

Lecture 13: Gradient Descent for linear regression

W

Gradient descent for linear regression

- For linear regression, we have $w^* = \arg \min_{w \in \mathbb{R}^d} \underbrace{\|\mathbf{y} - \mathbf{X}w\|_2^2}_{f(w)}$
- Gradient Descent:
 - Initialize: $w_0 = 0$
 - For $t=0,1,2,\dots$
 - $w_{t+1} \leftarrow w_t - \eta \cdot \nabla_w f(w_t)$

$$\nabla f(w_t) = -2\mathbf{X}^T(\mathbf{y} - \mathbf{X}w_t)$$

$$w_{t+1} = w_t + \eta 2\mathbf{X}^T(\mathbf{y} - \mathbf{X}w_t) = (\mathbf{I} - 2\eta\mathbf{X}^T\mathbf{X})w_t + 2\eta\mathbf{X}^T\mathbf{y}$$

Let the least-squares solution be $w^* = (\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T\mathbf{y}$

$$\begin{aligned} w_{t+1} - w^* &= (\mathbf{I} - 2\eta\mathbf{X}^T\mathbf{X})w_t + 2\eta\mathbf{X}^T\mathbf{y} - w^* \\ &= (\mathbf{I} - 2\eta\mathbf{X}^T\mathbf{X})(w_t - w^*) + 2\eta\mathbf{X}^T\mathbf{y} - 2\eta\mathbf{X}^T\mathbf{X}w^* \\ &= (\mathbf{I} - 2\eta\mathbf{X}^T\mathbf{X})(w_t - w^*) \end{aligned}$$

Gradient Descent (GD) for Linear Regression (LR)

- We use this analytical derivation of GD for LR to understand how the choice of step size η impacts the algorithm

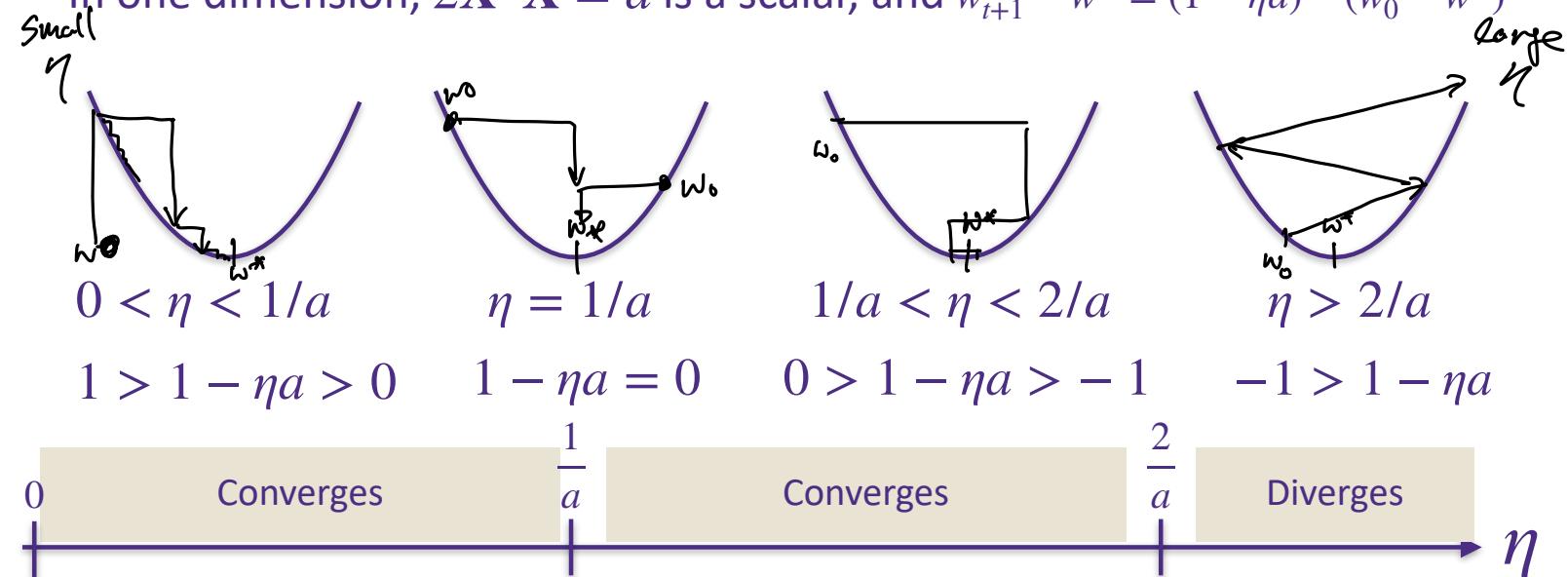
$$\begin{aligned} w_{t+1} = w_t - \eta \nabla f(w_t) \implies w_{t+1} - w^* &= (\mathbf{I} - 2\eta \mathbf{X}^T \mathbf{X})(w_t - w^*) \\ &= (\mathbb{I} - 2\eta X^T X) (\mathbb{I} - 2\eta X^T X)(w_{t-1} - w^*) \\ &\quad \vdots \\ &= (\mathbb{I} - 2\eta X^T X)^{t+1} (w_0 - w^*) \end{aligned}$$

\uparrow evolution \uparrow initial error.

Gradient descent for linear regression

$$\begin{aligned} w_{t+1} &= w_t - \eta \nabla f(w_t) \implies w_{t+1} - w^* &= (\mathbf{I} - 2\eta \mathbf{X}^T \mathbf{X})(w_t - w^*) \\ &&= (\mathbf{I} - 2\eta \mathbf{X}^T \mathbf{X})^2(w_{t-1} - w^*) \\ &&\vdots \\ &&= (\mathbf{I} - 2\eta \mathbf{X}^T \mathbf{X})^{t+1}(w_0 - w^*) \end{aligned}$$

In one dimension, $2\mathbf{X}^T \mathbf{X} = a^D$ is a scalar, and $w_{t+1} - w^* = (1 - \eta a)^{t+1}(w_0 - w^*)$



Gradient descent for linear regression

$$w_{t+1} = w_t - \eta \nabla f(w_t) \implies w_{t+1} - w^* = (\mathbf{I} - 2\eta \mathbf{X}^T \mathbf{X})^{t+1} (w_0 - w^*)$$

- In multi dimensions, **eigenvalues** of $\mathbf{I} - 2\eta \mathbf{X}^T \mathbf{X}$ are important
- Let the eigenvalue decomposition of $\mathbf{I} - 2\eta \mathbf{X}^T \mathbf{X}$ be $Q^{-1}DQ$,
 - Where D is a diagonal matrix with Eigen values $\{D_{ii}\}_{i=1}^d$ in the diagonal
 - And Q is an orthogonal matrix, with each Eigen vector as in a row

$$\mathbf{I} - 2\eta \mathbf{X}^T \mathbf{X} = Q^{-1} \cdot D \cdot Q , \quad D = \begin{bmatrix} D_{11} & & \\ & D_{22} & \\ & & \ddots & \\ & & & D_{dd} \end{bmatrix} \quad Q = \begin{bmatrix} \mathbf{\xi}_1^T \\ \mathbf{\xi}_2^T \\ \vdots \\ \mathbf{\xi}_d^T \end{bmatrix}$$

$$w_{\text{eff}, t+1} - w^* = \underbrace{Q^{-1} D Q \cdot Q^{-1} D Q \cdot Q^{-1} D Q \cdots}_{t+1} \cdot (w_0 - w^*)$$

$$= Q^{-1} \underbrace{D^{t+1}}_{\begin{bmatrix} D_{11}^{t+1} & & \\ & D_{22}^{t+1} & \\ & & \ddots \end{bmatrix}} \cdot Q (w_0 - w^*)$$

$$\begin{aligned} \xrightarrow{\text{defn}} \underbrace{Q \cdot (w_{\text{eff}, t+1} - w^*)}_{\text{drift}} &= \underbrace{D^{t+1} \cdot Q (w_0 - w^*)}_{\mathbf{\xi}_1^T (w_{\text{eff}, t+1} - w^*)} \\ \rightarrow \mathbf{\xi}_1^T (w_{\text{eff}, t+1} - w^*) &= \underbrace{D_{11}^{t+1} \cdot \mathbf{\xi}_1^T (w_0 - w^*)}_{\{(x_i, D_{11})\}} \end{aligned}$$

Gradient descent for linear regression

$$w_{t+1} = w_t - \eta \nabla f(w_t) \implies w_{t+1} - w^* = (\mathbf{I} - 2\eta \mathbf{X}^T \mathbf{X})^{t+1} (w_0 - w^*)$$

- In multi dimensions, **eigenvalues** of $\mathbf{I} - 2\eta \mathbf{X}^T \mathbf{X}$ are important
- Let the eigenvalue decomposition of $\mathbf{I} - 2\eta \mathbf{X}^T \mathbf{X}$ be $Q^{-1}DQ$

$$\begin{aligned} \text{Then, } w_{t+1} - w^* &= (Q^{-1}DQ)^{t+1} (w_0 - w^*) \\ &= \underbrace{Q^{-1}DQ Q^{-1}DQ \cdots Q^{-1}DQ}_{t+1 \text{ times}} (w_0 - w^*) \\ &= Q^{-1} D^{t+1} Q (w_0 - w^*) \\ Q(w_{t+1} - w^*) &= D^{t+1} Q (w_0 - w^*) \end{aligned}$$

- This defines a series of equations capturing how the error evolves in Directions defined by the rows of Q , which are the Eigen vectors of $\mathbf{I} - 2\eta \mathbf{X}^T \mathbf{X}$

Gradient descent for linear regression

$$Q(w_{t+1} - w^*) = D^{t+1} Q (w_0 - w^*)$$

- Where eigenvalue decomposition of $\mathbf{I} - 2\eta \mathbf{X}^T \mathbf{X}$ is $Q^{-1} D Q$

- Let $Q = \begin{bmatrix} - & q_1^T & - \\ - & q_2^T & - \\ & \vdots & \end{bmatrix}$, then the above multi-dimensional dynamics

of GD can be decomposed into multiple 1-d dynamics we saw before

- The eigenvector-eigenvalue pairs $\{(q_i, D_{ii})\}_{i=1}^d$ of the matrix $\mathbf{I} - 2\eta \mathbf{X}^T \mathbf{X}$ determines the behavior of gradient descent
- In direction q_1 , the error decreases multiplicatively according to D_{11}

$$\text{Error in direction } q_1 \longrightarrow q_1^T (w_{t+1} - w^*) = D_{11}^{t+1} q_1^T (w_0 - w^*)$$

$$q_2^T (w_{t+1} - w^*) = D_{22}^{t+1} q_2^T (w_0 - w^*)$$

⋮
⋮
⋮

Gradient descent for linear regression

$$w_{t+1} = w_t - \eta \nabla f(w_t) \implies w_{t+1} - w^* = (\mathbf{I} - 2\eta \mathbf{X}^T \mathbf{X})^{t+1} (w_0 - w^*)$$

$$\implies Q(w_{t+1} - w^*) = D^{t+1} Q(w_0 - w^*)$$

$$q_1^T (w_{t+1} - w^*) = D_{11}^{t+1} q_1^T (w_0 - w^*)$$

$$q_2^T (w_{t+1} - w^*) = D_{22}^{t+1} q_2^T (w_0 - w^*)$$

- For example suppose, the step size η is chosen such that

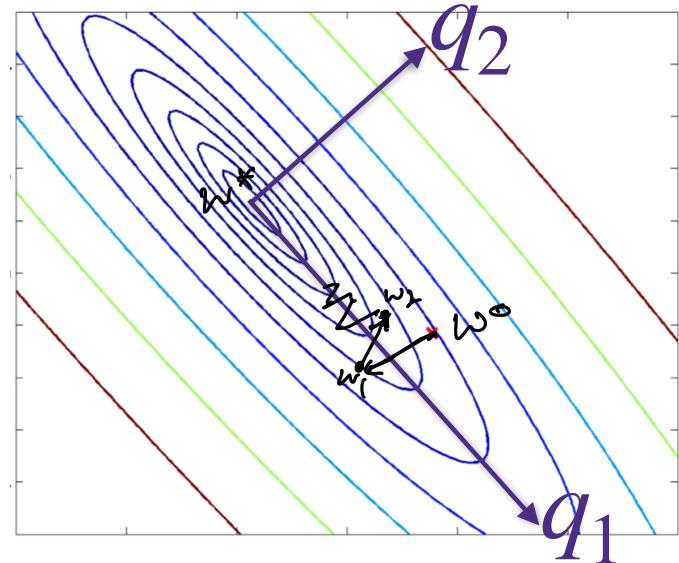
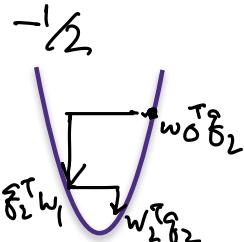
In direction q_1

$$0 < D_{11} < 1$$

0.9

In direction q_2

$$-1 < D_{22} < 0$$



Gradient descent for logistic regression

- Now we know how to find the global minimum of a logistic regression problem, numerically

Loss function: Conditional Likelihood

$$\{(x_i, y_i)\}_{i=1}^n \quad x_i \in \mathbb{R}^d, \quad y_i \in \{-1, 1\}$$

$$\widehat{w}_{MLE} = \arg \max_w \prod_{i=1}^n P(y_i | x_i, w) \quad P(Y = y | x, w) = \frac{1}{1 + \exp(-y w^T x)}$$

$$= \arg \min_w \sum_{i=1}^n \log(1 + \exp(-y_i x_i^T w))$$

$$\nabla f(w) = \sum_{i=1}^n \frac{1}{1 + \exp(-y_i x_i^T w)} \exp(-y_i x_i^T w) (-y_i x_i)$$

What is known for Gradient descent for convex fcs)

- $f(\cdot)$ is L -smooth if $\|\nabla f(w) - \nabla f(v)\|_2 \leq L\|w - v\|_2$ for all $w, v \in \mathbb{R}^d$
- $f(\cdot)$ is μ -strongly convex if $f(w) \geq f(v) + \nabla f(v)^T(w - v) + \frac{\mu}{2}\|w - v\|_2^2$
- For L -smooth functions, with a fixed step size $\eta < 1/L$

- if $f(w)$ is convex,

$$f(w_t) - f(w^*) \leq \frac{\|w_0 - w^*\|_2^2}{2\eta t} \leq \frac{L \cdot \cancel{\mathbb{R}^2}}{2t} = O\left(\frac{1}{t}\right)$$

- if $f(w)$ is μ -strongly convex,

$$f(w_t) - f(w^*) \leq (1 - \eta\mu)^t (f(w_0) - f(w^*)) \cong O(e^{-\frac{1}{\mu} \cdot t})$$

- Gradient Descent is oftentimes called full-batch gradient descent to differentiate it from stochastic gradient descent, which uses only a (randomly chosen) subset of training data at each iteration
- In practice, people use Stochastic Gradient Descent (SGD).

Questions?

Stochastic Gradient Descent

-What do we use in practice?

W

Machine Learning Problems

- Given data: $\{(x_i, y_i)\}_{i=1}^n \quad x_i \in \mathbb{R}^d \quad y_i \in \mathbb{R}$
- Learning a model's parameters: $\frac{1}{n} \sum_{i=1}^n \ell_i(w) = \frac{1}{n} \sum_{i=1}^n \ell(f_w(x_i), y_i)$

• Gradient Descent (GD):

one update takes cdn operations/time for some constant $c > 0$

$$w_{t+1} \leftarrow w_t - \eta \frac{1}{n} \sum_{i=1}^n \nabla \ell_i(w_t)$$

• Stochastic Gradient Descent (SGD): one update takes cd operations/time

$$w_{t+1} \leftarrow w_t - \eta \boxed{\nabla \ell_{I_t}(w_t)}$$

I_t drawn uniform at random from $\{1, \dots, n\}$

• SGD is an unbiased estimate of the GD

$$\mathbb{E}[\nabla \ell_{I_t}(w)] = \sum_{i=1}^n \Pr(I_t=i) \cdot \nabla \ell_i(w) = \frac{1}{n} \sum_{i=1}^n \nabla \ell_i(w).$$

Stochastic Gradient Descent

Theorem

Let $w_{t+1} = w_t - \eta \nabla_w \ell_{I_t}(w) \Big|_{w=w_t}$ I_t drawn uniform at random from $\{1, \dots, n\}$ so that

$$\mathbb{E}[\nabla \ell_{I_t}(w)] = \frac{1}{n} \sum_{i=1}^n \nabla \ell_i(w) =: \nabla \ell(w)$$

If $\|w_0 - w_*\|_2^2 \leq R$ and $\sup_w \max_i \|\nabla \ell_i(w)\|_2^2 \leq G$ then

after T iterations we will show

$$\mathbb{E}[\ell(\bar{w}) - \ell(w_*)] \stackrel{\textcircled{1}}{\leq} \underbrace{\frac{R}{2T\eta} + \frac{\eta G}{2}}_{\text{RHS}} \stackrel{\textcircled{2}}{\leq} \sqrt{\frac{RG}{T}}$$

$$\eta = \sqrt{\frac{R}{GT}}$$

$$\bar{w} = \frac{1}{T} \sum_{t=1}^T w_t$$

(In practice use last iterate)

Convergence rate: $O\left(\frac{1}{\sqrt{T}}\right)$

$$-\frac{R}{2T\eta^2} + \frac{G}{2} = 0$$

Taking the derivative of RHS to zero

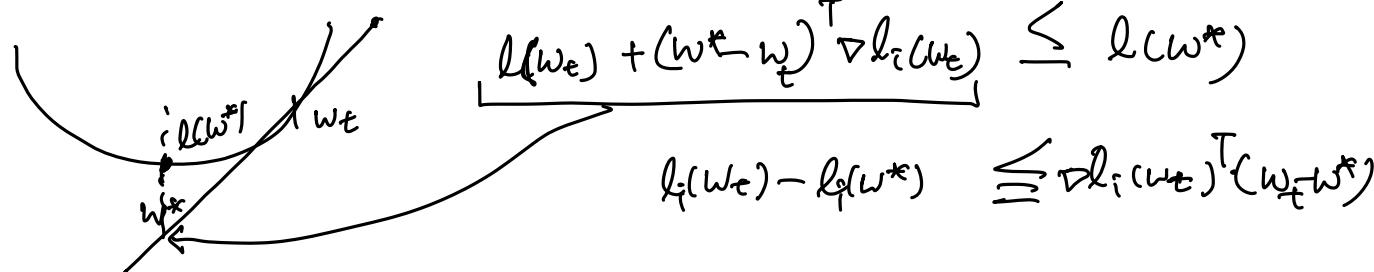
We want to show that

$$\begin{aligned}
 \mathbb{E} \left[\ell \left(\frac{1}{T} \sum_{t=1}^T w_t \right) - \ell(w_*) \right] &\leq \mathbb{E} \left[\frac{1}{T} \sum_{i=1}^T \ell(w_t) - \ell(w_*) \right] && \text{Follows from convexity of } \ell(\cdot) \\
 &\stackrel{\text{and Jensen's inequality}}{\leq} \frac{1}{T} \sum_{i=1}^T \mathbb{E} [\ell(w_t) - \ell(w_*)] && \text{(3 slides later)} \\
 &\stackrel{\text{Follows from linearity of expectation}}{\leq} \frac{R}{2T\eta} + \frac{\eta G}{2} && \text{We are left to show this}
 \end{aligned}$$

Proof

$$\begin{aligned}
 \mathbb{E}[\|w_{t+1} - w_*\|_2^2] &= \mathbb{E}[\|w_t - \eta \nabla \ell_{I_t}(w_t) - w_*\|_2^2] && \leftarrow w_{t+1} = w_t - \eta \nabla \ell_{I_t}(w_t) \\
 &= \mathbb{E}[\|w_t - w_*\|_2^2] + \eta^2 \cdot \underbrace{\mathbb{E}[\|\nabla \ell_{I_t}(w_t)\|^2]}_{\leq G} - 2\eta \underbrace{\mathbb{E}[\nabla \ell_{I_t}(w_t)^T (w_t - w_*)]}_{\ell(w_t) - \ell(w_*)} \\
 &\leq \mathbb{E}[\|w_t - w_*\|_2^2] + \eta^2 \cdot \underbrace{\mathbb{E}[\ell_{I_t}(w_t) - \ell_{I_t}(w^*)]}_{\ell(w_t) - \ell(w^*)}
 \end{aligned}$$

*Convexity $\ell_i(w)$



Stochastic Gradient Descent

Proof

$$\mathbb{E}[||w_{t+1} - w_*||_2^2] = \mathbb{E}[||w_t - \eta \nabla \ell_{I_t}(w_t) - w_*||_2^2]$$

$$\leq \mathbb{E}[||w_t - w_*||_2^2] + \eta^2 G - 2\eta (\ell(w_t) - \ell(w_*))$$

$$\sum_{t=0}^{T-1} (\ell(w_t) - \ell(w_*)) \leq \frac{1}{2\eta} \left(\underbrace{\mathbb{E}[||w_t - w_*||_2^2] - \mathbb{E}[||w_{t+1} - w_*||_2^2]}_{\text{Towering}} + \eta^2 G \right)$$
$$\sum_{t=1}^T (\ell(w_t) - \ell(w_*)) \leq \frac{1}{2\eta} \left(\mathbb{E}[||w_0 - w_*||_2^2] - \underbrace{\mathbb{E}[||w_T - w_*||_2^2]}_{\leq R} + T\eta^2 G \right)$$
$$\leq \frac{R}{2\eta} + \frac{T\eta^2 G}{2}$$

Stochastic Gradient Descent

Proof

$$\begin{aligned}\mathbb{E}[||w_{t+1} - w_*||_2^2] &= \mathbb{E}[||w_t - \eta \nabla \ell_{I_t}(w_t) - w_*||_2^2] \\ &= \mathbb{E}[||w_t - w_*||_2^2] - 2\eta \mathbb{E}[\nabla \ell_{I_t}(w_t)^T (w_t - w_*)] + \eta^2 \mathbb{E}[||\nabla \ell_{I_t}(w_t)||_2^2] \\ &\leq \mathbb{E}[||w_t - w_*||_2^2] - 2\eta \mathbb{E}[\ell(w_t) - \ell(w_*)] + \eta^2 G\end{aligned}$$

$$\begin{aligned}\mathbb{E}[\nabla \ell_{I_t}(w_t)^T (w_t - w_*)] &= \mathbb{E}[\mathbb{E}[\nabla \ell_{I_t}(w_t)^T (w_t - w_*) | I_1, w_1, \dots, I_{t-1}, w_{t-1}]] \\ &= \mathbb{E}[\nabla \ell(w_t)^T (w_t - w_*)] \\ &\geq \mathbb{E}[\ell(w_t) - \ell(w_*)]\end{aligned}$$

$$\begin{aligned}\sum_{t=1}^T \mathbb{E}[\ell(w_t) - \ell(w_*)] &\leq \frac{1}{2\eta} (\mathbb{E}[||w_1 - w_*||_2^2] - \mathbb{E}[||w_{T+1} - w_*||_2^2] + T\eta^2 G) \\ &\leq \frac{R}{2\eta} + \frac{T\eta G}{2}\end{aligned}$$

We have:

$$\sum_{t=1}^T \mathbb{E}[\ell(w_t) - \ell(w_*)] \leq \frac{R}{2\eta} + \frac{T\eta G}{2}$$

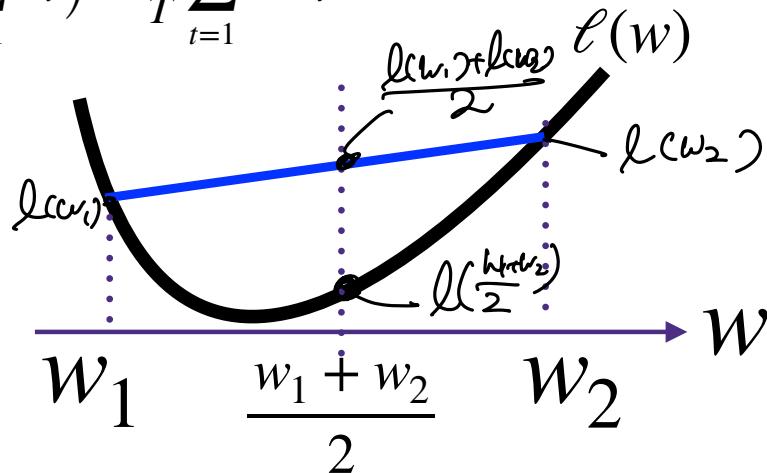
$$\begin{aligned}\mathbb{E}[\ell(\bar{w}) - \ell(w_*)] &\stackrel{\textcircled{1}}{\leq} \frac{1}{T} \sum_{t=1}^T \mathbb{E}[\ell(w_t) - \ell(w_*)] \\ &\leq \frac{R}{2\eta T} + \frac{\eta G}{2} \quad \longrightarrow (\star)\end{aligned}$$

$$\bar{w} = \frac{1}{T} \sum_{t=1}^T w_t$$

① Jensen's inequality:

For any $\{w_1, \dots, w_T\}$ and a convex function $\ell(\cdot)$, we have

$$\ell\left(\frac{1}{T} \sum_{t=1}^T w_t\right) \leq \frac{1}{T} \sum_{t=1}^T \ell(w_t)$$



Mini-batch SGD

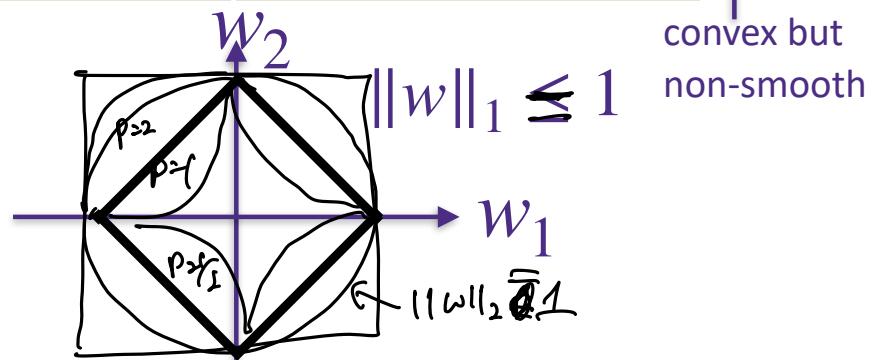
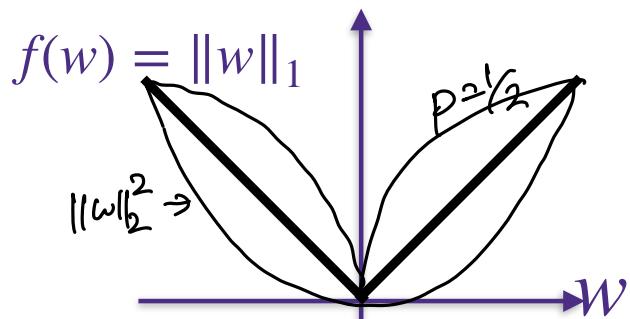
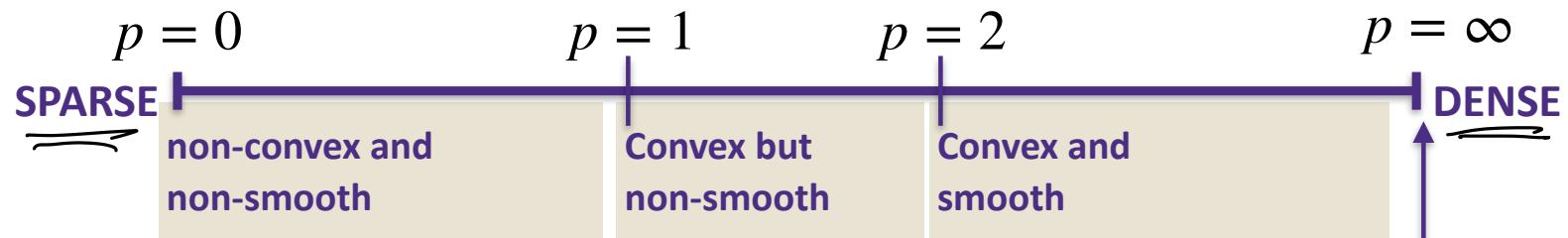
- Instead of one iterate, average B stochastic gradient together
- Advantages:
 - Smaller variance: the variance of the stochastic gradient is smaller by a factor of $1/\sqrt{B}$
 - Parallelization: each gradient in the mini-batch can be computed in parallel
- If you have regularizer, $\frac{1}{n} \sum_{i=1}^n \ell_i(w) + r(w)$, then update with the stochastic gradient of the loss and gradient of the regularizer

Sparsity/Complexity tradeoff

- ℓ_p -norm of a vector is defined as $\|w\|_p \triangleq \left(w_1^p + w_2^p + \dots + w_d^p \right)^{1/p}$
- Consider regularized least squares problem of minimizing
$$\mathcal{L}(w) = \sum_{i=1}^n (y_i - w^T x_i)^2 + \lambda \|w\|_p^p$$
- This is ridge regression for $p = 2$ and Lasso for $p = 1$

$$\|w\|_0 = \# \text{ of non-zero entries}$$

$$\|w\|_\infty = \max\{w_i\}$$



Questions?
