

01 Armenian Alphabet Entropy

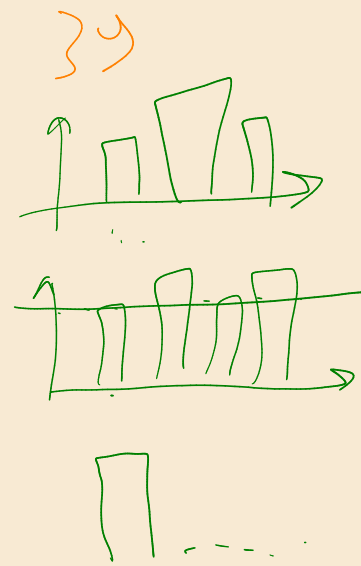
The Armenian alphabet has **39 letters** (Mashtots' original 36 + Ռ, Օ, Ֆ).

- a) Compute the **uniform entropy**: how many bits per letter if every letter were equally likely?
- b) English's uniform bound is $\log_2 26 \approx 4.70$ bits/letter, yet real English text has empirical entropy around 1.3 bits/letter (Shannon, 1951). Explain *why* the empirical entropy is so much smaller than the uniform bound. Name the **two** distinct effects.
- c) (Python) Estimate the *real* per-letter entropy of Armenian. Use the provided file [assets/panir_hy.txt](#) (the «Պանիր» Wikipedia article), or download fresh text yourself with the snippet below. Keep only Armenian letters, lowercase them, and compute the empirical letter-frequency distribution and its entropy.

```
# Option A - download the provided text file:
import urllib.request
url = ("https://raw.githubusercontent.com/HaykTarkhanyan/"
      "python_math_ml_course/main/math/assets/panir_hy.txt")
text = urllib.request.urlopen(url).read().decode("utf-8")

# Option B - download any Armenian article's plain text yourself:
import urllib.request, urllib.parse, json
title = "Պանիր" # try any article title
url = "https://hy.wikipedia.org/w/api.php?" + urllib.parse.urlencode({
    "action": "query", "prop": "extracts", "explaintext": 1,
    "titles": title, "format": "json", "redirects": 1,
})
req = urllib.request.Request(url, headers={"User-Agent": "course/1.0"})
page = next(iter(json.load(urllib.request.urlopen(req))["query"]["pages"].values()))
text = page["extract"]
```

- d) Compare your empirical entropy to the uniform bound from (a). By roughly how many bits per letter does the frequency skew save you?
- e) A 100-page Armenian e-book has roughly 200,000 characters. Estimate the entropy **lower bound** on its size (KB). Then explain the catch: you'll get ≈ 4.5 bits/letter, but a real code must assign a **whole** number of bits to each letter — so how do you actually get close to the entropy?



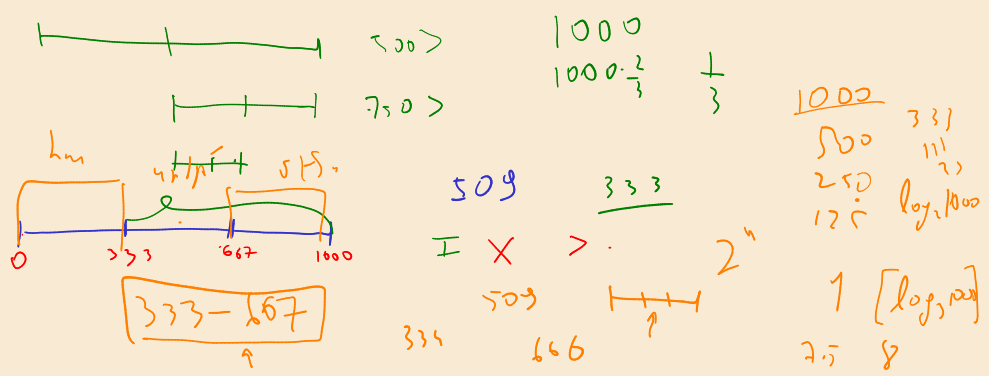
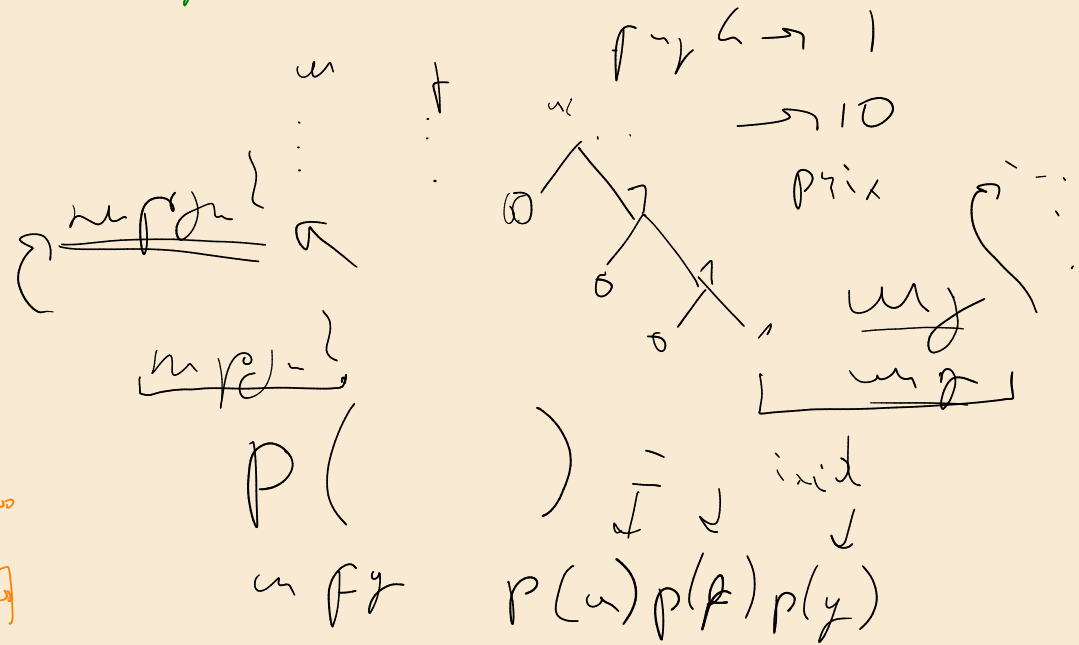
u 00000
k
y
}

unif

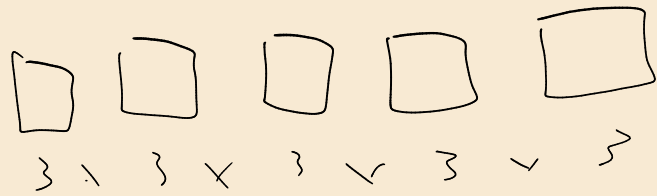
$$p(u) = p(y) = \frac{1}{39}$$

$$-\sum_{i=1}^{39} \frac{1}{39} \cdot \log \frac{1}{39} = -\log \frac{1}{39} = \log 39 = 3.2 - 64$$

$p(u) \neq p(y)$ 7.3



2.

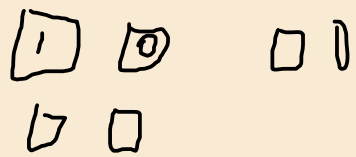
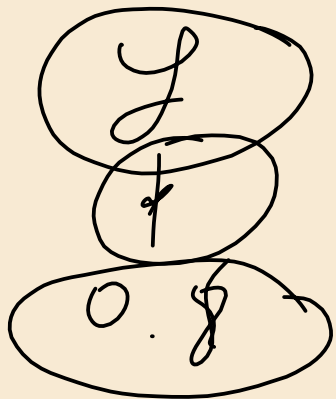


$$3^5 = 243$$

$$\rightarrow \sum_{i=1}^5 \frac{1}{243} \log_2 243 \approx 8 \quad \boxed{8}$$

$$\log_2 2315 \approx 11$$

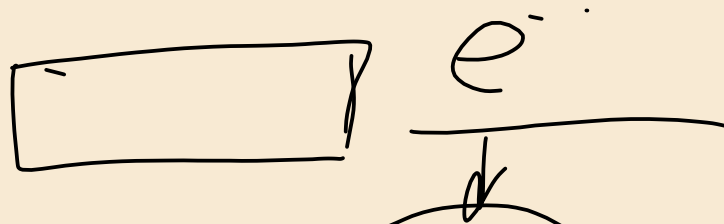
X



$$m \sim \frac{1}{n}$$

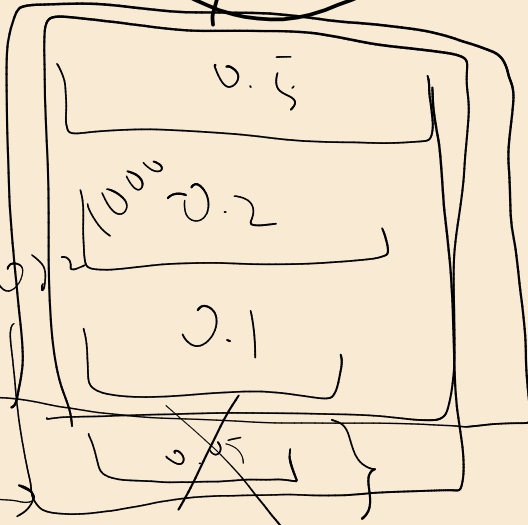
θ_1





t
to 32

80%
 $h(1000 \times 0.4) =$
 $\Rightarrow h(400) =$



Top



t
0.8

