

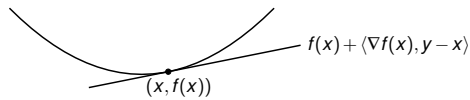
Convex optimization

Slides: Thanks to Di Wang.

$$\min_{x \in Q} f(x)$$

$$f(x) - f(y) \leq \langle \nabla f(x), x - y \rangle$$

Q : feasible space, convex.



First-order Iterative Methods

- ▶ Query $x \in Q$, update using $\nabla f(x)$
- ▶ Low per-iteration cost, $\text{poly}(\frac{1}{\epsilon})$ convergence.
- ▶ Methods of choice in large-scale regime.

Gradient Descent: one dimensional intuition.

Convexity:

$$f(x^*) \geq f(x) + \nabla f(x)^T(x^* - x). \implies f(x) \leq f(x^*) + \nabla f(x)^T(x - x^*)$$

Also: $f(x) - f(x^*) \leq \nabla f(x)^T(x - x^*) = gR$

L-Lipschitz, $R = \|x_0 - x^*\|$:

$$x^+ = x - \frac{1}{L} \nabla f(x) \quad f(x) - f(x^+) \geq \frac{1}{2L} \|\nabla f(x)\|_2^2$$

In one dimension: $\nabla f(x) = g$.

Gap: gR . Progress/step: Roughly $g^2/2$.

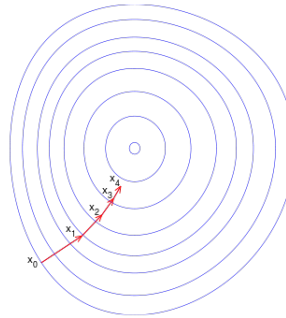
Idea: Gap/(progress/step) \implies roughly $2LR/g$ steps.

Convexity: $g \geq (f(x) - f(x^*)) / R \implies 2LR^2 / (f(x) - f(x^*))$ steps.

While gap $f(x) - f(x^*) \geq \epsilon$ we have $g \geq \epsilon / R$.

$$\implies O(LR^2/\epsilon) \text{ steps reduce gap by } 1/2.$$

Gradient Descent



- ▶ Moves in down-hill direction.
- ▶ Improve objective function value each iteration.
- ▶ Output final point.

Gradient Descent: convergence in ℓ_2

Convexity:

$$f(x^*) \geq f(x) + \nabla f(x)^T(x^* - x). \implies f(x) \leq f(x^*) + \nabla f(x)^T(x - x^*)$$

L-Lipschitz, $R = \|x_0 - x^*\|$:

$$x^+ = x - \frac{1}{L} \nabla f(x) \quad f(x) - f(x^+) \geq \frac{1}{2L} \|\nabla f(x)\|_2^2$$

$$f(x^+) \leq f(x^*) + \nabla f(x)^T(x - x^*) - \frac{1}{2L} \|\nabla f(x)\|_2^2$$

$$\implies f(x^+) - f(x^*) \leq \frac{1}{2} \left(\frac{2}{L} \nabla f(x)^T(x - x^*) - \frac{1}{L^2} \|\nabla f(x)\|_2^2 \right)$$

$$\leq \frac{1}{2} \left(\frac{2}{L} \|\nabla f(x)\|_2 \|x - x^*\| - \frac{1}{L^2} \|\nabla f(x)\|_2^2 + \|x - x^*\|_2^2 \right) \text{ Add 0}$$

$$\leq \frac{1}{2} (\|x - x^*\|_2^2 - \|(x - x^*) - \frac{1}{L} \nabla f(x)\|_2^2)$$

$$\leq \frac{1}{2} (\|x - x^*\|_2^2 - \|x^+ - x^*\|_2^2)$$

$$\sum_k^T f(x_k) - f(x^*) \leq \sum_k^T \frac{1}{2} (\|x_{k-1} - x^*\|_2^2 - \|x_k - x^*\|_2^2)$$

$$\leq \frac{1}{2} (\|x_0 - x^*\|_2^2 - \|x_T - x^*\|_2^2) \leq \frac{1}{2} \|x_0 - x^*\|_2^2$$

$f(x_k)$ is decreasing, we have $f(x_T) \leq \frac{1}{T} \sum_k f(x_k)$.

$$\implies f(x_T) - f(x^*) \leq \frac{LR^2}{2T} \text{ where } R = \|x_0 - x^*\|.$$

Also: $T = O(LR^2/\epsilon)$ iterations for $f(x_T) - f(x^*) \leq \epsilon$.

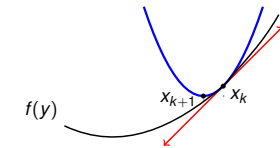
Gradient Descent

L-Lipschitz continuous

$$\|\nabla f(x) - \nabla f(y)\|_* \leq L\|x - y\| \quad \forall x, y \in Q$$

- ▶ Global linear lower bound and quadratic upper bound:

$$\forall y \quad f(x) + \langle \nabla f(x), y - x \rangle \leq f(y) \leq f(x) + \langle \nabla f(x), y - x \rangle + \frac{L}{2} \|y - x\|^2$$



$f(x)$ rises faster when going up, and falls slower when going down.

- ▶ Minimize using quadratic bound

Gradient Descent

Primal progress

$$f(x_k) - f(x_{k+1}) \geq \frac{1}{2L} \|\nabla f(x_k)\|_2^2$$

Convergence

L-Lipschitz, $R = \max_{x: f(x) \leq f(x_0)} \|x - x^*\|$:

$$f(x_T) - f(x^*) \leq O\left(\frac{LR^2}{T}\right)$$

To get ϵ -approximation, need

$$T = O\left(\frac{LR^2}{\epsilon}\right)$$

Relationship?

What is relationship to move closer to feasible?

If wrong side of hyperplane by at least something.
Move to other side.

What is the "hyperplane" here?

$\nabla f(x)$ Maybe.

Mirror Descent: Regret Minimization

- ▶ Regularized average regret

$$\begin{aligned}\tilde{R}_k(u) &= \frac{1}{\alpha k} (-w(u) + \alpha \sum_{i=0}^{k-1} \langle \xi_i, z_i - u \rangle) \\ &= R_k(u) - \frac{w(u)}{\alpha k}\end{aligned}$$

Distance Generating Function

$$w : \mathbb{R}^n \rightarrow \mathbb{R}$$

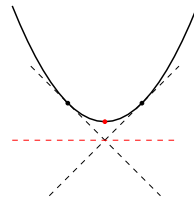
1-strongly convex for norm $\|\cdot\|$:

$$w(y) \geq w(x) + \langle \nabla w(x), y - x \rangle + \frac{1}{2} \|x - y\|^2$$

For ℓ_2 -norm, simply $w(x) = \frac{1}{2} \|x\|_2^2$.

(For distributions: $w(x) = -\sum_i x_i \log x_i$.)

Mirror Descent



- ▶ Each point gives a linear lower bound.
- ▶ Average of the lower bounds becomes flatter.
- ▶ Add the point with current worst regret.
- ▶ Output average of queried points.
- ▶ $x = \alpha x_1 + (1 - \alpha)x_2$
 $\implies f(x) \leq \alpha f(x_1) + (1 - \alpha)f(x_2)$.

Analysis doesn't require L -Lipschitz.

Bregman divergence

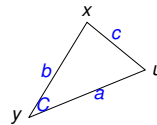
$$V_x(y) = w(y) - \langle \nabla w(x), y - x \rangle - w(x) \geq \frac{1}{2} \|x - y\|^2$$

Standard three point property of Bregman divergence:

$$\forall x, y \geq 0 \quad \langle -\nabla V_x(y), y - u \rangle = V_x(u) - V_y(u) - V_x(y),$$

For ℓ_2 -norm, $V_x(y) = \frac{1}{2} \|x - y\|_2^2$, $\nabla V_x(y) = (x - y)$

Three point property \leftrightarrow Law of cosines



$$c^2 = a^2 + b^2 - 2ab \cos(C) \quad \text{or} \quad 2ab \cos(C) = c^2 - a^2 - b^2$$

$$a^2 = V_y(u), \quad b^2 = V_x(u), \quad c^2 = V_x(u), \quad 2ab \cos(C) = -(x - y) \cdot (y - u)$$

Mirror Descent: Regret Minimization

- ▶ Average **Regret** with loss vector ξ_i 's

$$R_k(u) = \frac{1}{k} \sum_{i=0}^{k-1} \langle \xi_i, z_i - u \rangle$$

Why care about average regret? Bounds gap to OPT:

With $\xi_i = \nabla f(z_i)$, $\bar{z} = \frac{1}{k} \sum_{i=0}^{k-1} z_i$,

$$f(\bar{z}) - f(u) \leq \frac{1}{k} \sum_{i=0}^{k-1} f(z_i) - f(u) \leq \frac{1}{k} \sum_{i=0}^{k-1} \langle \nabla f(z_i), z_i - u \rangle = R_k(u)$$

$$f(\bar{z}) - \text{OPT} \leq \max_u R_k(u)$$

Mirror Descent

$$z_{k+1} = \text{Mirr}(z_k, \alpha \xi_k) = \arg \min_{z \in Q} \{V_{z_k}(z) + \alpha \langle \xi_k, z - z_k \rangle\}$$

Equivalent to regret minimization when $Q = \mathbb{R}^n$:

- ▶ Optimality condition of MD step:

$$\nabla V_{z_k}(z_{k+1}) = -\alpha \xi_k$$

$$z_{k+1} - z_k = -\alpha \xi_k$$

$$z_{k+1} = z_0 - \sum_i \alpha \xi_i$$

- ▶ Regret Minimization:

$$z_{k+1} = \arg \max_z \{-w(z) + \alpha \sum_{i=0}^k \langle \xi_i, z_i - z \rangle\}$$

Optimality condition:

$$z_{k+1} = -\sum_i \alpha \xi_i$$

Recall: $\nabla V_{z_k}(z_{k+1}) = -\alpha \xi_k$ $z_{k+1} - z_k = -\alpha \xi_k$

Mirror Descent

$$\alpha \langle \xi_k, z_k - u \rangle \leq \frac{\alpha^2}{2} \|\xi_k\|_*^2 + V_{z_k}(u) - V_{z_{k+1}}(u)$$

Telescoping T iterations, and **width** $\|\xi_k\|_*^2 \leq \rho^2$

$$1 \cdot \frac{\alpha^2 \rho^2}{2} + V_{z_0}(u) - V_{z_1}(u)$$

$$2 \cdot \frac{\alpha^2 \rho^2}{2} + V_{z_0}(u) - \cancel{V_{z_1}(u)} + \cancel{V_{z_1}(u)} - V_{z_2}(u)$$

$$3 \cdot \frac{\alpha^2 \rho^2}{2} + V_{z_0}(u) - \cancel{V_{z_1}(u)} + \cancel{V_{z_1}(u)} - \cancel{V_{z_2}(u)} + \cancel{V_{z_2}(u)} - V_{z_3}(u) \dots$$

$$\alpha \sum_{i=0}^{T-1} \langle \xi_i, z_i - u \rangle \leq \frac{\alpha^2 \rho^2 T}{2} + V_{z_0}(u) - V_{z_T}(u)$$

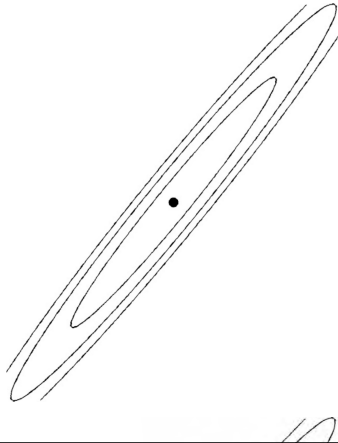
► $\alpha = \frac{\epsilon}{\rho^2}$, **diameter** $V_{z_0}(u) \leq \Theta$, in $T = \frac{2\rho^2\Theta}{\epsilon^2}$ iterations

$$\forall u, f(\bar{z}) - f(u) \leq \frac{\alpha \rho^2}{2} + \frac{V_{z_0}(u)}{\alpha T} \leq \epsilon$$

► Regret terms $\frac{\alpha^2}{2} \|\xi_k\|_*^2$ accumulate, bound step size α .

Linear Coupling

Momentum View:



Linear Coupling

Intuition: If $\|\nabla f(x_k)\|_*^2$ large

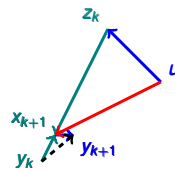
► GD can make large primal progress $\frac{1}{2L} \|\nabla f(x_k)\|_*^2$

► MD suffers large regret $\frac{\alpha^2}{2} \|\nabla f(x_k)\|_*^2$

► Use **primal progress** to cover **regret**.

► Regret terms no longer accumulates, telescope as the primal progress.

Bound $\alpha(f(x_{k+1}) - f(u)) \leq \alpha \langle \nabla f(x_{k+1}), x_{k+1} - u \rangle$



$$\begin{aligned} & \alpha \langle \nabla f(x_{k+1}), x_{k+1} - u \rangle \\ &= \alpha \langle \nabla f(x_{k+1}), z_k - u \rangle \\ & \quad + \alpha \langle \nabla f(x_{k+1}), x_{k+1} - z_k \rangle \\ & \leq \alpha^2 L (f(x_{k+1}) - f(y_{k+1})) \\ & \quad + V_{z_k}(u) - V_{z_{k+1}}(u) \\ & \leq \alpha^2 L (f(x_{k+1}) - f(y_{k+1})) \\ & \quad + V_{z_k}(u) - V_{z_{k+1}}(u) \\ & \quad + \alpha \langle \nabla f(x_{k+1}), \frac{1-\tau}{\tau} (y_k - x_{k+1}) \rangle \\ & \quad + \frac{1-\tau}{\tau} \alpha (f(y_k) - f(x_{k+1})) \\ & \quad + \alpha^2 L (f(y_k) - f(x_{k+1})) \\ & \quad + \alpha^2 L (f(y_k) - f(x_{k+1})) \end{aligned}$$

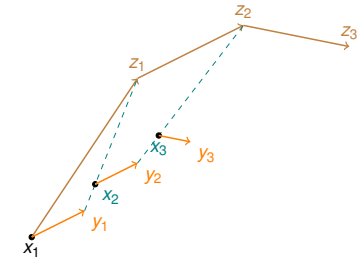
Linear Coupling

► $x_0 = y_0 = z_0$.

► **Coupling**: $x_{k+1} = \tau z_k + (1-\tau)y_k$.

► **MD**: $z_{k+1} = \text{Mirr}(z_k, \alpha \nabla f(x_{k+1}))$

► **GD**: $y_{k+1} = \text{Grad}(x_{k+1})$.



Linear Coupling

► Summing over $0, \dots, T-1$, with $\bar{x} = \frac{1}{T} \sum_i x_i$

$$f(\bar{x}) - f(u) \leq \frac{\alpha L}{T} (f(y_0) - f(y_T)) + \frac{V_{z_0}(u)}{\alpha T}$$

► If $f(y_0) - \text{OPT} \leq d$, diameter $V_{z_0}(u) \leq \Theta$,

$$\alpha = \sqrt{\frac{\Theta}{Ld}}, T = 4\sqrt{\frac{L\Theta}{d}}$$

$$f(\bar{x}) - f(u) \leq \frac{\alpha L d + \Theta / \alpha}{T} \leq \frac{d}{2}$$

► In $T = 4\sqrt{\frac{L\Theta}{d}}$ iterations,

$$f(x_0) - \text{OPT} \leq d \rightarrow f(\bar{x}) - \text{OPT} \leq \frac{d}{2}$$

To get ϵ -approximation:

$$T = O\left(\sqrt{\frac{L\Theta}{\epsilon}} + \sqrt{\frac{L\Theta}{2\epsilon}} + \dots\right) = O\left(\sqrt{\frac{L\Theta}{\epsilon}}\right)$$

Linear Coupling

- ▶ With $\alpha_k = \frac{k+1}{2L}$, can remove phases, and have $f(y_T) - f(u) \leq \epsilon$ after $T = O(\sqrt{\frac{L\Theta}{\epsilon}})$ iterations.
Almost the same as Nesterov's.
- ▶ GD: $O(\frac{LR^2}{\epsilon})$ v.s. MD: $O(\frac{\rho^2\Theta}{\epsilon^2})$ v.s. AGD: $O(\sqrt{\frac{L\Theta}{\epsilon}})$