

Revisiting Sentiment Analysis from Neural Network Perspective

Jiaxi Tang @ Simon Fraser University
jiaxit@sfu.ca

Agenda

- Part 1: Sentiment Analysis from past (5min)
 - Representation: Bag-of-words model
 - Method: KNN
 - Issues
- Part 2: Neural network(15min)
 - Representation: Word Embedding
 - Method: CNN
- Part 3: Practice(15min)
 - Code flow
 - Forward / Backward propagation
- Part 4: FAQ(15min)

Representation: Bag-of-words Model [1]

(1) John likes to watch movies. Mary likes movies too.

(2) John also likes to watch football games.

original sentence



(1) [1, 2, 1, 1, 2, 0, 0, 0, 1, 1]

(2) [1, 1, 1, 1, 0, 1, 1, 1, 0, 0]

bag-of-words representation

["John",
"likes",
"to",
"watch",
"movies",
"also",
"football",
"games",
"Mary",
"too"]

vocabulary

Representation: N-gram

(1) John likes to watch movies. Mary likes movies too.

(2) John also likes to watch football games.

Conceptually, we can view bag-of-word model as a special case of the n-gram model, with $n=1$

["John",	["John likes",
"likes",	"likes to",
"to",	"to watch",
"watch",	"watch
"movies",	movies",
"also",	"Mary likes",
"football",	"likes movies",
"games",	"movies too",]
"Mary",	
"too"]	

bigram feature

unigram feature

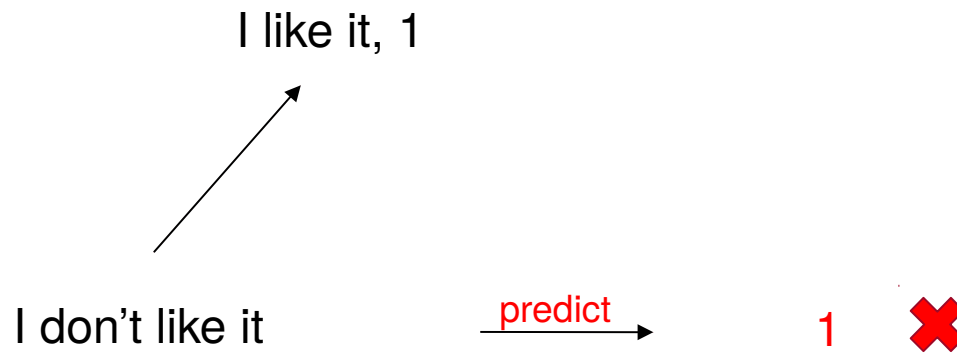
Method: K-Nearest Neighbor

Data:	Representation	Label:
(1) It is fantastic	[1,1,1,0,0,0,0,0]	1
(2) It is ultimately unsatisfying	[1,1,0,1,1,0,0,0]	0
(3) A mediocre film	[0,0,0,0,0,1,1,1]	0

It is ultimately mediocre film $\xrightarrow{\text{predict}}$ 0

[1,1,0,1,0,0,1,0]

Issue:

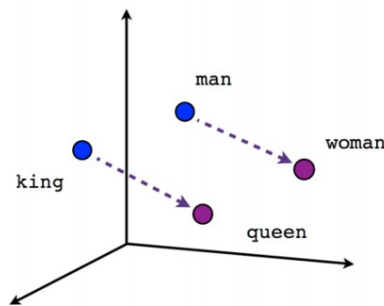


sparsity will lead to terrible performance!

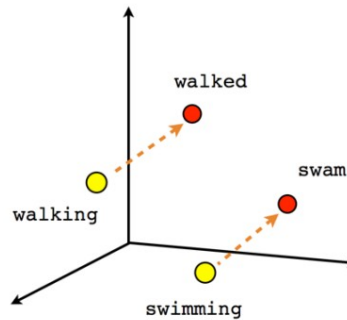
Representation: Word Embedding

Each word is associated with a vector(embedding) of latent values, e.g.

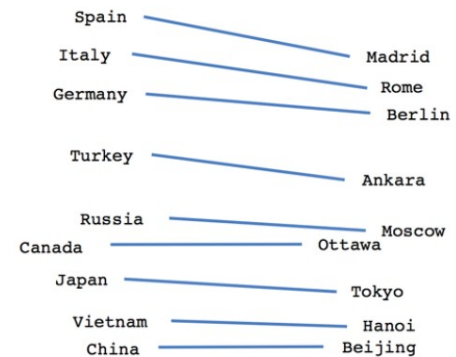
“like” = [0.7, 0.1, -0.5, 0.4, 1.2]



Male-Female



Verb tense



Country-Capital

Bag-of-words vs. Word Embedding

Bag-of-words

Each value has explicit meaning

Very sparse

Usually don't have
contextual information

Word embedding

Values are latent, don't know the meaning

Dense real values

Have contextual information

Representation: Word Embedding

two ways to get word embedding:

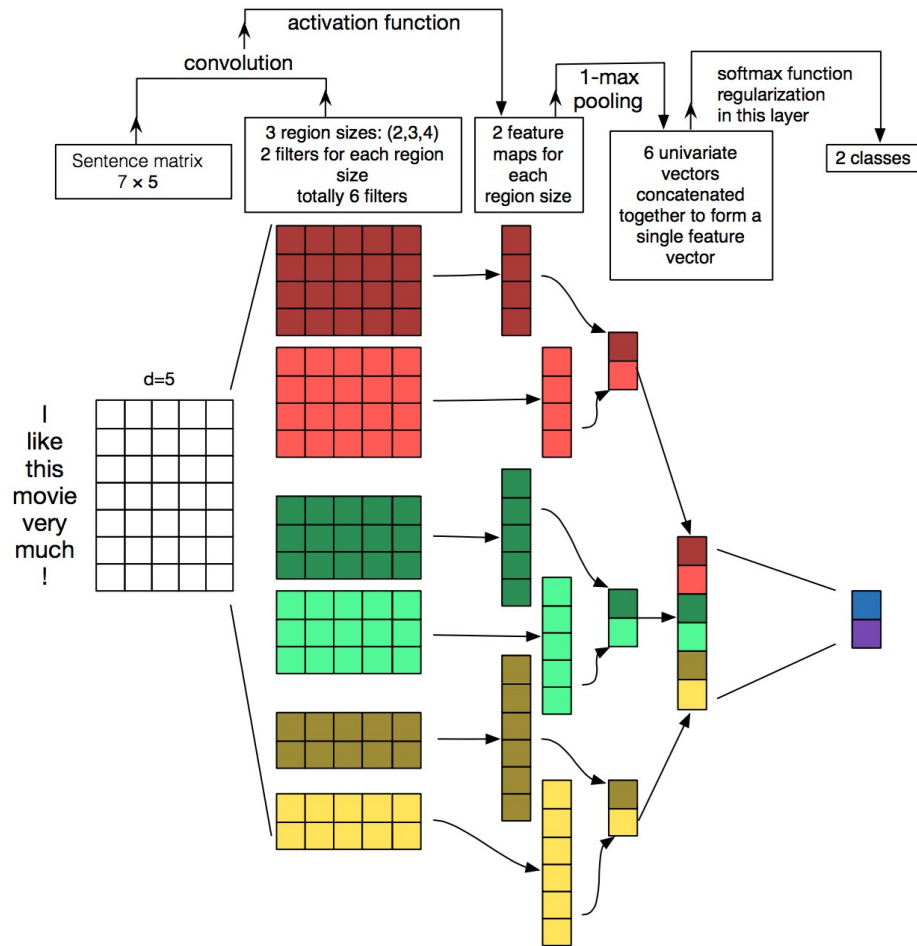
- (1) Initialize word embedding with random value, train them on the given dataset, and learn their optimal value.
- (2) Use existing word embedding trained from other datasets.

Google word2vec: <https://code.google.com/archive/p/word2vec/>

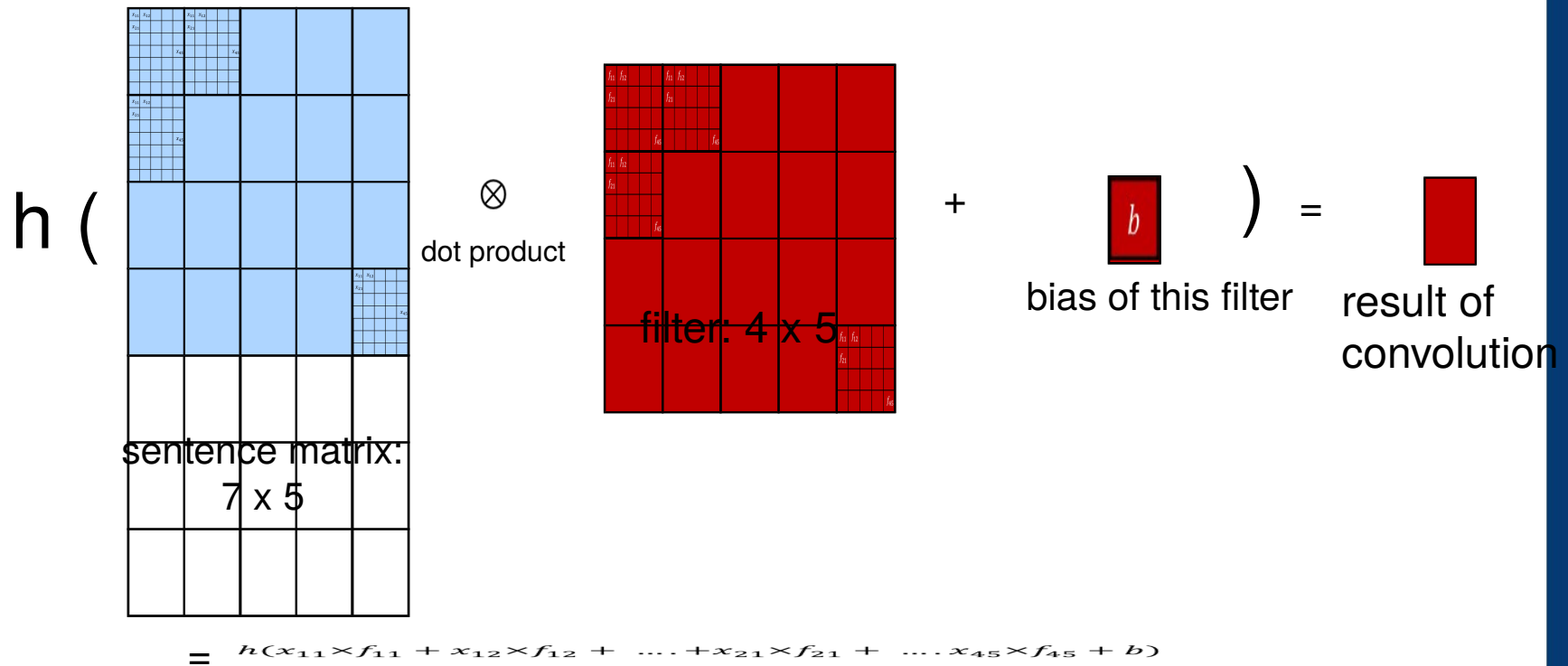
GloVe: <https://github.com/stanfordnlp/GloVe>

Try (1) first, then use (2)

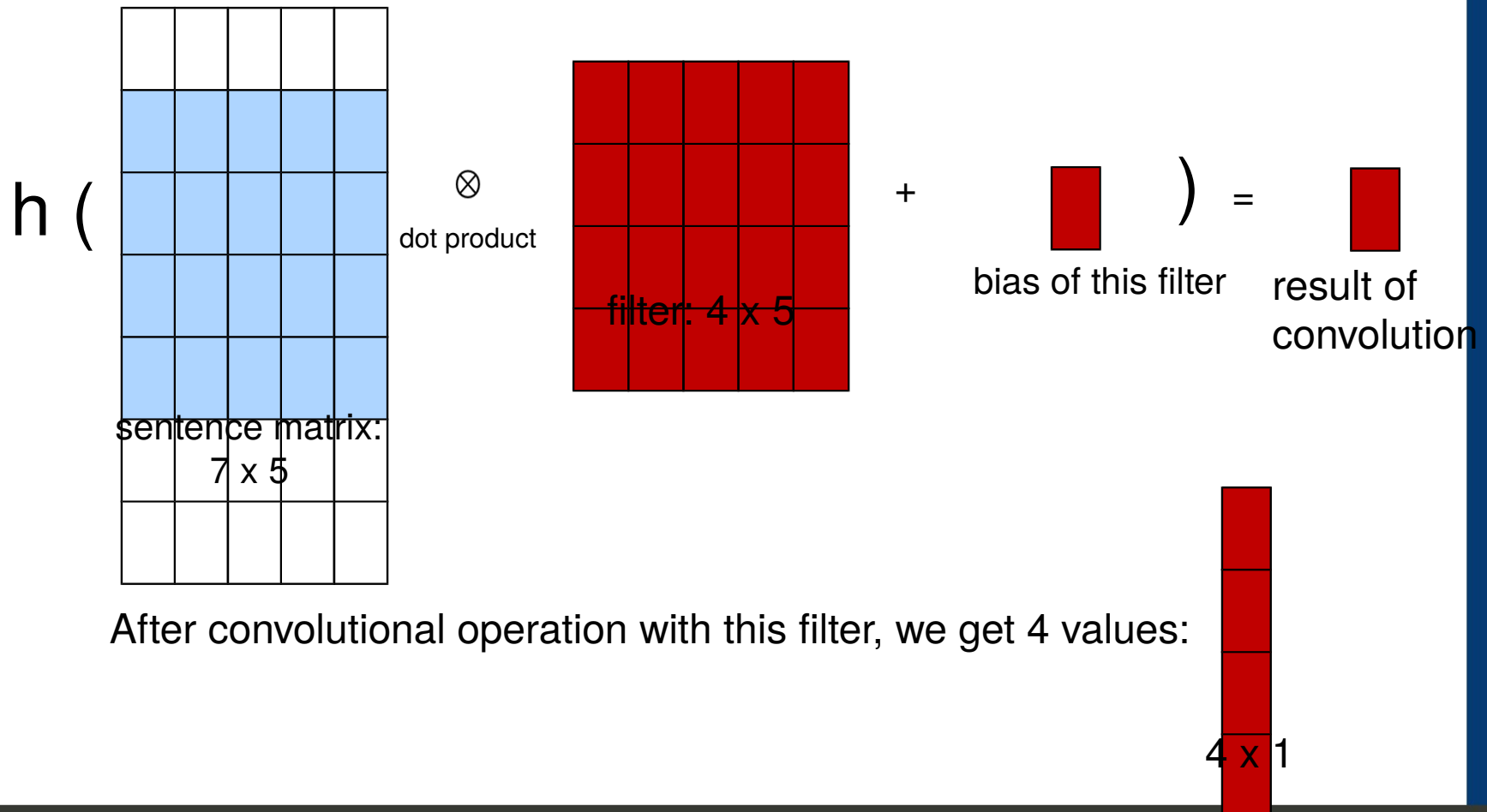
Method: CNN model



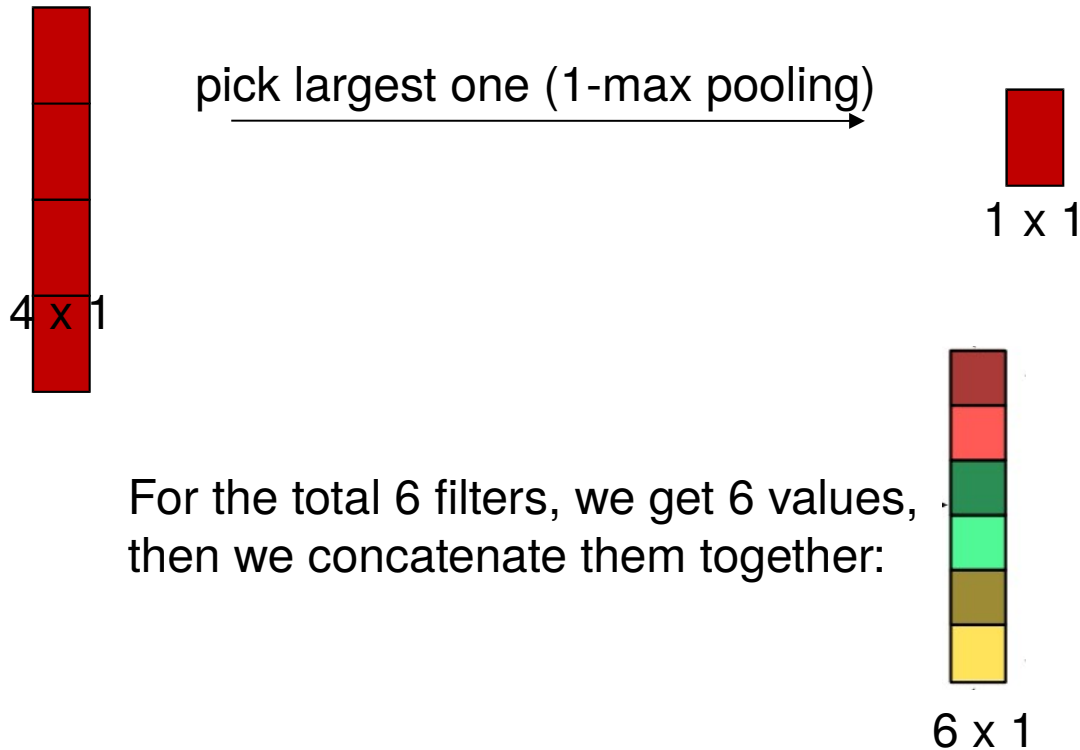
Method: CNN model – convolution



Method: CNN model – convolution



Method: CNN model - pooling



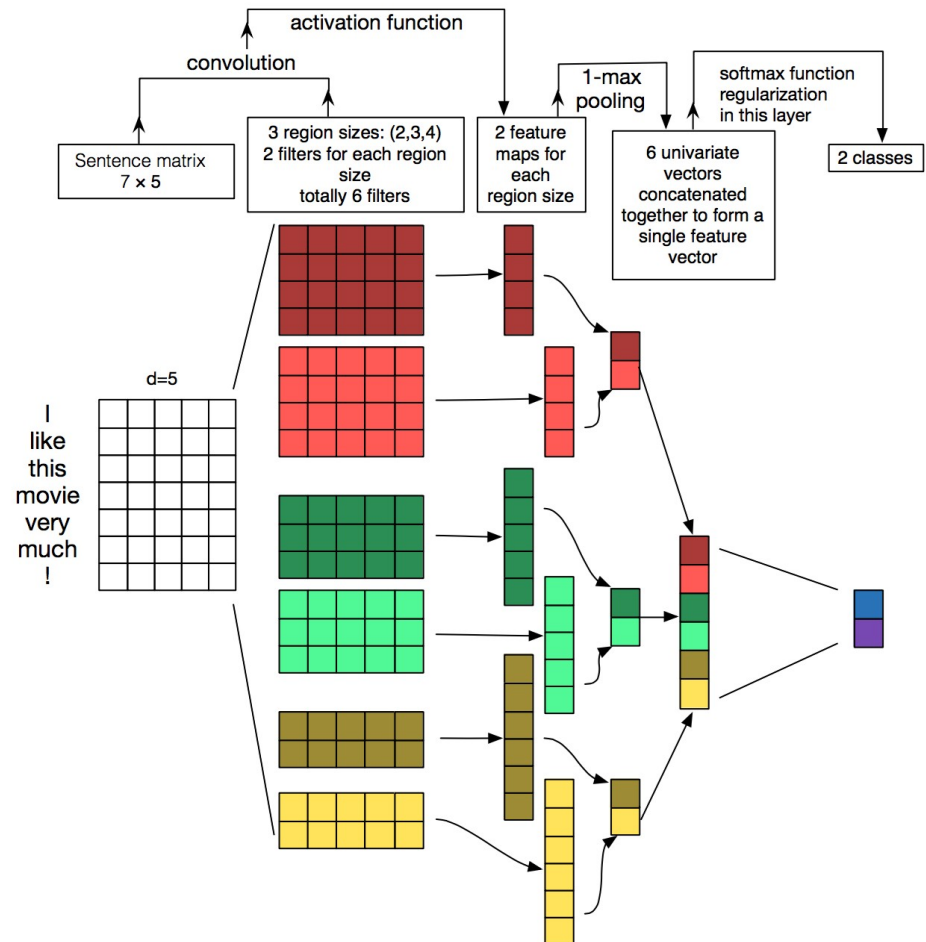
Method: CNN model

a filter with size = n

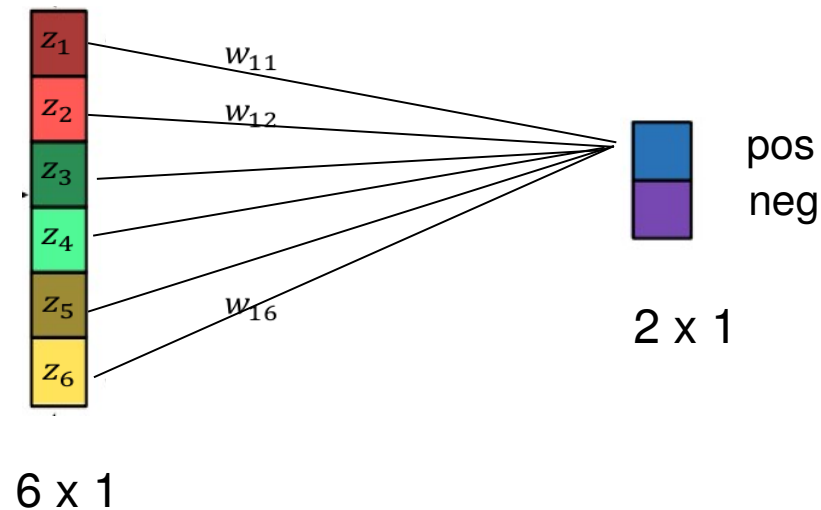
detect a specific n -gram feature

feature vector

embedding of this sentence



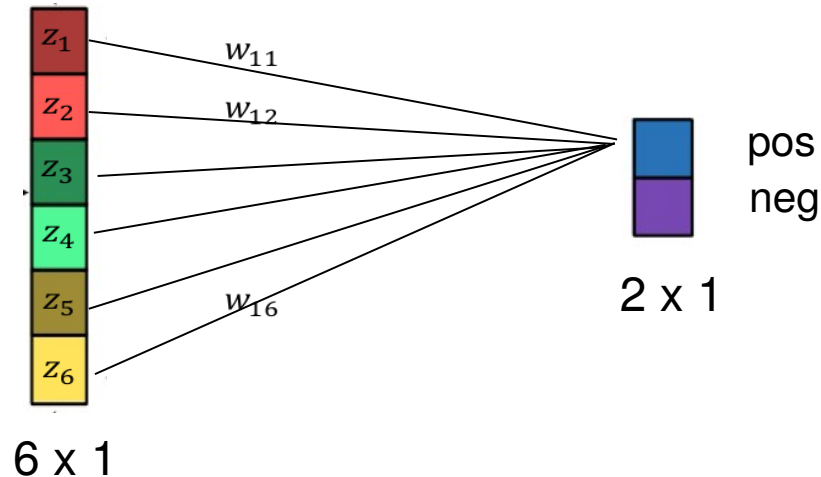
Method: CNN model – output layer, softmax loss



$$pos = z_1 \times w_{11} + z_2 \times w_{12} + \dots + z_6 \times w_{16} + b_{pos}$$

$$neg = z_1 \times w_{21} + z_2 \times w_{22} + \dots + z_6 \times w_{26} + b_{neg}$$

Method: CNN model – softmax loss



$$L_i = -\log\left(\frac{e^{s_{y_i}}}{\sum_j e^{s_j}}\right)$$

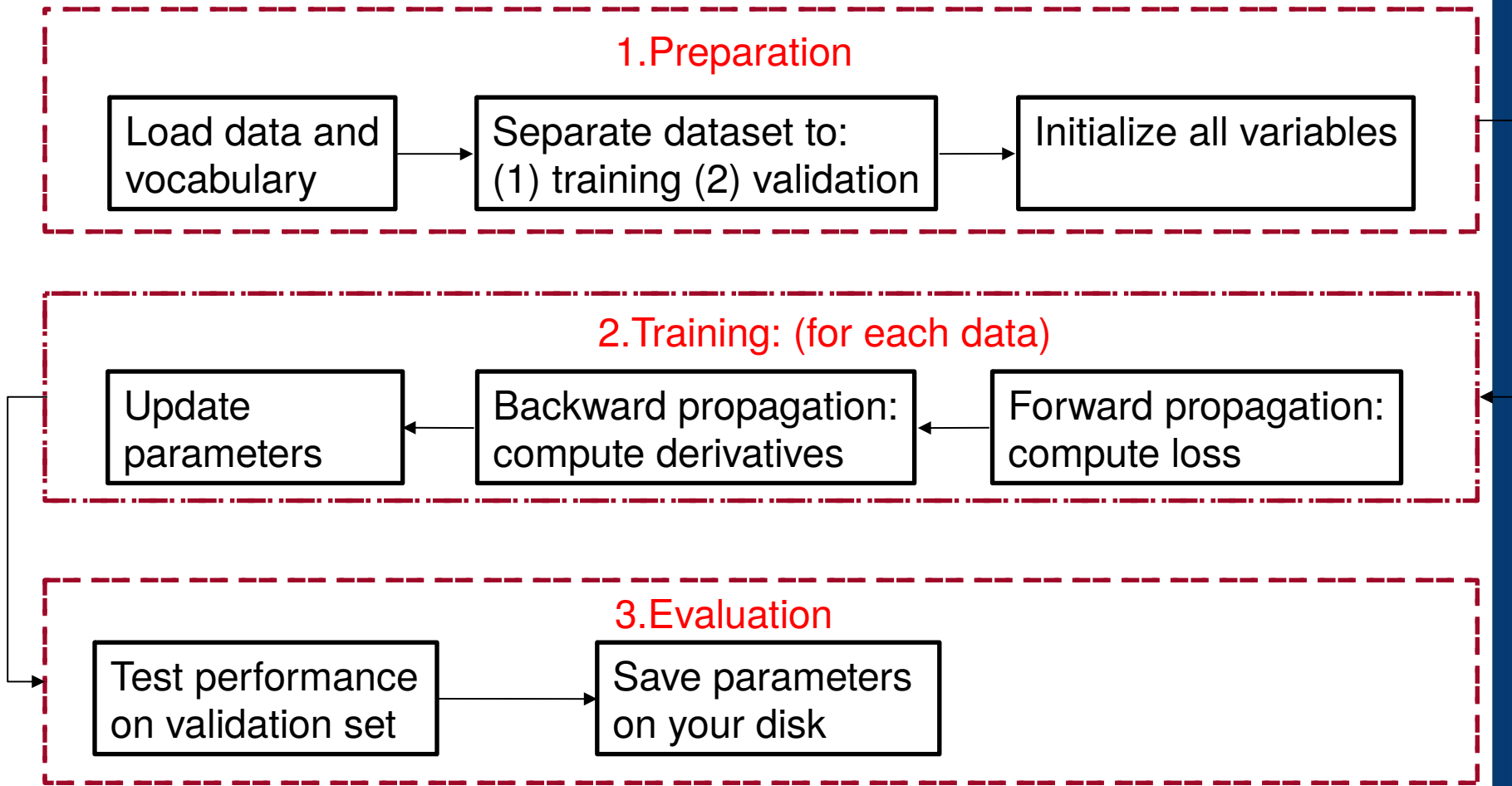
scores = unnormalized log probabilities of the classes

pos = -0.2 (ground truth)
neg = 1.3

exp → 0.8187
3.6693 → normalize → 0.1824
0.8176
probabilities

$$L_i = -\log(0.1824) = 1.7016$$

Practice: Code flow



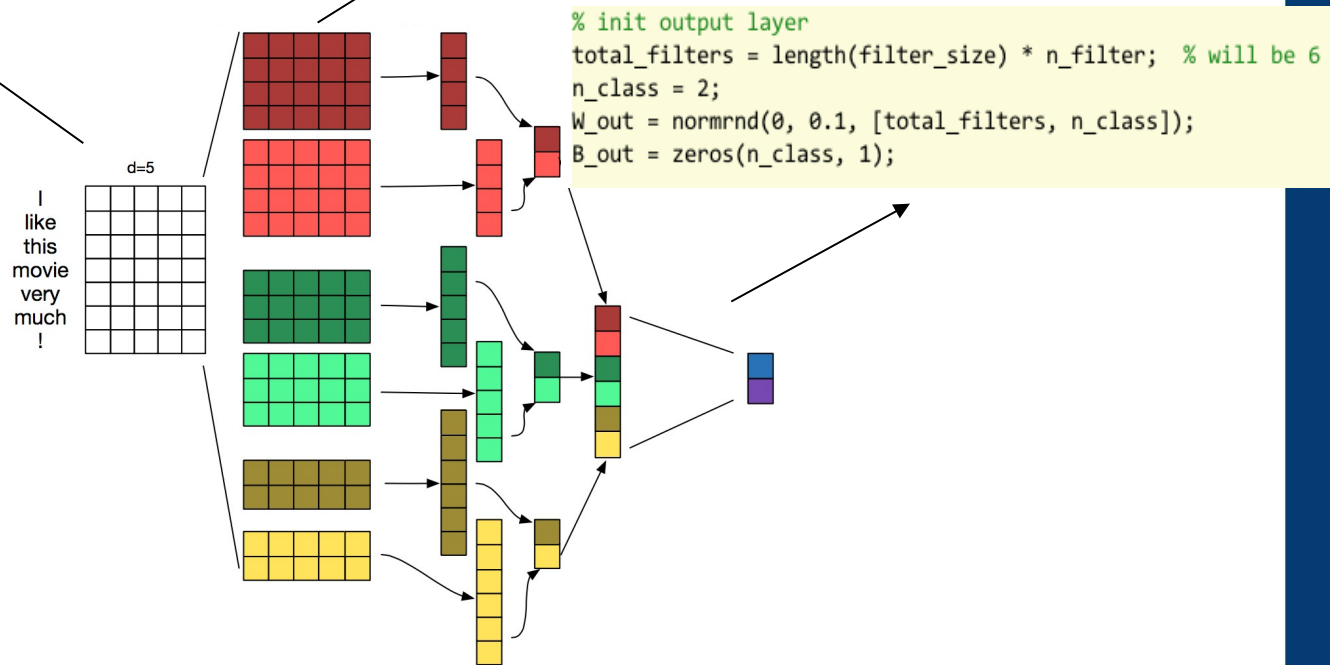
Practice: Initialization

```
% init embedding
d = 5;
total_words = length(wordMap);
% random sample from normal distribution
% with mean=0, variance=0.1
T = normrnd(0, 0.1, [total_words, d]);
```

```
% init filters
filter_size = [2,3,4];
n_filter = 2;

W_conv = cell(length(filter_size), 1);
B_conv = cell(length(filter_size), 1);

for i = 1: length(filter_size)
    % get filter size
    f = filter_size(i);
    % init W with: FW x FH x FC x K
    W_conv{i} = normrnd(0, 0.1, [f, d, 1, n_filter]);
    B_conv{i} = zeros(n_filter, 1);
end
```



Practice: Forward propagation

```
% get sentence matrix
% words_indexes = [wordMap('i'), wordMap('like'),
% .., wordMap('!')]
X = T(word_indexes, :);
```

```
pool_res = cell(1, length(filter_size));
cache = cell(2, length(filter_size));
for i = 1: length(filter_size)
    % convolutional operation
    conv = vl_nnconv(X, W_conv{i}, B_conv{i});

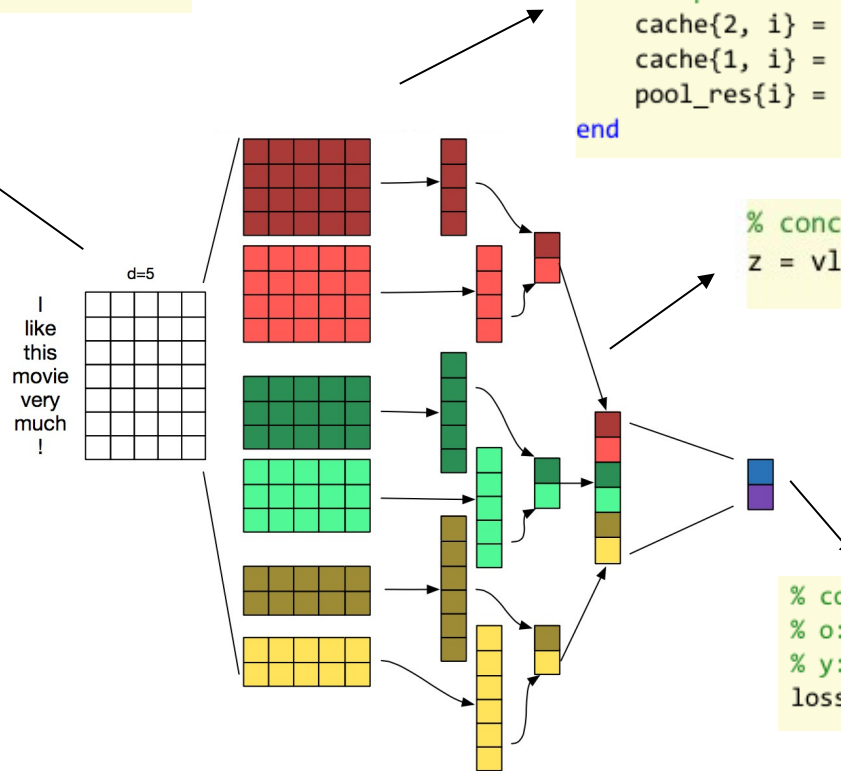
    % apply activation function: relu
    relu = vl_nnrelu(conv);

    % 1-max pooling operation
    sizes = size(conv);
    pool = vl_nnpool(relu, [sizes(1), 1]);

    % important: keep these values for back-prop
    cache{2, i} = relu;
    cache{1, i} = conv;
    pool_res{i} = pool;
end
```

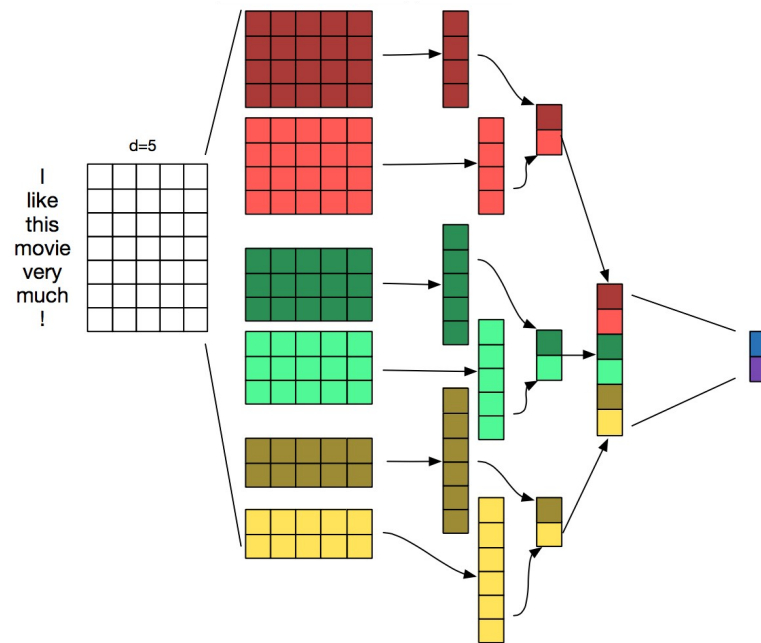
```
% concatenate
z = vl_nnconcat(pool_res, 3);
```

```
% compute loss
% o: value of output layer
% y: ground truth label (1 or 2)
loss = vl_nnloss(o, y);
```



Practice: Backward propagation

Part2 of <http://www.robots.ox.ac.uk/~vgg/practicals/cnn/>



Challenge:

1. split by ratio ✉ cross-validation
2. use existing word2vec to initialize word embedding
3. regularization (l2 norm, dropout)
4. SGD ✉ RMSProp or Adam
5. model ensembles
5. use mini-batch to accelerate
6. ...

FAQ

Thanks