

A decorative graphic consisting of two horizontal blue lines and two vertical blue lines. The top horizontal line starts from the left edge and ends with a small blue circle. The bottom horizontal line starts from the right edge and ends with a small blue circle. The two vertical lines intersect these horizontal lines, forming a frame-like structure.

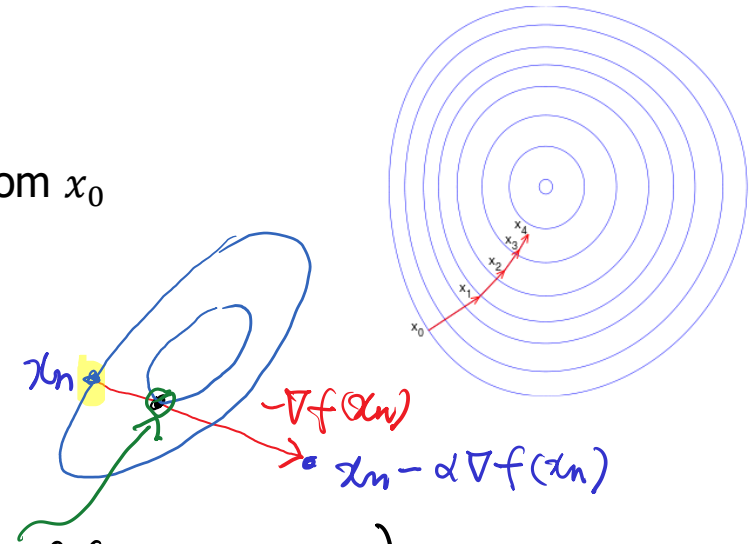
Optimization for ML

KAIST ISE
Hayong Shin

Gradient descent (GD)

- Local optimization problem

- find local minimum point of $f(x)$ starting from x_0
- GD : $x_{t+1} \leftarrow x_t - \alpha_t \nabla f(x_t)$
 - $\nabla f(x_t)$: **steepest descent** direction
 - TAQ. Is $\nabla f(x_t)$ the best direction?
- α_t : step size
 - How to set α_t ? \rightarrow non-trivial



① Exact line search $\alpha_n = \underset{\alpha}{\operatorname{argmin}} f(x_n - \alpha \nabla f(x_n))$

② Backtracking line search (Armijo rule)

$\alpha = s$

while $f(x_n - \alpha \nabla f(x_n)) \geq f(x_n)$, $\alpha \leftarrow \alpha/2$

i.e. find the first $\alpha \in [s, \frac{s}{2}, \frac{s}{4}, \dots]$ s.t. $f(x_n - \alpha \nabla f(x_n)) < f(x_n)$

③ 2nd order methods

Use $\nabla^2 f$ (Hessian) to determine α

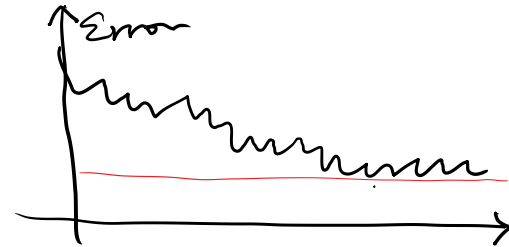
- $\nabla^2 f$ requires too much efforts to estimate

- approximate $\nabla^2 f$: BFGS, L-BFGS

Gradient descent : step size

- ④ constant $\alpha_n = \alpha$
- bigger α : may diverge
 - smaller α : slow

— may not converge in SGD

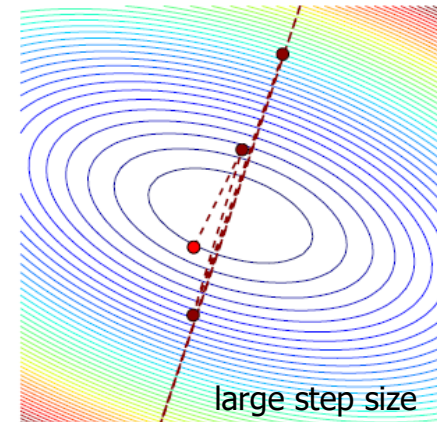
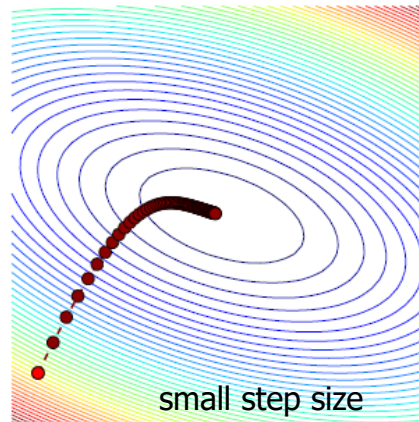
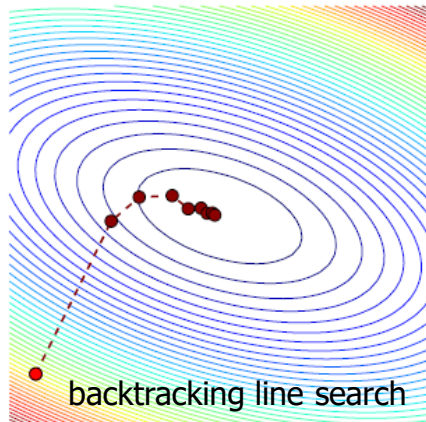


- ⑤ diminishing $\alpha_n \rightarrow 0$

— Essential for SGD

— Robbins - Monro condition

$$\sum_n \alpha_n = \infty, \quad \sum_n \alpha_n^2 < \infty$$



Stochastic Gradient

Local optimization problem : find x s.t. $\nabla f(x) = 0$

Let $g(x) = \nabla f(x) + \varepsilon$ $E[\varepsilon] = 0$

If we are to find (local) minimum point of $f(x)$,

we can assume $f(x)$ is (locally) convex, $\Rightarrow \nabla f(x)$: increasing

$\sum_{n=1}^{\infty} \alpha_n \rightarrow \infty$	and	$\sum_{n=1}^{\infty} \alpha_n^2 < \infty$
①		②

$\Rightarrow x_{n+1} = x_n - \alpha_n g(x_n),$

where $g(x_n)$ is an unbiased estimator of $\nabla f(x_n)$

→ stochastic gradient

* $f(x) = \sum_{i=1}^N f_i(x)$: $f_i(x)$: loss function for i^{th} data, $i \in D = \{1, \dots, N\}$

$\nabla f(x) = \sum_{i \in D} \nabla f_i(x)$: (full) batch gradient

let $M \subset D$ be a randomly chosen subset of size m ($|M| = m$)

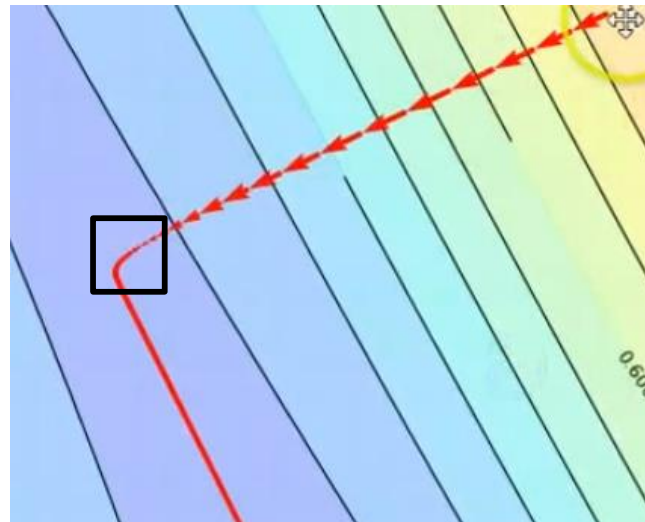
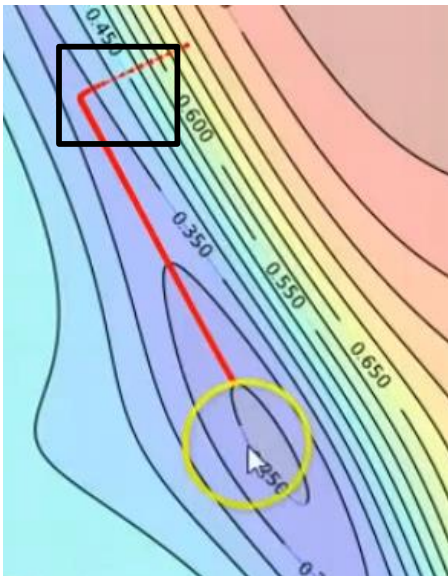
$g(x) = \sum_{i \in M} \nabla f_i(x) \Rightarrow E_M[g(x)] = \nabla f(x)$: unbiased

\Rightarrow stochastic gradient using minibatch

\Rightarrow extreme case : $m=1$ small $m \Rightarrow \text{Var}[g(x)]$ gets bigger

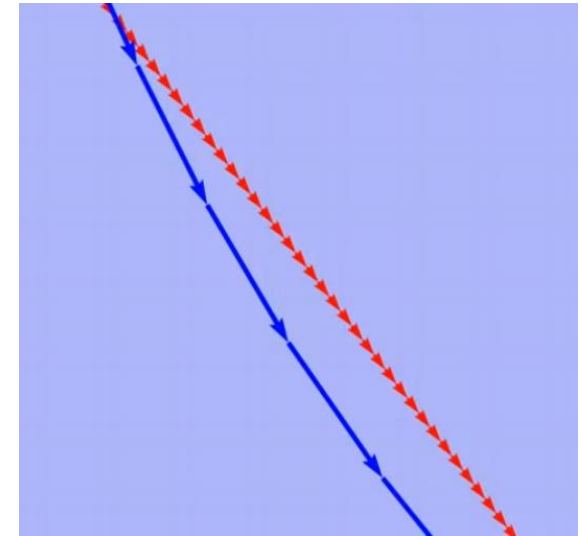
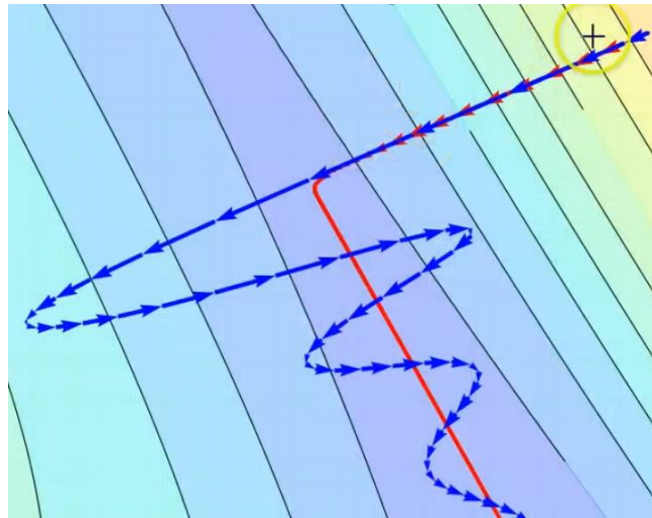
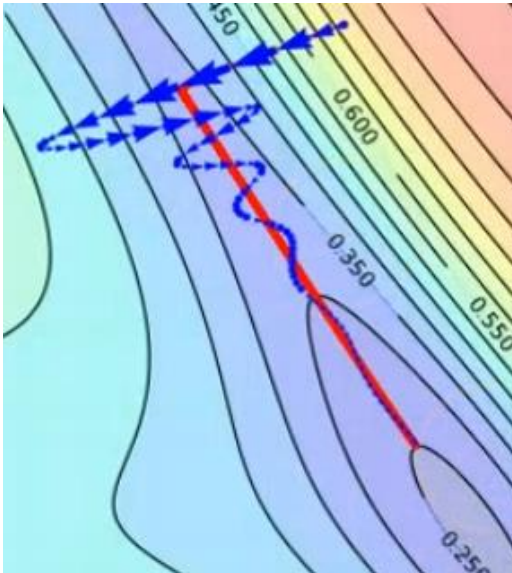
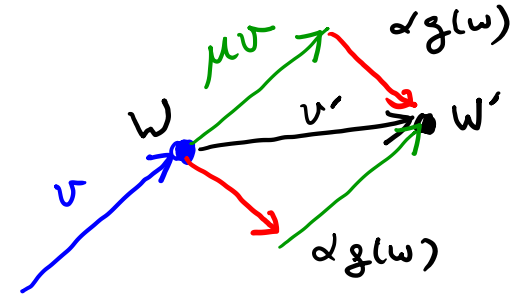
Enhancements in Optimization

- GD (Plain vanilla version)
 - $l(\mathbf{w})$: cost/loss/error function
 - \mathbf{g} : (an estimate of) $-\nabla l(\mathbf{w})$ (gradient at \mathbf{w})
 - α : learning rate, step size
 - $\mathbf{v} = \alpha \mathbf{g}$
 - $\mathbf{w} \leftarrow \mathbf{w} + \mathbf{v}$
 - very hard to set right step size

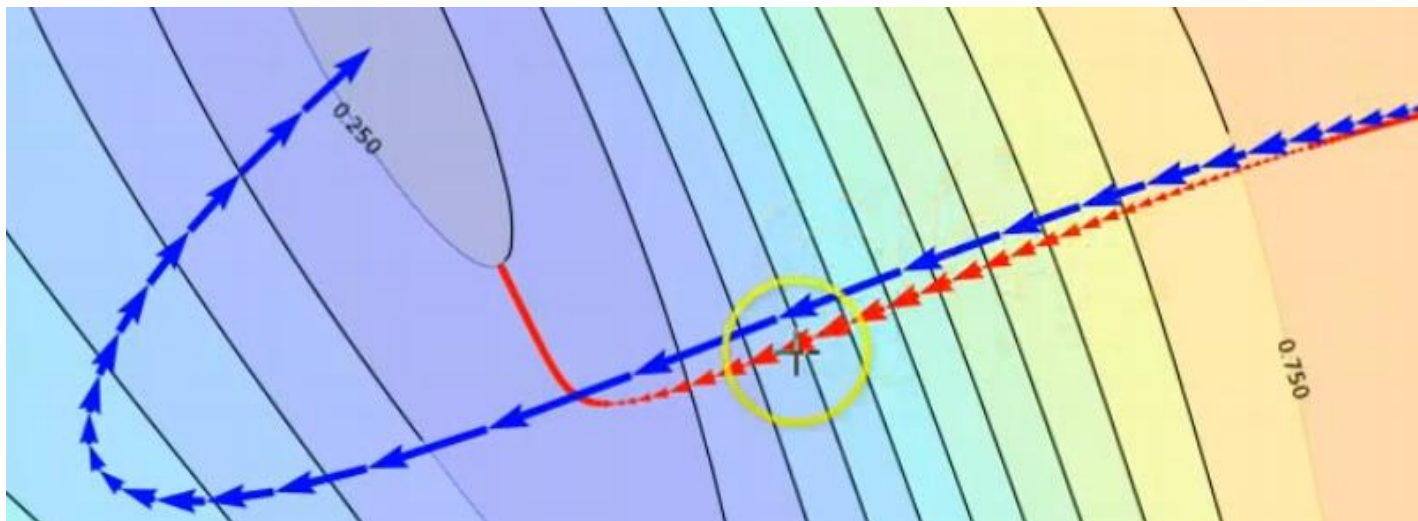
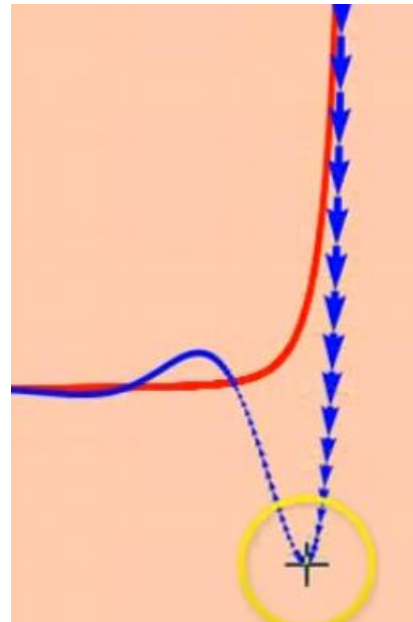
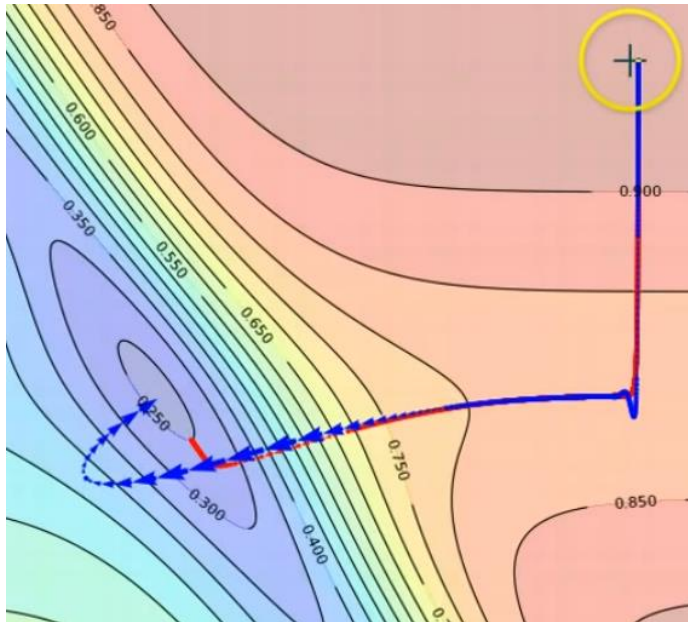


Enhancing GD : Momentum

- Momentum : keep velocity of (heavy) ball
 - $v \leftarrow \mu v + \alpha g$: velocity vector (initially 0)
 - μ : decaying factor, e.g. $\mu = 0.9$
 - $w \leftarrow w + v$
 - if g is constant, $w = \alpha[g + \mu g + \mu^2 g + \dots] = \frac{\alpha}{1-\mu} g$
 - if $\mu = 0.9$, the ball will move about 10 times faster



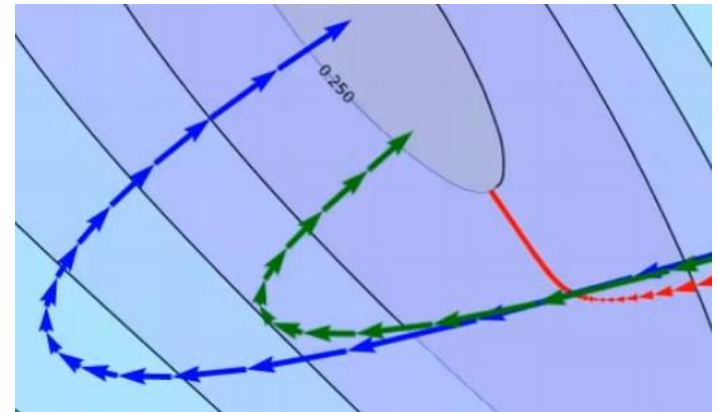
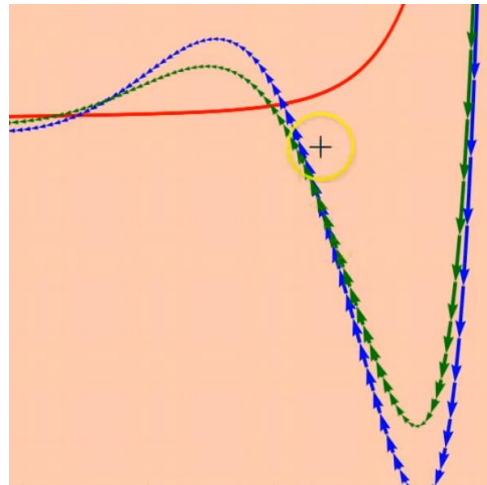
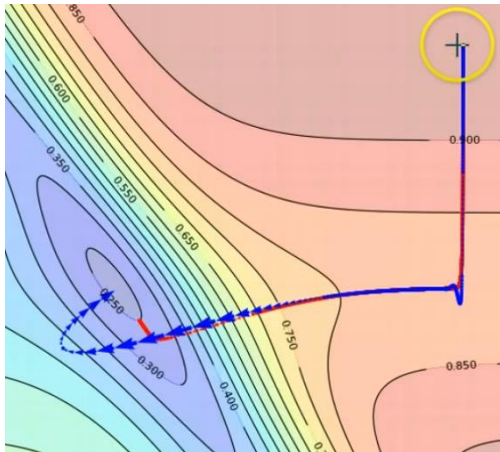
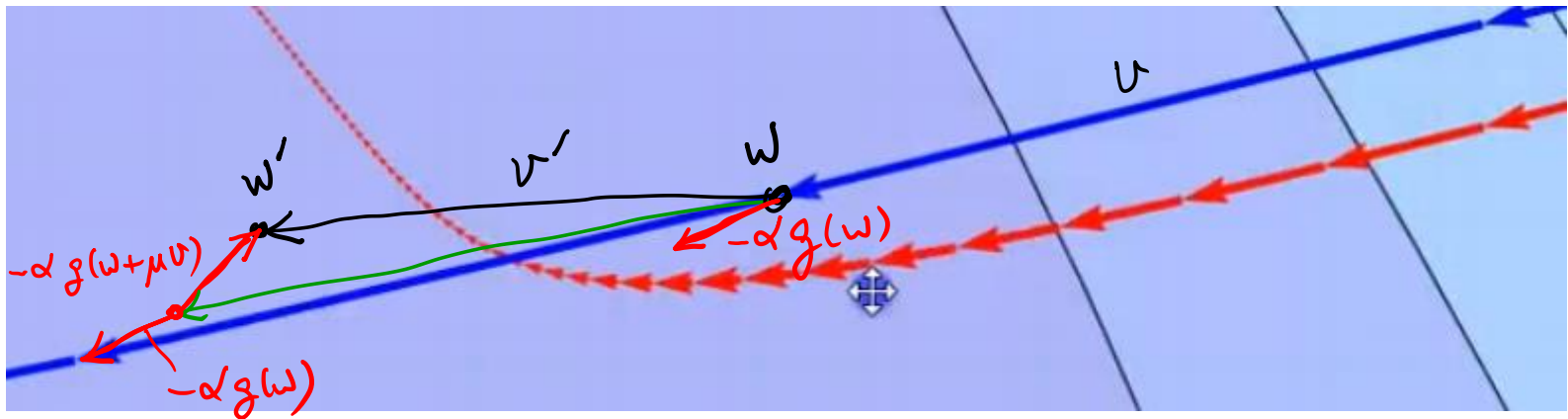
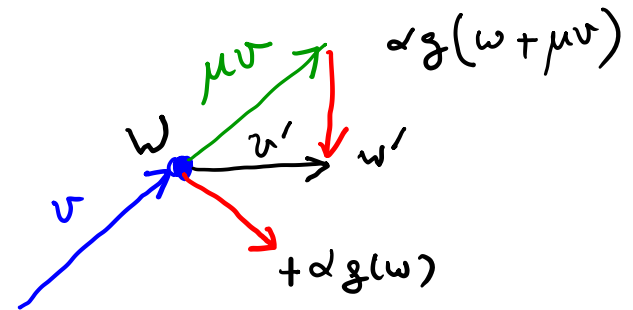
Enhancing GD : Momentum



Nesterov Accelerated Gradient (Sutskever, 2013)

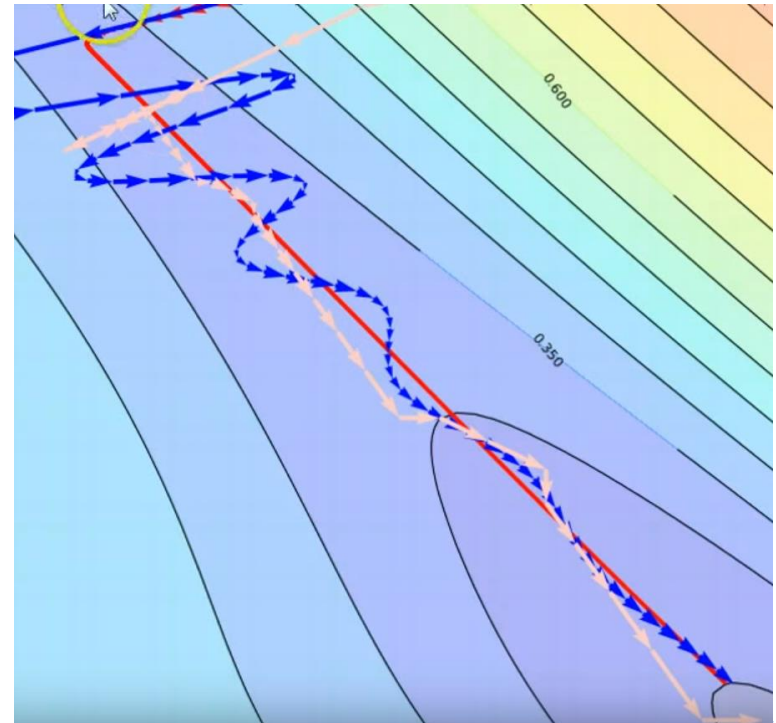
- NAG : Gamble, then correct

- $\mathbf{g} = -\nabla l(\mathbf{w} + \mu \mathbf{v})$
- $\mathbf{v} \leftarrow \mu \mathbf{v} + \alpha \mathbf{g}$
- $\mathbf{w} \leftarrow \mathbf{w} + \mathbf{v}$



RProp (Resilient Propagation, Riedmiller & Braun, 1993)

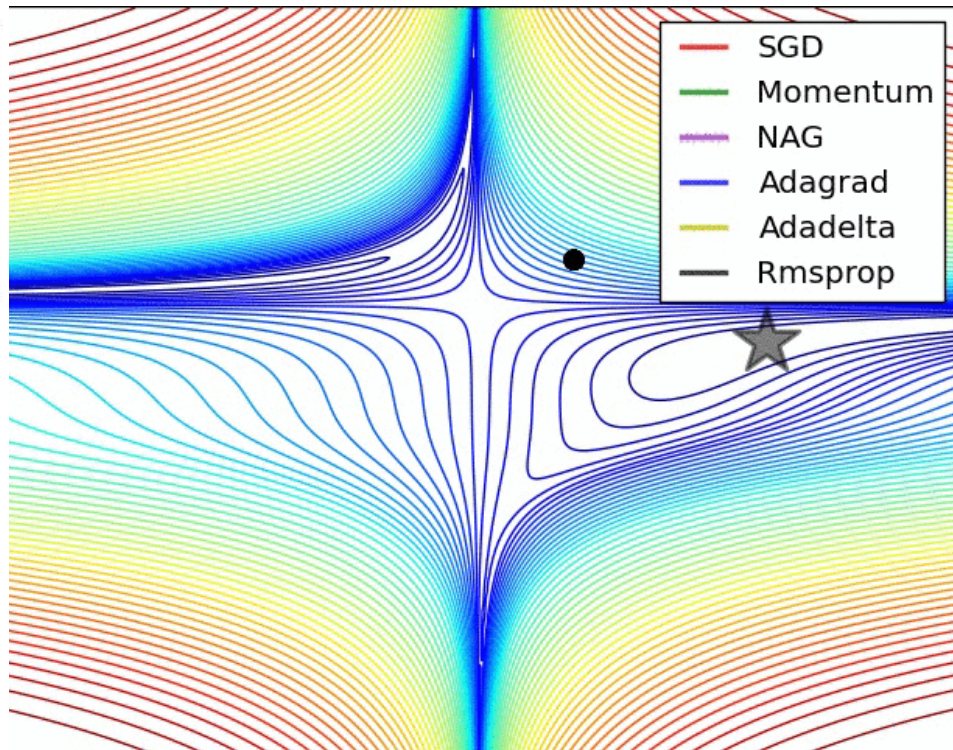
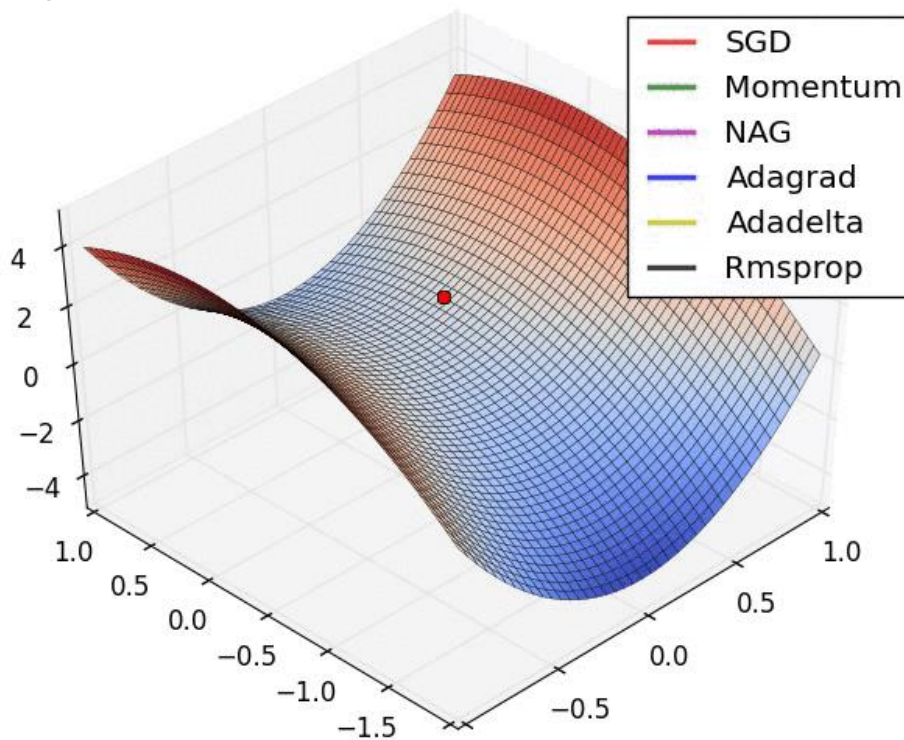
- Ignore gradient size, use only sign of gradient for each component
 - g_k : current gradient component k
 - $c_k = \text{sign}(g_k)$
 - p_k : previous sign, initially 0
- Keep step size α_k per component
 - if $p_k \cdot c_k > 0$ (same sign)
 - $\alpha_k \leftarrow 1.2$ (increase step size)
 - $\Delta w_k = \alpha_k \cdot c_k$
 - $p_k = c_k$
 - if $p_k \cdot c_k < 0$ (sign changes)
 - $\alpha_k \leftarrow \alpha_k / 2$ (decrease step size)
 - $\Delta w_k = 0$
 - $p_k = 0$
 - else (sign changed)
 - $\Delta w_k = \alpha_k \cdot c_k$
 - $p_k = c_k$



Other methods

- So far, we covered : GD, SGD, Momentum, NAG, RProp
- Recently, there are many other attempts on adaptive gradients
 - AdaGrad (Duchi et al, JMLR 2011)
 - AdaDelta (Zeiler, 2012)
 - RMSProp (Hinton, 2012)

$$p = \sum g^2, w \leftarrow w + \alpha \frac{g}{\sqrt{p}}$$



ADAM (Kingma & Ba, 2015)

- ADAM : ADaptive Moment estimation

- Combines benefits of AdaGrad (Duchi 2011), RMSProp (Hinton 2012)

- Note element-wise operations below : \mathbf{g}^2 , $\sqrt{\hat{\mathbf{p}}}$, and $\frac{\hat{\mathbf{v}}}{\sqrt{\hat{\mathbf{p}}} + \epsilon}$

- Adaptive step size per component

- Initialize : $\mathbf{v} = \mathbf{0}, \mathbf{p} = \mathbf{0}, \beta_1 = 0.9, \beta_2 = 0.999, \epsilon = 0.000001$

$\mathbf{g} = -\nabla l(\mathbf{w})$: opposite of gradient

$\mathbf{v} \leftarrow \beta_1 \mathbf{v} + (1 - \beta_1) \mathbf{g}$: update velocity (1st moment) estimate

- \mathbf{v} : momentum + gradient ("signal")
- bigger for components with consistent sign, smaller for sign changing component

$\mathbf{p} \leftarrow \beta_2 \mathbf{p} + (1 - \beta_2) \mathbf{g}^2$: update 2nd moment estimate

- \mathbf{p} : gradient variance ("noise")
- penalize components with bigger variance
- compensate scale difference between components

$\hat{\mathbf{v}} \leftarrow \frac{\mathbf{v}}{1 - \beta_1^t}, \hat{\mathbf{p}} \leftarrow \frac{\mathbf{p}}{1 - \beta_2^t}$: initial size bias correction

- \mathbf{v} and \mathbf{p} are small at the beginning
- $\beta_1^t = \beta_1 \beta_1^{t-1} \rightarrow 0$ and $\beta_2^t = \beta_2 \beta_2^{t-1} \rightarrow 0$

$\mathbf{w} \leftarrow \mathbf{w} + \alpha \frac{\hat{\mathbf{v}}}{\sqrt{\hat{\mathbf{p}}} + \epsilon}$: correct velocity with 2nd moment

Comparison

