

A Review of Gaussian Random Matrices

Department of Mathematics, Linköping University

Kasper Andersson

LITH-MAT-EX-2020/08-SE

Credits: **16hp**

Level: **G2**

Supervisor: **Martin Singull**,
Department of Mathematics, Linköping University

Examiner: **Xiangfeng Yang**,
Department of Mathematics, Linköping University

Linköping: **2020**

Abstract

While many university students get introduced to the concept of statistics early in their education, random matrix theory (RMT) usually first arises (if at all) in graduate level classes. This thesis serves as a friendly introduction to RMT, which is the study of matrices with entries following some probability distribution. Fundamental results, such as Gaussian and Wishart ensembles, are introduced and a discussion of how their corresponding eigenvalues are distributed is presented. Two well-studied applications, namely neural networks and PCA, are discussed where we present how RMT can be applied.

Keywords:

Random Matrix Theory, Gaussian Ensembles, Covariance, Wishart Ensembles, PCA, Neural Networks

URL for electronic version:

<http://urn.kb.se/resolve?urn=urn:nbn:se:liu:diva-171649>

Sammanfattning

Medan många stöter på statistik och sannolikhetslära tidigt under sina universitetsstudier så är det sällan slumpmatristeori (RMT) dyker upp förän på forskarnivå. RMT handlar om att studera matriser där elementen följer någon sannolikhetsfördelning och den här uppsatsen presenterar den mest grundläggande teorin för slumpmatriser. Vi introducerar Gaussian ensembles, Wishart ensembles samt fördelningarna för dem tillhörande egenvärdena. Avslutningsvis så introducerar vi hur slumpmatriser kan användas i neruonnät och i PCA.

Nyckelord:

Slumpmatristeori, Gaussian Ensembles, Kovarians, Wishart Ensembles, PCA, Neuronnät

URL för elektronisk version:

<http://urn.kb.se/resolve?urn=urn:nbn:se:liu:diva-171649>

Acknowledgements

I would like to thank my supervisor Martin Singull who introduced me to the topics of Random Matrix Theory, which I would never look upon by myself.

I also want to thank my mother, who always supported me throughout my studies.

Contents

1	Introduction	1
1.1	Aims and Outline	2
1.2	Preliminaries	3
2	Random Matrices	9
2.1	Gaussian Ensembles	9
2.2	Eigenvalues	11
2.3	Covariance and Wishart Ensembles	17
3	Applications	23
3.1	Random Matrices in Neural Networks	23
3.1.1	Neural Networks in General	23
3.1.2	RMT Model for Neural Networks	25
3.2	Dimensionality Reduction - PCA	29
3.2.1	Minimizing Mean Squared Error of Projection	29
3.2.2	Maximizing the Variance	30
3.2.3	Statistical inference of PCA	34
4	Ending discussion	37

Chapter 1

Introduction

A random matrix is a matrix \mathbf{M} with its entries m_{ij} being random variables. Two observations of 3×3 random matrices, both with elements being sampled uniformly on the interval $[0,1]$ can be

$$\begin{pmatrix} 0.9575 & 0.9706 & 0.8003 \\ 0.9649 & 0.9572 & 0.1419 \\ 0.1576 & 0.4854 & 0.4218 \end{pmatrix}, \begin{pmatrix} 0.9157 & 0.6557 & 0.9340 \\ 0.7922 & 0.0357 & 0.6787 \\ 0.9595 & 0.8491 & 0.7577 \end{pmatrix}.$$

We can do the same thing for random matrices with its elements being sampled from $\mathcal{N}(0,1)$

$$\begin{pmatrix} 0.7269 & -0.7873 & -1.0689 \\ -0.3034 & 0.8884 & -0.8095 \\ 0.2939 & -1.1471 & -2.9443 \end{pmatrix}, \begin{pmatrix} 1.4384 & 1.3703 & -0.2414 \\ 0.3252 & -1.7115 & 0.3192 \\ -0.7549 & -0.1022 & 0.3129 \end{pmatrix}.$$

Again, the elements of the two matrices differ from each other, which will be the general case when we are working with random matrices.

When we are studying random matrices, we are not interested in a specific sample, but rather a model of all the instances it can take. These models are known as *matrix ensembles*.

As always when we are working with matrices, eigenvalues are of great interest. Many of the results presented in this thesis will therefore be about how the eigenvalues are distributed for different kind of random matrices. Since the entries of random matrices are random themselves, the corresponding eigenvalues will be random as well. A major result in random matrix theory can be observed

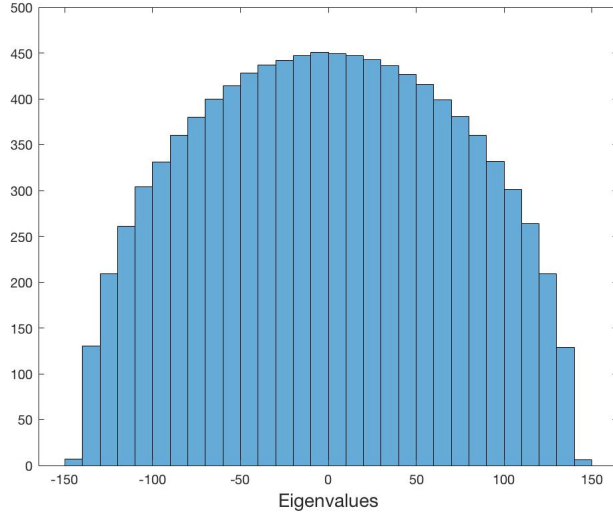


Figure 1.1: Empirical distribution of the eigenvalues for a 10000×10000 random matrix.

by plotting the eigenvalues corresponding to a symmetric matrix with entries being sampled from a standard normal distribution. The occurring elliptical shape, known as *Wigner's semicircle*, is demonstrated in Figure 1.1.

1.1 Aims and Outline

Most of the existing literature associated with random matrices is for graduate studies or research, despite the fact that many results can be well understood and interpreted by students on an undergraduate level. This thesis aims to give the reader a friendly introduction to the topics associated with random matrix theory.

Chapter 1 starts off by introducing the concept of matrix ensembles and briefly generalizes random variables into random matrices. Mathematical results which will be needed to develop the theory of random matrices will be presented here as well.

Chapter 2 collects the core of random matrix theory. The reader is introduced to Gaussian ensembles and to their corresponding eigenvalue distribution. We follow up with the definition of covariance matrices and their properties. Focusing on covariance matrices built upon Gaussian entries, a discussion of how these matrices and their eigenvalues are distributed is represented.

Chapter 3 aims to introduce the reader to how random matrix theory can be applied. First off, the idea of standard neural networks and how they are used is represented. We follow up with the concept of extreme learning machines and summarizes a framework, established by Z. Liao, of how random matrices can be used in neural networks. Lastly, we derive the concept of standard PCA and represent how statistical inference may be worked out when the underlying data follows a normal distribution.

Chapter 4 serves as a final discussion of what we have accomplished in this thesis and how the topics of random matrix could be further developed.

1.2 Preliminaries

Definition 1.1. A matrix \mathbf{M} with entries m_{ij} is **Hermitian** if $m_{ij} = \overline{m_{ji}}$ where \overline{z} is the complex conjugate of z .

Remark: If all entries in the matrix are real, then the matrix being Hermitian is equivalent to it being symmetric.

Definition 1.2. The **spectrum** of a matrix \mathbf{M} is given by the set $\{\lambda \in \mathbb{C} : \det(\mathbf{M} - \lambda \mathbf{I}) = 0\}$ where λ is an **eigenvalue** of \mathbf{M} and \mathbf{I} is the identity matrix.

Theorem 1.1. If X_1, \dots, X_n are i.i.d. random variables with corresponding probability density functions (pdf) $f(x_i)$ then their joint pdf is given by

$$f(x_1, \dots, x_n) = \prod_{i=1}^n f(x_i).$$

Remark: We will simply denote the joint pdf for the matrix entries m_{11}, \dots, m_{nn} of a $n \times n$ matrix \mathbf{M} as $f(\mathbf{M})$.

Definition 1.3. The **expected value** (mean) for a continuous random variable X with corresponding pdf $f(x)$ is given by

$$E[X] = \int_{-\infty}^{\infty} xf(x)dx.$$

Remark: For a discrete random variable, replace integral with proper sum and pdf with probability mass function (pmf).

Definition 1.4. The **variance** for a random variable X is given by

$$\text{Var}[X] = E[(X - E[X])^2].$$

Definition 1.5. A random variable X is **Gaussian** (normal) **distributed**, with mean μ and variance σ^2 if its pdf is given by

$$f(x|\mu, \sigma) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}, \quad x \in (-\infty, \infty)$$

and we denote this by $X \sim \mathcal{N}(\mu, \sigma^2)$.

Definition 1.6. The **gamma function** $\Gamma(\cdot)$ is given by

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

Definition 1.7. We say that X is **chi squared distributed** with n degrees of freedom if its pdf is given by

$$f(x|n) = \frac{1}{\Gamma(\frac{n}{2})2^{\frac{n}{2}}} x^{\frac{n}{2}-1} e^{-\frac{x}{2}}, \quad x \in (0, \infty),$$

which we denote as $X \sim \chi^2(n)$.

Theorem 1.2. If X_1, X_2, \dots, X_n are all i.i.d. $\mathcal{N}(0, 1)$ then it follows that

$$\sum_{i=1}^n X_i^2 \sim \chi^2(n).$$

Definition 1.8. X is **Rayleigh distributed** with scale parameter $b > 0$ if its pdf is given by

$$f(x|b) = \frac{x}{b^2} e^{-\frac{x^2}{2b^2}}, \quad x \in (0, \infty),$$

which we denote as $X \sim \text{Rayleigh}(b)$.

Lemma 1.1. If $X \sim \mathcal{N}(0, \sigma^2)$ and $Y \sim \mathcal{N}(0, \sigma^2)$ are independent, then it follows that

$$\sqrt{X^2 + Y^2} \sim \text{Rayleigh}(\sigma).$$

Proof: With Theorem 1.2 we get

$$\begin{aligned}
 P(\sqrt{X^2 + Y^2} \leq t) &= P(X^2 + Y^2 \leq t^2) \\
 &= P\left(\frac{X^2}{\sigma^2} + \frac{Y^2}{\sigma^2} \leq \frac{t^2}{\sigma^2}\right) \\
 &= \int_0^{\frac{t^2}{\sigma^2}} \frac{1}{2} e^{-\frac{x}{2}} dx \\
 &= 1 - e^{-\frac{t^2}{2\sigma^2}}.
 \end{aligned}$$

Taking derivative with the respect to t finally gives us

$$f(t) = \frac{t}{\sigma^2} e^{-\frac{t^2}{2\sigma^2}}$$

and we are done. \square

Definition 1.9. The *trace* for a $n \times n$ matrix \mathbf{M} with elements m_{ij} is given by

$$\text{tr}(\mathbf{M}) = \sum_{i=1}^n m_{ii}$$

Lemma 1.2. Given a matrix \mathbf{M} with eigenvalues $\lambda_1, \lambda_2, \dots, \lambda_n$ we have that

$$\text{tr}(\mathbf{M}^n) = \sum_{i=1}^n \lambda_i^n.$$

Definition 1.10. A matrix $\mathbf{M} \in \mathbb{R}^{n \times n}$ is *positive semi-definite* if $\mathbf{x}^T \mathbf{M} \mathbf{x} \geq 0$ for all vectors $\mathbf{x} \in \mathbb{R}^{n \times 1}$.

Lemma 1.3. A symmetric matrix is positive semi-definite if and only if all of its eigenvalues are non-negative.

Lemma 1.4. The eigenvectors corresponding to different eigenvalues of a symmetric matrix are orthogonal to each other.

Proof: Consider a symmetric matrix \mathbf{X} with eigenvalues $\lambda_i \neq \lambda_j$ and corresponding eigenvectors $\mathbf{x}_i, \mathbf{x}_j$. We have that $\mathbf{x}_j^T \mathbf{X} \mathbf{x}_i = \mathbf{x}_j^T \lambda_i \mathbf{x}_i$ and $\mathbf{x}_i^T \mathbf{X} \mathbf{x}_j = \mathbf{x}_i^T \lambda_j \mathbf{x}_j$ and since \mathbf{X} is symmetric, $\mathbf{x}_j^T \mathbf{X} \mathbf{x}_i = \mathbf{x}_i^T \mathbf{X} \mathbf{x}_j$. Together, this gives us the equality

$$(\lambda_i - \lambda_j) \mathbf{x}_i^T \mathbf{x}_j = 0,$$

and thus, $\mathbf{x}_i^T \mathbf{x}_j = 0$. \square

Definition 1.11. A sequence of random variables X_n **converges almost surely (a.s.)** to the random variable X as $n \rightarrow \infty$ if

$$P(\{\omega : X_n(\omega) \rightarrow X(\omega) \text{ as } n \rightarrow \infty\}) = 1.$$

Definition 1.12. A sequence of random variables X_n **converges in probability** to the random variable X as $n \rightarrow \infty$ if $\forall \varepsilon > 0$,

$$P(|X_n - X| > \varepsilon) \rightarrow 0 \quad \text{as } n \rightarrow \infty.$$

Definition 1.13. A real valued function f is said to be **Lipschitz continuous** if there exists a constant C for all $x \neq y$ such that

$$\left| \frac{f(x) - f(y)}{x - y} \right| \leq C.$$

Theorem 1.3. Let \mathbb{V} be a linear subspace of \mathbb{R}^n , $\mathbf{u} \in \mathbb{R}^n$ and $\mathbf{u}_{||\mathbb{V}}$ the orthogonal projection of \mathbf{u} onto \mathbb{V} . Then, for all $\mathbf{v} \in \mathbb{V}$, we have that

$$\|\mathbf{u}_{||\mathbb{V}} - \mathbf{v}\| \leq \|\mathbf{u} - \mathbf{v}\|.$$

Proof: We have

$$\mathbf{u} - \mathbf{v} = (\mathbf{u} - \mathbf{u}_{||\mathbb{V}}) + (\mathbf{u}_{||\mathbb{V}} - \mathbf{v}) \Rightarrow \|\mathbf{u} - \mathbf{v}\|^2 = \|(\mathbf{u} - \mathbf{u}_{||\mathbb{V}}) + (\mathbf{u}_{||\mathbb{V}} - \mathbf{v})\|^2$$

expanding the right hand side gives us

$$\begin{aligned} \|(\mathbf{u} - \mathbf{u}_{||\mathbb{V}}) + (\mathbf{u}_{||\mathbb{V}} - \mathbf{v})\|^2 &= ((\mathbf{u} - \mathbf{u}_{||\mathbb{V}}) + (\mathbf{u}_{||\mathbb{V}} - \mathbf{v}))^T \\ &\quad \cdot ((\mathbf{u} - \mathbf{u}_{||\mathbb{V}}) + (\mathbf{u}_{||\mathbb{V}} - \mathbf{v})) \\ &= (\mathbf{u} - \mathbf{u}_{||\mathbb{V}})^T (\mathbf{u} - \mathbf{u}_{||\mathbb{V}}) \\ &\quad + (\mathbf{u} - \mathbf{u}_{||\mathbb{V}})^T (\mathbf{u}_{||\mathbb{V}} - \mathbf{v}) \\ &\quad + (\mathbf{u}_{||\mathbb{V}} - \mathbf{v})^T (\mathbf{u} - \mathbf{u}_{||\mathbb{V}}) \\ &\quad + (\mathbf{u}_{||\mathbb{V}} - \mathbf{v})^T (\mathbf{u}_{||\mathbb{V}} - \mathbf{v}) \\ &= (\mathbf{u} - \mathbf{u}_{||\mathbb{V}})^T (\mathbf{u} - \mathbf{u}_{||\mathbb{V}}) \\ &\quad + (\mathbf{u}_{||\mathbb{V}} - \mathbf{v})^T (\mathbf{u}_{||\mathbb{V}} - \mathbf{v}) \\ &= \|\mathbf{u} - \mathbf{u}_{||\mathbb{V}}\|^2 + \|\mathbf{u}_{||\mathbb{V}} - \mathbf{v}\|^2, \end{aligned}$$

as a result of orthogonality between $(\mathbf{u} - \mathbf{u}_{||\mathbb{V}})$ and $(\mathbf{u}_{||\mathbb{V}} - \mathbf{v})$. Finally, due to the non-negative euclidean norm, we obtain

$$\|\mathbf{u} - \mathbf{v}\|^2 \geq \|\mathbf{u}_{||\mathbb{V}} - \mathbf{v}\|^2, \quad (1.1)$$

which is equivalent to

$$\|\mathbf{u} - \mathbf{v}\| \geq \|\mathbf{u}_{||\mathbb{V}} - \mathbf{v}\|, \quad (1.2)$$

and we are done. \square

Definition 1.14. *The **Frobenius norm**, denoted $\|\cdot\|_F$, of a matrix \mathbf{A} is given by*

$$\|\mathbf{A}\|_F = \sqrt{\text{tr}(\mathbf{A}\mathbf{A}^T)}.$$

Chapter 2

Random Matrices

2.1 Gaussian Ensembles

A cornerstone in random matrix theory (RMT) is to study the spectrum of random matrices. It's difficult to say anything about the eigenvalues for an arbitrary matrix with random entries, so we will focus on symmetric random matrices.

Definition 2.1. A $n \times n$ matrix \mathbf{M} with entries m_{ij} is a **Wigner matrix** if $m_{ij} = m_{ji}$ and m_{ij} are randomly i.i.d. up to symmetry.

E. P. Wigner used symmetric random matrices to study nuclear energy levels [23]. Another assumption which gives room for more calculations is to consider the entries being Gaussian.

Definition 2.2. Let \mathbf{M} be a real Wigner matrix with diagonal entries $m_{ii} \sim \mathcal{N}(0, 1)$ and $m_{ij} \sim \mathcal{N}(0, 1/2)$, $i \neq j$. \mathbf{M} is said to belong to the **Gaussian orthogonal ensemble (GOE)**.

One can easily achieve the required properties for a matrix to be considered as a GOE matrix. Consider a random $n \times n$ matrix \mathbf{A} with all entries being $\mathcal{N}(0, 1)$. Then

$$\mathbf{M} = \frac{\mathbf{A} + \mathbf{A}^T}{2} \tag{2.1}$$

is a GOE matrix.

We can derive a more practical way to describe the GOE.

Theorem 2.1. *The joint pdf for the Gaussian orthogonal ensemble is given by*

$$f(\mathbf{M}) = \frac{1}{2^{n/2}} \frac{1}{\pi^{n(n+1)/4}} e^{-\frac{1}{2} \text{tr}(\mathbf{M}^2)}.$$

Proof: Let m_{ij} be the entries for a $n \times n$ matrix \mathbf{M} generated by GOE. Recall, all entries in \mathbf{M} are independent so their joint pdf is given by

$$f(\mathbf{M}) = \left(\frac{1}{\sqrt{2\pi}} \right)^n \prod_{i=j} e^{-\frac{m_{ij}^2}{2}} \left(\frac{1}{\sqrt{\pi}} \right)^{\frac{n(n-1)}{2}} \prod_{1 \leq i < j \leq n} e^{-m_{ij}^2}. \quad (2.2)$$

Using exponential rules and the fact that \mathbf{M} is symmetric, we can rewrite the joint pdf as

$$\begin{aligned} C_n \prod_{i=j} e^{-\frac{m_{ij}^2}{2}} \prod_{1 \leq i < j \leq n} e^{-m_{ij}^2} &= C_n e^{-\frac{1}{2} \sum_{i=j} m_{ij}^2 - \sum_{1 \leq i < j \leq n} m_{ij}^2} \\ &= C_n e^{-\frac{1}{2} \left(\sum_{i=j} m_{ij}^2 + 2 \sum_{1 \leq i < j \leq n} m_{ij}^2 \right)} \\ &= C_n e^{-\frac{1}{2} \text{tr}(\mathbf{M}^2)}, \end{aligned} \quad (2.3)$$

where

$$\begin{aligned} C_n &= \left(\frac{1}{\sqrt{2\pi}} \right)^n \left(\frac{1}{\sqrt{\pi}} \right)^{\frac{n(n-1)}{2}} = \frac{1}{2^{n/2}} \frac{1}{\pi^{n/2}} \frac{1}{\pi^{n(n-1)/4}} \\ &= \frac{1}{2^{n/2}} \frac{1}{\pi^{n(n+1)/4}} \end{aligned} \quad (2.4)$$

is a normalization constant. □

Remark: The matrices sampled from the GOE are *not* orthogonal. The orthogonal in GOE comes from the fact that GOE matrices are invariant under orthogonal transformations, which means that \mathbf{M} and $\mathbf{Q}^T \mathbf{M} \mathbf{Q}$ has the same distribution, where \mathbf{Q} is an orthogonal matrix [1].

Matrices sampled from the GOE only contain real entries. Two other common Gaussian ensembles are the *Gaussian unitary ensemble* (GUE) and the *Gaussian symplectic ensemble* (GSE).

A random Hermitian $n \times n$ matrix \mathbf{M} is said to belong to the GUE if the diagonal entries $m_{jj} \sim \mathcal{N}(0, 1)$ and the upper triangular entries are given by

$m_{jk} = u_{jk} + iv_{jk}$ where $u_{jk}, v_{jk} \sim \mathcal{N}(0, 1/2)$. The pdf for GUE is given by

$$f(\mathbf{M}) = \prod_{j=1}^n \frac{1}{\sqrt{\pi}} e^{-m_{jj}^2} \prod_{1 \leq j < k \leq n} \frac{2}{\pi} e^{-2|m_{jk}^2|} = C_n e^{-\text{tr}(\mathbf{M})}, \quad (2.5)$$

where

$$C_n = \left(\frac{1}{\sqrt{\pi}} \right)^n \left(\frac{2}{\pi} \right)^n = \frac{2^n}{\pi^{3n/2}} \quad (2.6)$$

is a normalization constant. The U stands for GUE matrices being invariant under unitary transformations [8].

We will not discuss the GSE beside stating that its matrix entries are real quaternions.

2.2 Eigenvalues

We shall now look at the spectrum for random matrices with Gaussian entries. We will mostly limit ourselves to the GOE case but some more general results are presented as well.

Theorem 2.2. (Section 1.2 [8]) *The joint pdf for the eigenvalues of a $n \times n$ GOE matrix is given by*

$$f(\lambda_1, \dots, \lambda_n) = \frac{1}{C_n} e^{-\frac{1}{2} \sum_{i=1}^n \lambda_i^2} \prod_{1 \leq i < j \leq n} |\lambda_j - \lambda_i|,$$

where C_n is a normalization constant.

We shall only show that Theorem 2.2 holds for the simplest case, that is, we want to show that

$$f(\lambda_1, \lambda_2) = k e^{-\frac{1}{2}(\lambda_1^2 + \lambda_2^2)} |\lambda_1 - \lambda_2| \quad (2.7)$$

is the correct eigenvalue distribution for a 2×2 GOE matrix, where k is normalization constant.

Proof (Following from [7]): We will accomplish (2.7) by transforming the GOE pdf in terms of its eigenvalues. Recall that GOE matrices are invariant under orthogonal transformation, so consider $\mathbf{M} = \mathbf{Q}^T \mathbf{M}_\lambda \mathbf{Q}$, where $\mathbf{M} = \begin{pmatrix} m_{11} & m_{12} \\ m_{21} & m_{22} \end{pmatrix}$

is a GOE matrix, $\mathbf{Q} = \begin{pmatrix} \cos \theta & \sin \theta \\ -\sin \theta & \cos \theta \end{pmatrix}$ is an orthogonal matrix and $\mathbf{M}_\lambda = \begin{pmatrix} \lambda_1 & 0 \\ 0 & \lambda_2 \end{pmatrix}$. We have that

$$\begin{aligned} \begin{pmatrix} m_{11} & m_{12} \\ m_{21} & m_{22} \end{pmatrix} &= \begin{pmatrix} \cos \theta & \sin \theta \\ -\sin \theta & \cos \theta \end{pmatrix}^T \begin{pmatrix} \lambda_1 & 0 \\ 0 & \lambda_2 \end{pmatrix} \begin{pmatrix} \cos \theta & \sin \theta \\ -\sin \theta & \cos \theta \end{pmatrix} \\ &= \begin{pmatrix} \lambda_1 \cos^2 \theta + \lambda_2 \sin^2 \theta & (\lambda_1 - \lambda_2) \sin \theta \sin \theta \\ (\lambda_1 - \lambda_2) \sin \theta \sin \theta & \lambda_1 \sin^2 \theta + \lambda_2 \cos^2 \theta \end{pmatrix} \end{aligned} \quad (2.8)$$

that is

$$\begin{cases} m_{11} = \lambda_1 \cos^2 \theta + \lambda_2 \sin^2 \theta, \\ m_{22} = \lambda_1 \sin^2 \theta + \lambda_2 \cos^2 \theta, \\ m_{12} = \frac{1}{2}(\lambda_1 - \lambda_2) \sin 2\theta. \end{cases} \quad (2.9)$$

Note that $m_{12} = m_{21}$ due to symmetry, so we only need to consider one of the entries. To ensure that our function in the terms of its new variables (the eigenvalues) is a valid pdf, we need to take the Jacobian determinant into account. The Jacobian is given by

$$\mathcal{J} = \begin{pmatrix} \cos^2 \theta & \sin^2 \theta & (\lambda_2 - \lambda_1) \sin 2\theta \\ \sin^2 \theta & \cos^2 \theta & (\lambda_1 - \lambda_2) \sin 2\theta \\ \frac{1}{2} \sin 2\theta & -\frac{1}{2} \sin 2\theta & (\lambda_1 - \lambda_2) \cos 2\theta \end{pmatrix} \quad (2.10)$$

and

$$\begin{aligned} \det(\mathcal{J}) &= \det \begin{pmatrix} \cos^2 \theta & \sin^2 \theta & (\lambda_2 - \lambda_1) \sin 2\theta \\ \sin^2 \theta & \cos^2 \theta & (\lambda_1 - \lambda_2) \sin 2\theta \\ \frac{1}{2} \sin 2\theta & -\frac{1}{2} \sin 2\theta & (\lambda_1 - \lambda_2) \cos 2\theta \end{pmatrix} \\ &= \lambda_1 \cos^4 \theta \cos 2\theta - \lambda_2 \cos^4 \theta \cos 2\theta - \lambda_2 \sin^2 2\theta \\ &\quad + \lambda_1 \sin^2 2\theta + x \sin^2(2\theta) - \lambda_1 \sin^4 \theta \cos 2\theta \\ &\quad + \lambda_2 \sin^4 \theta \cos 2\theta \\ &= (\lambda_1 - \lambda_2)(\cos^2 2\theta + \sin^2 2\theta) \\ &= \lambda_1 - \lambda_2 \end{aligned} \quad (2.11)$$

Now, the change of variable can be done as $f(m_{11}, m_{22}, m_{12}) \rightarrow f(\lambda_1, \lambda_2, \theta) |\det(\mathcal{J})| =: \hat{f}(\lambda_1, \lambda_2, \theta)$, where $f(m_{11}, m_{22}, m_{12})$ is the GOE pdf. We have that

$$\hat{f}(\lambda_1, \lambda_2, \theta) = k' e^{-\frac{1}{2}(\lambda_1^2 + \lambda_2^2)} |\lambda_1 - \lambda_2|, \quad (2.12)$$

where k' is a constant. To get rid of θ , all we need to do is integrate with respect to θ which only contributes to the constant value. We have

$$\begin{aligned}\hat{f}(\lambda_1, \lambda_2) &= \int_{\theta} k' e^{-\frac{1}{2}(\lambda_1^2 + \lambda_2^2)} |\lambda_1 - \lambda_2| d\theta \\ &= k e^{-\frac{1}{2}(\lambda_1^2 + \lambda_2^2)} |\lambda_1 - \lambda_2|.\end{aligned}\tag{2.13}$$

□

Next thing we are going to look at is the distribution of the distance between two eigenvalues from a 2×2 GOE matrix.

Theorem 2.3. *The pdf for the distance $d = |\lambda_1 - \lambda_2|$ between the two eigenvalues λ_1, λ_2 in a 2×2 GOE matrix is given by*

$$f(d) = \frac{d}{2} e^{-\frac{d^2}{4}}.$$

Proof: Let $\mathbf{M} = \begin{pmatrix} m_{11} & m_{12} \\ m_{21} & m_{22} \end{pmatrix}$ be GOE matrix with eigenvalues λ_1, λ_2 such that $\lambda_1 > \lambda_2$ and consider $d = \lambda_1 - \lambda_2$. Solving $\det \begin{pmatrix} m_{11} - \lambda & m_{12} \\ m_{21} & m_{22} - \lambda \end{pmatrix} = 0$ gives us

$$\begin{aligned}\lambda_1 &= \frac{m_{11} + m_{22} + \sqrt{(m_{11} - m_{22})^2 + 4m_{12}m_{21}}}{2}, \\ \lambda_2 &= \frac{m_{11} + m_{22} - \sqrt{(m_{11} - m_{22})^2 + 4m_{12}m_{21}}}{2}\end{aligned}$$

and therefore

$$d = \sqrt{(m_{11} - m_{22})^2 + 4m_{12}^2},$$

where we used that $m_{12} = m_{21}$ due to symmetry. Now let $X = m_{11} - m_{22}$, $Y = 2m_{12}$. We have that $X \sim \mathcal{N}(0, 2)$, $Y \sim \mathcal{N}(0, 2)$ so

$$d = \sqrt{X^2 + Y^2} \sim \text{Rayleigh}(\sqrt{2})\tag{2.14}$$

according to Lemma 1.1. The pdf for d is therefore given by

$$f(d) = \frac{d}{2} e^{-\frac{d^2}{4}}.\tag{2.15}$$

□

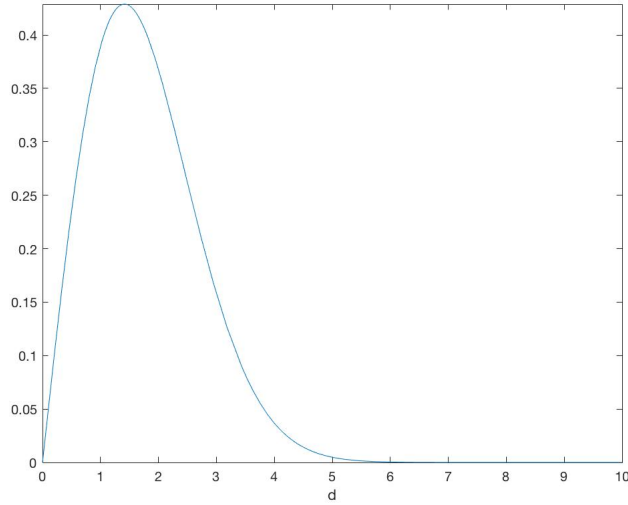


Figure 2.1: Pdf for distance between the two eigenvalues for a 2×2 GOE matrix.

What Theorem 2.3 actually says is that the eigenvalues do not want to be too close to each other, neither do they want to be too far apart. For d close to 0, the linear factor will dominate while for large d , the exponential factor will take over, see Figure 2.1.

One may rescale (2.15) to

$$f(s) = \frac{\pi s}{2} e^{-\frac{\pi s^2}{4}}, \quad (2.16)$$

which is known as *Wigner's surmise* [18].

Wigner suggested (2.16) to calculate the energy level spacings in the atomic nucleus where $s = \frac{\lambda_{i+1} - \lambda_i}{D}$ and D equals the mean distance for the energy levels. While (2.15) do gives us an correct pdf for the spacings between the eigenvalues, Wigner's surmise only works as an approximation for true value of the energy level spacings. The error may be up to 2% [9].

As we mentioned earlier there exist more Gaussian ensembles than the orthogonal one. Each of these ensembles have their own eigenvalue distribution as

well.

Theorem 2.4. (Section 1.3 [8]) *The joint pdf for the eigenvalues of the GOE ($\beta = 1$), the GUE ($\beta = 2$) and the GSE ($\beta = 4$) is given by*

$$f(\lambda_1, \dots, \lambda_n) = \frac{1}{G_{n,\beta}} e^{-\frac{\beta}{2} \sum_{i=1}^n \lambda_i^2} \prod_{1 \leq i < j \leq n} |\lambda_j - \lambda_i|^\beta,$$

where β is known as the Dyson index and

$$G_{n,\beta} = \beta^{-n/2-n\beta(n-1)/4} (2\pi)^{n/2} \prod_{j=0}^{n-1} \frac{\Gamma(1 + (j+1)\beta/2)}{\Gamma(1 + \beta/2)}$$

is a normalization constant.

So how do we actually calculate the probability that an eigenvalue ends up within a given interval \mathcal{I} ? Normally, this would be done as

$$P(\lambda \in \mathcal{I}) = \int_{\mathcal{I}} f(\lambda) d\lambda. \quad (2.17)$$

But how do we get the pdf for the eigenvalue? A good guess could be to integrate away the other variables in the joint pdf given in Theorem 2.4, that is

$$f(\lambda) = \int_{-\infty}^{\infty} \dots \int_{-\infty}^{\infty} f(\lambda, \lambda_2, \dots, \lambda_n) d\lambda_2 \dots d\lambda_n. \quad (2.18)$$

However, computing this multiple integral is far from trivial, even if we limit ourselves to the GOE. As we saw earlier in Figure 1.1, it seems like the eigenvalue distribution takes on a pileptic shape.

Theorem 2.5. (Section 1.4 [8]) *As the size n of a matrix \mathbf{M} from the Gaussian ensemble approaches infinity, the eigenvalue distribution for the scaled matrix $\frac{\mathbf{M}}{\sqrt{2n}}$ has a limit*

$$f(\lambda) = \frac{2}{\pi} \sqrt{1 - \lambda^2}, \quad \lambda \in [-1, 1].$$

Theorem 2.5 is known as *Wigner's semicircle law*. Wigner was first to prove it with the help of combinatorics in the 1950's. Later on, M. L. Mehta and M. Gaudin [18] proved it by computing the integral given in (2.18) and then taking the limit $n \rightarrow \infty$.

The curve generated for the semicircle law won't be a half-circle, but rather a half-ellipse, see Figure 2.2. Depending on how you choose to scale your matrix you will get different intervals for your elliptic curve.

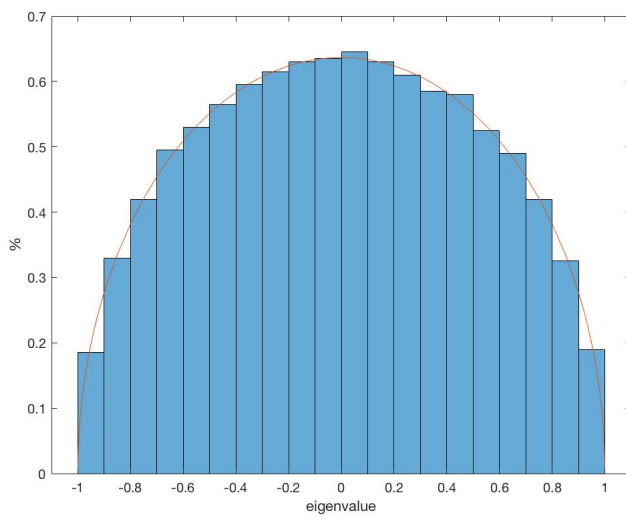


Figure 2.2: Eigenvalue distribution of 1000×1000 GOE matrix together with semicircle law.

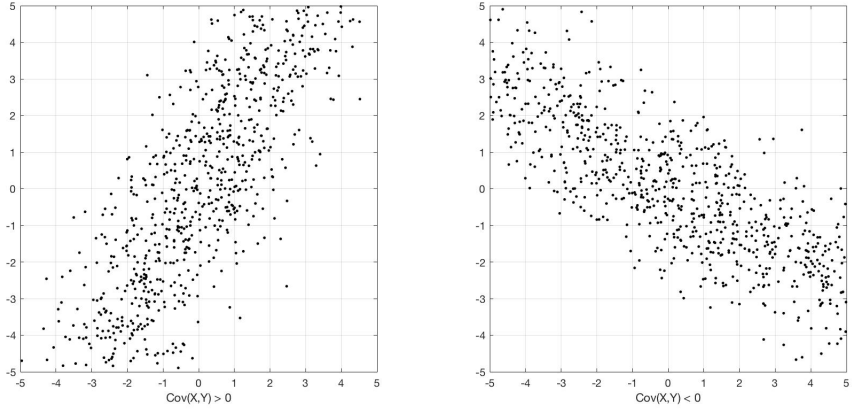


Figure 2.3: Plots of two data sets with positive and negative covariance respectively.

2.3 Covariance and Wishart Ensembles

So far have we only discussed theory regarding Wigner matrices and their natural symmetric counterpart $(\mathbf{A} + \mathbf{A}^T)/2$, which mostly appear when modelling, for example, physical systems [23]. However, when we are studying applications of multivariate statistics, such as variance analysis, regression analysis, dimension reduction etc. we are often interested in variation and dependency in our input data. A more occurring matrix in these topics is the *covariance matrix*, which will be the focus of this section.

Definition 2.3. The *covariance* between two random variables X and Y , where $E[X] = \mu_X, E[Y] = \mu_Y$, is given by

$$\text{Cov}[X, Y] = E[(X - \mu_X)(Y - \mu_Y)].$$

If X and Y are independent then $\text{Cov}(X, Y) = 0$. However, $\text{Cov}(X, Y) = 0$ does **not** necessary mean that X and Y are independent. The covariance acts as a metric of how two random variables vary compared to each other. Positive/negative covariance means that X and Y tend to be large/small simultaneously. This is illustrated in Figure 2.3.

Definition 2.4. The *covariance matrix* of a random vector $X \in \mathbb{R}^{p \times 1}$ is

given by

$$\mathbf{C} = E[(\mathbf{X} - \boldsymbol{\mu})(\mathbf{X} - \boldsymbol{\mu})^T],$$

where $\boldsymbol{\mu} = E[\mathbf{X}] = (E[X_1], E[X_2], \dots, E[X_p])^T$ is the mean vector of \mathbf{X} .

Each element c_{ij} in the covariance matrix \mathbf{C} denotes the covariance between X_i and X_j . Since $c_{ii} = \text{Cov}(X_i, X_i) = E[(X_i - \mu_{X_i})(X_i - \mu_{X_i})] = \sigma_i^2$ we have that the variance of X_i can be found on the diagonal of \mathbf{C} , as

$$\mathbf{C} = \begin{pmatrix} \sigma_1^2 & c_{12} & \dots \\ c_{21} & \sigma_2^2 & \dots \\ \vdots & \vdots & \ddots \end{pmatrix}. \quad (2.19)$$

Note that $\text{Cov}(X, Y) = \text{Cov}(Y, X)$ so \mathbf{C} is always symmetric. Also, for $\mathbf{v} \in \mathbb{R}^{n \times 1}$ we have that

$$\begin{aligned} \mathbf{v}^T \mathbf{C} \mathbf{v} &= \mathbf{v}^T E[(\mathbf{X} - \boldsymbol{\mu})(\mathbf{X} - \boldsymbol{\mu})^T] \mathbf{v} \\ &= E[\mathbf{v}^T (\mathbf{X} - \boldsymbol{\mu})(\mathbf{X} - \boldsymbol{\mu})^T \mathbf{v}] \\ &= E[(\mathbf{X} - \boldsymbol{\mu})^T \mathbf{v}]^T (\mathbf{X} - \boldsymbol{\mu})^T \mathbf{v}] \\ &= E[\|(\mathbf{X} - \boldsymbol{\mu})^T \mathbf{v}\|^2] \geq 0 \end{aligned} \quad (2.20)$$

so the covariance is positive semi-definite and therefore, due to Lemma 1.3, the eigenvalues of \mathbf{C} are non-negative. What (2.20) actually tells us is that the variance (see Definition 1.4) for a linear combination $\mathbf{a}^T \mathbf{X}$ of random variables is given by $\text{Var}[\mathbf{a}^T \mathbf{X}] = \mathbf{a}^T \mathbf{C} \mathbf{a}$.

With the covariance matrix defined we may expand the single variable Gaussian distribution into its multivariate case.

Definition 2.5. A random vector $\mathbf{X} \in \mathbb{R}^{p \times 1}$ has a **multivariate normal distribution** if its pdf is given by

$$f(\mathbf{X} | \boldsymbol{\mu}, \mathbf{C}) = \frac{1}{(2\pi)^{p/2} \det(\mathbf{C})} e^{-\frac{1}{2}(\mathbf{X} - \boldsymbol{\mu})^T \mathbf{C}^{-1} (\mathbf{X} - \boldsymbol{\mu})},$$

where $\boldsymbol{\mu}$ is the mean vector and \mathbf{C} is the covariance matrix. This will be denoted as $\mathbf{X} \sim \mathcal{N}(\boldsymbol{\mu}, \mathbf{C})$.

In most real-world scenarios, we don't know the true values of our parameters so we have to estimate them. For example, in the univariate case we, we estimate the mean μ with the sample mean $\bar{x} = \frac{1}{n} \sum x_i$ and variance is mostly estimated with the sample variance $s^2 = \frac{1}{n-1} \sum (x_i - \bar{x})^2$. What we often desire from our estimators is that they are *unbiased*, which means that their expected value is the true value.

Definition 2.6. Suppose we have n independent samples (x_i, y_i) . The **sample covariance** is given by

$$\hat{c} = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})$$

For the sample covariance, we have that

$$\begin{aligned} \hat{c} &= \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}) \\ &= \frac{1}{n-1} \left(\sum_{i=1}^n x_i y_i - \bar{y} \sum_{i=1}^n x_i - \bar{x} \sum_{i=1}^n y_i + n\bar{x}\bar{y} \right) \\ &= \frac{1}{n-1} \left(\sum_{i=1}^n x_i y_i - \frac{1}{n} \sum_{i=1}^n x_i \sum_{i=1}^n y_i \right) \\ &= \frac{1}{n-1} \left(\sum_{i=1}^n x_i y_i - \frac{1}{n} \left(\sum_{i=1}^n x_i y_i + \sum_{i \neq j} x_i y_j \right) \right) \end{aligned} \tag{2.21}$$

and calculating the corresponding expected value of \hat{C} gives us

$$\begin{aligned} E[\hat{C}] &= \frac{1}{n-1} (E[\sum_{i=1}^n X_i Y_i] - \frac{1}{n} (E[\sum_{i=1}^n X_i Y_i] + E[\sum_{i \neq j} X_i Y_i])) \\ &= \frac{1}{n-1} (nE[XY] - \frac{1}{n} (nE[XY] + n(n-1)E[X]E[Y])) \\ &= E[XY] - E[X]E[Y] \\ &= Cov(X, Y) \end{aligned} \tag{2.22}$$

and thus, the sample covariance is an unbiased estimator of the true covariance. The last step can be derived by expanding the covariance formula in Definition 2.3.

Definition 2.7. Suppose we have n independent samples $\mathbf{x}_i \in \mathbb{R}^{p \times 1}$ $i = 1, \dots, n$, then the **sample (empirical) covariance matrix** is given by

$$\mathbf{S} = \frac{1}{n-1} \sum_{i=1}^n (\mathbf{x}_i - \bar{\mathbf{x}})(\mathbf{x}_i - \bar{\mathbf{x}})^T.$$

By evaluating the sum for the sample covariance matrix \mathbf{S} as we did in (2.21), we get the more practical matrix notation

$$\mathbf{S} = \frac{1}{n-1} (\mathbf{X}\mathbf{X}^T - n\bar{\mathbf{x}}\bar{\mathbf{x}}^T), \tag{2.23}$$

where $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_n]$ is a $p \times n$ data matrix with each column being a sample and $\bar{\mathbf{x}}$ is the sample mean.

We shall now consider the distribution for matrices on the form $\mathbf{X}\mathbf{X}^T$ where the entries \mathbf{x}_i of \mathbf{X} are normal distributed. In general, when $\mathbf{x}_i \sim \mathcal{N}(\boldsymbol{\mu}, \mathbf{C})$, we have the *non-central Wishart distribution* [13]. For the sample covariance matrix \mathbf{S} , we indirectly assume $\boldsymbol{\mu} = \mathbf{0}$ as a result of subtracting the sample mean from each \mathbf{x}_i . Assuming zero mean explicitly is rather common in RMT and statistical literature, since most interesting results will be the same as with $\boldsymbol{\mu} \neq \mathbf{0}$ [20].

Theorem 2.6. (Section 7.2 [2]) *Let $\mathbf{W} = n\mathbf{S}$ where \mathbf{S} is the sample covariance matrix built upon $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_n] \in \mathbb{R}^{p \times n}$, $\mathbf{x}_i \sim \mathcal{N}(\boldsymbol{\mu}, \mathbf{C})$. \mathbf{W} is then called a **Wishart matrix** and has a (central) **Wishart Distribution**, denoted $W_p(n, \mathbf{C})$, with the pdf*

$$f(\mathbf{W}) = \frac{\det(\mathbf{W})^{\frac{n-p-1}{2}} e^{-\frac{1}{2}\text{tr}(\mathbf{C}^{-1}\mathbf{W})}}{2^{\frac{np}{2}} \pi^{\frac{p(p-1)}{4}} \det(\mathbf{C})^{\frac{n}{2}} \prod_{i=1}^p \Gamma\left(\frac{1}{2}(n+1-i)\right)}$$

Remark: To be consistent with the statistical literature we are referring to, we define the sample covariance matrix with a factor n^{-1} instead of $(n-1)^{-1}$. This change will be of minor importance for the theory we are representing and the empirical results will be sufficiently equal as n grows large.

Theorem 2.7. [6] *Let \mathbf{X} be a $p \times n$ matrix ($n \geq p$) with i.i.d entries x_{ij} being standard normal distributed. Then the matrix $\mathbf{W} = \mathbf{X}\mathbf{X}^T$ belongs to the **Wishart Ensemble** with the pdf*

$$f_{\beta}(\mathbf{W}) = \frac{1}{C_{p,n,\beta}} \det(\mathbf{W})^{\beta(p-n+1)/2-1} e^{-\frac{1}{2}\text{Tr}(\mathbf{W})},$$

where $C_{p,n,\beta}$ is a normalization constant and $\beta = 1, 2, 4$ corresponds to the orthogonal (real), unitary (complex) and symplectic (quaternion) ensembles, respectively.

Remark: We have only defined covariance for real matrices in this chapter but one may easily get the complex case by changing the real transpose operator $(\cdot)^T$ to the complex one $(\cdot)^*$. Note that for $\beta = 1$ we have the Wishart distribution $W_p(n, \mathbf{I})$.

Like the GOE case, matrices from the Wishart orthogonal ensemble are not orthogonal, they are invariant under an orthogonal transformation. Also, since the

entries of a Wishart matrix are not independent in general, it is not as trivial to prove Theorem 2.7 than, for example, Theorem 2.1. A proof by A. Edelman for the real and complex case, with help of matrix factorization, can be found in [5].

Just as eigenvalues tells us how eigenvectors are stretched when multiplied with their corresponding matrix, eigenvalues of covariance matrices do inform us how the data is spread (more about this in section 3.2). Like the Gaussian ensembles, we have a joint eigenvalue distribution for each of the Wishart ensembles.

Theorem 2.8. [6] *The joint eigenvalue pdf for a matrix generated by the Wishart ensemble is given by*

$$f(\lambda_1, \dots, \lambda_p) = c_{\beta,p} \prod_{i < j} |\lambda_i - \lambda_j|^\beta \prod_{i=1}^n \lambda_i^{\frac{\beta}{n}(p-n+1)-1} e^{-\frac{1}{2} \sum_{i=1}^n \lambda_i},$$

where $c_{\beta,p}$ is a normalization constant.

Since we only defined the Wishart ensemble for standard Gaussian entries, we do expect our covariance matrix to take the form of an identity matrix. One may mistakenly conclude assume that the eigenvalues will concentrate around 1. However, it turns out that $\lambda_1, \dots, \lambda_p \approx 1$ only holds when the sample size is much larger than the dimension of our data.

Theorem 2.9. [20] *Consider a $p \times n$ matrix \mathbf{X} ($n \geq p$) with i.i.d. entries x_{ij} being standard normal distributed and let $\gamma = \frac{p}{n} \in (0, 1]$. As the sample size n approaches infinity, the eigenvalue pdf for $\mathbf{S} = \frac{1}{n} \mathbf{X} \mathbf{X}^T$ has a limit*

$$f(\lambda) = \frac{1}{2\pi\gamma\lambda} \sqrt{(b_+ - \lambda)(\lambda - b_-)}, \quad \lambda \in [b_-, b_+], \quad (2.24)$$

where $b_{\pm} = (1 \pm \sqrt{\gamma})^2$.

Theorem 2.9 is known as the *Marčenko–Pastur law*. We observe in Figure 2.4 that our prediction about the eigenvalues being close to 1 gets less accurate when the dimension p approaches the sample size n .

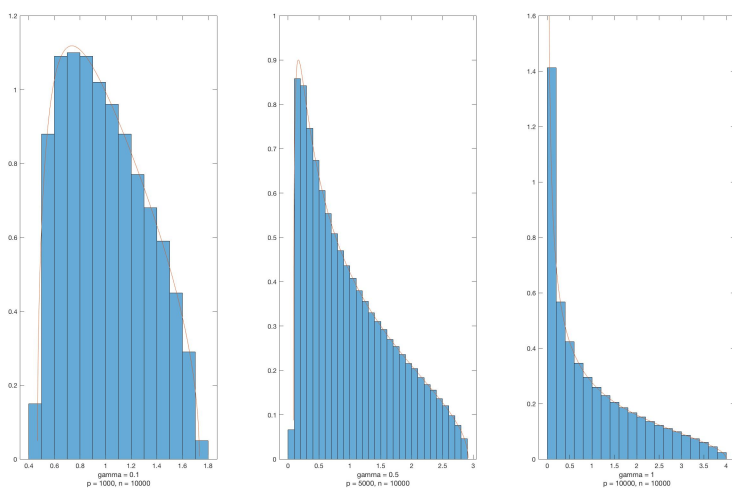


Figure 2.4: Empirical eigenvalue distribution of covariance matrices together with Marčenko–Pastur law for different $\gamma = \frac{p}{n}$ values.

Chapter 3

Applications

3.1 Random Matrices in Neural Networks

3.1.1 Neural Networks in General

Neural networks (NNs) are trained to predict/classify new data. The most basic approach is to train NNs with training samples $\{\mathbf{x}_i, y_i\}$ $i = 1, \dots, n$ where $\mathbf{x}_i \in \mathbb{R}^p$ is a feature vector and y_i is the correct class label. The output is a function $f(\mathbf{x}; \mathbf{w})$ which is a prediction of the true value y , where \mathbf{w} is a weight vector that determines how important each of the components in \mathbf{x} is.

The key idea is to find a function f and to adjust the parameters w such that the NN predicts class labels y with a high accuracy. We measure the accuracy of our NN with a *loss function* L , which we wish to minimize. The most elementary NNs uses $f(\mathbf{x}; \mathbf{w}) = \mathbf{w}^T \mathbf{x}$, which results in the loss function

$$L = \sum_{i=1}^n (y_i - \mathbf{w}^T \mathbf{x}_i)^2. \quad (3.1)$$

This is a least-square problem with an optimal solution $\mathbf{w} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X} \mathbf{y}$ where $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_n]$ is our data matrix and \mathbf{y} contains all our class labels y_i . An illustration of a very simple NN is given in Figure 3.1.

This NN structure can only solve linearly-separable classification problems which puts heavy limitations on the practical uses of it and a classic example of its limitations is the XOR problem [4]. A solution to this, which is based on Cover's

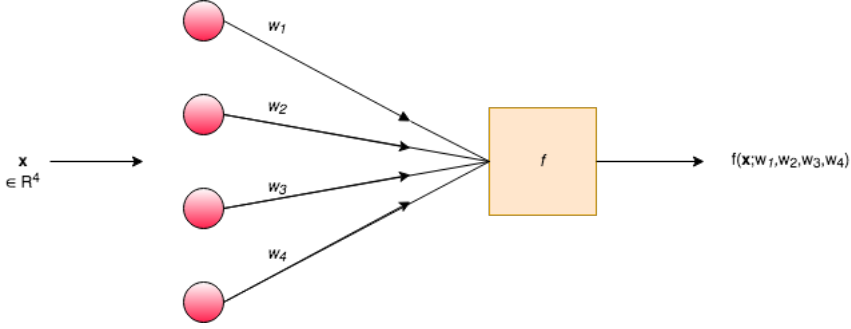


Figure 3.1: A very basic neural network with 4 dimensional feature vectors.

theorem [10], is to apply entry-wise a non-linear activation function $\sigma : \mathbb{R} \mapsto \mathbb{R}$ on our data. Our new loss function is then given by

$$L = \sum_{i=1}^n (y_i - \sigma(\mathbf{w}^T \mathbf{x}_i))^2. \quad (3.2)$$

Due to the non-linear mapping, our new optimization problem is non-linear and our least-square solution will not work. In general, a NN with a non-linear activation function will result in a non-convex optimization problem which is considered hard to solve. NNs are in general built of hidden-layer(s) which are used to apply our activation function and the training time for these networks can be up to days [14]. This is due to the massive amounts of parameters in large NNs, which in general results in many local minima while seeking an optimal solution. See Figure 3.2 for an illustration of this¹.

A way to tackle this problem is to consider networks with a single hidden-layer where the input weights are randomly selected, so called *Extreme Learning Machines* (ELMs) [12]. ELMs has shown to perform really good for large non-complex data sets compared to other popular classification algorithms (such as support vector machines)[11]. The fast performance comes from the fact that we do not train the first layer in ELMs, so our non-linear optimization problem turns into a regression problem with a closed form solution.

Due to the random weights we may use results from RMT to model our ELMs. A framework for this will be represented in the next subsection, solely based on

¹Figure is taken from <https://github.com/tomgoldstein/loss-landscape>

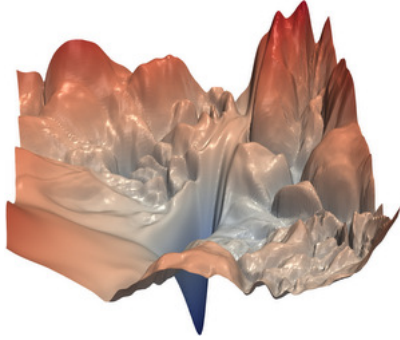


Figure 3.2: Illustration of a loss surface with many local minima .

the work by Z. Liao, C. Louart and R. Couillet [24] [19].

3.1.2 RMT Model for Neural Networks

Consider a data matrix $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_n] \in \mathbb{R}^{p \times n}$ with a target matrix² $\mathbf{Y} = [\mathbf{y}_1, \dots, \mathbf{y}_n] \in \mathbb{R}^{d \times n}$, a random weight matrix $\mathbf{H} \in \mathbb{R}^{N \times p}$ with i.i.d. Gaussian standard entries and a weight matrix $\boldsymbol{\beta} \in \mathbb{R}^{N \times d}$ for the second layer. We shall denote the output from the hidden layer as $\boldsymbol{\Sigma} = \sigma(\mathbf{H}\mathbf{X}) \in \mathbb{R}^{N \times n}$, where $\sigma : \mathbb{R} \mapsto \mathbb{R}$ is our activation function applied entry-wise, see figure Figure 3.3. As mentioned earlier, only the output layer $\boldsymbol{\beta}$ is to be trained, while \mathbf{H} is static but randomly chosen. We seek weights $\boldsymbol{\beta}$ which minimizes the loss function

$$L(\boldsymbol{\beta}) = \frac{1}{n} \sum_{i=1}^n \|\mathbf{y}_i - \boldsymbol{\beta}^T \sigma(\mathbf{H}\mathbf{x}_i)\|^2 + \lambda \|\boldsymbol{\beta}\|_F^2,$$

for some regularization parameter $\lambda \geq 0$ to ensure our model is not overfitting. An explicit solution for this given by

$$\boldsymbol{\beta} = \frac{1}{n} \left(\frac{1}{n} \boldsymbol{\Sigma} \boldsymbol{\Sigma}^T + \lambda \mathbf{I}_N \right)^{-1} \boldsymbol{\Sigma} \mathbf{Y}^T = \frac{1}{n} \boldsymbol{\Sigma} \left(\frac{1}{n} \boldsymbol{\Sigma}^T \boldsymbol{\Sigma} + \lambda \mathbf{I}_n \right)^{-1} \mathbf{Y}^T$$

which, per sample, results in a mean square training error

$$E_{train} = MSE_{train}(\lambda) = \frac{1}{n} \|\boldsymbol{\beta}^T \boldsymbol{\Sigma} - \mathbf{Y}\|_F^2 = \frac{\lambda^2}{n} \text{tr}(\mathbf{Y} \mathbf{Q}^2 \mathbf{Y}^T), \quad (3.3)$$

²This is a more general notation of the one we used in (3.1), which allows us to have multiple output nodes. However, when we represent the results in the end of this chapter, we will use unidimensional output columns ($d = 1$) as in (3.1).

1) $\frac{p}{n} \rightarrow c_1$ and $\frac{N}{n} \rightarrow c_2$ where $c_1, c_2 \in (0, \infty)$, 2) Input matrix \mathbf{X} has a bounded spectral norm, and 3) Output matrix \mathbf{Y} has bounded entries.

Determining a deterministic equivalent to \mathbf{Q} often boils down to calculating $E[\mathbf{Q}]$ which unfolds from the fact that $\|E[\mathbf{Q}] - \bar{\mathbf{Q}}\| \rightarrow 0$. Given the model we represented, Z. Liao shows that

$$\bar{\mathbf{Q}} = \bar{\mathbf{Q}}(\lambda) = (\check{\mathbf{K}} - \lambda \mathbf{I}_n)^{-1}, \quad \check{\mathbf{K}} = \frac{N}{n(1+\delta)} \mathbf{K},$$

where $\mathbf{K} = E_{\mathbf{w}}[\sigma(\mathbf{w}^T \mathbf{X})^T \sigma(\mathbf{w}^T \mathbf{X})]$ is an *equivalent kernel matrix* for $\mathbf{w} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_p)$ and δ implicitly defined as the solution to $\delta = \frac{1}{n} \text{tr}(\mathbf{K} \bar{\mathbf{Q}})$ do result in

$$\|E[\mathbf{Q}] - \bar{\mathbf{Q}}\| \rightarrow 0,$$

as $n, p, N \rightarrow \infty$. This, together with (3.3) and (3.4), gives us the deterministic training and test mean-squared errors

$$\begin{aligned} \bar{E}_{train} &= \frac{\lambda^2}{n} \text{tr} \left(\mathbf{Y} \bar{\mathbf{Q}} \left(\frac{\frac{1}{N} \text{tr}(\bar{\mathbf{Q}} \check{\mathbf{K}} \bar{\mathbf{Q}})}{1 - \frac{1}{N} \text{tr}(\check{\mathbf{K}} \bar{\mathbf{Q}} \check{\mathbf{K}} \bar{\mathbf{Q}})} \check{\mathbf{K}} + \mathbf{I}_n \right) \bar{\mathbf{Q}} \mathbf{Y}^T \right), \\ \bar{E}_{test} &= \frac{1}{\hat{n}} = \left\| \hat{\mathbf{Y}} - \hat{\mathbf{Y}} \bar{\mathbf{Q}} \check{\mathbf{K}}_{\mathbf{X} \hat{\mathbf{X}}} \right\|_F^2 \\ &\quad + \frac{\frac{1}{N} \text{tr}(\mathbf{Y} \bar{\mathbf{Q}} \check{\mathbf{K}} \bar{\mathbf{Q}} \mathbf{Y}^T)}{1 - \frac{1}{N} \text{tr}(\check{\mathbf{K}} \bar{\mathbf{Q}} \check{\mathbf{K}} \bar{\mathbf{Q}})} \left(\frac{1}{\hat{n}} \text{tr}(\check{\mathbf{K}}_{\hat{\mathbf{X}} \hat{\mathbf{X}}}) - \frac{1}{\hat{n}} \text{tr} \left(\mathbf{I}_n + \frac{N}{n} \lambda \bar{\mathbf{Q}} \right) \check{\mathbf{K}}_{\mathbf{X} \hat{\mathbf{X}}} \check{\mathbf{K}}_{\hat{\mathbf{X}} \mathbf{X}} \bar{\mathbf{Q}} \right), \end{aligned}$$

which results in the desired convergence

$$E_{train} - \bar{E}_{train} \rightarrow 0, \quad E_{test} - \bar{E}_{test} \rightarrow 0,$$

understood almost surely, where we used the notations

$$\mathbf{K}_{AB} = E_{\mathbf{w}}[\sigma(\mathbf{w}^T \mathbf{A})^T \sigma(\mathbf{w}^T \mathbf{B})], \quad \check{\mathbf{K}}_{AB} = \frac{N}{n} \frac{\mathbf{K}_{AB}}{1+\delta}, \quad \mathbf{K} = \mathbf{K}_{\mathbf{X} \mathbf{X}}, \quad \check{\mathbf{K}} = \check{\mathbf{K}}_{\mathbf{X} \mathbf{X}}.$$

We represent the result³ on a classification task of a 2-class MNIST⁴ dataset (7 and 9) for two different choices of activation function σ in Figure 3.4. We can observe an almost perfect match between theory and simulation for relatively small n, \hat{n}, p and N .

³The results are represented with a modified version of Z. Liao's code. Z. Liao's original code can be found on <https://github.com/Zhenyu-LIAO/RMT4ELM>

⁴<http://yann.lecun.com/exdb/mnist/>

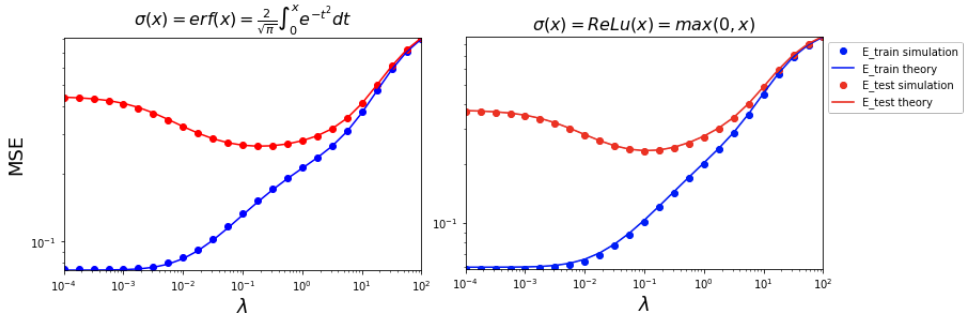


Figure 3.4: Comparison between predicted and simulated MSE (log-log scale) for sample size $n = \hat{n} = 1024$, variables $p = 784$ and hidden neurons $N = 512$.

3.2 Dimensionality Reduction - PCA

In modern technology it is very common that data sets are very large and problems may arise when we try to draw statistical conclusions [3, 22]. A common strategy to tackle this is to transform the data such that our transformed data contains fewer variables but most of the relevant information remains. A very popular method to accomplish this is the *Principal Component Analysis* (PCA), which may traced back to 1901 [21].

3.2.1 Minimizing Mean Squared Error of Projection

Consider a sampled data matrix $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_n] \in \mathbb{R}^{p \times n}$ that we want to project onto a lower dimensional subspace spanned by the orthogonal unit vectors $[\mathbf{v}_1, \dots, \mathbf{v}_n] = \mathbf{V} \in \mathbb{R}^{d \times n}$, where $d < p$. It is safe to assume that the sample mean \bar{x}_j for each column in \mathbf{X} equals to zero (on a real data set, this can be achieved by replacing our variables x_{ij} with $\hat{x}_{ij} = x_{ij} - \bar{x}_j$). This projection will result in loss of information and our goal is to minimize it. Due to Theorem 1.3, we know that the optimal projection (in terms of least squares) is given by the orthogonal projection. This gives us the mean squared error (MSE)

$$\begin{aligned} \frac{1}{n} \sum_{i=1}^n \|\mathbf{x}_i - \mathbf{V}\mathbf{V}^T \mathbf{x}_i\|_2^2 &= \frac{1}{n} \sum_{i=1}^n \left(\|\mathbf{x}_i\|_2^2 - \|\mathbf{V}^T \mathbf{x}_i\|_2^2 \right) \\ &= \frac{1}{n} \sum_{i=1}^n \|\mathbf{x}_i\|_2^2 - \frac{1}{n} \sum_{i=1}^n \|\mathbf{V}^T \mathbf{x}_i\|_2^2. \end{aligned} \quad (3.5)$$

As we are minimizing the MSE in terms of \mathbf{V} , we can simply ignore the first term and focusing on maximizing the second term. With help of the trace identities $\mathbf{a}^T \mathbf{a} = \text{tr}(\mathbf{a}\mathbf{a}^T)$ and $\text{tr}(\mathbf{ABC}) = \text{tr}(\mathbf{BCA})$ we get

$$\begin{aligned} \frac{1}{n} \sum_{i=1}^n \|\mathbf{V}^T \mathbf{x}_i\|_2^2 &= \frac{1}{n} \sum_{i=1}^n \text{tr}(\mathbf{V}^T \mathbf{x}_i \mathbf{x}_i^T \mathbf{V}) = \frac{1}{n} \text{tr}(\mathbf{V}^T \mathbf{X} \mathbf{X}^T \mathbf{V}) \\ &= \text{tr}(\mathbf{V}^T \mathbf{S} \mathbf{V}) = \text{tr}(\mathbf{S} \mathbf{V} \mathbf{V}^T) \\ &= \sum_{i=1}^n \text{tr}(\mathbf{S} \mathbf{v}_i \mathbf{v}_i^T) = \sum_{i=1}^n (\mathbf{S} \mathbf{v}_i)^T \mathbf{v}_i \\ &= \sum_{i=1}^n \mathbf{v}_i^T \mathbf{S} \mathbf{v}_i. \end{aligned} \quad (3.6)$$

Now let $y_i = \mathbf{v}^T \mathbf{x}_i$, we have that

$$\begin{aligned}
 \mathbf{v}^T \mathbf{S} \mathbf{v} &= \frac{1}{n-1} \left(\sum_{i=1}^n \mathbf{v}^T \mathbf{x}_i \mathbf{x}_i^T \mathbf{v} - n \mathbf{v}^T \bar{\mathbf{x}} \bar{\mathbf{x}}^T \mathbf{v} \right) \\
 &= \frac{1}{n-1} \left(\sum_{i=1}^n (\mathbf{v}^T \mathbf{x}_i)^2 - n (\mathbf{v}^T \bar{\mathbf{x}})^2 \right) \\
 &= \frac{1}{n-1} \left(\sum_{i=1}^n (\mathbf{v}^T \mathbf{x}_i)^2 - \frac{1}{n} \left(\sum_{i=1}^n \mathbf{v}^T \mathbf{x}_i \right)^2 \right) \\
 &= \frac{1}{n-1} \sum_{i=1}^n (y_i - \bar{y})^2 = s_y^2
 \end{aligned} \tag{3.7}$$

which tells us that $\mathbf{v}^T \mathbf{S} \mathbf{v}$ is the sample variance of y_1, \dots, y_n . This means that minimizing the MSE of the orthogonal projection is equivalent to maximizing the variance of our data. We shall exploit this in the next subsection to acquire a rather simple model for reducing the dimensionality while preserving as much information as possible.

3.2.2 Maximizing the Variance

We have that the key idea behind PCA is to maximize the variance of our data which means that we seek a linear combination $\mathbf{v}^T \mathbf{x} = \sum_{i=1}^p v_i x_i$ so that the variance is maximized. The sample variance for such a linear combination is given by $\mathbf{v}^T \mathbf{S} \mathbf{v}$ where \mathbf{S} is the sample covariance matrix. To proceed with a non-trivial optimization problem, we shall assume that the vectors spanning the subspace \mathbf{V} are of unit norm and orthogonal to each other, that is, $\|\mathbf{v}\|_2^2 = \mathbf{v}_i^T \mathbf{v}_j = 1$ for $i = j$ and $\mathbf{v}_i^T \mathbf{v}_j = 0$ otherwise. This results in the following optimization problem

$$\begin{aligned}
 \max_{\mathbf{v} \in \mathbb{R}^{p \times 1}} \quad & \mathbf{v}^T \mathbf{S} \mathbf{v} \\
 \text{s.t.} \quad & \mathbf{v}^T \mathbf{v} = 1
 \end{aligned} \tag{3.8}$$

Note that (3.8) is not a convex optimization problem. However, we can still solve it by looking at the spectral decomposition of \mathbf{S}

$$\mathbf{S} = \mathbf{E} \mathbf{D} \mathbf{E}^T = \begin{pmatrix} | & & | \\ \mathbf{e}_1 & \cdots & \mathbf{e}_p \\ | & & | \end{pmatrix} \begin{pmatrix} \lambda_1 & & 0 \\ & \ddots & \\ 0 & & \lambda_p \end{pmatrix} \begin{pmatrix} - & \mathbf{e}_1 & - \\ & \vdots & \\ - & \mathbf{e}_p & - \end{pmatrix},$$

where, without loss of generality, we can assume that the eigenvalues are ordered as $\lambda_1 \geq \dots \geq \lambda_p$. Due to the symmetry of \mathbf{S} , we also know that the eigenvectors $\mathbf{e}_1, \mathbf{e}_2, \dots$ will be orthogonal to each other (see Lemma 1.4). Multiplying \mathbf{S} from both sides with \mathbf{v} gives us

$$\mathbf{v}^T \mathbf{S} \mathbf{v} = \mathbf{v}^T \mathbf{E} \mathbf{D} \mathbf{E}^T \mathbf{v} = (\mathbf{E}^T \mathbf{v})^T \mathbf{D} \mathbf{E}^T \mathbf{v} = \mathbf{z}^T \mathbf{D} \mathbf{z} = \sum_{i=1}^p \lambda_i z_i^2 \leq \lambda_1 \sum_{i=1}^p z_i^2, \quad (3.9)$$

where we used the notation $\mathbf{z} = \mathbf{E}^T \mathbf{v}$. Note that

$$\sum_{i=1}^p z_i^2 = \|\mathbf{z}\|_2^2 = (\mathbf{E}^T \mathbf{v})^T \mathbf{E}^T \mathbf{v} = \mathbf{v}^T \mathbf{E} \mathbf{E}^T \mathbf{v} = \mathbf{v}^T \mathbf{v} = 1, \quad (3.10)$$

which tells us that $\mathbf{v}^T \mathbf{S} \mathbf{v} \leq \lambda_1$. Choosing \mathbf{v} to be the eigenvector \mathbf{e}_1 corresponding to λ_1 results in

$$\mathbf{v}^T \mathbf{S} \mathbf{v} = \mathbf{e}_1^T \mathbf{E} \mathbf{D} \mathbf{E}^T \mathbf{e}_1 = (\mathbf{E}^T \mathbf{e}_1)^T \mathbf{D} \mathbf{E}^T \mathbf{e}_1 = \lambda_1, \quad (3.11)$$

since

$$\mathbf{E}^T \mathbf{e}_1 = \begin{pmatrix} - & \mathbf{e}_1 & - \\ & \vdots & \\ - & \mathbf{e}_p & - \end{pmatrix} \begin{pmatrix} | \\ \mathbf{e}_1 \\ | \end{pmatrix} = \begin{pmatrix} \mathbf{e}_1^T \mathbf{e}_1 \\ \mathbf{e}_2^T \mathbf{e}_1 \\ \vdots \end{pmatrix} = \begin{pmatrix} 1 \\ 0 \\ \vdots \end{pmatrix}. \quad (3.12)$$

This can be done for all the eigenvalues of \mathbf{S} , that is,

$$\lambda_1 = \mathbf{e}_1^T \mathbf{S} \mathbf{e}_1 = \text{sample variance for linear combinations } \mathbf{e}_1^T \mathbf{x}_i \ \forall i = 1, \dots, n,$$

$$\vdots$$

$$\lambda_p = \mathbf{e}_p^T \mathbf{S} \mathbf{e}_p = \text{sample variance for linear combinations } \mathbf{e}_p^T \mathbf{x}_i \ \forall i = 1, \dots, n,$$

which tells us that the largest eigenvalue and its corresponding eigenvector matches the direction in which the data varies most. This means that reducing the dimensionality of a dataset while preserving maximum variance is obtained by projecting the data onto the eigenvector(s) corresponding to the largest eigenvalue(s).

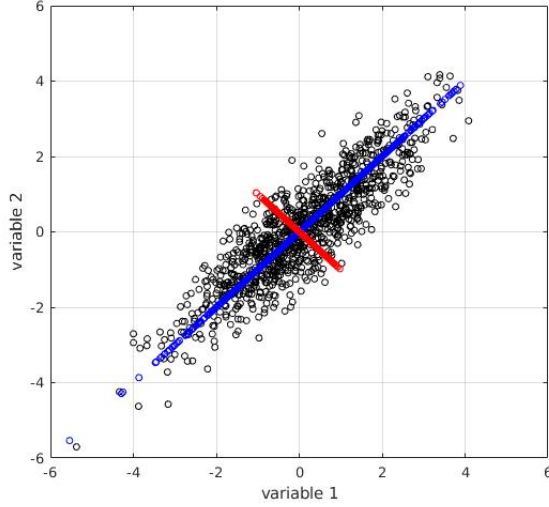


Figure 3.5: Blue/red dots represent the 2-dimensional data projected onto the first/second eigenvector.

We can measure how much of the variance is accounted for in a direction by calculating

$$\frac{\lambda_k}{\text{tr}(\mathbf{S})}, \quad (3.13)$$

where λ_k is the eigenvalue corresponding to the eigenvector \mathbf{e}_k we project our data onto. As an illustrating example, consider the 2-dimensional dataset in Figure 3.5. Projecting the data onto the first eigenvector (blue dots) will take about 94% of the the variance into account.

In general, one direction isn't enough to describe larger datasets sufficiently good. We can take more eigenvalues into consideration in our calculation and evaluate the accounted variances as

$$\Lambda_{\mathcal{S}} = \frac{\sum_{i \in \mathcal{S}} \lambda_i}{\text{tr}(\mathbf{S})}, \quad (3.14)$$

where \mathcal{S} is the set of eigenvalues we want to consider in our analysis. Nevertheless, correlation often do exist between sampled variables and then it's usually

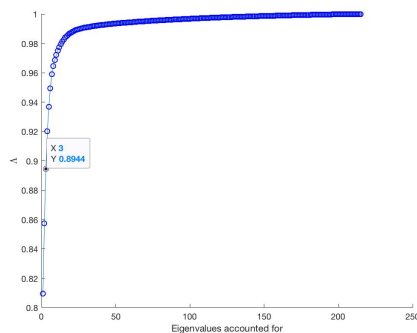


Figure 3.6: Cumulative sum of the eigenvalues corresponding to the sample covariance matrix of the data set 'ovariancancer'.

enough with a small set of eigenvalues to get a satisfying Λ_S . To illustrate this, we perform PCA on the dataset 'ovariancancer'⁵ which contains $n = 216$ samples⁶ and $p = 4000$ variables that has been accounted for. Observing the result in Figure 3.6, we can see that we only need 3 eigenvalues (out of 216 unique ones) to explain about 90% of the variance. Another interesting outcome is that most of the eigenvalues tend to be very small and lump together while only representing a small portion of the variance. Plotting the data projected onto the first three eigenvectors gives us a noticeable clustering of the data (see Figure 3.7), which illustrates how PCA can be used to represent high dimensional data.

⁵Can be loaded in MATLAB with the command "load ovariancancer"

⁶121 with ovarian cancer and 95 without.

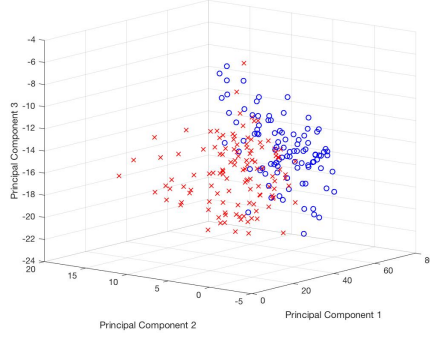


Figure 3.7: Red/Blue markers representing cancer/healthy patients in 3 dimension.

3.2.3 Statistical inference of PCA

So far, we treated PCA as a linear algebra problem without any assumptions about the underlying distribution of the data, and for many applications, this is good enough. For example, PCA may work as first step for other machine learning algorithms so intractable problems can be solved in a satisfactory time, or as we demonstrated in Figure 3.7, to visualize high-dimensional data. However, if we want to perform some kind of statistical inference, we need to consider some distribution for our data samples in \mathbf{X} .

Consider $\mathbf{x}_i \sim \mathcal{N}(\boldsymbol{\mu}, \mathbf{C})$ $i = 1, \dots, n$ and let λ_j $j = 1, \dots, p$ be unique distinct eigenvalues for \mathbf{C} and let $\hat{\lambda}_j$ be the eigenvalues for a sample covariance matrix built upon \mathbf{x}_i . D. N. Lawley [17] proves that the mean and variance for the random variable $\hat{\lambda}_j$ is given by

$$\begin{aligned} E[\hat{\lambda}_j] &= \lambda_j + \frac{\lambda_j}{n} \sum_{i=1, i \neq j}^p \left(\frac{\lambda_i}{\lambda_j - \lambda_i} \right) + \mathcal{O}\left(\frac{1}{n^2}\right), \\ \text{Var}[\hat{\lambda}_j] &= \frac{2\lambda_j^2}{n} \left(1 - \frac{1}{n} \sum_{i=1, i \neq j}^p \left(\frac{\lambda_i}{\lambda_j - \lambda_i} \right)^2 \right) + \mathcal{O}\left(\frac{1}{n^3}\right). \end{aligned} \tag{3.15}$$

T. W. Anderson [2] shows that the eigenvalues $\hat{\lambda}_j$ in fact are asymptotically

normal distributed. This means that as n grows large, we have

$$\hat{\lambda}_j \sim \mathcal{N}\left(\lambda_j, \frac{2\lambda_j^2}{n}\right) \quad (3.16)$$

approximately. Normalizing (3.16) results in

$$\frac{\hat{\lambda}_j - \lambda_j}{\lambda_j \sqrt{2/n}} \sim \mathcal{N}(0, 1), \quad (3.17)$$

which can be used for classic statistical inference such as hypothesis testing and interval estimation. For example, a confidence interval (CI) for λ_j with a confidence level $1 - \alpha$ may be constructed as

$$\begin{aligned} -z_{\alpha/2} &\leq \frac{\hat{\lambda}_j - \lambda_j}{\lambda_j \sqrt{2/n}} \leq z_{\alpha/2} \\ &\Leftrightarrow \\ -z_{\alpha/2} \lambda_j \sqrt{2/n} + \lambda_j &\leq \hat{\lambda}_j \leq z_{\alpha/2} \sqrt{2/n} + \lambda_j \\ &\Leftrightarrow \\ \frac{\hat{\lambda}_j}{1 + z_{\alpha/2} \sqrt{2/n}} &\leq \lambda_j \leq \frac{\hat{\lambda}_j}{1 - z_{\alpha/2} \sqrt{2/n}} \end{aligned} \quad (3.18)$$

and thus, we have a CI

$$I_{\lambda_j} = \left(\frac{\hat{\lambda}_j}{1 + z_{\alpha/2} \sqrt{2/n}}, \frac{\hat{\lambda}_j}{1 - z_{\alpha/2} \sqrt{2/n}} \right), \quad (3.19)$$

where $z_{\alpha/2}$ is the z -score which can be found in a normal table. Recall that we are only interested in I_{λ_j} for large n , so it is safe to assume that $z_{\alpha/2} \sqrt{2/n} < 1$.

For hypothesis testing, we may want to test

$$H_0 : \lambda_j = \lambda_{j0} \quad \text{vs} \quad H_1 : \lambda_j \neq \lambda_{j0} \quad (3.20)$$

where we reject H_0 on a significance level α if

$$\left| \frac{\hat{\lambda}_j - \lambda_{j0}}{\lambda_j \sqrt{2/n}} \right| \geq z_{\alpha/2}. \quad (3.21)$$

A perhaps more interesting test would be to check if the smallest $(p - k)$ eigenvalues are all the same, that is

$$H_0 : \lambda_{k+1} = \lambda_{k+2} = \dots = \lambda_p$$

vs

$$(3.22)$$

H_1 : At least two of the smallest $(p - k)$ eigenvalues are not equal.

Accepting H_0 as in (3.22) justifies that the k largest PCs are measuring some significant variation in our input data while the remaining $(p - k)$ PCs effectively measure noise. To test (3.22), one may use

$$Q = \left(\frac{\prod_{j=k+1}^p \hat{\lambda}_j}{\left(\sum_{j=k+1}^p \frac{\hat{\lambda}_j}{p-k} \right)^{p-k}} \right)^{n/2} \quad (3.23)$$

as a test statistic, where $-2 \ln(Q) \sim \chi^2(f)$ approximately under H_0 and $f = \frac{1}{2}(p - k + 2)(p - k - 1)$ [15]. This lets us reject the null hypothesis at a significance level α if

$$-2 \ln(Q) \geq \chi_{\alpha/2}^2(f), \quad (3.24)$$

where $\chi_{\alpha/2}^2(f)$ can be found in a χ^2 table.

Chapter 4

Ending discussion

We have introduced, discussed and applied the results from random matrix theory. The focus has solely been on matrices with entries being normal distributed, or even stricter, $\mathcal{N}(0, 1)$ distributed. Nevertheless, many interesting results, such as *Wigner's semicircle law* and the *Marčenko–Pastur law*, still emerges.

As we mentioned in the beginning, RMT mostly appears on research level. The aim of this thesis has been to introduce and perhaps simplify some results of this rather advanced topic, so that someone with basic probability and linear algebra skills can get a solid grasp of the properties related to random matrices.

For further research, it would be interesting to generalize some of the results. We have mostly focused on the spectrum of large matrices with an underlying $\mathcal{N}(0, 1)$ distribution. Having explicit non-asymptotically eigenvalue distributions for matrices with other distributions than standard normal would perhaps shed some light on the topics of RMT, such that it would be introduced much earlier in a mathematical oriented education.

Finally, as for applying the theory, we have only introduced RMT in already well established applications. Starting from the other side, that is, given an arbitrary matrix ensemble and trying to answer 'what does this model?' would perhaps result in new applications beside those we already have today.

Bibliography

- [1] G. W. Anderson, A. Guionnet, and O. Zeitouni. *An Introduction to Random Matrices*. Cambridge Studies in Advanced Mathematics. Cambridge University Press, 2009.
- [2] T. W. Anderson. *An Introduction to Multivariate Statistical Analysis*. John Wiley & Sons, 3rd edition, 2003.
- [3] S. Ayesha, M. K. Hanif, and R. Talib. Overview and comparative study of dimensionality reduction techniques for high dimensional data. *Information Fusion*, 59:44–58, 2020.
- [4] A. Cyr, F. Thériault, and S. Chartier. Revisiting the xor problem: a neuro-robotic implementation. *Neural Computing and Applications*, 32(14):9965–9973, 2020.
- [5] A. Edelman. *Eigenvalues and Condition Numbers of Random Matrices*. PhD thesis, Massachusetts Institute of Technology, 1989.
- [6] A. Edelman and N. Raj Rao. Random matrix theory. *Acta Numerica*, 14:233–297, 2005.
- [7] G. Ergün. *Random Matrix Theory*, pages 2549–2563. Springer New York, New York, NY, 2012.
- [8] P. J. Forrester. *Log-gases and random matrices*. Princeton University Press, Princeton, 2010.
- [9] P. J. Forrester and N. S. Witte. Exact Wigner surmise type evaluation of the spacing distribution in the bulk of the scaled random matrix ensembles. *arXiv e-prints*, pages math-ph/0009023, September 2000.
- [10] M. Hoogendoorn and B. Funk. *Machine Learning for the Quantified Self: On the Art of Learning from Sensory Data*. Cognitive Systems Monographs. Springer International Publishing, 2017.

- [11] G. Huang, H. Zhou, X. Ding, and R. Zhang. Extreme learning machine for regression and multiclass classification. *IEEE Transactions on Systems, Man, and Cybernetics, Part B (Cybernetics)*, 42(2):513–529, 2012.
- [12] G. Huang, Q-Y. Zhu, and C-K. Siew. Extreme learning machine: Theory and applications. *Neurocomputing*, 70(1):489 – 501, 2006. Neural Networks.
- [13] A. T. James and R. A. Fisher. The non-central wishart distribution. *Proceedings of the Royal Society of London. Series A. Mathematical and Physical Sciences*, 229(1178):364–366, 1955.
- [14] A. H. Jiang, D. L. K. Wong, G. Zhou, D. G. Andersen, J. Dean, G. R. Ganger, G. Joshi, M. Kaminsky, M. Kozuch, Z. C. Lipton, and P. Pillai. Accelerating deep learning by focusing on the biggest losers. *arXiv e-prints*, page arXiv:1910.00762, October 2019.
- [15] I. T. Jolliffe. *Principal Component Analysis*. Springer Series in Statistics. Springer, 2002.
- [16] I. T. Jolliffe and J. Cadima. Principal component analysis: a review and recent developments. *Philosophical Transactions of the Royal Society A: Mathematical, Physical and Engineering Sciences*, 374(2065):20150202, April 2016.
- [17] D. N. Lawley. Tests of significance for the latent roots of covariance and correlation matrices. *Biometrika*, 43:128, June 1956.
- [18] G. Livan, M. Novaes, and P. Vivo. Introduction to random matrices. *SpringerBriefs in Mathematical Physics*, 2018.
- [19] C. Louart, Z. Liao, and R. Couillet. A Random Matrix Approach to Neural Networks. *arXiv e-prints*, page arXiv:1702.05419, February 2017.
- [20] D. Paul. and A. Aue. Random matrix theory in statistics: A review. *Journal of Statistical Planning and Inference*, 150:1–29, 2014.
- [21] K. Pearson. On lines and planes of closest fit to systems of points in space. *The London, Edinburgh, and Dublin Philosophical Magazine and Journal of Science*, 2(11):559–572, 1901.
- [22] S. Velliangiri, S. Alagumuthukrishnan, and S Iwin Thankumar joseph. A review of dimensionality reduction techniques for efficient computation. *Procedia Computer Science*, 165:104–111, 2019.
- [23] E. P. Wigner. Random matrices in physics. *SIAM Review*, 9(1):1–23, 1967.

-
- [24] L. Zhenyu. *A random matrix framework for large dimensional machine learning and neural networks*. Theses, Université Paris-Saclay, September 2019.

Copyright

The publishers will keep this document online on the Internet – or its possible replacement – from the date of publication barring exceptional circumstances. The online availability of the document implies permanent permission for anyone to read, to download, or to print out single copies for his/her own use and to use it unchanged for non-commercial research and educational purpose. Subsequent transfers of copyright cannot revoke this permission. All other uses of the document are conditional upon the consent of the copyright owner. The publisher has taken technical and administrative measures to assure authenticity, security and accessibility.

According to intellectual property law the author has the right to be mentioned when his/her work is accessed as described above and to be protected against infringement.

For additional information about the Linköping University Electronic Press and its procedures for publication and for assurance of document integrity, please refer to its www home page: <http://www.ep.liu.se/>.

Upphovsrätt

Detta dokument hålls tillgängligt på Internet – eller dess framtida ersättare – från publiceringsdatum under förutsättning att inga extraordinära omständigheter uppstår.

Tillgång till dokumentet innebär tillstånd för var och en att läsa, ladda ner, skriva ut enstaka kopior för enskilt bruk och att använda det oförändrat för ickekommersiell forskning och för undervisning. Överföring av upphovsrätten vid en senare tidpunkt kan inte upphäva detta tillstånd. All annan användning av dokumentet kräver upphovsmannens medgivande. För att garantera äktheten, säkerheten och tillgängligheten finns lösningar av teknisk och administrativ art. Upphovsmannens ideella rätt innefattar rätt att bli nämnd som upphovsman i den omfattning som god sed kräver vid användning av dokumentet på ovan beskrivna sätt samt skydd mot att dokumentet ändras eller presenteras i sådan form eller i sådant sammanhang som är kränkande för upphovsmannens litterära eller konstnärliga anseende eller egenart.

För ytterligare information om Linköping University Electronic Press se förlagets hemsida <http://www.ep.liu.se/>.