## Vanishing and exploding gradient problem

Suppose that a computational graph contains a path that consists of repeatedly multiplying by a matrix $W$. After $t$ steps, this is equivalent to multiplying by $W^t$. Assuming $W$ has an eigendecomposition $W = V \, \text{diag}(\lambda) \, V^{-1}$ then.
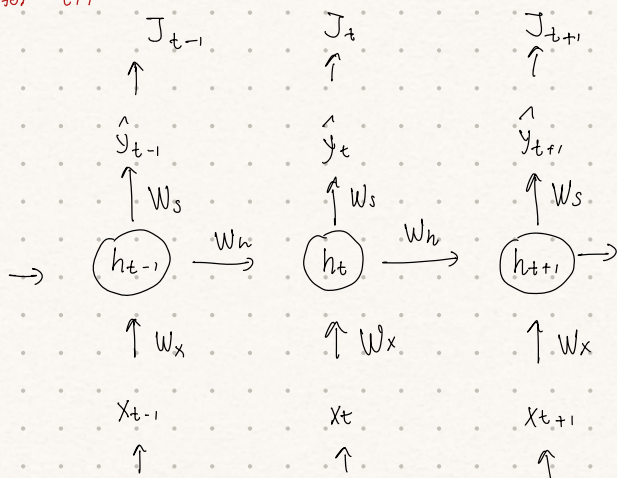
$$W^t = V \, \text{diag}(\lambda)^t \, V^{-t}$$

Therefore, any eigenvalues $\lambda_i$ that are not near $1$ in absolute value will either explode if they are greater than $1$ or vanish if they are less than $1$ in magnitude.

Vanishing gradients make it difficult to know which direction the parameters should move to improve the cost function, while exploding gradients can make learning unstable.

This problem is more relevant in RNNs than FFN because RNNs use the same matrices of parameters at each time step but FFN don't.

### RNN for LM



where $\quad h_t = \sigma\left( W_h \, h_{t-1} + W_x \, x_t \right)$

$\hat{y}_t = \text{softmax}\left( W_s \, h_t \right)$

Cross-Entropy loss at $t$: $\qquad J^{(t)}(\theta) = -\sum_{j=1}^{|V|} y_{t,j} \times \log \hat{y}_{t,j} \qquad \overset{OHE}{=} -\log \tilde{\hat{y}}_{t,j}$

Total loss: $\qquad J(\theta) = \dfrac{1}{T} \sum_{t=1}^{T} J^{(t)}(\theta) \qquad , \qquad \theta$ contains all parameters.

$$= -\dfrac{1}{T} \sum_{t=1}^{T} \sum_{j=1}^{|V|} y_{t,j} \times \log \hat{y}_{t,j}$$

### Backpropagation through time (BPTT)

Notation $\dagger$

$$\dfrac{\partial J}{\partial \theta} = \sum_{t=1}^{T} \dfrac{\partial J^{(t)}}{\partial \theta} \qquad , \qquad \dfrac{\partial J^{(t)}}{\partial \theta} = \sum_{k=1}^{t} \dfrac{\partial J^{(t)}}{\partial \hat{y}_t} \dfrac{\partial \hat{y}_t}{\partial h_t} \dfrac{\partial h_t}{\partial h_k} \dfrac{\partial h_k}{\partial \theta} \qquad \text{for all previous } k \text{ time-steps}$$

$$\dfrac{\partial J}{\partial \theta} = \sum_{t=1}^{T} \sum_{k=1}^{t} \dfrac{\partial J^{(t)}}{\partial \hat{y}_t} \dfrac{\partial \hat{y}_t}{\partial h_t} \dfrac{\partial h_t}{\partial h_k} \dfrac{\partial h_k}{\partial \theta}$$

where $\quad \dagger \; \dfrac{\partial h_t}{\partial h_k} = \overset{t}{\underset{j=k+1}{\prod}} \dfrac{\partial h_j}{\partial h_{j-1}} = \overset{t}{\underset{j=k+1}{\prod}} \dfrac{\partial h_j}{\partial \tilde{h}_j} \dfrac{\partial \tilde{h}_j}{\partial h_{j-1}} \quad , \quad \tilde{h}_j = W_h \, h_{j-1} + W_x \, x_j$

$$\frac{\partial h_t}{\partial h_k} = \prod_{j=k+1}^{t} W_h^T \, diag\left(\sigma'\left(h_{j-1}\right)\right)$$ and transport error from step $t$ back to step $k$

Since $h \in \mathbb{R}^{d_h}$, hence $\frac{\partial h_j}{\partial h_{j-1}}$ is the Jacobian matrix for $h_j$ wrt all the units of $h_{j-1}$

Notice: product of $(t-k)$ Jacobian matrices → exploding grads

vanishing

Gradient of loss:

$$\frac{\partial J}{\partial \theta} = \sum_{t=1}^{T} \sum_{k=1}^{t} \frac{\partial J^{(t)}}{\partial \hat{y}_t} \frac{\partial \hat{y}_t}{\partial h_t} \left( \prod_{j=k+1}^{t} \frac{\partial h_j}{\partial h_{j-1}} \right) \frac{\partial h_k}{\partial \theta}$$

Assuming $W^T$ and $diag(\sigma'(h_{j-1}))$ have bounded norms given by $\beta_w$ and $\beta_h$ respectively.

$$\left\| \frac{\partial h_j}{\partial h_{j-1}} \right\| \leq \| W_h^T \|_2 \, \| diag(\sigma'(h_{j-1})) \|_2 \leq \beta_w \beta_h$$

$\beta_w$ and $\beta_h$ bound the largest eigenvalue of the matrices.

↑ 2-norm definition.

Therefore,

$$\left\| \frac{\partial h_t}{\partial h_k} \right\| \leq (\beta_w \beta_h)^{t-k} \begin{array}{c} \nearrow 0 \\ \searrow \infty \end{array}$$ provided $t \gg k$   long term

Mikolov: sufficient $\lambda_1(w) < \frac{1}{\gamma} \Rightarrow$ vanishing $\nabla$.

Parameters update:

$$\theta \leftarrow \theta - \eta \nabla_\theta L$$

necessary $\lambda_1(w) > \frac{1}{\gamma} \Rightarrow$ exploding $\nabla$

→ $\theta$ too large

⇒ no learning

because $\frac{\partial J^{(t)}}{\partial h_t} \frac{\partial h_t}{\partial h_k} = \sum_i c_i \lambda_i^\ell q_i^T$

$\approx c_j \lambda_j^\ell q_j^T$

$\left| \frac{\lambda_i}{\lambda_j} \right| < 1 \quad \forall i \neq j$

$\lambda_t = \lambda_j$

Gradient Explosion: detectable at run time → overflow

Vanishing Gradients: undetectable and reduces the learning quality of the model for far away words

→ model unable to capture relation between words at steps $t$ and $t+k$

(Mikolov) Gradient Clipping

Clip the gradients to a small number whenever they explode

$$\hat{g} = \frac{\partial L}{\partial \theta} \longrightarrow \frac{\varepsilon}{\| \hat{g} \|} \hat{g}$$

Solution for Vanishing Gradients

1) $W_h^{(o)} = I$ instead of a random matrix
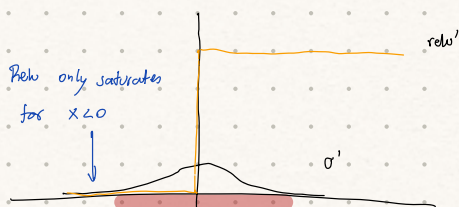
2) $\sigma = ReLU$ instead of sigmoid.

$$ReLU(x) = \begin{cases} 0 & , \ x < 0 \\ x & , \ x > 0 \end{cases}$$

$$\frac{d \, sigmoid(x)}{dx} = \sigma(x)(1 - \sigma(x))$$ vanishes gradients for $\sigma(x)$ close to 0 or 1

$$\frac{d \, ReLU(x)}{dx} = \begin{cases} 0 & , \ x < 0 \\ 1 & , \ x > 0 \end{cases}$$

$$\frac{d \, ReLU(x)}{dx} = \begin{cases} 0 & , \ x < 0 \\ 1 & , \ x > 0 \end{cases} = \mathbb{1}_{\mathbb{R}^+}(x)$$

Obs: training w/ ReLU is faster than sigmoid.



ReLU only saturates for $x < 0$

relu'

$\sigma'$

limited input interval w/ non-zero gradients → meaningless BP.

Recall: $h^{(+)} = \sigma(W_h h_{t-1} + W_x x_t)$

If $\sigma(x) = x$ then:

$$\frac{\partial h_t}{\partial h_{t-1}} = W_h^T \, \text{diag}\left(\sigma'(W_h h_{b-1} + W_x x_t)\right)$$

where $h_t$
$$\shortparallel$$

$$= W_h^T \cdot I = W_h^T$$

$$\uparrow$$
$$\sigma(x) = x$$

Hence,

$$\frac{\partial J^{(i)}(\theta)}{\partial h_j} = \frac{\partial J^{(i)}(\theta)}{\partial h_i} \prod_{t=j+1}^{i} \frac{\partial h_t}{\partial h_{t-1}}$$

$$= \frac{\partial J^{(P)}(\theta)}{\partial h_\ell} \left(W_h^T\right)^\ell \quad , \quad \ell = i - j$$

$$= \sum c_i \lambda_i^\ell q_i \quad \begin{array}{c} \nearrow 0 \\ \searrow \\ \infty \end{array}$$

$\lambda_i$ eigenvalues of $W_h$
$q_i$ eigenvectors.

Mikolov: in the linear case it is sufficient for the largest eigenvalue $\lambda_1$ to be smaller than 1 for long-term components to vanish ($t \to \infty$) and necessary for it to be larger than 1 for gradients to explode.