

Driving progress towards specific tasks and capabilities

$$\arg \max_{\theta} \prod_{(w,c) \in \mathcal{D}} p((w,c) \in \mathcal{D}) \prod_{(w,c) \in \mathcal{D}'} p((w,c) \notin \mathcal{D})$$

CCG: Combinatory Categorical Grammar

Propositions suggested by an utterance, but not explicitly expressed

Gricean Maxims; Common Ground

General principles we believe we mutually hold about how what kinds of utterances we should add to conversation given what's been said so far:

During interaction, we maintain some representation of what we believe is mutually known by conversation participants

Elements of Scenario Design

Interaction dynamics also depend heavily on the properties of the context itself:

- Incentive structure
- Environment design — perception and action, novelty
- Participants — how many, any existing structures among them, roles, a priori asymmetries
- Communication channel

Work in computational linguistics, psycholinguistics, and cognitive science aims to characterize the relationship between scenario design and linguistic behavior

Multilingual MLP and its challenges

For any task we expect out of language technologies, they should work for any language

Data Modality, Data Scarcity, Dialectical Variation, Speech System, Morphology, Lexical Semantics, Syntax, Semantics, Idioms, Difference in Language Use, Change

Autoregressive language modeling:

- The probability of a sequence is a product of local token probabilities
- The probability of a token depends on the ones that came before it

Masked language modeling:

- The probability of a sequence is a product of local token probabilities
- The probability of a token depends on the ones that came before and after it

$$p(\bar{x}) = \prod_{i=1}^n p(x_i \mid x_1, \dots, x_{i-1}, x_{i+1} \dots x_n)$$

Let's make a Markov assumption:
The probability of word at index i only depends on the $n - 1$ words that came before it

$$p(X_i = x) = p(x \mid x_1, \dots, x_{i-1})$$

$$\approx p(x \mid \text{Preceding (n-1)-gram})$$

$$\approx \frac{C(\text{green})}{C(\text{blue})}$$

For an n -gram language model, we need to store counts for:

- All sequences of length n (\mathcal{V}^n)
- All sequences of length $n - 1$ (\mathcal{V}^{n-1})

What if the count of the target n -gram is 0?

- Solution: add a small number to the count for every n -gram (aka “smoothing”)

What if our $n-1$ -gram prefix has a count of 0?

- Solution: condition on a shorter n -gram prefix (e.g., the previous $n-2$, or $n-3$, etc.) instead (aka “backoff”)

Can't learn anything from the counts of n -grams containing similar words

Without a big n , cannot handle long-distance dependencies

Measure of Fit (拟合优度量)

Likelihood: probability of the data under our model

$$\prod_{i=1}^M p_{\theta}(\bar{x}_i)$$

Negative log likelihood (fixes float underflow)

$$-\sum_{i=1}^M \log p_{\theta}(\bar{x}_i) = -\sum_{i=1}^M \sum_{j=1}^N \log p_{\theta}(x_j^i \mid x_1^i, \dots, x_{j-1}^i)$$

Perplexity: inverse probability of data, normalized by number of tokens in the dataset

Adjust the Temperature when generating

Operation: modify the logits before computing probabilities

$$s(w) = f(w \mid x_1, \dots, x_{i-1}; \theta) \leftarrow \text{logits}$$

$$p(X_i = w \mid x_1, \dots, x_{i-1}) = \frac{\exp(s(w)/\tau)}{\sum_{w' \in \mathcal{V}} \exp(s(w')/\tau)}$$

A better approximation for global argmax: **beam search**

- During generation, we maintain a “beam” of n sequences instead of just one
- At each generation step i ,
 - We select the n most likely next tokens \mathcal{X}_i for each prefix, and create n more sequences
- Then we look at all the n^2 sequences so far, and discard all but the n most likely sequences
- At the end, we select the sequence that has the highest probability among the set

Masking Out Wordtypes: (&constrained)

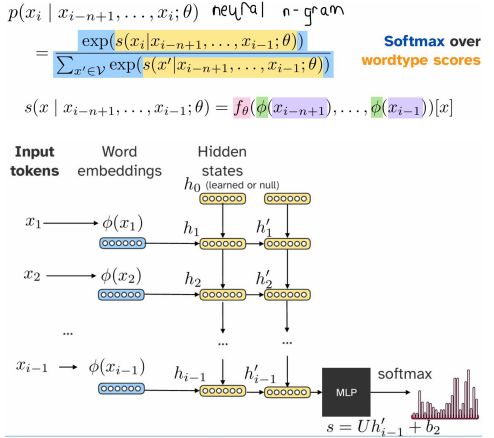
e sampling top-k sampling top-p (nucleus) sampling

Identify the set of n tokens \mathcal{E} such that $\forall x \in \mathcal{E}$, $p(x \mid x_1, \dots, x_{i-1}) \geq \epsilon$

Identify the set of k tokens \mathcal{K} that have the highest probabilities under $p(X_i \mid x_1, \dots, x_{i-1})$

Identify the set of n tokens \mathcal{P} that have the highest probabilities under $p(X_i \mid x_1, \dots, x_{i-1})$ and their cumulative probability is p

Similar to before: given a set of possible continuations $\mathcal{C} \subseteq \mathcal{V}$ we will set the probabilities of all other tokens to 0, then renormalize using $\sum p(X_i \mid x_1, \dots, x_{i-1})$

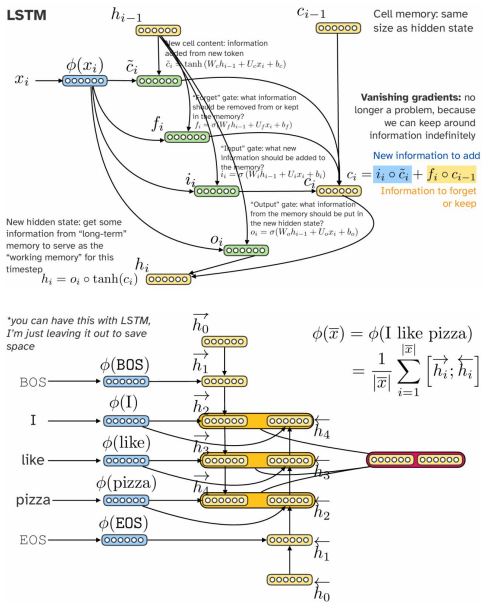


$$\theta^* = \arg \min_{\theta} \left(- \sum_{\bar{x} \in \mathcal{D}} \sum_{i=1}^{|\bar{x}|} \log p(x_i \mid x_1, \dots, x_{i-1}; \theta) \right)$$

Vanishing gradients: gradient signal from many timesteps in the future is negligible compared to gradient signal from nearby tokens

Exploding gradients: if the activations are too big, the gradients become too big, and the values of the weights become too big, eventually getting values closer to ∞

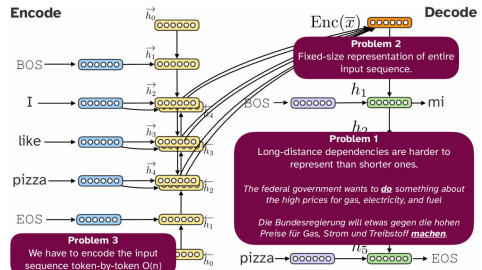
Can solve using gradient clipping (if the norm of the gradient is greater than a threshold, clamp its value or scale it down)



$$\text{Enc}(\bar{x}) = \text{Pool} (h_1, \dots, h_{|\bar{x}|}) \in \mathbb{R}^d$$

“Pooling” functions map from a sequence of items of type t to a single item of type t

Seq2seq (With challenges); Attention



We encoded our input sequence into hidden states: $h_1, \dots, h_{|\bar{x}|} = \mathbf{h} \in \mathbb{R}^{d \times |\bar{x}|}$
Now we want to predict the next word y_{i+1}
We have access to the previous decoder hidden state g_i

First, compute attention scores for each input hidden state using similarity between query (g_i) and keys (\mathbf{h}): $s_i = a(g_i, \mathbf{h}) \in \mathbb{R}^{|\bar{x}|}$
Then, take softmax of attention scores to get a distribution over keys: $\alpha_i = \text{softmax}(s_i) \in \Delta^{N_{1:|\bar{x}|}}$
Finally, compute a weighted sum of values (\mathbf{h}) using this distribution

$$c_i = \sum_{j=1}^{|\bar{x}|} \alpha_{i,j} h_j \in \mathbb{R}^d$$

Use the weighted sum to predict the next word:
 $p(Y_{i+1} \mid \bar{x}, y_1, \dots, y_i) = \text{softmax}(f(c_i, g_i))$
Use the weighted sum to update the decoder hidden state: $g_i = g(g_{i-1}, c_{i-1}, y_i)$

$$h_1, \dots, h_{|\bar{x}|} = \mathbf{h} \in \mathbb{R}^{d \times |\bar{x}|}$$

$$s_i = a(g_i, \mathbf{h}) \in \mathbb{R}^{|\bar{x}|}$$

$$\alpha_i = \text{softmax}(s_i) \in \Delta^{N_{1:|\bar{x}|}}$$

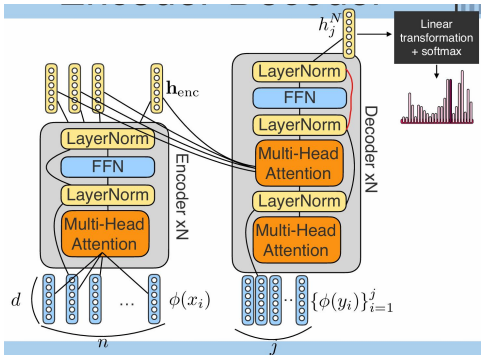
$$c_i = \sum_{j=1}^{|\bar{x}|} \alpha_{i,j} h_j \in \mathbb{R}^d$$

$$p(Y_{i+1} \mid \bar{x}, y_1, \dots, y_i) = \text{softmax}(f(c_i, g_i))$$

$$g_i = g(g_{i-1}, c_{i-1}, y_i)$$

• Sinusoidal embeddings

$$p_i = \begin{bmatrix} \sin(i/10000^{2\pi/d_1}) \\ \cos(i/10000^{2\pi/d_1}) \\ \vdots \\ \sin(i/10000^{2\pi/d_d}) \\ \cos(i/10000^{2\pi/d_d}) \end{bmatrix}$$



But we don't want the model to learn to rely on “future” words in the sequence!

Solution: during training, set $s_{ij} = -\infty$ if i (query index) $< j$ (key/value index)

BERT: Instead: randomly sample some tokens to replace with a MASK placeholder token in the input, then train the model to predict those words