

# High Dimensional Data Matrices and Random Matrices

Craig A. Tracy  
UC Davis

YITP @ 40

May 3, 2007

# DATA MATRICES AND PRINCIPAL COMPONENT ANALYSIS (PCA)

Suppose in some measurement (experiment) we have  $p$  variables

$$X_1, X_2, \dots, X_p$$

Theoretically, the  $X_k$  are random variables. The observed data on  $X_k$  can be viewed as a vector  $x_k \in \mathbb{R}^n$ . We form a  $p \times n$  data matrix

$$X = \begin{pmatrix} \longleftarrow & x_1 & \longrightarrow \\ & \vdots & \\ \longleftarrow & x_p & \longrightarrow \end{pmatrix}$$

- Astrophysics example: “In Orion A we have mapped 32 molecular transitions in the 3 mm wavelength band . . . .” Each “map” contains 360 pixels. Here  $p = 32$  and  $n = 360$ .

- Human Gene Structure: Properties of  $p = 38$  genes at  $n = 400$  locations in Europe.
- Finance: Twenty years of stock returns in the S&P 500. Here  $p = 500$  and  $n$  is the number of data points on an individual stock.

The idea of PCA (Hotelling, 1933; see Johnstone [6])

- Reduce dimensionality:  $W = \sum_k v_k X_k$  by requiring

$$\text{var}(W) = \sum_{k,k'} v_k \text{cov}(X_k, X_{k'}) v_{k'}$$

have *maximum variance*. Vector  $v$  is the 1st principal component vector. Then choose successive linear combinations that are orthogonal to previously chosen and maximize variance.

$$\ell_j = \max \{ v^T \Sigma v : v^T v_{j'} = 0, j' < j, |v| = 1 \}$$

## SAMPLE COVARIANCE MATRIX

Use data matrix  $X$  to get  $p \times p$  sample covariance matrix

$$S = \frac{1}{n} X^T X$$

and look for *sample principle components*:

$$S\hat{v}_j = \hat{\ell}_j \hat{v}_j$$

### Spreading of sample eigenvalues

Take  $n = p = 10$  for Wishart distribution:  $X_1, \dots, X_p$  follow a  $p$ -variate Gaussian distribution with  $\Sigma = 1$ :

$$\hat{\ell}_j = .003, .036, .095, .16, .30, .51, .78, 1.12, 1.40, 3.07$$

On the basis of this data, might (erroneously!) conclude population eigenvalues are quite different from each other (they all equal 1).

This spreading of the eigenvalues is the statistics version of

Wigner semicircle

or as its called here

Marčenko-Pastur limit density

$$g^{MP}(t) = \frac{\sqrt{(b_+ - t)(t - b_-)}}{2\pi\gamma t}, \quad b_{\pm} = (1 \pm \sqrt{\gamma})^2$$

where  $p/n \rightarrow \gamma$  as  $n, p \rightarrow \infty$ . When  $n = p$  the density is supported on the interval  $[0, 4]$ .

**Question:** Suppose one sees a largest sample eigenvalue of 4.25. Is this consistent with an identity covariance matrix? (It lies outside the M-P support.)

- $H_0 : \Sigma = I$ . The *null hypothesis*.
- $H_A : \Sigma \neq I$ . The *alternative hypothesis*.

Want

$$\mathbb{P} \left( \hat{\ell}_1 > t | H_0 = W_p(n, I) \right)$$

**Theorem** (Johnstone [5]):

$$\mathbb{P} \left( \hat{\ell}_1 \leq \mu_{np} + \sigma_{np}s | H_0 = W_p(n, I) \right) \longrightarrow F_1(s)$$

where

$$\mu_{np} = \left( \sqrt{n - 1/2} + \sqrt{p - 1/2} \right)^2$$

$$\sigma_{np} = \left( \sqrt{n - 1/2} + \sqrt{p - 1/2} \right) \left( \frac{1}{\sqrt{n - 1/2}} + \frac{1}{\sqrt{p - 1/2}} \right)^{1/3}$$

## WHAT IS $F_1$ ?

$F_1$  is one of three distributions first discovered by **Harold Widom** and **C.T.** [12, 13] in the context of the distribution of the largest eigenvalue in

GOE, GUE, and GSE

$$F_2(s) = \exp\left(-\int_s^\infty (x-s)q(x)^2 dx\right)$$

$$F_1(s)^2 = F_2(s) \exp\left(-\int_s^\infty q(x) dx\right)$$

$$q'' = sq + 2q^3, \text{ Painlevé II equation}$$

$$q(s) \sim \text{Ai}(s) \text{ as } s \rightarrow \infty$$

Appearance of  $F_1$  also in **double Wishart** : Two independent Wishart matrices  $A \sim W_p(n_1, I)$  and  $B \sim W_p(n_2, I)$ . Appears in

- Canonical correlation analysis

**William Chen**, of the IRS, computed tables of the exact distribution in the double Wishart

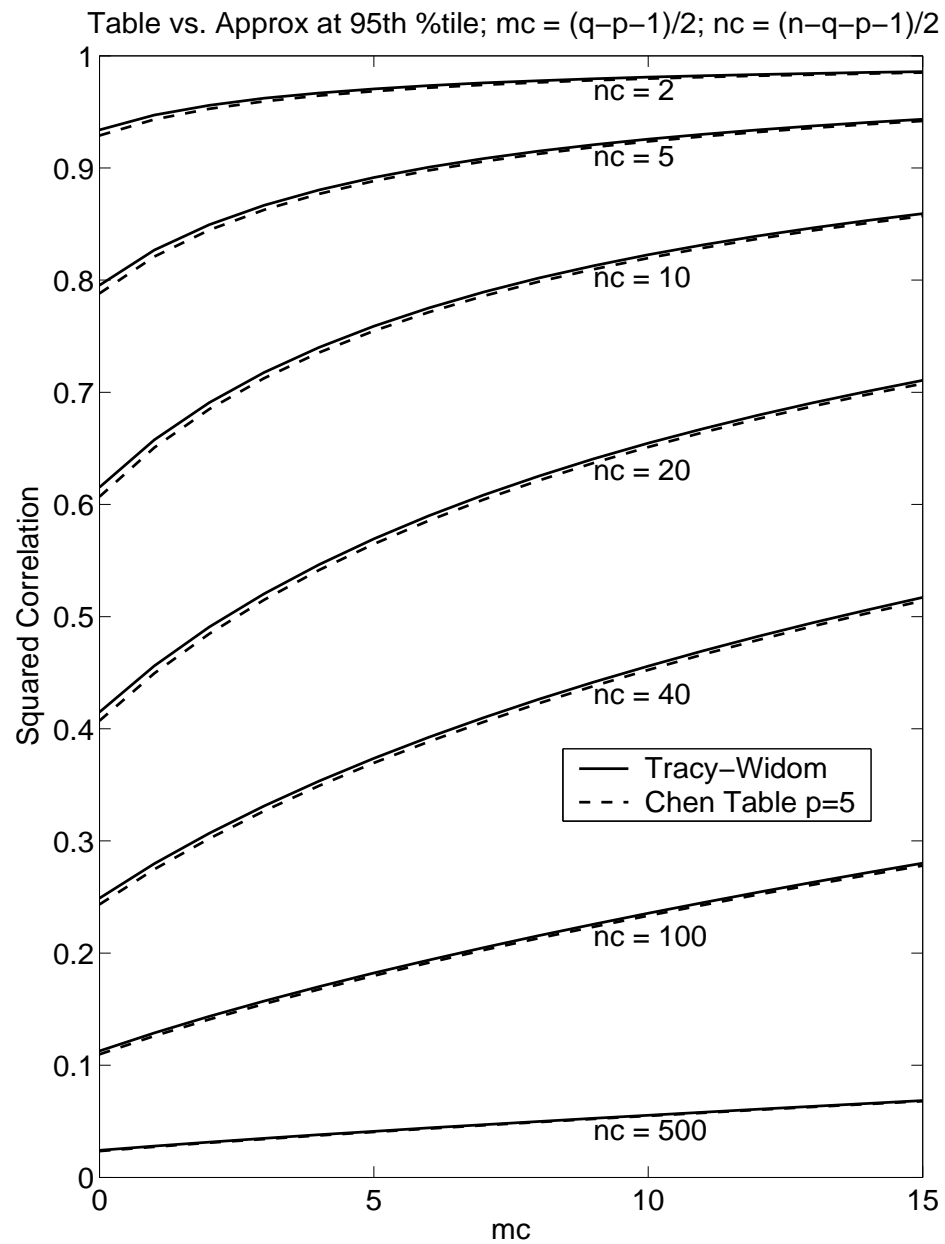
$$m_c = (n_1 - p - 1)/2, \quad n_c = (n_2 - p - 1)/2$$

Johnstone compared this with the TW approximation which is a limit theorem

$$n, p \rightarrow \infty, \quad \frac{p}{n} \rightarrow \gamma < \infty$$

The agreement is good. (See next slide.)

In the earlier example, the TW approximation yields a 6% chance of seeing a value more extreme than 4.25 even if “no structure” is present.



## FURTHER DEVELOPMENTS FOR WISHART DISTRIBUTION

- Johnstone [5], El Karoui [7, 8], Choup [3]: **Second-order accuracy**

$$\left| P \left( n\hat{\ell}_1 \leq \mu_{np} + \sigma_{np}s | H_0 \right) - F_{\beta}(s) \right| \leq Cp^{-2/3}$$

- El Karoui in null case for the largest eigenvalue, proves the limit law for  $0 \leq \gamma \leq \infty$ . This requires additional estimates to allow  $\gamma = \infty$ .
- Soshnikov [11] **removes Gaussian assumption** on the distribution of the matrix elements of  $X$  and only requires odd moments are zero and even moments satisfy a Gaussian type bound. Then for  $\Sigma = I$  with under the restriction that as  $n, p \rightarrow \infty$  that

$$n - p = O(p^{1/3})$$

we get the same limit law described by  $F_1$ .

- **Beyond the Null Hypothesis:** If  $A \sim W_p(n, \Sigma)$ , then joint eigenvalue density

$$c_{p,n,\Sigma} \prod_{j=1}^p l_j^{(n-p-1)/2} \prod_{j < k} |l_j - l_k| \times \int_{\mathcal{O}(p)} e^{-\frac{1}{2} \text{tr}(\Sigma^{-1} Q L Q^T)} dQ,$$

where  $L = \text{diag}(l_1, \dots, l_p)$  and  $dQ$  is normalized Haar measure.

Difficulty is the integral

$$\int_{\mathcal{O}(p)} e^{-\frac{1}{2} \text{tr}(\Sigma^{-1} Q L Q^T)} dQ$$

For **complex data** above integral is replaced by

$$\int_{\mathcal{U}(p)} e^{-\frac{1}{2} \text{tr}(\Sigma^{-1} U L U^*)} dU$$

This integral can be evaluated in terms of determinants:  
HARISH CHANDRA–ITZYKSON–ZUBER integral.

- **Spiked population data:**

$$\text{Eigenvalues of } \Sigma : \lambda_1 > 1 = \lambda_2 = \cdots = \lambda_p$$

When can we detect  $\lambda_1$  from the data? It depends! Baik, Ben Arous, P  ch   [1, 2] describe a PHASE TRANSITION for population covariance matrices of above form (special case of their theorem). For real data, see conjectures of Patterson, Price and Reich [10].

C.N. YANG INSTITUTE FOR  
THEORETICAL PHYSICS

HAPPY 40TH BIRTHDAY!



# References

- [1] J. Baik, G. Ben Arous, and S. Péché, Phase transition of the largest eigenvalue for nonnull complex sample covariance matrices, *Ann. Probab.* **33** (2005), 1643–1697.
- [2] J. Baik, Painlevé formulas of the limiting distributions for nonnull complex sample covariance matrices, *Duke Math. J.* **133** (2006), 205–235.
- [3] L. Choup, Edgeworth expansion of the largest eigenvalue distribution function of GUE and LUE, *Int. Math. Res. Not.* (2006), Art. ID 61049, 32 pp.
- [4] M. Dieng and C. A. Tracy, Application of random matrix theory to multivariate statistics, preprint, arXiv:math.PR/0603543.
- [5] I. M. Johnstone, On the distribution of the largest eigenvalue in principal component analysis, *Ann. Stat.* **29** (2001), 295–327.
- [6] I. M. Johnstone, High dimensional statistical inference and

random matrices, arXiv: math.ST/0611589.

- [7] N. El Karoui, On the largest eigenvalue of Wishart matrices with identity covariance when  $n$ ,  $p$  and  $p/n$  tend to infinity, arXiv: math.ST/0309355.
- [8] N. El Karoui, Tracy-Widom limit for the largest eigenvalue of a large class of complex Wishart matrices, arXiv: math.PR/0503109.
- [9] R. J. Muirhead, *Aspects of Multivariate Statistical Theory*, John Wiley & Sons, Inc., 1982.
- [10] N. Patterson, A. L. Price and D. Reich, Population structure and eigenanalysis, *PLoS Genetics* **2** (2006), 2074–2093.
- [11] A. Soshnikov, A note on universality of the distribution of the largest eigenvalue in certain sample covariance matrices, *J. Statistical Physics* **108** (2002), 1033–1056.
- [12] C. A. Tracy and H. Widom, Level-spacing distributions and the Airy kernel, *Phys. Letts. B* **305** (1993), 115–118; *Commun.*

*Math. Phys.* **159** (1994), 151–174.

- [13] C. A. Tracy and H. Widom, On orthogonal and symplectic matrix ensembles, *Commun. Math. Phys.* **177** (1996), 727–754.