# 1   Introduction

If $A$ is an $m \times n$, real-valued matrix recall that the *discrepancy* of $A$ is

$$\operatorname{disc}(A) = \min_{x \in \{\pm 1\}^n} \|Ax\|_\infty.$$

We have seen in the last lecture that matrix discrepancy generalizes combinatorial discrepancy (if $\mathcal{S}$ is a set system with $m$ sets, take $A$ to be the 0-1 incidence matrix of $\mathcal{S}$, and the $\pm 1$ vector $x$ is the colouring of the underlying elements), which means that the study of $\operatorname{disc}(A)$ is particularly interesting when the entries of $A$ are bounded in the interval $[-1, 1]$. In this regime we can get an upper bound of $\operatorname{disc}(A) = O(\sqrt{n \log m})$ by taking a uniformly random vector $x \in \{\pm 1\}^n$.

**Exercise 1.** *Prove that for any $A \in [-1, 1]^{m \times n}$ a uniformly random $x\{-1, 1\}^n$ has discrepancy $\|Ax\|_\infty = O(\sqrt{n \log m})$ with constant probability.*

The seminal "Six Standard Deviations Suffice" theorem improves this when $m$ is small:

**Theorem 1** (Six Standard Deviations Suffice [3]). *Let $A$ be any $m \times n$ matrix with entries bounded by $[-1, 1]$. Then $\operatorname{disc}(A) = O(\sqrt{n \log(m/n)})$. In particular, any set system $(U, \mathcal{S})$ of $m = |\mathcal{S}|$ sets over a universe of size $n = |U|$ has discrepancy $O(\sqrt{n \log(m/n)})$.*

This is especially interesting when $m = O(n)$ and we have $\operatorname{disc}(A) = O(\sqrt{n})$. Then the theorem shows that you can make the discrepancy of every set bounded by a constant times the standard deviation of a random coloring. With more careful analysis the constant can be shown to be smaller than 6, hence the name of the theorem, and of Spencer's famous paper.

In this lecture we give a constructive proof (by way of a simple randomized algorithm) due to Lovett and Meka [1] of the previous theorem.

**Exercise 2.** *For $n = 2^k$ for some positive integer $k$, let $H_n$ be the $n \times n$ Hadamard matrix, i.e. a matrix with $\pm 1$ entries so that any two of its rows are orthogonal. Show that its discrepancy is $\operatorname{disc}(H) = \Omega(\sqrt{n})$. Use this to show that there exists a set system with discrepancy $\Omega(\sqrt{n})$.*

# 2   The Algorithm

**Theorem 2.** *Let $m, n$ be positive integers with $m \geq n$ and let $A$ be any $m \times n$ matrix with entries from $[-1, 1]$. There is a randomized algorithm running in time polynomial in $m, n$ which, when given $A$ as input, outputs an $x \in \{\pm 1\}^n$ such that $\|Ax\|_\infty = O(\sqrt{n \log m/n})$ with high probability.*

In specifying and analyzing the algorithm we take a geometric viewpoint. Let $m, n$ be positive integers with $m \geq n$, and let $A$ be any $m \times n$ matrix with entries from the interval $[-1, 1]$. Let $C$ be a universal constant to be specified later, and consider the convex polytope

$$K = \{x \in \mathbf{R}^n \mid \|Ax\|_\infty \leq C\sqrt{n \log(8m/n)}\}.$$

If $a_1, a_2, \ldots, a_m$ are the rows of $A$ we can re-write the previous definition as

$$K = \{x \in \mathbf{R}^n \mid \forall i \in [m], |\langle a_i, x \rangle| \leq C\sqrt{n \log(8m/n)}\},$$

which will be more useful. Intuitively, $K$ is a convex polytope defined as the intersection of a set of "slabs" of the form

$$|\langle a_i, x \rangle| \leq C\sqrt{n \log(8m/n)}.$$

The algorithm will work as follows: we perform (an approximation of) a continuous random walk, starting from the origin, in the convex polytope $K \cap [-1, 1]^n$. When the random walk intersects a facet of $K \cap [-1, 1]^n$ we restrict further steps of the walk to remain on that facet, until we end up at a vertex of $K \cap [-1, 1]^n$. Ideally we would like this vertex to be a point in $\{-1, +1\}^n$, and this would be clearly enough to prove the theorem. This is too much to hope for, but we will show that the resulting vertex has a constant fraction of its coordinates in $\{-1, 1\}^n$.

We reduce Theorem 2 to the following theorem.

**Theorem 3.** *Let $m \geq n$ be positive integers and let $\delta = 1/\sqrt{n}$. There is a randomized, polynomial-time algorithm and a constant $C$ such that, when given an $m \times n$ matrix $A$ and a vector $x(0) \in [-1, 1]^n$, finds an $x \in [-1, 1]^n$ such that the following holds, with probability at least $1/6 - \varepsilon$ for any $1/6 > \varepsilon > 0$.*

1. *For each $i = 1, 2, \ldots, m$ we have $|\langle a_i, x - x(0) \rangle| \leq C\|a_i\|_2 \sqrt{\log(8m/n)}$*

2. *$|x_i| > 1 - \delta$ for at least $n/10$ indices $i$.*

*Proof of Theorem 2 from Theorem 3.* Start with $x(0) = \mathbf{0}$, run the algorithm from Theorem 3 and obtain an $x \in [-1, 1]^n$. Let $x'$ be the vector obtained by choosing all indices from $x$ for which (2) fails, and apply the algorithm recursively on $x(0) = x'$ and on the matrix $A'$ obtained by deleting the columns corresponding to the indices satisfying (2). At each recursive step we fix a constant fraction of the coordinates of $x$, and so we end up with a vector $x^*$ for which all indices satisfy (2) after $S = 10 \log n$ recursive steps. The discrepancy of the resulting vector is

$$\|Ax^*\|_\infty < C\sqrt{n}\sqrt{\log(8m/n)} + C\sqrt{n/10}\sqrt{\log(8m/(n/10))} + \cdots + C\sqrt{n/10^S}\sqrt{\log(8m/(n/10^S))}$$

$$< \sqrt{n} \sum_{s=0}^{\infty} \frac{C\sqrt{\log 8m \cdot 10^s/n}}{10^{s/2}} < C'\sqrt{n \log(m/n)}$$

for some constant $C'$. Finally, we round each coordinate in $x^*$ to the nearest integer. It is easy to see that this can change the discrepancy by at most $O(n/\delta) = O(\sqrt{n})$. $\square$

Let $\mathcal{N}(\mu, \sigma^2)$ denote the mean $\mu$ Gaussian distribution with variance $\sigma^2$. The algorithm is formally described in Algorithm 1. As stated, the algorithm includes several scalar parameters $\delta, \gamma, T, C$ that we fix during the analysis: for now, think of $\delta, \gamma$ as being small reals with, say, $1/\sqrt{n} \geq \delta \gg \gamma > 0$

---
**Algorithm 1:** Main Algorithm
---

    **Input**   : An $m \times n$ matrix $A$. A vector $x(0) \in [-1, 1]^n$.
    **Output**: A vector $x \in [-1, 1]^n$ satisfying the properties in Theorem 3.
    **for** $t = 1, 2, \ldots, T$ **do**
        |   Set $\mathcal{D}_t = \{i \in [m] \mid |\langle a_i, x(t-1) - x(0)\rangle| \geq C\|a_i\|_2\sqrt{\log(8m/n)} - \delta\}$;
        |   Set $\mathcal{V}_t = \{j \in [n] \mid |x_j(t-1)| > 1 - \delta\}$;
        |   Let $\mathcal{W}_t = \{y \in \mathbf{R}^n \mid \forall i \in \mathcal{D}_t, \langle a_i, y\rangle = 0 \text{ and } \forall j \in \mathcal{V}_t, y_j = 0\}$;
        |   Let $w_1, w_2, \ldots, w_k$ be an orthonormal basis of the subspace $\mathcal{W}_t$;
        |   Let $g_1, g_2, \ldots, g_k \sim \mathcal{N}(0, 1)$ be sampled i.i.d.;
        |   Set $\Delta x(t) = \sum_{i=1}^{k} g_i w_i$;
        |   Set $x(t) = x(t-1) + \gamma \Delta x(t)$;
    **end**
    **return** $x(T)$

---

and $T$ being some large integer on the order of $1/\gamma^2$. The set $\mathcal{D}_t$ contains the facets of $K$ for which the vector $x(t-1)$ is "almost tight", and the set $\mathcal{V}_t$ contains the set of facets of $[-1, 1]^n$ for which $x(t-1)$ is "almost tight". The subspace $\mathcal{W}_t$ contains all vectors orthogonal to the $\mathcal{D}_t$ facets and the $\mathcal{V}_t$ facets. In each iteration of the algorithm, we take the vector $x(t-1)$ and perturb it by Gaussian random noise in the subspace $\mathcal{W}_t$. By moving in the subspace $\mathcal{W}_t$, we never increase the discrepancy with respect to the facets in $\mathcal{D}_t$ and we never modify any coordinates that are sufficiently close to $\pm 1$.

The parameters $\gamma, \delta$ should be viewed as tolerance parameters that we must introduce since we are approximating a continuous random walk. The parameter $\delta$ defines a small region around the facets of $K \cap [-1, 1]^n$ which we use to define when a vector $x(t)$ is tight with respect to the facet. By choosing $\delta$ to be small enough we are guaranteed that the coordinates are close enough to the $\pm 1$ constraints so that we do not introduce too much extra discrepancy when rounding the fractional coordinates. The parameter $\gamma$ controls the step-size of our discretized walk — we choose a vector $\Delta x(t)$ of variance-1 Gaussian random noise, projected to the subspace $\mathcal{W}_t$, and make a $\gamma$-length step in that direction.

Gaussian random variables enjoy a number of useful properties (e.g. exponential tail bounds), but key to the analysis is the next property that states that a linear combination of samples of Gaussian noise is again Gaussian.

**Stability of Gaussians.** Let $g_1, g_2, \ldots, g_k$ be i.i.d. samples from $\mathcal{N}(0, 1)$, and let $g = (g_1, g_2, \ldots, g_k)$. Then $\langle a, g\rangle \sim \mathcal{N}(0, \|a\|_2^2)$ for any $a \in \mathbf{R}^k$.

To prove our Main Lemma we will also use following Azuma-type martingale concentration inequality.

**Lemma 4.** *Let $\sigma \in \mathbf{R}$ satisfy $0 < \sigma \leq \tau$. Suppose $y_1, y_2, \ldots, y_\ell$ are random variables where $y_1 \sim \mathcal{N}(0, \sigma^2)$, and for all $i > 1$, the conditional distribution of $y_i - y_{i-1}$ given the values of $y_1, y_2, \ldots, y_{i-1}$ is $\mathcal{N}(0, \sigma_i^2)$ for some random variable $0 < \sigma_i \leq \tau$ depending on $y_1, y_2, \ldots, y_{i-1}$. Then*

$$\Pr[|y_\ell| > t\tau\sqrt{\ell}] \leq 2e^{-t^2/2}$$

*for any $t > 0$.*

The main lemma now follows from the martingale concentration bound and the Stability of Gaussians. It says that if the random walk is not "too long" then, on average, the number of tight discrepancy constraints will be small.

**Lemma 5.** *If $T = O(1/\gamma^2)$ then there exists a constant $C$ such that*

$$\mathbb{E}\,|\mathcal{D}_{T+1}| = \mathbb{E}\,|\{i \in [m] \mid |\langle a_i, x(T) - x(0)\rangle| \geq C\|a_i\|_2\sqrt{\log(8m/n)}\}| \leq n/4.$$

*Proof.* Let $C$ be a constant that will be fixed later. By linearity of expectation we can write

$$\mathbb{E}\,|\mathcal{D}_{T+1}| = \sum_{i=1}^{m} \Pr[|\langle a_i, x^T - x(0)\rangle| \geq C\|a_i\|_2\sqrt{\log(8m/n)}].$$

We expand the inner product as

$$\langle a_i, x^T - x(0)\rangle = \sum_{t=1}^{(T)} \gamma\langle a_i, \Delta x(t)\rangle$$

where $\Delta x(t) = \sum_{j=1}^{k} g_j w_j$ for i.i.d. Gaussian samples $g_1, g_2, \ldots, g_k \sim \mathcal{N}(0,1)$. The Stability of Gaussians implies that

$$\gamma\langle a_i, \Delta x(t)\rangle = \gamma\sum_{j=1}^{k} g_j\langle a_i, w_j\rangle \sim \mathcal{N}(0, \sigma^2)$$

where $\sigma^2 = \gamma^2\sum_{j=1}^{k}\langle a_i, w_j\rangle^2 \leq \gamma^2\|a_i\|_2^2$ since the basis $w_1, \ldots, w_k$ is an orthonormal basis of a subspace. For each $t = 1, 2, \ldots, T$ let $y_t = \langle a_i, x(t) - x(0)\rangle$. It follows that the sequence of variables $y_1, y_2, \ldots, y_T$ satisfy the conditions of Lemma 4 with $\tau = \gamma\|a_i\|_2$, thus

$$\Pr[|\langle a_i, x(T) - x(0)\rangle| > t\gamma\|a_i\|_2\sqrt{T}] = \Pr[|y_T| > t\gamma\|a_i\|_2\sqrt{T}] \leq 2e^{-t^2/2}.$$

Since $T = O(1/\gamma^2)$, choosing $t = C\sqrt{\log 8m/n}$ and $C$ any constant such that $C \geq 1/\gamma\sqrt{T}$ yields

$$\Pr[|\langle a_i, x(T) - x(0)\rangle| > C\|a_i\|_2\sqrt{\log 8m/n}] \leq 2e^{-\log 8m/n} = \frac{n}{4m}.$$

By summing this inequality over all $i \in [m]$ we get $\mathbb{E}\,|\mathcal{D}_{T+1}| \leq n/4$. $\qquad\square$

With this lemma we can prove Theorem 3.

*Proof of Theorem 3.* Let $x = x(T)$ be the output of Algorithm 1. It is not hard to show that for $\gamma$ much smaller than $\delta$, $x(t) \in K \cap [-1,1]^n$ for all $t$ with high probability. We leave this detail as an exercise. This immediately gives that $x(T) \in [-1,1]^n$ and that the first propety in Theorem 3 is satisfied.

It remains to prove the second property in Theorem 3. We do so by estimating $\mathbb{E}\,|\mathcal{V}_T|$ and using Markov's inequality. By the definition of the algorithm, $\mathbb{E}\,\|x(T) - x(0)\|_2^2 \leq n$ and $x(T) - x(0) = \sum_{t=1}^{(T)} \gamma\Delta x(t)$. For any $t = 1, 2, \ldots, T$, if $w_1, w_2, \ldots, w_k$ is the orthonormal basis of $\mathcal{W}_t$ we have

$$\mathbb{E}\,\|\Delta x(t)\|_2^2 = \mathbb{E}\langle\sum_{j=1}^{k} g_j w_j, \sum_{j=1}^{k} g_j w_j\rangle = \sum_{j=1}^{k}\mathbb{E}\,g_j^2 = k = \dim\mathcal{W}_t \geq n - \mathbb{E}\,|\mathcal{D}_t| - \mathbb{E}\,|\mathcal{V}_t|,$$

where $\mathbb{E}\, g_j^2 = 1$ for any $j$ since the Gaussian samples have variance 1. Since the Gaussian samples used by the algorithm are independent and have mean 0, we have

$$n \geq \mathbb{E}\, \|x(T) - x(0)\|_2^2 = \mathbb{E}\langle \sum_{t=1}^{T} \gamma \Delta x(t), \sum_{t=1}^{T} \gamma \Delta x(t)\rangle$$

$$= \gamma^2 \sum_{t=1}^{T} \langle \Delta x(t), \Delta x(t)\rangle = \gamma^2 \sum_{t=1}^{T} \mathbb{E}\, \|x(t)\|_2^2.$$

For each $t = 1, 2, \ldots, T$ we have $|\mathcal{D}_t| \leq |\mathcal{D}_{t+1}|$ and $|\mathcal{V}_t| \leq |\mathcal{V}_{t+1}|$, since whenever the random walk is tight to a facet in $\mathcal{D}_t$ or $\mathcal{V}_t$ it remains tight to that facet for all further steps. Using this fact we continue the calculation:

$$n \geq \gamma^2 \sum_{t=1}^{T} \mathbb{E}\, \|\Delta x(t)\|_2^2 \geq \gamma^2 \sum_{t=1}^{T} n - \mathbb{E}\,|\mathcal{D}_t| - \mathbb{E}\,|\mathcal{V}_t| \geq \gamma^2(T)(n - \mathbb{E}\,|\mathcal{D}_T| - \mathbb{E}\,|\mathcal{V}_T|).$$

Choose $T = 2/\gamma^2$ and using the fact that $|\mathcal{D}_T| \leq |\mathcal{D}_{T+1}|$ rearrange to get

$$\mathbb{E}\,|\mathcal{V}_T| \geq \frac{1}{2}(n - 2\,\mathbb{E}\,|\mathcal{D}_T|) \geq n/4.$$

Applying Markov's inequality to the random variable $n - |\mathcal{V}_T|$ we get

$$\Pr[n - |\mathcal{V}_T| > 9n/10] = \Pr[n/10 > |\mathcal{V}_T|] \leq \frac{n - \mathbb{E}\,|\mathcal{V}_T|}{9n/10} \leq \frac{3n/4}{9n/10} = \frac{5}{6}$$

and thus $|\mathcal{V}_T| \geq n/10$ with probability at least $1/6$.

It is easy to verify that the algorithm runs in polynomial time. $\qquad\square$

In fact the proof above gives the following powerful statement which can be used to give many other discrepancy upper bounds. Some applications are indicated in the exercises that follow.

**Theorem 6.** *Let $m \geq n$ be positive integers and let $\delta = 1/n$. Let $\lambda_1, \ldots, \lambda_m \geq 0$ be such that $\sum_{i=1}^{m} e^{-\lambda_i^2/2} \leq \frac{n}{4}$. There is a randomized, polynomial-time algorithm such that, when given an $m \times n$ matrix $A$ and a vector $x(0) \in [-1,1]^n$, finds an $x \in [-1,1]^n$ such that the following holds, with probability at least $1/6 - \varepsilon$ for any $1/6 > \varepsilon > 0$.*

1. *For each $i = 1, 2, \ldots, m$ we have $|\langle a_i, x - x(0)\rangle| \leq \lambda_i \|a_i\|_2$;*

2. *$|x_i| > 1 - \delta$ for at least $n/10$ indices $i$.*

**Exercise 3.** *Let $A \in \{0,1\}^{m \times n}$ be a binary matrix with at most $d$ ones per column.*

   **a.** *Bound (in terms of $n$ and $d$) the numbers of rows of $A$ with at least $s$ ones, for any integer $1 \leq s \leq n$.*

   **b.** *Use Theorem 6 to show that for any $x(0) \in [-1,1]^n$ there exists an $x$ such that $\|A(x - x(0))\|_\infty = O(\sqrt{d})$ and for at least $n/10$ indices $i$ we have $|x_i| > 1 - \delta$.*

   **c.** *Show that $\mathrm{disc}(A) = O(\sqrt{d}\log n)$.*

**Exercise 4.** *Fix $k$ permutations $\pi_1, \ldots, \pi_k$ of $[n] = \{1, \ldots, n\}$. Let $\mathcal{S}$ be the set system consisting of the $kn$ sets $S_{ij} = \{\pi_i(1), \pi_i(2), \ldots, \pi_i(j)\}$, and let $A \in \{0,1\}^{kn \times n}$ be its incidence matrix.*

    **a.** *Use Theorem 6 to show that for any $x(0) \in [-1,1]^n$ there exists an $x$ such that $\|A(x - x(0))\|_\infty = O(k)$ and for at least $n/10$ indices $i$ we have $|x_i| > 1 - \delta$.*

        HINT*: Consider a different set system, derived by breaking each permutation into consecutive intervals of size $Ck$ for a large enough constant $C$*

    **b.** *Show that $\mathrm{disc}(\mathcal{S}) = \mathrm{disc}(A) = O(k \log n)$.*

    **c.** *(More challenging) Imrpove the bound from the first subproblem to $O(\sqrt{k})$ and the bound from the second subproblem to $O(\sqrt{k} \log n)$.*

# References

[1] Shachar Lovett and Raghu Meka. *Constructive discrepancy minimization by walking on the edges.* In the proceedings of FOCS 2012.

[2] Zbyněk Šidák. *Rectangular confidence regions for the means of multivariate normal distributions.* J. Amer. Statist. Assoc. 62: 626-633 (1967).

[3] Joel Spencer. *Six standard deviations suffice.* Transactions of the American Mathematical Society 289(2):679-706 (1985).