bat, con$^\alpha$          bat → i

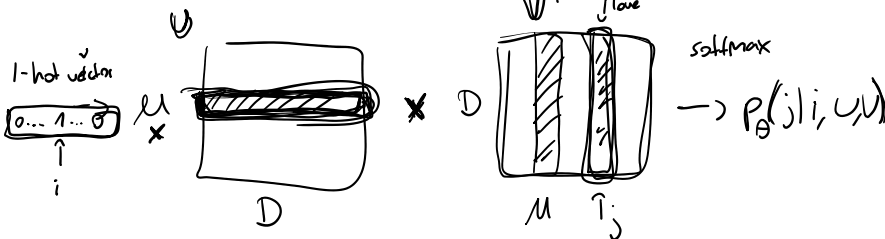$$h_\theta(x^{(t)}, x^{(t-1)}, x^{(t+1)} \ldots)$$

RNN

**Problem 1:** Word2vec defines a mapping from a single word to a single fixed vector. Explain and provide an example why this will not be expressive enough regarding homographs (i.e., words with the same spelling but having more than one meaning). Propose an alternative solution.
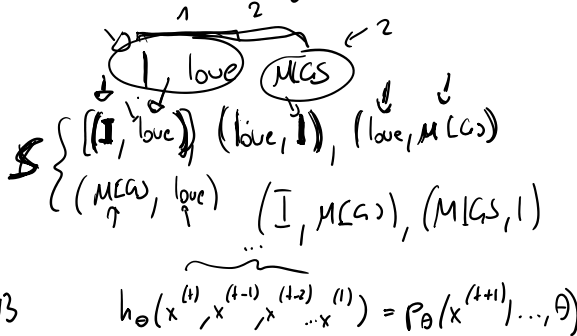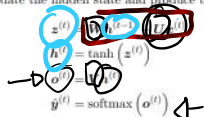
$W^T$    i love

1-hot vector

$\boxed{0 \ldots 1 \ldots 0}$    $\underset{x}{u}$    $\times$    D    softmax

$\uparrow$
$i$

$D$          $M$   $T_j$          $\longrightarrow P_\theta(j \mid i, u, v)$

Training on pairs:          1    2          2
$\overset{\smile}{(I \; love)}$   $(M[CS)$

objective:

$\$\begin{cases} [(I, love)], (love, I), (love, M[C_3)) \\ (M[C_3), love) \quad (I, M[C_3)), (M[C_3, I) \end{cases}$

$\max \prod_{i \in S} P(j \mid i, \theta)$          $\ldots$

$\theta = \{U, V\}$          $h_\theta(x^{(t)}, x^{(t-1)}, x^{(t-2)} \ldots x^{(1)}) = P_\theta(x^{(t+1)} \mid \ldots, \theta)$

**Problem 2:** Given a previous hidden state $h^{(t-1)} \in \mathbb{R}^D$ and a current input $x^{(t)} \in \mathbb{R}^N$, the recurrent neural network equations to update the hidden state and produce the output are:

$$z^{(t)} = U x^{(t)} + W h^{(t-1)}$$
$$h^{(t)} = \tanh(z^{(t)})$$
$$o^{(t)} = V h^{(t)}$$
$$\hat{y}^{(t)} = \mathrm{softmax}(o^{(t)})$$

where parameters $W \in \mathbb{R}^{D \times D}$, $U \in \mathbb{R}^{D \times N}$ and $V \in \mathbb{R}^{M \times D}$ are shared at every step.

To train an RNN we need gradients of loss w.r.t. the parameters: $\partial L/\partial W$, $\partial L/\partial U$ and $\partial L/\partial V$. Your task is to arrive at the equations given on slide 17 in the lecture.

Use the fact that $\partial L/\partial o^{(t)} = \hat{y}^{(t)} - y^{(t)}$, where $y^{(t)}$ is the true output.

*Hint: Since parameters are shared, the total gradient is the sum of the contributions over all the steps. Because of that, it might be easier to introduce copies of parameters, e.g. $W^{(t)}$ – a copy of $W$ at step $t$, calculate $\partial L/\partial W^{(t)}$ and sum over all $t$.*

$$\frac{\partial L}{\partial \theta} = \left( \frac{\partial L}{\partial o^{(t)}} \cdot \frac{\partial o^{(t)}}{\partial \theta} \right)$$
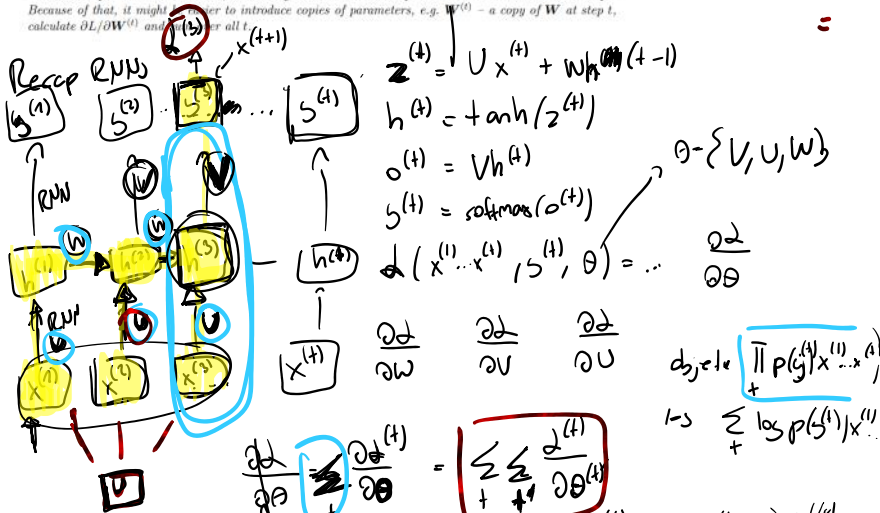
$$z^{(t)} = W h^{(t-1)} + (\text{const in } W)$$

$$L^{(t)} = \hat{y}^{(t)} - y^{(t)} = \frac{\partial L}{\partial o^{(t)}}$$

$=$

Recap RNNs

$s^{(1)}$  $s^{(2)}$  $\boxed{?}$ ... $s^{(t)}$          $x^{(t+1)}$

RNN

$W$        $W$

$h$   $A$  $h^{(2)}$  $h$          $\boxed{h^{(t)}}$

RNN
$U$      $U$      $U$          $\boxed{x^{(t)}}$

$x^{(1)}$  $x^{(2)}$  $x^{(3)}$

$\boxed{U}$

$$z^{(t)} = U x^{(t)} + W h^{(t-1)}$$
$$h^{(t)} = \tanh(z^{(t)})$$
$$o^{(t)} = V h^{(t)}$$
$$s^{(t)} = \mathrm{softmax}(o^{(t)})$$          $\theta = \{V, U, W\}$

$$L(x^{(1)} \ldots x^{(t)}, s^{(t)}, \theta) = \ldots$$          $\frac{\partial L}{\partial \theta}$

$\frac{\partial L}{\partial W}$    $\frac{\partial L}{\partial V}$    $\frac{\partial L}{\partial U}$          objective $\prod_t P(s^{(t)} \mid x^{(1)} \ldots x^{(t)})$

$\to \sum_t \log P(s^{(t)} \mid x^{(1)})$

$$\frac{\partial L}{\partial \theta} = \sum_t \frac{\partial L^{(t)}}{\partial \theta} = \boxed{\sum_t \sum_{t'} \frac{\partial L^{(t)}}{\partial \theta^{(t')}}}$$

$$\frac{\partial L^{(t)}}{\partial \theta} = \quad \frac{\partial L^{(t)}}{\partial U} = \sum \frac{\partial L^{(t)}}{\partial U(\text{path})} = \sum_{t'} \frac{\partial L^{(t)}}{\partial U^{(t')}} = \frac{\partial L^{(t)}}{\partial U^{(t)}} \quad \frac{\partial L^{(t)}}{\partial U^{(3)}} = 0$$

$\mathbb{R}^n$ paths

(1) $V$:  $\frac{\partial L}{\partial o^{(t)}} \frac{\partial o^{(t)}}{\partial V}$          $\mathbb{R}^{n \times m}$

before:  $\frac{\partial o^{(t)}}{\partial V}$

$$= \frac{\partial}{\partial V}(V h^{(t)})$$

Recap:  $S: \mathbb{R}^M \to \mathbb{R}$   $S(x)$          $f: \mathbb{R}^{m \times o} \to \mathbb{R}^n$   $f(X)$          $\uparrow$  $x$

$$\frac{\partial S}{\partial x} = J_S^T = \begin{pmatrix} \frac{\partial S_1}{\partial x_1} & \cdots & \frac{\partial S_1}{\partial x_m} \\ \vdots & & \vdots \\ \frac{\partial S_n}{\partial x_1} & & \frac{\partial S_n}{\partial x} \end{pmatrix}^T \quad \left[ \frac{\partial f}{\partial X} \right]_{ijh} = \frac{\partial f_i}{\partial X_{jh}}$$          equals?

$$\frac{\partial g}{\partial x} = J_g^T = \begin{pmatrix} \frac{\partial g_1}{\partial x_1} & \cdots & \frac{\partial g_1}{\partial x_m} \\ \frac{\partial g_n}{\partial x_1} & \cdots & \frac{\partial g_n}{\partial x_m} \end{pmatrix} \qquad \left[\frac{\partial f}{\partial X}\right]_{ijh} = \frac{\partial f_i}{\partial X_{jh}}$$

$$\frac{\partial f_i}{\partial X_{jh}} \qquad X_{jk}: \quad \frac{\partial f}{\partial X_{jh}} = \sum_i \frac{\partial f_i}{\partial X_{jh}} \qquad \sum_i \boxed{\frac{\partial o^{(t)}_i}{\partial V_{jk}}}$$

$$\boxed{\frac{\partial o^{(t)}_k}{\partial V_{ij}}} = \frac{\partial}{\partial V_{ij}}(Vh^{(t)})_k = \frac{\partial}{\partial V_{ij}}\left[\sum_l V_{kl} h^{(t)}_l\right]_k = \sum_l \frac{\partial}{\partial V_{ij}} V_{kl} h^{(t)}_l$$

$$= \sum_l h^{(t)}_l \frac{\partial}{\partial V_{ij}} V_{kl} = \sum_{l\neq j} h^{(t)}_l \frac{\partial}{\partial V_{ij}} V_{kl} + h^{(t)}_j \frac{\partial}{\partial V_{ij}} V_{kj}$$

$$= \begin{cases} 0 & \text{if } i\neq k \text{ or } j\neq l \\ 1 & \text{if } i=h \text{ and } j=l \end{cases} \qquad = h^{(t)}_j \frac{\partial}{\partial V_{ij}} V_{ij}$$

$$\frac{\partial o^{(t)}}{\partial V_{ij}} = \sum_h \frac{\partial o^{(t)}_h}{\partial V_{ij}} = \sum_h h^{(t)}_j \frac{\partial}{\partial V_{ij}} V_{ij}$$

$$= \sum_{h\neq i} h^{(t)}_j \frac{\partial}{\partial V_{ij}} V_{hl} + h^{(t)}_j \frac{\partial}{\partial V_{ij}} V_{ij} = \boxed{h^{(t)}_j} \cdot \boxed{\frac{\partial o^{(t)}}{\partial V_{ij}}}$$

$$=0 \qquad =1 \qquad \boxed{= 1 \ \text{if } i=h}$$

$$\boxed{\frac{\partial d}{\partial o^{(t)}_i}} = \boxed{\frac{\partial d}{\partial o^{(t)}}} \quad \boxed{\frac{\partial o^{(t)}}{\partial V_{ij}}} = \sum_h \frac{\partial d}{\partial o^{(t)}_h} \frac{\partial o^{(t)}_h}{\partial V_{ij}} = \sum_h \delta^{(t)}_h h^{(t)}_j \frac{\partial}{\partial V_{ij}} V_{ij}$$

$$\frac{\partial d^{(t)}}{\partial V} = \delta^{(t)}(h^{(t)})^T \qquad \boxed{\delta^{(t)}_i h^{(t)}_j}$$

$$\frac{\partial d}{\partial V} = \sum_t \frac{\partial d^{(t)}}{\partial V} = \sum_t \delta^{(t)}(h^{(t)})^T \qquad \delta' \qquad J = \hat{\delta}^{(t)} = \delta^{(t)} = \frac{\partial d}{\partial o^{(t)}}$$

$$z^{(t)} = W h^{(t-1)}$$

$$\boxed{\frac{\partial d}{\partial W}} = \sum_t \frac{\partial d}{\partial W^{(t)}} = \sum_t \boxed{\frac{\partial d}{\partial h^{(t)}_i}} \; \boxed{\frac{\partial h^{(t)}}{\partial z^{(t)}}} \; \boxed{\frac{\partial z^{(t)}}{\partial W_{hk}}}$$

i) $\frac{\partial z^{(t)}}{\partial W} = \frac{\partial}{\partial W}(Wh^{(t-1)} + \cancel{\cdots} W) = \frac{\partial}{\partial W}(Wh^{(t-1)})$ \qquad but we know what "effect this has"

$$\sum_t = \sum_t \boxed{\frac{\partial d}{\partial h^{(t)}} \frac{\partial h^{(t)}}{\partial z^{(t)}}} (h^{(t-1)})^T$$

ii) Now look at: $\boxed{\frac{\partial d}{\partial h^{(t)}}} \boxed{\frac{\partial h^{(t)}}{\partial z^{(t)}}} \xrightarrow{\mathbb{R}^n}$ \qquad as before: $\frac{\partial d}{\partial h^{(t)}_i} \boxed{\frac{\partial h^{(t)}_i}{\partial z^{(t)}_j}}$ \qquad $= 1 \text{ if } j=i$
\qquad $= 0 \text{ otherwise}$

$$\boxed{\frac{\partial h^{(t)}_i}{\partial z^{(t)}_j}} = \frac{\partial}{\partial z^{(t)}_j}\left[\tanh(z^{(t)})_i\right] = \frac{\partial}{\partial z^{(t)}_j}\left[\tanh(z^{(t)}_i)\right] = \left[\frac{\partial}{\partial x}\tanh(x)\right](z^{(t)}_i)\left(\frac{\partial}{\partial z^{(t)}_j} z^{(t)}_i\right)$$

$$= \begin{cases} \left[\frac{\partial}{\partial x}\tanh(x)\right](z^{(t)}_i) \cdot 1 & \text{if } i=j \\ 0 & \text{otherwise} \end{cases}$$

equis: $\hat{B}^{-1} x$

$\delta((\cdots)^{(t)}) \leftarrow \frac{\partial d}{\partial o^{(t)}}$

If smth of te form

| $\frac{\partial d}{\partial a}$ | $\frac{\partial a}{\partial B}$ |
|---|---|
| $\frac{\partial d}{\partial a}$ | $X$ |

$\boxed{\frac{\partial}{\partial a} = B \frac{\partial}{\partial x}}$

$$\frac{d}{dx}\tanh x = 1 - (\tanh x)^2$$
$$f \qquad 1 - f^2$$

$$\begin{cases} 0 & \text{otherwise} \end{cases}$$

$$= \begin{cases} 1 - \left(\tanh z_i^{(t)}\right)^2 & \text{if } i=j \\ 0 & \text{othw} \end{cases} = \begin{cases} 1 - \left(h^{(t)}_i\right)^2 & \text{if } i=j \\ 0 & \text{othw} \end{cases}$$

$$\top \quad 1 - f^2$$

$$\frac{\partial l}{\partial z_j^{(t)}} \quad \frac{\partial l}{\partial h}^{(t)} \quad \frac{\partial h^{(t)}}{\partial z_j^{(t)}} = \sum_i \frac{\partial l}{\partial h_i^{(t)}} \frac{\partial h_i^{(t)}}{\partial z_j^{(t)}} = \sum_{i \neq j} \frac{\partial l}{\partial h_i^{(t)}} \frac{\partial h_i^{(t)}}{\partial z_j^{(t)}} + \frac{\partial l}{\partial h_j^{(t)}} \frac{\partial h_j^{(t)}}{\partial z_j^{(t)}}$$

$$= 0 \qquad \underbrace{\phantom{xxx}}_{=0} \qquad \underbrace{\phantom{xxx}}_{= 1 - \left(h_j^{(t)}\right)^2}$$

$$= 1 - \left(h_j^{(t)}\right)^2 \cdot \frac{\partial l}{\partial h_j^{(t)}}$$

$$\frac{\partial l}{\partial z^{(t)}} = \text{diag}\left(1 - \left(h^{(t)}\right)^2\right) \frac{\partial l}{\partial h^{(t)}}$$

$$\sum_t \frac{\partial l}{\partial h^{(t)}} \frac{\partial h^{(t)}}{\partial z^{(t)}} \left[h^{(t-1)}\right]^\top = \sum_t \text{diag}\left(1 - \left(h^{(t)}\right)^2\right) \frac{\partial l}{\partial h^{(t)}} \left[h^{(t-1)}\right]^\top = \frac{\partial l}{\partial W}$$

$$\frac{\partial l}{\partial U} : \text{ Is the same derivation as for } \frac{\partial l}{\partial W} \text{ but we interchange } W \text{ and } U$$
$$h^{(t-1)} \text{ and } x^{(t)}$$

$$= \sum_t \text{diag}\left(1 - h^{(t)}\right)^2 \frac{\partial l}{\partial h^{(t)}} \left[x^{(t)}\right]^\top$$

**Problem 3:** What do you need to change in the equations that you got in the previous exercise if the output $o^{(t)}$ is used as an input to another neural network?

So, then $\exists$ $f$ as another Layer ($\cup$) in between

$$\frac{\partial l}{\partial \theta} = \left(\frac{\partial l}{\partial f}\right) \frac{\partial f}{\partial \theta} =$$