

Lecture 2

HWO is due
10:59pm Friday.

Optimization

- ① Implicit Regularization of Gradient Descent
- ② Stochastic Gradient Descent
- ③ Momentum

Last time: Gradient Descent converges for linear regression.

$$\eta < \frac{1}{\lambda_{\max}(X^T X)}$$

Key takeaway: Geometric properties of the loss matter, and dictate choice of learning rate.

Regularization + Ridge : Why all the singular values + spread of singular values matter.

$$\vec{w}_* = (X^T X + \lambda I)^{-1} X^T \vec{y} = X^T (X X^T + \lambda I)^{-1} \cdot \vec{y}$$

$$X = U \Sigma V^T$$

$$\vec{w}_* = (V \Sigma^T U^T (U \Sigma V^T + \lambda I)^{-1} V \Sigma U^T \cdot \vec{y})$$

$$= V (\Sigma^T \Sigma + \lambda I)^{-1} V^T \Sigma^T U^T \vec{y}$$

$$= V \underbrace{(\Sigma^T \Sigma + \lambda I)^{-1}}_{\text{Diagonal matrix}} \Sigma^T U^T \vec{y}.$$

$$= V \begin{bmatrix} \sigma_1^2 + \lambda & & & & \\ & \sigma_2^2 + \lambda & & & \\ & & \ddots & & \\ & & & 0 & \\ & & & & \ddots & \lambda & \lambda & \lambda \end{bmatrix}^{-1} \Sigma^T U^T \vec{y}.$$

$$\vec{w}_* = \sum_{i=1}^n \vec{u}_i \left[\frac{\sigma_i}{\sigma_i^2 + \lambda} \right] \vec{u}_i^T \cdot \vec{y}$$

$\frac{\sigma_i}{\sigma_i^2 + \lambda}$ if $\lambda \gg \sigma_i$
 $\rightarrow 0$ as $\lambda \rightarrow \infty$

$$\frac{\sigma_i}{\sigma_i^2 + \lambda} \rightarrow \frac{\sigma_i}{\sigma_i^2} \text{ if } \sigma_i \gg \lambda$$

GD update: $\vec{w}_{t+1} = \vec{w}_t - 2\eta \underbrace{x^T(x\vec{w}_t - \vec{y})}_{\text{residual}}$.

always remains in the span of the data.

GD update for ridge

$$\vec{w}_{t+1} = (1 - 2\eta\lambda) \vec{w}_t - 2\eta \cdot (x^T) (x\vec{w}_t - \vec{y})$$

weight decay.

Implicit Regularization

$$V^T \vec{w}_{t+1} = V^T \vec{w}_t - 2\eta V^T (V \Sigma^T V^T) (x\vec{w}_t - \vec{y})$$

$$\begin{aligned} \vec{\tilde{w}}_{t+1} &= \vec{\tilde{w}}_t - 2\eta \Sigma^T (U^T \Sigma V^T \vec{w}_t - U^T \vec{y}) \\ &= \vec{\tilde{w}}_t - 2\eta \Sigma^T (\Sigma \vec{\tilde{w}}_t - U^T \vec{y}) \end{aligned}$$

$$\begin{aligned} \tilde{w}_{t+1,i} &= \tilde{w}_{t,i} - 2\eta \cdot [\sigma_i^{-2} \tilde{w}_{t,i} - \sigma_i^{-1} \tilde{y}_i] \\ &= \tilde{w}_{t,i} - 2\eta \cdot \sigma_i^{-1} [\sigma_i^{-2} \tilde{w}_{t,i} - \tilde{y}_i] \end{aligned}$$

"Early Stopping"

Implicit Regularization

Ridge $\rightarrow \lambda$ is large \rightarrow we don't have large components
of directions of small singular values in
Ridge solution.

GD \rightarrow early stopping \rightarrow also does not move in
directions of small singular values.
or moves very little

Stochastic Gradient Descent

- GD is great for convex functions.
- Neural network fitting requires optimization over non-convex landscapes
- GD is also expensive.

SGD

- Adds some noise to gradient direction. In expectation, still descends.

$\min_{\theta} f(\theta)$

$$\theta_{t+1} = \theta_t - \eta \cdot \frac{1}{B} \sum_{i=1}^B \nabla f_i(\theta_t)$$

where $f(\theta) = \frac{1}{n} \sum_{i=1}^n f_i(\theta)$

$$E[\nabla f_i(\theta)] = \frac{1}{n} \cdot \sum_{i=1}^n \nabla f_i(\theta) = \nabla f(\theta)$$

draw i uniformly iid

B: batch size
"mini-batch"

Gradient descent is
"full-batch"

- Many standard loss functions have this form.

e.g. squared loss functions:

$$\begin{aligned} \operatorname{argmin}_{\vec{x}} \|\vec{A}\vec{x} - \vec{y}\|_2^2 &= \operatorname{argmin}_{\vec{x}} \frac{1}{m} \|\vec{A}\vec{x} - \vec{y}\|_2^2 = \\ &= \operatorname{argmin}_{\vec{x}} \frac{1}{m} \sum_{i=1}^m (\vec{q}_i^\top \vec{x} - y_i)^2. \end{aligned}$$

Loss on individual data point.

- Though noisy, improves fit to subset of data each step.
- Computationally less expensive
- Noise helps us (can help us) escape local minima and saddle points.

Question: Does it converge?

Optimization : Must use a learning-rate schedule with decreasing
step size for SGD

In practice: can see it converge with constant step size.

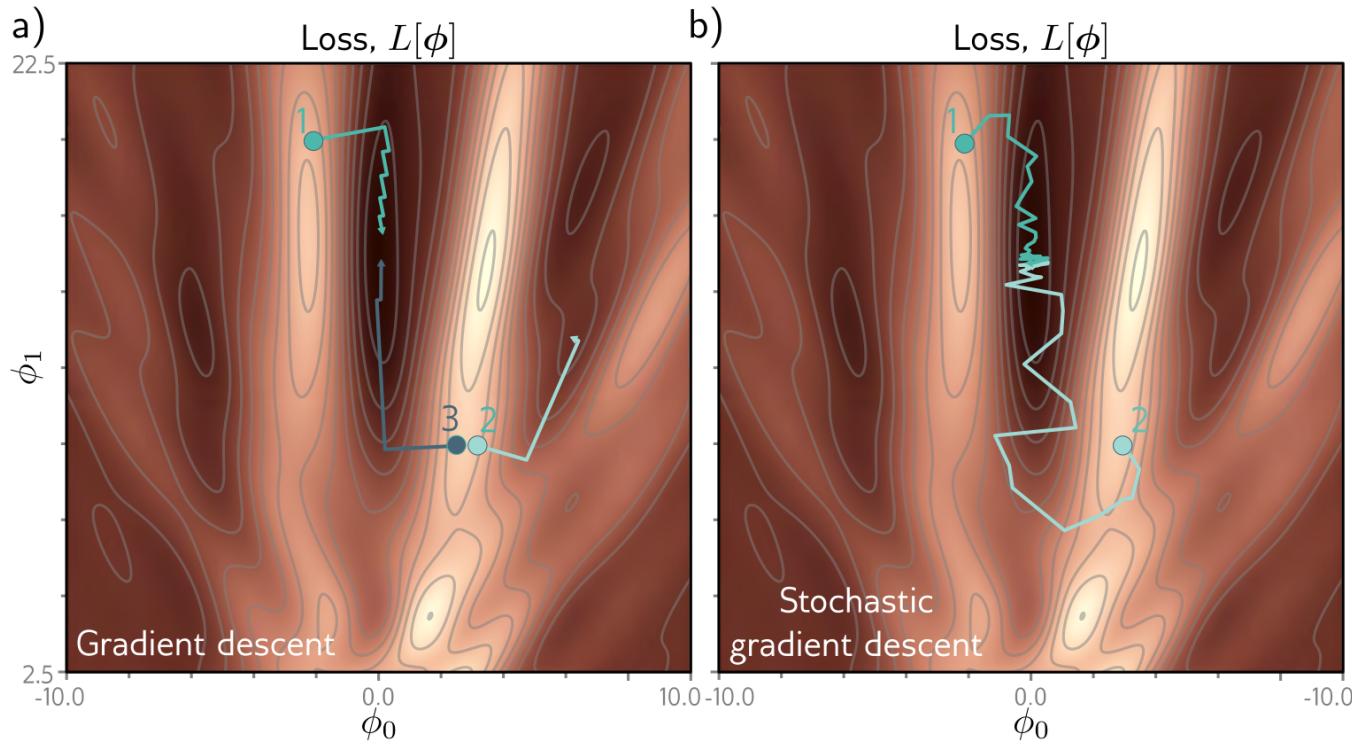


Figure 6.5 Gradient descent vs. stochastic gradient descent. a) Gradient descent with line search. As long as the gradient descent algorithm is initialized in the right “valley” of the loss function (e.g., points 1 and 3), the parameter estimate will move steadily toward the global minimum. However, if it is initialized outside this valley (e.g., point 2), it will descend toward one of the local minima. b) Stochastic gradient descent adds noise to the optimization process, so it is possible to escape from the wrong valley (e.g., point 2) and still reach the global minimum.

Prince 2025.

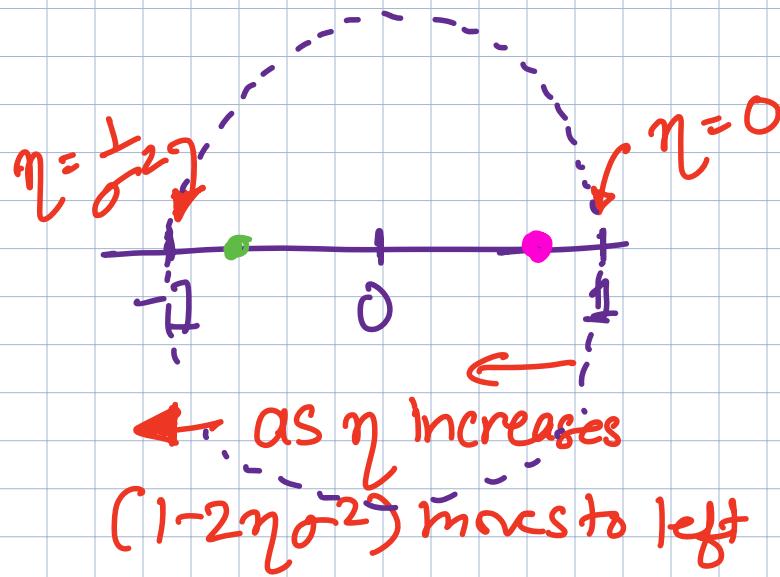
Momentum:

Gradient Descent.

Scalar case: $\sigma w = y$
 (i.e. on direction of SVD basis)

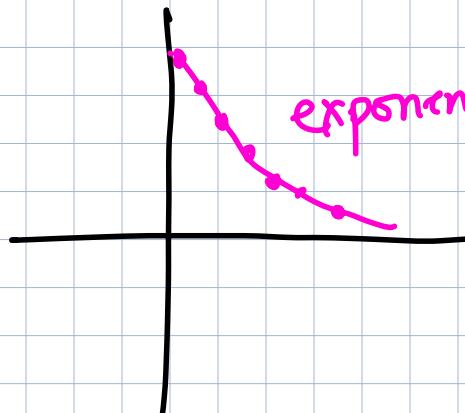
$$w_{t+1} = w_t + \eta \cdot 2\sigma(y - \sigma w_t) \\ = (1 - 2\eta\sigma^2)w_t + 2\eta\cdot\sigma\cdot y.$$

$$\Rightarrow \eta < \frac{1}{\sigma^2}$$

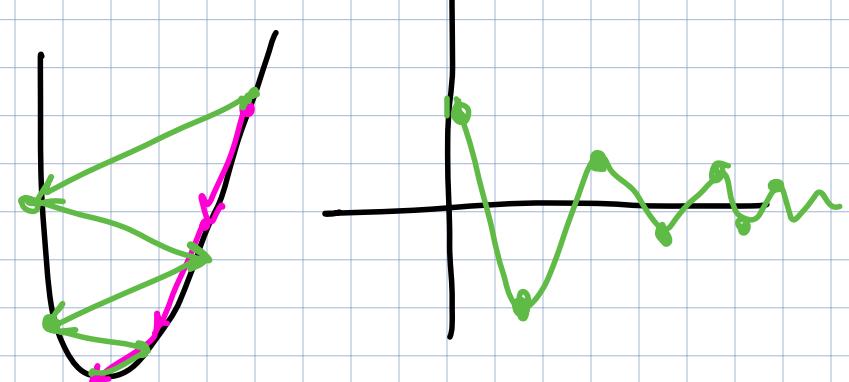


i.e. $\min_w \|y - \sigma w\|_2^2 = f(w)$
 $\frac{d}{dw} f(w) = -2\sigma(y - \sigma w)$

GD learning rate gated by max singular value.
 So smaller SVs may still converge slowly.



exponential decay



So how can we smooth out the oscillations?

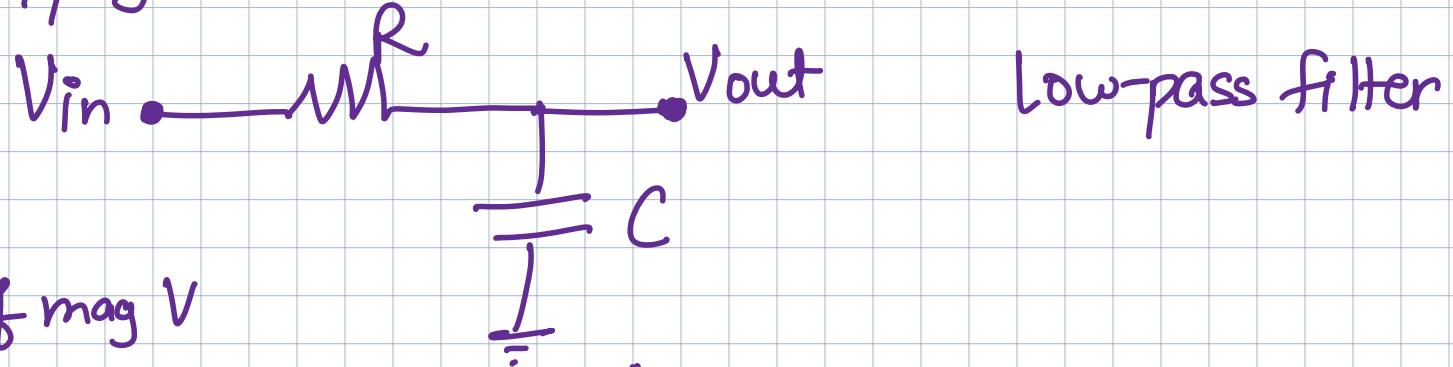
→ Averaging , Low-pass filtering

One strategy : $\frac{1}{t} \sum_{i=1}^t \nabla f(w_i) \rightarrow \text{straight average.}$

→ Requires lots of memory ... have to remember all the past gradients.

→ In physics) state-space systems
a state can be used for averaging.

A circuit /physics Intuition:



If V_{in} = step of mag V

$$V_{out} = V(1 - e^{-\frac{1}{RC}t})$$

Momentum : Another tweak on Gradient Descent. || How to make optimization go faster?
 Add memory.

$$\vec{w}_{k+1} = \vec{w}_k - \eta \cdot \vec{z}_{k+1}$$

$$\underline{\vec{z}_{k+1}} = \beta \cdot \vec{z}_k + (1-\beta) \nabla f(\vec{w}_k), z_0 = 0. \quad \beta \in [0, 1]$$

If $\beta = 0$, $z_{k+1} = \nabla f(w_k)$. i.e same old gradient descent -

$$z_1 = (1-\beta) \nabla f(w_0)$$

$$z_2 = \beta \cdot (1-\beta) \cdot \nabla f(w_0) + (1-\beta) \nabla f(w_1).$$

$$z_3 = \beta^2 \cdot (1-\beta) \cdot \nabla f(w_0) + \beta \cdot (1-\beta) \nabla f(w_1) + (1-\beta) \nabla f(w_2).$$

:

Add history of the gradients. Exponentially decaying weight on history.

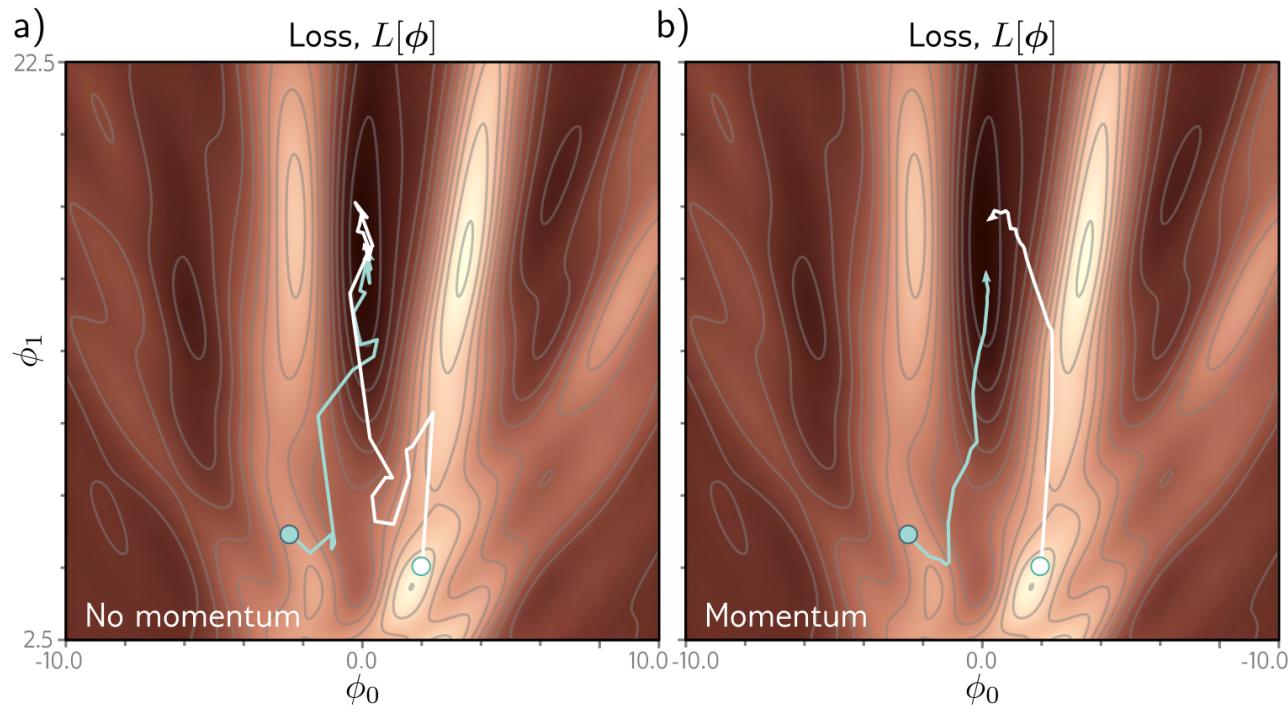


Figure 6.7 Stochastic gradient descent with momentum. a) Regular stochastic gradient descent takes a very indirect path toward the minimum. b) With a momentum term, the change at the current step is a weighted combination of the previous change and the gradient computed from the batch. This smooths out the trajectory and increases the speed of convergence.

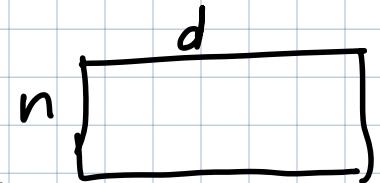
SGD Convergence

SGD can provably converge, if loss can go to zero even with constant stepsize

$$X\vec{w} = \vec{y} . \quad X \in \mathbb{R}^{n \times d} \\ d > n .$$

full-row rank.

under-determined.



Let \vec{w}_* be the min-norm solution.

$$\vec{w}_* = X^T (X X^T)^{-1} \vec{y} .$$

We do SGD, batch size 1.

$$\vec{w}_0 = 0 .$$

error

$$X(\vec{w} - \vec{w}^*) = \vec{y} - \vec{y} = 0$$

\vec{q}

$X\vec{q} = 0$ is equivalent equation with

$$\vec{q}_0 = \vec{w}_0 - \vec{w}_* \\ = -\vec{w}_*$$

We know from the SVD that only directions with non-zero singular values can be moved by SGD/GD.

$$\begin{aligned}
 X\vec{q} &= U\Sigma V^T \vec{q} \\
 &= U \left[\sum_{i=1}^{n \times n} \begin{matrix} \circ \\ \circ \end{matrix} \right] V^T \vec{q} \\
 &= \begin{bmatrix} \tilde{X} & 0 \end{bmatrix} \vec{z} \quad \begin{aligned} V^T \vec{q} &= \vec{z} \\ U\Sigma_1 &= \tilde{X} \end{aligned}
 \end{aligned}$$

$$\begin{aligned}
 z_0 &= V^T q_0 = V^T (-\vec{w}_x) \\
 &= - \begin{bmatrix} \sum_{i=1}^{d-n} V^T \vec{y} \\ 0 \end{bmatrix} \} \text{d-n zero entries.}
 \end{aligned}$$

So what matters is

$$\underbrace{\tilde{X} \vec{z}}_{n \times n} = 0$$

ith row of this is the same as ith row of

$$\tilde{X} \vec{q}$$

(we have only dropped zeros)

Now GD.

$$\vec{z}_{t+1} = \vec{z}_t - 2\gamma \cdot \tilde{X}^T \tilde{X} z_t$$

$$\text{loss} = \|\tilde{X} \vec{z}\|_2^2$$

$$\nabla \text{loss} = 2 \tilde{X}^T \tilde{X} z$$

$$\begin{aligned} \tilde{X}^T \tilde{X} &= \Sigma_1^T U^T U \Sigma_1 \\ &= \Sigma_1^T \Sigma_1 \end{aligned}$$

For SGD with batch size 1, this is.

$$\vec{z}_{t+1} = \vec{z}_t - 2\gamma \cdot \tilde{X}^T \tilde{X} z_t$$



point sampled at time I_t .

For this SGD update can we argue that

$$\text{loss} = \|\tilde{\vec{X}}\tilde{\vec{Z}}\|_2^2 = \tilde{\vec{Z}}^T \tilde{\vec{X}}^T \tilde{\vec{X}} \tilde{\vec{Z}}$$
 goes to zero?

How to show?

Argue that it decreases^{in expectation} with every step

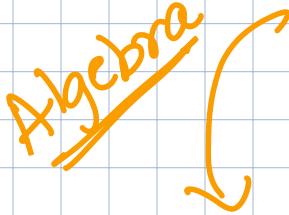
and is non-negative. Furthermore $\mathcal{L}(\vec{Z})=0 \Rightarrow \vec{Z}=0.$

"Lyapunov function"

since $\tilde{\vec{X}}$ is full rank.

So how can we connect $\mathcal{L}(z_{t+1})$ to $\mathcal{L}(z_t)$?

Write $\mathcal{L}(\vec{z}_{t+1}) = \mathcal{L}(\underline{\vec{z}_t} + (\vec{z}_{t+1} - \underline{\vec{z}_t}))$

Algebra 
 $= \mathcal{L}(\vec{z}_t) + A + B.$

$$A = 2\vec{z}_t^T \tilde{\vec{X}}^T \tilde{\vec{X}} (\vec{z}_{t+1} - \vec{z}_t)$$

$$B = (\vec{z}_{t+1} - \vec{z}_t)^T \tilde{\vec{X}}^T \tilde{\vec{X}} (\vec{z}_{t+1} - \vec{z}_t)$$

The linear term looks like a gradient step given \tilde{z}_t .

Can we show a contraction?

$$E[A|\tilde{z}_t] \leq -\eta\gamma \cdot L(\tilde{z}_t)$$

Recall only randomness here
comes from SGD.

$$\begin{aligned} E[A|\tilde{z}_t] &= 2\tilde{z}_t^T \tilde{X}^T \tilde{X} E[\tilde{z}_{t+1}, -\tilde{z}_t | \tilde{z}_t] \\ &= 2\tilde{z}_t^T \tilde{X}^T \tilde{X} (-2\eta) E[\tilde{x}_{I_t} \tilde{x}_{I_t}^T] \cdot \tilde{z}_t \end{aligned}$$

$$E[\tilde{x}_{I_t} \tilde{x}_{I_t}^T] = \frac{1}{n} \sum_{i=1}^n \tilde{x}_i \tilde{x}_i^T = \frac{1}{n} \tilde{X}^T \tilde{X}$$



$$\Rightarrow E[A|\tilde{z}_t] = -\frac{4\eta}{n} \cdot \tilde{z}_t^T \tilde{X}^T \tilde{X} \boxed{\tilde{X}^T \tilde{X} \cdot \tilde{z}_t}$$

max scaling is σ_{\min}^2

$$\geq -\frac{4\eta}{n} \cdot \sigma_{\min} \cdot \tilde{z}_t^T \tilde{X}^T \tilde{X} \tilde{z}_t$$

$$= -\frac{4\eta}{n} \cdot \sigma_{\min} L(\tilde{z}_t)$$

Now for the B term, quadratic in $(\vec{z}_{t+1} - \vec{z}_t)$.

$$E[B | \vec{z}_t] \leq C \cdot \eta \cdot L(\vec{z}_t)$$

$$\begin{aligned} E[B | \vec{w}_t] &= E\left[(\vec{z}_{t+1} - \vec{z}_t)^T \tilde{\vec{x}}^T \tilde{\vec{x}} (\vec{z}_{t+1} - \vec{z}_t) \right] \\ &\leq 4\eta^2 \underbrace{\sigma_{\max}^2}_{\text{max scaling.}} \vec{z}_t^T E\left[\tilde{\vec{x}}_{I_t} \tilde{\vec{x}}_{I_t}^T \tilde{\vec{x}}_{I_t} \tilde{\vec{x}}_{I_t}^T \right] \tilde{\vec{z}}_t \\ &\quad \cdot \|\vec{x}_{I_t}\|^2 \leq \gamma \end{aligned}$$

$$\vec{z}_{t+1} - \vec{z}_t$$

$$= -2\eta \tilde{\vec{x}}_{I_t} \tilde{\vec{x}}_{I_t}^T \vec{z}_t$$

$$\leq 2\eta^2 \underbrace{\sigma_{\max}^2}_{\text{max scaling.}} \cdot \gamma^2 \cdot L(\vec{w}_t)$$

Putting it all together.

$$E[L(\vec{z}_{t+1}) | z_t] \leq (1 - c_1\eta + c_2\eta^2) L(w_t).$$

smaller than 1 if η is small enough.

Takeaway:

- Interesting dependence on norms of data points.
- Decoupling to SVD helps us see what is going on.
- SGD shares many properties of GD.
e.g. can only move in span of data.