

MID-TERM EXAM

YangyangLi

li002252@umn.edu

1 Problem 1

1.1 a

Answers:

$g(z) = |z|$ that is convex function but non-differentiable at point 0. So using subdifferential to generalize gradient of differential functions at $z = 0$.

$$\partial g(z_0) = \{x \in \mathbb{R} : g(z) \geq g(z_0) + x(z - z_0)\}$$

Where $z_0 = 0$ $g(z_0) = 0$.

So,

$$\begin{aligned}\partial g(z_0) &= \{x \in \mathbb{R} : g(z) \geq g(z_0) + x(z - z_0)\} \\ &= \{x \in \mathbb{R} : g(z) \geq g(z_0) + xz\} \\ &= \{x \in \mathbb{R} : |z| \geq xz\} \\ &= \{x \in \mathbb{R} : -1 \leq x \leq 1\}\end{aligned}$$

Also, $\partial_z |z| = 1$ when $z > 0$. $\partial_z |z| = -1$ when $z < 0$. Hence,

$$\partial_z |z| = \begin{cases} 1 & z > 0 \\ -1 & z < 0 \\ [-1, 1] & z = 0 \end{cases}.$$

1.2 b

Answers:

As,

$$f(x) = \frac{1}{2} \|y - Ax\|_2^2 + \lambda \|x\|_1$$

Let $g(x) = \|y - Ax\|_2^2$, its gradient is $A^T(Ax - y)$, which is shown in previous HWs. Concretely,

$$\begin{aligned}\|y - Ax\|_2^2 &= \|Ax - y\|_2^2 \\ &\rightarrow \|A(x + \delta) - y\|_2^2 \\ &= \|Ax - y\|_2^2 + 2 \langle A^T(Ax - y), \delta \rangle + o(\|\delta\|_2)\end{aligned}$$

So, $\partial_x g(x) = 2A^T(Ax - y)$. As $A \in \mathbb{R}^{N \times d}$, $y \in \mathbb{R}^N$ and $x \in \mathbb{R}^{d \times 1}$

$$\partial_{\mathbf{x}} g(\mathbf{x}) = \begin{bmatrix} \frac{\partial g}{\partial x_1} \\ \frac{\partial g}{\partial x_2} \\ \vdots \\ \frac{\partial g}{\partial x_d} \end{bmatrix}$$

Let $h(\mathbf{x}) = \|\mathbf{x}\|_1$,

$$\partial_{\mathbf{x}} h(\mathbf{x}) = \begin{bmatrix} \frac{\partial h}{\partial x_1} \\ \frac{\partial h}{\partial x_2} \\ \vdots \\ \frac{\partial h}{\partial x_d} \end{bmatrix}$$

As shown in 1.1, for every scalar x_i , $\frac{\partial h}{\partial x_i} = \text{sign}(x_i)$. So,

$$\partial_{\mathbf{x}} h(\mathbf{x}) = \begin{bmatrix} \text{sign}(x_1) \\ \text{sign}(x_2) \\ \vdots \\ \text{sign}(x_d) \end{bmatrix}$$

Consequently,

$$\begin{aligned} \partial_{\mathbf{x}} f(\mathbf{x}) &= \partial_{\mathbf{x}} g(\mathbf{x}) + \partial_{\mathbf{x}} h(\mathbf{x}) \\ &= \mathbf{A}^T (\mathbf{A}\mathbf{x} - \mathbf{y}) + \lambda \text{sign}(\mathbf{x}) \end{aligned}$$

Where $\text{sign}(\mathbf{x})$ means applying the sign function elementwise.

1.3 c

Answers:

Let $f(\mathbf{z}) = \|\mathbf{z}\|_2$. It is not differentiable when $\mathbf{z} = \mathbf{0}$. So,

$$\begin{aligned} \partial_{\mathbf{z}} f(\mathbf{z}_0) &= \left\{ \mathbf{x} \in \mathbb{R}^N : f(\mathbf{z}) \geq f(\mathbf{z}_0) + \langle \mathbf{x}, \mathbf{z} - \mathbf{z}_0 \rangle \right\} \\ &= \left\{ \mathbf{x} \in \mathbb{R}^N : \|\mathbf{z}\| \geq \mathbf{x}^T \mathbf{z} \right\} \\ &= \left\{ \mathbf{x} \in \mathbb{R}^N : \|\mathbf{z}\| \leq 1 \right\} \end{aligned}$$

Where $\mathbf{z}_0 = \mathbf{0}$.

When \mathbf{z} does not equal to $\mathbf{0}$,

$$\partial_{\mathbf{z}} f(\mathbf{z}) = \begin{bmatrix} \frac{\partial f(\mathbf{z})}{\partial z_1} \\ \frac{\partial f(\mathbf{z})}{\partial z_2} \\ \vdots \\ \frac{\partial f(\mathbf{z})}{\partial z_n} \end{bmatrix}$$

Where $\frac{\partial f(z)}{\partial z_i} = \frac{z_i}{\|z\|}$.
Hence,

$$\partial_z \|z\|_2 = \begin{cases} \frac{z}{\|z\|} & z \neq \mathbf{0} \\ \left\{ \mathbf{x} \in \mathbb{R}^N : \|\mathbf{x}\| \leq 1 \right\} & z = \mathbf{0} \end{cases}$$

1.4 d

1.4.1 i

Answers:

I implement Lasso to estimate x on the [Colab Notebook](#) that is uploaded as well. Please feel free to check that.

1.4.2 ii

Answers:

Let $g(x) = \|y - Ax\|_2$. As proved in the questions 1.3, its subgradient is shown below in terms of chain rule.

$$g(x) = \begin{cases} \frac{(1^T(A \otimes (Ax - y)))^T}{\|y - Ax\|_2} & y - Ax \neq 0 \\ \{\|y - Ax\|_2 \leq 1\} & y - Ax = 0 \end{cases}$$

Where $A \otimes (Ax - y)$ is broadcasting operation.

I implement Lasso to estimate x on the [Colab Notebook](#) that is uploaded as well. Please feel free to check that.

1.5 iii

Answers:

I fix the $\lambda = 1e - 2$, and try $\sigma^2 = 0.7, 2, 4$. I find that both Lasso and square-root Lasso have similar recovery performance of groundtruth x_0 . Also, The higher the variance is, the worse the performance will be.

I finish this part on the [Colab Notebook](#) that is uploaded as well. Please feel free to check that.

2 Problem 2

2.1 a

Answers:

$$\min_{w \in \mathbb{H}} \sum_{i=1}^N (\langle w, \Phi(x_i) \rangle - y_i)^2 + \lambda \|w\|_{\mathbb{H}}^2 = \min_{w \in \mathbb{H}} \|\Phi(X)w - y\|^2 + \lambda \|w\|_{\mathbb{H}}^2$$

Where $\Phi(\mathbf{X}) = \begin{bmatrix} \Phi(\mathbf{x}_1^T) \\ \Phi(\mathbf{x}_2^T) \\ \vdots \\ \Phi(\mathbf{x}_N^T) \end{bmatrix}$. the gradient is,

$$\nabla_{\mathbf{w}} = 2\Phi(\mathbf{X})^T (\Phi(\mathbf{X}) \mathbf{w} - \mathbf{y}) + 2\lambda \mathbf{w}$$

the global minimizer is \mathbf{w}_* ,

$$(\Phi(\mathbf{X})^T \Phi(\mathbf{X}) + \lambda \mathbf{I}) \mathbf{w}_* = \Phi(\mathbf{X})^T \mathbf{y}$$

For any nonzero vector \mathbf{u} ,

$$\mathbf{u}^T (\Phi(\mathbf{X})^T \Phi(\mathbf{X}) + \lambda \mathbf{I}) \mathbf{u} = \|\Phi(\mathbf{X}) \mathbf{u}\|^2 + \lambda \|\mathbf{u}\|^2 > 0$$

So $\Phi(\mathbf{X})^T \Phi(\mathbf{X}) + \lambda \mathbf{I}$ is positive definite, which implies that it is invertible.

$$\mathbf{w}_* = (\Phi(\mathbf{X})^T \Phi(\mathbf{X}) + \lambda \mathbf{I})^{-1} \Phi(\mathbf{X})^T \mathbf{y} = \Phi(\mathbf{X})^T (\Phi(\mathbf{X}) \Phi(\mathbf{X})^T + \lambda \mathbf{I})^{-1} \mathbf{y}$$

Rewriting $\mathbf{w}_* = \sum_{i=1}^N \alpha_i \Phi(\mathbf{x}_i)$, $\alpha = (\Phi(\mathbf{X}) \Phi(\mathbf{X})^T + \lambda \mathbf{I})^{-1} \mathbf{y} = (\mathbf{G} + \lambda \mathbf{I})^{-1} \mathbf{y}$ where \mathbf{G} is the Gram matrix generated from K on our training data points $\{\mathbf{x}_i\}_{i=1}^N$.

Hence, the regularized regression problem can be solved by find the optimal α . Also, α only relates to \mathbf{G} and \mathbf{y}

2.2 b

Answers:

The new optimization problem need to find the optimal α . As proved in question 2.1, $\alpha = (\Phi(\mathbf{X}) \Phi(\mathbf{X})^T + \lambda \mathbf{I})^{-1} \mathbf{y} = (\mathbf{G} + \lambda \mathbf{I})^{-1} \mathbf{y}$. Obviously, it is unique global minimizer.

2.3 c

Answers:

the find predictor is $\langle \mathbf{w}_*, \Phi(\mathbf{x}) \rangle$. The matrix form is,

$$\Phi(\hat{\mathbf{X}}) \mathbf{w}_* = \Phi(\hat{\mathbf{X}}) \Phi(\mathbf{X})^T (\Phi(\mathbf{X}) \Phi(\mathbf{X})^T + \lambda \mathbf{I})^{-1} \mathbf{y} = \hat{\mathbf{G}} (\mathbf{G} + \lambda \mathbf{I})^{-1} \mathbf{y}$$

Where train data $\Phi(\mathbf{X}) = \begin{bmatrix} \Phi(\mathbf{x}_1^T) \\ \Phi(\mathbf{x}_2^T) \\ \vdots \\ \Phi(\mathbf{x}_N^T) \end{bmatrix}$. test data $\Phi(\hat{\mathbf{X}}) = \begin{bmatrix} \Phi(\mathbf{x}_1^T) \\ \Phi(\mathbf{x}_2^T) \\ \vdots \end{bmatrix}$. \mathbf{G} is the Gram matrix

generated from K on our training data points $\{\mathbf{x}_i\}_{i=1}^N$. $\hat{\mathbf{G}}$ is the Gram matrix generated from K on our testing data points and training data points.

3 Problem 3

3.1 a

Answers:

The problem is convex because $\frac{1}{2} \|\mathbf{w}\|_2^2$, $-v\rho$ and $\sum_{i=1}^N \xi_i$ are all convex. Their Hessian matrix $\nabla^2 \geq 0$. The positive combination $\frac{1}{2} \|\mathbf{w}\|_2^2 - v\rho + \sum_{i=1}^N \xi_i$ is convex. Also, the constraints is linear. So, the problem is convex.

3.2 b

Answers:

For the soft-margin SVM, there exist \mathbf{w}_0 , b_0 and ξ_i satisfying,

$$y_i (\langle \mathbf{w}_0, \mathbf{x}_i \rangle + b_0) \geq 1 - \xi_i \quad \forall i$$

So we can find $\lambda > 0$ so that,

$$y_i (\langle \lambda \mathbf{w}_0, \mathbf{x}_i \rangle + \lambda b_0) > \rho - \xi_i \quad \forall i$$

Where $\xi_i > 0$, $\rho > 0$.

So strict feasibility can be verified and we can invoke the KKT optimality condition. The Lagrangian is,

$$\mathcal{L}(\mathbf{w}, b, \rho, \boldsymbol{\xi}, \boldsymbol{\lambda}, \mathbf{u}, \pi) = \frac{1}{2} \|\mathbf{w}\|_2^2 - v\rho + \sum_{i=1}^N \xi_i + \sum_{i=1}^N \lambda_i (\rho - \xi_i - y_i (\langle \mathbf{w}, \mathbf{x}_i \rangle + b)) - \sum_{i=1}^N u_i \xi_i - \pi \rho$$

The KKT condition is,

stationarity:

$$\mathbf{w} = \sum_{i=1}^N \lambda_i y_i \mathbf{x}_i \quad \sum_{i=1}^N \lambda_i y_i = 0 \quad \mathbf{1} = \boldsymbol{\lambda} + \mathbf{u} \quad \sum_{i=1}^N \lambda_i = \pi + v$$

feasibility:

$$y_i (\langle \mathbf{w}, \mathbf{x}_i \rangle + b) \geq \rho - \xi_i, \quad \xi_i \geq 0, \quad \lambda_i \geq 0, \quad u_i \geq 0 \quad \forall i \quad \rho \geq 0, \quad \pi \geq 0$$

complementary slackness:

$$\lambda_i (\rho - \xi_i - y_i (\langle \mathbf{w}, \mathbf{x}_i \rangle + b)) = 0, \quad u_i \xi_i = 0 \quad \forall i \quad \pi \rho = 0$$

3.3 c

Answers:

the support vectors satisfy $y_i (\langle \mathbf{w}, \mathbf{x}_i \rangle + b) = \rho - \xi_i$. So, $\lambda_i \neq 0$ for the support vectors. As $\rho > 0$, $\pi = 0 \Rightarrow \sum_{i=1}^N \lambda_i = v$ in terms of the KKT conditions.

For the outliers, $\xi_i \geq 0 \Rightarrow u_i = 0 \Rightarrow \lambda_i = 1$. $\lambda_i = 1$ for all the outliers. Let N_o denote the number of outliers. $N_o \leq N$

$$N_o = \sum_{i=1}^{N_o} 1 \leq \sum_{i=1}^N \lambda_i = v$$

Where λ_i denote value of λ of all data points.

The support vectors contain points on the hyperplanes and outliers. For the points on the hyperplanes, $\xi_i = 0 \Rightarrow u_i > 0 \Rightarrow \lambda_i < 1$ in terms of KKT conditions.

Hence, $\lambda_i \leq 1$ for the support vectors. $\lambda_i = 0$ for other points. Let N_s denote the number of support vectors.

$$N_s = \sum_{i=1}^{N_s} 1 \geq \sum_{i=1}^N \lambda_i = v$$

Where λ_i denote value of λ of all data points.

3.4 d

Answers:

As the training set $\{(x_i, y_i)\}_{i=1}^N$ is linearly separable, there exists w_0 and b_0 so that

$$y_i (\langle w, x_i \rangle + b) \geq 0 \quad \forall i$$

So $\rho \neq 0$ in that $\rho - \xi_i \leq 0$ if $\rho = 0$. The constraint will be loose so that we cannot get right and optimal hyperplane.

So $\rho > 0$. Also, v is an upper bound on the number of outliers bounded as proved above 3.3. If $v = \frac{1}{2}$, there exists no outliers. Hence, $\xi_i = 0$ for all data points.

The objective function will be,

$$\min_{w, b, \rho} \frac{1}{2} \|w\|_2^2 - v\rho \quad \text{s.t. } y_i (\langle w, x_i \rangle + b) \geq \rho$$

Rewriting that,

$$\min_{w, b, \rho} \frac{1}{2} \left\| \frac{w}{\rho} \right\|_2^2 - \frac{v}{\rho} \quad \text{s.t. } y_i \left(\left\langle \frac{w}{\rho}, x_i \right\rangle + \frac{b}{\rho} \right) \geq 1$$

Let $w' = \frac{w}{\rho}$, $v' = \frac{v}{\rho}$, $b' = \frac{b}{\rho}$.

$$\min_{w', b'} \frac{1}{2} \|w'\|_2^2 - v' \quad \text{s.t. } y_i (\langle w', x_i \rangle + b') \geq 1$$

We can remove v' , which does not influence minimization,

$$\min_{w', b'} \frac{1}{2} \|w'\|_2^2 \quad \text{s.t. } y_i (\langle w', x_i \rangle + b') \geq 1 \quad \forall i$$

Now the objective function is consistent with hard-margin SVM.

Consequently,

If $\rho > 0$, this problem yields the same binary classifier as that of hard-margin SVM.

3.5 e

Answers:

Assume $(\mathbf{w}_*, b_*, \xi_*, \rho_*)$ is a global minimizer of problem,

$$\min_{\mathbf{w}, b, \xi, \rho} \frac{1}{2} \|\mathbf{w}\|_2^2 - \nu \rho + \sum_{i=1}^N \xi_i \quad \text{s.t.} \quad y_i (\langle \mathbf{w}, \mathbf{x}_i \rangle + b) \geq \rho - \xi_i, \quad \xi_i \geq 0 \quad \forall i \quad \rho \geq 0$$

$$y_i (\langle \mathbf{w}_*, \mathbf{x}_i \rangle + b_*) \geq \rho_* - \xi_i^* \quad \forall i$$

So,

$$y_i \left(\left\langle \frac{\mathbf{w}}{\rho}, \mathbf{x}_i \right\rangle + \frac{b}{\rho} \right) \geq 1 - \frac{\xi_i}{\rho} \quad \forall i$$

$$\text{Let } \mathbf{w}' = \frac{\mathbf{w}}{\rho}, b' = \frac{b}{\rho}, \xi' = \frac{\xi_i}{\rho}$$

$$y_i (\langle \mathbf{w}', \mathbf{x}_i \rangle + b') \geq 1 - \xi' \quad \forall i$$

As $\rho_* > 0$, The objective function is,

$$\min_{\mathbf{w}, b, \xi, \rho} \frac{1}{2} \|\mathbf{w}\|_2^2 - \nu \rho + \sum_{i=1}^N \xi_i = \min_{\mathbf{w}, b, \xi, \rho} \frac{1}{2} \left\| \frac{\mathbf{w}}{\rho} \right\|_2^2 + \rho^{-1} \sum_{i=1}^N \frac{\xi_i}{\rho} - \frac{\nu}{\rho}$$

Let $v' = \nu/\rho$, the objective function is,

$$\min_{\mathbf{w}, b, \xi, \rho} \frac{1}{2} \|\mathbf{w}'\|_2^2 + \rho^{-1} \sum_{i=1}^N \xi' - v'$$

Assume ρ is constant. Removing v' , which does not influence minimization,

$$\min_{\mathbf{w}, b, \xi} \frac{1}{2} \|\mathbf{w}'\|_2^2 + \rho^{-1} \sum_{i=1}^N \xi' \quad y_i (\langle \mathbf{w}', \mathbf{x}_i \rangle + b') \geq 1 - \xi' \quad \forall i$$

This is soft-margin SVM with $C = \rho^{-1}$. As $(\mathbf{w}_*, b_*, \xi_*, \rho_*)$ is a global minimizer of original problem, the global minimizer $\{\mathbf{w}, b, \xi\}$ of soft-margin SVM with $C = \rho_*^{-1}$ are $\frac{\mathbf{w}_*}{\rho_*}, \frac{b_*}{\rho_*}$ and $\frac{\xi_*}{\rho_*}$ respectively.

4 Problem 4

4.1 a

The ℓ_2 norm case is,

$$\mathcal{H} \doteq \{\mathbf{x} \rightarrow \langle \mathbf{w}, \mathbf{x} \rangle : \|\mathbf{w}\|_2 \leq 1\}$$

$$\begin{aligned}
N * G(\mathbf{X}) &= \mathbb{E}_{\mathbf{g} \sim iid \mathcal{N}(0,1)} \sup_{\|\mathbf{w}\|_2 \leq 1} \langle \mathbf{X} \mathbf{w}, \mathbf{g} \rangle \\
&= \mathbb{E}_{\mathbf{g} \sim iid \mathcal{N}(0,1)} \sup_{\|\mathbf{w}\|_2 \leq 1} \sum_{i=1}^N g_i \langle \mathbf{w}, \mathbf{x}_i \rangle \\
&= \mathbb{E}_{\mathbf{g} \sim iid \mathcal{N}(0,1)} \sup_{\|\mathbf{w}\|_2 \leq 1} \left\langle \mathbf{w}, \sum_{i=1}^N g_i \mathbf{x}_i \right\rangle \\
&\leq \mathbb{E}_{\mathbf{g} \sim iid \mathcal{N}(0,1)} \left\| \sum_{i=1}^N g_i \mathbf{x}_i \right\|_2
\end{aligned}$$

using Jensen's inequality,

$$\begin{aligned}
\mathbb{E}_{\mathbf{g} \sim iid \mathcal{N}(0,1)} \left\| \sum_{i=1}^N g_i \mathbf{x}_i \right\|_2 &= \mathbb{E}_{\mathbf{g} \sim iid \mathcal{N}(0,1)} \left(\left\| \sum_{i=1}^N g_i \mathbf{x}_i \right\|_2^2 \right)^{1/2} \\
&\leq \left(\mathbb{E}_{\mathbf{g} \sim iid \mathcal{N}(0,1)} \left\| \sum_{i=1}^N g_i \mathbf{x}_i \right\|_2^2 \right)^{1/2}
\end{aligned}$$

g_1, g_2, \dots, g_n are sampled from normal Gaussian distribution and independent.

$$\begin{aligned}
\mathbb{E}_{\mathbf{g} \sim iid \mathcal{N}(0,1)} \left\| \sum_{i=1}^N g_i \mathbf{x}_i \right\|_2^2 &= \mathbb{E}_{\mathbf{g} \sim iid \mathcal{N}(0,1)} \sum_{i,j} g_i g_j \langle \mathbf{x}_i, \mathbf{x}_j \rangle \\
&= \sum_{i \neq j} \langle \mathbf{x}_i, \mathbf{x}_j \rangle \mathbb{E}_{\mathbf{g} \sim iid \mathcal{N}(0,1)} g_i g_j + \sum_{i=1}^N \langle \mathbf{x}_i, \mathbf{x}_i \rangle \mathbb{E}_{\mathbf{g} \sim iid \mathcal{N}(0,1)} g_i^2 \\
&= \sum_{i=1}^N \|\mathbf{x}_i\|_2^2 \leq N \max_i \|\mathbf{x}_i\|_2^2
\end{aligned}$$

Hence,

$$G(\mathbf{X}) \leq \frac{\max_i \|\mathbf{x}_i\|_2}{\sqrt{N}}$$

4.2 b

The ℓ_∞ norm case is,

$$\mathcal{H} \doteq \{\mathbf{x} \rightarrow \langle \mathbf{w}, \mathbf{x} \rangle : \|\mathbf{w}\|_\infty \leq 1\}$$

using Holder's inequality,

$$\begin{aligned}
N * G(\mathbf{X}) &= \mathbb{E}_{\mathbf{g} \sim iid \mathcal{N}(0,1)} \sup_{\|\mathbf{w}\|_\infty \leq 1} \langle \mathbf{X} \mathbf{w}, \mathbf{g} \rangle \\
&= \mathbb{E}_{\mathbf{g} \sim iid \mathcal{N}(0,1)} \sup_{\|\mathbf{w}\|_\infty \leq 1} \sum_{i=1}^N g_i \langle \mathbf{w}, \mathbf{x}_i \rangle \\
&= \mathbb{E}_{\mathbf{g} \sim iid \mathcal{N}(0,1)} \sup_{\|\mathbf{w}\|_2 \leq 1} \left\langle \mathbf{w}, \sum_{i=1}^N g_i \mathbf{x}_i \right\rangle \\
&\leq \mathbb{E}_{\mathbf{g} \sim iid \mathcal{N}(0,1)} \left\| \sum_{i=1}^N g_i \mathbf{x}_i \right\|_1
\end{aligned}$$

For each $j \in [d]$, let $v_j = (x_1, j \dots x_N, j) \in \mathbb{R}^N$. Note that $\|v_j\|_2 \leq \sqrt{N} \max_i \|\mathbf{x}_i\|_1$. Let $V = \{v_1, \dots, -v_1, \dots, -v_n\}$. The right-hand side is $N G(V)$. Using Massart lemma,

$$G(V) \leq \max_i \|\mathbf{x}_i\|_1 \sqrt{2 \log(2d) / N}$$