

TIIC31020 HW1 : Regression

Part I : Linear Regression

Collaborated with Haichuan Wang, Lily Zhu

Problem 1 [6 points]

Suppose we take the data set from which we learned (estimated) w^* , and for every example (x_i, y_i) , change y_i to $y'_i = ay_i + b$, for some constants a and b (same a and b for all the i s). Now you fit least squares model to the new data set $\{x_i, y'_i\}$, yielding the parameter vector w' .

Can w' be computed directly from w^* , without looking at the data again? If yes, how exactly? If not, why not?

• Yes.

$$w^* = (X^T X)^{-1} X^T y$$

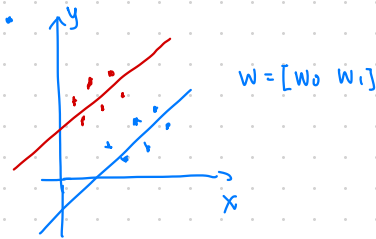
$$w' = (X^T X)^{-1} X^T (ay + b \begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix})$$

$$= a(X^T X)^{-1} X^T y + b(X^T X)^{-1} X^T \begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix}$$

$$\text{We can write } \begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix} = \underbrace{\begin{pmatrix} 1 & x_{11} & \dots & x_{1p} \\ \vdots & \vdots & & \vdots \\ 1 & x_{n1} & \dots & x_{np} \end{pmatrix}}_X \cdot \begin{pmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{pmatrix}$$

$$\text{Then, } w' = a(X^T X)^{-1} X^T y + b(X^T X)^{-1} X^T X \cdot \begin{pmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{pmatrix}$$

$$= a \cdot w^* + b \cdot \begin{pmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{pmatrix}$$



Problem 2 [7 points]

Now, instead of modifying y s, we will modify x s. For every feature (dimension of x) we will change x_{ij} to $\tilde{x}_{ij} = c_j x_{ij}$ for some constants c_j (same set of c_1, \dots, c_d for all i s). Again, we fit least squares model to the new data set $\{\tilde{x}_i, y_i\}$, and get w' .

Can w' be computed directly from w^* , without looking at the data again? If yes, how exactly? If not, why not?

• Yes.

$$\text{Let } X' = \begin{pmatrix} x_{11} & \dots & x_{1p} \\ \vdots & & \vdots \\ x_{n1} & \dots & x_{np} \end{pmatrix} \quad D = \begin{pmatrix} 1 & 0 & \dots & 0 \\ 0 & c_1 & & \\ \vdots & & \ddots & \\ 0 & & & c_p \end{pmatrix}$$

$$\text{Now, } w' = (X' D^T X' D)^{-1} (X' D)^T y$$

$$= (D^T X^T X D)^{-1} D^T X^T y$$

$$= D^{-1} \cdot (X^T X)^{-1} (D^T)^T \cdot D^T X^T y$$

$$= D^{-1} (X^T X)^{-1} X^T y$$

$$= D^{-1} \cdot w^*$$

$$= \begin{pmatrix} 1 & 0 & \dots & 0 \\ \frac{1}{c_1} & & & \\ \vdots & & \ddots & \\ 0 & & & \frac{1}{c_p} \end{pmatrix} \cdot w^*$$

In Lecture 2 we saw two necessary conditions for the optimal w as given below.

$$\sum^N (y_i - \mathbf{w} \cdot \mathbf{x}_i) = 0 \quad (1)$$

$$\forall j = 1, \dots, d: \sum_{i=1}^N (y_i - \mathbf{w} \cdot \mathbf{x}_i) x_{ij} = 0 \quad (2)$$

We are going to confirm that these indeed are necessary conditions for optimal linear model \mathbf{w} in a different way, without reasoning about derivatives, instead proving it by explicit contradiction.

First, consider a linear model \mathbf{w} which violates (1),

$$\sum_{i=1}^N (y_i - \mathbf{w} \cdot \mathbf{x}_i) = p \neq 0. \quad (3)$$

Problem 3 [6 points]

Show a model \mathbf{w}_p which satisfies (1) and has a lower empirical squared loss than \mathbf{w} .

• Let $W_p = W + \begin{pmatrix} p \\ 0 \\ 0 \\ 0 \\ 0 \end{pmatrix}$,

$$\begin{aligned} L(w_p) &= \frac{1}{N} \sum_{i=1}^N (y_i - w_p x_i)^2 \\ &= \frac{1}{N} \sum_{i=1}^N (y_i - w x_i - \begin{pmatrix} \frac{p}{N} \\ 0 \\ 0 \end{pmatrix} \cdot \begin{pmatrix} 1 & x_{i1} & \dots & x_{ip} \end{pmatrix})^2 \\ &= \frac{1}{N} \sum_{i=1}^N (y_i - w x_i - \frac{p}{N})^2 \\ &= \frac{1}{N} \sum_{i=1}^N \left((y_i - w x_i)^2 - 2 (y_i - w x_i) \cdot \frac{p}{N} + \frac{p^2}{N^2} \right) \\ &= L(w) - \frac{2}{N} \cdot p \cdot \frac{p}{N} + \frac{1}{N} \cdot N \cdot \frac{p^2}{N^2} \\ &= L(w) - \frac{p^2}{N^2} \end{aligned}$$

Therefore, $L(W_p) < L(W)$.

Now do the same for the other condition: consider a linear model \mathbf{w} which violates (2),

$$\forall j = 1, \dots, d: \sum_{i=1}^N (y_i - \mathbf{w} \cdot \mathbf{x}_i) x_{ij} = s \neq 0. \quad (4)$$

Problem 4 [7 points]

Show a model \mathbf{w}_s which satisfies (2) and has a lower empirical squared loss than \mathbf{w} .

- Pick $k \in [1, d]$, let $a = \frac{S}{\sum_{i=1}^d x_{ik}^2}$
- Define $W_s = \begin{pmatrix} w_0 \\ \vdots \\ w_{k+a} \\ \vdots \\ w_d \end{pmatrix} = W + \begin{pmatrix} 0 \\ \vdots \\ a \\ \vdots \\ 0 \end{pmatrix}$

$$\begin{aligned} \text{Then, } L(w_s) &= \frac{1}{N} \sum_{i=1}^N (y_i - w_s x_i)^2 \\ &= \frac{1}{N} \sum_{i=1}^N (y_i - (w + \frac{a}{b}) x_i)^2 \\ &= \frac{1}{N} \sum_{i=1}^N (y_i - w x_i - a \cdot x_{ik})^2 \\ &= \frac{1}{N} \sum_{i=1}^N (y_i - w x_i)^2 - 2 \cdot \frac{1}{N} \sum_{i=1}^N (y_i - w x_i) \cdot x_{ik} \cdot a + \frac{1}{N} \sum_{i=1}^N a^2 x_{ik}^2 \\ &= L(w) - 2 \cdot \frac{1}{N} s \cdot a + \frac{1}{N} a^2 \sum_{i=1}^N x_{ik}^2 \\ &= L(w) - 2 \cdot \frac{1}{N} s \cdot a + \frac{1}{N} a \cdot \frac{s}{N \cdot \sum x_{ik}} \cdot \sum x_{ik}^2 \\ &= L(w) - \frac{1}{N} s (2a - a) \\ &= L(w) - \frac{\frac{1}{N} s a}{1} \rightarrow = \frac{s^2}{N^2 \cdot \sum x_{ik}} > 0 \end{aligned}$$

Therefore, $L(w_s) < L(w)$

Part 2: Asymmetric Loss

In class we have discussed the idea of learning (via empirical risk minimization) by finding the parameter values that minimize the loss function on the training data. In the case of least squares regression there is a closed form solution. We will now consider a modification, and develop a gradient descent solution for minimizing it.

Consider an *asymmetric loss* function:

$$\ell_\alpha(y, \hat{y}) = \begin{cases} \alpha (y - \hat{y})^2 & \text{if } y \geq \hat{y} \\ (y - \hat{y})^2 & \text{if } y < \hat{y} \end{cases} \quad (5)$$

Problem 5 [9 points]

Assuming $\alpha > 1$, what kind of desired properties of the predictor does this loss function capture? When would you want to use it instead of the standard least squares? Feel free to provide a specific (possibly made up!) application scenario, or describe such a scenario in general terms.

• penalties under-predictions ($\hat{y} \leq y$)

• For example, when predicting the chance of a person getting infected, we can afford overpredicting a little, i.e. providing them with a harmless vaccine. But underpredicting the risk will cause the person not protected by vaccine.

Problem 6 [10 points]

Write down the expression for the gradient of the loss ℓ_α with respect to the parameters of a linear predictor $\hat{y}(x) = w \cdot x$.

$$\hat{y}(x) = w \cdot x$$

$$L(w) = \frac{1}{N} \sum_{i=1}^N (y_i - w \cdot x_i)^2$$

$$\frac{\partial L}{\partial w_0} = \begin{cases} \frac{\partial (\alpha (y - w_0 - w_1 x_1 - \dots - w_p x_p)^2)}{\partial w_0} & , y \geq \hat{y} \\ \frac{\partial (y - w_0 - w_1 x_1 - \dots - w_p x_p)^2}{\partial w_0} & , y < \hat{y} \end{cases}$$

$$= \begin{cases} 2\alpha (y - w \cdot x) & , y \geq \hat{y} \\ 2(y - w \cdot x) & , y < \hat{y} \end{cases}$$

$$\frac{\partial L}{\partial w_j}, j \neq 0 = \begin{cases} 2\alpha (y - w \cdot x) x_j & , y \geq \hat{y} \\ 2(y - w \cdot x) x_j & , y < \hat{y} \end{cases}$$

• Therefore,

$$\nabla L = \begin{cases} \begin{pmatrix} 2\alpha (y - w \cdot x) \\ 2\alpha (y - w \cdot x) x_1 \\ \vdots \\ 2\alpha (y - w \cdot x) x_p \end{pmatrix} & , y \geq \hat{y} \\ \begin{pmatrix} 2(y - w \cdot x) \\ 2(y - w \cdot x) x_1 \\ \vdots \\ 2(y - w \cdot x) x_p \end{pmatrix} & , y < \hat{y} \end{cases}$$

Problem 8 [20 bonus points]

In this problem, we will flesh out the analysis of the convergence rate of the Perceptron algorithm. It's a bonus problem: it's optional, and if you solve it you may get more than 100 points for the homework assignment.

Let's make the following assumptions. Suppose that our data lies in the unit ball and is linearly separable by a *unit-norm, origin-containing* hyperplane \mathbf{w}^* with margin γ . This means that for every training data point \mathbf{x}_i and corresponding label y_i , we have:

$$\begin{aligned} \|\mathbf{x}_i\| &\leq 1 \\ y_i \mathbf{w}^* \cdot \mathbf{x}_i &\geq \gamma \end{aligned}$$

Note that we can always ensure (by manipulating the data) that $\|\mathbf{x}_i\| \leq 1$, and $b = 0$ (think why this is true).

Let M denote the number of mistakes Perceptron makes. Recall from lecture that we want to show that $M \leq 1/\gamma^2$.

- Without loss of generality, we'll assume that all our datapoints have label +1. Explain why this is a valid assumption.

• We can transform the dataset as

$$\mathbf{x}_i = \begin{cases} \mathbf{x}_i, & \text{if } y_i = 1 \\ -\mathbf{x}_i, & \text{if } y_i = -1 \end{cases}$$

Then, we would still have the same prediction, because

$$\text{sign}(-\mathbf{x}_i \cdot \mathbf{w}) = -\text{sign}(\mathbf{x}_i \cdot \mathbf{w})$$

- In showing the convergence of Perceptron, we'll need two lemmas. The first roughly says, "we make decent progress towards a good solution after each mistake." The second roughly says, "the norm of our solution vector isn't too big." Think about why these might be good things to show (no need to write anything here).

- OK, let's formalize the statement of the first lemma. Let \mathbf{w}_i denote our weight vector after we've made i mistakes. Show that:

$$\mathbf{w}_{i+1} \cdot \mathbf{w}^* - \mathbf{w}_i \cdot \mathbf{w}^* \geq \gamma$$

Then, think about how this mathematical statement fits with the English description above.

• After making the $i+1^{\text{th}}$ mistake,

$$\mathbf{w}_{i+1} = \mathbf{w}_i + y_i \mathbf{x}_i$$

• now we get

$$\mathbf{w}^* (\mathbf{w}_{i+1} - \mathbf{w}_i) = \mathbf{w}^* y_i \mathbf{x}_i$$

• as $y_i \mathbf{w}^* \mathbf{x}_i > \gamma$, we obtained the desired inequality.

- Before we move on to the next lemma, let's try a naive bound to capture the idea of the second lemma. Using the triangle inequality, bound $\|\mathbf{w}_M\|$.

Hint: The triangle inequality states that for any two vectors \mathbf{a} and \mathbf{b} :

$$\|\mathbf{a} + \mathbf{b}\| \leq \|\mathbf{a}\| + \|\mathbf{b}\|$$

• Let $\alpha(i)$ be the index of $\mathbf{x}_{\alpha(i)}$ at the i^{th} mistake.

$$\text{then } \|\mathbf{w}_M\| = \left\| \mathbf{w}_0 + \sum_{i=1}^M y_{\alpha(i)} \cdot \mathbf{x}_{\alpha(i)} \right\| \quad \downarrow y_i = 1$$

$$= \left\| \sum_{i=1}^M \mathbf{x}_{\alpha(i)} \right\|$$

$$\leq \underbrace{\|\mathbf{x}_{\alpha(1)}\|}_{\leq 1} + \dots + \underbrace{\|\mathbf{x}_{\alpha(M)}\|}_{\leq 1}$$

$$\leq M$$

5. Let's now strengthen this. In general, the following "squared triangle inequality" doesn't hold (if you're unconvinced, find some counterexample):

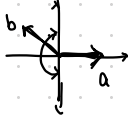
$$\|a+b\|^2 \leq \|a\|^2 + \|b\|^2$$

However, this is true for the iterates of our Perceptron algorithm! To see this, show that if $a \cdot b \leq 0$, then the above "squared triangle inequality" is actually true. Use this to show that:

$$\|w_M\| \leq \sqrt{M}$$

Contrast this with the bound we saw in the previous part.

As $a \cdot b = \|a\| \cdot \|b\| \cdot \cos \theta \leq 0$,
It must be that $\cos \theta \leq 0$,
meaning $\theta \in [\frac{\pi}{2}, \frac{3\pi}{2}]$.



- Therefore, a and b are pointing in "opposing directions". Thus, $a+b$ will be "shorter" than both a and b :

$$\begin{aligned} \|a+b\| &\leq \|a\| \text{ and } \|a+b\| \leq \|b\|, \\ \text{so } \|a+b\|^2 &\leq \|a\|^2 \text{ and } \|a+b\|^2 \leq \|b\|^2, \\ \text{so } \|a+b\|^2 &\leq \|a\|^2 + \|b\|^2. \end{aligned}$$

Last question:

$$\|w_M\| = \|x_{d(1)} + \dots + x_{d(M)}\|,$$

now

$$\begin{aligned} \|x_{d(1)} + \dots + x_{d(M)}\|^2 &\leq \underbrace{\|x_{d(1)}\|^2}_{\leq 1} + \dots + \underbrace{\|x_{d(M)}\|^2}_{\leq 1} \\ &\leq M \end{aligned}$$

$$\text{Thus, } \|w_M\| \leq \sqrt{M}.$$

6. Let's combine everything. Get an upper and lower bound on the quantity $w^* \cdot w_M$, simplify appropriately, and conclude that $M \leq 1/\gamma^2$. Here, it should become clear why the bound from part (4) is not sufficient to tell us anything interesting.
Hint: You might find it useful to know the *Cauchy-Schwarz inequality*, which states that for any two vectors a and b :

$$a \cdot b \leq \|a\| \cdot \|b\|$$

The proof of this is straightforward:

$$\begin{aligned} a \cdot b &= \|a\| \cdot \|b\| \cdot \cos(\theta_{a,b}) \\ &\leq \|a\| \cdot \|b\| \end{aligned}$$

- Use telescoping:

$$\begin{aligned} w_M \cdot w^* &= \overbrace{w_M w^* - w_{M-1} w^*}^{\geq \gamma} \\ &\quad + \underbrace{w_{M-1} w^* - w_{M-2} w^*}_{\vdots} \\ &\quad + \underbrace{w_1 w^* - w_0 w^*}_{} \end{aligned} \quad \left. \vphantom{\begin{aligned} w_M w^* &= \end{aligned}} \right\} M$$

$$\geq \boxed{M \cdot \gamma} \quad \text{lower bound}$$

- In addition, $\overset{=1, \text{ from class slide}}{\|w^*\|}$

$$\begin{aligned} w_M \cdot w^* &\leq \|w_M\| \cdot \|w^*\| \\ &= \|w_M\| \\ &\leq \boxed{\sqrt{M}} \quad \text{upper bound} \end{aligned}$$

- Combine:

$$\begin{aligned} M \cdot \gamma &\leq \sqrt{M} \\ M &\leq \frac{1}{\gamma^2} \end{aligned}$$

Problem 9 [10 points]

Write down a modified algorithm for the multi-way Perceptron, as pseudo-code. Explain the motivation for your proposed modification.

pseudo-code:

while epoch < max_epoch:

for i in 1...N:

each-class-score = $x_i \cdot W + b$

predicted_class = each-class-score.argmax()

if predicted_class != actual_class:

Weight vector for predicted class = $Lr * x_i$

Weight vector for actual class = $Lr * x_i$

$b[\text{predicted_class}] = Lr$

$b[\text{actual_class}] = Lr$

if number of wrong predicts < stop_threshold:

return

else:

epoch += 1 # continue to next epoch

$Lr = Lr * \text{decay}$

★ for Lr in 1 0.5 0.3:

decay = 0.95 (call it ξ)

⇒ results same:

• because: take learning rate η_1 and η_2 ,

say $\eta_2 = K \eta_1$, then

$$W_2 = \sum_{e=1}^T \xi^{(e)} \eta_2 \cdot y_{(e)} \cdot x_{(e)}$$

$$= \underbrace{\xi^T \cdot \eta_2}_{\text{a constant}} \cdot \sum_{e=1}^T y_{(e)} \cdot x_{(e)}$$

so W_2 is proportional to W_1 .

Thus, we need to change the decay rate.

$$x_i \cdot W + b = \begin{bmatrix} x_{i1} & \dots & x_{in} \end{bmatrix} \begin{bmatrix} w_1 \\ \vdots \\ w_n \end{bmatrix} + [b_0 \dots b_n]$$

\uparrow
 each-class-score

Motivation:

• old_score = $x_i W_c + b_c$ predicted class

new_score = $x_i (W_c - (Lr \cdot x_i^T)) + b_c - Lr$

= $x_i W_c - Lr x_i \cdot x_i^T + b_c - Lr$

= old_score - $Lr (||x_i||^2 + 1)$

< old_score

Therefore, the update will make

new_score < old_score.

• Similarly, for the actual class:

new_score = $x_i (W_c + (Lr \cdot x_i^T)) + b_c + Lr$

= old_score + $Lr (||x_i||^2 + 1)$

> old_score.

which means that after the update, the actual class will be more likely to chosen.