

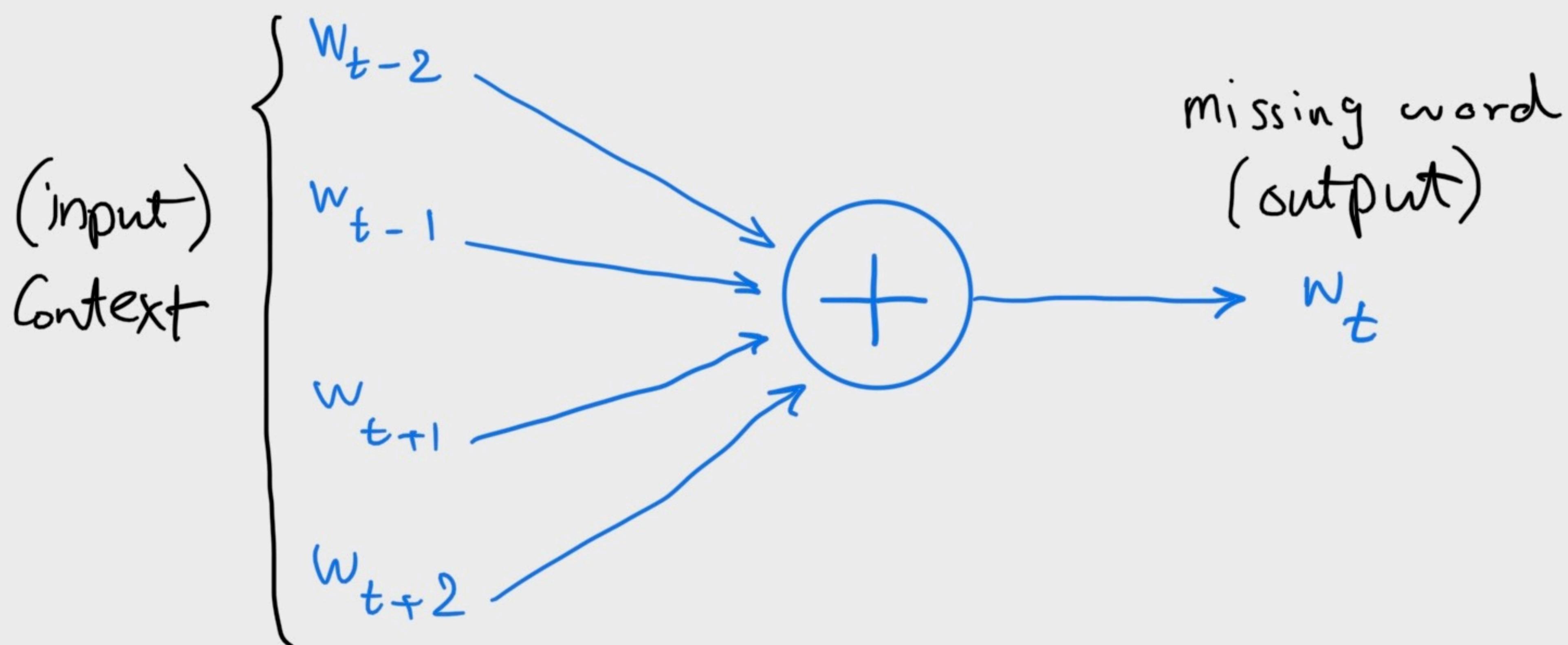
## Word2Vec

Any word that occurs in the same context,  
should have similar vector.

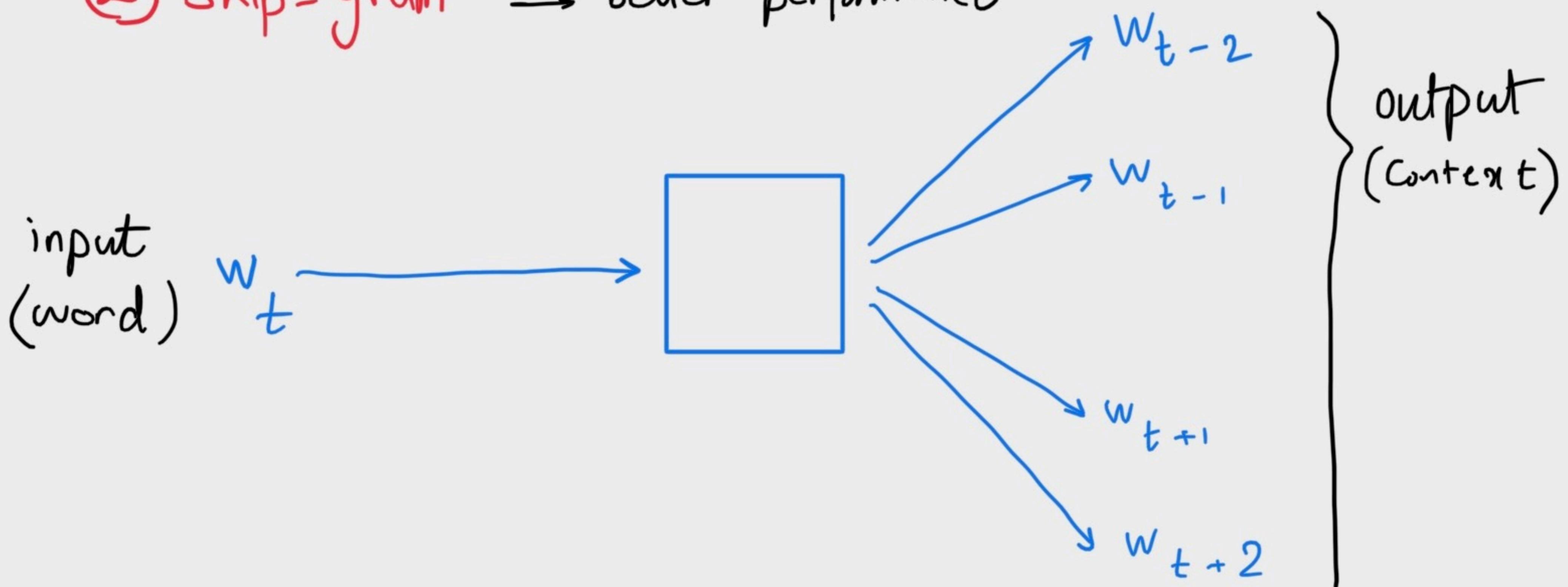
words before  
and after

- We can implement word2vec  
in 2 ways

### ① Continuous Bag-of-words (CBOW)



### ② Skip-gram → better performance



# Implementing Word2Vec

## ● Load the Data

with open('text.txt', 'r') as f:

```
text = f.read()
```

## ● Preprocess

```
import re  
from collections import Counter  
  
def preprocess(text):  
    text = text.lower().replace('.', '<PERIOD>')  
        .replace(':', '<COLON>'). ... .replace('...')  
    note the space
```

```
words = text.split()
```

```
c = Counter(words) → Counter for word frequency
```

```
ref_words = [word for word in words if c[word] > 5]
```

```
return ref_words
```

```
words = preprocess(text) → Create a big list of words
```

## • Create a Dictionary

```
def lookup_table(words):  
    c = Counter(words)  
    vocab = Sorted(c, k=c.get, reverse=True)  
    int2vocab = {i: word for i, word in enumerate(vocab)}  
    vocab2int = {word: i for i, word in int2vocab.items()}  
    return vocab2int, int2vocab
```

Sort by frequency

vocab2int, int2vocab = lookup\_table(words)

words\_idx = [vocab2int[w] for w in words]

## • Subsample

we don't want to weigh words that are too frequent, very much. this is basically noise

the, and, a, ...

### Mikolov Subsampling

$$P(w_i) = 1 - \sqrt{\frac{t}{f(w_i)}} \quad \begin{array}{l} \xrightarrow{\text{some threshold}} \\ \text{↓ for word } w_i \\ \text{↳ frequency of } w_i \text{ in dataset} \end{array}$$

$P$  is the probability of discarding  $w_i$ .

Example:  $P(0) = 1 - \sqrt{\frac{1 \times 10^{-5}}{1 \times 10^6 / 16 \times 10^6}} = 0.98735$

word "the"

we want to discard "the" most of the time; but still keeping enough of it to make an embedding.

```
from collections import Counter  
import random  
import numpy as np
```

$t = 1e-5$

$c = Counter(\\text{words\_idx})$

$text\_size = len(\\text{words})$

$freqs = \\{\\text{word} : \\text{count}/text\\_size \\text{ for } \\text{word}, \\text{count} \\text{ in } c.items()\\}$

$pdrop = \\{\\text{word} : 1 - np.sqrt(t/freqs[\\text{word}]) \\text{ for } \\text{word} \\text{ in } c\\}$

$train\\_words = [\\text{word} \\text{ for } \\text{word} \\text{ in } words\\_idx \\text{ if }$

$\\underbrace{random.random()}_{\\text{between } 0 \\& 1} < \\underbrace{(1 - pdrop[\\text{word}])}_{\\text{probability of keeping a word.}}$

$\\downarrow$   
discarding too frequent words from the text

### • Create batches of Data

```
def get_target(words, idx, wsize = 5):
```

$R = random.choice(range(1, wsize + 1))$  select a  $wsize$ ,

$frm = max(0, idx - R)$

randomly select a

$to = min(len(words), idx + R)$

number  $R$  from  $[1:wsize]$ ,

select  $R$  token in

future & past from  $words[idx]$

$return words[frm:idx] + words[idx + 1: to + 1]$

$words[idx]$

```

def batch(words, bsize, wsize=5):
    n = len(words) // bsize
    words = words[:n * bsize]

    for idx in range(0, len(words), bsize):
        x, y = [], []
        b = words[idx:idx+bsize]
        for i in range(len(b)):
            bx = b[i]
            by = get_target(b, i, wsize)
            x.extend([bx] * len(by))
            y.extend(by) → put each by on a
                           new row
        yield x, y

```

### • Define a Similarity Metric

```

def cos_Similarity(embed, vsize=16, vwindow=100,
                   device='CPU'):

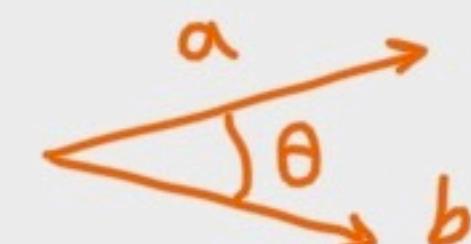
```

$\vec{v}_w \rightarrow \text{embed\_vecs} = \text{embed.weight}$

$\vec{v} \rightarrow \text{mags} = \text{embed\_vecs}.\text{pow}(2).\text{sum(dim}=1)$

$\cdot \text{sqrt}() \cdot \text{unsqueeze}(0)$

$$\begin{aligned} \text{Similarity} &= \cos \theta \\ &= \frac{\vec{a} \cdot \vec{b}}{\|\vec{a}\| \|\vec{b}\|} \end{aligned}$$



pick N  
words from examples = np.array(random.sample(range(vwindow),  
(0, window)  
and examples = np.append(examples,  
(1000, widow + 1000) random.sample(range(1000, 1000+vwindow),  
vsize // 2))

examples = torch.LongTensor(examples).to(device)

vvecs = embed(examples) embedding layer

Similarity → sims = torch.mm(vvecs, embed\_vecs.t()) / mags  
dotprod

return examples, sims

## • Skip Gram

`torch.nn.Embedding (num_embeddings, → # of rows in the  
embedding_dim, → # of columns in the  
...)`

↑ aka vocab size  
↑ lookup matrix

```
import torch  
from torch import nn  
import torch.optim as optim
```

```
class SkipGram(nn.Module) :
```

```
    def __init__(self, n_vocab, n_embed) :  
        super().__init__()
```

one row for  
each entry in vocab

```
        self.embedding = nn.Embedding(n_vocab, n_embed)
```

This FC  
layer takes  $\leftarrow$  self.output = nn.Linear(n\_embed, n\_vocab)  
the embedding self.norm = nn.LogSoftmax(dim=1)  
dim as input and output size of vocab length.

This is because the output is a series of word scores.

```
def forward(self, x) :  
    x = self.embedding(x)  
    Scores = self.output(x)  
    log_ps = self.norm(Scores)
```

return log-ps

## • Train the Model

```
device = 'cuda' if torch.cuda.is_available() else 'cpu'
```

embed\_dim=300

```
model = SkipGram(len(vocab2int), embed_dim).to(device)
```

```
criterion = nn.NLLLoss()
```

```
optimizer = optim.Adam(model.parameters(), lr=0.003)
```

Step = 0

epoch = 5

```
for e in range(epoch):
```

for input, target in batch(train\_words, 512):

Step += 1

```
input, target = torch.LongTensor(input),
```

`torch.LongTensor(target)`

```
input, target = input.to(device),
```

target · to (device)

`bg_ps = model(input)`

```
loss = criterion(log_ps, targets)
```

optimizer. zero\_grad()

loss.backward()

optimizer.step()

# doing some validation

if step % 20 == 0:

examples, sims = cos\_similarity(model.embedding,  
device=device)

- close = sims.topk(6)

examples, close = examples.to('cpu'),  
close.to('cpu')

for i, j in enumerate(examples):

so, the closest ← close\_words = [int2vocab[idx.item()]]  
words to int2vocab[j.item()] for idx in close[i][1:]  
would be: close\_words

## ● Visualize

import matplotlib.pyplot as plt

from sklearn.manifold import TSNE → t-distributed  
stochastic

neighbor embeddings

embeds = model.embedding.weight.to('cpu').data.numpy()

viz = 600

tsne = TSNE()

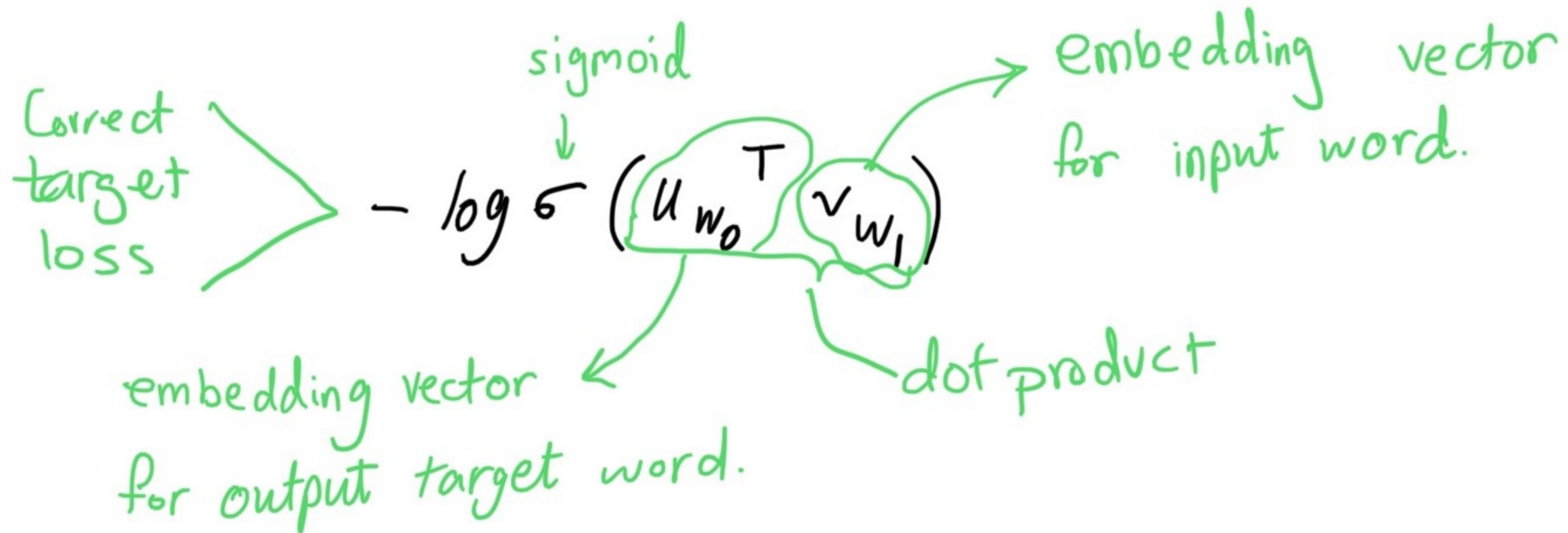
embed\_tsne = tsne.fit\_transform(embeds[:viz, :])

## ● Negative Sampling.

The implementation above can be very slow. We can use negative sampling to speed things up.

① do not take softmax over all the words. instead, we can add an embedding layer for the output.

② use a different loss function that only cares about true examples and a small subset of noise examples.



$\sum_i^N \mathbb{E}_{w_i \sim p_n(w)}$   
 sum over  
 all words  $w_i$   
 that are drawn from a noise distribution  
 ↓  
 Vocabulary of words that are  
 not in the content of input word  
 (i.e., irrelevant words)

Putting things together :

$$\text{loss} = -\log \sigma(u_{w_0}^\top v_{w_I}) - \sum_i^N \mathbb{E}_{w_i \sim p_n(w)} \log \sigma(-u_{w_i}^\top v_{w_I})$$

↓  
 prob  
 ↓  
 Correct target

↑  
 prob  
 ↑  
 Noisy target

SkipGram with Negative Sampling :

class SkipGramNeg(nn.Module) :

```
def __init__(self, n_vocab, n_embed, noise_dist=None):
    super().__init__()
```

`self.n_vocab = n_vocab`

`self.n_embed = n_embed`

`self.noise_dist = noise_dist`

embed layer  
for input

→ self.in\_embed = nn.Embedding(n\_vocab, n\_embed)

embedding layer for output → self.out\_embed = nn.Embedding(n\_vocab, n\_embed)

Initialization using uniform dist {  
self.in\_embed.weight.data.uniform\_(-1, 1)  
self.out\_embed.weight.data.uniform\_(-1, 1)}

def forward\_input (self, input\_words):  
return self.in\_embed (input\_words)

def forward\_output (self, output\_words):  
return self.out\_embed (output\_words)

def forward\_noise (self, bsize, n\_samples):  
} if self.noise\_dist is None:  
noise\_dist = torch.ones (self.n\_vocab)  
else:  
noise\_dist = self.noise\_dist

Set noise dist

Sampling words } noise\_words = torch.multinomial (noise\_dist,  
bsize \* n\_samples,  
replacement = True)

```

device = 'cuda' if model.out_embed.weight.is_cuda
else 'cpu'

noise_words = noise_words.to(device)

noise_vecs = self.out_embed(noise_words)
    • view(bsize, n-samples,
          self.n-embed)

```

return noise\_vecs

● We then need to define custom loss class

```
class NegativeSamplingLoss(nn.Module):
```

```
def __init__(self):
```

```
super().__init__()
```

```
def forward(self, input_vecs, output_vecs, noise_vecs):
```

bsize, embed\_size = input\_vecs.shape

*a batch of  
Column vectors* → input\_vecs = input\_vecs.view(bsize, embed\_size, 1)

*a batch of  
row vectors* → output\_vecs = output\_vecs.view(bsize, 1, embed\_size)

First  
Term of  
Loss Function

out\_loss = torch.bmm(output\_vecs, input\_vecs)
 • sigmoid() • log() • squeeze()

so that no empty dims  
are left in the output

Second  
Term of  
Loss Function

noise\_loss = torch.bmm (noise\_vecs.neg(), input\_vecs)  
• sigmoid () • log () • squeeze () • sum (1)

return - (out\_loss + noise\_loss).mean()

- There are some changes to the training loop

word\_freqs = np.array(sorted(Freqs.values()), reverse=True)

unigram\_dist = word\_freqs / word\_freqs.sum()

noise\_dist = torch.from\_numpy((unigram\_dist \*\* .75 /  
np.sum(unigram\_dist \*\* .75))

embed\_dim = 300

model = SkipGram Neg (len(vocab2int), embed\_dim,  
noise\_dist = noise\_dist).to(device)

criterion = Negative Sampling Loss

optimizer = optim.Adam (model.parameters(), lr=0.003)

the changes in the inner training loop:

Instead of just calling `model()`, we call the forward functions.

`input_vecs = model.forward_input(input)`

`output_vecs = model.forward_output(target)`

`noise_vecs = model.forward_noise(inputs.shape[0], 5)`

`loss = Criterion(input_vecs, output_vecs, noise_vecs)`

rest of the code is similar to before.