

# 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

<b>1</b>	<b>Abstract</b>	<b>2</b>
<b>2</b>	<b>Introduction</b>	<b>2</b>
<b>3</b>	<b>Background and Related Work</b>	<b>2</b>
3.1	MNIST . . . . .	2
3.2	Related Work . . . . .	3
<b>4</b>	<b>Methodology</b>	<b>3</b>
4.1	Traditional Machine Learning . . . . .	3
4.2	Our Bayes TDA Methodology . . . . .	5
4.3	TDA & ML . . . . .	7
<b>5</b>	<b>Experiments</b>	<b>8</b>
5.1	Data Summary . . . . .	8
5.2	Model Perfomance . . . . .	10
<b>6</b>	<b>Discussion</b>	<b>15</b>
<b>7</b>	<b>Conclusion</b>	<b>15</b>
<b>8</b>	<b>References</b>	<b>17</b>
<b>9</b>	<b>Code Appendix</b>	<b>18</b>

# 1 Abstract

This paper compares traditional machine learning models with a Bayesian Topological Data Analytic model to predict the numerical labels of a handwritten digits from the MNIST dataset. It also considers a Topological Data Analysis approach paired with machine learning models, to achieve comparable accuracy to the traditional models with a fraction of the features and offer insight to how the features influence classification. We found our Bayesian model to severely underperform, whereas our machine learning models with topological features performed with high accuracy. Models that accurately classifying digits can be expanded and adapted to other character or image recognition systems.

# 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. We also consider the ability of TDA-based models to assess the importance of individual predictors, providing insights into which features most influence MNIST classification.

# 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  $16 \times 16$  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.

## 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.

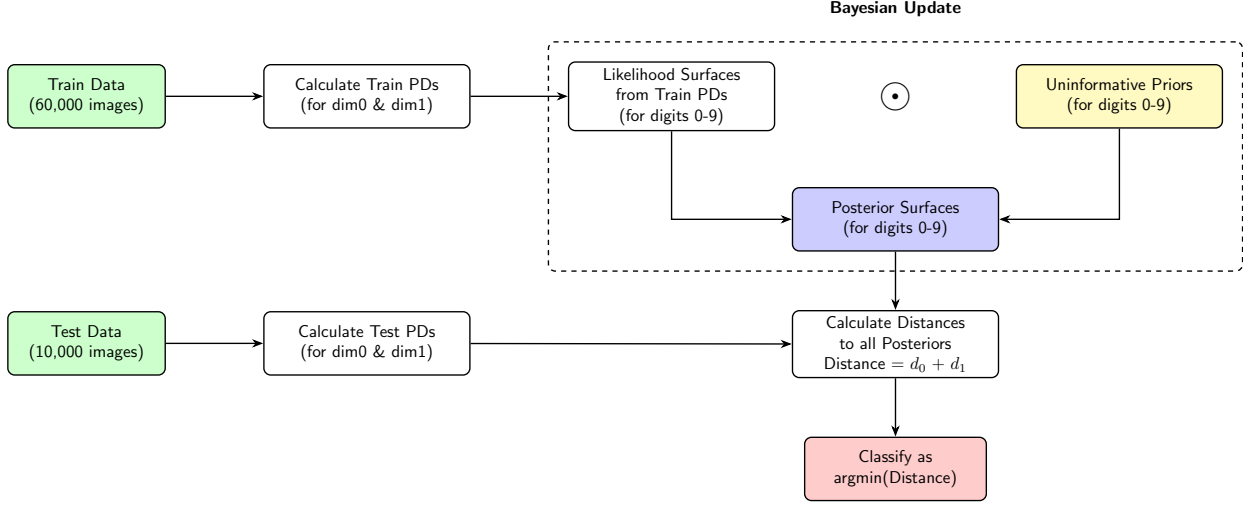
**4.1.2 Table 1. Model architectures**

Model	Hidden Layers (units)	Activation Functions	Regularization	Output Layer
<b>Multinomial</b>	None	–	None	10
<b>Dropout NN</b>	256, 128	ReLU (hidden), Softmax (out)	Dropout (0.4, 0.3)	10
<b>Ridge NN (L2)</b>	256, 128	ReLU (hidden), Softmax (out)	L2 penalty ( $\lambda = 0.01$ )	10
<b>Lasso NN (L1)</b>	256, 128	ReLU (hidden), Softmax (out)	L1 penalty ( $\lambda = 0.01$ )	10

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.

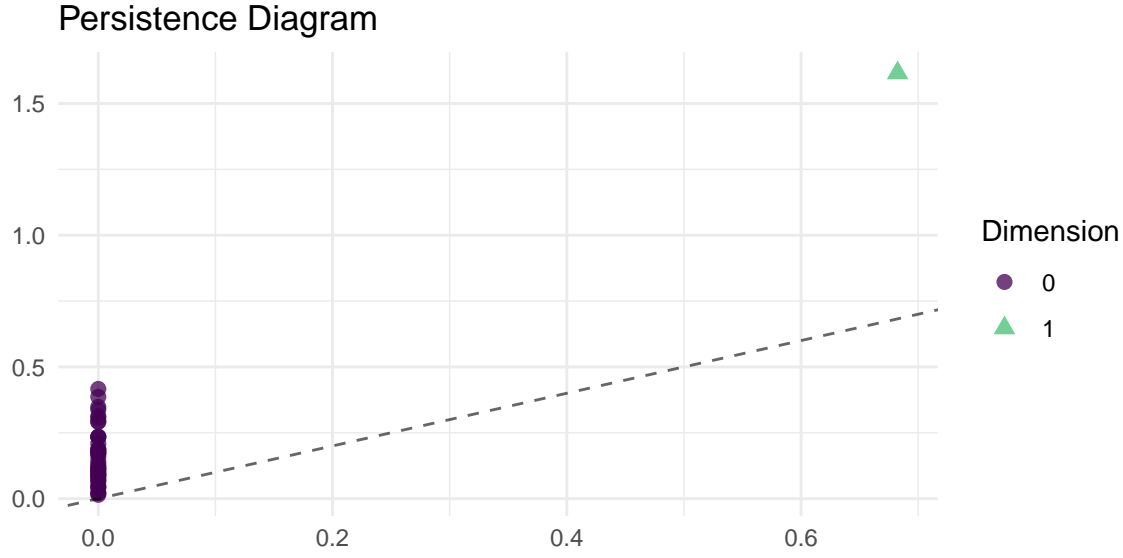


### 4.2.1 Cubical Complexes & Persistent Homology

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  $\Pi$  allows us to model a collection of random points  $\{x_1, \dots, x_n\}$  in a space  $\mathbb{X}$ , with an intensity measure  $\Lambda$ . 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:

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

Without going into all of the details in Maroulas' paper, this is achieved 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 distance<sup>1</sup> 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.

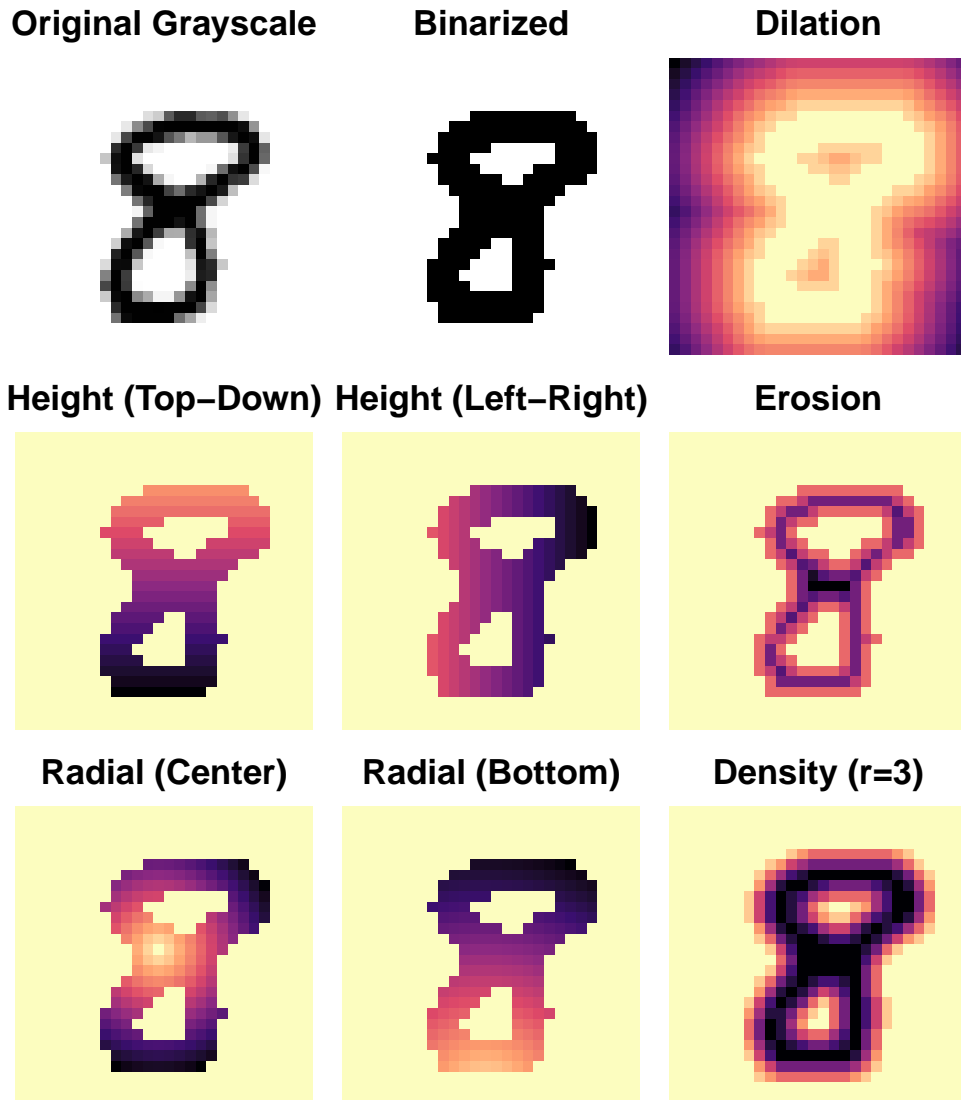
#### 4.3.1 Filtering

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.

<sup>1</sup>Also 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

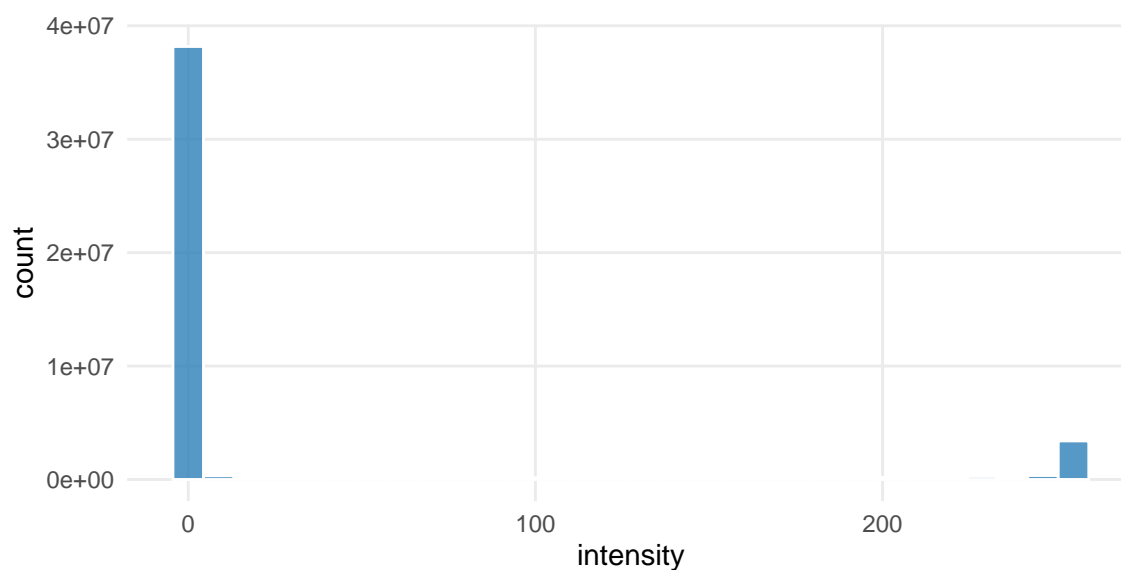
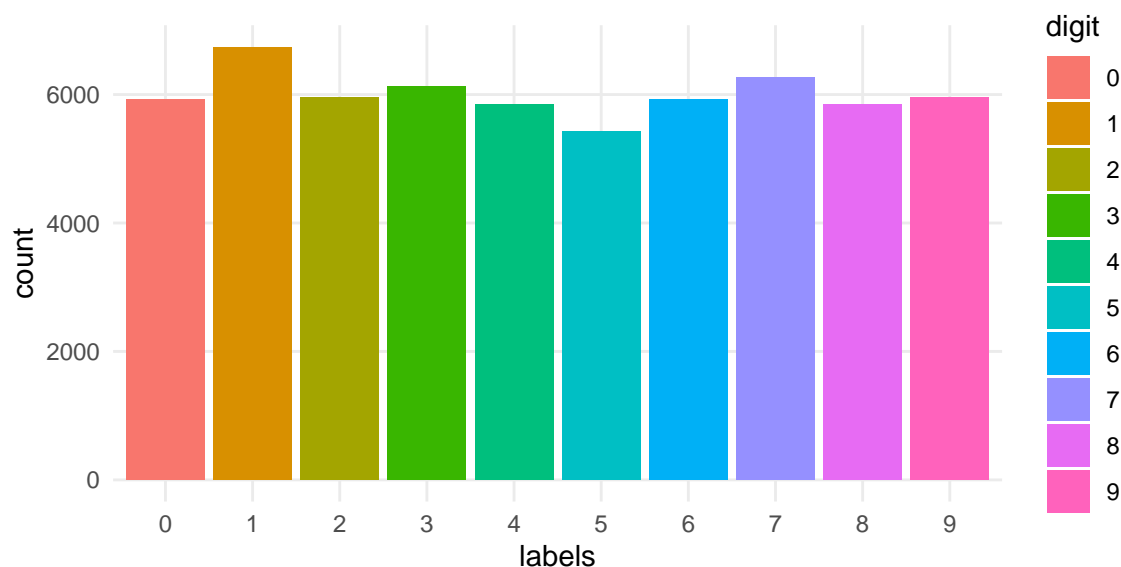
The MNIST data set is composed of the training set of 60,000 images (50% SD-3, 50% SD-7) and Test set of 10,000 images. Each image is 28x28 pixels, and are normalized to center of mass alignment.

#### 5.1.1 EDA

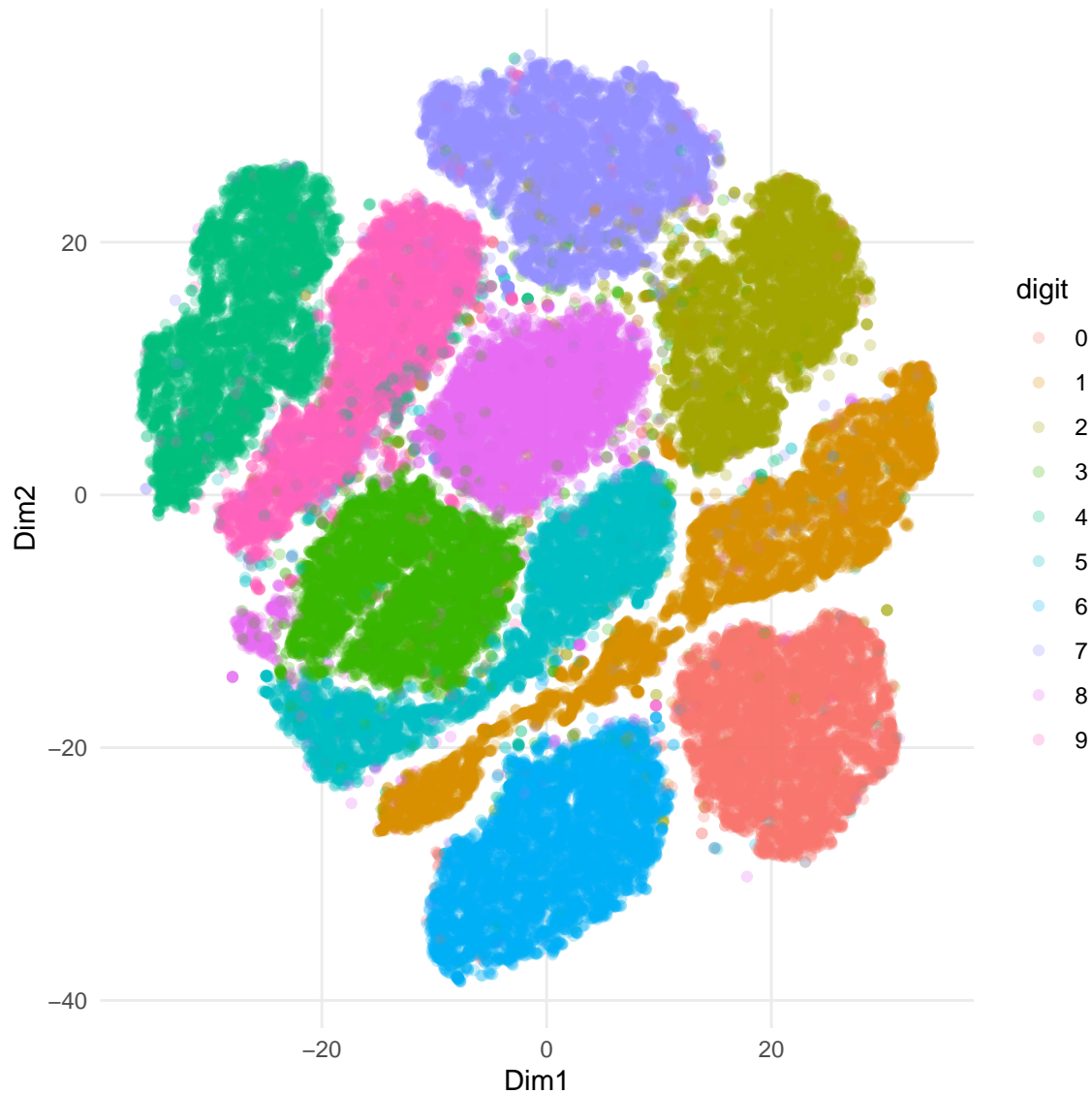
We examined the MNIST training dataset in terms of label distribution, pixel intensity, and overall structure. Label distribution is balanced, with approximately 6,000 examples per digit. We can



expect our classifiers not to be biased toward any particular class. Pixel intensity analysis reveals that most pixels are zero, reflecting the black background of the images. We also use t-Distributed Stochastic Neighbor Embedding(t-SNE)<sup>2</sup> as developed by (Krijthe 2015) to provide visualization of the 784-dimensional pixel space in two dimensions. The figure shows well-separated clusters for each digit, suggesting that the classes are distinguishable and that supervised models can capture these patterns effectively.



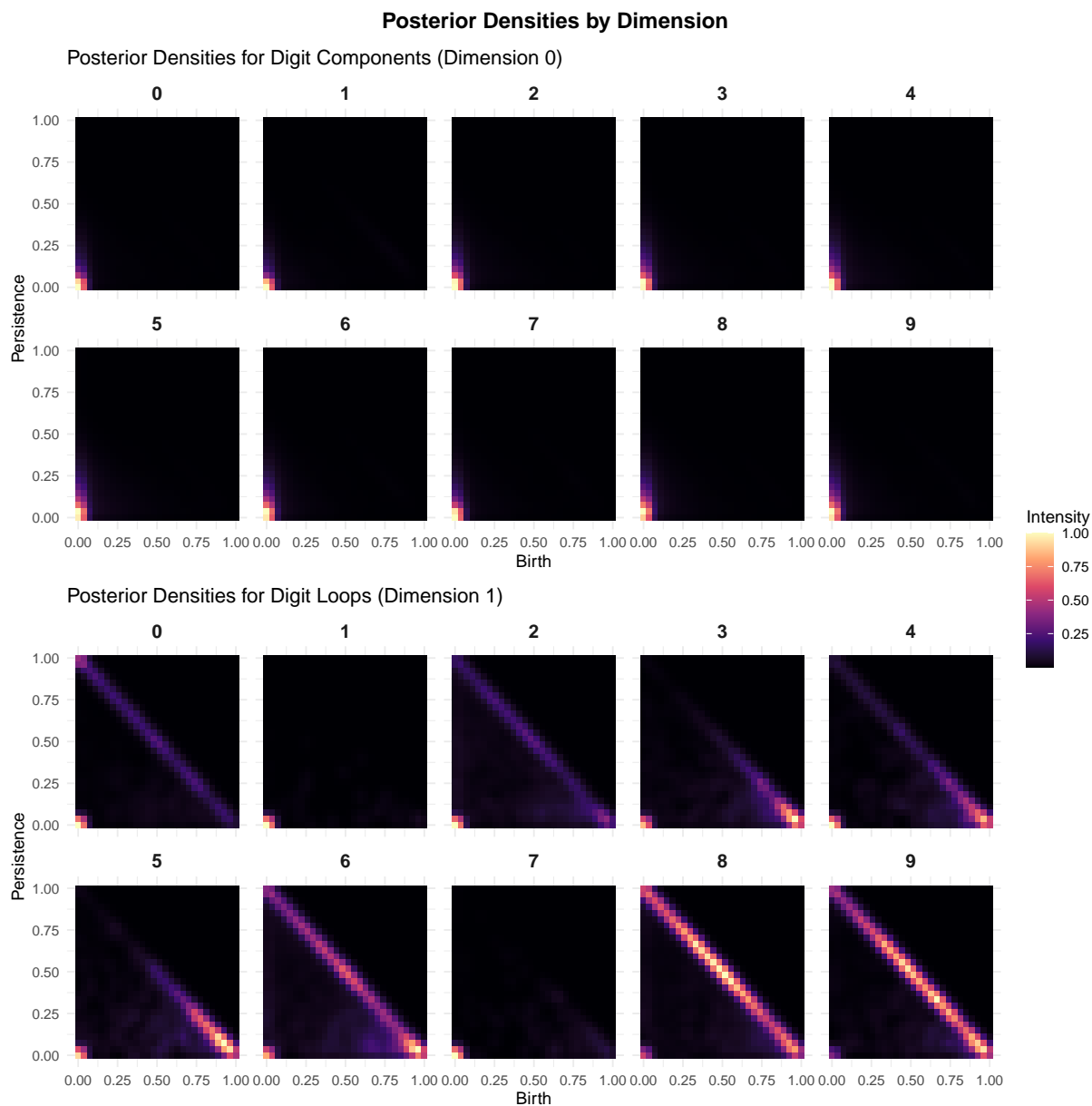
<sup>2</sup>See [T-SNE Exploration by TusVasMit](#) for more details.



## 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 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

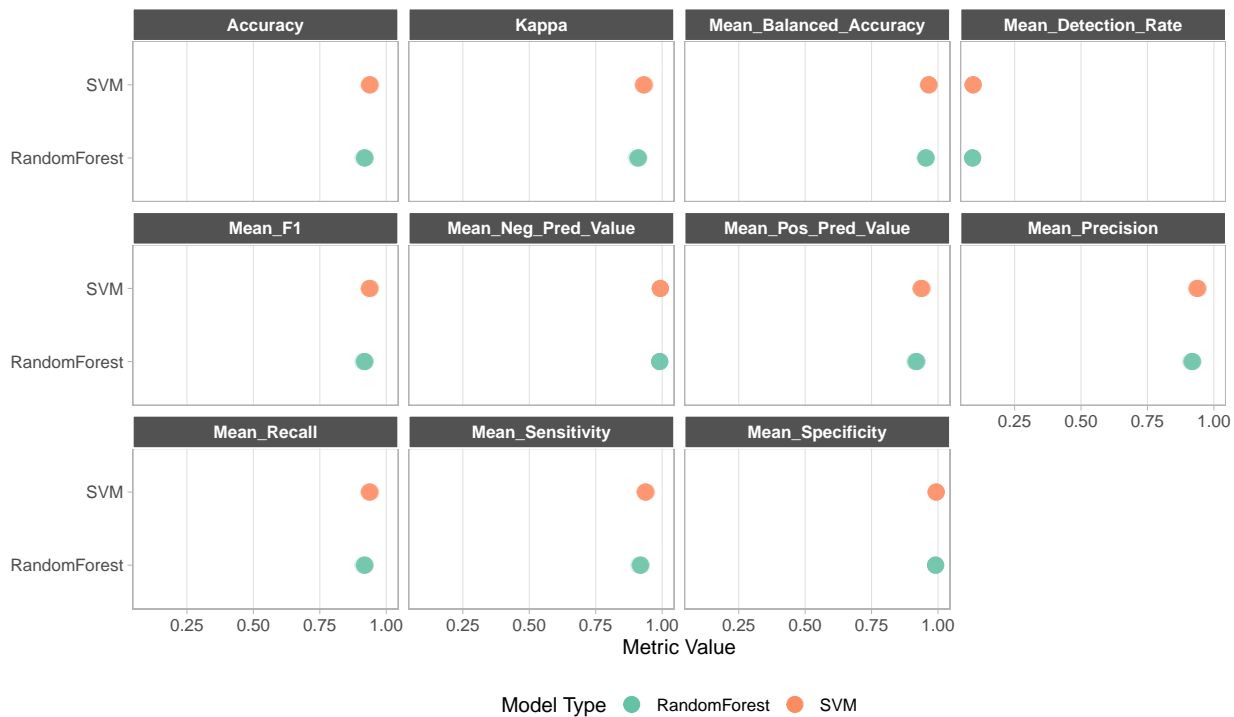
This appears to be much worse, and seeing the accuracy of each model confirms this.

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

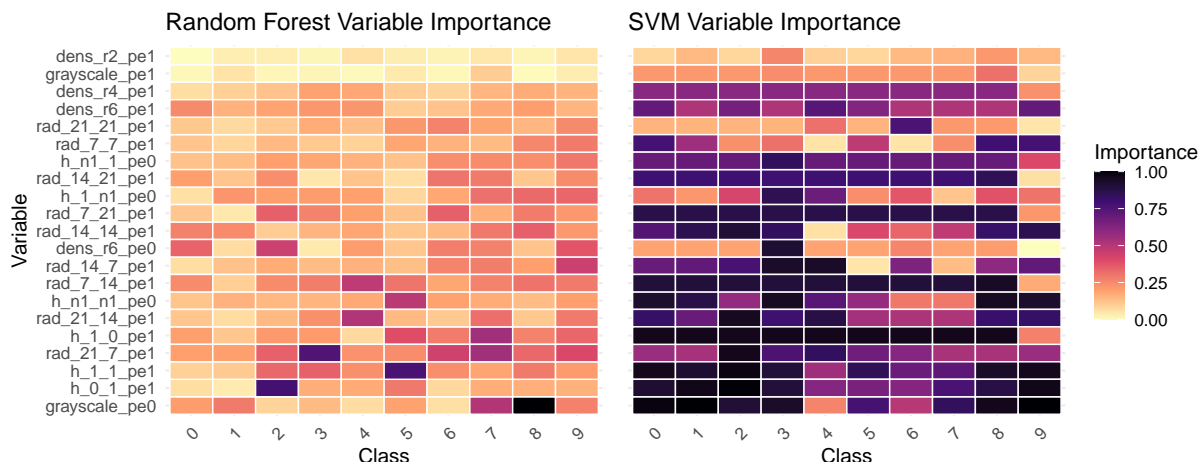
For our TDA with machine learning approach. We used 1 grayscale, 8 height, 1 dilation, 1 erosion, 9 radial, and 3 density features. We then used SVM with a radial basis function(RBF) kernel and a Random Forest model. Below is a comparison between the two models for various performance metrics. SVM is 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

Traditional machine learning methods performed very well, achieving 99% accuracy using all 784 pixels as predictors. However, our Bayesian TDA approach did not perform as well as we had hoped. The posterior densities for the persistence diagrams were not distinct enough to allow for accurate classification. The main issue is that the posterior densities for the 0 dimension were very similar across digits, making it difficult to discriminate between them. The 1st dimension was better, but not enough for discrimination. Another issue is that for a test image, if no features were detected, it was the null model, which was closest to the digit ‘1’. This is something we would expect for {1, 2, 3, 4, 5, 7}.

A better approach likely would have been not using the Bayesian framework and start with a simple model. For example, use persistence diagrams to obtain the Betti number<sup>3</sup> for each dimension, which should split up the digits into 3 sub classes, {1, 2, 3, 4, 5, 7}, {0, 6, 9}, {8}. Then find another feature to discriminate within each subclass.

Using TDA incorporated with machine learning models, we were able to achieve 93.1% accuracy with only 23 features. This was a significant improvement to our methodology, and gives similar results to the paper it was based on (Garin and Tauzin 2019). The most important variable in both the random forest and SVM model was the grayscale 0 dimensional persistent entropy.

## 7 Conclusion

In this project, we used a Bayesian topological data analytic approach to digit classification in the MNIST handwritten dataset and compared it with traditional machine learning models. We also used TDA as a dimension reduction technique to pair with machine learning models. We

<sup>3</sup>Informally, the kth Betti number refers to the number of k-dimensional holes on a topological surface. See more [here](#)

found that the Bayesian TDA approach did not perform well. However, using TDA with machine learning approach achieved 93.1% accuracy with only 23 features.

Future work could involve using a simpler model with Betti numbers or other topological features. For the TDA & machine learning model, other machine learning models could be explored. Better fine tuning on the parameters may bump up the accuracy. Including more features and performing feature selection may also help model performance.



## 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>.
- Krijthe, Jesse H. 2015. *Rtsne: T-Distributed Stochastic Neighbor Embedding Using Barnes-Hut Implementation*. <https://github.com/jkrijthe/Rtsne>.
- 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.”

## 9 Code Appendix

```
library(quarto)
library(knitr)
library(tidyverse)
library(conflicted)
library(janitor)
library(ggtda)
# library(TDAvis)
library(patchwork)
library(gganimate)
library(ggforce)
library(simpextree)
library(gifski)
library(magick)
library(ripserr)
library(reshape2)
# remotes::install_github("maroulaslab/BayesTDA") Use this if package 'BayesTDA' is not available
library(BayesTDA)
library(TDAstats)
library(mvtnorm)
library(kableExtra)
library(plotly)
library(DiagrammeR)
library(transport)
library(TDA)
library(RColorBrewer)
library(Rtsne)
library(keras)
library(furrr)
library(yardstick)
library(caret)
library(imager)
conflicted::conflict_prefer("filter", "dplyr")
conflicted::conflict_prefer("select", "dplyr")
conflicted::conflicts_prefer(ggtda::geom_simplicial_complex)
conflicted::conflicts_prefer(plotly::layout)
conflicts_prefer(magrittr::set_names)
knitr::opts_chunk$set(
  comment = "#>",
  message = FALSE,
  warning = FALSE,
  cache = FALSE,
  echo = FALSE,
  tidy.opts = list(width.cutoff = 100),
  tidy = FALSE,
```

```

  fig.align = "center"
)
ggplot2::theme_set(ggplot2::theme_minimal())
ggplot2::theme_update(panel.grid.minor = ggplot2::element_blank())

#-----#
#-----
#-----Load & Preprocess Data-----
#-----

mnist <- readRDS(file = "mnist_dataset")

train <- mnist$train
test <- mnist$test

train_images <- train$images
train_labels <- as.factor(train$labels)

test_images <- test$images
test_labels <- as.factor(test$labels)

train_images <- train_images / 255
test_images <- test_images / 255

train_images_list <- lapply(1:nrow(train_images), function(i) {
  matrix(train_images[i, ], nrow = 28) |> t()
})

test_images_list <- lapply(1:nrow(test_images), function(i) {
  matrix(test_images[i, ], nrow = 28) |> t()
})

plot_digit <- \(image_list = train_images_list, image_index = NULL, image_df = NULL, melted = FALSE) {
  if(!melted){
    image_df <- melt(image_list[image_index])
    colnames(image_df) <- c("y", "x", "value")
  }
  ggplot(image_df, aes(x = x, y = y, fill = value)) +
    geom_raster() +
    scale_fill_gradient(low = "white", high = "black") +
    scale_y_reverse() +
    coord_equal() +
    theme_void() +
    theme(legend.position = "none")
}

# plot_digit(image_index = 8)

```

```

# paste0("Label: ", train$labels[8])

binarize_images <- function(images_list, threshold = 0.5) {
  lapply(images_list, function(mat) {
    ifelse(mat < threshold, 0, 1)
  })
}

train_images_binarized <- binarize_images(train_images_list)
test_images_binarized <- binarize_images(test_images_list)

# plot_digit(image_index = 8) + plot_digit(train_images_binarized, image_index = 8)

#-----#
\usetikzlibrary{
  positioning,
  arrows.meta,
  shapes.geometric,
  fit,
  calc
}

\begin{tikzpicture}[
  % Adjusted node distances for better spacing
  node distance = 1.2cm and 2cm,
  every node/.style={
    draw,
    thick,
    rounded corners,
    align=center,
    minimum height=1.3cm,
    font=\sffamily
  },
  data/.style={fill=green!20, text width=3cm},
  prior/.style={fill=yellow!30, text width=4cm},
  posterior/.style={fill=blue!20, text width=4cm},
  result/.style={fill=red!20, text width=3.5cm},
  process/.style={text width=4cm},
  arrow/.style={->, >=Stealth, thick},
  connector/.style={draw=none, font=\sffamily\Huge},
  % A dedicated style for labels on arrows (edges)
  edge_label/.style={draw=none, midway, fill=none, font=\sffamily}
]

% == Column 1 & 2: Data and PD Calculation ==
% Position nodes in the first two columns

```

```

\node[data] (train) {Train Data \\ (60,000 images)};
\node[process, right=of train] (calc_pd_train) {Calculate Train PDs \\ (for dim0 \& dim1)};

% Increased vertical distance for a clearer separation of train/test paths
\node[data, below=3.75cm of train] (test) {Test Data \\ (10,000 images)};
\node[process, right=of test] (calc_pd_test) {Calculate Test PDs \\ (for dim0 \& dim1)};

% == Column 3: Bayesian Model Training ==
% Position this block relative to the training data processing nodes
\node[process, right=of calc_pd_train] (likelihoods) {Likelihood Surfaces from Train PDs \\ (for
%\node[connector, right=of likelihoods] (update_op) {\$\\otimes\$};
\node[connector, right=of likelihoods] (update_op) {\$\\odot\$};
\node[prior, right=of update_op] (priors) {Uninformative Priors \\ (for digits 0-9)};
\node[posterior, below=of update_op] (posteriors) {Posterior Surfaces \\ (for digits 0-9)};

% Bounding box for the Bayesian update process
\node[draw, dashed, inner sep=0.4cm, fit=(priors) (likelihoods) (update_op) (posteriors), label={

% == Column 4: Classification ==
% Position the classification node vertically centered between its inputs for a balanced look
\node[process, below=of posteriors] (calc_dist) {Calculate Distances to all Posteriors \\ Distance
\node[result, below=of calc_dist] (classify) {Classify as \\ argmin(Distance)};

% == Arrows ==
% Connect nodes with clearer, non-overlapping paths
\draw[arrow] (train) -- (calc_pd_train);
\draw[arrow] (test) -- (calc_pd_test);

% Bayesian model flow
\draw[arrow] (calc_pd_train) -- (likelihoods);
\draw[arrow] (priors) |- (posteriors);
\draw[arrow] (likelihoods) |- (posteriors);

% Classification flow
% Use |- routing to different anchors (north west and south west) to keep lines clean
\draw[arrow] (posteriors) -- (calc_dist);
\draw[arrow] (calc_pd_test) -- (calc_dist);

% Arrow with a nicely placed label for the distance formula
\draw[arrow] (calc_dist) -- (classify);
    \node[edge_label, right=0.2cm] {Distance = \\ \$(1-\lambda)d_{0} + \lambda d_{1}\$};
    \node[edge_label, right=0.2cm] {Distance = \$d_{0} + d_{1}\$};

\end{tikzpicture}
# example persistence diagram

```

```

set.seed(5926720)

df <- tdaunif::sample_circle(n = 40L, sd = .05)
colnames(df) <- c("x", "y")
df |> ripser::vietoris_rips() |>
  ggplot() +
  ggtda::stat_persistence(
    aes(start = birth, end = death,
        color = factor(dimension), shape = factor(dimension)),
    size = 2.5,
    alpha = .75
  ) +
  geom_abline(intercept = 0, slope = 1, linetype = "dashed", alpha = 0.6) +
  scale_color_viridis_d(end = .7, name = "Dimension") +
  scale_shape_discrete(name = "Dimension") +
  labs(title = "Persistence Diagram") +
  theme_minimal()

image_index <- which(mnist$train$labels == 8)[100]
image_matrix <- matrix(mnist$train$images[image_index, ], nrow = 28, byrow = TRUE)

binary_image <- ifelse(image_matrix > 0.4, 1, 0)
img_c <- as.cimg(binary_image)
img_c_inv <- as.cimg(1 - binary_image)

coords <- as.matrix(expand.grid(y = 1:28, x = 1:28))

height_filt_td <- matrix(coords %*% c(0, 1), nrow = 28, byrow = TRUE) * binary_image
height_filt_lr <- matrix(coords %*% c(1, 0), nrow = 28, byrow = TRUE) * binary_image

radial_filt_c <- matrix(sqrt(rowSums(sweep(coords, 2, c(14, 14), "-")^2)), nrow = 28, byrow = TRUE)
radial_filt_b <- matrix(sqrt(rowSums(sweep(coords, 2, c(14, 28), "-")^2)), nrow = 28, byrow = TRUE)

dilation_filt <- as.matrix(distance_transform(img_c, 1))
erosion_filt <- as.matrix(distance_transform(img_c_inv, 1))
density_filt <- as.matrix(boxblur(img_c, 3))

plot_filtration <- function(data_matrix, title) {
  as.data.frame(data_matrix) |>
  mutate(y = row_number()) |>
  pivot_longer(-y, names_to = "x", values_to = "value") |>
  mutate(x = as.integer(gsub("V", "", x))) |>
  ggplot(aes(x = x, y = y, fill = value)) +
  geom_tile() +
  scale_y_reverse() +
  coord_equal() +

```

```

    scale_fill_viridis_c(option = "magma", direction = -1) +
    # scale_fill_gradient(low = "white", high = "black", na.value = "white") +
    theme_void() +
    labs(title = title) +
    theme(plot.title = element_text(hjust = 0.5, face = "bold"), legend.position = "none")
}

p_orig <- plot_filtration(image_matrix, "Original Grayscale") + scale_fill_gradient(low = "white"
p_bin <- plot_filtration(binary_image, "Binarized") + scale_fill_gradient(low = "white", high = "
p_height_td <- plot_filtration(height_filt_td, "Height (Top-Down)")
p_height_lr <- plot_filtration(height_filt_lr, "Height (Left-Right)")
p_radial_c <- plot_filtration(radial_filt_c, "Radial (Center)")
p_radial_b <- plot_filtration(radial_filt_b, "Radial (Bottom)")
p_dilation <- plot_filtration(dilation_filt, "Dilation")
p_erosion <- plot_filtration(erosion_filt, "Erosion")
p_density <- plot_filtration(density_filt, "Density (r=3)")

(p_orig + p_bin + p_dilation) /
(p_height_td + p_height_lr + p_erosion) /
(p_radial_c + p_radial_b + p_density)
#-----
#-----Distribution of training labels-----
#-----
#first run the qmd file
train_df <- train$images |> as.data.frame()
train_df$labels <- train$labels |> as.factor()

ggplot(train_df, aes(labels, fill = labels))+
  geom_bar()+
  labs(fill = "digit")

#-----
#-----Pixel intensity representation-----
#-----
# train_df_long <- train_df |>
#   pivot_longer(-labels, names_to = "covariate", values_to = "intensity")
# saveRDS(train_df_long, "train_df_long.rds")
train_df_long <- readRDS("train_df_long.rds")

ggplot(train_df_long) +
  geom_histogram(
    aes(intensity),
    bins = 30,
    fill = "#2C7FB8",
    color = "white",

```

```

    alpha = 0.8
  )
# ggplot(train_df_long) +
#   geom_histogram(
#     aes(intensity),
#     bins = 30,
#     fill = "#2C7FB8",
#     color = "white",
#     alpha = 0.8
#   ) +
#   scale_y_log10() +
#   labs(y = "Count (log scale)")

#-----
#-----distribution of mnist training-----
#-----

#normalized pixel values and apply t-SNE
# tsne_results <- Rtsne(
#   train_images,
#   dim = 2,
#   perplexity = 30,
#   max_iter = 1000
# )
# df_tsne <- tibble(Dim1 = tsne_results$Y[, 1],
#                   Dim2 = tsne_results$Y[, 2],
#                   digit = train_df$labels)
#
# saveRDS(df_tsne, "df_tsne")
df_tsne <- readRDS("df_tsne")
ggplot(df_tsne, aes(Dim1, Dim2, color = digit)) +
  geom_point(alpha = .25)
posterior_list_dim0 <- readRDS("btda/posterior_list_dim0.rds")
posterior_list_dim1 <- readRDS("btda/posterior_list_dim1.rds")
posterior_dim1_df <- posterior_list_dim1 |>
  set_names(0:9) |>
  list_rbind(names_to = "digit")
posterior_dim0_df <- posterior_list_dim0 |>
  set_names(0:9) |>
  list_rbind(names_to = "digit")

(ggplot(posterior_dim0_df, aes(x = birth, y = persistence, fill = intensity)) +
  geom_tile() +
  facet_wrap(~digit, ncol = 5) +
  scale_fill_viridis_c(option = "magma") +
  labs(
    title = "Posterior Densities for Digit Components (Dimension 0)",

```



```

    x = "Birth",
    y = "Persistence",
    fill = "Intensity"
) +
theme_minimal() +
theme(strip.text = element_text(size = 12, face = "bold"))) / (ggplot(posterior_dim1_df, aes(x
geom_tile() +
facet_wrap(~digit, ncol = 5) +
scale_fill_viridis_c(option = "magma") +
labs(
  title = "Posterior Densities for Digit Loops (Dimension 1)",
  x = "Birth",
  y = "Persistence",
  fill = "Intensity"
) +
theme_minimal() +
theme(strip.text = element_text(size = 12, face = "bold"))) +
plot_layout(guides = "collect") +
plot_annotation(
  title = "Posterior Densities by Dimension",
  theme = theme(plot.title = element_text(hjust = 0.5, face = "bold", size = 14))
)
# --- Load Data ---
nn_dropout_conf_mat <- readRDS("nn_dropout_conf_mat")
nn_ridge_conf_mat <- readRDS("nn_ridge_conf_mat")
nn_lasso_conf_mat <- readRDS("nn_lasso_conf_mat")
mlogit_conf_mat <- readRDS("mlogit_conf_mat")

# --- Prepare data frames for all plots ---
df_dropout <- as.data.frame.matrix(nn_dropout_conf_mat$table) %>%
  mutate(true_label = fct_rev(rownames(.))) |>
  pivot_longer(cols = -true_label, names_to = "predicted_label", values_to = "count")

df_ridge <- as.data.frame.matrix(nn_ridge_conf_mat$table) %>%
  mutate(true_label = fct_rev(rownames(.))) |>
  pivot_longer(cols = -true_label, names_to = "predicted_label", values_to = "count")

df_lasso <- as.data.frame.matrix(nn_lasso_conf_mat$table) %>%
  mutate(true_label = fct_rev(rownames(.))) |>
  pivot_longer(cols = -true_label, names_to = "predicted_label", values_to = "count")

df_mlogit <- as.data.frame.matrix(mlogit_conf_mat$table) %>%
  mutate(true_label = fct_rev(rownames(.))) |>
  pivot_longer(cols = -true_label, names_to = "predicted_label", values_to = "count")

# --- Determine the global scale limits for the fill color ---

```

```

global_limits <- range(
  bind_rows(df_dropout, df_ridge, df_lasso, df_mlogit) |> pull(count)
)

# --- Define a plotting function to reduce code repetition ---
create_heatmap <- function(data, title, limits) {
  data |>
    mutate(text_color = if_else(count < (0.6 * max(limits)), "black", "white")) |>
    ggplot(aes(x = predicted_label, y = true_label, fill = count)) +
    geom_tile(color = "gray50", linewidth = 0.5) +
    geom_text(aes(label = count, color = text_color), size = 2) +
    scale_fill_viridis_c(
      option = "magma",
      direction = -1,
      limits = limits # Apply global limits here
    ) +
    scale_color_manual(values = c("black" = "black", "white" = "white"), guide = "none") +
    coord_fixed() +
    theme_minimal(base_size = 12) +
    labs(
      x = "Predicted Label",
      y = "True Label",
      fill = "Count",
      title = title
    ) +
    theme(
      axis.text.x = element_text(angle = 45, hjust = 1, face = "bold"),
      axis.text.y = element_text(face = "bold"),
      panel.grid = element_blank()
    )
}

# --- Create each plot using the function ---
p_dropout <- create_heatmap(df_dropout, "Neural Net (Dropout)", global_limits)
p_ridge <- create_heatmap(df_ridge, "Neural Net (Ridge)", global_limits)
p_lasso <- create_heatmap(df_lasso, "Neural Net (Lasso)", global_limits)
p_mlogit <- create_heatmap(df_mlogit, "Multinomial Logistic", global_limits)

# --- Combine Plots with Patchwork, collecting guides ---
(p_dropout + p_ridge) / (p_lasso + p_mlogit) +
  plot_annotation(
    theme = theme(plot.title = element_text(hjust = 0.5, size = 18, face = "bold"))
  ) &
  theme(legend.position = "none")

```

```

full_tda_01_results <- readRDS("btda/full_tda_01_results.rds")

conf_matrix <- table(
  true_label = full_tda_01_results$true_label,
  predicted_label = full_tda_01_results$predicted_label
)

as.data.frame.matrix(conf_matrix) %>%
  mutate(true_label = fct_rev(rownames(.))) |>
  pivot_longer(
    cols = -true_label,
    names_to = "predicted_label",
    values_to = "count"
  ) |>
  mutate(text_color = if_else(count < 450, "black", "white")) |>
  ggplot(aes(x = predicted_label, y = true_label, fill = count)) +
  geom_tile(color = "gray50", linewidth = 0.5) +
  geom_text(aes(label = count, color = text_color), size = 3.5) +
  scale_fill_viridis_c(option = "magma", direction = -1, guide = "none") +
  scale_color_manual(values = c("black" = "black", "white" = "white"), guide = "none") +
  coord_fixed() +
  theme_minimal(base_size = 12) +
  labs(
    x = "Predicted Label",
    y = "True Label",
    fill = "Count"
  ) +
  theme(
    axis.text.x = element_text(angle = 45, hjust = 1, face = "bold"),
    axis.text.y = element_text(face = "bold"),
    panel.grid = element_blank()
  )

# accuracy <- mean(full_results$true_label == full_results$predicted_label, na.rm = TRUE)

multi_metrics <- metric_set(
  yardstick::accuracy,
  yardstick::precision,
  yardstick::recall,
  yardstick::f_meas
)

results_for_metrics <- full_tda_01_results |>
  select(true_label, predicted_label) |>
  mutate(
    true_label = factor(true_label, levels = 0:9),

```

```

    predicted_label = factor(predicted_label, levels = 0:9)
  )

proposed_accuracy <- multi_metrics(results_for_metrics, truth = true_label, estimate = predicted_label,
  janitor::clean_names() |>
  filter(metric == "accuracy") |> pull(estimate)

# ml_results <- tibble(
#   method = c("multinomial", "dropout nn", "ridge nn", "lasso no"),
#   accuracy = c(mlogit_acc, nn_dropout_accu, nn_ridge_accu, nn_lasso_accu)
# )
# saveRDS(ml_results, "ml_results")

ml_results <- readRDS("ml_results") #still a tibble so easy to edit
ml_results <- ml_results %>%
  mutate(method = str_replace(method, "lasso no", "lasso nn"))
ml_results |>
  add_row(method = "proposed", accuracy = proposed_accuracy) |>
  kable(digits = 4)
all_features_train_data <- readRDS("rdss/all_features_train_data.rds")
all_features_test_data <- readRDS("rdss/all_features_test_data.rds")
all_features_rf_model <- readRDS("rdss/all_features_rf_model.rds")
all_features_svm_model <- readRDS("rdss/all_features_svm_model.rds")
model_comparison <- resamples(list(RandomForest = all_features_rf_model, SVM = all_features_svm_model))
# dotplot(model_comparison,
#   scales = list(x = list(rot = 45)))
tidy_comparison <- model_comparison$values |>
  pivot_longer(
    cols = -Resample,
    names_to = "key",
    values_to = "value"
  ) |>
  separate(key, into = c("Model", "Metric"), sep = "~")

ggplot(tidy_comparison, aes(x = value, y = Model, color = Model)) +
  geom_point(size = 4, alpha = 0.2) +
  scale_color_brewer(palette = "Set2") +
  scale_x_continuous(limits = c(NA, 1.0)) +
  facet_wrap(~ Metric, ncol = 4) +
  labs(
    x = "Metric Value",
    y = NULL,
    color = "Model Type"
  ) +
  theme_light(base_size = 12) +
  theme(

```

```

    plot.title = element_text(face = "bold", size = 18, margin = margin(b = 5)),
    plot.subtitle = element_text(size = 13, margin = margin(b = 15)),
    strip.text = element_text(face = "bold", size = 10, color = "white"),
    strip.background = element_rect(fill = "#525252", color = "white"),
    legend.position = "bottom",
    panel.grid.major.y = element_blank(),
    panel.grid.minor.x = element_blank()
  ) +
  guides(color = guide_legend(override.aes = list(alpha = 1)))
test_indices <- 1:10000
predictions <- predict(all_features_svm_model, newdata = all_features_test_data)
results <- tibble(
  true_label = factor(mnist$test$labels[test_indices]),
  predicted_label = predictions
)

conf_matrix_tda_full <- table(
  true_label = results$true_label,
  predicted_label = results$predicted_label
)

as.data.frame.matrix(conf_matrix_tda_full) %>%
  mutate(true_label = fct_rev(rownames(.))) |>
  pivot_longer(
    cols = -true_label,
    names_to = "predicted_label",
    values_to = "count"
  ) |>

  mutate(text_color = if_else(count < 450, "black", "white")) |>
  ggplot(aes(x = predicted_label, y = true_label, fill = count)) +
  geom_tile(color = "gray50", linewidth = 0.5) +

  geom_text(aes(label = count, color = text_color), size = 3.5, fontface = "bold") +
  scale_fill_viridis_c(option = "magma", direction = -1, guide = "none") +

  scale_color_manual(values = c("black" = "black", "white" = "white"), guide = "none") +
  coord_fixed() +
  theme_minimal(base_size = 12) +
  labs(
    x = "Predicted Label",
    y = "True Label",
    fill = "Count"
  ) +
  theme(
    axis.text.x = element_text(angle = 45, hjust = 1, face = "bold"),

```

```

    axis.text.y = element_text(face = "bold"),
    panel.grid = element_blank()
  )

# final_accuracy <- mean(results$true_label == results$predicted_label, na.rm = TRUE)
# print(paste("Final Test Set Accuracy (Random Forest):", scales::percent(final_accuracy, accurac
rf_model_importance <- varImp(all_features_rf_model, scale = TRUE)
svm_model_importance <- filterVarImp(
  x = all_features_train_data[, -which(names(all_features_train_data) == "label")],
  y = all_features_train_data$label
)

rf_plot <- rf_model_importance$importance |>
  as.data.frame() %>%
  mutate(variable = rownames(.)) |>
  pivot_longer(
    cols = -variable,
    names_to = "class",
    values_to = "importance"
  ) |>
  mutate(
    importance = (importance - min(importance)) / (max(importance) - min(importance))
  ) |>
  mutate(variable = fct_reorder(variable, importance, .fun = max, .desc = TRUE)) |>
  ggplot(aes(x = class, y = variable, fill = importance)) +
  geom_tile(color = "white", linewidth = 0.4) +
  scale_fill_viridis_c(
    option = "magma",
    direction = -1,
    limits = c(0, 1)
  ) +
  labs(
    x = "Class",
    y = "Variable",
    fill = "Importance",
    title = "Random Forest Variable Importance"
  ) +
  theme_minimal(base_size = 12) +
  theme(axis.text.x = element_text(angle = 45, hjust = 1))

svm_plot <- as.data.frame(svm_model_importance) %>%
  mutate(variable = rownames(.)) |>
  pivot_longer(
    cols = -variable,
    names_to = "class",
    values_to = "importance"
  )

```

```

) |>
mutate(
  importance = (importance - min(importance)) / (max(importance) - min(importance))
) |>
mutate(
  variable = fct_reorder(variable, importance, .fun = max, .desc = TRUE),
  class = gsub("X", "", class)
) |>
ggplot(aes(x = class, y = variable, fill = importance)) +
geom_tile(color = "white", linewidth = 0.4) +
scale_fill_viridis_c(option = "magma", direction = -1) +
labs(
  x = "Class",
  y = NULL,
  fill = "Importance",
  title = "SVM Variable Importance"
) +
theme_minimal(base_size = 12) +
theme(
  axis.text.x = element_text(angle = 45, hjust = 1),
  axis.text.y = element_blank()
)

rf_plot + svm_plot +
plot_layout(guides = 'collect')

```