

Deep Learning 101

CSCI 1460: Computational Linguistics

Lecture 8

Ellie Pavlick

Fall 2023

Announcements

- Assignment 2 fixes (tf-idf vs. BOW)
- NLP talks in the department!
 - **Oct 13th, 12pm Panel/Discussion** Strong vs. Weak Compositionality in Humans and Machines!

Topics

- More Followup on Word Embeddings from SVD
- Logistic Regression and Gradient Descent
- Multilayer Perceptrons
- Word Embeddings from NNs

Topics

- **More Followup on Word Embeddings from SVD**
- Logistic Regression and Gradient Descent
- Multilayer Perceptrons
- Word Embeddings from NNs

SVD Revisited

The below figure shows the following: (Part of) a term-document matrix M; A V matrix that results when running LSA on M; An embedding (i.e., row of the U matrix) associated with a document d. Which of the below represents the most likely content of document d? (Note that document d is not supposed to be one of the docs doc1, doc2, doc3, doc4, doc5 along the rows of M. You can assume d is a different document that also occurred in M but is not shown in the below figure.) *

	red	green	apple	kiwi
doc1	1	1	1	1
doc2	1	0	1	0
doc3	1	1	0	1
doc4	1	0	1	1
doc5	0	0	1	1

M

10	1	0
8	0	8
1	6	1
-1	7	11

V

d	1	2	12
---	---	---	----

SVD Revisited

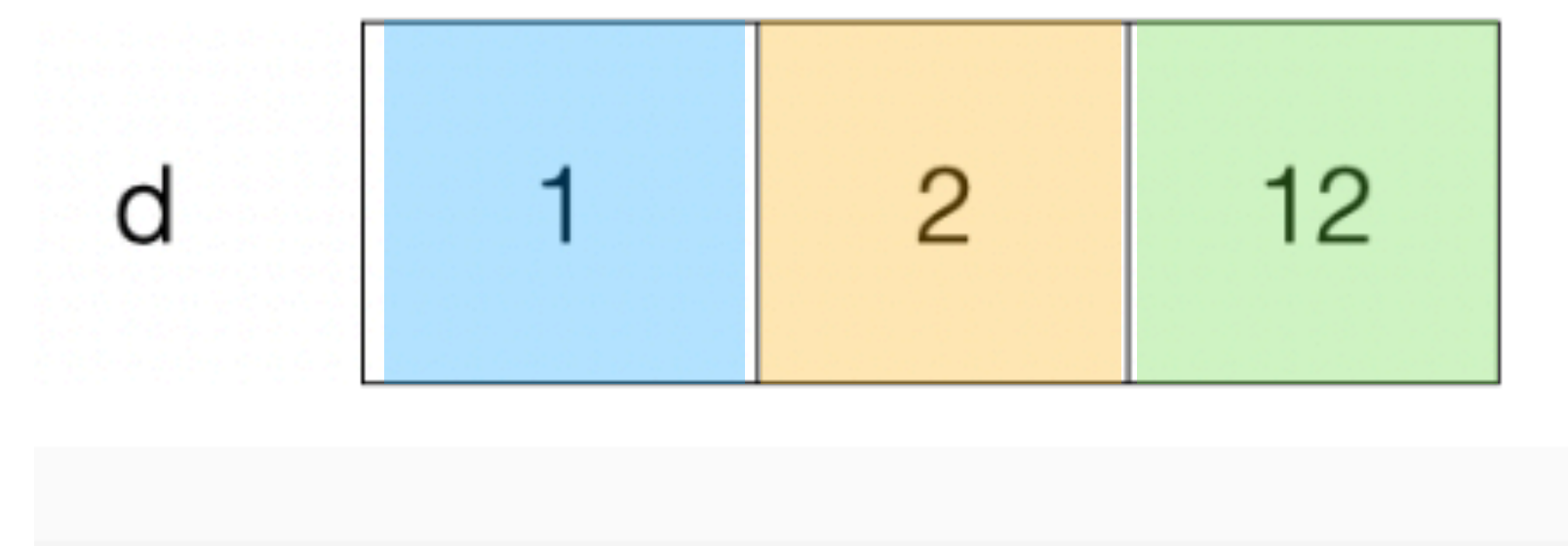
The below figure shows the following: (Part of) a term-document matrix M; A V matrix that results when running LSA on M; An embedding (i.e., row of the U matrix) associated with a document d. Which of the below represents the most likely content of document d? (Note that document d is not supposed to be one of the docs doc1, doc2, doc3, doc4, doc5 along the rows of M. You can assume d is a different document that also occurred in M but is not shown in the below figure.) *

	red	green	apple	kiwi
doc1	1	1	1	1
doc2	1	0	1	0
doc3	1	1	0	1
doc4	1	0	1	1
doc5	0	0	1	1

M

10	1	0
8	0	8
1	6	1
-1	7	11

V



SVD Revisited

The below figure shows the following: (Part of) a term-document matrix M; A V matrix that results when running LSA on M; An embedding (i.e., row of the U matrix) associated with a document d. Which of the below represents the most likely content of document d? (Note that document d is not supposed to be one of the docs doc1, doc2, doc3, doc4, doc5 along the rows of M. You can assume d is a different document that also occurred in M but is not shown in the below figure.) *

	red	green	apple	kiwi
doc1	1	1	1	1
doc2	1	0	1	0
doc3	1	1	0	1
doc4	1	0	1	1
doc5	0	0	1	1

M

10	1	0
8	0	8
1	6	1
-1	7	11

V

d	1	2	12
---	---	---	----

SVD Revisited

The below figure shows the following: (Part of) a term-document matrix M; A V matrix that results when running LSA on M; An embedding (i.e., row of the U matrix) associated with a document d. Which of the below represents the most likely content of document d? (Note that document d is not supposed to be one of the docs doc1, doc2, doc3, doc4, doc5 along the rows of M. You can assume d is a different document that also occurred in M but is not shown in the below figure.) *

	red	green	apple	kiwi
doc1	1	1	1	1
doc2	1	0	1	0
doc3	1	1	0	1
doc4	1	0	1	1
doc5	0	0	1	1

M

10	1	0
8	0	8
1	6	1
-1	7	11

V

d	1	2	12
---	---	---	----

SVD Revisited

The below figure shows the following: (Part of) a term-document matrix M; A V matrix that results when running LSA on M; An embedding (i.e., row of the U matrix) associated with a document d. Which of the below represents the most likely content of document d? (Note that document d is not supposed to be one of the docs doc1, doc2, doc3, doc4, doc5 along the rows of M. You can assume d is a different document that also occurred in M but is not shown in the below figure.) *

	red	green	apple	kiwi
doc1	1	1	1	1
doc2	1	0	1	0
doc3	1	1	0	1
doc4	1	0	1	1
doc5	0	0	1	1

M

10	1	0
8	0	8
1	6	1
-1	7	11

V

d	1	2	12
---	---	---	----

SVD Revisited

The below figure shows the following: (Part of) a term-document matrix M; A V matrix that results when running LSA on M; An embedding (i.e., row of the U matrix) associated with a document d. Which of the below represents the most likely content of document d? (Note that document d is not supposed to be one of the docs doc1, doc2, doc3, doc4, doc5 along the rows of M. You can assume d is a different document that also occurred in M but is not shown in the below figure.) *

	red	green	apple	kiwi
doc1	1	1	1	1
doc2	1	0	1	0
doc3	1	1	0	1
doc4	1	0	1	1
doc5	0	0	1	1

M

10	1	0
8	0	8
1	6	1
-1	7	11

V

d	1	2	12
---	---	---	----

SVD Revisited

The below figure shows the following: (Part of) a term-document matrix M; A V matrix that results when running LSA on M; An embedding (i.e., row of the U matrix) associated with a document d. Which of the below represents the most likely content of document d? (Note that document d is not supposed to be one of the docs doc1, doc2, doc3, doc4, doc5 along the rows of M. You can assume d is a different document that also occurred in M but is not shown in the below figure.) *

	red	green	apple	kiwi
doc1	1	1	1	1
doc2	1	0	1	0
doc3	1	1	0	1
doc4	1	0	1	1
doc5	0	0	1	1

M

10	1	0
8	0	8
1	6	1
-1	7	11

V

d	1	2	12
---	---	---	----

SVD Revisited

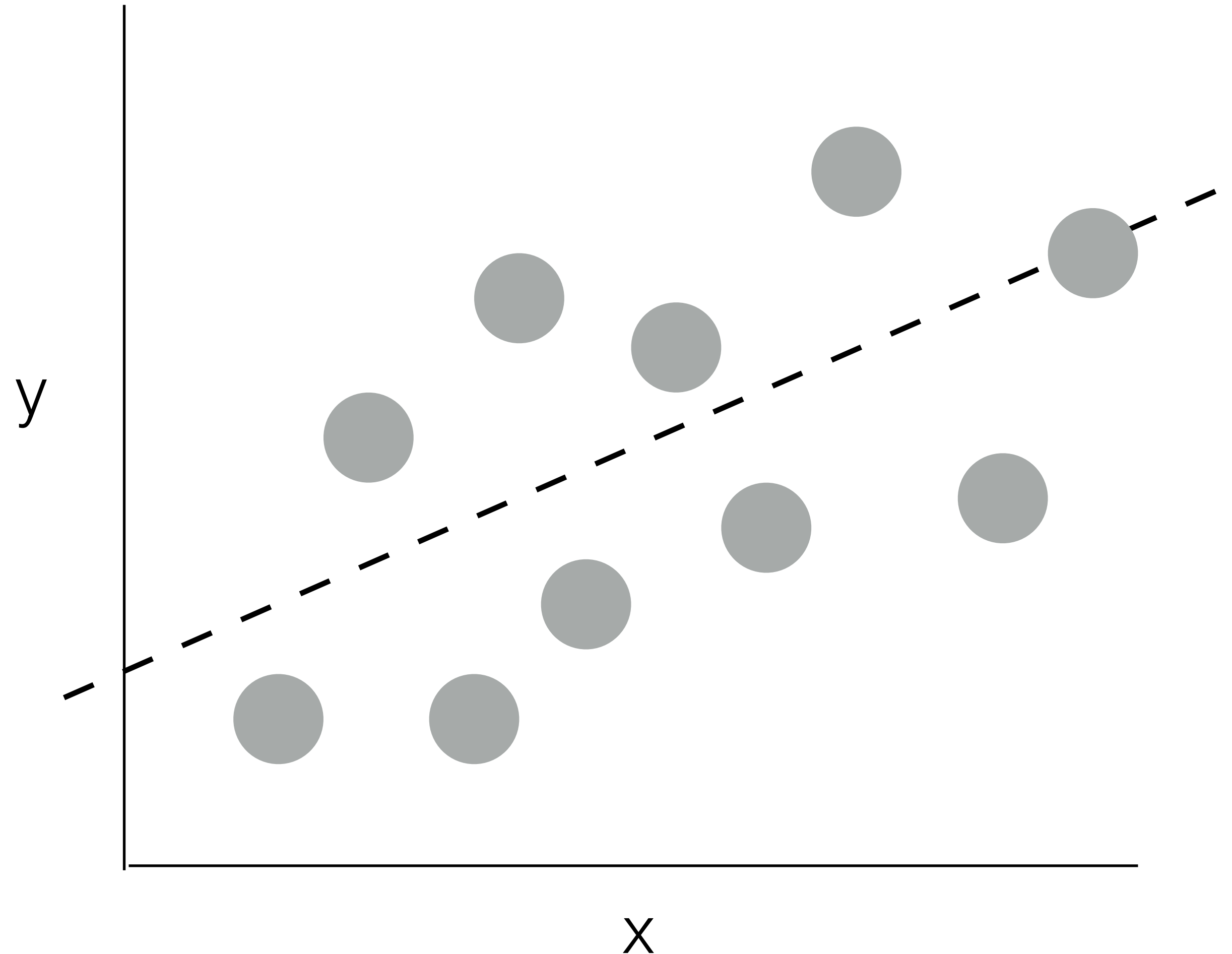
Colab Notebook

Topics

- More Followup on Word Embeddings from SVD
- **Logistic Regression and Gradient Descent**
- Multilayer Perceptrons
- Word Embeddings from NNs

Logistic Regression Classifiers

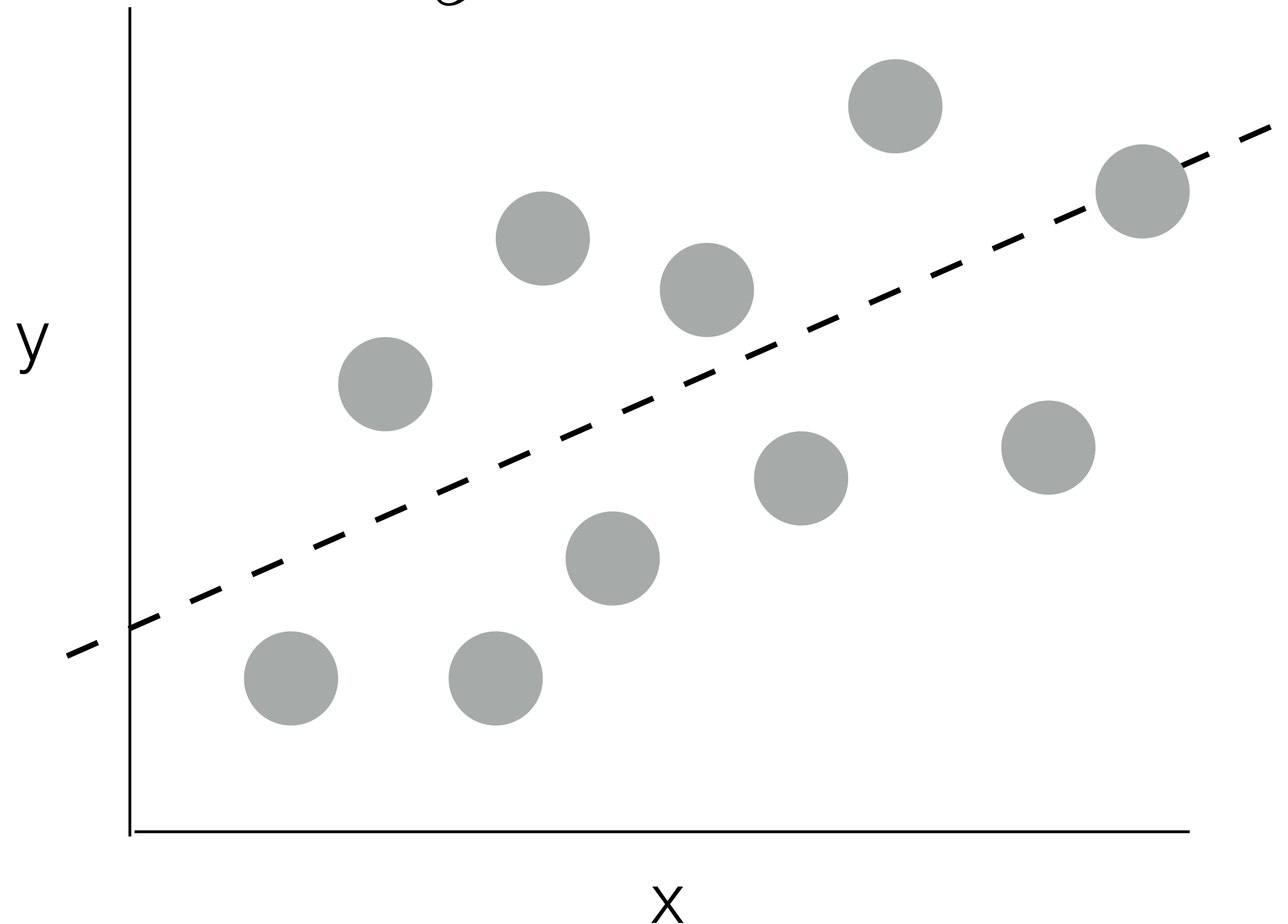
Making Predictions



Logistic Regression Classifiers

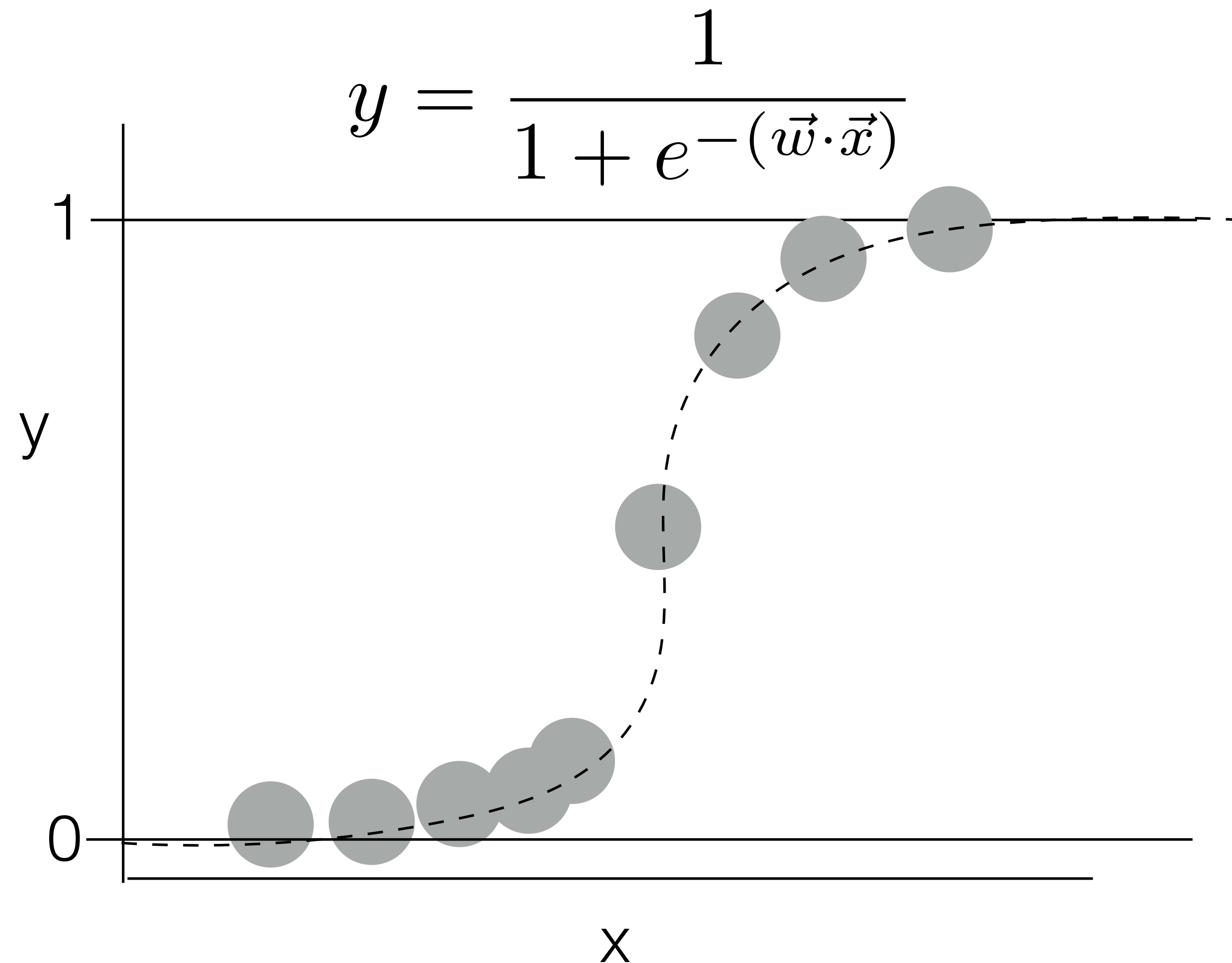
Making Predictions

$$y = \vec{w} \cdot \vec{x}$$



Logistic Regression Classifiers

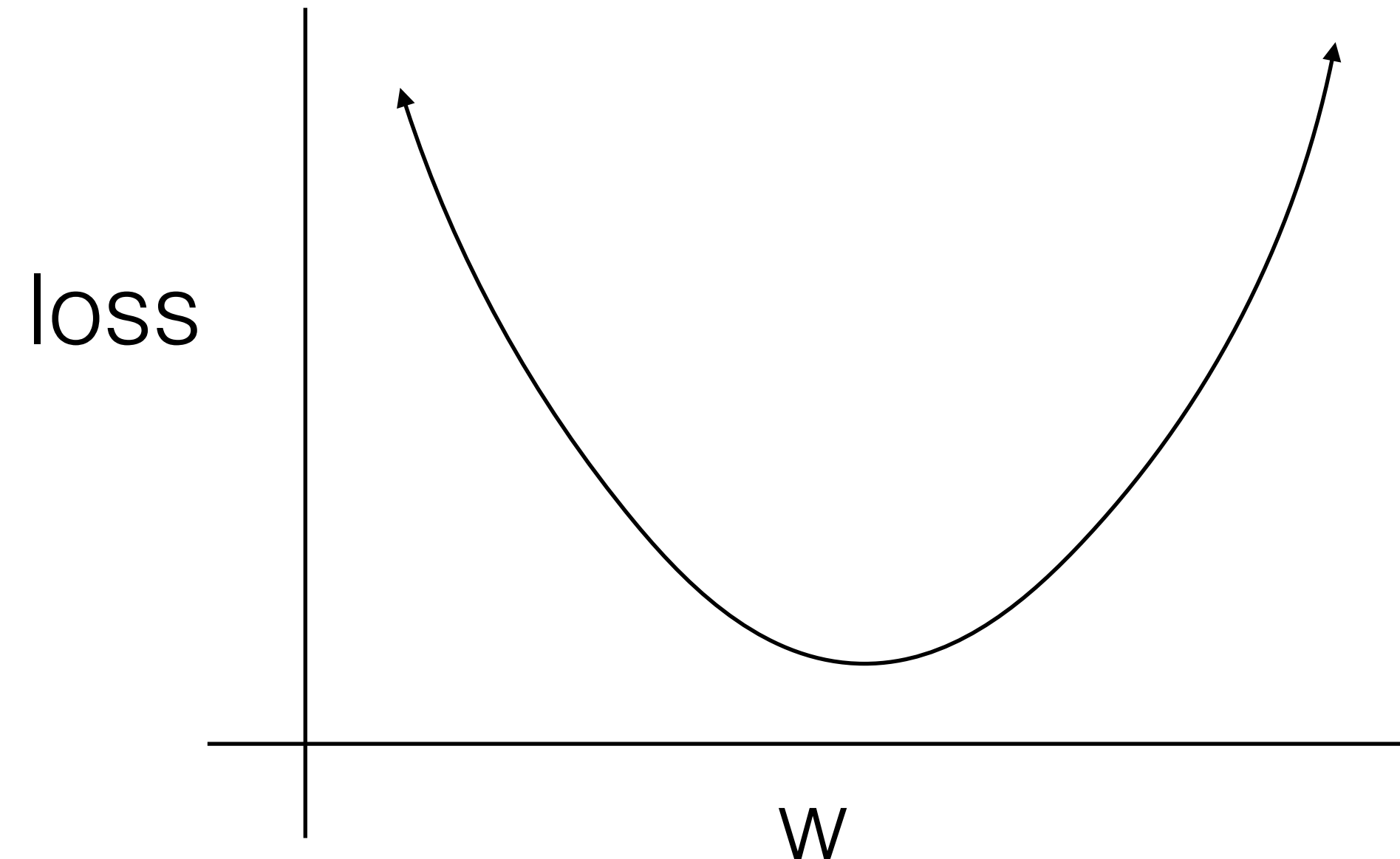
Making Predictions



Logistic Regression Classifiers

Training with Gradient Descent

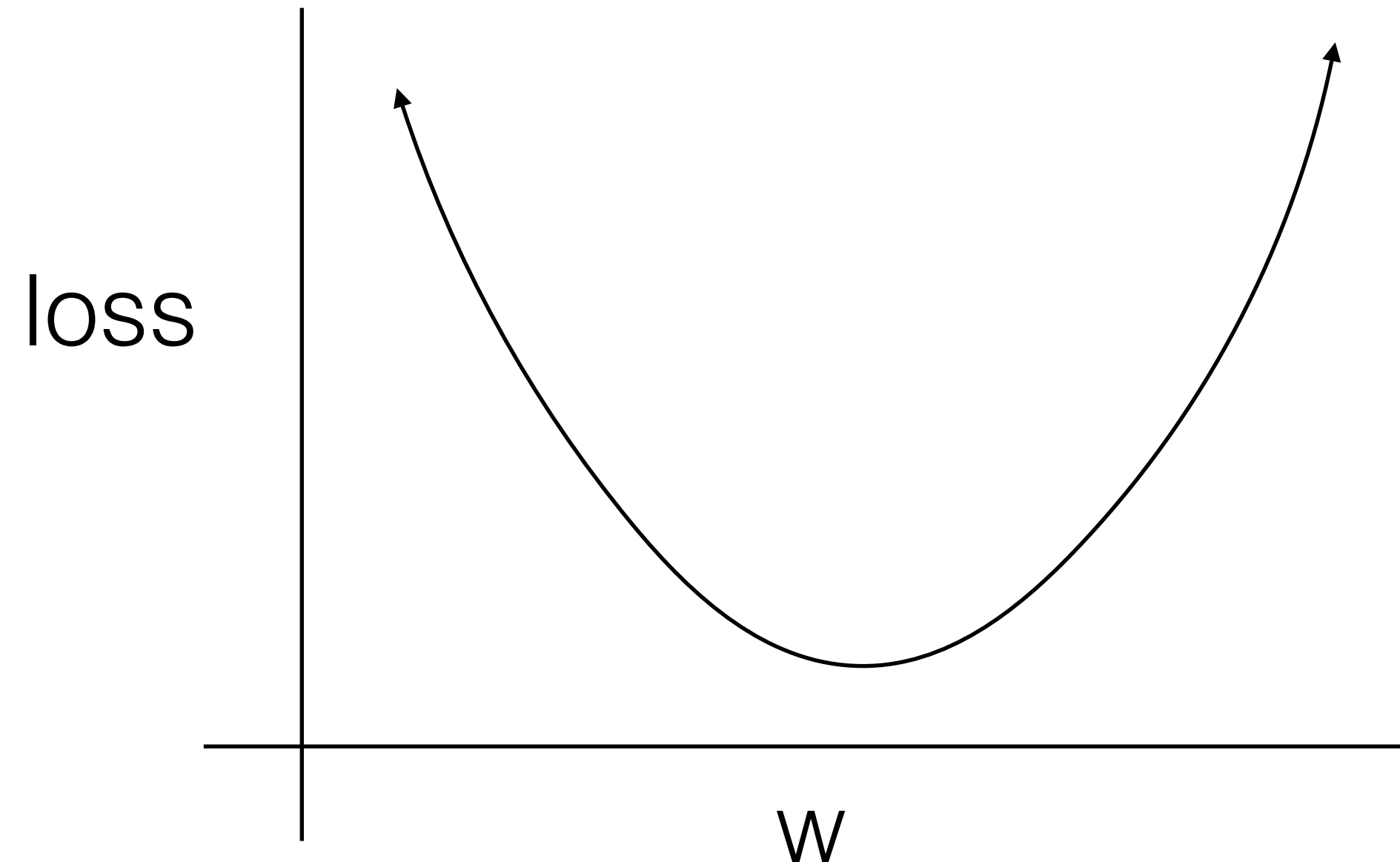
minimize $\text{loss}(\text{data}, w)$



Logistic Regression Classifiers

Training with Gradient Descent

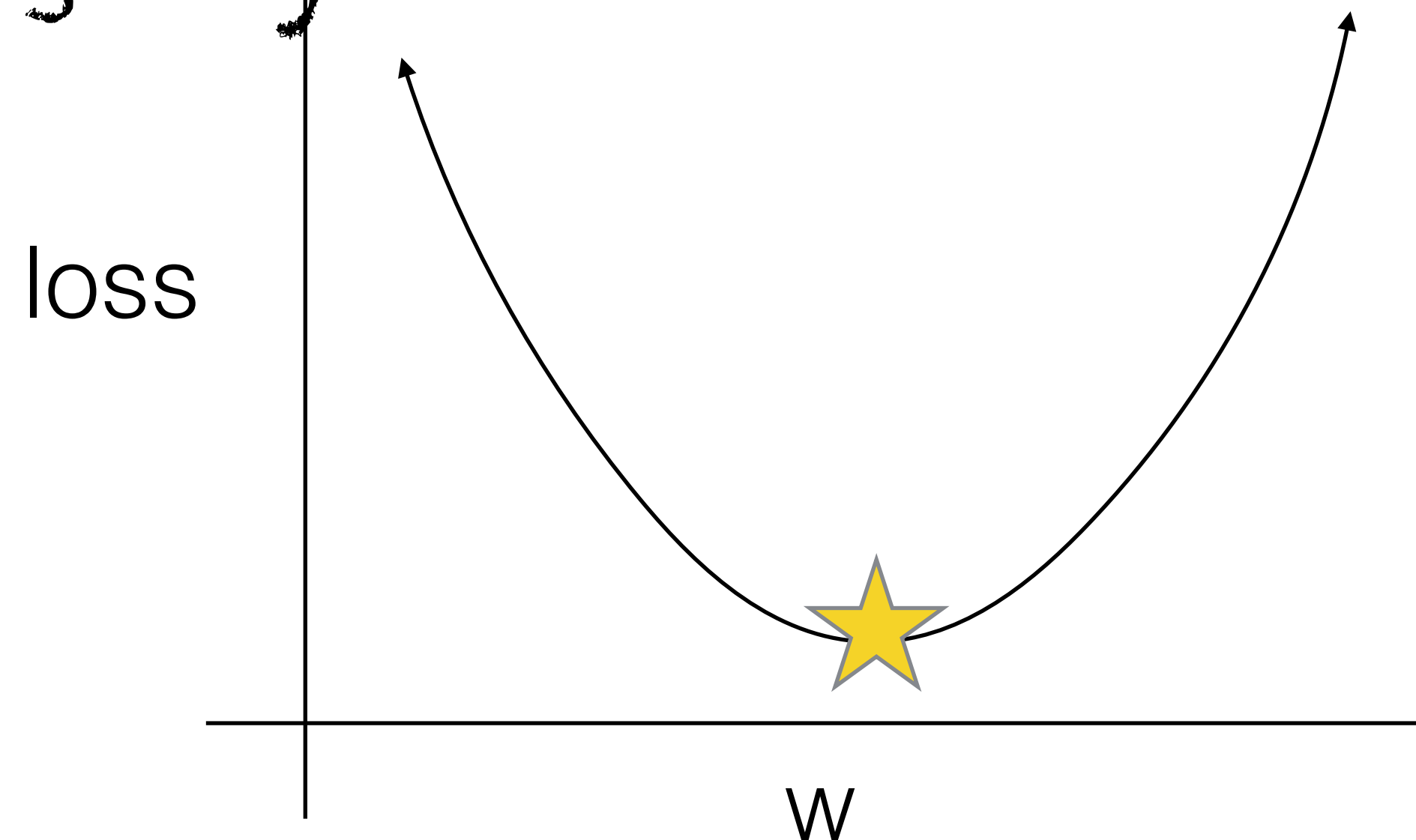
$$-Y \log \hat{Y} + (1 - Y) \log(1 - \hat{Y})$$



Logistic Regression Classifiers

Training with Gradient Descent

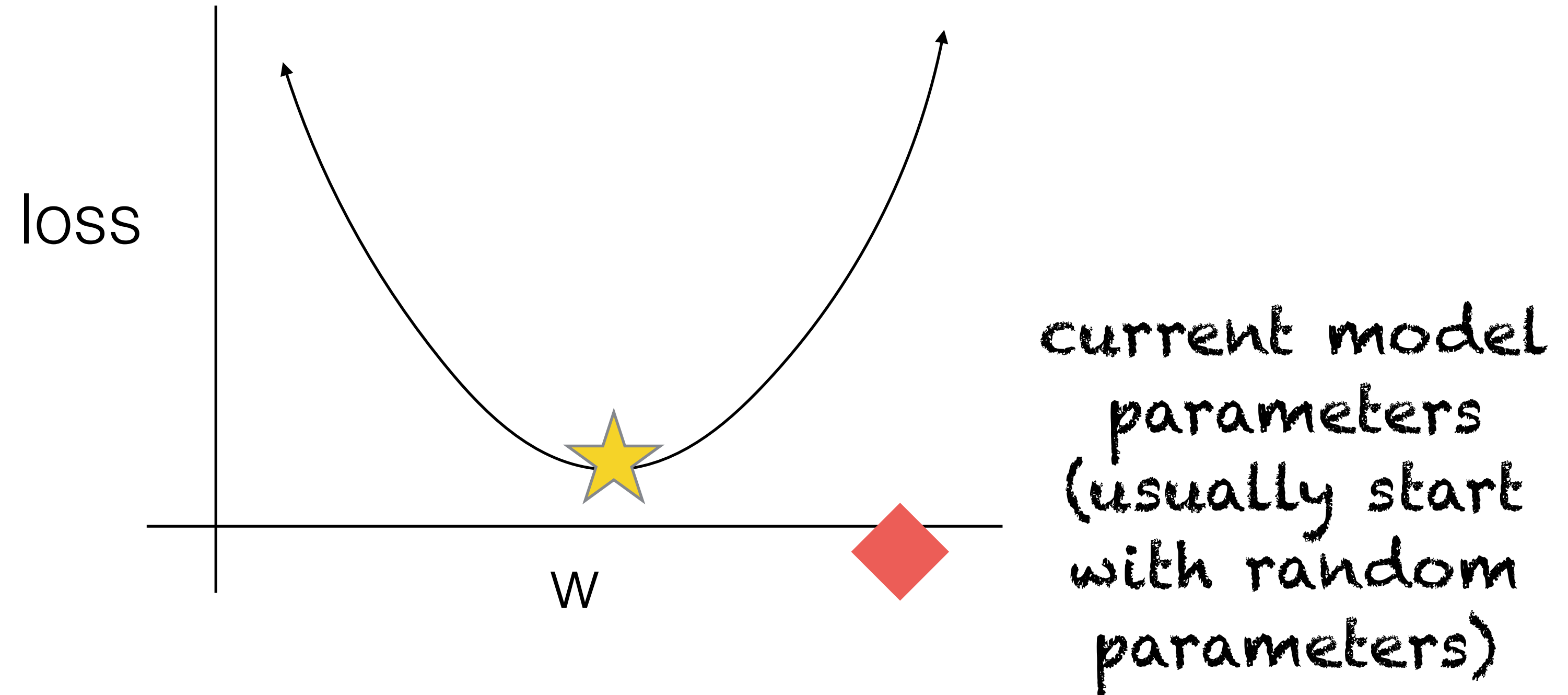
Goal (lowest achievable value for loss given data). You don't know what value of parameters will give you this.



Logistic Regression Classifiers

Training with Gradient Descent

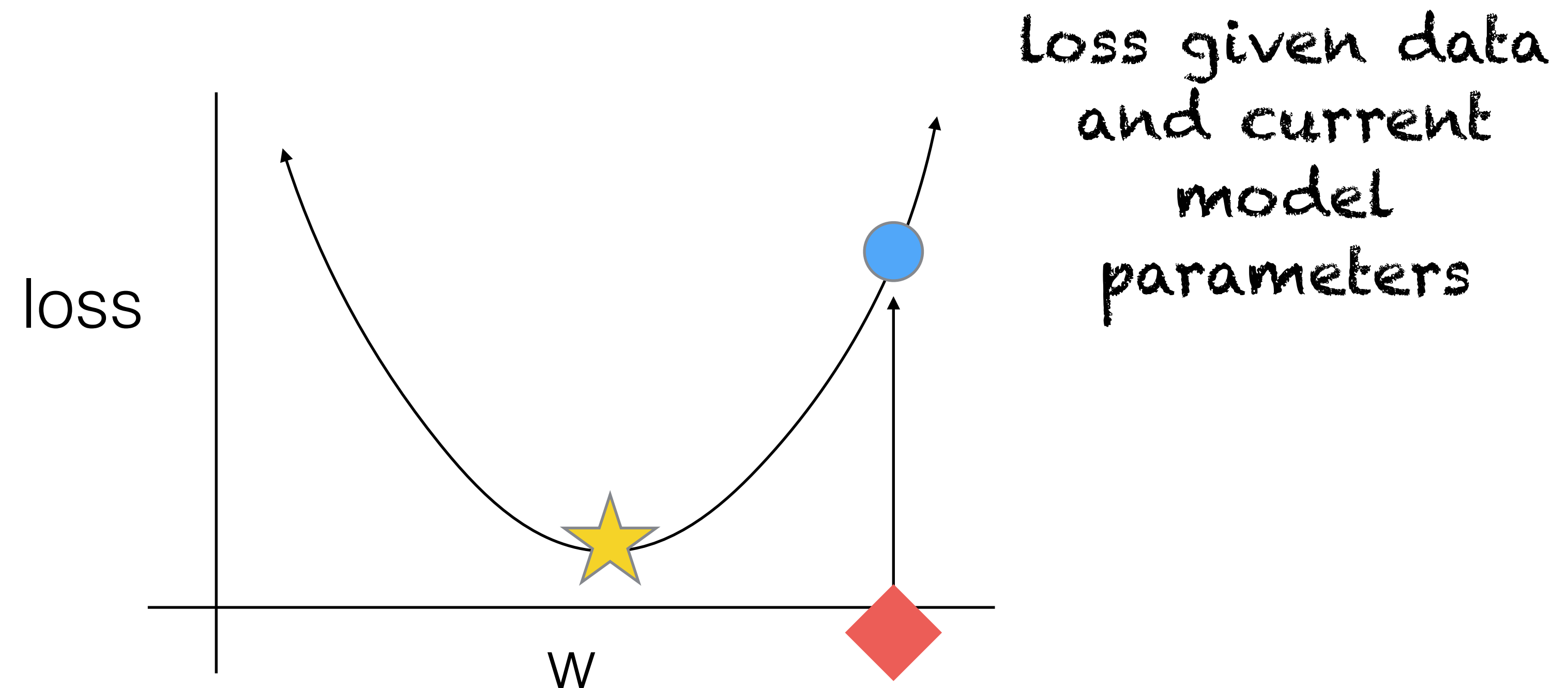
$$-Y \log \hat{Y} + (1 - Y) \log(1 - \hat{Y})$$



Logistic Regression Classifiers

Training with Gradient Descent

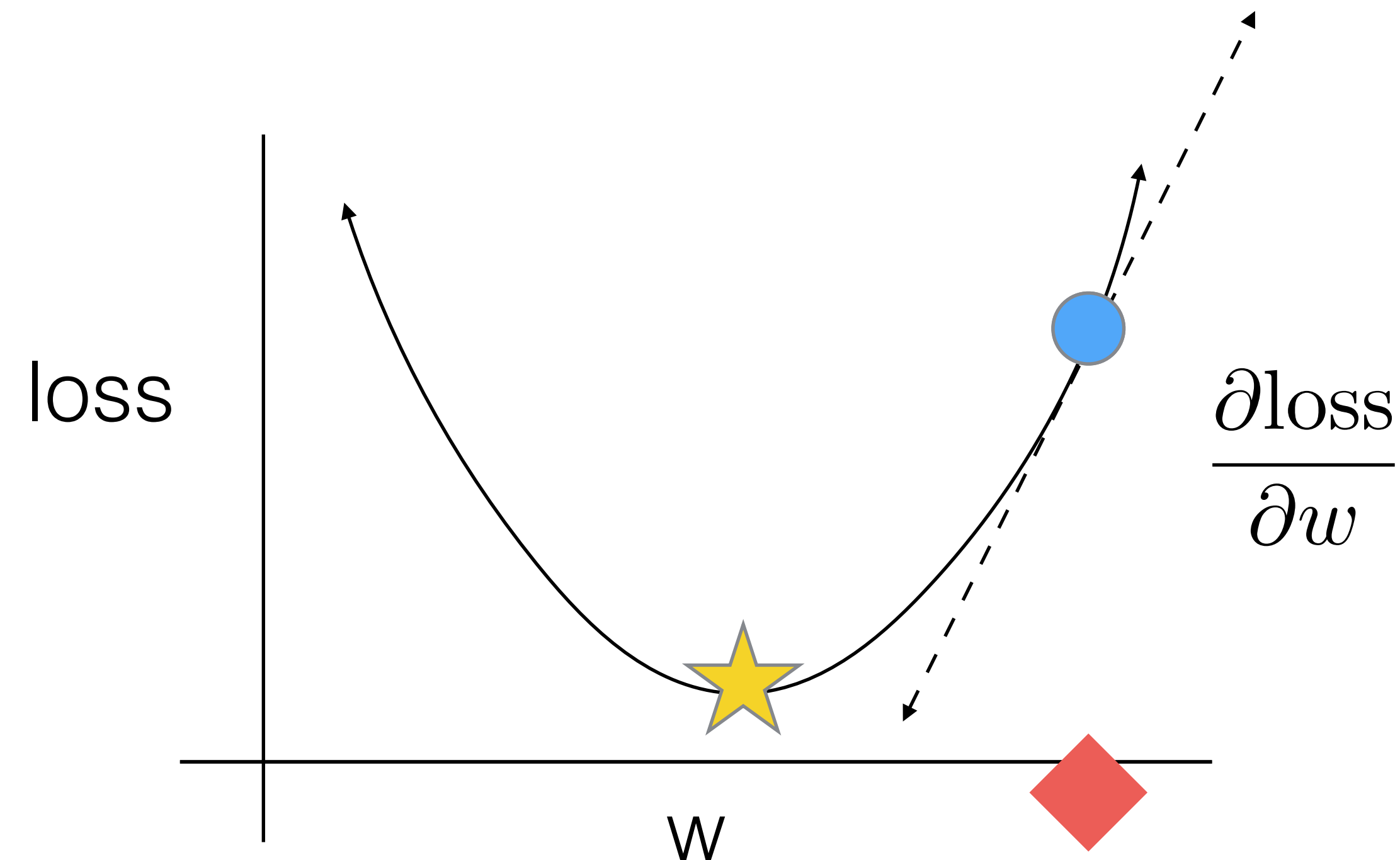
$$-Y \log \hat{Y} + (1 - Y) \log(1 - \hat{Y})$$



Logistic Regression Classifiers

Training with Gradient Descent

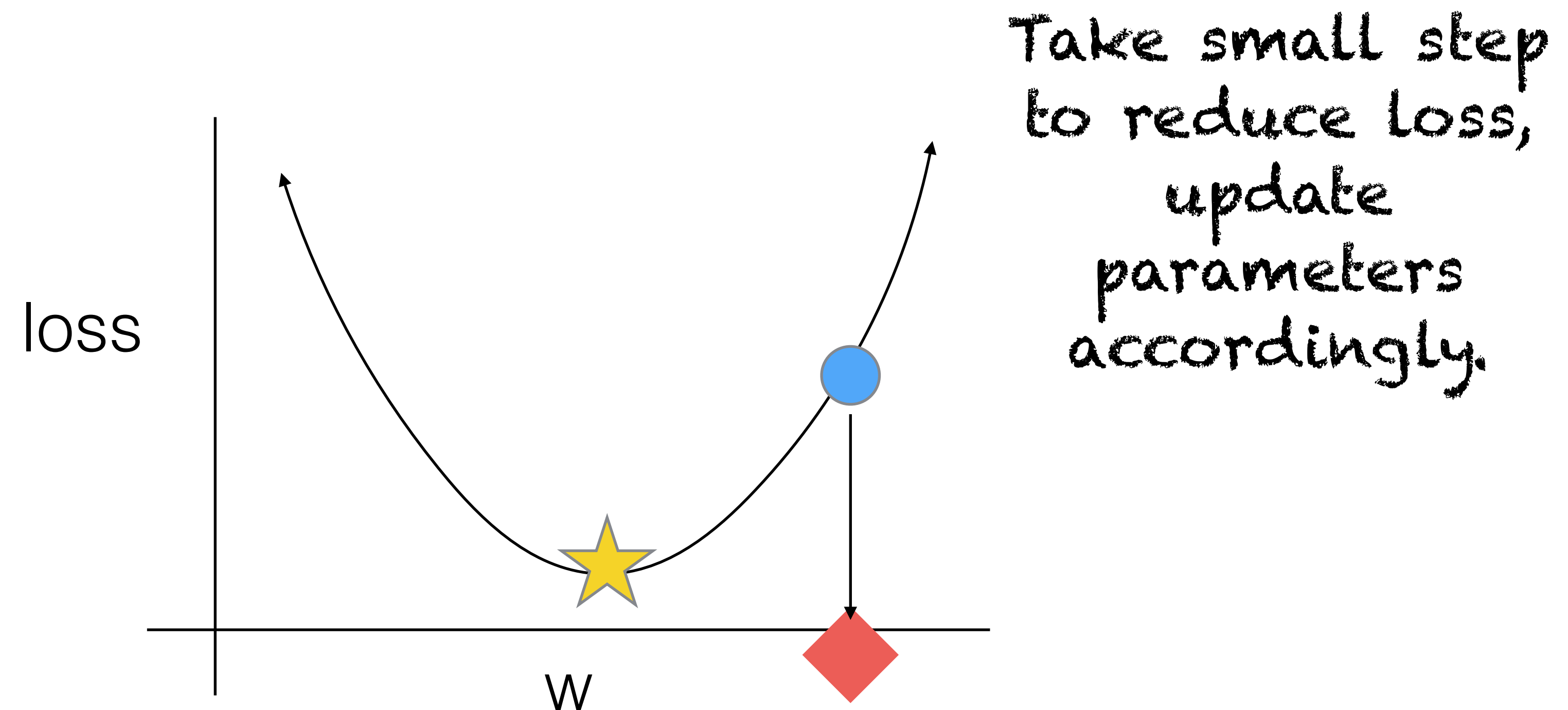
$$-Y \log \hat{Y} + (1 - Y) \log(1 - \hat{Y})$$



Logistic Regression Classifiers

Training with Gradient Descent

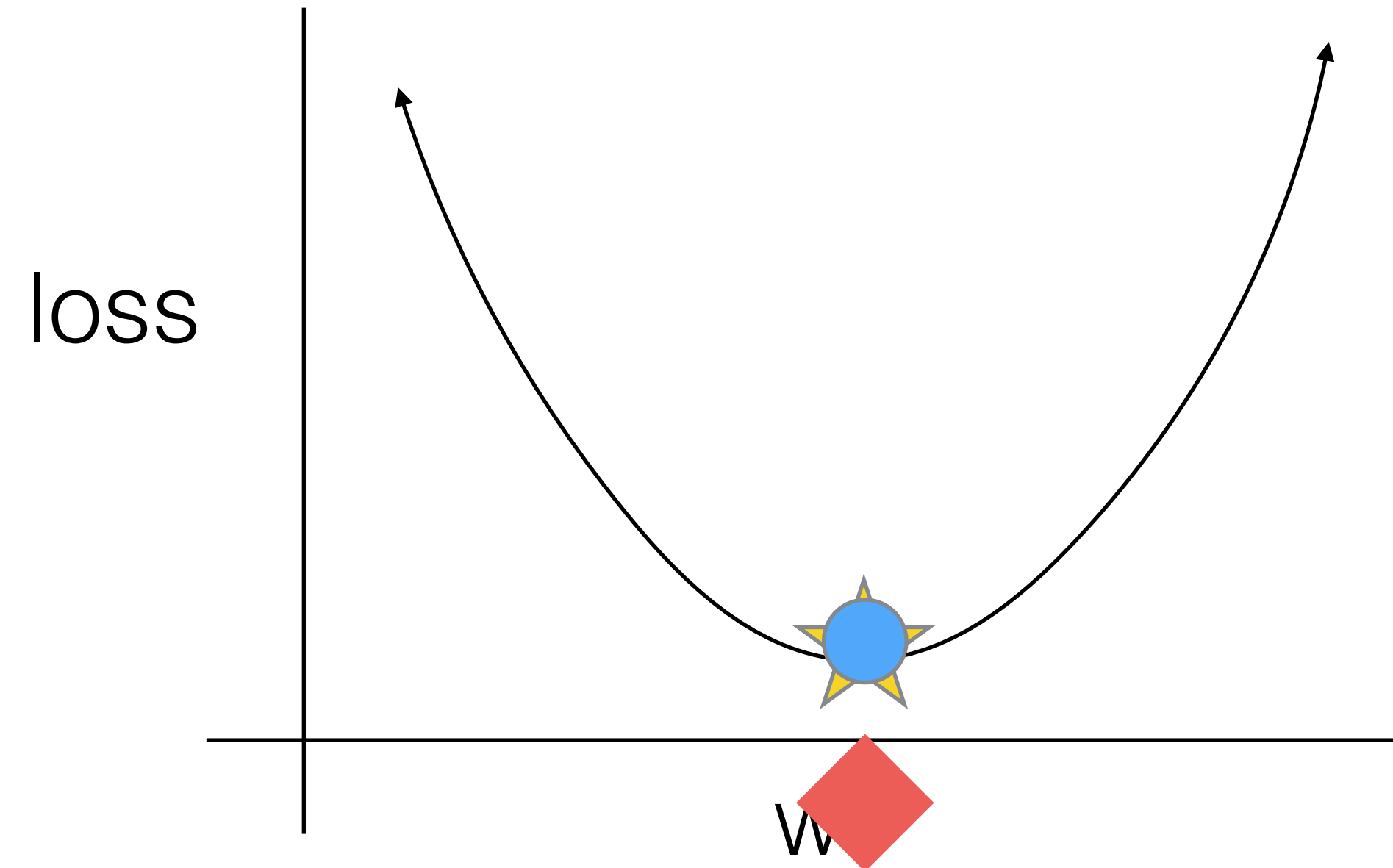
$$-Y \log \hat{Y} + (1 - Y) \log(1 - \hat{Y})$$



Logistic Regression Classifiers

Training with Gradient Descent

$$-Y \log \hat{Y} + (1 - Y) \log(1 - \hat{Y})$$

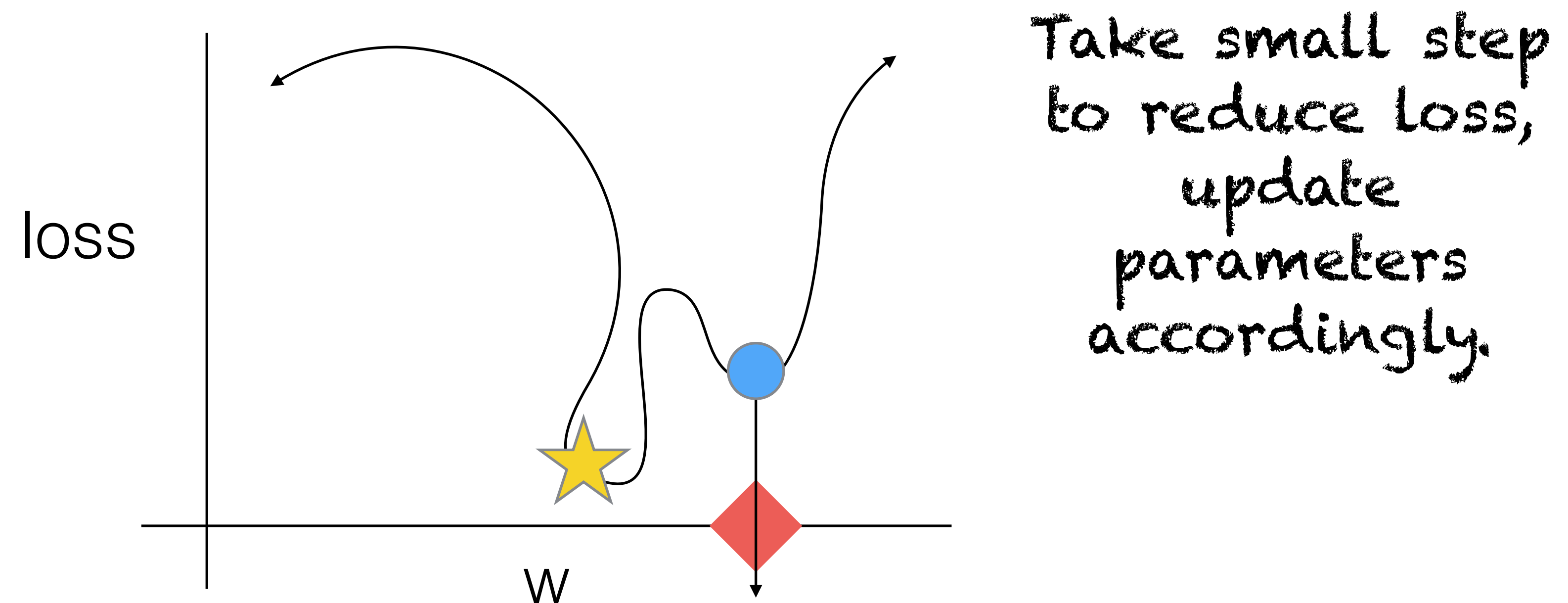


Repeat until
you converge,
i.e., loss can't
be decreased,
or you time
out (like in
kmeans).

Logistic Regression Classifiers

Training with Gradient Descent

$$-Y \log \hat{Y} + (1 - Y) \log(1 - \hat{Y})$$



Topics

- More Followup on Word Embeddings from SVD
- Logistic Regression and Gradient Descent
- **Multilayer Perceptrons**
- Word Embeddings from NNs

Language Modeling Task

Running Example

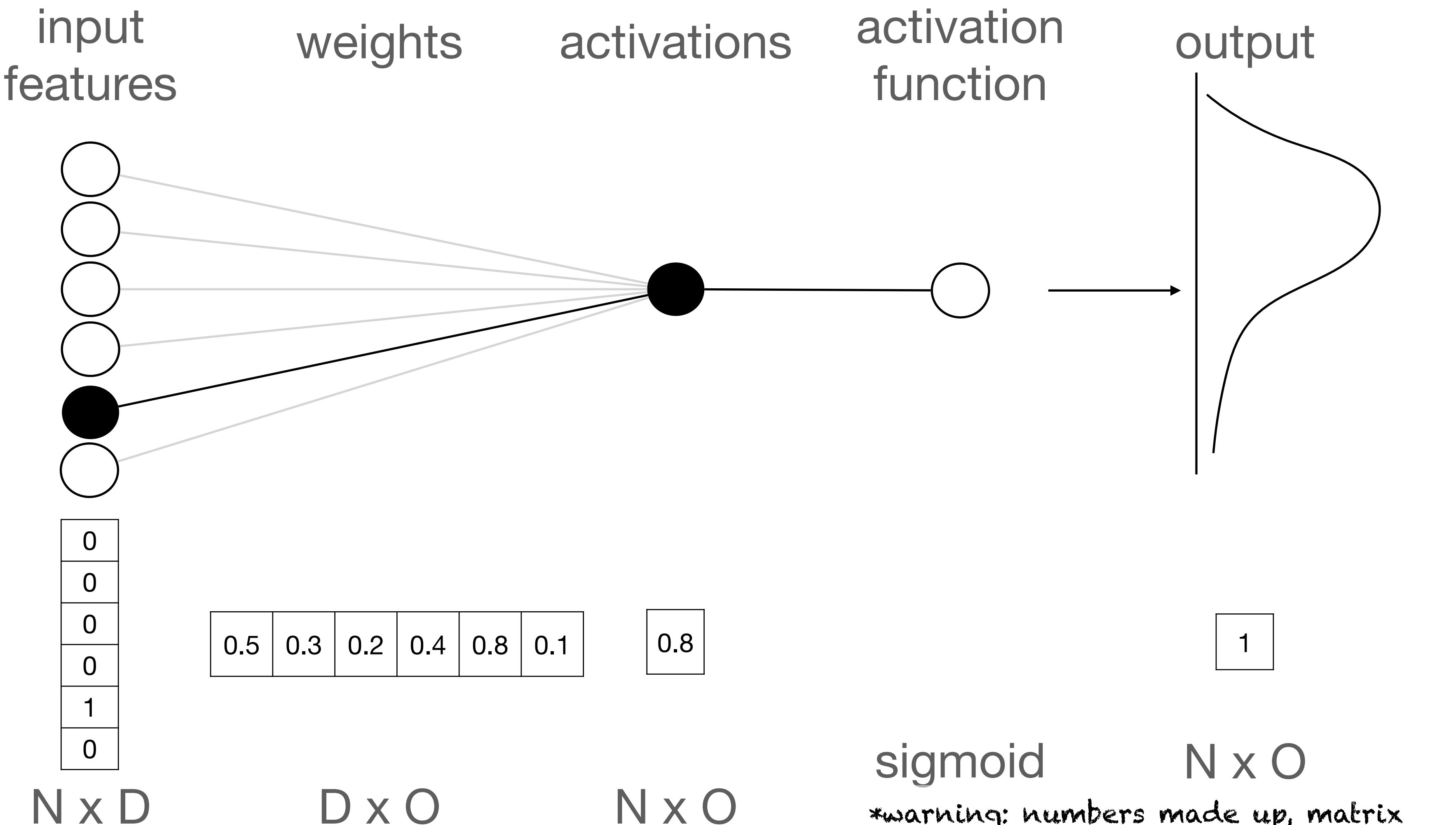
Task: Predict the next word in a sentence.

The cat sat on the ____

Basic Perceptron

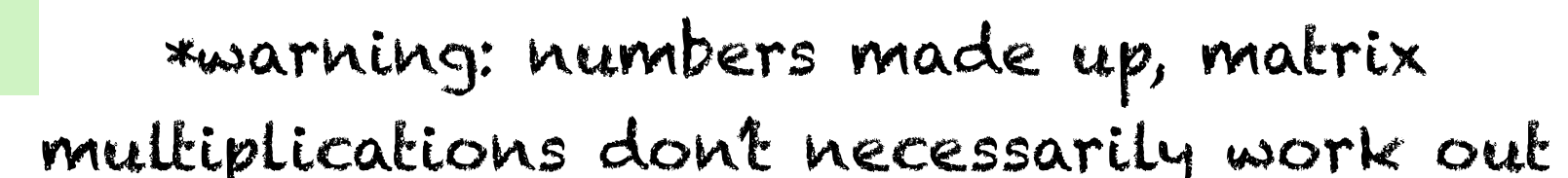
Same as Logistic Regression

Task: Predict the next word
Input: the
Expected: cat



*warning: numbers made up, matrix multiplications don't necessarily work out

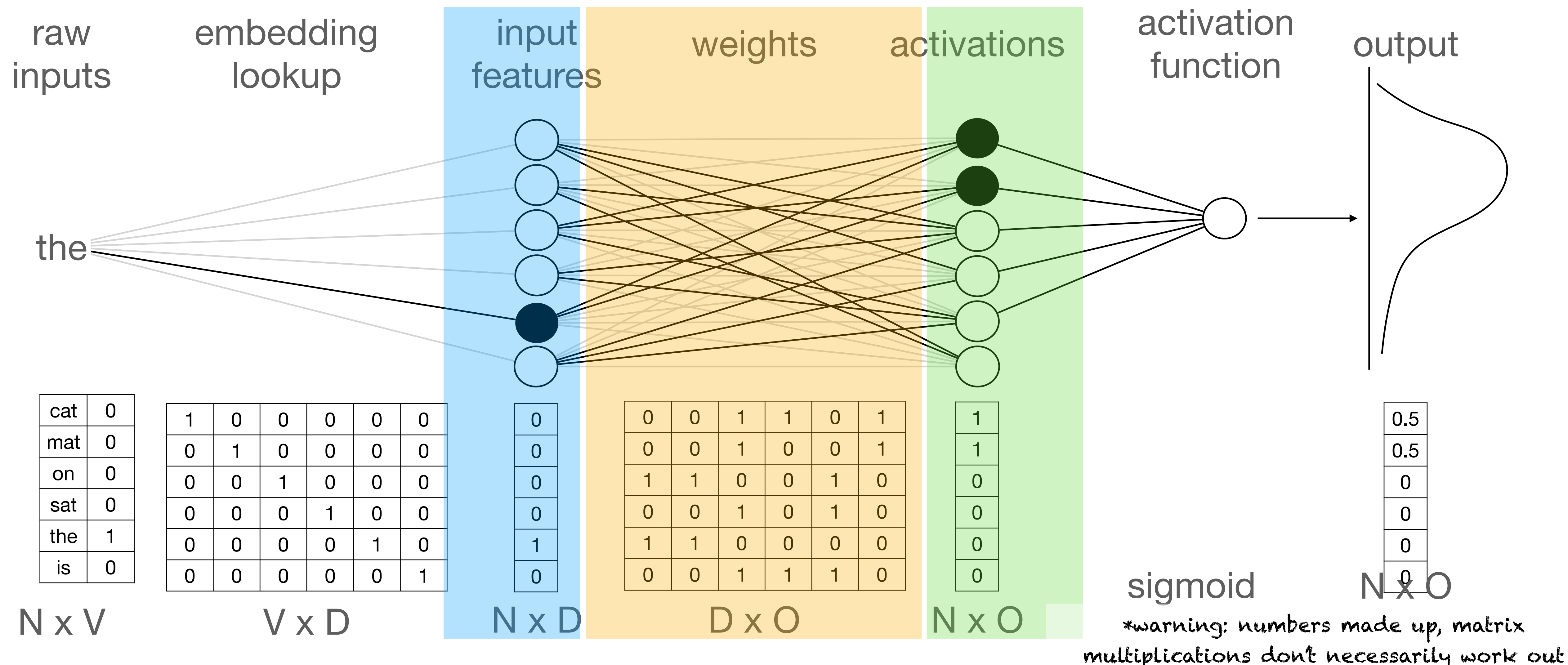
Task: Predict the next word
Input: the
Expected: cat



Basic Perceptron

Forward Pass

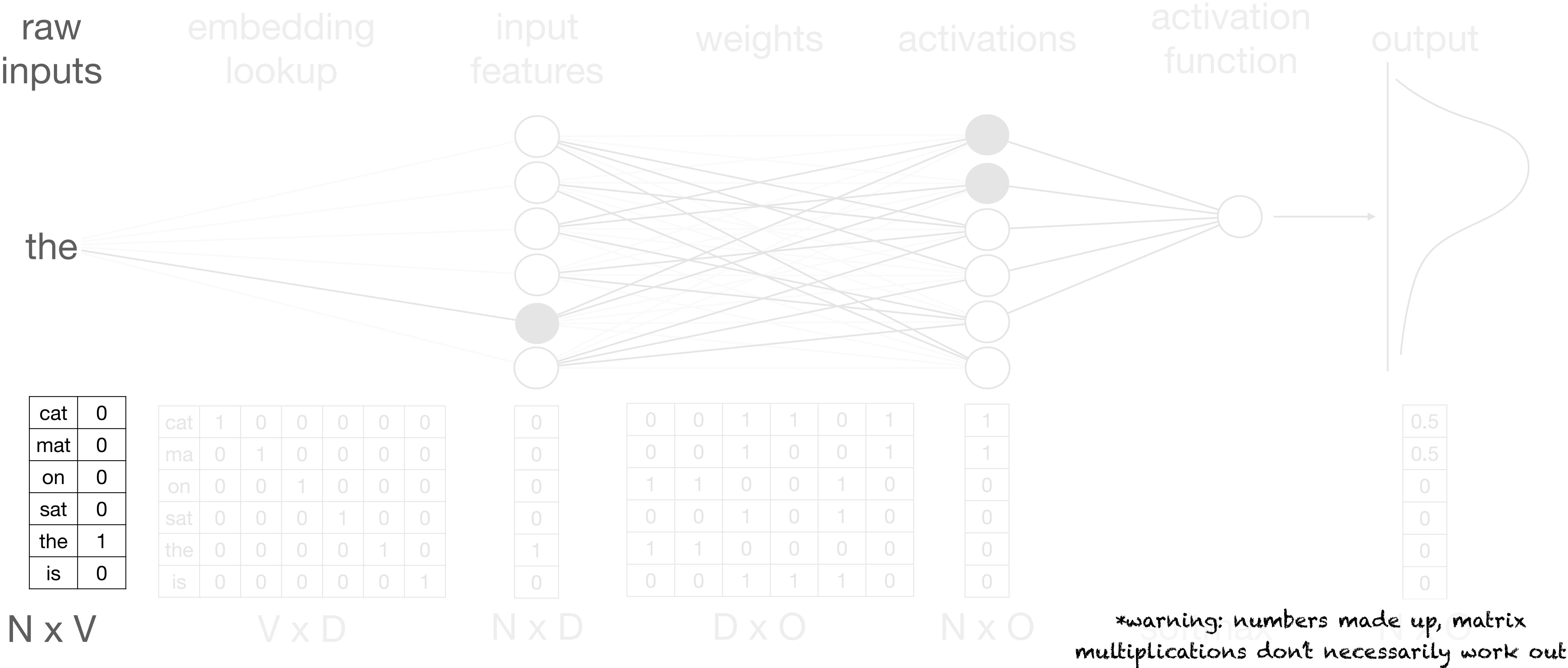
Task: Predict the next word
Input: the
Expected: cat



Basic Perceptron

Forward Pass

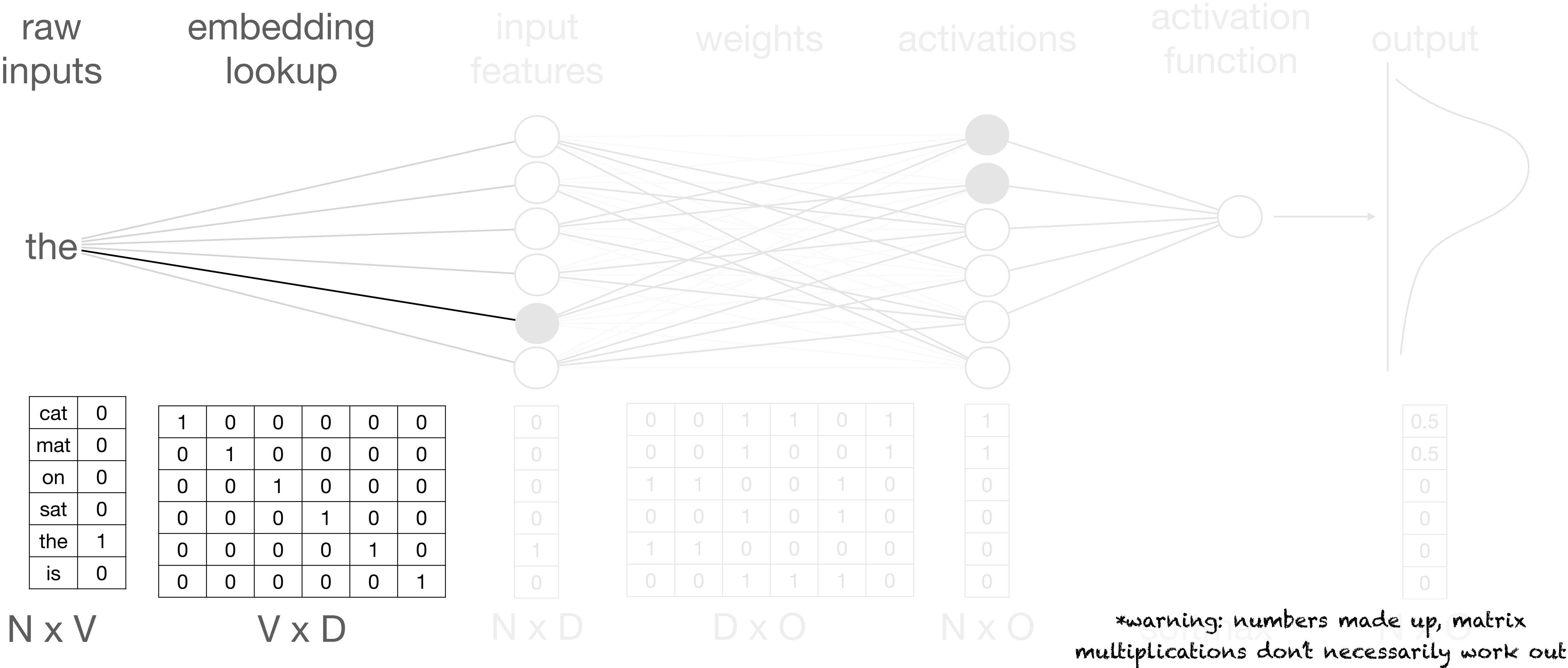
Task: Predict the next word
Input: the
Expected: cat



Basic Perceptron

Forward Pass

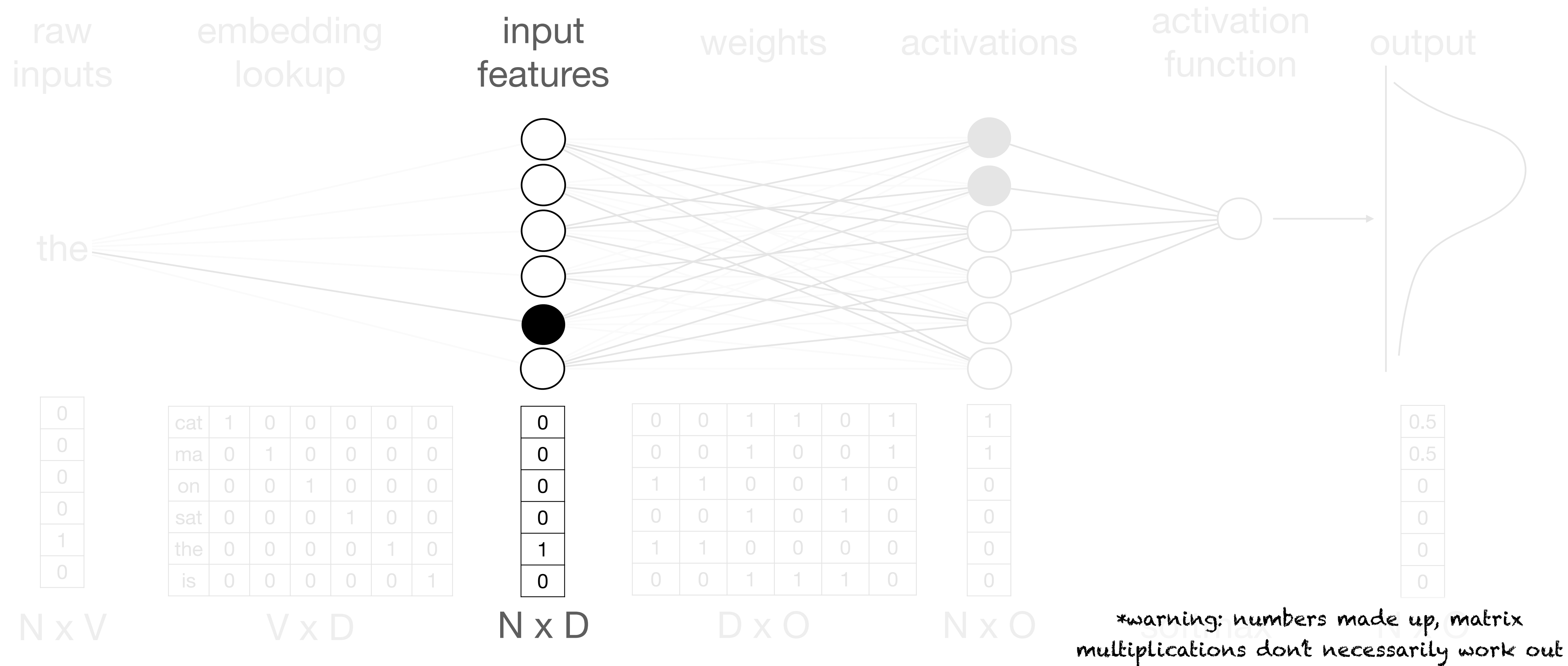
Task: Predict the next word
Input: the
Expected: cat



Basic Perceptron

Forward Pass

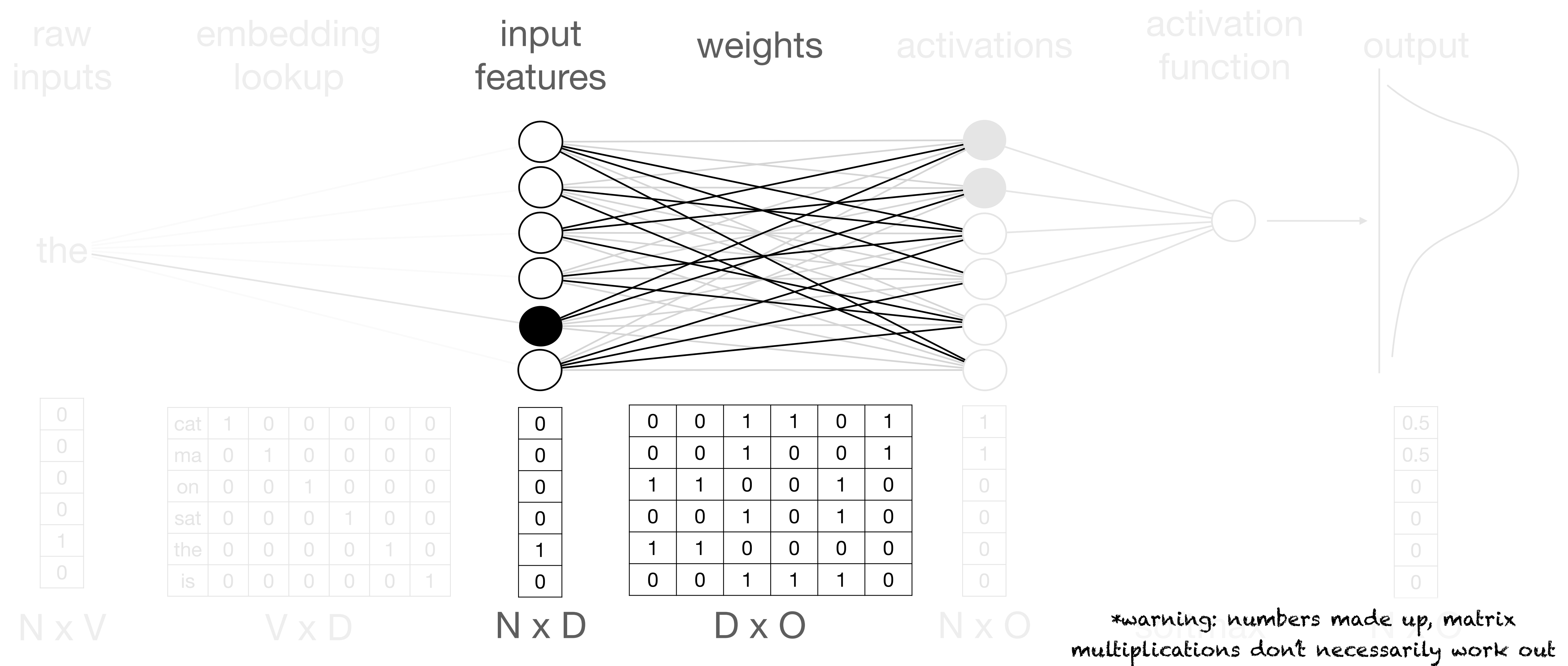
Task: Predict the next word
Input: the
Expected: cat



Basic Perceptron

Forward Pass

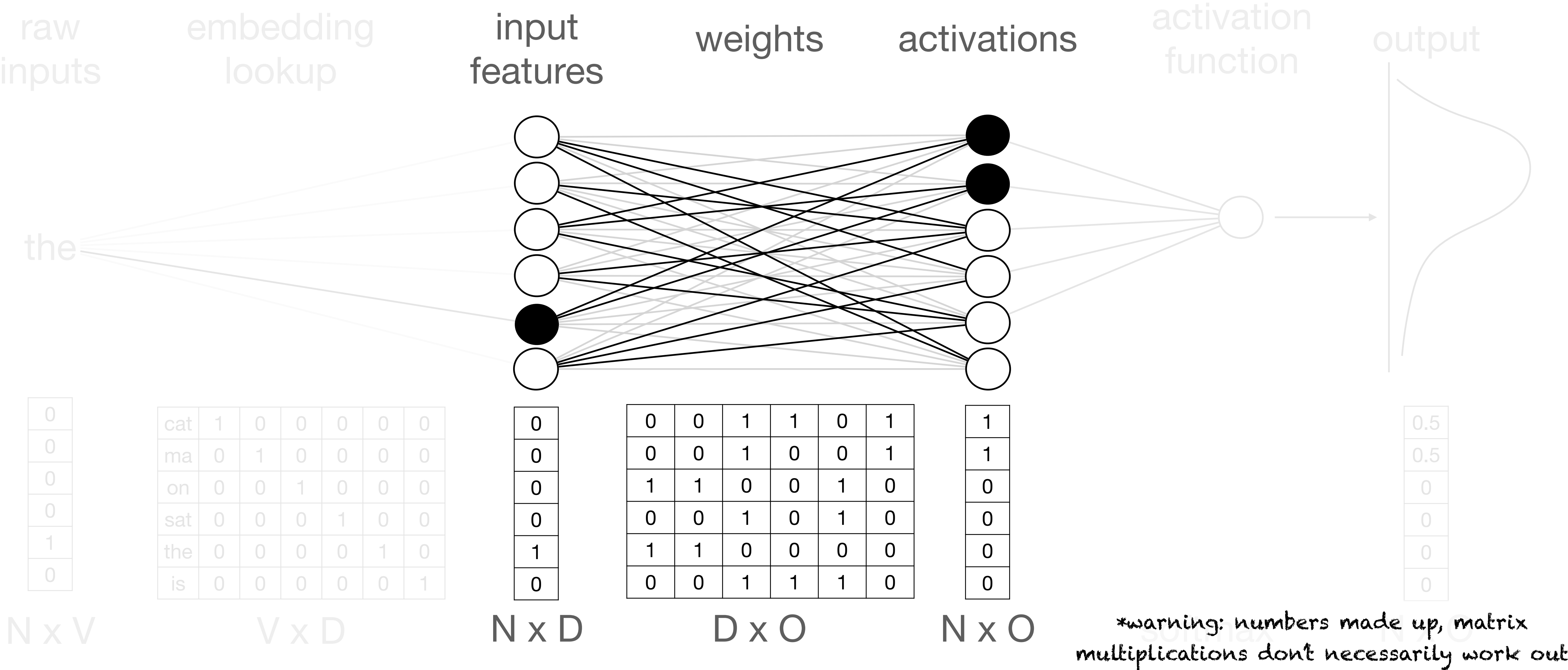
Task: Predict the next word
Input: the
Expected: cat



Basic Perceptron

Forward Pass

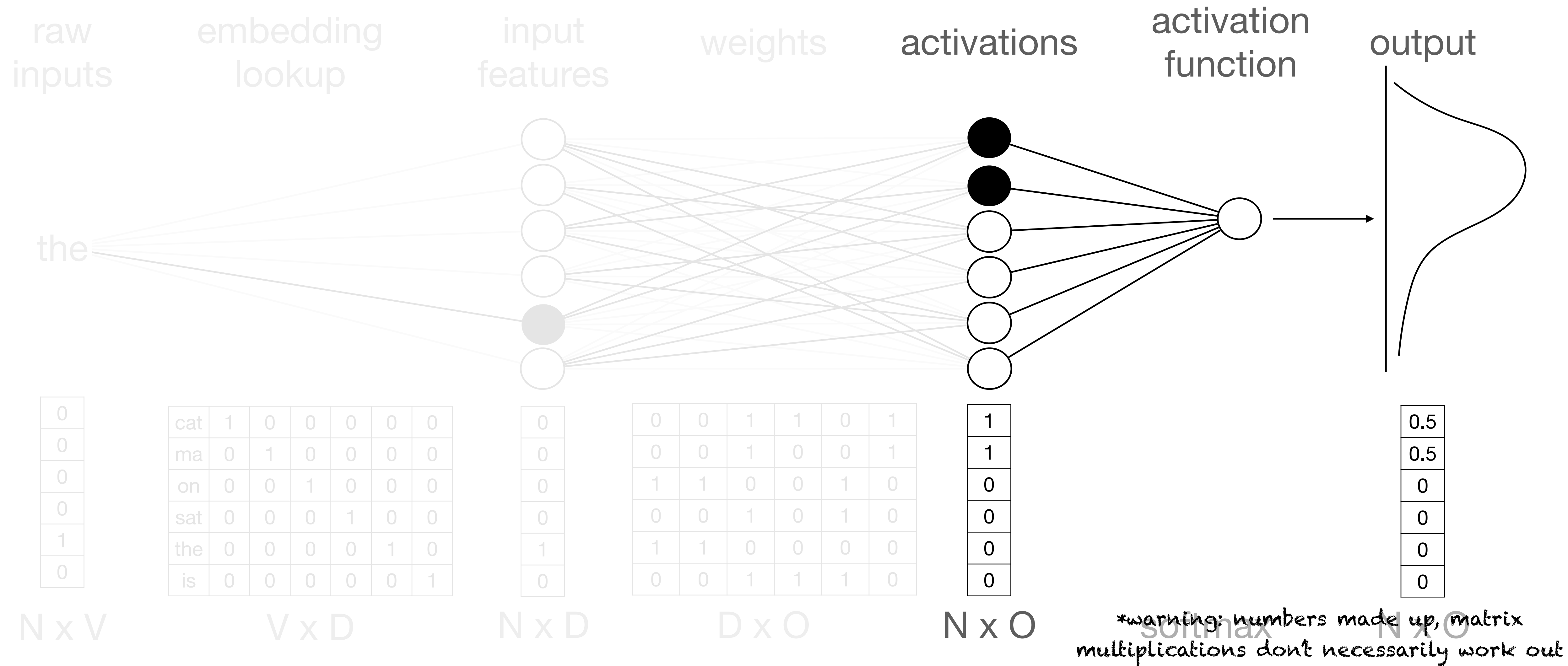
Task: Predict the next word
Input: the
Expected: cat



Basic Perceptron

Forward Pass

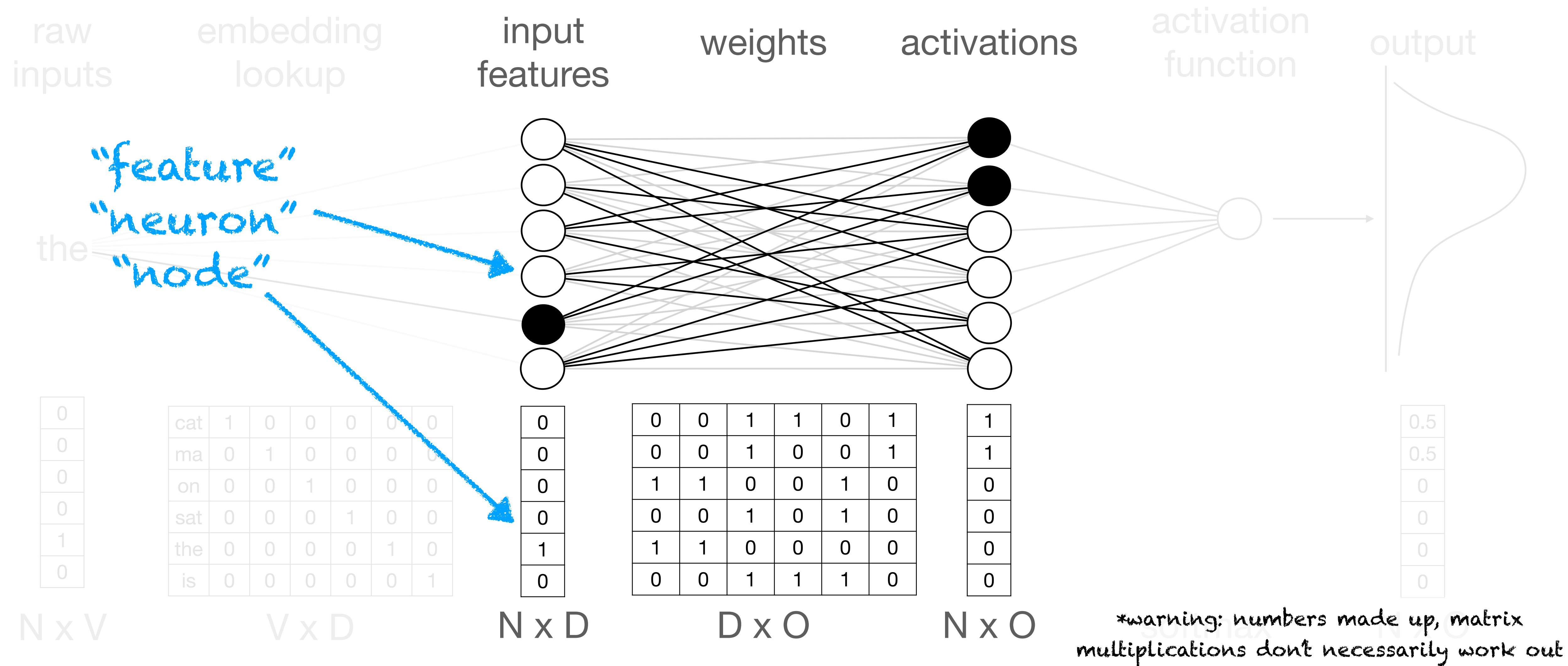
Task: Predict the next word
Input: the
Expected: cat



Basic Perceptron

Forward Pass

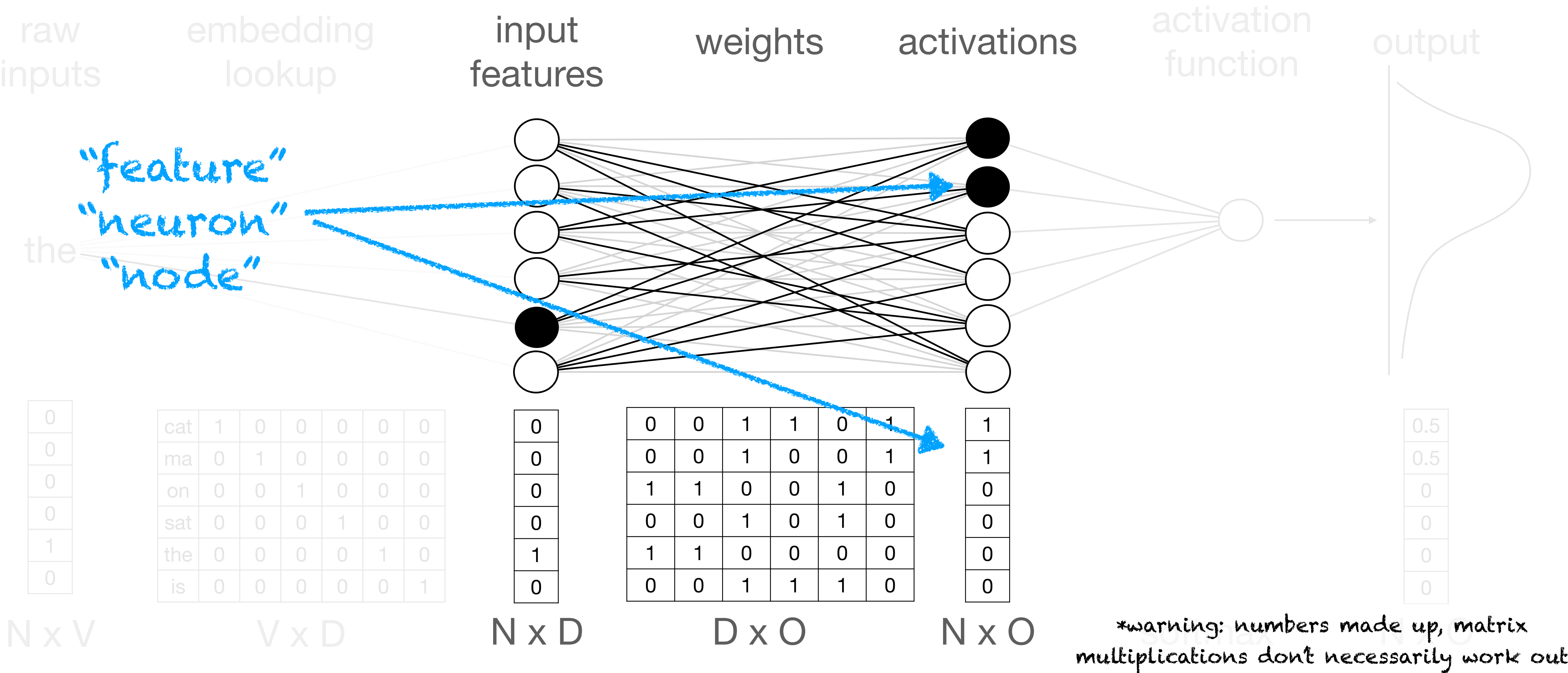
Task: Predict the next word
Input: the
Expected: cat



Basic Perceptron

Forward Pass

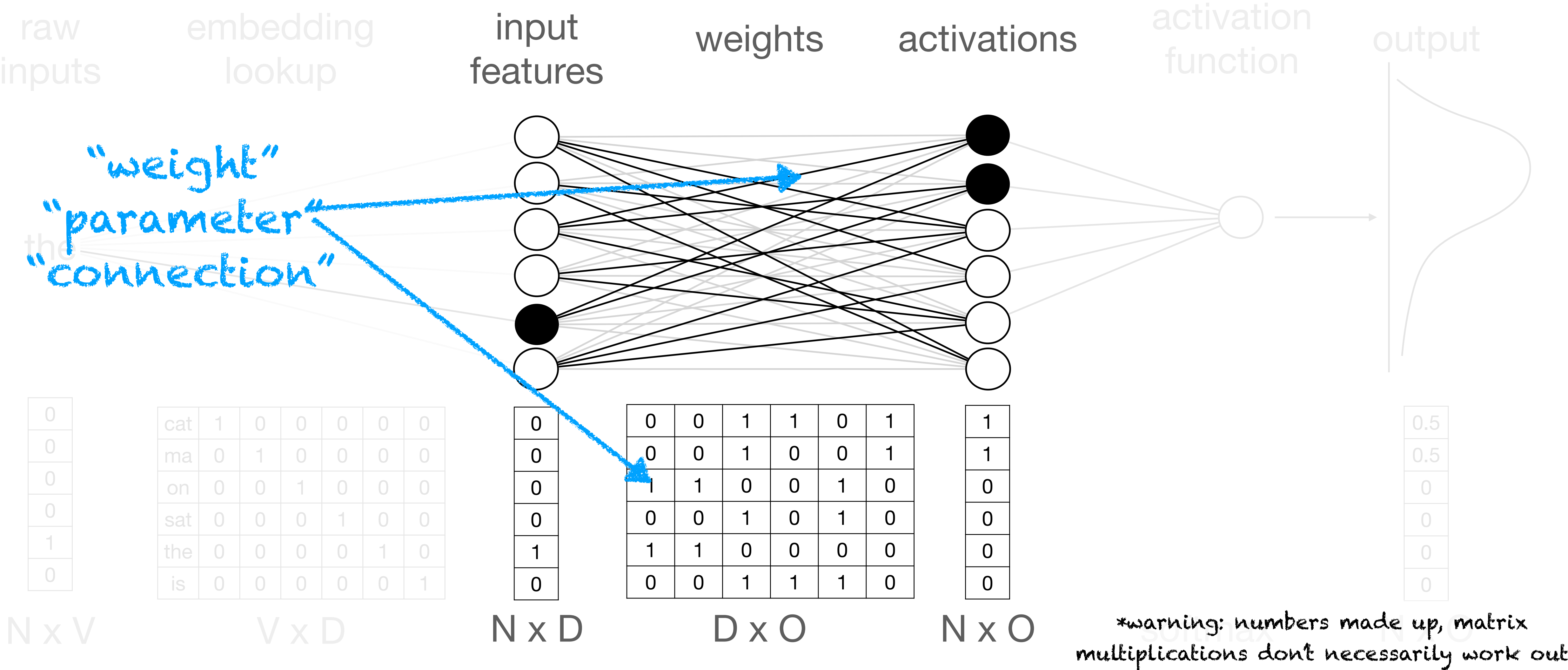
Task: Predict the next word
Input: the
Expected: cat



Basic Perceptron

Forward Pass

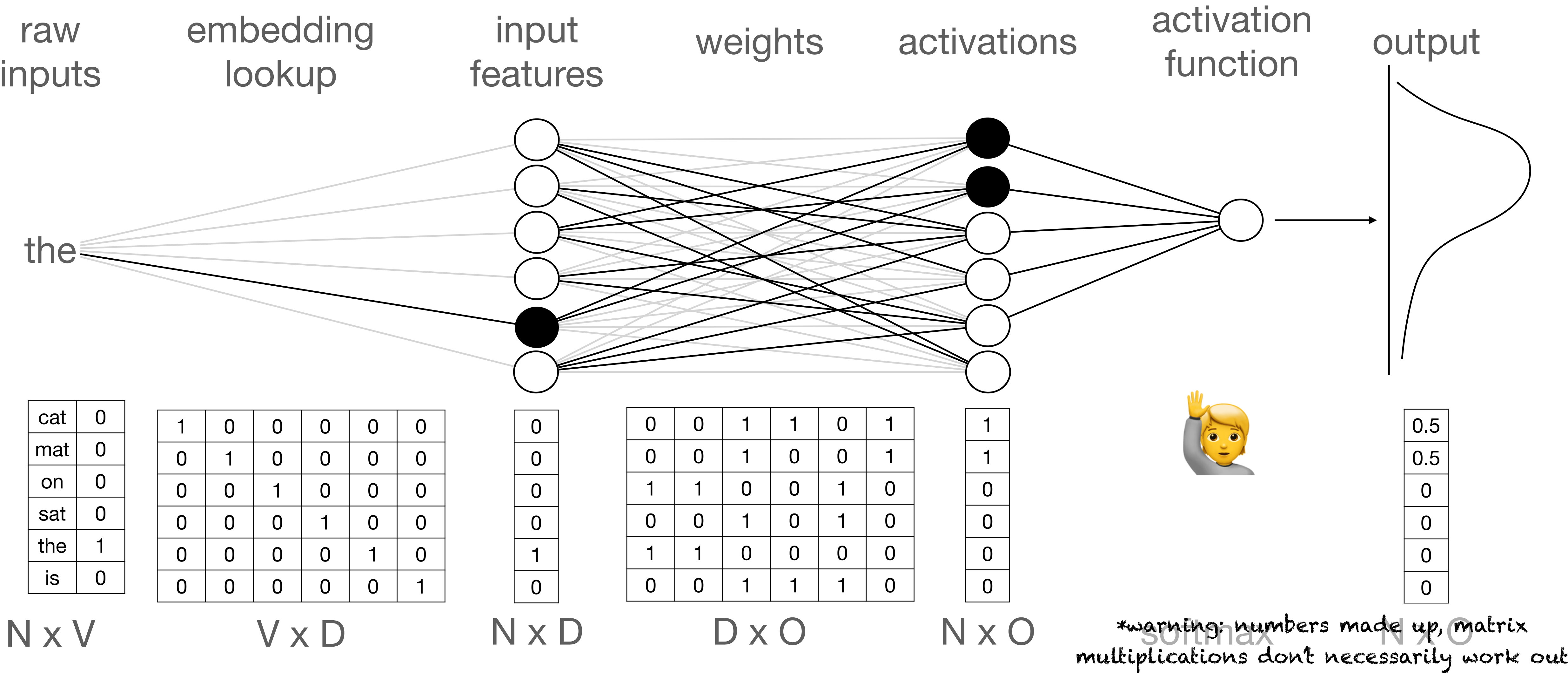
Task: Predict the next word
 Input: the
 Expected: cat



Basic Perceptron

Forward Pass

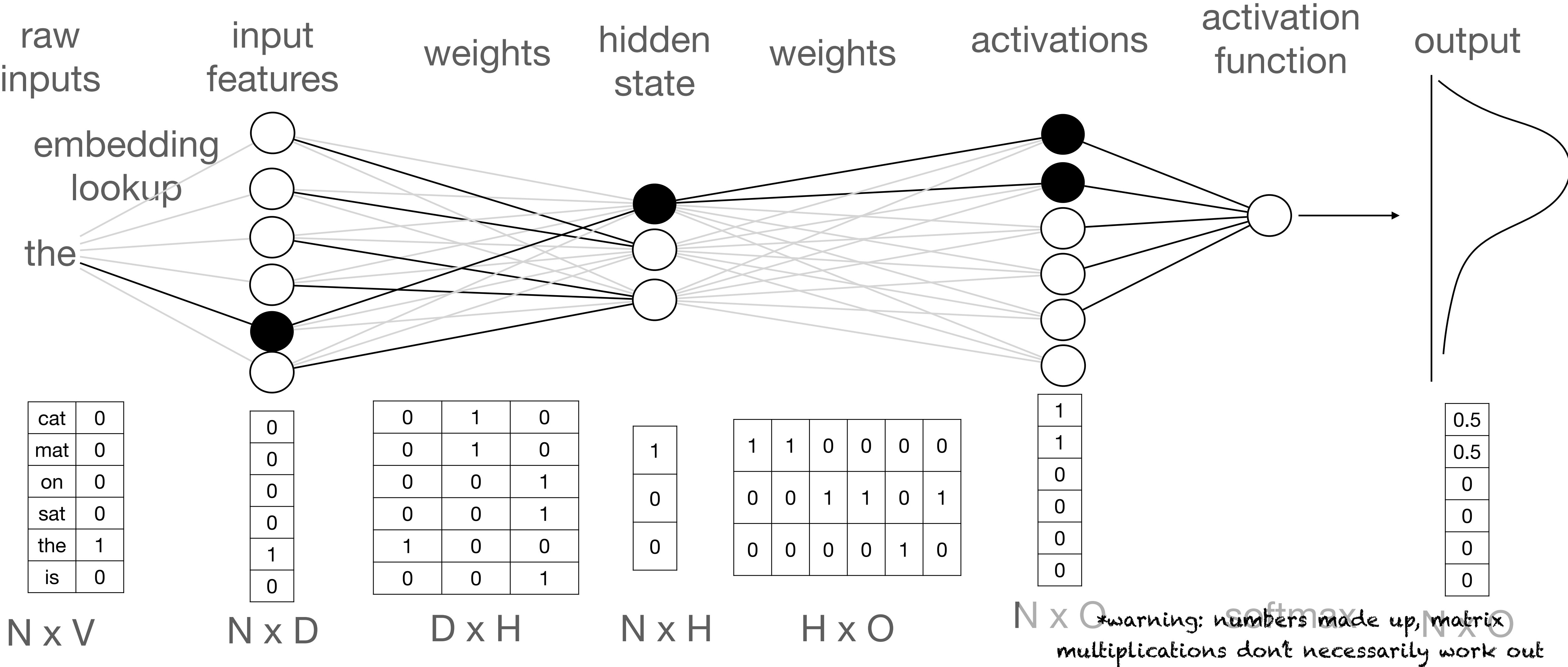
Task: Predict the next word
 Input: the
 Expected: cat



Multilayer Perceptron

Forward Pass

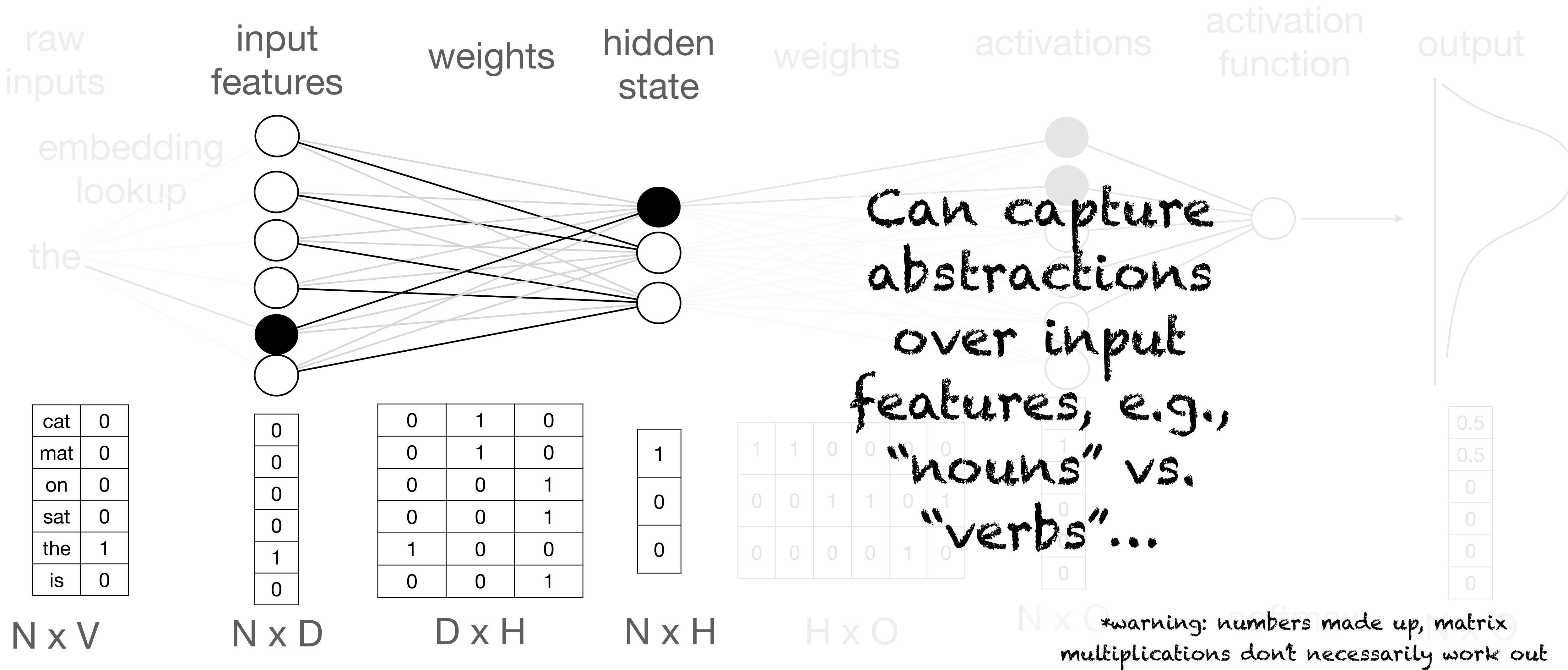
Task: Predict the next word
Input: the
Expected: cat



Multilayer Perceptron

Forward Pass

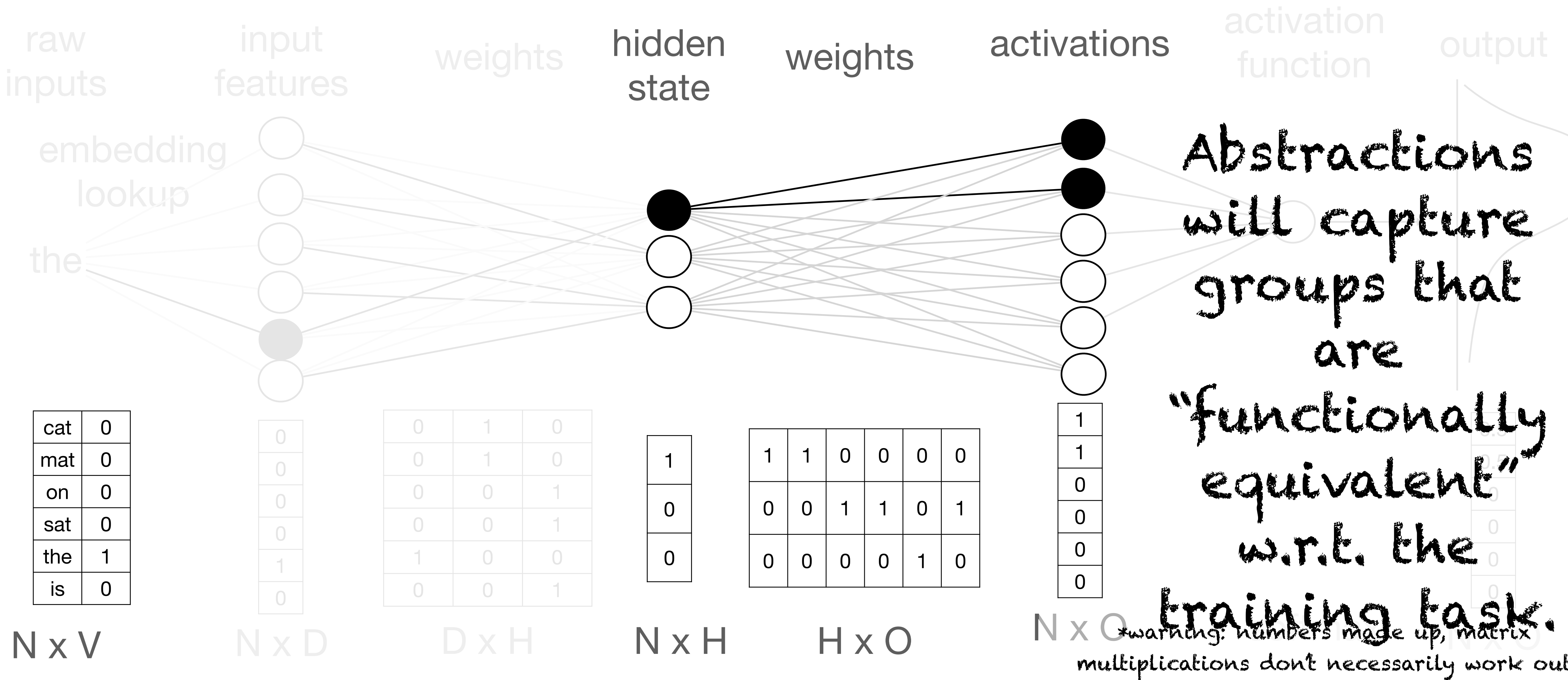
Task: Predict the next word
Input: the
Expected: cat



Multilayer Perceptron

Forward Pass

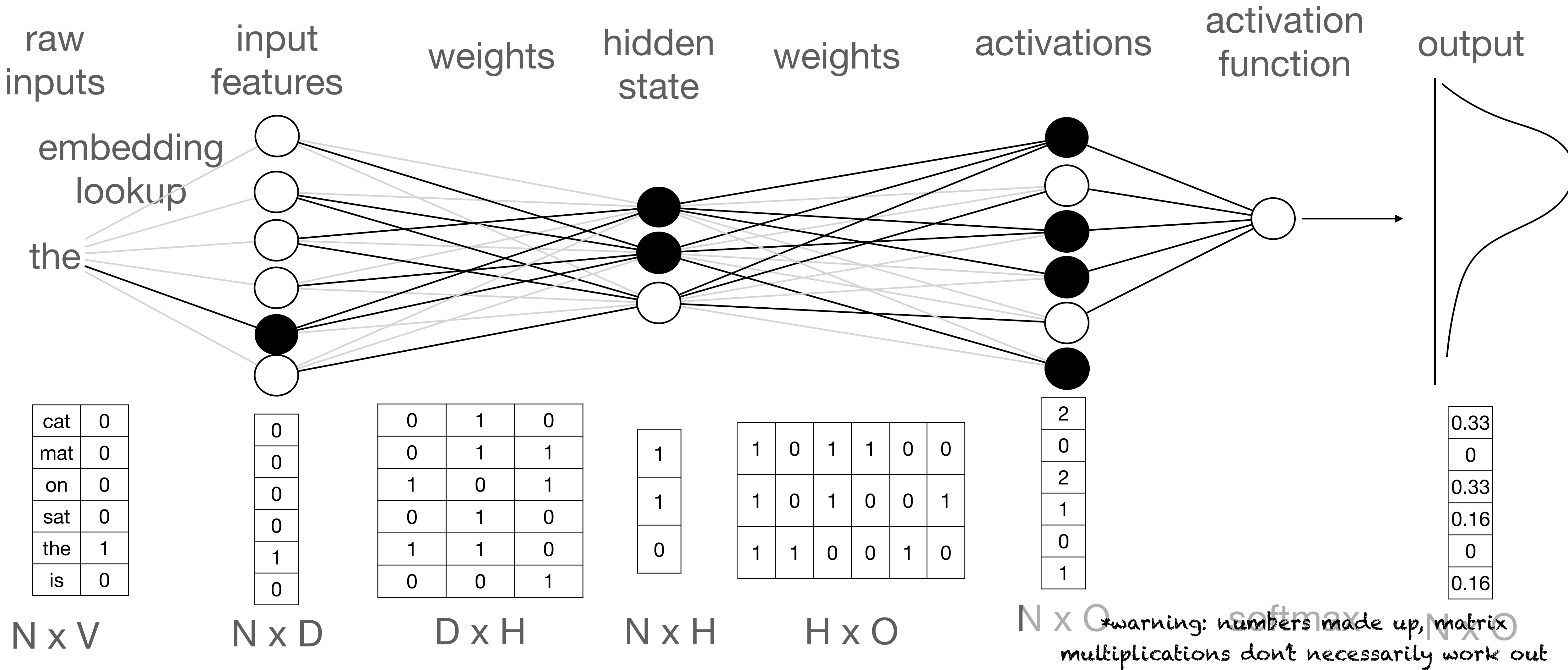
Task: Predict the next word
Input: the
Expected: cat



Multilayer Perceptron

Training with Backpropagation

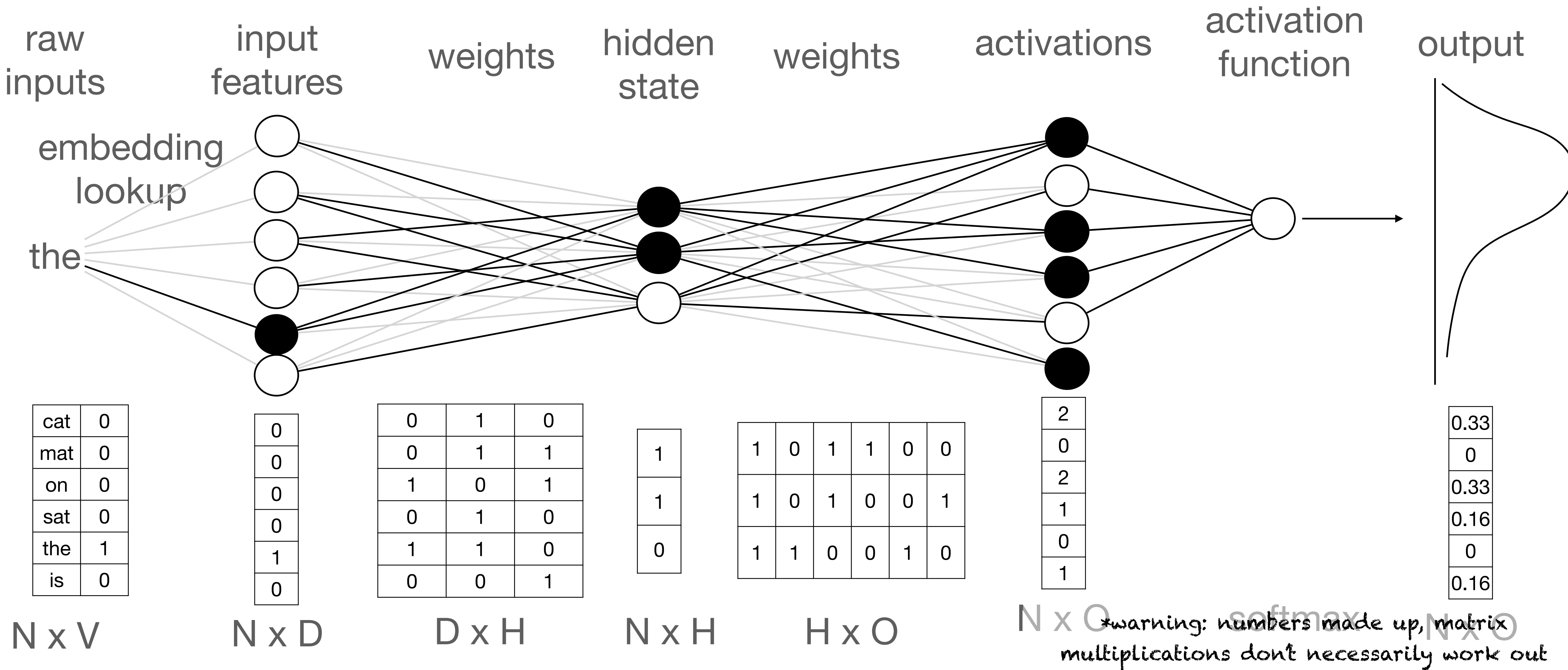
Parameters are randomly initialized.



Multilayer Perceptron

I.e., predictions are random

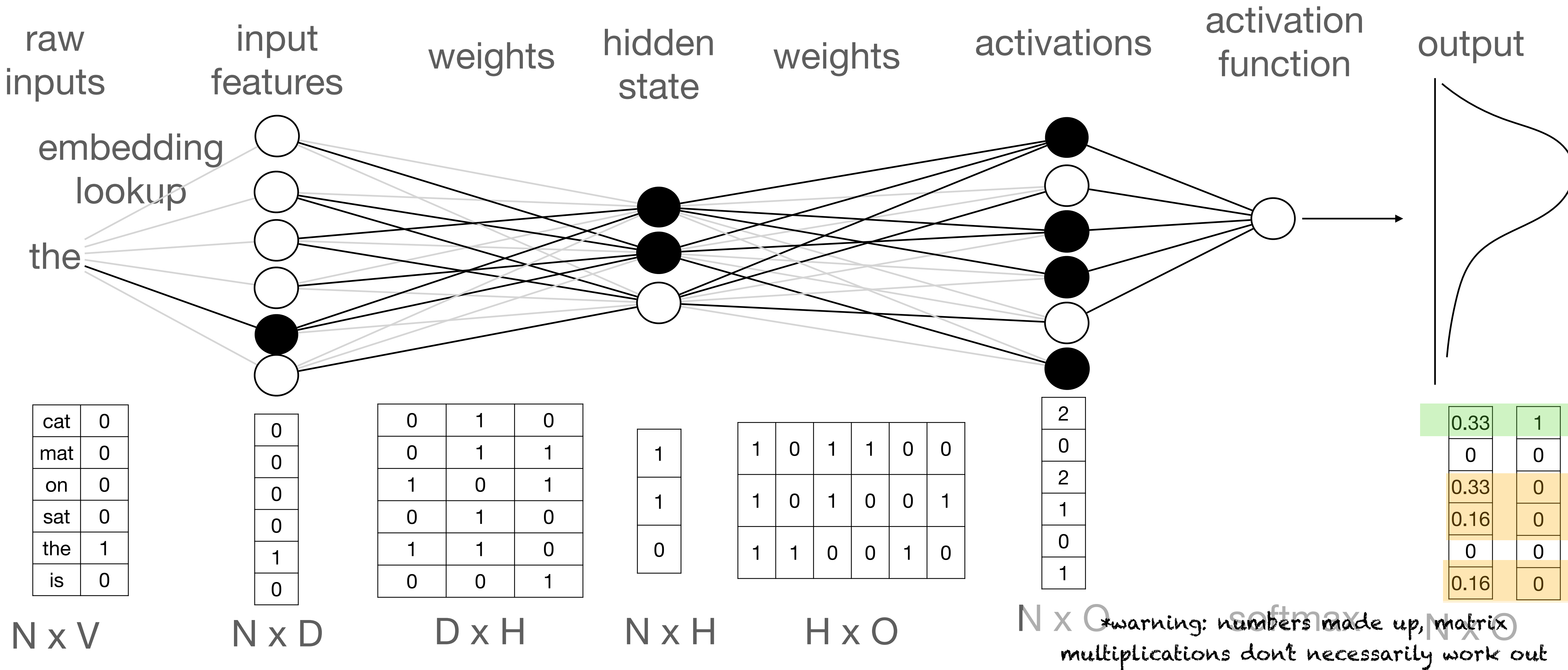
Training with Backpropagation



Multilayer Perceptron

Training with Backpropagation

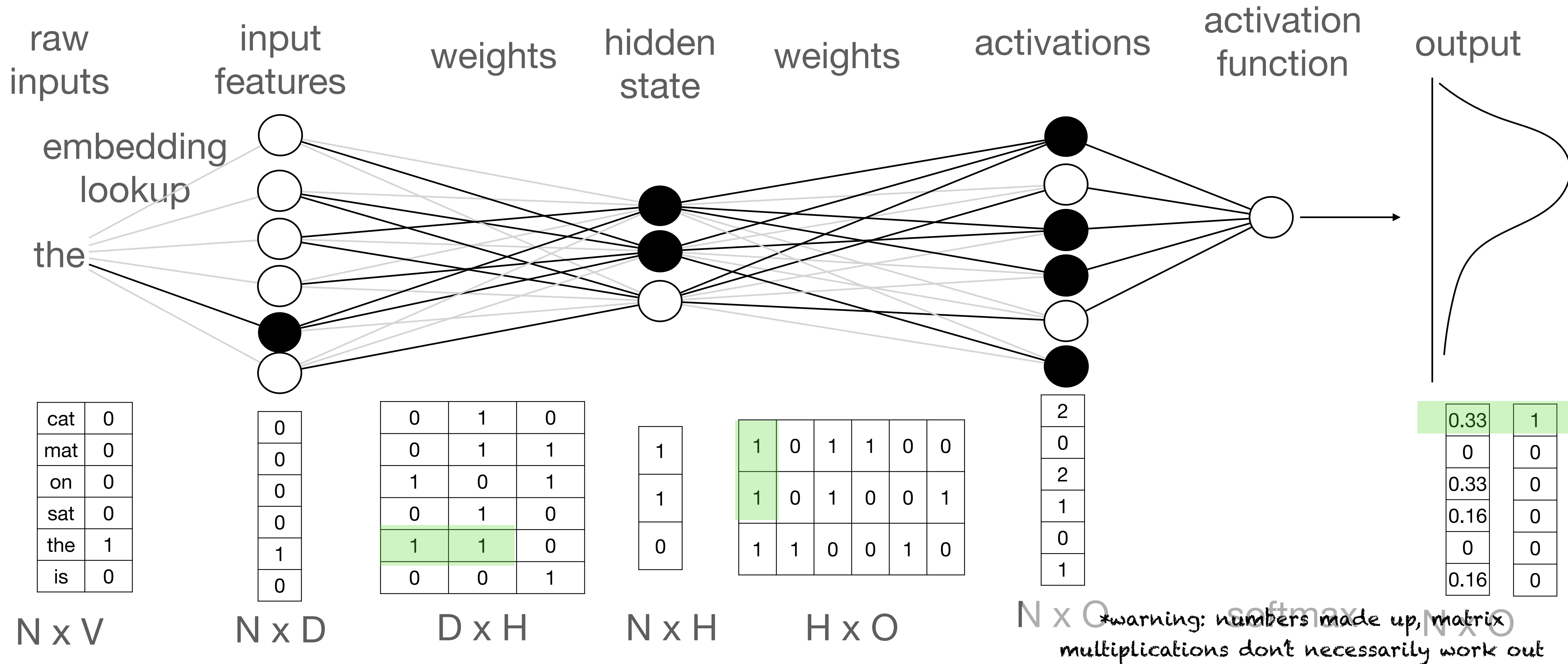
Compare predictions to ground truth output...



Multilayer Perceptron

Training with Backpropagation

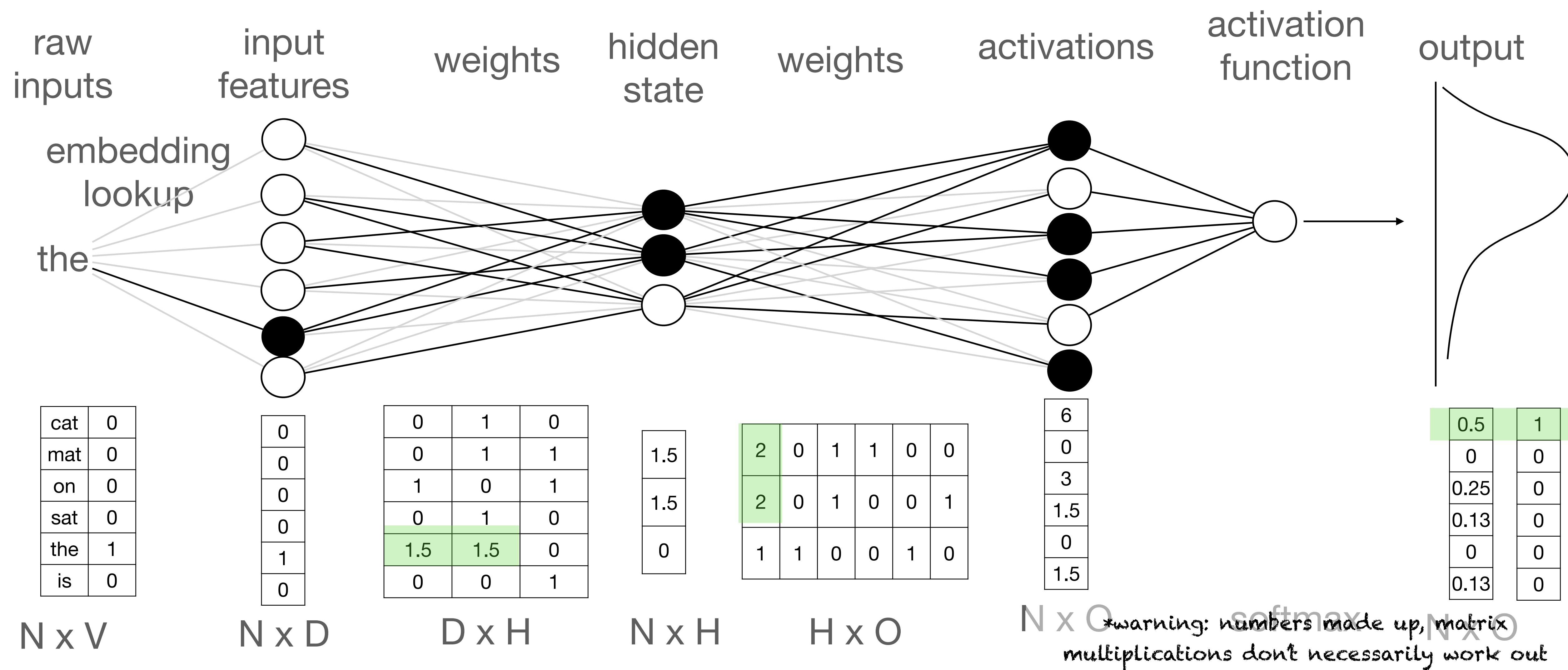
Adjust each weight
(using gradient descent
and chain rule)



Multilayer Perceptron

Training with Backpropagation

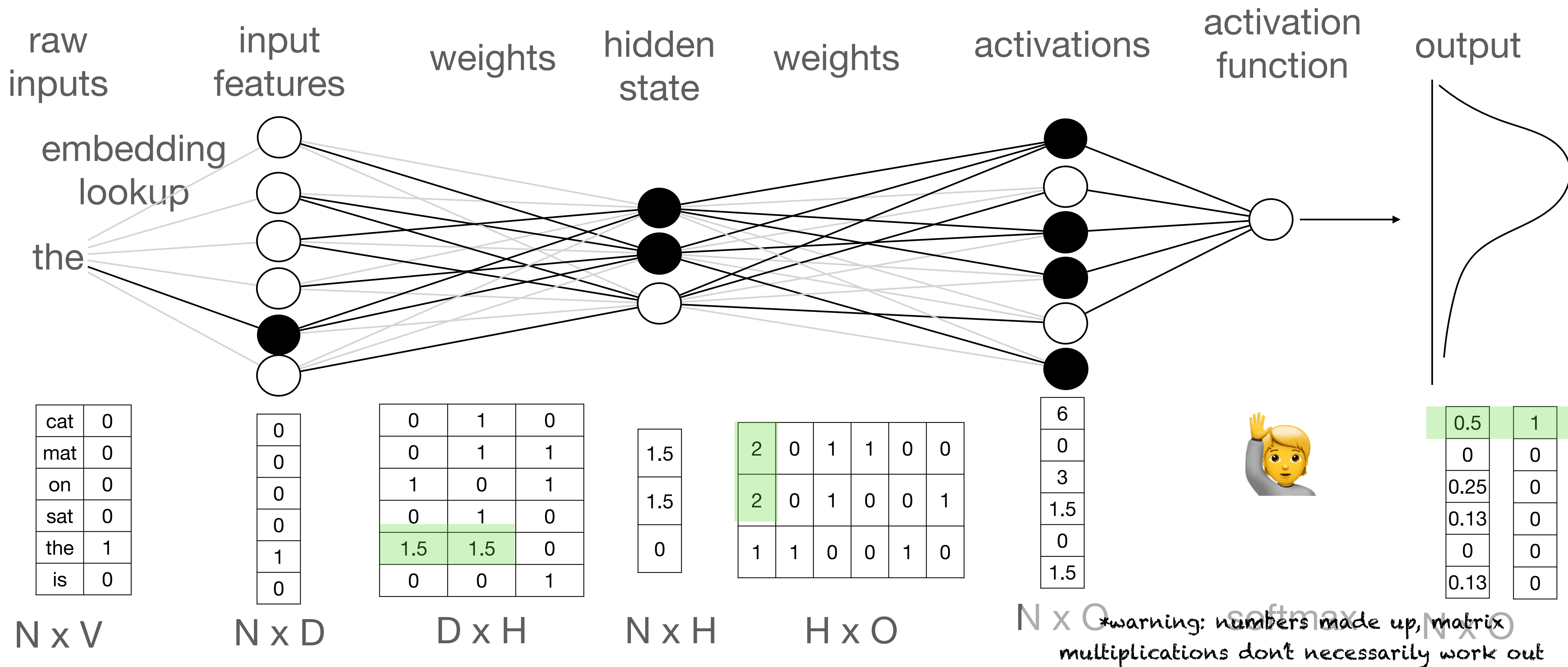
Adjust each weight
(using gradient descent
and chain rule)



Multilayer Perceptron

Training with Backpropagation

Adjust each weight
(using gradient descent
and chain rule)

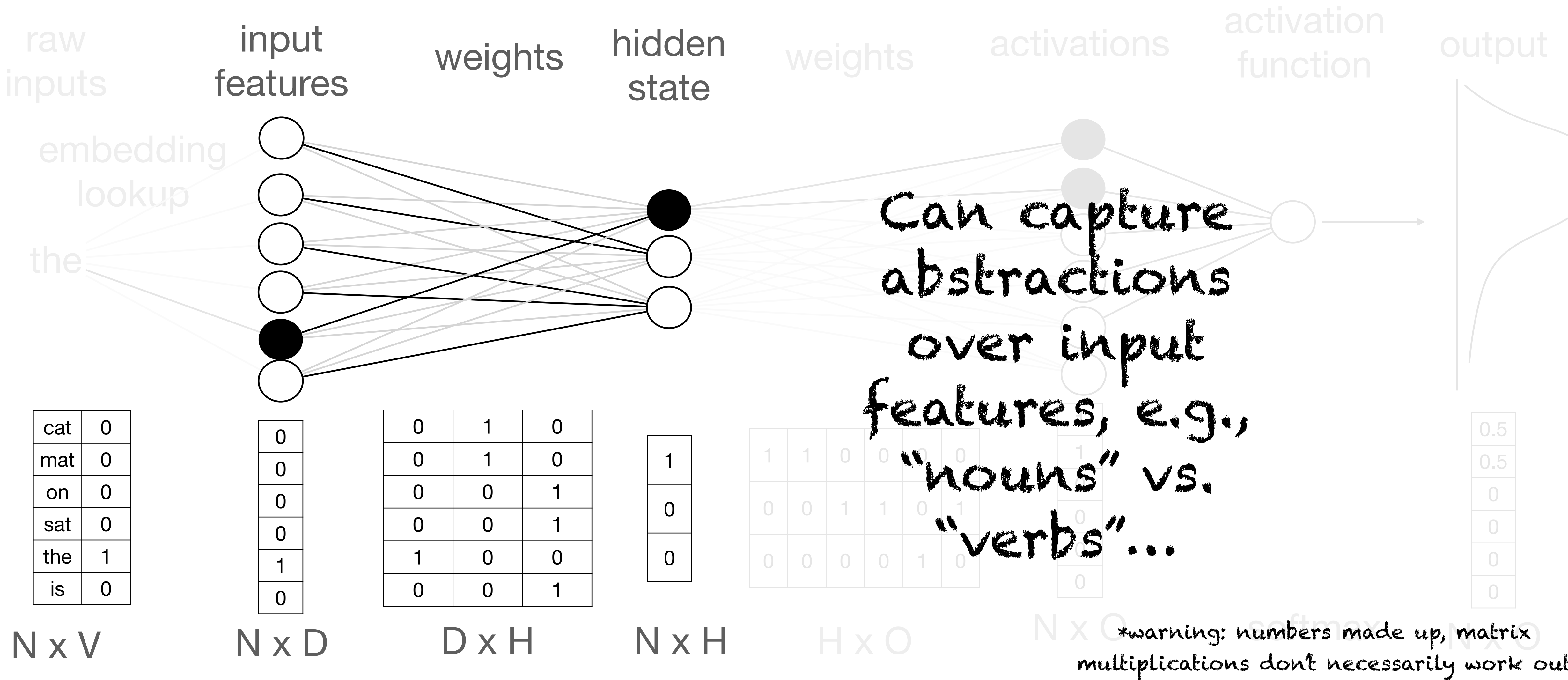


Topics

- More Followup on Word Embeddings from SVD
- Logistic Regression and Gradient Descent
- Multilayer Perceptrons
- **Word Embeddings from NNs**

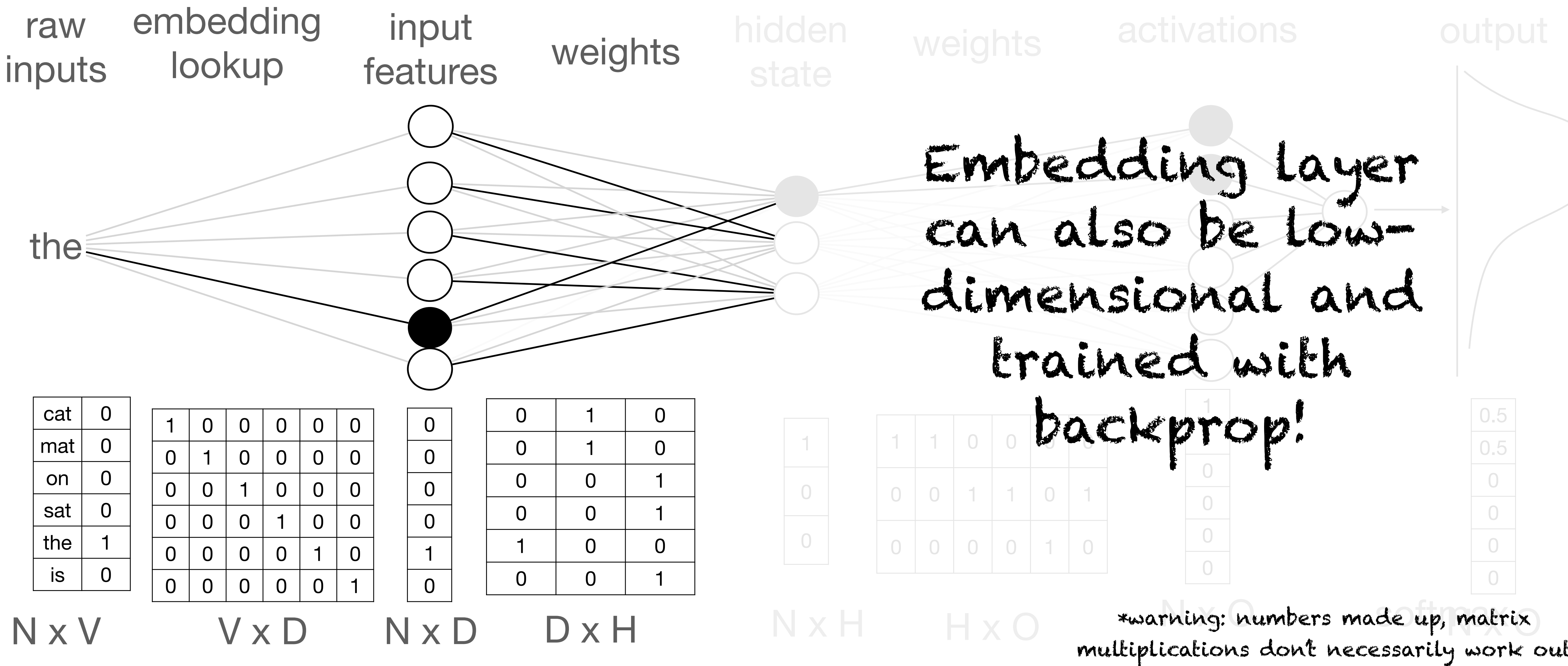
Word Embeddings from Neural Networks

Task: Predict the next word
Input: the
Expected: cat



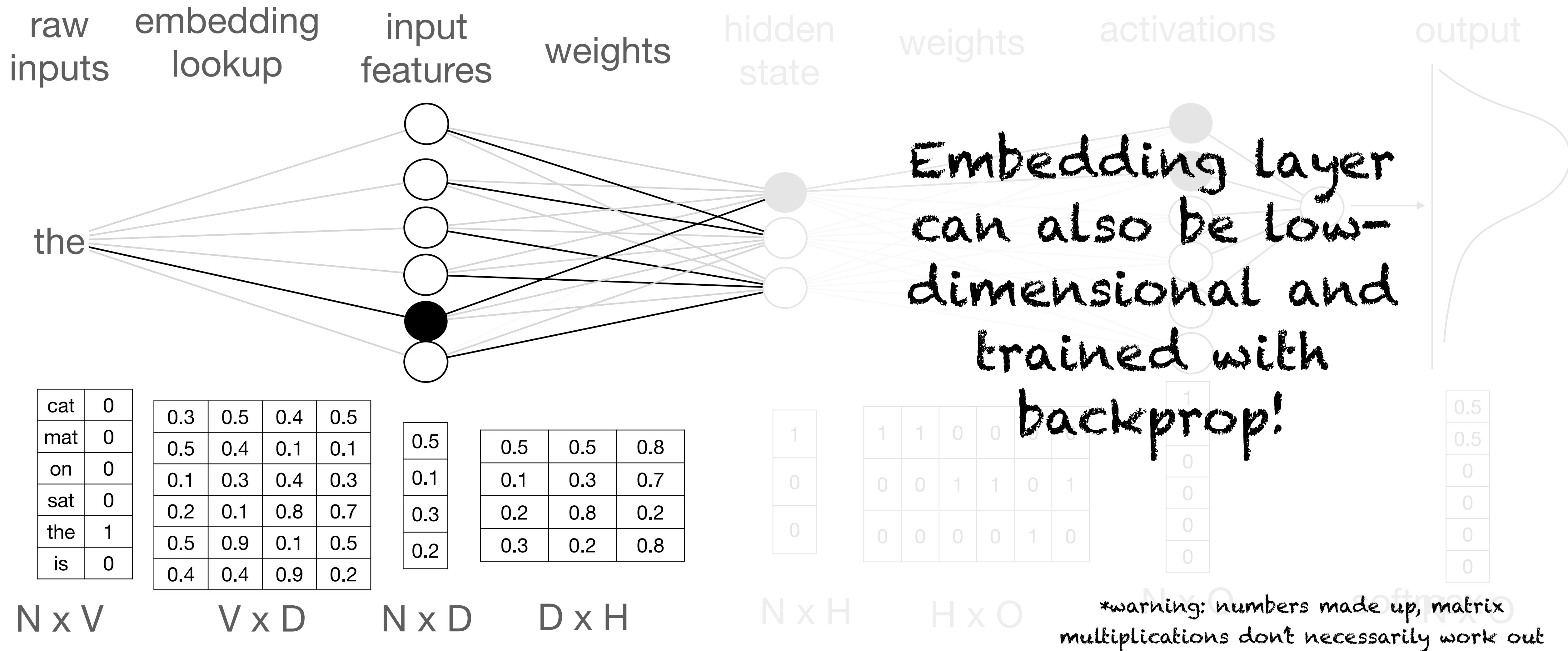
Word Embeddings from Neural Networks

Task: Predict the next word
Input: the
Expected: cat



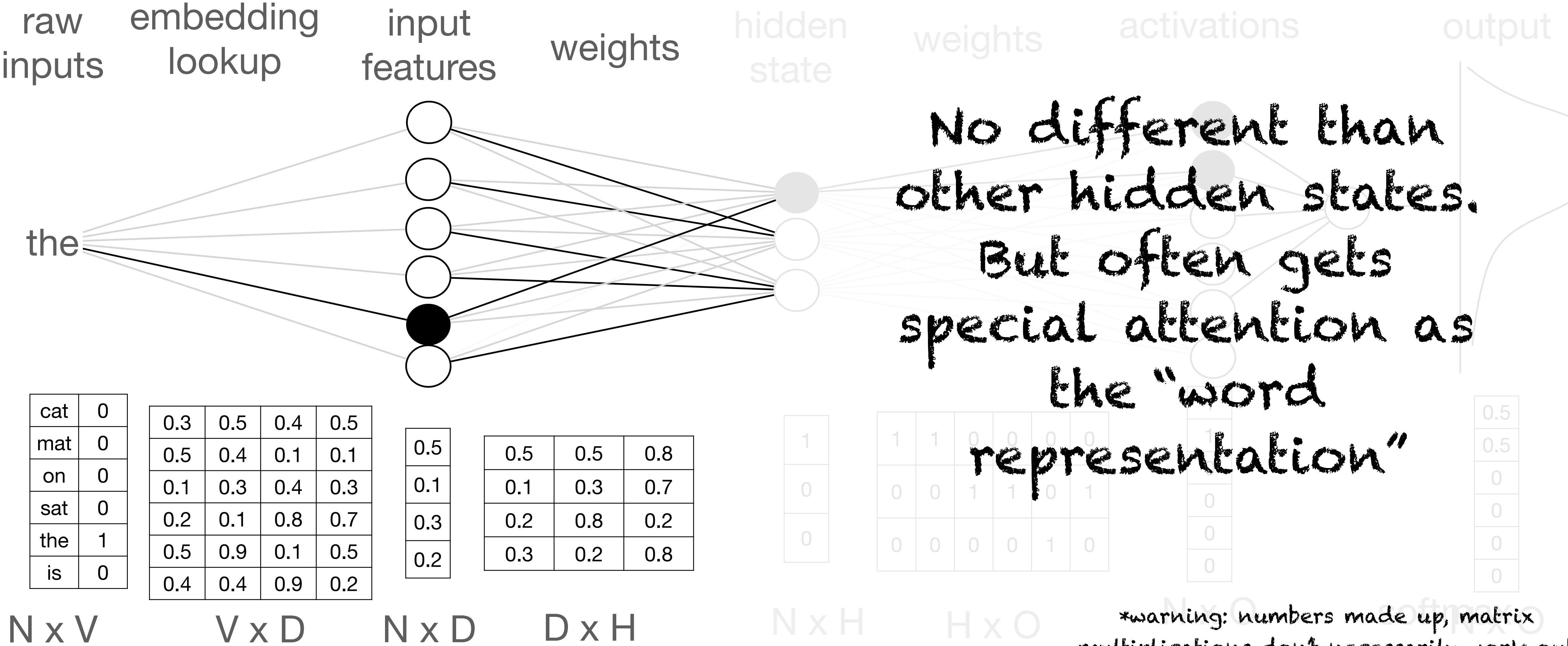
Word Embeddings from Neural Networks

Task: Predict the next word
Input: the
Expected: cat



Word Embeddings from Neural Networks

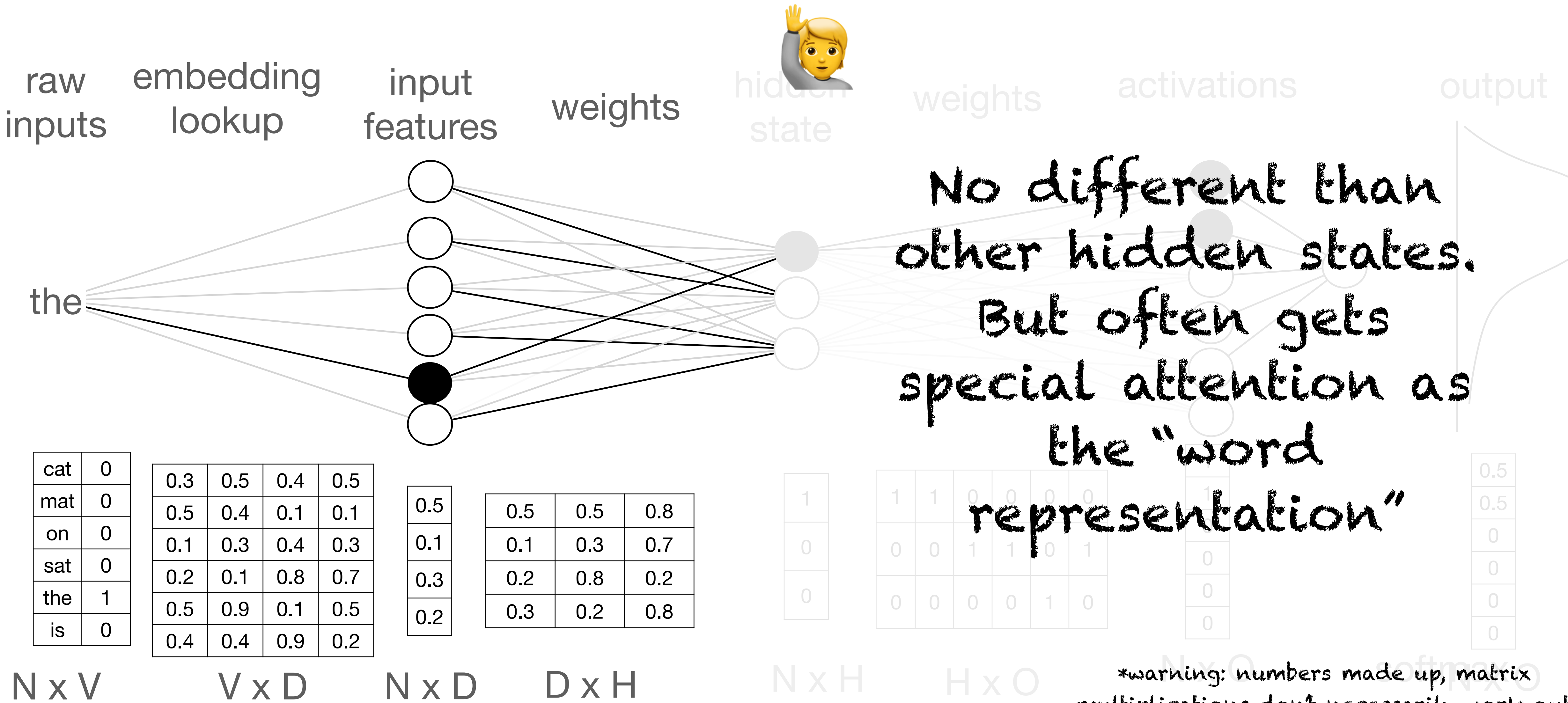
Task: Predict the next word
Input: the
Expected: cat



*warning: numbers made up, matrix multiplications don't necessarily work out

Word Embeddings from Neural Networks

Task: Predict the next word
Input: the
Expected: cat



*warning: numbers made up, matrix multiplications don't necessarily work out

Word Embeddings from Neural Networks

Efficient Estimation of Word Representations in Vector Space

Tomas Mikolov

Google Inc., Mountain View, CA
tmikolov@google.com

Kai Chen

Google Inc., Mountain View, CA
kaichen@google.com

Greg Corrado

Google Inc., Mountain View, CA
gcorrado@google.com

Jeffrey Dean

Google Inc., Mountain View, CA
jeff@google.com

Distributed Representations of Words and Phrases and their Compositionality

Tomas Mikolov

Google Inc.
Mountain View
mikolov@google.com

Ilya Sutskever

Google Inc.
Mountain View
ilyasu@google.com

Kai Chen

Google Inc.
Mountain View
kai@google.com

Greg Corrado

Google Inc.
Mountain View
gcorrado@google.com

Jeffrey Dean

Google Inc.
Mountain View
jeff@google.com

Word Embeddings from Neural Networks

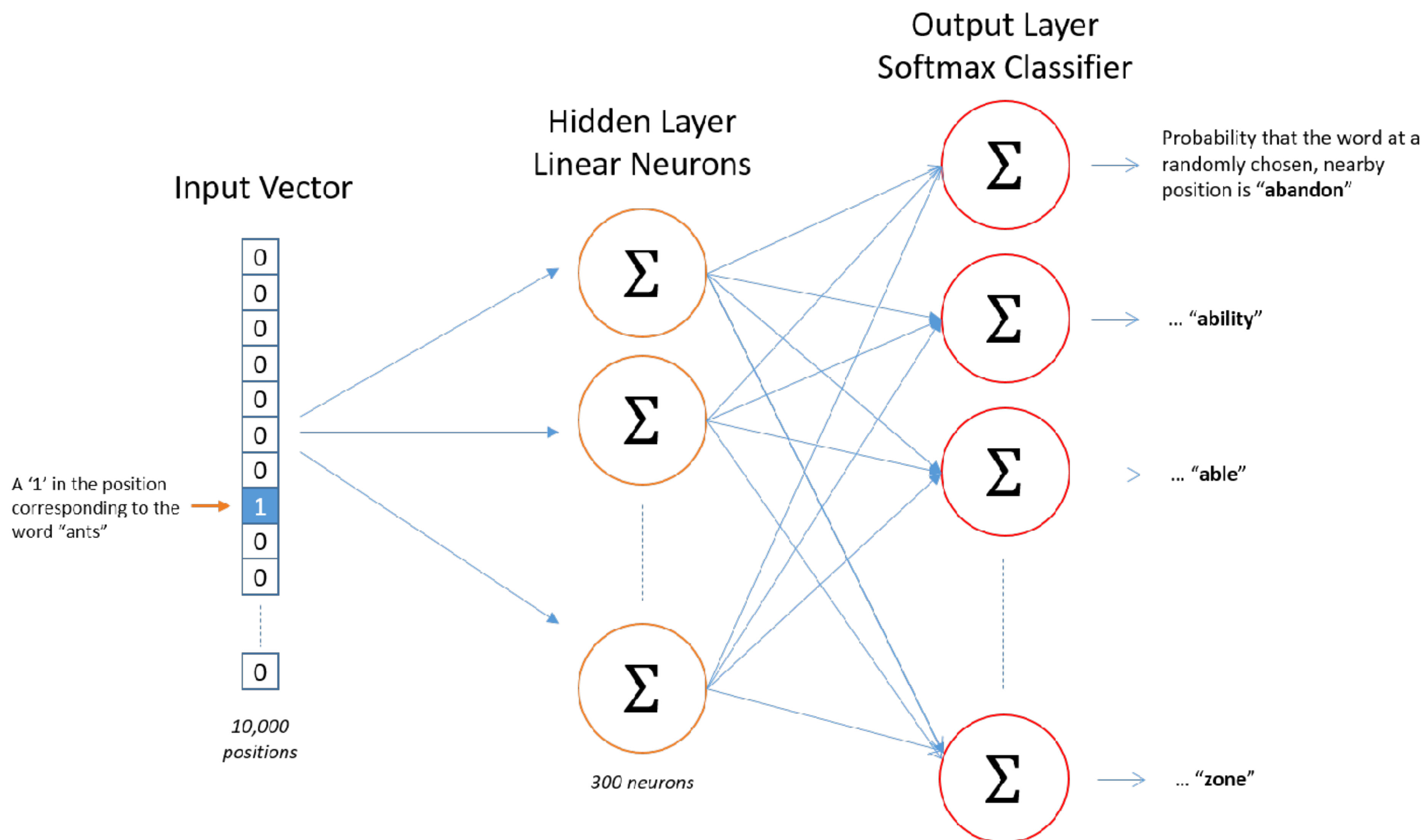
Distribution of words in w's context

Word w being represented

	cat	kitten	cute	adorable	gradients
cat	0	0	1	1	0
kitten	0	0	1	1	0
cute	1	1	0	0	0
adorable	1	1	1	0	0
gradients	0	0	1	1	1

Word Embeddings from Neural Networks

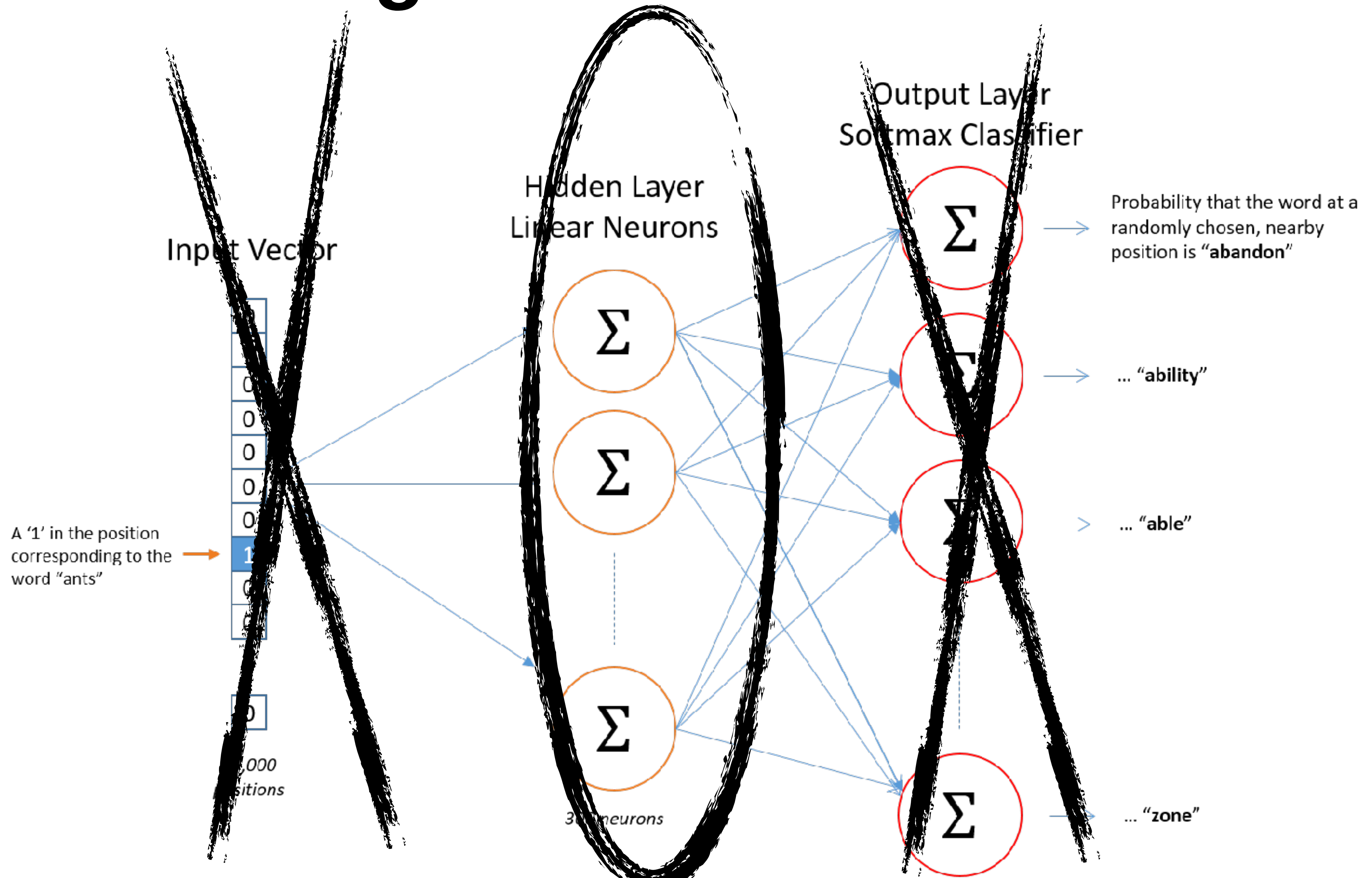
Word w being represented



Distribution of words in w 's context

Word Embeddings from Neural Networks

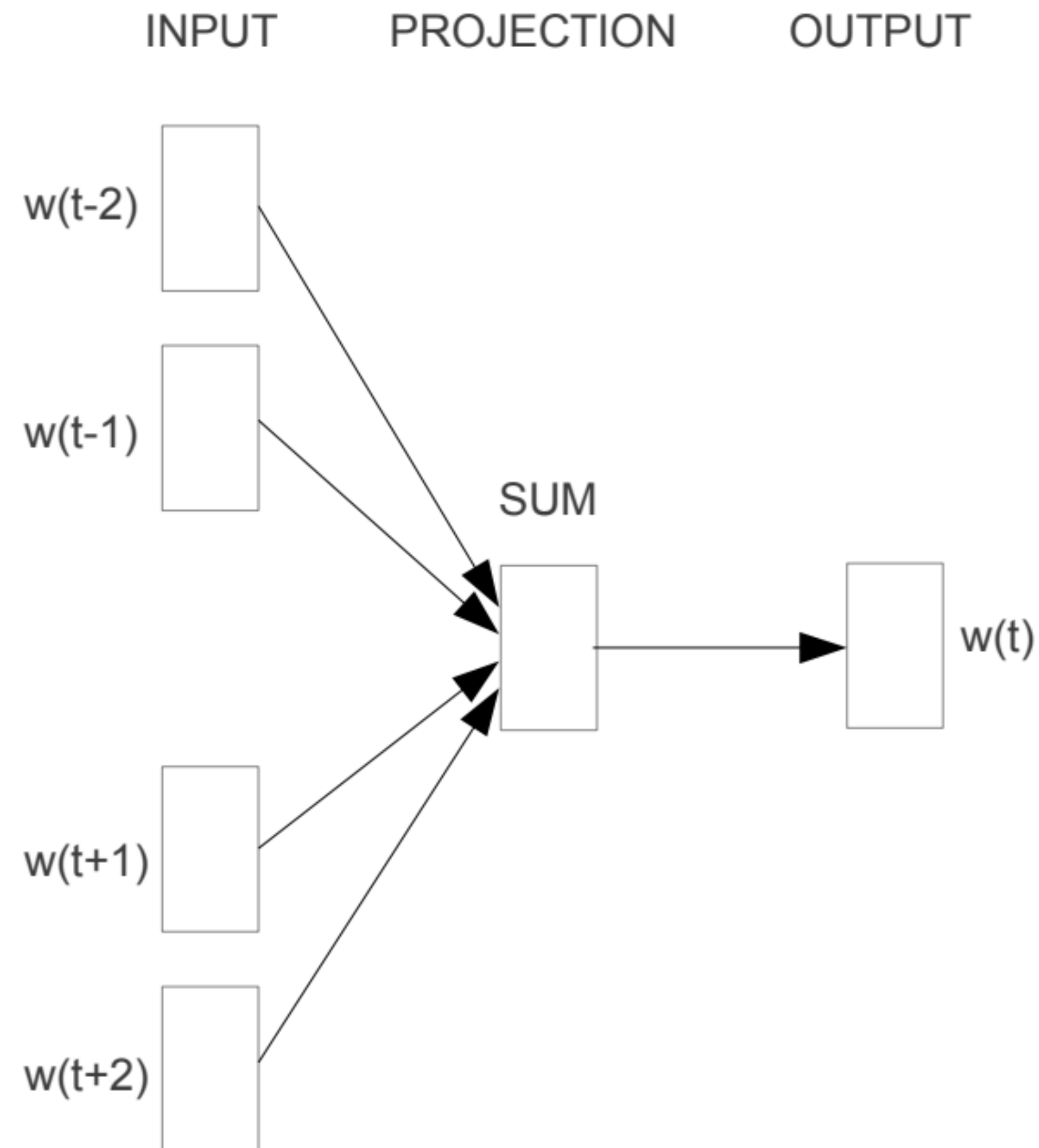
Word w being represented



Distribution of words in w 's context

Word Embeddings from Neural Networks

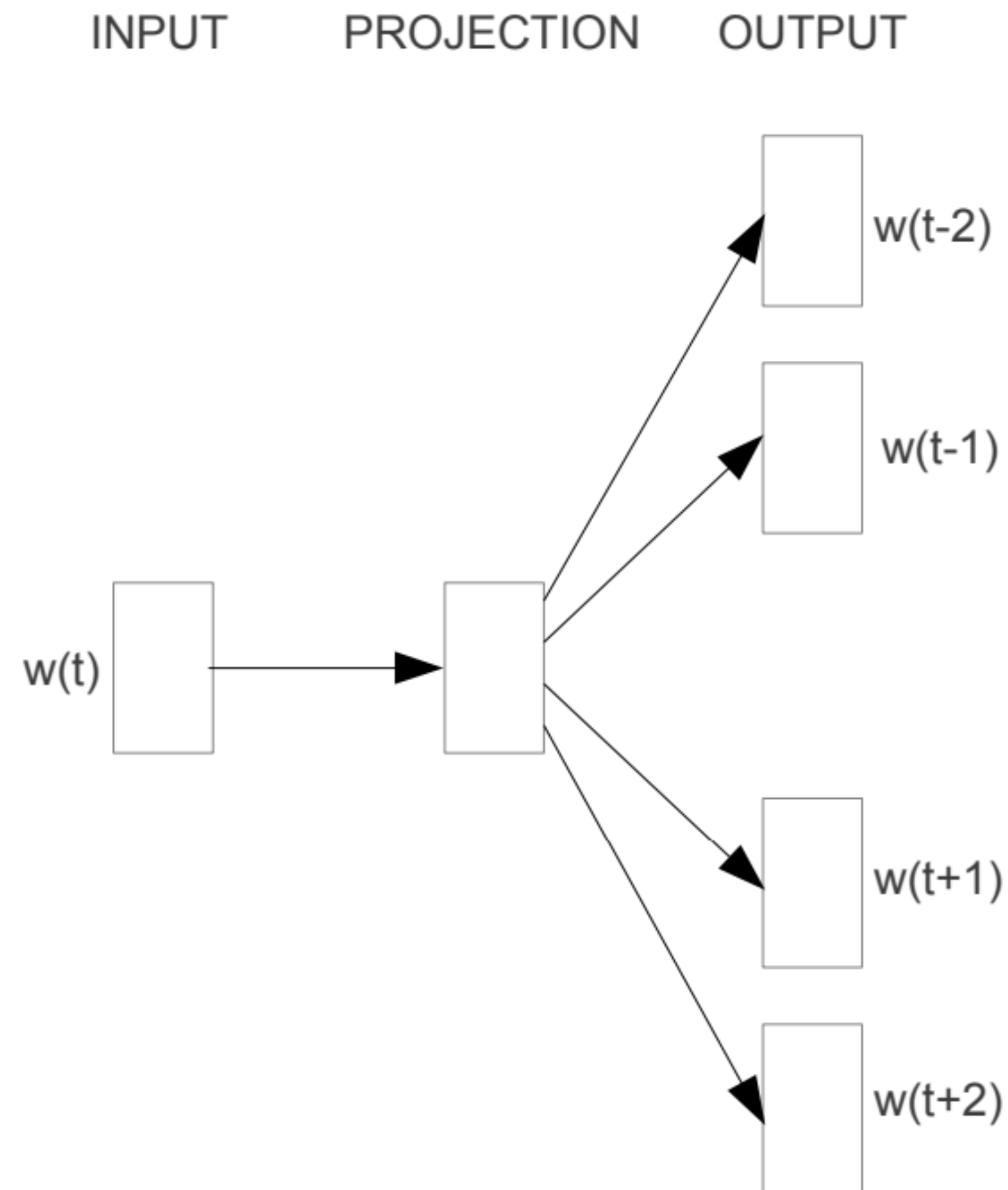
Continuous Bag of Words (CBOW)



Given context,
predict the word

Word Embeddings from Neural Networks

SkipGram



Given the word,
predict the
context

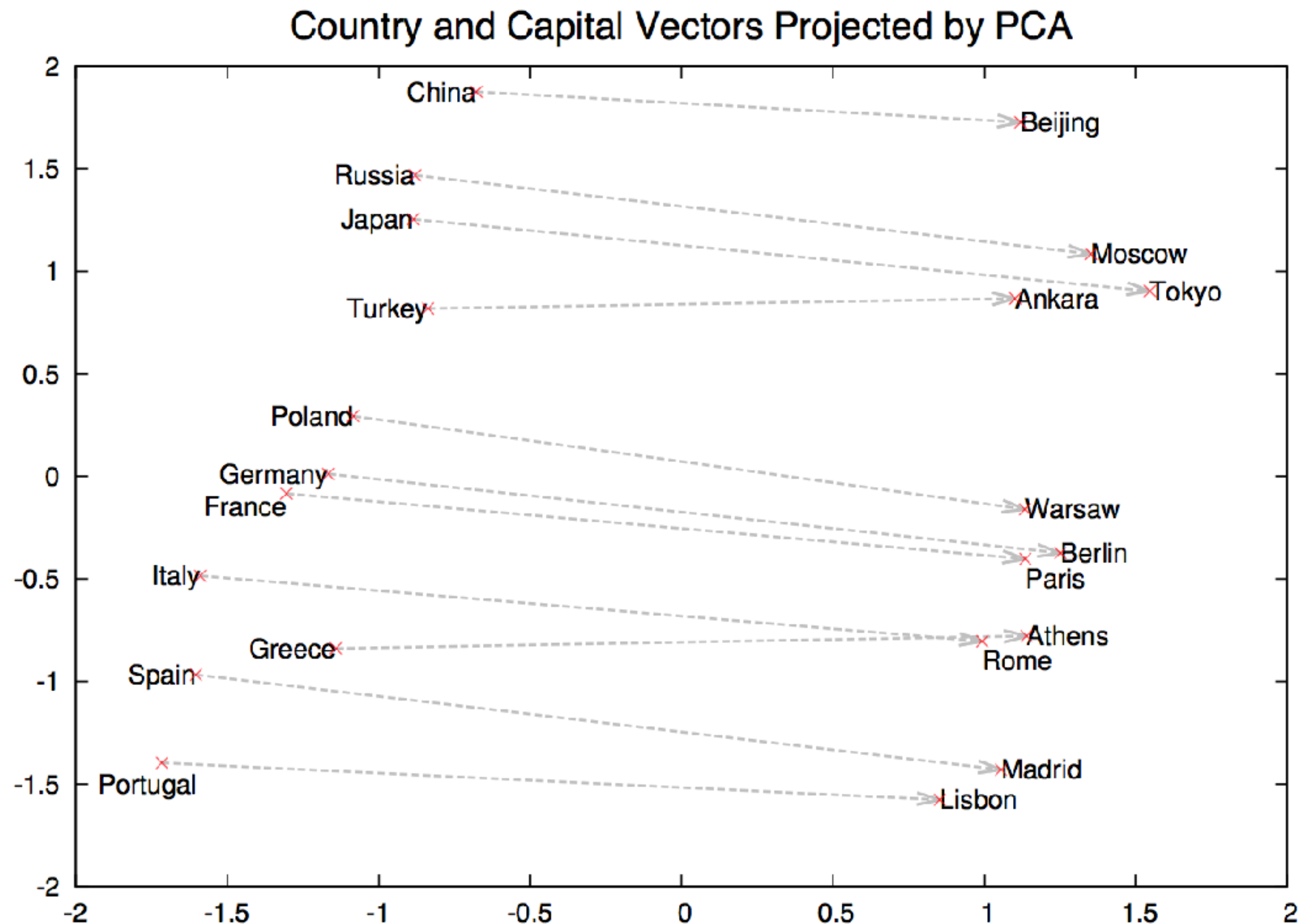
Pro

Eva



Pretrained Word Embeddings

Evaluations of word2vec embeddings



Pretrained Word Embeddings

Evaluations of word2vec embeddings

Table 8: *Examples of the word pair relationships, using the best word vectors from Table 4 (Skip-gram model trained on 783M words with 300 dimensionality).*

Relationship	Example 1	Example 2	Example 3
France - Paris	Italy: Rome	Japan: Tokyo	Florida: Tallahassee
big - bigger	small: larger	cold: colder	quick: quicker
Miami - Florida	Baltimore: Maryland	Dallas: Texas	Kona: Hawaii
Einstein - scientist	Messi: midfielder	Mozart: violinist	Picasso: painter
Sarkozy - France	Berlusconi: Italy	Merkel: Germany	Koizumi: Japan
copper - Cu	zinc: Zn	gold: Au	uranium: plutonium
Berlusconi - Silvio	Sarkozy: Nicolas	Putin: Medvedev	Obama: Barack
Microsoft - Windows	Google: Android	IBM: Linux	Apple: iPhone
Microsoft - Ballmer	Google: Yahoo	IBM: McNealy	Apple: Jobs
Japan - sushi	Germany: bratwurst	France: tapas	USA: pizza

Pretrained Word Embeddings

NNs vs. SVD

Pretrained Word Embeddings

NNs vs. SVD

- Same basic idea! Dimensionality reduction leads to good abstractions

Pretrained Word Embeddings

NNs vs. SVD

- Same basic idea! Dimensionality reduction leads to good abstractions
 - In fact, the two methods are provably equivalent in the simplest case

Neural Word Embedding as Implicit Matrix Factorization

Omer Levy

Department of Computer Science

Bar-Ilan University

omerlevy@gmail.com

Yoav Goldberg

Department of Computer Science

Bar-Ilan University

yoav.goldberg@gmail.com

Pretrained Word Embeddings

NNs vs. SVD

- Same basic idea! Dimensionality reduction leads to good abstractions
 - In fact, the two methods are provably equivalent in the simplest case
- But embeddings from NNs can become more powerful (and harder to interpret) as:

Neural Word Embedding as Implicit Matrix Factorization

Omer Levy

Department of Computer Science

Bar-Ilan University

omerlevy@gmail.com

Yoav Goldberg

Department of Computer Science

Bar-Ilan University

yoav.goldberg@gmail.com

Pretrained Word Embeddings

NNs vs. SVD

- Same basic idea! Dimensionality reduction leads to good abstractions
 - In fact, the two methods are provably equivalent in the simplest case
- But embeddings from NNs can become more powerful (and harder to interpret) as:
 - We add more layers

Neural Word Embedding as Implicit Matrix Factorization

Omer Levy
Department of Computer Science
Bar-Ilan University
omerlevy@gmail.com

Yoav Goldberg
Department of Computer Science
Bar-Ilan University
yoav.goldberg@gmail.com

Pretrained Word Embeddings

NNs vs. SVD

- Same basic idea! Dimensionality reduction leads to good abstractions
 - In fact, the two methods are provably equivalent in the simplest case
- But embeddings from NNs can become more powerful (and harder to interpret) as:
 - We add more layers
 - We add more non-linearity

Neural Word Embedding as Implicit Matrix Factorization

Omer Levy
Department of Computer Science
Bar-Ilan University
omerlevy@gmail.com

Yoav Goldberg
Department of Computer Science
Bar-Ilan University
yoav.goldberg@gmail.com

Pretrained Word Embeddings

NNs vs. SVD

- Same basic idea! Dimensionality reduction leads to good abstractions
 - In fact, the two methods are provably equivalent in the simplest case
- But embeddings from NNs can become more powerful (and harder to interpret) as:
 - We add more layers
 - We add more non-linearity
 - We invent more complex training objectives

Neural Word Embedding as Implicit Matrix Factorization

Omer Levy
Department of Computer Science
Bar-Ilan University
omerlevy@gmail.com

Yoav Goldberg
Department of Computer Science
Bar-Ilan University
yoav.goldberg@gmail.com

Pretrained Word Embeddings

NNs vs. SVD

- Same basic idea! Dimensionality reduction leads to good abstractions
 - In fact, the two methods are provably equivalent in the simplest case
- But embeddings from NNs can become more powerful (and harder to interpret) as:
 - We add more layers
 - We add more non-linearity
 - We invent more complex training objectives
 - More next lecture(s)!

Neural Word Embedding as Implicit Matrix Factorization

Omer Levy
Department of Computer Science
Bar-Ilan University
omerlevy@gmail.com

Yoav Goldberg
Department of Computer Science
Bar-Ilan University
yoav.goldberg@gmail.com