# Reminders

↳ Quiz

↳ Pset +Self Grade!

↳ Office Hours

↳ Discord

↳ Events

- Conference    Oct 3-5
- Redistricting    Sep 16   Lunch @ Ath
- Trees         Sep 16   7pm @ Mudd

# Review

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

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

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

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

① Model    Linear! $f(x) = w^T x$

② Loss    MSE! $\mathcal{L}(w) = \sum\limits_{i=1}^{n} \left[ f(x^{(i)}) - y^{(i)} \right]^2 \frac{1}{n}$

③ Optimizer  Exact! $w^* = (X^T X)^+ X^T y$

# Why pseudo inverse?

"Invert" non-invertible matrices

e.g. $X \in \mathbb{R}^{n \times d}$

$$X = \sum_{i=1}^{d} \sigma_i \, u_i \, v_i^T$$

$$X^+ X = \sum_{j=1}^{d} \frac{1}{\sigma_j} v_j \, u_j^T \sum_{i=1}^{d} \sigma_i \, u_i \, v_i^T$$
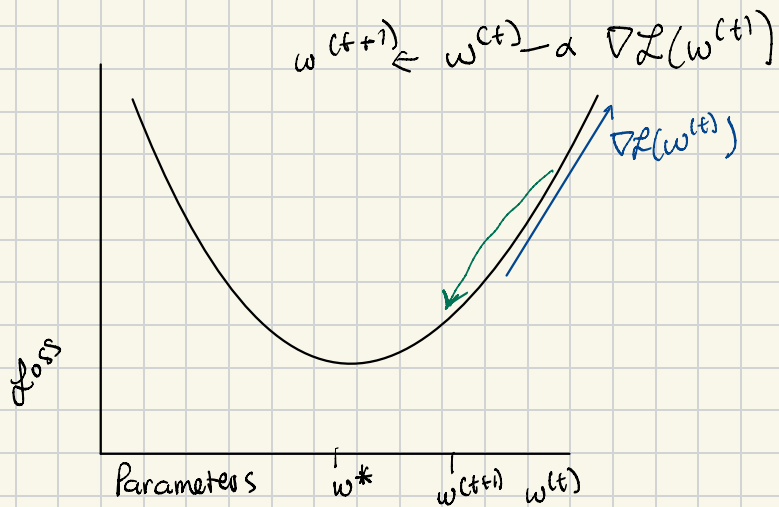
$$= \sum_{i=1}^{d} v_i \, v_i^T = I_d$$

Greatest advice: "Go sit on your rock"

Two issues:

↳ Time to compute $w^*$

↳ What if data is <u>not</u> linear?

## Gradient Descent

$$w^{(t+1)} \leftarrow w^{(t)} - \alpha \nabla \mathcal{L}(w^{(t)})$$



$\nabla \mathcal{L}(w^{(t)})$

loss

Parameters    $w^*$    $w^{(t+1)}$   $w^{(t)}$

Intuition : Move away from steepest ascent

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

## The Math   (as promised)

$\boxed{\text{GR}}$

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

$$\Rightarrow \quad \mathcal{L}(w + \Delta) - \mathcal{L}(w) \simeq \frac{\partial \mathcal{L}(w)}{\partial w} \Delta$$

Choose $\Delta$ so $\mathcal{L}(w+\Delta) - \mathcal{L}(w)$ is

negative ....        $\Delta = - \partial \frac{\mathcal{L}(w)}{\partial w}$

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

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

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

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

Choose $\quad \Delta v = -\alpha \nabla_w \mathcal{L}(w)$

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

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

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

For linear regression,

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

Time to compute:

# Reminders

↳ Notes!

↳ Self-grade due Friday

↳ Pset 3 due Monday

↳ Quiz Tuesday

  ↳ exercises + notes = fair game

↳ Events!

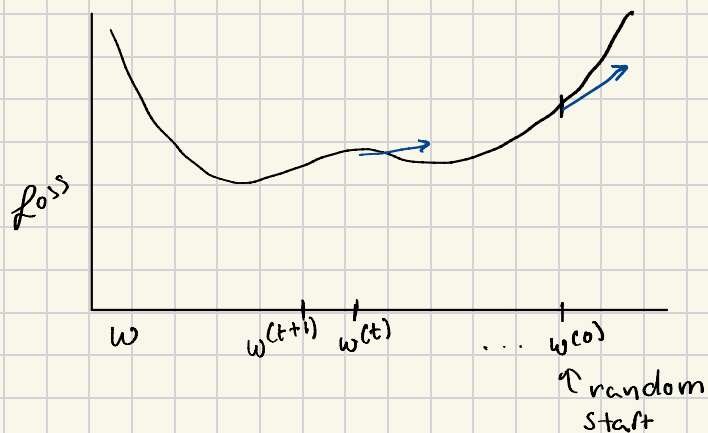↳ Please ask specific
  questions , e.g.,
  "Explain concept A"
  "Why is B true"

↳ | Sit    on    your    rock |
    ↰ most    important thing I can offer

<u>Review</u> :    Gradient   Descent



$w^{(t+1)} \leftarrow w^{(t)} - \alpha \, \nabla \mathcal{L}(w^{(t)})$

Linear  regression :

$$\nabla \mathcal{L}(w) = \underset{d \times n}{X^T} (\underset{n \times d}{X} \underset{d \times 1}{w} - \underset{n \times 1}{y})$$

Time :

<u>Even    Faster</u>

n  and  d  can  be  prohibitively large

Enter:        STOCHASTIC    gradient descent
                          ↑
                     "stokhos"
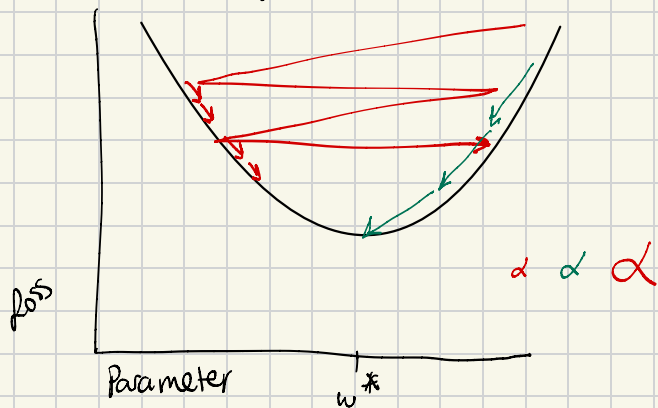
$S \subseteq \{1, \dots, n\}$        s.t.      $|S| = m$

$$\mathcal{L}_S(w) = \frac{1}{m} \sum_{i \in S} \left[ f(x^{(i)}) - y^{(i)} \right]^2$$

For    linear   regression,

$$\nabla_w \mathcal{L}_S(w) = \frac{2}{|S|} X_S^T (X_S w - y_S)$$

Time :

# Learning Rate
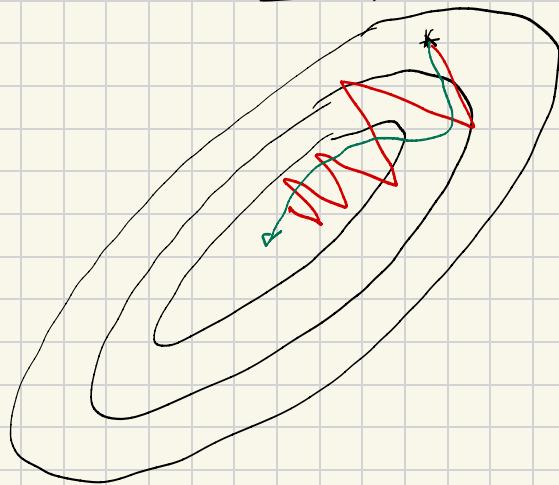


loss

Parameter    $w^*$

$\textcolor{red}{\alpha}\ \ \alpha\ \ \textcolor{red}{\alpha}$

## heuristic strategies:

↳ If loss unstable, decrease $(\alpha/2)$
   If loss slowly decreasing, increase $(2\alpha)$

↳ "scheduler"  start big → small, cycle

↳ Adaptive  $\alpha^{(t+1)} \leftarrow \dfrac{\alpha^{(t)}}{\left(\nabla_w \mathcal{L}(w^{(t)})\right)^2}$

# Momentum



$\uparrow w_2$
$\hookrightarrow w_1$

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

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

# Polynomial Regression

Q: We covered how to speed up optimization, how do we address model mis-specification?

A: Add more *expressive* features

Linear regression e.g.,

$$f(x) = w_1 x_1 + w_2 x_2$$

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

Polynomial regression e.g.,

$$f(x) = w_1 x_1 + w_2 x_2 + w_3 x_1 x_2 + w_4 x_1^2 + w_5 x_2^2$$

$$X_{aug} = \begin{bmatrix} X & Z \end{bmatrix}$$

$$\begin{bmatrix} \phantom{X} \\ \phantom{X} \end{bmatrix} \begin{bmatrix} \phantom{X} \\ \phantom{X} \end{bmatrix}$$
$$\underbrace{\phantom{XX}}_{d} \underbrace{\phantom{XX}}_{d'}$$

$$w^*_{aug} = X^+_{aug} y$$

Q: What changes in 3 step recipe?

# Expressivity

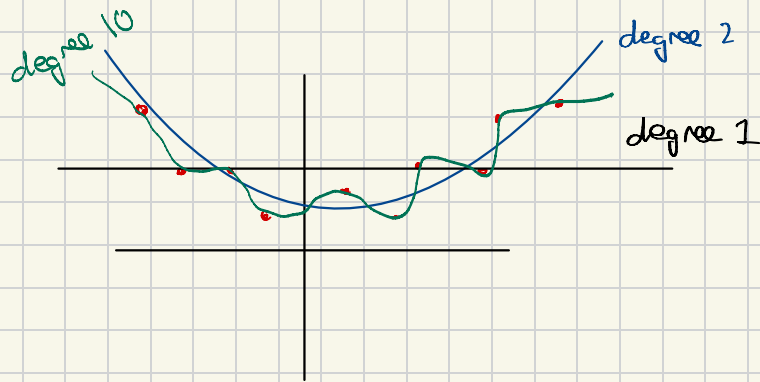$$w^* = X^+ y \qquad w^*_{aug} = X^+_{aug} y$$

Claim: More expressive

$$\min_{w \in \mathbb{R}^{d+d'}} \| X_{aug} w - y \|_2^2 \leq \min_{w \in \mathbb{R}^d} \| X w - y \|_2^2$$

Proof: choose $\hat{w} = \begin{bmatrix} w^* \\ 0 \\ 0 \\ 0 \end{bmatrix}$. Then $\| X_{aug} \hat{w} - y \|_2^2 = \| X w^* - y \|_2^2$

Since $\hat{w}$ is always a choice, $w^*_{aug}$ must be at least as good

# Which degree to choose?



degree 10

degree 2

degree 1

overfits (memorize)

$x^1, x^2, x^{10}$

not expressive enough

just right :)

# Generalization Error

| Training | Validation | Test |
|---|---|---|

Data

$\left\{\begin{array}{l}\end{array}\right.$ repeat

1. Train model

2. Check loss on validation set

3. Update hyperparameters
   (learning rate, momentum, etc)

Pernicious issue: model performs better on validation even though never "trained" on it

Solution: Use test loss <u>once</u> on final model

# Regularization

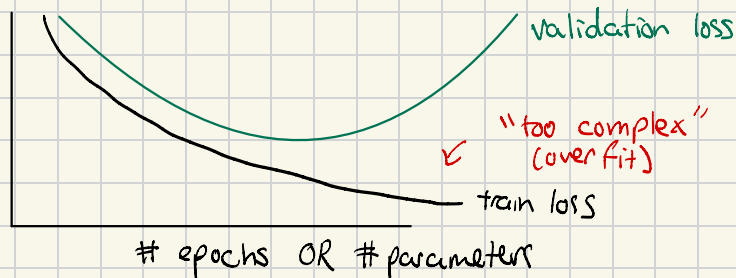**Idea:** Use simplest model

Weight magnitude $\approx$ complexity

↳ More output variation per input change

**New loss:**

$$\mathcal{L} + \lambda \|w\|_2^2$$

# Double Descent

## Traditional View:



validation loss

"too complex" (overfit)

train loss

# epochs OR # parameters

## Modern View:



validation

implicit regularization

train loss

# epochs OR # parameters