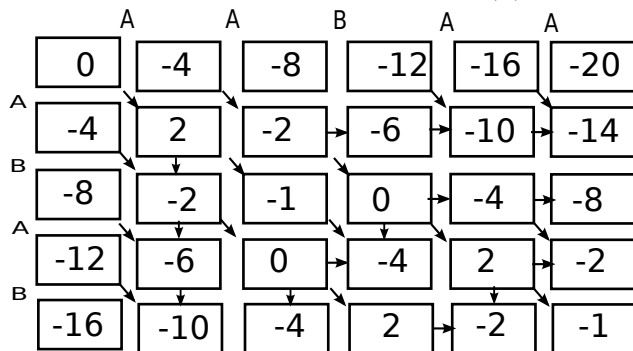


- Proszę nie edytować tego pliku. Tutaj są tylko treści zadań.
- Odpowiedzi proszę umieszczać w pliku o nazwie 'mbi-nazwisko.txt'. Proszę pobrać ten plik, zmienić mu nazwę wpisując własne nazwisko bez użycia polskich znaków. Dodatkowo proszę wpisać swoje imię i nazwisko w pierwszej linii pliku.
- Na koniec proszę pobrać i wypełnić oświadczenie (plik 'oświadczenie-nazwisko.doc').
- Egzamin składa się z 5 zadań. Treść każdego zadania jest na oddzielnej stronie.

Zadanie 1 (8 pkt)

Poniżej przedstawiono macierz dla algorytmu, który dostarcza uliniowienia globalnego (algorytm Needlemana-Wunscha) dwóch sekwencji. Macierz podobieństwa pokazano obok. Stosujemy liniową karę za przerwę, $\gamma(n) = n * d, d = -4$.

	A	B
A	2	-3
B	-3	2

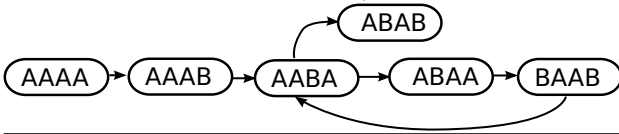


Pytania:

- algorytm uruchomiono dla sekwencji AABAA i ...
- rozwiązań optymalnych jest ...
- rozwiązania optymalne to ...

Zadanie 2 (8 pkt)

Dla zestawu odczytów zbudowano graf de Bruijna 5 rzędu pokazany na rysunku (nie jest to multi-graf, ani A-graf, nie ma wag w krawędziach).



Pytania:

- Czy istnieje rozwiązanie (ścieżka Eulera)?
- Jeżeli tak, podaj sekwencję odpowiadającą jednemu z rozwiązań:
- Zbuduj graf de Bruijna 4 rzędu (bez wag w krawędziach). Ile ten graf ma wierzchołków?
- Czy istnieje (dla grafu de Bruijna 4 rzędu) rozwiązanie (ścieżka Eulera)?
- Jeżeli tak, to podaj sekwencję odpowiadającą jednemu z rozwiązań:

Zadanie 3 (8 pkt)

Cząsteczka jest biopolimerem i może być reprezentowana przez napis nad alfabetem $\{A, B\}$. Obserwujemy oddziaływanie poszczególnych symboli z naszą sondą i oznaczamy je jako '+' lub '-'. Chcemy opisać naszą cząsteczkę oraz oddziaływanie za pomocą ukrytego modelu Markowa. Stanem będzie symbol, więc $Q = \{A, B\}$, obserwacją oddziaływanie, więc $V = \{+, -\}$. Dla cząsteczki $AAAAAAAAAA$ obserwujemy oddziaływanie $+ - + + + - - - -$, dla cząsteczki $BBBB$ obserwujemy $+ - + +$. Nasze sekwencje zaczynają się od z takim samym prawdopodobieństwem od A i od B, więc $P_A = \frac{1}{2}, P_B = \frac{1}{2}$. Znamy macierz przejść:

	A	B
A	$\frac{4}{5}$	$\frac{1}{5}$
B	$\frac{1}{5}$	$\frac{4}{5}$

Pytania:

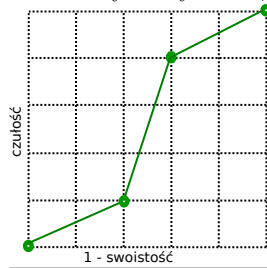
- podaj macierz emisji:

	+	-
A		
B		

- podaj najbardziej prawdopodobną sekwencję cząsteczki dla oddziaływania $+ + +$:

Zadanie 4 (8 pkt)

Opracowano test binarny X, dla którego krzywa ROC jest pokazana na rysunku. Test zwraca wynik 'tak' lub 'nie'. Używamy testu, gdzie obiekty mają w 50% stan 'tak' i w 50% stan 'nie'.



Pytania:

- czy test może zwracać zawsze wynik 'nie' ?
- jaka jest wtedy czułość?
- jaka jest czułość dla swoistości 0.4?
- Jaki będzie błąd przy czułości 0.2?
- Jaki będzie błąd przy czułości 0.8?

Zadanie 5 (8 pkt)

Mamy dane pokazane poniżej i zakładamy pierwszą ramkę odczytu (pierwszy kodon odpowiada pierwszemu aminokwasowi).

a) fragment sekwencji kodującej genu CLP1 w genomie referencyjnym:

```
>hg38_dna range=chr11:57659477-57659493 5'pad=0 3'pad=0 strand=+
ATGGGAGAAGAGGCT
```

b) sekwencję pacjenta:

```
> patient1
ATGGGAGTAGAGGCC
```

Tablica kodowania aminokwasów:

		Second base					
		U	C	A	G		
First base	U	UUU } Phenylalanine F UUC } UUA } Leucine L UUG }	UCU } Serine S UCC } UCA } UCG }	UAU } Tyrosine Y UAC } UAA } Stop codon UAG } Stop codon	UGU } Cysteine C UGC } UGA } Stop codon UGG } Tryptophan W	U	C
	C	CUU } Leucine L CUC } CUA } CUG }	CCU } Proline P CCC } CCA } CCG }	CAU } Histidine H CAC } CAA } Glutamine Q CAG }	CGU } Arginine R CGC } CGA } CGG }	C	A
	A	AUU } Isoleucine I AUC } AUA } AUG } Methionine M start codon	ACU } Threonine T ACC } ACA } ACG }	AAU } Asparagine N AAC } AAA } Lysine K AAG }	AGU } Serine S AGC } AGA } Arginine R AGG }	A	G
	G	GUU } Valine V GUC } GUA } GUG }	GCU } Alanine A GCC } GCA } GCG }	GAU } Aspartic acid D GAC } GAA } Glutamic acid E GAG }	GGU } Glycine G GGC } GGA } GGG }	G	
						U	C
						A	G
						G	

Polecenia:

- Dokonaj translacji sekwencji referencyjnej DNA na sekwencję aminokwasów
- Dokonaj translacji sekwencji DNA pacjenta na sekwencję aminokwasów
- Znajdź warianty SNV na poziomie sekwencji DNA i zapisz je w formacie:
 - CHR:POS_REF>ALT, np: chr11:57659477_A>T
- Określ efekt każdego wariantu jak będzie on miał dla sekwencji białkowej, nadając im jedną z 4 etykiet: 'missense', 'synonymous', 'stopgain', 'stoploss'