

Consider the Lasso Problem

$$\min_{\beta} \frac{1}{2} \|X\beta - y\|_2^2 + \gamma \|\beta\|_1$$

for $\beta \in \mathbb{R}^P$, X as an $n \times p$ matrix. This is a normalized matrix X , usually $\|X\delta_i\| \leq 1$. For classical ML setups, we might have loss functions corresp. to regression or classification. E.g.,

$$\min_{\beta} \frac{1}{n} \sum_i \phi_i(x_i^T \beta) + \gamma \|\beta\|_1$$

where $x_i = \delta_i^T X$. Convex conjugates are:

Squared $\frac{\phi(z)}{2} = \frac{1}{2} (z - y_i)^2$ $\phi^*(z) = \frac{1}{2} (z + y_i)^2 - \frac{1}{2} y_i^2$

Logistic $\log(1 + \exp(-y_i z))$ $\frac{-z}{y_i} \log\left(\frac{-z}{y_i}\right) + \left(1 + \frac{z}{y_i}\right) \log\left(1 + \frac{z}{y_i}\right)$

See Blitz: A Pivoted Meta
Algorithm Appendix E
for derivation

But the original problem is unconstrained.

So the dual of $\min_{\beta} L(\beta) + \lambda \|\beta\|_1$

is just $\max_{\beta} p^{\star \beta}$, which is trivial and unhelpful.

So we need to add constraints which will expose some perturbations, such that duals provide a useful optimization view. There will be 2 distinct duals (I tried messing with them, but they seem to be distinct). All of them will have Strong Duality by relaxed Slater

n-parameter dual.

$$\min_{\substack{\beta \in \mathbb{R}^p \\ r \in \mathbb{R}^n}} \frac{1}{2} \|r\|_2^2 + \lambda \|\beta\|_1$$

$$\text{st } X\beta - y = r$$

introduce the dual variable $\nu \in \mathbb{R}^n$

$$\max_{\nu} \inf_{\beta, r} \frac{1}{2} \|r\|_2^2 + \lambda \|\beta\|_1 + \nu^T (\chi\beta - y - r)$$

$\rightarrow \inf_{\beta} \lambda \|\beta\|_1 + \nu^T \chi\beta = \begin{cases} 0 & \|X^T \nu\|_{\infty} \leq 1 \\ -\infty & \text{o/w} \end{cases}$
 by looking at $-\beta$ & the convex conjugate of $\|\cdot\|_1$

$$= \max_{\nu} -\nu^T y - \frac{1}{2} \|\nu\|_2^2 + \inf_r \underbrace{\frac{1}{2} \|r\|_2^2 - \nu^T r + \frac{1}{2} \|\nu\|_2^2}_{\frac{1}{2} \|\nu - r\|_2^2}$$

st $\|X^T \nu\|_{\infty} \leq 1$

$0 @ r = \nu.$

$$= \min_{\nu} -\nu^T y - \frac{1}{2} \|\nu\|_2^2$$

$$\text{st } \|X^T \nu\|_{\infty} \leq 1$$

Now lets generalize this for GLMs.

$$\min_{\beta} \frac{1}{n} \sum_i \phi_i(x_i^T \beta) + \lambda \|\beta\|_1$$

Similarly, consider constraints:

$$\begin{aligned} \min_{\beta, z} \quad & \frac{1}{n} \sum_i \phi_i(z_i) + \lambda \|\beta\|_1 \\ \text{s.t.} \quad & X\beta = z \end{aligned}$$

$$= \max_v \min_{\beta, z} \underbrace{v^T X\beta + \lambda \|\beta\|_1}_{\text{as before}} + \underbrace{\frac{1}{n} \sum_i \phi_i(z_i) - v_i z_i}_{\text{negative of cvx conj.}}$$

$$= \max_v \frac{1}{n} \sum_i -\phi_i^*(v_i) \quad \text{s.t.} \quad \|X^T v\|_{\infty} \leq 1$$

$$= - \min_v \frac{1}{n} \sum_i \phi_i^*(v_i) \\ \text{s.t.} \quad \|X^T v\|_{\infty} \leq 1$$

P-parameter dual

$$\min_{\beta} \frac{1}{2} \|X\beta - y\|_2^2 + \lambda \|\beta\|_1$$

$$= \min_{\beta, \gamma} \frac{1}{2} \|X\beta - y\|_2^2 + \lambda \|\gamma\|_1$$

$$\text{s.t. } \gamma = \beta \quad (\gamma \in \mathbb{R}^p)$$

$$= \max_w \min_{\beta, \gamma} \frac{1}{2} \|X\beta - y\|_2^2 + \lambda \|\gamma\|_1 + w^T(\beta - \gamma)$$

$$= \max_{w: \|w\|_\infty \leq \lambda} \min_{\beta} \frac{1}{2} \|X\beta - y\|_2^2 + w^T \beta$$

At this point, let's pause for a minute.

Note the Moore-Penrose Pseudoinverse X^+ .

If the SVD is $U \Sigma V^T = X$, then

$X^+ = V \Sigma^{-1} U^T$, where the SVD here

is reduced i.e., Σ is $(\text{rank } X) \times (\text{rank } X)$.

Note the last-normal least squares solution
is then $X^+ y = \beta_{LS}$ for minimizing $\|X\beta - y\|_2^2$.
This does not require X to be full rank.

A critical property of X^+ is projection:

$X^+ X$ is an orthogonal projection
onto $\text{span}(X^T)$ from \mathbb{R}^P

This implies $I - X^+ X = \text{proj}_{\text{ker } X}$ so

$$X(I - X^+ X) = 0. \text{ Thus we}$$

can write any $\beta^* = X^+ \alpha + (I - X^+ X) \delta$

for some α , δ arbitrarily defined.

Notice that then we can write our old
optimization in terms of $\alpha \in \mathbb{R}^n$, $\delta \in \mathbb{R}^P$

$$\min_{\beta} \frac{1}{2} \|X\beta - y\|_2^2 + w^T \beta$$

$$= \min_{\alpha, \beta} \frac{1}{2} \|XX^T \alpha - y\|_2^2 + w^T (X^T \alpha + (I - X^T X)\beta)$$

if $w \notin \text{span}(X^T)$ then $\text{rank } X < p$

and $I - X^T X \neq 0$. Choose

$$\beta = n \left(-(I - X^T X) \right) w$$

results in the obj $\rightarrow -\infty$ as $n \rightarrow \infty$.

thus this minimum creates an explicit constraint on w .

This leaves $\min_{\alpha} \|XX^T \alpha - y\|_2^2 + w^T (X^T \alpha)$

$$\text{as } w \in \text{span}(X^T) \Rightarrow w^T (I - X^T X) = 0$$

$$\min_{\alpha} \frac{1}{2} \|X X^T \alpha\|_2^2 - y^T (X X^T \alpha) + \frac{1}{2} \|y\|_2^2 + \omega^T (X^T \alpha).$$

By projection $X^+ = X^T X X^T$, so

$$= \min_{\alpha} \frac{1}{2} \|X X^T \alpha\|_2^2 - (y - (X^+)^T \omega)^T X X^T \alpha + \frac{1}{2} \|y - (X^+)^T \omega\|_2^2 + (X^+ y)^T \omega - \frac{1}{2} \|(X^+)^T \omega\|_2^2$$

$$= \min_{\alpha} \frac{1}{2} \|X X^T \alpha - (y - (X^+)^T \omega)\|_2^2 + (X^+ y)^T \omega - \frac{1}{2} \|(X^+)^T \omega\|_2^2$$

And now we can change the optimisation back!

as $X X^T \alpha = X \beta$ for some $\beta \in \mathbb{R}^p$
 but then we see $\alpha^* = (y - (X^+)^T \omega)$

Then using the identity $XX^+(X^+)^T = (X^+)^T$,

$$\begin{aligned} & \|XX^+ \alpha^* - (y - (X^+)^T w)\|_2^2 \\ &= \|XX^+ y - y\|_2^2 \end{aligned}$$

which simplifies the original problem to

$$\begin{aligned} \max_w & (X^+ y)^T w - \|(X^+)^T w\|_2^2 \\ & + \frac{1}{2} \|(I - XX^+)y\|_2^2 \end{aligned}$$

$$\text{s.t. } \|w\|_\infty \leq \lambda$$

$$w \in \text{span } X^T$$

With full-rank columns, i.e. $\text{rank } X = p$,

$$X^+ = (X^T X)^{-1} X^T \text{ and } w \in \text{span } X^T \text{ holds, so we may simplify}$$

$$\max_w \beta_{LS}^T w + \frac{w^T (X^T X)^{-1} w}{2} +$$

$$s.t. \quad \|w\|_\infty \leq \lambda, \quad \frac{y^T (I-P)y}{2}$$

$$\text{where } \beta_{LS} = (X^T X)^{-1} X^T y$$

$$P = X (X^T X)^{-1} X$$

Since $I-P$ is idempotent.

$$\begin{aligned} \text{Note } \frac{y^T (I-P)y}{2} &= \frac{1}{2} \|(I-P)y\|_2^2 \\ &= \frac{RSS(\beta_{LS})}{2} \end{aligned}$$

Extending to ϕ :

generally seems harder, I got an np-harder dual that way which seems useless.