

6 week Final Project lets build Bayesian classifier

Learning Objectives:

- 6.1 Know definition of multinomial distribution.
- 6.2 Apply notions of conditional probability, random variables and their distributions to design classifier algorithm.
- 6.3 Apply Bayes' rule to select the most plausible hypothesis.

Оглавление:

1. Problem statement and approach
2. Mathematical model & multinomial distribution
3. Explicit formula of multinomial distribution
4. Language detection: formalization & maximum likelihood principle
5. Bayesian inference and interpretation of results
6. SGA: Final proof

Problem statement and approach

Final project: language classification

Problem statement

- We have a text in unknown language
- This text uses Latin alphabet
- We assume that the text is written using a language from some fixed set, e.g. {English, German, French, Spanish}
- Our solution should be fully automatic and robust
- Length of the text is arbitrary
- It is possible that the text is too small to detect its language confidently — in this case our algorithm should report that

Final project will be devoted to development of algorithm of text classifications. Our problem: we have a text that is written on some unknown language for us. We have to find the algorithm which analyse this text and say to which language it belongs. If there are slight mistakes in the text like typos or unusual language(slang), it should anyway detect the language correctly at least as far as it's possible.

Language detection: possible approaches

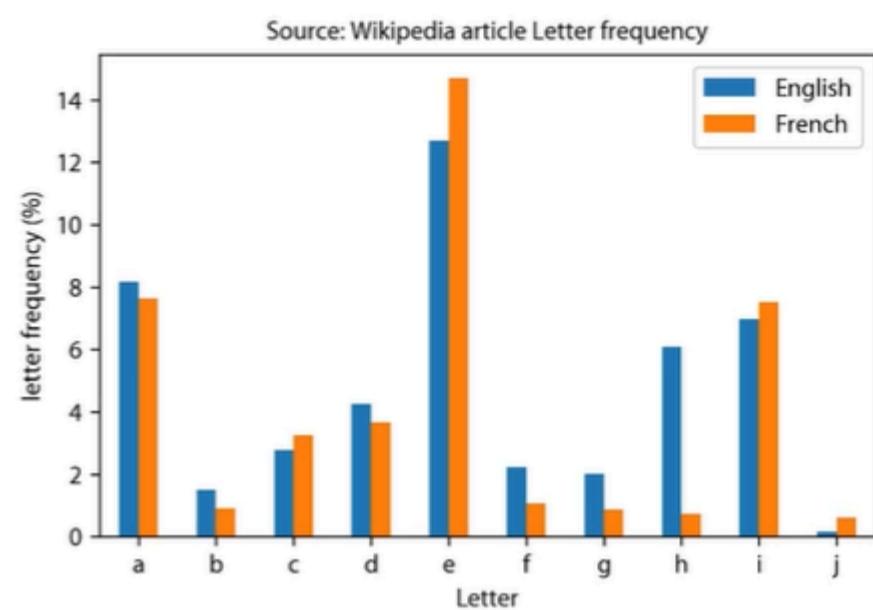
- We can leverage our knowledge about languages and make rule-based algorithm.
- For example, if there is word “**the**” in the text, it’s probably in English.
- This is error-prone: for example, “**The**” can be used as an alternative transcription of Chinese name **Zheng** or possible spelling of word **Te** (Tea) in Norwegian. In this case, our rule “The ⇒ English” will fail.

It is also not scalable: if we have new language, we should add new rules manually.

Language detection: statistical approach

- Different languages have different statistical properties.
- The simplest characteristics is letter frequencies.

Letter frequencies



We can see that we can meet h letter more frequently in English than in French.

Our approach

- Find frequencies of letters in all languages that we consider using some large corpora.
- Find frequencies of letters in our text.
- Compare our frequencies with language's frequencies.

Our approach: pro and contra

- (+) It's fully automatic. To add a new language, we just have to add its frequency table.
- (+) It's robust. For example, if a text contains typos or slang, it will not affect letter frequencies too much.
- (-) If our text is too small, letter frequencies can become unreliable.
- (?) We have to define what does it mean for frequency tables to be "similar" to each other.

Use relative frequency which can't be affected by large length of text.

Our plan

- Introduce probabilistic model that is based on *multinomial distribution* to answer the following question: for a given language, what is the probability to get a text with given frequency table? This probability is called *likelihood*.
- Use this model to estimate each language's likelihood for a given text.
- Take into account relative popularity of different languages using Bayes' rule.
- Write a code that will calculate likelihoods and so-called *posterior probabilities* of each language.

Mathematical model & multinomial distribution

In our language detection problem, we will use frequencies of characters in the text to determine the language of the text. To make our reasoning rigorous, we have to consider some mathematical model of this distribution. Under some assumption this distribution can be modelled by so called multinomial distribution. Let's discuss now:

In multinomial distribution, we assume that we have a kind of dice, but not actual dice, but some kind of generalized dice with K sides or K faces. And we assume that this dice can be non-symmetric. Probabilities of each face can be not equal to each other, but are given by some distribution.

Multinomial distribution

Dice with K faces
enumerated by numbers 1, ..., K

$$P(\text{get face } j) = p_j, \quad j=1, \dots, K$$
$$p_1 + \dots + p_K = 1 \quad p_j \geq 0$$

Roll our dice n times

Let $X_j = \text{"count of face } j"$

$$(X_1, \dots, X_K) \sim \text{Multinomial}(K, n, (p_1, \dots, p_K))$$

1. In this case multinomial distribution becomes binomial distribution

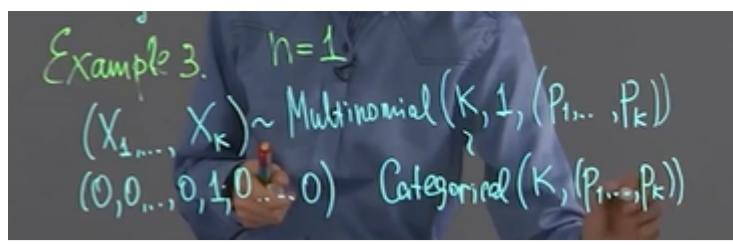
Example 1. $K=2$

$$(X_1, X_2) \sim \text{Multinomial}(2, n, (p, 1-p))$$
$$\Rightarrow X_1 \sim \text{Binomial}(n, p)$$

Example 2. $K=6$ $p_1 = \dots = p_6 = \frac{1}{6}$

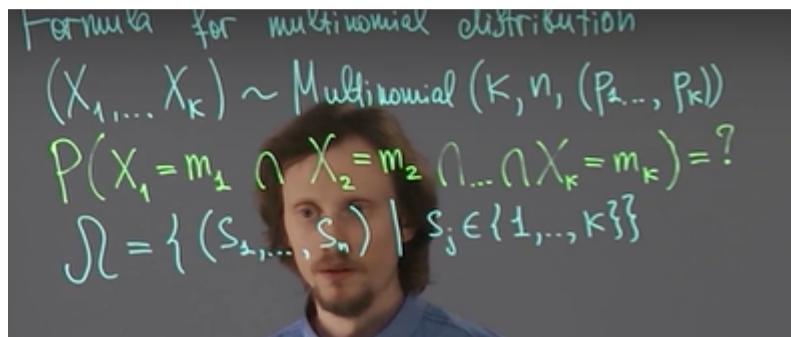
$$(X_1, \dots, X_6) \sim \text{Multinomial}(6, n, (\frac{1}{6}, \dots, \frac{1}{6}))$$

$X_j = \text{count of } j \text{ points}$

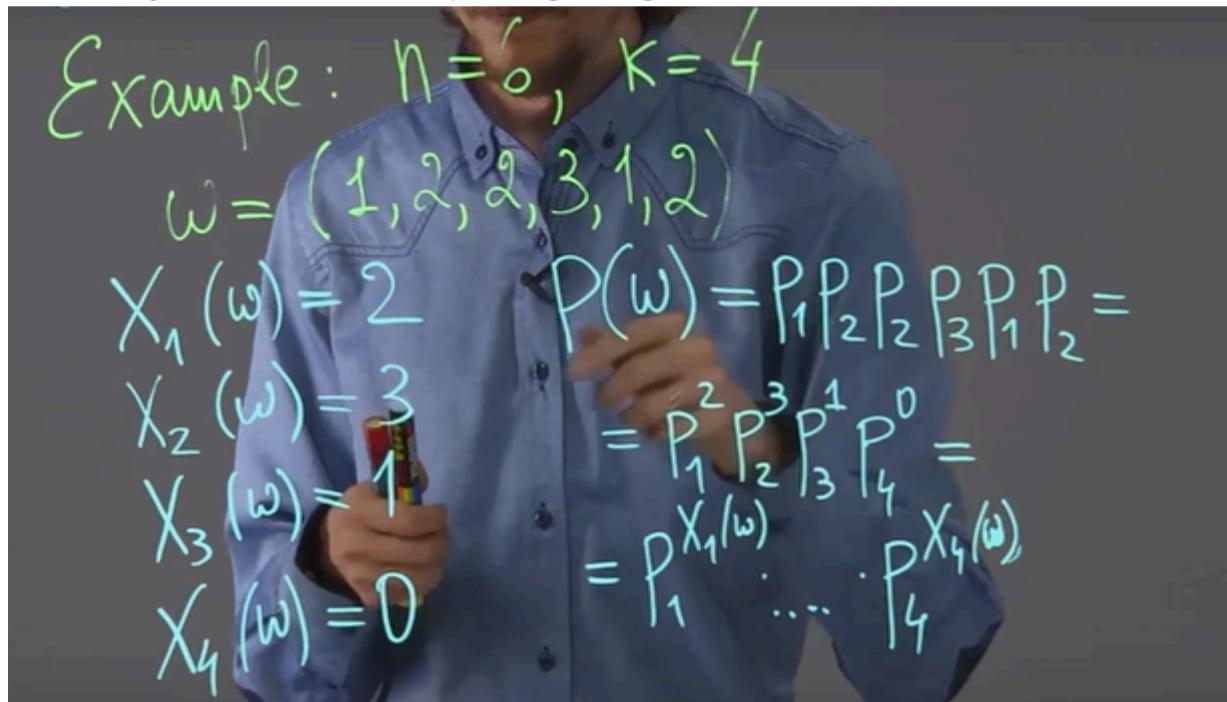


This kind vectors is used in so-called **one-hot encoding** when we have to encode a one of several categories using some numeric vector. And also this distribution called **categorical distribution**.

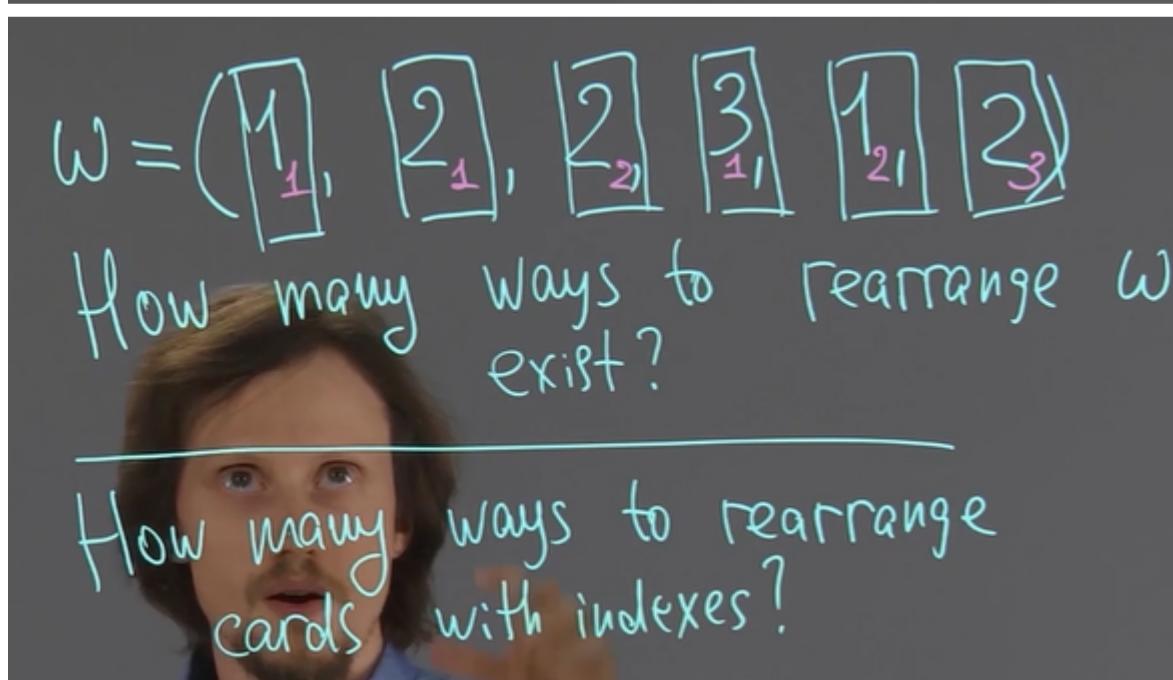
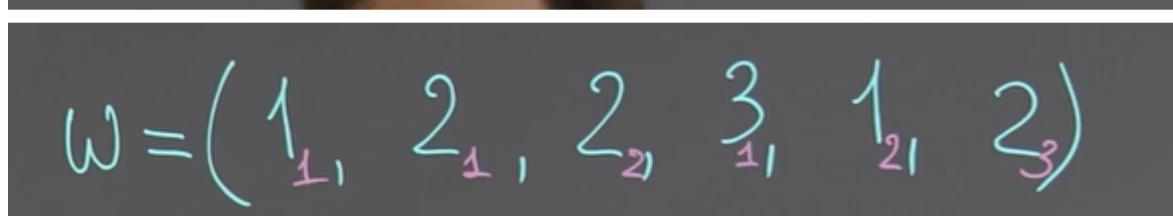
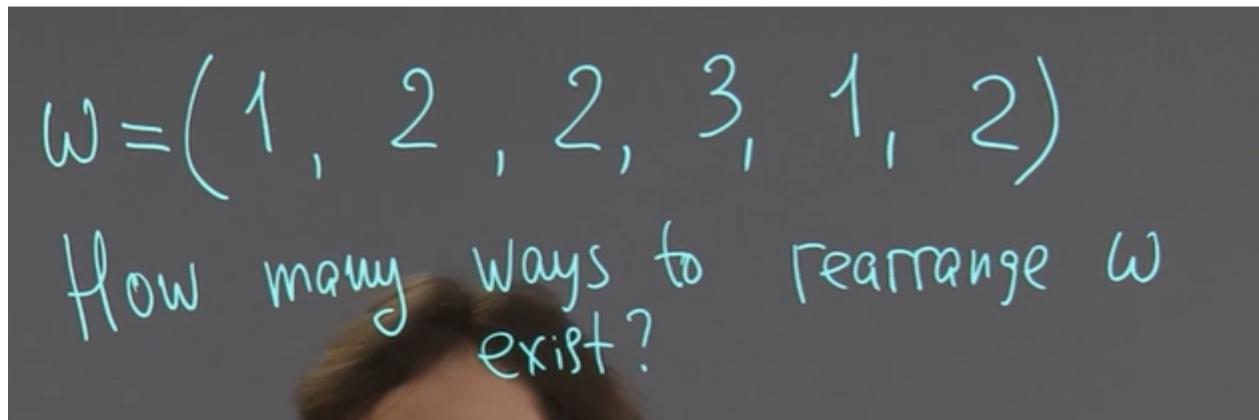
Explicit formula of multinomial distribution



Each S is just a result of corresponding rolling dice.



How many possible ways satisfy ω sequence? If we arrange it we get many of them. Let's count by combinatorics:



All this card are different to each other, so number of ways to rearrange them is equal to $n! = 6!$

But no of this rearranges are different if we will ignore indexes. Now let us try to remove some indexes and check how will affect this number($6!$).

We can 1 rearrange 2 times = 2!

We can 2 rearrange 3 times = 3!

We can 3 rearrange 3 times = 1!

Let's return to original question:

$$P(X_1=m_1 \cap X_2=m_2 \cap \dots \cap X_k=m_k) = ?$$

We found the number of outcomes that produced from this outcome w by rearranging its elements. And we seen that for each such outcome the values $X_i(w)$ are the same as before. It means that to find the probability that we need, we just have to multiply the probability of all such outcomes(w) which are all equal to the same number p by the number of these outcomes.

That's is universal formula

This formula gives us multinomial distribution in general case.

Language detection: formalization & maximum likelihood principle

Let's us formalize our language detection problem using multinomial distribution.

We believe that our text, which is also given and is just a sequence of characters, is just a result of random experiment that's similar to the dice rolling with k faces and on each face we have a letter from our alphabet. We throw this dice several times, number of times is equal n - size of text, and the text is generated just by this dice. Of course, this model is not correct. No real texts are generated by this random generation. But this formula is useful for our purpose.

So now we want to answer on : Which language came according to our model?

Likelihood is a probability to observe the data that we observed under assumption that some model for our data is valid. In this case, likelihood is the probability that some r.v. that is multinomially distributed with some properties, have values that we get from our text.

$$\underbrace{(f_1(T), \dots, f_k(T))}_{f(T)} \sim \text{Multinomial}(n, k, (p_{i1}, \dots, p_{ik}))$$

Size of T

Part 4

Like likelihood

$$(X_{i1}, \dots, X_{ik}) \sim \text{Multinomial}(n, k, (p_{i1}, \dots, p_{ik}))$$

$$P(f(T) | L_i) = P(X_{i1} = f_1(T) \cap \dots \cap X_{ik} = f_k(T))$$

So the likelihood is a probability to obtain the values of $f(T)$ that we really obtained provided that our model L_i takes place.

How can we use it to select the language?

In maximum likelihood principle, we just select the most likelihood model. So we just find i such that $P(f(T), L_i)$ is maximum.

$$i_{\max} = \underset{i=1, \dots, l}{\operatorname{argmax}} P(f(T) | L_i)$$

So in maximum likelihood settings, we just select the model for which likelihood of our data is maximum. However, this is not the only possible solution.

In fact, it's possible that we higher priority know that the probability of one language is less than probability of other language. Let us assume, for example, that we have some rare languages, and it's quite unlikely to meet it in these languages. In these case, we need to have much more evidence in favor of these languages to accept them, at least from naive point of view. To formalize this approach we have to switch maximum likelihood perspective to Bayesian perspective. Now, let us discuss how to do it.

Bayesian inference and interpretation of results

Bayesian approach allows us to take into account some prior information about the probability to get a text of some language. So let us discuss how Bayes' rule allows us to do it.

$f(T)$ = vector of frequencies of each character

L_1, \dots, L_e - languages
 T - text $f(T) = (f_1(T), \dots, f_k(T))$
 $P(f(T) | L_i)$ - likelihood
 $P(L_i) = q_i$ - a prior probabilities of languages
 $i = 1, \dots, l$
 $P(L_i | f(T)) = ?$

What we really want to select this language L_i .

In Bayesian approach, we believe that L_i is also chosen randomly before we start generating our text. So we make our model a little bit more complicated by adding one step to it. We believe, that first we select the language and then we use this language to generate out text $(f(T))$ and we some a prior distribution on the set of languages.

$P(L_i) = q_i$ - a prior probabilities of languages

What do we want to estimate in Bayesian approach is a probability that our text belongs to some language after we looked at this text and calculated $P(L_i)$ and took into account likelihood.

$$P(L_i | f(T)) = ?$$

Recall: Bayes' theorem

$$H_1, \dots, H_n \quad H_a \cap H_b = \emptyset \quad a \neq b$$

$$A \quad H_1 \cup \dots \cup H_n = \Omega$$

$$P(H_i | A) = \frac{P(A|H_i) \times P(H_i)}{\sum_{s=1}^n P(A|H_s) \times P(H_s)}$$

Rewrite for our settings:

$$P(L_i | f(T)) = \frac{P(f(T) | L_i) \times P(L_i)}{\sum_{s=1}^n P(f(T) | L_s) \times P(L_s)}$$

So we can use this to find posterior probability that we need to select the language with highest probability.

$$i_B = \arg \max_{i=1, \dots, l} P(L_i | f(T))$$

So if you want to choose select the language which maximizes $P(L_i | f(T))$, then you can ignore the denominator and compare only numerators, because for all values $P(L_i | f(T))$ we have the same denominator and it does not depend on i . But if you want this full probability $P(L_i | f(T))$ you have to take into account this denominator.

You can interpret $P(L_i | f(T))$ as a measure of how certain algorithm, how confident it's in its prediction. For example , if you see that the probability of even the best language is not so high, it means that model (your algorithm) isn't too much confident about this answer and probably you have to investigate it by different tools.

SGA: Final proof

The log-likelihood function

$$L((f_1, \dots, f_k), (p_1, \dots, p_k)) = \sum_{i=1}^k f_i \log p_i$$

that you (hopefully) obtained in final project is also known as cross entropy. It is extremely popular in machine learning, namely, in classification problems. It allows you to measure how good probability distribution (p_1, \dots, p_k) fits the actual absolute frequencies obtained from the data (f_1, \dots, f_k) . Assume that frequencies (f_1, \dots, f_k) are fixed. What is the best distribution (p_1, \dots, p_k) from likelihood's perspective? Intuitively, it seems that we have to put relative frequencies

$$r_i = \frac{f_i}{\sum_{j=1}^k f_j}$$

as p_i to get best fit. In fact, it is true. To prove it, let us use Jensen's inequality for logarithms. It is stated as follows:

For any values $\alpha_1, \dots, \alpha_k$, such that $\sum_{j=1}^k \alpha_j = 1$ and $\alpha_j \geq 0$ for all $j = 1, \dots, k$, and any positive values x_1, \dots, x_k , the following inequality holds:

$$\log \sum_{j=1}^k \alpha_j x_j \geq \sum_{j=1}^k \alpha_j \log(x_j).$$

Use this inequality to prove that

$$\sum_{j=1}^k r_j \log p_j - \sum_{j=1}^k r_j \log r_j \leq 0,$$

then prove that to obtain maximum log-likelihood (and therefore maximum likelihood) for fixed (f_1, \dots, f_k) we have to put $p_i = r_i$, $i = 1, \dots, k$.

Hint. Use properties of logarithm to transform left-hand part of the last inequality to right-hand part of the previous inequality.

Solution:

We need to prove inequality:

$$\sum_j^k r_j \log p_j - \sum_{j=1}^k r_j \log r_j \leq 0$$

To prove that to obtain maximum log-likelihood (and therefore maximum likelihood) for fixed (f_1, \dots, f_k) we have to put $p_i = r_i$, $i = 1, \dots, k$.

We can Jensen's inequality to prove it. For any values $\sum a_i = 1$, $a_i > 0$ and for all $i = 1, \dots, k$ any positive values x_1, \dots, x_k the inequality looks like:

$$\log \sum_{j=1}^k a_j x_j \geq \sum_{j=1}^k a_j \log(x_j)$$

Let's apply our inequality to our log likelihood function:

$$L(f_1, \dots, f_k, (p_1, \dots, p_k)) = \sum_{i=1}^k f_i \log p_i$$

So, we get:

$$\sum_{i=1}^k f_i \log(p_i) \geq \sum_{i=1}^k f_i \log(r_i)$$

We need to prove inequality:

$$\sum_j^k r_j \log p_j - \sum_{j=1}^k r_j \log r_j \leq 0$$

To prove that to obtain maximum log-likelihood (and therefore maximum likelihood) for fixed (f_1, \dots, f_k) we have to put $p_i = r_i$, $i = 1, \dots, k$.

We can Jensen's inequality to prove it. For any values $\sum a_i = 1$, $a_i > 0$ and for all $i = 1, \dots, k$ any positive values x_1, \dots, x_k the inequality looks like:

$$\log \sum_{j=1}^k a_j x_j \geq \sum_{j=1}^k a_j \log(x_j)$$

Let's apply our inequality to our log likelihood function:

$$L(f_1, \dots, f_k, (p_1, \dots, p_k)) = \sum_{i=1}^k f_i \log p_i$$

So, we get:

$$\sum_{i=1}^k f_i \log(p_i) \geq \sum_{i=1}^k f_i \log(r_i)$$

We need to show that $\sum_{i=1}^k f_i \log(r_i) \leq 0$. In order to prove that obtain maximum likelihood for fixed $(p_i = r_i, i = 1, \dots, k)$.

Let's simplify:

$$\sum f_i p_i \log(r_i) = \sum f_i (p_i \log(p_i))$$

Since $\sum f_i \cdot p_i = 1$:

$$\sum f_i p_i \log(r_i) - \sum f_i p_i \log(p_i) = \sum f_i p_i \log\left(\frac{r_i}{p_i}\right)$$

$$\sum f_i p_i \log(r_i) \leq \sum f_i p_i \log(p_i)$$

$$\sum f_i p_i \log(r_i) - \sum f_i p_i \log(p_i) \geq 0$$

$$\sum f_i p_i \log\left(\frac{r_i}{p_i}\right) \geq 0$$

Since $L((f_i)(p_i))$ is the log likelihood for the given probabilities is the likelihood function for the relative frequencies. Therefore we have shown that $\sum_{i=1}^k f_i p_i \log(r_i) \leq 0$. This means that in order to maximize the log likelihood function for the fixed frequency (f_1, \dots, f_k) we need to

set $p_i = r_i$.

