

Handout 3: Stochastic gradient descent

Lecturer & author: Georgios P. Karagiannis

georgios.karagiannis@durham.ac.uk

Aim. To introduce the stochastic gradient descent (motivation, description, practical tricks, analysis in the convex scenario, and implementation).

Reading list & references:

- Shalev-Shwartz, S., & Ben-David, S. (2014). Understanding machine learning: From theory to algorithms. Cambridge university press.
– Ch. 14.3 Stochastic Gradient Descent (SGD), 14.5 Variants, 14.5 Learning with SGD

Further reading

- Bottou, L. (2012). Stochastic gradient descent tricks. In Neural networks: Tricks of the trade (pp. 421-436). Springer, Berlin, Heidelberg.

1. MOTIVATIONS FOR STOCHASTIC GRADIENT DESCENT

Problem 1. Consider a learning problem $(\mathcal{H}, \mathcal{Z}, \ell)$. Learning may involve the computation of the minimizer $w^* \in \mathcal{H}$, where \mathcal{H} is a class of hypotheses, of the risk function (RF) $R(w) = \mathbb{E}_{z \sim g}(\ell(w, z))$ given an unknown data generating model $g(\cdot)$ and using a known tractable loss $\ell(\cdot, \cdot)$; that is

$$(1.1) \quad w^* = \arg \min_{w \in \mathcal{H}} (R_g(w)) = \arg \min_{w \in \mathcal{H}} (\mathbb{E}_{z \sim g}(\ell(w, z)))$$

Remark 2. Gradient descent (GD) cannot be directly utilized to address Problem 1 (i.e., minimize the Risk function) because g is unknown, and because (1.1) involves an integral which may be computationally intractable. Instead it aims to minimize the ERF $\hat{R}(w) = \frac{1}{n} \sum_{i=1}^n \ell(w, z_i)$ which ideally is used as a proxy when data size n is big (big-data).

Remark 3. The implementation of GD may be computationally impractical even in problems where we need to minimize an ERF $\hat{R}_n(w)$ if we have big data ($n \approx \text{big}$). This is because GD requires the recursive computation of the exact gradient $\nabla \hat{R}_n(w) = \frac{1}{n} \sum_{i=1}^n \nabla \ell(w, z_i)$ using all the data $\{z_i\}$ at each iteration. That may be too slow.

Remark 4. Stochastic gradient descent (SGD) aims at solving (1.1), and overcoming the issues in Remarks 2 & 3 by using an unbiased estimator of the actual gradient (or some sub-gradient) based on a sample properly drawn from g .

2. STOCHASTIC GRADIENT DESCENT

2.1. Description.

Notation 5. For the sake of notation simplicity and generalization, we present Stochastic Gradient Descent (SGD) in the following minimization problem

$$(2.1) \quad w^* = \arg \min_{w \in \mathcal{H}} (f(w))$$

where here $f : \mathbb{R}^d \rightarrow \mathbb{R}$, and $w \in \mathcal{H} \subseteq \mathbb{R}^d$; $f(\cdot)$ is the unknown function to be minimized, e.g., $f(\cdot)$ can be the risk function $R_g(w) = \mathbb{E}_{z \sim g}(\ell(w, z))$.

Algorithm 6. *Stochastic Gradient Descent (SGD) with learning rate $\eta_t > 0$ for the solution of the minimization problem (2.1)*

For $t = 1, 2, 3, \dots$ iterate:

(1) *compute*

$$(2.2) \quad w^{(t+1)} = w^{(t)} - \eta_t v_t,$$

where v_t is a random vector such that $E(v_t | w^{(t)}) \in \partial f(w^{(t)})$

(2) *terminate if a termination criterion is satisfied, e.g.*

If $t \geq T_{\max}$ then STOP

Remark 7. If f is differentiable at $w^{(t)}$, it is $\partial f(w^{(t)}) = \{\nabla f(w^{(t)})\}$. Hence v_t is such as $E(v_t | w^{(t)}) = \nabla f(w^{(t)})$ in Algorithm 6 step 1.

Note 8. Assume f is differentiable (for simplicity). To compare SGD with GD, we can re-write (2.2) in the SGD Algorithm 6 as

$$(2.3) \quad w^{(t+1)} = w^{(t)} - \eta_t [\nabla f(w^{(t)}) + \xi_t],$$

where

$$\xi_t := v_t - \nabla f(w^{(t)})$$

represents the (observed) noise introduced in (2.2) by using a random realization of the exact gradient.

Remark 9. Given T SGD algorithm iterations, the output of SGD can be (but not a exclusively)

(1) the average (after discarding the first few iterations of $w^{(t)}$ for stability reasons)

$$(2.4) \quad w_{\text{SGD}}^{(T)} = \frac{1}{T} \sum_{t=1}^T w^{(t)}$$

(2) or the best value discovered

$$w_{\text{SGD}}^{(T)} = \arg \min_{w_t} (f(w^{(t)}))$$

(3) or the last value discovered

$$w_{\text{SGD}}^{(T)} = w^{(T)}$$

Note 10. SGD output converges to a local minimum, $w_{\text{SGD}}^{(T)} \rightarrow w_*$ (in some sense), under different sets of regularity conditions. Section 4 has a brief analysis. To achieve this, Conditions 11 on the learning rate are rather inevitable and should be satisfied.

Condition 11. Regarding the learning rate (or gain) $\{\eta_t\}$ should satisfy conditions

- (1) $\eta_t \geq 0$,
- (2) $\sum_{t=1}^{\infty} \eta_t = \infty$
- (3) $\sum_{t=1}^{\infty} \eta_t^2 < \infty$

Remark 12. The popular learning rates $\{\eta_t\}$ in Remark 9 in Handout 2 satisfy Condition 11 and hence can be used in SGD too. Once parametrized, η_t can be tuned based on pilot runs using a reasonably small number of data.

Remark 13. Intuition on Condition 11. Assume that v_t is bounded. Condition 11(3) aims at reducing the effect of the randomness in v_t (introduced noise ξ_t) because it implies $\eta_t \searrow 0$ as $t \rightarrow \infty$; if this was not the case then

$$w^{(t+1)} - w^{(t)} = -\eta_t v_t \rightarrow 0$$

may not be satisfied and the chain $\{w^{(t)}\}$ may not converge. Condition 11(2) prevents η_t from reducing too fast and allows the generated chain $\{w^{(t)}\}$ to be able to converge. E.g., after t iterations

$$\begin{aligned} \|w^{(t)} - w^*\| &= \|w^{(t)} \pm w^{(0)} - w^*\| \geq \|w^{(0)} - w^*\| - \|w^{(t)} - w^{(0)}\| \\ &\geq \|w^{(0)} - w^*\| - \sum_{t=0}^{\infty} \|w^{(t+1)} - w^{(t)}\| = \|w^{(0)} - w^*\| - \sum_{t=0}^{T-1} \|\eta_t v_t\| \end{aligned}$$

However if it was $\sum_{t=1}^{\infty} \eta_t < \infty$ it would be $\sum_{t=0}^{\infty} \|\eta_t v_t\| < \infty$ and hence $w^{(t)}$ would never converge to w^* if the seed $w^{(0)}$ is far enough from w^* .

3. STOCHASTIC GRADIENT DESCENT WITH PROJECTION

Remark 14. Consider the scenario in Problem 1 where the learning problem requires to discover w^* in the restricted/bounded set \mathcal{H} . Assume the function to be minimized is convex in the restricted hypothesis set \mathcal{H} , e.g. $\mathcal{H} = \{w : \|w\| \leq B\}$, but non-convex in \mathbb{R}^d . Direct implementation of vanilla SGD (Algorithm 6) may produce a chain stepping out \mathcal{H} and hence an output $w_{\text{SGD}} \notin \mathcal{H}$. SGD can be modified to address this issue, as in Algorithm 15, by including a projection step guarantying $w \in \mathcal{H}$.

Algorithm 15. *Stochastic Gradient Descent with learning rate $\eta_t > 0$ and with projection in \mathcal{H} for problem in (2.1)*

For $t = 1, 2, 3, \dots$ iterate:

(1) compute

$$(3.1) \quad w^{(t+\frac{1}{2})} = w^{(t)} - \eta_t v_t,$$

where v_t is a random vector such that $E(v_t | w^{(t)}) \in \partial f(w^{(t)})$

(2) compute

$$(3.2) \quad w^{(t+1)} = \arg \min_{w \in \mathcal{H}} \left(\|w - w^{(t+\frac{1}{2})}\| \right)$$

(3) terminate if a termination criterion is satisfied

4. ANALYSIS OF SGD (ALGORITHM 6)

Note 16. Recall that the stochasticity of SGD comes from the stochastic sub-gradients $\{v_t\}_{t=1}^T$; hence the expectations below are under these random vectors' distributions.

Theorem 17. *Let $f(\cdot)$ be a convex function. If we run SGD algorithm of f with learning rate $\eta_t > 0$ for T steps, the output $w_{\text{GD}}^{(T)} = \frac{1}{T} \sum_{t=1}^T w^{(t)}$ satisfies*

$$(4.1) \quad E \left(f \left(w_{\text{GD}}^{(T)} \right) \right) - f(w^*) \leq \frac{\|w^*\|^2}{2\eta T} + \frac{\eta}{2} \frac{1}{T} \sum_{t=1}^T E \|v_t\|^2$$

Proof. Let $v_{1:t} = (v_1, \dots, v_t)$. By Jensens' inequality (or see (4.3) in Handout 2)

$$(4.2) \quad E \left(f \left(w_{\text{GD}}^{(T)} \right) - f(w^*) \right) \leq E \left(\frac{1}{T} \sum_{t=1}^T \left(f(w^{(t)}) - f(w^*) \right) \right) = \frac{1}{T} \sum_{t=1}^T E \left(f(w^{(t)}) - f(w^*) \right)$$

I will try to use Lemma 22 from Handout 2, hence I need to show

$$(4.3) \quad E \left(f(w^{(t)}) - f(w^*) \right) \leq E \left(\langle w^{(t)} - w^*, v_t \rangle \right)$$

where the expectation is under $v_{1:T}$. It is

$$\begin{aligned} E_{v_{1:T}} \left(\langle w^{(t)} - w^*, v_t \rangle \right) &= E_{v_{1:t}} \left(\langle w^{(t)} - w^*, v_t \rangle \right) \\ &= E_{v_{1:t-1}} \left(E_{v_{1:t}} \left(\langle w^{(t)} - w^*, v_t \rangle | v_{1:t-1} \right) \right) \quad (\text{law of total expectation}) \end{aligned}$$

But $w^{(t)}$ is fully determined by $v_{1:t-1}$, (see (2.2)) so

$$E_{v_{1:t-1}} \left(E_{v_{1:t}} \left(\langle w^{(t)} - w^*, v_t \rangle | v_{1:t-1} \right) \right) = E_{v_{1:t-1}} \left(\langle w^{(t)} - w^*, E_{v_{1:t}}(v_t | v_{1:t-1}) \rangle \right)$$

As $w^{(t)}$ is fully determined by $v_{1:t-1}$ then $\mathbb{E}_{v_{1:t}}(v_t|v_{1:t-1}) = \mathbb{E}_{v_{1:t}}(v_t|w^{(t)}) \in \partial f(w^{(t)})$, hence $\mathbb{E}_{v_{1:t}}(v_t|v_{1:t-1})$ is a sub-gradient. By sub-gradient definition

$$(4.4) \quad \begin{aligned} \mathbb{E}_{v_{1:t-1}} \left(\langle w^{(t)} - w^*, \mathbb{E}_{v_{1:t}}(v_t|v_{1:t-1}) \rangle \right) &\geq \mathbb{E}_{v_{1:t-1}} \left(f(w^{(t)}) - f(w^*) \right) \\ &= \mathbb{E}_{v_{1:T}} \left(f(w^{(t)}) - f(w^*) \right) \end{aligned}$$

Hence combining (4.4), (4.3), and (4.2)

$$\mathbb{E} \left(f(w_{\text{SGD}}^{(T)}) - f(w^*) \right) \leq \frac{1}{T} \sum_{t=1}^T \mathbb{E} \left(\langle w^{(t)} - w^*, v_t \rangle \right)$$

Then Lemma 22 in Handout 2 implies

$$\mathbb{E} \left(f(w_{\text{SGD}}^{(T)}) - f(w^*) \right) \leq \mathbb{E} \left(\frac{1}{T} \frac{\|w^*\|^2}{2\eta} + \frac{\eta}{2} \frac{1}{T} \sum_{t=1}^T \|v_t\|^2 \right) = \frac{\mathbb{E} \|w^*\|^2}{2\eta T} + \frac{\eta}{2} \frac{1}{T} \sum_{t=1}^T \mathbb{E} \|v_t\|^2$$

□

Remark 18. The upper bound in (4.1) depends on the variation of v_t as

$$(4.5) \quad \mathbb{E} \|v_t\|^2 = \sum_{j=1}^d \text{Var}(v_{t,j}) + \sum_{j=1}^d (\mathbb{E}(v_{t,j}))^2$$

where d is the dimension of $v_t = (v_{t,1}, \dots, v_{t,d})^\top$. The second term on the right hand side of (4.5) is constant as $v_{t,j}$ is the unbiased estimator of the sub-gradient by construction.

Proposition 19. (*Cont. Theorem 17*) Let $f(\cdot)$ be a convex function, and let $\mathcal{H} = \{w \in \mathbb{R} : \|w\| \leq B\}$. Let $\mathbb{E} \|v_t\|^2 \leq \rho^2$. Assume we run SGD algorithm of $f(\cdot)$ with learning rate $\eta_t = \sqrt{\frac{B^2}{\rho^2 T}}$ for T steps, and output $w_{\text{SGD}}^{(T)} = \frac{1}{T} \sum_{t=1}^T w^{(t)}$. Then

(1) upper bound on the sub-optimality is

$$(4.6) \quad \mathbb{E} \left(f(w_{\text{SGD}}^{(T)}) \right) - f(w^*) \leq \frac{B\rho}{\sqrt{T}}$$

(2) a given level off accuracy ε such that $\mathbb{E} \left(f(w_{\text{SGD}}^{(T)}) \right) - f(w^*) \leq \varepsilon$ can be achieved after T iterations

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

Proof. It follows from 17. □

Remark 20. Lemma 22 in Handout 2 holds even when projection steps are added in the vanilla SGD algorithm, hence the above analysis holds for the SGD with projections (Algorithm 15) too.

5. IMPLEMENTATION OF SGD IN THE LEARNING PROBLEM 1

Note 21. Below we show SGD can be implement in the learning Problem 1. Proposition 22 allows a convenient implementation of SGD to the learning problem (Problem 1).

Proposition 22. For a randomly drawn $z \sim g(\cdot)$, the sub-gradient v of $\ell(w, z)$ at point w is an unbiased estimator of the sub-gradient of the risk $R_g(w)$ at point w .

Proof. Let v be a sub-gradient of $\ell(w, z)$ at point w , then

$$(5.1) \quad \ell(u, z) - \ell(w, z) \geq \langle u - w, v \rangle$$

It is

$$\begin{aligned} R_g(u) - R_g(w) &= \mathbb{E}_{z \sim g} (\ell(u, z) - \ell(w, z) | w) \geq \mathbb{E}_{z \sim g} (\langle u - w, v \rangle | w) \\ &= \langle u - w, \mathbb{E}_{z \sim g} (v | w) \rangle \end{aligned}$$

Hence, by definition, v is such that $\mathbb{E}_{z \sim g} (v | w)$ is a sub-gradient of $R_g(w)$. \square

Example 23. Consider that loss ℓ is differentiable wrt w . If $v = \nabla_w \ell(w, z)$, then

$$\mathbb{E}_{z \sim g} (v | w) = \nabla_w R_g(w)$$

Note 24. Below we show how SGD can be implement in the learning Problem 1.

Note 25. Assume there is available a finite dataset $\mathcal{S}_n = \{z_i; i = 1, \dots, n\}$ of size n which consists of independent realizations z_i from the data generating distribution g ; $z_i \stackrel{\text{ind}}{\sim} g$. Batch sampling (Algorithm 26) is an implementation of the SGD (Algorithm 6) in the learning Problem 1.

Algorithm 26. batch Stochastic Gradient Descent with learning rate $\eta_t > 0$, and batch size m , for Problem 1.

For $t = 1, 2, 3, \dots$ iterate:

(1) get a random sub-sample $\{\tilde{z}_j^{(t)}; j = 1, \dots, m\}$ of size m with or without replacement from the complete data-set \mathcal{S}_n .

(2) compute

$$(5.2) \quad w^{(t+1)} = w^{(t)} - \eta_t v_t,$$

where $v_t = \frac{1}{m} \sum_{j=1}^m \tilde{v}_{t,j}$ and $\mathbb{E}(\tilde{v}_{t,j} | w^{(t)}) \in \partial_w \ell(w^{(t)}, \tilde{z}_j^{(t)})$.

(3) terminate if a termination criterion is satisfied

Remark 27. Step 1 can be presented equivalently as

(1) Randomly generate a set $\mathcal{J}^{(t)} \subseteq \{1, \dots, n\}^m$ of m indices from 1 to n with or without replacement, and set a $\tilde{\mathcal{S}}_m = \{z_i; i \in \mathcal{J}^{(t)}\}$.

Hence $v_t = \frac{1}{m} \sum_{j \in \mathcal{J}^{(t)}} v_{t,j}$ where $\mathbb{E}(\tilde{v}_{t,j} | w^{(t)}) \in \partial_w \ell(w^{(t)}, z_j)$.

Remark 28. A projection step, such as (3.2) can be added right after step 2 in Algorithm 26, if needed.

Remark 29. If it is possible to sample anytime fresh examples z_i directly from the data generation model g instead of just having access to only a given finite dataset of examples \mathcal{S}_n , then step 1 in Algorithm 26 can become

(1) sample $\tilde{z}_j^{(t)} \sim g(\cdot)$ for $j = 1, \dots, m$.

Definition 30. Online Stochastic gradient descent is the special case of Algorithm 26 using sub-samples of size one ($m = 1$), namely, only one example is randomly chosen in Step 1 of Algorithm 26.

Remark 31. In theory, using larger batch size m has the benefit that reduces the variance $\text{Var}(v_t|w^{(t)})$ of v_t at iteration t due to averaging effect, stabilizes the SGD algorithm, and reduces the error bound (4.1); see Remark 18.

Remark 32. In practice, for a given fixed computational time, using smaller batch size m has the benefit that the algorithm iterates faster as each iteration processes less number of examples. E.g., consider the extreme cases GD vs online SGD in a scenario with big-data: if the dataset consists of several replications of the same values, GD (using all the data) has to process the same information multiple times, while the online SGD (using only one example at a time) would avoid this issue.

Remark 33. In practice, for a given fixed computational time, it is possible for a batch SGD with smaller batch size m to present better generalization properties (wrt the theoretical assumptions) than those with larger m (GD is included). It is observed for the former to be often less prone to getting stuck in shallow local minima because of the additional amount of “noise” E.g., consider the extreme cases GD vs online SGD in a scenario with non-convex risk function (e.g. our theoretical assumptions as violated): if the Risk function presents local minima, considering less examples randomly chosen each time may cause fluctuations in the gradient that allow the chain to accidentally jump/escape to an area with a lower minimum.

6. STOCHASTIC VARIANCE REDUCED GRADIENT (SVRG)

Remark 34. Recall the upper bound of the error (4.1) in SGD depends on the variance of the stochastic gradient, as shown in Remark 4.5. Hence the algorithm may be improved by reducing the variance of each element of v_t .

Remark 35. Control variate is a general way to perform variance reduction. Let random variables $v \in \mathbb{R}$, and $y \in \mathbb{R}$. Let $z = v + c(y - \mathbb{E}(y))$ for some constant $c \in \mathbb{R}$. It is $\mathbb{E}_c(z) = \mathbb{E}(v)$ and

$$\text{Var}_c(z) = \text{Var}(v) + c^2 \text{Var}(y) + 2c \text{Cov}(v, y)$$

which is minimized for

$$c^* = -\frac{\text{Cov}(v, y)}{\text{Var}(y)}$$

hence

$$\text{Var}_{c^*}(z) = \text{Var}(v) - \frac{(\text{Cov}(v, y))^2}{\text{Var}(y)}$$

Note 36. For simplicity, SVRG is presented in the case where $\kappa = 1$ and the loss function $\ell(w, z)$ is differentiable wrt w hence its gradient exists. However it is applicable to the more general cases introduced.

Remark 37. Every m iterations, SVRG keeps a snapshot \tilde{w} , and computes the gradient using all data i.e. $\frac{1}{n} \sum_{i=1}^n \ell(\tilde{w}, z_i)$. At each iteration t , the update is

$$w^{(t+1)} = w^{(t)} - \eta_t \left[\underbrace{\underbrace{\nabla \ell(w^{(t)}, \tilde{z}^{(t)})}_{=v} - \underbrace{1}_{=c}}_z \underbrace{\left(\nabla \ell(\tilde{w}, \tilde{z}^{(t)}) - \frac{1}{n} \sum_{i=1}^n \nabla \ell(\tilde{w}, z_i) \right)}_y \right]$$

given that a random $\tilde{z}^{(t)}$ has been collected from the sample. The symbols below the brackets are given with reference to Remark 35.

Algorithm 38. *Stochastic Variance Reduced Gradient with learning rate $\eta_t > 0$ for Problem 1.*

For $t = 1, 2, 3, \dots$ iterate:

- (1) randomly get an example $\tilde{z}^{(t)}$ from S_n .
- (2) compute

$$(6.1) \quad w^{(t+1)} = w^{(t)} - \eta_t \left[\nabla \ell(w^{(t)}, \tilde{z}^{(t)}) - \nabla \ell(\tilde{w}, \tilde{z}^{(t)}) + \frac{1}{n} \sum_{i=1}^n \nabla \ell(\tilde{w}, z_i) \right]$$

- (3) if modulo $(t, \kappa) = 0$, set $\tilde{w} = w^{(t)}$
- (4) if a termination criterion is satisfied STOP

Remark 39. Iterations of SVRG are computationally faster than those of full GD, but SVRG can still match the theoretical convergence rate of GD.

Remark 40. How often we get snapshots, κ , in Algorithm 38 is specified by the researcher. The **smaller** the κ , the more frequent snapshots, and the more correlated the baseline y will be with the objective x and hence the bigger the performance improvement; however the iterations will be slower.

7. PRECONDITIONED SGD; THE ADAGRAD ALGORITHM

Remark 41. All SGD (introduced earlier) are first order methods in the sense they consider only the gradient. Their advantage is each iteration is fast. The disadvantage is that they ignore the curvature of the space and hence can be slower to converge in cases the curvature changes eg. among dimensions of w .

Remark 42. To address this, SGD (Algorithm 6) can be modified in the update step as in Algorithm 43 by using a preconditioner P_t that accounts for the curvature (or geometry in general of f).

Algorithm 43. *Preconditioned Stochastic Gradient Descent with learning rate $\eta_t > 0$, and preconditioner P_t for the solution of the minimization problem (2.1)*

For $t = 1, 2, 3, \dots$ iterate:

(1) *compute*

$$(7.1) \quad w^{(t+1)} = w^{(t)} - \eta_t P_t v_t,$$

where v_t is a random vector such that $E(v_t | w^{(t)}) \in \partial f(w^{(t)})$, and P_t a preconditioner

(2) *terminate if a termination criterion is satisfied, e.g.*

If $t \geq T$ then STOP

Remark 44. A natural choice of P_t can be $P_t := [H_t + \epsilon I_d]^{-1}$, where H_t is the Hessian matrix $[H_t]_{i,j} = \frac{\partial^2}{\partial w_i \partial w_j} f(w) \Big|_{w=w^{(t)}}$ (ie. the gradient of the gradient's elements), and ϵ is a tiny $\epsilon > 0$ to mitigate machine error when Hessian elements are close to zero.

Remark 45. If the preconditioner P_t is set to be the inverse of the full Hessian, it may be too expensive to perform matrix operations in (7.1) with the full Hessian, and such operations can be too unstable/inaccurate due to the random error induced by the stochasticity of the gradient.

7.1. Adaptive Stochastic Gradient Decent (AdaGrad).

Remark 46. AdaGrad dynamically incorporates knowledge of the geometry of function $f(\cdot)$ (to be minimized) in earlier iterations to perform more informative gradient-based learning.

Remark 47. AdaGrad aims to perform larger updates (i.e. high learning rates) for those dimensions of w that are related to infrequent features (largest partial derivative) and smaller updates (i.e. low learning rates) for frequent ones (smaller partial derivative).

Remark 48. Hence, this strategy often improves convergence performance over standard stochastic gradient descent in settings where data is sparse and sparse features w 's are more informative.

Definition 49. Adaptive Stochastic Gradient Decent (AdaGrad) can be presented in terms of preconditioned SGD (Algorithm 43) with preconditioner $P_t = [I_d \text{diag}(G_t) I_d + \epsilon I_d]^{-1/2}$ in (7.1)

$$(7.2) \quad w^{(t+1)} = w^{(t)} - \eta_t [I_d \text{diag}(G_t) I_d + \epsilon I_d]^{-1/2} v_t,$$

where $\text{diag}(A) := (A_{1,1}, A_{2,2}, \dots, A_{d,d})^\top$ is the vector of the diagonal of a matrix A , $G_t = \sum_{\tau=1}^t v_\tau^\top v_\tau$ is the sum of the outer products of the gradients $\{v_\tau; \tau \leq t\}$ up to the state t , and $\epsilon > 0$ is a tiny value (eg, 10^{-6}) set for computational stability in case the gradient becomes too close to zero.

Remark 50. AdaGrad algorithm individually adapts the learning rate of each dimension of w_t by scaling them inversely proportional to the square root of the sum of all the past squared values of the gradient $\{v_\tau; \tau \leq t\}$. This is because

$$(7.3) \quad [G_t]_{j,j} = \sum_{\tau=1}^t (v_{\tau,j})^2$$

where j denotes the j -th dimension of w . Hence (7.2) implies that the j -th dimension of w is updated as

$$w_j^{(t+1)} = w_j^{(t)} - \eta_t \frac{1}{\sqrt{[G_t]_{j,j} + \epsilon}} v_{t,j}.$$

Remark 51. The accumulation of positive terms in (7.3) makes the sum keep growing during training and causes the learning rate to shrink and becoming infinitesimally small. This offers an automatic way to choose a decreasing learning rate simplifying setting the learning rate; however it may result in a premature and excessive decrease in the effective learning rate. This can be mitigated by still considering in (7.2) a (user specified rate) $\eta_t \geq 0$ and tuning it properly via pilot runs.

8. EXAMPLE¹

Note 52. We continue Example in Section 5 in Handout 2. Recall, we considered a hypothesis space \mathcal{H} of linear functions $h : \mathbb{R}^2 \rightarrow \mathbb{R}$ with $h(w) = w_1 + w_2 x$, $w = (w_1, w_2)^\top$, and $\ell(w, z = (x, y)^\top) = (y_i - w_1 - w_2 x)^2$. Here we consider a big dataset $\mathcal{S}_n = \{z_1, \dots, z_n\}$ with $n = 10^6$ examples.

Example 53. The batch SGD algorithm (Algorithm 26) with learning rate η_t and batch size $m = 10$ is

- For $t = 1, 2, 3, \dots$ iterate:
 - (1) Randomly generate a set $\mathcal{J}^{(t)}$ by drawing $m = 10$ numbers from $\{1, \dots, n = 10^6\}$, and set $\tilde{\mathcal{S}}_m = \{z_i; i \in \mathcal{J}^{(t)}\}$
 - (2) compute

$$w^{(t+1)} = w^{(t)} - \eta_t v_t,$$

where

$$v_t = \begin{pmatrix} 2w_1^{(t)} + 2w_2^{(t)} \frac{1}{m} \sum_{j \in \mathcal{J}^{(t)}} x_j - 2 \frac{1}{m} \sum_{j \in \mathcal{J}^{(t)}} y_j \\ 2w_1^{(t)} \bar{x} + 2w_2^{(t)} \frac{1}{m} \sum_{j \in \mathcal{J}^{(t)}} x_j^2 - 2 \sum_{j \in \mathcal{J}^{(t)}} y_j x_j \end{pmatrix},$$

- (3) if $t \geq T = 1000$ STOP

because

$$v_t = \frac{1}{m} \sum_{j \in \mathcal{J}^{(t)}} \nabla \ell(w^{(t)}, \tilde{z}^{(t)}) = \dots = \begin{pmatrix} 2w_1^{(t)} + 2w_2^{(t)} \frac{1}{m} \sum_{j \in \mathcal{J}^{(t)}} x_j - 2 \frac{1}{m} \sum_{j \in \mathcal{J}^{(t)}} y_j \\ 2w_1^{(t)} \bar{x} + 2w_2^{(t)} \frac{1}{m} \sum_{j \in \mathcal{J}^{(t)}} x_j^2 - 2 \sum_{j \in \mathcal{J}^{(t)}} y_j x_j \end{pmatrix}$$

In Figures 8.1a & 8.1d, we observe that increasing the batch size has improved the convergence however this is not a panacea. Also it had reduced the oscillations of chain $\{w^{(t)}\}$.

Example 54. The SVRG with learning rate $\eta_t > 0$ and batch size $m = 1$ is

- For $t = 1, 2, 3, \dots$ iterate:
 - (1) randomly generate a set $\mathcal{J}^{(t)} = \{j^*\}$ by drawing one number j^* from $\{1, \dots, n = 10^6\}$, and set $\tilde{\mathcal{S}}_1 = \{z_{j^*}\}$

¹Code can be found in https://github.com/georgios-stats/Machine_Learning_and_Neural_Networks_III_Epiphany_2023/tree/main/Lecture_handouts/code/03.Stochastic_gradient_descent

(2) compute

$$(8.1) \quad w^{(t+1)} = \begin{bmatrix} w_1^{(t)} \\ w_2^{(t)} \end{bmatrix} - \eta_t \left(\begin{bmatrix} 2(w_1^{(t)} - \tilde{w}_1^{(t)}) + 2w_2^{(t)}x_{j^*} \\ 2(w_2^{(t)} - \tilde{w}_2^{(t)})\bar{x} + 2(w_2^{(t)}(x_{j^*})^2 - \tilde{w}_2^{(t)}(x_{j^*})^2) \end{bmatrix} + \nabla \hat{R}_{\mathcal{D}}(\tilde{w}) \right)$$

(3) if modulo $(t, \kappa) = 0$,

(a) set $\tilde{w} = w^{(t)}$

(b) compute

$$(8.2) \quad \frac{1}{n} \sum_{i=1}^n \ell(\tilde{w}, z_i) = \begin{pmatrix} 2\tilde{w}_1^{(t)} + 2\tilde{w}_2^{(t)}\bar{x} - 2\bar{y} \\ 2\tilde{w}_1^{(t)}\bar{x} + 2\tilde{w}_2^{(t)}\bar{x}^2 - 2\bar{y}^\top \bar{x} \end{pmatrix} = \nabla \hat{R}_{\mathcal{D}}(\tilde{w})$$

(4) if a termination criterion is satisfied STOP

Because (8.2) is actually the gradient of the Risk function at \tilde{w} and

$$\nabla \ell(w^{(t)}, z_{j^*}) - \nabla \ell(\tilde{w}, z_{j^*}) = \begin{pmatrix} 2(w_1^{(t)} - \tilde{w}_1^{(t)}) + 2w_2^{(t)}x_{j^*} \\ 2(w_2^{(t)} - \tilde{w}_2^{(t)})\bar{x} + 2(w_2^{(t)}(x_{j^*})^2 - \tilde{w}_2^{(t)}(x_{j^*})^2) \end{pmatrix}$$

In Figure 8.1c we observe the frequency of the snapshots has improved the convergence; however this is not a panacea as seen in Figure 8.1f.

Example 55. The AdaGrad with $\eta_t = 1$, $\epsilon = 10^{-6}$, and batch size $m = 1$ is

- Set $G_0 = (0, 0)^\top$.
- For $t = 1, 2, 3, \dots$ iterate:
 - (1) randomly generate a set $\mathcal{J}^{(t)} = \{j^*\}$ by drawing one number from $\{1, \dots, n = 10^6\}$, and set $\tilde{\mathcal{S}}_1 = \{z_{j^*}\}$
 - (2) compute

$$v_t = \begin{pmatrix} 2w_1^{(t)} + 2w_2^{(t)}x_{j^*} - 2y_{j^*} \\ 2w_1^{(t)}\bar{x} + 2w_2^{(t)}\bar{x}^2 - 2y_{j^*}\bar{x} \end{pmatrix}$$

$$G_t = G_{t-1} + \begin{pmatrix} v_{t,1}^2 \\ v_{t,2}^2 \end{pmatrix}$$

$$w^{(t+1)} = \begin{pmatrix} w_1^{(t)} - \frac{\eta_t}{\sqrt{G_{t,1} + \epsilon}} v_{t,1} \\ w_2^{(t)} - \frac{\eta_t}{\sqrt{G_{t,2} + \epsilon}} v_{t,2} \end{pmatrix},$$

(3) if $t \geq T = 1000$ STOP

because

$$v_t = \nabla \ell(w^{(t)}, z_{j^*}) = \dots = \begin{pmatrix} 2w_1^{(t)} + 2w_2^{(t)}x_{j^*} - 2y_{j^*} \\ 2w_1^{(t)}\bar{x} + 2w_2^{(t)}\bar{x}^2 - 2y_{j^*}\bar{x} \end{pmatrix}$$

In Figures 8.1g & 8.1h, we see that AdaGrad with $\eta = 1$ works (I did not try to tune it), however to make vanilla SGD to work I have to tune $\eta = 0.03$ otherwise for $\eta = 1.0$ it did not work.

Example 56. Consider a (rather naive) loss function $\ell(w, z = (x, y)) = -\cos(0.5(y - w_1 - w_2x))$, a hypothesis class $\mathcal{H} = \{w \in \mathbb{R}^2 : \|w\| \leq 1.5\}$, and assume that inputs x in dataset \mathcal{D} are such that

$x \in [-1, 1]$. Note that $-\cos(\cdot)$ is convex in $[-1.5, 1.5]$ and non-convex in \mathbb{R} . Consider learning rate $\eta_t = 50/t$ reducing to zero, and seed $w^{(0)} = (1.5, 1.5)$. An unconstrained SGD may produce a minimizer/solution outside \mathcal{H} because η_t is too large at the first few iterations. We can design the online SGD (batch size $m = 1$) with projection to \mathcal{H} as

- For $t = 1, 2, 3, \dots$ iterate:

- (1) randomly generate a set $\mathcal{J}^{(t)} = \{j^*\}$ by drawing one number from $\{1, \dots, n = 10^6\}$, and set $\tilde{\mathcal{S}}_1 = \{z_{j^*}\}$

- (2) compute

$$w^{(t+1/2)} = w^{(t)} - \frac{50}{t} \begin{pmatrix} \sin(y_{j^*} - w_1^{(t)} - w_2^{(t)} x_{j^*}) \\ \sin(y_{j^*} - w_1^{(t)} - w_2^{(t)} x_{j^*}) x_{j^*} \end{pmatrix}$$

$$w^{(t+1/2)} = \arg \min_{\forall w: \|w\| \leq 1.5} (\|w - w^{(t+1/2)}\|)$$

- (3) if $t \geq T = 1000$ STOP

because

$$v_t = \nabla \ell(w^{(t)}, z_{j^*}) = \dots = \begin{pmatrix} \sin(y_{j^*} - w_1^{(t)} - w_2^{(t)} x_{j^*}) \\ \sin(y_{j^*} - w_1^{(t)} - w_2^{(t)} x_{j^*}) x_{j^*} \end{pmatrix}$$

In Figures 8.1b & 8.1e, we observe that the SGD got trapped outside \mathcal{H} due to the unreasonably large learning rate at the beginning of the iterations, while the SGD with projection step managed to stay in \mathcal{H} and converge.



(A) batch SGD Example 53 w_1



(B) online SGD with projection Example 56 w_1



(C) SVRG Example 54 w_1



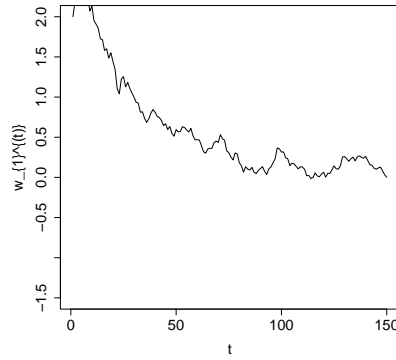
(D) batch SGD Example 53 w_2



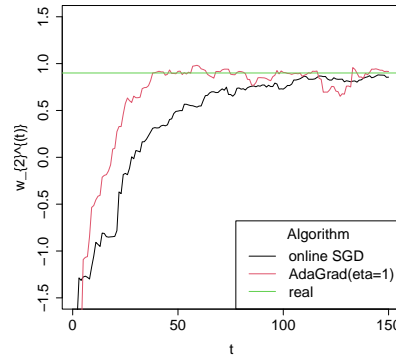
(E) online SGD with projection Example 56 w_2



(F) SVRG Example 54 w_2



(G) online AdaGrad Example 55 w_1



(H) online AdaGrad Example 55 w_2

FIGURE 8.1. Simulations of the Examples