

MAP 534

Introduction to machine learning

Linear models for classification I

Alain Durmus

A quick clarification regarding homework and Labs

- No libraries can be used for the homeworks.
- Extension for the first one to Monday.
- Students which would have submitted their homework in due date will have a bonus for the next one +2 points.
- Correction for some exercises/labs will be provided.

Introduction

Bayes classifier

Deterministic discriminative learning

Generative probabilistic models

Naive Bayes

Discriminant analysis (linear and quadratic)

Introduction to classification: setting

- We consider a supervised setting.
- Decide how to encode inputs and outputs: this defines the input space X , and the output space Y .
- Here we consider specifically the classification problem: Y is a finite set, $Y = \{1, \dots, K\}$, in most of this lecture even $Y = \{0, 1\}$ (or $\{-1, 1\}$). (1)
- In this lecture, we apply the three machine learning paradigms to address:

$$y_{\text{pred}} = \mathcal{C}(x_{\text{new}}), \quad (2)$$

and aim to quantify if possible the uncertainties of our predictions.

- Here \mathcal{C} is called a classifier.
- Recall that the three paradigms are:
 - deterministic discriminative learning;
 - probabilistic discriminative learning;
 - probabilistic generative learning;
- We will use these three paradigms to learn particular classifiers.
- A first question we deal with is the existence of an optimal classifier, also called Bayes classifier.

Introduction

Bayes classifier

Deterministic discriminative learning

Generative probabilistic models

Naive Bayes

Discriminant analysis (linear and quadratic)

An optimal classifier?

- Consider $Y = \{0, 1\}$ and $X = \mathbb{R}^d$.
- What would be an optimal classifier $\mathcal{C} : X \rightarrow \{0, 1\}$.
- To this end, we take a decision theoretical approach.
- Consider a probabilistic model for the data (X, Y) and denote by $(x, y) \mapsto p(x, y)$ the corresponding density.
- Namely,

$$\mathbb{P}(X \in A, Y = i) = \int_A p(x, i) dx, \quad (3)$$

where p_0, p_1 are probability densities on \mathbb{R}^d .

- We define the risk of a classifier $\mathcal{C} : X \rightarrow \{0, 1\}$ as

$$\text{Risk}(\mathcal{C}) = \mathbb{E}[\mathbb{1}\{\mathcal{C}(X) \neq Y\}]. \quad (4)$$

- We can show that a classifier which minimizes the risk is given by $x \mapsto \mathcal{C}^*(x)$ where $\mathcal{C}^*(x) = \mathbb{1}\{q^*(x) \geq 1/2\}$ with

$$q^*(x) = \mathbb{E}[Y|X = x] = \mathbb{P}(Y = 1|X = x) = p(x, 1)/(p(x, 0) + p(x, 1)). \quad (5)$$

- Proof: see in small classes.

Introduction to classification: hypothesis class

- Take $Y = \{0, 1\}$. Recall $\mathcal{C}^*(x) = \mathbb{1} \{q^*(x) \geq 1/2\}$.
- Here we aim to model somehow $q^*(x) = \mathbb{P}(Y = 0|X = x)$ ($\mathbb{P}(Y = 1)$ is then completely determined...)
- To this end, we approximate $x \mapsto q^*(x)$ by functions

$$h_w : x \mapsto \sigma \circ f_w(x) \quad \sigma : \mathbb{R} \rightarrow [0, 1] , \quad (6)$$

- $f_w \in \mathcal{F}$ where \mathcal{F} is a hypothesis class of functions from X to \mathbb{R} parametrized by w .
- For example,

$$\mathcal{F}_\phi = \left\{ x \mapsto \phi(x)^T w + w_0 : w \in \mathbb{R}^{d+1} \right\} , \quad \phi(x) = (\phi_1(x), \dots, \phi_d(x))^T . \quad (7)$$

Introduction to classification: hypothesis class

- Take $Y = \{0, 1\}$. Recall $\mathcal{C}^*(x) = \mathbb{1} \{q^*(x) \geq 1/2\}$.
- Here we aim to model somehow $q^*(x) = \mathbb{P}(Y = 1|X = x)$ ($\mathbb{P}(Y = 0|X = x)$ is then completely determined...)
- To this end, we approximate $x \mapsto q^*(x) - 1/2$ by functions

$$h_w : x \mapsto \sigma \circ f_w(x) \quad \sigma : \mathbb{R} \rightarrow [0, 1] , \quad (8)$$

with

- $f_w \in \mathcal{F}$ where \mathcal{F} is a hypothesis class of functions from X to \mathbb{R} parametrized by w .
- Here σ is called an activation function.
- Ideally, in discriminative deterministic learning, we would like to choose $\sigma_1 : t \mapsto \mathbb{1} \{t - 1/2 \geq 0\}$ but non-smooth. Hence we use a regularization of this function (in general).
- In many cases, σ will come from distribution function from probabilistic models that we will use.

- Take $Y = \{0, 1\}$. Recall $\mathcal{C}^*(x) = \mathbb{1} \{q^*(x) \geq 1/2\}$.
- Here we aim to model somehow $q^*(x) = \mathbb{P}(Y = 1|X = x)$
($\mathbb{P}(Y = 0|X = x)$ is then completely determined...)
- To this end, we approximate $x \mapsto q^*(x) - 1/2$ by functions

$$h_w : x \mapsto \sigma \circ f_w(x) \quad \sigma : \mathbb{R} \rightarrow [0, 1], \quad (9)$$

with $f_w \in \mathcal{F}$ where \mathcal{F} is a hypothesis class of functions from X to \mathbb{R} parametrized by w .

- This defines a hypothesis class of functions from X to $[0, 1]$,

$$\mathcal{H} = \{\sigma \circ f : f \in \mathcal{F}\} = \{\sigma \circ f_w : w \in \Theta\}. \quad (10)$$

Introduction to classification: hypothesis class

- Take $Y = \{0, 1\}$. Recall $\mathcal{C}^*(x) = \mathbb{1} \{q^*(x) \geq 1/2\}$.
- Here we aim to model somehow $q^*(x) = \mathbb{P}(Y = 0|X = x)$ ($\mathbb{P}(Y = 1)$ is then completely determined...)
- To this end, we approximate $x \mapsto q^*(x)$ by functions

$$h_w : x \mapsto \sigma \circ f_w(x) \quad \sigma : \mathbb{R} \rightarrow [0, 1], \quad (11)$$

$f_w \in \mathcal{F}$ where \mathcal{F} is a hypothesis class of functions from X to \mathbb{R} parametrized by w .

- This defines a hypothesis class of functions from X to $[0, 1]$,

$$\mathcal{H} = \{\sigma \circ f : f \in \mathcal{F}\} = \{\sigma \circ f_w : w \in \Theta\}. \quad (12)$$

- Given a learned $h_{\hat{w}}$, the predictor/classifier is

$$\mathcal{C}_{\hat{w}} : x \mapsto \mathbb{1} \{h_{\hat{w}}(x) \geq \eta\}, \text{ where } \eta \text{ is a threshold, most of the times } 1/2 \text{ or } 0. \quad (13)$$

- The set

$$\{x \in X : h_{\hat{w}}(x) = \eta\}. \quad (14)$$


is called the decision boundary associated with $h_{\hat{w}}$.

Introduction to classification: hypothesis class

- Take $Y = \{0, 1\}$. If we change the risk with an asymmetric risk

$$\text{Risk}_\eta(\mathcal{C}) = \eta \mathbb{P}(Y = 1, \mathcal{C}(X) = 0) + (1 - \eta) \mathbb{P}(Y = 0, \mathcal{C}(X) = 1) , \quad (15)$$

the optimal classifier is no more $\mathcal{C}^*(x) = \mathbb{1} \{q^*(x) \geq 1/2\}$ but $\mathcal{C}^*(x) = \mathbb{1} \{q^*(x) \geq \eta\}$.

- Here recall $q^*(x) = \mathbb{P}(Y = 0|X = x)$.
- The extreme cases are easy... 
- Given a learned $h_{\hat{w}}$ approximating q^* , the predictor/classifier is then in that case

$$\mathcal{C}_{\hat{w}} : x \mapsto \mathbb{1} \{h_{\hat{w}}(x) \geq \eta\} , \text{ where } \eta \text{ is a threshold, most of the times } 1/2 \text{ or } 0 . \quad (16)$$

- The set

$$\{x \in X : h_{\hat{w}}(x) = \eta\} . \quad (17)$$

is called the decision boundary associated with $h_{\hat{w}}$.

- Take $Y = \{0, 1\}$.
- We consider first deterministic discriminative learning;
- and the activation function $\sigma(t) = \mathbb{1} \{t \geq 1/2\}$.
- In this simple case, the previous discussion simplifies and the classifiers have the form

$$\mathcal{C}_w : x \mapsto \mathbb{1} \{f_w(x) \geq 0\} , \quad f_w \in \mathcal{F} , \quad (18)$$

with \mathcal{F} hypothesis class of functions from $X \rightarrow \mathbb{R}$.

- In the following, we consider only linear classifiers:

$$\mathcal{C}_w(x) = \mathbb{1} \{f_w(x) > 0\} , \quad f_w(x) = w^T \Phi(x) , \quad \Phi : X \rightarrow \mathbb{R}^d . \quad (19)$$

Introduction

Bayes classifier

Deterministic discriminative learning

Generative probabilistic models

Naive Bayes

Discriminant analysis (linear and quadratic)

Some reminders

- Choose a class of hypotheses/representations $\mathcal{C} = \{\mathcal{C}_w : X \rightarrow Y : w \in \Theta\}$.
- Choose a loss function $\ell : Y \times X \times \Theta$.
- Define the error function

$$E_\ell(w) = \sum_{i=1}^n \ell(y_i, x_i, w) , \quad \text{equivalently} \quad E_\ell(w) = n^{-1} \sum_{i=1}^n \ell(y_i, x_i, w) . \quad (20)$$

- Choose an algorithm to solve

$$\text{minimize } E_\ell . \quad (21)$$

- Here, we start by specifying \mathcal{C} and the form of ℓ .

- First, choose the class $\mathcal{C} = \{\mathcal{C}_w : X \rightarrow Y : w \in \mathbb{R}^d\}$.
- Consider $Y = \{0, 1\}$ and $X = \mathbb{R}^d$.
- The most simple predictors are linear predictor

$$\mathcal{C}_w : x \mapsto \mathbb{1}\{f_w(x) \geq 0\} , \quad f_w(x) = w^T x + w_0 . \quad (22)$$

- Here the decision boundary is a the affine hyperplan:

$$H = \{(w_0, w) \in \mathbb{R}^{d+1} : w^T x + w_0 = 0\} . \quad (23)$$

- w is called a weight vector;
- w_0 a bias (not in statistical sense) or sometimes a negative threshold (why?).

Linear discriminant functions

- First, choose the class $\mathcal{C} = \{\mathcal{C}_w : X \rightarrow Y : w \in \mathbb{R}^d\}$.
- Consider $Y = \{0, 1\}$ and $X = \mathbb{R}^d$.
- The most simple predictors are linear predictors

$$\mathcal{C}_w : x \mapsto \mathbb{1} \{f_w(x) \geq 0\} , \quad f_w(x) = w^T x + w_0 . \quad (24)$$

- Adding a component equal to 1 to x (dummy variable), we can simply consider

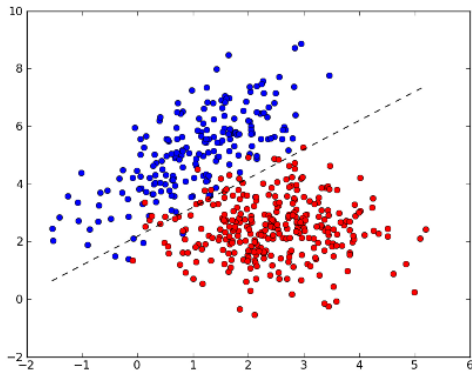
$$\mathcal{C}_w : x \mapsto \mathbb{1} \{f_w(x) \geq 0\} , \quad f_w(x) = w^T x . \quad (25)$$

- In what follows, if not specify, I always consider that a dummy variable is added to the data...
- Define a loss ℓ and the error function

$$E_\ell(w) = n^{-1} \sum_{i=1}^n \ell(y_i, f_w(x_i)) . \quad (26)$$

- Choose an algorithm to solve minimize E_ℓ .

Geometrically



- Learn a boundary to separate two “groups” of points.

Perceptron algorithm

- Here we consider $Y = \{-1, 1\}$ which allow a more simple presentation of the perceptron algorithm.
- If data $\in \{0, 1\}$, apply the simple transformation

$$y \mapsto 2y - 1, \quad \text{with inverse } y \mapsto (y + 1)/2. \quad (27)$$

- We consider here a slight generalization of the previous class function using basis functions $\{\phi_j : X \rightarrow \mathbb{R}\}$.
- Denote by $\phi(x) = (\phi_1(x), \dots, \phi_d(x))^T$ and consider the class of linear discriminant functions associated with:

$$\mathcal{C} = \left\{ x \mapsto \mathcal{C}(x) = 2 \times \mathbb{1} \left\{ \phi(x)^T w > 0 \right\} - 1 : w \in \mathbb{R}^d \right\} \quad (28)$$

- Now which loss and error functions to estimate the weight w ?

Perceptron algorithm

- Here we consider $Y = \{-1, 1\}$ which allow a more simple presentation of the perceptron algorithm.
- Denote by $\phi(x) = (\phi_1(x), \dots, \phi_d(x))^T$ and consider the class of linear discriminant functions associated with:


$$\mathcal{C} = \left\{ x \mapsto \mathcal{C}_w(x) = 2 \times \mathbb{1} \{f_w(x) > 0\} - 1 : w \in \mathbb{R}^d \right\}, \quad f_w(x) = \phi(x)^T w. \quad (29)$$

- Here we consider the error function:

$$E_{\text{perc}}(w) = - \sum_{i=1}^N y_i f_w(x_i) \mathbb{1} \{y_i f_w(x_i) < 0\}. \quad (30)$$

- Justification: we aim to fit the data:
 - if the data i well-classified, no contribution since $y_i f_w(x_i) \geq 0$;
 - if the data i not well-classified, linear contribution $-y_i \phi_x(x_i)^T w \geq 0$.
- Exercise: what is the corresponding loss?

The perceptron algorithm: introduction to stochastic gradient descent

- We aim to minimize $E(w) = \sum_{i=1}^N E_i(w)$, $E_i = y_i f_w(x_i) \mathbb{1} \{y_i f_w(x_i) < 0\}$.
- $\nabla E_i(w) = ??$ 
- First option gradient descent (GD):

$$w_{k+1} = w_k - \eta \nabla E(w) . \quad (31)$$

- The parameter $\eta > 0$ is called the learning rate in ML or stepsize in optimization.
- Problem of GD: ∇E not always accessible
 - $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^N$ not always entirely available;
 - data may arrive sequentially/ in a continuous stream;
 - yet you want to make some predictions and adjust them with new data...
 - N is too large...

The perceptron algorithm: introduction to stochastic gradient descent

- We aim to minimize $E(w) = \sum_{i=1}^N E_i(w)$, $E_i = y_i f_w(x_i) \mathbb{1} \{y_i f_w(x_i) < 0\}$.

Algorithm 1: Online SGD

- Initialize w_0
- Given a stream of data $\{(x_i, y_i)\}_{i=1}^N$, for $k = 0 \dots$
 - $w_{k+1} = w_k - \eta_k \nabla E_k(w_k) = w_k - \eta_k y_k \phi(x_k) \mathbb{1} \{y_k \phi(x_k)^T w_k < 0\}$
- $(\eta_k)_{k \in \mathbb{N}}$ is a sequence of stepsize which either $\rightarrow 0$ and is constant $\eta_k \equiv \eta$.
- The data is used online.
- Once the k -th pair (x_k, y_k) has been used, it can be forgotten.
- Number of iterations: number of data which can be collected.

The perceptron algorithm: introduction to stochastic gradient descent

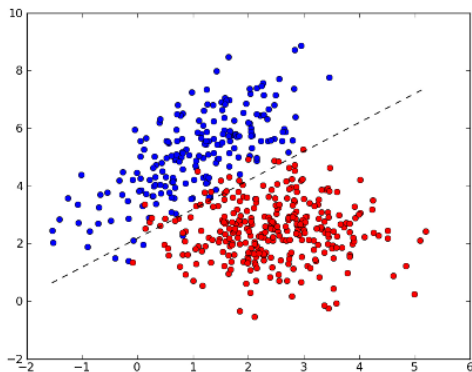
- We aim to minimize $E(w) = \sum_{i=1}^N E_i(w)$, $E_i = y_i f_w(x_i) \mathbb{1} \{y_i f_w(x_i) < 0\}$.
- When N is too large

Algorithm 2: Batch SGD

- Initialize w_0 , size of a batch b ;
 - Given a dataset $\{(x_i, y_i)\}_{i=1}^N$, for $k = 0 \dots$
 - Take a random batch $B_k \subset \{1, \dots, N\}$ with $\text{Card}(B_k) = b$;
 - $w_{k+1} = w_k - \eta_k \sum_{i \in B_k} \nabla E_i(w_k) = w_k - \eta_k \sum_{i \in B_k} y_i \phi(x_i) \mathbb{1} \{y_i \phi(x_i)^T w_k < 0\}$
-
- The data is used by batches.
 - Number of iterations: potentially infinite.
 - Under appropriate condition $(w_k)_{k \in \mathbb{N}}$ converges to a minimizer of E .
 - An alternative to $(w_k)_{k \in \mathbb{N}}$ is the corresponding Polyak-Ruppert averaging.

- Cons
 - Perceptron does not generalize well for $K > 2$;
 - does not provide probabilistic outputs;
- In the sequel, we present probabilistic models to circumvent the last issue.

Geometrically



- Learn a boundary to separate two “groups” of points.
- ? : what about Y with $K = \text{Card}(Y) > 2$?

Case $K = \text{Card}(Y) > 2$

- Several strategies for $Y = \{1, \dots, K\}$.
- One-versus-the rest: consider the $K - 1$ classifiers \mathcal{C}_k associated to the $\{0, 1\}$ -prediction problems for $k \in \{1, \dots, K - 1\}$:

$$y = k, \quad y \neq k. \quad (32)$$

- The one-versus-the rest classifier:

$$\mathcal{C}^{\text{OR}} = \sum_{k=1}^{K-1} (k+1) \mathcal{C}_k. \quad (33)$$

- This makes sense for $K = 2$?
- However, this leads to ambiguous classification, i.e., classifiers give non-consistent results
- Example: $\mathcal{C}_{K-1}(x) = \mathcal{C}_{K-2}(x) = 1$ or $\mathcal{C}^{\text{OR}}(x) \notin Y = \{1, \dots, K\}$.

Case $K = \text{Card}(Y) > 2$

- Several strategies for $Y = \{1, \dots, K\}$.
- One-versus-one: consider the $K(K-1)/2$ classifiers $\mathcal{C}_{k,i}$ associated to the $\{0,1\}$ -prediction problems for $k, i \in \{1, \dots, K\}$, $i \neq k$:

$$y = k, \quad y = i. \quad (34)$$

- The one-versus-the rest classifier:

$$\mathcal{C}^{\text{OR}} = \sum_{\substack{k,i=1 \\ k \neq i}}^K (k+1) \mathcal{C}_{k,i}. \quad (35)$$

- However, this leads to ambiguous classification again.

Case $K = \text{Card}(Y) > 2$

- Several strategies for $Y = \{1, \dots, K\}$.
- A constant solution is to consider K discriminant functions $\{f_k : X \rightarrow \mathbb{R}\}_{k=1}^K$ and set

$$\mathcal{C}^{K, \text{lin}}(x) = \underset{k \in Y}{\operatorname{argmax}} f_k(x) . \quad (36)$$

- It is in line with the formal classification problem.
- Defining a classifier is in fact equivalent to define a partition of $X = \sqcup_{i=1}^K X_i$.
- Indeed, if we have a classifier \mathcal{C} , we can define $X_i = \{x \mid \mathcal{C}(x) = i\} = \mathcal{C}^{-1}(\{i\})$.
- If we have a partition, we can define $\mathcal{C}(x) = \sum_{i=1}^K i \mathbb{1}_{X_i}(x)$.

Case $K = \text{Card}(Y) > 2$

- Several strategies for $Y = \{1, \dots, K\}$.
- A constant solution is to consider K discriminant functions $\{f_k : X \rightarrow \mathbb{R}\}_{k=1}^K$ and set

$$\mathcal{C}^{K, \text{lin}}(x) = \min \arg\max_{k \in Y} f_k(x) . \quad (37)$$

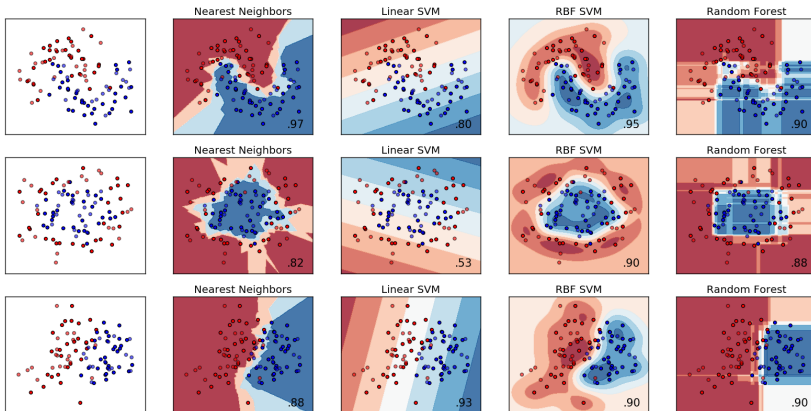
- It is in line with the formal classification problem.
- The \min in the definition (39) is to make a choice in the case $\arg\max_{k \in Y} f_k(x)$ admits strictly more than one element.
- For $A \subset Y$, $A = \{k_1, \dots\}$,

$$D_A = \{x : \arg\max_k f_k(x) = A\} , \quad (38)$$

is called the decision boundary associated with the labels in A .

Classification


...many ways to separate points!




Case $K = \text{Card}(Y) > 2$

- Several strategies for $Y = \{1, \dots, K\}$.
- A consistent solution is to consider K discriminant functions $\{f_k : X \rightarrow \mathbb{R}\}_{k=1}^K$ and set

$$\mathcal{C}^{K, \text{lin}}(x) = \min_{k \in Y} \arg \max_{k \in Y} f_k(x) . \quad (39)$$

- Example: for $X = \mathbb{R}^d$, $f_k(x) = w_k^T x$
-  For the moment $\{w_k\}_{k=1}^K$ are arbitrary and define the hypothesis class of our classifiers.
- They have to be learned to fit the data; up coming!
- Then $D_{\{1,2\}} = \{x : (w_1 - w_2)^T x\}$.
- For any $A \subset Y$, D_A is in fact convex (polygon):

$$x_1, x_2 \in D_A \implies tx_1 + (1-t)x_2 \in D_A \text{ for any } t \in [0, 1] . \quad (40)$$

-  : how to find $\{f_k\}_{k=1}^K$ or equivalently here $\{w_k\}_{k=1}^K$?

Least square for classification

- We consider a 1 of K binary coding scheme

$$y = k \mapsto y = \mathbf{e}_k = \begin{pmatrix} 0 \\ \vdots \\ 1 \text{ position } k \\ \vdots \\ 0 \end{pmatrix}. \quad (41)$$

- And then K linear discriminant (bias included...):

$$f_{w^{(k)}}(x) = [w^{(k)}]^T x. \quad (42)$$

- This gives the model

$$F_W(x) = Wx, \quad \text{where } W = [w^{(1)} \dots w^{(K)}] \in \mathbb{R}^{d \times K}. \quad (43)$$

- We expect that if $y_i = k$, $(Wx_i)_k > (Wx_i)_{k'}$, $k' \neq k$.
- This motivates the introduction of the error function:

$$E(W) = \frac{1}{2} \sum_{i=1}^N \|y_i - W^T w_i\|^2. \quad (44)$$

- E can be written of the of the form

$$E(W) = \frac{1}{2} \text{trace} \left((Y - XW)^T (Y - XW) \right) , \quad (45)$$

with

$$Y = [y_1 \cdots y_N]^T \in \mathbb{R}^{N \times K} , \quad X = [x_1 \cdots x_N]^T \in \mathbb{R}^{N \times d} . \quad (46)$$

- E is quadratic ! and admits then a unique minimizer:

$$\hat{W} = X^\dagger Y = (X^T X)^{-1} X^T Y , \quad (47)$$

if $X^T X$ is invertible.

- Drawbacks
 - lack of robustness to outliers;
 - does not corresponds to any data: Solution to E corresponds to the maximum likelihood estimator for Gaussian conditional likelihood.
- In the sequel, we explore more models to addresse these issues.

Fisher's linear discriminant

- We consider the two classes problem $Y = \{0, 1\}$.
- Projection can lead to significant loss of information in the data.
- However, we can seek a vector w such that the projection onto $\text{Span}(w)$ which maximizes the separation between the two classes.
- The simplest measure of separation is the separation between the projected means:

$$\text{maximize } w \mapsto w^T \{m_1 - m_0\}, \text{ subject to } \|w\| = 1, \quad (48)$$

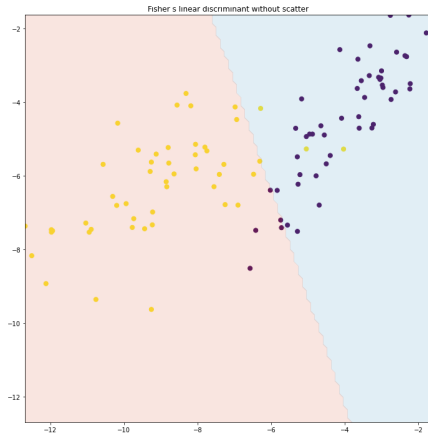
$$m_i = \sum_{j=1}^N x_j \mathbb{1}\{y_j = i\} / \sum_{j=1}^N \mathbb{1}\{y_j = i\}. \quad (49)$$

- Solution

$$\hat{w} = (m_1 - m_0) / \|m_1 - m_0\|. \quad (50)$$

- However this criteria is too sensible to correlated data: does not catch second order information.

Example: Fisher first attempt



- We consider the two classes problem $Y = \{0, 1\}$.
- Projection can lead to significant loss of information in the data.
- However, we can seek a vector w such that the projection onto $\text{Span}(w)$ which maximizes the separation between the two classes.
- To catch second order information, we aim to take into account the variance of the data once projected to w :

$$s^2 = w^T \Sigma w, \quad \mathbf{S}_W = \sum_{j=1}^N (x_j - m_i)(x_j - m_i)^T \mathbb{1}\{y_j = i\}. \quad (51)$$

- Remarks:
 - $N^{-1}\mathbf{S}_W$: empirical covariance for the $\{x_i\}_{i=1}^N$ within class.
 - It is implicitly assumed that the covariances of two classes are the same.
 - s is sometimes called scatter.

Fisher's linear discriminant

- We consider the two classes problem $Y = \{0, 1\}$.
- Projection can lead to significant loss of information in the data.
- However, we can seek a vector w such that the projection onto $\text{Span}(w)$ which maximizes the separation between the two classes.
- To catch second order information, we aim to take into account the variance of the data once projected to w :

$$s^2 = w^T \Sigma w, \quad \mathbf{A} = \sum_{j=1}^N (x_j - m_i)(x_j - m_i)^T \mathbb{1}_{\{y_j = i\}}. \quad (52)$$

- Based on these quantities, we can seek the projection which maximizes the separation between the projected means and with smallest scatter:

$$J(w) = (w^T(m_1 - m_0))^2 / s^2 = \frac{w^T \mathbf{S}_B w}{w^T \mathbf{S}_W w}, \quad (53)$$

with $N^{-1} \mathbf{S}_B = (m_1 - m_0)(m_1 - m_0)^T / N$ between-class covariance matrix.

- Based on these quantities, we can seek the projection which maximizes the separation between the projected means and with smallest scatter:

$$J(w) = (w^T(m_1 - m_0))^2 / s^2 = \frac{w^T \mathbf{S}_B w}{w^T \mathbf{S}_W w}, \quad (54)$$

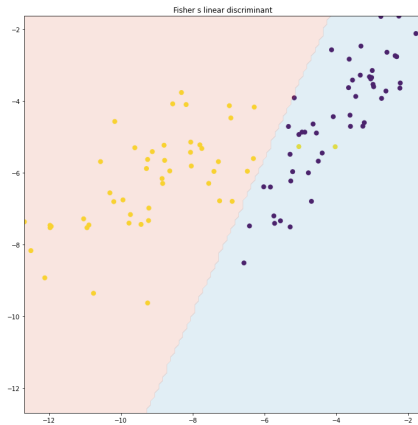
with $N^{-1} \mathbf{S}_B = (m_1 - m_0)(m_1 - m_0)^T / N$ between-class covariance matrix.

- $\nabla J(w) = 0$ implies that there exist $a_1, a_2, a_3 \in \mathbb{R}$ such that 

$$\text{Conclusion: } \hat{w} \propto \mathbf{S}_W^{-1}(m_1 - m_0). \quad (55)$$

- Maximizer by Cauchy-Schwarz inequality!

Example: Fisher DA



Introduction

Bayes classifier

Deterministic discriminative learning

Generative probabilistic models

Naive Bayes

Discriminant analysis (linear and quadratic)

Formal presentation

- We consider now a probabilistic model for the data $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^N$.
- We always assume if not specified that the pairs are i.i.d observations of a pair of random variable (X, Y) !
- In the generative probabilistic model setting:
 - we choose a prior on Y : $\mathbb{P}(Y = k) = q_k$;
 - a parametric model for $\{\mathbb{P}_w(X \in \cdot | Y = k) : w \in \mathbb{R}^d\}$;
 - we assume that $\mathbb{P}_w(X \in \cdot | Y = k)$ has a density $p_w(\cdot | k)$.
- This gives rise to the joint model and the likelihood


$$L(w; \mathcal{D}) = \prod_{i=1}^N \{p_w(x_i | y_i) q_{y_i}\} . \quad (56)$$

- In what follows, w is then estimated by maximum likelihood estimation:

$$\hat{w} \in \operatorname{argmax}_w L(w; \mathcal{D}) . \quad (57)$$

- The prediction is then:

$$y_{\text{pred}} = \operatorname{argmax}_{k \in Y} \mathbb{P}_{\hat{w}}(Y = k | X = x_{\text{new}}) = \operatorname{argmax}_{k \in Y} [q_k p_{\hat{w}}(x_{\text{new}} | k)] . \quad (58)$$

- We consider now a probabilistic model for the data $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^N$.
- We always assume if not specified that the pairs are i.i.d observations of a pair of random variable (X, Y) !
- In the generative probabilistic model setting:
 - we estimate $\mathbb{P}(Y = k) = q_k$;
 - choose a parametric model for $\{\mathbb{P}_w(X \in \cdot | Y = k) : w \in \mathbb{R}^d\}$;
 - we assume that $\mathbb{P}_w(X \in \cdot | Y = k)$ has a density $p_w(\cdot | k)$.
- Estimation of $\mathbb{P}(Y = k) = q_k$? 

- We consider now a probabilistic model for the data $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^N$.
- We always assume if not specified that the pairs are i.i.d observations of a pair of random variable (X, Y) !
- In the generative probabilistic model setting:
 - we choose a prior on Y : $\mathbb{P}(Y = k) = q_k$ or we can use the MLE...;
 - a parametric model for $\{\mathbb{P}_w(X \in \cdot | Y = k) : w \in \mathbb{R}^d\}$;
 - we assume that $\mathbb{P}_w(X \in \cdot | Y = k)$ has a density $p_w(\cdot | k)$.
- It remains to specify the model $\{\mathbb{P}_w(X \in \cdot | Y = k) : w \in \mathbb{R}^d\}$.

Introduction

Bayes classifier

Deterministic discriminative learning

Generative probabilistic models

Naive Bayes

Discriminant analysis (linear and quadratic)

- The Naive Bayes assumption assume that $X = (X^{(1)}, \dots, X^{(d)})$ have independent component given Y :

$$\mathbb{P}_w(X|Y) = \prod_{i=1}^d \mathbb{P}_w(X^{(i)}|Y) .$$

- Crude modeling for $\mathbb{P}_w(X|Y)$;
- **Feature independence** assumption;
- **Simple featurewise model**: binomial if binary, multinomial if finite and Gaussian if continuous.
- If all features are continuous, the law of X given Y is Gaussian with a **diagonal covariance matrix**!

Gaussian Naive Bayes

- The Naive Bayes assumption assume that $X = (X^{(1)}, \dots, X^{(d)})$ have independent component given Y :

$$\mathbb{P}_w(X|Y) = \prod_{i=1}^d \mathbb{P}_w(X^{(i)}|Y).$$

- For $k \in Y = \{1, \dots, K\}$, the conditional density of $X^{(i)}$ given $\{Y = k\}$ is

$$p_w(x^{(i)}|k) = (2\pi\sigma_{i,k}^2)^{-1/2} \exp \left\{ -(x^{(i)} - \mu_{i,k})^2 / (2\sigma_{i,k}^2) \right\}.$$

- The complete conditional distribution of X given $\{Y = k\}$ is then

$$p_w(x|y = k) = (\det(2\pi\Sigma_k))^{-1/2} \exp \left\{ -(x - \mu_k)^T \Sigma_k^{-1} (x - \mu_k) / 2 \right\},$$

- The parameter w consists in $w = \{(\mu_k, \Sigma_k)\}_{k=1}^K$ with

$$\Sigma_k = \text{diag}(\sigma_{1,k}^2, \dots, \sigma_{d,k}^2) \text{ and } \mu_k = (\mu_{1,k}, \dots, \mu_{d,k})^T. \quad (59)$$

- MLE? 

Gaussian Naive Bayes

- The Naive Bayes assumption assume that $X = (X^{(1)}, \dots, X^{(d)})$ have independent component given Y :

$$\mathbb{P}_w(X|Y) = \prod_{i=1}^d \mathbb{P}_w(X^{(i)}|Y).$$

- For $k \in Y = \{1, \dots, K\}$, the conditional density of $X^{(i)}$ given $\{Y = k\}$ is

$$p_w(x^{(i)}|k) = (2\pi\sigma_{i,k}^2)^{-1/2} \exp \left\{ -(x^{(i)} - \mu_{i,k})^2 / (2\sigma_{i,k}^2) \right\}.$$

- The complete conditional distribution of X given $\{Y = k\}$ is then

$$p_w(x|y = k) = (\det(2\pi\Sigma_k))^{-1/2} \exp \left\{ -(x - \mu_k)^T \Sigma_k^{-1} (x - \mu_k) / 2 \right\},$$

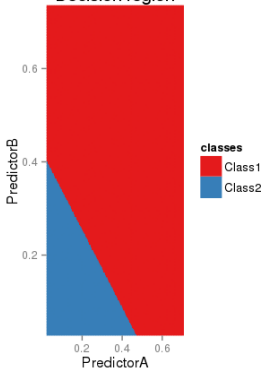
- The parameter w consists in $w = \{(\mu_k, \Sigma_k)\}_{k=1}^K$ with

$$\Sigma_k = \text{diag}(\sigma_{1,k}^2, \dots, \sigma_{d,k}^2) \text{ and } \mu_k = (\mu_{1,k}, \dots, \mu_{d,k})^T. \quad (60)$$

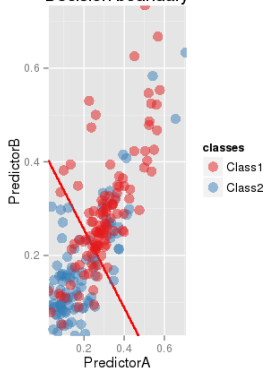
- MLE? 

Naive Bayes with Gaussian model

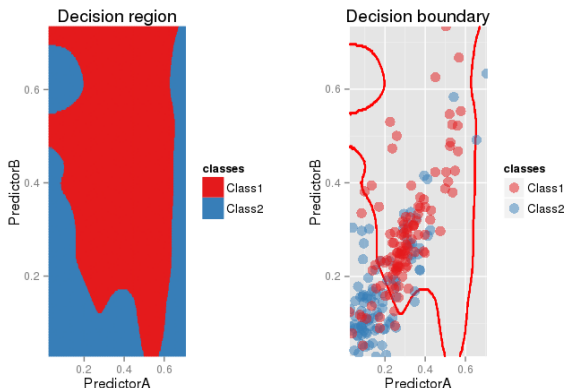
Decision region



Decision boundary



Naive Bayes with kernel density estimates



- In this experiment, $\mathbb{P}(X^{(i)} \in \cdot | Y = k)$ or more exactly the densities $p(\cdot)$ for $k \in Y$ are estimated by kernel density estimation and used to obtain a classifier:

$$y_{\text{pred}} = \underset{k \in Y}{\operatorname{argmax}} [q_k \hat{p}(x_{\text{new}} | k)] . \quad (61)$$

Introduction

Bayes classifier

Deterministic discriminative learning

Generative probabilistic models


Naive Bayes

Discriminant analysis (linear and quadratic)

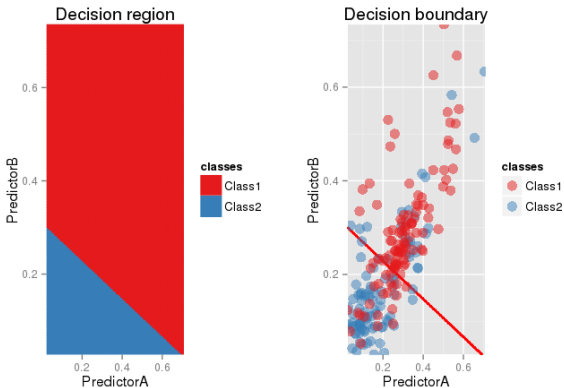
Discriminant Analysis (Gaussian model)

- Other models for $\{\mathbb{P}_w(X \in \cdot | Y = k) : w \in \mathbb{R}^d\}$ in the case $X = \mathbb{R}^d$.
- For any $k \in Y$, linear discriminant analysis (QDA):

$$\mathbb{P}_w(X \in \cdot | Y = k) = N(\mu_k, \Sigma) . \quad (62)$$

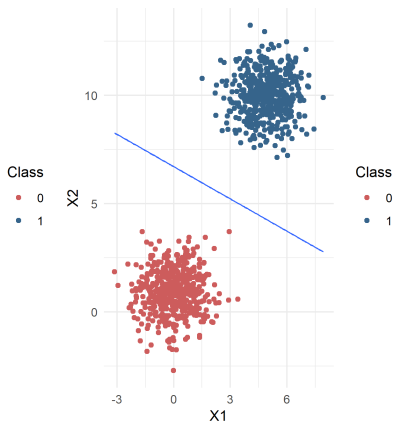
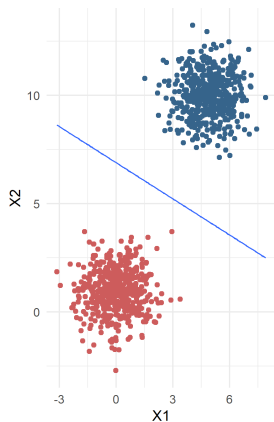
- Here the parameter $w = (\{\mu_k\}_{k=1}^K, \Sigma)$.
- MLE? 
- $K = 2$, then the **decision boundary is an affine hyperplane**.

Linear Discriminant Analysis




Example: LDA

```
boundary_true_parameters = function(x, mu0, mu1, Sigma, pi0){  
  u = t(as.matrix(mu1-mu0)) %*% inv(Sigma)  
  v = (u %*% (matrix(x - ((mu1+mu0)/2)) )) - log(pi0/(1-pi0))  
  return(as.numeric(v))  
}
```

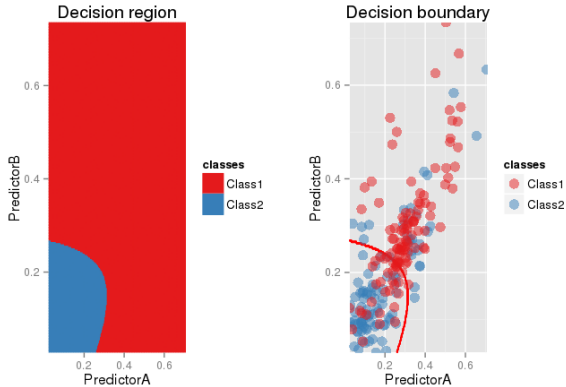


- Other models for $\{\mathbb{P}_w(X \in \cdot | Y = k) : w \in \mathbb{R}^d\}$ in the case $X = \mathbb{R}^d$.
- For any $k \in Y$, quadratic discriminant analysis (QDA):

$$\mathbb{P}_w(X \in \cdot | Y = k) = N(\mu_k, \Sigma_k) . \quad (63)$$

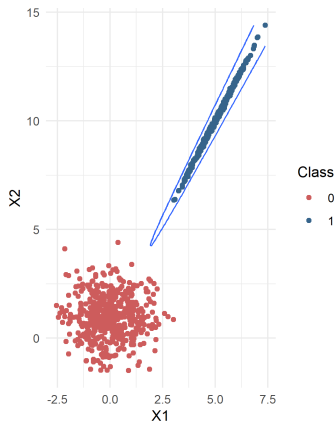
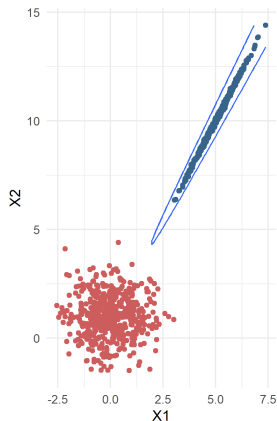
- Here the parameter $w = \{(\mu_k, \Sigma_k)\}_{k=1}^K$
- MLE? 

Quadratic Discriminant Analysis



Example: QDA

```
boundary_true_parameters_quadratic = function(x, mu0, mu1, Sigma0, Sigma1, pi0){  
  u1 = -0.5*(t(as.matrix(x-mu1))) %*% inv(Sigma1) %*% as.matrix((x - mu1))  
  u0 = 0.5*(t(as.matrix(x-mu0))) %*% inv(Sigma0) %*% as.matrix((x - mu0))  
  cste = - log(pi0/(1-pi0))  
  bonus = -0.5*log(abs(det(Sigma1))) + 0.5*log(abs(det(Sigma0)))  
  
  return(as.numeric(u1+u0+cste+bonus))  
}
```



- Relation between Fisher's DA and least squares [Bis07, Section 4.1.5].
- Fisher's discriminant for multiple classes [Bis07, Section 4.1.6].
- Probabilistic generative models, discrete features/Exponential family [Bis07, Section 4.2.3, 4.2.4].



Christopher M. Bishop. *Pattern Recognition and Machine Learning (Information Science and Statistics)*. 1st ed. Springer, 2007. ISBN: 0387310738.