

# TALLER 29 may - 5 jun

## CV Clasificación LFW:

### 1) Naive Bayes Gaussian:

Red Bayesiana simple, donde  $Y(\text{Clase})$  es el nodo padre de un conjunto  $X_i$  que representa una evidencia. Si tenemos un vector de evidencia  $x = (x_1, x_2, \dots, x_d)$  determina la prob. de que pertenezcan a una de las clases  $y$ , valores que puede tomar  $y$ .

Distribución Gaussiana de datos. Por medio de train estimamos  $\mu_y$  y  $\sigma_y^2$  de cada clase.

Se empieza guardando Priors en una clase, preestablecidos o estimados con clases. Los Priors responden a las tablas de prob. para  $Y, p(Y)$ . También se estima tabla condicional para cada  $X_i$ , que tienen padre en  $Y$ . Así  $P(X_i|Y)$ , gaussianas por eso:  $\mu_y$  y  $\sigma_y^2$ ; así:

$$p(X_i = x_i | Y = y) = \frac{1}{\sqrt{2\pi}\sigma_y^2} e^{-\frac{(x_i - \mu_y)^2}{2\sigma_y^2}}$$

Probabilidad conjunta:

$$P(Y = y, X_1 = x_1, \dots, X_d = x_d) = P(Y = y) \prod_{i=1}^d P(X_i = x_i | Y = y)$$

Decidir sobre clase, se elegirá la clase que maximice la probabilidad estimada:

$$\hat{y} = \underset{y}{\operatorname{argmax}} P(Y = y, X_1 = x_1, \dots, X_d = x_d)$$

### SGD Classifier:

Un método de optimización. Para minimizar un loss, SGD actualiza los parámetros de forma iterativa.

Dataset:  $\{(x_i, y_i)\}_{i=1}^N$ ,  $x_i \in \mathbb{R}^d$ ,  $y_i \in \{-1, 1\}$

Objetivo encontrar un vector de pesos  $w \in \mathbb{R}^d$  y bias  $b \in \mathbb{R}$  que minimice la pérdida logística regularizada:

- Logistic Loss (Clasificación binaria):

$$P(y_i | x_i) = \frac{1}{1 + e^{(-y_i(w^T x_i + b))}}$$

$$\mathcal{L}(y_i, w^T x_i + b) = \log(1 + e^{(-y_i(w^T x_i + b))})$$

- Función objetivo con regularización: Prevenir sobreajuste

$$J(w, b) = \frac{1}{N} \sum_{i=1}^N \log(1 + e^{(-y_i(w^T x_i + b))}) + \frac{\lambda}{2} \|w\|^2$$

- Cálculo de gradiente:

$$\nabla_w J = -\frac{1}{N} \sum_{i=1}^N \frac{y_i x_i}{1 + \exp(y_i(w^T x_i + b))} + \lambda w$$

Para el bias  $b$ :

$$\frac{\partial J}{\partial b} = -\frac{1}{N} \sum_{i=1}^N \frac{y_i}{1 + \exp(y_i(w^T x_i + b))}$$

- Descenso de gradiente estocástico:

No calcula el gradiente completo, 1 ejemplo a la vez: para  $(x_+, y_+)$ , es:

$$\begin{aligned} w &\leftarrow w - \eta_+ \left( \frac{-y_+ x_+}{1 + \exp(y_+(w^T x_+ + b))} + \lambda w \right) \\ b &\leftarrow b - \eta_+ \left( \frac{-y_+}{1 + \exp(y_+(w^T x_+ + b))} \right) \end{aligned}$$

← tasa de aprendizaje en t

SGD resuelve el siguiente problema de regresión logística regularizada:

$$\min_{w, b} \frac{1}{N} \sum_{i=1}^N \log(1 + \exp(-y_i(w^T x_i + b))) + \frac{\lambda}{2} \|w\|^2$$

Optimiza iterativamente, actualizando pesos después de cada muestra

### Logistic Regression: Casando SGD

Dataset:  $\{(x_i, y_i)\}_{i=1}^N$ ; input  $x_i \in \mathbb{R}^d$ ; labels  $y_i \in \{0, 1\}$

- Logistic model (Sigmoid Activation)

Probabilidad de class 1:

$$P(y = 1 | x) = \sigma(w^T x + b) = \frac{1}{1 + \exp(-w^T x - b)}$$

$$\hat{y} = \begin{cases} 1, & \sigma(w^T x + b) > 0.5 \\ 0, & \text{e.o.c} \end{cases}$$

- Log-likelihood (Forma maximizada)

Likelihood para dataset:

$$\mathcal{L}(w, b) = \prod_{i=1}^N \sigma(w^T x_i + b)^{y_i} (1 - \sigma(w^T x_i + b))^{1-y_i}$$

Log-Likelihood:

$$\log(\mathcal{L}(w, b)) = \sum_{i=1}^N (y_i \log(\sigma(w^T x_i + b)) + (1 - y_i) \log(1 - \sigma(w^T x_i + b)))$$

- Forma minimizada (Negative Log-likelihood):

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

- Vector con Regularización:

$y_i \in \{-1, 1\}$ , simplifica la expresión loss

$$J(w, b) = \sum_{i=1}^N \log(1 + \exp(-y_i(w^T x_i + b)))$$

$$\min_{w, b} \frac{1}{N} \sum_{i=1}^N \log(1 + \exp(-y_i(w^T x_i + b))) + \frac{\lambda}{2} \|w\|^2$$

←  $\lambda$  regularización

### Linear Discriminant Analysis (LDA):

- Distribuciones de clase:

$$x | y = c \sim N(\mu_c, \Sigma)$$

Covarianza en común

Priors de clase  $P(y = c) = \pi_c$

• • •

Regla de Bayes y función discriminante:

$$P(y=c|x) = \frac{P(x|y=c)\pi_c}{P(x)}$$

Log-posterior lleva a la función discriminante:

$$\delta_c(x) = x^T \Sigma^{-1} \mu_c - \frac{1}{2} \mu_c^T \Sigma^{-1} \mu_c + \log \pi_c$$

$$\hat{y} = \arg \max_c \delta_c(x)$$

Estimación de parámetros:

Media de clase:

$$\hat{\mu}_c = \frac{1}{N_c} \sum_{i: y_i=c} x_i$$

Covarianza:

$$\hat{\Sigma} = \frac{1}{N-K} \sum_{c=1}^K \sum_{i: y_i=c} (x_i - \hat{\mu}_c)(x_i - \hat{\mu}_c)^T$$

Priors de clase:

$$\hat{\pi}_c = \frac{N_c}{N}$$

Optimización:

sobre direcciones  $E_n \in \mathbb{R}^d$  que maximizan la separación de clases:

$$\max_W \frac{W^T S_B W}{W^T S_W W}$$

Dispersión entre clases:

$$S_B = (\mu_1 - \mu_0)(\mu_1 - \mu_0)^T$$

Dispersión dentro de la clase:

$$S_W = \Sigma_0 + \Sigma_1$$

Problema generalizado de valores propios: Cociente de Rayleigh, que maximiza en:

$$W \propto S_W^{-1} (\mu_1 - \mu_0)$$

Para multiclases:

$$S_B W = \lambda S_W W$$

$$\max_W \frac{W^T (\mu_1 - \mu_0)(\mu_1 - \mu_0)^T W}{W^T \Sigma W} \rightarrow W = \Sigma^{-1} (\mu_1 - \mu_0)$$

LDA multiclase generaliza esto mediante vectores propios

$$\Sigma_W^{-1} S_B$$

## K Neighbors Classifier:

Busca predecir  $\hat{y}$  encontrando la etiqueta mayoritaria entre los k vecinos más cercanos de x en train.

$$\{(x_i, y_i)\}_{i=1}^N; x_i \in \mathbb{R}^d; y_i \in \{1, \dots, C\}$$

$$\text{test } x \in \mathbb{R}^d$$

...

Calculo de distancia:

$$D(x, x_i) = \|x - x_i\|_2 = \sqrt{(x - x_i)^T (x - x_i)}$$

k vecinos cercanos  $\{N_k(x) \subset \{1, \dots, N\}\}$

$N_k(x) \rightarrow$  Más pequeño  $D(x, x_i)$  de k

Regla predictiva (mayoría)

$$\hat{y} = \arg \max_{c \in \{1, \dots, C\}} \sum_{i \in N_k(x)} 1(y_i = c)$$

KNN soportando (opcional):

$$\hat{y} = \arg \max_c \sum_{i \in N_k(x)} \frac{1}{D(x, x_i)} \cdot 1(y_i = c)$$

Cálculo de costo (al predecir):

$$\arg \min_{N_k(x)} \sum_{i \in N_k(x)} D(x, x_i)$$

## Linear SVC:

Busca un hiperplano  $w^T x + b = 0$  que separe los 2 clases con el máximo margen.

Loss function hinge:

Penaliza los puntos mal clasificados y que violan el margen:

$$\mathcal{L}(y_i, f(x_i)) = \max(0, 1 - y_i(w^T x_i + b))$$

Loss es 0 cuando  $y_i(w^T x_i + b) \geq 1$

Riesgo empírico regularizado (Primordial): función de costo:

$$J(w, b) = \frac{1}{2} \|w\|^2 + C \sum_{i=1}^N \max(0, 1 - y_i(w^T x_i + b))$$

$\|w\|^2$  impone la maximización de márgenes (pesos pequeños  $\rightarrow$  Margen grande)

C es un hiperparámetro que controla el equilibrio entre errores de margen y clasificación.

Optimización (primordial):

$$\min_{w, b} \frac{1}{2} \|w\|^2 + C \sum_{i=1}^N \max(0, 1 - y_i(w^T x_i + b))$$

## SVC:

función de decisión:

$$f(x) = \text{Sign} \left( \sum_{i=1}^N \alpha_i y_i \overset{\text{kernel function}}{K(x_i, x)} + b \right)$$

multiplicadores lagrange

Problema primordial (SVM margen suave):

$$\min_{w, b, \xi} \frac{1}{2} \|w\|^2 + C \sum_{i=1}^N \xi_i$$

Parámetro para misclasificación

Característica de transformación via kernel

$$\text{Sujeto a } y_i(w^T \phi(x_i) + b) \geq 1 - \xi_i, \xi_i \geq 0$$

- Problema de derivación dual =

$$\max_{\alpha} \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i,j=1}^N \alpha_i \alpha_j y_i y_j K(x_i, x_j)$$

Sujeto a  $0 \leq \alpha_i \leq C$

$$\sum_{i=1}^N \alpha_i y_i = 0$$

- Función de decisión final =

$$f(x) = \text{sign} \left( \sum_{i=1}^N \alpha_i^* y_i K(x_i, x) + b \right)$$

## Random Forest Classifier =

Tenemos  $M$  árboles de decisión =

$$\hat{y} = \text{majority\_vote}(T_1(x), T_2(x), \dots, T_M(x))$$

Cada árbol =

- Muestreo de datos con sustitución
- Seleccionar un subconjunto aleatorio de entidades en cada división
- Cultivar un árbol usando criterios codiciosos (Gini, entropía)

- Bootstrap Sample:

Para training  $D = \{(x_i, y_i)\}_{i=1}^N$ , genera un sample  $D_m$  de tamaño  $N$

- Crecimiento de árboles =

Seleccionar atributo  $j$  y dividir  $S$ , tal que se divida en left/right maximizando la separación de clase.

- Criterio de división =

$$\max_{j,s} \Delta I = I(D) - \left( \frac{|D_{\text{left}}|}{|D|} I(D_{\text{left}}) + \frac{|D_{\text{right}}|}{|D|} I(D_{\text{right}}) \right)$$

- $I()$ : Medida de impurezas (Gini, Entropía)
- $D_{\text{left}}, D_{\text{right}}$ : divide
- Intenta todos los  $s$  en  $j$  seleccionados

Gini impureza =

$$I(D) = 1 - \sum_c p_c^2, \quad p_c = \frac{1}{|D|} \sum_{i \in D} 1(y_i = c)$$

Entropía:

$$I(D) = - \sum_c p_c \log p_c$$

- Repetir recursivamente:

Repetir hasta que:

- Máximo de profundidad
- Mínimas samples por nodo
- Hojas puras (todos los samples de misma clase)

$$\hat{y}(x) = \arg \max_c \sum_{m=1}^M 1(T_m(x) = c)$$