

Lecture 5

①

Now that we have some integration tools at our disposal, we can consider some integrals for moments and statistics of the MP distribution.

Moments of the MP distribution

Proposition: For the standard MP distribution F_y with index $y > 0$ and $\sigma^2 = 1$, it holds for any analytic function f on a domain containing the interval $[a, b] = [(1 - \sqrt{y})^2, (1 + \sqrt{y})^2]$

$$\int f(x) dF_y(x) = -\frac{1}{4\pi i} \oint_{|z|=1} \frac{f((1+\sqrt{y}z)^2)(1-z^2)^2}{z^2(1+\sqrt{y}z)(z+\sqrt{y})} dz$$

Proof: (We will prove a stronger case later).

Let's look at some applications.

Example 1: Logarithms of eigenvalues are often used in multivariate analysis. Set

$$f(x) = \log(x).$$

Assume $0 < y < 1$ so that we don't get zero eigenvalues.

$$\int \log(x) dF_y(x) = \oint_{|z|=1} \frac{\log(1 + \sqrt{y}z) (1 - z^2)^2}{z^2 (1 + \sqrt{y}z)(z + \sqrt{y})} dz \quad \frac{1}{4\pi i}$$

$z \in \mathbb{C}$.

$$= -\frac{1}{4\pi i} \oint_{|z|=1} \frac{\log(1 + \sqrt{y}z) (1 - z^2)^2}{z^2 (1 + \sqrt{y}z)(z + \sqrt{y})} dz$$

$$|1 + \sqrt{y}z|^2$$

$$= (1 + \sqrt{y}z)(1 + \sqrt{y}\bar{z})$$

$$= (1 + \sqrt{y}z)(1 + \sqrt{y}\bar{z})$$

$$= -\frac{1}{4\pi i} \oint_{|z|=1} \frac{\log(1 + \sqrt{y}\bar{z}) (1 - z^2)^2}{z^2 (1 + \sqrt{y}z)(z + \sqrt{y})} dz$$

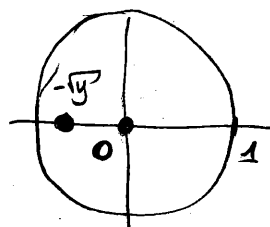
$$= I_1 + I_2$$

Q: When do these integrals have singularities?

There is one at the point $z=0$, due to the $\frac{1}{z^2}$ term.

Another at $z = -\sqrt{y}$.

both within contour $|z|=1$



"order 2 pole"

By Cauchy residue theorem,

$$\int_C f(z) dz = 2\pi i \sum_{a \in C} \text{Res}(f; a)$$

where a are points of singularity.

We need to find the residues at the points $z=0$ and $z=-\sqrt{y}$. We could expand and find the Laurent series but there is an easier way.

Proposition: If f has a pole of order $n \geq 1$ at a . Define $g(z) = (z-a)^n f(z)$ then

$$\text{Res}(f; a) = \frac{1}{(n-1)!} \lim_{z \rightarrow a} g^{(n-1)}(z).$$

Proof: Remember that the residue is the term c_{-1} in the Laurent series expansion of $f(z)$:

$$f(z) = \frac{c_{-n}}{(z-a)^n} + \dots + \frac{c_{-1}}{z-a} + a_0 + \dots$$

$$\text{So } g(z) = c_{-n} + \dots + c_{-1}(z-a)^{n-1} + c_0(z-a)^n + \dots$$

$$\text{and } g^{(n-1)}(z) = (n-1)! c_{-1} + n(n-1) \dots 2 \cdot c_0(z-a) + \dots$$

$$\text{Hence, } \lim_{z \rightarrow a} g^{(n-1)}(z) = g^{(n-1)}(a) = (n-1)! c_{-1}. \quad \square$$

Applying this proposition at $a = -\sqrt{y}$

$$\lim_{z \rightarrow -\sqrt{y}} \frac{\log(1+\sqrt{y}z)(1-z^2)^2}{z^2(1+\sqrt{y}z)(z+\sqrt{y})} \cdot \cancel{(z-\sqrt{y})} = \frac{\log(1-y)(1-y)^2}{y(1-y)} \\ = \log(1-y) \frac{(1-y)}{y}$$

The singularity at $a=0$ is of order 2, so

$$g(z) = \cancel{z^2} \frac{\log(1+\sqrt{y}z)(1-z^2)^2}{\cancel{z^2}(1+\sqrt{y}z)(z+\sqrt{y})} \\ = \frac{\log(1+\sqrt{y}z)(1-z^2)^2}{(1+\sqrt{y}z)(z+\sqrt{y})}$$

$$g'(z) = \frac{\sqrt{y}(1-z^2)^2}{(\sqrt{y}+z)(1+\sqrt{y}z)^2} - \frac{4z(1-z^2)\log(1+\sqrt{y}z)}{(\sqrt{y}+z)(1+\sqrt{y}z)} - \\ \frac{\sqrt{y}(1-z^2)^2\log(1+\sqrt{y}z)}{(\sqrt{y}+z)(1+\sqrt{y}z)^2} - \frac{(1-z^2)^2\log(1+\sqrt{y}z)}{(\sqrt{y}+z)^2(1+\sqrt{y}z)}$$

$$g'(0) = \frac{\sqrt{y}}{\sqrt{y}} - 0 - 0 - 0 = 1.$$

So by the residue theorem $I_1 = -\frac{1}{4\pi i} \left[2\pi i \cdot \left(\log(1-y) \frac{(1-y)}{y} + 1 \right) \right]$

$$= -\frac{1}{2} \left(\log(1-y) \frac{(1-y)}{y} + 1 \right)$$

Now for I_2 we have

$$I_2 = -\frac{1}{4\pi i} \oint_{|z|=1} \frac{\log(1+\sqrt{y}\bar{z})(1-z^2)^2}{z^2(1+\sqrt{y}z)(z+\sqrt{y})} dz.$$

we shall make the change of variable $s = \bar{z}$ and notice that since $|z|=1$, we have.

$$\frac{1}{z} = \frac{1}{e^{i\theta}} = e^{-i\theta} = \bar{z}$$

$$\text{So } I_2 = -\frac{1}{4\pi i} \oint_{|s|=1} \frac{\log(1+\sqrt{y}s)(1-(\frac{1}{s})^2)^2}{(\frac{1}{s})^2(1+\sqrt{y}(\frac{1}{s}))(\frac{1}{s}+\sqrt{y})} (-\frac{1}{s^2}) ds.$$

and this can be shown to be

$$I_2 = I_1.$$

$$\text{hence. } I = -\log(1-y) \frac{(1-y)}{y} - 1.$$

Example 2. We can calculate the mean of the MP distribution.⁶

For all $y > 0$,

$$\int x dF_y(x) = 1.$$

Proof: This can be shown in the same way as Example 1.

For any monomial function $f(x) = x^k$ for $k \in \mathbb{N}$, the residue approach becomes tedious. There is a direct proof as well.

(Bai & Silvestri 2010; Lemma 3.1).

Proposition: The moments of the standard MP distribution

$$\beta_k := \int x^k dF_y(x) = \sum_{r=0}^{k-1} \frac{1}{r+1} \binom{k}{r} \binom{k-1}{r} y^r.$$

Proof: $p_y(x) = \begin{cases} \frac{1}{2\pi xy} \sqrt{(b-x)(x-a)}, & a \leq x \leq b \\ 0, & \text{otherwise} \end{cases}$

density.

$$a = (1 - \sqrt{y})^2 \\ b = (1 + \sqrt{y})^2$$

$$\beta_k = \frac{1}{2\pi y} \int_a^b x^{k-1} \sqrt{(b-x)(x-a)} dx$$

set $x = 1 + y + z$, $dx = dz$

$$x = a \Rightarrow (1 - \sqrt{y})^2 = 1 + y + z \\ z = (1 - \sqrt{y})^2 - 1 - y \\ = -2\sqrt{y}.$$

$$(b-x)(x-a) = (2\sqrt{y} - z)(2\sqrt{y} + z) \\ = (4y - z^2)$$

$$x = b \Rightarrow z = 2\sqrt{y}.$$

$$\text{So } \beta_k = \frac{1}{2\pi y} \int_{-2\sqrt{y}}^{2\sqrt{y}} (1+y+z)^{k-1} \sqrt{4y-z^2} dz$$

Recall $(a+b)^\alpha = \sum_{k=0}^{\alpha} \binom{\alpha}{k} a^k b^{\alpha-k}$ $\binom{\alpha}{k} = \frac{\alpha!}{k!(\alpha-k)!}$

$$\begin{aligned} a &= 1+y \\ b &= z \\ \alpha &= k-1 \end{aligned}$$

$$\begin{aligned} &= \frac{1}{2\pi y} \int_{-2\sqrt{y}}^{2\sqrt{y}} \sum_{\ell=0}^{k-1} \binom{k-1}{\ell} (1+y)^{k-1-\ell} z^\ell \sqrt{4y-z^2} dz \\ &= \frac{1}{2\pi y} \sum_{\ell=0}^{k-1} \binom{k-1}{\ell} (1+y)^{k-1-\ell} \int_{-2\sqrt{y}}^{2\sqrt{y}} z^\ell \sqrt{4y-z^2} dz. \end{aligned}$$

set $z = 2\sqrt{y}u$, $dz = 2\sqrt{y}du$. $z = -2\sqrt{y} \Rightarrow u = -1$.
 $z = 2\sqrt{y} \Rightarrow u = 1$.
 $\Rightarrow 4y - z^2 = 1 - u^2$

$$\begin{aligned} &= \frac{1}{2\pi y} \sum_{\ell=0}^{k-1} \binom{k-1}{\ell} (1+y)^{k-1-\ell} 2\sqrt{y} (2\sqrt{y})^\ell \int_{-1}^1 u^\ell \sqrt{1-u^2} du \\ &= \frac{1}{2\pi y} \sum_{\ell=0}^{k-1} \binom{k-1}{\ell} (1+y)^{k-1-\ell} (4y)^{\frac{\ell+1}{2}} \int_{-1}^1 u^\ell \sqrt{1-u^2} du. \\ &= \frac{1}{2\pi y} \sum_{\ell=0}^{[(k-1)/2]} \binom{k-1}{2\ell} (1+y)^{k-1-2\ell} (4y)^{\ell+1} \int_{-1}^1 u^{2\ell} \sqrt{1-u^2} du. \end{aligned}$$

set $u = \sqrt{w}$, $du = \frac{1}{2} \frac{1}{\sqrt{w}} dw$. $u = -1, w = 1$

$$= \frac{1}{2\pi y} \sum_{\ell=0}^{[(k-1)/2]} \binom{k-1}{2\ell} (1+y)^{k-1-2\ell} (4y)^{\ell+1} \int_0^1 w^{\ell-1/2} \sqrt{1-w} dw.$$

$$\text{As } \int_0^1 w^{l-\frac{1}{2}} \sqrt{1-w} dw = \frac{\sqrt{\pi} \Gamma(l+\frac{1}{2})}{2 \Gamma(2+l)}$$

if $l > -\frac{1}{2}$.

$$\Gamma(l+\frac{1}{2}) = \frac{(2l)!}{4^l l!} \sqrt{\pi}$$

$$\Gamma(n) = (n-1)!$$

$$\Gamma(t) = \int_0^\infty x^{t-1} e^{-x} dx.$$

$$= \sum_{l=0}^{[(k-1)/2]} \frac{1}{\cancel{2\pi y}} \frac{(k-1)!}{\cancel{(2l)!} ((k-1)-2l)!} \frac{\sqrt{\pi} \cancel{(2l)!} \sqrt{\pi}}{\cancel{2} \cancel{l!} (l+1)!} \cancel{y^{l+1}} y^{l+1} (1+y)^{k-1-2l}.$$

$$= \sum_{l=0}^{[(k-1)/2]} \frac{(k-1)!}{l! (l+1)! (k-1-2l)!} y^l (1+y)^{k-1-2l}.$$

$$\text{As } (1+y)^{k-1-2l} = \sum_{s=0}^{k-1-2l} \binom{k-1-2l}{s} y^s = \sum_{s=0}^{k-1-2l} \frac{(k-1-2l)!}{s! (k-1-2l-s)!} y^s$$

$$= \sum_{l=0}^{[(k-1)/2]} \frac{(k-1)!}{l! (l+1)! (k-1-2l)!} y^l \sum_{s=0}^{k-1-2l} \frac{(k-1-2l)!}{s! (k-1-2l-s)!} y^s$$

$$= \sum_{l=0}^{[(k-1)/2]} \sum_{s=0}^{k-1-2l} \frac{(k-1)!}{l! (l+1)! s! (k-1-2l-s)!} y^{l+s}$$

Subst. $r = l+s$

$$s=0 \Rightarrow r=l \quad s=k-1-2l \Rightarrow r=k-1-l$$

$$= \sum_{l=0}^{[(k-1)/2]} \sum_{r=l}^{k-1-l} \frac{(k-1)!}{l! (l+1)! s! (k-1-r-l)!} y^r$$

$$= \frac{1}{k} \sum_{r=0}^{k-1} \binom{k}{r} y^r \sum_{l=0}^{\min(r, k-r)} \binom{r}{l} \binom{k-r}{k-r-l-1}$$

$$= \frac{1}{k} \sum_{r=0}^{k-1} \binom{k}{r} \binom{r}{r+1} y^r = \sum_{r=0}^{k-1} \frac{1}{r+1} \binom{k}{r} \binom{k-1}{r} y^r.$$

Fubini theorem for sequences: If $\sum_{n=0}^{\infty} \sum_{m=0}^{\infty} |a_{nm}| < \infty$ then

$$\sum_{m=0}^{\infty} \sum_{n=0}^{\infty} a_{nm} = \sum_{n=0}^{\infty} \sum_{m=0}^{\infty} a_{nm}.$$

How did I use that?

$$\sum_{l=0}^{\lfloor (k-1)/2 \rfloor} \sum_{r=l}^{k-1-l} \frac{(k-1)!}{l!(l+1)!(r-l)!(k-1-r-l)!} y^r$$

$$= \sum_{l=0}^{\infty} \mathbb{1}(l \leq \lfloor (k-1)/2 \rfloor) \sum_{r=0}^{\infty} \mathbb{1}(r \geq l) \mathbb{1}(r \leq k-1-l)$$

$$= \sum_{l=0}^{\infty} \sum_{r=0}^{\infty} \mathbb{1}(l \leq \lfloor (k-1)/2 \rfloor) \mathbb{1}(l \leq r) \mathbb{1}(l \leq k-1-r) \mathbb{1}(r \leq k-1).$$

$$= \sum_{r=0}^{k-1} \sum_{l=0}^{\min(r, k-1-r)} \left\{ \frac{k!}{r!(k-r)!} y^r \frac{(k-1)! r! (k-r)!}{k! l! (l+1)! (r-l)! (k-1-r-l)!} \right\}$$

$$= \binom{k}{r} y^r \frac{1}{k} \cdot \frac{r!}{l!(r-l)!} \cdot \frac{(k-r)!}{(l+1)!(k-1-r-l)!}$$

$$k-r-(k-1-r-l) = l+1.$$

$$= \frac{1}{k} \binom{k}{r} y^r \binom{r}{l} \binom{k-r}{k-1-r-l}$$

$$= \sum_{r=0}^{k-1} \frac{1}{k} \binom{k}{r} y^r \sum_{l=0}^{\min(r, k-1-r)} \binom{r}{l} \binom{k-r}{k-1-r-l}$$

Generalised MP distribution

Previously, we've seen the case where the population covariance matrix has the simple form $\Sigma = \sigma^2 I_p$.

We can consider a slightly more general case if we make the assumption that the observation vectors $\{y_k\}_{1 \leq k \leq n}$ can be represented as

$$y_k = \Sigma^{1/2} x_k \quad x_k \text{ iid}, \quad \Sigma^{1/2} \text{ nonneg. sqroot of } \Sigma.$$

This gives the associated covariance matrix

$$\begin{aligned} \tilde{B}_n &= \frac{1}{n} \sum_{k=1}^n y_k y_k^* = \Sigma^{1/2} \left(\frac{1}{n} \sum_{k=1}^n x_k x_k^* \right) \Sigma^{1/2} \\ &= \Sigma^{1/2} S_n \Sigma^{1/2}. \end{aligned}$$

S_n is the sample covariance matrix with iid components

The eigenvalues of \tilde{B}_n are the same as $S_n \Sigma$.

11

The following result holds for $B_n = S_n T_n$ for general nonnegative definite matrix T_n . ($T_n = \Sigma$ is a special case)

Theorem Let S_n be the sample covariance matrix $S_n = \frac{1}{n} \sum_{i=1}^n x_i x_i^*$ with IID components and let (T_n) be a sequence of nonnegative definite Hermitian matrices of size $p \times p$.

Define $B_n = S_n T_n$ and assume:

- (1) The entries (x_{jk}) of the data matrix $X = (x_1, \dots, x_n)$ are IID with mean zero and variance 1.
- (2) The data dimension to sample size ratio $p/n \rightarrow y > 0$ as $n \rightarrow \infty$.
- (3) The sequence (T_n) is either deterministic or independent of (S_n) .
- (4) Almost surely, the sequence $(H_n = F^{T_n})$ of the ESD of (T_n) weakly converges to a non-random probability measure H .

Then, almost surely, F^{B_n} weakly converges to a non-random probability measure $F_{y,H}$. Its Stieltjes transform is given

by
$$S(z) = \int \frac{1}{t(1-y-yzs(z))-z} dH(t), \quad z \in \mathbb{C}_+. \quad (*)$$

Notice that the ST of $F_{y,H}$ is implicitly defined.

It can be shown that a unique solution exists but, unfortunately, no closed-form solution exists.

(see Silverstein & Combettes 1992.)

There is a better way to present the ST of $F_{y,H}$.

Consider for B_n a companion matrix

$$\underline{B_n} = \frac{1}{n} X^* T X$$

Size $n \times n$.

Both matrices share the same nonzero eigenvalues so

Their ESD satisfy

$$n F^{\underline{B_n}} - p F^{B_n} = (n-p) \delta_0$$

Note: Given two matrices $A_{p \times q}$ and $B_{q \times p}$ where $p \geq q$, eigenvalues of AB is that of BA augmented by $p-q$ zeros.

$$B_n = S_n T_n = \frac{1}{n} X X^* T_n. \quad \underline{B_n} = \frac{1}{n} X^* T X$$

X $p \times n$ matrix

When $p/n \rightarrow y > 0$, F^{Bn} has limit $F_{C,H}$ if and only if \underline{F}^{Bn} has limit $\underline{F}_{C,H}$. In this case, the limit satisfies

$$\underline{F}_{C,H} - y F_{C,H} = (1-y) S_0.$$

and their ST are related by

$$\underline{S}(z) = -\frac{1-y}{z} + y S(z).$$

Now substituting \underline{S} for S in (*) yields

$$\underline{S}(z) = \left(z - y \int \frac{t}{1+t\underline{S}(z)} dH(t) \right)^{-1}$$

solving in z gives

$$z = -\frac{1}{\underline{S}(z)} + y \int \frac{t}{1+t\underline{S}(z)} dH(t). \quad (**)$$

which defines the inverse function of \underline{S} .

(*) is called the Marcenko-Pastur equation and

(**) is the Silverstein equation.

Limiting spectral distribution for Random Fisher matrices

In the univariate case, when we need to test equality between the variances of 2 Gaussian populations, a Fisher statistic of the form S_1^2/S_2^2 is used where S_i^2 are estimators of the unknown variances in the two populations.

The equivalent for the multivariate setting is:

Take two independent samples $\{X_1, X_2, \dots, X_{n_1}\}$ and $\{Y_1, Y_2, \dots, Y_{n_2}\}$ both from p -dimensional population with iid components and finite second moment.

$$S_1 = \frac{1}{n_1} \sum_{k=1}^{n_1} X_k X_k^*$$

$$S_2 = \frac{1}{n_2} \sum_{k=1}^{n_2} Y_k Y_k^*$$

Then $F_n = S_1 S_2^{-1}$ is called a Fisher matrix. $n = (n_1, n_2)$
(Note: need $p \leq n_2$ so that S_2 invertible)

Let $s > 0$ and $0 < t < 1$. The Fisher LSD $F_{s,t}$ is the distribution with density function

$$p_{s,t}(x) = \frac{1-t}{2\pi x(s+tx)} \sqrt{(b-x)(x-a)} \quad a \leq x \leq b.$$

with $a = a(s,t) = \frac{(1-h)^2}{(1-t)^2}$, $b = b(s,t) = \frac{(1+h)^2}{(1-t)^2}$,

$$h = h(s,t) = (s+t-st)^{1/2}.$$

When $s > 1$, $F_{s,t}$ has a mass at $x=0$ of value $1-1/s$ with the total mass of the rest of the distribution for $x > 0$ is equal to $1/s$.

The Fisher LSD has many similarities to the standard MP distribution. This is not a coincidence as the MP LSD F_y is the Fisher LSD $F_{y,0}$ (ie. $s,t = y,0$)

Also note $t \rightarrow 1_-$, $a(s,t) \rightarrow \frac{1}{2}(1-s)^2$, $\underbrace{b(s,t) \rightarrow \infty}_{\text{supp}(F_{s,t}) \text{ becomes unbounded.}}$

Theorem: For an analytic function f on a domain containing $[a, b]$ (as above). We have.

$$\int_a^b f(x) dF_{s,t}(x) = - \frac{h^2(1-t)}{4\pi i} \oint_{|z|=1} \frac{f\left(\frac{1+hz}{(1-t)^2}\right) (1-z^2)^2 dz}{z(1+hz)(z+h)(tz+h)(t+hz)}$$

Proof: Using the density $\beta_{s,t}(x)$

$$I = \int_a^b f(x) dF_{s,t}(x) = \int_a^b f(x) \frac{1-t}{2\pi x(s+xt)} \sqrt{(x-a)(b-x)}' dx.$$

Make change of variable $x = \frac{1+h^2+2h\cos(\theta)}{(1-t)^2} \quad \theta \in (0, \pi)$

$$x-a = \frac{1+h^2+2h\cos(\theta)}{(1-t)^2} - \frac{(1-h)^2}{(1-t)^2} = \frac{2h+2h\cos(\theta)}{(1-t)^2}$$

$$b-x = \frac{(1+h)^2}{(1-t)^2} - \frac{1+h^2+2h\cos(\theta)}{(1-t)^2} = \frac{2h-2h\cos(\theta)}{(1-t)^2}$$

$$\sqrt{(x-a)(b-x)} = \sqrt{\frac{(2h)^2}{(1-t)^4} (1-\cos(\theta))(1+\cos(\theta))}$$

$$= \frac{2h}{(1-t)^2} \sin(\theta)$$

$$x=a \Rightarrow \cos(\theta) = \frac{a(1-t)^2 - (1+h^2)}{2h}$$

$$= 0.$$

$$\Rightarrow \theta = 0$$

$$x=b \Rightarrow \theta = \pi$$

$$dx = \frac{-2h\sin(\theta)}{(1-t)^2} d\theta.$$

hence,

$$I = \frac{2h^2(1-t)}{\pi} \int_0^\pi \frac{f\left(\frac{1+h^2+2h\cos(\theta)}{(1-t)^2}\right) \sin^2(\theta) d\theta}{(1+h^2+2h\cos(\theta))(s(1-t)^2+t(1+h^2+2h\cos(\theta)))}$$

$$= \frac{2h^2(1-t)}{\pi} \frac{1}{2} \int_0^{2\pi} \frac{f\left(\frac{1+h^2+2h\cos(\theta)}{(1-t)^2}\right) \sin^2(\theta) d\theta}{(1+h^2+2h\cos(\theta))(s(1-t)^2+t(1+h^2+2h\cos(\theta)))}$$

Now let $z = e^{i\theta}$,

$$1+h^2+2h\cos(\theta) = |1+hz|^2$$

$$\sin(\theta) = \frac{z - z^{-1}}{2i}$$

$$\log(z) = i\theta \Rightarrow \theta = \frac{1}{i} \log(z) \Rightarrow d\theta = \frac{1}{iz} dz.$$

$$I = - \frac{h^2(1-t)}{4\pi i} \oint_{|z|=1} \frac{f\left(\frac{|1+hz|^2}{(1-t)^2}\right) (1-z^2)^2 dz}{z^3 |1+hz|^2 (s(1-t)^2+t|1+hz|^2)}$$

On $|z|=1$, we have $|1+hz|^2 = (1+hz)(1+hz^{-1})$.

so expanding denominator and simplifying we have result.

Example: Take $(s, t) = (y, 0)$ and we get
the result for MP distribution

Example 2: The first two moments are

$$\int x dF_{s,t}(x) = \frac{1}{1-t}, \quad \int x^2 dF_{s,t}(x) = \frac{h^2 + 1 - t}{(1-t)^3}.$$

Hence the variance equals $h^2/(1-t)^3$.
