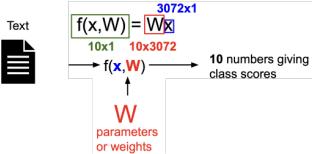


Linear Classifier

- Consider a 10 way classification problem
- 10 classes, 3072 features
- $f(x, W) = Wx$



Softmax classifier:

- Scores (vector):** $s = Wx_i$
- Softmax function for probability:** $P(Y = k) = \frac{e^{s_k}}{\sum_j e^{s_j}}$

Given: (many) input / output examples (x, y)

Find mapping f s.t. $y \approx f(x)$

$$\text{Define } L(W) = \sum_{i=1}^N L(f_{w_i}, \{x_i, y_i\}) + R(W)$$

learn by optimizing $W = \arg \min_W L(W)$

How does one pick the best W ?

Sigmoid classifier (binary):

- Score (scalar):** $s = \mathbf{w} \cdot \mathbf{x}_i$
- Sigmoid function for probability:**

- Positive class:

$$P(Y = 1) = \sigma(s) = \frac{1}{1+e^{-s}}$$

- Negative class:

$$P(Y = 0) = 1 - \sigma(s)$$

$$\frac{\partial L_i(\mathbf{W})}{\partial \mathbf{W}_{y_i}} = -x_i + p(y_i | x_i; \mathbf{W})x_i$$

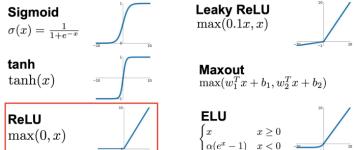
$$W_{t+1} = W_t - \alpha \nabla L(W_t) \quad \text{Gradient Descent Update Rule}$$

where:

α : step size (learning rate)

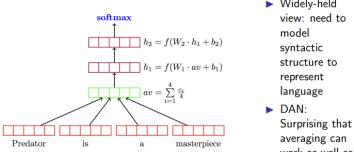
$\nabla L(W_t)$: gradient of L at W_t

W_t ; W_{t+1} : current and new value of W



ReLU is a good default choice for most problems

Deep Averaging Networks: feedforward neural network on average of word embeddings from input text



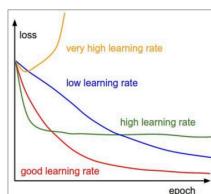
- Widely-held view: need to model syntactic structure to represent language
- DAN: Surprising that averaging can work as well as it does

$s = f(x; W_1, W_2) = W_2 g(W_1 x)$ non-linear score function

$$L_i = -\log \left(\frac{e^{s_{y_i}}}{\sum_j e^{s_j}} \right) \text{ loss on predictions}$$

$$L = \frac{1}{N} \sum_{i=1}^N L_i + \lambda R(W_1) + \lambda R(W_2) \quad \text{Total loss: data loss + regularization}$$

- Learning rate:** Step size in gradient descent. Too small: slow convergence. Too large: overshoots the minimum



Batching: Training on a single example at a time is slow, batching data gives speedups due to more efficient matrix operations

- Stochastic Gradient Descent (SGD):** Mini-batch size of 1, noisy updates
- Mini-Batch Gradient Descent:** Training on a small batch of examples at a time

Batch size is a hyperparameter (32, 64, 128, ...)

Skip Gram: Learning Word Embedding (Mikolov et al. 2013)

- Input:** a corpus of text (e.g. all of Wikipedia, News articles, Books, etc.)
- Output:** a set of embeddings: a real-valued vector for each word in the vocabulary
- We are going to learn these by setting up a fake prediction problem: predict a word's context from that word

the dog bit the man

(word = bit, context = dog)
(word = bit, context = the)

For each position $t = 1, \dots, T$, predict context words within a window of fixed size m , given center word w_j .

$$\text{Data Likelihood} = \prod_{t=1}^T \prod_{\substack{m \leq j \leq m \\ j \neq t}} P(w_{t+j} | w_t; \theta)$$

where θ are all the parameters of the model

We can write the loss as:

$$L(\theta) = -\frac{1}{T} \sum_{t=1}^T \sum_{\substack{m \leq j \leq m \\ j \neq t}} \log P(w_{t+j} | w_t; \theta)$$

For a center word c and a context word

o :

$$P(o | c) = \frac{\exp(u_o^T v_c)}{\sum_{w \in V} \exp(u_w^T v_c)}$$

$$\frac{\partial}{\partial v_c} \log(p(o | c)) = u_o - \frac{1}{\sum_{w=1}^V \exp(u_w^T v_c)} \sum_{x=1}^V \exp(u_x^T v_c) u_x$$

Distribute term across the sum:

$$\begin{aligned} &= u_o - \sum_{x=1}^V \frac{\exp(u_x^T v_c)}{\sum_{w=1}^V \exp(u_w^T v_c)} u_x \\ &= u_o - \sum_{x=1}^V P(x | c) u_x \\ &= \text{observed context vector} - \text{expected context vector} \end{aligned}$$

Thus center word is pulled towards words that are observed in its context, and away from those that are not. i.e.

$$v_o^{\text{new}} = v_o^{\text{old}} + \text{observed} - \text{expected}$$

Key Limitation of Word2Vec: Capturing cooccurrences inefficiently

Subword tokenization: break words into multiple word pieces

- Generally want more frequent words to be represented by fewer tokens, e.g. "the" should be a single token
- Can handle rare words better than word-level tokenization
- Can share parameters between related words (e.g., "opened", "opens", "opening")
- Reduce vocabulary size \rightarrow reduce num parameters, compute + memory

Developed for machine translation by Sennrich et al., ACL 2011, Later used in BERT, T5, RoBERTa, GPT, etc.

Relies on a simple algorithm called Byte Pair Encoding (Gage, 1994)

- Iteratively merges the most frequent pair of consecutive characters, until a fixed vocabulary size is reached

Language Modeling is the task of predicting what word comes next

$$P(X_1 = x_1, X_2 = x_2, \dots, X_n = x_n)$$

$$\underset{\text{assumption}}{=} P(X_1 = x_1) \times P(X_2 = x_2 | X_1 = x_1)$$

$$\times \prod_{i=3}^n P(X_i = x_i | X_{i-2} = x_{i-2}, X_{i-1} = x_{i-1})$$

$$= \prod_{i=1}^n P(X_i = x_i | X_{i-2} = x_{i-2}, X_{i-1} = x_{i-1})$$

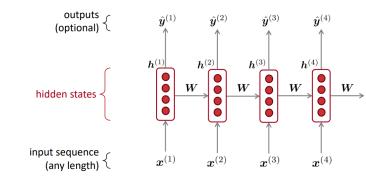
Unrolled Recurrent Neural Networks (RNNs)

hidden states

$$h^{(t)} = \sigma(W_h h^{(t-1)} + W_e e^{(t)} + b_1)$$

$h^{(0)}$ is the initial hidden state

Core idea: Apply the same weights W repeatedly



RNN Language Model

output distribution

$$\hat{y}^{(t)} = \text{softmax}(U h^{(t)} + b_2) \in \mathbb{R}^{|V|}$$

$$\hat{y}^{(4)} = P(\text{books}/\text{the students opened their laptops})$$

hidden states

$$h^{(1)} = \sigma(W_h h^{(0)} + W_e e^{(1)} + b_1)$$

$h^{(0)}$ is the initial hidden state

word embeddings

$$e^{(t)} = E x^{(t)}$$

words / one-hot-vectors

$$x^{(t)} \in \mathbb{R}^{|V|}$$

Training RNNs not much different from FFN, by unrolling the RNN we effectively create a feedforward network

Key Difference is Weight Sharing in RNNs: same weight is used at every time step; in FFN, each layer has its own weights.

Gradient clipping: if the norm of the gradient is greater than some threshold, scale it down before applying SGD update

Algorithm 1 Pseudo-code for norm clipping
Input: gradient g , threshold threshold
Output: clipped gradient \hat{g}
end if

If the gradient becomes too big, then the SGD update becomes too big:

$$\theta^{new} = \theta^{old} - \alpha \nabla_\theta J(\theta)$$

This can cause bad updates: we take too large a step and reach a weird and bad parameter configuration (with large loss)

If the gradient becomes too small, then the SGD update becomes too small:

If FFN, layers have different weights, some numbers might be small, some large, there is a balance

If numbers are < one, multiplying them together gets us a tiny product, eventually converges towards zero (**vanishing gradients**)

If numbers > one, we get the opposite problem, the product will increase exponentially in the length of the sequence (**exploding gradients**)

If FFN, layers have different weights, some numbers might be small, some large, there is a balance

If numbers are < one, multiplying them together gets us a tiny product, eventually converges towards zero (**vanishing gradients**)

If numbers > one, we get the opposite problem, the product will increase exponentially in the length of the sequence (**exploding gradients**)

If FFN, layers have different weights, some numbers might be small, some large, there is a balance

If numbers are < one, multiplying them together gets us a tiny product, eventually converges towards zero (**vanishing gradients**)

If numbers > one, we get the opposite problem, the product will increase exponentially in the length of the sequence (**exploding gradients**)

If FFN, layers have different weights, some numbers might be small, some large, there is a balance

If numbers are < one, multiplying them together gets us a tiny product, eventually converges towards zero (**vanishing gradients**)

If numbers > one, we get the opposite problem, the product will increase exponentially in the length of the sequence (**exploding gradients**)

If FFN, layers have different weights, some numbers might be small, some large, there is a balance

If numbers are < one, multiplying them together gets us a tiny product, eventually converges towards zero (**vanishing gradients**)

If numbers > one, we get the opposite problem, the product will increase exponentially in the length of the sequence (**exploding gradients**)

If FFN, layers have different weights, some numbers might be small, some large, there is a balance

If numbers are < one, multiplying them together gets us a tiny product, eventually converges towards zero (**vanishing gradients**)

If numbers > one, we get the opposite problem, the product will increase exponentially in the length of the sequence (**exploding gradients**)

If FFN, layers have different weights, some numbers might be small, some large, there is a balance

If numbers are < one, multiplying them together gets us a tiny product, eventually converges towards zero (**vanishing gradients**)

If numbers > one, we get the opposite problem, the product will increase exponentially in the length of the sequence (**exploding gradients**)

If FFN, layers have different weights, some numbers might be small, some large, there is a balance

If numbers are < one, multiplying them together gets us a tiny product, eventually converges towards zero (**vanishing gradients**)

If numbers > one, we get the opposite problem, the product will increase exponentially in the length of the sequence (**exploding gradients**)

If FFN, layers have different weights, some numbers might be small, some large, there is a balance

If numbers are < one, multiplying them together gets us a tiny product, eventually converges towards zero (**vanishing gradients**)

If numbers > one, we get the opposite problem, the product will increase exponentially in the length of the sequence (**exploding gradients**)

If FFN, layers have different weights, some numbers might be small, some large, there is a balance

If numbers are < one, multiplying them together gets us a tiny product, eventually converges towards zero (**vanishing gradients**)

If numbers > one, we get the opposite problem, the product will increase exponentially in the length of the sequence (**exploding gradients**)

If FFN, layers have different weights, some numbers might be small, some large, there is a balance

If numbers are < one, multiplying them together gets us a tiny product, eventually converges towards zero (**vanishing gradients**)

If numbers > one, we get the opposite problem, the product will increase exponentially in the length of the sequence (**exploding gradients**)

If FFN, layers have different weights, some numbers might be small, some large, there is a balance

If numbers are < one, multiplying them together gets us a tiny product, eventually converges towards zero (**vanishing gradients**)

If numbers > one, we get the opposite problem, the product will increase exponentially in the length of the sequence (**exploding gradients**)

If FFN, layers have different weights, some numbers might be small, some large, there is a balance

If numbers are < one, multiplying them together gets us a tiny product, eventually converges towards zero (**vanishing gradients**)

If numbers > one, we get the opposite problem, the product will increase exponentially in the length of the sequence (**exploding gradients**)

If FFN, layers have different weights, some numbers might be small, some large, there is a balance

If numbers are < one, multiplying them together gets us a tiny product, eventually converges towards zero (**vanishing gradients**)

If numbers > one, we get the opposite problem, the product will increase exponentially in the length of the sequence (**exploding gradients**)

If FFN, layers have different weights, some numbers might be small, some large, there is a balance

If numbers are < one, multiplying them together gets us a tiny product, eventually converges towards zero (**vanishing gradients**)

If numbers > one, we get the opposite problem, the product will increase exponentially in the length of the sequence (**exploding gradients**)

If FFN, layers have different weights, some numbers might be small, some large, there is a balance

If numbers are < one, multiplying them together gets us a tiny product, eventually converges towards zero (**vanishing gradients**)

If numbers > one, we get the opposite problem, the product will increase exponentially in the length of the sequence (**exploding gradients**)

If FFN, layers have different weights, some numbers might be small, some large, there is a balance

If numbers are < one, multiplying them together gets us a tiny product, eventually converges towards zero (**vanishing gradients**)

If numbers > one, we get the opposite problem, the product will increase exponentially in the length of the sequence (**exploding gradients**)

If FFN, layers have different weights, some numbers might be small, some large, there is a balance

If numbers are < one, multiplying them together gets us a tiny product, eventually converges towards zero (**vanishing gradients**)

If numbers > one, we get the opposite problem, the product will increase exponentially in the length of the sequence (**exploding gradients**)

If FFN, layers have different weights, some numbers might be small, some large, there is a balance

If numbers are < one, multiplying them together gets us a tiny product, eventually converges towards zero (**vanishing gradients**)

If numbers > one, we get the opposite problem, the product will increase exponentially in the length of the sequence (**exploding gradients**)

If FFN, layers have different weights, some numbers might be small, some large, there is a balance

If numbers are < one, multiplying them together gets us a tiny product, eventually converges towards zero (**vanishing gradients**)

If numbers > one, we get the opposite problem, the product will increase exponentially in the length of the sequence (**exploding gradients**)

If FFN, layers have different weights, some numbers might be small, some large, there is a balance

If numbers are < one, multiplying them together gets us a tiny product, eventually converges towards zero (**vanishing gradients**)

If numbers > one, we get the opposite problem, the product will increase exponentially in the length of the sequence (**exploding gradients**)

If FFN, layers have different weights, some numbers might be small, some large, there is a balance

If numbers are < one, multiplying them together gets us a tiny product, eventually converges towards zero (**vanishing gradients**)

If numbers > one, we get the opposite problem, the product will increase exponentially in the length of the sequence (**exploding gradients**)

If FFN, layers have different weights, some numbers might be small, some large, there is a balance

If numbers are < one, multiplying them together gets us a tiny product, eventually converges towards zero (**vanishing gradients**)

If numbers > one, we get the opposite problem, the product will increase exponentially in the length of the sequence (**exploding gradients**)

If FFN, layers have different weights, some numbers might be small, some large, there is a balance

If numbers are < one, multiplying them together gets us a tiny product, eventually converges towards zero (**vanishing gradients**)

If numbers > one, we get the opposite problem, the product will increase exponentially in the length of the sequence (**exploding gradients**)

If FFN, layers have different weights, some numbers might be small, some large, there is a balance

If numbers are < one, multiplying them together gets us a tiny product, eventually converges towards zero (**vanishing gradients**)

If numbers > one, we get the opposite problem, the product will increase exponentially in the length of the sequence (**exploding gradients**)

If FFN, layers have different weights, some numbers might be small, some large, there is a balance

If numbers are < one, multiplying them together gets us a tiny product, eventually converges towards zero (**vanishing gradients**)

If numbers > one, we get the opposite problem, the product will increase exponentially in the length of the sequence (**exploding gradients**)

If FFN, layers have different weights, some numbers might be small, some large, there is a balance

If numbers are < one, multiplying them together gets us a tiny product, eventually converges towards zero (**vanishing gradients**)

If numbers > one, we get the opposite problem, the product will increase exponentially in the length of the sequence (**exploding gradients**)

If FFN, layers have different weights, some numbers might be small, some large, there is a balance

If numbers are < one, multiplying them together gets us a tiny product, eventually converges towards zero (**vanishing gradients**)

If numbers > one, we get the opposite problem, the product will increase exponentially in the length of the sequence (**exploding gradients**)

If FFN, layers have different weights, some numbers might be small, some large, there is a balance

If numbers are < one, multiplying them together gets us a tiny product, eventually converges towards zero (**vanishing gradients**)

If numbers > one, we get the opposite problem, the product will increase exponentially in the length of the sequence (**exploding gradients**)

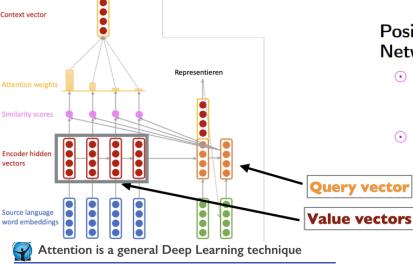
If FFN, layers have different weights, some numbers might be small, some large, there is a balance

If numbers are < one, multiplying them together gets us a tiny product, eventually converges towards zero (**vanishing gradients**)

If numbers > one, we get the opposite problem, the product will increase exponentially in the length of the sequence (**exploding gradients**)

If FFN, layers have different weights, some numbers might be small, some large, there is a balance

If numbers are < one, multiplying them together gets us a tiny product, eventually converges towards zero (**vanishing gradients**)



We have some **values** $h_1, \dots, h_N \in \mathbb{R}^{d_h}$ and a **query** $s \in \mathbb{R}^{d_q}$

Attention always involves:

1. Computing the **attention scores** $e \in \mathbb{R}^N$ There are multiple ways to do this
2. Taking softmax to get **attention distribution** a : $a = \text{softmax}(e) \in \mathbb{R}^N$

3. Using attention distribution to take weighted sum of values:

$$a = \sum_{i=1}^N \alpha_i h_i \in \mathbb{R}^{d_h}$$

thus obtaining the **attention output** a (sometimes called the context vector)

here are several ways you can compute $e \in \mathbb{R}^N$ from $h_1, \dots, h_N \in \mathbb{R}^{d_h}$ and $s \in \mathbb{R}^{d_q}$:

Basic dot-product attention: $e_i = s^T h_i \in \mathbb{R}$

- Note: this assumes $d_s = d_h$.

Multiplicative attention: $e_i = s^T W h_i \in \mathbb{R}$ [Luong, Pham, and Manning 2015]

- Where $W \in \mathbb{R}^{d_q \times d_h}$. W is a weight matrix. Perhaps better called "bilinear attention"

Reduced-rank multiplicative attention: $e_i = s^T (U^T V) h_i = (Us)^T (Vh_i)$

- For low rank matrices $U \in \mathbb{R}^{d_q \times d_q}$, $V \in \mathbb{R}^{k \times d_h}$, $k \ll d_h$, d_q

Additive attention: $e_{ij} = v^T \tanh(W_1 h_i + W_2 v) \in \mathbb{R}$ [Bahdanau, Cho, and Bengio 2014]

- Where $W_1 \in \mathbb{R}^{d_q \times d_h}$, $W_2 \in \mathbb{R}^{d_q \times d_v}$ are weight matrices and $v \in \mathbb{R}^{d_v}$ is a weight vector.
- d_v (the attention dimensionality) is a hyperparameter
- "Additive" is a weird/bad name. It's really using a feed-forward neural net layer.

Attention significantly improves NMT performance

- It's very useful to allow decoder to focus on certain parts of the source
- Attention solves the bottleneck problem
- Attention allows decoder to look directly at source; bypass bottleneck

Attention provides some interpretability

- By inspecting attention distribution, we can see what the decoder was focusing on
- We get (soft) alignment for free!
- This is cool because we never explicitly trained an alignment system
- The network just learned alignment by itself

In greedy decoding, usually we decode until the model produces an <END> token

- For example: <START> he hit me with a pie <END>

In beam search decoding, different hypotheses may produce <END> tokens on different timesteps

- When a hypothesis produces <END>, that hypothesis is complete.
- Place it aside and continue exploring other hypotheses via beam search.

$$\frac{1}{t} \sum_{i=1}^t \log P_{\text{LM}}(y_i | y_1, \dots, y_{i-1}, x)$$

Attention is All you Need

► Attention

- Self-attention
- Interspersing of Attention and Feedforward layers
- Multiple attention heads in parallel
- Encoder-Decoder Attention (cross-attention)
- Positional Encoding (attention is order-agnostic)
- Optimization Tricks
 - Residual network (ResNet) structure
 - LayerNorms
 - Scaled dot-product attention

Given: **Query vector** q and a set of **Key-Value** (k - v) vector pairs.

Output: a weighted sum of the **values**, where the weights are determined by the similarity of the **query** to the **keys**.

- A selective summary of the **values**, with respect to the **query**.

$$k = W_k \cdot x \quad v = W_v \cdot x$$

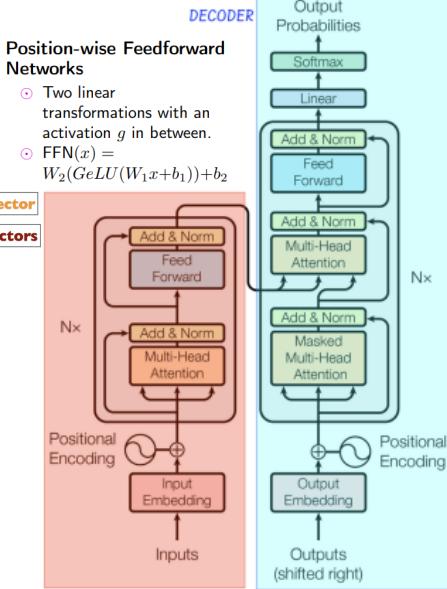
The query q can come from a separate input y :

$$q = W_q \cdot y$$

Or from the same input x , then we call it "self-attention":

$$q = W_q \cdot x$$

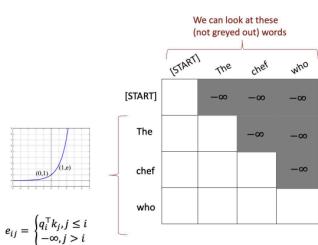
Attention is a general Deep Learning technique



$$A(Q, K, V) = \text{softmax}\left(\frac{QK^\top}{\sqrt{d_k}}\right) \cdot V$$

We want to mask with upper triangular matrix. In PyTorch use: `torch.triu`

Setting attention score to $-\infty$ ensures that attention probability will be zero.

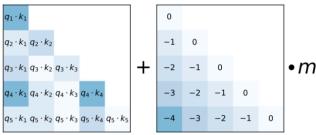


Encoder-Decoder Attention

- Each decoded token can "look at" the encoder's output (as in seq2seq w/ attention)
- This allows the decoder to focus on relevant parts of the encoder's output.

Absolute vs Relative Positional Encodings

- Absolute:** add an encoding to the input in hope that relative position will be captured
 - Vaswani et al. 2017: Positional Encodings
 - Fails to extrapolate to longer sequences
- Relative:** Explicitly encode relative position e.g., "I" is 3 words away from "run"
 - Press et al. 2021: "Attention with Linear Biases" (ALiBi), parameter-free linear bias added to dot-product similarity. Extrapolates to longer sequences.



Residual connections: training converges faster, and allows for deeper networks.



Layer normalization: Normalizes the outputs to be within a consistent range, preventing too much variance in scale of outputs

Challenges with RNNs

- Gradient vanishing and explosion**
 - Long-range dependencies are not captured
- Recurrence prevents parallel computation**
 - Slow training

Transformer Networks

- Facilitate long range dependencies**
 - No gradient vanishing and explosion
 - No recurrence: facilitate parallel computation

Even though most parameters and FLOPs are in feedforward layers, Transformers are still limited by quadratic $O(n^2)$ complexity of self-attention

Pretraining: The model learns general knowledge from a large, diverse dataset.

Finetuning: The model adapts this general knowledge to a specific task with a smaller, task-specific dataset.

ELMo is built using a deep Bi-LSTM model, where each word token's embedding is a function of:

$$h_k^{LM} = \overrightarrow{h_k^{LM}}, \overleftarrow{h_k^{LM}}$$

$\overrightarrow{h_k^{LM}}$ and $\overleftarrow{h_k^{LM}}$ represent the hidden states of forward and backward LSTMs for word k .

$$\text{ELMo}_k = \gamma \sum_{j=1}^L s_j h_{k,j}^{LM}$$

re:

$$h_{k,j}^{LM} \text{ is the } j\text{-th layer of the LSTM for word } k.$$

s_j : trainable scalar weights.

γ : trainable scalar parameter to scale the final embedding.

d: Nucleus Sampling

atic Decoding Method using Top-p Vocabulary

Given a probability distribution $P(x | x_{1:i-1})$, the **top-p vocabulary** $V^{(p)} \subset V$ is defined as the smallest set of tokens satisfying:

$$\sum_{x \in V^{(p)}} P(x | x_{1:i-1}) \geq p$$

Focus on (the nucleus) ensures sampling focuses on high-probability tokens while discarding the unreliable tail

Let $p' = \sum_{x \in V^{(p)}} P(x | x_{1:i-1})$. Original distribution is re-scaled to a new distribution, from which the next word is sampled:

$$P'(x | x_{1:i-1}) = \begin{cases} P(x | x_{1:i-1}) / p' & \text{if } x \in V^{(p)} \\ 0 & \text{otherwise} \end{cases}$$

At each time step, top-k sampling considers only the k token with the highest probabilities, re-scaling these probabilities and sampling from them

Sampling with Temperature

High temperature: more randomness; **Low temperature:** less randomness

Given the scores $u_{1:|V|}$ and temperature t , adjust softmax as:

$$p(x | V_i | x_{1:i-1}) = \frac{\exp(u_i/t)}{\sum_{j \in V} \exp(u_j/t)}$$

Negative Sampling

- Purpose:
- Efficiently approximate the softmax function by simplifying the computation.

Method:

- For each positive example (target-context pair), sample k negative examples (words not in the context).
- Update weights for both positive and negative samples.

Loss Function:

- Uses a binary logistic regression objective:
- $\mathcal{L} = - \left(\log \sigma(u_{w_i}^\top v_{w_i}) + \sum_{k=1}^k \log \sigma(-u_{w_k}^\top v_{w_k}) \right)$

where:

- σ is the sigmoid function.
- u_{w_i} : Output embedding of context word.
- v_{w_i} : Input embedding of target word.
- w_i : Negative sample words.

Teacher Forcing During Training

Definition:

- During training, the decoder is provided with the ground truth previous token y_{i-1} as input rather than its own previous prediction.

Benefits:

- Stability:**
 - Reduces the accumulation of errors over time steps.

Faster Convergence:

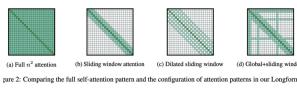
- The model learns the correct mapping between inputs and outputs more efficiently.

Limitations:

- Exposure Bias:**
 - At inference time, the model relies on its own predictions, which may lead to errors if it hasn't learned to handle them during training.

Alternatives:

- Scheduled Sampling:**
 - Gradually replaces ground truth inputs with model predictions during training.



part 2: Comparing the full self-attention pattern and the configuration of attention patterns in our Longformer

► Several pre-specified self-attention patterns that limit the number of operations

-training BERT: Masked Language Model

Masked Language Model (MLM) is one of the two pre-training tasks in BERT.

► Given an input sentence, BERT uses the context of surrounding words to predict the masked word:

$$P(\text{word}_i | \text{context}) = \text{softmax}(\mathbf{W} \cdot \mathbf{h}_i)$$

► Where:

- \mathbf{h}_i : The hidden state of the Transformer for word i .
- \mathbf{W} : Output projection matrix.

LMs place a distribution $P(x_i | x_1, \dots, x_{i-1})$

How do we generate text from these?

- Greedy decoding:** take the word with the highest probability at each step

- Beam search:** find the sequence with the highest probability

- Sample from the model:** draw x_i from that distribution

LMs place a distribution $P(x_i | x_1, \dots, x_{i-1})$. How do we generate text from these?

Maximization-based decoding:

- Greedy decoding:** pick the most likely word at each step
- Beam search:** keep track of the top k most likely sequences

Sampling-based decoding:

- Random sampling:** sample from the full distribution
- Nucleus sampling:** sample from the top p fraction of the distribution
- Top-k sampling:** sample from the top k most likely words

Encoder : $P(x_i | x_1, \dots, x_{i-1})$

- Bidirectional attention:** trained with masked language modeling
- Unidirectional attention:** trained to predict the next word
- To use in practice:** Ignore this probability distribution. Fine-tune the model for some other task $P(y | x)$

Decoder: $P(x_i | x_1, \dots, x_{i-1})$

- Unidirectional attention:** trained to predict the next word
- Alternative:** Can ignore this probability distribution and train a model for $P(y | x)$.
- To use in practice:** Generate text by sampling from this distribution

Momentum:

- Incorporates past gradients to smooth updates.
- Update rule:

$$v_t = \gamma v_{t-1} + \eta \nabla_\theta \mathcal{L}(\theta)$$

$$\theta_{\text{old}} = \theta_{\text{old}} - v_t$$

where γ is the momentum coefficient.

Sinusoidal Positional Encoding:

- Formulas:**

$$PE_{(\text{pos}, 2i)} = \sin\left(\frac{\text{pos}}{10000^{2i/d_{\text{model}}}}\right)$$

$$PE_{(\text{pos}, 2i+1)} = \cos\left(\frac{\text{pos}}{10000^{2i+1/d_{\text{model}}}}\right)$$

- pos : Position in the sequence.
- i : Dimension index.
- d_{model} : Model dimensionality.

Properties:

- Provides unique encoding for each position.
- Allows the model to generalize to sequences longer than those seen during training.
- Alternative:
 - Learned Positional Embeddings:**
 - Position embeddings are learned during training.
 - May not generalize to longer sequences.