

LLaMa 2

→ we copy the kv-heads to match the number of query heads

SwiGLU Function

$$\text{LLaMa 2 FFn SwiGLU}(x, W, V, W_2) = (\text{Swish}_\beta(xW) \otimes xV) W_2$$

$$\beta = 1, \text{Swish} = \text{SiLU}$$

INFERENCE :-

Logits are the output of the last linear layer in

→ transformer model. They are unscaled probabilities, but not clearly as they don't sum up to 1.

→ The softmax scales them to make them sum up to 1.

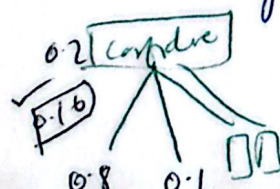
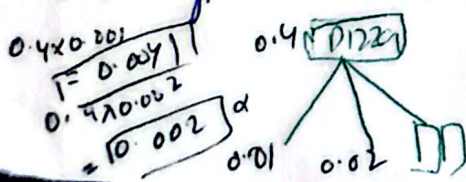
1) Greedy

→ select the highest probability

→ easy to implement
→ performs poor
→ initially wrong
guess can lead to bad results

2) Beam-Search ($k=2$)

→ select Top(2) or k with highest probability, we make 2 prompt in next step and ask the model to ^{output} generate next token and calculate cumulative reward for both paths and select the highest cumulative reward.



→ Top-k=2

and same for

increase inference time
perform better.

Temperature:

- The idea is to scale the logits before applying softmax
- low temperature → more confident model → high gap b/w L&H probabilities
- high temperature → less confident model → low gap b/w L&H probabilities
- probabilities / temperature → amplify & deamplify probs

Random Sampling

- after softmax, as it is distribution we can take the number (sample) ^{high-probability} from it.

issues e.g. [0.12, 0.07, 0.80]

(first we choose 12%, then 7%, then 80%, the higher the probability the more likely to be chosen).

Problem: the very little probability it may happen that we choose tokens that are bullshit. so we end up choosing low probabilities.

- Top-K = to avoid bad tokens, we remove them
- sort all the logits, and keep the highest K, and apply softmax to rest. But still low probabilities can make biased on distribution. we want a bit randomness at mid level.

Top-P: we keep only the items with highest probabilities such that their cumulative probability is greater than or equal to P (parameter). This let us deal with all types of distribution (flat or \mathcal{L}_1 ($P=0.5$), area under distribution.