

Why GBDT?

The new tree corrects the error of previous trees, but not fully.

$$l(y_i, \hat{y}_i^{k-1} + f_k(x_i)) \approx \underline{l(y_i, \hat{y}_i^k)} + \underline{g_i f_k(x_i) + \frac{1}{2} h_i f_k^2(x_i)}$$

previous loss      correction of the new tree.

only find local min in every tree

decrease the loss

$$\hat{y}_i^k = \hat{y}_i^0 + \underline{\text{something}}$$

$$\hat{y}_i^k = \hat{y}_i^0 + \eta \frac{\partial L(y_i, \hat{y}_i^0)}{\partial \hat{y}_i^0}$$

we fit a new model to generate  $g_i^k$

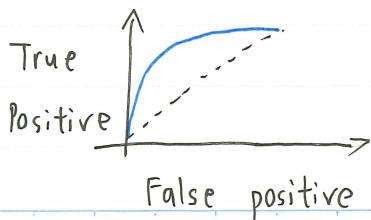
$$\therefore \hat{y}_i^k = \hat{y}_i^0 + f_1(x) + f_2(x) + \dots + f_K(x)$$

$$w_j^k = \frac{-g_j}{h_j + 1}$$

It's essentially the average of  $g_i$ 's inside every leaf

Second order term can be treated as a scaling factor that impact how far we go in the gradient direction based on curvature

ROC curve (receiver operating characteristic)



A good classifier stays as far away from diagonal line as possible

## Deep Neural Network (DNN)

$$\frac{\partial L}{\partial w_{ij}^{[2]}} = \frac{\partial L}{\partial a^{[3]}} \frac{\partial a^{[3]}}{\partial z^{[3]}} \frac{\partial z^{[3]}}{\partial a^{[2]}} \frac{\partial a^{[2]}}{\partial z^{[2]}} \frac{\partial z^{[2]}}{\partial w_{ij}^{[2]}}$$

$$(a^{[3]} - y) \quad W^{[3]} \quad \text{diag}(g'(z^{[2]})) \quad a_j^{[2]} e_i \\ 1 \times 1 \quad 1 \times 2 \quad 2 \times 2 \quad 2 \times 1$$

$$= (a^{[3]} - y) W^{[3]} \circ g'(z^{[2]}) a_j^{[2]} e_i$$

$$= [(a^{[3]} - y) W^{[3]} \circ g'(z^{[2]})]_j a_j^{[2]}$$

hadamard product (elementwise product)

In matrix form:

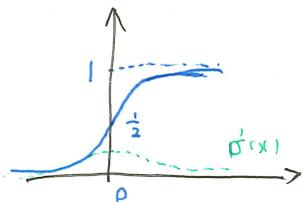
$$f(x) = \sigma^{[L]}(W^{[L]} \sigma^{[L-1]}(W^{[L-1]} \sigma^{[L-2]}(W^{[L-2]} \dots \sigma^{[2]}(W^{[1]} + b^{[1]}) \dots + b^{[L-2]}) + b^{[L-1]}) + b^{[L]}$$

a linear transformation followed by a non-linear transformation

if activation function is linear, then  $f(x) = w'x + b$  is a linear model

## Activation Functions

## Sigmoid



- ✓ Maps to probability
- ✓ Logistic regression

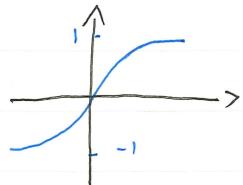
- not zero centered → bias

- saturate locally at extremes (grad = 0) → no flow and vanishing gradients

- Verdict: don't use in non-output layers

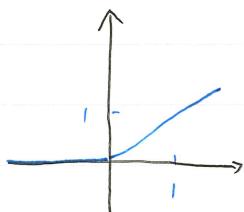
## Tangens hyperbolics ( $\tanh$ )

- ✗ Maps to probability
- ✗ logistic regression
- ✗ saturates ( $\text{grad} = 0$ ), no flow
- ✓ zero-centered



**Verdict:** can use in non-output layers, but not recommended

## ReLU



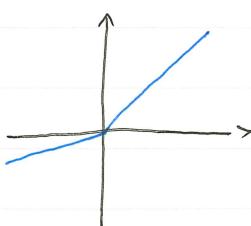
- non-linear due to clipping
- no expensive functions
- doesn't saturate on right side
- ✗ not zero-centered
- ✗ neurons can get stuck: requires increased tuning of learning rate

**Verdict:** default

## Leaky ReLU:

$$g(z) = \max(\epsilon z, z)$$

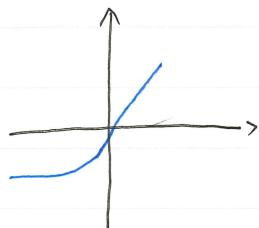
with  $\epsilon \ll 1$



## ELU

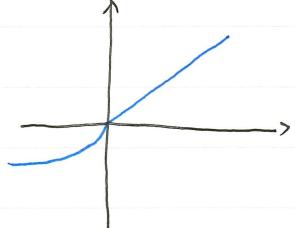
$$g(z) = \max(\alpha(e^z - 1), z)$$

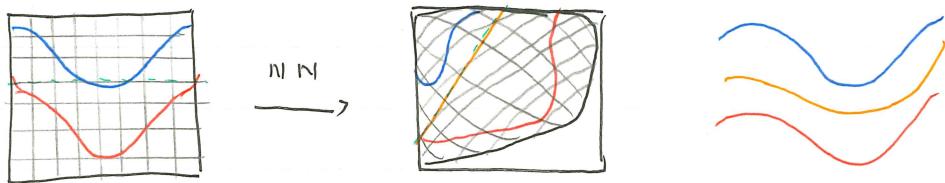
with  $\alpha \ll 1$



## SELU

$$\begin{cases} f(x) = \lambda x & x > 0 \\ f(x) = \lambda \alpha(e^x - 1) & x \leq 0 \end{cases}$$



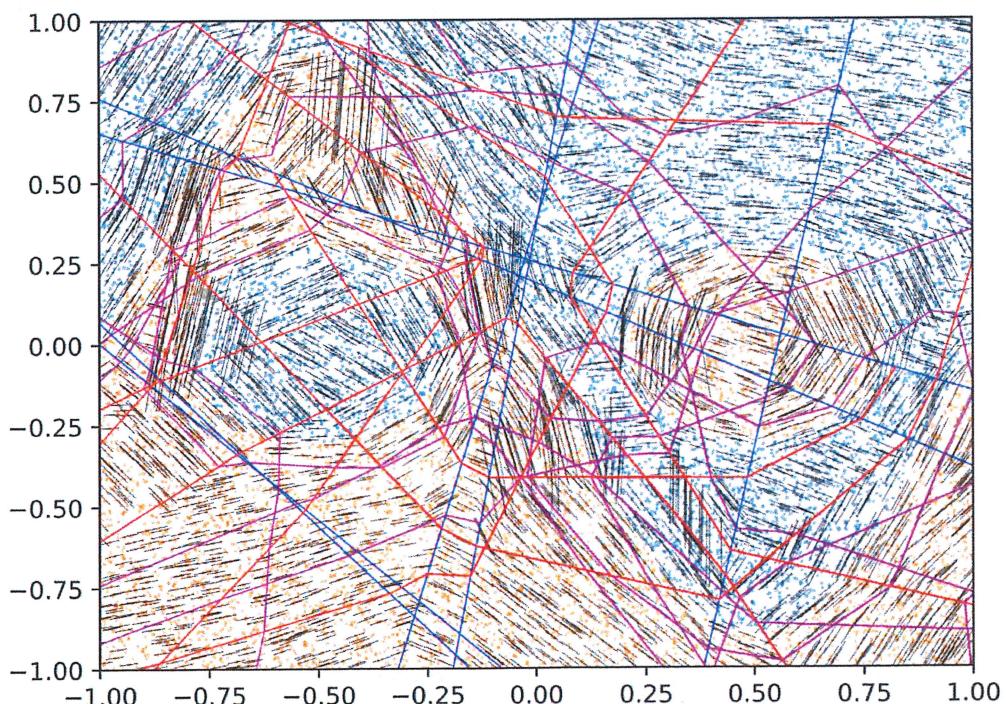


Universality Theorem

Any continuous function  $f$   
 $f: \mathbb{R}^n \rightarrow \mathbb{R}^m$

can be realized by a network with one hidden layer  
 given enough hidden neurons

ReLU DNNs Partitions Input Space      3 layers x 8 neurons



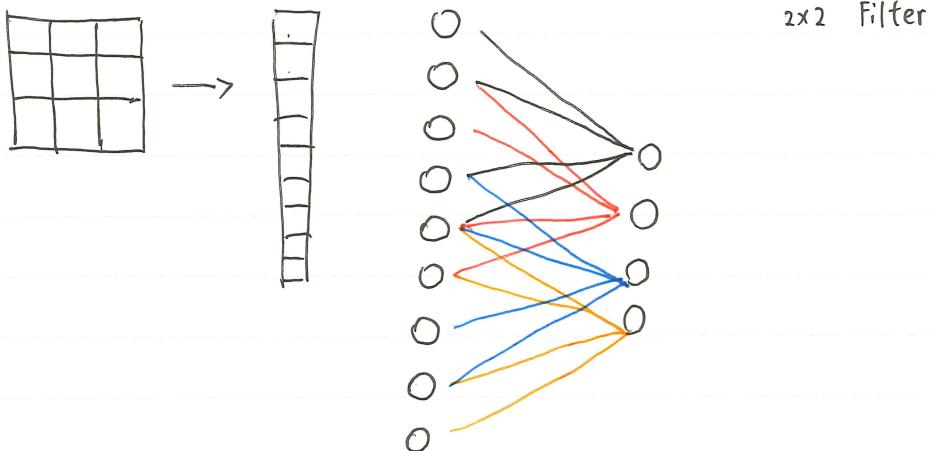
Blue: Linear regions (convex polyhedrons) generated by 1<sup>st</sup> layer weights

Red: 2<sup>nd</sup>

Purple: 3<sup>rd</sup>

### Convolutional Layer:

CNN is a sparsely connected DNN.



$$\text{output size} = \frac{N + 2P - F}{\text{stride}} + 1$$

Image  $\rightarrow$  CONV  $\rightarrow$  Max Pooling  $\rightarrow$  Flatten  $\rightarrow$  FC

$\underbrace{\hspace{10em}}$   
 $\times N$

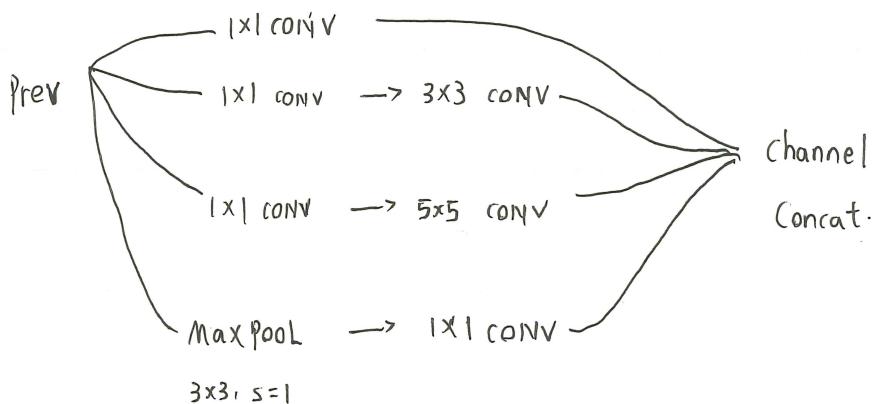
- In CNN:
  - first layers usually acts as edge detectors
  - deeper feature maps encode high level features

## VGG Net

smaller filters are cheaper.

- stack of  $3 \times 3$  filters  $\times 3$  has same receptive field as  $7 \times 7$

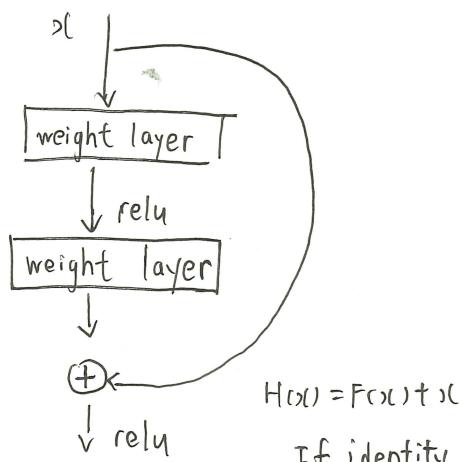
## Inception Net.



## ResNet - Motivation

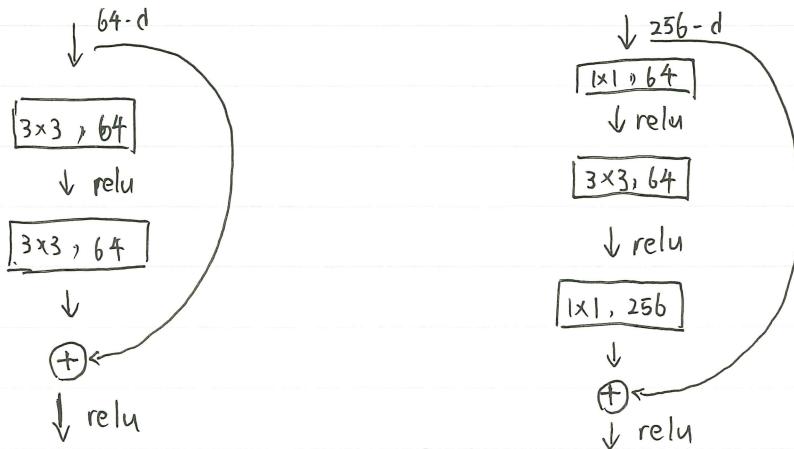
Degradation problem: adding more layers to a suitably deep model  $\rightarrow$  higher training error

Optimization difficulties: optimizers can't find the solution when going deeper



$$H(x) = F(x) + x$$

If identity were optimal, easy to set weights as 0



all  $3 \times 3$   $\leftarrow$  similar complexity  $\rightarrow$  bottleneck for Res 50/101/152

## Useful tricks

- LR scheduler : (multi)step LR, cosine Anneal LR

Dropout: During training, randomly remove neurons, by a factor  $\frac{1}{1-p}$

During inference ① Use all neurons

② still dropout, but do ensembling over multiple runs

prevent co-adaptation of neurons, (similar to random forest)

## Data Augmentation

- ## - Mixup

Batch norm

Address covariance shift: Layers' scales change based on previous layers.

↑  
M samples

Features →

+	-	+	-	+	-	+	-	+	-
+	-	+	-	+	-	+	-	+	-
+	-	+	-	+	-	+	-	+	-
+	-	+	-	+	-	+	-	+	-
+	-	+	-	+	-	+	-	+	-

$\mu_i = \frac{1}{M} \sum A_i$   
 $\sigma_i^2 = \sqrt{\frac{1}{M} \sum (A_i - \mu_i)^2}$

→  $\hat{A}_i = \frac{A_i - \mu_i}{\sigma_i} \rightarrow \tilde{B}_i \tilde{N}_i = \gamma \odot \hat{A}_i + \beta \rightarrow$

↓

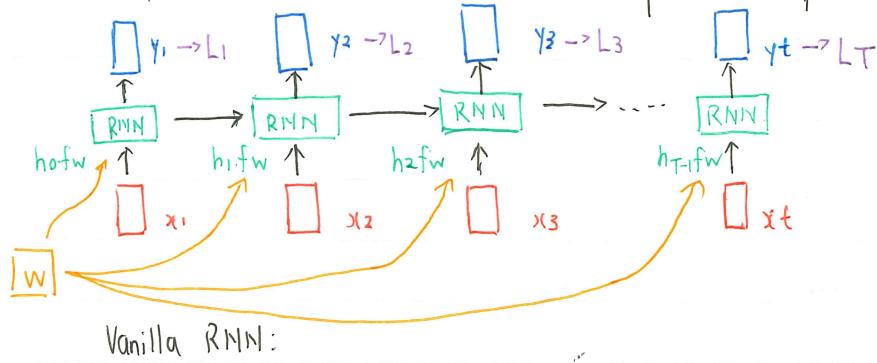
Moving Average

$\mu_{movi} = \alpha \mu_{movi} + (1-\alpha) \mu_i$   
 $\sigma_{movi}^2 = \alpha \sigma_{movi}^2 + (1-\alpha) \sigma_i^2$

KOKUYO

## Recurrent Neural Network

update an internal state as a sequence is processed



$$h_t = \tanh(W_{hh}h_{t-1} + W_{xh}x_t)$$

$$y_t = W_{hy}h_t$$

Many to one: Encode input sequence in a single vector

One to many: Produce output sequence from single vector

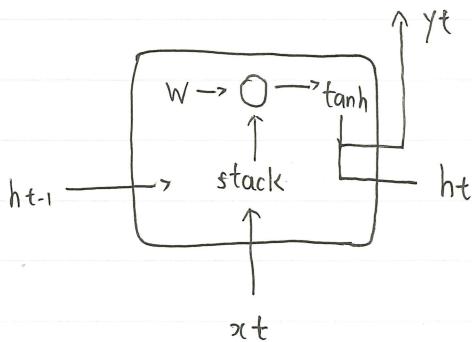
sequence 2 sequence

Advantage:

- Can process any length input.
- Computation for step  $t$  can use information from many steps back
- Model size doesn't increase for larger input.

Disadvantages:

- Recurrent computation is slow.



Backpropagation through time

$$\frac{\partial J_T}{\partial w} = \sum_{i=1}^T \frac{\partial J_T}{\partial w} \Big|_{(i)}$$

$$\approx h_t \approx \tanh(W_{hh}h_{t-1} + W_{xh}x_t)$$

$$\therefore \frac{\partial h_t}{\partial h_{t-1}} = \text{diag}[\tanh'(W_{hh}h_{t-1} + W_{xh}x_t)] W_{hh}$$

$$\therefore \frac{\partial J_T}{\partial h_t} = \frac{\partial J_T}{\partial h_T} W_{hh}^{T-t} \prod_{t < i \leq T} \text{diag}[\tanh'(W_{hh}h_{i-1} + W_{xh}x_i)]$$

To prevent gradient explosion: Gradient Clipping

Before SGD scale down gradient if its norm is greater than some threshold  
 $\uparrow$   
 update

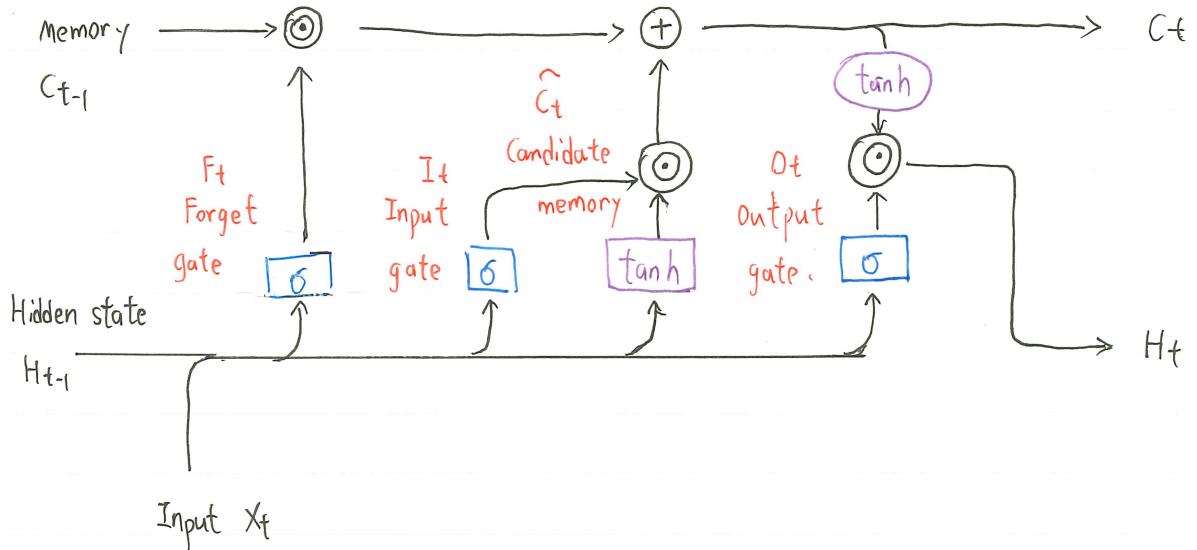
if  $\|g\| \geq \text{threshold}$  then

$$g \leftarrow \frac{\text{threshold}}{\|g\|} g$$

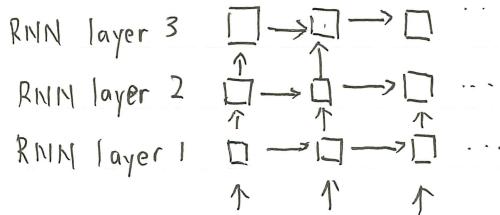
Simple RNNs have bad long-term memory

Solution: LSTM. GRU (gated recurrent unit)

## Long short-term memory (LSTM)



## Multi-layer RNNs

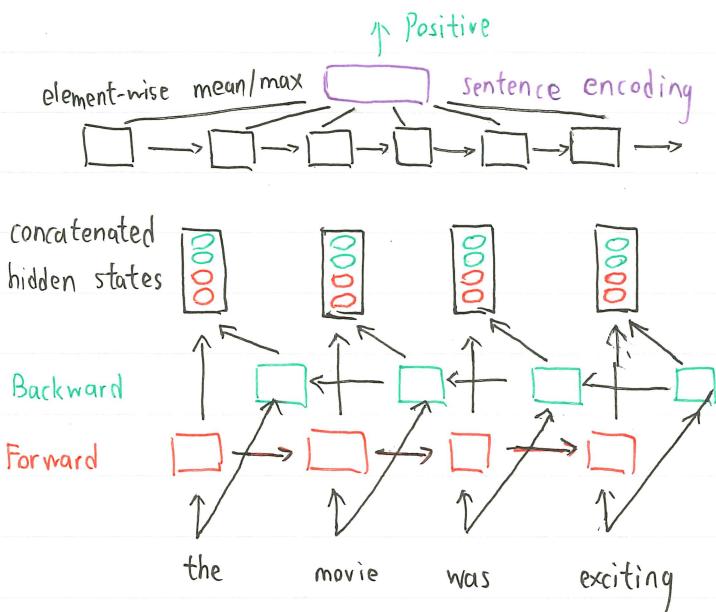


the movie was

High-performing RNNs are often multilayer (not as deep as CNN )

However, this only contain information about the left context

Bidirectional RNN :



$$\begin{aligned} \text{Forward RNN} \quad \vec{h}_t &= \text{RNN}_{\text{FW}}(\vec{h}_{t-1}, x_t) \\ \underline{h}_t &= \text{RNN}_{\text{BW}}(\underline{h}_{t-1}, x_t) \end{aligned}$$

$$h_t = [\vec{h}_t, \underline{h}_t]$$

Encoder  $\xrightarrow{\text{hidden state}}$  Decoder

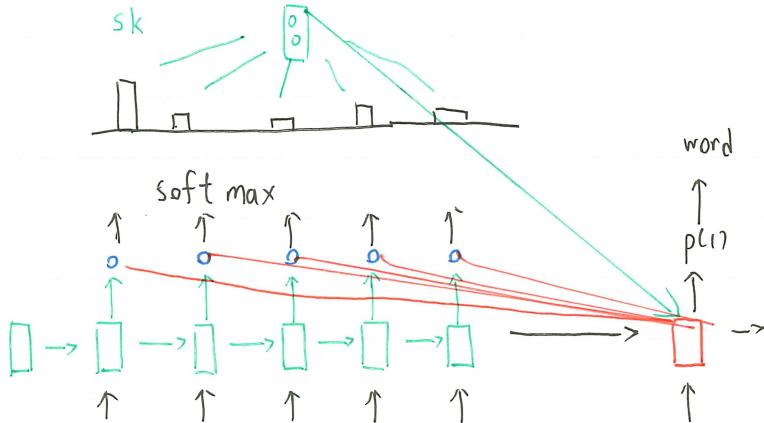
Bottleneck Problem  $\xrightarrow{\quad}$  solved by attention

Attention:

score( $h_t, s_k$ )

scalar out

Attention function:

How relevant is source token  $k$  for target step  $t$ ?Decoder state at step  $t$ :  $h_t$ Encoder state for token  $k$ :

Source sentence

$$a_K^{(t)} = \frac{\exp(\text{score}(h_t, s_k))}{\sum_{i=1}^m \exp(\text{score}(h_t, s_i))}, \quad (k=1, \dots, m)$$

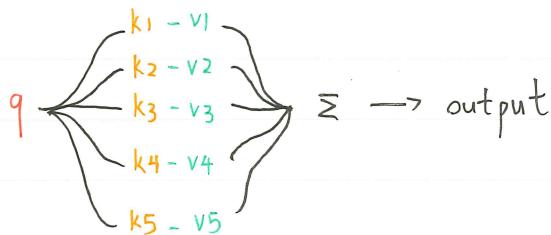
$$c^{(t)} = a_1 s_1^{(t)} + a_2 s_2^{(t)} + \dots + a_m s_m^{(t)} = \sum_{i=1}^m a_i^{(t)} s_i$$

Source context for decoder step  $t$

Finally, concatenate with decoder state and pass on to output layer

$$\tilde{h}_t = \tanh h_t^c w_c [ct; h_t^{dec}] \in \mathbb{R}^h, w_c \in \mathbb{R}^{2h \cdot h}$$

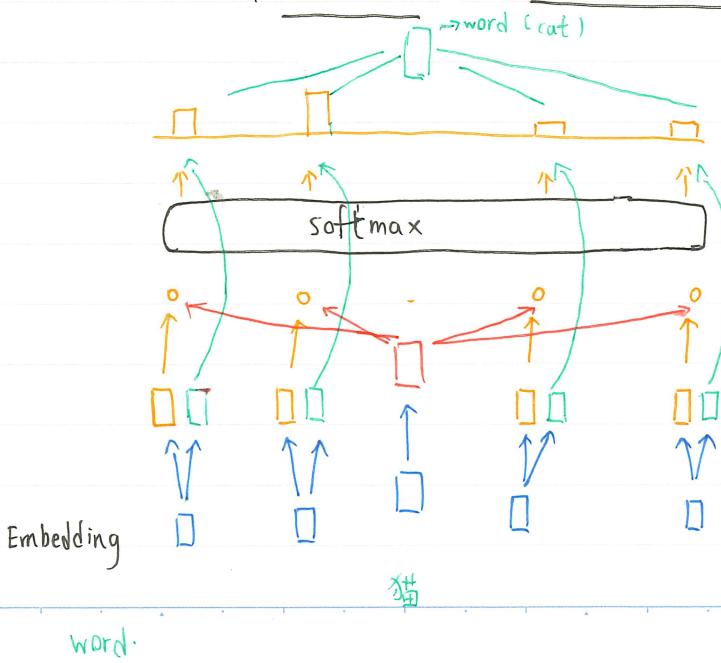
Intuition: The **query** matches to all **keys** softly to a weight  $\in [0, 1]$   
 The **key's value** are multiplied by weights and summed



Self attention:

$$\text{Attention}(q, k, v) = \underbrace{\text{softmax}\left(\frac{q k^T}{\sqrt{K}}\right)}_{\text{Attention weights}} V$$

From each state to All other tokens in the sequence



## Recommendation Systems

A set of tuples  $(u, i, r, t), r_{u,i}^t$

For a linear model.

$$\text{rating (user, movie)} = f(\phi_u(u), \theta_u) + f(\phi_i(i), \theta_i)$$

To recommend a movie to a user

$$\underset{i \in \text{candidates}}{\operatorname{argmax}} : f(u) + \boxed{f(i)} \quad \text{not personalized}$$

The key is to model interactions between users and items