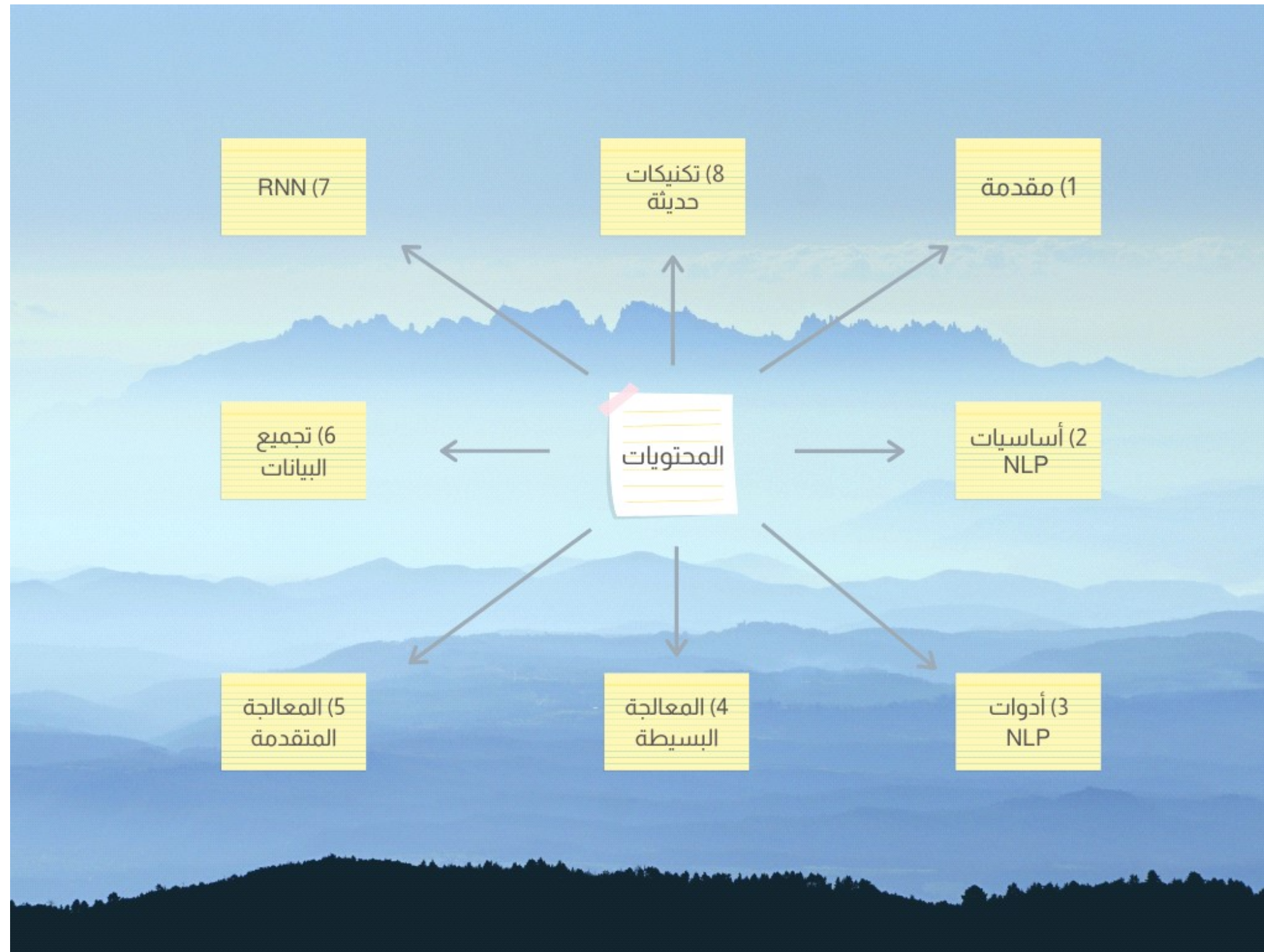


NATURAL LANGUAGE PROCESSING

المعالجة اللغوية الطبيعية



المحتويات

				التطبيقات	العقبات و التحديات	تاريخ NLP	ما هو NLP	المحتويات	1) مقدمة
					البحث في النصوص	ملفات pdf	الملفات النصية	المكتبات	2) أساسيات NLP
T.Visualization	Syntactic Struc.	Matchers	Stopwords	NER	Stem & Lemm	POS	Sent. Segm.	Tokenization	3) أدوات NLP
	Dist. Similarity	Text Similarity	TF-IDF	BOW	Word2Vec	T. Vectors	Word embed	Word Meaning	4) المعالجة البسيطة
T. Generation	NGrams	Lexicons	GloVe	L. Modeling	NMF	LDA	T. Clustering	T. Classification	5) المعالجة المتقدمة
	Summarization & Snippets		Ans. Questions		Auto Correct	Vader	Naïve Bayes	Sent. Analysis	
Search Engine	Relative Extraction		Information Retrieval		Information Extraction		Data Scraping	Tweet Collecting	6) تجميع البيانات
					Rec NN\TNN	GRU	LSTM	Seq to Seq	7) RNN
Chat Bot	Gensim	FastText	Bert	Transformer	Attention Model	T. Forcing	CNN	Word Cloud	8) تكنيكات حديثة

القسم الرابع : المعالجة البسيطة للنصوص

الجزء الثالث : Text Vectors

=====

نتناول في هذا الجزء , العمليات الحسابية المعتمدة علي قيم word embedding لكل كلمة , وذلك لتحديد كلمات مطلوبة , بناء علي عدد آخر من الكلمات . .

مثال بسيط : لو كان لدينا عاصمة الولايات المتحدة و نريد استنتاج عاصمة روسيا , فكيف يتم هذا ؟



USA



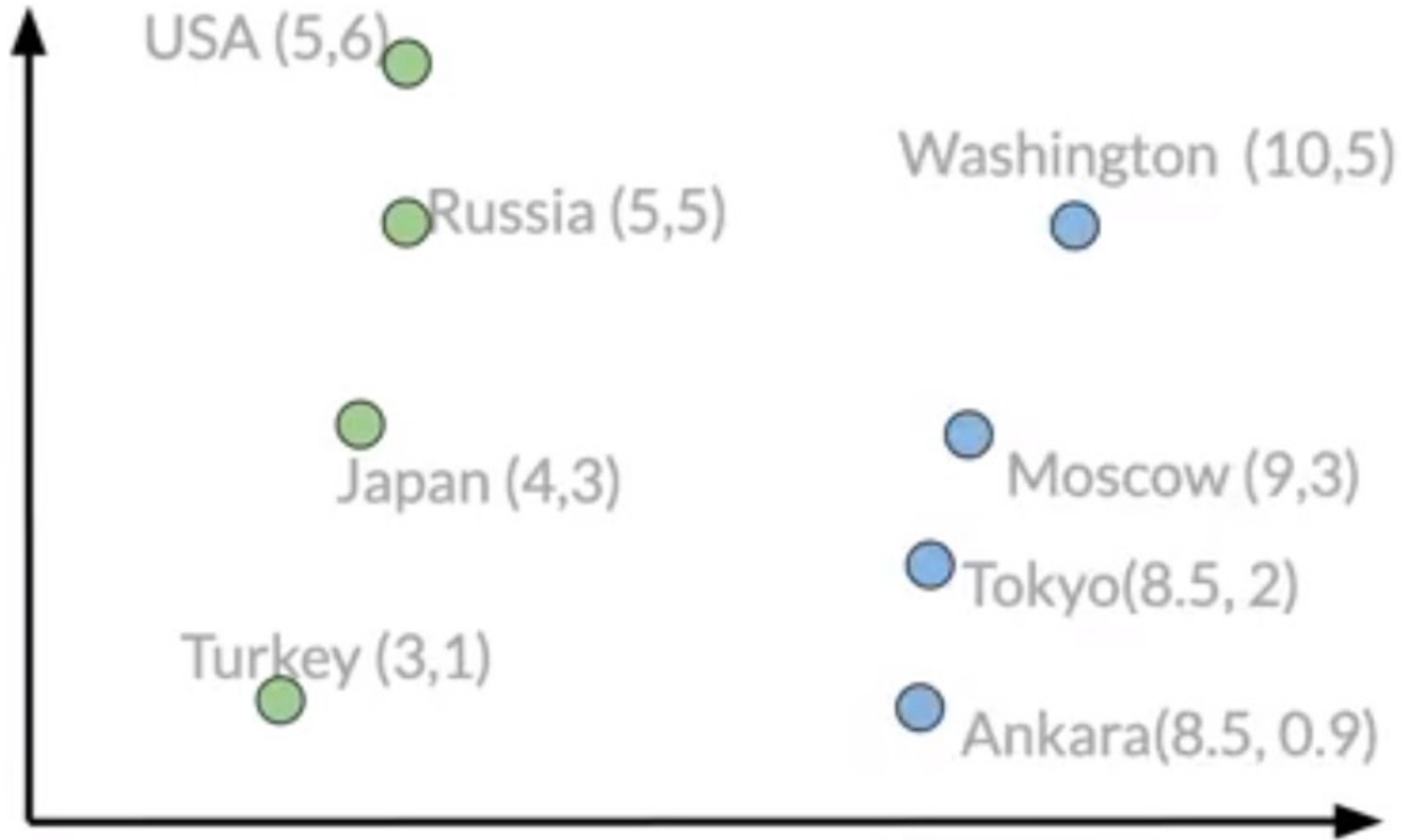
Washington
DC

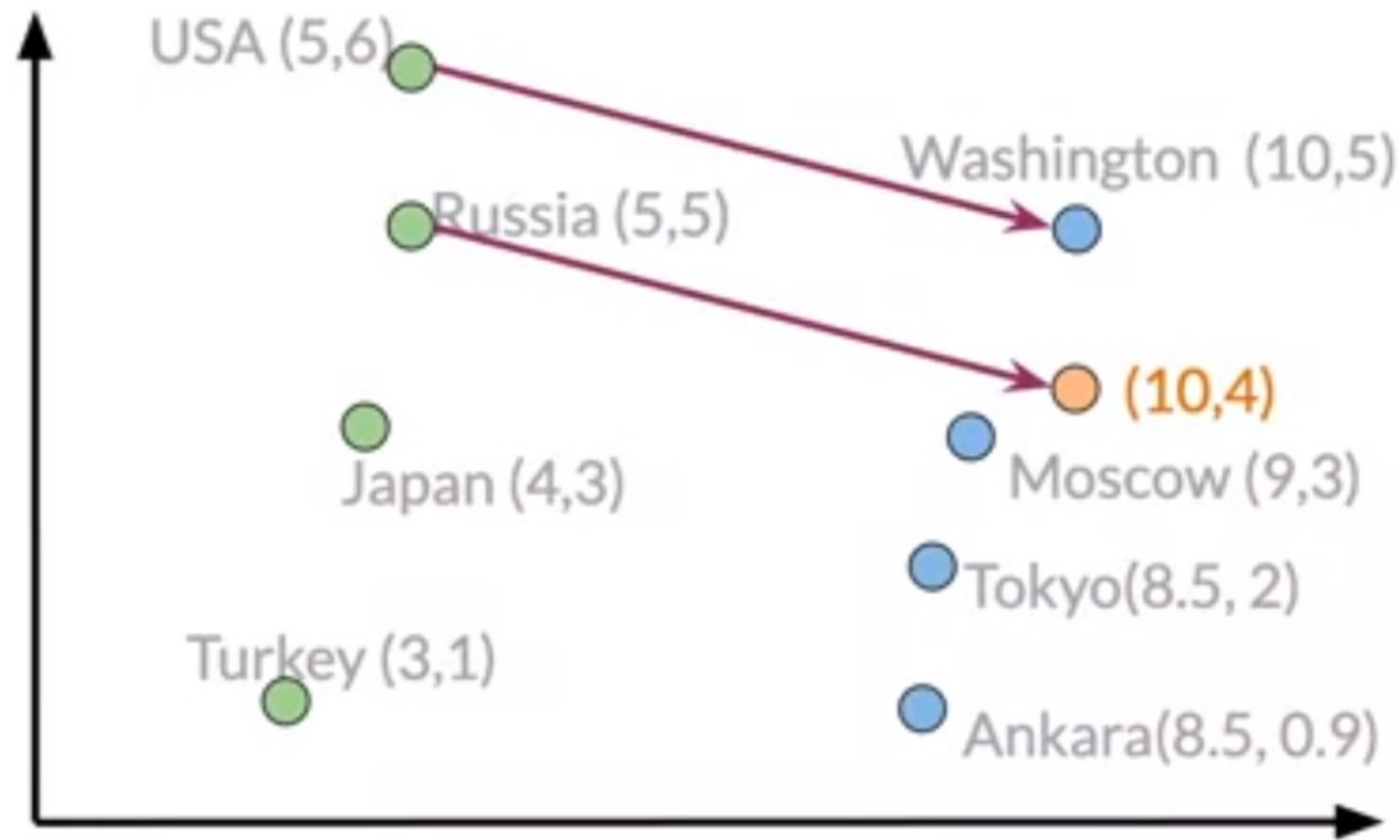


Russia



يتم عبر رسم الدول و العواصم في جراف معين , فلو تخيلنا ان قيم WE هي قيمتين فقط مثل هذا الرسم , وتم تحديد اماكن الدول و المدن الهمة هنا , و تم تطبيق عملية طرح بين امريكا و عاصمتها , التي ستكون -1 , 5 ثم مع هذا الفارق مع قيم روسيا , والتي ستكون بقيمة معينة 10,4





$$\text{Washington} - \text{USA} = \begin{bmatrix} 5 & -1 \end{bmatrix}$$

$$\text{Russia} + \begin{bmatrix} 5 & -1 \end{bmatrix} = \begin{bmatrix} 10 & 4 \end{bmatrix}$$

و بالبحث , نجد ان اقرب مدينة لهذه القيمة هو موسكو , و يمكن تطبيق هذا الامر علي عدد اخر من البلاد

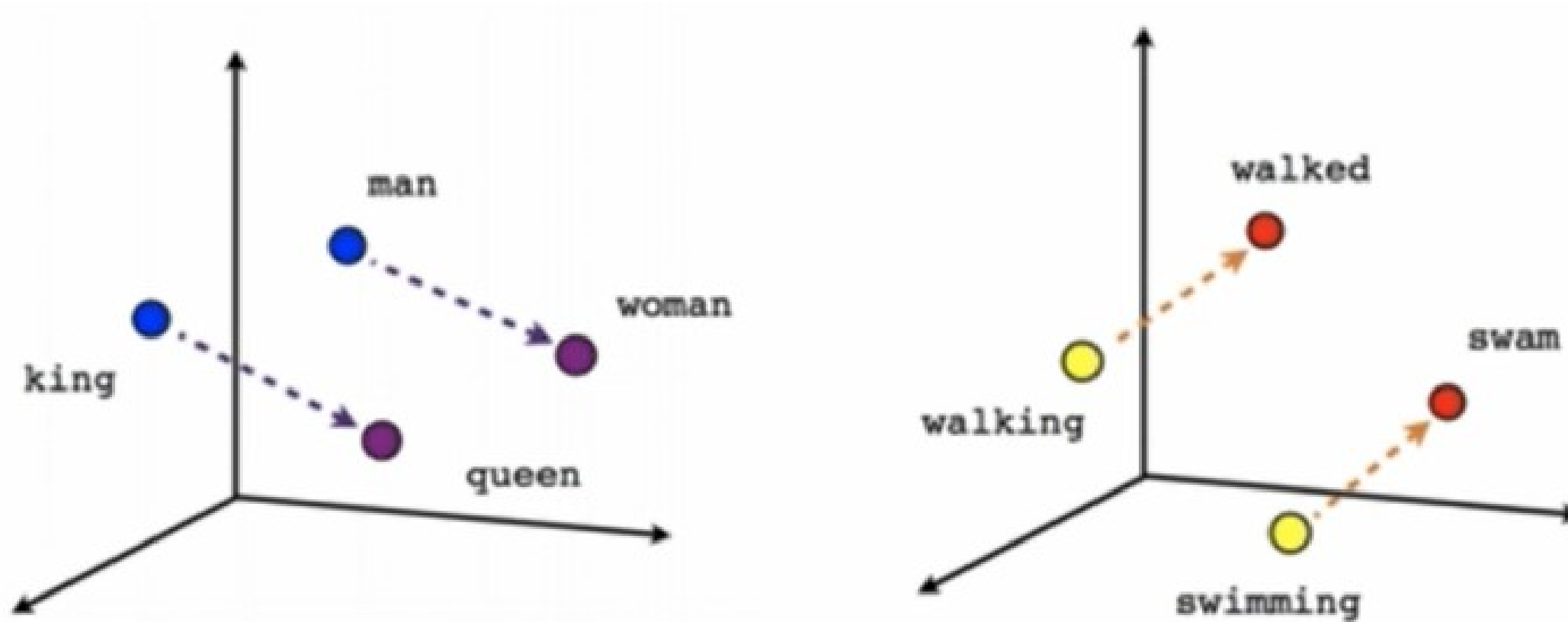
و هذا الأمر لو تم علي مقياس صغير مثل بعدين فقط , فيمكن أن يتم علي حجم أكبر , مثل 300 بعد الخاص بالـ WE

كذلك يمكن عمل جمع او طرح للجملة المشهورة, فيمكن مثلا ان نقول ان :

$$X = \text{King} - \text{Man} + \text{Woman}$$

قيمة X هنا هي ان نقوم بايجاد الفارق بين الملك و الرجل , والتي ستكون صفات الملوك , واذا اضعناها الي المرأة , ستكون هي قيمة الملكة . . و عمليات الجمع و الطرح هي في الحقيقة لكل الفيكتورز الـ 300

كذلك نفس الامر في الازمنة في الافعال وغيرها هكذا



و هذا الأمر يسمى word analogy او تشابه الكلمات

و هي الخاصة بمدي اقتراب او ابتعاد الكلمات عن بعضها البعض , و مدي الاختلافات او التشابهات بينها , وهناك امثلة اخري
مثل :

King - Queen ~= Prince - Princess

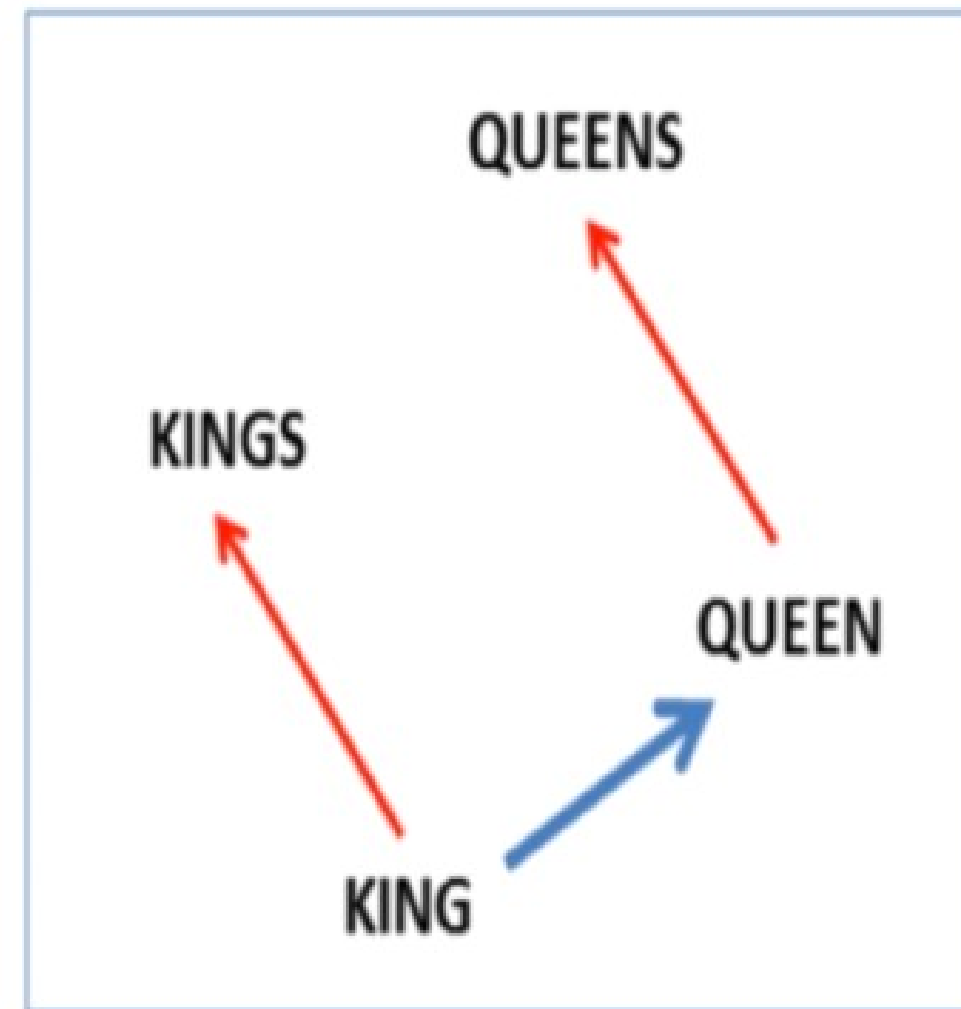
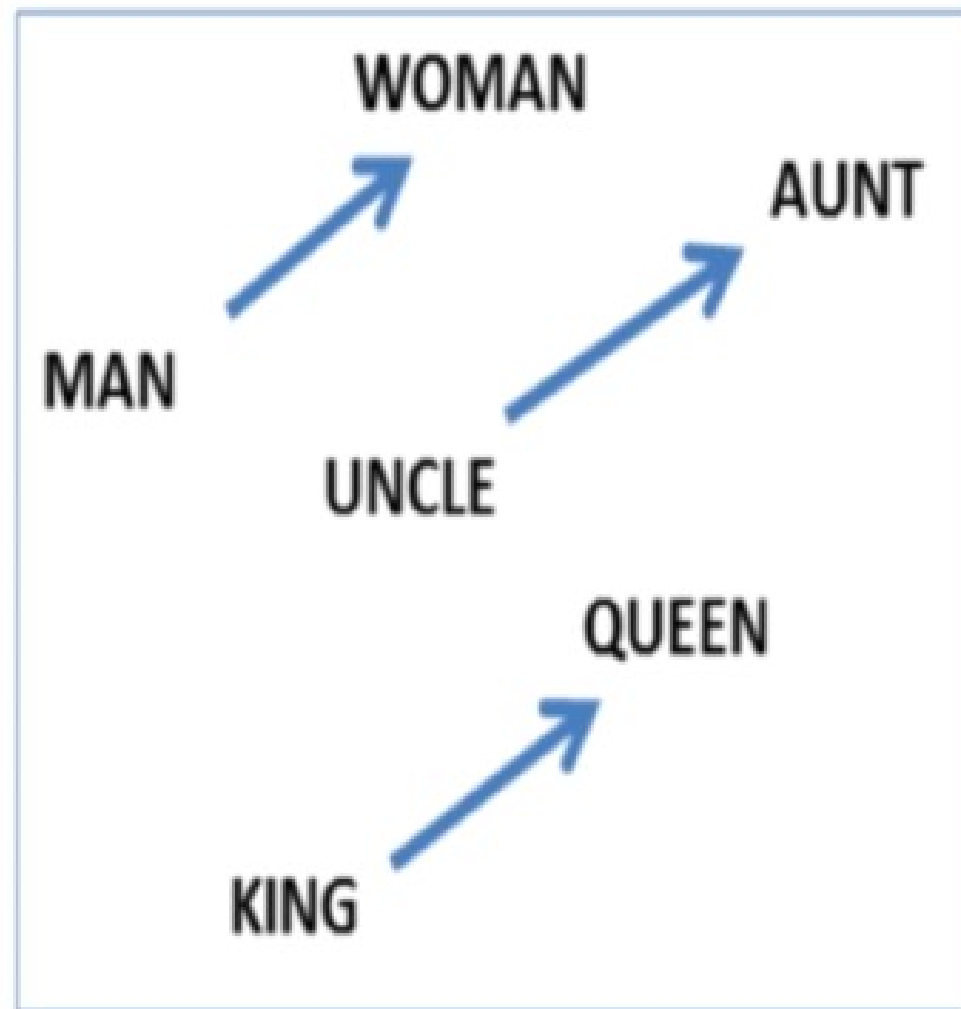
France - Paris ~= Germany - Berlin

Japan - Japanese ~= China - Chinese

Brother - Sister ~= Uncle - Aunt

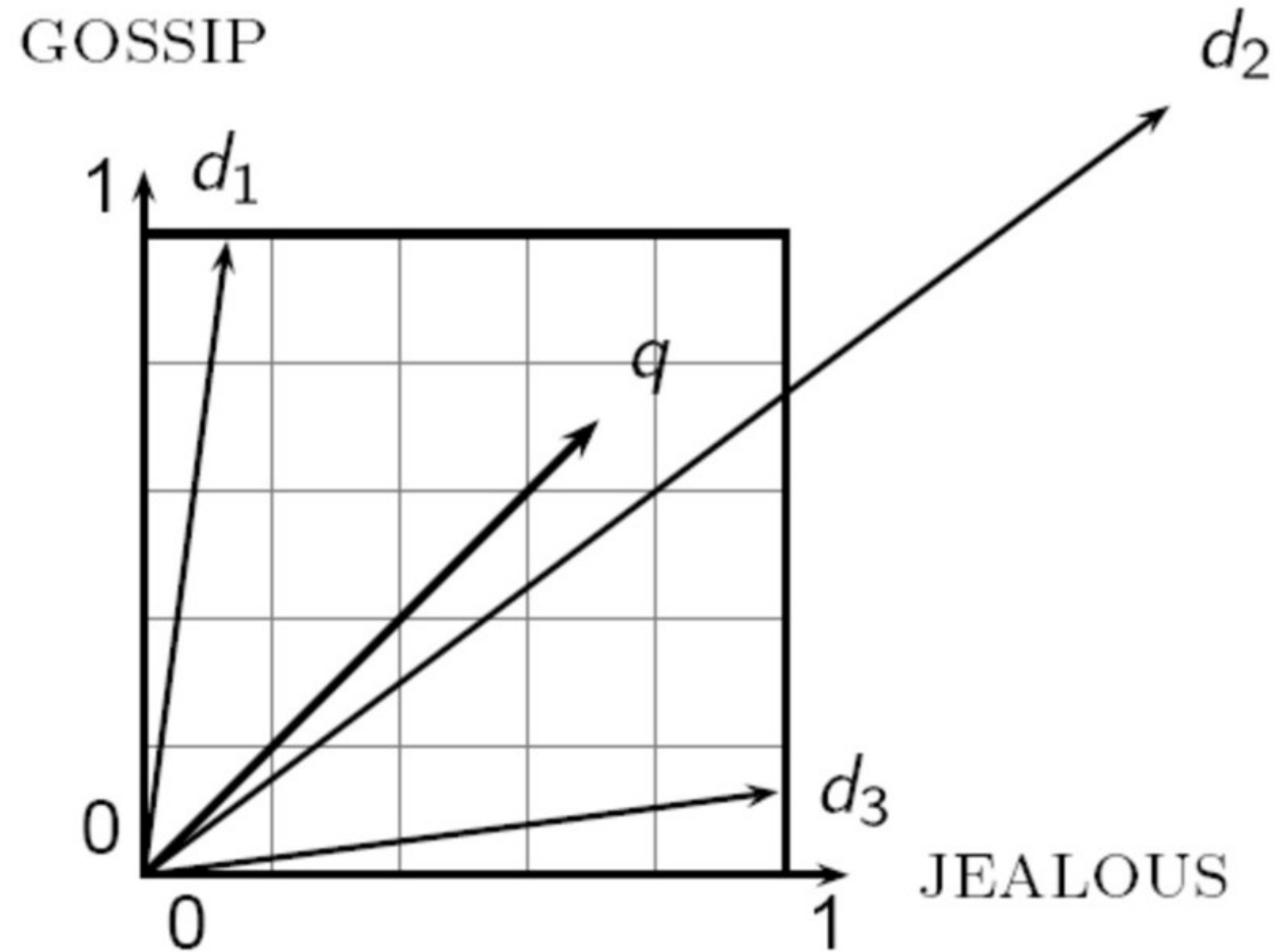
Walk - Walking ~= Swim - Swimming

و التي تكون بهذا الشكل



الان , كيف يمكن حساب مدي اقتراب او ابتعاد كلمتين عن بعضهما البعض ؟ ؟

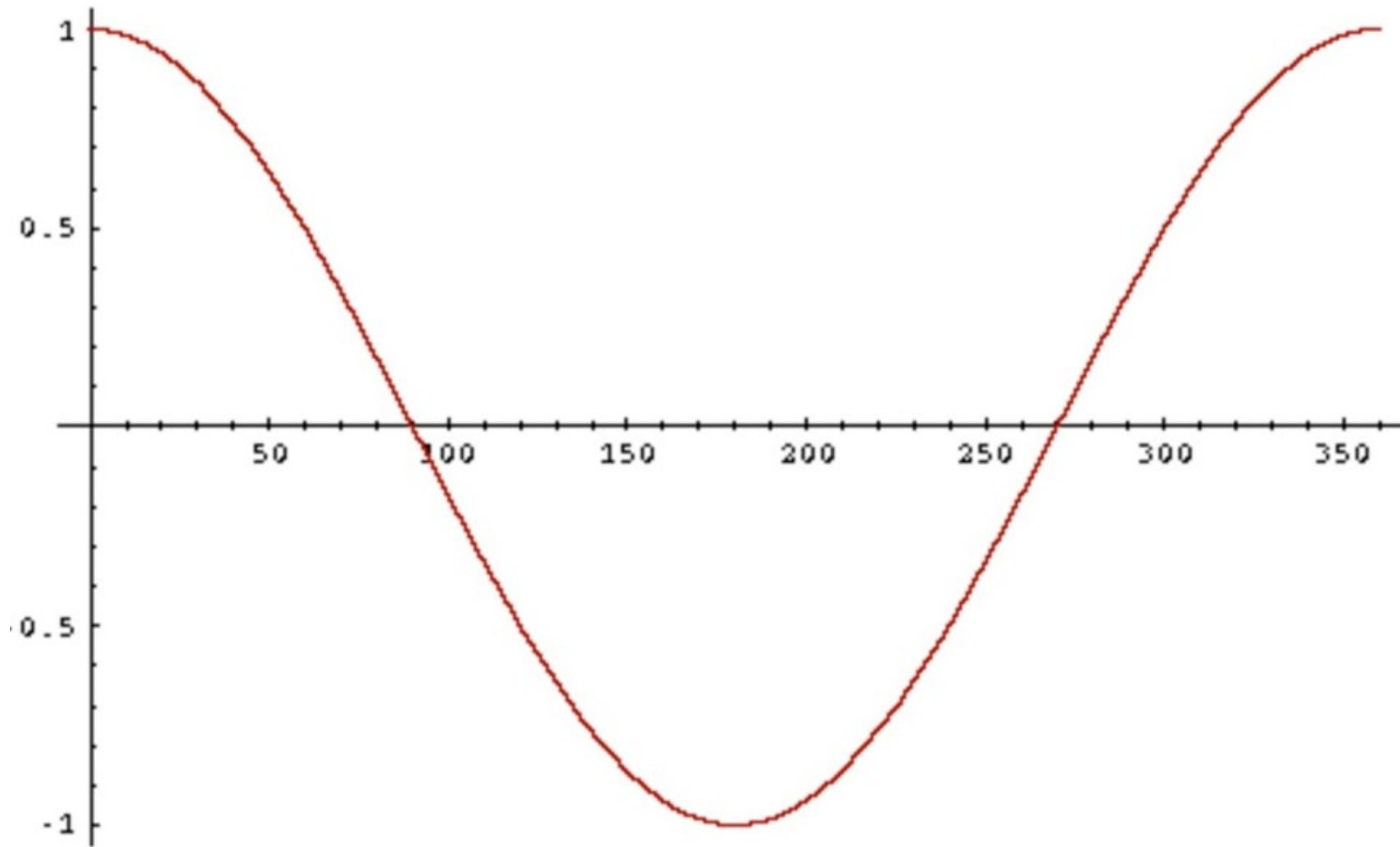
التناول الأول هو المسافة بأسلوب Euclidean distance , وهي أن يتم حساب المسافة بين النقطتين بعضها , بحيث تكون المسافة بين كل متجه و الثاني , يدل علي العلاقة بين الكلمتين هكذا:



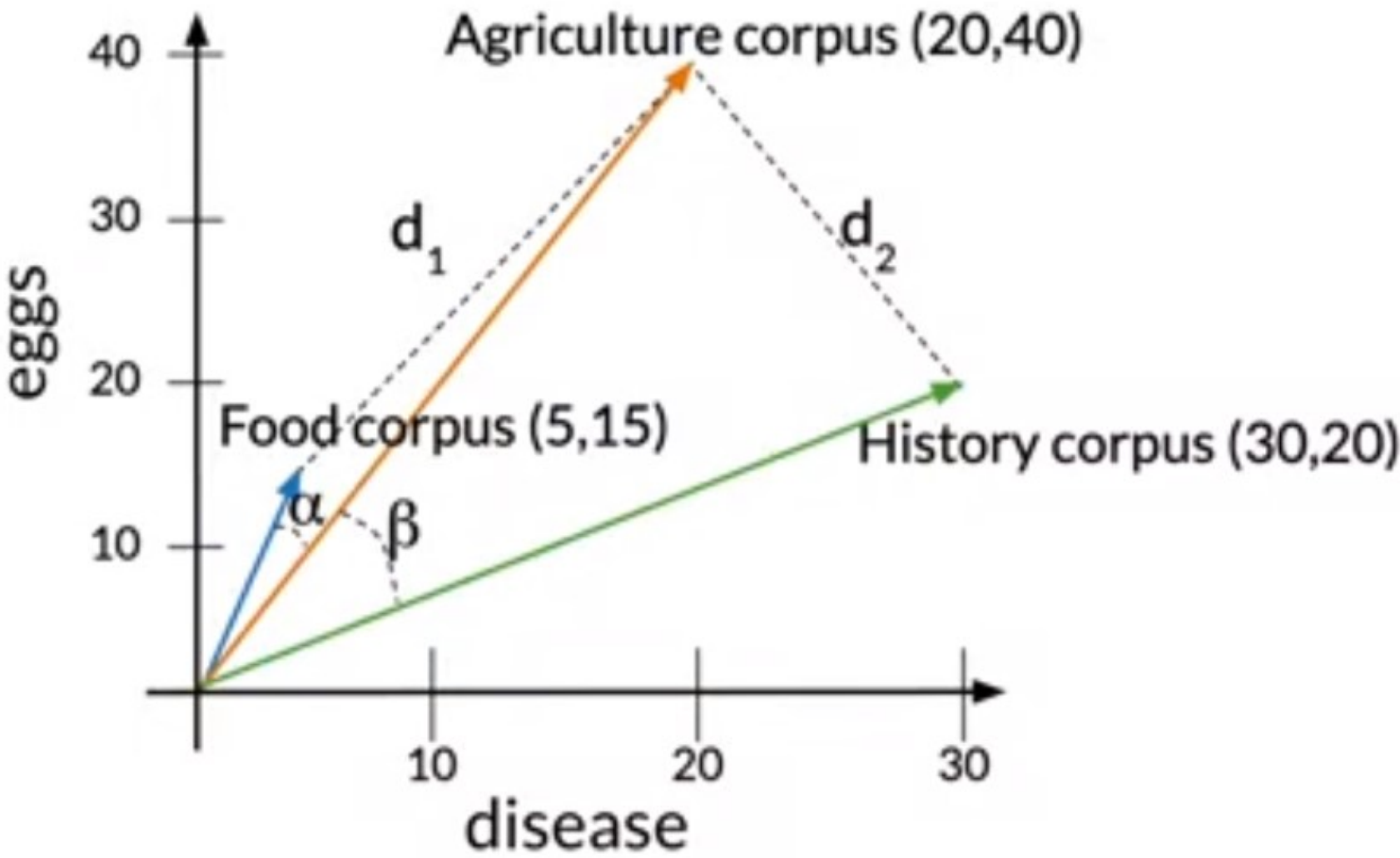
لكن المشكلة ان المسافة بين نقطتين لا يعبر عن الفرق الحقيقي بينهما , لأن المسافة بين q و $d2$ اكبر من مسافتها مع $d1$ او $d3$, علي الرغم من ان التطابق يكاد يكون بينها و بين $d2$

لذا فاننا نعتمد علي الزاوية بينهم , و بالتحديد قيمة \cos الزاوية فيما يسمى cosine similarity

حيث أن الكلمتين لو كانا متطابقتين فتكون الزاوية بينهم 0 فتكون القيمة 1 , وكلما ابتعدا لما زادت قيمة \cos حتي تصل ل 0 مع زاوية 90 و سالب 1 مع الزاوية العكسية , فيكون الترتيب تنازليا حسب المسافة , وتصاعديا حسب قيمة \cos هكذا :



و هنا مثال آخر للفارق بين cosine similarity و Euclidean distance



Euclidean distance: $d_2 < d_1$
Angles comparison: $\beta > \alpha$

The cosine of the angle between the vectors

* * * * *

اذن كيف يتم حساب قيمة cosine similarity خاصة اننا ليس لدينا قيمة للزاوية بين الكلمات بالفعل

اولا نتذكر معلومة هامة و هي vector norm & dot products و التي يتم حسابها هكذا

Vector norm

$$\|\vec{v}\| = \sqrt{\sum_{i=1}^n v_i^2}$$

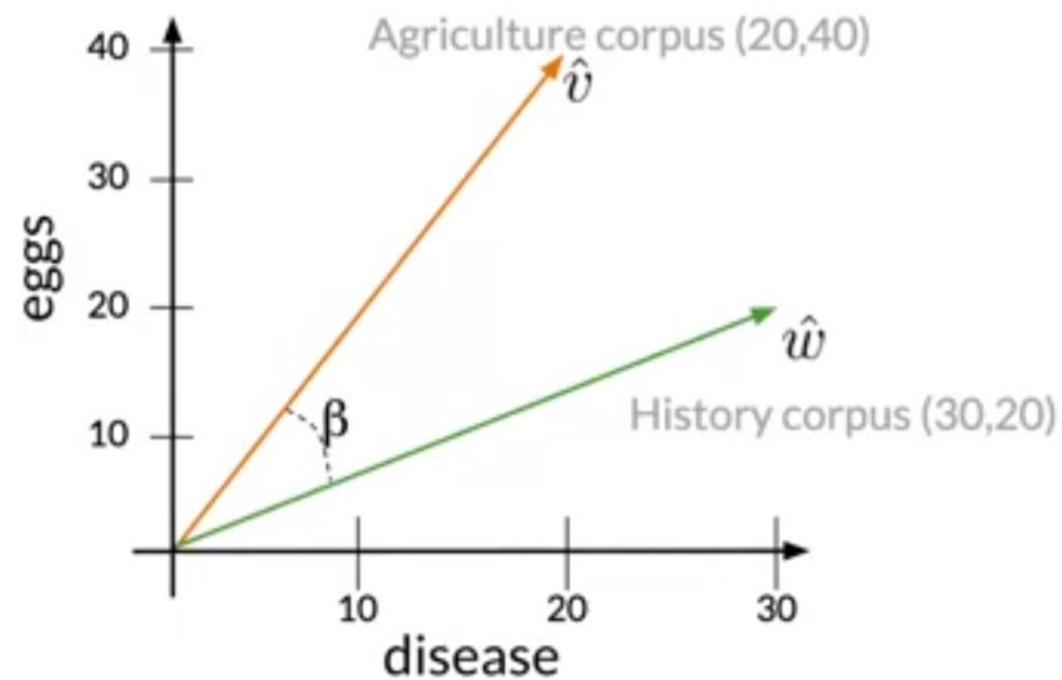
Dot product

$$\vec{v} \cdot \vec{w} = \sum_{i=1}^n v_i \cdot w_i$$

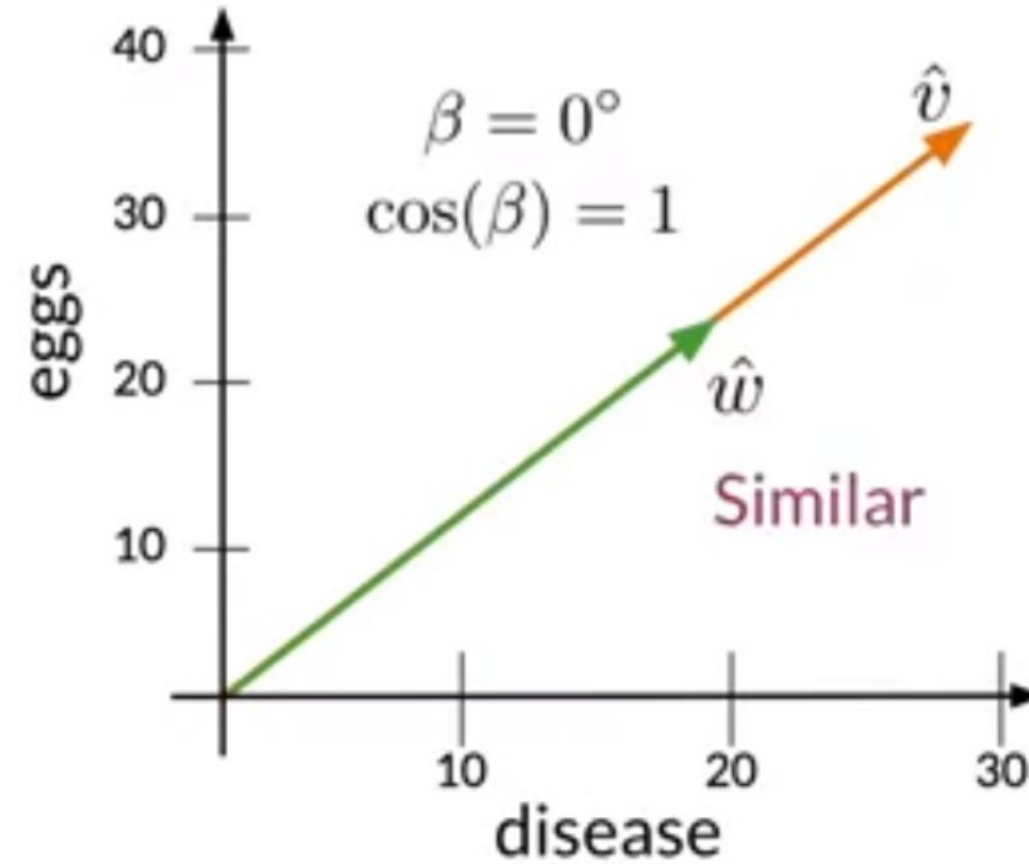
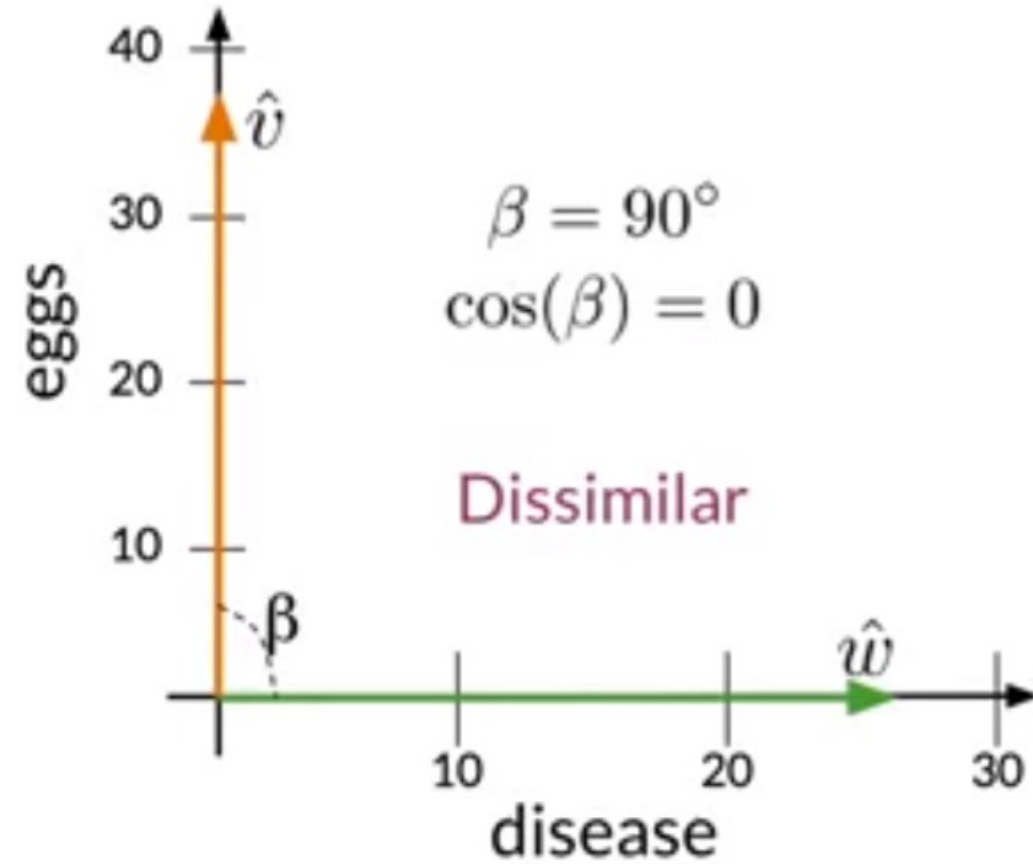
و نستخدم هذه القاعدة المعروفة لحساب قيمة ال COS

$$\hat{v} \cdot \hat{w} = \|\hat{v}\| \|\hat{w}\| \cos(\beta)$$

$$\cos(\beta) = \frac{\hat{v} \cdot \hat{w}}{\|\hat{v}\| \|\hat{w}\|}$$



و لا تنس ان لو تعاملنا مع كلمتين بعيدتين عن بعضهما فتكون الزاوية 90 و يكون ال \cos يساوي 0 اي مختلفين
و لو تعاملنا مع كلمتين متطابقين عن بعضهما فتكون الزاوية 0 و يكون ال \cos يساوي 1 اي متشابهين
و اي رقم بينهما يشير الي مدي التشابه بينهم



و لن تزيد الزاوية بينهما عن 90 او تقترب من 180 لان هذا يشير الي تضاد و لا توجد قيم سالبة في عدد الكلمات