

ML Kit: A Machine Learning Library for Rust

Owen Wetherbee (ocw6), Ethan Ma (em834), Sylvan Martin (slm338)

Keywords: Machine learning, Gradient Descent, PCA, Rust

Application Setting

Rust is a relatively new programming language with a dearth of machine learning infrastructure. We aimed to create a library for Rust programmers to make machine learning convenient, much like NumPy or TensorFlow in Python.

Project Description

Initially, we wanted to create a comprehensive machine learning library for Rust, which would use Rust's safety and speed to implement many algorithms used in data science. We had the original (lofty) goal of being able to run a diffusion model by the end of the semester, or be able to do anything that NumPy/SciKitLearn could do. As we began implementation we recognized that we lacked the time and resources to implement the sheer amount of algorithms we set out to. So, we shifted our focus to what we believe are some of the core algorithms in the field of Machine Learning.

In the end, we created the pure-Rust library `ml_kit`, which implements, from scratch, the following:

- Neural Network based learning, consisting of
 - the basic neural network model with user-defined network shape and activation functions for each layer
 - functionality for handling large datasets for training and testing
 - a stochastic gradient descent trainer, with whatever batch size and epochs a user may want
 - various “gradient update schedules,” such as fixed learning rates, time-decay learning rates, AdaGrad, etc.
 - Convolutional Neural Networks (Owen/Ethan talk more on this?)

- Principle Component Analysis, consisting of
 - an implementation of Singular Value Decomposition (SVD),
 - using SVD to obtain a k -dimensional plane of best fit for a set of points in \mathbb{R}^n ,
 - compressing images (or any data) using SVD by truncating low-variance dimensions

Neural Networks

Neural network based learning (i.e. deep learning) is a type of supervised learning that involves using black-box models which depend on a large number of tunable parameters to classify complicated inputs. The models are trained by iteratively evaluating their performance on training examples and updating their parameters accordingly, usually through gradient descent. We call these black-box models ‘networks’ because they are generally conceptualized as a series of composed functions, or ‘layers’, where each layer linearly maps some input vector (or, more generally, tensor) to a possibly differently-sized output vector, and then applies a simple element-wise non-linear ‘activation function’. This sequence of highly inter-connected linear functions interspersed with element-wise non-linearities enables the model to identify complex features only using parameters that have relatively simple and computable gradients, ensuring efficient training.

Fully-connected neural networks

The traditional neutral network (NN) is a series of fully-connected layers, meaning each input element of a layer has a parameter through which it can linearly affect each output element. In particular, the function associated with each layer is of the form ¹

$$\vec{v}_{\text{out}} = f(\vec{v}_{\text{in}}) = \sigma \left(\mathbf{W} \vec{v}_{\text{in}} + \vec{\mathbf{b}} \right) ,$$

where \mathbf{W} is a matrix of ‘weight’ parameters, $\vec{\mathbf{b}}$ is a vector of ‘bias’ parameters, and σ is the non-linear activation function. The derivative of the output vector with respect to the input vector and these parameters can then be simply calculated as

$$\frac{\partial v_{\text{out},i}}{\partial W_{i,j}} = v_{\text{in},j} \cdot \sigma'(\vec{x}) , \quad \frac{\partial v_{\text{out},i}}{\partial b_i} = \sigma'(\vec{x}) , \text{ and} \quad \frac{\partial v_{\text{out},i}}{\partial v_{\text{in},j}} = W_{i,j} \cdot \sigma'(\vec{x}) ,$$

where $\vec{x} = \mathbf{W} \vec{v}_{\text{in}} + \vec{\mathbf{b}}$ and all unspecified derivatives are zero. Thus, the output of the entire neural-network model can be simply computed as the composition of each

1. Michael A. Nielsen, “Neural Networks and Deep Learning,” 2015, <http://neuralnetworksanddeeplearning.com>.

layer function, and the gradient of the loss with respect to the model parameters (to use for parameter updates) can be determined by iteratively applying the chain rule to the above derivatives.

Convolutional neural networks

Fully-connected neural networks are very general, but they do not take advantage of any potential structure of the input to minimize the number of required parameters or improve efficiency or performance. Enter, convolutional neural networks (CNNs). CNNs, primarily designed for image recognition and classification, use the spatial structure of an image to allow for shared weights and biases, reducing the number of required parameters. In particular, CNNs scan a series of filters across the image, each of which conceptually is intended to identify a certain local feature of the image (e.g. identify straight or curved lines). This scanning process is specifically done by convolving each filter with the image, hence the ‘convolutional’ in CNN. As described in ², this convolution step can be implemented efficiently via matrix multiplication by first reshaping the input matrices and filters. Although we do not give the explicit formulas here (see ³), the derivatives of the output of a convolution layer with respect to the inputs and filters can also be computed via a convolution.

CNNs often also contain pooling layers, which further reduce the dimensionality of the data. These layers scan a ‘window’ along the input matrix, combining all the elements in each window into a single output element. This combination process can be done in a number of ways, including taking the maximum, average, or sum of the window elements, which are referred to as max pooling, average pooling, and sum pooling, respectively. These pooling layers do not have any parameters, and their outputs are simply related to the elements of each window, so the derivatives of this layer are straightforward to compute.

A CNN usually involves a sequence of several of these convolution and pooling layers. However, once the dimensionality of the input data has been sufficiently reduced, the final classification is usually performed by a fully-connected neural network, which constitutes the last few layers of the CNN. As with the fully-connected NN, training proceeds by iteratively passing the training inputs through each layer to get the model output, and then moving back through the layers via the chain rule to compute the loss gradient.

2. December 2019, https://leonardoaraujosantos.gitbook.io/artificial-intelligence/machine-learning/deep_learning/convolution_layer/making_faster.

3. Pavithra Solai, *Convolutions and Backpropagations*, April 2018, <https://pavisj.medium.com/convolutions-and-backpropagations-46026a8f5d2c>.

Stochastic gradient descent

The training of these neural networks is often powered by stochastic gradient descent (SGD). SGD is a scheme for updating the parameters of a network using the gradient of the loss for small sampled batches of training examples. In particular, to perform SGD, a set of training examples is provided, along with a specification of batch size and number of epochs. In each epoch, the whole set of training examples is sampled into different smaller sets specified by the batch size. The network parameters are then updated based on the loss incurred by each batch of training examples, rather than the entire training set. This greatly reduces the computational cost in each update step, while not changing the performance of the model in expectation. Each batch step specifically involves iterating through each training example in the batch. Each training input is passed forward through the layers of the network to compute the output and loss, which is then passed back through the layers to compute the derivative of the loss with respect to each parameter via the chain rule. These derivatives are summed for each example from the batch, the resulting gradient is normalized, and then each parameter p is updated as according to

$$p \longmapsto p - \alpha \frac{\partial L}{\partial p},$$

where $\frac{\partial L}{\partial p}$ is the normalized derivative of the loss L with respect to parameter p and α is the step size, or ‘learning rate’. This process is repeated for each batch in the training set, yielding an epoch. Many epochs are usually required to train a network to achieve reasonable classification accuracies.

What we implemented

Principle Component Analysis and SVD

Singular Value Decomposition (SVD) is a factorization of a matrix $A \in \mathbb{R}^{m \times n}$ into

$$A = U\Sigma V^\top$$

where $U \in \mathbb{R}^{m \times m}$ and $V \in \mathbb{R}^{n \times n}$ are orthogonal matrices with columns $\vec{u}_1, \dots, \vec{u}_m$ and $\vec{v}_1, \dots, \vec{v}_n$ respectively, and Σ is a diagonal matrix of the form (assuming $n \leq m$)

$$\Sigma = \begin{bmatrix} \sigma_1 & & & \\ & \sigma_2 & & \\ & & \ddots & \\ & & & \sigma_n \end{bmatrix} \in \mathbb{R}^{m \times n}$$

where $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_n \geq 0$ are the *singular values* of A . This allows us to write A as the linear combination

$$A = \sum_{i=1}^n \sigma_i \vec{u}_i \vec{v}_i^\top$$

Note that because the singular values are sorted in decreasing order, we can effectively “save data” in representing A by truncating all but the first $r \leq n$ singular vectors, as the last items in the sum do not contribute as much to the overall product. The immediate application is identifying the most significant axes of correlation in data, allowing dimensionality reduction of datasets by re-writing each item in the basis $\vec{v}_1, \dots, \vec{v}_r$. Commonly, the “data” of concern is written into the rows of A .

What we implemented

Our library implements the Golub-Kahan SVD algorithm (described in “Matrix Computations,” by Golub and van Loan.⁴) which begins by bi-diagonalizing A , then performing SVD on the bidiagonalization, as the numerical stability and performance are better. Once we were able to compute the SVD of any matrix, we could use SVD as a subroutine in other useful techniques.

The (subjectively) coolest application of SVD we implemented is image compression. By taking an $m \times n$ image and writing it as a 4-tuple of matrices (R, G, B, A) representing the red, green, blue, and alpha channels of the pixels, we can perform SVD on each color channel, store only the SVD representation after truncating insignificant singular vectors, and then de-compress the image later by computing $U\Sigma V^\top$ for each color channel. In practice, one can discard roughly half the singular values and still obtain a recognizable image. Examples of this and discussion of runtime and compression rates will be left to the evaluation section.

A common goal in statistics is finding the “line of best fit” of a set of points, or in higher dimensions, a k -dimensional plane of best fit. For a cluster of data centered at the origin, the first singular vector, \vec{v}_1 , is the line of best fit. This is because SVD is equivalently defined as an optimization procedure where

$$\vec{v}_1 = \underset{\|x\|=1}{\operatorname{argmax}} \|Ax\|$$

Since \vec{v}_1 is maximizing the sum of squares of the inner products with all the data points, it is minimizing the sum of squared distances from each data point to the line spanned by \vec{v}_1 . More generally, taking the first $\vec{v}_1, \dots, \vec{v}_k$ vectors, we obtain a basis for a k -dimensional “plane of best fit.”

4. Gene H. Golub and Charles F. van Loan, *Matrix Computations*, Fourth (JHU Press, 2013), ISBN: 1421407949 9781421407944, <http://www.cs.cornell.edu/cv/GVL4/golubandvanloan.htm>.

We implement a procedure which takes a data matrix $D \in \mathbb{R}^{n \times m}$ whose columns are each \mathbb{R}^n data points, and returns $\vec{\mu} \in \mathbb{R}^n$ and $V \in \mathbb{R}^{n \times k}$ such that the plane defined by

$$\left\{ \vec{\mu} + \sum_{i=1}^k \alpha_i \vec{v}_i \mid \alpha_i \in \mathbb{R} \right\}$$

is the plane of best fit for the data. Examples of this will be left to the evaluations section.

Relationship to Other Work

Linear algebra is the foundation to machine learning, so we needed to use a good linear algebra library. Sylvan had previously spent winter break working on `matrix_kit`, which is a pure-Rust linear algebra package that implemented incredibly basic matrix-vector operations. We continued developing this library in parallel with `ml_kit` over the semester as we recognized more features that were needed from the linear algebra library. So, the sum of our work for the course can be thought of as the entirety of `ml_kit`, as well as significant improvement to the functionality of `matrix_kit`.

The `matrix_kit` library can be found on GitHub at https://github.com/SylvanM/matrix_kit.

Evaluation

Neural Networks

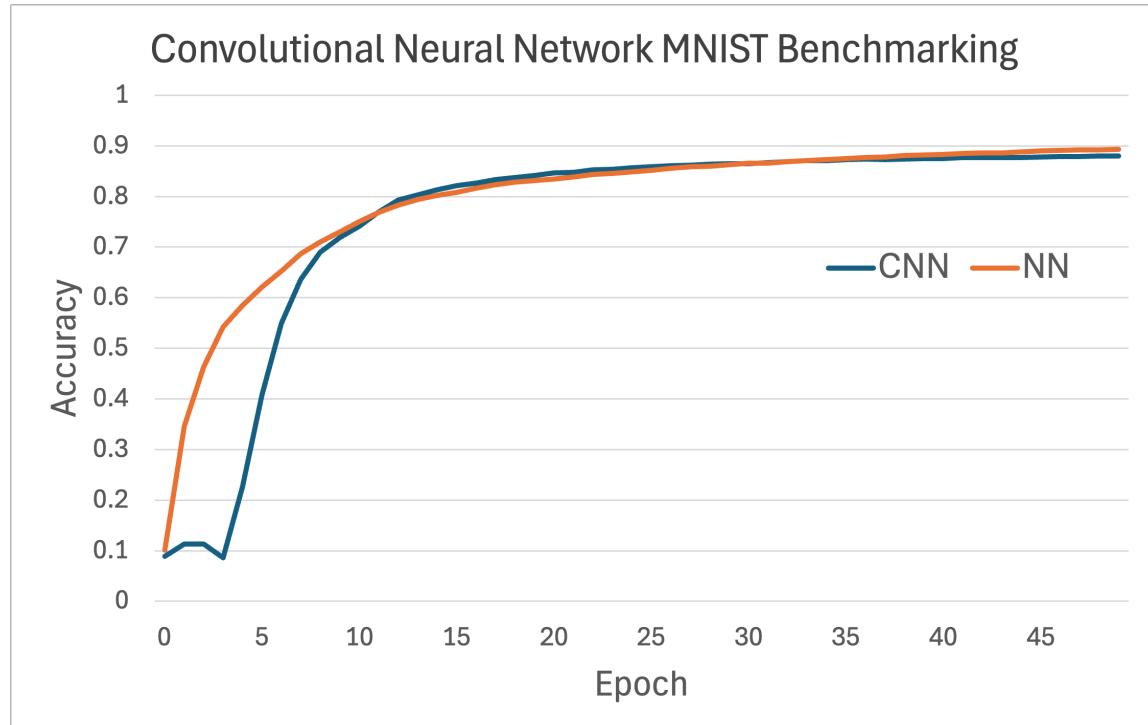
We performed unit tests of various functions of our fully-connected neural network (NN), convolutional neural network (CNN), and stochastic gradient descent (SGD) implementations to ensure their correctness. Apart from this, the primary way we verified these implementations in a more integrated way was by testing toy network models on the MNIST database of handwritten digits⁵. After parsing the MNIST data into rust and creating data items for easy integration with our SGD trainer methods, we constructed toy NN and CNN models to train on this data.

Each MNIST input consisted of a 28×28 pixel image, and each output consisted of one of the ten labels $\{0, 1, \dots, 9\}$. We therefore constructed the fully-connected NN to have 3 layers, with input and output sizes of $(784, 16)$, $(16, 16)$, and $(16, 10)$, respectively, thus taking the $28 \times 28 = 784$ pixels to one of two possible label probabilities. For the CNN, we used the following layer structure: a convolutional layer

5. <https://www.kaggle.com/datasets/hojjatk/mnist-dataset>.

with four 5×5 filters; a max pooling layer with a 2×2 window with stride 2; a convolutional layer with two $5 \times 5 \times 4$ filters; another max pooling layer with a 2×2 window with stride 2; and finally two fully-connected layers with input and output sizes of (32, 16) and (16, 10), respectively. Both convolutional layers had zero padding and stride 1, and all the layers used a sigmoid activation function.

We trained these two networks using the MNIST data and the SGD trainer methods, with the learning rate set to 0.05, a squared loss function, and a batch size of 32. The resulting accuracy after each epoch is shown in the plot below.



As can be seen, the two networks have a very similar performance curve, both achieving a classification accuracy of about 85% after 20 epochs, continuing more gradually up to 90% by the 50-th epoch. Note that the CNN is specified by 1,004 total parameters while the fully-connected NN is specified by 13,002 total parameters. Despite this, the two networks achieve a similar performance, demonstrating the advantage of the CNN. However, it is important to note that the CNN took considerably longer to train for each epoch. This is likely due to the overhead in the convolutional layers, stemming from some inefficient implementations in both CNN and matrix kit.

Regardless, this MNIST benchmarking shows that our implementations of NN, CNN, and SGD indeed yield trainable networks which can achieve good classification accuracies. While not guaranteeing the correctness or efficiency of our implementations, this serves as a good integrated test to show that our implementations are behaving

as intended.

Principle Component Analysis

In this section, we discuss the performance of our SVD-based algorithms in terms of their speed and correctness.

Singular Value Decomposition

Correctness

In numerical methods such as SVD, the accumulation of floating point error prevents us from measuring exact correctness. Instead, we define some small ε , and consider two matrices $A \approx_\varepsilon B$ to be “basically equal” if for all i, j , $|a_{ij} - b_{ij}| < \varepsilon$. When testing on large matrices, we chose $\varepsilon = 1 \times 10^{-5}$ to be our goal, though in practice on smaller matrices ($m, n \leq 50$) we attain correctness up to $\varepsilon = 1 \times 10^{-9}$.

To test correctness of our SVD implementation, we start with a randomly generated matrix $A \in \mathbb{R}^{m \times n}$ where each entry is sampled from the normal distribution $\mathcal{N}(0, 1)$, and compute its SVD into U , V , and Σ . We then check that U and V are orthogonal by testing $U^\top U \approx_\varepsilon I_m$ and $V^\top V \approx_\varepsilon I_n$. We then check that $A \approx_\varepsilon U\Sigma V^\top$.

Sometimes, our SVD implemtion does not converge on a particular matrix, and so we abort the execution. The main loop of the Golub-Kahan algorithm should in theory take $\mathcal{O}(n)$ steps, so we introduced a check that if the main loop runs more than $30n$ times, we declare the instance a failure, and return empty 0×0 matrices.

With m and n sampled randomly in the range 2 through 20, we ran the above testing procedure 10,000 times and found no correctness failures, though we observe a divergence rate of about 2.49%. That is, on small matrices, our algorithm failed to converge on 249 out of the 10,000 random matrices.

The issue of convergence seems to become better with larger matrices. When we sample m and n randomly from the range 2 through 200, we found that out of 300 trials, only 1 instance failed to converge, a rate of .33%.

For matrices with $200 \leq m, n \leq 300$, we found no matrices with failed to converge out of 30 trials.

Runtime Performance

Our SVD implementation is absurdly slow. It became impractical to test at a large scale for matrices with more than $500 \times 500 = 250,000$ entries. This is unfortunate,

as most real-world applications would require much larger matrices. When we first tried image compression, we used a 4K image taken by an iPhone, and running SVD on each color channel took so long that even when let to run for 12 hours overnight, it didn't terminate. It is hard to know if this is because it genuinely just took so long, or if it fell into the small fraction of cases that caused divergence. Because of the decaying rate of divergence for larger dimensions, we believe it is the former. We also tried running SVD on matrices of comparable size (with m and n in the low thousands) and were unable to get it to terminate after several hours.

This is almost certainly due to the fact that at this stage of development we did *not* fully optimize our Golub-Kahan implementation. The first step, which relies on bi-diagonalizing the matrix, uses a sequence of Householder reflections, which if implemented efficiently can take linear time as they only affect two rows of a matrix. Further optimization would include only applying the reflection to a sub-matrix of A , and not the entire matrix. Instead, we applied the householder reflection matrix to the entire matrix A , taking $\mathcal{O}(n^3)$ time for each step. This easily balloons the runtime, since this is done inside an inner loop. At first, we tried implementing the fully optimized Golub-Kahan algorithm, but were unable to get it to work correctly. So, we focused on just trying to get a *correct* implementation instead of premature optimization. The natural next steps after this project would be to get this running faster!

Image Compression

As discussed above, we were unable to run our image compression on large images, so we started with low-resolution images. We ran our image compression algorithm on two images, one of some sheep⁶, and one of a cat.

When a matrix A is decomposed into U, Σ, V , and all but the first k singular values are ditched, we are left with three matrices of sizes $m \times k$, k and $n \times k$ (we need only store the diagonal of Σ). So, the number of floating point values we must store as a fraction of the amount we would need to fully represent A is what we call the “compression ratio,” computed as

$$r = \frac{r + rm + rn}{mn}$$

For small matrices and large values of r , this doesn't save much. However, this scales incredibly well for large matrices, as we can pick modestly small values of r and end up with impressive compression ratios while still being able to recognize the original image. A small caveat here is that typically a single pixel's channel for one color only uses one byte to represent it, whereas we are using 64-bit floating points, so one

6. The sheep are in the teaching barn by the Vet School. I highly recommend that you go visit and pet them! It's free and you can just walk right up to them!

could argue that we need to include a factor of 8 in our compression ratio. However, our compression would work just as well if each color channel was represented as a floating point value between 0 and 1 (since this is what we convert the image to anyway when we begin compression), so we think it is fair to compare the sizes when thinking of the original channel as being a decimal between 0 and 1 as opposed to an integer between 0 and 255.

Sheep

We used a 240×320 pixel image of some sheep laying in a barn, and truncated the singular values up to the first $k \in \{100, 50, 10\}$. Running SVD to decompose the image took **1m 15s**.



Original



$k = 100$, 73.0% compression



$k = 50$, 36.5% compression



$k = 10$, 7.3% compression

Panini (the cat)

On a 640×480 image of Panini, SVD took **20m 50s** to finish. Here are the results using the first $k \in \{150, 50, 7\}$ singular values.



Original



$k = 150$, 54.7% compression



$k = 50$, 18.2% compression



$k = 7$, 2.6% compression

Planes of Best Fit

References

December 2019. [https://leonardoaraujosantos.gitbook.io/artificial-intelligence/machine_learning/deep_learning/convolution_layer/making_faster](https://leonardoaraujosantos.gitbook.io/artificial-intelligence/machine-learning/deep-learning/convolution_layer/making_faster).

<https://www.kaggle.com/datasets/hojjatk/mnist-dataset>.

Golub, Gene H., and Charles F. van Loan. *Matrix Computations*. Fourth. JHU Press, 2013. ISBN: 1421407949 9781421407944. <http://www.cs.cornell.edu/cv/GVL4/golubandvanloan.htm>.

Nielsen, Michael A. “Neural Networks and Deep Learning,” 2015. <http://neuralnetworksanddeeplearning.com>.

Solai, Pavithra. *Convolutions and Backpropagations*, April 2018. <https://pavisj.mediump.com/convolutions-and-backpropagations-46026a8f5d2c>.