Midterm Review

3 Linear Regression

① Replacing least squares error with sum of absolute values as error:

$$\sum_{i=1}^{n} (y_i - \hat{y}_i)^2 \implies \sum_{i=1}^{n} |y_i - \hat{y}_i|$$

where $\hat{y}_i = w^T x_i + b$

Advantage : In general, the sum of absolute error doesn't penalize outliers as much as MSE does (due to the square). Thus, this error function is more robust to outliers.
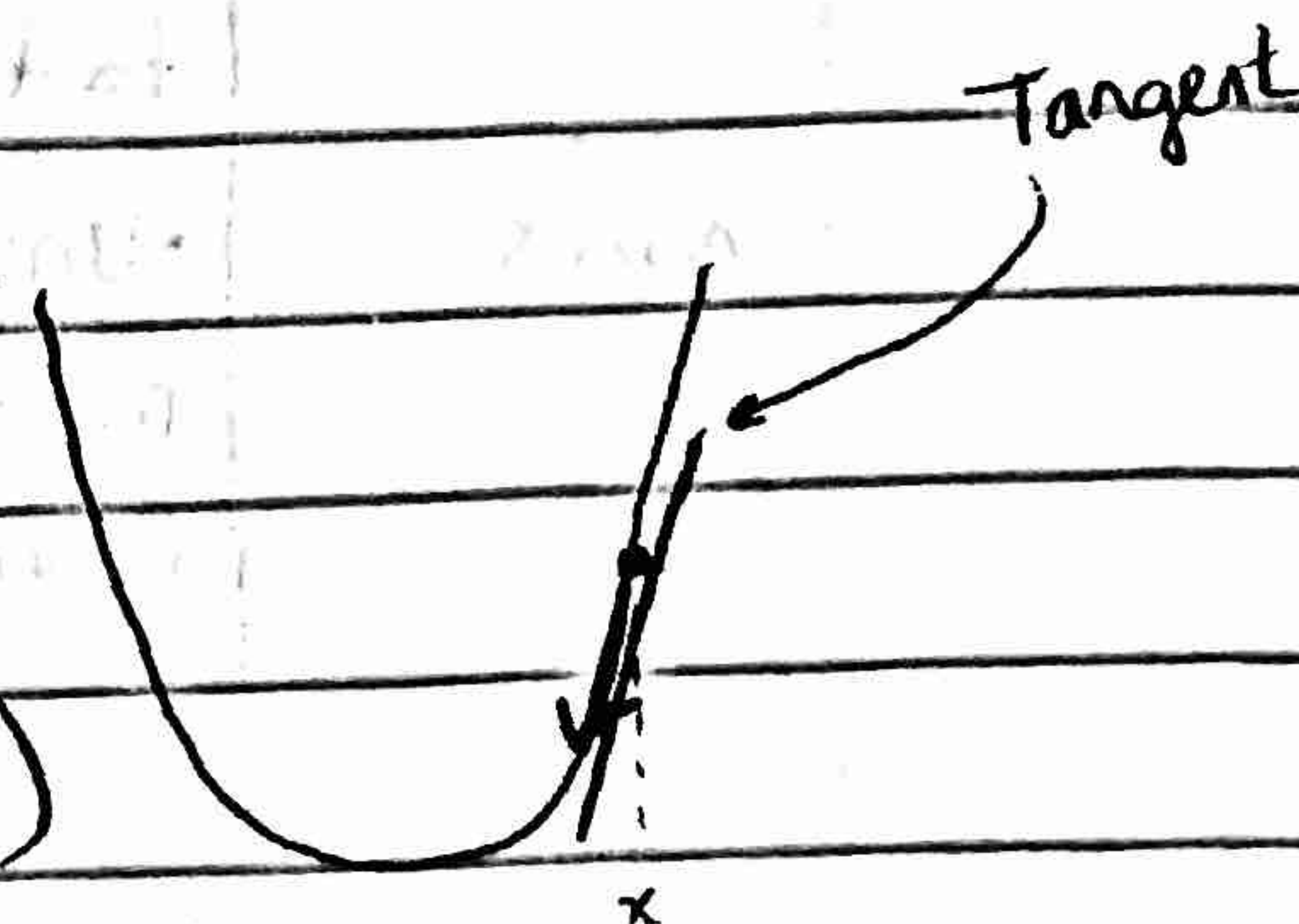
Disadvantage: If data has samples that appear like outliers (but are not), this error function won't penalize a poor fit as much as desired.

② SGD on least squares objective (batch size = 1)
$$\hookrightarrow \mathcal{L}(y, \hat{y}) = (y - \hat{y})^2 = (y - w^T x + b)^2$$

- Batch size $\Rightarrow$ how many samples do we train on & collect loss before updating weights?
- Update weights via gradient descent (SGD in this case, where samples are randomized & batch size = 1)

Gradient of $\mathcal{L}(y, \hat{y})$ – partial derivative of loss. Subtract from current weight parameters (we want to go opposite gradient – ball rolling down hill)


Tangent

**Gradients:**

$$\ell(y,\hat{y}) = (y-\hat{y})^2 = (y - w^T x + b)^2$$

$$\frac{\partial \ell}{\partial w} = 2(y-\hat{y})x \qquad \Bigg\}\ \text{Chain rule!}$$

$$\frac{\partial \ell}{\partial b} = 2(y-\hat{y})$$

**Code:**

```
for i=1 to n-epochs
    for j = 1 to n          # shuffled!
        ŷⱼ = wᵀxⱼ + b        # make prediction
        diff = yⱼ - ŷⱼ

        w = w - 2Ɛ diff x  ⎫  # update
        b = b - 2Ɛ diff    ⎭
```

where  n-epochs = # of times we iterate over training set

n = # of training examples  (shuffled for SGD)

Ɛ = learning rate

**4 Classification**

| ① | Decision Tree | LDA |
|---|---|---|
| Pros | • Interpretable, fast to test points | • Simple computations, closed form |
| cons | • Unstable to new data, may require restructuring (complex) | • Linear decision boundary (not suited for all problems) |

*I replaced the a in the question with x

② Derivative of $\sigma(x) = \dfrac{1}{1+e^{-x}} = (1+e^{-x})^{-1}$

$$\sigma'(x) = -1 \cdot (1+e^{-x})^{-2} \cdot -e^{-x} \qquad \text{Chain rule}$$

$$= \frac{e^{-x}}{(1+e^{-x})^2} \qquad \text{Consolidate}$$

$$= \frac{1}{1+e^{-x}} \cdot \frac{e^{-x}}{1+e^{-x}} \qquad \text{Separate terms}$$

$$= \frac{1}{1+e^{-x}} \cdot \frac{(1+e^{-x})-1}{1+e^{-x}} \qquad \begin{array}{l}\text{Add } 1/\text{subtract 1 to} \\ \text{numerator of second term}\end{array}$$

$$= \frac{1}{1+e^{-x}} \cdot \left(\frac{1+e^{-x}}{1+e^{-x}} - \frac{1}{1+e^{-x}}\right) \qquad \text{Factor out}$$

$$= \frac{1}{1+e^{-x}} \cdot \left(1 - \frac{1}{1+e^{-x}}\right) \qquad \text{Simplify}$$

$$= \sigma(x) \cdot (1 - \sigma(x))$$

③ Fisher Linear Discriminant weights $(\hat{w}_{Fisher})$ for a 2 class problem

We know from class that $\hat{w}_{Fisher}$ can be computed by the following —

$$\hat{w}_{Fisher} = \underbrace{S_w^{-1}}_{\substack{\text{within} \\ \text{class variance}}} \underbrace{(\mu_2 - \mu_1)}_{\substack{\text{difference of} \\ \text{means}}}$$

Where $S_w = \sum\limits_{c} \sum\limits_{i \in C} (x_i - \mu_c)(x_i - \mu_c)^T$

To compute: $\mu_1, \mu_2, \Sigma_1, \Sigma_2, \Sigma_w (S_w)$

$\boxed{\mu_1}$ 1x2

$$\nu_1 = \begin{bmatrix} \dfrac{0 + (-1) + 1}{3} \\ \dfrac{0 + 0 + 1}{3} \end{bmatrix} = \begin{bmatrix} 0 \\ 1/3 \end{bmatrix}$$

$\boxed{\mu_2}$ 1x2

$$\mu_2 = \begin{bmatrix} \dfrac{10 + 11}{2} \\ \dfrac{10 + 10}{2} \end{bmatrix} = \begin{bmatrix} 21/2 \\ 10 \end{bmatrix}$$

$\boxed{\Sigma_1}$ 2x2

$$\Sigma_1 = \left( \begin{bmatrix} 0 \\ 0 \end{bmatrix} - \begin{bmatrix} 0 \\ 1/3 \end{bmatrix} \right) \left( \begin{bmatrix} 0 \\ 0 \end{bmatrix} - \begin{bmatrix} 0 \\ 1/3 \end{bmatrix} \right)^T + \left( \begin{bmatrix} -1 \\ 0 \end{bmatrix} - \begin{bmatrix} 0 \\ 1/3 \end{bmatrix} \right) \left( \begin{bmatrix} -1 \\ 0 \end{bmatrix} - \begin{bmatrix} 0 \\ 1/3 \end{bmatrix} \right)^T$$

$$+ \left( \begin{bmatrix} 1 \\ 1 \end{bmatrix} - \begin{bmatrix} 0 \\ 1/3 \end{bmatrix} \right) \left( \begin{bmatrix} 1 \\ 1 \end{bmatrix} - \begin{bmatrix} 0 \\ 1/3 \end{bmatrix} \right)^T$$

$$= \begin{bmatrix} 0 \\ -1/3 \end{bmatrix} \begin{bmatrix} 0 & -1/3 \end{bmatrix} + \begin{bmatrix} -1 \\ -1/3 \end{bmatrix} \begin{bmatrix} -1 & -1/3 \end{bmatrix} + \begin{bmatrix} 1 \\ 2/3 \end{bmatrix} \begin{bmatrix} 1 & 2/3 \end{bmatrix}$$

$$= \begin{bmatrix} 0 & 0 \\ 0 & 1/9 \end{bmatrix} + \begin{bmatrix} 1 & 1/3 \\ 1/3 & 1/9 \end{bmatrix} + \begin{bmatrix} 1 & 2/3 \\ 2/3 & 4/9 \end{bmatrix} = \begin{bmatrix} 2 & 1 \\ 1 & 5/9 \end{bmatrix} = \Sigma_1$$

$\boxed{\Sigma_2} \; 2\times 2$

$$\Sigma_2 = \left( \begin{bmatrix} 10 \\ 10 \end{bmatrix} - \begin{bmatrix} 2\frac{1}{2} \\ 10 \end{bmatrix} \right) \left( \begin{bmatrix} 10 \\ 10 \end{bmatrix} - \begin{bmatrix} 2\frac{1}{2} \\ 10 \end{bmatrix} \right)^T + \left( \begin{bmatrix} 11 \\ 10 \end{bmatrix} - \begin{bmatrix} 2\frac{1}{2} \\ 10 \end{bmatrix} \right) \left( \begin{bmatrix} 11 \\ 10 \end{bmatrix} - \begin{bmatrix} 2\frac{1}{2} \\ 10 \end{bmatrix} \right)^T$$

$$= \begin{bmatrix} -\frac{1}{2} \\ 0 \end{bmatrix} \begin{bmatrix} -\frac{1}{2} & 0 \end{bmatrix} + \begin{bmatrix} \frac{1}{2} \\ 0 \end{bmatrix} \begin{bmatrix} \frac{1}{2} & 0 \end{bmatrix}$$

$$= \begin{bmatrix} \frac{1}{4} & 0 \\ 0 & 0 \end{bmatrix} + \begin{bmatrix} \frac{1}{4} & 0 \\ 0 & 0 \end{bmatrix} = \begin{bmatrix} \frac{1}{2} & 0 \\ 0 & 0 \end{bmatrix} = \Sigma_2$$

$\boxed{\Sigma_{SW}} \; 2\times 2 = \Sigma_1 + \Sigma_2$

$$= \begin{bmatrix} 2 & 1 \\ 1 & 5/4 \end{bmatrix} + \begin{bmatrix} \frac{1}{2} & 0 \\ 0 & 0 \end{bmatrix} = \begin{bmatrix} 2\frac{1}{2} & 1 \\ 1 & 5/4 \end{bmatrix}$$

$\hat{W}_{Fisher} = S_W^{-1} (\mu_2 - \mu_1)$

Inverse of $2\times 2 \Rightarrow \begin{bmatrix} a & b \\ c & d \end{bmatrix}^{-1} = \dfrac{1}{ad-bc} \begin{bmatrix} d & -b \\ -c & a \end{bmatrix}$

$$\begin{bmatrix} 2\frac{1}{2} & 1 \\ 1 & 5/4 \end{bmatrix}^{-1} = \frac{1}{\frac{5}{2}\cdot\frac{5}{4}-1} \begin{bmatrix} 5/4 & -1 \\ -1 & 5/2 \end{bmatrix} = \frac{18}{7} \begin{bmatrix} 5/4 & -1 \\ -1 & 5/2 \end{bmatrix}$$

(Not simplified)

$$\hat{W}_{Fisher} = \frac{18}{7} \begin{bmatrix} 5/4 & -1 \\ -1 & 5/2 \end{bmatrix} \left( \begin{bmatrix} \frac{21}{2} \\ 10 \end{bmatrix} - \begin{bmatrix} 0 \\ 1/3 \end{bmatrix} \right)$$

## 5 General

When would you use a validation set in addition to training & testing sets?

A validation set is a portion of your data not used for training or testing sets, used to tune hyperparameters

In Fisher's LDA, the decision rule is

$$L_{Fisher}(\vec{x}) = \mathbb{I}\{\hat{w}_{Fisher}^T \vec{x} > \gamma\}$$

where $\gamma$ acts as the bias (more literally, $\gamma$ = "negative bias", since $\hat{w}_{Fisher}^T x + b = 0$

$\hat{w}_{Fisher}^T x = -b$, so $-b = \gamma$).

The threshold $\gamma$, the point where we determine whether a test point will be classified as class 1 or 2, is determined by cross-validation using a validation set.