

BIOINFORMATIKA V PRAXI – CVIČENÍ 2

SEQUENCE ALIGNMENT

STUDIJNÍ MATERIÁLY

Studijní materiály předmětu **C2130 Úvod do chemoinformatiky a bioinformatiky**, přednáška **Sequence alignment**.

VYUŽITÍ SEKVENČNÍHO PŘILOŽENÍ PRO IDENTIFIKACI GENU V ONLINE DATABÁZÍCH

Pro vyhledávání v internetových databázích lze použít několik přístupů (viz. Bioinformatika v praxi – cvičení 1). Pokud máme jako vstupní údaj sekvenci genu/proteinu, využíváme hledání na základě podobnosti sekvencí. Typickou ukázkou je aplikace **BLAST** na serveru NCBI (<http://blast.ncbi.nlm.nih.gov/Blast.cgi>).

ÚKOL 1

Pomocí aplikace BLAST identifikujte následující sekvence:

Sekvence 1:

```
MNGTEGPNFYVPFSNKTGVVRS PF EY P Q Y Y L A E P W Q F S M L A A Y M F L L I V L G F P I N F L T L Y V T V Q H K N V R T P L N Y I  
L L N L A V A N H F M V F G G F T T T L Y T S L H G Y F V F G S T G C N L E G F F A T L G G E I A L W S L V L A I E R Y V V C K P M S N F R F G E  
N H A I M G V A F T W M A L A C A A P P L V G W S R Y I P E G M Q C S C G I D Y Y T L K P E V N N E S F V I Y M F V V H F T I P M T I I F F C Y G Q  
L V F T V K E A A A Q Q Q E S A T T Q K A E K E V T R M V I I M V I A F L I C W V P Y A S V A F Y I F T H Q G S D F G P I L M T L P A F F A K S S A I  
Y N P V I Y I M M N K Q F R N C M L T T I C C G K N P F G E E E G S T T A S K T E T S Q V A P A
```

Název genu/proteinu:

Číslo záznamu v databázi NCBI:

Míra shody zadané a nalezené sekvence:

Sekvence 2:

```
atgtcgagcgttcagaccgctgccacttcgtggggaaccgtcccgtcgatccgtgtgtacacggccaataatggc  
aagatcaccgagcgttgctgggacgggaaggggtggtacacgggtgccttcaacgagcccggcgataacgtctcc  
gtgaccagctggctggtcggcagcgcgatccatatccgcgtctatgcaagcaccggcaccacgaccacagagtgg  
tgctgggacgggaacggctggaccaagggcgccctacaccgccaccaactga
```

Název genu/proteinu:

Číslo záznamu v databázi NCBI:

Míra shody zadané a nalezené sekvence:

Sekvence 3:

```
MGLSDGEWQMVLNIWGKVEGDLAGHGQEVLI SLFKAHPETLEKFDKFKNLKSEEEEMKSS E D L K K H G C T V L T A L G T  
I L K K K G Q H A A E I Q P L A Q S H A T K H K I P V K Y L E F I S E V I I Q V L K K R Y S G D F G A D A Q G A M S K A L E L F R N D I A A K Y K E L  
G F Q G
```

Název genu/proteinu:

Číslo záznamu v databázi NCBI:

Míra shody zadané a nalezené sekvence:

VYHLEDÁNÍ PODOBNÝCH SEKVENCÍ A URČENÍ PŘÍBUZNOSTI

Výhodou použití sequence alignmentu je schopnost nalezení nejen shodného záznamu, ale i záznamů podobných. Tak lze na základě podobných sekvencí identifikovat i dosud neznámou sekvenci a odhadnout její „příbuzenské“ vztahy.

ÚKOL 2

Nalezněte 4 nejpodobnější sekvence k sekvenci zadané. Použijte aplikaci BLAST.

Sekvence:

```
MNTRSFHRIDVHKARELLQRPDTVLLDCRHPSTDFRAGHIAGASPLGDYNADDHVLNIAKHRPVLIYCYHGNASQM
RAQLFADFGFAEVYSLDGGYEAWRKVHTPANSQLEALQCWLMAQEFPAADIHARTRDGVTPLMRAAGEGDPARV
AELLAAGADPHQRNNDGNQALWFACVSENLDTLDLLVAVGAHLNHQNDNGATCLMYAASAGKTAVVERLLAFGAD
RSLLSLDDFTALDMAANLECLNLLRETPRIKAVT
```

Číslo záznamu	Protein	Organismus	Score

VLIV POUŽITÉ MATICE NA VÝSLEDEK ALIGNMENTU

Dalším parametrem, který může ovlivnit výsledek alignmentu je použitá matice. Většina programů detekuje automaticky nukleotidovou sekvenci a použije příslušnou matici, v případě proteinových sekvencí je však situace komplikovanější.

ÚKOL 3

Následující sekvence identifikujte a přiložte v programu **ClustalW** (<http://www.ebi.ac.uk/Tools/clustalw2/index.html>). Poté proveďte alignment s použitím matice PAM, BLOSUM, Gonnet a ID a výsledky porovnejte. Která matice je použita při základním nastavení?

Sekvence 1:

```
MPGIRLRYALLALVFAYIYYYIVLSYRDQFSDIKKCFSSIRAKIEDYDSSKKGQPKLASNSYLEADMLYRD
RTQVG IENATMVMLVRNRELEGALSSMRSLED RFN RQYKYPWVFLNDEPFTEEFIEQTMLMASSQTFYEL
IPSSDWNMPDFIDNERVEQNIANSTDV IYGFSKSYRNMCHFNSGYFYKQKRLNLDWYFRVEPDVEYMCD
FQYDPFTLLRTNNKIYGFVIAIHEYENTIPTLWPTVEKFMQTYPDLIHANNSLKFITTTNESSLNHWVTP I
QSSSDYNLCHFWSNFEIGNLNFFRGEAYNKYFDFLDRAGGFYERWGDAPVHSIGLNLLADKNSIHHFED
IGYYHPPYLACPSSKDVI AAKRCVCRKRGNDGEVIDSAIDVNVFSCLSRWWRYGAGKRFLNEIDYTFNN
```

Sekvence 2:

```
MPESGVPAGRRRLAPLLLVTSA AALLAVTMANPPDLVDFHVYMLGGVALDRPDTLYSFAYS DQSPDQPL
PFVYPPFAAILFRPLTALPFV VAGVLWQLGILAAVYGVIRISQRLGGGSHRTAMLW TAGLIWLEPVRVA
LDLGQVGI FLTLAVLYAVCSRWWLSGLLVGLAAGVKLTPAITGLYFLGVRRWTA AAFSAVVFVATIGLS
YLVVGDQVRNYFTRVMGDT SINPIGIALNQSWRGAISRFLGHDAGGSALVIAA IAGTAVLALFAWAALGA
GSRPRDRLGSVLVVQLFGLLMSPI SWVHHWAWVVPLVLWLLSGPWRNEPGARVLGWGWLALTFVGVPSQL
SLLEPSLWEISR PWYLAWAGLAYVVA AVATLGWIVLTGRRNAAPT PPVRRRFRARVVGTRGKQIACEEHRA
GR
```

Sekvence 3:

MELQSLIDTVSLQKLLLLGALLRLILIAAFFHDQWFRVKYTDIDYMIVVDGARHMWNGGSPFDRTTFRY
TPLLAALVMPISIWIANPMGKLIFASSDLGAAWYCYGVLKSFAKERSAKWMVSLFILFNPIVLSVSTRGNS
DMLVTFMSLMVLSKFFARRKCYQAAAVLGFVAVHFKIYPIIYALPLTLGVWEQSVAASTNTWRRVVKTAVVV
SICALMAAISFAVPTVLCYMKYQQYLNEAFIYHVYREDHRHNFSPYWLLMYLNMARRHLGQGVDFSPRL
VAFAPQAVVLSFVSYKLRNRTAHACCVQTVLFAFNKVCTVQYFVWFIPFLAFLFCEPKEVEDEDESGGSG
AFKFFSWVKALGVVLMWAATIPLWVTTAVPLEFHGYSDFACLWIVSCLFFLAMVVLASMLARIAYRVQCT
KCSAKSIKVA

Sekvence 4:

MDASTPNPPTLGTDAVASTIVFFHPDLGIGGAERLVVDAAVGLQTRGHKVVIFTNHCDPTHCFDECRDGT
LDVVRVGRNSIVPPSIFSRILTILCAILRHIHLLLTIHILTGELAALSPRAFIVDQLSAGLPLMRFLAPDVPV
LFYCHFPDLLLAQGRQSLVKRLYRVFPDRLEEWSMGFAHAVAVNSKFTRGIVGNTWPALQNKVIPINVVYP
CVDTHTTTHETAPDEAKLAAGKKLILSINRFERKKDIGLAIRAFQAQIPEEQRRGARLVLAGGYDARVSENV
LYHAELQALATSLSLAHHTLTPAELGSAAPPDAQHFGIVPLEAMLARVPVLAANTGGPVETVADRETGW
LRDPADAPAWTDVMARCLALPDDQLAAMGDAGRRRVRELFGRDKMAQTLDESLVQIAGLAEERRVSGAAG
FGVLAAFIAACAVLAAWFAF

VLIV MEZER NA VÝSLEDEK ALIGNMENTU

Možnost vložení mezer významně zvyšuje šance na úspěšný alignment sekvencí. Při změně nastavení parametrů se mění skóre určující podobnost sekvencí a ty tak mají vliv i na určení vzájemné příbuznosti sekvencí. Při špatném nastavení pak umožňují provést alignment i u naprosto nepodobných sekvencí.

ÚKOL 4

Proveďte multiple alignment následujících sekvencí pomocí programu ClustalW. V prvním případě nastavte parametr GAP OPEN na 1, ve druhém případě na 100 a výsledky porovnejte.

Sekvence 1:

trypsin [*Homo sapiens*] – AAZ40216.1

RIQVRLGEHNIEVLEGNEQFINAAKIRHHPKYDRKTLNNDIMLIKLSRAVINARVSTISLPTAPPATGTKCLIS
GWGNTASSG

Sekvence 2:

FlgA [*Agrobacterium tumefaciens*] – AAB71791.1

MRFGRNSSCRTALVRMCLASAFSLGALAPALQAPMALVPVRTIYPGEAISPEQVKSVEVTNPNISAGYASDIS
EVEGMISKQTLTPGRTIPIAALREPSLVVRGTSVKLVFHIIGNMTLMASGTPMSDGLGEVVRVRNIDSGVMVSGT
VMKDGTIQVMAK

Sekvence 3:

collagen [*Caenorhabditis elegans*] – CAA35955.1

MSEDLKQIAQETESLRKVAFFFGIAVSTIATLTAI IAVPMLYNYMQHVQSSSQSEVEFCQHRNGLWDEYKRFQGV
SGVEGRIKRDAYHRSLGVSGASRKARRQSYGNDAVGGFGGSSGGSCCGSGAAGPAGSPGQDAGPNDGAPGA
PGNPGQDASEDQTAGPDSFCFDCPAGPPGPGSAPGQKGPSGAPGAPGQSGGAALPGPPGAGPPGAGQPGSNGN
AGAPGAPGQVVDVPGTPGPAGPPGSPGPAGAPGQPGQAGSSQPGGPGPQGDAGAPGAPGAPGQAGAPGQDGESGS
EGACDHCPPPRTAPGY

Gap open = 1

Vzájemně příbuznější sekvence:

Gap open = 100

Vzájemně příbuznější sekvence:

ALIGNMENT NA GENOVÉ vs. PROTEINOVÉ ÚROVNI

Často se setkáváme se situací, kdy alignment na genové úrovni není pro naše potřeby vhodný. Je tedy zapotřebí výsledné sekvence porovnat i na úrovni proteinu.

ÚKOL 5

U následujících dvojic sekvencí proved'te sequence alignment na genové úrovni (program **lalign** – http://www.ch.embnet.org/software/LALIGN_form.html). Tyto sekvence přeložte do sekvence aminokyselin programem **Translate** – server ExPassy (<http://www.expasy.ch/tools/dna.html>) a proved'te alignment těchto – přeložených sekvencí. Porovnejte množství nespárovaných nukleotidů/aminokyselin (resp. procento identity) v obou případech.

Sekvence a1

```
atg tgc gca gtg cgc agg gcc ggc tcg aag cgc aag cag gaa gcg ttt gcg gtg atc
ccg gcg act gcg ctg gct aat gcg gta ccg gct agc gtg gct tct gca ccc cgc act
gcc cag cac ctt ccg ctc agt cct cgg cgc ccg cct gcc gcc tcc gga gcg cct gtg
tgg ttt cca aaa aaa gac tta cag caa aat gaa tac tcc agc cat caa gag aat agg
aaa tca cat tac caa gtc tcc tga
```

Sekvence a2

```
atg tgt gca gtg cgc cga gcc ggc tcc aag agg aag caa gaa gcg ttt gcg gtt atc
ccg gcg act gct ctg gct aat gca gta ccg gct agc gtg gct tct gca ccg cgc act
gcc cag cat tta ccg ctg agt cct cgc cgg ccg cct gca gct tcc gga gcg cca gtg
tgg ttc cca aaa aaa gat ttg cag caa aat gaa tat tcc agc cac cag gag aat agg
aag tcc cat tac caa gtc tca tga
```

Identita nt sekvencí a1-a2:

Identita ak sekvencí a1-a2:

Sekvence b1

```
atg tgc gca gtg cgc agg gcc ggc tcg aag cgc aag cag gaa gcg ttt gcg gtg atc
ccg gcg act gcg ctg gct aat gcg gta ccg gct agc gtg gct tct gca ccc cgc act
gcc cag cac ctt ccg ctc agt cct cgg cgc ccg cct gcc gcc tcc gga gcg cct gtg
tgg ttt cca aaa aaa gac tta cag caa aat gaa tac tcc agc cat caa gag aat agg
aaa tca cat tac caa gtc tcc tga
```

Sekvence b2

```
atg tgc gca gtg cgc agg gcc ggc tcg aag cgc aag cag gaa gcg tgt gcg gtg atc
ccg gcg act gcg ctg gct aat gcg gaa ccg gct agc gtg gct tct gca ccc cgc act
gcc cag cac ctt ccg ctc agt cct cgg cgc ccg cct gcc gcc tcc cga gcg cct gtg
tgg ttt cca aaa aaa gac tta cag caa aat gaa tac tcc agc cat taa gag aat agg
aaa tca cat tac caa gtc tcc tga
```

Identita nt sekvencí b1-b2:

Identita ak sekvencí b1-b2:

VYUŽITÍ ALIGNMENTU PRO INTERPRETACI VÝSLEDKŮ SEKVENACE

Běžným užitím sequence alignmentu je analýza výstupu po sekvenaci. Detekujeme tak mutace (inzerce, delece, substitute), které mohou mít vliv na sekvenci kódovaného proteinu – záměna aminokyseliny, posunutí čtecího rámce, vytvoření nebo odstranění STOP kodonu, atd. Můžeme aplikovat pairwise alignment nebo u více sekvencí multiple alignment.

ÚKOL 6

Následující sekvence obsahují inzerce. Určete, která z obou sekvencí je vhodnější pro budoucí práci s proteinem a proč. Pro alignment použijte vámi zvolený program (lalign, ClustalW, případně jiný).

Původní gen:

```
atggctgattctcaaacgcatccaaccgcgccggcgaattctcgattccgcccgaataccgatttccgcgcgatt
ttcttcgcgaatgccgcccagcaacagcacatcaaaattgttcacgcccagccaggaacccgcccgcgtatcac
aagctgacgacgcgcgacggcccgcgcaagccacgctgaattccggcaacggcaagatccgtttcgaggtgtcg
gtgaacggcaagccgctcggcgaccgacgcgcgctctcgcgccgatcaacggcaagaagtcggacggctcgcggtc
acggtcaacttcgggatcgtcgtgtcgggaagacggccacgacagcgactacaacgacggcatcgtcgtgctccag
tggccgatcggctga
```

Sekvence 1:

```
atggctgattctcaaacgcatccaaccgcgccggcgaattctcgattccgcccgaataccgatttccgcgcgatt
ttcttcgcgaatgccgcccagcaacagcacatcaaaaattgttcacgcccagccaggaacccgcccgcgtatcac
cacaagctgacgacgcgcgacggcccgcgcaagccacgctgaattccggcaacggcaagatccgtttcgaggtg
tcggtgaacggcaagccgctcggcgaccgacgcgcgctctcgcgccgatcaacggcaagaagtcggacggctcgcg
ttcacggtaacttcgggatcgtcgtgtcgggaagacggccacgacagcgactacaacgacggcatcgtcgtgctc
cagtggccgatcggctga
```

Sekvence 2:

```
atggctgattctcaaacgcatccaaccgcgccggcgaattctcgattccgcccgaataccgatttccgcgcgatt
ttcttcgcgaatgccgcccagcaacagcacatcaaaattgttcacgcccagccaggaacccgcccgcgtatcac
aagctgacgacgcgcgacggcccgcgcaagccacgctgaattccggcaaacggcaagatccgtttcgaggtgtc
cgggtgaacggcaagccgctcggcgaccgacgcgcgctctcgcgccgatcaacggcaagaagtcggacggctcgcg
tcacggtaacttcgggatcgtcgtgtcgggaagacggccacgacagcgactacaacgacggcatcgtcgtgctcc
agtggccgatcggctga
```

Vhodnější sekvence:

Důvod:

ÚKOL 7

Následující sekvence obsahují různé mutace. Určete, které z těchto sekvencí jsou použitelné pro budoucí práci s proteinem a proč. Označte nejvhodnější sekvenci.

Originální sekvence:

```
atggctgattctcaaacgcatccaaccgcgccggcgaattctcgattccgcccgaataccgatttccgcgcgatt
ttcttcgcgaatgccgcccagcaacagcacatcaaaattgttcacgcccagccaggaacccgcccgcgtatcac
aagctgacgacgcgcgacggcccgcgcaagccacgctgaattccggcaacggcaagatccgtttcgaggtgtcg
gtgaacggcaagccgctcggcgaccgacgcgcgctctcgcgccgatcaacggcaagaagtcggacggctcgcggtc
acggtcaacttcgggatcgtcgtgtcgggaagacggccacgacagcgactacaacgacggcatcgtcgtgctccag
tggccgatcggctga
```

Sekvence 1:

atggctgattctcaaacgtcatccaaccgcgccggcgagttctcgattccgccaataaccgatttccgcgcgatt
 ttcttcgcgaatgccgcccagcaacagcacatcaaattggtcatcggcgacagccaggaaccagccgcgtatcac
 aagctgacgacgcgcgacggcccgcgcgaagccacggtgaattccggcaacggcaagatccggttccgaggtgctcg
 gtgaacggcaagccgctcggcgaccgacgcgcgtctcgcgccgatcaacggcaagaagtcggacggctcgcggttc
 acggtcaacttcgggatcgtcgtgctcgggaagacggccacgacagcgcgactacaacgacggcatcgtcgtgctccag
 tggccgatcggctga

Sekvence 2:

atggctgattctcaaacgtcatccaaccgcgccggcggaattctcgattccgccaataaccgatttccgcgcgatt
 ttcttcgcgaatgccgcccagcaacagcacatcaaattggtcatcggcgacagccaggaaccagccgcgtatcac
 aagctgacgacgcgcgacggcccgcgcgtagccacgctgaattccggcaacggcaagatccggttccgaggtgctcg
 gtgaacggcaagccgctcggcgaccgacgcgcgtctcgcgccgatcaacggcaagaagtcggacggctcgcggttc
 acggtcaacttcgggatcgtcgtgctcgggaagacggccacgacagcgcgactacaacgacggcatcgtcgtgctccag
 tggccgatcggctga

Sekvence 3:

atggctgattctcaaacgtcatccaaccgcgccggcggaattctcgattccgccaataaccgatttccgcgcgatt
 ttcttcgcgaatgccgcccagcaacagcacatcaaattggtcatcggcgacagccaggaaccagccgcgtatcac
 aagctgacgacgcgcgacggcccgcgctaagccacgctgaattccggcaacggcaagatccggttccgaggtgctcg
 gtgaacggcaagccgctcggcgaccgacgcgcgtctcgcgccgatcaacggcaagaagtcggacggctcgcggttc
 acggtcaacttcgggatcgtcgtgctcgggaagacggccacgacagcgcgactacaacgacggcatcgtcgtgctccag
 tggccgatcggctga

Sekvence 4:

aaggctgattctcaaacgtcatccaaccgcgccggcggaattctcgattccgccaataaccgatttccgcgcgatt
 ttcttcgcgaatgccgcccagcaacagcacatcaaattggtcatcggcgacagccaggaaccagccgcgtatcac
 aagctgacgacgcgcgacggcccgcgcgaagccacgctgaattccggcaacggcaagatccggttccgaggtgctcg
 gtgaacggcaagccgctcggcgaccgacgcgcgtctcgcgccgatcaacggcaagaagtcggacggctcgcggttc
 acggtcaacttcgggatcgtcgtgctcgggaagacggccacgacagcgcgactacaacgacggcatcgtcgtgctccag
 tggccgatcggctga

Sekvence	Charakter mutace z hlediska genu	Charakter mutace z hlediska proteinu	Použitelná pro další práci (ANO/NE)
1			
2			
3			
4			

PROBLÉM REPETIC

Při porovnávání dvou celkově podobných sekvencí užíváme zpravidla metody globálního alignmentu. V případě sekvencí, které jsou podobné jen v určité své části (např. jedné z domén), je vhodnější použít lokální alignment. Ten má svůj význam i v případě proteinů s tzv. repetitivy.

ÚKOL 8

Proveďte alignment následujících dvou sekvencí programem Align (<http://www.ebi.ac.uk/Tools/emboss/align>) s použitím algoritmu needle (globální alignment) a water (lokální alignment). V obou případech nastavte parametr Gap open na 15.0 a výsledky porovnejte.

Sekvence 1

PTEFLYTSKIAAISWAATGGRQQRVYFQDLNGKIREAQRGGDNPWTGGSSQNVIGEAKLFSPLAAVTWKSAQGIQ
 IRVYCVNKDNLSEFVYDGSKWITGQLGSGVGVKVGSNKLAALQWGGSESAPPNIRVYYQKSNNGSGSSIHEYVWS

GKWTAGASFGSTVPGTGIGATAIGPGRLRIYYQATDNKIREHCWDSNSWYVGGFSASASAGVSIAAISWGSTPNI
RVYWQKGREELYEAAYGGSWNTPGQIKDASRPTPSLPDTFIAANSSGNIDISVFFQASGVSLQQWQWISGKGWSI
GAVVPTGTPAGW

Sekvence 2

SSVQTAATSWGTVPSIRVYTANNNGKITERCWDGKGWYTGAFNPEPGDNVSVTSWLVGSAIHIRVYASTGTTTTEWC
WDGNGWTKGAYTATN

	Identické ak	Podobné ak	Mezery
Needle			
Water			

Výše uvedené sekvence jsou příkladem repetice, tj. opakujících se podobných (homologních) úseků v rámci jedné sekvence. Přítomnost repetice lze zjistit/ověřit programem **RADAR** (<http://www.ebi.ac.uk/Tools/Radar/>).

ÚKOL 9

V sekvencích z úkolu 8 detekujte repetice pomocí programu Radar. Uveďte počet repetice zjištěných u každé sekvence:

Sekvence 1:

Sekvence2:

Sekvenci s více repeticemi rozdělte na jednotlivé repetice a proveďte multiple alignment pomocí programu ClustalW. Která z residuí jsou v repeticích konzervována (zcela, částečně)? Využijte tzv. consensus.