

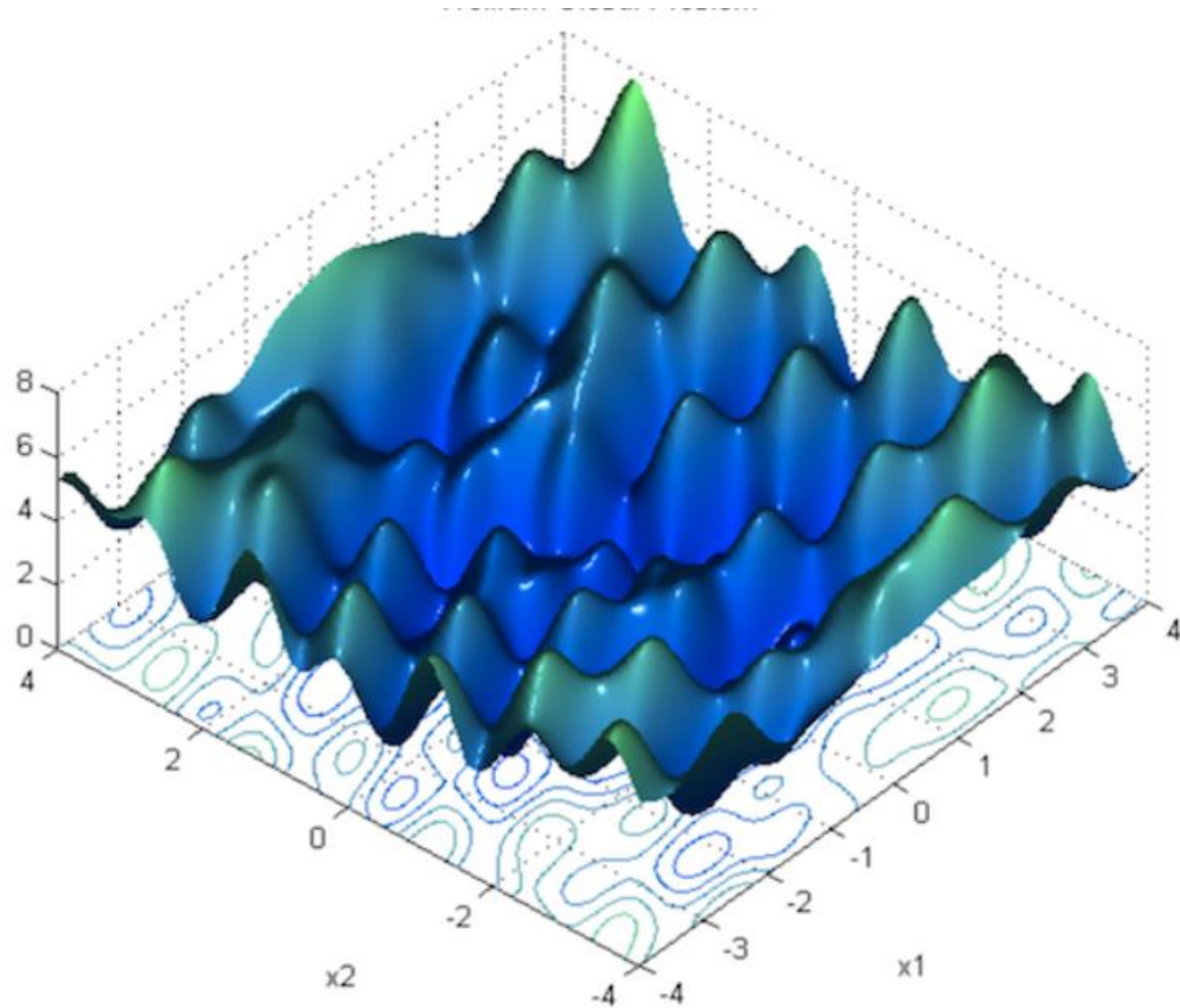
OPTIMIZATION FOR DNN

- **Difficulty in finding best parameters' values**
- **SGD**
- **SGD + Momentum**
- **Adaptive learning rates (Adaptive methods)**
- **Adam (Momentum + Adaptive rates)**
- **Second order methods**
- **More on gradient vanishing and exploding**
- **Training of LSTM**

- Combine materials from Fei-Fei Lee ppt, Chap. 8 of Deep Learning (Goodfellow, etc.) and [how neural networks are trained](#)

Difficulty in finding optimal parameter set in DNN

- Too many parameters, hundreds of thousand to million
- Curse of dimensionality— random search or grid search are impossible
- The convexity of the loss function is very complex, many local optimums, saddle points—for d parameters, when gradients are all zero, there is only $\frac{1}{3^d}$ possibility to be an optimum, and it could be a local minimum, or a local maximum, not the global minimum
- You may find a local minimum



Example of non-convex loss surface with two parameters. Note that in deep neural networks, we're dealing with millions of parameters, but the basic principle stays the same. Source: [Yoshua Bengio](#).

SGD (Stochastic Gradient Descent)

- Mini-batch gradient descent (MB-GD), in which the whole dataset is randomly subdivided into N equally-sized mini-batches of K samples each. K may be a small positive number, or it can be in the dozens or hundreds; it depends on the specific architecture and application. Note that if $K=1$, then you have SGD, and if K is the size of the whole dataset, it is **batch** gradient descent. Note also that confusingly, sometimes people say “SGD” to refer to both MB-GD and **one sample at a time**.

Mini-batch SGD

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

Require: Learning rate schedule $\epsilon_1, \epsilon_2, \dots$

Require: Initial parameter θ

$k \leftarrow 1$

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: $\theta \leftarrow \theta - \epsilon_k \hat{\mathbf{g}}$

$k \leftarrow k + 1$

end while

Problem with SGD

- What if loss changes quickly in one direction and slowly in another?
- What does gradient descent do?

```
# Vanilla Gradient Descent
```

```
while True:
```

```
    weights_grad = evaluate_gradient(loss_fun, data, weights)
```

```
    weights += - step_size * weights_grad # perform parameter update
```

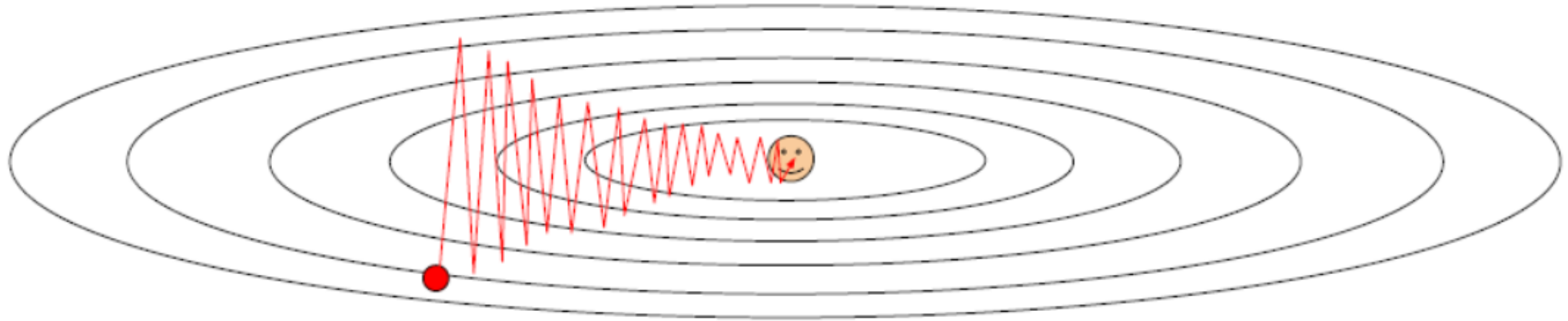
Lee's

Optimization: Problems with SGD

What if loss changes quickly in one direction and slowly in another?

What does gradient descent do?

Very slow progress along shallow dimension, jitter along steep direction



Lee's

Optimization: Problems with SGD

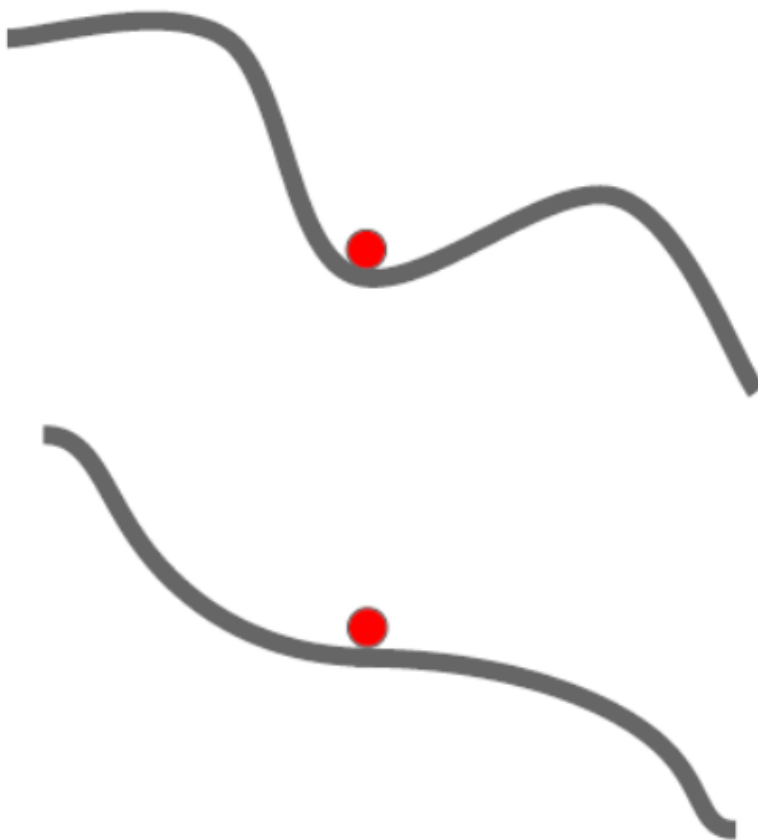
What if the loss function has a **local minima** or **saddle point**?



Optimization: Problems with SGD

What if the loss function has a **local minima** or **saddle point**?

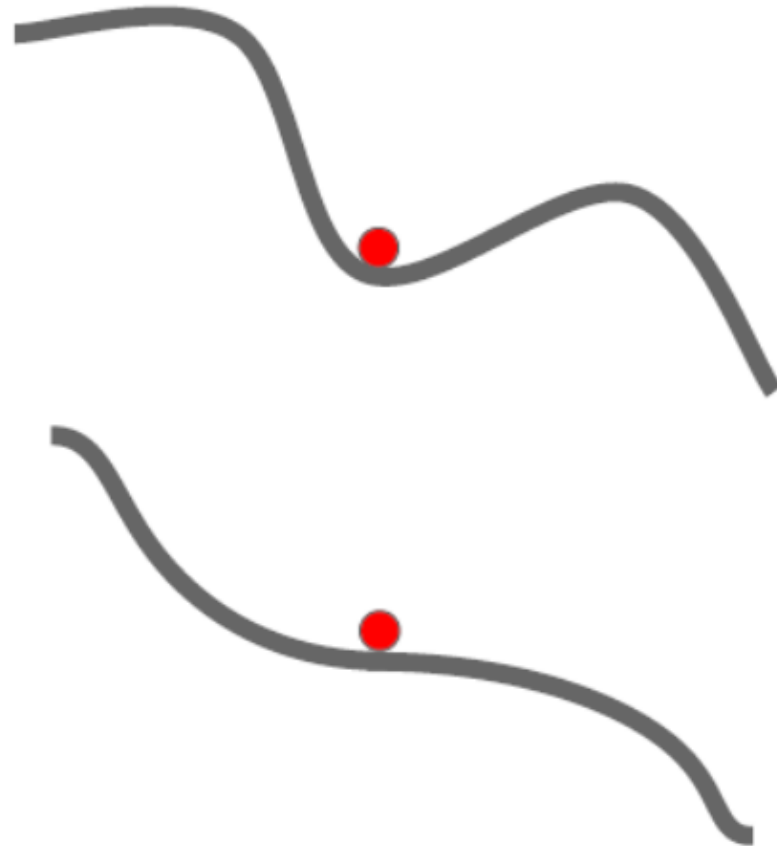
Zero gradient,
gradient descent
gets stuck



Optimization: Problems with SGD

What if the loss function has a **local minima** or **saddle point**?

Saddle points much more common in high dimension



Solution 1--add momentum

SGD + Momentum

SGD

$$x_{t+1} = x_t - \alpha \nabla f(x_t)$$

```
while True:
    dx = compute_gradient(x)
    x -= learning_rate * dx
```

SGD+Momentum

$$v_{t+1} = \rho v_t + \nabla f(x_t)$$

$$x_{t+1} = x_t - \alpha v_{t+1}$$

```
vx = 0
while True:
    dx = compute_gradient(x)
    vx = rho * vx + dx
    x -= learning_rate * vx
```

- Build up “velocity” as a running mean of gradients
- Rho gives “friction”; typically rho=0.9 or 0.99

SGD + Momentum

SGD+Momentum

$$v_{t+1} = \rho v_t - \alpha \nabla f(x_t)$$

$$x_{t+1} = x_t + v_{t+1}$$

```
vx = 0
while True:
    dx = compute_gradient(x)
    vx = rho * vx - learning_rate * dx
    x += vx
```

SGD+Momentum

$$v_{t+1} = \rho v_t + \nabla f(x_t)$$

$$x_{t+1} = x_t - \alpha v_{t+1}$$

```
vx = 0
while True:
    dx = compute_gradient(x)
    vx = rho * vx + dx
    x -= learning_rate * vx
```

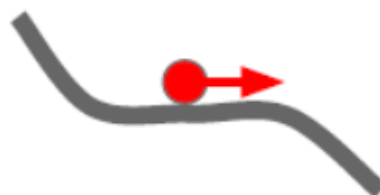
You may see SGD+Momentum formulated different ways,
but they are equivalent - give same sequence of x

SGD + Momentum

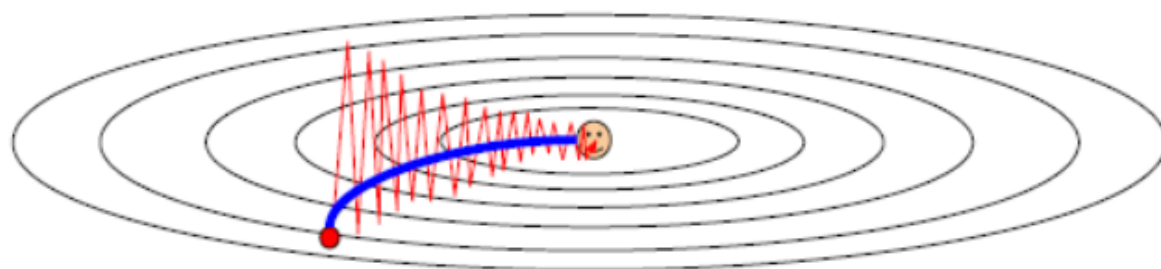
Local Minima



Saddle points



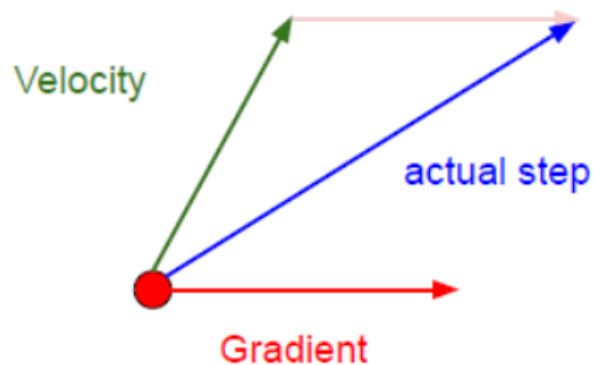
Poor Conditioning



Nesterov's momentum accelerated gradient(NAG)

Nesterov Momentum

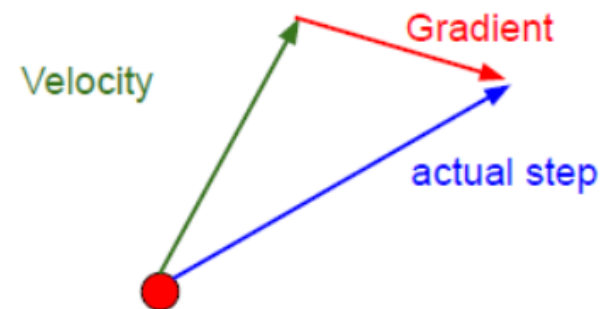
Momentum update:



Combine gradient at current point with velocity to get step used to update weights

Nesterov, "A method of solving a convex programming problem with convergence rate $O(1/k^2)$ ", 1983
Mertens, "Introduction to convex optimization: a basic course", 2024

Nesterov Momentum



"Look ahead" to the point where updating using velocity would take us; compute gradient there and mix it with velocity to get actual update direction

Use mini-batch

$$\boldsymbol{v} \leftarrow \alpha \boldsymbol{v} - \epsilon \nabla_{\boldsymbol{\theta}} \left[\frac{1}{m} \sum_{i=1}^m L \left(\boldsymbol{f}(\boldsymbol{x}^{(i)}; \boldsymbol{\theta} + \alpha \boldsymbol{v}), \boldsymbol{y}^{(i)} \right) \right], \quad (8.21)$$

$$\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} + \boldsymbol{v}, \quad (8.22)$$

Algorithm 8.3 Stochastic gradient descent (SGD) with Nesterov momentum

Require: Learning rate ϵ , momentum parameter α

Require: Initial parameter θ , initial velocity v

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 labels $\mathbf{y}^{(i)}$.

 Apply interim update: $\tilde{\theta} \leftarrow \theta + \alpha v$.

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

 Compute velocity update: $\mathbf{v} \leftarrow \alpha \mathbf{v} - \epsilon \mathbf{g}$.

 Apply update: $\theta \leftarrow \theta + v$.

end while

Adaptive learning rate--AdaGrad

- AdaGrad stands for **Adaptive subGradient**. [how neural networks are trained](#)
- Adaptive methods tend to scale gradient components differently. It will decrease the component which are always big.

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 $\mathbf{r} = \mathbf{0}$

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: $\mathbf{r} \leftarrow \mathbf{r} + \mathbf{g} \odot \mathbf{g}$.

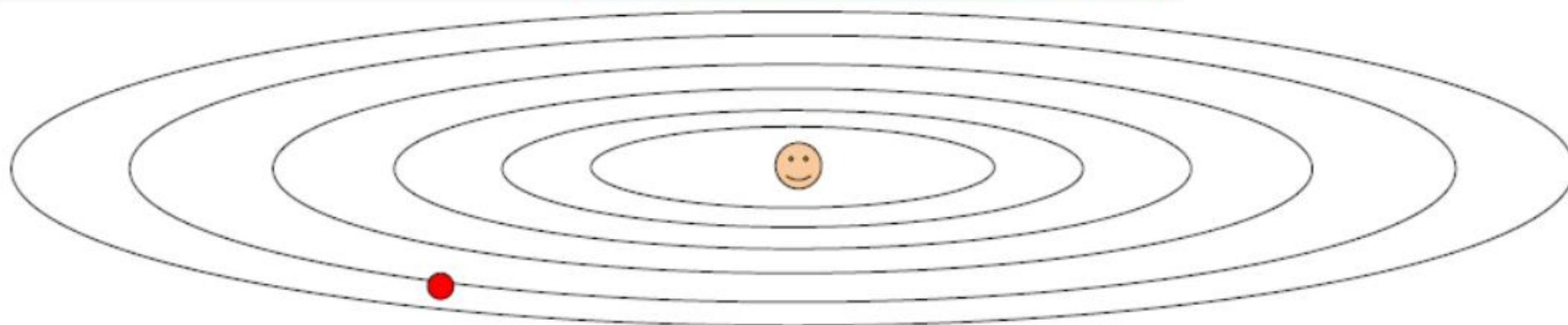
 Compute update: $\Delta \theta \leftarrow \frac{\epsilon}{\delta + \sqrt{\mathbf{r}}} \odot \mathbf{g}$. (Division and square root applied element-wise)

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

end while

AdaGrad

```
grad_squared = 0
while True:
    dx = compute_gradient(x)
    grad_squared += dx * dx
    x -= learning_rate * dx / (np.sqrt(grad_squared) + 1e-7)
```



Q2: What happens to the step size over long time? Decays to zero

RMSProp: “Leaky AdaGrad”

AdaGrad

```
grad_squared = 0
while True:
    dx = compute_gradient(x)
    grad_squared += dx * dx
    x -= learning_rate * dx / (np.sqrt(grad_squared) + 1e-7)
```



RMSProp

```
grad_squared = 0
while True:
    dx = compute_gradient(x)
    grad_squared = decay_rate * grad_squared + (1 - decay_rate) * dx * dx
    x -= learning_rate * dx / (np.sqrt(grad_squared) + 1e-7)
```

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 $\mathbf{r} = 0$

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

ρ controls the length scale of the moving window. RMSProp is effective and practical.

Finally, the Adam-consider momentum and adaptive rate scaling

- Adam stands for “Adaptive Moment”
- First moment for momentum (+current gradient)
- Second moment for weight scaling

Adam

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 $\mathbf{s} = \mathbf{0}$, $\mathbf{r} = \mathbf{0}$

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: $\mathbf{s} \leftarrow \rho_1 \mathbf{s} + (1 - \rho_1) \mathbf{g}$

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

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

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

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

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

end while

Reference: Kingma, Diederik and Jimmy Ba, “Adam: A method for Stochastic Optimization

$$v_i = \beta_2 v_{i-1} + (1 - \beta_2) g_i \odot g_i$$

$$v_i = \beta_2^{i-1} (1 - \beta_2) g_1 \odot g_1 + \beta_2^{i-2} (1 - \beta_2) g_2 \odot g_2 + \cdots + (1 - \beta_2) g_i \odot g_i$$

$$v_i = (1 - \beta_2) \sum_{k=1}^i \beta_2^{i-k} g_k \odot g_k$$

$$E(v_i) = E[(1 - \beta_2) \sum_{k=1}^i \beta_2^{i-k} (g_k \odot g_k)]$$

Assuming that, $E[g_k \odot g_k] = E[g_i \odot g_i]$

We have,

Note that, $1 - x^n = (1 - x)(1 + x + x^2 + \dots + x^{n-1})$

$$\begin{aligned} E(v_i) &= E(1 - \beta_2) \sum_{k=1}^i \beta^{i-k} g_k \odot g_k = E[g_i \odot g_i] (1 - \beta_2) \sum_{k=1}^i \beta^{i-k} \\ &= E[g_i \odot g_i] (1 - \beta_2^i) \end{aligned}$$

Therefore,

$$\tilde{v}_i = \frac{v_i}{(1 - \beta_2^i)}$$

So which one is the best?

- Unfortunately, there is currently no consensus on this point. Schaul et al. suggested that the family of algorithms with adaptive learning rates (RMSProp and AdaDelta, and ADAM) perform fairly robustly.
- No single best algorithm has emerged.
- Currently, the most popular optimization algorithms actively in use include SGD, SGD with momentum, RMSProp, RMSProp with momentum, AdaDelta, and Adam.
- The choice depend one the user's familiarity with the algorithm(for ease of hyperparameter tuning)

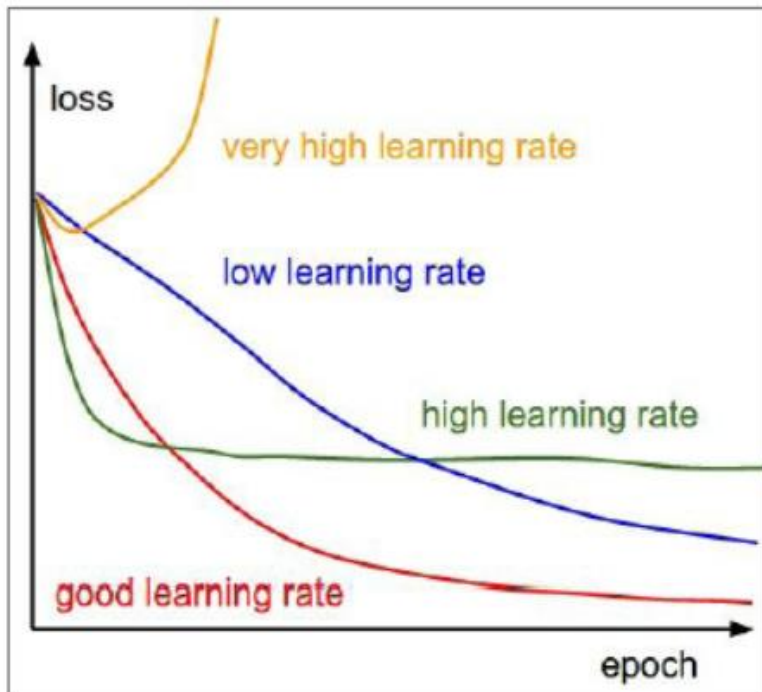
Lee's

In practice:

- **Adam** is a good default choice in many cases; it often works ok even with constant learning rate
- **SGD+Momentum** can outperform Adam but may require more tuning of LR and schedule
 - Try cosine schedule, very few hyperparameters!
- If you can afford to do full batch updates then try out **L-BFGS** (and don't forget to disable all sources of noise)

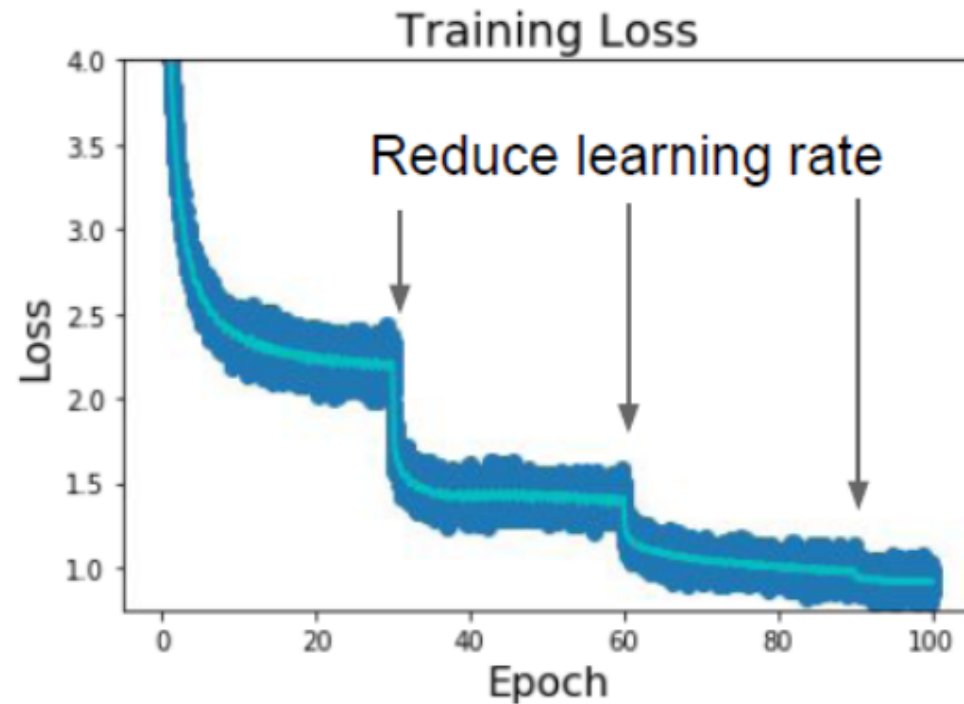
LR: learning rate tuning

SGD, SGD+Momentum, Adagrad, RMSProp, Adam all have **learning rate** as a hyperparameter.



Q: Which one of these learning rates is best to use?

Learning Rate Decay



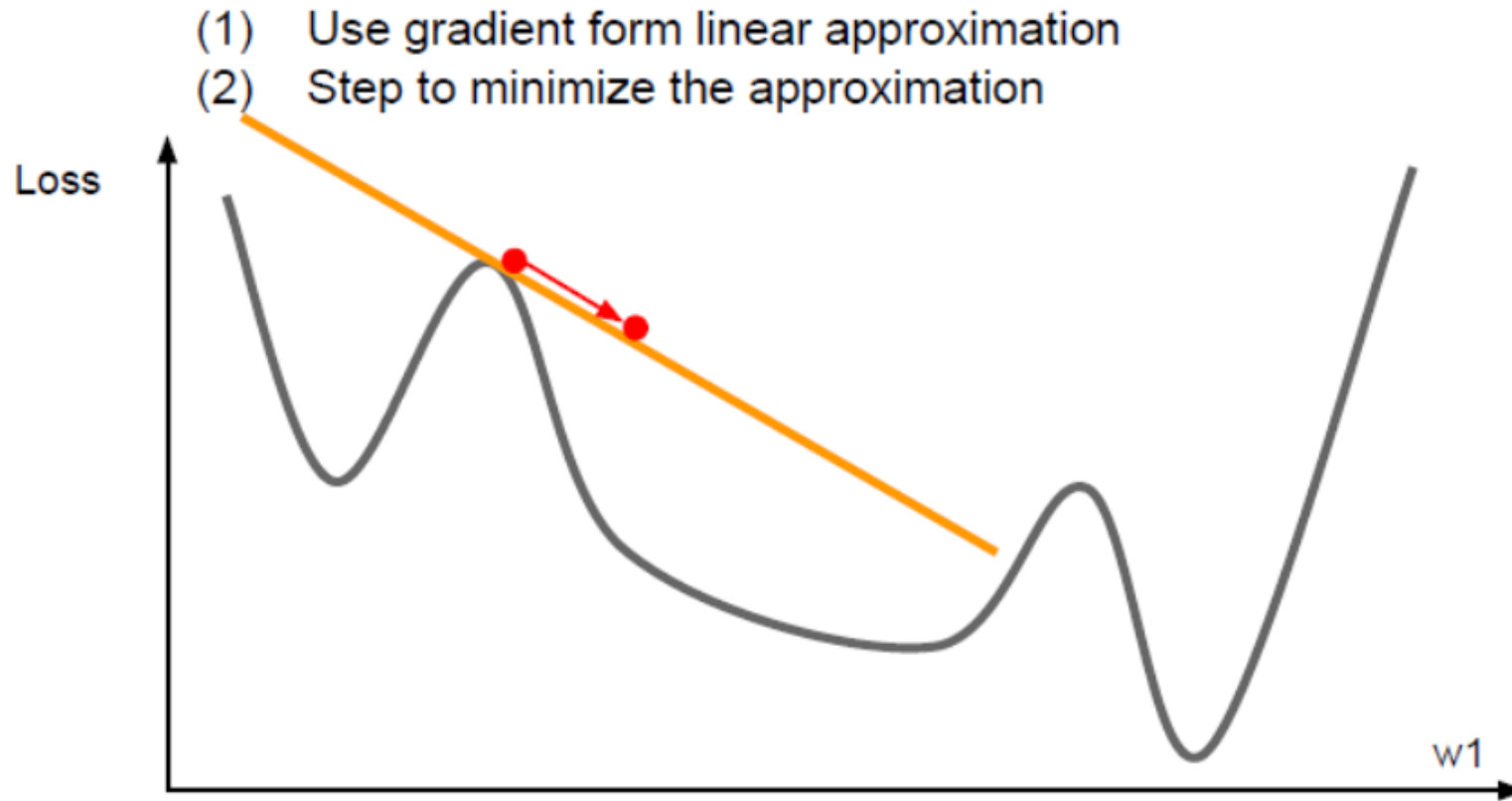
Step: Reduce learning rate at a few fixed points. E.g. for ResNets, multiply LR by 0.1 after epochs 30, 60, and 90.

See the animation!!

- <http://colah.github.io/posts/2015-08-Understanding-LSTMs/>
- <https://stackoverflow.com/questions/44273249/in-keras-what-exactly-am-i-configuring-when-i-create-a-stateful-lstm-layer-wi>
- <http://www.deeplearningbook.org/>

Lee's

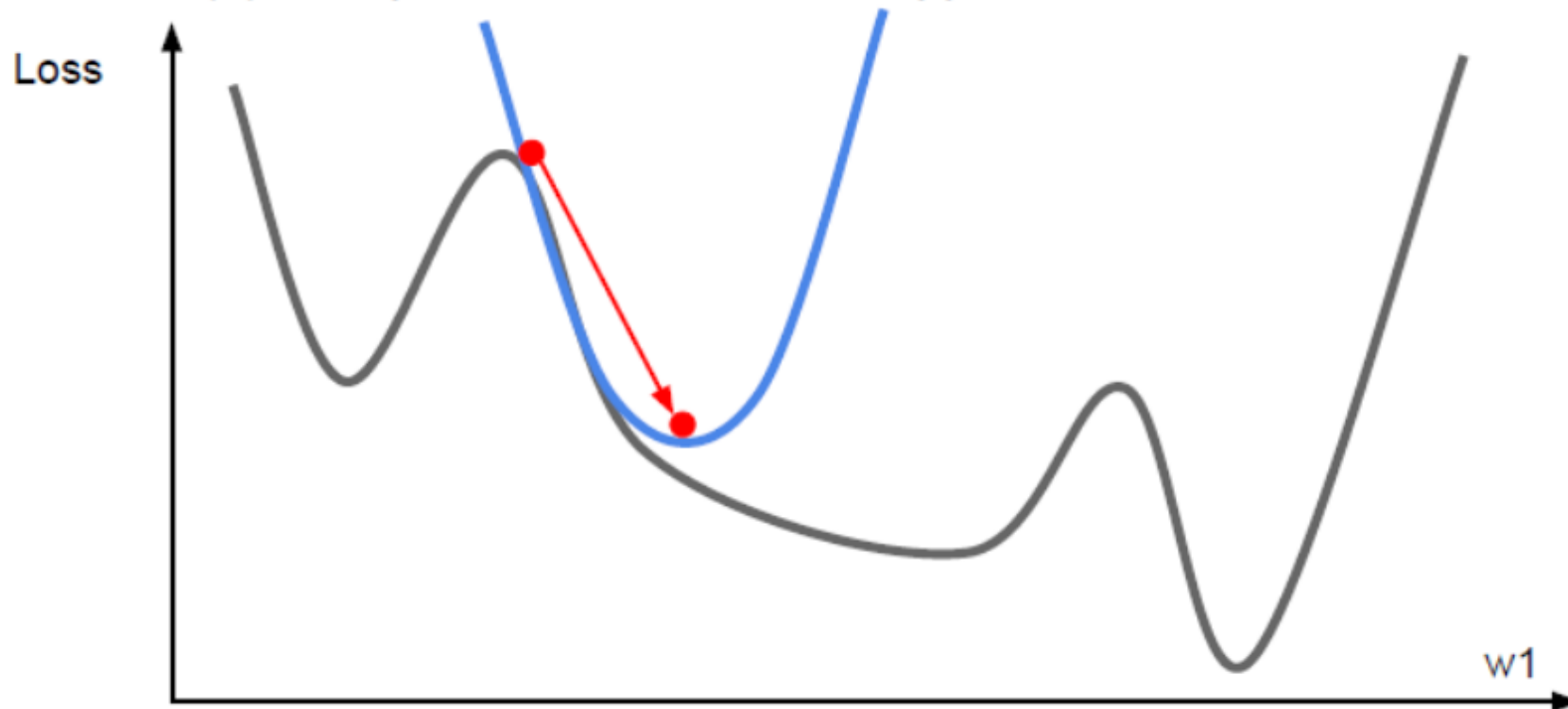
First-Order Optimization



Lee's

Second-Order Optimization

- (1) Use gradient and **Hessian** to form **quadratic** approximation
- (2) Step to the **minima** of the approximation



Note

- The first order method considers the tangent of a target point, while the second order method considers the curvature of the function near the target point

Second order methods

- Single variable Taylor's Expansion
- $f(x) \cong f(x_0) + (x - x_0)f'(x_0) + \frac{(x-x_0)^2}{2}f''(x_0) + \dots$
- Take derivative of the equation and let it equals 0.
- $f'(x) = f'(x_0) + (x - x_0)f''(x_0) = 0$
- $x = x_0 - \frac{f'(x_0)}{f''(x_0)}$ Newton's method for optimization

- Algorithm Newton's method

1: let x_0 be the initial point

2: **while** $|f'(x_0)| > \epsilon$ **do**

3: $x = x_0 - f'(x_0)/f''(x_0)$

4: $x_0 = x$

5: **end while**

6: **return** x

Multivariate Taylor's expansion

$$J(\theta) \cong J(\theta_0) + (\theta - \theta_0)^t \nabla_{\theta} J(\theta_0) + \frac{1}{2} (\theta - \theta_0)^t H (\theta - \theta_0)$$

Take derivative and equate to 0, we can find

$$\theta^* = \theta_0 - H^{-1} \nabla_{\theta} J(\theta_0)$$

Where d-dimensional Hessian matrix:

$$H(x) = \begin{bmatrix} \frac{\partial^2 f}{\partial_{x1} \partial_{x1}} & \frac{\partial^2 f}{\partial_{x1} \partial_{x2}} & \cdots & \frac{\partial^2 f}{\partial_{x1} \partial_{xd}} \\ \vdots & \vdots & & \vdots \\ \frac{\partial^2 f}{\partial_{xd} \partial_{x1}} & \frac{\partial^2 f}{\partial_{xd} \partial_{x2}} & \cdots & \frac{\partial^2 f}{\partial_{xd} \partial_{xd}} \end{bmatrix}$$

From Deep learning book

Algorithm 8.8 Newton's method with objective $J(\boldsymbol{\theta}) = \frac{1}{m} \sum_{i=1}^m L(f(\boldsymbol{x}^{(i)}; \boldsymbol{\theta}), y^{(i)})$

Require: Initial parameter $\boldsymbol{\theta}_0$

Require: Training set of m examples

while stopping criterion not met **do**

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

 Compute Hessian: $\boldsymbol{H} \leftarrow \frac{1}{m} \nabla_{\boldsymbol{\theta}}^2 \sum_i L(f(\boldsymbol{x}^{(i)}; \boldsymbol{\theta}), y^{(i)})$

 Compute Hessian inverse: \boldsymbol{H}^{-1}

 Compute update: $\Delta\boldsymbol{\theta} = -\boldsymbol{H}^{-1} \boldsymbol{g}$

 Apply update: $\boldsymbol{\theta} = \boldsymbol{\theta} + \Delta\boldsymbol{\theta}$

end while

Lee's

Second-Order Optimization

second-order Taylor expansion:

$$J(\boldsymbol{\theta}) \approx J(\boldsymbol{\theta}_0) + (\boldsymbol{\theta} - \boldsymbol{\theta}_0)^\top \nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta}_0) + \frac{1}{2} (\boldsymbol{\theta} - \boldsymbol{\theta}_0)^\top \mathbf{H} (\boldsymbol{\theta} - \boldsymbol{\theta}_0)$$

Solving for the critical point we obtain the Newton parameter update:

$$\boldsymbol{\theta}^* = \boldsymbol{\theta}_0 - \mathbf{H}^{-1} \nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta}_0)$$

Hessian has $O(N^2)$ elements

Inverting takes $O(N^3)$

N = (Tens or Hundreds of) Millions

Q: Why is this bad for deep learning?

Lee's

Second-Order Optimization

$$\theta^* = \theta_0 - H^{-1} \nabla_{\theta} J(\theta_0)$$

- Quasi-Newton methods (**BGFS** most popular):
instead of inverting the Hessian ($O(n^3)$), approximate inverse Hessian with rank 1 updates over time ($O(n^2)$ each).
- **L-BFGS** (Limited memory BFGS):
Does not form/store the full inverse Hessian.

Lee's

L-BFGS

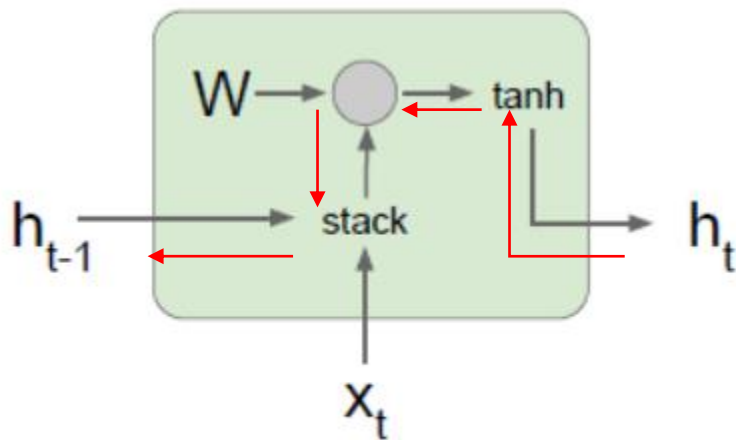
- **Usually works very well in full batch, deterministic mode** i.e. if you have a single, deterministic $f(x)$ then L-BFGS will probably work very nicely
- **Does not transfer very well to mini-batch setting.** Gives bad results. Adapting second-order methods to large-scale, stochastic setting is an active area of research.

Gradient vanishing & exploding in RNN

- Recap

Vanilla RNN Gradient Flow

Backpropagation from h_t
to h_{t-1} multiplies by W
(actually W_{hh}^T)

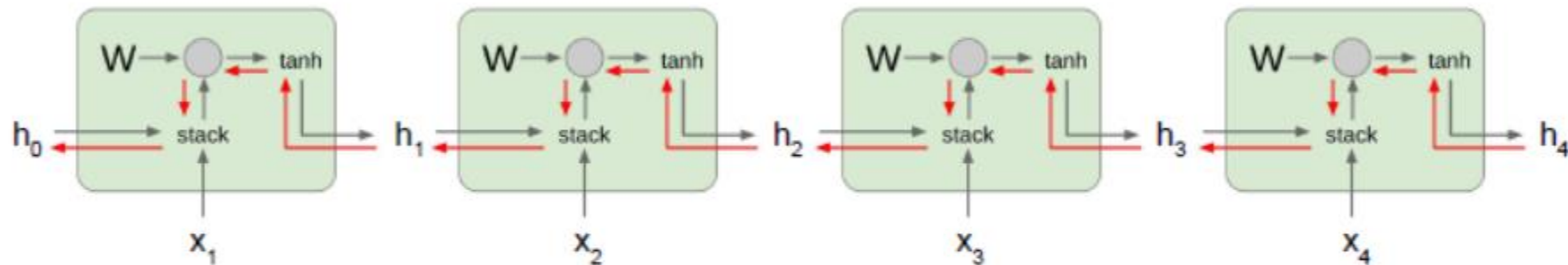


Bengio et al, "Learning long-term dependencies with gradient descent is difficult", IEEE Transactions on Neural Networks, 1994
Pascanu et al, "On the difficulty of training recurrent neural networks", ICML 2013

$$\begin{aligned} h_t &= \tanh(W_{hh}h_{t-1} + W_{hx}x_t) \\ &= \tanh\left((W_{hh} \quad W_{hx}) \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix}\right) \\ &= \tanh\left(W \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix}\right) \end{aligned}$$

Vanilla RNN Gradient Flow

Bengio et al, "Learning long-term dependencies with gradient descent is difficult", IEEE Transactions on Neural Networks, 1994
Pascanu et al, "On the difficulty of training recurrent neural networks", ICML 2013



Computing gradient
of h_0 involves many
factors of W
(and repeated \tanh)

Consider only W_{hh} , $h_t = W_{hh} \cdot h_{t-1}$

- Let ξ be the cost function: the gradient of m steps before the current step t is:
-

$$\frac{\partial \xi}{\partial h_{t-m}} = \frac{\partial \xi}{\partial h_t} * \frac{\partial h_t}{\partial h_{t-m}}$$

Where $\frac{\partial h_t}{\partial h_{t-m}} = \frac{\partial h_t}{\partial h_{t-1}} * \dots * \frac{\partial h_{t-m+1}}{\partial h_{t-m}} = W^m$

$$\frac{\partial \xi}{\partial h_{t-m}} = \frac{\partial \xi}{\partial h_t} * W^m$$

Deep learning book

- Suppose W has an eigendecomposition $W=V \text{diag}(\lambda)V^{-1}$
- $W^m = (V \text{diag}(\lambda)V^{-1})^m = V \text{diag}(\lambda)^m V^{-1}$
- Any eigenvalues λ_i that are not near an absolute value of 1 will either explode if they are greater than 1 in magnitude or vanish if they are less than 1 in magnitude.
- Vanishing gradients make it difficult to know which direction the parameters should move to improve the cost function
- Exploding gradients can make learning unstable

- Recurrent networks use **the same matrix W** at each time step, but **feedforward networks do not**, so even very deep feedforward networks can largely avoid the vanishing and exploding gradient problem.
- ---from deep learning book of Ian Goodfellow etc.

A diagonalization example

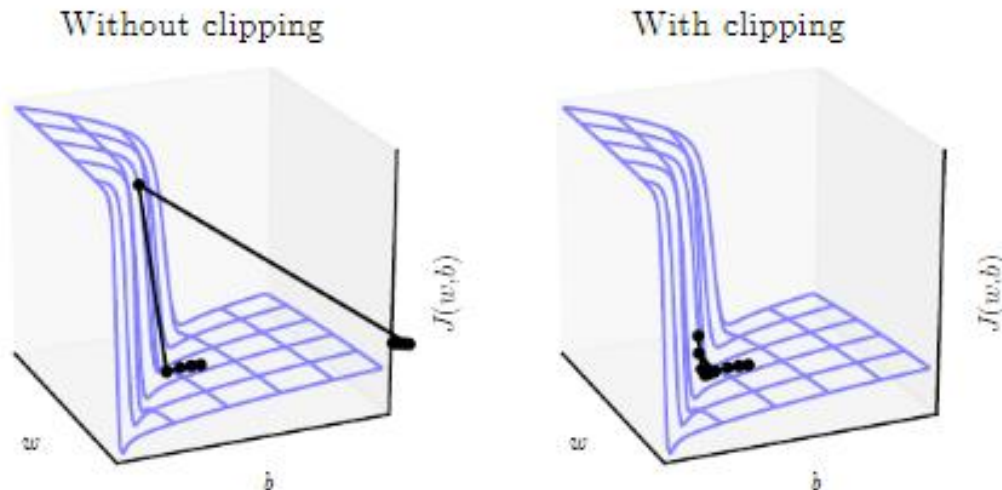
- $A = \begin{bmatrix} -4 & -5 \\ 10 & 11 \end{bmatrix}$ has eigen value $\lambda_1 = 1$, $x_1 = \begin{bmatrix} 1 \\ -1 \end{bmatrix}$,
 $\lambda_2 = 6$, $x_2 = \begin{bmatrix} -1 \\ 2 \end{bmatrix}$

$$\begin{aligned} A^k &= (S\Lambda S^{-1})^k = S\Lambda^k S^{-1} = \begin{bmatrix} 1 & -1 \\ -1 & 2 \end{bmatrix} \begin{bmatrix} 1^k & 0 \\ 0 & 6^k \end{bmatrix} \begin{bmatrix} 2 & 1 \\ 1 & 1 \end{bmatrix} \\ &= \begin{bmatrix} 2 - 6^k & 1 - 6^k \\ -2 + 2 \cdot 6^k & -1 + 2 \cdot 6^k \end{bmatrix} \end{aligned}$$

To deal with the gradient exploding problem

1. Gradient clipping, where we threshold the maximum value a gradient can get

$$\text{if } \|g\| > \beta, \quad g \leftarrow \frac{\beta}{\|g\|} g$$



Training LSTM

$$i_t = \sigma(W_{ii}x_t + b_{ii} + W_{hi}h_{(t-1)} + b_{hi})$$

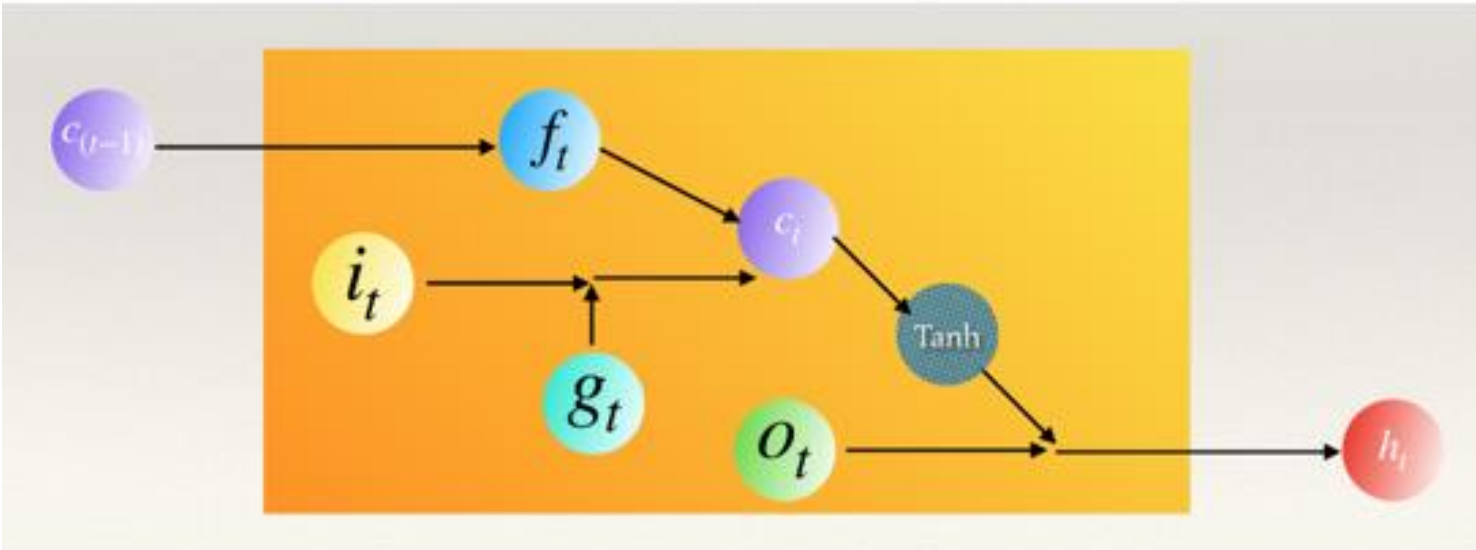
$$f_t = \sigma(W_{if}x_t + b_{if} + W_{hf}h_{(t-1)} + b_{hf})$$

$$g_t = \tanh(W_{ig}x_t + b_{ig} + W_{hg}h_{(t-1)} + b_{hg})$$

$$o_t = \sigma(W_{io}x_t + b_{io} + W_{ho}h_{(t-1)} + b_{ho})$$

$$c_t = f_t c_{(t-1)} + i_t g_t$$

$$h_t = o_t \tanh(c_t)$$

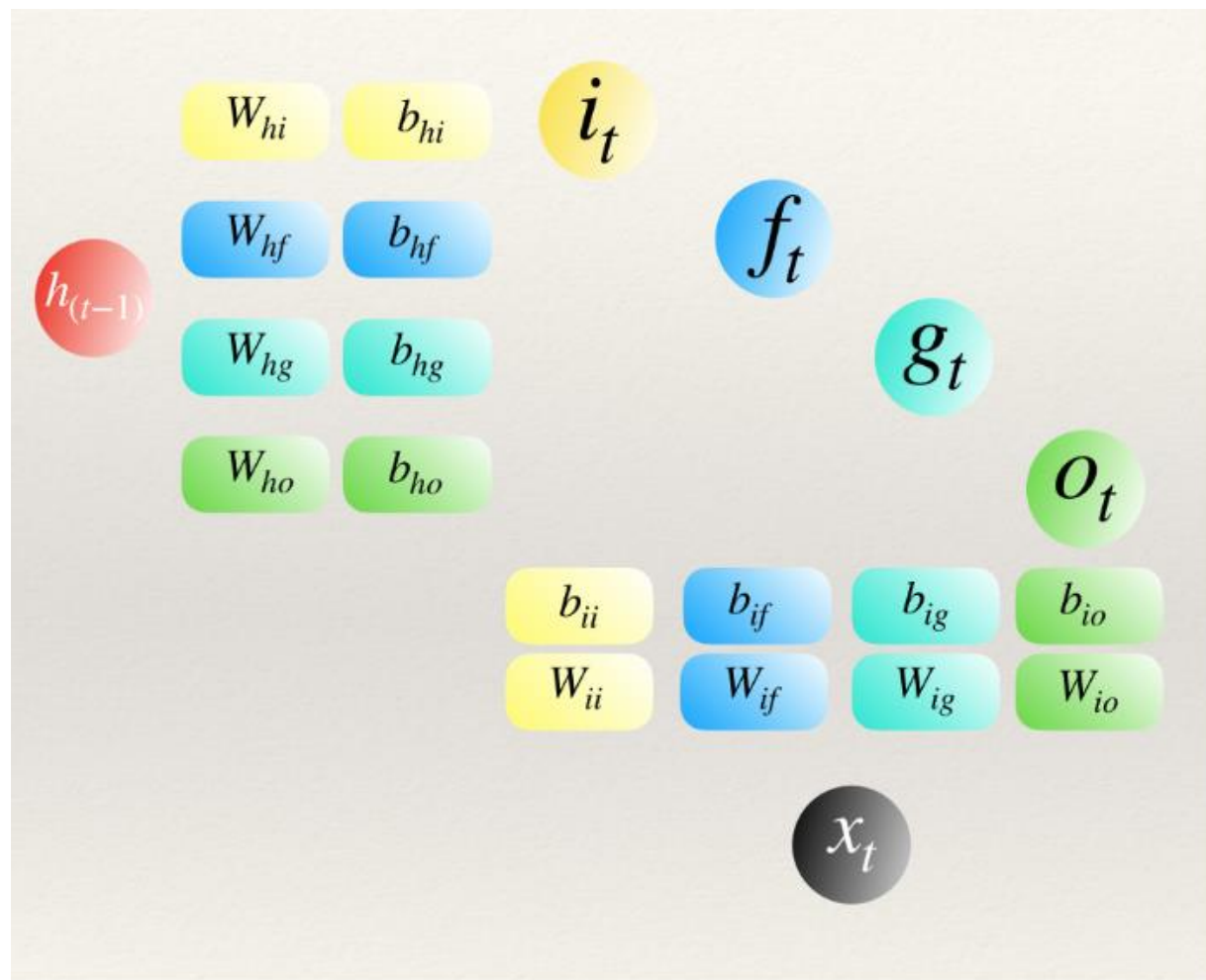


f: forget gate

i: input gate

g: gate gate (or cell gate)

o: output gate



To simplify, set x_t and h_t

$$x_t = [x_1, x_2, \dots, x_5]$$

輸入之向量

$$h_{(t-1)} = [h_1, h_2, \dots, h_{10}]$$

上一刻的輸出

$$W_{ii} = [W_{ii1}, \dots, W_{ii10}] = \begin{bmatrix} W_{ii1,1}, \dots, W_{ii10,1} \\ \dots \\ W_{ii1,5}, \dots, W_{ii10,5} \end{bmatrix}$$

W_{ii} is 5*10; W_{hi} is 10*10

$$\begin{aligned} W_{hi} &= [W_{hi1} \dots W_{hi10}] \\ &= \begin{bmatrix} W_{hi1,1} \dots W_{hi10,1} \\ W_{hi1,2} \dots W_{hi10,2} \\ \dots \dots \dots \\ W_{hi1,10} \dots W_{hi10,10} \end{bmatrix} \end{aligned}$$

X is 1*5 matrix, weight W_{ij} is 5*10 matrix to generate a vector of 1*10
To be added to the hidden state
Let ignore the bias.

$$y_i = h_{(t-1)} W_{hi} + x_t W_{ii} + b_{ii} + b_{hi}$$

$$i_t = \sigma(y_i) = [i_{t1}, \dots, i_{t10}]$$

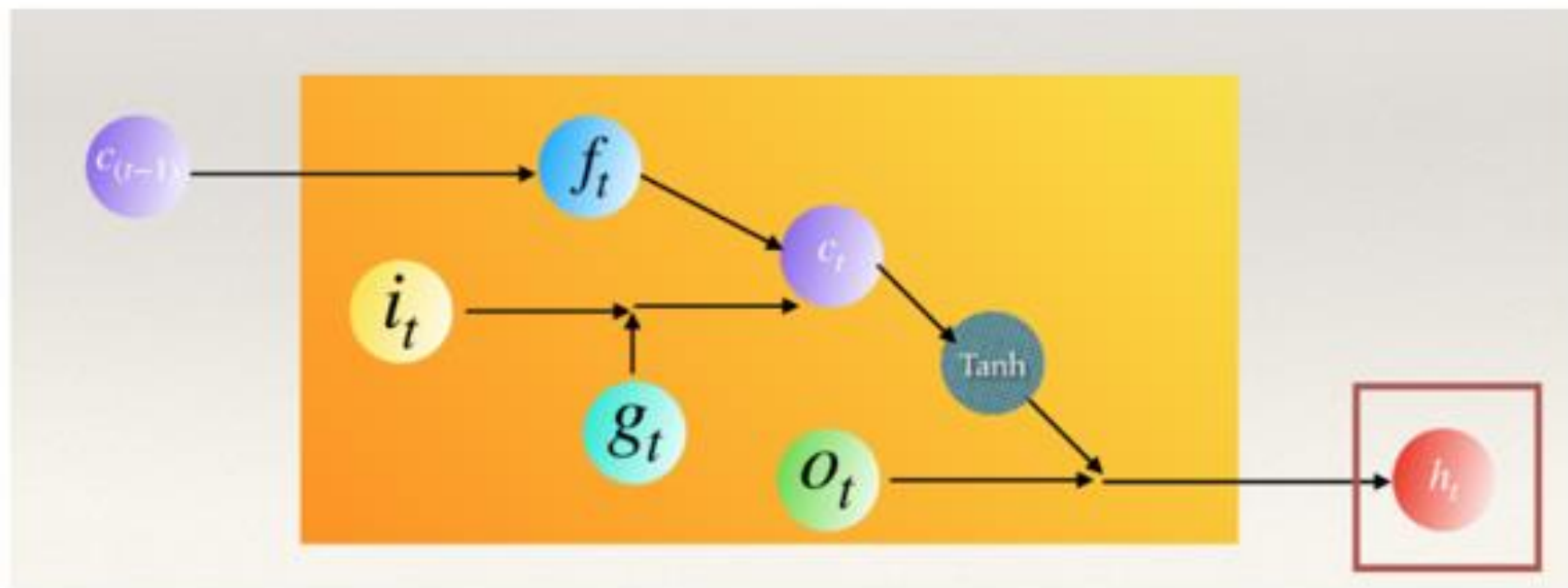
Element-wise operation

$$i_t = \sigma(y_i) = [i_{t1}, \dots, i_{t10}]$$

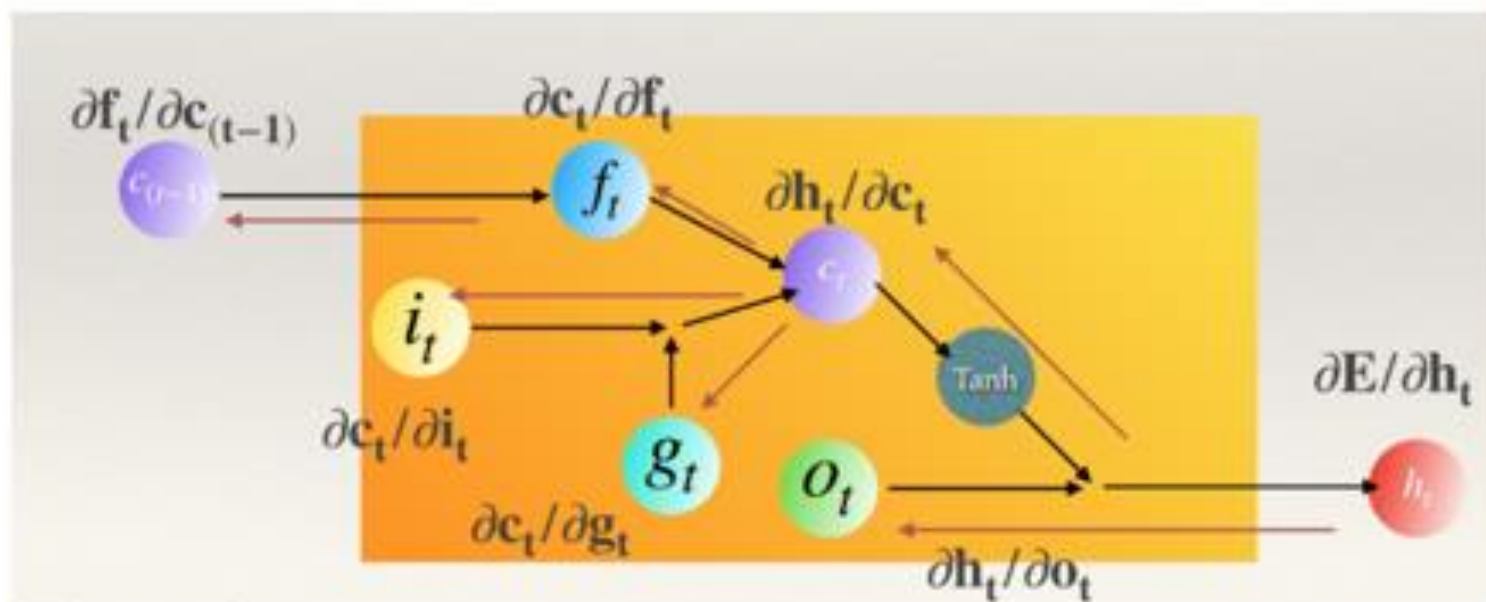
$$f_t = \sigma(y_f) = [f_{t1}, \dots, f_{t10}]$$

$$g_t = \tanh(y_g) = [g_{t1}, \dots, g_{t10}]$$

$$o_t = \sigma(y_o) = [o_{t1}, \dots, o_{t10}]$$



LSTM backward propagation (forget gate)



Chain rule

$$\partial \mathbf{E} / \partial \mathbf{f}_t = (\partial \mathbf{E} / \partial \mathbf{h}_t)(\partial \mathbf{h}_t / \partial \mathbf{c}_t)(\partial \mathbf{c}_t / \partial \mathbf{f}_t)$$

$$\partial \mathbf{h}_t / \partial \mathbf{c}_t = \mathbf{o}_t(1 - \tanh^2(\mathbf{c}_t))$$

$$\partial \mathbf{c}_t / \partial \mathbf{f}_t = \mathbf{c}_{(t-1)}$$

Training LSTM

$$i_t = \sigma(W_{ii}x_t + b_{ii} + W_{hi}h_{(t-1)} + b_{hi})$$

$$f_t = \sigma(W_{if}x_t + b_{if} + W_{hf}h_{(t-1)} + b_{hf})$$

$$g_t = \tanh(W_{ig}x_t + b_{ig} + W_{hg}h_{(t-1)} + b_{hg})$$

$$o_t = \sigma(W_{io}x_t + b_{io} + W_{ho}h_{(t-1)} + b_{ho})$$

$$c_t = f_t c_{(t-1)} + i_t g_t$$

$$h_t = o_t \tanh(c_t)$$

$$\partial \mathbf{f}_t / \partial \mathbf{W}_{if} = \begin{bmatrix} \partial f_{t1} / \partial W_{fi1,1}, \dots, \partial f_{t10} / \partial W_{fi10,1} \\ \dots \vdots \\ \partial f_{t1} / \partial W_{fi1,5}, \dots, \partial f_{t10} / \partial W_{fi10,5} \end{bmatrix}$$

where $\partial f_{ti} / \partial W_{fii,j} = f_{ti}(1 - f_{ti})x_{ti}$

Put them together

$$\partial \mathbf{E} / \partial \mathbf{W}_{\text{if}} = (\partial \mathbf{E} / \partial \mathbf{h}_t)(\partial \mathbf{h}_t / \partial \mathbf{c}_t)(\partial \mathbf{c}_t / \partial \mathbf{f}_t)(\partial \mathbf{f}_t / \partial \mathbf{W}_{\text{if}})$$

$$\partial \mathbf{h}_t / \partial \mathbf{c}_t = \mathbf{o}_t(1 - \tanh^2(\mathbf{c}_t))$$

$$\partial \mathbf{c}_t / \partial \mathbf{f}_t = \mathbf{c}_{(t-1)}$$

$$\partial \mathbf{f}_t / \partial \mathbf{W}_{\text{if}} = \begin{bmatrix} \partial f_{t1} / \partial W_{fi1,1}, \dots, \partial f_{t10} / \partial W_{fi10,1} \\ \vdots \\ \partial f_{t1} / \partial W_{fi1,5}, \dots, \partial f_{t10} / \partial W_{fi10,5} \end{bmatrix}$$

$$\text{where } \partial f_{ti} / \partial W_{fii,j} = f_{ti}(1 - f_{ti})x_{ti}$$

$$\therefore \mathbf{W}_{\text{if,new}} = \mathbf{W}_{\text{if}} - \gamma(\partial \mathbf{E} / \partial \mathbf{W}_{\text{if}})$$

Reference

1. <https://www.bioing.jku.at/publications/older/2604.pdf>
2. <https://pytorch.org/docs/stable/nn.html?highlight=lstm#torch.nn.LSTM>

Number of parameters

- Hidden state weight matrix has $(\text{hidden state length})^2$ parameters
- Input matrix has $(\text{hidden state length}) * \text{input vector length}$
- We have 4 sets of gate's matrices to train
- Many parameters to train