

Chapter 5 Coordinate Descent

- **Key Idea** Update only one coordinate
- **Results:** Worse case d times more of iteration, under suitable condition results may improve.

Polyak-Łojasiewicz inequality

- *Goal* only to prove function values can converge to optimal, no care of complexity.
- **Definition 5.1** Let $f : \mathbb{R}^d \rightarrow \mathbb{R} \in C^1$ has a global minimum \mathbf{x}^* . We say that f satisfies the *Polyak-Łojasiewicz inequality* (PL inequality) if $\exists \mu > 0$ s.t. $\forall \mathbf{x} \in \mathbb{R}^d, \frac{1}{2} \|\nabla f(\mathbf{x})\|^2 \geq \mu (f(\mathbf{x}) - f(\mathbf{x}^*))$.
- **Lemma 5.2 (Strong Convexity \rightarrow PL ineq)** The proof is in *Lemma E* of lecture 03.
- PL ineq is strictly weaker than strong convexity.
 - E.g. $f(x_1, x_2) = x_1^2$ is not strongly convex, but satisfies PL ineq.
 - Example of non-convex function that satisfies PL ineq: $f(x) = x^2 + \text{Sigmoid}(\frac{x}{20}) \times 5x^2$.
- **Theorem 5.3 (Exponential decay)** Let $f : \mathbb{R}^d \rightarrow \mathbb{R} \in C^1$ has a global minimum \mathbf{x}^* . Suppose that f is L -smooth satisfies μ -PL ineq. Choosing stepsize $\gamma = L^{-1}$, then gradient descent starting with arbitrary \mathbf{x}_0 satisfies
$$f(\mathbf{x}_T) - f(\mathbf{x}^*) \leq \left(1 - \frac{\mu}{L}\right)^T (f(\mathbf{x}_0) - f(\mathbf{x}^*))$$
 - **Proof**
 - By sufficient descent $f(\mathbf{x}_{t+1}) \leq f(\mathbf{x}_t) - \frac{1}{2L} \|\nabla f(\mathbf{x}_t)\|^2$,
 - by PL ineq **RHS** $\leq f(\mathbf{x}_t) - \frac{\mu}{L} (f(\mathbf{x}_t) - f(\mathbf{x}^*))$, then we get what we want.

Coordinate Smoothness

- **Definition 5.4** Let $f : \mathbb{R}^d \rightarrow \mathbb{R} \in C^1$ and $\mathcal{L} = (L_1, L_2, \dots, L_d) \in \mathbb{R}^d$. Function f is called *coordinate-wise smooth (with parameter \mathcal{L})* if $\forall i \in [d], \forall \mathbf{x} \in \mathbb{R}^d, \lambda \in \mathbb{R}, f(\mathbf{x} + \lambda \mathbf{e}_i) \leq f(\mathbf{x}) + \lambda \nabla_i f(\mathbf{x}) + \frac{L_i}{2} \lambda^2$
 - If $\forall i, L_i = L$, f is said to be coordinate-wise smooth with parameter L .
- Coordinate-wise smoothness is more fine grained.
 - $f(x_1, x_2) = x_1^2 + 10x_2^2$ is $(2, 20)$ -coordinate-smooth, but only **20**-smooth.
 - $f(x_1, x_2) = x_1^2 + x_2^2 + Mx_1x_2$ is $(2, 2)$ -coordinate-smooth, but only $(M + 2)$ -smooth.
- **Algorithm**
 - (i) choose an active coordinate $i \in [d]$
 - (ii) $\mathbf{x}_{t+1} := \mathbf{x}_t - \gamma_i \nabla_i f(\mathbf{x}_t) \mathbf{e}_i$
- **Lemma 5.5 (Sufficient Descent)** A stepsize of $\gamma_i = L_i^{-1}$ gives $f(\mathbf{x}_{t+1}) \leq f(\mathbf{x}_t) - \frac{1}{2L_i} |\nabla_i f(\mathbf{x}_t)|^2$.

Randomized CD

- **Algorithm Change** (i) sample $i \in [d]$ uniformly at random.
- **Theorem 5.6** Let $f : \mathbb{R}^d \rightarrow \mathbb{R} \in C^1$ has a global minimum \mathbf{x}^* . Suppose that f is coordinate-wise L -smooth and satisfies the μ -PL inequality. Choosing stepsize $\gamma_i = L^{-1}$, then randomized coordinate descent with arbitrary start satisfies
$$\mathbb{E}[f(\mathbf{x}_T) - f(\mathbf{x}^*)] \leq \left(1 - \frac{\mu}{dL}\right)^T (f(\mathbf{x}_0) - f(\mathbf{x}^*)).$$
 - **Proof**
 - Take expectation over sufficient descent with each coordinate of probability $1/d$,
$$\mathbb{E}[f(\mathbf{x}_{t+1})|\mathbf{x}_t] \leq f(\mathbf{x}_t) - \frac{1}{2L} \sum_{i=1}^d \frac{1}{d} |\nabla_i f(\mathbf{x}_t)|^2 = f(\mathbf{x}_t) - \frac{1}{2dL} \|\nabla f(\mathbf{x}_t)\|^2$$
 - With PL ineq, **RHS** $\leq f(\mathbf{x}_t) - \frac{\mu}{dL} (f(\mathbf{x}_t) - f(\mathbf{x}^*))$
 - Take expectation over condition $\mathbb{E}[\cdot|\mathbf{x}_t]$, we arrive at our conclusion.
 - **Comment** $\left(1 - \frac{\mu}{L}\right) \approx \left(1 - \frac{\mu}{dL}\right)^d$, this is nearly the same as vanilla GD, need improvement.

Importance Sampling

- **Algorithm Change** (i) when L_i not the same, sample $i \in [d]$ with probability $\frac{L_i}{\sum_{j=1}^d L_j}$
- **Theorem 5.7** Let $f : \mathbb{R}^d \rightarrow \mathbb{R} \in C^1$ has a global minimum \mathbf{x}^* . Suppose that f is coordinate-wise \mathcal{L} -smooth and satisfies the μ -PL inequality. Let $\bar{L} = \frac{1}{d} \sum_{i=1}^d L_i$ be the average of all coordinate-wise smoothness constants. Then coordinate descent with importance sampling with arbitrary start satisfies
$$\mathbb{E}[f(\mathbf{x}_T) - f(\mathbf{x}^*)] \leq \left(1 - \frac{\mu}{d\bar{L}}\right)^T (f(\mathbf{x}_0) - f(\mathbf{x}^*))$$

- **Proof**
 - $\mathbb{E}[f(\mathbf{x}_{t+1})|\mathbf{x}_t] \leq f(\mathbf{x}_t) - \sum_{i=1}^d \frac{1}{2L} \frac{L_i}{\sum_{j=1}^d L_j} |\nabla_i f(\mathbf{x}_t)|^2 = f(\mathbf{x}_t) - \frac{1}{2d\bar{L}} \|\nabla f(\mathbf{x}_t)\|^2$
- **Comment** \bar{L} can be much smaller than $L = \max_{i=1}^d L_i$, this is a improvement in constant. When all L are equal, no improvement.

Steepest Coordinate Descent

- **Algorithm Change** (i) Choose $i = \underset{i \in [d]}{\operatorname{argmax}} |\nabla_i f(\mathbf{x}_t)|$, also called *Gauss-Southwell* rule.
- Since $\max_i |\nabla_i f(\mathbf{x})|^2 \geq \frac{1}{d} \sum_{i=1}^d |\nabla_i f(\mathbf{x})|^2$, for a coordinate- L -smooth function, we have
 - $f(\mathbf{x}_{t+1}) \leq f(\mathbf{x}_t) - \frac{1}{2L} \max_i |\nabla_i f(\mathbf{x})|^2 \leq f(\mathbf{x}_t) - \frac{1}{2L} \sum_{i=1}^d \frac{1}{d} |\nabla_i f(\mathbf{x}_t)|^2$
 - Then we arrive at **Corollary 5.8**, still $f(\mathbf{x}_T) - f(\mathbf{x}^*) \leq (1 - \frac{\mu}{dL})^T (f(\mathbf{x}_0) - f(\mathbf{x}^*))$

Better Convergence of Steepest Descent with ℓ_1 norm strong convexity (Also known as Steeper)

- **Lemma 5.9** (ℓ_1 norm μ -strong convexity $\rightarrow \ell_\infty$ norm μ -PL inequality) This is a natural result for *Lemma E*, as the dual norm of $\|\cdot\|_1$ is $\|\cdot\|_\infty$.
- Proof of dual norm $\|\cdot\|_{p^*} = \|\cdot\|_q$ where $p^{-1} + q^{-1} = 1$ is by Hölder inequality $\sum_{k=1}^n |x_k y_k| \leq (\sum_{k=1}^n |x_k|^p)^{\frac{1}{p}} (\sum_{k=1}^n |y_k|^q)^{\frac{1}{q}}$ (equality holds if \mathbf{x}^p and \mathbf{y}^q are proportional.)
 - Proof of special case of $(1, \infty)$ is simple.
- **Theorem 5.10** Let $f: \mathbb{R}^d \rightarrow \mathbb{R} \in C^1$ has a global minimum \mathbf{x}^* . Suppose that f is coordinate-wise L -smooth and satisfies the μ_1 -PL inequality under ℓ_1 norm. Then *steepest coordinate descent* with step size $\gamma_i = L^{-1}$ arbitrary start satisfies

$$f(\mathbf{x}_T) - f(\mathbf{x}^*) \leq (1 - \frac{\mu_1}{L})^T (f(\mathbf{x}_0) - f(\mathbf{x}^*)).$$
 - **Key Idea** $f(\mathbf{x}_{t+1}) \leq f(\mathbf{x}_t) - \frac{1}{2L} \max_i |\nabla_i f(\mathbf{x})|^2 = f(\mathbf{x}_t) - \frac{1}{2L} \|\nabla f(\mathbf{x}_t)\|_\infty^2$
- **Comment** Since $\|\mathbf{y} - \mathbf{x}\|_1 \geq \|\mathbf{y} - \mathbf{x}\| \geq \|\mathbf{y} - \mathbf{x}\|_1 / \sqrt{d}$
 - If f is ℓ_1 μ -strong convex $\rightarrow f$ is ℓ_2 μ -strong convex.
 - If f is ℓ_2 μ -strong convex $\rightarrow f$ is ℓ_1 μ/d -strong convex.
 - Seems no better than ℓ_2 case, but cases like $f(x) = \mathbf{x}^\top \operatorname{diag}\{\lambda_i\} \mathbf{x} / 2$ shows ℓ_1 gives much better results. ([Appendix C of this](#))
 - Has sth to do with *convex conjugates*.

Greedy coordinate descent

- **Algorithm Change** (ii) $\mathbf{x}_{t+1} := \underset{\lambda \in \mathbb{R}}{\operatorname{argmin}} f(\mathbf{x}_t + \lambda \mathbf{e}_i)$
 - Might be easier when f is analytic.
- **Problem** May stuck at non-minimum.
 - Example: $f(\mathbf{x}) := \|\mathbf{x}\|^2 + |x_1 - x_2|$, (x, x) is minimal for all single coordinate in range $|x| \leq 1/2$.
 - The following form is guaranteed to not be.
- **Theorem 5.11** Let $f: \mathbb{R}^d \rightarrow \mathbb{R}$ be of the form $f(\mathbf{x}) := g(\mathbf{x}) + h(\mathbf{x})$ where $h(\mathbf{x}) = \sum_i h_i(x_i)$, g, f_i all convex, $g \in C^1$, then whenever *greedy coordinate descent* makes no improvement, its in minimum.
 - **Proof**
 - No improvement means $\forall \lambda, i \in [d], f(\mathbf{x} + \lambda \mathbf{e}_i) \geq f(\mathbf{x})$,
 - then $f(\mathbf{x} + (y_i - x_i) \mathbf{e}_i) \geq f(\mathbf{x})$
 - let $p(z) := g(z, \mathbf{x}_{-i}) - \partial_i g(x_i)(z - x_i)$ and $q(z) := h_i(z) + \partial_i g(x_i)(z - x_i)$, since g convex and differentiable, $p(z)$ reach global minimum at $z = x_i$, $f_i := f(z; \mathbf{x}_{-i}) = p(z) + q(z)$, so f_i also reach global minimum at $z = x_i$.
 - Since $p(z)$ differentiable, $\partial p(x_i) = 0$, if q does not reach global minimum at $z = x_i$, then $\exists y$ s.t. $q(y) < q(x_i)$, then by convexity $q(x_i + \lambda(y - x_i)) \leq q(x_i) + \lambda(q(y) - q(x_i))$ is bounded by a negative slope of $-(q(x_i) - q(y))$ at $z = x_i$.
 - This leads to contradiction where $f = p + q$ reaches minimum at $z = x$.
 - Therefore, $q(z)$ also reaches minimal at $z = x_i$.
 - The above discussion means $\forall y_i, q(y_i) \geq q(x_i)$, equivalently $h_i(y_i) + \partial_i g(\mathbf{x})(y_i - x_i) - h_i(x_i) \geq 0$
 - Since g convex,

$$f(\mathbf{y}) \geq g(\mathbf{x}) + \nabla g(\mathbf{x})^\top (\mathbf{y} - \mathbf{x}) + h(\mathbf{y}) - h(\mathbf{x}) = f(\mathbf{x}) + \underbrace{\sum_{i=1}^d (\nabla_i g(\mathbf{x})(y_i - x_i) + h_i(y_i) - h_i(x_i))}_{\geq 0} \geq f(\mathbf{x})$$
 - **Comment** LASSO $\min_{\mathbf{x} \in \mathbb{R}^n} \|\mathbf{A}\mathbf{x} - \mathbf{b}\|^2 + \lambda \|\mathbf{x}\|_1$ is of this form
 - Under mild regularity on g , convergence is affirmative.