

Procesamiento de Lenguaje Natural

Traducción automática y búsqueda



Dra. Helena Gómez Adorno

helena.gomez@iimas.unam.mx

Dra. Gemma Bel

gbele@iingen.unam.mx



Correo del curso:

pln.cienciadedatos@gmail.com

Asistente:

Luis Ramon Casillas

Contenido

Traducción automática

“hola” → “bonjour”

Búsqueda de documentos

“Puedo obtener”
un reembolso?”

“¿Cuál es tu política
de devoluciones?”

¿Puedo recuperar mi
dinero?

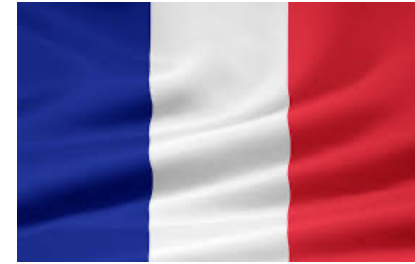
Contenido

- Transformación de vectores
- K-vecinos mas cercanos
- Tablas Hash
- Dividir espacio de vectores en regiones
- Hash sensible a la localidad
- Vecinos mas cercanos aproximados

Revisión de traducciones



?



Español $\left\{ \begin{array}{l} [-2, -4, -8] \\ \dots \\ [2, 1, 0] \end{array} \right.$

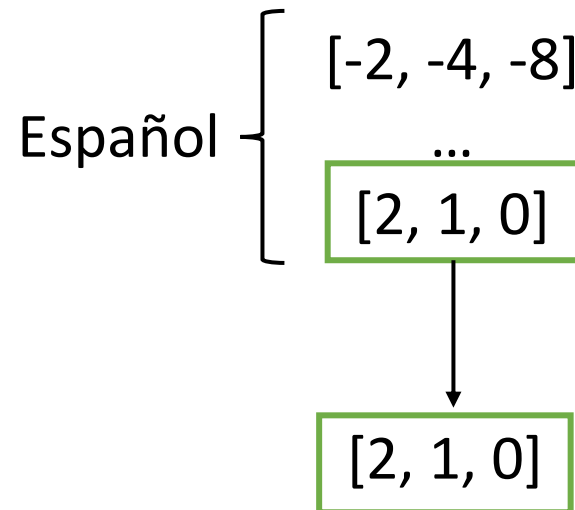
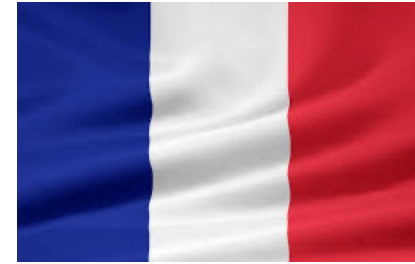
Francés $\left\{ \begin{array}{l} [3, 2, 7] \\ \dots \\ [9, 33, 12] \end{array} \right.$

Revisión de traducciones

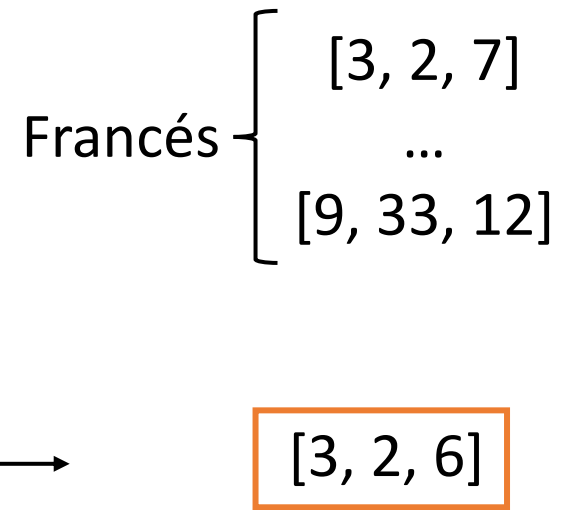


Gato

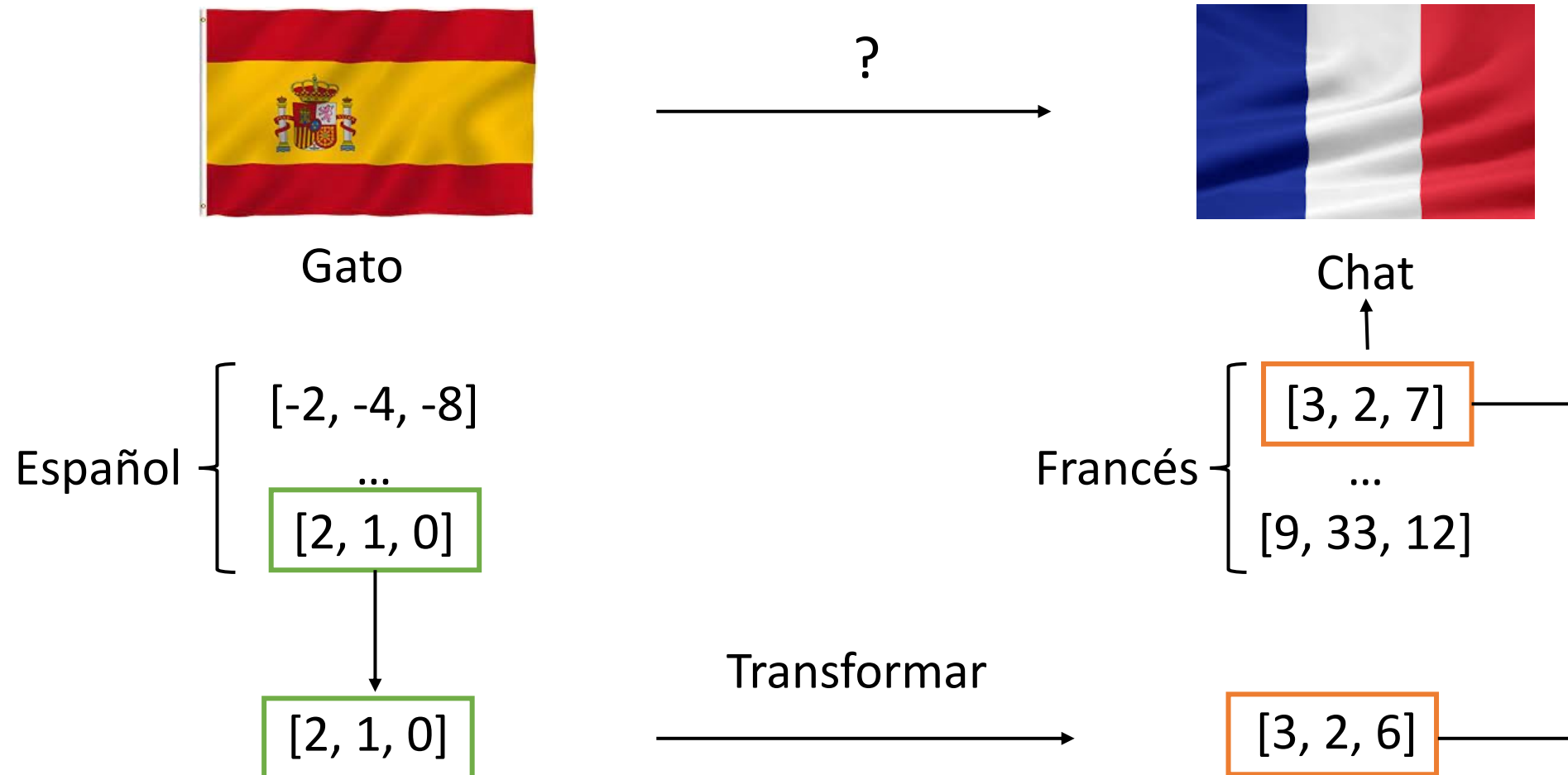
?



Transformar



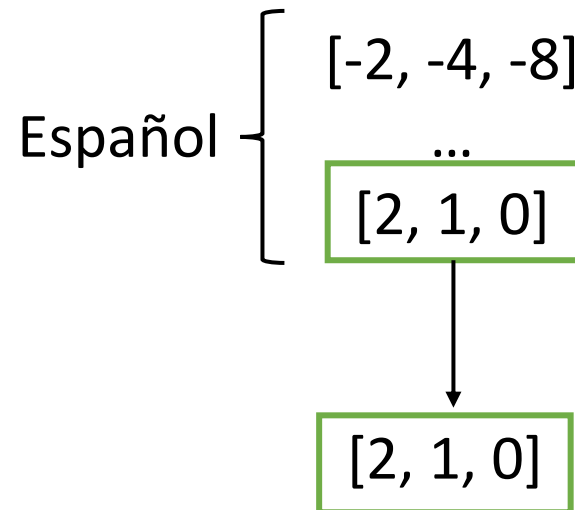
Revisión de traducciones



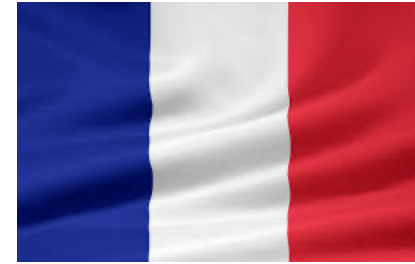
Revisión de traducciones



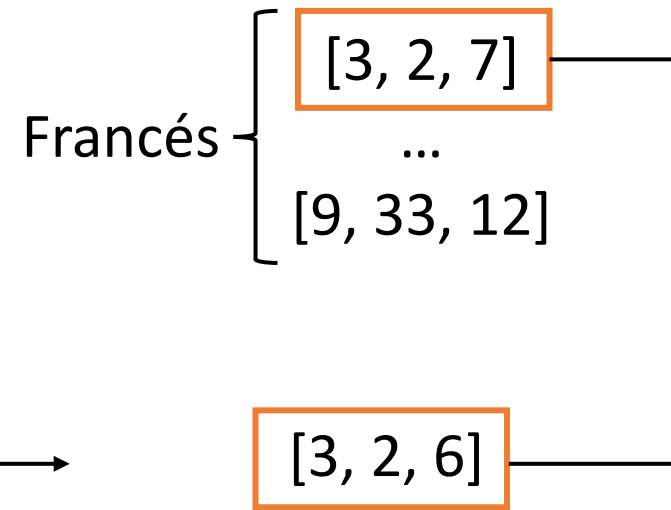
Gato



?



Chat



Transformar



Con una matriz

Transformación de vectores

```
R = np.array([[2, 0],  
              [0, -2]])
```

```
x = np.array([[1, 1]])
```

```
np.dot(x, R)
```

```
array([[2, -2]])
```


Alineación de vectores de palabras

$$XR \approx Y$$

$$\begin{pmatrix} \text{[vector "gato"]} \\ \text{[... vector]} \end{pmatrix}$$

X

$$\begin{pmatrix} \text{[vector "chat"]} \\ \text{[... vector]} \end{pmatrix}$$

Y

Alineación de vectores de palabras

$$XR \approx Y$$

$$\begin{pmatrix} [\text{vector "gato"}] \\ [\dots \text{vector}] \\ [\text{vector "zebra"}] \end{pmatrix}$$

X

$$\begin{pmatrix} [\text{vector "chat"}] \\ [\dots \text{vector}] \\ [\text{vector "zébresse"}] \end{pmatrix}$$

Y

Alineación de vectores de palabras

$$XR \approx Y$$

$$\begin{bmatrix} \text{[vector "gato"]} \\ \text{[... vector]} \\ \text{[vector "zebra"]} \end{bmatrix}$$

X

$$\begin{bmatrix} \text{[vector "chat"]} \\ \text{[... vector]} \\ \text{[vector "zébrasse"]} \end{bmatrix}$$

Y

Subconjunto del vocabulario total

Pasos para aprender R

$$Loss = \| \mathbf{XR} - \mathbf{Y} \|_F$$

Pasos para aprender R

initialize R

$$Loss = \| \mathbf{XR} - \mathbf{Y} \|_F$$

Pasos para aprender R

initialize R

in a loop:

$$Loss = \| \mathbf{XR} - \mathbf{Y} \|_F$$

Pasos para aprender R

initialize R

in a loop:

$$Loss = \| \mathbf{XR} - \mathbf{Y} \|_F$$

$$g = \frac{d}{dR} Loss$$

gradient

Pasos para aprender R

initialize R

in a loop:

$$Loss = \| \mathbf{X}\mathbf{R} - \mathbf{Y} \|_F$$

$$g = \frac{d}{dR} Loss$$

gradient

$$R = R - \alpha g$$

update

Norma Frobenius

$$\| \mathbf{XR} - \mathbf{Y} \|_F$$

$$\mathbf{A} = \begin{pmatrix} 2 & 2 \\ 2 & 2 \end{pmatrix}$$

$$\|\mathbf{A}_F\| = \sqrt{2^2 + 2^2 + 2^2 + 2^2}$$

Norma Frobenius

$$\| \mathbf{XR} - \mathbf{Y} \|_F$$

$$\mathbf{A} = \begin{pmatrix} 2 & 2 \\ 2 & 2 \end{pmatrix}$$

$$\|\mathbf{A}_F\| = \sqrt{2^2 + 2^2 + 2^2 + 2^2}$$

$$\|\mathbf{A}_F\| = 4$$

$$\|\mathbf{A}\|_F \equiv \sqrt{\sum_{i=1}^m \sum_{j=1}^n |a_{ij}|^2}$$

Norma Frobenius al cuadrado

$$\|\mathbf{XR} - \mathbf{Y}\|_F^2$$

$$\mathbf{A} = \begin{pmatrix} 2 & 2 \\ 2 & 2 \end{pmatrix}$$

$$\|\mathbf{A}\|_F^2 = \left(\sqrt{2^2 + 2^2 + 2^2 + 2^2} \right)^2$$

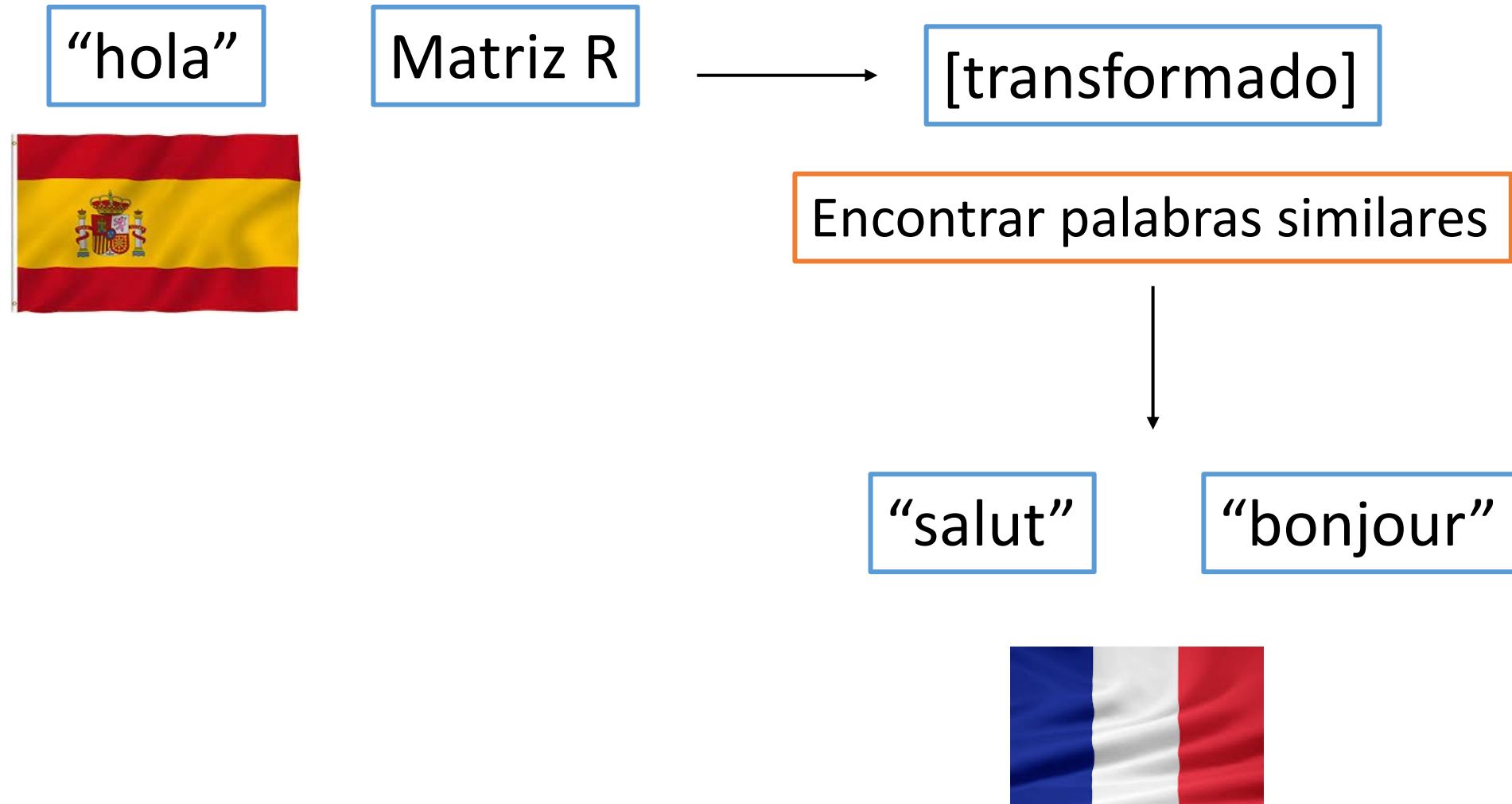
$$\|\mathbf{A}\|_F^2 = 16$$

Gradiente

$$Loss = \|\mathbf{XR} - \mathbf{Y}\|_F^2$$

$$g = \frac{d}{dR} Loss = \frac{2}{m} (\mathbf{X}^T (\mathbf{XR} - \mathbf{Y}))$$

Encontrar la traducción



Vecinos más cercanos



You



San Francisco

Friend



Vecinos más cercanos



You
San Francisco

Friend



Location

Shanghai



Bangalore



Los Angeles

Vecinos más cercanos



You
San Francisco



Friend



Location

Shanghai

Nearest

2



Bangalore

3



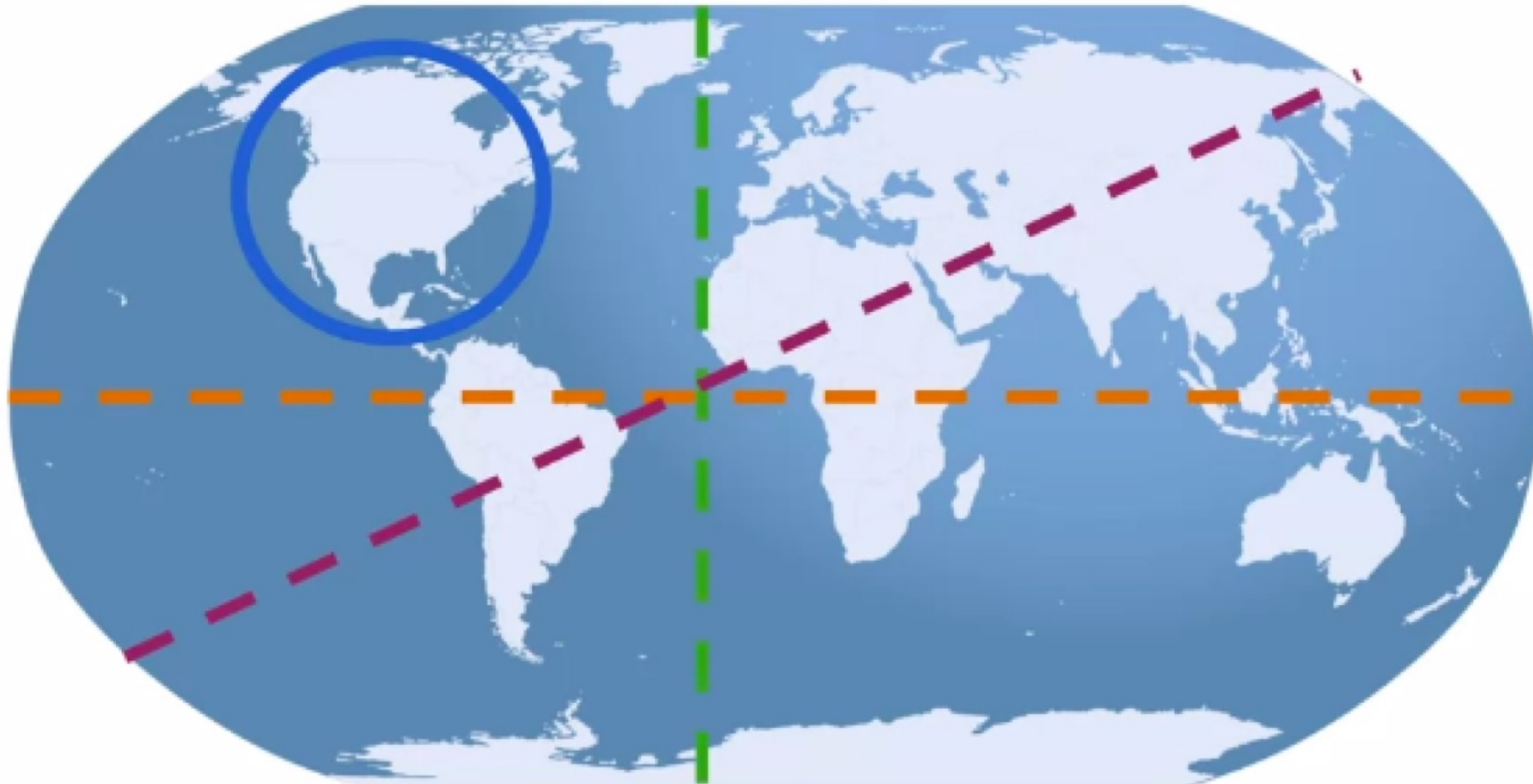
Los Angeles

1

Vecinos más cercanos



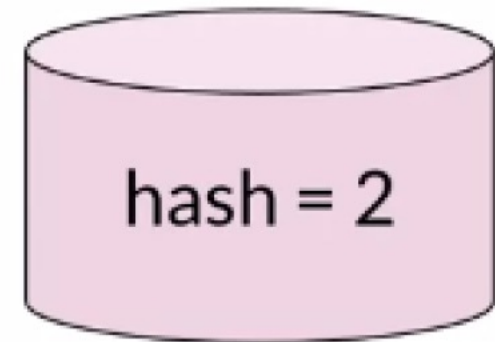
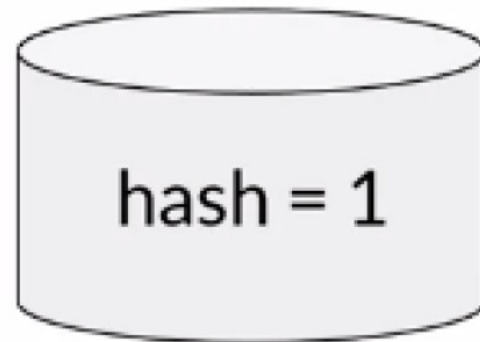
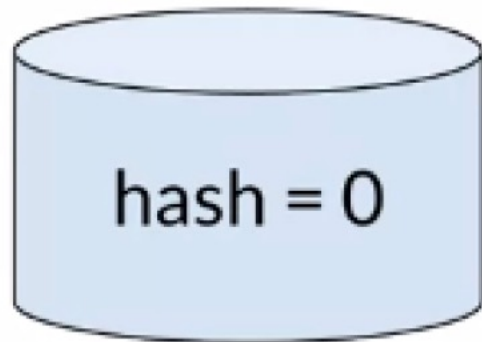
Vecinos más cercanos



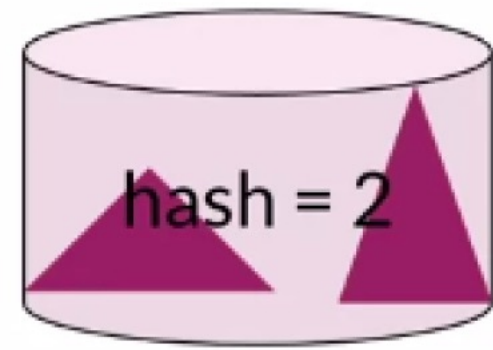
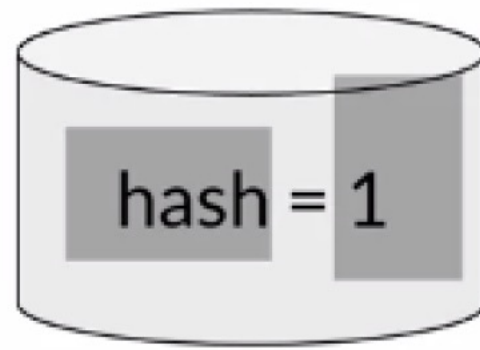
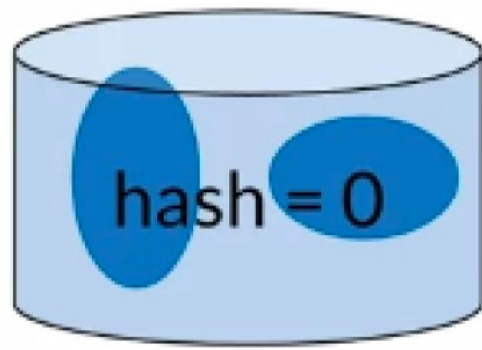
Tablas hash



Tablas hash



Tablas hash



Tablas hash



Hash function (vector) \longrightarrow Hash value

Tablas hash

0	1	2	3	4	5	6	7	8	9
100				14			17		
10							97		

Función hash (vector)= → valor hash

Tablas hash

0	1	2	3	4	5	6	7	8	9
100				14			17		
10							97		

Función hash (vector)= → valor hash

Valor hash= vector % número de cubetas

Creación de tabla hash básica

```
def basic_hash_table(value_l, n_buckets):  
    def hash_function(value_l, n_buckets):  
        return int(value) % n_buckets  
    hash_table = {i:[] for i in range(n_buckets)}  
    for value in value_l:  
        hash_value = hash_function(value, n_buckets)  
        hash_table[hash_value].append(value)  
    return hash_table
```


Función hash

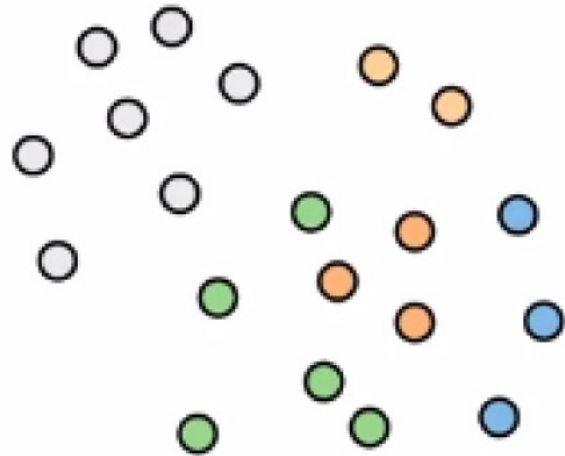
0	1	2	3	4	5	6	7	8	9
100				14			17		
10							97		

Función hash por ubicación?

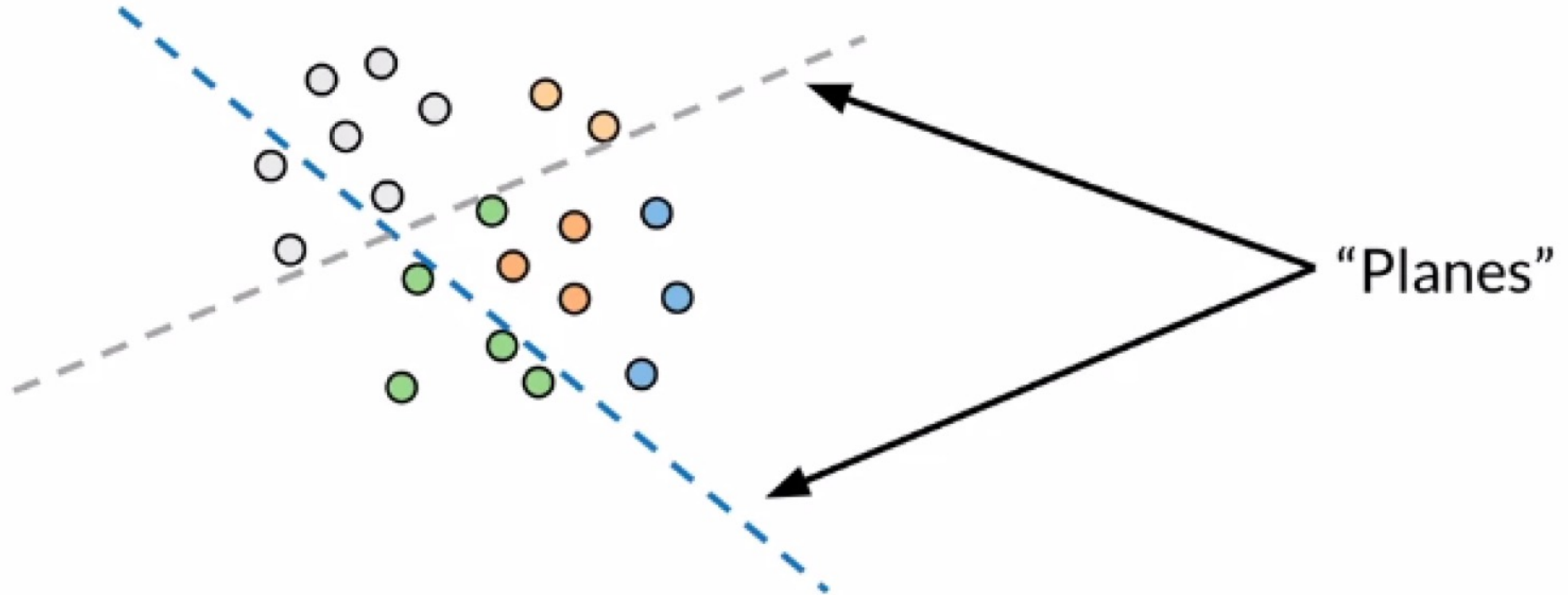
0	1	2	3	4	5	6	7	8	9
14									100
10									97
17									

Idealmente queremos que las palabras (números) cercanas queden en la misma cubeta.
Solución: Hashing sensible a la localidad (locality sensitive)

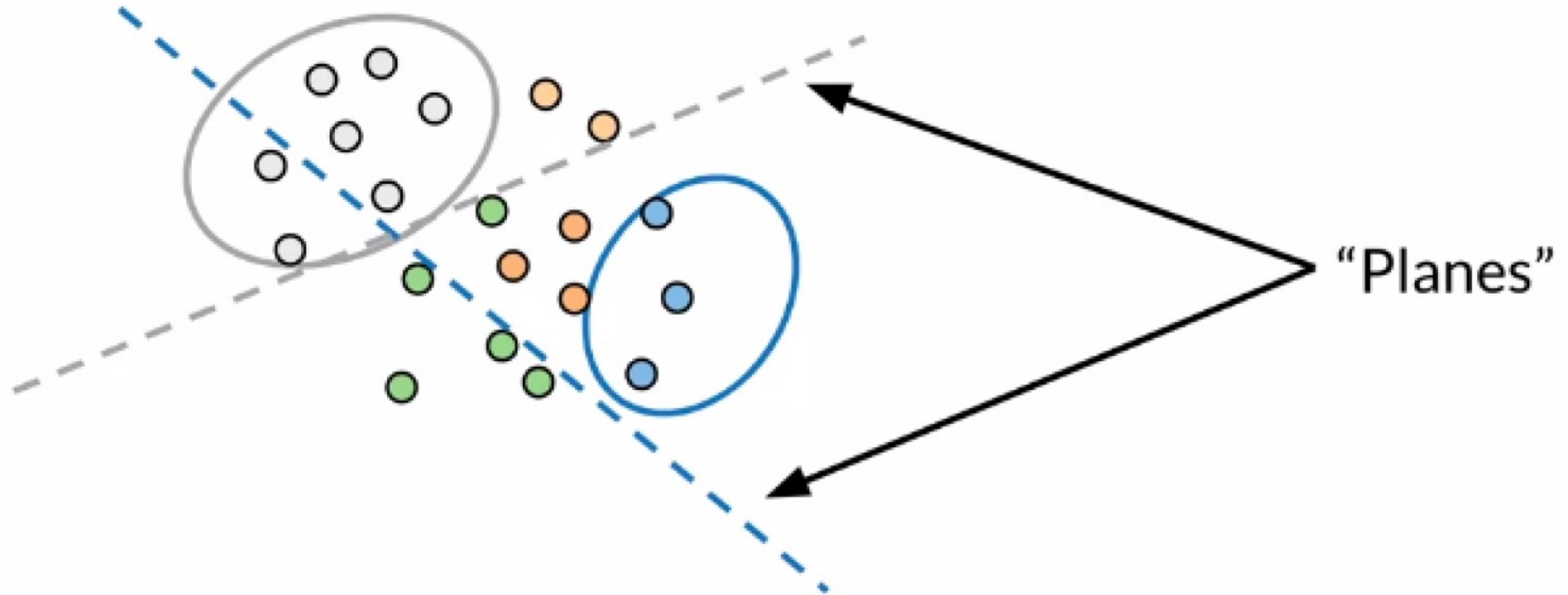
Locality sensitive hashing



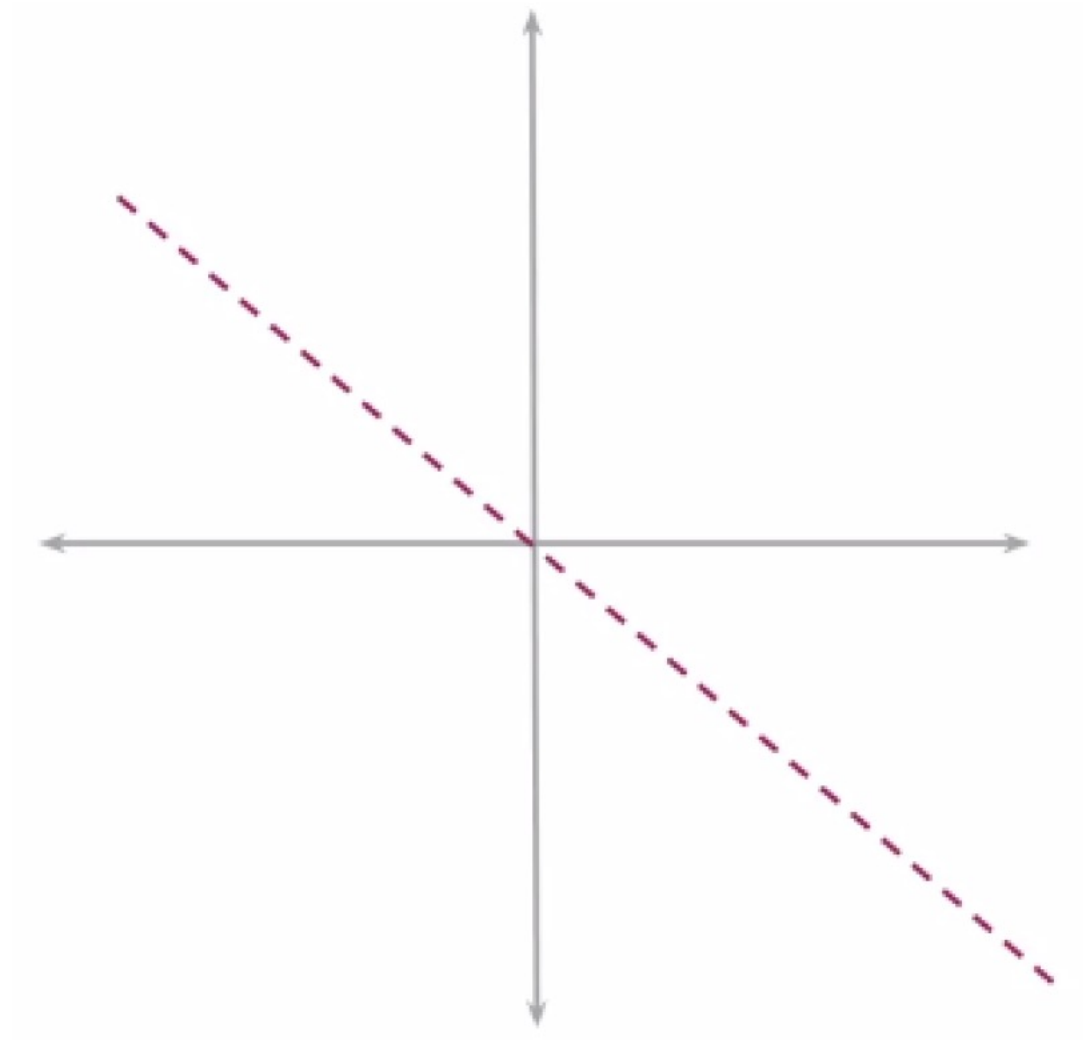
Locality sensitive hashing



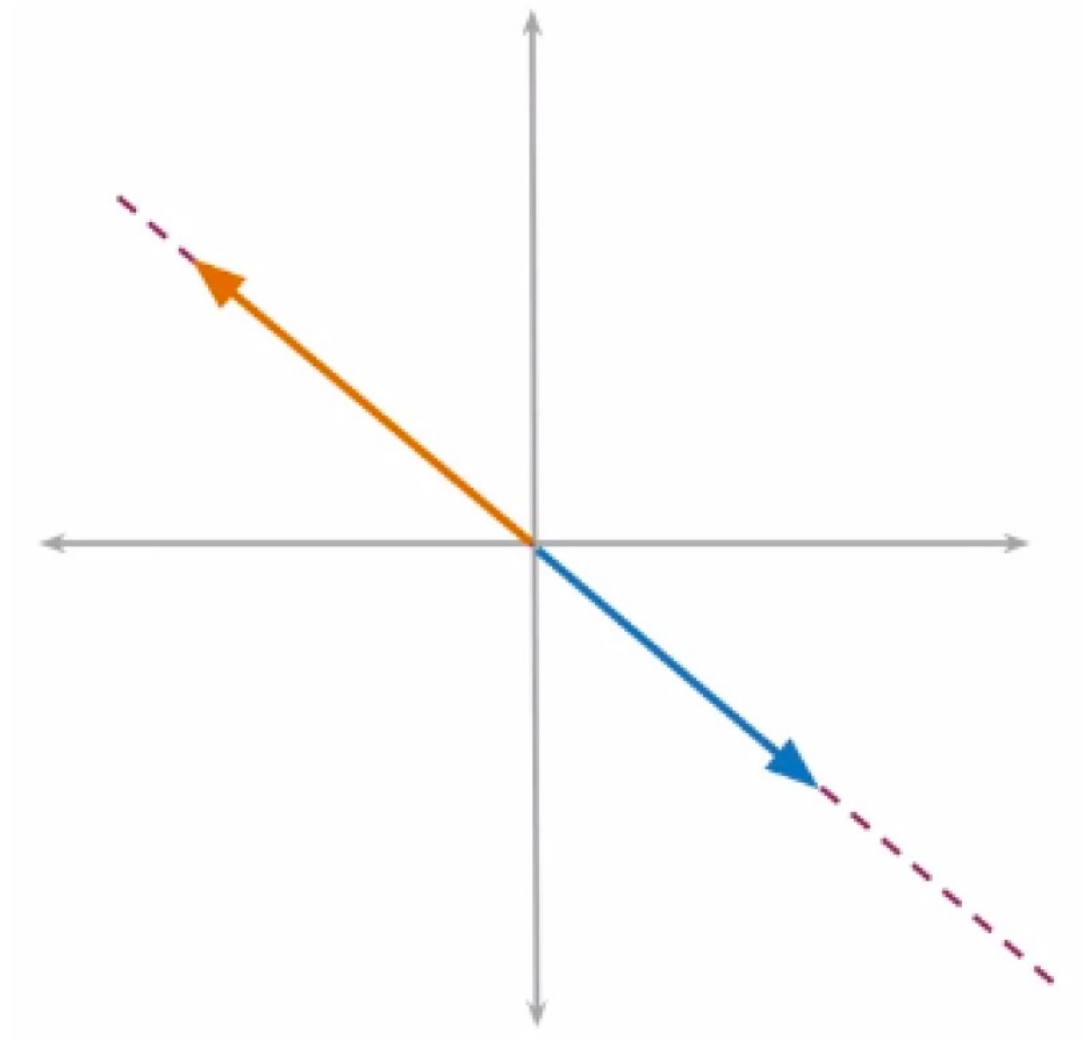
Locality sensitive hashing



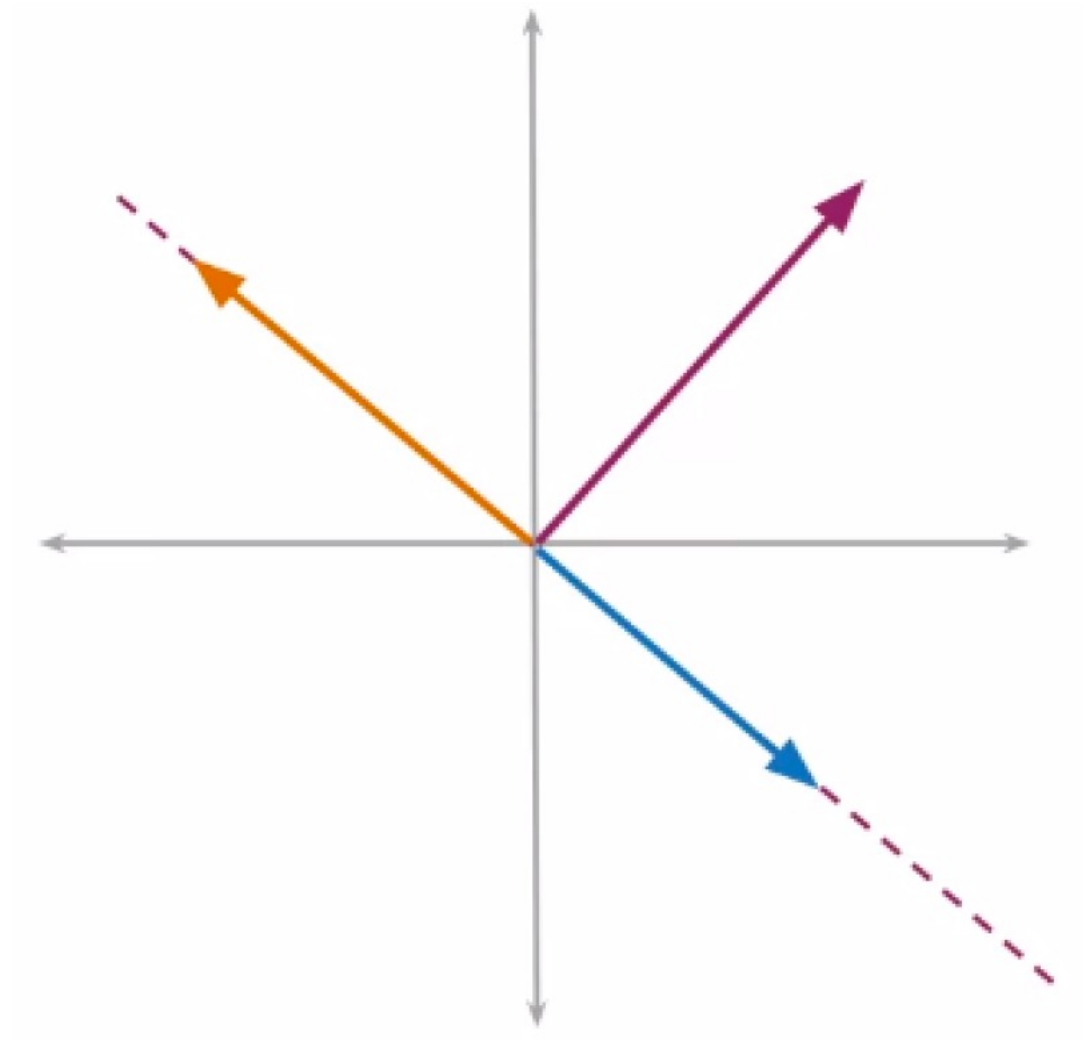
Planos



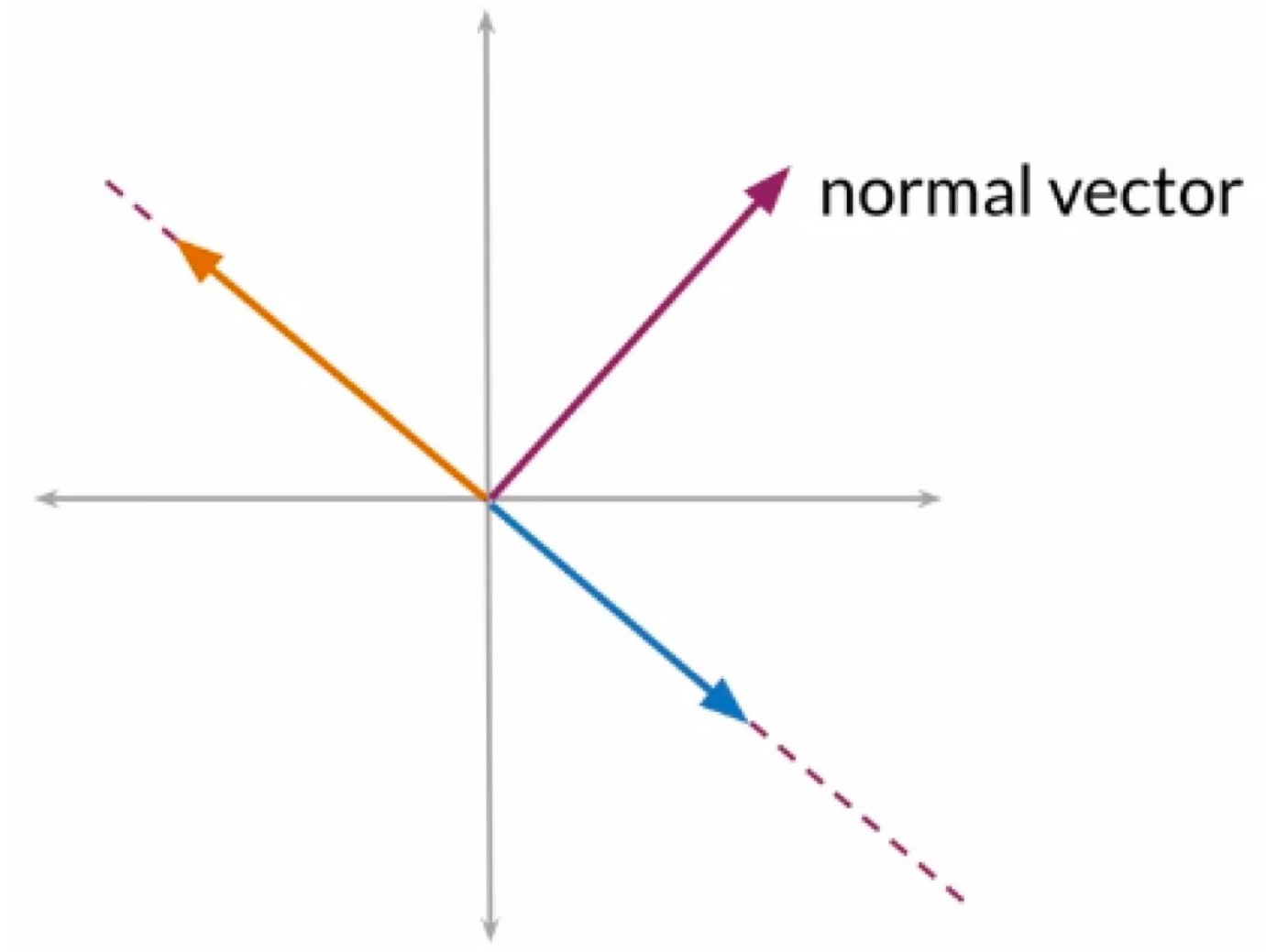
Planos



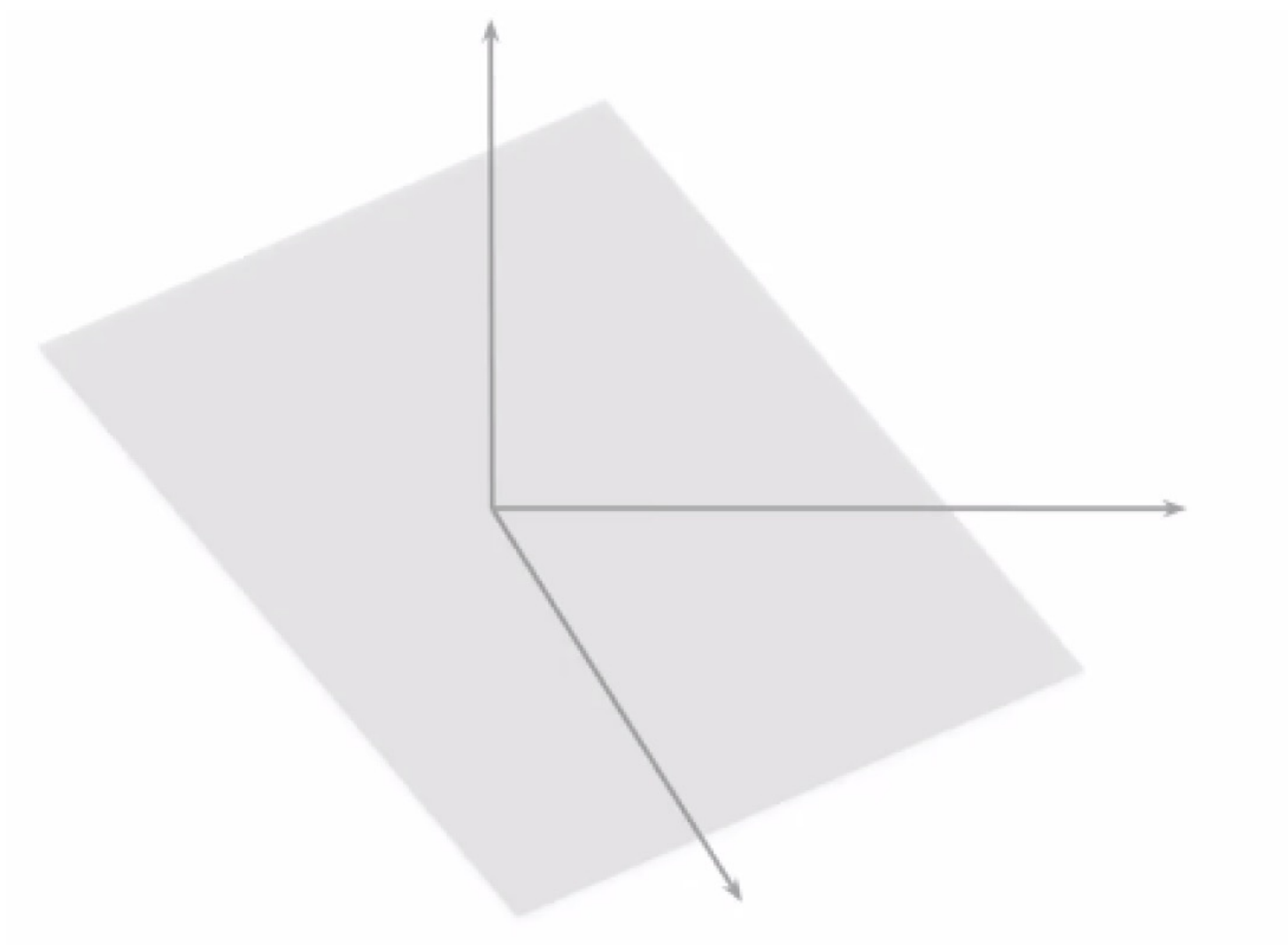
Planos



Planos

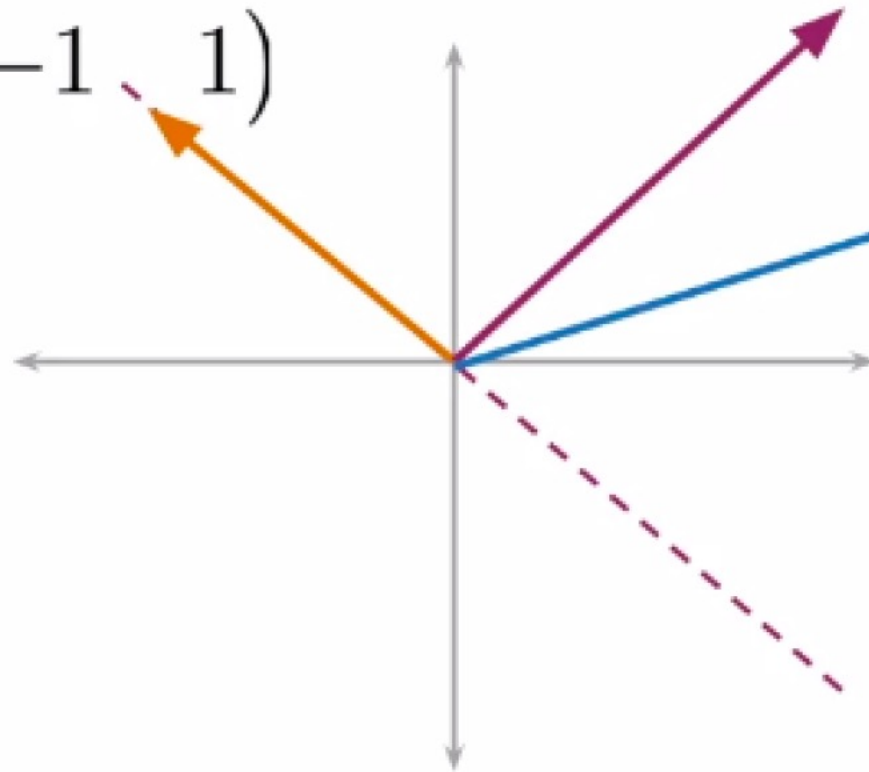


Planos



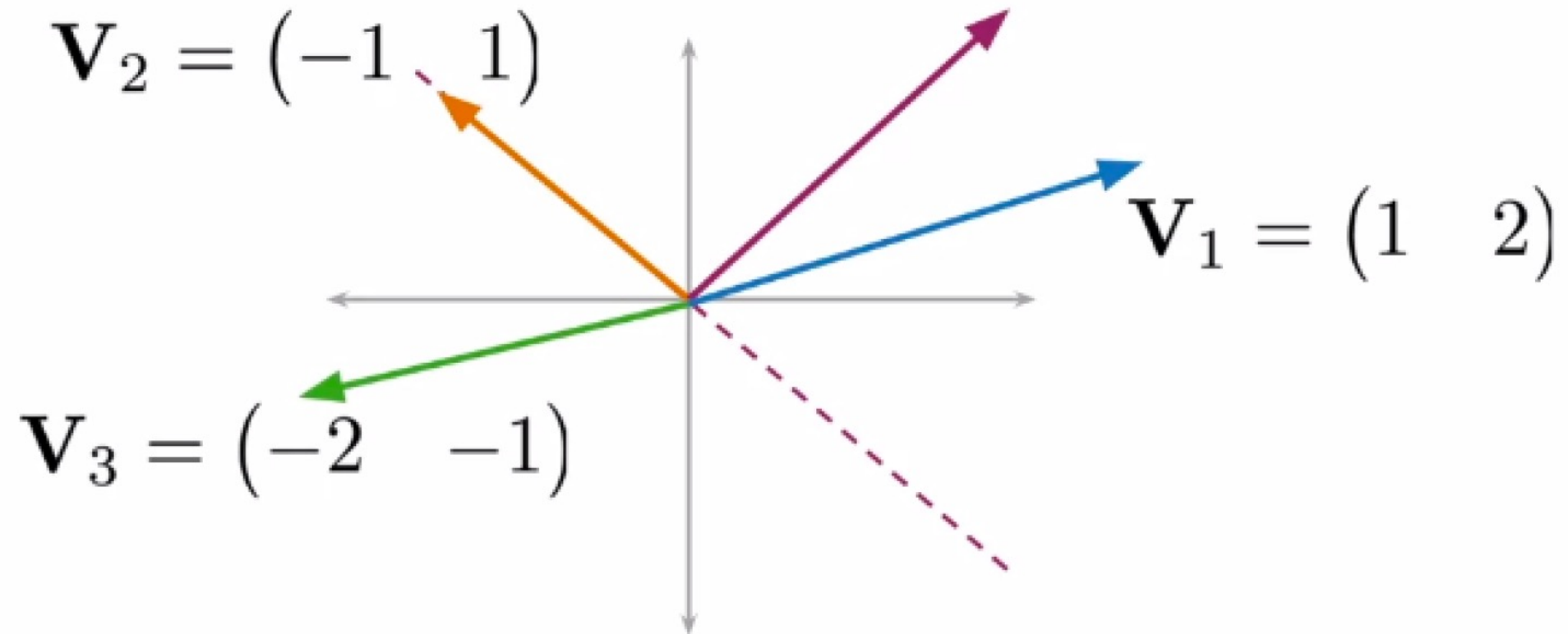
Cuál lado del plano?

$$\mathbf{V}_2 = (-1 \ 1)$$

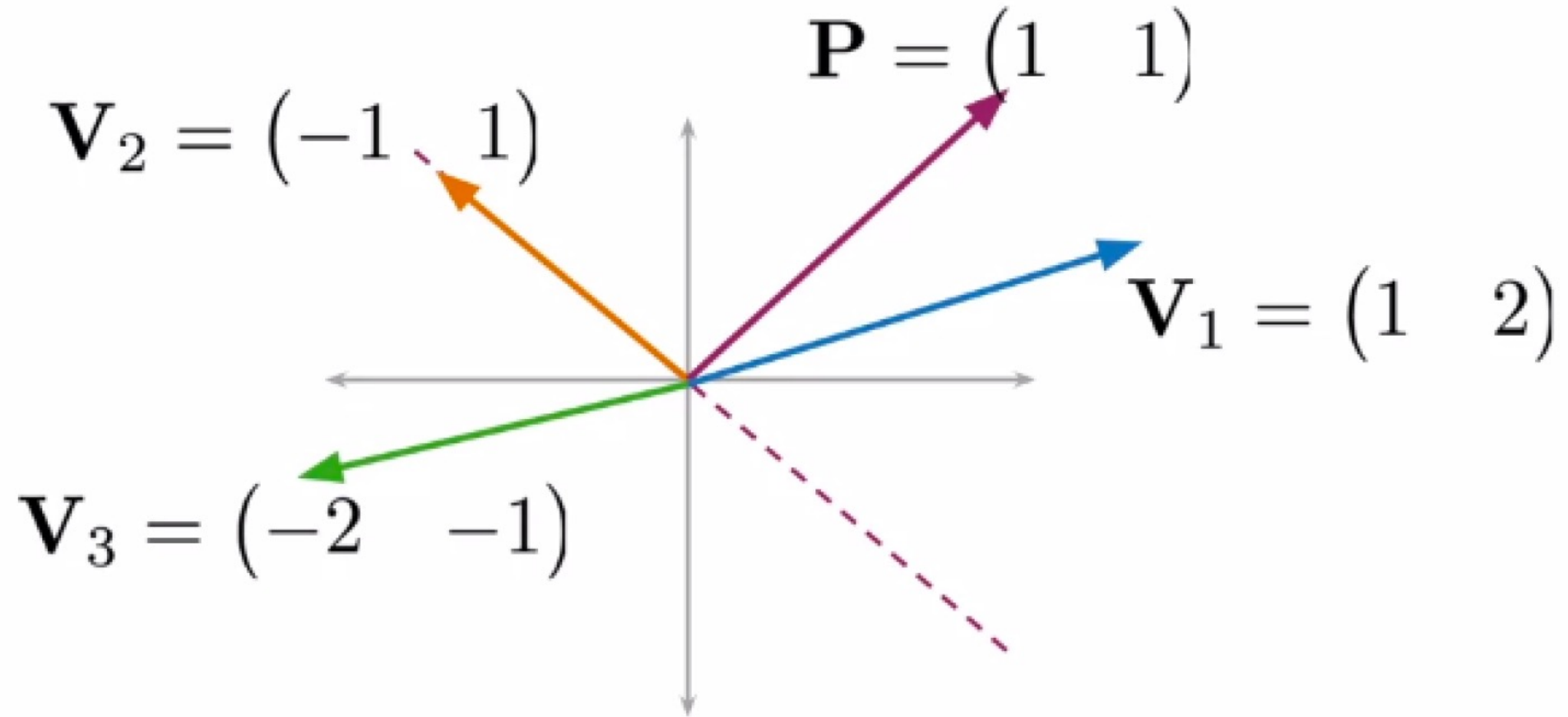


$$\mathbf{V}_1 = (1 \ 2)$$

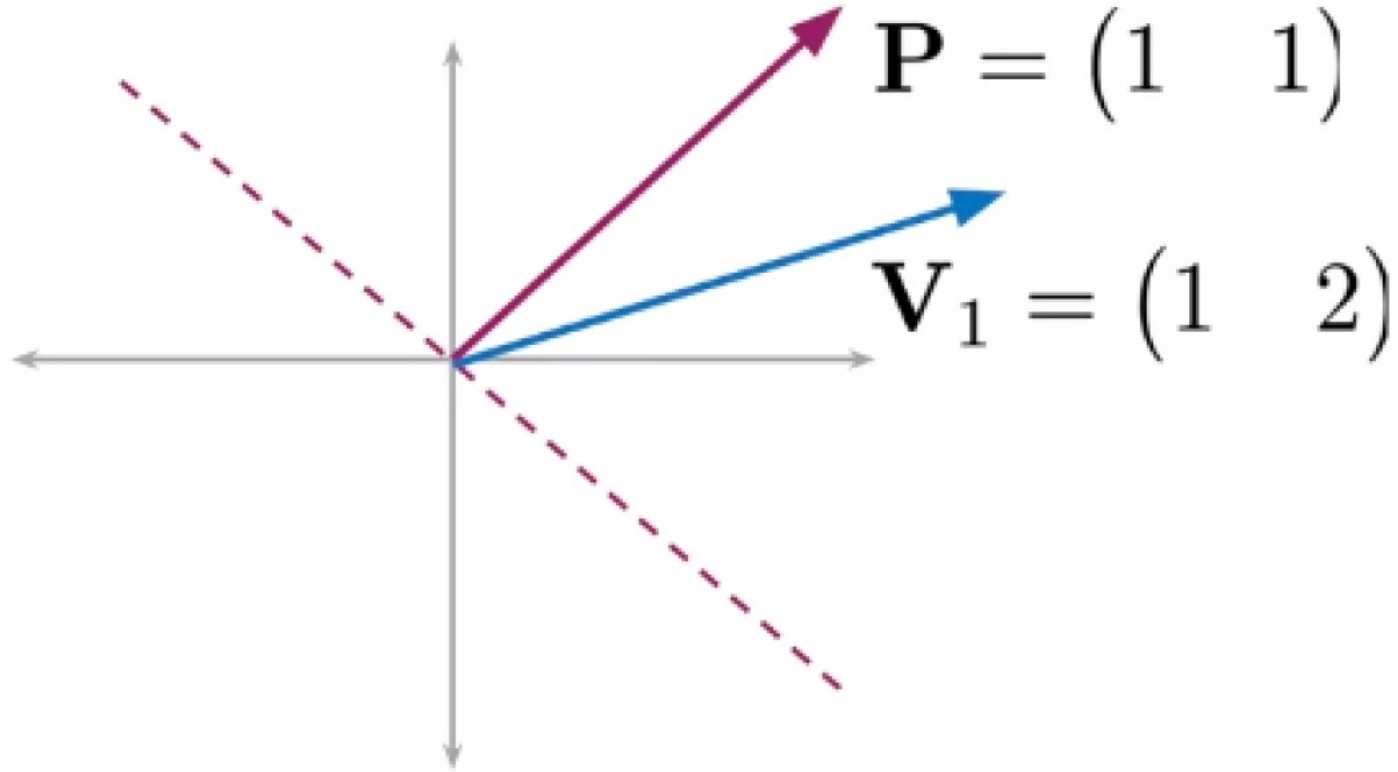
Cuál lado del plano?



Cuál lado del plano?



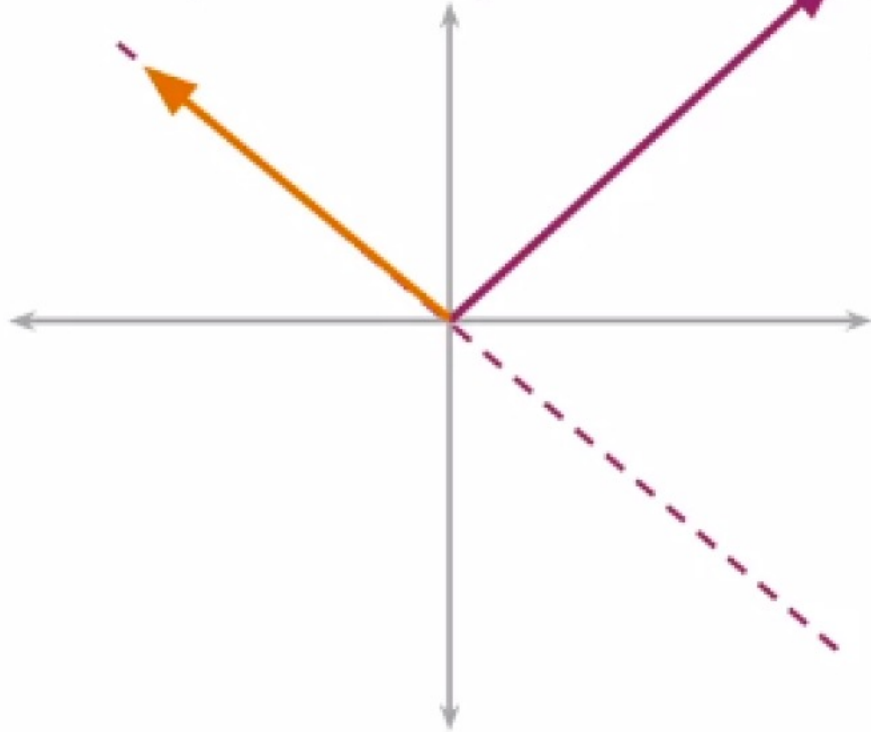
Cuál lado del plano?



$$\mathbf{P}\mathbf{V}_1^T = 3$$

Cuál lado del plano?

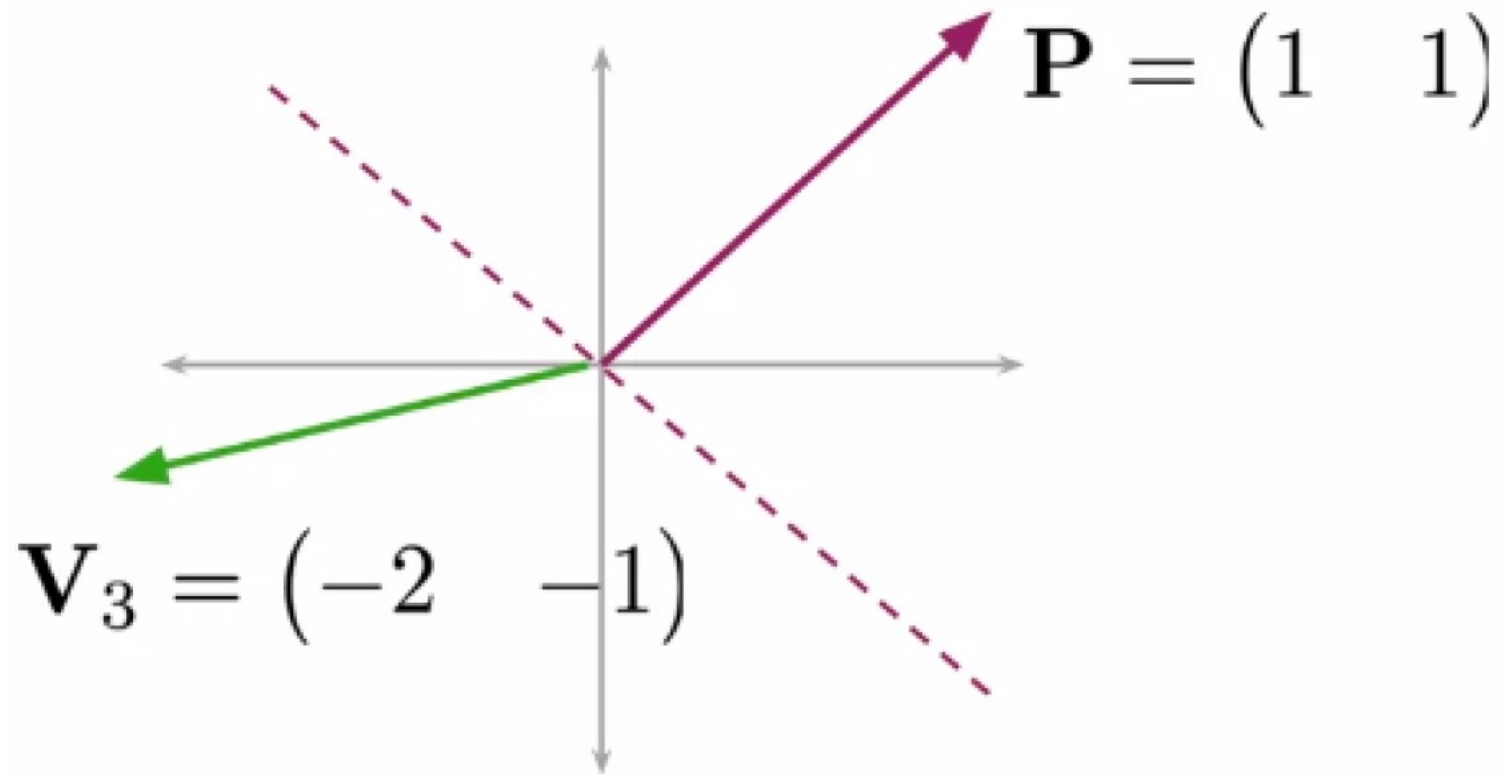
$$\mathbf{V}_2 = (-1 \ 1)$$



$$\mathbf{P} = (1 \ 1)$$

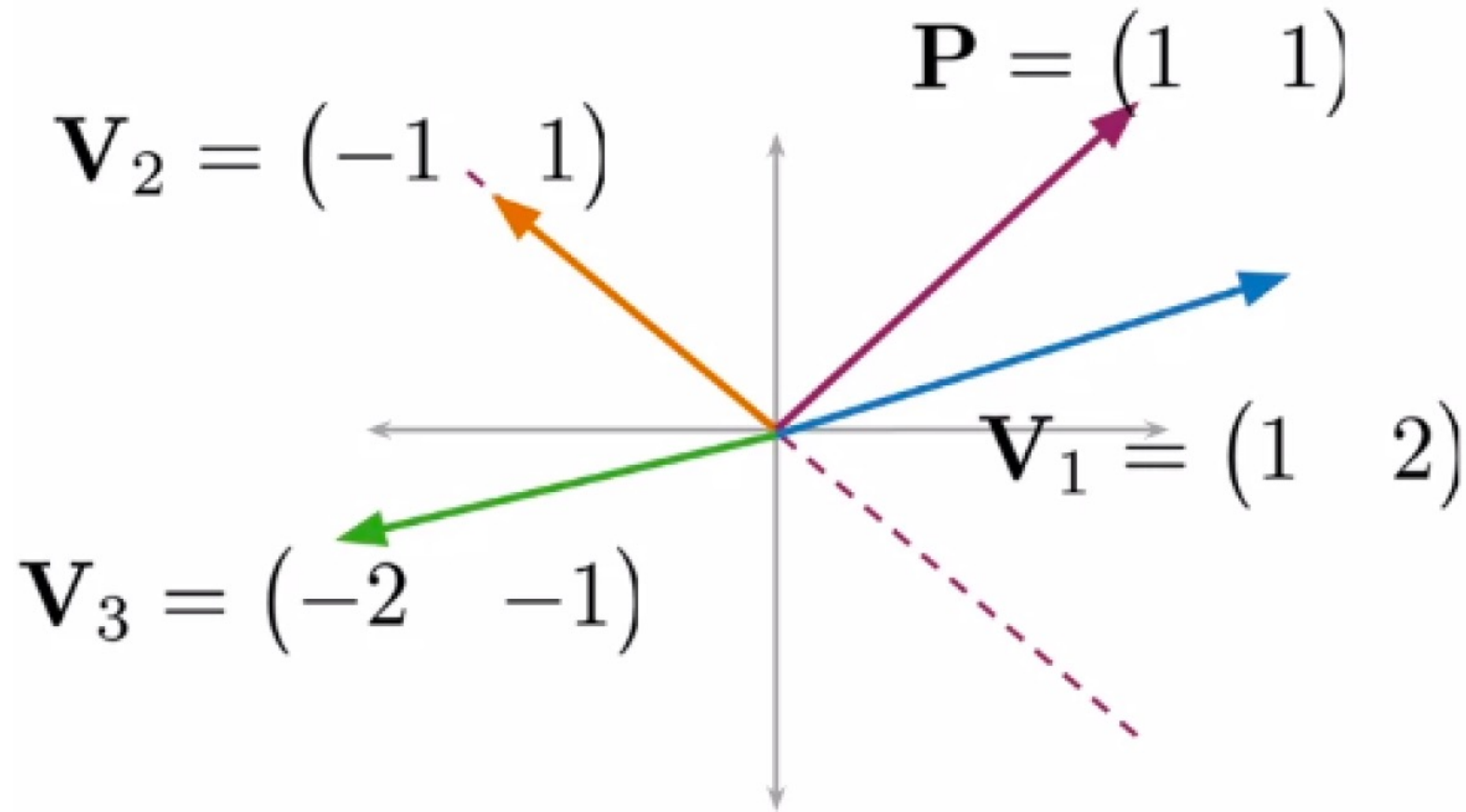
$$\mathbf{P}\mathbf{V}_2^T = 0$$

Cuál lado del plano?



$$\mathbf{P}\mathbf{V}_3^T = -3$$

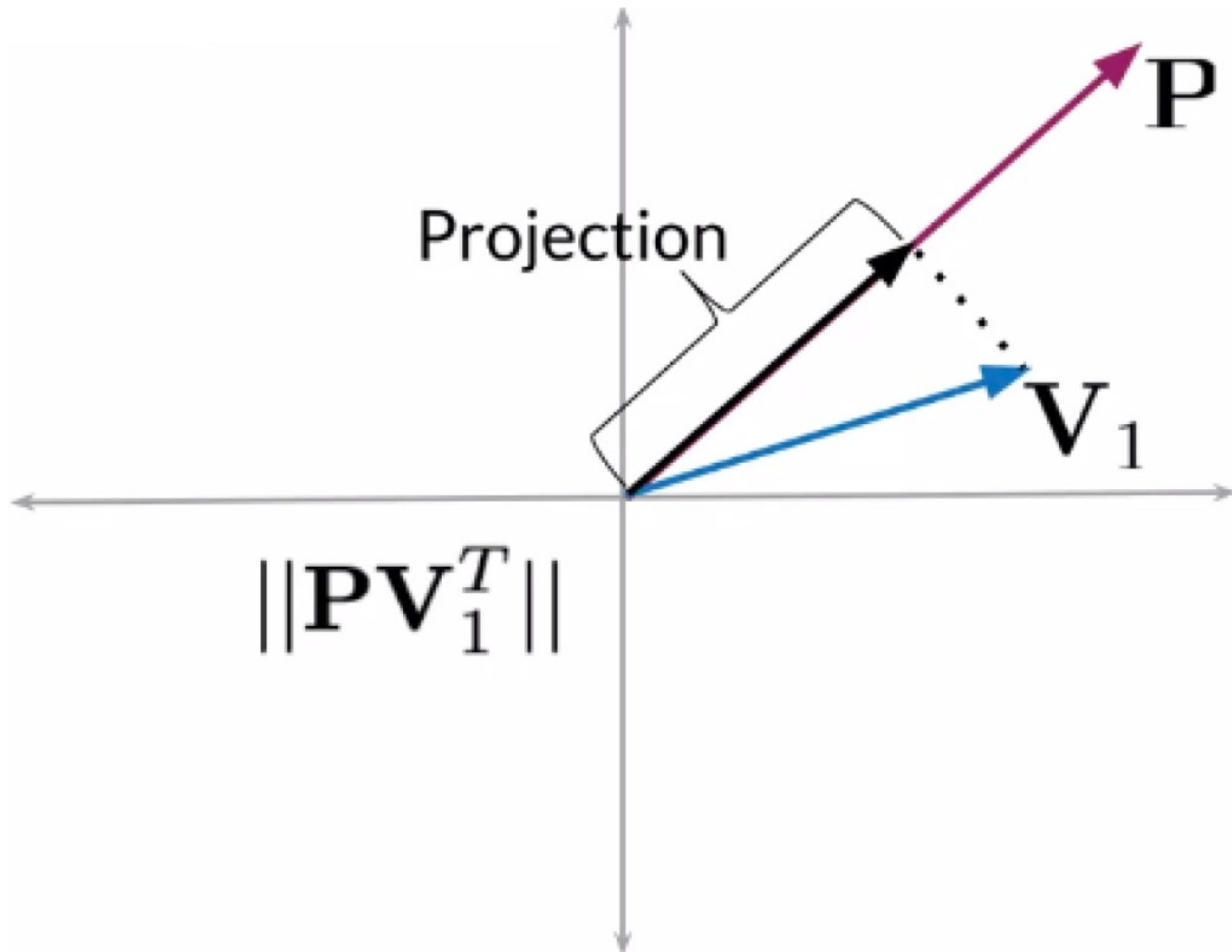
Cuál lado del plano?



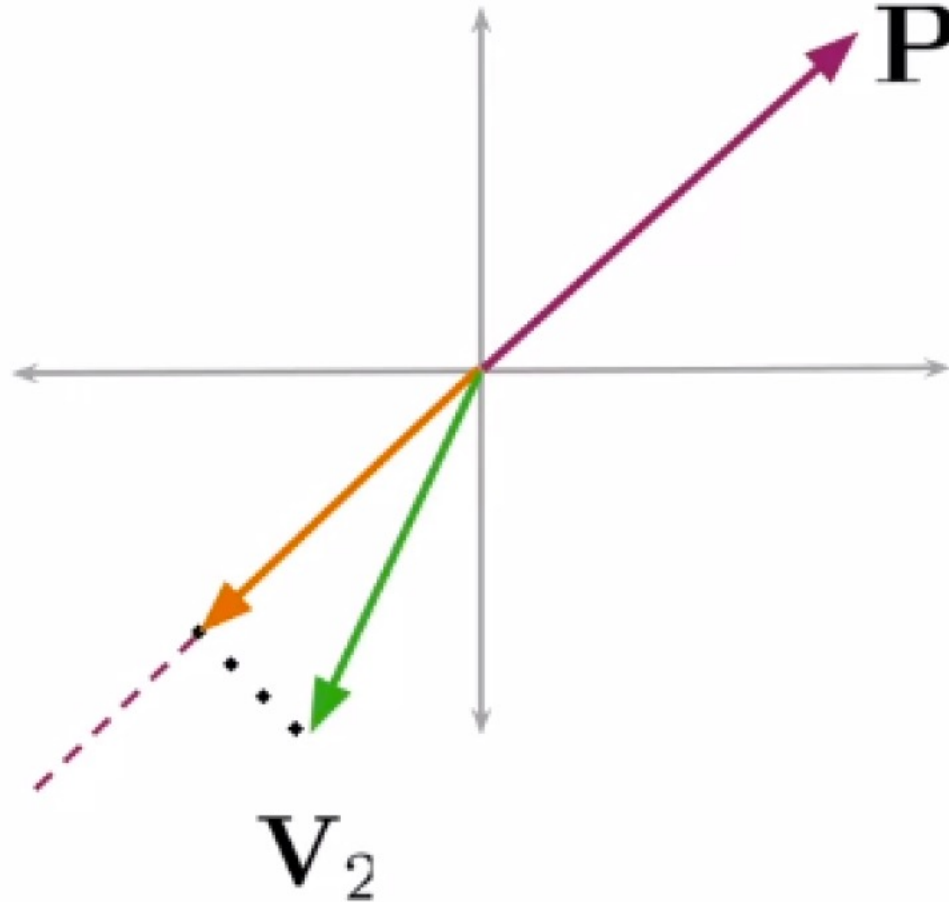
$$\begin{aligned}\mathbf{P}\mathbf{V}_1^T &= 3 \\ \mathbf{P}\mathbf{V}_2^T &= 0 \\ \mathbf{P}\mathbf{V}_3^T &= -3\end{aligned}$$

Nótese los signos

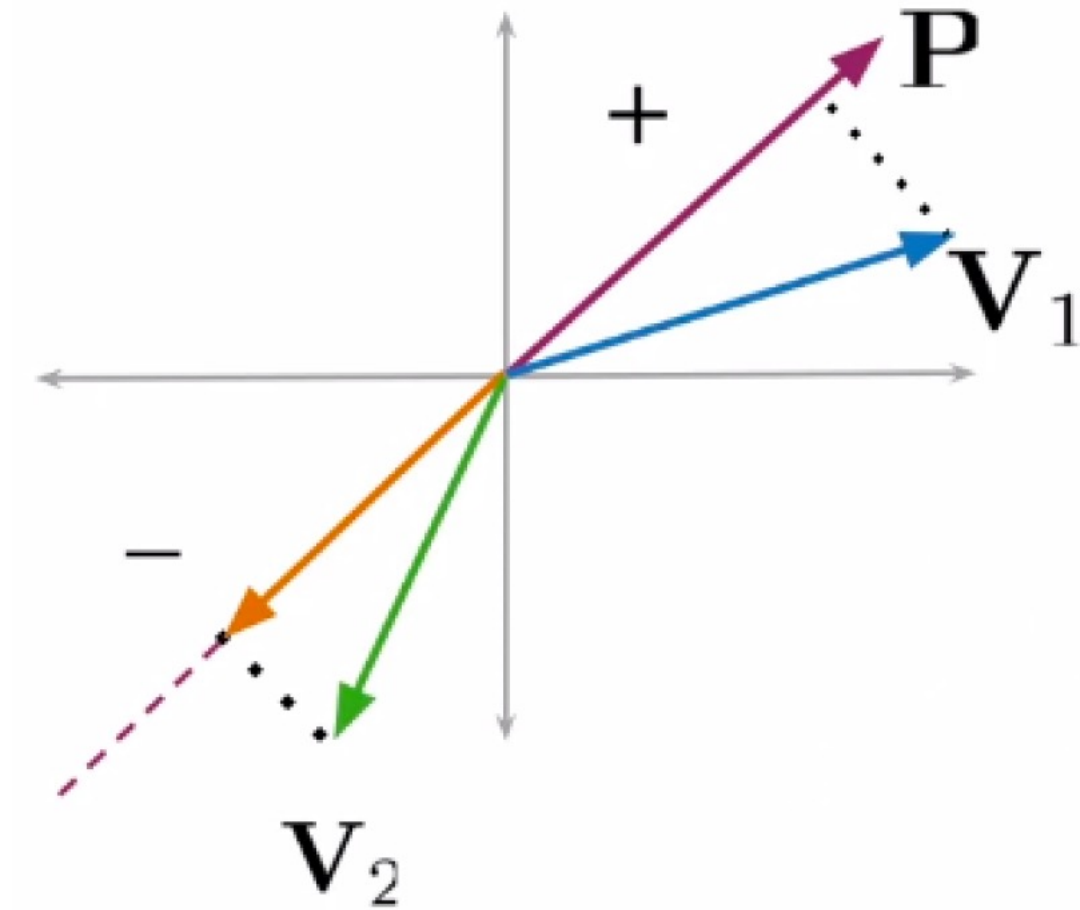
Visualizando un producto punto



Visualizando un producto punto



Visualizando un producto punto



El signo indica dirección

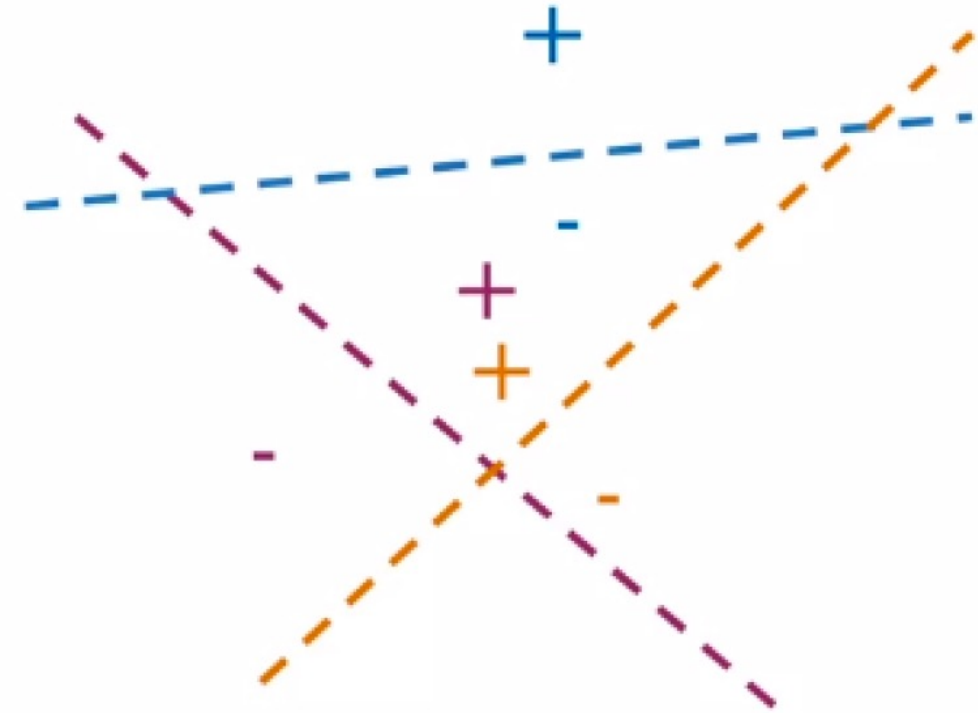
¿Cuál lado del plano?

```
def side_of_plane(P,v):  
    dotproduct = np.dot(P,v.T)  
    sign_of_dot_product = np.sign(dotproduct)  
    sign_of_dot_product_scalar= np.asscalar(sign_of_dot_product)  
    return sign_of_dot_product_scalar
```

Múltiples planos

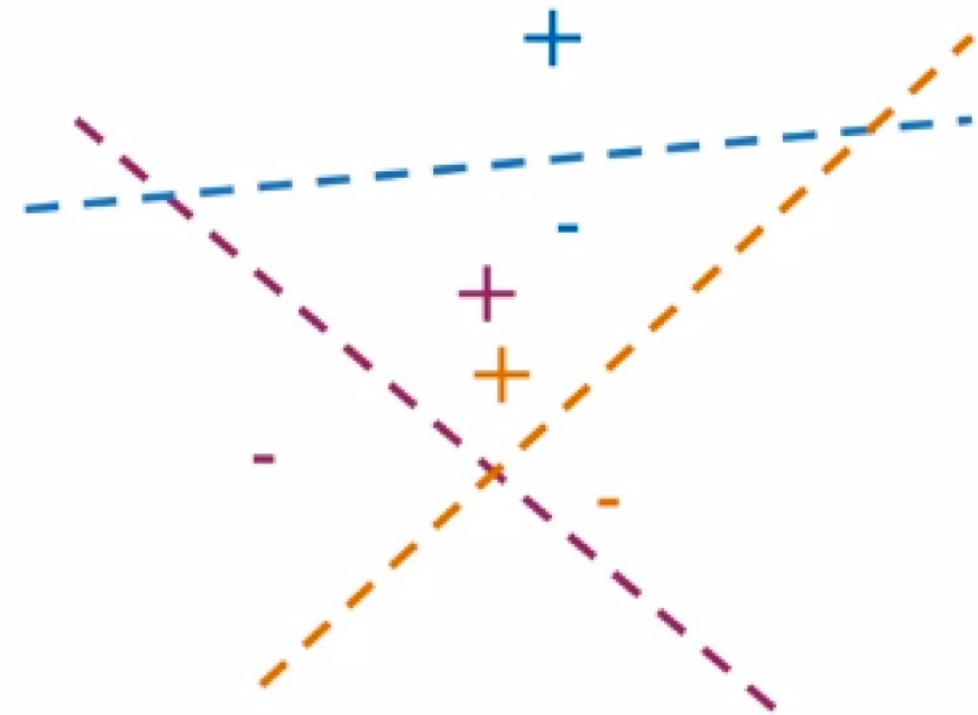
Múltiples planos → Productos puntos → Valores Hash

Múltiples planos, valor hash único?



$$\mathbf{P}_1 \mathbf{v}^T = 3, \text{sign}_1 = +1, h_1 = 1$$

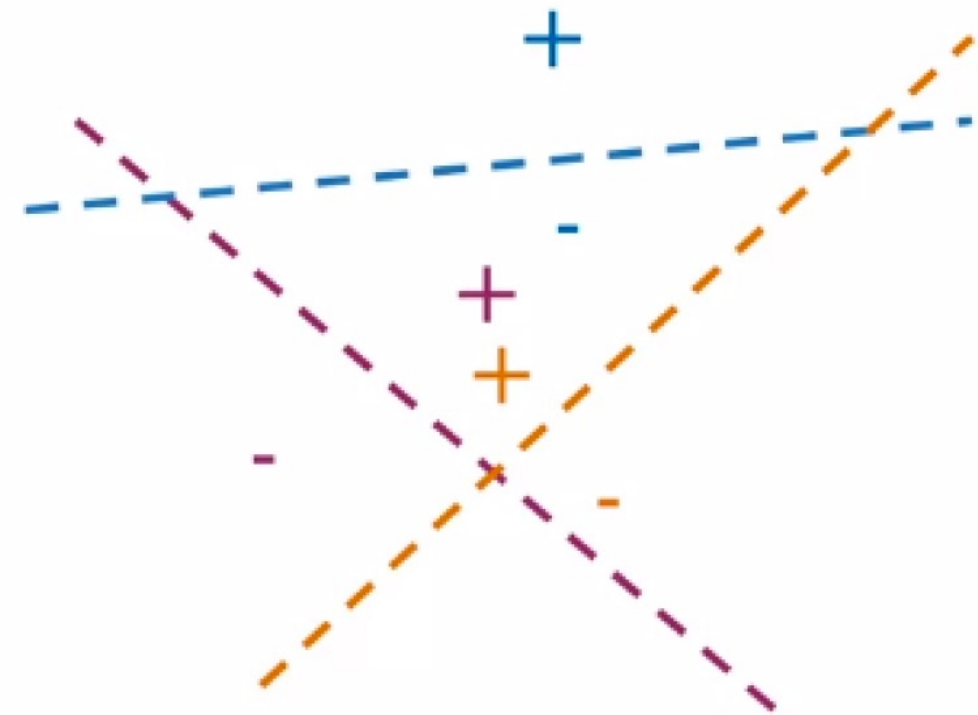
Múltiples planos, valor hash único?



$$\mathbf{P}_1 \mathbf{v}^T = 3, \text{sign}_1 = +1, h_1 = 1$$

$$\mathbf{P}_2 \mathbf{v}^T = 5, \text{sign}_2 = +1, h_2 = 1$$

Múltiples planos, valor hash único?

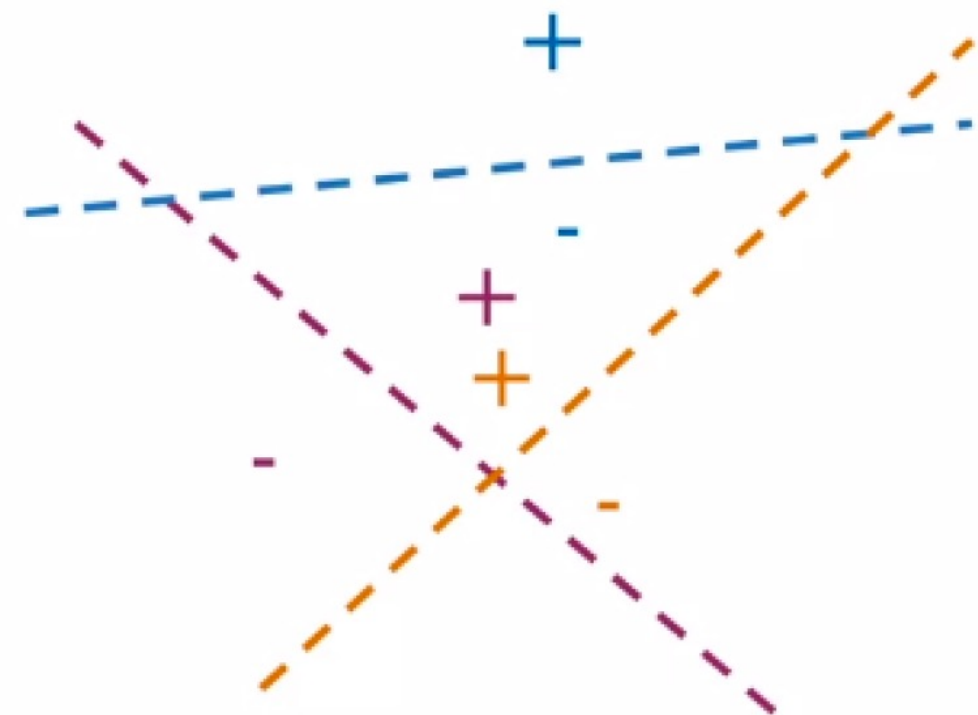


$$\mathbf{P}_1 \mathbf{v}^T = 3, \text{sign}_1 = +1, h_1 = 1$$

$$\mathbf{P}_2 \mathbf{v}^T = 5, \text{sign}_2 = +1, h_2 = 1$$

$$\mathbf{P}_3 \mathbf{v}^T = -2, \text{sign}_3 = -1, h_3 = 0$$

Múltiples planos, valor hash único?



$$\mathbf{P}_1 \mathbf{v}^T = 3, \text{sign}_1 = +1, h_1 = 1$$

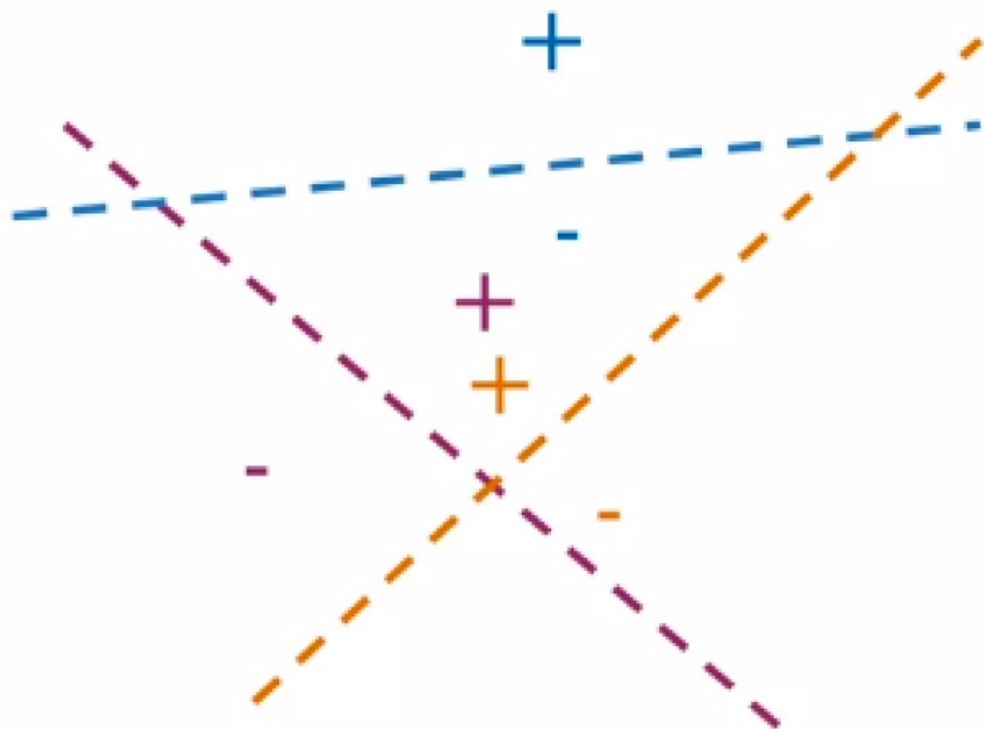
$$\mathbf{P}_2 \mathbf{v}^T = 5, \text{sign}_2 = +1, h_2 = 1$$

$$\mathbf{P}_3 \mathbf{v}^T = -2, \text{sign}_3 = -1, h_3 = 0$$

$$\begin{aligned} \text{hash} &= 2^0 \times h_1 + 2^1 \times h_2 + 2^2 \times h_3 \\ &= 1 \times 1 + 2 \times 1 + 4 \times 0 \end{aligned}$$

$$= 3$$

Múltiples planos, valor hash único!



$$\text{sign}_i \geq 0, \rightarrow h_i = 1$$

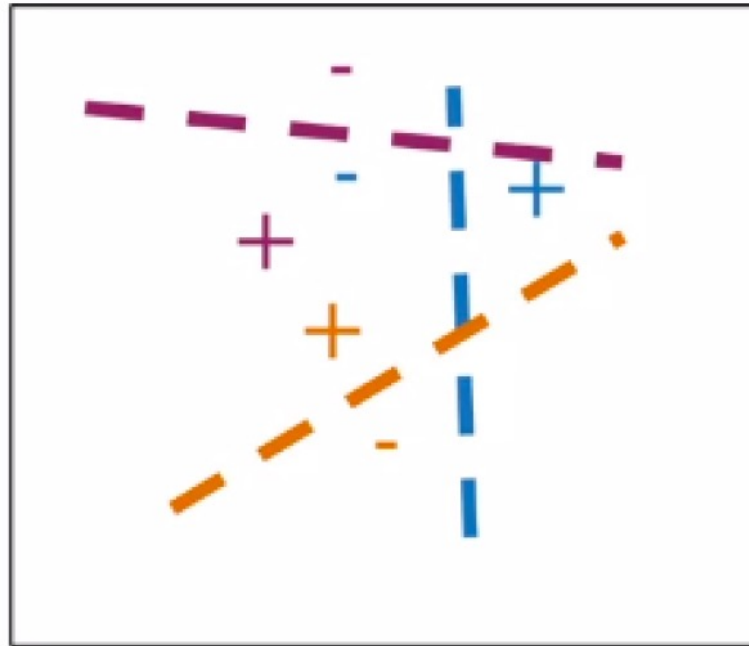
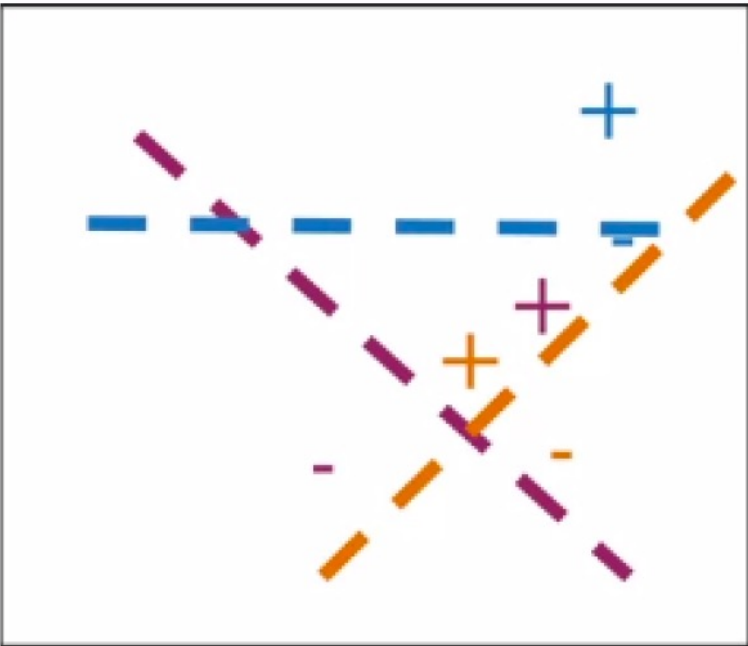
$$\text{sign}_i < 0, \rightarrow h_i = 0$$

$$\text{hash} = \sum_i^H 2^i \times h_i$$

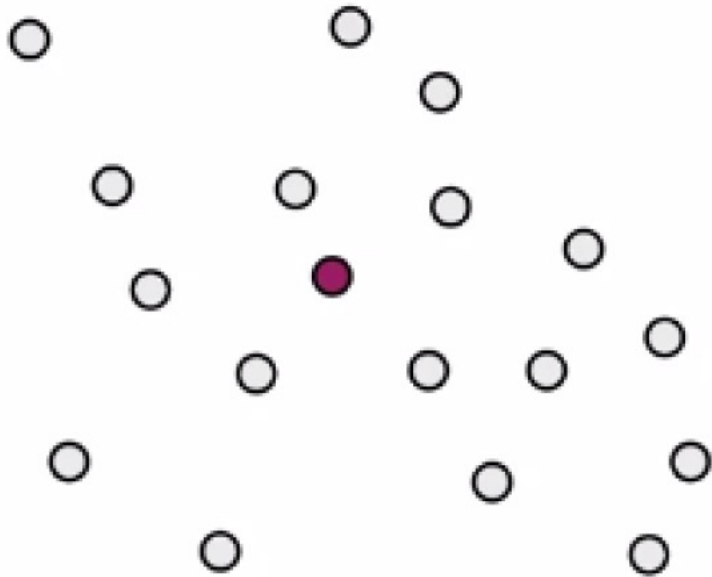
Múltiples planos, valor hash único!!

```
def hash_multiple_plane(P_l, v):  
    hash_value = 0  
  
    for i, P in enumerate(P_l):  
        sign = side_of_plane(P, v)  
        hash_i = 1 if sign >= 0 else 0  
        hash_value += 2**i * hash_i  
  
    return hash_value
```

Planos aleatorios

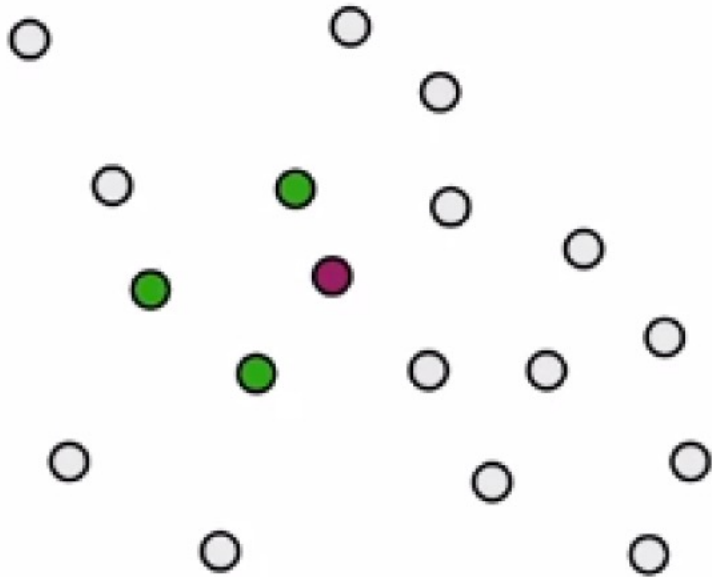


Múltiples conjuntos de planos aleatorios



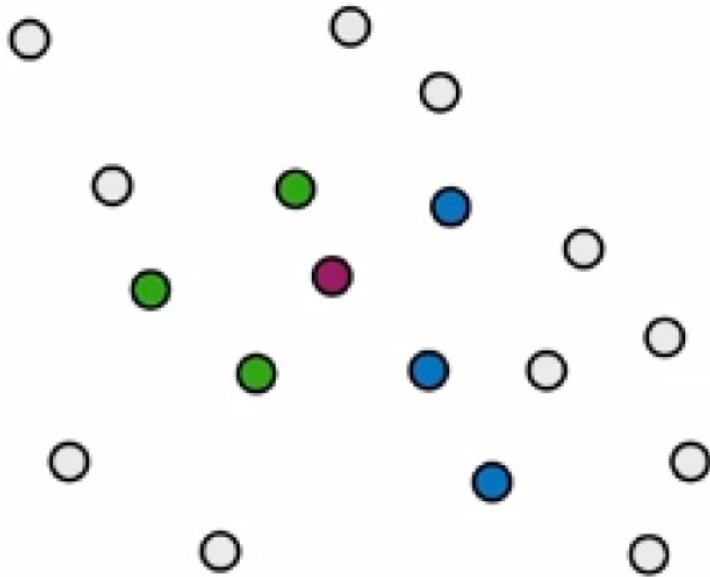
Supongamos que estamos tratando de encontrar el vecino más cercano para el vector rojo (punto).

Múltiples conjuntos de planos aleatorios



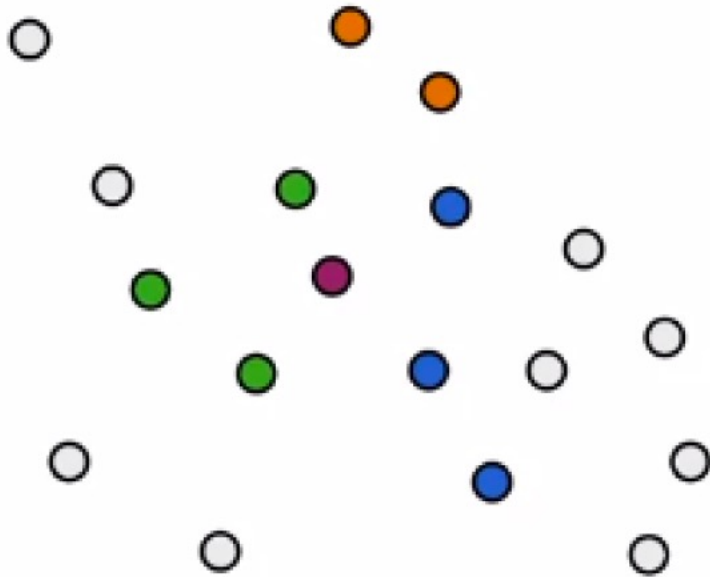
La primera vez, el plano nos da los puntos verdes

Múltiples conjuntos de planos aleatorios



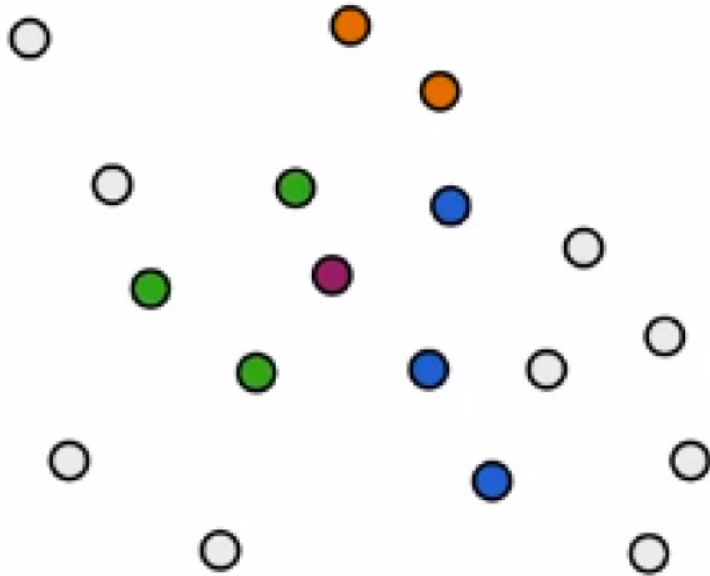
Luego, volvemos a ejecutar el lado del plano con otro plano aleatorio y el algoritmo me da los puntos azules

Múltiples conjuntos de planos aleatorios



La tercera vez, obtuvimos los puntos naranjas

Múltiples conjuntos de planos aleatorios



Si seguimos realizando este proceso, es muy probable que obtengamos todos los vecinos cercanos

↓
Vecinos cercanos
aproximados

Hacer un conjunto de planos aleatorios

```
num_dimensions = 2 #300 in assignment
num_planes = 3 #10 in assignment

random_planes_matrix = np.random.normal(
    size=(num_planes,
          num_dimensions))
```

```
array([[ 1.76405235  0.40015721]
       [ 0.97873798  2.2408932 ]
       [ 1.86755799 -0.97727788]])
```

```
v = np.array([[2,2]])
```

Hacer un conjunto de planos aleatorios

```
num_dimensions = 2 #300 in assignment
num_planes = 3 #10 in assignment

random_planes_matrix = np.random.normal(
    size=(num_planes,
          num_dimensions))
```

```
array([[ 1.76405235  0.40015721]
       [ 0.97873798  2.2408932 ]
       [ 1.86755799 -0.97727788]])
```

```
v = np.array([[2,2]])
```

```
def side_of_plane_matrix(P,v):
    dotproduct = np.dot(P,v.T)
    sign_of_dot_product = np.sign(dotproduct)
    return sign_of_dot_product
```

```
num_planes_matrix = side_of_plane_matrix(
    random_planes_matrix,v)
```

```
array([[1.]
       [1.]
       [1.]])
```

Representación de documentos

I love learning!

[?, ?, ?]

I

[1, 0, 1]

love

[-1, 0, 1]

learning

[1, 0, 1]

Representación de documentos

I love learning!	$[?, ?, ?]$
I	$[1, 0, 1]$
	+
love	$[-1, 0, 1]$
	+
learning	$[1, 0, 1]$
	=
I love learning!	$[1, 0, 3]$