

# Feature Selection and Image Classification

Landon Carter and Courtney Guo

December 2015

## 1 Introduction

In this work, we wish to implement supervised image clustering and classification which combines a feature extraction step with a clustering step. For the feature extraction stage, we will use two different methods: autoencoders and convolutional neural networks. For the clustering stage, we will use softmax layers in the respective neural nets, SVM, and k-means. We will then show the results of these algorithms on the CIFAR-10 image classification dataset.

## 2 Feature Selection

### 2.1 Autoencoder

We used `MATLAB`'s built in autoencoder implementation on the CIFAR-10 dataset. An autoencoder is a neural network used for learning efficient codings: it is trained such that its output matches its input, but it has a hidden layer in which the number of hidden nodes is less than the number of input nodes. Therefore, the autoencoder will learn to compress the data efficiently, and then decode the data to match the original input. The hidden layer represents the encoding scheme, so we can use the activations of the hidden layer as the features. Training the autoencoder is the same as training a neural network: for each data point we can perform back propagation to calculate the gradients of the error with respect to the weights, and then update the weights based on that, until the error converges or we reach the maximum number of iterations. In particular on the autoencoder, we do not need any labels - we simply try to reconstruct the input. In some implementations, the decoding weights can be inverted from the encoding weights such that the autoencoder is symmetric, though ours were not.

Autoencoders can also be stacked to produce deep neural nets, which can be significantly more effective than one-layer autoencoders, since they can learn higher-order features more easily. One nice feature is that each layer of the deep autoencoder can be trained independently relatively quickly to give approximately-optimal results. The entire deep neural net can then undergo full backpropagation to converge quickly to final optimal values. This method of training is significantly faster than traditional full backpropagation, since each iteration of optimizing a layer is much faster than full backpropagation, and the final full backpropagation step undergoes relatively few iterations since all weights are already close to optimal values.

### 2.2 Convolutional Neural Network

Convolutional neural networks (ConvNet's, or CNN's) are a particular type of neural net which is useful in image analysis and recognition. The distinguishing factor of CNN's are the way they can account for proximity in data, taking "convolutions" of small contiguous subsets of the data. This

is of particular importance in image processing, where large amounts of information and context can be gained by considering neighboring pixels. The CNN we constructed was based off a standard LeNet topology, consisting of 5 distinct blocks: 3 blocks of convolutional layer, Rectified Linear Unit (ReLU) activation layer, and pooling layer. After these 3 blocks (9 total layers), we had a final 4 layers which computed another convolution, ReLU, convolution, and then softmax in order to do a final bit of convolution before extracting labels. When paired with clustering methods, the output of the final convolution layer was sent to  $k$ -means and SVM instead of the softmax layer.

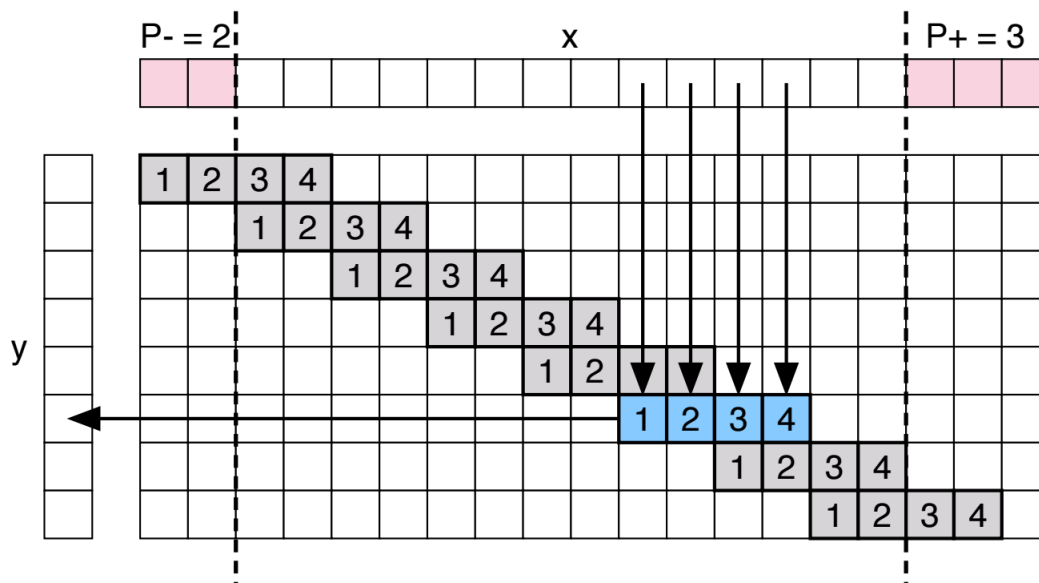
After constructing the CNN and initializing it with random values, the CNN is trained via stochastic gradient descent (SGD).

### 2.2.1 Convolutional Layer

The convolution layer is the most important part of the CNN, and is what distinguishes it most distinctly from other neural network topologies. Each convolutional neuron takes some small contiguous subset of the input layer and computes

$$y_{i''j''d''} = \sum_{i'=1}^{H'} \sum_{j'=1}^{W'} \sum_{d'=1}^D f_{i'j'd} \times x_{i''+i'-1,j''+j'-1,d',d''}$$

In other words, for the coordinate  $(i'', j'', d'')$ , the convolution is determined by applying some filter  $f_{i'j'd}$  to each of the inputs within  $H', W', D$  of the coordinate of interest (where the original input size is  $\mathbf{x} \in \mathbb{R}^{H \times W \times D}$ ). In this formulation, only inputs to the right and below the desired coordinate are considered (all depth dimensions - usually corresponding to different color channels - are used), though this does not affect the results - the output nominally for coordinate  $(i'', j'', d'')$  really corresponds to the output for coordinate  $(i'' + H'/2, j'' + W'/2, d'')$  in the traditional formulation, which considers inputs in the range  $H'/2$  to the left,  $H'/2$  to the right, and so on. The input is padded with a reasonable value such that edge effects are minimized. This can be demonstrated by considering the 1-dimensional case:



Here we can see the operation of a 1-d convolution, with  $H' = 4$ . Some filter is applied to boxes 1-4 at each  $x$ , which is then stored in the output  $y$ . In this case, the padding is made explicit as  $P_- = 2, P_+ = 3$ .

The filters are equivalent to the weights of normal neural networks, and are initialized randomly within a set of reasonable values. This implementation is relatively fast, because the number of weights is actually independent of input size - there are  $d''$  filters defined over an  $H' \times W' \times D$  volume, and each filter is applied across the entire input and trained across the entire input. For example, one basic filter may learn edge-detection. One artificial filter that this can be compared to is the Sobel filter, a convolution kernel which represents edge detection in traditional image processing:

$$G_x = \begin{bmatrix} -1 & 0 & +1 \\ -2 & 0 & +2 \\ -1 & 0 & +1 \end{bmatrix}, G_y = \begin{bmatrix} -1 & -2 & -1 \\ 0 & 0 & 0 \\ +1 & +2 & +1 \end{bmatrix}, \mathbf{G} = \sqrt{G_x^2 + G_y^2}$$

The equivalent convolution filter learned in the CNN might represent just vertical or horizontal edges. With enough filters, many unique image features can be learned. However, as is obvious, the use case of a single convolution layer is very limited, which is why we pair it with a ReLU and pooling layer, and use this trio to construct a deep NN.

### 2.2.2 Rectified Linear Unit Layer

The ReLU layer is a simple activation layer, which computes

$$y_{ijd} = \max\{0, x_{ijd}\}$$

This performs the same function as any other activation layer - adding nonlinearity to the NN. In this case, a linear activation is selected due to computational speed - it is much faster to compute a linear activation than a sigmoid, while the performance is only minimally decreased.

### 2.2.3 Pooling Layer

The pooling layer computes the maximum input in an  $H' \times W'$  patch:

$$y_{i''j''d} = \max_{1 \leq i' \leq H', 1 \leq j' \leq W'} x_{i'+i''-1, j'+j''-1, d}$$

This is useful for selecting the dominant features within each patch and essentially sparsify the data. Importantly, the pooling layer patches are non-overlapping, while the convolutional layers are not. This is also useful for providing robustness to image translation - if an image is translated by less than  $H'$  or  $W'$ , the output of the pooling layer will still be the same.

### 2.2.4 Softmax Layer

The softmax layer is the standard classification layer discussed in class:

$$y_{ijk} = \frac{e^{x_{ijk}}}{\sum_{t=1}^D e^{x_{ikt}}}$$

As discussed in class, this is useful for classification purposes. In other variations we tested, we replaced this layer with  $k$ -means clustering and multiclass SVM.

## 3 Clustering

### 3.1 $k$ -means

$k$ -means clustering is an unsupervised learning algorithm that takes a dataset and attempts to group the data into  $k$  clusters such that the sum of the distances from each point to the centroid of the cluster that it belongs to is minimized. The algorithm works by first initializing the centroids of  $k$  clusters, and at each step, it classifies the points based on which cluster it is closest to, and then updates the centroids of the clusters. One of the benefits of  $k$ -means clustering is that it will return the coordinates of the centroids of all the clusters, so for each cluster, we have a point that is “most representative” of the cluster.

We implemented the  $k$ -means algorithm ourselves. We first tested it on the MNIST raw data, which is a database of images of handwritten digits 0-9. Each image is grayscale and is 28 pixels by 28 pixels, and there are 5000 images in the dataset. So each data point can be represented as a feature vector of size  $28^2 = 784$ . We ran our  $k$ -means implementation on the MNIST dataset with  $k = 10$ , and initial centroids set equal to 10 random data points. To assign classes to clusters, for each cluster we looked at the label with the majority vote in each cluster, and assigned that label to that cluster. Therefore, to classify a new point, we first find out which cluster it is closest to, and then assign it the label corresponding to that cluster. To measure performance of the  $k$ -means clustering, we counted the number of misclassified points in the test set. For the MNIST dataset, we got an error rate of  $\frac{3695}{5000} = 73.9\%$ , suggesting that  $k$ -means does not work well on this dataset. However, it is slightly better than random, because a random classifier would give a 90% error rate. However, once we consider the 5 most likely labels for each data point, we have a 21.2% error rate, so our classifier is nontrivially better than random. We compared our implementation with MATLAB’s built-in  $k$ -means algorithm and they got almost identical results, so it is not a problem with our implementation.

We thought one possible reason  $k$ -means doesn’t work well might be because the dimensionality is too high, so we used an autoencoder to reduce the 784 features to 50, and then used  $k$ -means on the 50 features. However, we still got a high error rate of  $\frac{3570}{5000} = 71.4\%$ . Therefore, this suggests that the MNIST data is not suited for  $k$ -means. Possible reasons for this might be because there is a lot of overlap between clusters, so  $k$ -means cannot distinguish between them well. One other possible reason for poor performance is the highly correlated dependence between various dimensions, such that applying a kernel or a neural network to them will separate the clusters much better.

### 3.2 Support Vector Machine

Support Vector Machines are supervised learning models that construct a hyperplane to separate the dataset, classifying all points on one side of the hyperplane as belonging to one class, and all points on the other side as belonging to the other class. SVM aims to maximize the margin of the classification, which is the distance from the closest training point to the hyperplane. Finding the hyperplane with maximum margin can be expressed as a quadratic optimization problem (minimizing an objective while satisfying some constraints), so we can take the dual form as shown in homework 2 to solve the optimization problem. SVM is a binary classifier, so to extend it to  $n$  classes, we can train  $n$  one-vs-all classifiers: the  $i^{\text{th}}$  classifier’s output represents whether or not the point is in the  $i^{\text{th}}$  class. We adapted our code from homework 2 to extend it to the multiclass case by training one-vs-all classifiers. However, the code ran too slowly for our purposes, so we used an existing implementation of SVM along with our same multiclass extension.

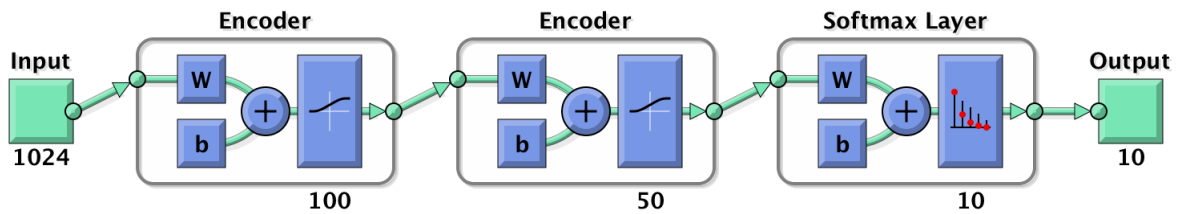
## 4 Results

The CIFAR-10 dataset consists of images corresponding to 10 different classes: airplane, automobile, bird, cat, deer, dog, frog, horse, ship, and truck. Each image is a  $32 \times 32$  color image, and there are 10000 examples in the training dataset, and 10000 examples in the test dataset. We can find the results of each combination of feature selection and clustering on the CIFAR-10 dataset.

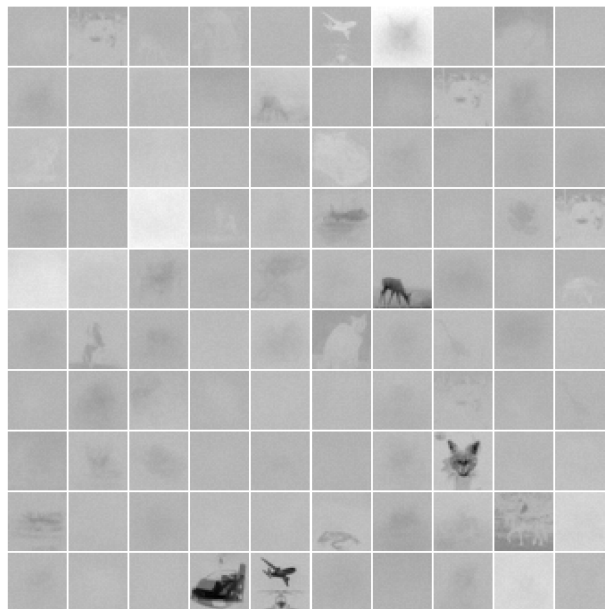
### 4.1 Autoencoder

#### 4.1.1 Autoencoder with softmax

We added a softmax layer to the autoencoder, so that the diagram of the neural network is given below:



We can also plot the weights of the first layer of the autoencoder, to see what kind of features it learned. We see that there is a little overfitting, because some of the nodes clearly resemble some of the original pictures, but otherwise it seems like the autoencoder has learned some important features.



Then, we can plot the confusion matrix, which shows how similar our classification is to the true labels of the data. This plot shows all of the counts for which a data point was classified by our classifier as belonging to class  $i$ , but it actually belongs to class  $j$ . Therefore the entries on the diagonal are when  $i = j$ , which are the cases when our classifier was correct. We see that our accuracy was 34% (therefore error rate was 66%), which is significantly better than random, but it is not very good.

		Confusion Matrix										
		1	2	3	4	5	6	7	8	9	10	
Output Class	1	387 3.9%	51 0.5%	84 0.8%	60 0.6%	81 0.8%	51 0.5%	46 0.5%	59 0.6%	163 1.6%	67 0.7%	36.9% 63.1%
	2	38 0.4%	367 3.7%	27 0.3%	42 0.4%	27 0.3%	21 0.2%	48 0.5%	36 0.4%	87 0.9%	180 1.8%	42.0% 58.0%
	3	123 1.2%	39 0.4%	294 2.9%	92 0.9%	165 1.7%	109 1.1%	114 1.1%	85 0.9%	40 0.4%	33 0.3%	26.9% 73.1%
	4	46 0.5%	56 0.6%	110 1.1%	236 2.4%	85 0.9%	185 1.8%	105 1.1%	94 0.9%	47 0.5%	57 0.6%	23.1% 76.9%
	5	81 0.8%	29 0.3%	141 1.4%	88 0.9%	274 2.7%	74 0.7%	105 1.1%	112 1.1%	52 0.5%	39 0.4%	27.5% 72.5%
	6	42 0.4%	26 0.3%	108 1.1%	164 1.6%	81 0.8%	307 3.1%	82 0.8%	93 0.9%	41 0.4%	34 0.3%	31.4% 68.6%
	7	36 0.4%	66 0.7%	99 1.0%	130 1.3%	138 1.4%	79 0.8%	368 3.7%	53 0.5%	21 0.2%	52 0.5%	35.3% 64.7%
	8	57 0.6%	40 0.4%	73 0.7%	77 0.8%	79 0.8%	93 0.9%	48 0.5%	374 3.7%	43 0.4%	66 0.7%	39.4% 60.6%
	9	153 1.5%	140 1.4%	42 0.4%	48 0.5%	56 0.6%	43 0.4%	38 0.4%	35 0.4%	433 4.3%	114 1.1%	39.3% 60.7%
	10	37 0.4%	186 1.9%	22 0.2%	63 0.6%	14 0.1%	38 0.4%	46 0.5%	59 0.6%	73 0.7%	358 3.6%	40.0% 60.0%
		38.7% 61.3%	36.7% 63.3%	29.4% 70.6%	23.6% 76.4%	27.4% 72.6%	30.7% 69.3%	36.8% 63.2%	37.4% 62.6%	43.3% 56.7%	35.8% 64.2%	34.0% 66.0%
		1	2	3	4	5	6	7	8	9	10	
		Target Class										

#### 4.1.2 Autoencoder with SVM

We used the features extracted by the autoencoder as inputs to the one-vs-all SVM multiclass method. On the CIFAR-10 dataset, all of the one-vs-all classifiers outputted 0 for all the data points, so none of the points were actually classified. This meant that none of the classes were significantly likely. Since we could not get any results for the CIFAR-10 dataset, we tried this method on the MNIST dataset. Around half of the data points were not classified by any of the 10 SVM classifiers, but some of the data points were classified, and the accuracy was high for those cases. If a data point was not classified by any of the 10 SVM classifiers, it gets assigned the 10th class by default.

	1	2	3	4	5	6	7	8	9	10	
1	349 7.0%	0 0.0%	0 0.0%	2 0.0%	0 0.0%	0 0.0%	26 0.5%	0 0.0%	0 0.0%	0 0.0%	92.6% 7.4%
2	3 0.1%	172 3.4%	4 0.1%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	96.1% 3.9%
3	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	NaN% NaN%
4	0 0.0%	0 0.0%	9 0.2%	277 5.5%	0 0.0%	6 0.1%	0 0.0%	1 0.0%	0 0.0%	0 0.0%	94.5% 5.5%
5	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	NaN% NaN%
6	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	NaN% NaN%
7	8 0.2%	0 0.0%	0 0.0%	10 0.2%	0 0.0%	0 0.0%	254 5.1%	0 0.0%	1 0.0%	0 0.0%	93.0% 7.0%
8	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	NaN% NaN%
9	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	12 0.2%	2 0.0%	240 4.8%	0 0.0%	94.5% 5.5%
10	166 3.3%	331 6.6%	480 9.6%	202 4.0%	517 10.3%	486 9.7%	207 4.1%	500 10.0%	247 4.9%	488 9.8%	13.5% 86.5%
	66.3% 33.7%	34.2% 65.8%	0.0% 100%	56.4% 43.6%	0.0% 100%	0.0% 100%	50.9% 49.1%	0.0% 100%	49.2% 50.8%	100% 0.0%	35.6% 64.4%
	1	2	3	4	5	6	7	8	9	10	

This plot is a confusion matrix, explained earlier. We can see that our classifier only predicted classes of 1, 2, 4, 7, 9, and 10. When it did not predict 10, its accuracy was 93.9%, which is pretty remarkable, but it predicted 10 most of the time which meant that the classifier was unsure on most of the data points. Therefore, we see that the one-vs-all strategy for multiclass SVM does not yield the best results, because if there is no class is more than 50% likely for a given point, SVM will not classify the point into any class.

### 4.1.3 Autoencoder with $k$ -means

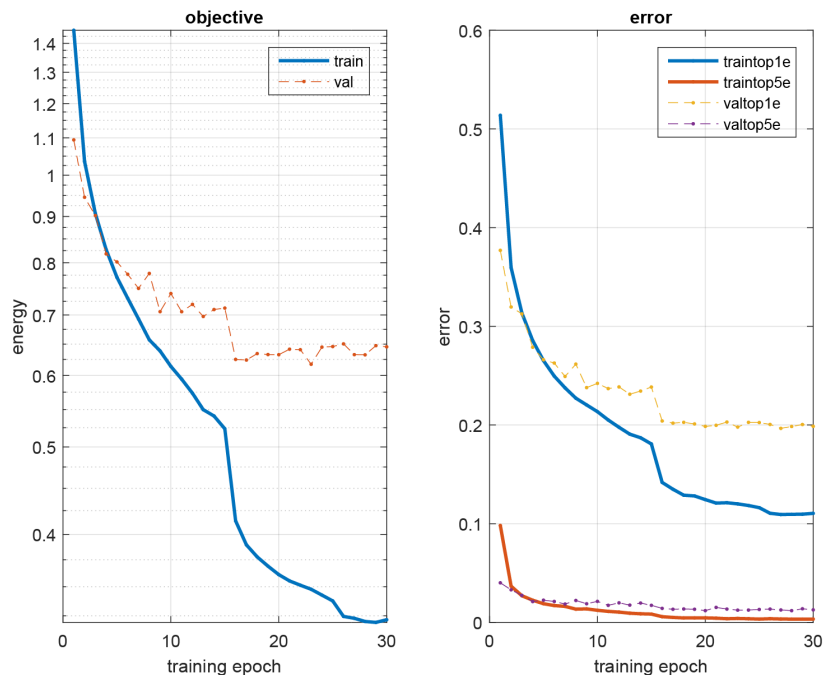
We ran the  $k$ -means algorithm on features extracted using an autoencoder, but the algorithm converged to only have one cluster. This probably means that the features extracted were not distinct enough, so the  $k$ -means algorithm could not find two separate clusters that would be more optimal than one cluster. From the confusion matrix below, we see that the accuracy of this is 10.3%, which is basically equivalent to random guessing. From our discussion of  $k$ -means above, this low accuracy is expected, because  $k$ -means did not work well with MNIST, and classification on the CIFAR-10 dataset is harder than classification on MNIST.

	1	2	3	4	5	6	7	8	9	10	
1	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	NaN% NaN%
2	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	NaN% NaN%
3	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	NaN% NaN%
4	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	NaN% NaN%
5	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	NaN% NaN%
6	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	NaN% NaN%
7	1005 10.1%	974 9.7%	1032 10.3%	1016 10.2%	999 10.0%	937 9.4%	1030 10.3%	1001 10.0%	1025 10.2%	981 9.8%	10.3% 89.7%
8	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	NaN% NaN%
9	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	NaN% NaN%
10	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	NaN% NaN%
	0.0% 100%	0.0% 100%	0.0% 100%	0.0% 100%	0.0% 100%	0.0% 100%	0.0% 100%	0.0% 100%	0.0% 100%	0.0% 100%	10.3% 89.7%
	1	2	3	4	5	6	7	8	9	10	

## 4.2 Convolutional Neural Network

### 4.2.1 Convolutional Neural Network with Softmax

The training of the CNN proceeded via SGD as mentioned, attempting to minimize an energy function which essentially captured the error in training. The error and objective function for the training data and validation data was plotted over the course of 30 epochs of training. After about 15 epochs, it can be seen that the validation error rate has stabilized to around 20%. When including the top 5 results of the CNN, validation error rate stabilized to less than 3% (`valtop5e`). In this case, each epoch represents passing through the entire training set, or 10,000 individual SGD steps.





This final error of around 20% is notable, since it is in the same ballpark as the best results obtained for CIFAR-10, which report as low as 15% error without data augmentation (we also don't use data augmentation, so this is a valid comparison). This was achieved with a very similar structure to the CNN - softmax combo we used to classify CIFAR-10, but with optimized weights for hyperparameters such as learning rate. The same approximate structure is found in GoogLeNet, which won the most recent ImageNet classification challenge.

### 4.2.2 Convolutional Neural Network with SVM

Multiclass SVM was much more successful in clustering the feature vectors from CNN's than autoencoders, reaching a final accuracy of a very disappointing 11.5% (hardly better than random guessing). Again, label 10 was assigned by default if no other label had probability  $\geq 50\%$ . Interestingly, labels 3 and 9 seemed popular for the SVM to predict, though it is unclear why this might be the case, and merits future analysis.

**Confusion Matrix**

1	653 6.5%	5 0.1%	61 0.6%	9 0.1%	7 0.1%	4 0.0%	4 0.0%	5 0.1%	25 0.3%	4 0.0%	84.0%
2	14 0.1%	1 0.0%	510 5.1%	19 0.2%	6 0.1%	20 0.2%	19 0.2%	10 0.1%	0 0.0%	1 0.0%	0.2%
3	43 0.4%	6 0.1%	188 1.9%	169 1.7%	888 8.9%	123 1.2%	186 1.9%	249 2.5%	19 0.2%	14 0.1%	10.0%
4	12 0.1%	34 0.3%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	1 0.0%	594 5.9%	0.0%
5	2 0.0%	1 0.0%	19 0.2%	26 0.3%	9 0.1%	1 0.0%	638 6.4%	0 0.0%	1 0.0%	3 0.0%	1.3%
6	12 0.1%	6 0.1%	0 0.0%	1 0.0%	1 0.0%	0 0.0%	0 0.0%	1 0.0%	649 6.5%	2 0.0%	0.0%
7	46 0.5%	376 3.8%	5 0.1%	3 0.0%	1 0.0%	1 0.0%	4 0.0%	2 0.0%	99 1.0%	83 0.8%	0.6%
8	0 0.0%	390 3.9%	0 0.0%	0 0.0%	1 0.0%	0 0.0%	0 0.0%	0 0.0%	1 0.0%	13 0.1%	0.0%
9	15 0.1%	12 0.1%	62 0.6%	624 6.2%	15 0.1%	741 7.4%	41 0.4%	74 0.7%	12 0.1%	7 0.1%	0.7%
10	203 2.0%	169 1.7%	155 1.6%	149 1.5%	72 0.7%	110 1.1%	108 1.1%	659 6.6%	193 1.9%	279 2.8%	13.3%
	65.3%	0.1%	18.8%	0.0%	0.9%	0.0%	0.4%	0.0%	1.2%	27.9%	11.5%
	34.7%	99.9%	81.2%	100%	99.1%	100%	99.6%	100%	98.8%	72.1%	88.5%
	1	2	3	4	5	6	7	8	9	10	
											Target Class

### 4.2.3 Convolutional Neural Network with $k$ -means

By taking the activations of the last convolution layer, right before the softmax layer, we obtain a lower-dimensionality representation of our input. Rather than send this to softmax, we instead used this as the input data for  $k$ -means clustering, where again the most common label within each cluster was assigned to the cluster. This yielded quite decent results, with an overall 63.9% accuracy, significantly better than any other feature extraction/classification pair except CNN with softmax, which achieved an 80% overall accuracy.

**Confusion Matrix**

	1	2	3	4	5	6	7	8	9	10	
1	659 6.6%	6 0.1%	54 0.5%	7 0.1%	5 0.1%	2 0.0%	4 0.0%	4 0.0%	25 0.3%	3 0.0%	85.7% 14.3%
2	2 0.0%	642 6.4%	0 0.0%	0 0.0%	1 0.0%	0 0.0%	0 0.0%	0 0.0%	5 0.1%	23 0.2%	95.4% 4.6%
3	27 0.3%	1 0.0%	545 5.5%	28 0.3%	8 0.1%	21 0.2%	20 0.2%	11 0.1%	3 0.0%	1 0.0%	82.0% 18.0%
4	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	NaN% NaN%
5	71 0.7%	6 0.1%	270 2.7%	201 2.0%	881 8.8%	163 1.6%	188 1.9%	310 3.1%	20 0.2%	16 0.2%	41.4% 58.6%
6	26 0.3%	10 0.1%	84 0.8%	694 6.9%	22 0.2%	788 7.9%	68 0.7%	87 0.9%	17 0.2%	9 0.1%	43.7% 56.3%
7	2 0.0%	2 0.0%	27 0.3%	30 0.3%	13 0.1%	1 0.0%	704 7.0%	0 0.0%	3 0.0%	3 0.0%	89.7% 10.3%
8	2 0.0%	0 0.0%	1 0.0%	5 0.1%	49 0.5%	16 0.2%	0 0.0%	565 5.7%	0 0.0%	0 0.0%	88.6% 11.4%
9	18 0.2%	5 0.1%	0 0.0%	1 0.0%	1 0.0%	0 0.0%	0 0.0%	1 0.0%	665 6.7%	3 0.0%	95.8% 4.2%
10	193 1.9%	328 3.3%	19 0.2%	34 0.3%	20 0.2%	9 0.1%	16 0.2%	22 0.2%	262 2.6%	942 9.4%	51.1% 48.9%
	65.9% 34.1%	64.2% 35.8%	54.5% 45.5%	0.0% 100%	88.1% 11.9%	78.8% 21.2%	70.4% 29.6%	56.5% 43.5%	66.5% 33.5%	94.2% 5.8%	63.9% 36.1%
	1	2	3	4	5	6	7	8	9	10	

Target Class

## 5 Discussion

Convolutional Neural Networks with a softmax classifier were by far the most effective pairing of feature extraction and classification. This makes sense - they're the state of the art for image classification, due to convolutional neural networks' ability to efficiently extract contextual image information and softmax's kernel-like trainability. It is understandable that the more optimized topology of the CNN would easily surpass the accuracy of the autoencoder.

Most interestingly, taking the final 64-dimensional activations of the final convolutional layer in the CNN and using them as inputs to a  $k$ -means clustering algorithm yielded very favorable results, achieving a very respectable 64% accuracy on CIFAR-10. This also holds potential interest, since  $k$ -means is an unsupervised clustering algorithm, it could theoretically be used in a semi-supervised classification task, where it would know to create a separate cluster for previously-unseen data. The effectiveness of this strategy is tremendously amplified when feeding it outputs from a CNN than raw data, and clustering can be done much more quickly due to the significantly decreased dimensionality. In order to be able to tell when a new cluster is added, a regularization function on the number of clusters would have to be introduced and the number of clusters allowed to vary.

Overall, the strategy of coupling a simple clustering algorithm to a feature extraction algorithm was less successful than anticipated, most likely due to the very high covariance of the extracted features, which prevents standard non-kernelized clustering algorithms from performing well. Manually kernelizing a clustering algorithm was not pursued due to the difficulty in coming up with a reasonable kernel for the artificial data. Softmax instead served as a substitute, allowing something similar to a kernelized clustering algorithm to be learned by the NN's.

## 6 Works Cited

- Learning Multiple Layers of Features from Tiny Images, Alex Krizhevsky, 2009.
- "MatConvNet - Convolutional Neural Networks for MATLAB", A. Vedaldi and K. Lenc,

Proc. of the ACM Int. Conf. on Multimedia, 2015.