

18-847F: Special Topics in Computer Systems

Foundations of Cloud and Machine Learning Infrastructure

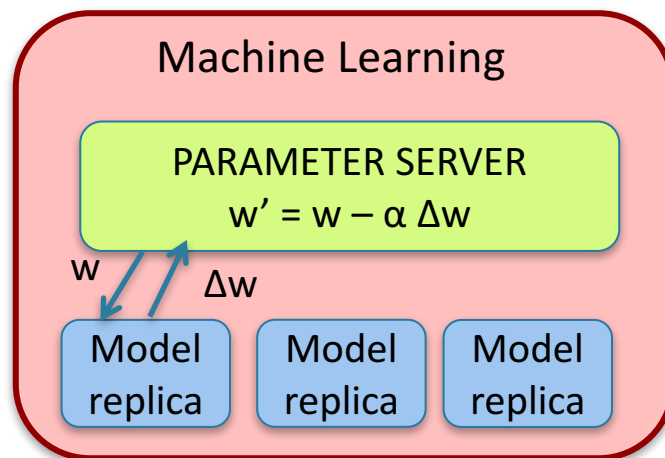
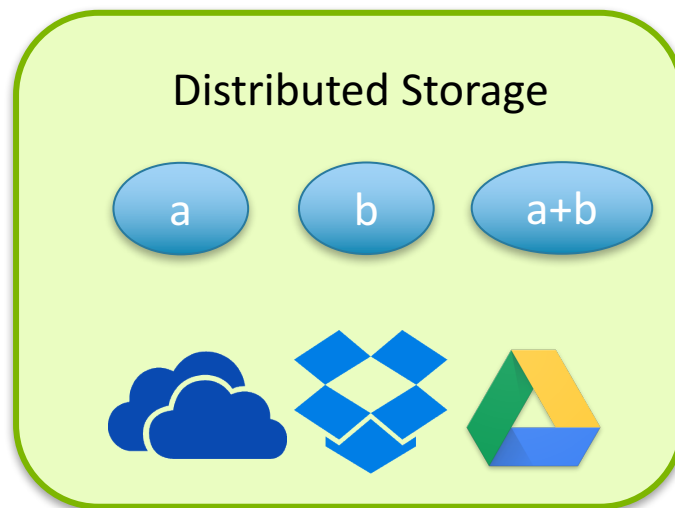
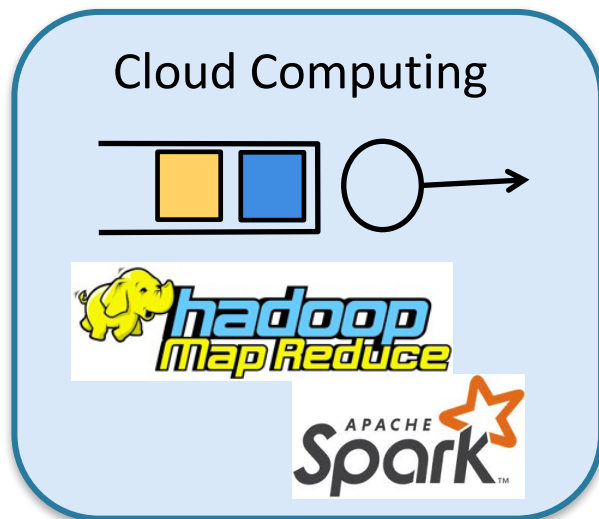


Lecture 8: Intro to Coding Theory

Foundations of Cloud and Machine Learning Infrastructure

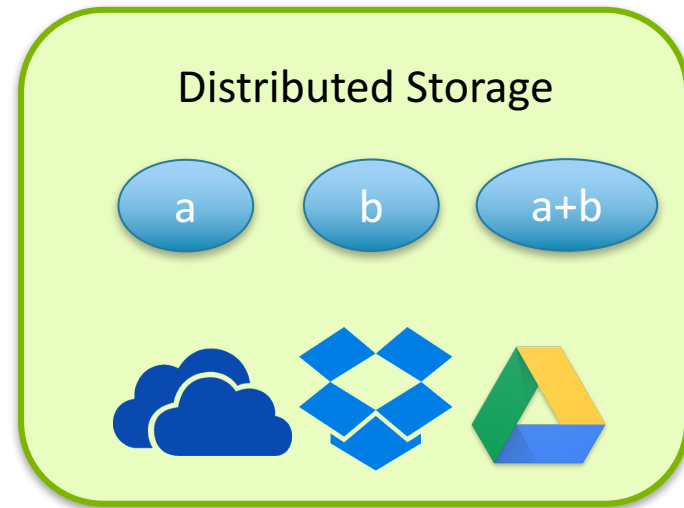


Topics Covered



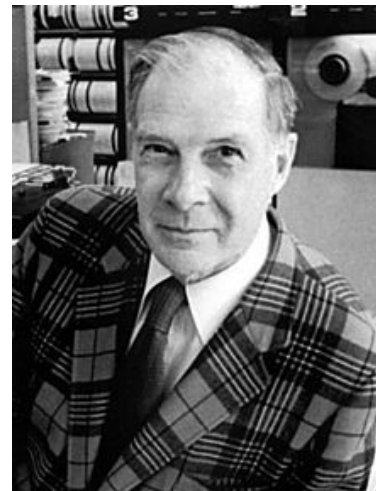
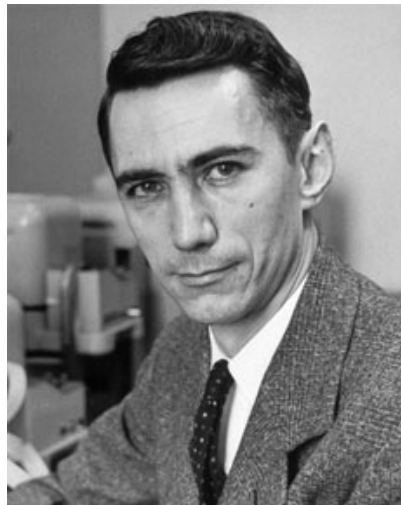
Topics Covered

- Coding for locality/repair
- Reducing latency in content download
- Coded Computing



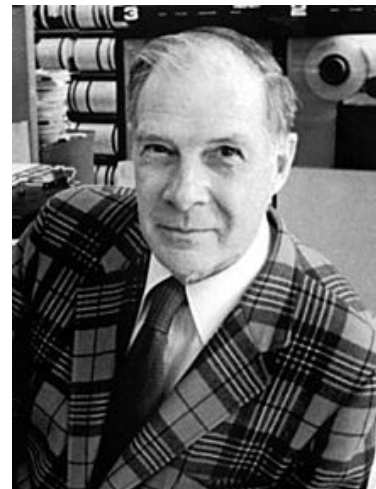
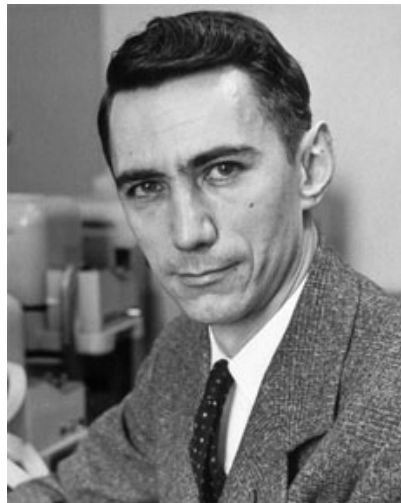
Coding Theory

- For reliable communication in presence of noise
- Bell Labs was one of the leaders in 1950's
- Key figures: Claude Shannon and Richard Hamming



Coding Theory

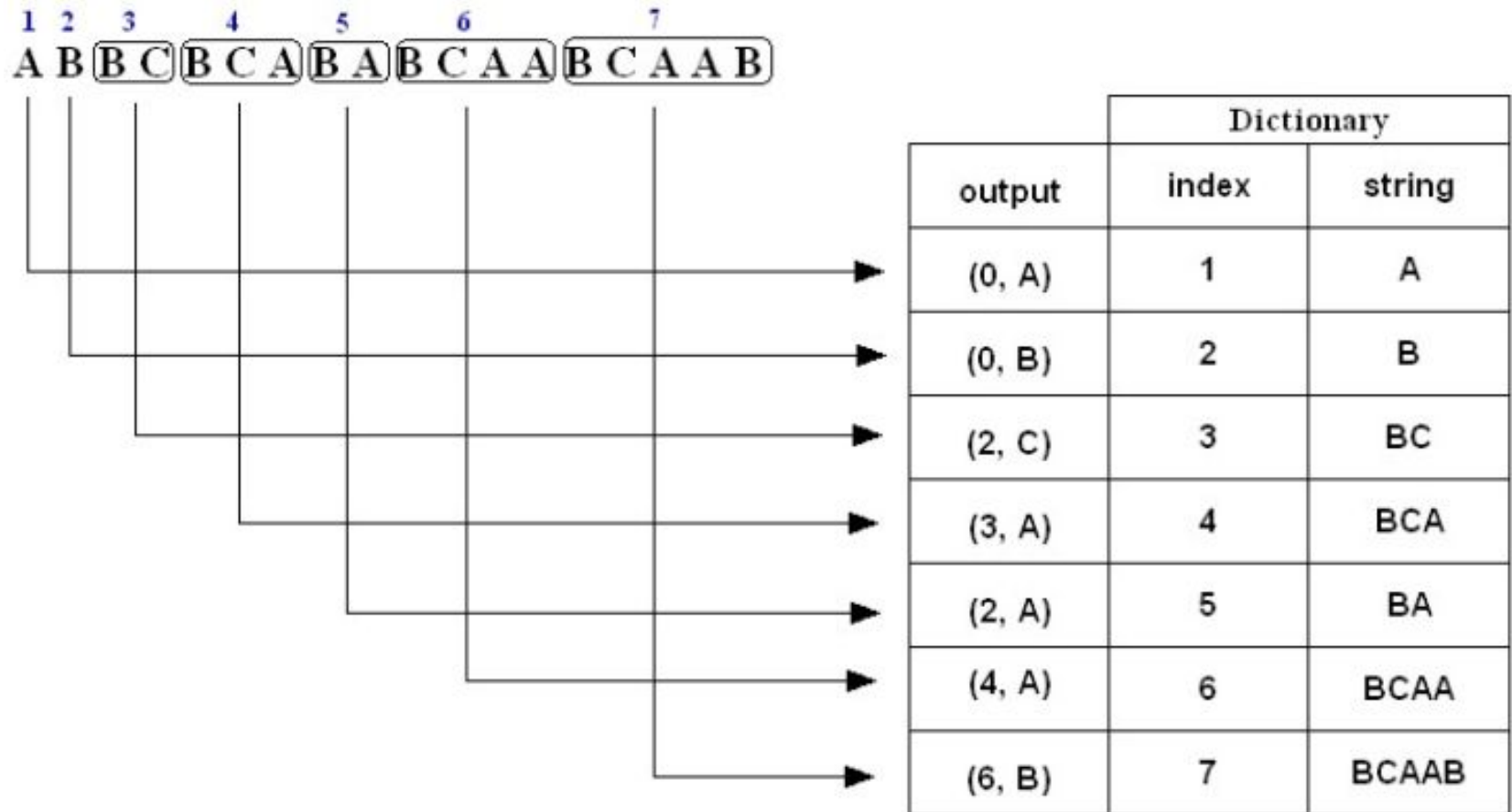
- Two types of Coding:
 - Source Coding: Data Compression
 - Channel Coding: Error Correction



Source Coding – Removing Redundancy

- Huffman Coding
- Zip Data Compression: Lempel-Ziv Coding
- Image/Video Compression: JPEG, MPEG
- Modern applications: Gradient & Model Compression

Source Coding: Lempel-Ziv Coding



Source Coding – Removing Redundancy

- Huffman Coding
- Zip Data Compression: Lempel-Ziv Coding
- Image/Video Compression: JPEG, MPEG
- Modern applications: Gradient & Model Compression

Channel Coding – Adding Redundancy

- Repetition Code
 - $0 \rightarrow 000$: Rate: $1/3$
 - If receive $0??$ we can recover from 2 erasures
- $(3,2)$ code: Data bits: a, b Parity bit: $(a \text{ XOR } b)$
 - Example: $011, 110$: Rate $2/3$
 - If we receive $0?1$ or $?10$ we can correct the failed bit
 - 2 bit symbols: $(01) \text{ ? } (11)$

Simplest Channel Codes

- Repetition Code
 - $0 \rightarrow 000$: Rate: $1/3$
 - If receive $0??$ we can recover from 2 erasures
- $(3,2)$ code: Data bits: a, b Parity bit: $(a \text{ XOR } b)$
 - Example: $011, 110$: Rate $2/3$
 - If we receive $0?1$ or $?10$ we can correct the failed bit
 - 2 bit symbols: (01) **$(1,0)$** (11)

Linear Codes: Definition and Notation

Linear Codes

An (n,k) linear code C is a dimension- k subspace of F_q^n , where F_q is a finite field of q elements

Generator Matrix

G is an $k \times n$ matrix for code C , if its k rows span C

For an $(7,4)$
binary ($q=2$) code

$$G = \begin{pmatrix} 1 & 0 & 0 & 0 & 1 & 1 & 0 \\ 0 & 1 & 0 & 0 & 0 & 1 & 1 \\ 0 & 0 & 1 & 0 & 1 & 1 & 1 \\ 0 & 0 & 0 & 1 & 1 & 0 & 1 \end{pmatrix}$$

Linear Codes: Definition and Notation

With an $(7,4)$ code, we encode a 4-bit string (a,b,c,d) as

The code is said to be systematic if $G = [I_k \mid A]$

$$\begin{array}{l} a \\ b \\ c \\ d \end{array} \begin{pmatrix} 1 & 0 & 0 & 0 & 1 & 1 & 0 \\ 0 & 1 & 0 & 0 & 0 & 1 & 1 \\ 0 & 0 & 1 & 0 & 1 & 1 & 1 \\ 0 & 0 & 0 & 1 & 1 & 0 & 1 \end{pmatrix}$$

$a, \quad b, \quad c, \quad d, \quad a+c+d, \quad a+b+c, \quad b+c+d$

Linear Codes: Definition and Notation

Rate of the Code

An (n,k) code has code rate $r = k/n$

For an $(7,4)$
binary ($q=2$) code

$$G = \begin{pmatrix} 1 & 0 & 0 & 0 & 1 & 1 & 0 \\ 0 & 1 & 0 & 0 & 0 & 1 & 1 \\ 0 & 0 & 1 & 0 & 1 & 1 & 1 \\ 0 & 0 & 0 & 1 & 1 & 0 & 1 \end{pmatrix}$$

Linear Codes: Definition and Notation

Distance

Minimum Hamming distance between any two codewords. For linear codes, it is the minimum Hamming weight of a non-zero codeword.

Distance = $d = 3$

For an $(7,4)$

binary ($q=2$) code

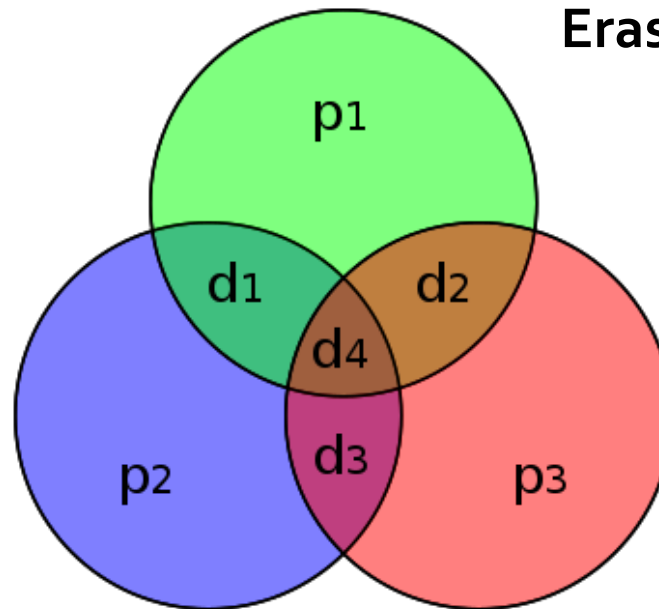
$$G = \begin{pmatrix} 1 & 0 & 0 & 0 & 1 & 1 & 0 \\ 0 & 1 & 0 & 0 & 0 & 1 & 1 \\ 0 & 0 & 1 & 0 & 1 & 1 & 1 \\ 0 & 0 & 0 & 1 & 1 & 0 & 1 \end{pmatrix}$$

Codes with $d = n - k + 1$ are called maximum-distance separable (MDS) codes

Hamming Codes

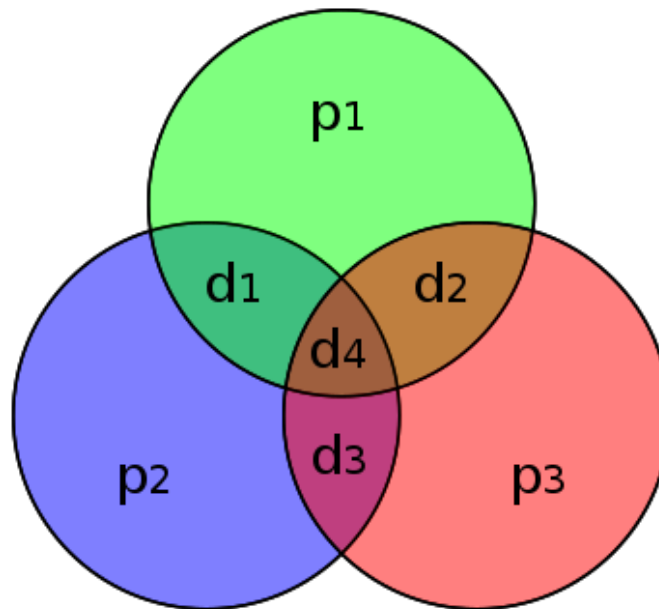
- (7,4) Hamming Code: 4 data bits, 3 parity bits
- Parity $p_1 = d_1 \oplus d_2 \oplus d_4$
- Can correct 1-bit errors or 2-bit erasures
- Can detect 1 or 2-bit errors

Errors: Bit Flips $1 \rightarrow 0$ or $0 \rightarrow 1$
Erasures: You receive a ?



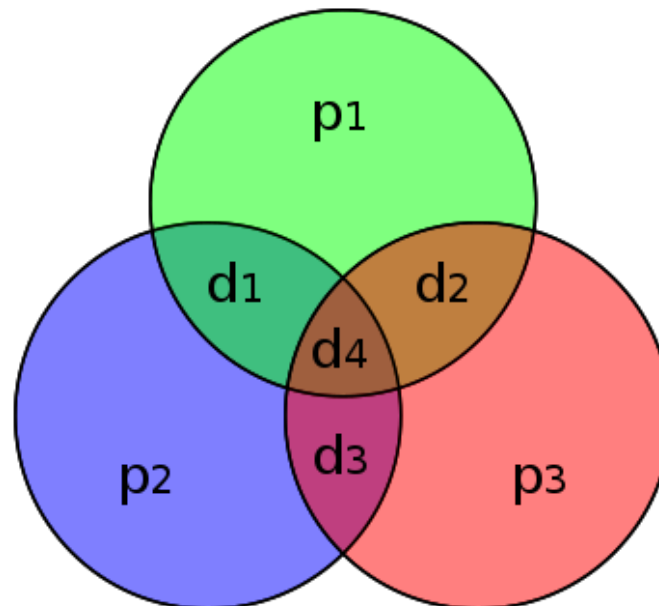
Concept Check: Erasure Codes

- What is the rate and distance of this code?
- Correct the 2 erasures
 - $(d_1, d_2, d_3, d_4, p_1, p_2, p_3) = (0, ?, 1, ?, 1, 0, 0)$



Concept Check: Answer

- What is the rate of the code? $r = 4/7, d = 3$
- Correct the 2 erasures
 - $(d_1, d_2, d_3, d_4, p_1, p_2, p_3) = (0, 0, 1, 1, 1, 0, 0)$

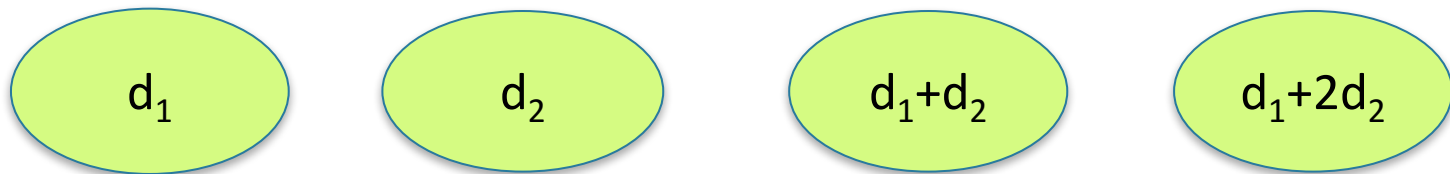


(n,k) Reed-Solomon Codes: 1960

- Data: $d_1, d_2, d_3, \dots, d_k$
- Polynomial: $d_1 + d_2 x + d_3 x^2 + \dots + d_k x^{k-1}$
- Parity bits: Evaluate at $n-k$ points:
 - $x=1:$ $d_1 + d_2 + d_3 + d_4$
 - $x=2:$ $d_1 + 2 d_2 + 4 d_3 + 8 d_4$
 - $x=3:$
 - $x=4:$
 - $x=n:$...
- Can solve for the coefficients from any k coded symbols

Example: (4,2) Reed-Solomon Code

- Data: $d_1, d_2 \rightarrow$ Polynomial: $d_1 + d_2 x + d_3 x^2 + \dots d_k x^{k-1}$



- Can solve for the coefficients from any k coded symbols
- Microsoft uses (7, 4) code
- Facebook uses (14,10) code