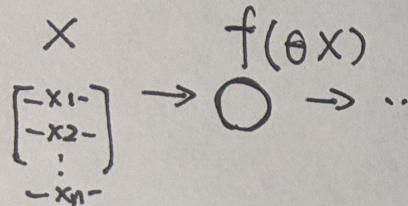


Deep Learning

• Neuron:



X : data matrix

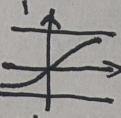
$$\theta = \begin{bmatrix} 1 \\ \theta_1 \\ \theta_2 \\ \vdots \end{bmatrix} \text{ Params}$$

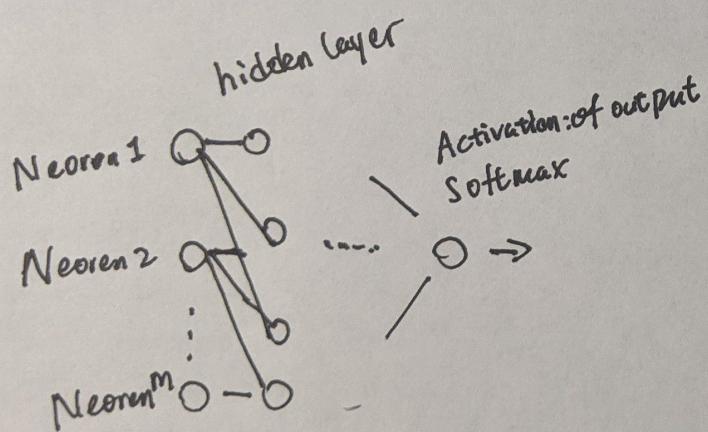
θX : linear comb of param * data

f : activation function:

ReLU: 

Sigmoid: 

tanh: 

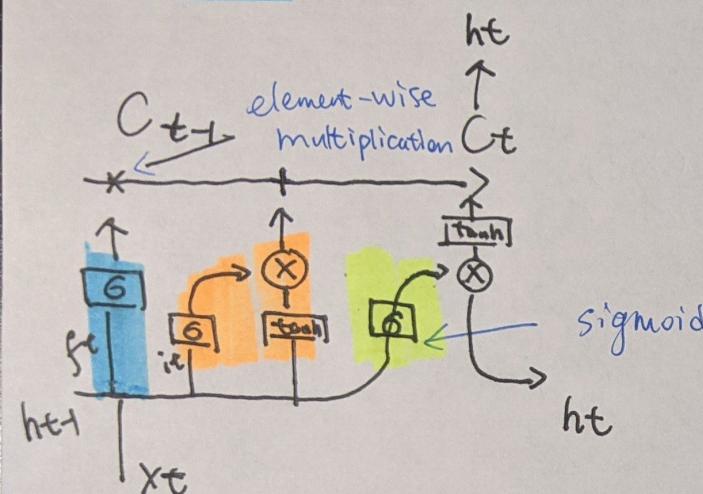


Softmax:

output a array of prob.
sum of these is 1

• Backpropagation: update params by gradient descent & chain rule.

LSTM



Computation:

$$\tanh \left[\sigma(w_f [x_t^{h_{t-1}}])^+ \sigma(w_i [x_t^{h_{t-1}}]) * \right. \\ \left. \tanh (w_c [x_t^{h_{t-1}}]) \right] \\ + \sigma(w_o [x_t^{h_{t-1}}]) \text{ new Value} \\ \text{Output gate}$$

params:

$$4 \times \text{Shape}(h) = (\text{Shape}(h) + \text{shape}(x))$$

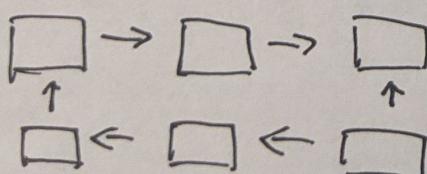
\nearrow \swarrow $\overbrace{\quad\quad\quad}$

$w_f, w_i, w_c,$
 $w [$

 $h+x$
 $]$

Bidirectional LSTM:

forward pass



Back Ward pass

Word Embedding

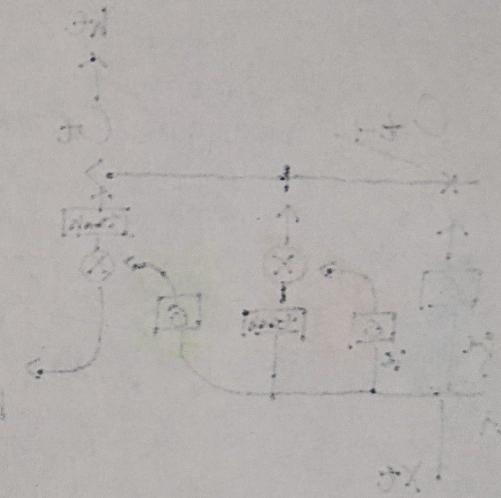
MTBD

$$\begin{array}{c} \text{Word in dict} \\ \boxed{\text{A} \quad \text{B} \quad \text{C} \quad \text{D}} \end{array} \times \begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \\ 0 \end{bmatrix} \rightarrow \boxed{\text{A}}$$

one hot encoding

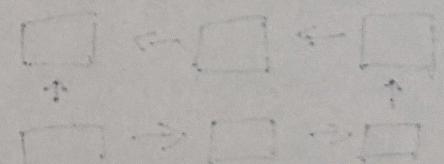
Encoding for one word,

Length is customized.



word embedding

word vector

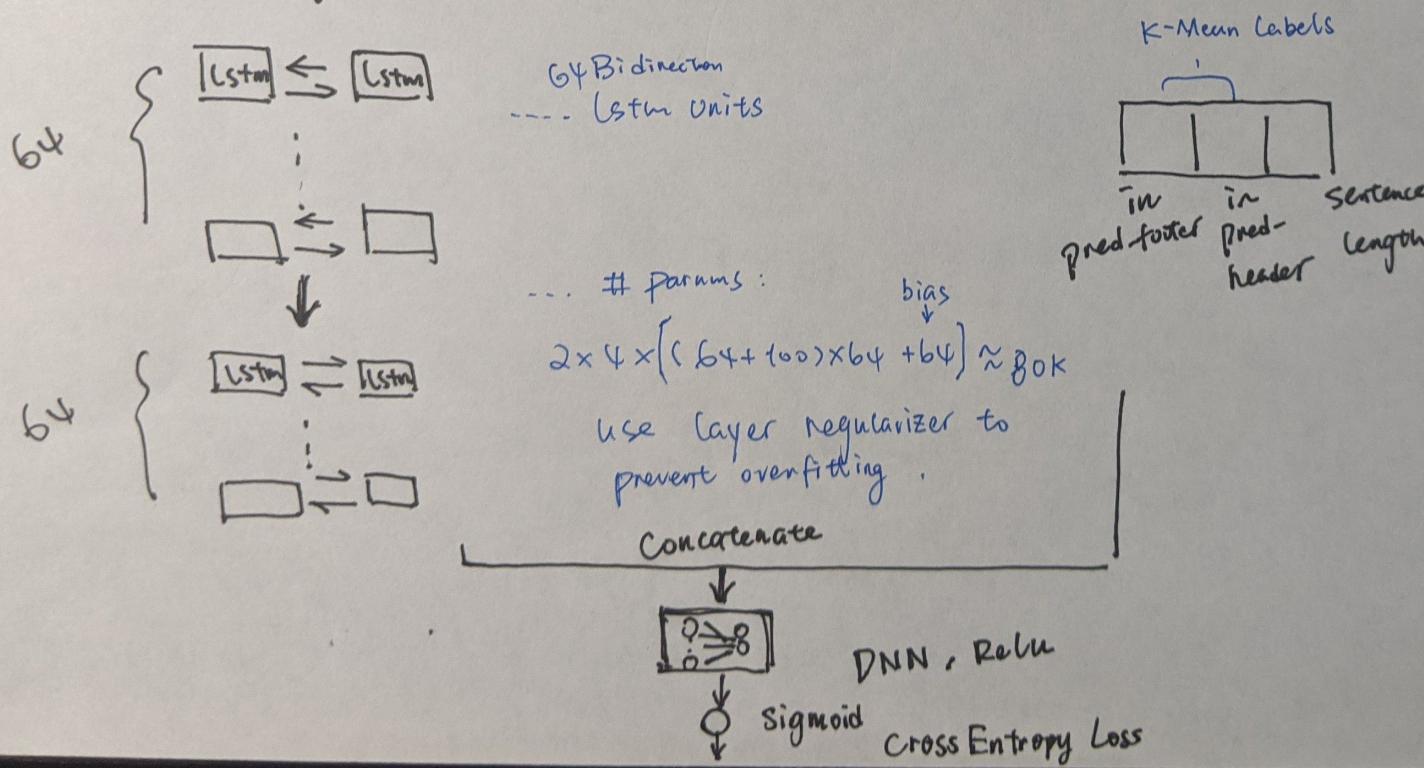
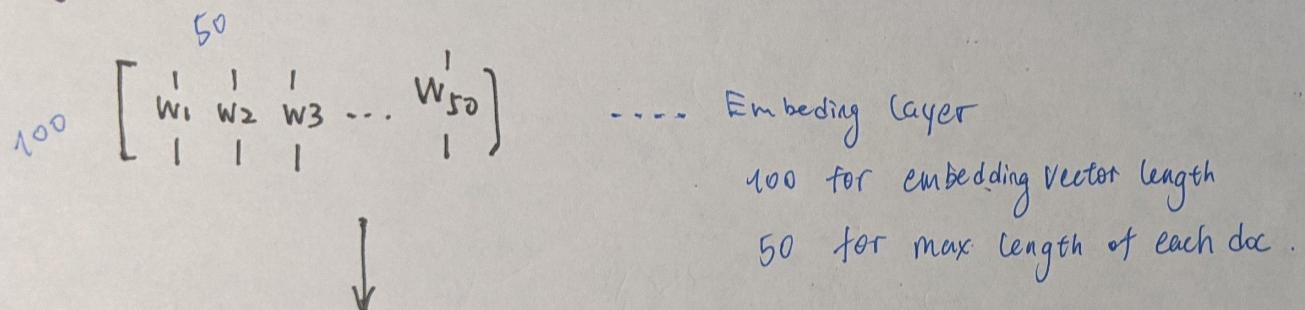


word vector

Model From Internship.

text classification for
cleaning text.

- each piece of doc are sep by " " (triple space)
 - 0 as need to keep
 - 1 as need to remove



Transformer Networks

Vector-Seg Model :

Dog img \rightarrow "Dog in a cup"

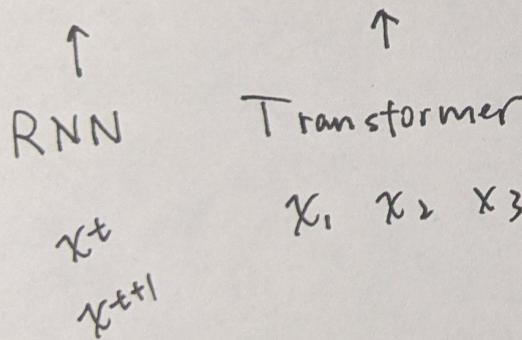
→
mark
[mark]

Seq-vec Model :

"Text..." \rightarrow [0~1]

Seg-seg Model :

Translation



-- Positional Encoder:

vector that gives context based on position of word in sentence

Seg. Model

Attention Layer : What part of the input should we focus?

E-C Attention

Enc

$[x_1 \ x_2 \ x_3 \ \dots \ x_m]$

$$\begin{aligned} & \rightarrow K_i = W_k x_i \quad \dots \\ & \rightarrow V_i = W_v x_i \quad \dots \end{aligned}$$

Structure:

Dec

$[x'_1 \ x'_2 \ \dots \ x'_m]$

$$\rightarrow q_j = W_q x'_j \ \dots$$

Formulas

$$C = \text{Attention}(x, x') \quad \alpha_i = \text{Softmax}(K^T q_i) \in \mathbb{R}^m$$

$$c_i = \alpha_1 v_1 + \dots + \alpha_m v_m = V \alpha_i$$

$$\alpha_2 = \dots$$

$$\alpha_m = \dots$$

Notes:

C can be used to generate next word (Translation Task)

Self-Attention Layer

$$\text{Query } q_i = W_Q x_i$$

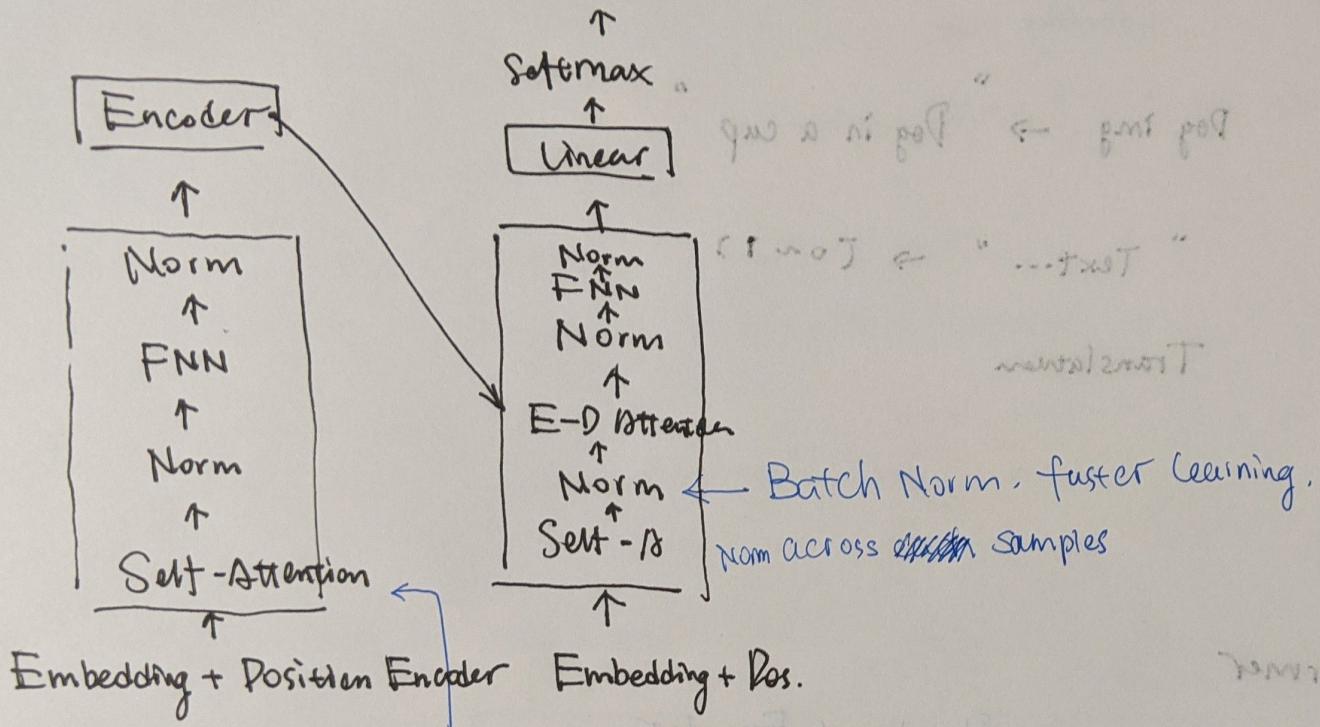
$$\text{Key } k_i = W_K x_i$$

$$\text{Value } v_i = W_V x_i$$

$$\alpha_j = \text{Softmax}(K^T q_i : j) \\ c_1 = \alpha_{11} v_1 + \dots + \alpha_{m1} v_m = V d_2 \quad \left. \right\} \Rightarrow c_j = V \cdot \text{Softmax}(K^T q_i : j)$$

$$\text{Self-attention layer: } C = \text{Attn}(X, X)$$

Encoder - Decoder



Learned representation.

: word - seq - relav

: word - seq - pos

: word - pos - pos

different attention vectors,

avg. for each word vector

$$[x_1 \ x_2 \ \dots \ x_n]$$

$$[w_1 \ w_2 \ \dots \ w_n]$$

$$\dots = \max_{i=1}^n \dots$$

$$\dots = \max_{i=1}^n \dots$$

BERT

- is for pre-training Transformer's encoders
- Predict masked word, next sentence.

① mask 15% words randomly.

output a prob @ masked position.

calculate $\text{loss} = \text{CrossEntropy}(e, p)$

ground truth Prob

Perform GD to update parameters

② ~~mask~~ use 50% true linked 2 sentences,
50% faked

[CLS] Sent 1 [SEP] Sent 2



token for classification for separating sentences

- Combing Two Tasks

input

example:

"[CLS] calculus is a [MASK]
of math,

[SEP] it [MASK] developed
by newton and leibniz."

Target: true, "branch", "was"

LOSS:

next-sent:

binary clf: cross

Entropy

Masked word:

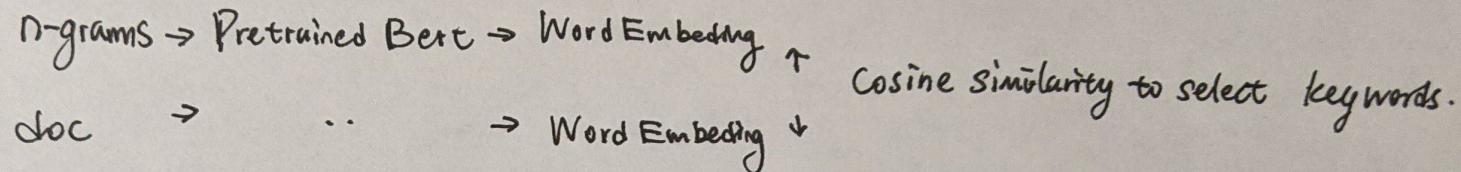
multi-class ~~clf~~
Classification

~~clf~~

Sum of three loss

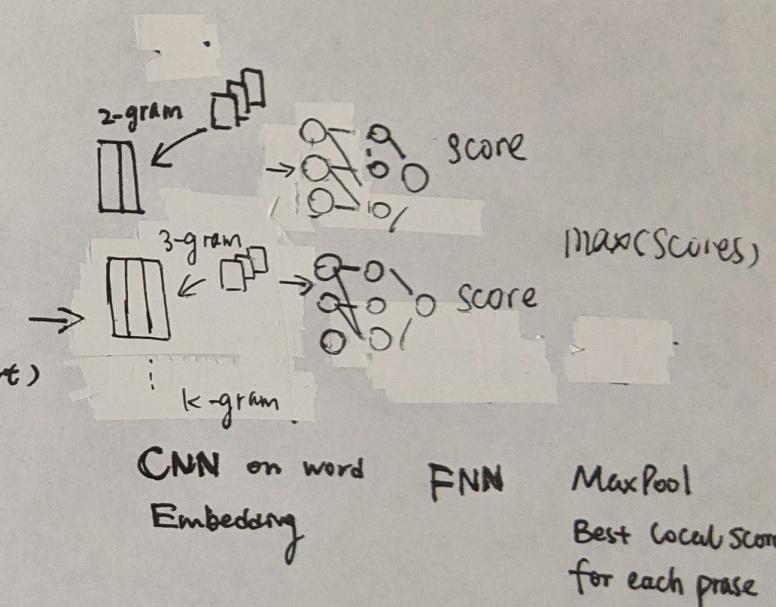
Model From Internship., Keyphrase extraction task

Unsupervised Keyphrase Extraction:



Supervised Keyphrase Extraction. (BERT-KPE)

n -grams \rightarrow Pretrained Bert
(Base/SpanBert/Roberta)



Loss :

Hinge Loss : pairwise ranking loss

$$\sum_{\text{phase-doc}} \max(0, 1 - f^*(p_+, D) + f^*(p_-, D))$$

CE Loss : is or not true keyphrase

$\angle \text{chunk} = \text{Cross Entropy } (P(c_i^k = y_i^k))$

Multi-task training

$$\text{Loss} = \angle \text{Rank} + \angle \text{chunk}$$

Pytorch:

Pros : Many things are modularizable

Cons : May not be good for launching model into production

Code example for train model (Stats 507),

for epoch in range(n-epochs),

model.train(),

for idx, batch in enumerate(train-loader):

forward passing

logits, probs = model(...)
input \rightarrow output of softmax

Cost = F. cross_entropy(logit, true_labels)

update gradients
optimizer.zero_grad()

Cost.backward()

update model parameters

optimizer.step()

evaluation

model.eval(),

Cost = Compute_epoch_Cost(model, train-loader)

:

More codes from lecture —>

PyTorch Usage: Step 1 (Definition)

```
class MultilayerPerceptron(torch.nn.Module):

    def __init__(self, num_features, num_classes):
        super(MultilayerPerceptron, self).__init__()

        ### 1st hidden layer
        self.linear_1 = torch.nn.Linear(num_feat, num_h1)

        ### 2nd hidden layer
        self.linear_2 = torch.nn.Linear(num_h1, num_h2)

        ### Output layer
        self.linear_out = torch.nn.Linear(num_h2, num_classes)

    def forward(self, x):
        out = self.linear_1(x)
        out = F.relu(out)
        out = self.linear_2(out)
        out = F.relu(out)
        logits = self.linear_out(out)
        probas = F.log_softmax(logits, dim=1)
        return logits, probas
```

Define model parameters that will be instantiated when created an object of this class

Define how and it what order the model parameters should be used in the forward pass

PyTorch Usage: Step 2 (Creation)

```
torch.manual_seed(random_seed)
model = MultilayerPerceptron(num_features=num_features,
                             num_classes=num_classes)    | Instantiate model
                                                               (creates the model parameters)
model = model.to(device)                                | Optionally move model to GPU, where
                                                               device e.g. torch.device('cuda:0')
optimizer = torch.optim.SGD(model.parameters(),
                           lr=learning_rate)      | Define an optimization method
```

PyTorch Usage: Step 3 (Training)

```
for epoch in range(num_epochs):
    model.train()
    for batch_idx, (features, targets) in enumerate(train_loader):

        features = features.view(-1, 28*28).to(device)
        targets = targets.to(device)

        ### FORWARD AND BACK PROP
        logits, probas = model(features)
        cost = F.cross_entropy(probas, targets)
        optimizer.zero_grad()

        cost.backward()

        ### UPDATE MODEL PARAMETERS
        optimizer.step()

    model.eval()
    with torch.no_grad():
        # compute accuracy
```

Run for a specified number of epochs

Iterate over minibatches in epoch

If your model is on the GPU, data should also be on the GPU

PyTorch Usage: Step 3 (Training)

```
for epoch in range(num_epochs):
    model.train()
    for batch_idx, (features, targets) in enumerate(train_loader):

        features = features.view(-1, 28*28).to(device)
        targets = targets.to(device)

        ### FORWARD AND BACK PROP
        logits, probas = model(features) ← This will run the forward() method
        loss = F.cross_entropy(logits, targets) ← Define a loss function to optimize
        optimizer.zero_grad() ← Set the gradient to zero
                                (could be non-zero from a previous forward pass)

        loss.backward() ← Compute the gradients, the backward is automatically
                        constructed by "autograd" based on the forward()
                        method and the loss function

        ### UPDATE MODEL PARAMETERS
        optimizer.step() ← Use the gradients to update the weights according to
                           the optimization method (defined on the previous slide)
                           E.g., for SGD,  $w := w + \text{learning\_rate} \times \text{gradient}$ 

    model.eval()
    with torch.no_grad():
        # compute accuracy
```

PyTorch Usage: Step 3 (Training)

```
for epoch in range(num_epochs):
    model.train()
    for batch_idx, (features, targets) in enumerate(train_loader):

        features = features.view(-1, 28*28).to(device)
        targets = targets.to(device)

        ### FORWARD AND BACK PROP
        logits, probas = model(features)
        loss = F.cross_entropy(logits, targets)
        optimizer.zero_grad()

        loss.backward()

        ### UPDATE MODEL PARAMETERS
        optimizer.step()

    model.eval()
    with torch.no_grad():
        # compute accuracy
```

For evaluation, set the model to eval mode (will be relevant later when we use DropOut or BatchNorm)

This prevents the computation graph for backpropagation from automatically being build in the background to save memory

PyTorch Usage: Step 3 (Training)

```
for epoch in range(num_epochs):
    model.train()
    for batch_idx, (features, targets) in enumerate(train_loader):

        features = features.view(-1, 28*28).to(device)
        targets = targets.to(device)

        ### FORWARD AND BACK PROP
        logits, probas = model(features)
        loss = F.cross_entropy(logits, targets)
        optimizer.zero_grad()

        loss.backward()

        ### UPDATE MODEL PARAMETERS
        optimizer.step()

    model.eval()
    with torch.no_grad():
        # compute accuracy
```

logits because of computational efficiency.
Basically, it internally uses a `log_softmax(logits)` function
that is more stable than `log(softmax(logits))`.
More on logits coming up.
