Matrix Completion comes from the Netflix Prize '07-'09
Given ratings of movies by people (person, movie, score)
Try to predict ratings by those people on other movies.

Data was very sparse
  $\approx$ 500k people, 18k movies, on average 200 ratings/person
And, very irregular: many people rate few movies.
                     some movies get few ratings

Idea: Consider matrix people $\begin{pmatrix} & \text{ratings} & \end{pmatrix}$

— movies —

Knowing only a few entries in the matrix,
  can we guess the rest.

Need to assume matrix has some structure,
  like approximately low rank (low-rank + noise).

Matrix Completion problem: given $m \times n$ matrix $M$,
  $\Omega \subseteq \{1..m\} \times \{1..,n\}$ and an integer $r$
    $\min_X \sum_{(i,j) \in \Omega} (M(i,j) - X(i,j))^2$ s.t. rank$(X) \leq r$.

or, given $\varepsilon > 0$,
    $\min_X$ rank$(X)$ s.t. $\|M - X\|_\Omega^2 \leq \varepsilon$

Both forms are NP-hard, even with $\varepsilon = 0$.

So, replace rank($X$) with something convex.
If $\sigma_1, ..., \sigma_m$ are singular values of $X$,
$$\text{rank}(X) = \|(\sigma_1, ..., \sigma_m)\|_0$$
So, by analogy with $\ell_1$, use the <u>nuclear norm</u> of $X$
$$\|X\|_* = \|(\sigma_1, ..., \sigma_m)\|_1 = \sum_i \sigma_i \qquad (\text{Fazel})$$

Contrast with $\quad \|X\|_F = \|(\sigma_1, ..., \sigma_m)\|_2 = \left(\sum \sigma_i^2\right)^{1/2}$

$$\|X\|_2 = \max_i \sigma_i$$

| vector | matrix |
|---|---|
| $\|\cdot\|_0$ | rank |
| $\|\cdot\|_\infty$ | $\|\cdot\|_2$ |
| $\|\cdot\|_2$ | $\|\cdot\|_F$ |
| $\|\cdot\|_1$ | $\|\cdot\|_*$ |

Exact recovery:
$$\min_X \|X\|_* \quad \text{s.t.} \quad \|M - X\|_\Omega = 0$$

Approximate recovery:
$$\min_X \|X\|_* \quad \text{s.t.} \quad \|M - X\|_\Omega \leq \varepsilon$$

Work well when $\Omega$ is random (big assumption)
$\quad X$ is nice: no big entries,
$\qquad\qquad$ singular vectors incoherent.
$\quad$ If $M$ is low rank, can actually find $M$.

# The nuclear norm    (also called trace norm)

Is the dual of the operator norm, $\|X\|_2$.

__Thm__  $\|X\|_* = \max\limits_{\|Y\|_2 = 1} Tr\left(X^T Y\right) = \sum\limits_{(i,j)} X(i,j) Y(i,j)$

__proof__  Let $X = USV^T$ be the SVD of $X$.
  If $Y = UV^T$
  $Tr\left(X^T Y\right) = Tr\left[VSU^TUV^T\right] = Tr\left(VSV^T\right)$
$$= Tr\left(SV^TV\right)$$
$$= Tr(S)$$
$$= \sum \sigma_i$$

  If $Y = \hat{U}\hat{V}^T$, where $\hat{U}, \hat{V}$ are orthogonal,
  $Tr\left(X^T Y\right) = Tr\left(VSU^T\hat{U}\hat{V}^T\right)$
$$= Tr\left(S \underbrace{U^T\hat{U}\hat{V}^TV}_{Q}\right)$$
$$= Tr(SQ) \text{ where } Q \text{ is orthogonal}$$
$$= \sum_i \sigma_i Q(i,i) \leq \sum_i \sigma_i$$

  ces  $Q(i,i) \leq 1$ for all $i$.

Now we know $\|X\|_*$ is convex

$$\|X\|_* = \left\{ \begin{array}{l} \min\limits_{W_1, W_2} \quad \frac{1}{2}\left(\text{Tr}(W_1) + \text{Tr}(W_2)\right) \\[2ex] \text{s.t.} \quad \begin{pmatrix} W_1 & X^T \\ X & W_2 \end{pmatrix} \succcurlyeq O \end{array} \right.$$

So, can minimize nuclear norms using semidefinite programming.

---

①riginal Story: approx low-rank by nuclear norm.

Alternate Story: Use complexity to guide problem and algorithm selection.
 Exactly solve tractable problem
 Approximately solve hard problem
 or just try to solve hard problem. (heuristic).

Alternating minimization:
 0. pick $L \in \mathbb{R}^{m \times r}$ at random
 1. Until $\|M - LR^T\|_\Omega$ is small
   $R = \arg\min\limits_{R \in \mathbb{R}^{n \times r}} \|M - LR^T\|_\Omega$

   $L = \arg\min\limits_{L \in \mathbb{R}^{m \times r}} \|M - LR^T\|_\Omega$

How hard are NP-hard problems?

For any, time less than $2^{n^\varepsilon}$ would be shocking

For most, time $\leq 2^{n^{1+\varepsilon}}$, any $\varepsilon$, would be surprising

Circuit SAT: time $\in 2^{(1-\varepsilon)n}$ unlikely.

Strong Exponential time hypothesis (SETH)
    $\Rightarrow$ lower bounds for many problems in P.

Are other notions of hard, from
    very hard — uncomputable, to pretty hard — Unique Games Conj.

When we can solve NP-hard problems,
    they are probably special:
        random
        structured
        structure + noise    (smoothed analysis)
        well conditioned
        other nice properties

Probably Certifiably Correct — Bandeira
    Random problems that can usually solve
        and certify optimality.

    Ex. relax, and certify by computing solution to dual
        of relaxation

$$\max_{\{\pm 1\}^n} \quad \leq \quad \max_{\|v_1\|=1,\dots,\|v_n\|=1} \quad \leq \text{ dual of that}$$

can provide certificate

In nice cases, can quickly solve dual from solution to original — complementary slackness.

---

Other notions of complexity :
   Information complexity of Black-Box algorithms.

See: Bubeck's monograph, Convex Optimization: Algorithms and Complexity, or Lectures on Modern Convex Optimization by Ben-Tal and Nemirovski

Consider an optimization algorithm that
   minimizes $f$ only using
      · evaluations of $f$
      · a few properties of $f$, like convexity.

Think of accessing $f$ through a black-box
      or an oracle.
Can prove lower bounds on the # of calls
      necessary

Can even allow evaluations of gradients of $f$
      (or sub-gradients : $f(y) \geq f(x) + g_x^T(y-x), \forall y$ )

For example, consider minimizing convex functions over $B(0,1) = \{x : \|x\|_2 \leq 1\} \subseteq \mathbb{R}^n$

Can prove that for every (deterministic) algorithm, $A$, that makes $k < n$ oracle calls, $\exists$ a convex $f$ s.t.
$$\|x_k - x_*\| \geq \text{const}/\sqrt{k}$$
where $x_* =$ minimum of $f$ and $x_k$ is output of algorithm.

Are many results like this for many function classes.
To prove faster convergence of an algorithm, must use more about the function.

---

Black-Box solvers usually arise in non-convex optimization.

These are heuristics people use to try to minimize or maximize functions.

I use them to find counter-examples.

For example, here's a conjecture:

Conj For all orthogonal matrices $Q$
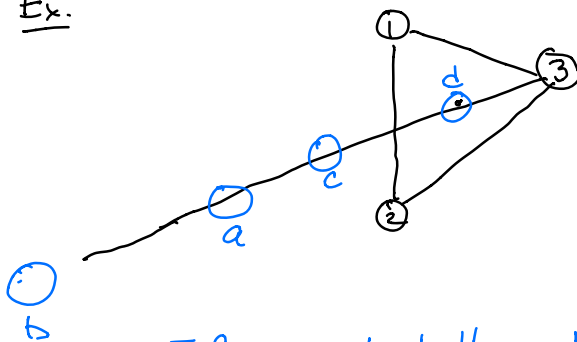$\exists x \in \{\pm 1\}^n$ such that $\|Qx\|_\infty \leq 2$.

It's false. I know because I found a $Q$.
(might be true with 2.1)

Typical approaches:
1. Gradient descent from many random starts.

2. Nelder - Mead. (the other simplex method)
keep $n+1$ vectors in $\mathbb{R}^n$, $x_1,\ldots,x_{n+1}$
adjust them until all approach the solution.

Order so that $f(x_1) \leq f(x_2) \leq \cdots \leq f(x_{n+1})$
$x_{n+1}$ is worst, so try to improve it by
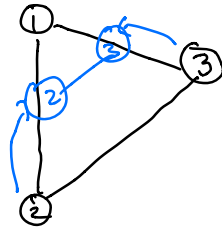moving on line through centroid $= \frac{1}{n+1} \sum_i x_i$

Ex.



Animate.

First try a.
If better, check if
b better than $x_1$

If a not better, try c and d

Get a smaller simplex.
Is the most popular black-box optimizer.

3. Line searches in random directions.

4. Simulated annealing: $g$ decreasing with $t$.
   $\hat{x} = x + \tau$, $\tau$ random   (usually discrete)
   if  $f(\hat{x}) \leq f(x) + g(t)$
       $x = \hat{x}$

5. Differential Evolution ⭐
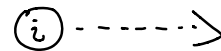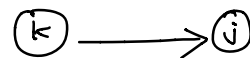   Generate a population $x_1, \ldots, x_{cn}$   $c > 0$
   Pick $i, j, k$ at random
      make sure $f(x_j) \leq f(x_k)$ (swap otherwise)
      $\hat{x} = x_i + (x_j - x_k)$
      if  $f(\hat{x}) \leq f(x_i)$
          $x_i = \hat{x}$

   $\underset{k}{\textcircled{k}} \longrightarrow \textcircled{j}$

   $\textcircled{i} \dashrightarrow$

Integer Linear Programming
$$\min c^T x \text{ st. } Ax \le b, \quad x \in \mathbb{Z}^n$$
$$\text{or } x \in \{0,1\}^n$$

Is NP-hard, but is good code
that often works in practice.
( CPLEX, Gurobi, Knitro, Mosek)

Might take a long time,
but can get certificates of optimality.

SAT solvers.
Solve problems like circuit satisfiability.
(usually SAT, which is slightly simpler)
Are competitions.

Conclusions of unsatisfiability are
key technology in program verification
and secure kernels.

SOS  (Sums of Squares)   <span style="color:magenta">Hot topic</span>

Any polynomial inequality like $f(\bar{x}) \geq 0$
has a proof  $h(x)^2 f(x) = \sum_i g_i(x)^2$

where $h$ and $g_i$ are polynomials.

Ex.  $x^2 + y^2 - 2xy \geq 0$
$x^2 + y^2 - 2xy = (x-y)^2$

Find them by Semidefinite Programming.

# Modern Machine Learning

Fit models to data by non-convex programming

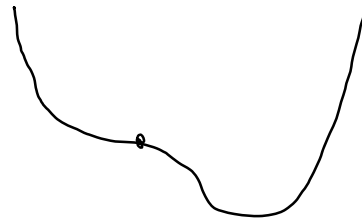Given vectors $x_1, ..., x_n$ and labels $y_1, ..., y_n$
try to choose parameters $\Theta$ (e.g. weights in neural net)
to minimize $F(\Theta) = \sum_i [f_\Theta(x_i) - y_i]^2$

or $F(\Theta) = \sum_i \ell(f_\Theta(x_i), y_i)$ some loss function $\ell$.

Often use gradient-based methods to find
local minima.
Need to escape saddle points!

Many theorems for special problems
Improvements in Black-Box algorithms, assuming more.

But, are often many local minima,
different algorithms produce different $\Theta$.
$\Theta$ MATTERS MORE THAN $F(\Theta)$