

## Reminders

- ↳ Quiz
- ↳ Pset
- ↳ Office Hours
- ↳ Discord
- ↳ Events
  - Conference Oct 3-5
  - Redistricting Sep 16 Lunch @ Ath
  - Trees Sep 16 7pm @ Mudd

## Review

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

$$y^{(i)} \in \mathbb{R}$$

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

$$\begin{bmatrix} -x^{(i)T} \end{bmatrix}$$

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

$$\begin{bmatrix} y^{(i)} \end{bmatrix}$$

Goal:  $f(\underline{x}^{(i)}) \approx y^{(i)}$

① Model Linear!  $f(\underline{x}) = \underline{\omega}^T \underline{x}$

② Loss MSE!  $\mathcal{L}(\underline{\omega}) = \sum_{i=1}^n [f(\underline{x}^{(i)}) - y^{(i)}]^2 \frac{1}{n}$

③ Optimizer Exact!  $\underline{\omega}^* = (\underline{X}^T \underline{X})^+ \underline{X}^T \underline{y}$

Why pseudo inverse?

"Invert" non-invertible matrices  
e.g.  $\underline{X} \in \mathbb{R}^{n \times d}$

$$\underline{X} = \sum_{i=1}^d \sigma_i \underline{u}_i \underline{v}_i^T$$

$$\begin{aligned} \underline{X}^+ \underline{X} &= \sum_{j=1}^d \frac{1}{\sigma_j} \underline{v}_j \underline{u}_j^T \sum_{i=1}^d \sigma_i \underline{u}_i \underline{v}_i^T \\ &= \sum_{i=1}^d \underline{v}_i \underline{v}_i^T = \underline{I}_d \end{aligned}$$

Greatest advice: "Go sit on your rock"

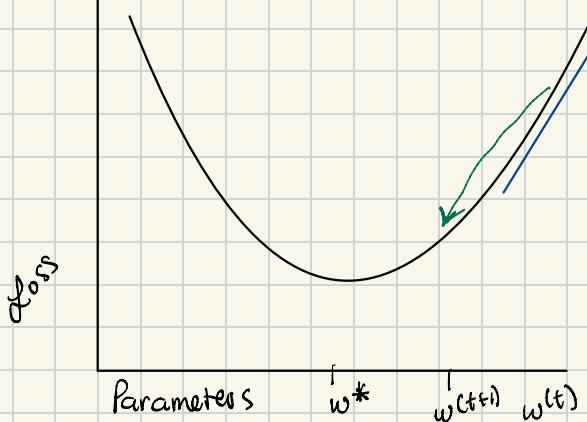
Two issues:

↳ Time to compute  $\underline{\omega}^*$

↳ What if data is not linear?

## Gradient Descent

$$\omega^{(t+1)} \in \omega^{(t)} - \alpha \nabla \mathcal{L}(\omega^{(t)})$$



Intuition: Move away from steepest ascent

$\alpha$  = "step size" or "learning rate"

## The Math (as promised)

$$[\omega \in \mathbb{R}]$$

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

$$\Rightarrow \mathcal{L}(\omega + \Delta) - \mathcal{L}(\omega) \approx \frac{\partial \mathcal{L}(\omega)}{\partial \omega} \Delta$$

Choose  $\Delta$  so  $\mathcal{L}(\omega + \Delta) - \mathcal{L}(\omega)$  is negative ...

$$\Delta = -\frac{\partial \mathcal{L}(\omega)}{\partial \omega}$$

$$w \in \mathbb{R}^d$$

$$\frac{\partial \mathcal{L}(w)}{\partial w_i} = \lim_{\Delta \rightarrow 0} \frac{\mathcal{L}(w + \Delta e_i) - \mathcal{L}(w)}{\Delta}$$

$$\Rightarrow \mathcal{L}(w + \Delta e_i) - \mathcal{L}(w) \approx \Delta \langle \nabla_w \mathcal{L}(w), e_i \rangle$$

$$\mathcal{L}(w + v) - \mathcal{L}(w) \approx \Delta \langle \nabla_w \mathcal{L}(w), v \rangle$$

$$\text{choose } \Delta v = -\alpha \nabla_w \mathcal{L}(w)$$

$$\begin{aligned} \mathcal{L}(w + v) - \mathcal{L}(w) &\approx -\alpha \|\nabla_w \mathcal{L}(w)\|_2^2 \\ &= -\alpha \|\nabla_w \mathcal{L}(w)\|_2 \|\nabla_w \mathcal{L}(w)\|_2 \end{aligned}$$

$$\text{Recall } \langle a, b \rangle = \|a\|_2 \|b\|_2 \cos(\theta)$$

$$\max_{\theta} \cos \theta = 1 \Rightarrow \text{achieve "best" update!}$$

For linear regression,

$$\nabla_w \mathcal{L}(w) = \frac{2}{n} X^T(Xw - y)$$

Time to compute:

Even Faster

n and d can be prohibitively large

Enter: STOCHASTIC gradient descent

↑  
"stokhos"

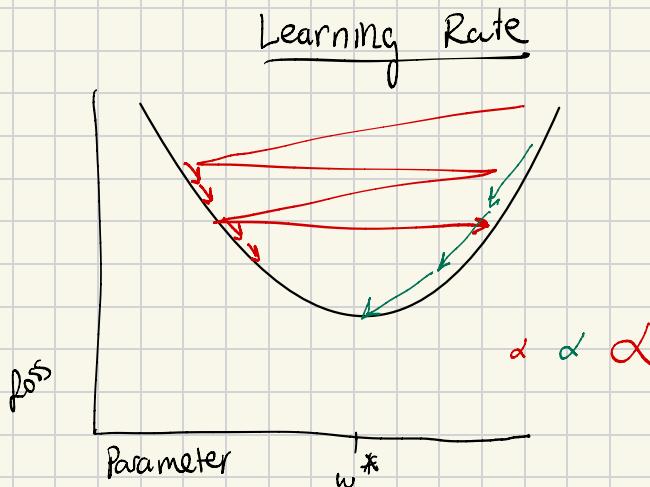
$$S \subseteq \{1, \dots, n\}$$

$$\mathcal{L}_s(\omega) = \frac{1}{|S|} \sum_{i \in S} [f(x^{(i)}) - y^{(i)}]^2$$

For linear regression,

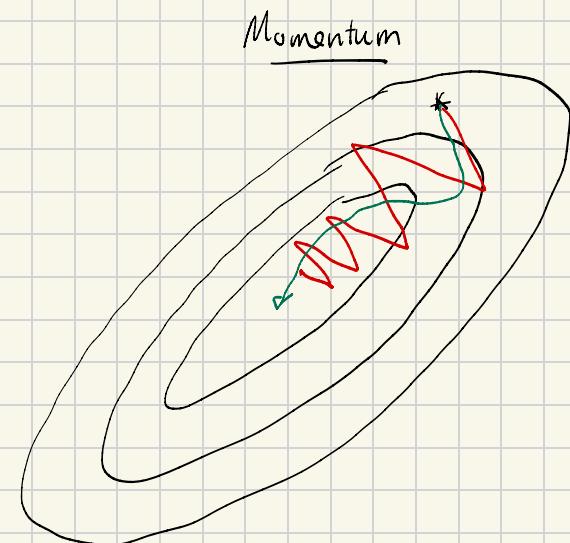
$$\nabla_{\omega} \mathcal{L}_s(\omega) = \frac{2}{|S|} x_s^T (x_s \omega - y_s)$$

Time to compute:



heuristic strategies:

- ↳ If loss unstable, decrease ( $\alpha/2$ )
- If loss slowly decreasing, increase ( $2\alpha$ )
- ↳ "scheduler" start big  $\rightarrow$  small, cycle
- ↳ Adaptive  $\alpha^{(t+1)} \leftarrow \frac{\alpha^{(t)}}{(\nabla_w \mathcal{L}(w^{(t)}))^2}$



$$m^{(t+1)} \leftarrow \beta m^{(t)} + (1-\beta) \nabla_w \mathcal{L}_s(w^{(t)})$$

$$w^{(t+1)} \leftarrow w^{(t)} - \alpha m^{(t+1)}$$