**Example:** Linear model with additively corrupted covariates

Consider

$$y_i = x_i^T \underset{\text{ground truth}}{\theta^*} + \varepsilon_i \quad , \quad y_i \in \mathbb{R}, \quad x_i \in \mathbb{R}^d, \quad \theta^* \in \mathbb{R}^d, \quad \varepsilon_i \in \mathbb{R} \quad - (L)$$

- response variable
- covariate vector
- noise  $\mathbb{E}[\varepsilon_i] = 0$

More compactly,

$$y = X\theta^* + \varepsilon$$

The estimation problem

$$\hat{\theta} \in \underset{\theta \in \mathbb{R}^d}{\arg\min} \left\{ \frac{1}{2n} \| y - X\theta \|_2^2 + R(\theta) \right\} \qquad (LS)$$

Can be viewed as a sample average version of the following idealized problem:

$$\hat{\theta} \in \underset{\theta \in \mathbb{R}^d}{\arg\min} \left\{ \frac{1}{2} \theta^T \underset{\frac{1}{n}X^TX}{\Sigma_x} \theta - \theta^T \underset{\frac{1}{n}X^Ty}{\Sigma_x \theta^*} + R(\theta) \right\}, \qquad (\mathcal{I})$$

where $\Sigma_x > 0$ is the covariance of $\{x_i\}$, which are assumed to be iid mean-zero. This follows from

$$\Sigma_x = \mathbb{E}[xx^T], \qquad \underset{\varepsilon_i}{\mathbb{E}}[y_i x_i] = x_i x_i^T \theta^*,$$

$$X^TX = \sum_{i=1}^{n} x_i x_i^T, \qquad X^Ty = X^T(X\theta^* + \varepsilon)$$

Now, suppose that $x_i$ is not observed directly, but rather we observe a $z_i \in \mathbb{R}^d$ that is related to $x_i$ via

$$z_i = x_i + w_i,$$

$$x_i \in \mathbb{R}^d, \quad z_i \in \mathbb{R}^d,$$
$$i = 1, \cdots, n \quad (n \ll d)$$

where $w_i \in \mathbb{R}^d$ is the noise vector, assumed to be zero mean and have covariance $\Sigma_w$ (known), and $w_i$ is independent of $x_i$ and $\varepsilon_i$. Following our previous argument, we compute

$$\frac{1}{n} \underset{w}{\mathbb{E}}[z^Tz] = \frac{1}{n} X^TX + \Sigma_w, \qquad \frac{1}{n} \underset{w}{\mathbb{E}}[z^Ty] = \frac{1}{n} X^Ty$$

Then, we see that in this setting, we can use

$$\hat{\Gamma} = \frac{1}{n} z^Tz - \Sigma_w \quad , \quad \hat{\gamma} = \frac{1}{n} z^Ty$$

$$\hat{\Gamma} = \frac{1}{n}\underbrace{Z^T Z}_{d\times d} - \Sigma_w \quad , \quad \hat{\gamma} = \frac{1}{n}Z^T y$$

and this gives rise to

$$\hat{\Theta} \in \underset{\|\Theta\|_1 \leq R}{\arg\min} \left\{ \frac{1}{2}\Theta^T \hat{\Gamma}\Theta - \hat{\gamma}^T\Theta \right\}$$

<span style="color:red">make sure $\hat{\Theta}$ exists</span>

$$\mathbb{R}^{n\times d} \ni Z = \begin{bmatrix} - z_1^T - \\ \vdots \\ - z_n^T - \end{bmatrix}$$

$$\frac{1}{n}\sum_{i=1}^{n} z_i z_i^T$$

$$=$$

Note that $\hat{\Gamma}$ need not be p.s.d. Indeed, $\frac{1}{n}Z^T Z$ has rank at most $n$ but $\Sigma_w \in S_+^d$ could have rank $d$ (e.g., $\Sigma_w = I_d$). Hence, the above formulation is non-convex in general. We can also consider the regularized version
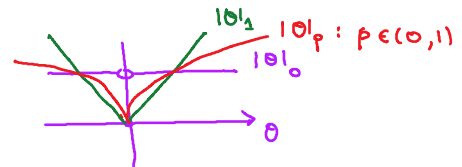
$$\hat{\Theta}' \in \underset{\|\Theta\|_1 \leq R'}{\arg\min} \left\{ \frac{1}{2}\Theta^T \hat{\Gamma}\Theta - \hat{\gamma}^T\Theta + \lambda\|\Theta\|_1 \right\}$$

<span style="color:red">make sure $\hat{\Theta}$ exists</span>

Example: Non-Convex regularizer

When dealing with sparse model, we want to use $\ell_0$-quasi-norm to measure sparsity.
This function is challenging, so we replace it by a convex surrogate:

$\ell_1$-norm. To get a better approximation of $\|\cdot\|_0$, one can consider non-convex (but continuous) surrogates

e.g. $\|\Theta\|_p^p = \sum_{i=1}^{d} |\Theta_i|^p \quad \ell_p$-(quasi)-norm if $p \in (0,1)$

(Bridge penalty)

However, this is not quite desirable, because

$$\lim_{\Theta\searrow 0}(|\Theta|^p)' = +\infty$$

Consider the following family of regularizers:

$$\mathbb{R}^d \ni \theta \mapsto R_\lambda(\theta) = \sum_{i=1}^{d} R_\lambda(\theta_i)$$

<span style="color:red">↑ Separable</span>

And the function $R_\lambda : \mathbb{R} \to \mathbb{R}_+$ satisfies the following properties:

(1) $R_\lambda(0) = 0$ and $R_\lambda(t) = R_\lambda(-t)$ $\forall t \in \mathbb{R}$,

(2) $R_\lambda$ is non-decreasing on $\mathbb{R}_+$

(3) For $t > 0$, the function $t \mapsto \dfrac{R_\lambda(t)}{t}$ is non-increasing in $t$.

(4) $R_\lambda$ is differentiable at any $t \neq 0$ and <u>subdifferentiable</u>

<span style="color:red">needs definition b/c $R_\lambda$ may not be convex</span> <span style="color:orange">✓ (see below)</span>

at $t = 0$ with $\lim\limits_{t \searrow 0} R_\lambda'(t) = \lambda L$ for some $L > 0$.

(5) There exists a $\mu > 0$ s.t. $t \mapsto R_{\lambda,\mu}(t) \triangleq R_\lambda(t) + \frac{\mu}{2} t^2$ is convex (in other words, $R_\lambda$ is $\mu$-weakly convex).

With (5), we can write

$$R_\lambda(t) = R_{\lambda,\mu}(t) - \frac{\mu}{2} t^2.$$

Hence, formally, we can write

$$\partial R_\lambda(t) = \partial\left( R_{\lambda,\mu}(t) - \frac{\mu}{2} t^2 \right)$$

$$= \partial R_{\lambda,\mu}(t) - \mu t$$

<span style="color:red">Well-defined b/c $R_{\lambda,\mu}$ is convex</span>

Hence, we can define

$$\partial R_\lambda(t) \triangleq \partial R_{\lambda,\mu}(t) - \mu t.$$