

INF8953 CE - Fall 2020

Machine Learning

- Sarath Chandar

5. Classification



5. Classification

Consider a k -class classification problem.

$$c \in \{c_1, \dots, c_k\}$$

Two stages in classification :-

- Inference stage : Use training data to learn a model for $p(c_k|x)$
- Decision stage : Use posterior probabilities to make optimal class assignments.

$$\hat{c}(n) = \underset{c \in C}{\operatorname{argmin}} \sum_{k=1}^K L(c_k, c) \Pr(c_k | x=n)$$

↑
 Price for classifying an obs.
 belonging to class C_k to class c .

For 0/1 loss fn:

$$\hat{c}(n) = \underset{c \in C}{\operatorname{argmin}} [1 - \Pr(c | x=n)]$$

$$= \underset{c \in C}{\operatorname{argmax}} \Pr(c | x=n)$$

↳ Bayes classifier.

3 approaches to solving classification problems:-

① First solve the inference problem of determining the class-conditional densities $P(x|C_k)$ for each class C_k individually.

→ Also infer prior class probabilities $P(C_k)$

→ Then use Bayes theorem:

$$P(C_k|x) = \frac{P(x|C_k) P(C_k)}{P(x)}$$

Posterior class probabilities.

where $P(x) = \sum_k P(x|C_k) P(C_k)$

(or) we can also learn the joint distribution

$P(x, C_k)$ and then normalize to obtain

$$P(C_k|x).$$

Once you find $P(C_k|x)$, use decision theory

to determine the class.

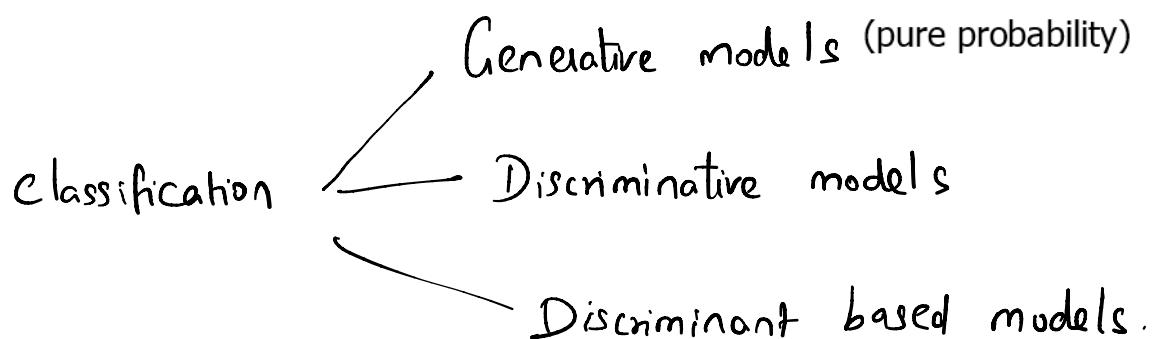
— This is known as "generative model".
↓
we can generate synthetic
data points in the input space
by using the learnt distribution.

② Directly model $P(C_k|x)$ and then use decision theory to determine the class.

Approaches that model $P(C_k|x)$ directly
are called "discriminative models".

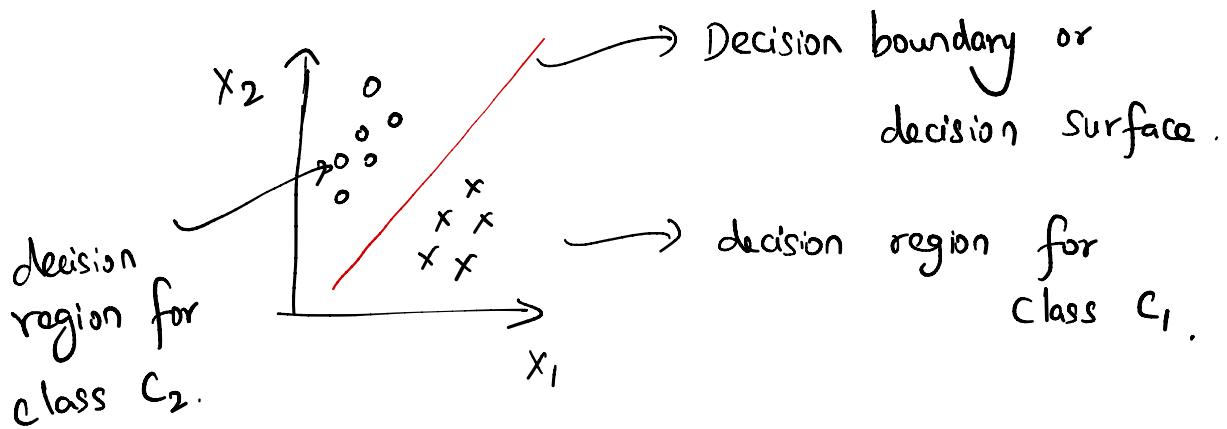
③ Find a function $f(x)$ — called discriminant function —
which maps x directly to class label.

- There is no probability here.
 - discriminant based models.
-



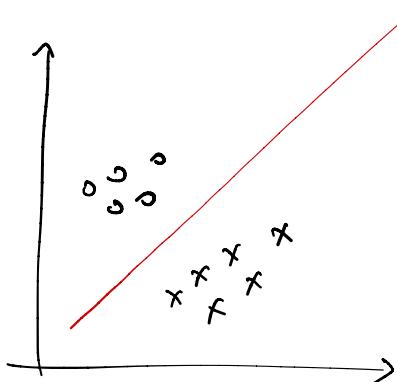
Linear classification :-

Consider a 2-class problem : $x - C_1$
 $o - C_2$

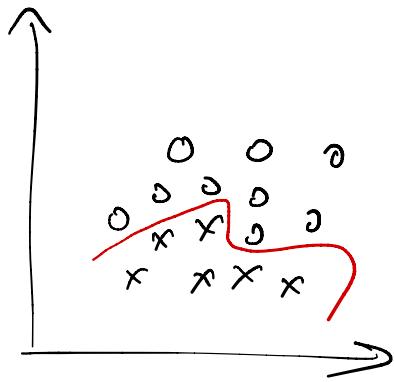


Linear model — decision surfaces are linear functions of the input vector x and hence are defined by ($D-1$) dimensional hyperplanes with the D -dimensional input space.

Linearly separable problems :-



linearly separable



Not linearly separable.

Target Variable for classification :-

For regression, target variable $t \in \mathbb{R}$

For classification, target class $c \in \{c_1, \dots, c_k\}$

2. class problem : (c_1, c_2) Label encoding

$t = 1$ class c_1

$t = 0$ class c_2

$t \in \{0, 1\}$

t - can be interpreted as the probability
that the class is c_1 .

Probability takes extreme values of 0 and 1.

$k > 2$ classes :-

① Can we consider $t \in \{1, 2, \dots, k\}$?

Problem: this representation creates artificial distance
between classes. c_1 is close to c_2 than c_5 .

② 1-of-k Coding scheme:- one hot encoding

for $k=5$ classes,

$$\text{class 1 : } t = (1, 0, 0, 0, 0)^T$$

$$\text{class 2 : } t = (0, 1, 0, 0, 0)^T$$

and so on.

Each t_k can be interpreted as the probability
that the class is C_k .

Discriminant functions:-

- takes an input vector x and assigns it to
one of k classes, denoted as C_k .
- linear discriminants \rightarrow decision surfaces are hyperplanes.

2-class scenario:-

Simple linear discriminant function:-

$$y(n) = \omega^T x + \omega_0$$

↑
weight vector ↑ bias

if $y(x) \geq 0$ then C_1

if $y(x) < 0$ then C_2

decision boundary : $y(x) = 0$

w_0 - can be considered as a threshold.

For example, if $w_0 = -7$, then $w^T x$ has to be at least $+7$ to be classified as C_1 .

Consider two points x_A and x_B both of which lie on the decision surface.

$$y(x_A) = 0$$

$$w^T x_A + w_0 = 0$$

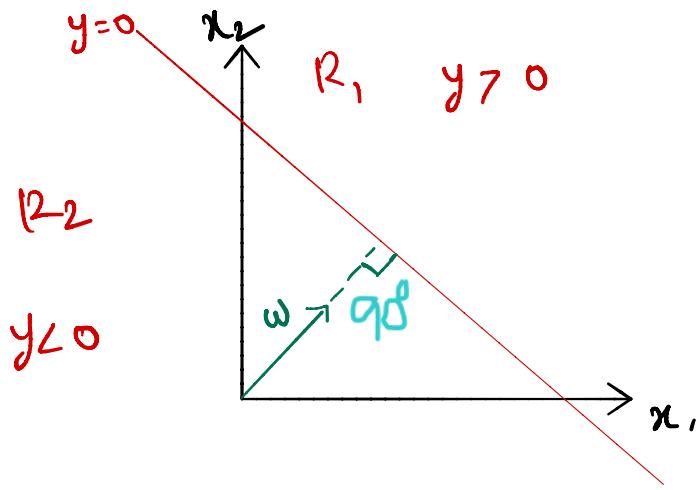
$$y(x_B) = 0$$

$$w^T x_B + w_0 = 0$$

$$w^T (x_A - x_B) = 0$$

\Rightarrow w vector is orthogonal to every vector lying within the decision surface.

\Rightarrow w determines the orientation of the decision surface.



to find the shortest distance to the decision line, ie which is orthogonal is the shortest

We know that w is orthogonal to the decision Surface.

Any point x' on this decision surface that is closest to the origin can be represented as

$$x' = \alpha w \quad \text{for some scalar } \alpha.$$

x' is on the decision surface.

$$\hookrightarrow w^T x' + w_0 = 0$$

$$\alpha w^T w + w_0 = 0$$

$$\alpha = \frac{-w_0}{\|w\|^2}$$

The distance from x' to the origin is

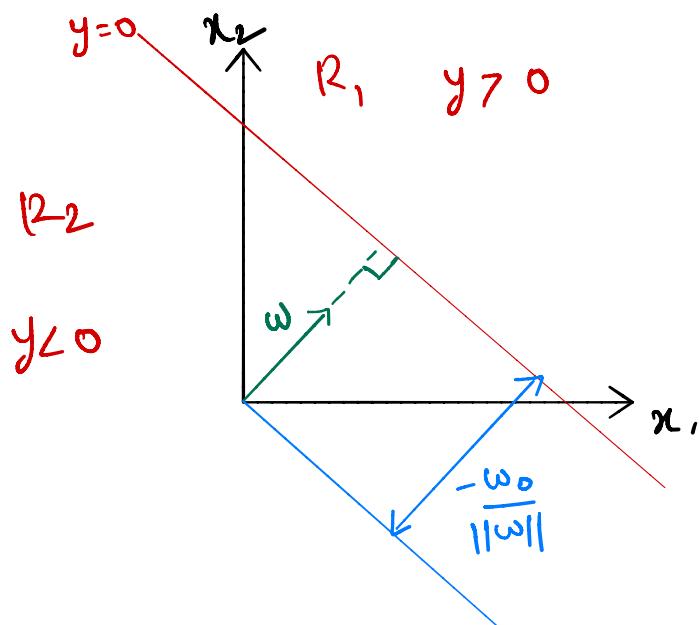
$$\begin{aligned} \|x'\| &= \|\alpha w\| \\ &= \alpha \|w\| \end{aligned}$$

$$= -\frac{w_0}{\|w\|^2} \|w\|$$

$$= -\frac{w_0}{\|w\|}$$

w_0 determines the location of the decision

function.



Value of $y(x)$ gives a signed measure of the perpendicular distance r of the point x from the decision surface.

Consider an arbitrary point x and let x_{\perp} be its orthogonal projection onto the decision surface, so that

$$x = x_{\perp} + \gamma \frac{w_0}{\|w\|}$$

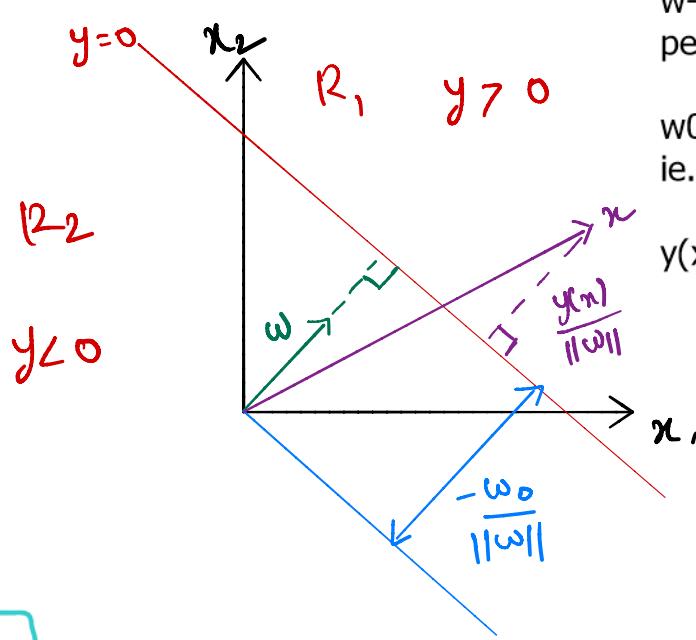
$$\omega^T x = \omega^T x_{\perp} + r \frac{\omega^T \omega}{\|\omega\|}$$

$$\underline{\omega^T x + w_0} = \underline{\omega^T x_{\perp}} + w_0 + r \cdot \frac{\omega^T \omega}{\|\omega\|}$$

$$y(x) = y(x_{\perp}) + r \frac{\omega^T \omega}{\|\omega\|}$$

$$y(x) = 0 + r \cdot \|\omega\|$$

$$r = \frac{y(x)}{\|\omega\|}$$



w- decides the orientation ie perpendicular

w0--> how far it is from the origin ie. displacement

y(x)--> sign, to determine the class

Summary:-

- ① Decision surface is perpendicular to ω .
- ② Displacement of decision surface is controlled by the bias parameter.

③ Signed orthogonal distance of a general point x from the decision surface

is given by $\frac{y(n)}{\|\omega\|}$

Magnitude--> tells how far from decision surface

Sign--> tells what is the class

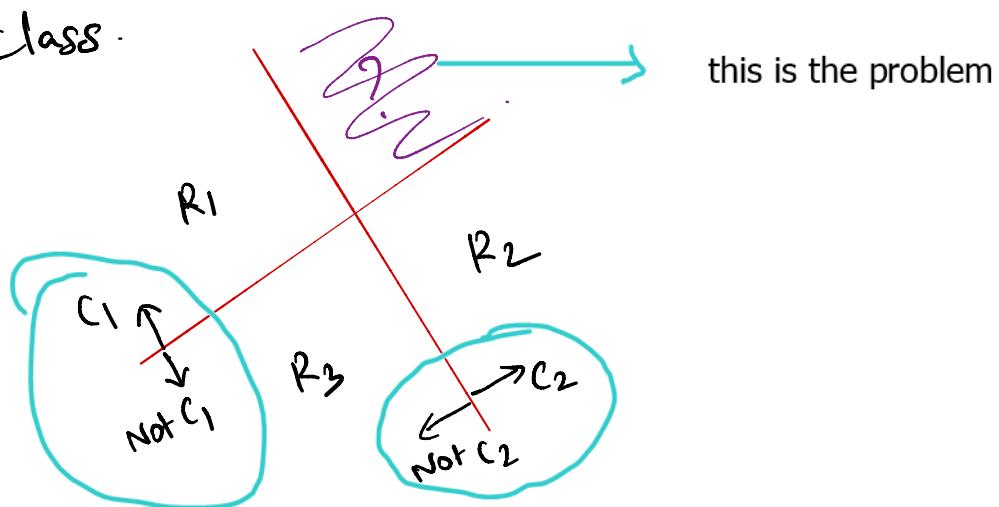
Multiple classes :-

① One - versus - the - rest classifier :-

$k-1$ classifiers each of which

Solves a two-class problem of separating points in class C_k from points not in that

class.

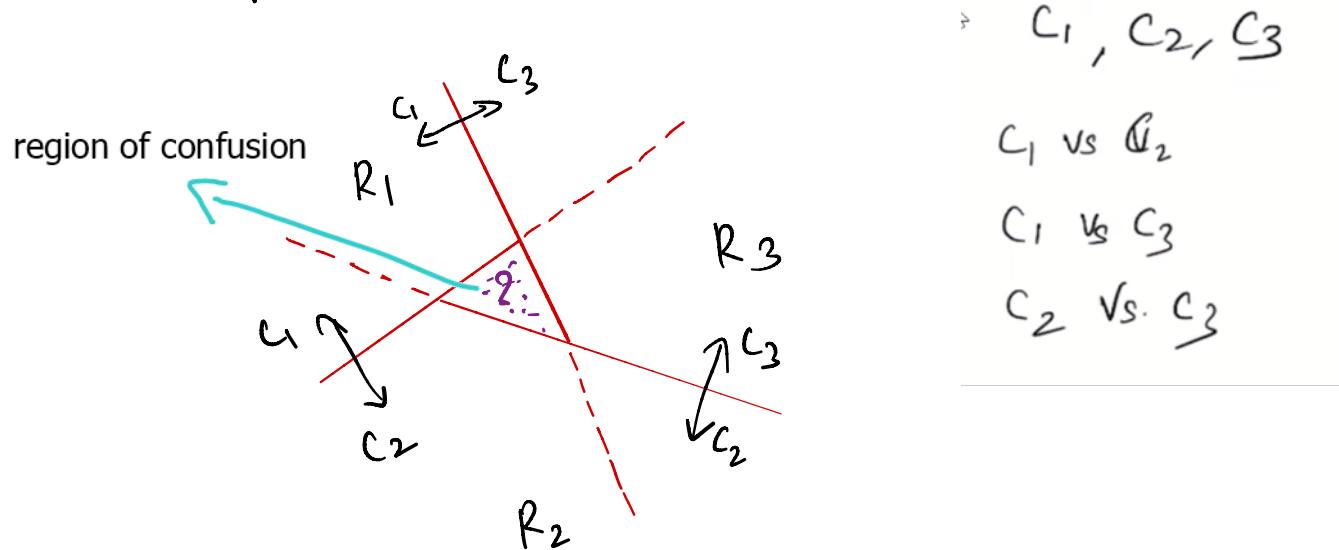


② One - versus - one classifier:-

Consider $k(k-1)/2$ binary discriminant

functions, one for every possible pair of classes.

- Majority voting among the discriminant functions for classification.



- ③ Consider a single k-class discriminant comprising k-linear functions of the form

$$y_k(x) = \omega_k^T x + \omega_{k0}$$

if $y_k(x) > y_j(x)$ for all $j \neq k$, then x belongs to C_k .

Decision boundary between C_k and C_j :

$$y_k(x) = y_j(x)$$

$$\omega_k x + \omega_{k0} = \omega_j x + \omega_{j0}$$

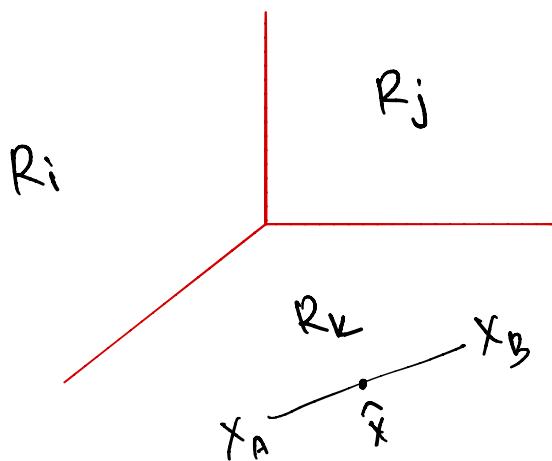
$$(\omega_k - \omega_j)^T x + (\omega_{k0} - \omega_{j0}) = 0$$

→ Similar to decision boundary for two classes.

→ So analogous geometrical properties apply.

Note: Decision regions of such a discriminant

are always singly connected and convex.



it is better than the previous 2 approaches because there is no confusion region

Consider $x_A, x_B \in R_k$

Any point \hat{x} that lies on the line connecting x_A and x_B can be expressed as

then \hat{x} lies in that same region as x_A and x_B

$$\hat{x} = \lambda x_A + (1-\lambda) x_B \quad 0 \leq \lambda \leq 1$$

$$y_k(\hat{x}) = \lambda y_k(x_A) + (1-\lambda) y_k(x_B)$$

[by linearity of
discriminant]

we know that $y_k(x_A) > y_j(x_A)$ and

$y_k(x_B) > y_j(x_B)$ for all $j \neq k$.

$$\Rightarrow y_k(\hat{x}) > y_j(\hat{x}) \text{ for all } j \neq k.$$

So $\hat{x} \in R_k$.

Thus R_k is singly connected and convex.

Least Squares for classification:-

- classification problem with 'k' classes.
- 1-of-k binary coding scheme for the target vector t .
- We know that least squares approximates the conditional expectation $E[b|x]$ of the target values given the input vector.
For the binary coding scheme, this conditional expectation is given by the vector of posterior probabilities. So it makes sense to use least squares.

- However, there is no guarantee for the values to be in the $(0,1)$ range!

Each class C_k is described by its own
linear model:

$$y_k(n) = \omega_k^T x + \omega_{k_0} \quad \text{for } k=1\dots K$$

Vector notation:-

$$y(x) = \tilde{w}^T \tilde{x}$$

\tilde{w} is a matrix whose k^{th} column comprises
the $D+1$ dimensional vector $\tilde{\omega}_k = (\omega_{k_0}, \omega_k^T)^T$.

\tilde{x} - is the augmented input vector $(1, x^T)^T$.

→ Consider training data set: $\{x^{(n)}, t^{(n)}\}_{n=1}^N$

→ define a matrix T whose n^{th} row is
the vector $t^{(n)T}$, together with a matrix
 \tilde{X} whose n^{th} row is $\tilde{x}^{(n)T}$

→ Sum of Square error fn:

$$E_D(\tilde{w}) = \frac{1}{2} \operatorname{Tr} \left\{ (\tilde{x} \tilde{w} - T)^T (\tilde{x} \tilde{w} - T) \right\}$$

$$\text{Solution :- } \tilde{w} = (\tilde{x}^T \tilde{x})^{-1} \tilde{x}^T T$$

$$= \tilde{x}^+ T$$

 \rightarrow Pseudo inverse.

discriminant function : $y(x) = \tilde{w}^T x = T^T (\tilde{x}^+)^T x$.

Note 1: If every target vector in the training set satisfies some linear constraint

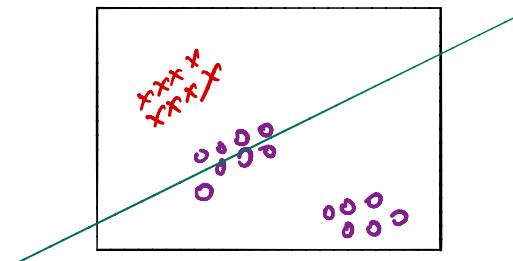
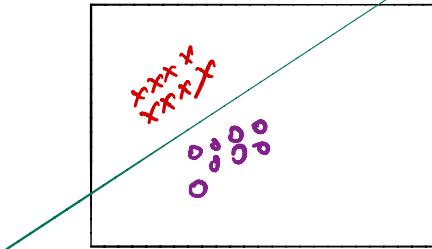
$$a^T t^{(n)} + b = 0$$

for some constants a, b , then the least squares model prediction for any value of x will satisfy the same constraint so that

$$a^T y(x) + b = 0$$

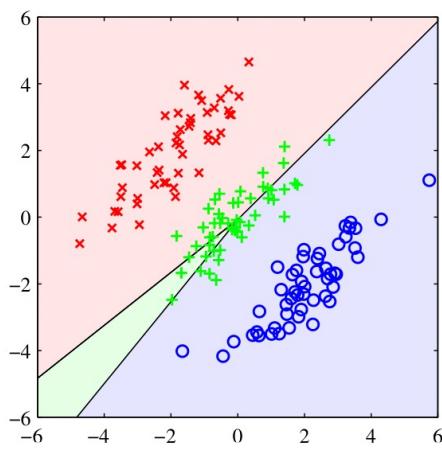
Since we use 1-0k-k encoding, model's predictions will sum to 1 for any value of x . But the values are not constrained to be in the range (0,1) ! ie there can be negative too.

Note 2: Least-squares solutions lack robustness to outliers.



The sum of Square error function penalizes predictions that are "too correct" in that they lie a long way on the correct side of the decision boundary.

Note 3: Doesn't work well in many cases.



only small region assigned to green class!

You should know!

1. 2 stages of classification
 - Inference
 - Decision.
2. 3 approaches to classification
 - Generative model
 - Discriminative model
 - Discriminant-based model.
3. Linear classification, linear decision boundary,
linear model, linearly separable problems.
4. 1-of-k Coding scheme / 1-hot vector
5. Discriminant functions.
6. Geometric interpretation of $y(n) = \mathbf{w}n + w_0$.
7. Multiple classes
 - 1 vs. rest classifier
 - 1 vs. 1 classifier
 - k-class discriminants