



NPTEL ONLINE CERTIFICATION COURSES

Course Name: Deep Learning

Faculty Name: Prof. P. K. Biswas

Department : E & ECE, IIT Kharagpur

Topic

Lecture 41: Popular CNN Models V

CONCEPTS COVERED

Concepts Covered:

☐ CNN

☐ AlexNet

☐ VGG Net

☐ Transfer Learning

☐ Challenges in Deep Learning

☐ GoogLeNet

☐ ResNet

☐ etc.

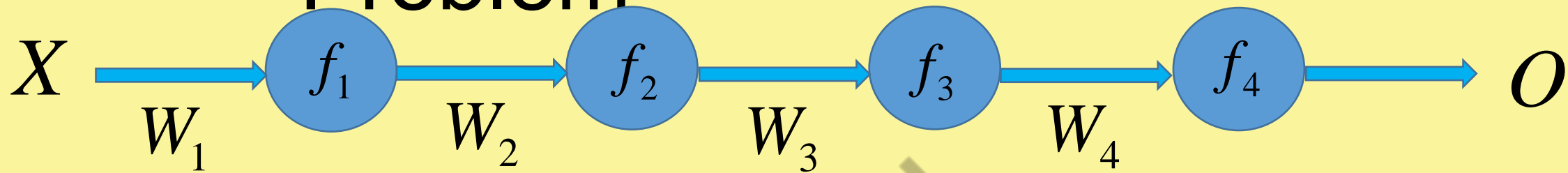


Challenges

- ❑ Deep learning is data hungry.
- ❑ Overfitting or lack of generalization.
- ❑ Vanishing/Exploding Gradient Problem.
- ❑ Appropriate Learning Rate.
- ❑ Covariate Shift.
- ❑ Effective training.



Vanishing Gradient Problem



$$\frac{\partial O}{\partial W_1} = X \cdot f'_1 \cdot W_2 \cdot f'_2 \cdot W_3 \cdot f'_3 \cdot W_4 \cdot f'_4$$



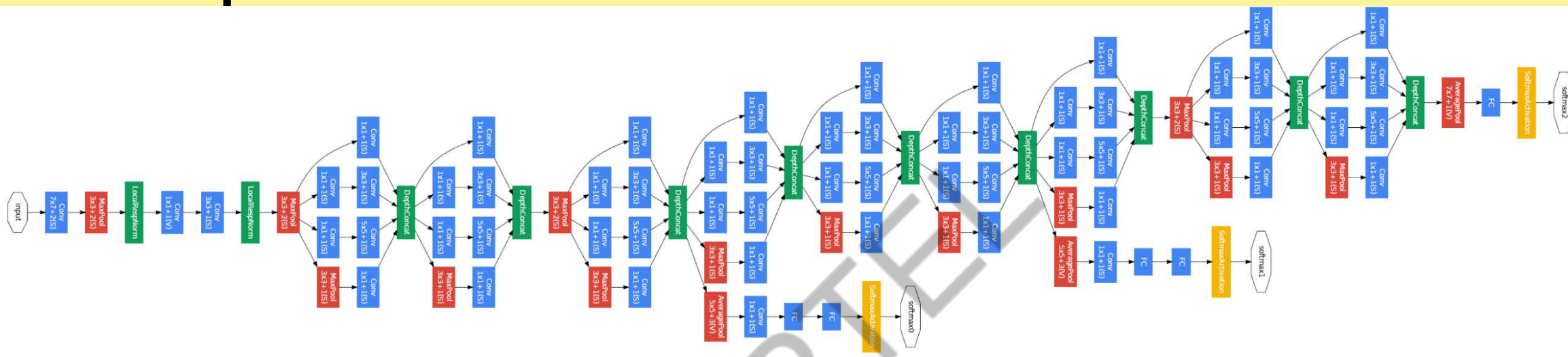
Vanishing Gradient Problem

- ❑ Choice of activation function: ReLU instead of Sigmoid.
- ❑ Appropriate initialization of weights.
- ❑ Intelligent Back Propagation Learning Algorithm.



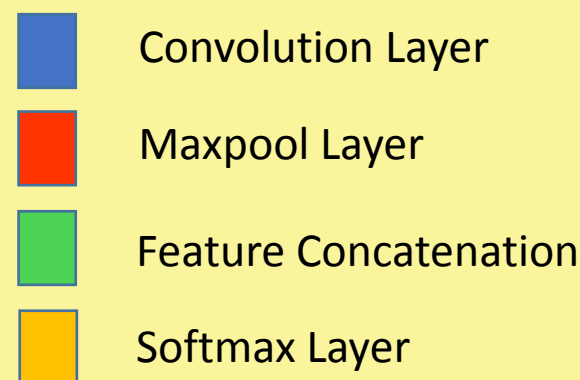
GoogLeNet ILSVRC 2014 Winner





❖ 22 Layers with parameters

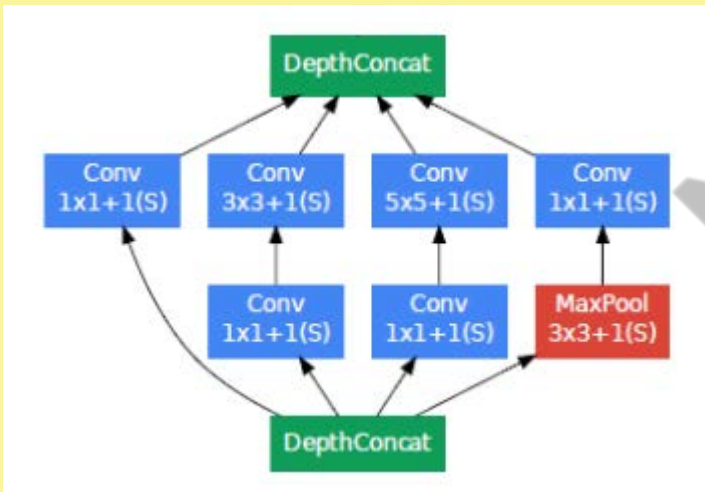
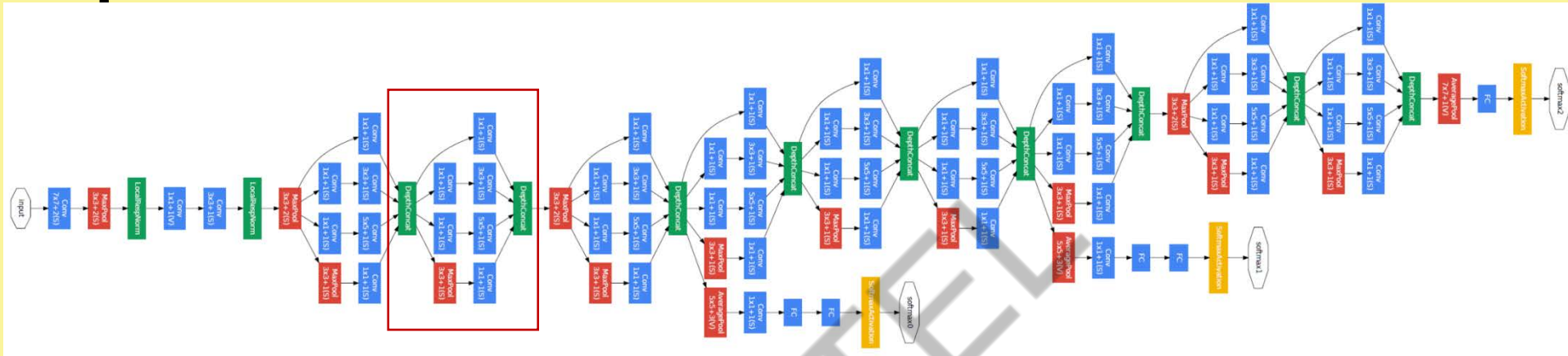
❖ 27 Layer including Maxpool layers



GoogLeNe

13

†



Inception Module

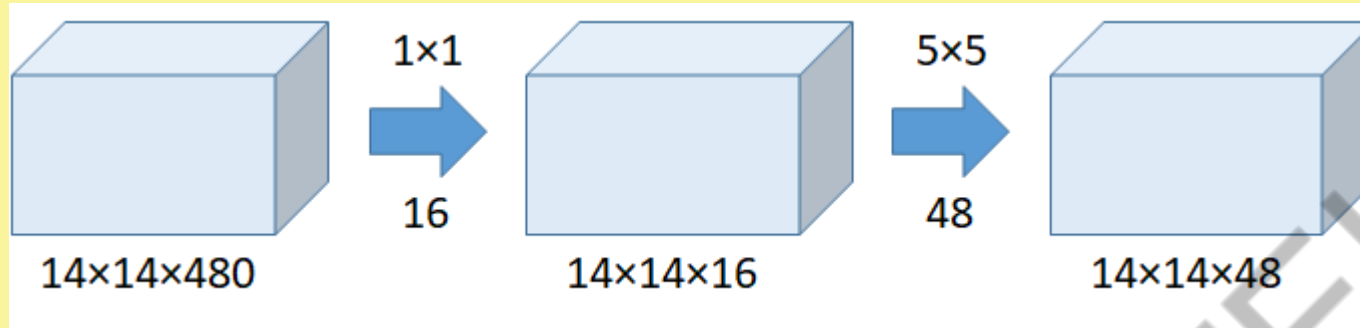


Inception Module

- ❑ Computing 1×1 , 3×3 , and 5×5 convolutions within the same module of the network.
- ❑ Covers a bigger area, at the same time preserves fine resolution for small information on the images.
- ❑ Use different convolution kernels of different sizes in parallel from the most accurate detailing (1×1) to a bigger one (5×5).
- ❑ 1×1 convolution also reduces computation.



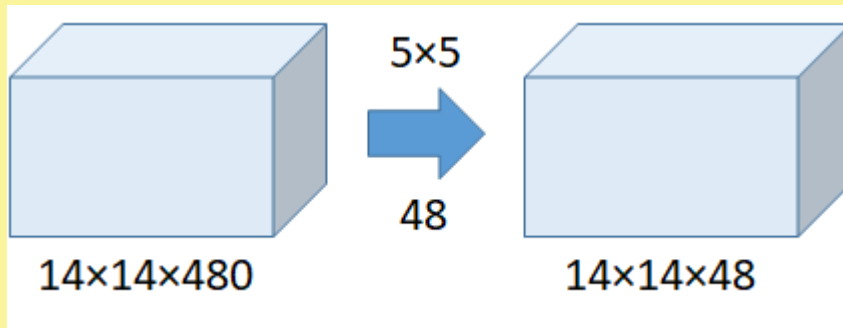
Inception Module



Number of operations for 1×1 = $(14 \times 14 \times 16) \times (1 \times 1 \times 480) = 1.5\text{M}$

Number of operations for 5×5 = $(14 \times 14 \times 48) \times (5 \times 5 \times 16) = 3.8\text{M}$

Total number of operations = $1.5\text{M} + 3.8\text{M} = 5.3\text{M}$



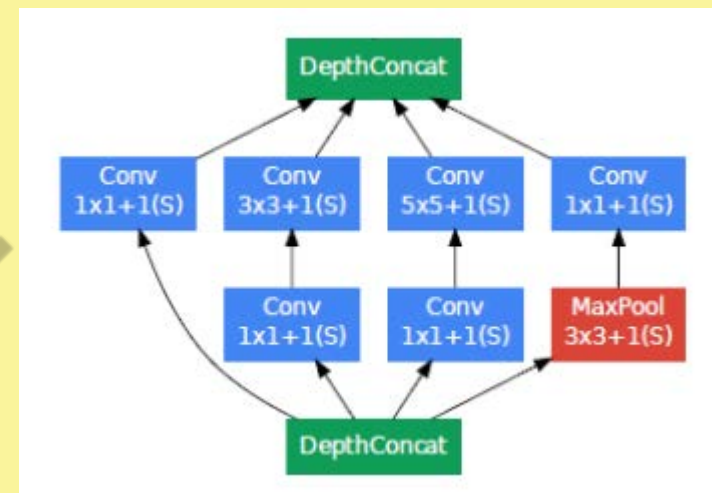
Number of operations = $(14 \times 14 \times 48) \times (5 \times 5 \times 480) = 112.9\text{M}$



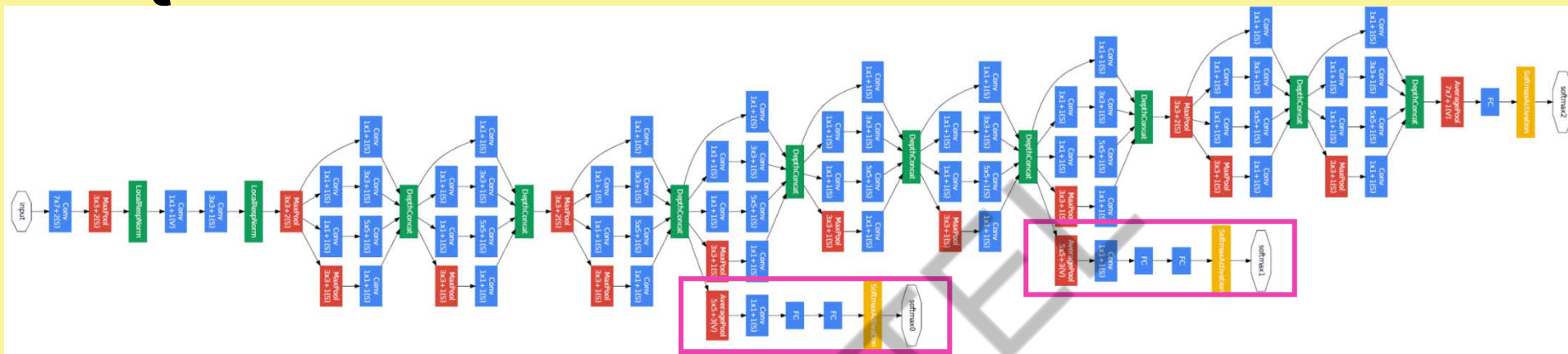
<https://medium.com/coinmonks/paper-review-of-googlenet-inception-v1-winner-of-ilsvlc-2014-image-classification-c2b3565a64e7>

Inception Module

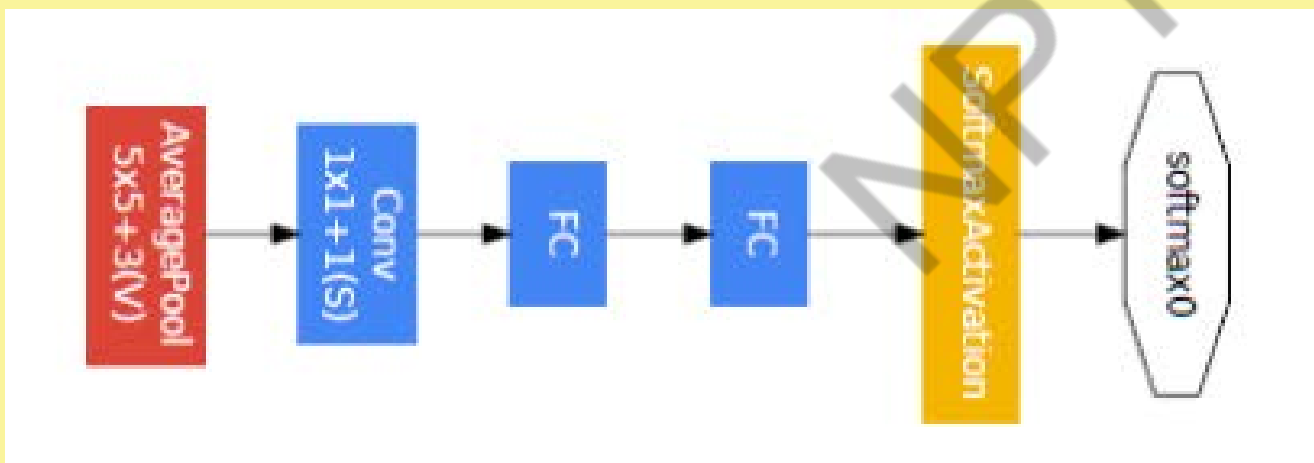
- ❑ Outputs of these filters are then stacked along the channel dimension.
- ❑ Multi-level feature extractor.
- ❑ There are 9 such inception modules.
- ❑ Top-5 error rate of less than 7 %.



GoogLeNet



Auxiliary Classifier



Auxiliary Classifier

- ❑ Due to large depth of the network, ability to propagate gradient back through all the layers was a concern.
- ❑ Auxiliary Classifiers are smaller CNNs put on top of middle Inception modules.
- ❑ Addition of auxiliary classifiers in the middle exploits the discriminative power of the features produced by the layers in the middle.



Auxiliary Classifier

- ❑ During training, loss of Auxiliary classifiers are added to the total loss of the network.
- ❑ Losses from Auxiliary classifiers were weighted by 0.3.
- ❑ Auxiliary classifiers are discarded at Inference time.





NPTEL ONLINE CERTIFICATION COURSES

*Thank
you*





NPTEL ONLINE CERTIFICATION COURSES

Course Name: Deep Learning

Faculty Name: Prof. P. K. Biswas

Department : E & ECE, IIT Kharagpur

Topic

Lecture 42: Popular CNN Models VI

CONCEPTS COVERED

Concepts Covered:

- ❑ CNN

- ❑ Challenges in Deep Learning

- ❑ GoogLeNet

- ❑ ResNet

- ❑ Momentum Optimizer



Challenges

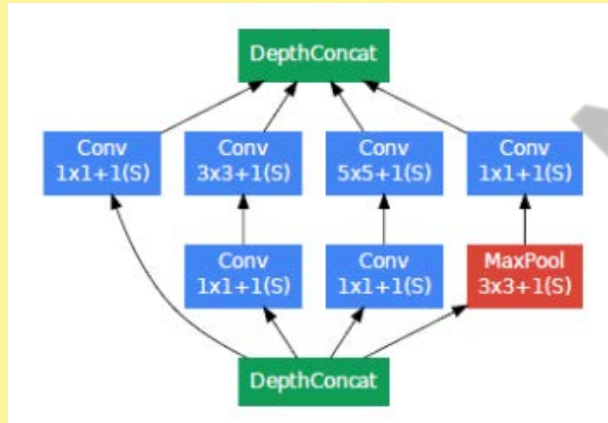
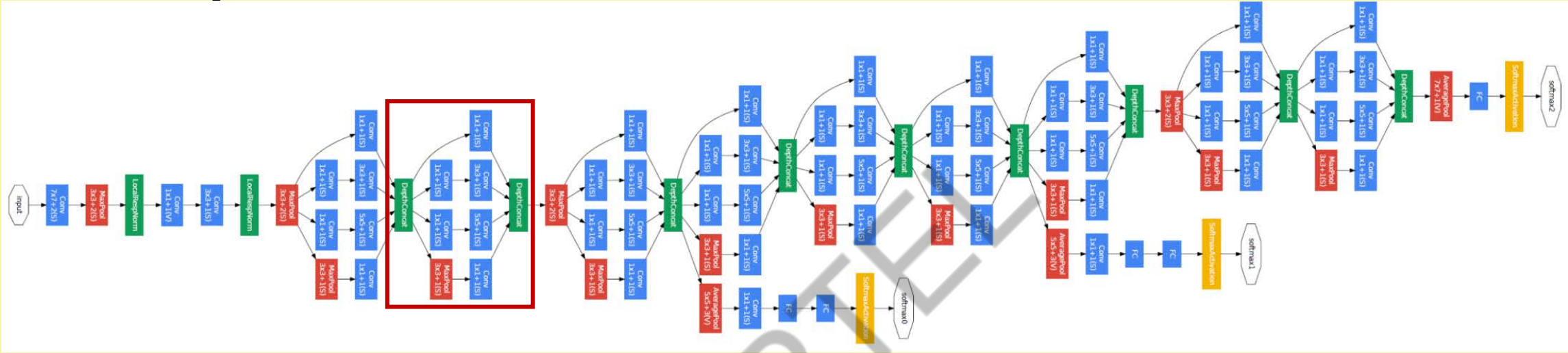
- ❑ Deep learning is data hungry.
- ❑ Overfitting or lack of generalization.
- ❑ Vanishing/Exploding Gradient Problem.
- ❑ Appropriate Learning Rate.
- ❑ Covariate Shift.
- ❑ Effective training.



Vanishing Gradient Problem

- ❑ Choice of activation function: ReLU instead of Sigmoid.
- ❑ Appropriate initialization of weights.
- ❑ Intelligent Back Propagation Learning Algorithm.

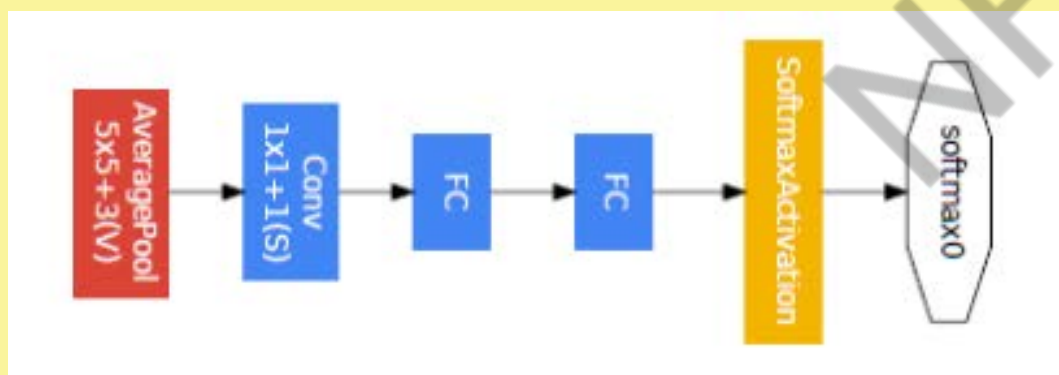
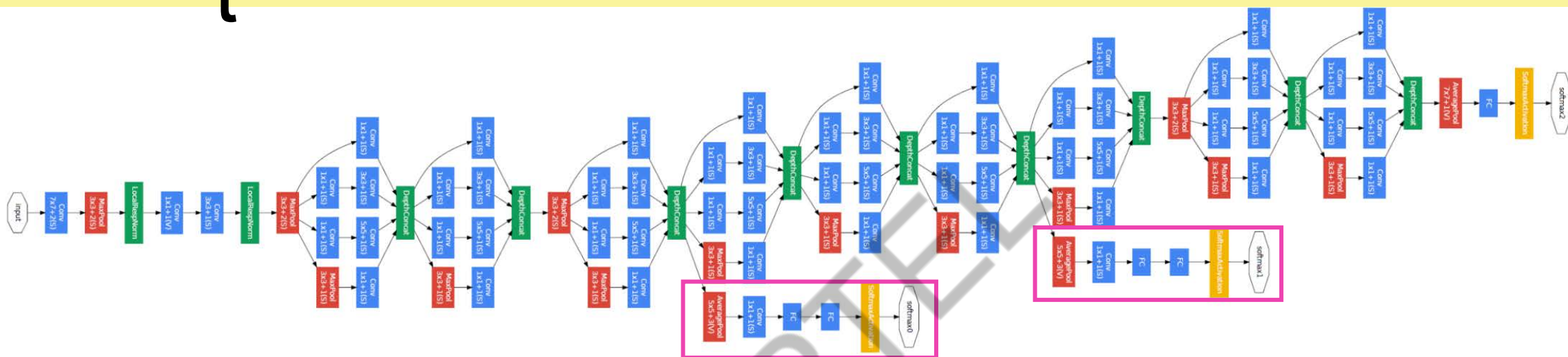




Inception Module



GoogLeNet



Auxiliary Classifier



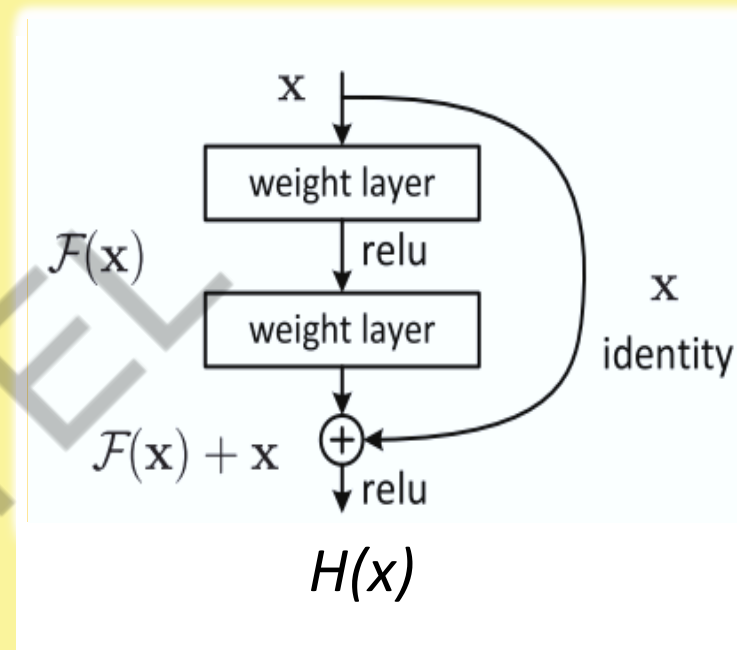
ResNet



ResNe

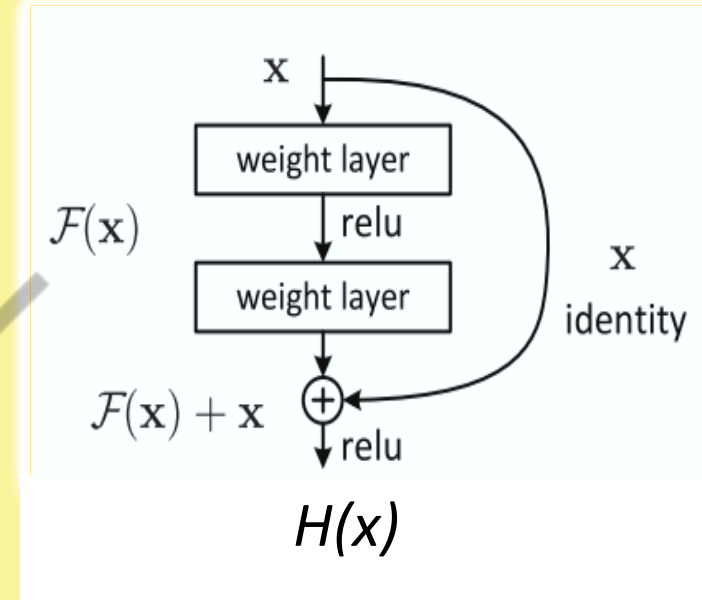
t

- ❑ Core idea is: introduction of Skip Connection/ Identity Shortcut Connection that skips one or more layers.
- ❑ Stacking layers should not degrade performance compared to its shallow counterpart.
- ❑ Weight layer learns $F(x)=H(x)-x$



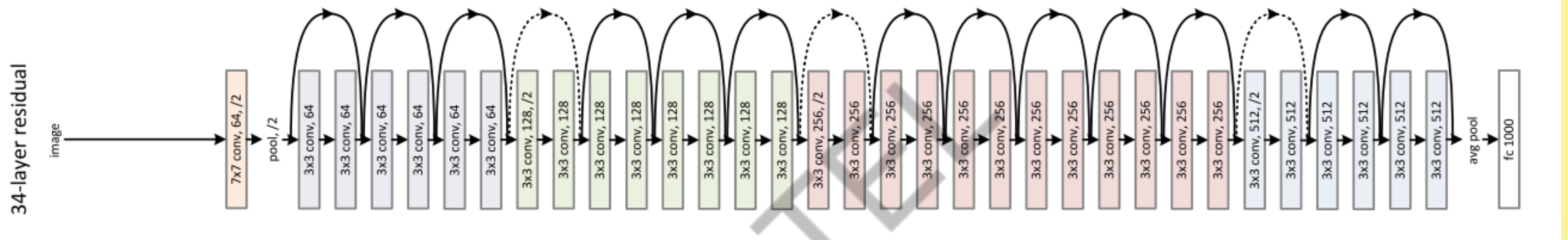
ResNe

- By stacking identity mappings the resultant deep network should give at least same performance as its shallow counterpart.
- Deeper network should not give higher training error than shallow network.
- During learning the gradient can flow to any earlier network through shortcut connections alleviating vanishing gradient problem.



ResNet

3



<https://towardsdatascience.com/an-overview-of-resnet-and-its-variants-5281e2f56035>

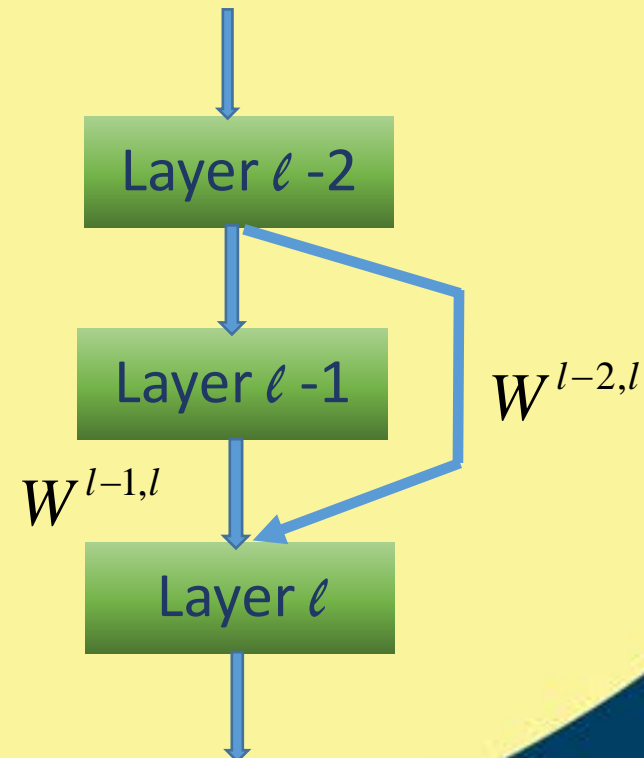
ResNe

t

Forward flow:

$$\begin{aligned} a^l &= f(W^{l-1,l} \cdot a^{l-1} + b^l + W^{l-2,l} \cdot a^{l-2}) \\ &= f(Z^l + W^{l-2,l} \cdot a^{l-2}) \end{aligned}$$

$$a^l = f(Z^l + a^{l-2}) \quad \text{if same dimension}$$

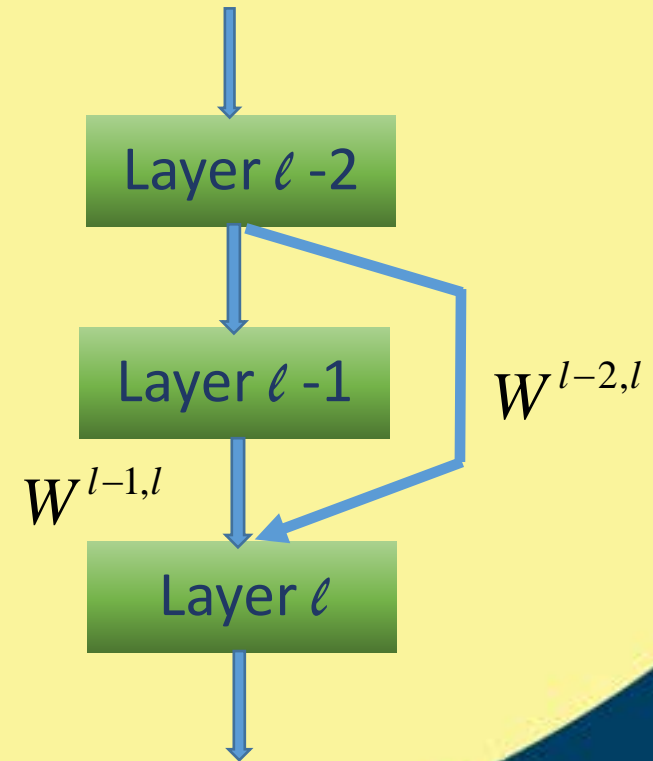


Backward Propagation:

$$\nabla W^{l-1,l} = -a^{l-1} \cdot \delta^l \quad \text{normal path}$$

$$\nabla W^{l-2,l} = -a^{l-2} \cdot \delta^l \quad \text{skip path}$$

If the skip path has fixed weights, identity matrix, then they are not updated.



Challenges

- ☐ Deep learning is data hungry.
- ☐ Overfitting or lack of generalization.
- ☐ Vanishing/Exploding Gradient Problem.
- ☐ Appropriate Learning Rate.
- ☐ Covariate Shift.
- ☐ Effective training.



Optimizing Gradient Descent



Gradient Descent Challenges

Challenges of Mini-batch Gradient Descent

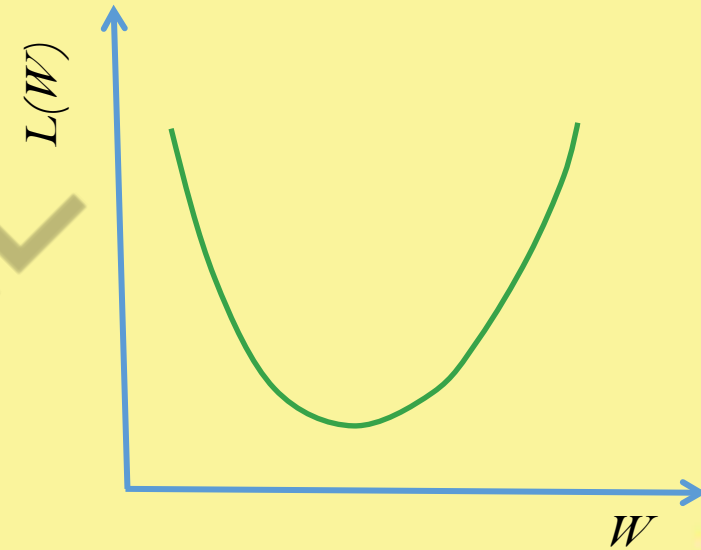
- ☐ Choice of Proper Learning Rate:
 - ☐ Too small a learning rate leads to slow convergence.
 - ☐ A large learning rate may lead to oscillation around the minima or may even diverge.



Gradient Descent Challenges

Challenges of Mini-batch Gradient Descent

- ❑ Choice of Proper Learning Rate:
 - ❑ Too small a learning rate leads to slow convergence.
 - ❑ A large learning rate may lead to oscillation around the minima or may even diverge.



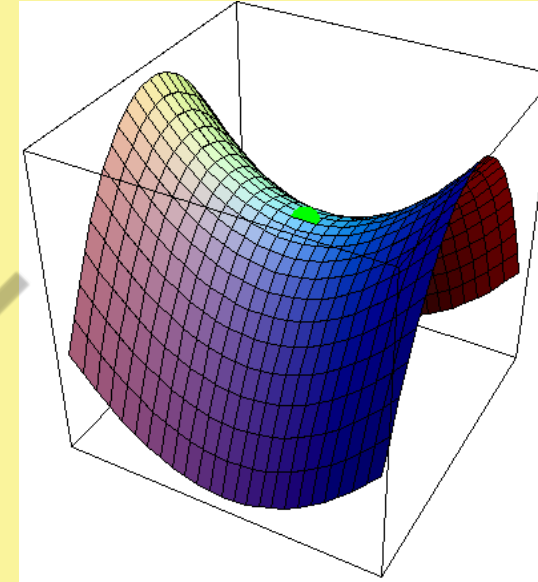
Gradient Descent Challenges

- ☐ Learning Rate Schedules: changing learning rate according to some predefined schedule.
- ☐ The same learning rate applies to all parameter updates.
- ☐ The data may be sparse and different features have very different frequencies.
- ☐ Updating all of them to the same extent might not be proper.
- ☐ Larger update for rarely occurring features might be a better choice.

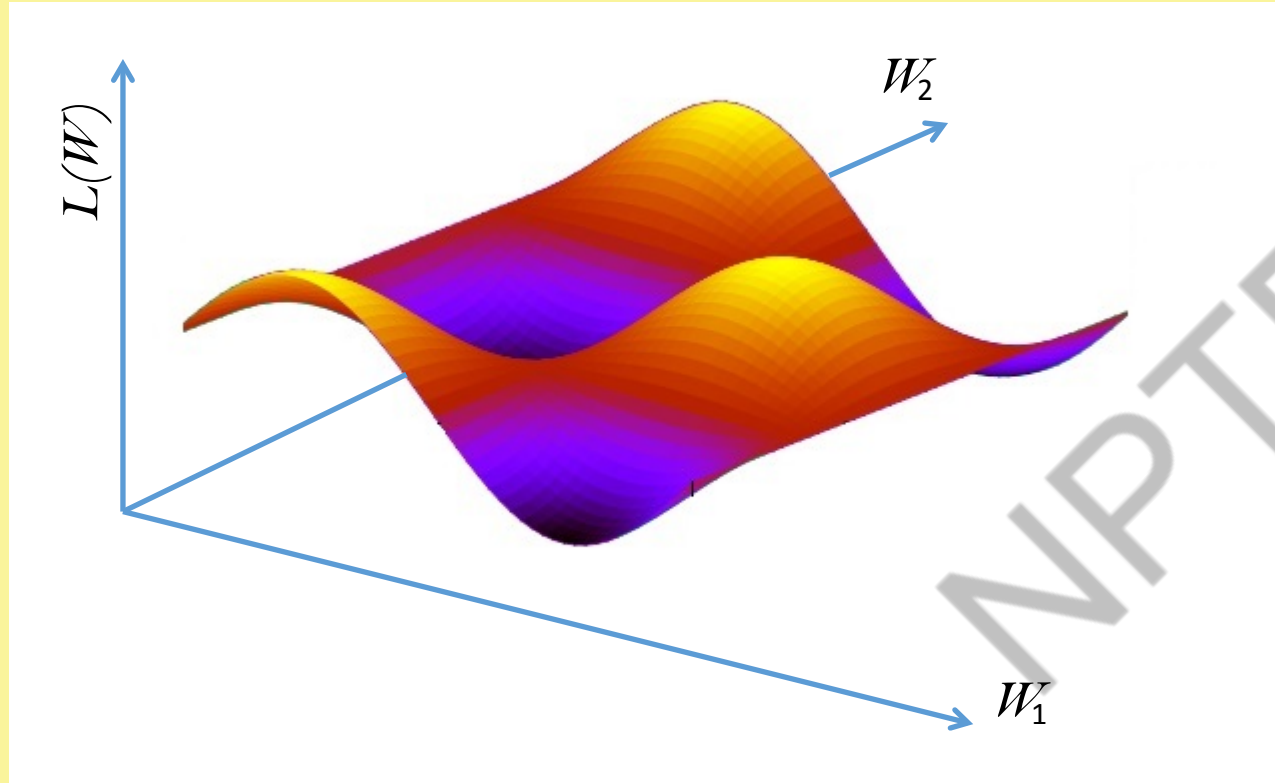


Gradient Descent Challenges

- ❑ Avoiding getting trapped in suboptimal local minima.
- ❑ Difficulty arises in from saddle points, i.e. points where one dimension slopes up and another slopes down.
- ❑ These saddle points are usually surrounded by a plateau of the same error, which makes it hard for SGD to escape, as the gradient is close to zero in all dimensions.



Momentum Optimizer





NPTEL ONLINE CERTIFICATION COURSES

*Thank
you*





NPTEL ONLINE CERTIFICATION COURSES

Course Name: Deep Learning

Faculty Name: Prof. P. K. Biswas

Department : E & ECE, IIT Kharagpur

Topic

Lecture 43: Popular Optimizing Gradient Descent

Challenges

- ☐ Deep learning is data hungry.
- ☐ Overfitting or lack of generalization.
- ☐ Vanishing/Exploding Gradient Problem.
- ☐ Appropriate Learning Rate.
- ☐ Covariate Shift.
- ☐ Effective training.



CONCEPTS COVERED

Concepts Covered:

- ❑ CNN

- ❑ ResNet

- ❑ Gradient Descent Challenges

- ❑ Momentum Optimizer

- ❑ Nesterov Accelerated Gradient

- ❑ Adagrad.

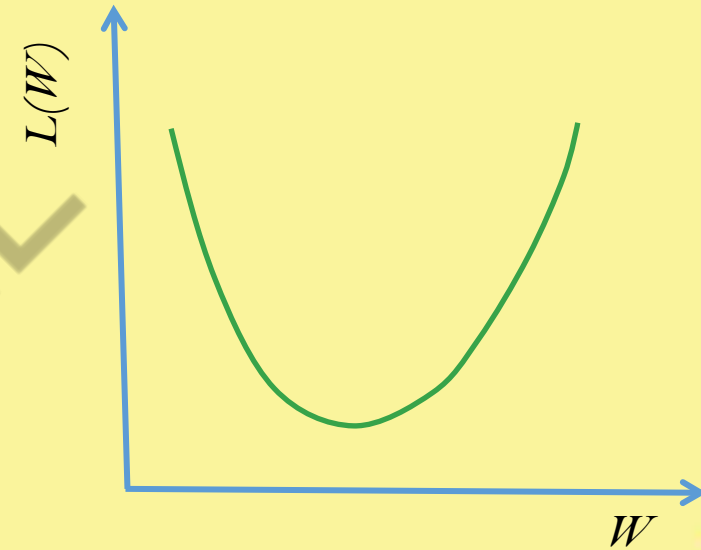
- ❑ etc.



Gradient Descent Challenges

Challenges of Mini-batch Gradient Descent

- ❑ Choice of Proper Learning Rate:
 - ❑ Too small a learning rate leads to slow convergence.
 - ❑ A large learning rate may lead to oscillation around the minima or may even diverge.



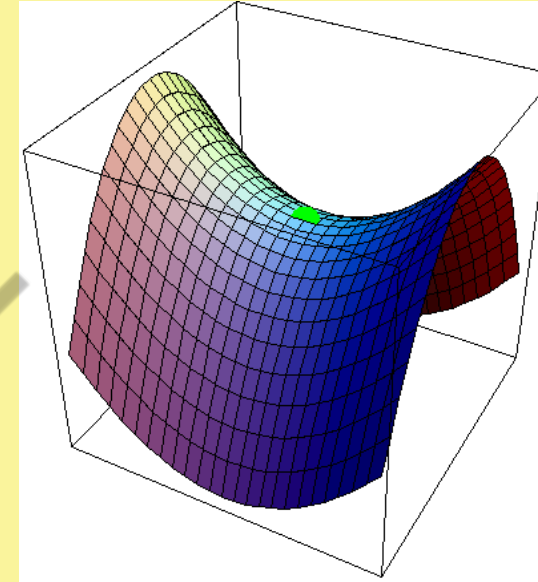
Gradient Descent Challenges

- ☐ Learning Rate Schedules: changing learning rate according to some predefined schedule.
- ☐ The same learning rate applies to all parameter updates.
- ☐ The data may be sparse and different features have very different frequencies.
- ☐ Updating all of them to the same extent might not be proper.
- ☐ Larger update for rarely occurring features might be a better choice.



Gradient Descent Challenges

- ❑ Avoiding getting trapped in suboptimal local minima.
- ❑ Difficulty arises from saddle points, i.e. points where one dimension slopes up and another slopes down.
- ❑ These saddle points are usually surrounded by a plateau of the same error, which makes it hard for SGD to escape, as the gradient is close to zero in all dimensions.



Optimizing Gradient Descent



CONCEPTS COVERED

Concepts Covered:

- ❑ CNN

 - ❑ ResNet

 - ❑ Gradient Descent Challenges

 - ❑ Momentum Optimizer

 - ❑ Adagrad.

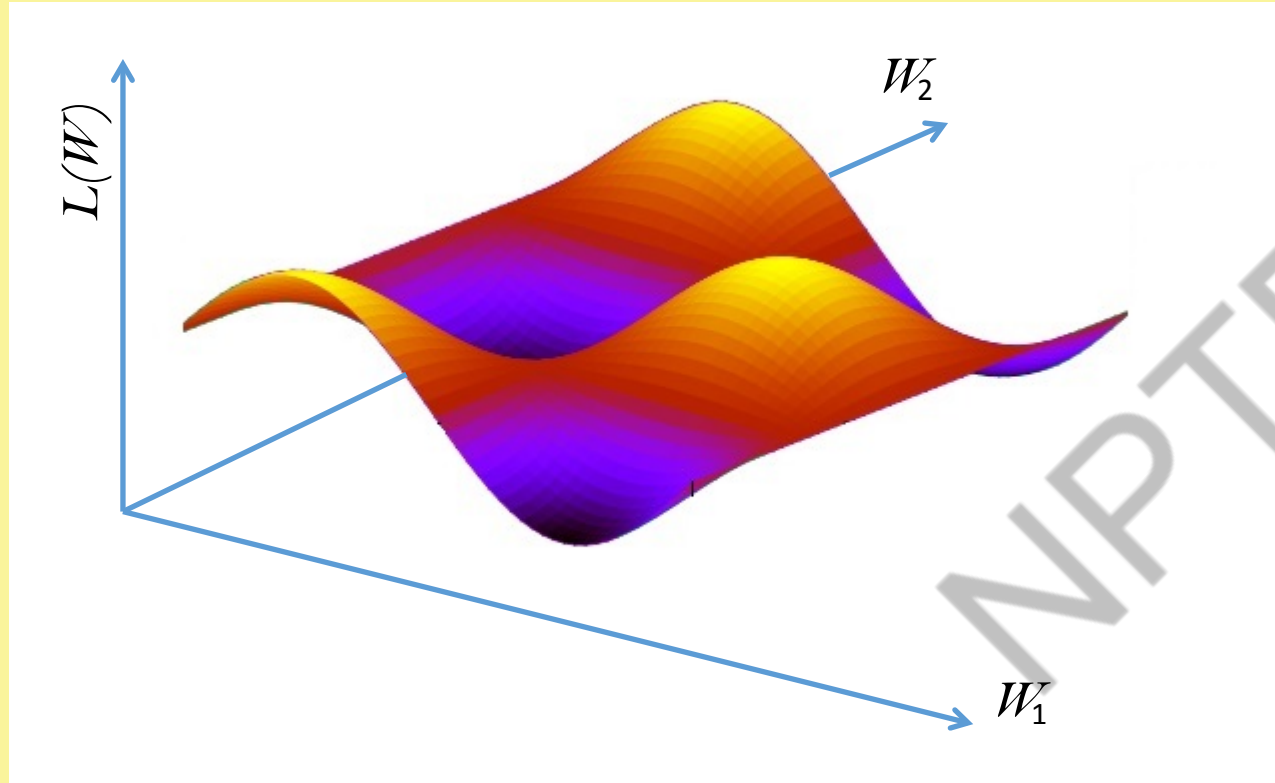
 - ❑ etc.



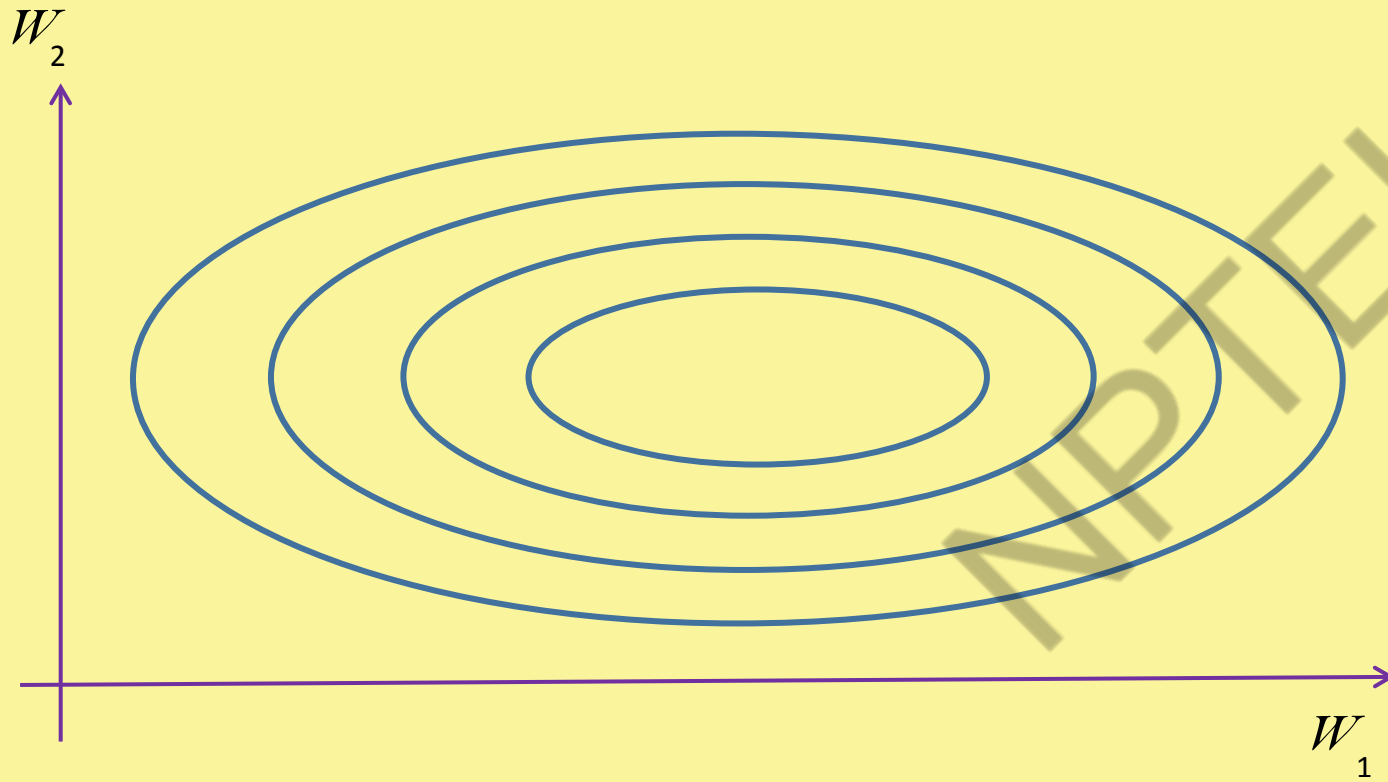
Momentum Optimizer



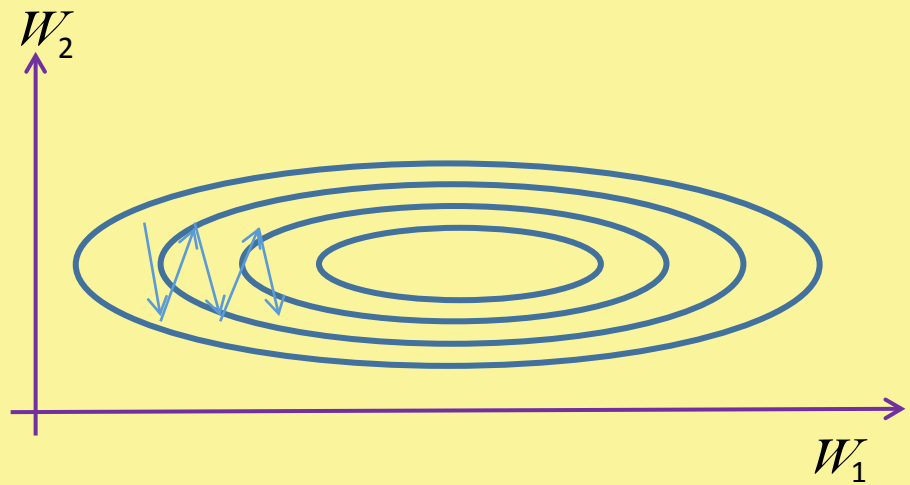
Momentum Optimizer



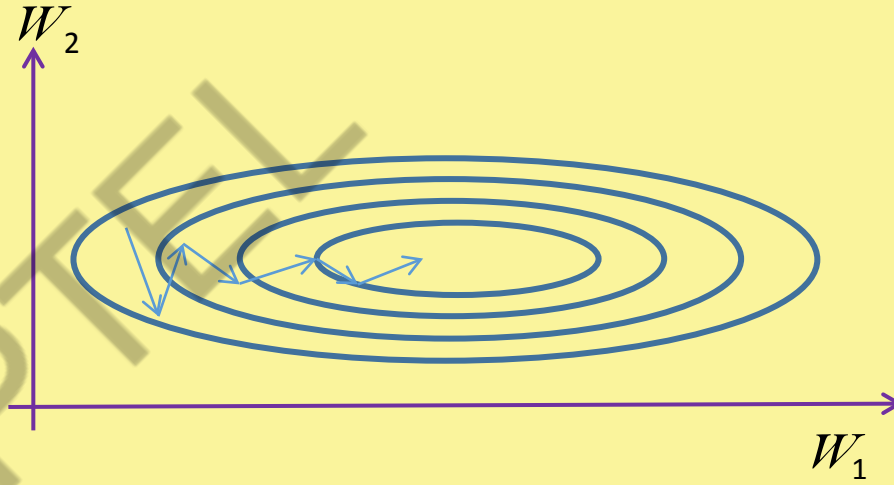
Momentum Optimizer



Momentum Optimizer



SGD



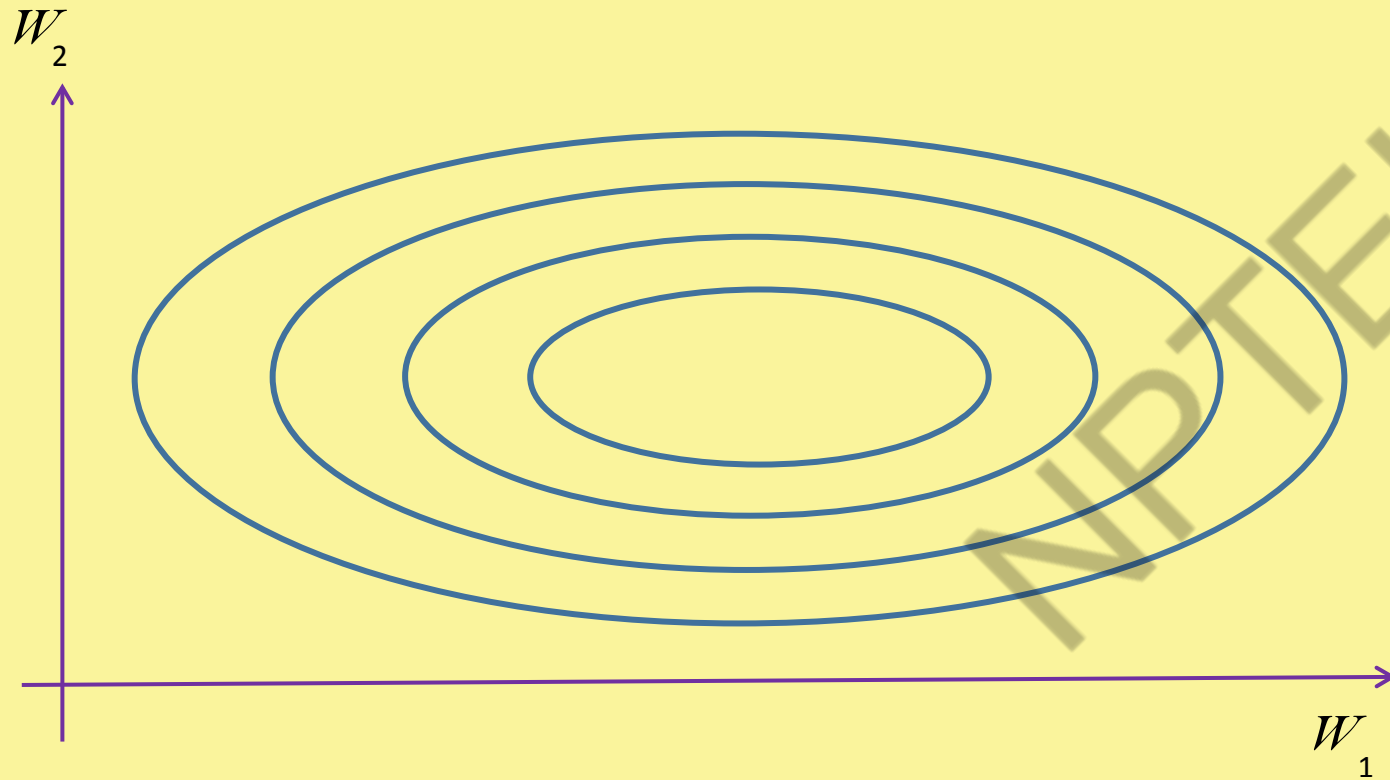
SGD with Momentum



Nesterov Accelerated Gradient (NAG)



Nesterov Accelerated Gradient (NAG)



Problem with Momentum Optimizer/NAG

- ❑ Both the algorithms require the hyper-parameters to be set manually.
- ❑ These hyper-parameters decide the learning rate.
- ❑ The algorithm uses same learning rate for all dimensions.
- ❑ The high dimensional (mostly) non-convex nature of loss function may lead to different sensitivity on different dimension.
- ❑ We may require learning rate could be small in some dimension and large in another dimension.





NPTEL ONLINE CERTIFICATION COURSES

*Thank
you*





NPTEL ONLINE CERTIFICATION COURSES

Course Name: Deep Learning

Faculty Name: Prof. P. K. Biswas

Department : E & ECE, IIT Kharagpur

Topic

Lecture 44: Optimizing Gradient Descent II

CONCEPTS COVERED

Concepts Covered:

☐ CNN

- ☐ Gradient Descent Challenges

- ☐ Momentum Optimizer

- ☐ Nesterov Accelerated Gradient

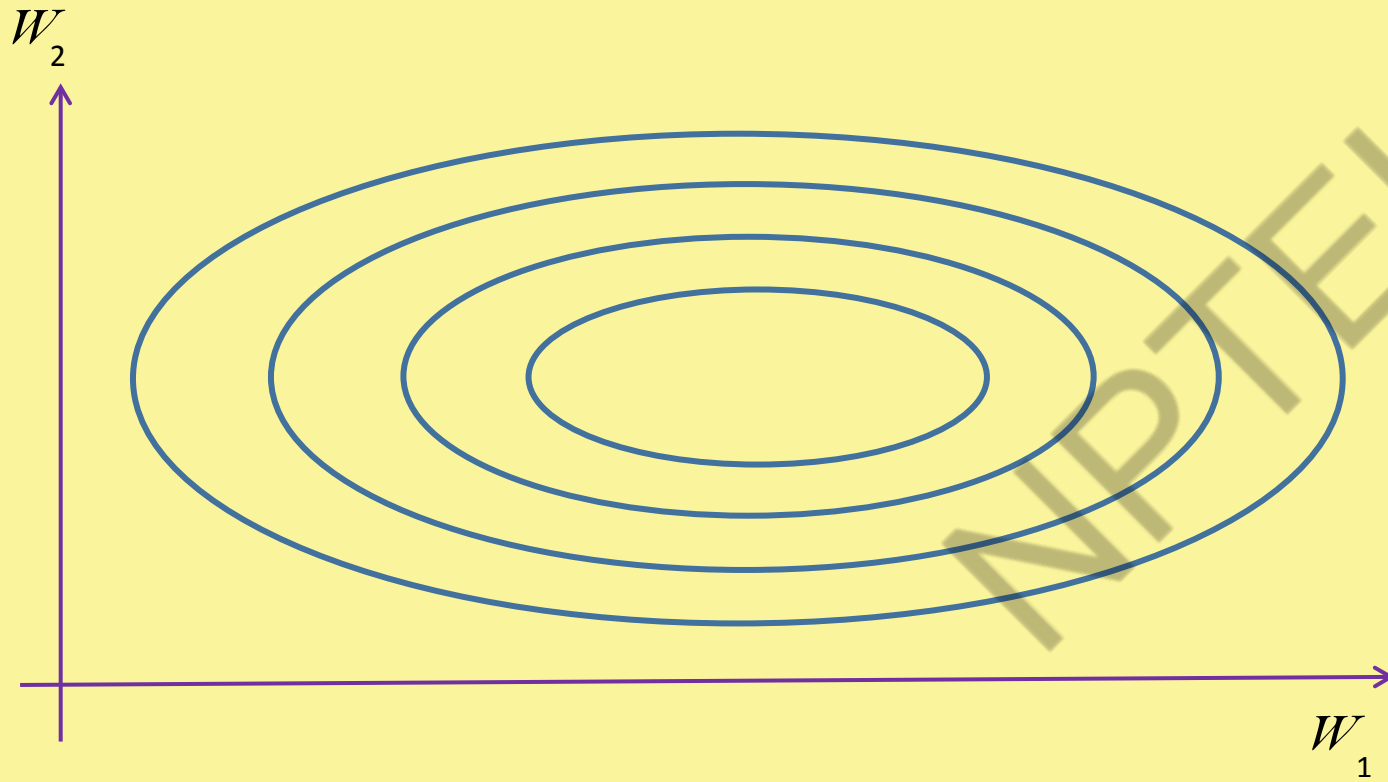
- ☐ Adagrad

- ☐ RMSProp

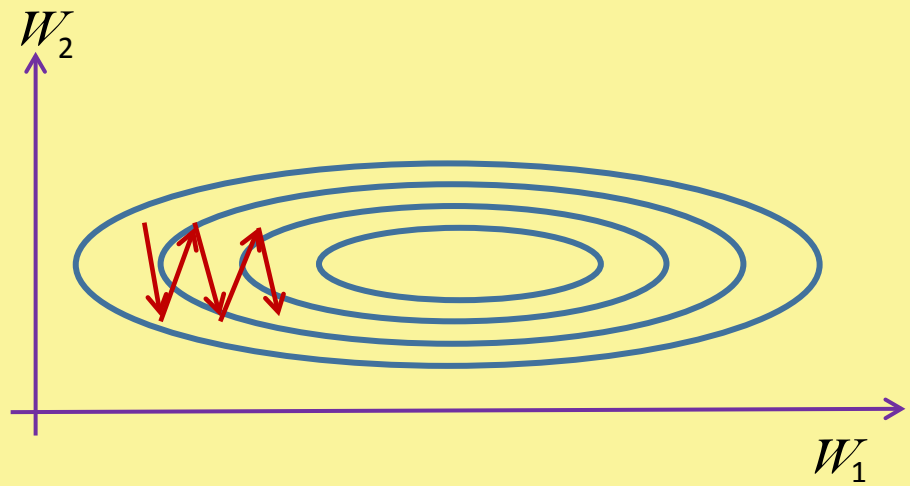
- ☐ etc.



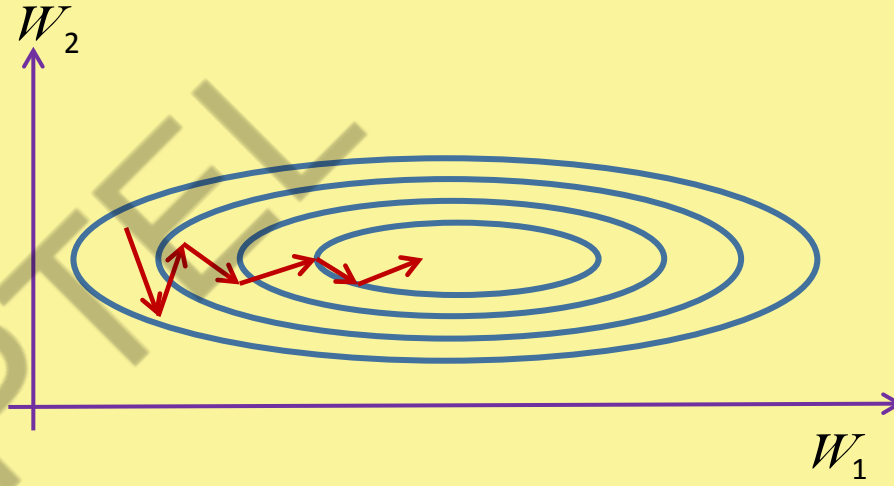
Momentum Optimizer



Momentum Optimizer



SGD



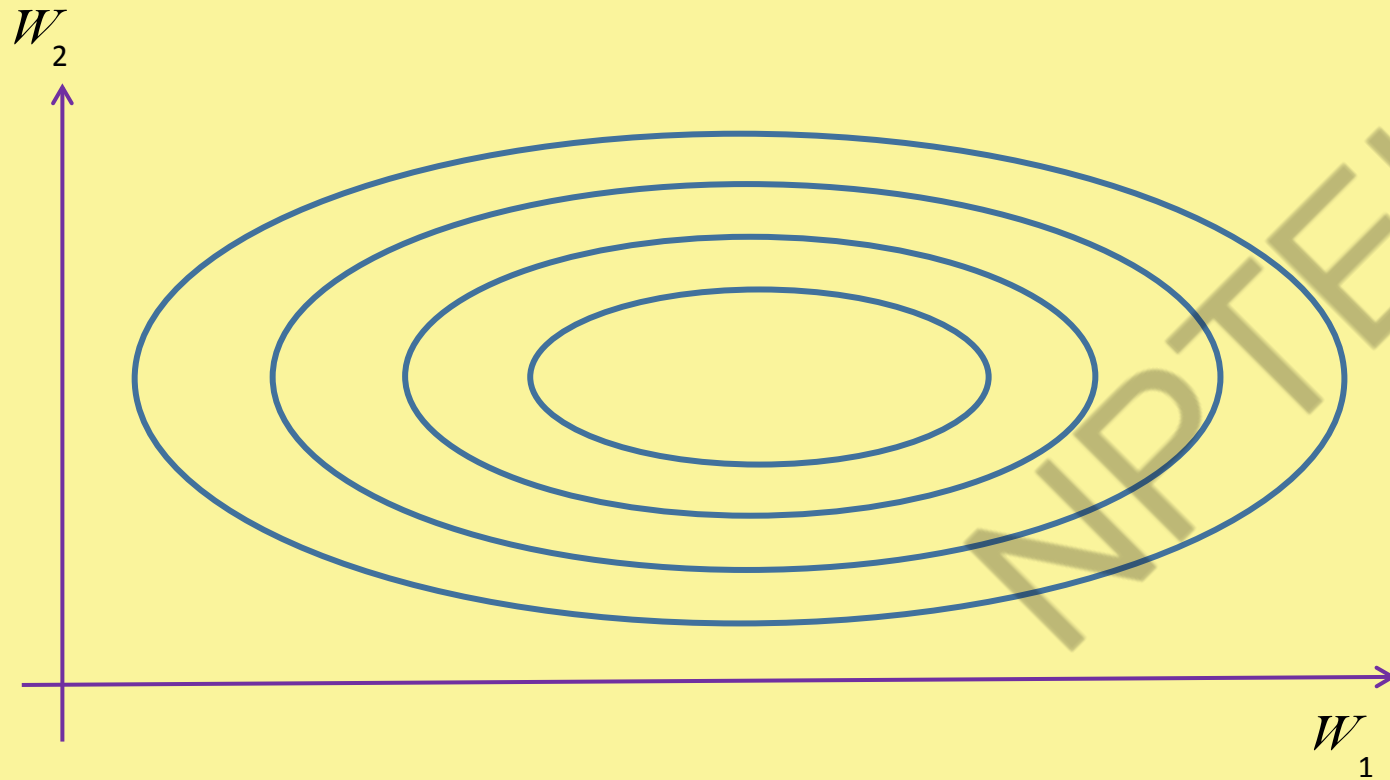
SGD with Momentum



Nesterov Accelerated Gradient (NAG)



Nesterov Accelerated Gradient (NAG)



Problem with Momentum Optimizer/NAG

- ❑ Both the algorithms require the hyper-parameters to be set manually.
- ❑ These hyper-parameters decide the learning rate.
- ❑ The algorithm uses same learning rate for all dimensions.
- ❑ The high dimensional (mostly) non-convex nature of loss function may lead to different sensitivity on different dimension.
- ❑ We may require learning rate be small in some dimension and large in another dimension.



Adagrad



Adagrad

- ☐ Adagrad adaptively scales the learning rate for different dimensions.
- ☐ Scale factor of a parameter is inversely proportional to the square root of sum of historical squared values of the gradient.
- ☐ The parameters with the largest partial derivative of the loss will have rapid decrease in their learning rate.
- ☐ Parameters with small partial derivatives will have relatively small decrease in learning rate.



Adagrad

$$g_t = \frac{1}{n} \sum_{\forall X \in \text{Minibatch}} \nabla_W L(W_t, X) \quad r_t = \sum_{\tau=1}^t g_\tau \circ g_\tau$$

$$W_{t+1} = W_t - \frac{\eta}{\sqrt{\epsilon I + r_t}} \circ g_t$$

$\circ \rightarrow$ element - wise product



Adagrad

$$\begin{bmatrix} W_{t+1}^{(1)} \\ W_{t+1}^{(2)} \\ \vdots \\ W_{t+1}^{(d)} \end{bmatrix} = \begin{bmatrix} W_t^{(1)} \\ W_t^{(2)} \\ \vdots \\ W_t^{(d)} \end{bmatrix} - \begin{bmatrix} \frac{\eta}{\sqrt{\epsilon + r_t^{(1)}}} \cdot g_t^{(1)} \\ \frac{\eta}{\sqrt{\epsilon + r_t^{(2)}}} \cdot g_t^{(2)} \\ \vdots \\ \frac{\eta}{\sqrt{\epsilon + r_t^{(d)}}} \cdot g_t^{(d)} \end{bmatrix}$$



Adagrad

Positive Side:

- ❑ Adagrad adaptively scales the learning rate for different dimensions by normalizing with respect to the gradient magnitude in the corresponding dimension.
- ❑ Adagrad eliminates the need to manually tune the learning rate.
- ❑ Reduces learning rate faster for parameters showing large slope and slower for parameters giving smaller slope.
- ❑ Adagrad converges rapidly when applied to convex functions.



Adagrad

Negative side:

- ☐ If the function is non-convex:- trajectory may pass through many complex terrains eventually arriving at a locally region.
- ☐ By then learning rate may become too small due to the accumulation of gradients from the beginning of training.
- ☐ So at some point the model may stop learning.





NPTEL ONLINE CERTIFICATION COURSES

*Thank
you*





NPTEL ONLINE CERTIFICATION COURSES

Course Name: Deep Learning

Faculty Name: Prof. P. K. Biswas

Department : E & ECE, IIT Kharagpur

Topic

Lecture 45: Optimizing Gradient Descent III

CONCEPTS COVERED

Concepts Covered:

☐ CNN

☐ Gradient Descent Challenges

☐ Momentum Optimizer

☐ Nesterov Accelerated Gradient

☐ Adagrad

☐ RMSProp

☐ etc.



Adagrad

$$g_t = \frac{1}{n} \sum_{\forall X \in \text{Minibatch}} \nabla_W L(W_t, X) \quad r_t = \sum_{\tau=1}^t g_\tau \circ g_\tau$$

$$W_{t+1} = W_t - \frac{\eta}{\sqrt{\epsilon I + r_t}} \circ g_t$$

$\circ \rightarrow$ element - wise product



Adagrad

Positive Side:

- ❑ Adagrad adaptively scales the learning rate for different dimensions by normalizing with respect to the gradient magnitude in the corresponding dimension.
- ❑ Adagrad eliminates the need to manually tune the learning rate.
- ❑ Reduces learning rate faster for parameters showing large slope and slower for parameters giving smaller slope.
- ❑ Adagrad converges rapidly when applied to convex functions.



Adagrad

Negative side:

- ☐ If the function is non-convex:- trajectory may pass through many complex terrains eventually arriving at a locally region.
- ☐ By then learning rate may become too small due to the accumulation of gradients from the beginning of training.
- ☐ So at some point the model may stop learning.



RMSProp



RMSPro

p

- ❑ RMSProp uses exponentially decaying average of squared gradient and discards history from the extreme past.
- ❑ Converges rapidly once it finds a locally convex bowl.
- ❑ Treats this as an instance of Adagrad algorithm initialized within that bowl.



RMSPro

$$g_t = \frac{1}{n} \sum_{\forall X \in \text{Minibatch}} \nabla_W L(W_t, X)$$

$$r_t = \beta r_{t-1} + (1 - \beta) g_t \circ g_t \rightarrow \text{Exponentially decaying average}$$

$$W_{t+1} = W_t - \frac{\eta}{\sqrt{\epsilon I + r_t}} \circ g_t$$



RMSProp with Nesterov Momentum

$$\tilde{W} = W_t + \alpha v \quad g_t = \frac{1}{n} \sum_{\forall X \in \text{Minibatch}} \nabla_W L(\tilde{W}, X)$$

$$r_t = \beta r_{t-1} + (1 - \beta) g_t \circ g_t$$

$$v_{t+1} = \alpha v_t - \frac{\eta}{\sqrt{\epsilon I + r_t}} \circ g_t \quad W_{t+1} = W_t + v_t$$



Adaptive Moments (Adam)



Adam

- ❑ Variant of the combination of RMSProp and Momentum.
- ❑ Incorporates first order moment (with exponential weighting) of the gradient (Momentum term).
- ❑ Momentum is incorporated in RMSProp by adding momentum to the rescaled gradients.
- ❑ Both first and second moments are corrected for bias to account for their initialization to zero.



Adam

$$g_t = \frac{1}{n} \sum_{\forall X \in \text{Minibatch}} \nabla_w L(W, X)$$

Biased first and second moments

$$s_t = \beta_1 s_{t-1} + (1 - \beta_1) g_t$$

$$r_t = \beta_2 r_{t-1} + (1 - \beta_2) g_t \circ g_t$$



Adam

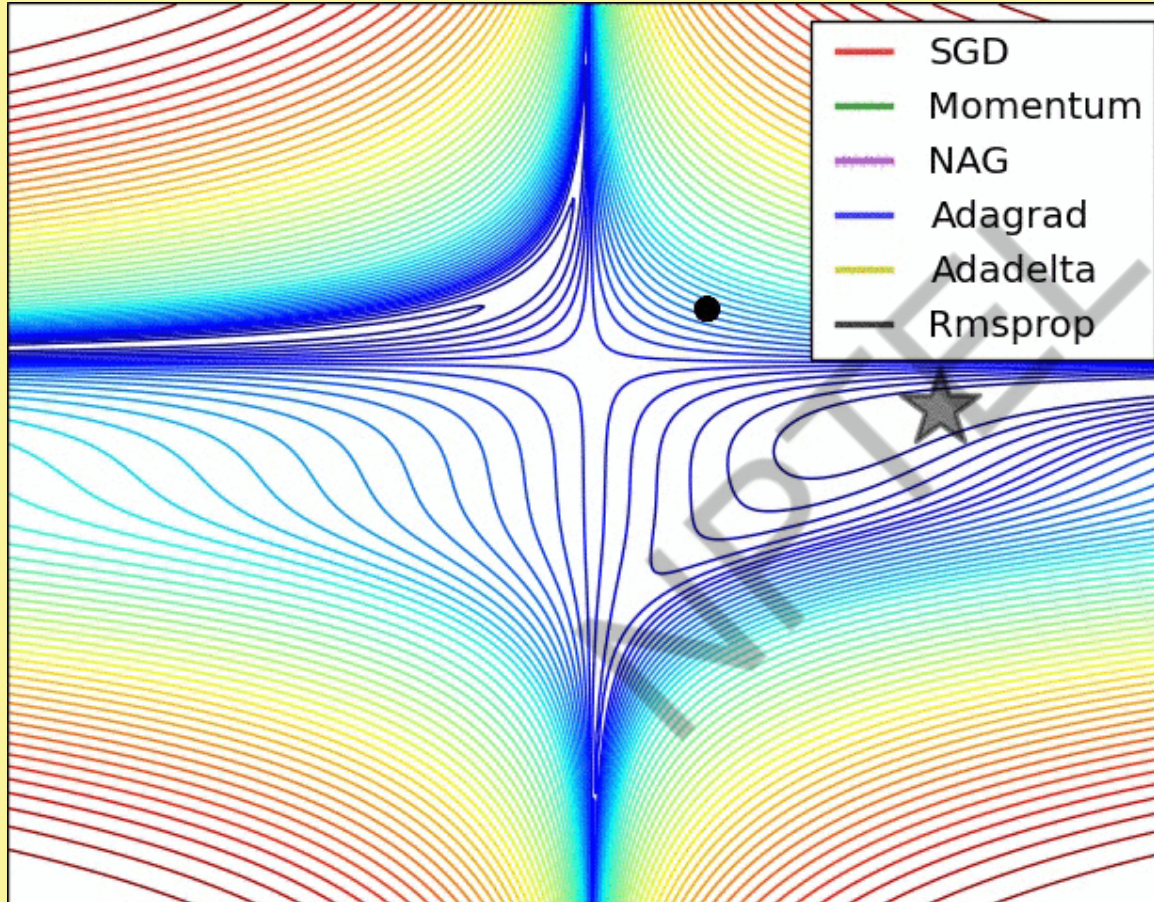
Bias corrected first and second moments

$$\hat{s}_t = \frac{s_t}{1 - \beta_1} \quad \hat{r}_t = \frac{r_t}{1 - \beta_2}$$

$$W_{t+1} = W_t - \eta \frac{\hat{s}_t}{\sqrt{\epsilon I + \hat{r}_t}}$$



Momentum Optimizer



Animation Source:-
<https://imgur.com/a/Hqolp>



NPTEL ONLINE CERTIFICATION COURSES

*Thank
you*

