

LECTURE 7

Note Title

2/25/2006

TODAY: DATA COMPRESSION UPPER BOUNDS

- TIGHTNESS OF KRAFT
 - ENTROPY BASED ENCODING (SHANNON)
 - = HUFFMAN CODING (OPTIMAL)
 - Other schemes
-

Review of last lecture

Definitions

- Lossless or non-singular compressions
- Uniquely Decodable
- Prefix code

Limitations

① C is prefix-free of lengths $l_1, \dots, l_n \Rightarrow \sum D^{-l_i} \leq 1$

uniquely-decodable

Mu-Millan

② Entropy lower bound :

if C is prefix-free

$$\text{then } E_x [|C(x)|] \geq \frac{H(x)}{\log D}$$

③ if C is any code with expected length

L , then \exists prefix code C'

with expected length $L + O(\sqrt{L})$.

(so non-prefix-codes are not much better anyway)

Today: first see Kraft's tight.

Theorem: if l_1, \dots, l_n satisfy $\sum 2^{-l_i} \leq 1$
then \exists prefix free code C with $|C(i)| = l_i$.

Proof: Back to infinite D-ary tree.

Assume $l_1 \leq l_2 \leq \dots \leq l_n$.

Assign codewords as follows.

For $C(1)$ - assign left most node at level l_1 & delete subtree below it.

$C(2)$ - assign left most node in
remaining tree at level l_2
& delete subtree below it.

\vdots

To prove we don't get stuck need to prove that if

$$\sum_{i=1}^K D^{-l_i} \leq 1 - D^{-l_{i+1}}$$

then \exists free node at level l_{i+1} .

(Simple exercise: left to class).

Entropy upper bound:

Recall proof of Entropy lower bound.

set $D^{-l_i} = q_i$ & q_i look like

probabilities. Resulting length is

$$\frac{H(x)}{\log D} + D(p \parallel q).$$

To minimize compression length should

make $q_i \sim p_i$

$$\Rightarrow D^{-l_i} \sim p_i$$

$$\Rightarrow l_i = \frac{-\log p_i}{\log D}$$

But l_i needs to be an integer?

$$\text{so set } l_i = \left\lceil \frac{-\log p_i}{\log D} \right\rceil$$

$$\leq \frac{-\log p_i}{\log D} + 1$$

By tightness of Kraft, there exists
(prefix) code C with $|C(i)| = l_i$

Expected length of C

$$= \sum_i P_i l_i$$

$$\leq - \sum \frac{P_i \log P_i}{\log D} + \sum P_i$$

$$= \frac{H(x)}{\log D} + 1$$

Conclude

for prefix code C

$$\frac{H(x)}{\log D} \leq E[|C(x)|] \leq \frac{H(x)}{\log D} + 1$$

[Code known as Shannon code]

Notes about Shannon Code

- Not "optimal" for every X .

(Example in Cover & Thomas)

- But nearly optimal if $H(X)$ large

- So can apply to $\bar{X} = (X_1, \dots, X_n)$ & then

Inefficiency is very small ($\leq \frac{1}{n}$).

- But this introduces other complications.

Need to maintain encoding & decoding

table. Their lengths $\sim 2^{H(\bar{X})} = 2^{H(X) \cdot n}$.

- So if we can avoid appealing to n ,

it is good....

- Can do optimal "effectively".

Huffman Coding

$$C: \Omega \rightarrow \{0,1\}^*$$

$$\Omega = \{1, \dots, N\}$$

with p_1, \dots, p_N

Coding algorithm Huff (p_1, \dots, p_N);

if $N=2$ then $C(1)=0$; $C(2)=1$;

else if $N \geq 3$

[Sort so that $p_1 \geq \dots \geq p_N$;

$$C' \leftarrow \text{Huff}(p_1, \dots, p_{N-2}, p_{N-1} + p_N)$$

let $C(i) = C'(i) \quad \forall i \leq N-2$

$$\leftarrow C(N-1) = C'(N-1) \cdot 0$$

$$C(N) = C'(N) \cdot 1$$

Return C ;

]

Optimality of Huffman Code

Other coding schemes

Shannon-Elias-Fano

Arithmetic

} Described well
in Cover &
Thomas.