

CS60021: Scalable Data Mining

Large Scale Machine Learning

Sourangshu Bhattacharya

# Much of ML is optimization

## Linear Classification

$$\begin{aligned} \arg \min_w \sum_{i=1}^n \|w\|^2 + C \sum_{i=1}^n \xi_i \\ \text{s.t. } 1 - y_i x_i^T w \leq \xi_i \\ \xi_i \geq 0 \end{aligned}$$

## Maximum Likelihood

$$\arg \max_{\theta} \sum_{i=1}^n \log p_{\theta}(x_i)$$

## K-Means

$$\arg \min_{\mu_1, \mu_2, \dots, \mu_k} J(\mu) = \sum_{j=1}^k \sum_{i \in C_j} \|x_i - \mu_j\|^2$$

# Stochastic optimization

- Goal of machine learning :
  - Minimize expected loss

$$\min_h L(h) = \mathbf{E} [\text{loss}(h(x), y)]$$

given samples  $(x_i, y_i) \ i = 1, 2 \dots m$

- This is Stochastic Optimization
  - Assume loss function is convex

# Batch (sub)gradient descent for ML

- Process all examples together in each step

$$w^{(k+1)} \leftarrow w^{(k)} - \eta_t \left( \frac{1}{n} \sum_{i=1}^n \frac{\partial L(w, x_i, y_i)}{\partial w} \right)$$

where  $L$  is the regularized loss function

- Entire training set examined at each step
- Very slow when  $n$  is very large

# Stochastic (sub)gradient descent

- “Optimize” one example at a time
- Choose examples randomly (or reorder and choose in order)
  - Learning representative of example distribution

for  $i = 1$  to  $n$ :

$$w^{(k+1)} \leftarrow w^{(k)} - \eta_t \frac{\partial L(w, x_i, y_i)}{\partial w}$$

where  $L$  is the regularized loss function

# Stochastic (sub)gradient descent

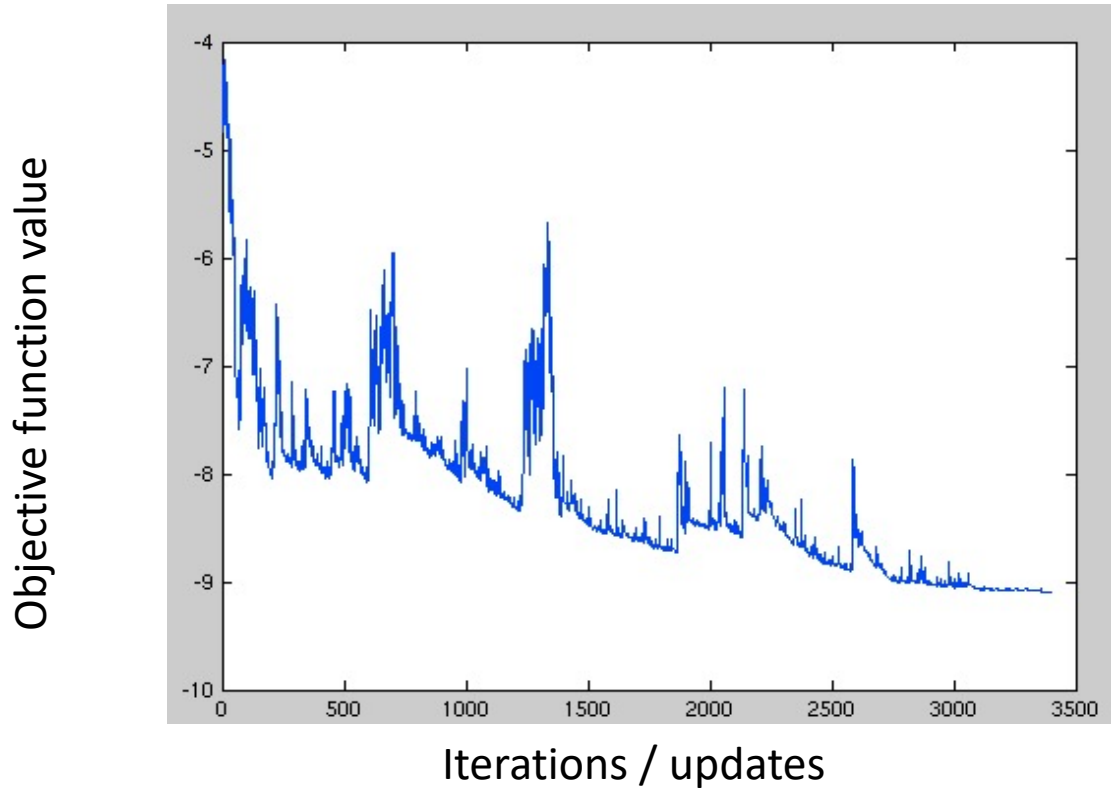
for  $i = 1$  to  $n$ :

$$w^{(k+1)} \leftarrow w^{(k)} - \eta_t \frac{\partial L(w, x_i, y_i)}{\partial w}$$

where  $L$  is the regularized loss function

- Equivalent to online learning (the weight vector  $w$  changes with every example)
- Convergence guaranteed for convex functions (to local minimum)

# SGD convergence



# Stochastic gradient descent

- Given dataset  $D = \{(x_1, y_1), \dots, (x_m, y_m)\}$
- Loss function:  $L(\theta, D) = \frac{1}{N} \sum_{i=1}^N l(\theta; x_i, y_i)$
- For linear models:  $l(\theta; x_i, y_i) = l(y_i, \theta^T \phi(x_i))$
- Assumption  $D$  is drawn IID from some distribution  $\mathcal{P}$ .
- Problem:

$$\min_{\theta} L(\theta, D)$$



# Stochastic gradient descent

- Input:  $D$
- Output:  $\bar{\theta}$

## Algorithm:

- Initialize  $\theta^0$
- For  $t = 1, \dots, T$   
$$\theta^{t+1} = \theta^t - \eta_t \nabla_{\theta} l(y_t, \theta^T \phi(x_t))$$
- $$\bar{\theta} = \frac{\sum_{t=1}^T \eta_t \theta^t}{\sum_{t=1}^T \eta_t}.$$

# SGD convergence

- Expected loss:  $s(\theta) = E_{\mathcal{P}}[l(y, \theta^T \phi(x))]$
- Optimal Expected loss:  $s^* = s(\theta^*) = \min_{\theta} s(\theta)$
- Convergence:

$$E_{\bar{\theta}}[s(\bar{\theta})] - s^* \leq \frac{R^2 + L^2 \sum_{t=1}^T \eta_t^2}{2 \sum_{t=1}^T \eta_t}$$

- Where:  $R = \|\theta^0 - \theta^*\|$
- $L = \max \nabla l(y, \theta^T \phi(x))$

# SGD convergence proof

- Define  $r_t = \|\theta^t - \theta^*\|$  and  $g_t = \nabla_{\theta} l(y_t, \theta^T \phi(x_t))$
- $r_{t+1}^2 = r_t^2 + \eta_t^2 \|g_t\|^2 - 2\eta_t (\theta^t - \theta^*)^T g_t$
- Taking expectation w.r.t  $\mathcal{P}, \bar{\theta}$  and using  $s^* - s(\theta^t) \geq g_t^T (\theta^* - \theta^t)$ , we get:

$$E_{\bar{\theta}}[r_{t+1}^2 - r_t^2] \leq \eta_t^2 L^2 + 2\eta_t (s^* - E_{\bar{\theta}}[s(\theta^t)])$$

- Taking sum over  $t = 1, \dots, T$  and using

$$E_{\bar{\theta}}[r_{T+1}^2 - r_0^2] \leq L^2 \sum_{t=0}^{T-1} \eta_t^2 + 2 \sum_{t=0}^{T-1} \eta_t (s^* - E_{\bar{\theta}}[s(\theta^t)])$$

# SGD convergence proof

- Using convexity of  $s$ :

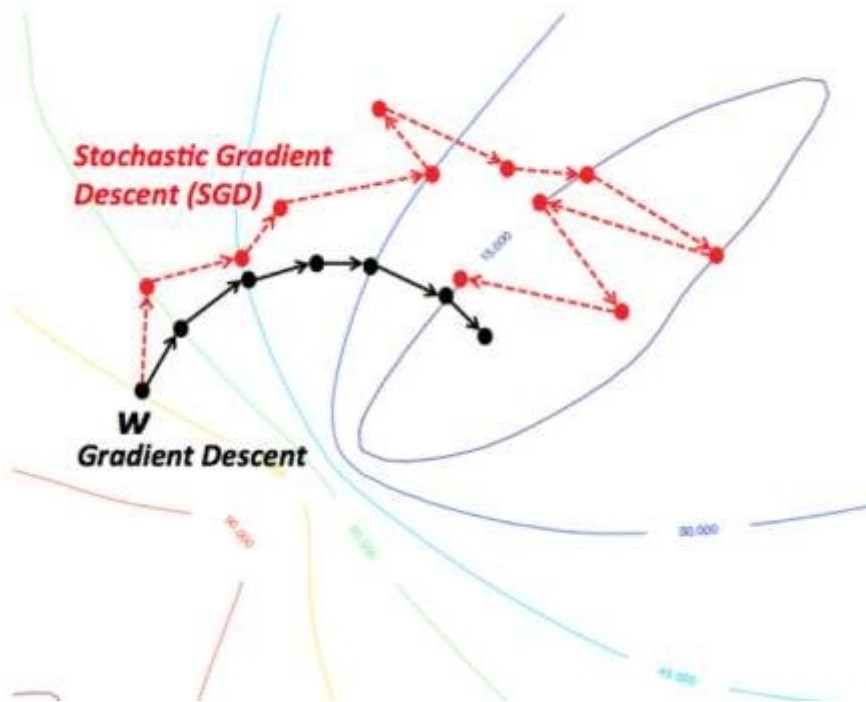
$$\left( \sum_{t=0}^{T-1} \eta_t \right) E_{\bar{\theta}} [s(\bar{\theta})] \leq E_{\bar{\theta}} \left[ \sum_{t=0}^{T-1} \eta_t s(\theta^t) \right]$$

- Substituting in the expression from previous slide:

$$E_{\bar{\theta}} [r_{t+1}^2 - r_0^2] \leq L^2 \sum_{t=0}^{T-1} \eta_t^2 + 2 \sum_{t=0}^{T-1} \eta_t (s^* - E_{\bar{\theta}} [s(\bar{\theta})])$$

- Rearranging the terms proves the result.

# The fluctuation : Batch vs SGD



<https://wikidocs.net/3413>

Batch gradient descent converges to the minimum of the basin the parameters are placed in and the **fluctuation is small**.

**SGD's fluctuation is large** but it enables to jump to new and potentially better local minima.

However, this ultimately complicates convergence to the exact minimum, as SGD will keep overshooting

# SGD - Issues

- Convergence very sensitive to learning rate ( $\eta_t$ ) (oscillations near solution due to probabilistic nature of sampling)
  - Might need to decrease with time to ensure the algorithm converges eventually
- Basically – SGD good for machine learning with large data sets!

# Mini-batch SGD

- Stochastic – 1 example per iteration
- Batch – All the examples!
- Mini-batch SGD:
  - Sample  $m$  examples at each step and perform SGD on them
- Allows for parallelization, but choice of  $m$  based on heuristics

# Example: Text categorization

- **Example by Leon Bottou:**
  - **Reuters RCV1** document corpus
    - Predict a category of a document
      - One **vs.** the rest classification
  - $n = 781,000$  training examples (documents)
  - 23,000 test examples
  - $d = 50,000$  features
    - One feature per word
    - Remove stop-words
    - Remove low frequency words



# Example: Text categorization

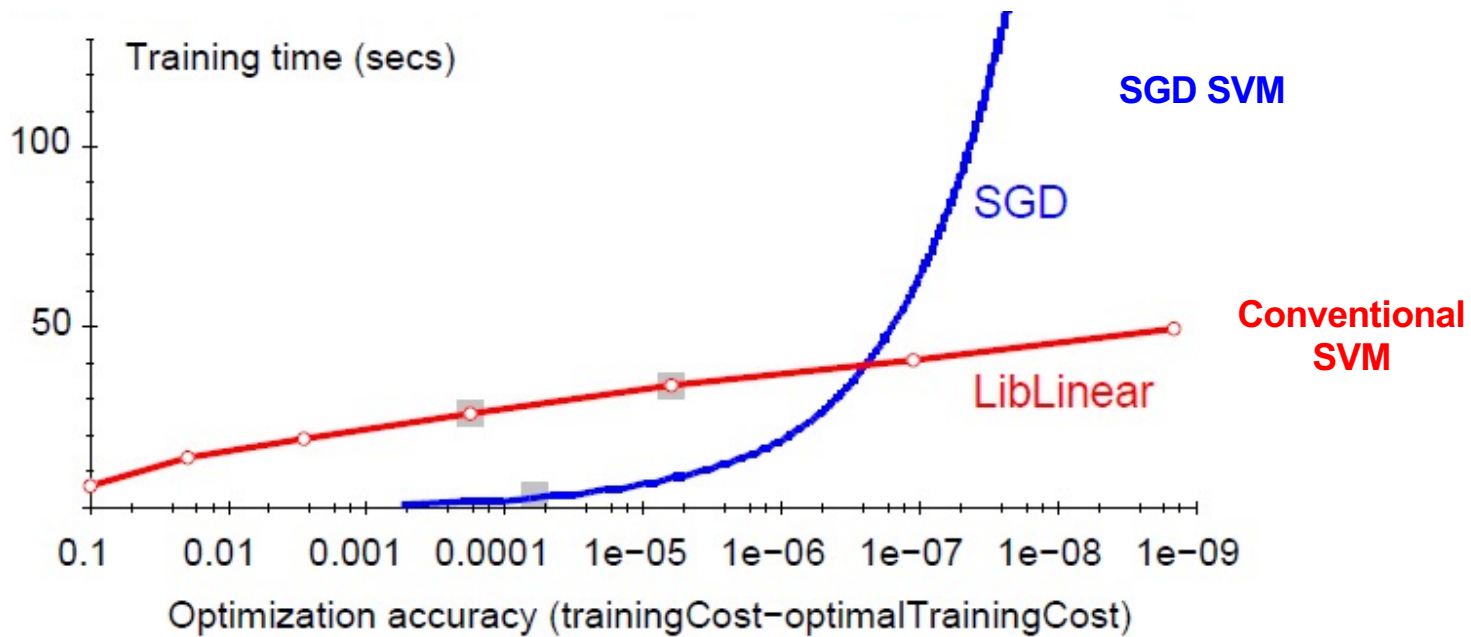
- **Questions:**

- (1) Is **SGD** successful at minimizing  $f(\mathbf{w}, \mathbf{b})$ ?
- (2) How quickly does **SGD** find the min of  $f(\mathbf{w}, \mathbf{b})$ ?
- (3) What is the error on a test set?

	<i>Training time</i>	<i>Value of <math>f(\mathbf{w}, \mathbf{b})</math></i>	<i>Test error</i>
Standard SVM	23,642 secs	0.2275	6.02%
“Fast SVM”	66 secs	0.2278	6.03%
<b>SGD SVM</b>	1.4 secs	0.2275	6.02%

- (1) SGD-SVM is successful at minimizing the value of  $f(\mathbf{w}, \mathbf{b})$
- (2) SGD-SVM is super fast
- (3) SGD-SVM test set error is comparable

# Optimization “Accuracy”



Optimization quality:  $| f(w,b) - f(w^{opt},b^{opt}) |$

For optimizing  $f(w,b)$  within reasonable quality SGD-SVM is super fast

# Practical Considerations

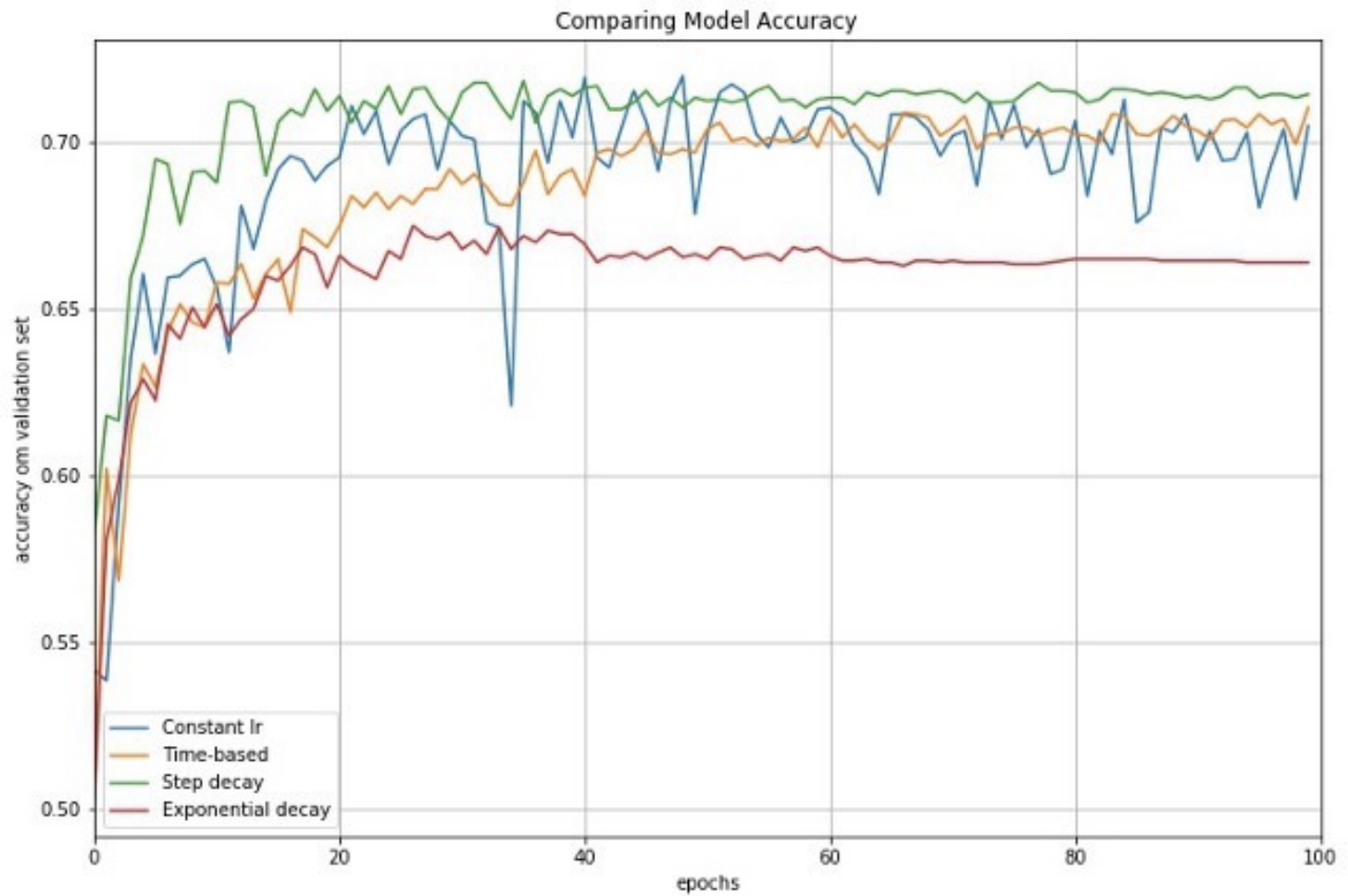
- **Need to choose learning rate  $\eta$  and  $t_0$**

$$w_{t+1} \leftarrow w_t - \frac{\eta_0}{t + t_0} \left( w_t + C \frac{\partial L(x_i, y_i)}{\partial w} \right)$$

- **Leon suggests:**

- Choose  $t_0$  so that the expected initial updates are comparable with the expected size of the weights
- Choose  $\eta$ :
  - Select a **small subsample**
  - Try various rates  $\eta$  (e.g., 10, 1, 0.1, 0.01, ...)
  - Pick the one that most reduces the cost
  - Use  $\eta$  for next 100k iterations on the full dataset

# Learning rate comparison



# **ACCELERATED GRADIENT DESCENT**

# Stochastic gradient descent

Idea: Perform a parameter update for each training example  $x(i)$  and label  $y(i)$

Update:  $\theta = \theta - \eta \cdot \nabla_{\theta} J(\theta; x(i), y(i))$

Performs redundant computations for large datasets

# Momentum gradient descent

- Idea: Overcome ravine oscillations by momentum

- 

Update:

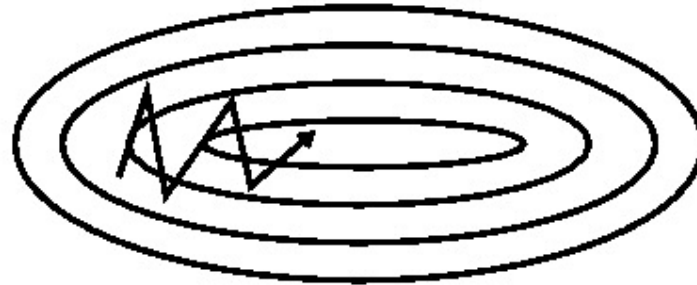
- $v_t = \gamma v_{t-1} + \eta \cdot \nabla_{\theta} J(\theta)$

- $\theta = \theta - v_t$

SGD

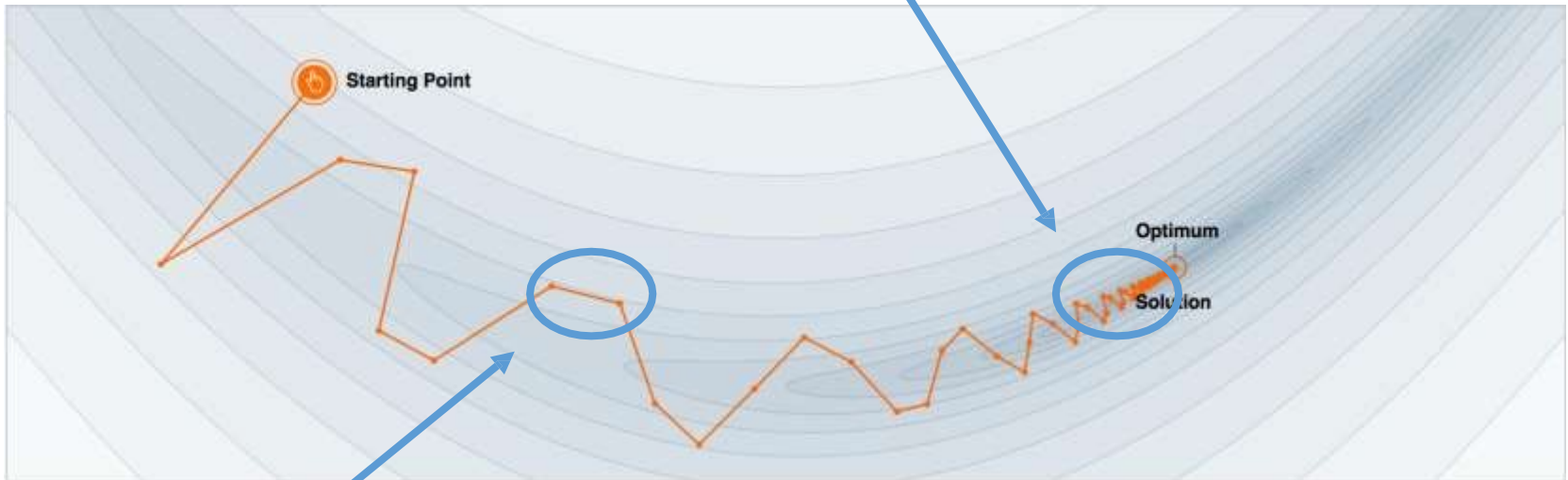


SGD with momentum



# Why Momentum Really Works

The momentum term **reduces updates for dimensions whose gradients change directions.**



The momentum term **increases for dimensions whose gradients point in the same directions.**



# Nesterov accelerated gradient

- However, a ball that rolls down a hill, blindly following the slope, is highly unsatisfactory.
- We would like to have a smarter ball that has a notion of where it is going so that it knows to slow down before the hill slopes up again.
- **Nesterov accelerated gradient** gives us a way of it.

# Nesterov accelerated gradient

$$v_t = \gamma v_{t-1} + \eta \nabla_{\theta} J(\theta - \gamma v_{t-1})$$
$$\theta = \theta - v_t$$

Approximation of the next position of the parameters(predict)

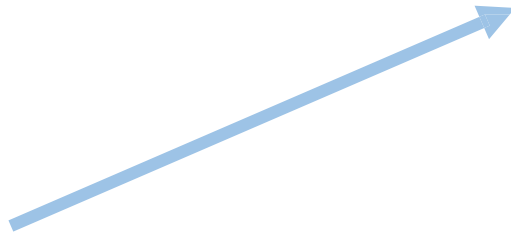
# Nesterov accelerated gradient

Approximation of the next position of the parameters' gradient(**correction**)

$$v_t = \gamma v_{t-1} + \eta \nabla_{\theta} J(\theta - \gamma v_{t-1})$$
$$\theta = \theta - v_t$$

Approximation of the next position of the parameters(**predict**)

# Nesterov accelerated gradient



Blue line : predict

Approximation of the next position of the parameters' gradient (correction)

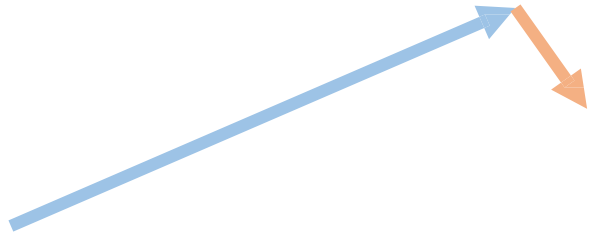
Red line : correction

$$v_t = \gamma v_{t-1} + \eta \nabla_{\theta} J(\theta - \gamma v_{t-1})$$
$$\theta = \theta - v_t$$

Green line : accumulated gradient

Approximation of the next position of the parameters (predict)

# Nesterov accelerated gradient



Blue line : predict

Approximation of the next position of the parameters' gradient (correction)

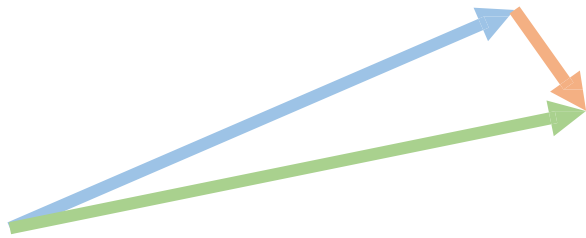
Red line : correction

$$v_t = \gamma v_{t-1} + \eta \nabla_{\theta} J(\theta - \gamma v_{t-1})$$
$$\theta = \theta - v_t$$

Green line : accumulated gradient

Approximation of the next position of the parameters (predict)

# Nesterov accelerated gradient



Blue line : predict

Red line : correction

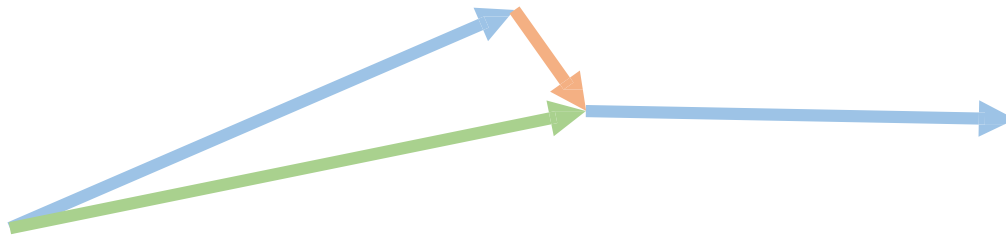
Green line : accumulated gradient

Approximation of the next position of the parameters' gradient (correction)

$$v_t = \gamma v_{t-1} + \eta \nabla_{\theta} J(\theta - \gamma v_{t-1})$$
$$\theta = \theta - v_t$$

Approximation of the next position of the parameters (predict)

# Nesterov accelerated gradient



Blue line : predict

Red line : correction

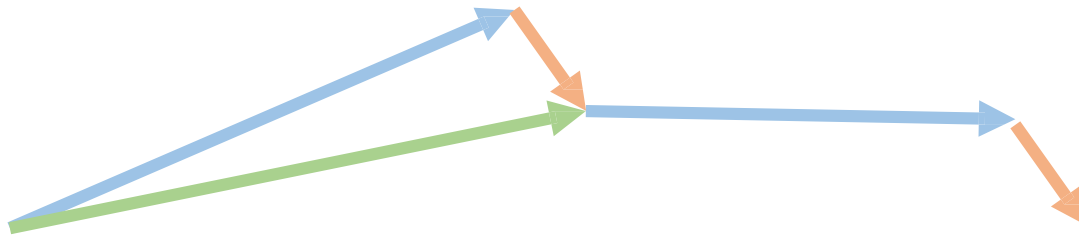
Green line : accumulated gradient

Approximation of the next position of the parameters' gradient (correction)

$$v_t = \gamma v_{t-1} + \eta \nabla_{\theta} J(\theta - \gamma v_{t-1})$$
$$\theta = \theta - v_t$$

Approximation of the next position of the parameters (predict)

# Nesterov accelerated gradient



Blue line : predict

Red line : correction

Green line : accumulated gradient

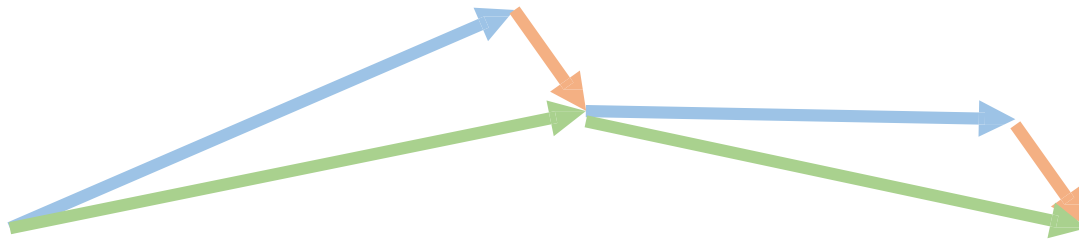
Approximation of the next position of the parameters' gradient (correction)

$$v_t = \gamma v_{t-1} + \eta \nabla_{\theta} J(\theta - \gamma v_{t-1})$$
$$\theta = \theta - v_t$$

Approximation of the next position of the parameters (predict)



# Nesterov accelerated gradient



Blue line : predict

Red line : correction

Green line : accumulated gradient

Approximation of the next position of the parameters' gradient (correction)

$$v_t = \gamma v_{t-1} + \eta \nabla_{\theta} J(\theta - \gamma v_{t-1})$$
$$\theta = \theta - v_t$$

Approximation of the next position of the parameters (predict)

# Nesterov accelerated gradient

- This anticipatory update **prevents** us from **going too fast** and **results in increased responsiveness**.
- Now , we can adapt our updates to the slope of our error function and **speed up SGD** in turn.

# What's next...?

- We also want to adapt our updates to each individual parameter to perform larger or smaller updates **depending on their importance.**
  - Adagrad
  - Adadelta
  - RMSprop
  - Adam

# Adagrad

- Adagrad adapts the learning rate to the parameters
  - Performing larger updates for infrequent
  - Performing smaller updates for frequent parameters.
- Ex.
  - Training large-scale neural nets at Google that learned to recognize cats in Youtube videos.

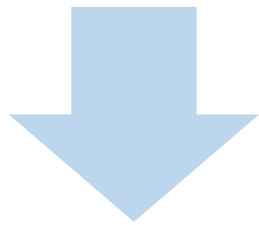
# Different learning rate for every parameter

- Previous methods :
  - **we used the same learning rate  $\eta$  for all parameters  $\theta$**
- Adagrad :
  - It uses a different learning rate for every parameter  $\theta_i$  at every time step  $t$

# Adagrad

SGD

$$\theta_{t+1,i} = \theta_{t,i} - \eta \cdot g_{t,i}$$



Adagrad

$$\theta_{t+1,i} = \theta_{t,i} - \frac{\eta}{\sqrt{G_{t,ii} + \epsilon}} \cdot g_{t,i}$$

$$G_t = \begin{matrix} \mathbb{R}^{d \times} \\ d \end{matrix} \left( \begin{array}{ccc} \square & \dots & \circ & \dots & \circ \\ \vdots & & \vdots & & \vdots \\ \circ & \dots & \square & \dots & \circ \\ \vdots & & \vdots & & \vdots \\ \circ & \dots & \circ & \dots & \square \end{array} \right)$$

$$g_{t,i} = \nabla_{\theta} J(\theta_i)$$

Vectorize

$$\theta_{t+1} = \theta_t - \frac{\eta}{\sqrt{G_t + \epsilon}} \odot g_t.$$

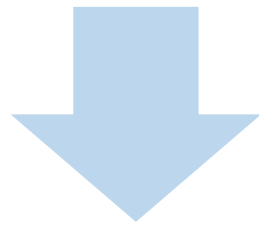
# Adagrad

SGD

$$\theta_{t+1,i} = \theta_{t,i} - \eta \cdot g_{t,i}$$

$\mathbb{R}^{d \times d}$

$$G_t = \begin{pmatrix} \square & \dots & \circ & \dots & \circ \\ \vdots & & \vdots & & \vdots \\ \circ & \dots & \square & \dots & \circ \\ \vdots & & \vdots & & \vdots \\ \circ & \dots & \circ & \dots & \square \end{pmatrix}$$



Adagrad modifies the general learning rate  $\eta$  based on the **past gradients** that have been computed for  $\theta_i$

Adagrad

$$\theta_{t+1,i} = \theta_{t,i} - \frac{\eta}{\sqrt{G_{t,ii} + \epsilon}} \cdot g_{t,i}$$

$$g_{t,i} = \nabla_{\theta} J(\theta_i)$$

Vectorize

$$\theta_{t+1} = \theta_t - \frac{\eta}{\sqrt{G_t + \epsilon}} \odot g_t.$$

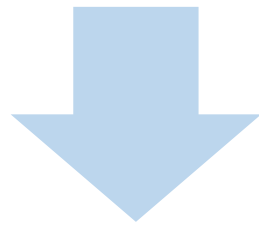
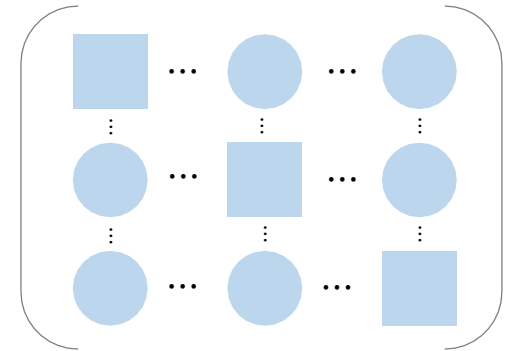
# Adagrad

SGD

$$\theta_{t+1,i} = \theta_{t,i} - \eta \cdot g_{t,i}$$

$\mathbb{R}^{d \times d}$

$G_t =$



$G_t$  is a diagonal matrix where each diagonal element  $(i,i)$  is the sum of the squares of the gradients  $\theta_i$  up to time step  $t$ .

Adagrad

$$\theta_{t+1,i} = \theta_{t,i} - \frac{\eta}{\sqrt{G_{t,ii} + \epsilon}} \cdot g_{t,i}$$

$$g_{t,i} = \nabla_{\theta} J(\theta_i)$$

Vectorize

$$\theta_{t+1} = \theta_t - \frac{\eta}{\sqrt{G_t + \epsilon}} \odot g_t.$$

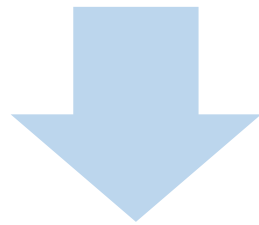


# Adagrad

SGD

$$\theta_{t+1,i} = \theta_{t,i} - \eta \cdot g_{t,i}$$

$$G_t = \begin{matrix} \mathbb{R}^{d \times d} \\ d \end{matrix} \left( \begin{array}{ccc} \square & \dots & \circ & \dots & \circ \\ \vdots & & \vdots & & \vdots \\ \circ & \dots & \square & \dots & \circ \\ \vdots & & \vdots & & \vdots \\ \circ & \dots & \circ & \dots & \square \end{array} \right)$$



$\epsilon$  is a smoothing term that avoids division by zero (usually on the order of  $1e - 8$ ).

Adagrad

$$\theta_{t+1,i} = \theta_{t,i} - \frac{\eta}{\sqrt{G_{t,ii} + \epsilon}} \cdot g_{t,i}$$

$$g_{t,i} = \nabla_{\theta} J(\theta_i)$$

Vectorize

$$\theta_{t+1} = \theta_t - \frac{\eta}{\sqrt{G_t + \epsilon}} \odot g_t.$$

# Adagrad's advantages

- Advantages :
  - It is well-suited for dealing with sparse data.
  - It greatly improved the robustness of SGD.
  - It eliminates the need to manually tune the learning rate.

# Adagrad's disadvantage

- Disadvantage :
  - Main weakness is its accumulation of the squared gradients in the denominator.

# Adagrad's disadvantage

- The disadvantage causes the learning rate to shrink and become infinitesimally small. The algorithm can no longer acquire additional knowledge.
- The following algorithms aim to resolve this flaw.
  - Adadelta
  - RMSprop
  - Adam

# Adadelta : extension of Adagrad

- Adadelta is an extension of Adagrad.
- Adagrad :
  - It accumulate all past squared gradients.
- Adadelta :
  - It restricts the window of accumulated past gradients to some fixed size  $w$ .

# Adadelta

- Instead of inefficiently storing, the sum of gradients is recursively defined as a decaying average of all past squared gradients.

$$E[g^2]_t = \gamma E[g^2]_{t-1} + (1 - \gamma)g_t^2$$

- $E[g^2]_t$  : The running average at time step  $t$ .
- $\gamma$  : A fraction similarly to the Momentum term, around 0.9

# Adadelta

Adagrad

$$\Delta\theta_t = -\frac{\eta}{\sqrt{G_t + \epsilon}} \odot g_t$$

SGD

$$\begin{aligned}\Delta\theta_t &= -\eta \cdot g_{t,i} \\ \theta_{t+1} &= \theta_t + \Delta\theta_t\end{aligned}$$



Adadelta

$$\Delta\theta_t = -\frac{\eta}{\sqrt{E[g^2]_t + \epsilon}} g_t$$

# Adadelta

Adagrad

$$\Delta\theta_t = -\frac{\eta}{\sqrt{G_t + \epsilon}} \odot g_t$$

SGD

$$\begin{aligned}\Delta\theta_t &= -\eta \cdot g_{t,i} \\ \theta_{t+1} &= \theta_t + \Delta\theta_t\end{aligned}$$



Replace the diagonal matrix  $G_t$  with the decaying average over past squared gradients  $E[g^2]_t$

Adadelta

$$\Delta\theta_t = -\frac{\eta}{\sqrt{E[g^2]_t + \epsilon}} g_t$$



# Adadelta

Adagrad

$$\Delta\theta_t = -\frac{\eta}{\sqrt{G_t + \epsilon}} \odot g_t$$

SGD

$$\begin{aligned}\Delta\theta_t &= -\eta \cdot g_{t,i} \\ \theta_{t+1} &= \theta_t + \Delta\theta_t\end{aligned}$$



Replace the diagonal matrix  $G_t$  with the decaying average over past squared gradients  $E[g^2]_t$

Adadelta

$$\Delta\theta_t = -\frac{\eta}{\sqrt{E[g^2]_t + \epsilon}} g_t$$



Adadelta

$$\Delta\theta_t = -\frac{\eta}{\text{RMS}[g]_t} g_t$$

# Update units should have the same hypothetical units

- The units in this update do not match and the update should have the same hypothetical units as the parameter.
  - As well as in SGD, Momentum, or Adagrad
- To realize this, first defining another exponentially decaying average

$$E[\Delta\theta^2]_t = \gamma E[\Delta\theta^2]_{t-1} + (1 - \gamma)\Delta\theta_t^2$$

# Adadelta

$$E[g^2]_t = \gamma E[g^2]_{t-1} + (1 - \gamma)g_t^2$$

$$E[\Delta\theta^2]_t = \gamma E[\Delta\theta^2]_{t-1} + (1 - \gamma)\Delta\theta_t^2$$



$$RMS[\Delta\theta]_t = \sqrt{E[\Delta\theta^2]_t + \epsilon}$$

Adadelta

$$\Delta\theta_t = -\frac{\eta}{\sqrt{E[g^2]_t + \epsilon}}g_t$$

Adadelta

$$\Delta\theta_t = -\frac{\eta}{RMS[g]_t}g_t$$

# Adadelta

$$E[g^2]_t = \gamma E[g^2]_{t-1} + (1 - \gamma)g_t^2$$

$$E[\Delta\theta^2]_t = \gamma E[\Delta\theta^2]_{t-1} + (1 - \gamma)\Delta\theta_t^2$$



We approximate RMS with the RMS of parameter updates until the previous time step.

$$RMS[\Delta\theta]_t = \sqrt{E[\Delta\theta^2]_t + \epsilon}$$

**Adadelta**

$$\Delta\theta_t = -\frac{\eta}{\sqrt{E[g^2]_t + \epsilon}}g_t$$

**Adadelta**

$$\Delta\theta_t = -\frac{\eta}{RMS[g]_t}g_t$$

# Adadelta update rule

- Replacing the learning rate  $\eta$  in the previous update rule with  $RMS[\Delta\theta]_{t-1}$  finally yields the Adadelta update rule:

$$\Delta\theta_t = -\frac{RMS[\Delta\theta]_{t-1}}{RMS[g]_t} g_t$$
$$\theta_{t+1} = \theta_t + \Delta\theta_t$$

- Note : **we do not even need to set a default learning rate**

# RMSprop

RMSprop and Adadelta have both been developed independently around the same time to resolve Adagrad's radically diminishing learning rates.

## RMSprop

$$E[g^2]_t = 0.9E[g^2]_{t-1} + 0.1g_t^2$$
$$\theta_{t+1} = \theta_t - \frac{\eta}{\sqrt{E[g^2]_t + \epsilon}}g_t$$

# RMSprop

RMSprop as well divides the learning rate by an exponentially decaying average of squared gradients.

## RMSprop

$$E[g^2]_t = 0.9E[g^2]_{t-1} + 0.1g_t^2$$
$$\theta_{t+1} = \theta_t - \frac{\eta}{\sqrt{E[g^2]_t + \epsilon}}g_t$$

Hinton suggests  $\gamma$  to be set to 0.9, while a good default value for the learning rate  $\eta$  is 0.001.

# Adam

- Adam's feature :
  - Storing an exponentially decaying average of past squared gradients  $v_t$  like Adadelta and RMSprop
  - Keeping an exponentially decaying average of past gradients  $m_t$ , similar to momentum.

$$m_t = \beta_1 m_{t-1} + (1 - \beta_1) g_t \longrightarrow \text{The first moment (the mean)}$$

$$v_t = \beta_2 v_{t-1} + (1 - \beta_2) g_t^2 \longrightarrow \text{The second moment (the uncentered variance)}$$



# Adam

- As  $m_t$  and  $v_t$  are initialized as vectors of 0's, they are biased towards zero.
  - Especially during the initial time steps
  - Especially when the decay rates are small
    - (i.e.  $\beta_1$  and  $\beta_2$  are close to 1).
- Counteracting these biases in Adam

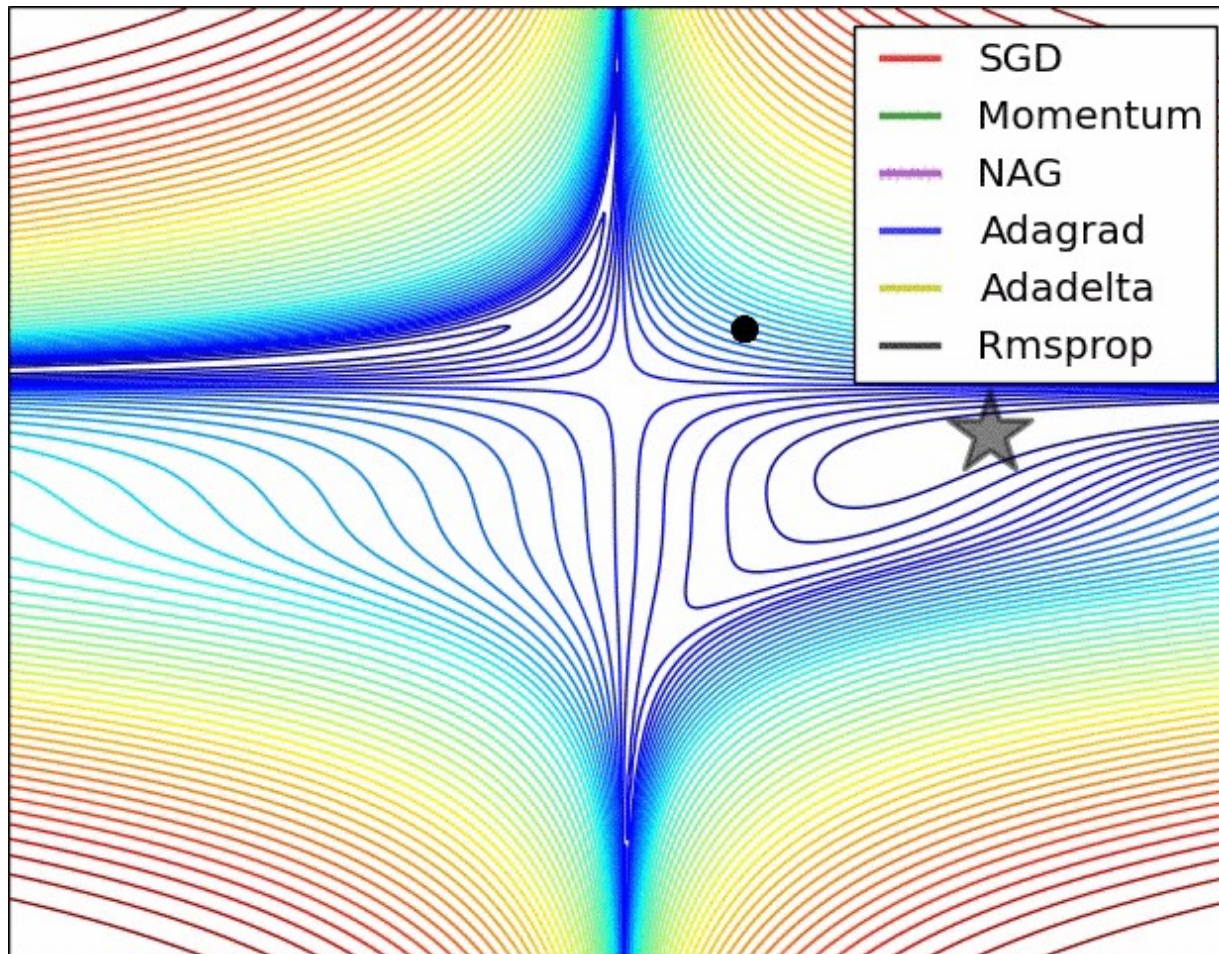
$$\hat{m}_t = \frac{m_t}{1 - \beta_1^t}$$
$$\hat{v}_t = \frac{v_t}{1 - \beta_2^t}$$

**Adam**

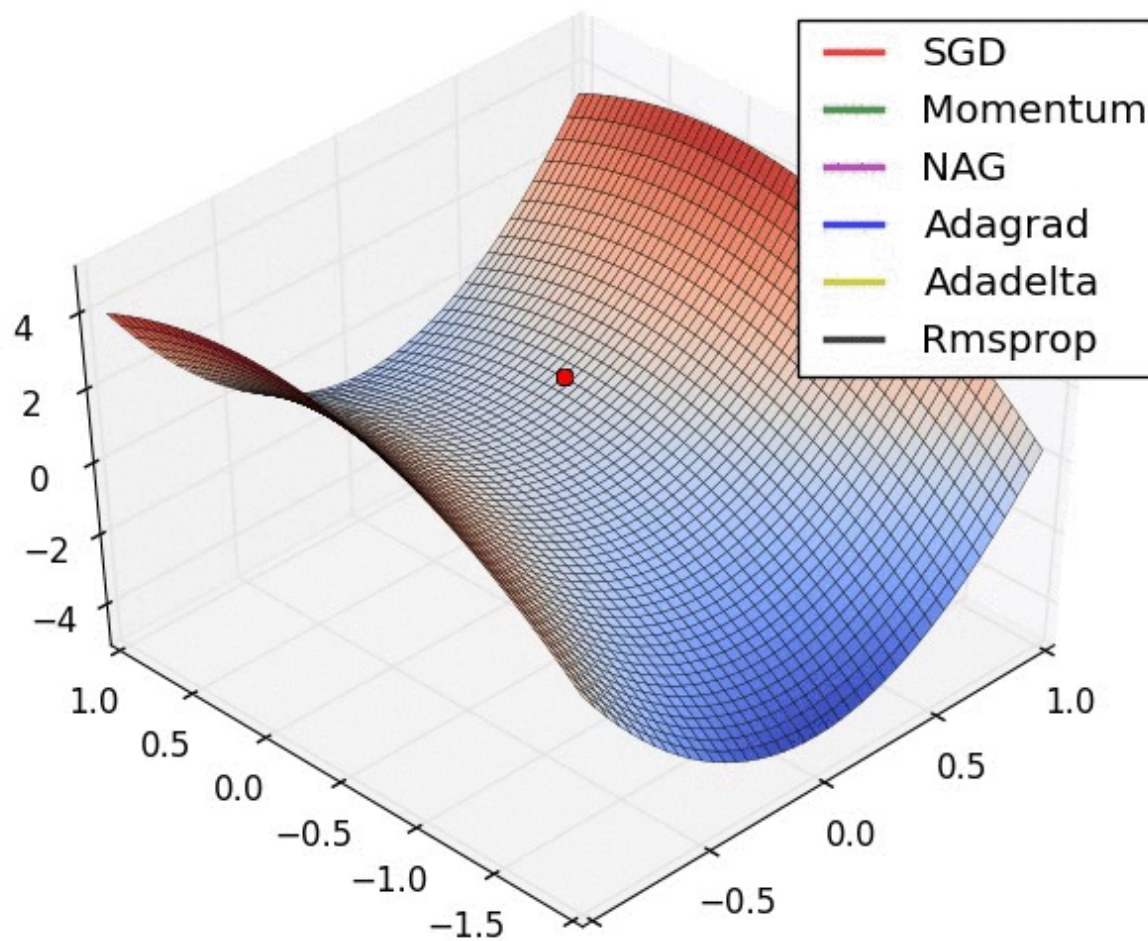
$$\theta_{t+1} = \theta_t - \frac{\eta}{\sqrt{\hat{v}_t} + \epsilon} \hat{m}_t$$

Note : default values of 0.9 for  $\beta_1$ , 0.999 for  $\beta_2$ , and  $10^{-8}$  for  $\epsilon$

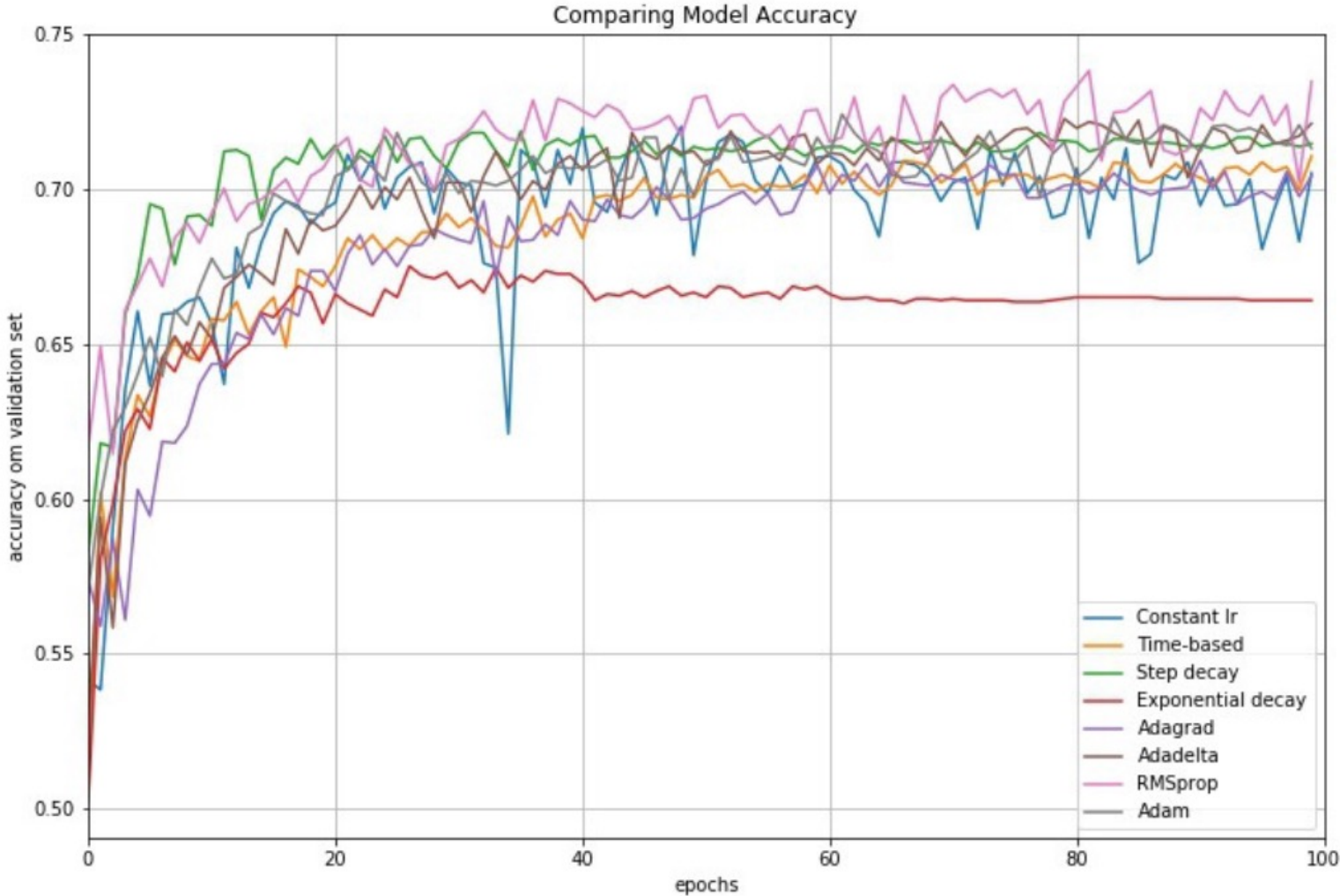
# Visualization



# Visualization



# Enhancements comparison



# Summary

- There are two main ideas at play:
  - **Momentum** : Provide consistency in update directions by incorporating past update directions.
  - **Adaptive gradient** : Scale the scale updates to individual variables using the second moment in that direction.
  - This also relates to adaptively altering step length for each direction.

# **THEORETICAL GUARRANTEES**



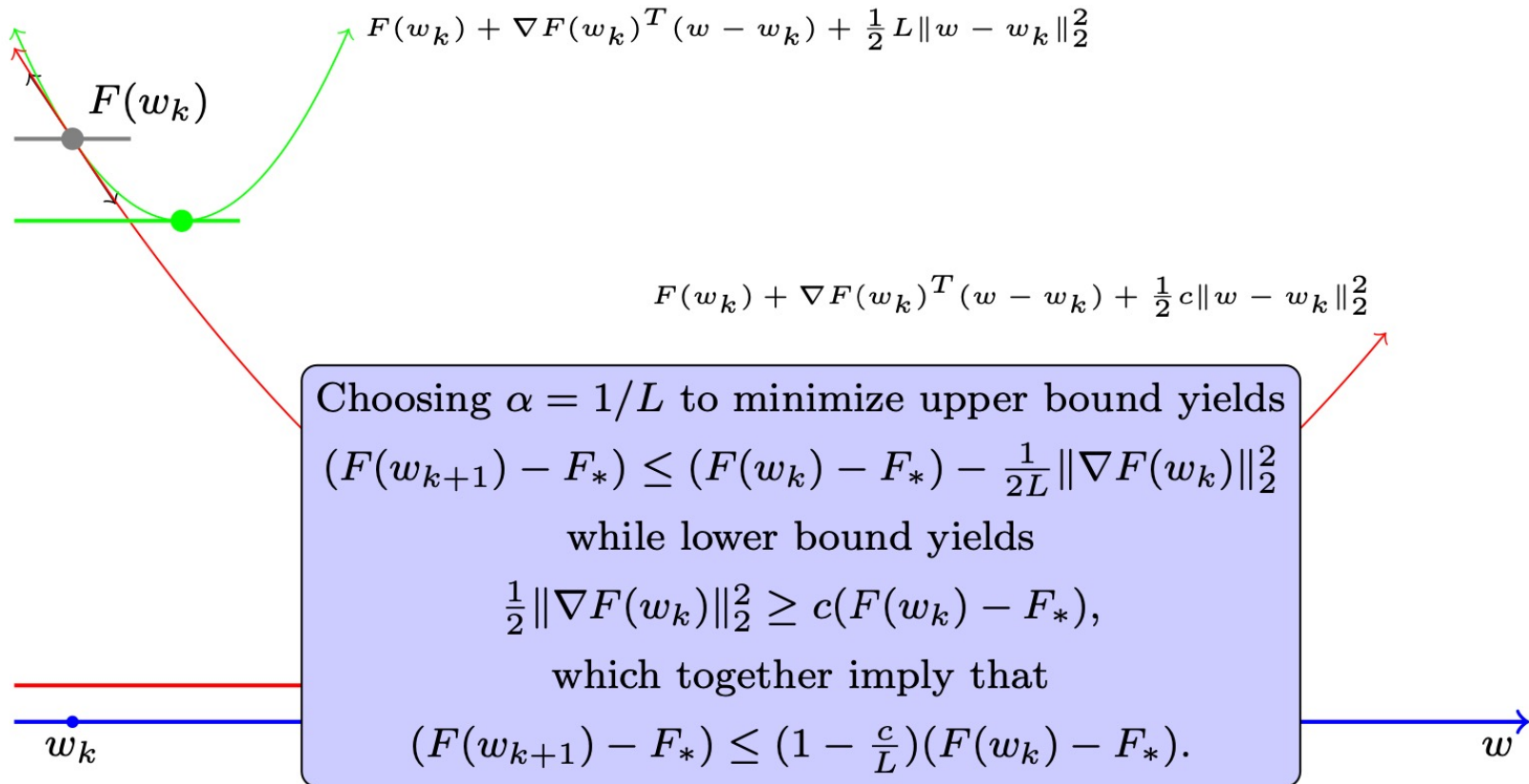
# Gradient Descent Convergence

## Assumption $\langle L/c \rangle$

*The objective function  $F : \mathbb{R}^d \rightarrow \mathbb{R}$  is*

- ▶  *$c$ -strongly convex ( $\Rightarrow$  unique minimizer) and*
- ▶  *$L$ -smooth (i.e.,  $\nabla F$  is Lipschitz continuous with constant  $L$ ).*

# Gradient Descent Convergence





# Convergence Rate and Computational Complexity

Overall Complexity ( $\epsilon$ ) = Convergence Rate<sup>-1</sup>( $\epsilon$ ) \* Complexity of each iteration

	Strongly Convex + Smooth			Convex + Smooth		
	Convergence Rate	Complexity of each iteration	Overall Complexity	Convergence Rate	Complexity of each iteration	Overall Complexity
GD	$O\left(\exp\left(-\frac{t}{Q}\right)\right)$	$O(n \cdot d)$	$O\left(nd \cdot Q \cdot \log\left(\frac{1}{\epsilon}\right)\right)$	$O\left(\frac{\beta}{t}\right)$	$O(n \cdot d)$	$O\left(nd \cdot \beta \cdot \left(\frac{1}{\epsilon}\right)\right)$
SGD	$O\left(\frac{1}{t}\right)$	$O(d)$	$O\left(\frac{d}{\epsilon}\right)$	$O\left(\frac{1}{\sqrt{t}}\right)$	$O(d)$	$O\left(\frac{d}{\epsilon^2}\right)$

# SGD Analysis

**THEOREM 14.8** *Let  $B, \rho > 0$ . Let  $f$  be a convex function and let  $\mathbf{w}^* \in \operatorname{argmin}_{\mathbf{w}: \|\mathbf{w}\| \leq B} f(\mathbf{w})$ . Assume that SGD is run for  $T$  iterations with  $\eta = \sqrt{\frac{B^2}{\rho^2 T}}$ . Assume also that for all  $t$ ,  $\|\mathbf{v}_t\| \leq \rho$  with probability 1. Then,*

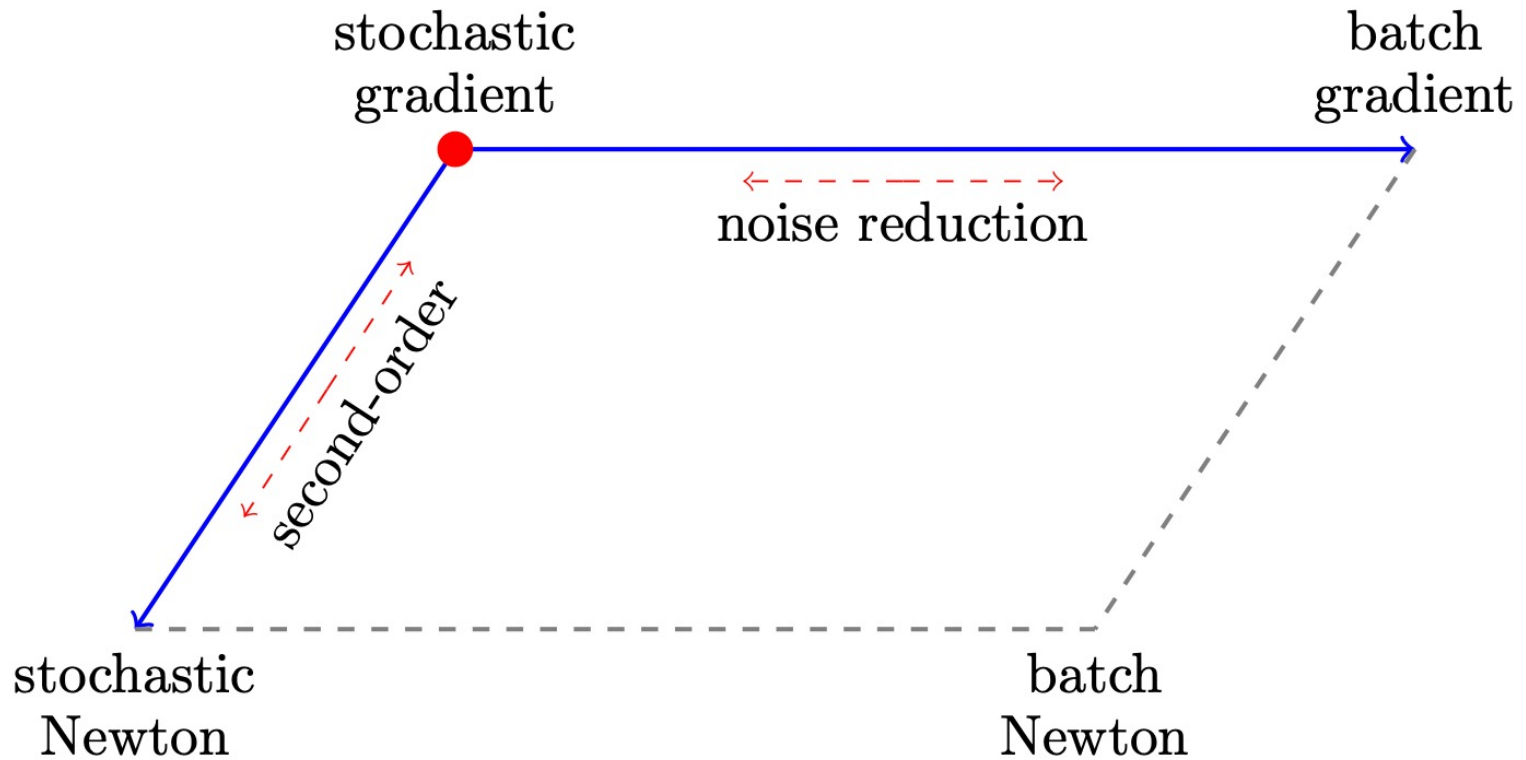
$$\mathbb{E}[f(\bar{\mathbf{w}})] - f(\mathbf{w}^*) \leq \frac{B\rho}{\sqrt{T}}.$$

*Therefore, for any  $\epsilon > 0$ , to achieve  $\mathbb{E}[f(\bar{\mathbf{w}})] - f(\mathbf{w}^*) \leq \epsilon$ , it suffices to run the SGD algorithm for a number of iterations that satisfies*

$$T \geq \frac{B^2 \rho^2}{\epsilon^2}.$$

# **LINEAR RATE METHODS**

# Improving SGD



Slides taken from Jorge Nocedal

# Stochastic Averaged Gradient

- Can we have a rate of  $O(\rho^t)$  with only 1 gradient evaluation per iteration?
  - YES! The **stochastic average gradient (SAG)** algorithm:
    - Randomly select  $i_t$  from  $\{1, 2, \dots, N\}$  and compute  $f'_{i_t}(x^t)$ .

$$x^{t+1} = x^t - \frac{\alpha^t}{N} \sum_{i=1}^N y_i^t$$

- **Memory:**  $y_i^t = \nabla f_i(x^t)$  from the **last  $t$**  where  $i$  was selected.  
[Le Roux et al., 2012]
  - **Stochastic** variant of increment average gradient (IAG).  
[Blatt et al., 2007]
  - Assumes gradients of non-selected examples don't change.
  - Assumption becomes accurate as  $\|x^{t+1} - x^t\| \rightarrow 0$ .

# SAG Convergence Rate

- If each  $f'_i$  is  $L$ -continuous and  $f$  is strongly-convex, with  $\alpha_t = 1/16L$  SAG has

$$\mathbb{E}[f(x^t) - f(x^*)] \leq \left(1 - \min\left\{\frac{\mu}{16L}, \frac{1}{8N}\right\}\right)^t C,$$

where

$$C = [f(x^0) - f(x^*)] + \frac{4L}{N} \|x^0 - x^*\|^2 + \frac{\sigma^2}{16L}.$$

- **Linear convergence rate but only 1 gradient per iteration.**
  - For well-conditioned problems, constant reduction per pass:

$$\left(1 - \frac{1}{8N}\right)^N \leq \exp\left(-\frac{1}{8}\right) = 0.8825.$$

- For ill-conditioned problems, almost same as deterministic method (but  $N$  times faster).

# SAG Convergence Rate

- Assume that  $N = 700000$ ,  $L = 0.25$ ,  $\mu = 1/N$ :
  - Gradient method has rate  $\left(\frac{L-\mu}{L+\mu}\right)^2 = 0.99998$ .
  - Accelerated gradient method has rate  $(1 - \sqrt{\frac{\mu}{L}}) = 0.99761$ .
  - **SAG ( $N$  iterations)** has rate  $(1 - \min\{\frac{\mu}{16L}, \frac{1}{8N}\})^N = 0.88250$ .
  - *Fastest possible* first-order method:  $\left(\frac{\sqrt{L}-\sqrt{\mu}}{\sqrt{L}+\sqrt{\mu}}\right)^2 = 0.99048$ .
- **SAG beats two lower bounds:**
  - Stochastic gradient bound (of  $O(1/t)$ ).
  - Deterministic gradient bound (for typical  $L$ ,  $\mu$ , and  $N$ ).
- Number of  $f'_i$  evaluations to reach  $\epsilon$ :
  - Stochastic:  $O(\frac{L}{\mu}(1/\epsilon))$ .
  - Gradient:  $O(N\frac{L}{\mu}\log(1/\epsilon))$ .
  - Accelerated:  $O(N\sqrt{\frac{L}{\mu}}\log(1/\epsilon))$ .
  - **SAG:**  $O(\max\{N, \frac{L}{\mu}\}\log(1/\epsilon))$ .

# SAG Convergence Rate

- Use SGD for well conditioned problems.
- Use Accelerated SGD for ill-conditioned problems where  $N$  is lower than  $O(\sqrt{C})$ .
- Otherwise use SAG.



# SAG Implementation

- Basic SAG algorithm:
  - while(1)
  - Sample  $i$  from  $\{1, 2, \dots, N\}$ .
  - Compute  $f'_i(x)$ .
  - $d = d - y_i + f'_i(x)$ .
  - $y_i = f'_i(x)$ .
  - $x = x - \frac{\alpha}{N}d$ .
- Practical variants of the basic algorithm allow:
  - Regularization.
  - Sparse gradients.
  - Automatic step-size selection.
    - Common to use adaptive step-size procedure to estimate  $L$ .
  - Termination criterion.
    - Can use  $\|x^{t+1} - x^t\|/\alpha = \frac{1}{n}d \approx \|\nabla f(x^t)\|$  to decide when to stop.
  - Acceleration [Lin et al., 2015].
  - Adaptive non-uniform sampling [Schmidt et al., 2013].

# SAG Implementation

- Does **re-shuffling** and doing full passes work better?
  - For classic SG: **Maybe?**
    - Noncommutative arithmetic-geometric mean inequality conjecture.
  - For SAG: **NO.**
  - Performance is intermediate between IAG and SAG.
- Can **non-uniform** sampling help?
  - For classic SG methods, can only improve constants.
  - For SAG, **bias sampling towards Lipschitz constants  $L_i$ ,**

[Recht & Ré, 2012]

$$\|\nabla f_i(x) - \nabla f_i(y)\| \leq L_i \|x - y\|.$$

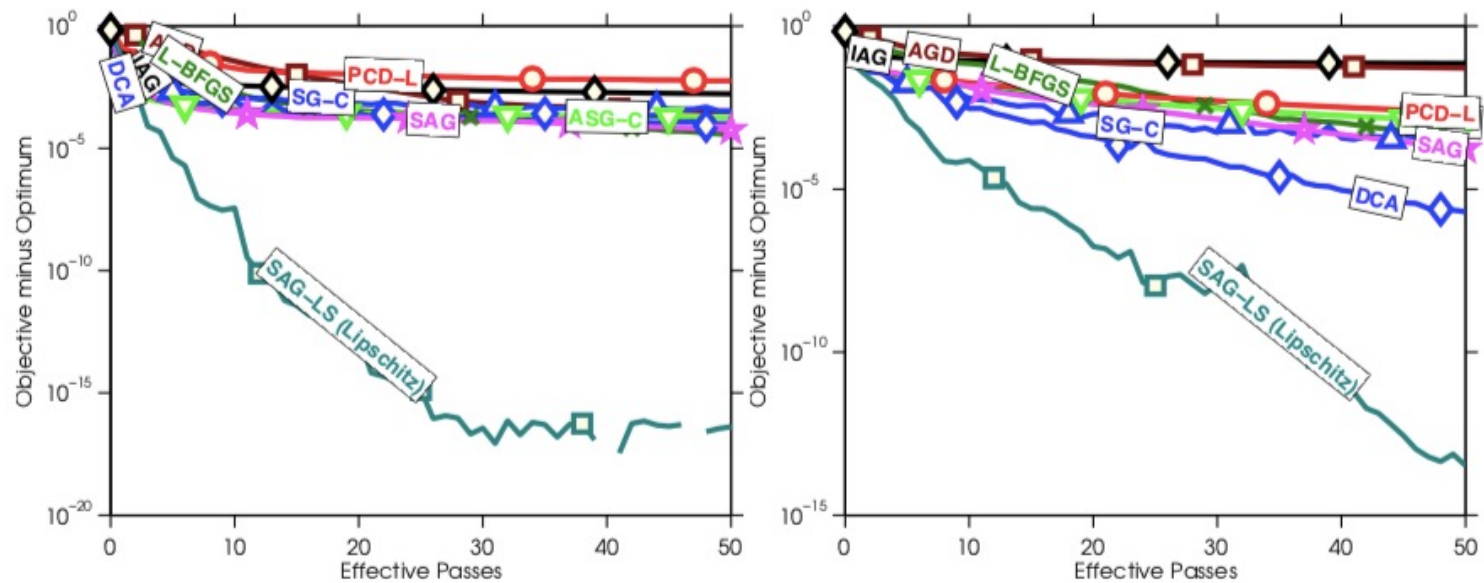
improves rate to depend on  $L_{\text{mean}}$  instead of  $L_{\text{max}}$ .

(with bigger step size)

- **Adaptively estimate  $L_i$  as you go.** (see paper/code).
- Slowly learns to **ignore well-classified examples.**

## SAG with Non-Uniform Sampling

- protein ( $n = 145751$ ,  $p = 74$ ) and sido ( $n = 12678$ ,  $p = 4932$ )



- Adaptive non-uniform sampling helps a lot.

# Stochastic Variance Reduced GD

SVRG algorithm:

- Start with  $x_0$
- for  $s = 0, 1, 2 \dots$ 
  - $d_s = \frac{1}{N} \sum_{i=1}^N f'_i(x_s)$
  - $x^0 = x_s$
  - for  $t = 1, 2, \dots, m$ 
    - Randomly pick  $i_t \in \{1, 2, \dots, N\}$
    - $x^t = x^{t-1} - \alpha_t (f'_{i_t}(x^{t-1}) - f'_{i_t}(x_s) + d_s)$ .
  - $x_{s+1} = x^t$  for random  $t \in \{1, 2, \dots, m\}$ .

Requires **2 gradients per iteration and occasional full passes**,  
but **only requires storing  $d_s$  and  $x_s$** .

Practical issues similar to SAG (acceleration versions, automatic step-size/termination, handles sparsity/regularization, non-uniform sampling, mini-batches).

# BATCH NORMALIZATION

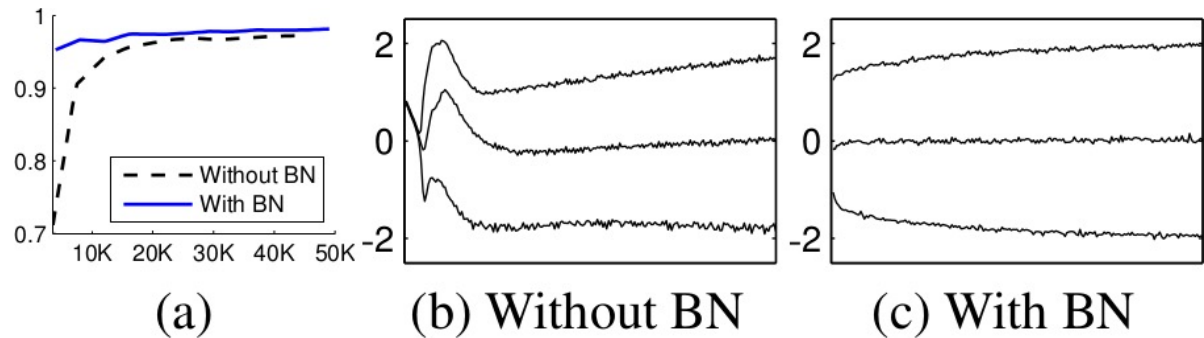
Slides taken from Jude Shavlik: [http://pages.cs.wisc.edu/~shavlik/cs638\\_cs838.html](http://pages.cs.wisc.edu/~shavlik/cs638_cs838.html)

# Batch normalization:

## Other benefits in practice

- BN reduces training times. (Because of less Covariate Shift, less exploding/vanishing gradients.)
- BN reduces demand for regularization, e.g. dropout or L2 norm.
  - Because the means and variances are calculated over batches and therefore every normalized value depends on the current batch. I.e. the network can no longer just memorize values and their correct answers.)
- BN allows higher learning rates. (Because of less danger of exploding/vanishing gradients.)
- BN enables training with saturating nonlinearities in deep networks, e.g. sigmoid. (Because the normalization prevents them from getting stuck in saturating ranges, e.g. very high/low values for sigmoid.)

# Batch normalization: Better accuracy , faster.



*BN applied to MNIST (a), and activations of a randomly selected neuron over time (b, c), where the middle line is the median activation, the top line is the 15th percentile and the bottom line is the 85th percentile.*

# Why the naïve approach Does not work?

- Normalizes layer inputs to zero mean and unit variance. *whitening*.
- Naive method: Train on a batch. Update model parameters. Then normalize. **Doesn't work:** Leads to exploding biases while distribution parameters (mean, variance) don't change.
  - If we do it this way gradient always ignores the effect that the normalization for the next batch would have
  - i.e. : **“The issue with the above approach is that the gradient descent optimization does not take into account the fact that the normalization takes place”**



# Doing it the “correct way” Is too expensive!

- A proper method has to include the current example batch *and* somehow all previous batches ( all examples) in the normalization step.
- This leads to calculating in covariance matrix and its inverse square root. That's expensive. The authors found a faster way!

The issue with the above approach is that the gradient descent optimization does not take into account the fact that the normalization takes place. To address this issue, we would like to ensure that, for any parameter values, the network *always* produces activations with the desired distribution. Doing so would allow the gradient of the loss with respect to the model parameters to account for the normalization, and for its dependence on the model parameters  $\Theta$ . Let again  $x$  be a layer input, treated as a vector, and  $\mathcal{X}$  be the set of these inputs over the training data set. The normalization can then be written as a transformation

$$\hat{x} = \text{Norm}(x, \mathcal{X})$$

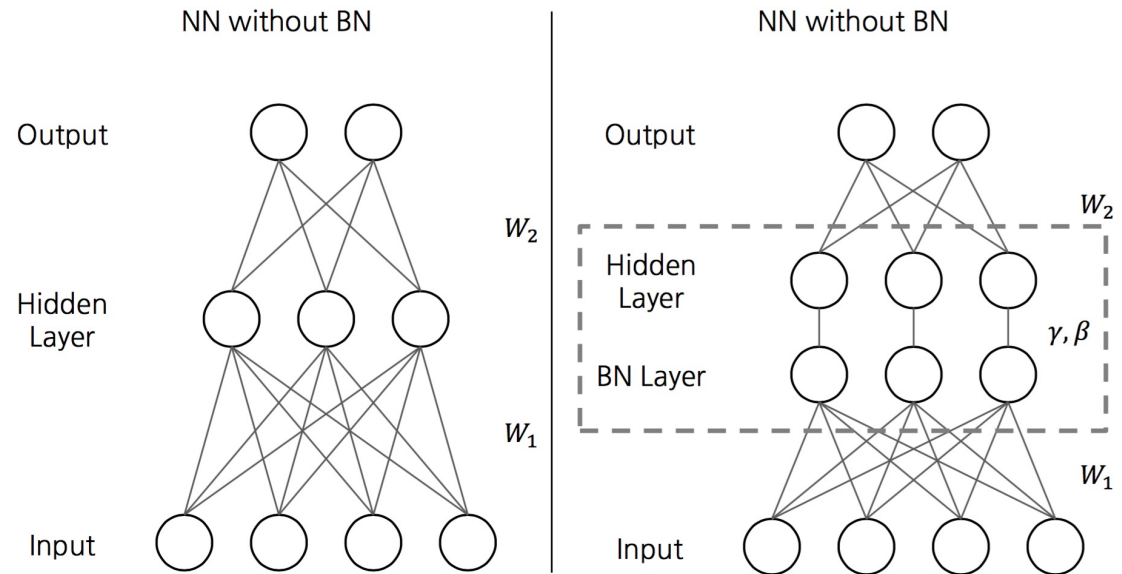
which depends not only on the given training example  $x$  but on all examples  $\mathcal{X}$  – each of which depends on  $\Theta$  if  $x$  is generated by another layer. For backpropagation, we would need to compute the Jacobians  $\frac{\partial \text{Norm}(x, \mathcal{X})}{\partial x}$  and  $\frac{\partial \text{Norm}(x, \mathcal{X})}{\partial \mathcal{X}}$ ; ignoring the latter term would lead to the ex-

ploration described above. Within this framework, whitening the layer inputs is expensive, as it requires computing the covariance matrix  $\text{Cov}[x] = E_{x \in \mathcal{X}}[xx^T] - E[x]E[x]^T$  and its inverse square root, to produce the whitened activations  $\text{Cov}[x]^{-1/2}(x - E[x])$ , as well as the derivatives of these transforms for backpropagation. This motivates us to seek an alternative that performs input normalization in a way that is differentiable and does not require the analysis of the entire training set after every parameter update.

# The proposed solution: To add an extra reg

we introduce, for each activation  $x^{(k)}$ , a pair of parameters  $\gamma^{(k)}, \beta^{(k)}$ , which scale and shift the normalized value:

$$y^{(k)} = \gamma^{(k)} \hat{x}^{(k)} + \beta^{(k)}.$$



A new layer is added so the gradient can “see” the normalization and make adjustments if needed.

# Algorithm Summary:

## Normalization via Mini-Batch Statistics

- Each feature (component) is normalized individually
- Normalization according to:
  - $\text{componentNormalizedValue} = (\text{componentOldValue} - E[\text{component}]) / \sqrt{\text{Var}(\text{component})}$
- A new layer is added so the gradient can “see” the normalization and made adjustments if needed.
  - The new layer has the power to learn the identity function to de-normalize the features if necessary!
  - Full formula:  $\text{newValue} = \gamma * \text{componentNormalizedValue} + \beta$  (gamma and beta learned per component)
- E and Var are estimated for each mini batch.
- BN is fully differentiable.

# The Batch Transformation: formally from the paper.

**Input:** Values of  $x$  over a mini-batch:  $\mathcal{B} = \{x_{1\dots m}\}$ ;

Parameters to be learned:  $\gamma, \beta$

**Output:**  $\{y_i = \text{BN}_{\gamma, \beta}(x_i)\}$

$$\mu_{\mathcal{B}} \leftarrow \frac{1}{m} \sum_{i=1}^m x_i \quad // \text{ mini-batch mean}$$

$$\sigma_{\mathcal{B}}^2 \leftarrow \frac{1}{m} \sum_{i=1}^m (x_i - \mu_{\mathcal{B}})^2 \quad // \text{ mini-batch variance}$$

$$\hat{x}_i \leftarrow \frac{x_i - \mu_{\mathcal{B}}}{\sqrt{\sigma_{\mathcal{B}}^2 + \epsilon}} \quad // \text{ normalize}$$

$$y_i \leftarrow \gamma \hat{x}_i + \beta \equiv \text{BN}_{\gamma, \beta}(x_i) \quad // \text{ scale and shift}$$

# The full algorithm as proposed in the paper

**Input:** Network  $N$  with trainable parameters  $\Theta$ ;  
subset of activations  $\{x^{(k)}\}_{k=1}^K$

**Output:** Batch-normalized network for inference,  $N_{\text{BN}}^{\text{inf}}$

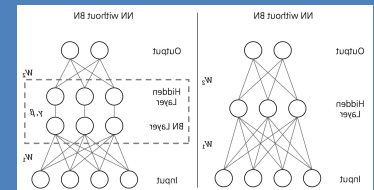
- 1:  $N_{\text{BN}}^{\text{tr}} \leftarrow N$  // Training BN network
- 2: **for**  $k = 1 \dots K$  **do**
- 3: Add transformation  $y^{(k)} = \text{BN}_{\gamma^{(k)}, \beta^{(k)}}(x^{(k)})$  to  $N_{\text{BN}}^{\text{tr}}$  (Alg. 1)
- 4: Modify each layer in  $N_{\text{BN}}^{\text{tr}}$  with input  $x^{(k)}$  to take  $y^{(k)}$  instead
- 5: **end for**
- 6: Train  $N_{\text{BN}}^{\text{tr}}$  to optimize the parameters  $\Theta \cup \{\gamma^{(k)}, \beta^{(k)}\}_{k=1}^K$
- 7:  $N_{\text{BN}}^{\text{inf}} \leftarrow N_{\text{BN}}^{\text{tr}}$  // Inference BN network with frozen // parameters
- 8: **for**  $k = 1 \dots K$  **do**
- 9: // For clarity,  $x \equiv x^{(k)}$ ,  $\gamma \equiv \gamma^{(k)}$ ,  $\mu_{\mathcal{B}} \equiv \mu_{\mathcal{B}}^{(k)}$ , etc.
- 10: Process multiple training mini-batches  $\mathcal{B}$ , each of size  $m$ , and average over them:
 
$$E[x] \leftarrow E_{\mathcal{B}}[\mu_{\mathcal{B}}]$$

$$\text{Var}[x] \leftarrow \frac{m}{m-1} E_{\mathcal{B}}[\sigma_{\mathcal{B}}^2]$$
- 11: In  $N_{\text{BN}}^{\text{inf}}$ , replace the transform  $y = \text{BN}_{\gamma, \beta}(x)$  with  $y = \frac{\gamma}{\sqrt{\text{Var}[x] + \epsilon}} \cdot x + \left(\beta - \frac{\gamma E[x]}{\sqrt{\text{Var}[x] + \epsilon}}\right)$
- 12: **end for**

**Algorithm 2:** Training a Batch-Normalized Network

Alg 1 (previous slide)

Architecture modification



**Note that  $\text{BN}(x)$  is different during test...**

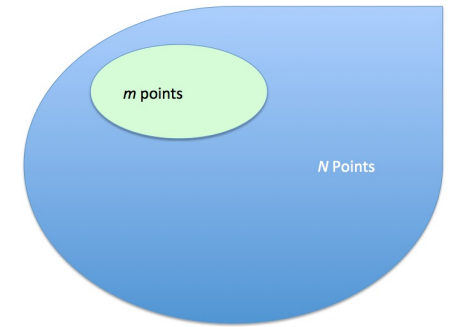
$$\sigma_{\mathcal{B}}^2 \leftarrow \frac{1}{m} \sum_{i=1}^m (x_i - \mu_{\mathcal{B}})^2$$

**Vs.**

$$\text{Var}[x] \leftarrow \frac{m}{m-1} E_{\mathcal{B}}[\sigma_{\mathcal{B}}^2]$$

# Populations stats vs. sample stats

- In algorithm 1, we are estimating the true mean and variance over the entire population for a given batch.
- When doing inference you're minibatching your way through the entire dataset, you're calculating statistics on a per sample/batch basis. We want our sample statistics to be *unbiased* to population statistics.



	Population Statistics over $N$	Sample/Batch Statistics over $m$
Mean Estimate	$\mu = \frac{1}{N} \sum_i x_i$	$\bar{x} = \frac{1}{m} \sum_i x_i$
Variance Estimate	$\sigma = \frac{1}{N} \sum_i (x_i - \mu)^2$	$\sigma_B = \frac{1}{m-1} \sum_i (x_i - \bar{x})^2$

## ACCELERATING BN NETWORKS

Batch normalization only not enough!

- Increase learning rate.
- Remove Dropout.
- Shuffle training examples more thoroughly
- Reduce the L2 weight regularization.
- Accelerate the learning rate decay.
- Reduce the photometric distortions.

# References:

- SGD proof by Yuri Nesterov.
- MMDS <http://www.mmds.org/>
- *Blog of Sebastian Ruder* <http://ruder.io/optimizing-gradient-descent/>
- *Learning rate comparison* <https://towardsdatascience.com/learning-rate-schedules-and-adaptive-learning-rate-methods-for-deep-learning-2c8f433990d1>