

LEARNING:

→ Reinforcement, Classification etc.

MODULE - I

LEARNING AS OPTIMIZATION:

Input: $x_1, x_2 \dots x_n \in D$
↳ Distribution

$y_1, y_2 \dots y_n \rightarrow$ Labels such that

$$H: x \rightarrow y$$

(There is a hypothesis function)

LOSS FUNCTION:

$$l: \text{Label} \times \text{Label} \rightarrow \mathbb{R}$$

Output: h such that loss is least but to
avoid overfitting we will constrain it by
parameter θ : Parameterized Hypothesis Class.
 \mathcal{H}

GOAL: EMPIRICAL RISK MINIMIZATION (ERM)

$$\min_{\theta \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^n l(h_\theta(x_i), y_i)$$

GD:

$$x_{i+1} = x_i - \eta \nabla f(x_i)$$

FUNCTIONS THAT ARE GOOD/BAD FOR GD:

LIPSCHITZNESS:

f is L -Lipschitz if

$\forall x, y$

$$|f(x) - f(y)| \leq L \|x - y\|_2$$

L_2 norm $\sqrt{\sum x_i^2}$

SMOOTHNESS:

f is β -smooth if

$\forall x, y$

$$\|\nabla f(x) - \nabla f(y)\|_2 \leq \beta \|x - y\|_2$$

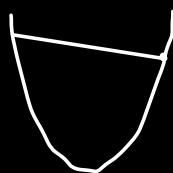
[Smoothness is a stronger constraint]

MONOTONICITY OF GD:

f is β -smooth function, if $\eta \leq 1/\beta$. Then

$$f(x_{i+1}) \leq f(x_i) - \frac{\eta}{2} \|\nabla f(x_i)\|^2$$

[Idea: So β gives an idea of how much to jump without gradient being too different]



avoids such jumps.

So monotonic!

Note: * If β is small, we say that gradient change is small \rightarrow better for SGD.
 So η can be large. Learn fast.

* Stopped when gradient = 0.

SMOOTHNESS UPPERBOUND:

* f is β -smooth $\forall a, b \quad [f: \mathbb{R} \rightarrow \mathbb{R}]$

$$f(b) \leq f(a) + f'(a) \cdot (b-a) + \frac{\beta}{2} (b-a)^2$$

* $f: \mathbb{R}^d \rightarrow \mathbb{R}$

$$\boxed{f(b) \leq f(a) + \langle \nabla f(a), b-a \rangle + \frac{\beta}{2} \|b-a\|_2^2}$$

TAYLOR'S THEOREM:

$$f(x+h) = f(x) + f'(x) \cdot h + \underbrace{\int_0^1 (f'(x+th) - f'(x)) th dt}_{\text{error term}}$$

$$f(x+h) = f(x) + f'(x) \cdot h + f''(x) \cdot \frac{h^2}{2} + \dots$$

MONOTONICITY PROOF:

UNIVARIATE:

Given:

$$* \|\nabla f(x) - \nabla f(y)\|_2 \leq \beta \|x-y\|_2 \quad [\text{Smoothness}]$$

$$* x_{i+1} = x_i - \eta f'(x_i) \quad [\text{GD}]$$

$$* f(b) \leq f(a) + f'(a) \cdot (b-a) + \frac{\beta}{2} (b-a)^2$$

[Smoothness]

upper bound]

PROOF:

$$f(x_{i+1}) = f(x_i - \eta f'(x_i))$$

$$b = x_i - \eta f'(x_i)$$

$$a = x_i$$

$$\leq f(x_i) + f'(x_i) \cdot (-\eta f'(x_i))$$

$$+ \frac{\beta}{2} (-\eta f'(x_i))^2$$

$$\leq f(x_i) - \eta f'(x_i)^2 + \frac{\eta^2 \beta}{2} f'(x_i)^2$$

$$\eta \leq 1/\beta \quad \beta \geq \frac{1}{\eta}$$

$$\leq f(x_i) - \frac{\eta}{2} \|f'(x_i)\|^2$$

MULTIVARIABLE:

$$f(x_{i+1}) = f(x_i - \eta \nabla f(x_i))$$

$$\leq f(x_i) + \langle \nabla f(x_i), -\eta \nabla f(x_i) \rangle +$$

$$\frac{\beta}{2} \|(-\eta \nabla f(x_i))\|_2^2$$

$$\leq f(x_i) - \eta \|\nabla f(x_i)\|_2^2 + \frac{\eta^2 \beta}{2} \|\nabla f(x_i)\|_2^2$$

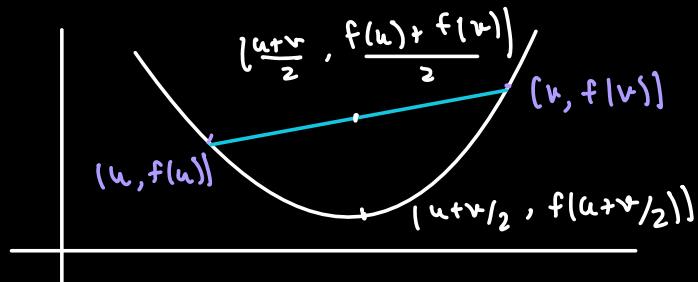
$$\leq f(x_i) - \frac{\eta}{2} \|\nabla f(x_i)\|_2^2$$

note: $\langle x, x \rangle = \|x\|_2^2$

Convex Functions: [Only one minima]

$f: \mathbb{R}^d \rightarrow \mathbb{R}$ is convex if the tangent plane

at any point is below the curve.



$\forall u, v$

$$f\left(\frac{u+v}{2}\right) \leq \frac{f(u) + f(v)}{2}$$

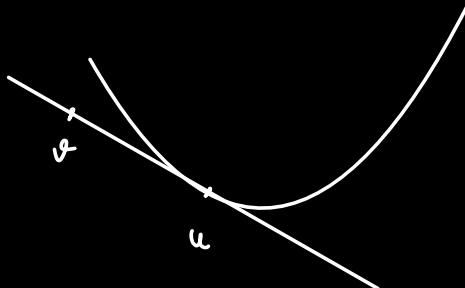
$\lambda \in (0, 1)$

$$f(\lambda u + (1-\lambda)v) \leq \lambda f(u) + (1-\lambda)f(v)$$

$$f(u) + \langle \nabla f(u), v-u \rangle \leq f(v)$$

$\underbrace{\quad \quad \quad}_{\downarrow}$

value at v along the tangent line



$$\text{value at } v = f(u) + \text{slope}(v-u)$$

$$\begin{aligned} f, g : \text{convex} &\rightarrow f+g \\ a > 0 &\qquad\qquad\qquad a.f \end{aligned} \quad \left. \begin{array}{c} \text{Convex} \\ \text{Convex} \end{array} \right\}$$

$f-g$ - need not be convex

$\Rightarrow g : \mathbb{R} \rightarrow \mathbb{R}$: convex

$$\omega \in \mathbb{R}^d$$

$$g_\omega : \mathbb{R}^d \rightarrow \mathbb{R} \quad g_\omega(x) = g(\langle \omega, x \rangle)$$

\hookrightarrow convex!

Eg: e^x , x^2 , $(x-a)^2$, $|x|$

$e^{<\omega, x>}$, $<\omega, x>^2$, $(<\omega, x> - a)^2$, $|<\omega, x>|$

ERM: $L(\theta) = \frac{1}{n} \sum_{i=1}^n l(h_\theta(x_i), y_i)$

$\underbrace{\qquad\qquad\qquad}_{a > 0} \hookrightarrow \text{Sum of convex}$

So only prove $l(h_\theta(x_i), y_i)$ is convex.

* Least Squares Regression:

Loss: $(<\theta, x_i> - y_i)^2$

\downarrow
 $(x-a)^2$

\downarrow
Convex.

PROVE CONVEX:

- * Remove positive values multiplied.
- * If sum, prove individually convex.
- * If $<\theta, x_i>$, prove for just x_i .
- * Know e^x , x^2 , $(x-a)^2$, $|x|$ are convex.

Convex - GD CONVERGENCE:

If f is β -smooth and convex, then

if $\eta \leq \gamma_\beta$

$$f(x_k) \leq f(x_*) + \frac{2\beta \cdot \|x_0 - x_*\|}{k}$$

Note: If β is low, good function,

so converge faster

If starting point close to optimum,
converge faster.

x_* → Global Optimum

x_0 → Starting point

$$f(x_k) - f(x_*) \leq \frac{2\beta \cdot \|x_0 - x_*\|}{k}$$

So Distance to optimum $\propto 1/k$

To get with ϵ of $x_k \rightarrow 1/\epsilon$

- If f was L -Lipschitz steps -

$$f(x_k) \leq f(x_*) + \frac{L \cdot \|x_0 - x_*\|}{\sqrt{k}}$$

Distance to optimum $\propto 1/\sqrt{k}$

$1/\epsilon^2$ steps

PROOF :

To prove:

$$f(x_k) \leq f(x_*) + \frac{2\beta \cdot \|x_0 - x_*\|}{k}$$

Given:

$$\star \quad f(x_{i+1}) \leq f(x_i) - \frac{\eta}{2} \|\nabla f(x_i)\|_2^2 \quad [\text{MONOTONICITY}]$$

$$\star \quad f(u) + \langle \nabla f(u), v-u \rangle \leq f(v) \quad [\text{CONVEX}]$$

\star Properties of vectors:

$$\rightarrow \|u-v\|^2 = \|u\|^2 + \|v\|^2 - 2\langle u, v \rangle$$

$$2\langle u, v \rangle = \|u\|^2 + \|v\|^2 - \|u-v\|^2$$

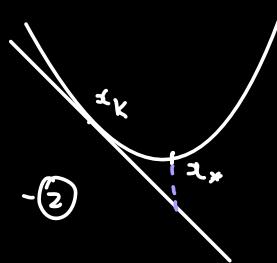
$$\rightarrow |\langle u, v \rangle| \leq \|u\| \cdot \|v\| \quad [\text{Cauchy-Schwarz}]$$

$$\star \quad \text{GD: } x_{k+1} = x_k - \frac{1}{\beta} \nabla f(x_k)$$

$$f(x_{k+1}) \leq f(x_k) - \frac{1}{2\beta} \|\nabla f(x_k)\|_2^2 \quad (1) \quad (\eta = 1/\beta)$$

$$f(x_*) \geq f(x_k) + \langle \nabla f(x_k), x_* - x_k \rangle$$

$$\Rightarrow f(x_k) \leq f(x_*) - \langle \nabla f(x_k), x_* - x_k \rangle \quad (2)$$



Combine ①, ②

$$f(x_{k+1}) \leq f(x_*) - \langle \nabla f(x_k), x_* - x_k \rangle - \frac{1}{2\beta} \|\nabla f(x_k)\|_2^2$$

$$\nabla f(x_k) = \beta (x_k - x_{k+1})$$

$$f(x_{k+1}) \leq f(x_*) + \beta \langle x_k - x_{k+1}, x_k - x_* \rangle -$$

$$\frac{\beta}{2} \| (x_k - x_{k+1}) \|_2^2$$

$$\leq f(x_*) + \frac{\beta}{2} \| x_k - x_{k+1} \|^2 + \frac{\beta}{2} \| x_k - x_* \|^2$$

$$- \frac{\beta}{2} \| x_* - x_{k+1} \|^2 - \frac{\beta}{2} \| x_k - x_{k+1} \|^2$$

$$f(x_{k+1}) \leq f(x_*) + \frac{\beta}{2} \| x_k - x_* \|^2 - \frac{\beta}{2} \| x_* - x_{k+1} \|^2$$

$$\Rightarrow f(x_1) - f(x_*) \leq \frac{\beta}{2} \| x_0 - x_* \|^2 - \frac{\beta}{2} \| x_* - x_1 \|^2$$

$$f(x_2) - f(x_*) \leq \frac{\beta}{2} \| x_1 - x_* \|^2 - \frac{\beta}{2} \| x_* - x_2 \|^2$$

⋮

$$f(x_{k+1}) - f(x_*) \leq \frac{\beta}{2} \| x_k - x_* \|^2 - \frac{\beta}{2} \| x_* - x_{k+1} \|^2$$

Add them

$$\sum_{i=1}^{k+1} [f(x_i) - f(x_*)] \leq \frac{\beta}{2} \left[\|x_0 - x_*\|^2 - \|x_{k+1} - x_*\|^2 \right]$$

As GO is monotonic

$$f(x_{k+1}) \leq f(x_i) \quad \forall i \leq k+1$$

$$f(x_{k+1}) - f(x_*) \leq f(x_i) - f(x_*)$$

$$(k+1)(f(x_{k+1}) - f(x_*)) \leq \frac{\beta}{2} \left[\|x_0 - x_*\|^2 - \|x_{k+1} - x_*\|^2 \right]$$
$$\leq \frac{\beta}{2} \|x_0 - x_*\|^2$$

$$f(x_{k+1}) - f(x_*) \leq \frac{\beta}{z(k+1)} \|x_0 - x_*\|^2 //$$

NESTEROV'S ACCELERATED GRADIENT DESCENT (NAGD):

$$f(x_k) \leq f(x_*) + \frac{2\beta \|x_0 - x_*\|^2}{k^2}$$

$\eta \leq 1/\beta$
 $\eta_i = \frac{(i+1)\eta}{2}$
 $\alpha_i = 2/(i+3)$

To be within $\epsilon = 0.01$ of optimum:

Convex

$$\left\{ \begin{array}{l} * \text{Lipschitz + GD} \Rightarrow 1/\epsilon^2 = 1/10^4 = 10000 \text{ iterations} \\ * \text{Smooth + GD} \Rightarrow 1/\epsilon = 100 \text{ iterations} \\ * \text{Smooth + NAGD} \Rightarrow 1/\sqrt{\epsilon} = 10 \text{ iterations.} \end{array} \right.$$

What about time taken for each iteration!

For LSR:

$$L(\omega) = \frac{1}{n} \sum_{i=1}^n (\langle \omega; x_i \rangle - y_i)^2$$

$$\nabla L(\omega) = \frac{2}{n} \sum_{i=1}^n (\langle \omega; x_i \rangle - y_i) x_i$$

↓
iterate over d dimensions
iterate over n examples.

So each iteration takes $O(nd)$ time.

So to be within ϵ of answer

$1/\epsilon$ iterations $\times O(nd)$ $\Rightarrow O(nd/\epsilon)$ time.
per iteration

STOCHASTIC GRADIENT DESCENT:

k elements \rightarrow Gradient

$$\text{So } \mathbf{x}_{t+1} = \mathbf{x}_t - \eta \mathbf{b}(\mathbf{x}_t)$$

\hookrightarrow gradient estimator

* b should have low variance

$$\text{var}(b) := \sum_{i=1}^k \text{var}(b_i) \quad \text{and } E[b(x)] = \nabla f(x)$$

SGD CONVERGENCE:

f : convex and β -smooth

$$\eta \leq 1/\beta, \text{ var}(b(x)) \leq \sigma^2$$

$$E[f(\tilde{\mathbf{x}}_k)] \leq f(\mathbf{x}_*) + \frac{\|\mathbf{x}_0 - \mathbf{x}_*\|_2^2}{2\eta k} + \eta \sigma^2$$



\sim
only difference.

$$\underbrace{\frac{1}{k} (\mathbf{x}_1 + \dots + \mathbf{x}_k)}$$

average of all the iteration values.

equate

$$\eta = \frac{\|\mathbf{x}_0 - \mathbf{x}_*\|_2^2}{\sigma \sqrt{k}}$$

\Rightarrow So take $\eta_k = 1/\sqrt{k}$, not $1/\beta$.

MONOTONICITY OF SGD:

Given smoothness

$$\boxed{E[f(x_{i+1})] \leq E[f(x_i)] - \frac{\eta}{2} E[\|\nabla f(x_i)\|_2^2] + \frac{\eta\sigma^2}{2}}$$

PROOF:

$$\begin{aligned} f(x_{i+1}) &= f(x_i - \eta h(x_i)) \\ &\leq f(x_i) + \langle \nabla f(x_i), -\eta h(x_i) \rangle + \frac{\beta}{2} \|\eta h(x_i)\|_2^2 \\ [f(b) &\leq f(a) + \langle \nabla f(a), b-a \rangle + \frac{\beta}{2} \|b-a\|_2^2 - \\ &\quad \text{Smoothness Upperbound}] \end{aligned}$$

$$\leq f(x_i) - \eta \langle \nabla f(x_i), h(x_i) \rangle + \frac{\eta^2 \beta}{2} \|h(x_i)\|_2^2$$

Apply expectation

$$\begin{aligned} E[f(x_{i+1})] &\leq E[f(x_i)] - \eta \langle E[\nabla f(x_i)], E[h(x_i)] \rangle \\ &\quad + \frac{\eta^2 \beta}{2} E\|E[h(x_i)]\|_2^2 \end{aligned}$$

$$\leq E[f(x_i)] - \eta E[\|\nabla f(x_i)\|_2^2] +$$

$$\frac{\eta^2 \beta}{2} E\|\nabla f(x_i)\|_2^2$$

$$+ \frac{\eta\sigma^2}{2}$$

$$\begin{aligned}
\text{var}(\hat{h}(x)) &= E\left\{\left\|h(x) - E[h(x)]\right\|_2^2\right\} \\
&= E\left[\|h(x)\|_2^2\right] - \|E[h(x)]\|_2^2 \\
&= E\left[\|h(x)\|_2^2\right] - \|\nabla f(x)\|_2^2
\end{aligned}$$

SUMMARY :

1. GD: $x_{i+1} = x_i - \eta \nabla f(x_i)$ •

2. LIPSCHITZNESS: $|f(x) - f(y)| \leq L \|x - y\|_2$

3. SMOOTHNESS: $\|\nabla f(x) - \nabla f(y)\|_2 \leq \beta \|x - y\|_2$

UPPER- BOUND:

$$f(b) \leq f(a) + \langle \nabla f(a), b - a \rangle + \frac{\beta}{2} \|b - a\|_2^2$$

4. MONOTONICITY OF GD GIVEN SMOOTH:

$$f(x_{i+1}) \leq f(x_i) - \frac{\eta}{2} \|\nabla f(x_i)\|^2$$

Proof : GD + S. Upper Bound

$$f(x_k) \leq f(x_k) \leq f(x_i) \quad \forall i \leq k$$

5. CONVEX FUNCTIONS:

$$f\left(\frac{u+v}{2}\right) \leq \frac{f(u) + f(v)}{2}$$

$$f(\lambda u + (1-\lambda)v) \leq \lambda f(u) + (1-\lambda)f(v)$$

$$\star \quad f(u) + \langle \nabla f(u), v - u \rangle \leq f(v)$$

6. PROVE CONVEX:

- * Remove positive values multiplied.
- * If sum, prove individually convex.
- * If $\langle \theta, x_i \rangle$, prove for just x_i .
- * Know $e^x, x^2, (x-a)^2, |x|$ are convex.

7. GD CONVERGENCE: Convex + Smooth

$$f(x_k) \leq f(x_*) + \frac{2\beta \cdot \|x_0 - x_*\|}{k} .$$

Proof: Monotonicity - ①

Convex Definition (x_k, x_*) - ②

① + ② + GD Replace $\nabla f(x) + \gamma_2 + \text{sum trick}$

$$f(x_k) \leq f(x_i) \quad \forall i \leq k$$

8. NGD CONVERGENCE: Convex + Smooth

$$f(x_k) \leq f(x_*) + \frac{2\beta \|x_0 - x_*\|}{k^2}$$

9. SGD: $x_{t+1} = x_t - \eta \nabla h(x_t)$

$$\mathbb{E}[h(x)] = \nabla f(x), \quad \text{var}(h(x)) \leq \sigma^2$$

10. MONOTONICITY (Smooth):

$$E[f(x_{i+1})] \leq E[f(x_i)] - \frac{\eta}{2} E[\|\nabla f(x_i)\|_2^2] + \frac{\eta\sigma^2}{2}$$

11. CONVERGENCE:

$$E[f(\tilde{x}_k)] \leq f(x_*) + \frac{\eta \|x_0 - x_*\|_2^2}{2\eta k} + \eta\sigma^2$$

12. SGD: \downarrow Iteration time

PGD: Constrained Optimization, New point outside convex set.

Project it.

$$x_{t+1} = \text{Proj}_C(x_t - \eta \nabla f(x_t))$$

\Rightarrow VECTOR PROPERTIES:

$$\forall 1: \langle u, v \rangle = \|u\|_2^2$$

$$\forall 2: 2 \langle u, v \rangle = \|u\|^2 + \|v\|^2 - \|u - v\|^2$$

$$\forall 3: |\langle u, v \rangle| \leq \|u\| \cdot \|v\| \quad [\text{Cauchy-Schwarz}]$$

HOMEWORK

$$\|\nabla f(\omega)\| \leq \epsilon$$

$$r \quad f(\omega_{i+1}) \leq f(\omega_i) - \frac{\eta}{2} \|\nabla f(\omega_i)\|^2$$

$$* \quad f(b) \leq f(a) + \langle \nabla f(a), b-a \rangle + \frac{\beta}{2} \|b-a\|_2^2$$

$$* \quad \omega_{i+1} = \omega_i - \eta \nabla f(\omega_i)$$

$$f(\omega_{k+1}) \leq f(\omega_0) - \frac{1}{2\beta} \left(\sum_{i=0}^k \|\nabla f(\omega_i)\|_2^2 \right)$$

$$\sum_{i=0}^{k+1} \|\nabla f(\omega_i)\|_2^2 \leq 2\beta (f(\omega_0) - f(\omega_k))$$

$$\leq 2\beta (f(\omega_0) - f(\omega_x))$$

We have:

$$\rightarrow f(x_{i+1}) \leq f(x_i) - \frac{1}{2\beta} \|\nabla f(x_i)\|^2$$

$$\rightarrow x_{i+1} = x_i - \eta \nabla f(x_i)$$

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

Monotonicity - ①

Convex Definition (x_k, x^*) - ③

① + ② + GD Replace $\nabla f(x) + \nabla_2 + \text{sum trick}$

$$f(x_k) \leq f(x_i) \quad i \leq k$$

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

$$f(x^*) \geq f(x_k) + \langle \nabla f(x_k), x^* - x_k \rangle$$

$$+ \frac{\alpha}{2} \|x^* - x_k\|_2^2$$

$$f(x_{k+1}) \leq f(x^*) + \beta \langle x_k - x_{k+1}, x_k - x^* \rangle$$

$$- \frac{\alpha}{2} \|x^* - x_k\|_2^2 - \frac{\beta}{2} \|x_k - x_{k+1}\|_2^2$$

$$f(x_{k+1}) \leq f(x^*) + \frac{\beta}{2} \|x_k - x_{k+1}\|^2$$

$$+ \frac{\beta}{2} \|x_k - x^*\|_2^2 - \frac{\beta}{2} \|x^* - x_{k+1}\|_2^2$$

$$+ \frac{\alpha}{2} \|x^* - x_k\|_2^2 - \frac{\beta}{2} \|x_k - x_{k+1}\|_2^2$$

$$\begin{aligned} f(x_{k+1}) &\leq f(x^*) + \underbrace{\frac{(\beta-\alpha)}{2} \|x_k - x^*\|_2^2}_{-} \\ &\quad - \frac{\beta}{2} \|x^* - x_{k+1}\|_2^2 \end{aligned}$$

$$\begin{aligned} f(x_{k+1}) &\geq f(x^*) + \underbrace{\langle \nabla f(x^*), x_{k+1} - x^* \rangle}_{0} + \frac{\alpha}{2} \|x_{k+1} - x^*\|_2^2 \end{aligned}$$

$$\geq f(x^*) + \frac{\alpha}{2} \|x_{k+1} - x^*\|_2^2$$

$$f(x^*) + \frac{\alpha}{2} \|x_{k+1} - x^*\|_2^2 \leq f(x^*) + \underbrace{\frac{(\beta-\alpha)}{2} \|x_k - x^*\|_2^2}_{- \frac{\beta}{2} \|x_{k+1} - x^*\|_2^2}$$

$$- \frac{\beta}{2} \|x_{k+1} - x^*\|_2^2$$

$$\frac{(\alpha + \beta)}{2} \|x_{k+1} - x^*\|^2 \leq \frac{(\beta - \alpha)}{\sum} \|x_k - x^*\|^2$$

$$\frac{\beta - \alpha}{\beta}$$