

Self grades are due at 11 PM on April 26, 2024.

1. Median Versus Mean

For a given vector $\vec{v} \in \mathbb{R}^n$, the mean can be found as the solution to the optimization problem

$$\min_{x \in \mathbb{R}} \|\vec{v} - x\vec{1}\|_2^2, \tag{1}$$

where $\vec{1}$ is the vector of ones in \mathbb{R}^n . Similarly, the median (any value x such that there is an equal number of values in \vec{v} above or below x) can be found via

$$\min_{x \in \mathbb{R}} \|\vec{v} - x\vec{1}\|_1. \tag{2}$$

We consider a robust version of the mean problem (1):

$$\min_x \max_{\vec{u} : \|\vec{u}\|_\infty \leq \lambda} \|\vec{v} + \vec{u} - x\vec{1}\|_2^2, \tag{3}$$

in which we assume that the components of \vec{v} can be independently perturbed by a vector \vec{u} whose magnitude is bounded by a given number $\lambda \geq 0$.

- (a) Is the robust problem (3) convex? Justify your answer precisely, based on expression (3), and without further manipulation.

Solution: The robust problem is convex, since the objective function is the pointwise maximum (over \vec{u}) of convex functions, $x \rightarrow \|\vec{v} + \vec{u} - x\vec{1}\|_2^2$.

- (b) Show that problem (3) can be expressed as

$$\min_{x \in \mathbb{R}} \sum_{i=1}^n (|v_i - x| + \lambda)^2. \tag{4}$$

Solution: For any vector $\vec{z} \in \mathbb{R}^n$, we have

$$\max_{\vec{u} : \|\vec{u}\|_\infty \leq \lambda} \|\vec{z} + \vec{u}\|_2^2 = \max_{\vec{u} : \|\vec{u}\|_\infty \leq \lambda} \sum_{i=1}^n (z_i + u_i)^2 \tag{5}$$

$$= \max_{|u_1| \leq \lambda, |u_2| \leq \lambda, \dots, |u_n| \leq \lambda} \sum_{i=1}^n (z_i + u_i)^2 \tag{6}$$

$$= \sum_{i=1}^n \max_{u_i : |u_i| \leq \lambda} (z_i + u_i)^2. \tag{7}$$

The last equality follows since the problem is decomposable and it is optimal to maximize each of these terms. Next we have,

$$\sum_{i=1}^n \max_{u_i : |u_i| \leq \lambda} (z_i + u_i)^2 = \sum_{i=1}^n \max_{\eta : |\eta| \leq \lambda} (z_i + \eta)^2 \tag{8}$$

$$= \sum_{i=1}^n (|z_i| + \lambda)^2, \tag{9}$$

the last line resulting from

$$\forall \eta, |\eta| \leq \lambda : |z_i + \eta| \leq |z_i| + \lambda, \quad (10)$$

with upper bound attained with $\eta = \lambda \text{sign}(z_i)$. Taking $\vec{z} = \vec{v} - x\vec{1}$ we get the desired form.

- (c) Express the problem (4) as a QP. State precisely the variables, and constraints if any.

Solution: A QP formulation is

$$\min_{x,t} \sum_{i=1}^n (t_i + \lambda)^2 : t_i \geq \pm(v_i - x), \quad i = 1, \dots, n. \quad (11)$$

- (d) Show that when λ is large, the solution set approaches that of the median problem (2). *HINT: Given variable a , constants b, c , where $c \gg 1$, and the optimization problem $\min_a \frac{1}{c}(b - a)^2 + |b - a|$. The minimizer a^* tends to minimize the second term only.*

Solution: The objective function takes the form

$$\sum_{i=1}^n (|v_i - x| + \lambda)^2 = n\lambda^2 + 2\lambda \|\vec{v} - x\vec{1}\|_1 + \|\vec{v} - x\vec{1}\|_2^2, \quad (12)$$

The corresponding optimization problem has the same minimizers as the problem

$$\min_x \|\vec{v} - x\vec{1}\|_1 + \frac{1}{2\lambda} \|\vec{v} - x\vec{1}\|_2^2, \quad (13)$$

When λ is large, the minimizer will tend to minimize the first term only, which implies the desired result.

- (e) It is often said that the median is a more robust notion of “middle” value than the mean, when noise is present in \vec{v} . Based on the previous part, justify this statement.

Solution: The median problem can be interpreted as a robust version of the mean problem, when the uncertainty is large.

2. Sphere Enclosure

Let $B_i, i = 1, \dots, m$, be m Euclidean balls in \mathbb{R}^n , with centers \vec{x}_i , and radii $\rho_i \geq 0$. We wish to find a ball B with center $\vec{c} \in \mathbb{R}^n$ of minimum radius $r \geq 0$ that contains all the $B_i, i = 1, \dots, m$. Cast this problem as an SOCP.

Solution: Let $\vec{c} \in \mathbb{R}^n$ and $r \geq 0$ denote the center and radius of the enclosing ball B , respectively. We express the given balls B_i as

$$B_i = \{\vec{x} : \vec{x} = \vec{x}_i + \vec{\delta}_i, \|\vec{\delta}_i\|_2 \leq \rho_i\}, \quad i = 1, \dots, m. \quad (14)$$

We have that $B_i \subseteq B$ if and only if

$$\max_{\vec{x} \in B_i} \|\vec{x} - \vec{c}\|_2 \leq r. \quad (15)$$

Note that

$$\max_{\vec{x} \in B_i} \|\vec{x} - \vec{c}\|_2 = \max_{\|\vec{\delta}_i\|_2 \leq \rho_i} \|\vec{x}_i - \vec{c} + \vec{\delta}_i\|_2 = \|\vec{x}_i - \vec{c}\|_2 + \rho_i. \quad (16)$$

The last step follows by choosing $\vec{\delta}_i$ in the direction of $\vec{x}_i - \vec{c}$. The problem is then cast as the following SOCP

$$\min_{\vec{c}, r} \quad r \quad (17)$$

$$\text{s.t.} \quad \|\vec{x}_i - \vec{c}\|_2 + \rho_i \leq r, \quad i = 1, \dots, m. \quad (18)$$

3. LASSO vs. Ridge

Say that we have the data set $\{(\vec{x}^{(i)}, y^{(i)})\}_{i=1, \dots, n}$ of samples $\vec{x}^{(i)} \in \mathbb{R}^d$ and values $y^{(i)} \in \mathbb{R}$.

Define $X = [\vec{x}^{(1)} \ \dots \ \vec{x}^{(n)}]^\top$ and $y = [y^{(1)} \ \dots \ y^{(n)}]^\top$.

For the sake of simplicity, assume that each feature of the data has mean 0 and variance 1 and the features are uncorrelated, i.e. $X^\top X = nI$. Consider the linear least squares regression with regularization in the ℓ_1 -norm, also known as LASSO:

$$\vec{w}^* = \operatorname{argmin}_{\vec{w} \in \mathbb{R}^d} \|X\vec{w} - \vec{y}\|_2^2 + \lambda \|\vec{w}\|_1. \quad (19)$$

This problem will compare ℓ_1 -regularization with ℓ_2 -regularization (ridge regression) to understand their similarities and differences. We will do this by looking at the elements of \vec{w}^* in the solution to each problem.

- (a) First, we decompose this optimization problem into d univariate optimization problems over each element of \vec{w} . Let $X = [\vec{x}_1 \ \dots \ \vec{x}_d]$ and recall that $X^\top X = nI$.

Solution:

$$\|X\vec{w} - \vec{y}\|_2^2 + \lambda \|\vec{w}\|_1 = \sum_{i=1}^d [nw_i^2 - 2\vec{y}^\top \vec{x}_i w_i + \lambda |w_i|] + \vec{y}^\top \vec{y}. \quad (20)$$

Hence the original problem becomes

$$\min_{\vec{w} \in \mathbb{R}^d} \sum_{i=1}^d [nw_i^2 - 2\vec{y}^\top \vec{x}_i w_i + \lambda |w_i|] + \vec{y}^\top \vec{y}. \quad (21)$$

Since the objective is separable in w_i the problem decomposes into d univariate optimization problems and hence we have

$$\sum_{i=1}^d \min_{w_i \in \mathbb{R}} [nw_i^2 - 2\vec{y}^\top \vec{x}_i w_i + \lambda |w_i|] + \vec{y}^\top \vec{y}. \quad (22)$$

- (b) If $w_i^* > 0$, then what is the value of w_i^* ? What is the condition on $\vec{y}^\top \vec{x}_i$ for this to be possible? *HINT: Use the first order condition.*

Solution: If $w_i^* > 0$, then the first order optimality conditions for w_i^* write

$$2nw_i^* - 2\vec{y}^\top \vec{x}_i + \lambda = 0, \quad (23)$$

from which we obtain

$$w_i^* = \frac{2\vec{y}^\top \vec{x}_i - \lambda}{2n}, \quad (24)$$

which is positive when

$$\vec{y}^\top \vec{x}_i > \frac{\lambda}{2}. \quad (25)$$

- (c) If $w_i^* < 0$, then what is the value of w_i^* ? What is the condition on $\vec{y}^\top \vec{x}_i$ for this to be possible?

Solution: If $w_i^* < 0$, then the first order optimality conditions for w_i^* write

$$2nw_i^* - 2\vec{y}^\top \vec{x}_i - \lambda = 0, \quad (26)$$

from which we obtain

$$w_i^* = \frac{2\vec{y}^\top \vec{x}_i + \lambda}{2n}, \quad (27)$$

which is negative when

$$\bar{y}^\top \bar{x}_i < -\frac{\lambda}{2}. \quad (28)$$

- (d) What can we conclude about w_i^* if $|\bar{y}^\top \bar{x}_i| \leq \frac{\lambda}{2}$? How does the value of λ impact the individual entries w_i^* ?

Solution: From the previous parts we have $w_i^* \neq 0 \Rightarrow |\bar{y}^\top \bar{x}_i| > \frac{\lambda}{2}$. Hence if $|\bar{y}^\top \bar{x}_i| \leq \frac{\lambda}{2}$ then we must have $w_i^* = 0$. This means that a larger value of λ will force more entries of \bar{w} to be zero — i.e. larger λ will imply higher sparsity.

- (e) Now consider the case of ridge regression, which uses the the ℓ_2 regularization $\lambda \|\bar{w}\|_2^2$.

$$\bar{w}^* = \operatorname{argmin}_{\bar{w} \in \mathbb{R}^d} \|X\bar{w} - \bar{y}\|_2^2 + \lambda \|\bar{w}\|_2^2. \quad (29)$$

Write down the new condition for w_i^* to be 0. How does this differ from the condition obtained in part (d) and what does this suggest about LASSO?

Solution: In the case of ridge regression the optimal weight vector \bar{w} is given by

$$w_i^* = \frac{\bar{y}^\top \bar{x}_i}{n + \lambda}, \quad i = 1, \dots, d. \quad (30)$$

So w_i^* is only zero when $\bar{y}^\top \bar{x}_i = 0$, in contrast to LASSO where w_i^* is zero when $\bar{y}^\top \bar{x}_i \in [-\frac{\lambda}{2}, \frac{\lambda}{2}]$. This suggests that LASSO forces a lot of coordinates to be zero, i.e. induces sparsity to the optimal weight vector.

4. More Fun with Lasso and Ridge

Complete the Jupyter notebook `ridge_vs_lasso.ipynb` which demonstrates differences between ridge regression and LASSO.

5. Connecting Ridge Regression, LASSO, and Constrained Least Squares

This question aims to help you develop an understanding of how a constraint in an optimization problem has the same effect as a penalty term in the objective, and apply it to the context of regularized least squares. More formally, let $f: \mathbb{R}^n \rightarrow \mathbb{R}$ be strictly convex and such that $\lim_{t \rightarrow \infty} f(\vec{x}_t) = \infty$ for any sequence $(\vec{x}_t)_{t=0}^{\infty}$ such that $\lim_{t \rightarrow \infty} \|\vec{x}_t\|_2 = \infty$. Let $g: \mathbb{R}^n \rightarrow \mathbb{R}_+$ be convex and take non-negative values. Further, suppose that there exists $\vec{x}_0 \in \mathbb{R}^n$ such that $g(\vec{x}_0) = 0$. For $\lambda \geq 0$ and $k \geq 0$, define the “penalty” and “constraint” programs

$$P(\lambda) \doteq \operatorname{argmin}_{\vec{x}} \{f(\vec{x}) + \lambda g(\vec{x})\} \quad (31)$$

$$C(k) \doteq \operatorname{argmin}_{\vec{x}: g(\vec{x}) \leq k} f(\vec{x}). \quad (32)$$

We will show that for every λ there exists k such that $P(\lambda) = C(k)$, and vice versa.

- (a) Show that, for $k \geq 0$ and $\lambda \geq 0$, both $P(\lambda)$ and $C(k)$ have exactly one element, i.e., each problem has exactly one optimal solution.

HINT: You may use without proof that $P(\lambda)$ and $C(k)$ have at least one element each (this is true from assumptions but requires some analysis to show). Thus, you just need to show that there are not multiple optimal solutions to each problem. For this, use strict convexity of the objectives.

Solution: Since f is strictly convex and g is convex, both problems are convex with strictly convex objective. Thus both problems have at most one solution. Since we know from the assumptions that both problems have at least one solution, $P(\lambda)$ and $C(k)$ have exactly one element each.

We give some clarity on the fact that $P(\lambda)$ and $C(k)$ have at least one element each. This requires some analysis to prove, and so **it is all out of scope of the course**. We know that convex functions (hence also strictly convex functions) are continuous (see [this MathStackExchange link](#)). The assumption that $f(\vec{x}) \rightarrow \infty$ as $\|\vec{x}\|_2 \rightarrow \infty$ is called “coercivity”. We know that continuous coercive functions on closed sets (such as \mathbb{R}^n and $\{\vec{x} \in \mathbb{R}^n \mid g(\vec{x}) \leq k\}$) attain their global minima, i.e., their problems have at least one solution (see [this MathStackExchange link](#)). This is how we are able to say that $P(\lambda)$ and $C(k)$ have at least one element each.

- (b) Prove that for all $\lambda \geq 0$ there exists $k \geq 0$ such that $P(\lambda) = C(k)$.

HINT: Let $\vec{x}^ \in P(\lambda)$ and show that $\vec{x}^* \in C(k)$ for $k = g(\vec{x}^*)$. Use the fact, from part 5(a), that $P(\lambda)$ and $C(k)$ have exactly one element.*

Solution: Let $\lambda \geq 0$ and let $\vec{x}^* \in P(\lambda)$, so that $P(\lambda) = \{\vec{x}^*\}$. Set $k = g(\vec{x}^*)$. Certainly \vec{x}^* is feasible for the constraint problem with this k . We claim that $\vec{x}^* \in C(k)$, so that $C(k) = \{\vec{x}^*\}$. Suppose for the sake of contradiction that $C(k) = \{\vec{z}\}$ where $\vec{z} \neq \vec{x}^*$. Then \vec{z} must be feasible for the constraint problem, so $g(\vec{z}) \leq k = g(\vec{x}^*)$, and it must be better than \vec{x}^* , so $f(\vec{z}) < f(\vec{x}^*)$. Then we have

$$f(\vec{z}) + \lambda g(\vec{z}) < f(\vec{x}^*) + \lambda g(\vec{z}) \quad (33)$$

$$\leq f(\vec{x}^*) + \lambda g(\vec{x}^*) \quad (34)$$

so $\vec{x}^* \notin P(\lambda)$, a contradiction. Thus $\vec{x}^* \in C(k)$, i.e., $C(k) = \{\vec{x}^*\}$.

- (c) Prove that for all $k > 0$ there exists $\lambda \geq 0$ such that $P(\lambda) = C(k)$.

HINT: Prove that strong duality holds for the constraint problem, let $\vec{x}^ \in C(k)$ and μ^* be optimal primal and dual variables for the constraint problem and show that $\vec{x}^* \in P(\lambda)$ for $\lambda = \mu^*$.*

Solution: We know that there is a point \vec{x}_0 such that $g(\vec{x}_0) = 0 < k$. Thus the constraint problem has a strictly feasible point. Since the constraint problem is a convex problem with a strictly feasible point, strong duality holds for it. The Lagrangian of the constraint problem is given by

$$L_k(\vec{x}, \mu) = f(\vec{x}) + \mu(g(\vec{x}) - k). \quad (35)$$

Since strong duality holds, let (\vec{x}^*, μ^*) be optimal primal and dual variables for the constraint problem (so that $\vec{x}^* \in C(k)$, i.e., $C(k) = \{\vec{x}^*\}$) such that

$$p^* = L_k(\vec{x}^*, \mu^*) = d^*. \quad (36)$$

Since the constraint problem is convex, and strong duality holds, we have

$$\vec{x}^* \in \underset{\vec{x}}{\operatorname{argmin}} L_k(\vec{x}, \mu^*) \quad (37)$$

$$= \underset{\vec{x}}{\operatorname{argmin}} \{f(\vec{x}) + \mu^*(g(\vec{x}) - k)\} \quad (38)$$

$$= \underset{\vec{x}}{\operatorname{argmin}} \{f(\vec{x}) + \mu^*g(\vec{x}) - \mu^*k\} \quad (39)$$

$$= \underset{\vec{x}}{\operatorname{argmin}} \{f(\vec{x}) + \mu^*g(\vec{x})\} \quad (40)$$

$$= P(\mu^*), \quad (41)$$

so $\vec{x}^* \in P(\lambda)$, i.e., $P(\lambda) = \{\vec{x}^*\}$, with $\lambda = \mu^*$.

Now we apply our findings to regularized least squares, in order to understand why LASSO promotes sparsity more than ridge regression. Let $A \in \mathbb{R}^{m \times n}$ have full column rank, and let $\vec{y} \in \mathbb{R}^m$. In the course, we have looked at LASSO:

$$\text{LASSO}(\lambda) \doteq \underset{\vec{x}}{\operatorname{argmin}} \left\{ \|A\vec{x} - \vec{y}\|_2^2 + \lambda \|\vec{x}\|_1 \right\} \quad (42)$$

and ridge regression:

$$\text{Ridge}(\lambda) \doteq \underset{\vec{w}}{\operatorname{argmin}} \left\{ \|A\vec{w} - \vec{y}\|_2^2 + \lambda \|\vec{w}\|_2^2 \right\} \quad (43)$$

which add an ℓ^1 and ℓ^2 norm penalty to the least squares objective, respectively. The analogous constraint programs are the ℓ^1 - and ℓ^2 -constrained least squares problems:

$$\ell^1\text{CLS}(k) \doteq \underset{\vec{x}: \|\vec{x}\|_1 \leq k}{\operatorname{argmin}} \|A\vec{x} - \vec{y}\|_2^2 \quad (44)$$

$$\ell^2\text{CLS}(k) \doteq \underset{\vec{x}: \|\vec{x}\|_2 \leq k}{\operatorname{argmin}} \|A\vec{x} - \vec{y}\|_2^2. \quad (45)$$

- (d) Show that the result from part 5(b) and part 5(c) can be used to show the equivalence of LASSO with $\ell^1\text{CLS}$ and the equivalence of ridge regression with $\ell^2\text{CLS}$. Namely, for each pair of equivalent formulations, find f and g , prove that f is strictly convex, prove that g is convex, and prove that there is an \vec{x}_0 such that $g(\vec{x}_0) = 0$.

Solution: For all formulations,

$$f(\vec{x}) \doteq \|A\vec{x} - \vec{y}\|_2^2 \quad (46)$$

has Hessian

$$\nabla_{\vec{x}}^2 f(\vec{x}) = 2A^\top A \quad (47)$$

which is PD since A has full column rank. Thus f is strictly convex (and in fact is strongly convex), and in particular $f(\vec{x}_t) \rightarrow \infty$ as $t \rightarrow \infty$ for any $(\vec{x}_t)_{t=0}^{\infty}$ such that $\|\vec{x}_t\|_2 \rightarrow \infty$.

In the LASSO- ℓ^1 CLS formulations, $g(\vec{x}) = \|\vec{x}\|_1$, which is convex since it is a norm. In the LASSO- ℓ^2 CLS formulations, $g(\vec{x}) = \|\vec{x}\|_2^2$, which is (strictly) convex since it has Hessian $2I$ which is PD.

For both g , $g(\vec{0}) = 0$.

- (e) Complete the Jupyter notebook `ridge_lasso_constrained.ipynb`, which will use this equivalence to show geometrically why LASSO solutions tend to be sparse (i.e. have many zeros) while ridge regression doesn't, and attach a PDF printout of your answers.

Solution: See the Jupyter notebook solutions.

6. Homework Process

With whom did you work on this homework? List the names and SIDs of your group members.

NOTE: If you didn't work with anyone, you can put "none" as your answer.