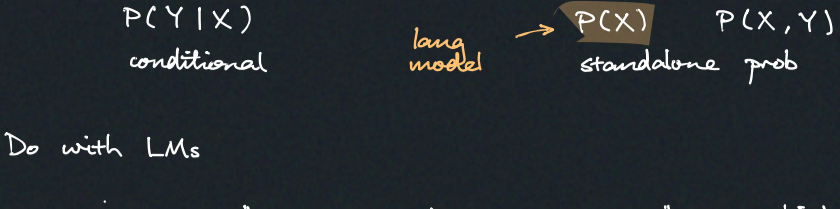


Lec 3 Language Modelling (LM)

The task



Do with LMs

- scoring "Jane went to the store" \rightarrow high score
"store to Jane went the" \rightarrow low
- generate - sample by $\tilde{x} \sim P(X)$
- find most probable sentence in prob space

Applying the LM

- \rightarrow QA \hookrightarrow prediction } see open-llm for rankings on these
- "Q: where is CMU located? A: Pittsburgh." \rightarrow high score
- "Q: where is CMU located? A: Birmingham." \rightarrow low score?

\rightarrow Prompt question & gen continuation

"Q: where is CMU located? A: _____"

\rightarrow Categorise

"This is great! Star rating _____"
"This is great! Star rating 5" \rightarrow score

\rightarrow Correct grammar

replace words, paraphrase etc.

Autoregressive LMs

\leftarrow GPT, LLaMA, etc.

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

next token context

So the problem is predicting $P(x_i | x_1, \dots, x_{i-1})$

Why not just do $P(X)$?
 \rightarrow the input space becomes $|V|^n$ where V is vocab and n is sentence length

Alternative: [MASK] LMs.
 problems:
 - doesn't give prob of sequence
 - hard to do generation

Alternative: energy-based
 very advanced

Unigram LMs

\hookrightarrow potential sols

Assume $x_i \perp x_1, \dots, x_{i-1}$
 $\Rightarrow P(x_i) \approx P(x_i | x_1, \dots, x_{i-1})$

$$P_{un}(x_i) = \frac{\text{count}(x_i)}{\sum_x \text{count}(x)}$$

\leftarrow count

- Problems
- $P(\text{unknown word}) = 0$
 \hookrightarrow segment word
 \hookrightarrow model unknown words by characters

Log space multiplication

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

\uparrow ends to underflow \uparrow better

Notation: $\theta_{x_i} := \log P(x_i)$ so unigram has $|V|$ params

n-gram model

Generalised

$$P_{LM}(x_i | x_{i-n+1}, \dots, x_{i-1}) := \frac{c(x_{i-n+1}, \dots, x_i)}{c(x_{i-n+1}, \dots, x_{i-1})}$$

- Problems
- Unseen context seq \Rightarrow prob is 0
 \hookrightarrow fall back to shorter context
 interpolate with shorter gram
 so we ensemble different-gram models
- } Idea of ensembling

$$P(x_i | x_{i-n+1}, \dots, x_{i-1}) = \lambda P_{3\text{-gram}} \dots + (1-\lambda) P_{4\text{-gram}} \text{ etc.}$$

- different varieties
- Goodman 1998 smoothing (for count-based models)
 - Additive / Dirichlet
 - Discounting
 - Kneser-Ney \leftarrow state-of-the-art when n-gram was popular.

fallback	bigram	
the 0.5	the the 0	
box 0.3	the box 0 \rightarrow 1	0.33
is 0.2	the is 0	

\exists solutions for individual issues but combining is hard and ppl used neural instead

\uparrow
 if more evidence for these, we weight them more

- Doesn't share among similar words
- Not condition on interleaving word
- Long distance dependencies not captured

n-gram

fast
learns low-freq better
scalable

neural

usually better performance

} kenlm for n-gram

LM evaluation

Log-likelihood (LL)

$$LL(X_{\text{test}}) = \sum_{x \in X_{\text{test}}} \log P(x)$$

WLL (per word LL)

LL norm. by # words

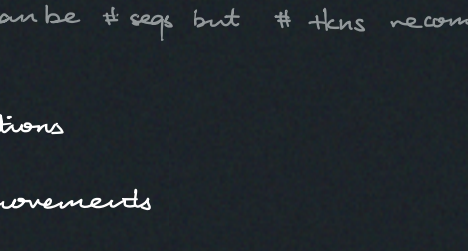
NLL

LL but negative

Entropy

$$H(X_{\text{test}}) = \frac{1}{\sum_{x \in X_{\text{test}}} |X|} - \sum_{x \in X_{\text{test}}} \log_2 P(x)$$

} related to info theory? prob model can be used to do compression



Perplexity

$T := X_{\text{test}}$

$$PPL(T) = 2^{H(T)} = e^{-WLL(T)}$$

Intuitively: # of samples needed to sample from probs to get next token right } i.e. guess

Note: make sure to divide to account for vocab size etc. before comparing

Feedforward neural model

1. Calculate some feature of context } SGD-optimised

2. Predict probs for next token

3. Gradient descent

Solves?:

- x Doesn't share among similar words
- ✓ Not condition on interleaving word
- x Long distance dependencies not captured

Feed forward neural model

Solves?:

- ✓ Doesn't share among similar words
- ✓ Not condition on interleaving word
- x Long distance dependencies not captured

Other LM desiderata

- ▷ Calibration: "model knows when it knows"
 \hookrightarrow want $P(a \text{ is correct} | P_{\text{model}} = p) = p$
 i.e. if model says 7% right then 7% right
 \hookrightarrow in practice, exact probs are hard, so use buckets
 \hookrightarrow note: calibration & accuracy not necessarily high correlation

- To compute prob of output i.e. confidence
- prob of answer - if single answer \rightarrow easy
 - if multi answer \rightarrow ?
 - prob of answer + paraphrases of the answer
 i.e. sum multiple acceptable answer
 - sample multiple outputs, count # of correct answers
 \hookrightarrow useful if LM doesn't return prob
 - ask the model what it thinks

\hookrightarrow Xiong + 2023 good comparison paper

▷ Efficiency

- Parameter count
- Memory usage
 - \hookrightarrow peak
 - \hookrightarrow model load
- Latency \approx encode time
 - \hookrightarrow until first token
 - \hookrightarrow until finished generating
- Throughput

Efficient NN implementation

- Batching
 \hookrightarrow batch size estimate can be # seqs but # tokens recommended
- GPU
- Avoid
- Avoid repeating operations
- Reduce operations
- Reduce CPU \leftrightarrow GPU movements