

Задачи оценивания значимости выравнивания при помощи скрытых марковских моделей

Власенко Даниил Владимирович, гр.19.Б04-мм

Научный руководитель: к.ф.-м.н. Коробейников А.И.

Санкт-Петербургский государственный университет
Прикладная математика и информатика
Вычислительная стохастика и статистические модели

Отчет по производственной практике

Санкт-Петербург, 2022

Оценивание значимости выравнивания

Задачи оценивания значимости выравнивания
при помощи скрытых марковских моделей

Власенко Даниил Владимирович, гр.19.Б04-мм
Научный руководитель: к.ф.-м.н. Коробейников А.И.

Санкт-Петербургский государственный университет
Прикладная математика и информатика
Вычислительная стохастика и статистические модели

Отчет по производственной практике
Санкт-Петербург, 2022

Научный руководитель к.ф.-м.н., Коробейников А.И.,
кафедра статистического моделирования

Введение

Пусть дан алфавит символов Σ .

Определение

Последовательностью длины L над алфавитом Σ будем называть такой элемент X , что $X \in \Sigma^L$. Последовательностью X над алфавитом Σ будем называть такой X , что $X \in \bigcup_{L=0}^{L=\infty} \Sigma^L$.

Определение

Парное выравнивание последовательностей называется отображение $Q : (\bigcup_{L_1=0}^{L_1=\infty} \Sigma^{L_1} \times \bigcup_{L_2=0}^{L_2=\infty} \Sigma^{L_2}) \rightarrow (\Sigma^{\max(L_1, L_2)} \times \Sigma^{\max(L_1, L_2)})$, такое что:

1. Возможны вставки символа — в последовательностях.
2. Вставка — на одинаковых позициях в обеих последовательностях запрещена.
3. Порядок изначальных символов внутри последовательностей сохраняется.

2/4

Власенко Д.В.

Оценивание значимости выравнивания

Оценивание значимости выравнивания

— Введение

Введение

Пусть дан алфавит символов Σ .

Определение

Последовательностью длины L над алфавитом Σ будем называть такой элемент X , что $X \in \Sigma^L$. Последовательностью X над алфавитом Σ будем называть такой X , что $X \in \bigcup_{L=0}^{L=\infty} \Sigma^L$.

Определение

Парное выравнивание последовательностей называется отображение $Q : (\bigcup_{L_1=0}^{L_1=\infty} \Sigma^{L_1} \times \bigcup_{L_2=0}^{L_2=\infty} \Sigma^{L_2}) \rightarrow (\Sigma^{\max(L_1, L_2)} \times \Sigma^{\max(L_1, L_2)})$, такое что:

1. Возможны вставки символа — в последовательностях.
2. Вставка — на одинаковых позициях в обеих последовательностях запрещена.
3. Порядок изначальных символов внутри последовательностей сохраняется.

Пусть дан алфавит символов Σ .

Определение

Последовательностью длины L над алфавитом Σ будем называть такой элемент X , что $X \in \Sigma^L$. Последовательностью X над алфавитом Σ будем называть такой X , что $X \in \bigcup_{L=0}^{L=\infty} \Sigma^L$.

Определение

Парное выравнивание последовательностей называется отображение $Q : (\bigcup_{L_1=0}^{L_1=\infty} \Sigma^{L_1} \times \bigcup_{L_2=0}^{L_2=\infty} \Sigma^{L_2}) \rightarrow (\Sigma^{\max(L_1, L_2)} \times \Sigma^{\max(L_1, L_2)})$, такое что:

1. Возможны вставки символа — в последовательностях.
2. Вставка — на одинаковых позициях в обеих последовательностях запрещена.
3. Порядок изначальных символов внутри последовательностей сохраняется.

A	C	E	A	A	F	A	E
C	E	A	F	D	C	E	

A	C	E	A	A	F	A	—	E
—	C	E	A	—	F	D	C	E

Рис. 1: Последовательности до и после парного выравнивания.

Примем множество $\bigcup_{L_1=0}^{L_1=\infty} \Sigma^{L_1} \times \bigcup_{L_2=0}^{L_2=\infty} \Sigma^{L_2}$ за пространство элементарных исходов Ω . Область значений выравнивания Q обозначим как $\bar{\Omega}$.

Определение

Оценкой парного выравнивания называется случайная величина $s : \bar{\Omega} \rightarrow \mathbb{R}$.

Оценивание значимости выравнивания

— Введение

Введение

A	C	E	A	A	F	A	E
C	E	A	F	D	C	E	

A	C	E	A	A	F	A	—	E
—	C	E	A	—	F	D	C	E

Рис. 1: Последовательности до и после парного выравнивания.

Примем множество $\bigcup_{L_1=0}^{L_1=\infty} \Sigma^{L_1} \times \bigcup_{L_2=0}^{L_2=\infty} \Sigma^{L_2}$ за пространство элементарных исходов Ω . Область значений выравнивания Q обозначим как $\bar{\Omega}$.

Определение

Оценкой парного выравнивания называется случайная величина $s : \bar{\Omega} \rightarrow \mathbb{R}$.

Примем множество $\bigcup_{L_1=0}^{L_1=\infty} \Sigma^{L_1} \times \bigcup_{L_2=0}^{L_2=\infty} \Sigma^{L_2}$ за пространство элементарных исходов Ω . Область значений выравнивания Q обозначим как $\bar{\Omega}$.

Определение

Оценкой парного выравнивания называется случайная величина $s : \bar{\Omega} \rightarrow \mathbb{R}$.

Способом вычисления оценки выравнивания s может быть, например, увеличение оценки на 1 при совпадении символов, стоящих друг под другом, и уменьшение на $\frac{1}{2}$ при несовпадении. Тогда оценка s приведенного на слайде выравнивания будет равна 3.

Сходство последовательностей может отражать функциональные, структурные или эволюционные взаимосвязи объектов, которые описывают эти последовательности. Таким образом вычисление оценки выравнивания последовательностей может быть полезно в задаче определения степени родства биологических организмов путем сравнения их ДНК или РНК, нуклеотидных последовательностей, задаче анализа свойств белков, аминокислотных последовательностей, задаче распознавания речи человека или письменного языка и многих других приложениях.

Определение

Множественным выравнением N последовательностей называется отображение $Q: \times_{i=1}^N (\bigcup_{L_i=0}^{L_i=\infty} \Sigma^{L_i}) \rightarrow \times_{i=1}^N (\Sigma^{\max_{L \in L_i} (L)})$, такое что:

1. Возможны вставки символа — в последовательностях.
2. Вставка — на одинаковых позициях во всех последовательностях запрещена.
3. Порядок изначальных символов внутри последовательностей сохраняется.

Примем множество $\times_{i=1}^N (\bigcup_{L_i=0}^{L_i=\infty} \Sigma^{L_i})$ за пространство элементарных исходов Ω . Область значений выравнения Q обозначим как $\bar{\Omega}$.

Определение

Оценкой парного выравнения называется случайная величина $s: \bar{\Omega} \rightarrow \mathbb{R}$.

Оценивание значимости выравнения

— Введение

Введение

Определение

Множественным выравнением N последовательностей называется отображение $Q: \times_{i=1}^N (\bigcup_{L_i=0}^{L_i=\infty} \Sigma^{L_i}) \rightarrow \times_{i=1}^N (\Sigma^{\max_{L \in L_i} (L)})$, такое что:

1. Возможны вставки символа — в последовательностях.
2. Вставка — на одинаковых позициях во всех последовательностях запрещена.
3. Порядок изначальных символов внутри последовательностей сохраняется.

Примем множество $\times_{i=1}^N (\bigcup_{L_i=0}^{L_i=\infty} \Sigma^{L_i})$ за пространство элементарных исходов Ω . Область значений выравнения Q обозначим как $\bar{\Omega}$.

Определение

Оценкой парного выравнения называется случайная величина $s: \bar{\Omega} \rightarrow \mathbb{R}$.

На предыдущем слайде приведен пример попарного выравнения двух строк, но если сходство последовательностей слабое, то через такое выравнение может не выйти идентифицировать взаимосвязь описываемых последовательностями объектов. В таких случаях может помочь множественное выравнение, обобщим имеющиеся определения.

Определение

Множественным выравнением N последовательностей называется отображение $Q: \times_{i=1}^N (\bigcup_{L_i=0}^{L_i=\infty} \Sigma^{L_i}) \rightarrow \times_{i=1}^N (\Sigma^{\max_{L \in L_i} (L)})$, такое что:

1. Возможны вставки символа — в последовательностях.
2. Вставка — на одинаковых позициях во всех последовательностях запрещена.
3. Порядок изначальных символов внутри последовательностей сохраняется.

Определение

Множественным выравниванием N последовательностей называется отображение $Q: \times_{i=1}^N (\bigcup_{L_i=0}^{L_i=\infty} \Sigma^{L_i}) \rightarrow \times_{i=1}^N (\Sigma^{\max_{L \in L_i} (L)})$, такое что:

- 1 Возможны вставки символа — в последовательностях.
- 2 Вставка — на одинаковых позициях во всех последовательностях запрещена.
- 3 Порядок изначальных символов внутри последовательностей сохраняется.

Примем множество $\times_{i=1}^N (\bigcup_{L_i=0}^{L_i=\infty} \Sigma^{L_i})$ за пространство элементарных исходов Ω . Область значений выравнивания Q обозначим как $\bar{\Omega}$.

Определение

Оценкой парного выравнивания называется случайная величина $s: \bar{\Omega} \rightarrow \mathbb{R}$.

Оценивание значимости выравнивания

└ Введение

Примем множество $\times_{i=1}^N (\bigcup_{L_i=0}^{L_i=\infty} \Sigma^{L_i})$ за пространство элементарных исходов Ω . Область значений выравнивания Q обозначим как $\bar{\Omega}$.

Определение

Оценкой парного выравнивания называется случайная величина $s: \bar{\Omega} \rightarrow \mathbb{R}$.

Введение

Определение

Множественным выравниванием N последовательностей называется отображение $Q: \times_{i=1}^N (\bigcup_{L_i=0}^{L_i=\infty} \Sigma^{L_i}) \rightarrow \times_{i=1}^N (\Sigma^{\max_{L \in L_i} (L)})$, такое что:

- 1 Возможны вставки символа — в последовательностях.
- 2 Вставка — на одинаковых позициях во всех последовательностях запрещена.
- 3 Порядок изначальных символов внутри последовательностей сохраняется.

Примем множество $\times_{i=1}^N (\bigcup_{L_i=0}^{L_i=\infty} \Sigma^{L_i})$ за пространство элементарных исходов Ω . Область значений выравнивания Q обозначим как $\bar{\Omega}$.

Определение

Оценкой парного выравнивания называется случайная величина $s: \bar{\Omega} \rightarrow \mathbb{R}$.