

CS60021: Scalable Data Mining

Large Scale Machine Learning

Sourangshu Bhattacharya

Much of ML is optimization

Linear Classification

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

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

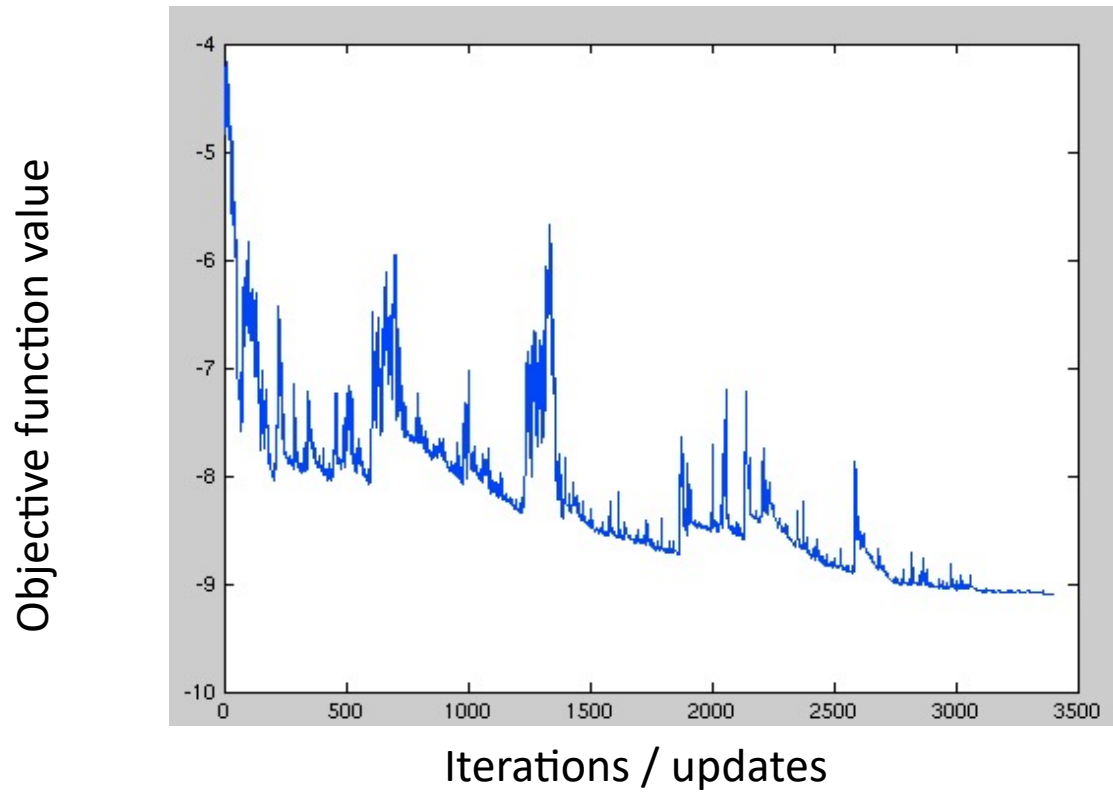
for $i = 1$ to n :

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

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

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

Algorithm:

- Initialize θ^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:

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

- 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

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

$E_{\bar{\theta}}[s]$

$\leq K$

how far in space

$\sum \eta_t$

$2\eta_t$

SGD convergence proof

- Using convexity of s :

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

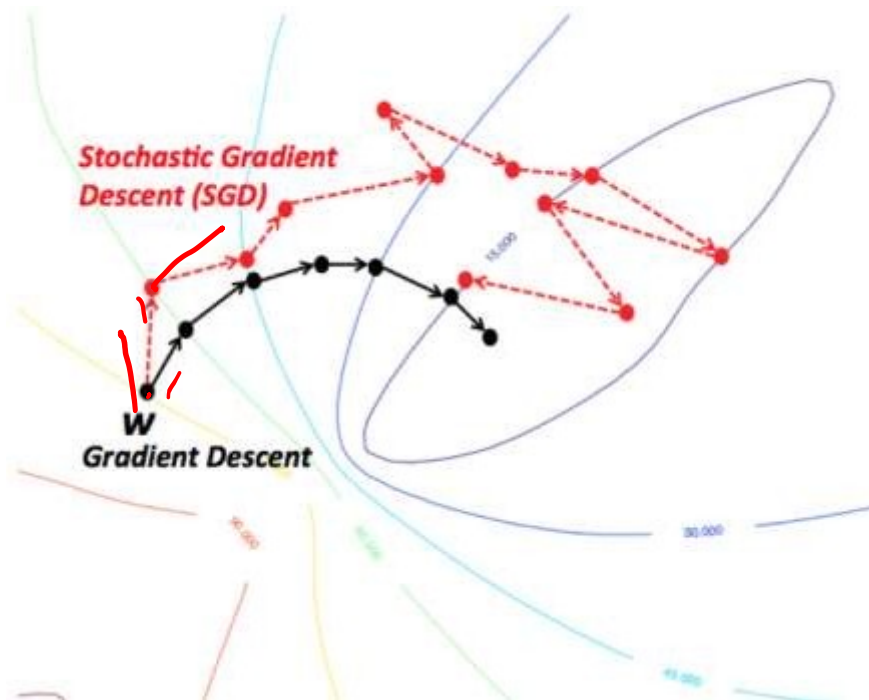
- Substituting in the expression from previous slide:

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

- Rearranging the terms proves the result.

$$\begin{aligned} E_{\bar{\theta}} [s(\bar{\theta})] &\leq s^* \\ L^2 \sum \eta_t^2 + 2 \sum \eta_t &\leq 2 \sum \eta_t \end{aligned}$$

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 (η_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

$$\theta^{t+1} = \theta^t - \eta_t g(B_{i_t})$$

- 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

$$g(B_i) = \frac{1}{m} \sum_{n_j \in B_i} g(n_j)$$

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

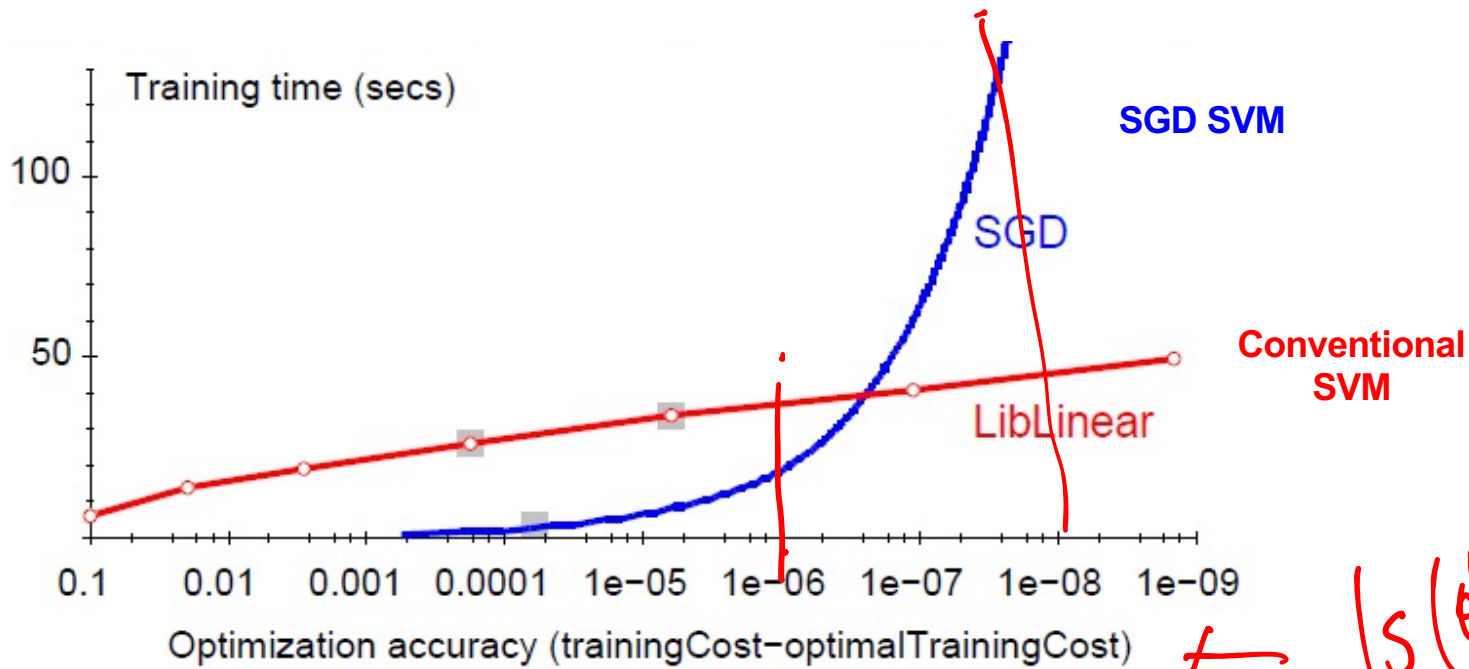
- **Questions:**

- (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 $f(\mathbf{w}, \mathbf{b})$</i> | <i>Test error</i> |
|----------------|----------------------|--|-------------------|
| Standard SVM | 23,642 secs | 0.2275 | 6.02% |
| “Fast SVM” | 66 secs | 0.2278 | 6.03% |
| SGD SVM | 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”

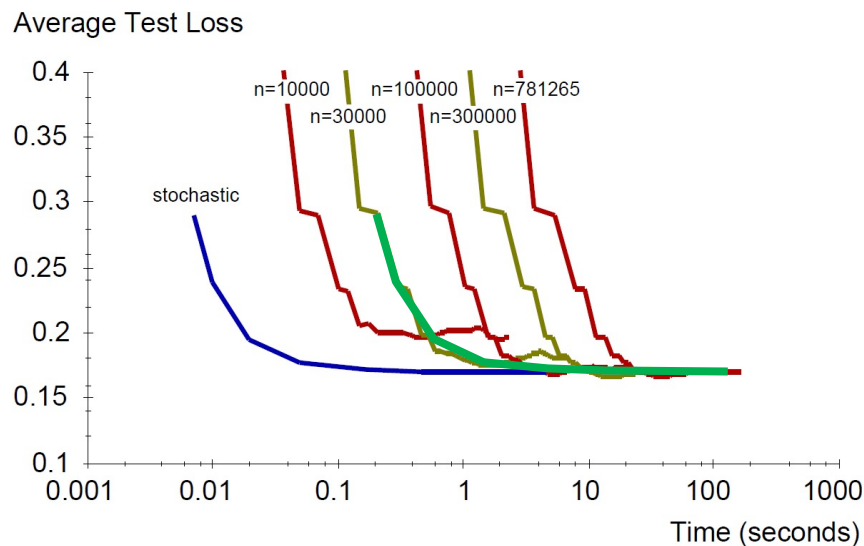


Optimization quality: $| \underline{f(w,b)} - \underline{f(w^{opt}, b^{opt})} |$

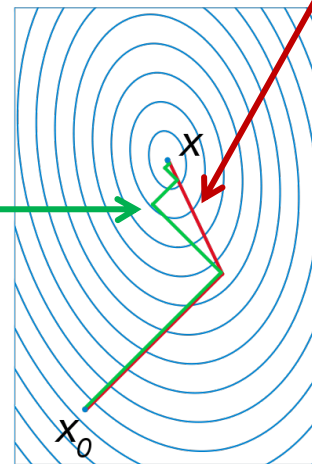
For optimizing $f(w,b)$ within reasonable quality SGD-SVM is super fast

SGD vs. Batch Conjugate Gradient

SGD on full dataset vs. **Conjugate Gradient** on a sample of n training examples



Theory says: Gradient descent converges in linear time k . Conjugate gradient converges in \sqrt{k} . k ... condition number



Bottom line: Doing a simple (but fast) SGD update many times is better than doing a complicated (but slow) CG update a few times

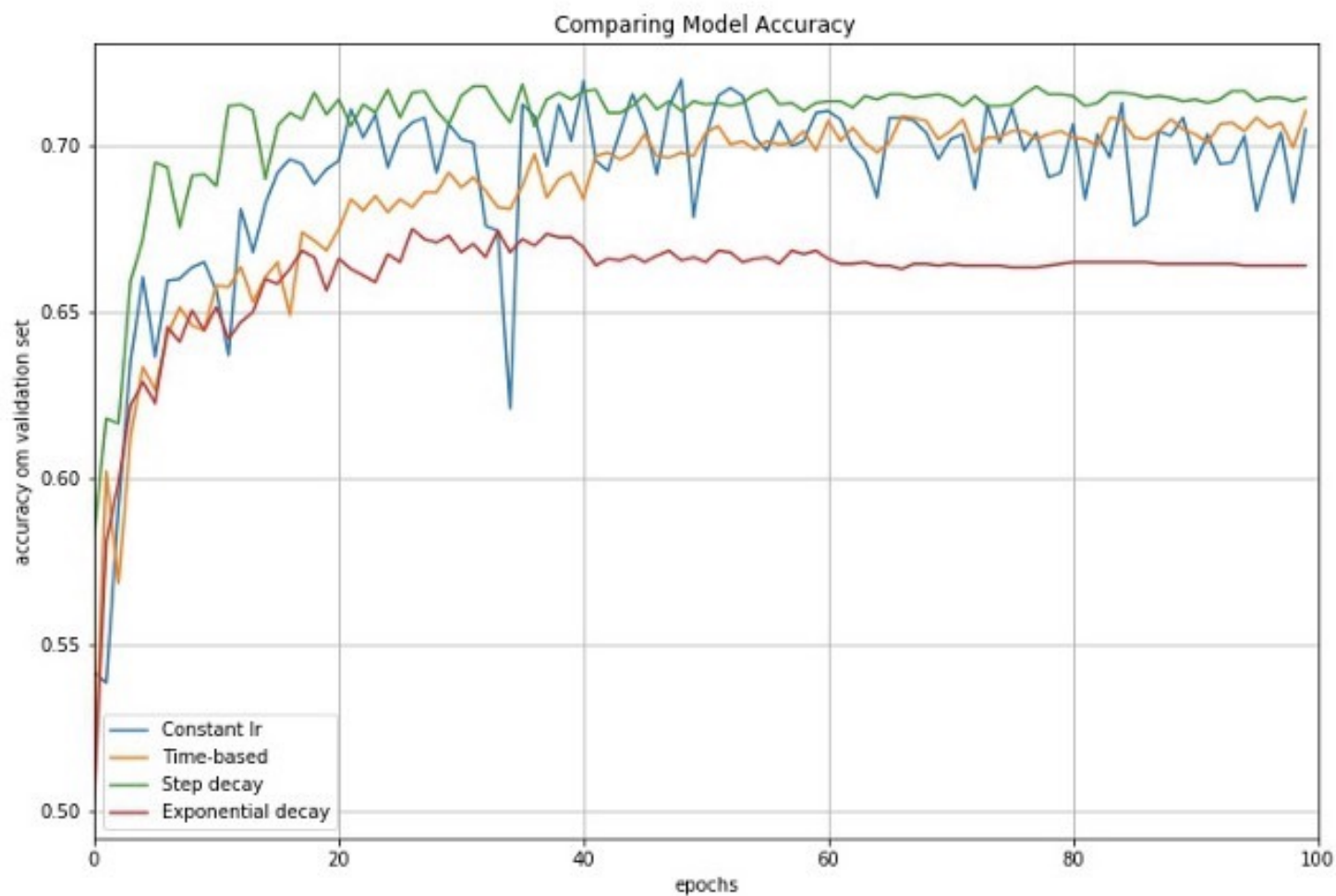
Practical Considerations

- Need to choose learning rate η 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 η :
 - Select a **small subsample**
 - Try various rates η (e.g., 10, 1, 0.1, 0.01, ...)
 - Pick the one that most reduces the cost
 - Use η for next 100k iterations on the full dataset

Learning rate comparison



Practical Considerations

- **Sparse Linear SVM:**

- **Feature vector \mathbf{x}_i is sparse (contains many zeros)**

- Do not do: $\mathbf{x}_i = [0, 0, 0, 1, 0, 0, 0, 0, 5, 0, 0, 0, 0, 0, \dots]$
- But represent \mathbf{x}_i as a sparse vector $\mathbf{x}_i = [(4, 1), (9, 5), \dots]$

- **Can we do the SGD update more efficiently?**

$$\mathbf{w} \leftarrow \mathbf{w} - \eta \left(\mathbf{w} + C \frac{\partial L(\mathbf{x}_i, y_i)}{\partial \mathbf{w}} \right)$$

- **Approximated in 2 steps:**

$$\mathbf{w} \leftarrow \mathbf{w} - \eta C \frac{\partial L(\mathbf{x}_i, y_i)}{\partial \mathbf{w}} \quad \text{cheap: } \mathbf{x}_i \text{ is sparse and so few coordinates } \mathbf{j} \text{ of } \mathbf{w} \text{ will be updated}$$

$$\mathbf{w} \leftarrow \mathbf{w}(1 - \eta) \quad \text{expensive: } \mathbf{w} \text{ is not sparse, all coordinates need to be updated}$$

Practical Considerations

■ **Solution 1:** $\mathbf{w} = \mathbf{s} \cdot \mathbf{v}$

- Represent vector \mathbf{w} as the product of scalar \mathbf{s} and vector \mathbf{v}
- Then the update procedure is:
 - $(1) \mathbf{v} = \mathbf{v} - \eta C \frac{\partial L(x_i, y_i)}{\partial \mathbf{w}}$
 - $(2) \mathbf{s} = \mathbf{s}(1 - \eta)$

• **Solution 2:**

- Perform only step **(1)** for each training example
- Perform step **(2)** with lower frequency and higher η

Two step update procedure:

$$(1) w \leftarrow w - \eta C \frac{\partial L(x_i, y_i)}{\partial w}$$

$$(2) w \leftarrow w(1 - \eta)$$

Practical Considerations

- **Stopping criteria:**

How many iterations of SGD?

- **Early stopping with cross validation**

- Create a validation set
- Monitor cost function on the validation set
- Stop when loss stops decreasing

- **Early stopping**

- Extract two disjoint subsamples **A** and **B** of training data
- Train on **A**, stop by validating on **B**
- Number of epochs is an estimate of k
- Train for k epochs on the full dataset

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

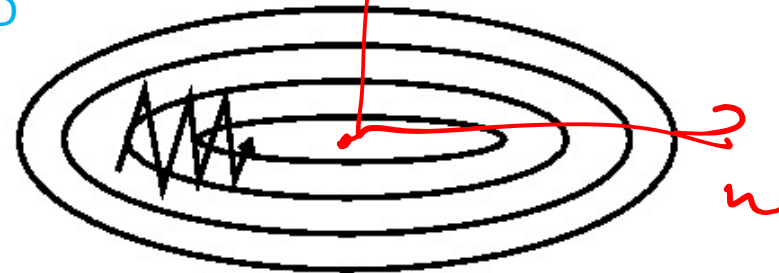
- Idea: Overcome ravine oscillations by momentum

Update:

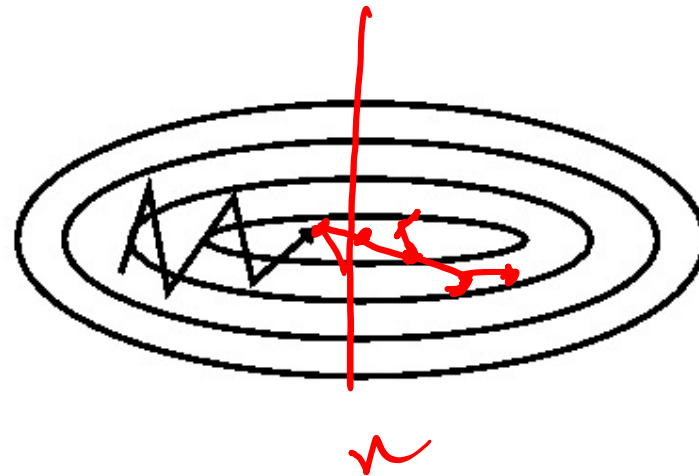
- $v_t = \gamma v_{t-1} + \eta \cdot \nabla_{\theta} J(\theta)$

- $\theta = \theta - v_t$

SGD



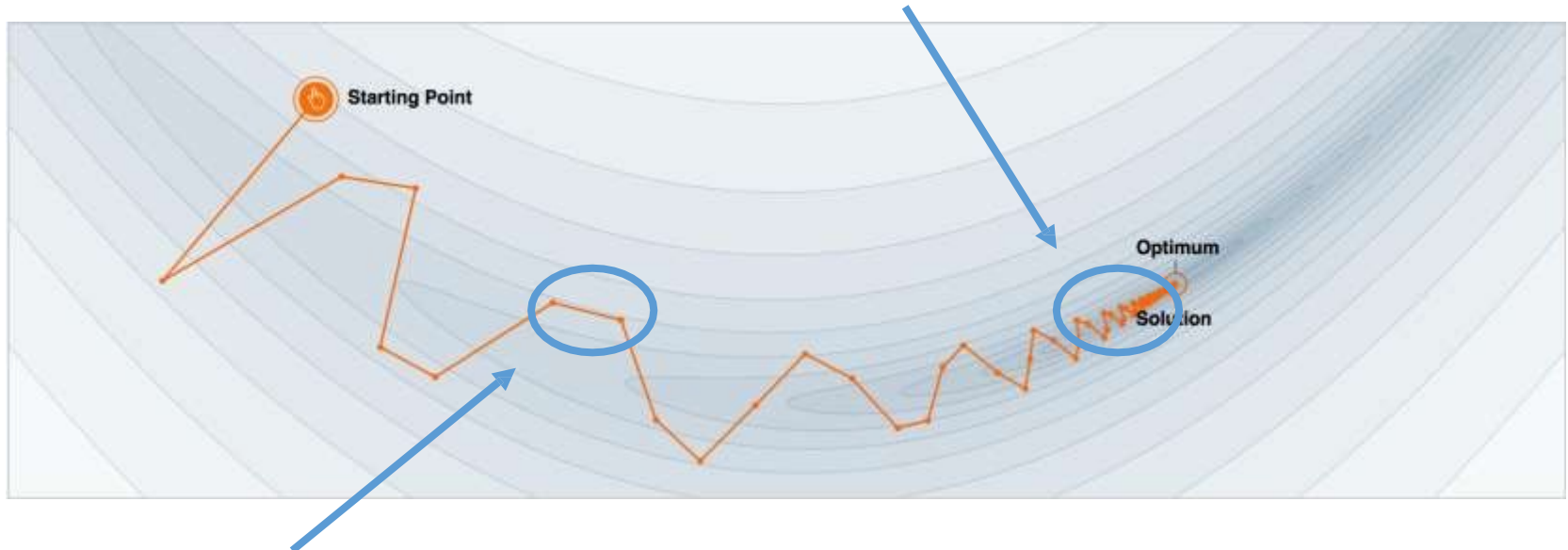
SGD with momentum



$$\frac{\eta}{10}$$

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.**

Demo : <http://distill.pub/2017/momentum/>

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.

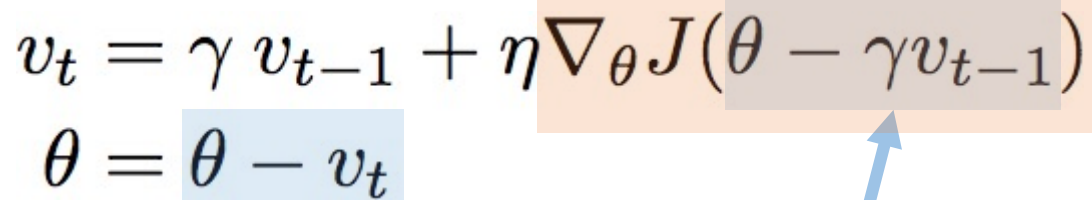
Nesterov accelerated gradient

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


Approximation of the next position of the parameters(predict)

Nesterov accelerated gradient

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

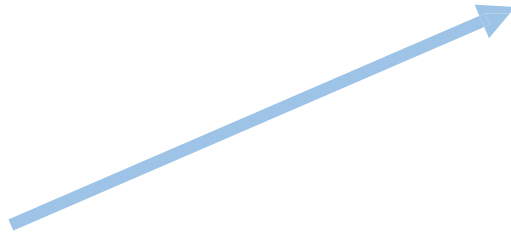


The diagram illustrates the Nesterov accelerated gradient algorithm. It features two equations: $v_t = \gamma v_{t-1} + \eta \nabla_{\theta} J(\theta - \gamma v_{t-1})$ and $\theta = \theta - v_t$. An orange arrow points from the text 'Approximation of the next position of the parameters' gradient(**correction**) to the term $\nabla_{\theta} J(\theta - \gamma v_{t-1})$ in the first equation. A blue arrow points from the text 'Approximation of the next position of the parameters(**predict**)' to the term $\theta - v_t$ in the second equation. The first equation is highlighted with an orange background, and the second equation is highlighted with a blue background.

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

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

Nesterov accelerated gradient



Blue line : predict

Red line : correction

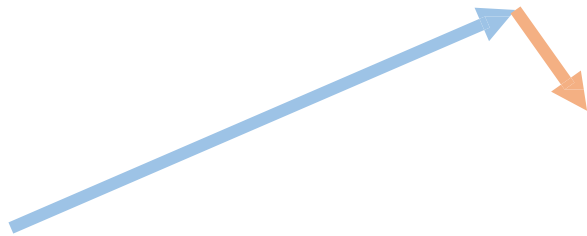
Green line : accumulated gradient

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

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

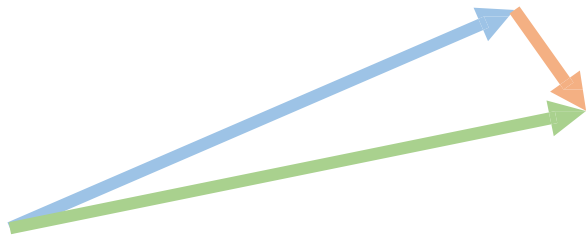
Red line : correction

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

Green line : accumulated gradient

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

Nesterov accelerated gradient



Blue line : predict

Red line : correction

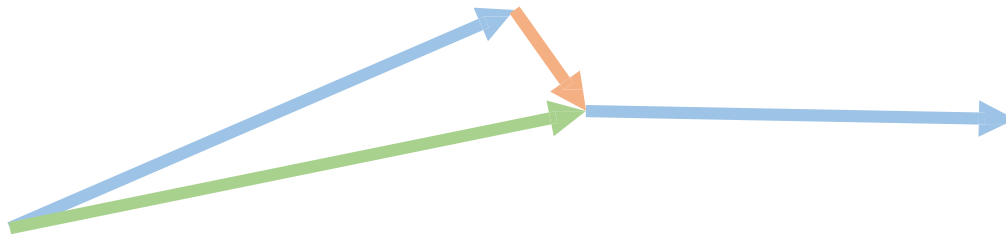
Green line : accumulated gradient

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

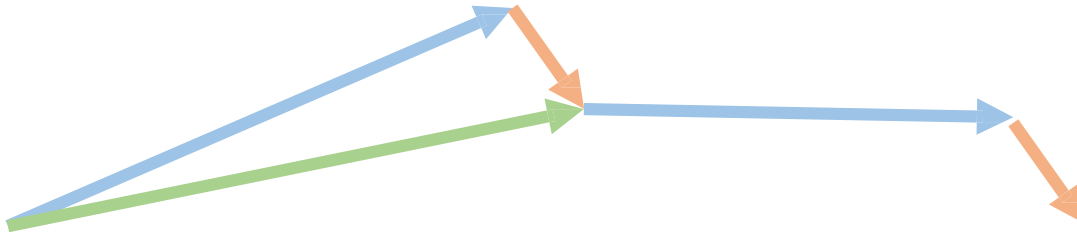
Green line : accumulated gradient

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

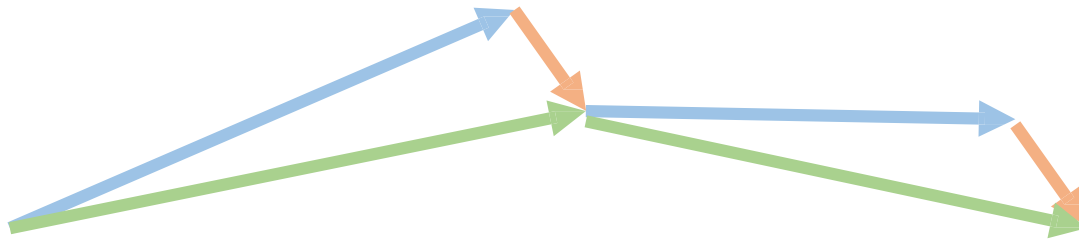
Green line : accumulated gradient

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...?

Adaptive
gradient
methods.

- 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

Newton's

$$\begin{pmatrix} -1 \\ H & g \end{pmatrix}$$

$$f(u, v) = au^2 + bv^2$$
$$\begin{bmatrix} \frac{1}{2a} & 0 \\ 0 & \frac{1}{2b} \end{bmatrix}^+ \begin{bmatrix} g_u \\ g_v \end{bmatrix}$$

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 η for all parameters θ
- Adagrad :
 - It uses a different learning rate for every parameter θ_i at every time step t

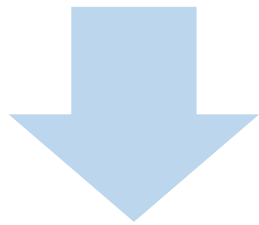
Adagrad

$\nabla^2 L \cdot (g_i)^2$

SGD

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

$$G_t = \begin{matrix} \mathbb{R}^{d \times d} \\ \begin{pmatrix} \text{square with red dot} & \dots & \text{circle} & \dots & \text{circle} \\ \vdots & & \vdots & & \vdots \\ \text{circle} & \dots & \text{square with red dot} & \dots & \text{circle} \\ \vdots & & \vdots & & \vdots \\ \text{circle} & \dots & \text{circle} & \dots & \text{square with red dot} \end{pmatrix} \end{matrix}$$



Adagrad

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

$g_{t,i}$

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

Vectorize

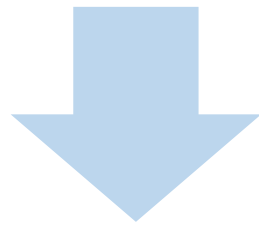
$$\theta_{t+1} = \theta_t - \frac{\eta}{\sqrt{G_t + \epsilon}} \odot g_t$$

Adagrad

SGD

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

$$G_t = \begin{matrix} \mathbb{R}^{d \times d} \\ \begin{pmatrix} \square & \dots & \circ & \dots & \circ \\ \vdots & & \vdots & & \vdots \\ \circ & \dots & \square & \dots & \circ \\ \vdots & & \vdots & & \vdots \\ \circ & \dots & \circ & \dots & \square \end{pmatrix} \end{matrix}$$



Adagrad modifies the general learning rate η based on the **past gradients that have been computed for θ_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

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

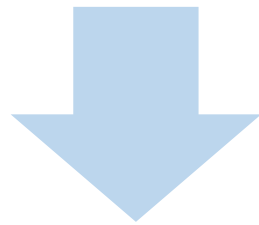
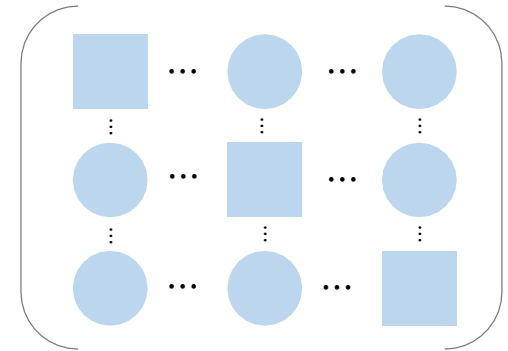
Adagrad

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 θ_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

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

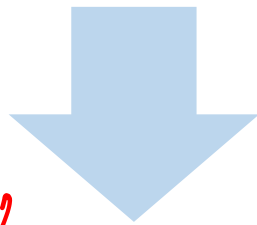
Adagrad

SGD

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

$$G_t = \begin{matrix} \mathbb{R}^{d \times d} \\ \begin{pmatrix} \square & \dots & \circ & \dots & \circ \\ \vdots & & \vdots & & \vdots \\ \circ & \dots & \square & \dots & \circ \\ \vdots & & \vdots & & \vdots \\ \circ & \dots & \circ & \dots & \square \end{pmatrix} \end{matrix}$$

Ug!



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

Adagrad

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




$\sqrt{g_{t,i}}$

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

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


Adadelata : extension of Adagrad

- Adadelata is an extension of Adagrad.
- Adagrad :
 - It accumulate all past squared gradients.
- Adadelata :
 - 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.

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

- $E[g^2]_t$: The running average at timestep t .
- γ : 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 η in the previous update rule with $RMS[\Delta\theta]_{t-1}$ finally yields the Adadelta update rule:

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

- 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.

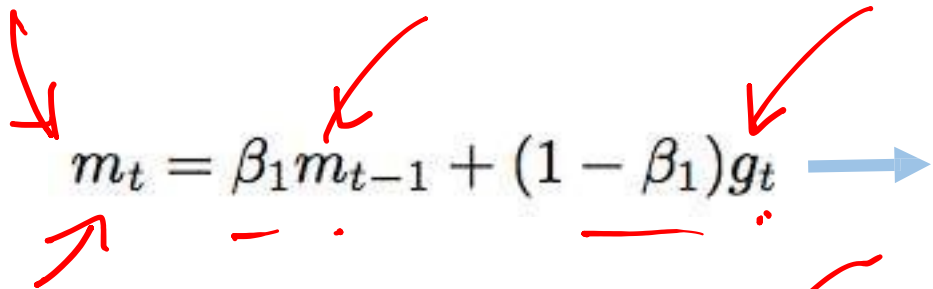
RMSprop

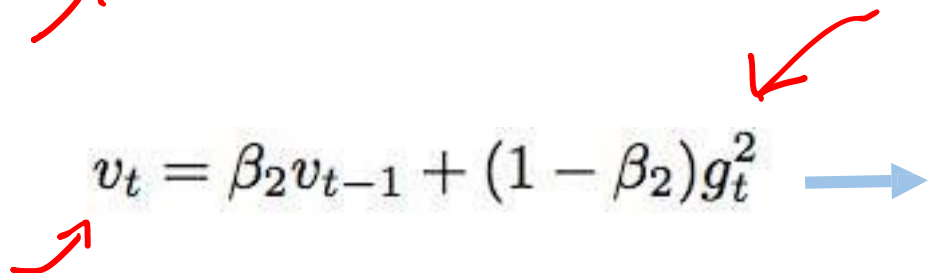
$$E[g^2]_t = 0.9E[g^2]_{t-1} + 0.1g_t^2$$
$$\theta_{t+1} = \theta_t - \frac{\eta}{\sqrt{E[g^2]_t + \epsilon}} g_t$$

Hinton suggests γ to be set to 0.9, while a good default value for the learning rate η 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.


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


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

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. β_1 and β_2 are close to 1).
- Counteracting these biases in Adam

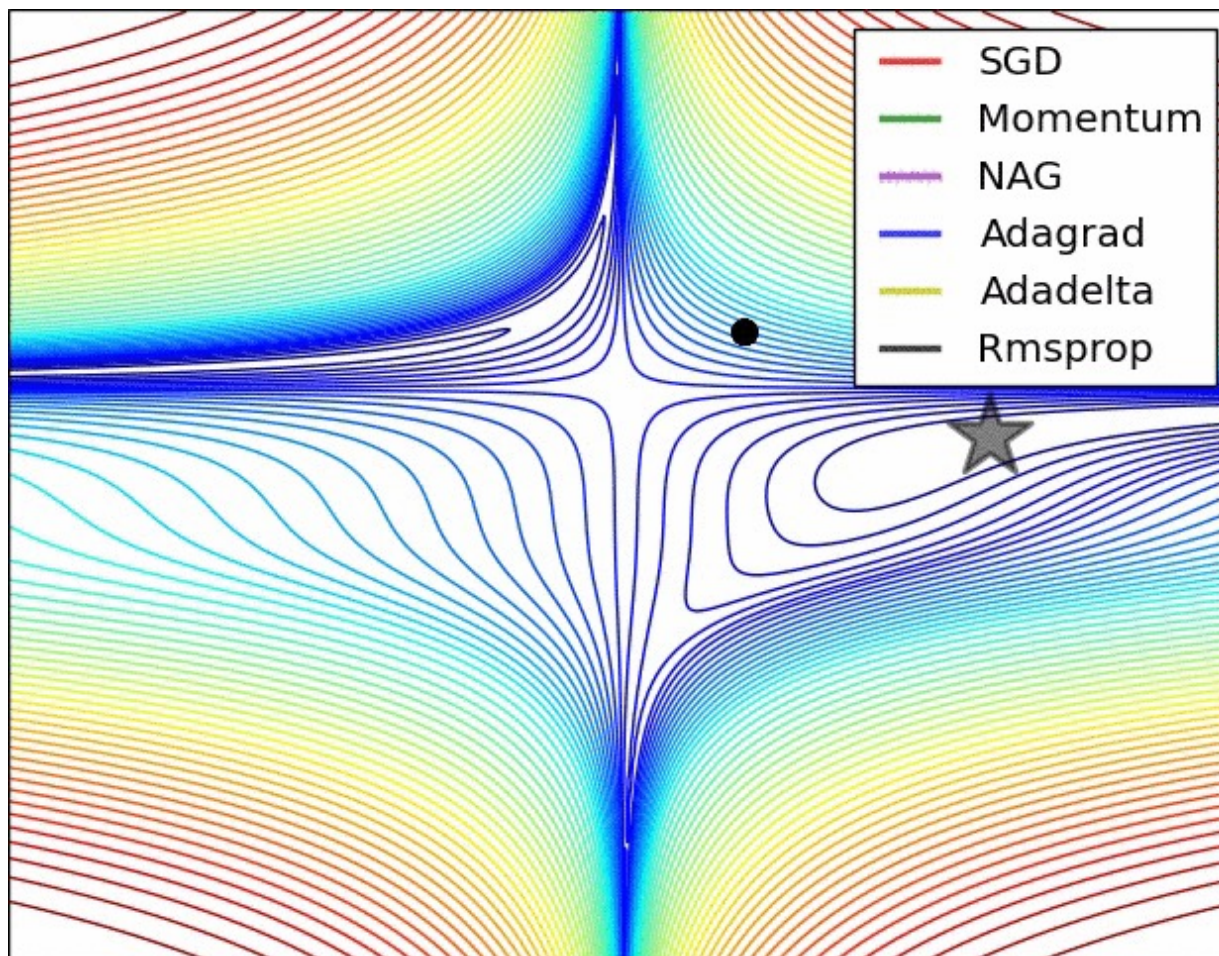
$$\hat{m}_t = \frac{m_t}{1 - \beta_1^t}$$
$$\hat{v}_t = \frac{v_t}{1 - \beta_2^t}$$

Adam

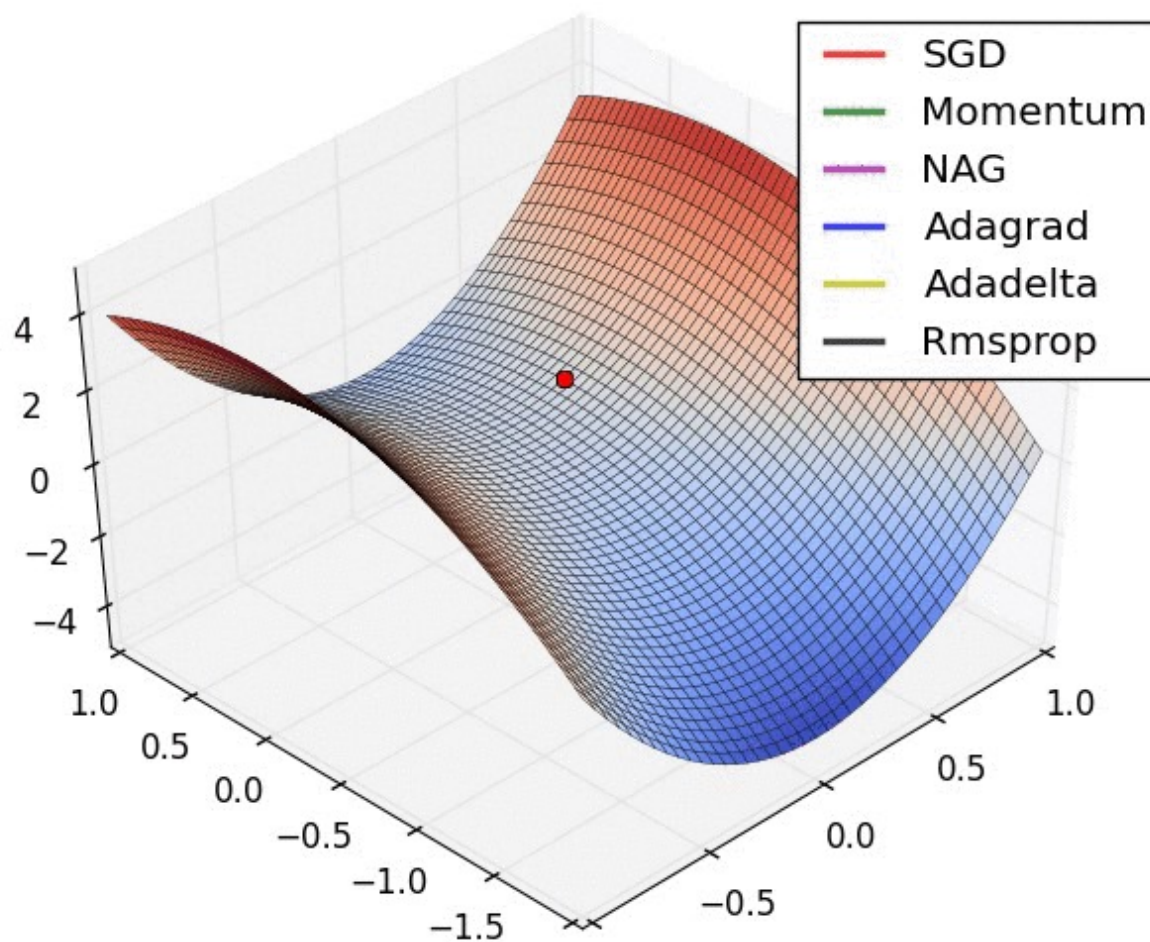
$$\theta_{t+1} = \theta_t - \frac{\eta}{\sqrt{\hat{v}_t} + \epsilon} \hat{m}_t$$

Note : default values of 0.9 for β_1 , 0.999 for β_2 , and 10^{-8} for ϵ

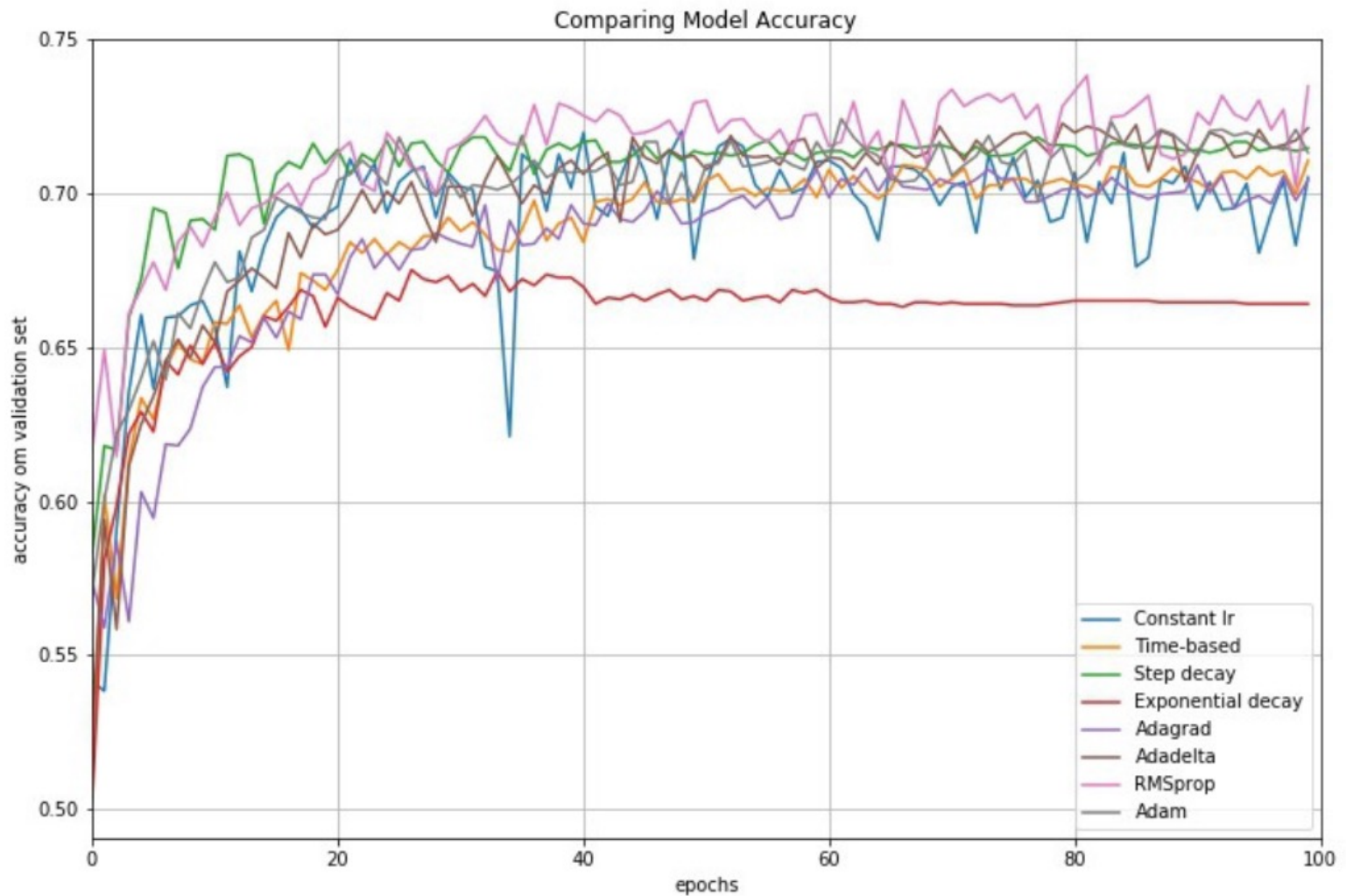
- SGD - Simple
 - RMSprop
 - Adam
- # Visualization
- Bias correction.



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.

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>