

Lecture notes 6: Variations of stochastic gradient descent

Lecturer & author: Georgios P. Karagiannis

georgios.karagiannis@durham.ac.uk

Aim. To introduce the stochastic gradient descent variations for restricted sets, adaptive step functions, and reduced errors.

Reading list & references:

- (1) Johnson, Rie, and Tong Zhang. "Accelerating stochastic gradient descent using predictive variance reduction." Advances in neural information processing systems 26 (2013).
- (2) Duchi, John, Elad Hazan, and Yoram Singer. "Adaptive subgradient methods for online learning and stochastic optimization." Journal of machine learning research 12, no. 7 (2011).

1. STOCHASTIC GRADIENT DESCENT WITH PROJECTION

Problem 1. Consider the minimization problem

$$(1.1) \quad w^* = \arg \min_w (f(w))$$

where $f : \mathbb{R}^d \rightarrow \mathbb{R}$, $w \in \mathcal{W}$, and assume $\mathcal{W} \subset \mathbb{R}^d$ is a restricted domain.

Note 2. Consider Problem 1 requiring to discover w^* in a restricted/bounded set $\mathcal{W} \subset \mathbb{R}^d$, e.g. $\mathcal{W} = \{w \in \mathbb{R}^d : \|w\| \leq B\}$ for $B > 0$. If f is convex over \mathcal{W} , but non-convex in \mathbb{R}^d/\mathcal{W} , direct implementation of (standard) SGD (Algorithm 6) may produce a chain stepping out \mathcal{W} and hence an output $w_{\text{SGD}} \notin \mathcal{W}$.

Note 3. Stochastic gradient descent with projection (prjSGD) is a modified SGD suitable for applications with restricted sets and relies on the consideration of a projection step guaranteeing $w \in \mathcal{W}$.

Algorithm 4. *Stochastic Gradient Descent with learning rate $\eta_t > 0$ and with projection in \mathcal{W} for Problem1*

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

(1) compute

$$(1.2) \quad w^{(t+\frac{1}{2})} = w^{(t)} - \eta_t v_t,$$

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

(2) compute

$$(1.3) \quad w^{(t+1)} = \arg \min_{w \in \mathcal{W}} \left(\|w - w^{(t+\frac{1}{2})}\| \right)$$

(3) terminate if a termination criterion is satisfied

Note 5. The analysis of SGD with projection is given as a Exercise 13 from the Exercise sheet.

Example 6.¹ Consider a (rather naive) loss function $\ell(w, z = (x, y)) = -\cos(0.5(y - w_1 - w_2 x))$, a hypothesis class $\mathcal{H} = \{h_w(x) = w^\top x : w \in \mathcal{W}\}$ and $\mathcal{W} = \{w \in \mathbb{R}^2 : \|w\| \leq 1.5\}$, and assume that inputs x in dataset S 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 ($m = 1$) 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 $\tilde{S}_1 = \{z_{j^*}\}$

(2) compute

$$w^{(t+1/2)} = w^{(t)} - \frac{25}{t} \begin{pmatrix} \sin\left(0.5\left(y_{j^*} - w_1^{(t)} - w_2^{(t)} x_{j^*}\right)\right) \\ \sin\left(0.5\left(y_{j^*} - w_1^{(t)} - w_2^{(t)} x_{j^*}\right)\right) x_{j^*} \end{pmatrix}$$

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

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

because

$$v_t = \nabla \ell(w^{(t)}, z_{j^*}) = \begin{pmatrix} \frac{\partial}{\partial w_1^{(t)}} \ell(w^{(t)}, z_{j^*}) \\ \frac{\partial}{\partial w_2^{(t)}} \ell(w^{(t)}, z_{j^*}) \end{pmatrix} = \begin{pmatrix} 0.5 \sin\left(0.5\left(y_{j^*} - w_1^{(t)} - w_2^{(t)} x_{j^*}\right)\right) \\ 0.5 \sin\left(0.5\left(y_{j^*} - w_1^{(t)} - w_2^{(t)} x_{j^*}\right)\right) x_{j^*} \end{pmatrix}$$

In Figures A.1b & A.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.

¹https://github.com/georgios-stats/Machine_Learning_and_Neural_Networks_III_Epiphany_2025/tree/main/Lecture_notes/code/example_projSGD.R

2. PRECONDITIONED SGD AND THE ADAGRAD ALGORITHM

Problem 7. Consider the minimization problem

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

where $f : \mathcal{W} \rightarrow \mathbb{R}$, and assume $\mathcal{W} \equiv \mathbb{R}^d$.

Note 8. All SGD (introduced earlier) are first order methods, in the sense they consider only the gradient, and hence they ignore the curvature of the space. Consequently they may present slow converge in cases the curvature changes, eg. among dimensions of w .

Note 9. Preconditioned Stochastic Gradient Descent (Algorithm 10) is a modified SGD aiming to expedite convergence of SGD by using a preconditioner P_t that accounts for the curvature (or geometry in general) of f .

Algorithm 10. *Preconditioned Stochastic Gradient Descent with learning rate $\eta_t > 0$, and preconditioner P_t for the solution of the minimization problem (3.1)*

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

(1) *compute*

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

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

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

If $t \geq T$ then STOP

Note 11. 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.

Note 12. Using the inverse of the full Hessian as preconditioner P_t may be too expensive to perform matrix operations in (2.2) and operations can be too unstable/inaccurate due to the random error induced by the stochasticity of the gradient.

2.1. Adaptive Stochastic Gradient Decent (AdaGrad).

Note 13. AdaGrad (Algorithm 15) is a modified SGD algorithm that aims to adaptively adjust the learning rate of SGD, perform more informative gradient-based learning, and by improve convergence performance over standard SGD by dynamically incorporating knowledge of the geometry of function $f(\cdot)$ in earlier iterations.

Note 14. It is particularly suitable in cases where the data are sparse and sparse features w 's are more informative. It 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).

Algorithm 15. *Adaptive Stochastic Gradient Decent (AdaGrad) is a preconditioned SGD (Algorithm 10) with preconditioner $P_t = [\text{diag}(G_t) + \epsilon I_d]^{-1/2}$ where notation $\text{diag}(A)$ denotes a $d \times d$ matrix whose diagonal is the d dimensional diagonal vector $(A_{1,1}, A_{2,2}, \dots, A_{d,d})$ of $d \times d$ matrix A and whose off-diagonal elements are zero, $G_t = \sum_{\tau=1}^t v_\tau v_\tau^\top$ 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. I.e.*

$$(2.3) \quad w^{(t+1)} = w^{(t)} - \eta_t [\text{diag}(G_t) + \epsilon I_d]^{-1/2} v_t.$$

Note 16. AdaGrad 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

$$(2.4) \quad [G_t]_{j,j} = \sum_{\tau=1}^t (v_{\tau,j})^2$$

where j denotes the j -th dimension of w . Hence (2.3) and (2.4) imply that the j -th dimension of w is updated as

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

Note 17. The accumulation of positive terms in (2.4) makes the sum keep growing during training and causes the learning rate to shrink and become 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 (2.3) a (user specified rate) $\eta_t \geq 0$ and tuning it properly via pilot runs.

Example 18. ² 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 $j^{(t)}$ from $\{1, \dots, n = 10^6\}$
 - (2) compute

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

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

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

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

²https://github.com/georgios-stats/Machine_Learning_and_Neural_Networks_III_Epiphany_2025/tree/main/Lecture_notes/code/example_AdaGrad.R

because

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

In Figures A.1g & A.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.

3. STOCHASTIC VARIANCE REDUCED GRADIENT (SVRG)

Note 19. Recall, in Note 16 of Lecture Notes 5, the upper bound of the error ((3.1); in L.N. 5) in SGD depends on the variance of the stochastic sub-gradient i.e. $E \|v_t - E(v_t)\|^2$.

$$(3.1) \quad E \left(f(w_{\text{SGD}}^{(T)}) \right) - f(w^*) \leq \frac{\|w^*\|^2}{2\eta T} + \frac{\eta}{2} \frac{1}{T} \sum_{t=1}^T E \|v_t\|^2$$

$$(3.2) \quad E \|v_t\|^2 = \underbrace{E \|v_t - E(v_t)\|^2}_{=\text{tr}(\text{Var}(v_t))} + \underbrace{\|E(v_t)\|^2}_{\text{const.}}$$

A way to improve the SGD may be by reducing / controlling the variances at each dimension of v_t .

3.1. Batch Stochastic gradient descent.

Note 20. A way to reduce variance in (3.2) is to consider a larger batch size $m = |\mathcal{J}^{(t)}|$. Note that

$$E_{z \sim g} \|v_t - E \|v_t\|\|^2 = E_{z \sim g} \left\| \nabla_w \frac{1}{|\mathcal{J}^{(t)}|} \sum_{j \in \mathcal{J}^{(t)}} \ell(w, z_j) - \nabla_w R_g(w) \right\|^2 = \frac{E_{z_j \sim g} \|\nabla_w \ell(w, z_j) - \nabla_w R_g(w)\|^2}{|\mathcal{J}^{(t)}|}$$

where the top part of the fraction is the gradient variance for SGD with $m = 1$.

3.2. Stochastic Variance Reduced Gradient (SVRG).

Note 21. Control variate is a general way to perform variance reduction. Let $x \in \mathbb{R}^d$ be an unbiased estimator of θ . The goal is to compute $v = x + c(y - E(y))$ by finding a suitable control variable $y \in \mathbb{R}^d$ and constant $c \in \mathbb{R}$ such that $\text{tr}(\text{Var}_c(v)) < \text{tr}(\text{Var}(x))$ with purpose to use it as an unbiased estimator of θ instead of x $E_c(v) = E(x) = \theta$.

Example 22. Let random variables $x \in \mathbb{R}^d$, and $y \in \mathbb{R}^d$. Let $v = x + c(y - E(y))$ for some constant $c \in \mathbb{R}$. It is $E_c(v) = E(x)$. Also

$$E_c \|v - E(v)\|^2 = E \|x - E(x)\|^2 + c^2 E \|y - E(y)\|^2 + 2c E \langle x, y \rangle$$

which is minimized for

$$0 = \frac{d}{dc} E_c \|v - E(v)\|^2 \Big|_{c=c^*} \implies c^* = - \frac{E \langle x, y \rangle}{E \|y - E(y)\|^2}$$

with

$$(3.3) \quad E_{c^*} \|v - E(v)\|^2 = E \|x - E(x)\|^2 - \frac{E \langle x, y \rangle^2}{E \|y - E(y)\|^2}$$

Algorithm 23. *Stochastic Variance Reduced Gradient with learning rate $\eta_t > 0$ for Problem 1.*

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

(1) randomly draw an example $z_{j(t)}$ indexed by $j(t) \in \{1, \dots, n\}$.

(2) compute

$$(3.4) \quad w^{(t+1)} = w^{(t)} - \eta_t v_t$$

where

$$(3.5) \quad v_t = \underbrace{\nabla \ell(w^{(t)}, z_{j(t)})}_{=x} - \underbrace{c}_{=c} \left(\underbrace{\nabla \ell(\tilde{w}, z_{j(t)})}_{=y} - \frac{1}{n} \sum_{i=1}^n \nabla \ell(\tilde{w}, z_i) \right)$$

(3) if modulo $(t, \kappa) = 0$, set $\tilde{w} = w^{(t)}$

(4) if a termination criterion is satisfied STOP

Note 24. Prepend, we do not know c in (3.3). The variance of gradient v_t^{SVRG} is

$$\begin{aligned} \mathbb{E}_{z \sim g} \|v_t^{\text{SVRG}} - \mathbb{E}(v_t^{\text{SVRG}})\|^2 &= \mathbb{E}_{z \sim g} \left\| \overbrace{\nabla_w \ell(w^{(t)}, z_{j(t)}) + c \left(\nabla_w \ell(\tilde{w}, z_{j(t)}) - \frac{1}{n} \sum_{i=1}^n \nabla \ell(\tilde{w}, z_i) \right)}^{v_t^{\text{SVRG}} =} - \nabla_w R_g(w^{(t)}) \right\|^2 \\ &\leq \mathbb{E}_{z \sim g} \left\| \nabla_w \ell(w^{(t)}, z_{j(t)}) + c \nabla_w \ell(\tilde{w}, z_{j(t)}) \right\|^2 \end{aligned}$$

As optimal c is often intractable, we sub-optimally set $c = -1$ and get

$$\mathbb{E}_{z \sim g} \|v_t^{\text{SVRG}} - \mathbb{E}(v_t^{\text{SVRG}})\|^2 \leq \mathbb{E}_{z \sim g} \left\| \nabla_w \ell(w^{(t)}, z_{j(t)}) - \nabla_w \ell(\tilde{w}, z_{j(t)}) \right\|^2$$

If $\ell(\cdot, z_j)$ is β_j -smooth, then

$$\mathbb{E}_{z \sim g} \|v_t^{\text{SVRG}} - \mathbb{E}(v_t^{\text{SVRG}})\|^2 \leq \max_j (\{\beta_j\}) \|w^{(t)} - \tilde{w}\|^2$$

Hence to bound/control variance, point \tilde{w} should be chosen to be as close as possible to $w^{(t)}$. Perhaps by taking a snapshot $\tilde{w} \leftarrow w^{(t)}$ every κ -th iteration.

Note 25. The frequency of the snapshots κ 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 in (3.5) and hence the bigger the performance improvement (See (3.3)); however the iterations will be slower.

Note 26. Iterations of SVRG are computationally faster than those of full GD but SVRG can still match the theoretical convergence rate of GD.

Example 27. ³ 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 $\mathcal{J}^{(t)} = \{j^{(t)}\}$ by drawing one number $j^{(t)}$ from $\{1, \dots, n = 10^6\}$.
- (2) compute

$$(3.6) \quad w^{(t+1)} = \begin{bmatrix} w_1^{(t)} \\ w_2^{(t)} \end{bmatrix} - \eta_t \left(\begin{bmatrix} 2(w_1^{(t)} - \tilde{w}_1) + 2(w_2^{(t)} - \tilde{w}_2)x_{j^{(t)}} \\ 2(w_2^{(t)} - \tilde{w}_2)x_{j^{(t)}} + 2(w_2^{(t)} - \tilde{w}_2)(x_{j^{(t)}})^2 \end{bmatrix} + \nabla \hat{R}_{\mathcal{S}}(\tilde{w}) \right)$$

- (3) if modulo $(t, \kappa) = 0$, set $\tilde{w} = w^{(t)}$
- (4) compute

$$(3.7) \quad \frac{1}{n} \sum_{i=1}^n \ell(\tilde{w}, z_i) = \begin{pmatrix} 2\tilde{w}_1 + 2\tilde{w}_2\bar{x} - 2\bar{y} \\ 2\tilde{w}_1\bar{x} + 2\tilde{w}_2\bar{x}^2 - 2\bar{y}^\top x \end{pmatrix} = \nabla \hat{R}_{\mathcal{S}}(\tilde{w})$$

- (5) if a termination criterion is satisfied STOP

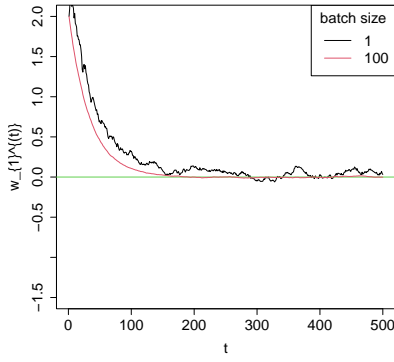
Because (3.7) is actually the gradient of the Risk function at \tilde{w} and

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

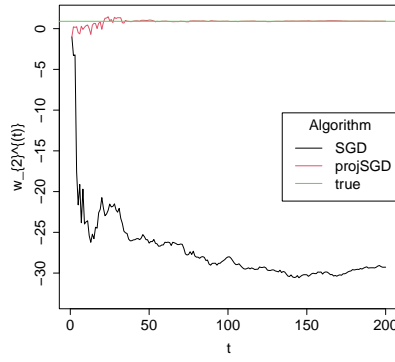
In Figure A.1c we observe the frequency of the snapshots has improved the convergence; however this is not a panacea as seen in Figure A.1f.

³https://github.com/georgios-stats/Machine_Learning_and_Neural_Networks_III_Epiphany_2025/tree/main/Lecture_notes/code/example_SVRG.R

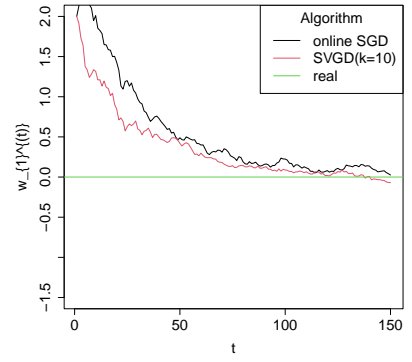
APPENDIX A. PLOTS



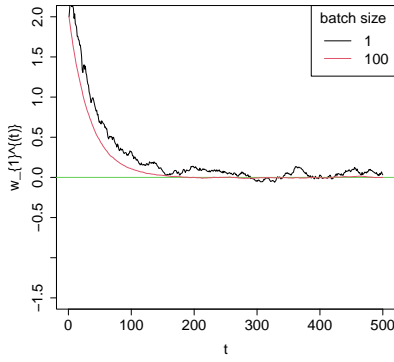
(A) batch SGD Example 30 w_1



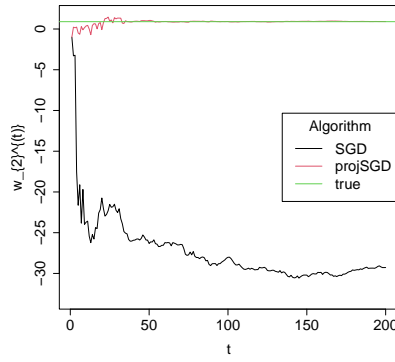
(B) online SGD with projection Example 6 w_1



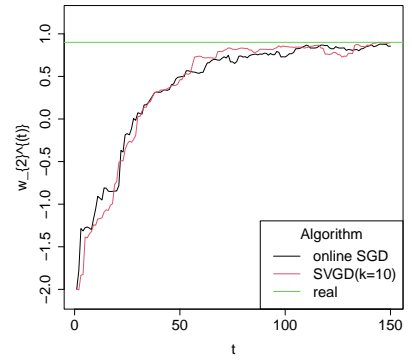
(C) SVRG Example 27 w_1



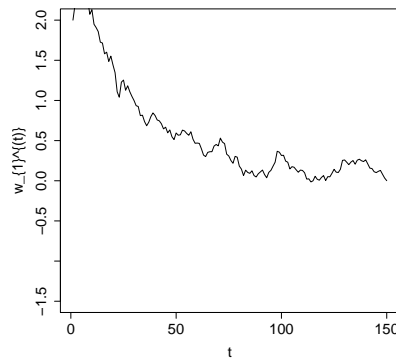
(D) batch SGD Example 30 w_2



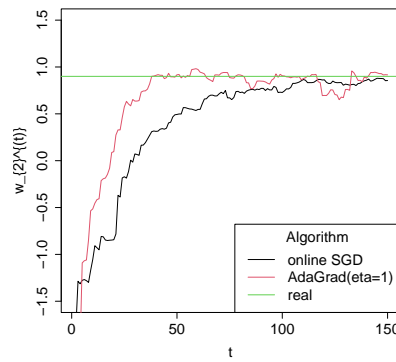
(E) online SGD with projection Example 6 w_2



(F) SVRG Example 27 w_2



(G) online AdaGrad Example 18 w_1



(H) online AdaGrad Example 18 w_2

FIGURE A.1. Simulations of the Examples