# Machine learning

| A | B | C | D | label |
|---|---|---|---|-------|
| $a_1$ | $b_1$ | $c_1$ | $d_1$ | $l_1$ |
| $a_2$ | $b_2$ | $c_2$ | $d_2$ | $l_2$ |
| $a_3$ | $b_3$ | $c_3$ | $d_3$ | $l_3$ |
| $a_4$ | $b_4$ | $c_4$ | $d_4$ | $l_1$ |
| $a_5$ | $b_5$ | $c_5$ | $d_5$ | $l_2$ |
| $a_6$ | $b_6$ | $c_6$ | $d_6$ | $l_2$ |

← $n$ tuple (vertical)   ← $k$ feature →

let us say there are '$l$' labels
'$k$' feature
'$n$' tuple



$I_1$, $I_2$ ... $I_k$ → $L_1$ → label 1

$L_{2l}$ → label '$l$'

weight label means
each of these

for each tuple we will read if
weight can be varied & xxx
is xxx upon input

| A | B | C | D | | Labels |
|---|---|---|---|---|---|
| $a_1$ | $b_1$ | $c_1$ | $d_1$ | | $l_1$ |
| $a_2$ | $b_2$ | $c_2$ | $d_2$ | | $l_2$ |
| $a_3$ | $b_3$ | $c_3$ | $d_3$ | | $l_3$ |
| $a_4$ | $b_4$ | $c_4$ | $d_4$ | | $l_1$ |
| $a_5$ | $b_5$ | $c_5$ | $d_5$ | | $l_2$ |
| $a_6$ | $b_6$ | $c_6$ | $d_6$ | | $l_2$ |

← n tuples ↓

← k features →

let us say there are '$l$' labels
'$k$' features
'$n$' tuples



$I_1$
$I_2$
$I_k$

$L_1$ → label 1
$L_l$ → label '$l$'

weighted label nodes
each of these
has weights for
each feature

for each tuple we will
have '$k$' features so
input nodes have it

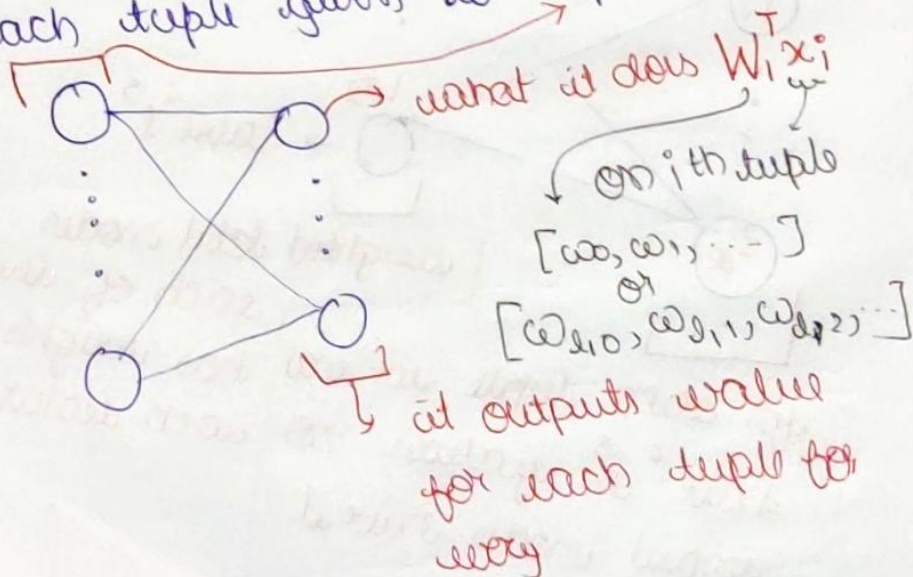→ each $\bigcirc$ $L_1$ label node has weights for each feature which are randomly initialized so

$L_1 \to$ outputs $W^T = [\omega_1, \omega_2 \ldots \omega_k]$

* can have constant or '0' weight.

$$W^T = [\omega_0, \omega_1 \ldots \omega_k]$$

|←————— $k$ features + 1 ——————→|
constant

* $X:$ $\begin{bmatrix} 1 & a_1 & b_1 & c_1 \\ 1 & a_2 & b_2 & c_2 \\ 1 & & & \\ \vdots & \vdots & & \end{bmatrix}$ $\Big\}$ $n$ tuples (leave out label)

|←——— $k$ feature + 1 ———→|

* each tuple given to input node



→ what it does $W_i^T x_i$

on $i$th tuple

$[\omega_0, \omega_1, \ldots]$
or
$[\omega_{i,0}, \omega_{i,1}, \omega_{i,2} \ldots]$

↳ it outputs value for each tuple for every

so each row will have value for each label

$$
\begin{bmatrix} \omega_{l,0} \\ \omega_{l,1} \\ \vdots \\ \omega_{l,k} \end{bmatrix}_{(k+1) \times 1}
\begin{bmatrix} x_0 & x_1 & \dots & x_k \end{bmatrix}_{1 \times (k+1)}
$$

$$
\begin{bmatrix} \omega_{l,0} & \omega_{l,1}, & \dots & \omega_{l,k} \end{bmatrix}_{(1 \times (k+1))}
*
\begin{bmatrix} x_0 \\ \vdots \\ x_k \end{bmatrix}_{(k+1) \times 1}
$$

↓ this is for $l_i$ label

* each node will have values for all  label

$[P_{i1}, P_{i2}, \dots P_{il}]$

values that label nodes return for $i^{th}$ row of the dataset

$$
\begin{bmatrix} c \\ c \end{bmatrix}
$$

$\dfrac{0.18P6}{+0.18P6}$  $\cdot (\overline{5}) \partial$      $0.18P6$

$0.18P6 =$

$0.45$

\* Now what is the use of
$$[\vec{\phi_{i1}}, \vec{\phi_{i2}}, \ldots \vec{\phi_{il}}] ?$$

whichever label has the higher probability will have more chance to belong to it

→ here comes the role of softmax

why softmax?

→ softmax function is a function that turns a vector of k real values into a k real vector which sums to 1

→ Input value might be positive, negative or zero. But transforms b/w 0 and 1

→ so it can be treated as probabilities

$$\sigma(\vec{z_q})_i = \frac{e^{z_i}}{\sum_{j=1}^{k} e^{z_j}}$$

1 example worth than 100 words

$$\begin{bmatrix} 8 \\ 5 \\ 0 \end{bmatrix} \text{ be } z$$

$e^8$, $e^{z_0}$ = 2981.0

$e^5$, $e^{z_1}$ = 148.4

$e^0$, $e^{z_2}$ = 1

softmax

$$\sigma(\vec{z})_0 = \frac{2981.0}{2981.0 + 148.4 + 1}$$

$$= \frac{2981.0}{3130.6}$$

$$= 0.95$$

$$\sigma(\vec{z})_1 = \frac{148.4}{3130.4} = 0.047$$

$$\sigma(\vec{z})_2 = \frac{1}{3130.4} = 0.0003$$

so label 0 has probability of 95%.

→ so now $[\overset{S}{\underset{}{P_{i1}}}, \overset{S}{\underset{}{P_{i2}}}, \dots \overset{S}{\underset{}{P_{il}}}]$ are values

we got for ith row for each label

give this to softmax.

$$P_{i1} = \frac{e^{S_{i1}}}{Z} \qquad \underrightarrow{} \qquad \sum_{j=1}^{l} e^{S_{ij}}$$

returns probability
that ith row belongs
to label 1

we have 'l' labels
so their summation

$$P_i = [\,P_{i1} \quad P_{i2} \dots \quad P_{il}\,]$$

probabilities of ith de Row belonging
to each label.

$$Y_i = [\,y_{i1} \quad y_{i2} \dots \quad y_{il}\,]$$

this is what actually
are label
if ith tuple has l labels

$$Y_i = [1 \; 0 \; 0 \dots 0]$$

Now comes KL divergence

it quantifies or shows how
much are the probabilities
differ in both the vectors

Kind of loss function which are
used to reduce

$$KL(P||Q) = -\sum_{j=0}^{l} \cdot Pij \; log\left(\frac{Qij}{Pij}\right)$$

i th row
ij is label
probability

now we have ground truth as

$\vec{Y_i}$ for i th row

$\vec{P_i}$ for ith row as prods generated?

$$KL(\vec{Y_i}||\vec{P_i}) = -\sum_{j=0}^{l} Pij \; log\left(\frac{Yij}{Pij}\right)$$

Q) $Y_i = [1, 0, 0 \ldots]$

something just like this
so we can show that
class label (prob where it
is not '0'

suppose say label is 1) $P_{i1}$ is all

we want

$$KL(\vec{Y_i} \| \vec{P_{i1}})^2 \quad -\ln P_{ij} + d\sum_{t=1}^{L} \|W_t\|^2$$

regularization

(CLOSS FONCTION)   where $\hat{j}$ is actual
label of the
tuple

Now minimize the class
and adjust weights $W_i^T \ldots W_L^T$
$[\omega_{i1}, \omega_{i21} \ldots \omega_{iL}]$

## Story Boils Down

$$\sum_{\substack{i=1 \\ i=\frac{1}{2}}}^{2} KL\left(\vec{Y_i} \| \vec{P_i}\right) + d\sum_{t=1}^{L} \| W_t \|^2 \qquad \left. \boxed{\phantom{XXXX}}^{\frac{1}{2}} \right)$$

for all tuples

$$\sum_{i=1}^{N} \sum_{j=1}^{L} I(Y_{ij}==1)\left(-\ln P_{ij}\right) + \left( \phantom{XXXX} \right)$$

don't worry it only means Identification function to recognize only the labels where $Y_{ij}$ is 1

$$-\ln P_{ij} = \frac{e^{W_j^T x_i}}{Z} \begin{array}{l} \rightarrow x_i \text{ ith row of the} \\ \text{set.} \\ \rightarrow j \text{ th label weight} \end{array}$$

$$\sum_{R=1}^{L} e^{W_R^T x_i}$$

NOW reduce loss $W_1^T, W_2^T \dots W_L^T$

let us first write derivation for

$$W_j^T = [\omega_{j1} \quad \omega_{j2} \dots \quad \omega_{jk}]$$

for som $\omega_{jk}$ as all are variables.

$$\frac{\partial L}{\partial W_{jk}} = -\sum_{i=1}^{n} \sum_{j'=1}^{L} I(y_{ij'}=1) \frac{1}{P_{ij'}} \boxed{\frac{\partial (P_{ij'})}{\partial W_{jk}}}$$

$$+ \lambda \frac{\partial}{\partial W_{jk}} \left[ \sum_{m=1}^{K} \omega_{jm}^2 \right] = 2\omega_{jk}$$

$$\frac{\partial (P_{ij'})}{\partial W_{jk}} = I(j=j') \left[ \underbrace{e^{W_j^T x_i}} \cdot \frac{x_{ik}}{Z} - \right. \overbrace{\qquad}$$

this only
appears when
this

$$\frac{e^{W_j^T x_i} \cdot e^{W_j^T x_i} \cdot x_{ik}}{Z^2}$$

$$= I(j=j') \cdot \frac{e^{W_j^T x_i} \cdot x_{ik}}{Z} - \frac{\left(e^{W_j^T x_i}\right)^2 \cdot x_{ik}}{Z^2}$$

$$= \frac{P_{ij} \cdot x_{ik} \cdot I(j=j')}{Z} - \frac{\left(P_{ij}\right)^2 \cdot x_{ik}}{Z^2}$$

$$= \frac{P_{ij} \cdot x_{ik}}{Z} \left( I(j=j') - \frac{P_{ij}}{Z} \right)$$

$$P_{ij} = \frac{e^{W_j^T x_i}}{\sum\limits_{k=1}^{L} e^{W_k^T x_i}} \qquad [W_{j1} \; W_{j2} \ldots W_{jk}]$$

$i \rightarrow$ row

$j \rightarrow$ label

features total are 'k'

$$\frac{1}{P_{ij}} \cdot \frac{\partial}{\partial W_{j'k}} \left( \frac{e^{W_j^T x_i}}{z} \right)$$

$j = j'$

some $\left[ \quad \bigcirc \quad \right] \rightarrow W_j$

$k^{th}$ feature weight

$$\left( \frac{x_{ik} \, e^{W_j^T x_i} \cdot z - e^{W_j^T x_i} \cdot}{z^2} \right)$$

$$\sum\limits_{t=1}^{L} \frac{\partial \, e^{W_t^T x_i}}{\partial W_{j'k}} \qquad \frac{- e^{W_j^T x_i} \cdot x_{jk} e}{z^2}$$

$$\frac{\partial L}{\partial W_{j'k}} = -\sum\limits_{i=1}^{N} \sum\limits_{j=1}^{L} \frac{1}{P_{ij}} \left[ \mathbb{I}(j = j') \left( \frac{x_{ik} \, e^{W_j^T x_i} \cdot z - e^{W_j^T x_i}}{} \right) \right.$$

$$\left( \frac{x_{ik} e^{W_j^T x_i} \, z - e^{W_j^T x_i} \cdot e^{W_{j'}^T x_i} x_{jk}}{z^2} \right)$$

$$+ \mathbb{I}(j \neq j') \left( 0 - \frac{e^{W_j^T x_i} \cdot x_{jk} \cdot e^{W_{j'}^T x_i}}{z^2} \right)$$

$$\left( x_{ik} \cdot P_{ij} - (P_{ij})^2 \, x_{ik} \right) I(j = j')$$

$$- I(j \neq j') \left( P_{ij'} \cdot P_{ij} \right) x_{ik}$$

$$\Rightarrow P_{ij} \, x_{ik} \left( I(j = j')(1 - P_{ij}) - I(j \neq j')(P_{ij'}) \right)$$

$$\frac{\partial P_{ij'}}{\partial \omega_{jk}} \Rightarrow \begin{cases} P_{ij} \, x_{ik}(1 - P_{ij}) & j = j' \\[2mm] P_{ij} \, x_{ik} \, P_{ij'} & j \neq j' \end{cases}$$

$$\omega_{jk} - \alpha \, \frac{\partial \, LOSS}{\partial \omega_{jk}}$$

| A | B | C |
|---|---|---|
|   |   |   |
|   |   |   |
|   |   |   |

| P | N |
|---|---|
|   |   |
|   |   |
|   |   |

NOW Version 2



Features "K" →          Labels L →

                    Feat
                    ures
                         ↓
↓                        F        PRIOR
data
denote   DATASET                  WORDS
by
(es)

(N × K)                  (F × L)

there are no              what are features?
labels                    features are nothing
(here also                but words
means the
same)

$I(g_f \in x_i) \rightarrow$ identity function

feature tuple i     $0 \rightarrow$ if printuple has no
                         $g_f$ feature
                    $1 \rightarrow$ if it has

Suppose

| Support   | Pos | Neutral | Neg  | $P_f$ |
|-----------|-----|---------|------|-------|
| Excellent | 0.9 | 0.05    | 0.05 | → all num |
|           |     |         |      | up to one |
| feature   |     |         |      | this what |
|           |     |         |      | F×L has. |

$\tilde{P}_f \rightarrow$ denotes probability vector for feature 'f' of the F×L Matrix or prior knowledge



$$C_f = \sum_{i=1}^{N} I(g_f \in r_i)$$

count of all the files having 'f' feature

we have got $P_i \begin{bmatrix} \quad \end{bmatrix}$ having probabilities of each tuple for each label

$$\tilde{P}_f = \frac{1}{C_f} \sum_{i=1}^{N} P_i \; I(g_f \in r_i)$$
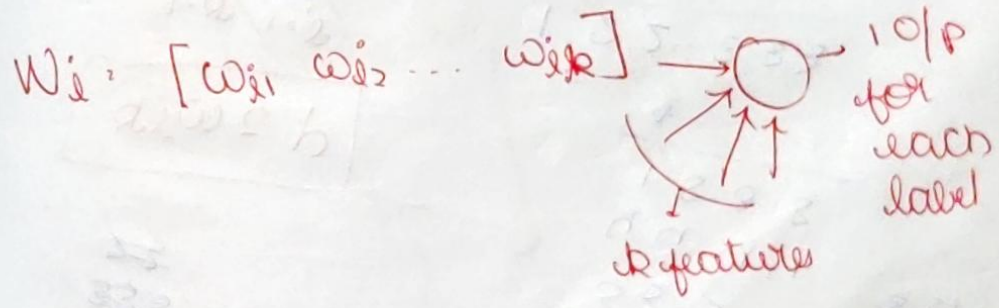
Sum of all probabilities vectors which has feature (f)

─────────────────────

total files which have feature 'f'

$\hat{P}$

| pos | neg | neu |
|-----|-----|-----|
|  |  |  |
|  |  |  |

ith

bias (N×L)

$\tilde{P}$

| pos | neg | neu |
|-----|-----|-----|
|  |  |  |
|  |  |  |

(F×L)

$$\sum_{f=1}^{F} KL(\hat{P}_f \| \tilde{P}_f) + \cancel{\text{....}}$$

$$\alpha \sum_{\ell=1}^{L} \|W_\ell\|^2$$

| | f | |
|--|--|--|
| | 1 | |

ith row

(N×K)

weights init randomly

loss function

$W_i = [\omega_{i1} \quad \omega_{i2} \ldots \quad \omega_{ik}]$ → ○ → o/p for each label

k features

We do differentiation for one
label and 1 feature weight keep in
mind repeat for $L \times k$ times

$$\frac{\partial O}{\partial \omega_{lk}} = \frac{\partial}{\partial \omega_{lk}}\left(\sum_{f=1}^{F}\left(\sum_{l'=1}^{L} \hat{P}_{fl'} \ln \frac{\hat{P}_{fl'}}{\tilde{P}_{fl'}}\right)\right)$$

subscript

for each label
for feature f

$$+ \frac{\partial}{\partial \omega_{lk}} \, d \sum_{l'=1}^{L} \sum_{k'=1}^{k} \omega^2_{l'k'}$$

$$= \sum_{f=1}^{F}\sum_{l'=1}^{L}\left(\ln \frac{\hat{P}_{fl'}}{\tilde{P}_{fl'}} + 1\right) \frac{\partial}{\partial \omega_{lk}} \hat{P}_{fl'}$$

$$d \sum_{l'=1}^{l}\sum_{k'=1}^{k}$$

$$d \, 2 \, \omega_{lk}$$

(upside-down text at bottom)

$$\frac{1\ 2\ 2\ 3\ 3\ 3}{\phantom{x}}$$

$$\frac{1\ 8\ 3\ 3}{4\ 4\ 4\ 4}$$

$$\frac{1\ 5\ 5\ 5\ 5\ 5}{1\ 4\ 4\ 4\ 4}$$

$$\frac{3\ 3}{1\ 2\ 2}$$

$$2\ 2$$

$$1$$

$$2\ 2\ 3\ 3\ 3$$

$$5\ 5\ 5\ 5\ 5$$

$8\ 5+1$

$5+0$

$2+0$

$2+2 = $

$2+3$

$n+3$

$$\hat{P}_{\theta'} = \frac{e^{W^T x_j}}{z}$$

---

# WORD2VEC

$M_1$ ~~~~ Movies $M_j$ ~~~~ $M_N$

may have
may not have
each movie
so some values
missing

$$r_{ij} \rightarrow \text{rating of } M_j$$
$$\text{by } u_i$$

users $\{ u_1, u_i, u_m$

$r_{ij}$

$M \times n$
matrix

$$u_i \longrightarrow P_{i}$$
$$\underset{\text{user i transformed}}{\downarrow}$$
$$\text{vector}$$

$$\frac{\partial p}{\partial a} \quad M_j \rightarrow Q_j$$
$$\underset{\text{movie j}}{\downarrow}$$

$P$

$M$

$Q$

$N$

$$\hat{r}_{ij} = P_i \, \omega_j^T + \underbrace{(b_i + b_j)}_{\text{biases}}$$

$\gamma_{ij}$ — predicted

$\gamma_{ij}$ — true

$$E_{ij} = [\gamma_{ij} - \hat{\gamma}_{ij}]^2$$

error$_{ij}$

$$\text{LOSS} = \frac{1}{N_D} \sum_{(i,j) \in D} E_{ij}$$

$\frac{1}{N_D}$ → no of entries in $D$

$D$ → sub matrix

$$= \frac{1}{N_D} \sum_{(i,j) \in D} [\gamma_{ij} - (P_i \omega_j^T + b_i + b_j)]^2$$



$D$

Machine Translation

source sentence $\xrightarrow{\text{paired}}$ target translation

□ □

learning Representation of words wring
unlabeled data

$[\omega, \qquad]$ $\omega T$

↓
make substrings of sizes
⟩ (examples)

$\omega_1 \; \omega_2 \; \underline{\omega_3} \; \omega_4 \; \omega_5$ $\longrightarrow$ PIVOT

$P(\underline{\omega_3} | \; \underline{\omega_1 \; \omega_2 \; \omega_4 \; \omega_5})$

② given these words

appearing
pivot

$$\frac{\partial}{\partial \omega_{jk\theta}} \left( \sum_{f=1}^{F} \sum_{\ell'=1}^{L} \hat{P}_{f\ell'} \ln \frac{\hat{P}_{f\ell'}}{\tilde{P}_{f\ell'}} \right)$$



$$\frac{\partial \hat{P}_{f\ell'}}{\partial \omega_{jk\theta}} \ln \frac{\hat{P}_{f\ell'}}{\tilde{P}_{f\ell'}}$$

$$+ \quad \frac{\tilde{P}_{f\ell'}}{\hat{P}_{f\ell'}} \hat{P}_{f\ell'} \frac{\partial}{\partial \omega_{jk\theta}} \left( \frac{\hat{P}_{f\ell'}}{\tilde{P}_{f\ell'}} \right)$$

$$\left( \ln \frac{\hat{P}_{f\ell'}}{\tilde{P}_{f\ell'}} + 1 \right) \frac{\partial}{\partial \omega_{jk\theta}} \left( \hat{P}_{f\ell'} \right)$$

$$\frac{1}{C_f} \sum_{j=1}^{2} I_\theta (f \in x_i) \, P_{i\ell'}$$

$$x_{ik} \, P_{i\ell'} \left( [\, I(\theta = \ell') - P_{ij} \,] \right)$$

# CBOW (Continuous bag of words)

skipgram model     (center word given return context |v|)

|V| {



vector of word
$W$

n {



$\vartheta$   for of word

$h = W x^T$

$h^{th}$ row of $W$   $|V| \times 1$

$(n \times 1)$ matrix

$uc = \vartheta^T \cdot h$

    $n \times |V| \quad n \times 1$

$uc = \vartheta^T h$

$$uc = \vartheta^T W^T x$$

    $n \times |V| \quad n \times |V| \quad |V| \times 1$

    $|V| \times n \quad\quad n \times 1$

input and some $x_b$ is 1 remaining are 0

$\leftarrow x$

    $1 \times n \quad 1 \times n \quad (|V| \times 1)$

     $|V| \times n \quad |V| \times 1$

     $n \times 1$

loss =

$|V| \times 1$

$$uc_j = \vartheta'^T_{w_j} h$$

$j^{th}$ word in vocabulary

one hot vector for content words

$$E = -\sum_{c=1}^{C} u_{j_c}^*$$

value at $|V| \times 1$ $j$th value or simply value of the word.

$$+ \ C \log \sum_{j'=1}^{|V|} e^{-u_{j'}} \ )$$

$$\Big[ \ \underline{\hspace{4cm}} \ \Big]_{|V| \times 1}$$

$$\frac{\partial E}{\partial u_{c,j}} = y_{c,j} - t_{c,j}$$

output that we get    ground truth

$E \cdot I_j$ column vector represent row wise sum of prediction errors across context word panel for current word.

glove (global vector)

co-occurance matrix is written

| $P_{ik}$ | | $k = solid$ | $k = gas$ | $k = water$ | $k = f$ |
|---|---|---|---|---|---|
| $P_{jk}$ | $P(k|ice)$ | $1.9 \times 10^4$ | $6.6 \times 10^{-5}$ | $3.3 \times 10^{-3}$ | $1.7 \times 10^{-5}$ |
| | $P(k)$ steam | $2.2 \times 10^{-5}$ | $7.8 \times 10^{-4}$ | $2.2 \times 10^{-3}$ | $1.8 \times 10^{-5}$ |

$\dfrac{P(k|ice)}{P(k|steam)}$  ......  $8.9$   $8.5 \times 10^{-2}$   $1.36$   $0.96$

$P_{ik} \to$ probably that word $k$ cocurring given $i$
$P_{jk} \to$ probability that word $k$ cocurring $j$
$P_{ik} = \dfrac{X\_ik}{X\_i} \to$ using words $i, k$ together
$\phantom{P_{ik} = } \dfrac{X\_ik}{X\_i} \to$ using $i$ alone in corpus

$k$ is norm coord

$i, j$ content words

$F(wi \; wj \; \bar{wk}) = \dfrac{P_{ik}}{P_{jk}}$

$\boxed{F(wi - wj, \bar{wk}) = \dfrac{P_{ik}}{P_{jk}} \to scalar}$

linear difference between the words.
vector value:

$F((wi - wj)^T \cdot \bar{wk}) = \dfrac{F(wi^T \cdot \bar{wk})}{F(wj^T \cdot \bar{wk})}$

$P_{ik} \searrow \dfrac{X\_ik}{X\_i}$

$\searrow P_{jk} \searrow \dfrac{X\_jk}{X\_j}$

replace $F$ with exponential

$e^{wi^T \cdot \bar{wk}} = P\_ik$

$\boxed{\log(xi) + wi^T \bar{w}k = \log(X\_ik)}$

bias term for $\omega_i$, $\omega_b$

$$\omega_i^T \cdot \bar{\omega}_b + b_i + \bar{b}_b = \log(x_{ib})$$

$$J = \sum_{i,j} f(X_{ij}) \left( \omega_i^T \cdot \bar{\omega}_b + b_i + \bar{b}_b - \log(x_{ib}) \right)$$

$$\underbrace{\qquad}_{\text{weight function}}$$

## heirarchical softmax

$$\omega_1 \quad \omega_2 \quad \omega_3 \quad \omega_4 \, \omega_5$$

$$\sum\sum \ln P(\omega_1|\omega_3) + \ln P(\omega_2|\omega_3) + \ln P(\omega_4|\omega_3) + \ln P(\omega_3|\omega_3)$$

$$P\left(\frac{\omega_i}{\omega_j}\right) \quad \frac{e^{V_{\omega_i}^T V_{\omega_j}}}{\sum\limits_{i=1}^{|V|} e^{V_{\omega_i}^T k_{\omega_j}}}$$



$\omega_1$
$\omega_2$

$|V| \times d_R$

$$\ln P = e^{V_{\omega_i}^T V_{\omega_j}} - \ln \left( \sum\limits_{i=1}^{|V|} e^{V_{\omega_i}^T V_{\omega_j}} \right)$$

$$\frac{a \cdot x}{1-x}$$

$$(a \bar{\omega}^T \omega) \, 7$$

$$\frac{(a \bar{\omega}^T \omega) \, 7}{7}$$

$$(a \bar{\omega}^T \cdot (\omega \cdot \omega)) \, 7$$

$$\frac{d_i \cdot x}{a \cdot x}$$

$$d_i \cdot 9 = a \bar{\omega}^T \omega$$

$$\frac{(a \cdot x) \, a \omega}{9}$$