

HW2 Handwritten Assignment Solution

Lecturer: Pei-Yuan Wu

TAs: Yuan-Chia Chang(Problem 5), Chun-Lin Huang(Problem 1, 2, 3, 4)

November 2023, First Edition

Problem 1 (Convolution)(0.5%)

As mentioned in class, image size may change after convolution layers. Consider a batch of image data with shape $(B, W, H, input_channels)$, how will the shape change after the following convolution layer:

$Conv2D(input_channels, output_channels, kernel_size = (k_1, k_2), stride = (s_1, s_2), padding = (p_1, p_2))$

For simplicity, the padding tuple means that p_1 pixels are padded on both left and right sides, and p_2 pixels are padded on both top and bottom sides.

Solution: The output is $(B, W', H', output_channels)$, where

$$W' = \left\lfloor \frac{W + 2 * p_1 - k_1}{s_1} + 1 \right\rfloor$$
$$H' = \left\lfloor \frac{H + 2 * p_2 - k_2}{s_2} + 1 \right\rfloor$$

Note that you should give some explanation to get all points.

Problem 2 (Batch Normalization)(1%)

Batch normalization [1] is a popular trick for training deep networks nowadays, which aims to preserve the distribution within hidden layers and avoids vanishing gradient issue. The algorithm can be written as below:

Algorithm 1 Batch Normalization

Input Feature from data points over a mini-batch $B = (x_i)_{i=1}^m$
Output $y_i = BN_{\gamma, \beta}(x_i)$

```

1: procedure BATCHNORMALIZE( $B, \gamma, \beta$ )
2:    $\mu_B \leftarrow \frac{1}{m} \sum_{i=1}^m x_i$  ▷ mini-batch mean
3:    $\sigma_B^2 \leftarrow \frac{1}{m} \sum_{i=1}^m (x_i - \mu_B)^2$  ▷ mini-batch variance
4:   for  $i \leftarrow 1$  to  $m$  do
5:      $\hat{x}_i \leftarrow \frac{x_i - \mu_B}{\sqrt{\sigma_B^2 + \epsilon}}$  ▷ normalize
6:      $y_i \leftarrow \gamma \hat{x}_i + \beta$  ▷ scale and shift
7:   end for
8:   return
9: end procedure

```

During training we need to backpropagate the gradient of loss ℓ through this transformation, as well as compute the gradients with respect to the parameters γ, β . Towards this end, please write down the close form expressions for $\frac{\partial \ell}{\partial x_i}$, $\frac{\partial \ell}{\partial \gamma}$, $\frac{\partial \ell}{\partial \beta}$ in terms of $x_i, \mu_B, \sigma_B^2, \hat{x}_i, y_i$ (given by the forward pass) and $\frac{\partial \ell}{\partial y_i}$ (given by the backward pass).

- Hint: You may first write down the close form expressions of $\frac{\partial \ell}{\partial \hat{x}_i}$, $\frac{\partial \ell}{\partial \sigma_B^2}$, $\frac{\partial \ell}{\partial \mu_B}$, and then use them to compute $\frac{\partial \ell}{\partial x_i}$, $\frac{\partial \ell}{\partial \gamma}$, $\frac{\partial \ell}{\partial \beta}$.

Solution: We use the gradient descent to update γ and β :

$$\begin{aligned}\gamma &\leftarrow \gamma - \eta \frac{\partial \ell}{\partial \gamma} \\ \beta &\leftarrow \beta - \eta \frac{\partial \ell}{\partial \beta}\end{aligned}$$

where η is a learning rate and

$$\begin{aligned}\frac{\partial \ell}{\partial \gamma} &= \sum_{i=1}^m \left(\frac{\partial \ell}{\partial y_i} \frac{\partial y_i}{\partial \gamma} \right) = \sum_{i=1}^m \frac{\partial \ell}{\partial y_i} \hat{x}_i \\ \frac{\partial \ell}{\partial \beta} &= \sum_{i=1}^m \left(\frac{\partial \ell}{\partial y_i} \frac{\partial y_i}{\partial \beta} \right) = \sum_{i=1}^m \frac{\partial \ell}{\partial y_i}\end{aligned}$$

Note that we sum from 1 to m because we are working with mini-batches. Now,

we derive some important terms by chain rule:

$$\begin{aligned}
\frac{\partial l}{\partial \hat{x}_i} &= \frac{\partial l}{\partial y_i} \frac{\partial y_i}{\partial \hat{x}_i} = \frac{\partial l}{\partial y_i} \gamma \\
\frac{\partial l}{\partial \sigma_B^2} &= \frac{\partial l}{\partial \hat{x}_i} \frac{\partial \hat{x}_i}{\partial \sigma_B^2} = -\frac{1}{2} \sum_{i=1}^m \frac{\partial l}{\partial \hat{x}_i} (x_i - \mu_B) (\sigma_B^2 + \epsilon)^{-\frac{3}{2}} \\
\frac{\partial l}{\partial \mu_B} &= \frac{\partial l}{\partial \hat{x}_i} \frac{\partial \hat{x}_i}{\partial \mu_B} \\
&= \sum_{i=1}^m \frac{\partial l}{\partial \hat{x}_i} \frac{\partial}{\partial \mu_B} (x_i - \mu_B) (\sigma_B^2 + \epsilon)^{-\frac{1}{2}} \\
&= \sum_{i=1}^m \frac{\partial l}{\partial \hat{x}_i} \left[\frac{\partial (x_i - \mu_B)}{\partial \mu_B} (\sigma_B^2 + \epsilon)^{-\frac{1}{2}} + (x_i - \mu_B) \frac{\partial (\sigma_B^2 + \epsilon)^{-\frac{1}{2}}}{\partial \mu_B} \right] \\
&= \sum_{i=1}^m \frac{\partial l}{\partial \hat{x}_i} \left[\frac{-1}{\sqrt{\sigma_B^2 + \epsilon}} + (x_i - \mu_B) \frac{\partial (\sigma_B^2 + \epsilon)^{-\frac{1}{2}}}{\partial \mu_B} \right]
\end{aligned}$$

We know $\sigma_B^2 = \frac{1}{m} \sum_{i=1}^m (x_i - \mu_B)^2$, so the second term in bracket is

$$\begin{aligned}
(x_i - \mu_B) \frac{\partial (\sigma_B^2 + \epsilon)^{-\frac{1}{2}}}{\partial \mu_B} &= (x_i - \mu_B) \frac{-1}{2} (\sigma_B^2 + \epsilon)^{-\frac{3}{2}} \frac{\partial (\sigma_B^2 + \epsilon)}{\partial \mu_B} \\
&= \frac{-1}{2} (x_i - \mu_B) (\sigma_B^2 + \epsilon)^{-\frac{3}{2}} \frac{\partial \left(\frac{1}{m} \sum_{i=1}^m (x_i - \mu_B)^2 + \epsilon \right)}{\partial \mu_B} \\
&= \frac{-1}{2} (x_i - \mu_B) (\sigma_B^2 + \epsilon)^{-\frac{3}{2}} \left(\frac{-2}{m} \sum_{i=1}^m (x_i - \mu_B) \right)
\end{aligned}$$

Hence,

$$\begin{aligned}
\frac{\partial l}{\partial \mu_B} &= \sum_{i=1}^m \frac{\partial l}{\partial \hat{x}_i} \frac{-1}{\sqrt{\sigma_B^2 + \epsilon}} + \underbrace{\sum_{i=1}^m \frac{\partial l}{\partial \hat{x}_i} \frac{-1}{2} (x_i - \mu_B) (\sigma_B^2 + \epsilon)^{-\frac{3}{2}} \left(\frac{-2}{m} \sum_{i=1}^m (x_i - \mu_B) \right)}_{\frac{\partial l}{\partial \sigma_B^2}} \\
&= \sum_{i=1}^m \frac{\partial l}{\partial \hat{x}_i} \frac{-1}{\sqrt{\sigma_B^2 + \epsilon}} + \frac{\partial l}{\partial \sigma_B^2} \left(\frac{-2}{m} \sum_{i=1}^m (x_i - \mu_B) \right)
\end{aligned}$$

To derive $\frac{\partial l}{\partial x_i}$, we use the chain rule $\frac{\partial l}{\partial x_i} = \frac{\partial l}{\partial \hat{x}_i} \frac{\partial \hat{x}_i}{\partial x_i} + \frac{\partial l}{\partial \sigma_B^2} \frac{\partial \sigma_B^2}{\partial x_i} + \frac{\partial l}{\partial \mu_B} \frac{\partial \mu_B}{\partial x_i}$. Now, calculate the remaining term:

$$\begin{aligned}
\frac{\partial \hat{x}_i}{\partial x_i} &= \frac{1}{\sqrt{\sigma_B^2 + \epsilon}} \\
\frac{\partial \mu_B}{\partial x_i} &= \frac{1}{m} \\
\frac{\partial \sigma_B^2}{\partial x_i} &= \frac{2(x_i - \mu)}{m}
\end{aligned}$$

That is,

$$\frac{\partial l}{\partial x_i} = \frac{\partial l}{\partial \hat{x}_i} \frac{1}{\sqrt{\sigma_B^2 + \epsilon}} + \frac{\partial l}{\partial \sigma_B^2} \frac{2(x_i - \mu)}{m} + \frac{1}{m} \frac{\partial l}{\partial \mu_B}$$

Problem 3 (Constrained Mahalanobis Distance Minimization Problem)(1.5%)

1. Let $\Sigma \in R^{m \times m}$ be a symmetric positive semi-definite matrix, $\mu \in R^m$. Please construct a set of points $x_1, \dots, x_n \in R^m$ such that

$$\frac{1}{n} \sum_{i=1}^n (x_i - \mu)(x_i - \mu)^T = \Sigma, \quad \frac{1}{n} \sum_{i=1}^n x_i = \mu$$

- Find the relation between set of points and (μ, Σ) and (μ, Σ) is known
Solution: WLOG (Without Loss of Generality), let $\mu = 0$. Since Σ is a symmetric positive semi-definite matrix, we can perform eigen decomposition as follows:

$$\Sigma = UDU^T = \sum_{i=1}^m (d_i u_i u_i^T).$$

where U and U^T are orthogonal matrix. Let $n = 2m$ and construct a set of points $x_1, \dots, x_m, \dots, x_{2m}$ where $x_i = \sqrt{d_i} u_i$ and $x_{m+i} = -\sqrt{d_i} u_i \forall 1 \leq i \leq m$. Then,

$$\begin{aligned} \frac{1}{n} \sum_{i=1}^n x_i &= \mu = 0 \\ \frac{1}{n} \sum_{i=1}^n x_i x_i^T &= \sum_{i=1}^m (d_i u_i u_i^T) = UDU^T = \Sigma \end{aligned}$$

Note that lots of students just use the covariance of eigen decomposition to construct $\{x_i\}_{i=1}^n$. However, It should satisfy the condition that $\frac{1}{n} \sum_{i=1}^n x_i = \mu$

2. Let $1 \leq k \leq m$, solve the following optimization problem (and justify with proof):
 minimize $\text{Trace}(\Phi^T \Sigma \Phi)$
 subject to $\Phi^T \Phi = I_k$
 variables $\Phi \in R^{m \times k}$

Solution: Let ϕ_1, \dots, ϕ_k be the columns of Φ . Then

$$\text{Trace}(\Phi^T \Sigma \Phi) = \sum_{i=1}^k \phi_i^T \Sigma \phi_i = \sum_{i=1}^k \phi_i^T \left(\sum_{j=1}^m (d_j u_j u_j^T) \right) \phi_i = \sum_{j=1}^m d_j \sum_{i=1}^k \langle u_j, \phi_i \rangle^2 = \sum_{j=1}^m c_j d_j$$

where $\langle \cdot, \cdot \rangle$ is standard inner product in Euclidean space, $c_j := \sum_{i=1}^k \langle u_j, \phi_i \rangle^2$ for each $j = 1, \dots, m$ and $d_1 \geq d_2 \geq \dots \geq d_m \geq 0$ Claim: $0 \leq c_j \leq 1$ and $\sum_{j=1}^m c_j = k$ Clearly, $c_j \geq 0$. Extending ϕ_1, \dots, ϕ_k to $\phi_1, \dots, \phi_k, \phi_{k+1}, \dots, \phi_m$ for \mathbb{R}^m . Then, for each $j = 1, \dots, m$

$$c_j = \sum_{i=1}^k \langle u_j, \phi_i \rangle^2 \leq \sum_{i=1}^m \langle u_j, \phi_i \rangle^2 = 1$$

Finally, since u_1, \dots, u_m is an orthonormal basis for \mathbb{R}^m ,

$$\sum_{j=1}^m c_j = \sum_{j=1}^m \sum_{i=1}^k \langle u_j, \phi_i \rangle^2 = \sum_{i=1}^k \sum_{j=1}^m \langle u_j, \phi_i \rangle^2 = \sum_{i=1}^k \|\phi_i\|_2^2 = k$$

Hence, the minimum value of $\sum_{j=1}^m c_j d_j$ over all choice of $c_1, c_2, \dots, c_m \in [0, 1]$ with $\sum_{j=1}^m c_j = k$ is d_{m-k+1}, \dots, d_m . This is achieved when $c_1, \dots, c_{m-k} = 0$ and $c_{m-k+1} = \dots = c_m = 1$

References

- [1] Sergey Ioffe and Christian Szegedy (2015), “Batch Normalization: Accelerating Deep Network Training by Reducing Internal Covariate Shift”, Arxiv:1502.03167

Problem 4 (Convergence of K-means Clustering) (1.5%)

In the K-means clustering algorithm, we are given a set of n points $x_i \in \mathbb{R}^d, i \in \{1, \dots, n\}$ and we want to find the centers of k clusters $\boldsymbol{\mu} = (\mu_1, \dots, \mu_k)$ by minimizing the average distance from the points to the closest cluster center. In general, $n \geq k$. Define function $\mathcal{C} : \{1, \dots, n\} \rightarrow \{1, 2, \dots, k\}$ assigns one of k clusters to each point in the data set such that $\mathcal{C}(i) = q$ if the i -th data point x_i is assigned to the q -th cluster where $i \in \{1, 2, \dots, n\}$ and $q \in \{1, 2, \dots, k\}$

Formally, we want to minimize the following loss function

$$L(\mathcal{C}, \boldsymbol{\mu}) = \sum_{i=1}^n \|x_i - \mu_{\mathcal{C}(i)}\|_2^2 = \sum_{q=1}^k \sum_{i:\mathcal{C}(i)=q} \|x_i - \mu_q\|_2^2$$

The K-means algorithm:

Algorithm 2 K-means algorithm

Initialize cluster center $\mu_j, j = 1, 2, \dots, k$ (k random x_n from data set)

Repeat:

1. Fix $\boldsymbol{\mu}$, update $\mathcal{C}(i)$ for each i that minimizes L . Formally, consider a data point x_i , and let $\mathcal{C}(i)$ be the assignment from the previous iteration and $\mathcal{C}^*(i)$ be the new assignment obtained as: $\mathcal{C}^*(i) = \arg \min_{j=1, \dots, k} \|x_i - \mu_j\|_2^2$
2. Fix \mathcal{C} , update the centers μ_j which satisfies

$$|\{i : \mathcal{C}(i) = j\}| \mu_j = \sum_{i:\mathcal{C}(i)=j} x_i,$$

for each j , where $|\{i : \mathcal{C}(i) = j\}|$ is the number of elements of set $\{i : \mathcal{C}(i) = j\}$. (i.e. Set the cluster centres to be the means of the points in each cluster.)

The algorithm stops when no change in loss function occurs during the assignment step.

Suppose that the algorithm proceeds from iteration t to $t + 1$.

1. Consider the points z_1, z_2, \dots, z_m , where $m \geq 1$. and for $i \in \{1, 2, \dots, m\}, z_i \in \mathbb{R}^d$. Let $\bar{z} = \frac{1}{m} \sum_{i=1}^m z_i$ be the mean of these points, and let $z \in \mathbb{R}^d$ be an arbitrary point in the same (d -dimensional) space. Then

$$\sum_{i=1}^m \|z_i - z\|_2^2 \geq \sum_{i=1}^m \|z_i - \bar{z}\|_2^2$$

Solution:

$$\begin{aligned}
\sum_{i=1}^m \|z^i - z\|^2 &= \sum_{i=1}^m \|(z^i - \bar{z}) + (\bar{z} - z)\|^2 \\
&= \sum_{i=1}^m \left(\|z^i - \bar{z}\|^2 + \|\bar{z} - z\|^2 + 2(z^i - \bar{z}) \cdot (\bar{z} - z) \right) \\
&= \sum_{i=1}^m \|z^i - \bar{z}\|^2 + \sum_{i=1}^m \|\bar{z} - z\|^2 + 2 \sum_{i=1}^m (z^i \cdot \bar{z} - z^i \cdot z - \bar{z} \cdot \bar{z} + \bar{z} \cdot z) \\
&= \sum_{i=1}^m \|z^i - \bar{z}\|^2 + m\|\bar{z} - z\|^2 + 2(m\bar{z} \cdot \bar{z} - m\bar{z} \cdot z - m\bar{z} \cdot \bar{z} + m\bar{z} \cdot z) \\
&= \sum_{i=1}^m \|z^i - \bar{z}\|^2 + m\|\bar{z} - z\|^2 \\
&\geq \sum_{i=1}^m \|z^i - \bar{z}\|^2.
\end{aligned}$$

2. Show that $L(\mathcal{C}^{t+1}, \mu^t) \leq L(\mathcal{C}^t, \mu^t)$ i.e. The first step in K-means clustering

Solution: It follows directly from the logic of the algorithm: \mathcal{C}^t and \mathcal{C}^{t+1} are different only if there is a point that finds a closer cluster center in μ^t than the one assigned to it by \mathcal{C}^t :

$$L(\mathcal{C}^{t+1}, \mu^t) = \sum_{i=1}^n \|x^i - \mu_{\mathcal{C}^{t+1}(i)}^t\|^2 < \sum_{i=1}^n \|x^i - \mu_{\mathcal{C}^t(i)}^t\|^2 = L(\mathcal{C}^t, \mu^t)$$

3. Show that $L(\mathcal{C}^{t+1}, \mu^{t+1}) \leq L(\mathcal{C}^{t+1}, \mu^t)$ i.e. The second step in K-means clustering. (Hint: Use the result of (a)) *Solution:* Use the result in (1):

$$\begin{aligned}
L(\mathcal{C}^{t+1}, \mu^{t+1}) &= \sum_{i=1}^n \|x^i - \mu_{\mathcal{C}^{t+1}(i)}^{t+1}\|^2 \\
&= \sum_{k'=1}^k \sum_{i \in \{1, 2, \dots, n\}, \mathcal{C}^{t+1}(i)=k'} \|x^i - \mu_{\mathcal{C}^{t+1}(i)}^{t+1}\|^2 \\
&\leq \sum_{k'=1}^k \sum_{i \in \{1, 2, \dots, n\}, \mathcal{C}^{t+1}(i)=k'} \|x^i - \mu_{\mathcal{C}^{t+1}(i)}^t\|^2 \text{ (by 1)} \\
&= \sum_{i=1}^n \|x^i - \mu_{\mathcal{C}^{t+1}(i)}^t\|^2 \\
&= L(\mathcal{C}^{t+1}, \mu^t).
\end{aligned}$$

4. Use the result in (b) and (c) to show that the loss of k -means clustering algorithm is monotonic decreasing. (Hint: Show that the sequence $\{l_t\}$, where $l_t = L(\mathcal{C}^t, \boldsymbol{\mu}^t)$, which is monotone decreasing ($l_{t+1} \leq l_t, \forall t$) and bounded below ($l_t \geq 0$). Then, we use monotone convergence theorem for sequences, $\{l_t\}$ converges.) *Solution:* Define the sequence $\{l_t\}$, where $l_t = L(\mathcal{C}^t, \boldsymbol{\mu}^t)$. By previous results, we have

$$l_t = L(\mathcal{C}^t, \boldsymbol{\mu}^t) \leq L(\mathcal{C}^{t+1}, \boldsymbol{\mu}^{t+1}) = l_{t+1}$$

for all t . Hence, $\{l_t\}$ is a monotonic decreasing sequence. Note that we apply **monotonic convergence theorem of sequence** to prove the sequence is convergence, which does not guarantee this algorithm could find the **global** minimum, just a **local** minimum.

5. Show that the k -means clustering algorithm converges in finitely many steps. *Solution:* There are at most k^N ways to partition N data points into k clusters. Then, this algorithm converges in finitely many steps. Note that the upper bound (k^N) may not tight.

Problem 5 (Gradient Descent Convergence) (1.5%)

Suppose the function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is differentiable. Also, f is β -smoothness and α -strongly convex.

$$\beta\text{-smoothness} : \beta > 0, \forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^n, \|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\|_2 \leq \beta \|\mathbf{x} - \mathbf{y}\|_2$$

$$\alpha\text{-strongly convex} : \alpha > 0, \forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^n, f(\mathbf{x}) - f(\mathbf{y}) - \nabla f(\mathbf{y})^T(\mathbf{x} - \mathbf{y}) \geq \frac{\alpha}{2} \|\mathbf{x} - \mathbf{y}\|_2^2$$

Then we propose a gradient descent algorithm

1. Find a initial $\boldsymbol{\theta}^0$.
2. Let $\boldsymbol{\theta}^{n+1} = \boldsymbol{\theta}^n - \eta \nabla_{\boldsymbol{\theta}} f(\boldsymbol{\theta}^n) \quad \forall n \geq 0$, where $\eta = \frac{1}{\beta}$.

The following problems lead you to prove the gradient descent convergence.

1. Prove the property of β -smoothness function

$$\forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^n, f(\mathbf{x}) - f(\mathbf{y}) - \nabla f(\mathbf{y})^T(\mathbf{x} - \mathbf{y}) \leq \frac{\beta}{2} \|\mathbf{x} - \mathbf{y}\|_2^2$$

- (a) Define $g : \mathbb{R} \rightarrow \mathbb{R}, g(t) = f(\mathbf{y} + t(\mathbf{x} - \mathbf{y}))$. Show that $f(\mathbf{x}) - f(\mathbf{y}) = \int_0^1 g'(t) dt$.

$$\text{Solution: } f(\mathbf{x} - \mathbf{y}) = g(1) - g(0) = \int_0^1 g'(t) dt.$$

- (b) Show that $g'(t) = \nabla f(\mathbf{y} + t(\mathbf{x} - \mathbf{y}))^T(\mathbf{x} - \mathbf{y})$.

$$\text{Solution: Let } \mathbf{z} = \mathbf{y} + t(\mathbf{x} - \mathbf{y}).$$

$$g'(t) = \frac{dg(t)}{dt} = \frac{df(\mathbf{z})}{dt} = \sum_{i=1}^n \frac{df(\mathbf{z})}{dz_i} \frac{dz_i}{dt} = \sum_{i=1}^n \frac{df(\mathbf{z})}{dz_i} \frac{d(\mathbf{y}_i + t(\mathbf{x}_i - \mathbf{y}_i))}{dt} = \sum_{i=1}^n \frac{df(\mathbf{z})}{dz_i} (\mathbf{x}_i - \mathbf{y}_i) = \nabla f(\mathbf{z})^T(\mathbf{x} - \mathbf{y}) = \nabla f(\mathbf{y} + t(\mathbf{x} - \mathbf{y}))^T(\mathbf{x} - \mathbf{y})$$

- (c) Show that $|f(\mathbf{x}) - f(\mathbf{y}) - \nabla f(\mathbf{y})^T(\mathbf{x} - \mathbf{y})| \leq \int_0^1 |(\nabla f(\mathbf{y} + t(\mathbf{x} - \mathbf{y})) - \nabla f(\mathbf{y}))^T(\mathbf{x} - \mathbf{y})| dt$.

Solution:

$$\begin{aligned} |f(\mathbf{x}) - f(\mathbf{y}) - \nabla f(\mathbf{y})^T(\mathbf{x} - \mathbf{y})| &= \left| \int_0^1 g'(t) dt - \nabla f(\mathbf{y})^T(\mathbf{x} - \mathbf{y}) \right| \\ &= \left| \int_0^1 \nabla f(\mathbf{y} + t(\mathbf{x} - \mathbf{y}))^T(\mathbf{x} - \mathbf{y}) dt - \int_0^1 \nabla f(\mathbf{y})^T(\mathbf{x} - \mathbf{y}) dt \right| \\ &= \left| \int_0^1 (\nabla f(\mathbf{y} + t(\mathbf{x} - \mathbf{y}))^T - \nabla f(\mathbf{y})^T)(\mathbf{x} - \mathbf{y}) dt \right| \end{aligned}$$

- (d) By Cauchy-Schwarz inequality and the definition of β -smoothness, show that $|f(\mathbf{x}) - f(\mathbf{y}) - \nabla f(\mathbf{y})^T(\mathbf{x} - \mathbf{y})| \leq \frac{\beta}{2} \|\mathbf{x} - \mathbf{y}\|_2^2$, hence we get

$$f(\mathbf{x}) - f(\mathbf{y}) - \nabla f(\mathbf{y})^T(\mathbf{x} - \mathbf{y}) \leq \frac{\beta}{2} \|\mathbf{x} - \mathbf{y}\|_2^2$$

Solution: By Cauchy-Schwarz inequality,

$$\begin{aligned} & \left| \int_0^1 (\nabla f(\mathbf{y} + t(\mathbf{x} - \mathbf{y}))^T - \nabla f(\mathbf{y})^T)(\mathbf{x} - \mathbf{y}) dt \right| \\ & \leq \int_0^1 |(\nabla f(\mathbf{y} + t(\mathbf{x} - \mathbf{y}))^T - \nabla f(\mathbf{y})^T)(\mathbf{x} - \mathbf{y})| dt \\ & \leq \int_0^1 |(\nabla f(\mathbf{y} + t(\mathbf{x} - \mathbf{y})) - \nabla f(\mathbf{y}))^T(\mathbf{x} - \mathbf{y})| dt \\ & \leq \int_0^1 \|(\nabla f(\mathbf{y} + t(\mathbf{x} - \mathbf{y})) - \nabla f(\mathbf{y}))\|_2 \|(\mathbf{x} - \mathbf{y})\|_2 dt \quad \text{By Cauchy-Schwarz inequality} \\ & \leq \int_0^1 \beta \|t(\mathbf{x} - \mathbf{y})\|_2 \|(\mathbf{x} - \mathbf{y})\|_2 dt \leq \int_0^1 t \beta \|(\mathbf{x} - \mathbf{y})\|_2^2 dt = \beta \|(\mathbf{x} - \mathbf{y})\|_2^2 \int_0^1 t dt = \frac{\beta}{2} \|(\mathbf{x} - \mathbf{y})\|_2^2 \end{aligned}$$

2. Let $\mathbf{y} = \mathbf{x} - \frac{1}{\beta} \nabla f(\mathbf{x})$ and use 1., prove that

$$f(\mathbf{x} - \frac{1}{\beta} \nabla f(\mathbf{x})) - f(\mathbf{x}) \leq -\frac{1}{2\beta} \|\nabla f(\mathbf{x})\|_2^2$$

and

$$f(\mathbf{x}^*) - f(\mathbf{x}) \leq -\frac{1}{2\beta} \|\nabla f(\mathbf{x})\|_2^2,$$

where $\mathbf{x}^* = \arg \min_{\mathbf{x}} f(\mathbf{x})$.

Solution: By the fact that $f(\mathbf{y}) - f(\mathbf{x}) - \nabla f(\mathbf{x})^T(\mathbf{y} - \mathbf{x}) \leq \frac{\beta}{2} \|\mathbf{x} - \mathbf{y}\|_2^2$,

$$f(\mathbf{x}) - f(\mathbf{x} - \frac{1}{\beta} \nabla f(\mathbf{x})) - \nabla f(\mathbf{x})^T \frac{1}{\beta} \nabla f(\mathbf{x}) \leq \frac{1}{2\beta} \|\nabla f(\mathbf{x})\|_2^2.$$

Also, $\nabla f(\mathbf{x})^T \frac{1}{\beta} \nabla f(\mathbf{x}) = \frac{1}{\beta} \|\nabla f(\mathbf{x})\|_2^2$. Hence,

$$f(\mathbf{x} - \frac{1}{\beta} \nabla f(\mathbf{x})) - f(\mathbf{x}) \leq -\frac{1}{2\beta} \|\nabla f(\mathbf{x})\|_2^2$$

Owing to the fact that $f(\mathbf{x}^*) \leq f(\mathbf{x} - \frac{1}{\beta} \nabla f(\mathbf{x}))$, we get

$$f(\mathbf{x}^*) - f(\mathbf{x}) \leq -\frac{1}{2\beta} \|\nabla f(\mathbf{x})\|_2^2.$$

3. Show that $\forall n \geq 0$,

$$\|\theta^{n+1} - \theta^*\|_2^2 = \|\theta^n - \theta^*\|_2^2 + \eta^2 \|\nabla_{\theta} f(\theta^n)\|_2^2 - 2\eta \nabla_{\theta} f(\theta^n)^T (\theta^n - \theta^*),$$

where $\theta^* = \arg \min_{\theta} f(\theta)$.

Solution:

$$\begin{aligned} \|\theta^{n+1} - \theta^*\|_2^2 &= \|\theta^n - \eta \nabla_{\theta} f(\theta^n) - \theta^*\|_2^2 \\ (\theta^n - \eta \nabla_{\theta} f(\theta^n) - \theta^*)^T (\theta^n - \eta \nabla_{\theta} f(\theta^n) - \theta^*) &= \|\theta^n - \theta^*\|_2^2 + \eta^2 \|\nabla_{\theta} f(\theta^n)\|_2^2 - 2\eta \nabla_{\theta} f(\theta^n)^T (\theta^n - \theta^*) \end{aligned}$$

4. Use 2. and the definition of α -strongly convex to prove $\forall n \geq 0$

$$\|\theta^{n+1} - \theta^*\|_2^2 \leq (1 - \frac{\alpha}{\beta}) \|\theta^n - \theta^*\|_2^2,$$

where $\theta^* = \arg \min_{\theta} f(\theta)$.

Solution: By rearranging the inequality we get in 2., $\eta^2 \|\nabla f(\theta^n)\|_2^2 \leq 2\beta \eta^2 (f(\theta^*) - f(\theta^n))$.

Also, by the definition of α -strongly convex, $-2\eta \nabla_{\theta} f(\theta^n)^T (\theta^n - \theta^*) \leq -\alpha \eta \|\theta^n - \theta^*\|_2^2 - 2\eta (f(\theta^*) - f(\theta^n))$ Combined the above results and $\eta = \frac{1}{\beta}$, we get

$$\|\theta^{n+1} - \theta^*\|_2^2 \leq \|\theta^n - \theta^*\|_2^2 + \frac{2}{\beta} (f(\theta^*) - f(\theta^n)) - \frac{\alpha}{\beta} \|\theta^n - \theta^*\|_2^2 - \frac{2}{\beta} (f(\theta^*) - f(\theta^n)) = (1 - \frac{\alpha}{\beta}) \|\theta^n - \theta^*\|_2^2.$$

5. Use the above inequality to show that θ^n will converge to θ^* when n goes to infinity.

Solution: Because $\frac{\alpha}{\beta} > 0$, also, by the property of β -smoothness function and the definition of α strongly convex function, we get $\alpha \leq \beta$, which derives $\frac{\alpha}{\beta} < 1$ and $0 < 1 - \frac{\alpha}{\beta} < 1$. When n goes to infinity, $\|\theta^n - \theta^*\| = 0$. Hence, θ^n will converge to θ^*

Version Description

1. First Edition: Finish Problem Solution 1 to 5