

# CMPT 409/981: Optimization for Machine Learning

## Lecture 10

---

Sharan Vaswani

October 20, 2022

# Minimizing smooth, strongly-convex functions using SGD

For smooth, strongly-convex functions, SGD with an  $O(1/k)$  decreasing step-size converges to the minimizer at an  $\Theta(1/T)$  rate (we will prove this later today).

Similar to the convex setting, using SGD with a constant step-size results in convergence to the neighbourhood that depends on the noise in the stochastic gradients.

**Claim:** For  $L$ -smooth,  $\mu$ -strongly convex functions,  $T$  iterations of SGD with  $\eta_k = \eta = \frac{1}{L}$  returns iterate  $w_T$  such that,

$$\mathbb{E}[\|w_T - w^*\|^2] \leq \exp\left(\frac{-T}{\kappa}\right) \|w_0 - w^*\|^2 + \frac{\sigma^2}{\mu L}$$

Hence, SGD results in an exponential convergence to the neighbourhood of the minimizer.

Unlike the convex case for which we proved a guarantee on the average iterate  $\bar{w}_T$ , here we have a guarantee for the last iterate  $w_T$ .

# Minimizing smooth, strongly-convex functions using SGD

**Proof:** Following a proof similar to the convex case,

$$\begin{aligned}\|w_{k+1} - w^*\|^2 &= \|w_k - \eta_k \nabla f_{i_k}(w_k) - w^*\|^2 \\ &= \|w_k - w^*\|^2 - 2\eta_k \langle \nabla f_{i_k}(w_k), w_k - w^* \rangle + \eta_k^2 \|\nabla f_{i_k}(w_k)\|^2\end{aligned}$$

Taking expectation w.r.t  $i_k$  on both sides,

$$\begin{aligned}\mathbb{E}[\|w_{k+1} - w^*\|^2] &= \|w_k - w^*\|^2 - 2\mathbb{E}[\eta_k \langle \nabla f_{i_k}(w_k), w_k - w^* \rangle] + \mathbb{E}[\eta_k^2 \|\nabla f_{i_k}(w_k)\|^2] \\ \implies \mathbb{E}[\|w_{k+1} - w^*\|^2] &= \|w_k - w^*\|^2 - 2\eta_k \langle \nabla f(w_k), w_k - w^* \rangle + \eta_k^2 \mathbb{E}[\|\nabla f_{i_k}(w_k)\|^2] \\ &\quad \text{(Assuming } \eta_k \text{ is independent of } i_k \text{ and Unbiasedness)}\end{aligned}$$

# Minimizing smooth, strongly-convex functions using SGD

Recall that  $\mathbb{E}[\|w_{k+1} - w^*\|^2] = \|w_k - w^*\|^2 - 2\eta_k \langle \nabla f(w_k), w_k - w^* \rangle + \eta_k^2 \mathbb{E}[\|\nabla f_{ik}(w_k)\|^2]$ .

$$\begin{aligned} & \mathbb{E}[\|w_{k+1} - w^*\|^2] \\ &= \|w_k - w^*\|^2 - 2\eta_k \langle \nabla f(w_k), w_k - w^* \rangle + \eta_k^2 \mathbb{E}[\|\nabla f_{ik}(w_k) - \nabla f(w_k) + \nabla f(w_k)\|^2] \\ &\leq \|w_k - w^*\|^2 - 2\eta_k \langle \nabla f(w_k), w_k - w^* \rangle + \eta_k^2 \mathbb{E}[\|\nabla f(w_k)\|^2] + \eta_k^2 \sigma^2 \\ & \hspace{15em} \text{(Using the bounded variance assumption)} \end{aligned}$$

Using  $\mu$ -strong convexity,  $f(y) \geq f(x) + \langle \nabla f(x), y - x \rangle + \frac{\mu}{2} \|y - x\|^2$  with  $y = w^*$  and  $x = w_k$ ,

$$\begin{aligned} & \leq \|w_k - w^*\|^2 - 2\eta_k [f(w_k) - f(w^*)] - \mu\eta_k \|w_k - w^*\|^2 + \eta_k^2 \mathbb{E}[\|\nabla f(w_k)\|^2] + \eta_k^2 \sigma^2 \\ & \hspace{15em} \text{(Eq. (1))} \end{aligned}$$

$$\begin{aligned} \implies & \mathbb{E}[\|w_{k+1} - w^*\|^2] \\ & \leq (1 - \mu\eta_k) \|w_k - w^*\|^2 - 2\eta_k [f(w_k) - f(w^*)] + 2L\eta_k^2 \mathbb{E}[f(w_k) - f(w^*)] + \eta_k^2 \sigma^2 \\ & \hspace{15em} \text{(Using } L\text{-smoothness of } f\text{)} \end{aligned}$$

# Minimizing smooth, strongly-convex functions using SGD

$$\mathbb{E}[\|w_{k+1} - w^*\|^2] \leq (1 - \mu\eta_k) \|w_k - w^*\|^2 - 2\eta_k[f(w_k) - f(w^*)] + 2L\eta_k^2 \mathbb{E}[f(w_k) - f(w^*)] + \eta_k^2 \sigma^2.$$

Setting  $\eta_k = \eta = \frac{1}{L}$

$$\mathbb{E}[\|w_{k+1} - w^*\|^2] \leq \left(1 - \frac{\mu}{L}\right) \|w_k - w^*\|^2 + \frac{\sigma^2}{L^2}$$

Since the above inequality is true for all  $k$ , using it for  $k = T - 1$ ,

$$\mathbb{E}[\|w_T - w^*\|^2] \leq \left(1 - \frac{\mu}{L}\right) \|w_{T-1} - w^*\|^2 + \frac{\sigma^2}{L^2}$$

Taking expectation w.r.t the randomness from iterations  $k = 0$  to  $T - 1$ ,

$$\implies \mathbb{E}[\|w_T - w^*\|^2] \leq \rho \mathbb{E} \|w_{T-1} - w^*\|^2 + \frac{\sigma^2}{L^2} \quad (\text{Denoting } \rho := 1 - \mu/L)$$

# Minimizing smooth, strongly-convex functions using SGD

Recall that  $\mathbb{E}[\|w_T - w^*\|^2] \leq \rho \mathbb{E} \|w_{T-1} - w^*\|^2 + \frac{\sigma^2}{L^2}$ . Unrolling the recursion until  $k = 0$ ,

$$\begin{aligned}\mathbb{E}[\|w_T - w^*\|^2] &\leq \rho^T \|w_0 - w^*\|^2 + \frac{\sigma^2}{L^2} \sum_{k=0}^{T-1} \rho^k \leq \rho^T \|w_0 - w^*\|^2 + \frac{\sigma^2}{L^2} \sum_{k=0}^{\infty} \rho^k \\ &\leq \rho^T \|w_0 - w^*\|^2 + \frac{\sigma^2}{L^2} \frac{1}{1 - \rho} \quad (\text{Infinite geometric series}) \\ &= \left(1 - \frac{\mu}{L}\right)^T \|w_0 - w^*\|^2 + \frac{\sigma^2}{\mu L} \\ &\leq \exp\left(\frac{-T}{\kappa}\right) \|w_0 - w^*\|^2 + \frac{\sigma^2}{\mu L} \quad (1 - x \leq \exp(-x)) \\ \Rightarrow \mathbb{E}[\|w_T - w^*\|^2] &\leq \underbrace{\exp\left(\frac{-T}{\kappa}\right) \|w_0 - w^*\|^2}_{\text{bias}} + \underbrace{\frac{\sigma^2}{\mu L}}_{\text{neighbourhood}}\end{aligned}$$

Questions?

# Minimizing smooth, strongly-convex functions using SGD

Let us prove that SGD with an  $O(1/k)$  step-size results in  $O(1/T)$  convergence to the minimizer.

**Claim:** For  $L$ -smooth,  $\mu$ -strongly convex functions,  $T$  iterations of SGD (for  $T \geq 2\kappa$ ) with  $\eta_k = \frac{1}{\mu(k+1)}$  returns iterate  $\bar{w}_T = \frac{\sum_{k=0}^{T-1} w_k}{T}$  such that,

$$\mathbb{E}[\|\bar{w}_T - w^*\|^2] \leq \frac{\sigma^2 [1 + \log(T)]}{\mu T}$$

Three problems – the above result (i) requires knowledge of  $\mu$ , (ii) the guarantee only holds for  $T \geq 2\kappa$ , (iii) the guarantee only holds for the average iterate and not the last iterate.

Instead of bounded variance, [LJSB12] assume that  $\mathbb{E}[\|\nabla f_i(w)\|^2] \leq G$ . Solves (ii) (not (i) and (iii)) and requires additional assumptions about the boundedness of iterates.

[GLQ<sup>+</sup>19, Theorem 3.2] uses a constant, then  $O(1/k)$  step-size. Solves (iii) (not (i) and (ii))

[LZO21, VDTB21] use an  $O((1/T)^{k/T})$  step-size and solves all three problems. Also prove a noise-adaptive  $O\left(\exp\left(\frac{-T}{\kappa}\right) + \frac{\sigma^2}{T}\right)$  rate, but requires knowledge of  $T$ .



# Minimizing smooth, strongly-convex functions using SGD

**Proof:** Following the previous proof, we can recover (Eq. (1) on Slide 3),

$$\begin{aligned} & \mathbb{E} \|w_{k+1} - w^*\|^2 \\ & \leq \|w_k - w^*\|^2 - 2\eta_k [f(w_k) - f(w^*)] - \mu\eta_k \|w_k - w^*\|^2 + \eta_k^2 \mathbb{E} [\|\nabla f(w_k)\|^2] + \eta_k^2 \sigma^2 \\ & \leq \|w_k - w^*\|^2 - 2\eta_k [f(w_k) - f(w^*)] - \mu\eta_k \|w_k - w^*\|^2 + 2L\eta_k^2 \mathbb{E} [f(w_k) - f(w^*)] + \eta_k^2 \sigma^2 \\ & \quad \text{(Using } L\text{-smoothness of } f\text{)} \\ & \leq \|w_k - w^*\|^2 - 2\eta_k [f(w_k) - f(w^*)] - \mu\eta_k \|w_k - w^*\|^2 + \eta_k \mathbb{E} [f(w_k) - f(w^*)] + \eta_k^2 \sigma^2 \\ & \quad \text{(Using } \eta_k \leq \frac{1}{2L} \text{ for all } k\text{)} \\ \implies \mathbb{E}[f(w_k) - f(w^*)] & \leq \frac{\left[ \|w_k - w^*\|^2 (1 - \mu\eta_k) - \mathbb{E} \|w_{k+1} - w^*\|^2 \right]}{\eta_k} + \eta_k \sigma^2 \end{aligned}$$

Taking expectation w.r.t the randomness from iterations  $k = 0$  to  $T - 1$ ,

$$\mathbb{E}[f(w_k) - f(w^*)] \leq \frac{\mathbb{E} \left[ \|w_k - w^*\|^2 (1 - \mu\eta_k) - \|w_{k+1} - w^*\|^2 \right]}{\eta_k} + \eta_k \sigma^2$$

# Minimizing smooth, strongly-convex functions using SGD

Recall that  $\mathbb{E}[f(w_k) - f(w^*)] \leq \frac{\mathbb{E}[\|w_k - w^*\|^2 (1 - \mu \eta_k) - \|w_{k+1} - w^*\|^2]}{\eta_k} + \eta_k \sigma^2$ .

Summing from  $k = 0$  to  $T - 1$ ,

$$\begin{aligned} \sum_{k=0}^{T-1} \mathbb{E}[f(w_k) - f(w^*)] &\leq \sum_{k=0}^{T-1} \frac{\mathbb{E}[\|w_k - w^*\|^2 (1 - \mu \eta_k) - \|w_{k+1} - w^*\|^2]}{\eta_k} + \sigma^2 \sum_{k=0}^{T-1} \eta_k \\ &= \sum_{k=0}^{T-1} \frac{\mathbb{E}[\|w_k - w^*\|^2 (1 - \mu \eta_k) - \|w_{k+1} - w^*\|^2]}{\eta_k} + \sigma^2 \sum_{k=0}^{T-1} \frac{1}{\mu (k+1)} \\ &\leq \sum_{k=0}^{T-1} \frac{\mathbb{E}[\|w_k - w^*\|^2 (1 - \mu \eta_k) - \|w_{k+1} - w^*\|^2]}{\eta_k} + \frac{\sigma^2 [1 + \log(T)]}{\mu} \end{aligned}$$

Dividing by  $T$ , using Jensen's inequality for the LHS, and by definition of  $\bar{w}_T$ ,

$$\mathbb{E}[f(\bar{w}_T) - f(w^*)] \leq \frac{1}{T} \sum_{k=0}^{T-1} \frac{\mathbb{E}[\|w_k - w^*\|^2 (1 - \mu \eta_k) - \|w_{k+1} - w^*\|^2]}{\eta_k} + \frac{\sigma^2 [1 + \log(T)]}{\mu T}$$

# Minimizing smooth, strongly-convex functions using SGD

Recall that  $\mathbb{E}[f(\bar{w}_T) - f(w^*)] \leq \frac{1}{T} \sum_{k=0}^{T-1} \frac{\mathbb{E}[\|w_k - w^*\|^2 (1 - \mu \eta_k) - \|w_{k+1} - w^*\|^2]}{\eta_k} + \frac{\sigma^2 [1 + \log(T)]}{\mu T}$ .

$$\begin{aligned} & \sum_{k=0}^{T-1} \frac{\mathbb{E}[\|w_k - w^*\|^2 (1 - \mu \eta_k) - \|w_{k+1} - w^*\|^2]}{\eta_k} \\ &= \mathbb{E} \sum_{k=1}^{T-1} \left[ \|w_k - w^*\|^2 \left( \frac{1}{\eta_k} - \frac{1}{\eta_{k-1}} - \mu \right) \right] + \|w_0 - w^*\|^2 \left( \frac{1}{\eta_0} - \mu \right) - \frac{\|w_T - w^*\|^2}{\eta_{T-1}} \\ &\leq \mathbb{E} \sum_{k=1}^{T-1} \left[ \|w_k - w^*\|^2 (\mu(k+1) - \mu k - \mu) \right] + \|w_0 - w^*\|^2 (\mu - \mu) = 0 \end{aligned}$$

Putting everything together,

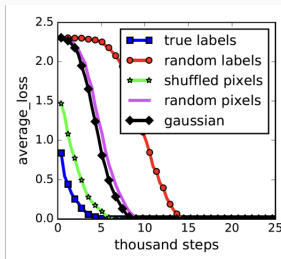
$$\mathbb{E}[f(\bar{w}_T) - f(w^*)] \leq \frac{\sigma^2 [1 + \log(T)]}{\mu T}$$

Since we used the fact that  $\eta_k \leq \frac{1}{2L}$  for all  $k$ , it implies that  $\frac{1}{\mu T} \leq \frac{1}{2L} \implies T \geq 2\kappa$ .

Questions?

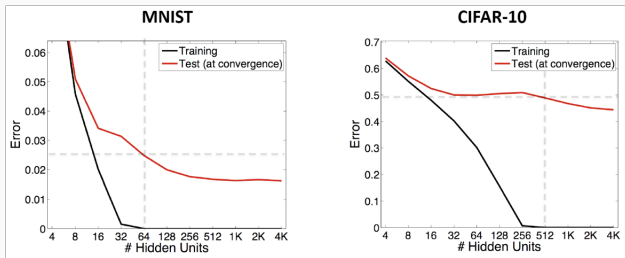
# Interpolation for over-parameterized models

**Interpolation:** Over-parameterized models (such as deep neural networks) are capable of exactly fitting the training dataset.



Zhang et al, "Understanding deep learning requires rethinking generalization", 2016.

Loss vs Training steps on CIFAR-10 dataset



[https://www.neyshabur.net/papers/inductive\\_bias\\_poster.pdf](https://www.neyshabur.net/papers/inductive_bias_poster.pdf)

Error vs Network size

Formally, when minimizing  $f(w) = \frac{1}{n} \sum_{i=1}^n f_i(w)$ , interpolation means that if  $\|\nabla f(w)\| = 0$ , then  $\|\nabla f_i(w)\| = 0$  for all  $i \in [n]$  i.e. the variance in the stochastic gradients becomes zero at a stationary point.

# SGD under Interpolation

Recall that SGD needs to decrease the step-size to counteract the noise (variance).

**Idea:** Under interpolation, since the noise is zero at the optimum, SGD does not need to decrease the step-size and can converge to the minimizer by using a *constant* step-size.

If  $f$  is strongly-convex and the model is expressive enough such that interpolation is satisfied (for example, when using kernels or least squares with  $d > n$ ), constant step-size SGD can converge to the minimizer at an  $O(\exp(-T/\kappa))$  rate.

In this setting, SGD matches the rate of deterministic (full-batch) GD, but compared to GD, each iteration is cheap.

Moreover, empirical results (and theoretical results on “benign overfitting”) suggest that interpolating the training dataset does not adversely affect the generalization error!

# Minimizing smooth, strongly-convex functions using SGD under interpolation

**Claim:** When minimizing  $f(w) = \frac{1}{n} \sum_{i=1}^n f_i(w)$  such that (i)  $f$  is  $\mu$ -strongly convex, (ii) each  $f_i$  is convex and  $L$ -smooth, (iii) interpolation is exactly satisfied i.e.  $\|\nabla f_i(w^*)\| = 0$ ,  $T$  iterations of SGD with  $\eta_k = \eta = \frac{1}{2L}$  returns iterate  $w_T$  such that,

$$\mathbb{E}[\|w_T - w^*\|^2] \leq \exp\left(\frac{-T}{\kappa}\right) \|w_0 - w^*\|^2$$

Before analyzing the convergence of SGD, let us first study the effect of interpolation on  $\sigma^2(w)$ .

$$\begin{aligned}\sigma^2(w) &:= \mathbb{E}_i \|\nabla f(w) - \nabla f_i(w)\|^2 = \|\nabla f(w)\|^2 + \mathbb{E}_i \|\nabla f_i(w)\|^2 - 2\mathbb{E}[\langle \nabla f(w), \nabla f_i(w) \rangle] \\ &= \mathbb{E}_i \|\nabla f_i(w)\|^2 + \|\nabla f(w)\|^2 - 2\|\nabla f(w)\|^2 \quad (\text{Unbiasedness}) \\ &\leq \mathbb{E}_i \|\nabla f_i(w)\|^2 \leq \mathbb{E}_i [2L[f_i(w) - f_i(w^*)]] \\ &\quad (\text{Using } L\text{-smoothness, convexity of } f_i \text{ and } \nabla f_i(w^*) = 0)\end{aligned}$$

$$\implies \sigma^2(w) \leq 2L[f(w) - f(w^*)] \quad (\text{Unbiasedness})$$

As  $w$  gets closer to the solution (in terms of the function values), the variance decreases becoming zero at  $w^*$ . Hence, under interpolation, we do not need to decrease the step-size.

# Minimizing smooth, strongly-convex functions using SGD under interpolation

**Proof:** Following the same proof as before, we get that,

$$\begin{aligned}\mathbb{E}[\|w_{k+1} - w^*\|^2] &= \|w_k - w^*\|^2 - 2\eta_k \langle \nabla f(w_k), w_k - w^* \rangle + \eta_k^2 \mathbb{E} [\|\nabla f_{ik}(w_k)\|^2] \\ &\leq \|w_k - w^*\|^2 - 2\eta_k \langle \nabla f(w_k), w_k - w^* \rangle + \eta_k^2 \mathbb{E}_i [2L [f_{ik}(w_k) - f_{ik}(w^*)]] \\ &\quad \text{(Using } L\text{-smoothness, convexity of } f_i \text{ and } \nabla f_i(w^*) = 0\text{)} \\ &= \|w_k - w^*\|^2 - 2\eta_k \langle \nabla f(w_k), w_k - w^* \rangle + 2L \eta_k^2 \mathbb{E} [f(w_k) - f(w^*)] \\ &\quad \text{(Unbiasedness)} \\ &= \|w_k - w^*\|^2 (1 - \mu\eta_k) - 2\eta_k [f(w_k) - f(w^*)] + 2L \eta_k^2 \mathbb{E} [f(w_k) - f(w^*)] \\ &\quad \text{(Strong-convexity)} \\ &= \left(1 - \frac{\mu}{2L}\right) \|w_k - w^*\|^2 \quad \text{(Since } \eta_k = \eta = \frac{1}{2L}\text{)}\end{aligned}$$

Taking expectation w.r.t the randomness from iterations  $k = 0$  to  $T - 1$  and recursing,

$$\mathbb{E}[\|w_T - w^*\|^2] \leq \left(1 - \frac{\mu}{2L}\right)^T \|w_0 - w^*\|^2 \leq \exp\left(\frac{-T}{\kappa}\right) \|w_0 - w^*\|^2$$



# Minimizing smooth, strongly-convex functions using SGD under interpolation

We can modify the proof in order to get an  $O\left(\exp\left(\frac{-T}{\kappa}\right) + \zeta^2\right)$  where  $\zeta^2 \propto \mathbb{E}_i \|\nabla f_i(w^*)\|^2$ .

Moreover, as before, if we use a mini-batch of size  $b$ , the effective noise is  $\zeta_b^2 \propto \frac{\mathbb{E}_i \|\nabla f_i(w^*)\|^2}{b}$ .

Hence, if the model is sufficiently over-parameterized so that it *almost* interpolates the data, and we are using a large batch-size, then  $\zeta_b^2$  is small, and constant step-size works well.

When minimizing convex functions under (exact) interpolation, constant step-size SGD results in  $O(1/T)$  convergence, matching deterministic GD, but with much smaller per-iteration cost (Need to prove this in Assignment 3!)

Questions?

# Minimizing smooth, non-convex functions using SGD under interpolation

When minimizing non-convex functions, interpolation is not enough to guarantee a fast (matching the deterministic)  $O(1/T)$  rate for SGD.

Can achieve this rate under the *strong growth condition* (SGC) on the stochastic gradients. Formally, there exists a constant  $\rho > 1$  such that for all  $w$ ,

$$\mathbb{E}_i \|\nabla f_i(w)\|^2 \leq \rho \|\nabla f(w)\|^2$$

Hence, SGC implies that  $\|\nabla f_i(w^*)\|^2 = 0$  for all  $i$  and hence interpolation.

As before, let us study the effect of SGC on the variance  $\sigma^2(w)$ .

$$\begin{aligned} \sigma^2(w) &:= \mathbb{E}_i \|\nabla f_i(w) - \nabla f(w)\|^2 \leq \mathbb{E}_i \|\nabla f_i(w)\|^2 - \|\nabla f(w)\|^2 && \text{(Unbiasedness)} \\ \implies \sigma^2(w) &\leq (\rho - 1) \|\nabla f(w)\|^2 && \text{(SGC)} \end{aligned}$$

Hence, SGC implies that as  $w$  gets closer to a stationary point (in terms of the gradient norm), the variance decreases and constant step-size SGD converges to a stationary point.

# Minimizing smooth, non-convex functions using SGD under interpolation

**Claim:** For (i)  $L$ -smooth functions lower-bounded by  $f^*$ , (ii) under  $\rho$ -SGC,  $T$  iterations of SGD with  $\eta_k = \frac{1}{\rho L}$  returns an iterate  $\hat{w}$  such that,

$$\mathbb{E}[\|\nabla f(\hat{w})\|^2] \leq \frac{2\rho L [f(w_0) - f^*]}{T}$$

**Proof:** Similar to the proof in Lecture 8, using the  $L$ -smoothness of  $f$  with  $x = w_k$  and  $y = w_{k+1} = w_k - \eta_k \nabla f_{ik}(w_k)$ ,

$$f(w_{k+1}) \leq f(w_k) + \langle \nabla f(w_k), -\eta_k \nabla f_{ik}(w_k) \rangle + \frac{L}{2} \eta_k^2 \|\nabla f_{ik}(w_k)\|^2$$

Taking expectation w.r.t  $i_k$  on both sides and using that  $\eta_k$  is independent of  $i_k$

$$\begin{aligned} \mathbb{E}[f(w_{k+1})] &\leq f(w_k) - \eta_k \mathbb{E}[\langle \nabla f(w_k), \nabla f_{ik}(w_k) \rangle] + \frac{L\eta_k^2}{2} \mathbb{E}[\|\nabla f_{ik}(w_k)\|^2] \\ \mathbb{E}[f(w_{k+1})] &\leq f(w_k) - \eta_k \|\nabla f(w_k)\|^2 + \frac{L\eta_k^2}{2} \mathbb{E}[\|\nabla f_{ik}(w_k)\|^2] \quad (\text{Unbiasedness}) \end{aligned}$$

# Minimizing smooth, non-convex functions using SGD under interpolation

Recall  $\mathbb{E}[f(w_{k+1})] \leq f(w_k) - \eta_k \|\nabla f(w_k)\|^2 + \frac{L\eta_k^2}{2} \mathbb{E}[\|\nabla f_{ik}(w_k)\|^2]$ . Using  $\rho$ -SGC,

$$\mathbb{E}[f(w_{k+1})] \leq f(w_k) - \eta_k \|\nabla f(w_k)\|^2 + \frac{L\rho\eta_k^2}{2} \|\nabla f(w_k)\|^2$$

$$\mathbb{E}[f(w_{k+1})] \leq f(w_k) - \frac{1}{2\rho L} \|\nabla f(w_k)\|^2 \quad (\text{Using } \eta_k = \eta = \frac{1}{\rho L})$$

Taking expectation w.r.t the randomness from iterations  $i = 0$  to  $k - 1$ , and summing





$$\sum_{k=0}^{T-1} \mathbb{E}[\|\nabla f(w_k)\|^2] \leq 2\rho L \sum_{k=0}^{T-1} \mathbb{E}[f(w_k) - f(w_{k+1})] \implies \frac{\sum_{k=0}^{T-1} \mathbb{E}[\|\nabla f(w_k)\|^2]}{T} \leq \frac{2\rho L \mathbb{E}[f(w_0) - f^*]}{T}$$

(Dividing by  $T$ )

Defining  $\hat{w} := \arg \min_{k \in \{0, 1, \dots, T-1\}} \mathbb{E}[\|\nabla f(w_k)\|^2]$ ,

$$\mathbb{E}[\|\nabla f(\hat{w})\|^2] \leq \frac{2\rho L [f(w_0) - f^*]}{T}$$

Questions?

-  Robert Mansel Gower, Nicolas Loizou, Xun Qian, Alibek Sailanbayev, Egor Shulgin, and Peter Richtárik, *Sgd: General analysis and improved rates*, International Conference on Machine Learning, PMLR, 2019, pp. 5200–5209.
-  Simon Lacoste-Julien, Mark Schmidt, and Francis Bach, *A simpler approach to obtaining an  $o(1/t)$  convergence rate for the projected stochastic subgradient method*, arXiv preprint arXiv:1212.2002 (2012).
-  Xiaoyu Li, Zhenxun Zhuang, and Francesco Orabona, *A second look at exponential and cosine step sizes: Simplicity, adaptivity, and performance*, International Conference on Machine Learning, PMLR, 2021, pp. 6553–6564.
-  Sharan Vaswani, Benjamin Dubois-Taine, and Reza Babanezhad, *Towards noise-adaptive, problem-adaptive stochastic gradient descent*, arXiv preprint arXiv:2110.11442 (2021).