

Algorithmen auf Sequenzen

Exakte Mustersuche

Sven Rahmann

Genominformatik
Universitätsklinikum Essen
Universität Duisburg-Essen
Universitätsallianz Ruhr

Das Pattern Matching Problem (Exakte Mustersuche)



Gegeben sei ein endliches Alphabet Σ , ein Text $T \in \Sigma^n$ und ein Muster (Pattern) $P \in \Sigma^m$ i.d.R. $m \ll n$.

Gesucht sind alle Positionen i , für die gilt: $T[i : i + m] = P$.

Das Pattern Matching Problem (Exakte Mustersuche)

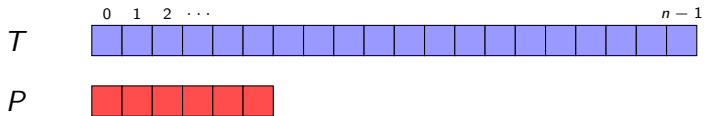
Gegeben sei ein endliches Alphabet Σ , ein Text $T \in \Sigma^n$ und ein Muster (Pattern) $P \in \Sigma^m$ i.d.R. $m \ll n$.

Gesucht sind alle Positionen i , für die gilt: $T[i : i + m] = P$.

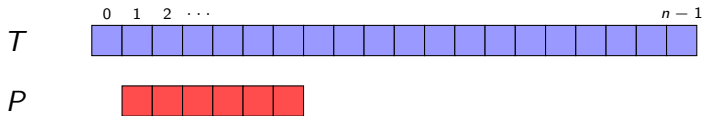
Varianten:

- Kommt P in T vor?
- Wie oft kommt P in T vor?

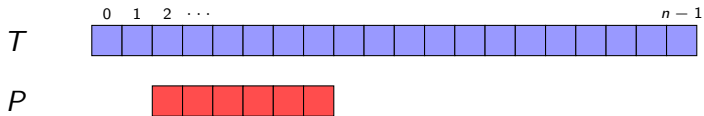
Naiver Algorithmus



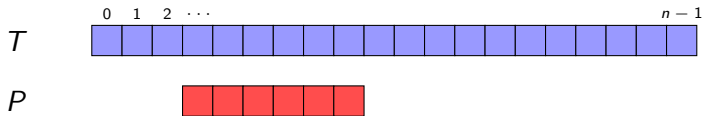
Naiver Algorithmus



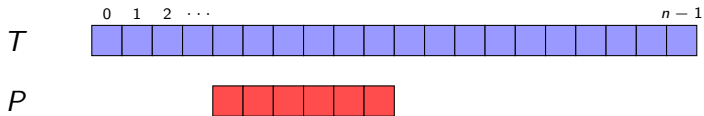
Naiver Algorithmus



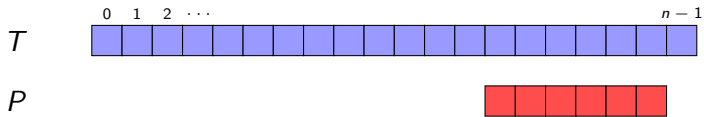
Naiver Algorithmus



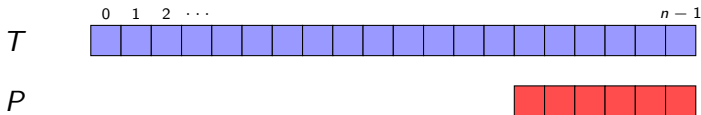
Naiver Algorithmus



Naiver Algorithmus

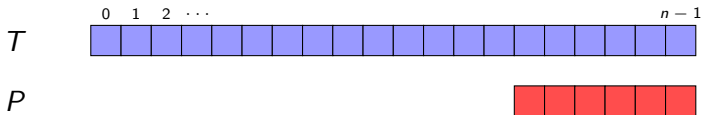


Naiver Algorithmus



Insgesamt $n - m + 1$ Verschiebungen.
Pro Verschiebung bis zu m Vergleiche.
Gesamte Laufzeit: $\mathcal{O}(mn)$.

Naiver Algorithmus



Insgesamt $n - m + 1$ Verschiebungen.
 Pro Verschiebung bis zu m Vergleiche.
 Gesamte Laufzeit: $\mathcal{O}(mn)$.

```

1 def naive(P, T):
2     m, n = len(P), len(T)
3     for i in range(n - m + 1):
4         if T[i:i+m] == P:
5             yield (i, i+m)
  
```

Theorem (Erwartete Laufzeit)

Sei $|\Sigma| \geq 2$. Seien ein Muster der Länge m und ein Text der Länge n zufällig gleichverteilt gewählt. Dann beträgt die Worst-case-Laufzeit des naiven Algorithmus $\mathcal{O}(mn)$, aber die erwartete Laufzeit lediglich $\mathcal{O}(n)$.

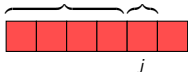
Die Wahrscheinlichkeit p , dass jeweils ein zufällig gezogenes Zeichen aus P und T übereinstimmen, beträgt

$$p := \frac{|\Sigma|}{|\Sigma|^2} = \frac{1}{|\Sigma|}.$$

Die Wahrscheinlichkeit p , dass jeweils ein zufällig gezogenes Zeichen aus P und T übereinstimmen, beträgt

$$p := \frac{|\Sigma|}{|\Sigma|^2} = \frac{1}{|\Sigma|}.$$

- Die W'keit, dass alle Zeichen des Musters mit dem Text übereinstimmen, beträgt p^m .
- Die W'keit, dass das Muster erst an Stelle j mit dem Text nicht übereinstimmt, beträgt $p^{j-1} \cdot (1 - p)$.



Die erwartete Anzahl an Vergleichen für ein Muster der Länge m beträgt demnach (Erfolg + Misserfolg an Stelle j für alle j)

$$E_m := mp^m + \sum_{j=1}^m jp^{j-1} \cdot (1 - p)$$

Die erwartete Anzahl an Vergleichen für ein Muster der Länge m beträgt demnach (Erfolg + Misserfolg an Stelle j für alle j)

$$E_m := mp^m + \sum_{j=1}^m jp^{j-1} \cdot (1-p)$$

Für beliebige Längen $m \rightarrow \infty$ ist

$$E_m < E_\infty := (1-p) \sum_{j=0}^{\infty} jp^{j-1}.$$

Die erwartete Anzahl an Vergleichen für ein Muster der Länge m beträgt demnach (Erfolg + Misserfolg an Stelle j für alle j)

$$E_m := mp^m + \sum_{j=1}^m jp^{j-1} \cdot (1-p)$$

Für beliebige Längen $m \rightarrow \infty$ ist

$$E_m < E_\infty := (1-p) \sum_{j=0}^{\infty} jp^{j-1}.$$

Der Term $\sum_{j=0}^{\infty} jp^{j-1}$ ist die Ableitung von $\sum_{j=0}^{\infty} p^j = 1/(1-p)$.
Somit gilt $\sum_{j=0}^{\infty} jp^{j-1} = 1/(1-p)^2$, also

$$E_\infty = \frac{1-p}{(1-p)^2} = \frac{1}{1-p}.$$

Mit $p = 1/|\Sigma|$ folgt

$$E_m < E_\infty = \frac{1}{1-p} = \frac{|\Sigma|}{|\Sigma|-1} \text{ für alle } m \geq 1.$$

Mit $p = 1/|\Sigma|$ folgt

$$E_m < E_\infty = \frac{1}{1-p} = \frac{|\Sigma|}{|\Sigma|-1} \text{ für alle } m \geq 1.$$

Für ein 2-buchstabiges Alphabet ist $E_m < 2$ für alle m .

Für $|\Sigma| \rightarrow \infty$ sogar $E_m \searrow 1$.

Die erwartete Laufzeit auf gleichverteilten zufälligen Texten ist also $\mathcal{O}(n \cdot E_m) = \mathcal{O}(n)$

