

Week 2

9/2/2025

- ↳ Quiz!
- ↳ How many read the notes? Important!
- ↳ Encourage practice, problem sets due every week
 - First 2 problems still 9/19
 - In general, due Mondays 11:59 pm
- ↳ Office hours change:
Monday + Wednesday 12:30 to 2
- ↳ Post questions, come to office hours :-)

Last week was warm up for machine learning!

Review

What about when M is not square?

$M \in \mathbb{R}^{d \times d}$

$$M = \sum_{i=1}^r \underline{v}_i \lambda_i \underline{w}_i^T$$

$\lambda_1, \dots, \lambda_r \in \mathbb{R}$ eigenvalues

$\underline{v}_1, \dots, \underline{v}_r \in \mathbb{R}^d$ right eigenvectors

$\underline{w}_1, \dots, \underline{w}_r \in \mathbb{R}^d$ left eigenvectors

Use: Page Rank where we powered up matrix multiplication!

$$\underline{A} \underline{v}_i = \lambda_i \underline{v}_i$$

$$\underline{w}_i^T \underline{A} = \lambda_i \underline{w}_i^T$$

$$\underline{v}_i^T \underline{v}_j = \mathbb{I}[i=j] = \underline{w}_i^T \underline{w}_j$$

$$\underline{w}_j^T \underline{A} \underline{v}_i = \underline{w}_j^T \lambda_i \underline{v}_i$$

\Leftrightarrow

$$\underline{w}_j^T \underline{v}_i = \mathbb{I}[i=j]$$

$$\lambda_j \underline{w}_j^T \underline{v}_i = \lambda_i \underline{w}_j^T \underline{v}_i$$

$$(\lambda_j - \lambda_i) \underline{w}_j^T \underline{v}_i = 0$$

Supervised Learning

Problems like:

- predicting temperature
- identifying objects in an image
- generating next word

Labelled data

$$\underline{x} \in \mathbb{R}^d \quad y \in \mathbb{R}$$

n points...

$$(\underline{x}^{(1)}, y^{(1)}), \dots, (\underline{x}^{(n)}, y^{(n)})$$

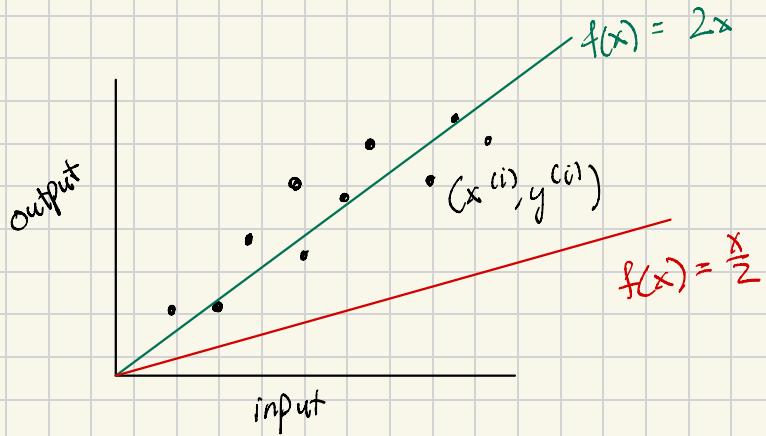
Goal: Find function $f: \mathbb{R}^d \rightarrow \mathbb{R}$
so that $f(\underline{x}^{(i)}) \approx y^{(i)} \quad \forall i$

Empirical risk minimization:

- ① Function class \mathcal{F} from which to select f
- ② Loss to measure how well f fits data
- ③ Optimizer, method to select f with low loss

Univariate Linear Regression

$$x^{(1)}, \dots, x^{(n)} \in \mathbb{R}$$



① Function class:

$$f(x) = \omega x \quad \text{for } \omega \in \mathbb{R}$$

② Loss

Attempt #1: $f(x) - y$

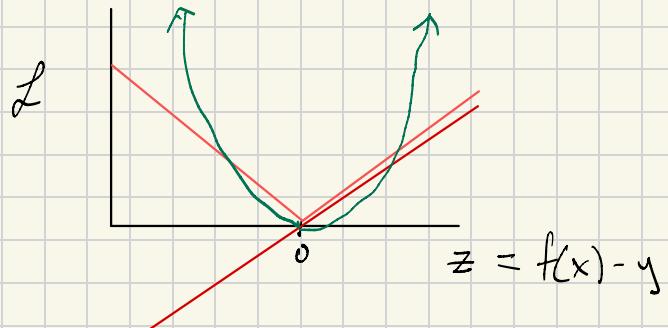
$\because y > f(x) \Rightarrow$ very negative

Attempt #2: $|f(x) - y|$

\because not differentiable at 0

Attempt #3: $(f(x) - y)^2$

\therefore far gets penalized more



Mean Squared Error (MSE) Loss

$$\mathcal{L}(\omega) = \frac{1}{n} \sum_{i=1}^n [f_\omega(x^{(i)}) - y^{(i)}]^2$$

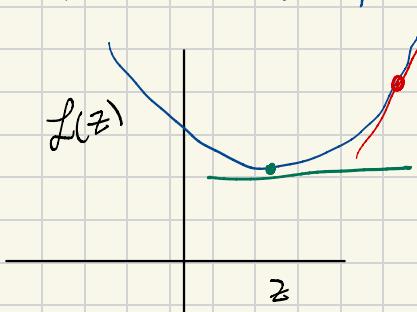
③ Exact Optimization

Key insight: Squared loss is convex

↳ differentiable

↳ single minima

Q: How do we find minima?



no improvement in any direction!

$$\begin{aligned}\frac{\partial}{\partial \omega} [\mathcal{L}(\omega)] &= \frac{1}{n} \sum_{i=1}^n \frac{\partial}{\partial \omega} [w x^{(i)} - y^{(i)}]^2 \\ &= \frac{1}{n} \sum_{i=1}^n 2[w x^{(i)} - y^{(i)}] \frac{\partial}{\partial \omega} (w x^{(i)} - y^{(i)}) \\ &\stackrel{\text{set}}{=} \frac{1}{n} \sum_{i=1}^n 2[w x^{(i)} - y^{(i)}] x^{(i)} \\ &= 0\end{aligned}$$

$$\begin{aligned}\stackrel{?}{=} \frac{1}{n} \sum_{i=1}^n w^* (x^{(i)})^2 &= \frac{2}{n} \sum_{i=1}^n y^{(i)} x^{(i)} \\ w^* &= \frac{\sum_{i=1}^n y^{(i)} x^{(i)}}{\sum_{i=1}^n [x^{(i)}]^2}\end{aligned}$$

In general, $\underline{x}^{(i)} \in \mathbb{R}^d$

Again, but with some linear algebra

Multivariate Linear Regression

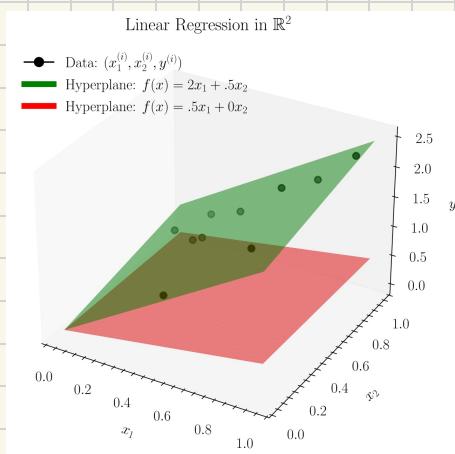
$(\underline{x}^{(1)}, y^{(1)}), \dots, (\underline{x}^{(n)}, y^{(n)})$ $\underline{x}^{(i)} \in \mathbb{R}^d$

① Function Class

Let weights $\underline{w} \in \mathbb{R}^d$

$$f(\underline{x}) = \langle \underline{x}, \underline{w} \rangle = \underline{x}^\top \underline{w} = \sum_{k=1}^d w_k x_k$$

coefficient
for each
feature



$$f(\underline{x}) = [2 \quad .5] \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$$

$$f(\underline{x}) = [.5 \quad 0] \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$$

② MSE

$$\mathcal{L}(\underline{w}) = \frac{1}{n} \sum_{i=1}^n (\langle \underline{x}^{(i)}, \underline{w} \rangle - y^{(i)})^2$$

$$\underline{X} \in \mathbb{R}^{n \times d}$$

$$\begin{bmatrix} \underline{x}^{(1)\top} \\ \underline{x}^{(2)\top} \\ \vdots \\ \underline{x}^{(n)\top} \end{bmatrix}$$

$$\underline{y} \in \mathbb{R}^n$$

$$\begin{bmatrix} y^{(1)} \\ y^{(2)} \\ \vdots \\ y^{(n)} \end{bmatrix}$$

$$\frac{1}{n} \|\underline{X}\underline{w} - \underline{y}\|_2^2 = \mathcal{L}(\underline{w})$$

$$\frac{1}{n} \left\| \begin{bmatrix} \langle \underline{x}^{(1)}, \underline{w} \rangle \\ \langle \underline{x}^{(2)}, \underline{w} \rangle \\ \vdots \\ \langle \underline{x}^{(n)}, \underline{w} \rangle \end{bmatrix} - \begin{bmatrix} y^{(1)} \\ y^{(2)} \\ \vdots \\ y^{(n)} \end{bmatrix} \right\|_2^2$$

(3)

Optimizer

Key insight:

$\underline{\omega}^*$ when no improvement
in any direction

$$\nabla_{\underline{\omega}} \mathcal{L}(\underline{\omega}) = \begin{bmatrix} \frac{\partial \mathcal{L}}{\partial \omega_1} \\ \frac{\partial \mathcal{L}}{\partial \omega_2} \\ \vdots \\ \frac{\partial \mathcal{L}}{\partial \omega_d} \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix}$$

$$\underline{e}_i = \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 0 \\ 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix} \leftarrow i^{\text{th position}}$$

$$\frac{\partial \mathcal{L}(\underline{\omega})}{\partial \omega_i} = \lim_{\Delta \rightarrow 0} \frac{\mathcal{L}(\underline{\omega} + \Delta \underline{e}_i) - \mathcal{L}(\underline{\omega})}{\Delta}$$

$$= \lim_{\Delta \rightarrow 0} \frac{\|\underline{X}(\underline{\omega} + \Delta \underline{e}_i) - \underline{y}\|_2^2 - \|\underline{X}\underline{\omega} - \underline{y}\|_2^2}{\Delta} = (*)$$

$$\begin{aligned} \|\underline{a} + \underline{b}\|_2^2 &= \sum_{i=1}^d (a_i + b_i)^2 = (\underline{a} + \underline{b})^T (\underline{a} + \underline{b}) \\ &= \underline{a}^T \underline{a} + \underline{a}^T \underline{b} + \underline{b}^T \underline{a} + \underline{b}^T \underline{b} \\ &= \|\underline{a}\|_2^2 + 2 \langle \underline{a}, \underline{b} \rangle + \|\underline{b}\|_2^2 \end{aligned}$$

$$\underline{a} = \underline{X}\underline{\omega} - \underline{y} \quad \underline{b} = \Delta \underline{X}\underline{e}_i$$

$$(*) = \lim_{\Delta \rightarrow 0} \frac{\|\underline{X}\underline{\omega} - \underline{y}\|_2^2 + 2 \langle \underline{X}\underline{\omega} - \underline{y}, \Delta \underline{X}\underline{e}_i \rangle + \|\Delta \underline{X}\underline{e}_i\|_2^2 - \|\underline{X}\underline{\omega} - \underline{y}\|_2^2}{\Delta}$$

$$= \lim_{\Delta \rightarrow 0} \frac{2 \Delta \langle \underline{X}\underline{\omega} - \underline{y}, \underline{X}\underline{e}_i \rangle + \Delta^2 \|\underline{X}\underline{e}_i\|_2^2}{\Delta}$$

$$= 2 \langle \underline{X}\underline{\omega} - \underline{y}, \underline{X}\underline{e}_i \rangle$$

$$\underline{X} \underline{e}_i = \underline{\underline{X}}_i$$

$$\begin{bmatrix} -\underline{\underline{x}}^{(1)\top} \\ \vdots \\ -\underline{\underline{x}}^{(n)\top} \end{bmatrix} \begin{bmatrix} 0 \\ \vdots \\ 0 \end{bmatrix}$$

$$\nabla_{\underline{\omega}} \mathcal{L}(\underline{\omega}) = \frac{1}{2} \begin{bmatrix} \underline{\underline{x}}_1^\top (\underline{\underline{x}} \underline{\omega} - \underline{y}) \\ \vdots \\ \underline{\underline{x}}_i^\top (\underline{\underline{x}} \underline{\omega} - \underline{y}) \end{bmatrix} = 2 \underline{\underline{X}}^\top (\underline{\underline{X}} \underline{\omega} - \underline{y})$$

$$\nabla_{\underline{\omega}} \mathcal{L}(\underline{\omega}) = 0 = 2 \underline{\underline{X}}^\top (\underline{\underline{X}} \underline{\omega^*} - \underline{y})$$

$$\underline{\underline{X}}^\top \underline{\underline{X}} \underline{\omega^*} = \underline{\underline{X}}^\top \underline{y}$$

$$\underline{\omega^*} = (\underline{\underline{X}}^\top \underline{\underline{X}})^+ \underline{\underline{X}}^\top \underline{y}$$