Expectation on a Single Iteration I

• Recall earlier we derived

$$\sum_{t=0}^{T-1} \alpha_t E[\|\nabla f(\boldsymbol{\theta}_t)\|^2]$$

$$\leq f(\boldsymbol{\theta}_0) - f^* + \sum_{t=0}^{T-1} \frac{\alpha_t^2 G^2 L}{2}, \quad (1)$$

- Let's get a better representation by taking the expectation on a single iteration
- Assume that we run up to a random number of iteration τ

Expectation on a Single Iteration II

• And τ follows the following probability distribution

$$P(\tau = t) = \frac{\alpha_t}{\sum_{k=0}^{T-1} \alpha_k}$$

• The expected squared-norm of the gradient is

$$E_{\tau,\tilde{i}_{0},...,\tilde{i}_{\tau-1}}[\|\nabla f(\boldsymbol{\theta}_{\tau})\|^{2}]$$

$$= \sum_{t=0}^{T-1} E_{\tilde{i}_{0},...,\tilde{i}_{t-1}}[\|\nabla f(\boldsymbol{\theta}_{t})\|^{2}]P(\tau = t)$$

$$= \left(\sum_{k=0}^{T-1} \alpha_{k}\right)^{-1} \sum_{t=0}^{T-1} \alpha_{t} E[\|\nabla f(\boldsymbol{\theta}_{t})\|^{2}] \quad (2)$$

Expectation on a Single Iteration III

• From (2) and (1),

$$E[\|\nabla f(\boldsymbol{\theta}_{\tau})\|^{2}] \leq \left(\sum_{k=0}^{T-1} \alpha_{k}\right)^{-1} \left(f(\boldsymbol{\theta}_{0}) - f^{*} + \frac{G^{2}L}{2} \sum_{t=0}^{T-1} \alpha_{t}^{2}\right) (3)$$

< □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > <

Expectation on a Single Iteration IV

• If the learning rate is a constant $\alpha_t = \alpha$, then

$$E[\|\nabla f(\boldsymbol{\theta}_{\tau})\|^{2}] \leq (T\alpha)^{-1} \left(f(\boldsymbol{\theta}_{0}) - f^{*} + \frac{G^{2}L}{2}T\alpha^{2} \right)$$
$$= \frac{f(\boldsymbol{\theta}_{0}) - f^{*}}{\alpha T} + \frac{\alpha L G^{2}}{2}$$

• The right-hand side does not go to zero as $T \to \infty$ due to the term

$$\frac{\alpha LG^2}{2}$$

(日)

Reducing the Learning Rate I

- We can make it zero by decreasing the learning rate over time
- To make the right-hand side of (3) go to zero, we need

$$\sum_{t=0}^{T-1} \alpha_t \text{ grows much faster than } \sum_{t=0}^{T-1} \alpha_t^2$$

• An example is

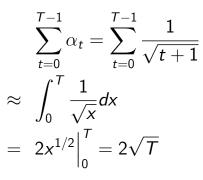
$$\alpha_t = \frac{1}{\sqrt{t+1}}$$

1

< ロ > < 同 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 > < 0 >

Reducing the Learning Rate II

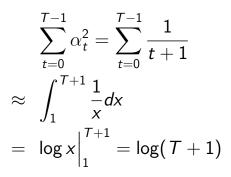
We have



< □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > <

Reducing the Learning Rate III

and



< □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > <

Reducing the Learning Rate IV

• Now (3) becomes

$$egin{aligned} & E[\|
abla f(m{ heta}_{ au})\|^2] \ &\lesssim \ (2\sqrt{T})^{-1}\left(f(m{ heta}_0)-f^*+rac{G^2L}{2}(\log(T+1))
ight) \ &= \ rac{2(f(m{ heta}_0)-f^*)+G^2L\log(T+1)}{4\sqrt{T}} = O(rac{\log T}{\sqrt{T}}) \end{aligned}$$

• The value goes to zero as T increases

イロト 不得 トイヨト イヨト