
1. Definición del problema

Definimos en primer lugar el conjunto A como el conjunto de caracteres definidos en la lengua española en minúscula, sin tener en cuenta signos de puntuación ni acentos. Es decir,

$$A = \{a, b, c, d, e, f, g, h, i, j, k, l, m, n, \tilde{n}, o, p, q, r, s, t, u, v, w, x, y, z\}$$

Consideramos además el conjunto B dado por los mismos caracteres pero en mayúscula. Es decir,

$$B = \{A, B, C, D, E, F, G, H, I, J, K, L, M, N, \tilde{N}, O, P, Q, R, S, T, U, W, X, Y, Z\}$$

Podemos definir la aplicación $toUpper : A \rightarrow B$ que asigna a cada caracter de A su correspondiente caracter en mayúsculas de B . Es fácil ver que esta aplicación es biyectiva y su inversa la denotaremos por $toLower : B \rightarrow A$. Podremos denotar también $B = toUpper(A)$.

Definimos también el conjunto $V = \{a, e, i, o, u\} \subset A$ de las vocales, el conjunto T de acentos¹ y consideramos el conjunto W el cual incluye todas las variaciones de dichas vocales mediante los acentos de T . Consideramos entonces la aplicación $tilde : T \times V \rightarrow W$ el cual añade el acento de T a la vocal de V resultando en un elemento de W .

Definimos el conjunto P como el conjunto de caracteres que no se encuentra en $A \cap B \cap V \cap W$, el cual incluye signos de puntuación y el resto de caracteres ASCII. Al conjunto de caracteres ASCII lo denotaremos por Ω y de esta forma tenemos que $\{A, B, V, W, P\}$ define una partición de Ω , es decir

1. $A \cup B \cup V \cup W \cup P = \Omega$
2. $X \cap Y = \emptyset \quad \forall X, Y \in \{A, B, V, W, P\}, \quad X \neq Y$

Definiremos una **palabra** p como un vector de elementos de Ω , es decir

$$p = (x_1, \dots, x_n) \text{ con } x_i \in \Omega \quad \forall i \in 1, \dots, n$$

Observación. Este concepto es más extenso que el de palabra que se entiende en el lenguaje ya que a priori puede contener cualquier signo de puntuación, incluyendo caracteres en blanco (espacios).

Diremos que una palabra es **propia** si no contiene espacios en blanco, es decir, si

$$p = (x_1, \dots, x_n) \text{ con } x_i \in \Omega \setminus \{ ' '\} \quad \forall i \in 1, \dots, n$$

Diremos que n es la **longitud** de la palabra. Además por la definición de p tenemos que p tiene longitud n si y solo si $p \in \Omega^n$.

De esta forma, el problema que se plantea es encontrar una aplicación

$$d : \Omega^n \times \Omega^m \rightarrow \mathbb{R}^+ \quad m, n \in \mathbb{N}$$

de forma que d sea una distancia en el espacio $\Omega^n \times \Omega^m$.

¹Cabe destacar que T no es un conjunto de caracteres sino de tildes (variaciones)

Para seguir trabajando con la notación adecuada definiremos unos cuantos conceptos que serán útiles.

Definiremos la aplicación $C : \Omega^n \rightarrow \Omega$ como los **caracteres** de una palabra y estará definida como

$$C(p) = \{x_1, \dots, x_n\}, \quad \text{para } p = (x_1, \dots, x_n)$$

Notemos que con esta aplicación se pierde la propiedad de orden que tenía p como vector.

Definimos así una **subpalabra** sp de una palabra $p \in \Omega^n$ como una proyección con respecto a las coordenadas i -ésima a la $(i+k)$ -ésima de p , es decir, sp será una subpalabra de p si y solo si

$$sp = \pi_{i,i+1,\dots,i+k}(p) \text{ con } 1 \leq i, \quad i+k \leq n$$

Notaremos por $SP(p)$ al conjunto de subpalabras de p .

Algunas propiedades inmediatas son

1. $C(sp) \subseteq C(p)$
2. $sp \in \Omega^{k+1}$, es decir, sp tiene longitud $k+1 \leq n$.

Podemos además definir una aplicación aditiva como

$$\begin{aligned} + : \Omega^n \times \Omega^m &\rightarrow \Omega^{m+n} \\ +(p, p') &= p + p' = (x_1, \dots, x_n, x'_1, \dots, x'_n) \end{aligned}$$

Notemos que con esta definición, la suma no es conmutativa pero sí es asociativa. Podemos entenderla como una concatenación de las palabras p y p' .

Definiremos tres aplicaciones bastante importantes para las soluciones propuestas. Definimos la **inserción** como la aplicación $Ins : \Omega^n \times \Omega \times \mathbb{N} \cap [0, n] \rightarrow \Omega^{n+1}$ dada por

$$Ins(p, x, i) = (x_1, \dots, x_{i-1}, x, x_i, \dots, x_n)$$

y diremos que insertamos en la palabra p el caracter x en la posición i -ésima.

Definiremos la **eliminación** como $Del : \Omega^n \times \mathbb{N} \cap [0, n] \rightarrow \Omega^{n-1}$ dada por

$$Del(p, i) = (x_1, \dots, x_{i-1}, x_{i+1}, \dots, x_n)$$

y diremos que eliminamos el caracter i -ésimo de la palabra p .

Definiremos la **sustitución** como $Sus : \Omega^n \times \Omega \times \mathbb{N} \cap [0, n] \rightarrow \Omega^n$ dada por

$$Sus(p, x, i) = (x_1, \dots, x_{i-1}, x, x_{i+1}, \dots, x_n)$$

y diremos que insertamos el sustituimos por caracter x la posición i -ésima de la palabra p .

A estas tres aplicaciones las llamaremos **aplicaciones elementales**.

2. Primera solución

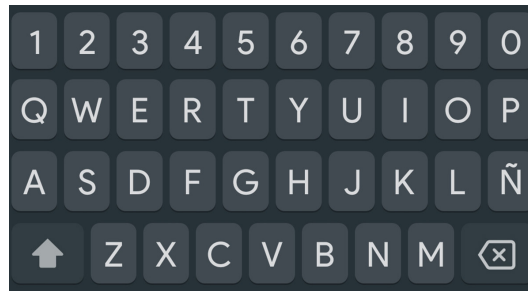
La primera distancia que se define que verifica todos los requisitos del problema podría ser la conocida como **distancia de Levenshtein**. Esta distancia consiste en encontrar el número de inserciones, eliminaciones y sustituciones mínimas necesarias para transformar una palabra p en otra p' , es decir

$$d_L(p, p') = k \iff \text{se puede pasar de } p \text{ a } p' \text{ con } k \text{ aplicaciones elementales}$$

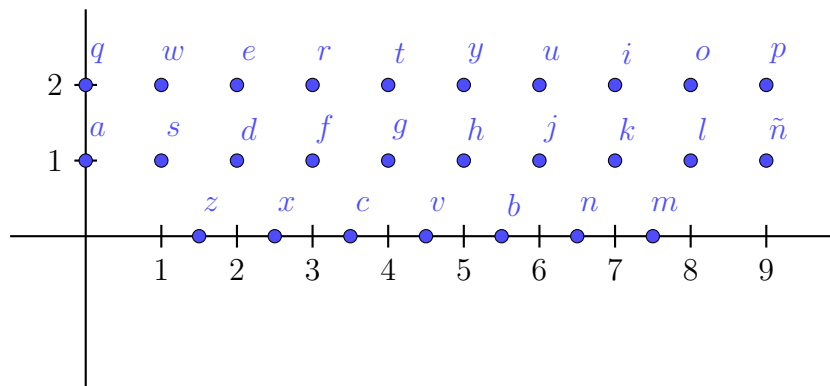
Notemos que esta distancia toma valores en \mathbb{N} y que no es muy precisa ya que tenemos que $d((c, a), (c, a, s, a)) = d((c, a), (x, \tilde{n}))$ lo cual parece poco intuitivo y poco práctico para comparar palabras.

3. Segunda solución

La segunda solución que se plantea está basada en la anterior pero añade un **peso** a cada aplicación elemental. Para ello definiremos una distancia geométrica sobre un teclado real considerando para ello cada caracter de A como la coordenada central de cada tecla y utilizando la distancia euclídea en \mathbb{R}^2 . Sabemos que un teclado estándar QWERTY (de móvil) tiene la siguiente distribución:



Por lo que podríamos considerar la siguiente representación en el plano:



De esta forma tendremos definida una distancia entre caracteres dada por $d_2(c_1, c_2) = \sqrt{(c_{11} - c_{21})^2 + (c_{12} - c_{22})^2}$ donde $c_1 = (c_{11}, c_{12})$, $c_2 = (c_{21}, c_{22})$. Sin embargo vamos a considerar una nueva distancia d dada por

$$d(c_1, c_2) = \frac{d_2(c_1, c_2)}{\sqrt{\max_{c \in A} \{d_2(c_1, c)\} \cdot \max_{c \in A} \{d_2(c_2, c)\}}}$$

y de esta forma tendremos que d sigue siendo una distancia y además $0 \leq d(c_1, c_2) \leq 1$ verificándose que

1. $d(c_1, c_2) = 0 \iff c_1 = c_2$ (por ser d_2 una distancia)
2. $d(c_1, c_2) = 1 \iff d_2(c_1, c_2) = \max_{c \in A} \{d_2(c_1, c)\} = \max_{c \in A} \{d_2(c_2, c)\}$, es decir, si c_1 y c_2 son los caracteres más alejados mutuamente.

De esta forma podemos definir una matriz $M \in \mathcal{M}_{27}(\mathbb{R})$ dada por

$$\{a_{ij}\}_{i,j} = \{d(c_i, c_j)\}_{i,j}$$

y M será simétrica y será la **matriz de distancias**.