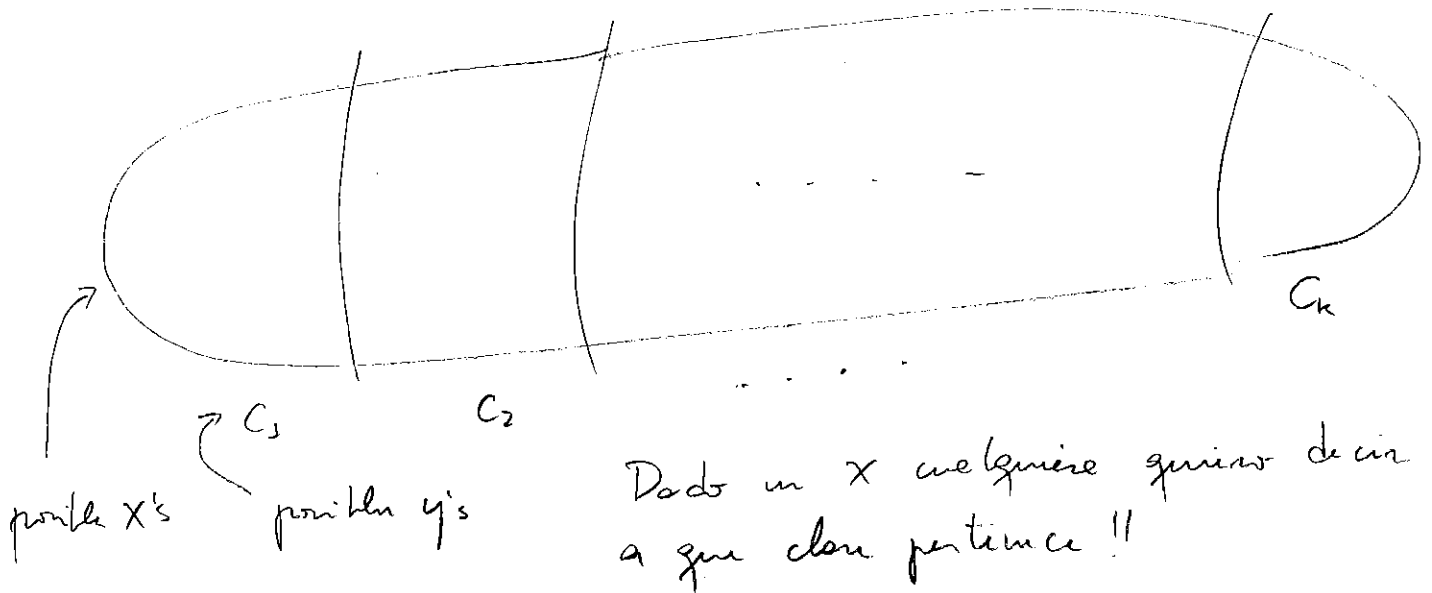


3a) Estimador de máxima verosimilitud, Información y Entropía. ①

Supongamos un universo de ejemplos, clasificados en k clases



Dado un x cualquiera quisiera decir a qué clase pertenece!!

Si no tengo acceso a todos los datos, sólo me queda mirar el experimento como si los datos son generados en un proceso estocástico. Quiero aproximar entonces a $P_{\text{real}}(C_i | x)$

Heure usado redes neuronales para hacer esto

$$P_{NN}(C_j | x) = \hat{y}_j = \text{forward}_{NN}(x)_j$$

todo esto depende de los parámetros θ de la red

\Rightarrow cuento con $P_{\theta}(C_j | x)$ computado por mi red para cada posible input.

Supongo que cuento con un conjunto de ejemplos

$$X = \{(x^{(1)}, s^{(1)}), (x^{(2)}, s^{(2)}), \dots, (x^{(N)}, s^{(N)})\} \quad s^{(i)} \in \{C_1, C_2, \dots, C_k\}$$

¿Cuál es la probabilidad que mi red le origine a observar exactamente estos datos? (si estos son generados i.i.d) la probabilidad es

$$\prod_{i=1}^N P_{\theta}(s^{(i)} | x^{(i)})$$

← probabilidad total de que un θ explique todos los datos observados. Me gustaría que esta probabilidad sea tan grande como sea posible.

Un buen (el mejor) estimador para los parámetros, es el que maximiza la anterior expresión

$$\theta_{ML} = \underset{\theta}{\operatorname{argmax}} \underbrace{\prod_{i=1}^N P_{\theta}(s^{(i)} | x^{(i)})}_{\text{likelihood}}$$

← Estimador de máxima verosimilitud para un parámetro, maximum likelihood

c) Por qué es bueno?

Supongamos que existe un $\hat{\theta}$ y $P_{\hat{\theta}} = P_{\text{real}}$ entonces

$$1) \text{ si } N \rightarrow \infty \Rightarrow \theta_{ML} \rightarrow \hat{\theta} \Rightarrow P_{\theta_{ML}} \rightarrow P_{\text{real}}$$

$$\lim_{N \rightarrow \infty} P_{\theta_{ML}} = P_{\text{real}} \quad \leftarrow \text{consistencia}$$

2) dentro de todos los estimadores consistentes θ_{ML} es el que converge más rápido a $\hat{\theta}$ ← eficiencia

(1) y (2) hacen que θ_{ML} sea el estimador preferido en machine learning.

Log likelihood: las siguientes son formulaciones equivalentes para θ_{ML}

$$\theta_{ML} = \underset{\theta}{\operatorname{argmax}} \prod_{i=1}^N P_{\theta}(s^{(i)} | x^{(i)}) = \underset{\theta}{\operatorname{argmax}} \log \left(\prod_{i=1}^N P_{\theta}(s^{(i)} | x^{(i)}) \right)$$

$$\theta_{ML} = \underset{\theta}{\operatorname{argmax}} \sum_{i=1}^N \log (P_{\theta}(s^{(i)} | x^{(i)}))$$

$$\theta_{ML} = \underset{\theta}{\operatorname{argmax}} \frac{1}{N} \sum_{i=1}^N \log (P_{\theta}(s^{(i)} | x^{(i)}))$$

$$\theta_{ML} = \underset{\theta}{\operatorname{argmax}} \mathbb{E}_{(s,x) \in X} (\log (P_{\theta}(s^{(i)} | x^{(i)})))$$

es como una pérdida (que queremos minimizar)

Me gustaría que de alguna forma capturemos esta forma de calcular los mejores parámetros. Para darle una interpretación usamos Teo. Información.

enfaticar después que esta formulación es el soporte teórico del entrenamiento por paquetes

Para darle una interpretación usamos Teo. Información.

3b) Teor. Información

Shannon (interés en comunicación):

- La información de un evento debe ser mayor mientras menos probable es el evento (evento con probabilidad 1 entrega 0 información).
- La información debe ser siempre positiva.
- La información de eventos independientes debe ser aditiva.

$$\Rightarrow \text{Inf}(x) = \log\left(\frac{1}{P(x)}\right) \quad \leftarrow \begin{array}{l} \text{información de un evento } x \\ \text{es el logaritmo del inverso de} \\ \text{su probabilidad.} \end{array}$$

La "entropía" de una distribución de probabilidad (o de un experimento) es el promedio ponderado de la información entregada por todos los eventos.

$$\text{Entropía de } P: H(P) = \sum_x P(x) \cdot \text{Inf}(x)$$

$$= \sum_x P(x) \cdot \log\left(\frac{1}{P(x)}\right)$$

menor log
mayor prob

$\text{Inf}(x)$ se puede interpretar como la cantidad mínima de bits que necesitamos para representar a x si quisieramos copiarlo de diferencia con los otros eventos $\Rightarrow H(P) = \mathbb{E}_{x \sim P}(\text{Inf}(x))$ cantidad promedio de bits que necesitamos para codificar los eventos de mi distribución

Imaginemos que los datos distribuyeron según P pero los codificamos según una distribución enorme Q ; ¿Cuántos bits estaríamos perdiendo?

$$\sum P(x) \log\left(\frac{1}{Q(x)}\right) - \sum P(x) \log\left(\frac{1}{P(x)}\right) = D_{\text{KL}}(P, Q)$$

Bits con la distribución enorme

Bits originales

(Kullback-Leibler)

divergencia

$D_{KL}(P, Q)$ es una buena medida de error y vemos que está íntimamente relacionada a \mathcal{O}_{ML}

$$D_{KL}(P, Q) = \underbrace{\sum P(x) \log\left(\frac{1}{Q(x)}\right)}_{H(P, Q)} - \underbrace{\sum P(x) \log\left(\frac{1}{P(x)}\right)}_{H(P)}$$

↖ Para minimizar, para Q , $D_{KL}(P, Q)$ tiene un mínimo cuando $Q = P$

↗ Entropía cruzada entre P y Q
 si tenemos P fijo y queremos minimizar $D_{KL}(P, Q)$
 \Rightarrow lo mismo con minimizar $H(P, Q) = \sum P(x) \log\left(\frac{1}{Q(x)}\right)$

Volviendo ahora a nuestra red con ejemplos $X = \{(x^{(i)}, s^{(i)})\}_{i=1}^N$

para el ejemplo $(x^{(i)}, s^{(i)})$ tengo dos distribuciones de probabilidad

$$P_{\theta}^{(i)}(c_j) = P_{\theta}(c_j | x^{(i)})$$

y la probabilidad empírica

$$P_{emp}^{(i)}(c_j) = P_{emp}(c_j | x^{(i)}) = \begin{cases} 1 & \text{si } c_j = s^{(i)} \\ 0 & \text{si } c_j \neq s^{(i)} \end{cases}$$

$$\mathcal{L}' = \frac{1}{N} \sum_{i=1}^N D_{KL}(P_{emp}^{(i)}, P_{\theta}^{(i)})$$

queremos minimizar \mathcal{L}' respecto a θ

\Rightarrow Para encontrar lo mismo que minimizar

$$\mathcal{L} = \frac{1}{N} \sum_{i=1}^N H(P_{emp}^{(i)}, P_{\theta}^{(i)})$$

$$\begin{aligned} \mathcal{O}_{ML} &= \arg \min_{\theta} \mathcal{L} \\ &= \arg \min_{\theta} \frac{1}{N} \sum_{i=1}^N H(P_{emp}^{(i)}, P_{\theta}^{(i)}) \end{aligned}$$

cross-entropy loss $\sum_j \underbrace{P_{emp}(c_j | x^{(i)}) \log\left(\frac{1}{P_{\theta}(c_j | x^{(i)})}\right)}_{\substack{= 0 \text{ si } c_j \neq s^{(i)} \\ = 1 \text{ si } c_j = s^{(i)}}}$

El optimizador de máxima verosimilitud obtiene de minimizar la cross-entropy

$$\Rightarrow \mathcal{L} = \frac{1}{N} \sum_{i=1}^N \log\left(\frac{1}{P_{\theta}(s^{(i)} | x^{(i)})}\right) = -\frac{1}{N} \sum_{i=1}^N \log(P_{\theta}(s^{(i)} | x^{(i)}))$$