# Introduction to Machine Learning: Linear and Logistic Regression and Neural Networks Using Python

# 6 Optimization, Validation, Accuracy & Regularization

# Optimization, Validation, Accuracy & Regularization

| Section | Title | Description |
| --- | --- | --- |
| 6.1 | Stochastic and Mini-Batch Gradient Descent | This section presents the stochastic Gradient Descent and mini-batch optimization algorithm |
| 6.2 | Momentum, RmsProp and Adam Optimization | This section presents optimizers that are alternatives to the Gradient Descent |
| 6.3 | Code Walkthrough Version 3.1 | Code walkthrough of optimization updates to the framework |
| 6.4 | Validation | This section describes using a validation dataset to measure performance |
| 6.5 | Performance Measures | This section describes metrics in addition to accuracy for measuring performance |
| 6.6 | Code Walkthrough Version 3.2 | This section presents the updates to the machine learning framework for validation and performance measures |
| 6.7 | Regularization | This section describes the purpose and approaches for regularization |
| 6.8 | Hyperparameter Search | This section describes the hyperparameter search process for tuning a machine learning system |
| 6.9 | Code Walkthrough Version 3.3 | This section presents the updates to the machine learning framework for regularization and hyperparameter searches |
| 6.10 | Underfitting and Overfitting | This section defines underfitting and overfitting and discusses how to diagnose and remedy these conditions |

# 6.1 Stochastic and Mini-Batch Gradient Descent

# Stochastic and Mini-Batch Gradient Descent

Goal of this Section:

- Present algorithms for Stochastic and Mini-Batch Gradient Descent

# Batch versus Stochastic Gradient Descent

- Gradient Descent algorithm used for Training presented earlier is often referred to as "Batch" Gradient Descent
  - Called Batch because at each epoch gradient calculation/update to W and b is based on all data samples at once
- Alternative: Stochastic Gradient Descent – at each epoch compute gradient/update using one training data sample at a time, looping over all data samples

# Training Algorithm: Stochastic Gradient Descent

Assume Neural Network with N layers

Input training data: feature matrix X and values Y

Make initial guess for parameters: $W^{[k]}$ and $b^{[k]}$ for k=1,…,N

Choose learning rate $\alpha > 0$

To simplify notation, remove subscript epoch=

1. Loop for epoch i = 1, 2, …
   - Randomly permute samples
   - Loop over samples j=0,…,m-1
     - Forward Propagate using feature vector $X_j$: compute $A_j^{[k]}$ for k=1,…,N
     - Back Propagate using $X_j$, $Y_j$: compute $\nabla_{W^{[k]}} L_j$, $\nabla_{b^{[k]}} L_j$ for k=1,…,N
     - Update parameters: for k=1,…,N
       - $W^{[k]} \leftarrow W^{[k]} - \alpha \nabla_{W^{[k]}} L_j$
       - $b^{[k]} \leftarrow b^{[k]} - \alpha \nabla_{b^{[k]}} L_j$
   - Forward Propagate using X: compute $A^{[k]}$ for k=1,…,N
   - Predict: compute $Y_{pred}$ and compute accuracy comparing with $Y$
   - Compute loss using $A^{[N]}$

Loop for fixed number of epochs

# Stochastic Gradient Descent - Example

- Consider Logistic Regression and data (2 features and 2 data points)

$$X = \begin{bmatrix} 1 & 2 \\ -2 & -5 \end{bmatrix} \qquad Y = [0 \quad 1]$$

$$X_{sample=0} = \begin{bmatrix} 1 \\ -2 \end{bmatrix} \qquad Y_{sample=0} = [0] \qquad X_{sample=1} = \begin{bmatrix} 2 \\ -5 \end{bmatrix} \qquad Y_{sample=1} = [1]$$

- Assume that initial parameter values are:

$$W = [0.1 \quad 0.1] \qquad b = [0.2]$$

- Assume binary cross entropy loss

- Assume training is performed using Gradient Descent with $\alpha = 0.1$

- Epoch 1 – Sample Index = 0

- Forward Propagation:

$$Z_{sample=0} = WX_{sample=0} + b = [0.1 \quad 0.1]\begin{bmatrix} 1 \\ -2 \end{bmatrix} + [0.2] = [0.1]$$

$$A_{sample=0} = f(Z_{sample=1}) = \left[\frac{1}{1+e^{-0.1}}\right] = [0.5250]$$

# Stochastic Gradient Descent - Example

- Back Propagation: Recall for binary cross entropy: $\nabla_A L = \begin{bmatrix} \frac{\partial L}{\partial A_0} & \frac{\partial L}{\partial A_1} & \cdots & \frac{\partial L}{\partial A_{m-1}} \end{bmatrix}, \frac{\partial L}{\partial A_j} = -\frac{1}{m}\left[\frac{Y_j}{A_j} - \frac{1-Y_j}{1-A_j}\right]$

- In this example m =1, and we use the subscript sample=0

$$\nabla_A L_{sample=0} = -\frac{1}{1}\left(\frac{Y_{sample=0}}{A_{sample=0}} - \frac{1-Y_{sample=0}}{1-A_{sample=0}}\right) = -\frac{1}{1}\left[-\frac{1}{1-0.5250}\right] = [2.1052]$$

$$\frac{\partial A_{sample=0}}{\partial Z_{sample=0}} = A_{sample=0} - A^2_{sample=0} = [0.2494]$$

$$\nabla_Z L_{sample=0} = \nabla_A L_{sample=0} * \frac{\partial A_{sample=0}}{\partial Z_{sample=0}} = [2.1052] * [0.2494] = [0.5250]$$

$$\nabla_W L_{sample=0} = \nabla_Z L_{sample=0} X^T_{sample=0} = [0.5250][1 \quad -2] = [0.5250 \quad -1.0500]$$

$$\nabla_b L_{sample=0} = \nabla_Z L_{sample=0} = [0.5250] \text{ (no sum here as } \nabla_Z L_{sample=0} \text{ has only 1 entry)}$$

- Update the parameters:

$$W = W - \alpha \nabla_W L_{sample=0} = [0.1 \quad 0.1] - 0.1*[0.5250 \quad -1.0500] = [0.0475 \quad 0.2050]$$

$$b = b - \alpha \nabla_b L_{sample=0} = [0.2] - 0.1*[0.5250] = [0.1475]$$

# Stochastic Gradient Descent - Example

- Epoch 1 – Sample Index=1
- Forward Propagation:

$$Z_{sample=1} = WX_{sample=1} + b = [0.0475 \quad 0.2050]\begin{bmatrix} 2 \\ -5 \end{bmatrix} + [0.1475] = [-0.7825]$$

$$A_{sample=1} = f(Z_{sample=1}) = \left[\frac{1}{1+e^{0.7825}}\right] = [0.3138]$$

- Back Propagation

$$\nabla_A L_{sample=1} = -\frac{1}{1}\left(\frac{Y_{sample=1}}{A_{sample=1}} - \frac{1-Y_{sample=1}}{1-A_{sample=1}}\right) = -\frac{1}{1}\left[\frac{1}{0.3138}\right] = [-3.1869]$$

$$\frac{\partial A_{sample=1}}{\partial Z_{sample=1}} = A_{sample=1} - A^2_{sample=1} = [0.2153]$$

$$\nabla_Z L_{sample=1} = \nabla_A L_{sample=1} * \frac{\partial A_{sample=1}}{\partial Z_{sample=1}} = [-3.1869] * [0.2153] = [-0.6862]$$

$$\nabla_W L_{sample=1} = \nabla_Z L_{sample=1} X^T_{sample=1} = [-0.6826][2 \quad -5] = [-1.3724 \quad 3.4311]$$

$$\nabla_b L_{sample=1} = \nabla_Z L_{sample=1} = [-0.6862] \text{ (no sum as } \nabla_Z L_{sample=1} \text{ has only 1 entry)}$$

- Update the parameters:

$$W = W - \alpha \nabla_W L_{sample=1} = [0.0475 \quad 0.2050] - 0.1*[-1.3724 \quad 3.4311] = [0.1847 \quad -0.1381]$$

$$b = b - \alpha \nabla_b L_{sample=1} = [0.1475] - 0.1*[-0.6862] = [0.2161]$$

# Stochastic Gradient Descent - Example

- Example shows 1 epoch of Stochastic Gradient Descent Algorithm
- In that epoch, W and b are updated twice (since there are 2 samples – there would be m updates in case of m samples)
- Process is repeated for epoch = 2, 3, 4, …

# Mini-Batch Gradient Descent

- Mini-Batch Gradient is compromise between Stochastic and Batch Gradient Descent

- Split training data into mini-batches (eg: 50 samples per mini-batch)

- At each epoch, loop over mini-batches and compute update to parameters one mini-batch at a time

- Stochastic Gradient Descent is equivalent to Mini-Batch Gradient Descent with 1 sample per batch

# Mini-Batch Creation Algorithm

Assume m training samples

Choose mini-batch size  b

Let n = number of mini-batches
- n = m/b if b goes into m evenly
- n = ceil(m/b) if b doesn't go into m evenly
  (ceil() always rounds up to next integer eg: ceil(4.1)=5)

1. Randomly permute samples
2. Split training samples into n mini-batches
   - 1$^{st}$ b samples for mini-batch 0,
   - 2$^{nd}$ b samples for mini-batch 1,
   - …
   - Final mini-batch may have less than b samples

# Mini-Batch Creation - Example

- Consider a case of 2 features and 5 data points

$$X = \begin{bmatrix} 1 & 2 & 4 & 1 & -1 \\ -2 & -5 & -8 & -2 & 2 \end{bmatrix} \qquad Y = \begin{bmatrix} 0 & 1 & 0 & 1 & 1 \end{bmatrix}$$

- Assume a mini-batch size of 3

- 3 doesn't go into 5 evenly so, there will be 2 mini-batches (3 & 2 data points)

- Split into 2 mini-batches:

$$X_{minibatch=0} = \begin{bmatrix} 1 & 2 & 4 \\ -2 & -5 & -8 \end{bmatrix} \qquad Y_{minibatch=0} = \begin{bmatrix} 0 & 1 & 0 \end{bmatrix}$$

$$X_{minibatch=1} = \begin{bmatrix} 1 & -1 \\ -2 & 2 \end{bmatrix} \qquad Y_{minibatch=1} = \begin{bmatrix} 1 & 1 \end{bmatrix}$$

# Training Algorithm: Mini-Batch Gradient Descent

Assume Neural Network with N layers

Input training data: feature matrix X and values Y

Choose batch size b

Make initial guess for parameters: $W^{[k]}$ and $b^{[k]}$ for k=1,…,N

Choose learning rate $\alpha > 0$

To simplify notation, remove subscript epoch=

1. Loop for epoch i = 1, 2, …
   - Generate mini-batches $(X_{minibatch=j}, Y_{minibatch=j})$, j=0,…,n-1
   - Loop over mini-batches j=0,…,n-1
     - Forward Propagate using $X_{minibatch=j}$: compute $A^{[k]}_{minibatch=j}$ for k=1,…,N
     - Back Propagate using $X_{minibatch=j}, Y_{minibatch=j}$: compute $\nabla_{W^{[k]}} L_{minibatch=j}, \nabla_{b^{[k]}} L_{minibatch=j}$ k=1,…,N
     - Update parameters for k=1,…,N
       - $W^{[k]} \leftarrow W^{[k]} - \alpha \nabla_{W^{[k]}} L_{minibatch=j}$
       - $b^{[k]} \leftarrow b^{[k]} - \alpha \nabla_{b^{[k]}} L_{minibatch=j}$
   - Forward Propagate using X: compute $A^{[k]}$ for k=1,…,N
   - Predict: compute $Y_{pred}$ and compute accuracy comparing with $Y$
   - Compute loss using $A^{[N]}$

Loop for fixed number of epochs

# Training with Mini-Batch Gradient Descent - Example

- Consider Logistic Regression and data (2 features and 5 data points)
- Assume mini-batches have been created as in Mini-Batch Creation example:

$$X_{minibatch=0} = \begin{bmatrix} 1 & 2 & 4 \\ -2 & -5 & -8 \end{bmatrix} \qquad Y_{minibatch=0} = \begin{bmatrix} 0 & 1 & 0 \end{bmatrix}$$

$$X_{minibatch=1} = \begin{bmatrix} 1 & -1 \\ -2 & 2 \end{bmatrix} \qquad Y_{minibatch=1} = \begin{bmatrix} 1 & 1 \end{bmatrix}$$

- Assume that initial parameter values are:

$$W = \begin{bmatrix} 0.1 & 0.1 \end{bmatrix} \qquad b = \begin{bmatrix} 0.2 \end{bmatrix}$$

- Assume binary cross entropy loss
- Assume training is performed using Gradient Descent with $\alpha=0.1$
- Epoch 1 - Mini-batch 0
- Forward Propagation:

$$Z_{minibatch=0} = WX_{minibatch=0} + b = \begin{bmatrix} 0.1 & 0.1 \end{bmatrix} \begin{bmatrix} 1 & 2 & 4 \\ -2 & -5 & -8 \end{bmatrix} + \begin{bmatrix} 0.2 \end{bmatrix} = \begin{bmatrix} 0.1 & -0.1 & -0.2 \end{bmatrix}$$

$$A_{minibatch=0} = f(Z_{minibatch=0}) = \begin{bmatrix} \frac{1}{1+e^{-0.1}} & \frac{1}{1+e^{0.1}} & \frac{1}{1+e^{0.2}} \end{bmatrix} = \begin{bmatrix} 0.5250 & 0.4750 & 0.4502 \end{bmatrix}$$

# Training with Mini-Batch Gradient Descent - Example

- Back Propagation: Recall: $\nabla_A L = \begin{bmatrix} \frac{\partial L}{\partial A_0} & \frac{\partial L}{\partial A_1} & \cdots & \frac{\partial L}{\partial A_{m-1}} \end{bmatrix}, \frac{\partial L}{\partial A_j} = -\frac{1}{m} \begin{bmatrix} \frac{Y_j}{A_j} - \frac{1-Y_j}{1-A_j} \end{bmatrix}$

$\nabla_A L_{minibatch=0} = -\frac{1}{3} \left( \frac{Y_{minibatch=0}}{A_{minibatch=0}} - \frac{1-Y_{minibatch=0}}{1-A_{minibatch=0}} \right) = -\frac{1}{3} \left[ -\frac{1}{1-0.5250} \quad \frac{1}{0.4750} \quad -\frac{1}{1-0.4502} \right]$

$= \begin{bmatrix} 0.7017 & -0.7017 & 0.6062 \end{bmatrix}$

$\frac{\partial A_j}{\partial Z_j} = A_j - A_j^2 \quad \begin{bmatrix} \frac{\partial A_0}{\partial Z_0} & \frac{\partial A_1}{\partial Z_1} & \frac{\partial A_2}{\partial Z_2} \end{bmatrix} = \begin{bmatrix} 0.2494 & 0.2494 & 0.2475 \end{bmatrix}$

$\nabla_Z L_{minibatch=0} = \nabla_A L_{minibatch=0} * \begin{bmatrix} \frac{\partial A_0}{\partial Z_0} & \frac{\partial A_1}{\partial Z_1} & \frac{\partial A_2}{\partial Z_2} \end{bmatrix} = \begin{bmatrix} 0.7017 & -0.7017 & 0.6062 \end{bmatrix} *$
$\begin{bmatrix} 0.2494 & 0.2494 & 0.2475 \end{bmatrix} = \begin{bmatrix} 0.1750 & -0.1750 & 0.1501 \end{bmatrix}$

$\nabla_W L_{minibatch=0} = \nabla_Z L_{minibatch=0} X^T_{minibatch=0} = \begin{bmatrix} 0.1750 & -0.1750 & 0.1501 \end{bmatrix} \begin{bmatrix} 1 & -2 \\ 2 & -5 \\ 4 & -8 \end{bmatrix}$

$= \begin{bmatrix} 0.4252 & -0.6755 \end{bmatrix}$

$\nabla_b L_{minibatch=0} = \sum_{j=0}^{m-1} \nabla_Z L_{minibatch=0,j} = 0.1501$ (sum of entries of $\nabla_Z L_{minibatch=0}$)

- Update the parameters:

$W = W - \alpha \nabla_W L_{minibatch=0} = \begin{bmatrix} 0.1 & 0.1 \end{bmatrix} - 0.1*\begin{bmatrix} 0.4252 & -0.6755 \end{bmatrix} = \begin{bmatrix} 0.0575 & 0.1675 \end{bmatrix}$

$b = b - \alpha \nabla_b L_{minibatch=0} = \begin{bmatrix} 0.2 \end{bmatrix} - 0.1*\begin{bmatrix} 0.1501 \end{bmatrix} = \begin{bmatrix} 0.1850 \end{bmatrix}$

# Training with Mini-Batch Gradient Descent - Example

- Epoch 1 - Mini-batch 1
- Forward Propagation:

$$Z_{minibatch=1} = WX_{minibatch=1} + b = [0.0575 \quad 0.1675]\begin{bmatrix} 1 & -1 \\ -2 & 2 \end{bmatrix} + [0.1850] = [-0.0926 \quad 0.4626]$$

$$A_{minibatch=1} = f(Z_{minibatch=1}) = \left[\frac{1}{1+e^{0.0926}} \quad \frac{1}{1+e^{-0.4626}}\right] = [0.4769 \quad 0.6136]$$

- Back Propagation

$$\nabla_A L_{minibatch=1} = -\frac{1}{2}\left(\frac{Y_{minibatch=1}}{A_{minibatch=1}} - \frac{1-Y_{minibatch=1}}{1-A_{minibatch=1}}\right) = -\frac{1}{2}\left[\frac{1}{0.4769} \quad \frac{1}{0.6136}\right] = [-1.0485 \quad -0.8148]$$

$$\frac{\partial A_j}{\partial Z_j} = A_j - A_j^2 \quad \left[\frac{\partial A_0}{\partial Z_0} \quad \frac{\partial A_1}{\partial Z_1}\right] = [0.2495 \quad 0.2371]$$

$$\nabla_Z L_{minibatch=1} = \nabla_A L_{minibatch=1} * \left[\frac{\partial A_0}{\partial Z_0} \quad \frac{\partial A_1}{\partial Z_1}\right] = [-1.0485 \quad -0.8148] * [0.2495 \quad 0.2371] = [-0.2616 \quad -0.1932]$$

$$\nabla_W L_{minibatch=1} = \nabla_Z L_{minibatch=1} X_{minibatch=1}^T = [-0.2616 \quad -0.1932]\begin{bmatrix} 1 & -2 \\ -1 & 2 \end{bmatrix} = [-0.0684 \quad 0.1368]$$

$$\nabla_b L_{minibatch=1} = \sum_{j=0}^{m-1} \nabla_Z L_{minibatch=1,j} = -0.4548 \quad \text{(sum of entries of } \nabla_Z L_{minibatch=1})$$
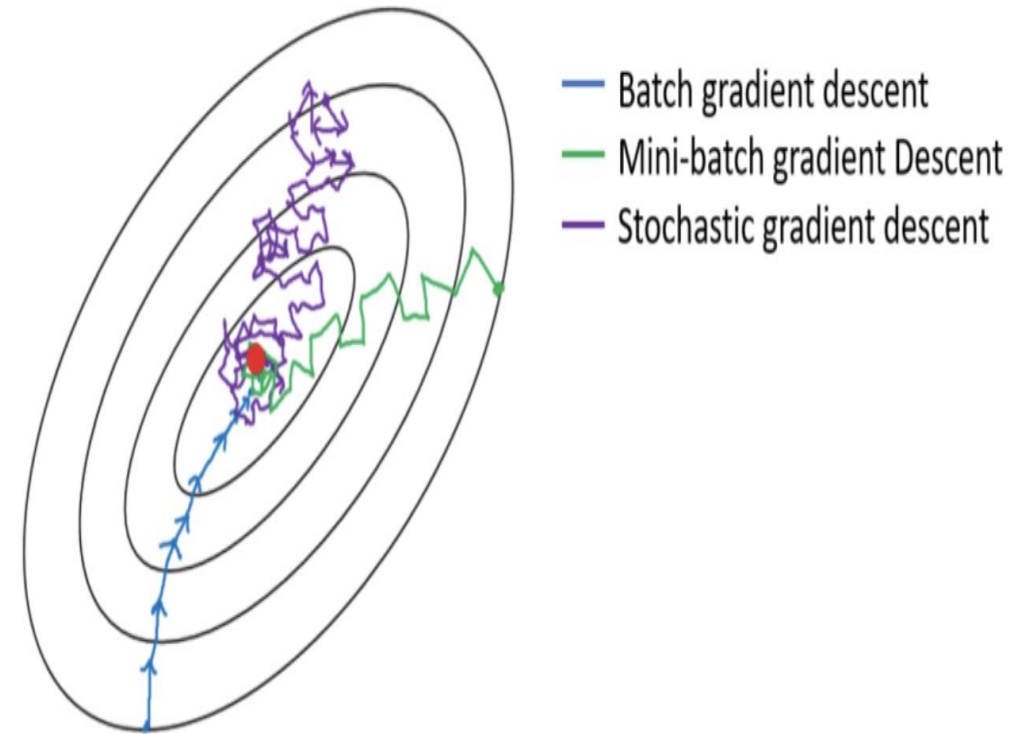
- Update the parameters:

$$W = W - \alpha \nabla_W L_{minibatch=1} = [0.0575 \quad 0.1675]\text{-}0.1*[-0.0684 \quad 0.1368] = [0.0643 \quad 0.1539]$$

$$b = b - \alpha \nabla_b L_{minibatch=1} = [0.1850]\text{-}0.1*[-0.4548] = [0.2305]$$

# Summary

| Batch Gradient Descent (BGD) | Stochastic Gradient Descent (SGD) | Mini-Batch Gradient Descent (MGD) |
|---|---|---|
| • Update $W^{[k]}$ and $b^{[k]}$ once per epoch<br>• Computationally efficient, in general, but may be issues with large datasets<br>• Usually stable learning path<br>• May converge to a local minimum if loss function not convex | • Update $W^{[k]}$ and $b^{[k]}$ m (number of samples) times per epoch<br>• Typically more computationally expensive than BGD<br>• Noisiest learning path<br>• Noisy learning process may help avoid local minimum | • Update $W^{[k]}$ and $b^{[k]}$ Ceil(m/batchSize) times per epoch<br>• Computational expense between SGD and BGD<br>• Noisy learning path<br>• Noisy learning process may help avoid local minimums |



— Batch gradient descent
— Mini-batch gradient Descent
— Stochastic gradient descent

Picture source (with permission):
https://towardsdatascience.com/gradient-descent-algorithm-and-its-variants-10f652806a3

# 6.1 Stochastic and Mini-Batch Gradient Descent - Jupyter Notebook DEMO

- See notebooks
  - IntroML/Examples/Chapter6/StochasticGradientDescent.ipynb
  - IntroML/Examples/Chapter6/MiniBatchGradientDescent.ipynb

Course Resources at:

- [https://github.com/satishchandrareddy/IntroML/](https://github.com/satishchandrareddy/IntroML/)

# 6.2 Momentum, RmsProp, and Adam Optimization

# Momentum, RmsProp, and Adam Optimization

Goal of this Section:

- Present additional optimization algorithms:
  - Momentum
  - RmsProp
  - Adam

# Optimization

- Gradient Descent is a rudimentary optimization algorithm
- Not guaranteed to converge to global minimum or may converge slowly or may get stuck near local minimum or not converge at all
- Has fixed learning rate $\alpha$ – may want to use an adaptive learning rate
- Enhancements to Gradient Descent have been developed to address these issues and improve convergence
- In this section discuss additional optimization approaches
- Can apply mini-batch optimization to these methods as well

# Optimization Algorithm

Let W = [$W_0$ $W_1$ $W_2$ $W_{d-1}$] denote parameter vector of variables. Let L(W)= L($W_0$,$W_1$,$W_2$ ,…, $W_{d-1}$) be a function of these parameters

Make initial guess $W_{epoch=0}$ for parameters

1. Loop epoch i = 1, 2, 3, …

- Compute gradient vector $\nabla_W L_{epoch=i-1}$ at $W_{epoch=i-1}$

- Compute $\boxed{\text{Update}_{epoch=i}}$ $\longleftarrow$ <span style="color:red">Depends on approach</span>

- Update W: $W_{epoch=i}$=$W_{epoch=i-1}$ + Update$_{epoch=i}$

- Compute $L(W_{epoch=i})$

Loop for fixed number of epochs

# Momentum

- Momentum algorithm has a 2-step update formula:

$$v_{epoch=i} = \beta v_{epoch=i-1} + \nabla_W L_{epoch=i-1} \qquad v_{epoch=0} = 0$$
$$Update_{epoch=i} = -\alpha v_{epoch=i}$$

where $0 \leq \beta < 1$ and $\alpha > 0$

- $v$ is a weighted average of gradients at all previous epochs
- As $\beta < 1$, weight of gradient at early epochs decreases as epochs increase
- $\beta = 0$ case corresponds to Gradient Descent
- Motivation for momentum is to avoid getting stuck in a local minimum
  - See blogposts on Momentum approach for more details

Reference:
Rumelhart, David E.; Hinton, Geoffrey E.; Williams, Ronald J. (8 October 1986). "Learning representations by back-propagating errors". Nature. 323 (6088): 533–536. Bibcode:1986Natur.323..533R. doi:10.1038/323533a0.

# Momentum - Example

- Consider simple function (has minimum at [0 0])

$L(W) = L(W_0, W_1) = 2W_0^2 + W_1^2$ with gradient $\nabla_W L = [4W_0 \quad 2W_1]$

- Choose $\alpha$=0.1, $\beta$=0.9, and set $v_{epoch=0}$=0 and make initial guess

$W_{epoch=0} = [2\ 2] \quad L(W_{epoch=0}) = 2 * 2^2 + 2^2 = 12$

- Epoch 1:

$\nabla_W L_{epoch=0} = \nabla_W L(W_{epoch=0}) = [4W_0 \quad 2W_1] = [8\ 4]$

$v_{epoch=1} = \beta v_{epoch=0} + \nabla_W L_{epoch=0} \qquad v_{epoch=1} = [8\ 4]$

$W_{epoch=1} = W_{epoch=0} - \alpha v_{epoch=1} = [2 \quad 2] - 0.1 * [8 \quad 4] = [1.2 \quad 1.6]$
$L(W_{epoch=1}) = 2 * 1.2^2 + 1.6^2 = 5.44$

- Epoch 2:

$\nabla_W L(W_{epoch=1}) = [4W_0 \quad 2W_1] = [4.8 \quad 3.2]$

$v_{epoch=2} = \beta v_{epoch=1} + \nabla_W L_{epoch=1} \qquad v_{epoch=2} = 0.9 * [8\ 4] + [4.8 \quad 3.2] = [12 \quad 6.8]$

$W_{epoch=2} = W_{epoch=1} - \alpha\, v_{epoch=2} = [1.2 \quad 1.6] - 0.1 * [12 \quad 6.8] = [0 \quad 0.92]$
$L(W_{epoch=2}) = +0.92^2 = 0.8464$

# RMSProp

- RMSProp = Root Mean Square Propagation
- RMSProp algorithm has a 2-step update formula:

$$v_{epoch=i} = \beta v_{epoch=i-1} + (1-\beta)\nabla_W L^2_{epoch=i-1} \qquad v_{epoch=0} = 0$$

$$Update_{epoch=i} = -\alpha \frac{\nabla_W L_{epoch=i-1}}{\sqrt{v_{epoch=i}} + \epsilon}$$

- $v$ represents moving average of square of gradient for each entry of W
- Effective learning rate $\frac{\alpha}{\sqrt{v_{epoch=i}}+\epsilon}$ is different for different entries of W
- $\varepsilon$ included to avoid division by 0
- Typical parameter values: $\beta$=0.9 and $\varepsilon = 10^{-8}$

Reference: unpublished work(Hinton, Srivastava & Swersky)
http://www.cs.toronto.edu/~tijmen/csc321/slides/lecture_slides_lec6.pdf

# RmsProp - Example

- Consider simple function (has minimum at [0 0])

$L(W) = L(W_0, W_1) = 2W_0^2 + W_1^2$ with gradient $\nabla_W L = [4W_0 \quad 2W_1]$

- Choose $\alpha = 0.1,\ \beta = 0.9,\ \varepsilon = 0$ and initial guess

$W_{epoch=0} = [2\ 2] \qquad L(W_{epoch=0}) = 2 * 2^2 + 2^2 = 12$

Epoch 1:

$\nabla_W L_{epoch=0} = \nabla_W L(W_{epoch=0}) = [4W_0 \quad 2W_1] = [8\ 4] \quad \nabla_W L^2_{epoch=0} = [64 \quad 16]$

$v_{epoch=1} = \beta v_{epoch=0} + (1 - \beta)\ \nabla_W L^2_{epoch=0} = 0.1 * [64 \quad 16] = [6.4 \quad 1.6]$

$W_{epoch=1} = W_{epoch=0} - \dfrac{\alpha}{\sqrt{v_{epoch=1}} + \epsilon} \nabla_W L_{epoch=0} = [2 \quad 2] - 0.1\left[\dfrac{8}{\sqrt{6.4}} \quad \dfrac{4}{\sqrt{1.6}}\right]$

Component-wise square

$= [1.6838 \quad 1.6838]$

Component-wise
square root and division

$L(W_{epoch=1}) = 2 * 1.6838^2 + 1.6838^2 = 8.51$

# RmsProp - Example

Epoch 2:

$$\nabla_W L_{epoch=1} = \nabla_W L(W_{epoch=1}) = [4W_0 \quad 2W_1] = [6.7351 \quad 3.3675] \quad \nabla_W L_{epoch=1}^2 = [45.3614 \quad 11.3406]$$

$$v_{epoch=2} = \beta v_{epoch=1} + (1-\beta)\nabla_W L_{epoch=1}^2 = 0.9 * [6.4. \quad 1.6] + 0.1[45.3614 \quad 11.3406] = [10.2961 \quad 2.5740]$$

$$W_{epoch=2} = W_{epoch=1} - \frac{\alpha}{\sqrt{v_{epoch=1}} + \epsilon} \nabla_W L_{epoch=1} = [1.6838 \quad 1.6838] - 0.1 \left[ \frac{6.7351}{\sqrt{10.2961}} \quad \frac{3.3675}{\sqrt{2.5740}} \right]$$

$$= [1.4739 \quad 1.4739]$$

$$L(W_{epoch=2}) = 2 * 1.4739^2 + 1.4739^2 = 6.5169$$

# Adam

- Adam = Adaptive Moment Estimation
- Adam algorithm has a 3-step update formula:

$$m_{epoch=i} = \beta_1 m_{epoch=i-1} + (1 - \beta_1)\nabla_W L_{epoch=i-1} \qquad m_{epoch=0} = 0$$

$$v_{epoch=i} = \beta_2 v_{epoch=i-1} + (1 - \beta_2)\nabla_W L^2_{epoch=i-1} \qquad v_{epoch=0} = 0$$

$$Update_{epoch=i} = -\alpha \frac{m_{epoch=i}}{\sqrt{v_{epoch=i}} + \varepsilon}$$

- $v$ represents moving average of square of gradient
- Like RmsProp, learning rate is divided by this moving average
- Effective learning rate $\frac{\alpha}{\sqrt{v_{epoch=i}} + \epsilon}$ is different for different entries of W
- $m$ represents moving average of gradient
- Typical parameter values: $\beta_1$=0.9, $\beta_2$=0.999, and $\varepsilon = 10^{-8}$

Reference: (Kingma & Ba)
https://arxiv.org/abs/1412.6980v8

# Adam - Example

- Consider simple function (has minimum at [0 0])

$L(W) = L(W_0, W_1) = 2W_0^2 + W_1^2$ with gradient $\nabla_W L = [4W_0 \quad 2W_1]$

- Choose $\alpha$=0.1, $\beta_1 = 0.9$, $\beta_2 = 0.999$, $\varepsilon = 0$ and set $v_{epoch=0}$=0, $m_{epoch=0}$=0 and make initial guess:

$W_{epoch=0} = [2\ 2] \quad L(W_{epoch=0}) = 2 * 2^2 + 2^2 = 12$

Epoch 1:

$\nabla_W L_{epoch=0} = \nabla_W L(W_{epoch=0}) = [4W_0 \quad 2W_1] = [8\ 4] \quad \nabla_W L^2_{epoch=0} = [64 \quad 16]$

$m_{epoch=1} = \beta_1 m_{epoch=0} + (1 - \beta_1)\nabla_W L_{epoch=0} = 0.1* [8 \quad 4] = [0.8 \quad 0.4]$

$v_{epoch=1} = \beta_2 v_{epoch=0} + (1 - \beta_2)\nabla_W L^2_{epoch=0} = 0.001* [64 \quad 16] = [0.064 \quad 0.016]$

$W_{epoch=1} = W_{epoch=0} - \dfrac{\alpha}{\sqrt{v_{epoch=1}} + \epsilon} m_{epoch=1} = [2 \quad 2] - 0.1\left[\dfrac{0.8}{\sqrt{0.064}} \quad \dfrac{0.4}{\sqrt{0.016}}\right] = [1.6838 \quad 1.6838]$

Component-wise square

Component-wise
square root and division

$L(W_{epoch=1}) = 2 * 1.6838^2 + 1.6838^2 = 8.51$

# Adam - Example

Epoch 2:

$$\nabla_W L_{epoch=1} = \nabla_W L(W_{epoch=1}) = [4W_0 \quad 2W_1] = [6.7351 \quad 3.3675] \quad \nabla_W L^2_{epoch=1} = [45.3614 \quad 11.3406]$$

$$m_{epoch=2} = \beta_1 m_{epoch=1} + (1 - \beta_1)\nabla_W L_{epoch=1} = 0.9* [0.8 \quad 0.4] + 0.1* [6.7351 \quad 3.3675] = [1.3935 \quad 0.6968]$$

$$v_{epoch=2} = \beta v_{epoch=1} + (1 - \beta)\nabla_W L^2_{epoch=1} = 0.999 * [0.064 \quad 0.016] + 0.001 * [45.3614 \quad 11.3406] = [0.1093 \quad 0.0273]$$

$$W_{epoch=2} = W_{epoch=1} - \frac{\alpha}{\sqrt{v_{epoch=2}} + \epsilon} m_{epoch=2} = [1.6838 \quad 1.6838] - 0.1 \left[ \frac{1.3935}{\sqrt{0.1093}} \quad \frac{0.6968}{\sqrt{0.0273}} \right] = [1.2623 \quad 1.2623]$$

$$L(W_{epoch=2}) = 2 * 1.2623^2 + 1.2623^2 = 4.7799$$

# Optimization: Summary of Algorithms

| Name | Description |
|------|-------------|
| Optimization | Optimization algorithms have general form:<br>$W_{epoch=i} = W_{epoch=i-1} + Update_{epoch=i}$<br>Each algorithm will have a different formula for Update<br>None of these algorithms is guaranteed to converge to local or global minimum |
| Gradient Descent | $Update_{epoch=i} = -\alpha \nabla_W L_{epoch=i-1}$ |
| Momentum | $v_{epoch=i} = \beta v_{epoch=i-1} + \nabla_W L_{epoch=i-1}$      start with $v_{epoch=0} = 0$<br>$Update_{epoch=i} = -\alpha v_{epoch=i}$ |
| RmsProp | $v_{epoch=i} = \beta v_{epoch=i-1} + (1-\beta)\nabla_W L^2_{epoch=i-1}$    start with $v_{epoch=0} = 0$<br><br>$Update_{epoch=i} = -\alpha \dfrac{\nabla_W L_{epoch=i-1}}{\sqrt{v_{epoch=i}} + \epsilon}$ |
| Adam | $m_{epoch=i} = \beta_1 m_{epoch=i-1} + (1-\beta_1)\nabla_W L_{epoch=i-1}$    start with $m_{epoch=0} = 0$<br>$v_{epoch=i} = \beta_2 v_{epoch=i-1} + (1-\beta_2)\nabla_W L^2_{epoch=i-1}$    start with $v_{epoch=0} = 0$<br><br>$Update_{epoch=i} = -\alpha \dfrac{m_{epoch=i}}{\sqrt{v_{epoch=i}} + \varepsilon}$ |

# 6.2 Optimization – Jupyter Notebook DEMO

- Notebooks are:
  - IntroML/Examples/Chapter6/Momentum.ipynb
  - IntroML/Examples/Chapter6/RmsProp.ipynb
  - IntroML/Examples/Chapter6/Adam.ipynb

Course Resources at:

- [https://github.com/satishchandrareddy/IntroML/](https://github.com/satishchandrareddy/IntroML/)

# 6.3 Code Walkthrough Version 3.1

# Code Walkthrough Version 3.1

Goal of this Section:

- Walkthrough addition of mini-batch optimization
- Walkthrough addition of Momentum, RmsProp, and Adam optimizers

# Coding Walkthrough: Version 3.1 To Do

| File/Component | To Do |
|---|---|
| NeuralNetwork_Base | Add method to create mini-batches |
| NeuralNetwork_Base | Update fit method to perform mini-batch optimization |
| Optimizer | Add derived classes for Momentum, RmsProp, and Adam optimizers |
| driver_neuralnetwork_binary<br>driver_neuralnetwork_multiclass | Update drivers to use mini-batch optimization and additional optimizers and measure time for training |

# NeuralNetwork_Base Class – Methods

| Method | Input | Description |
|---|---|---|
| create_minibatches | X (numpy array)<br>Y (numpy array)<br>batch_size (integer) | Method to create mini-batches of specified size<br>Return: List of tuples (Xminibatch,Yminibatch) one for each mini-batch |
| fit | X (numpy array)<br>Y (numpy array)<br>epochs (integer)<br>**kwargs<br>batch_size (integer) | Update the existing fit method to use mini-batches in training if batch_size is input<br><br>Return: dictionary containing loss and accuracy histories at each epoch |

# Momentum, RmsProp, Adam - Methods

| Momentum: methods | Input | Description |
|---|---|---|
| __init__ | learning_rate (float)<br>beta (float) | Takes in parameters and initializes v=0<br>Return: nothing |
| update | gradient (numpy array) | Performs 2 step process to calculate update<br>Return: update |

| RmsProp: methods | Input | Description |
|---|---|---|
| __init__ | learning_rate (float)<br>beta (float)<br>epsilon (float) | Takes in parameters and initializes v=0<br>Return: nothing |
| update | gradient (numpy array) | Performs 2 step process to calculate update<br>Return: update |

| Adam: methods | Input | Description |
|---|---|---|
| __init__ | learning_rate (float)<br>beta1 (float)<br>beta2 (float)<br>epsilon (float) | Takes in parameters and initializes m=0 and v=0<br>Return: nothing |
| update | gradient (numpy array) | Performs 3 step process to calculate update<br>Return: update |

# Why Use Multiple Instances of Optimizer?

- Why use separate instance of optimizer object for each parameter matrix: $W^{[1]}, b^{[1]}, W^{[2]}, b^{[2]}, \ldots, W^{[N]}, b^{[N]}$ in code?

- Momentum, RmsProp, and Adam methods use information from previous epochs to compute $Update_{epoch=i}$ not just the most recently computed gradients

- This info from previous epochs depends on the parameter and can't be mixed with info for other parameters – can't mix histories for $W^{[1]}$ and $b^{[1]}$ or $W^{[1]}$ and $W^{[2]}$

# 6.3 Code Version 3.1 Walkthrough DEMO

- Code located at: IntroML/Code/Version3.1

- How to proceed:
  - Using information from previous videos, see if you can make updates to framework using an original Version2.2 as starting point
  - Alternatively, watch this video, which walks through updates and then use as a guide to make updates
  - The examples in the Jupyter notebooks indicate how we will use numpy functionality to convert algorithms into Python code

Course Resources at:

- https://github.com/satishchandrareddy/IntroML/

# 6.4 Validation

# Validation

Goal of this Section:

- Review of validation and testing approaches

# Effectiveness of Neural Network Models

- In code examples, we have used loss and accuracy of prediction for training data as a measure of effectiveness

- This provides some information but is not a reliable measure of model effectiveness
  - In trivial case, can create a model that is simply a look-up table based on training data (looks up output label based on input features in training data)
  - Accuracy is 100% for predicting results for input data in training set
  - This model is useless for prediction if new input is not in training data

- Effectiveness of model measured by how well predictions match actual results for data not used in training

# Train, Validation, and Test Data Sets

| | | | |
|---|---|---|---|
| Train | | Validation | Test |

| Dataset | Description |
|---|---|
| Train | Data used to fit the model<br>(Feature Matrix X and Label Vector Y in notation of this course) |
| Validation | Data not used in training. Loss/Accuracy calculations are performed for this dataset to see how well model performs at prediction for data not used in training. As we will see later in this chapter use this dataset to measure performance as machine learning system set up is varied: including number of layers, number of units, optimization parameters, $\alpha, \beta, \beta_1, \beta_2, \varepsilon$, etc. (This is called hyperparameter tuning – see Section 6.8.) |
| Test | The sample of data not used in training or in tuning. This dataset is used to measure performance after machine learning set-up is finalized. |
| | Source: Towards Data Science<br>https://towardsdatascience.com/train-validation-and-test-sets-72cb40cba9e7 |
| Dataset split | No definitive answer in literature. In this course we will focus on train and validation and use roughly 80-90%/20-10% split |

# Training Algorithm With Validation

Assume Neural Network with N layers

Input training data: feature matrix Xtrain and values Ytrain

Input validation data: feature matrix Xvalid and values Yvalid

Choose batch size b

Make initial guess for parameters: $W^{[k]}$ and $b^{[k]}$ for k=1,…,N

Choose optimizer and parameters: $\alpha, \beta, \beta_1, \beta_2, \varepsilon$, etc

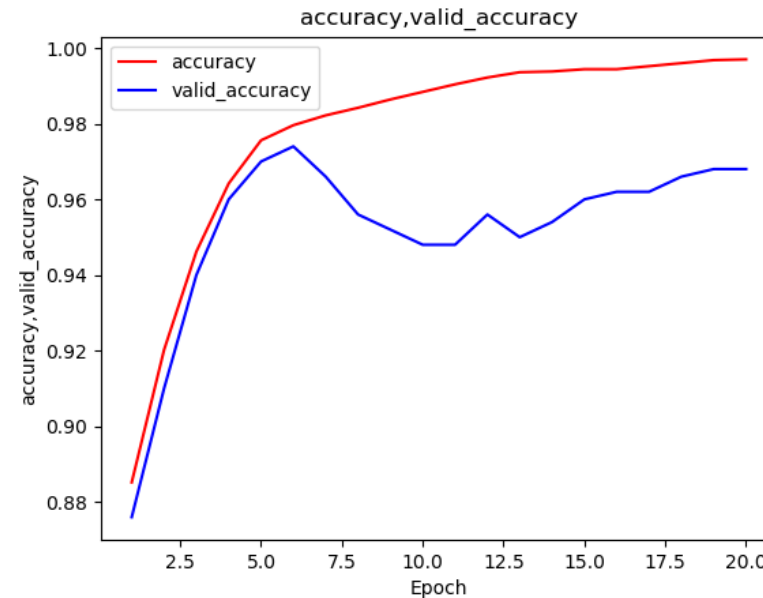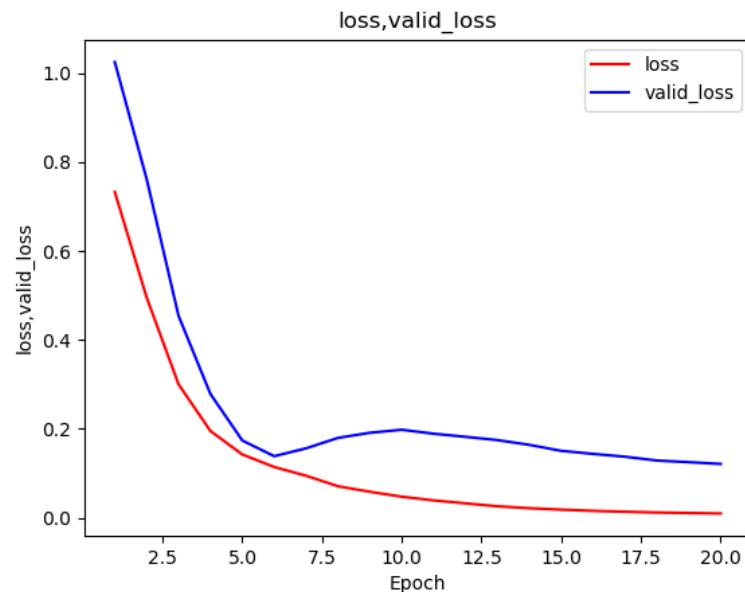1. Loop for epoch i = 1, 2, …
   - Generate mini-batches $(X_{minibatch=j}, Y_{minibatch=j})$, j=0,…,n-1
   - Loop over minibatches j=0,…,n-1
     - Forward Propagate using $Xtrain_{minibatch=j}$: compute $A^{[k]}_{minibatch=j}$ for k=1,…,N
     - Back Propagate using $Xtrain_{minibatch=j}, Ytrain_{minibatch=j}$: compute $\nabla_{W^{[k]}} L_{minibatch=j}$, $\nabla_{b^{[k]}} L_{minibatch=j}$ k=1,…,N
     - Update parameters: $W^{[k]}$ and $b^{[k]}$ for k=1,…,N
   - Forward Propagate using validation feature matrix Xvalid: compute $Avalid^{[k]}$ for k=1,…,N
   - Predict using validation data: compute $Yvalid_{pred}$ and compute validation accuracy comparing with $Yvalid$
   - Compute validation loss using $Avalid^{[N]}$
   - Forward Propagate using training feature matrix Xtrain: compute $Atrain^{[k]}$ for k=1,…,N
   - Predict using training data: compute $Ytrain_{pred}$ and compute training accuracy comparing with $Ytrain$
   - Compute training loss using $Atrain^{[N]}$

Loop for fixed number of iterations

# Loss and Accuracy Plots

- Loss and Accuracy plots will include histories for both training and validation

- Will discuss how to assess and improve performance using these plots later in chapter

# K-Fold Cross Validation

- K-Fold Validation
  - Randomize training and validation data
  - Split training and validation data into K "Folds" of roughly even size

| Fold 1 | Fold 2 | | | Fold K-1 | Fold K |

- Perform training with folds 1,2,…,K-1 and validation with fold K
- Repeat: train with folds 2,3,…,K and validation with fold 1
- Repeat with remaining combinations
- Validation accuracy is average over all possible combinations
- With this process each data point will be part of validation once and will be used in training K-1 times

# Stratified K-Fold Cross Validation

- Validation Data and Train Data should have same distribution
  - For binary classification, ensure proportion of 0 and 1 labels is same in train and validation data sets
  - For multiclass classification, ensure distribution of classes is same in train and validation data sets
- In K-Fold Validation, each fold should have similar distribution of outputs

# K-Fold Cross Validation Algorithm

Randomly shuffle training and validation data

Split data into K folds

Create K training/validation dataset combinations

1. Loop for combination i=1,...,K

    - Perform training using:
        - training dataset $Xtrain_{combination=i}, Ytrain_{combination=i}$
        - validation dataset $Xvalid_{combination=i}, Yvalid_{combination=i}$
    - Save validation accuracy at final epoch for combination=i

Final validation accuracy is average of validation accuracies for the K runs

# 6.5 Performance Measures

# Performance Measures

Goal of this Section:

- Review of accuracy measure

- Introduce precision, recall, and f1score measures

- Introduce Confusion Matrix

# Accuracy Calculation

- Accuracy by itself may not be best metric for measuring effectiveness of neural network

- Example:
  - Consider binary classification of x-ray images (normal versus abnormal)
  - Suppose that training data has few abnormal images (only 2%)
  - Without any neural network, we can create a trivial model: prediction of model is normal for all input images
  - This model will have 98% accuracy but is obviously a lousy model – it will never yield prediction of abnormal

- What are alternatives to accuracy metric for measuring effectiveness?

- References:

https://heartbeat.fritz.ai/evaluation-metrics-for-machine-learning-models-d42138496366

https://en.wikipedia.org/wiki/Confusion_matrix

# Precision, Recall, F1 Score for Binary Classification

- Accuracy is fraction of positives and negatives predicted correctly

$$Accuracy = \frac{\#TruePositive + \#TrueNegative}{\#Total} = \frac{\#Blue + \#Red}{\#Total}$$

- Precision is fraction of predicted positives that are correct

$$Precision = \frac{\#TruePositive}{\#PredictedPositive} = \frac{\#Blue}{\#Blue + \#Green}$$

- Recall is fraction of actual positives predicted correctly

$$Recall = \frac{\#TruePositive}{\#ActualPostive} = \frac{\#Blue}{\#Blue + \#Black}$$

- F1 score combines into a single figure

$$F1 = 2\frac{Precision * Recall}{Precision + Recall}$$

- For each of these measures, 1 is optimal score

|   | Actual | |
|---|---|---|
|  | Negative | Positive |
| **Predicted** Negative | True Negative (red) | False Negative (dark) |
| Positive | False Positive (green) | True Positive (blue) |

# Precision, Recall, F1 Score - Example

- Consider Binary Classification of X-Ray images
- 100 images total: actual and predictions in table
- Treat normal as negative and abnormal as positive

$$Accuracy = \frac{\#TruePositive + \#TrueNegative}{\#Total} = \frac{1+97}{100} = 98\%$$

$$Precision = \frac{\#TruePositive}{\#Predicted\ Positive} = \frac{1}{2} = 50\%$$

$$Recall = \frac{\#TruePositive}{\#ActualPositive} = \frac{1}{2} = 50\%$$

$$F1score = 2\frac{Precision * Recall}{Precision + Recall} = 50\%$$

|  | Actual Normal (Negative) | Actual Abnormal (Positive) |
|---|---|---|
| Predicted Normal (Negative) | 97 | 1 |
| Predicted Abnormal (Positive) | 1 | 1 |

- If relatively small number of actual positives, then accuracy alone may give a misleading picture of performance of machine learning system

# Target/Achievable Accuracy

- What accuracy should we expect from machine learning system?
  - For classification, perfect score is 100% accuracy
  - For regression, perfect score is 0 mean absolute error
  - Nearly impossible to achieve these perfect scores
- Machine learning models have prediction errors
  - Errors result from incomplete data – feature matrix may not include all relevant features
  - Noise in data
  - Stochastic nature of machine learning algorithm
    - Random choice of parameters W and b
  - Stop after fixed number of epochs so don't minimize loss function
- Target or achievable accuracy is problem dependent
  - For classification of cats and dogs based on photos, target should be close to 100%
  - In general, in case of noisy data, may not be able to achieve 100% accuracy
- In machine learning training and testing, goal should be to achieve target accuracy for training, validation, and test sets

# Target Accuracy - Example

- Figure on left:
  - Relatively sharp boundary between 0 and 1 training data points
  - Should expect target accuracy to be close to 100%

- Figure on right:
  - "Fuzzy" boundary between 0 and 1 training data points – noise has been added
  - Expect target/achievable to be less than 100% - not clear what label should be for new data points near boundary

# Confusion Matrix

- For classification (binary or multiclass), Confusion Matrix is a matrix/table listing counts of actual and predicted for each class

- Example: 3 classes
  - First Row indicates Predicted = 0 for 48 cases. Of these, 45 were actual class 0, 2 were actual class 1, and 1 was actual class 2

|  | Actual 0 | Actual 1 | Actual 2 |
|---|---|---|---|
| Predicted 0 | 45 | 2 | 1 |
| Predicted 1 | 3 | 51 | 2 |
| Predicted 2 | 1 | 1 | 44 |

- Useful for determining patterns in mis-classification

# Algorithm for Producing Confusion Matrix

Assume training has been performed

Assume prediction Y_pred  is made for which actual labels are known

Input actual labels Y, predicted labels Y_pred and number of classes c

1. For predicted class p= 0,1,…,c-1
   - Determine indices (idx_p) where Y_pred = p
   - For actual class a = 0,1,…,c-1
     - Count = number of entries where actual Y[idx_p] = a
     - Confusion Matrix (Row = p, Col = a) = Count

|  | Actual 0 | Actual  1 | Actual 2 |
|---|---|---|---|
| Predicted   0 | 45 | 2 | 1 |
| Predicted   1 | 3 | 51 | 2 |
| Predicted   2 | 1 | 1 | 44 |

# 6.6 Code Walkthrough Version 3.2

# Coding Walkthrough: Version 3.2

Goal of this Section:

- Add calculation of loss and accuracy for validation data set

- Add Confusion Matrix, Precision, Recall and F1 Score calculations

- Update drivers to input validation data for training

- Create K-fold cross validation example

# Coding Walkthrough: Version 3.2 To Do

| File/Component | To Do |
|---|---|
| NeuralNetwork_Base | Add calculation of loss/accuracy for validation data set for each epoch in fit method and add input to turn off printing of epoch level results |
| metrics | Add function to calculate precision, recall and F1 score |
| metrics | Add function to print confusion matrix |
| example_classification | Add functionality to generate a validation dataset and add noise to data |
| plot_results | Add functionality to plot validation data as well as training data in scatter plot |
| driver_neuralnetwork_ binary driver_neuralnetwork_ multiclass | Update to use validation dataset |
| driver_neuralnetwork_bi nary_kfold | Create example of K-fold cross validation |

# NeuralNetwork_Base Class – Methods

| Method/Function | Input | Description |
| --- | --- | --- |
| fit | X (numpy array)<br>Y (numpy array)<br>epochs (integer)<br>**kwargs<br>batch_size (integer)<br>verbose (boolean)<br>validation_data (tuple containing feature matrix and values for validation dataset) | Update the existing train method to compute accuracy/loss for validation dataset if provided, and turn off printing of loss and accuracy if verbose is False<br><br>Return: dictionary containing loss and accuracy histories for train and validation dataset (if provided) |

# Metrics

| Function | Input | Description |
|---|---|---|
| f1score | Y (numpy array)<br>Y_pred (numpy array) | Return: f1score, precision, and recall |
| confusion_matrix | Y (numpy array)<br>Y_pred (numpy array)<br>nclass (integer) | Prints confusion matrix<br><br>Return: nothing |

# example_classification

| Function | Input | Description |
|----------|-------|-------------|
| example | nfeature (integer)<br>m (integer)<br>case (string)<br>nclass (integer)<br>noise (boolean)<br>validpercent (float) | Add functionality to add noise to data (if noise = True).<br>Add functionality to create validation data – number of validation data points is m * validpercent<br><br>Return: Xtrain,Ytrain,Xvalid,Yvalid |

# plot_results

| Function | Input | Description |
|---|---|---|
| plot_results_classification | data_train (tuple (X,Y))<br>model (object)<br>nclass<br>**kwargs<br>validation_data (tuple (Xvalid,Yvalid)) | Umbrella function to plot data and heatmap for classification results. Update to include keyword argument for validation data |
| plot_results_data | data_train (tuple (X,Y))<br>nclass<br>**kwargs<br>validation_data (tuple (Xvalid,Yvalid)) | Update to plot validation data if keyword "validation_data" input is used.<br><br>Return: nothing |

# 6.6 Code Version 3.2 Walkthrough DEMO

- Code located at: IntroML/Code/Version3.2

- How to proceed:
  - Using information from previous videos, see if you can make updates to framework using an original Version3.1 as starting point
  - Alternatively, watch this video, which walks through updates and then use as a guide to make updates
  - The examples in the Jupyter notebooks indicate how we will use numpy functionality to convert algorithms into Python code

Course Resources at:

- https://github.com/satishchandrareddy/IntroML/

# 6.7 Regularization

# Regularization

Goal of this Section:

- Discuss purpose and approaches for regularization

# Purpose of Regularization and Approaches

Purpose of Regularization:

- Technique for solving ill-posed problems/avoid overfitting
- Tool to improve machine learning performance

Approaches for Regularization:

- Add "information" to problem
  - L1 regularization
  - L2 regularization
- Stop training algorithm early
- Dropout (remove some nodes in layers randomly)
- Good article for more information

https://www.analyticsvidhya.com/blog/2018/04/fundamentals-deep-learning-regularization-techniques/

# Underfitting and Overfitting - Definitions

- Underfitting
  - Occurs when function structure cannot capture the trend of the data
  - Example: Linear function attempting to fit "quadratic" data

- Overfitting
  - Occurs when training yields function that captures training data extremely well but does not generalize to data not seen in training – fits training data well but does not capture "trend" of data
  - Example: 11$^{th}$ degree polynomial used to fit 12 data points – fits training data exactly, but large discrepancy at end point does not capture trend of data

- Good/Robust Fit:
  - Occurs when training yields function that captures training data well and also generalizes to data not seen in training
  - Example: Model prediction captures training data well and captures trend of data

# L1 and L2 Regularization

- L1,L2 Regularization involve adding "penalty" term to loss function
$\hat{L} = regular\ loss + Penalty = L + P$

- L1 Regularization:

$$Penalty = \sum_{k=1}^{N} \lambda^{[k]} \sum_{i=0}^{n^{[k]}-1} \sum_{j=0}^{n^{[k-1]}-1} \left| W_{ij}^{[k]} \right| \qquad \text{(Lasso)}$$

- L2 Regularization:

$$Penalty = \sum_{k=1}^{N} \lambda^{[k]} \sum_{i=0}^{n^{[k]}-1} \sum_{j=0}^{n^{[k-1]}-1} \left( W_{ij}^{[k]} \right)^2 \qquad \text{(Ridge)}$$

- Constants $\lambda^{[k]} \geq 0$ for k=1,...,N are user specified and may be different for each layer of neural network

- Penalty terms are always non-negative

- Typically don't include parameters $b^{[k]}$ k=1,...,N in penalty term

# Impact of Regularization

- Penalty term forces training algorithm to smaller parameter values $W_{ij}^{[k]}$

- Why is this important?

- Recall forward propagation:

  Loop for k=1,…,N (number of layers)
  - Linear part: $Z^{[k]} = W^{[k]}A^{[k-1]} + b^{[k]} \quad A^{[0]} = X$
  - Activation: $A^{[k]} = f^{[k]}(Z^{[k]})$

- Large entries in $W^{[k]}$ can lead to large changes in $Z^{[k]}, A^{[k]}$ for small changes in input X leading to overfitting

- Section 6.10 has much more detail regarding use of regularization to address overfitting. In this section, I want to present the underlying math

# L2 Regularization - Example

- Consider binary classification with 2 features and 3 training data points (m=3)

$$X = \begin{bmatrix} 1 & 2 & 4 \\ -2 & -5 & -8 \end{bmatrix} \qquad Y = \begin{bmatrix} 0 & 1 & 0 \end{bmatrix}$$

- Assume 2 layer neural network: layer 1 (tanh) has 2 units, layer 2 (sigmoid) has 1unit

- Parameter matrices:

$$W^{[1]} = \begin{bmatrix} W_{00}^{[1]} & W_{01}^{[1]} \\ W_{10}^{[1]} & W_{11}^{[1]} \end{bmatrix} = \begin{bmatrix} 0.5 & 0.5 \\ 0.5 & -0.5 \end{bmatrix} \qquad b^{[1]} = \begin{bmatrix} 0.5 \\ 0.5 \end{bmatrix}$$

$$W^{[2]} = \begin{bmatrix} W_{00}^{[2]} & W_{01}^{[2]} \end{bmatrix} = \begin{bmatrix} -1 & 1 \end{bmatrix} \qquad b^{[2]} = \begin{bmatrix} -0.1 \end{bmatrix}$$

- Forward Propagation:

$$A^{[1]} = \begin{bmatrix} 0 & -0.7616 & -0.9051 \\ 0.9640 & 0.9993 & 1.0 \end{bmatrix}, \qquad A^{[2]} = \begin{bmatrix} 0.7035 & 0.8404 & 0.8588 \end{bmatrix}$$

- Binary cross entropy loss function with penalty is:

$$\hat{L} = -\frac{1}{3}\sum_{j=0}^{2} Y_j \ln\left(A_j^{[2]}\right) + (1 - Y_j)\ln\left(1 - A_j^{[2]}\right) + \lambda^{[1]}\left(\left(W_{00}^{[1]}\right)^2 + \left(W_{01}^{[1]}\right)^2 + \left(W_{10}^{[1]}\right)^2 + \left(W_{11}^{[1]}\right)^2\right) + \lambda^{[2]}\left(\left(W_{00}^{[2]}\right)^2 + \left(W_{01}^{[2]}\right)^2\right)$$

# Impact of L2 Regularization on Back Propagation

- L2 Regularization does not directly change Forward Propagation

- L2 Regularization does change Loss function and changes Back Propagation

$$\hat{L} = regular\ loss + Penalty = L + P$$

- Recall penalty term is:

$$Penalty = P = \sum_{k=1}^{N} \lambda^{[k]} \sum_{i=0}^{n^{[k]}-1} \sum_{j=0}^{n^{[k-1]}-1} \left(W_{ij}^{[k]}\right)^2$$

- Computing derivatives - each $W_{ij}^{[k]}$ appears quadratically, so

$$\frac{\partial P}{\partial W_{ij}^{[k]}} = 2\lambda^{[k]} W_{ij}^{[k]} \quad \text{which implies} \quad \nabla_{W^{[k]}} P = 2\lambda^{[k]} W^{[k]}$$

- In the back propagation algorithm, the only change to gradient formulas is:

$$\nabla_{W^{[k]}} \hat{L} = \nabla_{Z^{[k]}} L A^{[k-1]^T} + \boxed{2\lambda^{[k]} W^{[k]}}$$

# 6.8 Hyperparameter Search

# Hyperparameter Search

Goal of this Section:

• Describe hyperparameter searches/tuning

# What are Hyperparameters?

Parameters:

- Parameter matrices $W^{[k]}$ and $b^{[k]}$ for k=1,...,N
- Also refer to entries of $W^{[k]}$ and $b^{[k]}$ as parameters
- Learned during training process

Hyperparameters (not learned)

- Settings/quantities that define neural network structure, regularization, or are used to control learning/training process

| Component | Hyperparameters |
|---|---|
| Neural Network | Number of hidden layers<br>Units per hidden layer<br>Activation function (not strictly a hyperparameter) |
| fit method | Number of epochs<br>Batch size |
| Optimization | Gradient Descent: learning rate $\alpha$<br>Momentum: learning rate $\alpha$, $\beta$<br>RmsProp: learning rate $\alpha$, $\beta$, $\varepsilon$<br>Adam learning rate $\alpha$, $\beta_1$, $\beta_2$, $\varepsilon$ |
| L2 Regularization | Regularization multipliers $\lambda^{[k]}$ for layers k=1,...,N |

# Hyperparameter Search

Hyperparameter search process

- Define training and validation datasets

- Create Neural Network Model

- Define optimization method

- Search for "Best" Hyperparameter combination
    - Perform training for each hyperparameter combination
    - Measure accuracy for validation dataset

- Choose Hyperparameter combination with highest validation accuracy

# Hyperparameter Search - Example

- Example:  Dataset
  - 1000 samples in training dataset
  - 200 samples in validation dataset
  - Binary classification: 2 features "quadratic" example with noise
  - Feature matrix is 2x1000 for training and 2x200 for validation

- Neural Network
  - 5 layer neural network (15, 11, 8, 4, 1 units in layers)
  - tanh activation in layers 1,2,3,4 and sigmoid in final layer
  - L2 regularization with variable $\lambda^{[k]}$

- Optimization:
  - Binary cross entropy loss function
  - Momentum with fixed $\beta$=0.9  and variable learning rate $\alpha$
  - 300 epochs
  - Mini-batch with variable batch size

- Search parameters
  - L2 regularization parameter (same value for all layers) $\lambda^{[k]}$ = 0, 0.01, 0.02
  - Learning rate $\alpha$ = 0.01, 0.03, 0.1, 0.3
  - Mini-batch size = 64, 1000

# Hyperparameter Search – Example 1

| | A | B | C | D | E |
|---|---|---|---|---|---|
| 1 | Lambda | Learning Rate | batch_size | Train Accuracy | Valid Accuracy |
| 2 | 0 | 0.01 | 64 | 0.95 | 0.925 |
| 3 | 0 | 0.01 | 1000 | 0.945 | 0.95 |
| 4 | 0 | 0.03 | 64 | 0.953 | 0.935 |
| 5 | 0 | 0.03 | 1000 | 0.95 | 0.94 |
| 6 | 0 | 0.1 | 64 | 0.936 | 0.94 |
| 7 | 0 | 0.1 | 1000 | 0.948 | 0.935 |
| 8 | 0 | 0.3 | 64 | 0.592 | 0.57 |
| 9 | 0 | 0.3 | 1000 | 0.957 | 0.915 |
| 10 | 0.01 | 0.01 | 64 | 0.947 | 0.94 |
| 11 | 0.01 | 0.01 | 1000 | 0.951 | 0.94 |
| 12 | 0.01 | 0.03 | 64 | 0.947 | 0.94 |
| 13 | 0.01 | 0.03 | 1000 | 0.948 | 0.93 |
| 14 | 0.01 | 0.1 | 64 | 0.926 | 0.925 |
| 15 | 0.01 | 0.1 | 1000 | 0.949 | 0.935 |
| 16 | 0.01 | 0.3 | 64 | 0.911 | 0.9 |
| 17 | 0.01 | 0.3 | 1000 | 0.948 | 0.95 |
| 18 | 0.02 | 0.01 | 64 | 0.948 | 0.95 |
| 19 | 0.02 | 0.01 | 1000 | 0.948 | 0.93 |
| 20 | 0.02 | 0.03 | 64 | 0.935 | 0.935 |
| 21 | 0.02 | 0.03 | 1000 | 0.945 | 0.945 |
| 22 | 0.02 | 0.1 | 64 | 0.929 | 0.92 |
| 23 | 0.02 | 0.1 | 1000 | 0.946 | 0.935 |
| 24 | 0.02 | 0.3 | 64 | 0.755 | 0.785 |
| 25 | 0.02 | 0.3 | 1000 | 0.946 | 0.935 |

- 3 Cases have 95% Validation Accuracy
- Pick between Row 17 and 18 (as these have higher Training accuracy)

# Hyperparameter Search – Example 2

| | A | B | C | D | E |
|---|---|---|---|---|---|
| 1 | Lambda | Learning Rate | batch_size | Train Accuracy | Valid Accuracy |
| 2 | 0 | 0.01 | 64 | 0.95 | 0.929166667 |
| 3 | 0 | 0.01 | 1000 | 0.929 | 0.916666667 |
| 4 | 0 | 0.03 | 64 | 0.9465 | 0.923333333 |
| 5 | 0 | 0.03 | 1000 | 0.944 | 0.931666667 |
| 6 | 0 | 0.1 | 64 | 0.947333333 | 0.926666667 |
| 7 | 0 | 0.1 | 1000 | 0.950666667 | 0.9375 |
| 8 | 0 | 0.3 | 64 | 0.856833333 | 0.8425 |
| 9 | 0 | 0.3 | 1000 | 0.9545 | 0.929166667 |
| 10 | 0.01 | 0.01 | 64 | 0.945 | 0.943333333 |
| 11 | 0.01 | 0.01 | 1000 | 0.9505 | 0.944166667 |
| 12 | 0.01 | 0.03 | 64 | 0.939833333 | 0.940833333 |
| 13 | 0.01 | 0.03 | 1000 | 0.945833333 | 0.943333333 |
| 14 | 0.01 | 0.1 | 64 | 0.914666667 | 0.915833333 |
| 15 | 0.01 | 0.1 | 1000 | 0.946833333 | 0.945 |
| 16 | 0.01 | 0.3 | 64 | 0.901 | 0.895833333 |
| 17 | 0.01 | 0.3 | 1000 | 0.943833333 | 0.9425 |
| 18 | 0.02 | 0.01 | 64 | 0.9445 | 0.946666667 |
| 19 | 0.02 | 0.01 | 1000 | 0.944 | 0.9425 |
| 20 | 0.02 | 0.03 | 64 | 0.938 | 0.936666667 |
| 21 | 0.02 | 0.03 | 1000 | 0.9445 | 0.9425 |
| 22 | 0.02 | 0.1 | 64 | 0.9035 | 0.9 |
| 23 | 0.02 | 0.1 | 1000 | 0.945166667 | 0.945 |
| 24 | 0.02 | 0.3 | 64 | 0.754166667 | 0.740833333 |
| 25 | 0.02 | 0.3 | 1000 | 0.946 | 0.946666667 |

- Use k-fold cross validation
- In this example 6 folds (of size 200 samples)
- Train and Validation Accuracy are averages over the 6 folds
- 2 cases with validation accuracy of 94.67%
- Pick row 25 (has higher training accuracy)

# 6.9 Code Walkthrough Version 3.3

# Coding Walkthrough: Version 3.3

Goal of this Section:

- Walkthrough code changes for regularization
- Walkthrough code updates for hyperparameter search

# Coding Walkthrough: Version 3.3 To Do

| File/Component | To Do |
|---|---|
| NeuralNetwork_Base | Update compute_loss method to include L2 regularization penalty term |
| LRegression | Update __init__ method to allow input of lambda for regularization<br>Update back_propagation method to incorporate L2 regularization change |
| NeuralNetwork | Update add_layer method to allow input of lambda for regularization<br>Update back_propagation method to incorporate L2 regularization change |
| driver_linearregression<br>driver_logisticregression<br>driver_neuralnetwork_<br>binary<br>driver_neuralnetwork_mult<br>iclass<br>unittest_forwardback<br>propagation | Update these drivers and unittest file to include input of lambda value for L2 regularization |
| driver_neuralnetwork_<br>binary_search<br>driver_neuralnetwork_<br>binary_search_kfold | Create new drivers for performing hyperparameter search/tuning |
| write_csv | Function for writing results to csv |

# LRegression, NeuralNetwork, NeuralNetwork_Base – Methods

| Method | Input | Description |
|---|---|---|
| LRegression __init__ | nfeature (integer) activation (string) lamb (float) | Update: take input lamb for regularization constant. Save in info[0]["lambda"] variable Return: nothing |
| LRegression back_propagate | X (numpy array) Y (numpy array) | Update: add $2\lambda W$ in calculation of $\nabla_W L$ Return: nothing |
| NeuralNetwork add_layer | nunit (integer) activation (string) lamb (float) | Update: take input lamb for regularization constant. Save in info[layer]["lambda"] variable Return: nothing |
| NeuralNetwork back_propagate | X (numpy array) Y (numpy array) | Update: add $2\lambda^{[k]} W^{[k]}$ for calculation of $\nabla_{W^{[k]}} L$ Return: nothing |
| NeuralNetwork_Base compute_loss | Y (numpy array) | Update to include L2 regularization penalty term Return: loss (including penalty term) |

# write_csv

| File/Component | Input | Description |
|---|---|---|
| write_csv | filename (string)<br>list_out (list of lists) | This program uses the csv.writer function in the csv package to write to file. Variable list_out is a list of lists – list_out=[list1,list2, ….]. Each of list1,list2,… is written to separate row in filename.<br>Return: nothing |

# 6.9 Code Version 3.3 Walkthrough DEMO

- Code located at: IntroML/Code/Version3.3

- How to proceed:
  - Using information from previous videos, see if you can make updates to framework using an original Version3.2 as starting point
  - Alternatively, watch this video, which walks through updates and then use as a guide to make updates
  - The examples in the Jupyter notebooks indicate how we will use numpy functionality to convert algorithms into Python code

Course Resources at:

- https://github.com/satishchandrareddy/IntroML/

# 6.10 Underfitting and Overfitting

# Underfitting and Overfitting

Goal of this Section:

- Define and provide examples of Underfitting and Overfitting

- Show how to diagnose Underfitting and Overfitting

- Present remedies for addressing Underfitting and Overfitting

# Underfitting and Overfitting - Definitions

- Underfitting
  - Occurs when function structure cannot capture the trend of the data
  - Example: Linear function attempting to fit "quadratic" data

- Overfitting
  - Occurs when training yields function that captures training data extremely well but does not generalize to data not seen in training – fits training data well but does not capture "trend" of data
  - Example: 11$^{th}$ degree polynomial used to fit 12 data points – fits training data exactly, but large discrepancy at end point does not capture trend of data

- Good/Robust Fit:
  - Occurs when training yields function that captures training data well and also generalizes to data not seen in training
  - Example: Model prediction captures training data well and captures trend of data

# Bias Error and Variance

In Machine Learning literature terms Bias Error and Variance are used when discussing underfitting and overfitting

Bias Error

- Occurs from erroneous assumptions in learning algorithm. High bias (underfitting) can cause algorithm to miss relationship between input features and target outputs.

Variance

- Sensitivity to small fluctuations in training data. High variance (overfitting) can cause algorithm to capture random noise in training data as opposed to general trend.

Source

- https://en.wikipedia.org/wiki/Bias%E2%80%93variance_tradeoff

Bias/Variance Tradeoff

- Models with low bias often have high variance and vice versa
- Goal: find model with low bias and low variance

# Underfitting - Example

- Example: Binary Classification for case of "ring"
  - 2000 training data points chosen from uniform distribution in [-2,2]x[-2,2]

- Machine Learning approach:
  - Logistic Regression
  - Adam $\alpha$=0.005, $\beta_1$=0.9, $\beta_2 = 0.999$, $\varepsilon = 10^{-8}$
  - 100 epochs and batch size of 64

- Results
  - Heatmap shows logistic regression cannot capture "nonlinear" data in ring data set – prediction of model is 1 for all points

- Symptoms of Underfitting
  - Training Loss well above 0 – roughly 0.62 after 100 epochs
  - Training Accuracy well below target accuracy – stays at roughly 0.69 after 100 epochs

# Underfitting – Remedy

- Remedy for underfitting: use "better" model: Here we try models with more layers
  - Model1: 1 layer (1 unit)
  - Model2: 2 layers (4,1 units)
  - Model3: 3 layers (8,4,1 units)
  - Model4: 4 layers (12,8,4,1 units)
  - Model5: 5 layers (16,12,8,4,1 units)
  - Use tanh activation in all hidden layers and sigmoid in final layer
  - Train for 100 epochs using batch size of 64 and Adam optimizer $\alpha$=0.005, $\beta_1$=0.9, $\beta_2 = 0.999$, $\varepsilon = 10^{-8}$
- Model5 achieves loss <0.1 after 100 epochs
- Model5 achieves accuracy of 97.9% after 100 epochs
- Heatmap shows results for Model5

# Overfitting - Example

- Example: Binary Classification – Quadratic with Noise
  - 1000 training data points chosen from uniform distribution in [-2,2]x[-2,2]
  - 200 validation data points chosen from uniform distribution in [-2,2]x[-2,2]
- Machine Learning approach:
  - Neural Network with 5 layers (15,11,8,4,1 units)
  - tanh activation for all layers except sigmoid in last
  - Momentum Optimization with $\alpha$=0.3 and $\beta$=0.9
  - Train with 300 epochs and batchsize = 1000 (number training data points)
- Symptoms:
  - Loss for validation dataset greater than that for training dataset (plot below shows that loss increases with epoch for validation dataset)
  - Accuracy for validation dataset less than that for training dataset  (Training Accuracy = 0.957, Validation Accuracy = 0.915)
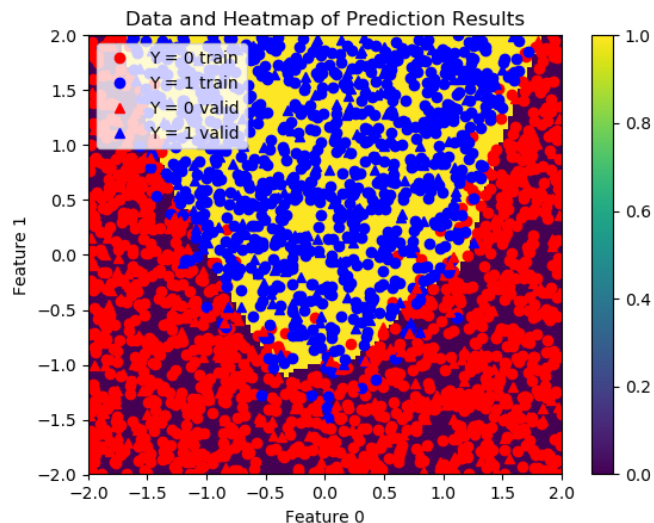
# Overfitting - Remedies

Overfitting Remedies

- Use more data points
  - This may not always be possible in real world setting

- Use simpler function structure
  - With more layers, neural network can achieve more complex function structure. Training tries to fit more complex function structure more precisely to training data. This may lead to missing the general trend of data
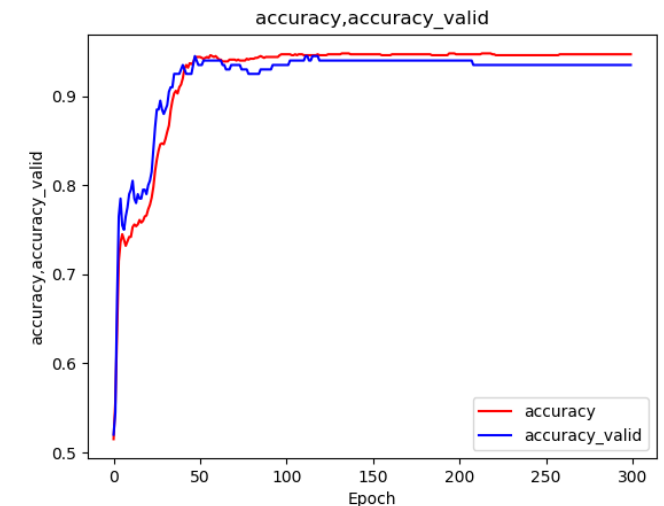  - Use fewer layers in neural network or fewer units in layers
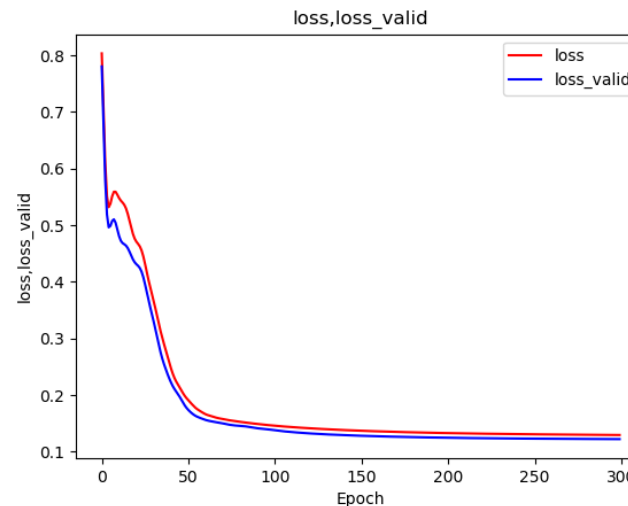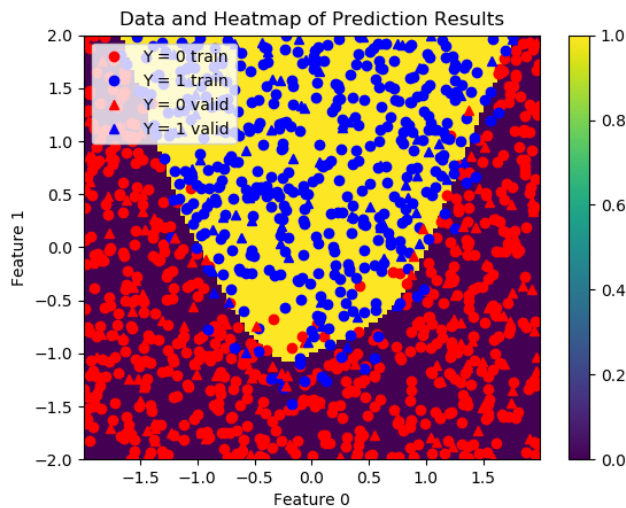
- Use Regularization
  - L2 regularization

# Overfitting Remedy – Use More Data

- Example: Binary Classification – Quadratic with Noise – Double Number of Data Points
  - 2000 training data points chosen from uniform distribution in [-2,2]x[-2,2]
  - 400 validation data points chosen from uniform distribution in [-2,2]x[-2,2]

- Machine Learning approach:
  - Neural Network with 5 layers (15,11,8,4,1 units)
  - tanh activation for all layers except sigmoid in last
  - Momentum Optimization with $\alpha$=0.3 and $\beta$=0.9
  - Train with 300 epochs and batch_size = 1000

- Results
  - Heatmap shows smoother boundary between 0 and 1 regions
  - Loss for validation dataset slightly greater than that for training dataset
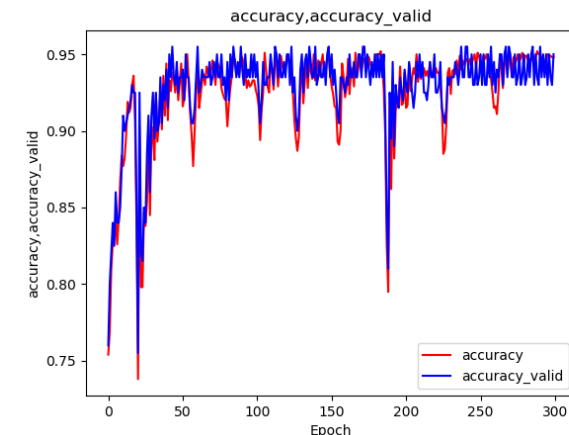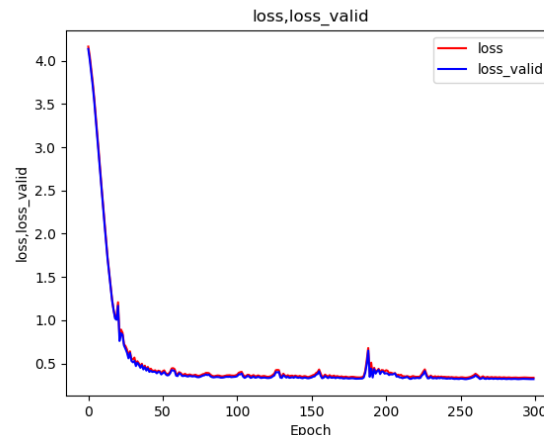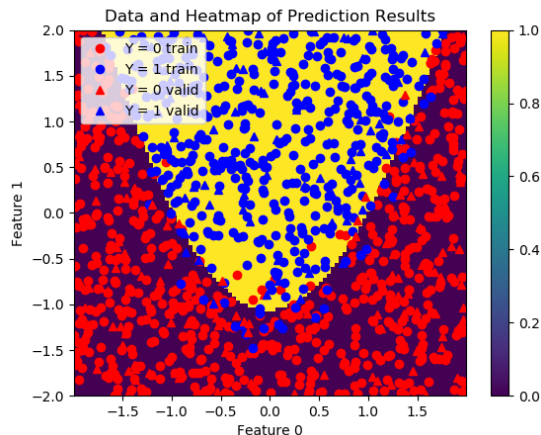  - At final epoch: Training Accuracy 0.9485, Validation Accuracy 0.965

# Overfitting Remedy – Use Simpler Neural Network

- Example: Binary Classification – Quadratic with Noise
  - 1000 training data points chosen from uniform distribution in [-2,2]x[-2,2]
  - 200 validation data points chosen from uniform distribution in [-2,2]x[-2,2]

- Machine Learning approach:
  - Neural Network with 2 layers (4,1 units)
  - tanh activation for all layers except sigmoid in last
  - Momentum Optimization with $\alpha$=0.3 and $\beta$=0.9
  - Train with 300 epochs and batch_size = 1000

- Results
  - Heatmap shows smoother boundary between 0 and 1 regions
  - Loss for validation dataset slightly greater than that for training dataset
  - At final epoch: Training Accuracy 0.947, Validation Accuracy 0.935

# Overfitting Remedy – Use Regularization

- Example: Binary Classification – Quadratic with Noise
  - 1000 training data points chosen from uniform distribution in [-2,2]x[-2,2]
  - 200 validation data points chosen from uniform distribution in [-2,2]x[-2,2]

- Machine Learning approach:
  - Neural Network with 5 layers (15,11,8,4,1 units)
  - tanh activation for all layers except sigmoid in last
  - Momentum Optimization with $\alpha$=0.3 and $\beta$=0.9
  - Train with 300 epochs and batch_size = 1000 (number training  data points)
  - Use L2 regularization with $\lambda^{[k]}$=0.01 for all layers k=1,2,3,4,5

- Results
  - Heatmap shows smoother boundary between 0 and 1 regions
  - Loss for validation dataset slightly greater than that for training dataset
  - At final epoch: Training Accuracy 0.948, Validation Accuracy 0.95 (validation accuracy oscillates between 0.93 and 0.95)
  - Maximum absolute value of entries of W for layers 1,2,3,4,5 ($\lambda^{[k]}$=0.0): 1.93, 2.94, 3.60, 2.98, 3.28
  - Maximum absolute value of entries of W for layers 1,2,3,4,5 ($\lambda^{[k]}$=0.01): 0.79, 0.41, 0.39, 0.44, 1.35

# Machine Learning Approach – Flowchart



Machine Learning System → Train Accuracy achieve target? → **Yes** → Validation Accuracy /Loss similar to Training Accuracy/Loss → **Yes** → Use Test Dataset to compute final accuracy score

Train Accuracy achieve target? → **No** → Underfit Use "better" NN

Validation Accuracy /Loss similar to Training Accuracy/Loss → **No** → Overfit Use simpler NN or Use more data or Use regularization or Use other technique