

Lo q' queda pendiente es calcular estimadores  $P(x_i/S)$  y  $P(x_i/S)$ .  
 Si tuviéramos una cantidad razonable de <sup>mails de</sup> entrenamiento, ~~una~~  
 es natural estimar, por ejemplo  $P(x_i/S)$  como <sup>el conteo</sup> ~~la frecuencia~~  
<sup>mails spam</sup> q' contienen la palabra  $w_i$  ~~adecuado~~ y  
 el total de mails <sup>q' tienen spam</sup> ~~q' contienen~~  $w_i$  =  $\frac{\# \text{ mails spam con } w_i}{\# \text{ mails totales con } w_i}$

$$P(x_i/S) \sim \frac{\# \text{ mails spam con } w_i}{\# \text{ mails totales con } w_i}$$

Problema: esta forma de estimar  $P(x_i/S)$  tiene un problema.  
 Supongamos q' la palabra "data" ocurre solo en mails q'  
 no son spam en nuestro conjunto de entrenamiento. Es  
 decir:  ~~$P(\text{"data"}/S) = 0$~~   $P(\text{"data"}/S) = 0$

Entonces el clasificador asignará probabilidad 0 a cualquier  
 mail q' contenga la palabra "data", aunque el texto sea

"data on cheap viagra and authentic rolex watches"

Solución: la estimación se hace utilizando <sup>una constante</sup> ~~el pseudocount~~  
 $K$  ( $K$  dice q' es chico,  $K_j = 1$ )  
 $K$ :

$$P(x_i/S) \sim \frac{(K + \# \text{ mails spam con } w_i)}{(2K + \# \text{ mails con spam})}$$

↓  
 se asigna un peso  
 menor a las 'observaciones' ficticias.