

기계학습
장르 예측

CNN

RNN

LSTM GRU.

transformer

→ motivation

actual computation (to GRU
RNN
Trans)

AE VAE GAN.

Convolution Neural Network

Convolution \Rightarrow spatial filtering

kernel. \Rightarrow input = m kernel size = k stride = s output size = $\frac{(m-k)}{s} + 1$

padding : zeroed frame. border effects

$$\text{Compute Output Shape} = \left(\frac{H+2P-F}{S} + 1 \right) \times \left(\frac{W+2P-F}{S} + 1 \right)$$

H: input height W: input width
P: padding size F: filter size S: stride

Conv properties preserve signal support structure

sparse interaction (local connectivity)

parameter sharing

equivariance to translation $\text{conv}(\text{shift}(\text{input})) = \text{shift}(\text{conv}(\text{input}))$

high K \Rightarrow larger receptive field (local region of input)

Pooling : replace output of layer at certain location with summary stat. of nearby outputs.

smaller representation, more manageable

\Rightarrow max pooling, average pooling, ...

invariant to small translations of the input. small change of input \Rightarrow most pooled output no change.

Parameter : by model in training. (W, b)

Hyperparameter : by user before training. (LR, Epoch, ...)

Calculate # params

1. Linear: $(I \times O) + O$ (I : input size, O : output size)

2. Conv: $((K_w \times K_h \times C_{in}) + 1) \times C_{out}$ (K : kernel C : channel)

3. Batch Norm: $C \times 2$ C : channel.

∇ pooling : 0

4. RNN: $H \times I + H \times H + H$ (H : hidden state, I : input)

6. LSTM: $4 \times (\text{RNN})$ ($\because 4$ gates)

7. GRU: $3 \times (\text{RNN})$ ($\because 3$ gates)

8. Transformer: (d : model dimension, V : vocab size.)

Each Encoder: Multi-head Attention: $4 \times d^m$ (K, Q, V, Out)

Feed Forward Network: $2(d \times 4d) = 8d^m$ ($d \rightarrow 4d \rightarrow d$)

Total: $(V \times d) + N \times (\quad)$

VGGNet

ResNet

Sequential Data: Sentences, video frames, speech waveform, ...
 bag of words: vocab set, no order. \Rightarrow one hot encode, ...

model for sequential data: deal with variable length sequence
 maintain sequence order
 keep track of long term dependencies
 share parameters across the sequence.

Vectorizing: transform text into numeric tensors.
 \Rightarrow break down text into (word, character, n-gram)

Tokenization: word : delimiter = space, most common, big vocab. (embedding matrix), OOV
 character : small vocab, no meaning, long seq length.
 subword : no split frequent words, split rare words into meaningful subword, zero-shot OK
 Byte-Pair Encoding : iteratively find most frequent pair of adjacent characters in data
 merge them into a new single token. \Rightarrow add new tokens until desired vocab size.

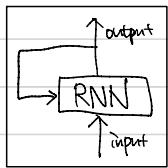
Vectorization:

one-hot encoding : size of vector = size of vocab. each token has unique index of vector
 binary, sparse, high-dimensional, hard-coded

token(word) embedding: use of dense word vectors. lower-dimension, learn from data

Why Embedding?
 efficiency: one-hot \Rightarrow large, mostly zero vectors. embedding \Rightarrow direct row indexing, O(1) MEM efficient.
 interpretability: directly storing token vector. \Rightarrow pretrained, visualized, reused.
 simplicity: each token has its own learnable vector.

Recurrent Neural Networks \Rightarrow process sequence by iterating through step elem and maintain state
 state: containing info relative to what it has seen so far



state = 0

for input in input-seq:

$$\text{output} = \text{activation}(\text{dot}(W, \text{input}) + \text{dot}(U, \text{state}) + b)$$

$$\text{state} = \text{output}$$

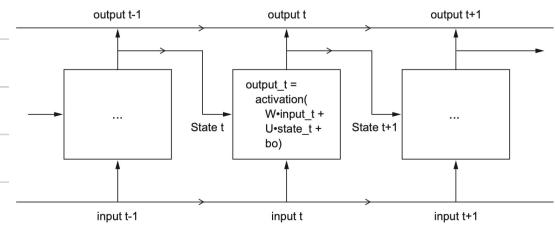


Figure 6.10 A simple RNN, unrolled over time

param of Simple RNN.

Embedding layer: Vocab size \times Embedding Dim (1000×64)

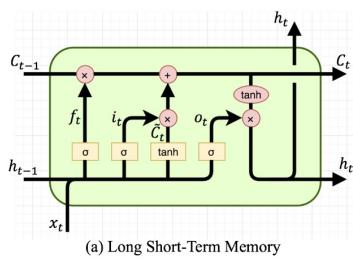
$$y = \sigma(\text{state} \cdot U + \text{input} \cdot W + b)$$

RNN layer: $H \times (H+I) + B$ $64(64+64) + (64 \times 2)$

Linear layer: $(I \times D) + B$ $(64 \times 1) + 1$

Packing \Rightarrow padding ~~out~~ ~~2D~~ NLP.

LSTM : carry info across many timesteps.
Simple RNN + additional data flow. (Carry)



$$i_t = \sigma(\mathbf{x}_t U^i + h_{t-1} W^i)$$

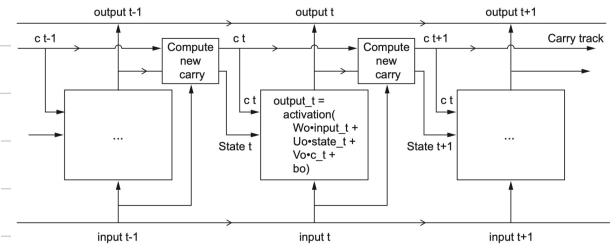
$$f_t = \sigma(\mathbf{x}_t U^f + h_{t-1} W^f)$$

$$o_t = \sigma(\mathbf{x}_t U^o + h_{t-1} W^o)$$

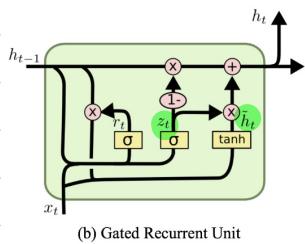
$$\tilde{C}_t = \tanh(\mathbf{x}_t U^g + h_{t-1} W^g)$$

$$C_t = f_t * C_{t-1} + i_t * \tilde{C}_t$$

$$h_t = \tanh(C_t) * o_t$$



GRU



$$z_t = \sigma(x_t U^z + h_{t-1} W^z)$$

$$r_t = \sigma(x_t U^r + h_{t-1} W^r)$$

$$\tilde{h}_t = \tanh(x_t U^h + (r_t * h_{t-1}) W^h)$$

$$h_t = (1 - z_t) * h_{t-1} + z_t * \tilde{h}_t$$

Bidirectional RNNs

forward backward } knowledge coming next

Attention \Rightarrow vector of importance weights. To predict one elem, using attention how strongly correlated with others,

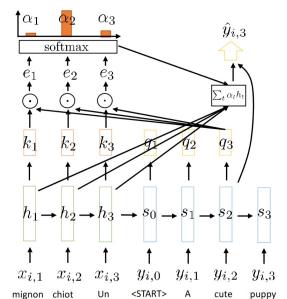
Seq2Seq \Rightarrow fixed size embedding. is incapable of remembering long seq. decode from last hidden state

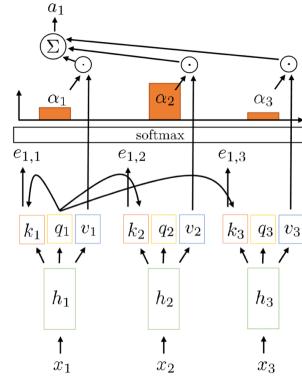
Attention \Rightarrow context vector (decoder state) from entire encoded vectors.

teacher forcing in training by feeding GT token.

$$\text{softmax}\left(\frac{\mathbf{QK}^T}{\sqrt{dk}}\right) \cdot \mathbf{V} = \text{attention}(\mathbf{Q}, \mathbf{K}, \mathbf{V})$$

Attention can access every time step. RNN \rightarrow pure attention (OK)
no temporal dependencies.





$$a_t = \sum_t \alpha_{l,t} v_t$$

$$\alpha_{l,t} = \exp(e_{l,t}) / \sum_{t'} \exp(e_{l,t'})$$

we'll see why this is important soon

$$e_{l,t} = q_t \cdot k_t$$

c.g., $k_t = W_k h_t$

$$q_t = q(h_t)$$

e.g., $q_t = W_q h_t$

this is *not* a recurrent model!
but still weight sharing:

$$h_t = \sigma(Wx_t + b)$$

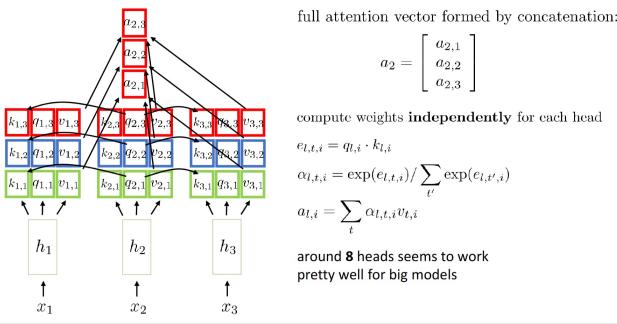
shared weights at all time steps

(or any other nonlinear function)

Self-attention limitation

1. Positional encoding: position of words has meaning. represent relative position by frequency-based sin/cos encoding.
 learning positional encoding \Rightarrow more flexible, complex. cannot generalize.
 $\text{input} = \text{emb}(1_t) + p_t$

2. Multi-head attention: softmax \Rightarrow dominated by one value. hard to specify 2 diff things.
 Then... multiple keys, queries, values for every time step.



full attention vector formed by concatenation:

$$a_2 = \begin{bmatrix} a_{2,1} \\ a_{2,2} \\ a_{2,3} \end{bmatrix}$$

compute weights independently for each head

$$e_{l,t,i} = q_{l,i} \cdot k_{l,i}$$

$$\alpha_{l,t,i} = \exp(e_{l,t,i}) / \sum_{t'} \exp(e_{l,t',i})$$

$$a_{l,i} = \sum_t \alpha_{l,t,i} v_{l,i}$$

around 8 heads seems to work pretty well for big models

3. Adding nonlinearity: every self attention layer is a linear transformation of previous layer

$$h_t^l = \sigma(W^l a_t^l + b^l)$$

4. Masked encoding: prevent attention lookup into the future.

problem: self attention can look at the value at next steps, based on inputs at next steps.

During test time \Rightarrow for output of step t requires input of $t+1$.

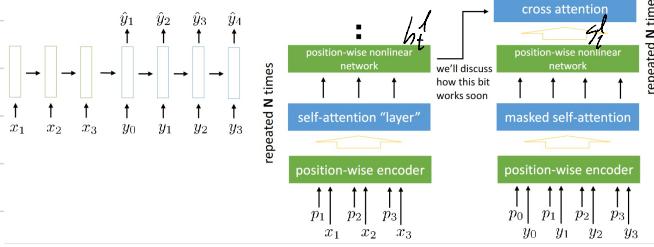
Self attention \Rightarrow allow past info, but not future.

Solution: $e_{l,t} = \begin{cases} g_l \cdot k_t & (\text{if } l > t) \\ -\infty & (\text{o.w.}) \end{cases} \Rightarrow \text{Zero value in softmax.}$

Transformer

The “classic” transformer

As compared to a sequence to sequence RNN model



Cross Attention

$$\text{query: } q_t^l = W_q^l h_t^l$$

$$\text{key: } k_t^l = W_k^l h_t^l$$

$$\text{value: } v_t^l = W_v^l h_t^l$$

$$c_{i,t}^l = \beta_i^l \cdot k_t^l, \quad \alpha_{i,t}^l = \exp(c_{i,t}^l) / \sum_t \exp(c_{i,t}^l), \quad C_i^l = \sum_t \alpha_{i,t}^l v_t^l$$

* Layer Normalization for variable length of text. (dimension-wise)

Each data get its own stat(μ, σ) \Rightarrow norm

Putting it all together

The Transformer

multi-head attention keys and values
 $k_{t,1}^l, \dots, k_{t,m}^l$ and $v_{t,1}^l, \dots, v_{t,m}^l$

6 layers, each with $d = 512$

$$\bar{h}_t^l = \text{LayerNorm}(\bar{a}_t^l + h_t^l)$$

passed to next layer $\ell + 1$

$$h_t^l = W_2^\ell \text{ReLU}(W_1^\ell \bar{a}_t^l + b_1^\ell) + b_2^\ell$$

2-layer neural net at each position

$$\bar{a}_t^l = \text{LayerNorm}(\bar{h}_{t-1}^{\ell-1} + a_t^l)$$

essentially a residual connection with LN

$$\text{input: } \bar{h}_{t-1}^{\ell-1}$$

output: a_t^l

concatenates attention from all heads

Decoder decodes one position at a time with masked attention

Drawback

$$\rightarrow O(n^2)$$

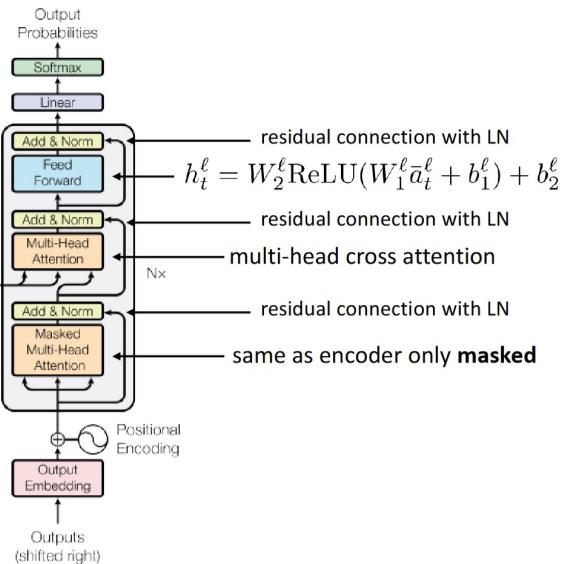
complex implementation

Benefits

\rightarrow long range connection

easy to parallelize

deeper network.



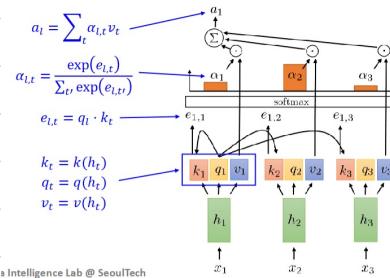
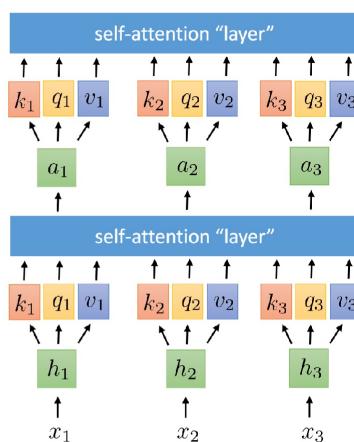
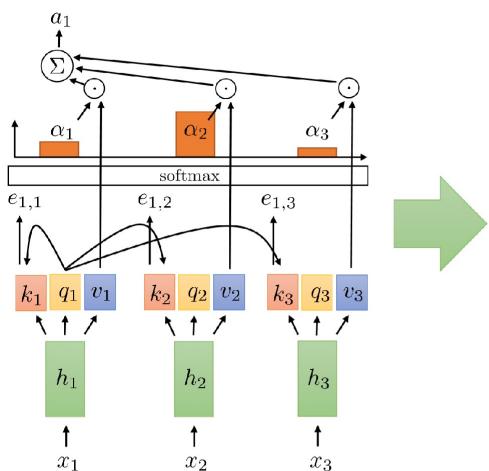
Vaswani et al. *Attention Is All You Need*. 2017.

Language Model.: probabilistic model that assigns a probability $P(w_1, w_2, \dots, w_n)$ to every finite seq w_1, w_2, \dots, w_n .

chain rule of conditional prob. $P(w_1, w_2, \dots, w_n) = P(w_1) P(w_2|w_1) \dots P(w_n|w_1, w_2, \dots, w_{n-1})$

\Rightarrow equivalent to being able to predict next word.

Attention can access every time step. RNN can do, Attention can do more.



Intelligence Lab @ SeoulTech

Normalization Position. LN $\Rightarrow \frac{x - \mu}{\sigma}$, $\sigma = \sqrt{\frac{1}{n} \sum (x_i - \mu)^2}$

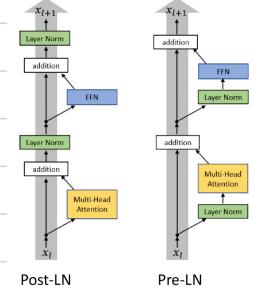
Post-LN : vanilla. between residual blocks. instable due to large gradients near the output layer.

Pre-LN : LN applied before each sub-layer.

RNSNorm : simplify LN by removing mean stat. $RMS(a) = \sqrt{\frac{1}{n} \sum a_i^2}$



Activation Function : ReLU, Gelu ...



Positional Embedding : APE : added to context representation,

RPE : add position info to model.

FLAN : data \Rightarrow Instruction, improve zero-shot.

Alignment & RLHF \Rightarrow ① Supervised Fine Tuning. ② Reward Model. ③ Proximal Policy Optimization

\hookrightarrow objective = follow user's instruction usefully safely

Discriminative vs. Generative Models. (Give observed variable x , target variable y)

Discriminative \Rightarrow estimate $P(y|x)$ conditional prob.

Generative \Rightarrow estimate $P(x|y)$ joint prob.

w/o assuming y , learn $p(x)$ from given data.

\hookrightarrow enable use to generate new data similar to training dataset

how to sample new plausible data?

Statistical perspective: generation/sampling \Rightarrow data point from probability distribution

Map simple dist $P(z)$ to data dist $P(x)$. z : latent variable, $p(z)$: latent dist.

Sample from $p(z)$, map it to data point ① $\xleftarrow{p(x|z)}$ ④

AutoEncoder. $x \xrightarrow[\text{enc}]{e_\theta(x)} z \xrightarrow[\text{dec}]{d_\theta(z)} \hat{x}$

\Rightarrow enable accurate reconstruction

decoder (latent \rightarrow input) \Rightarrow latent space is unstructured.

guarantee latent is mapped to data point

random sampling latent code produce meaningless outputs.

Variational AutoEncoder. (probabilistic structure on latent space)

$$x - e_\theta(x) - \frac{\mu}{\sigma} - z \sim N(\mu, \sigma^2) - z - d_\theta(z) - \hat{x} \quad \text{①} \xleftarrow[\text{④}]{p(x|z)} \text{②} \quad p(x|z) = N(x; D(z), \sigma^2 I), p(z) = N(z; 0, I)$$

Max ELBO. Max $P(x)$

$$\text{ELBO} = \mathbb{E}_{q_\theta(z|x)} [\log(p_\theta(x|z))] - D_{KL}(q_\theta(z|x) || p(z))$$

\Rightarrow reconstruction (max)
accurate recovery of x from z .

\Rightarrow prior matching term (min)
 $q_\theta(z|x)$ stay close to $p(z)$

$$p(x) = \int p(x,z) dz = \int p(x|z)p(z) dz, \quad p(x|z) = N(x; D(z), \sigma^2 I), \quad p(z) = N(z; 0, I) \quad : \text{decoder}$$

$$q_\theta(z|x) = N(z; M_\theta(x), C_\theta(x)) \quad : \text{encoder}$$

VAE training

1. Feed data point \mathbf{x} to encoder to predict $\mu_\phi(\mathbf{x})$, $\sigma_\phi^2(\mathbf{x})$
 2. Sample latent variable \mathbf{z} from $p_\theta(z|\mathbf{x}) = N(\mathbf{z}; \mu_\phi(\mathbf{x}), \sigma_\phi^2(\mathbf{x})\mathbf{I})$
 3. Feed \mathbf{z} to decoder to predict $\hat{\mathbf{x}} = D_\theta(\mathbf{z})$
 4. Perform gradient descent through (-ELBO) $\rightarrow \max (+\text{ELBO})$
- (2) $\Rightarrow \mu + \sigma \cdot \mathbf{e}$ \rightarrow differentiable. \mathbf{e} from $N(0, 1)$ by reparameterization trick. Randomness separation to deterministic

VAE Generation.

1. Sample latent value \mathbf{z} from $p(\mathbf{z}) = N(\mathbf{z}; 0, \mathbf{I})$
2. Feed \mathbf{z} to decoder to predict $\hat{\mathbf{x}} = D_\theta(\mathbf{z})$

Generative Adversarial Networks, \Rightarrow minimax game D vs. G .

$$\text{Objective: } \min_G \max_D V(D, G) = \mathbb{E}_{\mathbf{x} \sim p(\text{data})} [\log D(\mathbf{x})] + \mathbb{E}_{\mathbf{z} \sim p(\mathbf{z})} [\log (1 - D(G(\mathbf{z})))]$$

$D(\mathbf{x})$ = prob. input \mathbf{x} is real.

Heuristics \Rightarrow Saturating Gradients. $D(G(\mathbf{z}))$ instead of $1 - D(G(\mathbf{z}))$.

$$D-G \begin{cases} \text{sat. until } D=1 & D(\mathbf{x}) \rightarrow 0.5 \\ & D(G(\mathbf{z})) \rightarrow 0.5 \end{cases}$$

$$\ell(\mathbf{x}, y) = -[y \log \mathbf{x} + (1-y) \log (1-\mathbf{x})] \quad \text{real: 1} \quad \text{fake: 0}$$