

# Regularized Discriminant Analysis: A Large Dimensional Study

ISIT 2018

---

Xiaoke Yang

**Khalil Elkhail**

Abla Kammoun

Tareq Y. Al-Naffouri

Mohamed-Slim Alouini

June 19, 2018

King Abdullah University of Science and Technology



1. Gaussian discriminant analysis
2. Asymptotic Performance of RDA
3. Discussion
4. Conclusions

## Gaussian discriminant analysis

---

Classification is the task of assigning a label to an input data among a set of possible categories.

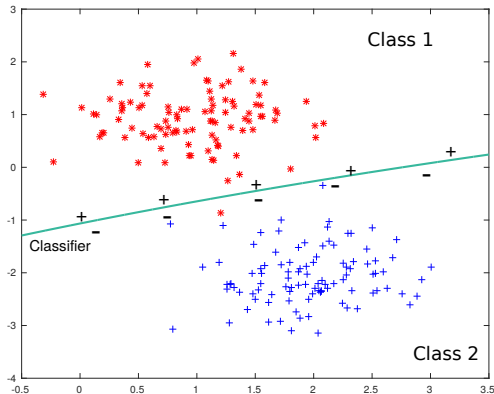


---

Yann LeCun, Leon Bottou, Yoshua Bengio, and Patrick Haffner. Gradient-Based Learning Applied To Document Recognition. Proceedings of the IEEE, 86(11):2278-2324, 1998.

# 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.

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

- 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 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) ( $\boldsymbol{\Sigma}_0 = \boldsymbol{\Sigma}_1$ )

$$W^{LDA}(\mathbf{x}) = \left( \mathbf{x} - \frac{\boldsymbol{\mu}_0 + \boldsymbol{\mu}_1}{2} \right)^T \boldsymbol{\Sigma}^{-1} (\boldsymbol{\mu}_0 - \boldsymbol{\mu}_1) - \log \frac{\pi_1}{\pi_0} \quad (1)$$

$$\begin{cases} \text{Assign } \mathbf{x} \text{ to class } 0 & \text{if } W^{LDA} > 0 \\ \text{Assign } \mathbf{x} \text{ to class } 1 & \text{otherwise} \end{cases}$$

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

### Quadratic discriminant analysis ( $\Sigma_0 \neq \Sigma_1$ )

$$W^{QDA}(\mathbf{x}) = -\frac{1}{2} \log \frac{|\Sigma_0|}{|\Sigma_1|} - \frac{1}{2}(\mathbf{x} - \mu_0)^T \Sigma_0^{-1}(\mathbf{x} - \mu_0) + \frac{1}{2}(\mathbf{x} - \mu_1)^T \Sigma_1^{-1}(\mathbf{x} - \mu_1)$$

(2)

$$\begin{cases} \text{Assign } \mathbf{x} \text{ to class 0} & \text{if } W^{QDA}(\mathbf{x}) > 0 \\ \text{Assign } \mathbf{x} \text{ to class 1} & \text{otherwise} \end{cases}$$

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

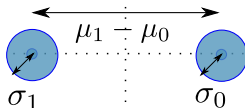


## LDA: finite regime (statistics are known)

- Assume  $\Sigma$ ,  $\mu_0$  and  $\mu_1$  known.
- Equal priors :  $\pi_0 = \pi_1 = 0.5$
- No asymptotic regime,  $p$  is fixed.

The total misclassification rate is given by <sup>1</sup>

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



<sup>1</sup>Jerome Friedman, Trevor Hastie, and Robert Tibshirani. The Elements of Statistical Learning. Springer, 2009.

For Gaussian data, LDA and QDA are respectively the *best* classifiers (with the smallest possible risk) in the case of equal covariances and distinct covariances!

How things will look like with noisy estimates of the statistics?

Consider the supervised setting

$$\bar{\mathbf{x}}_i = \frac{1}{n_i} \sum_{\mathbf{x}_j \in \mathcal{C}_i} \mathbf{x}_j$$

$$\hat{\Sigma}_i = \frac{1}{n_i - 1} \sum_{\mathbf{x}_j \in \mathcal{C}_i} (\mathbf{x}_j - \bar{\mathbf{x}}_i)(\mathbf{x}_j - \bar{\mathbf{x}}_i)^T$$

- Introduce regularization parameter  $\lambda \in [0, 1]$

$$\hat{\Sigma}_i(\lambda) = \frac{(1 - \lambda)n_i \hat{\Sigma}_i + \lambda n \hat{\Sigma}}{(1 - \lambda)n_i + \lambda n}, \quad i \in \{0, 1\}$$

- Introduce regularization parameter  $\gamma \in [0, 1]$  to avoid singularity

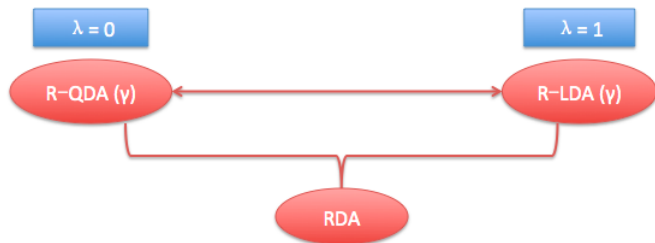
$$\hat{\Sigma}_i(\lambda, \gamma) = \gamma \frac{(1 - \lambda)n_i \hat{\Sigma}_i + \lambda n \hat{\Sigma}}{(1 - \lambda)n_i + \lambda n} + (1 - \gamma)\mathbf{I}_p.$$

Both R-LDA and R-QDA are special cases of RDA.

- $\lambda = 0 \rightarrow$  R-QDA
- $\lambda = 1 \rightarrow$  R-LDA
- Define  $\mathbf{H}_i = \hat{\Sigma}_i^{-1}$

---

<sup>2</sup>J. Friedman, Regularized discriminant analysis, Journal of the American Statistical Association, vol. 84, pp. 165{175, 1989}



The extreme cases (R-LDA and R-QDA) have been analyzed in <sup>3</sup> and <sup>4</sup>

<sup>3</sup>K. Elkhailil, A. Kammoun, R. Couillet, T. Y. Al-Naffouri, and M.-S. Alouini. A Large Dimensional Study of Regularized Discriminant Analysis Classifiers. ArXiv e-prints, Nov. 2017.

<sup>4</sup>A. Zollanvari and E. R. Dougherty, Generalized Consistent Error Estimator of Linear Discriminant Analysis, IEEE Transactions on Signal Processing, vol. 63, no. 11, pp. 28042814, June 2015

# LDA: Asymptotic regime (equal covariances)

## Asymptotic growth regime

Let  $n = n_0 + n_1$ .

- $n_0, n_1, p \rightarrow \infty$  such that  $\frac{n_0}{n_1} \rightarrow 1$  and  $\frac{p}{n} \rightarrow c < 1$
- $\Sigma_0 = \Sigma_1$
- $\mu \triangleq \mu_0 - \mu_1$  is such that  $\|\mu\| = O(1)$ .

Under these assumptions, the misclassification rate converges to <sup>5</sup>:

$$R_{LDA} = \Phi \left[ -\frac{\Delta}{2} \sqrt{1-c} \right] \rightarrow_{a.s.} 0. \quad (3)$$

→ 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  $c < 1$ .

---

<sup>5</sup>Cheng Wang and Binyan Jiang. On the dimension effect of regularized linear discriminant analysis, arXiv:1710.03136v1

## Asymptotic Performance of RDA

---

- Make some simplification

$$\mathbf{H}_0 = \left[ (1 - \gamma)\mathbf{I}_p + \alpha_0 \widehat{\boldsymbol{\Sigma}}_0 + \beta_0 \widehat{\boldsymbol{\Sigma}}_1 \right]^{-1}$$

$$\alpha_0 = \frac{n_0 \gamma}{n_0 + \lambda n_1}, \beta_0 = \frac{n_1 \gamma \lambda}{n_0 + \lambda n_1}$$

$$\mathbf{H}_1 = \left[ (1 - \gamma)\mathbf{I}_p + \alpha_1 \widehat{\boldsymbol{\Sigma}}_0 + \beta_1 \widehat{\boldsymbol{\Sigma}}_1 \right]^{-1}$$

$$\alpha_1 = \frac{n_0 \gamma \lambda}{n_1 + \lambda n_0}, \beta_1 = \frac{n_1 \gamma}{n_1 + \lambda n_0}$$

- The Discriminant score for the RDA is:

$$\widehat{\delta}_i^{RDA}(\mathbf{x}) = \frac{1}{2} \log |\mathbf{H}_i| - \frac{1}{2} (\mathbf{x} - \bar{\mathbf{x}}_i)^T \mathbf{H}_i (\mathbf{x} - \bar{\mathbf{x}}_i) + \log \pi_i.$$

- The conditional misclassification error rate associated to class  $C_i$

$$\epsilon_i^{RDA} = \mathbb{P} \left[ (-1)^i \widehat{\delta}_0^{RDA}(\mathbf{x}) < (-1)^i \widehat{\delta}_1^{RDA}(\mathbf{x}) \mid \mathbf{x} \in C_i \right]. \quad (4)$$

The conditional misclassification error can easily be shown to write as

$$\epsilon_i = \mathbb{P} \left[ \boldsymbol{\omega}^T \mathbf{B}_i \boldsymbol{\omega} + 2\boldsymbol{\omega}^T \mathbf{y}_i < \xi_i \right], \text{ where } \boldsymbol{\omega} \sim \mathcal{N}(\mathbf{0}_p, \mathbf{I}_p), \quad (5)$$

$$\mathbf{B}_i = \boldsymbol{\Sigma}_i^{1/2} (\mathbf{H}_1 - \mathbf{H}_0) \boldsymbol{\Sigma}_i^{1/2},$$

$$\mathbf{y}_i = \boldsymbol{\Sigma}_i^{1/2} [\mathbf{H}_1 (\boldsymbol{\mu}_i - \bar{\mathbf{x}}_1) - \mathbf{H}_0 (\boldsymbol{\mu}_i - \bar{\mathbf{x}}_0)],$$

$$\xi_i = -\log \left( \frac{|\mathbf{H}_0|}{|\mathbf{H}_1|} \right) + (\boldsymbol{\mu}_i - \bar{\mathbf{x}}_0)^T \mathbf{H}_0 (\boldsymbol{\mu}_i - \bar{\mathbf{x}}_0) - (\boldsymbol{\mu}_i - \bar{\mathbf{x}}_1)^T \mathbf{H}_1 (\boldsymbol{\mu}_i - \bar{\mathbf{x}}_1) + 2 \log \frac{\pi_1}{\pi_0}.$$

### Assumptions

1.  $p, n_1, n_2 \rightarrow \infty$  with  $\frac{n_i}{p} \rightarrow c(0, \infty)$ ,  $\frac{n_0}{n} = \frac{n_1}{n} \rightarrow 0.5$
2. The difference in means  $\boldsymbol{\mu} = \boldsymbol{\mu}_0 - \boldsymbol{\mu}_1$  satisfies  $\|\boldsymbol{\mu}\|^2 = O(\sqrt{p})$ .
3.  $\|\boldsymbol{\Sigma}_i\| = O(1)$ .
4. Matrix  $\boldsymbol{\Sigma}_0 - \boldsymbol{\Sigma}_1$  has at most  $O(\sqrt{p})$  eigenvalues of order 1 while the remaining ones decay at an order of  $O(1/\sqrt{p})$ .



### Central Limit Theorem (CLT)<sup>6</sup>

Assume  $\lambda \neq 1$ . Under assumptions 1-4,  $\omega^T \mathbf{B}_i \omega + 2\omega^T \mathbf{y}_i$  can be treated as a Gaussian random variable  $\mathbf{z} \sim \mathcal{N}(\text{tr} \mathbf{B}_i, 2 \text{tr} \mathbf{B}_i^2 + 4\mathbf{y}_i^T \mathbf{y}_i)$ , so the conditional classification error of RDA satisfies

$$\epsilon_i^{RDA} - \Phi \left( (-1)^i \frac{\xi_i - \text{tr} \mathbf{B}_i}{\sqrt{2 \text{tr} \mathbf{B}_i^2 + 4\mathbf{y}_i^T \mathbf{y}_i}} \right) \xrightarrow{\text{a.s.}} 0. \quad (6)$$

---

<sup>6</sup>K. Elkhail, A. Kammoun, R. Couillet, T. Y. Al-Naffouri, and M.-S. Alouini. A Large Dimensional Study of Regularized Discriminant Analysis Classifiers. ArXiv e-prints, Nov. 2017.

*Benaych and Couillet, 2016*

$$\mathbf{H}_i = \left( \frac{\mathbf{y}_0 \mathbf{y}_0^T}{p} + \frac{\mathbf{y}_1 \mathbf{y}_1^T}{p} + (1 - \gamma) \mathbf{I} \right)^{-1} : \text{Gram matrix of mixture models}^7$$

Define

$$\delta_i = \frac{1}{n_i} \text{tr} \Sigma_i \mathbf{Q}_i, \quad \tilde{\delta}_i = \frac{\alpha_i}{1 + \alpha_i \delta_i}$$

where

$$\mathbf{Q}_i = \left[ (1 - \gamma) \mathbf{I}_p + \frac{\alpha_i}{1 + \alpha_i \delta_i} \Sigma_0 + \frac{\beta_i}{1 + \beta_i \delta_i} \Sigma_1 \right]^{-1},$$

$$\frac{1}{p} \text{tr} \mathbf{A} (\mathbf{H}_i - \mathbf{Q}_i) \rightarrow_p 0, \quad \mathbf{u}^T (\mathbf{H}_i - \mathbf{Q}_i) \mathbf{v} \rightarrow_p 0. \quad (7)$$

<sup>7</sup>F. Benaych-Georges and R. Couillet, Spectral Analysis of the Gram Matrix of Mixture Models, ESAIM: Probability and Statistics, vol. 20, pp. 217237, 2016.

## RDA (Asymptotic result)

With assumptions 1-4 satisfied, we have

$$\epsilon_i - \Phi \left( (-1)^i \frac{\bar{\xi}_i - \bar{b}_i}{\sqrt{2\bar{B}_i}} \right) \xrightarrow{p} 0. \quad (8)$$

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

$$\begin{aligned} \bar{\xi}_i \triangleq & \frac{1}{\sqrt{\rho}} \log \left( \frac{1 + \alpha_0 \delta_0}{1 + \alpha_1 \delta_1} \right)^{n_0} \left( \frac{1 + \beta_0 \delta_0}{1 + \beta_1 \delta_1} \right)^{n_1} + \frac{1}{\sqrt{\rho}} \left[ \frac{\alpha_1 \delta_1 n_0 - \alpha_0 \delta_0 n_0}{(1 + \alpha_1 \delta_1)(1 + \alpha_0 \delta_0)} + \frac{\beta_1 \delta_1 n_0 - \beta_0 \delta_0 n_0}{(1 + \beta_1 \delta_1)(1 + \beta_0 \delta_0)} \right] \\ & + \frac{1}{\sqrt{\rho}} \log \frac{|\mathbf{Q}_1|}{|\mathbf{Q}_0|} + \frac{1}{\sqrt{\rho}} (-1)^{i+1} \boldsymbol{\mu}^T \mathbf{Q}_{1-i} \boldsymbol{\mu}. \end{aligned}$$

$$\bar{b}_i = \frac{1}{\sqrt{\rho}} \text{tr} \boldsymbol{\Sigma}_i (\mathbf{Q}_1 - \mathbf{Q}_0).$$

$$\bar{B}_i \triangleq \frac{1}{\rho} \frac{2n_1 \phi}{1 - (\tilde{\delta}_1^2 + \tilde{\delta}_0^2) \phi} - \frac{1}{\rho} \frac{2n_1 \phi}{1 - 2\tilde{\delta}_0 \tilde{\delta}_1 \phi}.$$

$$\phi = \frac{1}{n_1} \text{tr} \boldsymbol{\Sigma}_1 \mathbf{Q}_1 \boldsymbol{\Sigma}_1 \mathbf{Q}_1.$$

## Discussion

---

## What happens if

- $\|\mu_0 - \mu_1\| = O(1)$ ?

## What happens if

- $\|\mu_0 - \mu_1\| = O(1)$ ?

The difference in means will not be asymptotically used by RDA.

**What happens if**

- $\|\mu_0 - \mu_1\| = O(1)$ ?

The difference in means will not be asymptotically used by RDA.

- Matrix  $\Sigma_0 - \Sigma_1$  has more than  $O(\sqrt{p})$  eigenvalues of order 1?

### What happens if

- $\|\mu_0 - \mu_1\| = O(1)$ ?

The difference in means will not be asymptotically used by RDA.

- Matrix  $\Sigma_0 - \Sigma_1$  has more than  $O(\sqrt{p})$  eigenvalues of order 1?

RDA will perform asymptotically perfect classification.

- Matrix  $\Sigma_0 - \Sigma_1$  has less than  $O(\sqrt{p})$  eigenvalues of order 1 and  $\|\mu_0 - \mu_1\| = O(1)$ ?



## What happens if

- $\|\mu_0 - \mu_1\| = O(1)$ ?

The difference in means will not be asymptotically used by RDA.

- Matrix  $\Sigma_0 - \Sigma_1$  has more than  $O(\sqrt{p})$  eigenvalues of order 1?

RDA will perform asymptotically perfect classification.

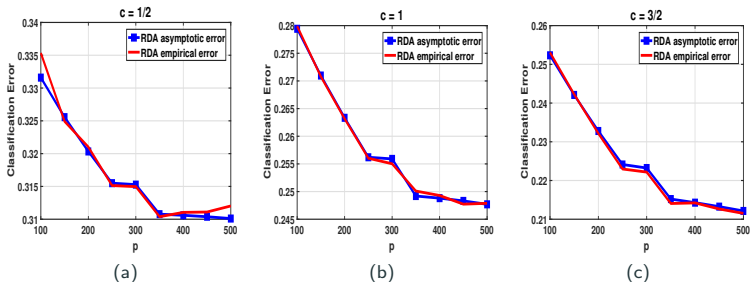
- Matrix  $\Sigma_0 - \Sigma_1$  has less than  $O(\sqrt{p})$  eigenvalues of order 1 and  $\|\mu_0 - \mu_1\| = O(1)$ ?

The misclassification rate of RDA will converge to 0.5. Classification is asymptotically impossible.

# RDA: Synthetic data Experiment (1)

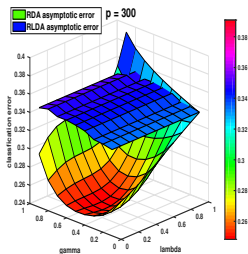
$$\{\Sigma_0\}_{i,j} = 0.6^{|i-j|}, \Sigma_1 = \Sigma_0 + 2 \begin{bmatrix} \mathbf{I}_k & \mathbf{0}_{k \times (p-k)} \\ \mathbf{0}_{(p-k) \times k} & \mathbf{0}_{(p-k) \times (p-k)} \end{bmatrix}, k = \lfloor \sqrt{p} \rfloor, \mu_0 = \mathbf{1}_{p \times 1},$$

$$\mu_1 = \mu_0 + 2\rho^{-\frac{1}{4}} \mathbf{1}_{p \times 1}, \lambda = 0.5, \gamma = 0.5, c = \frac{n_1}{p}$$

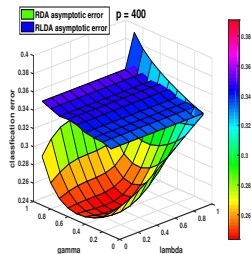


**Figure 3.1:** RDA classifier performance in terms of classification error with equal training,  $n_0 = n_1$ . The x axis is the number of the data dimension.

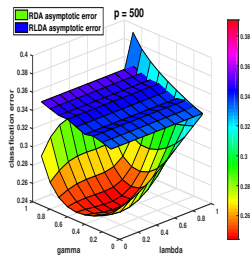
# RDA: Synthetic data Experiment (2)



(a)



(b)



(c)

- $\lambda$  approaching the minimum classification error is neither 1 nor 0

The optimal classifier minimizing the classification error is neither R-LDA nor R-QDA.

- RDA offers better classification performance than R-LDA and R-QDA with proper regularizers selection

## Conclusions

---

- Leveraging results from RMT, we identified the growth rate regime in which RDA results in a non-trivial classification risk.
- We derived a closed-form expression for the asymptotic classification risk reflecting the impact of the data statistics on the performance.
- Practical insights are drawn to help design the optimal classifier.

Thank you!