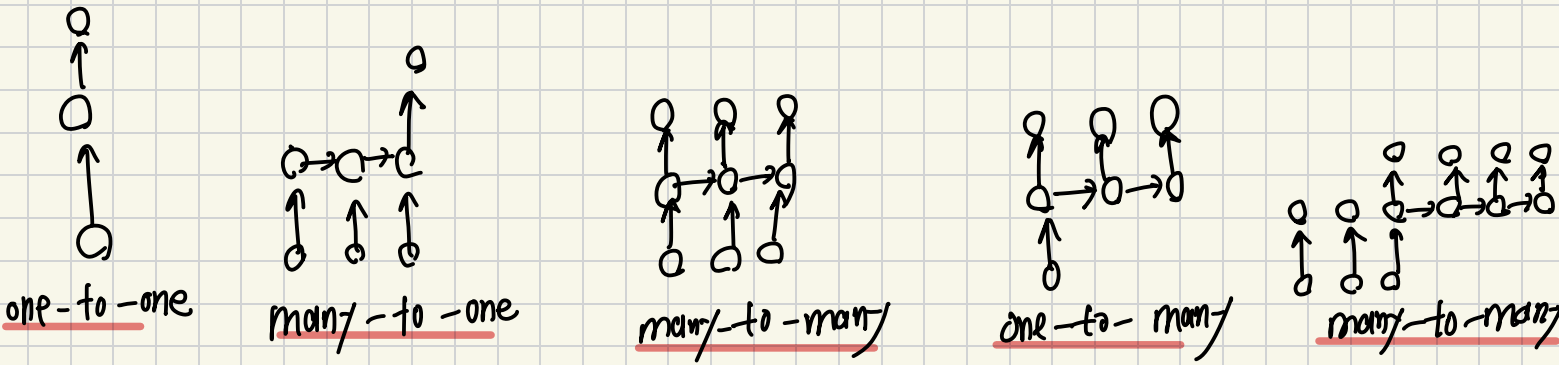


Sequential Data

순서가 중요한 data를



Word embedding

data를 vector로 만들어준다.

embedding : N-1을 유클리드 dimension에 나눈 값과 관련 (분포)가 다르



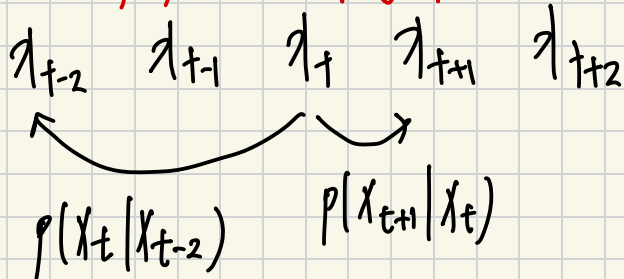
비슷한 단어들은 distance가 ↓
 다른 단어들은 distance가 ↑ 가 되도록
 word를 vector로 embedding 한다.

word2vec : 의미가 있는 최소 단위

data driven approach \Rightarrow 사전적 정의 사용X 어느 단어와 같이 사용되는지 통계적 추론.

< word2vec >

\Rightarrow 결국 sentence에서 단어 사용 likelihood $(L(\theta) = \prod_{i=1}^n \prod_{-m \leq j \leq m} p(w_{t+j} | w_t; \theta))$ 를 maximize.



< Glove >

Learn two separate embeddings, one as i and the other as j . Later, they are averaged.

$$w_i^T w_j = \log p(i|j) \quad J = \sum_{i,j=1}^V f(X_{ij}) (w_i^T \tilde{w}_j + b_i + \tilde{b}_j - \log X_{ij})^2$$

of times word j occurs in the context of word i .

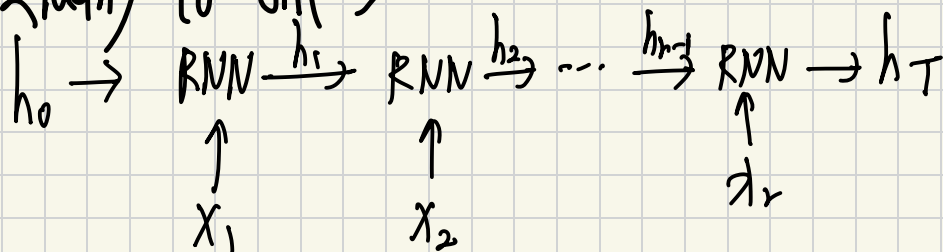
Dot-product similarity of two word embeddings
 \approx Order of word-pair frequency

: 자주 동시에 나오는 단어들은 비슷한 vector로 embedding 하자 라는 Idea 이다.

RNN for sequence classification (T/F를 맞추기)

* Image라 다르게 읽히길 갖는 data를 처리할 수 있어야 한다.

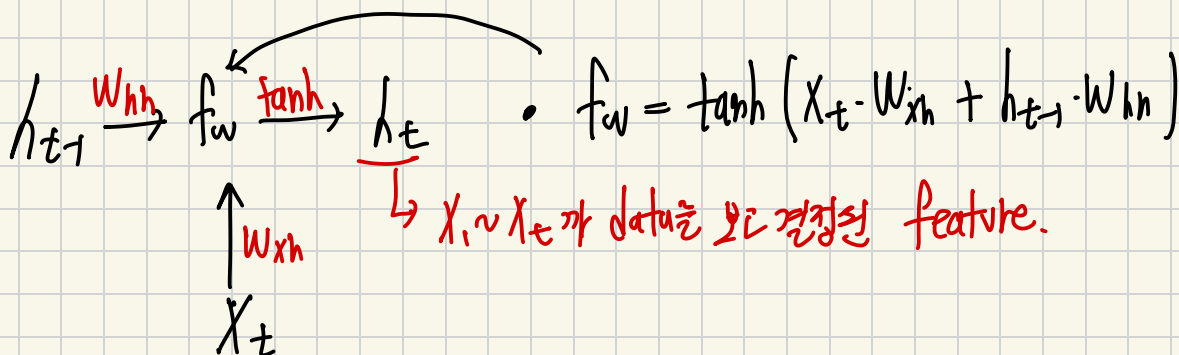
< Many to one >



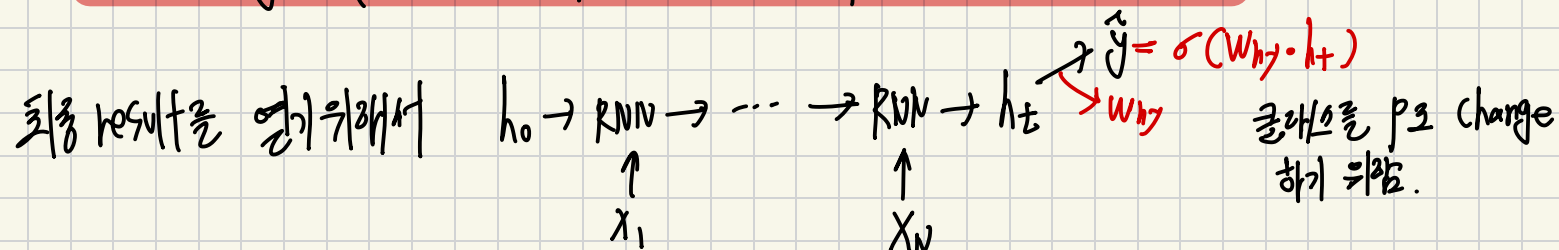
$$h_t = f_w(h_{t-1}, x_t)$$

이전 h 와 입력을 이용.

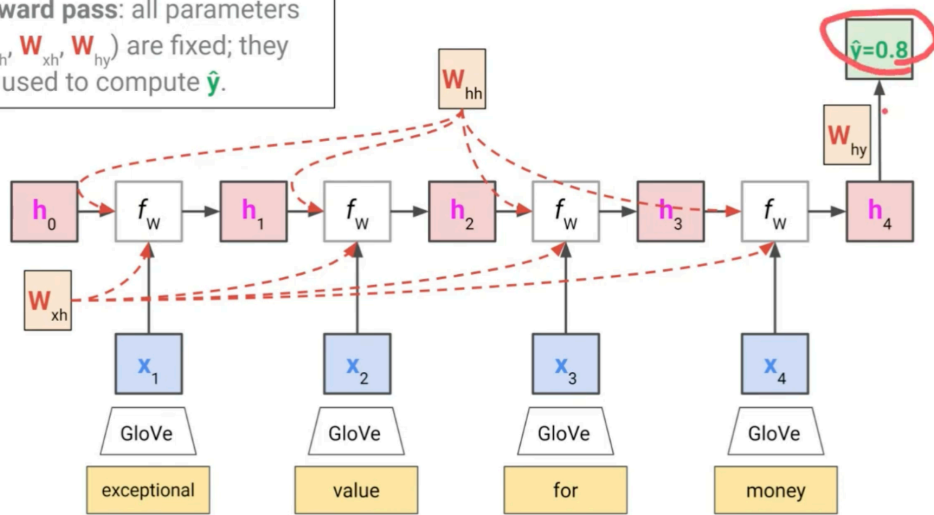
f_w 은 어느 step에나 same이다.
 Because 길이에 영향을 받으면 순서가 바뀌기 때문



모든 weight (W_{nh}, W_{xh})는 모든 step에 대해 same.



Forward pass: all parameters (W_{hh} , W_{xh} , W_{hy}) are fixed; they are used to compute \hat{y} .



이제 Backpass를 하면 되기 단계는 계속 동일한 W gradient를 구해야 한다.

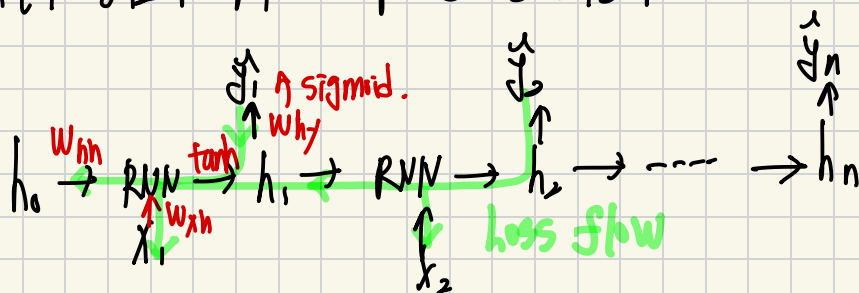
< Many to Many >

언어 = 단어들을 특정 규칙에 맞게 배열하는 것.

Traditional: n-grams.

모든 data에 대해 조건부 확률로 A라는 단어가 나올 때 B가 나올 확률은 모든 사전으로 등록
 \Rightarrow 결국 정확도 \uparrow 사건의 크기가 \uparrow 되게 된다!

X_t 가 들어오자마자 output을 출력해준다.

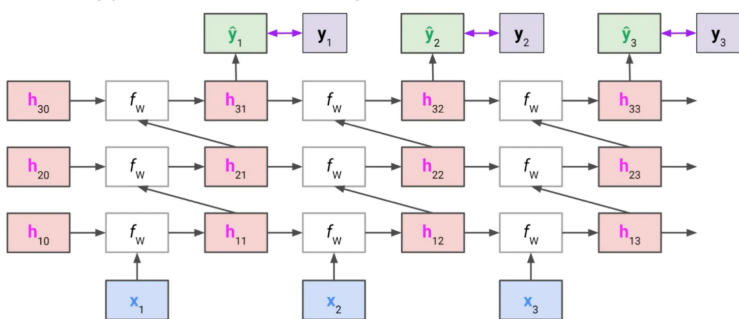


X 는 사전에 존재하는 모든 word 들 중 가장 높은 값으로 예측.

\Rightarrow train은 cross entropy 사용.

Multi-Layer RNN (RNN을 여러층 쌓아 올린다)

- We may put more than one hidden layers.



장) Input sequence 길이에 제약 X.
 \rightarrow 매 step 가중치를 동일하게 하여 해결.

단) 느리다 (순차적 처리)

vanishing gradient & long range dependence

