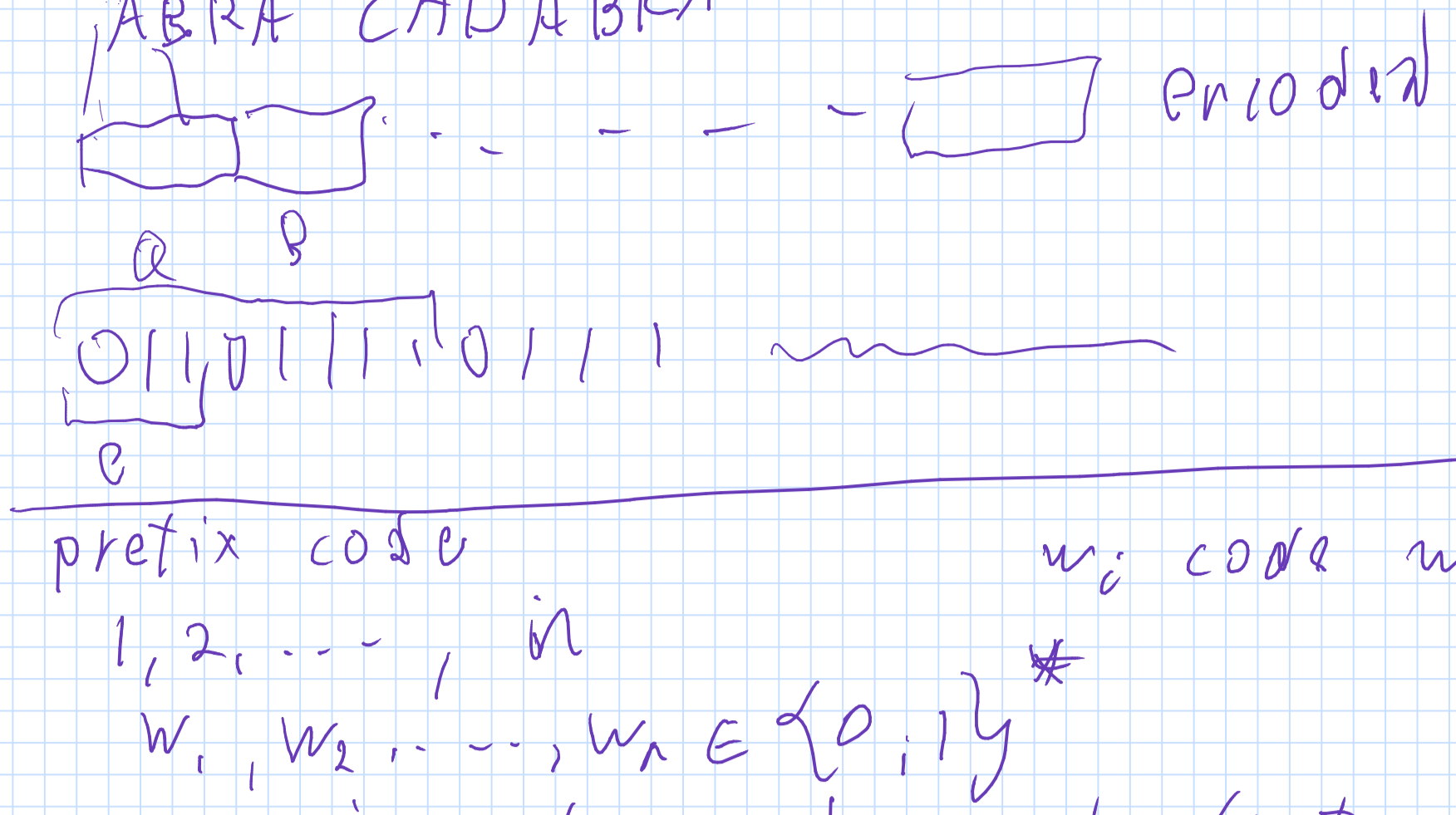


3/4/21

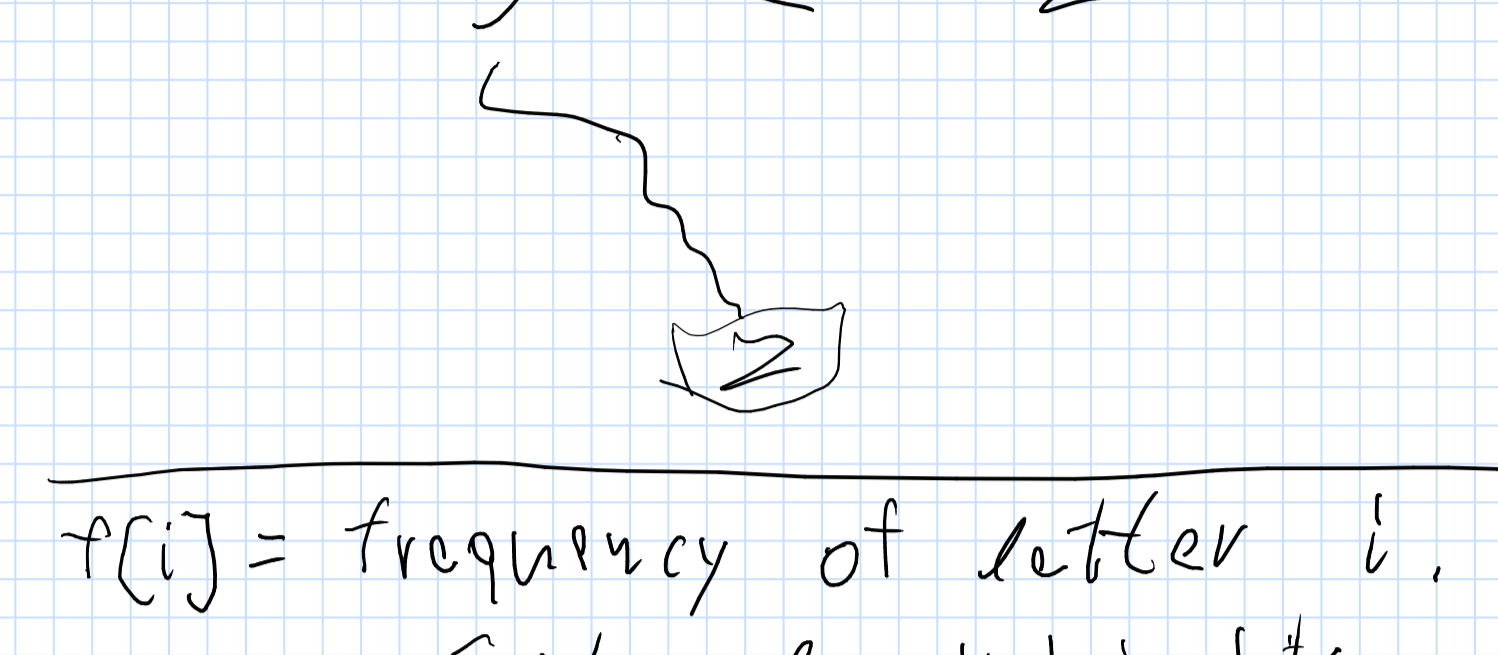
Huffman's encoding
 prefix codes binary
 ASCII 8 bits
 A = 64

English
 E, A > 10%
 X, Z, Q 1%

Z: 215
 E: 74809
 A: 48,000

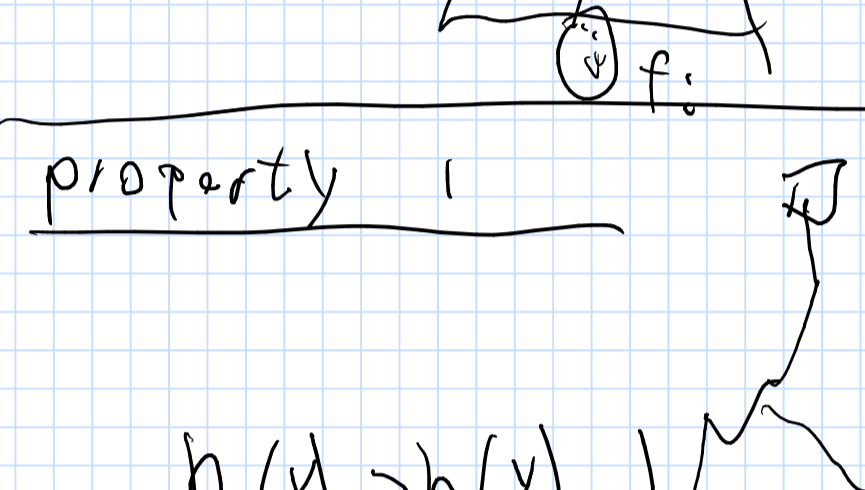


prefix code w_i code word
 $1, 2, \dots, n$
 $w_1, w_2, \dots, w_n \in \{0, 1\}^*$
 there is no code word that is prefix of another code word. } prefix free code



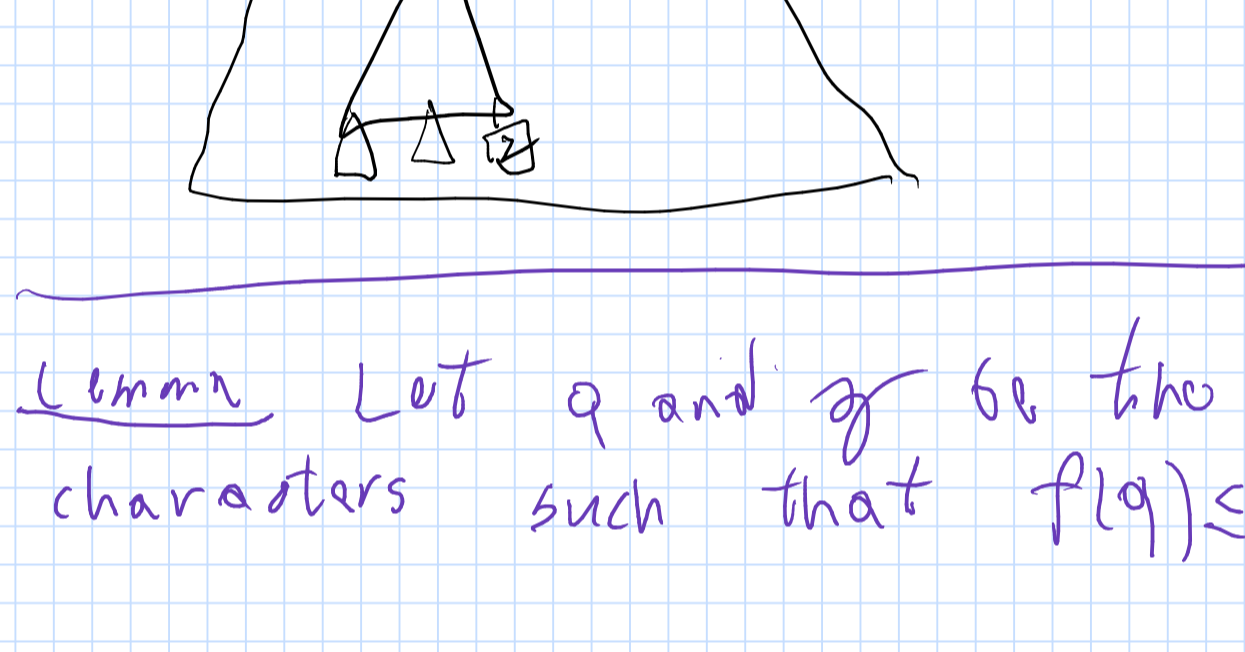
$f(i)$ = frequency of letter i .
 C : code $\begin{cases} w_1 & l_1 = |w_1| \text{ in bits} \\ w_2 & l_2 \\ \vdots & \vdots \\ w_n & l_n \end{cases}$ $F = \sum f_i \cdot l_i$

size compressed text $(F) \equiv \text{cost}(F)$
 $= \sum_{i=1}^n (w_i \cdot f_i)$
 compute prefix free code that min cost(F)

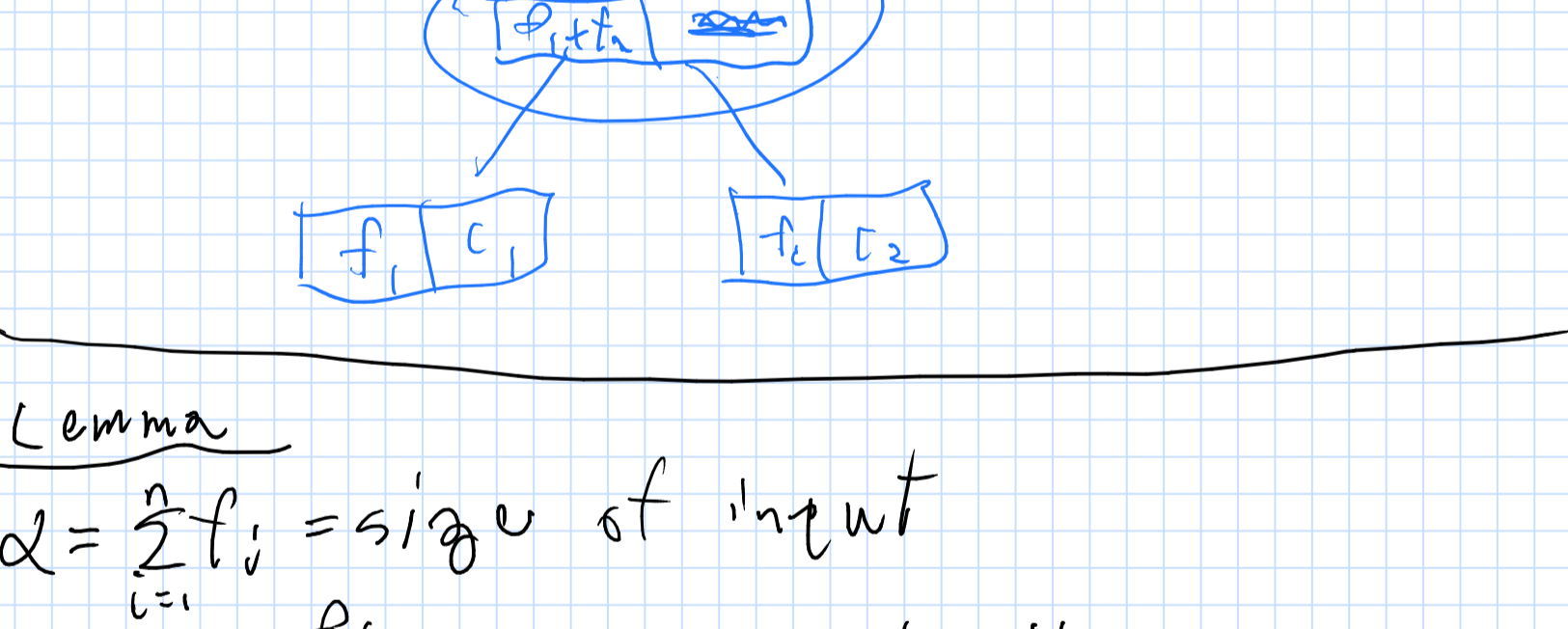
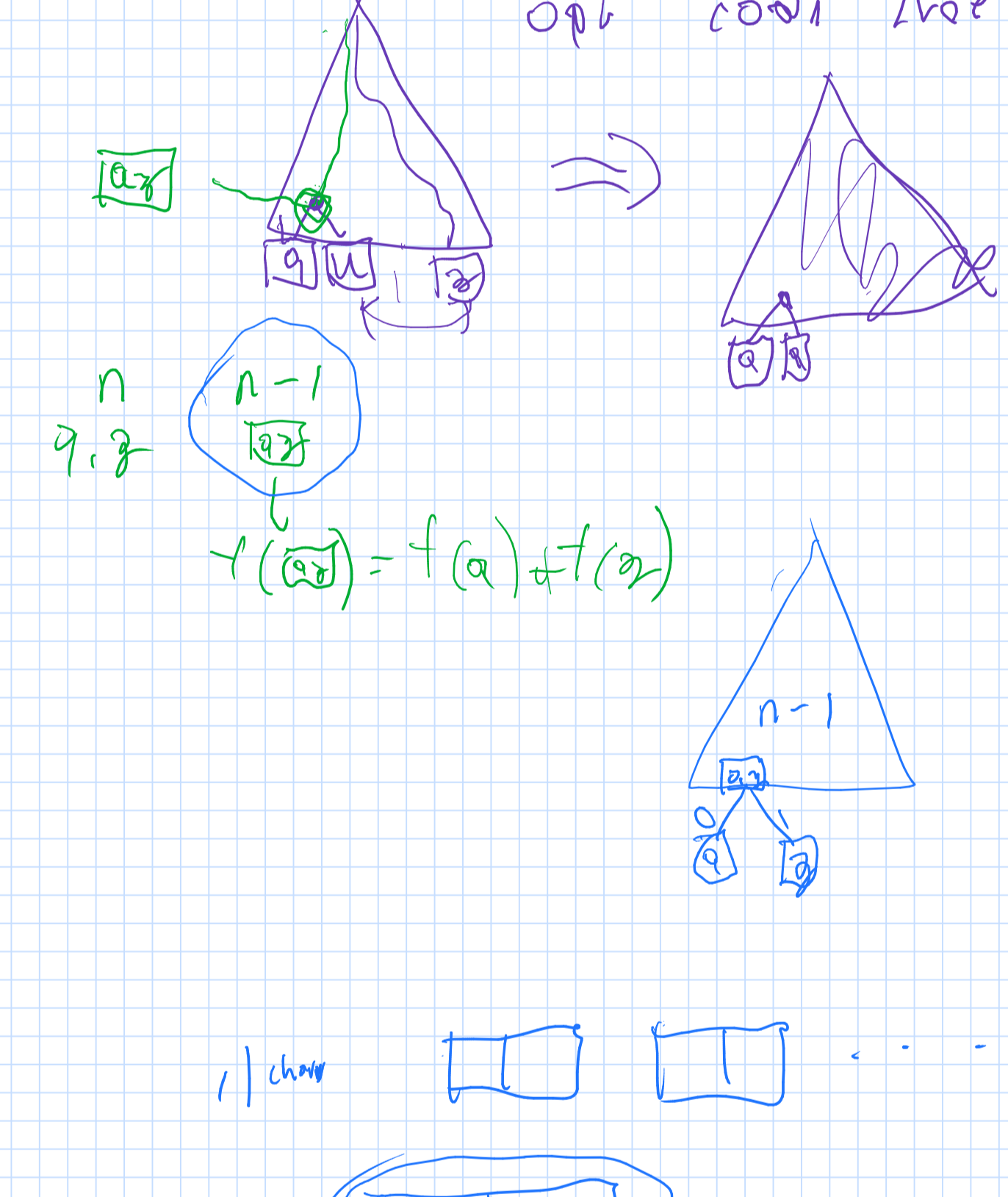


property 1
 $h(x) > h(y)$
 optimal code tree.
 $\Rightarrow f(x) < f(y)$
 if false $f(x) > f(y)$
 $f(x) \cdot h(x) > f(y) \cdot h(y)$
 $f(x)h(y) + f(y)h(x)$
 $h(x) > h(y)$
 $f(x) > f(y)$

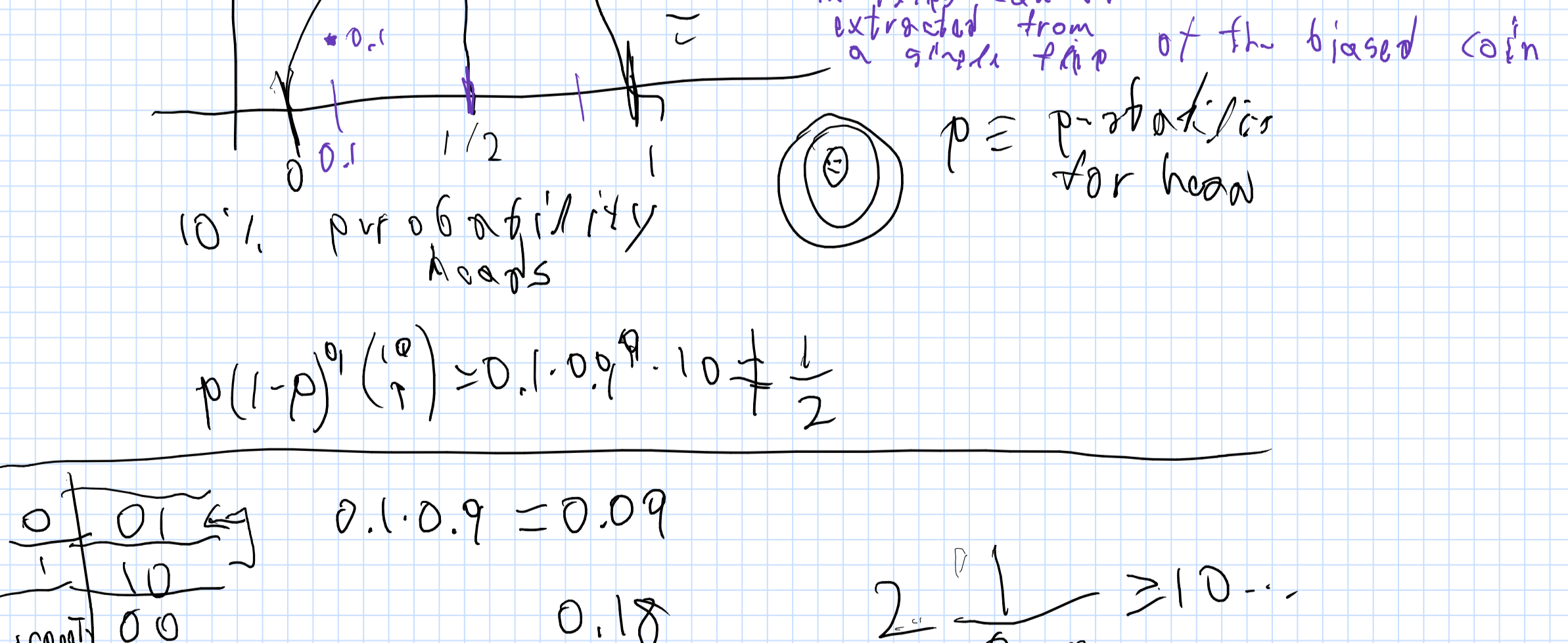
proved:
 claim In optimal as
 if $f(x) > f(y) \Rightarrow h(x) \leq h(y)$.



Lemma Let q and z be the two characters such that $f(q) \leq f(z) \leq$ frequency of any other character
 $\Rightarrow \exists$ optimal code tree T such that
 - q and z are siblings
 - q is the deepest leaf.
 - z is also deepest leaf.



Lemma
 $\alpha = \sum_{i=1}^n f_i$ = size of input
 $p_i = \frac{f_i}{\alpha}$ "probability base 2
 $H(p_1, p_2, \dots, p_n) = \sum_{i=1}^n p_i \log \frac{1}{p_i}$ Entropy
 $= \sum_{i=1}^n -p_i \log p_i$
 $H(p) = p \log \frac{1}{p} + (1-p) \log \frac{1}{1-p}$ binary entropy
 $H(p) = H(1-p)$ symmetric



0	01	0.1 * 0.9 = 0.09
1	10	0.18
repeat	00	
	11	

$2 \cdot \frac{1}{0.18} \geq 10 \dots$

Lemma
 Given n characters with probabilities p_1, p_2, \dots, p_n such that
 p_i is a power of 2
 $p_i = 1/2^{t_i}$

Then, the average cost of encoding the input using Huffman encoding is the entropy:
 $E[\text{expected len code}] = \sum_{i=1}^n p_i \log \frac{1}{p_i} = \sum_{i=1}^n p_i t_i$

\Leftrightarrow the i th character with probability $p_i = \frac{1}{2^{t_i}}$ has a word of length t_i .

