

# 10-605/805 – ML for Large Datasets

## Lecture 5: Distributed Linear Regression

Henry Chai

9/13/22

# Front Matter

- HW1 released 8/30, due 9/13 (today!) at 11:59 PM
  - Recitation 3 on 9/16 will go over HW1 solutions
- HW2 released 9/8, due 9/22 at 11:59 PM
- Mini-project details released 9/9

# Background: Big $O$ Notation

- Used to describe an algorithm's time or space (storage) complexity in terms of the input size
- Formally:
$$f(x) = O(g(x)) \iff \exists C, x_0 \text{ s.t. } f(x) \leq Cg(x) \forall x \geq x_0$$
  - $O(1)$  = constant time/space, i.e., a fixed number of operations or storage regardless of input
  - $O(\log(n))$  = logarithmic time/space
  - $O(n)$  = linear time/space
- An algorithm's time and space complexity can be different
  - Example: multiplying an  $a \times b$  matrix with an  $b \times c$  matrix takes  $O(abc)$  time ( $ac$  dot products between  $b$ -length vectors) but the result uses  $O(ac)$  storage

# Background: Empirical Risk Minimization

- A common framework for supervised learning
- Given:
  - some labelled training dataset  $\mathcal{D} = \{(\mathbf{x}^{(i)}, y^{(i)})\}_{i=1}^n$
  - a loss function  $\ell: \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}$
  - a hypothesis class or set of functions  $\mathcal{F}$

the goal is to find

$$\hat{f} = \operatorname{argmin}_{f \in \mathcal{F}} \sum_{i=1}^n \ell(f(\mathbf{x}^{(i)}), y^{(i)})$$

with the hope that

$$\mathbb{E}_{p(\mathbf{x}, y)}[\ell(f(\mathbf{x}), y)] \approx \sum_{i=1}^n \ell(f(\mathbf{x}^{(i)}), y^{(i)})$$

# Background: Empirical Risk Minimization

- A common framework for supervised learning
- Given:
  - some labelled training dataset  $\mathcal{D} = \{(\mathbf{x}^{(i)}, y^{(i)})\}_{i=1}^n$
  - a loss function  $\ell: \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}$
  - a hypothesis class or set of functions  $\mathcal{F}$

the goal is to find

$$\hat{f} = \operatorname{argmin}_{f \in \mathcal{F}} \sum_{i=1}^n \ell(f(\mathbf{x}^{(i)}), y^{(i)})$$

- Depending on the choice of  $\mathcal{F}$  and  $\ell$ , this objective function may be convex (easy to optimize) or non-convex (hard)
- Our focus will be solving this problem for large  $n$  and/or  $k$

# Background: Regression

- A type of supervised learning
- Given:
  - some labelled training dataset  $\mathcal{D} = \{(\mathbf{x}^{(i)}, y^{(i)})\}_{i=1}^n$
  - a loss function  $\ell: \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}$  where  $\mathcal{Y} = \mathbb{R}$
  - a hypothesis class or set of functions  $\mathcal{F}$

the goal is to find

$$\hat{f} = \operatorname{argmin}_{f \in \mathcal{F}} \sum_{i=1}^n \ell(f(\mathbf{x}^{(i)}), y^{(i)})$$

- Fun example: predicting the year a song was released based on (a representation of) its audio (see HW2)

# Background: Linear Regression (Ordinary Least Squares)

- A type of supervised learning
- Given:
  - some labelled training dataset  $\mathcal{D} = \{(\mathbf{x}^{(i)}, y^{(i)})\}_{i=1}^n$
  - $\ell(y, y') = (y - y')^2$
  - $\mathcal{F}$  = all functions of the form  $f(\mathbf{x}) = w_0 + \sum_{d=1}^k w_d x_d$

the goal is to find

$$\hat{f} = \operatorname{argmin}_{f \in \mathcal{F}} \sum_{i=1}^n (f(\mathbf{x}^{(i)}) - y^{(i)})^2$$

# Background: Linear Regression (Ordinary Least Squares)

- A type of supervised learning
- Given:
  - some labelled training dataset  $\mathcal{D} = \{(\mathbf{x}^{(i)}, y^{(i)})\}_{i=1}^n$
  - $\ell(y, y') = (y - y')^2$
  - $\mathcal{F}$  = all functions of the form  $f(\mathbf{x}) = \mathbf{w}^T [1 \ \mathbf{x}^T]^T$

the goal is to find

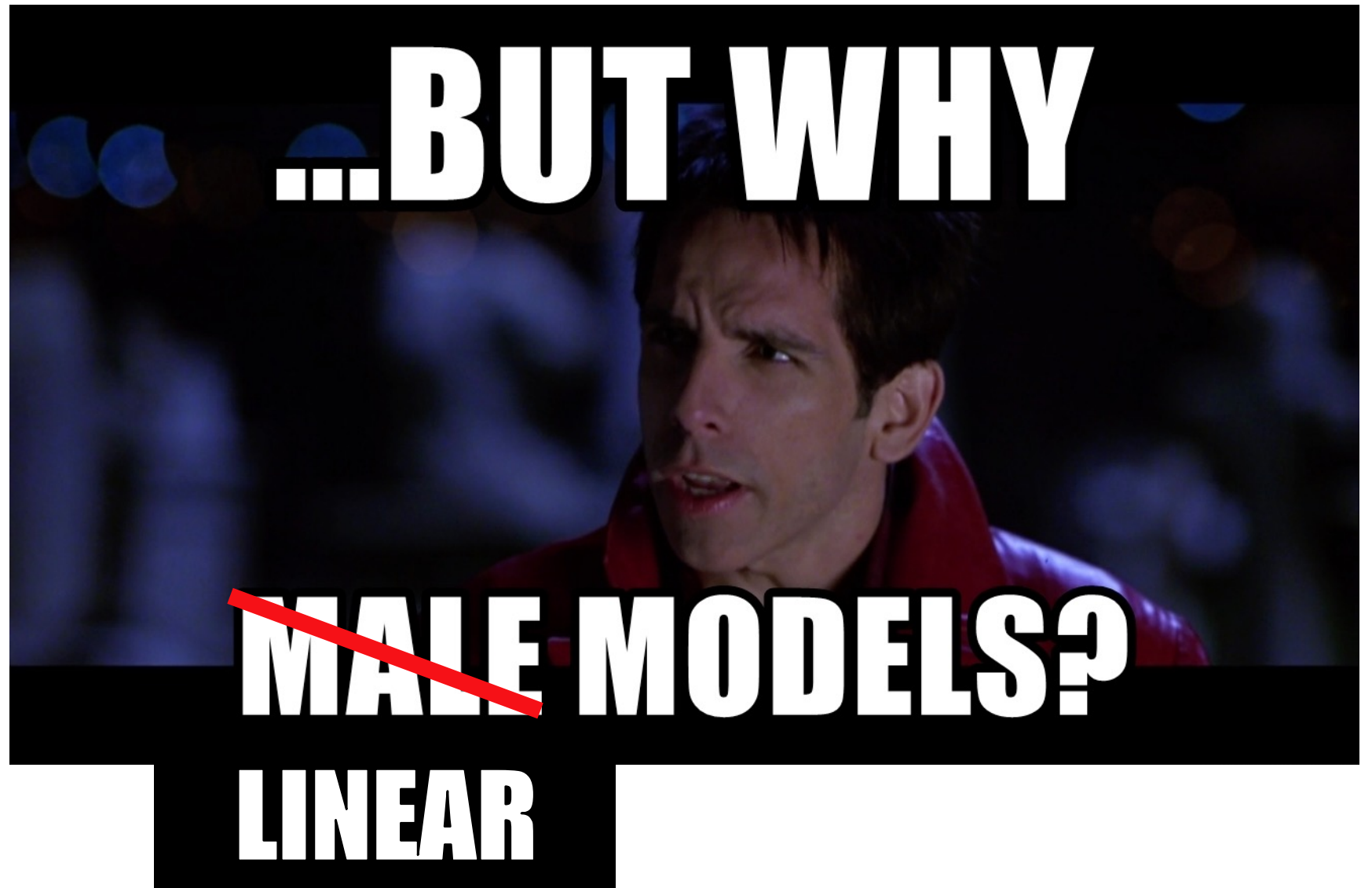
$$\hat{f} = \operatorname{argmin}_{f \in \mathcal{F}} \sum_{i=1}^n (f(\mathbf{x}^{(i)}) - y^{(i)})^2$$



# Background: Linear Regression (Ordinary Least Squares)

- A type of supervised learning
  - Given:
    - some labelled training dataset  $\mathcal{D} = \{(\mathbf{x}^{(i)}, y^{(i)})\}_{i=1}^n$
    - $\ell(y, y') = (y - y')^2$
    - $\mathcal{F}$  = all functions of the form  $f(\mathbf{x}) = \mathbf{w}^T \mathbf{x}$
- the goal is to find
- $\mathbf{w}$  implicitly prepended

$$\hat{\mathbf{w}} = \underset{\mathbf{w}}{\operatorname{argmin}} \sum_{i=1}^n (\mathbf{w}^T \mathbf{x}^{(i)} - y^{(i)})^2$$




# Background: Linear Regression (Ordinary Least Squares)

- A type of supervised learning
  - Given:
    - some labelled training dataset  $\mathcal{D} = \{(\mathbf{x}^{(i)}, y^{(i)})\}_{i=1}^n$
    - $\ell(y, y') = (y - y')^2$
    - $\mathcal{F}$  = all functions of the form  $f(\mathbf{x}) = \mathbf{w}^T \mathbf{x}$
- the goal is to find
- $\mathbf{w}$  implicitly prepended

$$\hat{\mathbf{w}} = \underset{\mathbf{w}}{\operatorname{argmin}} \sum_{i=1}^n (\mathbf{w}^T \mathbf{x}^{(i)} - y^{(i)})^2$$

# Background: Linear Regression (Ordinary Least Squares)

- A type of supervised learning
- Given:
  - some labelled training dataset  $\mathcal{D} = \{(\mathbf{x}^{(i)}, y^{(i)})\}_{i=1}^n$
  - $\ell(y, y') = (y - y')^2$
  - $\mathcal{F}$  = all functions of the form  $f(\mathbf{x}) = \mathbf{w}^T \mathbf{x}$   
the goal is to find  $\mathbf{w}$   1 implicitly prepended

$$\hat{\mathbf{w}} = \underset{\mathbf{w}}{\operatorname{argmin}} (X\mathbf{w} - \mathbf{y})^T (X\mathbf{w} - \mathbf{y})$$

$$\hat{\mathbf{w}} = \underset{\mathbf{w}}{\operatorname{argmin}} \|X\mathbf{w} - \mathbf{y}\|_2^2$$

- where  $X = \begin{bmatrix} \mathbf{x}^{(1)T} \\ \vdots \\ \mathbf{x}^{(n)T} \end{bmatrix}$  and  $\mathbf{y} = \begin{bmatrix} y^{(1)} \\ \vdots \\ y^{(n)} \end{bmatrix}$

Background:  
Linear  
Regression  
(Ordinary Least  
Squares)

$$L_{\mathcal{D}}(\mathbf{w}) = (\mathbf{X}\mathbf{w} - \mathbf{y})^T (\mathbf{X}\mathbf{w} - \mathbf{y})$$

# Background: Regularization

- A modification to empirical risk minimization that penalizes model complexity in order to combat overfitting
- Given:
  - some labelled training dataset  $\mathcal{D} = \{(\mathbf{x}^{(i)}, y^{(i)})\}_{i=1}^n$
  - a loss function  $\ell: \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}$
  - a hypothesis class or set of functions  $\mathcal{F}$
  - a regularizer  $R: \mathcal{W} \rightarrow \mathbb{R}$
  - a coefficient of regularization  $\lambda$

the goal is to find

$$\hat{\mathbf{w}} = \operatorname{argmin}_{\mathbf{w}} \sum_{i=1}^n \ell(f_{\mathbf{w}}(\mathbf{x}^{(i)}), y^{(i)}) + \lambda R(\mathbf{w})$$

# Background: Ridge Regression

- A modification to empirical risk minimization that penalizes model complexity in order to combat overfitting
- Given:
  - some labelled training dataset  $\mathcal{D} = \{(\mathbf{x}^{(i)}, y^{(i)})\}_{i=1}^n$
  - $\ell(y, y') = (y - y')^2$
  - $\mathcal{F}$  = all functions of the form  $f(\mathbf{x}) = \mathbf{w}^T \mathbf{x}$
  - $R(\mathbf{w}) = \|\mathbf{w}\|_2^2 = \mathbf{w}^T \mathbf{w}$
  - a coefficient of regularization  $\lambda$

the goal is to find

$$\hat{\mathbf{w}} = \underset{\mathbf{w}}{\operatorname{argmin}} (\mathbf{X}\mathbf{w} - \mathbf{y})^T (\mathbf{X}\mathbf{w} - \mathbf{y}) + \lambda \mathbf{w}^T \mathbf{w}$$

## Background: Ridge Regression

$$L_{\mathcal{D}}(\mathbf{w}) = (\mathbf{X}\mathbf{w} - \mathbf{y})^T (\mathbf{X}\mathbf{w} - \mathbf{y}) + \lambda \mathbf{w}^T \mathbf{w}$$

•  
•  
•

$$\rightarrow \hat{\mathbf{w}} = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I}_k)^{-1} \mathbf{X}^T \mathbf{y}$$

where  $\mathbf{I}_k$  is the  $k \times k$  identity matrix



Aside:  
How can we  
set  $\lambda$ ?

$$L_{\mathcal{D}}(\mathbf{w}) = (\mathbf{X}\mathbf{w} - \mathbf{y})^T (\mathbf{X}\mathbf{w} - \mathbf{y}) + \lambda \mathbf{w}^T \mathbf{w}$$

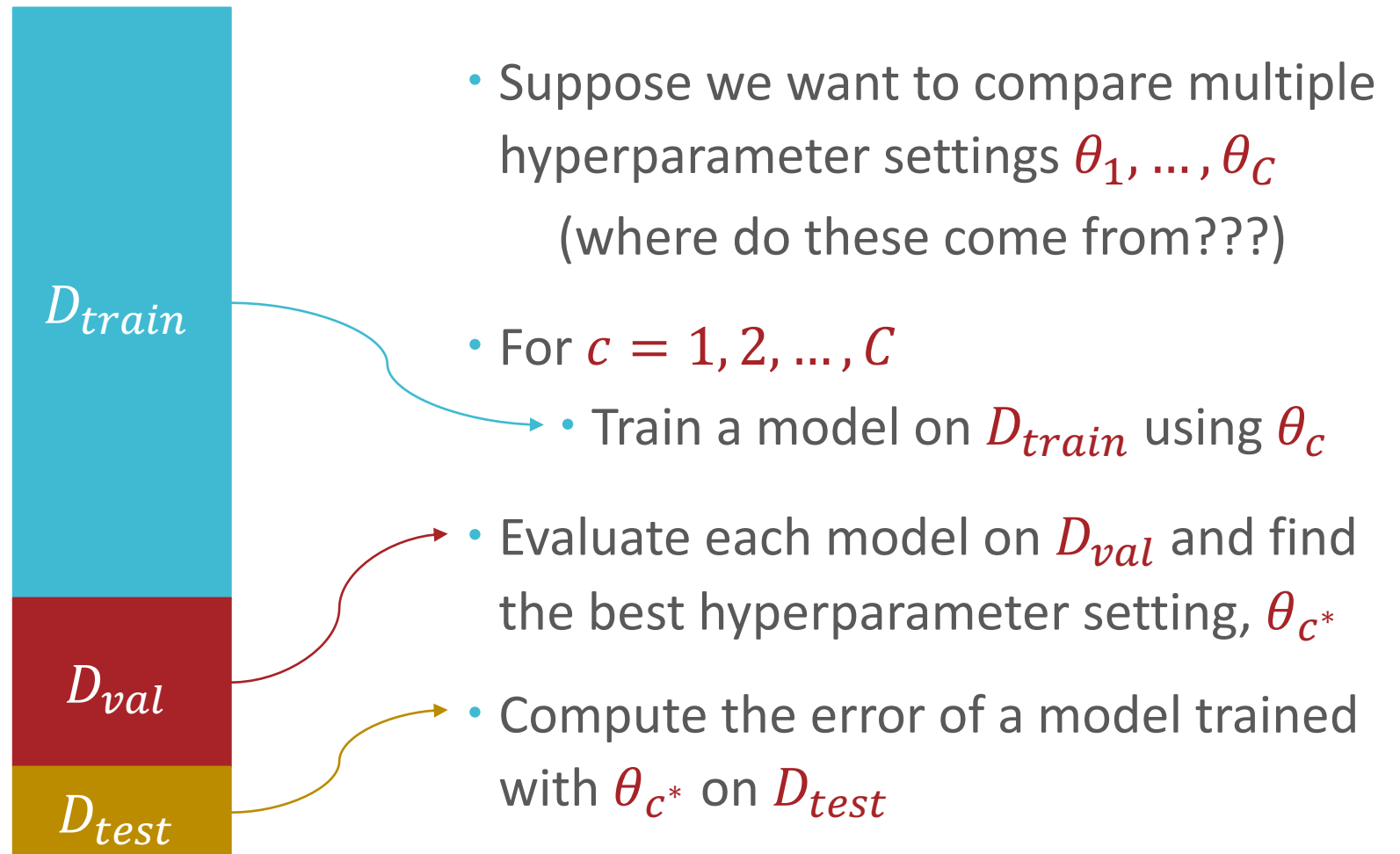
•  
•  
•

$$\rightarrow \hat{\mathbf{w}} = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I}_k)^{-1} \mathbf{X}^T \mathbf{y}$$

where  $\mathbf{I}_k$  is the  $k \times k$  identity matrix

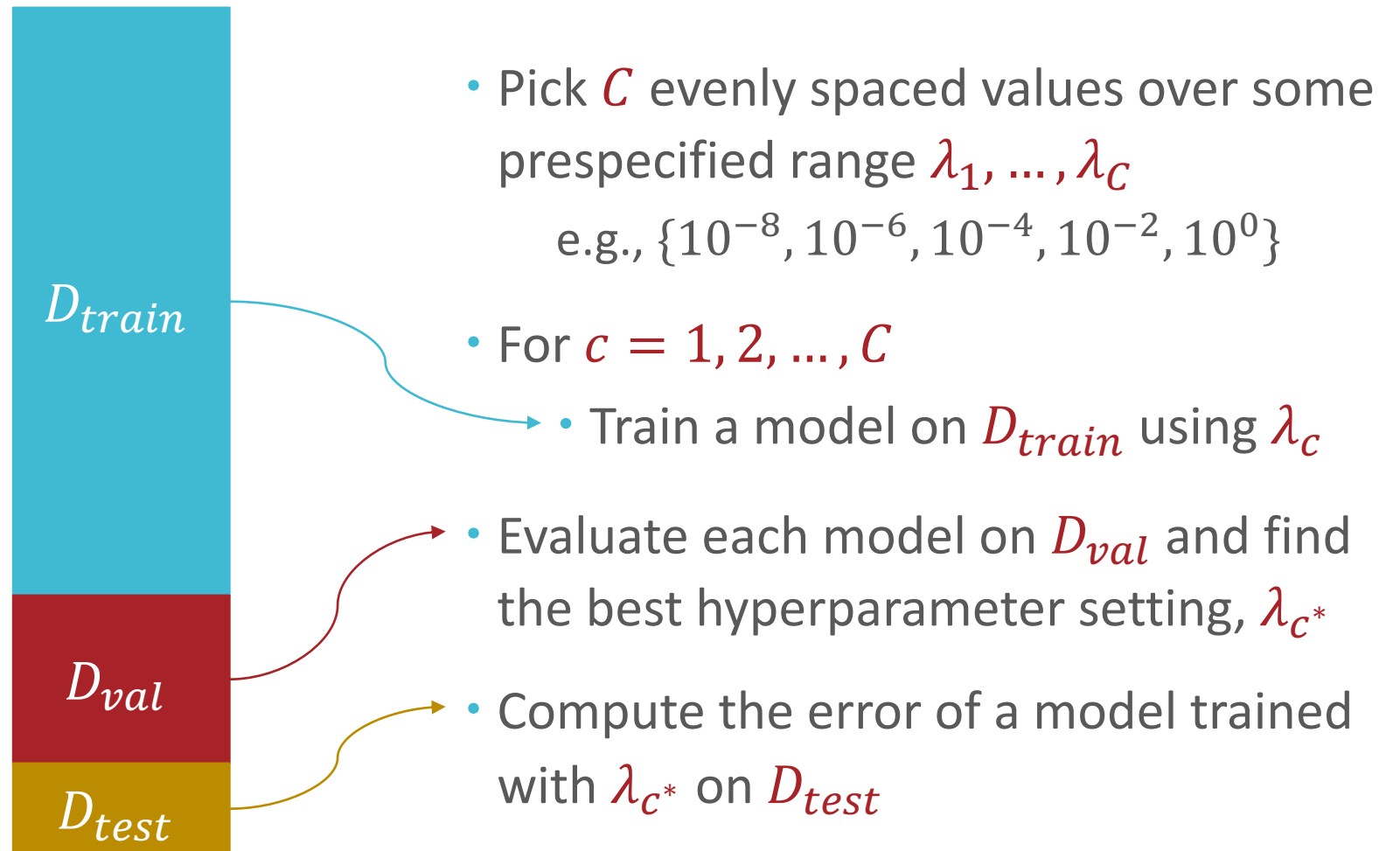
# Recall: Machine Learning Pipeline

- Hyperparameter optimization



# HW2 Preview: Grid Search

- Hyperparameter optimization



# Linear Regression: Computational Cost

$$\hat{\mathbf{w}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

1. Does this quantity exist, i.e., is  $\mathbf{X}^T \mathbf{X}$  invertible?

When  $n \gg k + 1$ ,  $\mathbf{X}^T \mathbf{X}$  is (almost always) full rank and; otherwise, we can compute a *pseudoinverse* via singular value decomposition

2. If so, how expensive is it to compute?
  - Computing  $\mathbf{X}^T \mathbf{X}$ :  $O(nk^2)$  time &  $O(k^2)$  space
  - Inverting  $\mathbf{X}^T \mathbf{X}$ :  $O(k^3)$  time &  $O(k^2)$  space
  - Storing  $\mathbf{X}$ :  $O(nk)$  space

Key takeaway: computational bottlenecks will change based on the relationship between  $n$  and  $k$

# Linear Regression: Large $n$ , Small $k$

- Assume  $O(k^3)$  computation and  $O(k^2)$  storage is possible on a single machine
  - ✓ We can store and invert  $X^T X$ 
    - We cannot compute  $X^T X$
    - We cannot store  $X$
- Idea: distribute storage of  $X$  and computation of  $X^T X$ 
  1. Store the rows of  $X$  across different machines
  2. Compute  $X^T X$  as the sum of outer products

# Matrix Multiplication via Inner Products

$$\begin{bmatrix} 1 & 4 & 5 \\ 3 & 1 & 3 \end{bmatrix} \begin{bmatrix} 2 \\ 2 \\ 3 \end{bmatrix} = \begin{bmatrix} \phantom{0} \\ \phantom{0} \end{bmatrix}$$

$$1 * 2 + 4 * 2 + 5 * 3 = 25$$

# Matrix Multiplication via Inner Products

$$\begin{bmatrix} 1 & 4 & 5 \\ 3 & 1 & 3 \end{bmatrix} \begin{bmatrix} 2 \\ 2 \\ 1 \\ 4 \end{bmatrix} = \begin{bmatrix} 25 \\ \end{bmatrix}$$

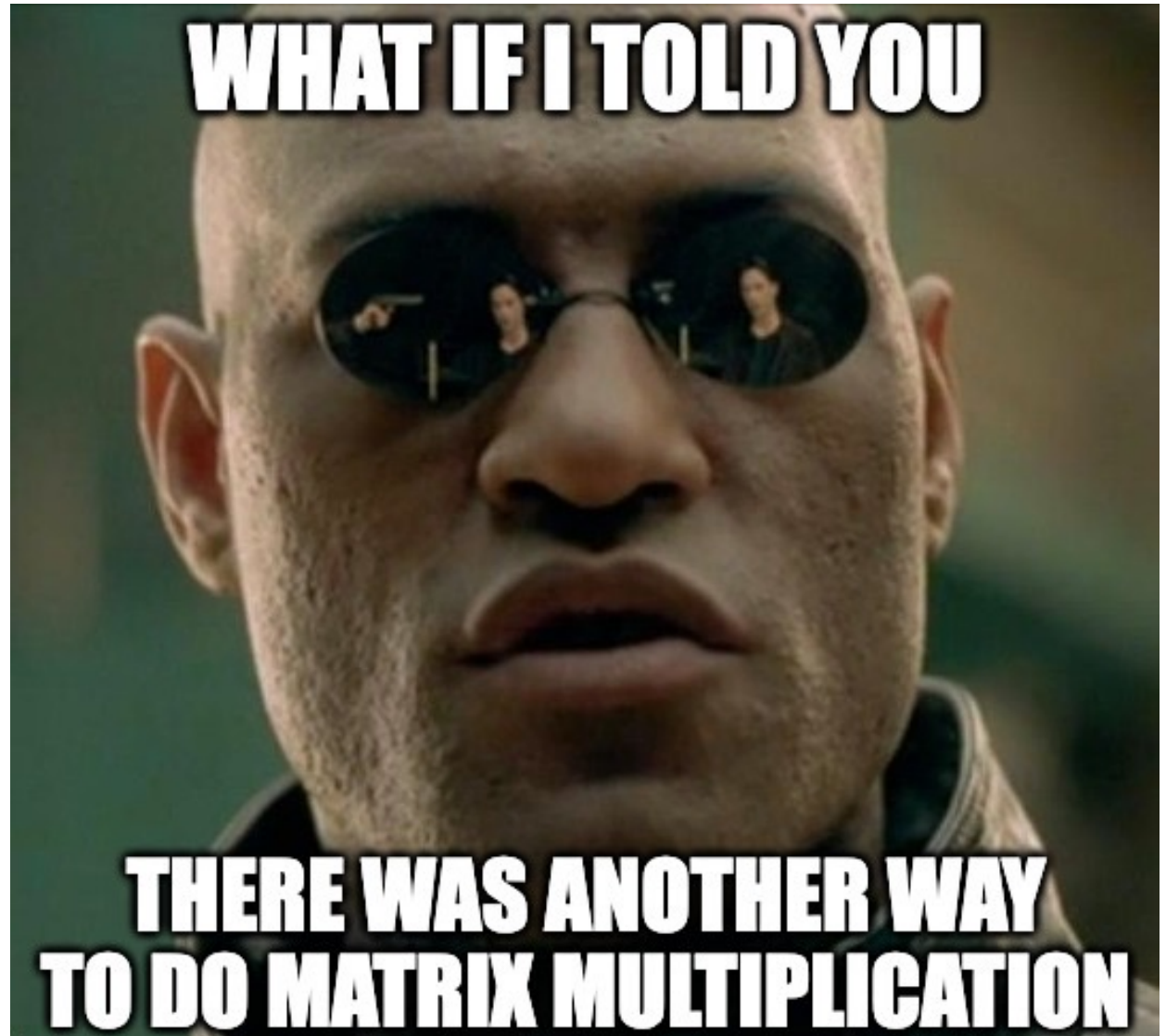
$$1 * 2 + 4 * 1 + 5 * 4 = 26$$

# Matrix Multiplication via Inner Products

$$\begin{bmatrix} 1 & 4 & 5 \\ 3 & 1 & 3 \end{bmatrix} \begin{bmatrix} 2 & 2 \\ 2 & 1 \\ 3 & 4 \end{bmatrix} = \begin{bmatrix} 25 & 26 \\ 17 & 19 \end{bmatrix}$$



# Matrix Multiplication via Inner Products



# Matrix Multiplication via Outer Products

$$\begin{bmatrix} 1 \\ 3 \end{bmatrix} \begin{matrix} 4 & 5 \\ 1 & 3 \end{matrix} \begin{bmatrix} 2 & 2 \\ 2 & 1 \\ 3 & 4 \end{bmatrix} = \begin{bmatrix} \phantom{0} & \phantom{0} \\ \phantom{0} & \phantom{0} \end{bmatrix}$$

$$\begin{bmatrix} 2 & 2 \\ 6 & 6 \end{bmatrix}$$

# Matrix Multiplication via Outer Products

$$\begin{bmatrix} 1 & 4 & 5 \\ 3 & 1 & 3 \end{bmatrix} \begin{bmatrix} 2 & 2 \\ 2 & 1 \\ 3 & 4 \end{bmatrix} = \begin{bmatrix} & & \\ & & \\ & & \end{bmatrix}$$

$$\begin{bmatrix} 2 & 2 \\ 6 & 6 \end{bmatrix} + \begin{bmatrix} 8 & 4 \\ 2 & 1 \end{bmatrix}$$

# Matrix Multiplication via Outer Products

$$\begin{bmatrix} 1 & 4 & 5 \\ 3 & 1 & 3 \end{bmatrix} \begin{bmatrix} 2 & 2 \\ 2 & 1 \\ 3 & 4 \end{bmatrix} = \begin{bmatrix} & \\ & \end{bmatrix}$$

$$\begin{bmatrix} 2 & 2 \\ 6 & 6 \end{bmatrix} + \begin{bmatrix} 8 & 4 \\ 2 & 1 \end{bmatrix} + \begin{bmatrix} 15 & 20 \\ 9 & 12 \end{bmatrix}$$

# Matrix Multiplication via Outer Products

$$\begin{bmatrix} 1 & 4 & 5 \\ 3 & 1 & 3 \end{bmatrix} \begin{bmatrix} 2 & 2 \\ 2 & 1 \\ 3 & 4 \end{bmatrix} = \begin{bmatrix} 25 & 26 \\ 17 & 19 \end{bmatrix}$$

$$\begin{bmatrix} 2 & 2 \\ 6 & 6 \end{bmatrix} + \begin{bmatrix} 8 & 4 \\ 2 & 1 \end{bmatrix} + \begin{bmatrix} 15 & 20 \\ 9 & 12 \end{bmatrix}$$

# Matrix Multiplication via Outer Products

$$\begin{bmatrix} A_{11} & A_{12} & \cdots & A_{1m} \\ A_{21} & A_{22} & \cdots & A_{2m} \\ \vdots & \vdots & \ddots & \vdots \\ A_{n1} & A_{n2} & \cdots & A_{nm} \end{bmatrix} \begin{bmatrix} B_{11} & \cdots & B_{1k} \\ B_{21} & \cdots & B_{2k} \\ B_{31} & \cdots & B_{3k} \\ \vdots & \ddots & \vdots \\ B_{m1} & \cdots & B_{mk} \end{bmatrix}$$

$$= \begin{bmatrix} \sum_{i=1}^m A_{1i}B_{i1} & \cdots & \sum_{i=1}^m A_{1i}B_{ik} \\ \vdots & \ddots & \vdots \\ \sum_{i=1}^m A_{ni}B_{i1} & \cdots & \sum_{i=1}^m A_{ni}B_{ik} \end{bmatrix} = \sum_{i=1}^m \begin{bmatrix} A_{1i}B_{i1} & \cdots & A_{1i}B_{ik} \\ \vdots & \ddots & \vdots \\ A_{ni}B_{i1} & \cdots & A_{ni}B_{ik} \end{bmatrix}$$

# Distributed Computation of $(X^T X)^{-1}$

$$X^T X = \begin{bmatrix} \uparrow & \uparrow & \dots & \uparrow \\ \mathbf{x}^{(1)} & \mathbf{x}^{(2)} & \dots & \mathbf{x}^{(n)} \\ \downarrow & \downarrow & \dots & \downarrow \end{bmatrix} \begin{bmatrix} \leftarrow & \mathbf{x}^{(1)T} & \rightarrow \\ \leftarrow & \mathbf{x}^{(2)T} & \rightarrow \\ \vdots & \vdots & \vdots \\ \leftarrow & \mathbf{x}^{(n)T} & \rightarrow \end{bmatrix} = \sum_{i=1}^n \mathbf{x}^{(i)} \mathbf{x}^{(i)T}$$

- Idea: distribute  $\mathbf{x}^{(i)}$  and compute summands in parallel

Workers	$\begin{bmatrix} \leftarrow & \mathbf{x}^{(1)T} & \rightarrow \\ \leftarrow & \mathbf{x}^{(4)T} & \rightarrow \\ \vdots & \vdots & \vdots \end{bmatrix}$	$\begin{bmatrix} \leftarrow & \mathbf{x}^{(2)T} & \rightarrow \\ \leftarrow & \mathbf{x}^{(3)T} & \rightarrow \\ \vdots & \vdots & \vdots \end{bmatrix}$	$\begin{bmatrix} \leftarrow & \mathbf{x}^{(5)T} & \rightarrow \\ \leftarrow & \mathbf{x}^{(7)T} & \rightarrow \\ \vdots & \vdots & \vdots \end{bmatrix}$	$O(nk)$ distributed storage (total)	
Map	$\mathbf{x}^{(i)} \mathbf{x}^{(i)T}$	$\mathbf{x}^{(i)} \mathbf{x}^{(i)T}$	$\mathbf{x}^{(i)} \mathbf{x}^{(i)T}$	$O(nk^2)$ distributed work (total)	$O(k^2)$ local storage
Reduce	$\left( \sum_{i=1}^n \mathbf{x}^{(i)} \mathbf{x}^{(i)T} \right)^{-1}$			$O(k^3)$ local work	$O(k^2)$ local storage

# Distributed Computation of $(X^T X)^{-1}$



Workers

$$\begin{bmatrix} \leftarrow & \mathbf{x}^{(1)T} & \rightarrow \\ \leftarrow & \mathbf{x}^{(4)T} & \rightarrow \\ \vdots & \vdots & \vdots \end{bmatrix} \quad \begin{bmatrix} \leftarrow & \mathbf{x}^{(2)T} & \rightarrow \\ \leftarrow & \mathbf{x}^{(3)T} & \rightarrow \\ \vdots & \vdots & \vdots \end{bmatrix} \quad \begin{bmatrix} \leftarrow & \mathbf{x}^{(5)T} & \rightarrow \\ \leftarrow & \mathbf{x}^{(7)T} & \rightarrow \\ \vdots & \vdots & \vdots \end{bmatrix}$$

$O(nk)$  distributed  
storage (total)

`trainData.map(compute_outer_prods)`

$O(nk^2)$   
distributed  
work (total)

$O(k^2)$   
local  
storage

`trainData.reduce(sum_and_invert)`

$O(k^3)$   
local work

$O(k^2)$   
local  
storage

# Distributed Computation of $(X^T X)^{-1}$

# Linear Regression: Large $n$ , Large $k$

- Now,  $O(k^3)$  computation and  $O(k^2)$  storage is *not* possible on a single machine

We cannot store and invert  $X^T X$

We cannot compute  $X^T X$

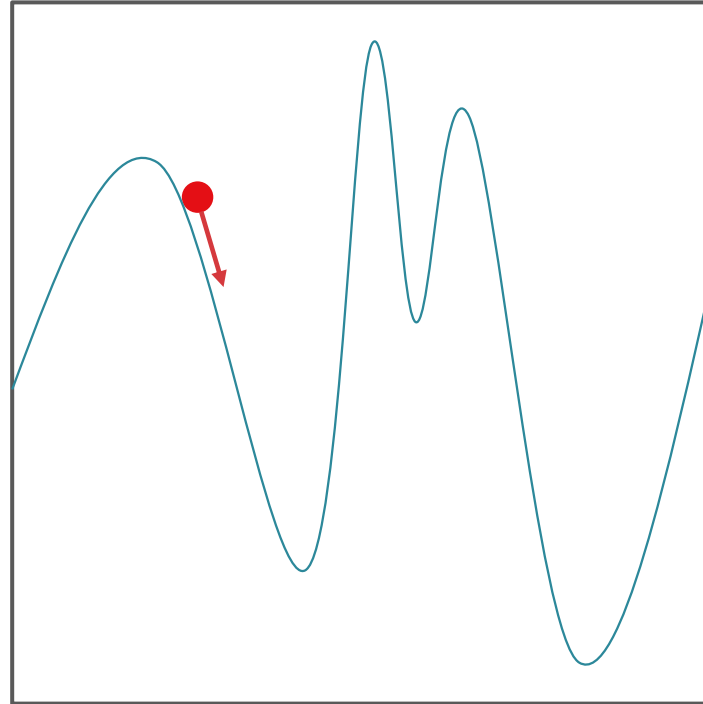
We cannot store  $X$

10-605/805 Principle #1: computation and storage should be at most linear in  $n$  and  $k$

- Idea: use a different algorithm!

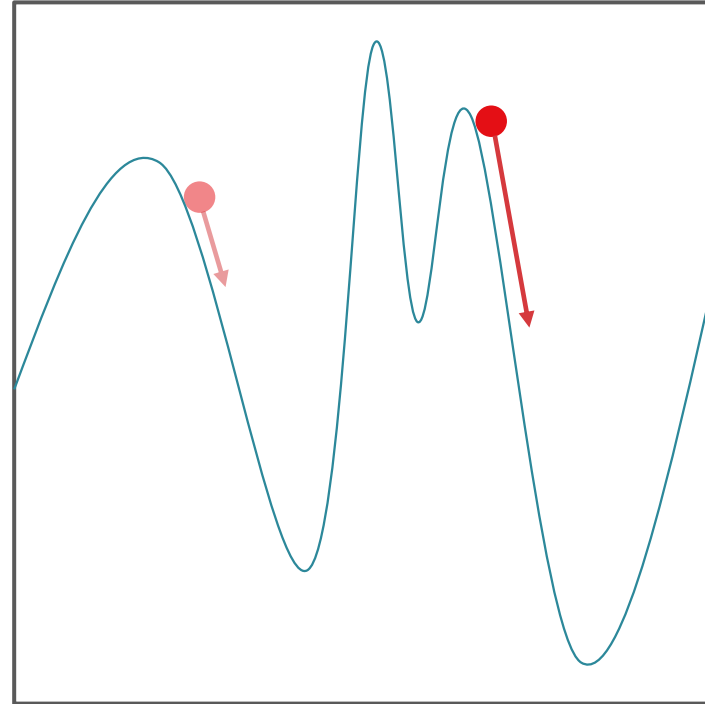
# Background: Gradient Descent

- An iterative method for minimizing functions
- Requires the gradient to exist everywhere



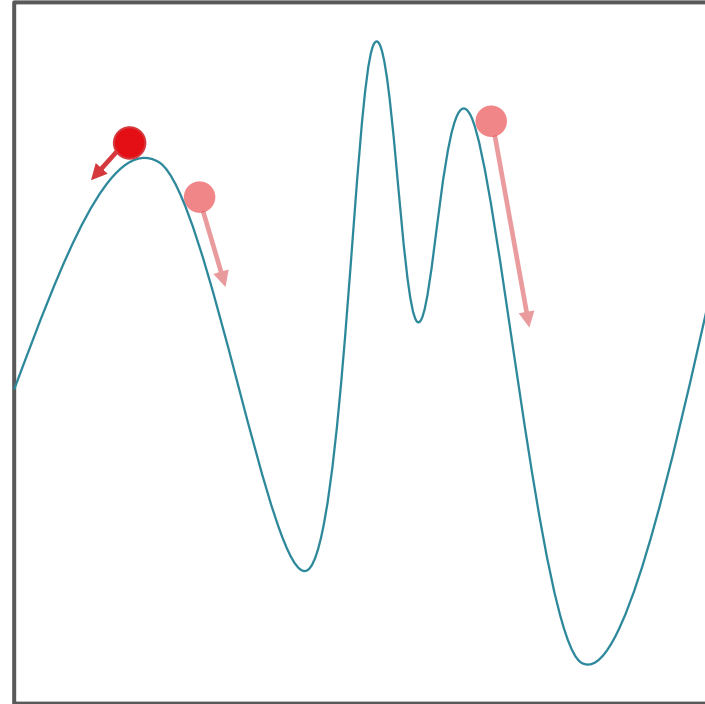
# Background: Gradient Descent

- An iterative method for minimizing functions
- Requires the gradient to exist everywhere



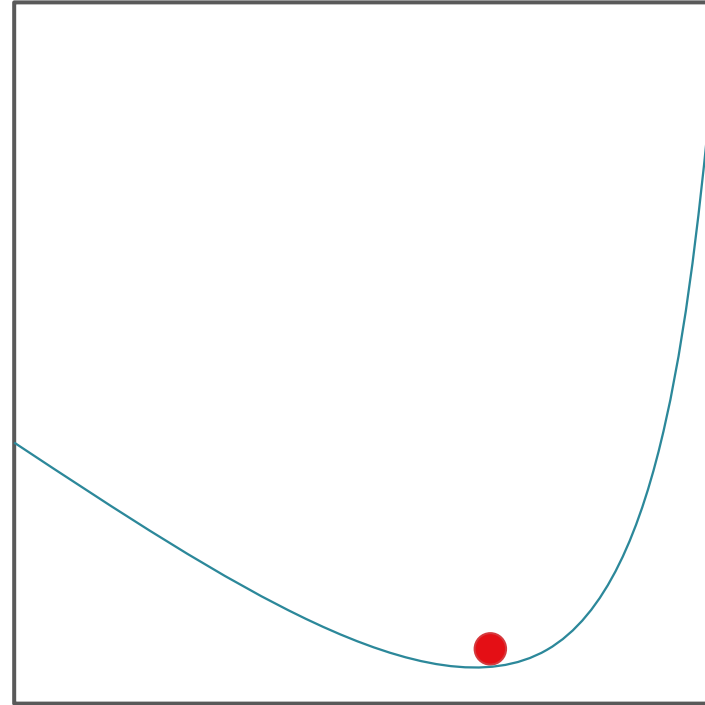
# Background: Gradient Descent

- An iterative method for minimizing functions
- Requires the gradient to exist everywhere



# Background: Gradient Descent

- An iterative method for minimizing functions
- Requires the gradient to exist everywhere



- Good news: the linear regression objective is convex so gradient descent will always converge to the global minimum

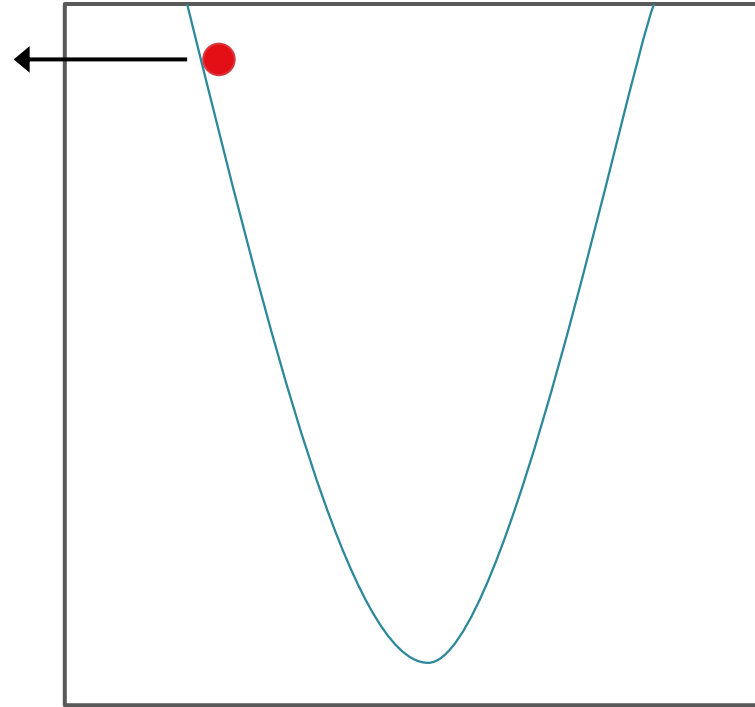
# Background: Gradient Descent

- Suppose we're trying to minimize some function  $L$  and we're currently at some location  $\mathbf{w}^{(t)}$
- Move some distance,  $\alpha$ , in the “most downhill” direction,  $\mathbf{v}$ :  
$$\mathbf{w}^{(t+1)} = \mathbf{w}^{(t)} + \alpha \mathbf{v}$$
- The gradient points in the direction of steepest *increase* ...
- ... so let's move in the opposite direction!

$$\mathbf{w}^{(t+1)} = \mathbf{w}^{(t)} - \eta \nabla_{\mathbf{w}} L(\mathbf{w}^{(t)})$$

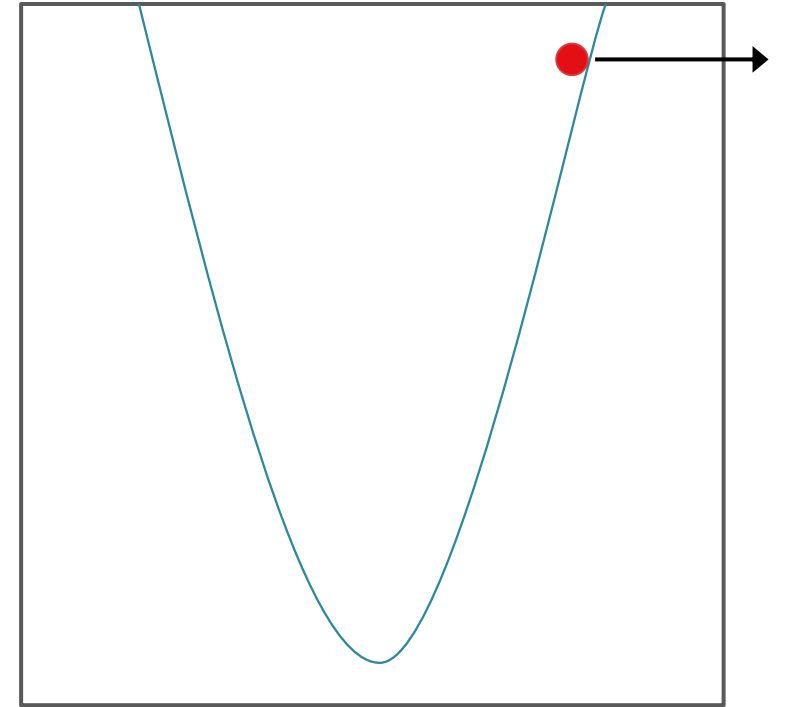
# Background: Gradient Descent

Direction of  
gradient



$$2x < 0 \text{ for } x < 0$$

Direction of  
gradient



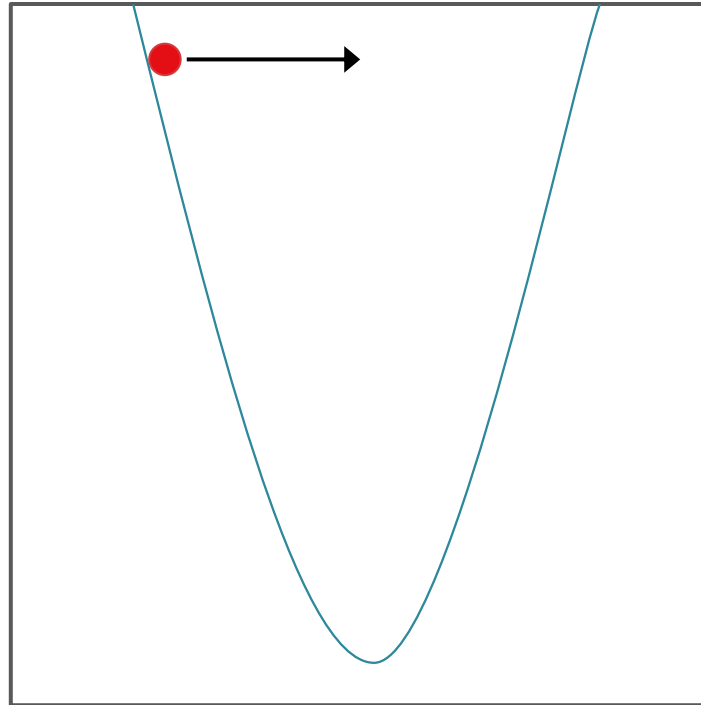
$$2x > 0 \text{ for } x > 0$$

$$\frac{\partial}{\partial x} x^2 = 2x$$



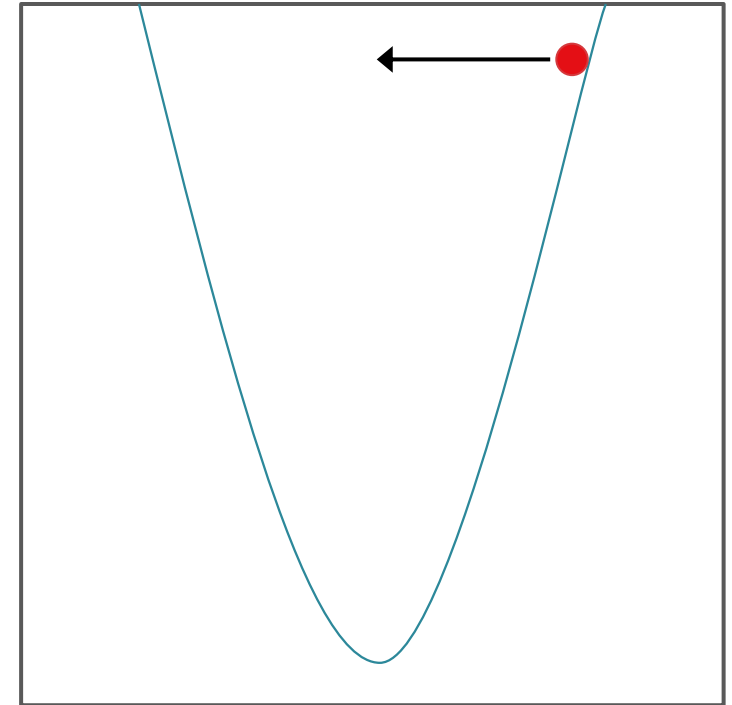
# Background: Gradient Descent

Direction of  
global minimum



$$2x < 0 \text{ for } x < 0$$

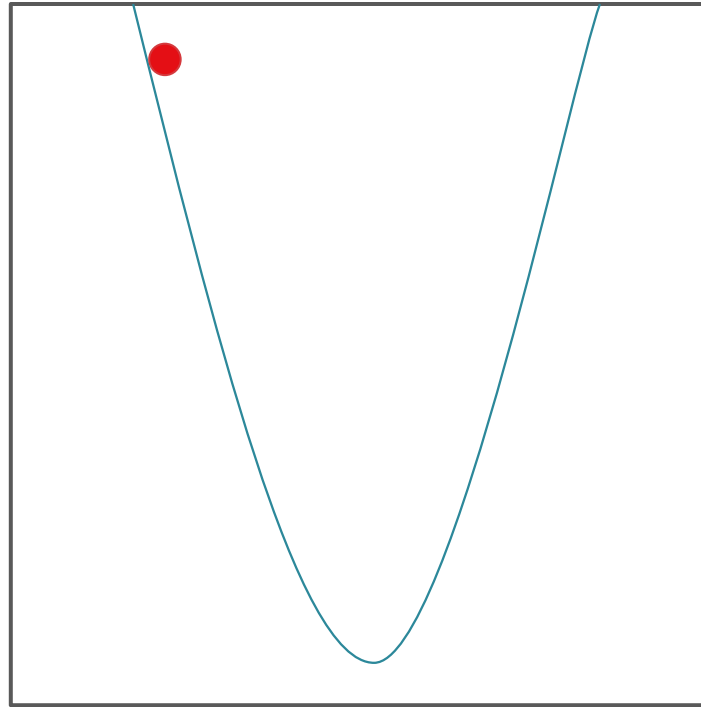
Direction of  
global minimum



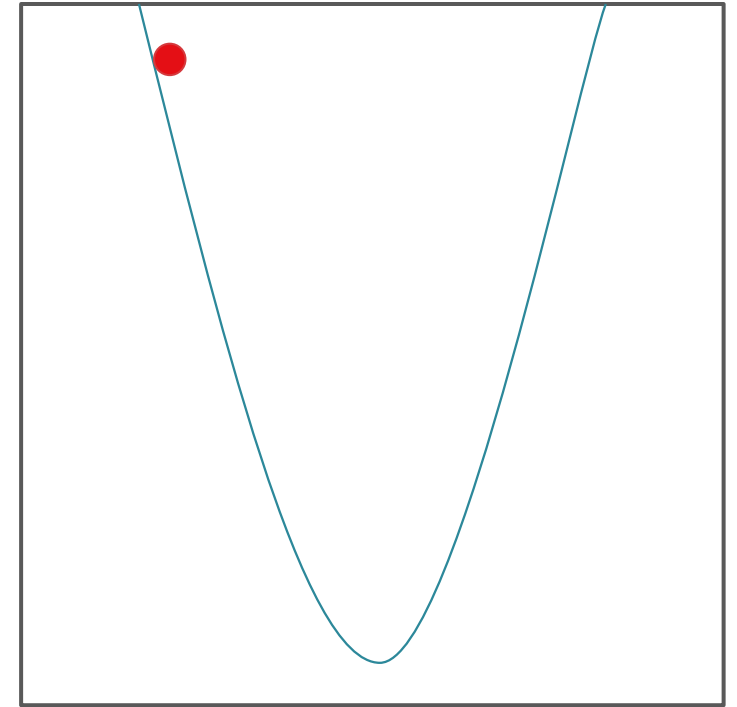
$$2x > 0 \text{ for } x > 0$$

$$\frac{\partial}{\partial x} x^2 = 2x$$

# Background: Gradient Descent Step Size

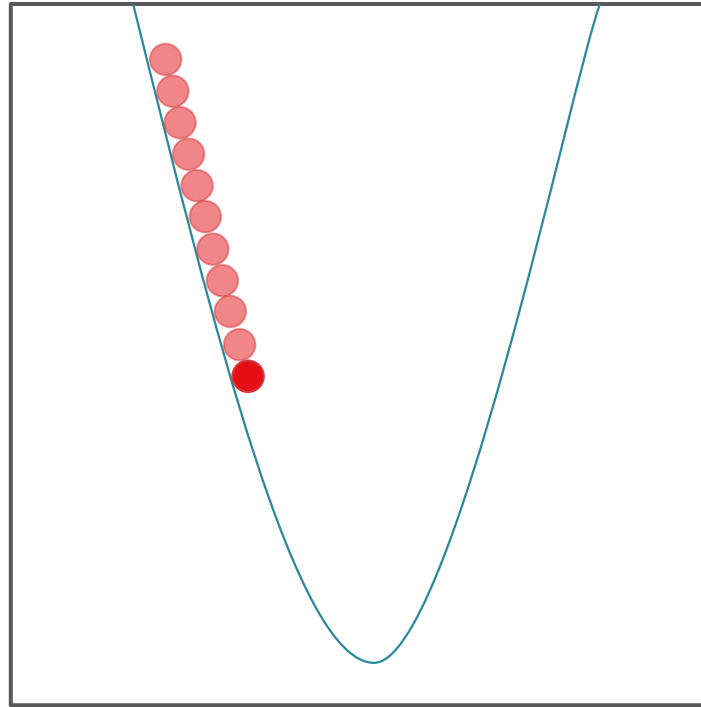


Small  $\alpha$

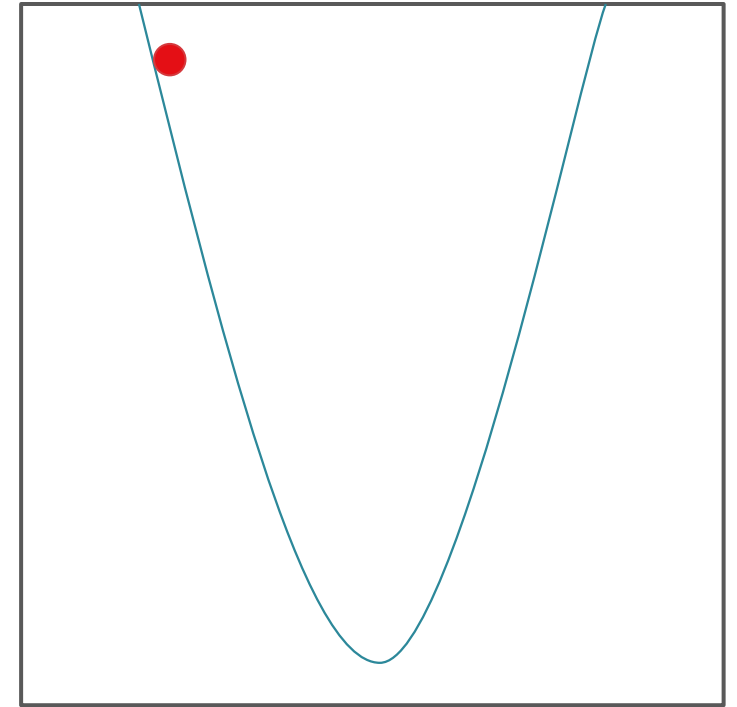


Large  $\alpha$

# Background: Gradient Descent Step Size

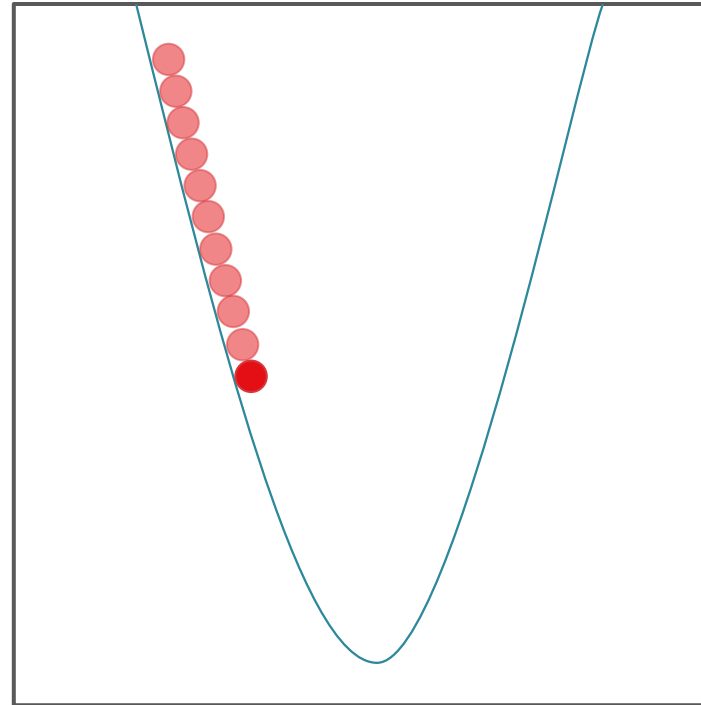


Small  $\alpha$

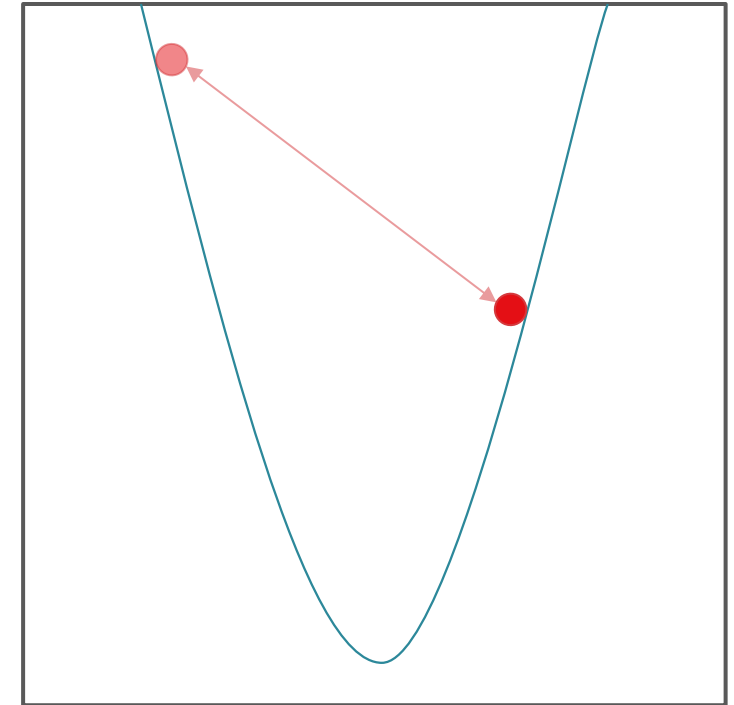


Large  $\alpha$

# Background: Gradient Descent Step Size



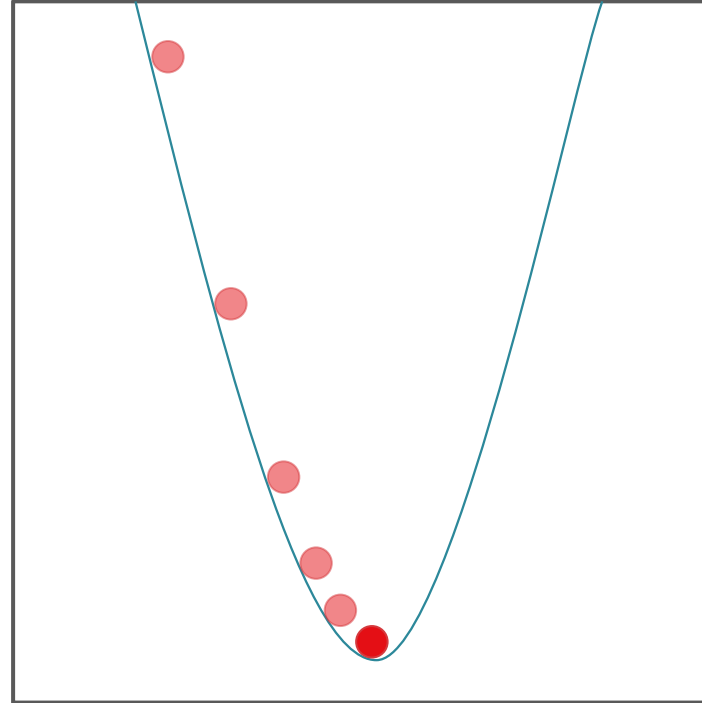
Small  $\alpha$



Large  $\alpha$

# Background: Gradient Descent Step Size

- Use a variable  $\alpha^{(t)}$  instead of a fixed  $\alpha$ !



- Example:  $\alpha^{(t)} = \frac{\alpha}{n\sqrt{t}}$

# Gradient Descent for Linear Regression

- Input:  $\mathcal{D} = \{(\mathbf{x}^{(i)}, y^{(i)})\}_{i=1}^n, \alpha$
- 1. Initialize  $\mathbf{w}^{(0)}$  to all zeros and set  $t = 0$
- 2. While TERMINATION CRITERION is not satisfied
  - a. Compute the gradient:

$$\nabla_{\mathbf{w}} L_{\mathcal{D}}(\mathbf{w}^{(t)}) = (2X^T X \mathbf{w}^{(t)} - 2X^T \mathbf{y})$$

- b. Update the weights:

$$\mathbf{w}^{(t+1)} \leftarrow \mathbf{w}^{(t)} - \frac{\alpha}{n\sqrt{t}} (2X^T X \mathbf{w}^{(t)} - 2X^T \mathbf{y})$$

- c. Increment  $t$ :  $t \leftarrow t + 1$

- Output:  $\mathbf{w}^{(t)}$

# Gradient Descent for Linear Regression

- Input:  $\mathcal{D} = \{(\mathbf{x}^{(i)}, y^{(i)})\}_{i=1}^n, \alpha$

1. Initialize  $\mathbf{w}^{(0)}$  to all zeros and set  $t = 0$
2. While TERMINATION CRITERION is not satisfied
  - a. Compute the gradient:

$$\nabla_{\mathbf{w}} L_{\mathcal{D}}(\mathbf{w}^{(t)}) = 2 \sum_{i=1}^n \left( \mathbf{w}^{(t)T} \mathbf{x}^{(i)} - y^{(i)} \right) \mathbf{x}^{(i)}$$

- b. Update the weights:

$$\mathbf{w}^{(t+1)} \leftarrow \mathbf{w}^{(t)} - \frac{\alpha}{n\sqrt{t}} \sum_{i=1}^n \left( \mathbf{w}^{(t)T} \mathbf{x}^{(i)} - y^{(i)} \right) \mathbf{x}^{(i)}$$

- c. Increment  $t$ :  $t \leftarrow t + 1$

- Output:  $\mathbf{w}^{(t)}$

Idea: distribute  $\mathbf{x}^{(i)}$  and compute summands in parallel

- Input:  $\mathcal{D} = \{(\mathbf{x}^{(i)}, y^{(i)})\}_{i=1}^n, \alpha$

1. Initialize  $\mathbf{w}^{(0)}$  to all zeros and set  $t = 0$
2. While TERMINATION CRITERION is not satisfied
  - a. Compute the gradient:

$$\nabla_{\mathbf{w}} L_{\mathcal{D}}(\mathbf{w}^{(t)}) = 2 \sum_{i=1}^n \left( \mathbf{w}^{(t)T} \mathbf{x}^{(i)} - y^{(i)} \right) \mathbf{x}^{(i)}$$


- b. Update the weights:

$$\mathbf{w}^{(t+1)} \leftarrow \mathbf{w}^{(t)} - \frac{\alpha}{n\sqrt{t}} \sum_{i=1}^n \left( \mathbf{w}^{(t)T} \mathbf{x}^{(i)} - y^{(i)} \right) \mathbf{x}^{(i)}$$

- c. Increment  $t$ :  $t \leftarrow t + 1$

- Output:  $\mathbf{w}^{(t)}$





Workers	$\begin{bmatrix} \leftarrow & \mathbf{x}^{(1)T} & \rightarrow \\ \leftarrow & \mathbf{x}^{(4)T} & \rightarrow \\ \vdots & \vdots & \vdots \end{bmatrix}$	$\begin{bmatrix} \leftarrow & \mathbf{x}^{(2)T} & \rightarrow \\ \leftarrow & \mathbf{x}^{(3)T} & \rightarrow \\ \vdots & \vdots & \vdots \end{bmatrix}$	$\begin{bmatrix} \leftarrow & \mathbf{x}^{(5)T} & \rightarrow \\ \leftarrow & \mathbf{x}^{(7)T} & \rightarrow \\ \vdots & \vdots & \vdots \end{bmatrix}$	$O(nk)$ distributed storage (total)	
Map	$(\mathbf{w}^{(t)T} \mathbf{x}^{(i)} - y^{(i)}) \mathbf{x}^{(i)}$	$(\mathbf{w}^{(t)T} \mathbf{x}^{(i)} - y^{(i)}) \mathbf{x}^{(i)}$	$(\mathbf{w}^{(t)T} \mathbf{x}^{(i)} - y^{(i)}) \mathbf{x}^{(i)}$	$O(nk)$ distributed work (total)	$O(k)$ local storage
Reduce	$\mathbf{w}^{(t+1)} \leftarrow \mathbf{w}^{(t)} - \frac{\alpha}{n\sqrt{t}} \sum_{i=1}^n (\mathbf{w}^{(t)T} \mathbf{x}^{(i)} - y^{(i)}) \mathbf{x}^{(i)}$			$O(k)$ local work	$O(k)$ local storage

Issue: all workers must have the latest weight vector

# Distributed Gradient Descent

Workers

$$\begin{bmatrix} \leftarrow & \mathbf{x}^{(1)T} & \rightarrow \\ \leftarrow & \mathbf{x}^{(4)T} & \rightarrow \\ \vdots & \vdots & \vdots \end{bmatrix}$$

$$\begin{bmatrix} \leftarrow & \mathbf{x}^{(2)T} & \rightarrow \\ \leftarrow & \mathbf{x}^{(3)T} & \rightarrow \\ \vdots & \vdots & \vdots \end{bmatrix}$$

$$\begin{bmatrix} \leftarrow & \mathbf{x}^{(5)T} & \rightarrow \\ \leftarrow & \mathbf{x}^{(7)T} & \rightarrow \\ \vdots & \vdots & \vdots \end{bmatrix}$$

$O(nk)$  distributed storage (total)

`trainData.map(compute_pointwise_grads( $\mathbf{w}^{(t)}$ ))`

$O(nk)$   
distributed  
work (total)

$O(k)$   
local storage

`trainData.sum()`

$O(k)$   
local work

$O(k)$   
local storage

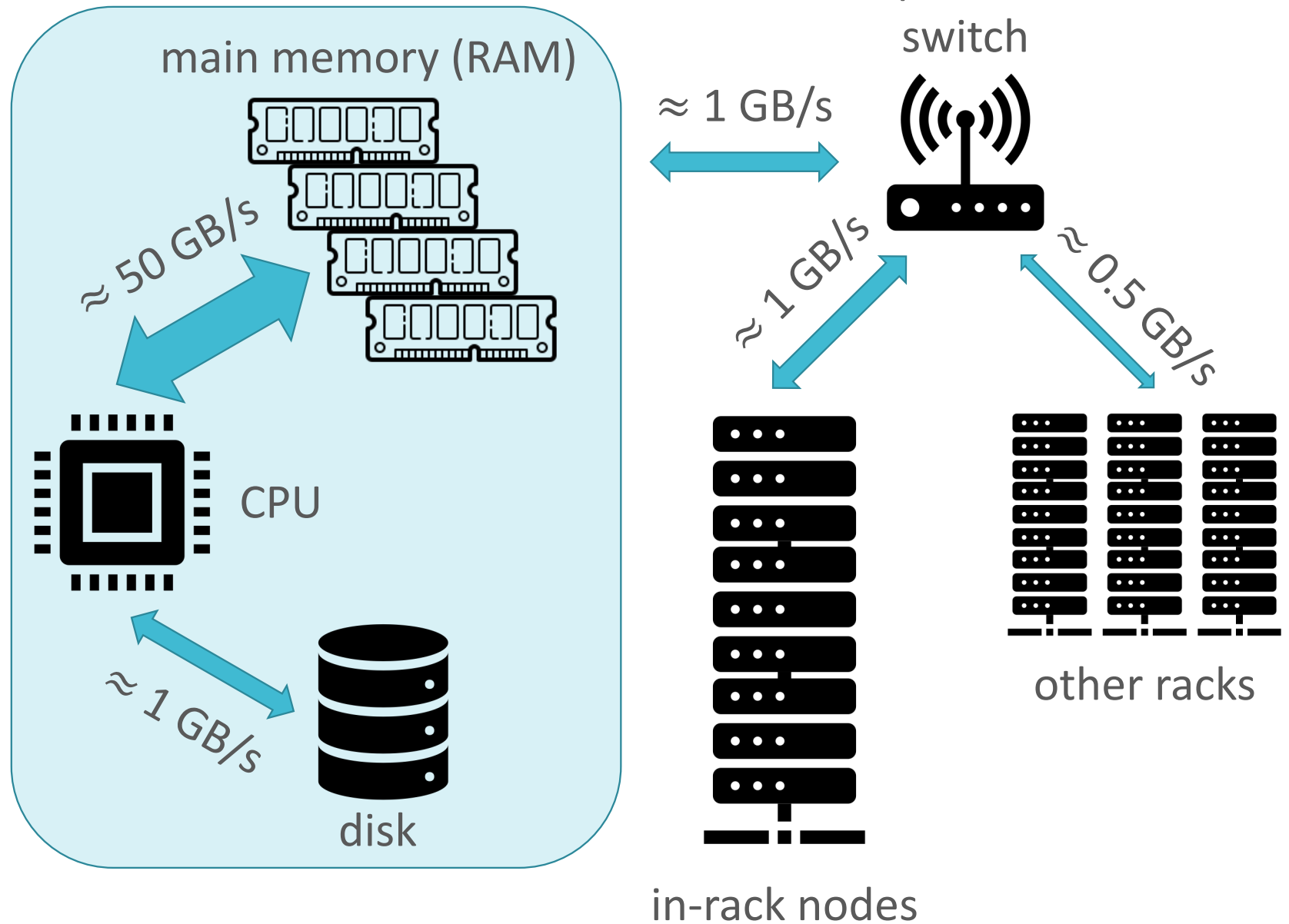
Issue: all workers must have the latest weight vector

# Distributed Gradient Descent

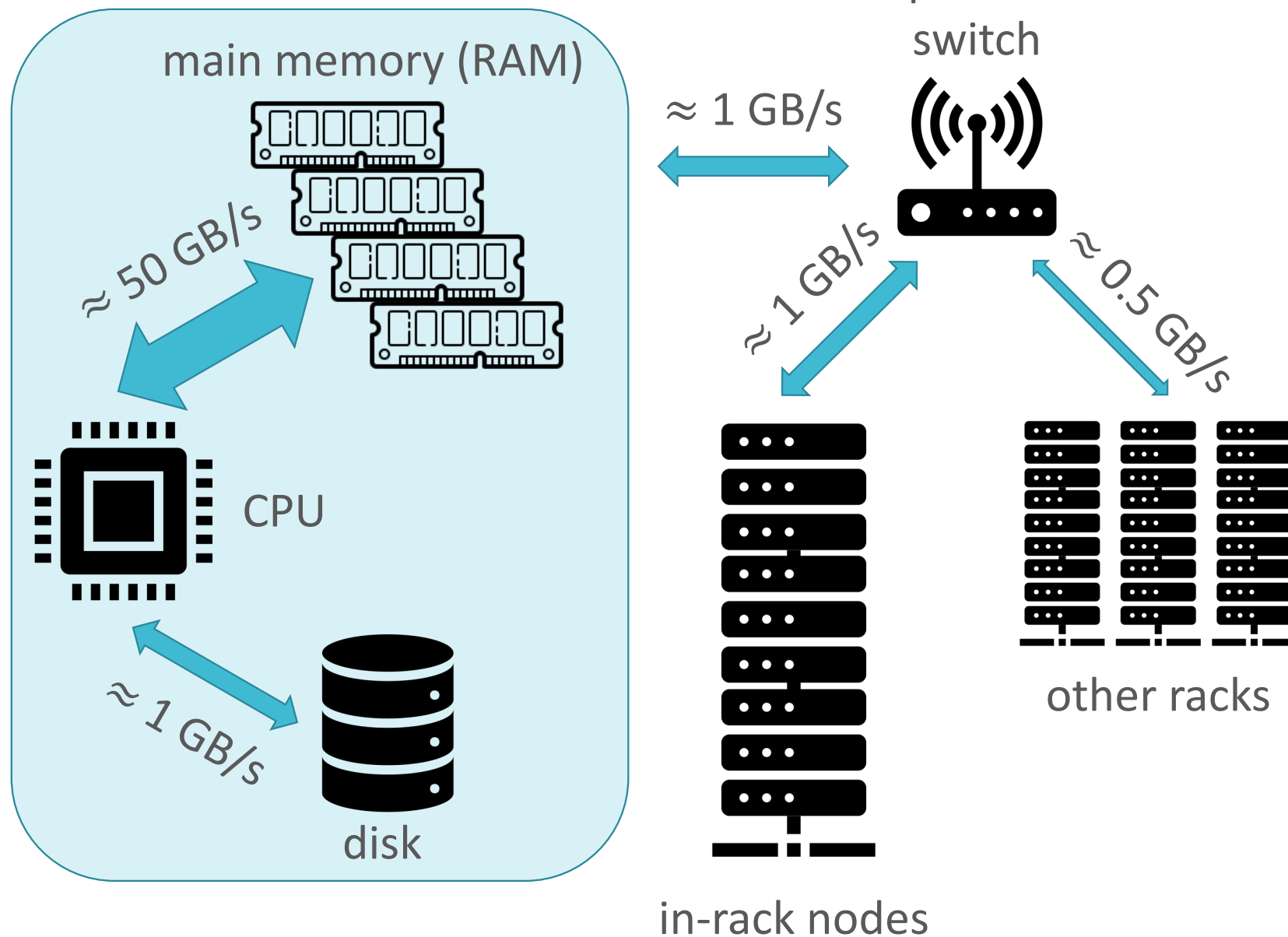
# Gradient Descent

- Pros:
  - Trivially parallelizable
  - Each individual iteration is cheap
    - Can be further improved using stochastic variants
  - Guaranteed to converge on convex objective functions
- Cons:
  - Potentially slow convergence
  - Introduction of a hyperparameter
  - Network communication in each iteration

# Recall: Communication Hierarchy



10-605/805  
Principle #2:  
Perform parallel  
and in-memory  
computation  
whenever  
possible



# 10-605/805

## Principle #2: Perform parallel and in-memory computation whenever possible

- Persisting data in-memory reduces communication, especially for iterative procedures

```
trainData.cache()
for t in range(num_iters):
    alpha_t = alpha / n * sqrt(t)

    grad = trainData.map(compute_pointwise_grads(w)).sum()

    w -= alpha_t * grad
```

# 10-605/805


## Principle #3: Minimize network communication

- Inherently at odds with Principle #2 → need to tradeoff between parallelism and network communication
- Three types of objects that may need to be communicated:
  - Data
  - Models
  - Intermediate objects
- Strategies:
  - Keep large objects local
  - Reduce the number of iterations

Workers	$\begin{bmatrix} \leftarrow & \mathbf{x}^{(1)T} & \rightarrow \\ \leftarrow & \mathbf{x}^{(4)T} & \rightarrow \\ \vdots & \vdots & \vdots \end{bmatrix}$	$\begin{bmatrix} \leftarrow & \mathbf{x}^{(2)T} & \rightarrow \\ \leftarrow & \mathbf{x}^{(3)T} & \rightarrow \\ \vdots & \vdots & \vdots \end{bmatrix}$	$\begin{bmatrix} \leftarrow & \mathbf{x}^{(5)T} & \rightarrow \\ \leftarrow & \mathbf{x}^{(7)T} & \rightarrow \\ \vdots & \vdots & \vdots \end{bmatrix}$	$O(nk)$ distributed storage (total)	
Map	$\mathbf{x}^{(i)} \mathbf{x}^{(i)T}$	$\mathbf{x}^{(i)} \mathbf{x}^{(i)T}$	$\mathbf{x}^{(i)} \mathbf{x}^{(i)T}$	$O(nk^2)$ distributed work (total)	$O(k^2)$ local storage
Reduce	$\left( \sum_{i=1}^n \mathbf{x}^{(i)} \mathbf{x}^{(i)T} \right)^{-1}$			$O(k^3)$ local work	$O(k^2)$ local storage

# Data Parallel: Compute outer products locally





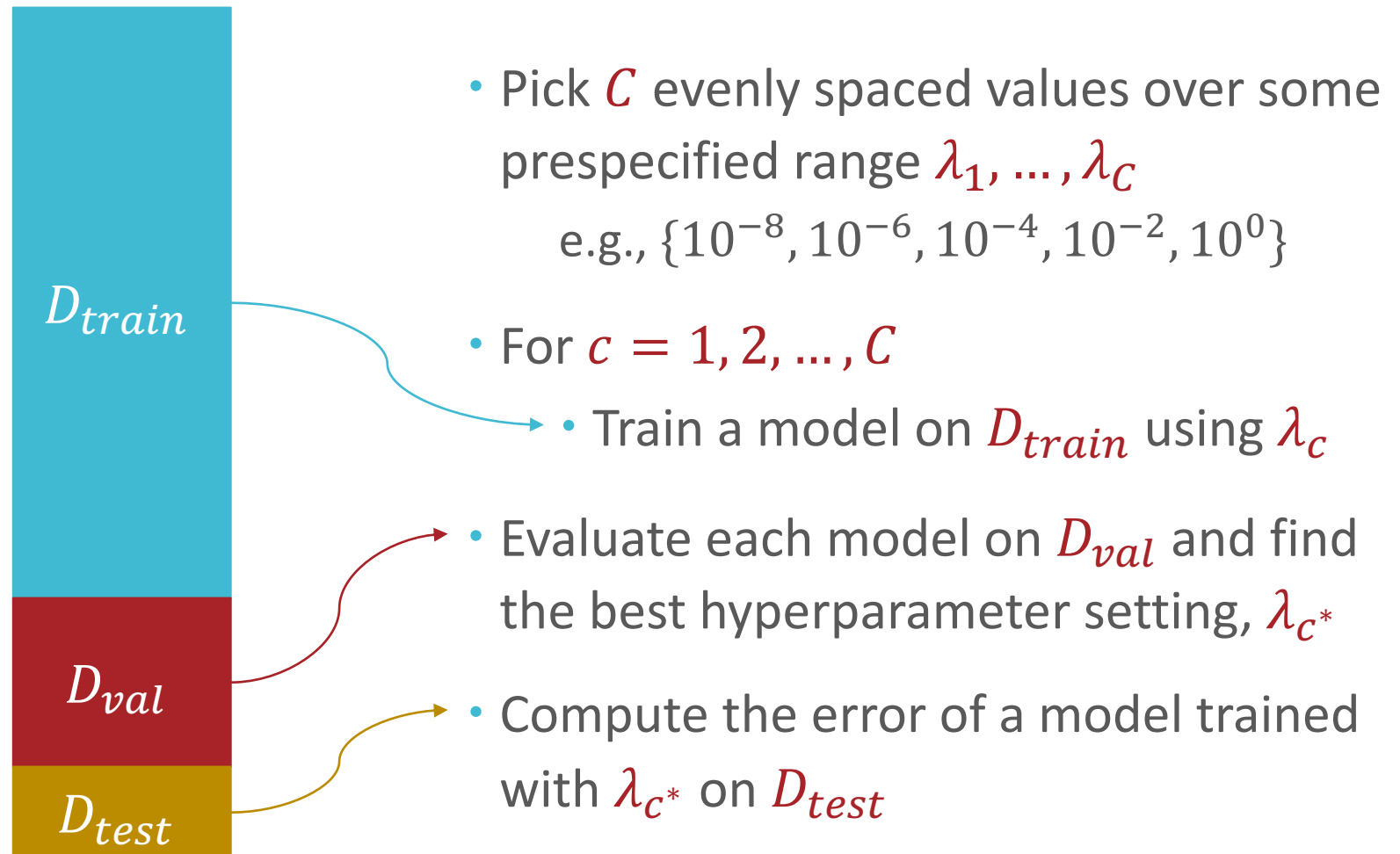
Workers	$\begin{bmatrix} \leftarrow & \mathbf{x}^{(1)T} & \rightarrow \\ \leftarrow & \mathbf{x}^{(4)T} & \rightarrow \\ \vdots & \vdots & \vdots \end{bmatrix}$	$\begin{bmatrix} \leftarrow & \mathbf{x}^{(2)T} & \rightarrow \\ \leftarrow & \mathbf{x}^{(3)T} & \rightarrow \\ \vdots & \vdots & \vdots \end{bmatrix}$	$\begin{bmatrix} \leftarrow & \mathbf{x}^{(5)T} & \rightarrow \\ \leftarrow & \mathbf{x}^{(7)T} & \rightarrow \\ \vdots & \vdots & \vdots \end{bmatrix}$	$O(nk)$ distributed storage (total)	
Map	$(\mathbf{w}^{(t)T} \mathbf{x}^{(i)} - y^{(i)}) \mathbf{x}^{(i)}$	$(\mathbf{w}^{(t)T} \mathbf{x}^{(i)} - y^{(i)}) \mathbf{x}^{(i)}$	$(\mathbf{w}^{(t)T} \mathbf{x}^{(i)} - y^{(i)}) \mathbf{x}^{(i)}$	$O(nk)$ distributed work (total)	$O(k)$ local storage
Reduce	$\mathbf{w}^{(t+1)} \leftarrow \mathbf{w}^{(t)} - \frac{\alpha}{n\sqrt{t}} \sum_{i=1}^n (\mathbf{w}^{(t)T} \mathbf{x}^{(i)} - y^{(i)}) \mathbf{x}^{(i)}$			$O(k)$ local work	$O(k)$ local storage

Issue: all workers must have the latest weight vector

# Data Parallel: Compute pointwise gradients locally

# Model Parallel: Train each hyperparameter setting on different machine(s)

- Hyperparameter optimization



# Key Takeaways

- 10-605/805 Principles:
  1. Computation and storage should be linear in  $n$  and  $k$ 
    - For linear regression:
      - When  $k$  is small, distribute covariance matrix computation using outer products
      - When  $k$  is large, minimize squared error via distributed gradient descent
  2. Perform parallel and in-memory computation whenever possible
  3. Minimize network communication
    - Data vs model parallelism