

## LECTURE 8: FREE ENTROPY AND LARGE DEVIATIONS

An important concept in classical probability theory is Shannon's notion of entropy. Having developed the analogy between free and classical probability theory, one hopes to find that a notion of "free entropy" exists in counterpart to the Shannon entropy. In fact there is a useful notion of free entropy. However, the development of this new concept is at present far from complete. The current state of affairs is that there are two distinct approaches to free entropy. These should give isomorphic theories, but at present we only know that they coincide in a limited number of situations.

The first approach to a theory of free entropy is via "microstates." This is rooted in the concept of large deviations. The second approach is "microstates free." This draws its inspiration from the statistical approach to classical entropy via the notion of Fisher information. The unification problem in free probability theory is to prove that these two theories of free entropy are consistent.

Let us return to the connection between random matrix theory and free probability theory which we have been developing. We know that a  $p$ -tuple

$$(A_N^{(1)}, \dots, A_N^{(p)})$$

of  $N \times N$  matrices chosen independently at random with respect to the density

$$(1) \quad P_N(A) = \text{const} \cdot e^{-\frac{N}{2} \text{Tr}(A^2)}$$

on the space of  $N \times N$  Hermitian matrices converges almost surely (in moments with respect to the normalized trace) to a freely independent family

$$(s_1, \dots, s_p)$$

of semi-circular elements lying in a non-commutative probability space. The von Neumann algebra generated by  $p$  freely independent semi-circulars is the von Neumann algebra  $L(\mathbb{F}_p)$  of the free group on  $p$  generators.

How likely is it to observe other distributions/operators for large  $N$ ?

---

*Date:* Lecture given on Nov. 22, 2007.

Let us consider the case  $p = 1$  more closely. For a random Hermitian matrix  $A = A^*$  (distribution as above) with real random eigenvalues

$$(2) \quad \lambda_1 \leq \cdots \leq \lambda_N,$$

denote by

$$(3) \quad \mu_A = \frac{1}{N}(\delta_{\lambda_1} + \cdots + \delta_{\lambda_N})$$

the eigenvalue distribution of  $A$  (also known as the “empirical eigenvalue distribution”), which is a random measure on  $\mathbb{R}$ . Wigner’s semicircle law states that as  $N \rightarrow \infty$

$$(4) \quad P_N(\mu_A \approx \mu_W) \rightarrow 1,$$

where  $\mu_W$  is the (non-random) semicircular distribution (see Lecture 1). What is the rate of decay of the probability

$$(5) \quad P_N(\mu_A \approx \nu),$$

where  $\nu$  is some measure (not necessarily the semicircle)? We expect that

$$(6) \quad P_N(\mu_A \approx \nu) \sim e^{-N^2 I(\nu)}$$

for some “rate function”  $I$  vanishing at  $\mu_W$ . By analogy with the classical theory of large deviations,  $I$  should correspond to a suitable notion of free entropy.

## 1. LARGE DEVIATION THEORY

Consider a sequence  $X_1, X_2, \dots$  of independent identically distributed random variables with distribution  $\mu$ , and put

$$(7) \quad S_n = \frac{X_1 + \cdots + X_n}{n}.$$

Let  $m = \mathbb{E}[X_i]$ . Then the law of large numbers asserts that  $S_n \rightarrow m$ , while the central limit theorem tells us that for large  $n$

$$(8) \quad S_n \approx m + \frac{\sigma^2}{\sqrt{n}} N(0, 1).$$

For example if  $\mu = N(0, 1)$  then  $S_n$  has distribution  $N(0, \frac{1}{n})$  and hence

$$(9) \quad P(S_n \approx x) = e^{-n \frac{x^2}{2}} \frac{\sqrt{n}}{\sqrt{2\pi}} \sim e^{-nI(x)}.$$

Thus the probability that  $S_n$  is near the value  $x$  decays exponentially in  $n$  at a rate determined by  $x$ , namely the “rate function”  $I(x) = \frac{x^2}{2}$ . Note that the convex function  $I(x)$  has a global minimum at  $x = 0$ ,

the minimum value there being 0, which corresponds to the fact that  $S_n$  approaches the mean 0 in probability.

This behavior is described in general by the following theorem of Cramér: Let  $X_1, X_2, X_3, \dots$  be a sequence of i.i.d random variables, with mean  $m$ , and set

$$(10) \quad S_n := \frac{1}{n}(X_1 + \dots + X_n),$$

the empirical mean. There exists a function  $I(x)$ , the rate function, such that

$$(11) \quad P(S_n > x) \approx e^{-nI(x)}, x > m$$

$$(12) \quad P(S_n < x) \approx e^{-nI(x)}, x < m.$$

How does one calculate the rate function  $I(x)$  for a given distribution  $\mu$ ? Suppose  $\mu$  has mean 0. For arbitrary  $x > 0$ , one has for all  $\lambda > 0$

$$\begin{aligned} P(S_n > x) &= P(nS_n > nx) \\ &\leq \mathbb{E}[e^{\lambda(nS_n - nx)}] \\ &= e^{-\lambda nx} \mathbb{E}[e^{\lambda(X_1 + \dots + X_n)}] \\ &= e^{-\lambda nx} \mathbb{E}[e^{\lambda X_i}]^n. \end{aligned}$$

Now put

$$(13) \quad \Lambda(\lambda) := \log \mathbb{E}[e^{\lambda X_i}],$$

the cumulant generating series of  $\mu$ . Then the above reads

$$(14) \quad P(S_n > x) \leq e^{-\lambda nx + n\Lambda(\lambda)} = e^{-n(\lambda x - \Lambda(\lambda))},$$

valid for all  $\lambda > 0$ . As  $\Lambda$  achieves its minimum at  $\lambda = 0$ , the above estimate is also valid for negative  $\lambda$ . Thus

$$(15) \quad P(S_n > x) \leq \inf_{\lambda} e^{-n(\lambda x - \Lambda(\lambda))} = e^{-n \sup_{\lambda} \{\lambda x - \Lambda(\lambda)\}}.$$

The function  $\lambda \mapsto \Lambda(\lambda)$  is convex, and the ‘‘Legendre transform’’ of  $\Lambda$  defined by

$$(16) \quad \Lambda^*(x) := \sup_{\lambda} \{\lambda x - \Lambda(\lambda)\}$$

is also convex. Thus we have proved that

$$(17) \quad P(S_n > x) \leq e^{-n\Lambda^*(x)},$$

where  $\Lambda^*$  is the Legendre transform of the cumulant generating function  $\Lambda$ . This gives  $\Lambda^*$  as a candidate for the rate function; however we also

have to check that  $e^{-n\Lambda^*(x)}$  is, at least asymptotically, a lower bound; more precisely, we need to verify that

$$(18) \quad \liminf \frac{1}{n} \log P(x - \delta < S_n < x + \delta) \geq -\Lambda^*(x)$$

for all  $x$  and all  $\delta > 0$ . By making an appropriate shift we can reduce this to the case  $x = 0$ . Then  $-\Lambda^*(0) = \inf_{\lambda} \Lambda(\lambda)$ . The idea of the proof is then to perturb the distribution  $\mu$  to  $\tilde{\mu}$  such that  $x = 0$  is the mean of  $\tilde{\mu}$ . Consider the case where  $\Lambda(\lambda)$  has a global minimum at  $\eta$ , and put

$$(19) \quad \frac{d\tilde{\mu}}{d\mu}(x) = e^{\eta x - \Lambda(\eta)}.$$

Note that

$$(20) \quad \int_{\mathbb{R}} d\tilde{\mu} = e^{-\Lambda(\eta)} \int e^{\eta x} dx$$

$$(21) \quad = e^{-\Lambda(\eta)} \mathbb{E}[e^{\eta X_i}]$$

$$(22) \quad = e^{-\Lambda(\eta)} e^{\Lambda(\eta)}$$

$$(23) \quad = 1,$$

which verifies that  $\tilde{\mu}$  is a probability measure. Consider now  $\tilde{X}_i$  i.i.d with distribution  $\tilde{\mu}$ , and put

$$(24) \quad \tilde{S}_n = \frac{\tilde{X}_1 + \cdots + \tilde{X}_n}{n}.$$

We have

$$(25) \quad \mathbb{E}[\tilde{X}_i] = \int x d\tilde{\mu}(x)$$

$$(26) \quad = e^{-\Lambda(\eta)} \int x e^{\eta x} dx$$

$$(27) \quad = e^{-\Lambda(\eta)} \frac{d}{d\lambda} e^{\Lambda(\lambda)} \Big|_{\lambda=\eta}$$

$$(28) \quad = e^{-\Lambda(\eta)} \Lambda'(\eta) e^{\Lambda(\eta)}$$

$$(29) \quad = \Lambda'(\eta)$$

$$(30) \quad = 0.$$

Now, for all  $\epsilon > 0$ , we have

$$(31) \quad P(-\epsilon < S_n < \epsilon) = \int_{|\sum_{i=1}^n X_i| < n\epsilon} \mu(dx_1) \dots \mu(dx_n)$$

$$(32) \quad \geq e^{-n\epsilon|\eta|} \int_{|\sum_{i=1}^n X_i| < n\epsilon} e^{n\sum X_i} \mu(dx_1) \dots \mu(dx_n)$$

$$(33) \quad = e^{-n\epsilon|\eta|} e^{n\Lambda(\eta)} \int_{|\sum_{i=1}^n X_i| < n\epsilon} \tilde{\mu}(dx_1) \dots \tilde{\mu}(dx_n)$$

$$(34) \quad = e^{-n\epsilon|\eta|} e^{n\Lambda(\eta)} P(-\epsilon < \tilde{S}_n < \epsilon).$$

By the law of large numbers,  $\tilde{S}_n \rightarrow \mathbb{E}_{\tilde{\mu}}[\tilde{X}_i] = 0$ , i.e.

$$(35) \quad \lim_{n \rightarrow \infty} P(-\epsilon < \tilde{S}_n < \epsilon) = 1$$

for all  $\epsilon > 0$ . Thus for all  $0 < \epsilon < \delta$

$$(36) \quad \liminf \frac{1}{n} \log P(-\delta < S_n < \delta) \geq \liminf \frac{1}{n} \log P(-\epsilon < S_n < \epsilon)$$

$$(37) \quad \geq \Lambda(\eta) - \epsilon|\eta|, \text{ for all } \epsilon$$

$$(38) \quad \geq \Lambda(\eta)$$

$$(39) \quad = \inf \Lambda(\lambda)$$

$$(40) \quad = -\Lambda^*(0).$$

This sketches the proof of Cramer's theorem for  $\mathbb{R}$ . A higher-dimensional form of Cramer's theorem is given in the next section.

## 2. CRAMER'S THEOREM FOR $\mathbb{R}^d$

Let  $X_1, X_2, \dots$  be a sequence of i.i.d random vectors, i.e. independent  $\mathbb{R}^d$ -valued random variables with common distribution  $\mu$  (a probability measure on  $\mathbb{R}^d$ ). Put

$$(41) \quad \Lambda(\lambda) := \mathbb{E}[e^{\langle \lambda, X_i \rangle}], \quad \lambda \in \mathbb{R}^d,$$

and

$$(42) \quad \Lambda^*(x) := \sup_{\lambda \in \mathbb{R}^d} \{\langle \lambda, x \rangle - \Lambda(\lambda)\}.$$

Assume that  $\Lambda(\lambda) < \infty$  for all  $\lambda \in \mathbb{R}^d$ , and put

$$(43) \quad S_n := \frac{1}{n}(X_1 + \dots + X_n).$$

Then the distribution  $\mu_{S_n}$  of the random variable  $S_n$  satisfies a large deviation principle with rate function  $\Lambda^*$ , i.e.

- $x \mapsto \Lambda^*(x)$  is lower semicontinuous (actually convex)

- $\Lambda^*$  is “good,” i.e.  $\{x \in \mathbb{R}^d : \Lambda^*(x) \leq \alpha\}$  is compact for all  $\alpha \in \mathbb{R}$
- For any closed set  $F \subset \mathbb{R}^d$ ,  $\limsup_n \frac{1}{n} \log P(S_n \in F) \leq -\inf_{x \in F} \Lambda^*(x)$
- For any open set  $G \subset \mathbb{R}^d$ ,  $\liminf_n \log P(S_n \in G) \geq -\inf_{x \in G} \Lambda^*(x)$ .

This is Cramer’s Theorem for  $\mathbb{R}^d$ , and in an informal way it says

$$(44) \quad P(S_n \approx x) \sim e^{-n\Lambda^*(x)}.$$

Actually, we are interested not in  $S_n$ , but in the empirical distribution

$$(45) \quad \frac{1}{n}(\delta_{X_1} + \cdots + \delta_{X_n}).$$

Let us consider this in the special case of random variables taking values in a finite alphabet  $A = \{a_1, \dots, a_d\}$ :

$$(46) \quad X_i : \Omega \rightarrow A,$$

with  $p_k := P(X_i = a_k)$ . As  $n \rightarrow \infty$ , the empirical distribution of the  $X_i$ ’s should converge to the “most likely” probability measure  $(p_1, \dots, p_d)$  on  $A$ .

Now define the vector of indicator functions  $Y_i : \Omega \rightarrow \mathbb{R}^d$  by

$$(47) \quad Y_i := (1_{a_1}(X_i), \dots, 1_{a_d}(X_i)),$$

so that in particular  $p_k$  is equal to the probability that  $Y_i$  will have a 1 in the  $k$ -th spot and 0’s elsewhere. Then

$$(48) \quad \frac{1}{n}(Y_1 + \cdots + Y_n)$$

gives the relative frequency of  $a_1, \dots, a_d$ :

$$(49) \quad \frac{1}{n}(\delta_{X_1} + \cdots + \delta_{X_n}),$$

i.e. the empirical distribution of  $(X_1, \dots, X_n)$ .

A probability measure on  $A$  is given by a  $d$ -tuple  $(q_1, \dots, q_d)$  of positive real numbers satisfying  $q_1 + \cdots + q_d = 1$ . By Cramer’s theorem,

$$(50) \quad P\left\{\frac{1}{n}(\delta_{X_1} + \cdots + \delta_{X_n}) \approx (q_1, \dots, q_d)\right\} = P\left(\frac{Y_1 + \cdots + Y_n}{n} \approx (q_1, \dots, q_d)\right) \sim e^{-n\Lambda^*(q_1, \dots, q_d)}.$$

Here

$$(51) \quad \Lambda(\lambda_1, \dots, \lambda_d) = \log \mathbb{E}[e^{\langle \lambda, Y_i \rangle}]$$

$$(52) \quad = \log(p_1 e^{\lambda_1} + \cdots + p_d e^{\lambda_d}).$$

Thus the Legendre transform is given by

$$(53) \quad \Lambda^*(q_1, \dots, q_d) = \sup_{(\lambda_1, \dots, \lambda_d)} \{\lambda_1 q_1 + \cdots + \lambda_d q_d - \Lambda(\lambda_1, \dots, \lambda_d)\}.$$

We compute the supremum over all tuples  $(\lambda_1, \dots, \lambda_d)$  by finding the partial  $\partial/\partial\lambda_i$  of

$$(54) \quad \lambda_1 q_1 + \dots + \lambda_d q_d - \Lambda(\lambda_1, \dots, \lambda_d)$$

to be

$$(55) \quad q_i - \frac{1}{p_1 e^{\lambda_1} + \dots + p_d e^{\lambda_d}} p_i e^{\lambda_i}.$$

Thus the max occurs when

$$(56) \quad \lambda_i = \log \frac{q_i}{p_i} + \log(p_1 e^{\lambda_1} + \dots + p_d e^{\lambda_d}),$$

and we compute

$$(57) \quad \Lambda^*(q_1, \dots, q_d) = q_1 \log \frac{q_1}{p_1} + \dots + q_d \log \frac{q_d}{p_d} + (q_1 + \dots + q_d)\Lambda - \Lambda$$

$$(58) \quad = q_1 \log \frac{q_1}{p_1} + \dots + q_d \log \frac{q_d}{p_d} + \Lambda - \Lambda$$

$$(59) \quad = q_1 \log \frac{q_1}{p_1} + \dots + q_d \log \frac{q_d}{p_d}$$

$$(60) \quad = H((q_1, \dots, q_d)|(p_1, \dots, p_d)),$$

the relative entropy of  $(q_1, \dots, q_d)$  with respect to  $(p_1, \dots, p_d)$ . Note that  $H((q_1, \dots, q_d)|(p_1, \dots, p_d)) \geq 0$ , with equality holding if and only if  $q_1 = p_1, \dots, q_d = p_d$ .

Thus  $(p_1, \dots, p_d)$  is the most likely realization, with other realizations exponentially unlikely; their unlikelihood is measured by the rate function  $\Lambda^*$ . And this rate function is indeed Shannon's relative entropy. This statement is *Sanov's theorem* for a finite alphabet; it also holds true for continuous distributions.

**Sanov's Theorem:** Let  $X_1, X_2, \dots$  be i.i.d real valued random variables with common distribution  $\mu$ , and let

$$(61) \quad \nu_n = \frac{1}{n}(\delta_{X_1} + \dots + \delta_{X_n})$$

be the empirical distribution of  $X_1, \dots, X_n$ , which is a random probability measure on  $\mathbb{R}$ . Then  $\{\nu_n\}$  satisfies a large deviation principle with rate function  $I(\nu) = S(\nu, \mu)$  (called the "relative entropy") given by

$$(62) \quad I(\nu) = \begin{cases} \int p(x) \log(p(x)) d\mu(x), & \text{if } \nu = p\mu \\ +\infty, & \text{otherwise.} \end{cases}$$

Concretely, this means the following. Consider the set  $\mathcal{M}$  of probability measures on  $\mathbb{R}$  with the weak topology (which is a metrizable

topology, e.g. by the Lévy metric). Then for closed  $F$ , open  $G$  in  $\mathcal{M}$  Sanov's theorem yields

$$(63) \quad \limsup_{n \rightarrow \infty} \frac{1}{n} \log P(\nu_n \in F) \leq - \inf_{\nu \in F} S(\nu, \mu)$$

$$(64) \quad \liminf_{n \rightarrow \infty} \frac{1}{n} \log P(\nu_n \in G) \geq - \inf_{\nu \in G} S(\nu, \mu).$$

### 3. BACK TO RANDOM MATRICES

Consider again the space  $\mathcal{H}_N$  of Hermitian matrices equipped with the probability measure  $P_N$  having density

$$(65) \quad dP_N(A) = \text{const} \cdot e^{-\frac{N}{2} \text{Tr}(A^2)} dA.$$

The eigenvalue distribution  $\tilde{P}_N$  on  $\mathbb{R}^N$  is defined by

$$(66) \quad \tilde{P}_N(B) := P_N\{A \in \mathcal{H}_N : (\lambda_1(A), \dots, \lambda_N(A)) \in B\}.$$

One knows that  $\tilde{P}$  is absolutely continuous with respect to Lebesgue measure on  $\mathbb{R}^N$  and has density

$$(67) \quad d\tilde{P}_N(\lambda_1, \dots, \lambda_N) = C_N \cdot e^{-\frac{N}{2} \sum_{i=1}^N \lambda_i^2} \prod_{i < j} (\lambda_i - \lambda_j)^2 \prod_{i=1}^N d\lambda_i,$$

where

$$(68) \quad C_N = \frac{N^{N^2/2}}{(2\pi)^{N/2} \prod_{j=1}^N j!}.$$

We want to establish a large deviation principle for the empirical eigenvalue distribution

$$(69) \quad \mu_A = \frac{1}{N} (\delta_{\lambda_1(A)} + \dots + \delta_{\lambda_N(A)})$$

of a random matrix in  $\mathcal{H}_N$ .

Heuristics for the rate function are as follows. We have

(70)

$$(71) \quad \begin{aligned} P_N(\mu_A \approx \nu) &= \tilde{P}_N\left(\frac{1}{N} (\delta_{\lambda_1(A)} + \dots + \delta_{\lambda_N(A)}) \approx \nu\right) \\ &= C_N \cdot \int_{\{\frac{1}{N} (\delta_{\lambda_1(A)} + \dots + \delta_{\lambda_N(A)}) \approx \nu\}} e^{-\frac{N}{2} \sum \lambda_i^2} \prod_{i < j} (\lambda_i - \lambda_j)^2 \prod_{i=1}^N d\lambda_i. \end{aligned}$$

Now for  $\frac{1}{N} (\delta_{\lambda_1(A)} + \dots + \delta_{\lambda_N(A)}) \approx \nu$ ,

$$(72) \quad -\frac{N}{2} \sum_{i=1}^N \lambda_i^2 = -\frac{N^2}{2} \frac{1}{N} \sum_{i=1}^N \lambda_i^2$$

is a Riemann sum for the integral  $\int x^2 d\nu(x)$ . Moreover

$$(73) \quad \prod_{i < j} (\lambda_i - \lambda_j)^2 = \exp\left(\sum_{i < j} \log |\lambda_i - \lambda_j|^2\right) = \exp\left(\sum_{i \neq j} \log |\lambda_i - \lambda_j|\right)$$

is a Riemann sum for  $N^2 \iint \log |x - y| d\nu(x) d\nu(y)$ .

Hence, heuristically, we expect that

$$(74) \quad P_N(\mu_A \approx \nu) \sim e^{-N^2 I(\nu)},$$

with

$$(75) \quad I(\nu) = - \iint \log |x - y| d\nu(x) d\nu(y) + \frac{1}{2} \int x^2 d\nu(x) - \lim_{N \rightarrow \infty} \frac{1}{N^2} \log C_N.$$

The value of the limit can be explicitly computed as  $6/8$ .

This argument is made rigorous in the following theorem of Ben-Arous and Guionnet from 1997.

Put

$$(76) \quad I(\nu) = - \iint \log |x - y| d\nu(x) d\nu(y) + \frac{1}{2} \int x^2 d\nu(x) - 6/8.$$

Then:

- (1)  $I : \mathcal{M} \rightarrow [0, \infty]$  is a well-defined, convex, good function on the space of real probability measures. It has unique minimum value 0 which occurs at the Wigner semicircle distribution  $\mu_W$ .
- (2) The empirical eigenvalue distribution satisfies a large deviation principle with respect to  $\tilde{P}_N$  with rate function  $I$ : for any open  $G$  and closed  $F$  in  $\mathcal{M}$ 
  - $\liminf_{N \rightarrow \infty} \frac{1}{N^2} \log \tilde{P}_N\left(\frac{\delta_{\lambda_1 + \dots + \delta_{\lambda_N}}}{N} \in G\right) \geq - \inf_{\nu \in G} I(\nu)$ .
  - $\limsup_{N \rightarrow \infty} \frac{1}{N^2} \log \tilde{P}_N\left(\frac{\delta_{\lambda_1 + \dots + \delta_{\lambda_N}}}{N} \in F\right) \leq - \inf_{\nu \in F} I(\nu)$ .