

## Gradient Descent (we never use)

$$\theta_j^{\text{new}} = \theta_j^{\text{old}} - \alpha \frac{\partial}{\partial \theta_j^{\text{old}}} J(\theta)$$

problem: Calculating our objective function over all of our data for our model (遍历所有训练数据) would take very very long time

## Stochastic Gradient Descent (always use for neural net) $\xrightarrow{\text{SGD}}$

also referred to Mini Batch Gradient Descent

Repeatingly sample windows, and update after each one

$\xrightarrow{\text{pick a very small subset of data}}$

this is a noisy, inaccurate estimate of the gradient

虽然 SGD 每一步看上去是在“瞎猜”，但在统计意义上是正确的，这些 noisy gradient 的“长期平均方向”，正好就是批量梯度下降的准确方向

Example: "the quick brown fox jumps" window size = 1 (center, context)

(quick, the) (quick, brown) (brown, quick) (brown, fox)  
(fox, brown) (fox, jumps) (jumps, fox)

## Traditional GD

1. input all the sample into model
2. calculate the average of the loss of all sample
3. calculate the gradient that loss to all parameters
4. update

## SGD

1. pick first sample (brown, fox) randomly, calculate this sample's loss to all parameters  
then update.

PS: 虽然依旧是针对所有参数，但在例子 like the 之类的参数上，梯度会非常大或为 0

main gradient is focus on the UBrown and Utok

2. pick second sample (quick, the) randomly

3. Continue like this - - -

(Bow)

Bag of words

Count word occurrences, ignore order



LSA

reduce dimension for Bow



Glove (2014)



Word2vec (2013)

→ SG (skip-gram)

Skip-gram

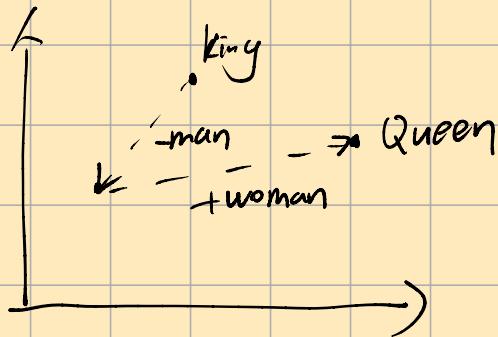
center word predict context

CBOW

context predict center word



(Continue Bag of words)



Analogy (类似)

Why have 2 vectors?

In math:

$$2 \text{ vecs} : P(O|C) = \frac{\exp(u^T v_c)}{\sum_{w=1}^V \exp(u_w^T v_c)}$$

$$1 \text{ vec} : P(O|C) = \frac{\exp(w^T v_c)}{\sum_{w=1}^V \exp(w_w^T v_c)}$$

Both 分子 & 分母 have "v<sub>c</sub>"

Obviously, it's very complex if we want to calculate the partial derivative

① skip-gram model with negative Sampling  
(负采样)

why use it?

Cause in Softmax, it takes long time to calculate the value of  $\sum_{w=1}^V \exp(u_w^T v_c)$ , The value of  $V$  is very big

what is the idea?

Maximize probability of real outside word

Minimize probability of random words

formula?

$$J_{\text{neg-sample}}(u_o, v_c, V) = -\log \sigma(u_o^T v_c) - \sum_{k \in \{\text{K sampled indices}\}} \log \sigma(-u_k^T v_c)$$

↑ make bigger  
make smaller

$\sigma(\cdot)$ : sigmoid

or logistic function



we learned from ML

$$\sigma(x) = \frac{1}{1+e^{-x}}$$

In this algorithm, we want Positive sample as big as possible

↓  
(probability of real outside word is "yes" bigger)

negative sample as big as possible too

(so we make  $u_k^T v_c$  smaller to ensure  $-u_k^T v_c$  bigger)

Unigram distribution (-元分布) 把各单个词的出现程度独立地考虑它们

We always sample with  $P(w) = U(w)^{\frac{3}{4}} / Z$ , the unigram distribution  $U(w)$  the <sup>raised to</sup> power



obviously, raising up the probability of the less frequent word

Why raise to  $\frac{3}{4}$  power?

make sampling more uniformly, do not overly lean towards high-frequency words (such as 'the', 'a') [why? in math, example  $\begin{cases} 10000^{0.75} \approx 5623 \\ 10^{-0.75} \approx 5.6 \end{cases}$ ]

## (2) Co-occurrence matrix (共發現矩阵)

Eg: Window length = 1

corpus: I like you  
I like him

	I	like	you	him
I	0	2	0	0
like	2	0	1	1
you	0	1	0	0
him	0	1	0	0

We find that the matrix is very high dimensional

So we want to reduce the dimensionality of that matrix (usually 25-1000 dimensions, similar to word2vec)

How to reduce?

Classic method: Singular Value Decomposition (奇异值分解) I've learned from linear algebra

$$\text{A matrix } \Rightarrow X = U \Sigma V^T$$

1.  $U$  &  $V$  are orthonormal

$$2. \Sigma = \begin{bmatrix} s_1 & & & \\ & s_2 & & \\ & & s_3 & \\ & & & \ddots \\ & & & s_{n-1} \end{bmatrix} \quad s_1 > s_2 > s_3 > \dots > s_n$$

(It rules that  $\Sigma$  order in size)

then, we set the smallest singular vectors to 0

## (3) Global Vectors for Word Representation (GloVe model)

Crucial insight: Ratios of co-occurrence probabilities can encode meaning components  
We want to capture them as linear meaning components in a word vector space

$$P(x|\text{ice}) \quad x = \text{solid} \quad x = \text{gas} \quad x = \text{water/random}$$

$$\text{eg: } P(x|\text{steam}) \quad \text{large} \quad \text{small} \quad 1 \quad 1$$

this means if the ratio is large or small,  $x$  can distinguish ice and steam  
if the ratio is 1,  $x$  can't distinguish them

How can we capture the ratios as linear meaning components?

using Log-bilinear model :  $w_i \cdot w_j = \log P(i|j)$

with vector differences :  $w_x \cdot (w_a - w_b) = \log \frac{P(x|a)}{P(x|b)}$  ← We assume it (come dot product also represent the probability)

$$\begin{cases} w_x \cdot w_a = \log P(x|a) \\ w_x \cdot w_b = \log P(x|b) \end{cases}$$

$$\text{Loss: } J = \sum_{i,j=1}^V f(X_{ij}) (w_i^T \tilde{w}_j + b_i + b_j - \log X_{ij})^2$$

(4) how to evaluate word vectors?

A general concept of evaluation (in NLP)

Intrinsic (内在)

Evaluation on a specific / intermediate task

word vector distances

to evaluate meaning similarity

human judgement

Extrinsic (外在)

Evaluation on real task

one example could help evaluate named entity recognition (NER)

(5) Word senses and word sense ambiguity

problem: most words have lots of meaning

how to solve?

Idea: Cluster word windows around words, retrain with each word assigned to multiple different clusters

but we don't usually do

what's we usually do?

⇒ Different senses of a word reside in a linear superposition (weighted sum)  
in standard word embeddings like word2vec

$$\text{Eg: } v_{\text{pike}} = \alpha_1 v_{\text{pike1}} + \alpha_2 v_{\text{pike2}} + \alpha_3 v_{\text{pike3}}$$

叠加  
↑  
not yet

where  $\alpha_i = \frac{f_i}{f_1 + f_2 + f_3}$   $f \leftarrow$  frequency

## (6) Cross entropy

a concept from information theory

$$H(p, q) = - \sum_{c=1}^C p(c) \log q(c)$$

$p$  = true, objective probability distribution

$q$  = the probability distribution your model've predicted

if  $p$  and  $q$  are same, the cross-entropy will be small

if differ significantly the cross-entropy will be large

A special case:

when  $p$  is One-hot coding ( $p = [0, 1, 0, 0, 1 \dots]$ )

$$H(p, q) = - \sum_{c=1}^C p(c) \log q(c) \quad \leftarrow p(c) \text{ will be } 0 \text{ except the right one}$$

$$= -\log q(C_r) \quad \leftarrow C_r \text{ is the right class which model predict}$$

the essence of cross entropy is:

To measure the difference between two probability distributions

## (7) neural net work

### NER: Binary classification for center word being location

- We do supervised training and want high score if it's a location

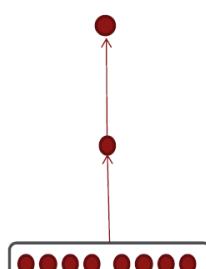
$$J_t(\theta) = \sigma(s) = \frac{1}{1 + e^{-s}}$$

↑  
predicted model  
probability of class



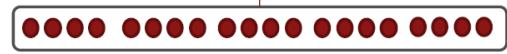
$$s = \mathbf{u}^T \mathbf{h}$$

$$\mathbf{h} = f(\mathbf{Wx} + \mathbf{b})$$



$f =$  Some element-wise non-linear function, e.g.,

$x$  (input)  $\in \mathbb{R}^{5d}$



$$x = [x_{\text{museums}} \quad x_{\text{in}} \quad x_{\text{Paris}} \quad x_{\text{are}} \quad x_{\text{amazing}}]$$

Embedding of  
1-hot words

