

Machine Learning

Lecture 06: Deep Feedforward Networks

Nevin L. Zhang
lzhang@cse.ust.hk

Department of Computer Science and Engineering
The Hong Kong University of Science and Technology

This set of notes is based on various sources on the internet and
Goodfellow, I., Bengio, Y., & Courville, A. (2016). *Deep Learning*. MIT press.
www.deeplearningbook.org

Introduction

- So far, probabilistic models for supervised learning

$$\{\mathbf{x}_i, y_i\}_{i=1}^N \rightarrow P(y|\mathbf{x})$$

- Next, **deep learning**:

$$\{\mathbf{x}_i, y_i\}_{i=1}^N \rightarrow \mathbf{h} = f(\mathbf{x}), P(y|\mathbf{h})$$

- $\mathbf{h} = f(\mathbf{x})$ is a feature transformation represented by a neural network,
- $P(y|\mathbf{h})$ is a probabilistic model on the transformed features.
- Regarded as **one whole model**: $P(y|\mathbf{x})$.
- This lecture: $\mathbf{h} = f(\mathbf{x})$ as a **feedforward neural network (FNN)**.
- Next lecture: $\mathbf{h} = f(\mathbf{x})$ as a **convolutional neural network (CNN or ConvNet)**.

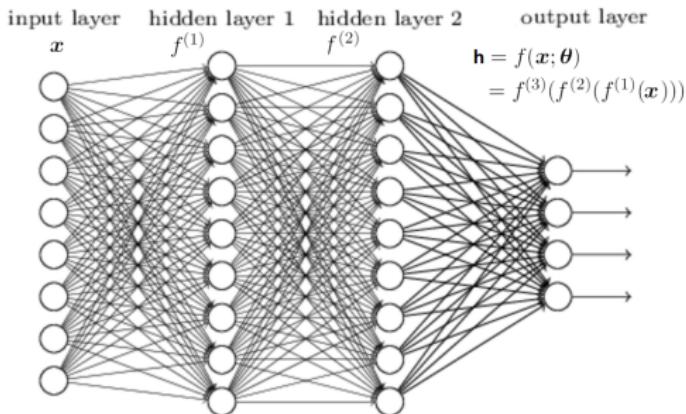
Outline

- 1 Feedforward Neural Network as Function Approximator
- 2 Feedforward Neural Network as Probabilistic Model
- 3 Backpropagation
- 4 Dropout
- 5 Optimization Algorithms

Deep Feedforward Networks

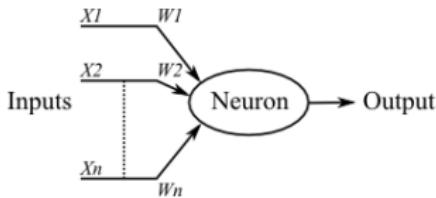
- Deep feedforward networks, also often called feedforward neural networks (FNNs), or multilayer perceptrons (MLPs), are the quintessential deep learning models
- A feedforward network defines a function $\mathbf{h} = f(\mathbf{x}, \theta)$.
- During learning, the parameters θ are optimized so that $f(\mathbf{x}, \theta)$ approximates some target function $f^*(\mathbf{x})$

Feedforward Neural Networks



- Networks of simple computing elements (**units, neurons**).
 - Each unit is connected to units on the previous layer and units on the next layer
 - Parameters include a **bias** for each unit and a **weight** for each link.
 - Units are divided into: **input units**, **hidden units**, and **output units**
- **Feedforward networks:**
 - Inputs enter the input units,
 - Propagate through the network to the output units.

The Units



- A unit accepts vector of inputs \mathbf{x} ,
- computes an affine transformation $z = \mathbf{W}^\top \mathbf{x} + b$, where $\mathbf{W} = (W_1, W_2, \dots, W_n)^\top$ are the **link weights** and b is the **bias** of the unit. z is sometimes called the **net input** of the unit.
- applies nonlinear **activation function** $g(z)$, and
- gives output $g(z) = g(\mathbf{W}^\top \mathbf{x} + b)$

Different types of units have different activation functions.

Initialize all weights to small random values, and biases to zero or to small positive values.

Weight Initialization

- The *Xavier* initialization method:

$$w \sim U\left[-\frac{1}{\sqrt{n}}, \frac{1}{\sqrt{n}}\right]$$

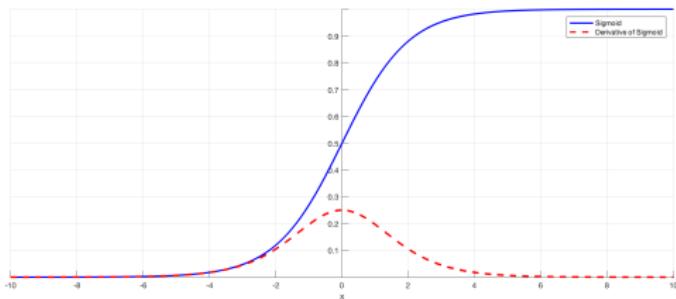
where n is the number of inputs to the neuron. Usually used for tanh.

- The *He* initialization or *Kaimin* initialization method:

$$w \sim \mathcal{N}\left[0, \frac{2}{\sqrt{n}}\right]$$

where n is the number of inputs to the neuron. Usually used for ReLU.

Sigmoid

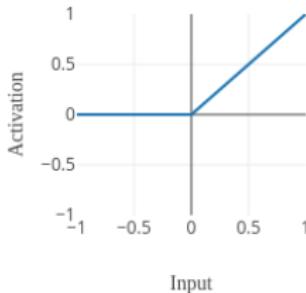
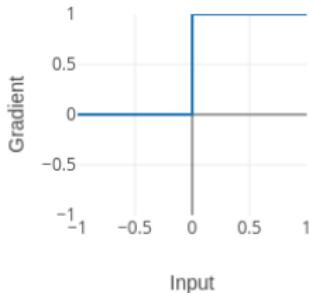


- **Sigmoid activation function:**

$$g(z) = \sigma(z) = \frac{1}{1 + \exp(-z)}$$

- Sigmoidal unit has small gradients across most of its domain (**vanishing gradient**), and hence leads to slow learning, especially in deep models.
 - It is said to saturate easily, where saturation refers to the state in which a neuron predominantly outputs values close to the asymptotic ends of the bounded activation function. Saturation implies slow learning.
- Not recommended as internal unit.

Rectified Linear Units (ReLU)

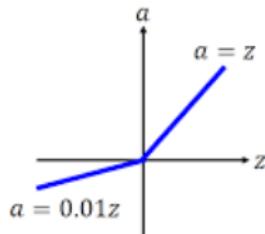


General form of activation function: $g(z) = g(\mathbf{W}^T \mathbf{x} + b)$

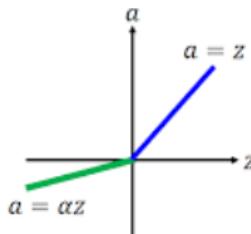
- **ReLU:** $g(z) = \max\{0, z\}$.
- Constant gradient when $z > 0$, which leads to faster learning. Probably the most commonly used activation function.
- Zero gradient when $z < 0$. The neuron is **dead** if $z < 0$ for all training examples. Dead neurons cannot be revived.
 - To mitigate the problem somewhat, initialize b to be small positive value, e.g. 0.1, so that the unit is initially active ($z > 0$).

Variations of ReLU

Leaky ReLU



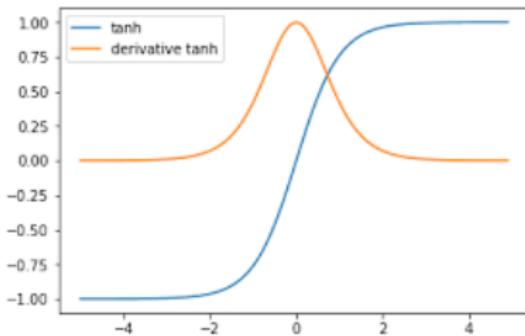
Parametric ReLU



$$g(z, \alpha) = \max\{0, z\} + \alpha \min\{0, z\}$$

- **Absolute value rectification:** $\alpha = -1$,
- **Leaky ReLU:** α is small value like 0.01,
- **Parametric ReLU:** Learns α
- Non-zero gradient even when $z < 0$, which mitigates the dead neuron problem.

Hyperbolic Tangent



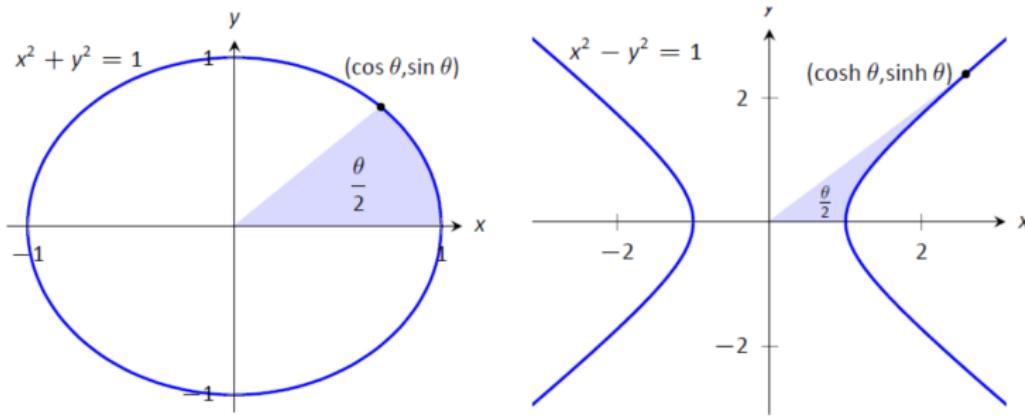
- Hyperbolic tangent activation function:

$$g(z) = \tanh(z) = \frac{\exp(z) - \exp(-z)}{\exp(z) + \exp(-z)}$$

- Tanh has larger gradients than sigmoid, and smoother than ReLU. However, it can still saturate.
- A popular choice in practice.

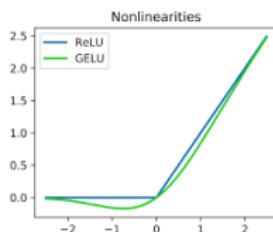
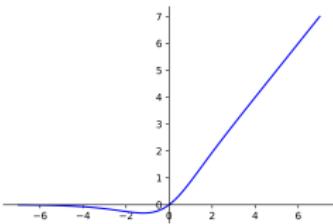
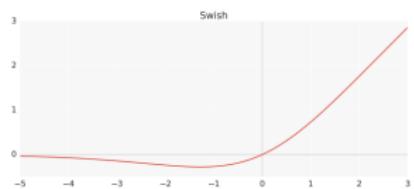
Hyperbolic Tangent

Just as cosine and sine are used to define points on the circle defined by $x^2 + y^2 = 1$, the functions hyperbolic cosine and hyperbolic sine are used to define points on the hyperbola $x^2 - y^2 = 1$:



Source: <https://math.libretexts.org/>

Swish, Mish and GELU



- **Swish activation:**

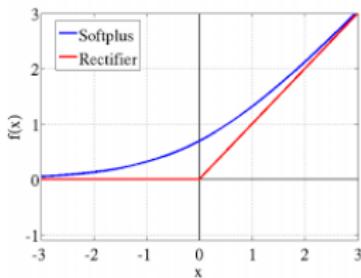
$$g(z) = z\sigma(\beta z), \text{ where } \beta \text{ is parameter fixed or learned}$$

- **Mish activation function:**

$$g(z) = z \tanh(\log(1 + \exp(z))) = z \tanh(\zeta(z))$$

- **GELU activation function:** $g(z) = z\Phi(z)$ Where $\Phi(z)$ is the standard Gaussian cumulative distribution function.
- Non-zero gradients when $z < 0$, smooth and non-monotonic. Unbounded above, which avoids saturation. Outperform ReLU in deep networks.

Softplus

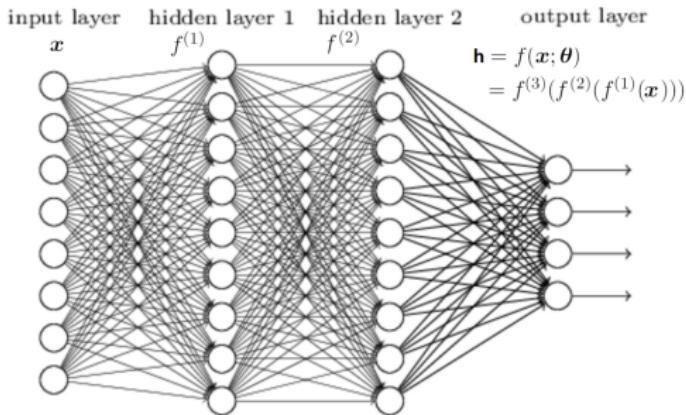


- **Softplus activation function:**

$$g(z) = \zeta(z) = \log(1 + \exp(z))$$

- Softplus is a smooth version of ReLU, empirically not as good as ReLU.

Computation by Feedforward Neural Network



- $\mathbf{h}^{(i)}$ — Column vector for units on layer i ; $\mathbf{b}^{(i)}$ — Biases for units on layer i ;
- $g^{(i)}$ — Activation functions for units on layer i .
- $\mathbf{W}^{(i)}$ — Matrix of weights for units on layer i , with weight for unit j at the j -th column.

$$\begin{aligned}\mathbf{h}^{(1)} &= g^{(1)}(\mathbf{W}^{(1)\top} \mathbf{x} + \mathbf{b}^{(1)}) \\ \mathbf{h}^{(2)} &= g^{(2)}(\mathbf{W}^{(2)\top} \mathbf{h}^{(1)} + \mathbf{b}^{(2)}) \\ \mathbf{h} &= g^{(3)}(\mathbf{W}^{(3)\top} \mathbf{h}^{(2)} + \mathbf{b}^{(3)})\end{aligned}$$

Universal Approximation

Theorem

Only one layer of sigmoid hidden units suffices to approximate any well-behaved function (e.g., bounded continuous functions) to arbitrary precision.

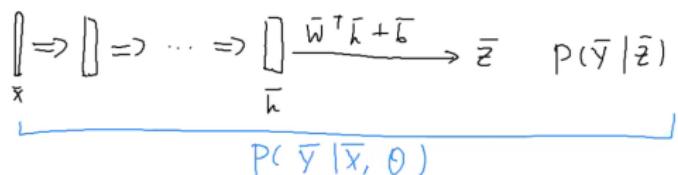
Deep learning useful when you:

- Need a complex function, and
- Have abundant data.

Outline

- 1 Feedforward Neural Network as Function Approximator
- 2 Feedforward Neural Network as Probabilistic Model
- 3 Backpropagation
- 4 Dropout
- 5 Optimization Algorithms

An FNN can be used to define a probabilistic model



- The first-second last layers define a feature transformation:

$$\mathbf{h} = f(\mathbf{x})$$

- The last layer defines a probabilistic model on the features:

$$P(\mathbf{y}|\mathbf{h})$$

Note: Here \mathbf{y} is a vector of output variables, while so far we have been talking about only one output variable y .

- The whole network defines a probabilistic model:

$$P(\mathbf{y}|\mathbf{x}, \theta),$$

where θ consists of weights in f and parameters in $P(\mathbf{y}|\mathbf{h})$.

Logits

- To define a probabilistic model on the features \mathbf{h} , first define a **logit vector**

$$\mathbf{z} = \mathbf{W}^\top \mathbf{h} + \mathbf{b},$$

where \mathbf{W} is a weight matrix and \mathbf{b} is a bias vector. They are the parameters of the last layer. $\mathbf{z} = (z_1, z_2, \dots)^\top$.

- Then, we can define various probabilistic models using \mathbf{z} , which are viewed as units for the last layer.

Linear-Gaussian Output Unit

$$\underbrace{[\![\dots]\!]}_{\mathbf{x}} \Rightarrow \underbrace{[\![\dots]\!]}_{\mathbf{z}} \dots \Rightarrow \underbrace{[\![\dots]\!]}_{\bar{\mathbf{z}}} \xrightarrow{\bar{w}^T \bar{\mathbf{z}} + \bar{b}} \bar{\mathbf{z}} \quad p(\bar{\mathbf{y}} | \bar{\mathbf{z}})$$

- When \mathbf{y} is real-valued vector, we can assume that \mathbf{y} follows a Gaussian distribution with mean \mathbf{z} and identity covariance matrix \mathbf{I} :

$$p(\mathbf{y}|\mathbf{x}) = \mathcal{N}(\mathbf{y}|\mathbf{z}, \mathbf{I})$$

- In this case, the **per-sample loss** is:

$$L(\mathbf{x}, \mathbf{y}, \theta) = -\log P(\mathbf{y}|\mathbf{x}, \theta) \propto \frac{1}{2} \|\mathbf{y} - \mathbf{z}\|_2^2 \quad (\text{See p12 of L02})$$

- Partial derivative of per-sample loss with respect to logit z_k :

$$\frac{\partial L}{\partial z_k} = z_k - y_k.$$

It is error.

Sigmoid Output Unit

- When there is only one binary output variable $y \in \{0, 1\}$, we can define a distribution $p(y|\mathbf{x})$ using a sigmoid unit:

$$P(y|\mathbf{x}) = \text{Ber}(y|\sigma(z)),$$

where z is a scalar.

- In this case, the **per-sample loss** is:

$$L(\mathbf{x}, \mathbf{y}, \theta) = -[y \log \sigma(z) + (1 - y) \log(1 - \sigma(z))] \quad (\text{See p10 of L03})$$

- Partial derivative of per-sample loss with respect to logit z :

$$\frac{\partial L}{\partial z} = \sigma(z) - y. \quad (\text{See p21 of L03})$$

It is error.

Sigmoid Output Unit

- We said earlier that sigmoid units are not recommended for hidden units because they saturates across most of their domains.
- They are fine as output units because the negative log-likelihood in the cost function helps to avoid the problem.
- In fact, $\sigma(z) - y \approx 0$ means:
 - $z \gg 0, y = 1$, or
 - $z \ll 0, y = 0$

In words, saturation occurs only when the model already has the right answer: When $y = 1$ and z is very positive, or $y = 0$ and z is very negative.

Softmax Output Unit

- When y is a variable with C possible values $\{1, 2, \dots, C\}$, we can define a distribution $p(y|\mathbf{x})$ using a **softmax unit**.
- Here $\mathbf{z} = \mathbf{W}^\top \mathbf{h} + \mathbf{b}$ is a vector of n components: $\mathbf{z} = (z_1, z_2, \dots, z_C)^\top$.
- The following formula gives a distribution over the domain of y :

$$P(y = k|\mathbf{x}) = \frac{\exp(z_k)}{\sum_{j=1}^C \exp(z_j)}$$

- Partial derivative of per-sample cross entropy loss with respect to logit z_k :

$$\frac{\partial L}{\partial z_k} = P(y = k) - \mathbf{1}(y = k). \quad (\text{See p40 of L03})$$

It is error.

Outline

- 1 Feedforward Neural Network as Function Approximator
- 2 Feedforward Neural Network as Probabilistic Model
- 3 Backpropagation
- 4 Dropout
- 5 Optimization Algorithms

Training Feedforward Neural Networks

- Given data $\{\mathbf{x}_i, \mathbf{y}_i\}_{i=1}^N$, we want to learn the parameters θ by minimizing the cross-entropy loss:

$$J(\theta) = \frac{1}{N} \sum_{i=1}^N L(\mathbf{x}_i, y_i, \theta) = \frac{1}{N} \sum_{i=1}^N (-\log P(\mathbf{y}_i | \mathbf{x}_i, \theta))$$

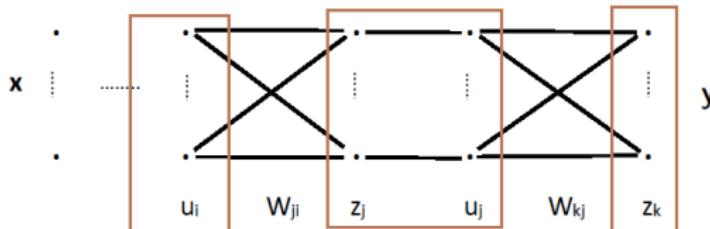
- Algorithm: Stochastic gradient descent (SGD)

- Initialize θ
- Repeat for a predetermined number of epochs
 - For each minibatch B ,

$$\theta \leftarrow \theta - \alpha \frac{1}{|B|} \sum_{i \in B} \nabla L(\mathbf{x}_i, y_i, \theta)$$

- Question: How to compute the gradient $\nabla L(\mathbf{x}_i, y_i, \theta)$ for each training example?
- Answer: [Backpropagation](#).

Backpropagation: Output Layer



- Let y_k , u_j and u_i each be the output of a unit at the last three layers, and z_k , z_j , z_i be the output before applying the activation function:

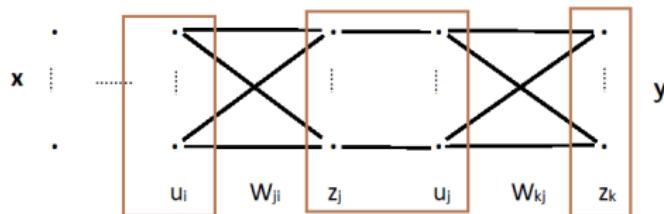
$$z_k = \sum_j u_j W_{kj}, \quad u_j = g(z_j), \quad z_j = \sum_i u_i W_{ji}$$

(Bias ignored for simplicity.)

- Let $L = L(\mathbf{x}, \mathbf{y}, \theta)$ for a particular training example (\mathbf{x}, \mathbf{y}) . The partial derivative of L w.r.t a weigh in the last layer:

$$\frac{\partial L}{\partial W_{kj}} = \frac{\partial L}{\partial z_k} \frac{\partial z_k}{\partial W_{kj}} = u_j \delta_k, \text{ where } \delta_k = \frac{\partial L}{\partial z_k}$$

Backpropagation: Hidden Layer



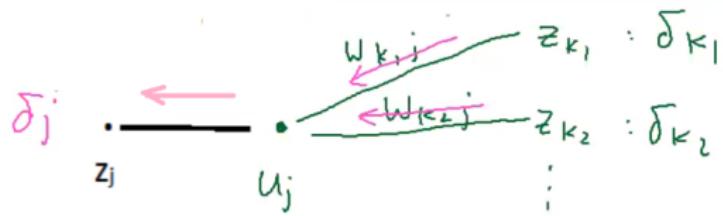
$$z_k = \sum_j u_j W_{kj}, \quad u_j = g(z_j), \quad z_j = \sum_i u_i W_{ji}$$

- The partial derivative of L w.r.t a weigh in the second last layer (similar for weight in other hidden layers):

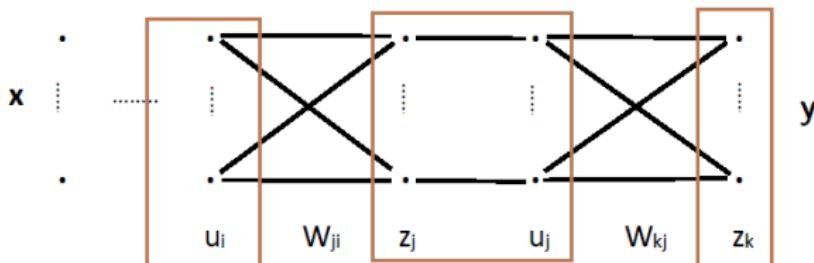
$$\begin{aligned} \frac{\partial L}{\partial W_{ji}} &= \sum_k \frac{\partial L}{\partial z_k} \frac{\partial z_k}{\partial W_{ji}} = \sum_k \delta_k \frac{\partial z_k}{\partial u_j} \frac{\partial u_j}{\partial z_j} \frac{\partial z_j}{\partial W_{ji}} \\ &= \sum_k \delta_k W_{kj} \frac{\partial u_j}{\partial z_j} u_i = u_i \delta_j, \text{ where } \delta_j = \frac{\partial u_j}{\partial z_j} \sum_k W_{kj} \delta_k \end{aligned}$$

Backpropagation of error

$$\delta_j = \frac{\partial u_j}{\partial z_j} \sum_k W_{kj} \delta_k$$



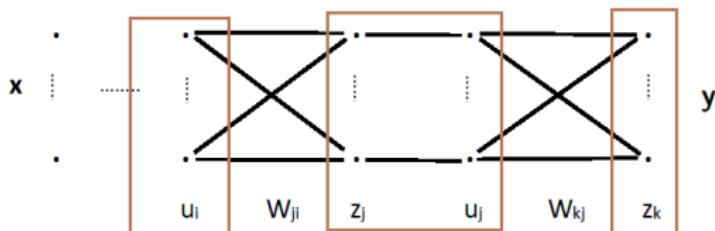
Backpropagation



Summary:

- For a weight W_{kj} the output layer: $\frac{\partial L}{\partial W_{kj}} = u_j \delta_k$, where $\delta_k = \frac{\partial L}{\partial z_k}$
 - u_j is the output of unit j from input x
 - δ_k is the “error value” for unit k obtained by comparing the model output $z = f(x)$ and the observed output y in L .
- For a weight W_{ji} in a hidden layer: $\frac{\partial L}{\partial W_{ji}} = u_i \delta_j$, where $\delta_j = \frac{\partial u_j}{\partial z_j} \sum_k W_{kj} \delta_k$
 - u_i is the output of unit i from input x
 - δ_j is the “error value” for unit j obtained by backpropagating the errors (δ_k) from the next layer .

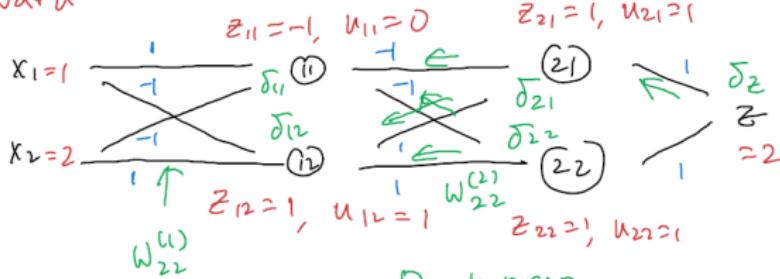
The Backpropagation Algorithm



- Objective: Compute $\frac{\partial L}{\partial W_{ji}}$ for all weights (and biases) for each training example:
 - 1 **Propagate forward** to compute the net input and output of each unit (i.e., u_i , u_j , ...).
 - 2 **Propagate “error” backward**
 - For each output unit k : $\delta_k \leftarrow \frac{\partial L}{\partial z_k}$
 - For each hidden unit j : $\delta_j \leftarrow \frac{\partial u_j}{\partial z_j} \sum_k W_{kj} \delta_k$
 - 3 Get the gradient for each weight: $\frac{\partial L}{\partial W_{ji}} \leftarrow u_i \delta_j$, $\frac{\partial L}{\partial b_j} \leftarrow \delta_j$

Backprop Example

Forward



hidden unit:

$$\text{ReLU } \text{PCY}(z) = \begin{cases} 0(z) & y=1 \\ 1-\sigma(z) & y=0 \end{cases}$$

$$\delta_{11} = 0$$

Back prop

$$\delta_{21} = \frac{\partial u_{21}}{\partial z_{21}} 1 \cdot \delta_2$$

$$= 0.88$$

$$\delta_2 = \sigma(2) - 0 \approx 0.88$$

$$\delta_{12} = \frac{\partial u_{12}}{\partial z_{12}}$$

$$(1 \cdot \delta_{21} + 1 \cdot \delta_{22})$$

$$= 1.76$$

$$\delta_{22} = 0.88$$

$$\frac{\partial L}{\partial w_{22}^{(2)}} = 2 \times 1.76$$

$$\frac{\partial L}{\partial w_{12}^{(1)}} = 1 \times 0.88$$

Example

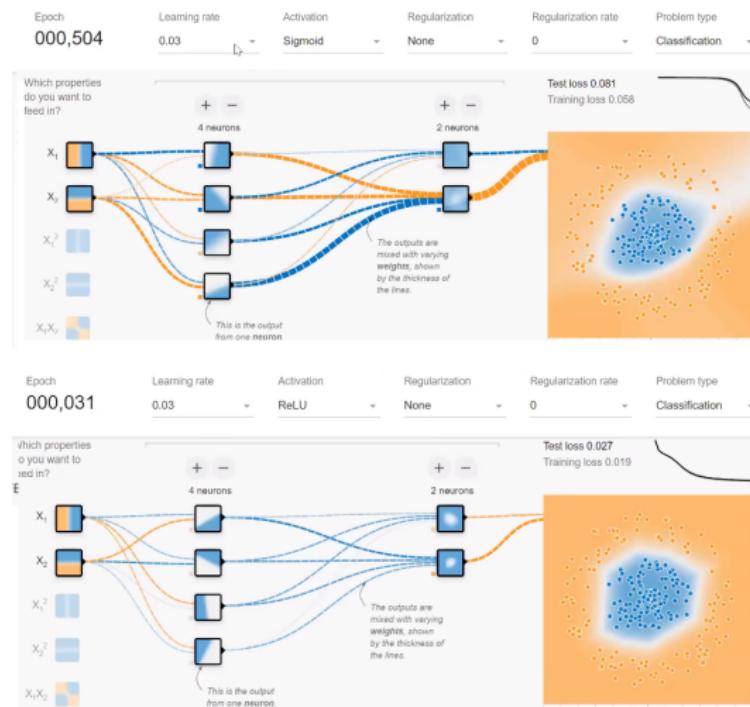
$$\frac{x_1 \quad x_2 \quad y}{1 \quad 2 \quad 0}$$

The Backpropagation Algorithm

- The backpropagation algorithm and SGD are implemented in deep learning packages such as Tensorflow.
- Weights are packed into **tensors** ($[a_{ijk}]$) for efficiency
 - A 1-D tensor is a vector
 - A 2-D tensor is a matrix
 - A 3-D tensor is several matrices stacked on top of each other.
- Although you do not need to know everything about backpropagation and SGD, some basic understanding will be very helpful when you work with deep learning.

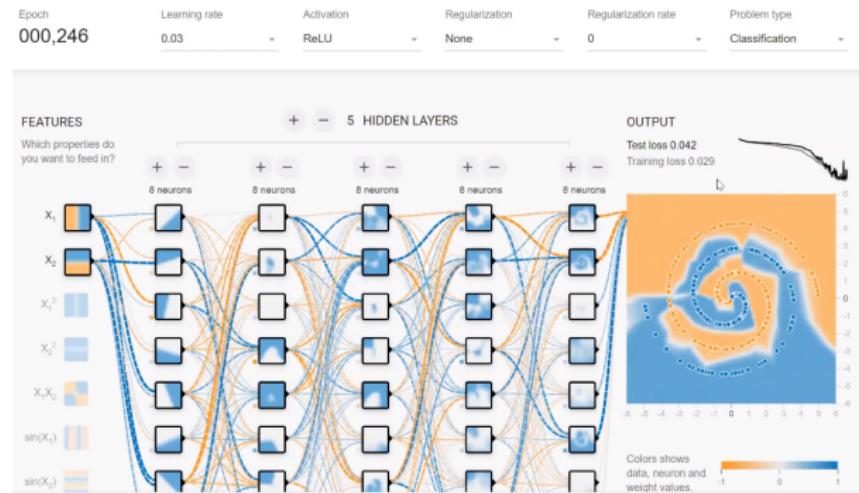
Demonstrations <http://playground.tensorflow.org/>

- ReLU units are easier to learn than Sigmoid and Tanh units



Demonstrations <http://playground.tensorflow.org/>

- Deep structure helps



- Need to adjust learning rate during training.

Outline

- 1 Feedforward Neural Network as Function Approximator
- 2 Feedforward Neural Network as Probabilistic Model
- 3 Backpropagation
- 4 Dropout
- 5 Optimization Algorithms

Introduction

- DL models are complex. Overfitting is an important issue.
- One way to do that is to regularize (weight decay):

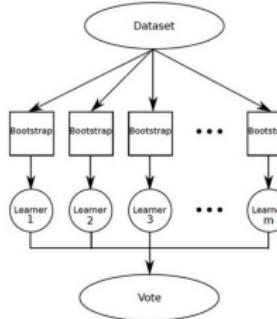
$$\tilde{J}(\theta) = J(\theta) + \lambda \|\theta\|_2^2$$

Almost always do this.

- Dropout is another technique for avoiding overfitting specifically proposed for deep learning.

Bagging (Bootstrap Aggregating)

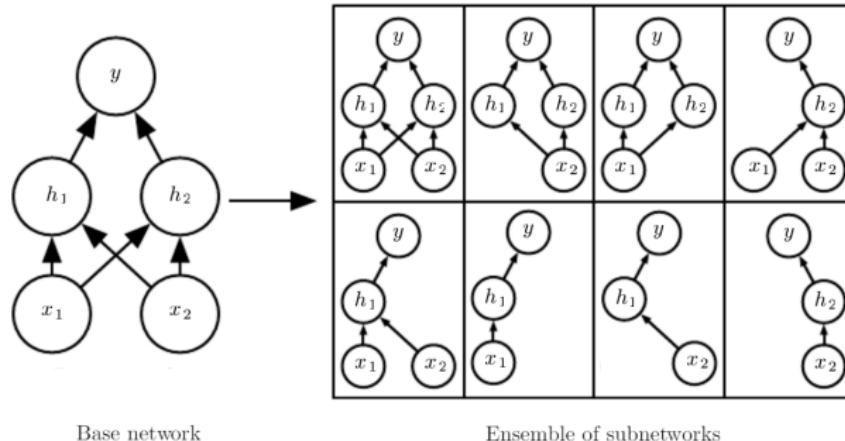
- A technique for reducing generalization error by combining several models
 - Train several different models separately on different randomly selected subsets of data called **bootstrap samples**,
 - Classification: Have all of the models vote on the output for test examples.



- Bagging reduces variance: Variance is error due to randomness. Different bootstrap samples contain different randomness, and hence unlikely to **most** make the same errors on the test set.

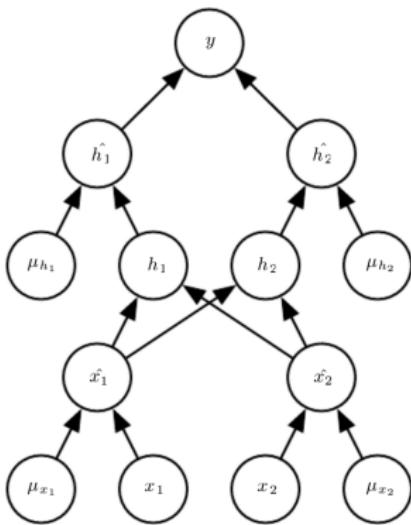
One Way to Use Bagging for Feedforward Networks

- Start with a large network.
- Obtain subnetworks and train them separately on different randomly selected subsets of data.
- To classify future examples, let the subnetworks vote.
- **Problem:** Too expensive.
- **Dropout** is an inexpensive approximation of the idea.



Subnetwork Selection via Masking

- Associate a **binary mask variable** with each input and hidden unit.
- Sample their values randomly and independently (e.g., with probability 0.5 for hidden units and 0.8 for input units).
- Multiply the output of each unit by its mask value before passing it to the next layer.
- Units with 0 mask values are effectively removed from the network.



Dropout

- To be used with minibatch-based algorithms such as SGD.
- For each minibatch,
 - Randomly sample values for mask variables
 - Carry out one step of gradient descent.
 - (Only parameters for the units with mask value 1 are updated. The others have gradient zero and hence not updated.)
- There is only **one network**, no subnetworks. At each step, training is conducted in a part of the network (a subnetwork).
- Dropout is a widely used regularization method for deep learning. It reduces overfitting by preventing complex co-adaptations of parameters on training data.

See example in CNN code.

Outline

- 1 Feedforward Neural Network as Function Approximator
- 2 Feedforward Neural Network as Probabilistic Model
- 3 Backpropagation
- 4 Dropout
- 5 Optimization Algorithms

- Our task is: Find θ to minimize the loss function $J(\theta)$.
- There are multiple algorithms that we can use. Stochastic gradient descent is one of them.

Algorithm 8.1 Stochastic gradient descent (SGD) update at training iteration k

Require: Learning rate ϵ_k .

Require: Initial parameter θ

while stopping criterion not met do

 Sample a minibatch of m examples from the training set $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}\}$ with corresponding targets $\mathbf{y}^{(i)}$.

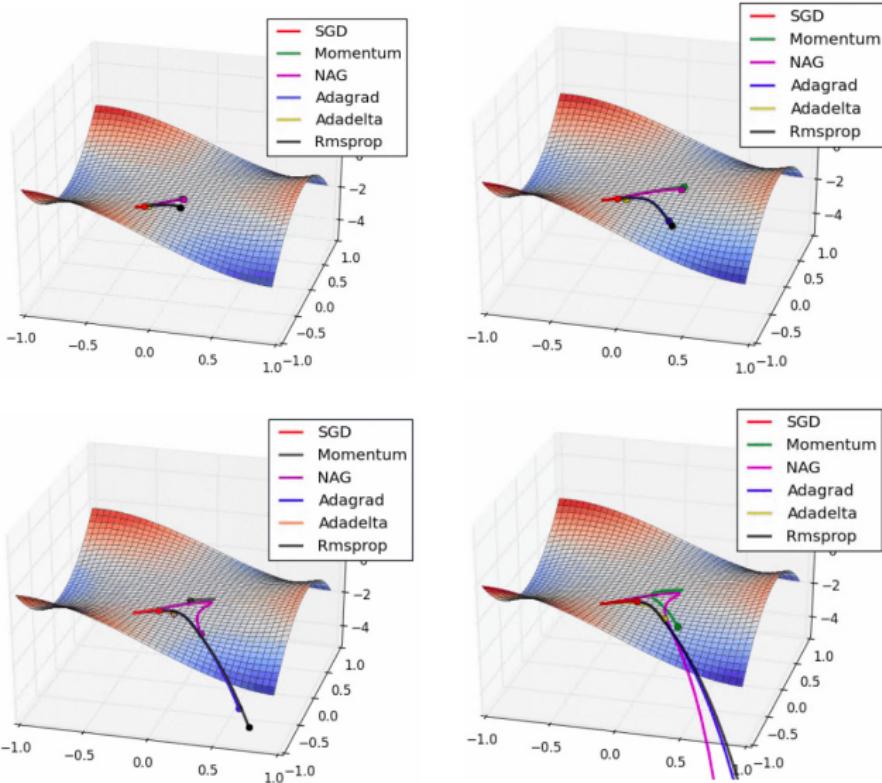
 Compute gradient estimate: $\hat{\mathbf{g}} \leftarrow +\frac{1}{m} \nabla_{\theta} \sum_i L(f(\mathbf{x}^{(i)}; \theta), \mathbf{y}^{(i)})$

 Apply update: $\theta \leftarrow \theta - \epsilon \hat{\mathbf{g}}$

end while

- There are others

Performances of Different Optimizers Vary a Lot



Momentum

SGD

Compute gradient estimate: $\hat{g} \leftarrow +\frac{1}{m} \nabla_{\theta} \sum_i L(f(\mathbf{x}^{(i)}; \boldsymbol{\theta}), \mathbf{y}^{(i)})$
Apply update: $\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} - \epsilon \hat{g}$

SGD with Momentum

Compute gradient estimate: $g \leftarrow \frac{1}{m} \nabla_{\theta} \sum_i L(f(\mathbf{x}^{(i)}; \boldsymbol{\theta}), \mathbf{y}^{(i)})$
Compute velocity update: $v \leftarrow \alpha v - \epsilon g$
Apply update: $\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} + v$

- **Momentum** is a technique to accelerate SGD.
- View θ as a “particle” that travels in the space of parameter values.
- In SGD, its movement is determined by gradient and learning rate.
- In SGD with momentum, the movement is additionally affected by its **velocity** v
 - v is an exponentially decaying average of past negative gradients.
 - The hyperparameter α is usually set at 0.5, 0.9 or 0.99.

Momentum

- The incorporation of momentum into stochastic gradient descent **reduces the variation in overall gradient directions** and **speeds up learning** (larger improvements to parameters at each step).



Online demo: <https://medium.com/@ramrajchandradevan/the-evolution-of-gradient-descend-optimization-algorithm-4106a6702d39>

Algorithms with Adaptive Learning Rates

- Learning rate should be reduced gradually.
- Several common algorithms:
 - AdaGrad
 - RMSProp
 - Adam
- Key idea: Parameters that have changed a lot should be allowed to change less in future.

AdaGrad (Adaptive Gradient Algorithm)

- AdaGrad scales the learning rates of all model parameters by inversely proportional to the square root of the sum of all of their historical squared gradients.
- The net effect is greater progress in the more gently sloped directions of parameter space.

Algorithm 8.4 The AdaGrad algorithm

Require: Global learning rate ϵ

Require: Initial parameter θ

Require: Small constant δ , perhaps 10^{-7} , for numerical stability

Initialize gradient accumulation variable $r = \mathbf{0}$

while stopping criterion not met **do**

 Sample a minibatch of m examples from the training set $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}\}$ with corresponding targets $\mathbf{y}^{(i)}$.

 Compute gradient: $\mathbf{g} \leftarrow \frac{1}{m} \nabla_{\theta} \sum_i L(f(\mathbf{x}^{(i)}; \theta), \mathbf{y}^{(i)})$

 Accumulate squared gradient: $\mathbf{r} \leftarrow \mathbf{r} + \mathbf{g} \odot \mathbf{g}$

 Compute update: $\Delta\theta \leftarrow -\frac{\epsilon}{\delta + \sqrt{\mathbf{r}}} \odot \mathbf{g}$. (Division and square root applied element-wise)

 Apply update: $\theta \leftarrow \theta + \Delta\theta$

end while

RMSProp (Root Mean Square Propagation)

- RMSProp uses an exponentially decaying average to discard history from the extreme past

Algorithm 8.5 The RMSProp algorithm

Require: Global learning rate ϵ , decay rate ρ .

Require: Initial parameter θ

Require: Small constant δ , usually 10^{-6} , used to stabilize division by small numbers.

Initialize accumulation variables $r = 0$

while stopping criterion not met **do**

 Sample a minibatch of m examples from the training set $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}\}$ with corresponding targets $\mathbf{y}^{(i)}$

 Compute gradient: $\mathbf{g} \leftarrow \frac{1}{m} \nabla_{\theta} \sum_i L(f(\mathbf{x}^{(i)}; \theta), \mathbf{y}^{(i)})$

 Accumulate squared gradient: $\mathbf{r} \leftarrow \rho \mathbf{r} + (1 - \rho) \mathbf{g} \odot \mathbf{g}$

 Compute parameter update: $\Delta\theta = -\frac{\epsilon}{\sqrt{\delta + \mathbf{r}}} \odot \mathbf{g}$. ($\frac{1}{\sqrt{\delta + \mathbf{r}}}$ applied element-wise)

 Apply update: $\theta \leftarrow \theta + \Delta\theta$

end while

AdaGrad vs RMSProp

$$\bar{\theta} = \begin{bmatrix} \theta_1 \\ \theta_2 \\ \vdots \\ \theta_n \end{bmatrix} \quad \bar{g} = \begin{bmatrix} g_1 \\ g_2 \\ \vdots \\ g_n \end{bmatrix} \quad \text{'sum of past changes': } \quad \bar{r} = \begin{bmatrix} r_1 \\ r_2 \\ \vdots \\ r_n \end{bmatrix}$$

AdaGrad

$$\bar{r} \leftarrow \bar{r} + \bar{g} \otimes \bar{g}$$

<u>t</u>	<u>r_i</u>
1	g_1^2
2	$g_1^2 + g_2^2$
3	$g_1^2 + g_2^2 + g_3^2$

No decaying

RMSProp

$$\bar{r} \leftarrow p \bar{r} + (1-p) \bar{g} \otimes \bar{g} \quad 0 < p < 1$$

$$\begin{aligned} t & \quad r_i \\ \hline 1 & \quad (1-p) g_1^2 \quad (g_{11}) \\ 2 & \quad p(1-p) g_1^2 + (1-p) g_2^2 \\ & = (1-p)(pg_1^2 + g_2^2) \\ 3 & \quad (1-p)(\underline{p^2} g_1^2 + \underline{pg_2^2} + g_3^2) \end{aligned}$$

Exponentially decaying weights

Adam (Adaptive Moments)

- Roughly a combination of RMSProp and momentum, with bias correction.
- “Insofar, Adam might be the best overall choice.”

Momentum

$$\bar{v} \leftarrow \alpha \bar{v} - \epsilon \bar{g} \quad \Delta \bar{\theta} = \bar{v}$$

$$\bar{\mu} = -\bar{v}$$

$$\bar{\mu} \leftarrow \alpha \bar{\mu} + \epsilon \bar{g} \quad \Delta \bar{\theta} = -\bar{\mu}$$

RMS Prop

$$\bar{r} \leftarrow p \bar{r} + (1-p) \bar{g} \odot \bar{g}$$

$$\Delta \theta = -\frac{\epsilon}{\bar{s} + \sqrt{\bar{r}}} \odot \bar{g}$$

Adam

$$\bar{s} \leftarrow p_1 \bar{s} + (1-p_1) \bar{g} \quad \bar{r} \leftarrow p_2 \bar{r} + (1-p_2) \bar{g} \odot \bar{g}$$

Bias

correction

t: iteration

number

$$\bar{s} \leftarrow \frac{\bar{s}}{1-p_1^t}$$

$$\bar{r} \leftarrow \frac{\bar{r}}{1-p_2^t}$$

$$\Delta \bar{\theta} = -\frac{\epsilon}{\bar{s} + \sqrt{\bar{r}}} \odot \bar{g}$$

Adam: Bias Correction

$$\bar{r} \leftarrow p_2 \bar{r} + (1-p_2) \bar{g} \circ \bar{g} \quad \hat{r} \leftarrow \frac{r}{1-p_2}$$

Adam (Adaptive Moments)

Algorithm 8.7 The Adam algorithm

Require: Step size ϵ (Suggested default: 0.001)

Require: Exponential decay rates for moment estimates, ρ_1 and ρ_2 in $[0, 1]$.
(Suggested defaults: 0.9 and 0.999 respectively)

Require: Small constant δ used for numerical stabilization. (Suggested default:
 10^{-8})

Require: Initial parameters θ

Initialize 1st and 2nd moment variables $s = \mathbf{0}$, $r = \mathbf{0}$

Initialize time step $t = 0$

while stopping criterion not met **do**

 Sample a minibatch of m examples from the training set $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}\}$ with
 corresponding targets $\mathbf{y}^{(i)}$.

 Compute gradient: $\mathbf{g} \leftarrow \frac{1}{m} \nabla_{\theta} \sum_i L(f(\mathbf{x}^{(i)}; \theta), \mathbf{y}^{(i)})$

$t \leftarrow t + 1$

 Update biased first moment estimate: $\hat{s} \leftarrow \rho_1 s + (1 - \rho_1) \mathbf{g}$ Momentum

 Update biased second moment estimate: $\hat{r} \leftarrow \rho_2 r + (1 - \rho_2) \mathbf{g} \odot \mathbf{g}$ RMSPROP

Correct bias in first moment: $\hat{s} \leftarrow \frac{\hat{s}}{1 - \rho_1^t}$

Correct bias in second moment: $\hat{r} \leftarrow \frac{\hat{r}}{1 - \rho_2^t}$

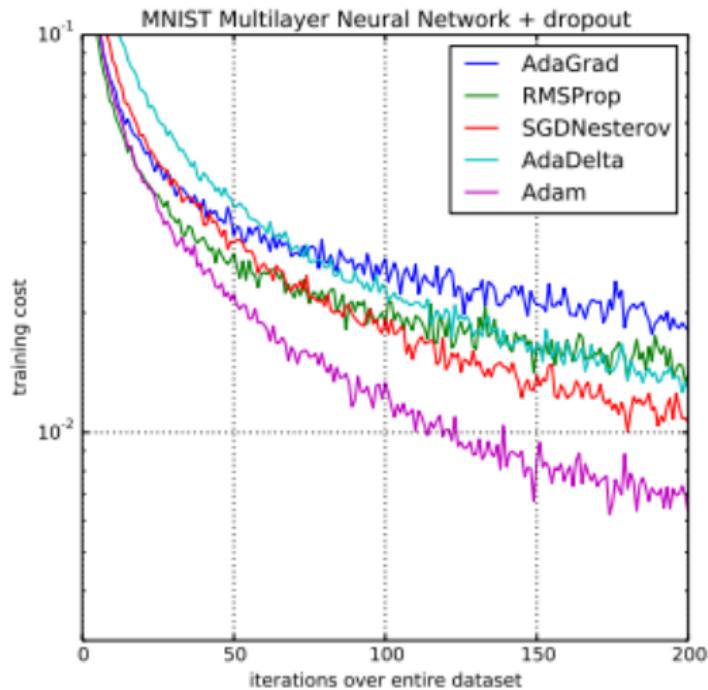
 Compute update: $\Delta\theta = -\epsilon \frac{\hat{s}}{\sqrt{\hat{r} + \delta}}$ (operations applied element-wise)

 Apply update: $\theta \leftarrow \theta + \Delta\theta$

end while

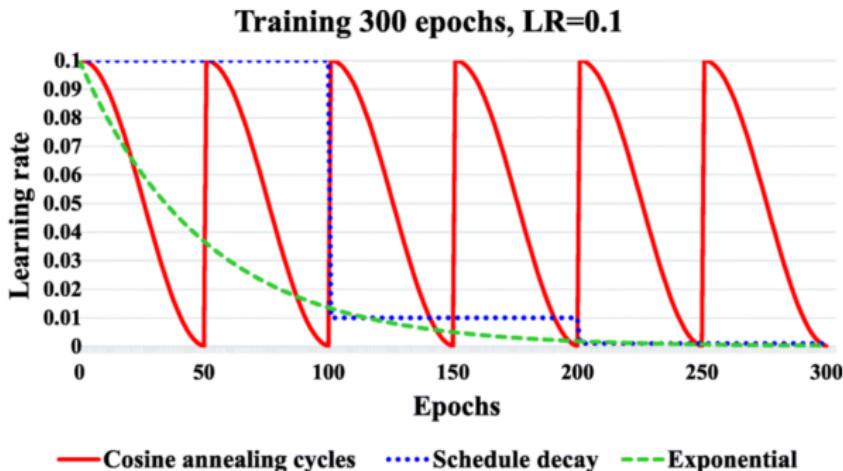
Adam (Adaptive Moments)

Results of different optimizers



Learning Rate Scheduling

- Momentum and Adam are about the use of gradients.
- Here are some ways to schedule the learning rate α



https://www.researchgate.net/figure/Schedule-decay-vs-Cyclic-Cosine-Annealing-vs-Exponential-decay_fig6_336989719

- Cosine Annealing: the resetting of the learning rate acts like a simulated restart of the learning process and the re-use of good weights as the starting point of the restart is referred to as a "warm restart".