

Example 2: Least Norm Solution

Least norm solution:

$$\min_{x \in \mathbb{R}^n} \frac{1}{2} x^T x$$

$A \in \mathbb{R}^{m \times n}$,
 $n \geq m$
 A is full-rank.

s.t. $Ax = b$: $\mu \in \mathbb{R}^m$

Find L and d .

$$L(x, \mu) = \frac{1}{2} x^T x + \mu^T (Ax - b)$$

$$\underline{d(\mu)} = \inf_{x \in \mathbb{R}^n} L(x, \mu) = \inf_{x \in \mathbb{R}^n} \left[\frac{1}{2} x^T x + \mu^T Ax - \mu^T b \right]$$

necessary condition

$$\begin{cases} \nabla_x L(x, \mu) = x + A^T \mu = 0 \\ \nabla_x^2 L(x, \mu) = I > 0 \end{cases}$$

↓

$$\bar{x} = -A^T \mu$$

minimizes $L(x, \mu)$.

then, $d(\mu) = \frac{1}{2} (-A^T \mu)^T (-A^T \mu) + \mu^T A (-A^T \mu) - \mu^T b$

$$= \frac{1}{2} \mu^T A A^T \mu - \mu^T A A^T \mu - \mu^T b$$

Dual optimization problem: $\max_{\mu \in \mathbb{R}^m} \left[-\frac{1}{2} \mu^T A A^T \mu - \mu^T b \right]$

$$\cong \min_{\mu \in \mathbb{R}^m} \frac{1}{2} \mu^T A A^T \mu + \mu^T b.$$

optimal dual solution μ^* satisfies $\nabla_{\mu} d(\mu) = 0$

$$\nabla_{\mu}^2 d(\mu) = -A A^T \text{ and } A A^T \text{ is positive definite.}$$

$$\Rightarrow A A^T \mu^* + b = 0$$

$$\Rightarrow \underline{\mu^* = -(A A^T)^{-1} b.}$$

then $x^* = -A^T (-(A A^T)^{-1} b)$

$$\underline{x^* = A^T (A A^T)^{-1} b.}$$

Towards Optimality Conditions

Consider the primal optimization problem:

$$\begin{aligned} \min_{x \in \mathbb{R}^n} \quad & f(x) \\ \text{s.t.} \quad & g_i(x) \leq 0, i \in [m] := \{1, 2, \dots, m\}, \\ & h_j(x) = 0, j \in [p]. \end{aligned}$$

Let the dual function be defined as

$$d(\lambda, \mu) := \inf_x L(x, \lambda, \mu) = \inf_x \left[f(x) + \sum_{i \in [m]} \lambda_i g_i(x) + \sum_{j \in [p]} \mu_j h_j(x) \right].$$

The corresponding dual optimization problem is:

$$\begin{aligned} \max_{\lambda \in \mathbb{R}^m, \mu \in \mathbb{R}^p} \quad & d(\lambda, \mu) \\ \text{s.t.} \quad & \lambda \geq 0, \\ & (\lambda, \mu) \in \text{dom}(d). \end{aligned}$$

Consequently, for any $(\bar{x}, \bar{\lambda}, \bar{\mu})$ with \bar{x} being primal feasible and $\bar{\lambda} \geq 0$, we have

$$\begin{aligned} d(\bar{\lambda}, \bar{\mu}) &= \inf_x \left[f(x) + \sum_{i \in [m]} \bar{\lambda}_i g_i(x) + \sum_{j \in [p]} \bar{\mu}_j h_j(x) \right] \\ &= \leq \left[f(\bar{x}) + \sum_{i \in [m]} \bar{\lambda}_i g_i(\bar{x}) + \sum_{j \in [p]} \bar{\mu}_j h_j(\bar{x}) \right] \\ &= \leq f(\bar{x}). \quad \underbrace{\sum_{i \in [m]} \bar{\lambda}_i g_i(\bar{x})}_{\leq 0} \quad \underbrace{\sum_{j \in [p]} \bar{\mu}_j h_j(\bar{x})}_{=0} = f(\bar{x}) \end{aligned}$$

When do both the above inequalities hold with equality?

- the last inequality holds with equality when

$$\sum_{i \in [m]} \bar{\lambda}_i g_i(\bar{x}) = 0 \Leftrightarrow \bar{\lambda}_i g_i(\bar{x}) = 0 \quad \forall i=1, 2, \dots, m$$

- the first inequality holds with equality when

- \bar{x} is the minimizer of $L(x, \bar{\lambda}, \bar{\mu})$.

- necessary condition $\nabla_x L(\bar{x}, \bar{\lambda}, \bar{\mu}) = 0$

- sufficient condition $\nabla_x^2 L(\bar{x}, \bar{\lambda}, \bar{\mu}) \succ 0$.

Necessary Conditions for Optimality

Theorem 3

Suppose strong duality holds. Let x^* be the optimal solution of the primal problem, and λ^*, μ^* be the optimal solution of the dual problem. Then, the following conditions are satisfied.

- **Primal Feasibility:** $g_i(x^*) \leq 0, i \in [m], h_j(x^*) = 0, j \in [p]$.
- **Dual Feasibility:** $\lambda^* \geq 0$.
- **Complementary Slackness:** $\lambda_i^* g_i(x^*) = 0$ for all $i \in [m]$.
- **Lagrangian Stationarity:**

$$\nabla_x L(x^*, \lambda^*, \mu^*) = \nabla_x f(x^*) + \sum_{i \in [m]} \lambda_i^* \nabla_x g_i(x^*) + \sum_{j \in [p]} \mu_j^* \nabla_x h_j(x^*) = 0.$$

The above four conditions are called Karush–Kuhn–Tucker (KKT) optimality conditions.

When are the above conditions sufficient for optimality?

Essentially, the Lagrangian stationarity conditions must imply that x^* is the optimal solution of the Lagrangian problem. When it is true?

Suppose $(\hat{x}, \hat{\lambda}, \hat{\mu})$ satisfying the KKT conditions. Then,

$$d(\hat{\lambda}, \hat{\mu}) = \inf_x L(x, \hat{\lambda}, \hat{\mu})$$

$$\leq L(\hat{x}, \hat{\lambda}, \hat{\mu}) = f(\hat{x}) + \sum_i \hat{\lambda}_i \underbrace{g_i(\hat{x})}_{=0} + \sum_j \mu_j \underbrace{h_j(\hat{x})}_{=0} = f(\hat{x})$$

Sufficient Condition for Optimality for Convex Problems

Theorem 4

Suppose the primal optimization problem is convex. Let \bar{x} , $\bar{\lambda}$ and $\bar{\mu}$ satisfy KKT conditions stated above. Then,

- $d(\bar{\lambda}, \bar{\mu}) = f(\bar{x})$ (strong duality holds).
- \bar{x} is optimal solution of primal problem.
- $(\bar{\lambda}, \bar{\mu})$ are optimal solution of dual problem.

It is still possible for a convex optimization problem to have an optimal solution but no KKT points. We need additional conditions to make sure that optimal solutions satisfy KKT conditions.

Let \bar{x} , $\bar{\lambda}$ and $\bar{\mu}$ satisfy KKT conditions stated above. From primal and dual feasibility we have

$$\begin{aligned} d(\bar{\lambda}, \bar{\mu}) &= \inf_x \left[f(x) + \sum_{i \in [m]} \bar{\lambda}_i g_i(x) + \sum_{j \in [p]} \bar{\mu}_j h_j(x) \right] \\ &\leq f(\bar{x}) + \sum_{i \in [m]} \bar{\lambda}_i g_i(\bar{x}) + \sum_{j \in [p]} \bar{\mu}_j h_j(\bar{x}) \leq f(\bar{x}). \end{aligned}$$

Further, both inequalities hold with equality.

Constraint Qualification and Strong Duality

Theorem 5

Suppose the primal optimization problem is convex which satisfies Slater's constraint qualification condition: there exists $\bar{x} \in \text{int}(\mathcal{D})$ in the domain of the optimization problem for which $g_i(\bar{x}) < 0$ for all $i \in [m]$ and $h_i(\bar{x}) = 0$ for all $i \in [p]$. Then, strong duality holds with $f^* = d^*$. Moreover, if $f^* > -\infty$, then, there exist (λ^*, μ^*) such that $g(\lambda^*, \mu^*) = d^* = f^*$.

- The point \bar{x} need not be an optimal solution. It is any arbitrary feasible point.
- **Relaxed Slater Condition:** If some of the inequality constraints are affine, then they need not hold with strict inequality. It is sufficient to find $\bar{x} \in \text{relint}(\mathcal{D})$ such that $g_i(\bar{x}) < 0$ for all g_i that are not affine.
- We now have the following result.

Proposition 1. *Suppose the primal problem is convex and satisfies Slater's condition. Then, if a feasible solution x^* is optimal, then there exist λ^*, μ^* such that (x^*, λ^*, μ^*) satisfy KKT conditions.*

- Note that sufficiency part holds even without Slater's condition.
- An alternative condition is **Linear Independence Constraint Qualification (LICQ)** which holds at a feasible solution x^* if the vectors

$$\begin{aligned} & \nabla h_j(x^*), \quad j \in [p], \\ & \nabla g_i(x^*), \quad i \in \{k \in [m] \mid g_k(x^*) = 0\} \end{aligned}$$

are linearly independent. If LICQ holds at a point x^* , then KKT conditions are necessary for the local optimality of x^* .

Convex Theorem of the Alternative

Consider the following general form of optimization problem:

$$\begin{aligned} \min_{x \in \mathbb{R}^n} \quad & f(x) \\ \text{s.t.} \quad & g_i(x) \leq 0, i \in [m] := \{1, 2, \dots, m\}, \end{aligned}$$

where f and g_i are convex functions.

Theorem 6

Let the constraint functions g_i satisfy Slater's condition: there exists \bar{x} such that $g_i(\bar{x}) < 0$ for all $i \in [m]$. Then, exactly one of the following two sets must be empty.

- $S = \{x \in \mathbb{R}^n \mid f(x) < 0, g_i(x) \leq 0, i \in [m]\}$
- $T = \{\lambda \in \mathbb{R}^m \mid \lambda \geq 0, \inf_{x \in \mathbb{R}^n} [f(x) + \sum_{i \in [m]} \lambda_i g_i(x)] \geq 0\}$.

Case 1: If T is non-empty, then S is empty. Suppose S is not empty.

$$\Rightarrow \bar{x} \in S \Leftrightarrow f(\bar{x}) < 0, g_i(\bar{x}) \leq 0 \quad \forall i.$$

Case 2: If S is empty, then T is non-empty. Can be shown via separating hyperplane theorem, but bit more involved.

Skipped

$$f(\bar{x}) + \sum_i \lambda_i g_i(\bar{x}) < 0 \quad \text{when } \lambda_i \geq 0$$

which means

$$\inf_{x \in \mathbb{R}^n} [f(x) + \sum_i \lambda_i g_i(x)] < 0 \quad \forall \lambda \geq 0$$

$\Rightarrow T$ is empty, contradicting initial hypothesis.

Strong Duality Theorem

weak duality:
 $d^* \leq f^*$

Consider the following general form of optimization problem:

$$\begin{aligned} \min_{x \in \mathbb{R}^n} \quad & f(x) \\ \text{s.t.} \quad & g_i(x) \leq 0, i \in [m] := \{1, 2, \dots, m\}, \end{aligned}$$

where f and g_i are convex functions satisfying Slater's condition.

Theorem 7

x^* is an optimal solution to the above problem if and only if there exists $\lambda^* \geq 0$ such that $\inf_{x \in \mathbb{R}^n} [f(x) + \sum_{i \in [m]} \lambda_i^* g_i(x)] \geq f(x^*)$. $\Rightarrow d^* = f^*$.

Since x^* is an optimal solution, the set

$$S = \{x \in \mathbb{R}^n \mid f(x) - f(x^*) < 0, g_i(x) \leq 0, i \in [m]\}$$

is infeasible.

It follows from the above theorem that the set

$$T = \{\lambda \in \mathbb{R}^m \mid \lambda \geq 0, \inf_{x \in \mathbb{R}^n} [f(x) - f(x^*) + \sum_{i \in [m]} \lambda_i g_i(x)] \geq 0\}$$

is feasible.

Lecture 17: 7th Feb 2025

Binary Classification

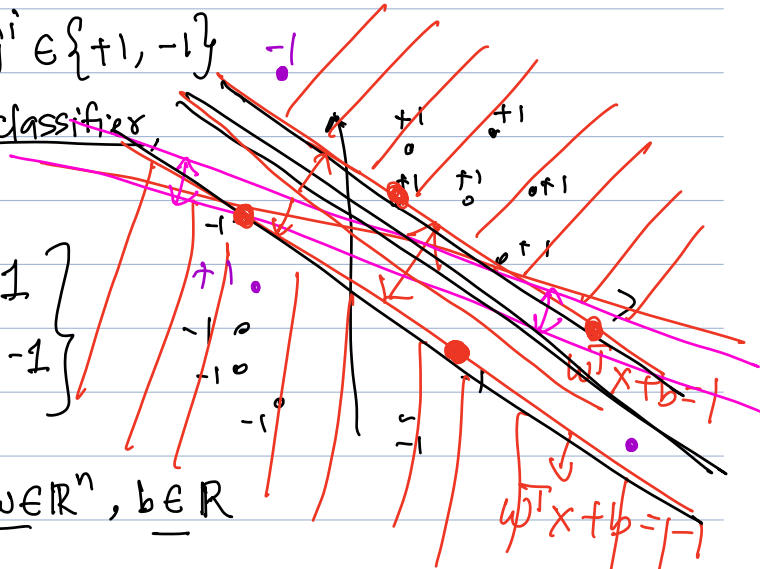
we are given $(\hat{x}^1, \hat{y}^1), (\hat{x}^2, \hat{y}^2), \dots, (\hat{x}^N, \hat{y}^N)$,

where $\hat{x}^i \in \mathbb{R}^n$, $\hat{y}^i \in \{+1, -1\}$

the goal is to find a linear classifier,

$w^T x + b$ such that

$$\begin{cases} w^T \hat{x}^i + b > 0 & \text{whenever } \hat{y}^i = +1 \\ w^T \hat{x}^i + b < 0 & \text{whenever } \hat{y}^i = -1 \end{cases}$$



decision variables: (w, b) , $w \in \mathbb{R}^n$, $b \in \mathbb{R}$

constraints:

$$\hat{y}^i (w^T \hat{x}^i + b) > 0 \quad \text{for all } i = 1, 2, \dots, N.$$

$$\Leftrightarrow \hat{y}^i (w^T \hat{x}^i + b) \geq \beta \quad \text{where } \beta > 0, \forall i$$

$$\Leftrightarrow \hat{y}^i \left(\left(\frac{w}{\beta} \right)^T \hat{x}^i + \frac{b}{\beta} \right) \geq 1 \quad \forall i$$

$$\Leftrightarrow \hat{y}^i (w^T \hat{x}^i + b) \geq 1 \quad \forall i = 1, 2, \dots, N$$

without loss of generality.

cost function:

$$w^T \hat{x}^i + b \geq 1 \quad \text{for all } i \text{ with } \hat{y}^i = +1$$

$$w^T \hat{x}^i + b \leq -1 \quad \text{for all } i \text{ with } \hat{y}^i = -1$$

For maximum margin classifiers, we try to maximize

the distance between the two hyperplanes

$$H_1 = \{x \in \mathbb{R}^n \mid w^T x + b = 1\}$$

$$H_2 = \{x \in \mathbb{R}^n \mid w^T x + b = -1\}$$

$$\text{dist}(H_1, H_2) = \frac{2}{\|w\|_2} \quad (\text{try to prove yourself}).$$

we choose the cost function $\frac{1}{2} \|w\|_2^2$ in order to maximize $\text{dist}(H_1, H_2)$.

optimization problem:

$$\min_{w \in \mathbb{R}^n, b \in \mathbb{R}^n} \frac{1}{2} w^T w$$

$$\text{s.t.} \quad 1 - \hat{y}^i (w^T \hat{x}^i + b) \leq 0, \quad i=1, 2, \dots, N.$$

$$L(w, b, \lambda) = \frac{1}{2} w^T w + \sum_{i=1}^N \lambda_i [1 - \hat{y}^i (w^T \hat{x}^i + b)]$$

$\lambda_i, \lambda \in \mathbb{R}^N$

$$= \left[\frac{1}{2} w^T w - \sum_{i=1}^N \lambda_i \hat{y}^i w^T \hat{x}^i \right] - \sum_{i=1}^N \lambda_i \hat{y}^i b + \sum_{i=1}^N \lambda_i$$

Let us find (\bar{w}, \bar{b}) that minimize $L(w, b, \lambda)$.

$$\nabla_w L(w, b, \lambda) = w - \sum_{i=1}^N \lambda_i \hat{y}^i \hat{x}^i$$

$$\nabla_w^2 L(w, b, \lambda) = I \geq 0$$

$$\Rightarrow \bar{w} = \sum_{i=1}^N \lambda_i \hat{y}^i \hat{x}^i \text{ minimizes } L(w, b, \lambda).$$

dual function: $d(\lambda) = -\infty$ if $\sum_{i=1}^N \lambda_i \hat{y}^i \neq 0$

and $d(\lambda) = \frac{1}{2} (\bar{w})^T \bar{w} - \sum_{i=1}^N \lambda_i \hat{y}^i (\bar{w})^T \hat{x}^i + \sum_{i=1}^N \lambda_i$ when $\sum_{i=1}^N \lambda_i \hat{y}^i = 0$

$$\Rightarrow d(\lambda) = \frac{1}{2} \left(\sum_{i=1}^N \lambda_i \hat{y}^i \hat{x}^i \right)^T \left(\sum_{j=1}^N \lambda_j \hat{y}^j \hat{x}^j \right) - \sum_{i=1}^N \lambda_i \hat{y}^i (\hat{x}^i)^T \left(\sum_{j=1}^N \lambda_j \hat{y}^j \hat{x}^j \right) + \sum_{i=1}^N \lambda_i$$

$$= \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \lambda_i \lambda_j \hat{y}^i \hat{y}^j (\hat{x}^i)^T \hat{x}^j$$

$$- \sum_{i=1}^N \sum_{j=1}^N \lambda_i \lambda_j \hat{y}^i \hat{y}^j (\hat{x}^i)^T (\hat{x}^j) + \sum_{i=1}^N \lambda_i$$

$$= -\frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \lambda_i \lambda_j (\hat{y}^i)^T (\hat{y}^j) (\hat{x}^i)^T (\hat{x}^j) + \sum_{i=1}^N \lambda_i$$

Dual optimization:

$$\begin{array}{ll} \max & d(\lambda) \\ \lambda \in \mathbb{R}^N & \\ \text{s.t.} & \lambda \geq 0, \sum_{i=1}^N \lambda_i \hat{y}^i = 0 \end{array}$$

Recovering optimal solution of the primal from optimal dual solution

Let λ^* be the optimal solution of the dual.

Let (w^*, b^*) be the optimal solution of the primal.

KKT conditions are given by:

$$\left[\begin{array}{l} \text{i)} \lambda^* \geq 0, \sum_{i=1}^N \lambda_i^* \hat{y}^i = 0 \\ \text{ii)} 1 - \hat{y}^i ((w^*)^T \hat{x}^i + b^*) \leq 0, \forall i = 1, 2, \dots, N \\ \text{iii)} \lambda_i [1 - \hat{y}^i ((w^*)^T \hat{x}^i + b^*)] = 0, \forall i = 1, 2, \dots, N \\ \text{iv)} \nabla_w L(w^*, b^*, \lambda^*) = 0 \Leftrightarrow w^* = \sum_{i=1}^N \lambda_i^* \hat{y}^i \hat{x}^i \end{array} \right.$$

Finding w^* is straightforward.

to find b^* , first find an index z s.t. $\lambda_z^* \neq 0$

then $1 - \hat{y}^z [(w^*)^T \hat{x}^z + b^*] = 0$

$$\Rightarrow (w^*)^T \hat{x}^z + b^* = \frac{1}{\hat{y}^z}$$

$$\Rightarrow b^* = \frac{1}{\hat{y}^z} - (w^*)^T \hat{x}^z$$

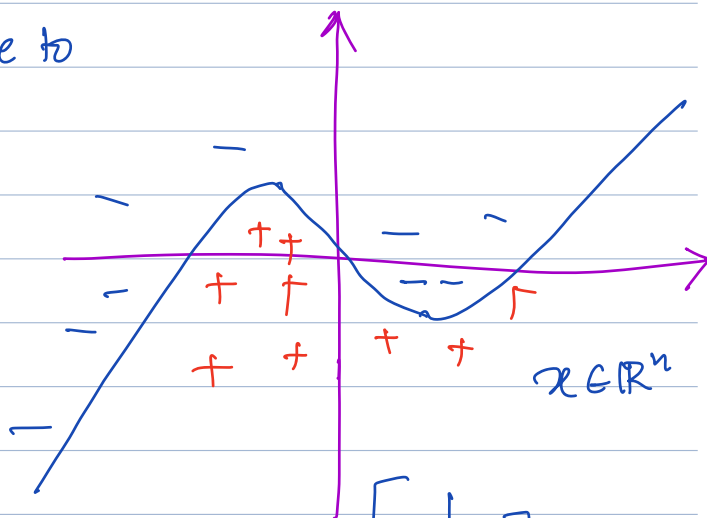
If the points are amenable to a linear classifier, we can try to define a feature vector

$$\phi: \mathbb{R}^n \rightarrow \mathbb{R}^k,$$

$k \gg n$. and try to find

(w, b) s.t.

$w^T \phi(x) + b$ classifies the points. $w \in \mathbb{R}^k$



$$\phi(x) = \begin{bmatrix} 1 \\ x_1 \\ x_2 \\ x_3 \\ x_4 \end{bmatrix}, \quad k=5$$

$w \in \mathbb{R}^5$

$$\text{If } x = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}, \quad \phi(x) = \begin{bmatrix} 1 \\ x_1 \\ x_2 \\ x_1^2 \\ x_2^2 \\ x_1 x_2 \end{bmatrix}$$

We can now formulate the optimization problem

Lecture -18
8th Feb.

$$(P) \begin{cases} \min_{\substack{\mathbf{w} \in \mathbb{R}^k, b \in \mathbb{R} \\ \boldsymbol{\varepsilon} \in \mathbb{R}^N}} & \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \sum_{i=1}^N \varepsilon_i \\ \text{s.t.} & 1 - \hat{y}^i (\mathbf{w}^T \phi(\hat{x}^i) + b) \leq \varepsilon_i; \quad i=1, 2, \dots, N \\ & \varepsilon_i \geq 0; \quad i=1, 2, \dots, N \end{cases}$$

$\lambda_{1,i}, \lambda_{1,i} \in \mathbb{R}^N$
 $\lambda_{2,i}, \lambda_{2,i} \in \mathbb{R}^N$

C : large constant that penalizes non-zero ε values.

Homework: [derive dual of (P).
from optimal dual solution, recover the optimal primal solution.

$$L(\mathbf{w}, b, \boldsymbol{\varepsilon}, \lambda_1, \lambda_2) = \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \sum_{i=1}^N \varepsilon_i + \sum_{i=1}^N \lambda_{1,i} (1 - \hat{y}^i (\mathbf{w}^T \phi(\hat{x}^i) + b) - \varepsilon_i) + \sum_{i=1}^N \lambda_{2,i} (-\varepsilon_i)$$

$$= \left[\frac{1}{2} \mathbf{w}^T \mathbf{w} - \sum_{i=1}^N \lambda_{1,i} \hat{y}^i \mathbf{w}^T \phi(\hat{x}^i) \right] - \sum_{i=1}^N \lambda_{1,i} \hat{y}^i b + \sum_{i=1}^N \lambda_{1,i} + \sum_{i=1}^N \varepsilon_i (C - \lambda_{1,i} - \lambda_{2,i})$$

$$\nabla_{\mathbf{w}} L = \mathbf{w} - \sum_{i=1}^N \lambda_{1,i} \hat{y}^i \phi(\hat{x}^i) = \mathbf{0} \Rightarrow \bar{\mathbf{w}} = \sum_{i=1}^N \lambda_{1,i} \hat{y}^i \phi(\hat{x}^i)$$

to obtain a finite value for $\inf_{\omega, b, \varepsilon} L(\omega, b, \varepsilon, \lambda_1, \lambda_2)$,
we need to set

$$c - \lambda_{1,i} - \lambda_{2,i} = 0 \quad \forall i = 1, 2, \dots, N$$

$$\sum_{i=1}^N \lambda_{1,i} \hat{y}^i = 0$$

The dual optimization problem is given by:

$$\max_{\lambda_1, \lambda_2} d(\lambda_1, \lambda_2)$$

$$\text{s.t. } \lambda_1 \geq 0, \lambda_2 \geq 0$$

$$c - \lambda_{1,i} = \lambda_{2,i} \geq 0$$

$$c - \lambda_{1,i} - \lambda_{2,i} = 0 \quad \forall i = 1, 2, \dots, N$$

$$\sum_{i=1}^N \lambda_{1,i} \hat{y}^i = 0$$

$$d(\lambda_1, \lambda_2) = \frac{1}{2} (\bar{\omega})^T \bar{\omega} - \sum_{i=1}^N \lambda_{1,i} \hat{y}^i (\bar{\omega})^T \phi(\hat{x}^i) + \sum_{i=1}^N \lambda_{2,i}$$

$$= -\frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \lambda_{1,i} \lambda_{1,j} \hat{y}^i \hat{y}^j \phi(\hat{x}^i)^T \phi(\hat{x}^j) + \sum_{i=1}^N \lambda_{1,i}$$

Equivalently:

(dual)

$$\begin{cases} \max_{\lambda_1 \in \mathbb{R}^N} d(\lambda_1) \\ \text{s.t. } \lambda_1 \geq 0, c - \lambda_{1,i} \geq 0 \quad \forall i, \\ \sum_{i=1}^N \lambda_{1,i} \hat{y}^i = 0 \end{cases}$$

Suppose the optimal dual solution λ_1^* is given. Let us try to recover the optimal primal solution $(\omega^*, b^*, \varepsilon^*)$.

$$\omega^* = \sum_{i=1}^N \lambda_{1,i}^* \hat{y}^i \phi(\hat{x}^i)$$

Let us state complementarity slackness conditions.

$$\lambda_{1,i}^+ (1 - \hat{y}^i ((\omega^*)^T \phi(\hat{x}^i) + b^*) - \varepsilon_i^+) = 0 \quad (CS)$$

$$\varepsilon_i^+ (\cancel{\lambda_{1,i}^+}) (C - \lambda_{1,i}^+) = 0$$

Simply finding an index z at which $\lambda_{1,z}^+ > 0$ is not enough since if $\lambda_{1,z}^+ = C$, we can't determine ε_i^+ .

Hence, we choose index z s.t. $0 < \lambda_{1,z}^+ < C$, which implies

$$\varepsilon_z^+ = 0 \text{ and } 1 - \hat{y}^z ((\omega^*)^T \phi(\hat{x}^z) + b^*) - \varepsilon_z^+ = 0$$

$$\Rightarrow b^* = \frac{1}{\hat{y}^z} - (\omega^*)^T \phi(\hat{x}^z).$$

After finding (ω^*, b^*) , we can find ε_i^+ from (CS) conditions.

Suppose we are given a new data point x^{new} .

We can predict its label by evaluating

$$(\omega^*)^T \phi(x^{\text{new}}) + b^*.$$

if $\text{sign} \left[\begin{array}{c} \downarrow \\ \phantom{(\omega^*)^T \phi(x^{\text{new}}) + b^*} \end{array} \right] > 0$, then assign label $+1$
 < 0 , " " " -1 .

observe that:

$$b^* + (\omega^*)^T \phi(x^{\text{new}}) = b^* + \left(\sum_{i=1}^N \lambda_{1,i}^+ \hat{y}^i \phi(\hat{x}^i) \right)^T \phi(x^{\text{new}})$$

$$= \left[\frac{1}{\hat{y}^z} - \sum_{i=1}^N \lambda_{1,i}^+ \hat{y}^i \underbrace{\phi(\hat{x}^i)^T \phi(\hat{x}^z)}_{K(\hat{x}^i, \hat{x}^z)} \right] + \sum_{i=1}^N \lambda_{1,i}^+ \hat{y}^i \underbrace{\phi(\hat{x}^i)^T \phi(x^{\text{new}})}_{K(\hat{x}^i, x^{\text{new}})}$$

- Note that the above expression does not depend on primal decision variables.

- In addition, the dependence on ϕ is in terms of $\phi(x)^T \phi(y)$ for various x, y .

- In practice, it could be challenging to try out different choice of ϕ , and a larger value of K ($\dim \phi(x)$), increases the complexity of the primal optimization problem.

- To avoid the above issues, often kernel methods are deployed.

$$K(x_1, x_2), \quad \text{e.g.: } K(x_1, x_2) = \exp\left(\frac{\|x_1 - x_2\|^2}{2M}\right),$$

- we replace $\phi(x_i)^T \phi(x_j)$ terms by $K(x_i, x_j)$.

- The dual optimization problem can be written as

$$\left[\begin{array}{l} \max_{\lambda_i \in \mathbb{R}^N} \quad -\frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \lambda_{i,i} \lambda_{j,j} \hat{y}^i \hat{y}^j K(\hat{x}_i, \hat{x}_j) + \sum_{i=1}^N \lambda_{i,i} \\ \text{s.t.} \quad \lambda_{i,i} \geq 0, \quad \lambda_{i,i} \leq C \quad \forall i \\ \sum_{i=1}^N \lambda_{i,i} \hat{y}_i = 0 \end{array} \right.$$

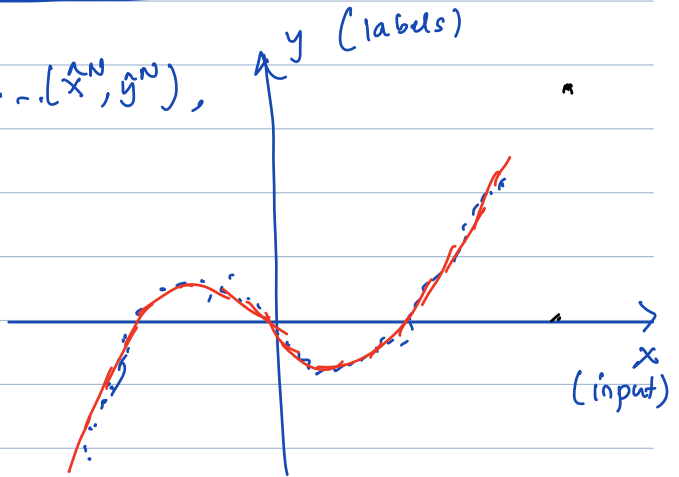
Regression Problems

As before, we are given $(\hat{x}^1, \hat{y}^1) \dots (\hat{x}^N, \hat{y}^N)$,
 $\hat{y}^i \in \mathbb{R}$, $\hat{x}^i \in \mathbb{R}^n$.

we wish to find $w \in \mathbb{R}^k$ s.t.

$$\underline{w^T \phi(\hat{x}^i) \approx \hat{y}^i} \quad \forall i=1, 2, \dots, N.$$

$\phi: \mathbb{R}^n \rightarrow \mathbb{R}^k$ feature map.



For a fixed w , the residual error for i th data point: $\underline{w^T \phi(\hat{x}^i) - \hat{y}^i}$

residual error vector:
$$\begin{bmatrix} w^T \phi(\hat{x}^1) - \hat{y}^1 \\ \vdots \\ w^T \phi(\hat{x}^N) - \hat{y}^N \end{bmatrix} \in \mathbb{R}^N = \underline{\underline{\phi(\hat{x})w - \hat{y}}}$$

where $\hat{y} \in \mathbb{R}^N$, $(\hat{y})_i = \hat{y}^i$,

$\underline{\phi(\hat{x})} \in \mathbb{R}^{N \times k}$, $\phi(\hat{x}^i)$ is the i -th row of $\phi(\hat{x})$ matrix.

Regression problem: $\min_{w \in \mathbb{R}^k} \|\phi(\hat{x})w - \hat{y}\|_{2/1/\infty}$

Least squares problem: $\min_{w \in \mathbb{R}^k} \|\phi(\hat{x})w - \hat{y}\|_2^2$

the above problem is convex
 \uparrow

$$\begin{aligned} &= (\phi(\hat{x})w - \hat{y})^T (\phi(\hat{x})w - \hat{y}) \\ &= w^T \phi(\hat{x})^T \phi(\hat{x}) w - 2\hat{y}^T \phi(\hat{x}) w + \hat{y}^T \hat{y} \end{aligned}$$

Hessian: $\underline{\phi(\hat{x})^T \phi(\hat{x})} \succeq 0$

If there are constraints on $w \in W$, then the problem remains a convex optimization problem if W is a convex set.

Now, let us tackle the Regression problem with respect to the ∞ -norm. Recall $\|x\|_\infty = \max_{1 \leq i \leq n} |x_i|$

$$\min_{w \in \mathbb{R}^k} \left[\|\phi(\hat{x})w - \hat{y}\|_\infty \right] = \min_{w \in \mathbb{R}^k} \left[\max_{1 \leq i \leq N} |\phi(\hat{x}^i)^T w - \hat{y}^i| \right]$$

$$= \min_{\substack{w \in \mathbb{R}^k \\ t \in \mathbb{R}}} t \\ \text{s.t.} \quad \max_{1 \leq i \leq N} |\phi(\hat{x}^i)^T w - \hat{y}^i| = t$$

equivalently:

$$\min_{\substack{w \in \mathbb{R}^k \\ t \in \mathbb{R}}} t \\ \text{s.t.} \quad |\phi(\hat{x}^i)^T w - \hat{y}^i| \leq t, \quad \forall i = 1, 2, \dots, N$$

$$= \min_{\substack{w \in \mathbb{R}^k \\ t \in \mathbb{R}}} t \\ \text{s.t.} \quad \begin{aligned} \phi(\hat{x}^i)^T w - \hat{y}^i &\leq t & \forall i = 1, 2, \dots, N \\ -\phi(\hat{x}^i)^T w + \hat{y}^i &\leq t & \forall i = 1, 2, \dots, N. \end{aligned}$$

is a Linear programming problem.

$$\|x\|_1 = \sum_{i=1}^N |x_i|$$

Regression problem w.r.t 1-norm

$$\min_{w \in \mathbb{R}^k} \|\phi(\hat{x})w - \hat{y}\|_1 = \min_{w \in \mathbb{R}^k} \sum_{i=1}^N \underbrace{|\phi(\hat{x}^i)^T w - \hat{y}^i|}_{t_i}$$

$$= \min_{\substack{w \in \mathbb{R}^k \\ t \in \mathbb{R}^N}} \sum_{i=1}^N t_i$$

$$\text{s.t.} \quad |\phi(\hat{x}^i)^T w - \hat{y}^i| = t_i \quad \forall i = 1, 2, \dots, N$$

equivalently:

$$\begin{aligned} & \min_{\substack{w \in \mathbb{R}^k \\ t \in \mathbb{R}^N}} \sum_{i=1}^N t_i \\ & \text{s.t.} \quad \phi(\hat{x}^i)^T w - \hat{y}^i \leq t_i \\ & \quad \quad -\phi(\hat{x}^i)^T w + \hat{y}^i \leq t_i \\ & \quad \quad t_i \geq 0 \end{aligned}$$

which is a Linear programming problem.

Class Test on Wednesday during class hours.

↳ Syllabus: all topics covered so far.
