

2.1 Hierarchical Softmax

Q. 가장 유용한 approximation은 무엇인가?
e.g) 왜 PLM가둠

A computationally efficient approximation of the full softmax is the hierarchical softmax. In the context of neural network language models, it was first introduced by Morin and Bengio [12]. The main advantage is that instead of evaluating W output nodes in the neural network to obtain the probability distribution, it is needed to evaluate only about $\log_2(W)$ nodes.

The hierarchical softmax uses a binary tree representation of the output layer with the W words as its leaves and, for each node, explicitly represents the relative probabilities of its child nodes. These define a random walk that assigns probabilities to words.

More precisely, each word w can be reached by an appropriate path from the root of the tree. Let $n(w, j)$ be the j -th node on the path from the root to w , and let $L(w)$ be the length of this path, so $n(w, 1) = \text{root}$ and $n(w, L(w)) = w$. In addition, for any inner node n , let $\text{ch}(n)$ be an arbitrary fixed child of n and let $\llbracket x \rrbracket$ be 1 if x is true and -1 otherwise. Then the hierarchical softmax defines $p(w_O | w_I)$ as follows:

$$p(w | w_I) = \prod_{j=1}^{L(w)-1} \sigma \left(\llbracket n(w, j+1) = \text{ch}(n(w, j)) \rrbracket \cdot v'_{n(w, j)}{}^\top v_{w_I} \right) \quad (3)$$

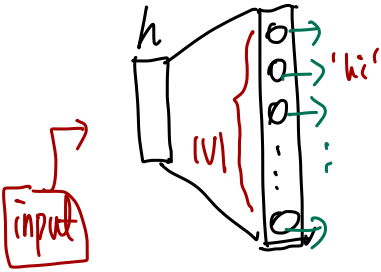
where $\sigma(x) = 1/(1 + \exp(-x))$. It can be verified that $\sum_{w=1}^W p(w | w_I) = 1$. This implies that the cost of computing $\log p(w_O | w_I)$ and $\nabla \log p(w_O | w_I)$ is proportional to $L(w_O)$, which on average is no greater than $\log W$. Also, unlike the standard softmax formulation of the Skip-gram which assigns two representations v_w and v'_w to each word w , the hierarchical softmax formulation has one representation v_w for each word w and one representation v'_n for every inner node n of the binary tree.

The structure of the tree used by the hierarchical softmax has a considerable effect on the performance. Mnih and Hinton explored a number of methods for constructing the tree structure and the effect on both the training time and the resulting model accuracy [10]. In our work we use a binary Huffman tree, as it assigns short codes to the frequent words which results in fast training. It has been observed before that grouping words together by their frequency works well as a very simple speedup technique for the neural network based language models [5, 8].

→ Q. check 1

Q. check 2.

① 1개 softmax



training time. The basic Skip-gram formulation defines $p(w_{t+j}|w_t)$ using the softmax function:

$$p(w_o|w_I) = \frac{\exp(v'_{w_o} \top v_{w_I})}{\sum_{w=1}^W \exp(v'_{w_o} \top v_{w_I})} \quad (2)$$

where v_w and v'_w are the "input" and "output" vector representations of w , and W is the number of words in the vocabulary. This formulation is impractical because the cost of computing $\nabla \log p(w_o|w_I)$ is proportional to W , which is often large (10^5 – 10^7 terms).

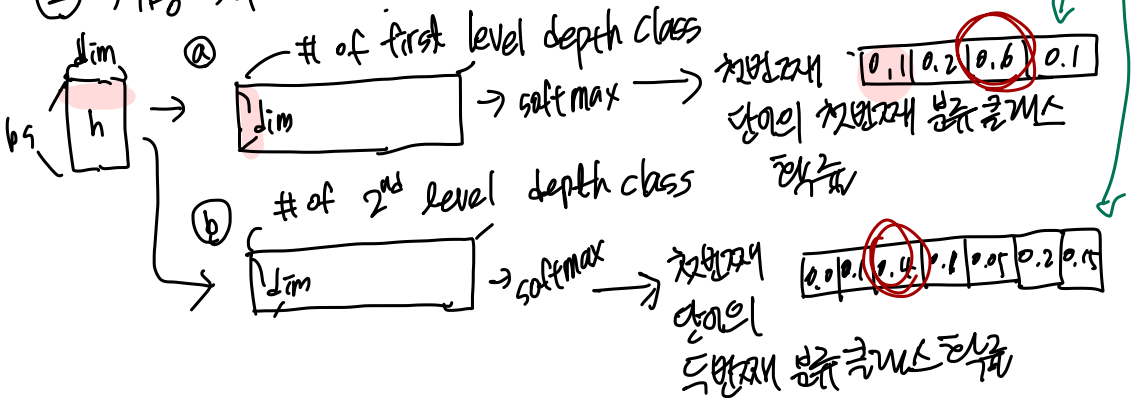
⇒ 'hi'라는 단어의 prob을 알려주려 hidden vector

h 는 output matrix $\in \mathbb{R}^{hidden-dim \times vocab-size}$

Vocab size 많음 softmax 수행 $\Rightarrow \log(Vocab\ size)$

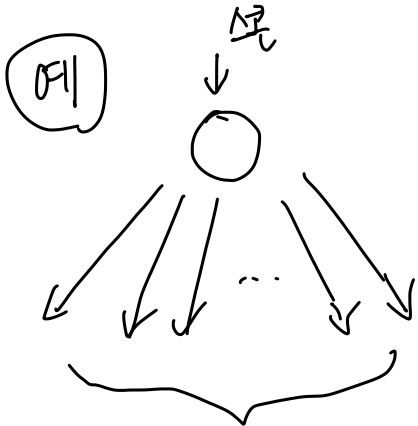
→ 과거 같은 이거 너무 computation 크다!

② 계층 softmax

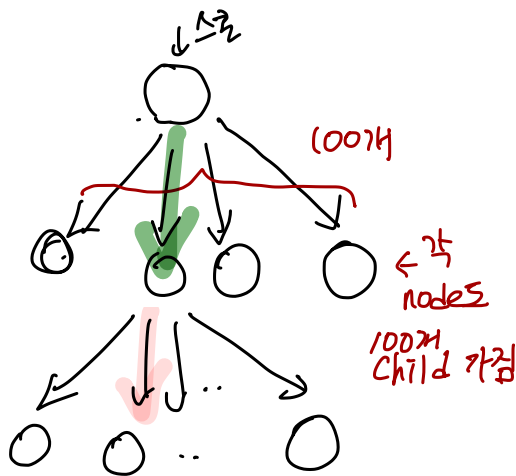


$$\therefore \text{최종 } P(W_{out} | h_{in}) = P(W_{첫번째계층클래스} | h_{in}) \times P(W_{두번째계층클래스} | h_{in})$$

W_{out} : 최종 prediction 단어
 h_{in} : input 단어 embedding



10,000개 문
10,000개 word 확률
다 계산

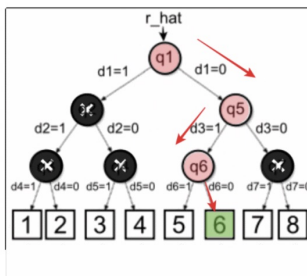
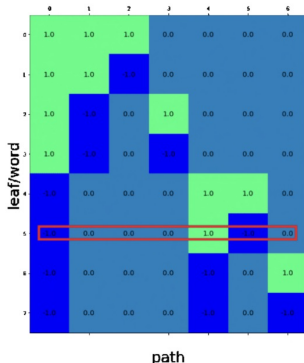


→ : 100번 수행
→ : 100번 수행

∴ output 확률 계산 위해 200번 수행

Q. 이렇게 계산해서 확률에 쓰려면?

⇒ training data 의 ground truth label 이
무엇인지 index 와 함께 graph traverse path도
주어야 함

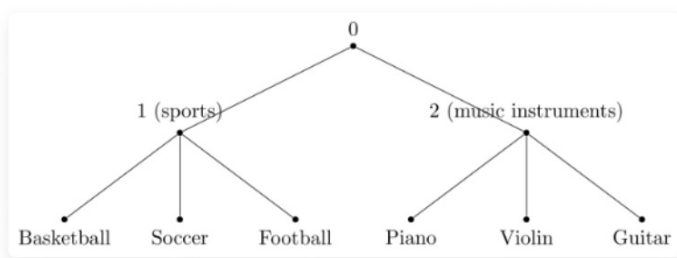


Q. 객 이진트리여야 하나?

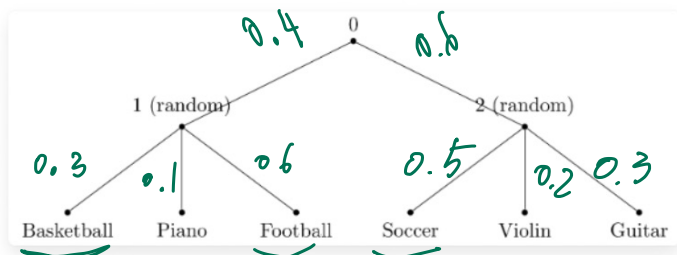
⇒ 대부분 하프만 tree 구조에서 이진 tree가 많은데,

객 그걸 필수는 없음 (루트 큰 레프랑 리프만 때..)

⇒ 단, 이렇게 계층을 쌓는지가 성능에 영향 많이 줌



Good Tree



• basketball, football, soccer 이 sports로 각자 높은 확률
속인데,

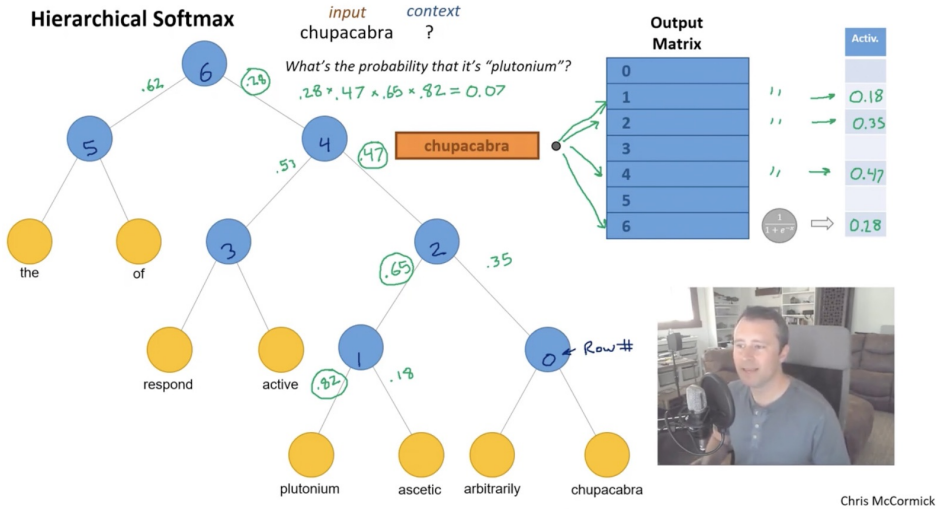
$$0.4 \times 0.3 = 0.12$$

$$0.4 \times 0.6 = 0.24$$

$$0.6 \times 0.5 = 0.3$$

모델이 임의이론하면 soccer이
있는쪽으로 첫 layer부터 높은
확률 주수 있게 해야 하는데,
유사한 클래스가 많다면
모델이 헷갈릴 수 있음.

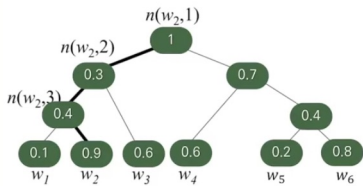
히프와 트리



4. Complexity Reduction

Efficient Estimation of Word Representations in Vector Space

• Hierarchical Softmax



$$\sum_{i=1}^6 P(w_i = w_0) = 1 \times 0.3 \times 0.4 \times 0.1 + 1 \times 0.3 \times 0.4 \times 0.9 + 1 \times 0.3 \times 0.6 + 1 \times 0.7 \times 0.6 + 1 \times 0.7 \times 0.4 \times 0.2 + 1 \times 0.7 \times 0.4 \times 0.8 = 1$$

$$P(n, \text{left}) = \sigma(v_n^T \cdot h)$$

: node n의 left로 갈 확률

$$P(n, \text{right}) = \sigma(-v_n^T \cdot h)$$

: node n의 right로 갈 확률

$$P(n, \text{left}) + P(n, \text{right}) = 1$$

: sigmoid이므로 위의 식이 성립

$$\sum_{i=1}^V P(w_i = w_0) = 1$$

: 모든 단어의 값에 대한 계산 없이

전체 합이 1이 되는 확률값을 구할 수 있다.

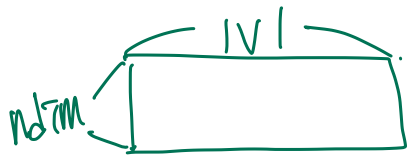
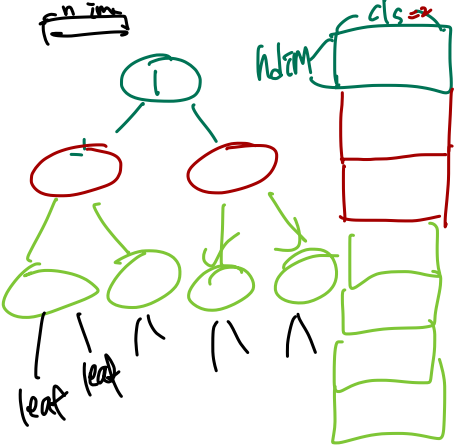
Computational complexity $\rightarrow D \times \ln(V)$

D: 각 step마다 길이 D의 vector 입력

$\ln(V)$: binary tree의 depth는 $\ln(V)$ 에 비례한다. 26

forward 계산 주고 backward도 주겠지만,

model complexity (모델 파라미터) 증가.. 크기가 커진다.



$$\underline{\text{ndim}} \times \overset{\text{Intermediate}}{\text{node}}_{\text{NT}} \times 2.$$

$$\text{intermediate node}_{\text{NT}} \leq \frac{|V|}{2}.$$