

# Regularized Quadratic Discriminant Analysis for Large Dimensional Data

Khalil Elkhail<sup>1</sup>, Abla Kammoun<sup>1</sup>, Romain Couillet<sup>2</sup>, Tareq Y. Al-Naffouri<sup>1</sup> and Mohamed Slim Alouini<sup>1</sup>



<sup>1</sup> CEMSE Division, King Abdullah University of Science and Technology (KAUST), KSA

<sup>2</sup> CNRS-CentraleSupélec-Université Paris-Sud, France



## Introduction

### Motivation

- ▶ Closed form expressions for the misclassification probability exists in the case where the training size  $n$  is larger than the number of predictors  $p$  (the signal dimension) using some properties of Wishart distribution (Exact analysis)
- ▶ Most analysis is only valid for Gaussian data limiting the scope of application.
- ▶ Most of asymptotic approaches rely on the assumption that  $p \gg n$  and assume some sparsity on the class' statistics.

### Contributions

- ▶ We consider the double asymptotic regime where both  $p, n \rightarrow \infty$  with  $p/n \rightarrow c \in (0, \infty)$ .
- ▶ We make mild assumptions on the statistical means and the covariance matrices just to achieve non trivial error rates.
- ▶ We derive an asymptotic limit for the misclassification probability that reveals the mathematical connection between the classification error and the statistical parameters associated with each class.
- ▶ We leverage this result to propose a more efficient design of the regularized QDA classifier by selecting the regularization parameter that minimizes the asymptotic classification error.

## Regularized QDA for Binary Classification

### QDA

An observation  $\mathbf{x} \in \mathbb{R}^{p \times p}$  belongs to class  $\mathcal{C}_i$ ,  $i = 0, 1$ , if and only if

$$\mathbf{x} = \boldsymbol{\mu}_i + \boldsymbol{\Sigma}_i^{1/2} \boldsymbol{\omega},$$

with  $\boldsymbol{\omega} \sim \mathcal{N}(\mathbf{0}_{p \times 1}, \mathbf{I}_p)$ . The Bayes rule classifier

$$W_i^{QDA}(\mathbf{x}) = -\frac{1}{2} \log |\boldsymbol{\Sigma}_i| - \frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_i)^T \boldsymbol{\Sigma}_i^{-1} (\mathbf{x} - \boldsymbol{\mu}_i) + \log \pi_i, \quad i \in \{0, 1\}. \quad (1)$$

The classification rule is given by

$$\begin{cases} \mathbf{x} \in \mathcal{C}_0 \\ \mathbf{x} \in \mathcal{C}_1 \end{cases} \quad \text{if } W_0^{QDA}(\mathbf{x}) > W_1^{QDA}(\mathbf{x}) \quad \text{otherwise.} \quad (2)$$

### Training

Statistics estimation for  $i \in \{0, 1\}$

$$\begin{aligned} \bar{\mathbf{x}}_i &= \frac{1}{n_i} \sum_{\mathbf{l} \in \mathcal{T}_i} \mathbf{x}_l, \\ \hat{\boldsymbol{\Sigma}}_i &= \frac{1}{n_i - 1} \sum_{\mathbf{l} \in \mathcal{T}_i} (\mathbf{x}_l - \bar{\mathbf{x}}_i) (\mathbf{x}_l - \bar{\mathbf{x}}_i)^T. \end{aligned}$$

The empirical discriminant analysis score becomes

$$\widehat{W}_i^{QDA}(\mathbf{x}) = -\frac{1}{2} \log |\hat{\boldsymbol{\Sigma}}_i| - \frac{1}{2} (\mathbf{x} - \bar{\mathbf{x}}_i)^T \hat{\boldsymbol{\Sigma}}_i^{-1} (\mathbf{x} - \bar{\mathbf{x}}_i) + \log \pi_i. \quad (3)$$

If  $n_i < p$ , we use a regularized estimate for the inverse covariance matrix

$$\mathbf{H}_i = \left( \mathbf{I}_p + \gamma \hat{\boldsymbol{\Sigma}}_i \right)^{-1}, \quad (4)$$

where  $\gamma > 0$  is a regularizer.

$$\widehat{W}_i^{RQDA} = \frac{1}{2} \log |\mathbf{H}_i| - \frac{1}{2} (\mathbf{x} - \bar{\mathbf{x}}_i)^T \mathbf{H}_i (\mathbf{x} - \bar{\mathbf{x}}_i) + \log \pi_i. \quad (5)$$

The conditional classification error can easily be shown to write as

$$\epsilon_i = \mathbb{P} [\boldsymbol{\omega}^T \mathbf{B}_i \boldsymbol{\omega} + 2\boldsymbol{\omega}^T \mathbf{y}_i < \xi_i], \quad \text{where } \boldsymbol{\omega} \sim \mathcal{N}(\mathbf{0}_p, \mathbf{I}_p), \quad (6)$$

$$\begin{aligned} \mathbf{B}_i &= \boldsymbol{\Sigma}_i^{1/2} (\mathbf{H}_1 - \mathbf{H}_0) \boldsymbol{\Sigma}_i^{1/2}, \\ \mathbf{y}_i &= \boldsymbol{\Sigma}_i^{1/2} [\mathbf{H}_1 (\boldsymbol{\mu}_1 - \bar{\mathbf{x}}_1) - \mathbf{H}_0 (\boldsymbol{\mu}_1 - \bar{\mathbf{x}}_0)], \\ \xi_i &= -\log \left( \frac{|\mathbf{H}_0|}{|\mathbf{H}_1|} \right) + (\boldsymbol{\mu}_1 - \bar{\mathbf{x}}_0)^T \mathbf{H}_0 (\boldsymbol{\mu}_1 - \bar{\mathbf{x}}_0) - (\boldsymbol{\mu}_1 - \bar{\mathbf{x}}_1)^T \mathbf{H}_1 (\boldsymbol{\mu}_1 - \bar{\mathbf{x}}_1) + 2 \log \frac{\pi_1}{\pi_0}. \end{aligned}$$

## Important Results from RMT

### Technical assumptions

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  $O(1)$ .

### Recent RMT tools (Hachem et al, 2008)

For any  $\|\mathbf{A}\| = O(1)$  and under assumptions 1-4:

$$\frac{1}{p} \text{tr} \mathbf{A} \mathbf{H}_i - \frac{1}{p} \text{tr} \mathbf{A} \left( \mathbf{I}_p + \frac{\gamma}{1 + \gamma \delta_i} \boldsymbol{\Sigma}_i \right)^{-1} \xrightarrow[p \rightarrow \infty]{\text{a.s.}} 0, \quad \text{where } \delta_i = \frac{1}{n_i} \text{tr} \boldsymbol{\Sigma}_i \left( \mathbf{I}_p + \frac{\gamma}{1 + \gamma \delta_i} \boldsymbol{\Sigma}_i \right)^{-1}.$$

## Asymptotic analysis using RMT tools

Define

$$\mathbf{T}_i = \left( \mathbf{I}_p + \frac{\gamma}{1 + \gamma \delta_i} \boldsymbol{\Sigma}_i \right)^{-1} \quad \text{and the scalars } \phi_i \text{ and } \tilde{\phi}_i \text{ as } \phi_i = \frac{1}{n_i} \text{tr} \boldsymbol{\Sigma}_i^2 \mathbf{T}_i^2, \quad \tilde{\phi}_i = \frac{1}{(1 + \gamma \delta_i)^2}.$$

Let  $\boldsymbol{\mu} = \boldsymbol{\mu}_0 - \boldsymbol{\mu}_1$ , and set  $\bar{\xi}_i$ ,  $\bar{b}_i$  and  $\bar{B}_i$  to

$$\bar{\xi}_i \triangleq \frac{1}{\sqrt{p}} \left[ -\log \frac{|\mathbf{T}_0|}{|\mathbf{T}_1|} + \log \frac{(1 + \gamma \delta_0)^{n_0}}{(1 + \gamma \delta_1)^{n_1}} + \gamma \left( \frac{n_1 \delta_1}{1 + \gamma \delta_1} - \frac{n_0 \delta_0}{1 + \gamma \delta_0} \right) + (-1)^{i+1} \boldsymbol{\mu}^T \mathbf{T}_{1-i} \boldsymbol{\mu} \right].$$

$$\bar{b}_i \triangleq \frac{1}{\sqrt{p}} \text{tr} \boldsymbol{\Sigma}_i (\mathbf{T}_1 - \mathbf{T}_0).$$

$$\bar{B}_i \triangleq c \left[ \frac{\phi_0}{1 - \gamma^2 \phi_0 \tilde{\phi}_0} + \frac{\phi_1}{1 - \gamma^2 \phi_1 \tilde{\phi}_1} \right] - \frac{2}{p} \text{tr} \boldsymbol{\Sigma}_i \mathbf{T}_1 \boldsymbol{\Sigma}_i \mathbf{T}_0.$$

### Asymptotic Conditional Probability of Misclassification

Under assumptions 1-4, the following convergence holds for  $i \in \{0, 1\}$

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

### Special cases

- ▶  $\|\boldsymbol{\mu}_0 - \boldsymbol{\mu}_1\| = O(1)$ . In this case, the classification error rate would still converge to a non trivial limit but would not asymptotically depend on the difference  $\|\boldsymbol{\mu}_0 - \boldsymbol{\mu}_1\|$ .
- ▶  $\|\boldsymbol{\Sigma}_0 - \boldsymbol{\Sigma}_1\| = O(p^{-\frac{1}{2}-\alpha})$ ,  $\alpha > 0$ .

$$\epsilon - \Phi \left( -\frac{\boldsymbol{\mu}^T \mathbf{T} \boldsymbol{\mu}}{2\sqrt{p}} \sqrt{\frac{1 - \gamma^2 \phi \tilde{\phi}}{c \gamma^2 \phi^2 \tilde{\phi}}} \right) \xrightarrow{p} 0, \quad (7)$$

which means that RLDA is better than RQDA (Zollanvari and Dougherty, 2015).

## Experiments

### Synthetic data

$$\begin{aligned} \{\boldsymbol{\Sigma}_0\}_{i,j} &= 0.6^{|i-j|}, \quad \boldsymbol{\Sigma}_1 = \boldsymbol{\Sigma}_0 + 2 \begin{bmatrix} \mathbf{I}_k & \mathbf{0}_{k \times (p-k)} \\ \mathbf{0}_{(p-k) \times k} & \mathbf{0}_{(p-k) \times (p-k)} \end{bmatrix}, \quad k = \lfloor \sqrt{p} \rfloor, \\ \boldsymbol{\mu}_0 &= [\mathbf{1}, \mathbf{0}_{(p-1) \times 1}], \quad \boldsymbol{\mu}_1 = \boldsymbol{\mu}_0 + p^{-\frac{1}{4}} \mathbf{1}_{p \times 1}. \end{aligned}$$

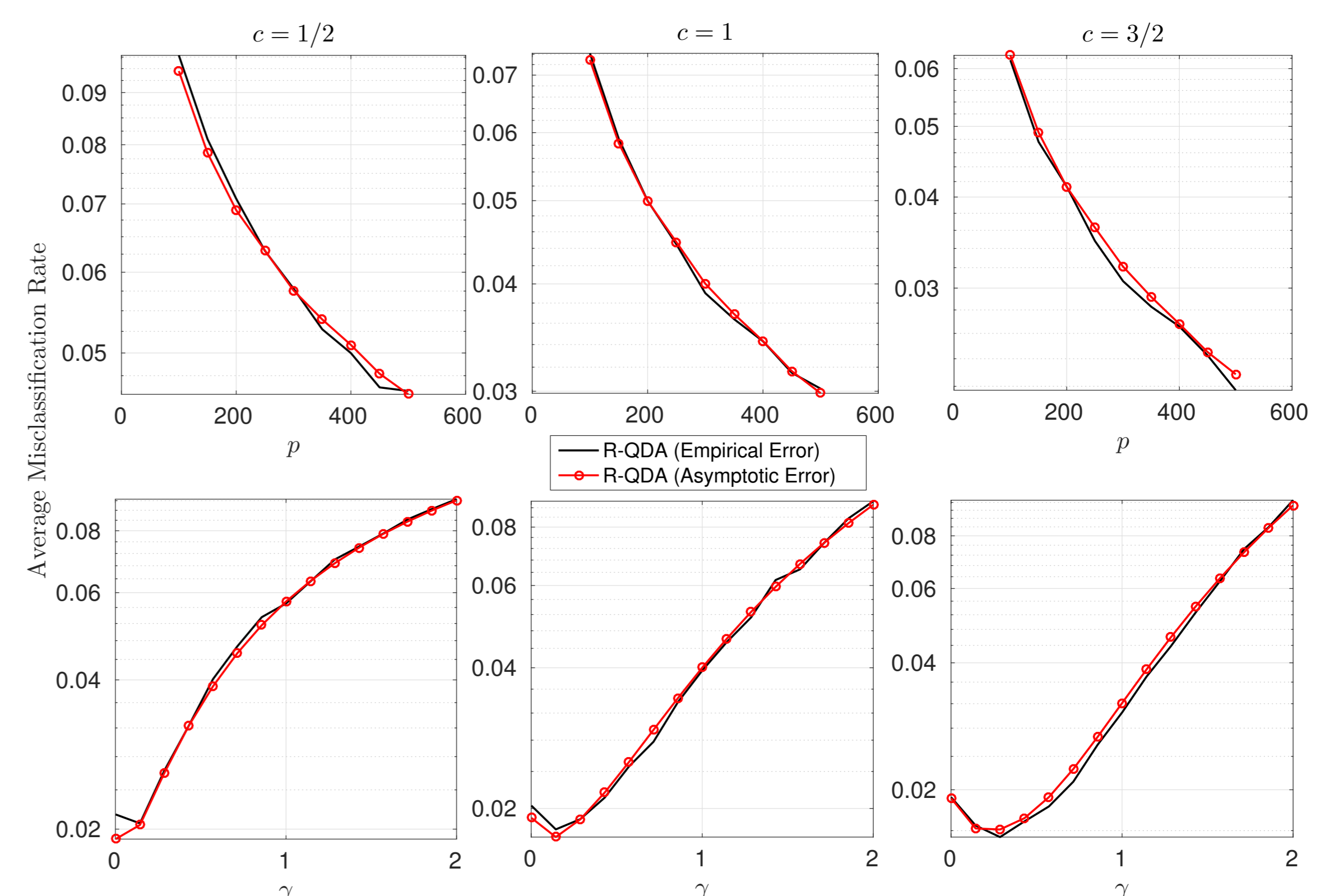


Figure:  $n_0 = n_1$  for both experiments with  $\gamma = 1$  in the first experiment and  $p = 300$  in the second experiment.

### Breast cancer data

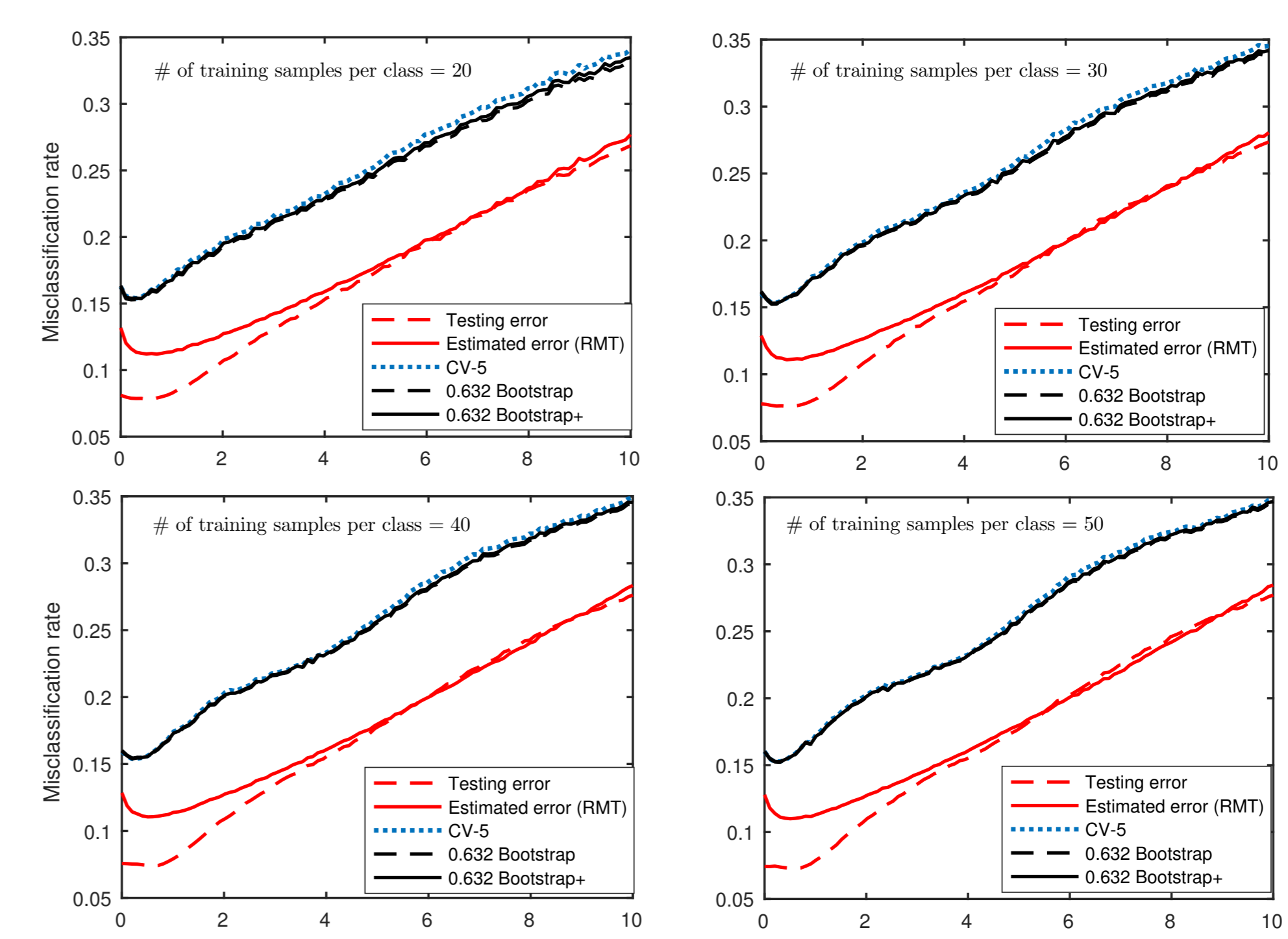


Figure: Average misclassification rate versus the regularization parameter  $\gamma$  with equal training size ( $n_0 = n_1$ ).