# Defining a new learning problem

I want a glass of orange juice to go along with my cereal.

| Context | word | target? |
|---------|------|---------|
| orange | juice | 1 |
| orange | king | 0 |
| orange | book | 0 |
| orange | the | 0 |
| orange | of | 0 |

$k = 5\text{-}20$    smaller. datasets

$k = 2\text{-}5$    larger dataset

[Mikolov et. al., 2013. Distributed representation of words and phrases and their compositionality]

Andrew Ng
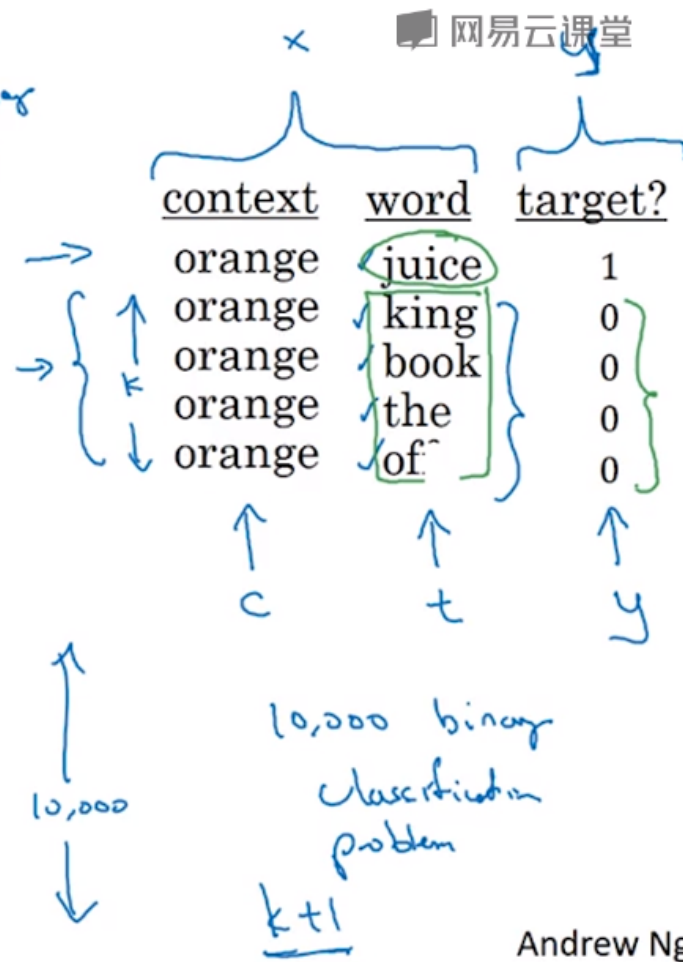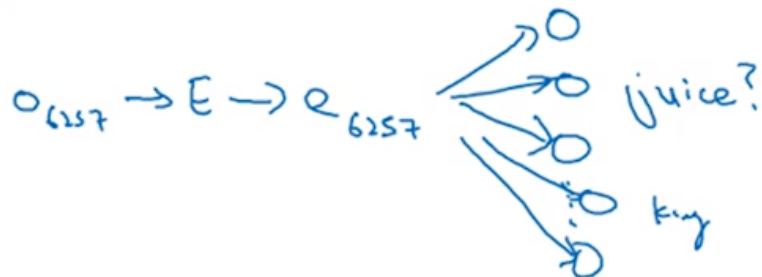
larger training dataset smaller k

It's really to try to distinguish between these two types of distributions from which you might sample a pair of words.

Next, let's describe the supervised learning model for learning a mapping from x to y.

# Model

Softmax: $\quad p(t|c) = \dfrac{e^{\theta_t^T e_c}}{\sum_{j=1}^{10,000} e^{\theta_j^T e_c}}$ $\bigg\}$ 10,000-way softmax

$$P(y=1 \mid c,t) = \sigma\left(\theta_t^T e_c\right) \leftarrow$$

Orange 6257

$O_{6257} \rightarrow E \rightarrow e_{6257}$ →→ juice?

→→ king

10,000

| context | word | target? |
|---------|------|---------|
| orange | juice | 1 |
| orange | king | 0 |
| orange | book | 0 |
| orange | the | 0 |
| orange | of | 0 |

c     t     y

10,000 binary classification problem

k+1

k to 1 ratio of negative to positive examples

And on every iteration, we're only going to train k+1 of them.

Andrew Ng

So the parameters are similar as before, you have one parameter vector theta for each possible target word.
And a separate parameter vector, really the embedding vector for each possible context word.
And we're going to use this formula to estimate the probability that y is equal to 1.

# Selecting negative examples

$\overbrace{\qquad\qquad}^{x} \quad \overbrace{\quad}^{y}$

| context | word | target? |
|---------|------|---------|
| orange | juice | 1 |
| orange | king | 0 |
| orange | book | 0 |
| orange | the | 0 |
| orange | of | 0 |

the, of, and, ...

$t$

$$P(w_i) = \frac{f(w_i)^{3/4}}{\sum_{j=1}^{10,000} f(w_j)^{3/4}} \qquad \frac{1}{|V|}$$

f empirical frequency of word in your corpus

Andrew Ng