# Least singular value of random matrices

# Lewis Memorial Lecture / DIMACS minicourse

# March 18, 2008

Terence Tao (UCLA)

1

## Extreme singular values

Let $M = (a_{ij})_{1 \leq i \leq n; 1 \leq j \leq m}$ be a square or rectangular matrix with $1 \leq m \leq n$, thus $M : \mathbf{C}^m \to \mathbf{C}^n$. The *least singular value* $\sigma_m \geq 0$ can be defined as

$$\sigma_m := \inf_{v \in \mathbf{C}^m : \|v\| = 1} \|Mv\|.$$

The *largest singular value* $\sigma_1 \geq \sigma_m$ is similarly given by

$$\sigma_1 := \sup_{v \in \mathbf{C}^m : \|v\| = 1} \|Mv\|.$$

Of course, if the $a_{ij}$ are real, one can replace $\mathbf{C}$ by $\mathbf{R}$ throughout.

As in the previous lecture, we are interested in the case when the $a_{ij}$ are independent random variables. Important model examples include

- **Gaussian ensemble:** $a_{ij} \equiv N(0, 1)$ for all $i, j$.

- **Bernoulli ensemble:** $a_{ij} \in \{-1, +1\}$ for all $i, j$, with equal probability of each.

- **Resolvent models:** $n = m$, and $a_{ij} \equiv a - z\delta_{ij}$ for some fixed distribution $a$ and a deterministic $z \in \mathbf{C}$.

We shall mostly focus on the case of square matrices $n = m$, but also consider the rectangular case $m = (1 - \delta)n$ for fixed eccentricity $0 < \delta < 1$.

3

In this lecture, we shall focus on the problem of <span style="color:red">tail estimates</span> for the singular values. In particular, we are interested in <span style="color:red">upper tail estimates</span>

$$\mathbf{P}(\sigma_1 \geq K) \leq \ldots$$

for the largest singular value, and <span style="color:red">lower tail estimates</span>

$$\mathbf{P}(\sigma_m \leq \varepsilon) \leq \ldots$$

for the least singular value.

[Upper and lower tail estimates for other singular values are also of interest, but these two are particularly useful for upper-bounding the <span style="color:red">condition number</span> $\sigma_1/\sigma_m$.]

4

## The largest singular value

We begin with the theory of $\sigma_1$, which is much simpler and better understood. Furthermore, upper tail estimates on $\sigma_1$ will be needed to obtain lower tail estimates on $\sigma_m$.

The techniques here are quite general, but for sake of concreteness we focus on the Bernoulli case.

One can obtain remarkably good control on $\sigma_1$ by the moment method. For instance, from the second moment identity

$$nm = \text{tr}(M^*M) = \sum_{j=1}^{m} \sigma_j^2$$

we obtain the bounds

$$\sqrt{n} \leq \sigma_1 \leq \sqrt{nm}.$$

By computing higher moments, one can obtain much better bounds, for instance one can show

$$\mathbf{P}(\sigma_1 \leq C\sqrt{n}) \ll \exp(-cn)$$

for some absolute constants $C, c > 0$. Thus $\sigma_1 \sim \sqrt{n}$ with exponentially high probability. (This result holds more generally when the coefficients are uniformly subgaussian; weaker bounds are also available assuming uniformly bounded fourth moment (Latala, 2005).)

Similar higher moment computations also show that $\sigma_1$ is distributed according to the Tracy-Widom law. (This result holds more generally when the coefficients are uniformly subgaussian (Soshnikov, 2002; recently improved by Ruzmaikina, 2008).)

The moment method does not generalise well to lower tail bounds for the least singular value (because negative moments such as $\operatorname{tr}(M^*M)^{-1}$ are hard to compute). So let us give another proof of the upper tail estimate that does not rely as heavily on the moment method.

The starting point is the identity

$$\sigma_1 = \sup_{v \in S^{m-1}} \|Mx\|$$

where $S^{m-1} := \{v \in \mathbf{R}^m : \|x\| = 1\}$ is the unit sphere in $\mathbf{R}^m$.

The sphere $S^{m-1}$ is uncountable, and so we have an uncountably infinite supremum. But we can use the <span style="color:red">$\varepsilon$-net argument</span> to replace this infinite supremum with a finite one.

Let $\Omega \subset S^{m-1}$ be a 1/2-net of $S^{m-1}$, i.e. a maximal 1/2-separated subset of $S^{m-1}$. Then every element of $S^{m-1}$ lies within 1/2 of an element of $\Omega$. From the triangle inequality we conclude

$$\sigma_1 = \sup_{v \in S^{m-1}} \|Mv\| \le \sup_{x \in \Omega} \|Mv\| + \frac{1}{2}\sigma_1$$

and thus

$$\sigma_1 \le 2 \sup_{v \in \Omega} \|Mv\|.$$

On the other hand, the spherical caps centred at elements of $\Omega$ of radius 1/4 are all disjoint. Standard volume packing arguments thus yield the cardinality bound $|\Omega| \ll \exp(O(m))$. We can thus apply the crude <span style="color:red">union bound</span>:

$$\begin{aligned}
\mathbf{P}(\sigma_1 \geq K) &\leq \mathbf{P}(\sup_{v \in \Omega} \|Mv\| \geq K/2) \\
&\leq \sum_{v \in \Omega} \mathbf{P}(\|Mv\| \geq K/2) \\
&\ll \exp(O(m)) \sup_{v \in \Omega} \mathbf{P}(\|Mv\| \geq K/2).
\end{aligned}$$

We have thus reduced the upper tail estimate for $\sigma_1$ to upper tail estimates for $\|Mv\|$ for various unit vectors $v = (v_1, \ldots, v_m)$, paying an "entropy cost" $\exp(O(m))$.

If we let $X_1, \ldots, X_n \in \{-1, +1\}^m$ be the rows of $M$, we can write

$$\|Mv\| = \left(\sum_{i=1}^{n} |X_i \cdot v|^2\right)^{1/2}.$$

Direct computation shows that

$$\mathbf{E}\|Mv\|^2 = \sum_{i=1}^{n} \mathbf{E}|X_i \cdot v|^2 = \sum_{i=1}^{n} |v|^2 = n$$

so $\|Mv\|$ has an average value of $O(\sqrt{n})$.

Indeed, from the <span style="color:red">Chernoff inequality</span> one has an exponential tail estimate

$$\mathbf{P}(\|Mv\| \geq K\sqrt{n}) \ll \exp(-cKn)$$

for any $K \geq 2$ and some absolute constant $c > 0$. For $K$ large enough, the exponential tail of $\exp(-cKn)$ overwhelms the entropy penalty of $\exp(O(m))$, and we obtain

$$\mathbf{P}(\sigma_1 \geq K\sqrt{n}) \ll \exp(-cn)$$

as claimed.

Moral: The $\varepsilon$-net argument can control tail probabilities for singular values, as long as one has very good tail probabilities for $\|Mv\|$ and very good entropy bounds for the $\varepsilon$-net.

For the least singular value problem, this argument works well for vectors $v$ which are "compressible" or otherwise "structured", as such sets of vectors tend to have small entropy. But for generic vectors $v$, it may be that there is too much entropy for the $\varepsilon$-net argument to work.

## The least singular value

Now we consider lower tail probabilities for the least singular value $\sigma_1$. We continue to work with the Bernoulli case for concreteness.

We begin with

> **Rectangular case.** (Litvak-Pajor-Rudelson-Tomczak-Jaegermann, 2005) If $m = (1 - \delta)n$ for some fixed $0 < \delta < 1$ (independent of $n$), then
>
> $$\mathbf{P}(\sigma_m \leq \varepsilon\sqrt{n}) \ll \exp(-cn)$$
>
> for some $\varepsilon, c > 0$ depending on $\delta$ but not on $n$.

We rely again on the $\varepsilon$-net argument.

From the preceding discussion, we already know that $\sigma_1 = O(\sqrt{n})$ with exponentially high probability. Thus we can assume that $\sigma_1 \leq K\sqrt{n}$ for some $K = O(1)$.

We now let $\Omega$ be an $\varepsilon/K$-net of $S^{m-1}$. Then

$$\sigma_m = \inf_{v \in S^{m-1}} \|Mv\| \geq \inf_{v \in \Omega} \|Mv\| - \frac{\varepsilon}{K} \sigma_1$$

and so

$$\mathbf{P}(\sigma_m \leq \varepsilon\sqrt{n}) \leq \mathbf{P}(\inf_{v \in \Omega} \|Mv\| \leq 2\varepsilon\sqrt{n}).$$

Volume-packing arguments show that $|\Omega| \ll O(1/\varepsilon)^m$. Thus we have

$$\mathbf{P}(\sigma_m \leq \varepsilon\sqrt{n}) \ll O(1/\varepsilon)^m \sup_{v \in \Omega} \mathbf{P}(\|Mv\| \leq 2\varepsilon\sqrt{n}).$$

Once again, we have $\|Mv\| = (\sum_{i=1}^n |X_i \cdot v|^2)^{1/2}$. For most unit vectors $v$, $X_i \cdot v$ is distributed much like a Gaussian $N(0,1)$ (thanks to the central limit theorem), and Chernoff-type estimates then give

$$\mathbf{P}(\|Mv\| \leq 2\varepsilon\sqrt{n}) \ll O(\varepsilon)^n.$$

Since $m = (1 - \delta)n$ for some fixed $\delta > 0$, we thus obtain the claim by taking $\varepsilon$ small enough.

Unfortunately, there are a few unit vectors $v$ which are too <span style="color:red">sparse</span> or <span style="color:red">compressed</span> for the central limit theorem to apply. For instance, if $v = e_j$ is a basis vector, then $X_i \cdot v$ has the Bernoulli distribution, and $\mathbf{P}(\|Mv\| \leq 2\varepsilon\sqrt{n})$ is at least as large as $2^{-n}$.

Fortunately, the space of compressed vectors has much smaller entropy, and a variant of the above argument suffices to treat this case.

These results have been extended to the intermediate eccentricity case $\delta = o(1)$ (Litvak-Pajor-Rudelson-Tomczak-Jaegermann, 2005; Rudelson-Vershynin 2008). But for the remainder of this lecture we focus on the square case $\delta = 0$, so $m = n$. Here, we expect the least singular value $\sigma_n$ to be much smaller than $\sqrt{n}$.

Indeed, from the moment method it is known that the limiting distribution of the $n$ singular values $\sigma_1, \ldots, \sigma_n$ is a continuous distribution on $[0, 2\sqrt{n}]$. This suggests (but does not prove) that $\sigma_n$ should be of size $n^{-1/2}$. (Moment methods only give the weaker bound of $O_\varepsilon(n^\varepsilon)$ for any $\varepsilon > 0$.)

For the Gaussian ensemble, one can explicitly compute the distribution of $\sigma_n$, and show that

$$\mathbf{P}(\sigma_n \leq \varepsilon n^{-1/2}) \ll \varepsilon$$

for any $\varepsilon > 0$ (Edelman 1988). In particular, we have $\sigma_n \gg n^{-1/2}$ with high probability.

Analogous bounds are now known for discrete ensembles such as the Bernoulli ensemble (Rudelson 2006; T.-Vu 2007; Rudelson-Vershynin 2007). For instance, for Bernoulli (or more generally, subgaussian) ensembles we have

$$\mathbf{P}(\sigma_n \leq \varepsilon n^{-1/2}) \ll \varepsilon + c^n$$

for some $0 < c < 1$ (Rudelson-Vershynin 2007).

Here we shall present a slightly different estimate:

> **Theorem.** (T.-Vu 2007) Let $M$ be (say) the Bernoulli ensemble. Then for every $A > 0$ there exists $B > 0$ such that
>
> $$\mathbf{P}(\sigma_n \leq n^{-B}) \ll_A n^{-A}.$$

For the Bernoulli case, the argument allows one to take $B = A + 1/2 + o(1)$, thus almost recovering the optimal bound in the polynomial regime. On the other hand, the method here also extends to matrices with much weaker moment conditions (one only needs uniform second moments), because one does not need strong control on $\sigma_1$ (crude polynomial-type bounds $\sigma_1 \ll n^{O(1)}$ suffice).

## Overview of argument

Once again, our starting points are the identities

$$\mathbf{P}(\sigma_n \leq n^{-B}) = \mathbf{P}(\|Mv\| \leq n^{-B} \text{ for some } v \in S^{n-1})$$

and

$$\|Mv\| = (\sum_{i=1}^{n} |X_i \cdot v|^2)^{-1/2}.$$

A key role is played by the <span style="color:red">small ball probabilities</span>

$$p_\varepsilon(v) := \mathbf{P}(|X \cdot v| \leq \varepsilon)$$

for various $v \in S^{n-1}$ and $\varepsilon > 0$, where $X$ is uniformly distributed on $\{-1, +1\}^n$.

Following the general strategy of (Rudelson, 2006), one breaks into three cases:

- In the compressible case in which $v$ is mostly concentrated on a few coefficients, an $\varepsilon$-net argument works well (due to the low entropy of the space of compressible vectors), combined with Littlewood-Offord inequalities of Erdős type.

- In the poor case in which the small ball probabilities are very small (less than $n^{-A-O(1)}$), a conditioning argument (similar to those used in the previous lecture) suffices.

23

- The most difficult case is the red rich incompressible case in which $v$ is spread out among many coefficients, but has large small ball probabilities. Here an $\varepsilon$-net argument eventually works, but the entropy estimate required is non-trivial. (In the case of integer-valued matrices, one can avoid entropy calculations by using a discretisation argument instead (T.-Vu, 2007).)

Actually, one can unify the compressible and rich cases, though for conceptual reasons it seems better to keep these cases separate.

## The conditioning argument

We sketch here the conditioning argument (essentially dating back to (Komlós, 1967)) used to establish the bound

$$\mathbf{P}(\|Mv\| \leq n^{-B} \text{ for some poor } v \in S^{n-1}) \ll n^{-A}$$

where the vectors $v$ are poor in the sense that $p_{n^{-B+O(1)}}(v) \leq n^{-A-O(1)}$. To simplify the exposition we shall be careless with the $n^{O(1)}$ factors.

The key point is that $M$ has the same least singular value as its adjoint $M^*$. In particular, if $\|Mv\| \leq n^{-B}$ for some (poor) unit vector $v$, then we must also have $\|M^*w\| \leq n^{-B}$ for some unit vector $w$, thus we have a near-linear dependence

$$\|w_1 X_1 + \ldots + w_n X_n\| \leq n^{-B}$$

between the rows of $M$. By symmetry (and paying a factor of $n$ in the probability bounds) we may assume that $|w_n| \geq 1/\sqrt{n}$, and so we can express $X_n$ as a linear combination of $X_1, \ldots, X_{n-1}$ (with coefficients $O(n^{O(1)})$) plus an error of $O(n^{-B+O(1)})$.

On the other hand, since $\|Mv\| \leq n^{-B}$, we have $X_i \cdot v = O(n^{-B})$ for all $1 \leq i \leq n$. In particular, we can find a poor vector $u$ depending only on $X_1, \ldots, X_{n-1}$ such that $X_i \cdot u = O(n^{-B})$ for all $1 \leq i \leq n-1$. (One can view $u$ as a kind of "unit normal" to the space spanned by $X_1, \ldots, X_{n-1}$.)

We now condition $X_1, \ldots, X_{n-1}$ to fix $u$. In order for $\|Mv\| \leq n^{-B}$ to hold, $X_n$ needs to be approximately a combination of $X_1, \ldots, X_{n-1}$, which forces $X_n \cdot u = O(n^{-B+O(1)})$. But as $u$ is poor, this only occurs with probability $n^{-A-O(1)}$, and the claim follows.

27

## The rich case

Now let us consider the problem of bounding the contribution

$$\mathbf{P}(\|Mv\| \leq n^{-B} \text{ for some rich } v \in S^{n-1})$$

where the vectors are $p$-rich in the sense that

$$p_{n^{-B+O(1)}}(v) \sim p$$

for some $n^{-A-O(1)} \ll p \leq 1$. In this case we will in fact obtain an exponentially good bound $\exp(-cn)$.

There are some important technicalities involving the $n^{O(1)}$ factors, but we shall temporarily ignore them.

By paying a logarithmic factor, we can fix $p$ (up to constants). The hardest case is when $p \sim n^{-A-O(1)}$. (In the other extreme case $p \sim 1$, the Littlewood-Offord inequality of Erdős forces $v$ to be compressed, and the $\varepsilon$-net argument works easily in this case.)

Using a crude bound $\sigma_1 \ll n^{O(1)}$ and the $\varepsilon$-net argument, we can replace the sphere $S^{n-1}$ by an $n^{-B+O(1)}$-net $\Omega$.

For each rich $v$, we morally have

$$\mathbf{P}(\|Mv\| \leq n^{-B}) \ll p^n.$$

So we need an entropy bound on $\Omega$ that is much better than $p^{-n}$.

The full net $\Omega$ has cardinality about $(n^{B+O(1)})^n$ - far too large! But the set of <span style="color:red">$p$-rich</span> vectors in $\Omega$ is much smaller. Indeed, modulo some technicalities involving $n^{O(1)}$ factors, we have

> **Proposition** (T.-Vu, 2008) The number of $p$-rich vectors in $\Omega$ is at most $n^{-n/2+o(n)}p^{-n}$.

Except for these technicalities, this proposition settles the rich case.

A similar bound was implicitly obtained in (Rudelson-Vershynin, 2007) which works for more general $p$, but requires more care with the $n^{O(1)}$ factors (in particular, one needs $\sigma_1 = O(n^{1/2})$.

A variant of the Rudelson-Vershynin bound has also been obtained recently in (Götze-Tikhomorov, 2008), which can handle $n^{O(1)}$ factors (thus allowing for $\sigma_1$ to be large) but needs $p$ to be close to 1, thus giving a $o(1)$ tail probability bound for the least singular value in this case.

What are the $n^{O(1)}$ technicalities?

- A problem arises because the small ball probabilities $p_{n^{-B+O(1)}}(v)$ vary with the $O(1)$ parameter in the exponent.

- On the other hand, in the rich case we have $n^{-A-O(1)} \ll p_{n^{-B+O(1)}}(v) \le 1$ for all values of $O(1)$ in a certain range. Also, $p_\varepsilon(v)$ is clearly increasing in $\varepsilon$.

- If $B$ is sufficiently large depending on $A$, we may apply the pigeonhole principle to find $B' \leq B$ such that

$$p_{n^{-B'}}(v) \leq p_{n^{-B'+10}}(v) \leq n^{0.1} p_{n^{-B'}}(v)$$

  (say). Using this range of scales $[n^{-B'}, n^{-B'+10}]$ instead of $n^{-B+O(1)}$, we can fix the argument.

Optimising this argument carefully in the case $\sigma_1 = O(n^{1/2})$ gives the near-sharp dependence $B = A + 1/2 + o(1)$. More generally, when $\sigma_1 = O(n^{\gamma})$ for some $\gamma \geq 1/2$, we obtain $B = (2A+1)\gamma + o(1)$ (T.-Vu, 2008).

## The entropy bound

The heart of the matter is the entropy bound counting the number of $p$-rich vectors. This is essentially the

> **Inverse Littlewood-Offord problem** If $v = (v_1, \ldots, v_n)$ is a unit vector such that
>
> $$p_\varepsilon(v) = \mathbf{P}(|a_1 v_1 + \ldots + a_n v_n| \leq \varepsilon) \geq p,$$
>
> where $a_1, \ldots, a_n \in \{-1, +1\}$ are iid Bernoulli signs, what does this tell us about $v$? (In particular, what entropy does the set of such $v$ have?)

For instance, the classical Littlewood-Offord inequality of Erdős, when rephrased in this language, tells us that at most $O(1/p^2)$ coefficients of $v$ can exceed $\varepsilon$ in magnitude. Unfortunately this bound is only non-trivial for $p \geq 1/\sqrt{n}$. Later Littlewood-Offord results of Moser, Halász, Kleitman, Sarkőzy-Szemerédi, and others reveal more information for smaller $p$, but not enough to get satisfactory entropy bounds for the inverse problem.

We do not yet have an absolutely sharp answer to this problem (as compared to, say, the inverse theorems of Freiman type in additive combinatorics). However, the bounds we have are sufficient to establish various non-trivial bounds for the least singular value.

One case in which $p_\varepsilon(v)$ is expected to be large, is when the coefficients of $v$ lie close to those of a (symmetric) generalised arithmetic progression (GAP)

$$P := \{n_1 w_1 + \ldots + n_r w_r : n_i \in \{-N_i, \ldots, N_i\}\}$$

for some rank $r \geq 1$, some dimensions $N_1, \ldots, N_r$, and some steps $w_1, \ldots, w_r \in \mathbf{R}$. Indeed, from the theory of random walks, we know that if $v_1, \ldots, v_n \in P$, then the random walk $a_1 v_1 + \ldots + a_n v_n$ mostly lies in $\sqrt{n}P$, which has cardinality $O(n^{r/2}|P|)$, and so the small ball probability should be at least $n^{-r/2}/|P|$ in this case. Similarly if $v_1, \ldots, v_n$ "almost" lie in $P$.

It turns out that this implication is partially reversible: if the small ball probability is large, then almost all the coefficients of $v$ lie near an arithmetic progression. Here is a typical statement:

> **Theorem** (T.-Vu, 2007) If $p_\varepsilon(v) \geq n^{-A}$, then all but $n^{0.01}$ of the coefficients of $n$ coefficients of $v$ lie within $O(n^{O_A(1)}\varepsilon)$ of a GAP $P$ of rank $O_A(1)$ and volume $|P| \ll_A n^{O_A(1)}$.

Other results of this type are in (Rudelson-Vershynin, 2007) and (T.-Vu, 2008). The latter result is what is needed to prove the entropy bound mentioned earlier.

The proofs of these inverse Littlewood-Offord theorems proceeds by expressing the small ball probability $p_\varepsilon(v)$ in a Fourier-analytic form (using the Esséen concentration inequality), and then relying heavily on the additive combiantorics of GAPs (i.e. Minkowski's geometry of numbers).

## A discretisation trick

In the case of integer-valued matrices (e.g. Bernoulli ensemble) there is an additional discretisation trick which can allow one to avoid detailed entropy calculations.

The key observation is this: if one has some concentration

$$|a_1 v_1 + \ldots + a_n v_n| \leq \varepsilon$$

for some small $0 < \varepsilon < 1/4$ and some signs $a_1, \ldots, a_n \in \{-1, +1\}$, and the $v_i$ are all very close (within $1/4n$) to an integer, then we can round off to obtain an exact identity

$$a_1 [v_1] + \ldots + a_n [v_n] = 0$$

where $[v]$ is the nearest integer to $v$.

Similar phenomena hold if the coefficients of $v$ lie close to some "coarse" GAP, modulo "fine" errors which are small with respect to the spacing of that GAP. In this case we say that $v$ is <span style="color:red">discretisable</span>.

In principle, this trick allows us to estimate a tail probability $\mathbf{P}(\sigma_n < \varepsilon)$ by the <span style="color:red">singularity probability</span> $\mathbf{P}(\sigma_n = 0)$, which is much better understood.

It turns out that an elementary analysis of GAPs, combined with the inverse Littlewood-Offord theorem mentioned earlier, implies that all (polynomially) rich vectors are discretisable (T.-Vu 2007; see Rudelson-Vershynin 2007 for a related result). This gives an alternate way to treat the Bernoulli case.