

Steps to follow in ANN:

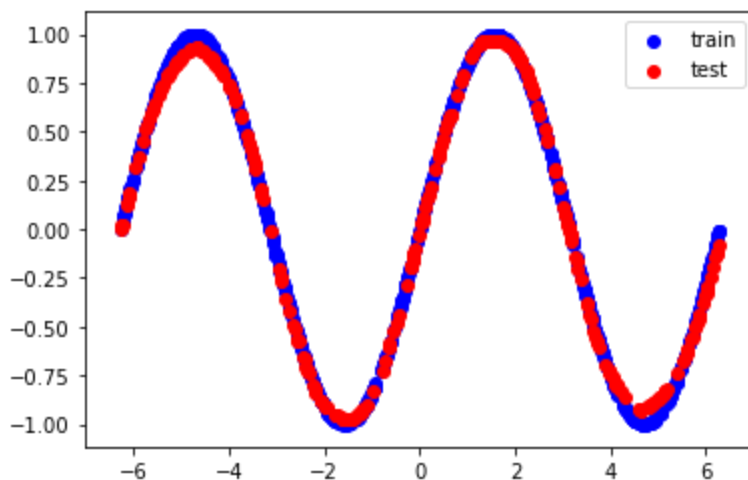
L_layer_model_minib() -- Main function

- `initialize_parameters_deep_he(layers_dims)` -- *initialise weights*
- `initialize_adam(parameters)` -- *initialise the optimiser*
- `random_mini_batches(X, Y, mini_batch_size, seed)` -- *defines randomised batches.*
- For i in mini batches
 - `L_model_forward(X, parameters, activation)` -- *computes the forward pass*
 - Loop over the layers and do forward prop on each
 - `linear_activation_forward(A_prev, W, b, activation)` -- *compute the linear forward and then apply activation*
 - `linear_forward(A, W, b):`
 - $Z = W \cdot A + b$
 - `A, activation_cache = activation(Z)`
 - `compute_cost(AL, Y, parameters, lambda, regularisation, cost_func='mse')` -- *computes the cost*
 - Mse : $\text{cost} = \text{np.mean}(\text{np.square}(\text{AL} - Y)) * 0.5$
 - Log : $\text{cost} = (1./m) * \text{np.sum}(-\text{np.dot}(Y, \text{np.log}(\text{AL} + \text{epsilon})).T) - \text{np.dot}(1 - Y, \text{np.log}(1 - \text{AL} + \text{epsilon})).T)$
 - `L_model_backward(AL, Y, caches, activation, regularisation, lambda, cost_func)` -- *perform backprop*
 - Mse : $dAL = (AL - Y)$
 - Log : $dAL = -(\text{np.divide}(Y, AL) - \text{np.divide}(1 - Y, 1 - AL))$
 - Loop over the layers and do back prop on each
 - `linear_activation_backward(dA, cache, regularisation, lambda, activation)` -- *compute back prop from activation to the weights*
 - $dZ = \text{activation_backward}(dA, \text{activation_cache})$
 - `linear_backward(dZ, cache, regularisation, lambda)`
 - $db = 1./m * \text{np.sum}(dZ, \text{axis} = 1, \text{keepdims} = \text{True})$
 - $dA_{\text{prev}} = \text{np.dot}(W.T, dZ)$
 - `update_parameters(parameters, grads, learning_rate)` -- *finally update all the weights in the network*
 - $\text{cost_avg} = \text{cost_total} / \text{batches}$ -- *compute the avg costs over the batch*
 - `predicterr(valid_x, valid_y, parameters, lambda, activation=activation, regularisation='none', cost_func=cost_func)` -- *predict the validation cost*
 - `plt.plot(costs)` -- *plot the costs*

Best Model for Sin Curve

```
layers_dims = [1,20,20,1]
```

```
parameters = L_layer_model_minib(  
train_x, train_y, layers_dims, valid=False,  
    num_iterations = 600,  
    learning_rate = 0.001,  
    print_cost = True,  
    lambd=0.1,  
    optimizer="adam",  
    beta = 0.9,beta1 = 0.9, beta2 = 0.999,epsilon = 1e-8,  
    activation='tanh',  
    regularisation='none',  
    mini_batch_size=64,  
    cost_func="mse",  
    he_init=True)
```



Architecture Comparisons

<p>layers_dims = [4,5,1]</p> <p>Training mape : 0.7941098689193549 Valid mape : 0.8063094442632108 Testing mape : 0.7990991185251256</p>	<p>Learning rate = 0.001</p> <p>This plot shows the training and validation cost over 100 epochs for the [4,5,1] architecture. The training cost (blue line) starts at approximately 0.175 and decreases to about 0.015. The validation cost (orange line) starts at approximately 0.14 and decreases to about 0.01. Both costs show a rapid initial drop followed by a slower decline.</p> <table><tr><th>Epochs</th><th>Train Cost</th><th>Validation Cost</th></tr><tr><td>0</td><td>0.175</td><td>0.140</td></tr><tr><td>20</td><td>0.050</td><td>0.025</td></tr><tr><td>40</td><td>0.030</td><td>0.018</td></tr><tr><td>60</td><td>0.020</td><td>0.015</td></tr><tr><td>80</td><td>0.018</td><td>0.014</td></tr><tr><td>100</td><td>0.015</td><td>0.013</td></tr></table>	Epochs	Train Cost	Validation Cost	0	0.175	0.140	20	0.050	0.025	40	0.030	0.018	60	0.020	0.015	80	0.018	0.014	100	0.015	0.013
Epochs	Train Cost	Validation Cost																				
0	0.175	0.140																				
20	0.050	0.025																				
40	0.030	0.018																				
60	0.020	0.015																				
80	0.018	0.014																				
100	0.015	0.013																				
<p>layers_dims = [4,10,1]</p> <p>Training mape : 0.7740370114496019 Valid mape : 0.779183617465357 Testing mape : 0.7870296006952487</p>	<p>Learning rate = 0.001</p> <p>This plot shows the training and validation cost over 100 epochs for the [4,10,1] architecture. The training cost (blue line) starts at approximately 0.25 and decreases to about 0.01. The validation cost (orange line) starts at approximately 0.18 and decreases to about 0.01. The training cost decreases more slowly than in the [4,5,1] architecture.</p> <table><tr><th>Epochs</th><th>Train Cost</th><th>Validation Cost</th></tr><tr><td>0</td><td>0.250</td><td>0.180</td></tr><tr><td>20</td><td>0.080</td><td>0.020</td></tr><tr><td>40</td><td>0.040</td><td>0.015</td></tr><tr><td>60</td><td>0.025</td><td>0.012</td></tr><tr><td>80</td><td>0.020</td><td>0.011</td></tr><tr><td>100</td><td>0.015</td><td>0.010</td></tr></table>	Epochs	Train Cost	Validation Cost	0	0.250	0.180	20	0.080	0.020	40	0.040	0.015	60	0.025	0.012	80	0.020	0.011	100	0.015	0.010
Epochs	Train Cost	Validation Cost																				
0	0.250	0.180																				
20	0.080	0.020																				
40	0.040	0.015																				
60	0.025	0.012																				
80	0.020	0.011																				
100	0.015	0.010																				
<p>layers_dims = [4,10,5,1]</p> <p>Training mape : 0.7692920370372286 Valid mape : 0.7763830442763537 Testing mape : 0.7821588242726298</p>	<p>Learning rate = 0.001</p> <p>This plot shows the training and validation cost over 100 epochs for the [4,10,5,1] architecture. The training cost (blue line) starts at approximately 0.11 and decreases to about 0.015. The validation cost (orange line) starts at approximately 0.04 and decreases to about 0.01. The training cost decreases more slowly than in the [4,10,1] architecture.</p> <table><tr><th>Epochs</th><th>Train Cost</th><th>Validation Cost</th></tr><tr><td>0</td><td>0.110</td><td>0.040</td></tr><tr><td>20</td><td>0.040</td><td>0.015</td></tr><tr><td>40</td><td>0.025</td><td>0.012</td></tr><tr><td>60</td><td>0.020</td><td>0.011</td></tr><tr><td>80</td><td>0.018</td><td>0.010</td></tr><tr><td>100</td><td>0.015</td><td>0.010</td></tr></table>	Epochs	Train Cost	Validation Cost	0	0.110	0.040	20	0.040	0.015	40	0.025	0.012	60	0.020	0.011	80	0.018	0.010	100	0.015	0.010
Epochs	Train Cost	Validation Cost																				
0	0.110	0.040																				
20	0.040	0.015																				
40	0.025	0.012																				
60	0.020	0.011																				
80	0.018	0.010																				
100	0.015	0.010																				
<p>layers_dims = [4,10,10,1]</p> <p>Training mape : 0.7716003838887064 Valid mape : 0.779187773992918 Testing mape : 0.7839908249893479</p>	<p>Learning rate = 0.001</p> <p>This plot shows the training and validation cost over 100 epochs for the [4,10,10,1] architecture. The training cost (blue line) starts at approximately 0.24 and decreases to about 0.01. The validation cost (orange line) starts at approximately 0.14 and decreases to about 0.01. The training cost decreases more slowly than in the [4,10,1] architecture.</p> <table><tr><th>Epochs</th><th>Train Cost</th><th>Validation Cost</th></tr><tr><td>0</td><td>0.240</td><td>0.140</td></tr><tr><td>20</td><td>0.080</td><td>0.015</td></tr><tr><td>40</td><td>0.040</td><td>0.012</td></tr><tr><td>60</td><td>0.025</td><td>0.011</td></tr><tr><td>80</td><td>0.020</td><td>0.010</td></tr><tr><td>100</td><td>0.015</td><td>0.010</td></tr></table>	Epochs	Train Cost	Validation Cost	0	0.240	0.140	20	0.080	0.015	40	0.040	0.012	60	0.025	0.011	80	0.020	0.010	100	0.015	0.010
Epochs	Train Cost	Validation Cost																				
0	0.240	0.140																				
20	0.080	0.015																				
40	0.040	0.012																				
60	0.025	0.011																				
80	0.020	0.010																				
100	0.015	0.010																				

Batch Size Comparisons

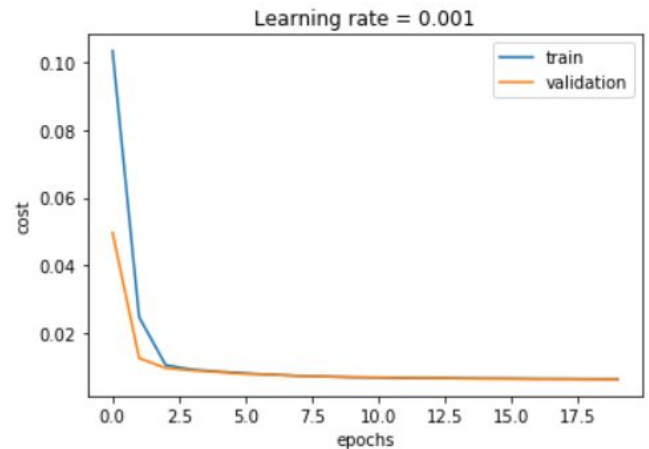
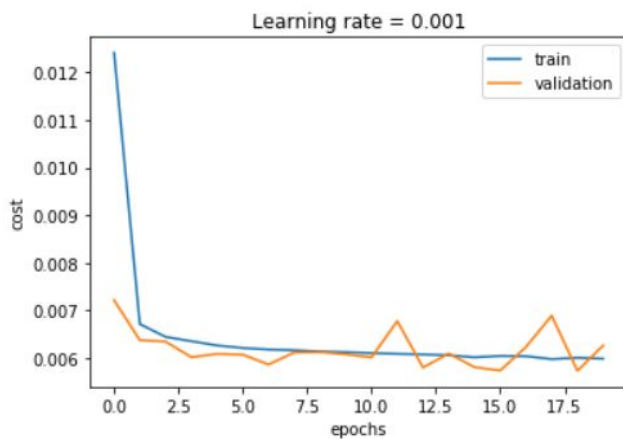
Freezing other parameters and changin batch sizes

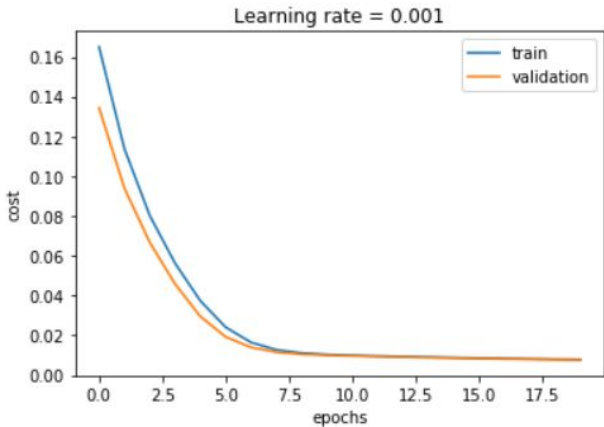
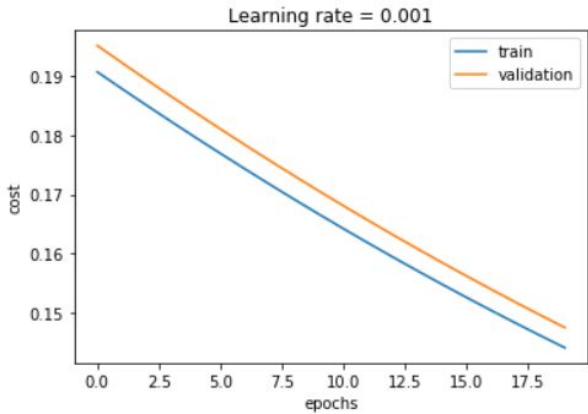
layers_dims = [4,10,10,1]
num_iterations = 20,
learning_rate = 0.001
mini_batch_size =1(SGD)
optimizer="adam"
activation=Tanh
he_init=True
cost_func='mse'
regularisation ='none'

Training mape : 0.6939322934163785
Valid mape : 0.7115727145429904
Testing mape : 0.7354052660571078

layers_dims = [4,10,10,1]
num_iterations = 20,
learning_rate = 0.001
mini_batch_size =64
optimizer="adam"
activation=Tanh
he_init=True
cost_func='mse'
regularisation ='none'

Training mape : 0.7226041626498847
Valid mape : 0.7243184792599733
Testing mape : 0.7423342115034387



<p> layers_dims = [4,10,10,1] num_iterations = 20, learning_rate = 0.001 mini_batch_size =256 optimizer="adam" activation=Tanh he_init=True cost_func='mse' regularisation ='none' </p> <p> Training mape : 0.787330891219502 Valid mape : 0.8000444029923501 Testing mape : 0.8021629604210103 </p>	<p> layers_dims = [4,10,10,1] num_iterations = 20, learning_rate = 0.001 mini_batch_size =6888(batch) optimizer="adam" activation=Tanh he_init=True cost_func='mse' regularisation ='none' </p> <p> Training mape : 3.600072792312424 Valid mape : 3.702112026264359 Testing mape : 3.6967978640657115 </p>
	

- As batch size increases learning becomes slow-*more iterations needed*
- As batch size increases learning per epoch becomes faster - *vectorisation comes into play*
- Even with sufficient epoch SGD tends to give slightly better results
- SGD graph for validation is usually a bit zig zag and not smooth(optimisers will help)

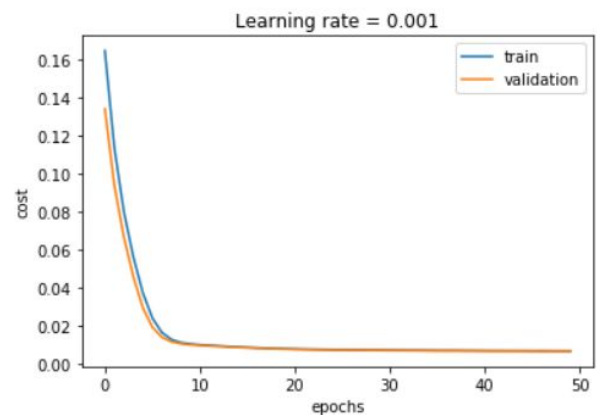
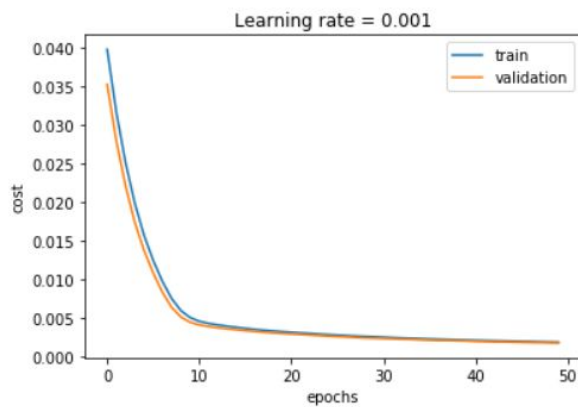
Output Activation Comparisons

layers_dims = [4,10,10,1]
num_iterations = 50,
learning_rate = 0.001
mini_batch_size =256
optimizer="adam"
activation="sigmoid"
he_init=True
cost_func='mse'
regularisation ='none'

Training mape : 0.7726484973627554
Valid mape : 0.7767228155762528
Testing mape : 0.7946247880970604

layers_dims = [4,10,10,1]
num_iterations = 50,
learning_rate = 0.001
mini_batch_size =256
optimizer="adam"
activation="tanh"
he_init=True
cost_func='mse'
regularisation ='none'

Training mape : 0.7264987907448855
Valid mape : 0.732542065744155
Testing mape : 0.7536587432898809



- Tanh converges faster

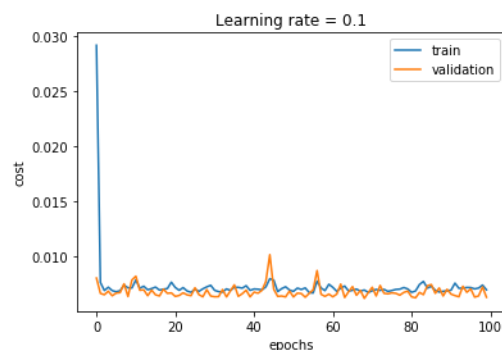
Learning Rate:

lr=0.1

Training mape : 0.7189152217350819

Valid mape : 0.7298132876871712

Testing mape : 0.7388300192620225

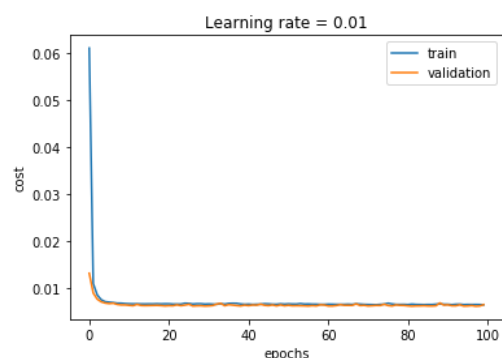


lr=0.01

Training mape : 0.7169999218682349

Valid mape : 0.7318793588526947

Testing mape : 0.7451245775566842

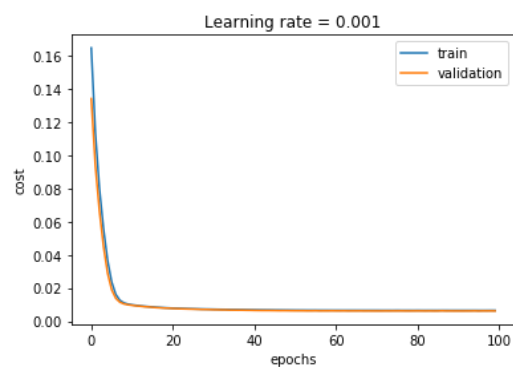


lr=0.001

Training mape : 0.7208036807521352

Valid mape : 0.7226891777193623

Testing mape : 0.7271498827770438

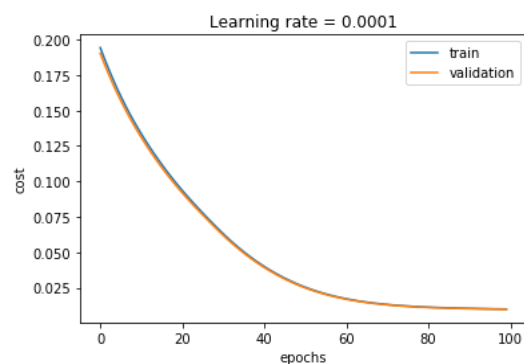


lr=0.0001

Training mape : 0.9026285294068407

Valid mape : 0.920032185594116

Testing mape : 0.9188982131125064



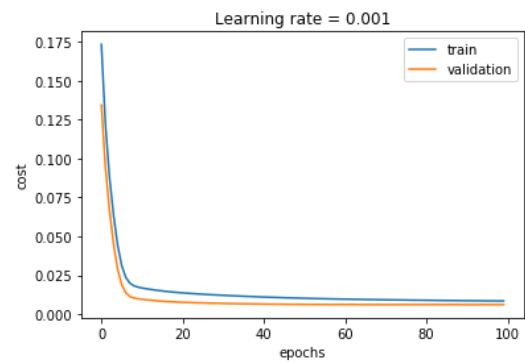
Regularisation

Lambd = 0.1

Training mape : 0.7208036807521352

Valid mape : 0.7226891777193623

Testing mape : 0.7271498827770438

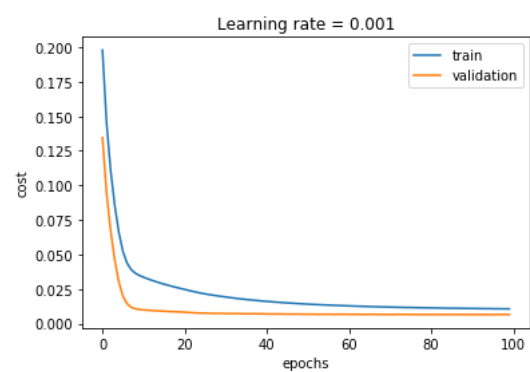


Lambd = 0.4

Training mape : 0.7367135963465501

Valid mape : 0.744504210476673

Testing mape : 0.7495242041036693

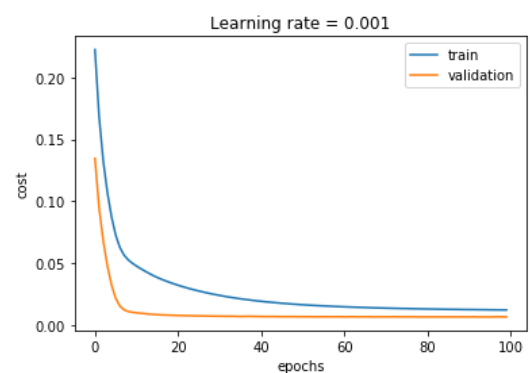


Lambd = 0.7

Training mape : 0.7581294783825394

Valid mape : 0.766588173560549

Testing mape : 0.7707128884610976

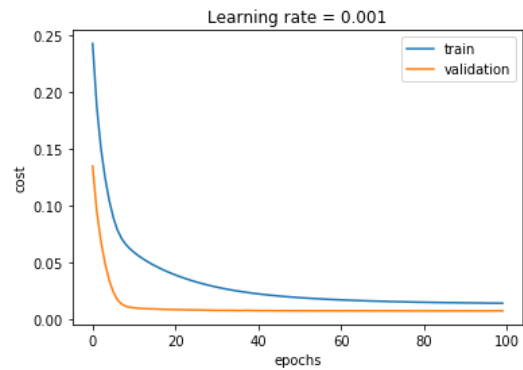


Lambda = 0.95

Training mape : 0.7716003838887064

Valid mape : 0.779187773992918

Testing mape : 0.7839908249893479



Learning Rate

- As Learning rate decreases we can observe the model learns slower
- More iterations will be required
- Higher the Learning rate higher the chances to encounter and fall into local minimas as we can see in the first graph. The graph is more noisy when dropping.
- Lower the learning rate smoother the graph

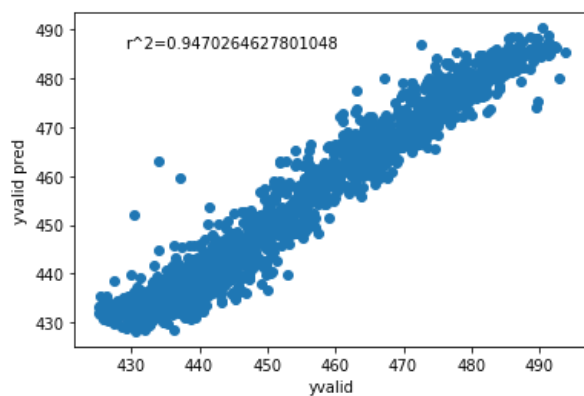
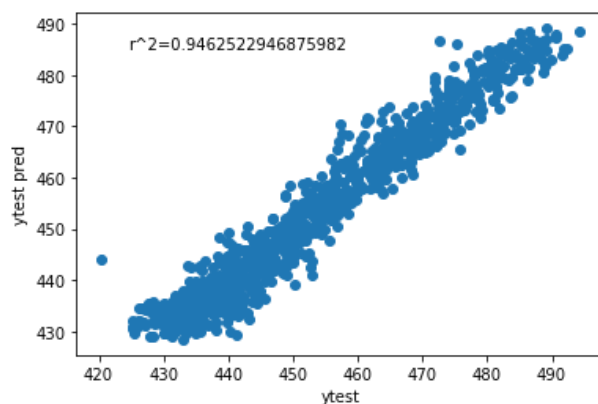
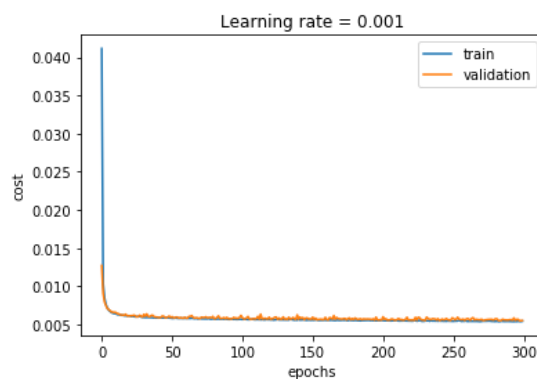
Regularisation

- As we increase lambda we decrease the rate of learning during training
- Validation loss is maintained constantly
- Convergence for higher lambda requires more iterations as we can see in maps vals.

Best Model for CCPP Dataset:

```
layers_dims = [4,20,20,1]
parameters = L_layer_model_minib(
X_train, y_train, layers_dims, valid=True, valid_x=X_valid, valid_y=y_valid,
num_iterations = 300,
he_init=True,
mini_batch_size = 64,
learning_rate = 0.001,
print_cost = True,
regularisation='none', lambd=0.1,
optimizer="adam",
beta = 0.9, beta1 = 0.9, beta2 = 0.999, epsilon = 1e-8,
activation='tanh',
cost_func='mse')
```

Training mape : 0.6552409400797429
Valid mape : 0.6688540464907443
Testing mape : 0.6783104901023306



Bonus :

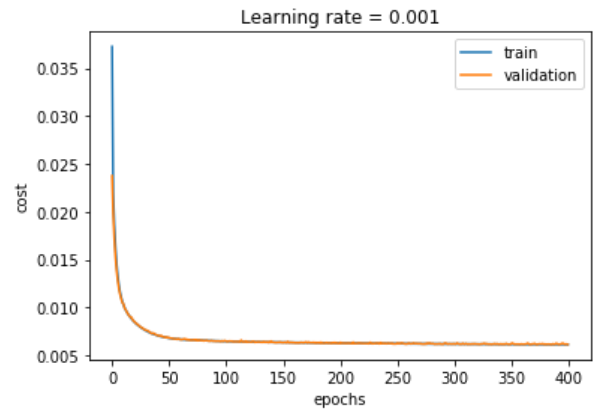
SGD with Momentum -

```
layers_dims = [4,10,10,1]
num_iterations = 400,
learning_rate = 0.001
mini_batch_size = 1
optimizer="momentum"
activation="tanh"
he_init=True
cost_func='mse'
regularisation ='none'
```

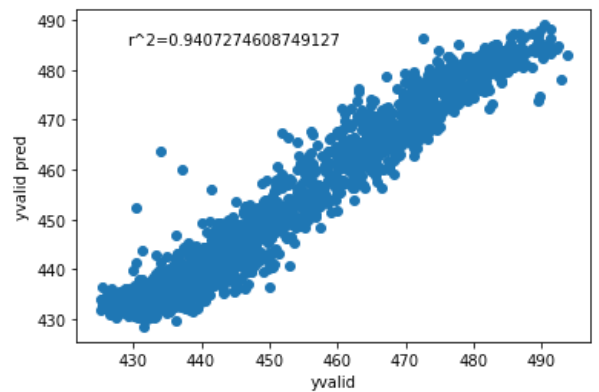
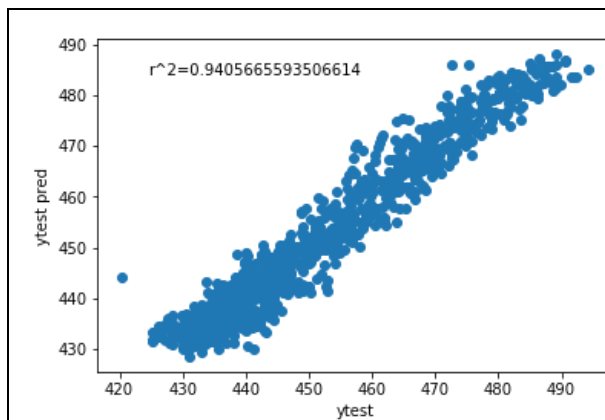
Training mape : 0.7107204433769398

Valid mape : 0.7157541114288213

Testing mape : 0.7271898861662556



R2-



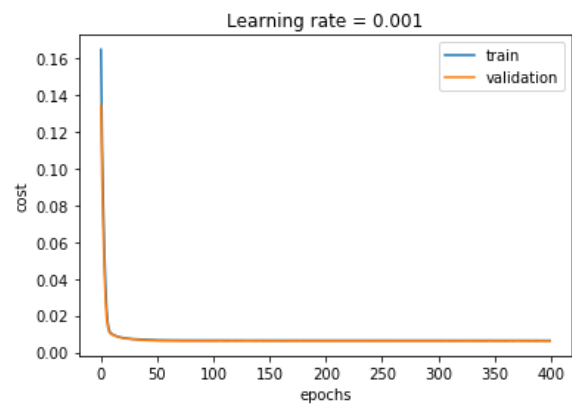
Adam with Batches-

```
layers_dims = [4,10,10,1]
num_iterations = 400,
learning_rate = 0.001
mini_batch_size = 256
optimizer="adam"
activation="tanh"
he_init=True
cost_func='mse'
regularisation ='L2'
```

Training mape : 0.7144425261228652

Valid mape : 0.7197032291132449

Testing mape : 0.7249401718028659



R2-

