

Language Model (LM):

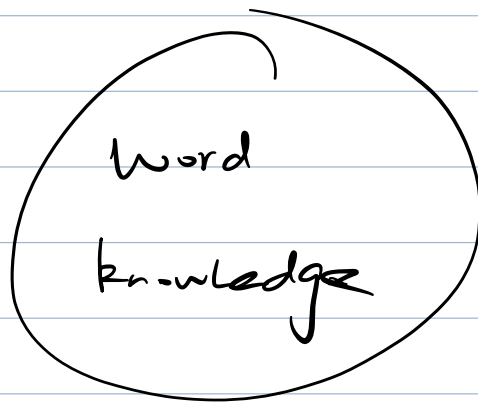
probability distribution P over sequences of tokens $X_{1:L}$

$$P(X_1, \dots, X_L)$$

sequence of tokens

(intuitively tells us how "good" a sequence of tokens is.)

Syntactic
knowledge



(linguistic capabilities)

Autoregressive language models

chain rule of probability

$$P(X_{1:L}) = P(X_1) P(X_2 | X_1) P(X_3 | X_1, X_2) \dots$$

$$P(X_L | X_{1:L-1})$$

$$= \prod_{i=1}^L P(x_i | x_{1:i-1})$$

Generation: generate an entire sequence

$x_{1:L}$ from p

for $i = 1, \dots, L$:

$$x_i \sim p(x_i | x_{1:i-1})^{\frac{1}{T}}$$

$T \geq 0$ temperature parameter

that controls how much randomness
(amount of variability in generation)

- $T = 0$: deterministically choose the most probable token x_i at each position i
- $T = 1$: sample "normally" from the pure language model
- $T = \infty$: sample from a uniform distribution over the entire vocabulary \square

Annealed conditional probability distribution

$$P_T(x_i | x_{1:i-1}) \propto p(x_i | x_{1:i-1})^{\frac{1}{T}}$$

Conditional Generation:

specify some prefix sequence $X_{1:i}$

(prompt)

sample the rest $X_{i+1:L}$

(completion)

the, mouse, ate $\xrightarrow{T=0}$ the, cheese

prompt

completion

$X_{1:i}$

$X_{i+1:L}$

Language Models:

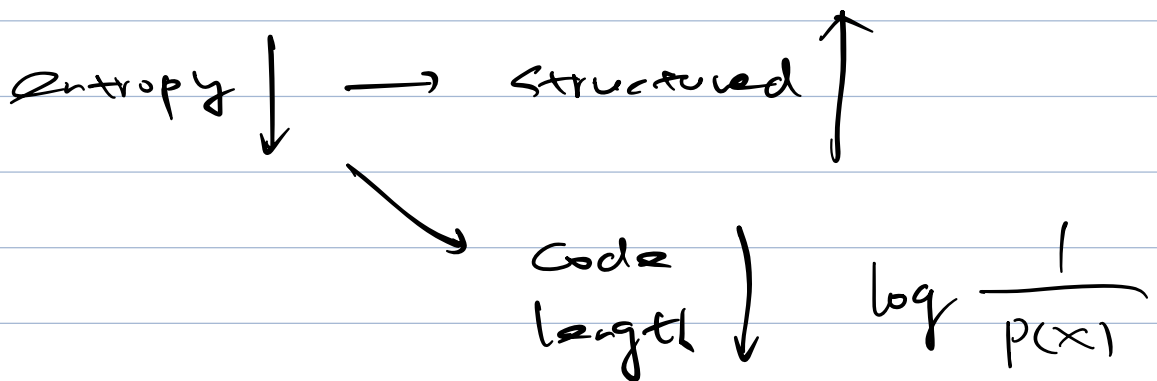
Claude Shannon

Information theory 1948

Entropy:

$$H(p) = \sum_x p(x) \log \frac{1}{p(x)}$$

Measure the expected number of bits any algorithm needs to encode (compress) a sample $x \sim p$ into a bitstring



Cross Entropy:

$$H(p, q) = \sum_x p(x) \log \frac{1}{q(x)}$$

Measure the expected number of bits (nats.) needed to encode a sample $x \sim p$ using the compression scheme given by the model q

(representing x with a code of length $\frac{1}{q(x)}$)

$$H(p, q) \geq H(p)$$

$\left\{ \begin{array}{l} p : \text{samples from the true data distribution} \\ q : \text{model} \end{array} \right.$

Shannon Game (human language model)

Shannon n-gram models 1948

Noisy channel model

$$p(\text{text} | \text{speech}) \propto \underbrace{p(\text{text})}_{\text{language model}} \underbrace{p(\text{speech} | \text{text})}_{\text{acoustic model}}$$

N-gram models

token X_i only depends on the last $n-1$

characters $X_{i-(n-1):i-1}$

$$P(X_i | X_{1:i-1}) = P(X_i | X_{i-(n-1):i-1})$$

Cons:

$n \downarrow \rightarrow$ capture long-range dependencies

$n \uparrow \rightarrow$ statistically infeasible to get
good estimates of the probabilities

local dependencies

Neural Language Models:

$$P(X_i | X_{i-(n-1):i-1})$$

RNNs, LSTMs:

Conditional distribution of a token X_i

to depend on the entire context $X_{1:i-1}$

Transformer:

GPT-3 : 2048

fixed context length n

increase in size

in-context learning

✓ Reliability

✓ Social bias

✓ Disinformation

✓ Security

✓ Legal consideration

Foundation Models: Key Points

(1) large-scale pre-trained language models
(PLMs, LLMs)

(2) can be instructed by prompts
to perform specific tasks