Part 3  Implicit Regularization / linear regression

## Question 1

Prove the following proposition.

### Proposition:

With the notations and setting established in class suppose we minimize $L_S(w)$ by initializing $w^{[0]} = a \in \mathbb{R}^d$ and producing iterates $w^{(1)}, w^{(2)}, \dots$ via iterative algorithm in which every update $w^{(t+1)} - w^{(t)} \in \text{span} \{ \nabla \ell_{(x_i, y_i)}(w) : i \in [m], w \in \mathbb{R}^d \}$.

Assume convergence to global min with zero loss.

Then the sub-optimality of the obtained norm (i.e. the extent to which it is larger than min norm across all global optima) is $\leq \| P_\perp a \|$ where $P_\perp : \mathbb{R}^d \to \mathbb{R}^d$ stands for projection onto the orthogonal complement of span $\{x_i\}_{i=1}^m$.

### Proof

Denote by

$$W^* = \arg \min_{\substack{w \in \mathbb{R}^d \text{ global min} \\ L_S(w) = 0}} \| w \|$$ the solution

with the min norm across all global optima.

### Lemma 1

The solution under the settings in the proposition is

$$\widetilde{W} = a + X(X^T X)^{-1} (y - X^T a)$$

With the markings above and as far as lemma 1 is true we have to prove that

$$\| \widetilde{W} \| \leq \| w^* \| + \| P_\perp a \|$$

in order to finish the question.

## Proof of lemma 1

First, we'll prove that for every $t \in \mathbb{N}$ $w^t \in \underbrace{a + \mathrm{span}\{x_i\}_{i=1}^m}_{\text{affine subspace}}$

by induction.

base: $W^0 = a = a + \sum_{i=1}^m 0 \cdot x_i \in a + \mathrm{span}\{x_i\}_{i=1}^m$

Step: Suppose $w^t \in a + \mathrm{span}\{x_i\}_{i=1}^m$. We'll prove that $w^{t+1} \in a + \mathrm{span}\{x_i\}_{i=1}^m$

$w^t \in a + \mathrm{span}\{x_i\}_{i=1}^m \implies w^t = a + \sum_{i=1}^m \alpha_i x_i$ for $\{\alpha_i\}_{i=1}^m$ scalars.

Note that for any $i \in [m]$ and $w \in \mathbb{R}^d$

$$\nabla \ell (w, x_i, y_i) = (x_i^T w - y_i) \cdot x_i$$

This implies that for every $t \in \mathbb{N}$

$$w^{t+1} - w^t \in \mathrm{span}\{x_i\}_{i=1}^m$$

(Since we know that $w^{t+1} - w^t \in \mathrm{span}\{\nabla \ell_{(x_i, y_i)}(w) | i \in [m] \; w \in \mathbb{R}^d\}$)

$\implies w^{t+1} - w^t = \sum_{i=1}^m \beta_i x_i$ for some $\{\beta_i\}_{i=1}^m$ scalars.

$\implies w^{t+1} - (a + \sum_{i=1}^m \alpha_i x_i) = \sum_{i=1}^m \beta_i x_i$

$\implies w^{t+1} = a + \sum_{i=1}^m (\alpha_i + \beta_i) x_i \implies w^{t+1} \in a + \mathrm{span}\{x_i\}_{i=1}^m$

$\square$

Hence we get for every $t \in \mathbb{N}$ $w^t \in a + \mathrm{span}\{x_i\}_{i=1}^m$.

Since the subspace $a + \mathrm{span}\{x_i\}_{i=1}^m$ is topologically closed

$\tilde{w} = \lim_{t \to \infty} w^t \in a + \mathrm{span}\{x_i\}_{i=1}^m$

$\implies \tilde{w} = a + \sum_{i=1}^m r_i x_i$ for some scalars $\{r_i\}_{i=1}^m$

$\implies \tilde{w} = a + X r$ for $r = \begin{pmatrix} r_1 \\ \vdots \\ r_m \end{pmatrix}$

We've got $\tilde{w} = a + Xt$.

Since the convergence is to global min with zero loss we get

$$X^T\tilde{w} = y \implies X^T(a + Xv) = y$$
$$\implies X^Ta + X^TXt = y$$
$$\implies X^TXt = y - X^Ta$$

$\text{rank}(X^TX) = m \implies t = (X^TX)^{-1}(y - X^Ta)$

$$\implies \tilde{w} = a + X(X^TX)^{-1}(y - X^Ta)$$

as required. $\square$

Now, remember that at class we proved (in class) that the solution

$$w^* = X(X^TX)^{-1}y \quad \text{is the one with minimal norm.}$$

Thus we have to show that

$$\|a + X(X^TX)^{-1}(y - X^Ta)\| = \|\tilde{w}\| \leq \|w^*\| + \|P_\perp a\| = \|X(X^TX)^{-1}y\| + \|P_\perp a\|$$

We'll show that:

$$\|a + X(X^TX)^{-1}(y - X^Ta)\| = \|a + X(X^TX)^{-1}y - X(X^TX)^{-1}X^Ta\|$$

$$\leq \|X(X^TX)^{-1}y\| + \|a - X(X^TX)^{-1}X^Ta\|$$

$\triangle$ inequality

$$= \|X(X^TX)^{-1}y\| + \|P_\perp a + P_{\|}a - X(X^TX)^{-1}X^T(P_\perp a + P_{\|}a)\|$$

$a = P_{\|}a + P_\perp a$
where $P_{\|}a \in \text{span}\{x_i\}$ and $P_\perp a$ is the projection as described in the question.

Notice that by definition $\langle x_i, P_\perp a\rangle = 0 \quad \forall i \implies X^T P_\perp a = \begin{pmatrix} 0 \\ \vdots \\ 0 \end{pmatrix}$

thus:

$\ominus \quad \|X(X^TX)^{-1}y\| + \|P_\perp a + P_{\|}a - X(X^TX)^{-1}X^T P_{\|}a\| \ominus$

$P_{\|}a \in \text{span}\{x_i\}_{i=1}^m$, thus $P_{\|}a = X \cdot b$ for some $b = \begin{pmatrix} b_1 \\ \vdots \\ b_m \end{pmatrix}$

$\ominus \quad \|X(X^TX)^{-1}y\| + \|P_\perp a + Xb - X(X^TX)^{-1}X^TXb\|$

$= \|X(X^TX)^{-1}y\| + \|P_\perp a\|$

Which means that the suboptimality of the obtained norm $\leq \|P_\perp a\|$

as required. $\square$