

Deep Learning for Language Modeling

Hung-yi Lee

Language model

Speech Recognition



$$\begin{aligned}\tilde{W} &= \mathit{arg} \max_W P(W|X) \\ &= \mathit{arg} \max_W \frac{P(X|W)P(W)}{P(X)} \\ &= \mathit{arg} \max_W P(X|W)P(W)\end{aligned}$$

$P(X|W)$:

Acoustic Model

$P(W)$:

Language Model

Language model

- Language model: Estimated the probability of word sequence
 - Word sequence: $w_1, w_2, w_3, \dots, w_n$
 - $P(w_1, w_2, w_3, \dots, w_n)$
- Useful in speech recognition
 - Different word sequence can have the same pronunciation



recognize speech
or
wreck a nice beach

If $P(\text{recognize speech}) > P(\text{wreck a nice beach})$

Output =
“recognize speech”

Language model

$$\begin{aligned} &P(\text{"wreck a nice beach"}) \\ &= P(\text{wreck} | \text{START})P(a | \text{wreck}) \\ &P(\text{nice} | a)P(\text{beach} | \text{nice}) \end{aligned}$$

- How to estimate $P(w_1, w_2, w_3, \dots, w_n)$
- Collect a large amount of text data as training data
 - However, the word sequence w_1, w_2, \dots, w_n may not appear in the training data
- N-gram language model: $P(w_1, w_2, w_3, \dots, w_n) = P(w_1 | \text{START})P(w_2 | w_1) \dots P(w_n | w_{n-1})$
- Estimate $P(\text{beach} | \text{nice})$ from training data

$$P(\text{beach} | \text{nice}) = \frac{C(\text{nice beach})}{C(\text{nice})}$$

Count of "nice beach" in the training data

Count of "nice" in the training data

Language model - Smoothing

- Training data:
 - The dog ran
 - The cat jumped

This is called **language model smoothing**.

$$P(\text{jumped} \mid \text{dog}) = \cancel{0} \quad 0.0001$$

$$P(\text{ran} \mid \text{cat}) = \cancel{0} \quad 0.0001$$

Give some small probability

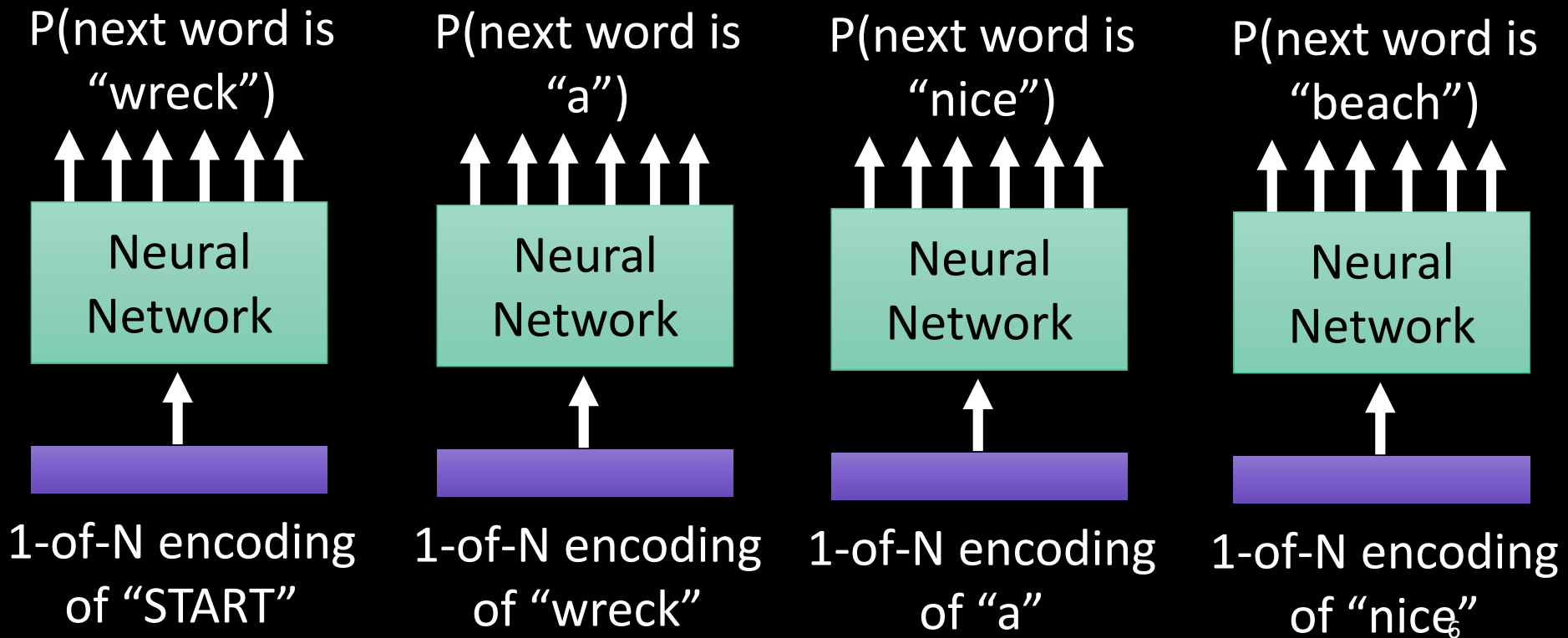
- The probability is not accurate.
- The phenomenon happens because we cannot collect all the possible text in the world as training data.

Neural-network based LM

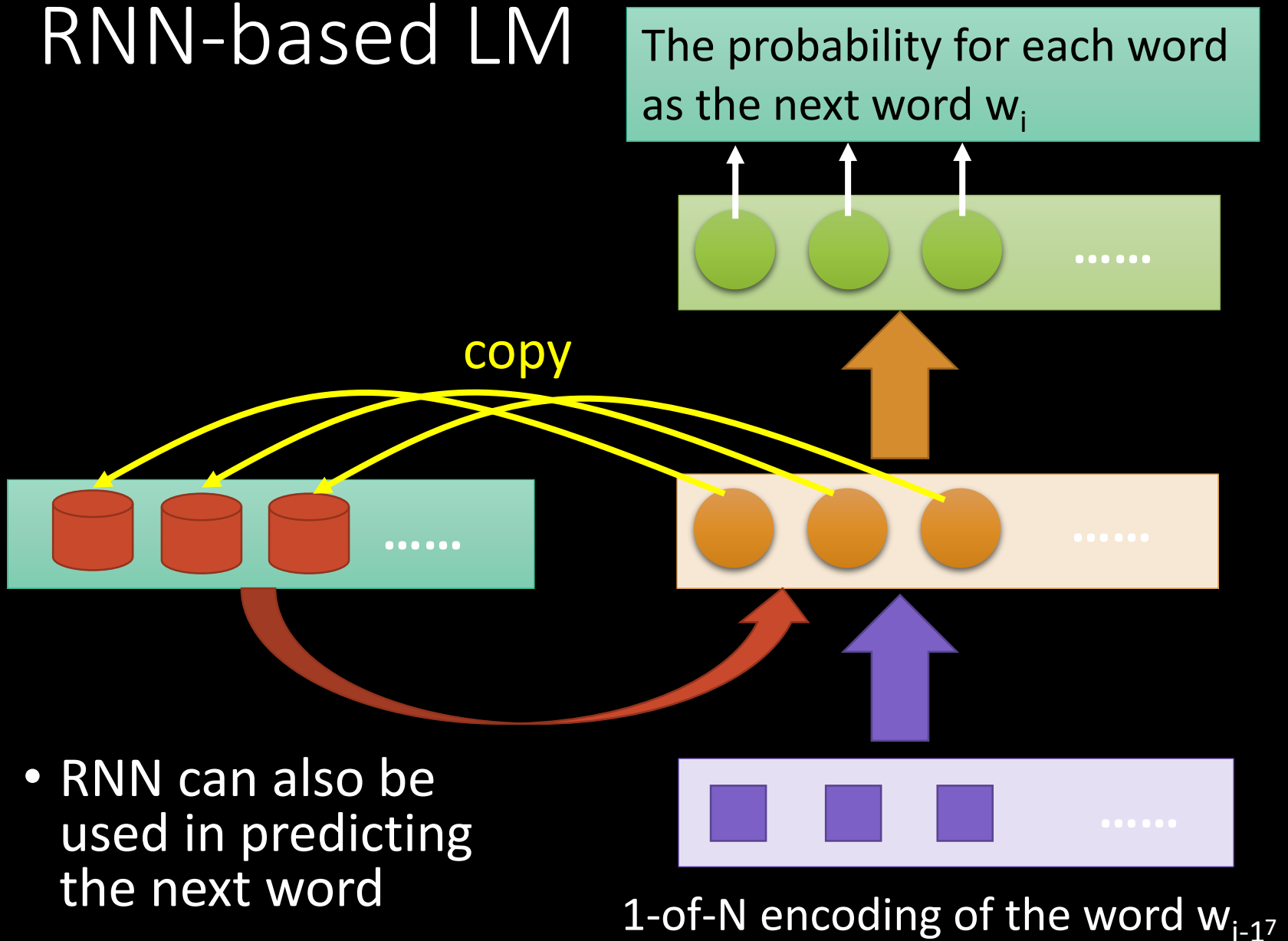
$P(\text{"wreck a nice beach"})$

$=P(\text{wreck} | \text{START})P(\text{a} | \text{wreck})P(\text{nice} | \text{a})P(\text{beach} | \text{nice})$

$P(b | a)$: not from counting, using NN to predict the next word.

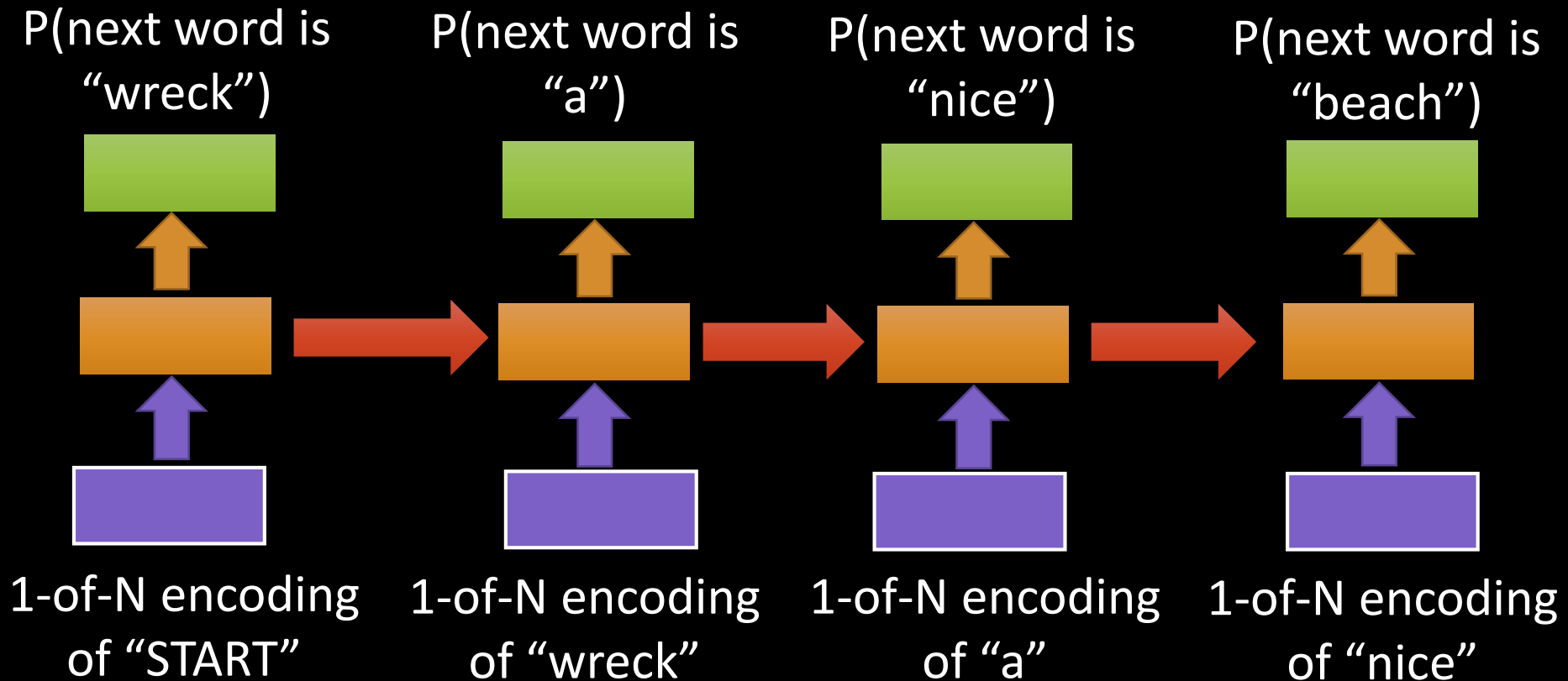


RNN-based LM



- RNN can also be used in predicting the next word

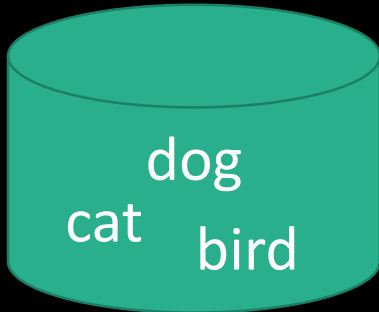
RNN-based LM



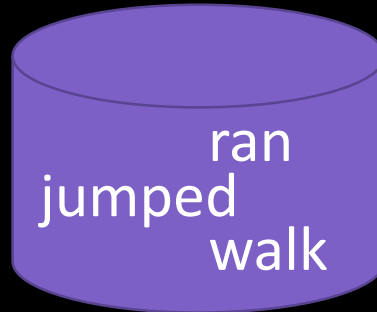
- Model long-term information
- People also used Deep RNN or LSTM

Why RNN works? – Class-based

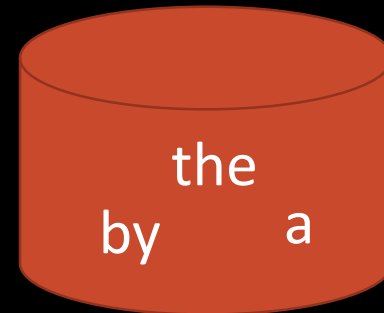
class 1: **A**nimal



class 2: **V**erb



class 3: **F**unction word



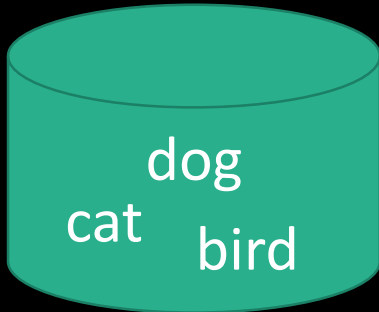
$W = "w_1 w_2 w_3"$ $C(w_i)$: class of word w_i

~~$$P(W) = P(w_1 | \text{START}) P(w_2 | w_1) P(w_3 | w_2)$$~~

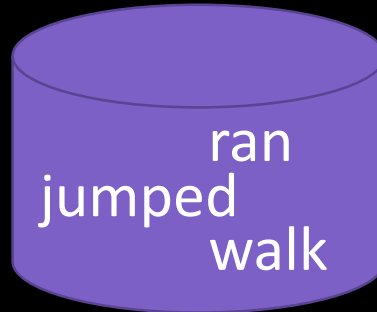
$$P(W) = P(C(w_1) | \text{START}) P(C(w_2) | C(w_1)) P(C(w_3) | C(w_2)) \\ \times P(w_1 | C(w_1)) P(w_2 | C(w_2)) P(w_3 | C(w_3))$$

Why RNN works? – Class-based

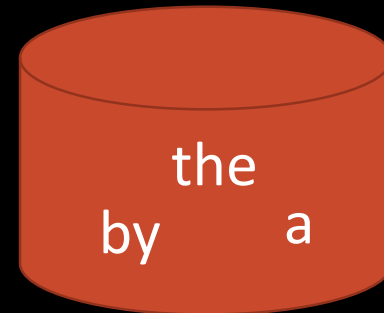
class 1: **A**nimal



class 2: **V**erb



class 3: **F**unction word



W = “the dog ran”
 F **A** **V**

$$P(W) = P(\mathbf{F}|\text{START}) P(\mathbf{A}|\mathbf{F}) P(\mathbf{V}|\mathbf{A})$$

$$\times P(\text{the}|\mathbf{F}) P(\text{dog}|\mathbf{A}) P(\text{ran}|\mathbf{V})$$

$P(\text{class } i | \text{class } j)$ and $P(\text{word } w | \text{class } i)$ are estimated from training data.

Why RNN works? – Class-based

$P(\text{class } i \mid \text{class } j)$ and $P(\text{word } w \mid \text{class } i)$ are estimated from training data.

Training data

the	dog	ran
<i>F</i>	<i>A</i>	<i>V</i>

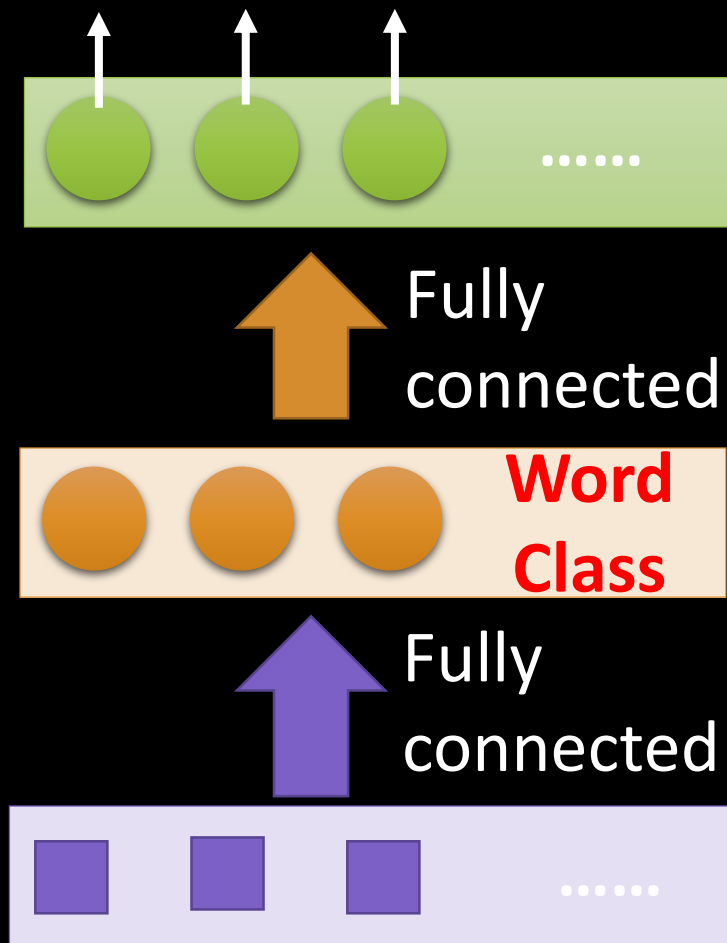
the	cat	jumped
<i>F</i>	<i>A</i>	<i>V</i>

$W =$ “the cat ran”
F *A* *V*

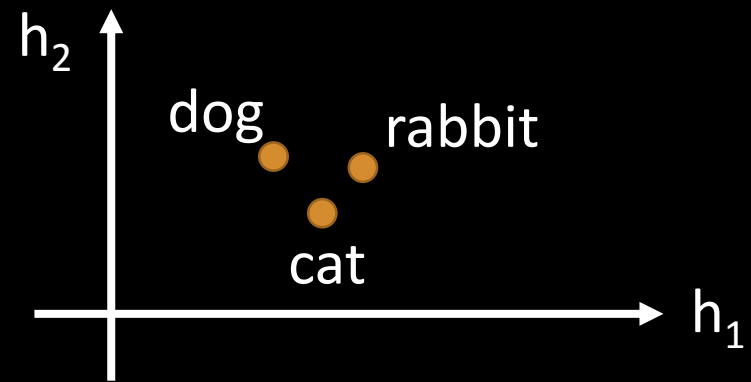
$P(\text{ran} \mid \text{cat})$ is zero given the training data

However, $P(\text{Verb} \mid \text{Animal})$ is not zero

Why RNN works? – Class-based



The hidden layer of the related words are close.

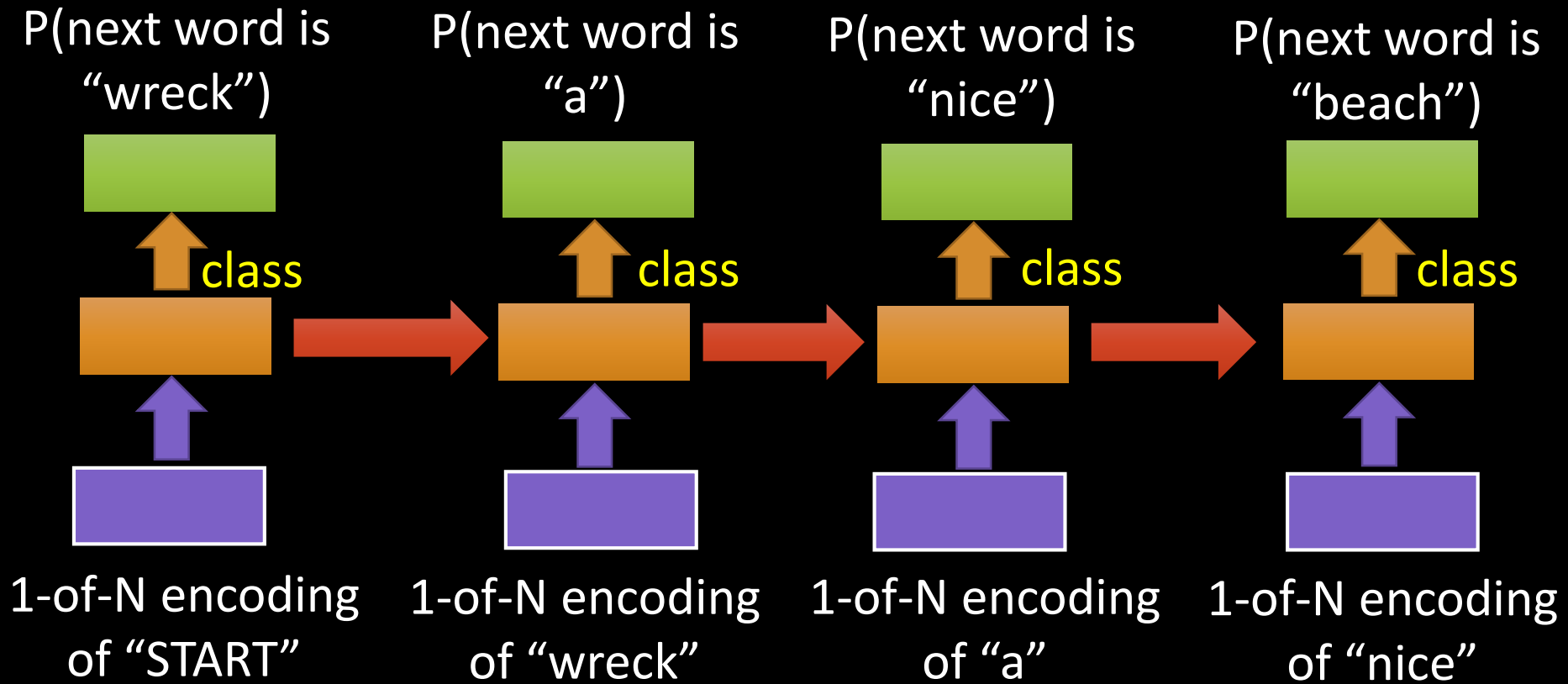


If $P(\text{jump} | \text{dog})$ is large, then $P(\text{jump} | \text{cat})$ increase accordingly. (even there is not "... cat jump ..." in the data)

Smoothing is automatically done.

Why RNN works? – Class-based

- Like class-based LM



Large Output Layer

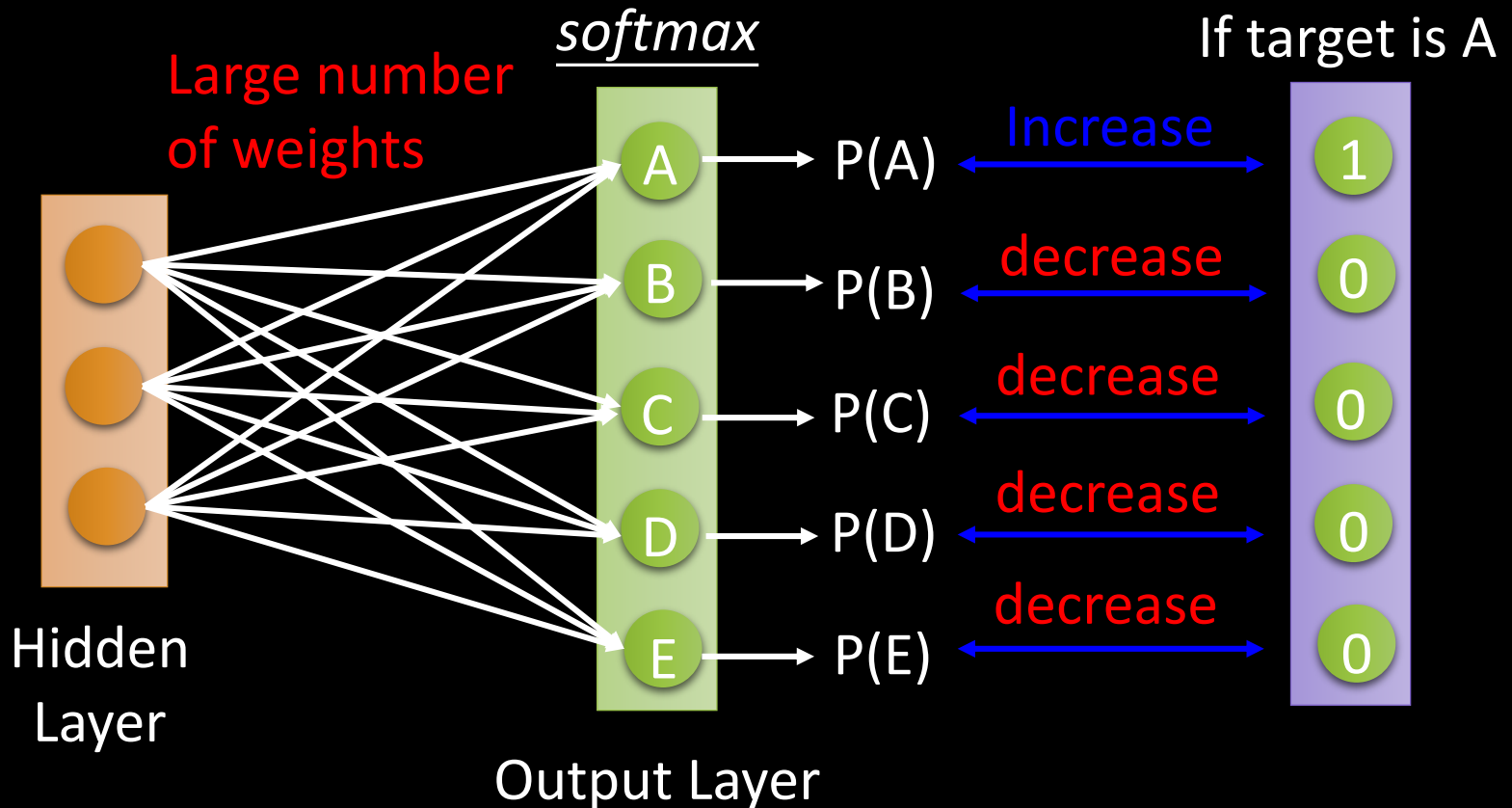
- Factorization of the Output Layer
 - Ref: Mikolov Tomáš: Statistical Language Models based on Neural Networks. PhD thesis, Brno University of Technology, 2012. (chapter 3.4.2)
- Noise Contrastive Estimation (NCE)
 - Ref: Ref: X. Chen, X. Liu, M. J. F. Gales and P. C. Woodland, "Recurrent neural network language model training with noise contrastive estimation for speech recognition," ICASSP, 2015
- More ways to deal with the large output layer
 - Hinton's course on Coursera: Ways to Deal with the Large Number of Possible Outputs
 - <https://www.youtube.com/watch?v=vLmgSo9LVMk>

Appendix

Some ways to deal with the
large output layer

Large Output Layer

Normal Training

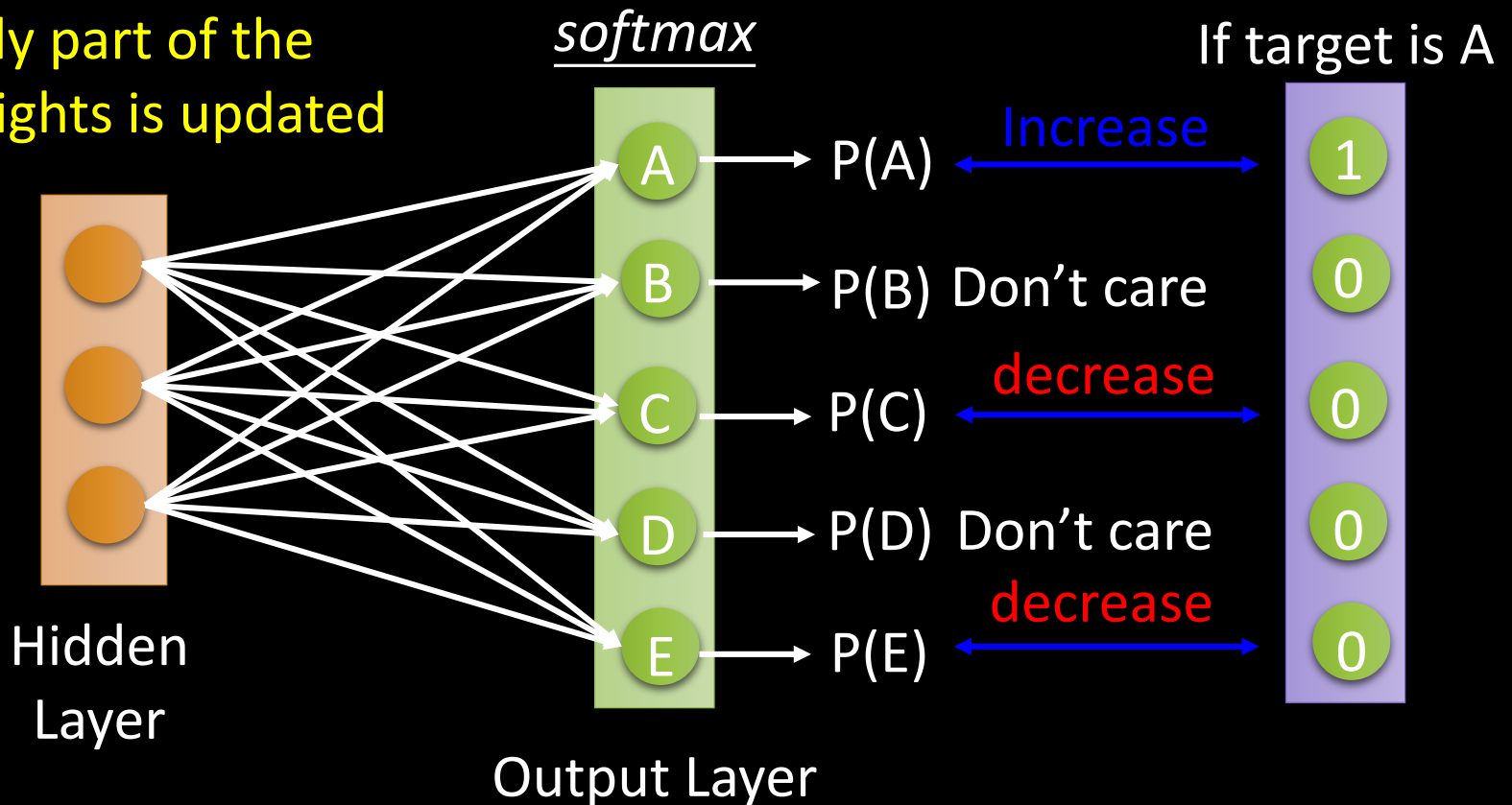


Large Output Layer – Noise Contrastive Estimation (NCE)

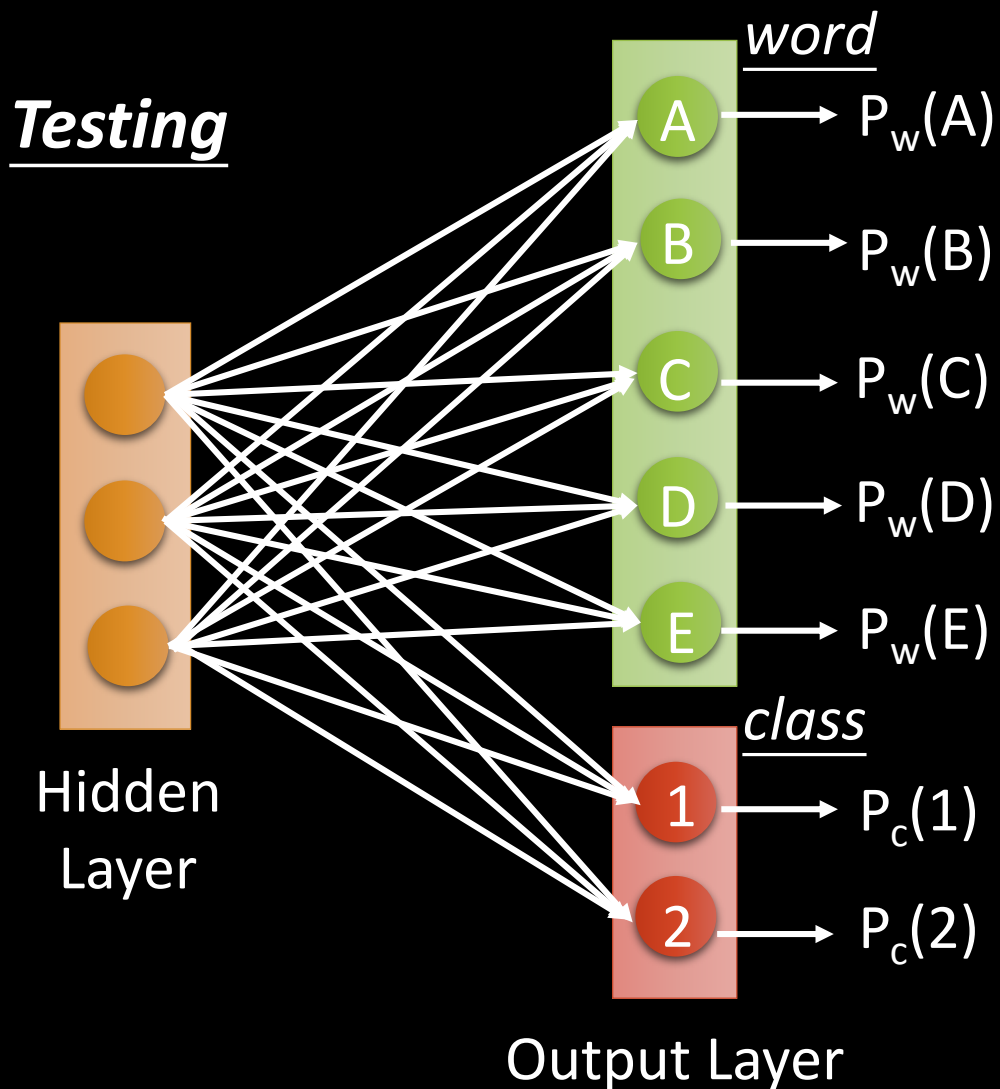
NCE

Randomly sample some words
to suppress the probability

Only part of the
weights is updated



Large Output Layer - Factorization



Word	Class
A	1
B	1
C	1
D	2
E	2

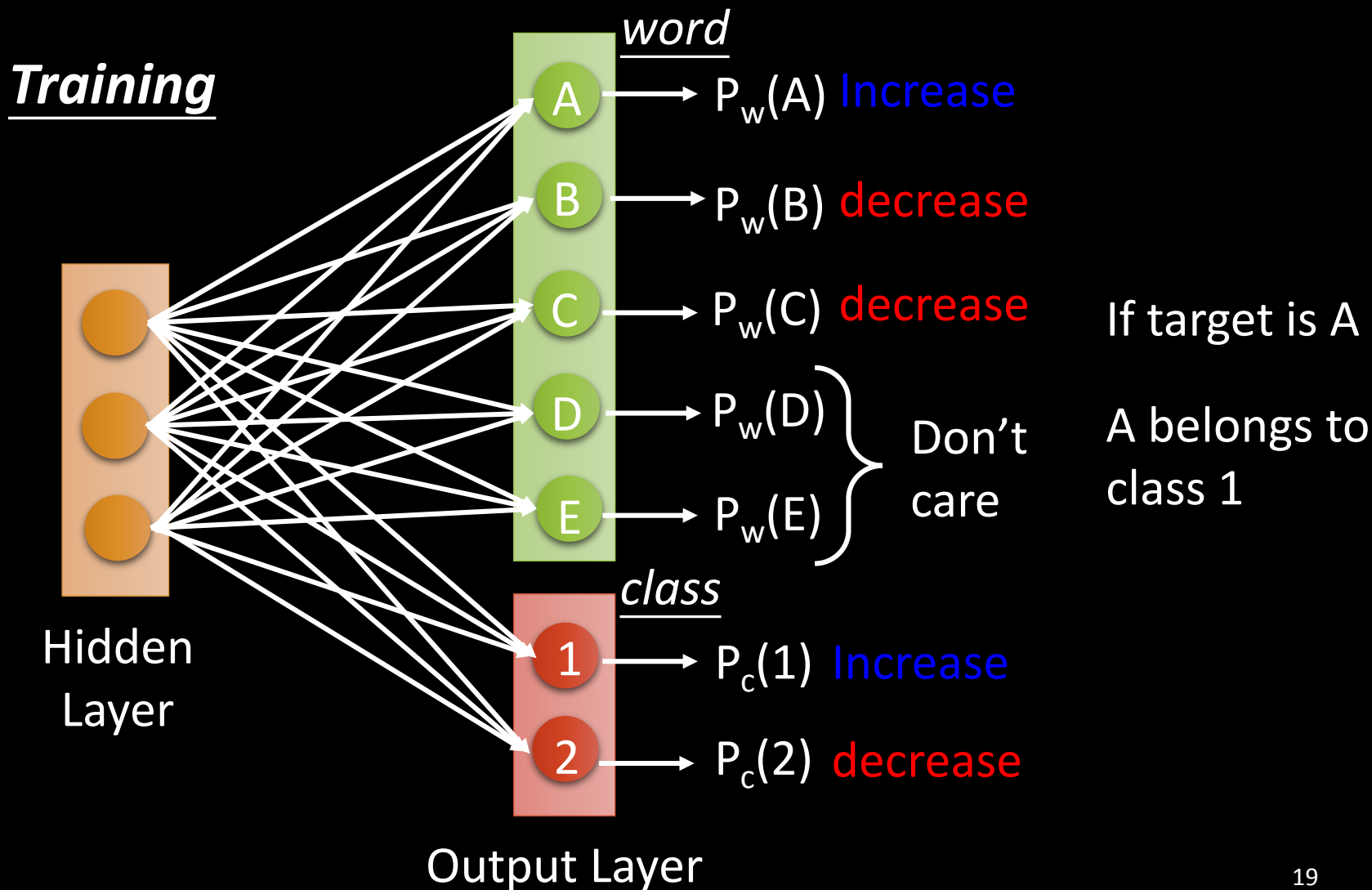
$$P(A) = P_w(A) \times P_c(1)$$

.....

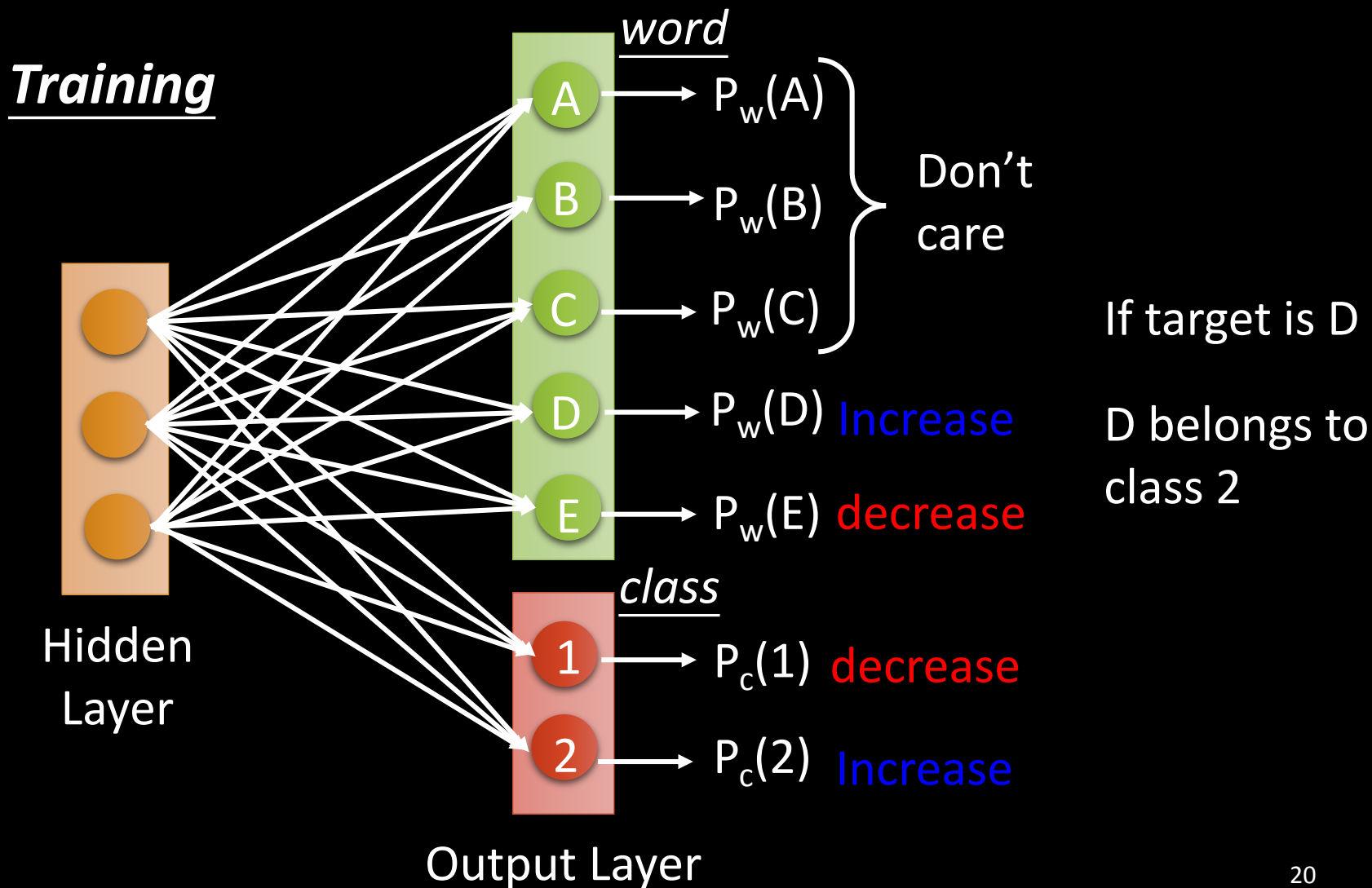
$$P(D) = P_w(D) \times P_c(2)$$

.....

Large Output Layer - Factorization



Large Output Layer - Factorization



Sentence Completion Task

- <https://docs.google.com/presentation/d/1K3nMJSczbdNTvp3GslgX59-drtVYJ558Lya2DXOeqeU/edit>