

# ④ Pre-train word and phrases vectors (Google News word embedding dataset)

100 Billion Words  $\rightarrow$  Train  $\rightarrow$  3 Million words and phrases with 300 dimensional vectors.

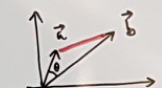
ex. 

word	1	2	...	300
word1	0.01	-0.02	...	0.3
...				
wordM				

 $\rightarrow$  300 dimensional vector

## • Cosine Similarity of $\vec{a}, \vec{b}$

$$\theta \in [0^\circ \div 180^\circ] \Rightarrow \cos \theta \in [-1, 1] = [-1 \quad 0 \quad 1]$$



$$\cos \theta = \frac{\vec{a} \cdot \vec{b}}{\|\vec{a}\| \|\vec{b}\|} = \frac{\sum_{i=1}^N a_i \cdot b_i}{\sqrt{\sum_{i=1}^N a_i^2} \sqrt{\sum_{i=1}^N b_i^2}} = \frac{\text{np.dot}(\vec{a}, \vec{b})}{\text{np.linalg.norm}(\vec{a}) \times \text{np.linalg.norm}(\vec{b})}$$

where  $\vec{a}, \vec{b}$  are embedding (word) vectors.

## • Euclidean Distance

$$d(\vec{a}, \vec{b}) = d(\vec{b}, \vec{a}) = \sqrt{\sum_{i=1}^N (a_i - b_i)^2} \quad \text{Ex in 2D: } d(\vec{a}, \vec{b}) = \sqrt{(a_1 - b_1)^2 + (a_2 - b_2)^2}$$

$= \text{np.sqrt}(\text{np.sum}(\vec{a} - \vec{b} ** 2))$   $\leftarrow$  Note: numpy broadcasting

• Predict: Ex. Canberra is to Australia as Hanoi is to ?

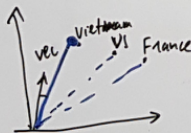
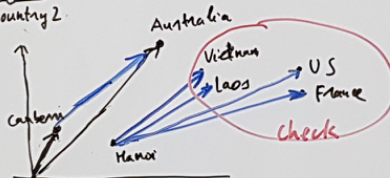
def get\_country(city1, country1, city2, embeddings):  
 embedding = {  
 "Canberra": [ ... ],  
 "Australia": [ ... ],  
 "Hanoi": [ ... ],  
 "US": [ ... ]  
 }

$$\text{Country1} - \text{city1} = \text{country2} - \text{city2}$$

$$\Rightarrow \text{vec} = \text{Country1} - \text{city1} + \text{city2}$$

$\Rightarrow$  Just choose the country that has maximum of cosine-similarity (vec, country)

• Test accuracy with given (known) set of (country, city)



## PCA

Eigen vector: Uncorrelated features of the data (orthogonal)

Eigen Value: The amount of information retained by each feature (variance)

Mean Normalized data (Z-score):  $\bar{x}_i = \frac{x_i - \mu_i}{\sigma_i}$

$\downarrow$   
Covariance Matrix  $\Sigma$  where  $\Sigma_{ij} = \text{cov}_{ij} = \frac{1}{N-1} \sum_{i=1}^N \frac{\sigma_i}{\sigma_j} (x_i - \bar{x}_i)(x_j - \bar{x}_j)$

$\downarrow$   
Perform SVD (Singular Value Decomposition):  $SVD(\Sigma)$

$\downarrow$   
Eigen Vector      Eigen Value

$$\begin{bmatrix} | & & | \\ | & & | \\ | & & | \end{bmatrix} \begin{bmatrix} \square & & \square \\ & \square & \\ & & \square \end{bmatrix} = \begin{bmatrix} \text{---} \\ \text{---} \\ \text{---} \end{bmatrix}$$