

EEE 551 Information Theory (Spring 2022)

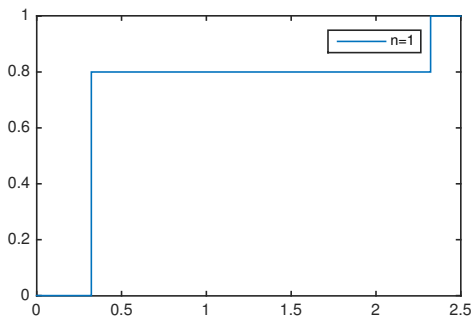
Chapter 3: Asymptotic Equipartition Property (AEP)

- Consider an **independent and identically distributed** (i.i.d.) sequence of random variables

$$(X_1, X_2, \dots, X_{n-1}, X_n) = X^n$$

also called a **discrete memoryless source** (DMS)

- How does $p(X^n)$ behave for large n ?
- For example, let $X_i \sim \text{Bern}(0.2)$
- Plot the CDF of $-\frac{1}{n} \log p(X^n)$:

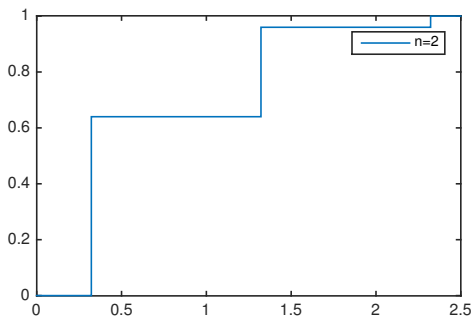


- Consider an **independent and identically distributed** (i.i.d.) sequence of random variables

$$(X_1, X_2, \dots, X_{n-1}, X_n) = X^n$$

also called a **discrete memoryless source** (DMS)

- How does $p(X^n)$ behave for large n ?
- For example, let $X_i \sim \text{Bern}(0.2)$
- Plot the CDF of $-\frac{1}{n} \log p(X^n)$:

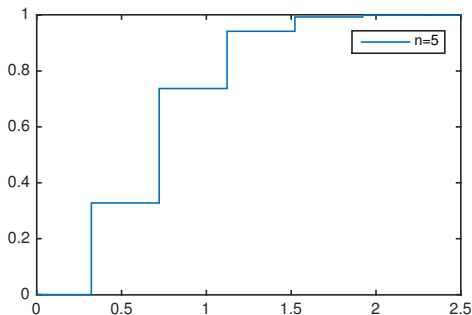


- Consider an **independent and identically distributed** (i.i.d.) sequence of random variables

$$(X_1, X_2, \dots, X_{n-1}, X_n) = X^n$$

also called a **discrete memoryless source** (DMS)

- How does $p(X^n)$ behave for large n ?
- For example, let $X_i \sim \text{Bern}(0.2)$
- Plot the CDF of $-\frac{1}{n} \log p(X^n)$:

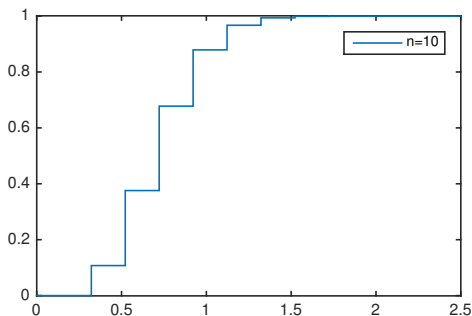


- Consider an **independent and identically distributed** (i.i.d.) sequence of random variables

$$(X_1, X_2, \dots, X_{n-1}, X_n) = X^n$$

also called a **discrete memoryless source** (DMS)

- How does $p(X^n)$ behave for large n ?
- For example, let $X_i \sim \text{Bern}(0.2)$
- Plot the CDF of $-\frac{1}{n} \log p(X^n)$:

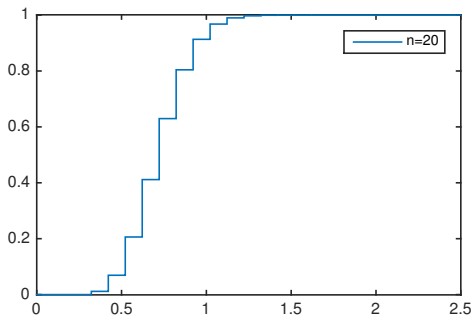


- Consider an **independent and identically distributed** (i.i.d.) sequence of random variables

$$(X_1, X_2, \dots, X_{n-1}, X_n) = X^n$$

also called a **discrete memoryless source** (DMS)

- How does $p(X^n)$ behave for large n ?
- For example, let $X_i \sim \text{Bern}(0.2)$
- Plot the CDF of $-\frac{1}{n} \log p(X^n)$:

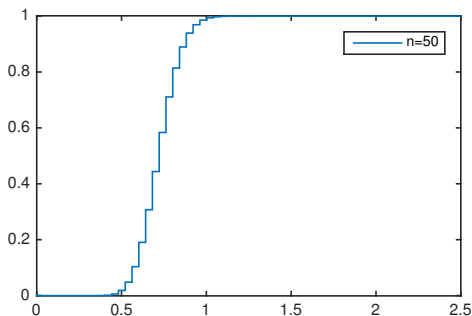


- Consider an **independent and identically distributed** (i.i.d.) sequence of random variables

$$(X_1, X_2, \dots, X_{n-1}, X_n) = X^n$$

also called a **discrete memoryless source** (DMS)

- How does $p(X^n)$ behave for large n ?
- For example, let $X_i \sim \text{Bern}(0.2)$
- Plot the CDF of $-\frac{1}{n} \log p(X^n)$:

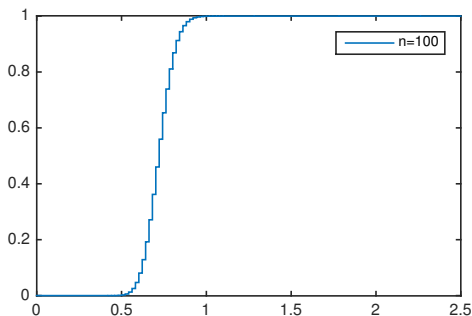


- Consider an **independent and identically distributed** (i.i.d.) sequence of random variables

$$(X_1, X_2, \dots, X_{n-1}, X_n) = X^n$$

also called a **discrete memoryless source** (DMS)

- How does $p(X^n)$ behave for large n ?
- For example, let $X_i \sim \text{Bern}(0.2)$
- Plot the CDF of $-\frac{1}{n} \log p(X^n)$:

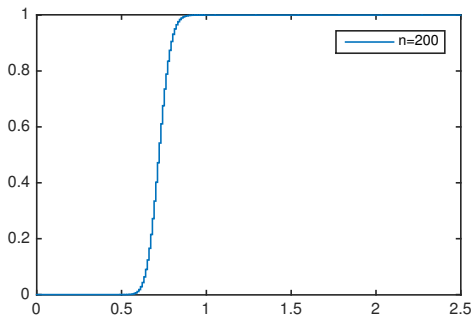


- Consider an **independent and identically distributed** (i.i.d.) sequence of random variables

$$(X_1, X_2, \dots, X_{n-1}, X_n) = X^n$$

also called a **discrete memoryless source** (DMS)

- How does $p(X^n)$ behave for large n ?
- For example, let $X_i \sim \text{Bern}(0.2)$
- Plot the CDF of $-\frac{1}{n} \log p(X^n)$:

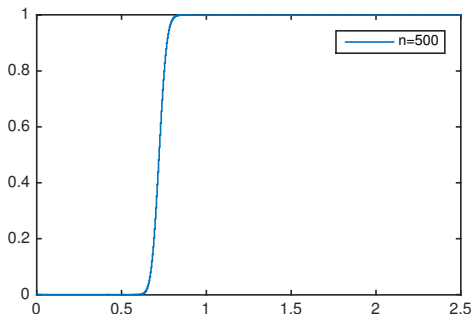


- Consider an **independent and identically distributed** (i.i.d.) sequence of random variables

$$(X_1, X_2, \dots, X_{n-1}, X_n) = X^n$$

also called a **discrete memoryless source** (DMS)

- How does $p(X^n)$ behave for large n ?
- For example, let $X_i \sim \text{Bern}(0.2)$
- Plot the CDF of $-\frac{1}{n} \log p(X^n)$:

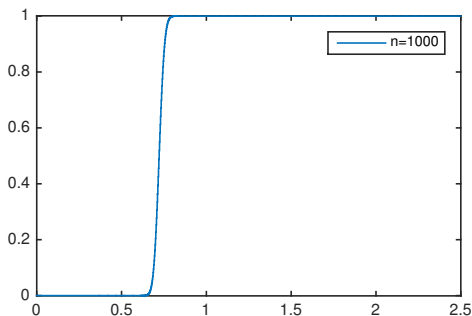


- Consider an **independent and identically distributed** (i.i.d.) sequence of random variables

$$(X_1, X_2, \dots, X_{n-1}, X_n) = X^n$$

also called a **discrete memoryless source** (DMS)

- How does $p(X^n)$ behave for large n ?
- For example, let $X_i \sim \text{Bern}(0.2)$
- Plot the CDF of $-\frac{1}{n} \log p(X^n)$:

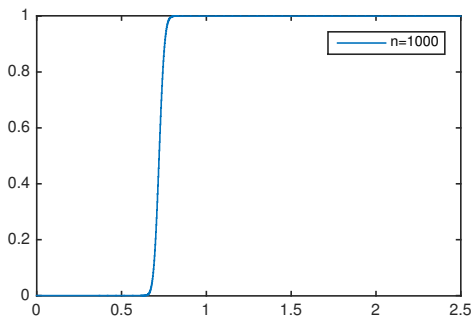


- Consider an **independent and identically distributed** (i.i.d.) sequence of random variables

$$(X_1, X_2, \dots, X_{n-1}, X_n) = X^n$$

also called a **discrete memoryless source** (DMS)

- How does $p(X^n)$ behave for large n ?
- For example, let $X_i \sim \text{Bern}(0.2)$
- Plot the CDF of $-\frac{1}{n} \log p(X^n)$:



- Concentrates around $H(X) = 0.7219$!

- Let $X^n \stackrel{\text{iid}}{\sim} \text{Bern}(p)$
- Given $x^n = (x_1, x_2, \dots, x_n)$,

$$p(x^n) = \prod_{i=1}^n p(x_i) = \left(\prod_{i:x_i=1} p \right) \left(\prod_{i:x_i=0} (1-p) \right) = p^{\sum_i x_i} (1-p)^{\sum_i (1-x_i)}$$

- By the law of large numbers, $\sum_i X_i \approx np$ and $\sum_i (1 - X_i) \approx n(1 - p)$

$$\begin{aligned} p(X^n) &\approx p^{np} (1-p)^{n(1-p)} \\ &= \left(2^{\log p} \right)^{np} \left(2^{\log(1-p)} \right)^{n(1-p)} \\ &= 2^{np \log p} 2^{n(1-p) \log(1-p)} \\ &= 2^{n(p \log p + (1-p) \log(1-p))} \\ &= 2^{-nH(X)} \end{aligned}$$

$$\text{i.e. } -\frac{1}{n} \log p(X^n) \approx H(X)$$

- “Almost all events are almost equally likely”

Weak Law of Large Numbers

A sequence of random variables Z_1, Z_2, Z_3, \dots converges to c **in probability** if, for all $\epsilon > 0$,

$$\lim_{n \rightarrow \infty} \Pr \{ |Z_n - c| > \epsilon \} = 0.$$

Theorem (Weak Law of Large Numbers)

Let Y_1, Y_2, \dots, Y_n be i.i.d. random variables with expectation μ , and finite variance.

Then $\frac{1}{n} \sum_{i=1}^n Y_i \rightarrow \mu$ in probability.

The proof is based on two **concentration inequalities**:

1. Markov's inequality: If $P(X < 0) = 0$, then for any $a > 0$, $P(X > a) \leq \frac{\mathbb{E}(X)}{a}$

Proof:

- Let $I(x > a) = \begin{cases} 1 & x > a \\ 0 & x \leq a \end{cases}$
- $I(x > a) \leq \frac{x}{a}$ for all $x \geq 0$
- $P(X > a) = \mathbb{E}(I(X > a)) \leq \mathbb{E}\left(\frac{X}{a}\right) = \frac{\mathbb{E}(X)}{a}$

Weak Law of Large Numbers, ctd

2. Chebyshev's inequality: For any random variable X , for any $a > 0$,

$$P(|X - \mathbb{E}(X)| > a) \leq \frac{\text{Var}(X)}{a^2}.$$

Proof: Apply Markov's inequality to $Y = (X - \mathbb{E}(X))^2$:

$$\begin{aligned} P(|X - \mathbb{E}(X)| > a) &= P((X - \mathbb{E}(X))^2 > a^2) = P(Y > a^2) \\ &\leq \frac{\mathbb{E}(Y)}{a^2} = \frac{\text{Var}(X)}{a^2} \end{aligned}$$

Proof of weak law of large numbers:

- Let $Z_n = \frac{1}{n} \sum_{i=1}^n Y_i$
- $\mathbb{E}(Z_n) = \frac{1}{n} \sum_{i=1}^n \mathbb{E}(Y_i) = \mu$
- $\text{Var}(Z_n) = \frac{1}{n^2} \sum_{i=1}^n \text{Var}(Y_i) = \frac{\text{Var}(Y_1)}{n}$
- Apply Chebyshev's inequality:

$$P(|Z_n - \mu| > \epsilon) \leq \frac{\text{Var}(Z_n)}{\epsilon^2} = \frac{\text{Var}(Y_1)}{\epsilon^2 n} \rightarrow 0$$

Thus $Z_n \rightarrow \mu$ in probability

AEP Theorem

If X_1, X_2, \dots, X_n are drawn i.i.d. from $p(x)$, then

$$-\frac{1}{n} \log p(X^n) \rightarrow H(X) \text{ in probability}$$

Proof:

$$\begin{aligned} -\frac{1}{n} \log p(X^n) &= -\frac{1}{n} \log \prod_{i=1}^n p(X_i) \\ &= -\frac{1}{n} \sum_{i=1}^n \log p(X_i) \\ &\rightarrow -\mathbb{E} \log p(X) && \text{by WLLN} \\ &= H(X) \end{aligned}$$

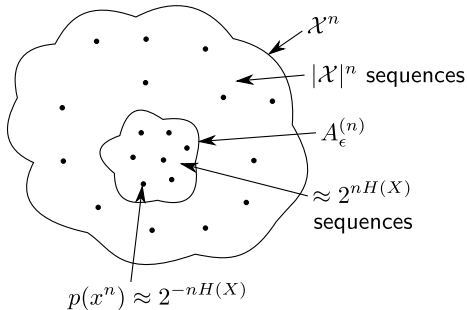
The Typical Set

Given a distribution $p(x)$, the **typical set** $A_\epsilon^{(n)}$ is the set of $(x_1, x_2, \dots, x_n) = x^n \in \mathcal{X}^n$ such that

$$2^{-n(H(X)+\epsilon)} \leq p(x^n) \leq 2^{-n(H(X)-\epsilon)}$$

Typical Set Properties

- $x^n \in A_\epsilon^{(n)}$ iff $\left| -\frac{1}{n} \log p(x^n) - H(X) \right| \leq \epsilon$
- $\lim_{n \rightarrow \infty} \Pr\{A_\epsilon^{(n)}\} = 1$
- $|A_\epsilon^{(n)}| \leq 2^{n(H(X)+\epsilon)}$
- For n sufficiently large,
 $|A_\epsilon^{(n)}| \geq (1 - \epsilon)2^{n(H(X)-\epsilon)}$



Proofs of Typical Set Properties

$$\blacksquare x^n \in A_\epsilon^{(n)} \text{ iff } \left| -\frac{1}{n} \log p(x^n) - H(X) \right| \leq \epsilon$$

Proof: From the definition of the typical set:

$$\begin{aligned} 2^{-n(H(X)+\epsilon)} &\leq p(x^n) \leq 2^{-n(H(X)-\epsilon)} \\ -n(H(X)+\epsilon) &\leq \log p(x^n) \leq -n(H(X)-\epsilon) \\ H(X)+\epsilon &\geq -\frac{1}{n} \log p(x^n) \geq H(X)-\epsilon \end{aligned}$$

$$\blacksquare \lim_{n \rightarrow \infty} \Pr\{A_\epsilon^{(n)}\} = 1$$

Proof: From the AEP theorem,

$$\Pr \left\{ \left| -\frac{1}{n} \log p(X^n) - H(X) \right| > \epsilon \right\} \rightarrow 0.$$

Thus $\Pr\{A_\epsilon^{(n)}\} \rightarrow 1$ as $n \rightarrow \infty$

- $|A_\epsilon^{(n)}| \leq 2^{n(H(X)+\epsilon)}$

Proof:

$$\begin{aligned}
 1 &\geq \Pr\{A_\epsilon^{(n)}\} \\
 &= \sum_{x^n \in A_\epsilon^{(n)}} p(x^n) \\
 &\geq \sum_{x^n \in A_\epsilon^{(n)}} 2^{-n(H(X)+\epsilon)} \\
 &= |A_\epsilon^{(n)}| 2^{-n(H(X)+\epsilon)}
 \end{aligned}$$

- For n sufficiently large, $|A_\epsilon^{(n)}| \geq (1 - \epsilon)2^{n(H(X)-\epsilon)}$

Proof: Since $\Pr\{A_\epsilon^{(n)}\} \rightarrow 1$, for sufficiently large n

$$\begin{aligned}
 1 - \epsilon &\leq \Pr\{A_\epsilon^{(n)}\} \\
 &= \sum_{x^n \in A_\epsilon^{(n)}} p(x^n) \\
 &\leq \sum_{x^n \in A_\epsilon^{(n)}} 2^{-n(H(X)-\epsilon)} \\
 &= |A_\epsilon^{(n)}| 2^{-n(H(X)-\epsilon)}
 \end{aligned}$$