

CS 260: Foundations of Data Science

Prof. Sara Mathieson

Spring 2025



Haverford
COLLEGE

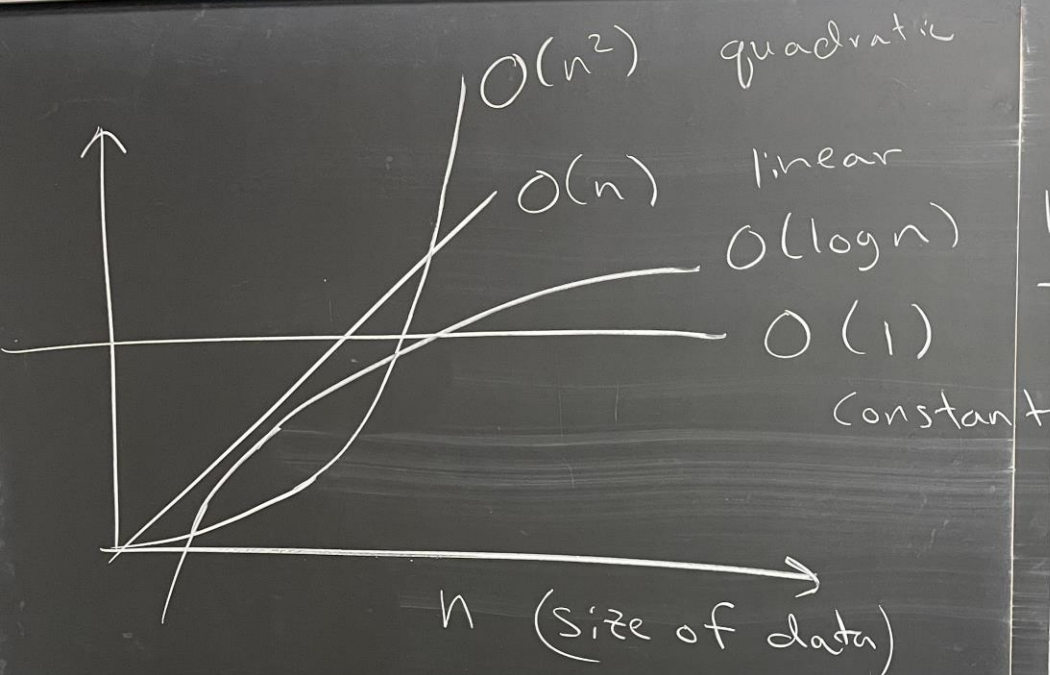
- Lab 5 released (partners optional), due after Spring Break
- Lab 3 grades will be up soon
- No class next Thursday (I'm at a conference)
- **Midterm 1: next Tuesday in class**
 - Study sheet created by 1 (one page front and back)
 - No other notes or resources

Outline

- Review runtime
- Handout 12
- Finish algorithmic bias
- Discussion: admissions at Haverford

Outline

- Review runtime
- Handout 12
- Finish algorithmic bias
- Discussion: admissions at Haverford



$$O(n) + O(\log n) \Rightarrow O(n)$$

logarithmic

$$O(\underset{\substack{\uparrow \\ \text{constant}}}{c}n) \Rightarrow O(n)$$

$$O(5n + 11n) \Rightarrow O(n)$$

$$O(n^k + n^j) = O(n^{\max(k, j)})$$

$$n^5 + n^7 \Rightarrow O(n^7)$$

① work inside out
(start with inner loop)

② what is changing?

③ what is constant?

④ log (see halving)

⑤ exp (see doubling)

Outline

- Review runtime
- **Handout 12**
- Finish algorithmic bias
- Discussion: admissions at Haverford

① $O(n)$

② $O(1)$

③ $O(n)$

④ $O(n^2)$

⑤

$$n + (n-1) + (n-2) + \dots + 3 + 2 + 1$$

$$\Rightarrow (n+1) \frac{n}{2} \Rightarrow O(n^2)$$



$$\frac{n^2}{2} \Rightarrow O(n^2)$$

$$(6) \quad O(n)$$

$$(7) \quad O(\log n)$$

$$(9) \quad O(n \log n)$$

$$(10) \quad O(2^n)$$

$$\begin{aligned} & \rightarrow 1 + 2 + 4 + 8 + \dots + 2^{n-1} \\ & 1 + \cancel{2} + \cancel{2^2} + \cancel{2^3} + \dots + \cancel{2^{n-1}} = S \\ & (-) \quad \cancel{2} + \cancel{2^2} + \cancel{2^3} + \dots + \cancel{2^{n-1}} + 2^n = 2S \\ & \hline & 1 - 2^n = S - 2S \\ & S = 2^n - 1 \\ & \Rightarrow O(2^n) \end{aligned}$$

⑧ $lst = [\dots]$

$n = \text{len}(lst)$

$\Rightarrow O(1)$ ★

\Rightarrow basic linked list
implementation

$\Rightarrow O(n)$
(see drawing)

can't assume

$p \neq n$ are similar

$O(1)$ ★

$$\textcircled{11} \quad \vec{w} = \underbrace{\left(\underbrace{X^T X}_{(a)} \right)^{-1}}_{(d)} \underbrace{X^T \vec{y}}_{(c)}$$

(b) $(p+1) \times (p+1) + (p+1) \times n$

(c) $X^T \Rightarrow O(np)$

(a) multiply $(p+1) \times n + n \times (p+1) \Rightarrow O(p^2 n)$

(b) invert $(p+1) \times (p+1)$ matrix $\Rightarrow O(p^3)$
inverse is cubic

(c) multiply $(p+1) \times n + (n \times 1) \Rightarrow O(pn)$

(d) multiply $(p+1) \times (p+1) + (p+1) \times 1 \Rightarrow O(p^2)$

$\Rightarrow O(p^2 n + p^3 + pn + p^2) \Rightarrow \boxed{O(p^2 n + p^3)} \star$

$\textcircled{8} \quad |st$

$n =$

$\Rightarrow \boxed{O(n^3)}$
 \Rightarrow bas
 imp

ca

matrix multiplication

$$\left. \begin{array}{l} A \Rightarrow (n, p) \\ B \Rightarrow (p, m) \end{array} \right\} \Rightarrow O(nmp)$$

$$\left. \begin{array}{l} C \Rightarrow (p, p) \\ D \Rightarrow (p, m) \end{array} \right\} \Rightarrow O(p^2 m)$$

inverse: matrix must be square

$$\Rightarrow (n \times n)$$

$$\Rightarrow \boxed{O(n^3)}$$

Outline

- Review runtime
- Handout 12
- **Finish algorithmic bias**
- Discussion: admissions at Haverford

How can we tell if an algorithm is biased?

D: dataset with attributes X , Y

- * X is protected
- * Y is unprotected (other features)

How can we tell if an algorithm is biased?

D: dataset with attributes X , Y

- * X is protected
- * Y is unprotected (other features)

Goal: determine outcome C (hired, admitted, etc)

How can we tell if an algorithm is biased?

D: dataset with attributes X , Y

- * X is protected
- * Y is unprotected (other features)

Goal: determine outcome C (hired, admitted, etc)

Direct discrimination: $C = f(X)$

- * Female instrumentalist not hired for orchestra
- * Some ethnic groups not allowed to eat at a restaurant

How can we tell if an algorithm is biased?

D: dataset with attributes X , Y

- * X is protected
- * Y is unprotected (other features)

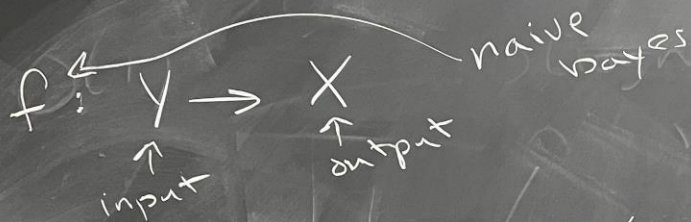
Goal: determine outcome C (hired, admitted, etc)

Indirect discrimination: $C = f(Y)$

- * but strong correlation between X and Y
- * Ex: housing loans
- * Ex: programming experience

Idea: if we can predict X from Y ,
could be disparate impact →

Predictor



metric

Balanced error rate (BER)

$$P(C=1|X=0) \leq 0.8 P(C=1|X=1)$$

hired
admitted

$C=0$ not hired

$$\epsilon = \text{BER} = \frac{1}{2} \left[P(f(Y)=0|X=1) + P(f(Y)=1|X=0) \right]$$

want high! i.e. confusion.

$$\frac{10}{1000}$$

$$\frac{25}{50}$$

① train a classifier $f(y) \rightarrow x$
(predict x)

② if BER is too low,
could be disparate impact

Normal error $\Rightarrow \frac{10 + 25}{1000 + 50} \Rightarrow \text{low}$

$$\text{BER} = \left(\frac{10}{1000} + \frac{25}{50} \right) \frac{1}{2}$$

\Rightarrow high

Outline

- Review runtime
- Handout 12
- Finish algorithmic bias
- **Discussion: admissions at Haverford**

Discussion: admissions at Haverford

- Haverford has suddenly started receiving 10x more applications than usual
- You are tasked with creating an algorithm to determine whether or not an applicant should be admitted
- Questions:
 - How would you encode features?
 - How would you use past admission data to train?
 - What cost function are you trying to optimize?