CS 331, Fall 2024    Today: - Kevin's research

Lecture 26  (12/9)
         - Fine-grained complexity
         - Popular conjectures

# Kevin's research

Algorithmic primitives for "big" data science

$\approx 1/2$ : continuous algos foundations (opt, samp, NLA)
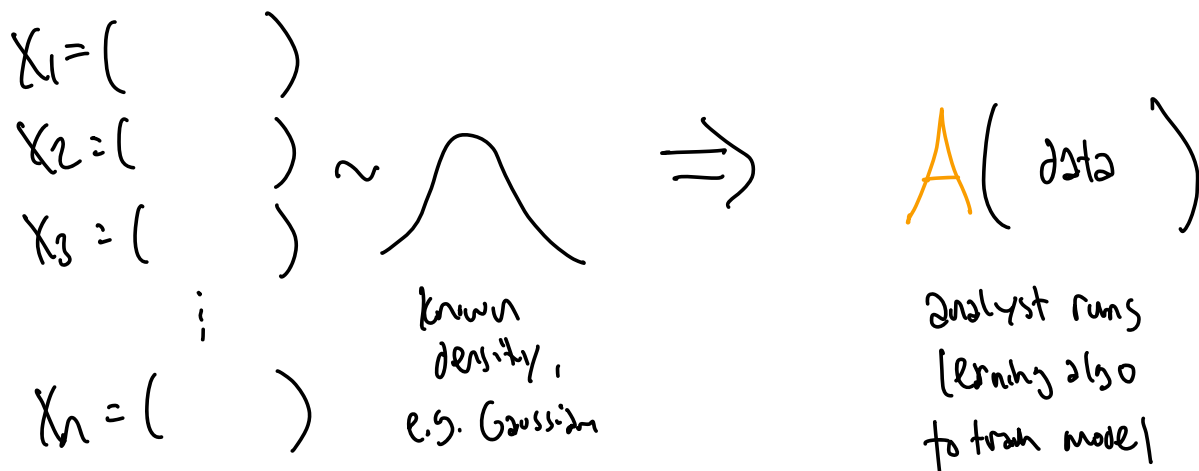
$\approx 1/2$ : trustworthy ML (robustness, privacy, fairness)
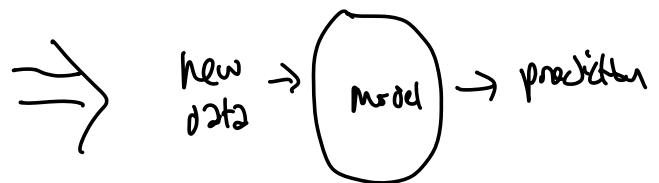
Key themes:

- Many problems hard in worst case. <u>Harness structure.</u>

- Where does the data come from? <u>Minimal assumptions.</u>

- Dataset sizes enormous, only growing. <u>Near-linear time?</u>

## Trustworthy ML

**Textbook setting** for statistics / learning:

$$X_1 = (\quad\quad)$$
$$X_2 = (\quad\quad)$$
$$X_3 = (\quad\quad)$$
$$\vdots$$
$$X_n = (\quad\quad)$$

$\sim$ known density, e.g. Gaussian

$\Rightarrow \quad A(\text{data})$

analyst runs learning algo to train model

$\Rightarrow$ new data $\to$ ( model ) $\to$ prediction

**Real life is harder.**

- What if we're wrong about the data? <u>Robustness</u>

- The data is coming from humans. <u>Privacy / fairness</u>

- Why do we believe the model's conclusion? <u>Interpretability</u>

- ... all needs to happen efficiently ...

## Continuous algos

What kind of tools are useful for modern algo design?

OPT: Minimize structured objectives.

e.g. minimax/stochastic optimization

Semidefinite programming ("matrix LP")

Structured nonconvex problems (sparsity, GLM, ...)

SAMP: Sample from structured densities.

e.g. logconcave sampling (basic tractable family)

Structured multimodal problems

NLA: numerical linear algebra primitives

e.g. preconditioning (solving linear systems, regression)

Sparsification (replace data w/ representatives)

2014: Google internship. Not good at it...

2015: Complexity research. Not good at it...

2016: Genomics research. Really fun! I liked algos best...

2017: Genomics / NLP / stats research. (Ph.D. rotations)

2018: Approximate maxflow.

2019: Nash equilibria, optimal transport, SDP.

2020: Sampling. SOTA for some logconcave families.

2021: Robust stats. PCA, regression, clustering in near-linear time.

2022-24: Privacy, interpretability, etc.
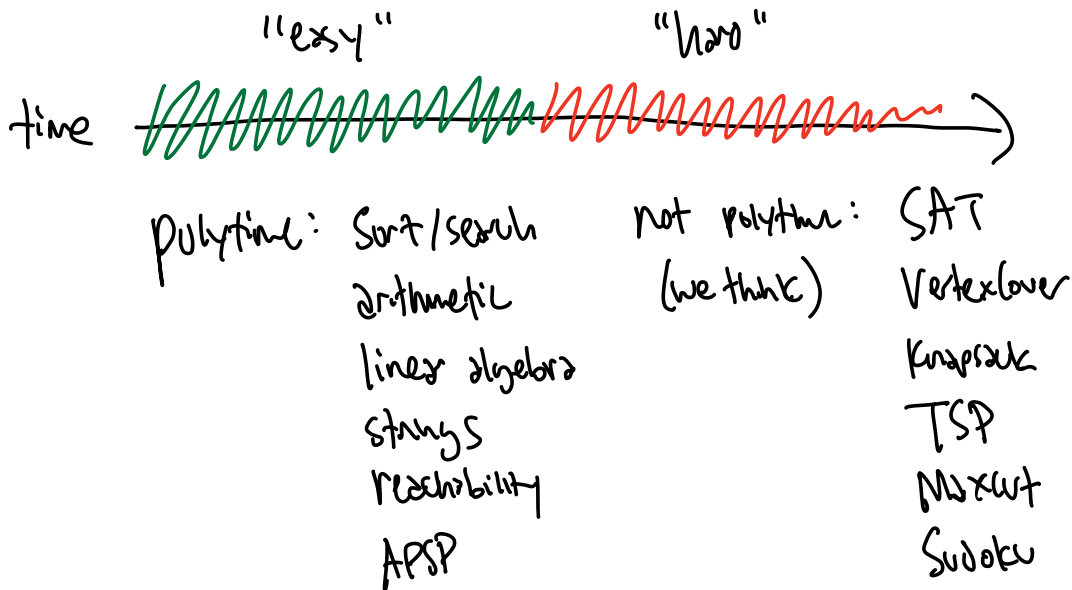
I am trying to learn more about modern ML...

Algorithms are cool and come in many flavors. There are so many connections.
Just keep learning and enjoy ☺

# Fine-grained Complexity

So far: Complexity theory for "small data"

$$n = 100, 1000, 10000 \ldots$$

"easy"          "hard"

time

polytime: Sort/search     not polytime: SAT
            arithmetic           (we think)     VertexCover
            linear algebra                    Knapsack
            strings                             TSP
            reachability                      Maxcut
            APSP                             Sudoku

It's 2024. We need complexity theory for "big data"

$$n = 10^9, 10^{12}, 10^{15} \ldots$$

"easy"          "medium"          "hard"

time

near-linear      polytime\         not polytime
time              near-linear

NP: a tool to prove problems hard.

Today: how to prove problems medium.

| Known easy | Suspected medium |
|---|---|
| FFT | 3-SUM |
| Shortest path | Allpairs shortest paths |
| Maxflow | Dynamic maxflow |
| Longest increasing subseq. | Longest common subseq. |
| Closest pair in $\mathbb{R}^2$ | Closest pair in $\mathbb{R}^d$ |
| Longest palindromic substring | Edit distance |

e.g. APSP     $n^3$    ...    ...    $n^{3-o(1)}$

Floyd-Warshall '62          Williams '14

Why are we so good at problems on LHS
   ... but bad at problems on RHS?

Goal in FGC: web of reductions, Common source of hardness
                                      must attack first!

# Popular Conjectures

Let $\varepsilon > 0$ be small constant.

3-SUM: Given list $L$ of #'s, $\exists\ a, b, c \in L$

$$\text{s.t. } a + b + c = 0?$$

... cannot be solved in time $O\left(|L|^{2-\varepsilon}\right)$

APSP: Given graph $G = (V, E, w)$ compute

$|V| \times |V|$ matrix encoding all-pairs shortest paths

... Cannot be solved in time $O\left(|V|^{3-\varepsilon}\right)$

SETH: $\exists$ constant $k$ s.t. $k$-SAT on

$\Phi$ w/ $m$ clauses, $n$ variables

... cannot be solved in time $O\left(2^{n(1-\varepsilon)} \text{poly}(m)\right)$

Rough intuition: we care about the _exponent_ now.

## 3-SUM

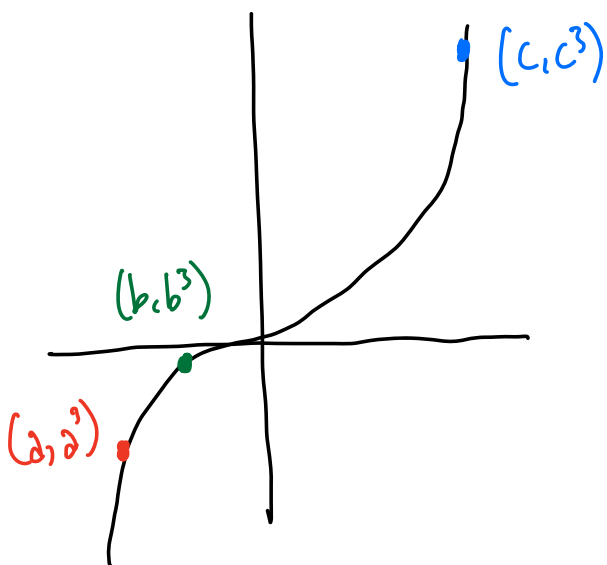Computational geometers Gajentaan–Overmars '95

Kickoff field of FGC by reducing to 3-SUM.

Example: Collinearity requires $\approx n^2$ time.

Input: $n$ points in $\mathbb{R}^2$. $\exists$ three on a line?

(This proof is crazy.) Reduce 3-SUM to collinearity.

Check collinearity $\left( \{ (x, x^3) \mid x \in L \} \right)$



$$\frac{b^3 - a^3}{b - a} = \frac{c^3 - b^3}{c - b}$$

$$a^2 + ab + b^2 = c^2 + cb + b^2$$

$$a^2 - c^2 = b(c - a)$$

$$-a - c = b \quad (3\text{-SUM})$$

More: visibility /reachability ⎫
      motion planning          ⎬ Several require
      polygon containment      ⎭ ingenuity to solve $\tilde{n} \approx n^2$.

## APSP

This is really about "Combinatorial" matrix multiplication.

Many amazing "algebraic" innovations:

$(n \times n)$ $(n \times n)$  takes time...   $n^3$ (duh)
     ↑          ↑                              $n^{2.807}$ Strassen
       multiply                                   ⋮
                                               $n^{2.3714}$ ADWXXZ

Use crazy cancellations on huge tensors...

What if we can only use more "baseline" techniques?

Recall <u>divide-and-conquer</u> + <u>DP</u> algo for APSP

$DP[s][t][\ell]$ : shortest $s-t$
path w/ $\leq 2^{\ell}$ edges

$$= \min_{u \in V} DP[s](u)[\ell-1]$$
$$+ \ DP(u)[t][\ell-1]$$

"guess the midpoint"



$\leq 2^{\ell-1}$  $\leq 2^{\ell-1}$

This is the same problem as "dot product"

except    sum $\longrightarrow$ min

   product $\longrightarrow$ sum

$$DP_{\ell-1} \otimes DP_{\ell-1} = DP_{\ell} \quad \forall \ell \in (O(\log(n)))$$

"min-plus" convolution

To solve APSP, need to solve combinatorial matmul $O(\log(n))$ times. Thus, APSP conjecture:

Combinatorial matmul needs $\approx n^3$ time :(

Good news: Structured matmul / inversion /... easier!

## SETH

Recall 3-SAT solvable in: $O(2^n m)$ time
Improvement: Try 7/8 for one clause, recurse.
$T(n) \leq 7 T(n-3) + O(m) \Rightarrow T(n) = O(1.913^n m)$

So, 3-SAT does not need $2^n$ time.

More general: k-SAT in $2^{n(1-O(\frac{1}{k}))}$ m time.

But we need _tight base_ (will later choose $n = \log$)
base becomes exponent.

Hence SETH: choose large enough k.

Almost all modern reductions use SETH thru:

$$\text{SETH} \le \text{OV} \qquad \text{Williams, 2005}$$
(orthogonal vectors)

OV implies FGC of so many problems:

- Diameter
- Dynamic reachability
- Single-source maxflow
- Edit distance
- Fréchet distance
- LCS
- Local alignment
- Stable matching
- Closest pair

2-OV problem:

Let $d = \omega(\log(n))$    "sparse subset"

$A, B \subset \{0,1\}^d$, size $n$

$\exists a \in A, b \in B$ s.t. $\underbrace{a^T b = 0}$ ?
                                      orthogonal
                                      vectors

Conjecture: no better than $\approx n^2 d$ possible.

k-OV: there are k sets

$A_1 \ A_2 \ A_3 \ \ldots \ \subset \{0,1\}^d$

$\exists a_1 \in A_1, \ a_2 \in A_2, \ a_3 \in A_3, \ldots$

s.t. $\sum_{i \in [d]} a_1[i] \, a_2[i] \, a_3[i] \ldots = 0$ ?

Conj: needs $\approx n^k d$ time.

Obs 1: 2-OV is very fundamental. FGC!

Obs 2: $2\text{-OV} \gtrsim 3\text{-OV} \gtrsim \dots \gtrsim k\text{-OV}$.

Obs 3: "$OV \gtrsim SETH$".

Suppose there's $k\text{-OV}$ in $O(N^{k(1-\varepsilon)})$.

Take $k$-SAT formula $\underline{\Phi}$, $n$ variables
$m$ clauses

Create $A_1, \dots, A_k \subset \{0,1\}^m$:

$$X \in \{0,1\}^n \rightarrow (\underbrace{x_1}_{n/k} \mid \underbrace{x_2}_{n/k} \mid \dots \mid \underbrace{x_k}_{n/k})$$

blocks

$A_i =$ index by $N = 2^{n/k}$ assignments to $i^{th}$ block

Faster $\Rightarrow$ SAT in time $N^{k(1-\varepsilon)} = 2^{n(1-\varepsilon)}$
$k$-OV
(violates SETH!)