# EE/Ma/CS 126a: Information Theory

Aman Bhargava

October-December 2022

# Contents

## 0.1   Introduction and Course Information

This document offers an overview of EE/Ma/CS 126a at Caltech. They comprise my condensed course notes for the course. No promises are made relating to the correctness or completeness of the course notes. These notes are meant to highlight difficult concepts and explain them simply, not to comprehensively review the entire course.

**Course Information**

- Professor: Michelle Effros

- Term: 2022 Fall

# Chapter 1

# Math Review

## 1.1  Combinatorics & Probability

**Binomial Distribution & Coefficient**

- **Bernoulli Process**: Repeated trials, each with one binary outcome. The probability of a positive outcome is $p \in [0, 1]$. Each trial is independent.

- **Binomial Distribution**: Let $x$ represent the number of successful trials in a Bernoulli process repeated $n$ times with success probability $p$. The binomial distribution gives the probability distribution on $x$:

$$b(x; n, p) = \binom{n}{k} p^x (1-p)^{n-x} \tag{1.1}$$

  Which has $\mu = np$, $\sigma^2 = npq$.

- **Intuition for Binomial Distribution**: The probability of observing a sequence with $x$ positive outcomes and $n - x$ negative outcomes is $p^x(1-p)^{n-x}$. There are $\binom{n}{k}$ different sequences (i.e., permutations) that have $x$ positive cases and $n$ negative cases. Thus the total probability of observing $x$ positive cases is given by Eq 1.1.

- **Binomial Coefficient**:

$$\binom{n}{k} = \frac{n!}{k!(n-k)!} \tag{1.2}$$

## 1.2 Logarithm Identities

Entropy calculations and manipulations involve a lot of logarithms. They're not so bad once you get to know them, though:

- **Definition**:
$$a = b^{\log_b a}$$

- **Sum-Product**:
$$\log_c(ab) = \log_c a + \log_c b$$

- **Difference-Quotient**:
$$\log(a/c) = \log a - \log c$$

$$\log \frac{1}{a} = -\log a$$

- **Product-Exponent**:
$$\log_c(a^n) = n \log_c(a)$$

- **Swapping Base**:
$$\log_b(a) = \log_a(b)$$

- **Swapping Exponential**:
$$a^{\log n} = n^{\log a}$$

- **Change of Base**;
$$\log_b(a) = \frac{\log_x(a)}{\log_x(b)}$$

# Chapter 2

# Entropy Definitions

*Chapter 2 of Elements of Information Theory.*

## 2.1 Entropy, Conditional Entropy, Joint Entropy

**Entropy Definition (Discrete):**

$$H(X) = \sum_{x \in \mathcal{X}} p(x) \log(\frac{1}{p(x)}) \qquad (2.1)$$

$$= -\sum_{x \in \mathcal{X}} p(x) \log p(x) \qquad (2.2)$$

$$= \mathbb{E}[\log \frac{1}{p(x)}] \qquad (2.3)$$

**Theorem 1** *Properties of Entropy*

1. ***Non-negativity:*** *$H(X) \geq 0$ – Reasoning: Entropy is the sum-product of non-negative terms.*

2. ***Change of base:*** *$H_b(X) = (\log_b a)H_a(X)$*

3. ***Bernoulli entropy:*** *$H(X) = -p \log p - q \log q \equiv H(p)$.*

    - *$H(p)$ is a concave function of $p$, peaks at $p = q = 0.5$.*

**Joint Entropy:**  Literally just entropy of vector $[X, Y]^\top$.

$$H(X, Y) = -\sum_{x \in \mathcal{X}} \sum_{y \in \mathcal{Y}} p(x, y) \log p(x, y) \qquad (2.4)$$

$$= \mathbb{E}[\log p(x, y)] \qquad (2.5)$$

$$\qquad (2.6)$$

**Conditional Entropy:**  $H(Y|X)$ is the expected entropy of $p(y|x)$ averaged across all $x$.

$$H(Y|X) = \sum_{x \in \mathcal{X}} p(x) \underbrace{\sum_{y \in \mathcal{Y}} p(y|x) \log p(y|x)}_{H(Y|X=x)} \qquad (2.7)$$

$$= -\mathbb{E}[\log p(Y|X)] \qquad (2.8)$$

Entropy can be thought of as the **uncertainty** in the value of a random variable. High entropy corresponds to a high degree of uncertainty. Conditional entropy $H(Y|X)$ can be thought of as the average **remaining uncertainty** in the value of $Y$ after learning the value of $X$.

**Theorem 2** *Chain Rule for Entropy*

$$H(X, Y) = H(X) + H(Y|X) \qquad (2.9)$$

$$= H(Y) + H(X|Y) \qquad (2.10)$$

*It also follows that*

$$H(X, Y|Z) = H(X|Z) + H(Y|X, Z) \qquad (2.11)$$

***Proof sketch:***

- *Recall that $H(X) = -\mathbb{E}[\log p(x)]$ and $H(Y|X) = -\mathbb{E}[\log p(y|x)]$.*

- *$\log p(x) + \log p(y|x) = \log(p(x) \cdot p(y|x)) = \log p(x, y)$.*

- *The proof follows from there. You can also write out the full sum form of $H(X, Y)$ and recover $H(X), H(Y|X)$ from there if you're feeling rigorous.*

## 2.2 Relative Entropy & Mutual Information

**Relative Entropy:** $D(p\|q)$ gives a *distance* between distributions $p(x)$ and $q(x)$. Also known as KL divergence.

$$D(p\|q) = \sum_{x \in \mathcal{X}} p(x) \log \frac{p(x)}{q(x)} \tag{2.12}$$

$$= \mathbb{E}_{p(x)}[\log \frac{p(x)}{q(x)}] \tag{2.13}$$

$$\tag{2.14}$$

This also corresponds to the **inefficiency** of using $q$ as a replacement for $p$ when generating codes for tokens drawn from $p(x)$.

- **Average code length with correct $p(x)$:** $H(p)$.

- **Average code length with incorrect $q(x)$:** $H(p) + D(p\|q)$.

**Theorem 3** *Properties of Relative Entropy*

1. ***Asymmetric:*** *In general, $D(p\|q) \neq D(q\|p)$.*

2. ***Non-negative:*** $D(p\|q) \geq 0$.

3. ***Identity:*** *If $D(p\|q) = 0$ then $p \equiv q$.*

**Conditional Relative Entropy/KL Divergence:** Distance between two distributions when conditioned on the same variable. Similar idea of averaging across all values of the conditioning variable.

$$D(p(y|x)\|q(y|x)) = \sum_{x \in \mathcal{X}} p(x) \Big[ \sum_{y \in \mathcal{Y}} p(y|x) \log \frac{p(y|x)}{q(y|x)} \Big] \tag{2.15}$$

$$= \mathbb{E}_{p(x,y)} \log \Big[ \frac{p(y|x)}{q(y|x)} \Big] \tag{2.16}$$

We now move onto **mutual information** – a measure of the dependence of two variables. As we will see, it is the **reduction in uncertainty** of $X$ due to knowing $Y$, on average.

**Mutual Information:**

$$I(X;Y) = \sum_{x\in\mathcal{X}}\sum_{y\in\mathcal{Y}} p(x,y) \log\left(\frac{p(x,y)}{p(x)p(y)}\right) \tag{2.17}$$

$$= D\big(p(x,y)\|p(x)p(y)\big) \tag{2.18}$$

$$= \mathbb{E}[\log \frac{p(X,Y)}{p(X)p(Y)}] \tag{2.19}$$

**Theorem 4** *Properties of Mutual Information*

- *It is the **divergence** between $p(x,y)$ and $p(x)p(y)$.*

- ***Symmetry:*** $I(X;Y) = I(Y;X)$.

- ***Relation to Entropy:*** *Mutual information is the reduction in uncertainty of each RV expected after discovering the other variable's value.*

$$I(X;Y) = H(X) - H(X|Y) \tag{2.20}$$

$$= H(Y) - H(Y|X) \tag{2.21}$$

$$\tag{2.22}$$

- ***Alternative Entropy Relation:***

$$I(X;Y) = H(X) + H(Y) - H(X,Y) \tag{2.23}$$

$$I(X;X) = H(X) - H(X|X) \tag{2.24}$$

$$= H(X) \tag{2.25}$$

***Proof Sketch for (3):***

- *Within the definition of $I(X;Y)$ there is a term $\log \frac{p(x,y)}{p(x)p(y)}$.*

- *Once you convert the argument of the log into $p(x|y)/p(x)$, you can separate out $H(X) - H(X|Y)$ using the quotient-difference logarithm rule.*

**Conditional Mutual Information:** $I(X,Y|Z)$ is the average reduction in uncertainty on the value of $X$ due to knowing $Y$ when $Z$ is given.

$$I(X;Y|Z) = H(X|Z) - H(X|Y,Z) \tag{2.26}$$

$$= \mathbb{E}_{p(x,y,z)} \log \Big[\frac{p(X,Y|Z)}{p(X|Z)p(Y|Z)}big] \tag{2.27}$$

$$\tag{2.28}$$

## 2.3  Chain Rules: $H(\cdot), I(\cdot; \cdot)$

These chain rules end up being very useful in a lot of proofs. Deeply under-standing them is a good idea.

**Theorem 5** *Entropy chain rule Let $X_1, X_2, \ldots, X_n \tilde{p}(x_1, x_2, \ldots, x_n)$. Then*

$$H(X_1, \ldots, X_n) = \sum_{i=1}^{n} H(X_i | X_{i-1}, X_{i-2}, \ldots, X_1) \qquad (2.29)$$

**Proof:**  *Repeatedly apply Equation 2.9.*

**Intuition of Chain Rule:**    It's important to note that the term in the sum is conditioned on elements $X_j$ with $j < i$. Conditioning always reduces entropy, so it's as though the "additional entropy" from the term must be reduced to account for the previous terms already having been added to the total. Also note that any order can suffice – there is no absolute order in the sum.

**Theorem 6** *Chain Rule for Mutual Information*

$$I(X_1, \ldots, X_n; Y) = \sum_{i=1}^{n} I(X_i; Y | X_{i-1}, X_{i-2}, \ldots, X_1) \qquad (2.30)$$

$$(2.31)$$

**Proof:**  *Start with the entropy definition of mutual information. Then apply the chain rule for entropy (Theorem 5).*

- *$I(X_{1:n}; Y) = H(X_{1:n}) - H(X_{1:n} | Y)$.*
- *$= \sum_{i=1}^{n} H(X_i | X_{i-1:1}) - \sum_{i=1}^{n} H(X_i | X_{i-1:1}, Y)$.*
- *$= \sum_{i=1}^{n} I(X_i; Y | X_{1:i-1})$.*

**Theorem 7** *Chain Rule for Relative Entropy*

$$D(p(x, y) \| q(x, y)) = D(p(x) \| q(x)) + D(p(y|x) \| q(y|x)) \qquad (2.32)$$

   **Proof sketch:** *Expand the LHS in log-sum form. Separate the term in the* log *into a sum of two* log *terms corresponding to the two divergences on the RHS.*

## 2.4   Jensen's Inequality & Consequences

**Theorem 8** *Jensen's Inequality (and Convexity) For any convex function f and random variable X,*

$$\mathbb{E}\big[f(x)\big] \geq f(\mathbb{E}[X]) \tag{2.33}$$

*Where "convex f in $(a,b)$"* $\iff$

$$f(\lambda x_1 + (1-\lambda)x_2) \leq \lambda f(x_1) + (1-\lambda)f(x_2) \tag{2.34}$$

*for all $x_1, x_2 \in (a,b)$ and $\lambda \in [0,1]$.* ***Strict*** *convexity has equality iff $\lambda \in \{0,1\}$.* ***Concave*** *f* $\iff$ *convex* $-f$.

   ***Proof sketch of Jensen's Inequality:***

- *Start with $|\mathcal{X}| = 2$. Then we can let a vector $\mathbf{p} \in \mathbb{R}^2$ represent f(x).*

- $\mathbb{E}[f(x)] = p_1 f(x_1) + p_2 f(x_2).$

- $f(\mathbb{E}[X]) = f(p_1 x_1 + p_2 x_2).$

- *Show equivalence of Jensen's inequality* $\iff$ $\mathbf{p}$ *is a valid PMF.*

- *Expand to $|\mathcal{X}| = k-1$ (induction proof).*

- *Use continuity arguments for continuous case.*

**Theorem 9** *Implications of Jensen's Inequality*

1. ***Information Inequality:*** $D(p\|q) \geq 0.$

2. ***MI Inequality:*** $I(X;Y) \geq 0.$

3. ***Maximum Entropy:*** $H(X) \leq \log|\mathcal{X}|.$ *Maximum is achieved by $X \tilde{} uniform(\mathcal{X}).$*

4. ***Information can't Hurt:*** $H(X|Y) \leq H(X).$

5. ***Entropy Sum Bound:*** $H(X_{1:n}) \leq \sum_{i=1}^{n} H(X_i)$ *with equality for independent $X_i$.*

## 2.5   Log-Sum Inequality & Consequences

**Theorem 10** *Log-Sum Inequality Let $\{a_i, b_i\}_{i=1}^n$ be **non-negative** numbers. Then*

$$\sum_{i=1}^n a_i \log \frac{a_i}{b_i} \geq \big( \sum_{i=1}^n a_i \big) \log \frac{\sum a_i}{\sum b_i} \tag{2.35}$$

*With equality iff $\frac{a_i}{b_i}$ is constant.*
   ***Proof sketch:***

1. *Introduce $\alpha_i = a_i / \sum a_j$ and $t_i = a_i/b_i$. $\alpha_i$ values comprise a probability distribution (sum to 1).*

2. *Apply **Jensen's** to $\sum \alpha_i f(t_i) \geq f(\sum \alpha_i t_i)$ where $f() = \log()$.*

**Theorem 11** *Applications of Log-Sum Inequality*

- ***Convexity of Relative Entropy:***  *$D(p\|q)$ is convex in pairs of $p, q$.*

$$D(\lambda p_1 + (1-\lambda)p_1 \| \lambda q_1 + (1-\lambda)q_2) \leq \lambda D(p_1\|q_1) + (1-\lambda)D(p_2\|q_2) \tag{2.36}$$

- ***Concavity of Mutual Information:***

    1. ***For fixed** $p(y|x)$**:** $I(X;Y)$ is concave in $p(x)$.*
    2. ***For fixed** $p(x)$**:** $I(X;Y)$ if concave in $p(y|x)$.*
    3. *This property is really handy when doing channel capacity calculations/bounds!*

## 2.6   Data Processing Inequality

**Theorem 12** *Data Processing Inequality Let $X \to Y \to Z$ form a markov chain. That is, $p(x,y,z) = p(x)p(y|x)p(z|y)$. Then*

$$I(X;Y) \geq I(X;Z) \tag{2.37}$$

*That is, the mutual information between $X, Z$ is bounded by the mutual information between $X, Y$.*
   ***Proof sketch:***

- *Expand $I(X;Y,Z)$ with the chain rule.*

- *($\star$) Observe that $I(X;Z|Y) = 0$ since $X, Z$ are **conditionally independent** given $Y$ (recall Markov theory – head-to-tail connection).*

- *Since $I(X;Y|Z) \geq 0$, we can conclude that $I(X;Y) \geq I(X;Z)$.*

**Sufficient Statistics Connection:** Assume that $X \to Y$ is a Markov chain. E.g., $X$ are some parameters of an underlying distribution and $Y$ are some data collected from the distribution. If $Z = f(Y)$, then we have $X \to Y \to Z$. Therefore any function $Z$ on the data $Y$ cannot have greater information on the underlying distribution $X$ than $Y$ did in the first place.

A **sufficient statistic** is one that has all the information that the data had about the underlying distribution. That is, $Z = f(Y)$ is a sufficient statistic on $Y$ if $I(Z; X) = I(Y; X)$.

**Minimal Sufficient Statistic:** Let $\theta \to T(X) \to U(X) \to X$. $T(X)$ is the *minimal sufficient statistic* iff $U(X)$ can be any sufficient statistic and still have $\theta \to T(X) \to U(X) \to X$ be a valid Markov chain.

## 2.7 Fano's Inequality

Fano's inequality concerns estimators for random variables. Given some correlated random variables $X, Y$, we want to understand the probability of error when using an estimator $\hat{X}(Y)$ to approximate $X$. The big reveal is that $\Pr(\text{error}) = \text{function}(H(X|Y))$.

**Theorem 13** *Fano's Inequality Let $X, Y$ be dependent random variables. $\hat{X}(Y)$ is an estimator on $X$, so $X \to Y \to \hat{X}$ is a valid Markov chain for the joint distribution. Then*

$$H(P_{err}) + P_{err} \log |\mathcal{X}| \geq H(X|\hat{X}) \geq H(X|Y) \tag{2.38}$$

*A weaker form of the same being:*

$$1 + P_{err} \log |\mathcal{X}| \geq H(X|Y) \tag{2.39}$$

$$P_{err} \geq \frac{H(X|Y) - 1}{\log |\mathcal{X}|} \tag{2.40}$$

$$\tag{2.41}$$

*   **Proof Sketch:**

*   *Represent the "error" event as random variable $E$. It takes value 1 if $\hat{X} \neq X$ and zero if $\hat{X} = X$.*

*   *$P_{err} = \mathbb{E}[E]$.*

*   *Expand $H(E, X|\hat{X}) = H(X|\hat{X}) + H(E|X, \hat{X})$, noting that the last term must be zero.*

- *Also note that $H(E|\hat{X}) \leq H(E)$ (conditioning reduces entropy).*

- *$H(X|E, \hat{X})$ can be expanded and shown to be bounded by $P_e \log |\mathcal{X}|$.*

- *Finally use Markov's inequality for $H(X|\hat{X}) \geq H(X|Y)$.*

- *Combine everything and you should be able to get the strongest version in Equation 2.38.*

You can also strengthen it to

$$H(P_{err}) + P_{err} \log(\underbrace{|\mathcal{X}| - 1}_{(\star)}) \geq H(X|Y)$$

since one element of $\mathcal{X}$ is eliminated from the error entropy calculation by random guessing.

**"Sharpness" of Fano's Inequality:** If you have no information $Y$ to inform $\hat{X}$, your best guess is $\hat{X} = \arg\max_x p(x)$. Equality in Fano's inequality (i.e., $H(P_{err}) + P_{err} \log(|\mathcal{X}| - 1) \geq H(X)$) is achieved by $p(\hat{x}) = 1 - P_{err}$ and $p(x \neq \hat{x}) = P_{err}/(|\mathcal{X}| - 1)$.

**Bound on Collision:** Let $X, X' \tilde{p}(x)$. Then $\Pr\{X = X'\} = \sum_{x \in \mathcal{X}} (p(x))^2$. Then

$$\Pr(X = X') \geq 2^{-H(X)} \tag{2.42}$$

with equality if $X \tilde{u}nif(\mathcal{X})$.

**Proof sketch:** Expand RHS with entropy in $\mathbb{E}[]$ form. Apply Jensen's inequality, and you'll recover $\sum p(x)^2$ as the upper bound.

**Collisions with Different Distributions:** Let $X \tilde{p}(x)$ and $\hat{X} \tilde{r}(x)$. Then

$$\Pr(X \neq \hat{X}) \geq 2^{-H(p) - D(p\|r)} \tag{2.43}$$

$$\Pr(X \neq \hat{X}) \geq 2^{-H(r) - D(r\|p)} \tag{2.44}$$

$$\tag{2.45}$$

**Proof sketch:** We know the LHS is $\sum_x p(x) r(x)$. Expanding the RHS, we can apply Jensen's inequality if $H$ and $D$ are in $\mathbb{E}$ form. Then we will have a bound equal to LHS.

# Chapter 3

# Asymptotic Equipartition Theorem (AEP)

## 3.1 Typicality and AEP Theorem

AEP is the application of the **weak law of large numbers** (WLLN) to entropy. Recall WLLN:

$$\lim_{n \to \infty} \frac{1}{n} \sum_{i=1}^{n} X_i = \mathbb{E}[X] \tag{3.1}$$

Applying to entropy,

**Theorem 14** *Asymptotic Equipartition Theorem Let $X_1 \ldots X_n$ be iid with PMF $p(x)$. Then*

$$\frac{1}{n} \log \frac{1}{p(X_1, \ldots, X_n)} \to H(x) \tag{3.2}$$

$$p(X_1, \ldots, X_n) \to 2^{-nH(x)} \tag{3.3}$$

*in probability. That is, for all $\epsilon > 0$, $\exists\, n$ such that the difference between the sample entropy (LHS) and the true entropy (RHS) is less than $\epsilon$.*

   ***Proof sketch:*** *Apply WLLN to the definition of entropy. LHS = RHS in limit $n \to \infty$..*

**Typical sequences:**   Sequences of $x_i$ with **sample entropy** $\approx nH(X)$.

**Typical set** $A_\epsilon^{(n)}$ **w.r.t.** $p(x)$:    A set of sequences $x^n \in \mathcal{X}^n$ that satisfy

$$2^{-n(H(X)+\epsilon)} \leq p(x_1, \ldots, x_n) \leq 2^{-n(H(X)-\epsilon)} \qquad (3.4)$$

**Theorem 15** *Properties of $A_\epsilon^{(n)}$*

1. *$x^n \in A_\epsilon^{(n)} \rightarrow$ sample entropy $-\frac{1}{n} \log p(x^n) \in [H(X) - \epsilon, H(X) + \epsilon]$.*

2. *$\Pr\{X^n \in A_\epsilon^{(n)}\} = \Pr\{A_\epsilon^{(n)}\} \geq 1 - \epsilon$ for large $n$.*

3. *$|A_\epsilon^{(n)}| \leq 2^{n(H(x)+\epsilon)}$ for all $n$.*

4. *$A_\epsilon^{(n)} \geq (1 - \epsilon)2^{n(H(x)-\epsilon)}$ for sufficiently large $n$.*

   *Proof Sketches:*

1. *Rearrangement of AEP/definition of $A_\epsilon^{(n)}$.*

2. *Apply the lower bound on $\Pr x^n : x^n \in A_\epsilon^{(n)}$.*

3. *$\Pr\{A_\epsilon^{(n)}\} \leq 1$. But we also know that $\Pr(x^n) \geq 2^{-n(H(x)+\epsilon)}$ if $x^n \in A_\epsilon^{(n)}$. Do the math.*

4. *$\Pr\{A_\epsilon^{(n)}\} \geq 1 - \epsilon$ (from 2). Since typical sequences $x^n$ have maximum probability $2^{-n(H(X)-\epsilon)}$, we can derive this bound.*

# Chapter 4

# Data Compression

*Tail end of EIT chapter 3 and the entirety of chapter 5.*

**Data Compression – Problem Statement:** We want to find the *shortest* codes for transmitting sequences $x^n$ where each token is iid with PMF $p(x)$.

- We can leverage notions of **typicality** since we know that $X^n \in A_\epsilon^{(n)}$ with arbitrarily high probability for large $n$.

- **Simple implementation:** Introduce some ordering to elements in $A_\epsilon^{(n)}$. Since the size of $A_\epsilon^{(n)} \approx 2^{nH(X)}$, a binary index would need $n(H + \epsilon) + 1$ bits. This is a pretty good code already!

- For items not in $A_\epsilon^{(n)}$, we can prepend a flag bit and use codes of length $n \log |\mathcal{X}| + 1$.

## 4.1  Definitions & Source Coding Theorem

**Compression/Source Coding Preliminaries:**

- **Source:** Random variable $X : x \in \mathcal{X}$.

- **Codeword Alphabet:** $\mathcal{D}$ is a $D$-ary alphabet.

- **Source Code** $C : \mathcal{X} \to \mathcal{D}^*$ maps from the source alphabet to strings of arbitrary length from the code alphabet.

- **Expected Length** $L(C) = \sum_{x \in \mathcal{X}} p(x)\ell(x)$ where $\ell(x) = |C(x)|$.

- **Assumption:** We tend to represent every $D$-ary alphabet with natural numbers $\{0, 1, \ldots, 1 - D\}$.

- **Non-singular:** Code $C$ is non-singular iff

$$x_1 \neq x_t \implies C(x_1) \neq C(x_2) \qquad (4.1)$$

- **Extension Code:** $C^*$ is the extension of $C$ –

$$C^*(x_1, x_2, \ldots) = C(x_1)C(x_2)\ldots \qquad (4.2)$$

  I.e., the concatenation of $C(x_i)$.

- **Uniquely Decodable:** $C$ is UD if $C^*$ is non-singular.

- **Prefix Code:** No codeword $c(x_1)$ is a prefix of another codeword $c(x_2)$. This also means one can decode without any future information – it is an **instantaneous code**.

**Theorem 16** *Source Coding Theorem Let $X^n$ be iid with each $X_i \tilde{p}(x)$. Then **there exists** a code with lengths satisfying:*

$$\mathbb{E}\big[\frac{1}{n}\ell(X^n)\big] \leq H(X) + \epsilon \qquad (4.3)$$

*for $n$ sufficiently large. In other words, we can represent sequences of length $n$ using $nH(X)$ bits on average!*

   ***Proof sketch:***

- *For large $n$, $\Pr\{A_\epsilon^{(n)}\} \to 1$.*

- *Then $\mathbb{E}[\ell(X^n)] = \sum_{x^n \in \mathbb{X}^n} p(x^n)\ell(x^n)$.*

- *Split the sum into $x^n \in A_\epsilon^{(n)}$ and the compliment.*

- *Apply codes of length $n(H+\epsilon)+2$ for the first sum and lengths $n \log |\mathcal{X}| + 2$ to the second (1 flag bit, 1 rounding bit).*

- *Simplify using $\Pr\{A_\epsilon^{(n)}\}$ and upper bounds on $A_\epsilon^{(n)}$ size.*

## 4.2　Kraft Inequality

The Kraft inequality answers the question, "how short can codes get?"

**Theorem 17** *Kraft Inequality Let $C : \mathcal{X}^n \to \mathcal{D}^*$ be a **prefix code**. Let $\{\ell_i\}_{i=1}^m$ be the codeword lengths for each $x^n \in \mathcal{X}^n$ (i.e., $m = |\mathcal{X}^n|$). Then*

$$\sum_{i=1}^m D^{-\ell_i} \leq 1 \qquad (4.4)$$

*And conversely: For any $\{\ell_i\}$ satisfying the inequality, **there exists a prefix code with those lengths**!*
　　**Proof sketch:**

- *Consider a tree structure for all possible $\mathcal{D}^*$. Each layer adds one character, etc.*

- *At the deepest layer $\ell$, there are $D^\ell$ leaf nodes.*

- *Each non-leaf node on layer $\ell_i$ has $D^{\ell-\ell_i}$ descendants. **To maintain prefix quality**, all descendants are eliminated!*

- *The number of descendants on the final layer must sum to $D^{\ell_{max}}$.*

- *Manipulate these equations and the inequality will pop out :)*

**Theorem 18** *Extended Kraft Inequality Even for an infinite prefix code (i.e., countably infinite codewords), the lengths $\{\ell_i\}_{i=1}^\infty$ will satisfy*

$$\sum_{i=1}^\infty D^{-\ell_i} \leq 1 \qquad (4.5)$$

　　**Proof sketch:** *Apply analogy to floating point numbers/place value. "Descendants" represent intervals, they must be disjoint, etc.*

## 4.3　Finding Optimal Codes

We learned from the Kraft Inequality (Equation 4.4) that the codeword lengths are constrained to follow $\sum D^{-\ell_i} \leq 1$. Optimal codes will there-

fore solve the following optimization problem:

$$\min_{\{\ell_i\}} \sum_{i=1}^{m} p_i \ell_i \tag{4.6}$$

$$s.t. \sum_{i=1}^{m} D^{-\ell_i} \leq 1 \tag{4.7}$$

**Lagrange Multiplier Solution:** $\ell_i^* = -\log_D p_i$ satisfies Kraft inequality and minimizes $\mathbb{E}[\ell_i]$ – expected code length converges to $H_D(x)$. However, we need to round off code lengths so they're integers!

**Theorem 19** *Expected code length inequality For a D-ary code and random variable $X : x \in \mathcal{X}$,*

$$L \geq H_D(X) \tag{4.8}$$

*with equality iff $D^{-\ell_i} = p_i$ where $p_i = p(x_i)$.*

***Proof sketch:*** *Start by expanding $L - H_D(X)$. Combine the sums using product-sum log rule and recover $L - H_D(X) = D(\mathbf{p}\|\mathbf{r})$ where $r_i = D^{-\ell_i}/\sum_j D^{-\ell_j}$. By non-negativity of $D(\|)$, the proof is done.*

**D-adic Distributions:** $p(x)$ is D-adic if $\exists \{n_i\}$ such that

$$p(x_i) = D^{-n_i} \tag{4.9}$$

where each $n_i$ is a natural number.

Generally, we want a $D$-adic distribution $p(x)$ so that our codes achieve expected length $L = \sum p_i \ell_i = H_D(X)$. Failing that, we want to **minimize** $D(d - adic\|p(x))$. We are approximating $p(x)$ with the closet $D$-adic distribution. That's really all that coding methods boil down to!

- Shannon-Fano: Offers good, easy, suboptimal codes.

- Huffman: Truly optimal codes based on $p(x)$ – we actually find the nearest $D$-adic distribution!

# 4.4 Sardinas-Patterson Test for Unique Decodability