

APRENDIZAJE DE MÁQUINA

CLASIFICACIÓN LINEAL

AGENDA

01 Introducción

Problema de clasificación y regresión lineal

O2 Clasificación binaria

Regresión logística y método de Newton

Modelos Lineales Generalizados
Familia exponencial, Construyendo GLMs, Softmax

O4 Modelos discriminativos y generativos

Diferencias, Análisis de Gaussianas, Clasificador de Bayes

O5 Métricas de evaluación binarias
Sensibilidad, Especificidad, Puntaje F1, Curva ROC





INTRODUCCIÓN PROBLEMA DE CLASIFICACIÓN



El **problema** es el **mismo** que el de **regresión**: **predecir** un conjunto de variables de **salidas** y dado un **conjunto** de **datos** de **entrada** X.

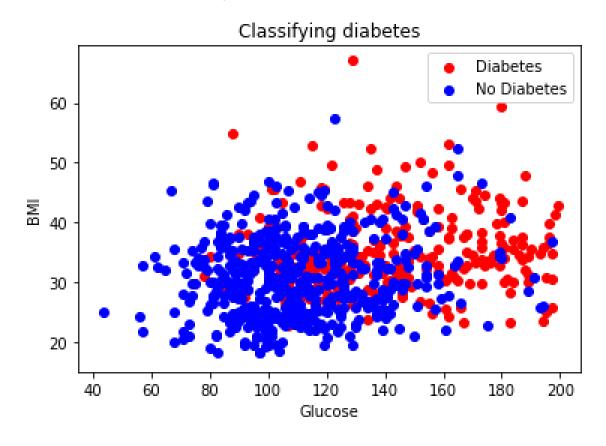
La única deferencia recae en que los valores de y toman un conjunto de valores discretos.

x_1 =Glucosa [mg/dl]	x_2 =IMC	$oldsymbol{y}=$ Presencia de Diabetes
148	33.6	1
85	26.6	0
183	23.3	1
89	28.1	0
137	43.1	1
:	:	:

INTRODUCCIÓN PROBLEMA DE CLASIFICACIÓN

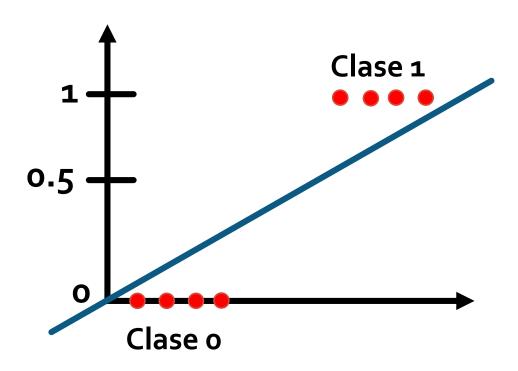


Graficamos los 752 datos con Python:



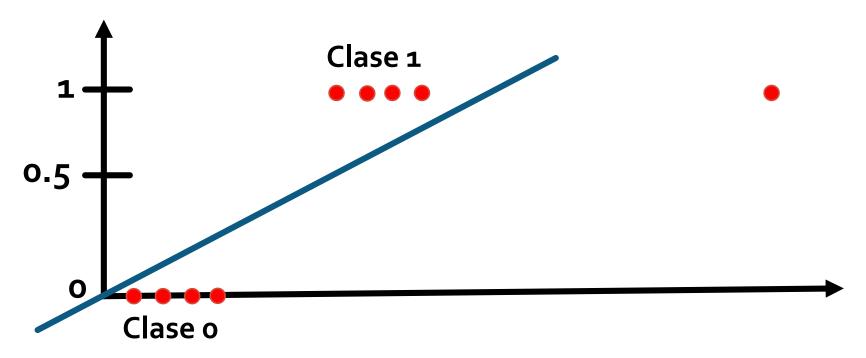
INTRODUCCIÓN REGRESIÓN LINEAL

¿PORQUÉ NO APLICAR LA REGRESIÓN LINEAL PARA CLASIFICAR?



I N T R O D U C C I Ó N R E S I Ó N L I N E A L

¿PORQUÉ NO APLICAR LA REGRESIÓN LINEAL PARA CLASIFICAR?



La función $w^T X$ no describe el valor esperado de $Y \in \{0, 1\}$



REGRESIÓN LOGÍSTICA H I P Ó T E S I S



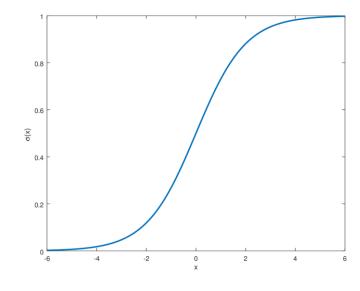
La **hipótesis** que utilizamos **para** la **regresión lineal no** es **adecuada** para resolver nuestro **problema** de **clasificación**.

Se propone una hipótesis diferente que se ajuste mejor al problema:

$$h_w(x) = g(w^T X) = \frac{1}{1 + e^{-w^T X}}$$

A esta función se le denomina función sigmoide o logística.

Se garantiza que $g(z) \in \{0, 1\}$.



REGRESIÓN LOGÍSTICA FUNCIÓN SIGMOIDE



TAREA

Demostrar que:

$$g'(\mathbf{z}) = g(\mathbf{z})(\mathbf{1} - g(\mathbf{z}))$$

REGRESIÓN LOGÍSTICA S U P O S I C I O N E S



Ahora se define el **mismo problema** de **encontrar** la **mejor combinación** de **pesos** que se **ajuste** al **problema** de **clasificación**.

Se aplican las siguientes suposiciones:

1. Suponemos que la hipótesis define una medida de probabilidad (Bernoulli):

$$P(y = 1/x; w) = h_w(x)$$
 $P(y = 0/x; w) = 1 - h_w(x)$
 $P(y/x; w) = h_w(x)^y (1 - h_w(x))^{1-y}$
 $L(w) = P(y/x; w)$

REGRESIÓN LOGÍSTICA S U P O S I C I O N E S



2. Suponiendo que se toman m muestras iid, se tiene que la verosimilitud se puede escribir como:

$$L(w) = \prod_{i=1}^{m} P(y^{(i)}/x^{(i)}; w)$$

$$L(w) = \prod_{i=1}^{m} h_{w}(x^{(i)})^{y^{(i)}} (1 - h_{w}(x^{(i)}))^{1-y^{(i)}}$$

REGRESIÓN LOGÍSTICA PÉRDIDA LOGARÍTMICA



Se puede escribir la pérdida logarítmica:

$$l(w) = log \prod_{i=1}^{m} P(y^{(i)}/x^{(i)}; w)$$

$$l(w) = \sum_{i=1}^{m} y^{(i)} log h_{w}(x^{(i)}) + (1 - y^{(i)}) log (1 - h_{w}(x^{(i)}))$$

$$l(w) = \sum_{i=1}^{m} y^{(i)} log \left(\frac{1}{1 + e^{-w^{T}X}} \right) + \left(1 - y^{(i)} \right) log \left(1 - \frac{1}{1 + e^{-w^{T}X}} \right)$$

REGRESIÓN LOGÍSTICA FUNCIÓN DE COSTO



Se puede escribir la función de costo:

$$J(w) = -l(w)$$

$$J(w) = -\sum_{i=1}^{m} y^{(i)} log \left(\frac{1}{1 + e^{-w^{T}X}} \right) + \left(1 - y^{(i)} \right) log \left(1 - \frac{1}{1 + e^{-w^{T}X}} \right)$$

REGRESIÓN LOGÍSTICA FUNCIÓN DE COSTO



Interpretando la función de costo:

$$J(w) = -\sum_{i=1}^{m} y^{(i)} \log h_w(x^{(i)}) + (1 - y^{(i)}) \log (1 - h_w(x^{(i)}))$$

Si y = 0	Siy = 1
$J(y=0,\widehat{y})=-log(1-\widehat{y})$	$J(y=1,\widehat{y})=-log(\widehat{y})$
$\lim_{\widehat{y}\to 0} \boldsymbol{log}(1-\widehat{y}) \to 0$	$\lim_{\widehat{y}\to0} \boldsymbol{log}(\widehat{y}) \to \infty$
$\lim_{\widehat{y}\to 1} log(1-\widehat{y}) \to \infty$	$\lim_{\widehat{y}\to 1} log(\widehat{y}) \to 0$

REGRESIÓN LOGÍSTICA O P T I M I Z A C I Ó N



¿CÓMO ENCONTRAR LOS MEJORES PESOS W?

REGRESIÓN LOGÍSTICA DESCENSO POR GRADIENTE



Encontrar la mejor combinación de pesos w por descenso de gradiente se ajuste al problema de clasificación.

$$\mathbf{w} \coloneqq \mathbf{w} - \alpha \, \nabla_{\mathbf{w}} \mathbf{J}(\mathbf{w})$$

REGRESIÓN LOGÍSTICA DESCENSO POR GRADIENTE



TAREA

Demostrar que el **gradiente** respecto a un **solo dato** de **entrenamiento** (x, y) y **peso** w_i está dado por:

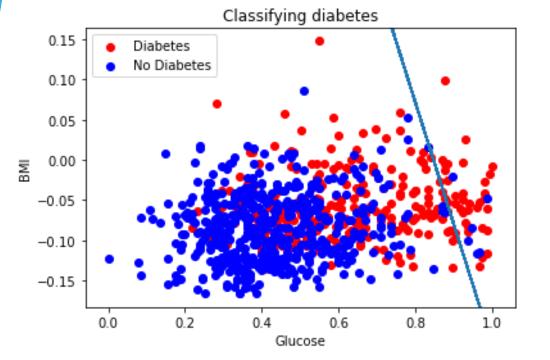
$$\frac{\partial}{\partial w_j}J(w)=(y-h_w(x))x_j$$

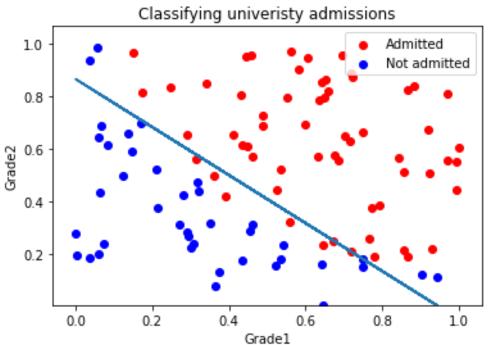
Donde $x, w_j \in \mathbb{R}^n$ y $y \in \{0, 1\}$

RECORDAR EL DESCENSO POR GRADIENTE EN REGRESIÓN LINEAL → MISMO RESULTADO

REGRESIÓN LOGÍSTICA E J E M P L O R E A L







MÉTODO DE NEWTON O P T I M I Z A C I Ó N



Recordamos de métodos numéricos, el **método** de **Newton**, en donde se encuentra el valor de w para el cual f(w) = 0:

$$w := w - \frac{f(w)}{f'(w)}$$

$$\text{slope} = f'(x_n)$$

$$f(x_n)$$

MÉTODO DE NEWTON O P T I M I Z A C I Ó N



Si queremos encontrar el mínimo de la función de costo J(w), esto correspondería a encontrar los puntos en donde $\nabla_w J(w) = 0$, por lo que se puede usar el **Método** de **Newton** (mucho **más rápido**):

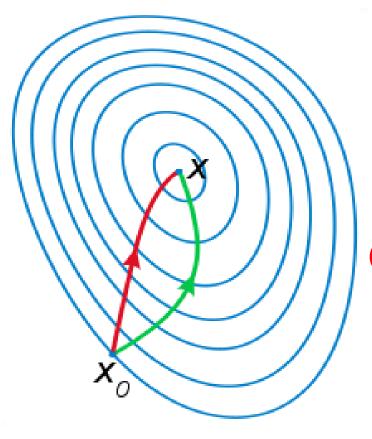
$$w \coloneqq w - H^{-1} \nabla_w J(w)$$

donde H^{-1} es la matriz Hessiana y ∇_w es el gradiente.

$$H = \begin{bmatrix} \frac{\partial^2 J(w)}{\partial w_1^2} & \dots & \frac{\partial^2 J(w)}{\partial w_1 \partial w_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial^2 J(w)}{\partial w_n \partial w_1} & \dots & \frac{\partial^2 J(w)}{\partial w_n^2} \end{bmatrix}$$

MÉTODO DE NEWTON O P T I M I Z A C I Ó N





Descenso por gradiente (Convergencia lineal)

Método de Newton (Convergencia cuadrática)



¿PORQUÉ NO USAR EL MÉTODO DE NEWTON EN LUGAR DE DESCENSO POR GRADIENTE?



MODELOS LINEALES GENERALIZADOS FAMILIA EXPONENCIAL



Recordamos que realizamos las siguientes suposiciones:

MÍNIMOS CUADRADOS

 $y \in \mathbb{R} \sim Gaussiana$ (Regresión lineal)

REGRESIÓN LOGÍSTICA

 $y \in \{0, 1\} \sim Bernoulli$ (Clasificación binaria)

MODELOS LINEALES GENERALIZADOS

FAMILIA EXPONENCIAL



Decimos que una distribución de probabilidad pertenece a la familia exponencial si se puede escribir de la siguiente forma:

$$p(y; \eta) = b(y)e^{\left(\eta^T T(y) - a(\eta)\right)}$$

Donde:

- η es denominado como el **parámetro natural** o **canónico** de la distribución.
- T(y) estadístico suficiente (estadístico que resume la información completa de una muestra).
- $a(\eta)$ la función logística de partición.

NOTA: $e^{-a(\eta)}$ funciona como constante de normalización para asegurar que $p(y;\eta)$ integre a 1.

MODELOS LINEALES GENERALIZADOS FAMILIA EXPONENCIAL



$$p(y; \eta) = b(y)e^{(\eta^T T(y) - a(\eta))}$$

Sí fijamos a, b y T, entonces podemos decir que tenemos una distribución parametrizada solo por η , donde al variar η se tienen diferentes distribuciones.

MODELOS LINEALES GENERALIZADOS EJEMPLOS DE FAMILIA EXPONENCIAL



Probamos que la distribución de Bernoulli es parte de la familia exponencial:

$$p(y; \phi) = \phi^y (1 - \phi)^{1-y}$$

Al variar ϕ se tienen diferentes distribuciones $p(y; \phi)$.

MODELOS LINEALES GENERALIZADOS

EJEMPLOS DE FAMILIA EXPONENCIAL

Encontramos:

$$T(y) = y$$

$$b(y) = 1$$

$$\eta = log\left(\frac{\phi}{1-\phi}\right)$$

$$a(\eta) = -log(1-\phi)$$

$$p(y; \phi) = e^{\left(log\left(\frac{\phi}{1-\phi}\right)y + log(1-\phi)\right)}$$

MODELOS LINEALES GENERALIZADOS EJEMPLOS DE FAMILIA EXPONENCIAL



Probamos que la distribución Gaussiana (donde σ no importa) es parte de la familia exponencial (por lo que $\sigma^2 = 1$):

$$p(y;\mu) = \frac{1}{\sqrt{2\pi}}e^{\frac{(y-\mu)^2}{2}}$$

TAREA 1

MODELOS LINEALES GENERALIZADOS CONSTRUYENDO LOS MODELOS



Vamos a realizar las siguientes suposiciones:

- 1. La variable de salida y / X; $w \sim FamExp(\eta)$
- 2. Objetivo: dada una matriz de características X, calcular $E(T(y); \eta)$ denominada como función canónica de respuesta.

$$h(\eta) = E(T(y); \eta)$$

3. Se define que la relación entre η , w y X es lineal (solo sí $\eta \in \mathbb{R}$):

$$\eta = w^T X$$

MODELOS LINEALES GENERALIZADOS CONSTRUYENDO LOS MODELOS



EJEMPLO: Bernoulli

Para valores fijos de *X* y *w* se tiene que el **objetivo** es **calcular**:

$$h(w) = E(T(y); \eta)$$

Pero sabemos que T(y) = y por lo que el valor esperado de una variable que se distribuye como Bernoulli es:

$$E(y; \eta) = P(y = 1; \eta)$$

De igual forma, sabemos de la familia exponencial que:

$$\phi = P(y = 1; \eta)$$

MODELOS LINEALES GENERALIZADOS

EJEMPLOS DE FAMILIA EXPONENCIAL



Se tiene de la definición de la familia exponencial que el parámetro natural η está expresado como:

$$\eta = log\left(\frac{\phi}{1-\phi}\right)$$

Al **despejar** ϕ de la igualdad anterior:

$$\eta = log\left(\frac{\phi}{1-\phi}\right)$$

Se tiene lo siguiente:

$$\phi = \frac{1}{1 + e^{-\eta}}$$

MODELOS LINEALES GENERALIZADOS CONSTRUYENDO LOS MODELOS



Por lo que el **valor esperado** es:

$$E(y; \eta) = \phi = \frac{1}{1 + e^{-\eta}}$$

Pero gracias a la tercera suposición $oldsymbol{\eta} = w^T X$

$$E(y; w, X) = \phi = \frac{1}{1 + e^{-W^T X}}$$

FUNCIÓN SIGMOIDE

MODELOS LINEALES GENERALIZADOS CONSTRUYENDO LOS MODELOS



EJEMPLO: Gaussiana

Probar que la función canónica de respuesta es equivalente a W^TX suponiendo que $y \sim Gaussiana$ y que $\eta = W^TX$

$$h(w) = E(T(y); \eta) = w^T X$$

TAREA 2

MODELOS LINEALES GENERALIZADOS

REGRESIÓN SOFTMAX



Se define el siguiente problema, en donde la variable de salida $y \in \{1, ..., k\}$ se distribuye como una función multinomial $y \sim Multinomial$.

- Parámetros: $\phi_1, \phi_2, ..., \phi_{k-1}$
- $P(y=i) = \phi_i$
- $\phi_k = 1 (\phi_1 + \phi_2 + \dots + \phi_{k-1})$

Se define el **estadístico suficiente** T(y) como **vector**:

$$T(1) = \begin{bmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix} T(2) = \begin{bmatrix} 0 \\ 1 \\ \vdots \\ 0 \end{bmatrix} \dots T(k-1) = \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 1 \end{bmatrix} T(k) = \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix} \in \mathbb{R}^{k-1}$$

Función indicador:

$$1\{True\} = 1$$
 $1\{False\} = 0$ $T(y)_i = 1\{y = i\}$

REGRESIÓN SOFTMAX



Se expresa la distribución multinomial en forma de la familia exponencial:

$$P(y) = \phi_1^{1\{y=1\}} \phi_2^{1\{y=2\}} \dots \phi_k^{1\{y=k\}}$$

$$P(y) = \phi_1^{T(y)_1} \phi_2^{T(y)_2} \dots \phi_k^{T(y)_{k-1}} \phi_k^{1 - \sum_{j=1}^{k-1} T(y)_j}$$

$$P(y) = e^{\left[T(y)_1 \log(\phi_1) + T(y)_2 \log(\phi_2) + \dots + (1 - \sum_{j=1}^{k-1} T(y)_j) \log(\phi_k)\right]}$$

$$P(y) = e^{\left[T(y)_{1} \log\left(\frac{\phi_{1}}{\phi_{k}}\right) + T(y)_{2} \log\left(\frac{\phi_{2}}{\phi_{k}}\right) + \dots + T(y)_{k-1} \log\left(\frac{\phi_{k-1}}{\phi_{k}}\right) + \log(\phi_{k})\right]}$$

REGRESIÓN SOFTMAX



Se compara con la forma de la familia exponencial:

$$P(y) = e^{\left[T(y)_1 \log\left(\frac{\phi_1}{\phi_k}\right) + T(y)_2 \log\left(\frac{\phi_2}{\phi_k}\right) + \dots + T(y)_{k-1} \log\left(\frac{\phi_{k-1}}{\phi_k}\right) + \log(\phi_k)\right]}$$

$$p(y; \eta) = b(y)e^{(\eta^T T(y) - a(\eta))}$$

Por lo tanto:

$$b(y)=1$$

$$a(\eta) = -\log(\phi_k)$$

REGRESIÓN SOFTMAX



$$P(y) = e^{\left[T(y)_{1} \log\left(\frac{\phi_{1}}{\phi_{k}}\right) + T(y)_{2} \log\left(\frac{\phi_{2}}{\phi_{k}}\right) + \dots + T(y)_{k-1} \log\left(\frac{\phi_{k-1}}{\phi_{k}}\right) + \log(\phi_{k})\right]}$$

$$p(y; \eta) = b(y)e^{(\eta^T T(y) - a(\eta))}$$

Finalmente:

$$oldsymbol{\eta} = egin{bmatrix} logigg(rac{oldsymbol{\phi}_1}{oldsymbol{\phi}_k}igg) \ logigg(rac{oldsymbol{\phi}_{k-1}}{oldsymbol{\phi}_k}igg) \end{bmatrix}$$

REGRESIÓN SOFTMAX



La función canónica de respuesta sería:

$$\eta_i = log\left(\frac{\phi_i}{\phi_k}\right)$$

$$\eta_k = log\left(\frac{\phi_k}{\phi_k}\right) = 0$$

$$\phi_k e^{\eta_i} = \phi_i$$

REGRESIÓN SOFTMAX



La suma de probabilidades debe ser igual a 1:

$$\sum_{i=1}^k \phi_k e^{\eta_i} = \sum_{i=1}^k \phi_i = 1$$

$$\phi_k \sum_{i=1}^k e^{\eta_i} = \sum_{i=1}^k \phi_i = 1$$

$$\phi_k = \frac{1}{\sum_{i=1}^k e^{\eta_i}}$$

REGRESIÓN SOFTMAX



Sustituyendo
$$oldsymbol{\phi}_k = rac{1}{\sum_{i=1}^k e^{\eta_i}}$$
 en la función canónica de respuesta:

$$\phi_k e^{\eta_i} = \phi_i$$

$$\phi_i = \frac{e^{\eta_i}}{\sum_{j=1}^k e^{\eta_j}}$$

FUNCIÓN SOFTMAX

REGRESIÓN SOFTMAX



Utilizando la **tercera suposición** de que $\eta_i = w_i^T X^*$

$$\phi_i = \frac{e^{w_i^T X}}{\sum_{j=1}^k e^{w_j^T X}}$$

Sabiendo que el objetivo es calcular $h(w) = E(T(y); \eta)$:

$$E\left(\begin{bmatrix} T(y)_1 \\ T(y)_2 \\ \vdots \\ T(y)_{k-1} \end{bmatrix}\right) = E\left(\begin{bmatrix} \mathbf{1}\{y=1\} \\ \mathbf{1}\{y=2\} \\ \vdots \\ \mathbf{1}\{y=k-1\} \end{bmatrix}\right) = \begin{bmatrix} \boldsymbol{\phi}_1 \\ \boldsymbol{\phi}_2 \\ \vdots \\ \boldsymbol{\phi}_{k-1} \end{bmatrix}$$

REGRESIÓN SOFTMAX



La hipótesis h(w) queda expresada como un vector:

$$h(w) = \begin{bmatrix} \frac{e^{w_1^T X}}{\sum_{j=1}^k e^{w_j^T X}} \\ \frac{e^{w_2^T X}}{\sum_{j=1}^k e^{w_j^T X}} \\ \vdots \\ \frac{e^{w_{k-1}^T X}}{\sum_{j=1}^k e^{w_j^T X}} \end{bmatrix} = \begin{bmatrix} P(y = 1) \\ P(y = 2) \\ \vdots \\ P(y = k-1) \end{bmatrix}$$

que tiene como componentes la probabilidad cada clase i: P(y = i)

REGRESIÓN SOFTMAX



Calculamos la **pérdida logarítmica** y **función de costo** para **un solo** dato de **entrenamiento**:

$$l(w) = \sum_{l=1}^{k} \mathbf{1}(y = l) log(p(y/x)) = \sum_{i=1}^{k} \mathbf{1}(y = l) log\left(\frac{e^{w_i^T X}}{\sum_{j=1}^{k} e^{w_j^T X}}\right)$$

$$J(w) = -l(w)$$

REGRESIÓN SOFTMAX



Calculamos la **pérdida logarítmica** y **función de costo** para m datos de **entrenamiento**:

$$l(w) = \sum_{i=1}^{m} \sum_{l=1}^{k} 1(y^{(i)} = l) log \left(\frac{e^{w_l^T x^{(i)}}}{\sum_{j=1}^{k} e^{w_j^T x^{(i)}}} \right)$$

$$J(w) = -l(w)$$

SE PUEDE USAR EL MÉTODO DE NEWTON O DESCENSO POR GRADIENTE





Modelos discriminativos:

- Los modelos que se han visto hasta ahora, se denominan como modelos discriminativos, en donde se realiza el aprendizaje sobre la distribución de probabilidad p(y/x) directamente.
- **Ejemplos**: regresión logística y mínimos cuadrados.

Modelos generativos:

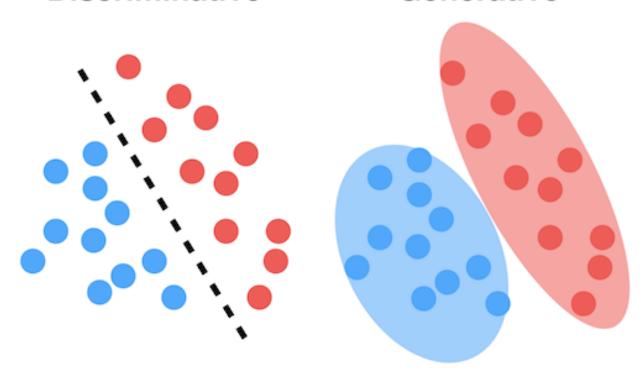
- Modelos que tratan de aprender p(x/y) y p(y). Para el caso de clasificación binaria el algoritmo aprendería tres distribuciones distintas p(x/y=0), p(x/y=1) y p(y).
- Entonces se clasificaría un nuevo dato comparándolo con ambas distribuciones p(x/y = 0) y p(x/y = 1).

MODELOS DISCRIMINATIVOS Y G E N E R A T I V O S D I F E R E N C I A S



Discriminative

Generative





D E S C R I P C I Ó N

Un modelo generativo modela las características que están condicionadas por la variable de respuesta p(x/y).

Utilizando el **Teorema** de **Bayes** se puede calcular p(y/x).

$$p(y/x) = \frac{p(x/y)p(y)}{p(x)}$$

Se tiene que el denominador se puede calcular de la siguiente forma (clasificación binaria):

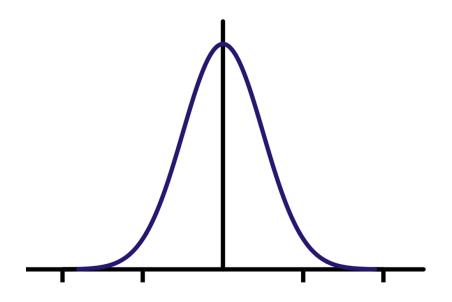
$$p(x) = p(x/y = 1)p(y = 1) + p(x/y = 0)p(y = 0)$$

MODELOS GENERATIVOS ANÁLISIS DISCRIMINATIVO GAUSSIANO



En un análisis discriminativo Gaussiano suponemos que:

$$p(x/y)\sim Gaussiana$$



MODELOS GENERATIVOS GAUSSIANAS MULTIVARIADAS



Una **distribución Gaussiana** de *d* **dimensiones** es parametrizada por:

- Vector de media: $\mu \in \mathbb{R}^d$
- Matriz de covarianza: $\Sigma \in \mathbb{R}^{dxd}$

donde Σ es simétrica y positiva semi-definida. Su **densidad** se puede **calcular** como:

$$p(x;\mu,\Sigma) = \frac{1}{(2\pi)^{d/2}|\Sigma|^{1/2}}e^{\left(-\frac{1}{2}(x-\mu)^T\Sigma^{-1}(x-\mu)\right)}$$

Propiedades:

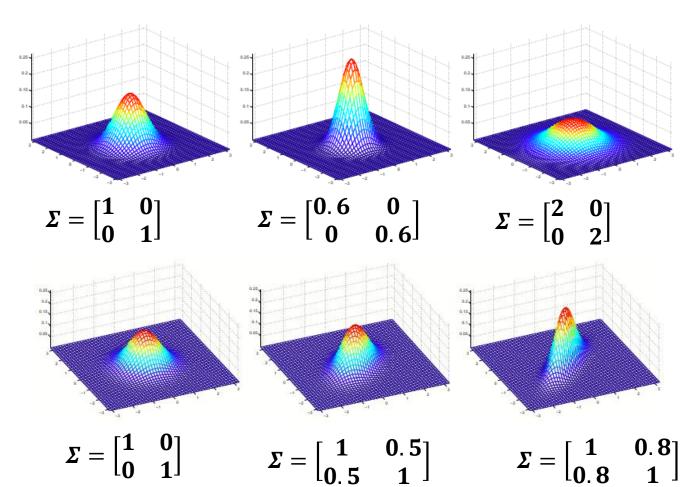
$$E[X] = \int_{\mathcal{X}} x \, p(x; \mu, \Sigma) \, dx = \mu$$

$$Cov(X) = E[(X - E[X])(X - E[X])^{T}] = \Sigma$$

GAUSSIANAS MULTIVARIADAS



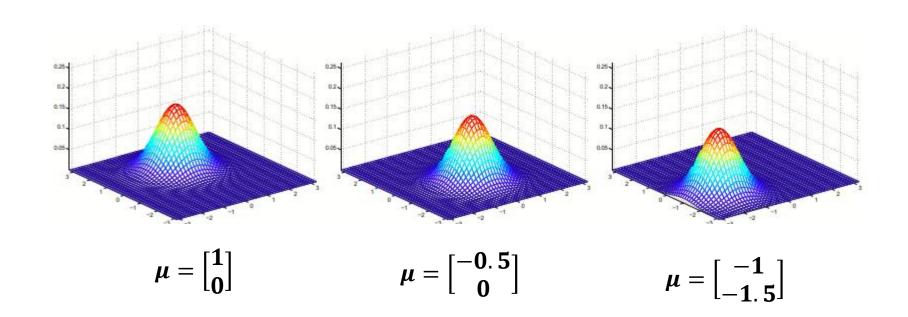
Ejemplo: Variando Σ .



GAUSSIANAS MULTIVARIADAS



Ejemplo: Variando μ .





ANÁLISIS DISCRIMINATIVO GAUSSIANO

El análisis discriminativo Gaussiano resuelve el problema de clasificación binaria donde:

$$x \in \mathbb{R}^n$$
$$y \in \{0, 1\}$$

Las **suposiciones** del **modelo son**:

$$x/y = 0 \sim N(\mu_0, \Sigma)$$

$$x/y = 1 \sim N(\mu_1, \Sigma)$$

MODELOS GENERATIVOS ANÁLISIS DISCRIMINATIVO GAUSSIANO



Las distribuciones están dadas por:

$$p(y) = \phi^{y} (1 - \phi)^{1 - y}$$

$$p(x|y = 0) = \frac{1}{(2\pi)^{d/2} |\Sigma|^{1/2}} \exp\left(-\frac{1}{2} (x - \mu_0)^T \Sigma^{-1} (x - \mu_0)\right)$$

$$p(x|y = 1) = \frac{1}{(2\pi)^{d/2} |\Sigma|^{1/2}} \exp\left(-\frac{1}{2} (x - \mu_1)^T \Sigma^{-1} (x - \mu_1)\right)$$

Los parámetros de las distribuciones serían: ϕ , μ_0 , μ_1 , Σ



ANÁLISIS DISCRIMINATIVO GAUSSIANO

La verosimilitud conjunta de los datos estaría dada por:

$$l(\phi, \mu_k, \Sigma) = \prod_{i=1}^m p(x^{(i)}, y^{(i)}; \phi, \mu_k, \Sigma)$$

$$l(\phi, \mu_k, \Sigma) = \prod_{i=1}^{m} p(x^{(i)}/y^{(i)}; \mu_k, \Sigma) p(y^{(i)}; \phi)$$

$$\log l(\phi, \mu_k, \Sigma) = \log \prod_{i=1}^{m} p(x^{(i)}/y^{(i)}; \mu_k, \Sigma) p(y^{(i)}; \phi)$$

$$\log l(\phi, \mu_k, \Sigma) = \sum_{i=1}^{m} \log(p(x^{(i)}/y^{(i)}; \mu_k, \Sigma) + \log p(y^{(i)}; \phi)$$

MODELOS GENERATIVOS ANÁLISIS DISCRIMINATIVO GAUSSIANO



La verosimilitud conjunta de los datos estaría dada por:

$$\log l(\phi, \mu_k, \Sigma) = \sum_{i=1}^m \log(p(x^{(i)}/y^{(i)}; \mu_k, \Sigma) + \log p(y^{(i)}; \phi)$$

$$log \ l(\phi, \mu_k, \Sigma) = \sum_{i=1}^{m} log \left(\frac{1}{(2\pi)^{\frac{d}{2}} |\Sigma|^{\frac{1}{2}}} \right) + log \left[e^{\left(-\frac{1}{2} (x^{(i)} - \mu_k)^T \Sigma^{-1} (x^{(i)} - \mu_k) \right)} \right] + log(\Phi^{y^{(i)}} (1 - \Phi)^{1 - y^{(i)}})$$

$$log \ l(\phi, \mu_k, \Sigma) = \sum_{i=1}^{m} log(\frac{1}{(2\pi)^{d/2} |\Sigma|^{1/2}}) - \frac{1}{2} (x^{(i)} - \mu_k)^T \Sigma^{-1} (x^{(i)} - \mu_k) + y^{(i)} \ log(\phi) + (1 - y^{(i)}) \log((1 - \Phi))$$



ANÁLISIS DISCRIMINATIVO GAUSSIANO

Maximizando respecto a ϕ :

$$\nabla_{\emptyset} \log l(\phi, \mu_k, \Sigma) = \nabla_{\phi} \sum_{i=1}^{m} y^{(i)} \log \left(\frac{\phi}{1 - \phi} \right) + \log((1 - \Phi)) = 0$$

$$\nabla_{\emptyset} \log l(\phi, \mu_k, \Sigma) = \sum_{i=1}^{m} \left(\frac{\mathbf{y}^{(i)}}{\phi(1-\phi)} - \frac{1}{1-\phi} \right) = 0$$

$$\sum_{i=1}^{m} \frac{y^{(i)}}{\phi(1-\phi)} = \frac{m}{1-\phi}$$

$$\nabla_{\emptyset} \log l(\phi, \mu_k, \Sigma) = \sum_{i=1}^{m} \frac{\mathbf{y}^{(i)}}{m} = \sum_{i=1}^{m} \frac{1(\mathbf{y}^{(i)} = \mathbf{k})}{m} = \phi = \frac{numero\ de\ datos\ de\ la\ clase\ k}{total\ de\ datos}$$

ANÁLISIS DISCRIMINATIVO GAUSSIANO



Maximizando respecto a μ_k :

$$\nabla_{\mu_k} \log l(\phi, \mu_k, \Sigma) = \nabla_{\mu_k} \sum_{i=1}^m -\frac{1}{2} (x^{(i)} - \mu_k)^T \Sigma^{-1} (x^{(i)} - \mu_k) = 0$$

Aplicando la propiedad $\nabla_x x^T A x = 2Ax$:

$$\nabla_{\mu_k} \log l(\phi, \mu_k, \Sigma) = \sum_{i=1}^m \Sigma^{-1} (x^{(i)} - \mu_k) = 0$$

$$\nabla_{\mu_k} \log l(\phi, \mu_k, \Sigma) = \Sigma^{-1} \sum_{i=1}^m (x^{(i)} - \mu_k) = 0$$

$$\nabla_{\mu_k} \log l(\phi, \mu_k, \Sigma) = \sum_{i=1}^m \mathbf{x}^{(i)} - \sum_{i=1}^m \mu_k = 0$$

MODELOS GENERATIVOS ANÁLISIS DISCRIMINATIVO GAUSSIANO



Solo nos **interesan las** $x^{(i)}$ que pertenezcan a la **clase k**:

$$\nabla_{\mu_k} \log l(\phi, \mu_k, \Sigma) = \sum_{i=1}^m x^{(i)} - \sum_{i=1}^m \mu_k = 0$$

$$\nabla_{\mu_k} \log l(\phi, \mu_k, \Sigma) = \sum_{i=1}^m 1(y^{(i)} = k) x^{(i)} - \sum_{i=1}^m 1(y^{(i)} = k) \mu_k = 0$$

$$\mu_k = \frac{\sum_{i=1}^m 1(y^{(i)} = k)x^{(i)}}{\sum_{i=1}^m 1(y^{(i)} = k)} = \frac{suma\ de\ x\ que\ pertenecen\ a\ clase\ k}{n\'umero\ de\ datos\ de\ clase\ k}$$

ANÁLISIS DISCRIMINATIVO GAUSSIANO



Maximizando respecto a Σ^{-1} :

$$\nabla_{\Sigma^{-1}} \log l(\phi, \mu_k, \Sigma) = \nabla_{\Sigma^{-1}} \sum_{i=1}^{m} log(\frac{1}{(2\pi)^{d/2} |\Sigma|^{1/2}}) - \frac{1}{2} (x^{(i)} - \mu_k)^T \Sigma^{-1} (x^{(i)} - \mu_k) = 0$$

$$\nabla_{\Sigma^{-1}} \log l(\phi, \mu_k, \Sigma) = \nabla_{\Sigma^{-1}} \sum_{i=1}^{m} log(\frac{1}{(2\pi)^{d/2} |\Sigma|^{1/2}}) - \frac{1}{2} (x^{(i)} - \mu_k)^T (x^{(i)} - \mu_k) \Sigma^{-T} = 0$$

Como Σ es simétrica a $\Sigma^{-1} = \Sigma^{-T}$:

$$\nabla_{\Sigma^{-1}} \log l(\phi, \mu_k, \Sigma) = \nabla_{\Sigma^{-1}} \sum_{i=1}^{m} -\frac{d}{2} log(2\pi) + \frac{1}{2} log(\left|\Sigma^{-1}\right|) - \frac{1}{2} \left(x^{(i)} - \mu_k\right)^T \left(x^{(i)} - \mu_k\right) \Sigma^{-1} = 0$$

$$\nabla_{\boldsymbol{\Sigma}^{-1}} \log l(\boldsymbol{\phi}, \boldsymbol{\mu}_k, \boldsymbol{\Sigma}) = \sum_{i=1}^{m} \nabla_{\boldsymbol{\Sigma}^{-1}} \frac{1}{2} log(\left|\boldsymbol{\Sigma}^{-1}\right|) - \nabla_{\boldsymbol{\Sigma}^{-1}} \frac{1}{2} \left(\boldsymbol{x}^{(i)} - \boldsymbol{\mu}_k\right)^T \left(\boldsymbol{x}^{(i)} - \boldsymbol{\mu}_k\right) \boldsymbol{\Sigma}^{-1} = 0$$

ANÁLISIS DISCRIMINATIVO GAUSSIANO



Aplicando la propiedad $\nabla_x b^T x = b$:

$$\nabla_{\boldsymbol{\Sigma}^{-1}} \log l(\boldsymbol{\phi}, \boldsymbol{\mu}_k, \boldsymbol{\Sigma}) = \frac{1}{2} \sum_{i=1}^{m} \nabla_{\boldsymbol{\Sigma}^{-1}} \left[log(\left|\boldsymbol{\Sigma}^{-1}\right|) \right] - \left(\boldsymbol{x}^{(i)} - \boldsymbol{\mu}_k\right) \left(\boldsymbol{x}^{(i)} - \boldsymbol{\mu}_k\right)^T = 0$$

Aplicando la propiedad $\nabla_x \log(|X|) = X^{-T}$:

$$\nabla_{\boldsymbol{\Sigma}^{-1}} \log l(\boldsymbol{\phi}, \boldsymbol{\mu}_k, \boldsymbol{\Sigma}) = \sum_{i=1}^{m} (\boldsymbol{\Sigma}^{-1})^{-T} - (\boldsymbol{x}^{(i)} - \boldsymbol{\mu}_k) (\boldsymbol{x}^{(i)} - \boldsymbol{\mu}_k)^T = 0$$

$$m\boldsymbol{\Sigma}^{T} - \sum_{i=1}^{m} (\boldsymbol{x}^{(i)} - \boldsymbol{\mu}_{k}) (\boldsymbol{x}^{(i)} - \boldsymbol{\mu}_{k})^{T} = 0$$

$$\Sigma^{T} = \frac{1}{m} \sum_{i=1}^{m} (x^{(i)} - \mu_{k}) (x^{(i)} - \mu_{k})^{T}$$

ANÁLISIS DISCRIMINATIVO GAUSSIANO



Como Σ es simétrica a $\Sigma^T = \Sigma$:

$$\Sigma = \frac{1}{m} \sum_{i=1}^{m} (x^{(i)} - \mu_k) (x^{(i)} - \mu_k)^T$$

VARIANZA DE LOS DATOS DE LA CLASE K

Recordando que **es un "outer"** product: $\Sigma \in \mathbb{R}^{n \times n}$

$$\Sigma = \begin{bmatrix} \frac{1}{m} \sum_{i=1}^{m} (x^{(i)}_{1} - \mu_{k})^{2} & \dots & \frac{1}{m} \sum_{i=1}^{m} (x^{(i)}_{1} - \mu_{k}) (x^{(i)}_{n} - \mu_{k}) \\ \vdots & \ddots & \vdots \\ \frac{1}{m} \sum_{i=1}^{m} (x^{(i)}_{n} - \mu_{k}) (x^{(i)}_{1} - \mu_{k}) & \dots & \frac{1}{m} \sum_{i=1}^{m} (x^{(i)}_{n} - \mu_{k})^{2} \end{bmatrix}$$

ANÁLISIS DISCRIMINATIVO GAUSSIANO



Una vez que hemos encontrado los parámetros:

$$\phi = \sum_{i=1}^{m} \frac{1(\mathbf{y}^{(i)} = \mathbf{k})}{m}$$

$$\mu_{k} = \frac{\sum_{i=1}^{m} 1(y^{(i)} = k)x^{(i)}}{\sum_{i=1}^{m} 1(y^{(i)} = k)}$$

$$\Sigma = \frac{1}{m} \sum_{i=1}^{m} (x^{(i)} - \mu_k) (x^{(i)} - \mu_k)^{T}$$

Ya podemos realizar **predicciones**.

5

ANÁLISIS DISCRIMINATIVO GAUSSIANO

Para realizar una predicción calculamos la probabilidad a posteriori P(y = 1/x) con el nuevo dato x:

$$p(y = 1/x) = \frac{p(x/y = 1)p(y = 1)}{p(x)}$$

Asignamos un umbral: si $P(y = 1/x) \ge 0.5$ pertenece a la clase 1, sino pertenece a la clase 0. Es decir estamos **maximizando** la **probabilidad**:

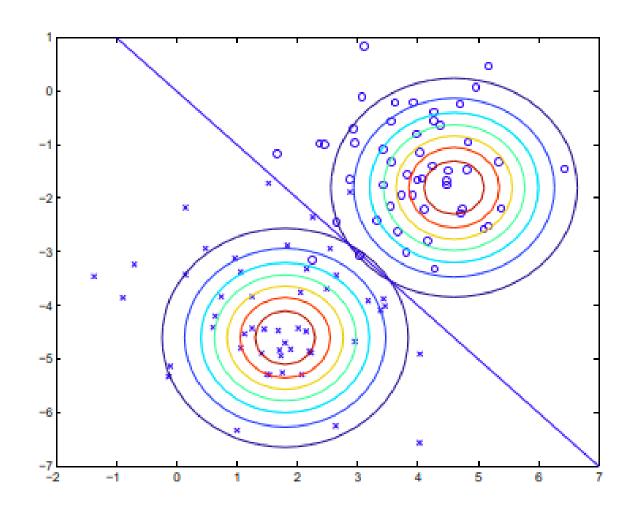
$$argmax_y p(y/x) = argmax_y \frac{p(x/y)p(y)}{p(x)}$$

Sí
$$p(y = 1) = p(y = 0) = 0.5$$
 se tiene:

$$argmax_{y}p(y/x) = argmax_{y} p(x/y)$$

MODELOS GENERATIVOS ANÁLISIS DISCRIMINATIVO GAUSSIANO





ANÁLISIS DISCRIMINATIVO GAUSSIANO



RESUMEN:

1. Obtengo los parámetros para cada Gaussiana (en clasificación binaria serían dos curvas) a partir de los datos.

$$\phi, \mu, \Sigma$$

2. Calculo las probabilidades

$$P(x/y = k) P(y) P(x)$$

3. Utilizo Teorema de Bayes para hacer una nueva predicción.

$$p(y = k/x) = \frac{p(x/y = k)p(y = k)}{p(x)}$$

COMPARACIÓN MODELOS DISCRIMINATIVOS

Si analizamos la probabilidad p(y = 1/x), en el Teorema de Bayes, podremos darnos cuenta que obtenemos la misma curva de regresión logística.

$$p(y = 1/x) = \frac{p(x/y = 1)p(y = 1)}{p(x)}$$

$$p(y = 1/x) = \frac{p(x/y = 1)p(y = 1)}{p(x/y = 1)p(y = 1) + p(x/y = 0)p(y = 0)} \left(\frac{\frac{1}{p(x/y = 1)p((y = 1))}}{\frac{1}{p(x/y = 1)p((y = 1))}}\right)$$

$$p(y = 1/x) = \frac{1}{1 + \frac{p(x/y = 0)p(y = 0)}{p(x/y = 1)p((y = 1))}}$$

INATIVOS

COMPARACIÓN MODELOS DISCRIMINATIVOS

$$\frac{p(x/y=0)p(y=0)}{p(x/y=1)p((y=1))} = e^{\left(-\frac{1}{2}(x^{(i)}-\mu_0)^T \Sigma^{-1}(x^{(i)}-\mu_0) + \frac{1}{2}(x^{(i)}-\mu_1)^T \Sigma^{-1}(x^{(i)}-\mu_1)\right)} \left(\frac{1-\phi}{\phi}\right)$$

$$\frac{p(x/y=0)p(y=0)}{p(x/y=1)p((y=1))} = e^{\left(-\frac{1}{2}(x^{(i)}-\mu_0)^T \Sigma^{-1}(x^{(i)}-\mu_0) + \frac{1}{2}(x^{(i)}-\mu_1)^T \Sigma^{-1}(x^{(i)}-\mu_1)\right)} \left(e^{\log\left(\frac{1-\phi}{\phi}\right)}\right)$$

$$\frac{p(x/y=0)p(y=0)}{p(x/y=1)p((y=1))} = e^{\left(-\frac{1}{2}\Sigma^{-1}\left[\left(x^{(i)}-\mu_0\right)^T\left(x^{(i)}-\mu_0\right)-\left(x^{(i)}-\mu_1\right)^T\left(x^{(i)}-\mu_1\right)\right] + \log\left(\frac{1-\phi}{\phi}\right)\right)}$$

$$\frac{p(x/y=0)p(y=0)}{p(x/y=1)p((y=1))} = e^{\left(-\frac{1}{2}\Sigma^{-1}\left[\sum_{j=1}^{n}(x^{(i)}_{j}-\mu_{0})^{2}-(x^{(i)}_{j}-\mu_{1})^{2}\right]+\log\left(\frac{1-\phi}{\phi}\right)\right)}$$

COMPARACIÓN MODELOS DISCRIMINATIVOS



$$\frac{p(x/y=0)p(y=0)}{p(x/y=1)p((y=1))} = e^{\left(-\frac{1}{2}\Sigma^{-1}\left[\sum_{j=1}^{n}-2\mu_0 x^{(i)}_{j}+\mu_0^2+2\mu_1 x^{(i)}_{j}-\mu_1^2\right]+\log\left(\frac{1-\phi}{\phi}\right)\right)}$$

Como $x^{(i)}_{0} = 1$

$$\frac{p(x/y=0)p(y=0)}{p(x/y=1)p((y=1))} = e^{\left(-\left[\sum_{j=1}^{n}\frac{1}{2}\Sigma^{-1}(\mu_{0}-\mu_{1})(x^{(i)}_{j})\right] + \left[\frac{1}{2}\Sigma^{-1}(\mu_{0}^{2}-\mu_{1}^{2}) + \log\left(\frac{1-\phi}{\phi}\right)\right](x^{(i)}_{0})}\right)}$$

Comparando $w^T X$:

$$w^{T} = \begin{bmatrix} -\frac{1}{2} \Sigma^{-1} (\mu_{0}^{2} - \mu_{1}^{2}) - \log \left(\frac{1 - \phi}{\phi}\right) \\ \frac{1}{2} \Sigma^{-1} (\mu_{0} - \mu_{1}) \\ \vdots \\ \frac{1}{2} \Sigma^{-1} (\mu_{0} - \mu_{1}) \end{bmatrix}$$

O S .

COMPARACIÓN MODELOS DISCRIMINATIVOS

Por lo tanto:

$$p(y = 1/x) = \frac{1}{1 + e^{-\left(\left[\sum_{j=1}^{n} \frac{1}{2} \Sigma^{-1} (\mu_0 - \mu_1)(x^{(i)}_j)\right] + \left[-\frac{1}{2} \Sigma^{-1} (\mu_0^2 - \mu_1^2) + \log\left(\frac{1-\phi}{\phi}\right)\right](x^{(i)}_0)\right)}} = \frac{1}{1 + e^{-W^T X}}$$

Donde w^T :

$$w^{T} = \begin{bmatrix} -\frac{1}{2} \Sigma^{-1} (\mu_{0}^{2} - \mu_{1}^{2}) - \log \left(\frac{1 - \phi}{\phi}\right) \\ \frac{1}{2} \Sigma^{-1} (\mu_{0} - \mu_{1}) \\ \vdots \\ \frac{1}{2} \Sigma^{-1} (\mu_{0} - \mu_{1}) \end{bmatrix}$$

COMPARACIÓN MODELOS DISCRIMINATIVOS



De manera más general:

Función de Verosimilitud $x/y \sim FamiliaExponencial(\eta)$



Distribución Posterior P(y/x)=función sigmoide

ESA ES LA RAZÓN PORQUE SE USA REGRESIÓN LOGÍSTICA → FUNCIONA PARA MUCHOS TIPOS DE SUPOSICIONES



COMPARACIÓN MODELOS DISCRIMINATIVOS

Un modelo generativo tendrá mejor rendimiento que uno discriminativo, sí la suposición que hicimos de la forma de la distribución P(x/y) se cumple para los datos reales (se provee de más información al algoritmo).

De otra forma, si los datos no se comportan como supusimos, el modelo discriminativo tendrá mejores resultados, porque aún cuando nuestras suposiciones no fueron tan certeras, la regresión logística funciona para muchas suposiciones distintas.

GENERATIVO	DISCRIMINATIVO
Mejor rendimiento cuando se conoce la distribución P(x/y)	Mejor rendimiento cuando se desconoce P(x/y) (robusto a suposiciones incorrectas)
Eficiencia asintótica	Puede existir un mejor modelo
Necesita pocos datos	Necesita muchos datos





MODELO DE EVENTOS MULTIVARIADO DE BERNOULLI:

Para el caso del análisis discriminativo Gaussiano, se tenía el caso de que $x^{(i)}_{j} \in \mathbb{R}$, es decir, eran valores continuos.

Para casos, como clasificación de textos ("correo deseado" y "correo no deseado") los vectores x pueden tener una dimensionalidad muy alta (el número de palabras es inmenso).

$$x \in \{0, 1\}^{5000}$$

Se necesitarían 2⁵⁰⁰⁰ parámetros.

$$x = \begin{bmatrix} 1 & a \\ 0 & aardvark \\ 0 & aardwolf \\ \vdots & \vdots \\ 1 & buy \\ \vdots & \vdots \\ 0 & zygmurgy \end{bmatrix}$$

MODELOS GENERATIVOS CLASIFICADOR INGENUO DE BAYES



SOLUCIÓN: suposición ingenua de Bayes (muy fuerte)

Suponemos que $x^{(i)}_{j}$ son condicionalmente independientes dado y.

Es decir, si se sabe que un texto es "correo no deseado" y = 1, el hecho de que aparezca la palabra $x_{2087} = "comprar"$ en el texto, no afecta las creencias sobre la aparición de cualquier otra palabra, como $x_{39831} = "precio"$.

$$p(x_1, \dots x_j, \dots, x_n/y) = p(x_1)p(x_2/x_1, y)p(x_3/x_1, x_2, y) \dots = p(x_1/y)p(x_2/y) \dots = \prod_{j=1}^{n} p(x_j/y)$$

AÚN CUANDO LA SUPOSICIÓN NO ES CIERTA, EL ALGORITMO TIENE BUEN RENDIMIENTO EN MUCHAS APLICACIONES

OS YES

CLASIFICADOR INGENUO DE BAYES

Como el clasificador de Bayes es un modelo generativo, es de interés modelar las distribuciones $p(x_j/y)$ y p(y). De manera particular se quiere maximizar la verosimilitud conjunta:

$$l(\phi_y, \phi_{j/y=0}\phi_{j/y=1}) = \prod_{i=1}^m p(x^{(i)}, y^{(i)})$$

Donde las **suposiciones son**:

$$p(x/y) = p(x_1, ... x_j, ..., x_n/y) = \prod_{j=1}^n p(x_j/y)$$
 $x_j / y = 0 \sim Bernoulli(\phi_{j/y=0})$
 $x_j / y = 1 \sim Bernoulli(\phi_{j/y=1})$
 $y \sim Bernoulli(\phi_j)$

) S :S

CLASIFICADOR INGENUO DE BAYES

El **resultado** de la **estimación** por **máxima verosimilitud** da:

$$\phi_{j/y=1} = \frac{\sum_{i=1}^{m} 1\{x_j^{(i)}, y^{(i)} = 1\}}{\sum_{i=1}^{m} 1\{y^{(i)} = 1\}} = \frac{\text{\# veces que se repite la palabra j en "no deseado"}}{\text{total de correos "no deseados"}}$$

$$\phi_{j/y=0} = \frac{\sum_{i=1}^m 1\{x_j^{(i)}, y^{(i)} = 0\}}{\sum_{i=1}^m 1\{y^{(i)} = 0\}} = \frac{\text{\# veces que se repite la palabra j en "deseado"}}{\text{total de correos "deseados"}}$$

$$\phi_y = \frac{\sum_{i=1}^m 1\{y^{(i)} = 1\}}{m} = \frac{total\ de\ correos\ "no\ deseados"}{total\ de\ correos}$$

Si las características $x_j^{(i)}$ toman más valores, se pueden modelar como multinomiales.

Si las características $x_i^{(i)}$ toman valores continuos, se discretiza en intervalos.

CLASIFICADOR INGENUO DE BAYES



Para realizar una nueva predicción usamos el teorema de Bayes:

$$p(y = 1/x) = \frac{p(x/y = 1)p(y = 1)}{p(x)}$$

Donde:

$$p(x/y = 1) = \prod_{j=1}^{n} p(x_j/y = 1) = \prod_{j=1}^{n} (\phi_j/y = 1)$$

$$p(y) = \phi_y = \frac{\sum_{i=1}^m 1\{y^{(i)} = 1\}}{m}$$

$$p(x) = \prod_{j=1}^{n} p(x_j/y = 1) p(y = 1) + \prod_{j=1}^{n} p(x_j/y = 0) p(y = 0)$$

MODELOS GENERATIVOS CLASIFICADOR INGENUO DE BAYES



¿Qué pasa sí aparece en un correo una palabra que nunca ha visto el algoritmo en otro correo?

CLASIFICADOR INGENUO DE BAYES



La probabilidad que asignará, de ver la palabra en cualquiera de los dos correos será 0.

$$\phi_{j/y=1} = \frac{\sum_{i=1}^{m} 1\{x_{j}^{(i)}, y^{(i)} = 1\}}{\sum_{i=1}^{m} 1\{y^{(i)} = 1\}} = \frac{\text{\# veces que se repite la palabra j en "no deseado"}}{\text{total de correos "no deseados"}} = \frac{0}{m_{negal}}$$

$$\phi_{j/y=0} = \frac{\sum_{i=1}^{m} 1\{x_j^{(i)}, y^{(i)} = 0\}}{\sum_{i=1}^{m} 1\{y^{(i)} = 0\}} = \frac{\text{\# veces que se repite la palabra j en "deseado"}}{\text{total de correos deseados"}} = \frac{0}{m_{pos}}$$

Al hacer la **predicción** tendremos una **incongruencia**:

$$p(y = 1/x) = \frac{p(x/y = 1)p(y = 1)}{p(x)} = \frac{0}{0}$$

MODELOS GENERATIVOS CLASIFICADOR INGENUO DE BAYES



ESTADÍSTICAMENTE NO ES ADECUADO DECIR QUE LA PROBABILIDAD DE UN EVENTO ES CERO SOLO PORQUE NO LO HAYAS VISTO EN TUS DATOS

NOTA: LAPLACE Y EL SOL



CLASIFICADOR INGENUO DE BAYES

SUAVIZADO DE LAPLACE:

La **solución** es **cambiar** nuestra **estimación** (caso general multinomial donde $y = \{1, 2, ..., k\}$):

$$p(y = 1) = \frac{(\# "1"s + 1)}{(\# "0"s + 1) + (\# "1"s + 1)}$$

$$p(y = j) = \phi_y = \frac{1 + \sum_{i=1}^m 1\{y^{(i)} = j\}}{m + k}$$

Al calcular los **otros estimadores** se tiene:

$$\phi_{j/y=1} = \frac{1 + \sum_{i=1}^{m} 1\{x_j^{(i)}, y^{(i)} = 1\}}{2 + \sum_{i=1}^{m} 1\{y^{(i)} = 1\}}$$

$$\phi_{j/y=0} = \frac{1 + \sum_{i=1}^{m} 1\{x_j^{(i)}, y^{(i)} = 0\}}{2 + \sum_{i=1}^{m} 1\{y^{(i)} = 0\}}$$

CLASIFICADOR INGENUO DE BAYES



MODELO DE EVENTOS MULTINOMIAL:

Vamos a ver el caso cuando $x_i = \{1, 2, ..., k\}$.

Por ejemplo, ahora se tiene que el valor de x_j representa la **posición** de la **palabra** en el **diccionario**, y el índice j indica la **posición** de la **palabra** en el **correo**.

Entonces n ahora representa la **longitud** del correo (y varía acorde a cada correo).

Las **suposiciones** serían:

$$p(x/y) = p(x_1, ... x_j, ..., x_n/y) = \prod_{j=1}^{n} p(x_j/y)$$
 $x_j / y = 0 \sim Multinomial(\phi_{j/y=0})$
 $x_j / y = 1 \sim Multinomial(\phi_{j/y=1})$
 $y \sim Bernoulli(\phi_j)$

CLASIFICADOR INGENUO DE BAYES



Los parámetros, al calcular la máxima verosimilitud serían (con suavizado de Laplace):

$$\phi_{k/y=1} = \frac{\sum_{i=1}^{m} \sum_{j=1}^{n_i} \mathbf{1}\{x_j^{(i)} = k, y^{(i)} = 1\} + 1}{\sum_{i=1}^{m} \mathbf{1}\{y^{(i)} = 1\} d_i + n_i + k} = \frac{\# \ de \ veces \ que \ aparece \ la \ palabra \ k \ en \ correos \ ND}{total \ de \ palabras \ en \ correos \ ND}$$

$$\phi_{k/y=0} = \frac{\sum_{i=1}^{m} \sum_{j=1}^{n_i} \mathbf{1} \left\{ x_j^{(i)} = k, y^{(i)} = 0 \right\} + \mathbf{1}}{\sum_{i=1}^{m} \mathbf{1} \left\{ y^{(i)} = 0 \right\} d_i + n_i + k} = \frac{\# \ de \ veces \ que \ aparece \ la \ palabra \ k \ en \ correos \ D}{total \ de \ palabras \ en \ correos \ D}$$

$$\phi_y = \frac{\sum_{i=1}^m \mathbf{1}\{y^{(i)} = \mathbf{1}\} + \mathbf{1}}{m}$$





Falso positivo (error tipo I):

- Estadística: Se rechaza incorrectamente una hipótesis nula verdadera.
- Aprendizaje máquina: el modelo predice que la clase de un dato de entrenamiento p(y/x) = 1 es positiva, cuando la realidad es que el dato pertenecía a la clase negativa y = 0.



Falso negativo (error tipo II):

- Estadística: Se acepta incorrectamente la hipótesis nula falsa.
- Aprendizaje máquina: el modelo predice que la clase de un dato de entrenamiento p(y/x) = 0 es negativa, cuando la realidad es que el dato pertenecía a la clase positiva y = 1.

Verdadero positivo:

- Estadística: Se acepta correctamente la hipótesis nula falsa.
- Aprendizaje máquina: el modelo predice que la clase de un dato de entrenamiento p(y/x) = 1 es positiva, cuando la realidad es que el dato pertenecía a la clase positiva y = 1.

Verdadero negativo:

- Estadística: Se acepta correctamente la hipótesis nula verdadera.
- Aprendizaje máquina: el modelo predice que la clase de un dato de entrenamiento p(y/x) = 0 es negativa, cuando la realidad es que el dato pertenecía a la clase negativa y = 0.

MATRIZ DE CONFUSIÓN



PREDICTIVE VALUES

POSITIVE (1) NEGATIVE (0)

TUAL VALUES

POSITIVE (1)

NEGATIVE (0)

TP	FN
FP	TN

S E N S I B I L I D A D



La sensibilidad ("true positive rate", "recall"): mide el rendimiento del modelo para predecir la clase positiva y = 1.

Es difícil que un algoritmo sensible se equivoque en predecir la clase positiva.

Aún así, una alta sensibilidad puede venir acompañada de muchos falsos positivos.

$$Sens = rac{TP}{TP + FN}$$

MÉTRICAS DE EVALUACION BINARIAS E S P E C I F I C I D A D



La especificidad ("true negative rate"): mide el rendimiento del modelo para predecir la clase negativa y=0.

Es difícil que un algoritmo negativo se equivoque en predecir la clase negativa.

Aún así, una alta especificidad puede venir acompañada de muchos falsos negativos.

$$Spec = rac{TN}{TN + FP}$$

P R E C I S I Ó N



La precisión o "accuracy" (no confundir con "positive predictive value" o "precisión") mide el rendimiento del modelo de manera general. Cuanto se desvían las predicciones de los valores reales.

$$Acc = \frac{TP + TN}{TP + TN + FN + FP}$$

VALOR POSITIVO PREDICHO



El "positive predictive value" o "precisión", mide la fracción de positivos que fueron correctamente predichos.

$$PPV = \frac{TP}{TP + FP}$$

PUNTAJE F1



$$Sens = \frac{TP}{TP + FN}$$

$$PPV = \frac{TP}{TP + FP}$$

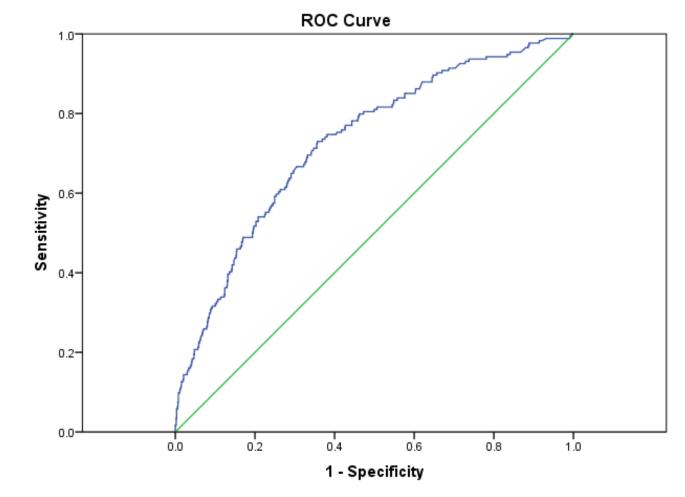
$$F1 = \frac{2}{Sens^{-1} + PPV^{-1}}$$

CURVAROC



FPR vs TPR

Se varía el umbral



Diagonal segments are produced by ties.

C U R V A

R O C



