

Solución

1-)

### • Clasificador Naive Bayes

Basado en el modelo de Bayes asume que las características son independientes entre sí y que siguen una distribución normal (o sea) gaussiana

$$P(X, A) = P(X|A)P(A) = \prod_{j=1}^p P(x_j|A)P(A)$$

Se asume una distribución gaussiana para

$$P(x_k|A_i) = \frac{1}{\sqrt{2\pi} \sigma_{k,i}} \exp\left(-\frac{(x_k - \mu_{k,i})^2}{2\sigma_{k,i}^2}\right)$$

La función de predicción se define como:

$$\hat{y} = f(X) = \operatorname{argmax}_A \prod_{j=1}^p P(x_j|A_c)P(A_c)$$

Naive Bayes no posee un problema de optimización específico sin embargo se basa en estimar medias y varianzas

• Parámetros

$$\mu_{ki} = \frac{1}{|N_A|} \sum_{i=1}^{N_A} x_i$$

$$\sigma_{A,i}^2 = \frac{1}{N_A} \sum_{j=1}^{N_A} (x_j^{(i)} - \mu_{A,j})^2$$

## • SGD (gradient descent) classifier

consideramos la optimización estocástica donde se minimiza el valor esperado

$$L(\theta) = E_{q(z)} [L(\theta, z)]$$

en cada iteración obtenemos observaciones  $L_+(\theta) = L(\theta, z_+)$

$z_+ \sim q$ , obtenemos la estimación del gradiente de  $L$ ;  $g(z)$   
es independiente de los parámetros como  $g_+ = \nabla_{\theta} L_+(\theta_+)$   
el algoritmo puede escribirse así

$$\theta_{t+1} = \theta_t - \eta_t \nabla L(\theta_t, z_t) = \theta_t - \eta_t g_+$$

Para regresión lineal modelo  $(x_n^T \theta - y_n)$

la función objetivo (costo) tiene la forma

$$L(\theta) = \frac{1}{2N} \sum_{n=1}^N (x_n^T \theta - y_n)^2 = \frac{1}{2N} \|X\theta - y\|_2^2$$

el gradiente

$$g_+ = \frac{1}{N} \sum_{n=1}^N (b_+^T x_n - y_n) x_n$$

la actualización de parámetros  $\theta$  basada en el gradiente  
(Así resolvemos el problema de minimización de la función de pérdida  
SGD)

$$\theta_{t+1} = \theta_t - \eta_t (b_+^T x_n - y_n) x_n$$

## • Regresión logística

Utilizamos función lineal por tener decisión

$$f(x) = w^T x + b$$

La regresión logística transforma la salida de la función lineal  $f(x)$  en una probabilidad utilizando la función sigmoide

$$\sigma(f(x)) = \frac{1}{1 + e^{-f(x)}}$$

La función de pérdida logística

$$J(w, b) = \frac{1}{N} \sum_{i=1}^N [y_i \log(\sigma(f(x_i))) + (1 - y_i) \log(1 - \sigma(f(x_i)))]$$

por lo que el problema de optimización consiste en minimizar la función de pérdida logística

$$\min_{w, b} (L(w, b) + \lambda R(w))$$

y utilizamos los algoritmos de optimización SGD, gradiente descendente o métodos como newton-raphson para minimizar  $w$  y  $b$  de la función logística

$$w \leftarrow w - \eta \nabla_w J(w, b)$$

$$b \leftarrow b - \eta \nabla_b J(w, b)$$

## • Linear Discriminant Analysis (LDA)

busca una proyección lineal de las características originales  
que maximice la separabilidad entre clases

$$P(y=k|X) = \frac{\exp(-1/2 (X - \mu_k)^T \Sigma^{-1} (X - \mu_k))}{(2\pi)^{d/2} |\Sigma|^{1/2}}$$

$\mu_k$  es la media de la clase  $k$ ,  $\Sigma$  es la matriz de covarianza  
común y  $d$  la dimensión del espacio

El objetivo es encontrar un hiperplano que maximice la separación  
entre clases

$$w^* = \underset{w}{\operatorname{argmax}} \frac{w^T S_B w}{w^T S_W w}$$

y su solución

$$S_W^{-1} S_B w = \lambda w$$

## • K-Neighbors Classifier

Dado un conjunto de datos  $\{(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)\}$   
Se calcula una distancia  $d(x_{\text{new}}, x_i)$ , por ejemplo usando  
métrica euclidiana

$$d(x_{\text{new}}, x_i) = \left( \sum_{j=1}^p |x_{\text{new}, j} - x_{i,j}|^q \right)^{1/q}$$

Se hacen  $k$  vecinos cercanos  $x_{\text{new}}$

$$N_k(x_{\text{new}}) = \{x_i \mid x_i \text{ entre los } k \text{ puntos más cercanos a } x_{\text{new}}\}$$



la clase  $y_{new}$  se asigna a la clase que nos da más votos entre los  $k$  vecinos

$$y_{new} = \underset{c \in \{1, 2, \dots, C\}}{\operatorname{argmax}} \sum_{i \in N_k(x_{new})} \mathbb{I}(y_i = c)$$

$\mathbb{I}$  → función indicadora 
$$\begin{cases} 1 & \text{si } y_i = c \\ 0 & \text{en otro caso} \end{cases}$$

el conjunto de error para un valor dado de  $k$

$$\operatorname{error}(k) = \frac{1}{n} \sum_{i=1}^n \mathbb{I}(y_i \neq \hat{y}_i(k))$$

donde  $\hat{y}_i(k)$  es la clase predicha el valor óptimo de  $k$  es

$$k^* = \underset{k \in \mathbb{N}}{\operatorname{argmin}} \operatorname{error}(k)$$

## • SVC

Se debe encontrar un hiper-plano definido por

$$w^T x + b = 0 \quad \text{y } y_n \in \{-1, +1\}$$

es un problema de optimización cuadrática (con restricciones) el objetivo es minimizar la norma del vector  $w$  para maximizar el margen

$$\min_{w, b} \frac{1}{2} \|w\|_2^2 \quad \text{s.t.}, \quad y_n (w^T x_n + b) \geq 1 \quad \forall n \in \{1, 2, \dots, n\}$$

la regularizaci3n de problemas presenta un trade-off entre (des)no ser posible encontrar la mayor direcci3n

$$\min_{w, b} = \frac{1}{2} \|w\|_2^2 + C \sum_{n=1}^N \xi_n$$

$$\text{s.t. } y_n (f(x_n) - 1) \geq 1 - \xi_n$$

$$\xi_n \geq 0$$

el problema dual se obtiene introduciendo multiplicadores de Lagrange

$$\max_{\delta} \sum_{i=1}^N \delta_i - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \delta_i \delta_j y_i y_j x_i^T x_j$$

- $\delta_i$  multiplicadores de Lagrange
- $x_i^T x_j$  es el producto interno

el problema no lineal por SVC se utiliza kernel trick con funciones de decisi3n

$$f(x_{\text{new}}) = \text{Sign} \left( \sum_{i=1}^N \delta_i y_i K(x_i, x_{\text{new}}) + b \right)$$

### • Random forest classifier

Se propuso se crea un le generaci3n de fronteras mediante unbolizaci3n sobre las caracteristicas

• consiste en un lenguaje arbol de decisi3n  $\{h_1, h_2, \dots, h_N\}$

el entrenamiento de  $h$  se realiza seleccionando un subconjunto de  $n'$  muestras ( $n' \leq n$ ) del conjunto  $D$ , en cada divisi3n del arbol la predicci3n final se obtiene por votaci3n mayoritaria

$$y_{\text{new}} = \arg \max_{c \in \{1, 2, \dots, C\}} \sum_{f=1}^F \mathbb{I}(h_f(x_{\text{new}}) = c)$$

lo creas de discurrir sobre su utilidad en la vida  
Oni imparty

$$G_i = 1 - \sum_{k=1}^n p_{jk}^2$$

### • Gaussian process classifier

funcion latente que sigue el proceso gaussiano

$$f(x) \sim GP(m(x), k(x, x'))$$

$m(x)$  funcion de media  $m(x) = 0$   
 $k(x, x')$  kernel

la probabilidad pertenece a la clase 1, se utiliza funcion de enlace link function

$$P(y=1|f(x)) = \sigma(f(x)) = \sigma = \frac{1}{1 + e^{-f(x)}}$$

el objetivo es encontrar distribucion posterior de  $f(x)$   
en tanto a maximice la verosimilitud

$$p(y|x) = \int p(y|f)p(f|x)df$$

lo hiper parametros del kernel  $k(x, x')$  se optimizan maximizando la verosimilitud marginal

$$\theta^* = \arg \max_{\theta} p(y|x, \theta)$$

### Post data

la IA fue utilizada para explicar de manera general los modelos  
y adicionalmente permitian elegir y plantear los modelos  
como Gaussian process y k-NN como el Random forest regressive  
ademas de implementarlos con herramientas nuevas

IA's usadas Deepseek, ChatGPT