

Section 2 (skip-gram with negative sampling (SGNS) setup)

skip-gram model : stream of words $w_i \in [N]$, contexts $c_i \in [N]^L := [M]$
 observed in n pairs (w_i, c_i) (so c_i is fixed window around w_i)
 and looking for word representations w_1, \dots, w_N ,
 context representations c_1, \dots, c_M

To maximize

$$\sum_{i=1}^n f(w_i, c_i; w_1, \dots, w_N, c_1, \dots, c_M)$$

where

$$f(t, l; \dots) = \log \left[\gamma(w_e^T c_t) + t \cdot \log \gamma(-w_e^T c_l) \right]$$

$$\text{where } L \sim \text{categorical} \left\{ \frac{n_1}{n}, \dots, \frac{n_m}{n} \right\}$$

and $\gamma(u) = \frac{1}{1+e^{-u}}$ is the standard logistic function Logit

which, by the way, satisfies $\delta(u) = 1 - \gamma(u)$, and $\gamma'(p) = \frac{1}{(1-p)^2}$

Binary Logistic regression

$$P(Y=1 | X) = \gamma(\theta^T X) ; P(Y=0 | X) = 1 - \gamma(\theta^T X) = \gamma(-\theta^T X)$$

Interpretation:

fix c_1, \dots, c_m , and replace $t \cdot \log \gamma(-w_e^T c_l)$

$$\text{with } \log \gamma(-w_e^T c_{l_1}) + \dots + \log \gamma(-w_e^T c_{l_t})$$

L_i : i.i.d. replicates of L

Then optimizing

$$\log \gamma(w_e^T c_{l_1}) + \log \gamma(-w_e^T c_{l_2}) + \dots + \log \gamma(-w_e^T c_{l_t})$$

for w_e is getting a logistic classifier to recognise c_{l_i} as valid (class 1)

and c_{l_i} as invalid (class 0)

Section 3 (Implicit matrix factorization)

Remarks we are trying to maximize

$$\sum_{i=1}^n f(w_i, c_i; w_1, \dots, w_N, c_1, \dots, c_M)$$

where

$$f(w_i, c_i; \dots) = \log \left\{ r(w_i^T c_i) \right\} + t \cdot \log \left\{ r(-w_i^T c_i) \right\}$$

where $L \sim \text{categorical} \left\{ \frac{n_1^{(c)}}{n}, \dots, \frac{n_M^{(c)}}{n} \right\}$

Let $M_{i,c} = w_i^T c_i$. Then can write objective as optimisation

of $\sum_{i,c} (\text{terms involving } M_{i,c})$

where each term in the sum is

$$n_{i,c}^{(wc)} \log(p_{i,c}) + t \cdot n_{i,c}^{(w)} \cdot \frac{n_{i,c}^{(c)}}{n} \cdot \log(p_{i,c})$$

Now look to maximise for each $p = r(M_{i,c})$:

$$n_{i,c}^{(wc)} \log(p) + t \cdot n_{i,c}^{(w)} \cdot \frac{n_{i,c}^{(c)}}{n} \log(1-p)$$

$$\hat{p} = \frac{n_{i,c}^{(wc)}}{n_{i,c}^{(wc)} + t \cdot n_{i,c}^{(w)} \cdot \frac{n_{i,c}^{(c)}}{n}}$$

$$\hat{m}_{i,c}^{(wc)} = \log(p) = \log\left(\frac{\hat{p}}{1-\hat{p}}\right) = \log\left(n \cdot \frac{n_{i,c}^{(wc)}}{t \cdot n_{i,c}^{(w)} \cdot \frac{n_{i,c}^{(c)}}{n}}\right)$$

$$= \log\left\{\frac{\frac{n_{i,c}^{(wc)}}{n}}{\left(n_{i,c}^{(w)}/n\right) \left(n_{i,c}^{(c)}/n\right)}\right\} - \log(t).$$

For discrete variables X, Y , the pointwise mutual information matrix is

$$\log \left\{ \frac{P(x, y)}{P(x) P(y)} \right\}, \text{ and so } \hat{M} \text{ can be viewed as a shifted empirical PMI.}$$

Obviously can find decomposition $\hat{M}_{kl} = \hat{W}_k^T \cdot \hat{C}_l$
 if \hat{W}_k, \hat{C}_l are high-dimensional enough (e.g. via SVD)

So the claim is SGDNS is factoring \hat{M} , but I don't know how to make this statement formal if you admit \hat{W}_k, \hat{C}_l aren't arbitrarily high-dimensional.

Just because:

$$\hat{M} = \hat{W} \hat{C}^T, \text{ where } \hat{w}, \hat{c} = \arg \max(\text{objective function}) \quad (\text{unconstrained})$$

doesn't mean:

$$\hat{M} \approx \tilde{W} \tilde{C}^T, \text{ when } \tilde{w}, \tilde{c} = \arg \min(\text{obj}) \quad (\text{constrained})$$

A weighted matrix factorization?

SVD of \hat{M} gives $\hat{M} = \tilde{W} \tilde{C}^T$ to minimize $\|\hat{M} - M\|_F^2 = \sum |\hat{M}_{ij} - M_{ij}|^2$

However, we see that objective

$$\underbrace{n_{kl}^{(wc)} \log \delta(M_{kl})}_{\text{fixed}} + t \cdot \underbrace{n_k^{(w)} \cdot \frac{n_l^{(c)}}{n} \cdot \log \delta(-M_{kl})}_{\text{fixed}}$$

is not a fixed function of $\hat{M}_{kl} = \log \left(\frac{a}{b} \right)$

and in fact for $\log \left(\frac{a}{b} \right)$ fixed, i.e. $\frac{n_{kl}}{n_k n_l}$ fixed,

contribution is α times higher if word k or context l appears α more times.

Performance metrics

- rank correlation between word pairs (human assigned similarity) (Spearman correlation)
and corresponding cosine similarity $s(u, v) = \frac{u^T v}{\|u\| \|v\|}$
- Analogy "Paris is to France as Tokyo is to Japan"
found using equation $s(b^*, a^*) s(b^*, b) / [s(b^*, a) + \epsilon]$
e.g. $s(\text{Japan}, \text{France}) s(\text{Japan}, \text{Tokyo}) / s(\text{Japan}, \text{Paris})$