

GD continued

Saturday, January 11, 2025 7:52 PM

prop $\nabla f(x) = 0 \Rightarrow x$ is global min for convex f .

Lera. $f(x) - f(y) \leq \|\nabla f(x)\|_2 \|x - y\|_2$ for convex f .

Pf. $f(y) - f(x) \geq \langle \nabla f(x), y - x \rangle$

$$\begin{aligned} f(x) &\leq f(y) - \langle \nabla f(x), y - x \rangle \\ &\leq f(y) + \|\nabla f(x)\|_2 \|y - x\|_2. \end{aligned}$$

Thm. $x^{(k+1)} = x^{(k)} - \frac{1}{L} \nabla f(x^{(k)})$ achieves

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

$$\begin{aligned} \text{Pf. } f(x^{(k+1)}) &= f(x^{(k)} - \frac{1}{L} \nabla f(x^{(k)})) \\ &\leq f(x^{(k)}) - \frac{1}{2L} \|\nabla f(x^{(k)})\|_2^2 \\ f(x^{(k+1)}) - f(x^*) &\leq \underbrace{f(x^{(k)}) - f(x^*)}_{\mathcal{E}_k} - \frac{1}{2L} \underbrace{\left(\frac{f(x^{(k)}) - f(x^*)}{\|x^{(k)} - x^*\|_2} \right)^2}_{\mathcal{E}_{k+1}} \end{aligned}$$

$$\mathcal{E}_{k+1} \leq \mathcal{E}_k - \frac{1}{2L} \left(\frac{\mathcal{E}_k}{R} \right)^2$$

$$\frac{1}{\mathcal{E}_{k+1}} - \frac{1}{\mathcal{E}_k} = \frac{\mathcal{E}_k - \mathcal{E}_{k+1}}{\mathcal{E}_{k+1} \mathcal{E}_k} \leq \frac{1}{2L} \frac{\mathcal{E}_k^2}{R^2} = \frac{1}{2LR^2}$$

$$\frac{1}{\varepsilon_{k+1}} - \frac{1}{\varepsilon_k} = \frac{\varepsilon_k - \varepsilon_{k+1}}{\varepsilon_{k+1} \varepsilon_k} \leq \frac{1}{2L} \frac{R^2}{\varepsilon_k^2} = \frac{R^2}{2L}$$

$$\begin{aligned}\mathcal{E}_0 &= f(x^{(0)}) - f(x^*) \leq \langle \nabla f(x^*), x^{(0)} - x^* \rangle + \frac{L}{2} \|x^{(0)} - x^*\|_2^2 \\ &= \frac{L}{2} \|x^{(0)} - x^*\|_2^2 \leq \frac{LR^2}{2}\end{aligned}$$

$$\begin{aligned}\frac{1}{\varepsilon_k} &\leq \frac{2}{LR^2} + \frac{k}{2LR^2} = \frac{k+4}{2LR^2} \\ \Rightarrow \varepsilon_k &\leq \frac{2LR^2}{k+4}\end{aligned}$$

$$R = \max_{\substack{x: f(x) \leq f(x^{(0)})}} \|x - x^*\|_2$$

In fact, we can set $R = \|x^{(0)} - x^*\|_2$

Lemma. For $h \leq \frac{2}{L}$, $\|x^{(k+1)} - x^*\|_2 \leq \|x^k - x^*\|_2$.

Pf.

$$\begin{aligned}\|x^{(k+1)} - x^*\|_2^2 &= \|x^{(k)} - h\nabla f(x^k) - x^*\|_2^2 \\ &= \|x^{(k)} - x^*\|_2^2 + h^2 \|\nabla f(x^k)\|_2^2 - 2h \langle \nabla f(x^k), x^{(k)} - x^* \rangle\end{aligned}$$

$$\nabla f(x^k) = \nabla f(x^*) + \int_0^1 \langle \nabla^2 f(x^* + t(x^k - x^*)), x^k - x^* \rangle dt$$

$$\langle \nabla f(x^k), x^k - x^* \rangle = (x^k - x^*)^\top \int_0^1 \nabla^2 f(x^* + t(x^k - x^*)) dt \quad (x^k - x^*)$$

$$H \succ 0$$

$$H^2 \preceq L H \geq \frac{1}{L} \|H(x^k - x^*)\|_2^2$$

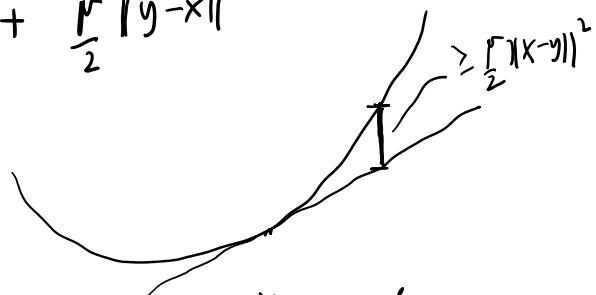
$$= \frac{1}{L} \|\nabla f(x^k)\|_2^2$$

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

$$h \geq \frac{2}{L} \Rightarrow \geq 0.$$

Strongly convex.

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



Lemma The following are equivalent.

$$\mu > 0$$

$$(1) \quad \forall t \quad f((1-t)x + ty) \leq (1-t)f(x) + t f(y) - \frac{\mu t(1-t)}{2} \|x - y\|^2$$

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

$$(3) \quad \nabla^2 f(x) \succ \mu I$$

Def

$$(3) \quad \nabla^2 f(x) \succcurlyeq \text{MI}$$

Pf. Define $g(x) = f(x) - \frac{\mu}{2} \|x\|^2$. Apply equivalent definitions of convexity to g .

Lemma For f μ -strongly convex,

$$\|\nabla f(x)\|^2 \geq 2\mu (f(x) - f(y)).$$

$$\begin{aligned} \text{Pf. } f(y) &\geq f(x) + \langle \nabla f(x), y-x \rangle + \frac{\mu}{2} \|y-x\|^2 \\ f(x) - f(y) &\leq \langle \nabla f(x), x-y \rangle - \frac{\mu}{2} \|y-x\|^2 \\ &\leq \max_z \langle \nabla f(x), z \rangle - \frac{\mu}{2} \|z\|^2 \end{aligned}$$

$$\text{At optimal } z, \quad \nabla f(x) = \text{M}z$$

$$\text{So } f(x) - f(y) \leq \frac{1}{2\mu} \|\nabla f(x)\|_2^2.$$

$$\text{Thm. } f(x^{(k)}) - f(x^*) \leq \left(1 - \frac{\mu}{L}\right)^k (f(x^0) - f(x^*))$$

$$\begin{aligned} \text{Pf. } f(x^{(k+1)}) - f(x^*) &\leq f(x^{(k)}) - f(x^*) - \frac{1}{2L} \|\nabla f(x^{(k)})\|_2^2 \\ &\leq f(x^{(k)}) - f(x^*) - \frac{2\mu}{2L} (f(x^{(k)}) - f(x^*)) \\ &\quad \dots \quad \text{repeat } k \rightarrow k+1 \quad (1, +1) \end{aligned}$$

$$\leq \left(1 - \frac{\mu}{L}\right) \left(f(x^*) - f(x^+)\right).$$

f		# steps	
general	$\ \nabla f(x)\ \leq \varepsilon$	$\frac{2L(f(x^0) - f(x^*))}{\varepsilon^2}$	} different measures of convergence.
convex	$f(x) \leq f^* + \varepsilon$	$\frac{2L \ x^0 - x^*\ ^2}{\varepsilon}$	
β , strongly convex	$f(x) \leq f^* + \varepsilon$	$\frac{L}{\beta} \log \left(\frac{f(x^0) - f(x^*)}{\varepsilon} \right)$	

Note that for $f(x) = \frac{1}{2}\|x\|^2$
 $\|\nabla f(x)\| \leq \varepsilon$ is $\|x\| \leq \varepsilon$
 vs. $f(x) \leq f^* + \varepsilon$ is $\|x\|^2 \leq \varepsilon$

So the rates are the same.

Sampling

recall goal: given access to f , sample according to e^{-f} .

Ex. 1. $f(x) = \frac{\|x\|^2}{2}$ $e^{-\frac{\|x\|^2}{2}}$ $N(0, I)$

Ex. 2

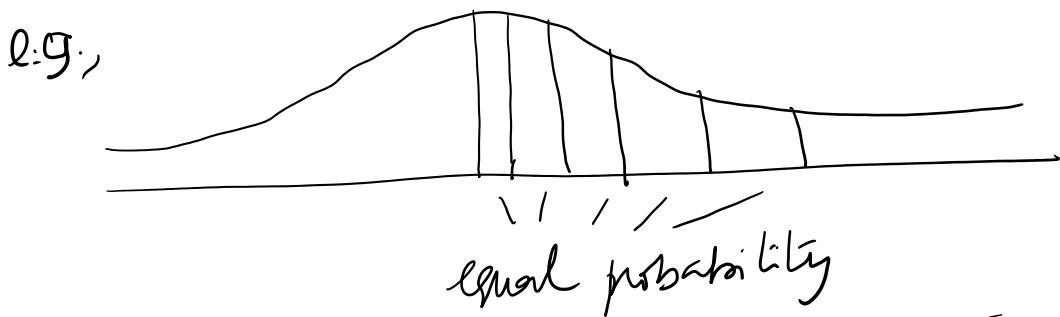
$$f(x) = \begin{cases} 0 & x \in K \\ \infty & x \notin K \end{cases}$$

uniformly sample K .

To sample $N(0, I)$

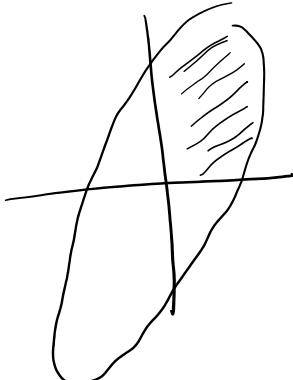
observe that this is independent $x_i \sim N(0, I)$.

Many numerical methods to sample $N(0, I)$



How about

$$e^{-\frac{x^T A x}{2}} \cdot \mathbb{1}_{x \geq 0} ?$$



without the constraint

$$e^{-\frac{x^T A x}{2}}$$
 is $N(0, \bar{A}^{-1})$

Gaussian with covariance \bar{A}^{-1}

$$y = \bar{A}^{\frac{1}{2}} x \quad y^T y = x^T A x$$

$$\mathbb{E}(y y^T) = \mathbb{E}(\bar{A}^{\frac{1}{2}} X X^T \bar{A}^{\frac{1}{2}}) = I$$

$$x = \bar{A}^{\frac{1}{2}} y \quad y \sim N(0, I).$$

How about with constraint $x \geq 0$? .. n+1

How about with constraint $x \geq 0$?
 Rejection sampling (sample $x \sim N(0, A^{-1})$
 is too expensive. output x if $x \geq 0$)

Need an algorithm.

Back to 6D. Can view as a continuous time
 algorithm.

Gradient flow

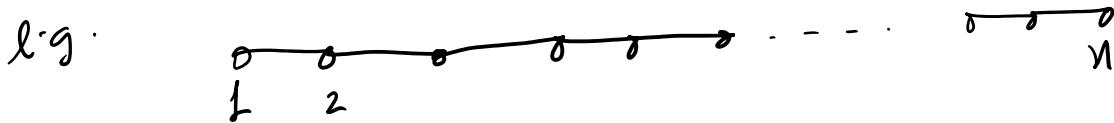
$$dx_t = -\nabla f(x_t) dt$$

$$\begin{aligned} d(f(x_t) - f(x_0)) &= \langle \nabla f(x_t), dx_t \rangle \\ &= -\|\nabla f(x_t)\|_2^2 dt \end{aligned}$$

Under γ -strong convexity,
 $d(f(x_t) - f(x^*)) \leq -2\gamma(f(x_t) - f(x^*))dt$

$$f(x_t) - f(x^*) \leq e^{-2\gamma t}(f(x_0) - f(x^*))$$

How about Sampling?
 need to introduce randomness.



To sample uniformly, start at i , $i-1$, ..., 1 in $\mathbb{C}[n]$

To sample uniformly, start at i

w.p. $\frac{1}{2} \rightarrow i-1$ | stay if $i-1 \notin \mathbb{N}$
 |
 $\frac{1}{2} \rightarrow i+1$

$p^{(k+1)}(i) = p^{(k)}(i-1) \cdot \frac{1}{2}$ $p(i) = \frac{1}{n}$ is stationary
 $\rightarrow p^{(k)}(i+1) \cdot \frac{1}{2}$

$dX_t = -\nabla f(X_t) dt$

(LD)
$$dX_t = -\nabla f(X_t) dt + \sqrt{2} dW_t$$

"Brownian motion". ↑
 infinitesimal Gaussian

$W_t \sim N(0, t)$

Langevin Diffusion.

Algorithm? Fix $h > 0$

$$X^{(k+1)} = X^{(k)} - h \nabla f(X^{(k)}) + \sqrt{2h} Z$$
 $Z \sim N(0, I)$
 $\sqrt{2h} Z \quad \bar{Z} \sim N(0, I)$

Thm1. For continuously differentiable f , with $\int e^{-f} < \infty$,
 density $\propto e^{-f}$ is stationary for LD.

More generally $dX_t = \mu(X_t) dt + \sigma(X_t) dW_t$ is a Stochastic Differential Equation (SDE).
 $X_t \in \mathbb{R}^n$
 $\sigma \in \mathbb{R}^{n \times m}$
 $W_t \in \mathbb{R}^m$

Thm [Fokker-Planck]

$$\frac{d p_t}{dt} = - \sum_{i=1}^n \frac{\partial}{\partial x_i} (r(x_t)_i p_t) + \frac{1}{2} \sum_{i,j=1}^n \frac{\partial^2}{\partial x_i \partial x_j} (D(x_t)_{ij} p_t)$$

$$\begin{aligned} D(x_t) &= \sigma(x_t) \sigma(x_t)^T \\ &= - \nabla \cdot (p_t r(x_t)) + \frac{1}{2} \nabla \cdot (\nabla \cdot (p_t D(x_t))) \\ &= - \nabla \cdot (p_t r(x_t)) + \frac{1}{2} \Delta (p_t D(x_t)). \end{aligned}$$

$$\text{Pf. (of Thm 1). } r(x_t)_i = - \nabla f(x_t)_i = - \frac{\partial f(x)}{\partial x_i}$$

$$\sigma(x_t) = \sqrt{2} I$$

$$D(x_t)_{ij} = \begin{cases} 2 & \text{if } i=j \\ 0 & \text{if } i \neq j \end{cases}$$

So at stationary p

$$0 = \frac{d p}{dt} = - \sum_{i=1}^n \frac{\partial}{\partial x_i} \left(- \frac{\partial f(x)}{\partial x_i} p(x) \right) + \frac{1}{2} \sum_i \frac{\partial^2}{\partial x_i^2} (2 p(x))$$

$$0 = \sum_i \frac{\partial}{\partial x_i} \left(p(x) \frac{\partial f(x)}{\partial x_i} + \frac{\partial p(x)}{\partial x_i} \right)$$

$$= \sum_i \frac{\partial}{\partial x_i} \left(p(x) \left(\frac{\partial f(x)}{\partial x_i} + \frac{\partial \log p(x)}{\partial x_i} \right) \right)$$

$$= \sum_i \frac{\partial}{\partial x_i} \left(p(x) \frac{\partial \log \left(\frac{p(x)}{e^{-f(x)}} \right)}{\partial x_i} \right)$$

and hence we see that $p(x) = C e^{-f(x)}$

And now we see that $p(x) = C e^{-\lambda x}$
is a solution.

Q. Rate of convergence?

How to measure it?

$$d_{TV}(p, q) = \frac{1}{2} \int_{\Omega} |p(x) - q(x)| dx = \sup_{A \subseteq \Omega} |p(A) - q(A)|$$

$$\chi^2(p, q) = \int \left(\frac{p(x)}{q(x)} - 1 \right)^2 q(x) dx = \mathbb{E}_q \left(\left(\frac{p}{q} - 1 \right)^2 \right)$$

$$d_{KL}(p, q) = \int p(x) \log \frac{p(x)}{q(x)} dx = \mathbb{E}_p \left(\log \frac{p}{q} \right)$$

$$W_2(p, q) = \min_{\pi: \text{coupling of } X \sim p, Y \sim q} \mathbb{E} (||x - y||^2)$$

$||\cdot||$ some norm on Ω .

Coupling is a 1-1 map $\pi(x, y)$ joint distribution
marginal $\pi(\cdot, y)$ is q
marginal $\pi(x, \cdot)$ is p .

$$\lim_{t \rightarrow \infty} W_2(p_t, \bar{e}^f) \leq e^{-2\mu t} W_2(p_0, \bar{e}^f)$$

for f -strongly convex f .

Ex. X, Y Lévy Diffusions for f N -strongly

Lemma: X_0, Y_0 . X_t, Y_t Langeri Diffusions for f μ -strongly convex
 be coupled s.t. $E(\|X_t - Y_t\|^2) \leq e^{-2\mu t} \|X_0 - Y_0\|^2$

Pf.

$dX_t = -\nabla f(X_t) dt + \sqrt{2} dW_t$ $dY_t = -\nabla f(Y_t) dt + \sqrt{2} dW_t$	couple with same W_t !
--	-----------------------------

$$d(X_t - Y_t) = -(\nabla f(X_t) - \nabla f(Y_t)) dt$$

$$d\|X_t - Y_t\|^2 = -2 \langle \nabla f(X_t) - \nabla f(Y_t), X_t - Y_t \rangle dt$$

By μ -strong convexity,

$$f(Y_t) - f(X_t) \geq \langle \nabla f(X_t), Y_t - X_t \rangle + \frac{\mu}{2} \|Y_t - X_t\|^2$$

$$f(X_t) - f(Y_t) \geq \langle \nabla f(Y_t), X_t - Y_t \rangle + \frac{\mu}{2} \|X_t - Y_t\|^2$$

$$\langle \nabla f(X_t) - \nabla f(Y_t), X_t - Y_t \rangle \geq \mu \|X_t - Y_t\|^2$$

$$\therefore d\|X_t - Y_t\|^2 \leq -2\mu \|X_t - Y_t\|^2 dt$$

$$d \log \|X_t - Y_t\|^2 \leq -2\mu dt$$

$$\|X_t - Y_t\|^2 \leq e^{-2\mu t} \|X_0 - Y_0\|^2$$

$$\begin{aligned} d \log Z &\leq \alpha dt \\ \log\left(\frac{Z}{Z_0}\right) &\leq \alpha t \\ Z &\leq e^{\alpha t} Z_0. \end{aligned}$$

Next time: Proof of F-P and hence of
limiting distribution.

Main tool: A chain rule for
stochastic variables!