

Optimization: compute predictions (at current prediction) → compute the loss → backpropagate the error (which direction to change) → update parameters → zero the gradients. Analogy of ∇ math notation: given data training steps

$L(\theta, x, y) = -\log p_{\text{ml}}(x_i | y)$, find params that minimize loss for data: $\hat{\theta} = \arg \max_{\theta} \sum_i L(\theta, x_i, y_i)$

\Rightarrow loss function: quantity the difference.

$\text{MSE} = \frac{1}{n} \sum_i (y_i - \hat{y}_i)^2$. \approx cross-entropy loss:

$$-\frac{1}{n} \sum_i y_i \log(\hat{y}_i).$$

\Rightarrow Gradient: vectors of partial derivatives of the loss function which suggests the rate & direction to update params (opposite to gradient, increase, \rightarrow decrease) in order to minimize loss function. ($\frac{dL}{d\theta}$ math.)

\Rightarrow Back-propagation: calculate gradients by propagating the errors backward through the net from output to input, updating each param along way.

\Rightarrow Stochastic GD: use random subsets (mini-batches) of data instead of the entire dataset to iteratively calculating loss, gradients & update params. (Vanilla-BatchGD)

$\text{ANN} \Rightarrow$ Unit score: $z = b + \sum_i w_i x_i = b + w \cdot x$.

\Rightarrow Anyfunc: $a = \delta(z) = \delta(w \cdot x + b)$.

\Rightarrow Activation function: Perceptron: $f(z) = \delta z > 0$.

\vee Sigmoid (0-1): $f(z) = \delta(z) = \frac{1}{1 + \exp(-z)}$

\vee Hyperbolic tangent: $f(z) = \tanh(z) = \frac{\exp(z) - \exp(-z)}{\exp(z) + \exp(-z)}$.

tan \Rightarrow hidden layers: intermediate layers that performs computations.

\Rightarrow Score: last layer output, often raw score or log-odds (that would transform to prob.)

\Rightarrow Loss: concept: a function that assigns to each input

X a probability distribution over S given param $\theta \in \Theta$ (Input conditions) \Rightarrow Neural LM: an LM realised as a neural network.

\Rightarrow Causal LM: a function that maps an initial token sequence to a next-token distribution.

$\text{LM: } \text{run} \rightarrow \Delta(T) \text{ for vocab-set of tokens.}$

\Rightarrow Surprisal: $-\log(p_{\text{ml}}(w_{i+1} | w_i:n))$ next-token prob

\Rightarrow Goodness of fit: 1) perplexity $p_{\text{ml}}(w_i:n) = p_{\text{ml}}(w_i:n) - \ln \frac{1}{n}$ more weights to learn, so \downarrow higher t-data demand.

2) Avg. surprisal: $\log \text{of pp: } -\frac{1}{n} \log p_{\text{ml}}(w_i:n)$.

\Rightarrow Autograd generation: generate data based on previous elements, left-to-right, no bidirectional.

(inputs → embeddings - RNN - Softmax - sampled word maintain a hidden state that capture contents)

\Rightarrow Denoising scheme: 1) pure sampling: pure NLP.

2) Greedy decoding: best word (highest NLP). 3) softmax sampling: k-most likely words. (repack and optimise wrt particular task. \Rightarrow Maintains memory: give enough info to contextualized embedding. \Rightarrow Contextualized embeddings between vectors is useful.)

\Rightarrow Auto-regressive LM: one that only have access to previous tokens (and output become inputs)

\Rightarrow Masked LM: predicts missing tokens in a sequence by masking some tokens and train the model to fill in the blank. Eval: downstream NLP tasks.

\Rightarrow Bidirectional ~: predicts tokens based on preceding and succeeding context in a sequence. Eval: all etc.

\Rightarrow Training regimes: 1) teacher forcing: LM is fed by masking some tokens and train the model to

the word seq. not predictions. 2) autoregressive forcing: i.e. free running. 3) curriculum mode: auto-regency generates a seq.

\Rightarrow RNN: NN designed to process sequential data where the output depends on the current & past input

\Rightarrow FFNN: info flow in 1D. \Rightarrow Convolutional NN: specialised for processing grid-like data (3, 1, 1).

\Rightarrow GRU: recurrent RNN: capture more context while avoiding sparsity. hidden layers represent word meaning. storage, compute issues. word window

\Rightarrow LSTM: can handle infinite memory instead of finite memory. instead of current hidden state, it's hidden state from previous time step, do every step (forward). \Rightarrow It's great for embedding learning & powerful too

\Rightarrow BiRNN (Bidirectional RNN): Bidirectional, out-train encoder, no decoder.

\Rightarrow S-Vanilla RNN: can handle infinite memory instead of finite memory. instead of current hidden state, it's hidden state from previous time step, do every step (backward). \Rightarrow It's great for embedding learning & powerful too

\Rightarrow Inhibition: At each time step t , we have a hidden state h_t and cell state C_t (store h_t into), LSTM erases, writes, and reads into

from the cell. \Rightarrow New undergoes a nonlinear activation, just add grad.

\Rightarrow LSTM short-L: LSTM respond to learn from the context it's been next token ($\text{LSTM}(\text{hidden} \times 3)$). 2) Prepped LM (classifiers): fine-tuned, predict to pleasure ($\text{LSTM}(\text{hidden} \times 3)$). 3) LSTM-based RNN (lagers): algorithm using LM (sophisticated pumping, cheap inputs, $\text{LSTM}(\text{hidden} \times 3)$).

\Rightarrow LSTM short-L: LSTM need to learn from the context it's been next token ($\text{LSTM}(\text{hidden} \times 3)$). 2) Prepped LM (classifiers): fine-tuned, predict to pleasure ($\text{LSTM}(\text{hidden} \times 3)$). 3) LSTM-based RNN (lagers): algorithm using LM (sophisticated pumping, cheap inputs, $\text{LSTM}(\text{hidden} \times 3)$).

\Rightarrow Summary: 1) LM need native genuine "learning" as structure of simple examples and target input. 2) k-shot examples and phone structures when part of cell state is used to compute output (like new sum).

5) RNN: H_t captures long-range dependencies well & TGRNN: \Rightarrow Intuition: 1) more weights to learn, so \downarrow higher t-data demand. 2) slower to train. 2) full batch vs. mini-batch. 3) loss function from exploding gradients.

\Rightarrow Bi-LSTM: we both tritute, starts content. \Rightarrow Intuition: 1) more top of each other, these stacking hidden layers provides increased abstractions (old memory). Bi-LSTM-FWD-BKWARD

\Rightarrow Intuition: at each step the decoder pays attention to all of the encoders' hidden states.

\Rightarrow Self-Attention: Within same seq., each word is transformed into a rich abstract rep (context embed) based on the weighted sums of other words. (How much should each word be influenced by its neighbors, the explicit positionality).

\Rightarrow Transformer: $\text{MLP} \times n$ that uses self-att mechanism to static embeddings, \Rightarrow slower to train. \Rightarrow Trained embeddings, \Rightarrow context-aware & more powerful than static embeddings. \Rightarrow Intuition: MLP that uses self-att mechanism to recurrent layers to process data.

\Rightarrow Decoder: component that receive inputs. It produces output sequence to optimise based on actual output.

\Rightarrow Encoder: component that receive inputs. It produces new seq. of They are identical to encoders, rather have an additional att head in b/w the self-att and FFN layers. This focuses on part of the encoder's repetitions.

\Rightarrow Decoders: ~ produce output. It generates new seq. of text. They are identical to encoders, rather have an additional att head in b/w the self-att and FFN layers. This focuses on part of the encoder's repetitions.

\Rightarrow BERET (Bidirectional Trans): Bidirectional, out-train encoder, no decoder.

\Rightarrow 2 training objectives: 1) Predict the masked word. 2) 2 sentences are fed in at a time. Predict if 1st. \Rightarrow Intuition: 1) Great for embedding learning & powerful too

\Rightarrow 2) transfer learning to other task. 2) no dev, not costly generalization

\Rightarrow 2) kind of LMs: 1) core-LMs (foundations): predict statically y_1, \dots, y_N is embedding for input x_1, \dots, x_N . 2) word embeddings: type-based hope, contextualized token embeddings. token-level next token ($\text{LSTM}(\text{hidden} \times 3)$). 3) Prepped LM (classifiers): predict to pleasure ($\text{LSTM}(\text{hidden} \times 3)$). 3) LSTM-based RNN (lagers): algorithm using LM (sophisticated pumping, cheap inputs, $\text{LSTM}(\text{hidden} \times 3)$).

\Rightarrow Summary: 1) LM need native genuine "learning" as structure of simple examples and target input. 2) k-shot examples and phone structures when part of cell state is used to compute output (like new sum).

