

# CS4220: Machine Learning

Yuhang Tian<sup>1</sup>

<sup>1</sup>Technology University of Delft

y.tian-13@student.tudelft.nl

## 1. Bayesian Classifier Theorem

### 1.1. Decision Boundary

$$p(x|w_1)p(w_1) \geq p(x|w_2)p(w_2) \quad (1)$$

### 1.2. Minimizing the Classification Error Probability (Bayes Error)

$$P_e = p(w_2) \int_{-\infty}^{x_0} p(x|w_2)dx + p(w_1) \int_{x_0}^{\infty} p(x|w_1)dx \quad (2)$$

where  $p(x_0|w_1)p(w_1) = p(x_0|w_2)p(w_2)$

### 1.3. Minimizing the Average Risk

Loss Matrix

$$L = \begin{bmatrix} 0 & \lambda_{12} \\ \lambda_{21} & 0 \end{bmatrix} \quad (3)$$

Risk Function

$$r = \lambda_{21}p(w_2) \int_{-\infty}^{x_0} p(x|w_2)dx + \lambda_{12}p(w_1) \int_{x_0}^{\infty} p(x|w_1)dx \quad (4)$$

### 1.4. Gaussian pdf in the l-dimensional space

$$p(x) = \frac{1}{(2\pi)^{l/2}|\Sigma|^{1/2}} \exp\left(-\frac{1}{2}(x - \mu)^T \Sigma^{-1}(x - \mu)\right) \quad (5)$$

where  $\Sigma = E[(x - \mu)(x - \mu)^T]$

### 1.5. Bayesian Classifier

discriminant functions

$$g_i(x) = \ln(p(x|w_i)p(w_i)) = \ln p(x|w_i) + \ln p(w_i) \quad (6)$$

Normally Distributed Classifier

$$g_i(x) = -\frac{1}{2}x^T \Sigma_i^{-1}x + \frac{1}{2}x^T \Sigma_i^{-1}\mu_i - \frac{1}{2}\mu_i^T \Sigma_i^{-1}\mu_i + \frac{1}{2}\mu_i^T \Sigma_i^{-1}x + \ln p(w_i) + c_i \quad (7)$$

where  $c_i = -(l/2)\ln 2\pi - (1/2)\ln(\det(\Sigma_i))$

- if l=2, correlation=0

$$g_i(x) = -\frac{1}{2\sigma_i^2}(x_1^2 + x_2^2) + \frac{1}{\sigma_i^2}(\mu_{i1}x_1 + \mu_{i2}x_2) - \frac{1}{2\sigma_i^2}(\mu_{i1}^2 + \mu_{i2}^2) + \ln p(w_i) + c_i \quad (8)$$

$g_i(x) - g_j(x) = 0$  are quadratics (i.e., ellipsoids, parabolas, hyperbolas, pairs of lines)

- if covariance matrix is the same in all classes

$$g_i(x) = w_i^T x + b \quad (9)$$

where  $w_i = \Sigma^{-1}\mu_i$  and  $b = \ln p(w_i) - \frac{1}{2}\mu_i^T \Sigma^{-1}\mu_i$

- if Diagonal covariance matrix with equal elements ( $\Sigma = \sigma^2 I$ )

$$g_i(x) = \frac{1}{\sigma^2} \mu_i^T x + b \quad (10)$$

## 2. ESTIMATION OF UNKNOWN PROBABILITY DENSITY FUNCTIONS

### 2.1. ML

we considered  $\theta$  as an unknown parameter.

$$\hat{\theta}_{ML} = \arg \max_{\theta} \prod_{k=1}^N p(x_k; \theta) \quad (11)$$

ML estimate of  $\sigma^2$

$$\hat{\sigma}_{ML}^2 = \frac{1}{N} \sum_{k=1}^N (x_k - \mu)^2 \quad (12)$$

ML estimate of  $\mu$

$$\hat{\mu}_{ML} = \frac{1}{N} \sum_{k=1}^N x_k \quad (13)$$

### 2.2. MAP

we considered  $\theta$  as a random vector.

$$p(\theta|X) = \frac{p(\theta)p(X|\theta)}{P(X)} \quad (14)$$

then,

$$\hat{\theta}_{MAP} : \frac{\partial}{\partial \theta} p(\theta|X) = 0 \text{ or } \frac{\partial}{\partial \theta} p(X|\theta)p(\theta) = 0 \quad (15)$$

### 2.3. Bayesian Inference

Given the set  $X$  of the  $N$  training vectors and the *a priori information* about the pdf  $p(\theta)$ , the goal is to compute the conditional pdf  $p(x|X)$ .

$$p(x|X) = \int p(x|\theta)p(\theta|X)d\theta \quad (16)$$

with

$$p(\theta|X) = \frac{p(X|\theta)p(\theta)}{p(X)} \quad (17)$$

$$p(X|\theta) = \prod_{k=1}^N p(x_k|\theta) \quad (18)$$

### 3. Normal-based Classifier: Quadratic Discriminant, Linear Discriminant and Nearest Mean

Let's assume that we have two classes:

#### 3.1. Quadratic Discriminant

by eq.7, the quadratic classifier,

$$f(x) = x^T W x + w^T + w_0 \quad (19)$$

with

$$W = \frac{1}{2}(\Sigma_2^{-1} - \Sigma_1^{-1}) \quad (20)$$

$$w = \mu_1^T \Sigma_1^{-1} - \mu_2^T \Sigma_2^{-1} \quad (21)$$

$$w_0 = -\frac{1}{2} \ln(\det(\Sigma_1)) - \frac{1}{2} \mu_1^T \Sigma_1^{-1} \mu_1 + \ln p(y_1) + \frac{1}{2} \ln(\det(\Sigma_2)) + \frac{1}{2} \mu_2^T \Sigma_2^{-1} \mu_2 - \ln p(y_2) \quad (22)$$

(i.e., ellipsoids, parabolas, hyperbolas, pairs of lines)

#### 3.2. Linear Discriminant

by eq.9, the linear classifier,

$$f(x) = w^T x + w_0 \quad (23)$$

with

$$w = \Sigma^{-1}(\mu_2 - \mu_1) \quad (24)$$

$$w_0 = \frac{1}{2} \mu_2^T \Sigma^{-1} \mu_2 - \frac{1}{2} \mu_1^T \Sigma^{-1} \mu_1 + \ln \frac{p(y_1)}{p(y_2)} \quad (25)$$

Therefore, the Linear Discriminant has a **strong assumption** on that  $\Sigma_1 = \Sigma_2$

#### 3.3. Nearest Mean Classifier

by eq.10, the nearest mean classifier,

$$f(x) = w^T x + w_0 = \sigma^2(g_1(x) - g_2(x)) \quad (26)$$

with

$$w = \mu_2 - \mu_1 \quad (27)$$

$$w_0 = \frac{1}{2} \mu_2^T \mu_2 - \frac{1}{2} \mu_1^T \mu_1 + \sigma^2 \ln \frac{p(y_1)}{p(y_2)} \quad (28)$$

Therefore, the Nearest Mean has a **strong assumption** on mutually uncorrelated and of the same variance ( $\Sigma_1 = \Sigma_2 = \sigma^2 I$ )

## 4. More Parametric Classifiers

### 4.1. Logistic Classifier

$$\begin{cases} p(y_1|x) = \frac{1}{e^{-(w^T x + w_0)} + 1} \\ p(y_2|x) = \frac{1}{e^{(w^T x + w_0)} + 1} \end{cases} \quad (29)$$

Maximize Log Likelihood

$$\ln p(y|x) = \sum_{i=1}^N \ln \left( \frac{1}{e^{-y_i(w^T x_i + w_0)} + 1} \right) \quad (30)$$

### 4.2. Fisher Classifier

$$y_i = w^T x_i \begin{cases} \geq 0 & \text{if class 1} \\ < 0 & \text{if class 2} \end{cases} \quad (31)$$

Minimize Square Loss

$$L(w) = \sum_{i=1}^N (w^T x_i - y_i)^2 \quad (32)$$

### 4.3. The Perception

Forward&backward propagation and update  $w$

$$w \leftarrow w + \eta y x \quad (33)$$

## 5. Non-parametric Classification

In the above sections, all the classifiers are based on **Parametric Classification** method, more precisely, based on normal distribution and Bayes Theorem. In this section, it will mainly demonstrate two non-parametric classifiers - Parzen Classifier and Nearest Neighbour Classifier (Both methods are sensitive to the scaling of the features).

### 5.1. Parametric vs. Non-parametric

- Parametric: Assumptions can greatly simplify the learning process, but can also limit what can be learned. Algorithms that simplify the function to a known form are called parametric machine learning algorithms.
- Non-parametric: Algorithms that do not make strong assumptions about the form of the mapping function are called non-parametric machine learning algorithms. By not making assumptions, they are free to learn any functional form from the training data.

### 5.2. Histogram method

it's easy, so I just skip it.

### 5.3. Parzen Density Estimation

- Kernel

$$f(x) = \begin{cases} 0 & \text{if } |x| > h \\ \frac{1}{V} & \text{if } |x| \leq h \end{cases} \quad (34)$$

- Parzen Classifier

$$p(z|h) = \frac{1}{n} \sum_{i=1}^n K(\|z - x_i\|, h) \quad (35)$$

- Parzen plugs in the Gaussian density:

$$p(x|w_i) = \frac{1}{n_i} \sum_{j=1}^{n_i} N(x|x_j^{(i)}, hI) \quad (36)$$

### 5.4. Nearest Neighbour Classification

$$p(x) = p(x|w_m)p(w_m) = \frac{k_m}{n_m V_k} \cdot \frac{n_m}{n} \quad (37)$$

where  $V_k$  is the volume of the sphere centered at  $x$  with radius  $r$  (the distance to the  $k$ -th nearest neighbor)

## 6. More Non-parametric Classifiers

GitHub Markdown

### 6.1. SVM

By putting some constraints on the linear classifier, the VC dimension can be reduced. Why do that? Ans: When  $h$  is small, the true error is close to the apparent error

$$\begin{cases} w^T x_i + b \geq +1 & \text{for } y_i = +1 \\ w^T x_i + b \leq -1 & \text{for } y_i = -1 \end{cases} \quad (38)$$

Core Idea of SVM: Find the decision boundary, while maximize the margin

1. The above two equation can be merged into one

$$y_i(w^T x_i + b) - 1 \geq 0 \quad (39)$$

2. The distance between the two boundaries

$$\text{maximize } \frac{2}{\|w\|} \rightarrow \text{minimize } \frac{1}{2} \|w\|^2 \quad (40)$$

3. by Lagrange Multiplier

$$L = \frac{1}{2} \|w\|^2 - \sum \alpha_i [y_i(w^T x_i + b) - 1] \quad (41)$$

4. by  $\frac{\partial L}{\partial w} = 0$

$$w = \sum \alpha_i y_i x_i \quad (42)$$

5. by  $\frac{\partial L}{\partial b} = 0$

$$\sum \alpha_i y_i = 0 \quad (43)$$

6. put eq.42 back to eq.41

$$L = \sum \alpha_i - \frac{1}{2} \sum_i \sum_j \alpha_i \alpha_j y_i y_j \mathbf{x}_i \cdot \mathbf{x}_j \quad (44)$$

7. put eq.42 back to decision rule

$$\begin{cases} \sum \alpha_i y_i \mathbf{x}_i \cdot \mathbf{u}_i + b - 1 \geq 0 & \text{Then, +} \\ \sum \alpha_i y_i \mathbf{x}_i \cdot \mathbf{u}_i + b - 1 \leq 0 & \text{Then, -} \end{cases} \quad (45)$$

8. kernelize

$$K(x_i, x_j) = \phi(x_i) \cdot \phi(x_j) \quad (46)$$

## 6.2. Agglomerative Hierarchical Clustering

1. Determine distances between all clusters
  - Two nearest objects in the clusters: single linkage
  - Two most remote objects in the clusters: complete linkage
  - Cluster centers: average linkage
2. Merge clusters that are closes
3. IF #clusters>1 THEN GOTO 1
  - Dendrogram: Cut at "largest jump" → Clustering
  - Fusion Graph: Cut at "largest drop" → Clustering

## 7. Regression

### 7.1. Intuitively Understanding

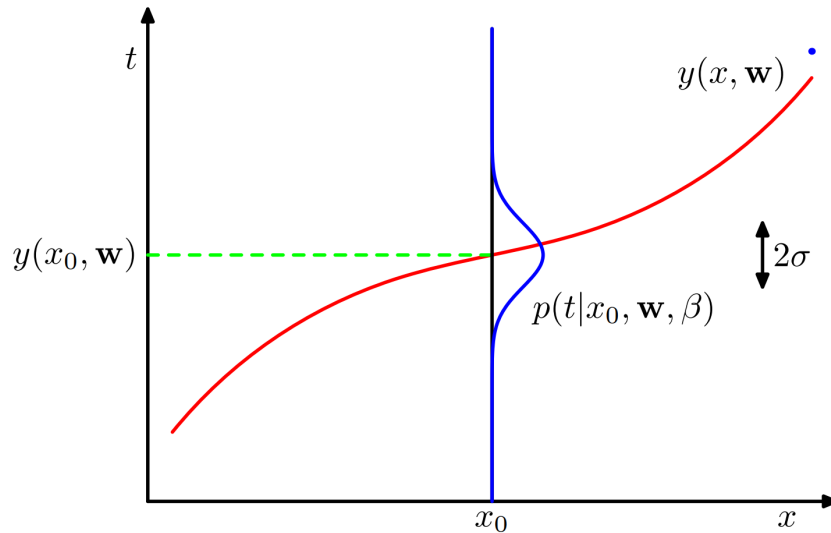


Figure 1. Regression Curve

Suppose that there is a fixed  $x_0$  and a bunch of choices of  $\theta$  (here are  $\mu$  and  $\sigma^2(\beta^{-1})$ ). If we expect the predicted value  $\hat{y}(x_0, w)$  to locate on the true value  $t_0$ , the probability of this occurrence  $p(t_0|x_0, w, \beta)$  should be maximized.

## 7.2. Maximum Likelihood Regression

$$p(t|x, w, \beta) = \prod_{n=1}^N N(t_n|y(x_n, w), \beta^{-1}) \quad (47)$$

Log Likelihood function

$$\ln p(t|x, w, \beta) = -\frac{\beta}{2} \sum_{n=1}^N y(x_n, w) - t_n^2 + \frac{N}{2} \ln \beta - \frac{N}{2} \ln(2\pi) \quad (48)$$

To solve

$$\frac{\partial}{\partial w} \ln p(t|x_0, w, \beta) = 0 \quad (49)$$

and

$$\frac{\partial}{\partial \beta^{-1}} \ln p(t|x_0, w, \beta) = 0 \quad (50)$$

If we use linear regression, the solutions of eq.49 and eq.50

$$w_{ML} = (X^T X)^{-1} X^T Y \quad \beta_{ML}^{-1} = \frac{1}{N} (w_{ML}^T X - Y) \quad (51)$$

## 7.3. Max a Posterior Regression

Suppose that we have some knowledge about  $w \sim N(0, \alpha I)$

$$w_{MAP} : \left( \prod_{i=1}^N p(y_i|w^T x_i, \sigma^2) \right) p(w|0, \alpha I) \quad (52)$$

$$\frac{\partial}{\partial w} \left( \prod_{i=1}^N p(y_i|w^T x_i, \sigma^2) \right) p(w|0, \alpha I) = 0 \rightarrow w_{MAP} = (X^T X + \frac{\sigma^2}{\alpha} I)^{-1} X^T Y \quad (53)$$

## 8. Regularization

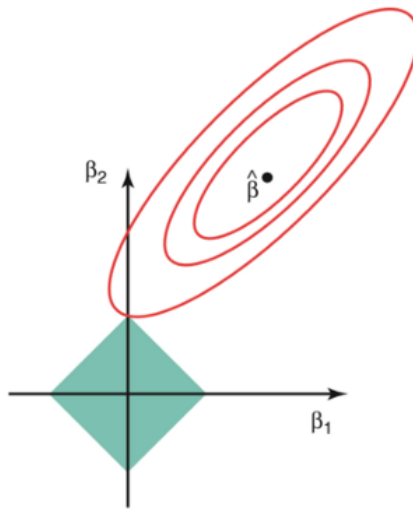
### 8.1. Keep Eigenvalues Away From Zero

Add identity to  $X X^T$

$$\hat{w} = (X^T X + \lambda I)^{-1} X^T Y \quad (54)$$

### 8.2. LASSO, L1 Norm

$$\min_{\beta} \frac{1}{N} \sum_{i=1}^n (y_i - w^T x_i)^2 + \lambda \|w\| \quad (55)$$

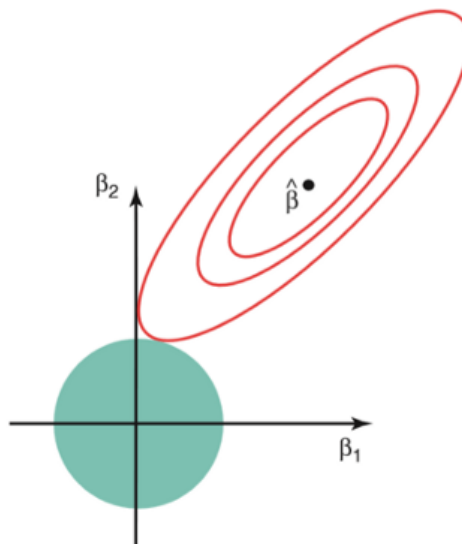


**Figura 2. LASSO/L1 Regularization**

$$\lambda \propto \frac{1}{\tau}$$

### 8.3. Ridge, L2 Norm

$$\min_{\beta} \frac{1}{N} \sum_{i=1}^n (y_i - w^T x_i)^2 + \lambda \|w\|^2 \quad (56)$$



**Figura 3. Ridge/L2 Regularization**

$$\lambda \propto \frac{1}{\tau}$$

### 8.4. L1 vs. L2

- L1 is for feature selection



- L2 is for avoiding overfitting

[Read More](#)

## 9. Curves

### 9.1. Bias-Variance Decomposition

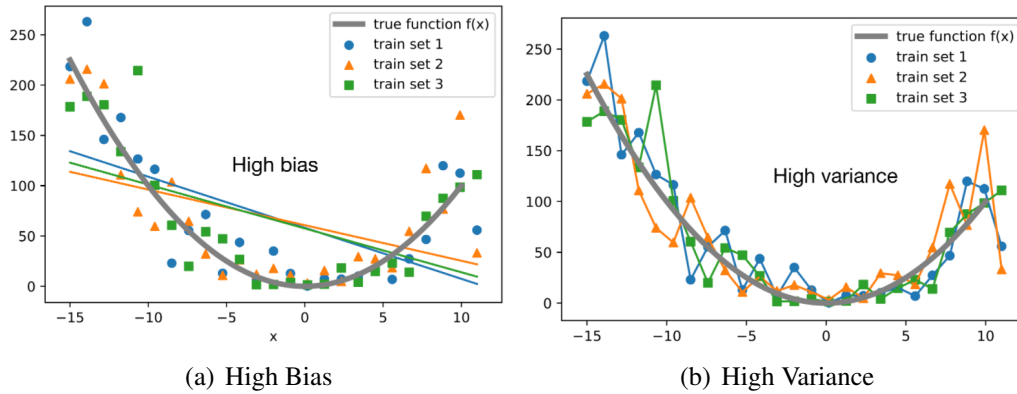


Figure 4. High Bias vs. High Variance

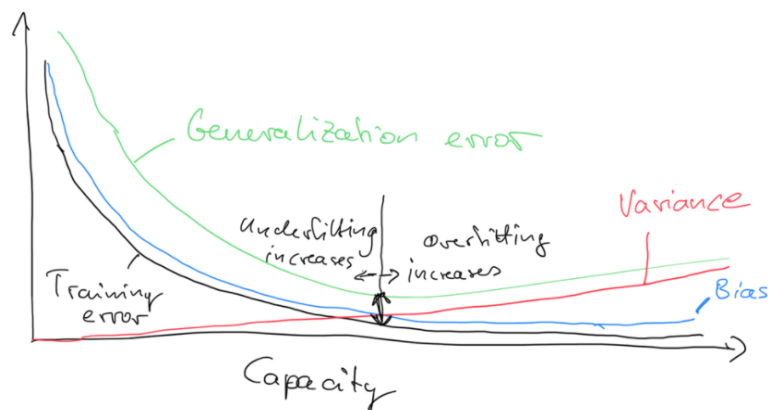


Figure 5. Decomposition Of Loss

### 9.2. Cross Validation Curve

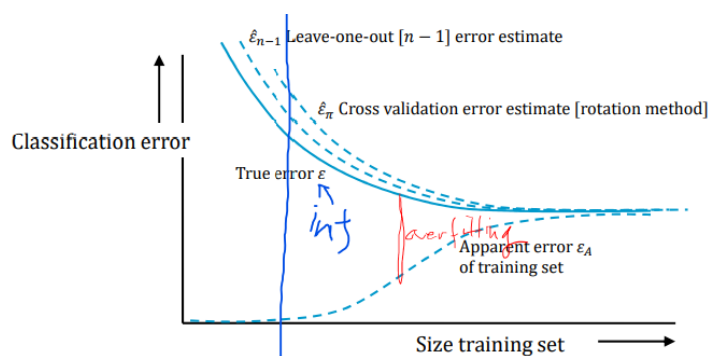


Figure 6. Cross Validation Curve

### 9.3. Learning Curve

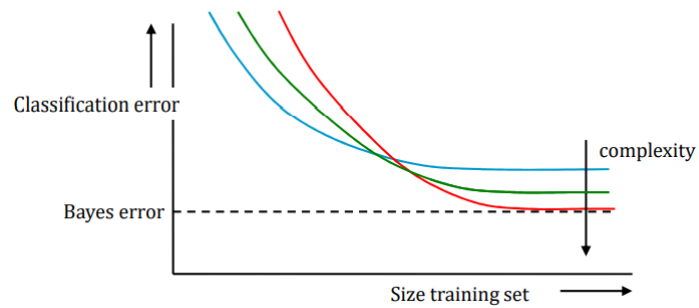


Figura 7. Learning Curve

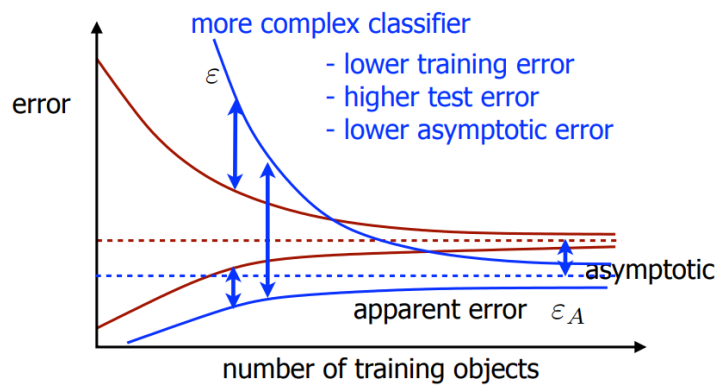


Figura 8. Learning Curve 2

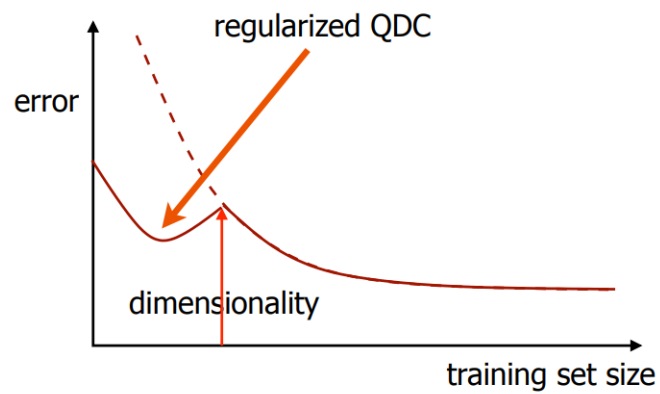
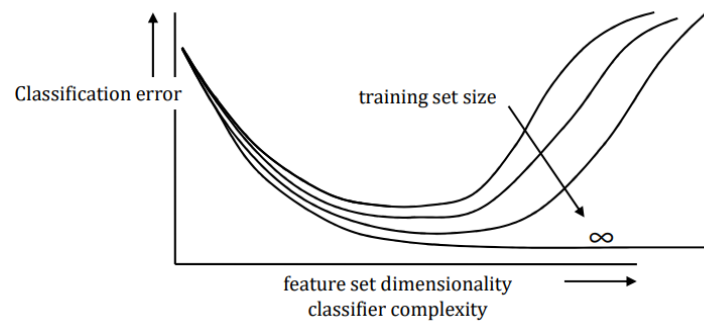
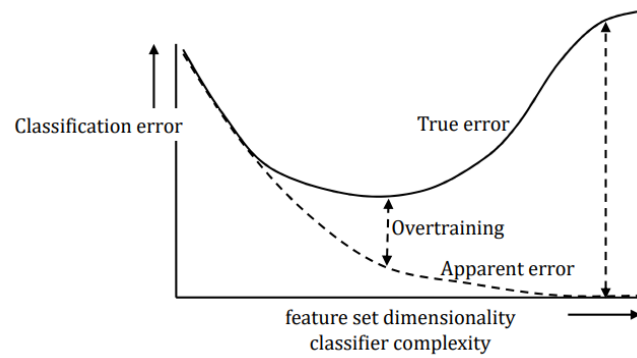


Figura 9. Regularized QDC

## 9.4. Feature Curve



**Figure 10. Feature Curve**



**Figure 11. Curse of Dimensionality**