# Sequence Modeling

**Week #4 - 2 Oct '24**

**NLP Lab**

# NLP Lab Project : Build your own GPT

Text Processing and Bag of Words Model

N-grams and TF-IDF

Theory of Word Embeddings : Word2Vec, GloVe

RNNs: LSTM, GRU; Attention Mechanisms

Transformer Model, Advanced Transformers: GPT, BERT

Fine-Tuning, SFT, Continued Pre-training

Custom Language Model using Prompt Engineering, RAG

# 1. Language Models Recap

# 1.1 Language Models

- Models that assign a probability to each possible next word.

- Language models can also assign a probability to an entire sentence, telling us that the following sequence has a much higher probability of appearing in a text.

- **Goal** : Compute  the probability   of a sentence or sequence   of words.

$$P(W) = P(w_1,w_2,w_3,w_4,w_5...w_n)$$

# 1.2 Discriminative vs Generative Models

- **Discriminative model:** a model that calculates the probability of a latent trait given the data.

$$P(Y \mid X) \text{ (Conditional)}$$

- **Generative model:** a model that calculates the probability of the input data itself.

$$P(X) \qquad P(X, Y)$$

stand-alone              joint

- Discriminative language models are designed to directly learn the boundary between different classes or outputs, typically in tasks involving classification or structured prediction.

# 1.3 Autoregressive Language Models

$$P(X) = \prod_{i=1}^{I} P(x_i \mid x_1, \ldots, x_{i-1})$$

Next Token     Context

- The big problem : How do we predict

$$P\,(x_i | x_i, x_2, \ldots, x_{i-1})$$

# 1.4 Examples of Language Models

- Count-based Language Model

  - Independence assumption: • Count-based maximum-likelihood estimation:

  $$P(x_i|x_1, \ldots, x_{i-1}) \approx P(x_i)$$

  - Count-based maximum-likelihood estimation.

  $$P_{\text{MLE}}(x_i) = \frac{c_{\text{train}}(x_i)}{\sum_{\tilde{x}} c_{\text{train}}(\tilde{x})}$$

- Higher-order n-gram Models

  $$P_{ML}(x_i \mid x_{i-n+1}, \ldots, x_{i-1}) := \frac{c(x_{i-n+1}, \ldots, x_i)}{c(x_{i-n+1}, \ldots, x_{i-1})}$$

- Featurized Log-Linear Models

# 1.5 Problems & Solutions

- Cannot share strength among **similar words**

  | | |
  |---|---|
  | she bought a car | she bought a bicycle |
  | she purchased a car | she purchased a bicycle |

  → solution: class based language models

- Cannot condition on context with **intervening words**

  Dr. Jane Smith    Dr. Gertrude Smith

  → solution: skip-gram language models

- Cannot handle **long-distance dependencies**

  for tennis class he wanted to buy his own racquet

  for programming class he wanted to buy his own computer

  → solution: cache, trigger, topic, syntactic models, etc.

# 2. Sequence Modeling

# 2.1 NLP and Sequential Data

- NLP is full of sequential data :

    ○ Characters in words

    ○ Words in sentences

    ○ Sentences in discourse

    ○ ……

# 2.2 Long-distance Dependencies in Language

- Agreement in number, gender, etc.

  **He** does not have very much confidence in **himself**.

  **She** does not have very much confidence in **herself**.


- Selectional preference.

  The **reign** has lasted as long as the life of the **queen**.

  The **rain** has lasted as long as the life of the **clouds**.

# 2.3 Can be Complicated !

- What is the referent of "**it**"?

The trophy would not fit in the brown suitcase because it was too **big**.

**Trophy**

The trophy would not fit in the brown suitcase because it was too **small**.

**Suitcase**

# 2.4 Types of Sequential Prediction

## Unconditioned vs Conditioned

- **Unconditioned Prediction:** Predict the probability of a single variable $P(X)$.

- **Conditioned Prediction:** Predict the probability of an output variable given an input $P(Y \mid X)$.

# 2.4 Types of Unconditioned Prediction



Left-to-right Autoregressive Prediction

$$P(X) = \prod_{i=1}^{|X|} P(x_i | x_1, \ldots, x_{i-1})$$

*(e.g. RNN or Transformer LM)*

Left-to-right Markov Chain (order n-1)

$$P(X) = \prod_{i=1}^{|X|} P(x_i | x_{i-n+1}, \ldots, x_{i-1})$$

*(e.g. n-gram LM, feed-forward LM)*

Independent Prediction

$$P(X) = \prod_{i=1}^{|X|} P(x_i)$$

*(e.g. unigram model)*

Bidirectional Prediction

$$P(X) \neq \prod_{i=1}^{|X|} P(x_i | x_{\neq i})$$

*(e.g. masked language model)*

# 2.4 Types of Conditioned Prediction



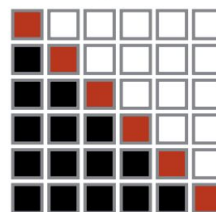Autoregressive Conditioned Prediction

$$P(Y|X) = \prod_{i=1}^{|Y|} P(y_i|X, y_1, \ldots, y_{i-1})$$
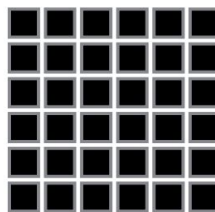
*(e.g. seq2seq model)*

Source X    Target Y

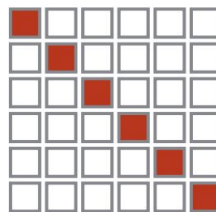Non-autoregressive Conditioned Prediction

$$P(Y|X) = \prod_{i=1}^{|Y|} P(y_i|X)$$

Source X    Target Y
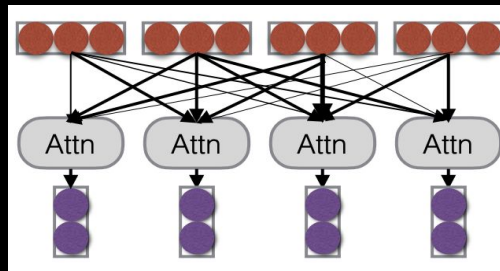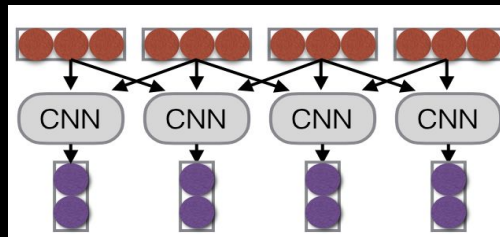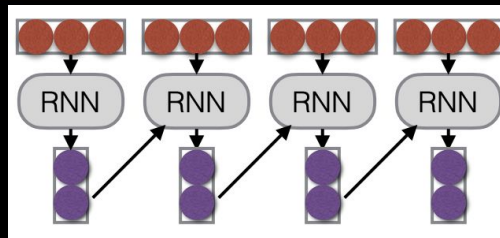
*(e.g. sequence labeling, non-autoregressive MT)*

# 2.5 Three Major Types of Sequence Models

- **Recurrence**: Condition representations on an encoding of the history.



- **Convolution**: Condition representations on local context.



- **Attention**: Condition representations on a weighted average of all tokens
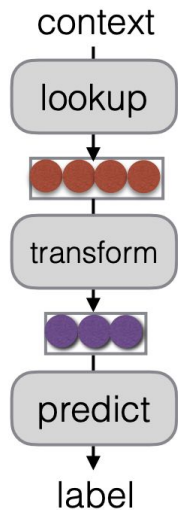
# 3. Recurrent Neural Networks

# 3.1 Recurrent Neural Networks

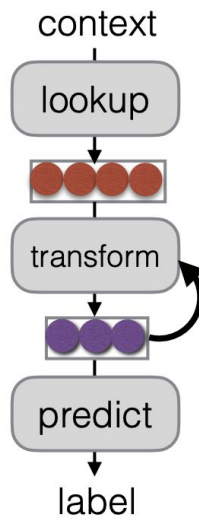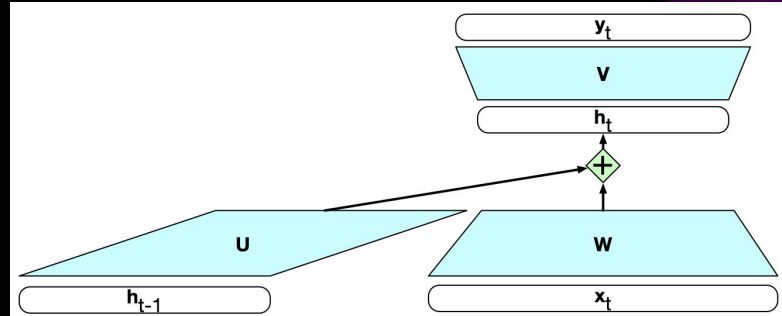- Tools to "remember" information.



Feed-forward NN

context
lookup
transform
predict
label

$$h_t = \mathrm{f}(W_x x_t + b)$$

Recurrent NN

context
lookup
transform
predict
label

$$h_t = \mathrm{f}(W_h h_{t-1} + W_x x_t + b)$$

# 3.2 Unrolling in Time

# 3.3 Training RNN Language Model



- **Self-Supervision** : the natural sequence of words is its own supervision!

- **Teacher Forcing** : model is always given the correct history sequence to predict the next word (rather than feeding the model its best case from the previous time step).

# 3.4 Generation with RNN-Based Language Models



- **Autoregressive Generation** : Sampling words conditioned on previous choices until a predetermined length is reached, or an end of sequence token is generated.

- The key to these approaches is to prime the generation component with an appropriate context. That is, instead of simply using <s> to get things started we can provide a richer task-appropriate context.

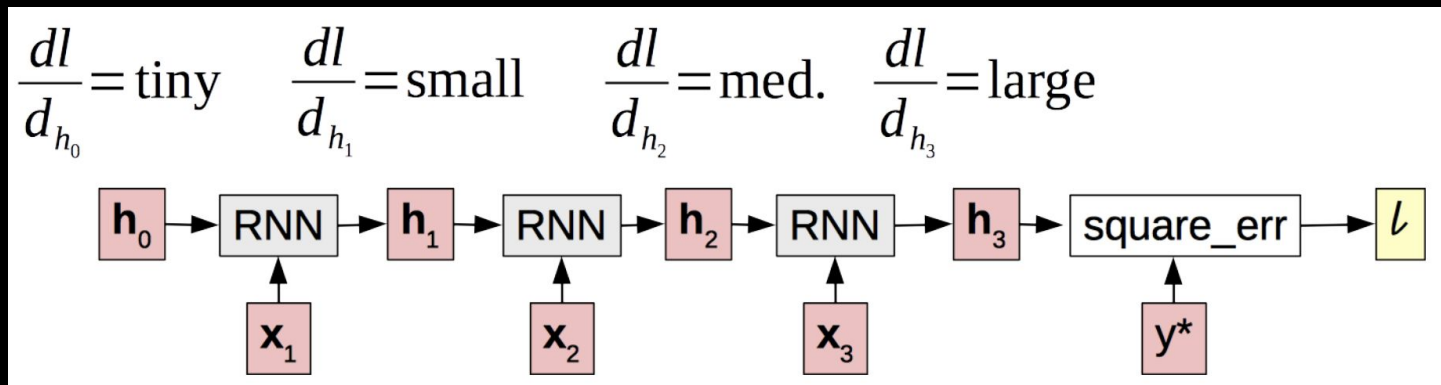# 3.5 Vanishing Gradient

- Gradients decrease as they get pushed back.

$$\frac{dl}{d_{h_0}} = \text{tiny} \qquad \frac{dl}{d_{h_1}} = \text{small} \qquad \frac{dl}{d_{h_2}} = \text{med.} \qquad \frac{dl}{d_{h_3}} = \text{large}$$

$h_0 \rightarrow$ RNN $\rightarrow h_1 \rightarrow$ RNN $\rightarrow h_2 \rightarrow$ RNN $\rightarrow h_3 \rightarrow$ square_err $\rightarrow l$

$x_1 \qquad x_2 \qquad x_3 \qquad y^*$

- Why? "Squashed" by non-linearities or small weights in matrices.

# 4. Long Short Term Memory

# 4.1 Basic Idea of LSTM



- **Basic idea:** make additive connections between time steps.
- Gates control the flow of the information and the gradients.

# 4.1 Equations of Gates in LSTM

- Forget Gate :  $\mathbf{f}_t = \sigma(\mathbf{U}_f\mathbf{h}_{t-1} + \mathbf{W}_f\mathbf{x}_t)$

- Update Gate :  $\mathbf{g}_t = \tanh(\mathbf{U}_g\mathbf{h}_{t-1} + \mathbf{W}_g\mathbf{x}_t)$

- Input/Add Gate :  $\mathbf{i}_t = \sigma(\mathbf{U}_i\mathbf{h}_{t-1} + \mathbf{W}_i\mathbf{x}_t)$

- Output Gate :  $\mathbf{o}_t = \sigma(\mathbf{U}_o\mathbf{h}_{t-1} + \mathbf{W}_o\mathbf{x}_t)$

# 5. Gated Recurrent Unit

# 4.1 GRU  Architecture



Gates:
$$o_t = \sigma(W_o s_{t-1} + U_o x_t + b_o)$$
$$i_t = \sigma(W_i s_{t-1} + U_i x_t + b_i)$$

States:
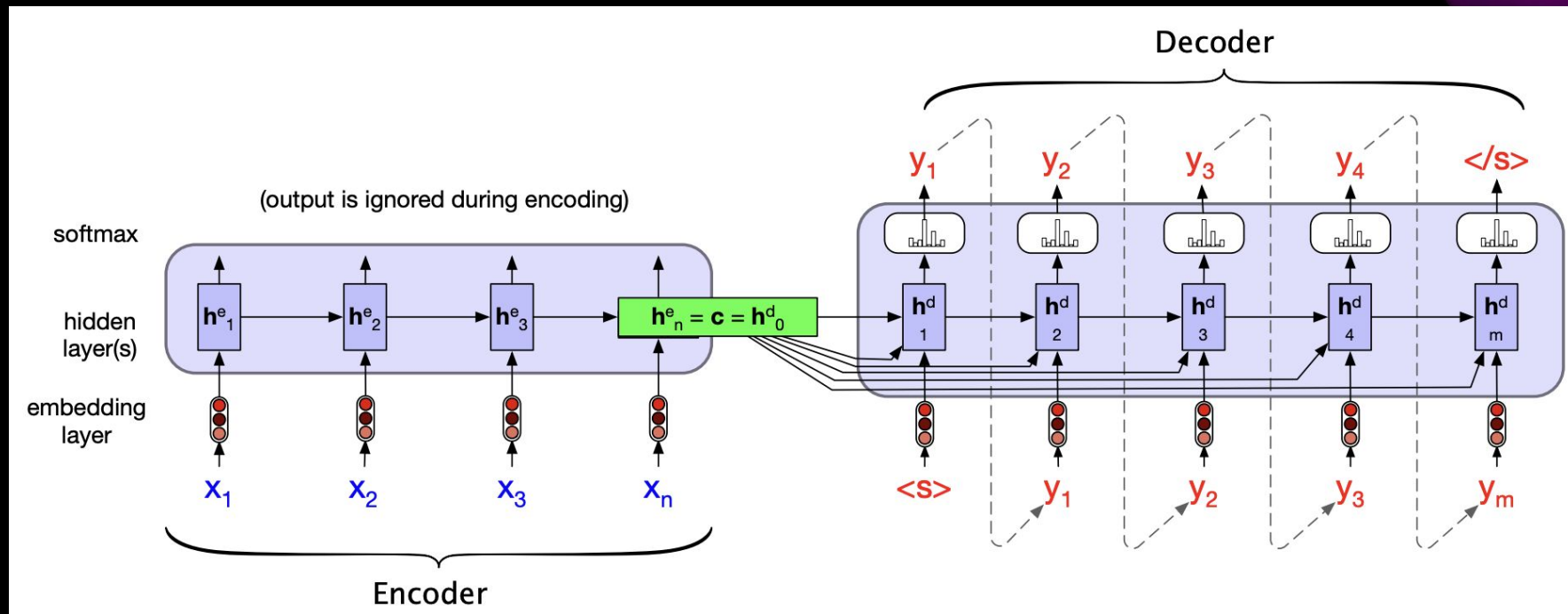$$\tilde{s}_l = \sigma(W(o_l \odot s_{l-1}) + U x_l + b)$$
$$s_l = (1 - i_l) \odot s_{l-1} + i_l \odot \tilde{s}_l$$

# 6. Encoder-Decoder

# 6.1 What is Encoder-Decoder ?

- Encoder-decoder networks, sometimes called sequence-to-sequence networks, encoder-decoder are models capable of generating contextually appropriate, arbitrary length, output sequences given an input sequence.

- The key idea underlying these networks is the use of an encoder network that takes an input sequence and creates a contextualized representation of it, often called the context.

- This representation is then passed to a decoder which generates a taskspecific output sequence.

# 6.2 Encoder-Decoder Models with RNN

# 7. Attention
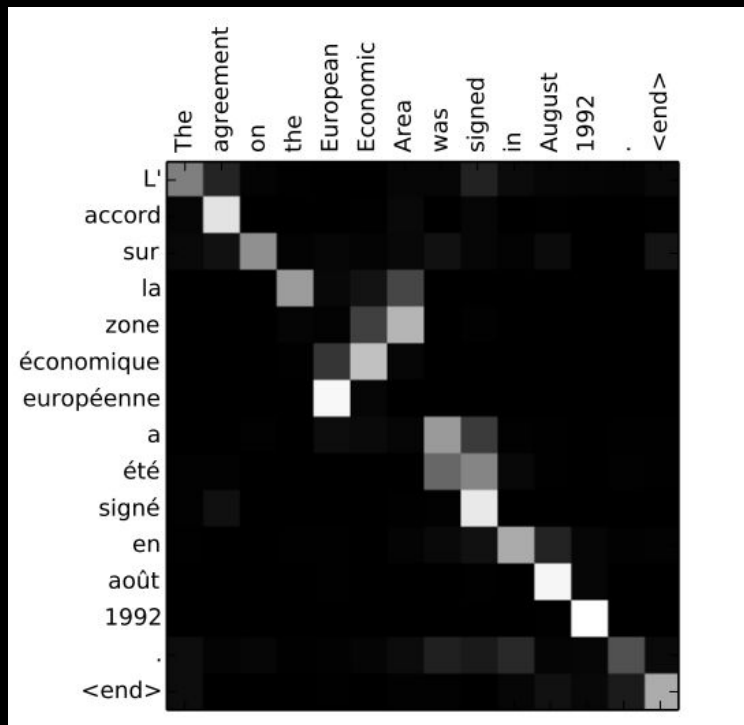
# 7.1 Basic Idea of Attention

- A mechanism that helps models focus on relevant parts of the input when making predictions.

- Overcomes limitations of traditional sequence models (like RNNs) that struggle with long dependencies.

- Two main types of Attention :
  - Cross-Attention
  - Self-Attention

# 7.2 Cross Attention

- Each element in a sequence attends to elements of another sequence.

# 7.3 Cross Attention : Bahdanau et al. 2015

- Encoder-Decoder Model:

  - The encoder processes the input sequence and generates hidden states.

  - The attention mechanism computes a weighted sum of these hidden states to form a context vector.

- Key Components:

  - Encoder Hidden States ($h_1$, $h_2$, ..., $h_\square$): Represent the input sequence.

  - Decoder Hidden State ($s_\square$): Represents the decoder's state at time step t.

  - Alignment Score ($e_{\square i}$): Measures the relevance between the decoder state at time t and the encoder hidden state $h_i$.

- Attention Weights ($\alpha_{\square i}$):

  - Computed using the alignment scores.

  - Determine how much "attention" should be paid to each input word.

# 7.3 Bahdanau Attention : Step-by-Step

- Alignment Scores:

    - Compute alignment scores between the current decoder hidden state and each encoder hidden state.

    - The alignment model is parametrized as a feedforward neural network and jointly trained with the remaining system components.

$$e_{ti} = v^T \tanh(W_1 s_{t-1} + W_2 h_i)$$

- Softmax Over Alignment Scores:

    - These attention weights tell how much the model should focus on each encoder hidden state.

- Context Vector (c□):

    - Compute a weighted sum of the encoder hidden states based on the attention weights

$$c_t = \sum_{i=1}^{n} \alpha_{ti} h_i$$

# 7.4 Luong Attention Mechanism

- Encoder Hidden States ($h_1$, $h_2$, ..., $h_\square$): Represent the input sequence.

- Decoder Hidden State ($s_\square$): Represents the decoder's state at time step t.

- Alignment Scores:

  - Compute alignment scores by taking the dot product of the decoder hidden state $s_\square$ and each encoder hidden state $h_i$.

$$e_{ti} = s_t \cdot h_i$$

- Attention Weights ($a_{\square i}$):

  - Apply softmax to the alignment scores to get attention weights

- Context Vector ($c_\square$):

  - Compute a weighted sum of the encoder hidden states based on the attention weights.

$$c_t = \sum_{i=1}^{n} \alpha_{ti} h_i$$

# Test you knowledge!