



General Consulate of the Republic of Sierra Leone

Teoria ML

Desigualdad de Hoeffding:

$$P(Z_i = 1) = \phi \quad P(Z_i \neq 1) = 1 - \phi$$

$$\hookrightarrow P(|\phi - \hat{\phi}| > \delta) \leq 2 \exp(-2\delta^2 m)$$

\hookrightarrow Conc. Valores

\hookrightarrow Se acerca al valor real de ϕ

II) Error de entrenamiento

$$\hat{\epsilon}(h) = \frac{1}{m} \sum \mathbb{1}\{h(x_i) \neq y_i\}$$

Error de generalización

$$\epsilon(h) = P(h(x) \neq y) \quad x \in D$$

III) Underfitting $\rightarrow \hat{\epsilon}$ alto
Overfitting $\rightarrow \hat{\epsilon}$ bajo y ϵ alto

IV) h def por ϕ

$$\hat{\phi} = \operatorname{argmin} \hat{\epsilon}(h_{\phi}) \rightarrow \text{minimizar } \hat{\epsilon}$$

V) $H = \{h_1, \dots, h_k\}$ hipótesis $h_i(x) \in \{0, 1\}$

$$\hat{\epsilon} \quad Z_j = \mathbb{1}\{h_i(x_j) \neq y_j\}$$

$$\hat{\epsilon}_{(hi)} = \frac{1}{m} \sum_{j=1}^m Z_j$$

$$\epsilon \quad Z = \mathbb{1}\{h_i(x) \neq y\}$$

\downarrow

$$\epsilon = P(Z=1)$$

Córdoba 1233-44 floor Tel. 40-9573-2022-1806 - 46-9782. Buenos Aires, Argentina

$$P(|\epsilon(h_i) - \hat{\epsilon}(h_i)| > \delta) \leq 2 \exp(-2\delta^2 m)$$

VI) Convergencia Uniforme



$$P(\exists h \in H / |\mathcal{E}(h) - \hat{\mathcal{E}}(h)| > \delta)$$

$$P(A_1 \cup A_2 \cup \dots \cup A_K)$$

$$\leq \sum_{i=1}^K P(A_i)$$

$$\leq \sum_{i=1}^K 2 \exp(-2\delta^2 m) = 2K \exp(-2\delta^2 m) \rightarrow \text{Prob } \hat{\mathcal{E}} \neq \mathcal{E} \geq \delta$$

VII) Despejamos m

$$m \gg \frac{1}{2\delta^2} \log \frac{2K}{\delta}$$

Complejidad de la muestra.

VIII)

$$h^* = \underset{h \in H}{\operatorname{argmin}} \mathcal{E}(h) \rightarrow \text{Hipotesis Optima}$$

$$\mathcal{E}(\hat{h}) \leq \hat{\mathcal{E}}(\hat{h}) + \delta$$

$$\hat{\mathcal{E}}(\hat{h}) \leq \hat{\mathcal{E}}(h^*)$$

$$\hookrightarrow \mathcal{E}(\hat{h}) \leq \hat{\mathcal{E}}(h^*) + \delta$$

$$\mathcal{E}(h^*) - \hat{\mathcal{E}}(h^*) \leq \delta$$

$$\hookrightarrow \hat{\mathcal{E}}(h^*) \leq \mathcal{E}(h^*) + \delta$$

$$\hookrightarrow \boxed{\mathcal{E}(\hat{h}) \leq \mathcal{E}(h^*) + 2\delta}$$

$$\text{IX) } \boxed{\mathcal{E}(\hat{h}) \leq \min_{h \in H} \mathcal{E}(h) + 2 \sqrt{\frac{1}{2m} \log \frac{2K}{\delta}}}$$

\hookrightarrow TEOREMA DE LA GENERALIZACIÓN

X) Dimension VC

$$VC(H)$$

$$VC(\mathcal{O}_x + \mathcal{O}_1) = 3$$

\hookrightarrow Cardinalidad máxima de puntos en 2D

Teo Radon: Si tenemos $d+2$ puntos en \mathbb{R}^d podemos particionarlos en 2 conjuntos tal que sus fronteras convexas se intersecten \rightarrow no separables linealmente

$$\hookrightarrow VC = d+1$$



General Consulate of the Republic of Sierra Leone

TEOREMA HART

↳ Si tenemos m puntos en d dims. Existe con prob alta una f que proyecta \mathbb{R}^d en \mathbb{R}^k $k \gg d$ de forma tal que los puntos sean linealmente inde separables.

TEOREMA VAPNIK: Sea H una flia de Hipotesis $VC(H) = d \Rightarrow$

con prob al menos $1 - \delta \quad \forall h \in H$

$$|\mathcal{E}(h) - \hat{\mathcal{E}}(h)| \leq O \sqrt{\frac{d}{m} \log \frac{m}{d} + \frac{1}{m} \log \frac{1}{\delta}}$$

Cor: $|\mathcal{E}(h) - \hat{\mathcal{E}}(h)| \leq \epsilon$

$\forall h \in H$ con prob $1 - \delta$ es suficiente que $m = O(d)$

$$d \approx |\mathcal{H}|$$



General Consulate of the Republic of Sierra Leone

CLUSTERING

• Clustering Jerarquico:

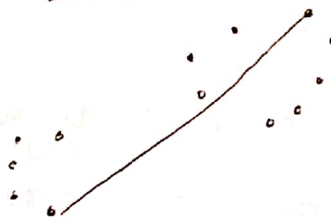
- 1) Partimos de que cada punto es un cluster.
- 2) En cada paso, unimos los dos cluster mas cercano.
- 3) Repetir hasta que quede uno solo o la cantidad buscada.

Distancias

- Single-linkage



- Complete-linkage



- Average-link



- Medion-link



• OPTIMIZACIÓN LSH Single Linkage

↳ Aproximación

↳ En cada bucket habra como maximo un punto de cada cluster

↳ Los clusters en un mismo bucket son candidatos

Cordoba 1233-4th floor - Tel. 40-9573-2022-1806 - 46-9782 - Buenos Aires, Argentina

K-MEANS

↳ Dado un set de datos en \mathbb{R}^d , particionar en K sets ($K \leq D$) de forma tal de minimizar la sumatoria de las distancias al centroide.

$$\sum_{i=1}^K \sum_{x \in S_i} \|x - \mu_i\| \rightarrow \text{NP-HARD}$$

algoritmo

- 1) Elegir K puntos, uno para cada cluster, centroides ^{iniciales}
- 2) Asignarle a cada punto el cluster con centroide ^{inicial} mas cercano
- 3) Recalcular el centroide como el centro de los puntos
- 4) Volver a 2

K-MEANS ++

- 1) Elegir un punto al azar como centroide
- 2) Calcular la distancia de cada punto contra los centroides y quedarse con la minima.
- 3) Calcular la probabilidad de cada punto como la distancia dividida por la suma de todas las distancias de los puntos
- 4) Elegir un punto al azar usando la nueva proba.
- 5) Repetir.

$$\|x_1 - y\|^2 + \|x_2 - y\|^2 + \|x_3 - y\|^2 = 0$$

$$\frac{x_1 - c_1}{\|x_1 - y\|} + \frac{x_2 - c_1}{\|x_2 - y\|} = \cancel{x_1 + x_2} \|x_1 - y\|$$



General Consulate of the Republic of Sierra Leone

STREAMING \leadsto FLUJO DE DATOS SIN FIN

Reservoir Sampling

- \rightarrow Guardar muestra de tamaño constante K
- $\rightarrow n$ cantidad de datos vistos
- \rightarrow Probabilidad de que un dato este en la muestra $\rightarrow \frac{K}{n}$
- \rightarrow Genero un $x \in [0, 1]$ aleatorio
- \rightarrow Si $x < \frac{K}{n}$, reemplazo uno aleatorio del stream

MOMENTOS DE STREAM

\rightarrow Decimos que la frecuencia de cada dato i es $M_i \rightarrow$ cant. veces visto

\rightarrow Momento de orden K del stream

$\rightarrow \sum M_i^K$

$K=0$ \rightarrow CANTIDAD DE ELEM DIFERENTES VISTOS

$K=1$ \rightarrow CANTIDAD DE ELEM OBSERVADOS

$K=2$ \rightarrow factor sorpresa \rightarrow indicador de distribución pareja.

FLAOLET-MARTIN $\rightarrow K=0$

\rightarrow APLICA FUNCIÓN DE HASH DE M bits, $M > \log_2 n$

\rightarrow OBSERVA CANTIDAD DE 0 a 129

$\rightarrow r = \text{MAX}(\uparrow)$

\rightarrow CANT ELEM DIFERENTES $\approx 2^r$

\rightarrow Requiere buen hash \rightarrow USAMOS VARIOS HASH EN GRUPOS.

Córdoba 1233-44 floor - Tel. 40-9573-2022-1806 - 46-9782 - Buenos Aires, Argentina

HyperLogLog

- ↳ USAMOS UN HASH DE 64 BITS
- ↳ Los primeros n bits son número de estimadores
- ↳ 64- n bits contamos 0 a izquierda y actualizamos el estimador si es mayor.
- ↳ La estimación es el promedio armónico de los estimadores

AMS → $k=2$

- ↳ se define un número k de estimadores
 - ↳ CV con 2 campos → VALOR
CANTIDAD
- ↳ Por cada elem
 - ↳ Si está en el valor de algún k estimador → cantidad ++
 - ↳ Si no se sortea si reemplaza a un estimador.
- ↳ Reservoir Sampling.

→ EL MOMENTO DE ORDEN 2 es el promedio de $n(2C_i - 1)$
↳ cantidad est. R_i

BLOOM FILTERS → SABER SI UN ELEM PERTENECE A UN SET

- ↳ Vector binario de m -bits y k funciones de hash de 0 a $m-1$
- ↳ Para construir el filtro, por cada elem del set le aplico los hash y enciendo en 1 las posiciones indicadas
- ↳ Para consultar, se aplican los hash y me fijo si todas las bits están en 1 → Pertenece con proba $(1-\delta)$ con δ P(Falso pos)
- ↳ Si no está encendido alguno, definitivamente no pertenece.

→ El valor OPTIMO de Hashes k para el vec de M bits en un universo de n elementos
↳ $k = \left(\frac{m}{n}\right) \log 2$

→ SABIENDO EL k , LA CANTIDAD DE M BITS
↳ P(FALSO POS)

$$m = \frac{n \ln p}{(\ln 2)^2}$$

→ ESTIMAR CUANTOS ELEM EN EL FILTRO $E = \frac{-(M \ln(1 - \frac{x}{M}))}{k}$
bits bits en 1



General Consulate of the Republic of Sierra Leone

Counting Filter

- ↳ en vez de bitmap, tengo un contador
- ↳ Nos permite eliminar elementos del filtro.
- ↳ Construimos como bloom, pero aumentando en 1 el contador.
- ↳ Consulta como bloom, > 0
- ↳ USA MAS ESPACIO

COUNT MIN SKETCH

- ↳ Basado en Counting Filter
- ↳ Se usan d filas de w bits asociado a una func. de hash.
- ↳ Observo elem, aplico d func. hash e incremento la pos en cada vez.
- ↳ Para estimar frecuencia, hashéo y me quedo con el contador min.
- ↳ Para mantener los i elementos mas frec, por cada elem actualizo y hago COUNT-MIN, si supera al menor de los i elementos me lo reemplazo
- ↳ Guardo los TOP- i en memoria

HEAVY HITTER PROBLEMS

- ↳ Encontrar valor mas repetido

CUCKOO FILTERS

- ↳ USA CUCKOO HASHING
- ↳ + RAPIDOS PARA CONSULTA
- ↳ NO ALMACENA CLAVES SI NO FINGERPRINT (6 a 8 bits)
- ↳ ALMACENAMOS TABLA DE CUCKOO
 - ↳ M buckets
 - ↳ b fingerprint por bucket.

SISTEMAS DE RECOMENDACIÓN \leadsto Recomendar Contenido al usuario

- \hookrightarrow PRECISIÓN
- \hookrightarrow Serendipity \rightarrow Cosas Nuevas
- \hookrightarrow Diversidad

SISTEMAS NO-PERSONALIZADOS \leadsto Para todos, no específicos.

- \hookrightarrow Recomendación de Comentarios
 - \hookrightarrow UPVOTES - DOWNVOTES \rightarrow NO PROPORCIÓN
 - \hookrightarrow $\frac{\text{UPVOTES}}{\text{UP} + \text{DOWN}}$ \leadsto ECUACIÓN MAS PELIGROSA

\hookrightarrow Intervalos de Confianza

- \hookrightarrow Dado los UP \uparrow y los DOWN \downarrow queremos saber con confianza del $x\%$ cual es el limite inferior para la proba de que el comentario sea POSITIVO.
- \hookrightarrow ORDENAMOS POR ESA PROBA.

Siendo U los UP \rightarrow Δ los down, ^{proporcion} proba upvote $P = \frac{U}{U+D}$ y $n = U+D$
USAMOS La formula de Wilson para intervalos de confianza:

z -VALUE PARA INTERVALO \hookrightarrow
$$\left(P + \frac{z^2}{2n} \pm z \sqrt{\left[P(1-P) + \frac{z^2}{4n} \right] / n} \right) \frac{1 + \frac{z^2}{n}}$$

\hookrightarrow No sirve para noticias

ORDENANDO NOTICIAS

\hookrightarrow función de utilidad U , P proba que le guste la noticia, q proba que no sea nuevo

$$u(P, q) = a * Pq + (1-P)q * b + P(1-q) * c + (1-P)(1-q) * d$$

Siendo $P = \frac{U+1}{U+D+2}$ y $q = e^{-\lambda \tau}$

$$\hookrightarrow u(U, D, A, B) = \log(U-D) + \frac{A-B}{45000}$$



General Consulate of the Republic of Sierra Leone

SISTEMAS BASADOS EN CONTENIDOS

- ↳ Recomendar basado en gustos previos
- ↳ Por cada item mantengo un profile
- ↳ Por cada usuario mantengo un profile

Calcular interces:

$$\cos \theta = \frac{x \cdot y}{\|x\| \|y\|}$$

Ej:

	1	2
ANIMATED	✓	x
MARVEL	x	✓
BECHDEL TEST	✓	✓

→ Perfil.

Collaborative Filtering

- ↳ Suponiendo n usuarios y m items, construyo matriz de utilidades.

	I1	I2	I3
U1	2	?	4
U2	?	3	4
U3	2	4	?

Valores a estimar

Dado el usuario X, busco los n usuarios mas similares y estimo base a eso

Calculo la semejanza entre usuarios como

$$\text{sim}(x, y) = \frac{\sum_{s \in S_{xy}} (r_{xs} - \bar{r}_x)(r_{ys} - \bar{r}_y)}{\sqrt{\sum_{s \in S_{xy}} (r_{xs} - \bar{r}_x)^2} \sqrt{\sum_{s \in S_{xy}} (r_{ys} - \bar{r}_y)^2}}$$

	P1	P2	P3	P4	P5	P6	PROMEDIO
U1	2		2	4	5		13/4
U2	5		4			1	10/3
U3			5		2		7/2
U4		1		5		4	16/3
U5			4			2	6/2
U6	4	5		1			10/3

→ Redo los PROMEDIOS a las calificaciones

Córdoba 1233-4th floor Tel. 40-9573-2022-1806 - 46-9782 Buenos Aires, Argentina

Para estimar la calificación de un usuario a un item uso los n vecinos más cercanos

$$r_{xi} = \frac{\sum_{y \in N} S_{xy} r_{yi}}{\sum_{y \in N} S_{xy}}$$

Cálculo por desviación:

$\hookrightarrow r_{ij} = \mu + \delta_i + \delta_j + \delta_{ij}$ donde:

μ = promedio de todas las calificaciones de todas las películas

δ_i = desviación usuario i respecto a μ

δ_j = desviación item j respecto a μ

δ_{ij} = desviación usuario i para película j calculado como:

$$\delta_{ij} = \frac{\sum_{j \in N} S_{ij} (r_{ij} - b_{xj})}{\sum_{j \in N} S_{ij}} \quad \text{y } b_{xj} = \mu + \delta_x + \delta_j$$

Machine Learning



General Consulate of the Republic of Sierra Leone

Page Rank

- ↳ Independiza el resultado del contenido del mismo
- ↳ Basado en la estructura de links de la web
- ↳ La importancia de un nodo viene de la importancia de b que lo apuntan

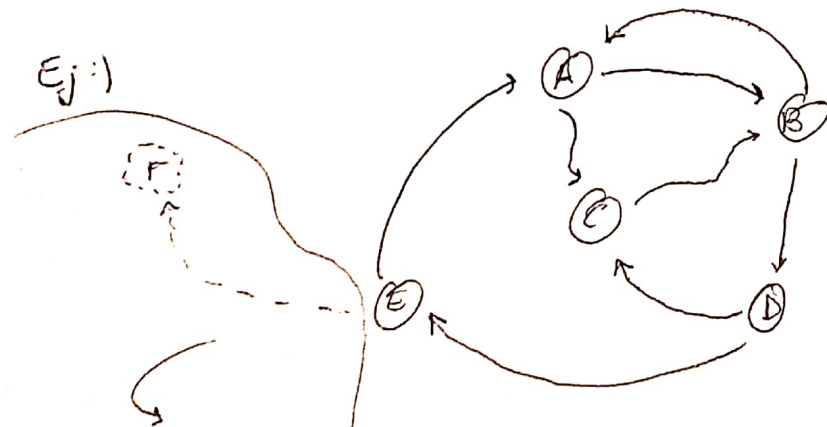
Algoritmo

- ↳ Asignar a cada nodo $PR = \frac{1}{n}$

Realizar k veces:

- ↳ Cada nodo reparte su PR en partes iguales a los que apunta
- ↳ El nuevo PR es la suma de los pedazos de PR recibidos

Ej:)



	A	B	C	D	E
<u>1</u>	1/5	1/5	1/5	1/5	1/5
<u>2</u>	3/10	3/10	2/10	1/10	1/10
	5/20	7/20	4/20	3/20	1/20
	6/25	8/25	5/25	4/25	2/25

Solución para grafos con dead-end

Repartimos en partes iguales la proba del nodo entre todos los demás
 (sumo $1/n$ de la probabilidad del nodo a cada nodo) → le reparto a todos, incluyendo

Otro Problema → Ciclos entre nodos que reparten entre ellos

↳ Introduzco parametro B tal que con proba B elegimos un link al azar y con proba $(1-B)$ se teletransporta a cualquier pagina al azar

En genl, $0,8 < B < 0,9$

TRUST RANK

- ↳ APLICAR TELETRANSPORTACIÓN ÚNICAMENTE HACIA PÁGINAS CONFIABLES
- ↳ CÁLCULO PR y TR
 - ↳ $MASA\ SPAM = (PR - PT) / PR$
 - ↳ ELIMINO LAS QUE SUPEREN UN LÍMITE