

Introducción a algoritmos probabilistas: El problema de "Selección de Secrétaire"

Una empresa quiere contratar a un asistente.
Pide a una empresa de selección de personal que
le envíe N candidatos, uno cada día. $t=1, \dots, N$
En el día t la empresa:

- (1) Paga el valor e de entrevistar al candidato y
- (2) Si el candidato es mejor que todos
los anteriormente vistos paga el valor K Suponga $K \gg e$
de contratar a este nuevo candidato
(y de despedir al anterior)

Asumimos que la empresa tiene un "test
de aptitud" que le permite ranquear totalmente
a los candidatos, sin repeticiones

1 2 3 ... N

$A(1)$ $A(2)$... $A(N)$

Los "scores" $A(i)$
no importa sino sólo
el orden relativo.

Sea $\text{Rank} : \{1, \dots, N\} \rightarrow \{1, \dots, N\}$

$\text{Rank}(i)$ = "ranking del candidato i según $A(\cdot)$ "

Por ejemplo $\text{Rank}(j)=1 \Leftrightarrow j$ es el mejor de todos los
candidatos

Problema: ¿Cuánto le cuesta a la compañía
llevar a cabo esta estrategia?

Sol: El costo depende mucho del ranking
en el que vengan. Para N candidatos

candidatos 1 2 3 4 5

R_{k_1} 1 2 3 4 5 $\rightarrow K + Ne$

R_{k_2} 5 4 3 2 1 $\rightarrow KN + Ne$

En el peor de los casos $KN + eN$

La empresa no sabe en qué Ranqueo van a venir los N candidatos.

Candidato	1	2	3	4	5
Rank	2	3	1	5	4

Permutaciones de n elementos.

$\sigma \in S_n$

Podemos representar nuestra ignorancia diciendo que " σ es una permutación **ALEATORIA** **UNIFORMEMENTE DISTRIBUIDA** en $\Omega = S_N$ "

Más precisamente:

$\Omega = S_n =$ "Permutaciones de N elementos"

$$\mathbb{P}(\{\sigma\}) = \frac{1}{|S_n|} = \frac{1}{N!} \leftarrow \text{medida de probabilidad uniforme en } S_N$$

Si los candidatos tienen rankeo σ el

$$\text{costo es } C(\sigma) = Ne + K \left| \left\{ i: \forall t < i \left[\sigma(t) > \sigma(i) \right] \right\} \right|$$

$$C: S_N \longrightarrow \mathbb{R}$$

es una variable aleatoria.

Cuando i cumple esto contamos a i

$$\sigma \quad 3 \quad 5 \quad 2 \quad 4 \quad 1 \quad 6 \rightarrow 6e + K \cdot 2$$

¿Cuál es el costo esperado?

** Si una variable aleatoria Z puede escribirse como suma de variables aleatorias X_1, \dots, X_N

$$Z = \sum_{i=1}^N X_i$$

entonces

$$\mathbb{E}[Z] \equiv \sum_{i=1}^N \mathbb{E}[X_i]$$

SIEMPRE,
Sin importar qué relaciones
hay entre las X_i 's

En nuestro caso definimos la variable aleatoria $Y^{(i)}: S_N \rightarrow \mathbb{R}$

$$Y^{(i)}(\sigma) = \begin{cases} 1 & \text{si deberíamos contratar al candidato } i \text{ según el ranqueo } \sigma. \\ 0 & \text{si no} \end{cases}$$

Como $Y^{(i)}$ solo asume dos valores

$$\mathbb{E}[Y^{(i)}] = \mathbb{P}\{\sigma \in S_n: \sigma(i) < \sigma(t) \forall t \leq i-1\}$$

Lema: $\mathbb{P}\{\sigma \in S_n : \sigma(i) < \sigma(t) \ \forall t \leq i-1\} = \frac{1}{i}$

$\mathbb{P}\{\sigma \in S_n : \sigma(i) \text{ es el número más pequeño entre } \sigma(1), \sigma(2), \dots, \sigma(i)\}$

Defina $B_j = \{\sigma \in S_n :$

Dem: El mínimo valor de $(\sigma(1), \dots, \sigma(i))$ se alcanza en la componente $j\}$

(1) B_1, B_2, \dots, B_i son disjuntos y cubren todo

(2) $S_n = B_1 \cup \dots \cup B_i$.

Precomponiendo σ con la transposición $(1j)$ construimos una biyección entre B_1 y B_j luego $\mathbb{P}(B_1) = \dots = \mathbb{P}(B_i)$ ③

Hay la misma cantidad de σ con mínimo en 1 que con mínimo en j .

①② y ③ implican que $\mathbb{P}(B_i) = \frac{1}{i}$

Ahora definamos $X = \sum_{i=1}^N Y^{(i)}$ y notemos que

$X(\sigma) = \#$ de candidatos que combatamos si vienen con rango σ

y que

$$C(\sigma) = Ne + K X(\sigma)$$

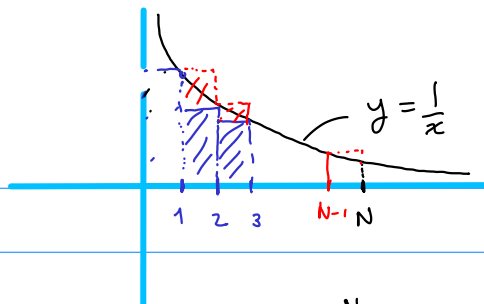
Teorema: El costo esperado de la estrategia es

$$\mathbb{E}[C] = \Theta(K \log(N)) + Ne$$

Dem: Del razonamiento anterior, $X = \sum_{i=1}^N Y^{(i)}$

$$\begin{aligned} \mathbb{E}[C] &= \mathbb{E}\left[Ne + K \sum_{i=1}^N Y^{(i)}\right] = \\ &= Ne + K \sum_{i=1}^N \mathbb{E}[Y^{(i)}] = Ne + K \left(\sum_{i=1}^N \frac{1}{i}\right) = \end{aligned}$$

Cómo se comporta?



Según dibujo:

$$|\text{Área Azul}| \leq \int_1^N \frac{1}{x} dx \leq \text{Área Rojo} = 1 + \frac{1}{2} + \frac{1}{3} + \dots + \frac{1}{N-1}$$

$$\sum_{i=1}^N \frac{1}{i} \leq 1 + \int_1^N \frac{1}{x} dx = 1 + \log(N) - \log(1)$$

$$\int_1^N \frac{1}{x} dx \leq \sum_{i=1}^N \frac{1}{i} \Rightarrow \sum_{i=1}^N \frac{1}{i} \in \Theta(\log(N))$$

$$\text{Luego } \mathbb{E}[C] = Ne + \Theta(K \log(N)).$$

② Aleatoriedad como herramienta:

En el problema anterior hicimos un análisis usando la **probabilidad uniforme** en S_N

como una representación de nuestra ignorancia.

Ahora vamos a utilizarla para mejorar

nuestros costos. Concretamente supongamos que

La empresa de selección nos da la lista de candidatos

① $\begin{bmatrix} A \\ B \\ C \\ \vdots \\ Z \end{bmatrix}$

② \rightarrow
nosotros
aleatorizamos
la lista

$\mathcal{J} \in S_N$
uniforme

$\begin{bmatrix} \mathcal{J}(A) \\ \mathcal{J}(B) \\ \vdots \\ \mathcal{J}(Z) \end{bmatrix}$

nosotros los entrevistamos
en ESTE orden.

Como \mathcal{J} es uniforme,
el Ranking resultante
también lo es
así que **NUESTRA**
Suposición de ranking
uniforme ahora

tendremos $\mathbb{E}[C] = ne + K \log(n)$ es verdadera y

Para poder hacer esto, tenemos que poder "generar" una permutación de S_N con la medida uniforme.

En la clase anterior vimos como generar un elemento de Ω finito cualquiera con una \mathbb{P} dada. El problema es que S_N es ENORME, con $n!$ elementos así que nuestro método general no sirve y requerimos algún mecanismo...

[PERMUTACIÓN UNIFORME MEDIANTE SORTING]

Algoritmo de re-ordenamiento quasi aleatorio:

NOTA: Todo algoritmo randomizado asume que podemos generar números aleatorios uniformes independientes

Random $(i, j) \leftarrow$ entero uniforme entre $\{i, i+1, \dots, j-1, j\}$

Es un tema interesante y profundo cómo construir tales generadores. Por ahora asumámoslo.

Algoritmo: INPUT: Vector A con N componentes
OUTPUT: Vector $C = \sigma(A)$ con N componentes
para $\sigma \in S_n$ aleatoria
con dist. uniforme

$P = [] \leftarrow$ Vector de "Prioridades"

for i in range(len(A)):

$P[i] = \text{Range}(1, N^3)$

$P_{\text{sorted}} = \text{numpy.argsort}(P)$

ORDENA A
de acuerdo a los
tamaños de P

$C = [A[P_{\text{sorted}}[i]] \text{ for } i \text{ in range}(N)]$

tiene un time-complexity $O(N \log(N))$
[Fast-FREE primitive]

* Lema: Si todas las prioridades son distintas entonces la permutación σ resultante tiene distribución uniforme.

$P(\sigma \mid \vec{P} \text{ tiene todas las prioridades distintas})$

Dem: La distribución de probabilidad de las prioridades es invariante respecto al orden de las componentes. Como resultado de esto

↑ Esto es suficiente.

la probabilidad de que el mínimo ocurra en cualquier localización i es la misma.

Usando esto definiremos:

$$X^{(i)} = \{ \text{ } i\text{-ésima prioridad se asigne en la } i\text{-ésima componente} \}$$

$$\begin{aligned} \mathbb{P}\{\text{generar la identidad}\} &= \mathbb{P}\{X^{(1)} \cap X^{(2)} \cap \dots \cap X^{(N)}\} \\ &= \mathbb{P}\{X^{(1)}\} \mathbb{P}\{X^{(2)} | X^{(1)}\} \cdot \mathbb{P}\{X^{(3)} | X^{(1)} \cap X^{(2)}\} \cdot \dots \cdot \mathbb{P}\{X^{(N)} | X^{(1)} \cap \dots \cap X^{(N-1)}\} \\ &= \frac{1}{N} \cdot \frac{1}{(N-1)} \cdot \frac{1}{(N-2)} \cdot \dots \cdot \frac{1}{2} \cdot 1 = \frac{1}{N!} \end{aligned}$$

Por invarianza ese es el valor de las demás también. \square

Aclaración: El algoritmo muestra las prioridades \vec{p} y luego chequea si hay alguna igualdad.

Si
hay alguna igualdad
Descarte todo y
vuelva a empezar

NO
Complete el algoritmo
y produzca σ

$\mathbb{P}\{\sigma \in S_n \mid \vec{p} \text{ NO tiene igualdad}\}$
es uniforme. Este es el
contenido de Lema(*)
preciso

Qué tan probable es que el algoritmo tenga que volver a empezar? Ver ejercicio

$$\mathbb{P}\{(p_1, \dots, p_n) \text{ no tiene repeticiones}\}.$$

③ "Selección de Secretaries .ONLINE "

El término ONLINE no quiere decir en internet sino quiere decir "en tiempo real", es decir es la siguiente variante del problema: (Ver página siguiente)

Vemos N candidatos en orden secuencial,
 Una empresa quiere contratar UN SOLO candidato
 tomando SOLO UNA decisión de contratación
 con el siguiente requerimiento extra: "Justamente
 después de la entrevista TENEMOS que DECIDIR
 inmediatamente si contratamos al candidato o
 no (e informarle)" (la decisión es instantánea = ON-LINE)
 y los datos aparecen en orden secuencial)

La empresa usará su "test de aptitud" A para comparar candidatos
 $i < j \Rightarrow A(i) > A(j)$.

Una estrategia posible es la siguiente:

- (1) Fijamos un número $K \leq N$
- (2) Entrevistamos y rechazamos los primeros k
 candidatos recordando su máximo score \bar{S}
- (3) Para $i \geq k+1$ si $\text{Score}(i) > \bar{S} \Rightarrow$ Contratamos a i ✓
 Si $\forall i \geq k+1$ $\text{Score}(i) < \bar{S} \Rightarrow$ Contratamos al último,
 el candidato N .

Esta estrategia depende de K ¿Cómo escogen
 K de manera óptima? Hay varios criterios
 de optimalidad. Nos enfocaremos en escoger
 al mejor candidato:

Teorema: Si S es el evento de escoger al mejor
 candidato entonces

$$\frac{k}{n} [\ln(n) - \ln(k)] \leq \mathbb{P}(S) \leq \frac{k}{n} (\ln(n-1) - \ln(k-1))$$

Si $k \approx \frac{n}{e}$ tendremos la/el mejor con

$$\mathbb{P}(S) \geq \frac{1}{e} \approx 0.367879$$

$$0.367 \cdot n$$

↑ más de $\frac{1}{3}$ de
 las veces!!

Dem: $S = \{\sigma : \text{escogemos el mejor}\}$

$S_i = \{\sigma : \text{escogemos el mejor y este está en posición } i\}$

$S = \bigcup_i S_i$ y $S_i = \emptyset$ para $i \leq k$ pues este nunca se escoge.

$$\mathbb{P}(S) = \sum_{i=k+1}^N \mathbb{P}(S_i)$$

$S_i = \text{Best applicant en posición } i$ y $\underbrace{y = B_i \cap O_i}_{\substack{\text{No escogemos } j \text{ en } k+1, \dots, i-1 \\ (\text{porque } A(i,j) < \bar{s}, j \in \{k+1, \dots, i-1\})}}$

O_i depende de los scores $1, \dots, i-1$ o mejor, de su ordenamiento relativo = Máximo está en posiciones $1, 2, \dots, k$

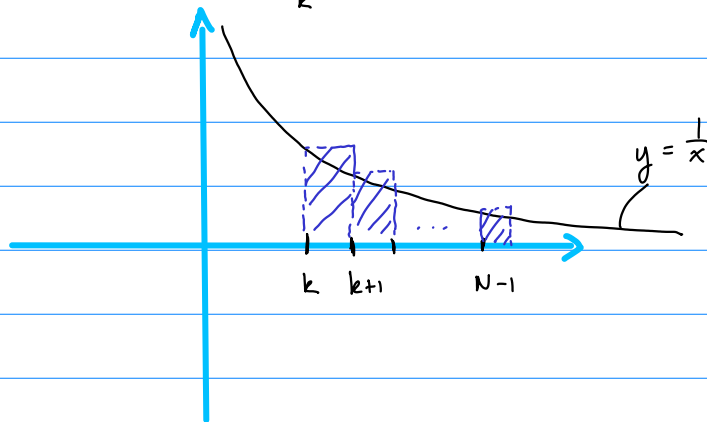
B_i depende de si el $A(i)$ es más grande que todos los demás.

Así que O_i y B_i son independientes
luego

$$\mathbb{P}(S_i) = \mathbb{P}(B_i) \mathbb{P}(O_i) = \frac{1}{N} \frac{k}{(i-1)}$$

$$\mathbb{P}(S) = \sum_{i=k+1}^N \frac{k}{N(i-1)} = \frac{k}{N} \sum_{i=k+1}^N \frac{1}{i-1}$$

$$= \frac{k}{N} \sum_{i=k}^{N-1} \frac{1}{i} \geq \frac{k}{N} \int_k^N \frac{1}{x} dx = \frac{k}{N} [\log(N) - \log(k)]$$



$$\frac{k}{N} (\log(N) - \log(k)) = \varphi(k)$$

$$\varphi'(k) = \frac{1}{N} (\log(N) - \log(k)) + \frac{k}{N} \left[-\frac{1}{k} \right] = 0$$

$$\frac{\log(N)}{N} - \frac{\log(k)}{N} - \frac{1}{N} = 0$$

$$\log(k) = \log(N) - 1$$

$k^* = \frac{N}{e}$ y reemplazando $\mathbb{P}(S) \geq \frac{1}{e}$