

MA4261 Information and Coding Theory

AY24/25 Semester 1

by Isaac Lai

Probability

- $\text{Var}(X) = \mathbb{E}[(X - \mathbb{E}[X])^2] = \mathbb{E}[X^2] - \mathbb{E}[X]^2$
- $\text{Var}(X_1 + X_2) = \text{Var}(X_1) + \text{Var}(X_2) + 2\text{Cov}(X_1, X_2)$
- **Union bound:** In a probability space with σ -algebra \mathcal{F} we have

$$\Pr\left(\bigcup_{i=1}^k A_i\right) \leq \sum_{i=1}^k \Pr(A_i)$$

This holds in the infinite case too.

- $\mathbb{E}[X] = \mathbb{E}_Y[\mathbb{E}_X[X | Y]]$
- Random variables X, Y, Z form a **Markov chain** in the order $X - Y - Z$ if their joint distribution P_{XYZ} satisfies for all $(x, y, z) \in \mathcal{X} \times \mathcal{Y} \times \mathcal{Z}$

$$P_{XYZ}(x, y, z) = P_X(x)P_{Y|X}(y | x)P_{Z|Y}(z | y)$$

This is equivalent to saying X and Z are **conditionally independent given Y** .

- **Markov's Inequality:** Let X be a real-valued non-negative random variable. Then for any $a > 0$ we have $\Pr(X > a) \leq \frac{\mathbb{E}[X]}{a}$.
- **Chebyshev's Inequality:** Let X be a real-valued random variable with mean μ and variance σ^2 . Then for any $a > 0$

$$\Pr(|X - \mu| > a\sigma) \leq \frac{1}{a^2}$$

- **Weak Law of Large Numbers:** For every $\epsilon > 0$,

$$\lim_{n \rightarrow \infty} \Pr\left(\left|\frac{1}{n} \sum_{i=1}^n X_i\right| > \epsilon\right) \leq \lim_{n \rightarrow \infty} \frac{\sigma^2}{n\epsilon^2} = 0$$

Information Quantities

Definition. The **entropy** $H(X)$ of a discrete random variable X is defined by

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

Properties of H

1. $H(X) \geq 0$
2. $H_b(X) = (\log_b a)H_a(X)$ (binary entropy)
3. (Conditioning does not increase entropy) For any two random variables X and Y , $H(X | Y) \leq H(X)$ with equality iff X and Y are independent.
4. $H(X_1, X_2, \dots, X_n) \leq \sum_{i=1}^n H(X_i)$ with equality iff all X_i are independent.

5. $H(X) \leq \log |\mathcal{X}|$ with equality iff X is distributed uniformly over \mathcal{X} .

6. $H(p)$ is concave in p .

7. Han's Inequality:

$$H(X_1, \dots, X_n) \leq \frac{1}{n-1} \sum_{i=1}^n H(X_1, \dots, X_{i-1}, X_{i+1}, \dots, X_n)$$

8. $H(g(X)) \leq H(X)$

Definition. The **relative entropy** $D(p \parallel q)$ of pmf p wrt pmf q is

$$D(p \parallel q) = \sum_x p(x) \log \frac{p(x)}{q(x)}$$

Definition. The **mutual information** between two random variables X and Y is defined as

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

Alternatively,

$$H(X) = E_p \log \frac{1}{p(X)}$$

$$H(X, Y) = E_p \log \frac{1}{p(X, Y)}$$

$$H(X | Y) = E_p \log \frac{1}{p(X | Y)}$$

$$I(X; Y) = E_p \log \frac{p(X, Y)}{p(X)p(Y)}$$

$$D(p \parallel q) = E_p \log \frac{p(X)}{q(X)}$$

Properties of D and I

1. $I(X; Y) = H(X) - H(X | Y) = H(Y) - H(Y | X) = H(X) + H(Y) - H(X, Y)$
2. $D(p \parallel q) \geq 0$ with equality iff $p(x) = q(x)$ for all $x \in \mathcal{X}$
3. $I(X; Y) = D(p(x, y) \parallel p(x)p(y)) \geq 0$ with equality iff $p(x, y) = p(x)p(y)$, i.e. X and Y are independent.
4. If $|\mathcal{X}| = m$ and u is the uniform distribution over \mathcal{X} , then $D(p \parallel q) = \log m - H(p)$.
5. $D(p \parallel q)$ is convex in the pair (p, q) .

Chain rules

- Entropy: $H(X_1, X_2, \dots, X_n) = \sum_{i=1}^n H(X_i | X_{i-1}, \dots, X_1)$
- Mutual information: $I(X_1, X_2, \dots, X_n; Y) = \sum_{i=1}^n I(X_i; Y | X_1, X_2, \dots, X_{i-1})$
- Relative entropy: $D(p(x, y) \parallel q(x, y)) = D(p(x) \parallel q(x)) + D(p(y | x) \parallel q(y | x))$

Important results

- **Jensen's Inequality:** If f is a convex function, then $\mathbb{E}f(X) \geq f(\mathbb{E}X)$
- **Log sum Inequality:** For n positive numbers, a_1, a_2, \dots, a_n and b_1, b_2, \dots, b_n

$$\sum_{i=1}^n a_i \log \frac{a_i}{b_i} \geq \left(\sum_{i=1}^n a_i \right) \log \frac{\sum_{i=1}^n a_i}{\sum_{i=1}^n b_i}$$

with equality iff $\frac{a_i}{b_i} = \text{constant}$.

- **Data-processing Inequality:** If $X \rightarrow Y \rightarrow Z$ forms a Markov chain, $I(X; Y) \geq I(X; Z)$.
- **Sufficient statistic:** $T(X)$ is sufficient relative to $\{f_\theta(x)\}$ iff $I(\theta; X) = I(\theta; T(X))$ for all distributions on θ .
- **Fano's Inequality:** Let $P_e = \Pr\{\hat{X}(Y) \neq X\}$. Then

$$H(P_e) + P_e \log |\mathcal{X}| \geq H(X | Y)$$

This can be loosened to

$$P_e \geq \frac{H(X | Y) - 1}{\log |\mathcal{X}|}$$

- If X and X' are i.i.d., then $\Pr(X = X') \geq 2^{-H(X)}$

Asymptotic Equipartition Property

Definition. The **typical set** of X , a discrete memoryless source (DMS) is defined as

$$A_\epsilon^{(n)}(X) := \left\{ x^n \in \mathcal{X}^n : \left| \frac{1}{n} \log \frac{1}{P_{X^n}(x^n)} - H(X) \right| \leq \epsilon \right\}$$

where for all $x^n \in \mathcal{X}^n$

$$P_{X^n}(x^n) = \Pr(X^n = x^n) = \prod_{i=1}^n P_X(x_i)$$

Theorem (AEP). 1. $\Pr(X^n \in A_\epsilon^{(n)}(X)) \geq 1 - \epsilon$ for all sufficiently large n .

2. The size of the typical set satisfies $(1 - \epsilon)2^{n(H(X) - \epsilon)} \leq |A_\epsilon^{(n)}(X)| \leq 2^{n(H(X) + \epsilon)}$.

Definition (Code). An $(n, 2^{nR})$ -fixed-to-fixed-length source code consists of an encoder f and a decoder φ where

1. $f : \mathcal{X}^n \rightarrow \{1, \dots, 2^{nR}\}$ and

2. $\varphi : \{1, \dots, 2^{nR}\} \rightarrow \mathcal{X}^n$

n is the blocklength of the code and R is the rate of the code.

Definition (Achievable rate). $R \geq 0$ is achievable if there exists a sequence of $(n, 2^{nR})$ -codes such that $\lim_{n \rightarrow \infty} \Pr(\hat{X}^n \neq X^n) = 0$ where $\hat{X}^n = \varphi(M)$ and $M = f(X^n)$ are the reconstructed source and compression index respectively.

Definition (Optimum Source Coding Rate). The optimum source coding rate for the DMS X is $R^*(X) = \inf\{R : R \text{ is achievable}\}$.

Theorem (Fixed-to-Fixed-Length Data Compression).

$$R^*(X) = H(X)$$

Theorem. If $R < H(X)$, then $P_e^{(n)} := \Pr(\hat{X}^n \neq X^n) \rightarrow 1$ as $n \rightarrow \infty$.

Theorem (Han-Verdu Lemma). Fix any $(n, 2^{nR})$ -code. Then $P_e = \Pr(\hat{X}^n \neq X^n)$ satisfies

$$P_e \geq \sup_{\gamma > 0} \Pr\left\{ \frac{1}{n} \log \frac{1}{P_{X^n}(X^n)} \geq R + \gamma \right\} - e^{-n\gamma}$$

Theorem. Let $B_\delta^{(n)} \subset \mathcal{X}^n$ be such that if $X_1, X_2, \dots \sim P_X$, then for every $\delta \in (0, 1)$, $\Pr(X^n \in B_\delta^{(n)}) \geq 1 - \delta$ for all n sufficiently large. Then for any $\delta' > 0$,

$$\frac{1}{n} \log |B_\delta^{(n)}| \geq H(X) - \delta'$$

for n sufficiently large. Here $H(X)$ is computed wrt PMF P_X

Entropy Rates of Stochastic Processes

A **stochastic process** $\{x_i\}_{i \in \mathbb{N}}$ is an indexed sequence of random variables where i is the time.

Definition. A stochastic process is **stationary** if $\Pr(X_1 = x_1, \dots, X_n = x_n) = \Pr(X_{1+\ell} = x_1, \dots, X_{n+\ell} = x_n)$ for all $n \in \mathbb{N}$ and every shift $\ell \in \mathbb{N}$, and for all $x_1, \dots, x_n \in \mathcal{X}$

Definition. A stochastic process is a **Markov chain** if $\forall n \geq 1, \Pr(X_{n+1} = x_{n+1} | X_1 = x_1, \dots, X_n = x_n) = \Pr(X_{n+1} = x_{n+1} | X_n = x_n) \forall x_1, \dots, x_{n+1} \in \mathcal{X}$

Definition. The Markov chain is **time-invariant** if $P(x_{n+1} | x_n)$ does not depend on n . Such a Markov chain is characterised by a transition probability matrix (TPM) $P = [P_{ij}]$, $i, j \in \mathcal{X}$, $P_{ij} = \Pr(X_{n+1} = j | X_n = i)$ for all time-invariant n . In other words, we have $p_{n+1} = p_n P$

If it is possible to go from any state to any other in a finite number of steps, the Markov chain is **irreducible**. If the GCD of the lengths of different paths from a state to itself is 1, the Markov chain is **aperiodic**.

Definition (Entropy rate). Two definitions:

$$H(X) = \lim_{n \rightarrow \infty} \frac{1}{n} H(X_1, X_2, \dots, X_n)$$

$$H'(X) = \lim_{n \rightarrow \infty} H(X_n | X_{n-1}, X_{n-2}, \dots, X_1)$$

For a stationary stochastic process, $H(\mathcal{X}) = H'(\mathcal{X})$

Theorem (Cesaro mean). If $a_n \rightarrow a$ and $b_n = \frac{1}{n} \sum_{i=1}^n a_i$, then $b_n \rightarrow a$.

Theorem (Shannon-McMillan-Breiman). For a stationary, ergodic (irreducible and aperiodic) process, the AEP holds: $\lim_{n \rightarrow \infty} -\frac{1}{n} \log p(X_1, \dots, X_n) = H(X)$

- **Entropy rate of an ergodic Markov chain:**

$$H(X) = H'(X) = H(X_2 | X_1)$$

- **Functions of a Markov chain:** If X_1, X_2, \dots, X_n form a stationary Markov chain and $Y_i = \phi(X_i)$, then

$$H(Y_n | Y^{n-1}, X_1) \leq H(Y) \leq H(Y_n | Y^{n-1})$$

$$\lim_{n \rightarrow \infty} H(Y_n | Y^{n-1}, X_1) = H(Y) = \lim_{n \rightarrow \infty} H(Y_n | Y^{n-1})$$

Fixed-to-Variable-Length Source Coding

Definition. A fixed-to-variable-length (F2V) source code for a random variable X is a map C for \mathcal{X} to $\{0, 1\}^*$. $C(x)$ is the codeword corresponding to $x \in \mathcal{X}$ and $l(w)$ is the length of the codeword corresponding to $x \in \mathcal{X}$.

Definition. The expected length $L(C)$ of a code $C : \mathcal{X} \rightarrow \{0, 1\}^*$ for a random variable $X \sim p_X$ is $L(C) = \sum_{x \in \mathcal{X}} p_X(x) l(x) = \mathbb{E}_{p_X} [l(X)]$.

Definition. A code C is **non-singular** if every $x \in \mathcal{X}$ gets mapped to a different codeword, i.e. for all $x, x' \in \mathcal{X}$ such that $x \neq x'$, we have $C(x) \neq C(x')$.

Definition. The **extension** C^* of a code C is the map from finite-length strings in \mathcal{X} to finite-length strings in $\{0, 1\}^*$.

Definition. A **uniquely decodable** code is one in which its extension is non-singular.

Definition. A code is called **prefix-free** or **instantaneous** if no codeword is a prefix of any other codeword.

Theorem (Kraft's Inequality). For any PF code over an alphabet of size 2, its codeword lengths l_1, l_2, \dots, l_m must satisfy $\sum_{i=1}^m 2^{-l_i} \leq 1$. Conversely, if the inequality is satisfied, then there exists a PF code with those lengths.

Theorem. The expected codeword length L^* of any binary PF code for a random variable X satisfies $L^* \geq H(X)$ with equality iff $2^{-l_i} = p_i$. Moreover, $L^* < H(X) + 1$.

Definition (Shannon code). For all $i \in \mathcal{X}$, $l_i = \left\lceil \log \frac{1}{p_i} \right\rceil$

Theorem (Coding over long blocks). We have $\frac{H(X^n)}{n} \leq L_n^* < \frac{H(X^n)}{n} + \frac{1}{n}$. If $\mathcal{X} = \{X_n\}_{n=1}^\infty$ is a stationary stochastic process, then $L_n^* \rightarrow H(\mathcal{X})$.

Theorem (Wrong code). For the code assignment $l(x) = \left\lfloor \log \frac{1}{q(x)} \right\rfloor$,

$$H(p) + D(p \| q) \leq E_p l(X) < H(p) + D(p \| q) + 1$$

Huffman codes

- We expect that an optimal PF code will have the longest codeword for the two least probable symbols. Otherwise we can delete one bit from the longer one and retain the PF property while decreasing the expected codeword length.

- Algorithm: rank symbols by probability. Combine the two least probable ones, and re-rank the probabilities. Repeat this process until we only have one symbol, then generate the codewords using the binary tree by branching with 0 and 1 at each node.

- WLOG $p_1 \geq p_2 \geq \dots \geq p_m$. The code is optimal iff $\sum p_i l_i$ is minimised.

- There exists an optimal code with $l_1 \leq l_2 \leq \dots \leq l_m$ where $C(m-1)$ and $C(m)$ are siblings that differ only in their last bits.

- The Huffman procedure yields an optimal code.

Channel Capacity

Definition. A **discrete** channel is a system consisting of (1) input alphabet \mathcal{X} , (2) output alphabet \mathcal{Y} , and (3) probability transition matrix $p_{Y|X}$.

Definition. The channel is **memoryless** if the probability distribution of the output at time i depends only on the input at time i , i.e. for all x^n, y^n ,

$$\Pr(Y^n = y^n | X^n = x^n) = \prod_{i=1}^n p_{Y|X}(y_i | x_i)$$

Definition. The **channel capacity** of a DMC $(\mathcal{X}, \mathcal{Y}, p_{Y|X})$ is

$$C = C(p_{Y|X}) = \max_{p_X} I(X; Y)$$

Definition. The n -th **extension of a DMC** $(\mathcal{X}, p(y | x), \mathcal{Y})$ without feedback is the channel $(\mathcal{X}^n, p(y^n | x^n), \mathcal{Y}^n)$ where $p(y^n | x^n) = \prod_{i=1}^n p(y_i | x_i)$ for $(x^n, y^n) \in \mathcal{X}^n \times \mathcal{Y}^n$.

Definition. An (M, n) -**code** for the DMC $(\mathcal{X}, p(y | x), \mathcal{Y})$ consists of (1) the message set $\{1, \dots, M\}$, (2) the encoder $f : \{1, \dots, M\} \rightarrow \mathcal{X}^n$, (3) the decoder $\varphi : \mathcal{Y}^n \rightarrow \{1, \dots, M\}$.

Definition. The **conditional probability of error** of a code (f, φ) of sending a message $w \in [M] = \{1, \dots, M\}$ is

$$\begin{aligned} \lambda_w &= \Pr(\varphi(Y^n \neq w | X^n = w^n(w))) \\ &= \sum_{y^n} p_{Y^n|X^n}(y^n | x^n(w)) \mathbf{1}\{\varphi(y^n) \neq w\} \end{aligned}$$

Definition. The **maximal probability of error** of a code (f, φ) is $\lambda_{\max}^{(n)} = \max_{w \in [M]} \lambda_w$.

Definition. The **average probability of error** of a code (f, φ) is $P_e^{(n)} = \lambda_{\text{ave}}^{(n)} = \frac{1}{M} \sum_{w=1}^M \lambda_w$.

Definition. The **rate** of an (M, n) -code is $R = \frac{1}{n} \log M$ bits per channel use, or $R = \frac{1}{n} \ln M$ nats per channel use.

Definition. A rate $R \geq 0$ is **achievable** for a DMC $p_{Y|X}$ if there exists a sequence of $(2^{nR}, n)$ -codes such that $\lambda_{\max}^{(n)} \rightarrow 0$ as $n \rightarrow \infty$. Note that if $R \geq 0$ is achievable, then $R' \leq R$ is achievable too.

Definition. The **capacity** of a DMC $p_{Y|X}$ is

$$\tilde{C} = \tilde{C}(p_{Y|X}) = \sup\{R \geq 0 : R \text{ is achievable}\}$$

Theorem. $\tilde{C} = C(p_{Y|X}) = \max_{p_X} I(p_X, p_{Y|X})$

Examples

- Noiseless binary channel, noisy channel with non-overlapping output: $C = 1$

- Binary symmetric channel: $C = 1 - H_b(p)$ when p_X is uniform on $\{0, 1\}$

- Binary erasure channel: $C = 1 - \alpha$ when p_X is uniform on $\{0, 1\}$

- Symmetric channels (doubly stochastic PTM, rows and columns are permutations of one another respectively): $C = \log |\mathcal{Y}| - H(\mathbf{r})$ where \mathbf{r} is the distribution on one row

Jointly Typical Sequences

Definition. The set $A_\epsilon^{(n)}(X, Y)$ of jointly typical sequences (x^n, y^n) wrt $p_{X, Y}$ is

$$A_\epsilon^{(n)} = A_\epsilon^{(n)}(X, Y) = \{(x^n, y^n) \in \mathcal{X}^n \times \mathcal{Y}^n : \begin{aligned} &\left| -\frac{1}{n} \log p_{X^n}(x^n) - H(X) \right| < \epsilon, \\ &\left| -\frac{1}{n} \log p_{Y^n}(y^n) - H(Y) \right| < \epsilon, \\ &\left| -\frac{1}{n} \log p_{X^n, Y^n}(x^n, y^n) - H(X, Y) \right| < \epsilon \end{aligned}\}$$

where $p_{X^n, Y^n}(x^n, y^n) = \prod_{i=1}^n p_{X, Y}(x_i, y_i)$.

Theorem (Joint AEP). 1. $\Pr((X^n, Y^n) \in A_\epsilon^{(n)}) \rightarrow 1$ as $n \rightarrow \infty$

2. $|A_\epsilon^{(n)}| \leq 2^{nH(X, Y)}$

3. If $(\tilde{X}^n, \tilde{Y}^n) \sim p_{X^n}(x^n) p_{Y^n}(y^n)$, then

$$\begin{aligned} (1 - \epsilon) 2^{-n(I(X; Y) + 3\epsilon)} &\leq \Pr((\tilde{X}^n, \tilde{Y}^n) \in A_\epsilon(x, y)) \\ &\leq 2^{-n(I(X; Y) - 3\epsilon)} \end{aligned}$$

Channel Coding Theorem

Theorem (Direct/Achievability). For a DMC, all rates $R < C$ are achievable. For all $R < C$, there exists a sequence (in $n \rightarrow \infty$) of $(2^{nR}, n)_{n \in \mathbb{N}}$ -codes with $\lambda_{\max}^{(n)} \rightarrow 0$.

Theorem (Converse/Impossibility). Conversely, any sequence of $(2^{nR}, n)_{n \in \mathbb{N}}$ -codes with $\lambda_{\text{ave}}^{(n)} \rightarrow 0$ satisfies $R \leq C$. (Proof using Fano's inequality)

Proof of achievability

1. Fix $p_X(x)$. Generate the codewords in iid fashion.
2. By symmetry of codebook generation and uniformity of codewords, WLOG choose $w = 1$. Calculate the error probability and bound it.
3. Choose a rate R such that the bound approaches 0 as $n \rightarrow \infty$.

Definition. The **feedback capacity** $C_{FB}(p_{Y|X})$ is the supremum of all achievable rates with feedback codes.

Theorem. $C_{FB} = C = \max_{p_X} I(X; Y)$

Theorem (Source-channel Separation Theorem). If $\{V_n\}$ is a finite-alphabet stationary stochastic process that satisfies AEP and $H(\mathcal{V}) < C$, then there exists a source-channel code (f_n, φ_n) with $P_e^{(n)} \rightarrow 0$. This code can be realised by a separation scheme. Conversely, for all stationary stochastic processes $\{V^n\}$ with $H(\mathcal{V}) > C$, $P_e^{(n)} \not\rightarrow 0$, i.e. $\limsup_{n \rightarrow \infty} P_e^{(n)} > 0$.

Gaussian Channels

Definition. The **differential entropy** $h(X)$ of a continuous RV X is

$$h(X) = - \int_{s_X} f_X(x) \log f_X(x) dx$$

- $h(\mathcal{N}(0, \sigma^2)) = \frac{1}{2} \log 2\pi e \sigma^2$
- $h(\mathcal{N}_n(\mu, K)) = \frac{1}{2} \log(2\pi e)^n |K|$
- $D(f \| g) = \int f \log \frac{f}{g} \geq 0$
- $h(X + c) = h(X)$
- $h(aX) = h(X) + \log |a|$

Theorem. Let X_1, X_2, \dots be a sequence of i.i.d. continuous RVs with common density f_X . Then $-\frac{1}{n} \log f_X(X_1, \dots, X_n) \rightarrow \mathbb{E}[-\log f_X(X)] = h(X)$.

Definition. $\forall \epsilon > 0$, the ϵ -**weakly typical set** wrt $f = f_X$ is

$$A_\epsilon^{(n)}(X) = \left\{ x_n \in \mathbb{R}^n : \left| -\frac{1}{n} \log f_X(x_1, \dots, x_n) - h(X) \right| < \epsilon \right\}$$

where $f_{X^n}(x_1, \dots, x_n) = \prod_{i=1}^n f_X(x_i)$. We have (1) $\Pr(X_n \in A_\epsilon^{(n)}(X)) \rightarrow 1$ (2) $\text{Vol}(A_\epsilon^{(n)}(X)) \leq 2^{n(h(X) + \epsilon)}$ (3) $\text{Vol}(A_\epsilon^{(n)}(X)) \geq (1 - \epsilon) 2^{n(h(X) - \epsilon)}$

Theorem. The information capacity of the Gaussian channel with $Y_i = X_i + Z_i$, $Z_i \sim \mathcal{N}(0, N)$, power constraint $\frac{1}{n} \sum_{i=1}^n x_i^2 \leq P$ is

$$C = \max_{f_X : \mathbb{E}[X^2] \leq P} I(X; Y) = \frac{1}{2} \log \left(1 + \frac{P}{N} \right)$$

bits per channel use.

Theorem (Water-filling). For k parallel Gaussian channels, find optimal power allocation P_1, \dots, P_k by differentiating the Lagrangian

$L = \sum_{i=1}^k \frac{1}{2} \log \left(1 + \frac{P_i}{N} \right) - \lambda \left(\sum_{i=1}^k P_i - P \right)$. Solution will be $P_i = (v - N_i)^+$ s.t. $\sum_{i=1}^k (v - N_i)^+ = P$.

Coding

Introduction

Definition. $\mathcal{A} = \{a_1, \dots, a_q\}$ is the **alphabet**; a_i are the **symbols** of the alphabet. A **block code** C of length n over \mathcal{A} is a subset of \mathcal{A}^n . Any vector $c \in C$ is called a **codeword**. $|C|$ is the **size** of the code. A code of length n and size M is an (n, M) -code.

Definition. Let C be an (n, M) -code over alphabet \mathcal{A} of size q . **Dimension** of code is $\log_q M$, **rate** of code is $\frac{1}{n} \log_q M$.

Definition. The **Hamming distance** is $d(x, y) = \sum_{i=1}^n \mathbf{1}\{x_i \neq y_i\}$.

Definition. Let C be a code of length n over alphabet \mathcal{A} . The **nearest neighbour** decoding rule $D(\cdot)$ states that $D(x) = c_x = \arg \min_{c \in C} d(x, c)$. If there exists more than one codeword achieving this minimum distance, D outputs \perp .

Definition. Let C be a code. The **distance** of the code is $d(C) = \min\{d(c_1, c_2) : c_1, c_2 \in C, c_1 \neq c_2\}$. An (n, M) -code of distance d is called an (n, M, d) -code.

Definition. Let $C \subset \mathcal{A}^n$ be a code. C **detects u errors** if $\forall c \in C, \forall x \in \mathcal{A}^n \setminus \{c\}$ it holds that $d(x, c) \leq u \implies x \notin C$. C **corrects v errors** if $d(x, c) \leq v \implies \text{NN decoding of } x \text{ outputs } c$.

Theorem. (1) C detects U errors $\iff d(C) > u$. (2) C corrects v errors $\iff d(C) \geq 2v + 1$. In other words, there exists a decoder (i.e. NN decoder) that can correct up to $\lfloor \frac{d(C)-1}{2} \rfloor$ errors.

Theorem. Let C be an (n, M, d) -code over \mathcal{A} and $v, u \in \mathbb{N}$ s.t. $2v + u \leq d(C) - 1$. Then there exists a decoder $D : \mathcal{A}^n \rightarrow C \cup \{\perp\}$ s.t. (i) if number of errors $\leq v$, errors can be corrected and (ii) if number of errors $\leq v + u$, they can be corrected.

Finite Fields

Definition. A group satisfies closure, associativity, identity, and inverse properties. Alternatively, replace closure and inverse properties with permutation property: $\forall a \in B, a \oplus G = \{a \oplus b : b \in G\}$ is a permutation of G .

Definition. A field \mathbb{F} has at least two elements, with two operations \oplus and \ast s.t. (1) (\mathbb{F}, \oplus) is an abelian group, (2) $\mathbb{F}^\ast = \mathbb{F} \setminus \{0\}$ is an abelian group under \ast , (3) $\forall a, b, c, \in \mathbb{F}, (a \oplus b) \ast c = (a \ast c) \oplus (b \ast c)$.

Theorem. $\{0, 1, \dots, p-1\}$ forms a field \mathbb{F}_p under mod- p addition and multiplication iff p is prime.

Theorem. If $g(x)$ is a prime polynomial of degree m over a prime field \mathbb{F}_p , then the set of remainder polynomials $R_{\mathbb{F}_p, m}$ with mod- $g(x)$ arithmetic forms a finite field $\mathbb{F}_{g(x)}$ with p^m elements.

Linear Codes

Definition. A **linear code** with length n over \mathbb{F}_q is a subspace of \mathbb{F}_q^n .

Definition. Let C be a linear code over \mathbb{F}_q . Then

- The **dual code** of C is $C^\perp = \{x \in \mathbb{F}_q^n : \forall c \in C, \langle x, c \rangle = 0\}$.
- The **dimension** of C is the dimension of C as a subspace of \mathbb{F}_q^n , denoted $\dim(C)$.

Theorem. Let C be a linear code over \mathbb{F}_q . Then $|C| = q^{\dim(C)}$, C^\perp is a linear code and $\dim(C) + \dim(C^\perp) = n$, $(C^\perp)^\perp = C$.

Definition. C is **self-orthogonal** if $C \subset C^\perp$. C is **self-dual** if $C = C^\perp$.

Theorem. C is a self-orthogonal linear code of length $n \implies \dim(C) \leq \frac{n}{2}$. C is a self-dual linear code of length $n \implies \dim(C) = \frac{n}{2}$.

Definition. The **Hamming weight** of $x \in \mathbb{F}_q^n$ $\text{wt}(x)$ is the number of nonzero elements of x , i.e. $\text{wt}(x) = d(x, 0)$. If $x, y \in \mathbb{F}_q^n$, we have $\text{wt}(x + y) = \text{wt}(x) + \text{wt}(y) - 2\text{wt}(x \ast y)$, $\text{wt}(x) + \text{wt}(y) \geq \text{wt}(x + y)$, $\text{wt}(x) + \text{wt}(y) \geq \text{wt}(x + y) \geq \max\{\text{wt}(x) - \text{wt}(y), 0\}$.

Theorem. For a (not necessarily linear) code C , $\text{wt}(C) = \min_{c \in C: c \neq 0} \text{wt}(C)$. If C is linear, then $d(C) = \text{wt}(C)$.

Definition. A **generator matrix** G for a linear code C is a matrix whose rows form a basis for C . A **parity check matrix** H for a linear code C is a generator matrix for C^\perp . If C is an $[n, k]$ -linear code over \mathbb{F}_q , then $G \in \mathbb{F}_q^{k \times n}$ and $H \in \mathbb{F}_q^{(n-k) \times n}$. To show a matrix is a generator matrix, it suffices to check that its rows are linearly independent.

Definition. A generator matrix G is in **standard form** if it is of the form $[I_k \mid X]$ where $I_k \in \mathbb{F}_q^{k \times k}$ is the identity and $X \in \mathbb{F}_q^{k \times (n-k)}$. A parity check matrix H is in **standard form** if it is of the form $[Y \mid I_{n-k}]$ for some $Y \in \mathbb{F}_q^{(n-k) \times k}$. In particular, $Y = -X^T$.

Lemma. Let C be a linear $[n, k]$ -code with generator G . $\forall v \in \mathbb{F}_q^n, v \in C^\perp$ iff $vG^T = 0 \in \mathbb{F}_q^k$. H is a parity check matrix iff its rows are LI and $HG^T = 0$.

Theorem. Let C be a linear code and H a parity-check matrix for C . Then (1) $d(C) \geq d$ iff every subset of $d-1$ columns of H is LI, (2) $d(C) \leq d$ iff there exists a subset of d columns of H that is LD.

Corollary. Let C be a linear code and H be its parity check matrix. Then $d(C) = d$ iff every subset of $d-1$ columns in H are LI and there exists a subset of d columns in H that is LD.

Bounds

Definition. For a linear code $[n, k]$ -code C , its **rate** is $R(C) = \frac{k}{n}$.

Definition. For a $[n, M, d]$ -code over \mathbb{F}_q , the **relative distance** of C is $\delta(C) = \frac{d-1}{n}$.

Definition. Let A be an alphabet of size $q > 1$ and fix blocklength n and distance d . Define $A_q(n, d) = \max\{M : \exists (n, M, d)\text{-code over } A\}$. An (n, M, d) -code for which $M = A_q(n, d)$ is an **optimal code**.

Definition. Let $q > 1$ be a prime power and fix blocklength n and distance d . Define $B_q(n, d) = \max\{q^k : \exists [n, k, d]\text{-linear code over } \mathbb{F}_q^n\}$. A linear $[n, k, d]$ -code for which $q^k = B_q(n, d)$ is an **optimal linear code**.

Theorem. Let $q \geq 2$ be a prime power. Then for every n , (1) $\forall 1 \leq d \leq n, B_q(n, d) \leq A_q(n, d) \leq q^n$, (2) $B_q(n, 1) = A_q(n, 1) = q^n$, (3) $B_q(n, n) = A_q(n, n) = q$.

Definition. Let A be an alphabet of size $q > 1$. Then $\forall u \in A^n$ and $\forall r \in \mathbb{N}$, a **sphere** with centre u and radius r is $S_A(u, r) = \{v \in A^n : d(u, v) \leq r\}$. The **volume** of $S_A(u, r)$ is $V_q^n(r) = |S_A(u, r)|$.

$$V_q^n(R) = \begin{cases} \sum_{i=0}^r \binom{n}{i} (q-1)^i & 0 \leq r \leq n \\ q^n & r > n \end{cases}$$

Theorem (Sphere Covering Lower Bound). For every natural number $q > 1$ and $n, d \in \mathbb{N}$ s.t. $1 \leq d \leq n$, $A_q(n, d) \geq \frac{q^n}{V_q^n(d-1)}$. Aka Gilbert-Varshamov Lower Bound

Theorem (Hamming (Sphere Packing) Upper Bound). For every natural number $q > 1$ and every $n, d \in \mathbb{N}$ s.t. $1 \leq d \leq n$, $A_q(n, d) \leq \frac{q^n}{V_q^n(\lfloor \frac{d-1}{2} \rfloor)}$

Definition. A code over an alphabet of size q with parameters (n, M, d) is **perfect** if it achieves the Hamming (sphere packing) upper bound. Perfect \implies optimal, but the converse is not true.

Definition (Binary Hamming Code). Let $r \geq 2$ and let C be a binary linear code with $n = 2^r - 1$ whose parity check matrix H is s.t. the columns are all of the nonzero vectors in \mathbb{F}_2^r . C is a binary Hamming code of length $2^r - 1$, denoted $\text{Ham}(r, 2)$.

Theorem. • All binary Hamming codes of a given length are equivalent.

- $\forall r \in \mathbb{N}, \dim(\text{Ham}(r, 2)) = k = 2^r - r - 1$.
- $\forall r \in \mathbb{N}, d(\text{Ham}(r, 2)) = d = 3 \implies$ code can correct 1 error.
- Hamming codes are perfect.

Theorem (Singleton Bound). For every $q \in \mathbb{N}$ and all $n, d \in \mathbb{N}$ with $1 \leq d \leq n, A_q(n, d) \leq q^{n-d+1}$. In particular, if C is a linear $[n, k, d]$ -code, then $k \leq n - d + 1$. This is only of interest for large q .

Definition. A linear code with parameters $[n, k, d]$ s.t. $k = n - d + 1$ is called a **maximum distance separable (MDS)** code.

Theorem. Let C be an $[n, k, d]$ -linear code over \mathbb{F}_q with generatr matrix G and parity check matrix H . Then C is MDS \iff every subset of $n - k$ columns of H is LI \iff every subset of k columns of G is LI $\iff C^\perp$ is MDS.

Lemma. Let $0 \leq \lambda \leq \frac{1}{2}$. Then

$$V_2^n(\lambda n) = \sum_{i=0}^{\lambda n} \binom{n}{i} \leq 2^{nH(\lambda)}$$
$$V_2^n(\lambda n) \geq \frac{1}{n+1} 2^{nH(\lambda)}$$

$$\lim_{n \rightarrow \infty} \frac{1}{n} \log V_2^n(n\lambda) = H(\lambda)$$

Theorem (Asymptotic Sphere Packing Bound). For every binary code C with asymptotic relative distance $\delta \leq \frac{1}{2}$ and rate R , we have $R \leq 1 - H(\delta/2)$.

Theorem (Asymptotic Gilbert-Varshamov Bound). Let n, k be s.t. $R \leq 1 - H(\delta)$ where $R = \frac{k}{n}$ and $\delta = \frac{d-1}{n} \leq \frac{1}{2}$. Then there exists a binary linear code C_n with rate R and distance at least d .

Theorem (Asymptotic Singleton Bound). The rate R and relative distance δ of a code over a q -ary alphabet satisfy $R \leq 1 - \delta$.

Reed-Solomon Codes

Let $n \in \mathbb{N}, 1 \leq k \leq n$, and q be a prime power s.t. $q \geq n$. Construct $\text{RS}_{q,n,k}$ as follows:

- Choose n different evaluation points $\alpha_1, \dots, \alpha_n \in \mathbb{F}_q$.
- Let $(m_0, \dots, m_{k-1}) \in \mathbb{F}_q^k$ be the message to be sent. Define a polynomial $C_m(x) = \sum_{i=0}^{k-1} m_i x^i$.
- Encode $m \in \mathbb{F}_q^k$ as $\text{RS}(m) = (C_m(\alpha_1), C_m(\alpha_2), \dots, C_m(\alpha_n)) \in \mathbb{F}_q^n$.

Theorem. $\text{RS}_{q,n,k}$ is a linear $[n, k, n - k + 1]$ -code and is thus MDS.

Berlekamp-Welch Algorithm

Problem: Given $y = (y_1, \dots, y_n) \in \mathbb{F}_q^n$, find $C(x) \in \mathbb{F}_q[x]$ s.t. (1) $\deg(C(x)) \leq k-1$ and (2) $C(\alpha_i) \neq y_i$ for at most $e \leq \lfloor \frac{n-k}{2} \rfloor$ values of $i \in [n]$, else return FAIL.

$$\text{Error locator polynomial: } E(x) = \prod_{i \in [n]: y_i \neq C(\alpha_i)} (x - \alpha_i)$$

$$\text{Key equation: } y_i \cdot E(\alpha_i) = C(\alpha_i)E(\alpha_i)$$

Let $Q(x) = C(x)E(x)$. Algorithm:

- Find
 - Monic degree e polynomial $E(x)$
 - Degree $\leq e + k - 1$ polynomial $Q(x)$ so that $y_i \cdot E(\alpha_i) = Q(\alpha_i)$
 - If $E(x)$ or $Q(x)$ do not exist, return FAIL
- Let $\tilde{C}(x) := Q(x)/E(x)$
 - If $d(\tilde{C}, y) = \sum_{i=1}^n d(\tilde{C}(\alpha_i), y_i) > e$ return FAIL
 - Else return \tilde{C}

Step 1 involves solving a linear system which takes $O(n^3)$ time, while step 2 involves polynomial division which takes $O(n^2)$ time. Hence the Berlekamp-Welch algorithm runs in $O(n^3)$ time.