

Convexity

Definition: $f(\alpha x + (1-\alpha)y) \leq \alpha f(x) + (1-\alpha)f(y)$ (No Assumption)
Alternative Notion:
1. $f(y) \geq f(x) + f'(x)(y-x)$ All in Taylor theory (twice differentiable)
2. $\nabla^2 f(x) \succeq 0$, PSD (twice differentiable)
3. $\nabla f(x)$ is monotone, $\langle \nabla f(x) - \nabla f(y), x-y \rangle \geq 0$ (once differentiable)

Convex function (convex set) \rightarrow Local min = Global min
Convex function $\nabla^2 f(x) \succeq 0$, as long as $\nabla f(x) = 0 \rightarrow$ Local/Global min

Convex set: $x \in C, y \in C \rightarrow \alpha x + (1-\alpha)y \in C$ (Use in Constraint)

Gradient Descent Finds Optimality

Deriving GD Comes from ① Pick ∇ Satisfy Taylor Theorem, since that $f(\bar{x} + \mu \nabla) - f(\bar{x}) < 0$
② Make simplified assumption (Locally L-smooth)

Finite Descent

Taylor theory: $f(\bar{x} + \mu \nabla) = f(\bar{x}) + \nabla f(\bar{x})^T (\mu \nabla) + \frac{1}{2} \mu^2 \nabla^2 f(\bar{x})^T \nabla$
Should be negative
We derived how we can't remove ∇^2
Descent on GD to satisfy Taylor theorem such that $f(\bar{x} + \mu \nabla) - f(\bar{x}) < 0$
Educated Guess Interpretation

Deriving GD from Local Convexity + traditional Calculus
We assume + combine local convexity by Aroniso property

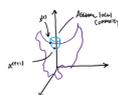
$f(\bar{x}) = f(x^{(k)}) + \nabla f(x^{(k)})^T (\bar{x} - x^{(k)}) + \frac{1}{2} (\bar{x} - x^{(k)})^T \nabla^2 f(\xi) (\bar{x} - x^{(k)})$
For seen the same reason by Taylor theory
We combine our bounds by saying that $\text{Curvature} = 1/\mu$, Locally L-smooth

Assume $\nabla^2 f(\xi) \leftarrow \frac{1}{\mu} I$ and thus

$g(z) = f(x^{(k)}) + \nabla f(x^{(k)})^T (\bar{x} - x^{(k)}) + \frac{1}{2\mu} \|\bar{x} - x^{(k)}\|^2$
We say $g(z)$ looks like $f(z)$, but importantly, it is convex, lets look with this convex function

With this assumption, we can derive GD
Just like in Calculus $\Rightarrow \nabla g(z) = 0$
 $\Rightarrow \nabla f(x^{(k)}) + \frac{1}{\mu} (\bar{x} - x^{(k)}) = 0$
 $\Rightarrow \bar{x}^* = x^{(k)} - \mu \nabla f(x^{(k)})$
Global min Can be found

This is GD Algorithm, we make an educated guess saying the function is Locally L-smooth
Which hides away complexity, then optimize on this simple function
This seems nice like EM



Optimality Guaranteed

Convex + Taylor theory + GD Protocol

① L-Lipschitz: $\|f(x) - f(y)\| \leq L \|x - y\|$
L-Lip bound (and): $\|\nabla f(x)\| \leq L$ (lip stands)
L-Lip convergence guarantee:
1. Convex
2. $\|x^{(k)} - x^*\| \leq R$
3. T Iteration
4. $\mu = \frac{R}{L}$
 $\Rightarrow f(\frac{R}{L} \nabla f(x^{(k)})) - f(x^*) \leq \frac{R^2}{2L}$
Average iteration errors bounded

② L-smooth: $\|\nabla^2 f(x)\| \leq L$
L-smooth bound (and): $0 \leq \nabla^2 f(x) \preceq L I$ if $x \in \mathcal{R}^n$ and $\forall \theta \in \mathbb{R}^n, \|\theta\| = 1$
L-smooth convergence guarantee: $f(x^{(k)}) \leq f(x^*) + \frac{L}{2} \|\nabla f(x^{(k)})\|^2$ At each step decrease
L-smooth (and) steps: $\|\nabla f(x^{(k)})\| \leq \sqrt{\frac{2(f(x^{(k)}) - f(x^*))}{L}}$ At least one $x^{(k)}$ meet satisfy this within T iterations
* strong, no convexity assumption
only differentiable + L-smooth (root makes a huge difference)

③ If we have $f(x)$ as a strongly convex and smooth
 $f(x^{(k)}) - f(x^*) \leq (1 - \frac{\mu}{L})^k (f(x^{(0)}) - f(x^*))$

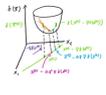
Twice Some Comments from Theoretical Perspective

Every time we check something, it is not just adding onto GD, but rather redefining complexity from a novel perspective, then finding connection to GD (solve the problem you have to solve)

Maintaining Local Convexity

How to Pick μ so we don't pass the optimality point
If we can make a projection line, we can maintain some local convexity

Aroniso Condition: $f(x^{(k+1)}) = f(x^{(k)} - \mu \nabla f(x^{(k)})) \leq f(x^{(k)}) - \mu \|\nabla f(x^{(k)})\|^2$
Convergence Guarantee: $f(x^{(k)}) - f(x^*) \leq \frac{\mu \|\nabla f(x^{(k)})\|^2}{2 \min(\mu, L)}$
and $\min(\mu, L) = \min(1, \frac{L}{\mu})$
Thus $f(x^{(k+1)}) - f(x^*) \leq \frac{\mu}{2} \|\nabla f(x^{(k)})\|^2$
At each step, error is bounded



Descent From Constraint's Perspective

We can describe a descent direction of GD with Norm in it
Taylor theorem first order expansion says that

$$\begin{cases} f(x^{(k+1)}) \approx f(x^{(k)}) + \nabla f(x^{(k)})^T (x^{(k+1)} - x^{(k)}) \\ \text{since } x^{(k+1)} = x^{(k)} - \mu \nabla f(x^{(k)}) \\ \rightarrow x^{(k+1)} - x^{(k)} = -\mu \nabla f(x^{(k)}) \\ \rightarrow f(x^{(k+1)}) + \mu \nabla f(x^{(k)})^T (-\mu \nabla f(x^{(k)})) = f(x^{(k)}) - \mu \nabla f(x^{(k)})^T \nabla f(x^{(k)}) \approx f(x^{(k)}) - \mu \|\nabla f(x^{(k)})\|^2 \end{cases}$$

This is dot product = norm

Thus, GD can be seen like:

$$f(x^{(k+1)}) = f(x^{(k)}) - \mu \|\nabla f(x^{(k)})\|^2$$

In addition, by adjusting the projection or norm that GD uses, we have different norms to GD
Norm forms a constraint on the descent

$$\begin{aligned} L_1 &\rightarrow \text{Sparse step direction (constraint distance)} \|\nabla f(x^{(k)})\|_1 & \mathcal{P}(\text{Descent direction}) &= -\text{Sign}\left(\frac{\partial f}{\partial x_i}\right) \\ L_{\infty} &\rightarrow \text{Uniform step in each dimension} \|\nabla f(x^{(k)})\|_{\infty} & \mathcal{P} &= \|\nabla f(x^{(k)})\|_{\infty} \begin{cases} \text{Sign}(f, 1) \\ \text{Sign}(f, 2) \end{cases} \end{aligned}$$

Constraint Projection: $x^{(k+1)} = x^{(k)} - \mu \nabla f(x^{(k)})$ s.t. $x \in \mathcal{R}^n$
 $y^{(k+1)} = x^{(k)} - \mu \nabla f(x^{(k)})$
then $x^{(k+1)} = \text{Proj}_{\mathcal{C}}(y^{(k+1)}) \leftarrow x^{(k+1)} = \text{argmin}_{\|y\|} \|y^{(k+1)} - z\|$

Newton's Method

Using distance way to derive GD
We treat $f(x) \approx f(x^{(k)}) + \nabla f(x^{(k)})^T (x - x^{(k)})$ as deriving GD
Now let's say $f(x) \approx f(x^{(k)}) + \nabla f(x^{(k)})^T (x - x^{(k)}) + \frac{1}{2} (x - x^{(k)})^T \nabla^2 f(x^{(k)}) (x - x^{(k)})$

Distance from GD minimization, lets set this to 0 directly with assumption that $\nabla f(x^*) = 0$
we will get a new method called Newton's method where

$$x^{(k+1)} = x^{(k)} - [\nabla^2 f(x^{(k)})]^{-1} \nabla f(x^{(k)})$$

Theoretical Support: ① $\|\nabla^2 f(x)\| \leq \frac{1}{\mu}$ $\forall \mu > 0 \rightarrow \lambda_{\min}(A) \geq \mu \rightarrow$ ensure Inverse exist
② $\|\nabla^2 f(x) - \nabla^2 f(x^*)\| \leq L \|x - x^*\| \rightarrow$ Hessian bounded

Guarantees: ① $\|x^{(k)} - x^*\| \leq \frac{2\eta}{L} \forall t$ ② $\|x^{(k)} - x^*\| \leq \frac{2\eta}{L} \|x^{(0)} - x^*\|^2$
① Convex exponentially fast if convex
② Non-Convex Mapping
 $\|x^{(k)} - x^{(k-1)}\|$ not bounded ≤ 1

Alternative Problem: $\nabla^2 f(x^{(k)}) (x - x^{(k)}) = -\nabla f(x^{(k)})$

$$\begin{matrix} \uparrow & \uparrow & \uparrow \\ A & \text{unknown } x & b \\ \text{matrix} & \text{vector} & \text{vector} \end{matrix} \Rightarrow Ax + b = 0$$
 Solving

Inverting matrix A has a extremely high computation cost:

Q.essi - Newton: Convergence $\nabla^2 f(x^{(k)}) \succeq \beta I^d$ where $\beta^{(k)}$ remains key insight (Dynamic learning rate for variable)
but Computationally cheap and blowrate
ADAM $\beta^{(k)} = \text{diag}(\nabla^2 f(x^{(k)}))$
General = $\sqrt{\text{diag}(L \text{Jacobian})}$

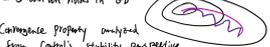
Momentum and Nesterov Acceleration

Momentum: If $-\nabla f(x^{(k)})$ is in the same direction as previous step ($x^{(k)} - x^{(k-1)}$), then move a bit further (Constructive Zeros)

If opposite direction, move a bit less (destructive Zeros)

Introducing Momentum manually
$$x^{(k+1)} = x^{(k)} - \mu \nabla f(x^{(k)}) + \beta (x^{(k)} - x^{(k-1)})$$

Smooth out noise in GD



Convergence property analyzed from Control's stability perspective

$$\begin{bmatrix} x^{(k+1)} \\ x^{(k)} \end{bmatrix} = \begin{bmatrix} 1 - \mu L + \beta & -\beta \\ \mu & 0 \end{bmatrix} \begin{bmatrix} x^{(k)} \\ x^{(k-1)} \end{bmatrix}$$

Use this to look at the eigenvalue thus deriving stability of GD w/m

For Quadratic Problems like $\frac{1}{2} x^T A x$, GD + m converges to x^* at rate of $(\frac{K-1}{K+1})^t$ where $K = \frac{\lambda_{\max}}{\lambda_{\min}}$

Nesterov Acceleration: Going a little bit more first, then take the gradient

$$\begin{aligned} y^{(k+1)} &= x^{(k)} + \beta (x^{(k)} - x^{(k-1)}) \\ x^{(k+1)} &= y^{(k+1)} - \mu \nabla f(y^{(k+1)}) \end{aligned}$$

N.A. is a less intuitive version of momentum, but it can go into continuous space and model as differential equation

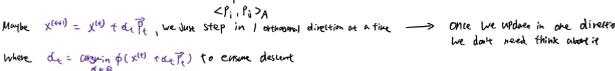
$$\text{GD} \left(\frac{K-1}{K+1} \right)^t \rightarrow \text{GD+M} \left(\frac{K-1}{K+1} \right)^t \rightarrow \left(\frac{N-1}{N} \right)^t \text{ For } \min_x f(x) \text{ where } f(x) = \frac{1}{2} x^T A x$$

Conjugate Gradient Descent

Reoptimizing Again, solve $Ax = b$, $\nabla f(x) = Ax - b$
Question: can we \odot not invert A (expensive + unstable)
① Do it step wise

Conjugate (General notion of orthogonality): $P_0^T A P_0 = 0$ $\forall i: P_0, P_1, \dots, P_{n-1}$ is the Conjugate of A (PD)
 $\langle P_i, P_j \rangle = \delta_{ij}$
Maybe $x^{(k+1)} = x^{(k)} + d_k \tilde{P}_k$, we use step in 1 orthonormal direction at a time \rightarrow Once we update in one direction, it's done for optimization, we don't need think about it

where $d_k = \text{argmin}_{d \in \mathbb{R}^n} f(x^{(k)} + d \tilde{P}_k)$ to ensure descent



Along this line, we can't find point (any part of the line to be minimized)

Let's solve this convex problem $d_k = x^* - \frac{(b - Ax^{(k)})^T P_k}{P_k^T A P_k} P_k$

Think about what this is: $b - Ax^{(k)} = -\nabla f(x^{(k)})$, then $d_k = \frac{-\nabla f(x^{(k)})^T P_k}{P_k^T A P_k} P_k$

Almost like a Gradient descent Scheme
Step we end up taking is just how P_k Project on our gradient. Take the Gradient in the Conjugate direction of P_k

How to find P_k then? $P_0^T A P_0 = 0$ is expensive to calculate
Need to dynamically Choose P_k
Compute P_k using Conjugate (next direction from only previous direction). Then along all P_0 to P_{k-2}

Like Ballon Equation idea, it is it can be conjugate to the previous vector and keep doing so, all are conjugate

Let's start with:

$$P_k = -\nabla f(x^{(k)}) + \beta_k P_{k-1}$$

Pick a new direction based on Gradient and previous direction

$$P_k^T A P_k = -P_{k-1}^T A \nabla f(x^{(k)}) + \beta_k P_{k-1}^T A P_{k-1}$$

Make this zero, then P_{k-1}^T and P_k is conjugate

$$\text{This makes } \beta_k = \frac{P_{k-1}^T A \nabla f(x^{(k)})}{P_{k-1}^T A P_{k-1}} \text{ to satisfy } P_{k-1}^T A P_k = 0$$

Start with $\beta_0 = 0$ and we pick a P_0 that we don't care what it is, then we conjugate at the next iteration

CGD also Strongly Convergent

① For $f(x)$ as quadratic, for any $x^{(k)}$, the set of step $\{x^{(k)}\}$ from CGD converge to x^* (solution) in at most n steps where $A \in \mathbb{R}^{n \times n}$

$$\|x^{(k+1)} - x^*\|_A^2 \leq \left(\frac{\lambda_{k+1} - \lambda_1}{\lambda_{k+1} + \lambda_1} \right)^2 \|x^{(k)} - x^*\|_A^2$$

At every step pick up the factor that is related to the eigenvalue

$$\|x^{(k)} - x^*\| \leq 2 \left(\frac{\sqrt{k} - 1}{\sqrt{k} + 1} \right)^k \|x^{(0)} - x^*\|_A$$

Similar Convergence Property to GD + M

Theoretical Reason: Strongly Convex

Properties of function that can make things easier and generalizable

Strongly Convex

Swap out $\nabla^2 f(x)$ in Taylor theorem with Smallest Eigenvalue $\nabla^2 f(x) \succeq cI$

$$f(x) \geq f(x^*) + \nabla f(x^*)^T (x - x^*) + \frac{c}{2} \|x - x^*\|^2 \leftarrow \text{Exponential!}$$

(Strongly Convex)

Instead of tangent line, we actually have a Tangent Curve

Size of the tangent tells you how close you are to the x^* when your condition is under strongly convex, and the relation is squared

$$\Rightarrow f(x) - f(x^*) \leq \frac{\|\nabla f(x)\|^2}{2c}$$

Strongly Convex yields very strong convergence rate

$$f(x^{(k+1)}) - f(x^*) \leq \left(1 - \frac{c}{L}\right) (f(x^{(k)}) - f(x^*)) = \left(1 - \frac{c}{L}\right)^k (f(x^{(0)}) - f(x^*))$$

When GD + ① Strong Convexity
② L-smooth $\|\nabla f(x)\| \leq L$
③ $\mu = \frac{c}{L}$

Strongest flavor of Strongly Convex is that with this condition satisfied, all previous method that we need quadratic form can generalize to any equation satisfying strongly convex

If $f(w)$ is convex and $R(w)$ is C-strongly convex, then $f(w) + R(w)$ is C-strongly convex
This give room for regularization to show its power (i.e. Ridge regression)

PL Condition

$f: \mathbb{R}^n \rightarrow \mathbb{R}$ satisfies μ -PL-condition if $\exists \alpha \in \mathbb{R}^n$

$$\frac{1}{2} \|\nabla f(x)\|_2^2 \geq \mu (f(x) - f(x^*)) \text{ Looks like Strong Convexity!}$$

When far from x^* , the gradient is big

SC \rightarrow PL, PL \nrightarrow SC

① PL-condition can hold for non-convex functions
Acts as a stand in for non-strongly-convex functions

② Importantly, if $f(x)$ is L-smooth and is μ -PL-condition, then gradient descent with $\alpha = \frac{\mu}{L}$ converges at a rate of

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

Exponentially fast, just like Strongly Convex

Notice that this does not need convex assumption, just L-smooth + μ -PL

With other Regularized Method Network, we can write the Mean Square Error as

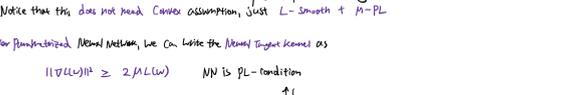
$$\|\nabla \ell(w)\|_2^2 \geq 2\mu \ell(w) \text{ NN is PL-condition}$$

And Mean Square Converges at a rate of $\left(1 - \frac{\mu}{L}\right)^k$

Overparameterized at some to converge is exponentially fast.

With PL-condition, even for non-convex function, all previous analysis can be done

Data Science is always Convex



For non-convex Problem, there are many w^* , Regularizer put a weight on the w^* with some requirement \rightarrow Small $\|w\|_2$
 \rightarrow Small $\|w\|_1$
 \rightarrow Small Entropy

More importantly, we can Reframe Constraint optimization into Regularized Unconstrained Problem

Optimal Constraint (convexness) \rightarrow optimization target f (Hinge Loss)
Optimal unconstrained (margin) \rightarrow Constraint

$$\min_{w, b} \frac{1}{2} \|w\|_2^2 \text{ Subject to Constraint of Correctness on labels}$$

Now the constraint is a function, the whole joint Minimization SUM can be reframed as:

$$\min_{w, b} \frac{1}{2} \sum_{i=1}^n (1 - y_i \langle w, x_i \rangle + b)^2 + \frac{\lambda}{2} \|w\|_2^2$$

$f(w; x)$

Convex Function \rightarrow every where you go $\nabla_w f(w; x) = 0$ PSD, no more w form at 2nd derivative