

## ASSIGNMENT 3

# SYSTEMS ENGINEERING FOR DEEP LEARNING

Submitted By: Sooryakiran (ME17B174)

### STEP 1

A neural network with the given architecture was trained on the FMNIST dataset. The network weights were initialized with Xavier uniform initialisation. Stochastic Gradient Descent optimizer with momentum was used with a learning rate of 0.001 and with a momentum weightage of 0.9.

Dataset	FMNIST  60000 Training samples 10000 Testing samples
Preprocessing	Normalization
Optimizer	SGD with momentum
Learning Rate	1e-3
Momentum ( $\beta$ )	0.9
Batch Size	32
Network initialization	Weights : Xavier uniform Biases : Zero
Network Architecture	Layer 1 : Conv 3x3, 16 Channels, ReLU  Layer 2: Conv 3x3, 16 Channels, ReLU  Layer 3: Dense 100 neurons, ReLU  Layer 4: Dense 10 neurons, Softmax

*Table 1.1:  
Model and training  
hyperparameters.*

The network was trained in a CPU only instance on Google Colaboratory. With a batch size of 32, I got 90.18 % training accuracy and a test accuracy of 88.33%. The loss vs batches plots and accuracy vs batches plots and the final results are given in the following page.

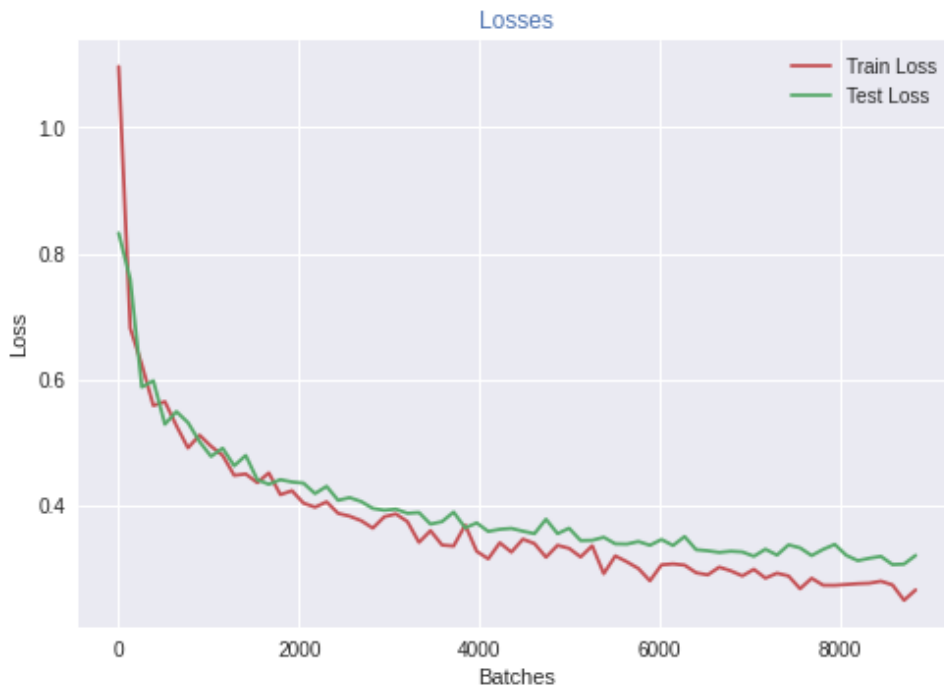


Figure 1.1.a:  
Loss vs Batches plot

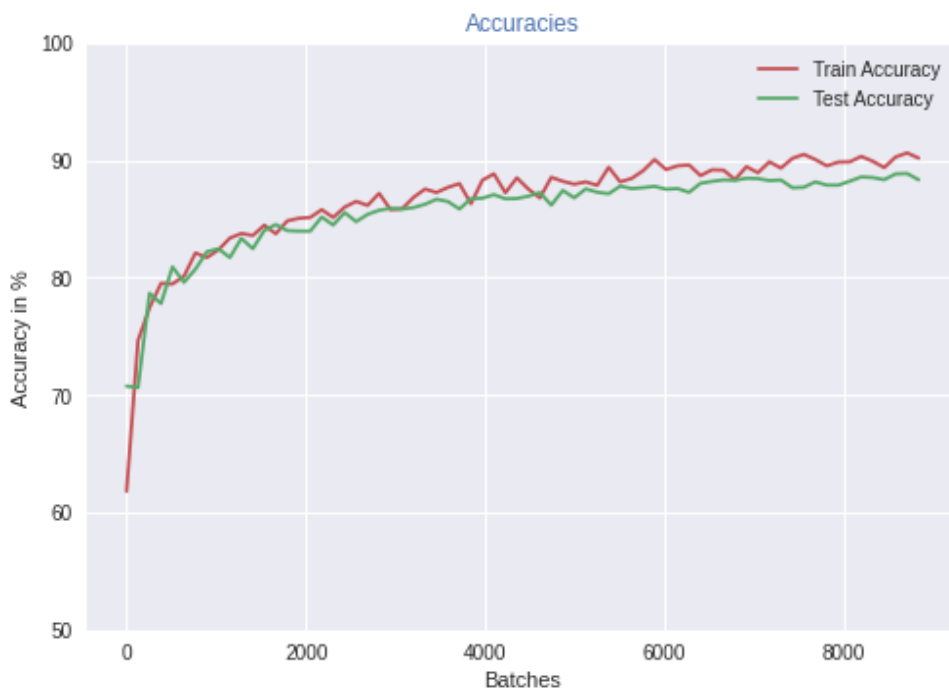


Figure 1.1.b:  
Accuracy vs Batches  
plot

Google Facets was used to analyse the prediction performance of the trained model. Google Facets allow inputs only in CSV format. So the first the model was used to generate predictions for the entire test dataset. Then a CSV file was generated with columns being the True and Predicted labels. This CSV file was uploaded to Google Facets website. Many visualizations like the confusion matrix, distribution of predictions were obtained.

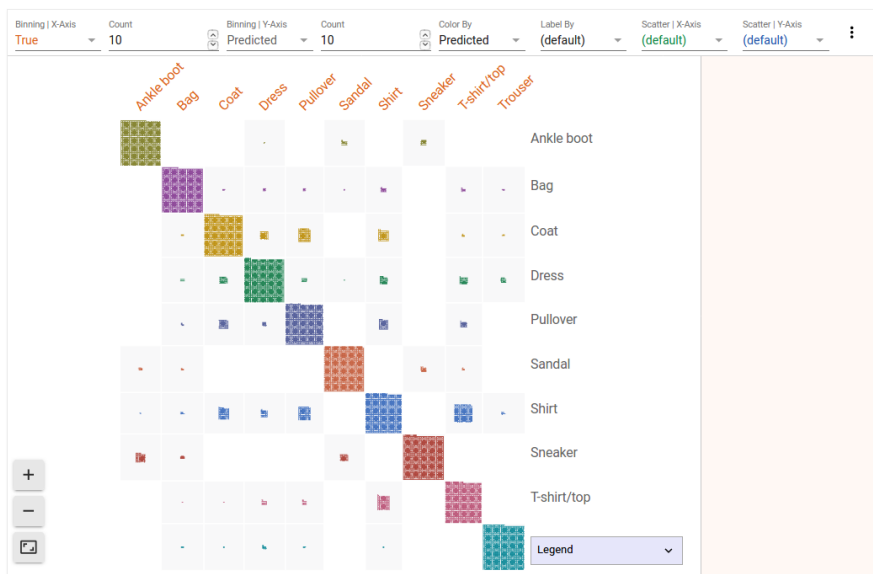


Figure 1.2.a:  
Confusion Matrix

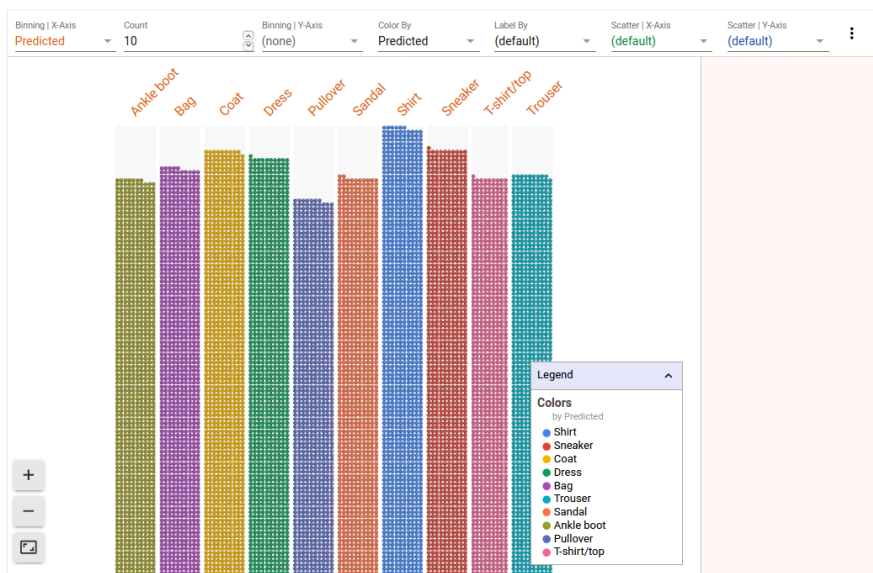


Figure 1.2.b:  
Distribution of  
predictions.



Figure 1.2.c:  
Predictions for each  
true label.

These visualizations helps us to understand where out model is failing to perform. For example, in Fig. 1.2.c, we can see that nearly 30% of T-shirts/tops are misclassified as shirts.

👉(ツ)👈	Loss	Accuracy
Training	0.267	90.186 %
Testing	0.321	88.338 %

Table 1.2:  
Train and test performances

STEP 2

MLflow was used to conduct grid search to find the optimal hyperparameters for the model. Two experiments were conducted, one to find the right Model Hyperparameters and another one to find the right Training Hyperparameters. These experiments were conducted one by one. The optimal model architecture obtained from the first experiment was used for the second experiment.

The first experiment search space is as follows:

Hyperparameter	Values
Conv 1 kernel sizes	3x3, 5x5
Conv 2 kernel sizes	3x3, 5x5
Conv 1 no of filters	8, 16, 32
Conv 2 no of filters	8, 16, 32
Dense 1 no of neurons	64, 100

Table 2.1.a:  
The search space for experiment 1

Results are as follows.

▼ Notes🔗

None

Search Runs: metrics.rmse < 1 and params.model = "tree" and tags.mlflow.source.type = "LOCAL" State: Active Search Clear

Showing 80 matching runs Compare Delete Download CSV📄 Columns

	Parameters <					Metrics	
<input type="checkbox"/>	Conv layer 1 channel size	Conv layer 1 filter size	Conv layer 2 channel size	Conv layer 2 filter size	Layer 3 size	↓ Test Accuracy	Training time
<input type="checkbox"/>	16	5	32	5	100	0.891473642172524	570.8142561912537
<input type="checkbox"/>	16	5	32	5	64	0.8889776357827476	606.5375547409058
<input type="checkbox"/>	8	3	8	3	64	0.8887779552715654	351.8503818511963
<input type="checkbox"/>	32	3	32	3	100	0.8885782747603834	829.9171538352966
<input type="checkbox"/>	32	5	16	5	64	0.8884784345047924	711.0066647529602
<input type="checkbox"/>	8	5	8	5	64	0.8882787539936102	441.177805185318
<input type="checkbox"/>	8	5	32	5	100	0.8880790734824281	486.8598906993866
<input type="checkbox"/>	32	5	8	5	100	0.8880790734824281	638.7947866916656
<input type="checkbox"/>	32	3	32	3	100	0.8876797124600639	855.222030878067
<input type="checkbox"/>	16	3	8	5	100	0.8876797124600639	501.46214294433594
<input type="checkbox"/>	16	3	32	5	100	0.8873801916932907	523.446713924408
<input type="checkbox"/>	16	3	32	5	64	0.8870806709265175	655.0486626625061
<input type="checkbox"/>	16	3	32	3	100	0.8866813099041534	667.5740876197815
<input type="checkbox"/>	16	5	32	3	100	0.8863817891373802	545.255553483963

Fig 2.1.a:  
Experiment results sorted by decreasing accuracy on the test dataset.

▼ Notes [🔗](#)

None

Search Runs:  State: Active ▾ Search Clear

Showing 80 matching runs Compare Delete Download CSV Columns

Parameters <						Metrics	
<input type="checkbox"/>	Conv layer 1 channel size	Conv layer 1 filter size	Conv layer 2 channel size	Conv layer 2 filter size	Layer 3 size	Test Accuracy	↑ Training time
<input type="checkbox"/>	8	3	8	3	100	0.8857827476038339	305.4713706970215
<input type="checkbox"/>	8	5	8	3	64	0.8798921725239617	331.3197486400604
<input type="checkbox"/>	8	3	8	5	100	0.8772963258785943	349.50655937194824
<input type="checkbox"/>	8	3	8	3	64	0.888779552715654	351.8503818511963
<input type="checkbox"/>	8	3	16	3	100	0.8828873801916933	352.47868728637695
<input type="checkbox"/>	8	5	8	5	100	0.8855830670926518	369.799165725708
<input type="checkbox"/>	8	5	16	3	100	0.8850838658146964	380.6902048587799
<input type="checkbox"/>	8	3	16	5	100	0.884185303514377	386.0200471878052
<input type="checkbox"/>	8	5	8	3	100	0.8799920127795527	392.9934973716736
<input type="checkbox"/>	8	3	16	5	64	0.8835862619808307	398.44480061531067
<input type="checkbox"/>	8	5	16	5	100	0.8810902555910544	402.0381429195404
<input type="checkbox"/>	16	3	8	3	64	0.8822883386581469	409.6416800022125
<input type="checkbox"/>	8	5	16	5	64	0.876797124600639	411.1483144760132
<input type="checkbox"/>	16	5	16	3	100	0.878694089456869	413.8496277332306

*Fig 2.1.b:  
Experiment results  
sorted by increasing  
training time required to  
reach 90% train  
accuracy.*

Top 3 results are as shown below:

	#1	#2	#3
Conv 1	5x5, 16 filters	5x5, 16 filters	3x3, 8 filters
Conv 2	5x5, 32 filters	5x5, 32 filters	3x3, 8 filters
Dense 3	100 neurons	64 neurons	64 neurons
Test Accuracy	89.14 %	88.89 %	88.87 %
Training Time	570.81s	606.53s	351.85s

*Table 2.1.a:  
Top 3 results based on  
test accuracy.*

	#1	#2	#3
Conv 1	3x3, 8 filters	5x5, 8 filters	3x3, 8 filters
Conv 2	3x3, 8 filters	3x3, 8 filters	5x5, 8 filters
Dense 3	100 neurons	64 neurons	100 neurons
Test Accuracy	88.57 %	87.98 %	87.72 %
Training Time	305s	331s	349s

*Table 2.1.b:  
Top 3 results based on  
training time.*

Smaller filter sizes have smaller training time because there are less parameters and less computation involved. But filters with larger receptive field performs better as they have higher representative power. This effect is more profound for shallow networks like the one we experimented with. Also, networks with higher number of filters can represent more features and are likely to have higher performance but compromises on training time. When the number of neurons on the dense layer is increased, it improves the performance but increases the training time.

I decided to choose the model with lowest training time for the next experiment because it is acceptable to sacrifice 0.64% of accuracy for 46% reduction in training time.

Experiment 2 was conducted on the model that had the smallest training time.

Hyperparameter	Values
Initialisation	Xavier Uniform, Xavier Normal, Kaiming Uniform
Batch sizes	8, 16, 32, 64
Optimizer params	1. Learning rate : 0.001 Momentum : 0.9 2. Learning rate : 0.01 Momentum : 0.9 3. Learning rate : 0.001 Momentum : 0.5 4. Learning rate : 0.01 Momentum : 0.5

Table 2.2:  
The search space for experiment 2

Results are as follows.

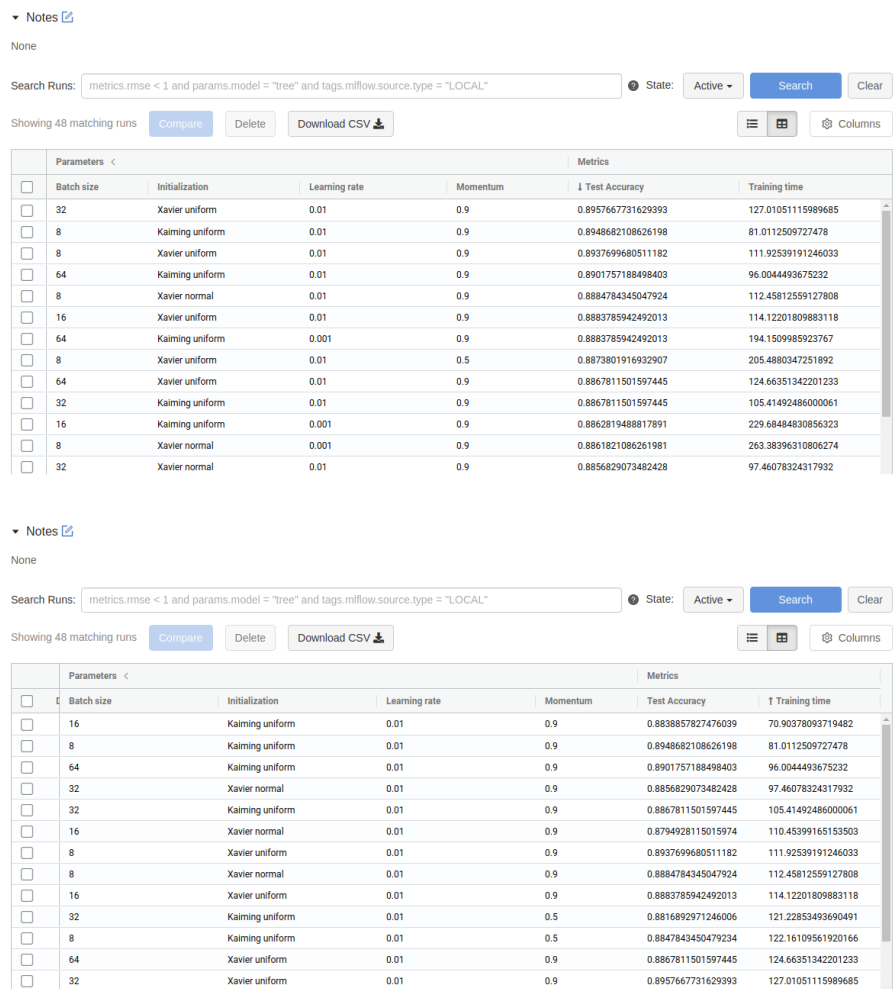


Fig 2.2.a:  
Experiment 2 results sorted by decreasing accuracy on the test dataset.

Fig 2.2.b:  
Experiment 2 results sorted by increasing training time required to reach 90% train accuracy.

Top 3 results are as shown below:

	#1	#2	#3
Batch size	32	8	8
Initialization	Xavier Uni.	Kaiming Uni.	Xavier Uni.
Learning rate	0.01	0.01	0.01
Momentum	0.9	0.9	0.9
Test Accuracy	89.57 %	89.48 %	89.37 %
Training Time	127s	81s	111s

*Table 2.2.a:  
Top 3 results based on  
test accuracy.*

	#1	#2	#3
Batch size	16	8	64
Initialization	Kaiming Uni.	Kaiming Uni.	Kaiming Uni.
Learning rate	0.01	0.01	0.01
Momentum	0.9	0.9	0.9
Test Accuracy	88.38 %	89.48 %	89.01 %
Training Time	70s	81s	96s

*Table 2.2.b:  
Top 3 results based on  
training time.*

Learning rate of 0.01 performed well in most cases. It converges faster and results in better performance on the test dataset. Also models with momentum of 0.9 outperformed models with momentum of 0.5 because it is better in escaping local minima traps. Kaiming initialization is found to converge faster. A smaller batch size give rise to gradients that are more randomized comparatively. This might help in escaping the local minima traps.

Here are some scatter plots that can help us choose even better models. However, not every one of these plots are conclusive. Plots like the one for learning rate, momentum, and number of filters can help us generalise how different models behave according to different hyperparameters. The MLflow logs are converted to CSV files and uploaded to Google Facets.

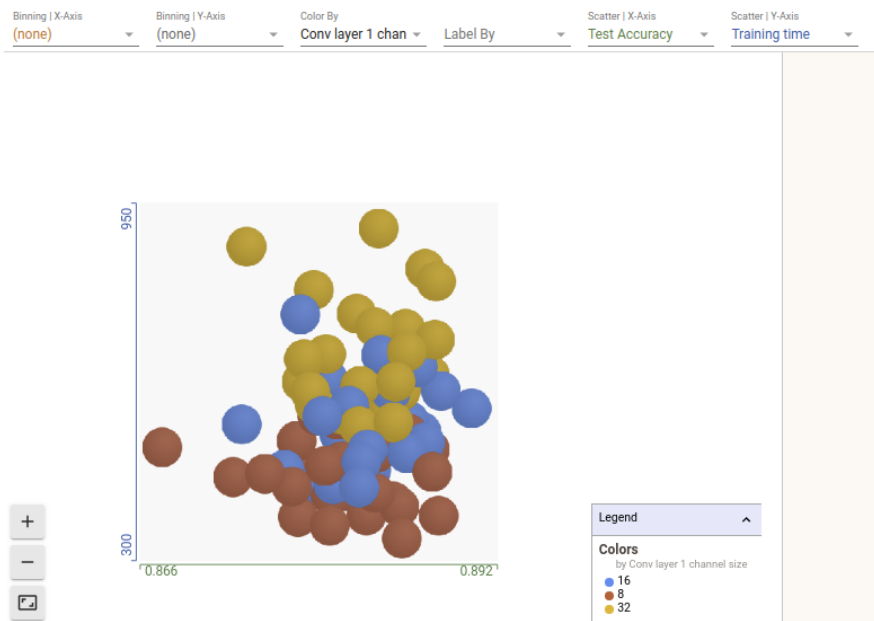


Fig 2.3.a:  
Conv 1 no of filters  
x-axis : Test accuracy  
y-axis: Training time

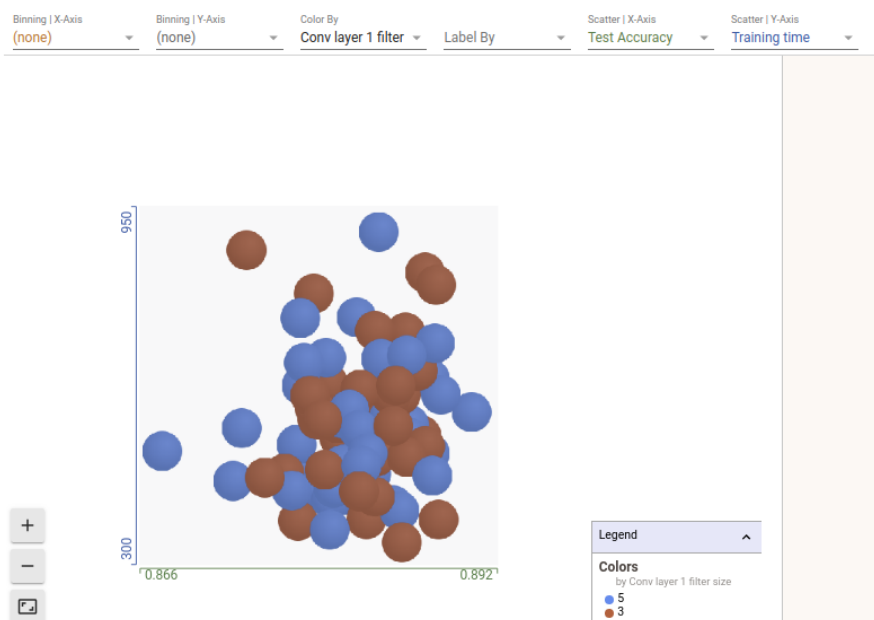


Fig 2.3.b:  
Conv 1 kernel size  
x-axis : Test accuracy  
y-axis: Training time

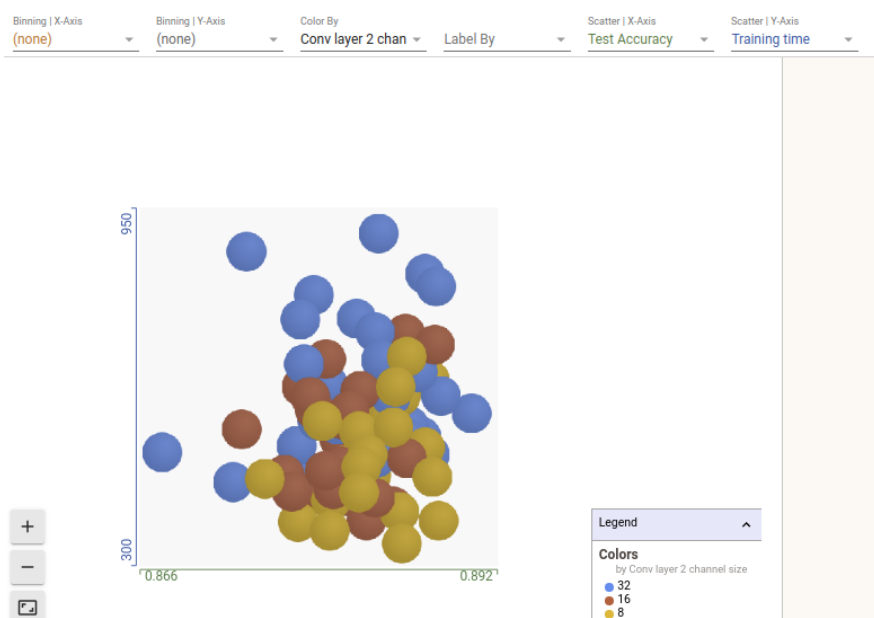


Fig 2.3.c:  
Conv 2 no of filters  
x-axis : Test accuracy  
y-axis: Training time



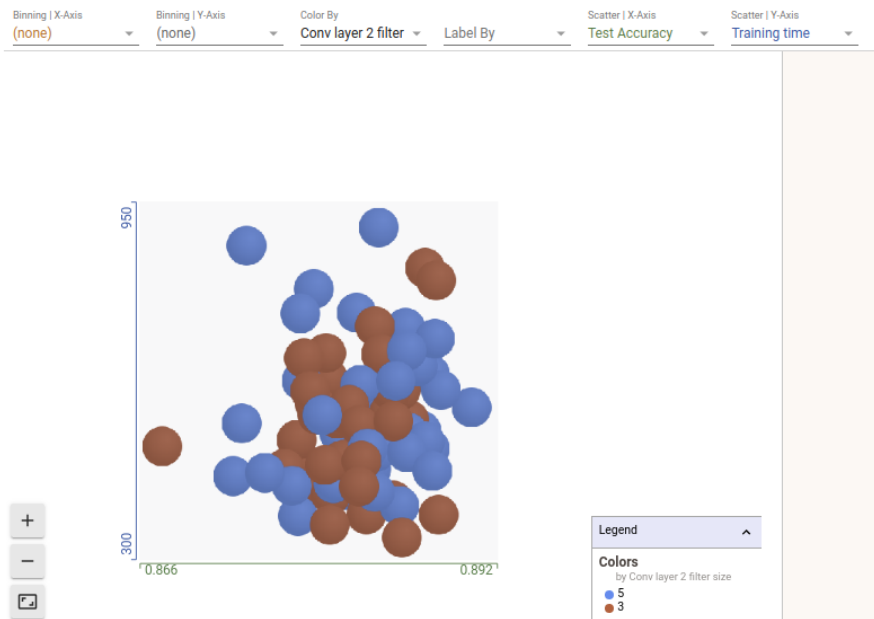


Fig 2.3.d:  
Conv 2 kernel size  
x-axis : Test accuracy  
y-axis: Training time

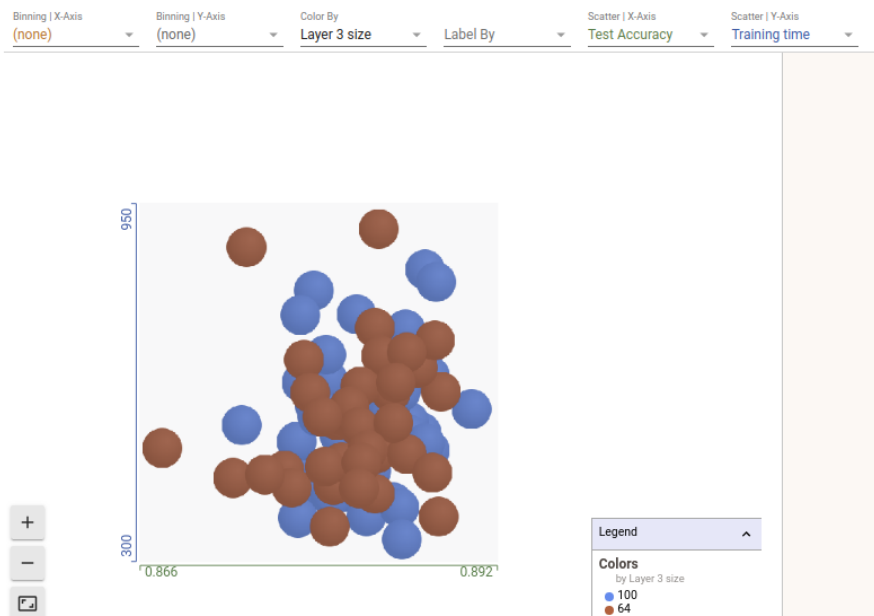


Fig 2.3.e:  
Number of neurons in  
layer 3  
x-axis : Test accuracy  
y-axis: Training time

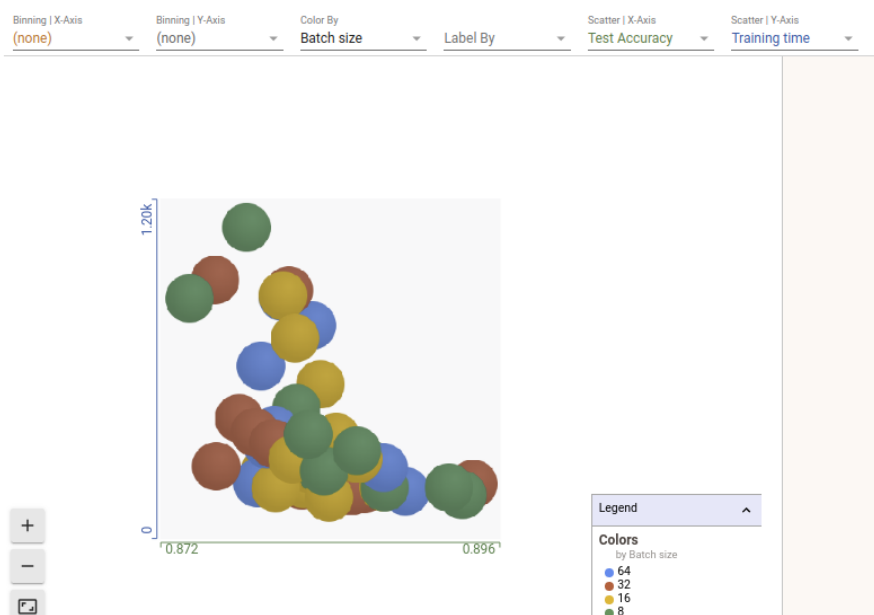


Fig 2.3.f:  
Batch size  
x-axis : Test accuracy  
y-axis: Training time

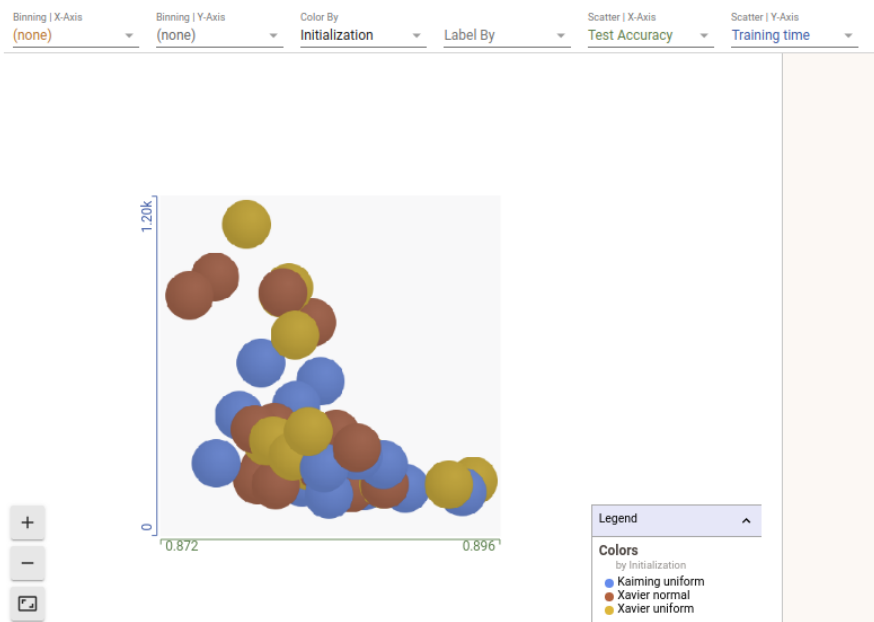


Fig 2.3.g:  
Initialization  
x-axis : Test accuracy  
y-axis: Training time

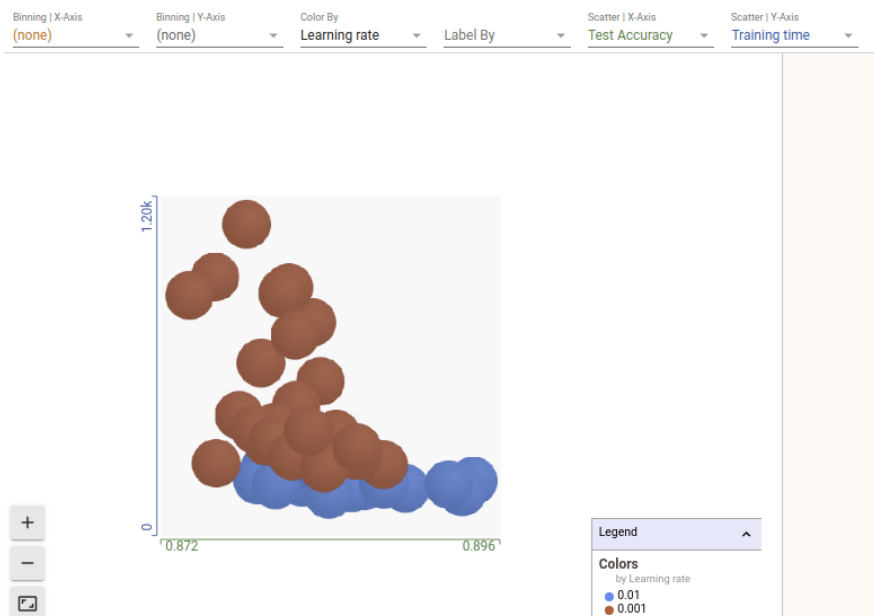


Fig 2.3.h:  
Learning rate  
x-axis : Test accuracy  
y-axis: Training time

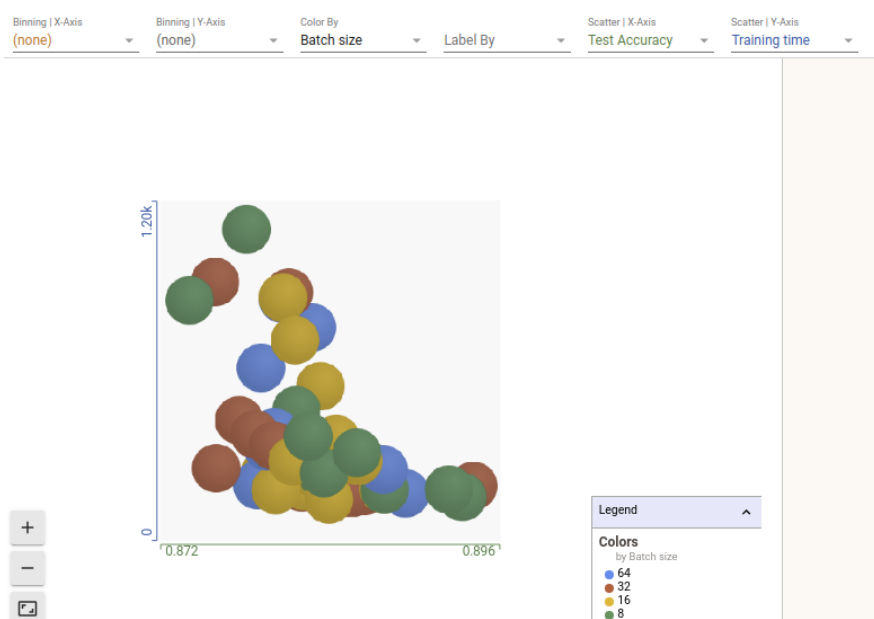


Fig 2.3.i:  
Momentum  
x-axis : Test accuracy  
y-axis: Training time

The pareto plots of all models with respect to testing accuracy and training time are as given below.

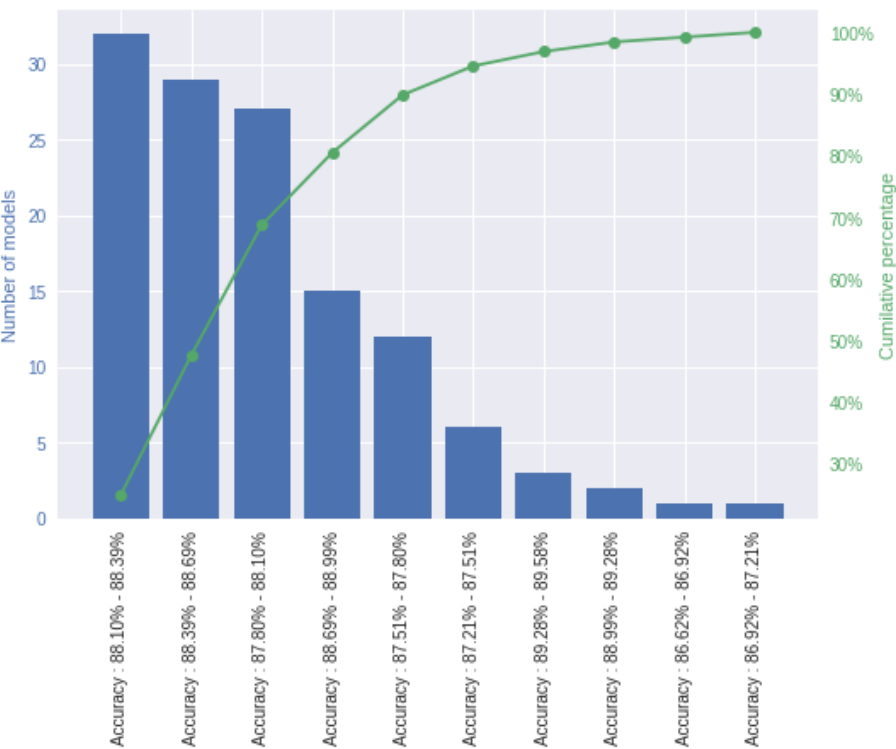


Fig 2.4.a:  
Pareto plot of all models  
with respect to testing  
accuracy.

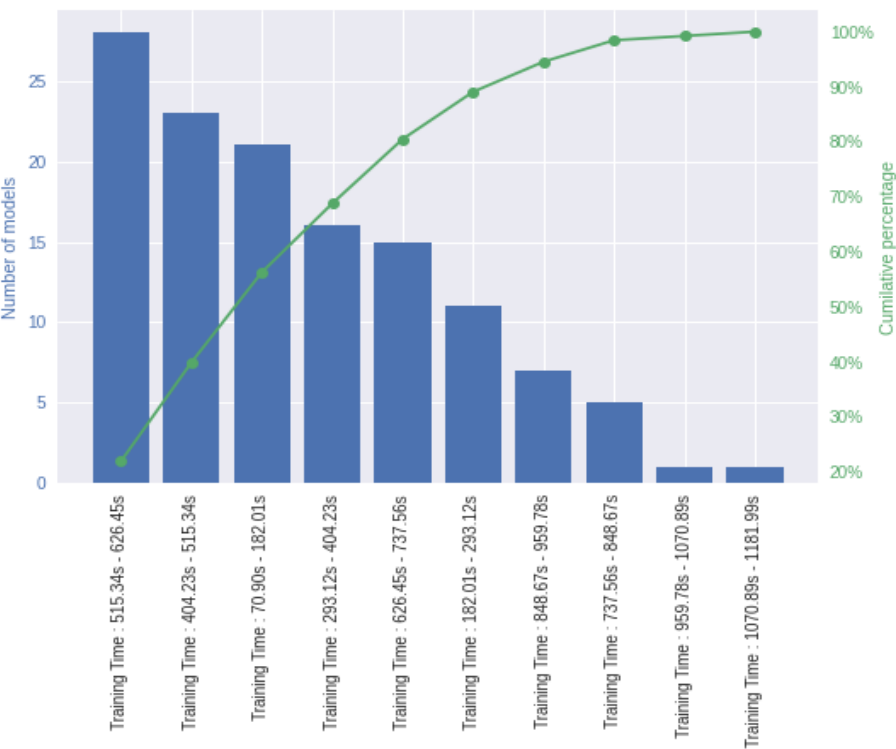


Fig 2.4.b:  
Pareto plot of all models  
with respect to training  
time.

Loss vs epoch/batches plots for some special cases asked in the question are given below.

Learning rate too high: The largest learning rate used for the experiment was 0.01. This is not a high learning rate. Therefore to produce the plot, I trained a model separately with learning rate. We can see that the loss plot has got many irregular jumps.



*Fig 2.5.a:  
High learning rate  
(0.05). Loss plot has  
irregular jumps.*

Learning rate too low: The MLflow experiment did not contain any examples for low learning rates. So I trained a separate model to demonstrate. The network is taking longer to converge.



*Fig 2.5.b:  
Low learning rate  
(0.0001). Loss plot  
takes very longer to  
converge.*

Momentum too high: This plot is also generated by running an experiment separately. We can see that the plot has got an oscillatory nature. The model oscillates around local minimas. Sometimes the network is able to escape local minimas, but sometimes it keeps on oscillating.



*Fig 2.5.c:  
High momentum (0.99,  
lr: 0.001). The loss  
oscillates too much.*

Batch size too low: When a network is trained using stochastic gradient descend on a very small batch size. The updates become more randomized. The plot below shows a network trained with batch size of 4. All other settings were similar to that of the above given networks.



*Fig 2.5.d:  
Low batch size. All the  
updates become more  
approximate and  
randomized, thus  
leading to a rugged loss  
plot.*

Overfitting due to large model: The experiments conducted did not have any good examples of overfitting. So a different larger model was trained to demonstrate the effect. This example was trained longer than others to make the effect more profound. We can see the test loss actually diverging.



*Fig 2.5.e:  
Overfitting on a large  
model.(Conv 5x5, 64;  
Conv 5x5, 64; 100  
neurons)*

