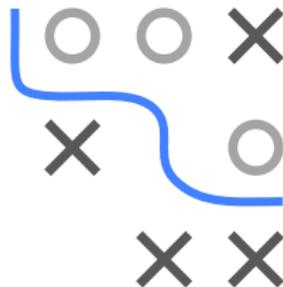


Optimization in Machine Learning

First order methods

Gradient descent



Learning goals

- Iterative Descent / Line Search
- Descent directions
- GD
- ERM with GD
- Pseudoresiduals

ITERATIVE DESCENT

Let f be the height of a mountain depending on the geographic coordinates (x_1, x_2)

$$f : \mathbb{R}^2 \rightarrow \mathbb{R}, \quad f(x_1, x_2) = y.$$

Goal: Reach the valley

$$\arg \min_{\mathbf{x}} f(\mathbf{x})$$

Central idea: Repeat

- 1 At current location $\mathbf{x} \in \mathbb{R}^d$ search for **descent direction** $\mathbf{d} \in \mathbb{R}^d$
- 2 Move along \mathbf{d} until f „sufficiently“ reduces (**step size control**) and update location



"Walking down the hill, towards the valley."

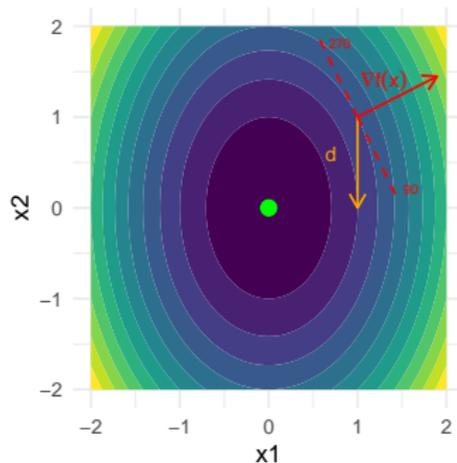
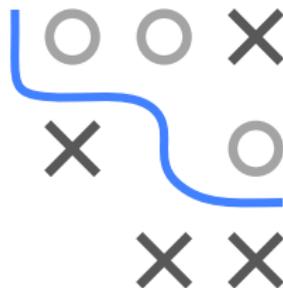


ITERATIVE DESCENT

Let $f : \mathbb{R}^d \rightarrow \mathbb{R}$ continuously differentiable.

Definition: $\mathbf{d} \in \mathbb{R}^d$ is a **descent direction** in \mathbf{x} if

$$D_{\mathbf{d}}f(\mathbf{x}) = \nabla f(\mathbf{x})^T \mathbf{d} < 0 \text{ (neg. directional derivative)}$$



Angle between $\nabla f(\mathbf{x})$ and \mathbf{d} must be $\in (90^\circ, 270^\circ)$.

ITERATIVE DESCENT

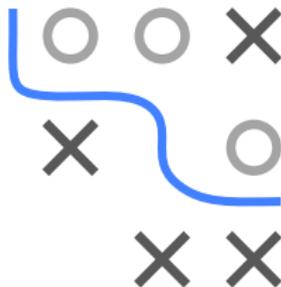
Algorithm Iterative Descent / Line search

- 1: Starting point $\mathbf{x}^{[0]} \in \mathbb{R}^d$
 - 2: **while** Stopping criterion not met **do**
 - 3: Calculate a descent direction $\mathbf{d}^{[t]}$ for current $\mathbf{x}^{[t]}$
 - 4: Find $\alpha^{[t]}$ s.t. $f(\mathbf{x}^{[t+1]}) < f(\mathbf{x}^{[t]})$ for $\mathbf{x}^{[t+1]} = \mathbf{x}^{[t]} + \alpha^{[t]}\mathbf{d}^{[t]}$.
 - 5: Update $\mathbf{x}^{[t+1]} = \mathbf{x}^{[t]} + \alpha^{[t]}\mathbf{d}^{[t]}$
 - 6: **end while**
-

NB: Terminology is sometimes ambiguous: "line search" can refer to Step 4 (selecting the step size that decreases $f(\mathbf{x})$) and can mean umbrella term for iterative descent algorithms.

Key questions:

- How to choose $\mathbf{d}^{[t]}$ (now)
- How to choose $\alpha^{[t]}$ (later)

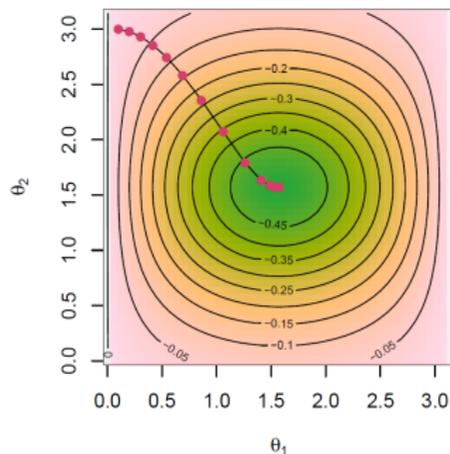
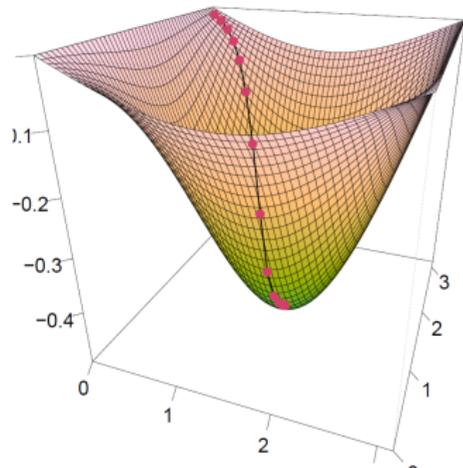
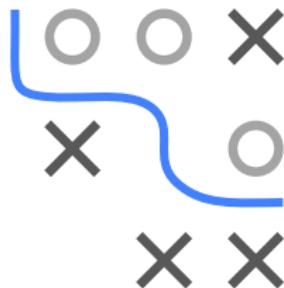


GRADIENT DESCENT

Properties of gradient:

- $\nabla f(\mathbf{x})$: direction of greatest increase
- $-\nabla f(\mathbf{x})$: direction of greatest decrease

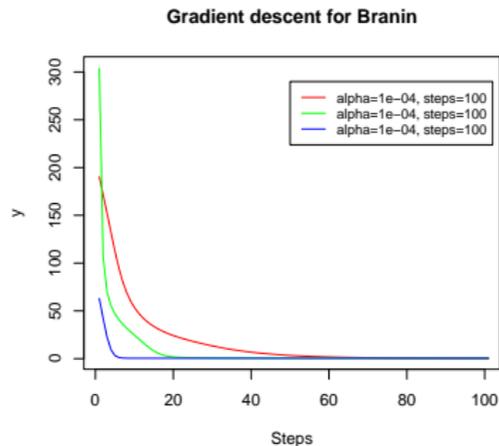
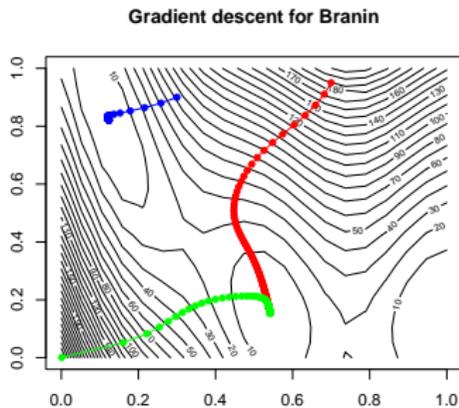
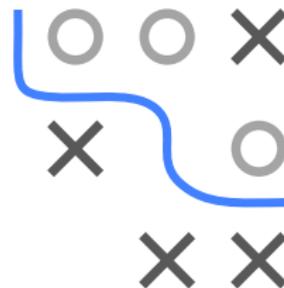
Using $\mathbf{d} = -\nabla f(\mathbf{x})$ is called **gradient descent**.



GD for $f(x_1, x_2) = -\sin(x_1) \cdot \frac{1}{2\pi} \exp((x_2 - \pi/2)^2)$ with sensibly chosen step size $\alpha^{[t]}$.

GD AND MULTIMODAL FUNCTIONS

Outcome will depend on start point.



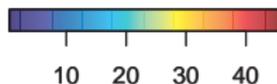
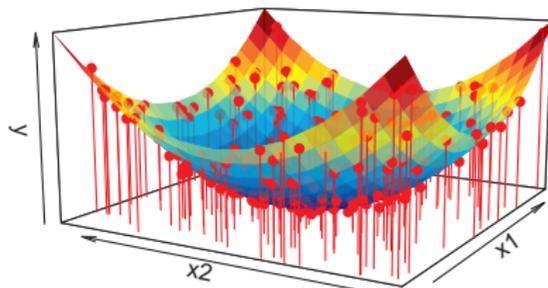
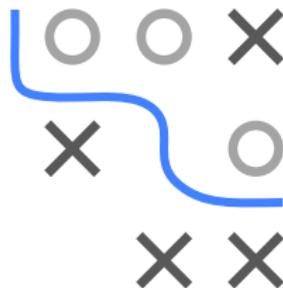
100 iters of GD with const $\alpha = 10^{-4}$.

ERM FOR NN WITH GD

Let $\mathcal{D} = ((\mathbf{x}^{(1)}, y^{(1)}), \dots, (\mathbf{x}^{(n)}, y^{(n)}))$, with $y = x_1^2 + x_2^2$ and minimize

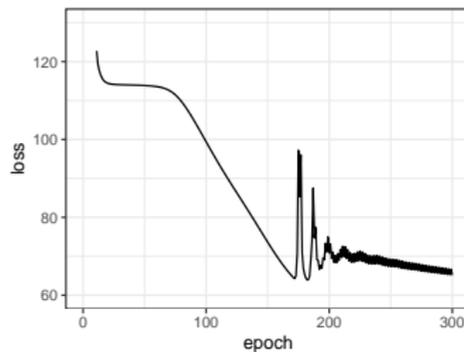
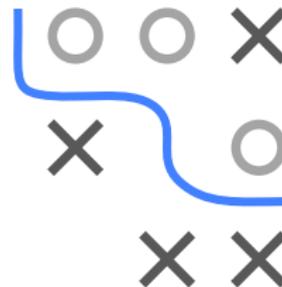
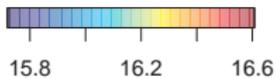
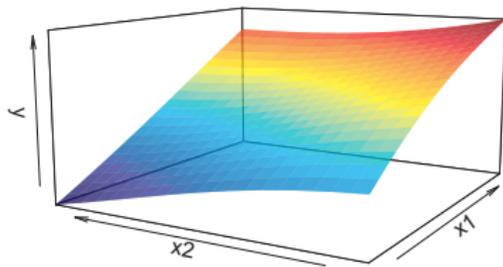
$$\mathcal{R}_{\text{emp}}(\boldsymbol{\theta}) = \sum_{i=1}^n \left(f(\mathbf{x} | \boldsymbol{\theta}) - y^{(i)} \right)^2$$

where $f(\mathbf{x} | \boldsymbol{\theta})$ is a neural network with 2 hidden layers (2 units each).



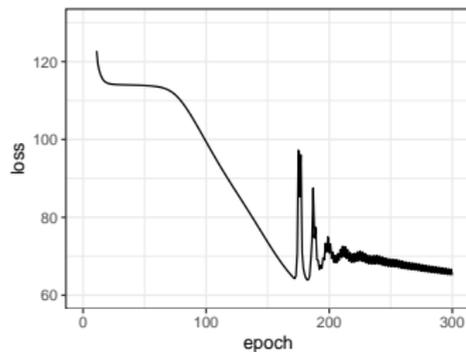
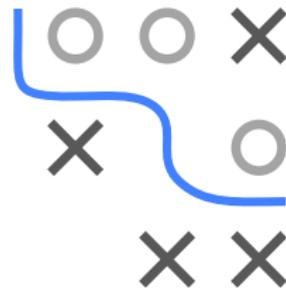
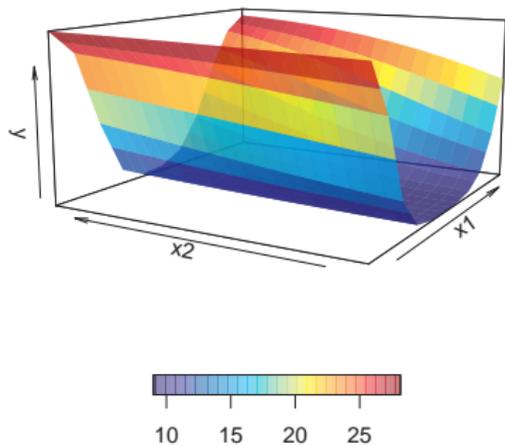
ERM FOR NN WITH GD

After 10 iters of GD:



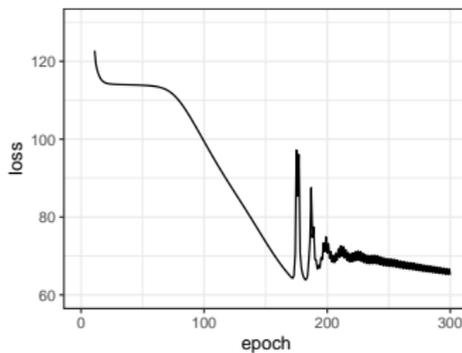
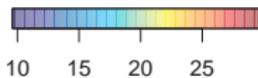
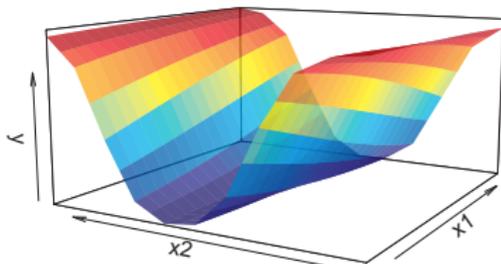
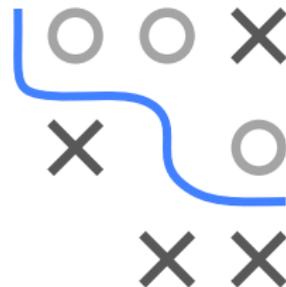
ERM FOR NN WITH GD

After 100 iters of GD:



ERM FOR NN WITH GD

After 300 iters of GD (note the zig-zag-behavior after iter. 200):



GD FOR ERM: PSEUDORESIDUALS

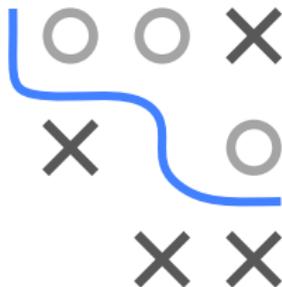
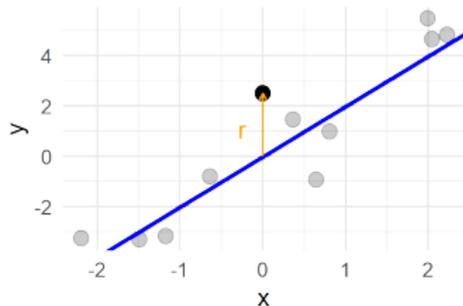
Gradient for ERM problem:

$$-\frac{\partial \mathcal{R}_{\text{emp}}(\theta)}{\partial \theta} = -\frac{\partial \sum_{i=1}^n L(y^{(i)}, f(\mathbf{x}^{(i)} | \theta))}{\partial \theta} = -\sum_{i=1}^n \underbrace{\frac{\partial L(y^{(i)}, f(\mathbf{x}^{(i)}))}{\partial f}}_{\text{pseudo residual } \tilde{r}^{(i)}(f)} \frac{\partial f(\mathbf{x}^{(i)} | \theta)}{\partial \theta}$$

- **pseudo residuals** tell us how to distort $f(\mathbf{x}^{(i)})$ to achieve greatest decrease of $L(y^{(i)}, f(\mathbf{x}^{(i)}))$ (best pointwise update)
- $\frac{\partial f(\mathbf{x}^{(i)} | \theta)}{\partial \theta}$ tells us how to modify θ accordingly and wiggle model output
- GD step sums up these modifications across all observations i

NB: Pseudo-residuals $\tilde{r}(f)$ match usual residuals for L2 loss:

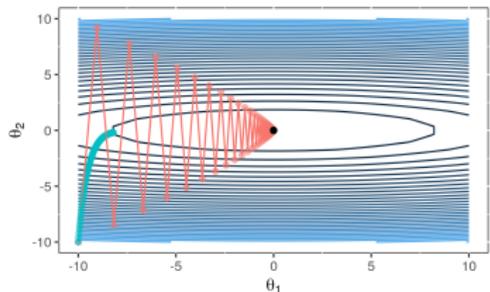
$$\begin{aligned} \frac{\partial L(y, f(\mathbf{x}))}{\partial f} &= \frac{1}{2} \frac{\partial (y - f)^2}{\partial f} \Bigg|_{f=f(\mathbf{x})} \\ &= y - f(\mathbf{x}) \end{aligned}$$



Optimization in Machine Learning

First order methods

Step size and optimality



Learning goals

- Impact of step size
- Fixed vs. adaptive step size
- Exact line search
- Armijo rule & Backtracking
- Bracketing & Pinpointing

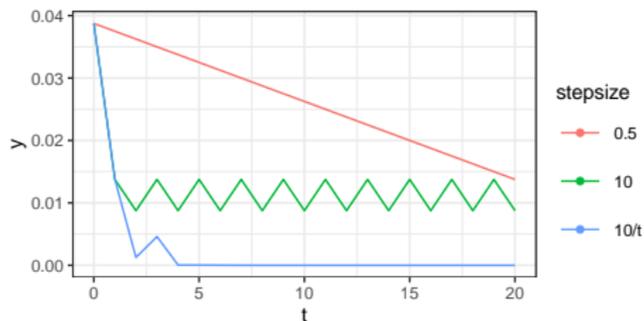
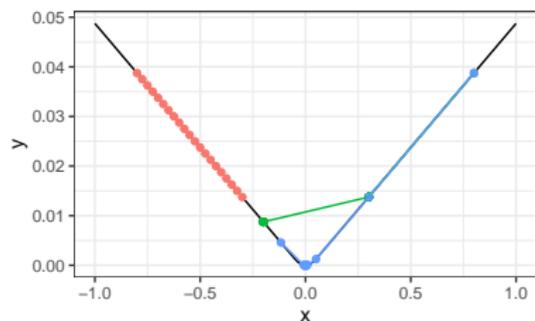
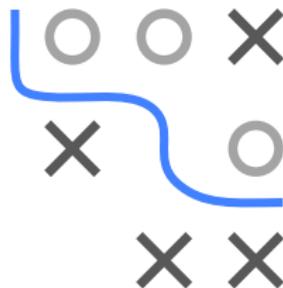
STEP SIZE CONTROL: DIMINISHING STEP SIZE

How can we adaptively control step size?

A natural way of selecting $\alpha^{[t]}$ is to decrease its value over time

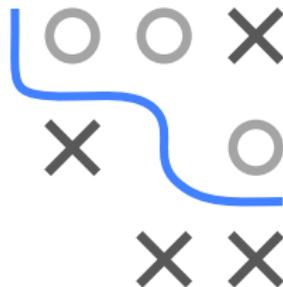
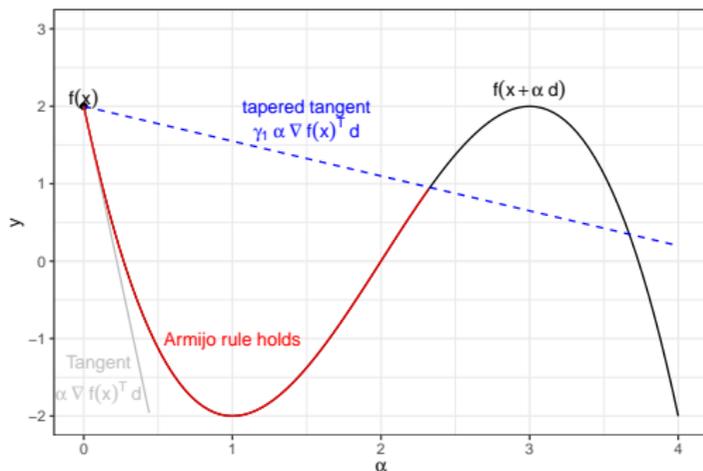
Example: GD on

$$f(x) = \begin{cases} \frac{1}{2}x^2 & \text{if } |x| \leq \delta, \\ \delta \cdot (|x| - 1/2 \cdot \delta) & \text{otherwise.} \end{cases}$$



GD with small constant (**red**), large constant (**green**), and diminishing (**blue**) step size

ARMIJO RULE



Inexact line search: Minimize objective “sufficiently” without computing optimal step size exactly

Common condition to guarantee “sufficient” decrease: **Armijo rule**

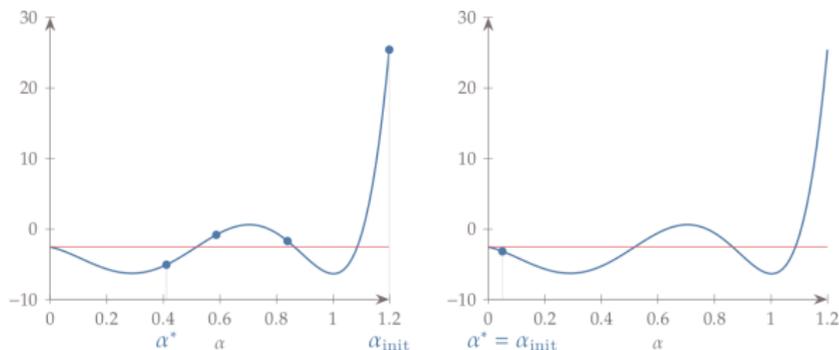
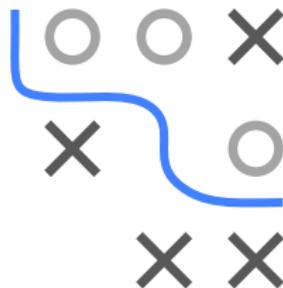
BACKTRACKING LINE SEARCH

Procedure to meet the Armijo rule: **Backtracking** line search

Idea: Decrease α until Armijo rule is met

Algorithm Backtracking line search

- 1: Choose initial step size $\alpha = \alpha_{\text{init}}$, $0 < \gamma_1 < 1$ and $0 < \tau < 1$
 - 2: **while** $f(\mathbf{x} + \alpha \mathbf{d}) > f(\mathbf{x}) + \gamma_1 \alpha \nabla f(\mathbf{x})^\top \mathbf{d}$ **do**
 - 3: Decrease α : $\alpha \leftarrow \tau \cdot \alpha$
 - 4: **end while**
-



(Source: Martins and Ning. *Engineering Design Optimization*, 2021.)

WOLFE CONDITIONS

Backtracking is simple and shows good performance in practice

But: Two undesirable scenarios

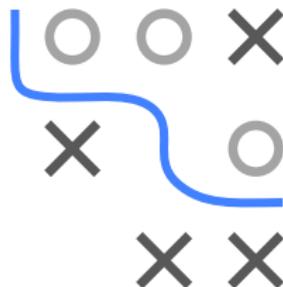
- 1 Initial step size α_{init} is too large \Rightarrow need multiple evaluations of f
- 2 Step size is too small with highly negative slopes

Solution for small step sizes:

- Fix γ_2 with $0 < \gamma_1 < \gamma_2 < 1$.
- α satisfies **sufficient curvature condition** in \mathbf{x} for \mathbf{d} if

$$|\nabla f(\mathbf{x} + \alpha\mathbf{d})^\top \mathbf{d}| \leq \gamma_2 |\nabla f(\mathbf{x})^\top \mathbf{d}|.$$

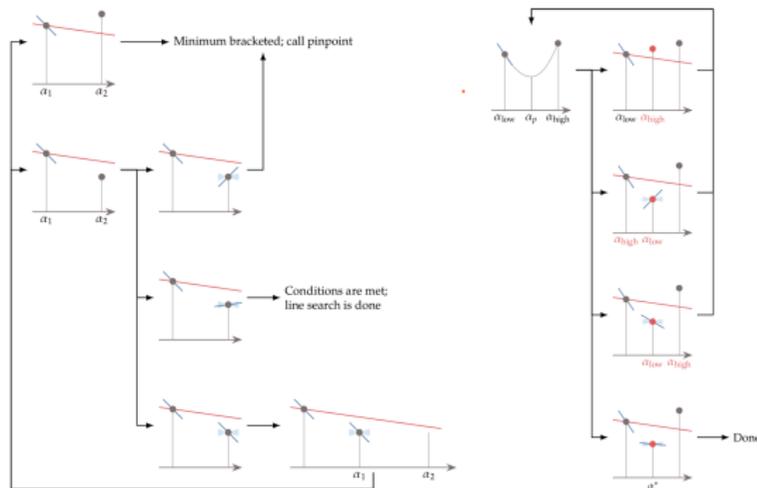
Armijo rule + sufficient curvature condition = **Wolfe conditions**



WOLFE CONDITIONS

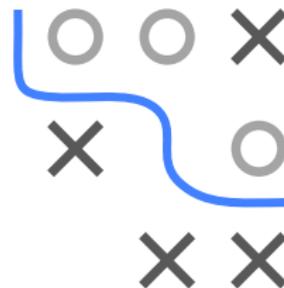
Algorithm for finding a Wolfe point (point satisfying Wolfe conditions):

- 1 Bracketing:** Find interval containing Wolfe point
- 2 Pinpointing:** Find Wolfe point in interval from bracketing



Left: Bracketing. Right: Pinpointing.

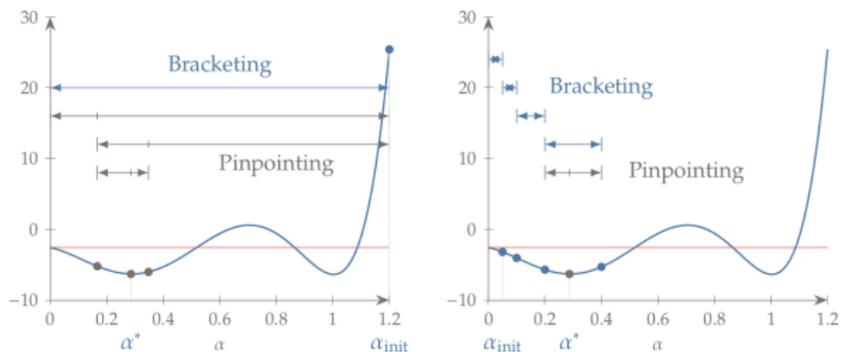
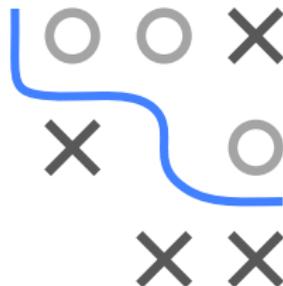
(Source: Martins and Ning. *EDO*, 2021.)



BRACKETING & PINPOINTING

Example:

- Large initial step size results in quick bracketing but multiple pinpointing steps (**left**).
- Small initial step size results in multiple bracketing steps but quick pinpointing (**right**).

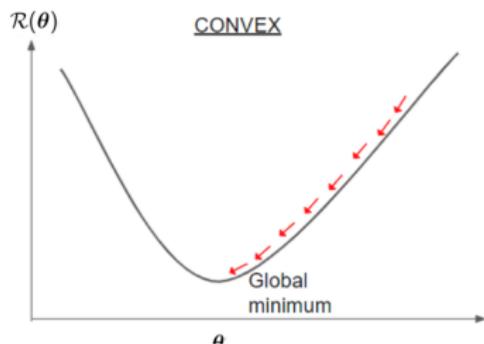
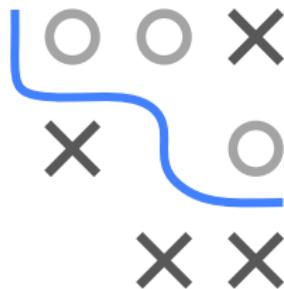


Source: Martins and Ning. *EDO*, 2021.

Optimization in Machine Learning

Deep dive

Gradient descent and optimality



Learning goals

- Convergence of GD
- Proof strategy and tools
- Descent lemma

SETTING

Assumptions:

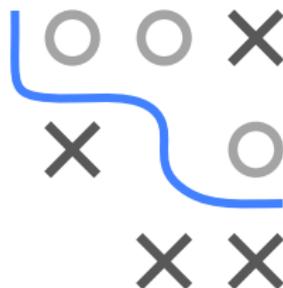
- f convex and differentiable
- Global minimum \mathbf{x}^* exists
- f has Lipschitz gradient (∇f does not change too fast)

$$\|\nabla f(\mathbf{x}) - \nabla f(\tilde{\mathbf{x}})\| \leq L\|\mathbf{x} - \tilde{\mathbf{x}}\| \quad \text{for all } \mathbf{x}, \tilde{\mathbf{x}}$$

Theorem (Convergence of GD). GD with step size $\alpha \leq 1/L$ yields

$$f(\mathbf{x}^{[k]}) - f(\mathbf{x}^*) \leq \frac{\|\mathbf{x}^{[0]} - \mathbf{x}^*\|^2}{2\alpha k}.$$

In other words: GD converges with rate $\mathcal{O}(1/k)$.



PROOF STRATEGY

- 1 Show that $f(\mathbf{x}^{[t]})$ **strictly decreases** with each iteration t

Descent lemma:

$$f(\mathbf{x}^{[t+1]}) \leq f(\mathbf{x}^{[t]}) - \frac{\alpha}{2} \|\nabla f(\mathbf{x}^{[t]})\|^2$$

at each iteration we at least gain



- 2 Bound **error of one step**

$$f(\mathbf{x}^{[t+1]}) - f(\mathbf{x}^*) \leq \frac{1}{2\alpha} \left(\|\mathbf{x}^{[t]} - \mathbf{x}^*\|^2 - \|\mathbf{x}^{[t+1]} - \mathbf{x}^*\|^2 \right)$$

we're not far more than by

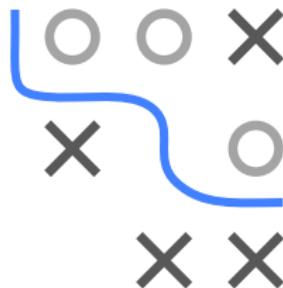
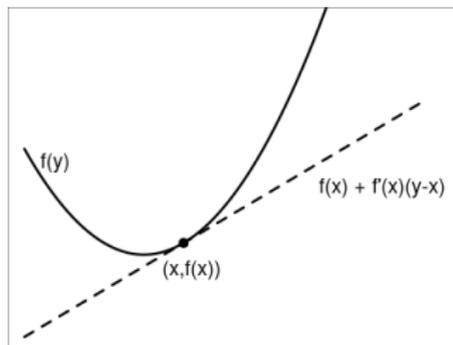
- 3 Finalize by **telescoping** argument

MAIN TOOL

Recall: First order condition of convexity

Every tangent line of f is always below f .

$$f(\mathbf{y}) \geq f(\mathbf{x}) + \nabla f(\mathbf{x})^\top (\mathbf{y} - \mathbf{x})$$



DESCENT LEMMA

Recall: ∇f Lipschitz $\implies \nabla^2 f(\mathbf{x}) \preceq L \cdot \mathbf{I}$ for all \mathbf{x}

This gives convexity of $g(\mathbf{x}) := \frac{L}{2} \|\mathbf{x}\|^2 - f(\mathbf{x})$ since

$$\nabla^2 g(\mathbf{x}) = L \cdot \mathbf{I} - \nabla^2 f(\mathbf{x}) \succeq 0.$$

First order condition of convexity of g yields

$$g(\mathbf{x}) \geq g(\mathbf{x}^{[t]}) + \nabla g(\mathbf{x}^{[t]})^\top (\mathbf{x} - \mathbf{x}^{[t]})$$

$$\Leftrightarrow \frac{L}{2} \|\mathbf{x}\|^2 - f(\mathbf{x}) \geq \frac{L}{2} \|\mathbf{x}^{[t]}\|^2 - f(\mathbf{x}^{[t]}) + (L\mathbf{x}^{[t]} - \nabla f(\mathbf{x}^{[t]}))^\top (\mathbf{x} - \mathbf{x}^{[t]})$$

\Leftrightarrow

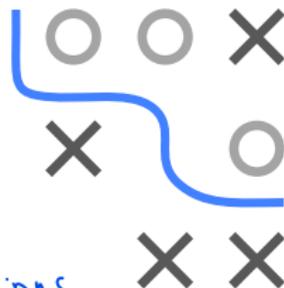
\vdots

$$\Leftrightarrow f(\mathbf{x}) \leq f(\mathbf{x}^{[t]}) + \nabla f(\mathbf{x}^{[t]})^\top (\mathbf{x} - \mathbf{x}^{[t]}) + \frac{L}{2} \|\mathbf{x} - \mathbf{x}^{[t]}\|^2$$

Now: One GD step with step size $\alpha \leq 1/L$:

$$\mathbf{x} \leftarrow \mathbf{x}^{[t+1]} = \mathbf{x}^{[t]} - \alpha \nabla f(\mathbf{x}^{[t]})$$

$$\nabla g(\mathbf{x}) = L\mathbf{x} - \nabla f(\mathbf{x})$$



plugging in definitions

$$\|\mathbf{x} - \mathbf{x}^{[t]}\|^2 = \|\mathbf{x}\|^2 - 2\mathbf{x}^\top \mathbf{x}^{[t]} + \|\mathbf{x}^{[t]}\|^2$$

$$L\mathbf{x}^{[t]\top} \mathbf{x} - L\mathbf{x}^{[t]\top} \mathbf{x}^{[t]}$$

DESCENT LEMMA

$$\begin{aligned} f(\mathbf{x}^{[t+1]}) &\leq f(\mathbf{x}^{[t]}) + \nabla f(\mathbf{x}^{[t]})^\top (\mathbf{x}^{[t+1]} - \mathbf{x}^{[t]}) + \frac{L}{2} \|\mathbf{x}^{[t+1]} - \mathbf{x}^{[t]}\|^2 \\ &= f(\mathbf{x}^{[t]}) + \nabla f(\mathbf{x}^{[t]})^\top (\mathbf{x}^{[t]} - \alpha \nabla f(\mathbf{x}^{[t]}) - \mathbf{x}^{[t]}) \\ &\quad + \frac{L}{2} \|\mathbf{x}^{[t]} - \alpha \nabla f(\mathbf{x}^{[t]}) - \mathbf{x}^{[t]}\|^2 \\ &= f(\mathbf{x}^{[t]}) - \nabla f(\mathbf{x}^{[t]})^\top \alpha \nabla f(\mathbf{x}^{[t]}) + \frac{L}{2} \|\alpha \nabla f(\mathbf{x}^{[t]})\|^2 \\ &= f(\mathbf{x}^{[t]}) - \alpha \|\nabla f(\mathbf{x}^{[t]})\|^2 + \frac{L\alpha^2}{2} \|\nabla f(\mathbf{x}^{[t]})\|^2 \\ &\leq f(\mathbf{x}^{[t]}) - \frac{\alpha}{2} \|\nabla f(\mathbf{x}^{[t]})\|^2 \end{aligned}$$

Note: $\alpha \leq 1/L$ yields $L\alpha^2 \leq \alpha$

- $\|\nabla f(\mathbf{x}^{[t]})\|^2 > 0$ unless $\nabla f(\mathbf{x}) = \mathbf{0}$
- f **strictly decreases** with each GD iteration until optimum reached
- Descent lemma yields bound on **guaranteed progress** if $\alpha \leq 1/L$
(explains why GD may diverge if step sizes too large)



ONE STEP ERROR BOUND

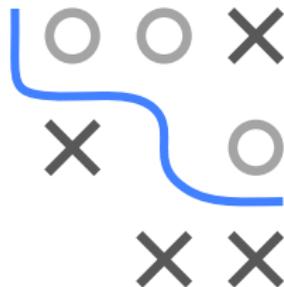
Again, first order condition of convexity gives

$$f(\mathbf{x}^{[t]}) - f(\mathbf{x}^*) \leq \nabla f(\mathbf{x}^{[t]})^\top (\mathbf{x}^{[t]} - \mathbf{x}^*).$$

This and the descent lemma yields

$$\begin{aligned} f(\mathbf{x}^{[t+1]}) - f(\mathbf{x}^*) &\leq f(\mathbf{x}^{[t]}) - \frac{\alpha}{2} \|\nabla f(\mathbf{x}^{[t]})\|^2 - f(\mathbf{x}^*) \\ &= \underbrace{f(\mathbf{x}^{[t]}) - f(\mathbf{x}^*)}_{\text{DB}} - \frac{\alpha}{2} \|\nabla f(\mathbf{x}^{[t]})\|^2 \\ &\leq \nabla f(\mathbf{x}^{[t]})^\top (\mathbf{x}^{[t]} - \mathbf{x}^*) - \frac{\alpha}{2} \|\nabla f(\mathbf{x}^{[t]})\|^2 \\ &= \frac{1}{2\alpha} \left(\|\mathbf{x}^{[t]} - \mathbf{x}^*\|^2 - \|\mathbf{x}^{[t]} - \mathbf{x}^* - \alpha \nabla f(\mathbf{x}^{[t]})\|^2 \right) \\ &= \frac{1}{2\alpha} \left(\|\mathbf{x}^{[t]} - \mathbf{x}^*\|^2 - \|\mathbf{x}^{[t+1]} - \mathbf{x}^*\|^2 \right) \end{aligned}$$

Note: Line 3 \rightarrow 4 is hard to see (just expand line 4).

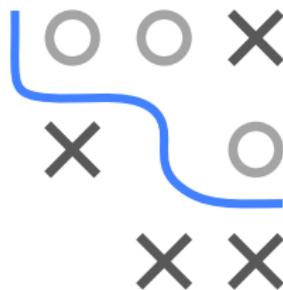


HW, expand and simplify

FINALIZATION

Summing over iterations yields *because each next iteration yields better approx.*

$$\begin{aligned}k(f(\mathbf{x}^{[k]}) - f(\mathbf{x}^*)) &\leq \sum_{t=1}^k [f(\mathbf{x}^{[t]}) - f(\mathbf{x}^*)] \\ &\leq \sum_{t=1}^k \frac{1}{2\alpha} \left[\|\mathbf{x}^{[t-1]} - \mathbf{x}^*\|^2 - \|\mathbf{x}^{[t]} - \mathbf{x}^*\|^2 \right] \\ &= \frac{1}{2\alpha} \left(\|\mathbf{x}^{[0]} - \mathbf{x}^*\|^2 - \|\mathbf{x}^{[k]} - \mathbf{x}^*\|^2 \right) \\ &\leq \frac{1}{2\alpha} \left(\|\mathbf{x}^{[0]} - \mathbf{x}^*\|^2 \right).\end{aligned}$$



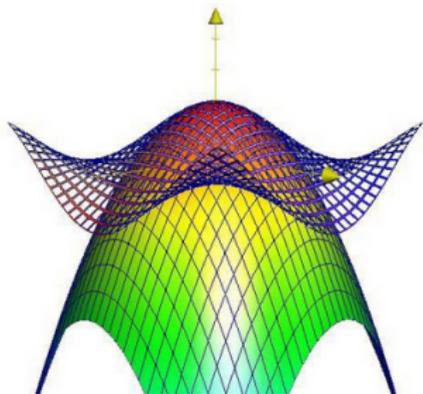
Arguments: Descent lemma (line 1). Telescoping sum (line 2 \rightarrow 3).

$$f(\mathbf{x}^{[k]}) - f(\mathbf{x}^*) \leq \frac{\|\mathbf{x}^{[0]} - \mathbf{x}^*\|^2}{2\alpha k}$$

Optimization in Machine Learning

First order methods

Weaknesses of GD – Curvature



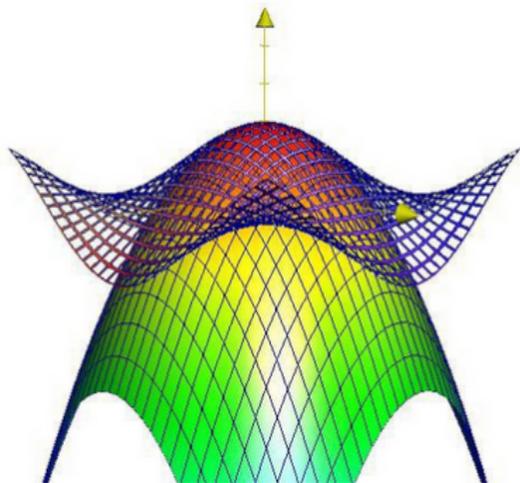
Learning goals

- Effects of curvature
- Step size effect in GD

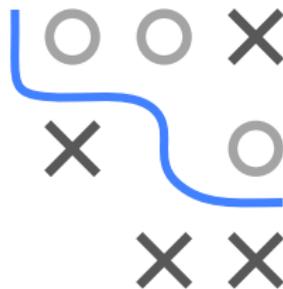
REMINDER: LOCAL QUADRATIC GEOMETRY

Locally approximate smooth function by quadratic Taylor polynomial:

$$f(\mathbf{x}) \approx f(\tilde{\mathbf{x}}) + \nabla f(\tilde{\mathbf{x}})^\top (\mathbf{x} - \tilde{\mathbf{x}}) + \frac{1}{2} (\mathbf{x} - \tilde{\mathbf{x}})^\top \nabla^2 f(\tilde{\mathbf{x}}) (\mathbf{x} - \tilde{\mathbf{x}})$$



Source: daniloroccatano.blog.



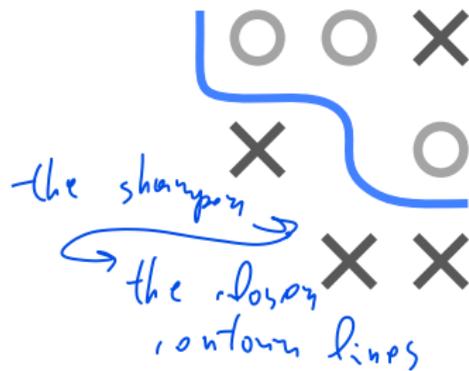
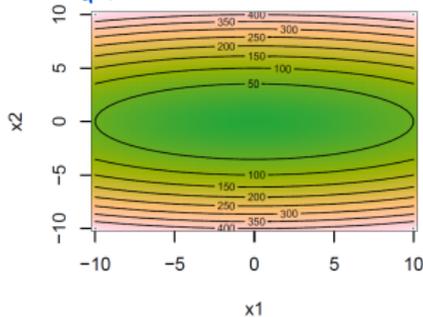
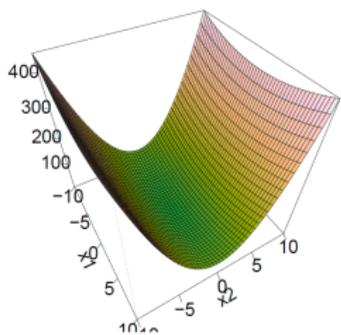
REMINDER: LOCAL QUADRATIC GEOMETRY

Study Hessian $\mathbf{H} = \nabla^2 f(\mathbf{x}^{[t]})$ in GD to discuss effect of curvature

Recall for quadratic forms:

- Eigenvector \mathbf{v}_{\max} (\mathbf{v}_{\min}) is direction of largest (smallest) curvature
- \mathbf{H} called ill-conditioned if $\kappa(\mathbf{H}) = |\lambda_{\max}|/|\lambda_{\min}|$ is large

condition number
sounding thing



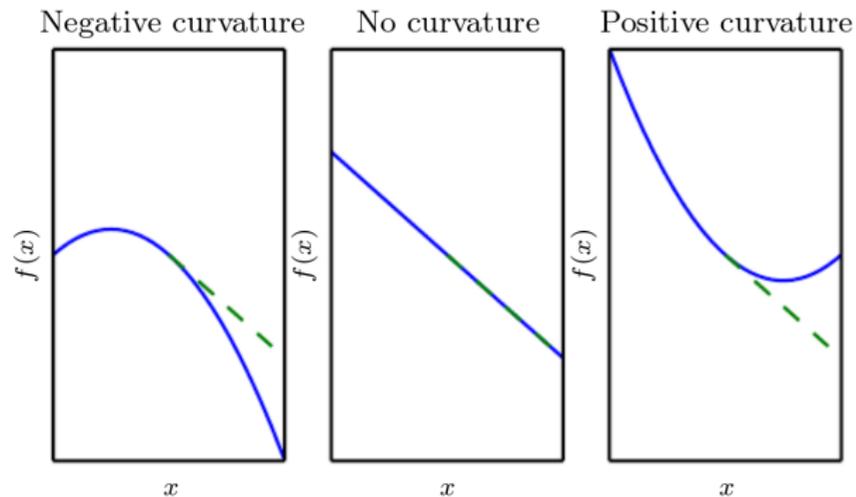
$$\hat{\Theta} = (X^T X)^{-1} X^T y$$
$$(X^T X) \hat{\Theta} = X^T y$$

Sys. lin equations
S.L.E.

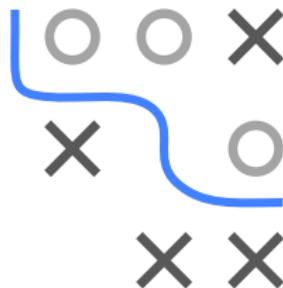
float - 4 byte
double - 8 byte

EFFECTS OF CURVATURE

Intuitively, curvature determines reliability of a GD step



Quadratic objective f (blue) with gradient approximation (dashed green).
Left: f decreases faster than ∇f predicts. **Center:** ∇f predicts decrease correctly. **Right:** f decreases more slowly than ∇f predicts.
(Source: Goodfellow et al., 2016)



CURVATURE AND STEP SIZE IN GD

Worst case: \mathbf{H} is ill-conditioned. What does this mean for GD?

- Quadratic Taylor polynomial of f around $\tilde{\mathbf{x}}$ (with gradient $\mathbf{g} = \nabla f$)

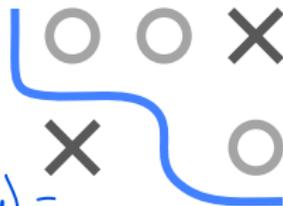
$$f(\mathbf{x}) \approx f(\tilde{\mathbf{x}}) + (\mathbf{x} - \tilde{\mathbf{x}})^\top \mathbf{g} + \frac{1}{2}(\mathbf{x} - \tilde{\mathbf{x}})^\top \mathbf{H}(\mathbf{x} - \tilde{\mathbf{x}})$$

- GD step with step size $\alpha > 0$ yields

$$f(\tilde{\mathbf{x}} - \alpha \mathbf{g}) \approx f(\tilde{\mathbf{x}}) - \alpha \mathbf{g}^\top \mathbf{g} + \frac{1}{2} \alpha^2 \mathbf{g}^\top \mathbf{H} \mathbf{g}$$

- If $\mathbf{g}^\top \mathbf{H} \mathbf{g} > 0$, we can solve for optimal step size α^* :

$$\alpha^* = \frac{\mathbf{g}^\top \mathbf{g}}{\mathbf{g}^\top \mathbf{H} \mathbf{g}}$$

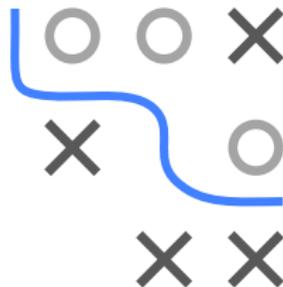
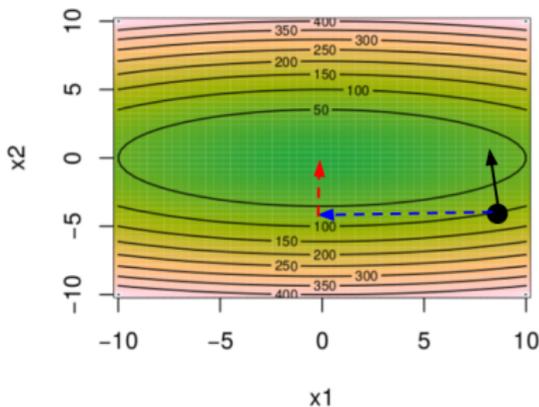


$\Delta_\alpha f(\tilde{\mathbf{x}} - \alpha \mathbf{g}) =$
 $= -\mathbf{g}^\top \mathbf{g} + \frac{1}{2} \alpha^2 \mathbf{g}^\top \mathbf{H} \mathbf{g}$
 $\frac{d}{d\alpha} \Delta_\alpha f(\tilde{\mathbf{x}} - \alpha \mathbf{g}) = -\mathbf{g}^\top \mathbf{g} + \alpha \mathbf{g}^\top \mathbf{H} \mathbf{g} = 0$

if $\mathbf{g}^\top \mathbf{H} \mathbf{g} < 0$ would α be the worst α ?

CURVATURE AND STEP SIZE IN GD

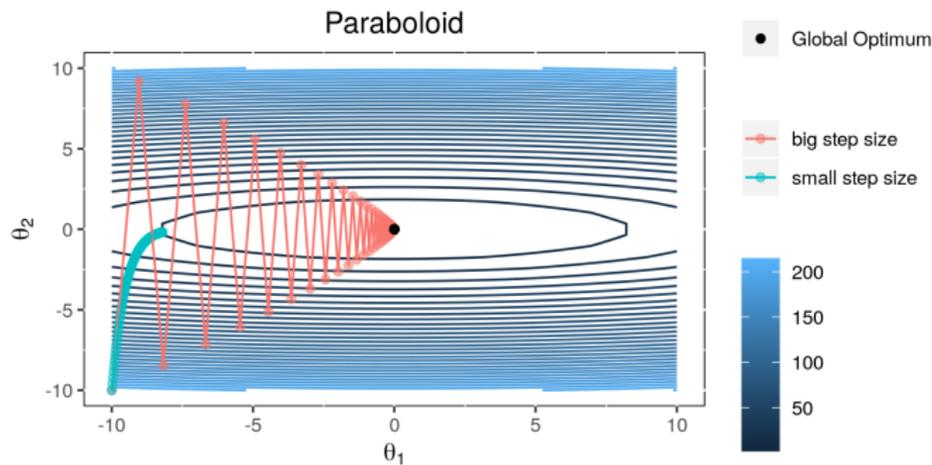
- What if \mathbf{g} is not aligned with eigenvectors?
- Consider 2D case: Decompose \mathbf{g} (black) into \mathbf{v}_{\max} and \mathbf{v}_{\min}



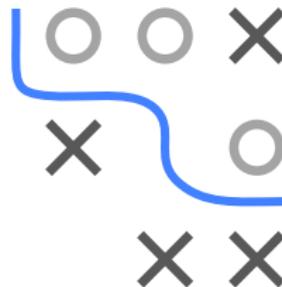
- Ideally, perform **large** step along \mathbf{v}_{\min} but **small** step along \mathbf{v}_{\max}
- However, gradient almost only points along \mathbf{v}_{\max}

CURVATURE AND STEP SIZE IN GD

- GD is not aware of curvatures and can only walk along \mathbf{g}
- Large step sizes result in “zig-zag” behaviour.
- Small step sizes result in weak performance.



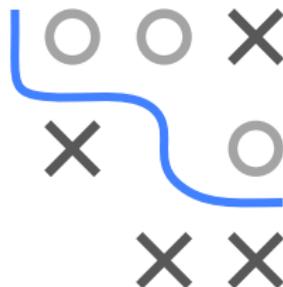
Poorly conditioned quadratic form. GD with large (red) and small (blue) step size. For both, convergence to optimum is slow.



CURVATURE AND STEP SIZE IN GD

- If gradient norms $\|\mathbf{g}\|$ increase, GD is not converging since $\mathbf{g} \neq 0$.
- Even if $\|\mathbf{g}\|$ increases, objective may stay approximately constant:

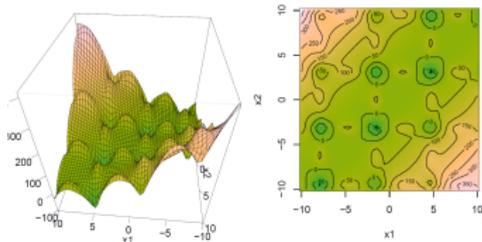
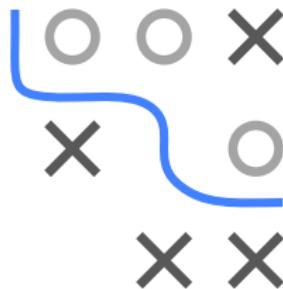
$$\underbrace{f(\tilde{\mathbf{x}} - \alpha \mathbf{g})}_{\approx \text{constant}} \approx f(\tilde{\mathbf{x}}) - \alpha \underbrace{\mathbf{g}^T \mathbf{g}}_{\text{increases}} + \frac{1}{2} \alpha^2 \underbrace{\mathbf{g}^T \mathbf{H} \mathbf{g}}_{\text{increases}}$$



Optimization in Machine Learning

First order methods

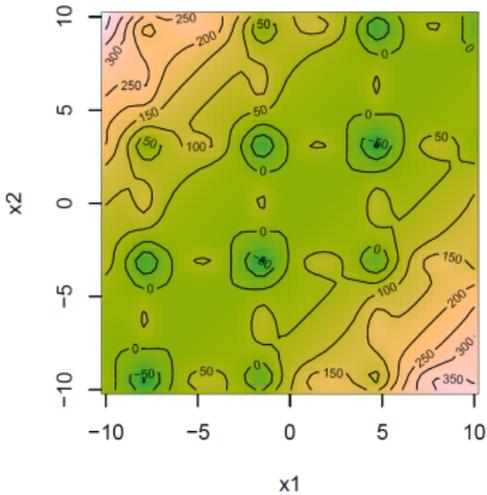
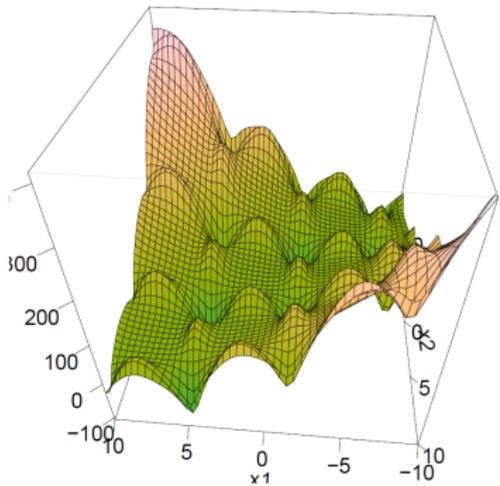
GD – Multimodality and Saddle points



Learning goals

- Multimodality, GD result can be arbitrarily bad
- Saddle points, major problem in NN error landscapes, GD can get stuck or slow crawling

UNIMODAL VS. MULTIMODAL LOSS SURFACES

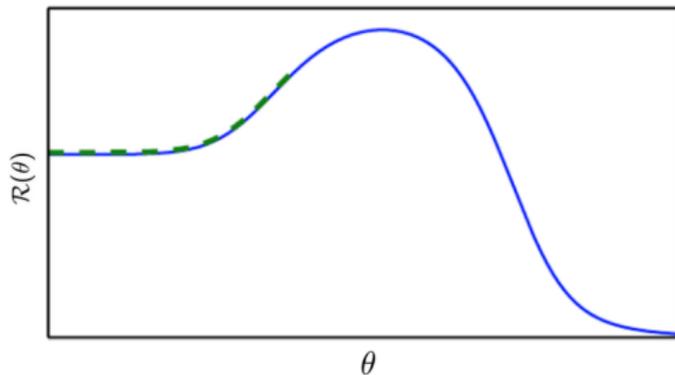


Snippet of a loss surface with many local optima



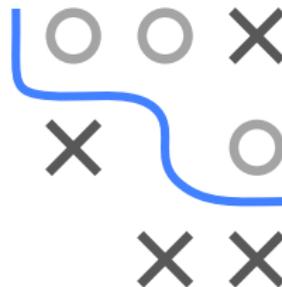
GD: ONLY LOCALLY OPTIMAL MOVES

- GD makes only **locally** optimal moves
- It may move away from the global optimum



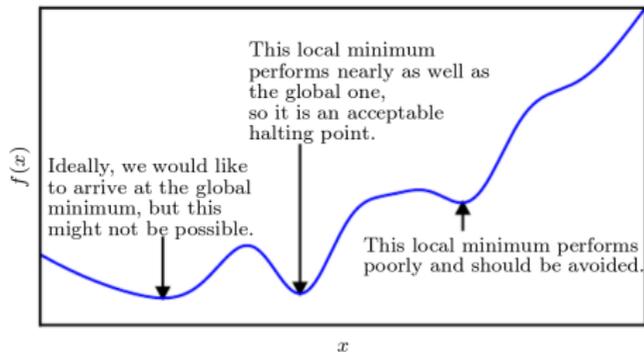
Source: Goodfellow et al., 2016

- Initialization on “wrong” side of the hill results in weak performance
- In higher dimensions, GD may move around the hill (potentially at the cost of longer trajectory and time to convergence)

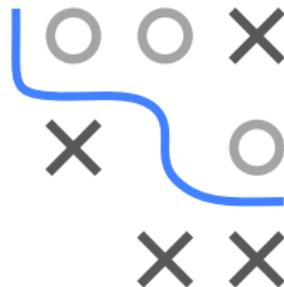


LOCAL MINIMA

- **In practice:** Only local minima with high value compared to global minimum are problematic.

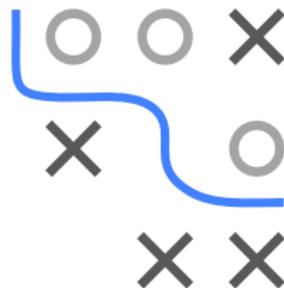
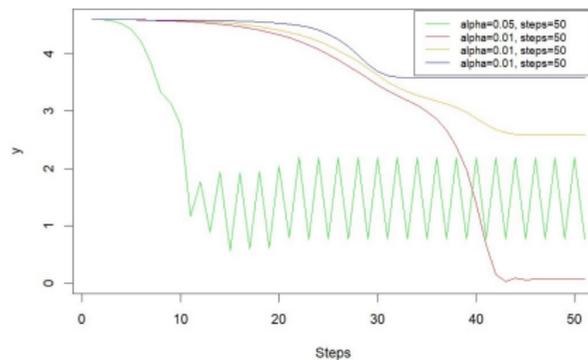
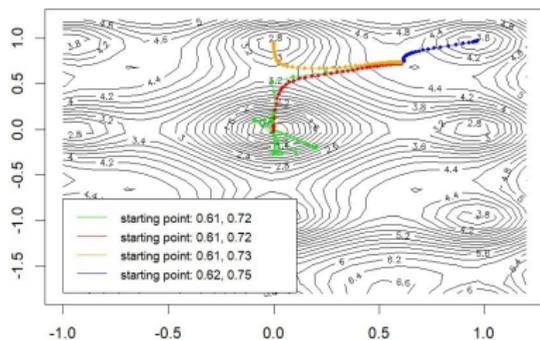


Source: Goodfellow et al., 2016



LOCAL MINIMA

- Small differences in starting point or step size can lead to huge differences in the reached minimum or even to non-convergence



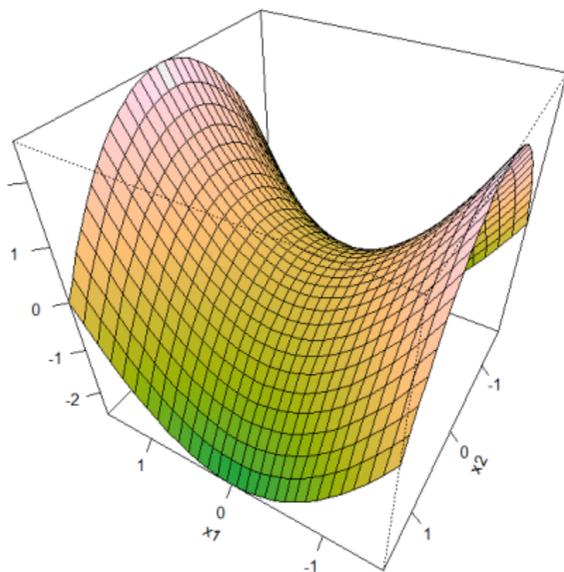
(Non-)Converging gradient descent for Ackley function

GD AT SADDLE POINTS

Example:

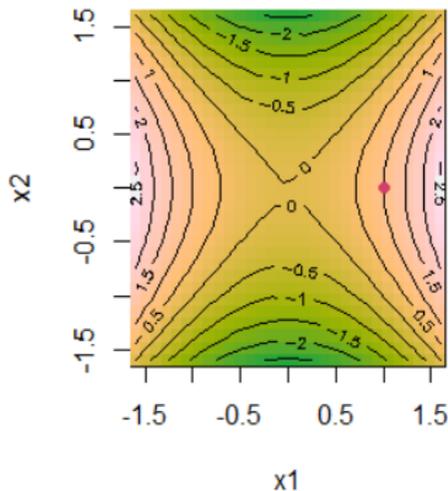
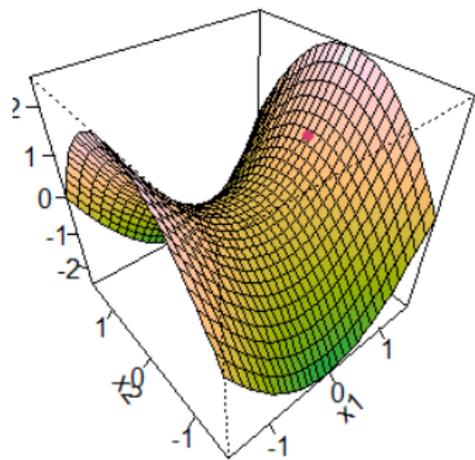
$$\begin{aligned}f(x_1, x_2) &= x_1^2 - x_2^2 \\ \nabla f(x_1, x_2) &= (2x_1, -2x_2)^\top \\ \mathbf{H} &= \begin{pmatrix} 2 & 0 \\ 0 & -2 \end{pmatrix}\end{aligned}$$

- Along x_1 , curvature is positive ($\lambda_1 = 2 > 0$).
- Along x_2 , curvature is negative ($\lambda_2 = -2 < 0$).



EXAMPLE: SADDLE POINT WITH GD

- How do saddle points impair optimization?
- Gradient-based algorithms **might** get stuck in saddle points

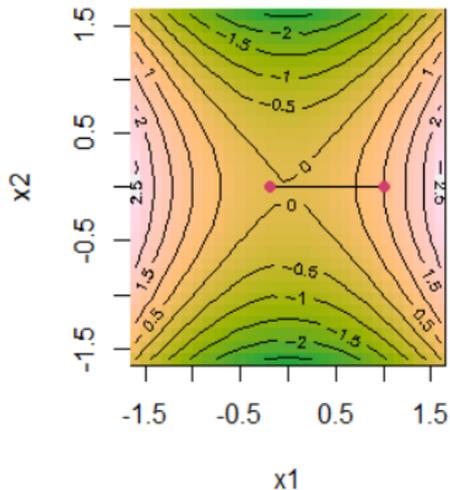
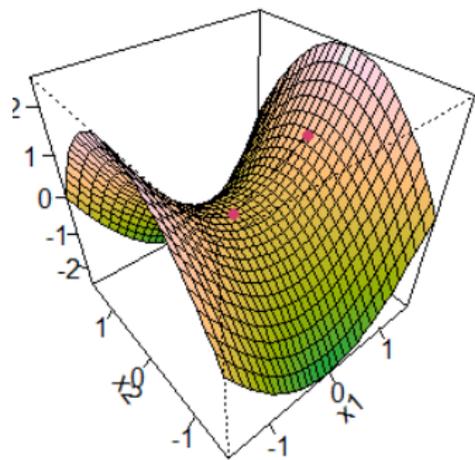


Red dot: Starting location



EXAMPLE: SADDLE POINT WITH GD

- How do saddle points impair optimization?
- Gradient-based algorithms **might** get stuck in saddle points

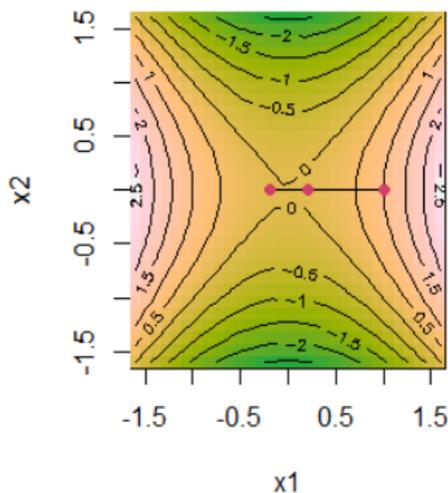
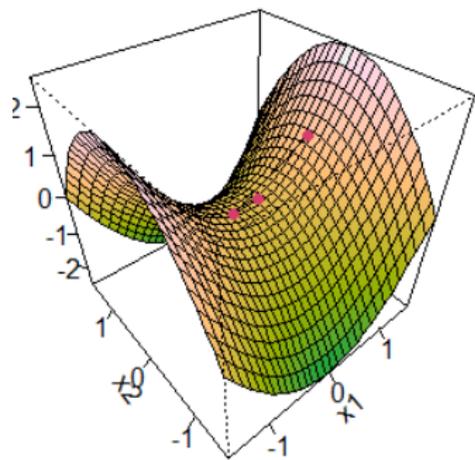


Step 1 ...



EXAMPLE: SADDLE POINT WITH GD

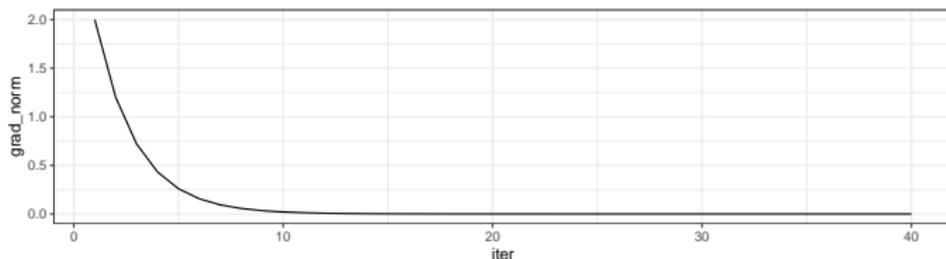
- How do saddle points impair optimization?
- Gradient-based algorithms **might** get stuck in saddle points



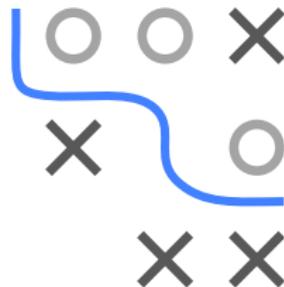
... Step 10 ... got stuck and cannot escape saddle point

EXAMPLE: SADDLE POINT WITH GD

- How do saddle points impair optimization?
- Gradient-based algorithms **might** get stuck in saddle points

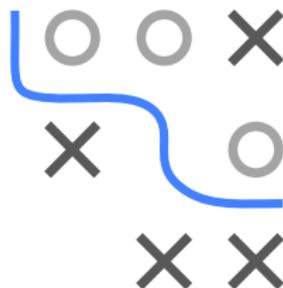


... Step 10 ... got stuck and cannot escape saddle point



SADDLE POINTS IN NEURAL NETWORKS

- For the empirical risk $\mathcal{R} : \mathbb{R}^d \rightarrow \mathbb{R}$ of a neural network, the expected ratio of the number of saddle points to local minima typically grows exponentially with d *→ Is there a proof?*
- In other words: Networks with more parameters (deeper networks or larger layers) exhibit a lot more saddle points than local minima
- **Reason:** Hessian at local minimum has only positive eigenvalues. Hessian at saddle point has positive and negative eigenvalues.



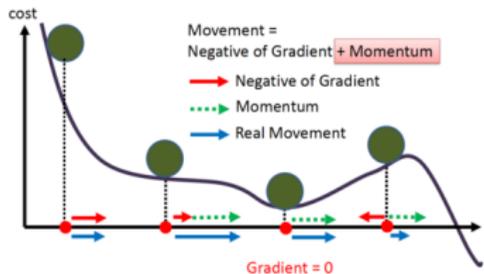
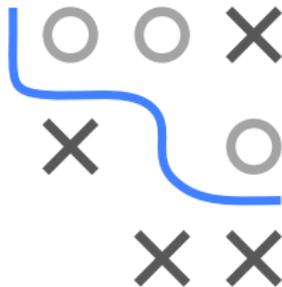
SADDLE POINTS IN NEURAL NETWORKS

- Imagine the sign of each eigenvalue is generated by coin flipping:
 - In a single dimension, it is easy to obtain a local minimum (e.g. “head” means positive eigenvalue).
 - In an m -dimensional space, it is exponentially unlikely that all m coin tosses will be head.
- A property of many random functions is that eigenvalues of the Hessian become more likely to be positive in regions of lower cost.
- For the coin flipping example, this means we are more likely to have heads m times if we are at a critical point with low cost.
- That means in particular that local minima are much more likely to have low cost than high cost and critical points with high cost are far more likely to be saddle points.
- “Saddle points are surrounded by high error plateaus that can dramatically slow down learning, and give the illusory impression of the existence of a local minimum” (Dauphin et al. (2014)).



Optimization in Machine Learning

First order methods GD with Momentum

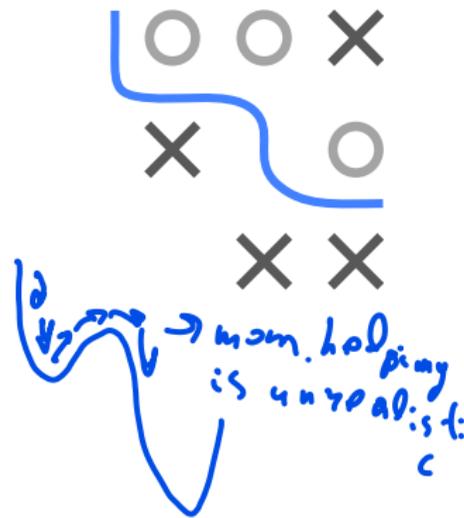


Learning goals

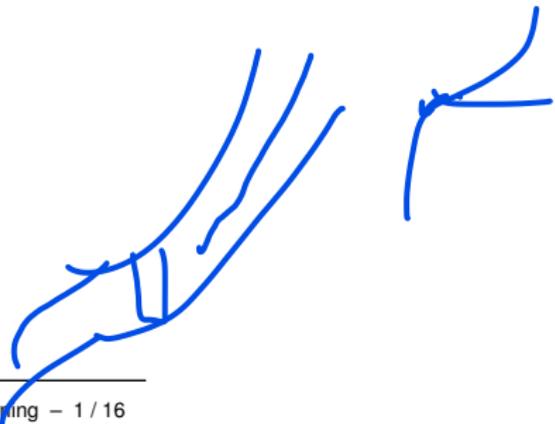
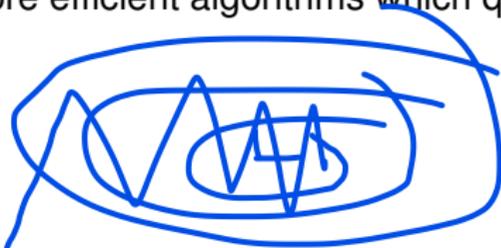
- Recap of GD problems
- Momentum definition
- Unrolling formula
- Examples
- Nesterov

RECAP: WEAKNESSES OF GRADIENT DESCENT

- **Zig-zagging behavior:** For ill-conditioned problems, GD moves with a zig-zag course to the optimum, since the gradient points approximately orthogonal in the shortest direction to the minimum.
- **Slow crawling:** may vanish rapidly close to stationary points (e.g. saddle points) and hence also slows down progress.
- **Trapped in stationary points:** In some functions GD converges to stationary points (e.g. saddle points) since gradient on all sides is fairly flat and the step size is too small to pass this flat part.



Aim: More efficient algorithms which quickly reach the minimum.



MOMENTUM: ANALYSIS

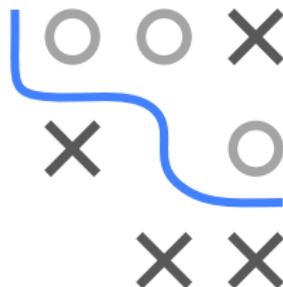
$$\boldsymbol{\nu}^{[1]} = \varphi \boldsymbol{\nu}^{[0]} - \alpha \nabla f(\mathbf{x}^{[0]})$$

$$\mathbf{x}^{[1]} = \mathbf{x}^{[0]} + \varphi \boldsymbol{\nu}^{[0]} - \alpha \nabla f(\mathbf{x}^{[0]})$$

$$\boldsymbol{\nu}^{[2]} = \varphi \boldsymbol{\nu}^{[1]} - \alpha \nabla f(\mathbf{x}^{[1]})$$

$$= \varphi(\varphi \boldsymbol{\nu}^{[0]} - \alpha \nabla f(\mathbf{x}^{[0]})) - \alpha \nabla f(\mathbf{x}^{[1]})$$

$$\mathbf{x}^{[2]} = \mathbf{x}^{[1]} + \varphi(\varphi \boldsymbol{\nu}^{[0]} - \alpha \nabla f(\mathbf{x}^{[0]})) - \alpha \nabla f(\mathbf{x}^{[1]})$$



MOMENTUM: ANALYSIS

$$\mathbf{v}^{[1]} = \varphi \mathbf{v}^{[0]} - \alpha \nabla f(\mathbf{x}^{[0]})$$

$$\mathbf{x}^{[1]} = \mathbf{x}^{[0]} + \varphi \mathbf{v}^{[0]} - \alpha \nabla f(\mathbf{x}^{[0]})$$

$$\begin{aligned} \mathbf{v}^{[2]} &= \varphi \mathbf{v}^{[1]} - \alpha \nabla f(\mathbf{x}^{[1]}) \\ &= \varphi(\varphi \mathbf{v}^{[0]} - \alpha \nabla f(\mathbf{x}^{[0]})) - \alpha \nabla f(\mathbf{x}^{[1]}) \end{aligned}$$

$$\mathbf{x}^{[2]} = \mathbf{x}^{[1]} + \varphi(\varphi \mathbf{v}^{[0]} - \alpha \nabla f(\mathbf{x}^{[0]})) - \alpha \nabla f(\mathbf{x}^{[1]})$$

$$\begin{aligned} \mathbf{v}^{[3]} &= \varphi \mathbf{v}^{[2]} - \alpha \nabla f(\mathbf{x}^{[2]}) \\ &= \varphi(\varphi(\varphi \mathbf{v}^{[0]} - \alpha \nabla f(\mathbf{x}^{[0]})) - \alpha \nabla f(\mathbf{x}^{[1]})) - \alpha \nabla f(\mathbf{x}^{[2]}) \end{aligned}$$

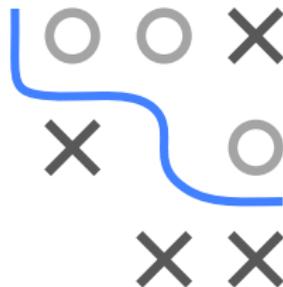
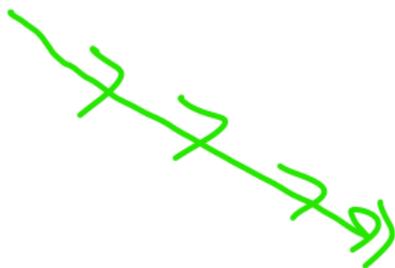
$$\mathbf{x}^{[3]} = \mathbf{x}^{[2]} + \varphi(\varphi(\varphi \mathbf{v}^{[0]} - \alpha \nabla f(\mathbf{x}^{[0]})) - \alpha \nabla f(\mathbf{x}^{[1]})) - \alpha \nabla f(\mathbf{x}^{[2]})$$

$$\begin{aligned} &= \mathbf{x}^{[2]} + \varphi^3 \mathbf{v}^{[0]} - \varphi^2 \alpha \nabla f(\mathbf{x}^{[0]}) - \varphi \alpha \nabla f(\mathbf{x}^{[1]}) - \alpha \nabla f(\mathbf{x}^{[2]}) \\ &= \mathbf{x}^{[2]} - \alpha(\varphi^2 \nabla f(\mathbf{x}^{[0]}) + \varphi \nabla f(\mathbf{x}^{[1]}) + \nabla f(\mathbf{x}^{[2]})) + \varphi^3 \mathbf{v}^{[0]} \end{aligned}$$

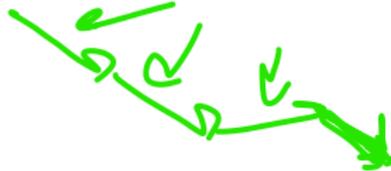
exponential smoothing

$$\mathbf{x}^{[t+1]} = \mathbf{x}^{[t]} - \alpha \sum_{j=0}^t \varphi^j \nabla f(\mathbf{x}^{[t-j]}) + \varphi^{t+1} \mathbf{v}^{[0]}$$

φ^j



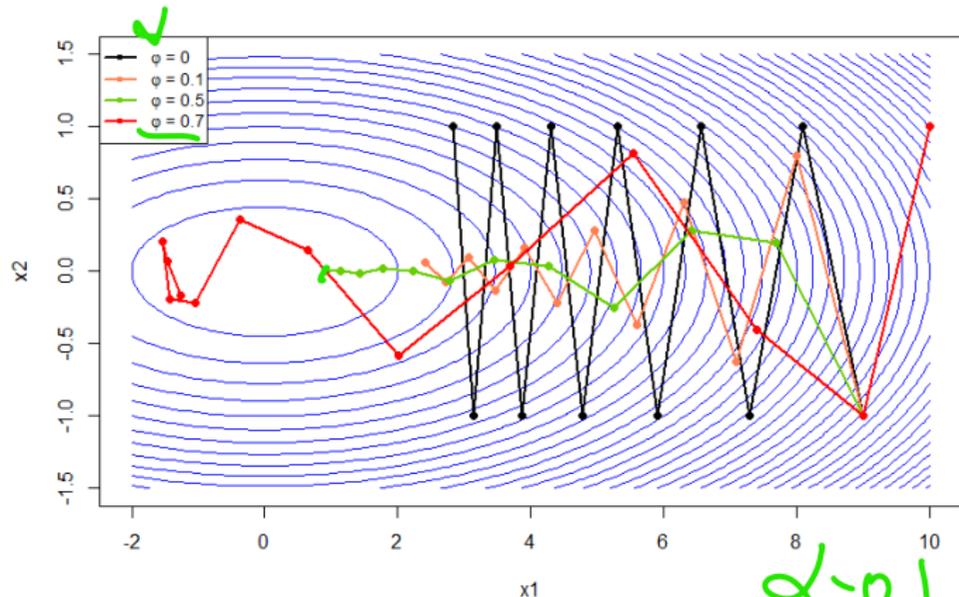
the older the step the less weight it gets



GD WITH MOMENTUM: ZIG-ZAG BEHAVIOUR

Caution:

- If momentum is too high, minimum is possibly missed
- We might go back and forth around or between local minima



α too close to 1
can't guarantee convergence



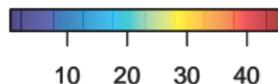
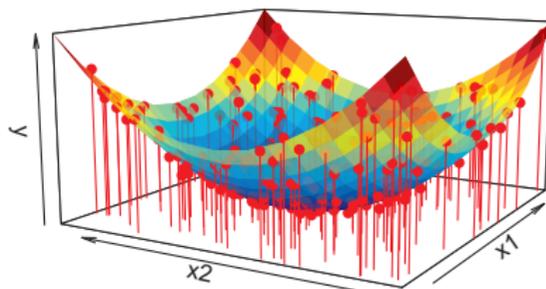
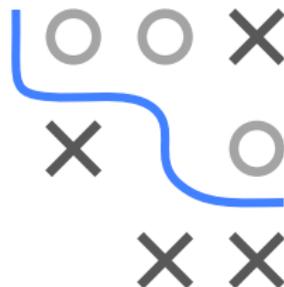
$$\underline{\theta}_0 + \underline{\alpha} x$$

ERM FOR NN WITH GD

Let $\mathcal{D} = ((\mathbf{x}^{(1)}, y^{(1)}), \dots, (\mathbf{x}^{(n)}, y^{(n)}))$, with $y = x_1^2 + x_2^2$ and minimize

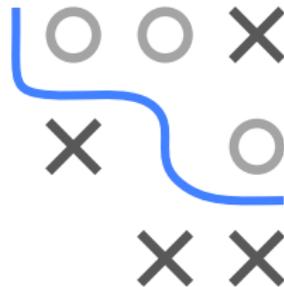
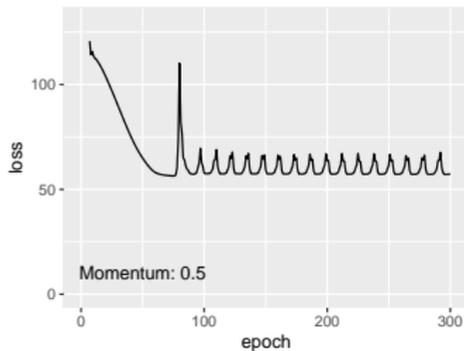
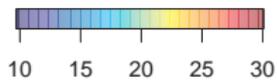
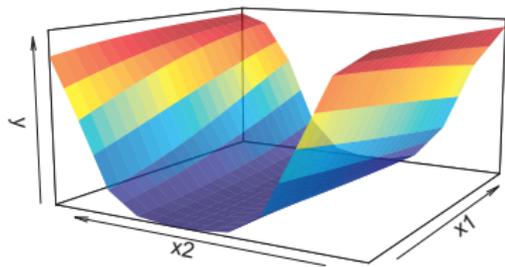
$$\mathcal{R}_{\text{emp}}(\boldsymbol{\theta}) = \sum_{i=1}^n \left(f(\mathbf{x} \mid \boldsymbol{\theta}) - y^{(i)} \right)^2$$

where $f(\mathbf{x} \mid \boldsymbol{\theta})$ is a neural network with 2 hidden layers (2 units each).



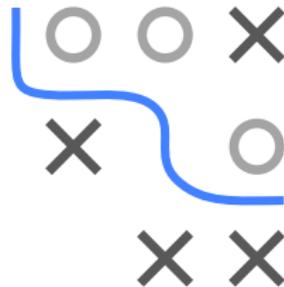
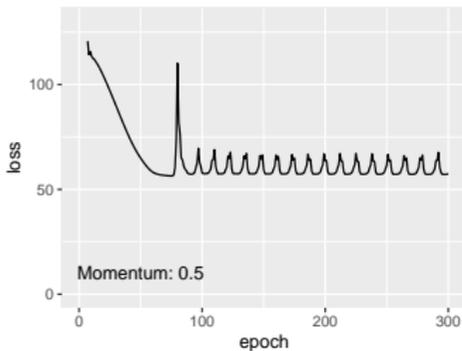
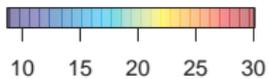
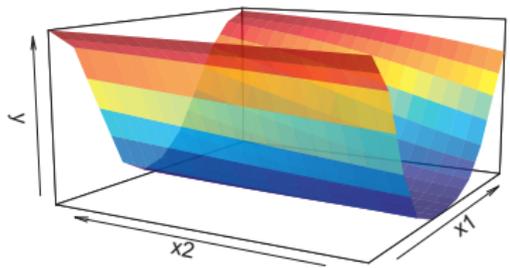
ERM FOR NN WITH GD

After 100 iters of GD:



ERM FOR NN WITH GD

After 300 iters of GD:



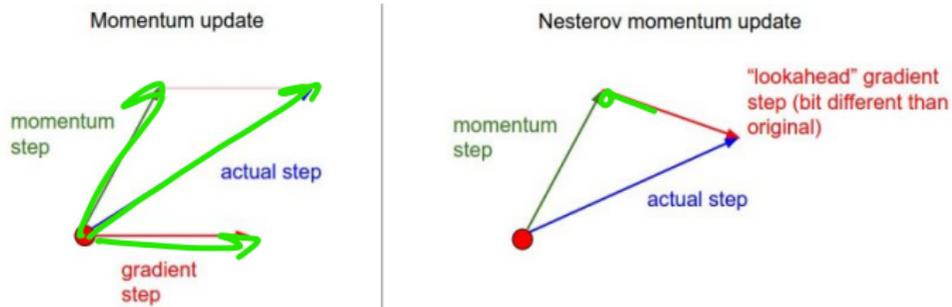
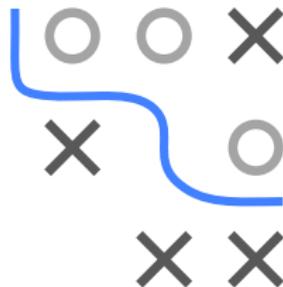
NESTEROV ACCELERATED GRADIENT

- Slightly modified version: **Nesterov accelerated gradient**
- Stronger theoretical convergence guarantees for convex functions
- Avoid moving back and forth near optima

$$\boldsymbol{v}^{[t+1]} = \varphi \boldsymbol{v}^{[t]} - \alpha \nabla f(\boldsymbol{x}^{[t]} + \varphi \boldsymbol{v}^{[t]})$$

$$\boldsymbol{x}^{[t+1]} = \boldsymbol{x}^{[t]} + \boldsymbol{v}^{[t+1]}$$

ToDo



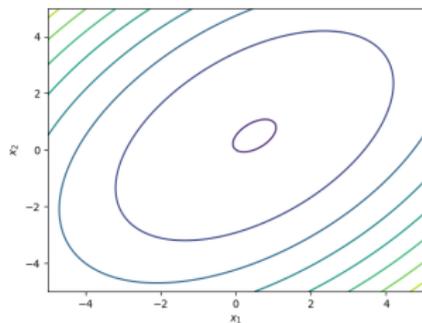
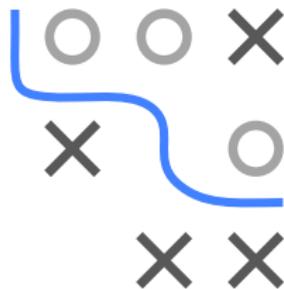
Nesterov momentum update evaluates gradient at the "look-ahead" position.

(Source: <https://cs231n.github.io/neural-networks-3/>)

Optimization in Machine Learning

First order methods

GD on quadratic forms



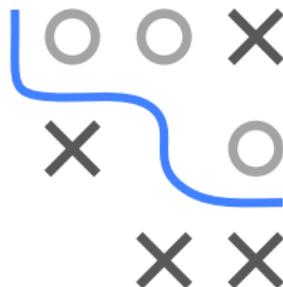
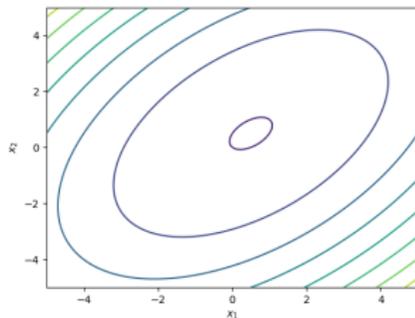
Learning goals

- Eigendecomposition of quadratic forms
- GD steps in eigenspace

QUADRATIC FORMS & GD

- We consider the quadratic function $q(\mathbf{x}) = \mathbf{x}^\top \mathbf{A} \mathbf{x} - \mathbf{b}^\top \mathbf{x}$.
- We assume that Hessian $\mathbf{H} = 2\mathbf{A}$ has full rank
- Optimal solution is $\mathbf{x}^* = \frac{1}{2} \mathbf{A}^{-1} \mathbf{b}$
- As $\nabla q(\mathbf{x}) = 2\mathbf{A} \mathbf{x} - \mathbf{b}$, iterations of gradient descent are

$$\mathbf{x}^{[t+1]} = \mathbf{x}^{[t]} - \alpha(2\mathbf{A}\mathbf{x}^{[t]} - \mathbf{b})$$

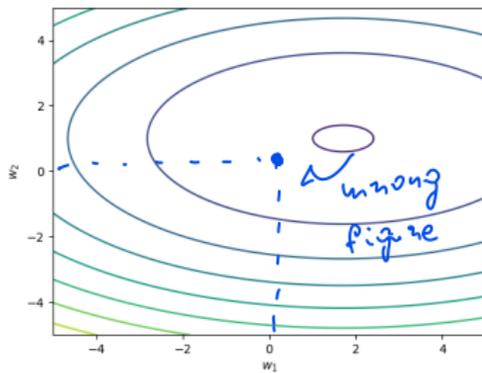
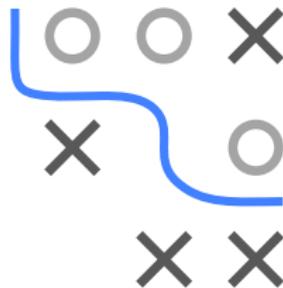


The following slides follow the blog post "Why Momentum Really Works", Distill, 2017.

<http://doi.org/10.23915/distill.00006>

EIGENDECOMPOSITION OF QUADRATIC FORMS

- We want to work in the coordinate system given by q
- **Recall:** Coordinate system is given by the eigenvectors of $\mathbf{H} = 2\mathbf{A}$
- Eigendecomposition of $\mathbf{A} = \mathbf{V}\mathbf{\Lambda}\mathbf{V}^T$
- \mathbf{V} contains eigenvectors \mathbf{v}_i and $\mathbf{\Lambda} = \text{diag}(\lambda_1, \dots, \lambda_n)$ eigenvalues
- Change of basis: $\mathbf{w}^{[t]} = \mathbf{V}^T(\mathbf{x}^{[t]} - \mathbf{x}^*)$



sym \Rightarrow \mathbf{V} rotations
orthogonal

- 1) move optimum to origin
- 2) rotate (because \mathbf{V} is orthogonal)

GD STEPS IN EIGENSPACE

With $\mathbf{w}^{[t]} = \mathbf{V}^\top (\mathbf{x}^{[t]} - \mathbf{x}^*)$, a single GD step

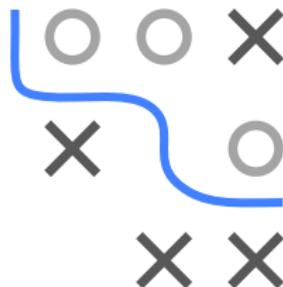
$$\mathbf{x}^{[t+1]} = \mathbf{x}^{[t]} - \alpha(2\mathbf{A}\mathbf{x}^{[t]} - \mathbf{b})$$

becomes

$$\mathbf{w}^{[t+1]} = \mathbf{w}^{[t]} - 2\alpha\mathbf{\Lambda}\mathbf{w}^{[t]}.$$

Therefore:

$$\begin{aligned}w_i^{[t+1]} &= w_i^{[t]} - 2\alpha\lambda_i w_i^{[t]} \\ &= (1 - 2\alpha\lambda_i)w_i^{[t]} \\ &= \dots \\ &= (1 - 2\alpha\lambda_i)^{t+1} w_i^{[0]}\end{aligned}$$



GD STEPS IN EIGENSPACE / 2

Proof (for $\mathbf{w}^{[t+1]} = \mathbf{w}^{[t]} - 2\alpha\Lambda\mathbf{w}^{[t]}$):

- A single GD step means

$$\mathbf{x}^{[t+1]} = \mathbf{x}^{[t]} - \alpha(2\mathbf{A}\mathbf{x}^{[t]} - \mathbf{b})$$

- Then:

$$\mathbf{V}^\top(\mathbf{x}^{[t+1]} - \mathbf{x}^*) = \mathbf{V}^\top(\mathbf{x}^{[t]} - \mathbf{x}^*) - \alpha\mathbf{V}^\top(2\mathbf{A}\mathbf{x}^{[t]} - \mathbf{b})$$

definition

$$\mathbf{w}^{[t+1]} = \mathbf{w}^{[t]} - \alpha\mathbf{V}^\top(2\mathbf{A}\mathbf{x}^{[t]} - \mathbf{b})$$

$$\mathbf{w}^{[t+1]} = \mathbf{w}^{[t]} - \alpha\mathbf{V}^\top(2\mathbf{A}(\mathbf{x}^{[t]} - \mathbf{x}^*) + \underbrace{2\mathbf{A}\mathbf{x}^* - \mathbf{b}}_{=0})$$

$$= \mathbf{w}^{[t]} - 2\alpha\Lambda\mathbf{V}^\top(\mathbf{x}^{[t]} - \mathbf{x}^*)$$

$$= \mathbf{w}^{[t]} - 2\alpha\Lambda\mathbf{w}^{[t]}$$

$$\Lambda = \mathbf{V}\Lambda\mathbf{V}^\top$$

$$\mathbf{V}^\top\mathbf{V}\Lambda\mathbf{V}^\top = \mathbf{H}$$



subt. and apply

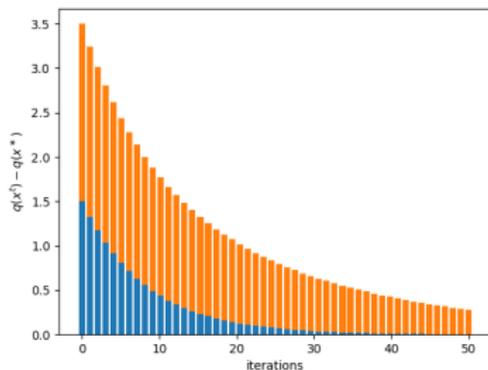
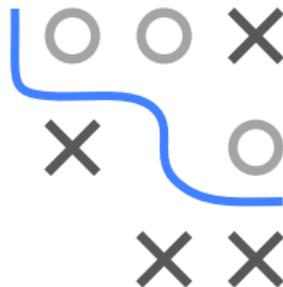
add and subtract $2\mathbf{A}\mathbf{x}^*$

because optimal solution also gradient at the optimum must be 0

GD ERROR IN ORIGINAL SPACE / 2

We now consider the contribution of each eigenvector to the total loss

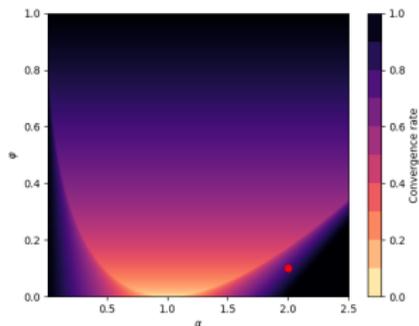
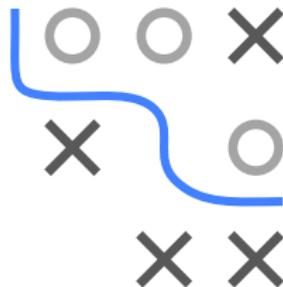
$$q(\mathbf{x}^{[t]}) - q(\mathbf{x}^*) = \frac{1}{2} \sum_i^d (1 - 2\alpha\lambda_i)^{2t} \lambda_i (w_i^{[0]})^2$$



Optimization in Machine Learning

First order methods

Momentum on quadratic forms



Learning goals

- Momentum update in Eigenspace
- Effect of φ

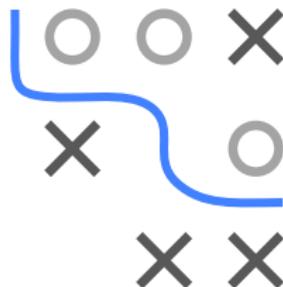
RECAP: MOMENTUM UPDATE

$$\begin{aligned}\boldsymbol{\nu}^{[t+1]} &= \varphi \boldsymbol{\nu}^{[t]} + \alpha \nabla f(\mathbf{x}^{[t]}) \\ \mathbf{x}^{[t+1]} &= \mathbf{x}^{[t]} - \boldsymbol{\nu}^{[t+1]},\end{aligned}$$

which simplifies to

$$\begin{aligned}\boldsymbol{\nu}^{[t+1]} &= \varphi \boldsymbol{\nu}^{[t]} + \alpha(\mathbf{A}\mathbf{x}^{[t]} - \mathbf{b}) \\ \mathbf{x}^{[t+1]} &= \mathbf{x}^{[t]} - \boldsymbol{\nu}^{[t+1]},\end{aligned}$$

for the quadratic form.



RECAP: DYNAMICS OF MOMENTUM

Changing the basis as before with $\mathbf{w}^{[t]} = \mathbf{V}^\top (\mathbf{x}^{[t]} - \mathbf{x}^*)$ and $\mathbf{u}^{[t]} = \mathbf{V}\nu^{[t]}$, we get the following set of equations, where each component acts independently, although $w_i^{[t]}$ and $u_i^{[t]}$ are coupled:

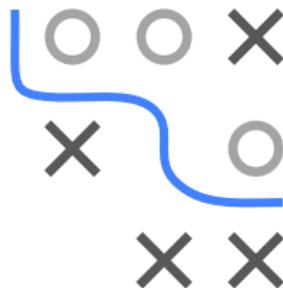
$$\begin{aligned}u_i^{[t+1]} &= \varphi u_i^{[t]} + \alpha \lambda_i w_i^{[t]}, \\w_i^{[t+1]} &= w_i^{[t]} - u_i^{[t+1]}\end{aligned}$$

We rewrite this:

$$\begin{pmatrix} 1 & 0 \\ 1 & 1 \end{pmatrix} \begin{pmatrix} u_i^{[t+1]} \\ w_i^{[t+1]} \end{pmatrix} = \begin{pmatrix} \varphi & \alpha \lambda_i \\ 0 & 1 \end{pmatrix} \begin{pmatrix} u_i^{[t]} \\ w_i^{[t]} \end{pmatrix}$$

and invert the matrix on the LHS:

$$\begin{pmatrix} u_i^{[t+1]} \\ w_i^{[t+1]} \end{pmatrix} = \begin{pmatrix} \varphi & \alpha \lambda_i \\ -\varphi & 1 - \alpha \lambda_i \end{pmatrix} \begin{pmatrix} u_i^{[t]} \\ w_i^{[t]} \end{pmatrix} = \mathbf{R}^{t+1} \begin{pmatrix} u_i^0 \\ w_i^0 \end{pmatrix}$$



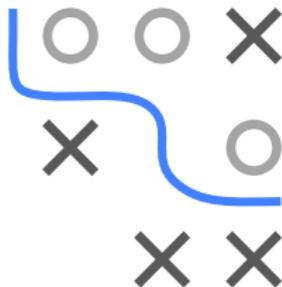
RECAP: DYNAMICS OF MOMENTUM / 2

Taking a 2×2 matrix to the t^{th} power reduces to a formula involving the eigenvalues of R , σ_1 and σ_2 :

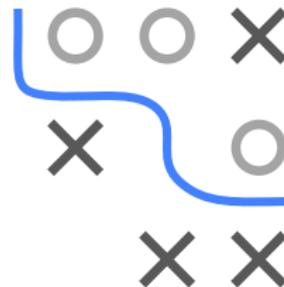
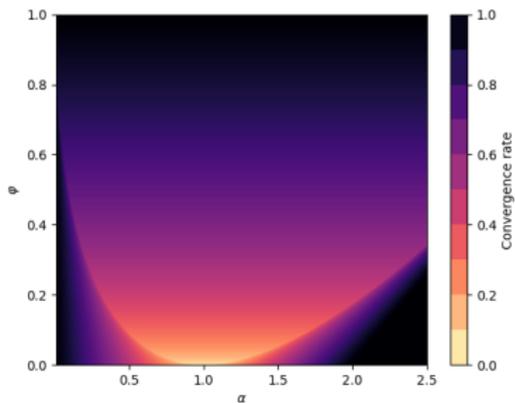
$$R^t = \begin{cases} \sigma_1^t R_1 - \sigma_2^t R_2, & \text{if } \sigma_1 \neq \sigma_2 \\ \sigma_1^t (tR/\sigma_1 - (t-1)I), & \text{if } \sigma_1 = \sigma_2 \end{cases}$$

where $R_j = \frac{R - \sigma_j I}{\sigma_1 - \sigma_2}$.

In contrast to GD, where we got one geometric series, we have two coupled series with real or complex values.



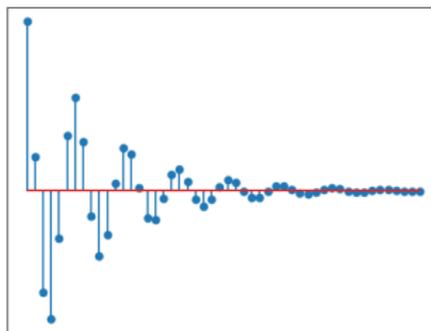
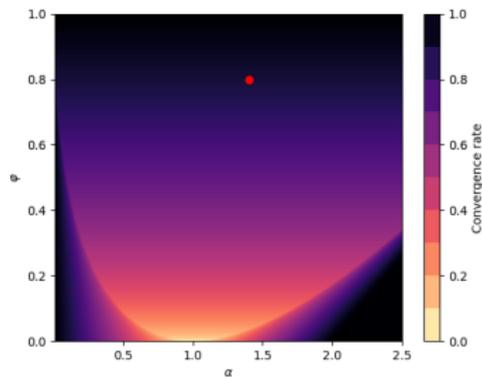
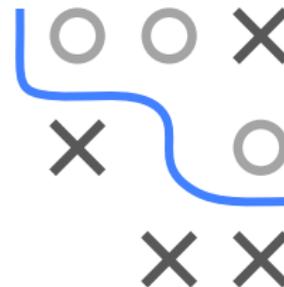
RECAP: DYNAMICS OF MOMENTUM / 3



The achieved convergence rate is therefore the slowest of the two, $\max\{|\sigma_1|, |\sigma_2|\}$.

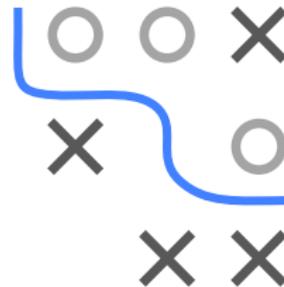
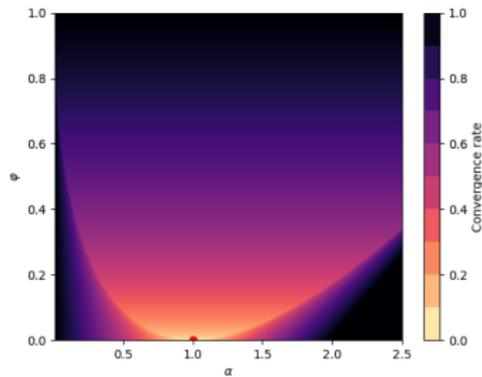
Each region shows a different convergence behavior.

RECAP: DYNAMICS OF MOMENTUM / 4



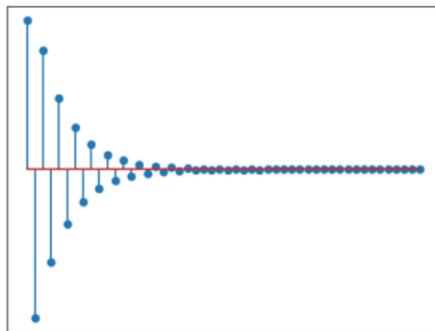
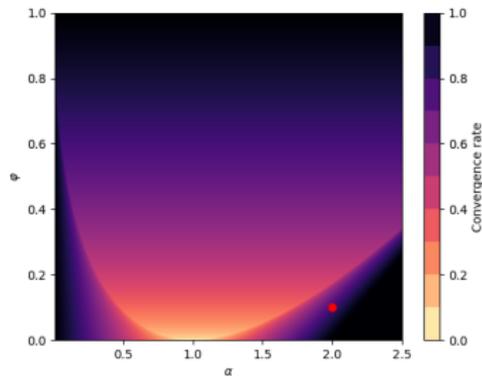
The eigenvalues of R are complex and we see low frequency ripples.

RECAP: DYNAMICS OF MOMENTUM / 6

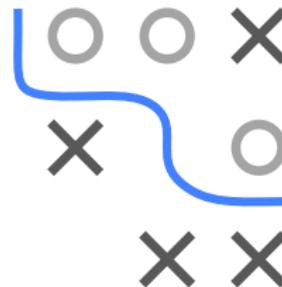


The step size is $\alpha = 1/\lambda_i$ and $\varphi = 0$ - we converge in one step.

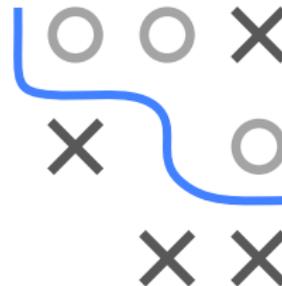
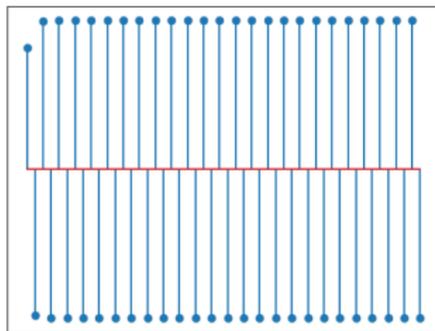
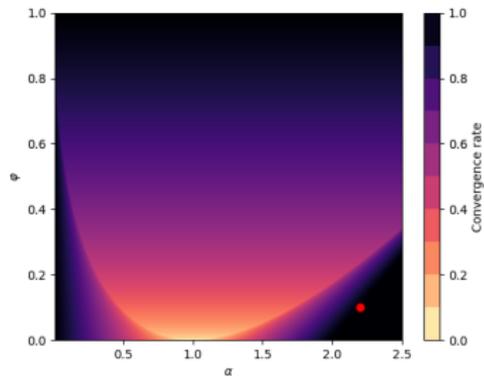
RECAP: DYNAMICS OF MOMENTUM / 7



When $\alpha > 1/\lambda_i$, the iterates flip sign every iteration.



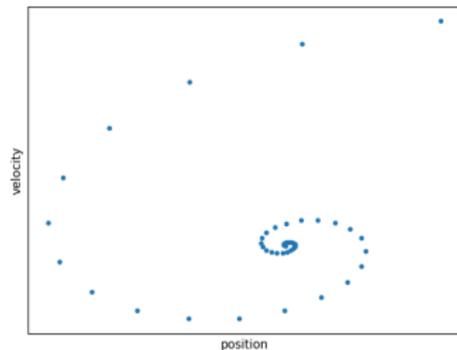
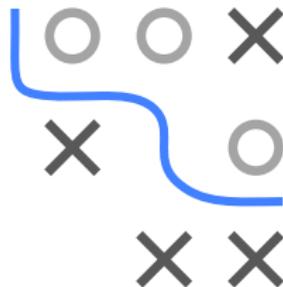
RECAP: DYNAMICS OF MOMENTUM / 8



If $\max\{|\sigma_1|, |\sigma_2|\} > 1$, the iterates diverge.

RECAP: DYNAMICS OF MOMENTUM / 9

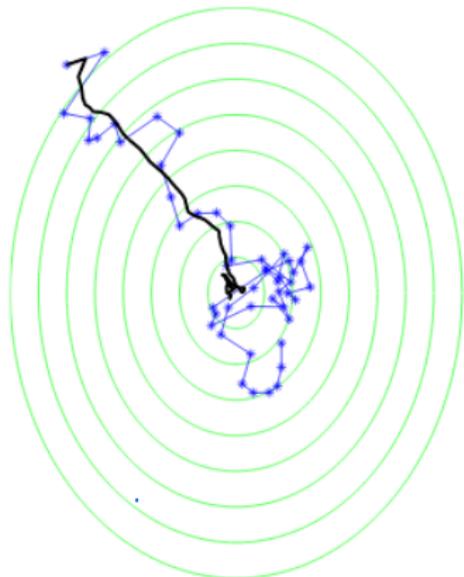
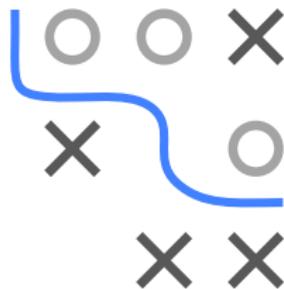
- Finally, we investigate the role of φ .
- We can think of gradient descent with momentum as a damped harmonic oscillator: a weight on a spring. We pull the weight down and study the path back to the equilibrium in phase space (looking at the position and the velocity).
- Depending on the choice of φ , the rate of return to the equilibrium position is affected.



Optimization in Machine Learning

First order methods

SGD



Learning goals

- SGD
- Stochasticity
- Convergence
- Batch size

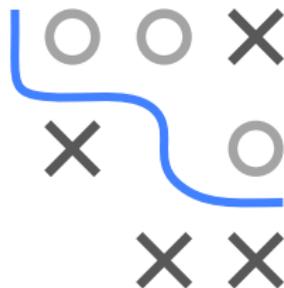
STOCHASTIC GRADIENT DESCENT

NB: We use g instead of f as objective, bc. f is used as model in ML.

$g : \mathbb{R}^d \rightarrow \mathbb{R}$ objective, g **average over functions**:

$$g(\mathbf{x}) = \frac{1}{n} \sum_{i=1}^n g_i(\mathbf{x}), \quad g \text{ and } g_i \text{ smooth}$$

loss at datapoint i



Stochastic gradient descent (SGD) approximates the gradient

$$\begin{aligned} \nabla_{\mathbf{x}} g(\mathbf{x}) &= \frac{1}{n} \sum_{i=1}^n \nabla_{\mathbf{x}} g_i(\mathbf{x}) := \mathbf{d} \quad \text{by} \\ &\frac{1}{|J|} \sum_{i \in J} \nabla_{\mathbf{x}} g_i(\mathbf{x}) := \hat{\mathbf{d}}, \end{aligned}$$

*in the limit we'll have
normally distrib grads.
CLT helps,
Didn't understand this but
sounds cool*

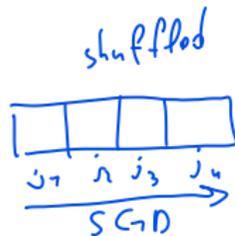
with random subset $J \subset \{1, 2, \dots, n\}$ of gradients called **mini-batch**.

This is done e.g. when computing the true gradient is **expensive**.

STOCHASTIC GRADIENT DESCENT / 2

Algorithm Basic SGD pseudo code

- 1: Initialize $\mathbf{x}^{[0]}$, $t = 0$
- 2: **while** stopping criterion not met **do**
- 3: Randomly shuffle indices and partition into minibatches J_1, \dots, J_K of size m
- 4: **for** $k \in \{1, \dots, K\}$ **do**
- 5: $t \leftarrow t + 1$
- 6: Compute gradient estimate with J_k : $\hat{\mathbf{d}}^{[t]} \leftarrow \frac{1}{m} \sum_{i \in J_k} \nabla_{\mathbf{x}} g_i(\mathbf{x}^{[t-1]})$
- 7: Apply update: $\mathbf{x}^{[t]} \leftarrow \mathbf{x}^{[t-1]} - \alpha \cdot \hat{\mathbf{d}}^{[t]}$
- 8: **end for**
- 9: **end while**



- Instead of drawing batches randomly we might want to go through the g_i sequentially (unless g_i are sorted in any way)
- Updates are computed faster, but also more stochastic:
 - In the simplest case, batch-size $m := |J_k|$ is set to $m = 1$
 - If n is a billion, computation of update is a billion times faster
 - **But** (later): Convergence rates suffer from stochasticity!

SGD IN ML

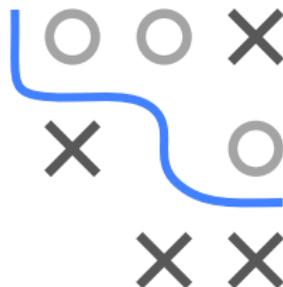
In ML, we perform ERM:

$$\mathcal{R}(\theta) = \frac{1}{n} \sum_{i=1}^n \underbrace{L(y^{(i)}, f(\mathbf{x}^{(i)} | \theta))}_{g_i(\theta)}$$

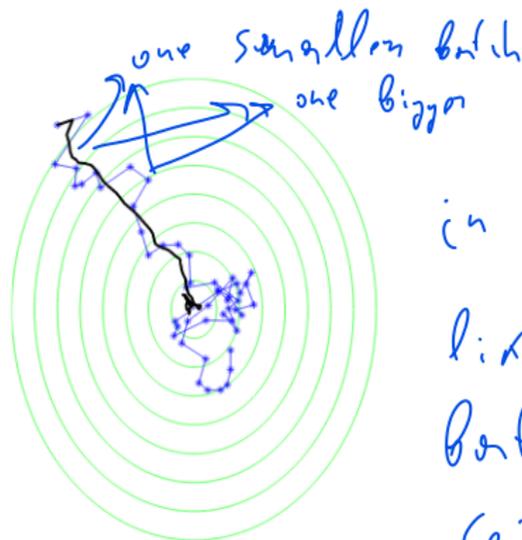
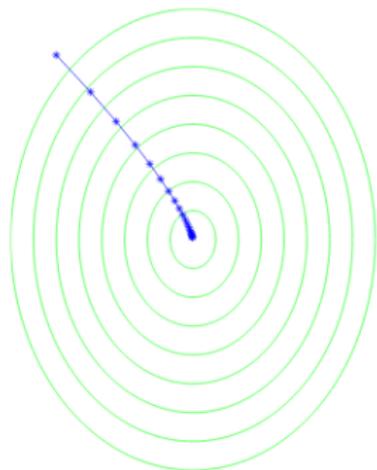
- for a data set

$$\mathcal{D} = \left((\mathbf{x}^{(1)}, y^{(1)}), \dots, (\mathbf{x}^{(n)}, y^{(n)}) \right)$$

- a loss function $L(y, f(\mathbf{x}))$, e.g., L2 loss $L(y, f(\mathbf{x})) = (y - f(\mathbf{x}))^2$,
- and a model class f , e.g., the linear model $f(\mathbf{x}^{(i)} | \theta) = \theta^\top \mathbf{x}$.



STOCHASTICITY OF SGD



$$\text{Minimize } g(x_1, x_2) = 1.25(x_1 + 6)^2 + (x_2 - 8)^2.$$

Left: GD. **Right:** SGD. Black line shows average value across multiple runs.

(Source: Shalev-Shwartz et al., Understanding Machine Learning, 2014.)

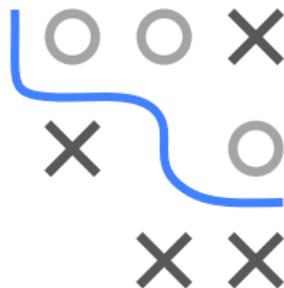
STOCHASTICITY OF SGD / 2

Assume batch size $m = 1$ (statements also apply for larger batches).

- **(Possibly) suboptimal direction:** Approximate gradient $\hat{\mathbf{d}} = \nabla_{\mathbf{x}}g_j(\mathbf{x})$ might point in suboptimal (possibly not even a descent!) direction
- **Unbiased estimate:** If J drawn i.i.d., approximate gradient $\hat{\mathbf{d}}$ is an unbiased estimate of gradient $\mathbf{d} = \nabla_{\mathbf{x}}g(\mathbf{x}) = \sum_{i=1}^n \nabla_{\mathbf{x}}g_i(\mathbf{x})$:

$$\begin{aligned}\mathbb{E}_j [\nabla_{\mathbf{x}}g_j(\mathbf{x})] &= \sum_{i=1}^n \nabla_{\mathbf{x}}g_i(\mathbf{x}) \cdot \mathbb{P}(i = j) \\ &= \sum_{i=1}^n \nabla_{\mathbf{x}}g_i(\mathbf{x}) \cdot \frac{1}{n} = \nabla_{\mathbf{x}}g(\mathbf{x}).\end{aligned}$$

momentum helps a lot



Conclusion: SGD might perform single suboptimal moves, but moves in “right direction” **on average**.

CONVERGENCE OF SGD

As a consequence, SGD has worse convergence properties than GD.

But: Can be controlled via **increasing batches** or **reducing step size**.

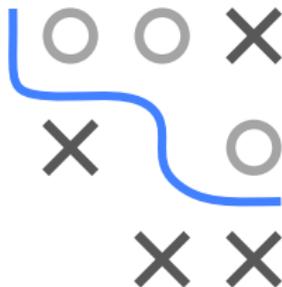
The larger the batch size m

- the better the approximation to $\nabla_{\mathbf{x}}g(\mathbf{x})$
- the lower the variance
- the lower the risk of performing steps in the wrong direction

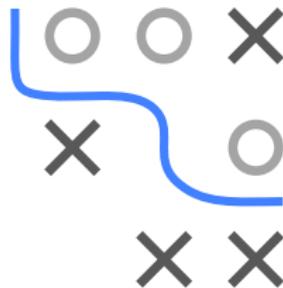
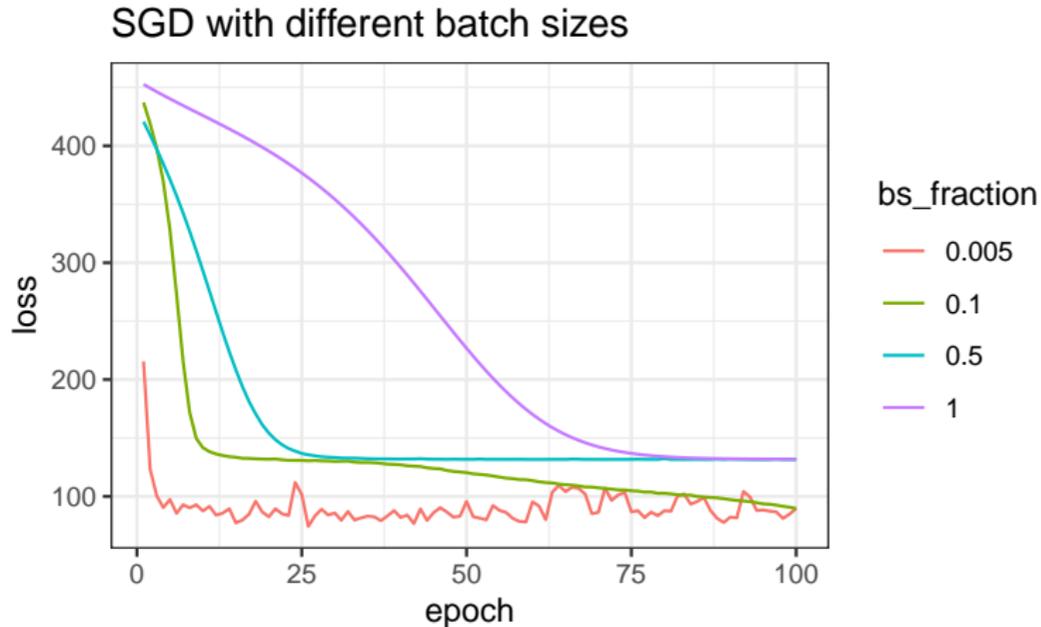
The smaller the step size α

- the smaller a step in a potentially wrong direction
- the lower the effect of high variance

As maximum batch size is usually limited by computational resources (memory), choosing the step size is crucial.



EFFECT OF BATCH SIZE

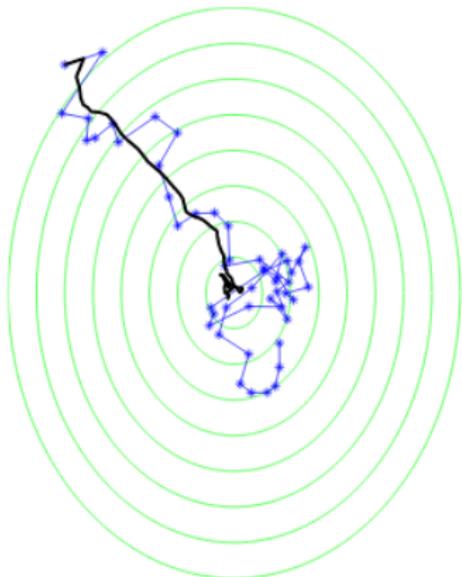


SGD for a NN with batch size $\in \{0.5\%, 10\%, 50\%\}$ of the training data.
The higher the batch size, the lower the variance.

Optimization in Machine Learning

First order methods

SGD Further Details

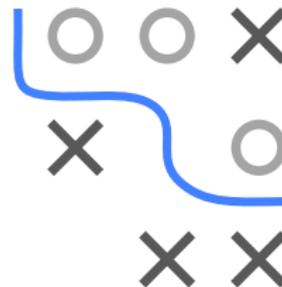
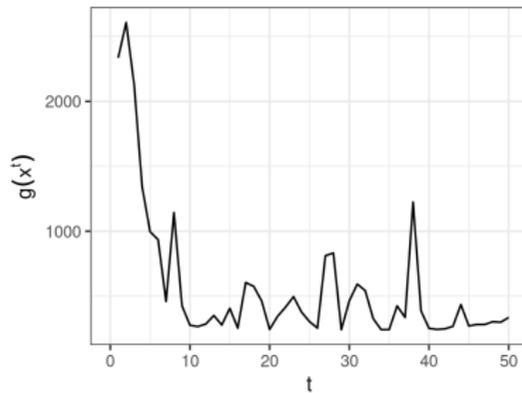
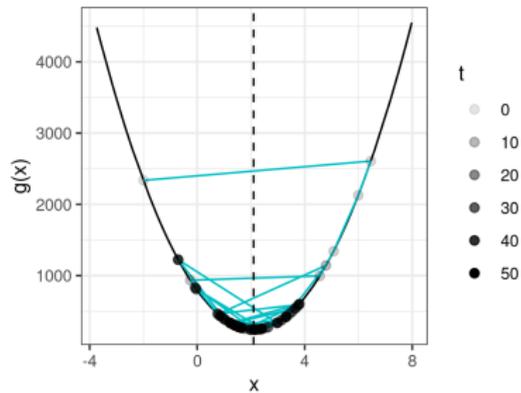


Learning goals

- Decreasing step size for SGD
- Stopping rules
- SGD with momentum

SGD WITH CONSTANT STEP SIZE

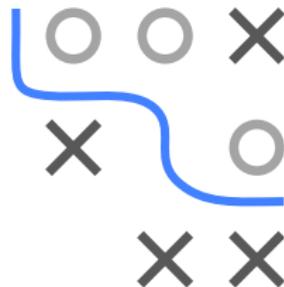
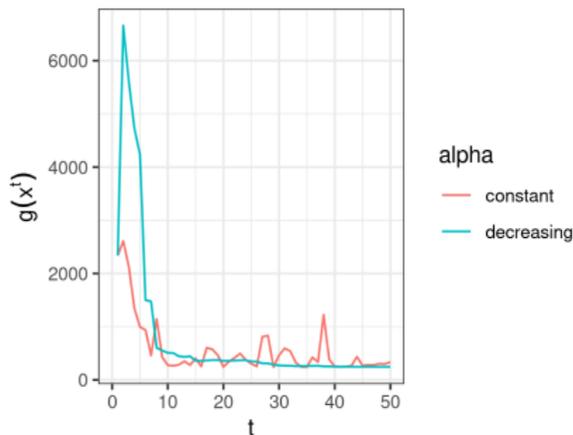
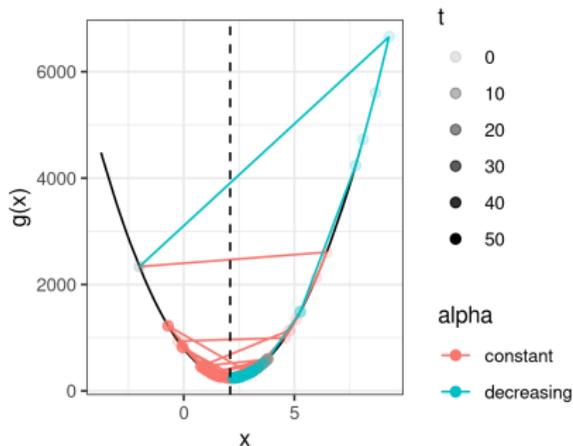
Example: SGD with constant step size.



Fast convergence of SGD initially. Erratic behavior later (variance too big).

SGD WITH DECREASING STEP SIZE / 2

- Popular solution: step size fulfilling $\alpha^{[t]} \in \mathcal{O}(1/t)$.



Example continued. Step size $\alpha^{[t]} = 0.2/t$.

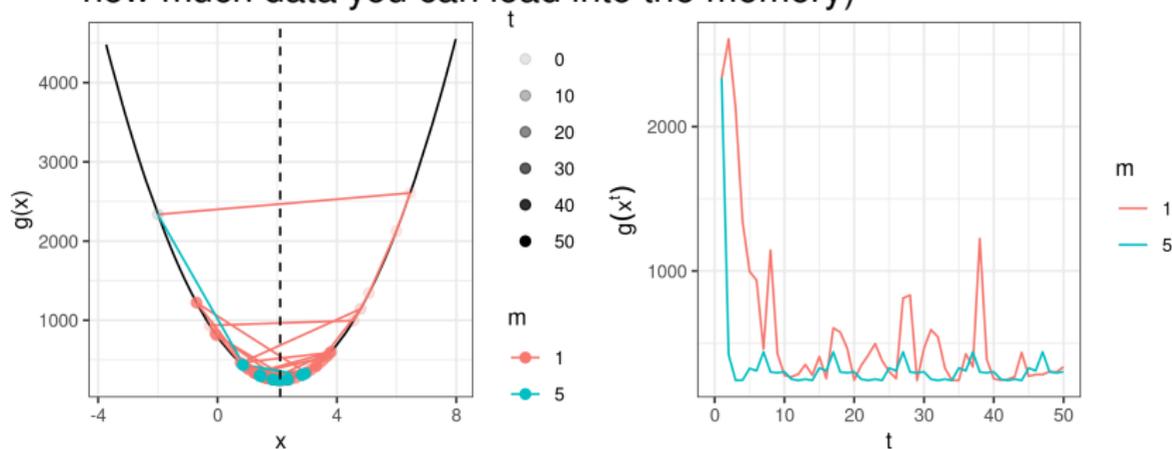
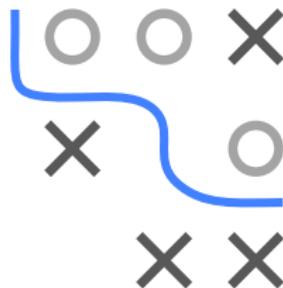
- Often not working well in practice: step size gets small quite fast.
- Alternative: $\alpha^{[t]} \in \mathcal{O}(1/\sqrt{t})$

MINI-BATCHES

- Reduce noise by increasing batch size m for better approximation

$$\hat{\mathbf{d}} = \frac{1}{m} \sum_{i \in J} \nabla_{\mathbf{x}} g_i(\mathbf{x}) \approx \frac{1}{n} \sum_{i=1}^n \nabla_{\mathbf{x}} g_i(\mathbf{x}) = \mathbf{d}$$

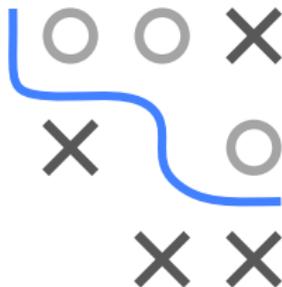
- Usually, the batch size is limited by computational resources (e.g., how much data you can load into the memory)



Example continued. Batch size $m = 1$ vs. $m = 5$.

STOPPING RULES FOR SGD

- **For GD:** We usually stop when gradient is close to 0 (i.e., we are close to a stationary point)
- **For SGD:** individual gradients do not necessarily go to zero, and we cannot access full gradient.
- Practicable solution for ML:
 - Measure the validation set error after T iterations
 - Stop if validation set error is not improving



SGD AND ML

General remarks:

- SGD is a variant of GD
- SGD particularly suitable for large-scale ML when evaluating gradient is too expensive / restricted by computational resources
- SGD and variants are the most commonly used methods in modern ML, for example:
 - Linear models
Note that even for the linear model and quadratic loss, where a closed form solution is available, SGD might be used if the size n of the dataset is too large and the design matrix does not fit into memory.
 - Neural networks
 - Support vector machines
 - ...

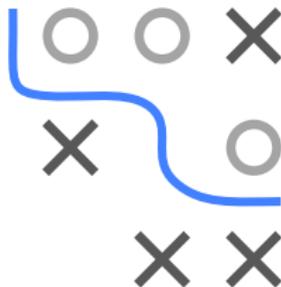


SGD WITH MOMENTUM

SGD is usually used with momentum due to reasons mentioned in previous chapters.

Algorithm Stochastic gradient descent with momentum

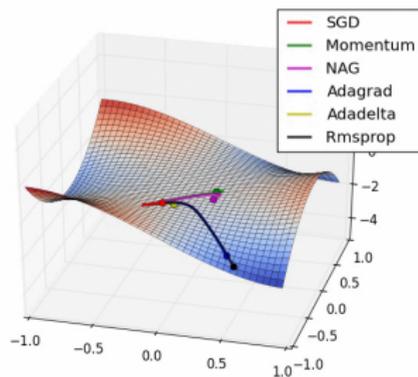
- 1: **require** step size α and momentum φ
 - 2: **require** initial parameter \mathbf{x} and initial velocity $\boldsymbol{\nu}$
 - 3: **while** stopping criterion not met **do**
 - 4: Sample mini-batch of m examples
 - 5: Compute gradient estimate $\nabla \hat{g}(\mathbf{x})$ using mini-batch
 - 6: Compute velocity update: $\boldsymbol{\nu} \leftarrow \varphi \boldsymbol{\nu} - \alpha \nabla \hat{g}(\mathbf{x})$
 - 7: Apply update: $\mathbf{x} \leftarrow \mathbf{x} + \boldsymbol{\nu}$
 - 8: **end while**
-



Optimization in Machine Learning

First order methods

Adam and friends

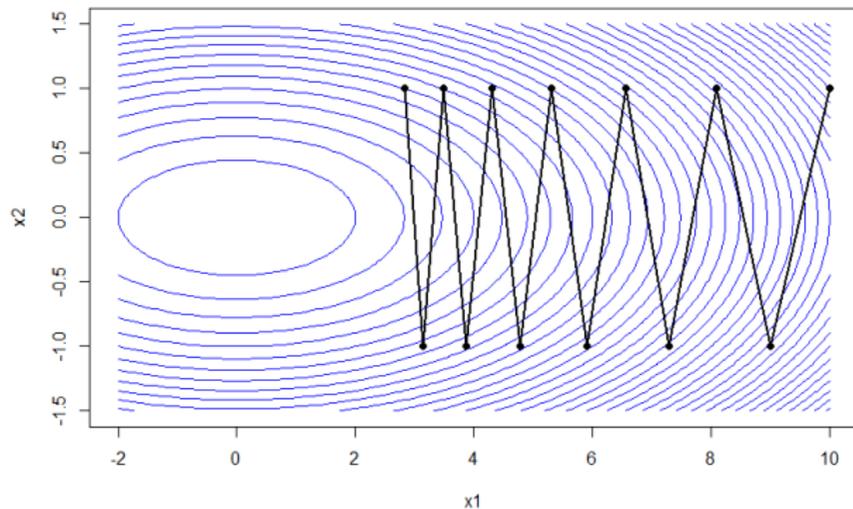
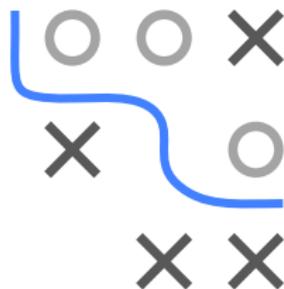


Learning goals

- Adaptive step sizes
- AdaGrad
- RMSProp
- Adam

ADAPTIVE STEP SIZES

- Step size is probably the most important control parameter
- Has strong influence on performance
- Natural to use different step size for each input individually and automatically adapt them

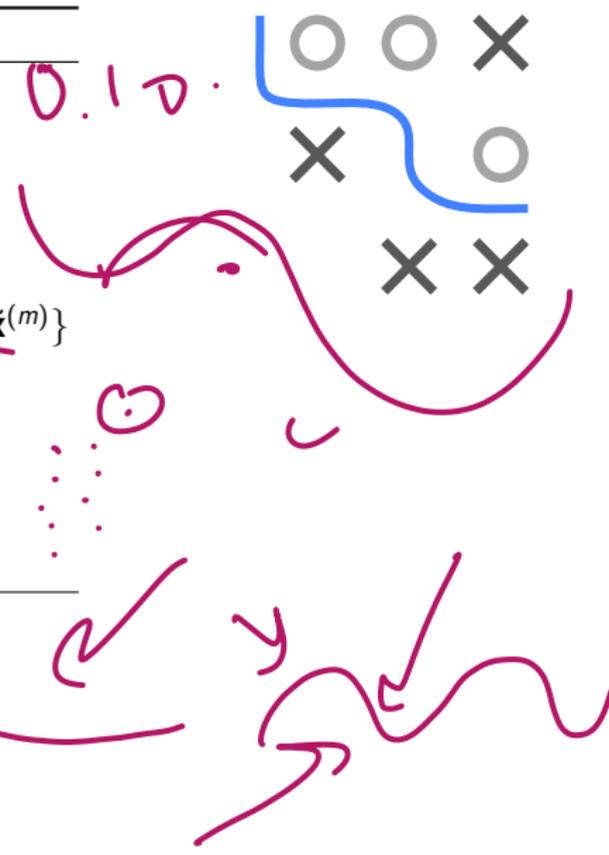


ADAGRAD / 2

Algorithm AdaGrad

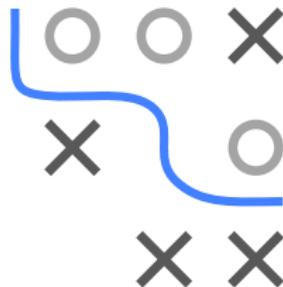
- 1: **require** Global step size α
- 2: **require** Initial parameter θ
- 3: **require** Small constant β , perhaps 10^{-7} , for numerical stability
- 4: **Initialize** gradient accumulation variable $\mathbf{r} = \mathbf{0}$
- 5: **while** stopping criterion not met **do**
- 6: Sample minibatch of m examples from the training set $\{\tilde{\mathbf{x}}^{(1)}, \dots, \tilde{\mathbf{x}}^{(m)}\}$
- 7: Compute gradient estimate: $\hat{\mathbf{g}} \leftarrow \frac{1}{m} \nabla_{\theta} \sum_{i=1}^m L(y^{(i)}, f(\tilde{\mathbf{x}}^{(i)} | \theta))$
- 8: Accumulate squared gradient $\mathbf{r} \leftarrow \mathbf{r} + \hat{\mathbf{g}} \odot \hat{\mathbf{g}}$
- 9: Compute update: $\nabla \theta = -\frac{\alpha}{\beta + \sqrt{\mathbf{r}}} \odot \hat{\mathbf{g}}$ (operations element-wise)
- 10: Apply update: $\theta \leftarrow \theta + \nabla \theta$
- 11: **end while**

\odot : element-wise product (Hadamard)



RMSPROP

- Modification of AdaGrad
- Resolves AdaGrad's radically diminishing step sizes.
- Gradient accumulation is replaced by exponentially weighted moving average
- Theoretically, leads to performance gains in non-convex scenarios
- Empirically, RMSProp is a very effective optimization algorithm. Particularly, it is employed routinely by DL practitioners.



$$\beta \cdot v^t + (1 - \beta) v^T$$

$$\beta^n$$

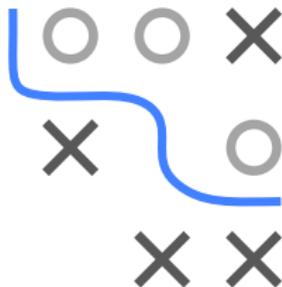
...

$$\beta$$

RMSPROP / 2

Algorithm RMSProp

- 1: **require** Global step size α and decay rate $\rho \in [0, 1)$
- 2: **require** Initial parameter θ
- 3: **require** Small constant β , perhaps 10^{-6} , for numerical stability
- 4: Initialize gradient accumulation variable $\mathbf{r} = \mathbf{0}$
- 5: **while** stopping criterion not met **do**
- 6: Sample minibatch of m examples from the training set $\{\tilde{\mathbf{x}}^{(1)}, \dots, \tilde{\mathbf{x}}^{(m)}\}$
- 7: Compute gradient estimate: $\hat{\mathbf{g}} \leftarrow \frac{1}{m} \nabla_{\theta} \sum_i L(y^{(i)}, f(\tilde{\mathbf{x}}^{(i)} | \theta))$
- 8: Accumulate squared gradient $\mathbf{r} \leftarrow \rho \mathbf{r} + (1 - \rho) \hat{\mathbf{g}} \odot \hat{\mathbf{g}}$
- 9: Compute update: $\nabla \theta = -\frac{\alpha}{\beta + \sqrt{\mathbf{r}}} \odot \hat{\mathbf{g}}$
- 10: Apply update: $\theta \leftarrow \theta + \nabla \theta$
- 11: **end while**



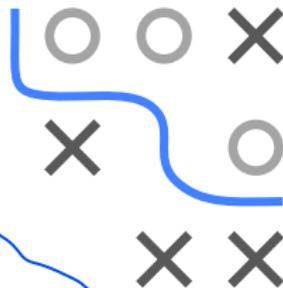
ADA is ma
RMSProp

we accumulate the squared grad differently, by decaying the previous ones.
grads of amount of emphasis on new and old gradients

ADAM / 2

Algorithm Adam

- 1: **require** Global step size α (suggested default: 0.001)
- 2: **require** Exponential decay rates for moment estimates, ρ_1 and ρ_2 in $[0, 1)$ (suggested defaults: 0.9 and 0.999 respectively)
- 3: **require** Small constant β (suggested default 10^{-8})
- 4: **require** Initial parameters θ
- 5: Initialize time step $t = 0$
- 6: Initialize 1st and 2nd moment variables $\mathbf{s}^{[0]} = \mathbf{0}, \mathbf{r}^{[0]} = \mathbf{0}$
- 7: **while** stopping criterion not met **do**
- 8: $t \leftarrow t + 1$
- 9: Sample a minibatch of m examples from the training set $\{\tilde{\mathbf{x}}^{(1)}, \dots, \tilde{\mathbf{x}}^{(m)}\}$
- 10: Compute gradient estimate: $\hat{\mathbf{g}}^{[t]} \leftarrow \frac{1}{m} \nabla_{\theta} \sum_i L(y^{(i)}, f(\tilde{\mathbf{x}}^{(i)} | \theta))$
- 11: Update biased first moment estimate: $\mathbf{s}^{[t]} \leftarrow \rho_1 \mathbf{s}^{[t-1]} + (1 - \rho_1) \hat{\mathbf{g}}^{[t]}$
- 12: Update biased second moment estimate: $\mathbf{r}^{[t]} \leftarrow \rho_2 \mathbf{r}^{[t-1]} + (1 - \rho_2) \hat{\mathbf{g}}^{[t]} \odot \hat{\mathbf{g}}^{[t]}$
- 13: Correct bias in first moment: $\hat{\mathbf{s}} \leftarrow \frac{\mathbf{s}^{[t]}}{1 - \rho_1^t}$
- 14: Correct bias in second moment: $\hat{\mathbf{r}} \leftarrow \frac{\mathbf{r}^{[t]}}{1 - \rho_2^t}$
- 15: Compute update: $\nabla \theta = -\alpha \frac{\hat{\mathbf{s}}}{\sqrt{\hat{\mathbf{r}} + \beta}}$
- 16: Apply update: $\theta \leftarrow \theta + \nabla \theta$
- 17: **end while**



More expensive but not really, we just have extra 2 vectors. and we don't have second-order stuff.

Momentum \rightarrow $\mathbf{s}^{[t]}$
RMS Prop \rightarrow $\hat{\mathbf{r}}^{[t]}$

$\mathbf{s}^{[0]} = \mathbf{0}, \mathbf{r}^{[0]} = \mathbf{0}$

this later

$$\frac{1}{10} \frac{1}{1 - \rho^t}$$

$t=10$
 $\rho=0.9$

$\{m, n\}$

ADAM / 3

- Initializes moment variables \mathbf{s} and \mathbf{r} with zero \Rightarrow Bias towards zero

$$\mathbb{E}[\mathbf{s}^{[t]}] \neq \mathbb{E}[\hat{\mathbf{g}}^{[t]}] \quad \text{and} \quad \mathbb{E}[\mathbf{r}^{[t]}] \neq \mathbb{E}[\hat{\mathbf{g}}^{[t]} \odot \hat{\mathbf{g}}^{[t]}]$$

(\mathbb{E} calculated over minibatches)

Own I, II moment estimates aren't unbiased estimates of gradient

- Indeed: Unrolling $\mathbf{s}^{[t]}$ yields

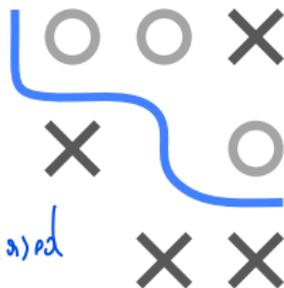
cancel out
cancel $\mathbf{s}^{[0]} = 0$

$$\mathbf{s}^{[1]} = \rho_1 \mathbf{s}^{[0]} + (1 - \rho_1) \hat{\mathbf{g}}^{[1]} = (1 - \rho_1) \hat{\mathbf{g}}^{[1]}$$

$$\mathbf{s}^{[2]} = \rho_1 \mathbf{s}^{[1]} + (1 - \rho_1) \hat{\mathbf{g}}^{[2]} = \rho_1 (1 - \rho_1) \hat{\mathbf{g}}^{[1]} + (1 - \rho_1) \hat{\mathbf{g}}^{[2]}$$

$$\mathbf{s}^{[3]} = \rho_1 \mathbf{s}^{[2]} + (1 - \rho_1) \hat{\mathbf{g}}^{[3]} = \rho_1^2 (1 - \rho_1) \hat{\mathbf{g}}^{[1]} + \rho_1 (1 - \rho_1) \hat{\mathbf{g}}^{[2]} + (1 - \rho_1) \hat{\mathbf{g}}^{[3]}$$

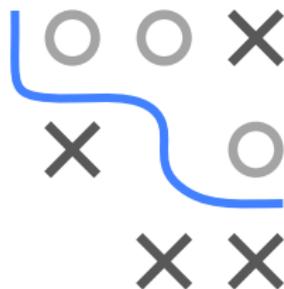
- Therefore: $\mathbf{s}^{[t]} = (1 - \rho_1) \sum_{i=1}^t \rho_1^{t-i} \hat{\mathbf{g}}^{[i]}$. *(can prove by induction)*
- **Note:** Contributions of past $\hat{\mathbf{g}}^{[i]}$ decreases rapidly



ADAM / 4

- We continue with

$$\begin{aligned}\mathbb{E}[\mathbf{s}^{[t]}] &= \mathbb{E}\left[\left(1 - \rho_1\right) \sum_{i=1}^t \rho_1^{t-i} \hat{\mathbf{g}}^{[i]}\right] \\ &= \mathbb{E}[\hat{\mathbf{g}}^{[t]}] \left(1 - \rho_1\right) \sum_{i=1}^t \rho_1^{t-i} + \zeta \\ &= \mathbb{E}[\hat{\mathbf{g}}^{[t]}] \left(1 - \rho_1^t\right) + \zeta,\end{aligned}$$



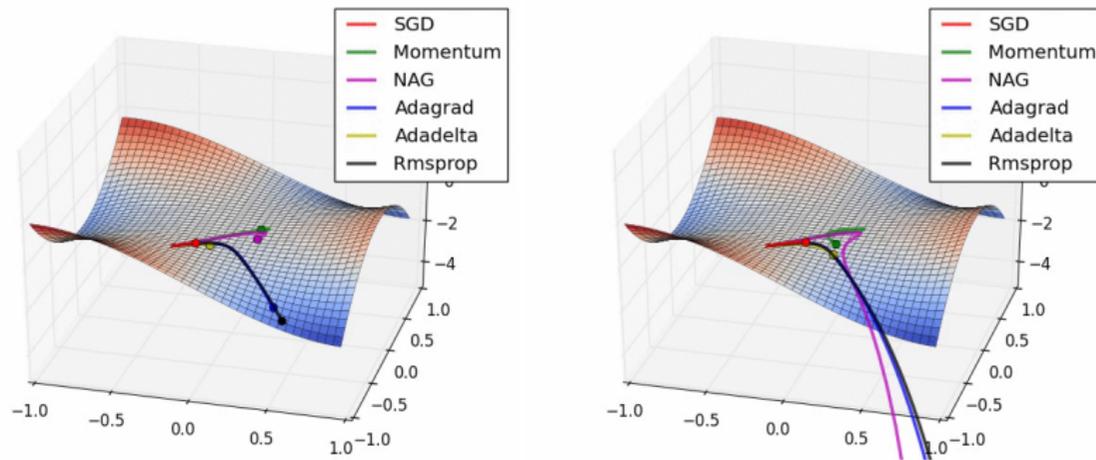
assume all grads are the same

where we approximated $\hat{\mathbf{g}}^{[i]}$ by $\hat{\mathbf{g}}^{[t]}$. The resulting error is put in ζ and be kept small due to the exponential weights of past gradients.

- Therefore: $\mathbf{s}^{[t]}$ is a biased estimator of $\hat{\mathbf{g}}^{[t]}$
- But bias vanishes for $t \rightarrow \infty$ ($\rho_1^t \rightarrow 0$)
- Ignoring ζ , we correct for the bias by $\hat{\mathbf{s}}^{[t]} = \frac{\mathbf{s}^{[t]}}{(1-\rho_1^t)}$
- Analogously: $\hat{\mathbf{r}}^{[t]} = \frac{\mathbf{r}^{[t]}}{(1-\rho_2^t)}$

*why we want unbiasedness
seems counterproductive...*

COMPARISON OF OPTIMIZERS: ANIMATION



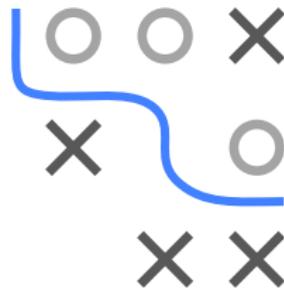
Credits: Dettmers (2015) and Radford

Comparison of SGD optimizers near saddle point.

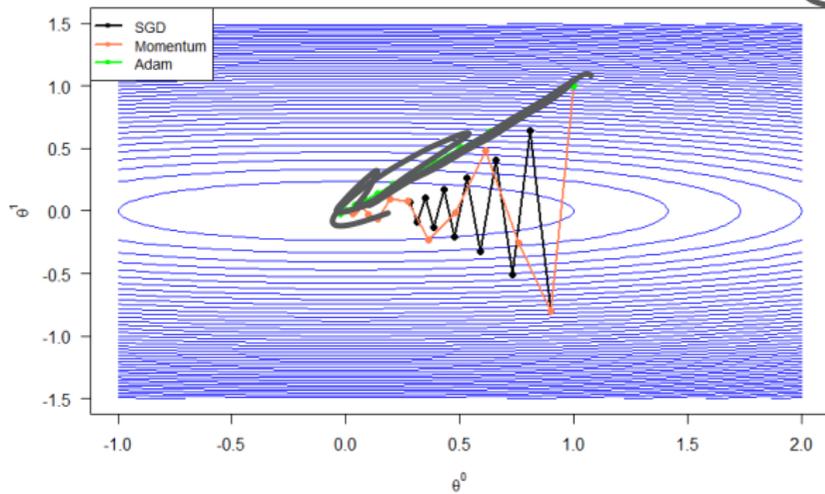
Left: After start. **Right:** Later.

All methods accelerate compared to vanilla SGD.

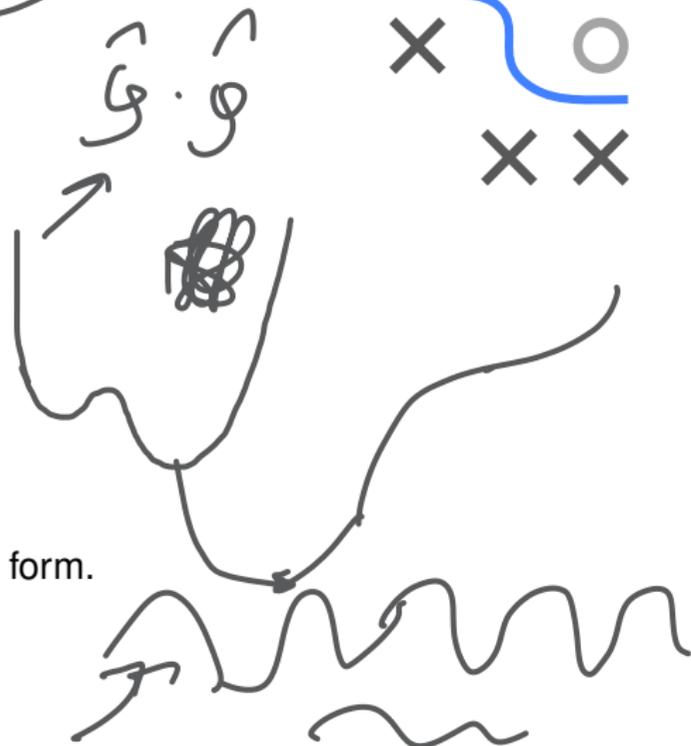
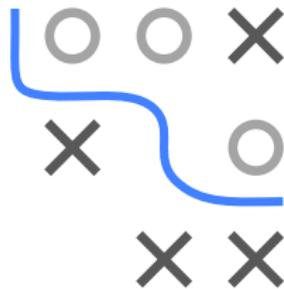
Best is RMSProp, then AdaGrad. (Adam is missing here.)



COMPARISON ON QUADRATIC FORM



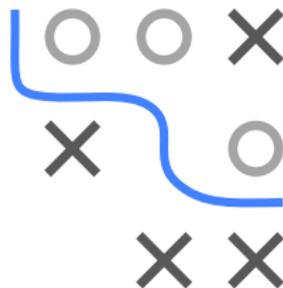
SGD vs. SGD with Momentum vs. Adam on a quadratic form.



Optimization in Machine Learning

First order methods

Comparison of first order methods

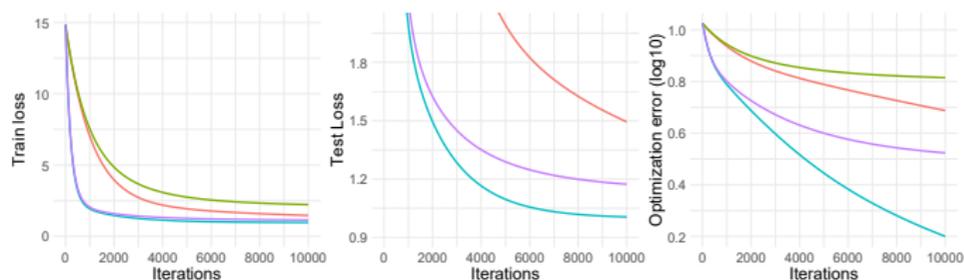


Learning goals

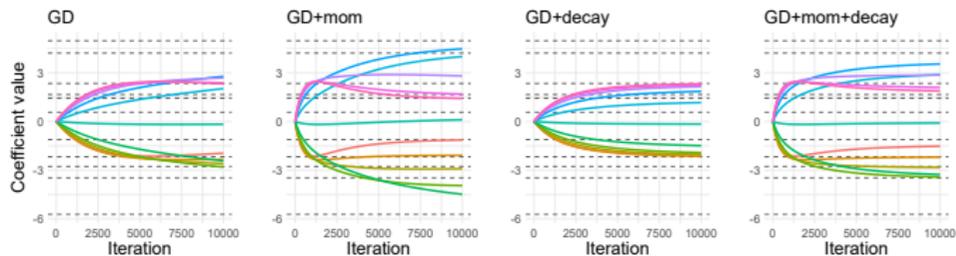
- Gradient Descent
- Stochastic Gradient Descent
- Momentum
- Step size decay

LIN-REG (GD + CORR. FEATURES)

GD with medium $\alpha = 2 \cdot 10^{-3}$ and bad conditioning (corr. features):



Method — GD — GD+decay — GD+mom — GD+mom+decay



→ $\hat{\theta}$ quite far from θ^*

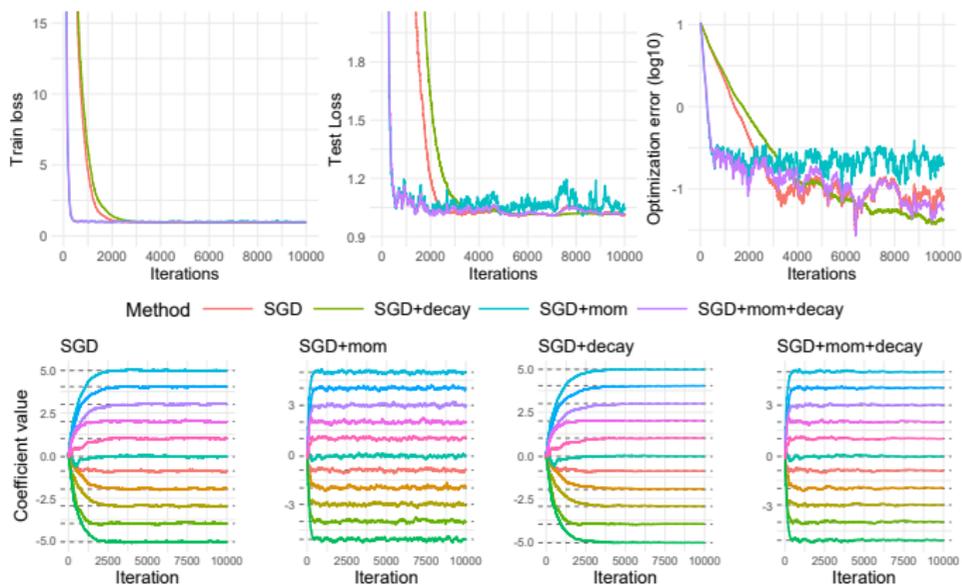
Irreducible loss due to additive noise is $\sigma = 1$. Dotted lines indicate global minimizers.

Bad **conditioning slows down optim** severely! Only **momentum** w/o decay comes close to global min. **Momentum** causes “overshooting” of some coefs. Corr. features cause corr. global minimizers.



LIN-REG (SGD + MED STEP SIZE)

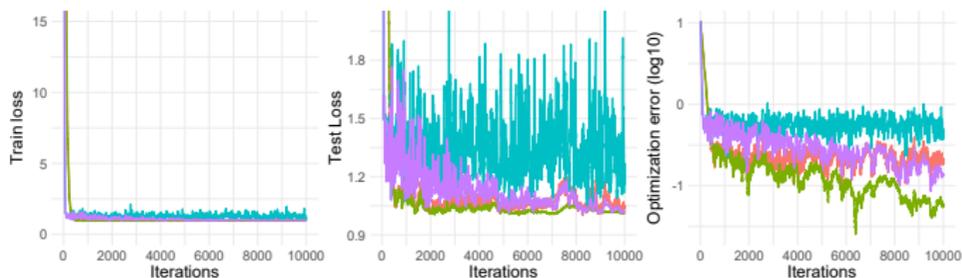
SGD with medium $\alpha = 2 \cdot 10^{-3}$ and indep. features:



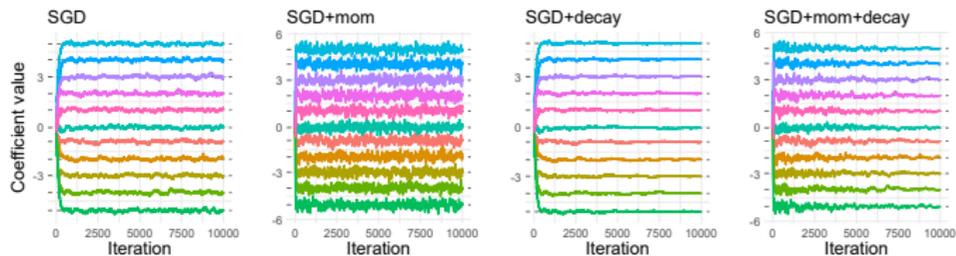
Momentum accelerates optim initially but is eventually outperformed by other variants. **Momentum+decay** is both fast initially and has small final error. **Decay** performs best overall but is slowest initially.

LIN-REG (SGD + LARGE STEP SIZE)

SGD with large $\alpha = 1 \cdot 10^{-2}$ and indep. features:

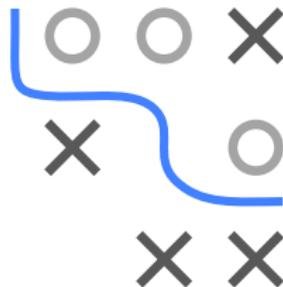


Method — SGD — SGD+decay — SGD+mom — SGD+mom+decay



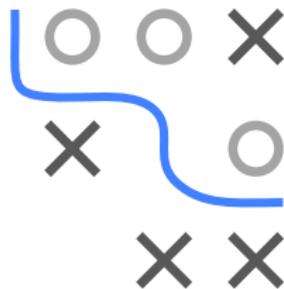
Irreducible error due to additive noise is $\sigma = 1$

High variance in SGD dynamics. **Momentum** becomes unstable while **decay** (necessary to eliminate noise) performs best and is fastest.



DIGRESSION: SOLVING OLS WITH QR DECOMP.

Solving linear least squares via (S)GD is rarely done in practice. But inversion of $\mathbf{X}^T \mathbf{X}$ is numerically unstable and to be avoided. A standard numerical approach applies **QR Decomposition**:



- Factorize $\mathbf{X} \in \mathbb{R}^{n \times p}$ as $\mathbf{X} = \mathbf{QR}$, where $\mathbf{Q} \in \mathbb{R}^{n \times p}$ (thin form) so that $\mathbf{Q}^T \mathbf{Q} = \mathbf{I}$ and $\mathbf{R} \in \mathbb{R}^{p \times p}$ is upper triangular.

- Purpose: replace solving $\hat{\theta} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$ with a more numerically stable method by avoiding direct inversion.

- The QR decomposition can be computed using Gram-Schmidt orthogonalization or Householder transformations

rewrite
 $(\mathbf{X}^T \mathbf{X}) \theta = \mathbf{X}^T \mathbf{y}$
Lin. eq. system,
because computing inverse
of $\mathbf{X}^T \mathbf{X}$ is awful

Steps:

- Decompose \mathbf{X} into \mathbf{Q} and \mathbf{R} .
- Compute $\mathbf{Q}^T \mathbf{y}$.
- Solve triangular system $\mathbf{R} \hat{\theta} = \mathbf{Q}^T \mathbf{y}$ via **back substitution**.

*- maybe we want
to avoid $\mathbf{X}^T \mathbf{X}$ mat.
multip.*

$$\mathbf{X} = \mathbf{Q} \mathbf{R} \quad \mathbf{Q} \text{-orth. to.}$$

$$\mathbf{X}^T \mathbf{X} = \mathbf{R}^T \mathbf{Q}^T \mathbf{Q} \mathbf{R}$$

To do $\hat{\theta}$

DIGRESSION: SOLVING OLS WITH QR DECOMP.

Why this system? Remember normal equation for least squares problem is $\mathbf{X}^\top \mathbf{X} \hat{\boldsymbol{\theta}} = \mathbf{X}^\top \mathbf{y}$. Now replace $\mathbf{X} = \mathbf{QR}$:

$$\begin{aligned}\mathbf{X}^\top \mathbf{X} \hat{\boldsymbol{\theta}} = \mathbf{X}^\top \mathbf{y} &\iff \mathbf{R}^\top (\mathbf{Q}^\top \mathbf{Q}) \mathbf{R}^\top \hat{\boldsymbol{\theta}} = \mathbf{R}^\top \mathbf{Q}^\top \mathbf{y} \\ \mathbf{R}^\top \mathbf{R} \hat{\boldsymbol{\theta}} = \mathbf{R}^\top \mathbf{Q}^\top \mathbf{y} &\iff \mathbf{R} \hat{\boldsymbol{\theta}} = \mathbf{Q}^\top \mathbf{y}\end{aligned}$$

be careful \mathbf{R}^\top isn't square

$\mathbf{R} \hat{\boldsymbol{\theta}} = \mathbf{Q}^\top \mathbf{y}$ is a triangular system easily solvable by back substitution:

$$\mathbf{R} = \begin{bmatrix} r_{11} & r_{12} & \dots & r_{1p} \\ 0 & r_{22} & \dots & r_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & r_{pp} \end{bmatrix}, \quad \hat{\boldsymbol{\theta}} = \begin{bmatrix} \hat{\theta}_1 \\ \hat{\theta}_2 \\ \vdots \\ \hat{\theta}_p \end{bmatrix}, \quad \mathbf{Q}^\top \mathbf{y} = \begin{bmatrix} b_1 \\ b_2 \\ \vdots \\ b_p \end{bmatrix}.$$



DIGRESSION: SOLVING OLS WITH QR DECOMP.

Steps for Back Substitution:

- 1 Start with the last equation (1 unknown):

$$\hat{\theta}_p = \frac{b_p}{r_{pp}}.$$

- 2 Move upwards to the $(p - 1)$ -th equation:

$$\hat{\theta}_{p-1} = \frac{b_{p-1} - r_{p-1,p}\hat{\theta}_p}{r_{p-1,p-1}}.$$

- 3 Continue this process up to the first row:

$$\hat{\theta}_i = \frac{b_i - \sum_{j=i+1}^p r_{ij}\hat{\theta}_j}{r_{ii}} \quad \text{for each } i = p - 1, \dots, 1.$$

Back substitution leverages triangular structure of \mathbf{R} , moving upward from the last row and inserting already known values of $\hat{\theta}$ as we go.



(S)GD FOR LOGISTIC REGRESSION

Comparison of (S)GD, (S)GD + momentum, and (S)GD + momentum + step size control on simulated data:

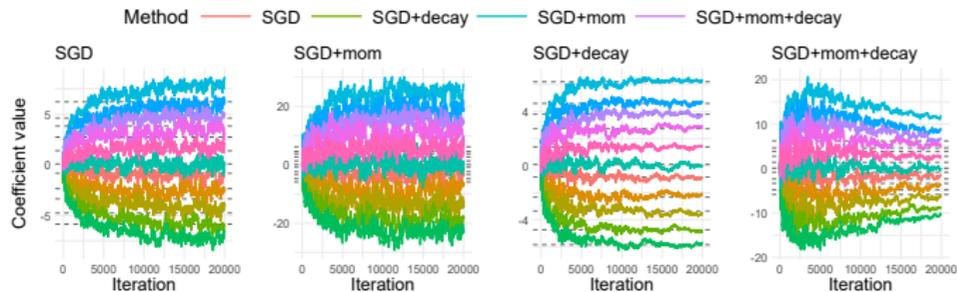
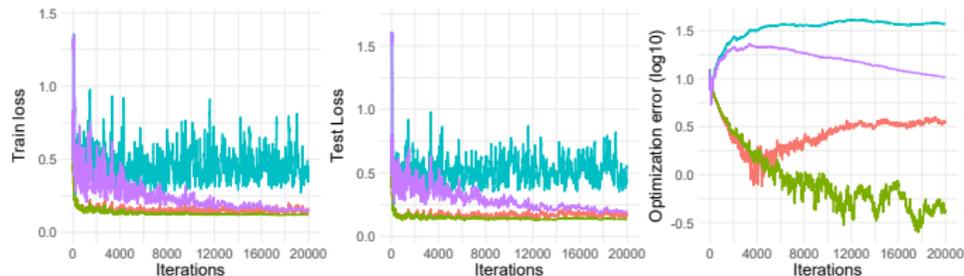
- Logistic regression (log loss) with same simulation setup as for linear regression
- To simulate response, we set $y^{(i)} \sim \mathcal{B}(\pi^{(i)})$, $\pi^{(i)} = \frac{1}{1 + e^{-(\mathbf{x}^{(i)})^\top \boldsymbol{\theta}^*}}$
- We set the momentum parameter to 0.8 and decay = 0.1
- We again use different step sizes to illustrate benefit of momentum and decay depends on step size
- ERM has unique global minimizer but no closed-form solution. We can approximate $\hat{\boldsymbol{\theta}}$ using the glm solution (second-order optim)
- We also track the optimization error $\|\boldsymbol{\theta} - \hat{\boldsymbol{\theta}}\|_2$

~ get from GLM (which uses Newton-Raphson)

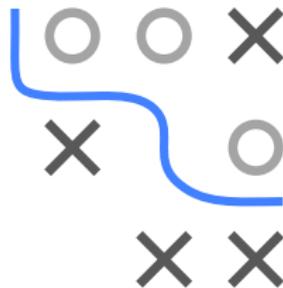


LOG-REG (SGD + LARGE STEP SIZE)

SGD with large $\alpha = 0.3$ and indep. features:



High variance in optim dynamics. Plain SGD and **Momentum** become unstable and oscillate at suboptimal loss values/diverge while **decay** (necessary to eliminate oscillations) performs best and is fastest. **Momentum+decay** initially diverges but recovers once step size reduces (would converge with many more steps).



u .

CLASSIFICATION ON MNIST (SGD)

For a more realistic application, we compare optimizers on training an NN on subset of MNIST image classification task.

- DNN has two hidden layers with 128 and 64 ReLU-activated units and Kaiming normal init. The loss is cross-entropy
Sample around 2040 from Normal
- Instead of SGD we use mini-batch SGD with 100 images per batch

- For the step size schedule, we use cosine decay

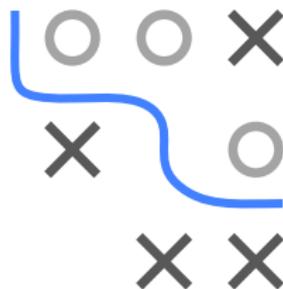
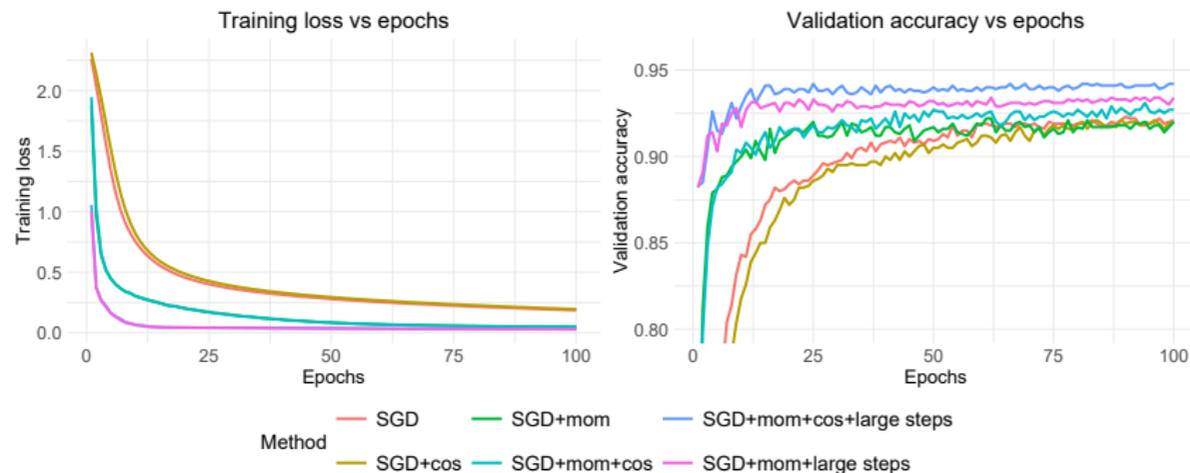
$$\alpha^{[t]} = \alpha^{[0]} \left[(1 - r_{\min}) \cdot \frac{1}{2} \left(1 + \cos \left(\pi \cdot \frac{t}{t_{\max}} \right) \right) + r_{\min} \right] \text{ with final step size fraction } r_{\min} = 0.01$$

- In case of momentum we set the parameter to 0.8
- Regular initial step size is 0.01 and to 0.1 for large step size



CLASSIFICATION ON MNIST (SGD)

Mini-batch SGD with $\alpha \in \{0.01, 0.1\}$ for 100 epochs (5000 iterations):



Observations (NB: green/cyan + blue/purple train losses overlap but not val. acc.):

1. **Momentum** drastically speeds up optimization in all settings.
2. **Plain SGD**: step size decay slows optimization slightly.
3. SGD, SGD+cos and SGD+mom achieve **same** val. acc.
 \Rightarrow no generalization benefit for medium α
4. **Cosine decay+momentum** improves generalization for medium α slightly.
5. **Large step size** with momentum and decay performs best

