

First-order methods

Convexity



10-725 Optimization
Geoff Gordon
Ryan Tibshirani

Administrivia



- HW1 out, due 9/20
 - ▶ in class, at **beginning** of class—no skipping lecture to keep working on it ;-)
 - ▶ if you use late days, hand in to course assistant before 1:30PM on due date + k days
- Reminder: think about project teams and project proposals (due 9/25)

Review



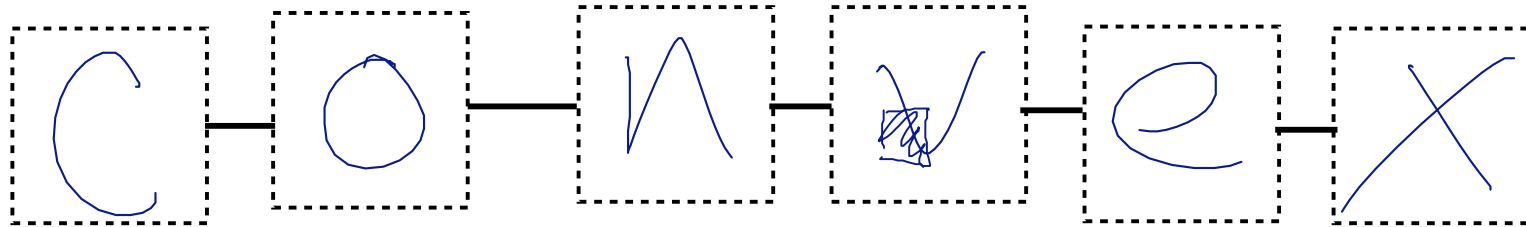
- Convex sets
 - ▶ primal (convex hull) vs. dual (intersect hyperplanes)
 - ▶ supporting, separating hyperplanes
 - ▶ operations that preserve convexity
 - ▶ affine fn; perspective
 - ▶ open/closed/compact

Review



- Convex functions
 - ▶ epigraph
 - ▶ domain
 - ▶ sublevel sets; quasiconvexity
 - ▶ first order, 2nd order conditions
 - ▶ operations that preserve convexity
 - ▶ perspective; minimization over one argument

Ex: structured classifier



x_i pixels of char i

$\phi_j(x_i)$ feature of a char

y_i a ... z

$\Psi_{ijk}(x_i, y_i)$

$\phi_j(x_i) \delta(y_i = k) \leftarrow w_{jk}$

$\chi_{ikl}(y_i, y_{i+1})$

$\delta(y_i = k) \delta(y_{i+1} = l) \leftarrow v_{kl}$

$$L(x, y; v, w) = \sum_{ijk} \phi_{ijk}(x_i, y_i) w_{ijk} + \sum_{ikl} \chi_{ikl}(y_i, y_{i+1}) v_{ikl}$$

Classifier:

$$y = \arg \max_y L(x, y; v, w)$$

Learning structured classifier

- Get it right if: $L(x, y; v, w) > L(x, y', v, w)$ $\forall y' \neq y$
- So, want: $L(x, y; v, w) \geq \max_{y'} (L(x, y'; v, w) + \pi(y, y'))$
- Where $\pi(y, y') = \begin{cases} 0 & y = y' \\ > 0 & y \neq y' \end{cases}$
- RHS: convex in v, w
- RHS – LHS: convex
- Train: lots of pairs (x^t, y^t) $\min_{v, w} \sum_t (RHS^t - LHS^t) + C(\|v\|^2 + \|w\|^2)$

Strict convexity; strong convexity

- Strictly convex:

$x, y \in \text{dom } f$ $t \in (0, 1)$
 $x \neq y$

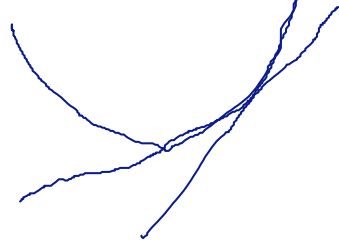
$f(tx + (1-t)y) < t f(x) + (1-t)f(y)$

$\ln(1+e^x)$



- k-strongly convex:

$f(y) \geq f(x) + (y-x)f'(x) + \frac{k}{2}(y-x)^2$



$|x| + x^2 \Rightarrow 2\text{-strongly}$

Extended reals

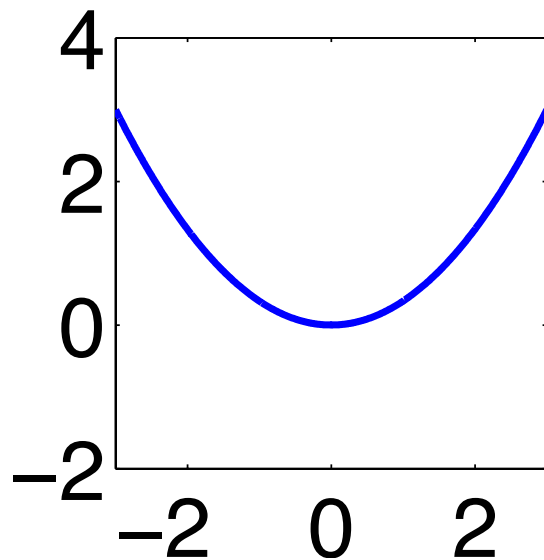
- Suppose $\text{dom } f \subset \mathbb{R}^n$
- Define $g(x) = \begin{cases} f(x) & x \in \text{dom } f \\ \infty & \text{o/w} \end{cases}$
- f convex \iff g convex
 - ▶ $f(tx+(1-t)y) \leq tf(x) + (1-t)f(y)$ ←
 - ▶ $g(tx+(1-t)y) \leq tg(x) + (1-t)g(y)$ ←
 - ▶ cases: $x, y \in \text{dom } f$
 $x, y \notin \text{dom } f$
 $x \in \text{dom } f \quad y \notin \text{dom } f$

Lipschitz

- Function $f(x)$ is Lipschitz (in norm $\|x\|$) if:

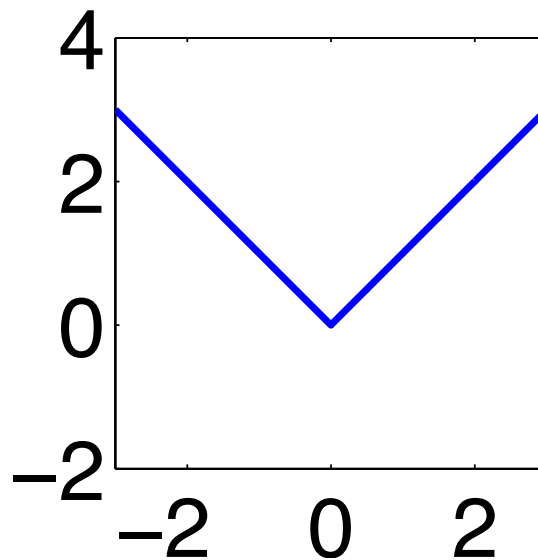
- ▶ $|f(x) - f(y)| \leq L \cdot \|x - y\|$

$x^2/3$ *no*



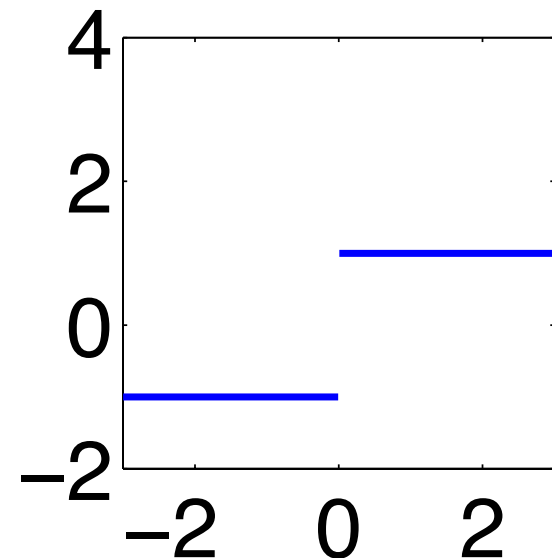
LCG

$|x|$ $L=1$



not LCG

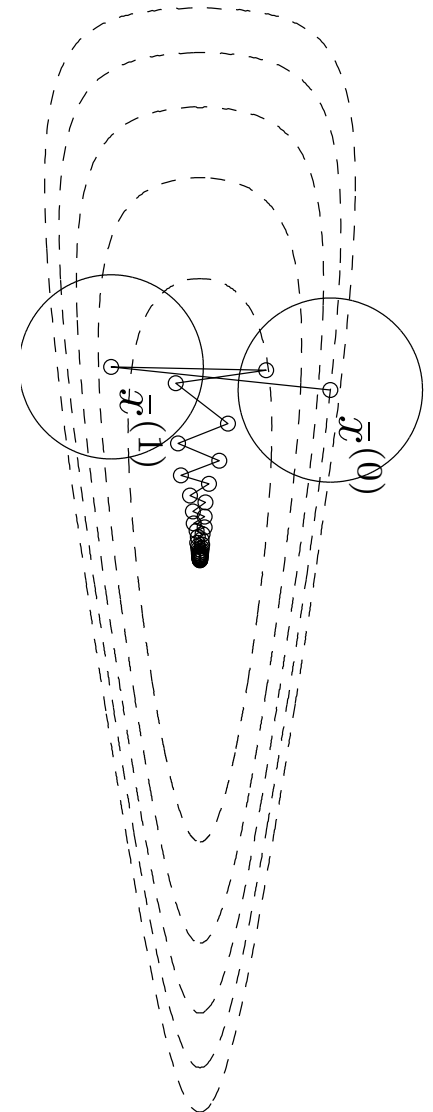
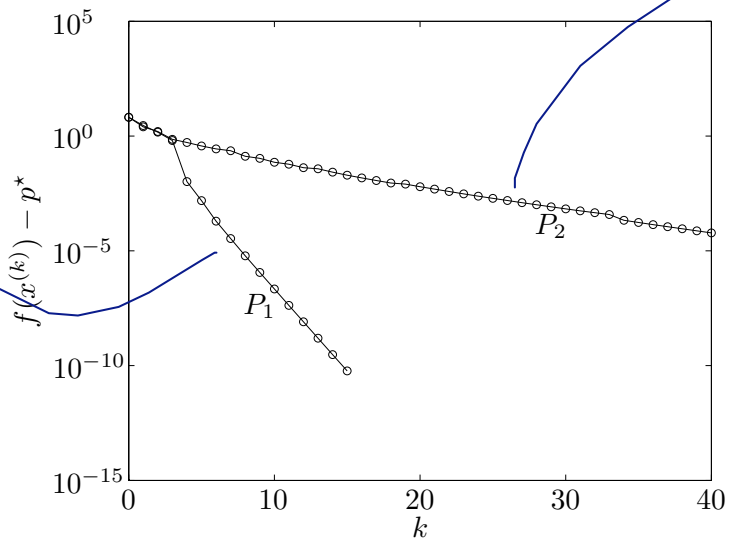
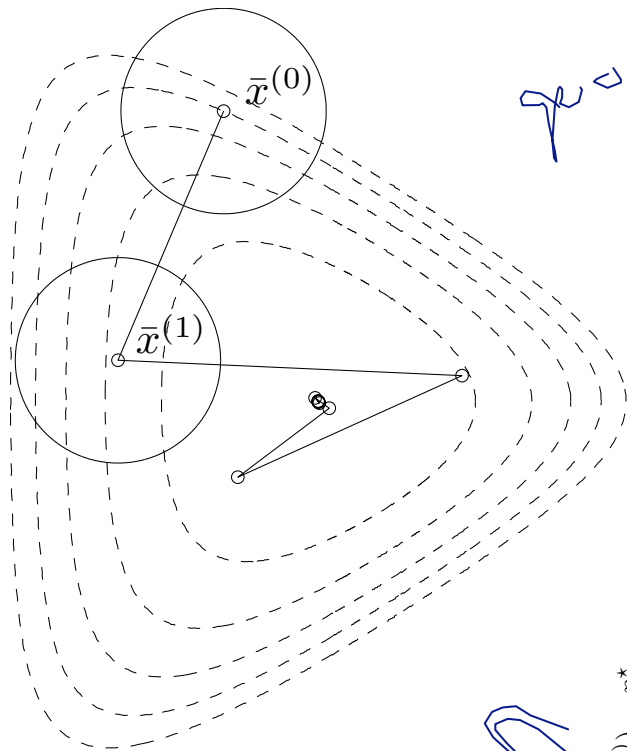
$\text{sgn}(x)$ *no*



Back to gradient descent

- Suppose $f(x)$ is convex, $\nabla f(x)$ exists
- Iterations to get to accuracy $\epsilon(f(x_0) - f(x^*))$: if
 - ▶ f Lipschitz: $O(1/\epsilon^2)$ — $t_k \approx 1/\sqrt{k}$
 - ▶ ∇f Lipschitz: $O(1/\epsilon)$ — t_k const
 - ▶ f strongly convex: $O(\ln(1/\epsilon))$
- Constant in $O(\dots)$: **conditioning** of f

Conditioning



Extensions

- Subgradient descent
- Prox operator (e.g., $g_t = \arg \min_g \|g\|_p^2 + \nabla f \cdot g$)
- FISTA, mirror descent, conjugate gradient
- Nesterov's smoothing
- Line search (BV sec 9.2)
- Stochastic GD (when $f(x) = E(f_i(x) \mid i \sim P)$)
 - ▶ sample one i on each iter, use $\nabla f_i(x)$
 - ▶ or, minibatches: sample a few i 's, use mean $\nabla f_i(x)$

Comparison: stochastic GD

- Iteration bounds for stochastic GD
 - ▶ f Lipschitz: $O(1/\epsilon^2)$ (worse const, same $O()$ as GD)
 - ▶ f strongly convex: $O(\ln(1/\epsilon)/\epsilon)$ (much worse)
- f Lipschitz: stochastic GD often wins big
- Even if f strongly convex:
 - ▶ plain GD: each iter $O(N)$, #iters $O(\ln(1/\epsilon))$
 - ▶ stochastic: each iter $O(1)$, #iters $O(\ln(1/\epsilon)/\epsilon)$
 - ▶ stochastic can win if lots of data, loose tolerance
 - ▶ could make sense to throw data away, use full GD