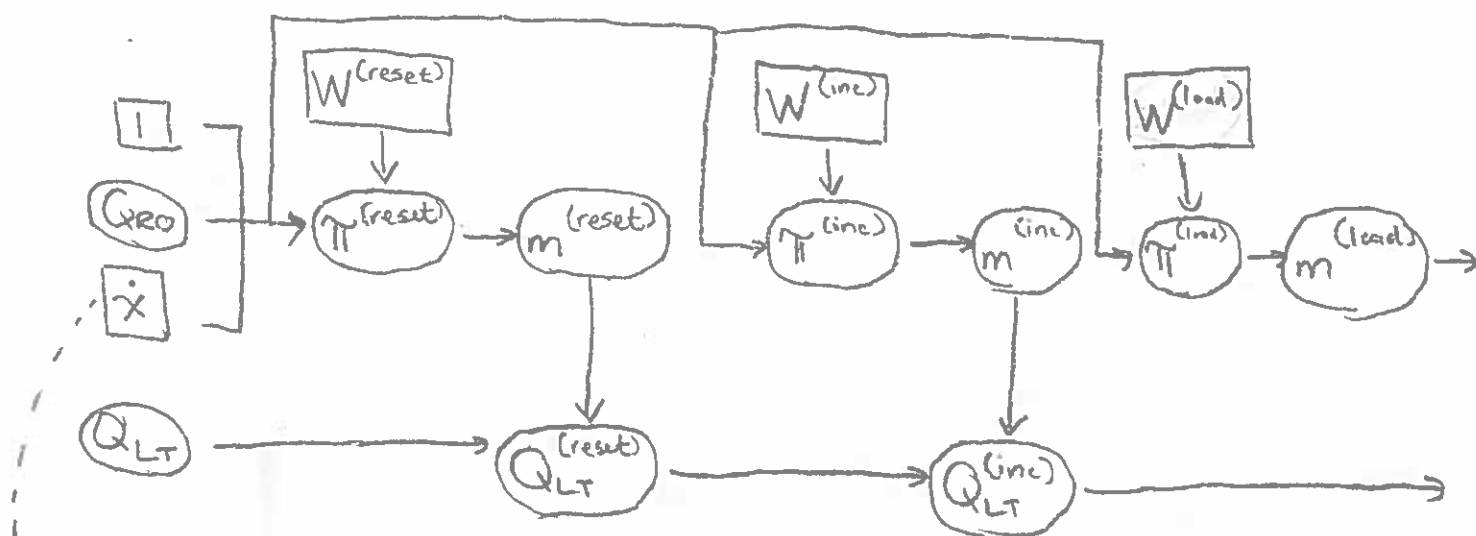


WORD VECTORS

- ① As we use increasingly complex architectures for natural language processing, a major problem manifests:



if our vocabulary contains V words,
then all the weight matrices will be $V \times H$
if we represent each word as a one-hot vector.

- ② This can make training difficult and slow. Suppose the vocabulary has 100K words. If our hidden state Q_{LT} has size 100, then $W^{(reset)}$ has 10 million parameters to learn.

With a large enough vocabulary, it's also entirely possible that these weight matrices won't fit in memory.

Word Vectors

- ③ It'd be nice to start with vector representations that more compactly encoded each word of the vocabulary.

One could imagine a representation that represents a word in terms of a selection of attributes, e.g.

$$\text{elephant} = \begin{bmatrix} 1.8 \\ 0.4 \end{bmatrix} \quad \text{dolphin} = \begin{bmatrix} 0.7 \\ 1.5 \end{bmatrix} \quad \text{truck} = \begin{bmatrix} 1.7 \\ -3.1 \end{bmatrix}$$

where each vector is:

$$\begin{bmatrix} \text{size} \\ \text{intelligence} \end{bmatrix}$$

Elephants are large and relatively intelligent, while trucks are large but rather dumb.

- ④ One way to learn these representations is with a simple neural network. It had better be simple, because we'll need to use one-hot vectors as its input (lest we enter some infinite recursion).

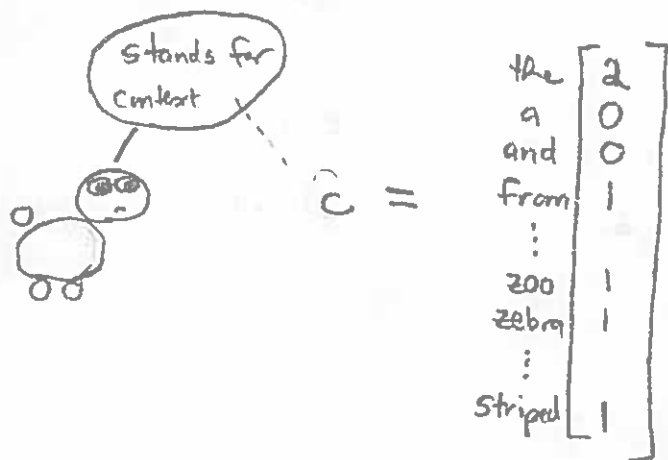
WORD VECTORS

- ⑤ A popular approach is to train a neural network to solve "fill-in-the-blank" statements, e.g.

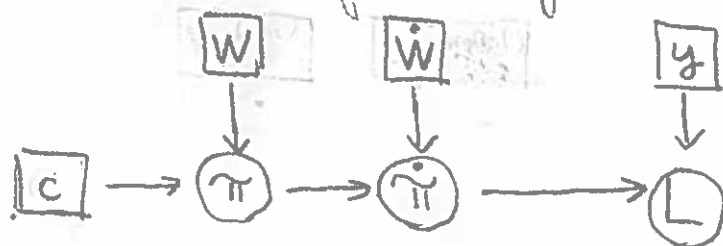
the striped zebra _____ from the zoo

A successful response could be "escaped", while other responses like "horse" or "delighted" should be considered less likely.

- ⑥ The input to the network is the sum of one-hot vectors for the words in the context, e.g.



- ⑦ The network itself is quite simple:



one-hot representation of the word that should fill in the blank

where: $\pi = W^T C$

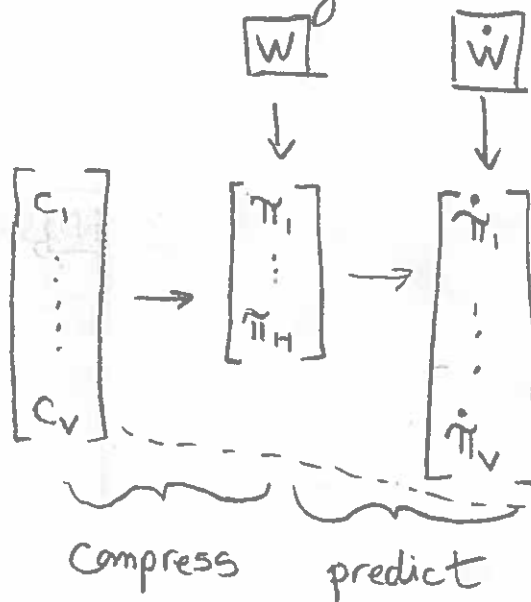
$$\hat{\pi} = \hat{W}^T \pi$$

$$L = -\log(y^T \text{softmax}(\hat{\pi}))$$

WORD VECTORS

- ⑧ Observe that the dimension of W is $V \times H$, π is $H \times 1$
 \dot{W} is $H \times V$, $\dot{\pi}$ is $V \times 1$.

Essentially this simple network needs to learn to compress large context vectors (of size V) into a compact vector (of size $H \ll V$), which can be used to predict the missing word.



this model is called the Continuous Bag-of-Words (CBOW) model, because the context vector doesn't represent the order of the context words, it just jumbles them all together in a "bag"



- ⑨ Interestingly, no activation function is ever applied.

This is a design decision to keep the model extremely simple so it can be trained on billions of word/context pairs.

WORD VECTORS

- ⑩ Billions? Why yes. Notice that it's easy to get training data for this task just by scraping documents from the web.

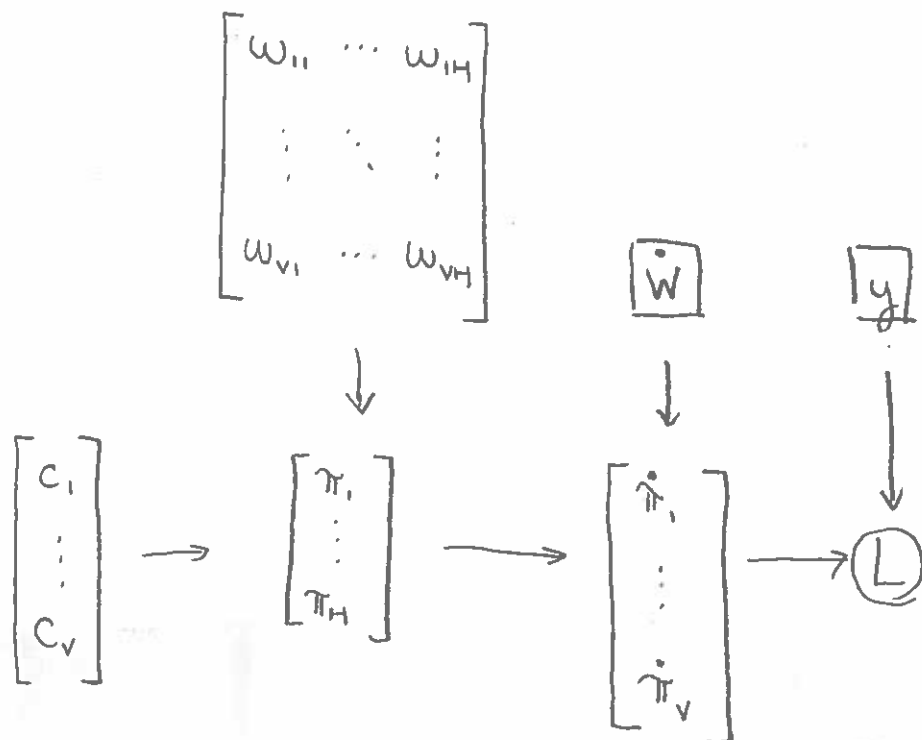
The unique stripes of zebras make them distinctive

- $b = \text{"the"}, c = \{\text{"of"}, \text{"stripes"}, \text{"unique"}\}$
- $b = \text{"unique"}, c = \{\text{"of"}, \text{"stripes"}, \text{"the"}, \text{"unique"}, \text{"zebras"}\}$
- $b = \text{"stripes"}, c = \{\text{"make"}, \text{"of"}, \text{"the"}, \text{"unique"}, \text{"zebras"}\}$
- $b = \text{"of"}, c = \{\text{"make"}, \text{"stripes"}, \text{"the"}, \text{"them"}, \text{"unique"}, \text{"zebras"}\}$
- etc.

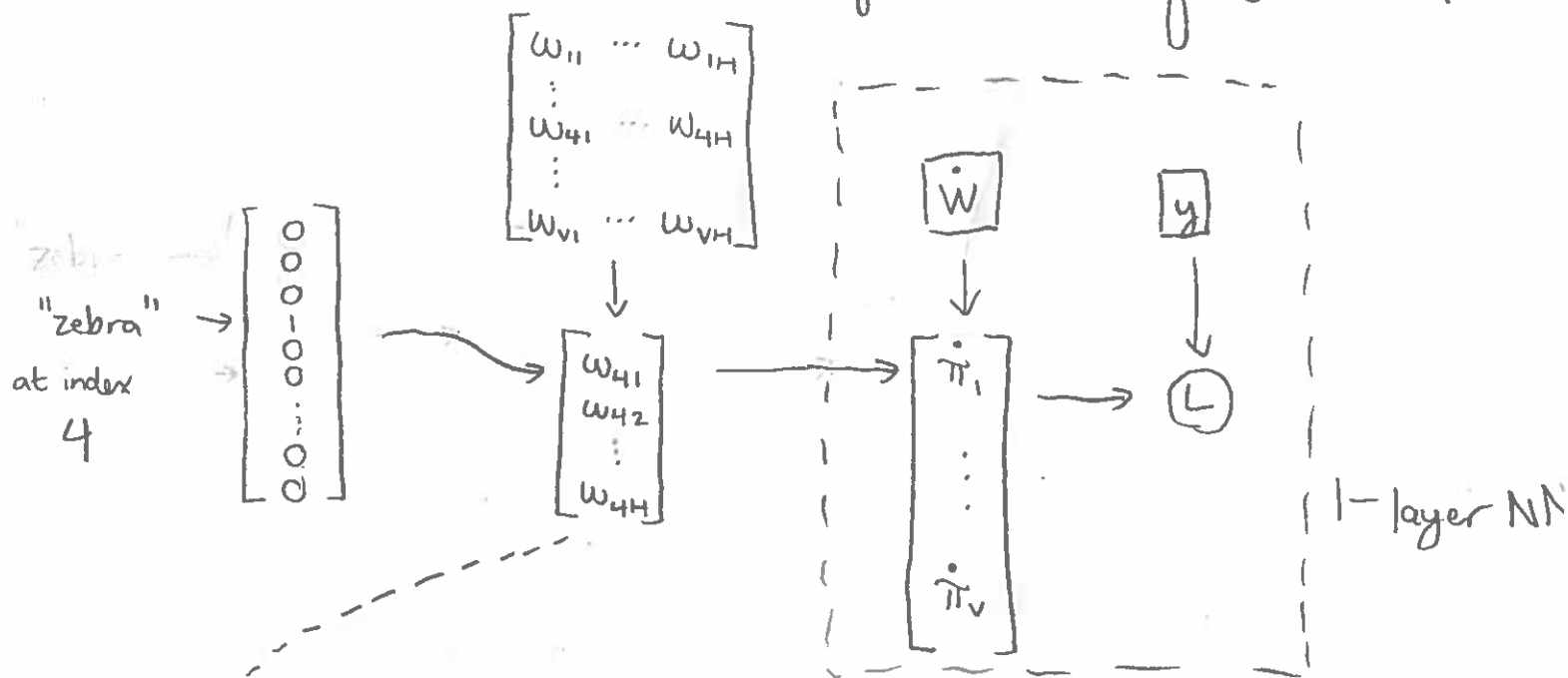
- ⑪ So we've trained a network. What then? I thought we wanted vector representations of words, not a classifier for fill-in-the-blank questions.

WORD VECTORS

⑫ Well, let's take a closer look at what the trained network looks like:



⑬ What if we provide a single word as an input context? Then c would be a one-hot representation of that word:



The vector π needs to capture H features of the context that are sufficient to predict nearby words with a 1-layer NN without activation (i.e. a logistic regression model).

WORD VECTORS

⑭ Observe that $\pi = \begin{bmatrix} w_{41} \\ w_{42} \\ \vdots \\ w_{4H} \end{bmatrix}$, which is just the 4th row

of matrix W . So effectively the 4th row is compressing word 4 of the vocabulary ("zebra") into a fixed size vector that is optimized to be able to predict nearby words.

And in general, the h^{th} row of W compresses the h^{th} word of the vocabulary into a set of predictive features.

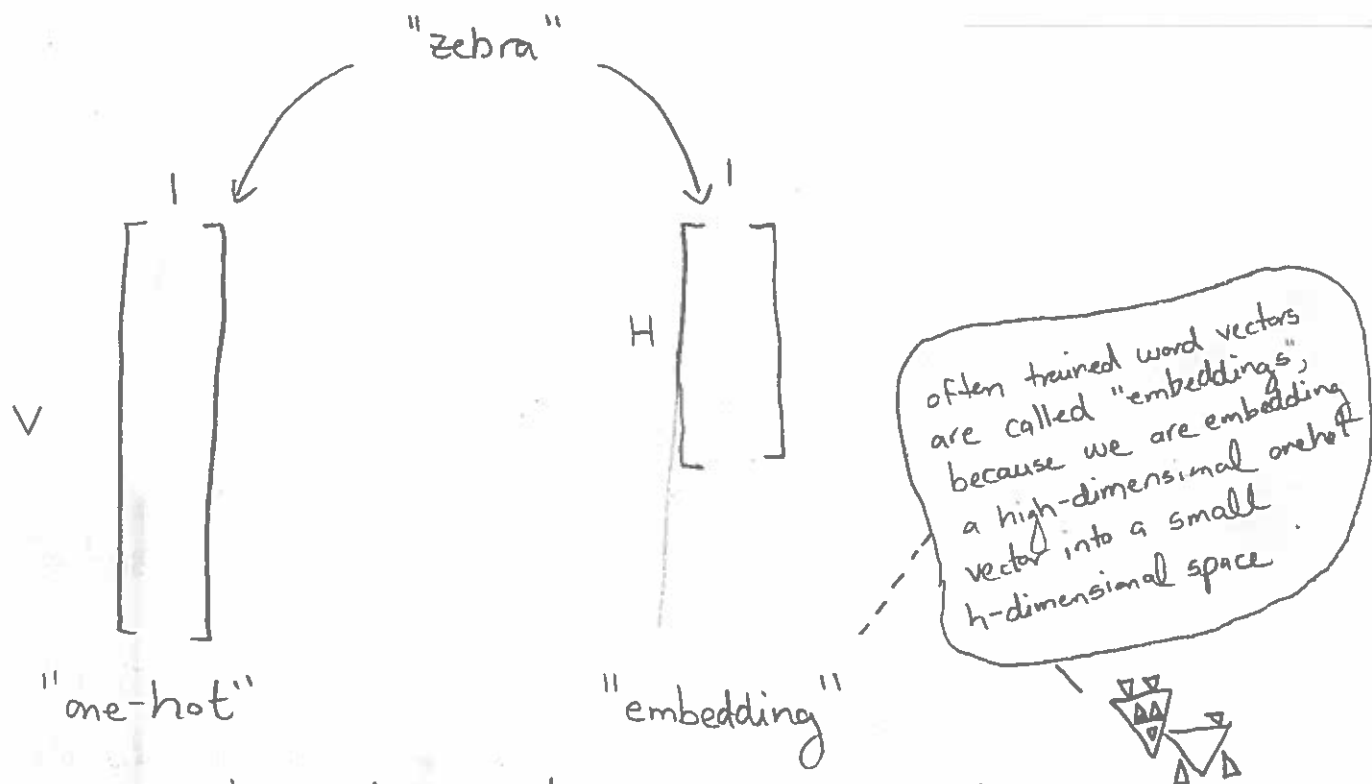
⑮ This can give us some very useful representations of words. Notice that dogs and cats are both often adopted, petted, fed, loved, and housebroken.

This gives the words "dog" and "cat" a distinct context from, say, "zebras", whom are rarely loved due to their anarchistic tendencies.

Therefore, the trained vectors ^{may} v encode some idea of "being a pet", in order to more accurately predict surrounding context.

WORD VECTORS

- ⑩ A practical upshot is that we can then use these word vector representations in more complex models instead of a one-hot representation.

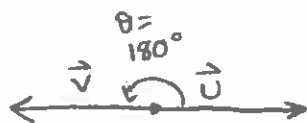
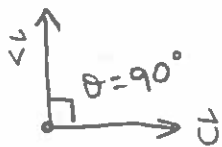
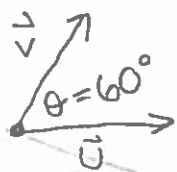


We can choose H to be whatever manageable size we like ($H=300$ is a common choice), whereas the vocab size V is out of our control.

WORD VECTORS

⑦ What do the trained word vectors end up looking like?

Well, there are challenges to visualizing them, since they are in something like a 300-dimensional space, but even in a high-dimensional space, two vectors are always separated by some angle θ .



If $\theta = 0^\circ$, then they're pointing in the same direction (and therefore very similar)

If $\theta = 90^\circ$, then they're orthogonal (and presumably also appear in unrelated contexts)

If $\theta = 180^\circ$, then they're pointing in opposing directions (and therefore may have opposing contexts)

⑧ Note that $\theta = 180^\circ$ doesn't suggest opposite meanings, but rather opposite contexts. "fast" and "slow" are opposites, but often appear in similar sentential contexts.

WORD VECTORS

- 19) Given two vectors \vec{u} and \vec{v} , how do we compute the angle between them? Well, the geometric formulation of dot product tells us:

$$\vec{u} \cdot \vec{v} = \|\vec{u}\| \|\vec{v}\| \cos \theta$$

where $\|\vec{x}\| = \sqrt{x_1^2 + \dots + x_n^2}$ is the length (L_2 -norm) of vector \vec{x} , and θ is the angle that separates \vec{u} and \vec{v} . Thus:

$$\cos \theta = \frac{\vec{u} \cdot \vec{v}}{\|\vec{u}\| \|\vec{v}\|}$$

-
- 20) This measure, $\frac{\vec{u} \cdot \vec{v}}{\|\vec{u}\| \|\vec{v}\|}$ is called the cosine similarity

of vectors \vec{u} and \vec{v} . Observe:

$$\text{when } \cos \theta = 1, \quad \text{then } \theta = 0^\circ$$

$$\cos \theta = -1, \quad \text{then } \theta = 180^\circ$$

$$\cos \theta = 0, \quad \text{then } \theta = 90^\circ$$

So a cosine similarity of 1 indicates maximal similarity, while a cosine similarity of -1 indicates maximal dissimilarity.

WORD VECTORS

- ② A commonly used set of word vectors was released by Google and trained on 100 billion words of Google News documents. These vectors have dimension 300. If we compute the cosine similarity of the word "king" with various other words, we get:

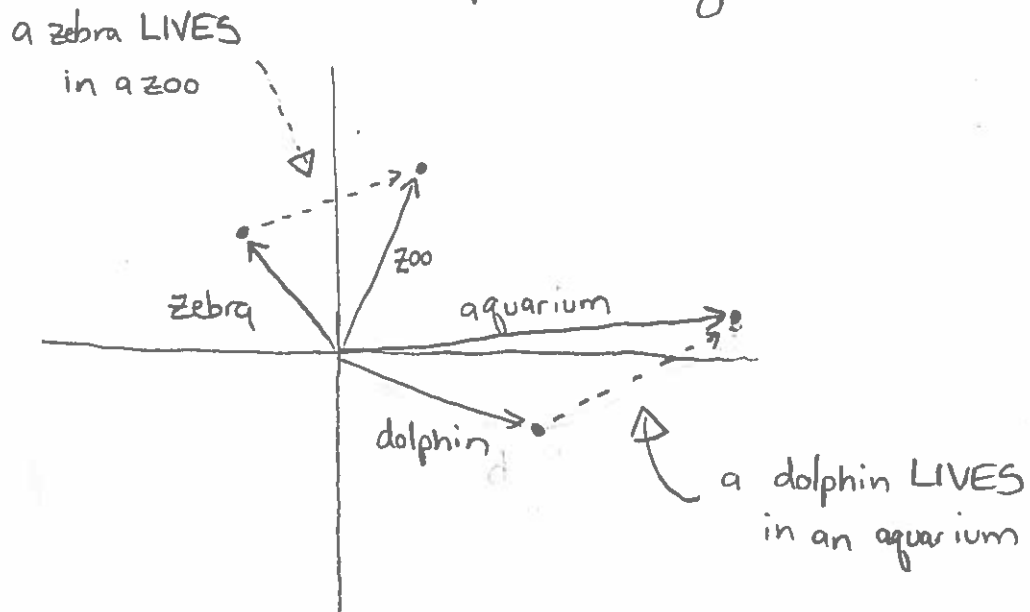
queen : 0.651

chess : 0.179

zebra : 0.098

farmed : 0.063

- ② But there are even deeper, more interesting structures that one can uncover with cosine similarity. One could hypothesize that perhaps word relationships are also captured in the vector space, e.g.

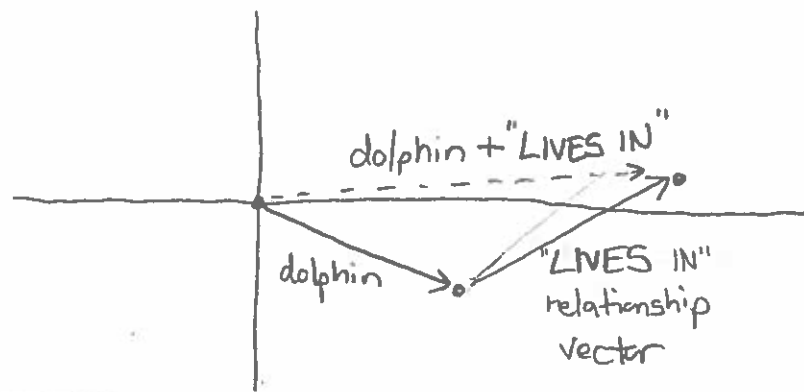


WORD VECTORS

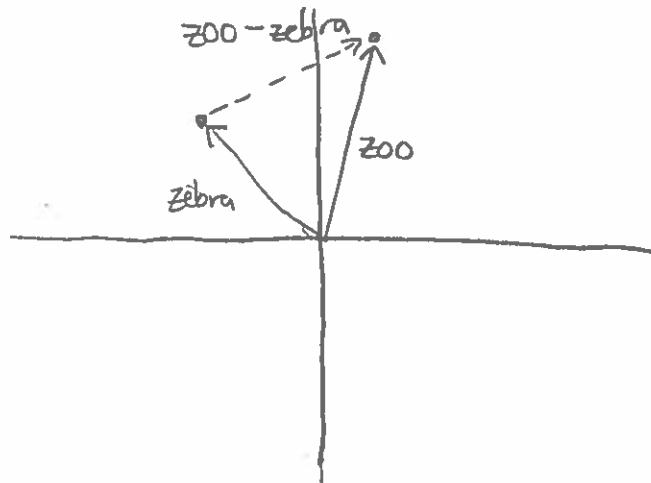
(23) If that were indeed the case, then we should be able to solve analogies, e.g.

zebra is to zoo AS dolphin is to _____

Namely, we look for the closest vector \vec{v} such that



(24) We can compute the relationship vector between zebra and zoo geometrically:



Thus, we look for the closest vector to:

$$\text{dolphin} + (\text{zoo} - \text{zebra})$$

WORD VECTORS

②⑤ Shockingly, this often works. Here are the closest five vectors using the Google News vectors:

zebra is to zoo AS dolphin is to

- 1. aquarium
- 2. Aquarium
- 3. dolphins
- 4. Vancouver Aquarium
- 5. Oceanarium

Seaworld comes in
at #7

