

Revisiting gradient descent : We will show convergence

only Lipschitz continuity

Lipschitz continuity

+ strong convexity

We have $\forall x, c, y \in \text{dom } f$,

- $f(y) = f(x) + \nabla^T f(x)(y - x) + \frac{1}{2}(y - x)^T \nabla^2 f(c)(y - x)$
- $\nabla^2 f(c) \preceq L I$
 - $\implies (y - x)^T \nabla^2 f(c)(y - x) \leq L \|y - x\|^2$
- Thus we get

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

$$x^{(k+1)} = x^{(k)} - t^{(k)} \nabla f(x^{(k)})$$

Generally, convergence analysis will require Lipschitz continuity

- Considering $x^k \equiv x$, and $x^{k+1} \equiv y$, and a fixed step size t , we get

$$f(x^{k+1}) \leq f(x^k) - t \nabla^\top f(x^k) \nabla f(x^k) + \frac{Lt^2}{2} \|\nabla f(x^k)\|^2$$

$$\implies f(x^{k+1}) \leq f(x^k) - \left(1 - \frac{Lt}{2}\right)t \|\nabla f(x^k)\|^2$$

- Taking $0 < t \leq \frac{1}{L} \implies 1 - \frac{Lt}{2} \geq \frac{1}{2}$, we have

$$f(x^{k+1}) \leq f(x^k) - \frac{t}{2} \|\nabla f(x^k)\|^2$$

Note: $t < \frac{2}{L}$

is necessary
for convergence

For a weaker
notion of convergence discussed here $t \leq \frac{2}{L}$ should be ok.
So you can replace t with $t/2$ on this & next pages if result
should still hold.

- Using convexity, we have $f(x^*) \geq f(x^k) + \nabla^\top f(x^k)(x^* - x^k)$
 $\implies f(x^k) \leq f(x^*) + \nabla^\top f(x^k)(x^k - x^*)$
- Thus,

$$f(x^{k+1}) \leq f(x^k) - \frac{t}{2} \|\nabla f(x^k)\|^2$$

$$\implies f(x^{k+1}) \leq f(x^*) + \nabla^\top f(x^k)(x^k - x^*) - \frac{t}{2} \|\nabla f(x^k)\|^2$$

$$\implies f(x^{k+1}) \leq f(x^*) - \frac{1}{2t} \|x^k - x^*\|^2 + \nabla^\top f(x^k)(x^k - x^*) - \frac{t}{2} \|\nabla f(x^k)\|^2 + \frac{1}{2t} \|x^k - x^*\|^2$$

$$\implies f(x^{k+1}) \leq f(x^*) + \frac{1}{2t} (\|x^k - x^*\|^2 - \|x^k - x^* - t\nabla f(x^k)\|^2)$$

$$\implies f(x^{k+1}) \leq f(x^*) + \frac{1}{2t} (\|x^k - x^*\|^2 - \|x^{k+1} - x^*\|^2)$$

$$\implies f(x^{k+1}) - f(x^*) \leq \frac{1}{2t} (\|x^k - x^*\|^2 - \|x^{k+1} - x^*\|^2)$$

- Over all iterations, we have

$$\sum_{i=0}^k (f(x^i) - f(x^*)) \leq \frac{1}{2t} \left(\|x^{(0)} - x^*\|^2 - \frac{1}{2\lambda} \|x^{(k+1)} - x^*\|^2 \right)$$

≤ 0

- Since $f(x^{k+1}) \leq f(x^k) \forall k = 0, 1, \dots$, we get

$$f(x^k) - f(x^*) \leq \frac{1}{k} \sum_{i=0}^{k-1} (f(x^i) - f(x^*)) \leq \frac{\|x^{(0)} - x^*\|^2}{2tk}$$

Q: What abt rate of convergence?

Question: Could we analyze Gradient descent more generally?

- Assume backtracking line search OR exact line search?
- Continue assuming Lipschitz continuity
 - ▶ Curvature is upper bounded: $\nabla^2 f(x) \preceq L I$
- Assume **strong convexity**
 - ▶ Curvature is lower bounded: $\nabla^2 f(x) \succeq m I$
 - ▶ For instance, we wouldn't want to use gradient descent for a linear function (no curvature)

④ How can we get Q or R linear convergence?

$$\Theta(\epsilon^{(0,1)}) \rightarrow f(x^k) - f(x^*) \leq \Theta^k (f(x^0) - f(x^*)) \xrightarrow{v^k} \text{R-linear}$$
$$\Theta(\epsilon^{(0,1)}) \rightarrow f(x^k) - f(x^*) \leq \Theta(f(x^{k-1}) - f(x^*)) \xrightarrow{\text{Or-linear}}$$

- Lipschitz continuity

$$\nabla^2 f(x) \preceq L I$$

$$\|\nabla f(x) - \nabla f(y)\| \leq L \|x - y\|$$

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

- Convexity

- ▶ Curvature should **not** be negative

$$\nabla^2 f(x) \succeq 0$$

$$f(y) \geq f(x) + \nabla^T f(x)(y - x)$$

- Strong convexity

$$\nabla^2 f(x) \succeq m I$$

$$f(y) \geq f(x) + \nabla^T f(x)(y - x) + \frac{m}{2} \|y - x\|^2$$

- ▶ For example, augmented Lagrangian is used to introduce strong convexity

Using strong convexity

- $f(y) \geq f(x) + \nabla^T f(x)(y - x) + \frac{m}{2} \|y - x\|^2$
 \geq minimum value the RHS can take as a function of y
- Minimum value of RHS
$$\nabla f(x) + my - mx = 0$$

$$\implies y = x - \frac{1}{m} \nabla f(x)$$
- Thus,
$$f(y) \geq f(x) + \nabla^T f(x) \left(-\frac{1}{m} \nabla f(x) \right) + \frac{m}{2} \left\| -\frac{1}{m} \nabla f(x) \right\|^2$$

$$\implies f(y) \geq f(x) - \frac{1}{2m} \|\nabla f(x)\|^2$$
 - ▶ Here, LHS is independent of x , and RHS is independent of y

$$f(x^*) \geq f(x) - \frac{1}{2m} \|\nabla f(x)\|^2$$

- If $\|\nabla f(x)\|$ is small, the point is nearly optimal

- ▶ If $\|\nabla f(x)\| \leq \sqrt{2m\epsilon}$, then:

$$f(x) - f(x^*) \leq \epsilon$$

- ▶ As the gradient $\|\nabla f(x)\|$ approaches 0, we get closer to the optimal solution x^*

} A knob provided by strong convexity!

Analysis for Backtracking Line Search

- Backtracking line search exits when

$$f\left(x^k - t \nabla f(x^k)\right) \leq f(x^k) - \frac{t}{2} \|\nabla f(x^k)\|^2$$

- ▶ where $t = (\beta)^r t_{orig}$
 - ★ t_{orig} was the initial step size before the invocation of backtracking line search
 - ★ r is the number of iterations before the loop terminated
- The margin of backtracking line search, $\frac{t}{2} \|\nabla f(x^k)\|^2$, is inspired by strong convexity

- Since f is strongly convex, and also Lipschitz continuous, we have for some L :

$$f(x^{k+1}) \leq f(x^k) + \left(\frac{Lt^2}{2} - t\right) \|\nabla f(x^k)\|^2$$

- We also consider

$$\begin{aligned} 0 < t \leq \frac{1}{L} &\implies t^2 \leq \frac{t}{L} \implies \frac{Lt^2}{2} \leq \frac{t}{2} \\ &\implies \frac{Lt^2}{2} - t \leq -\frac{t}{2} \end{aligned}$$

- Thus, we get the exit condition of backtracking line search

$$f(x^{k+1}) \leq f(x^k) - \frac{t}{2} \|\nabla f(x^k)\|^2$$

$$\implies f\left(x^k - t\nabla f(x^k)\right) \leq f(x^k) - \frac{t}{2} \|\nabla f(x^k)\|^2$$

- Convergence of gradient descent, given this condition, has been proved below