

Large Scale Optimization for Machine Learning:

How to Use Hessian?

Meisam Razaviyayn

razaviya@usc.edu

Using Hessian in Optimization Algorithms

- Using Hessian is beneficial if possible (why?)
- Hessian is huge when the dimension is large
- How should we use Hessian in large size problems?
 - (Block) diagonal Newton
 - Conjugate gradient method
 - Use Hessian-Vector product approximation

Conjugate Directions

$$\mathbf{Q} \in \mathbb{R}^{n \times n}$$

Start from simple convex positive definite quadratic optimization:

$$\min_{\mathbf{x}} \frac{1}{2} \mathbf{x}^T \mathbf{Q} \mathbf{x} - \mathbf{b}^T \mathbf{x}$$

Solution: $\mathbf{Q} \mathbf{x}^* = \mathbf{b}$

Definition: Directions $\mathbf{d}^0, \mathbf{d}^1, \dots, \mathbf{d}^k$ are \mathbf{Q} – conjugate if

$$(\mathbf{d}^i)^T \mathbf{Q} \mathbf{d}^j = \mathbf{0}, \quad \forall i, j \text{ with } i \neq j$$

Lemma: If a set of vectors are \mathbf{Q} -conjugate, then they are linearly independent

There are n \mathbf{Q} -conjugate directions and they form a (non-orthogonal) basis

Finding Solution Using Conjugate Directions

Idea: Expand using conjugate directions:

$$\mathbf{x}^* = \sum_{i=1}^n \alpha^i \mathbf{d}^i \quad \text{where } \alpha^k = \frac{(\mathbf{d}^k)^T \mathbf{b}}{(\mathbf{d}^k)^T \mathbf{Q} \mathbf{d}^k}, \forall k$$

$$\min_{\mathbf{x}} \frac{1}{2} \mathbf{x}^T \mathbf{Q} \mathbf{x} - \mathbf{b}^T \mathbf{x}$$

Solution: $\mathbf{Q} \mathbf{x}^* = \mathbf{b}$

Alternative approach for computing coefficients:

Start from $\mathbf{x}^0 = \mathbf{0}$

for $i = 0, 1, \dots, n - 1$

$$\alpha^i = \arg \min_{\alpha} f(\mathbf{x}^i + \alpha \mathbf{d}^i)$$

$$\mathbf{x}^{i+1} = \mathbf{x}^i + \alpha^i \mathbf{d}^i$$

return $\mathbf{x}^* = \mathbf{x}^n$

This iterative procedure finds the solution in finite number of iterations!

figure

How can we find \mathbf{Q} -conjugate directions?

Generating Q-conjugate Directions

Gram-Schmidt Procedure:

Given linearly independent vectors ξ^0, \dots, ξ^{n-1}

generate Q-conjugate \mathbf{d}^k 's such that $\xi^k = \sum_{i=0}^{k-1} \beta_k^i \mathbf{d}^i + \mathbf{d}^k, \forall k$

In other words,

$$\mathbf{d}^0 = \xi^0$$
$$\mathbf{d}^k = \xi^k - \sum_{i=0}^{k-1} \frac{(\mathbf{d}^i)^T \mathbf{Q} \xi^k}{(\mathbf{d}^i)^T \mathbf{Q} \mathbf{d}^i} \mathbf{d}^i$$

A Simple Recursion for Generating Q-conjugate Directions

Start with $\xi^0 = -\nabla f(\mathbf{x}^0), \dots, \xi^{n-1} = -\nabla f(\mathbf{x}^{n-1})$

Gram-Schmidt Procedure:

$$\mathbf{d}^0 = -\nabla f(\mathbf{x}^0)$$

$$\mathbf{d}^k = -\nabla f(\mathbf{x}^k) + \frac{(\nabla f(\mathbf{x}^k))^T (\nabla f(\mathbf{x}^k) - \nabla f(\mathbf{x}^{k-1}))}{(\nabla f(\mathbf{x}^{k-1}))^T (\nabla f(\mathbf{x}^{k-1}))} \mathbf{d}^{k-1}$$

for $k = 1, \dots, n - 1$

Conjugate Gradient Method for Quadratic Optimization

Start with $\mathbf{d}^0 = -\nabla f(\mathbf{x}^0)$

for $k = 1, 2, \dots, n$

$$\alpha^{k-1} = \arg \min_{\alpha} f(\mathbf{x}^{k-1} + \alpha \mathbf{d}^{k-1})$$

$$\mathbf{x}^k = \mathbf{x}^{k-1} + \alpha^{k-1} \mathbf{d}^{k-1}$$

$$\mathbf{d}^k = -\nabla f(\mathbf{x}^k) + \frac{(\nabla f(\mathbf{x}^k))^T (\nabla f(\mathbf{x}^k) - \nabla f(\mathbf{x}^{k-1}))}{(\nabla f(\mathbf{x}^{k-1}))^T (\nabla f(\mathbf{x}^{k-1}))} \mathbf{d}^{k-1}$$

return $\mathbf{x}^* = \mathbf{x}^n$

Solves quadratic problems
in exactly n iterations

For general optimization:

- It is similar to approximating the function with second order Taylor expansion
- The **Hessian is changing**
- Do not stop
- **Restart** at some $m < n$ (when not enough decrease)

Other ways to use Hessian information

For Newton step, we need to compute $(\nabla^2 f(\mathbf{x}))^{-1} \nabla f(\mathbf{x})$

- Approach 1: Approximate Hessian with diagonal or block diagonal matrices
- Approach 2: Use Hessian-vector product

$$\nabla^2 f(\mathbf{x})\mathbf{v} \approx \frac{\nabla f(\mathbf{x} + \alpha\mathbf{v}) - \nabla f(\mathbf{x})}{\alpha} \text{ for small } \alpha$$

Large Scale Optimization for Machine Learning:

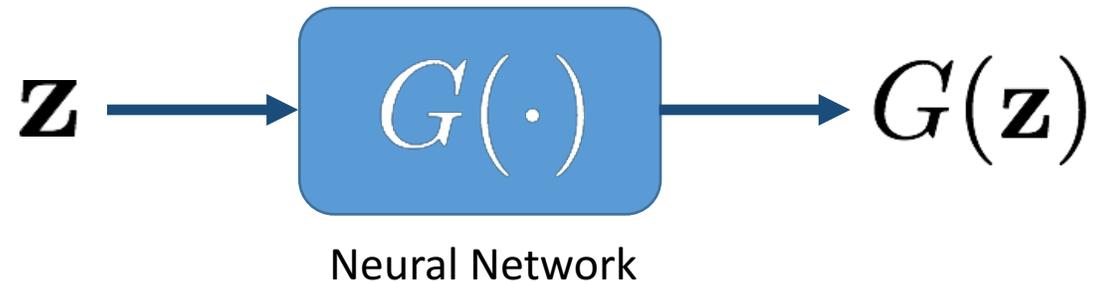
GANs and Min-Max Problems

Meisam Razaviyayn

razaviya@usc.edu

Generative Adversarial Networks

Goal: Generate samples that look like real samples $\mathbf{x}_1, \dots, \mathbf{x}_n \sim \mathbb{P}_x$

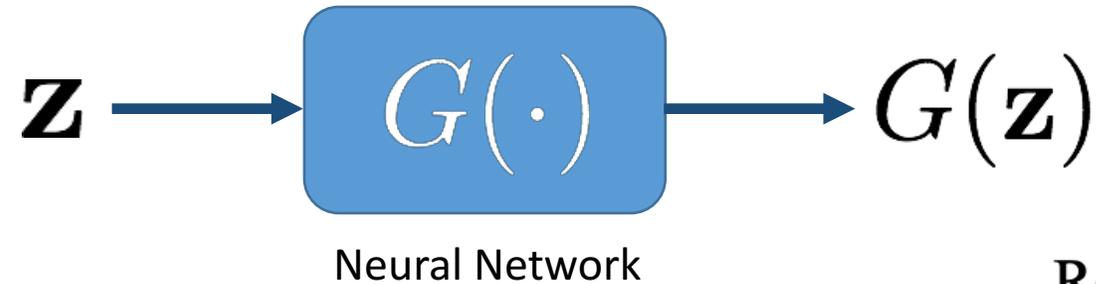


\mathbf{z} : input random variable, e.g., $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$

We need $G(\mathbf{z})$ to have the same distribution as \mathbb{P}_x

$$\min_G \text{dist}(G(\mathbf{z}), \mathbf{x})$$

Generative Adversarial Networks

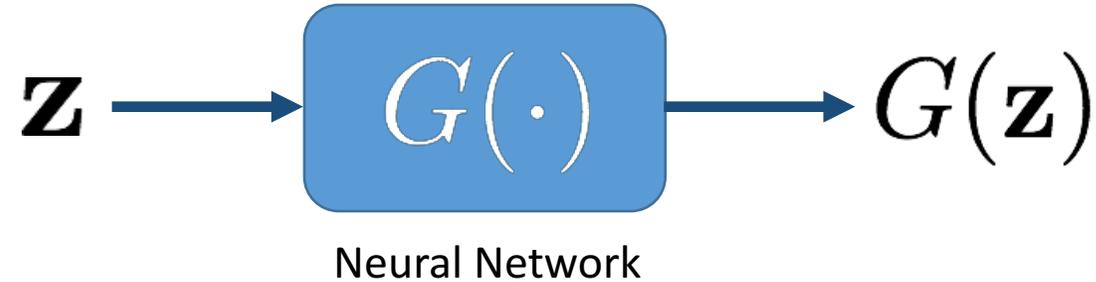


Applications of GANs:

- Simulated environments and training data
- Dealing with missing data [Dai et al. 2017]
- Super-resolution



Generative Adversarial Networks



Applications of GANs:

- Simulated environments and training data
- Dealing with missing data
- Super-resolution
- Text-to-image synthesis

This bird is white with some black on its head and wings, and has a long orange beak



This bird has a yellow belly and tarsus, grey back, wings, and brown throat, nape with a black face



This flower has overlapping pink pointed petals surrounding a ring of short yellow filaments



GAN Formulation



$$\min_G \text{dist}(G(\mathbf{z}), \mathbf{x})$$

How to define a distance between two distributions?

Many “distances” between distributions:

KL-divergence, Jensen-Shannon divergence, TV distance, Wasserstein divergence,...

Kullback-Leibler Divergence

Definition:
$$\text{KL}(P\|Q) = \int P(x) \log \left(\frac{P(x)}{Q(x)} \right) dx = \mathbb{E}_{x \sim P} \left[\log \left(\frac{P(x)}{Q(x)} \right) \right]$$

Discrete case:
$$\text{KL}(P\|Q) = \sum_i p_i \log \left(\frac{p_i}{q_i} \right)$$

This divergence is zero if and only if $P = Q$ almost everywhere.

Drawbacks:

- Not clear how to optimize it based on samples for GANs
- It is a non-symmetric divergence

Jensen-Shannon Divergence

Definition:
$$\text{JS}(P\|Q) = \frac{1}{2}\text{KL}(P\|M) + \frac{1}{2}\text{KL}(Q\|M)$$

where
$$M = \frac{P + Q}{2}.$$

Drawback: Not clear how to minimize it based on samples.

Lemma:
$$\text{JS}(P\|Q) = 1 + \frac{1}{2} \max_{D \in \mathbb{D}} \mathbb{E}_{x \sim P} [\log D(x)] + \mathbb{E}_{y \sim Q} [\log(1 - D(y))]$$

where $\mathbb{D} =$ set of all functions with range $(0, 1)$

Proof

More appropriate formulation for optimization

Wasserstein Distance

Definition: $W(P\|Q) = \min_{\pi \in \Pi} \int c(x, y) \pi(x, y) dy dx$

where $\Pi = \{ \pi(x, y) \mid \int \pi(x, y) dy = P(x), \int \pi(x, y) dx = Q(y) \}$

Discrete case: $W(P\|Q) = \min_{\pi} \sum_{x, y} c(x, y) \pi(x, y)$

s.t. $\pi(x, y) \geq 0, \sum_y \pi(x, y) = P(x), \forall x, \sum_x \pi(x, y) = Q(y), \forall y$

Intuition and earth mover's distance

Wasserstein Distance

Definition: $W(P\|Q) = \min_{\pi \in \Pi} \int c(x, y) \pi(x, y) dy dx$

where $\Pi = \{\pi(x, y) \mid \int \pi(x, y) dy = P(x), \int \pi(x, y) dx = Q(y)\}$

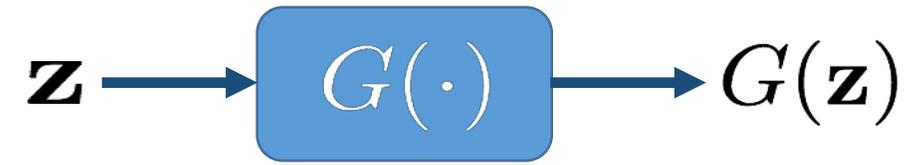
Issue: Need to work with probabilities directly and it also requires $O(n^2)$ variables

Dual formulation:
$$W(P\|Q) = \max_{\gamma, \lambda} \mathbb{E}_{x \sim P} [\gamma(x)] + \mathbb{E}_{y \sim Q} [\lambda(y)]$$
$$\text{s.t. } \gamma(x) + \lambda(y) \leq c(x, y), \forall x, y$$

Proof

When $c(x, y) = \|x - y\|_2$,
$$W(P\|Q) = \max_{\gamma} \mathbb{E}_{x \sim P} [\gamma(x)] - \mathbb{E}_{y \sim Q} [\gamma(y)]$$
$$\text{s.t. } \gamma(x) - \gamma(y) \leq \|x - y\|_2, \forall x, y$$

GAN Formulation



$$\min_G \text{dist}(G(\mathbf{z}), \mathbf{x})$$

Jensen-Shannon GANs:
$$\min_G \max_{D \in \mathbb{D}} \mathbb{E}_{\mathbf{x}} [\log D(\mathbf{x})] + \mathbb{E}_{\mathbf{z}} \log (1 - D(G(\mathbf{z})))$$

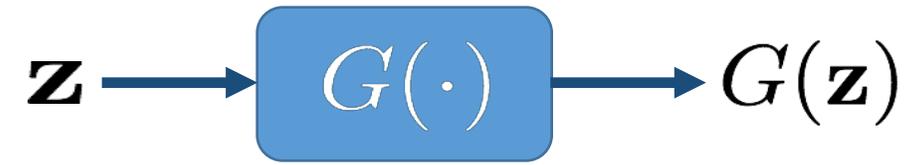
$\mathbb{D} =$ set of all functions with range $(0, 1)$

Wasserstein GANs:
$$\min_G \max_{\gamma} \mathbb{E}_{\mathbf{x}} [\gamma(\mathbf{x})] - \mathbb{E}_{\mathbf{z}} [\gamma(G(\mathbf{z}))]$$

s.t.
$$\gamma(\mathbf{x}) - \gamma(\mathbf{y}) \leq \|\mathbf{x} - \mathbf{y}\|_2, \forall \mathbf{x}, \mathbf{y}$$

- How to solve in practice over the set of functions?
 - Parameterized G, γ, D with neural network
- Why do we call it Generative Adversarial Networks? Discriminator vs. Generator

GAN Formulation



Parameterizing the functions:

$$\min_G \text{dist}(G(\mathbf{z}), \mathbf{x})$$

Jensen-Shannon GANs:
$$\min_{\theta} \max_{\beta} \mathbb{E}_{\mathbf{x}} [\log D_{\beta}(\mathbf{x})] + \mathbb{E}_{\mathbf{z}} \log (1 - D_{\beta}(G_{\theta}(\mathbf{z})))$$

$$\text{range}(D_{\beta}(\cdot)) = (0, 1)$$

Wasserstein GANs:
$$\min_{\theta} \max_{\beta} \mathbb{E}_{\mathbf{x}} [\gamma_{\beta}(\mathbf{x})] - \mathbb{E}_{\mathbf{z}} [\gamma_{\beta}(G_{\theta}(\mathbf{z}))]$$

s.t.
$$\gamma_{\beta}(\mathbf{x}) - \gamma_{\beta}(\mathbf{y}) \leq \|\mathbf{x} - \mathbf{y}\|_2, \forall \mathbf{x}, \mathbf{y}$$

How to make it unconstrained optimization?

Resulting Optimization Problem

$$\min_{\theta} \max_{\beta} f(\theta, \beta)$$

Many other GAN formulations as well

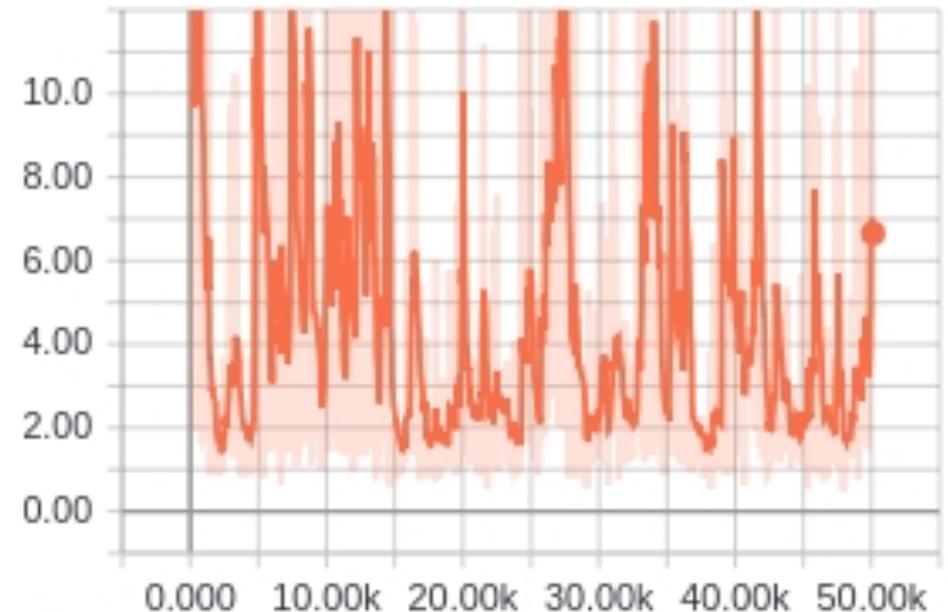
➤ How to solve?

➤ Gradient Descent/Ascent used in practice

$$\theta_{t+1} = \theta_t - \alpha \nabla_{\theta} f(\theta_t, \beta_t)$$

$$\beta_{t+1} = \beta_t + \alpha \nabla_{\beta} f(\theta_t, \beta_t)$$

➤ Very unstable in practice (even its variants)



Should GDA Work?

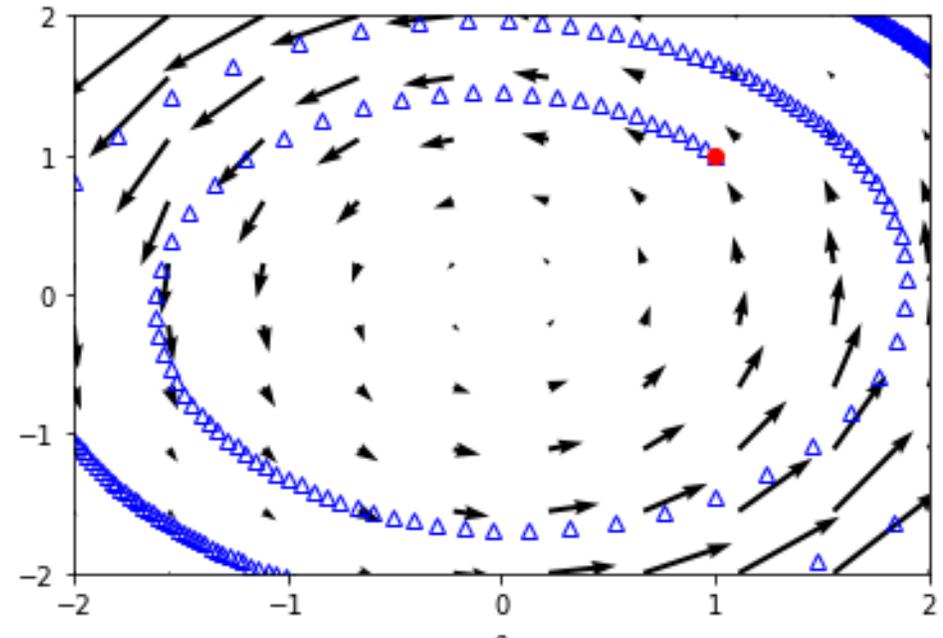
$$\min_{\theta} \max_{\beta} f(\theta, \beta)$$

Simple example: $f(\theta, \beta) = \theta\beta$

$$\theta_{t+1} = \theta_t - \alpha \nabla_{\theta} f(\theta_t, \beta_t)$$

$$\beta_{t+1} = \beta_t + \alpha \nabla_{\beta} f(\theta_t, \beta_t)$$

- Where should the algorithm converge to?
- Does the algorithm converge?



How to solve?

$$\min_{\theta} \max_{\beta} f(\theta, \beta)$$

- What should we hope for computing?
 - Finding the global min?
 - Finding a first order stationary point with respect to θ, β
 - **Open problem** in the general case.

Another Point of View

- Can we do gradient descent on $g(\theta)$?
- Find a point that $\|\nabla g(\theta)\|$ is small enough?
- Danskin's theorem requires solving the max problem exactly

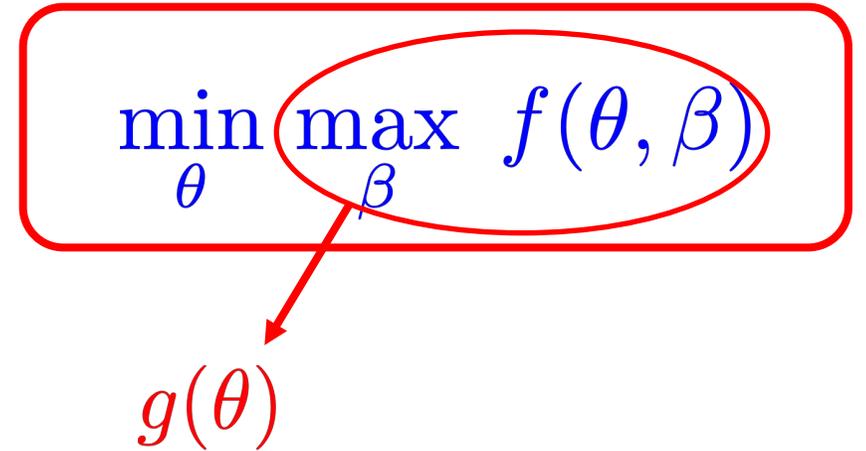
$$\theta_{t+1} = \theta_t - \alpha \nabla g(\theta_t)$$



$$\beta_{t+1} = \arg \max_{\beta} f(\theta_t, \beta)$$

$$\theta_{t+1} = \theta_t - \alpha \nabla_{\theta} f(\theta_t, \beta_{t+1})$$

- In practice, we may not be able to solve the max
 - No hope for this approach when non-concave w.r.t. β . Even under concavity, we cannot solve the inner problem exactly



Another Point of View

- Simpler case: what happens if $f(\theta, \beta)$ is strongly concave in β ?

$$\min_{\theta} \max_{\beta} f(\theta, \beta)$$

$g(\theta)$

$$\theta_{t+1} = \theta_t - \alpha \nabla g(\theta_t)$$

$$\beta_{t+1} = \arg \max_{\beta} f(\theta_t, \beta)$$
$$\theta_{t+1} = \theta_t - \alpha \nabla_{\theta} f(\theta_t, \beta_{t+1})$$

for $t = 1, 2, \dots,$

$$\tilde{\beta}_0 = \beta_t$$

for $i = 0, 1, \dots, T - 1$

$$\tilde{\beta}_{i+1} = \tilde{\beta}_i + \alpha \nabla_{\beta} f(\theta_t, \tilde{\beta}_i)$$

$$\beta_{t+1} = \tilde{\beta}_T$$

$$\theta_{t+1} = \theta_t - \alpha \nabla_{\theta} f(\theta_t, \beta_{t+1})$$

Theorem: Designing T and step-size carefully, the algorithm converges to ϵ -stationarity in polynomial time.

Can be extended to stochastic setting and using mini-batches

Thank You!