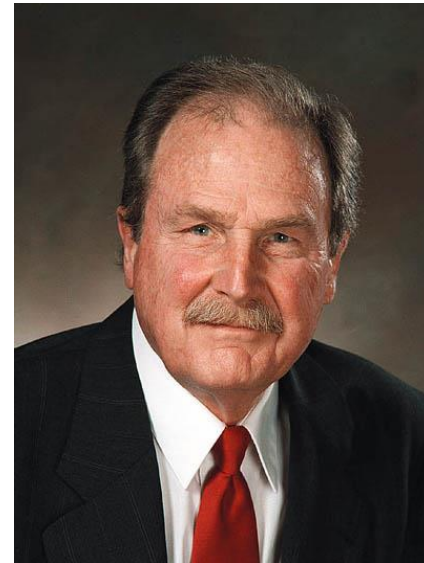


Huffman codes



David Huffman
(1925~1999)

現在使用的jpg檔案是失真的，有壓縮，壓縮的過程有使用這個huffman code

Coding

- The assignment of ***binary sequences*** to ***characters*** of an alphabet

alphabet	character / letter	code	codeword
	<i>a</i>		1
	–		001
	<i>b</i>		01100
	<i>f</i>		0100
	<i>n</i>		0111
	<i>r</i>		000
	<i>w</i>		01101
	<i>y</i>		0101

把左邊的字母轉換成0101組成的字串

Coding

- Fixed-length code

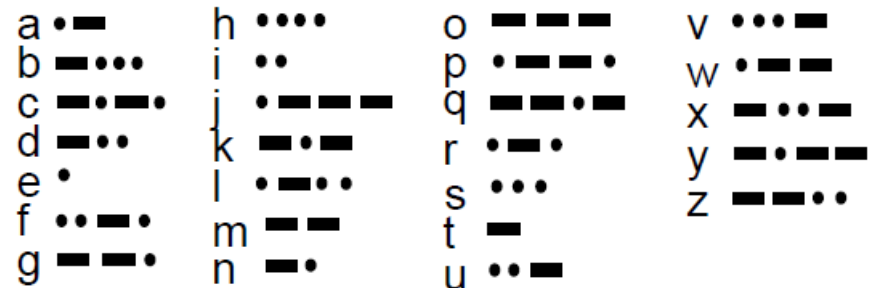
- Example: ASCII Code, Big5, UTF

- Variable-length code

- Example: Morse Code (1838)

早期是航海用的，用閃光或鳴笛來送出訊息。長度是2~4不等

Dec	Hx	Oct	Char	Dec	Hx	Oct	Html	Chr	Dec	Hx	Oct	Html	Chr	Dec	Hx	Oct	Html	Chr
0	0	000	NUL (null)	32	20	040	#32;	Space	64	40	100	#64;	@	96	60	140	#96;	^
1	1	001	SOH (start of heading)	33	21	041	#33;	!	65	41	101	#65;	A	97	61	141	#97;	a
2	2	002	STX (start of text)	34	22	042	#34;	"	66	42	102	#66;	B	98	62	142	#98;	b
3	3	003	ETX (end of text)	35	23	043	#35;	#	67	43	103	#67;	C	99	63	143	#99;	c
4	4	004	EOT (end of transmission)	36	24	044	#36;	\$	68	44	104	#68;	D	100	64	144	#100;	d
5	5	005	ENO (enquiry)	37	25	045	#37;	%	69	45	105	#69;	E	101	65	145	#101;	e
6	6	006	ACK (acknowledge)	38	26	046	#38;	&	70	46	106	#70;	F	102	66	146	#102;	f
7	7	007	BEL (bell)	39	27	047	#39;	'	71	47	107	#71;	G	103	67	147	#103;	g
8	8	010	BS (backspace)	40	28	050	#40;	(72	48	110	#72;	H	104	68	150	#104;	h
9	9	011	TAB (horizontal tab)	41	29	051	#41;)	73	49	111	#73;	I	105	69	151	#105;	i
10	A	012	LF (NL line feed, new line)	42	2A	052	#42;	*	74	4A	112	#74;	J	106	6A	152	#106;	j
11	B	013	VT (vertical tab)	43	2B	053	#43;	+	75	4B	113	#75;	K	107	6B	153	#107;	k
12	C	014	FF (NP form feed, new page)	44	2C	054	#44;	,	76	4C	114	#76;	L	108	6C	154	#108;	l
13	D	015	CR (carriage return)	45	2D	055	#45;	-	77	4D	115	#77;	M	109	6D	155	#109;	m
14	E	016	SO (shift out)	46	2E	056	#46;	.	78	4E	116	#78;	N	110	6E	156	#110;	n
15	F	017	SI (shift in)	47	2F	057	#47;	/	79	4F	117	#79;	O	111	6F	157	#111;	o
16	10	020	DLE (data link escape)	48	30	060	#48;	0	80	50	120	#80;	P	112	70	160	#112;	p
17	11	021	DC1 (device control 1)	49	31	061	#49;	1	81	51	121	#81;	Q	113	71	161	#113;	q
18	12	022	DC2 (device control 2)	50	32	062	#50;	2	82	52	122	#82;	R	114	72	162	#114;	r
19	13	023	DC3 (device control 3)	51	33	063	#51;	3	83	53	123	#83;	S	115	73	163	#115;	s
20	14	024	DC4 (device control 4)	52	34	064	#52;	4	84	54	124	#84;	T	116	74	164	#116;	t
21	15	025	NAK (negative acknowledge)	53	35	065	#53;	5	85	55	125	#85;	U	117	75	165	#117;	u
22	16	026	STN (synchronous idle)	54	36	066	#54;	6	86	56	126	#86;	V	118	76	166	#118;	v
23	17	027	ETB (end of trans. block)	55	37	067	#55;	7	87	57	127	#87;	W	119	77	167	#119;	w
24	18	030	CAN (cancel)	56	38	070	#56;	8	88	58	130	#88;	X	120	78	170	#120;	x
25	19	031	EM (end of medium)	57	39	071	#57;	9	89	59	131	#89;	Y	121	79	171	#121;	y
26	1A	032	SUB (substitute)	58	3A	072	#58;	:	90	5A	132	#90;	Z	122	7A	172	#122;	z
27	1B	033	ESC (escape)	59	3B	073	#59;	;	91	5B	133	#91;	[123	7B	173	#123;	{
28	1C	034	FS (file separator)	60	3C	074	#60;	<	92	5C	134	#92;	\	124	7C	174	#124;	
29	1D	035	GS (group separator)	61	3D	075	#61;	=	93	5D	135	#93;]	125	7D	175	#125;	}
30	1E	036	RS (record separator)	62	3E	076	#62;	>	94	5E	136	#94;	^	126	7E	176	#126;	~
31	1F	037	US (unit separator)	63	3F	077	#63;	?	95	5F	137	#95;	_	127	7F	177	#127;	DEL



Coding

站在儲存或傳送資料的角度，會覺得檔案越小越好

- Which one is more effective in terms of *compression*?

	a	b	c	d	e	f
Frequency (in thousands)	45	13	12	16	9	5
Fixed length	000	001	010	011	100	101
Variable length	0	101	100	111	1101	1100

3：每個長度是3碼

Fixed-length code: $3 * (45 + 13 + 12 + 16 + 9 + 5) = 300$ k bits

用fixed length來存會需要300k bits



Variable-length code: $1 * 45 + 3 * (13 + 12 + 16) + 4 * (9 + 5) = 224$ k bits

Assign *short* codeword to *frequent* characters!

所以站在壓縮的角度variable length是比較好的

Coding

Characters	Code 1	Code 2	Code 3
A	0	0	0
C	0	1	10
T	1	00	110
G	10	11	111

Ambiguous

Not uniquely
decodable

100 → C T
→ C A A

code 2的編碼法會解出兩種可能，所以不行

Prefix code

在壓縮上是很有效的編碼法

Prefix code

- No codeword is a prefix of some other codeword.

- Example:

prefix code : 找不到任何code word是別的code word的prefix的編碼法

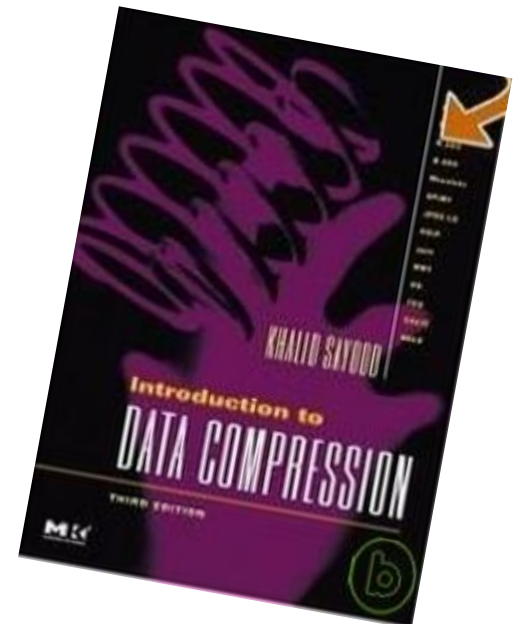
– Which one is prefix code? **Code3 !**

Letters	Code 2	Code 3	Code 4
a_1	0	0	0
a_2	1	10	01
a_3	00	110	011
a_4	11	111	0111



Prefix code

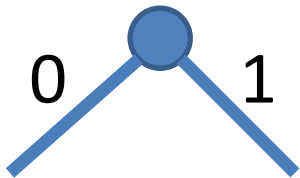
- Prefix code can always achieve the ***optimal*** data compress among any code!



Prefix (Optimal) code

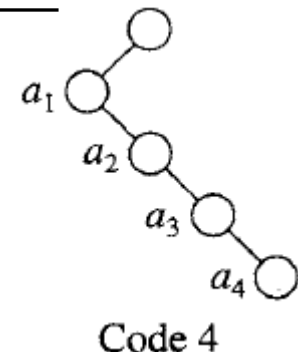
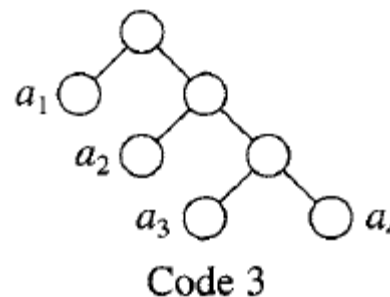
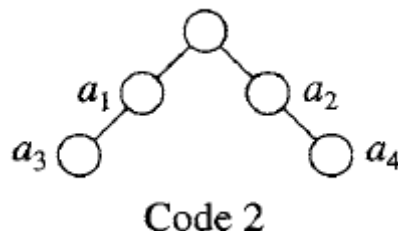
Q. 如何判斷一組code是不是prefix code呢？

- Can always be represented by a **binary tree**
- **Characters can only appear at leaf nodes**



Letters	Code 2	Code 3	Code 4
a_1	0	0	0
a_2	1	10	01
a_3	00	110	011
a_4	11	111	0111

A. 讓這一組code長出一顆樹，最後再去檢查是不是每個letters都在樹葉端上



code 3 所有的letters都只會出現在樹葉端

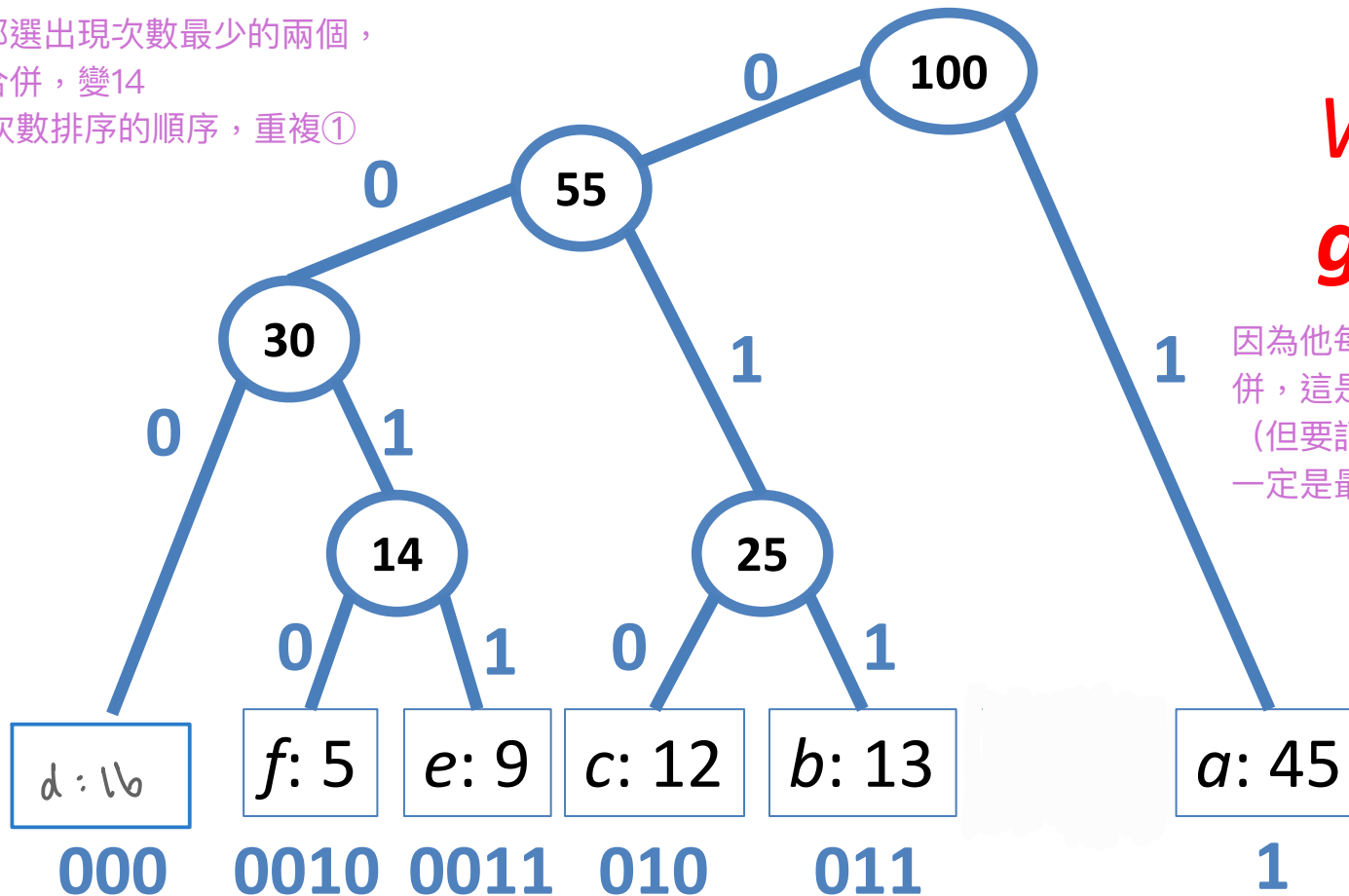
code 4都是別人的prefix

Huffman code

每個字母出現的次數（千次）

- $A = \{f: 5, e: 9, c: 12, b: 13, d: 16, a: 45\}$

1. 每次都選出現次數最少的兩個，
然後合併，變14
2. 按照次數排序的順序，重複①



Why is it greedy?

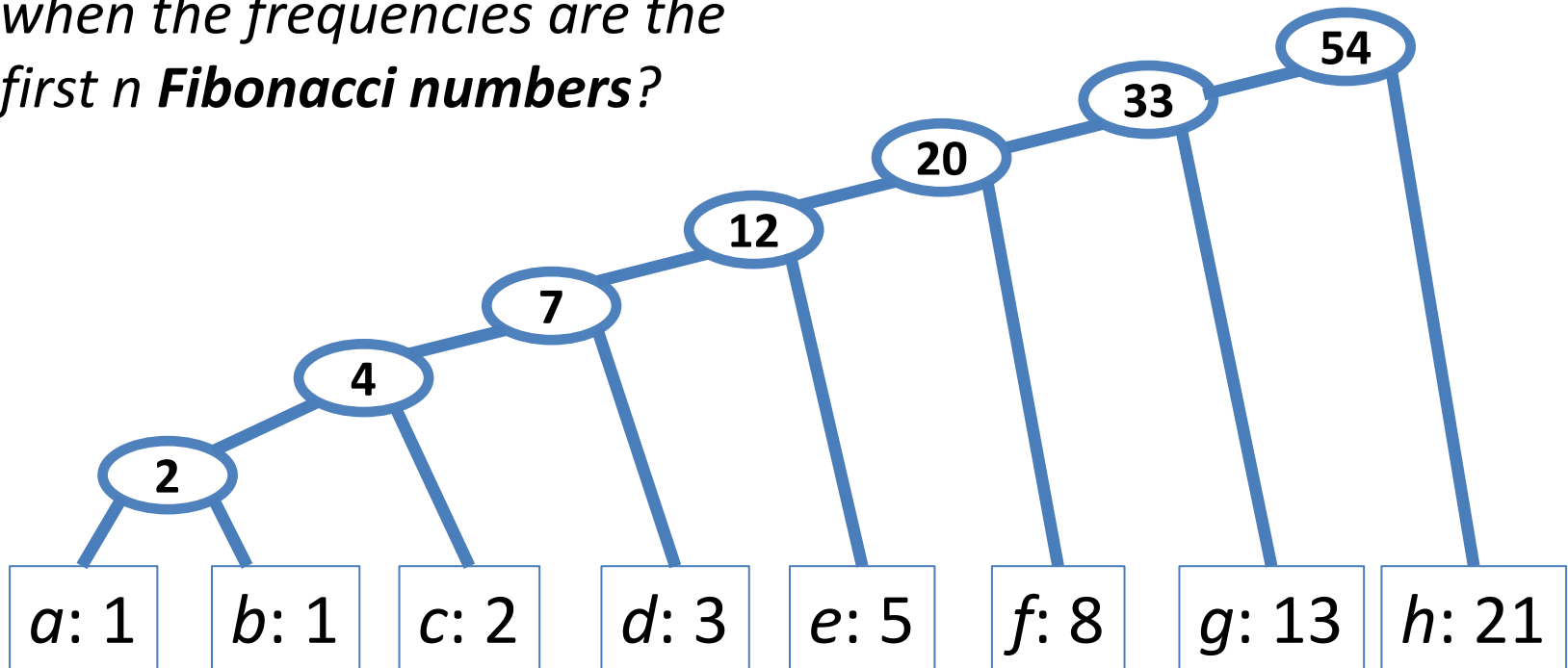
因為他每次都把兩個最小的合併，這是他的greedy choice
(但要記得貪婪法並沒有保證一定是最佳解)

有先排序過

Practice

- $A = \{a: 1, b: 1, c: 2, d: 3, e: 5, f: 8, g: 13, h: 21\}$

*Can you find the optimal code when the frequencies are the first n **Fibonacci numbers**?*



$$3 * 0.16 + 4 * 0.05 + 4 * 0.09 + 3 * 0.12 + 3 * 0.13 + 1 * 0.45 = 2.24$$

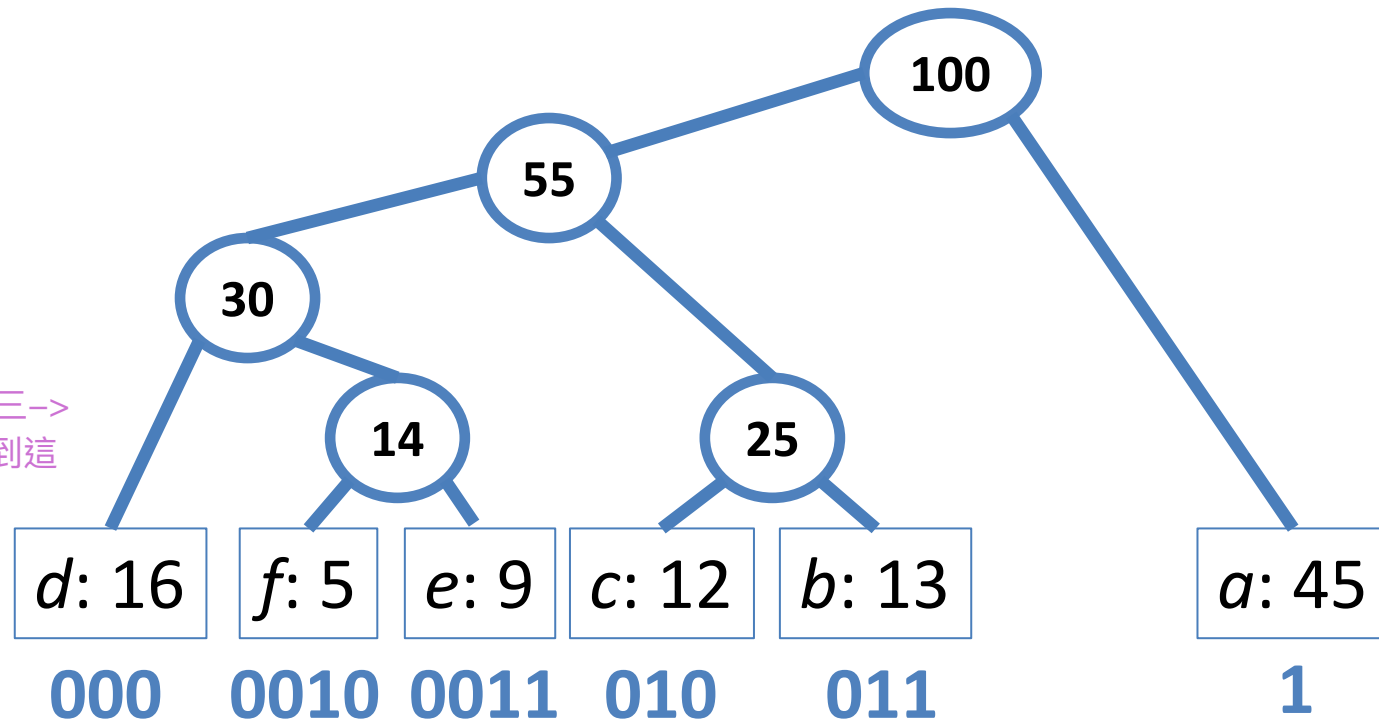
Cost (average bit rate)

平均來講每個字會需要2.24的bit，可以把文件傳出去



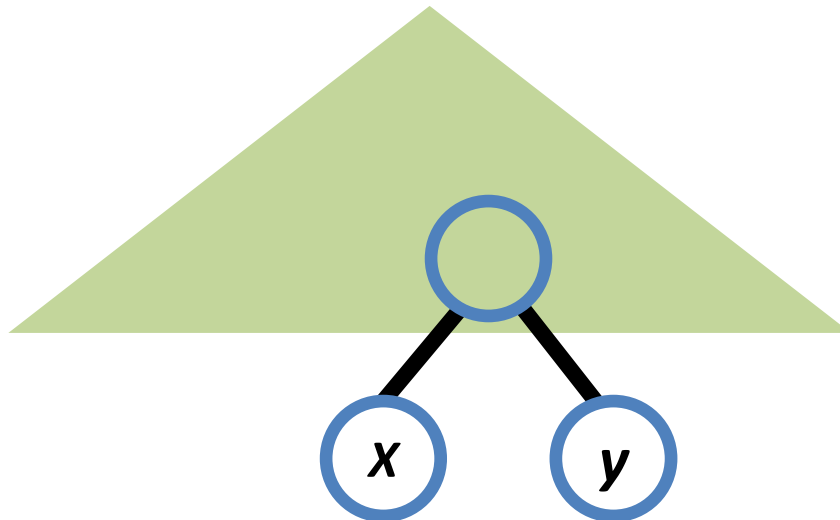
為什麼用最小的兩個字合併會是最佳解？

- The cost of a tree T : the sum of weighted depth of leaves
- An optimal tree has the *minimal* cost.

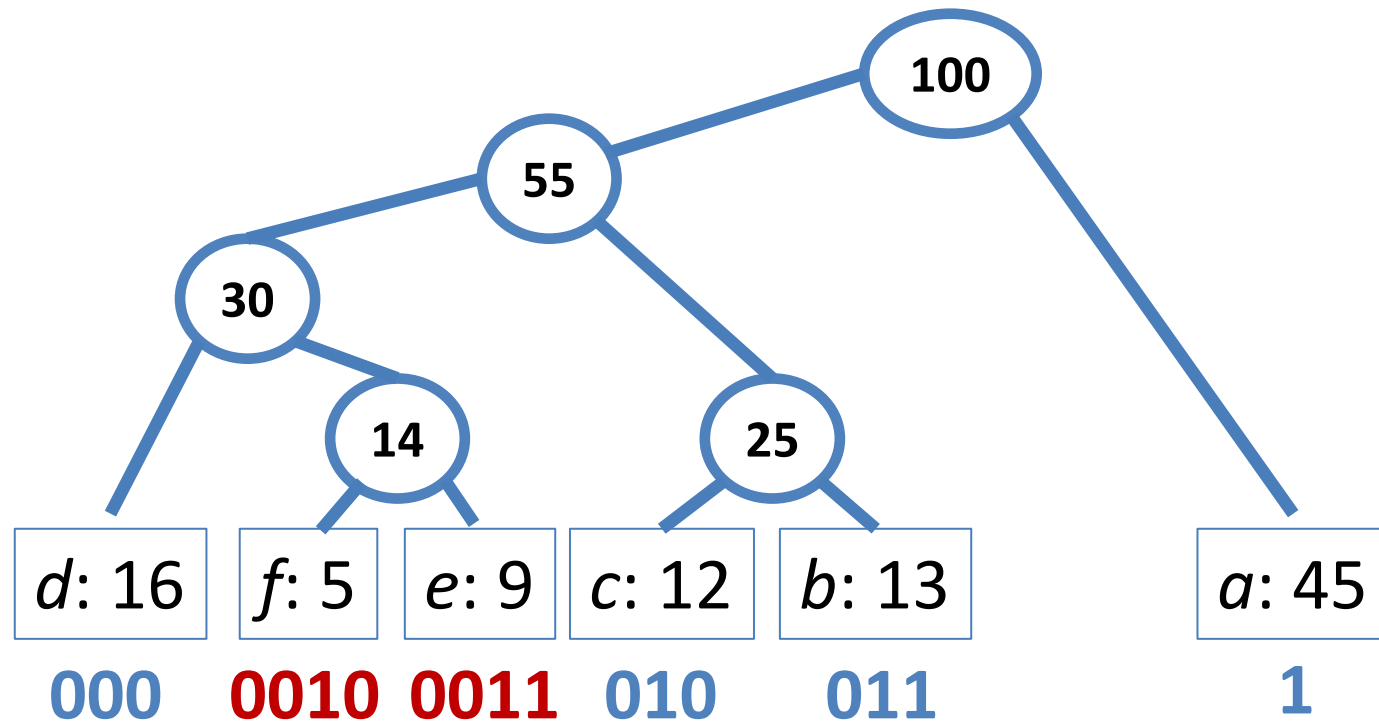


Correctness

- Claim: Let x and y be the lowest frequencies. There exists an optimal prefix code in which the code-words for x and y have the same length and differ only in the last bit.



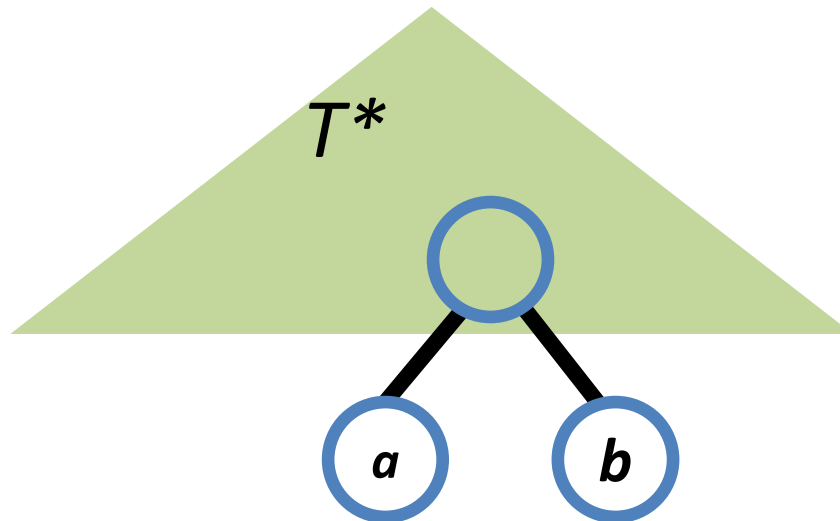
- Claim: Let x and y be the lowest frequencies. There exists an optimal prefix code in which the code-words for x and y have the same length and differ only in the last bit.



Correctness

- Proof: Let T^* be an optimal tree. Let a and b be two sibling leaves that are deepest in T^* .

假設 T^* 是optimal tree， a, b 是最深的兩個節點



x, y : the lowest frequencies
 a, b : the deepest nodes in T^*

Correctness

- Proof (continued):

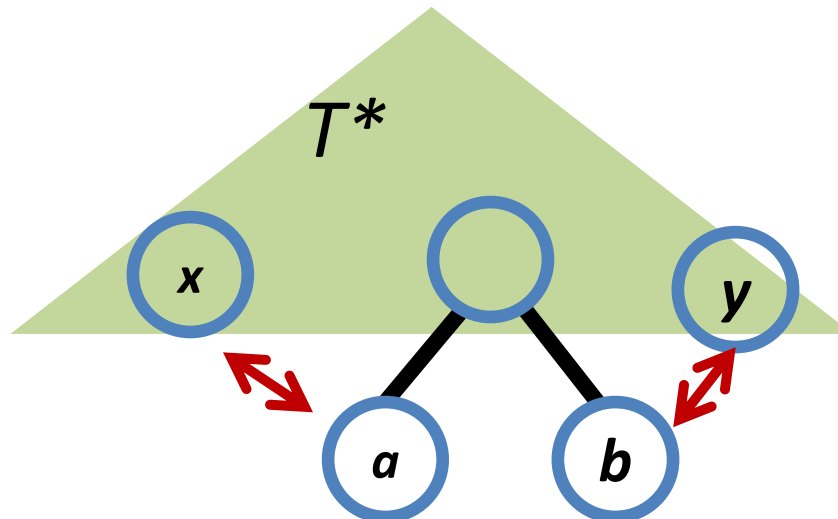
- Create T'' by swapping a and x , b and y

- $\text{cost}(T^*) \geq \text{cost}(T'')$ why?

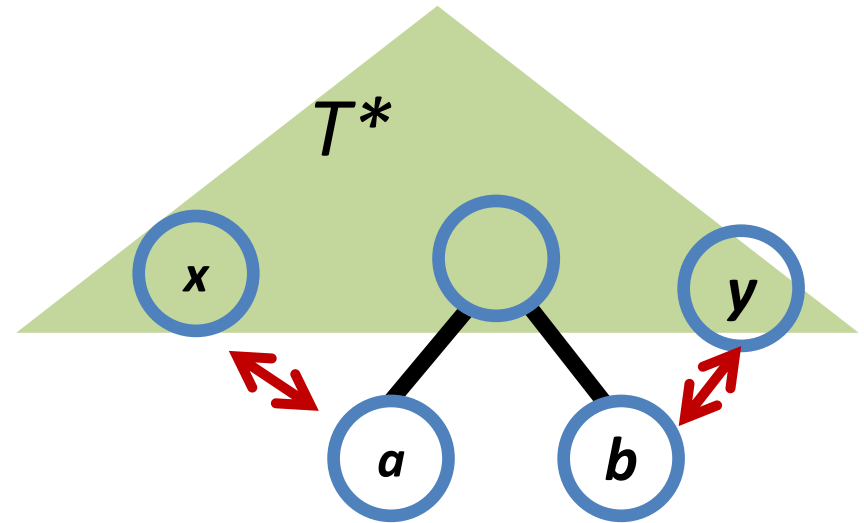
ab跟xy換位子後cost會往下掉

Q. why?

- Since T^* is optimal, which implies $\text{cost}(T^*) = \text{cost}(T'')$



x, y : the lowest frequencies
 a, b : the deepest nodes in T^*



$$\text{cost}(T^*) - \text{cost}(T'')$$

$$= \sum_{c \in C} c.\text{freq} \cdot d_{T^*}(c) - \sum_{c \in C} c.\text{freq} \cdot d_{T''}(c)$$

每一個樹的節點的頻率*深度

$$= (a.\text{freq} - x.\text{freq})(d_{T^*}(a) - d_{T^*}(x)) \\ + (b.\text{freq} - y.\text{freq})(d_{T^*}(b) - d_{T^*}(y))$$

p. 434

這四個()都會大於等於0。

因為我們已經是知道xy是最低頻率了，也知道ab的深度是最深的。（左上角有寫）

$$\geq 0$$

Conclusion

- An optimal solution for the coding problem can be found by a greedy algorithm.