

1 Model Framework

Empirical analysis of *natural language* reveals an interesting relation between the distance of characters in a text and their mutual information.

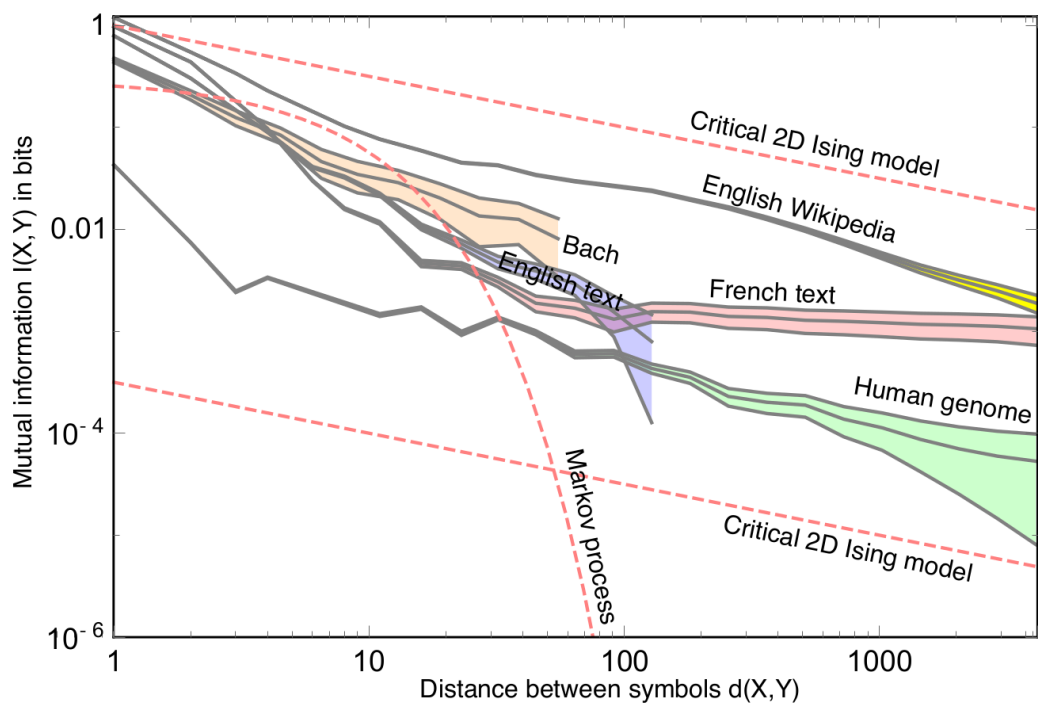


Figure 1.1: Decay of mutual information with separation. Source: [?]

The data show that mutual information tends to decay with distance, but it does so slowly. To be precise, *mutual information in natural language seems to follow a power-law*.

Thus, models of natural language should show a similar behavior. In order to filter potential models for natural language by their ability for *power-law behavior*, we need a precise definition.

1.1 The Framework

The tokens, represented by random variables X_t , are elements of a finite alphabet Σ . For every t and every separation $\tau > 0$, we want to bound

$$I(X_t; X_{t+\tau}) \in \Omega(\tau^{-\alpha}), \quad I(X_t; X_{t+\tau}) \in \mathcal{O}(\tau^{-\beta}) \quad ,$$

for some fixed $\alpha, \beta \in \mathbb{R}_{>0}$.

The first condition ensures that the mutual information does not decay too quickly, while the latter ensures that $I(X_t; X_{t+\tau}) \xrightarrow{\tau \rightarrow \infty} 0$, just like the data show. We also may replace the latter condition by this implication.

The importance of the second condition is further emphasized when considering models like the following time-homogenous Markov chain consisting of the two states A and B : The probability to transition from A to A is one, and hence the transition probability from A to B is zero. Similarly, starting at state B , we can only remain at state B .

The pairwise mutual information of this simple model equals a constant. Thus, we can lower bound it by a power-law. However, the problem is that the mutual information does not *decay with distance*. Thus, we conclude that high mutual information alone is not a good indicator for model quality.

1.1.1 Model Definition

In sequence modeling, a model is typically built from a finite set of rules or parameters defined over an alphabet Σ . These rules allow us to assign a probability to any finite string. For instance, a time-homogeneous Markov chain uses a fixed transition matrix to define a probability measure over the set of all strings of a given length n (i.e., over Σ^n) for any $n \in \mathbb{N}$.

This central idea of a family of finite domains leads to our general definition:

Definition 1.1.1 (Model over Σ^*). Let Σ be a finite alphabet. A model S over Σ^* is a function $S : \Sigma^* \mapsto [0, 1]$ that assigns a probability to each finite string $w \in \Sigma^*$, subject to the constraint that for any length $n \in \mathbb{N}$, the probabilities of all strings of that length sum to one:

$$\sum_{w \in \Sigma^n} S(w) = 1 \quad .$$

Furthermore, we denote the restriction of S to strings of length n as $S_n := S|_{\Sigma^n}$. The function $S_n : \Sigma^n \rightarrow [0, 1]$ is thus a probability measure over Σ^n .

Another strength of this definition is that we don't constrain models by their parameterization. How the models are defined, and how they compute the probabilities is up to them.

Additionally, we might want to restrain S in order to have reasonable time and space complexity, and to ensure the model is *consistent*, which means that the language of S_n should look *similar* to S_{n+d} , whatever this might mean. We also write w_i for X_i . We can think of w as a 1-indexed String of random variables.

We present one strict definition for this *similarity* in the following definition:

Definition 1.1.2. We say S has the *bulk marginal property* iff for every $n \in \mathbb{N}$, $w \in \Sigma^n$ it holds true that

$$\sum_{c \in \Sigma} S_{n+1}(wc) = S_n(w) \quad .$$

Lemma 1.1.1. For every $d \in \mathbb{N}$, let $I := [n+d] \setminus [n] = \{n+1, \dots, n+d\}$. Then, if S has the bulk marginal property, we have for every $w \in \Sigma^n$:

$$\sum_{s \in \Sigma^d} S_{n+d}(ws) = S_n(w) \quad .$$

Proof. We use induction over d . The base case directly follows from the definition of the bulk marginal property. Thus, assume the claim holds for some $d := k$. Then we have

$$\begin{aligned} \sum_{s \in \Sigma^{k+1}} S_{n+k+1}(ws) &= \sum_{v \in \Sigma^k} \sum_{c \in \Sigma} S_{n+k+1}(wvc) \\ &= \sum_{v \in \Sigma^k} \left(\sum_{c \in \Sigma} S_{n+k+1}((wv)c) \right) \\ &\stackrel{\text{bulk marginal property}}{=} \sum_{v \in \Sigma^k} S_{n+k}(wv) \\ &\stackrel{\text{inductive hypothesis}}{=} S_n(w) \quad , \end{aligned}$$

which concludes the induction. □

Definition 1.1.3 (Induced Bulk Marginal Model). Based on the model S , we can construct an *induced bulk marginal* model S^* by defining S_n^* recursively as

- $S_1^* := S_1$,
- $S_{n+1}^*(wc) := S_n^*(w) \frac{S_{n+1}(wc)}{\sum_{c' \in \Sigma} S_{n+1}(wc')} \quad ,$

where $w \in \Sigma^n, c \in \Sigma$.

Remark 1.1.1. If $\sum_{c' \in \Sigma} S_{n+1}(wc') = 0$, we might set $S_{n+1}^*(wc) := S_n^*(w) \frac{1}{|\Sigma|}$.

Lemma 1.1.2. *The induced bulk marginal model S^* indeed has the bulk marginal property.*

Proof. We have:

$$\begin{aligned}
\sum_{c \in \Sigma} S_{n+1}^*(wc) &\stackrel{\text{def of } S^*}{=} \sum_{c \in \Sigma} S_n^*(w) \frac{S_{n+1}(wc)}{\sum_{c' \in \Sigma} S_{n+1}(wc')} \\
&= S_n^*(w) \cdot \frac{1}{\sum_{c' \in \Sigma} S_{n+1}(wc')} \sum_{c \in \Sigma} S_{n+1}(wc) \\
&= S_n^*(w) \cdot \frac{\sum_{c \in \Sigma} S_{n+1}(wc)}{\sum_{c' \in \Sigma} S_{n+1}(wc')} \\
&\stackrel{\check{}}{=} S_n^*(w) \quad .
\end{aligned}$$

□

1.1.2 Restricting the Model

We can restrict our general model S by specifying that for each length n , the probability distribution S_n is computed by a parameterized **inference function**, $f_{n, \theta_n} : \Sigma^n \rightarrow [0, 1]$. Each function is identified by a parameter vector θ_n from a corresponding **parameter space** Θ_n .

For the framework to be valid, each function must define a proper probability distribution. We denote this distribution by S_{n, θ_n} and require that:

$$S_{n, \theta_n}(w) := f_{n, \theta_n}(w) \quad \text{and} \quad \sum_{w \in \Sigma^n} f_{n, \theta_n}(w) = 1 \quad .$$

This approach defines a **model class** \mathcal{S}_n , which is the set of all distributions that can be generated by the family of inference functions with parameters in Θ_n :

$$\mathcal{S}_n := \{S_{n, \theta_n} \mid \theta_n \in \Theta_n\} \quad .$$

A complete model $S \equiv (S_n)_{n \in \mathbb{N}}$ is thus specified by an inference function and a corresponding sequence of chosen parameters $(\theta_n)_{n \in \mathbb{N}}$.

Remark 1.1.2. The distinction between the inference function f and the distribution S is crucial. Two different functions, $f_{n, \theta}$ and $f_{n, \theta'}$, might have vastly different time and space complexities even if they compute the exact same distribution (i.e., $S_{n, \theta} = S_{n, \theta'}$). The complexity is therefore a property of the specific algorithmic implementation and parameterization of f_{n, θ_n} .

1.1.3 Power-Law Behavior

Now, we formalize what it means for a model to exhibit power-law behavior. Specifically, we require power-law decay in mutual information with respect to τ between *any* two variables X_t and $X_{t+\tau}$. This condition must hold for all t and for samples from *any* distribution S_n where $t + \tau \leq n$.

Definition 1.1.4. We define $i_{S_n}(\tau)$ and $I_{S_n}(\tau)$ to be the minimal and maximal mutual information between any two variables of S_n with distance τ . Formally, let $X_t, X_{t+\tau}$ ($t + \tau \leq n$) be random variables sampled from S_n . Then:

- $i_{S_n}(\tau) := \min_{t \in [n-\tau]} I(X_t; X_{t+\tau})$,
- $I_{S_n}(\tau) := \max_{t \in [n-\tau]} I(X_t; X_{t+\tau})$.

Definition 1.1.5 (Strong Power-Law Behavior). A model S has *strong lower bound power-law behavior* iff there exist constants $c_\alpha, \alpha \in \mathbb{R}_{>0}$ s.t. for every $n \in \mathbb{N}$ it holds true that $i_{S_n}(\tau) \geq c_\alpha \tau^{-\alpha}$. Similarly, S has *upper bound power-law behavior* iff there exist constants $c_\beta, \beta \in \mathbb{R}_{>0}$ s.t. for every $n \in \mathbb{N}$ it holds true that $I_{S_n}(\tau) \leq c_\beta \tau^{-\beta}$. Furthermore, S has *decaying behavior* iff for every $n \in \mathbb{N}$ we have $I_{S_{n+\tau}}(\tau) \xrightarrow{\tau \rightarrow \infty} 0$. Lastly, S has *strong power-law behavior* iff it has strong lower bound and upper bound power-law behavior (alternatively decaying behavior instead of upper bound power-law behavior).

Remark 1.1.3. For a model S^* with the bulk marginal property we can replace "for every $n \in \mathbb{N}$ " in definition 0.5 with "for $n \rightarrow \infty$ " thanks to lemma 0.1.

Proposition 1.1.1. *Upper bound power-law behavior implies decaying behavior.*

Proof. Assume model S has upper bound power-law behavior. Then there exist constants $c_\beta, \beta \in \mathbb{R}_{>0}$ s.t. for every $n \in \mathbb{N}$ it holds true that $I_{S_n}(\tau) \leq c_\beta \tau^{-\beta}$, especially for $n := n' + \tau$. Thus, for every $n' \in \mathbb{N}$:

$$I_{S_{n'+\tau}}(\tau) \leq c_\beta \tau^{-\beta} \xrightarrow{\tau \rightarrow \infty} 0 \quad .$$

□

Definition 1.1.6. We define $\overline{i_{S_n}}$ to be the minimal mutual information between any two variables over S_n with arbitrary distance τ . Formally, let X_i, X_j ($1 \leq i < j \leq n$) be random variables with distributions defined by S_n . Then:

$$\overline{i_{S_n}} := \min_{(i,j) \in [n]^2, i < j} I(X_i; X_j) = \min_{\tau \in [n-1]} i_{S_n}(\tau) \quad .$$

Definition 1.1.7 (Weak Power-Law Behavior). A model S has *weak lower bound power-law behavior* iff $\overline{i_{S_n}} \in \Omega(n^{-\alpha})$ for some $\alpha \in \mathbb{R}_{>0}$. Additionally, S has *weak power-law behavior* iff it has weak lower bound and upper bound power-law behavior (alternatively decaying behavior instead of upper bound power-law behavior).

Theorem 1.1.1 (Power-Law Decay for All Tokens in Models with the Bulk Marginal Property). *Let S be a model that satisfies the bulk marginal property and exhibits weak lower bound power-law behavior. Then, there exists an $\alpha \in \mathbb{R}_{>0}$ s.t. for every X_t , $I(X_t; X_{t+\tau}) \in \Omega(\tau^{-\alpha})$ (where X_t and $X_{t+\tau}$ are sampled over $S_{t+\tau}$, or, equivalently, any $S_{t+\tau+k}$).*

Proof. Since S has weak lower bound power-law behavior, there exist $\alpha', c' \in \mathbb{R}_{>0}$ s.t. $\overline{i_{S_n}} \geq c'n^{-\alpha'}$. Then, for every $t \in \mathbb{N}$, we have for $n := t + \tau$ by the definition of $\overline{i_{S_n}}$:

$$\begin{aligned} I(X_t; X_{t+\tau}) &\geq \overline{i_{S_{t+\tau}}} \\ &\geq c'(t + \tau)^{-\alpha'} \\ &= c'\tau^{-\alpha'}\left(\frac{t}{\tau} + 1\right)^{-\alpha'} \\ &\geq c'\tau^{-\alpha'}(t + 1)^{-\alpha'} \quad . \end{aligned}$$

Since S has the bulk marginal property, this inequality holds when sampling over any $S_{t+\tau+k}$, $k \in \mathbb{N}$. Now, set $\alpha := \alpha'$ and $c := c'(t + 1)^{-\alpha'}$. Note that α does not depend on t . Finally, we see that $I(X_t; X_{t+\tau}) \geq c\tau^{-\alpha}$. Thus, we get $I(X_t; X_{t+\tau}) \in \Omega(\tau^{-\alpha})$. \square

In other words, mutual information decays polynomially with distance for all tokens, and the constant α is always the same. This sounds like strong power-law behavior, but the issue is that the scalar c depends on t and we cannot lower bound it by constant greater than zero. Hence, the *starting threshold* $I(X_t; X_{t+1})$ can decay to zero for $t \rightarrow \infty$.

Remark 1.1.4. If additionally S had decaying behavior, then of course we would also have $I(X_t; X_{t+\tau}) \xrightarrow{\tau \rightarrow \infty} 0$.

Remark 1.1.5. The importance of this implication might depend on the context. However, this theorem proves to be very useful when considering its contraposition. In fact, we will use this contraposition later to disprove weak power-law behavior of hidden Markov models (and hence also strong power-law behavior).

Remark 1.1.6. It is crucial for S to have the bulk marginal property in theorem 0.1, or else $I(X_t; X_{t+\tau})$ might depend on S_n , and we cannot exclude $I(X_t; X_{t+\tau}) \xrightarrow{n \rightarrow \infty} 0$.

Proposition 1.1.2. *Strong lower bound power-law behavior implies weak lower bound power-law behavior.*

Proof. Assume model S has strong lower bound power-law behavior. Thus, it follows that there exist $c_\alpha, \alpha \in \mathbb{R}_{>0}$ s.t. for all $n \in \mathbb{N}$ we have that $i_{S_n}(\tau) \geq c_\alpha \tau^{-\alpha}$. Hence:

$$\begin{aligned}
\overline{i_{S_n}} &= \min_{\tau \in [n-1]} i_{S_n}(\tau) \\
&\geq \min_{\tau \in [n-1]} c_\alpha \tau^{-\alpha} \\
&\geq c_\alpha (n-1)^{-\alpha} \\
&= c_\alpha n^{-\alpha} \left(1 - \frac{1}{n}\right)^{-\alpha} \\
&\geq c_\alpha n^{-\alpha} 1^{-\alpha} \\
&= c_\alpha n^{-\alpha} .
\end{aligned}$$

It follows that $\overline{i_{S_n}} \in \Omega(n^{-\alpha})$, and hence S has weak lower bound power-law behavior. \square

Remark 1.1.7. Weak lower bound power-law behavior does *not* imply strong lower bound power-law behavior, not even for models with the bulk marginal property. To see this, note that we might have $i_{S_n}(1) \xrightarrow{n \rightarrow \infty} 0$ for some models with weak lower bound power-law behavior. (S_n may force $i_{S_n}(1)$ to decay to 0 for $n \rightarrow \infty$ because of weak correlations of consecutive tokens very late in the sequence.) The proof of theorem 0.1 fails when defining c , as it depends on t .

Remark 1.1.8. If S has decaying behavior, we cannot prove that S has strong lower bound power-law behavior by bounding $\overline{i_{S_{t+\tau}}}$ (using $\overline{i_{S_{t+\tau}}} \leq i_{S_{t+\tau}}(\tau)$), as we have for every $\tau \in \mathbb{N}$:

$$0 \leq \overline{i_{S_{t+\tau}}} \leq I_{S_{t+\tau}}(t) \xrightarrow{t \rightarrow \infty} 0 \quad .$$