

Reed Ballesteros
MSDS 458-DL: Artificial Intelligence & Deep Learning
Spring 2023
Prof. Syamala Srinivasan, Ph.D.
June 4, 2023

A.4: Fourth Research/Programming Assignment: Convolutional Neural Networks (CNNs) and the Fashion-MNIST Dataset

Abstract

This study details the development of several convolutional neural networks (CNNs) trained and tested for their performance accuracy using the Fashion-MNIST image dataset. The most optimal model would be integrated as part of a fashion retailer's recommendation system to offer other products to customers browsing their website. While we find a combination of regularization techniques applied on a CNN can yield significant gains in accuracy, we also compare them against the performance of various pre-trained models in which we find do not classify the dataset as well as our own models. Our most optimal CNN containing 9 total layers in 3 clustered groups was able to achieve a 93.7% accuracy against the Fashion-MNIST test dataset.

Introduction

This report documents the development of a convolutional neural network (CNN) with the goal of classifying clothing articles. The resulting model would be integrated into a fashion retailer's website in which it would scan and classify images that are displayed while a customer is browsing products on the website and its recommendation system would refer other items for them as well. Image classification accuracy for the system is important so the website can properly recommend to a visitor similar or related products that are popular, in-stock/available, are part of a current marketing campaign or promotion, or upsell to higher-end products. Displaying unrelated products could be jarring to the user experience and could possibly discourage them from continuing shopping on the website and thus a missed opportunity on a completed online sales transaction and collecting user information. We will choose the model configuration that yields the highest accuracy in image classification using the Fashion-MNIST dataset for model training, validation, and testing. We will also use this opportunity to build and compare several pre-trained CNNs on the same dataset.

Literature Review

There have been several studies on convolutional neural network (CNN) development using the Fashion-MNIST dataset. One study compared several CNNs trained with the dataset to create a model which

could be integrated into service robot technologies that could assist the elderly by identifying objects in their home such as clothing (Nocentini et al. 2022). The research led to the creation of the Multiple Convolutional Network (MCNN15) containing 15 CNNs achieving a test accuracy of 94% on the Fashion-MNIST dataset, surpassing other previous pre-trained models such as AlexNet and ResNet. Inspired by MCNN15, we created CNN models using sets of grouped convolutional layers and are able to achieve a 93.7% accuracy with the Fashion-MNIST test dataset.

Another study compared several pre-trained models using the Fashion-MNIST dataset in which the authors were able to yield a 96.2% test accuracy using the PyramidNet pre-trained model (Tang et al. 2020). The study also used a version of the ResNet pre-trained model where its optimal model was able to yield 94.2% accuracy. While our study includes the use of pre-trained models, including a version of ResNet, we were not able to achieve similar results with our ResNet pre-trained model.

Methods

Research Design

We will conduct our research by first downloading the Fashion-MNIST dataset available in the Keras framework and perform exploratory data analysis (EDA) on it to understand the data and its corresponding labels we will be using to build, train, and test our models. Each model is represented as an experiment. Our first section of experiments will involve building CNN models using 3 isolated convolution layers followed by 2 fully connected layers and a final softmax classification layer. Model variations will be created by adding regularization components such as batch normalization, dropout, and ridge regression (also called 'L2 Regularization') along with hyperparameter tuning of these components. We will also create a 4-layer CNN model and compare its performance to the 3-layer CNN models.

The second section of models contains a more complex architecture utilizing 3 groups of 3 convolutional layers in which each group is separated by a max pooling layer. Model variations in this section will be based on tuning regularization hyperparameters. The third section of models will utilize pre-trained models such as VGG16, VGG19, Inception V3 and ResNet152 V3, each of which will use a transformed Fashion-MNIST dataset to meet their respective image format requirements. The models in both of these sections will also be followed by 2 fully connected layers and a final softmax classification layer.

For each model we will record their accuracy and loss against the test subset, as well as precision, recall, F1 score, and RMSE. We will also plot each model's training and validation loss and accuracy over their respective number of epochs of training. Training for each model will automatically stop after validation accuracy does not improve after a given number of epochs, in which the best training epoch will be saved before the model overfits to the training data. Confusion matrices for each model will also be charted to find how well they perform for each image classification.

Implementation

Our models will be created using the Keras and TensorFlow frameworks available in Python. These are the most common tools used in developing such models based on their ease of use, customization options available, and their scalability for each convolutional or dense layer. Each convolutional and dense hidden layer will use the ReLU activation function as its simplicity in calculation can make the feedforward and backpropagation process fast, which will be helpful when we train deep models that contain a large number of parameters. The final output 10-way classification layer will use a softmax activation function.

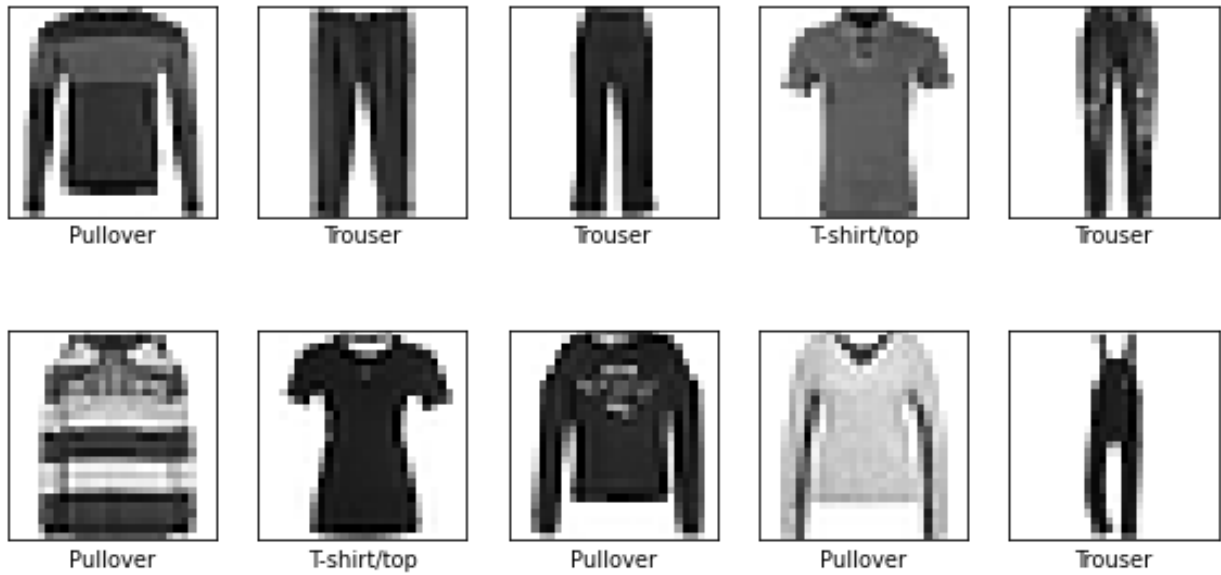
The Keras model compiler will be configured to use the Adam optimizer, calculate cross entropy loss using the SparseCategoricalCrossentropy class, and model performance will be based on accuracy metrics. For reproducibility and stabilize the random nature of the model training process, we will set a system seed of 43.

Dataset: Fashion-MNIST

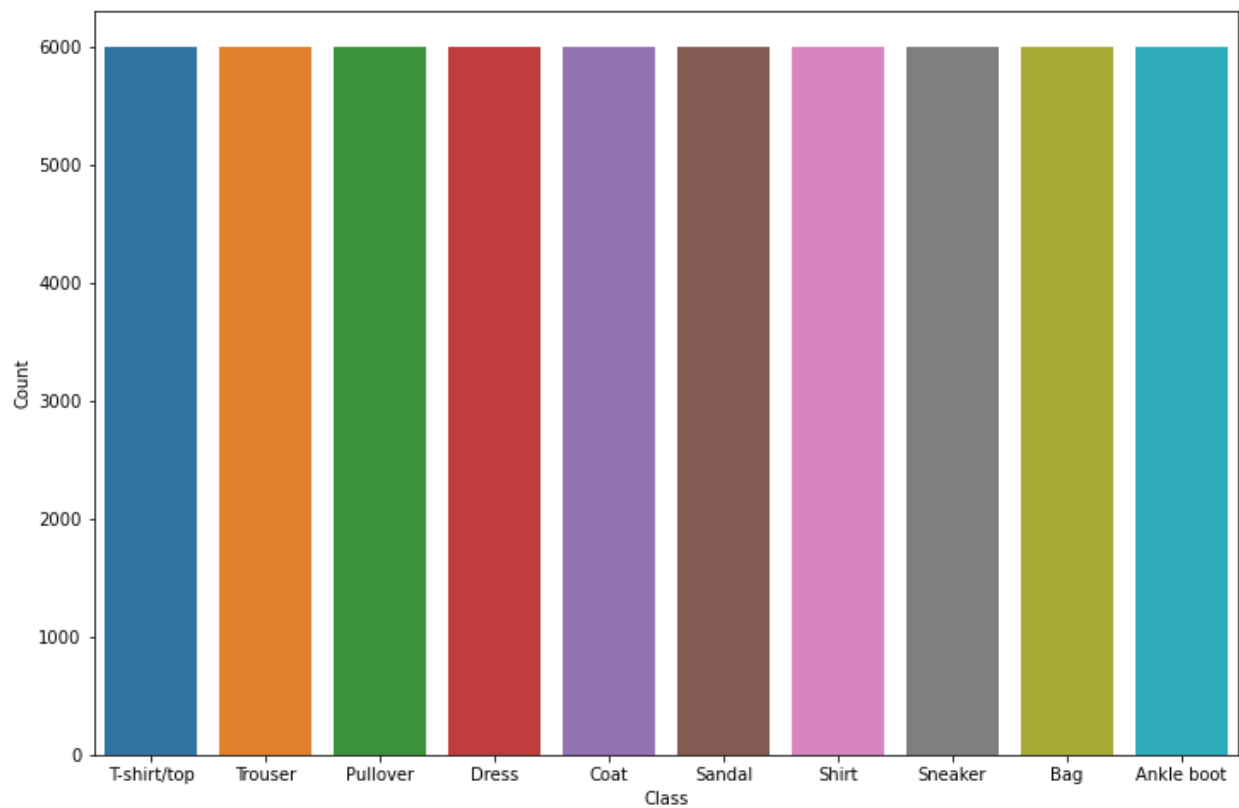
The Fashion-MNIST dataset provided by German online retailer Zalando is a collection of images (60000 for training, 10000 for testing), each categorized as 1 of 10 classes (t-shirt/top, trouser, pullover, dress, coat, sandal, shirt, sneaker, bag, ankle boot). The dataset is commonly used for benchmarking machine learning models. Each image is represented as a 28x28 matrix over a single greyscale channel, where each value in the matrix represents a pixel integer value between 0 to 255. We will transform the integer pixel data to a standardized value between 0 and 1 before feeding the image into the models. We will also transform the dataset into three-channel color images and resize them to a larger scale to meet the minimum image size requirements for each respective pre-trained model.

Exploratory Data Analysis (EDA)

As described above, the Fashion-MNIST dataset is a collection of 70000 28x28 color images, each classified as one of 10 categories: t-shirt/top, trouser, pullover, dress, coat, sandal, shirt, sneaker, bag, ankle boot. Below is a sample of images from the dataset:



The training dataset contains 6000 samples of each classification.

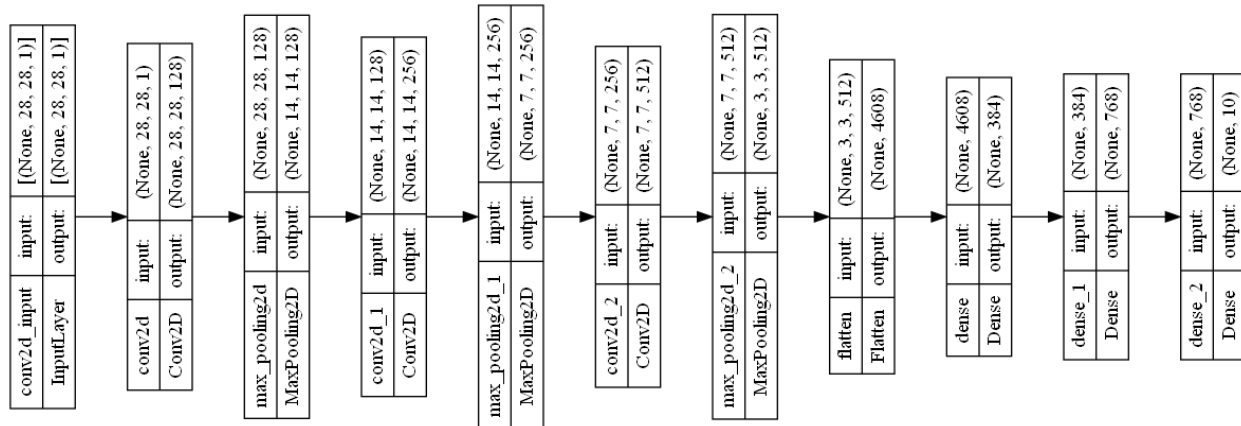


Results

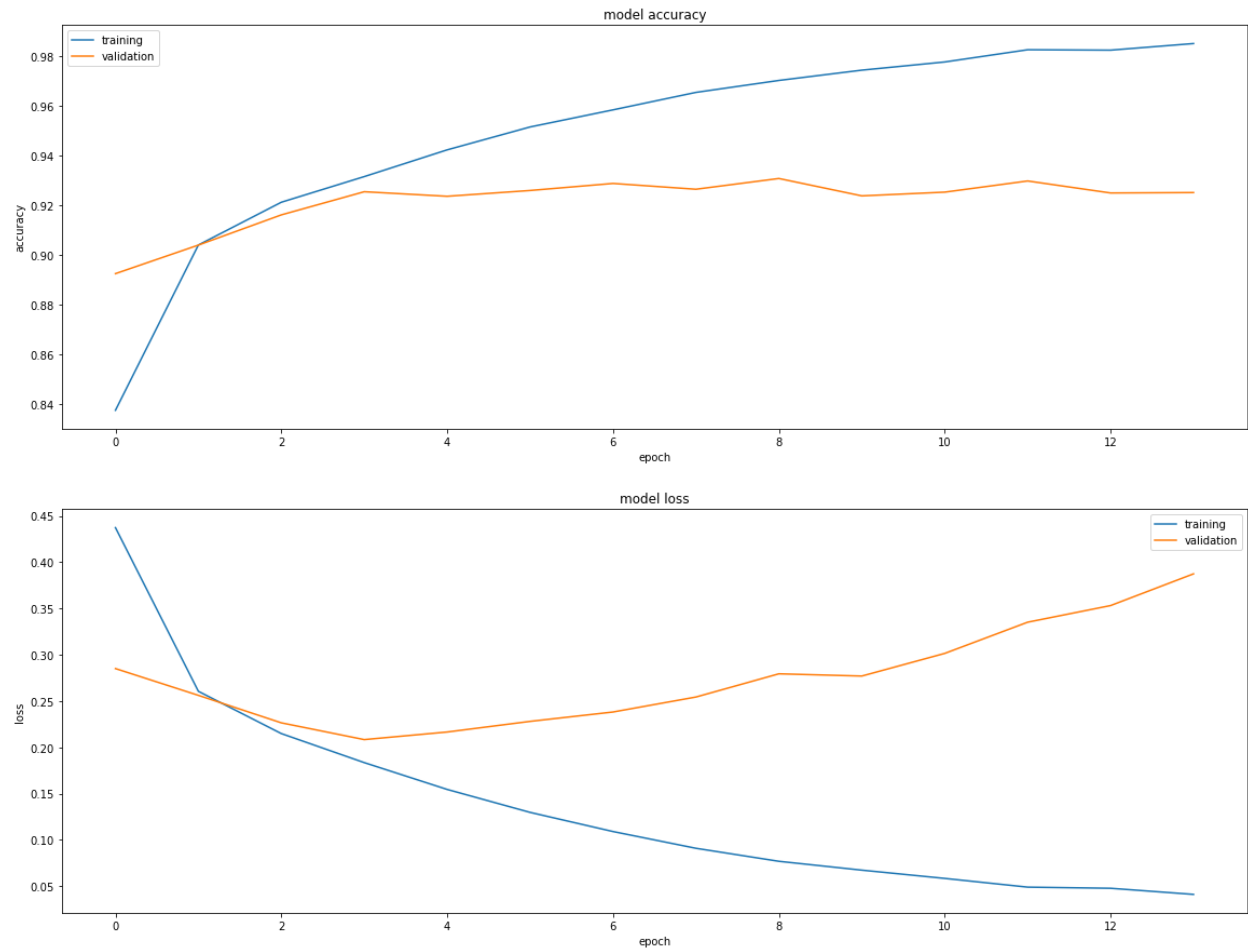
Section One: Convolutional Models

Experiment 1: CNN with 3 convolution/max pooling layers & 2 fully connected layers – no regularization

In this experiment we use a CNN consisting of three convolutional layers using a kernel size of 3x3 with the first layer containing 128 filters, the second layer containing 256 filters, and the third layer containing 512 filters, each separated by a 2x2 max pooling layer. The convolutional base then feeds into a DNN containing 2 fully connected layers, with the first layer containing 384 units and the second layer containing 768 units, followed by the 10-way softmax classification output layer.



Training and validation accuracy and cross entropy loss charts are below.



We see validation results start to diverge from the training results within 5 epochs, indicating signs of the model overfitting to the training data. The resulting test data accuracy was 91.9%, which is a high start in what we would expect our other models to be in within a similar range.

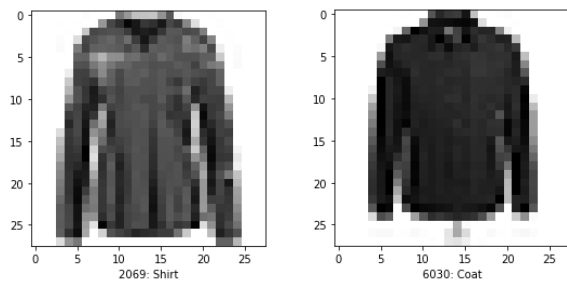
The experiment yields the following confusion matrix:

T-shirt/top	912	0	12	7	1	1	63	0	4	0
Trouser	0	986	0	10	2	0	1	0	1	0
Pullover	21	1	902	5	37	0	34	0	0	0
Dress	25	1	9	928	18	0	16	0	3	0
Coat	0	1	57	25	895	0	22	0	0	0
Sandal	0	0	0	0	0	984	0	14	0	2
Shirt	149	1	63	20	88	0	673	0	6	0
Sneaker	0	0	0	0	0	1	0	994	0	5
Bag	11	0	1	1	2	1	6	2	976	0
Ankle boot	0	0	0	0	0	6	1	53	0	940
	T-shirt/top	Trouser	Pullover	Dress	Coat	Sandal	Shirt	Sneaker	Bag	Ankle boot

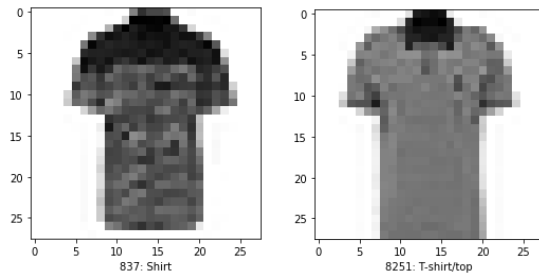
predicted label

From the matrix above we can see how well the model can correctly classify images from the Fashion-MNIST dataset, with the exception of shirts. Shirt misclassification would be a common theme shown in all confusion matrices for each model in this study, as images in the dataset would understandably look very similar to other clothing categories like tops, pullovers, and coats.

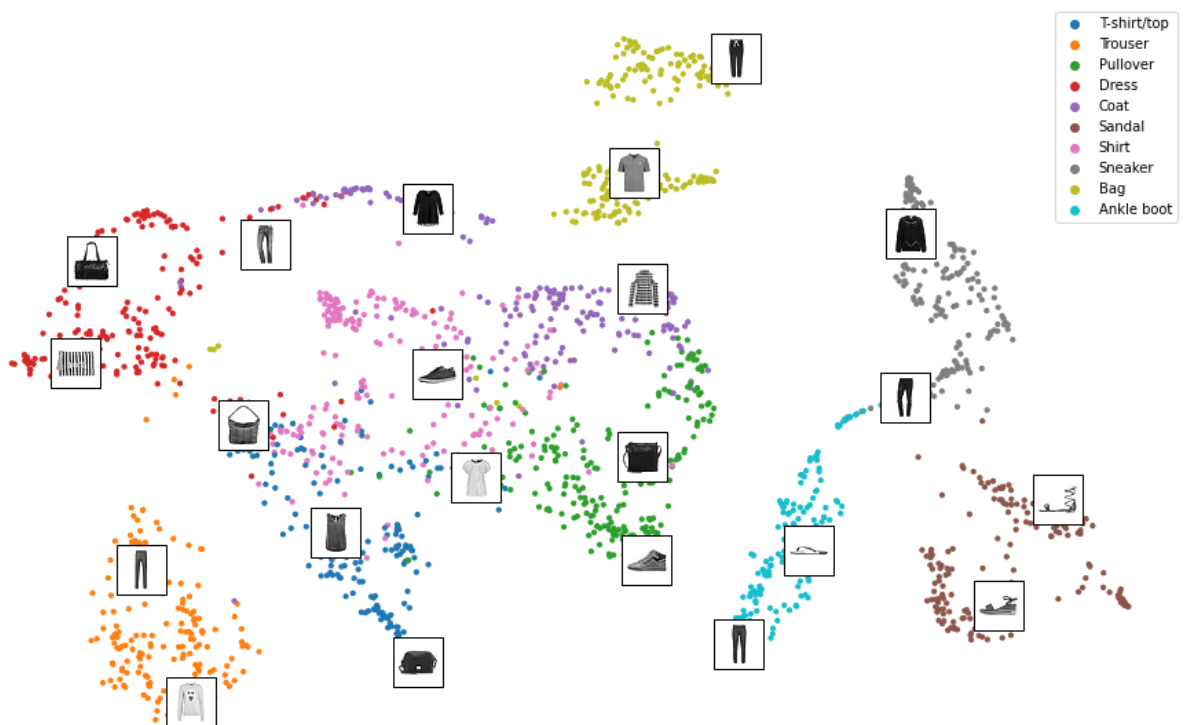
Images labeled 'shirt' and 'coat' in the dataset can visibly look similar such as these examples:



As well as between images labeled 'shirt' and 't-shirt/top':



We can reduce the model to 2 features and create a scatterplot to observe the distribution of a sample of model classifications.



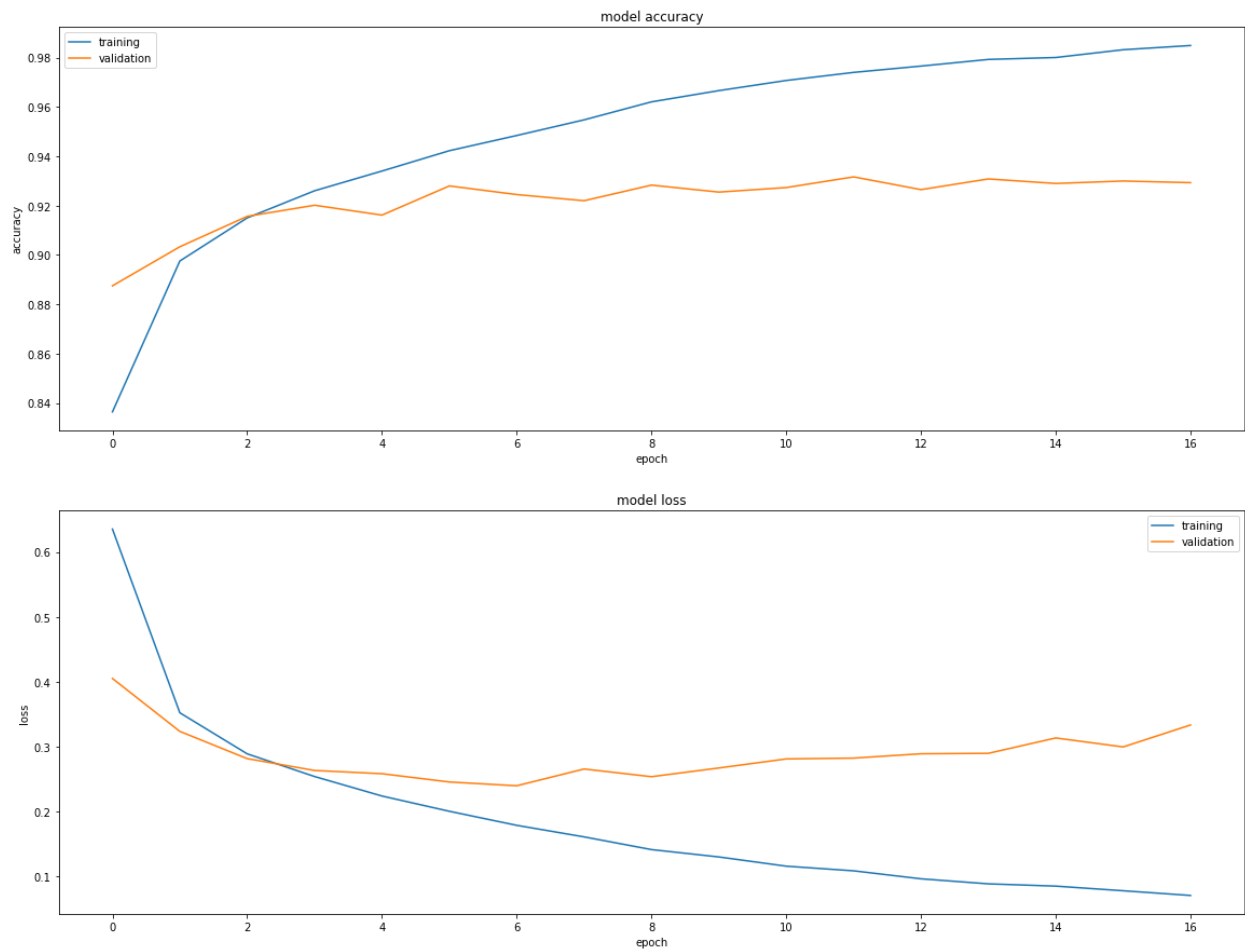
We can see distinct groups form in this scatterplot sample due to the model's high test accuracy. We've observed similar clustering in our other models as well as they're also able to achieve high test accuracy of about 90%.

Experiment 2: CNN with 3 convolution/max pooling layers & 2 fully connected layers – L2 Regularization (0.001)

Experiment 2 is similar to Experiment 1, containing a CNN consisting of three 3x3 convolutional layers with the first layer containing 128 filters, the second layer containing 256 filters, and the third layer containing 512 filters. The convolutional base then feeds into a DNN containing 2 fully connected layers,

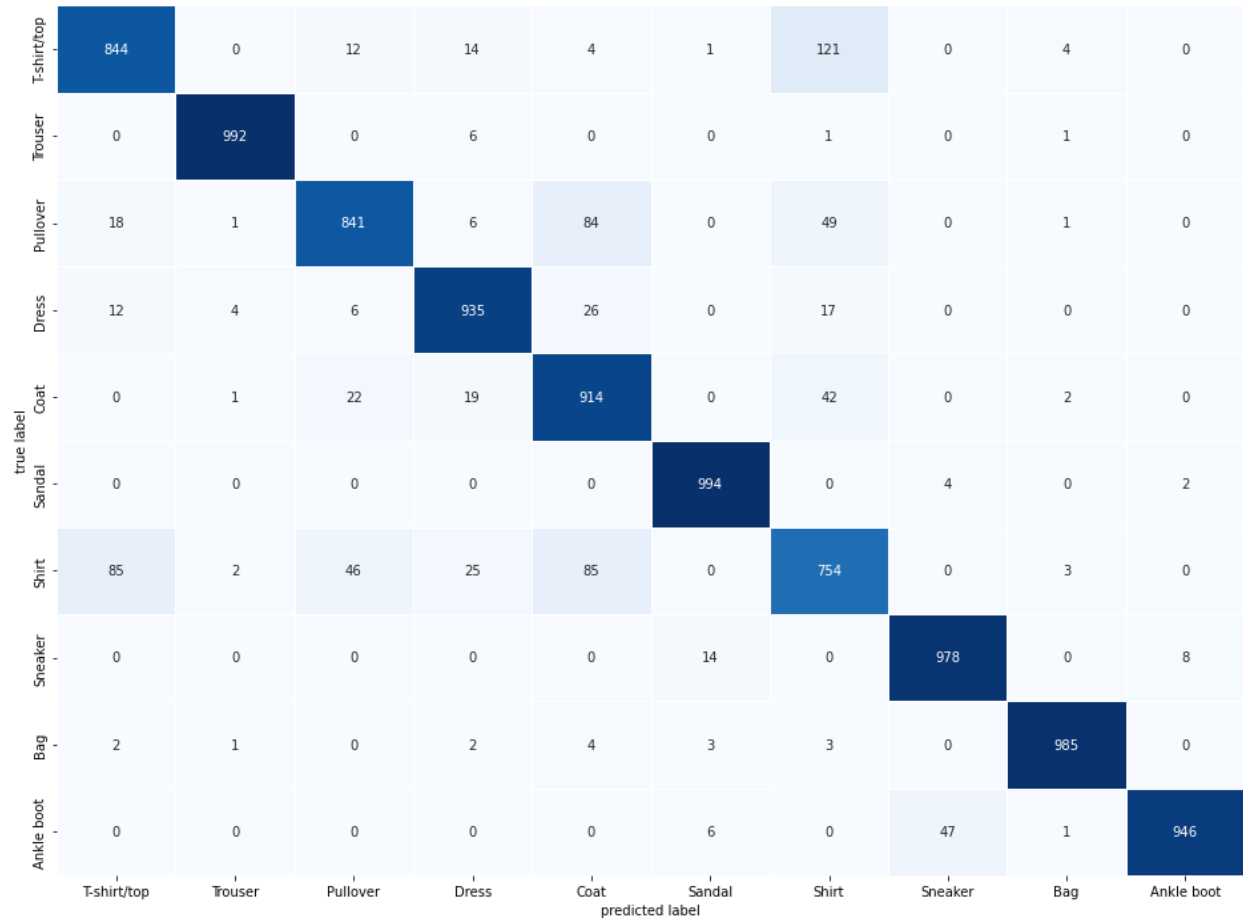
with the first layer containing 384 units and the second layer containing 768 units, followed by the 10-way softmax classification output layer. All fully connected layers utilize L2 regularization (0.001).

Training and validation accuracy and cross entropy loss charts are below.



We see validation results start to diverge from the training results within 4 epochs. The resulting test data accuracy is 91.83%, just 0.07% less than Experiment 1's accuracy.

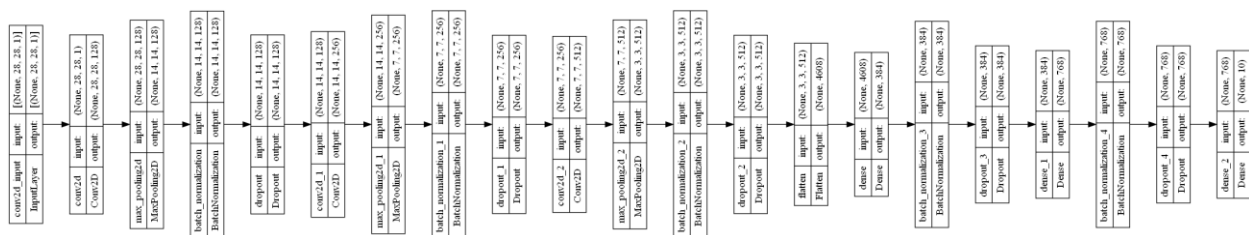
The experiment yields the following confusion matrix:



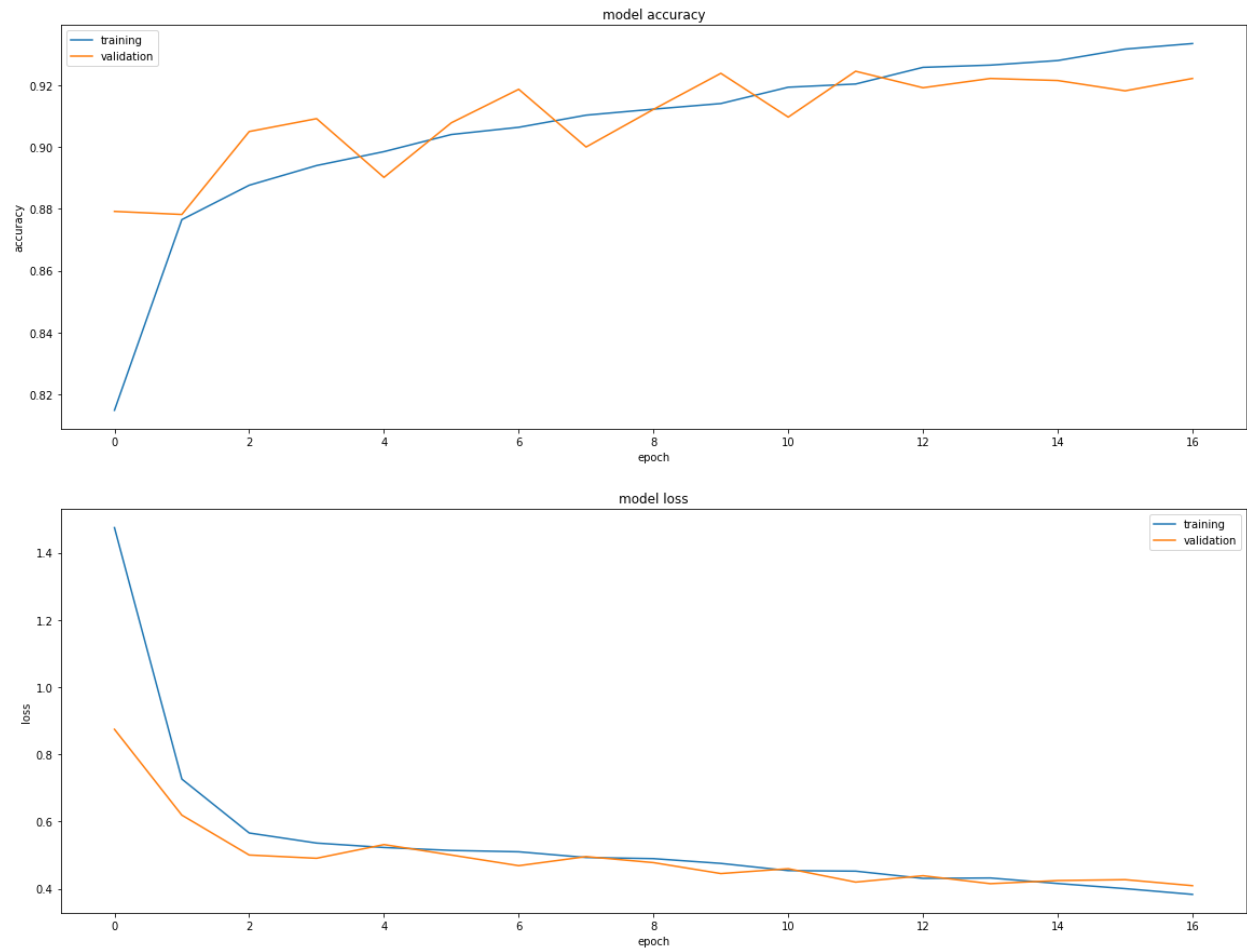
As expected, we see that this experiment also has issues classifying shirts as well.

Experiment 3: CNN with 3 convolution/max pooling layers & 2 fully connected layers – Add Batch Normalization and Dropout (0.3)

Experiment 3 adds batch normalization and dropout (0.3) regularization processes to all convolution and fully connected layers. A diagram of the full architecture is below.



Training and validation accuracy and cross entropy loss charts are below.



We see validation results start to diverge from the training results within 12 epochs. The resulting test data accuracy is 92.12%, exceeding Experiment 1's test accuracy. We see that adding regularization processes can improve model accuracy.

The experiment yields the following confusion matrix:

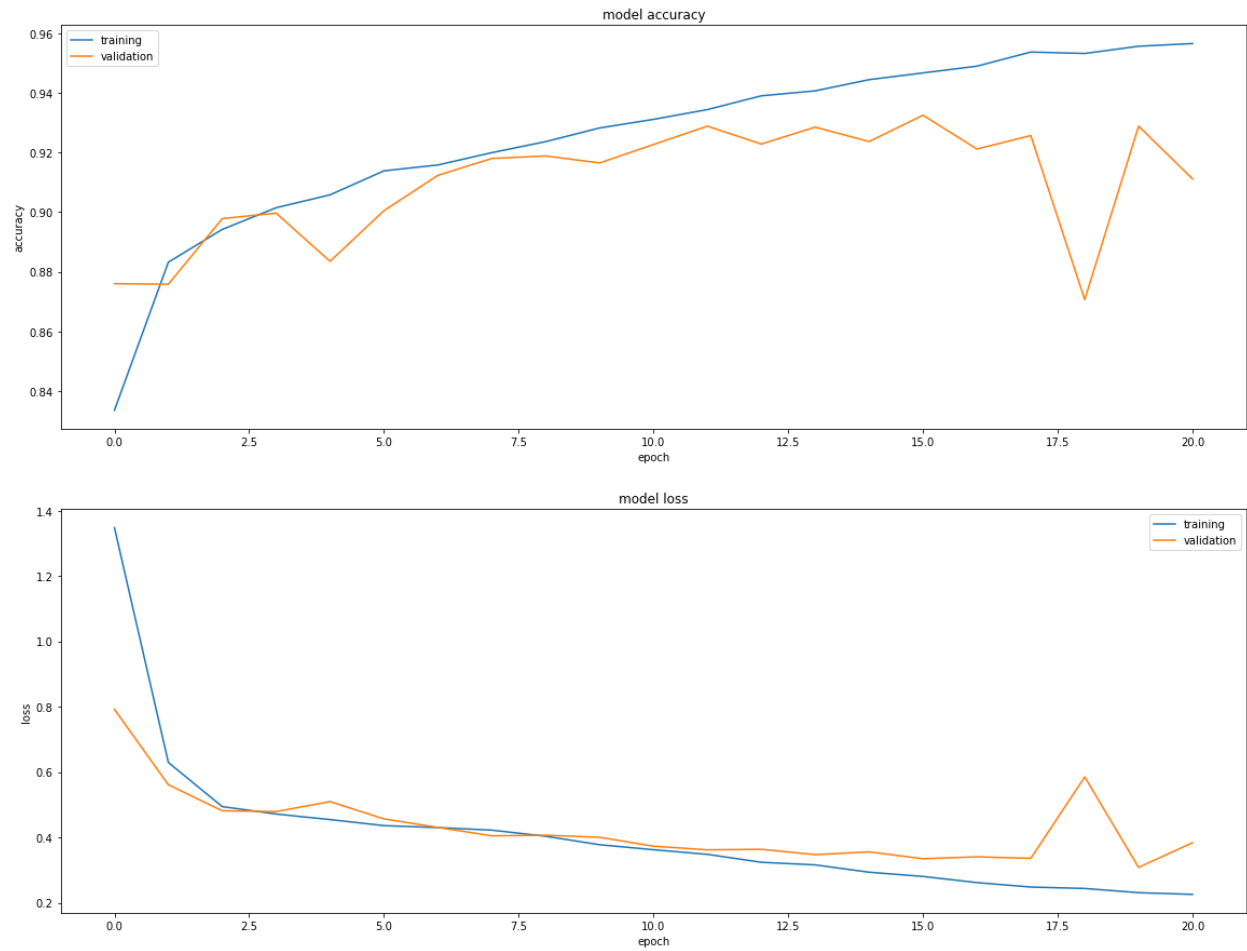
true label	T-shirt/top	911	0	19	8	3	1	51	0	7	0
	Trouser	1	973	0	22	1	0	1	0	2	0
	Pullover	17	1	938	11	23	0	9	0	1	0
	Dress	20	0	9	941	10	0	20	0	0	0
	Coat	0	0	80	35	848	0	36	0	1	0
	Sandal	0	0	0	0	0	978	0	19	0	3
	Shirt	115	0	97	25	58	0	702	0	3	0
	Sneaker	0	0	0	0	0	1	0	992	0	7
	Bag	3	0	0	6	0	0	1	3	986	1
	Ankle boot	0	0	0	0	0	4	0	52	1	943
		predicted label									

While the model can classify most images fairly well, it too has issues in classifying shirts.

Experiment 4: CNN with 3 convolution/max pooling layers & 2 fully connected layers – Change Dropout (0.2)

Experiment 4 is similar to Experiment 3 but changes the dropout rate to 0.2.

Training and validation accuracy and cross entropy loss charts are below.



We see validation results start to diverge from the training results within 7 epochs. The resulting test data accuracy is 92.45%, exceeding Experiment 3's test accuracy.

The experiment yields the following confusion matrix:

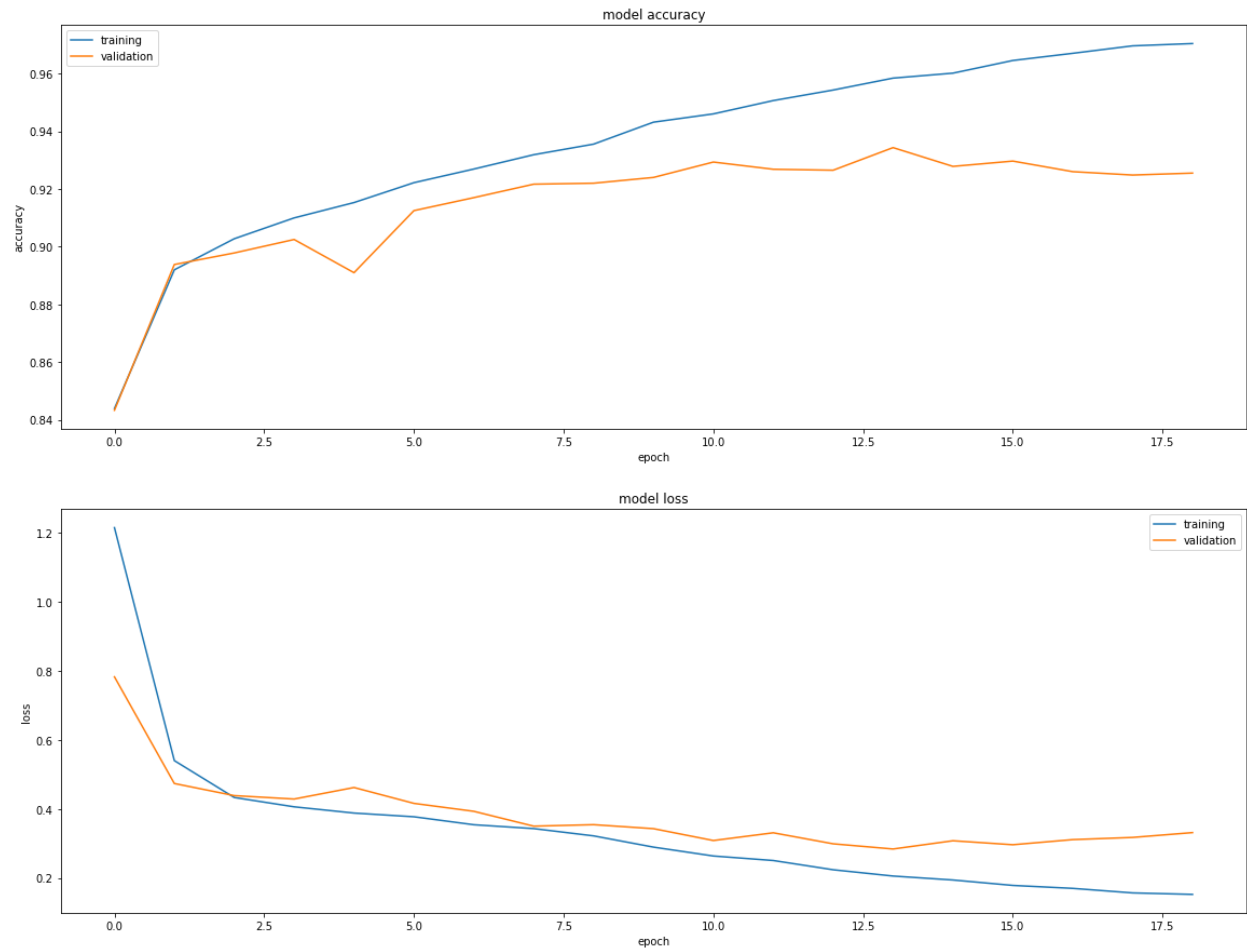
true label	T-shirt/top	874	0	14	25	4	1	78	0	4	0
	Trouser	1	977	0	19	1	0	2	0	0	0
	Pullover	19	1	869	9	58	0	44	0	0	0
	Dress	8	0	7	943	28	0	14	0	0	0
	Coat	1	1	13	13	901	0	71	0	0	0
	Sandal	0	0	0	0	0	974	0	23	0	3
	Shirt	97	0	44	32	44	0	780	0	3	0
	Sneaker	0	0	0	0	0	0	0	979	0	21
	Bag	0	1	1	5	2	4	4	1	981	1
	Ankle boot	0	0	0	0	0	6	1	26	0	967
		predicted label									

Shirt misclassification is still present compared to other categories.

Experiment 5: CNN with 3 convolution/max pooling layers & 2 fully connected layers – Change Dropout (0.1)

Experiment 5 is similar to Experiment 3 but changes the dropout rate to 0.1.

Training and validation accuracy and cross entropy loss charts are below.



We see validation results start to diverge from the training results by 10 epochs. The resulting test data accuracy is 92.54%, exceeding Experiment 4's test accuracy.

The experiment yields the following confusion matrix:

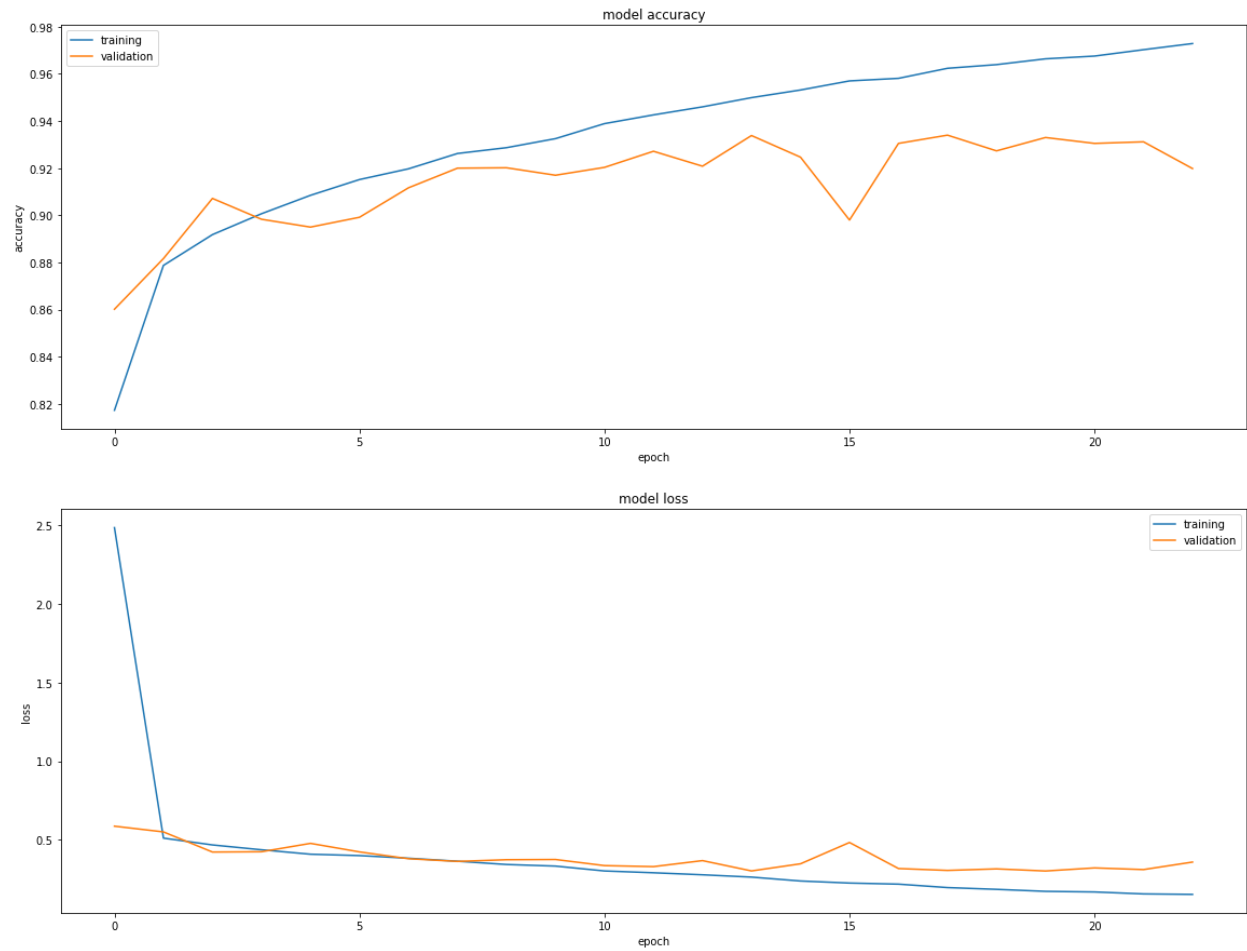
T-shirt/top	885	0	22	12	2	1	75	0	3	0
Trouser	1	981	0	14	1	0	2	0	1	0
Pullover	13	1	909	6	49	0	22	0	0	0
Dress	12	1	11	938	19	0	19	0	0	0
Coat	0	0	36	23	894	0	47	0	0	0
Sandal	0	0	0	0	0	970	0	20	0	10
Shirt	94	0	67	28	63	0	745	0	3	0
Sneaker	0	0	0	0	0	3	0	989	0	8
Bag	4	1	0	3	1	1	6	1	983	0
Ankle boot	0	0	0	0	0	2	0	37	1	960
	T-shirt/top	Trouser	Pullover	Dress	Coat	Sandal	Shirt	Sneaker	Bag	Ankle boot

As usual, the model shows signs of shirt misclassification.

Experiment 6: CNN with 3 convolution/max pooling layers & 2 fully connected layers – Change Dropout (0.2), L2 Regularization (0.01)

Experiment 6 is similar to Experiment 4 but changes L2 regularization learning rate to 0.01.

Training and validation accuracy and cross entropy loss charts are below.



We see validation results start to diverge from the training results after 14 epochs. The resulting test data accuracy is 92.8%, our best-performing model so far.

The experiment yields the following confusion matrix:

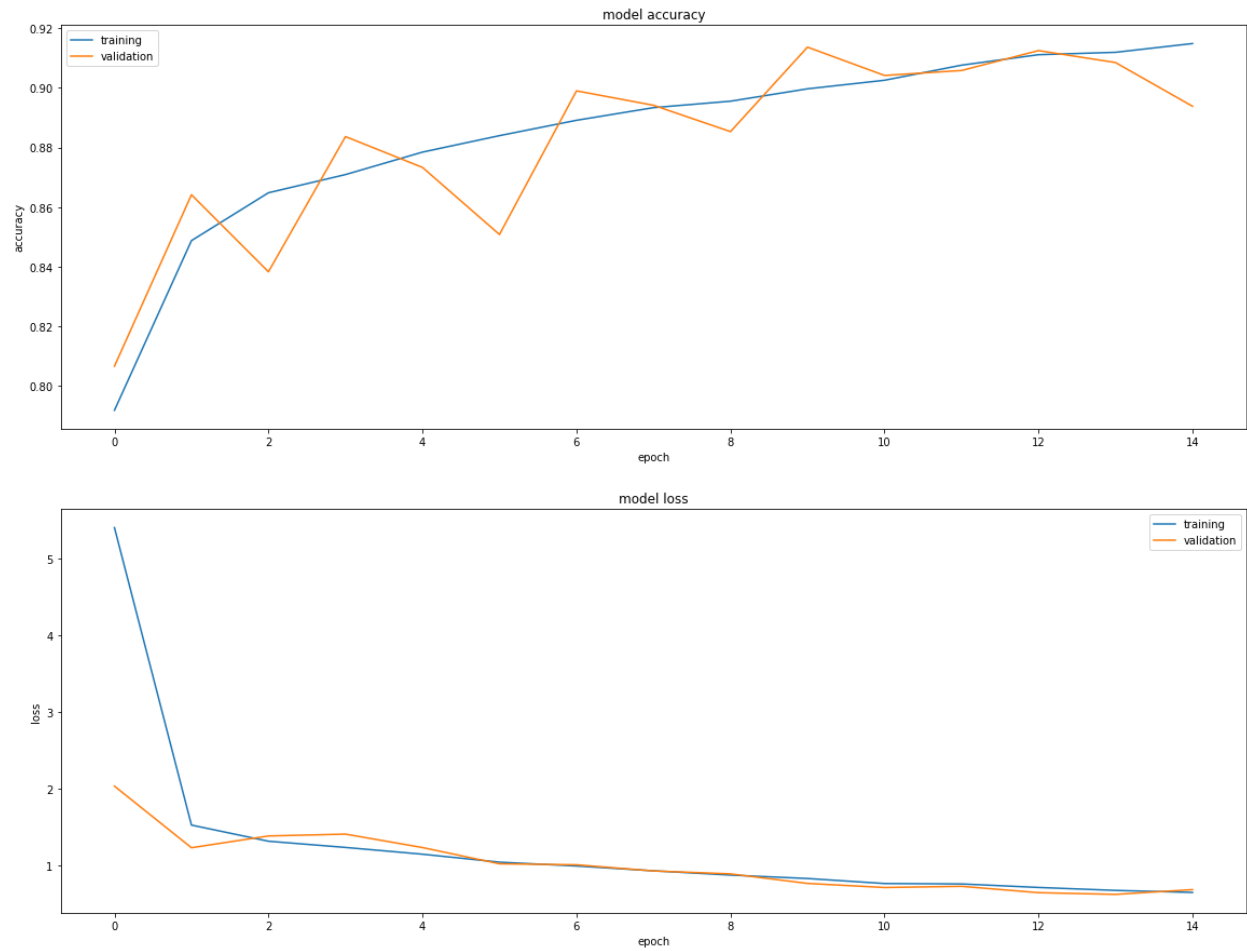
T-shirt/top	865	0	14	34	1	1	82	0	3	0
Trouser	0	985	0	10	2	0	2	0	1	0
Pullover	16	2	875	6	40	0	61	0	0	0
Dress	9	1	4	944	21	1	20	0	0	0
Coat	0	0	18	19	889	0	74	0	0	0
Sandal	0	0	0	0	0	993	0	5	0	2
Shirt	86	1	33	35	33	1	809	0	2	0
Sneaker	0	0	0	0	0	4	0	957	0	39
Bag	1	0	1	5	1	6	5	0	981	0
Ankle boot	0	0	0	0	0	5	1	12	0	982
	T-shirt/top	Trouser	Pullover	Dress	Coat	Sandal	Shirt	Sneaker	Bag	Ankle boot

Shirt classification is somewhat improved compared to previous models, but still underperforms compared to other categories.

Experiment 7: CNN with 3 convolution/max pooling layers & 2 fully connected layers – Change L2 Regularization (0.1)

Experiment 7 is similar to Experiment 6 but changes the L2 regularization learning rate to 0.1.

Training and validation accuracy and cross entropy loss charts are below.



We see validation results start to diverge from the training results after 12 epochs. The resulting test data accuracy is 91.9%. Compared to Experiment 6, our best-performing model so far, Experiment 7's accuracy is off by about 1%.

The experiment yields the following confusion matrix:

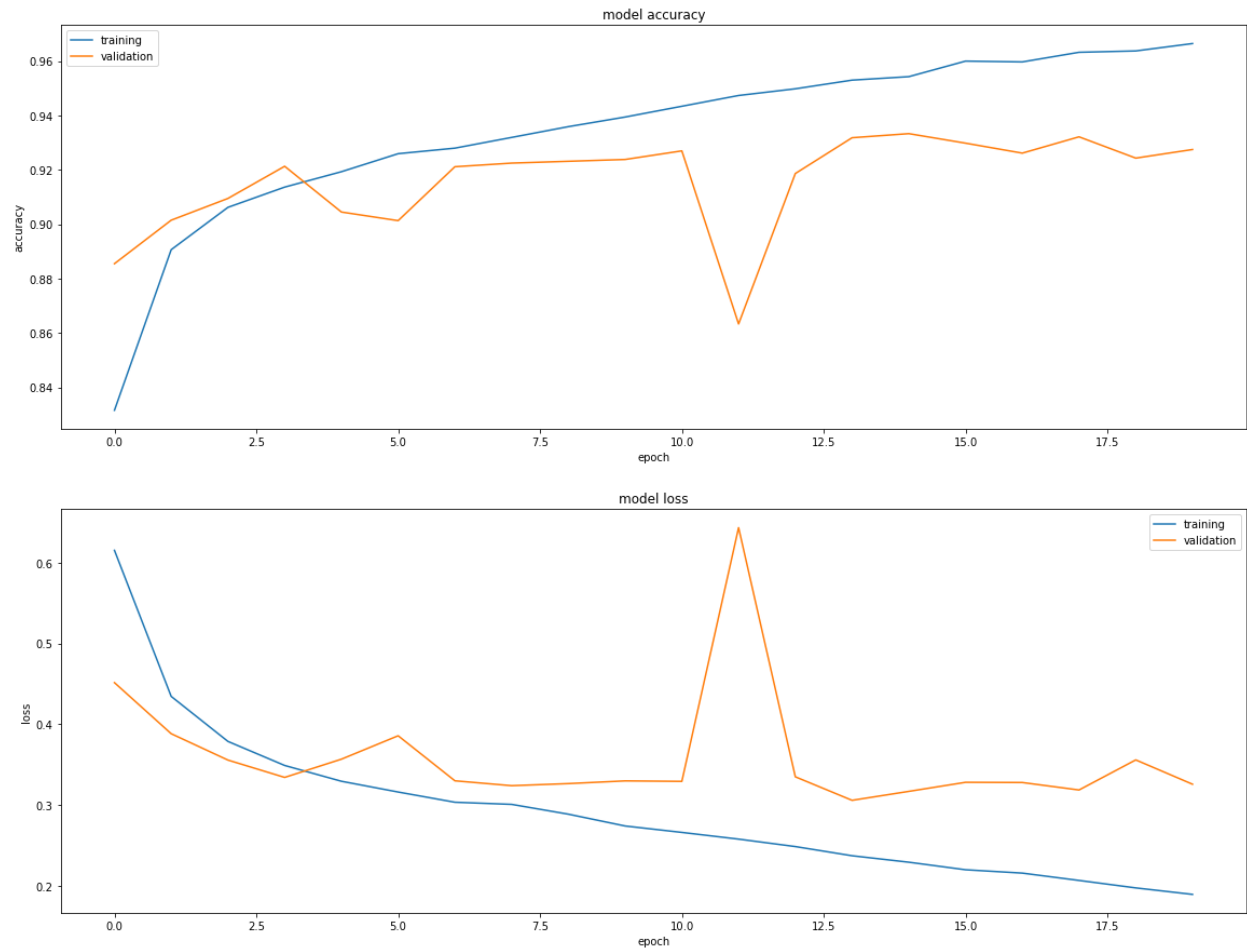
T-shirt/top	852	0	13	27	13	1	91	0	3	0
Trouser	0	979	0	13	6	0	1	0	1	0
Pullover	12	1	801	7	131	0	48	0	0	0
Dress	12	1	8	914	49	0	15	0	1	0
Coat	0	0	28	12	935	0	25	0	0	0
Sandal	0	0	0	0	0	986	0	12	0	2
Shirt	110	2	61	28	119	0	679	0	1	0
Sneaker	0	0	0	0	0	6	0	990	0	4
Bag	0	0	4	3	11	3	1	0	978	0
Ankle boot	0	1	0	0	0	11	0	60	0	928
	T-shirt/top	Trouser	Pullover	Dress	Coat	Sandal	Shirt	Sneaker	Bag	Ankle boot

Experiment 7 seems to have more difficulty classifying shirts compared to previous models.

Experiment 8: CNN with 3 convolution/max pooling layers & 2 fully connected layers – Change L2 Regularization (0.0001)

Experiment 8 is similar to Experiment 6 but changes the L2 regularization learning rate to 0.0001.

Training and validation accuracy and cross entropy loss charts are below.



We see validation results start to diverge from the training results after about 9 epochs. The resulting test data accuracy is 92.75%, which misses out Experiment 6's accuracy by just 0.05%.

The experiment yields the following confusion matrix:

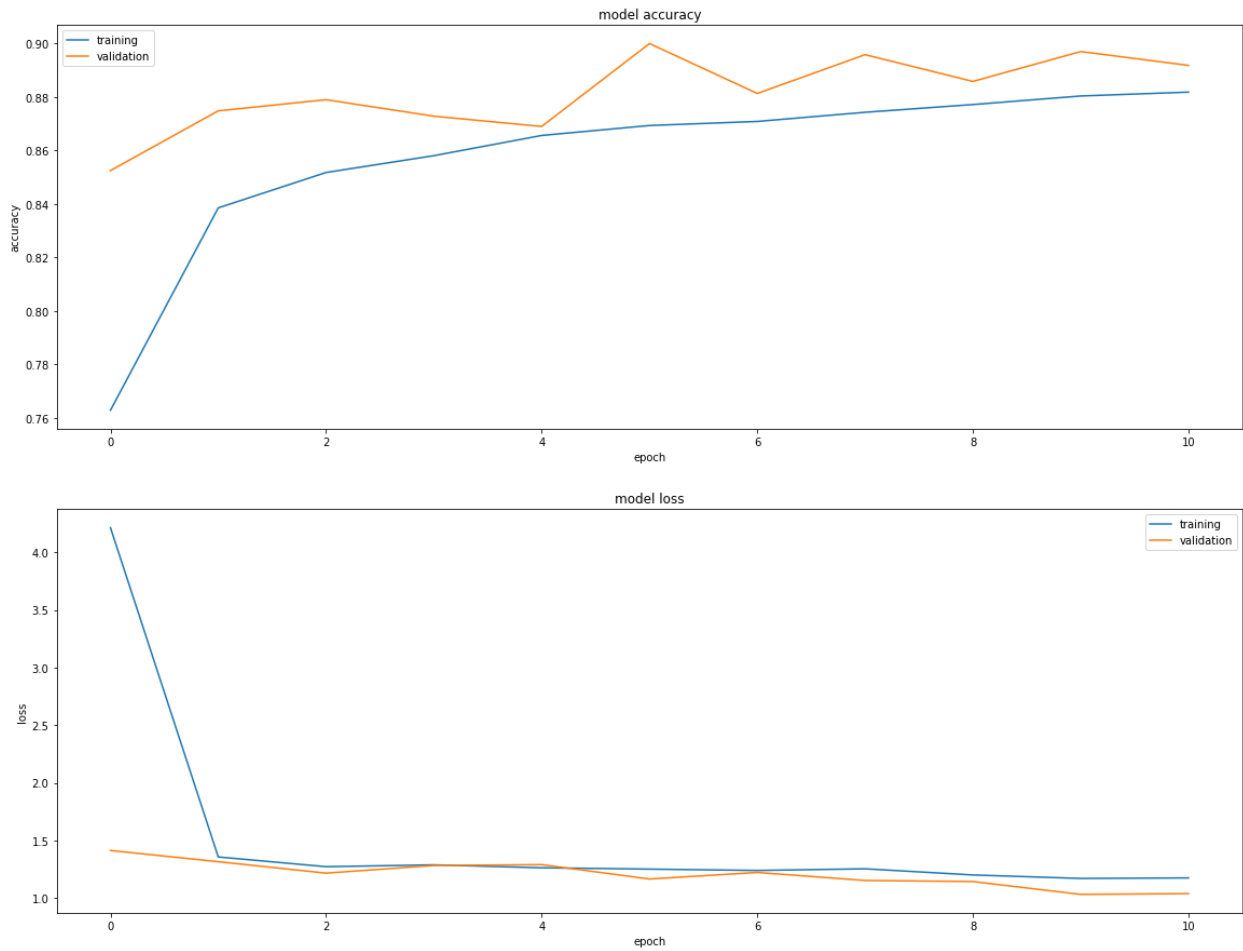
true label	T-shirt/top	898	0	21	8	4	1	68	0	0	0
	Trouser	0	986	0	11	1	0	1	0	1	0
	Pullover	12	2	893	4	37	0	51	0	1	0
	Dress	21	1	12	911	33	0	22	0	0	0
	Coat	0	1	40	8	884	0	67	0	0	0
	Sandal	0	0	0	0	0	978	0	16	0	6
	Shirt	90	1	46	13	50	0	797	0	3	0
	Sneaker	0	0	0	0	0	1	0	991	0	8
	Bag	3	1	0	3	1	4	2	1	984	1
	Ankle boot	0	0	0	0	0	5	0	42	0	953
		predicted label									

Shirt misclassification is still somewhat present in this model, but not as bad compared to most of the previous models.

Experiment 9: CNN with 3 convolution/max pooling layers & 2 fully connected layers – Change L2 Regularization (0.01), Dropout (0.5)

Experiment 9 is similar to Experiment 6, but changes the dropout rate to 0.5.

Training and validation accuracy and cross entropy loss charts are below.



While we see validation results keep up fairly well with training results, the whole training process only lasts for 10 epochs. The resulting test accuracy is 89.3%, well off our best test accuracy of 92.8% by 2.5%. We find that using larger dropout rates in this CNN model architecture does not result in better test accuracy.

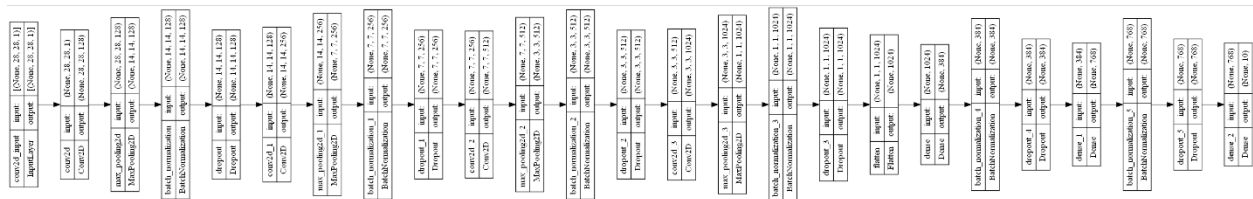
The experiment yields the following confusion matrix:

true label	T-shirt/top	958	1	7	6	3	1	18	0	6	0
	Trouser	2	981	0	10	4	0	1	0	2	0
	Pullover	29	1	876	7	35	0	52	0	0	0
	Dress	40	13	6	870	62	0	7	0	2	0
	Coat	1	1	88	16	830	0	63	0	1	0
	Sandal	0	0	0	0	0	967	0	28	0	5
	Shirt	274	1	61	25	101	0	532	0	6	0
	Sneaker	0	0	0	0	0	3	0	950	0	47
	Bag	4	1	0	2	1	1	4	1	984	2
	Ankle boot	0	0	0	0	0	4	1	17	0	978
		predicted label									

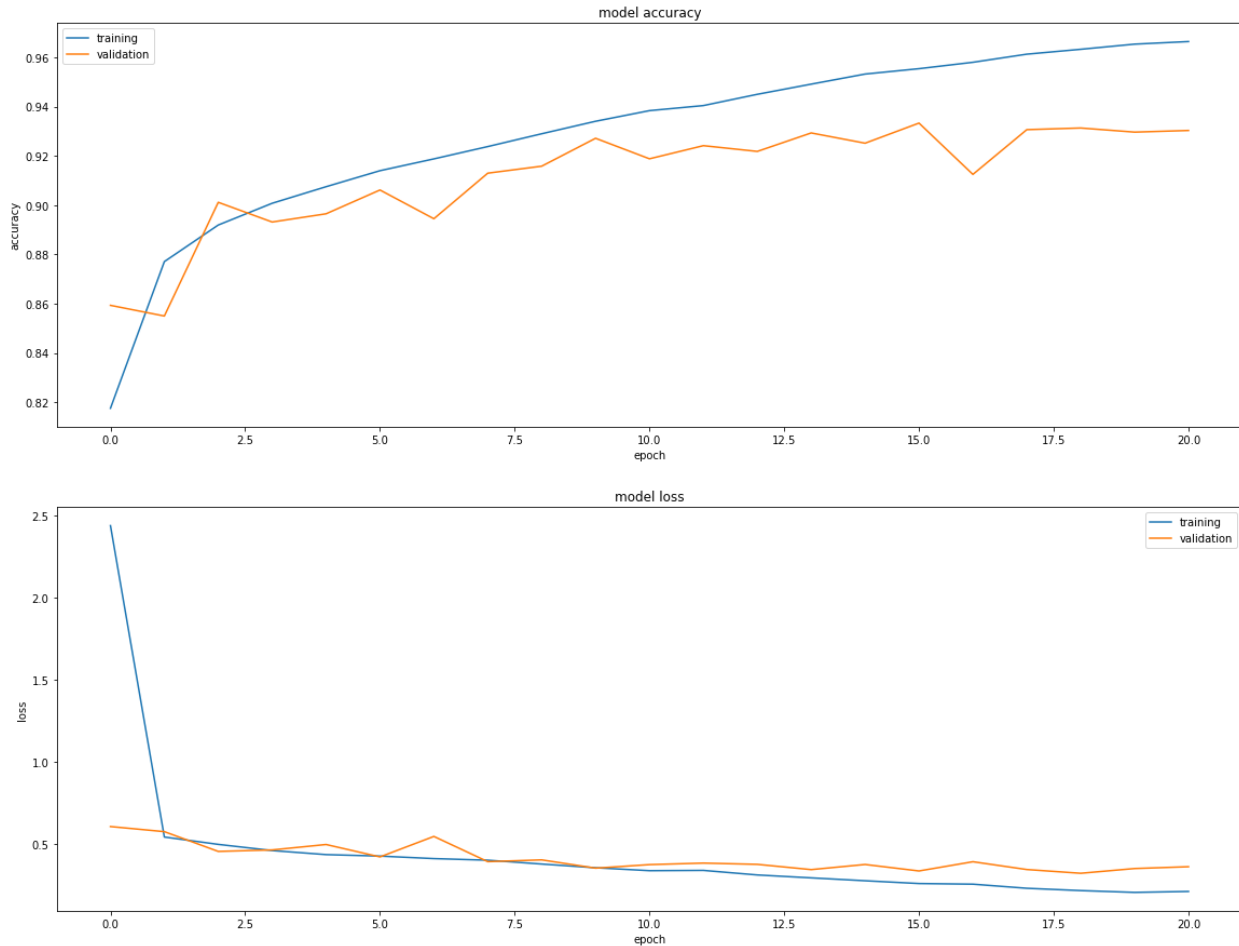
We see that using a larger dropout rate results in a worse shirt categorization process.

Experiment 10: CNN with 4 convolution/max pooling layers & 2 fully connected layers

Experiment 10 is similar to Experiment 6, our best-performing model so far, but adds another 3x3 convolutional layer which uses 1024 filters. A diagram of the full architecture is below.



Training and validation accuracy and cross entropy loss charts are below.



We see validation results start to diverge from the training results after 5 epochs. The resulting test data accuracy is 92.6%, which doesn't exceed the accuracy performance of Experiment 6, despite the added convolutional layer. Since Fashion-MNIST is a dataset which contains small images, there could be a ceiling of convolutional or max pooling layers, and if exceeded, could potentially lose important feature information (Nocentini et al. 2022). Adding an additional convolutional and max pooling layer in this case shows diminishing returns despite the added model complexity.

The experiment yields the following confusion matrix:

T-shirt/top	828	0	20	34	3	0	112	0	3	0
Trouser	0	979	0	15	3	0	1	0	2	0
Pullover	8	1	883	7	66	0	35	0	0	0
Dress	6	0	10	945	27	0	12	0	0	0
Coat	0	0	25	16	925	0	34	0	0	0
Sandal	0	0	0	0	0	970	0	18	1	11
Shirt	55	0	54	34	62	0	791	0	4	0
Sneaker	0	0	0	0	0	1	0	987	0	12
Bag	0	0	1	6	1	2	2	0	987	1
Ankle boot	0	0	1	0	0	0	0	31	0	968
	T-shirt/top	Trouser	Pullover	Dress	Coat	Sandal	Shirt	Sneaker	Bag	Ankle boot

Like in previous experiments, this model shows some signs of shirt misclassification.

The results of our convolutional model experiments in this section are summarized in the table below, with the best one highlighted (Experiment 6).

Experiment	Conv2d Layers	DNNs	L2 Reg	Batch Norm	Dropout	Test Accuracy	Test Loss	Precision	Recall	F1-Score	RMSE
1	3	2	None	None	None	0.919	0.2288	0.92	0.92	0.92	1.124
2	3	2	<u>0.001</u>	None	None	0.9183	0.2666	0.92	0.92	0.92	1.096
3	3	2	0.001	<u>Yes</u>	<u>0.3</u>	0.9212	0.4229	0.92	0.92	0.92	1.045
4	3	2	0.001	Yes	<u>0.2</u>	0.9245	0.331	0.92	0.92	0.92	1.03
5	3	2	0.001	Yes	<u>0.1</u>	0.9254	0.3061	0.93	0.93	0.93	1.022
6	3	2	<u>0.01</u>	Yes	0.2	0.928	0.3182	0.93	0.93	0.93	1.019
7	3	2	<u>0.1</u>	Yes	0.2	0.9188	0.4982	0.91	0.9	0.9	1.132
8	3	2	<u>0.0001</u>	Yes	0.2	0.9275	0.332	0.93	0.93	0.93	0.995
9	3	2	0.01	Yes	<u>0.5</u>	0.8926	1.0441	0.89	0.89	0.89	1.272
10	<u>4</u>	2	0.01	Yes	0.2	0.9263	0.3276	0.93	0.93	0.93	0.983

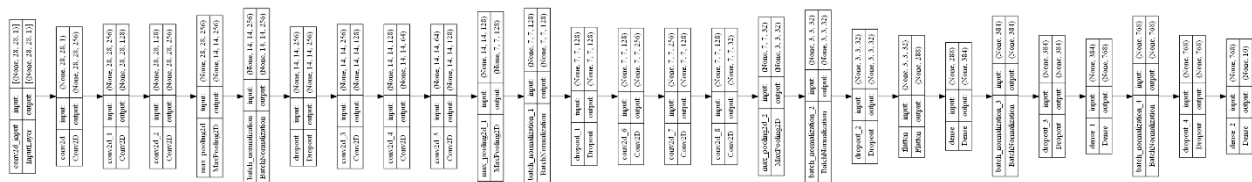
We were able to achieve a test accuracy of 92.8% in our group of convolutional models by tuning various hyperparameters, in which Experiment 6 was our most optimal model. While the last experiment in this section added another convolution layer we started to see diminishing returns as its test accuracy did not exceed Experiment 6.

Section Two: Grouped Convolutional Models (3x3 Stacks)

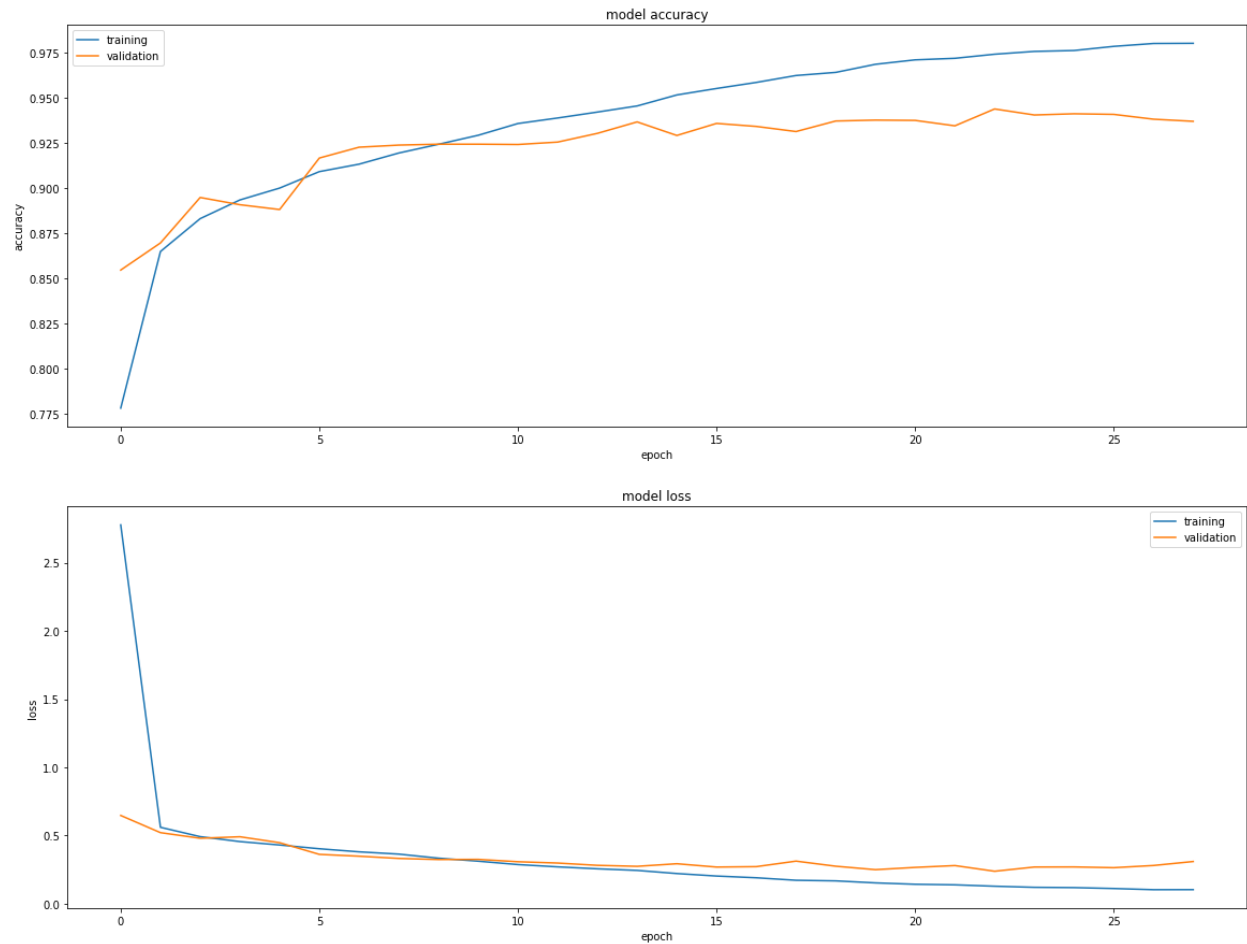
Experiment 1: CNN with 3 groups of 3 convolution layers & 2 fully connected layers – Add Batch Normalization, Dropout (0.2), L2 Regularization (0.01)

We model a CNN consisting of 3 groups of 3x3 convolutional layers, with each group separated by a 2x2 max pooling layer. The first group contains layers with 256, 128, and 256 filters respectively, the second group contains layers with 128, 64, and 128 filters respectively, and the third group contains layers with 256, 128, and 32 filters respectively. The convolutional base then feeds into a DNN containing 2 fully connected layers, with the first layer containing 384 units and the second layer containing 768 units, followed by the 10-way softmax classification output layer. Batch normalization and dropout (0.2) processes follow all max pooling layers and fully connected layers. All fully connected layers utilize L2 regularization configured with a learning rate of 0.01.

The model diagram below shows the full architecture used by all models in this section.



Training and validation accuracy and cross entropy loss charts are below.



We see validation results start to diverge from the training results after 7 epochs. The resulting test data accuracy is 93.7%, our best-performing model so far, outperforming all models in the first section.

The experiment yields the following confusion matrix:

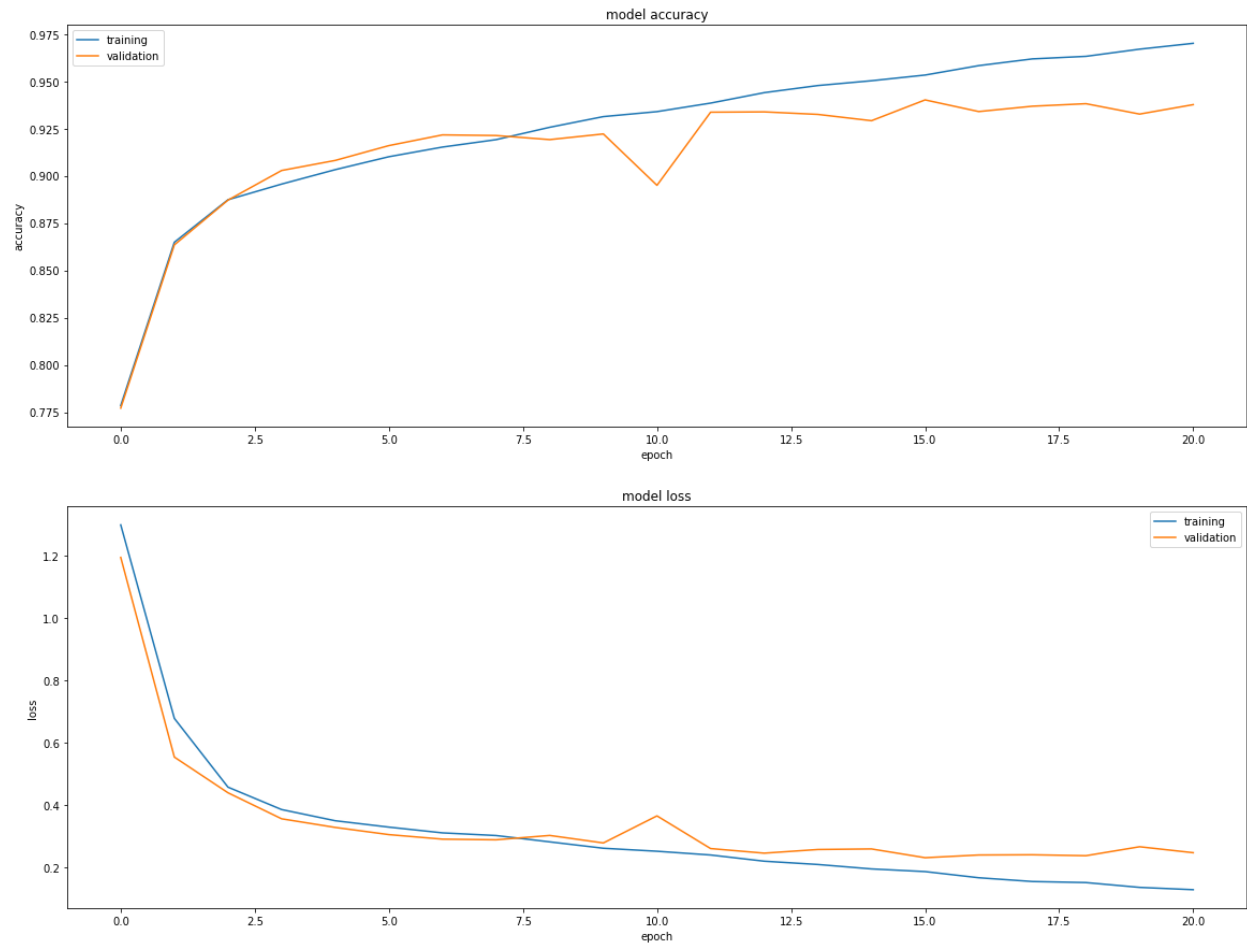
T-shirt/top	865	1	15	12	1	0	105	0	1	0
Trouser	0	992	0	5	2	0	0	0	1	0
Pullover	18	1	900	6	32	0	43	0	0	0
Dress	8	0	6	946	16	0	24	0	0	0
Coat	0	0	18	21	923	0	38	0	0	0
Sandal	0	0	0	0	0	981	0	11	1	7
Shirt	60	1	37	30	46	0	824	0	2	0
Sneaker	0	0	0	0	0	2	0	988	0	10
Bag	3	0	0	5	3	1	3	0	985	0
Ankle boot	0	0	0	0	0	3	0	30	1	966
	T-shirt/top	Trouser	Pullover	Dress	Coat	Sandal	Shirt	Sneaker	Bag	Ankle boot

While the model is able to classify most garments in the test dataset with well over 90% precision, it still has issues with classifying shirts with 79% precision.

Experiment 2: CNN with 3 groups of 3 convolution layers & 2 fully connected layers – Change L2 Regularization (0.001)

Experiment 2 follows the same architecture as Experiment 1 but changes L2 Regularization to a learning rate of 0.001.

Training and validation accuracy and cross entropy loss charts are below.



We see validation results start to diverge from the training results after about 8 epochs. The resulting test data accuracy is 93.6%, which is just lower than Experiment 1's test accuracy by 0.1%.

The experiment yields the following confusion matrix:

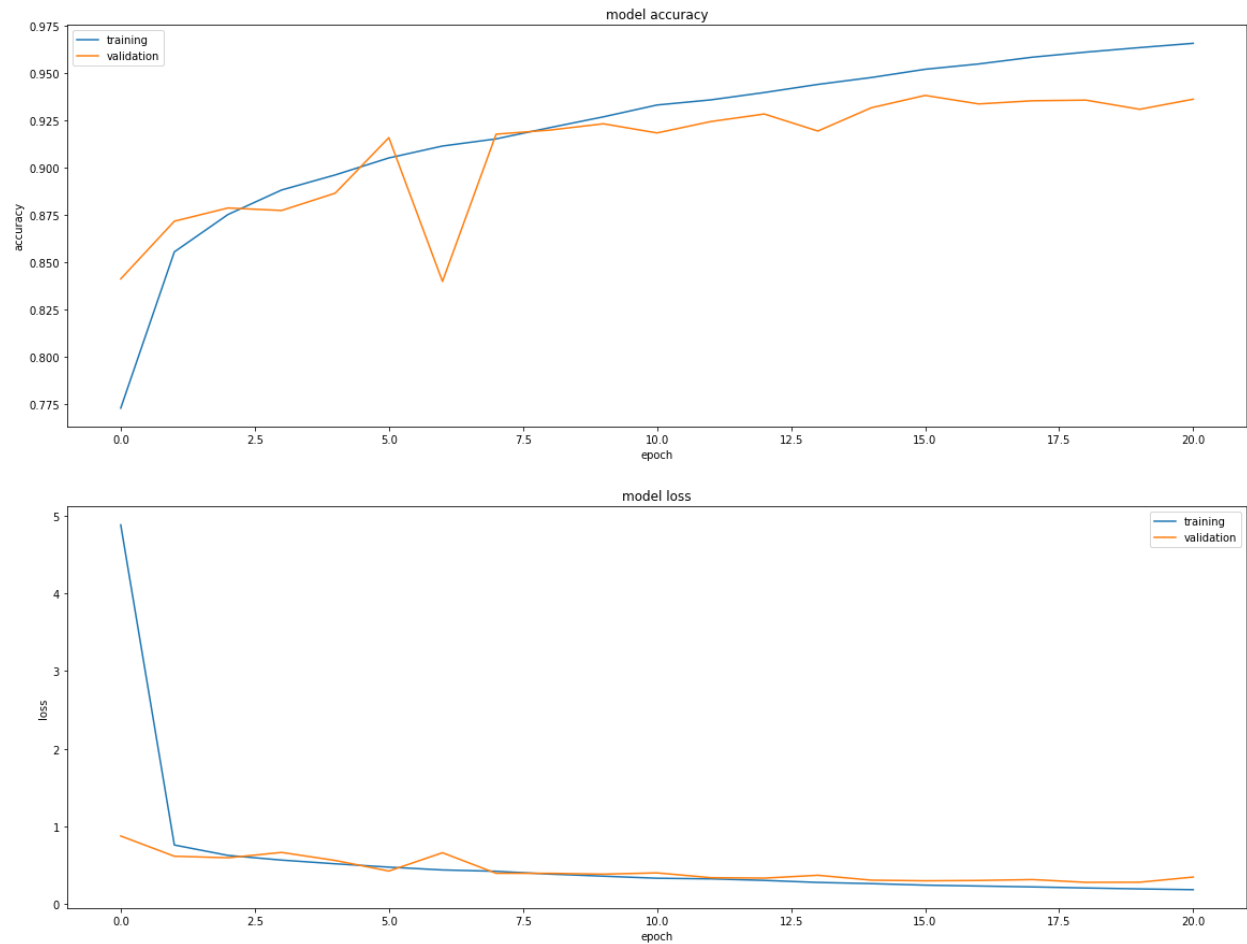
T-shirt/top	855	0	20	11	4	0	105	0	5	0
Trouser	0	987	0	4	2	0	5	0	2	0
Pullover	16	1	909	11	35	0	27	0	1	0
Dress	12	1	6	948	19	0	13	0	1	0
Coat	1	0	14	16	927	0	41	0	1	0
Sandal	0	0	0	0	0	995	0	2	0	3
Shirt	60	0	46	23	56	0	810	0	5	0
Sneaker	0	0	0	0	0	10	0	970	0	20
Bag	2	0	0	3	2	1	0	0	992	0
Ankle boot	0	0	1	0	0	6	0	24	0	969
	T-shirt/top	Trouser	Pullover	Dress	Coat	Sandal	Shirt	Sneaker	Bag	Ankle boot

While the model can correctly classify most images well, like Experiment 1 it too doesn't classify tops and shirt categories compared to other categories.

Experiment 3: CNN with 3 groups of 3 convolution layers & 2 fully connected layers – Change L2 Regularization (0.1)

Experiment 3 follows the same architecture as Experiment 1 but changes L2 Regularization to a learning rate of 0.1.

Training and validation accuracy and cross entropy loss charts are below.



We see validation results start to diverge from the training results after about 9 epochs. The resulting test data accuracy is 92.9%, which is just lower than Experiment 2's test accuracy by about 0.7%.

The experiment yields the following confusion matrix:

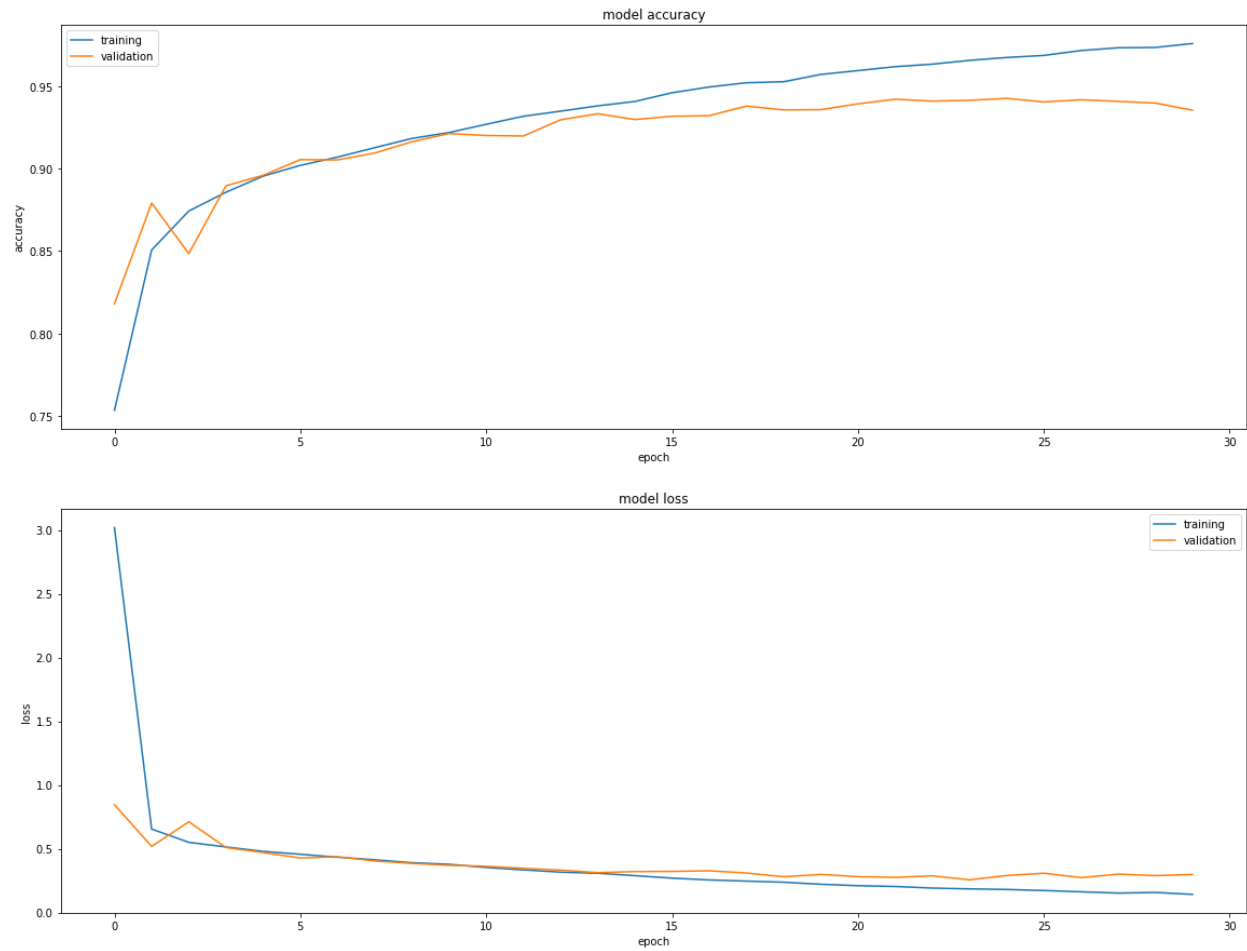
T-shirt/top	916	1	22	11	3	0	40	0	7	0
Trouser	1	983	0	13	2	0	1	0	0	0
Pullover	11	1	897	7	62	0	21	0	1	0
Dress	14	0	6	943	23	0	14	0	0	0
Coat	0	0	13	17	944	0	26	0	0	0
Sandal	0	0	0	0	0	987	0	6	0	7
Shirt	136	1	55	24	63	0	717	0	4	0
Sneaker	0	0	0	0	0	3	0	993	0	4
Bag	1	0	0	4	5	3	3	0	984	0
Ankle boot	0	0	0	0	1	2	1	68	0	928
	T-shirt/top	Trouser	Pullover	Dress	Coat	Sandal	Shirt	Sneaker	Bag	Ankle boot

The theme of shirt misclassification continues in this experiment despite its high overall test accuracy.

Experiment 4: CNN with 3 groups of 3 convolution layers & 2 fully connected layers – Change L2 Regularization (0.2), Dropout (0.3)

Experiment 4 follows the same architecture as Experiment 1 but changes dropout rate to 0.3.

Training and validation accuracy and cross entropy loss charts are below.



We see validation results start to diverge from the training results after about 12 or 13 epochs. The resulting test data accuracy is 93.66%, which is slightly higher than Experiment 2's test accuracy by about 0.04% but still lower than Experiment 1 by 0.04%.

The experiment yields the following confusion matrix:

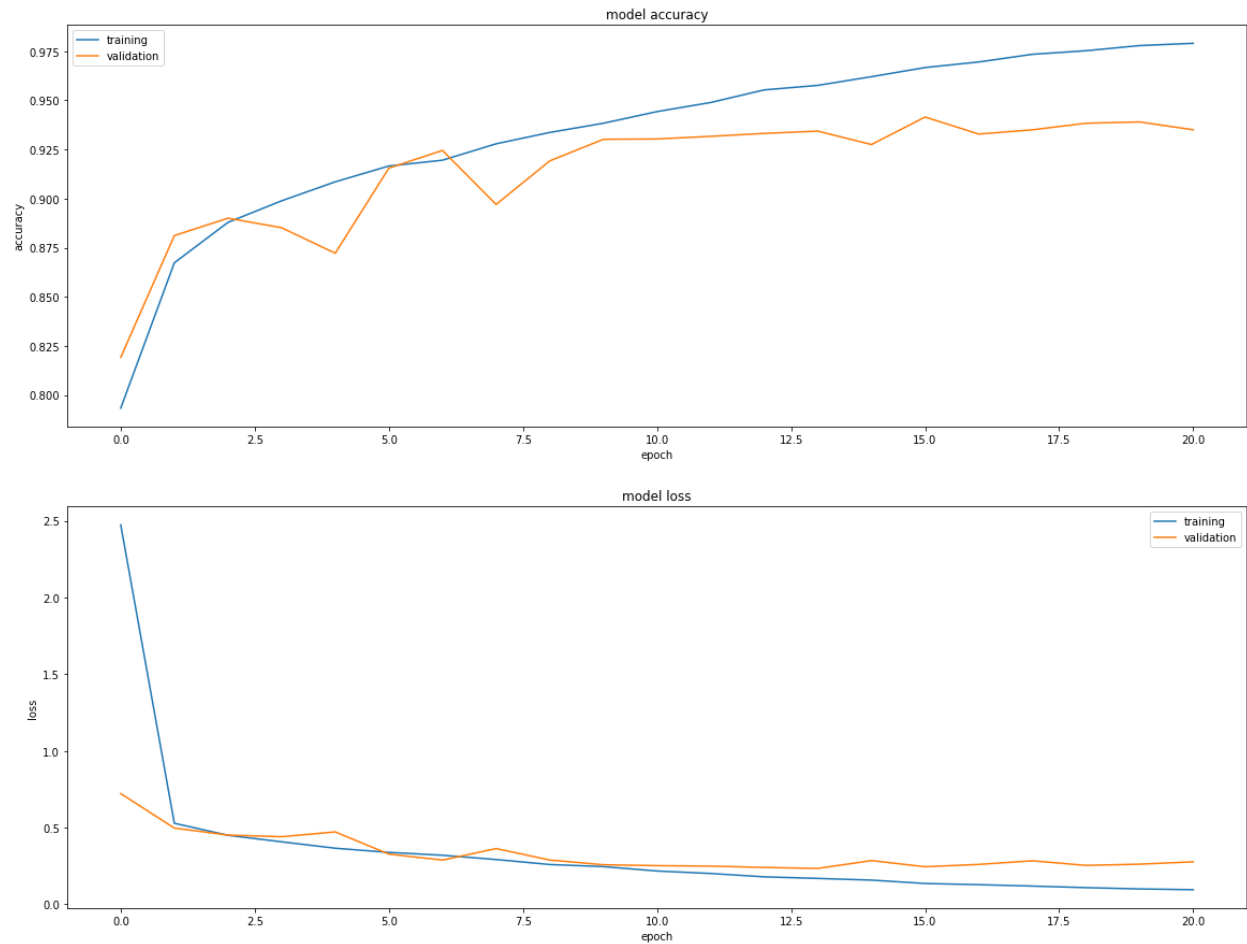
T-shirt/top	885	1	20	8	3	1	81	0	1	0
Trouser	0	989	0	6	2	0	2	0	1	0
Pullover	12	1	935	5	27	0	20	0	0	0
Dress	14	0	14	930	19	0	23	0	0	0
Coat	0	0	35	15	900	0	50	0	0	0
Sandal	0	0	0	0	0	990	0	7	0	3
Shirt	80	0	58	13	42	0	805	0	2	0
Sneaker	0	0	0	0	0	1	0	973	0	26
Bag	4	0	1	3	4	1	3	0	984	0
Ankle boot	0	0	1	0	0	4	0	20	0	975
	T-shirt/top	Trouser	Pullover	Dress	Coat	Sandal	Shirt	Sneaker	Bag	Ankle boot

The model can classify most categories well but still has some trouble categorizing shirts.

Experiment 5: CNN with 3 groups of 3 convolution layers & 2 fully connected layers – Change Dropout (0.1)

Experiment 5 follows the same architecture as Experiment 1 but changes dropout rate to 0.1.

Training and validation accuracy and cross entropy loss charts are below.



We see validation results start to diverge from the training results after 7 epochs. The resulting test data accuracy is 93%, which is the worst-performing model in this section but still out-performs all models in Section One.

The experiment yields the following confusion matrix:

T-shirt/top	864	0	29	26	1	1	76	0	3	0
Trouser	0	983	0	11	2	0	2	0	2	0
Pullover	8	1	923	6	30	0	32	0	0	0
Dress	5	3	8	953	15	0	16	0	0	0
Coat	0	0	38	23	877	0	62	0	0	0
Sandal	0	0	0	1	0	988	0	7	0	4
Shirt	80	0	50	29	39	0	797	0	5	0
Sneaker	0	0	0	0	0	4	0	974	0	22
Bag	1	1	3	2	3	3	4	1	981	1
Ankle boot	0	0	1	0	0	4	0	31	0	964
	T-shirt/top	Trouser	Pullover	Dress	Coat	Sandal	Shirt	Sneaker	Bag	Ankle boot

Like the previous models, it too has some trouble classifying shirt images.

The results of our grouped convolutional model experiments in this section are summarized in the table below, with the best one highlighted (Experiment 1).

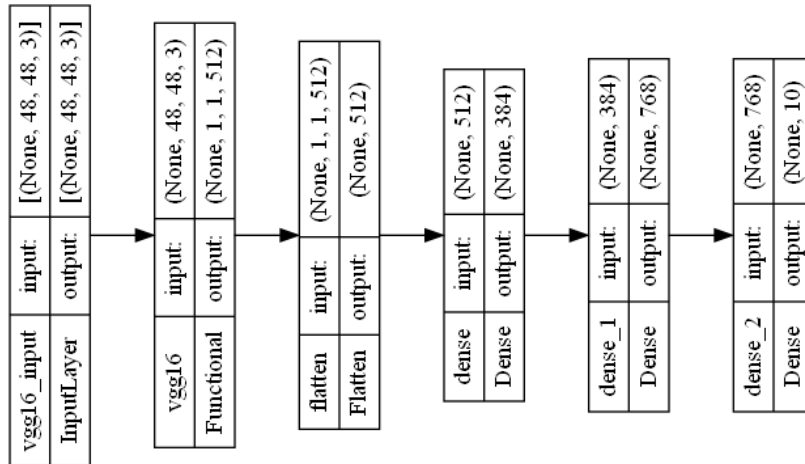
Experiment	Conv2d Layers	DNNs	L2 Reg	Batch Norm	Dropout	Test Accuracy	Test Loss	Precision	Recall	F1-Score	RMSE
1	9 (3x3 Groups)	2	0.01	Yes	0.2	0.937	0.2698	0.94	0.94	0.94	0.978
2	9 (3x3 Groups)	2	<u>0.001</u>	Yes	0.2	0.9362	0.2472	0.94	0.94	0.94	0.992
3	9 (3x3 Groups)	2	<u>0.1</u>	Yes	0.2	0.9292	0.3083	0.93	0.93	0.93	1.021
4	9 (3x3 Groups)	2	0.01	Yes	<u>0.3</u>	0.9366	0.2799	0.94	0.94	0.94	0.974
5	9 (3x3 Groups)	2	0.01	Yes	<u>0.1</u>	0.9304	0.2545	0.93	0.93	0.93	0.988

While we tuned dropout and L2 regularization hyperparameters above and below the values set in Experiment 1, none of the models were able to exceed its test accuracy performance. We found Experiment 1 to have the most optimal hyperparameter tuning settings and was able to achieve the best-performing model so far. All models in this section outperformed all models in the first section.

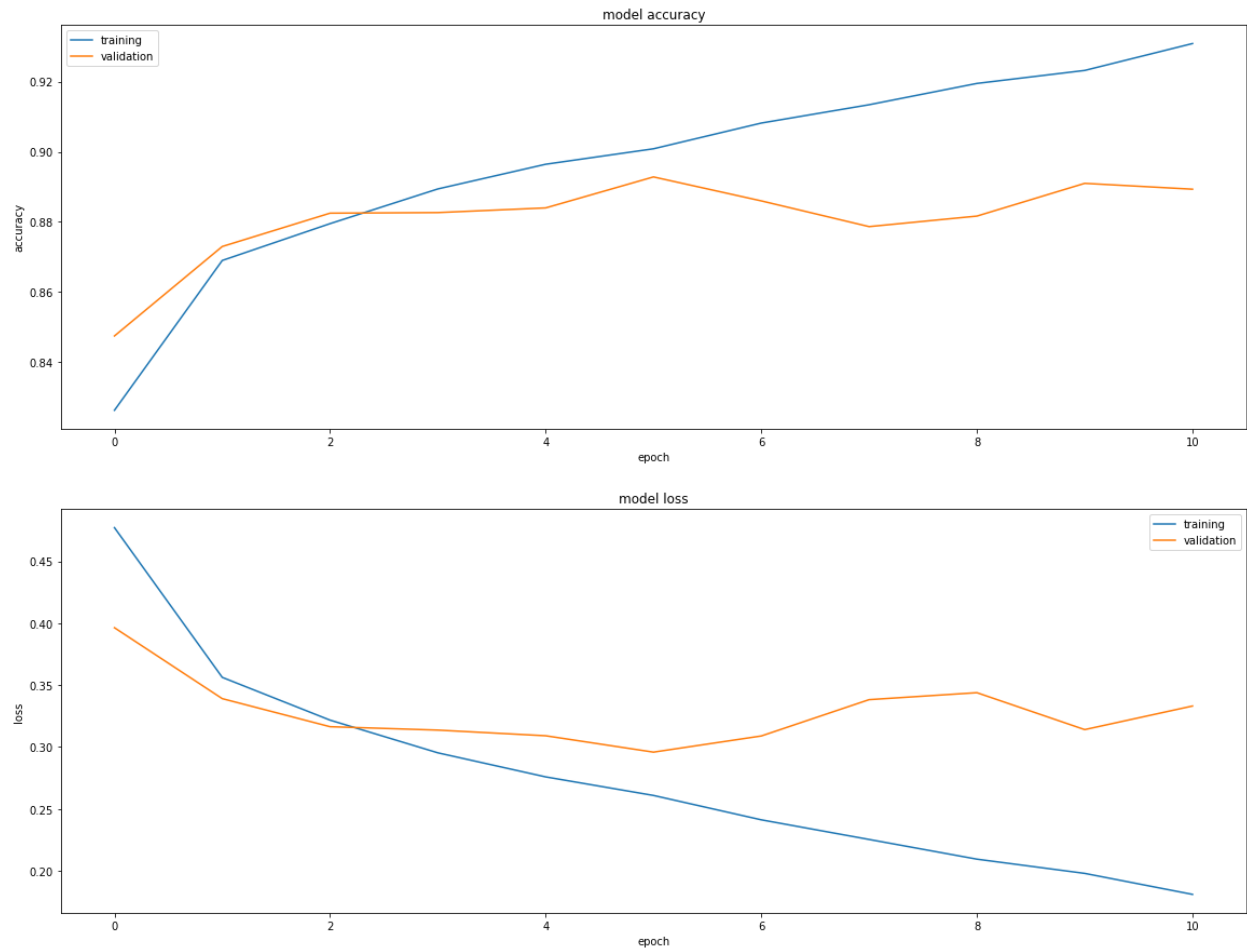
Section 3: Pre-Trained Models

Experiment 1: VGG16 with 2 fully connected layers

The VGG16 pre-trained model base feeds into a DNN containing 2 fully connected layers, with the first layer containing 384 units and the second layer containing 768 units, followed by the 10-way softmax classification output layer.



Training and validation accuracy and cross entropy loss charts are below.



We see validation results start to diverge from the training results after 5 epochs. The resulting test data accuracy is 88.1%, over 5% worse than our best-performing model that does not use a pre-trained model base.

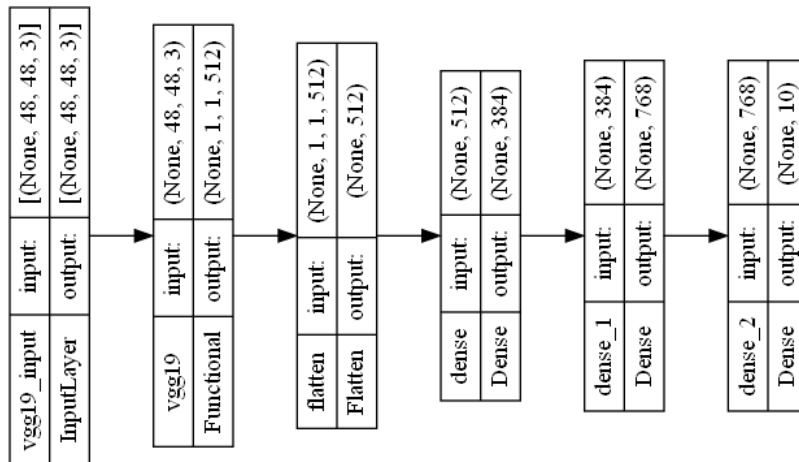
The experiment yields the following confusion matrix:

true label	T-shirt/top	847	0	6	32	7	0	97	0	11	0
	Trouser	1	966	3	25	1	0	1	0	3	0
	Pullover	18	2	791	16	93	0	77	0	3	0
	Dress	20	4	12	903	27	0	34	0	0	0
	Coat	6	1	62	42	819	0	69	0	1	0
	Sandal	0	0	0	0	0	939	0	40	4	17
	Shirt	137	1	60	46	90	0	655	0	11	0
	Sneaker	0	0	0	0	0	14	0	959	1	26
	Bag	4	0	2	2	3	3	7	1	977	1
	Ankle boot	0	0	1	0	0	3	0	45	1	950
		predicted label									

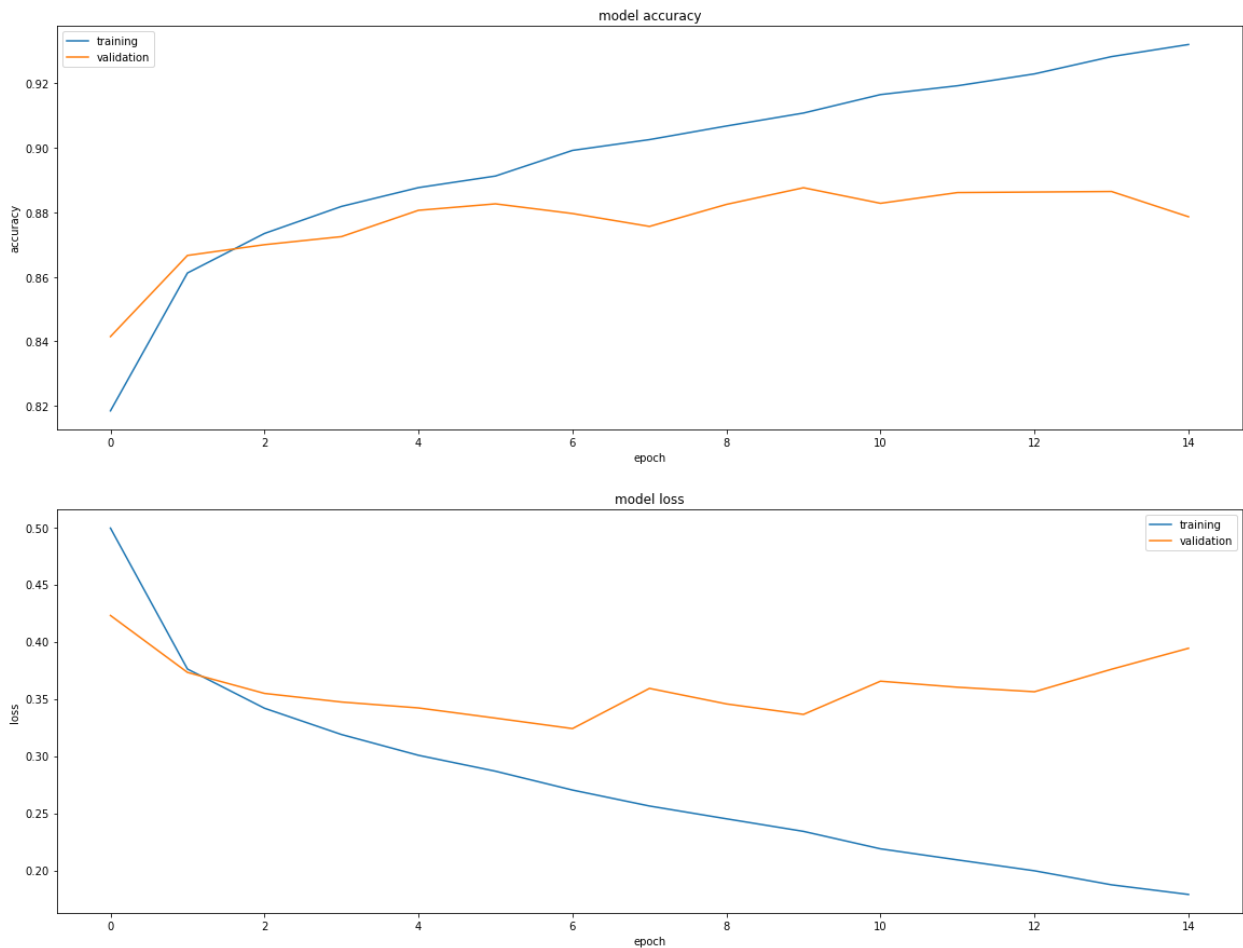
While the model exhibits the common theme of shirt misclassification, we see the model also has difficulties classifying between coats and pullovers.

Experiment 2: VGG19 with 2 fully connected layers

Like the VGG16 model above, the VGG19 pre-trained base feeds into a DNN containing 2 fully connected layers, with the first layer containing 384 units and the second layer containing 768 units, followed by the 10-way softmax classification output layer.



Training and validation accuracy and cross entropy loss charts are below.



We see validation results start to diverge from the training results after 5 epochs. The resulting test data accuracy is 87.3%, 0.7% worse than the less complex VGG16-based model.

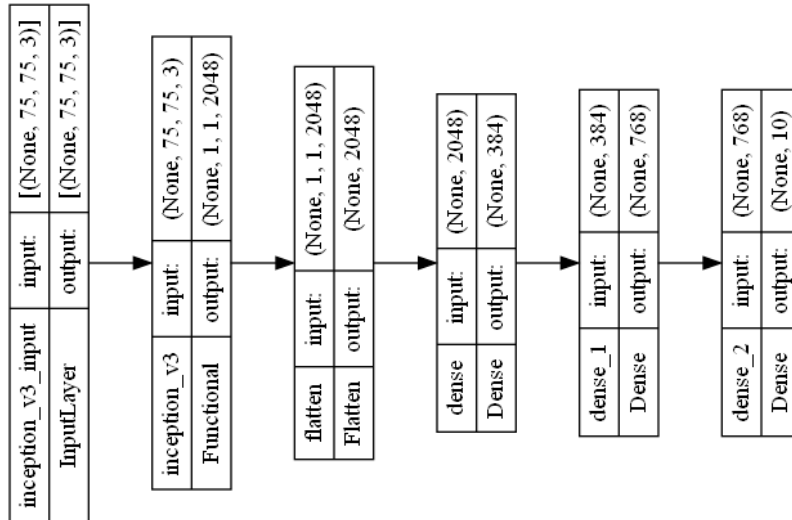
The experiment yields the following confusion matrix:

T-shirt/top	818	2	11	47	8	2	101	0	11	0
Trouser	0	967	2	26	1	0	1	0	3	0
Pullover	13	2	774	10	117	0	81	0	3	0
Dress	20	10	6	865	56	0	39	0	4	0
Coat	4	1	56	32	840	0	64	0	3	0
Sandal	0	0	0	1	0	937	0	48	0	14
Shirt	137	0	44	43	100	0	665	0	11	0
Sneaker	0	0	0	0	0	11	0	978	0	11
Bag	2	1	2	10	2	10	7	2	957	7
Ankle boot	0	0	1	0	0	6	0	60	0	933
	T-shirt/top	Trouser	Pullover	Dress	Coat	Sandal	Shirt	Sneaker	Bag	Ankle boot

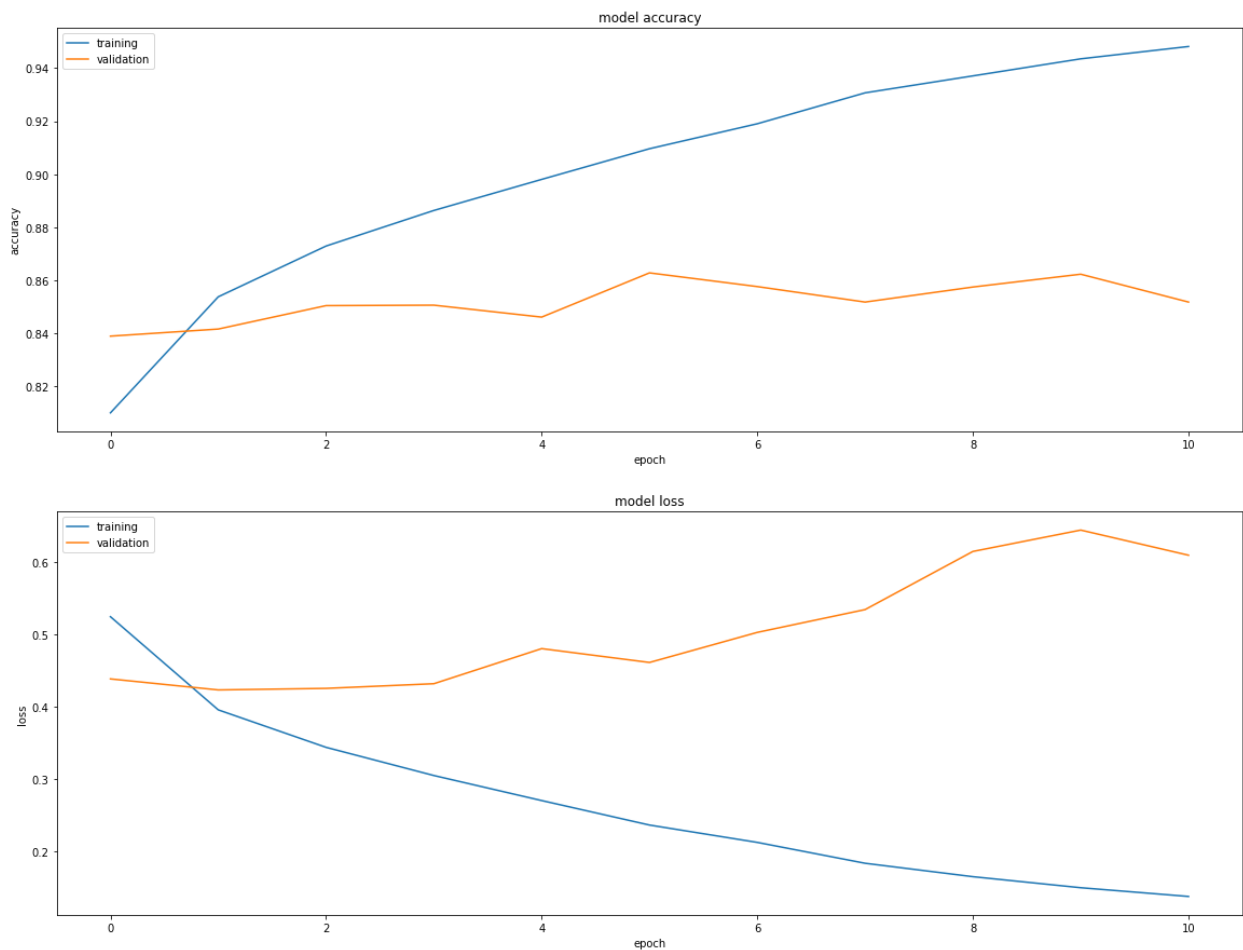
Along with shirt misclassification, VGG19 also has difficulties classifying between coats, dresses, and pullovers, having more trouble than its smaller VGG16 counterpart.

Experiment 3: Inception V3 with 2 fully connected layers

The Inception V3 pre-trained base feeds into a DNN containing 2 fully connected layers, with the first layer containing 384 units and the second layer containing 768 units, followed by the 10-way softmax classification output layer.



Training and validation accuracy and cross entropy loss charts are below.



We see validation results already start to diverge from the training results just after 2 or 3 epochs. The resulting test data accuracy is 84.3%, 3% worse than the VGG19 model.

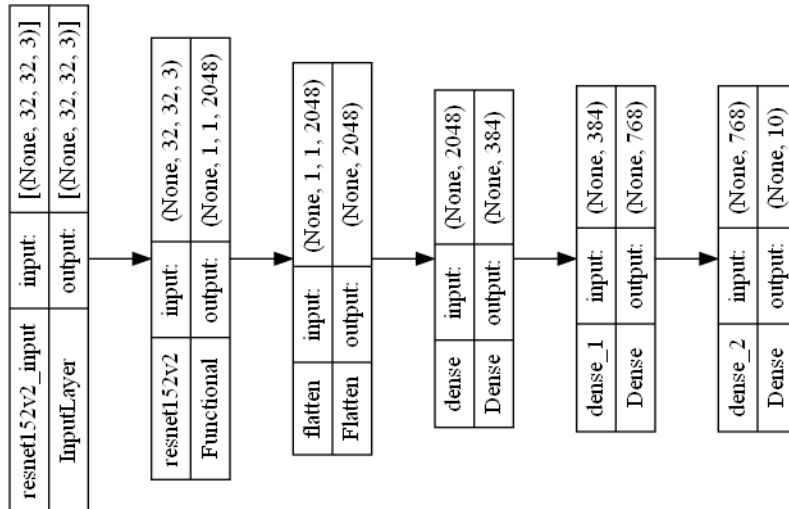
The experiment yields the following confusion matrix:

true label	T-shirt/top	Trouser	Pullover	Dress	Coat	Sandal	Shirt	Sneaker	Bag	Ankle boot
	744	2	21	70	24	2	119	0	18	0
	1	963	3	23	3	2	3	0	2	0
	11	4	758	8	123	0	86	0	10	0
	20	13	9	781	127	0	48	0	1	1
	3	0	57	21	857	2	54	0	6	0
	1	0	0	0	0	943	0	42	5	9
	92	2	50	59	218	2	564	0	13	0
	0	0	0	1	0	39	0	929	1	30
	2	1	8	4	5	5	7	0	965	3
	2	0	0	1	0	19	0	46	9	923
predicted label										

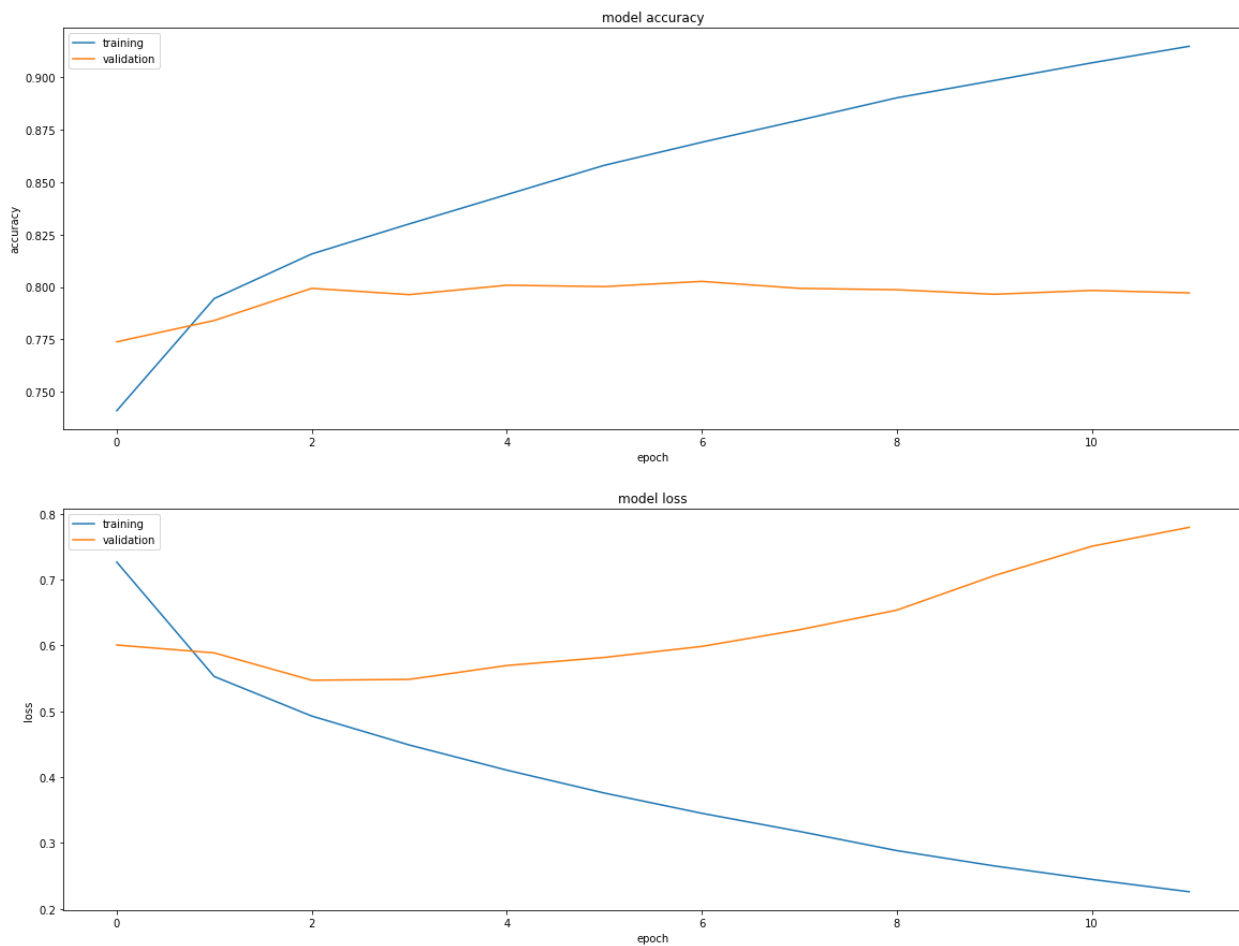
The Inception-based model has some trouble classifying any article of clothing worn on the upper body, performing worse than the VGG-based models.

Experiment 4: ResNet152 V3 with 2 fully connected layers

The ResNet152 V3 pre-trained base feeds into a DNN containing 2 fully connected layers, with the first layer containing 384 units and the second layer containing 768 units, followed by the 10-way softmax classification output layer.



Training and validation accuracy and cross entropy loss charts are below.



We see validation results already start to diverge from the training results just after 2 or 3 epochs. The resulting test data accuracy is 79.3%, our worst overall model.

The experiment yields the following confusion matrix:

T-shirt/top	770	3	18	69	11	3	92	0	33	1
Trouser	13	922	7	43	6	0	6	0	3	0
Pullover	47	5	603	20	155	0	148	0	22	0
Dress	43	21	29	770	63	2	50	2	20	0
Coat	10	4	86	47	675	0	169	0	9	0
Sandal	0	0	0	3	1	898	0	69	10	19
Shirt	185	8	64	59	135	2	516	0	31	0
Sneaker	0	0	0	0	0	30	1	945	5	19
Bag	13	0	5	6	18	4	30	10	913	1
Ankle boot	0	0	0	0	0	18	2	66	1	913
	T-shirt/top	Trouser	Pullover	Dress	Coat	Sandal	Shirt	Sneaker	Bag	Ankle boot

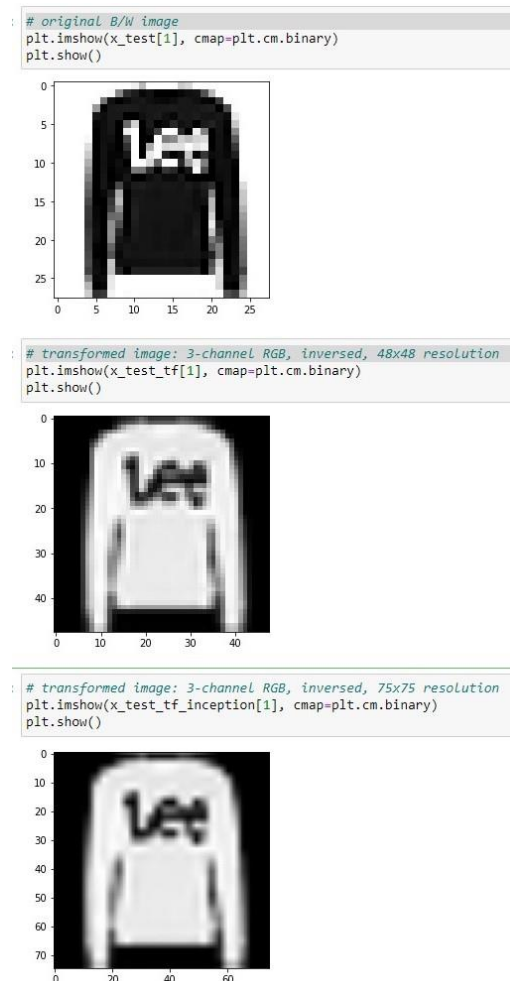
The ResNet-based pre-trained model correctly classifies shirts just over half of the time, which is not great. It performs worst overall out of all the models we've tested, as it has trouble categorizing upper-body garments in general.

The results of our pre-trained model experiments are summarized in the table below, with the best one highlighted (Experiment 1, using VGG16).

Experiment	Model	DNNs	L2 Reg	Batch Norm	Dropout	Test Accuracy	Test Loss	Precision	Recall	F1-Score	RMSE
1	VGG16	2	None	None	None	0.8806	0.3338	0.88	0.88	0.88	1.267
2	<u>VGG19</u>	2	None	None	None	0.8734	0.3533	0.87	0.87	0.87	1.288
3	<u>Inception V3</u>	2	None	None	None	0.8427	0.4356	0.85	0.84	0.84	1.356
4	<u>ResNet152V3</u>	2	None	None	None	0.7925	0.564	0.79	0.79	0.79	1.616

The reason why these pre-trained models did not perform as well as our own CNNs in this study could be due to the image transformation performed on the Fashion-MNIST dataset, as they require images to be in a 3-channel color format and need to be larger than their 28x28 resolution. While we are able to perform such transformations, the overall image quality might have been compromised in the process.

The figure below shows a comparison of an original Fashion-MNIST image with its transformed counterparts used for training and testing the pre-trained models. While the image colors are inverted, we see a blurrier image as its resolution is increased, thus a degradation in image quality.



With ResNet requiring a 75x75 resolution image and our transformation process results in a blurrier image, it is understandable why the model would perform worse than other models as it is more difficult for it to distinguish features, despite its complex architecture.

Overall, we found the best performing model to be Experiment 1 in the Grouped Convolutional Models section, yielding a 93.7% accuracy on the Fashion-MNIST test dataset.

Conclusion

We were able to build a convolutional neural network model that can classify images from the Fashion-MNIST test dataset with 93.7% accuracy. Incorporating regularization techniques to our models such as batch normalization, dropout, and ridge regression along with hyperparameter tuning can improve accuracy by almost 1%. We also saw diminishing returns when we created a CNN using 4 isolated convolutional layers did not improve test accuracy. Utilization of pre-trained models such as ResNet did

not achieve similar results compared to other studies that used them on the Fashion-MNIST dataset, but it could be due to the image quality of the transformed dataset. Our most optimal CNN-based model uses a complex architecture using multiple groups of convolutional layers. Given the 93.7% accuracy, we feel this model is a good first step towards integrating a classifier into the retailer's recommendation system. While we feel the model's accuracy is enough to be deployed for A/B testing on a beta version of the retail website, there is still more room for improvement. We can expect to optimize model accuracy over time, either with more nuanced tuning of our own developed models, creating a more complex architecture of grouped convolution layers, deep-diving into the details behind pre-trained models, or improving our image transformation processes.

Resources

Nocentini, O., Kim, J., Bashir, M. Z., & Cavallo, F. (2022). Image Classification Using Multiple Convolutional Neural Networks on the Fashion-MNIST Dataset. *Sensors*, 22(23), 9544.

<https://doi.org/10.3390/s22239544>

Tang, Y., Hanguo ., & Shuyong L. (2020). Optimal Design of Deep Residual Network Based on Image Classification of Fashion-MNIST Dataset. *Journal of Physics: Conference Series*, 1624.5 (2020), 52011.

<http://doi.org/10.1088/1742-6596/1624/5/052011>