

535520: Optimization Algorithms

Lecture 14: Optimizers for Neural Networks

Ping-Chun Hsieh (謝秉均)

December 18, 2024

This Lecture

1. AdaGrad

2. RMSProp, Adam, and Beyond

- Reading Material:

- Duchi and Hazan, “Adaptive subgradient methods for online learning and stochastic optimization,” COLT 2011.
- Kings and Ba, “Adam: A Method for Stochastic Optimization,” ICLR 2015.
- Part of the material is adapted from Prof. Gabriele Farina’s lecture notes (https://www.mit.edu/~gfarina/2024/67220s24_L13_adagrad/L13.pdf)

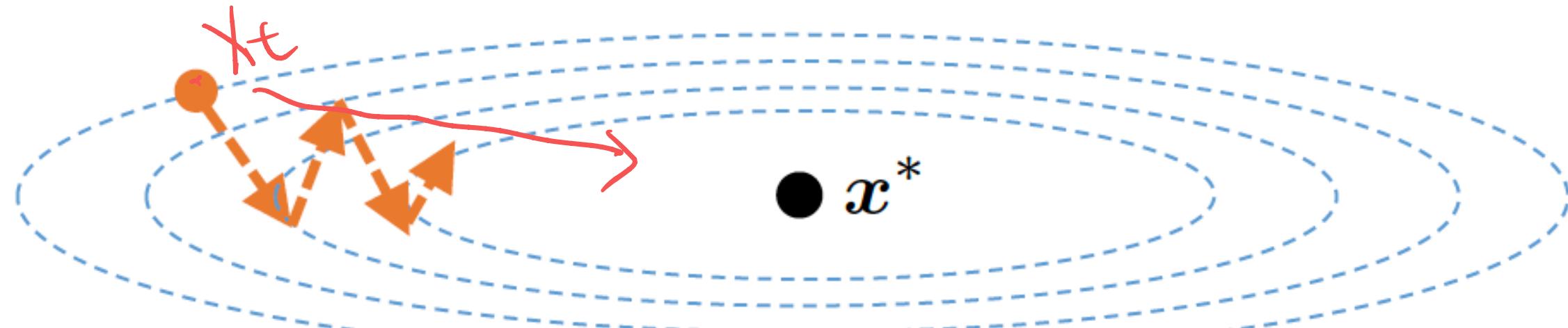
A Whole Family of Optimizers



- A lot of variants of SGD for optimization in practice
 - Learning rate
 - Momentum
- Today: We focus on the most widely used variants
 - AdaGrad
 - RMSProp
 - Adam
 - NAdam

Recall: Non-Homogeneity and “Scaled” Gradients

$$\text{minimize}_{x \in \mathbb{R}^2} \quad f(x) := \frac{1}{2}(x - x^*)^\top Q(x - x^*)$$



- Suppose $Q = [Q_{11}, 0; 0, Q_{22}]$ is a diagonal matrix with $Q_{11} \gg Q_{22}$

- Idea: Accelerate GD by “scaling the gradient”

$$x_{t+1} = x_t - \eta_t Q^{-1} \nabla f(x_t) = x_t - \eta_t (x_t - x^*)$$

$\tilde{Q}^{-1} = \begin{bmatrix} \frac{1}{Q_{11}} & 0 \\ 0 & \frac{1}{Q_{22}} \end{bmatrix}$

- Another interpretation: Scale the learning rate of each input variable separately
(in other words, $\eta_t Q^{-1}$ shall be jointly chosen to tackle non-homogeneity)
- Question: In practice, how to jointly choose $\eta_t Q^{-1}$?

AdaGrad: Adaptive Gradients

Adaptive Subgradient Methods for Online Learning and Stochastic Optimization*

John Duchi

*Computer Science Division
University of California, Berkeley
Berkeley, CA 94720 USA*

JDUCHI@CS.BERKELEY.EDU

Elad Hazan

*Technion - Israel Institute of Technology
Technion City
Haifa, 32000, Israel*

EHAZAN@IE.TECHNION.AC.IL

Yoram Singer

*Google
1600 Amphitheatre Parkway
Mountain View, CA 94043 USA*

SINGER@GOOGLE.COM

(COLT 2011)



John Duchi

(Professor @ Stanford)

Elad Hazan

(Professor @ Princeton)

This work has been cited for more
than 14,000 times

AdaGrad: Adaptive Gradients (Formally)

$t=2$

$$\nabla f(x_0) \quad \nabla f(x_1) \quad \nabla f(x_2)$$

$$x_t, s_t \in \mathbb{R}^d$$

- **Main idea:** Scale the learning rate of each input variable (i.e., per-dimension learning rate) based on the sum of squared gradients accumulated over time

AdaGrad

$$x_{t+1} = x_t - \eta M_t^{-1} \nabla f(x_t), \quad \text{where } M_t = \text{diag} \left([s_t]_i := \sqrt{\sum_{\tau=0}^t [\nabla f(x_\tau)]_i^2} \right)$$

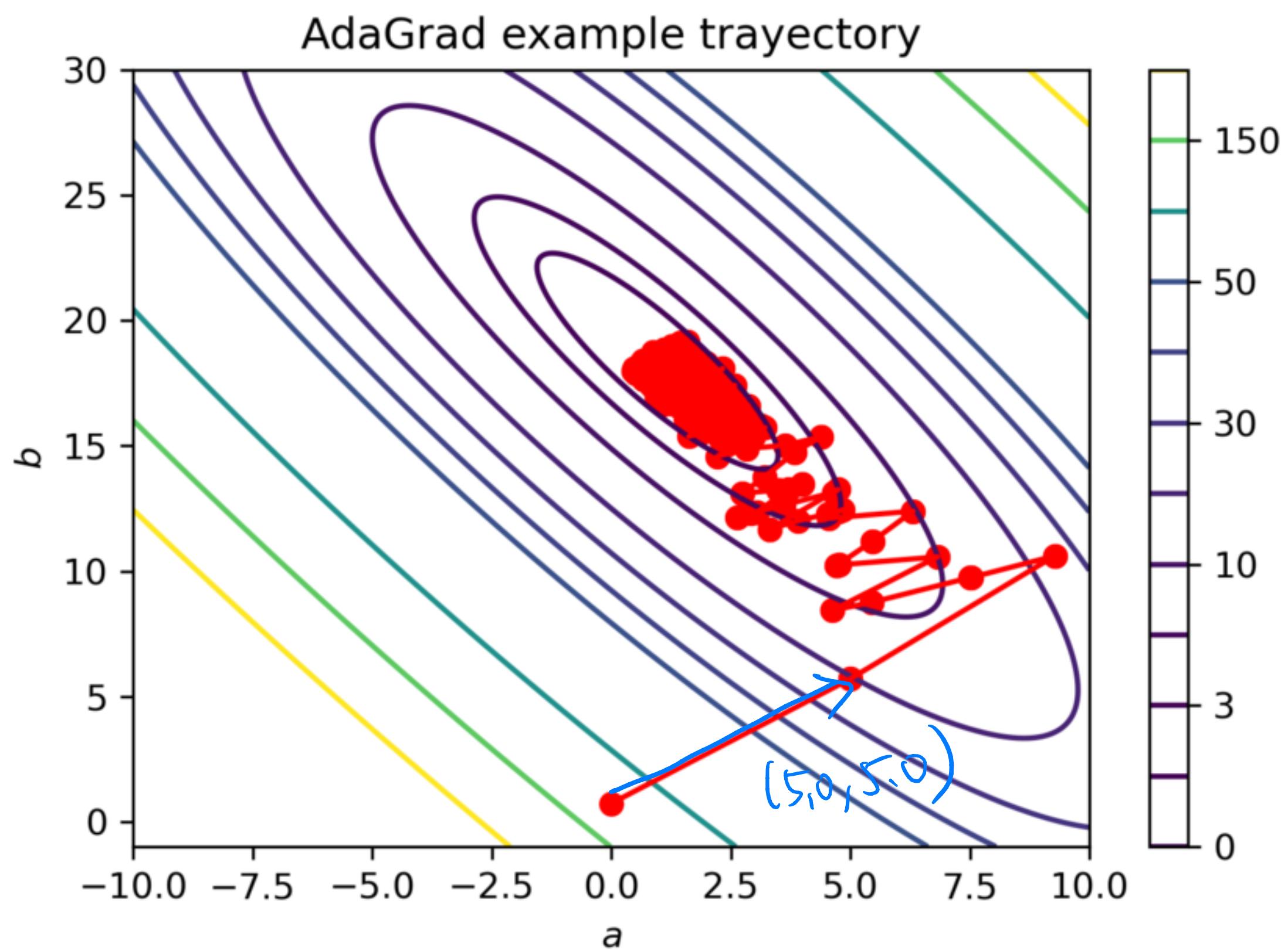
$x \in \mathbb{R}^3$

$M_t = \begin{bmatrix} \bullet & \bullet & \bullet \\ \bullet & \bullet & \bullet \\ \bullet & \bullet & \bullet \end{bmatrix}$

$i\text{-th component of the gradient at time } \tau$

- **Remark:** We assume $[\nabla f(x_0)]_i \neq 0$ for all i , and hence M_t is invertible
- **Remark:** In the stochastic setting, we can replace the true gradient $\nabla f(x_\tau)$ with the stochastic gradient $\hat{\nabla} f(x_\tau)$
- **Remark:** Clearly $[s_t]_i$ is monotonically increasing w.r.t. iteration t

AdaGrad: Numerical Example



Training dataset: Randomly generated data samples as $z = 20y + \varepsilon$

Learned predictor: $z = a + by$

Loss function: $f(a, b) = \frac{1}{n} \sum_{i=1}^n [(a + by_i) - z_i]^2$

Initial parameters: $x_0 \equiv (a, b) = (0, 0)$

- If the first sample is $(y_1, z_1) = (0.39, 9.84)$, what is x_1 under AdaGrad?

$$x_{t+1} = x_t - \eta M_t^{-1} \nabla f(x_t)$$

$$M_t = \text{diag} \left([s_t]_i := \sqrt{\sum_{\tau=0}^t [\nabla f(x_\tau)]_i^2} \right)$$

Intuition Behind AdaGrad

- Imagine a high-dimensional problem of $d \gg 1$
- Each gradient update can be “sparse” (only a few entries dominate the magnitude)
- Question: Is it good to use only one learning rate for all input variables?

$$\nabla f(x_0) = \begin{bmatrix} \textcircled{10} \\ \textcircled{10} \\ \vdots \\ \textcircled{10} \end{bmatrix} \quad 10^6$$

AdaGrad uses separate counters for different variables

1. If a variable $x^{(i)}$ has been updated many times, then the learning rate shall be made smaller

2. To “count” the number of updates, AdaGrad considers $[s_t]_i := \sqrt{\sum_{\tau=0}^t [\nabla f(x_\tau)]_i^2}$

AdaGrad is a Mirror Descent Algorithm

AdaGrad:

$$x_{t+1} = x_t - \eta \cdot M_t^{-1} \nabla f(x_t)$$

Theorem: AdaGrad's update rule is equivalent to the mirror descent update

with Bregman divergence $D_\phi(y||x) = \frac{1}{2}(x - y)^\top M_t(x - y)$, where $M_t := \text{diag}(s_t)$

Proof: Recall the “proximal viewpoint” of MD

$$x_{t+1} = \arg \min_x \left\{ \eta \nabla f(x_t)^\top x + D_\phi(x||x_t) \right\} = \arg \min_x \left\{ \eta \nabla f(x_t)^\top x + \frac{1}{2}(x - x_t)^\top M_t(x - x_t) \right\}$$

By setting the gradient of the above to 0 and solving for x :

$$\underbrace{\eta \nabla f(x_t) + M_t(x - x_t)}_{} = 0 \Leftrightarrow \underbrace{x = x_t - \eta \cdot M_t^{-1} \nabla f(x_t)}_{}$$

Remark: AdaGrad captures local geometry via M_t

Review: Mirror Descent (MD)

Key Idea: Generalize the L_2 proximal term to other distance measures

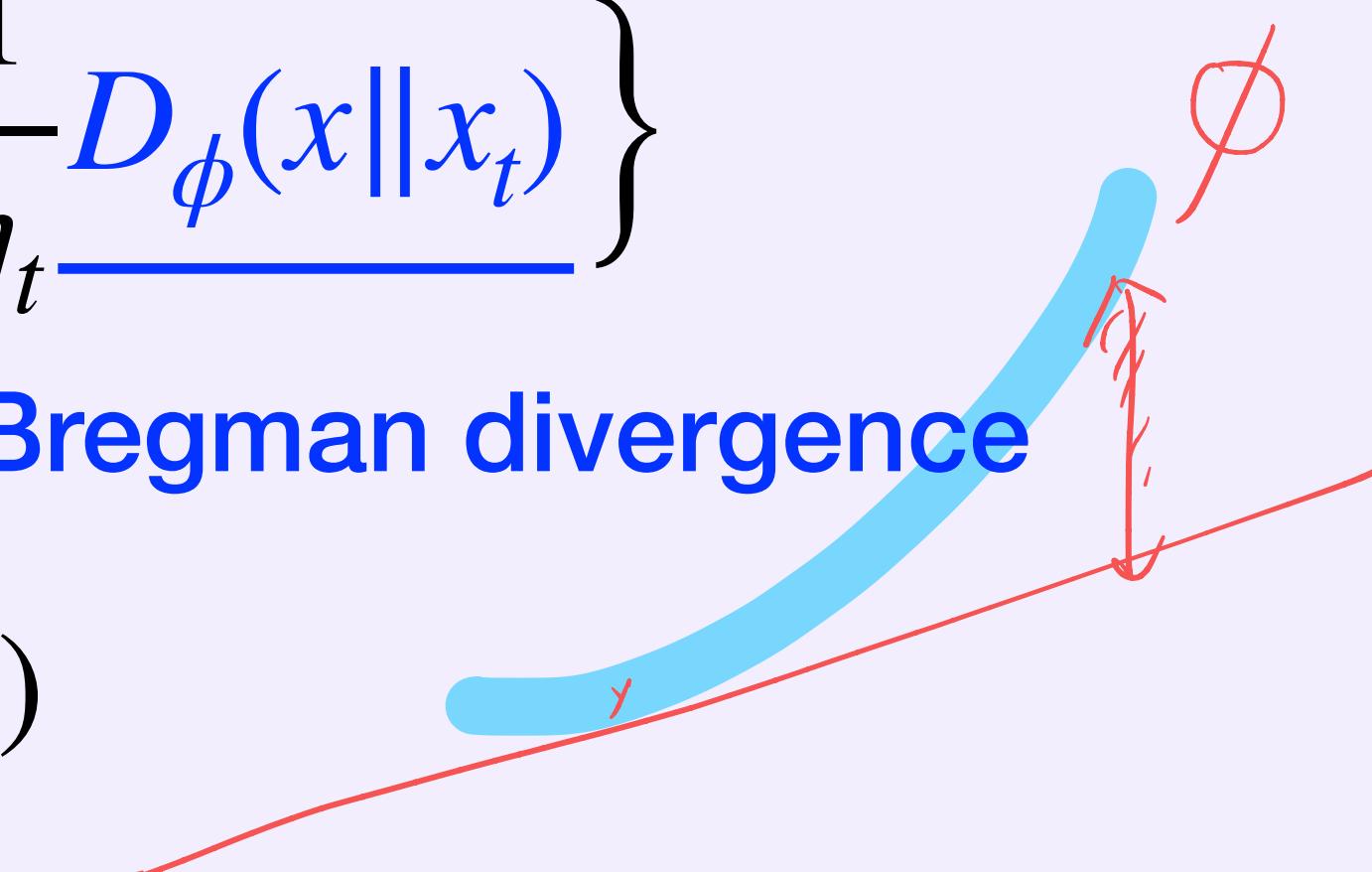
$$x_{t+1} = \arg \min_{x \in C} \left\{ f(x_t) + \nabla f(x_t)^\top (x - x_t) + \frac{1}{\eta_t} D_\phi(x \| x_t) \right\}$$

first-order approximation

Bregman divergence

$$\text{where } D_\phi(y \| x) := \phi(y) - \phi(x) - \nabla \phi(x)^\top (y - x)$$

($\phi(\cdot)$ is a strictly convex and differentiable “distance-generating” function)



- **Remark:** Bregman divergence is meant to capture "local geometry" of objective function

Review: Bregman Divergence as Mahalanobis Distance

Suppose $\phi(x) = x^\top Ax$ (where A is a positive definite $d \times d$ matrix)

- ▶ The resulting Bregman divergence is $D_\phi(y\|x) = (x - y)^\top A(x - y)$

(aka “Mahalanobis distance”)

$$D_\phi(y\|x) := \phi(y) - (\phi(x) + \nabla\phi(x)^\top(y - x))$$

$$= y^\top Ay - (x^\top Ax + 2(Ax)^\top(y - x))$$

$$= y^\top Ay + x^\top Ax - 2x^\top Ay$$

=

Performance Guarantee of AdaGrad

Theorem (Duchi et al., 2011): Let f be a convex and differentiable function.

For any choice of coefficients $\lambda_i \geq 0$, AdaGrad with step size $\eta = D/\sqrt{2}$ satisfies

$$\frac{1}{T} \sum_{t=0}^{T-1} f(x_t) \leq f(x_*) + \frac{\sqrt{2d}D}{T}$$

$$\min_{\lambda \in \mathbb{R}_{\geq 0}^n, \|\lambda\|_1=n} \sum_{t=0}^{T-1} \nabla f(x_t)^\top \text{diag}(\lambda)^{-1} \nabla f(x_t), \quad = O(T)$$

where $D := \max_{t \in [0, T]} \|x_t - x^*\|_\infty$ is the maximum distance from the optimizer

Moreover, with bounded gradients, AdaGrad achieves a regret bound of

$$R(T) := \sum_{t=1}^T (f(x_t) - f(x^*)) = O(\sqrt{T})$$

Cumulative regret

Regret

Recall: Mirror Descent Lemma

Lemma: Let f be a convex function. Under Mirror descent update, we have

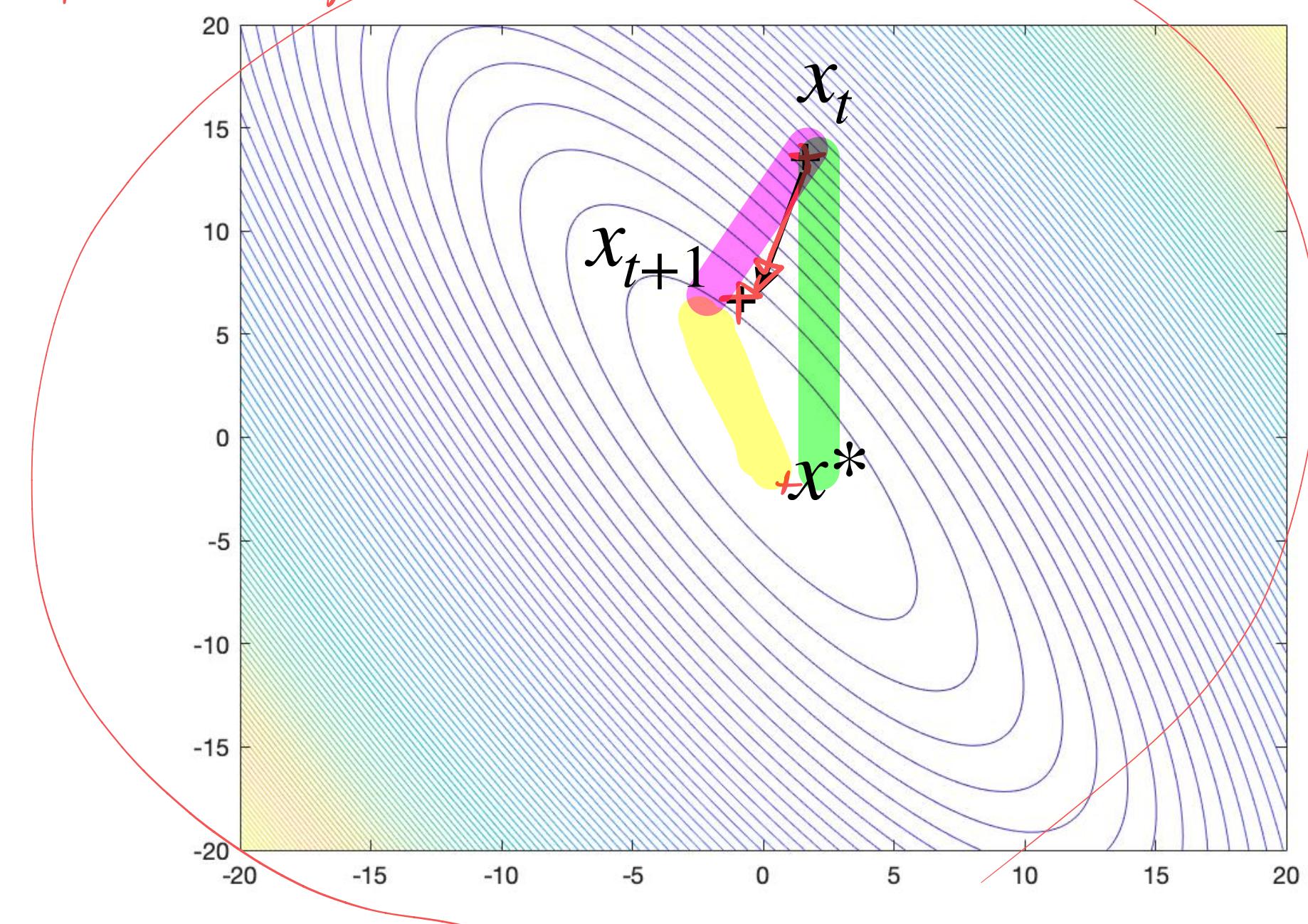
$$f(x_t) \leq f(x^*) + \frac{1}{\eta} \left(D_\phi(x^* \| x_t) - D_\phi(x^* \| x_{t+1}) + D_\phi(x_{t+1} \| x_t) \right).$$

$$\approx \|x^* - x_t\|^2$$

$$\approx \|x^* - x_{t+1}\|^2$$

$$\approx \|x_{t+1} - x_t\|^2$$

- What's the intuition?



Convergence Analysis for AdaGrad

$$f(x_t) \leq f(x^*) + \frac{1}{\eta} \left(D_{\phi_t}(x^* \| x_t) - D_{\phi_t}(x^* \| x_{t+1}) + D_{\phi_t}(x_{t+1} \| x_t) \right).$$

By taking the summation of the above over $t = 0, 1, \dots, T - 1$, we have

$$\sum_{t=0}^{T-1} f(x_t) \leq T \cdot f(x^*) + \frac{1}{\eta} \left(\sum_{t=0}^{T-2} \left(D_{\phi_{t+1}}(x^* \| x_{t+1}) - D_{\phi_t}(x^* \| x_{t+1}) \right) + \sum_{t=0}^{T-1} D_{\phi_t}(x_{t+1} \| x_t) \right).$$

(A)

(B)

Bounding the "Almost-telescopic" term (A)

Lemma: Let T be arbitrary and assume that $D := \max_{0 \leq t \leq T} \|x_t - x^*\|_\infty$ is finite.

For any $T \in \mathbb{N}$, the term (A) satisfies the inequality

$$\sum_{t=0}^{T-2} \left(D_{\phi_{t+1}}(x^* \| x_{t+1}) - D_{\phi_t}(x^* \| x_{t+1}) \right) \leq \frac{D^2}{2} \|s_{T-1}\|_1$$

Proof: Step 1: It is easy to verify that

$$\begin{aligned} D_{\phi_{t+1}}(x^* \| x_{t+1}) - D_{\phi_t}(x^* \| x_{t+1}) &= \frac{1}{2} \underbrace{(x_{t+1} - x^*)^T}_{\text{by Hölder's inequality}} \underbrace{(M_{t+1} - M_t)}_{\text{by the definition of } D} \underbrace{(x_{t+1} - x^*)}_{\text{by the definition of } D} \\ &\leq \frac{1}{2} \|x_{t+1} - x^*\|_\infty^2 \|s_{t+1} - s_t\|_1 \\ &\leq \frac{1}{2} D^2 \|s_{t+1} - s_t\|_1 \end{aligned}$$

(Cont.).

Step 2: By the fact that $\underline{[S_{t+1}]_i \geq [S_t]_i \geq 0}$ for all $i \in \{1, \dots, d\}$,

then we have

$$\|S_{t+1} - S_t\|_1 = \|S_{t+1}\|_1 - \|S_t\|_1,$$

This implies that

$$\begin{aligned} \sum_{t=0}^{T-2} \left(D_{\phi_{t+1}}(x^* \| x_{t+1}) - D_{\phi_t}(x^* \| x_{t+1}) \right) &\leq \frac{D^2}{2} \sum_{t=0}^{T-2} (\|S_{t+1}\|_1 - \|S_t\|_1) \\ &= \frac{D^2}{2} (\|S_{T-1}\|_1 - \|S_0\|_1) \\ &\leq \frac{D^2}{2} \|S_{T-1}\|_1 \end{aligned}$$

Bounding the Summation (B)

Lemma: For any $T \in \mathbb{N}$, we have

$$\sum_{t=0}^{T-1} D_{\phi_t}(x_t \| x_{t+1}) \leq \eta^2 \| s_{T-1} \|_1$$

Pf: Recall that $D_{\phi_t}(x_t \| x_{t+1}) = \frac{1}{2} (x_{t+1} - x_t)^T M_t (x_{t+1} - x_t)$

$x_{t+1} = x_t - \eta \cdot M_t^{-1} \nabla f(x_t)$

$$= \frac{\eta^2}{2} \nabla f(x_t)^T M_t^{-1} \nabla f(x_t) \quad \dots \text{ (by MD update.)}$$

By the definition of $M_t := \text{diag}(s_t)$ and $[s_t]_i = \sqrt{\sum_{t=0}^T [\nabla f(x_t)]_i^2}$

$$\begin{aligned} \nabla f(x_t)^T M_t^{-1} \nabla f(x_t) &= \sum_{i=1}^d \frac{[\nabla f(x_t)]_i^2}{[s_t]_i} \\ &= \sum_{i=1}^d \frac{[s_t]_i^2 - [s_{t-1}]_i^2}{[s_t]_i} \leq \sum_{i=1}^d \frac{[s_t]_i - [s_{t-1}]_i}{[s_t]_i} \end{aligned}$$

By taking the summation over $t=0, \dots, T-1$,
we get the desired result.

$$= 2(\|s_t\|_1 - \|s_{t-1}\|_1)$$

Bounding the norm $\| S_{T-1} \|_1$

Lemma: $\| S_{T-1} \|_1 \leq \sqrt{n} \cdot \sqrt{\min_{\lambda \in \mathbb{R}^d, \|\lambda\|_1=d} \sum_{t=0}^{T-1} \nabla f(x_t)^T \text{diag}(\lambda)^{-1} \nabla f(x_t)}$

Pf: $\| S_{T-1} \|_1^2 = \left(\sum_{i=1}^d \sqrt{\sum_{t=0}^{T-1} [\nabla f(x_t)]_i^2} \right)^2 \dots (1)$

Pick any $\lambda \in \mathbb{R}^d$
with $\|\lambda\|_1=d$ $= \left(\sum_{i=1}^d \left[\sqrt{\lambda_i} \cdot \left(\frac{1}{\sqrt{\lambda_i}} \sqrt{\sum_{t=0}^{T-1} [\nabla f(x_t)]_i^2} \right) \right]^2 \right)^2 \dots (2)$

$\leq \left(\sum_{i=1}^d (\sqrt{\lambda_i})^2 \right) \cdot \left(\sum_{i=1}^d \left(\frac{1}{\sqrt{\lambda_i}} \cdot \sqrt{\sum_{t=0}^{T-1} [\nabla f(x_t)]_i^2} \right)^2 \right) \dots (3)$

$= n \cdot \left(\sum_{t=0}^{T-1} \nabla f(x_t)^T \text{diag}(\lambda)^{-1} \nabla f(x_t) \right).$

Pros and Cons of AdaGrad

Pros:

1. No manual tuning of learning rate
2. Memory-efficient
3. Provable regret guarantee
4. Originally designed for “sparse gradient” regime (and hence good for this setting)

Cons:

1. AdaGrad can decrease the learning rate too fast in practice (why?)
2. Sensitive to the initialization (e.g., if $\nabla f(x_0)$ is large, then the learning rate can be get too small very quickly at later iterations)

RMSProp

- RMSProp = Root Mean Squared Propagation
- Main idea: RMSProp uses a less aggressive learning rate decay by considering the Root Mean Square of recent (stochastic) gradients

RMSProp

$$x_{t+1} = x_t - \eta M_t^{-1} \nabla f(x_t),$$

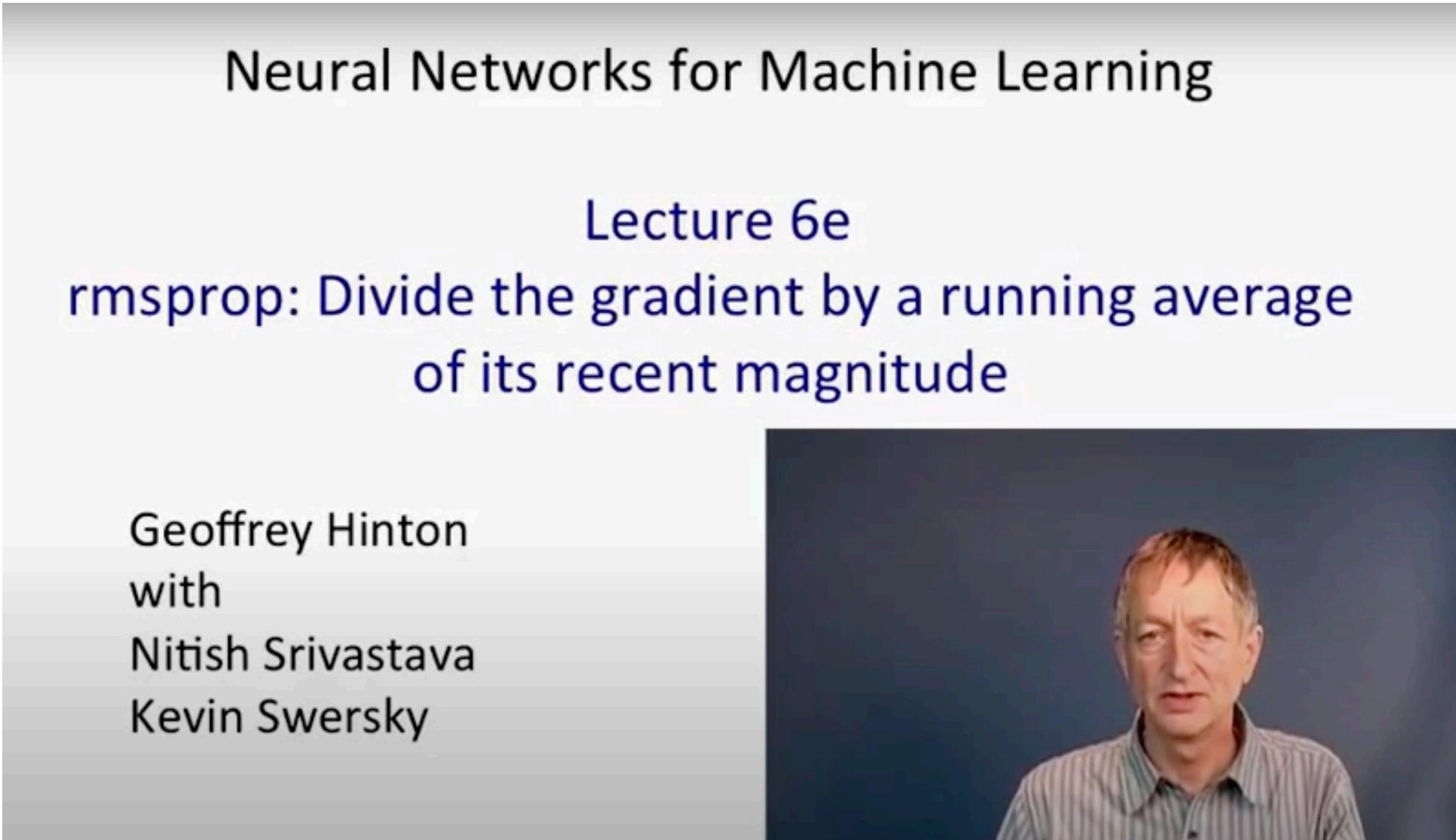
$$\sum_i = \frac{(1-\beta) \cdot (1-\beta^{t+1})}{1-\beta}$$
$$(1-\beta)[\nabla f(x_t)]_i^2 + (1-\beta)\cdot\beta \cdot [\nabla f(x_{t-1})]_i^2 + \dots + (1-\beta)\cdot\beta^t \cdot [\nabla f(x_0)]_i^2$$

where $[s_t]_i^2 = \beta [s_{t-1}]_i^2 + (1-\beta)[\nabla f(x_t)]_i^2$

$$M_t = \text{diag}([s_t])$$

- Remark: Now $[s_t]_i$ is no longer monotonically increasing

RMSProp



- Interestingly, RMSProp has never been formally published in any conference or journal
- In fact, RMSProp came from a Coursera lecture by Jeff Hinton!

Lecture Video: <https://www.youtube.com/watch?v=defQQqkXEfE>

Jeff Hinton's slides: https://www.cs.toronto.edu/~tijmen/csc321/slides/lecture_slides_lec6.pdf

Adadelta Optimizer

- Adadelta = Adaptive with Unit Correction via Δx_t
- **Main idea:** Adadelta enhances AdaGrad by incorporating two techniques
 1. Use exponential moving averages to estimate M_t (same as RMSProp)
 2. Apply exponential moving averages of Δx_t for unit correction

Adadelta

$$x_{t+1} = x_t - \eta M_t^{-1} \nabla f(x_t),$$
$$\text{where } [s_t]^2_i = \beta[s_{t-1}]_i^2 + (1 - \beta)[\nabla f(x_t)]_i^2$$
$$[\Delta_t]^2_i = \beta[\Delta_{t-1}]_i^2 + (1 - \beta)[(x_t - x_{t-1})]_i^2$$
$$M_t = \text{diag}([s_t / \Delta_t])$$

Adam Optimizer

- Adam = Adaptive Moment Estimation
- Main idea: Adam combines the adaptive learning rate of AdaGrad and RMSProp with the idea of momentum
 - 1. Use exponential moving averages to estimate gradients (aka momentum)
 - 2. Use exponential moving averages to estimate $\text{diag}(s_t)$

Adam

$$\text{where } \mathbf{g}_t = \gamma \mathbf{g}_{t-1} + (1 - \gamma) \nabla f(\mathbf{x}_t)$$

$$\mathbf{x}_{t+1} = \mathbf{x}_t - \eta M_t^{-1} \mathbf{g}_t,$$

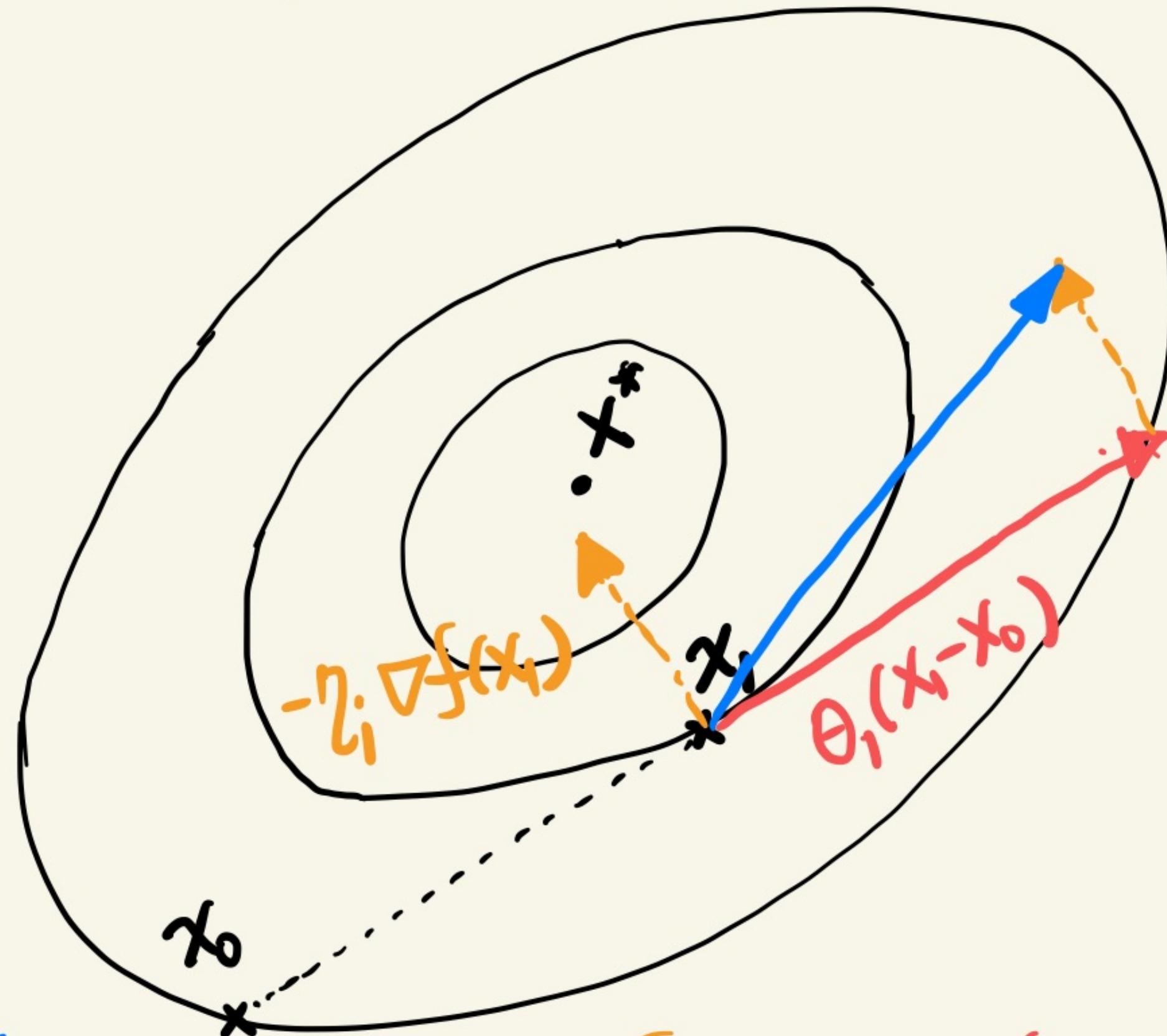
$$[s_t]_i^2 = \beta [s_{t-1}]_i^2 + (1 - \beta) [\nabla f(\mathbf{x}_t)]_i^2$$

$$M_t = \text{diag}([s_t]) + \varepsilon \mathbb{I}_d$$

- Remark: Typically, $\gamma = 0.9$ and $\beta = 0.999$ (e.g., PyTorch)
- Question: What method does Adam reduce to if $\gamma = 0$? RMSprop

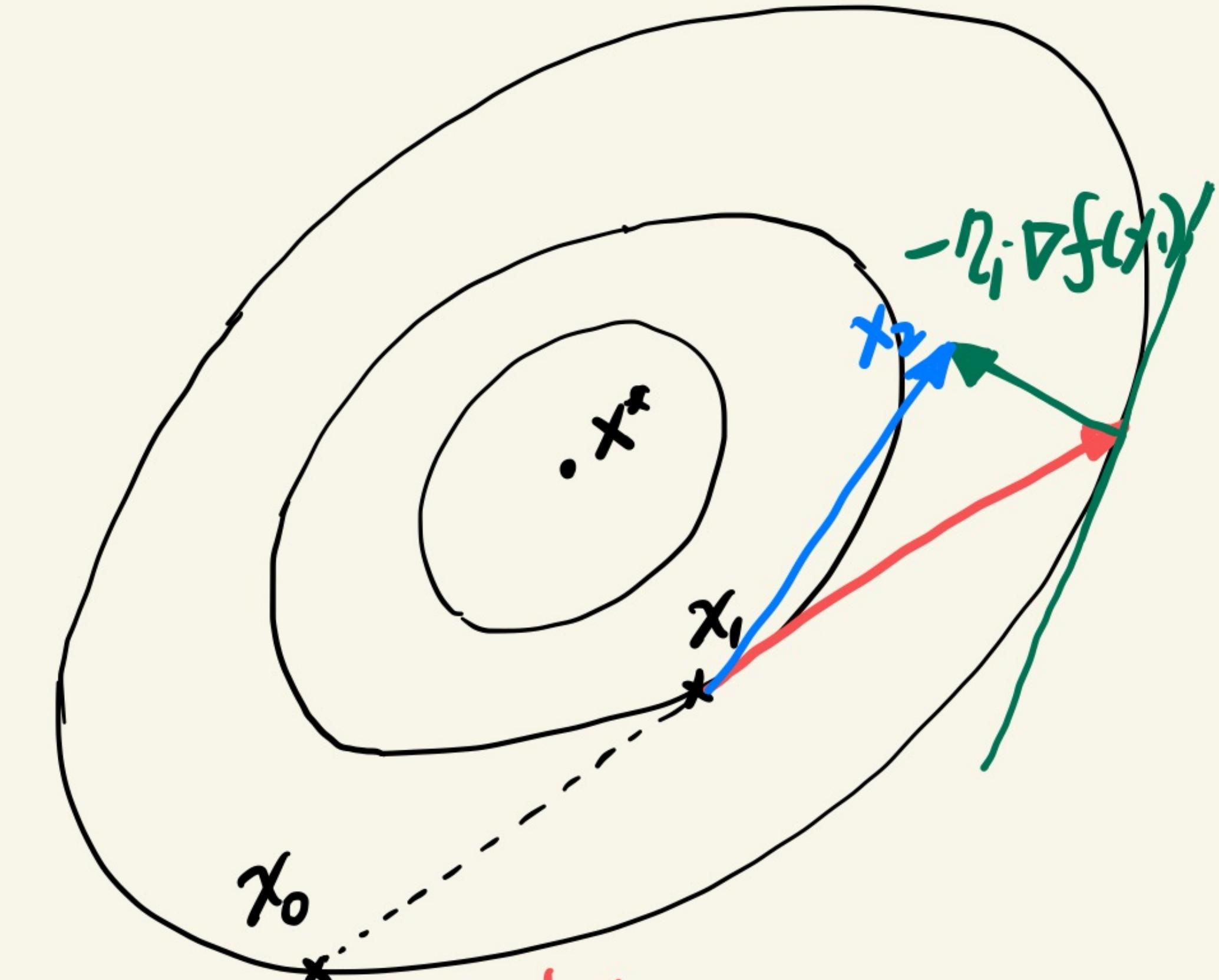
Visualization: Polyak's Momentum vs Nesterov's Method

Polyak's Momentum



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

Nesterov's Method



$$x_{t+1} = x_t + \underbrace{\frac{t-1}{t+2}(x_t - x_{t-1})}_{=: y_t} - \eta_t \nabla f(y_t)$$

Adam Optimizer

ADAM: A METHOD FOR STOCHASTIC OPTIMIZATION

Diederik P. Kingma*
University of Amsterdam, OpenAI
dpkingma@openai.com

Jimmy Lei Ba*
University of Toronto
jimmy@psi.utoronto.ca

ABSTRACT

(ICLR 2015)

We introduce *Adam*, an algorithm for first-order gradient-based optimization of stochastic objective functions, based on adaptive estimates of lower-order moments. The method is straightforward to implement, is computationally efficient, has little memory requirements, is invariant to diagonal rescaling of the gradients, and is well suited for problems that are large in terms of data and/or parameters. The method is also appropriate for non-stationary objectives and problems with very noisy and/or sparse gradients. The hyper-parameters have intuitive interpretations and typically require little tuning. Some connections to related algorithms, on which *Adam* was inspired, are discussed. We also analyze the theoretical convergence properties of the algorithm and provide a regret bound on the convergence rate that is comparable to the best known results under the online convex optimization framework. Empirical results demonstrate that Adam works well in practice and compares favorably to other stochastic optimization methods. Finally, we discuss *AdaMax*, a variant of *Adam* based on the infinity norm.



Diederik Kingma
(Anthropic)



Jimmy Ba
(University of Toronto)

This work has been cited for more than 200,000 times

Issue #1: Adam Optimizer with Initialization Bias Correction

- Adam usually initializes $g_0 = 0$ and $s_0 = 0$
- This can cause a little bias due to the exponential moving average

Adam with bias correction

where $g_t = \gamma g_{t-1} + (1 - \gamma) \nabla f(x_t)$

$$x_{t+1} = x_t - \eta M_t^{-1} \mathbf{g}_t,$$

$$[s_t]_i^2 = \beta [s_{t-1}]_i^2 + (1 - \beta) [\nabla f(x_t)]_i^2$$

$$\begin{aligned} g_t &\leftarrow \mathbf{g}_t / (1 - \gamma^t) \\ s_t &\leftarrow s_t / (1 - \beta^t) \end{aligned}$$

$$M_t = \text{diag}([s_t])$$

Issue #2: Adam Can Have “Non-Monotonic” Learning Rates

Adam

$$x_{t+1} = x_t - \eta M_t^{-1} g_t,$$

where $g_t = \gamma g_{t-1} + (1 - \gamma) \nabla f(x_t)$

$$[s_t]_i^2 = \beta [s_{t-1}]_i^2 + (1 - \beta) [\nabla f(x_t)]_i^2$$

$$M_t = \text{diag}([s_t])$$

- **Question:** Does Adam always have “non-increasing” learning rates?
(i.e., the learning rate strictly decays?)
- **Question:** How about AdaGrad and RMSProp?



Adam Can Possibly Suffer From Linear Regret

ON THE CONVERGENCE OF ADAM AND BEYOND

Sashank J. Reddi, Satyen Kale & Sanjiv Kumar
Google New York
New York, NY 10011, USA
`{sashank, satyenkale, sanjivk}@google.com`

ABSTRACT

Several recently proposed stochastic optimization methods that have been successfully used in training deep networks such as RMSPROP, ADAM, ADADELTA, NADAM are based on using gradient updates scaled by square roots of exponential moving averages of squared past gradients. In many applications, e.g. learning with large output spaces, it has been empirically observed that these algorithms fail to converge to an optimal solution (or a critical point in nonconvex settings). We show that one cause for such failures is the exponential moving average used in the algorithms. We provide an explicit example of a simple convex optimization setting where ADAM does not converge to the optimal solution, and describe the precise problems with the previous analysis of ADAM algorithm. Our analysis suggests that the convergence issues can be fixed by endowing such algorithms with “long-term memory” of past gradients, and propose new variants of the ADAM algorithm which not only fix the convergence issues but often also lead to improved empirical performance.

- $[s_t]_i^2 = \beta[s_{t-1}]_i^2 + (1 - \beta)[\nabla f(x_t)]_i^2$ is not monotonically increasing with t
- One can construct an online learning problem such that Adam can have a regret $\Theta(\sqrt{T})$

AMSGrad: Adam With Monotonic Learning Rates

Adam with monotonicity

where $g_t = \gamma g_{t-1} + (1 - \gamma) \nabla f(x_t)$

$$x_{t+1} = x_t - \eta M_t^{-1} \mathbf{g}_t,$$

$$[s_t]_i^2 = \beta [s_{t-1}]_i^2 + (1 - \beta) [\nabla f(x_t)]_i^2$$

$$[s_t]_i \leftarrow \max\{[s_t]_i, [s_{t-1}]_i\}$$

$$M_t = \text{diag}([s_t])$$

- **Intuition:** Use a simple max operation to ensure monotonicity in $[s_t]_i$

PyTorch Implementation of Adam Optimizer

```
[docs]class Adam(Optimizer):
    def __init__(self,
                 params: ParamsT,
                 lr: Union[float, Tensor] = 1e-3,
                 betas: Tuple[float, float] = (0.9, 0.999),
                 eps: float = 1e-8,
                 weight_decay: float = 0,
                 amsgrad: bool = False,
                 *,
                 foreach: Optional[bool] = None,
                 maximize: bool = False,
                 capturable: bool = False,
                 differentiable: bool = False,
                 fused: Optional[bool] = None,
                 ):
        if isinstance(lr, Tensor):
            if foreach and not capturable:
                raise ValueError(
                    "lr as a Tensor is not supported for capturable=False and
foreach=True")
            if lr.numel() != 1:
                raise ValueError("Tensor lr must be 1-element")
        if not 0.0 <= lr:
            raise ValueError(f"Invalid learning rate: {lr}")
        if not 0.0 <= eps:
            raise ValueError(f"Invalid epsilon value: {eps}")
        if not 0.0 <= betas[0] < 1.0:
            raise ValueError(f"Invalid beta parameter at index 0: {betas[0]}")
        if not 0.0 <= betas[1] < 1.0:
            raise ValueError(f"Invalid beta parameter at index 1: {betas[1]}")
        if not 0.0 <= weight_decay:
            raise ValueError(f"Invalid weight_decay value: {weight_decay}")

        defaults = dict(
            lr=lr,
            betas=betas,
            eps=eps,
            weight_decay=weight_decay,
            amsgrad=amsgrad,
            maximize=maximize,
            foreach=foreach,
            capturable=capturable,
            differentiable=differentiable,
            fused=fused,
        )
```

- A quite involved implementation (~1000 lines of python code)
- Can you identify some practical considerations not mentioned here?
(Hint: “AMSGrad”)

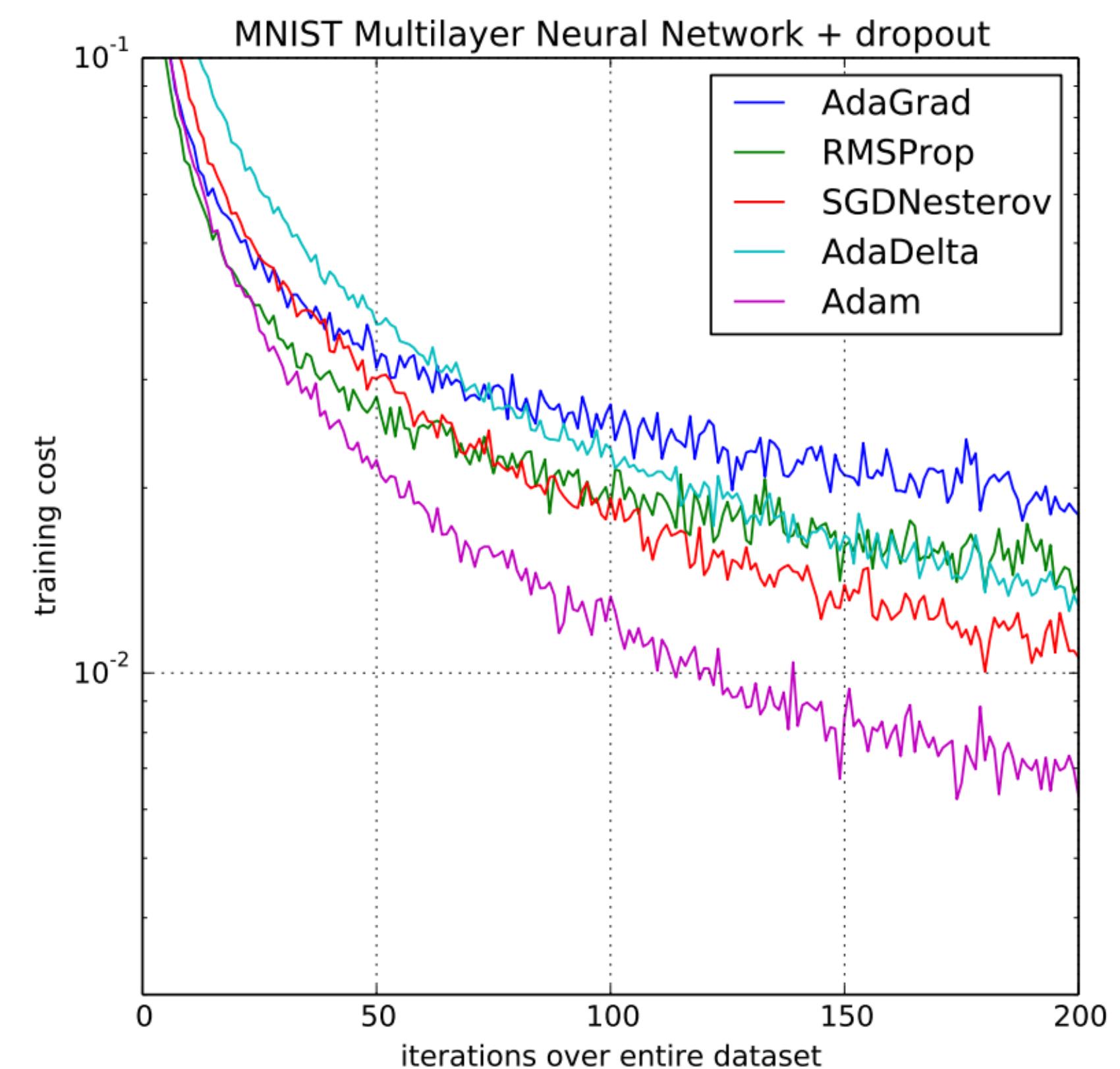
Some Remarks on Adam

1. Adam enjoys a nice $R(T) := \sum_{t=1}^T (f(x_t) - f(x^*)) = O(\sqrt{T})$ regret

bound in the convex setting (similar to AdaGrad)

2. Adam is widely used mainly due to its strong empirical convergence and a fairly low complexity

3. Adam can also be combined with Nesterov momentum (and hence called “**Nadam**”)



A General Form of Adaptive First-Order Methods

Adaptive First-Order Method

$$x_{t+1} = x_t - \eta V_t^{-1/2} m_t,$$

where

$$(1) \ g_t \leftarrow \nabla f(x_t) \text{ or } \hat{\nabla} f(x_t)$$

$$(2) \ m_t = \phi(g_1, \dots, g_t)$$

$$(3) \ V_t = \psi(g_1, \dots, g_t)$$

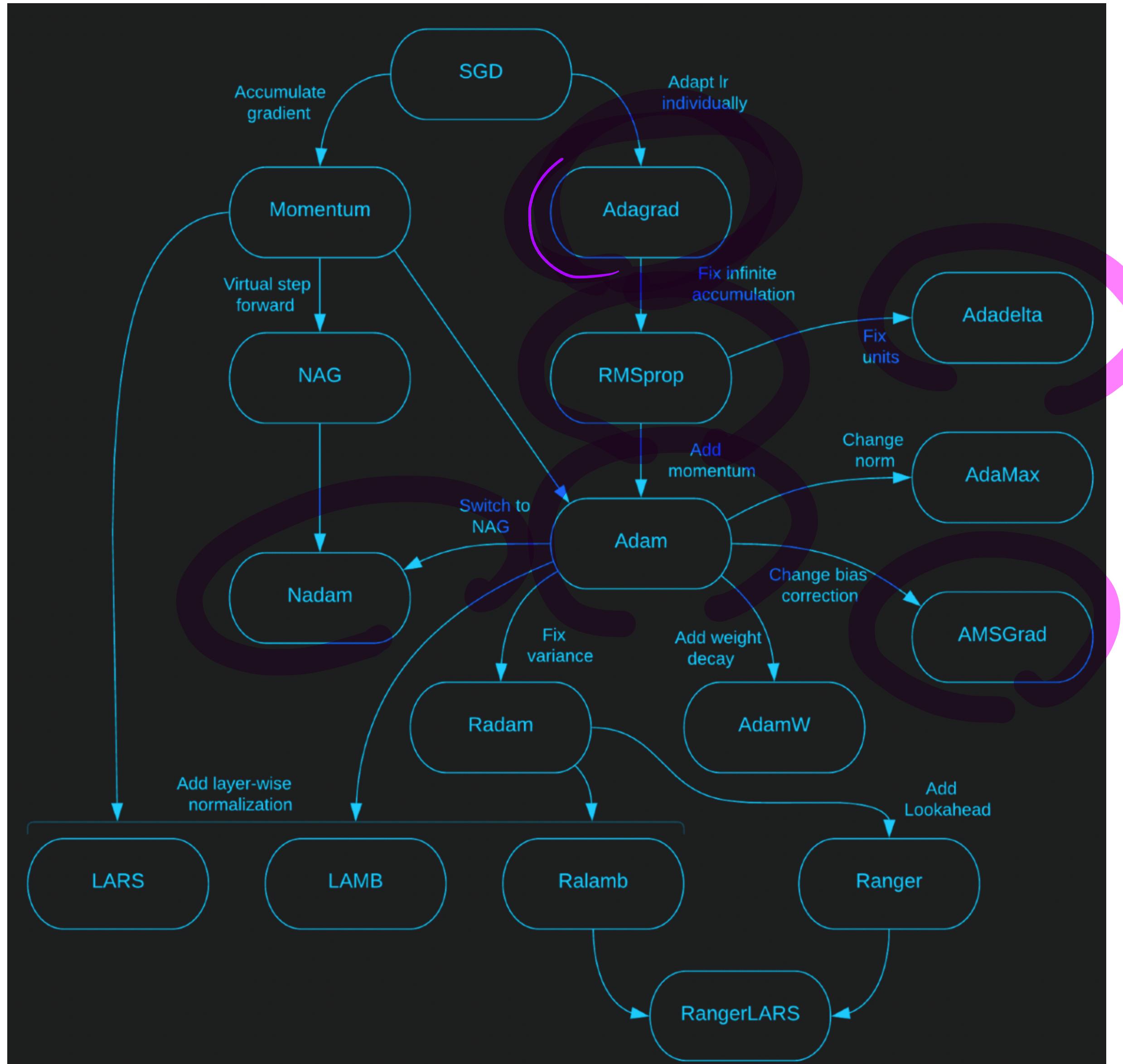
- What are the ϕ, ψ of GD (and SGD)?

- What are the ϕ, ψ of AdaGrad?

$$\text{Diag}\left(\sum_{t=1}^t g_i g_i^T\right)$$

- What are the ϕ, ψ of Adam?

Summary: Adaptive First-Order Optimizers



- AdaGrad (2011)
- RMSProp (2012)
- Adadelta (2012)
- Adam (2015)
- Nadam (2016)
- AMSGrad (2018)