



**How to prove tightness for the size of strange
random sets? Based on [GPS 2008].**

Gábor Pete

<http://www.math.toronto.edu/~gabor>

Recall the Fourier spectral sample

The space $L^2(\Omega, \mu)$, where $\Omega = \{\pm 1\}^V$, μ uniform probability measure, inner product $\mathbf{E}[fg]$, has a nice **orthonormal basis**:

For $S \subset V$, let $\chi_S(\omega) := \prod_{v \in S} \omega(v)$, the parity inside S .

Any function $f \in L^2(\Omega, \mu)$ decomposes in this basis (**Fourier-Walsh series**):

$$\hat{f}(S) := \mathbf{E}[f\chi_S]; \quad f(\omega) = \sum_{S \subset V} \hat{f}(S) \chi_S(\omega).$$

By Parseval, $\sum_S \hat{f}(S)^2 = \mathbf{E}[f^2]$. So can define probability measure $\mathbf{P}[\mathcal{S}_f = S] := \hat{f}(S)^2 / \mathbf{E}[f^2]$, the **spectral sample** $\mathcal{S}_f \subset V$.

$$\frac{\mathbf{E}[f(\omega^\epsilon)f(\omega)] - \mathbf{E}[f]^2}{\mathbf{E}[f^2]} = \sum_{S \neq \emptyset} \frac{\hat{f}(S)^2}{\mathbf{E}[f^2]} (1-\epsilon)^{|S|} = \mathbf{E}[(1-\epsilon)^{|\mathcal{S}_f|}; |\mathcal{S}_f| > 0],$$

hence small $\mathbf{P}[0 < |\mathcal{S}_f| < K/\epsilon]$ means small covariance after ϵ -noise.

Although goal is to understand size, **Gil Kalai** suggested trying to understand entire distribution. A strange random set of bits.

Effective sampling? If f is an effectively computable Boolean function, then there is an effective **quantum** algorithm for \mathcal{S}_f [**Bernstein-Vazirani 1993**].

For **critical planar percolation**, [**Smirnov '01**] + [**Tsirelson '04**] + [**Schramm-Smirnov**] implies that $\mathcal{S}_{Q,n}$ (left-right crossing in a conformal rectangle Q , mesh $1/n$) has a **conformally invariant scaling limit**.

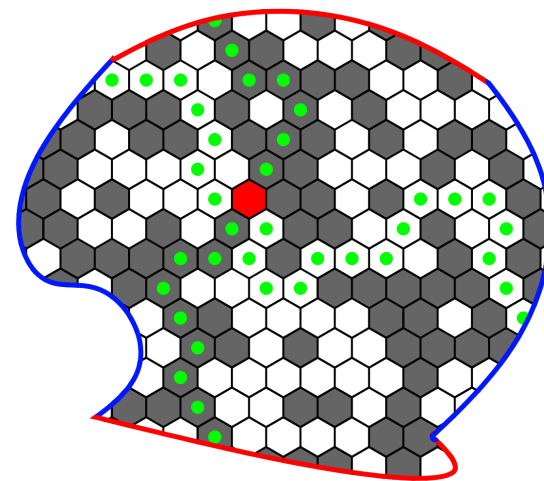
For ± 1 -valued f , can consider **pivotal bits**.

$\mathbf{P}[x, y \in \text{Piv}_f] = \mathbf{P}[x, y \in \mathcal{S}_f]$, but not for more points.

Both random subsets measure the “influence” or “relevance” of bits.

$\mathbf{P}[\mathcal{S}_{Q,n} \cap B \neq \emptyset] \asymp \mathbf{P}[B \text{ is pivotal for crossing } Q]$
 $= \alpha_4(B, Q)$, the 4-arm event.

$\mathbf{P}[\emptyset \neq \mathcal{S}_{Q,n} \subseteq B] \asymp \alpha_4(B, Q)^2$. But $\mathbf{P}[\emptyset \neq \text{Piv}_{Q,n} \subseteq B] \asymp \alpha_6(B, Q)$.



Three very simple examples

Dictator $_n(x_1, \dots, x_n) := x_1$.

Here $\text{Cov}[\text{Dic}_n(x), \text{Dic}_n(x^\epsilon)] = 1 - \epsilon$, so noise-stable.

And $\mathbf{P}[\mathcal{S}_n = \{x_1\}] = 1$.

Majority $_n(x_1, \dots, x_n) := \text{sgn}(x_1 + \dots + x_n) \approx \frac{1}{\sqrt{n}}(x_1 + \dots + x_n)$.

Here $\text{Cov}[\text{Maj}_n(x), \text{Maj}_n(x^\epsilon)] = 1 - O(\epsilon)$, so noise-stable.

And $\mathbf{P}[\mathcal{S}_n = \{x_i\}] \asymp 1/n$, most of the weight is on singletons.

On the other hand, $\mathbf{E}|\mathcal{S}_n| = \mathbf{E}|\text{Piv}_n| \asymp \frac{1}{\sqrt{n}} n \asymp \sqrt{n}$.

Parity $_n(x_1, \dots, x_n) := x_1 \cdots x_n$

Here $\text{Cov}[\text{Par}_n(x), \text{Par}_n(x^\epsilon)] = (1 - \epsilon)^n$, the most sensitive to noise.

And $\mathbf{P}[\mathcal{S}_n = \{x_1, \dots, x_n\}] = 1$.

Self-similarity for left-right crossing of $n \times n$ square

$$\mathbf{E}|\mathcal{S}_n| = \mathbf{E}|\text{Piv}_n| \asymp n^2 \alpha_4(1, n) \stackrel{\Delta}{\asymp} n^{3/4+o(1)},$$

$$\mathbf{E}|\mathcal{S}_n(r)| := \mathbf{E}\left[\#\{r\text{-boxes } \mathcal{S}_n \cap B_r \neq \emptyset\}\right] \asymp \frac{n^2}{r^2} \alpha_4(r, n) \asymp \mathbf{E}|\mathcal{S}_{n/r}|,$$

$$\mathbf{E}\left[|\mathcal{S}_n \cap B_r| \mid \mathcal{S}_n \cap B_r \neq \emptyset\right] \asymp r^2 \alpha_4(1, r) \asymp \mathbf{E}|\mathcal{S}_r|.$$

Of course, $r^2 \alpha_4(1, r) \cdot \frac{n^2}{r^2} \alpha_4(r, n) \asymp n^2 \alpha_4(1, n)$, by **quasi-multiplicativity**.

Self-similarity for left-right crossing of $n \times n$ square

$$\mathbf{E}|\mathcal{S}_n| = \mathbf{E}|\text{Piv}_n| \asymp n^2 \alpha_4(1, n) \stackrel{\Delta}{\asymp} n^{3/4+o(1)},$$

$$\mathbf{E}|\mathcal{S}_n(r)| := \mathbf{E}\left[\#\{r\text{-boxes } \mathcal{S}_n \cap B_r \neq \emptyset\}\right] \asymp \frac{n^2}{r^2} \alpha_4(r, n) \asymp \mathbf{E}|\mathcal{S}_{n/r}|,$$

$$\mathbf{E}\left[|\mathcal{S}_n \cap B_r| \mid \mathcal{S}_n \cap B_r \neq \emptyset\right] \asymp r^2 \alpha_4(1, r) \asymp \mathbf{E}|\mathcal{S}_r|.$$

Of course, $r^2 \alpha_4(1, r) \cdot \frac{n^2}{r^2} \alpha_4(r, n) \asymp n^2 \alpha_4(1, n)$, by **quasi-multiplicativity**.

Similar to the **zero-set of simple random walk**: $\mathbf{E}|\mathcal{Z}_n| \asymp n n^{-1/2} = n^{1/2}$,

$$\mathbf{E}|\mathcal{Z}_n(r)| := \mathbf{E}\left[\#\{r\text{-intervals } \mathcal{Z}_n \cap I_r \neq \emptyset\}\right] \asymp \frac{n}{r} (n/r)^{-1/2} \asymp \mathbf{E}|\mathcal{Z}_{n/r}|,$$

$$\mathbf{E}\left[|\mathcal{Z}_n \cap I_r| \mid \mathcal{Z}_n \cap I_r \neq \emptyset\right] \asymp r r^{-1/2} \asymp \mathbf{E}|\mathcal{Z}_r|.$$

The $\mathcal{S}_n(r)$ and $\mathcal{Z}_n(r)$ results are related to the existence of scaling limits.

What concentration can we expect?

\mathcal{S}_n is very different from **uniform set** of similar density:
i.i.d. $\mathbf{P}[x \in \mathcal{U}_n] = n^{-5/4}$. Hence $\mathbf{E}|\mathcal{U}_n| = n^{3/4}$.

For large r ($\gg n^{5/8}$), this \mathcal{U}_n intersects every r -box;
for small r , if it intersects one, there is just one point there.

Concentration of size: roughly within $\sqrt{\mathbf{E}|\mathcal{U}_n|} = n^{3/8}$.

A bit more similar: for $i = 1, \dots, (n/r)^2$, i.i.d. $\mathbf{P}[X_i = r^{3/4}] = (n/r)^{-5/4}$,
 $X_i = 0$ otherwise. Then $S_{n,r} := \sum_i X_i$. Hence $\mathbf{E}|S_{n,r}| = n^{3/4}$.

For $r = n^\gamma$, size $|S_{n,r}|$ is concentrated within $n^{3/8(1+\gamma)}$, still $o(\mathbf{E}|S_{n,r}|)$.

For self-similar sets, we expect only **tightness around the mean**:
 $\mathbf{P}[0 < |\mathcal{S}_n| < \lambda \mathbf{E}|\mathcal{S}_n|] \rightarrow 0$ as $\lambda \rightarrow 0$, uniformly in n .

Proving tightness with a lot of independence

Assume we have the following ingredients, true for the zeroes:

$$(1) \quad \mathbf{P} \left[|\mathcal{Z}_n \cap I_r| > c \mathbf{E}|\mathcal{Z}_r| \mid \mathcal{Z}_n \cap I_r \neq \emptyset, \mathcal{F}_{[n] \setminus I_r} \right] \geq c > 0.$$

$$(2) \quad \mathbf{P} \left[|\mathcal{Z}_n(r)| = k \right] \leq g(k) \mathbf{P} \left[|\mathcal{Z}_n(r)| = 1 \right], \text{ with sub-exponential } g(k):$$

when the r -intervals intersected are scattered, have to pay k times to get to and leave them, and this cost is not balanced by combinatorial entropy.

Proving tightness with a lot of independence

Assume we have the following ingredients, true for the zeroes:

$$(1) \quad \mathbf{P} \left[|\mathcal{Z}_n \cap I_r| > c \mathbf{E}|\mathcal{Z}_r| \mid \mathcal{Z}_n \cap I_r \neq \emptyset, \mathcal{F}_{[n] \setminus I_r} \right] \geq c > 0.$$

$$(2) \quad \mathbf{P} \left[|\mathcal{Z}_n(r)| = k \right] \leq g(k) \mathbf{P} \left[|\mathcal{Z}_n(r)| = 1 \right], \text{ with sub-exponential } g(k):$$

when the r -intervals intersected are scattered, have to pay k times to get to and leave them, and this cost is not balanced by combinatorial entropy.

$$\mathbf{P} \left[0 < |\mathcal{Z}_n| < c \mathbf{E}|\mathcal{Z}_r| \right] = \sum_{k \geq 1} \mathbf{P} \left[0 < |\mathcal{Z}_n| < c \mathbf{E}|\mathcal{Z}_r|, |\mathcal{Z}_n(r)| = k \right]$$

$$\text{by (1):} \quad \leq \sum_{k \geq 1} (1 - c)^k \mathbf{P} \left[|\mathcal{Z}_n(r)| = k \right]$$

$$\text{by (2):} \quad \leq O(1) \mathbf{P} \left[|\mathcal{Z}_n(r)| = 1 \right] \asymp (n/r)^{1-3/2},$$

which, using $\lambda = \frac{c \mathbf{E}|\mathcal{Z}_r|}{\mathbf{E}|\mathcal{Z}_n|} \asymp (r/n)^{1/2}$, reads as $\mathbf{P} \left[0 < |\mathcal{Z}_n| < \lambda \mathbf{E}|\mathcal{Z}_n| \right] \asymp \lambda$.

But we know much less independence for \mathcal{S}_n

$$(1') \quad \mathbf{P} \left[|\mathcal{S}_n \cap B_r/3| > c \mathbf{E}|\mathcal{S}_r| \mid \mathcal{S}_n \cap B_r \neq \emptyset = \mathcal{S}_n \cap W \right] \geq c > 0,$$

for any W that is not too close to B_r .

Why only this negative conditioning? **Inclusion formula:**

$$\mathbf{P}[\mathcal{S}_f \subset U] = \sum_{S \subset U} \hat{f}(S)^2 = \mathbf{E} \left[\left(\sum_{S \subset U} \hat{f}(S) \chi_S \right)^2 \right] = \mathbf{E} \left[\mathbf{E}[f \mid \mathcal{F}_U]^2 \right].$$

From this, for disjoint subsets A and B ,

$$\begin{aligned} \mathbf{P}[\mathcal{S}_f \cap B \neq \emptyset = \mathcal{S}_f \cap A] &= \mathbf{P}[\mathcal{S}_f \subseteq A^c] - \mathbf{P}[\mathcal{S}_f \subseteq (A \cup B)^c] \\ &= \mathbf{E} \left[\mathbf{E}[f \mid \mathcal{F}_{A^c}]^2 - \mathbf{E}[f \mid \mathcal{F}_{(A \cup B)^c}]^2 \right] \\ &= \mathbf{E} \left[\left(\mathbf{E}[f \mid \mathcal{F}_{A^c}] - \mathbf{E}[f \mid \mathcal{F}_{(A \cup B)^c}] \right)^2 \right]. \end{aligned}$$

So, what are we going to do?

With quite a lot of work for both items,

$$(1') \quad \mathbf{P} \left[|\mathcal{S}_n \cap B_r/3| > c \mathbf{E}|\mathcal{S}_r| \mid \mathcal{S}_n \cap B_r \neq \emptyset = \mathcal{S}_n \cap W \right] \geq c > 0.$$

$$(2) \quad \mathbf{P} \left[|\mathcal{S}_n(r)| = k \right] \leq g(k) \mathbf{P} \left[|\mathcal{S}_n(r)| = 1 \right], \text{ with sub-exponential } g(k).$$

We could repeat (1') for many r -boxes only if “not enough points in one box” meant “we found nothing in that box”.

So, take an **independent random dilute sample**: $\mathbf{P}[x \in \mathcal{R}] = 1/\mathbf{E}|\mathcal{S}_r|$ i.i.d.

Then, $|\mathcal{S}_n \cap B_r/3|$ is small $\implies \mathcal{R} \cap \mathcal{S}_n \cap B_r/3 = \emptyset$ is likely,

and $|\mathcal{S}_n \cap B_r/3|$ is large $\implies \mathcal{R} \cap \mathcal{S}_n \cap B_r/3 \neq \emptyset$ is likely.

So, what are we going to do?

With quite a lot of work for both items,

$$(1') \quad \mathbf{P} \left[|\mathcal{S}_n \cap B_{r/3}| > c \mathbf{E}|\mathcal{S}_r| \mid \mathcal{S}_n \cap B_r \neq \emptyset = \mathcal{S}_n \cap W \right] \geq c > 0.$$

$$(2) \quad \mathbf{P} \left[|\mathcal{S}_n(r)| = k \right] \leq g(k) \mathbf{P} \left[|\mathcal{S}_n(r)| = 1 \right], \text{ with sub-exponential } g(k).$$

We could repeat (1') for many r -boxes only if “not enough points in one box” meant “we found nothing in that box”.

So, take an **independent random dilute sample**: $\mathbf{P}[x \in \mathcal{R}] = 1/\mathbf{E}|\mathcal{S}_r|$ i.i.d.
Then, $|\mathcal{S}_n \cap B_{r/3}|$ is small $\implies \mathcal{R} \cap \mathcal{S}_n \cap B_{r/3} = \emptyset$ is likely,
and $|\mathcal{S}_n \cap B_{r/3}|$ is large $\implies \mathcal{R} \cap \mathcal{S}_n \cap B_{r/3} \neq \emptyset$ is likely.

But $\mathbf{P} \left[\mathcal{S}_n \neq \emptyset = \mathcal{R} \cap \mathcal{S}_n \mid |\mathcal{S}_n(r)| = k \right]$ is still problematic conditioning.

Or, if we **scan sequentially** the r -boxes until $\mathcal{R} \cap \mathcal{S}_n \cap B_{r/3} \neq \emptyset$, how would (2) imply that we had a good chance of success several times? We don't know how $\mathbf{P} \left[\mathcal{S}_n \cap B_r(t) \neq \emptyset \mid \mathcal{S}_n \cap W(t) = \emptyset \right]$ changes with the steps t .

Oded's first solution: a filtered Markov inequality

If \mathcal{F}_k is a monotone increasing filtration, X_k are non-negative variables, and $Y_k := \mathbf{E}[X_k \mid \mathcal{F}_k]$, then, for any $s, t \geq 0$,

$$\mathbf{P}\left[\sum_k Y_k \leq s, \sum_k X_k \geq t\right] \leq s/t.$$

In the application, \mathcal{F}_k is the σ -algebra generated by the random sets $\{\mathcal{R} \cap \mathcal{S}_n \cap B_j : j \leq k-1\}$, and $X_k = 1_{\{\mathcal{S}_n \cap B_k \neq \emptyset\}}$.

Since (2) says that $\sum_k X_k$ is probably large, we get the same for $\sum_k Y_k$. Hence, with large probability, there are several boxes where the scanning has a positive chance to succeed, so it is unlikely that it remains unsuccessful.

However, the Markov-type bound is too weak, we don't get the sharp result.

Oded's second solution: a large deviation lemma

Suppose $X_i, Y_i \in \{0, 1\}$, $i = 1, \dots, n$, and that $\forall J \subset [n]$ and $\forall i \in [n] \setminus J$

$$\mathbf{P}[Y_i = 1 \mid \forall_{j \in J} Y_j = 0] \geq c \mathbf{P}[X_i = 1 \mid \forall_{j \in J} Y_j = 0].$$

Then

$$\mathbf{P}[\forall_i Y_i = 0] \leq c^{-1} \mathbf{E}\left[\exp\left(-\frac{c}{e} \sum_i X_i\right)\right].$$

We use this with $X_j := 1_{\{\mathcal{S} \cap B_j \neq \emptyset\}}$ and $Y_j := 1_{\{\mathcal{S} \cap B_j \cap \mathcal{R} \neq \emptyset\}}$.

Proof: Instead of sequential scan, average everything together.

Choose $J \subset [n]$ randomly, Bernoulli($1-p$). Get $\mathbf{E}[Y p^Y] \geq c \mathbf{E}[X p^{Y+1}]$.

So, $\mathbf{E}[Z] \geq 0$, where $Z := (Y - cpX) p^Y$. Choose $p := e^{-1}$. Maximize Z over Y , and get the bound $Z \leq \exp(-1 - cX/e)$. Altogether, $ce^{-1} \mathbf{P}[Y = 0 < X] \leq \mathbf{E}[1_{X>0} \exp(-1 - cX/e)]$, and done.

Final result for the spectral sample

If $r \in [1, n]$, then $\{|\mathcal{S}_n| < \mathbf{E}|\mathcal{S}_r|\}$ is basically equivalent to being contained inside some $r \times r$ sub-square:

$$\mathbf{P}[0 < |\mathcal{S}_n| < \mathbf{E}|\mathcal{S}_r|] \asymp \alpha_4(r, n)^2 \left(\frac{n}{r}\right)^2.$$

In particular, on the triangular lattice Δ ,

$$\mathbf{P}[0 < |\mathcal{S}_n| < \lambda \mathbf{E}|\mathcal{S}_n|] \asymp \lambda^{2/3}.$$

The *scaling limit* of \mathcal{S}_n is a conformally invariant Cantor-set with Hausdorff-dimension $3/4$.

Remark. The same strategy gives $\mathbf{P}[0 < |\text{Piv}_n| < \lambda \mathbf{E}|\text{Piv}_n|] \asymp \lambda^{11/9}$, but it's an overkill, given all the independence in Piv_n .

Some related questions

Question 1: Can one build similar proofs for other Boolean functions?

Question 2: Self-similarity of Piv_n and \mathcal{S}_n is a lot of restriction.

Conjecture [Gil Kalai]: The **entropy** of such random sets X_n is at most $\mathbf{E}|X_n|$, i.e., there is no log factor as in uniform.

In particular, **Influence-Entropy conjecture [Friedgut-Kalai 1996]:** For some universal constant C , for any Boolean function f ,

$$\begin{aligned} \text{SpecEnt}(f) &:= \sum_{S \subset [n]} \hat{f}(S)^2 \log \frac{1}{\hat{f}(S)^2} \leq C \times \\ &\times \text{Influence}(f) := \mathbf{E}|\mathcal{S}_f| = \mathbf{E}|\text{Piv}_f| = \sum_{S \subset [n]} \hat{f}(S)^2 |S|. \end{aligned}$$

I think I can do it for Piv_n , but have no idea about \mathcal{S}_n .