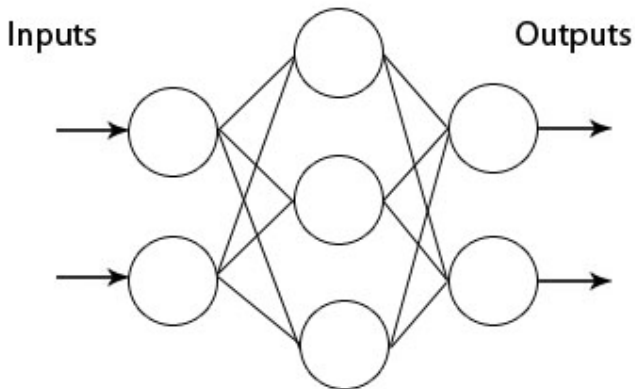


# Neural Networks (Part 2)

Matthew Galbraith & Mitchell Corbett

April 7, 2015

# Multilayer Perceptron Networks:



# Overview

## MLP: A type of feed-forward neural network

- Network composed of layers of neurons
- Each unit (neuron) in the a layer is directly connected to every unit in the following layer, starting from the input layer and ending at the output layer.
- The layers between the input and output layer are hidden layers

# Multilayer Perceptron Networks

## Overview (Continued)

- Whereas a true Perceptron is a binary classifier, Multilayer Perceptron networks are suited for both regression and classification problems.
- Clarification: MLP is not one perceptron, (although it's name makes it sound like that) but, a network composed of layers of perceptrons which are free to take any arbitrary activation function.

# Formal definition of a MLP with one hidden layer

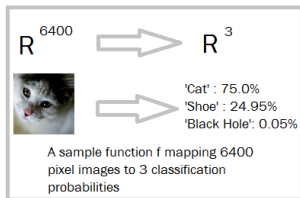
A one-hidden-layer MLP can be represented as a function

$$f : R^D \rightarrow R^L$$

$D$ : size of input vector  $x$

$L$ : size of the output vector  $f(x)$

Continued on chalkboard...



# So... What can we model with one-hidden-layer MLPs?

## EVERYTHING!

Well, not quite everything.. But when it comes to continuous functions...

### Universal Approximation Theorem

Summarized: A feed-forward network with a single hidden layer containing a finite number of neurons (Sound familiar?) can approximate continuous functions on compact subsets of  $R^n$ .

This means simple MLP networks such as the one I've been talking about can be used to model a wide variety of problems!

# Training the network

## High level pseudocode:

1. Initialize the weights
2. For a number of epochs, we:
  - 2.1. Forward-pass: We compute the output of the network
  - 2.1 We compute the loss (MSE) wrt our output
  - 2.2 Using backpropagation and an optimization method we work backwards to update the weights to minimize the error.

# Backpropagation

```
initialize network weights
loop:
forEach training example ex
  #forward-pass:
  prediction = network-output(ex)
  actual = correct-output(ex)
  compute loss at output units
  #backwards propagation:
  compute change in weights from hidden to output layer
  compute change in weights from input to hidden layer
  update network weights
if (stopping criterion satisfied):
  return network
```



# Update methods

## First order optimization algorithms

Stochastic Gradient Descent

Gradient Descent with Momentum

Gradient Descent with Nesterov Momentum

# Update methods - SGD

## Stochastic Gradient Descent

```
1 def sgd(loss, all_params, learning_rate):
2     all_grads = theano.grad(loss, all_params)
3     updates = []
4     for param_i, grad_i in zip(all_params, all_grads):
5         updates.append((param_i,
6                         param_i - learning_rate * grad_i))
7     return updates
```

# Update methods - Momentum

## Gradient Descent with Momentum

```
1 def momentum(loss, all_params, learning_rate, momentum=0.9):
2     all_grads = theano.grad(loss, all_params)
3     updates = []
4
5     for param_i, grad_i in zip(all_params, all_grads):
6         mparam_i = theano.shared(np.zeros(param_i.get_value().shape,
7                                             dtype=theano.config.floatX),
8                                   broadcastable=param_i.broadcastable)
9         v = momentum * mparam_i - learning_rate * grad_i
10        updates.append((mparam_i, v))
11        updates.append((param_i, param_i + v))
12
13    return updates
```

# Update methods - Nesterov Momentum

## Accelerated Gradient Descent with Nesterov Momentum

```
1  # using the alternative formulation of nesterov momentum described at
2  # https://github.com/lisa-lab/pylearn2/pull/136
3  # such that the gradient can be evaluated at the current parameters.
4  def nesterov_momentum(loss, all_params, learning_rate, momentum=0.9):
5      all_grads = theano.grad(loss, all_params)
6      updates = []
7
8      for param_i, grad_i in zip(all_params, all_grads):
9          mparam_i = theano.shared(np.zeros(param_i.get_value().shape,
10                                           dtype=theano.config.floatX),
11                                  broadcastable=param_i.broadcastable)
12          v = momentum * mparam_i - learning_rate * grad_i # new momentum
13          w = param_i + momentum * v - learning_rate * grad_i # new param values
14          updates.append((mparam_i, v))
15          updates.append((param_i, w))
16
17  return updates
```

# Moving on

## MLP applications

Now that I've introduced MLP Networks..

It's applications time!

# Facial Keypoint Recognition

## Problem:

Kaggle competition

## Observations:

7049 samples of 96x96 dimension images

Up to 15  $(x, y)$  pairs of coordinates for each image (Some missing)

## Model

We want to model some function

$$f : R^{9216} \rightarrow R^{30}$$

mapping our input image data to 15 pairs of  $(x, y)$  coordinates

## Python Packages used

```
from pandas.io.parsers import read_csv
import numpy as np
import cPickle as pickle
import os.path
import matplotlib.pyplot as plt
import theano
from lasagne import layers
from lasagne.nonlinearities import sigmoid
from lasagne.updates import nesterov_momentum
from nolearn.lasagne import NeuralNet
from sklearn.utils import shuffle
```

# Processing the data

## Load function courtesy of Daniel Nouri's tutorial

```
def load(test=False, cols=None, drop_missing=True):  
    #...  
    df = read_csv(os.path.expanduser(fname)) # load pandas dataframe  
  
    df['Image'] = df['Image'].apply(lambda im: np.fromstring(im, sep=' '))  
    if cols: # get a subset of columns  
        df = df[list(cols) + ['Image']]  
    if (drop_missing):  
        df = df.dropna() # drop all rows that have missing values  
    X = np.vstack(df['Image'].values) / 255. # scale pixel values to [0, 1]  
    X = X.astype(np.float32)  
  
    if not test: # only FTRAIN has any target columns  
        y = df[df.columns[:-1]].values  
        y = (y - 48) / 48 # scale target coordinates to [-1, 1]  
        X, y = shuffle(X, y, random_state=65539) # shuffle train data  
    else:  
        y = None
```



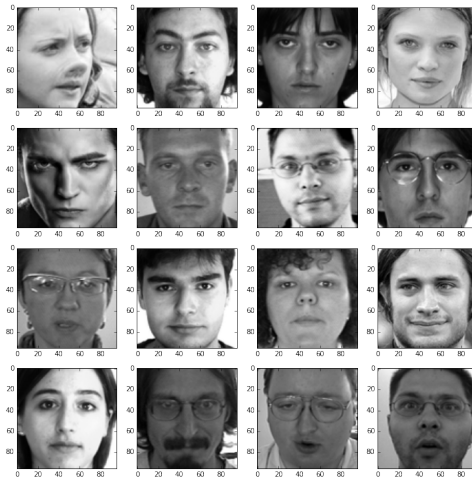
# Processing the data

Here's what loads..

```
>>>X, y = load()
>>>print summary(X,y)
X shape: (2140L, 9216L)
y shape: (2140L, 30L)
X: min, max = (0.0,1.0)
y: min, max = (-0.920286595821,0.996020495892)
Using 30.3589161583% of the provided training data
```

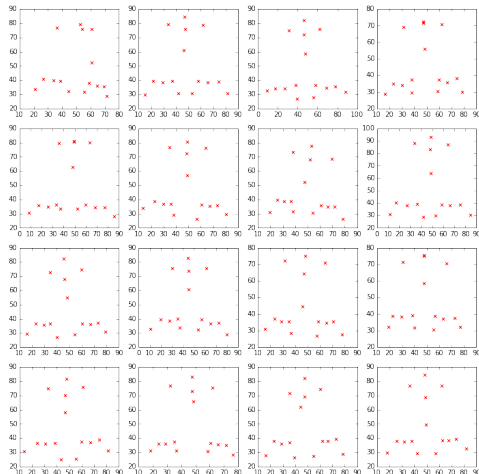
# Exploring the data

Facial Keypoints



# Exploring the data

Facial Keypoints



# Exploring the data

Facial Keypoints



# Creating NNets in Python

Here's the general workflow I followed

```
#Initialize the network
net = NeuralNet( #layers = ...
                 #more params = ...
                 #max_epochs = n
                )
X, y = load() #Load data
net.fit(X,y) #Train network
#analysis can go here...
X = load(test=false)[0] #Load test data
y_pred = net.predict(X) #Predict
#code to make Kaggle submission goes here...
```

# Using Early Stopping to reduce overfitting

```
1 class EarlyStopping(object):
2     def __init__(self, patience=100):
3         self.patience = patience
4         self.best_valid = np.inf
5         self.best_valid_epoch = 0
6         self.best_weights = None
7
8     def __call__(self, nn, train_history):
9         current_valid = train_history[-1]['valid_loss']
10        current_epoch = train_history[-1]['epoch']
11        if current_valid < self.best_valid:
12            self.best_valid = current_valid
13            self.best_valid_epoch = current_epoch
14            self.best_weights = [w.get_value() for w in nn.get_all_params()]
15        elif self.best_valid_epoch + self.patience < current_epoch:
16            print("Early stopping.")
17            print("Best valid loss was {:.6f} at epoch {}".format(
18                self.best_valid, self.best_valid_epoch))
19            nn.load_weights_from(self.best_weights)
20            raise StopIteration()
```

# Model 1

```
1  #Neural Network # 1 - A MLP with 1 hidden layer
2  net1 = NeuralNet(
3      layers=[ # three layers: one hidden layer
4          ('input', layers.InputLayer),
5          ('hidden', layers.DenseLayer),
6          ('output', layers.DenseLayer),
7      ],
8      # layer parameters:
9      input_shape=(None, 9216), # 96x96 input pixels
10     hidden_num_units=100, # number of units in hidden layer
11     hidden_nonlinearity=sigmoid,
12     output_nonlinearity=None, # output layer uses identity function
13     output_num_units=30, # 30 target values
14     # optimization params
15     update=nesterov_momentum,
16     update_learning_rate=0.01,
17     update_momentum=0.9,
18     regression=True,
19     max_epochs=1000,
20     on_epoch_finished=[
21         EarlyStopping(patience=50)
22     ],
23     verbose=1 #0 to not print anything
24 )
```

# Using gridsearchCV to determine appropriate number of hidden units

```
1  #parameter grid for gridsearch cv
2  param_grid = {
3  'more_params': [{'hidden_num_units': 100}, {'hidden_num_units':150}]...,
4  {'hidden_num_units': 300}]
5  }
6
7  #find net with best params , and then refit with all data on best net
8  gs = GridSearchCV(net1, param_grid, cv=2, refit=True, verbose=4)
9  X,y=load() # Load our data
10 gs.fit(X,y) #Fit our gridsearchCV object
11 with open('net1_gridsearch.pickle', 'wb') as f:
12     # we serialize the gridsearch model, as it also has the fitted
13     #neural network with best hidden unit number:
14     pickle.dump(gs, f, -1)
15     #if you just want to save the best estimator - ie - the neural network
16     #itself - without the gridsearch info:
17     pickle.dump(gs.best_estimator_, f, -1)
```



# Analyzing results of the gridsearchCV

```
1  >>>#Load the serialized gridsearchCV object which contains our network.
2  >>>net1_gridsearch = pickle.load(open('./net1_gridsearch.pickle', 'rb'))
3  >>>#This just shows the params the gridsearch considered when finding the best m
4  >>>net1_gridsearch.param_grid
5  {'more_params': [{'hidden_num_units': 100},
6    {'hidden_num_units': 150},
7    {'hidden_num_units': 200},
8    {'hidden_num_units': 250},
9    {'hidden_num_units': 300}]}
10 >>>#Here are the best parameters:
11 >>>print net1_gridsearch.best_params_
12 {'more_params': {'hidden_num_units': 100}}
```

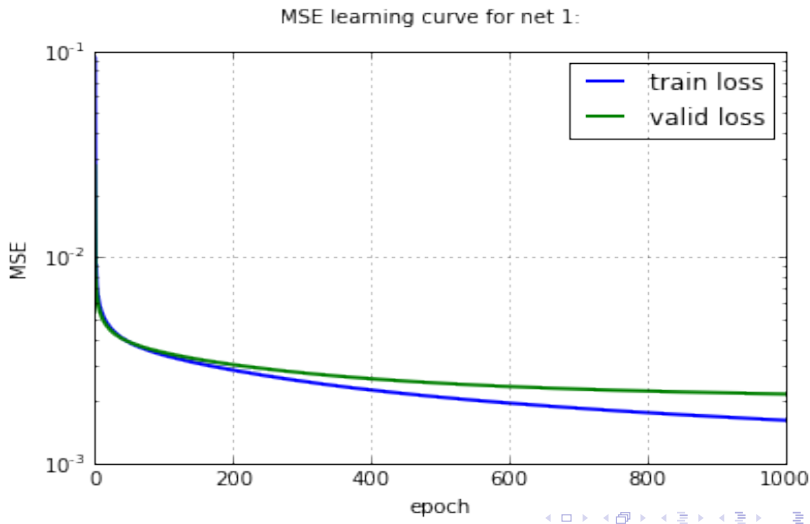
# Grid Search CV results

```
1 >>>print "The best score (Lowest MSE) was: " + str(net1_gridsearch.best_score_)
2 >>>print "This translates to a RMSE (Kaggle's evaluation) score of: " +
3 str(np.sqrt(net1_gridsearch.best_score_) * 48.)
4 The best score (Lowest MSE) was: 0.0025870305397
5 This translates to a RMSE (Kaggle evaluation) score of: 2.44141728581
```

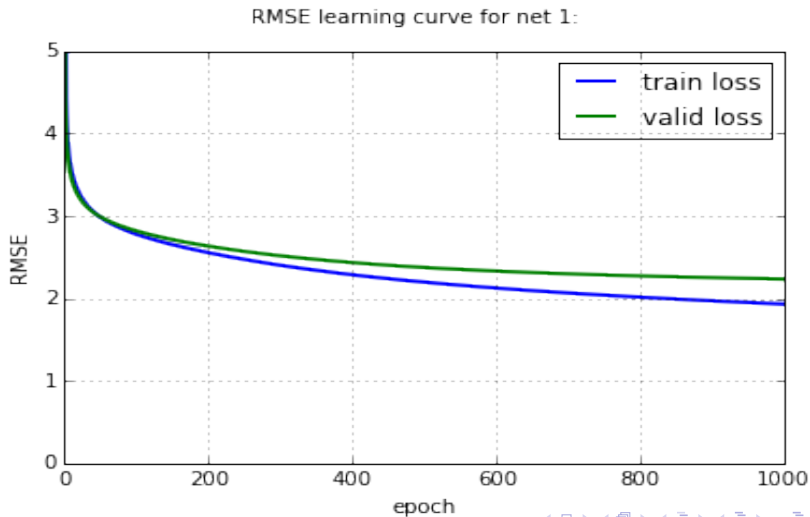
# Analyzing the first Neural Network model

```
1  #Let's load our neural net from the gridsearch object!
2  >>>net1 = net1_gridsearch.best_estimator_
3  >>>net1
4  NeuralNet(X_tensor_type=<function matrix at 0x...>,
5            batch_iterator_test=<nolearn.lasagne.BatchIterator object at 0x...>,
6            batch_iterator_train=<nolearn.lasagne.BatchIterator object at 0x...>,
7            eval_size=0.2,
8            hidden_nonlinearity=<theano.tensor.elemwise.Elemwise object at 0x...>,
9            hidden_num_units=100, input_shape=(None, 9216),
10           layers=[('input', <class 'lasagne.layers.input.InputLayer'>),
11                  ('hidden', <class 'lasagne.layers.dense.DenseLayer'>),
12                  ('output', <class 'lasagne.layers.dense.DenseLayer'>)],
13
14           loss=<function mse at 0x...>, max_epochs=1000,
15           more_params={'hidden_num_units': 100},
16           on_epoch_finished=[<__main__.EarlyStopping object at 0x...>],
17           on_training_finished=(), output_nonlinearity=None,
18           output_num_units=30, regression=True,
19           update=<function nesterov_momentum at 0x...>,
20           update_learning_rate=0.01, update_momentum=0.9,
21           use_label_encoder=False, verbose=1,
22           y_tensor_type=TensorType(float32, matrix))
```

# Model 1 - Learning Curve



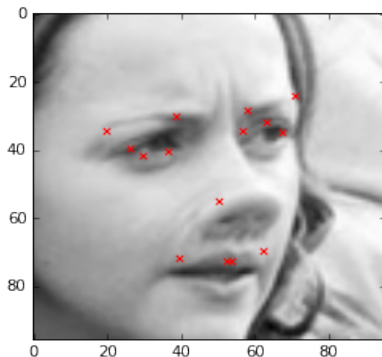
# Model 1 - Learning Curve



# Plotting a sample prediction

Taking a sample face (from the training set), here are predicted

net1: Facial Keypoints - Predicted

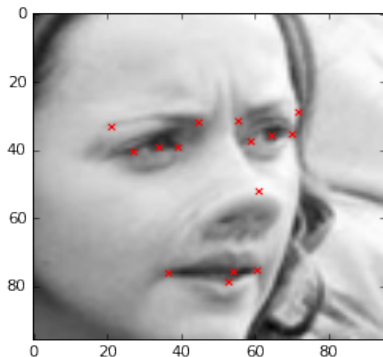


coordinates:

## Plotting a sample prediction –pt.2

Taking a sample face (from the training set), here are the actual

net1: Facial Keypoints - Actual

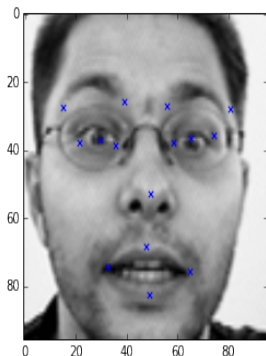


coordinates:

# Kaggle submission –pt.1

Here's a plot of predicted keypoints for a Kaggle test sample:

net1: Facial Keypoints - Prediction on Kaggle submission data (sample)





# Kaggle submission –pt.2

## Submitting to kaggle!

Not as simple as predicting and writing to csv...

Submission only wants specific coordinates from specific images in the test set

Need to filter our predictions by the coordinates Kaggle wants

# Kaggle submission –pt.3

## Submission made!

30	↓5	ashok k harnal	3.55854	22	Sun, 22 Mar 2015 12:49:35 (-2.1d)
31	new	rockers 🧑	3.58801	1	Thu, 02 Apr 2015 09:03:43
32	↓6	Jia Chen	3.58884	9	Sat, 14 Mar 2015 03:40:25 (-0.2h)
33	↓6	dscrimager	3.60199	1	Fri, 20 Feb 2015 04:16:55
34	↓6	kainster	3.62655	4	Sat, 14 Mar 2015 02:44:21 (-0.8h)
35	↑16	<b>mpg317</b>	<b>3.66788</b>	<b>7</b>	<b>Mon, 06 Apr 2015 03:37:16</b>

### Your Best Entry ↑

You improved on your best score by 0.49130.

You just moved up 25 positions on the leaderboard.

 Tweet this!

36	↓7	Suresh Kumar	3.75814	4	Fri, 13 Feb 2015 13:47:47 (-26.2h)
37	↓7	Elixir	3.77399	2	Fri, 13 Mar 2015 09:08:22 (-0.2h)
38	↓7	Manohar Swamynathan	3.80685	1	Thu, 12 Feb 2015 13:24:59

# Model 2

## Specialist Networks

- Idea: Train multiple networks on subsets of the 15 facial keypoints!
- Allows us to use more training data
- Result: 6 Specialized networks

# Dynamically updating learning rate and momentum coefficients

```
1  #http://danielnouri.org/notes/2014/12/17/using-convolutional-neural-nets-to-dete
2  class AdjustVariable(object):
3      def __init__(self, name, start=0.03, stop=0.001):
4          self.name = name
5          self.start, self.stop = start, stop
6          self.ls = None
7
8      def __call__(self, nn, train_history):
9          if self.ls is None:
10             self.ls = np.linspace(self.start, self.stop, nn.max_epochs)
11
12             epoch = train_history[-1]['epoch']
13             new_value = float32(self.ls[epoch - 1])
14             getattr(nn, self.name).set_value(new_value)
```

# The basis for our specialized neural networks

```
1 net = NeuralNet(  
2     layers=[ # three layers: one hidden layer  
3         ('input', layers.InputLayer),  
4         ('hidden', layers.DenseLayer),  
5         ('output', layers.DenseLayer),  
6     ],  
7     # layer parameters:  
8     input_shape=(None, 9216), # 96x96 input pixels per batch  
9     hidden_num_units=100, # number of units in hidden layer  
10    hidden_nonlinearity=sigmoid,  
11    output_nonlinearity=None, # output layer uses identity function  
12    output_num_units=30, # 30 target values for original, but we're changing  
13    # optimization params  
14    update=nesterov_momentum,  
15    update_learning_rate=theano.shared(float32(0.01)),  
16    update_momentum=theano.shared(float32(0.9)),  
17    regression=True,  
18    max_epochs=500,  
19    on_epoch_finished=[AdjustVariable('update_learning_rate', start=0.01, s  
20        AdjustVariable('update_momentum', start=0.9, stop=0.999),  
21        EarlyStopping(patience=200)  
22    ],  
23    verbose=1  
24 )
```

# the function to fit our specialized networks

```
1  for setting in SPECIALIST_SETTINGS:
2      cols = setting['columns']
3      X, y = load(cols=cols)
4
5      model = clone(net)
6      model.output_num_units = y.shape[1]
7      # set number of epochs relative to number of training examples:
8      model.max_epochs = int(1e7 / y.shape[0])
9      if 'kwargs' in setting:
10         # an option 'kwargs' in the settings list may be used to
11         # set any other parameter of the net:
12         vars(model).update(setting['kwargs'])
13
14     print("Training model for columns {} for up to {} epochs".format(
15         cols, model.max_epochs))
16     model.fit(X, y)
17     specialists[cols] = model
18
19 with open('net-specialists.pickle', 'wb') as f:
20     # we persist a dictionary with all models:
21     pickle.dump(specialists, f, -1)
```

# Analysis of the results - 1

```
1  #Load our 6 networks
2  net2_specialists = pickle.load(open('./net-specialists.pickle', 'rb'))
3  net2 = [x for x in net2_specialists.iteritems()]
4  #load networks into list
5  net2_nets = [x[1] for x in net2]
```

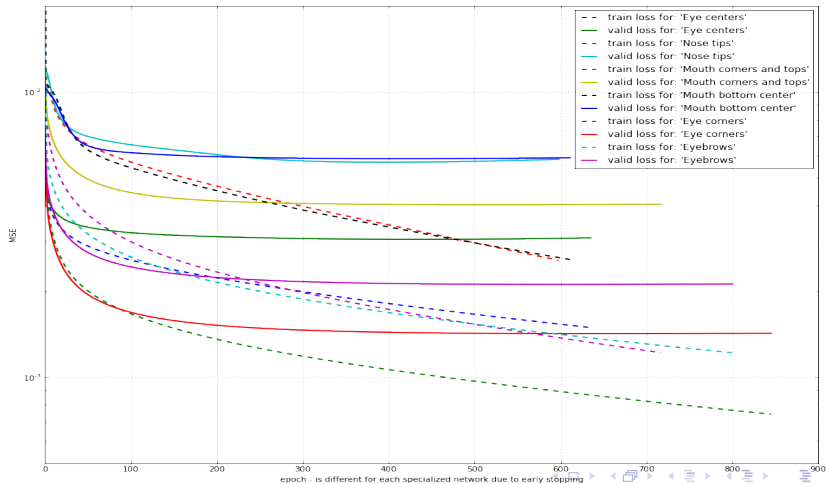
## Analysis of the results - 2

```
1  >>>for x,net in enumerate(net2):
2  >>>...print summary(x, net[0])
3  Specialized network #1 is trained on the keypoints:
4  ('left_eye_center_x', ...,
5   'right_eye_center_y')
6
7  Specialized network #2 is trained on the keypoints:
8  ('nose_tip_x', 'nose_tip_y')
9
10 Specialized network #3 is trained on the keypoints:
11 ('mouth_left_corner_x',...,
12 'mouth_right_corner_y', ..., 'mouth_center_top_lip_y')
13
14 Specialized network #4 is trained on the keypoints:
15 ('mouth_center_bottom_lip_x', 'mouth_center_bottom_lip_y')
16
17 Specialized network #5 is trained on the keypoints:
18 ('left_eye_inner_corner_x', ..., 'right_eye_outer_corner_y')
19
20 Specialized network #6 is trained on the keypoints:
21 ('left_eyebrow_inner_end_x', ..., 'right_eyebrow_outer_end_y')
```



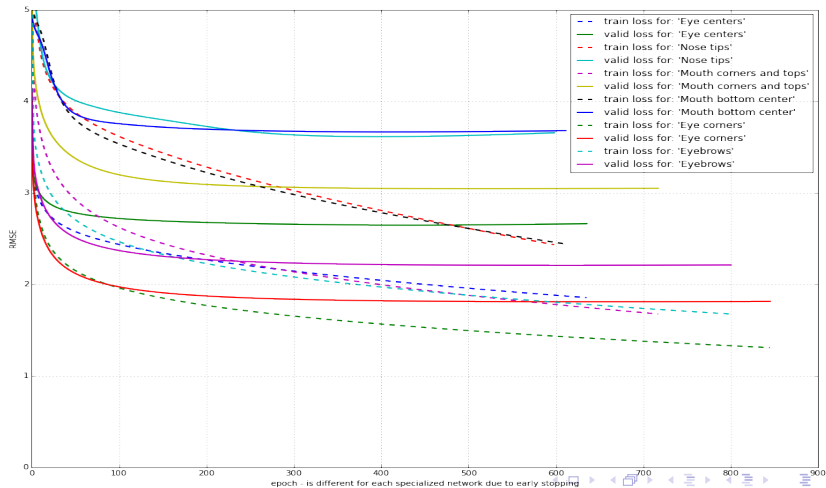
# Analysis of the results - MSE Learning Curve

MSE learning curve for specialized networks



# Analysis of the results - RMSE Learning Curve

RMSE learning curve for specialized networks



# Analysis of the results:

```
labels|train score|valid score
'Eye centers'|0.00174424204346|0.00304428217138
'Nose tips'|0.00344257359455|0.005985269282
'Mouth corners and tops'|0.00159532899415|0.00399073439391
'Mouth bottom center'|0.00331330332695|0.00583614665805
'Eye corners'|0.000931936667851|0.00145246363777
'Eyebrows'|0.00142408664489|0.00218300545633
```

# Submitting to kaggle

```
1 preds = []
2 For each specialized network, net:
3     Load X data for cols corresponding to net preds.append([net.predict(X)])
4     #rearrange indices to follow same order as kaggle data
5     preds[i] = preds[i][:, 0, 1, 4, 5, 2, 3, 6, 7]
6     #stack predictions into one large array
7     preds = np.hstack(preds)
8     #create submission the same way as submission 1
```

# Submitting to kaggle: results

## Submission made!

18	↓4	Carlos Mattoso	3.03034	5	Mon, 16 Mar 2015 14:26:15 (-0.3h)
19	↓4	Florian Muellerklein	3.04421	3	Fri, 13 Mar 2015 11:14:20
20	↑31	<b>mpg317</b>	<b>3.07987</b>	<b>9</b>	<b>Mon, 06 Apr 2015 04:20:15</b>

### Your Best Entry ↑

You improved on your best score by 0.58801.

You just moved up 15 positions on the leaderboard.



21	↓5	Il UWr 🤖	3.15571	5	Mon, 16 Feb 2015 18:28:53
22	↓5	Steven Durand	3.16545	4	Sun, 22 Mar 2015 15:10:52 (-6.7d)
23	↓5	nitinag	3.21959	4	Tue, 17 Mar 2015 06:02:40 (-3d)
24	↓2	0xBBBB	3.26921	3	Tue, 31 Mar 2015 04:23:06

# Conclusion

## Neural Networks

Complex! A whole course could be devoted to neural network theory

Can be applied to many different problems

# References

Useful code snippets:

<http://danielnouri.org/notes/2014/12/17/using-convolutional-neural-nets-to-detect-facial-keypoints-tutorial/>

1-hidden-layer MLP Mathematical definition -

<http://deeplearning.net/tutorial/mlp.html> Universal Approximation Theorem - <http://deeplearning.cs.cmu.edu/notes/SoniaHornik.pdf>

Theano - <http://deeplearning.net/software/theano/>

Lasagne - <http://lasagne.readthedocs.org/en/latest/>

Nolearn - <https://pythonhosted.org/nolearn/>

# references...

Nesterov Momentum -

<http://www.cs.toronto.edu/~fritz/absps/momentum.pdf>

Gradient Descent w/ Momentum -

<http://brahms.cpmc.columbia.edu/publications/momentum.pdf>

Early Stopping -

<http://papers.nips.cc/paper/1895-overfitting-in-neural-nets-backpropagation-conjugate-gradient-and-early-stopping.pdf>

Initial Weights -

<http://stats.stackexchange.com/questions/47590/what-are-good-initial-weights-in-a-neural-network>