

Porównywanie sekwencji

Homologia, podobieństwo i analogia

Homologi

- **Ortologi** – homologiczne geny, których rozdzielenie nastąpiło na skutek specjacji, czyli rozdzielenia gatunków, lub rzadziej horyzontalnego transferu genu. Geny ortologiczne mają zwykle taką samą, albo zbliżoną funkcję.
- **Paralogi** – geny pochodzące od wspólnego przodka, rozdzielone w wyniku duplikacji genu. Paralogi mają często różne funkcje w organizmie. Przykładem mogą być mioglobina i hemoglobina u człowieka.

dopasowanie sekwencji

- Dopasowanie/porównywanie
- Uliniowanie
- Alignment

W bioinformatyce, dopasowanie sekwencji jest sposobem dopasowania struktur pierwszorzędowych DNA, RNA, lub białek do zidentyfikowania regionów wykazujących podobieństwo, mogące być konsekwencją funkcjonalnych, strukturalnych, lub ewolucyjnych powiązań pomiędzy sekwencjami. Zestawione sekwencje nukleotydów lub aminokwasów są zazwyczaj przedstawiane jako wiersze macierzy. Pomiędzy reszty wprowadzane są przerwy, tak że reszty zbliżonych do siebie sekwencji tworzą kolejne kolumny.

Jeśli dwie dopasowywane sekwencje mają wspólne pochodzenie, niedopasowania mogą być interpretowane jako mutacje punktowe, a przerwy jako indela (mutacje polegające na delecji lub insercji), które zaszły w jednej lub obu liniach od czasu, kiedy obie sekwencje uległy rozdzieleniu. W przypadku dopasowywania sekwencji białek, stopień podobieństwa pomiędzy aminokwasami zajmującymi konkretną pozycję, może stanowić zgrubną miarę tego, jak konserwatywny jest dany region lub motyw. Brak substytucji lub obecność jedynie konserwatywnych substytucji (tj. zamiany reszty na inną, ale o podobnych właściwościach chemicznych) w określonym regionie sekwencji sugeruje, że jest on ważny strukturalnie lub funkcjonalnie. Dopasowywanie sekwencji może być także stosowane dla sekwencji pochodzenia poza biologicznego, np. danych finansowych lub sekwencji występujących w językach naturalnych.

Masuri i inni. Dopasowanie sekwencji. Wikipedia 11.2009

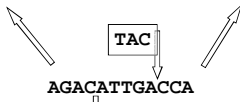
Skąd te różnice

różnice między sekwencjami świadczą o mutacjach, które zaszły po rozdzieleniu się sekwencji od wspólnego przodka

```

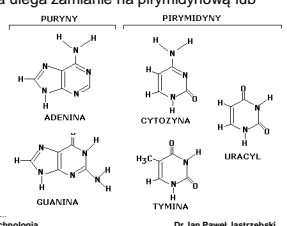
AGA--TTGATACCCA
|||  |||  |||  |||
AGACATTAA---CTA
    
```

AGA--TTGATACCCA Insercja +TAC Delecja -CA	AGACATTAA---CTA G->A C->T substytucje
---	---



Substytucje nukleotydowe

- Tranzycja** - okres przejściowy między systemem politycznym, który był, a tym który nastąpi. Proces ten jest krótszy i łatwiejszy od konsolidacji systemu politycznego. Tranzycja kończy się gdy pojawiają się ogólne ramy funkcjonowania nowego systemu. Przykładem są wszystkie państwa byłego bloku wschodniego, w tym Polska. (Czy o to chodzi?)
- Transwersja** - mutacja genu, punktowa zmiana chemiczna w obrębie nici DNA, w której zasada purynowa ulega zamianie na pirymidynową lub odwrotnie. Mutacja taka może nie spowodować żadnej zmiany lub zmianę kodu genetycznego (UUU -> UUA) albo też skróconą syntezę białka (UCG -> UCA).



Zastosowanie alignmentu

- poszukiwaniu oraz określaniu funkcji i struktury (białek) dla „nowych” sekwencji (nieznanych nam do tej pory)
- określaniu powiązań filogenetycznych między sekwencjami - homologii między sekwencjami oraz w analizach ewolucyjnych

Metody dopasowania

dopasowanie par sekwencji (*pairwise alignment*)

- **Macierz punktowe** - dot matrix, dotplot
- **Programowanie dynamiczne (DP)**
- **Metody słów (k - tuple methods)** - szybkie metody stosowane przy przeszukiwaniu baz danych sekwencji z wykorzystaniem programów FASTA i BLAST

- **dopasowanie wielu sekwencji (*multiple alignment*)**

Wstęp do bioinformatyki
Wykład 3

Biotechnologia
UWM

Dr Jan Paweł Jastrzębski

Etapy dopasowywania sekwencji

1 zestawienie (0 identycznych, 0% podobieństwa)

x = długość sekwencji (30)

MHSSIVLATVLFVAIASASKTRELCKMSLV

MHVSIVLATVLFVAIASAS

y = długość sekwencji (20)

2 zestawienie (0 identycznych, 0% podobieństwa)

MHSSIVLATVLFVAIASASKTRELCKMSLV

MHVSIVLATVLFVAIASAS

3 zestawienie (1 identyczna, 33% podobieństwa)

MHSSIVLATVLFVAIASASKTRELCKMSLV

MHVSIVLATVLFVAIASAS

Wstęp do bioinformatyki
Wykład 3

Biotechnologia
UWM

Dr Jan Paweł Jastrzębski

4 zestawienie (0 identycznych, 0% podobieństwa)

MHSSIVLATVLFVAIASASKTRELCKMSLV

MHVSIVLATVLFVAIASAS

5 zestawienie (0 identycznych, 0% podobieństwa)

MHSSIVLATVLFVAIASASKTRELCKMSLV

MHVSIVLATVLFVAIASAS

6 zestawienie (2 identyczne, 33% podobieństwa)

MHSSIVLATVLFVAIASASKTRELCKMSLV

MHVSIVLATVLFVAIASAS

7 zestawienie (0 identycznych, 0% podobieństwa)

MHSSIVLATVLFVAIASASKTRELCKMSLV

MHVSIVLATVLFVAIASAS

Wstęp do bioinformatyki
Wykład 3

Biotechnologia
UWM

Dr Jan Paweł Jastrzębski

Kilka możliwych rozwiązań:

```
AAGCTGAATTCGAA
AGGCTCATTCTGA
```

```
A-AGCTGAATTC--GAA
AG-GCTCA-TTTCTGA-
```

```
AAGCTGAATT-C-GAA
AGGCT-CATTCTGA-
```

Który alignment jest lepszy?

scoring system:

- Perfect match: +1
- Mismatch: -2
- Indel (gap): -1 (*kara za przerwy*)

```
AAGCTGAATT-C-GAA
AGGCT-CATTCTGA-
```

```
A-AGCTGAATTC--GAA
AG-GCTCA-TTTCTGA-
```

Score: = (+1)x10 + (-2)x2 + (-1)x4 = **2** Score: = (+1)x9 + (-2)x2 + (-1)x6 = -1

Wyższy score → Lepszy alignment

Zadanie 1

- Jaki jest **score** tego alignmentu??

dopasowanie: +1
 niedopasowanie: -1
 przerwa: -2

```
---bardzo---lubiebioinformatyke
| | | | | | | | | * | | | | | | | | | *
niebardzonielubiębioinformatyki
```

Metody dopasowania

dopasowanie par sekwencji (*pairwise alignment*)

1. **Metody słów (k - tuple methods)** - szybkie metody stosowane przy przeszukiwaniu baz danych sekwencji z wykorzystaniem programów FASTA i BLAST
2. **Macierz punktowe** - dot matrix, dotplot
3. **Programowanie dynamiczne (DP)**

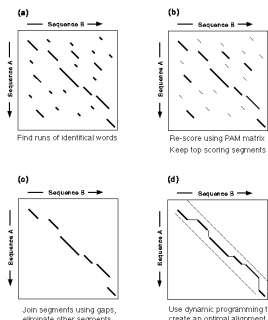
Wstęp do bioinformatyki
Wykład 3

Biotechnologia
UWM

Dr Jan Paweł Jastrzębski

1. „słowa” - FASTA

FASTA Algorithm



Wstęp do bioinformatyki
Wykład 3

Biotechnologia
UWM

Dr Jan Paweł Jastrzębski

1. „słowa” - BLAST vs. FASTA

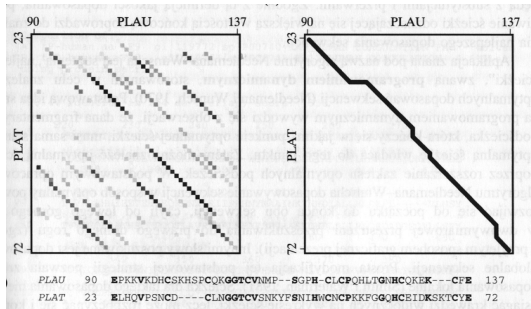
BLAST	FASTA
może podawać więcej niż jeden region o wysokiej punktacji	podaje tylko jedno najlepsze dopasowanie
lepszy dla sekwencji białek niż DNA	lepszy dla sekwencji DNA niż białek
szybszy niż FASTA	wolniejszy niż BLAST
mniej czuły niż FASTA przy użyciu domyślnych ustawień	bardziej czuły niż BLAST
daje gorsze rozróżnienie między prawdziwymi i fałszywymi homologami	daje lepsze rozróżnienie między prawdziwymi i fałszywymi homologami

Wstęp do bioinformatyki
Wykład 3

Biotechnologia
UWM

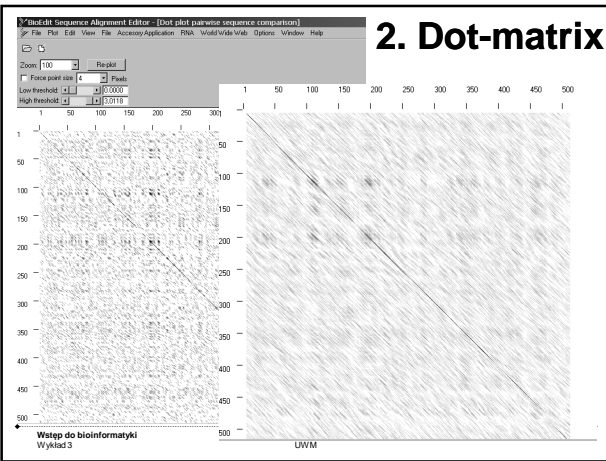
Dr Jan Paweł Jastrzębski

2. Macierze punktowe



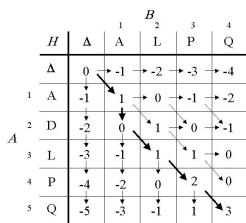
Wstęp do bioinformatyki Wykład 3 Biotechnologia UWM Dr Jan Paweł Jastrzębski

2. Dot-matrix

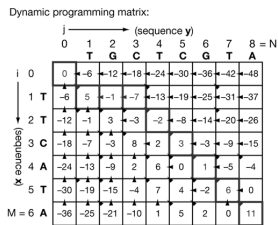


3. Programowanie dynamiczne

opiera się na podziale rozwiązywanego problemu na podproblemy względem kilku parametrów.



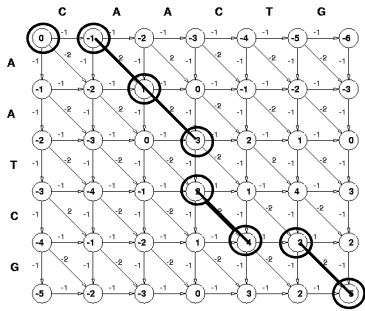
Alignment A... A D L P Q
 B... A A L P Q
 Similarity Score +1 -1 +1 +1 -1 = 3



Optimum alignment scores 11:
 T - - T C A T A
 T G C T C G T A
 +5 -6 -6 +5 +5 -2 +5 +5

Wstęp do bioinformatyki Wykład 3 Biotechnologia UWM Dr Jan Paweł Jastrzębski

3. Programowanie dynamiczne



Wstęp do bioinformatyki Wykład 3 Biotechnologia UWM Dr Jan Paweł Jastrzębski

Scoring matrix

- Reprezentuje system punktowania jako tabela lub macierz $n \times n$ (n jest liczbą liter, które zawiera alfabet. $n=4$ dla DNA, $n=20$ dla białek)
- Macierz punktowania jest symetryczna

	A	G	C	T
A	2			
G	-6	2		
C	-6	-6	2	
T	-6	-6	-6	2

Mismatch

Match

Wstęp do bioinformatyki Wykład 3 Biotechnologia UWM Dr Jan Paweł Jastrzębski

Podobieństwa biochemiczne i biofizyczne aminokwasów

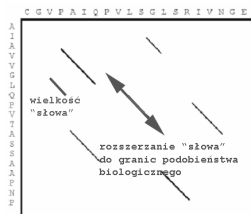
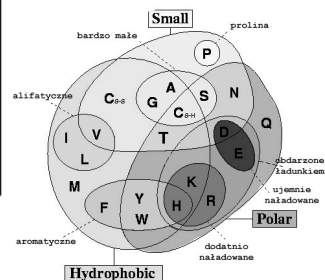


Diagram Venn-a



Wstęp do bioinformatyki Wykład 3 Biotechnologia UWM Dr Jan Paweł Jastrzębski

Macierze substytucji (podstawień)

- Jak za pomocą liczby określić podobieństwa biochemiczne i biofizyczne poszczególnych aminokwasów tak, aby liczba ta wyrażała jednocześnie realny wpływ na całe białko podstawienia danego aminokwasu w łańcuchu polipeptydowym i była uniwersalna dla wszystkich sekwencji?
- Przede wszystkim należy bazować na danych empirycznych
- Należy stworzyć alignment bardzo wielu blisko spokrewnionych sekwencji – na tyle podobnych, aby bez wątpliwości można było jednoznacznie i precyzyjnie określić częstotliwość substytucji poszczególnych aminokwasów w konkretnych pozycjach.

M	G	Y	D	E
M	G	Y	D	E
M	G	Y	E	E
M	G	Y	D	E
M	G	Y	E	E
M	G	Y	D	E
M	A	Y	E	E
M	A	Y	E	E

W kolumnie 4 E i D występują z częstotliwością w 4/8

Wstęp do bioinformatyki
Wykład 3

Biotechnologia
UWM

Dr Jan Paweł Jastrzębski

PAM Matrix – Point/Percent Accepted Mutations

*n*PAM (*n* Percent Accepted Mutations)

S_1, S_2 różnią się o jednostkę n PAM, jeśli S_2 można otrzymać z S_1 w ciągu akceptowalnych mutacji punktowych takich, że średnia liczba nieletalnych mutacji na 100 wynosi n . Najpopularniejsza jest 250PAM.

- Based on a database of 1,572 changes in 71 groups of closely related proteins (85% identity)
 - Alignment was easy

Wstęp do bioinformatyki
Wykład 3

Biotechnologia
UWM

Dr Jan Paweł Jastrzębski

PAM Matrices

- Family of matrices PAM 80, PAM 120, PAM 250
- The number on the PAM matrix represents evolutionary distance
- Larger numbers are for larger distances

Wstęp do bioinformatyki
Wykład 3

Biotechnologia
UWM

Dr Jan Paweł Jastrzębski

BLOSUM: Blocks Substitution Matrix

- Based on BLOCKS database
 - ~2000 blocks from 500 families of related proteins
 - Families of proteins with identical function
- Blocks are short conserved patterns of 3-60 aa long without gaps

AABCDA	---	BBBCDA
DABCDA	---	BBBCBB
BBBCDA	AA-	BCCAA
AAACDA	A-	CBDCB
CCBADA	---	DBDCC
AAACAA	---	BBCCC

Wstęp do bioinformatyki
Wykład 3

Biotechnologia
UWM

Dr Jan Paweł Jastrzębski

BLOSUM

- Each block represent sequences alignment with different identity percentage
- For each block the amino-acid substitution rates were calculated to create BLOSUM matrix

Wstęp do bioinformatyki
Wykład 3

Biotechnologia
UWM

Dr Jan Paweł Jastrzębski

BLOSUM Matrices

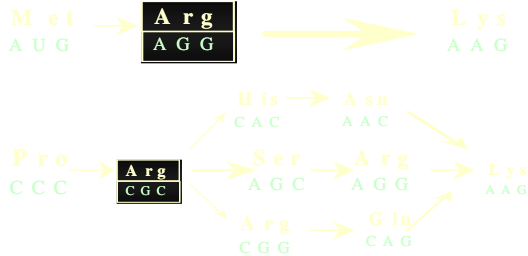
- BLOSUM n is based on sequences that shared at least n percent identity
- BLOSUM62 represents closer sequences than BLOSUM45

Wstęp do bioinformatyki
Wykład 3

Biotechnologia
UWM

Dr Jan Paweł Jastrzębski

Uwarunkowania genetyczne substytucji aminokwasowych



Wstęp do bioinformatyki
Wykład 3

Biotechnologia
UWM

Dr Jan Paweł Jastrzębski

Podstawy genetyczne algorytmów do zestawień aminokwasów?

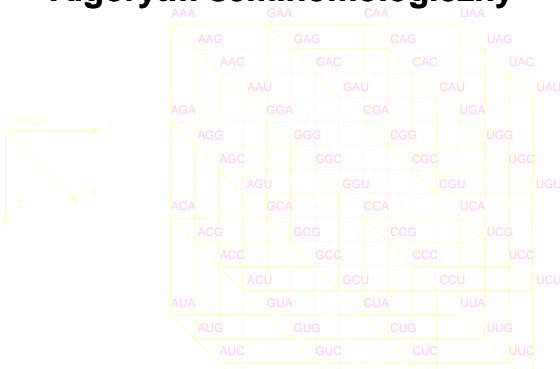
Replacement	PAM250	BLOSUM62
Arg/Lys	3	2
Lys/Gln	1	1
Arg/Gln	1	1
Lys/Glu	0	1
Arg/Glu	-1	0

Wstęp do bioinformatyki
Wykład 3

Biotechnologia
UWM

Dr Jan Paweł Jastrzębski

Algorytm semihomologiczny



Wstęp do bioinformatyki
Wykład 3

Biotechnologia
UWM

Dr Jan Paweł Jastrzębski

MSA & Evolution

MSA może dawać obraz sił kształtujących ewolucję !!!

- Ważne aminokwasy lub nukleotydy (pozycje w sekwencjach) mutują „niechętnie”
- Mniej ważne pozycje dla struktury i funkcji mogą wykazywać większą zmienność w kolumnach porównywanych sekwencji

Wstęp do bioinformatyki
Wykład 3

Biotechnologia
UWM

Dr Jan Paweł Jastrzębski

Pozycje konserwatywne

- Kolumny, gdzie wszystkie sekwencje zawierają takie same aminokwasy lub nukleotydy (lub w większości takie same – pozycje konserwatywne) są bardzo ważne (kluczowe) dla funkcji lub struktury.

```
VTISCTGSSSNIGAG-NHVKWYQQLPG
VTISCTGSSSNIGS--ITVNWYQQLPG
LRLSCTGSGFIFSS--YAMYWYQQAPG
LSLTC TGS GTSFDD-QYYSTWYQQPPG
```

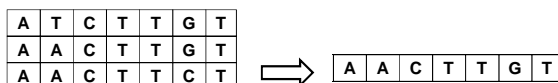
Wstęp do bioinformatyki
Wykład 3

Biotechnologia
UWM

Dr Jan Paweł Jastrzębski

Sekwencja konsensusowa

- W **sekwencji konsensusowej** zachowane są pozycje o największej częstotliwości występowania w każdej z kolumn alignmentu (The consensus sequence holds the most frequent character of the alignment at each column)
- Jest to sposób reprezentowania wyników multiple alignment, gdzie pokrewne sekwencje są porównywane każda do każdej, aby odnaleźć funkcjonalnie podobne motywy sekwencji (domeny białek). Sekwencja konsensusowa obrazuje które pozycje są konserwatywne, a które zmienne.



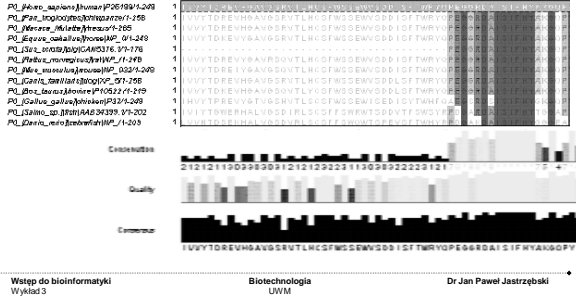
Wstęp do bioinformatyki
Wykład 3

Biotechnologia
UWM

Dr Jan Paweł Jastrzębski

Sekwencja konsensusowa

...***** ***** ..
.....



Alignment methods

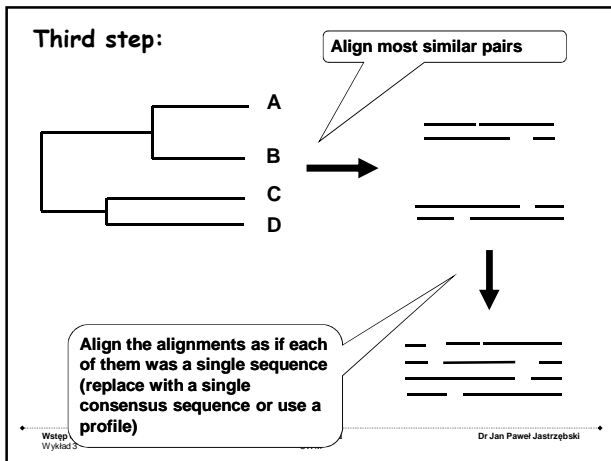
- Progressive alignment (Clustal)
- Iterative alignment (mafft, muscle)
- All methods today are an approximation strategy (**heuristic algorithm**), yield a possible alignment, but not necessarily the best one

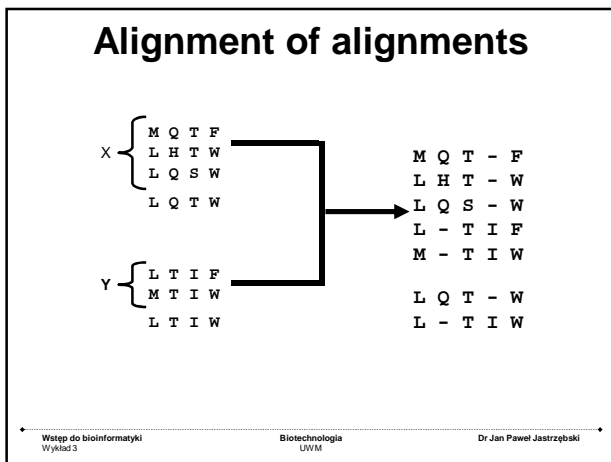
Wstęp do bioinformatyki Wykład 3 Biotechnologia UWM Dr Jan Paweł Jastrzębski

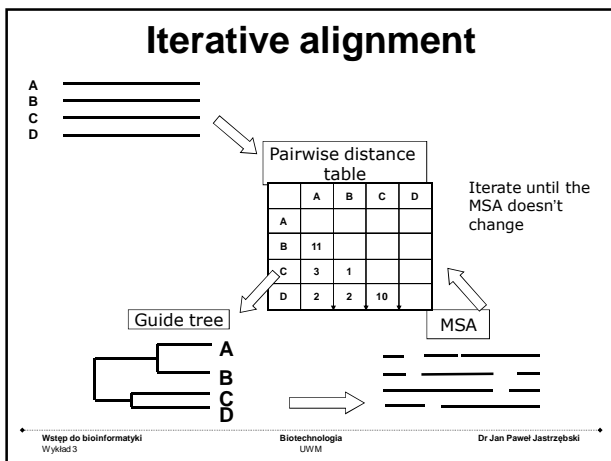
Praca domowa

- iteracja (np. pętle w programowaniu)
- heurestyka (głównie w informatyce)
- Alignment progresywny

Wstęp do bioinformatyki Wykład 3 Biotechnologia UWM Dr Jan Paweł Jastrzębski







Searching for remote homologs

- Sometimes BLAST isn't enough.
- Large protein family, and BLAST only gives close members. We want more distant members
- PSI-BLAST
- Profile HMMs

Wstęp do bioinformatyki Wykład 3 Biotechnologia UWM Dr Jan Paweł Jastrzębski

Profile

	1	2	3	4	5	6
A	1	0.67	0	0	.	.
T	0	0.33	1	1	.	.
C	0	0	0	0	.	.
G	0	0	0	0	.	.

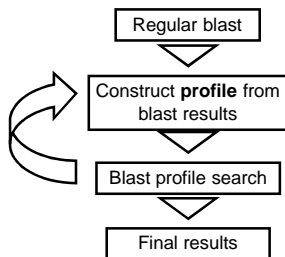
A T C T T G T
 A A C T T G T
 A A C T T C T

Profile =
 PSSM – Position Specific Score Matrix

Wstęp do bioinformatyki Wykład 3 Biotechnologia UWM Dr Jan Paweł Jastrzębski

PSI-BLAST

- Position Specific Iterated BLAST



Wstęp do bioinformatyki Wykład 3 Biotechnologia UWM Dr Jan Paweł Jastrzębski

PSI-BLAST

- zalety: PSI-BLAST looks for seq.s that are close to ours, and learns from them to extend the circle of friends
- wady: if we found a WRONG sequence, we will get to unrelated sequences (contamination). This gets worse and worse each iteration

Wstęp do bioinformatyki
Wykład 3

Biotechnologia
UWM

Dr Jan Paweł Jastrzębski

Profile HMM

- Similar to PSI-BLAST: also uses a profile
- Takes into account:
 - Dependence among sites (if site n is conserved, it is likely that site $n+1$ is conserved → part of a domain
 - The probability of a certain column in an alignment

Wstęp do bioinformatyki
Wykład 3

Biotechnologia
UWM

Dr Jan Paweł Jastrzębski

PSI BLAST vs profile HMM

PSI BLAST

Profile HMM

**Less exact
Faster**

**More exact
Slower**

Wstęp do bioinformatyki
Wykład 3

Biotechnologia
UWM

Dr Jan Paweł Jastrzębski
