

Random Matrix Theory: Selected Applications from Statistical Signal Processing and Machine Learning

Ph.D. Thesis Defense

Khalil Elkhail

Committee Chairperson: Dr. Tareq Y. Al-Naffouri

Committee Co-Chair: Dr. Mohamed-Slim Alouini

Committee Members: Dr. Xiangliang Zhang and Dr. Abla Kammoun

External Examiner: Dr. Alfred Hero

June 24, 2019

King Abdullah University of Science and Technology
Computer, Electrical and Mathematical Sciences and Engineering

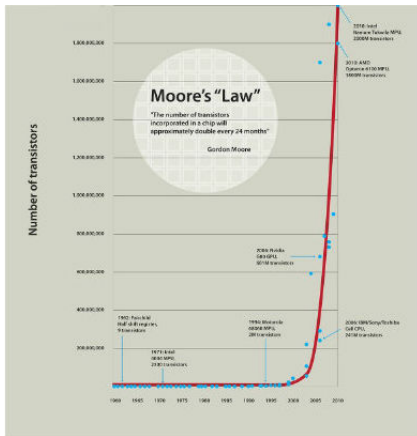


1. Introduction
2. Moments of Correlated Gram matrices
3. Regularized discriminant analysis with large dimensional data
4. Centered Kernel Ridge Regression (CKRR)
5. Conclusion
6. Future research directions

Introduction

Moore's law

- The # of transistors that you can fit into a piece of silicon doubles every couple of years.



1

¹C. M. Bishop, Microsoft research, Cambridge.

Big data

- ⦿ Large sample size
- ⦿ High dimensional data
- ⦿ High variability



New processing techniques

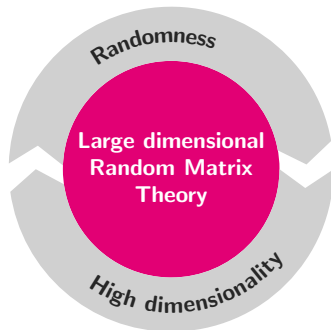
- ⦿ Optimization
- ⦿ Statistics

We need a tool that embraces these challenges

Random matrix theory (RMT)

Study the behavior of large random matrices

- Allow the prediction of the behavior of random quantities depending on large random matrices
- Key of success: **Randomness** + **High dimensionality**



Statistical Signal Processing

- Large number of antenna arrays vs large number of observations

→ Improved signal processing techniques

Wireless Communications

- Large # of antennas, Large # of users

→ Improved transmission and detection strategies

→ Low complexity design

Machine learning

- Supervised²/semi-supervised³/unsupervised learning⁴.

→ A better fundamental understanding

→ Improved classification performance

²Z. Liao, R. Couillet, "A Large Dimensional Analysis of Least Squares Support Vector Machines", submitted

³X. Mai, R. Couillet, "A random matrix analysis and improvement of semi-supervised learning for large dimensional data", submitted.

⁴R. Couillet, F. Benaych-Georges, "Kernel Spectral Clustering of Large Dimensional Data", Electronic Journal of Statistics, vol. 10, no. 1, pp. 1393-1454, 2016.

How does this work ?

- Self-averaging effect mechanism similar to that met in the law of large numbers
- $\mathbf{h}_1, \dots, \mathbf{h}_n \in \mathbb{C}^p$ with i.i.d entries with zero mean and variance $\frac{1}{n}$.
- $\mathbf{H}\mathbf{H}^H$ is an estimator of the cov. matrix with $\mathbf{H} = [\mathbf{h}_1, \dots, \mathbf{h}_n]$.

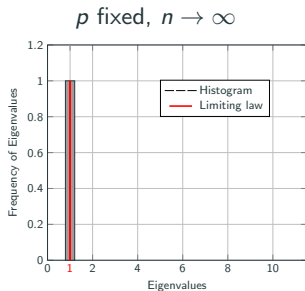


Figure 1.1: Histogram of eigenvalues of $\mathbf{H}\mathbf{H}^H$

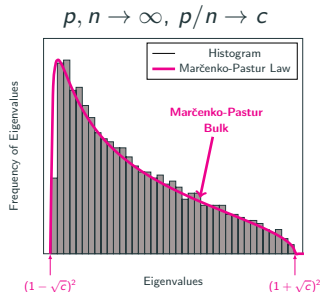


Figure 1.2: Histogram of eigenvalues of $\mathbf{H}\mathbf{H}^H$

Why is this useful?

The same result can be extended in the correlated case⁵

Certain functionals of $\mathbf{H}\mathbf{H}^H$ can be evaluated when $p, n \rightarrow \infty, p/n \rightarrow c$.

$$f(\mathbf{H}\mathbf{H}^H)$$

- $\frac{1}{n} \text{tr}(\mathbf{H}\mathbf{H}^H)$
- $\frac{1}{n} \text{tr}(\mathbf{H}\mathbf{H}^H)^{-k}$: performance of linear est. techniques
- $\frac{1}{n} \log \det(\mathbf{H}\mathbf{H}^H)$: MIMO systems, linear estimation (LCE)
- $\lambda_{\min}(\mathbf{H}\mathbf{H}^H), \lambda_{\max}(\mathbf{H}\mathbf{H}^H), \dots$: WEV in linear estimation

What happens for the moments in the finite regime?

$$\mathbb{E}_{\mathbf{H} \sim \mathcal{D}} f(\mathbf{H}\mathbf{H}^H)$$

⁵J. W. Silverstein and Z. D. Bai, On the Empirical Distribution of Eigenvalues of a Class of Large Dimensional Random Matrices, Journal of Multivariate Analysis, vol. 54, pp. 175192, May 2002.

Moments of Correlated Gram matrices

Linear estimation

Let $m < n$ and $\mathbf{H} \in \mathbb{C}^{n \times m}$ with i.i.d zero mean unit variance Gaussian entries and $\mathbf{\Lambda}$ is positive definite matrix with **distinct** eigenvalues $\theta_1, \theta_2, \dots, \theta_n$.

$$\mathbf{y}_{n \times 1} = \mathbf{H}_{n \times m} \mathbf{v}_{m \times 1} + \mathbf{\Lambda}^{-\frac{1}{2}} \mathbf{z}_{n \times 1}. \quad (1)$$

Define the **correlated Gram** matrix

$$\mathbf{G} = \mathbf{H}^* \mathbf{\Lambda} \mathbf{H}. \quad (2)$$

	LS	LMMSE
MSE	$\mathbb{E} \operatorname{tr} \mathbf{G}^{-1}$	$\mathbb{E} \operatorname{tr} (\mathbf{G} + \mathbf{R}_v^{-1})^{-1}$

Sample covariance matrix (SCM): $\mathbf{u}(k) = \mathbf{R}^{\frac{1}{2}} \mathbf{h}(k)$

$$\hat{\mathbf{R}}(n) = (1 - \lambda) \sum_{k=1}^n \lambda^{n-k} \mathbf{u}(k) \mathbf{u}^*(k) = \mathbf{R}^{\frac{1}{2}} \mathbf{H} \mathbf{\Lambda}(n) \mathbf{H}^* \mathbf{R}^{\frac{1}{2}}, \quad (3)$$

$$\begin{aligned} \text{Loss}(n) &\triangleq \mathbb{E} \left\| \mathbf{R}^{\frac{1}{2}} \hat{\mathbf{R}}^{-1}(n) \mathbf{R}^{\frac{1}{2}} - \mathbf{I}_m \right\|_F^2 \\ &= m + \mathbb{E} \operatorname{tr} \mathbf{G}_n^{-2} - 2 \mathbb{E} \operatorname{tr} \mathbf{G}_n^{-1} \end{aligned}$$

Negative moments of correlated Gram matrices

Define the **negative moments** of \mathbf{G} as

$$\mu_{\Lambda}(-k) \triangleq \mathbb{E} \operatorname{tr}(\mathbf{G}^{-k}), \quad k \in \mathbb{N}.$$

Then,

	LS (Exact)	LMMSE ($\mathbf{R}_x = \sigma_x^2 \mathbf{I}$, $\sigma_x^2 \gg 1$)
MSE	$\mu_{\Lambda}(-1)$	$\sum_{k=0}^l \frac{(-1)^k}{\sigma_x^{2k}} \mu_{\Lambda}(-k-1) + o(\sigma_x^{-2l})$

$$\text{Loss}(n) = m + \mu_{\Lambda(n)}(-2) - 2\mu_{\Lambda(n)}(-1).$$

Theorem (Negative moments)^a Let $p = \min(m, n - m)$, then for $1 \leq k \leq p$, we have

$$\mu_{\Lambda}(-k) = L \sum_{j=1}^k \sum_{i=1}^m \mathcal{D}(i, j) \frac{(-1)^{k-j}}{(k-j)!} \mathbf{b}_i^t \boldsymbol{\Psi}^{-1} \mathbf{D}_i \mathbf{a}_{j,k}.$$

^aK. Elkhailil, A. Kammoun, T. Al-Naffouri and M.-S. Alouini. Analytical Derivation of the Inverse Moments of One-sided Correlated Gram Matrices with Applications. IEEE Trans. Signal Processing, 2016.

$$\boldsymbol{\Psi} = \begin{bmatrix} 1 & \theta_1 & \cdots & \theta_1^{n-m-1} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & \theta_{n-m} & \cdots & \theta_{n-m}^{n-m-1} \end{bmatrix}$$

$$\mu_{\Lambda}(-k) = L \sum_{j=1}^k \sum_{i=1}^m \mathcal{D}(i, j) \frac{(-1)^{k-j}}{(k-j)!} \mathbf{b}_i^t \Psi^{-1} \mathbf{D}_i \mathbf{a}_{j,k}.$$

$$\Psi = \begin{bmatrix} 1 & \theta_1 & \cdots & \theta_1^{n-m-1} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & \theta_{n-m} & \cdots & \theta_{n-m}^{n-m-1} \end{bmatrix}$$

- So complicated formula
- Not useful if the eigenvalues of Λ are close to each other (We treat this issue for positive moments) ⁶.
- Not numerically stable if the dimensions are large.
- Not insightful!
- Not universal: the result will be different if we change the distribution from Gaussian.

⁶K. Elkhailil, A. Kammoun, T. Y. Al-Naffouri and M.-S. Alouini: Numerically Stable Evaluation of Moments of Random Gram Matrices With Applications. IEEE Signal Process. Lett. 24(9): 1353-1357 (2017)

Theorem (Silverstein and Bai ⁷)

Consider the Gram matrix $\mathbf{G} = \mathbf{H}^* \mathbf{\Lambda} \mathbf{H}$ with the following assumptions

- $m, n \rightarrow \infty$ with $\frac{m}{n} \rightarrow c \in (0, \infty)$
- $\|\mathbf{\Lambda}\| = O(1)$ with $\text{rank}(\mathbf{\Lambda}) = O(m)$.

$$\frac{1}{m} \text{tr} \mathbf{G}^{-1} - \delta \xrightarrow{a.s.} 0, \quad \delta = \left[\frac{1}{m} \text{tr} \mathbf{\Lambda} (\mathbf{I}_n + \delta \mathbf{\Lambda})^{-1} \right]^{-1}.$$

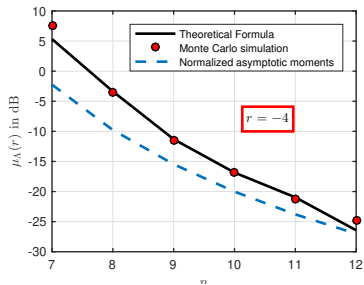
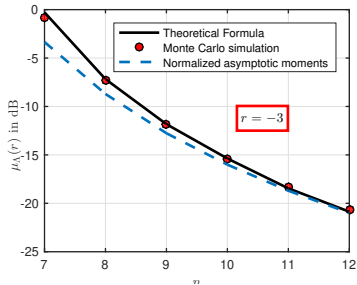
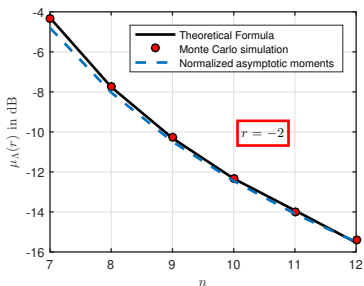
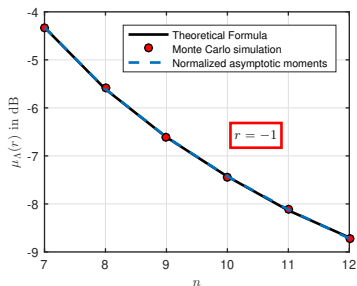
Higher inverse moments can be computed using an iterative process^a

^aKhalil Elkhail, Abla Kammoun, Tareq Y. Al-Naffouri, Mohamed-Slim Alouini: Analytical Derivation of the Inverse Moments of One-Sided Correlated Gram Matrices With Applications. IEEE Trans. Signal Processing 64(10): 2624-2635 (2016)

⁷J. W. Silverstein and Z. D. Bai, On the Empirical Distribution of Eigenvalues of a Class of Large Dimensional Random Matrices, Journal of Multivariate Analysis, vol. 54, pp. 175192, May 2002.

Validation of the inverse moments

$m = 3$



Optimal λ for SCM estimation

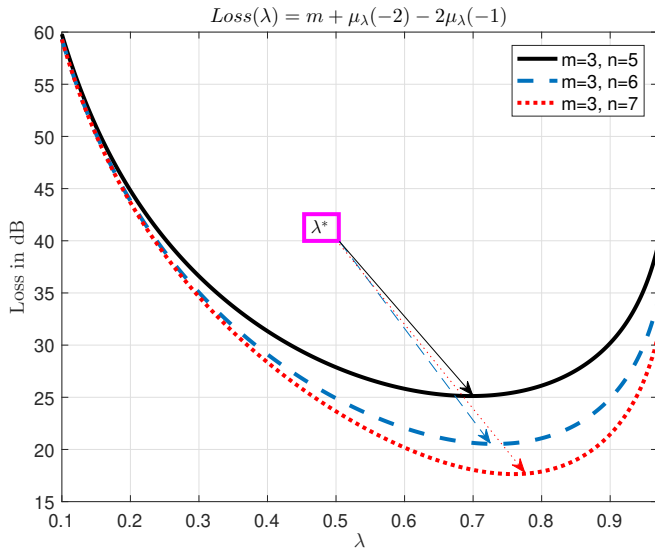
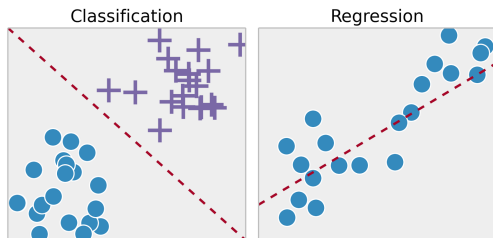


Figure 2.1: The estimation loss as a function of λ (Exact formula).

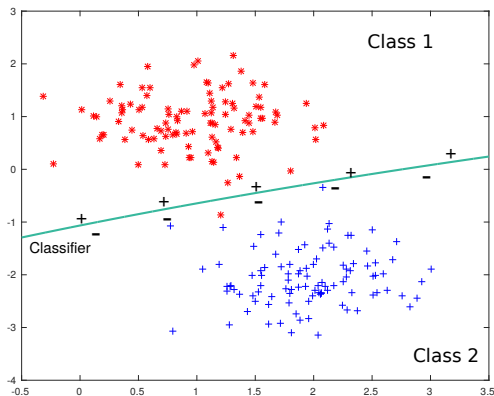
Regularized discriminant analysis with large dimensional data

- We are provided with labeled data $(\text{features}_i, \text{response/label}_i)_{1 \leq i \leq n}$.
- Fit a model to the data.



Classification

- Principle: Build a classification rule that allows to assign for an unseen observation its corresponding class.



Let \mathbf{x} be the input data and f be the classification rule.

$$\text{Classifier} \triangleq \begin{cases} \text{Assign class 1} & \text{if } f(\mathbf{x}) > 0 \\ \text{Assign class 2} & \text{if } f(\mathbf{x}) \leq 0 \end{cases}$$

- Data is assumed to be sampled from a certain dist.
- The decision rule is constructed based on that.
- The MAP rule is considered in the design

$$\hat{k} = \arg \max_{k: \text{classes}} \mathbb{P}[C_k | \mathbf{x}]$$

The classifier is designed to satisfy this rule.

Gaussian discriminant analysis

Gaussian mixture model for binary classification (2 classes)

- $\mathbf{x}_1, \dots, \mathbf{x}_n \in \mathbb{R}^p$
- Class k is formed by $\mathbf{x} \sim \mathcal{N}(\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$, $k = 0, 1$

Linear discriminant analysis (LDA): $\boldsymbol{\Sigma}_0 = \boldsymbol{\Sigma}_1 = \boldsymbol{\Sigma}$

$$W^{LDA}(\mathbf{x}) = \left(\mathbf{x} - \frac{\boldsymbol{\mu}_0 + \boldsymbol{\mu}_1}{2} \right)^T \boldsymbol{\Sigma}^{-1} (\boldsymbol{\mu}_0 - \boldsymbol{\mu}_1) - \log \frac{\pi_1}{\pi_0} > 0. \quad \begin{matrix} C_0 \\ < \\ C_1 \end{matrix}$$

→ Decision rule is linear in \mathbf{x} .

Quadratic discriminant analysis: $\boldsymbol{\Sigma}_0 \neq \boldsymbol{\Sigma}_1$

$$W^{QDA}(\mathbf{x}) = -\frac{1}{2} \log \frac{|\boldsymbol{\Sigma}_0|}{|\boldsymbol{\Sigma}_1|} - \frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_0)^T \boldsymbol{\Sigma}_0^{-1} (\mathbf{x} - \boldsymbol{\mu}_0) + \frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_1)^T \boldsymbol{\Sigma}_1^{-1} (\mathbf{x} - \boldsymbol{\mu}_1) > 0. \quad \begin{matrix} C_0 \\ < \\ C_1 \end{matrix}$$

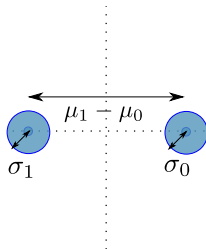
→ Decision rule is quadratic in \mathbf{x} .

How does this perform ?

- Assume Σ , μ_0 and μ_1 known.
- Equal priors : $\pi_0 = \pi_1 = 0.5$
- No asymptotic regime, p is fixed.

The total misclassification rate is given by ⁸

$$\epsilon = \Phi\left(-\frac{\Delta}{2}\right), \quad \Delta = \|\mu_0 - \mu_1\|_{\Sigma^{-1}}$$



What happens when the statistics are not known ?

⁸Jerome Friedman, Trevor Hastie, and Robert Tibshirani. The Elements of Statistical Learning. Springer, 2009.

LDA: Asymptotic regime (equal covariances)

Asymptotic growth regime

Let $n = n_0 + n_1$.

- $n_0, n_1, p \rightarrow \infty$ such that $\frac{p}{n} \rightarrow c < 1$.
- μ_0 and μ_1 are known.
- Σ is replaced by its sample estimate $\hat{\Sigma} = \frac{1}{n-2} \sum_{k=1}^2 \sum_{i=1}^{n_k} (\mathbf{x}_{k,i} - \bar{\mathbf{x}}_k)(\mathbf{x}_{k,i} - \bar{\mathbf{x}}_k)^T$.

Wang et al. 2018 ^a

$$\epsilon_{LDA} - \Phi \left[-\frac{\Delta}{2} \sqrt{1-c} \right] \rightarrow_{prob.} 0$$

^aCheng Wang and Binyan Jiang. On the dimension effect of regularized linear discriminant analysis, arXiv:1710.03136v1

- When $c \rightarrow 1$, the misclassification rate tends to 0.5.
- For the LDA to result in acceptable performance, we need c close to 0.
- Because its use of the inverse of the pooled covariance matrix, the LDA applies only when $c < 1$.

What happens if $p > n$?

Regularization

$$\mathbf{H} = (\mathbf{I}_p + \gamma \hat{\Sigma})^{-1}.$$

Optimal γ ?

Dimensionality reduction

$$\text{data}_{(d)} = \mathbf{W}_{d \times p} \times \text{data}_{(p)}$$

Best d ?

Random projections

$$\mathbb{R}^p \longrightarrow \mathbb{R}^d$$

$$\mathbf{x} \longmapsto \mathbf{W}\mathbf{x}$$

Projection matrix

We shall assume that the projection matrix \mathbf{W} writes as $\mathbf{W} = \frac{1}{\sqrt{p}}\mathbf{Z}$, where the entries $Z_{i,j}$ ($1 \leq i \leq d$, $1 \leq j \leq p$) of \mathbf{Z} are centered with unit variance and independent identically distributed random variables satisfying the following moment assumption. There exists $\epsilon > 0$, such that $\mathbb{E}|Z_{i,j}|^{4+\epsilon} < \infty$.

Johnson-Lindenstrauss Lemma

For a given n data points $\mathbf{x}_1, \dots, \mathbf{x}_n$ in \mathbb{R}^p , $\epsilon \in (0, 1)$ and $d > \frac{8 \log n}{\epsilon^2}$, there exists a linear map $f : \mathbb{R}^p \rightarrow \mathbb{R}^d$ such that

$$(1 - \epsilon) \|\mathbf{x}_i - \mathbf{x}_j\|^2 \leq \|\mathbf{W}\mathbf{x}_i - \mathbf{W}\mathbf{x}_j\|^2 \leq (1 + \epsilon) \|\mathbf{x}_i - \mathbf{x}_j\|^2, \quad (4)$$

Conditional risk after projection

$$\epsilon_i^{\text{P-LDA}} = \Phi \left[-\frac{1}{2} \sqrt{\boldsymbol{\mu}^\top \mathbf{W}^\top (\mathbf{W} \boldsymbol{\Sigma} \mathbf{W}^\top)^{-1} \mathbf{W} \boldsymbol{\mu}} + \frac{(-1)^{i+1} \log \frac{\pi_0}{\pi_1}}{\sqrt{\boldsymbol{\mu}^\top \mathbf{W}^\top (\mathbf{W} \boldsymbol{\Sigma} \mathbf{W}^\top)^{-1} \mathbf{W} \boldsymbol{\mu}}} \right].$$

Performance of LDA with random projections

Asymptotic Performance^a

$$\epsilon_i^{\text{P-LDA}} - \Phi \left[\frac{-\frac{1}{2} \boldsymbol{\mu}^\top (\boldsymbol{\Sigma} + \delta_d \mathbf{I}_p)^{-1} \boldsymbol{\mu} + (-1)^{i+1} \log \frac{\pi_0}{\pi_1}}{\sqrt{\boldsymbol{\mu}^\top (\boldsymbol{\Sigma} + \delta_d \mathbf{I}_p)^{-1} \boldsymbol{\mu}}} \right] \xrightarrow{\text{prob.}} 0, \quad (5)$$

$$\delta_d \text{tr}(\boldsymbol{\Sigma} + \delta_d \mathbf{I}_p)^{-1} = p - d. \quad (6)$$

δ_d can be seen as a penalty on projection.

^aK. Elkhalil, A. Kammoun, R. Calderbank, T. Al-Naffouri and M.-S. Alouini. Asymptotic Performance of Linear Discriminant Analysis with Random Projections. ICASSP 2019.

LDA

equal priors: $\Phi \left[-\frac{1}{2} \sqrt{\boldsymbol{\mu}^\top \boldsymbol{\Sigma} \boldsymbol{\mu}} \right]$

$\boldsymbol{\Sigma} = \mathbf{I}_p$: $\Phi \left[-\frac{1}{2} \|\boldsymbol{\mu}\| \right]$

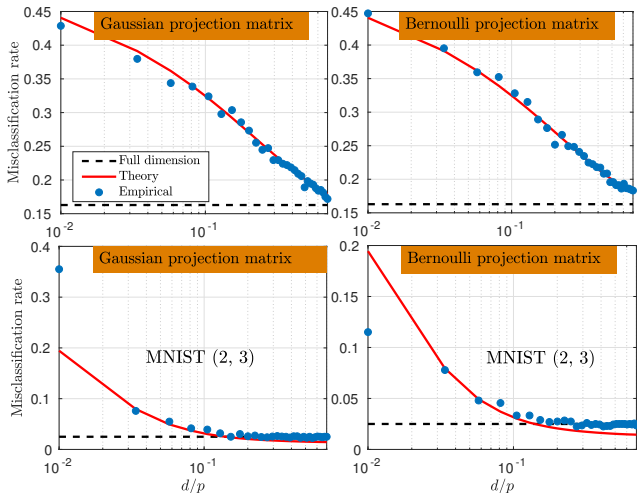
P-LDA

$$\Phi \left[-\frac{1}{2} \sqrt{\boldsymbol{\mu}^\top (\boldsymbol{\Sigma} + \delta_d \mathbf{I}_p)^{-1} \boldsymbol{\mu}} \right]$$

$$\Phi \left[-\frac{1}{2} \sqrt{d/p} \|\boldsymbol{\mu}\| \right]$$

P-LDA: Experiments

- $p = 800$.
- $\mu_0 = \mathbf{0}_p$ and $\mu_1 = \frac{3}{\sqrt{p}} \mathbf{1}_p$.
- $\Sigma = \{0.4^{|i-j|}\}_{i,j}$.



R-LDA: Asymptotic regime (equal covariances)

Asymptotic growth regime

- $n_0, n_1, p \rightarrow \infty$ such that $\frac{p}{n} \rightarrow c \in (0, \infty)$.
- $\boldsymbol{\mu}_k$ are replaced by $\bar{\mathbf{x}}_k = \frac{1}{n_k} \sum_{\mathbf{x}_i \in \mathcal{C}_k} \mathbf{x}_i$.
- $\boldsymbol{\Sigma}^{-1}$ is replaced by its ridge estimate $\mathbf{H} = (\mathbf{I}_p + \gamma \widehat{\boldsymbol{\Sigma}})^{-1}$.

Hachem et al 2008. ^a

$$\mathbf{H} \sim \mathbf{T} = (\mathbf{I}_p + \rho \boldsymbol{\Sigma})^{-1},$$

in the sense that $\mathbf{a}^T (\mathbf{H} - \mathbf{T}) \mathbf{b} \rightarrow_{prob.} 0$ and $\frac{1}{n} \text{tr} \mathbf{A} (\mathbf{H} - \mathbf{T}) \rightarrow_{prob.} 0$.

^aW. Hachem, O. Khorunzhiy, P. Loubaton, J. Najim, L. Pastur: A New Approach for Mutual Information Analysis of Large Dimensional Multi-Antenna Channels. IEEE Trans. Information Theory 54(9): 3987 - 4004 (2008).

Zollanvari and Dougherty 2015 ^a

$$\epsilon_{R-LDA}^{equal} - \Phi \left[\frac{-\boldsymbol{\mu}^T (\mathbf{I}_p + \rho \boldsymbol{\Sigma})^{-1} \boldsymbol{\mu}}{\sqrt{D}} \right] \rightarrow_{prob.} 0$$

^aAmin Zollanvari and Edward R. Dougherty: Generalized Consistent Error Estimator of Linear Discriminant Analysis. IEEE Trans. Signal Processing 63(11): 2804-2814 (2015)

R-LDA: Asymptotic regime (dist. covariances)

Asymptotic growth regime

- $n_0, n_1, p \rightarrow \infty$ such that $\frac{p}{n} \rightarrow c \in (0, \infty)$.
- $\boldsymbol{\mu}_k$ are replaced by $\bar{\mathbf{x}}_k = \frac{1}{n_k} \sum_{\mathbf{x}_i \in \mathcal{C}_k} \mathbf{x}_i$.
- $\boldsymbol{\Sigma}^{-1}$ is replaced by its ridge estimate $\mathbf{H} = (\mathbf{I}_p + \gamma \widehat{\boldsymbol{\Sigma}})^{-1}$.

Benaych and Couillet 2016 ^a

$$\mathbf{H} \sim \mathbf{T}_{0,1} \propto (\mathbf{I}_p + \rho_0 \boldsymbol{\Sigma}_0 + \rho_1 \boldsymbol{\Sigma}_1)^{-1}$$

^aF. Benaych-Georges and R. Couillet, Spectral Analysis of the Gram Matrix of Mixture Models, ESAIM: Probability and Statistics, vol. 20, pp. 217237, 2016.

Elkhalil et al. 2018 ^a

$$\epsilon_{R-LDA}^{dist.} - \left\{ \frac{1}{2} \Phi \left[\frac{-\boldsymbol{\mu}^T \mathbf{T}_{0,1} \boldsymbol{\mu} + \beta}{\sqrt{D_0}} \right] + \frac{1}{2} \Phi \left[\frac{-\boldsymbol{\mu}^T \mathbf{T}_{0,1} \boldsymbol{\mu} - \beta}{\sqrt{D_1}} \right] \right\} \rightarrow_{prob.} 0$$

^aK. Elkhalil, A. Kammoun, R. Couillet, T. Al-Naffouri and M.-S. Alouini. A Large Dimensional Study of Regularized Discriminant Analysis Classifiers. Under review in IEEE Trans. Information Theory.

How is this different from the case of equal covariances?

$$\epsilon_{R-LDA}^{dist.} = \left\{ \frac{1}{2} \Phi \left[\frac{-\boldsymbol{\mu}^T \mathbf{T}_{0,1} \boldsymbol{\mu} + \beta}{\sqrt{D_0}} \right] + \frac{1}{2} \Phi \left[\frac{-\boldsymbol{\mu}^T \mathbf{T}_{0,1} \boldsymbol{\mu} - \beta}{\sqrt{D_1}} \right] \right\} \rightarrow_{prob.} 0$$

Some insights

- If $\|\boldsymbol{\Sigma}_0 - \boldsymbol{\Sigma}_1\| = o(1)$

$$\epsilon_{R-LDA}^{dist.} = \epsilon_{R-LDA}^{equal} + o(1).$$

→ R-LDA is robust against small perturbations.

- Different misclassification rates across classes.
- The enhancement in the misclassification rate in one class is likely to be lost by the other class.
- R-LDA does not leverage well the information about the covariance differences.

What about R-QDA ?

R-LDA	R-QDA
n_0, n_1 samples	n_0, n_1 samples
$\hat{\Sigma} = \frac{1}{n-2} \sum_{k=1}^2 \sum_{i=1}^{n_k} (\mathbf{x}_{k,i} - \bar{\mathbf{x}}_k) (\mathbf{x}_{k,i} - \bar{\mathbf{x}}_k)^T$	$\hat{\Sigma}_0 = \frac{1}{n_0-1} \sum_{i=1}^{n_0} (\mathbf{x}_{0,i} - \bar{\mathbf{x}}_0) (\mathbf{x}_{0,i} - \bar{\mathbf{x}}_0)^T$ $\hat{\Sigma}_1 = \frac{1}{n_1-1} \sum_{i=1}^{n_1} (\mathbf{x}_{1,i} - \bar{\mathbf{x}}_1) (\mathbf{x}_{1,i} - \bar{\mathbf{x}}_1)^T$

$$\epsilon_i = \mathbb{P} \left[\boldsymbol{\omega}^T \mathbf{B}_i \boldsymbol{\omega} + 2\boldsymbol{\omega}^T \mathbf{y}_i < \xi_i \right], \text{ where } \boldsymbol{\omega} \sim \mathcal{N}(\mathbf{0}_p, \mathbf{I}_p),$$

Asymptotic growth regime

1. Data scaling: $\frac{n_i}{p} \rightarrow c \in (0, \infty)$, with $|n_0 - n_1| = o(p)$.
2. Mean scaling: $\|\boldsymbol{\mu}_0 - \boldsymbol{\mu}_1\|^2 = O(\sqrt{p})$.
3. Covariance scaling: $\|\boldsymbol{\Sigma}_i\| = O(1)$.
4. $\boldsymbol{\Sigma}_0 - \boldsymbol{\Sigma}_1$ has exactly $O(\sqrt{p})$ eigenvalues of $O(1)$.

CLT(Lyapunov)

$$\epsilon_i^{R-QDA} - \Phi \left[(-1)^i \frac{1/\sqrt{p}\xi_i - 1/\sqrt{p} \operatorname{tr} \mathbf{B}_i}{\sqrt{1/p^2 \operatorname{tr} \mathbf{B}_i^2 + 1/p^4 \mathbf{y}_i^T \mathbf{y}_i}} \right] \rightarrow_{\text{prob.}} 0.$$

Elkhalil et al. 2017/2018 ^{a b}

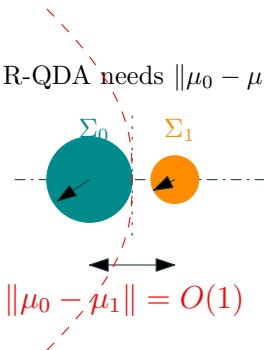
$$\epsilon_i^{R-QDA} - \left\{ \frac{1}{2} \Phi \left[\frac{\bar{\xi}_0 - \bar{b}_0}{\sqrt{2\bar{B}_0}} \right] + \frac{1}{2} \Phi \left[\frac{-\bar{\xi}_1 + \bar{b}_1}{\sqrt{2\bar{B}_1}} \right] \right\} \rightarrow_{prob.} 0.$$

$\bar{\xi}_i$, \bar{b}_i and \bar{B}_i depend on the classes' statistics.

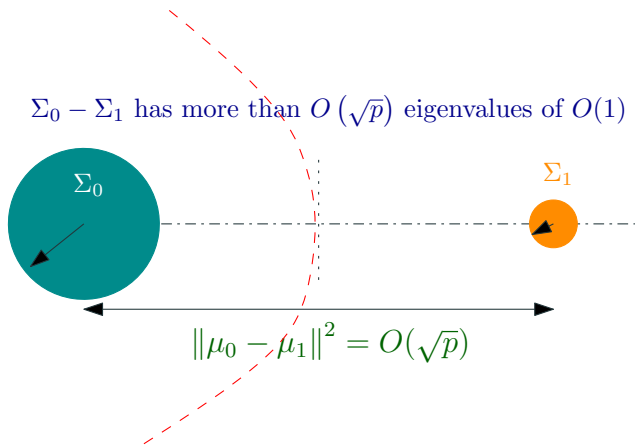
^aK. Elkhalil, A. Kammoun, R. Couillet, T. Y. Al-Naffouri, and M.-S. Alouini. Asymptotic Performance of Regularized Quadratic Discriminant Analysis-Based Classifiers. IEEE MLSP, Roppongi, Japan, Sept 2017.

^bK. Elkhalil, A. Kammoun, R. Couillet, T. Al-Naffouri and M.-S. Alouini. A Large Dimensional Study of Regularized Discriminant Analysis Classifiers. Under review in IEEE Trans. Information Theory.

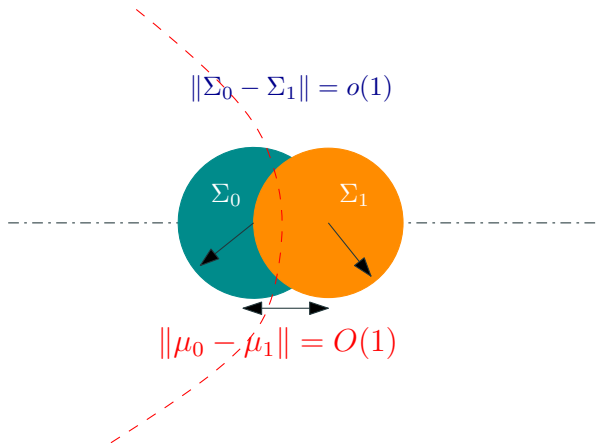
Recall that R-QDA needs $\|\mu_0 - \mu_1\|^2 = O(\sqrt{p})$



The information on the distance between the means is asymptotically useless!



R-QDA achieves asymptotic perfect classification.



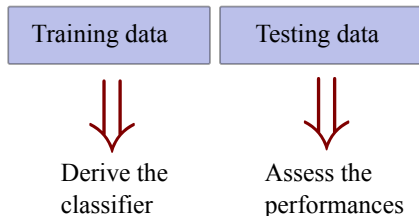
Classification is asymptotically impossible.

- Unbalanced training: $n_0 - n_1 = O(p)$.

R-QDA is equivalent to the naive classifier.

$$\epsilon \rightarrow \pi_0 \Phi(\infty) + \pi_1 \Phi(-\infty)$$

- Prone to estimation errors due to insufficiency in the number of observations.
- The tuning of the regularization parameter is very important



Model selection Given a set of candidate regularization factors

- Evaluate the performance using the test data for each regularization value ⁹
- Select the value that presents the lowest mis-classification rate

⁹J. Friedman. Regularized discriminant analysis. *Journal of the American Statistical Association*, 84:165175, 1989

Consistent estimator of the classification error

Exploiting the asymptotic equivalent

- If $p = O(1)$ and $n \rightarrow \infty$, then $\|\widehat{\Sigma}_i - \Sigma_i\| = o_p(1)$.
- When $p, n \rightarrow \infty$, then $\|\widehat{\Sigma}_i - \Sigma_i\| \neq o_p(1)$.

R-LDA (GE)

$$\widehat{\epsilon}_i^{R-LDA} - \Phi \left[\frac{(-1)^i G(\widehat{\mu}_i, \widehat{\mu}_0, \widehat{\mu}_1, \mathbf{H})}{\sqrt{D(\widehat{\mu}_0, \widehat{\mu}_1, \mathbf{H}, \widehat{\Sigma}_i)}} \right] \rightarrow_p 0$$

R-QDA (GE)

$$\widehat{\epsilon}_i^{R-QDA} - \Phi \left[(-1)^i \frac{\widehat{\xi}_i - \widehat{b}_i}{\sqrt{2\widehat{B}_i}} \right] \rightarrow_p 0$$

Optimal regularizer

$$\widehat{\gamma}^* = \arg \min_{\gamma > 0} \widehat{\epsilon}(\gamma).$$

- These results provides a glimpse on the region where the optimal γ is likely to belong.
- Perform a cross validation or testing in that region.

How well does this perform?

Benchmark estimation techniques:

- 5-fold cross-validation with 5 repetitions (5-CV).
- 0.632 bootstrap (B632).
- 0.632+ bootstrap (B632+)
- Plug-in estimator consisting of replacing the stats. in the DEs by their sample estimates.

Synthetic data

- $[\Sigma_0]_{i,j} = 0.6^{|i-j|}$.
- $\Sigma_1 = \Sigma_0 + 3 \begin{bmatrix} I_{\lceil \sqrt{p} \rceil} & \mathbf{0}_{k \times (p-k)} \\ \mathbf{0}_{(p-k) \times k} & \mathbf{0}_{(p-k) \times (p-k)} \end{bmatrix}$.
- $\mu_0 = [1, \mathbf{0}_{1 \times (p-1)}]^T$.
- $\mu_1 = \mu_0 + \frac{0.8}{\sqrt{p}} \mathbf{1}_{p \times 1}$.

Real data

- USPS dataset.
- $p = 256$ features (16×16) grayscale images.
- $n = 7291$ training examples.
- $n_{test} = 2007$ testing examples.



Performance: Synthetic data

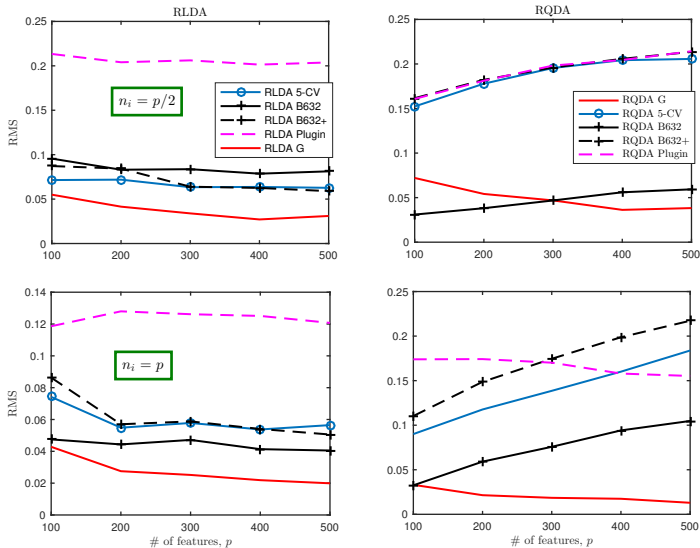


Figure 3.1: $n_0 = n_1$ and $\gamma = 1$.

Performance: Synthetic data

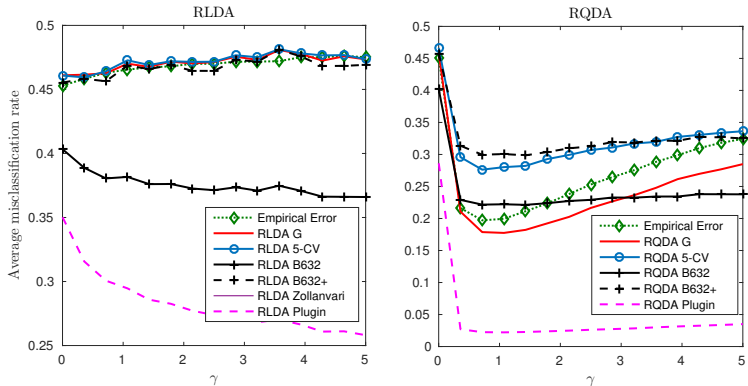


Figure 3.2: $p = 100$ features with equal training size ($n_0 = n_1 = p$).

Performance: USPS dataset

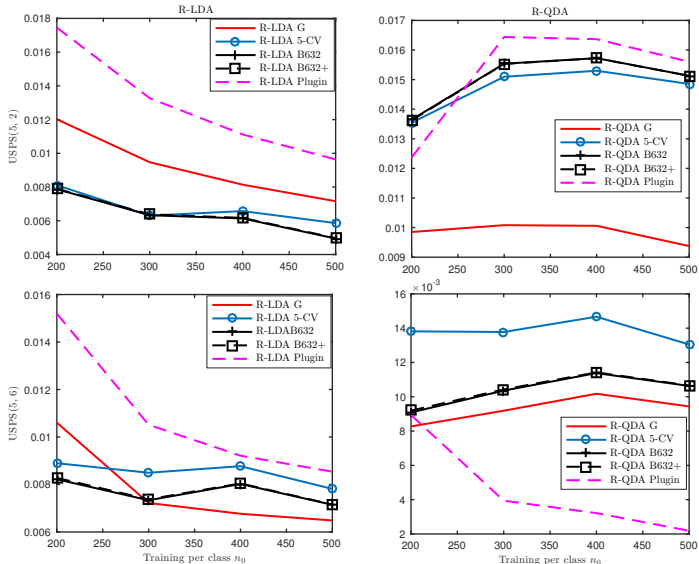
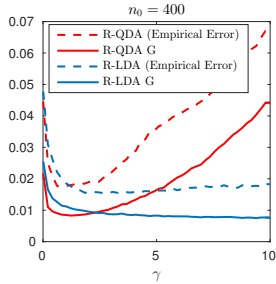
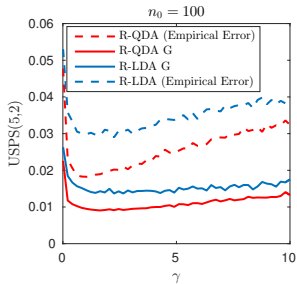
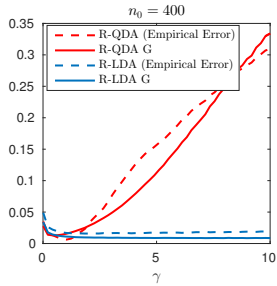
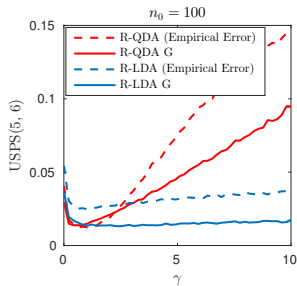
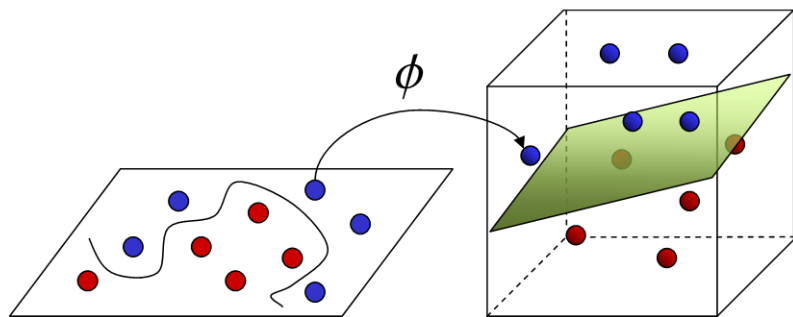


Figure 3.3: $n_0 = n_1$ and $\gamma = 1$. The first row gives the performance for the USPS data with digits (5, 2) whereas the second row considers the digits (5, 6).

Performance: USPS dataset



Centered Kernel Ridge Regression (CKRR)



Input Space

Feature Space

$$y = w^T x$$

$$y = w^T \phi(x)$$

KRR: Kernel trick

- $\{\mathbf{x}_i, y_i\}_{i=1}^n$ in $\mathcal{X} \times \mathcal{Y}$ s.t. $y_i = f(\mathbf{x}_i) + \sigma\epsilon_i$ with $\epsilon_i \sim_{\text{i.i.d}} \mathcal{N}(0, 1)$.
- Feature map: $\phi : \mathcal{X} \rightarrow \mathcal{H}$, with \mathcal{H} is a **RKHS**.
- Learning problem

$$\min_{\boldsymbol{\alpha}} \frac{1}{2} \|\mathbf{y} - \Phi\boldsymbol{\alpha}\|^2 + \frac{\lambda}{2} \|\boldsymbol{\alpha}\|^2$$

$$\Phi = [\phi(\mathbf{x}_1), \dots, \phi(\mathbf{x}_n)]^T \in \mathbb{R}^{|\mathcal{H}| \times n}$$

$$\boldsymbol{\alpha}^* = (\Phi^T \Phi + \lambda \mathbf{I}_{|\mathcal{H}|})^{-1} \Phi^T \mathbf{y} \in \mathbb{R}^{|\mathcal{H}|}$$

Woodbury

$$f^*(\mathbf{x}) = \phi(\mathbf{x})^T (\Phi^T \Phi + \lambda \mathbf{I}_{|\mathcal{H}|})^{-1} \Phi^T \mathbf{y}$$
$$f^*(\mathbf{x}) = \phi(\mathbf{x})^T \Phi^T (\Phi \Phi^T + \lambda \mathbf{I}_n)^{-1} \mathbf{y}$$

$$\{\Phi \Phi^T\}_{i,j} = \phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j)$$

$$f^*(\mathbf{x}) = \boldsymbol{\kappa}(\mathbf{x})^T (\mathbf{K} + \lambda \mathbf{I}_n)^{-1} \mathbf{y}$$

with $\boldsymbol{\kappa}(\mathbf{x})_i = \phi(\mathbf{x})^T \phi(\mathbf{x}_i)$ and $\mathbf{K}_{i,j} = \phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j)$.

Inner-product kernels

$$k(\mathbf{x}, \mathbf{x}') = \phi(\mathbf{x})^T \phi(\mathbf{x}') = g\left(\frac{\mathbf{x}^T \mathbf{x}'}{p}\right), \mathbf{x} \text{ and } \mathbf{x}' \in \mathcal{X}.$$

Asymptotic growth regime

Assumption 1.

- $p/n \rightarrow c(0, \infty)$.
- $\mathbb{E}\mathbf{x}_i = \mathbf{0}$ and $\text{cov}\mathbf{x}_i = \Sigma$ unif. bounded in p (e.g. $\mathbf{x}_i \sim \mathcal{N}(\mathbf{0}, \Sigma)$).

El Karoui 2010^a

$$\|\mathbf{K} - \mathbf{K}^\infty\| \rightarrow_{a.s.} 0$$

$$\text{with } \mathbf{K}^\infty = \underbrace{g(0) \mathbf{1}\mathbf{1}^T}_{\|\cdot\|=O(p)} + \underbrace{g'(0) \frac{\mathbf{X}\mathbf{X}^T}{p}}_{\|\cdot\|=O(1)} + \text{constant}(g, \Sigma).$$

^aN. El-Karoui, The Spectrum of Kernel Random Matrices, The Annals of Statistics, vol. 38, no. 1, pp. 150, 2010.

Inner-product kernels

$$k(\mathbf{x}, \mathbf{x}') = \phi(\mathbf{x})^T \phi(\mathbf{x}') = g\left(\frac{\mathbf{x}^T \mathbf{x}'}{p}\right), \quad \mathbf{x} \text{ and } \mathbf{x}' \in \mathcal{X}.$$

Asymptotic growth regime

Assumption 1.

- $p/n \rightarrow c(0, \infty)$.
- $\mathbb{E}\mathbf{x}_i = \mathbf{0}$ and $\text{cov}\mathbf{x}_i = \Sigma$ unif. bounded in p (e.g. $\mathbf{x}_i \sim \mathcal{N}(\mathbf{0}, \Sigma)$).

Centering with $\mathbf{P} = \mathbf{I}_n - \frac{\mathbf{1}\mathbf{1}^T}{n}$
 $\mathbf{K}_c = \mathbf{P}\mathbf{K}\mathbf{P}$

El Karoui 2010^a

$$\|\mathbf{K} - \mathbf{K}^\infty\| \xrightarrow{\text{a.s.}} 0$$

$$\text{with } \mathbf{K}^\infty = \underbrace{g(0)\mathbf{1}\mathbf{1}^T}_{\|\cdot\|=O(p)} + g'(0) \underbrace{\frac{\mathbf{X}\mathbf{X}^T}{p}}_{\|\cdot\|=O(1)} + \text{constant}(g, \Sigma).$$

^aN. El-Karoui, The Spectrum of Kernel Random Matrices, The Annals of Statistics, vol. 38, no. 1, pp. 150, 2010.

$$\mathbf{P} = \mathbf{I}_n - \frac{\mathbf{1}\mathbf{1}^T}{n}.$$

Learning problem

$$\min_{\alpha_0, \alpha} \frac{1}{2} \|\mathbf{y} - \Phi\alpha - \alpha_0\mathbf{1}_n\|^2 + \frac{\lambda}{2} \|\alpha\|^2 \Leftrightarrow \min_{\alpha} \frac{1}{2} \|\mathbf{P}(\mathbf{y} - \Phi\alpha)\|^2 + \frac{\lambda}{2} \|\alpha\|^2$$

$$\alpha^* = \Phi^T \mathbf{P} \left(\underbrace{\mathbf{P}\mathbf{K}\mathbf{P}}_{\mathbf{K}_c} + \lambda\mathbf{I}_n \right)^{-1} (\mathbf{y} - \bar{y}\mathbf{1}_n)$$

$$f_c^*(\mathbf{x}) = \kappa_c(\mathbf{x})^T (\mathbf{K}_c + \lambda\mathbf{I}_n)^{-1} \mathbf{P}\mathbf{y} + \bar{y}.$$

$$\kappa_c(\mathbf{x}) = \mathbf{P}\kappa(\mathbf{x}) - \frac{1}{n}\mathbf{P}\mathbf{K}\mathbf{1}_n, \quad \phi_c(\mathbf{x}) = \phi(\mathbf{x}) - \frac{1}{n} \sum_{i=1}^n \phi(\mathbf{x}_i).$$

Centered KRR \sim KRR with centered kernels

What about the performance ?

Performance metrics

$$\mathcal{R}_{train} = \frac{1}{n} \mathbb{E}_{\epsilon} \left\| \widehat{f}_{\epsilon}(X) - f(X) \right\|_2^2$$

$$\mathcal{R}_{test} = \mathbb{E}_{s \sim \mathcal{D}, \epsilon} \left| \widehat{f}_{\epsilon}(s) - f(s) \right|^2$$

Assumption 1. (Growth rate) As $p, n \rightarrow \infty$ we assume the following

- **Data scaling:** $p/n \rightarrow c \in (0, \infty)$.
- **Covariance scaling:** $\limsup_p \|\Sigma\| < \infty$.

Assumptions 2. (kernel function)

$$\mathbb{E} \left| g^{(3)} \left(\frac{1}{p} \mathbf{x}_i^T \mathbf{x}_j \right) \right|^k < \infty.$$

Assumption 3. (Data generating function)

-

$$\mathbb{E}_{x \sim \mathcal{N}(\mathbf{0}, \Sigma)} |f(\mathbf{x})|^k < \infty,$$

-

$$\mathbb{E}_{x \sim \mathcal{N}(\mathbf{0}, \Sigma)} \|\nabla_f(\mathbf{x})\|_2^k < \infty, \text{ where } \nabla_f(\mathbf{x}) = \left\{ \frac{\partial f(\mathbf{x})}{\partial x_l} \right\}_{l=1}^p.$$

Limiting risk^a Let $z = -\frac{\lambda + g(\tau) - g(0) - \tau g'(0)}{g'(0)}$ with $\tau = \frac{1}{p} \text{tr } \Sigma$.

$$\mathcal{R}_{\text{train}} - \mathcal{R}_{\text{train}}^{\infty} \rightarrow_{\text{prob.}} 0,$$

$$\mathcal{R}_{\text{test}} - \mathcal{R}_{\text{test}}^{\infty} \rightarrow_{\text{prob.}} 0,$$

$$\mathcal{R}_{\text{train}}^{\infty} = \left(\frac{c\lambda m_z}{g'(0)} \right)^2 \frac{n(1+m_z)^2 (\sigma^2 + \text{var}_f) - nm_z(2+m_z) \|\mathbb{E}[\nabla_f]\|^2}{n(1+m_z)^2 - pm_z^2} + \sigma^2 - 2\sigma^2 \frac{c\lambda m_z}{g'(0)}$$

$$\mathcal{R}_{\text{test}}^{\infty} = \frac{n(1+m_z)^2 (\sigma^2 + \text{var}_f) - nm_z(2+m_z) \|\mathbb{E}[\nabla_f]\|^2}{n(1+m_z)^2 - pm_z^2} - \sigma^2.$$

^aK. Elkhailil, A. Kammoun, X. Zhang, M.-S. Alouini and T. Al-Naffouri.

Risk Convergence of Centered Kernel Ridge Regression with Large Dimensional Data.
Submitted to IEEE Trans. Signal Processing.

Bad news ☹️

Minimum prediction risk is achieved by all kernels!! ~ Linear kernel

Good news 😊

kernel/regularizer can be jointly optimized!

Interesting relation between $\mathcal{R}_{train}^\infty$ and $\mathcal{R}_{test}^\infty$

$$\mathcal{R}_{test}^\infty = \left(\frac{c\lambda m_z}{g'(0)} \right)^{-2} \mathcal{R}_{train}^\infty - \sigma^2 \left(\frac{g'(0)}{c\lambda m_z} - 1 \right)^2.$$

Consistent estimator of $\widehat{\mathcal{R}}_{test}$

$$\widehat{\mathcal{R}}_{test} = \left(\frac{c\lambda \widehat{m}_z}{g'(0)} \right)^{-2} \widehat{\mathcal{R}}_{train} - \sigma^2 \left(\frac{g'(0)}{c\lambda \widehat{m}_z} - 1 \right)^2,$$
$$\widehat{m}_z = \frac{1}{p} \operatorname{tr} \left(\frac{XX^T}{p} - zI_n \right)^{-1}.$$

Issues with λ small ☹

Consistent estimator of $\widehat{\mathcal{R}}_{test}$

$$\widehat{\mathcal{R}}_{test} = \frac{1}{(cz\widehat{m}_z)^2} \left[\frac{1}{np} \mathbf{y}^T \mathbf{P} \mathbf{X} \left(z\widetilde{\mathbf{Q}}_z^2 - \widetilde{\mathbf{Q}}_z \right) \mathbf{X}^T \mathbf{P} \mathbf{y} + \text{var}(\mathbf{y}) \right] - \sigma^2.$$
$$\widetilde{\mathbf{Q}}_z = \left(\frac{\mathbf{X}^T \mathbf{P} \mathbf{X}}{p} - z\mathbf{I}_p \right)^{-1}.$$

More stable with respect to λ ☺

$$\widehat{\mathcal{R}}_{test}^* = \min_{z \notin \text{Supp}\{XX^T/p\}} \widehat{\mathcal{R}}_{test}(z), \quad z^* = -\frac{\lambda^* + g^*(\tau) - g^*(0) - \tau g'^*(0)}{g'^*(0)}.$$

Kernels

- Linear kernels: $k(\mathbf{x}, \mathbf{x}') = \alpha \mathbf{x}^T \mathbf{x}' / p + \beta$.
- Polynomial kernels: $k(\mathbf{x}, \mathbf{x}') = (\alpha \mathbf{x}^T \mathbf{x}' / p + \beta)^d$.
- Sigmoid kernels: $k(\mathbf{x}, \mathbf{x}') = \tanh(\alpha \mathbf{x}^T \mathbf{x}' / p + \beta)$.
- Exponential kernels: $k(\mathbf{x}, \mathbf{x}') = \exp(\alpha \mathbf{x}^T \mathbf{x}' / p + \beta)$.

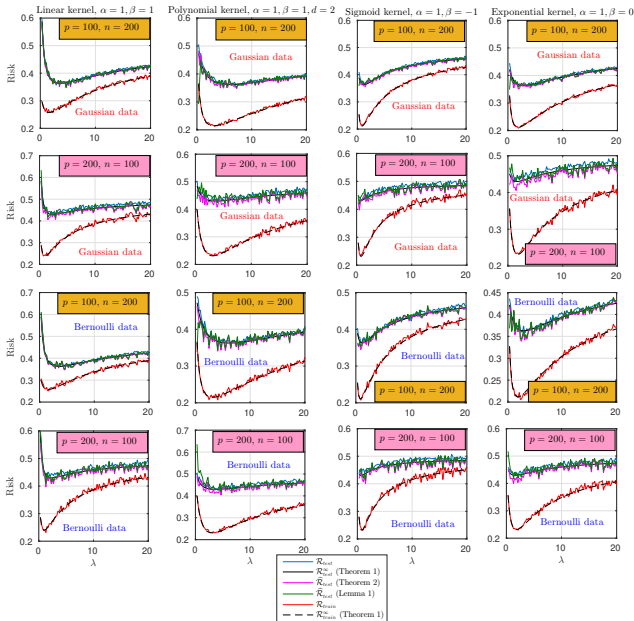
Synthetic data

- $\mathbf{x} \sim \Sigma^{\frac{1}{2}} \mathbf{z}$ with $\mathbf{z} = \{z_i\}_{i=1}^p$, $\mathbb{E} z_i = 0$, $\text{var} z_i = 1$ and $\mathbb{E} z_i^k = O(1)$.
- Generating function: $f(\mathbf{x}) = \sin\left(\frac{\mathbf{1}^T \mathbf{x}}{\sqrt{p}}\right)$.

Real data

- Communities and Crime dataset.
- $p = 122$, $n_{train} = 73$ and $n_{test} = 50$.
- Prediction risk is computed by averaging over 500 data shuffling.

CKRR: Synthetic data



CKRR: Synthetic data

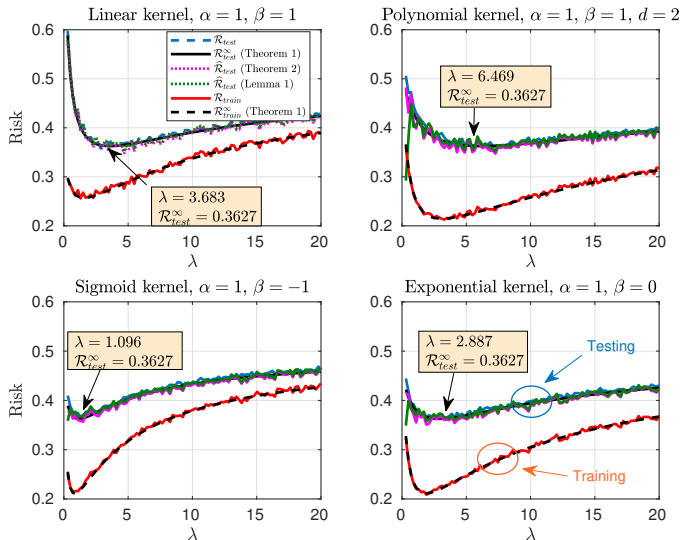


Figure 4.1: CKRR risk with respect to the regularization parameter λ on Gaussian data ($x \sim \mathcal{N}(0_p, \{0.4^{|i-j|}\}_{i,j})$), $n = 200$ training samples and $p = 100$ predictors.

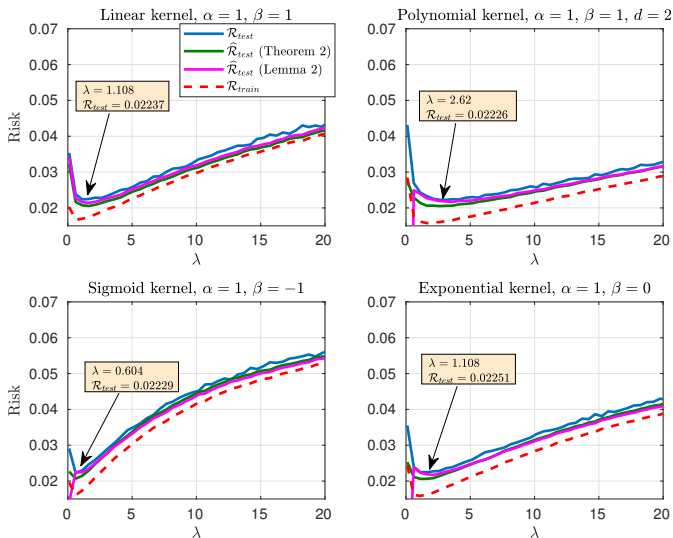


Figure 4.2: CKRR risk with respect to λ where independent zero mean Gaussian noise samples with variance $\sigma^2 = 0.05$ are added to the true response.

Conclusion

- Random matrix theory is a powerful tool that has been applied with success to the fields wireless communications and signal processing, providing solutions to very challenging problems
- High dimensionality along with stochasticity are the sole prerequisite of this tool
- Successful application of this tool has been demonstrated in the context of RDA.
- Fundamental limits of Centered kernel ridge regression.

Future research directions

- We can also consider the performance analysis of kernel LDA/QDA.
- Extend the analysis to *Homogenous* kernels.

Important results on Homogenous Kernel matrices

- $\phi(\mathbf{x})$ is a fixed non linear feature space mapping. The kernel function is given by

$$k(\mathbf{x}, \mathbf{x}') = \phi(\mathbf{x})^T \phi(\mathbf{x}').$$

- Homogeneous kernels

$$k(\mathbf{x}, \mathbf{x}') = f\left(\frac{\|\mathbf{x} - \mathbf{x}'\|^2}{p}\right).$$

- $\{\mathbf{K}\}_{i,j} = k(\mathbf{x}_i, \mathbf{x}_j)$.

Theorem (Spectrum of kernel random matrices, El-Karoui 2010)

¹⁰ [Informal statement]

$$\widehat{\mathbf{K}} = f(\tau) \mathbf{1}\mathbf{1}^T + f'(\tau) \mathbf{W} + f''(\tau) \mathbf{Q}, \quad \frac{\|\mathbf{x} - \mathbf{x}'\|^2}{p} \rightarrow_{a.s.} \tau.$$

$$\|\mathbf{K} - \widehat{\mathbf{K}}\| \xrightarrow{p} 0.$$

(7)

This might help to analyze the performance of some kernel methods in regression or classification.

¹⁰N. El Karoui, The spectrum of kernel random matrices, The annals of statistics, Volume 38, Number 1 (2010), 1-50.

That's it

Thank you for your time and attention!