dimension k . We first have our set of locations $\{x_i\}$, then for each x_i , a mark m_i is drawn conditionally & independently from a kernel $\ell(x_i, \cdot)$.

4.2.3 Bayes Update & Gaussian representation

Now we wish to incorporate a Bayesian framework to these persistence diagrams represented as point processes. Here we use the tilted representation of the diagram, so instead of (b, d) , we use $(b, d - b)$ for a persistence diagram D .

The first part we model is the latent or ‘true’ underlying persistence diagram D_x ; we model it by the intensity $\lambda_{D_x}(x)$. D_x is decomposed into two independent parts. D_{XO} are the points that can be observed with probability $\alpha(x)$: $\alpha(x)\lambda_{D_x}(x)$. The second part is for points that are missed or not observed: $(1 - \alpha(x))\lambda_{D_x}(x)$.

The second part we model is the observed persistence diagram D_Y , also two components. D_{YO} are the points generated from D_{XO} , which forms the pair $(\mathcal{D}_{XO}, \mathcal{D}_{YO})$. The connection is via the kernel $\ell(y|x)$. It gives the probability density of observing $y \in \mathcal{D}_{YO}$ given a latent point $x \in \mathcal{D}_{XO}$. The other component is that the points arise from noise: $\lambda_{D_{YS}}(y)$.

Now we bring the two parts together. The posterior intensity $\lambda_{\mathcal{D}_X|D_{Y^{1:m}}}(x)$ for the latent PD \mathcal{D}_X , given m independent observed PDs D_{Y^1}, \dots, D_{Y^m} . Let $D_{Y^{1:m}} = \cup_{i=1}^m D_{Y^i}$. The posterior intensity is:

$$\lambda_{\mathcal{D}_X|D_{Y^{1:m}}}(x) = \underbrace{(1 - \alpha(x))\lambda_{\mathcal{D}_X}(x)}_{\text{Prior Vanished Part}} + \underbrace{\frac{1}{m}\alpha(x) \sum_{i=1}^m \sum_{y \in D_{Y^i}} \frac{\ell(y|x)\lambda_{\mathcal{D}_X}(x)}{\lambda_{\mathcal{D}_{Y_S}}(y) + \int_W \ell(y|u)\alpha(u)\lambda_{\mathcal{D}_X}(u)du}}_{\text{Update from Observed Points } y} \quad \text{a.s.}$$

Without going into all of details, Maroulas achieves this is computationally using Gaussian mixtures for $\lambda_{\mathcal{D}_X}$, $\ell(y|x)$, and $\lambda_{\mathcal{D}_{Y_S}}(y)$.

4.2.4 Classification

Now that we are able to a posterior persistence diagram, we can do this for each digit 0-9. We can then calculate the Wasserstein distance¹ from the test persistence diagrams to the posterior diagrams for each dimension. Then we sum the distances for each dimension and classify the test image as the class associated with the minimum distance.

4.3 TDA & ML

With our Bayes TDA classification not being as successful as we would've liked, we were also interested in using TDA as a dimension reduction technique to pair with machine learning models. This approach uses the persistence diagram of an image as a predictor. Besides using the grayscale image, we can binarize the image, apply a filter that results in a different grayscale image for the digit. This allows us to generate multiple features for each observation.

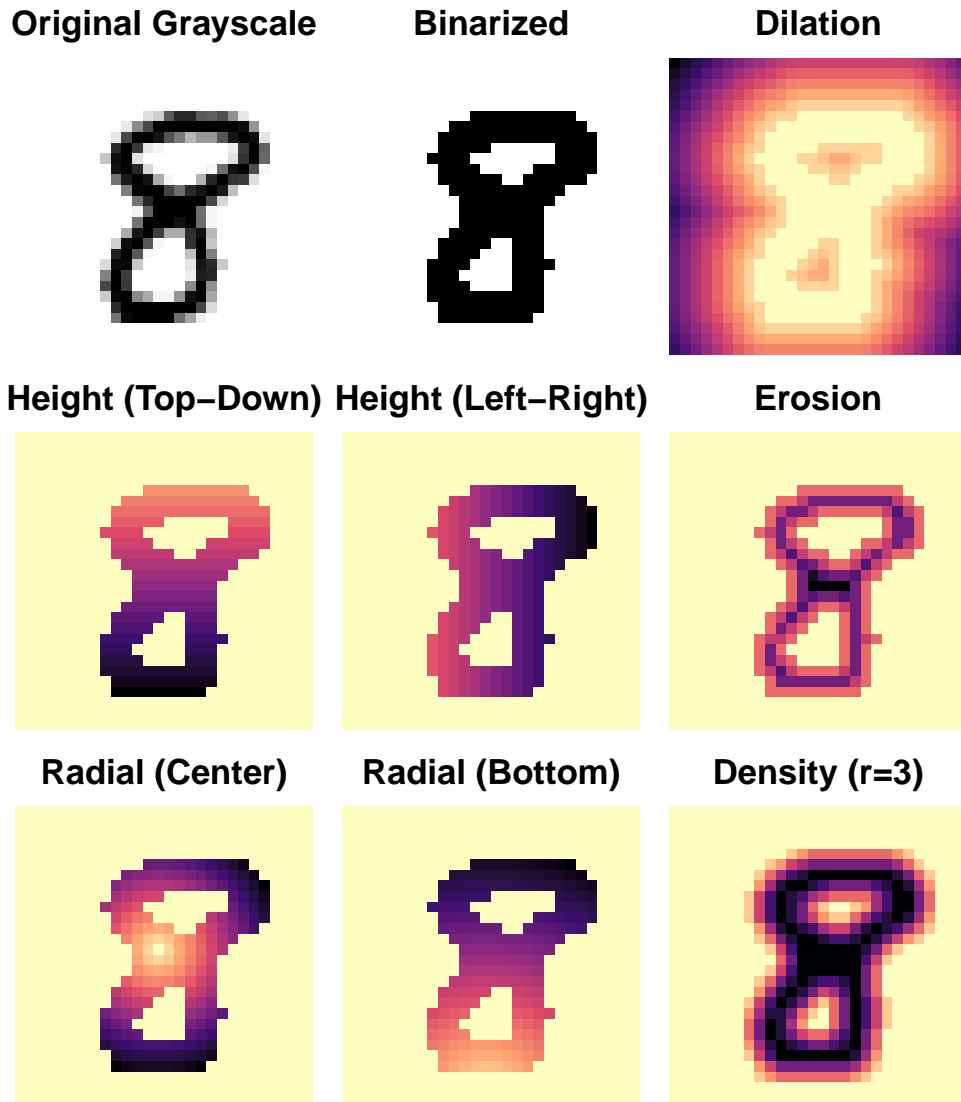
4.3.1 Filtering

The filters we used after binarizing the images were height, radial, dilation, erosion, and density.

The height filtration assigns a value to each pixel based on its projection onto a chosen direction vector, essentially measuring its “height” from a specific angle. Radial filtration works similarly, instead of assigning a value base on the distance to a vector, it assigns a value to each pixel based on its distance from a chosen center point. Dilation assigns each pixel a value corresponding to its shortest distance to a foreground pixel, where a foreground pixel is a pixel with an intensity of 1. This has the effect of “growing” or “dilating” the digit. Erosion is the inverse of dilation. It works by taking the dilation of the inverted image. This “shrinks” or “erodes” the digit. Lastly, density assigns each pixel a value based on the number of foreground neighbors within a given radius.

After applying the filter to the binarized image, we have a new grayscale image, which we can get a persistence diagram for. Examples of these filters are shown below. For purposes of visibility, a colored map is used to represent the grayscale values.

¹Also known as Kantorovich-Rubinstein metric, it is a distance function between probability distributions on a metric space \mathbb{M} . $W_q(X, Y) = \left(\inf_{\eta: X \rightarrow Y} \sum_{x \in X} \|x - \eta(x)\|_\infty^q \right)^{1/q}$. See more [here](#)

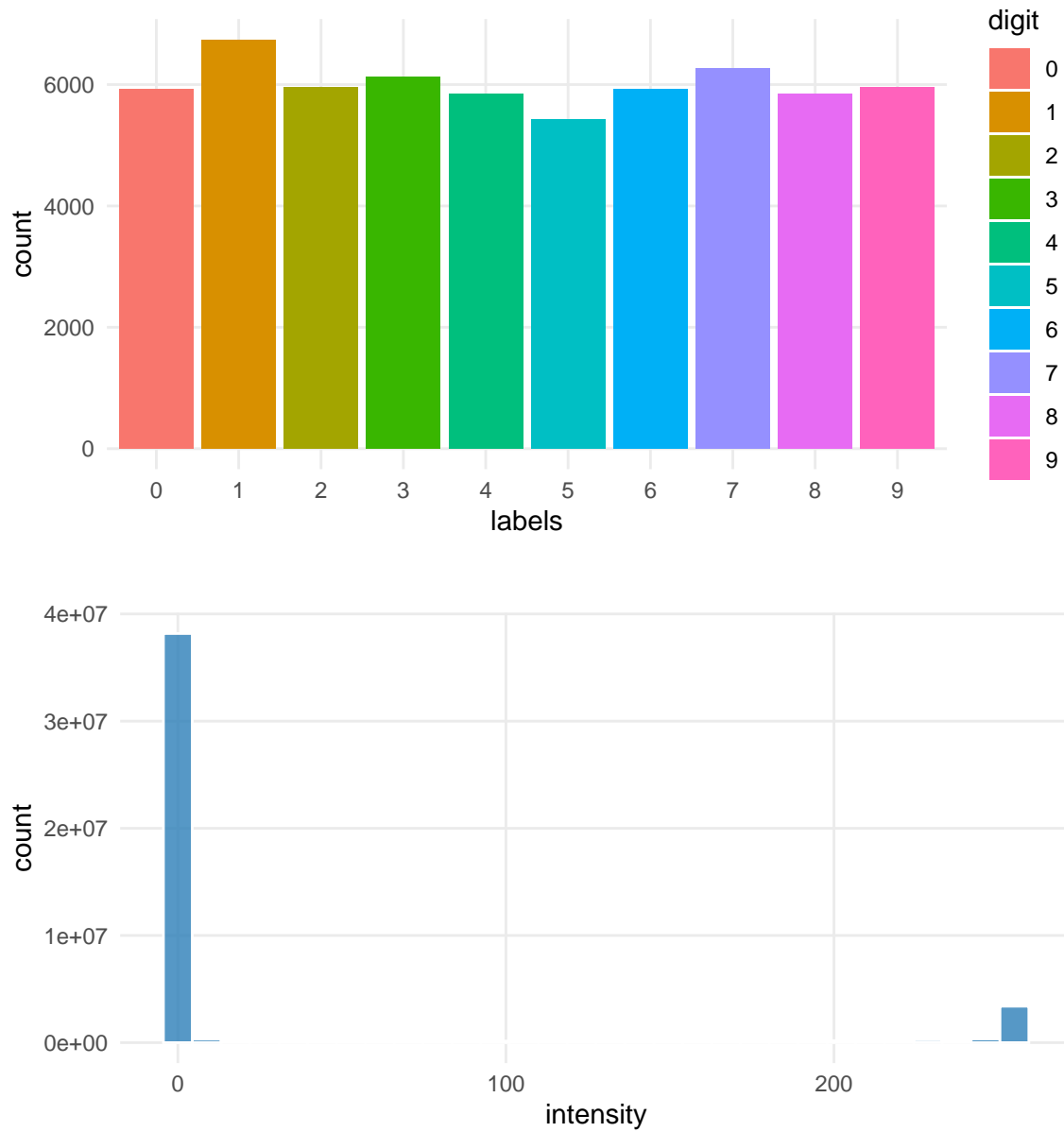


5 Experiments

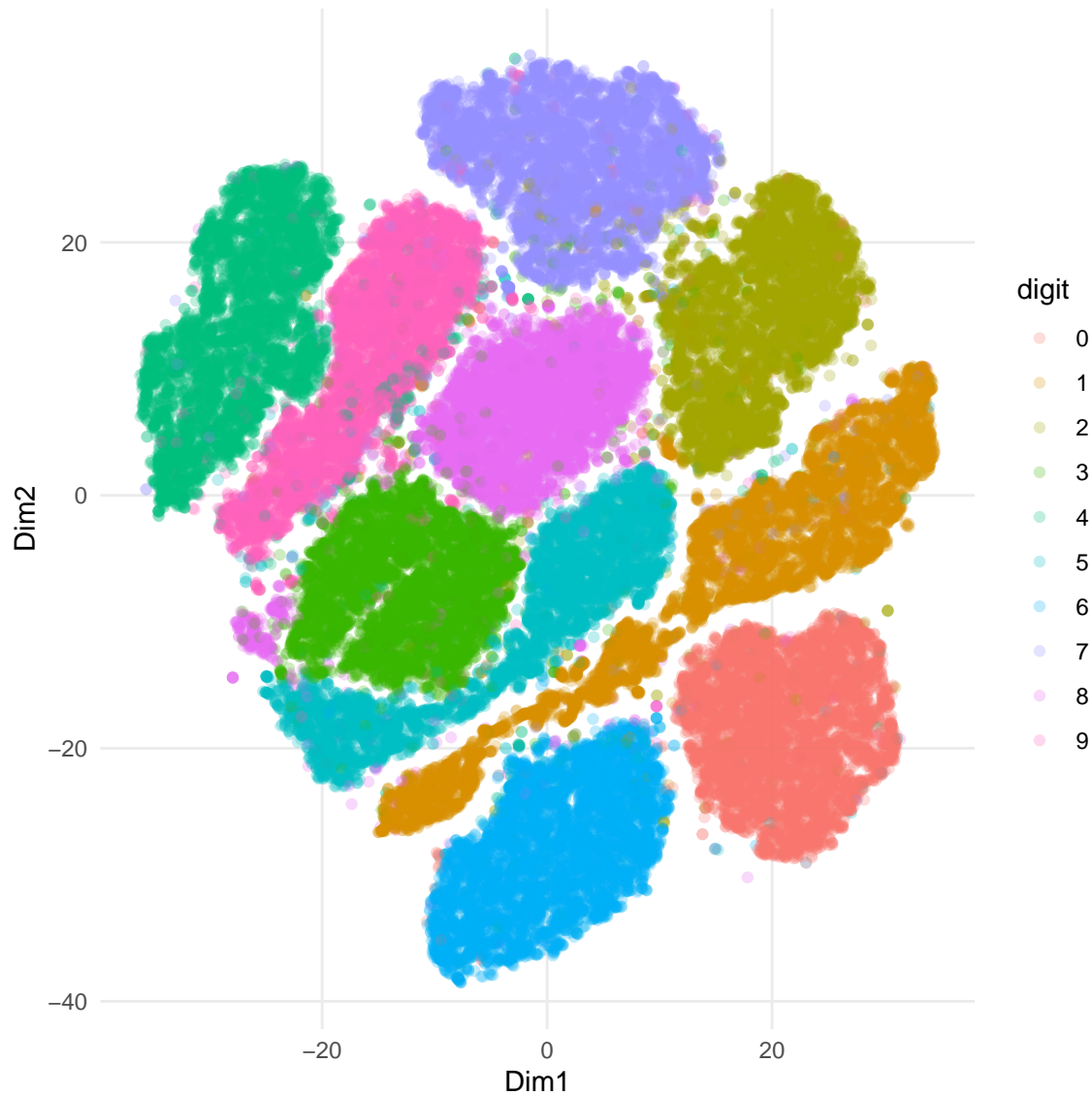
5.1 Data Summary

The MNIST data set is composed of the training set of 60,000 images (50% SD-3, 50% SD-7) and Test set of 10,000 images. Each image is 28x28 pixels, and are normalized to center of mass alignment.

5.1.1 EDA



We can see there are around 6000 digits in each class, so class imbalance is not an issue here. Most pixel values are zero, but out of the 784 pixels, 703 have an intensity greater than 0 at least once.



Using t-distributed Stochastic Neighbor Embedding(t-SNE), we can represent the data in 2D.² t-SNE is a technique for dimension reduction that is well-suited for visualizing high-dimensional data in a lower-dimensional space.

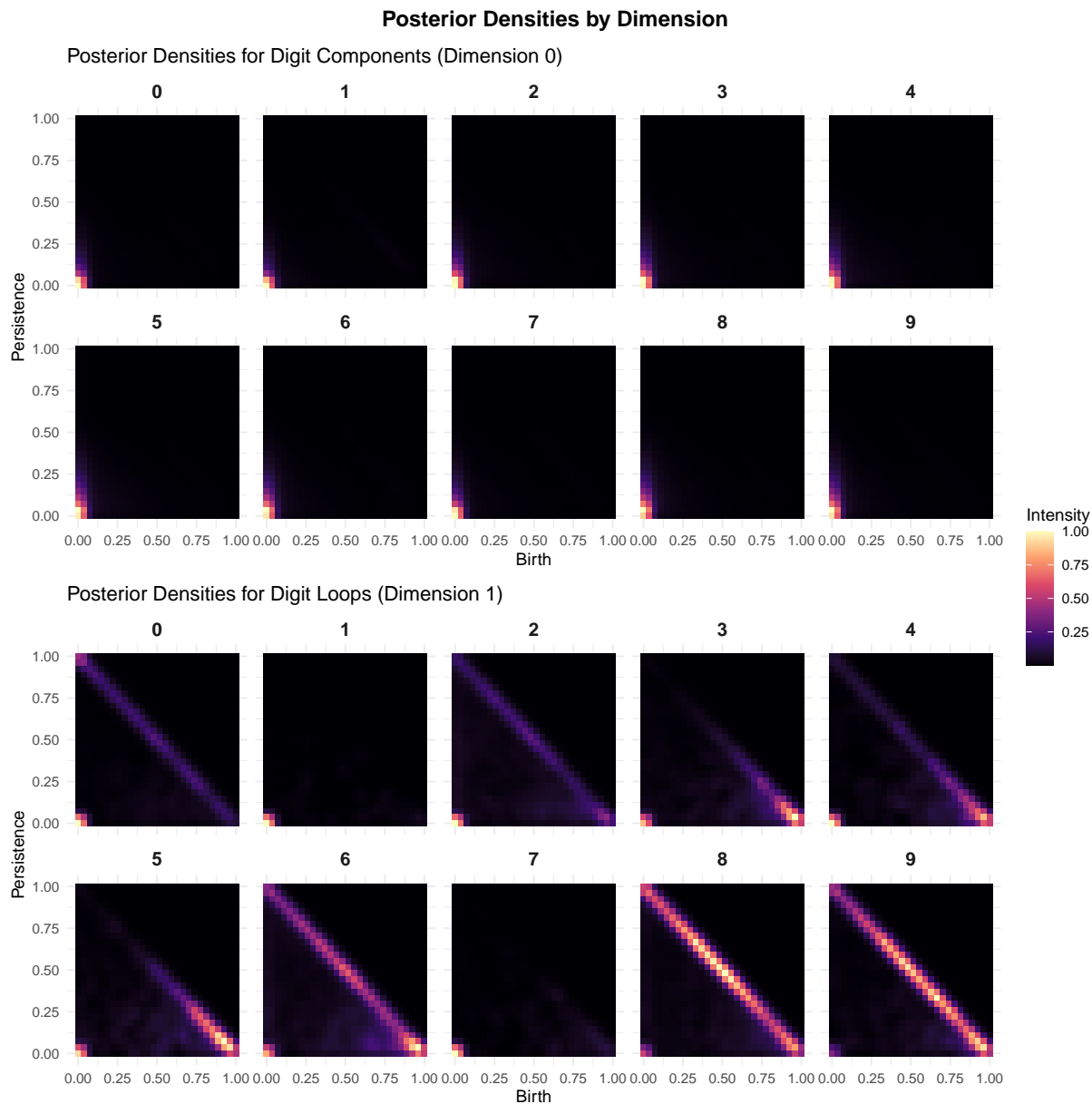
5.2 Model Performance

5.2.1 Bayes TDA Model

We can see the posterior densities for dimension 0 and 1 for our model. We see that the densities are essentially all the same for dimension 0, which is what we would expect, since there are no disconnected components in any of the digits. This would not be the case for letters, e.g., the dot

²See [T-SNE Exploration by TusVasMit](#) for more details.

over lowercase ‘i’ and ‘j’. For the 1st dimension, we see some similarities, but the posterior is about what we would expect. Unfortunately, some of the posteriors are fairly similar and it may be hard to discriminate between them.



For each of the neural network models, we can look at the confusion matrix to compare the predicted labels to the true labels.

Neural Net (Dropout)

True Label \ Predicted Label	0	1	2	3	4	5	6	7	8	9
0	971	0	5	0	1	2	2	1	4	4
1	1	1130	0	1	0	0	3	4	0	4
2	0	1	1012	3	1	0	0	12	4	0
3	1	0	2	990	0	0	1	1	2	5
4	0	0	1	0	959	0	3	1	4	8
5	0	0	0	7	0	882	3	0	4	7
6	5	2	3	0	5	4	946	0	4	1
7	1	0	5	5	0	0	0	1002	2	6
8	1	2	4	3	1	2	0	2	945	0
9	0	0	0	1	15	2	0	5	5	974

Neural Net (Ridge)

True Label \ Predicted Label	0	1	2	3	4	5	6	7	8	9
0	965	0	4	1	0	2	7	0	1	4
1	0	1113	0	0	0	0	2	3	0	2
2	4	2	1008	4	3	0	1	10	3	0
3	0	3	4	991	0	26	1	5	8	8
4	0	0	2	0	934	0	3	0	0	8
5	0	1	0	3	0	849	2	0	0	1
6	6	3	1	0	8	6	937	0	3	1
7	2	2	8	4	2	1	0	1002	3	11
8	3	11	5	7	3	7	5	3	955	8
9	0	0	0	0	32	1	0	5	1	966

Neural Net (Lasso)

True Label \ Predicted Label	0	1	2	3	4	5	6	7	8	9
0	956	0	7	0	1	6	12	0	3	3
1	0	1111	4	1	0	2	3	6	3	4
2	1	1	925	4	3	0	2	9	1	0
3	1	5	14	951	0	30	1	3	12	11
4	0	0	15	0	899	5	15	3	5	13
5	3	0	2	21	0	809	13	0	4	4
6	4	2	7	0	5	3	894	0	4	1
7	5	4	40	22	6	11	5	1000	16	31
8	5	12	12	10	4	20	13	4	924	14
9	5	0	6	1	64	6	0	3	2	928

Multinomial Logistic

True Label \ Predicted Label	0	1	2	3	4	5	6	7	8	9
0	961	0	5	3	1	9	9	2	7	11
1	0	1115	12	0	2	3	3	7	8	7
2	0	3	922	20	5	3	6	20	6	1
3	2	1	14	909	2	26	1	5	16	9
4	0	0	8	0	917	7	7	7	8	27
5	6	1	4	29	0	785	11	0	24	8
6	8	4	14	3	10	15	917	0	11	0
7	2	2	10	13	6	8	2	963	13	40
8	1	9	41	27	9	32	2	1	877	7
9	0	0	2	6	30	4	0	23	4	899

Now we compare with the confusion matrix for our method.

0	374	72	3	38	0	39	0	0	280	47
1	0	856	0	4	0	1	0	1	0	0
2	147	580	4	22	1	17	2	2	148	30
3	7	832	1	62	0	15	2	4	4	1
4	24	798	2	43	0	8	0	1	27	5
5	6	746	2	49	0	19	0	0	10	4
6	269	110	5	30	8	33	11	2	318	74
7	0	876	0	9	0	2	0	1	1	0
8	334	15	1	10	4	25	1	0	417	95
9	236	104	1	78	1	68	2	3	377	61
	0	1	2	3	4	5	6	7	8	9

True Label

Predicted Label

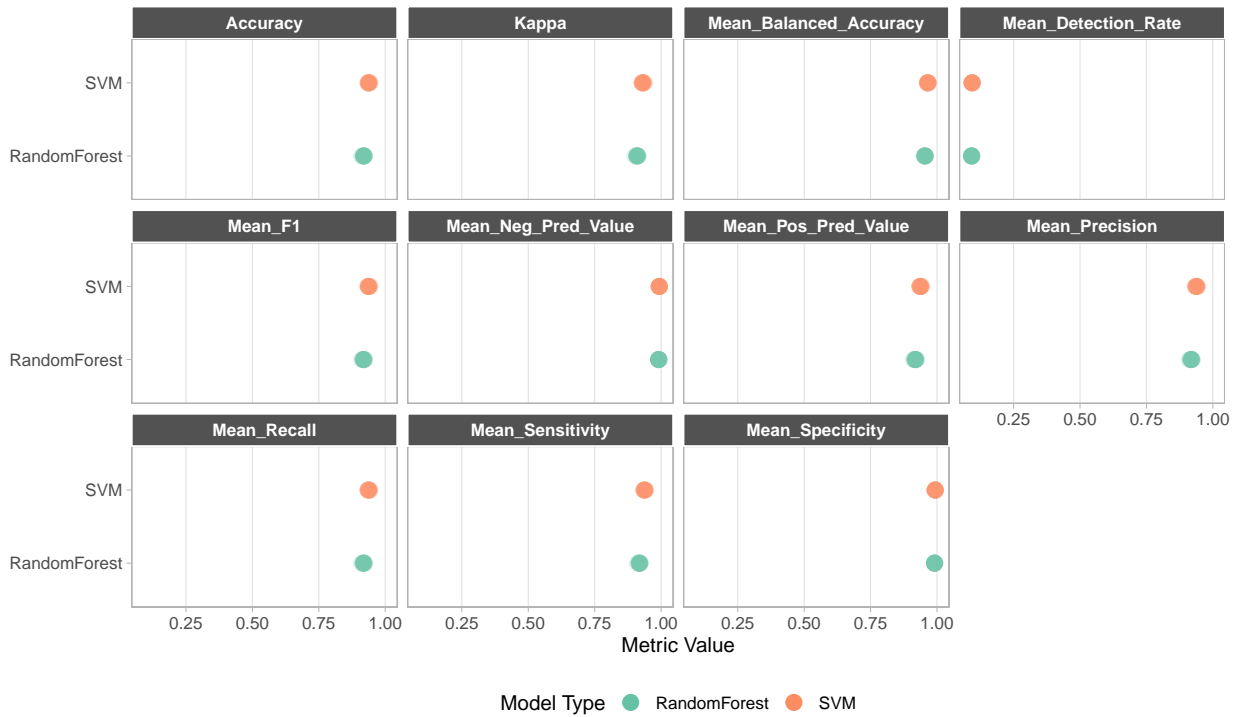
This appears to be much worse, and seeing the accuracy of each model confirms this.

method	accuracy
multinomial	0.9856
dropout nn	0.9962
ridge nn	0.9946
lasso nn	0.9946
proposed	0.2023

The neural network with dropout performs the best, achieving an accuracy over 99%. We discuss the likely failings of our model in the discussion section.

5.2.2 TDA & ML

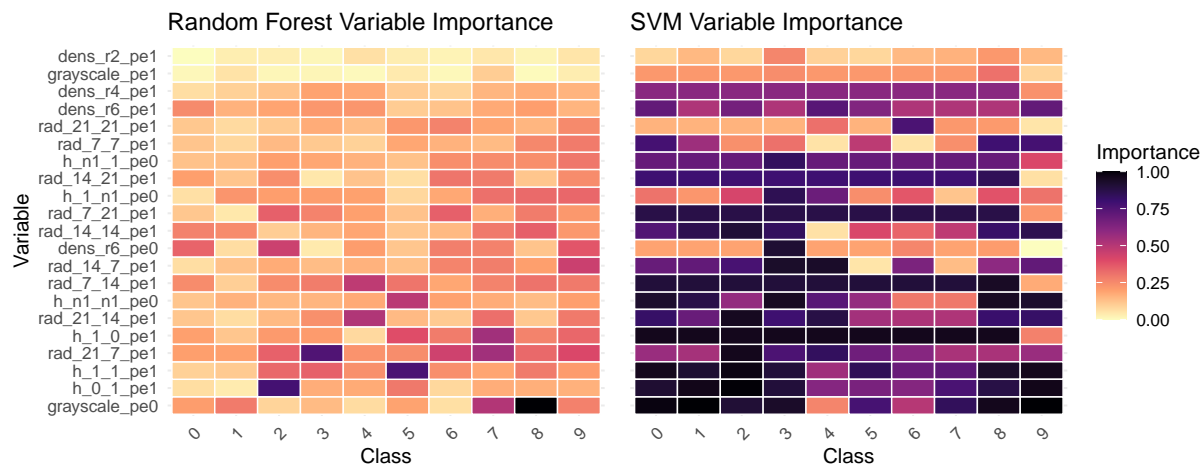
For our TDA with machine learning approach. We used 1 grayscale, 8 height, 1 dilation, 1 erosion, 9 radial, and 3 density features. We then used SVM with a radial basis function(RBF) kernel and a Random Forest model. Below is a comparison between the two models for various performance metrics. SVM is equal to or better than random forest in all metrics, so we will use it for prediction.



The confusion matrix below shows how well the SVM model performed. It achieved an accuracy of 93.1%. This isn't as good as the neural network, but performs comparably for having almost 1/40th the number of features.

True Label	0	934	0	4	1	2	1	11	1	19	7
	1	0	1109	6	0	4	0	5	3	5	3
	2	3	4	963	26	1	3	2	18	5	7
	3	1	0	47	911	2	16	0	12	12	9
	4	1	0	9	6	932	1	1	6	8	18
	5	2	0	11	18	0	831	4	4	14	8
	6	5	7	1	1	4	12	914	0	11	3
	7	1	3	17	11	13	3	1	942	6	31
	8	15	5	9	10	16	8	12	9	873	17
	9	9	6	8	6	15	6	3	29	26	901
		Predicted Label									

The benefit here is that the features here may have a better interpretation. The most important variable in both the random forest and SVM model was the grayscale 0 dimensional persistent entropy. Persistent entropy is a measure of the complexity of the persistence diagram.



6 Discussion

Comparing the models

7 Conclusion

8 References

- Garin, Adélie, and Guillaume Tauzin. 2019. “A Topological ”Reading” Lesson: Classification of MNIST Using TDA.” *CoRR* abs/1910.08345. <http://arxiv.org/abs/1910.08345>.
- James, Gareth, Daniela Witten, Trevor Hastie, and Robert Tibshirani. 2013. *An Introduction to Statistical Learning: With Applications in r*. Springer. <https://faculty.marshall.usc.edu/gareth-james/ISL/>.
- Kindelan, Rolando, José Frías, Mauricio Cerda, and Nancy Hitschfeld. 2021. “Classification Based on Topological Data Analysis.” *CoRR* abs/2102.03709. <https://arxiv.org/abs/2102.03709>.
- Maroulas, Vasileios, Farzana Nasrin, and Christopher Oballe. 2020. “A Bayesian Framework for Persistent Homology.” *SIAM J. Math. Data Sci.* 2 (1): 48–74.
- Nicolau, Monica, Arnold J. Levine, and Gunnar Carlsson. 2011. “Topology Based Data Analysis Identifies a Subgroup of Breast Cancers with a Unique Mutational Profile and Excellent Survival.” *Proceedings of the National Academy of Sciences* 108 (17): 7265–70. <https://doi.org/10.1073/pnas.1102826108>.
- Riihimäki, Henri, Wojciech Chachólski, Jakob Theorell, Jan Hillert, and Ryan Ramanujam. 2019. “A Topological Data Analysis Based Classification Method for Multiple Measurements.” *CoRR* abs/1904.02971. <http://arxiv.org/abs/1904.02971>.
- Yeboah, Felix. 2025. “Classification and Evaluation of Machine Learning Algorithms on the MNIST Dataset.”