Convex learning problems CPSC 532S: Modern Statistical Learning Theory 14 February 2022

cs.ubc.ca/~dsuth/532S/22/

Admin

- In hybrid mode now:
 - Thursday office hours available both in-person (ICICS X563) and on Zoom
- A2 due Friday night
 - Groups of up to three, allowed separate per question • If you don't have a group and want one, post on Piazza (asap)
- A1 grading: *allllllmost* done sorry

So far, only really talked about ERM (and variants like SRM)

- So far, only really talked about ERM (and variants like SRM)

• Not always practical: e.g. NP-hard to maximize accuracy of a linear binary classifier

- So far, only really talked about ERM (and variants like SRM)
- A scheme that usually is practical: **convex** learning problems

• Not always practical: e.g. NP-hard to maximize accuracy of a linear binary classifier

- So far, only really talked about ERM (and variants like SRM)
- Not always practical: e.g. NP-hard to maximize accuracy of a linear binary classifier
- A scheme that usually is practical: **convex** learning problems • Can get an ε -approximate ERM with gradient descent:

- So far, only really talked about ERM (and variants like SRM)
- Not always practical: e.g. NP-hard to maximize accuracy of a linear binary classifier
- A scheme that usually is practical: **convex** learning problems
 - Can get an ε -approximate ERM with gradient descent: • in $\mathcal{O}\left(\frac{1}{\epsilon}\right)$ steps, if loss is convex and has Lipschitz gradients

- So far, only really talked about ERM (and variants like SRM)
- Not always practical: e.g. NP-hard to maximize accuracy of a linear binary classifier
- A scheme that usually is practical: **convex** learning problems
 - Can get an ε -approximate ERM with gradient descent: • in $\mathcal{O}\left(\frac{1}{\varepsilon}\right)$ steps, if loss is convex and has Lipschitz gradients • in $\mathcal{O}\left(\log\frac{1}{\varepsilon}\right)$ steps, if loss is strongly convex with Lipschitz gradients

Convexity review

- Most of you probably already know most of this (based on the survey), but a quick reminder!
- For proofs and details, see e.g. chaps 2-3 of <u>Boyd and Vandenberghe</u> (free pdf)

Stephen Boyd and Lieven Vandenberghe

Convex Optimization



Convex sets

• If $x, y \in C$, then $\alpha x + (1 - \alpha)y \in C$ for all $\alpha \in [0, 1]$





 $f: \mathcal{X} \to \mathbb{R}$ is convex if, when \mathcal{X} is a convex set:

• the epigraph $\{(x, r) \in \mathcal{X} \times \mathbb{R} : r \ge f(x)\}$ is a convex set

 $f: \mathcal{X} \to \mathbb{R}$ is convex if, when \mathcal{X} is a convex set:



• the epigraph $\{(x, r) \in \mathcal{X} \times \mathbb{R} : r \ge f(x)\}$ is a convex set • it lies below its chords, $f(\alpha x + (1 - \alpha)y) \nearrow \alpha f(x) + (1 - \alpha)f(y)$ for $\alpha \in [0, 1]$

 $f: \mathcal{X} \to \mathbb{R}$ is convex if, when \mathcal{X} is a convex set:

- the epigraph $\{(x, r) \in \mathcal{X} \times \mathbb{R} : r \ge f(x)\}$ is a convex set
- it lies below its chords, $f(\alpha x + (1 \alpha)y) \ge \alpha f(x) + (1 \alpha)f(y)$ for $\alpha \in [0, 1]$
- if f is differentiable: convex iff f lies above its tangent planes $f(x) \ge f(y) + [\nabla f(y)](x - y)$ for all x, y



 $f: \mathcal{X} \to \mathbb{R}$ is convex if, when \mathcal{X} is a convex set:

- the epigraph $\{(x, r) \in \mathcal{X} \times \mathbb{R} : r \ge f(x)\}$ is a convex set
- it lies below its chords, $f(\alpha x + (1 \alpha)y) \ge \alpha f(x) + (1 \alpha)f(y)$ for $\alpha \in [0, 1]$
- if f is differentiable: convex iff f lies above its tangent planes $f(x) \ge f(y) + [\nabla f(y)](x - y)$ for all x, y
- if f is twice-differentiable: convex iff its Hessian $\nabla^2 f$ is positive semidefinite $\nabla^2 f \ge 0$ or all eigenvalues ≥ 0 or $v^\top [\nabla^2 f] v \ge 0$ for all v or $\nabla^2 f = A^\top A$



If f and g are convex functions, then so are

- If f and g are convex functions, then so are
- αf for any $\alpha \ge 0$

If f and g are convex functions, then so are

• αf for any $\alpha \ge 0$ • f + g, or even $\int_{\mathscr{A}} f_y(x) dw(y)$ if each f_y is convex and w a measure $\mathcal{L}(\cdot, \gamma_i) = (\cdot - \gamma_i)^2$

Ls (Ŷ) CR



If f and g are convex functions, then so are

- αf for any $\alpha \ge 0$ • f + g, or even $\int_{a} f_y(x) dw(y)$ if each f_y is convex and w a measure
- $x \mapsto f(Ax + b)$

Operations that preserve convexity

 $l_i(\hat{y}) = (\hat{y} - \gamma_i)^2$ $\widetilde{l_i}(w) = l_i(wx_i) = (wx - y_i)^2$

 $L_{S}(w) = L_{D}(w)$

- If f and g are convex functions, then so are
- αf for any $\alpha \ge 0$ • f + g, or even $\int_{\mathscr{A}} f_y(x) dw(y)$ if each f_y is convex and w a measure
- $x \mapsto f(Ax + b)$
- $x \mapsto g(f(x))$ if $f : \mathcal{X} \to \mathbb{R}$ and g is nondecreasing

- If f and g are convex functions, then so are
- αf for any $\alpha \ge 0$ • f + g, or even $\int_{-\pi}^{\pi} f_y(x) dw(y)$ if each f_y is convex and w a measure
- $x \mapsto f(Ax + b)$
- $x \mapsto g(f(x))$ if $f : \mathcal{X} \to \mathbb{R}$ and g is nondecreasing • $x \mapsto \max(f(x), g(x))$, or even $x \mapsto \sup f_v(x)$

 $y \in \mathscr{A}$

- If f and g are convex functions, then so are
- αf for any $\alpha \ge 0$ • f + g, or even $\int_{\mathcal{A}} f_y(x) dw(y)$ if each f_y is convex and w a measure
- $x \mapsto f(Ax + b)$
- $x \mapsto g(f(x))$ if $f : \mathcal{X} \to \mathbb{R}$ and g is nondecreasing
- $x \mapsto \max(f(x), g(x))$, or even $x \mapsto \sup f_v(x)$ $y \in \mathscr{A}$
- If f(x, y) is convex in (x, y), then $h(x) = \inf f(x, y)$ for nonempty convex C $y \in C$

- If f and g are convex functions, then so are
- αf for any $\alpha \ge 0$ • f + g, or even $\int_{\mathcal{A}} f_y(x) dw(y)$ if each f_y is convex and w a measure
- $x \mapsto f(Ax + b)$
- $x \mapsto g(f(x))$ if $f : \mathcal{X} \to \mathbb{R}$ and g is nondecreasing
- $x \mapsto \max(f(x), g(x)), \text{ or even } x \vdash$
- If f(x, y) is convex in (x, y), then $h(x) = \inf f(x, y)$ for nonempty convex C $v \in C$ • Perspective transform: $h(x, t) = tf\left(\frac{x}{t}\right)$ for t > 0

$$\rightarrow \sup_{y \in \mathscr{A}} f_y(x)$$

(pause)

- A learning problem $(\mathcal{H}, \mathcal{I}, \ell)$ is convex if
 - *H* is a convex set
 - for each $z \in \mathcal{Z}$, $\ell(\cdot, z)$ is a convex function

94 = EXH WTX : WERd3

91'= Ew: wERdz



af+bg = xf af(x)+bgG

 $l(h, 2) = (2h, x7 - y)^2$ entry (x, 7)



- A learning problem $(\mathcal{H}, \mathcal{I}, \ell)$ is convex if
 - \mathcal{H} is a convex set
 - for each $z \in \mathcal{Z}$, $\ell(\cdot, z)$ is a convex function
- Example: linear regression with square loss

- A learning problem $(\mathcal{H}, \mathcal{I}, \ell)$ is convex if
 - *H* is a convex set
 - for each $z \in \mathcal{Z}$, $\ell(\cdot, z)$ is a convex function
- Example: linear regression with square loss
- Non-example: linear classifiers with 0-1 loss



- A learning problem $(\mathcal{H}, \mathcal{I}, \ell)$ is convex if
 - *H* is a convex set
 - for each $z \in \mathcal{Z}$, $\ell(\cdot, z)$ is a convex function
- Example: linear regression with square loss
- Non-example: linear classifiers with 0-1 loss

• For convex learning problems, ERM is a convex optimization problem

- A learning problem $(\mathcal{H}, \mathcal{I}, \ell)$ is convex if
 - \mathcal{H} is a convex set
 - for each $z \in \mathcal{Z}$, $\ell(\cdot, z)$ is a convex function
- Example: linear regression with square loss
- Non-example: linear classifiers with 0-1 loss
- - Usually implies learnable in polynomial time (but not always)

For convex learning problems, ERM is a convex optimization problem

• Consider (homogeneous) linear regression on \mathbb{R} , $\ell(w, (x, y)) = (wx - y)^2$



- Consider (homogeneous) linear regression on \mathbb{R} , $\ell(w, (x, y)) = (wx y)^2$
- Suppose a deterministic A can (agnostically) PAC learn this problem



- Consider (homogeneous) linear regression on \mathbb{R} , $\ell(w, (x, y)) = (wx y)^2$
- Suppose a deterministic A can (agnostically) PAC learn this problem
- Take $\varepsilon = 0.01$, $\delta = 1/2$, *n* big enough that $L_{\Im}(A(S)) \inf L_{\Im}(w) \leq \varepsilon$



- Consider (homogeneous) linear regression on \mathbb{R} , $\ell(w, (x, y)) = (wx y)^2$
- Suppose a deterministic A can (agnostically) PAC learn this problem
- Take $\varepsilon = 0.01$, $\delta = 1/2$, *n* big enough that $L_{\mathcal{D}}(A(S)) \inf L_{\mathcal{D}}(w) \leq \varepsilon$ $\mathcal{D}_1(\{z\}) = \begin{cases} \mu & z = z_1 \\ 1 - \mu & z = z_2 \end{cases}; \quad \mathcal{D}_2(\{z_2\}) = 1$

• Let
$$\mu = \frac{1}{2n} \log \frac{100}{99}$$
; $z_1 = (1,0)$; $z_2 = (\mu, -1)$;

~ 0.01





- Consider (homogeneous) linear regression on \mathbb{R} , $\ell(w, (x, y)) = (wx y)^2$
- Suppose a deterministic A can (agnostically) PAC learn this problem
- Take $\varepsilon = 0.01$, $\delta = 1/2$, *n* big enough that $L_{\mathcal{D}}(A(S)) \inf L_{\mathcal{D}}(w) \leq \varepsilon$ • Let $\mu = \frac{1}{2n} \log \frac{100}{99}$; $z_1 = (1,0)$ $z_2 = (\mu, -1)$; $\mathfrak{D}_1(\{z\}) = \begin{cases} \mu & z = z_1 \\ 1 - \mu & z = z_2 \end{cases}$; $\mathfrak{D}_2(\{z_2\}) = 1$ • At least 99% prob to only see z_2 in S: $\mathfrak{D}_1^n((z_2, ..., z_2)) = (1 - \mu)^n \ge e^{-2\mu n} = 0.99$

• Let
$$\mu = \frac{1}{2n} \log \frac{100}{99}$$
; $z_1 = (1,0)$; $z_2 = (\mu, -1)$; $z_3 = (1,0)$





- Consider (homogeneous) linear regression on \mathbb{R} , $\ell(w, (x, y)) = (wx y)^2$
- Suppose a deterministic A can (agnostically) PAC learn this problem
- Take $\varepsilon = 0.01$, $\delta = 1/2$, *n* big enough that $L_{\mathcal{D}}(A(S)) \inf L_{\mathcal{D}}(w) \leq \varepsilon$ • Let $\mu = \frac{1}{2n} \log \frac{100}{99}$; $z_1 = (1,0)$ $z_2 = (\mu, -1)$; $\mathfrak{D}_1(\{z\}) = \begin{cases} \mu & z = z_1 \\ 1 - \mu & z = z_2 \end{cases}$; $\mathfrak{D}_2(\{z_2\}) = 1$ • At least 99% prob to only see z_2 in S: $\mathfrak{D}_1^n((z_2, ..., z_2)) = (1 - \mu)^n \ge e^{-2\mu n} = 0.99$

• Let
$$\mu = \frac{1}{2n} \log \frac{100}{99}$$
; $z_1 = (1,0)$; $z_2 = (\mu, -1)$; $z_3 = (1,0)$

- Let $\hat{w} = A((z_2, ..., z_2))$:





Convex problems aren't necessarily learnable • Consider (homogeneous) linear regression on \mathbb{R} , $\ell(w, (x, y)) = (wx - y)^2$ • Suppose a deterministic A can (agnostically) PAC learn this problem • Take $\varepsilon = 0.01$, $\delta = 1/2$, *n* big enough that $L_{\mathcal{D}}(A(S)) - \inf L_{\mathcal{D}}(w) \leq \varepsilon$ • Let $\mu = \frac{1}{2n} \log \frac{100}{99}$; $\begin{array}{l} z_1 = (1,0) \\ z_2 = (\mu, -1) \end{array}$; $\mathfrak{D}_1(\{z\}) = \begin{cases} \mu & z = z_1 \\ 1 - \mu & z = z_2 \end{cases}$; $\mathfrak{D}_2(\{z_2\}) = 1$ • At least 99% prob to only see z_2 in S: $\mathfrak{D}_1^n((z_2, \dots, z_2)) = (1 - \mu)^n \ge e^{-2\mu n} = 0.99$

• Let
$$\mu = \frac{1}{2n} \log \frac{100}{99}$$
; $z_1 = (1,0)$; $z_2 = (\mu, -1)$; $z_3 = (\mu, -1)$; $z_4 = (1,0)$

• Let
$$\hat{w} = A\left((z_2, ..., z_2)\right)$$
:
• If $\hat{w} < \frac{-1}{2\mu}$, have $L_{\mathcal{D}_1}(\hat{w}) \ge \mu \cdot (\hat{w} - 0)$

 $0)^2 \ge \frac{1}{4u} = \frac{8n}{\log \frac{100}{\log 2}} > 795n.$





Convex problems aren't necessarily learnable • Consider (homogeneous) linear regression on \mathbb{R} , $\ell(w, (x, y)) = (wx - y)^2$ • Suppose a deterministic A can (agnostically) PAC learn this problem • Take $\varepsilon = 0.01$, $\delta = 1/2$, *n* big enough that $L_{\mathscr{D}}(A(S)) - \inf L_{\mathscr{D}}(w) \leq \varepsilon$ $\mathcal{D}_1(\{z\}) = \begin{cases} \mu & z = z_1 \\ 1 - \mu & z = z_2 \end{cases}; \quad \mathcal{D}_2(\{z_2\}) = 1$ • At least 99% prob to only see z_2 in S: $\mathscr{D}_1^n((z_2, ..., z_2)) = (1 - \mu)^n \ge e^{-2\mu n} = 0.99$

• Let
$$\mu = \frac{1}{2n} \log \frac{100}{99}$$
; $z_1 = (1,0)$; $z_2 = (\mu, -1)$; $z_3 = (1,0)$

- Let $\hat{w} = A((z_2, ..., z_2))$:
 - If $\hat{w} < \frac{-1}{2u}$, have $L_{\mathcal{D}_1}(\hat{w}) \ge \mu \cdot (\hat{w} 0)^2 \ge \frac{1}{4u} = \frac{8n}{\log \frac{100}{99}} > 795n$.

• But $L_{\mathcal{D}_1}(0) = 1 - \mu$, so excess error on $\mathcal{D}_1 \ge \frac{1}{4\mu} - 1 + \mu > 795n - 1 > \varepsilon$




Convex problems aren't necessarily learnable • Consider (homogeneous) linear regression on \mathbb{R} , $\ell(w, (x, y)) = (wx - y)^2$ • Suppose a deterministic A can (agnostically) PAC learn this problem • Take $\varepsilon = 0.01$, $\delta = 1/2$, n big enough that $L_{\mathcal{D}}(A(S)) - \inf L_{\mathcal{D}}(w) \leq \varepsilon$

• Let
$$\mu = \frac{1}{2n} \log \frac{100}{99}$$
; $z_1 = (1,0)$
 $z_2 = (\mu, -1)$; $\mathscr{D}_1(\{z\}) = \begin{cases} \mu & z = z_1 \\ 1 - \mu & z = z_2 \end{cases}$; $\mathscr{D}_2(\{z_2\}) = (1 - \mu)^n \ge e^{-2\mu n} = 0.9$
• At least 99% prob to only see z_2 in S : $\mathscr{D}_1^n((z_2, \dots, z_2)) = (1 - \mu)^n \ge e^{-2\mu n} = 0.9$

- Let $\hat{w} = A\left((z_2, \dots, z_2)\right)$: If $\hat{w} < \frac{-1}{2u}$, have $L_{\mathcal{D}_1}(\hat{w}) \ge \mu \cdot (\hat{w} 0)^2 \ge \frac{1}{4\mu} = \frac{8n}{\log \frac{100}{99}} > 795n$.

• But $L_{\mathcal{D}_1}(0) = 1 - \mu$, so excess error on $\mathcal{D}_1 \ge \frac{1}{4\mu} - 1 + \mu > 795n - 1 > \varepsilon$ • If $\hat{w} \ge \frac{-1}{2\mu}$, then $L_{\mathcal{D}_2}(\hat{w}) \ge 1 \cdot (\mu \hat{w} + 1)^2 \ge (1 - \frac{1}{2})^2 \ge \frac{1}{4}$, but $L_{\mathcal{D}_2}\left(\frac{-1}{\mu}\right) = 0$



=



the norm of $h, \mathcal{H} = \{w \in \mathbb{R}^d : ||w|| \le B\}$

• Remember our bounds on \mathfrak{R}_n for linear classes depended on bounding

- Remember our bounds on \Re_n for linear classes depended on bounding the norm of $h, \mathcal{H} = \{ w \in \mathbb{R}^d : ||w|| \le B \}$
 - Optimal w on \mathcal{D}_2 is $-1/\mu\approx-200n,$ which is really big

- Remember our bounds on \Re_n for linear classes depended on bounding the norm of $h, \mathcal{H} = \{w \in \mathbb{R}^d : ||w|| \le B\}$
 - Optimal w on \mathscr{D}_2 is $-1/\mu \approx -200n$, which is really big

• Counterexample for learning $\mathcal{H} = \{x \mapsto wx : |w| \le 1\}$ with square loss:

- Remember our bounds on \mathfrak{R}_n for linear classes depended on bounding the norm of $h, \mathcal{H} = \{w \in \mathbb{R}^d : ||w|| \le B\}$
 - Optimal w on \mathscr{D}_2 is $-1/\mu \approx -200n$, which is really big

• Counterexample for learning $\mathscr{H} = \{x \mapsto wx : |w| \le 1\}$ with square loss: • Exactly the same as before, but scale by $1/\mu$: $z_1 = (1/\mu, 0)$ $z_2 = (1, -1)$

- Remember our bounds on \mathfrak{R}_n for linear classes depended on bounding the norm of $h, \mathcal{H} = \{w \in \mathbb{R}^d : ||w|| \le B\}$
 - Optimal w on \mathscr{D}_2 is $-1/\mu \approx -200n$, which is really big

• Counterexample for learning $\mathscr{H} = \{x \mapsto wx : |w| \le 1\}$ with square loss: • Exactly the same as before, but scale by $1/\mu$: $z_1 = (1/\mu, 0)$ $z_2 = (1, -1)$ • wx on the old problem corresponds to $\frac{w}{\mu}(\mu x)$ here; same loss

- Remember our bounds on \Re_n for linear classes depended on bounding the norm of $h, \mathcal{H} = \{w \in \mathbb{R}^d : ||w|| \le B\}$
 - Optimal w on \mathscr{D}_2 is $-1/\mu \approx -200n$, which is really big
- Counterexample for learning $\mathscr{H} = \{x \mapsto wx : |w| \le 1\}$ with square loss: • Exactly the same as before, but scale by $1/\mu$: $z_1 = (1/\mu, 0)$ $z_2 = (1, -1)$ • wx on the old problem corresponds to $\frac{w}{u}(\mu x)$ here; same loss
- - Now \mathcal{H} is bounded, but we have really big x values

- Remember our bounds on \Re_n for linear classes depended on bounding the norm of $h, \mathcal{H} = \{w \in \mathbb{R}^d : ||w|| \le B\}$
 - Optimal w on \mathscr{D}_2 is $-1/\mu \approx -200n$, which is really big
- Counterexample for learning $\mathscr{H} = \{x \mapsto wx : |w| \le 1\}$ with square loss: • Exactly the same as before, but scale by $1/\mu$: $z_1 = (1/\mu, 0)$ $z_2 = (1, -1)$ • wx on the old problem corresponds to $\frac{w}{u}(\mu x)$ here; same loss

 - Now \mathcal{H} is bounded, but we have really big x values • oh, yeah...our bound on \Re_n required bounding ||x|| too!

Some learnable classes

- A learning problem $(\mathcal{H}, \mathcal{I}, \ell)$ Convex ρ -Lipschitz *B*-Bounded if
 - $\mathcal{H} \subset \mathbb{R}^d$ is a convex set, with $||w|| \leq B$ for all $w \in \mathcal{H}$
 - for each $z \in \mathcal{X}$, $\ell(\cdot, z)$ is a convex, ρ -Lipschitz function

l: (w) = lw xi - Yi

 $\begin{array}{c}
 l_{c}(w) = (w \times (-\gamma)^{2})^{2} \\
 l_{c}^{c} = 2(w \times (-\gamma)^{2})^{2} \\
 l_{c}^{c} = 2(w \times (-\gamma)^{2})^{2} \\
 v_{c}^{c} & v_{c}^{c} \\
 v_{c}^{c} & v_{c}^$

Some learnable classes

- A learning problem $(\mathcal{H}, \mathcal{I}, \ell)$ Convex ρ -Lipschitz *B*-Bounded if
 - $\mathcal{H} \subset \mathbb{R}^d$ is a convex set, with $||w|| \leq B$ for all $w \in \mathcal{H}$
 - for each $z \in \mathcal{X}$, $\ell(\cdot, z)$ is a convex, ρ -Lipschitz function
- A learning problem $(\mathcal{H}, \mathcal{I}, \ell)$ Convex β -Smooth B-Bounded if
 - $\mathcal{H} \subset \mathbb{R}^d$ is a convex set, with $||w|| \leq B$ for all $w \in \mathcal{H}$
 - for each $z \in \mathcal{X}$, $\ell(\cdot, z)$ is a convex, nonnegative, β -smooth function • A function f is β -smooth if ∇f is β -Lipschitz



Some learnable classes

- A learning problem $(\mathcal{H}, \mathcal{I}, \ell)$ Convex ρ -Lipschitz *B*-Bounded if
 - $\mathcal{H} \subset \mathbb{R}^d$ is a convex set, with $||w|| \leq B$ for all $w \in \mathcal{H}$
 - for each $z \in \mathcal{Z}$, $\ell(\cdot, z)$ is a convex, ρ -Lipschitz function
- A learning problem $(\mathcal{H}, \mathcal{I}, \ell)$ Convex β -Smooth B-Bounded if
 - $\mathcal{H} \subset \mathbb{R}^d$ is a convex set, with $||w|| \leq B$ for all $w \in \mathcal{H}$
 - for each $z \in \mathcal{X}$, $\ell(\cdot, z)$ is a convex, nonnegative, β -smooth function • A function f is β -smooth if ∇f is β -Lipschitz
- We'll see soon that these classes are always learnable (efficiently!)

$$\begin{aligned} & \text{Key property of } \beta \text{-smooth } f \\ f(y) - f(x) &= \int_0^1 \nabla f \left(ty + (1-t)x \right) \cdot (y-x) \, dt \\ &= \int_0^1 \left[\nabla f(x) + \nabla f \left(ty + (1-t)x \right) - \nabla f(x) \right] \cdot (y-x) \, dt \end{aligned}$$

$$\leq \nabla f(x) \cdot (y-x) + \int_0^1 \left\| \nabla f \left(ty + (1-t)x \right) - \nabla f(x) \right\| \|y-x\| \, dt \\ &\leq \nabla f(x) \cdot (y-x) + \int_0^1 \beta \| ty + (1-t)x - x \| \|y-x\| \, dt \end{aligned}$$

$$= \nabla f(x) \cdot (y-x) + \beta \|y-x\|^2 \int_0^1 t \, dt \qquad A \succeq 0 \text{ means } f \text{ is } p \text{ id } A \not\equiv g \text{ means } A - B \succeq 0 \\ &= \nabla f(x) \cdot (y-x) + \frac{1}{2} \beta \|y-x\|^2 \end{aligned} \text{ (implies } \nabla^2 f(y) \leq \beta I \text{, if it exists)} \end{aligned}$$

(pause)

• Start at $w^{(0)}$, maybe 0 or sampled randomly. (SSBD calls this $w^{(1)}$)

- Start at $w^{(0)}$, maybe 0 or sampled randomly. (SSBD calls this $w^{(1)}$)

• Steps are $w^{(t+1)} = w^{(t)} - \eta \nabla f(w^{(t)})$, where $\eta > 0$ is a learning rate

- Start at $w^{(0)}$, maybe 0 or sampled randomly. (SSBD calls this $w^{(1)}$)
- Steps are $w^{(t+1)} = w^{(t)} \eta \nabla f(w^{(t)})$, where $\eta > 0$ is a learning rate
- Output last iterate $w^{(T)}$



 $\hat{f}_{i,j}(v) = f(w) + \nabla f(w) \cdot (v - w) + \frac{1}{2y} \|vw\|^{2}$





- Start at $w^{(0)}$, maybe 0 or sampled randomly. (SSBD calls this $w^{(1)}$)
- Steps are $w^{(t+1)} = w^{(t)} \eta \nabla f(w^{(t)})$, where $\eta > 0$ is a learning rate
- Output last iterate $w^{(T)}$
 - Not the only choice:

- Start at $w^{(0)}$, maybe 0 or sampled randomly. (SSBD calls this $w^{(1)}$)
- Output last iterate $w^{(T)}$
 - Not the only choice:

SSBD use average of iterate

Sometimes tail average $\frac{2}{T}$

- Sometimes best iterate: arg
- Or best on a validation set: a

• Steps are $w^{(t+1)} = w^{(t)} - \eta \nabla f(w^{(t)})$, where $\eta > 0$ is a learning rate

es,
$$\bar{w} = \frac{1}{T} \sum_{t=0}^{T-1} w^{(t)}$$

 $\sum_{t=0}^{T} w^{(t)}$
 $= T/2$
 $\min_{w^{(t)}:t\in[T]} f(w^{(t)})$
 $\operatorname{argmin}_{w^{(t)}:t\in[T]} L_V(w^{(t)})$

- SSBD section 14.1 analyzes:
 - average iterate, initialize at 0, Lipschitz f
 - very particular fixed η that depends on length of optimization T (also B, ρ)

- SSBD section 14.1 analyzes: average iterate, initialize at 0, Lipschitz fvery particular fixed η that depends on length of optimization T (also B, ρ)
- We'll do something else (more standard): last iterate, β -smooth f, fixed $\eta \leq 1/\beta$

- SSBD section 14.1 analyzes: average iterate, initialize at 0, Lipschitz fvery particular fixed η that depends on length of optimization T (also B, ρ)
- We'll do something else (more standard): last iterate, β -smooth f, fixed $\eta \leq 1/\beta$ We'll prove: $f(w^{(T)}) - f^* \leq \frac{\|w^{(0)} - w^*\|}{2\eta T}$

- SSBD section 14.1 analyzes: average iterate, initialize at 0, Lipschitz fvery particular fixed η that depends on length of optimization T (also B, ρ)
- We'll do something else (more standard): last iterate, β -smooth f, fixed $\eta \leq 1/\beta$ We'll prove: $f(w^{(T)}) - f^* \leq \frac{\|w^{(0)} - w^*\|}{2\eta T}$
 - so can get suboptimality ε in $\mathcal{O}(1/\varepsilon)$ steps

- SSBD section 14.1 analyzes: average iterate, initialize at 0, Lipschitz fvery particular fixed η that depends on length of optimization T (also B, ρ)
- $\left\| \left(\kappa \left(\kappa \right) w \left(\kappa \right) \right\| \right\| = \left\| \nabla f \left(w \kappa \right) \right\|$ • We'll do something else (more standard): last iterate, β -smooth *f*, fixed $\eta \leq 1/\beta$ • We'll prove: $f(w^{(T)}) - f^* \le \frac{\|w^{(0)} - w^*\|}{2\eta T}$
 - so can get suboptimality ε in $\mathcal{O}(1/\varepsilon)$ steps
 - In practice, can be hard to compute β , and $1/\beta$ usually too small



- SSBD section 14.1 analyzes: average iterate, initialize at 0, Lipschitz fvery particular fixed η that depends on length of optimization T (also B, ρ)
- We'll do something else (more standard): last iterate, β -smooth f, fixed $\eta \leq 1/\beta$ We'll prove: $f(w^{(T)}) - f^* \leq \frac{\|w^{(0)} - w^*\|}{2\eta T}$
 - so can get suboptimality ε in $\mathcal{O}(1/\varepsilon)$ steps
 - In practice, can be hard to compute β , and $1/\beta$ usually too small
 - Backtracking line search has a similar rate

 β , and $1/\beta$ usually too small nilar rate

- SSBD section 14.1 analyzes: average iterate, initialize at 0, Lipschitz fvery particular fixed η that depends on length of optimization T (also B, ρ)
- We'll do something else (more standard): last iterate, β -smooth f, fixed $\eta \leq 1/\beta$ We'll prove: $f(w^{(T)}) - f^* \leq \frac{\|w^{(0)} - w^*\|}{2\eta T}$
 - so can get suboptimality ε in $\mathcal{O}(1/\varepsilon)$ steps
 - In practice, can be hard to compute β , and $1/\beta$ usually too small
 - Backtracking line search has a similar rate

 β , and $1/\beta$ usually too small nilar rate

Springer Series in Operations Research

Jorge Nocedal Stephen J. Wright

Numerical Optimization Second Edition



2 Springer



• *f* is convex, ∇f is β -Lipschitz, learning rate $\eta \leq 1/\beta$

- f is convex, ∇f is β -Lipschitz, learning
- β -smooth functions have $f(v) \leq f(v)$

ng rate
$$\eta \leq 1/\beta$$

$$f(w) + \nabla f(w)^{\mathsf{T}}(v-w) + \frac{\beta}{2} \|v-w\|^2$$

- f is convex, ∇f is β -Lipschitz, learning
- β -smooth functions have $f(v) \leq$
- Iterate w goes to $w^+ = w \eta \nabla f(w)$; plugging in to above, get

ng rate
$$\eta \leq 1/\beta$$

$$f(w) + \nabla f(w)^{\mathsf{T}}(v-w) + \frac{\beta}{2} \|v-w\|^2$$

- f is convex, ∇f is β -Lipschitz, learning
- β -smooth functions have $f(v) \leq f(v)$
- Iterate w goes to $w^+ = w \eta \nabla f(w)$; plugging in to above, get

 $f(w^+) \le f(w) + \nabla f(w)^\top (w^+ - w)$

ng rate
$$\eta \leq 1/\beta$$

$$f(w) + \nabla f(w)^{\mathsf{T}}(v-w) + \frac{\beta}{2} \|v-w\|^2$$

$$(w) + \frac{\beta}{2} ||w^+ - w||^2$$

- f is convex, ∇f is β -Lipschitz, learning
- β -smooth functions have $f(v) \leq f(v)$
- Iterate w goes to $w^+ = w \eta \nabla f(w)$; plugging in to above, get
 - $f(w^+) \le f(w) + \nabla f(w)^\top (w^+ w)$ $= f(w) + \nabla f(w)^{\mathsf{T}}(-\eta \nabla f(v))^{\mathsf{T}}(-\eta \nabla f(v))$

ng rate
$$\eta \leq 1/\beta$$

$$f(w) + \nabla f(w)^{\mathsf{T}}(v-w) + \frac{\beta}{2} \|v-w\|^2$$

$$w) + \frac{\beta}{2} \|w^{+} - w\|^{2}$$
$$w)) + \frac{\beta}{2} \|-\eta \nabla f(w)\|^{2}$$

- f is convex, ∇f is β -Lipschitz, learning
- β -smooth functions have $f(v) \leq$
- Iterate w goes to $w^+ = w \eta \nabla f(w)$; plugging in to above, get
 - $f(w^+) \le f(w) + \nabla f(w)^\top (w^+ w)$ $= f(w) + \nabla f(w)^{\top} (-\eta \nabla f(w))^{\top}$ $= f(w) - \eta \|\nabla f(w)\|^2 + \frac{\beta \eta}{\beta}$

ng rate
$$\eta \leq 1/\beta$$

$$f(w) + \nabla f(w)^{\mathsf{T}}(v-w) + \frac{\beta}{2} \|v-w\|^2$$

$$\begin{split} w) + \frac{\beta}{2} \|w^{+} - w\|^{2} \\ (w)) + \frac{\beta}{2} \|-\eta \nabla f(w)\|^{2} \\ \frac{\beta \eta^{2}}{2} \|\nabla f(w)\|^{2} \end{split}$$

• f is convex, ∇f is β -Lipschitz, learning β -smooth functions have $f(v) \leq$ • Iterate w goes to $w^+ = w - \eta \nabla f(w)$; plugging in to above, get $f(w^+) \leq f(w) + \nabla f(w)^\top (w^+ - v)$ $= f(w) + \nabla f(w)^{\top} (-\eta \nabla f(w))^{\top}$ $= f(w) - \eta \|\nabla f(w)\|^2 + \frac{\beta \eta}{2}$ $= f(w) - \eta \left(1 - \frac{\beta \eta}{2}\right) \|\nabla f(w)\|^2$

ng rate
$$\eta \leq 1/\beta$$

$$f(w) + \nabla f(w)^{\mathsf{T}}(v-w) + \frac{\beta}{2} \|v-w\|^2$$

$$w) + \frac{\beta}{2} ||w^{+} - w||^{2}$$

$$(w)) + \frac{\beta}{2} ||-\eta \nabla f(w)||^{2}$$

$$\frac{\beta \eta^{2}}{2} ||\nabla f(w)||^{2}$$

$$7 f(w) ||^{2}$$

• f is convex, ∇f is β -Lipschitz, learning • β -smooth functions have $f(v) \leq$ • Iterate w goes to $w^+ = w - \eta \nabla f(w)$; plugging in to above, get $f(w^+) \leq f(w) + \nabla f(w)^\top (w^+ - v)$ $= f(w) + \nabla f(w)^{\top} (-\eta \nabla f(v))^{\top}$ $= f(w) - \eta \|\nabla f(w)\|^2 + \frac{\beta \eta}{-1}$ $= f(w) - \eta \left(1 - \frac{\beta \eta}{2}\right) \|\nabla f(w)\|^2$ $\leq f(w) - \frac{\eta}{2} \|\nabla f(w)\|^2$

ng rate
$$\eta \leq 1/\beta$$

$$f(w) + \nabla f(w)^{\mathsf{T}}(v-w) + \frac{\beta}{2} \|v-w\|^2$$

$$w) + \frac{\beta}{2} ||w^{+} - w||^{2}$$

$$(w)) + \frac{\beta}{2} ||-\eta \nabla f(w)||^{2}$$

$$\frac{\beta \eta^{2}}{2} ||\nabla f(w)||^{2}$$

$$7 f(w) ||^{2}$$

• f is convex, ∇f is β -Lipschitz, learning • β -smooth functions have $f(v) \leq$ • Iterate w goes to $w^+ = w - \eta \nabla f(w)$; plugging in to above, get $f(w^+) \leq f(w) + \nabla f(w)^\top (w^+ - v)$ $= f(w) + \nabla f(w)^{\top} (-\eta \nabla f(v))^{\top}$ $= f(w) - \eta \|\nabla f(w)\|^2 + \frac{\beta \eta}{-1}$ $= f(w) - \eta \left(1 - \frac{\beta \eta}{2}\right) \|\nabla f(w)\|^2$ $\leq f(w) - \frac{\eta}{\gamma} \|\nabla f(w)\|^2$

ng rate
$$\eta \leq 1/\beta$$

$$f(w) + \nabla f(w)^{\mathsf{T}}(v-w) + \frac{\beta}{2} \|v-w\|^2$$

$$\begin{split} w) + \frac{\beta}{2} \|w^{+} - w\|^{2} \\ (w)) + \frac{\beta}{2} \|-\eta \nabla f(w)\|^{2} \\ \frac{\beta \eta^{2}}{2} \|\nabla f(w)\|^{2} \\ \nabla f(w)\|^{2} \end{split}$$

"descent lemma": we're decreasing the objective! (note: didn't use convexity yet...)



• f is convex, ∇f is β -Lipschitz, learning • β -smooth functions have $f(v) \leq$ • Iterate w goes to $w^+ = w - \eta \nabla f(w)$; plugging in to above, get $f(w^+) \le f(w) + \nabla f(w)^\top (w^+ - w)^\top (w^$ $= f(w) + \nabla f(w)^{\mathsf{T}}(-\eta \nabla f(w))^{\mathsf{T}}(-\eta \nabla$ $= f(w) - \eta \|\nabla f(w)\|^2 + f'$ $\leq f(w) - \frac{\eta}{2} \|\nabla f(w)\|^2$

ng rate
$$\eta \leq 1/\beta$$

$$f(w) + \nabla f(w)^{\mathsf{T}}(v-w) + \frac{\beta}{2} \|v-w\|^2$$

$$w) + \frac{\beta}{2} ||w^{+} - w||^{2}$$

$$f(w)) + \frac{\beta}{2} ||-\eta \nabla f(w)||^{2}$$

$$\frac{\beta \eta^{2}}{2} ||\nabla f(w)||^{2}$$

 $= f(w) - \eta \left(1 - \frac{\beta \eta}{2} \right) \|\nabla f(w)\|^2 \qquad \text{first-order convexity condition:}$ $f(w^*) \ge f(w) + \nabla f(w)^{\mathsf{T}}(w^* - w)$ $\operatorname{so} f(w) \le f(w^*) + \nabla f(w)^{\mathsf{T}}(w - w^*)$ "descent lemma": we're decreasing the objective! (note: didn't use convexity yet...)



• f is convex, ∇f is β -Lipschitz, learning • β -smooth functions have $f(v) \leq$ • Iterate w goes to $w^+ = w - \eta \nabla f(w)$; plugging in to above, get $f(w^+) \leq f(w) + \nabla f(w)^\top (w^+ - v)$ $= f(w) + \nabla f(w)^{\mathsf{T}}(-\eta \nabla f(w))$ $= f(w) - \eta \|\nabla f(w)\|^2 + \frac{\beta}{2}$ $= f(w) - \eta \left(1 - \frac{\beta \eta}{2}\right) \|\nabla f(w)\|^2$ $\leq f(w) - \frac{\eta}{\gamma} \|\nabla f(w)\|^2$ $\leq f(w^*) + \nabla f(w)^{\mathsf{T}}(w - w^*) - \frac{\eta}{2} \|\nabla f(w)\|^2$

ng rate
$$\eta \leq 1/\beta$$

$$f(w) + \nabla f(w)^{\mathsf{T}}(v-w) + \frac{\beta}{2} \|v-w\|^2$$

$$w) + \frac{\beta}{2} ||w^{+} - w||^{2}$$

(w)) + $\frac{\beta}{2} ||-\eta \nabla f(w)||^{2}$
 $3\eta^{2} ||\nabla f(w)||^{2}$

first-order convexity condition: $f(w^*) \ge f(w) + \nabla f(w)^{\mathsf{T}}(w^* - w)$ $\operatorname{so} f(w) \le f(w^*) + \nabla f(w)^{\mathsf{T}}(w - w^*)$ "descent lemma": we're decreasing the objective! (note: didn't use convexity yet...)


ing rate
$$\eta \leq 1/\beta; w^+ = w - \eta \nabla f(w)$$

 $-w^*) - \frac{\eta}{2} \|\nabla f(w)\|^2$

 $2\eta \nabla f(w)^{\mathsf{T}}(w - w^*) - \eta^2 \|\nabla f(w)\|^2$

ing rate
$$\eta \leq 1/\beta; w^+ = w - \eta \nabla f(w)$$

 $-w^*) - \frac{\eta}{2} \|\nabla f(w)\|^2$

$$2\eta \nabla f(w)^{\mathsf{T}}(w - w^*) - \eta^2 \|\nabla f(w)\|^2$$
$$= \|w - w^*\|^2 - \|w - w^*\|^2$$

ing rate
$$\eta \leq 1/\beta; w^+ = w - \eta \nabla f(w)$$

 $-w^*) - \frac{\eta}{2} \|\nabla f(w)\|^2$

$+ 2\eta \nabla f(w)^{\mathsf{T}}(w - w^*) - \eta^2 \|\nabla f(w)\|^2$

$$2\eta \nabla f(w)^{\mathsf{T}}(w - w^*) - \eta^2 \|\nabla f(w)\|^2$$

= $\|w - w^*\|^2 - \|w - w^*\|^2 + 2\eta \nabla f(w)^{\mathsf{T}}(w - w^*) - \eta^2 \|\nabla f(w)\|^2$
= $\|w - w^*\|^2 - [\|w - w^*\|^2 - 2\eta \nabla f(w)^{\mathsf{T}}(w - w^*) + \eta^2 \|\nabla f(w)\|^2]$

ing rate
$$\eta \leq 1/\beta; w^+ = w - \eta \nabla f(w)$$

 $-w^*) - \frac{\eta}{2} \|\nabla f(w)\|^2$

$$2\eta \nabla f(w)^{\mathsf{T}}(w - w^*) - \eta^2 \|\nabla f(w)\|^2$$

= $\|w - w^*\|^2 - \|w - w^*\|^2 + 2\eta \nabla f(w)^{\mathsf{T}}(w - w^*) - \eta^2 \|\nabla f(w)\|^2$
= $\|w - w^*\|^2 - [\|w - w^*\|^2 - 2\eta \nabla f(w)^{\mathsf{T}}(w - w^*) + \eta^2 \|\nabla f(w)\|^2]$
= $\|w - w^*\|^2 - \|(w - w^*) - \eta \nabla f(w)\|^2$

ing rate
$$\eta \leq 1/\beta$$
; $w^+ = w - \eta \nabla f(w)$
 $-w^*) - \frac{\eta}{2} \|\nabla f(w)\|^2$

$$2\eta \nabla f(w)^{\mathsf{T}}(w - w^*) - \eta^2 \|\nabla f(w)\|^2$$

= $\|w - w^*\|^2 - \|w - w^*\|^2 + 2\eta \nabla f(w)^{\mathsf{T}}(w - w^*) - \eta^2 \|\nabla f(w)\|^2$
= $\|w - w^*\|^2 - [\|w - w^*\|^2 - 2\eta \nabla f(w)^{\mathsf{T}}(w - w^*) + \eta^2 \|\nabla f(w)\|^2]$
= $\|w - w^*\|^2 - \|(w - w^*) - \eta \nabla f(w)\|^2$
= $\|w - w^*\|^2 - \|w^* - w^*\|^2$

ing rate
$$\eta \leq 1/\beta; w^+ = w - \eta \nabla f(w)$$

 $-w^*) - \frac{\eta}{2} \|\nabla f(w)\|^2$

$$2\eta \nabla f(w)^{\mathsf{T}}(w - w^*) - \eta^2 \|\nabla f(w)\|^2$$

= $\|w - w^*\|^2 - \|w - w^*\|^2 + 2\eta \nabla f(w)^{\mathsf{T}}(w - w^*) - \eta^2 \|\nabla f(w)\|^2$
= $\|w - w^*\|^2 - [\|w - w^*\|^2 - 2\eta \nabla f(w)^{\mathsf{T}}(w - w^*) + \eta^2 \|\nabla f(w)\|^2]$
= $\|w - w^*\|^2 - \|(w - w^*) - \eta \nabla f(w)\|^2$
= $\|w - w^*\|^2 - \|w^+ - w^*\|^2$
Plugging back in, $f(w^+) \le f(w^*) + \frac{1}{2\eta} (\|w - w^*\|^2 - \|w^+ - w^*\|^2)$

ing rate
$$\eta \leq 1/\beta$$
; $w^+ = w - \eta \nabla f(w)$
 $-w^*) - \frac{\eta}{2} \|\nabla f(w)\|^2$

$$\begin{aligned} &2\eta \, \nabla f(w)^{\mathsf{T}}(w - w^*) - \eta^2 \| \nabla f(w) \|^2 \\ &= \|w - w^*\|^2 - \|w - w^*\|^2 + 2\eta \, \nabla f(w)^{\mathsf{T}}(w - w^*) - \eta^2 \| \nabla f(w) \|^2 \\ &= \|w - w^*\|^2 - \left[\|w - w^*\|^2 - 2\eta \, \nabla f(w)^{\mathsf{T}}(w - w^*) + \eta^2 \| \nabla f(w) \|^2 \right] \\ &= \|w - w^*\|^2 - \|(w - w^*) - \eta \, \nabla f(w)\|^2 \\ &= \|w - w^*\|^2 - \|w^+ - w^*\|^2 \end{aligned}$$
Plugging back in, $f(w^+) \leq f(w^*) + \frac{1}{2\eta} \left(\|w - w^*\|^2 - \|w^+ - w^*\|^2 \right)$
and so $f(w^{(k)}) - f(w^*) \leq \frac{1}{2n} \left(\|w^{(k-1)} - w^*\|^2 - \|w^{(k)} - w^*\|^2 \right)$

~//

ing rate
$$\eta \leq 1/\beta$$
; $w^+ = w - \eta \nabla f(w)$
 $-w^*) - \frac{\eta}{2} \|\nabla f(w)\|^2$

f is convex, ∇f is β -Lipschitz, learni now know $f(w^{(k)}) - f(w^*) \le \frac{1}{2\eta} \left(\parallel \frac{1}{2\eta} \right)$

ing rate
$$\eta \le 1/\beta$$
; $w^+ = w - \eta \nabla f(w)$
 $|w^{(k-1)} - w^*||^2 - ||w^{(k)} - w^*||^2)$

f is convex, ∇f is β -Lipschitz, learning $\operatorname{now}\operatorname{know} f(w^{(k)}) - f(w^*) \le \frac{1}{2n} \left(\|$ $f(w^{(T)}) - f(w^*) \le \frac{1}{T} \sum_{k=1}^{T} \left(f(w^{(k)}) - f(w^*) \right)$

ing rate
$$\eta \le 1/\beta$$
; $w^+ = w - \eta \nabla f(w)$
 $|w^{(k-1)} - w^*||^2 - ||w^{(k)} - w^*||^2)$

(because of the descent lemma)



f is convex, ∇f is β -Lipschitz, learning $\operatorname{now}\operatorname{know} f(w^{(k)}) - f(w^*) \le \frac{1}{2n} \left(\|$ $f(w^{(T)}) - f(w^*) \le \frac{1}{T} \sum_{k=1}^{T} \left(f(w^{(k)}) - f(w^*) \right)$ (because of the descent lemma) k = I $\leq \frac{1}{2\eta T} \sum_{k=1}^{T} \left(\| w^{(k-1)} \|_{k=1}^{\infty} \right)^{k-1}$

ing rate
$$\eta \le 1/\beta$$
; $w^+ = w - \eta \nabla f(w)$
 $|w^{(k-1)} - w^*||^2 - ||w^{(k)} - w^*||^2)$

$$^{1)} - w^* \|^2 - \|w^{(k)} - w^*\|^2)$$

$$- (4) - 2 + 2 - \cdots + 7 - 7 - 1$$



f is convex, ∇f is β -Lipschitz, learn $\operatorname{now\,know} f(w^{(k)}) - f(w^*) \le \frac{1}{2n} \left(\parallel \right)$ $f(w^{(T)}) - f(w^*) \le \frac{1}{T} \sum_{T}^{T} \left(f(w^{(k)}) - f(w^*) \right)$ (because of the descent lemma) k = 1 $\leq \frac{1}{2\eta T} \sum_{k=1}^{T} \left(\|w^{(k-1)}\|_{k=1}^{T} \right) \\ = \frac{1}{2\eta T} \left(\|w^{(0)} - w^{(0)}\|_{k=1}^{T} \right)$

ing rate
$$\eta \le 1/\beta$$
; $w^+ = w - \eta \nabla f(w)$
 $|w^{(k-1)} - w^*||^2 - ||w^{(k)} - w^*||^2)$

$$|w^{(k)} - w^{(k)}|^2 - ||w^{(k)} - w^{(k)}|^2$$

$$v^* \|^2 - \|w^{(T)} - w^*\|^2$$



f is convex, ∇f is β -Lipschitz, learn $\operatorname{now\,know} f(w^{(k)}) - f(w^*) \le \frac{1}{2n} \left(\parallel \right)$ $f(w^{(T)}) - f(w^*) \le \frac{1}{T} \sum_{T}^{T} \left(f(w^{(k)}) - f(w^*) \right)$ $\leq \frac{1}{2\eta T} \sum_{k=1}^{T} \left(\|w^{(k-1)} - w^{(k-1)}\| \right)$ $= \frac{1}{2\eta T} \left(\|w^{(0)} - w^{(k-1)}\| \right)$ $\leq \frac{1}{2nT} \|w^{(0)} - w^*\|^2$

ing rate
$$\eta \le 1/\beta$$
; $w^+ = w - \eta \nabla f(w)$
 $|w^{(k-1)} - w^*||^2 - ||w^{(k)} - w^*||^2)$

(because of the descent lemma)

$$(1) - w^* \|^2 - \|w^{(k)} - w^*\|^2$$

$$v^* \|^2 - \|w^{(T)} - w^*\|^2$$



f is convex, ∇f is β -Lipschitz, learn $\operatorname{now}\operatorname{know} f(w^{(k)}) - f(w^*) \le \frac{1}{2n} \left(\|$ $f(w^{(T)}) - f(w^*) \le \frac{1}{T} \sum_{k=1}^{T} \left(f(w^{(k)}) - f(w^*) \right)$ (because of the descent lemma) $\leq \frac{1}{2\eta T} \sum_{k=1}^{T} \left(\|w^{(k-1)}\|_{k=1}^{T} \right) \\ = \frac{1}{2\eta T} \left(\|w^{(0)} - w^{(0)}\|_{k=1}^{T} \right)$ $\leq \frac{1}{2nT} \|w^{(0)} - w^*\|^2 \leq \frac{B\beta}{2T} \text{ with } \eta = \frac{1}{\beta}$

ing rate
$$\eta \le 1/\beta$$
; $w^+ = w - \eta \nabla f(w)$
 $|w^{(k-1)} - w^*||^2 - ||w^{(k)} - w^*||^2)$

$$|w^{(k)} - w^{(k)}|^2 - ||w^{(k)} - w^{(k)}|^2$$

$$x^* \|^2 - \| w^{(T)} - w^* \|^2$$

$$< \frac{B\beta}{M} \text{ with } n = -\frac{1}{2}$$



f is convex, ∇f is β -Lipschitz, learn $\operatorname{now\,know} f(w^{(k)}) - f(w^*) \le \frac{1}{2n} \left(\parallel \right)$ $f(w^{(T)}) - f(w^*) \le \frac{1}{T} \sum_{T}^{T} \left(f(w^{(k)}) - f(w^*) \right)$ $\leq \frac{1}{2\eta T} \sum_{k=1}^{T} \left(\|w^{(k-1)}\|_{k=1}^{T} \right) \\ = \frac{1}{2\eta T} \left(\|w^{(0)} - w^{(0)}\|_{k=1}^{T} \right)$ $\leq \frac{1}{2nT} \| w^{(0)} - w^* \|$ note $f\left(\frac{1}{T}\sum_{k=1}^{T}w^{(k)}\right) \le \frac{1}{T}\sum_{k=1}^{T}f(w^{(k)})$: same for average iterate

ing rate
$$\eta \le 1/\beta$$
; $w^+ = w - \eta \nabla f(w)$
 $|w^{(k-1)} - w^*||^2 - ||w^{(k)} - w^*||^2)$

(because of the descent lemma)

$$|w^{(k)} - w^{(k)}|^2 - ||w^{(k)} - w^{(k)}|^2$$

$$\|w^*\|^2 - \|w^{(T)} - w^*\|^2$$

$$\|^{2} \leq \frac{\beta \rho}{2T} \text{ with } \eta = \frac{1}{\beta}$$



- Defined Convex-Lipschitz-Bounded and Convex-Smooth-Bounded classes
 - Will show they're learnable; didn't quite get there yet

nd Convex-Smooth-Bounded classes uite get there yet

- Defined Convex-Lipschitz-Bounded and Convex-Smooth-Bounded classes
 - Will show they're learnable; didn't quite get there yet

nd Convex-Smooth-Bounded classes uite get there yet

- Defined Convex-Lipschitz-Bounded and Convex-Smooth-Bounded classes • Will show they're learnable; didn't quite get there yet
- We can ~efficiently optimize empirical risk for Convex-Smooth-Bounded classes • SSBD section 14.1 shows (worse) rate for Convex-Lipschitz-Bounded, avg iterate Turns out can get much faster rate if objective is strongly convex Don't actually need differentiability everywhere: subgradient descent For more: check out <u>Bubeck (Convex Optimization: Algorithms and Complexity)</u>

 - Nocedal and Wright (Numerical Optimization), or take CPSC 536M





- Defined Convex-Lipschitz-Bounded and Convex-Smooth-Bounded classes
 - Will show they're learnable; didn't quite get there yet
- - Turns out can get much faster rate if objective is strongly convex
 - Don't actually need differentiability everywhere: subgradient descent
 - Nocedal and Wright (Numerical Optimization), or take CPSC 536M
- Next time:
 - How to use this / related stuff for learning guarantees
 - How to analyze classifiers (since 0-1 loss isn't Lipschitz or smooth)

 We can ~efficiently optimize empirical risk for Convex-Smooth-Bounded classes • SSBD section 14.1 shows (worse) rate for Convex-Lipschitz-Bounded, avg iterate For more: check out <u>Bubeck (Convex Optimization: Algorithms and Complexity)</u>



