## Lec 6    Generation Algorithms

Given model $M$, pretrained, lots of params
└ defines a conditional prob distribution
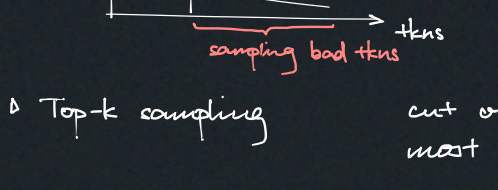
$M("2 + 2 = ")$         $M("favourite colour?")$


4                        green ← less confident

hallucination ← stars we always happen

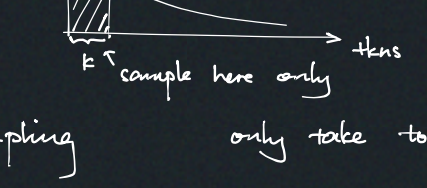Problem: how to get good output given hallucinations?

### * Sampling

▷ Ancestral sampling    $y_i \sim P(y_i | X, y_1, ..., y_{i-1})$
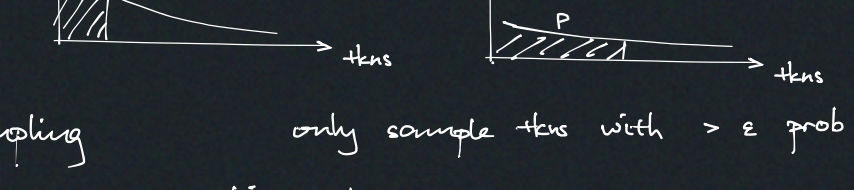└ sample according to dist at time step

Problem: long tail with 32,000 tokens


sampling bad this

▷ Top-k sampling    cut off tail and sample from the
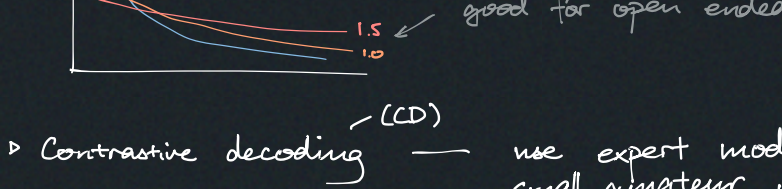most probable


k   sample here only

▷ Top-p (nucleus) sampling    only take top p prob



▷ ε - sampling    only sample tkns with > ε prob

▷ Temperature: morphing dist
before softmax, divide by temp


temp = 0.5
good for factual

1.5    good for open ended

▷ Contrastive decoding — (CD)    use expert model to "fix" dist from
ideas: if this is larger, probably    small, amateur model
because expert learnt sth the        then choose output the expert find more
amateur didn't                        likely than the amateur
also if both model have degeneracy   i.e. expert log prob − amateur log prob
case, CD cancels them out

### * Mode decoding method

When we want to search for the most likely seq

▷ Greedy decoding: $y_i = \text{argmax } P(y_i | X, y_1, ..., y_{i-1})$

Problem: local optimal not always global opt.

▷ Beam search

Better approximation than greedy

Problems — Less diversity

▷ Diverse beam search

→ score candidates differently st we get more diversity
eg if 2nd rank very similar to 1st rank, lower its
score and encourage sth different
→ Can define similarity func

▷ Stochastic beam search

→ keep scoring the same but sample from them
without replacement

### # Minimum Bayes Risk (MBR)

Do we really want the highest thing?

| | | |
|---|---|---|
| The cat sat down | 0.3 * | Really want k over a? |
| -- -- ran away | 0.001 | a similar & add up more |
| -- -- grew wings | 0.25 * | than a |
| -- -- spruced off | 0.2 * | Maybe * has lower risk in |
| | | terms of semantic |

Estimate prob: keep sampling, monte carlo

Estimate risk: eval agreement btwn candidates

MBR:  $\hat{y} = \text{argmax} \sum_{y' \in y_h} G(y, y')$    low risk
       $y \in y_h$          high prob
                            samp to pseudo-references

Sim metric $G$ ex: ROUGE, BEER, ...

Problems: sampling many things expensive

    Beam search known to
    sometimes degrade downstream
    metric perf.
    Because true mode sometimes
    bad, ex: "I don't know" may
    have very high prob
    Also, training may not lead
    to model being good

▷ Output ensembling

→ multiple models, rather than outputs MBR from same
model, we look at sim btwn models

▷ Self-consistency

1. Prompt for chain of thought
2. Sample many outs          ⎤
3. Extract answers from each  ⎬ MBR-ish
4. Return most frequent       ⎦

### # Constrained Generation

Eg  $M("Suggest activity")$  but don't want climbing
... doesn't quite work on LLMs

▷ Decode-side solution
└ set prob of "climbing" to 0
... but what about "bouldering" or other semantically similar
results

▷ Sample many, map check if about climbing, then filter

... expensive
also hard if model really likes climbing

▷ FUDGE — Future Discriminator       train additional
└ at decode step, modify dist by    classification model
estimation on whether
the seq will end up being formal vs informal, eg
estimator can be trained on labelled data

▷ Reinf. Learning from Human feedback  RLHF
                                       └ need more data source
Aligned LM combines original LM & reward model

▷ Reward-augmented decoding — no reinforcement learning
└ each decode step, predict future reward & adjust    involved
tkn probs, like FUDGE

→ Maybe just ask the model itself whether it violates the constraints

### * Human-in-the-loop Decoding

Decoding with some human intervention in between
- World Craft — story gen with human interaction
- Human edits to (fine grain replacement) and model continues
                 (specific edit request)       (apply own style)

- Regeneration button

- Multiple generations for human to choose

Tree of thought paper — use model to decide which option
to keep generating on

### * Practical Considerations

Decoding takes most time generally

▷ Speculative decoding
└ use small model to speculate next tkns
larger model rejects or accept
if usually accept, small model jumps forward quickly
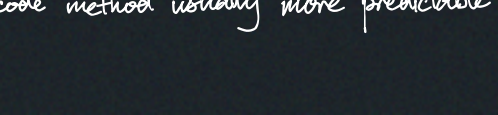while large model verifies block by block

▷ Fast inference / constraint inference libs

vLLM    disco    Outlines ...

### # Summary

- At decode step, adjust prob dist
- In decode process, choose btwn branches

Useful for  - Match constraints
            - Use data source not used in training
            - Compensate for worse model eg GPT2

Tradeoffs       diversity


             quality      speed

Decoding behaviour matters!

Decode method usually more predictable than fine tuning