

دانشکده آموزش های الکترونیکی دانشگاه شیراز

سیستم های چند رسانه ای

فشرده سازی متن

وحید اعتمادی

بهار 94

Entropy

- Which sentence contains the most information?
 - Shiraz temperature in Tir is over 30.
 - Shiraz temperature in Tir is under 20.
- The first statement does not tell us anything new
- The second statement is improbable, giving us a great deal of information
- Statement with unpredictability factor conveys the most new information and surprise

Entropy

- $I(A)$: Amount of information that event A provides
- The amount of information associated with an event should be inversely related to the probability of the event.

$$\text{Definition: } I(A) = \log(1/p_A) = -\log(p_A)$$

- Note: with this definition, an event that is sure to happen ($p_A = 1$) provides $I(A) = 0$ bits of information

مقدمات تئوري اطلاعات

- آنترپي η از يك منبع اطلاعاتي $S = \{s_1, s_2, \dots, s_n\}$ برابر است با:

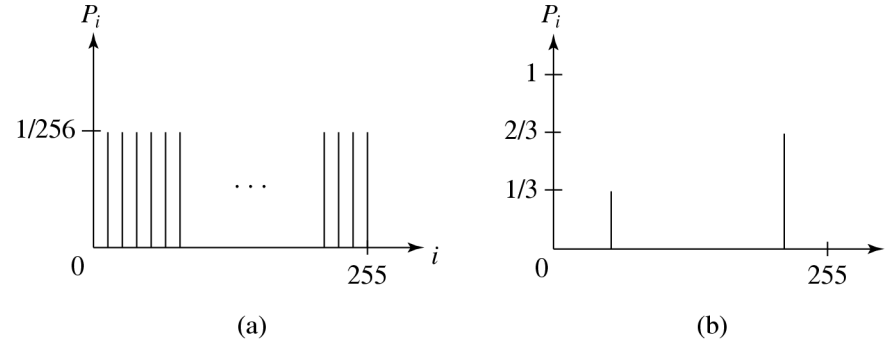
$$\eta = H(S) = \sum_{i=1}^n p_i \log_2 \frac{1}{p_i}$$

$$= - \sum_{i=1}^n p_i \log_2 p_i$$

p_i - احتمال رخداد s_i در S

$\log_2 \frac{1}{p_i}$ مشخص کننده ميزان اطلاعات موجود در s_i موجود است، بنابراین مي تواند تعيين کننده تعداد بيتهاي مورد نياز براي كد كردن s_i باشد.

مقدمات تئوري اطلاعات



هستوگرام سطح خاکستري دو تصوير

- (a) هستوگرام يك تصوير با شدت هاي توزيع سطح خاکستري يک نواخت و مساوي را نشان مي دهد، به ازاي هر i داريم:

$p_i = 1/256$ بنابراین آنروپي اين تصوير برابر است با:

$$\sum_{i=1}^{256} \frac{1}{256} \log_2 256 = 8 \quad \log_2 256 = 8$$

- (b) هستوگرام يك تصوير با دو مقدار ممکن را نشان مي دهد. آنروپي آن 0.92 است.

$$\frac{1}{3} \log_2 3 + \frac{2}{3} \log_2 \frac{3}{2}$$

کدینگ (VLC) Variable-Lenght

- الگوریتم شانون-فانو (روش بالا به پایین):

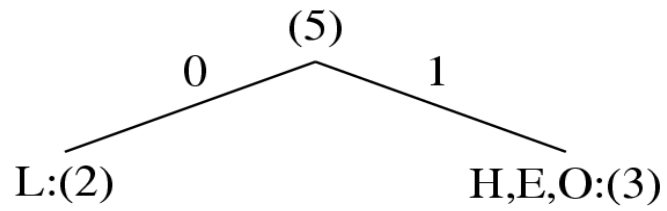
1. نمادها را بر اساس شمار فراوانی رخدادشان مرتب کنید.
2. نمادها را به صورت بازگشتی به دو قسمت تقسیم کنید، هر کدام با تعداد شماره‌های تقریباً برابر، تا جایی که هر قسمت فقط یک نماد داشته باشد.

- مثال: کدینگ کلمه HELLO

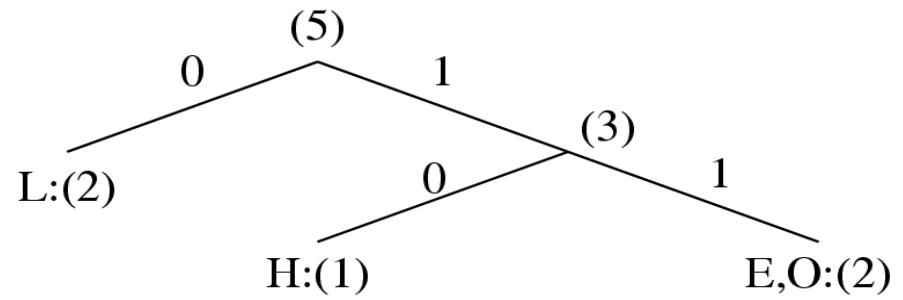
نماد	H	E	L	O
تعداد	1	1	2	1

شمار فراوانی نمادها در HELLO

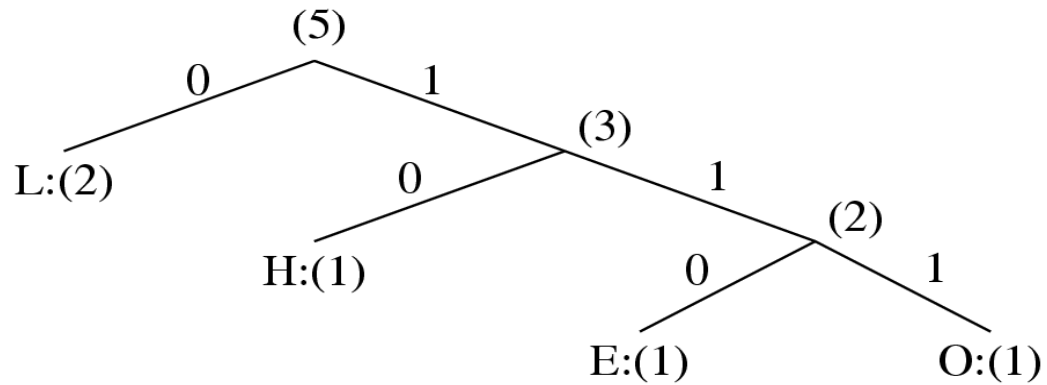
VLC



(a)



(b)



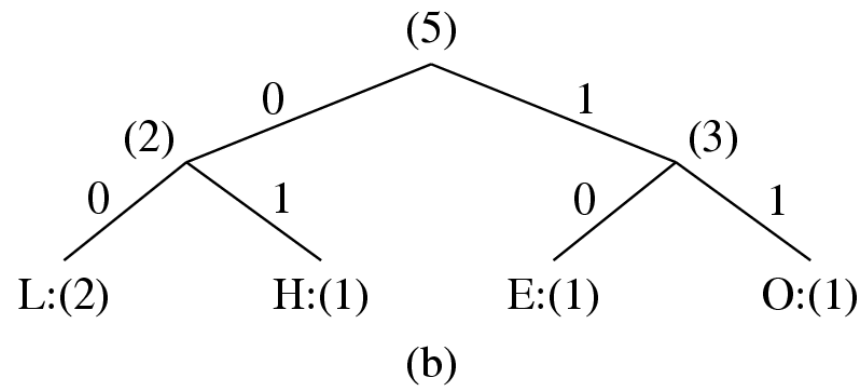
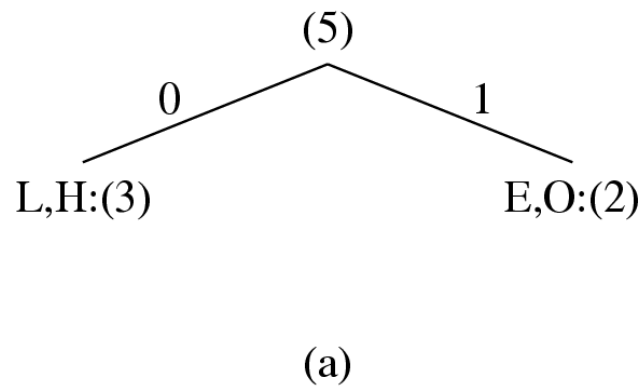
(c)

VLC

نتایج اجرای شانون-فانو بر روی HELLO

Symbol	Count	$\log_2 \frac{1}{p_i}$	Code	# of bits used
L	2	1.32	0	1
H	1	2.32	10	2
E	1	2.32	110	3
O	1	2.32	111	3
TOTAL # of bits:				10

VLC



• درخت کدینگ دیگری توسط شانون-فانو برای HELLO

VLC

• نتایج دیگر از اجرای شانون-فانو بر روی HELLO

Symbol	Count	$\log_2 \frac{1}{p_i}$	Code	# of bits used
L	2	1.32	00	4
H	1	2.32	01	2
E	1	2.32	10	2
O	1	2.32	11	2
TOTAL # of bits:				10

VLC

• الگوریتم کدینگ هافمن (Huffman) روند پایین به بالا:

1. تمام نمادها را بر اساس شمار فراوانی آنها در یک لیست مرتب کنید.

2. تا زمانی که فقط یک نماد باقی بماند، ادامه دهید:

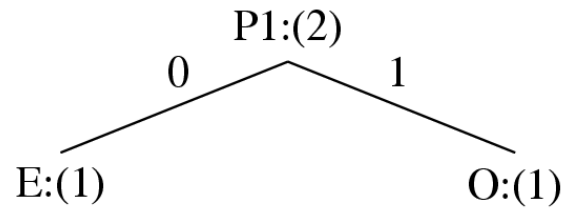
(1) دو نمادی که دارای کمترین فراوانی هستند را از لیست بردارید. یک زیردرخت هافمن ایجاد کنید که این دو نماد را به عنوان گره فرزند قرار داده و یک گره پدر ایجاد کنید.

(2) جمع تعداد فرزندان را به نود پدر اختصاص داده و در لیست به عنوان ترتیب قرارگیری وارد کنید.

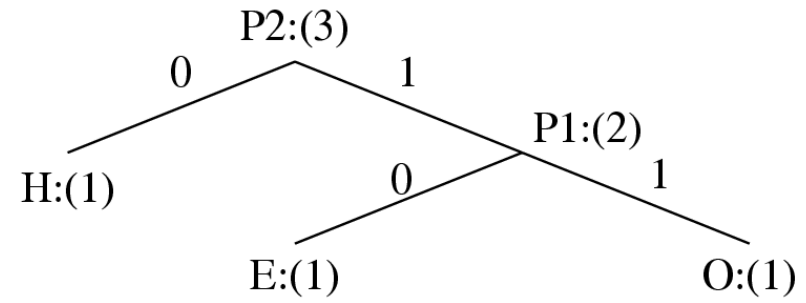
(3) فرزندان را از لیست حذف کنید.

(4) یک کلیدواژه برای هر برگ بر اساس مسیر از ریشه تعیین می شود.

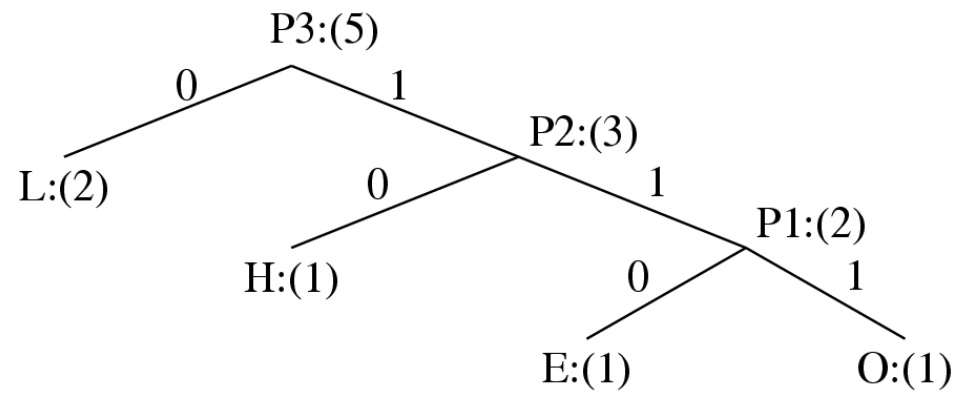
VLC



(a)



(b)



(c)

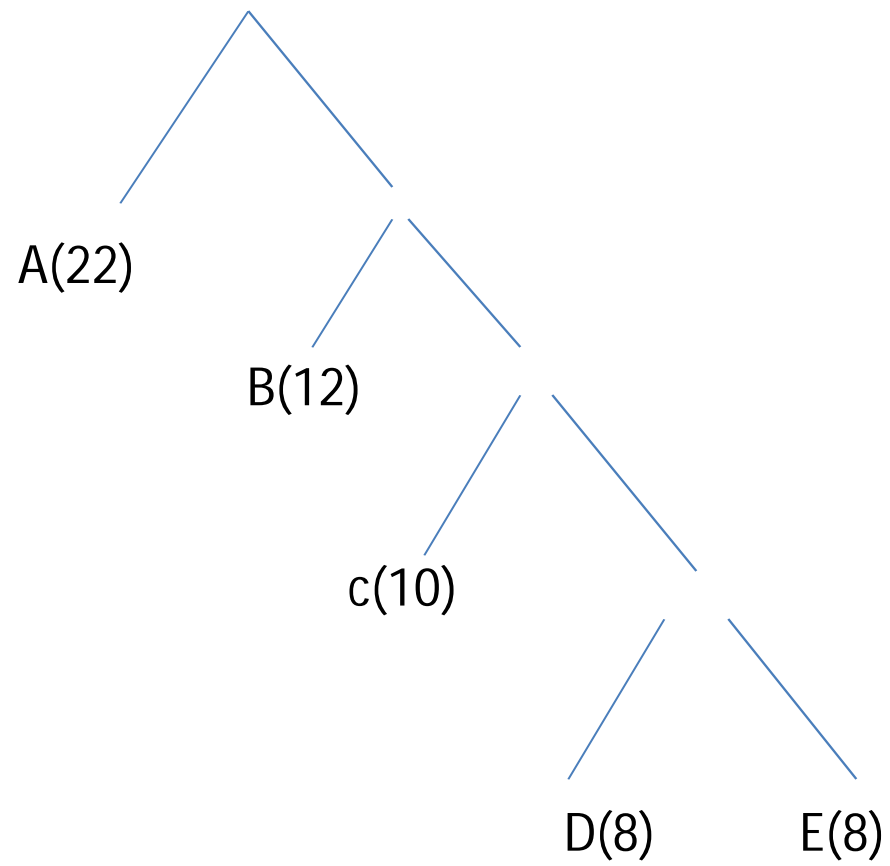
VLC

- نمادهای جدید $p1$ ، $p2$ ، $p3$ برای اشاره به نودهای پدر در درخت کدینگ هافمن تولید شده‌اند. محتویات موجود در لیست، در زیر مشخص شده است:

O(1)	E(1)	H(1)	L(2)	شروع:
	P1	H	L	(a): بعد از تکرار
	P2	L		(b): بعد از تکرار
		P3		(c): بعد از تکرار

سوال

- مجموعه سمبل‌های موجود در سیستم شامل پنج عنصر متفاوت A، B، C، D و E است. تعداد تکرار این سمبل‌ها در توالی ورودی به ترتیب 22، 12، 10، 8 و 8 مرتبه است.
- الف) با استفاده از روش Shannon-Fano کد هر یک از سمبل‌ها را مشخص نمایید. و طول میانگین کد برای این روش را حساب کنید.
- ب) با استفاده از روش Huffman کد هر یک از سمبل‌ها را مشخص نمایید. و طول میانگین کد برای این روش را حساب کنید.
- ج) انتروپی را برای این سیستم و مجموعه سمبل‌ها محاسبه کنید.



- Average # of bit: $22*1 + 12*2 + 10*3 + 4*8 + 4*8$
 $=140$ so:
 $140/60=2.33$

$$\text{Entropy}(S) = .36 * 1.47 + .2 * 2.32 + .17 * 2.55 + .13 * 2.94 + .13 * 2.94 = 2.19$$

Avg. # of bits \approx Entropy(S)

HW 1

- بخش ب و ج را محاسبه کنید.

خصوصیات کدینگ هافمن

- **1- خصوصیت منحصر به فرد پیشوندی:** هیچ کد هافمنی پیشوند کد هافمن دیگر نیست. این نکته مانع ایجاد هر گونه ابهامی در مرحله رمزگشایی می شود.

- **2- بهینگی:** حداقل افزونگی کد.

- دو نماد با کمترین فراوانی، دارای طول مساوی برای کد هافمنشان هستند و تنها در آخرین بیت تفاوت دارند.
- نمادهایی با فراوانی بیشتر دارای کد هافمن کوتاه تری نسبت به نمادهای با فراوانی کمتر هستند.
$$l < \eta + 1$$
- میانگین طول کد، برای یک منبع اطلاعات S حتماً از $\eta + 1$ کمتر است.