

UNIVERZITA KOMENSKÉHO, BRATISLAVA

FAKULTA MATEMATIKY, FYZIKY A INFORMATIKY

ZAROVNÁVANIE SEKVENCIÍ S POUŽITÍM METÓD
KLASIFIKÁCIE

Diplomová práca

2014

Bc. Michal Hozza

UNIVERZITA KOMENSKÉHO, BRATISLAVA

FAKULTA MATEMATIKY, FYZIKY A INFORMATIKY

ZAROVNÁVANIE SEKVENCIÍ S POUŽITÍM METÓD KLASIFIKÁCIE

Diplomová práca

Študijný program: Informatika
Študijný odbor: 2508 Informatika
Školiace pracovisko: Katedra Informatiky
Školiteľ: Mgr. Tomáš Vinař, PhD.
Konzultant: Mgr. Michal Nánási

Bratislava, 2014

Bc. Michal Hozza

Podakovanie...

Bc. Michal Hozza

Abstrakt

Zarovnávanie dvoch DNA sekvencií je jedným zo základných bioinformatických problémov. Obvykle takéto zarovnanie hľadáme pomocou jednoduchých párových skrytých Markovovských modelov (pHMM). V tejto práci sa zaoberáme možnosťami použitia prídavnej informácie o funkcii vstupných sekvencií na zlepšenie kvality takýchto zarovnaní. Informácie sme zakomponovali pomocou klasifikátorov, ktoré rozhodujú či dané pozície majú byť zarovnané k sebe alebo nie. Ako klasifikátor sme použili Random forest.

Ukázalo sa, že klasifikátor sa dokáže naučiť, ktoré okná majú byť zarovnané k sebe a ktoré nie.

Vyvinuli sme 2 modely pre zarovnanie sekvencií s anotáciami za pomoci klasifikátora, ktoré sú založené na párových skrytých Markovovských modeloch.

Model s klasifikátorom ako emisiou, kde sme nahradili emisné tabuľky stavov výstupom z klasifikátora.

Model s klasifikátorovou páskou, kde navyše modelujeme aj výstup z klasifikátora.

Kľúčové slová: zarovnávanie sekvencií, strojové učenie, Random forest, anotácie

Abstract

English abstract

Key words: ...

Obsah

Úvod	1
1 Zarovnávanie sekvencií	3
1.1 Podobnosť sekvencií, sekvenčná homológia a zarovnanie	3
1.2 Párové zarovnávanie	5
1.3 Typy zarovnaní	5
1.4 Skórovacie systémy	5
1.4.1 Skórovacie matice	6
1.5 Algoritmy na hľadanie zarovnaní	6
1.5.1 Algoritmus pre globálne zarovnanie: Needleman-Wunch	7
1.5.2 Algoritmus pre lokálne zarovnanie: Smith-Waterman	9
1.5.3 Afínne skórovanie medzier	10
1.6 Zarovnávanie pomocou skrytých Markovovských modelov	11
1.6.1 Zarovnávanie pomocou párových skrytých markovovských mode- lov (pHMM)	12
1.6.2 Viterbiho algoritmus na pHMM	12
1.6.3 Nastavenie parametrov pHMM	14
1.7 Štatistická významnosť zarovnania	14
2 Súvisiaca práca	17
3 Klasifikácia na základe lokálnej informácie	19
3.1 Náhodné lesy	19
3.1.1 Klasifikačné stromy	20
3.1.2 Náhodný les	22

3.2	Použitie náhodných lesov na klasifikáciu zarovnaní	23
3.2.1	Výber atribútov	23
3.2.2	Trénovanie a testovanie klasifikátorov	25
3.3	Výsledky experimentov	27
3.3.1	Dôležitosť atribútov	27
3.3.2	Úspešnosť klasifikátora	30
3.3.3	Nezávislosť typu dát na veľkosti okna	31
3.4	Zhrnutie	33
4	Modely	34
4.1	Integrácia párových HMM a výsledkov z klasifikátorov	34
4.1.1	Definície modelov	35
4.1.2	Diskretizácia skóre klasifikátorov	37
4.2	Trénovanie modelov	39
4.3	Výsledky experimentov	40
4.3.1	Vhodná veľkosť okna	42
4.3.2	Vplyv dodatočnej informácie na zarovnanie	42
4.3.3	Použitie jedného klasifikátora pre Match aj Inzert stav	43
4.4	Zhrnutie	44
5	Implementácia	45
5.1	Výber jazyka a použité knižnice	45
5.1.1	Realigner	45
5.1.2	Pythonové knižnice	46
5.2	Triedy na zarovnávanie s klasifikátorom	46
5.2.1	Predspracovanie dát	46
5.2.2	Abstrakcia klasifikátora	47
5.2.3	HMM stavy s klasifikátorom	48
5.3	Pomocné programy	49
5.3.1	Simulátor	49
5.3.2	Trénovanie modelov	51
5.4	Použitie	52
5.4.1	Konfigurácia a Formáty súborov	54

5.5 Rozšírenie programu	57
Záver	58
Literatúra	59

Zoznam obrázkov

1.1	Lokálne zarovnanie	4
1.2	Skórovacia matica	6
1.3	Tabuľka dyn. programovania pre globálne zarovnanie	8
1.4	Tabuľka dyn. programovania pre lokálne zarovnanie	9
1.5	Situácie pri afínnoom skórovaní	10
1.6	Stavový diagram pre zarovnanie sekvencií	11
1.7	Párový HMM pre zarovnávanie sekvencií	13
1.8	P-hodnota lokálneho zarovnaní	15
3.1	Klasifikačný strom	21
3.2	Okno klasifikátora	24
3.3	Dôležitosť atribútov pre typ dát D	29
3.4	Distribúcia výstupu z klasifikátora pri type D	30
3.5	Horšie distribúcie výstupov z klasifikátora	31
4.1	Použité klasifikátory v klasifikátorovej páske	35
4.2	Model s klasifikátorom ako emisiou	36
4.3	Model s klasifikátorovou páskou	37
4.4	Diskretizácia výstupu klasifikátora	38
5.1	Markovova reťaz použitá na generovanie informácie o génoch	50

Zoznam tabuliek

3.1	Porovnanie vlastností klasifikátorov	28
3.2	Porovnanie úspešnosti klasifikátorov pri rôznej veľkosti okna	32
4.1	Porovnanie úspešností pri rôznom spracovaní pásky	38
4.2	Porovnanie s existujúcimi zarovnávačmi	41
4.3	Porovnanie úspešností pri rôznej veľkosti okna	42
4.4	Vplyv dodatočnej informácie na zarovnanie	43
4.5	Porovnanie použitia jedného alebo dvoch typov klasifikátorov	43

Úvod

Najnovšie technológie sekvenovania DNA produkujú stále väčšie množstvo sekvencií rôznych organizmov. Spolu s tým stúpa aj potreba rozumieť týmto dátam, čomu napomáha *zarovnávanie sekvencií*. Zarovnávanie dvoch DNA sekvencií je jedným zo základných bioinformatických problémov. DNA sekvencia je reťazec, ktorý sa skladá z jednotlivých báz: adenín (A), guanín (G), cytozín (C) a tymín (T) [Wik14a]. Správne zarovnanie identifikuje časti sekvencie, ktoré vznikli z toho istého predka (zarovnané bázy), ako aj inzercie a delécie v priebehu evolúcie (medzery v zarovnaní). Je nápomocné pri zisťovaní ich štruktúry a následne funkcie jednotlivých častí.

Existujú rôzne algoritmy na zarovnávanie sekvencií. Algoritmy sú zvyčajne založené na dynamickom programovaní a pracujú v kvadratickom čase v závislosti od dĺžok sekvencií. Niekedy sa na urýchlenie použijú rôzne heuristické algoritmy, ktoré nie vždy nájdu najpravdepodobnejšie zarovnanie, ale pracujú oveľa rýchlejšie. Jeden z nich je aj algoritmus na zarovnávanie sekvencií pomocou jednoduchých párových skrytých Markovovských modelov (pHMM) [DEKM98], kde kvalita výsledného zarovnania je ovplyvnená len pravdepodobnostným modelom. Základný model berie do úvahy len jednotlivé bázy a pravdepodobnosti substitúcie (mutácie), inzercie a delécie.

My sme sa v práci zaoberali návrhom modelov pre tento algoritmus. Naše modely navyše uvažujú aj prídavné informácie (takzvané anotácie) získané z externých programov (napr. anotácie o génoch z vyhľadávača génov). Na zakomponovanie tejto informácie sme sa rozhodli využiť klasifikátory¹, ktoré sme trénovali na sekvenciách so známym zarovnaním a potom použili na zarovnanie nových sekvencií. Naše klasifikátory vracajú čísla z intervalu $\langle 0, 1 \rangle$, ktoré určujú, či dané dve bázy majú byť zarovnané spolu.

Ako klasifikátor sme použili *náhodný les* [Bre01], pretože aktuálne patrí medzi naj-

¹klasifikátor je program, ktorý na základe vopred natrénovaných parametrov klasifikuje vstupné dáta do niektorej triedy z danej množiny tried

lepšie klasifikátory.

V kapitole 1 popíšeme algoritmy na zarovnávanie a modely, z ktorých budeme vychádzať. V kapitole 2 si predstavíme súvisiace práce. Potom si predstavíme naše riešenie, ktoré sme v našej práci rozdelili do dvoch častí. V kapitole 3 si povieme niečo o náhodných lesoch a predstavíme si klasifikátory, ktoré použijeme v našich modeloch. Povieme si niečo voľbe atribútov a vlastnostiach klasifikátora. V kapitole 4 si predstavíme dva modely s klasifikátorom na zarovnávanie sekvencií s dodatočnou informáciou a porovnáme ich úspešnosť s referenčnými modelmi. Nakoniec si v kapitole 5 povieme niečo o tom ako sme implementovali náš zarovnávač.

1 Zarovnávanie sekvencií

V tejto kapitole si stručne popíšeme čo je to globálne a lokálne zarovnanie a ukážeme základné algoritmy na hľadanie globálneho a lokálneho zarovnania. Tieto algoritmy budeme neskôr modifikované používať pri našom riešení.

1.1 Podobnosť sekvencií, sekvenčná homológia a zarovnanie

V prírode vznikajú evolúciou nové sekvencie modifikáciou už existujúcich. Preto môžeme často spozorovať podobnosť medzi neznámou sekvenciou a sekvenciou o ktorej už niečo vieme. Ak zistíme podobnosti medzi sekvenciami, môžeme preniesť informácie o štruktúre a/alebo funkcii na novú sekvenciu.

Podobné sekvencie, ktoré sa vyvinuli mutáciami zo sekvencie v spoločnom predkovi sa nazývajú *homologické* a pod pojmom *hľadanie homológov* rozumieme hľadanie takých podobností, ktoré s veľkou pravdepodobnosťou vznikli práve spoločnou evolučnou históriou.

Počas evolúcie dvoch homologických sekvencií nastane veľa *inzercií*, *delécií* a *substitúcií*, preto predtým ako môžeme začať porovnávať sekvencie, ich musíme zarovnať tak, aby homologické časti sekvencií boli na rovnakom mieste v zarovnaní. [DEKM98, BV11]

1.1. PODOBNOSŤ SEKVENCÍÍ, SEKVENCIONÁLNA HOZADIVOSŤ A ZAROVNAVANIE

Sekvencia 1:

gagacccgcctaggtgaatatttagcagc
gattaaataccacgta**TATAAGGTGGACC**
GTTCTCGAGAGGTTCTTCCGGCAATGAC
GGCCAGAGCAAAAGCCACGTgtaggactg
catacgcctctacgcctccactgacgcga
tgatgtggcgtggatctgtttgctcttgg
tataggtcacggagacggctggtactgat
cccttcgggagtaaaaatataatgacat
ggcccaggcttcaggaggagttgtgcgg

Sekvencia 2:

tgtacagcactgcaacgagcatctggggg
ttggttattccgatggcgctggacagcta
gcggacagtagttctcaggccttagtaga
aaggtgggaaccccc**TATGAGGTGACCG**
TTTCAGCGTGACTATAGACGTCATTGAAG
CAATATACAGGAACACCACCTacttagga
agggagttcgggtgcagtaaagcattctta
cctcagggcacggtagagaacactacaac
cagaatagcaacgtgatgcggcgactctc

Lokálne zarovnanie:

TATAAGGTGGACCGTT-----CCTCGAGAGGTTCTTCCGGCAATGGCCACGAGAGCAAAAGCCACGT
TATGAGGTGACCGTTTTCAGCGTGACTATAGACGTCATTGAAGCAATATACAGG-----AACACCACCT

Obr. 1.1: Dve sekvencie a ich lokálne zarovnanie. Veľkými písmenami a červenou farbou sú vyznačené zarovnané časti. V zarovnaní sa nachádzajú zhody, nezhody a medzery v oboch sekvenciách

1.2 Párové zarovnávanie

Párové zarovnávanie je základná úloha zarovnávanie sekvencií, kde sa k sebe zarovnávajú dve sekvencie. V tejto práci sa budeme zaoberať len párovým zarovnávaním.

Kľúčové problémy sú:

1. Aké typy zarovnávanie by sme mali uvažovať
2. Skórovací systém, ktorý použijeme na ohodnotenie zarovnanie a tréovanie
3. Algoritmus, ktorý použijeme na hľadanie optimálneho alebo dobrého zarovnanie podľa skórovacieho systému
4. Štatistická významnosť zarovnanie.

[DEKM98]

1.3 Typy zarovnaní

Základné typy zarovnaní sú *Globálne zarovnanie* a *Lokálne zarovnanie*.

Definícia 1.3.1 (Globálne zarovnanie). Vstupom sú dve sekvencie $X = x_1x_2 \dots x_n$ a $Y = y_1y_2 \dots y_m$. Výstupom je zarovnanie celých sekvencií X a Y .

Definícia 1.3.2 (Lokálne zarovnanie). Vstupom sú dve sekvencie $X = x_1x_2 \dots x_n$ a $Y = y_1y_2 \dots y_m$. Výstupom je zarovnanie nejakých podreťazcov $x_i \dots x_j$ a $y_k \dots y_l$ sekvencií.

V praxi sa väčšinou snažíme nájsť zarovnanie s najvyšším skóre. [BV11]

1.4 Skórovacie systémy

Takmer všetky metódy zarovnanie hľadajú zarovnanie dvoch reťazcov na základe nejakej *skórovacej schémy*. Skórovacia schéma je metóda, ktorou priradíme zarovnanie skóre – zvyčajne čím väčšie skóre, tým realistejšie zarovnanie by to malo byť. Skórovacie schémy môžu byť veľmi jednoduché, napr. +1 za *zhodu* a −1 za *nezhodu* a

medzeru. Ak chceme mať schému, kde biologicky najkorektnejšie zarovnanie má najvyššie skóre, musíme vziať do úvahy, že biologické sekvencie majú evolučnú históriu, 3D štruktúru a mnohé ďalšie vlastnosti obmedzujúce ich evolučné procesy. Skórovací systém, ktorý by uvažoval všetky tieto faktory, by bol veľmi zložitý, preto sa používajú jednoduchšie, ktoré zanedbávajú veľa faktorov. [DEKM98]

1.4.1 Skórovacie matice

Skoro vždy však chceme rôzne zhody a nezhody skórovať rôzne - nie len všetky zhody $+1$ a nezhody -1 . Skóre môže závisieť od toho aké bázy sú v danom stĺpci zarovnania. Na to sa používa *skórovacia matica* (obr. 1.2), kde máme definované skóre pre každú dvojicu. Skórovacie matice sa využívajú aj pri zarovnávaní proteínov, kde niektoré dvojice majú podobné chemické vlastnosti. [DEKM98, BV11] Môžeme mať aj afinné skórovanie medzier, kde máme iné skóre za začatie medzery ako za predĺženie medzery. O tom si popíšeme v časti 1.5.3 Iným typom skórovacej schémy je napríklad *párový skrytý markvovský model* (viac v sekcii 1.6)

	A	C	G	T	-
A	6	-1	-2	-1	-3
C	-1	5	-3	-2	-4
G	-2	-3	5	-2	-2
T	-1	-2	-2	6	-1
-	-3	-4	-2	-1	∞

Obr. 1.2: Ukážka skórovacej matice. Všimnime si, že môžu byť rôzne skóre aj za jednotlivé zhody, nezhody alebo medzery.

1.5 Algoritmy na hľadanie zarovnaní

Pre danú skórovaciu schému potrebujeme algoritmus, ktorý nájde optimálne zarovnanie dvoch sekvencií. Budeme uvažovať zarovnávanie s medzerami. Medzery používame na znázornenie inzercie alebo delécie¹. Na označenie medzier používame pomlčku '-' a

¹vo všeobecnosti nevieme povedať, či ide o inzerciu v jednej sekvencii alebo deléciu v druhej sekvencii

medzeru do sekvencie medzi dve bázy pridáme tak, že medzi dve bázy napíšeme jednu alebo viac pomlčiek – počet pomlčiek zodpovedá počtu medzier a teda aj počtu báz, ktoré má na danom mieste jedna sekvencia navyše oproti druhej. Do sekvencie môžeme pridať ľubovoľne veľa medzier (pričom nemôžeme mať dve pomlčky v jednom stĺpci), aby sme dosiahli lepšie skóre. Pre dve sekvencie dĺžky n existuje

$$\binom{2n}{n} = \frac{(2n)!}{(n!)^2} \simeq \frac{2^{2n}}{\sqrt{\pi n}}$$

možných globálnych zarovnaní. [DEKM98] Pre dlhšie sekvencie nie je možné enumerovať všetky a vybrať tú s najlepším skóre.

Algoritmy na hľadanie zarovnaní využívajú *dynamické programovanie*. Často sa používajú aj rôzne heuristiky na urýchlenie výpočtu. My sa zaoberať algoritmami využívajúcimi heuristiky nebudeme. Pre rôzne typy zarovnaní a skórovacie schémy máme rôzne algoritmy zarovnávaní. [DEKM98, BV11]

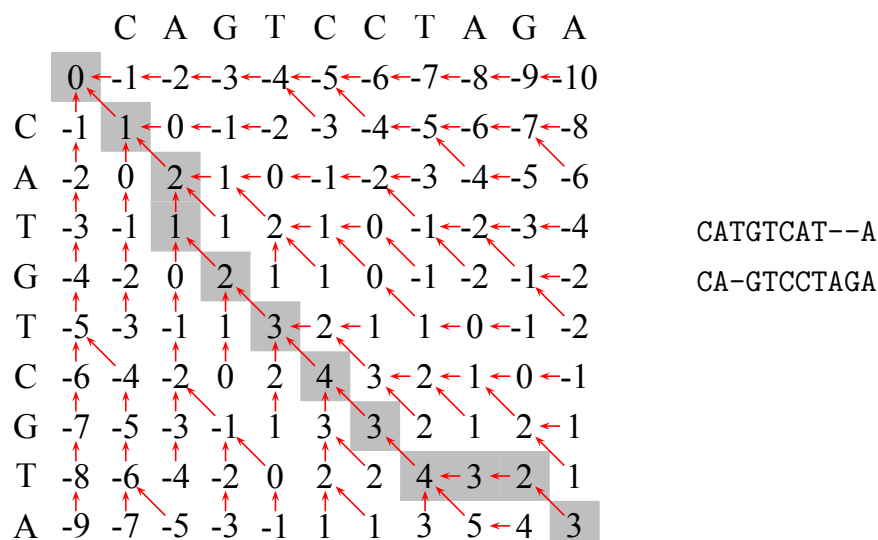
1.5.1 Algoritmus pre globálne zarovnanie: Needleman-Wunch

Máme dané dve sekvencie $X = x_1x_2 \dots x_n$ a $Y = y_1y_2 \dots y_m$, budeme zarovnávať všetky znaky sekvencie X a všetky znaky sekvencie Y . Definujeme si jednoduchú skórovaciu tabuľku kde $s(x, y)$ bude udávať skóre pre danú dvojicu báz (napr. +1 za zhodu, -1 za nezhodu) a za každú medzeru budeme dávať penaltu $-d$.

Algoritmus postupne vyplní 2-rozmernú maticu A . Riadky zodpovedajú bázam sekvencie X a stĺpce bázam Y . Na políčku $A[i, j]$ bude skóre najlepšieho zarovnania prvých i báz sekvencie X a prvých j báz Y .

Keď zarovnáваме sekvenciu s prázdnu sekvenciou, tak skóre bude $-n$, kde n je dĺžka sekvencie. Bude tam n pomlčiek, každá nám dá skóre -1 . Takto vyplníme riadky a stĺpce $A[i, 0]$ a $A[0, j]$.

Ak chceme vyplniť políčko $A[i, j]$, musíme si uvedomiť ako môže vyzeráť posledný stĺpec zarovnania $x_1x_2 \dots x_i$ a $y_1y_2 \dots y_j$. Máme iba 3 možnosti ako môže vyzeráť posledný stĺpec najlepšieho zarovnania. Buď obsahuje x_i alebo y_j alebo oboje. V prípade, že posledný stĺpec obsahuje oboje, cena tohto stĺpca je $s(x_i, y_j)$. Ak by sme posledný stĺpec zmazali, dostali by sme zarovnanie $x_1x_2 \dots x_{i-1}$ a $y_1y_2 \dots y_{j-1}$, pričom musí ísť o najlepšie zarovnanie. To už máme vypočítané v políčku $A[i-1, j-1]$, čiže výsledné skóre bude $A[i-1, j-1] + s(x_i, y_j)$.



Obr. 1.3: Tabuľka dynamického programovania pre globálne zarovnanie (vľavo) a výsledné zarovnanie (vpravo). Skóre je +1 za zhodu a -1 za neznhodu alebo medzeru.

V prípade, že posledný stĺpec obsahuje len x_i zarovnané s pomlčkou, skóre stĺpca bude -1 a po zmazaní dostávame zarovnanie $x_1x_2 \dots x_{i-1}$ a $y_1y_2 \dots y_j$, výsledné skóre bude teda $A[i-1, j] - 1$. V prípade, že posledný stĺpec obsahuje len y_i , tak skóre vypočítame analogicky.

Najlepšie skóre bude maximálne skóre pre všetky 3 prípady. Dostávame teda nasledujúci vzťah pre výpočet $A[i, j]$:

$$A[i, j] = \max \begin{cases} A[i-1, j-1] + s(x_i, y_j) \\ A[i-1, j] - d \\ A[i, j-1] - d \end{cases}$$

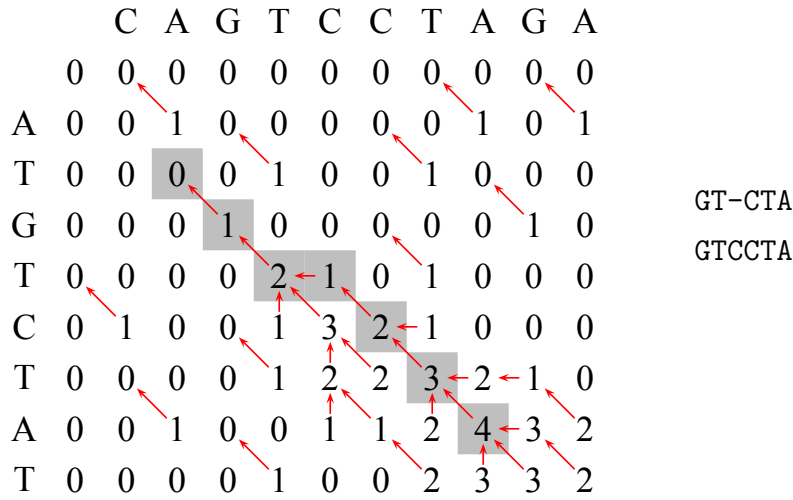
Maticu vieme vyplňať po riadkoch, pričom každé políčko vieme vypočítať z troch políčok, ktoré už sú vypočítané. Políčko $A[0, 0] = 0$ a krajné políčka potom vieme vypočítať ako $A[i, 0] = A[i-1, j] - d$ a $A[0, j] = A[i, j-1] - d$

Ak nás zaujíma aj zarovnanie – nie len jeho skóre – vieme si pre každé políčko zapamätať ktorá z troch možností dosiahla maximálnu hodnotu (červené šípky na Obr. 1.3). Na základe tejto informácie potom vieme zrekonštruovať zarovnanie tak, že postupne z posledného políčka ($A[n, m]$) budeme prechádzať na políčko, z ktorého sme vypočítali aktuálnu hodnotu.

Časová zložitosť je $O(nm)$, pretože vyplňame nm políček, každé v konštantnom čase. Zjavne aj pamäťová zložitosť je $O(nm)$.

Pamäťová zložitosť sa dá zredukovať na $O(n + m)$ za cenu zhruba dvojnásobného času výpočtu [Hir75].

1.5.2 Algoritmus pre lokálne zarovnanie: Smith-Waterman



Obr. 1.4: Tabuľka dynamického programovania pre lokálne zarovnanie (vľavo) a výsledné zarovnanie (vpravo). Skóre je +1 za zhodu a -1 za nezhodu alebo medzeru.

Algoritmus pre lokálne zarovnanie sa líši len v niekoľkých malých detailoch. Opäť vyplňame maticu A , s tým, že v $A[i, j]$ bude najvyššie skóre lokálneho zarovnanie medzi sekvenciami $x_1x_2 \dots x_i$ a $y_1y_2 \dots y_j$, ktoré buď obsahuje bázy x_i aj y_j , alebo je prázdne. Teda na ľubovoľnom mieste uvažujeme aj prázdne zarovnanie so skóre 0 (v matici nebudú záporné čísla). Vzťah pre výpočet $A[i, j]$ vyzerá takto:

$$A[i, j] = \max \begin{cases} 0 \\ A[i-1, j-1] + s(x_i, y_j) \\ A[i-1, j] - d \\ A[i, j-1] - d \end{cases}$$

V tomto prípade sú všetky krajné políčka nulové.

Zarovnanie potom nájdeme tak, že začneme z políčka s maximálnym skóre, a postupne budeme prechádzať na políčka, z ktorých sme získali maximálnu hodnotu, kým nenarazíme na nulu.

Časová aj pamäťová zložitosť sú, rovnako ako pri globálnom zarovnaní $O(nm)$.

1.5.3 Afínne skórovanie medzier

V jednoduchom skórovaní sme dávali za pomlčku vždy rovnaké skóre (-1) . Pri evolúcii sa však môže stať, že sa naraz zmaže (alebo vloží) niekoľko susedných báz, čo je pravdepodobnejšie ako to, že sa dané bázy zmazali (alebo vložili) nezávisle. Pri *afínnom skórovaní medzier* teda zavedieme dva typy skóre. Skóre za *začatie medzery* a skóre za *rozšírenie medzery*.

Algoritmus globálneho zarovnania vieme upraviť nasledovne: Namiesto matice A teraz budeme mať 3 matice M , I_x , I_y zodpovedajúce trom situáciám (Obr. 1.5).

ACT x_i	ACTTA x_i	ACT x_i --
AGTy j	AGTy j --	AGTATy j
(a) Mutácia(M)	(b) Inzercia v X (I_x)	(c) Inzercia v Y (I_y)

Obr. 1.5: Tri situácie pri afínnom skórovaní medzier

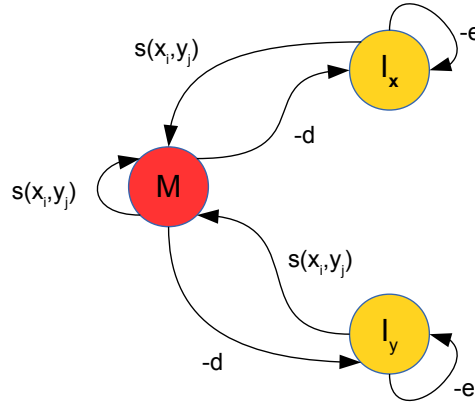
Nech $M[i, j]$ je najlepšie skóre prvých i báz zo sekvencie X a prvých j báz zo sekvencie Y , pričom x_i je zarovnané k y_j , $I_x[i, j]$ je najlepšie skóre ak x_i je zarovnané k medzere a $I_y[i, j]$ je najlepšie skóre ak y_j je zarovnané k medzere.

Označme si d penaltu za začatie medzery a e penaltu za rozšírenie medzery. Vzťahy pre výpočet políčok sú nasledovné:

$$\begin{aligned}
 M[i, j] &= \max \begin{cases} M[i-1, j-1] + s(x_i, y_j) \\ I_x[i-1, j-1] + s(x_i, y_j) \\ I_y[i-1, j-1] + s(x_i, y_j) \end{cases} \\
 I_x[i, j] &= \max \begin{cases} M[i-1, j] - d \\ I_x[i-1, j] - eI_y[i-1, j] - d \end{cases} \\
 I_y[i, j] &= \max \begin{cases} M[i, j-1] - d \\ I_x[i, j-1] - dI_y[i, j-1] - e \end{cases}
 \end{aligned}$$

V týchto rovniciach predpokladáme, že delécia nie je nasledovaná inzerciou. Toto platí v optimálnej sekvencii, ak $-d - s$ je menšie ako najmenšie skóre nezhody.

Tieto vzťahy vieme popísať stavovým diagramom na obrázku 1.6:



Obr. 1.6: Stavový diagram pre zarovnanie sekvencií - obsahuje *Match* (M), *InsertX* (I_x) a *InsertY* (I_y) stav a prechody medzi nimi spolu s ich cenou. Napríklad prechod z M do I_x znamená vloženie medzery do Y -ovej sekvencie a to penalizujeme $-d$

Časová zložitosť je $O(nm)$, pretože vyplníme $3nm$ políčok, každé v konštantnom čase. Pamäťová zložitosť je opäť $O(nm)$.

[DEKM98]

1.6 Zarovnávanie pomocou skrytých Markovovských modelov

Skrytý markovovský model (*Hidden Markov Model*, *HMM*) je generatívny pravdepodobnostný model, ktorý generuje náhodnú sekvenciu spolu s jej anotáciou (stavmi). HMM sa podobá na konečný automat. Skladá sa z konečného množstva stavov, prechodov medzi nimi a emisií. V každom kroku HMM vygeneruje a presunie sa do nejakého (aj toho istého) stavu. Generovanie symbolov aj presun medzi stavmi prebieha pravdepodobnostne. Konkrétny HMM je definovaný množinou stavov a nasledujúcimi distribúciami.

- distribúcia začiatočných stavov (pravdepodobnosť π_i , že HMM začne v stave i)

- distribúcia prechodov (pravdepodobnosť $a_{i,j}$, že HMM prejde zo stavu i do stavu j)
- distribúcia emisií (pravdepodobnosť $e_{i,x}$, že HMM v stave i vygeneruje symbol x)

Generovanie sekvencie teda vyzerá nasledovne: Na začiatku je HMM v stave i s pravdepodobnosťou π_i . Potom v každom kroku HMM emituje symbol x s pravdepodobnosťou $e_{i,x}$ a prejde do stavu j s pravdepodobnosťou $a_{i,j}$. Po n krokoch takto vygenerujeme sekvenciu dĺžky n , pričom každý symbol je anotovaný stavom, ktorý ho vygeneroval.

V takomto modeli vieme počítať pravdepodobnosť, že model vygeneruje sekvenciu x dĺžky n s anotáciou s ako súčin pravdepodobností prechodov a emisií. Výpočet vyzerá nasledovne:

$$P(X = x|S = s) = \pi_{s_1} \left(\prod_{i=1}^{n-1} e_{s_i, x_i} a_{s_i, s_{i+1}} \right) e_{s_n, x_n},$$

kde x je postupnosť symbolov a s je postupnosť stavov.

[BV11, DEKM98]

1.6.1 Zarovnávanie pomocou párových skrytých markovovských modelov (pHMM)

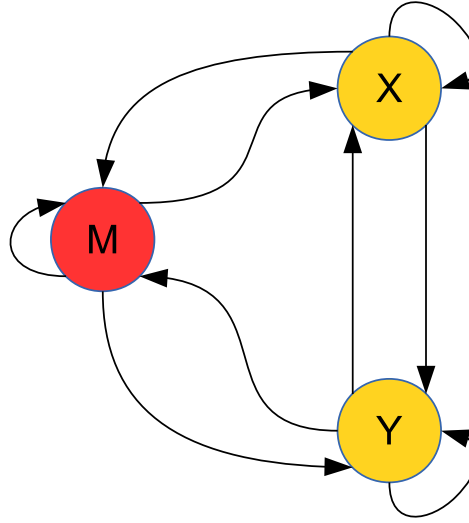
Párové skryté Markovovské modely (pHMM) sa od tých obyčajných líšia v tom, že namiesto jednej sekvencie generujú dvojicu sekvencií. Vedia generovať aj dvojicu symbolov aj jeden symbol v jednej sekvencii a nula symbolov v druhej. Na pHMM sa dajú použiť podobné algoritmy ako na HMM, ale treba ich upraviť tak, aby pracovali v dvoch rozmeroch.

V časti 1.5.3 sme si ukázali jednoduchý algoritmus na globálne zarovnávanie s afínnym skórovaním medzier. K tomuto algoritmu sme si uviedli aj jednoduchý stavový automat (Obr. 1.6). Tento automat vieme previesť na pHMM.

Na to aby sme automat previedli na pHMM, musíme urobiť niekoľko zmien – musíme nastaviť emisné a prechodové pravdepodobnosti, tak aby sčítavali do jedna.[DEKM98]

1.6.2 Viterbiho algoritmus na pHMM

Hľadáme najpravdepodobnejšiu postupnosť stavov A ak máme dané sekvencie X a Y , teda $\arg \max_A \Pr(A, X, Y)$. Úlohu budeme riešiť dynamickým programovaním.



Obr. 1.7: Párový HMM pre zarovnávanie sekvencií

Podproblém $V[i, j, u]$ je pravdepodobnosť najpravdepodobnejšej postupnosti stavov končiacej v $X[i]$ a $Y[j]$ v stave u .

Inicializácia (* označíme ľubovoľnú možnú hodnotu):

$$V[0, 0, M] = 1 \quad (1.1a)$$

$$V[i, 0, *] = 0 \quad \forall i \quad (1.1b)$$

$$V[0, j, *] = 0 \quad \forall j \quad (1.1c)$$

Rekurentné vzťahy:

$$V[i, j, M] = e_{M, (x_i, y_j)} \max \begin{cases} a_{M, M} V[i-1, j-1, M] \\ a_{I_x, M} V[i-1, j-1, I_x] \\ a_{I_y, M} V[i-1, j-1, I_y] \end{cases} \quad (1.2a)$$

$$V[i, j, I_x] = e_{I_x, x_i} \max \begin{cases} a_{M, I_x} V[i-1, j, M] \\ a_{I_x, I_x} V[i-1, j, I_x] \\ a_{I_y, I_x} V[i-1, j, I_y] \end{cases} \quad (1.2b)$$

$$V[i, j, I_y] = e_{I_y, y_j} \max \begin{cases} a_{M, I_y} V[i, j-1, M] \\ a_{I_y, I_y} V[i, j-1, I_y] \\ a_{I_x, I_y} V[i, j-1, I_x] \end{cases} \quad (1.2c)$$

Algoritmus funguje takto: Nech n a m sú dĺžky sekvencií X a Y . Na začiatku na-
inicializuje hodnoty pre $V[0, *, *]$ a $V[*, 0, *]$ podľa 1.1. Potom postupne pre všetky

$i = 1 \dots n$, $j = 1 \dots m$ a $u \in \{M, I_x, I_y\}$ dynamicky vypočíta $V[i, j, u]$ podľa 1.2.

Maximálne $V[n, m, u] \forall u \in \{M, I_x, I_y\}$ je pravdepodobnosť najpravdepodobnejšej postupnosti stavov. Aby sme stavy vedeli zrekonštruovať, pamätáme si pre každé $V[i, j, u]$ stav w , ktorý viedol k maximálnej hodnote. Táto postupnosť stavov nám udáva najpravdepodobnejšie zarovnanie.

Časová zložitosť tohto algoritmu je $O(nm)$.

Poznámka: pre dlhé sekvencie budú čísla $V[i, j, u]$ veľmi malé a môže dôjsť k podtečeniu. V praxi teda používame zlogaritmované hodnoty a namiesto násobenia súčet. [DEKM98]

1.6.3 Nastavenie parametrov pHMM

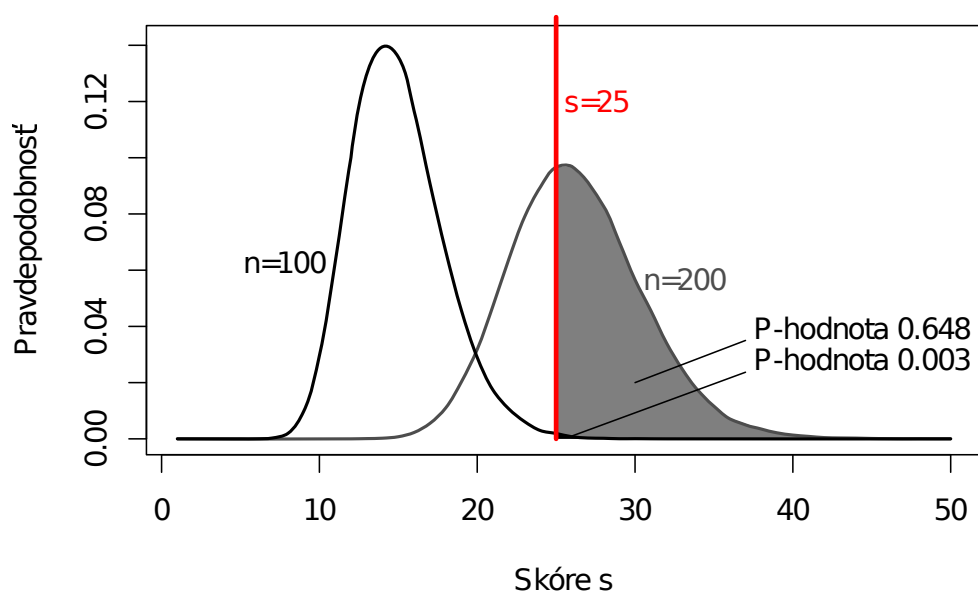
Ak máme anotované trénovacie sekvencie, môžeme HMM trénovať *metódou maximálnej vierohodnosti* (v angličtine *maximum likelihood estimation*). [DEKM98, Wik14c] Emisie získame tak, že vyfiltrujeme symboly s príslušným stavom a spočítame frekvencie pre každý stav zvlášť, pričom v Match stave počítame frekvencie dvojíc symbolov. Tranzície získame tak, že pre každý stav spočítame frekvencie nasledujúcich stavov. Parametre modelu môžeme teda ľahko natrénovať z existujúcich párových zarovnaní.

1.7 Štatistická významnosť zarovnaní

Smith-Watermanov algoritmus nájde najlepšie lokálne zarovnanie pre ľubovoľné dve sekvencie. Treba však rozhodnúť, či je zarovnanie dostatočne vierohodné nato, aby predstavovalo skutočnú podobnosť sekvencií a nie len najlepšie zarovnanie dvoch nesúvisiacich sekvencií. Ako vodítko pri rozhodnutí sa používajú identifikátory *štatistickej významnosti* zarovnaní: *P-hodnota* (*P-value*) alebo *E-hodnota* (*E-value*).

P-hodnota zarovnaní je pravdepodobnosť, že medzi náhodne generovanými sekvenciami tej istej dĺžky by sme našli zarovnanie s rovnakým skóre alebo vyšším. Keďže P-hodnota závisí od dĺžok sekvencií a skóre, musíme ju počítať pri každom zarovnaní. Je však časovo náročné robiť to generovaním veľkého množstva zarovnaní, preto sa používajú matematicky odvodené vzorce na odhad tejto hodnoty. ([KA90], [MB06]).

E-hodnota vyjadruje strednú hodnotu počtu zarovnaní so skóre aspoň takým ako má naše zarovnanie medzi náhodne generovanými sekvenciami. E-hodnota teda môže



Obr. 1.8: P-hodnota lokálneho zarovnania so skóre $s = 25$ medzi 2 sekvenciami dĺžky $n = 100$ alebo $n = 200$ (skórovanie $+1$ zhoda, -1 nezhoda alebo medzera). Rozdelenie bolo získané zarovnávaním 100000 párov náhodných sekvencií. Pri $n = 100$ je P-hodnota približne 0.003 a zodpovedá malej čiernej ploche pod krivkou napravo od zvislej čiary pre $s = 25$. Pri dlhších sekvenciách zodpovedá P-hodnota veľkej sivej ploche napravo od zvislej čiary. Pri takto dlhých sekvenciách očakávame skóre 25 alebo väčšie vo viac ako 60% prípadov čisto náhodou. Nejde teda o štatisticky významné zarovnanie.

byť aj väčšia ako jedna. Ak je E-hodnota väčšia ako jedna, tak čisto náhodou by sme očakávali aspoň jedno také silné zarovnanie, a teda zarovnania s takouto (a nižšou) E-hodnotou nebudeme považovať za štatisticky významné.

V štatistike sa v rôznych testoch štandardne používajú prahy na P-hodnotu 0.05 alebo 0.01. Pri zarovnávaní sekvencií však často používame ešte nižší prah, teda uvažujeme len zarovnania s P-hodnotou menšou ako napr. 10^{-5} . Pri malých hodnotách sú P-hodnota a E-hodnota približne rovnaké, teda taký istý prah môžeme použiť aj na E-hodnotu.

2 Súvisiaca práca

V tejto kapitole si uvedieme stručný prehľad modelov, ktoré zahŕňajú doplnkové informácie do zarovnania pomocou metód klasifikácie a stručne uvedieme v čom sa bude náš model líšiť.

V princípe môžeme rozlišovať dva typy modelov - *generatívny model* a *diskriminačný model*.

Konvenčné techniky odhadu pre zarovnávanie sa zakladajú na generatívnom modeli. Generatívny model (napr. HMM) sa snaží modelovať proces, ktorý generuje dáta ako pravdepodobnosť $P(X, Y, Z)$, kde $X = x_1x_2 \dots x_n$, $Y = y_1y_2 \dots y_m$ a Z je zarovnanie. Ak vieme vypočítať (alebo dobre odhadnúť) $P(X, Y, Z)$, hľadáme

$$\arg \max_z P(X = x, Y = y, Z = z)$$

čiže optimálne zarovnanie z z dvoch sekvencií x a y .

Aby sme zľahčili odhad pravdepodobnosti $P(X, Y, Z)$, rozložíme ju pomocou nezávislých pravdepodobností (napríklad v HMM na pravdepodobnosti prechodov a emisné pravdepodobnosti v stavoch). To síce vedie k efektívnym a jednoduchým výpočtom odhadu, ale obmedzuje to interakcie v rámci sekvencií, ktoré by sme mohli modelovať. Napríklad nemôžeme modelovať rôzne doplnkové informácie, ktoré sú závislé na sekvenciách [YJEP07].

Naproti tomu, v diskriminačných modeloch počítame

$$\arg \max_z P(Z = z | X = x, Y = y).$$

V tom prípade nemusíme modelovať proces, ktorý generuje x a y a môžeme sa zamerať na podstatnú časť problému [YJEP07].

Aktuálne existuje len niekoľko prístupov k diskriminačnému učeniu modelov zarovnávaní. Jeden z možných prístupov je riešiť *problém inverzného zarovnania* pomocou strojového učenia. [YJEP07]

Definícia 2.0.1 (Inverzné zarovnanie). Máme dané sekvencie a k nim zarovnanie. Inverzné zarovnanie nám vráti váhový model, s ktorým daný algoritmus na zarovnávanie vráti požadované zarovnanie k daným sekvenciám.

Problém inverzného zarovnania bol prvý krát formulovaný v [GS96]. Na tomto probléme je postavený aj model v [YJEP07], kde sa na trénovanie Support Vector Machine (SVM) dá pozerat ako na riešenie tohto problému. V článku sa zaoberajú použitím *Structural SVM* algoritmu na zarovnávanie proteínových sekvencií. Diskriminačné učenie umožňuje zahrnutie množstva dodatočnej informácie – státisíce parametrov. Navyše SVM umožňuje trénovanie pomocou rôznych účelových funkcií (loss functions). SVM algoritmus má lepšiu úspešnosť ako generatívna metóda SSALN, ktorá je veľmi presným generatívnym modelom zarovnaní zahrňajúcim informáciu o štruktúre.

Podobný prístup je aj v CONTRAlign [DGB06], kde sa používajú Conditional Random Fields (CRF). Tento prístup tiež ťaží z benefitov diskriminačného učenia, avšak na rozdiel od [YJEP07] neumožňuje použitie účelových funkcií.

3 Klasifikácia na základe lokálnej informácie

Cieľom našej práce bolo zahrnúť dodatočnú informáciu do zarovnávaní sekvencií, ktorá je poskytnutá formou anotácií k príslušným bázam. Keďže náš zarovnávač je založený na pHMM (sekcia 1.6.1), potrebujeme určiť emisné pravdepodobnosti v jednotlivých stavoch. Keďže chceme okrem jednotlivých báz zahrnúť aj iné informácie, rozhodli sme sa na tento účel využiť klasifikátory. Klasifikátory určujú, ktoré pozície sa majú zarovnať k sebe a túto informáciu zakomponovávame do výsledného zarovnávača pomocou modelov, ktoré si popíšeme v kapitole 4.

V tejto kapitole sa budeme venovať výberu klasifikátora a jeho vstupných atribútov. Najskôr si stručne popíšeme algoritmus klasifikátora, potom typy vstupných atribútov, ktoré sme navrhli a nakoniec sa budeme venovať vlastnostiam klasifikátora s danými typmi atribútov a výberu finálneho klasifikátora pre použitie v našich modeloch.

3.1 Náhodné lesy

Ako klasifikátor sme si vybrali *náhodný les* (*angl. Random forest*), pretože patrí v súčasnosti medzi najlepšie klasifikátory. V tejto sekcii si popíšeme základný algoritmus a v čom spočívajú jeho výhody a jeho najdôležitejšie vlastnosti.

Dôležité vlastnosti náhodných lesov

Klasifikátor náhodný les má niekoľko dôležitých vlastností, kvôli ktorým sme si ho vybrali:

- má vysokú presnosť (podľa [BC] je najpresnejší spomedzi všetkých vtedajších klasifikačných algoritmov)

- je efektívny aj na veľkých dátach
- dokáže obsiahnuť tisícky vstupných premenných
- dáva odhad, ktoré premenné sú dôležité na klasifikáciu.

[BC]

Navyše má ešte niekoľko ďalších užitočných vlastností, ktoré však pre nás neboli rozhodujúce:

- dokáže sa vysporiadať aj s väčším množstvom chýbajúcich dát
- produkuje interný nevychýlený odhad chyby aj počas tréovania lesa
- dokáže počítat vzťahy medzi atribútmi a klasifikáciou
- počítanie vzdialenosti môže byť použité aj na dáta bez informácie o triede, čo sa dá použiť na klastrovanie a detekciu outlierov.
- Ponúka experimentálnu metódu na detekciu interakcie medzi premennými.

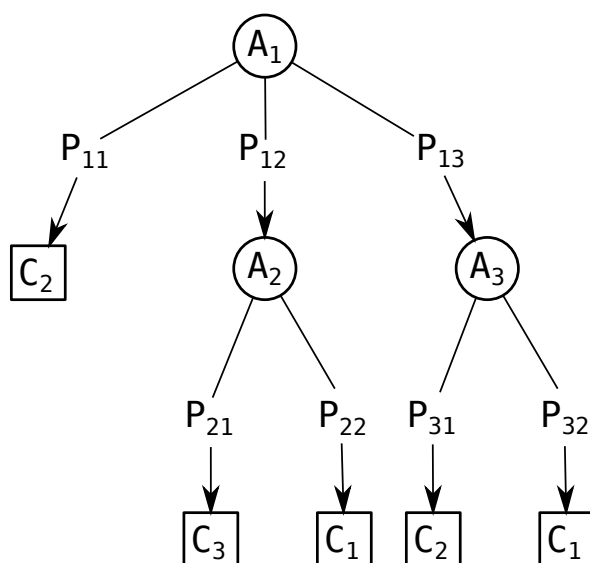
[BC]

3.1.1 Klasifikačné stromy

Keďže náhodný les je vlastne les *klasifikačných stromov* (niekedy označované ako *rozhodovacie stromy*), je nevyhnutné, aby sme si najskôr povedali niečo o nich.

Klasifikačný strom (obr 3.1) sa skladá z *uzlov* a *listov*. V uzle sa nachádza rozdeľovacie kritérium, podľa ktorého sa vieme rozhodnúť do ktorej vetvy máme pokračovať. V listoch sa nachádzajú triedy, do ktorých sa vstupný vektor klasifikuje. Klasifikačný strom je strom vybudovaný na základe trénovacej množiny, ktorý na základe vstupných údajov (vstupného vektora) predpovedá hodnotu výstupnej premennej.

Klasifikácia vstupného vektora vyzerá nasledovne: prechádzame strom zhora nadol a postupne sa zaraďujeme do vetiev podľa rozdeľovacích kritérií. Výstupná hodnota sa určí podľa toho, v ktorom liste skončíme.



Obr. 3.1: Klasifikačný strom. A sú atribúty, P sú konkrétne hodnoty atribútov a C sú triedy, do ktorých klasifikujeme dáta.

Trénovanie

Na tréovanie rozdovacích stromov sa používajú rôzne algoritmy. My si spomenieme najjednoduchší z nich – algoritmus ID3 [Wik14d]. Budeme rozlišovať dva typy listov: uzavretý a neuzavretý list. Uzavreté listy majú už priradenú triedu. Neuzavreté listy ešte nemajú priradenú triedu. Ku každému neuzavretému list prislúcha nejaká sada príkladov.

Algoritmus začína so stromom, ktorý má práve jeden vrchol – neuzavretý list. Kým existuje nejaký neuzavretý list l opakuje:

- Ak sú v jednom liste všetky príklady jednej triedy, označí list touto triedou a získame uzavretý list.
- Inak vyberie najinformatívnejší atribút A_l a vrcholu l pridá deti, rozdelené podľa hodnôt atribútu A_l .

Najinformatívnejší atribút vyberáme pomocou miery *zisk informácie* (angl. *information gain*). Zisk informácie je založený na rozdieli entrópie pred a po rozdelení podľa daného atribútu. Nech C je množina tried a p_c je pravdepodobnosť, že náhodný prvok

z S patrí do C , potom entropiu množiny S vypočítame ako:

$$E(S) = - \sum_{c \in C} p_c \log p_c$$

Zisk informácie potom vypočítame takto:

$$IG(S, A) = E(S) - \sum_{v \in A} \frac{|S_{\{A=v\}}|}{|S|} E(S_{\{A=v\}})$$

3.1.2 Náhodný les

Klasifikácia

Náhodný les sa skladá z mnohých klasifikačných stromov. Pri klasifikácii nejakého vstupného vektora sa predloží vstupný vektor každému stromu. Každý strom následne vráti klasifikáciu a hovoríme, že stromy *hlasujú*. Náhodný les následne zvolí klasifikáciu podľa väčšiny z hlasov všetkých stromov v lese.

Trénovanie

Pri trénovaní náhodných lesov sa využívajú dve hlavné techniky. Prvou z nich je takzvaný *bagging*, alebo tiež *bootstrapping*. Máme N trénovacích vektorov. Trénovaniu množinu pripravíme pre každý strom zvlášť tak, že z pôvodných N vektorov vyberieme náhodne s opakovaním N vektorov, ktoré budú použité na natrénovanie daného rozhodovacieho stromu.

Druhá technika spočíva v zmene trénovacieho algoritmu pre stromy. Nech trénovacie vektory majú M atribútov. Namiesto všetkých M atribútov sa v každom vrchole vyberie náhodne len $m \ll M$ z nich. Najlepší atribút na rozdelenie vrcholu sa potom vyberá z týchto m atribútov, pričom m je konštantné počas celého algoritmu.

Parameter m je jediným parametrom, na ktorý je náhodný les citlivý. Ako bolo ukázané v [Bre01], chyba náhodného lesa závisí na dvoch veciach: korelácii medzi dvoma stromami v lese (zvýšenie korelácie vedie k zväčšeniu chyby) a sily¹ každého stromu v lese (zvýšenie sily jednotlivých stromov vedie k zníženiu chyby). Ak znížime m znížime oboje – koreláciu aj silu stromov. Optimálny rozsah m je zvyčajne celkom veľký a v praxi sa zvykne často používať $m = \sqrt{M}$ alebo sa zvykne m zvoliť, tak aby sa

¹silný klasifikátor je taký, ktorý má malú chybu

minimalizovala *out of bag* (OOB) chyba. OOB chyba je počítaná na dátach, ktoré sa nenachádzajú v jednotlivých tréningových množinách pre dané stromy (viac v [BC]).

3.2 Použitie náhodných lesov na klasifikáciu zarovnaní

Keďže v našich modeloch máme dva typy stavov, použili sme dva typy klasifikátorov: *Match klasifikátor* pre Match stav a *Indel klasifikátor* pre Inzert stavy.

Klasifikátory dostanú na vstupe dáta asociované s pozíciami v daných sekvenciách a rozdeľujú ich do dvoch tried – nula a jedna. V Match klasifikátore trieda jedna znamená, že dané dve pozície majú byť zarovnané k sebe a nula znamená, že nie. V Indel klasifikátore trieda jedna označuje pozície, ktoré majú byť zarovnané k medzere a nula tie, ktoré nemajú. Ako výstupnú hodnotu pre naše modely berieme pravdepodobnosť, že dané dáta patria do triedy jedna. Túto pravdepodobnosť definujeme ako $P(C = 1|W = w)$, kde C je náhodná premenná pre výstupnú triedu, W je náhodná premenná pre množinu atribútov a w sú konkrétne hodnoty atribútov.

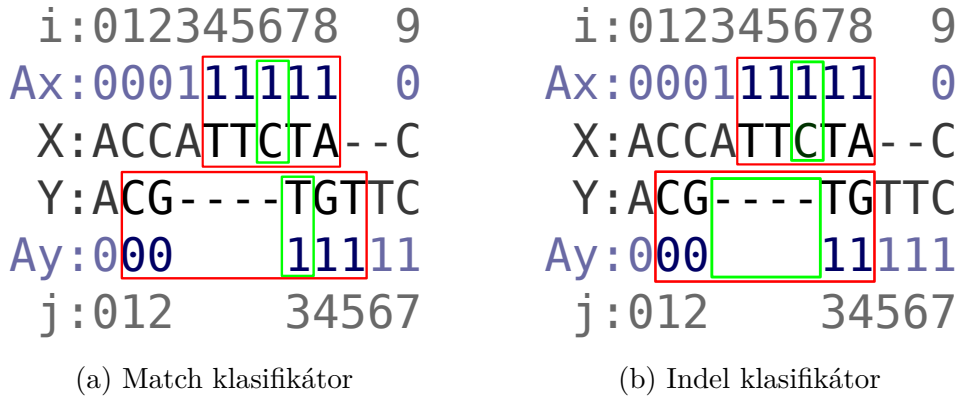
Keďže tieto klasifikátory sú nezávislé, nemusí platiť, že $P_M(C = 1|W_M = w_M) + P_I(C = 1|W_I = w_I) = 1$, kde P_M je výstup z Match klasifikátora a P_I je výstup z Indel klasifikátora a w_M a w_I sú atribúty pre tie isté pozície v Match a Indel klasifikátory. Môže sa napríklad stať, že oba klasifikátory zaradia atribúty prislúchajúce tým istým pozíciám do triedy jedna, alebo do triedy nula.

3.2.1 Výber atribútov

Klasifikátor dokáže pracovať len s poľom číselných atribútov, preto bolo treba najskôr dáta dať do tejto podoby. Navyše bolo treba rozhodnúť, aké dáta bude klasifikátor používať. Potrebovali sme, aby klasifikátor videl bázy a ich anotácie na daných pozíciách. Navyše sme sa rozhodli pridať ako pomôcku aj okolie daných báz spolu s anotáciami. Toto okolie budeme volať *okno*. Vyskúšali sme štyri rôzne formy okna a zisťovali sme, ktorá forma je pre klasifikátor najvhodnejšia. Na to sme mali dve miery: úspešnosť klasifikátora a dôležitosť atribútov. Pri úspešnosti sme sledovali rozdelenie výstupu klasifikátora pre pozitívne a negatívne atribúty. Dôležitosť atribútov sme mali hlavne na kontrolu a vizualizáciu atribútov.

Definícia okna

Okno veľkosti w pozostáva z $2w$ blokov veľkosti $k = (1 + \# \text{anotácií})$. Majme teda dve sekvencie, $X = x_1x_2 \dots x_n$ a $Y = y_1y_2 \dots y_n$ a pozície i a j . Pri Match klasifikátore okno veľkosti w obsahuje $x_{i-w/2} \dots x_i \dots x_{i+(1+w)/2}$, $y_{j-w/2} \dots y_j \dots y_{j+(1+w)/2}$ a všetky anotácie príslušných báz. (Obr. 3.2a) Pri Indel klasifikátore používame tiež dve pozície – prvá je pozícia bázy, na ktorú sa pýtame a druhá je pozícia medzery medzi dvoma bázami. Predpokladajme teraz, že X je sekvencia s bázou. Okno Indel klasifikátora veľkosti w obsahuje $x_{i-w/2} \dots x_i \dots x_{i+(1+w)/2}$, $y_{j-w/2} \dots y_j \dots y_{j+(1+w)/2-1}$ a všetky anotácie príslušných báz. (Obr. 3.2b)



Obr. 3.2: Okno klasifikátora veľkosti 5 pre pozície $i = 6$ a $j = 3$

Typ dát A - okno bez úpravy

Ako prvý typ dát sme zobrali okno tak, ako sme ho definovali v predošlej sekcii. Dáta obsahujú priamo všetky bázy a anotácie tak, ako sú v okne.

Typ dát B - zhody v stĺpcoch okna

Druhý typ dát obsahuje aktuálnu bázu spolu s jej anotáciami a navyše pole veľkosti $k * w$, ktoré má na i -tom mieste jedna ak $okno_X[i] = okno_Y[i]$, ináč nula (w je veľkosť okna, k je veľkosť bloku, $okno_X$ je časť okna zodpovedajúca sekvencii X a $okno_Y$ sekvencii Y). V Indel klasifikátore je jedna malá zmena: pozícia 0 aj 1 v x -ovej sekvencii sú porovnávané s pozíciou 1 v y -ovej a pozícia 2 v x -ovej sa porovnáva s pozíciou 2

v y -ovej. Pozíciu 1 v y -ovej sekvencii sme zopakovali preto, že sme experimentom zistili, že pre klasifikátor je dôležitá.

Typ dát C - matica zhôd v okne

Tretí typ dát je podobný ako typ B. Rozdiel je v tom, že teraz pole obsahuje nielen zhody po dvojiciach ale celú maticu zhôd. Teda opäť máme aktuálne bázy s anotáciami a pole má veľkosť $k * w^2$. Každý riadok sa skladá z jednotlivých blokov a v tabuľke v x -tom riadku, y -tom stĺpci a i -tom mieste v bloku je jedna práve vtedy, keď $okno_X[x+i] = okno_Y[y+i]$.

Typ dát D - kombinácia A a B

Posledný typ dát je kombináciou typov dát A a B. Dáta opäť obsahujú všetky bázy a anotácie tak ako v type A a navyše sme pridali pole zhôd z dát typu B. Táto informácia je síce redundantná a klasifikátor by si ju mal vedieť odvodiť aj sám, no experimenty ukázali, že to pomôže.

3.2.2 Trénovanie a testovanie klasifikátorov

Pre oba typy klasifikátorov sme sa snažili nájsť vhodné pozitívne a negatívne tréningové príklady. Z oboch sme do tréningovej množiny zahrnuli rovnako, aby bola tréningová množina vyvážená. Ak by sme nemali vyváženú tréningovú množinu, klasifikátory by zvýhodnili triedu, ktorej príkladov by bolo viac. Napríklad, keď sme pri Indel klasifikátore zobrali všetky pozície s medzerou ako pozitívne a všetky zarovnané pozície ako negatívne príklady, natrénovaný klasifikátor mal tendenciu dávať nižšie hodnoty, aby znížil tréningovú chybu, čomu sa chceme vyhnúť. Klasifikátor, ktorý by dával stále hodnotu nula, by mal tiež pomerne vysokú úspešnosť, pretože zarovnaných pozícií je rádovo viac ako pozícií s medzerou.

Výber pozitívnych príkladov bol v oboch prípadoch intuitívny. Pri výbere negatívnych príkladov sme sa museli zamyslieť nad vhodným protipólom k pozitívnym dátam.

Výber pozitívnych a negatívnych príkladov pre Match klasifikátor

Match klasifikátor sme chceli natrénovať tak, aby pre pozície, ktoré majú byť pri sebe, dával hodnoty blízke jednej a pre pozície ktoré k sebe byť zarovnané nemajú dával hodnoty blízke nule. Ako pozitívne príklady sme teda vybrali pozície z tréningových sekvencií, ktoré boli zarovnané k sebe.

Ako negatívne príklady sme vyskúšali dva prístupy: *náhodné dáta* a *dáta s posunom*.

Náhodné dáta: Náhodné dáta sme vybrali ako náhodné pozície, ktoré neboli zarovnané k sebe. Toto sa však ukázalo ako nedostatočné riešenie. Preto sme sa rozhodli pristúpiť k inému spôsobu výberu negatívnych vzoriek.

Dáta s posunom: Dáta s posunom sme vybrali posunutím jednej z dvoch pozícií v zarovnanom okne. Počas zarovňavania sa totiž väčšinou lokálne rozhodujeme, či zarovnať dané pozície k sebe, alebo vložiť medzeru. Ak by sme vložili medzeru, nastal by posun v jednej zo sekvencií. Rozhodli sme sa teda týmto inšpirovať a vyrábať negatívne príklady posunom zarovnaných pozícií. Pre každú zarovnanú pozíciu si najskôr rovnomerne náhodne vyberieme sekvenciu, ktorú budeme posúvať a potom z exponenciálneho rozdelenia náhodne vyberieme veľkosť posunu. Presný vzťah pre výpočet posunu Δ je

$$\Delta = (2D - 1) \cdot (1 + \lfloor S \rfloor) \quad D \sim \text{Alt}(0.5), S \sim \text{Exp}(0.75),$$

kde prvý člen nám určuje smer posunu (čo je to isté ako: ktorá sekvencia sa bude posúvať) a druhý člen určuje veľkosť posunu (chceme celočíselnú hodnotu ≥ 1).

Výber pozitívnych a negatívnych príkladov pre Indel klasifikátor

Indel klasifikátor sme chceli natrénovať tak, aby pre miesta, ktoré majú byť zarovnané s medzerou, dával čo najvyššie hodnoty a pre miesta, ktoré nemajú byť zarovnané k medzere, dával čo najnižšie hodnoty. Ako pozitívne príklady sme sa teda rozhodli vyberať pozície, ktoré boli v tréningových sekvenciách zarovnané k medzere. Ako negatívne príklady sme vybrali pozície, ktoré boli zarovnané k sebe, teda tie isté, čo sme pri tréningu Match klasifikátora považovali za pozitívne. Akurát v tomto prípade mala jedna

zo sekvencií okno skrátené o jedna, teda akoby bola medzera pred bázou, s ktorou je aktuálne uvažovaná pozícia zarovnaná.

Parametre trénovacej a testovacej množiny

Dáta pre naše experimenty pochádzajú z nášho simulátora, ktorý si popíšeme v 5.3.1. Trénovacia sada sa skladá z 20 zarovnaní sekvencií, pričom každé zarovnanie má dĺžku 10000. Celá trénovacia sada pre Match klasifikátor má 329428 príkladov. Pre Indel klasifikátor má sada 34150 príkladov. Testovacia sada sa skladá z jedného zarovnania s dĺžkou 10000. Pre Match klasifikátor obsahuje 16498 príkladov. Pre Indel klasifikátor obsahuje 1674 príkladov. Vo všetkých sadách je polovica príkladov pozitívnych a polovica negatívnych.

3.3 Výsledky experimentov

Použili sme klasifikátory tak, ako sme popísali v sekcii 3.2 a výsledky sa nachádzajú v tabuľke 3.1. Pre naše experimenty sme použili okno veľkosti päť, ale ako si ukážeme v časti 3.3.3 na veľkosti okna pri našich experimentoch nezáleží.

3.3.1 Dôležitosť atribútov

Atribúty sme označili takto: stredná pozícia je 0, pozície naľavo sú -1, -2 a pozície napravo 1 a 2, pričom v Indel klasifikátore v sekvencii s medzerou pozícia 0 chýba. Bázy v sekvencii X označujeme x_i a v Y y_i . Anotácie v sekvencii X označujeme a_i a v Y b_i . Zhody medzi bázami i, j budeme označovať $m_{i,j}$, pričom označenie zhôd v type dát B a D zjednodušíme na m_i . Podobne zhody medzi anotáciami označíme $l_{i,j}$ (resp. l_i).

V tabuľke 3.1 si môžeme všimnúť, že klasifikátory sa zamerali najmä na bázy a okrem prípadu C anotácie skoro nebrali do úvahy. Toto správanie zodpovedá tomu, že v praxi bázy majú podstatne väčší význam pri zarovnávaní sekvencií.

V prípade dát typu A a B sa Match klasifikátory nezamerali na stredné bázy ako by sme to čakali, ale všetky bázy brali približne s rovnakou váhou. Indel klasifikátory na tom boli podobne, hoci v prípade dát typu B je vidieť výraznejšiu preferenciu atribútu zhody báz x_0, y_1 (l_0), avšak atribúty l_1 a l_2 klasifikátoru neprišli dôležité.

Typ dát	chyba		atribúty	
	trénovacia	testovacia	najvýznamnejšie	najbezvýznamnejšie
A	93,07%	83,57%	$x_0, y_{-1}, y_1, x_*, y_*$	$\{a, b\}_*$
B	84,05%	84,31%	m_{-1}, m_0, m_*	l_*
C	80,44%	79,79%	$m_{-2,-2}, m_{2,2}, l_{-1,-1}, l_{1,1}, m_{0,0}$	ostatné
D	93,65%	84,32%	$m_0, m_{-1}, m_1, m_2, m_{-2}$	$l_*, \{a, b\}_*$

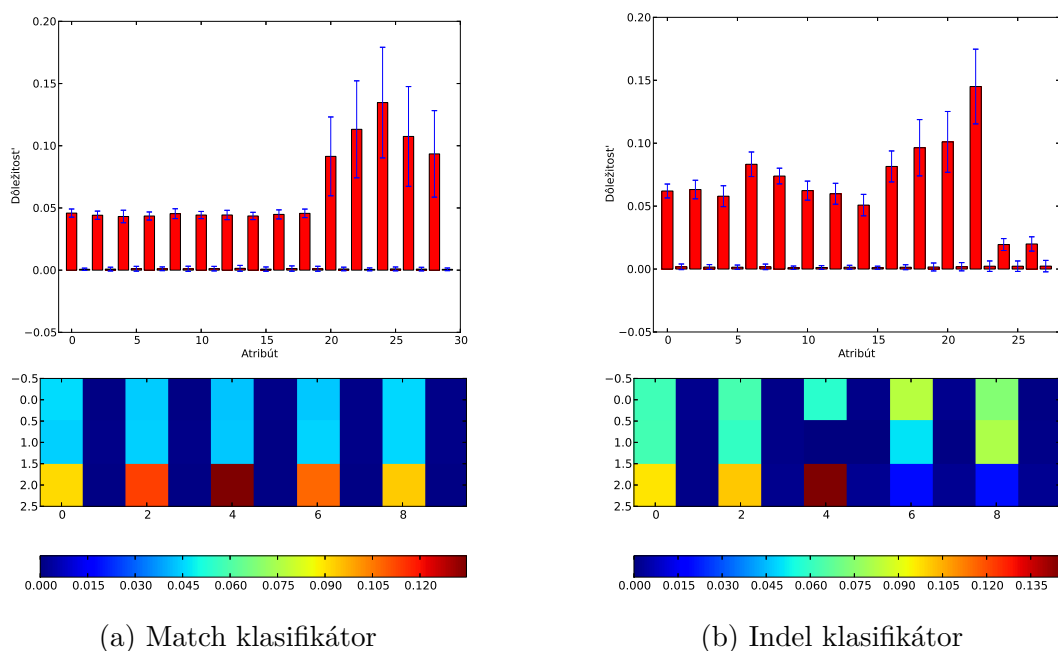
(a) Match klasifikátor

Typ dát	chyba		atribúty	
	trénovacia	testovacia	najvýznamnejšie	najbezvýznamnejšie
A	88,51%	74,13%	$x_0, x_1, x_{-1}, x_*, y_*$	$\{a, b\}_*$
B	77,17%	75,75%	m_0, m_{-1}, m_{-2}	l_*, m_1, m_2
C	78,03%	75,03%	$l_{2,2}, m_{-2,-2}, m_{0,1}, l_{-1,-1}, x_0, a_0$	ostatné
D	88,78%	76,46%	$m_0, m_{-1}, m_{-2}, x_1, \{x, y\}_*$	$\{l, a, b\}_*$

(b) Indel klasifikátor

Tabuľka 3.1: Porovnanie vlastností klasifikátorov pri rôznych typoch dát.

Pri type dát C si môžeme všimnúť, že najdôležitejšie atribúty sú na uhlopriečke, čo zodpovedá typu dát B, ibaže v tomto prípade sa nám vyskytli na niektorých pozíciách anotácie namiesto báz. Avšak ak sa nám tam vyskytla anotácia, bázu už klasifikátor nepovažoval za dôležitú a aj naopak. Indel klasifikátor navyše považoval za dôležitejšiu aj bázu na aktuálnej pozícii a jej anotáciu. Pri type dát C, na rozdiel od ostatných, klasifikátor viac berie do úvahy aj anotácie.

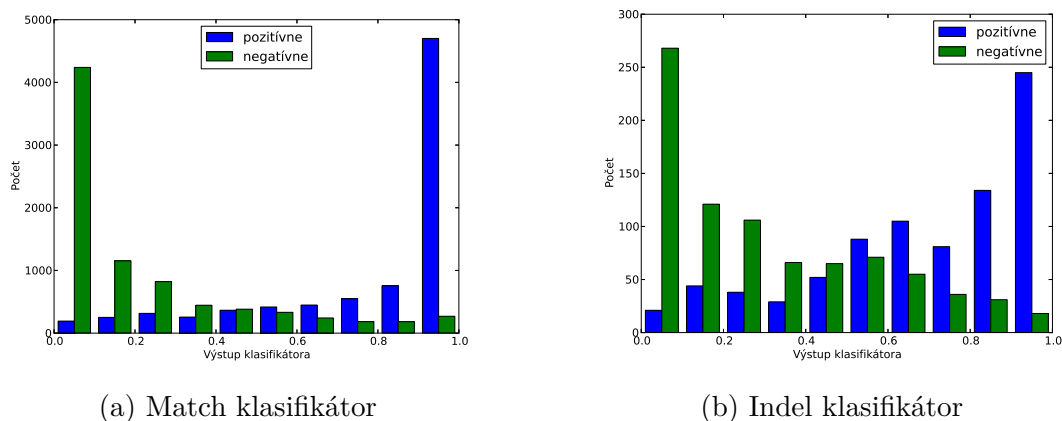


Obr. 3.3: **Dôležitosť atribútov pre typ dát D** - hodnoty sú normalizované, aby súčet bol jedna, modrý pásik označuje štandardnú odchýlku cez jednotlivé stromy v náhodnom lese. Pod grafom je tepelná mapa pre lepšiu vizualizáciu.

Pri dátach typu D sme dostali v prípade Match klasifikátora očakávanú distribúciu dôležitosti atribútov (obr. 3.3). Dáta typu B boli uprednostnené voči typu A, avšak aj bázy z typu A boli dôležité. Pri Indel klasifikátore sme dostali graf, ktorý je zložením grafov z dát typu A a B. Zaujímavosťou je, že pri tomto type dát má výraznejšiu úlohu práve zhoda na stredovej pozícii, teda pozície už nie sú také rovnocenné ako to bolo v type B.

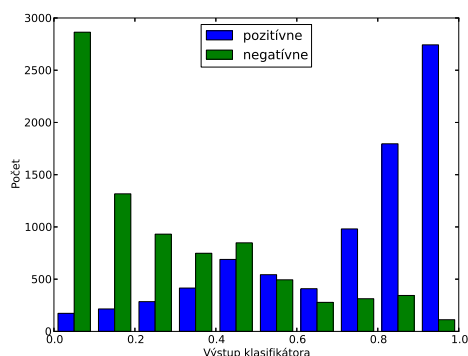
3.3.2 Úspešnosť klasifikátora

V tabuľke 3.1 vidíme, že hoci typ C obsahuje nadmnožinu atribútov B, jeho úspešnosť je nižšia a tento typ nemá zmysel používať. Ďalej si môžeme všimnúť, že kombinácia typov A a B využila klady oboch a mierne ich vylepšila.

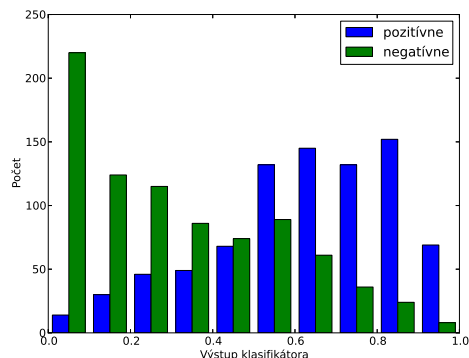


Obr. 3.4: Distribúcia výstupu z klasifikátora pri type dát D – modré sú pozitívne príklady a zelené sú negatívne. Na x -ovej osi je výstup klasifikátora a na y -je počet inštancií, pre ktoré výstup z klasifikátora padol do daného chlievika

Úspešnosť klasifikátorov s rôznym typom dát sme skúmali aj podrobnejšie, pozerali sme sa aj na distribúciu výstupov klasifikátora. Silný klasifikátor má totiž distribúciu takú, že väčšina pozitívnych príkladov má vysoké výstupné hodnoty a väčšina negatívnych klasifikátorov má nízke výstupné hodnoty. Najlepšiu distribúciu výstupov sa nám podarilo dosiahnuť práve pri type D (obr. 3.4). Na obrázku 3.5 sú pre porovnanie najhoršie distribúcie z ostatných typov dát.



(a) Typ C: Match klasifikátor



(b) Typ A: Indel klasifikátor

Obr. 3.5: Horšie distribúcie výstupov z klasifikátora – modré sú pozitívne príklady a zelené sú negatívne. Na x -ovej osi je výstup klasifikátora a na y -je počet inštancií, pre ktoré výstup z klasifikátora padol do daného chlievika

3.3.3 Nezávislosť typu dát na veľkosti okna

V tejto časti si ukážeme, že výber nášho typu dát nie je závislý na veľkosti okna. A teda môžeme najskôr vybrať typ dát a potom sa rozhodnúť, akú veľkosť okna použijeme.

Z tabuľky 3.2 vidíme, že klasifikátor s atribútmi typu D je, až na okno veľkosti tri, vždy najlepší. Rovnako zostalo zachované aj poradie ostatných typov dát. Z toho môžeme usúdiť, že pri voľbe typu dát nezáleží na veľkosti okna.

veľkosť okna	typ A	typ B	typ C	typ D
1	72,07	72,14	72,17	71,82
3	79,82	80,07	75,29	79,43
5	83,57	84,31	79,79	84,32
7	86,81	87,51	81,66	88,30
9	88,13	89,60	83,48	89,87
11	88,45	90,34	86,30	91,62
13	89,35	91,06	87,85	92,76
15	89,12	92,41	89,39	93,59

(a) Match klasifikátor

veľkosť okna	typ A	typ B	typ C	typ D
1	48,81	49,16	47,61	50,24
3	70,91	75,63	73,84	70,79
5	74,13	75,75	75,03	76,46
7	76,17	78,08	76,28	81,24
9	79,51	79,87	77,06	81,78
11	75,39	79,09	80,05	81,60
13	76,34	81,66	81,42	83,51
15	74,13	82,56	81,36	85,54

(b) Indel klasifikátor

Tabuľka 3.2: Porovnanie úspešnosti klasifikátorov pri rôznej veľkosti okna.

3.4 Zhrnutie

V tejto kapitole sme si predstavili klasifikátory, ktoré použijeme v našich modeloch. Uviedli sme dva typy klasifikátorov – pre Match stav a Inzert stav. Venovali sme sa výberu atribútov pre klasifikátory a ich trénovaniu. Predstavili sme si štyri množiny atribútov pre naše klasifikátory, pre ktoré sme porovnali vlastnosti klasifikátora.

Nakoniec sme sa rozhodli použiť ako množinu atribútov typ D, pretože s týmito atribútmi mali klasifikátory najlepšie vlastnosti a aj najvyššiu úspešnosť.

4 Modely

Hlavnou úlohou našej práce bolo preskúmať využitie dodatočnej informácie pri zarovnávaní sekvencií. Dodatočnú informáciu máme vo forme anotácií k sekvenciám. Anotácie sme sa rozhodli zakomponovať do zarovnania pomocou klasifikátorov, ktoré zosumarizujú informáciu obsiahnutú v anotáciách do jedného čísla. V predchádzajúcej kapitole sme riešili problém ako na základe lokálnej informácie rozhodnúť, či dané pozície majú byť zarovnané k sebe alebo k medzere. Navrhli sme dva klasifikátory – pre Match a Inzert stav, pričom prvý rozhoduje, či dané dve pozície majú byť zarovnané k sebe alebo nie. Druhý rozhoduje, či daná pozícia má byť zarovnaná k medzere.

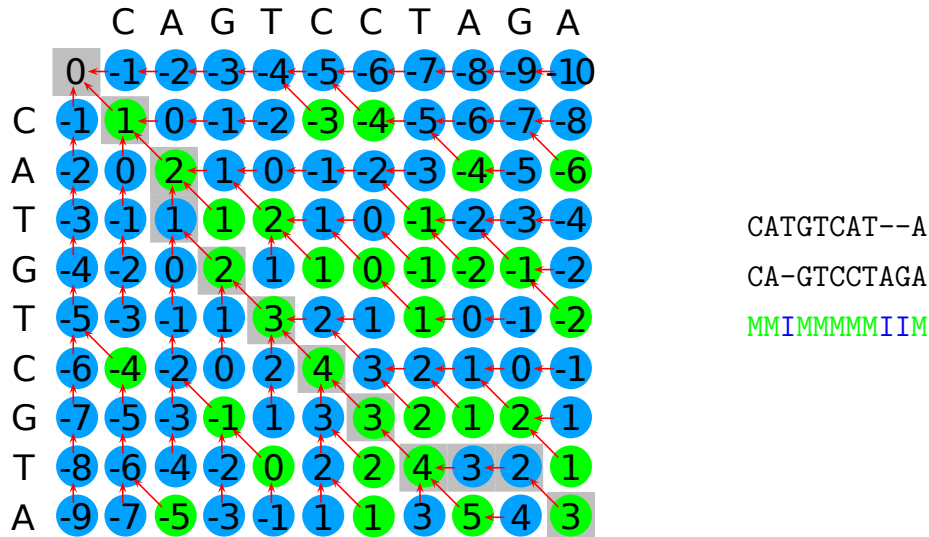
V tejto kapitole sa venujeme zakomponovaniu výsledkov z klasifikátorov do modelov pre zarovnávanie sekvencií. V sekcii 1.6 sme si zadefinovali pHMM pre zarovnávanie sekvencií (obr. 1.7). Najskôr si popíšeme architektúry našich dvoch modelov, ktoré sú založené na pHMM, pričom pHMM sme modifikovali tak, aby zahŕňal výstupy z klasifikátorov. Potom sa budeme venovať trénovaniu jednotlivých modelov a nakoniec si modely porovnáme a zhodnotíme ich úspešnosť.

4.1 Integrácia párových HMM a výsledkov z klasifikátorov

Pri zarovnávaní sekvencií si klasifikátor, z ktorého zoberieme výstup vyberieme podľa toho, v akom stave sa práve nachádzame. Pri Viterbiho algoritme konštruujeme tri dvojrozmerné tabuľky dynamického programovania, v ktorých sú na políčku (i, j) optimálne hodnoty pre zarovnanie končiace na i -tej pozícii v sekvencii X a na j -tej pozícii v sekvencii Y . Okrem toho si pamätáme aj odkiaľ sme do daného políčka prišli.

V jednom z našich modelov, ktoré definujeme v nasledujúcej sekcii, využijeme pásku s výstupom klasifikátora. Výstup z klasifikátora si môžeme pridať do tabuliek dynamic-

kého programovania. To, ktorý klasifikátor použijeme, závisí od smeru, ktorým sme sa v tabuľkách pohli. Ak sme sa pohli diagonálne, použijeme Match klasifikátor, ináč Indel klasifikátor. Ak teda pridáme k zarovnaniu pásku s výstupom klasifikátora, môžeme si ju predstaviť ako cestu v tabuľke, ktorá zodpovedá ceste zarovnania (obr. 4.1).

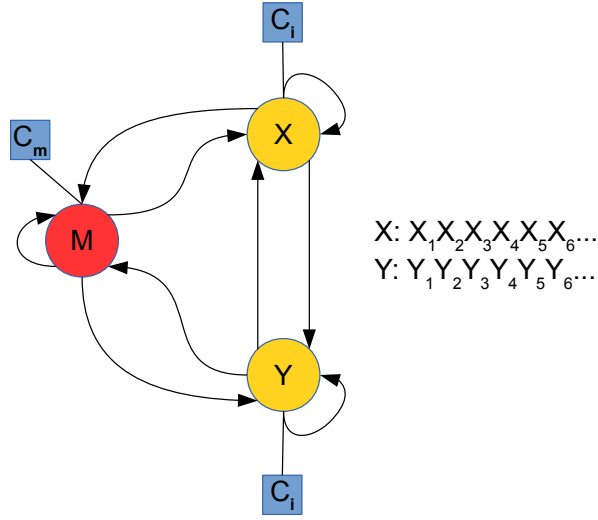


Obr. 4.1: Použité klasifikátory v klasifikátorovej páske pri zarovnávaní sekvencií. Pre jednoduchosť je použitá tabuľka z algoritmu Needleman-Wunch (kapitola 1.5.1) na globálne zarovnanie sekvencií. Viterbiho algoritmus používa tri tabuľky a cesta zarovnania vedie cez všetky tri. Zelenou farbou sú označené políčka, kde sa použije Match klasifikátor, modrou sú políčka, kde sa použije Indel klasifikátor. Šedou je vyznačené zarovnanie spolu s klasifikátorovou páskou pre zarovnanie. Vpravo je vypísané to isté zarovnanie spolu s klasifikátormi z klasifikátorovej pásky.

4.1.1 Definície modelov

Prvým modelom je *model s klasifikátorom ako emisiou (ďalej len model A)*. Tento model vyzerá rovnako ako základný model, aj pravdepodobnosti prechodov sú rovnaké, ale emisné pravdepodobnosti sú nahradené výstupom z klasifikátora.

V základnom modeli pre zarovnanie (kapitola 1.6) je pravdepodobnosť emisie v match stave $P(X = x \wedge Y = y)$, kde X, Y sú náhodné premenné pre bázy v sekvenciách



Obr. 4.2: Model s klasifikátorom ako emisiou

X a Y a x, y sú konkrétne bázy. Pre Inzert stavy analogicky. Na rozdiel od toho, výstup klasifikátora je definovaný ako $P(C = 1|W = w)$, kde $C \in \{0, 1\}$ je náhodná premenná pre výstupnú triedu, pričom nula znamená, že dáta nemajú byť zarovnané k sebe a 1 znamená, že majú byť, W je náhodná premenná pre okno klasifikátora, zahŕňajúce všetky atribúty a w je konkrétna hodnota okna. Pre porovnanie, ak by sme mali Match klasifikátor a ako atribúty len príslušné bázy, bolo by to $P(C = 1|X = x \wedge Y = y)$. Vidíme teda, že zatiaľ čo

$$\sum_{x \in X, y \in Y} P(X = x \wedge Y = y) = 1,$$

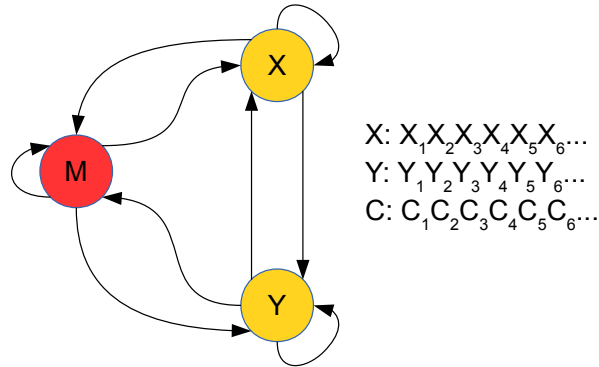
pre

$$\sum_{x \in X, y \in Y} P(C = 1|W = w)$$

takáto rovnosť neplatí. Viterbiho algoritmus však bude fungovať aj s hodnotami $P(C = 1|W = w)$ namiesto emisií a bude hľadať zarovnanie, ktoré bude maximalizovať takto definované skóre. O takomto modeli už však nemôžeme hovoriť ako o pravdepodobnostnom, ani ako o pHMM. Tento model je len inšpirovaný pHMM.

Keďže predošlý model nie je korektný pravdepodobnostný model, navrhli sme alternatívny model, *model s klasifikátorovou páskou (ďalej len model B)*, ktorý navyše modeluje aj výstup z klasifikátora, a teda je korektný pravdepodobnostný model. Nemodelujeme len dvojicu sekvencií, ale aj sekvenciu výstupov klasifikátora vo forme

pásky, ktorú sme si definovali na v sekcii 4. Pásku s výstupom z klasifikátora považujeme za akúsi pomôcku pre náš zarovnávač. Klasifikátor môže vrátiť ľubovoľnú



Obr. 4.3: Model s klasifikátorovou páskou

hodnotu z intervalu $\langle 0, 1 \rangle$. Klasické HMM robí iba s diskretnými hodnotami, a keďže výstup z klasifikátora je spojitý, museli sme sa s týmto problémom vysporiadať. Preto sme navrhli dve varianty tohto modelu – diskretnú a spojitú.

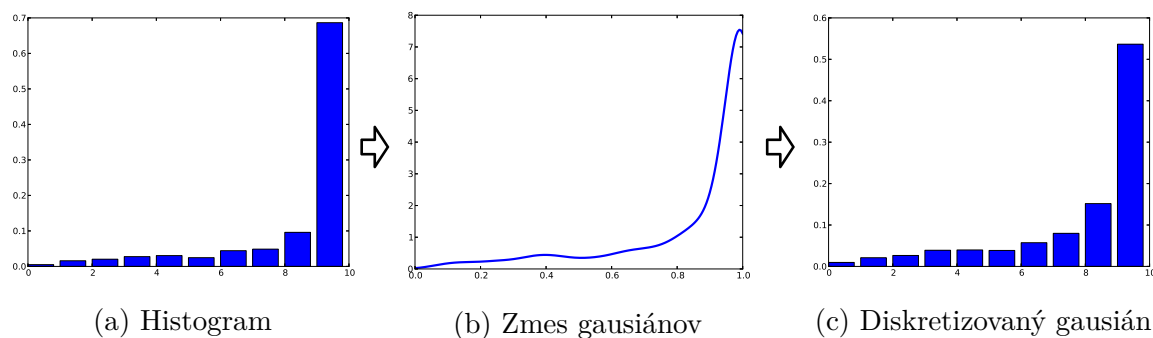
V diskretnom modeli sme hodnoty klasifikátora rovnomerne rozdelili na b hodnôt a výstup sme zaokrúhlili na najbližšiu z nich. V našich modeloch sme zvolili $b = 10$.

Alternatíva k predošlej metóde je použiť spojitý HMM. Hlavný rozdiel medzi spojitými a diskretnými HMM je, že spojitý HMM nepočítajú s pravdepodobnosťou, ale s hustotou. Hodnota hustoty v danom bode na rozdiel od pravdepodobnosti môže byť aj väčšia ako jedna. Podľa [HHL89], jedným zo štandardných spôsobov reprezentácie hustoty v spojitých HMM je zmes gausiánov, takže využijeme interpoláciu takouto distribúciou, a to budeme brať ako hustotu výstupu klasifikátora (obr. 4.4b).

4.1.2 Diskretizácia skóre klasifikátorov

Pre diskretnú verziu modelu B sme sa zaoberali spôsobom diskretizácie. Diskretizovať hodnoty môžeme dvoma spôsobmi. Buď jednoducho spočítame histogram pre trénovacie dáta, alebo interpolujeme vstupnú vzorku pomocou spojitej distribúcie, napr. pomocou zmesi gausiánov (tak ako v spojitej verzii) a potom toto rozdelenie diskretizujeme (obr. 4.4). Druhá spomenutá metóda má výhodu v tom, že vyhladí šum a doplní chýbajúce dáta, ale keďže je použitá aproximácia, zavádza určitú nepresnosť. V oboch prípadoch si

musíme zvoliť počet košov b , do ktorých dáta rozdelíme. Koše sme rozdelili rovnomerne a na základe experimentov (tabuľka 4.1) sme si zvolili $b = 10$.



Obr. 4.4: Dve možnosti diskretizácie výstupu klasifikátora: buď použijeme histogram 4.4a, alebo dáta najskôr aproximujeme pomocou zmesi gausiánov 4.4b a následne zdiskretizujeme ako na obr. 4.4c.

	počet košov	Úspešnosť
Histogram	10	81,78%
	20	81,35%
	50	59,97%
Diskretizovaný gausián	10	81,80%
	20	81,55%
	50	81,20%
	100	81,12%
Spojitý model	—	65,59%

Tabuľka 4.1: Porovnanie úspešností modelu B pri rôznom spracovaní pásky a rôznom počte košov pri diskretizácii.

Ako vidíme z tabuľky 4.1, spojitá verzia modelu sa nám neosvedčila. Má totiž nedostatok, že distribučná funkcia nedosahuje maximum pri výstupe klasifikátora rovnom jedna, ale kúsok predtým, čo znamená istú penalizáciu v prípade, že si je klasifikátor príliš istý. Môžeme si tiež všimnúť prudký pokles úspešnosti pri diskretizácii pomocou

histogramu s počtom košov 50. Je to spôsobené tým, že na trénovacej množine klasifikátor nevrátil žiadne hodnoty pre niektoré koše. Preto je rozumné použiť vyhladenie pomocou gausiánu, kde sa takéto problémy vyriešia. Okem toho bola úspešnosť s použitím gausiánu trochu vyššia ako bez neho takmer na všetkých testovaných sekvenciách.

4.2 Trénovanie modelov

Na testovanie úspešnosti klasifikátora sme použili tri typy množín, ktoré si označíme *simulované 1* (skr. *sim1*), *simulované 2* (*sim2*) a *biologické* (*bio*)

Každá dátová množina sa skladá z troch podmnožín – množina na trénovanie klasifikátora, množina na trénovanie modelu a testovacia množina. Množina na trénovanie modelu totiž musí byť nezávislá od množiny na trénovanie klasifikátora, pretože model budeme tiež používať na nezávislej množine a ak by sme použili tú istú množinu aj na trénovanie klasifikátora aj modelu, tak by model bol vychýlený.

Množiny pre simulované dáta pozostávajú z 20 zarovnaní s dĺžkou 10000 pre trénovanie klasifikátora, jedného zarovnania s dĺžkou 10000 pre trénovanie modelu a 6 zarovnaní s dĺžkou 1000 pre testovanie modelu. Množina biologických dát pre trénovanie klasifikátora sa skladá z 30 sekvencií s dĺžkami od 1000 do 1500, trénovacia množina pre modely obsahuje 41 sekvencií s dĺžkou od 1000 do 2500 a testovacia množina sa skladá zo 6 sekvencií s dĺžkami od 1000 do 1500.

V simulovaných dátach sme použili jednu anotáciu, ktorá označovala, či ide o gén, alebo nie. Simulované dáta 1 boli vygenerované našim simulátorom, ktorý je popísaný v kapitole 5.3.1. Simulované dáta 2 boli vygenerované tiež našim simulátorom, ktorý bol upravený tak, že v prípade génov sa sekvencie generovali z dvoch komplementárnych sekvencií a mimo génov z rovnakých sekvencií. Napríklad v prípade génov by sa namiesto AA vygenerovalo AT. To malo za úlohu posilniť význam anotácie. Biologické dáta boli prevzaté z databázy génov UCSC ([KBD⁺03]) a išlo o viacnásobné zarovnanie 7 sekvencií. Z nich sme si vybrali tri organizmy: človek (hg19), šimpanz (*Pan troglodytes*, panTro2) a pes (*Canis familiaris*, canFam2). Ako anotácie sme použili dáta z programu genscan¹ a anotáciu opakovaní, ktorá obsahuje informácie o tom, že daný kus sekvencie je repetitívny. Obe anotácie sme mali tiež z databázy génov UCSC.

¹program na hľadanie génov

Naše modely trénujeme pomocou metódy maximálnej vierohodnosti, tak ako sme si to popísali v kapitole 1.6.3. V HMM trénujeme prechodové a emisné pravdepodobnosti. Prechodové pravdepodobnosti sme trénovali vo všetkých modeloch rovnako a emisné pravdepodobnosti rôzne.

V modeli A sme trénovali iba prechodové pravdepodobnosti, emisie sme mali priamo z natrénovaného klasifikátora.

V modeli B sme trénovali aj prechodové aj emisné pravdepodobnosti. Označme si bázy v sekvenciách X, Y a výstup klasifikátora C . Veľkými písmenami X, Y, C budeme označovať náhodné premenné pre bázy a výstup klasifikátora, malými písmenami budeme označovať konkrétne hodnoty. V match stave generujeme trojice: (x, y, c) a v Inzert stavoch generujeme dvojice (x, c) alebo (y, c) . Popíšeme si len trénovanie emisií v Match stave, Inzert stavy trénujeme analogicky. V diskkrétnej verzii modelu si potrebujeme vyrobiť tabuľku $e[x, y, c] = P(X = x, Y = y, C = c) \quad \forall x, y, c$. V spojitnej verzii si potrebujeme zapamätať hustotu C pre každú možnú dvojicu (x, y) . Teda budeme mať tabuľku $e[x, y] = f_{x,y}(c) \quad \forall x, y$, kde $f_{x,y}(c)$ je hustota náhodnej premennej C za predpokladu, že $X = x$ a $Y = y$. Keďže distribúciu C chceme počítať pre každú možnú dvojicu báz zvlášť, musíme si ukázať, že pravdepodobnosť $P(X = x, Y = y, C = c)$ vieme rozložiť pomocou $P(C|X \cap Y)$. To vyplýva priamo z definície podmienenej pravdepodobnosti:

$$P(C = c \wedge X = x \wedge Y = y) = P(C = c|X = x \wedge Y = y) P(X = x \wedge Y = y)$$

4.3 Výsledky experimentov

V tejto sekcii porovnáme naše výsledky s výsledkami referenčného modelu a zároveň návača *muscle* [Edg04] na našich troch dátových množinách popísaných v sekcii 4.2.

Potom uvidíme výsledky niekoľkých experimentov, v ktorých sme porovnávali naše modely s rôznymi parametrami a v rôznych podmienkach.

Keďže k biologickým dátam nepoznáme správne zarovnanie, rozhodli sme sa zaviesť okrem percentuálnej zhody ďalšiu mieru. Túto mieru budeme označovať *tranzitivita* a budeme ju počítať z troch párových zarovnaní, ktoré dostaneme z troch sekvencií zarovnaním všetkých ich kombinácií. Označme si tri sekvencie ako A, B, C . Vyrobíme si tri zarovnania: AB, BC a AC . Tranzitivitu spočítame ako percentuálnu zhodu medzi

zložením prvých dvoch zarovnaní ($AB \circ BC$) a tretieho zarovnania AC . V ideálnom prípade by sme dostali 100% a mali by sme viacnásobné zarovnanie troch sekvencií. V praxi sú však hodnoty nižšie. Táto miera určuje konzistentnosť zarovnávača, čo je jedna z dôležitých vlastností kvalitného zarovnávača. V prípade biologických dát považujeme túto mieru za dôležitejšiu, pretože v tomto prípade náš zarovnávač môže nájsť biologicky správnejšie zarovnanie ako bolo to pôvodné. Avšak my sme si na naše testy vybrali viacnásobného zarovnanie 7 sekvencií, ktoré sú považované za veľmi kvalitné. Preto budeme požadovať aj vysoké percento zhody ako sekundárne kritérium.

Dáta	Model A		Model B		Ref. Model		Muscle	
	Zhoda	Tranz.	Zhoda	Tranz.	Zhoda	Tranz.	Zhoda	Tranz.
sim1	79,75%	44,97%	84,35%	56,5%	85,78%	61,03%	82,72%	58,76%
sim2	70,14%	—	71,47%	—	60,38%	—	61,47%	—
bio	91,40%	96,63%	91,24%	96,89%	91,34%	96,45%	91,28%	95,98%

Tabuľka 4.2: Porovnanie našich modelov s referenčným modelom a zarovnávačom muscle.

V tabuľke 4.2 sme porovnávali naše modely so základným modelom a zarovnávačom muscle. Typ dát v našich experimentoch je typ D a okno je veľkosti 9. Na simulovaných dátach, ktoré boli generované náhodne, naše modely nedokázali prekonať referenčné modely. Keď sme však posilnili význam anotácie (simulované dáta 2) referenčné modely, ktoré anotáciu k dispozícii nemali boli v nevýhode a naše modely dopadli lepšie. Na biologických dátach, dokázali dodatočné informácie pomôcť dosiahnuť lepšiu tranzitivitu (čo bol v tomto prípade náš primárny cieľ) a v prípade modelu A aj lepšiu úspešnosť.

Naše modely produkujú porovnateľné výsledky, avšak model B bol vo väčšine prípadov lepší. Model B je totiž kombináciou pôvodného pHMM a modelu A. Na rozdiel od modelu A sa nespolieha len na výsledky klasifikácie, pri ktorých sa občas môže stať, že nie sú dobré (úspešnosť je len okolo 80%), ale aj na samotné frekvencie výskytov dvojíc písmen. Klasifikátor potom na základe dodatočnej informácie upravuje túto pravdepodobnosť.

4.3.1 Vhodná veľkosť okna

V časti 3.3.3 sme si ukázali, že výber typu dát je nezávislý na výbere veľkosti okna. V tejto sekcii chceme zistiť aký vplyv na úspešnosť našich modelov má veľkosť okna. Vyskúšali sme viacero možností a výsledky sme zapísali do tabuľky 4.3.

Veľkosť okna	Model A	Model B
1	76,95%	57,86%
3	70,62%	81,42%
5	73,95%	82,98%
7	78,31%	83,05%
9	79,75%	84,35%
11	79,28%	83,57%
13	80,58%	83,98%
15	81,62%	83,98%

Tabuľka 4.3: Porovnanie úspešností modelov pri rôznej veľkosti okna.

Z tabuľky vidíme, že pre model B je optimálna veľkosť okna 7-9, pri väčších oknách už úspešnosť nerastie. Pri modeli A rástla úspešnosť až po veľkosť okna 15. Optimálnu hodnotu pre model A sme na základe tabuľky určili na 9-15. Pri vyšších hodnotách okna sa už úspešnosť nezvyšovala, tak sme ich do tabuľky nedali.

4.3.2 Vplyv dodatočnej informácie na zarovnanie

V tejto sekcii sme zisťovali, aký vplyv má dodatočná informácia na zarovnanie v simulovaných dátach 1. Porovnali sme štyri scenáre:

- okno veľkosti jedna bez anotácie – zodpovedá tomu, čo vidí aj základný pHMM
- okno veľkosti jedna s anotáciou – obsahuje navyše iba anotáciu
- okno veľkosti 9 bez anotácie – obsahuje navyše iba okolie v sekvencii
- okno veľkosti 9 s anotáciou – obsahuje navyše okolie v sekvencii aj anotáciu

Veľkosť okna	Anotácia	Model A	Model B
1	nie	74,42%	18,35%
1	áno	76,95%	57,86%
9	nie	78,43%	82,67%
9	áno	79,75%	84,35%

Tabuľka 4.4: Vplyv dodatočnej informácie na zarovnanie. V tabuľke je uvedená úspešnosť na simulovaných dátach 1.

V tabuľke 4.4 si môžeme všimnúť, že zatiaľ čo pri modeli A je úspešnosť veľmi podobná, hoci aj anotácie aj väčšie okno mierne pomohli. Naopak pri modeli B sú rozdiely výrazne väčšie, vidíme, že aj samotná anotácia výrazne pomohla a väčšie okno je pre tento model nutnosťou. Pri menšom počte atribútov totiž klasifikátor môže vrátiť len malý počet rôznych hodnôt (maximálne toľko, koľko je možných vstupov), čo je v modeli B nedostatočné. V modeli A to až tak nevádi, pretože v tomto prípade sa bude správať podobne ako pHMM – pre rovnaké bázy bude dávať vyššie hodnoty a pre rôzne nižšie.

4.3.3 Použitie jedného klasifikátora pre Match aj Inzert stav

V tejto sekcii si ukážeme, že voľba dvoch špecializovaných klasifikátorov mala svoje opodstatnenie. Namiesto dvoch špecifických klasifikátorov sme použili jeden klasifikátor aj pre Match aj pre Inzert stav. Tento klasifikátor bol taký istý ako náš Match klasifikátor a v Inzert stave sme prevrátili jeho hodnotu. Tieto možnosti sme porovnali na simulovaných dátach 1.

	Model A	Model B
1 klasifikátor	74,53%	82,72%
2 klasifikátory	79,75%	84,35%

Tabuľka 4.5: Porovnanie použitia jedného alebo dvoch typov klasifikátorov.

Z tabuľky 4.5 vidíme, že hoci naše modely dokážu fungovať aj s jedným klasifikáto-

rom, použitie dvoch typov klasifikátorov má svoje opodstatnenie a zvyšuje úspešnosť klasifikátorov.

4.4 Zhrnutie

V tejto kapitole sme si predstavili dva modely pre zarovnávanie sekvencií s dodatočnou informáciou pomocou klasifikátorov: model A a model B. Oba modely boli založené na pHMM. Prvý z nich priamo nahradil emisné pravdepodobnosti výstupom klasifikátora, a teda nemôže byť považovaný za pravdepodobnostný model. Druhý model sme navrhli ako pravdepodobnostný model, a okrem dvoch sekvencií mal k dispozícii pásku s výstupom z klasifikátorov. Uviedli sme si aj spôsob diskretizácie výstupu klasifikátora pre model B.

Ukázali sme, že na biologických dátach dokážu naše modely prekonať úspešnosť referenčných modelov. Okrem toho sme ukázali, že naše modely sú na rozdiel od referenčných modelov použiteľné aj v prípade, že vlastnosti zarovnania veľmi závisia od anotácie.

Nakoniec sme si ukázali aj opodstatnenie dvoch typov klasifikátorov namiesto jedného a aj to, že anotácia aj väčšie okolie daných báz majú pozitívny vplyv na zarovnávanie našimi modelmi.

5 Implementácia

V tejto kapitole sa budeme zaoberať implementáciou nášho riešenia. Popíšeme si dôvody pre výber jazyka a knižníc a niektoré časti nášho zarovnávača slúžiace na predspracovanie dát, klasifikáciu a zarovnávanie. Potom si popíšeme dôležité pomocné programy – simulátor a program na tréning modelov – a nakoniec si popíšeme použitie, konfiguráciu a možnosti rozšírenia programu.

5.1 Výber jazyka a použité knižnice

Ako jazyk pre písanie nášho zarovnávača sme si vybrali Python (vo verzii 2.7). Dôvodov sme mali niekoľko: V prvom rade zarovnávač sme nepísali od začiatku, ale ako základ sme zobrali zarovnávač od Michala Nánásiho (viď sekcia 5.1.1), ktorý už bol napísaný v jazyku Python 2.7. Ďalšími dôvodmi sú pohodlnosť písania v jazyku Python a dostupnosť veľkého množstva užitočných knižníc. Knižnice pre Python sa dajú písať aj v jazyku C alebo C++, čo umožňuje skĺbiť rýchlosť C a jednoduchosť jazyka Python.

5.1.1 Realigner

Realigner je program na zarovnávanie a prezarovnávanie¹, ktorý napísal Michal Nánási pre potreby [NVB13]. Umožňuje prezarovnanie sekvencií pomocou párových HMM a Viterbiho algoritmu. Prezarovnanie dokáže robiť aj v rámci lokálneho okolia okolo pôvodného zarovnania. Program je rozšíriteľný o nové modely, čo sme využili pre naše potreby.

¹modifikáciu existujúceho zarovnania

5.1.2 Pythonové knižnice

V našom programe sme sa okrem štandardnej knižnice jazyka Python rozhodli využiť niekoľko knižníc. Najdôležitejšími boli knižnice *numpy*, *scipy*, *scikit-learn*, *pyplot*, *track* a *pandas*.

Numpy je matematická knižnica, ktorá obsahuje množstvo užitočných matematických funkcií. Pre nás z nej boli dôležité hlavne štatistické funkcie, napr. generovanie náhodných čísel z rôznych rozdelení a počítanie štandardnej odchýlky. Okrem toho je táto knižnica základom aj pre väčšinu ostatných použitých knižníc.

Scipy obsahuje funkcie pre pokročilejšie vedecké výpočty. My sme z nej využili hlavne funkciu na odhad rozdelenia pomocou gausiánov.

Scikit-learn je knižnica zaoberajúca sa hlavne strojovým učením. Odtiaľ sme využili náhodné lesy.

Pyplot je knižnica na kreslenie grafov a využili sme ju na vygenerovanie všetkých našich vizualizácií.

Track je knižnica na prácu s '.bed' súbormi, v ktorých sme mali uložené naše anotácie.

Pandas je knižnica na analýzu a úpravu veľkých dát. Poslúžila nám najmä pri príprave anotácií pre biologické zarovnanie.

5.2 Triedy na zarovnávanie s klasifikátorom

V tejto sekcii si popíšeme najdôležitejšie triedy, ktoré sme doplnili do programu Realigner. Ide o triedy na prespracovanie dát, klasifikáciu a triedy pre stavy pHMM.

5.2.1 Predspracovanie dát

ToDo: obrazok classdiagramu

Na prespracovanie dát nám slúžili triedy *DataPreparer*, *ComparingDataPreparer*, *FullComparingDataPreparer*, *CombinedDataPreparer*, *IndelDataPreparer*, *ComparingIndelDataPreparer*, *FullComparingIndelDataPreparer* a *CombinedIndelDataPreparer* (obr ??).

Trieda *DataPreparer* je rodičovskou triedou pre ostatné triedy, takže si popíšeme

najmä ju. Trieda `IndelDataPreparer` je rodičom `Indel` verzií tried. Tieto triedy majú dve hlavné úlohy: Prvá je predpríprava dát na klasifikáciu a druhá je výroba trénovacej množiny zo sekvencií.

Pri predpríprave dát je úlohou dať dáta do formy vhodnej pre klasifikátor. To znamená získať okno a anotácie a ďalšie atribúty tak ako sme ich popísali v kapitole 3.2.1 a previesť ich do číselnej podoby. Naše štyri typy tried, ktoré zodpovedajú štyrom typom dát, sa líšia iba v tejto úlohe. Najdôležitejšia práca sa vykoná v metóde `_prepare_sequence`, ktorá pripraví dáta pre jednu sekvenciu, a potom sa zlepiť dáta z volania tejto funkcie pre obe sekvencie. `Indel` verzie tried obsahujú príslušné úpravy aby v sekvencii s medzerou použili kratšie okno.

Pri predspracovaní trénovacej množiny treba zo zarovnaní vytvoriť pozitívne a negatívne príklady. Táto metóda je implementovaná len v `DataPreparer` a `IndelDataPreparer` triede, ostatné ju zdedia. Pri výbere pozitívnych príkladov v triede `DataPreparer` označujeme pozície v sekvencii, ktoré sú zarovnané. To má na starosti metóda `prepare_positive_data`. Metóda `prepare_negative_data` vyberá negatívne príklady použitím zarovnaných pozícií z pozitívnych príkladov a náhodne ich poposúva podľa vzorca spomenutého v sekcii 3.2.2. Okrem toho sme implementovali aj náhodný výber negatívnych dát v metóde `prepare_negative_data_random`, kde negatívne príklady vyberáme náhodne z nezarovnaných pozícií s tým, že vyberieme rovnako veľa negatívnych ako pozitívnych príkladov. V `Indel` klasifikátore sa oba typy príkladov vyberajú v metóde `prepare_training_data`, pričom jedna sekvencia sa zvolí za medzerovú a za pozitívne príklady sa vezme všetko, čo je zarovnané k medzere v medzerovej sekvencii. Za negatívne sa vezme rovnako veľa náhodných zarovnaných pozícií.

5.2.2 Abstrakcia klasifikátora

V pre väčšiu modularitu sme sa rozhodli naprogramovať vrstvu medzi klasifikátorom a zvyškom programu. Trieda sprostredkujúca túto vrstvu sa volá *PairClassifier*. *PairClassifier* poskytuje niekoľko metód navyše oproti tým, ktoré poskytuje klasifikátor. Jej hlavnými výhodami sú:

- možnosť rýchlej výmeny klasifikátora bez úpravy zvyšku programu
- automatické predspracovanie pri klasifikácii a trénovaní pomocou príslušnej triedy

DataPreparer

- automatické tréovanie – ak klasifikátor nie je natrénovaný, tak sa automaticky natrénuje z prednastavenej tréovacej množiny
- hromadná klasifikácia s predspracovaním
- otočenie výstupu klasifikátora

5.2.3 HMM stavy s klasifikátorom

V zarovnávači bol implementovaný Viterbiho algoritmus, ktorý pracoval s párovými HMM. Tento algoritmus pracoval s triedou *GeneralizedPairState*, ktorá poskytovala metódu `emission` vracajúcu emisiu báz na daných pozíciách. Od tejto triedy dedia naše triedy *ClassifierState*, *ClassifierIndelState*, ktoré tvoria náš model A a preťažili sme metódu `emission` tak, aby vracala hodnoty z klasifikátora. Museli sme však urobiť niekoľko zmien. Viterbiho algoritmus si pýta emisie po jednej. Avšak volanie C-čkovej funkcie z knižnice `scikit-learn` z Pythonu je pomalé, takže pre dlhé sekvencie by zarovnanie trvalo veľmi dlho. Preto sme si najskôr naraz, počas jedného volania funkcie predpočítali emisie pre všetky pozície a tie sme potom zapamätali a vracali v momente, keď si ich Viterbiho algoritmus vypýtal. Na to sme využili možnosť hromadnej klasifikácie v našej triede *PairClassifier* (sekcii 5.2.2). Tieto dve triedy sa líšia iba v použitej triede na predspracovanie dát.

Od tried *ClassifierState* a *ClassifierIndelState* potom dedia triedy *ClassifierAnnotationState* a *ClassifierAnnotationIndelState*, ktoré tvoria diskretnú verziu nášho modelu B. Triedy mali v súbore uložené tabuľky emisií a podľa aktuálnych pozícií si zistili bázy a aktuálny výstup klasifikátora, ktoré slúžili ako indexy do tabuľky. Tieto triedy navyše obsahujú metódy na tréovanie emisií.

Od *ClassifierAnnotationState* a *ClassifierAnnotationIndelState* dedia triedy *ContinuousClassifierAnnotationState* a *ContinuousClassifierAnnotationIndelState*, ktoré stvárajú spojitú verziu modelu B. Tieto triedy majú v tabuľke emisií pre každú dvojicu báz uloženú serializovanú verziu objektu *gaussian_kde* (viac detailov je v sekcii 5.3.2), ktorý obsahuje distribúciu výstupu klasifikátora. Konkrétnu hodnotu potom získa volaním príslušného objektu k danej dvojici báz (resp. báze pri *IndelState*) s hodnotou klasifikátora. Tieto triedy takisto obsahujú aj metódy na tréovanie emisií.

Od triedy *GeneralizedPairState* navyše dedia aj naše triedy pre referenčný model – *SimpleMatchState* a *SimpleIndelState*, ktoré triedu *GeneralizedPairState* rozširujú o metódy na tréovanie emisných pravdepodobností.

ToDo: obrazok classdiagramu

5.3 Pomocné programy

V tejto sekcii sa budeme venovať dvom hlavným pomocným programom. Prvým z nich je náš simulátor, ktorý nám umožňuje generovať zarovnanie a druhým je program na tréovanie zarovnaní. Okrem toho sme implementovali množstvo programov na tvorbu vizualizácií pre túto prácu, alebo konverziu dát z rôznych zdrojov do formátov podporovaných našim zarovnávačom, o ktorých si však hovoriť nebudeme.

5.3.1 Simulátor

Simulátor slúži na generovanie zarovnaní pre naše experimenty. Náhodne vygeneruje tri² sekvencie, ktoré vznikli zo spoločného predka a vyrobí korektné spoločné zarovnanie všetkých troch sekvencií. Okrem toho vyrobí aj dodatočné informácie, ktoré majú pomôcť pri zarovnávaní.

Algoritmus

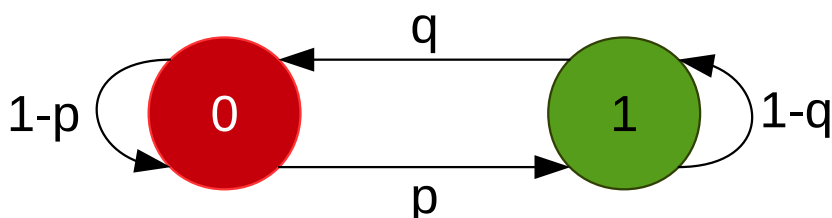
Ako prvú simulátor vygeneruje informáciu o tom, ktoré časti sekvencie prislúchajú génom a ktoré nie. Informácia má podobu binárneho vektora. Najskôr sa informácia vygeneruje pre *základnú (master) sekvenciu* a z nej sa potom skopíruje do troch dcérskych sekvencií, pričom s určitou pravdepodobnosťou sa niektoré gény zmažú (jednotky sa prepíšu na nulu). Táto informácia bude tvoriť pásku s anotáciou o géne.

Potom simulátor vygeneruje základnú sekvenciu a z nej následne pomocou mutácie a delécie (inzercia = delécia v zvyšných sekvenciách) odvodí tri ďalšie sekvencie. Mutácia má rôzne pravdepodobnosti v prípade, že ide o gén v jednej či druhej sekvencii, alebo oboch. Simulátor rovno generuje zarovnanie tým, že keď robí deléciu, nahradzuje bázy pomlčkami.

²aby sme vedeli počítať tranzitivitu

Generovanie informácie o génoch

ToDo: awk: Gény bývajú súvislé úseky, takže ich treba generovať tak, že niekedy začneme gén, potom generujeme jednotky, potom skončíme gén a generujeme nuly. Potom môžeme opäť začať gén atď.



Obr. 5.1: Markovova reťaz použitá na generovanie informácie o génoch. p je pravdepodobnosť, že začneme generovať gén, q je pravdepodobnosť, že prestaneme generovať gén.

Generovanie robíme pomocou *Markovovskej reťaze* (*Markov Chain*) obr. 5.1. Generujeme podľa aktuálneho stavu a v každom kroku sa náhodne rozhodneme, či sa prepne do iného stavu, alebo ostaneme v tom istom. Rozhodnutie robíme pomocou *falošnej mince* (*biased coin*), kde hlava padne s danou pravdepodobnosťou, ktorú nastavíme ako parameter. V našom prípade je to p pre stav nula a q pre stav jedna.

Simulácia mutácie a delécie

Ak už máme vygenerovanú základnú sekvenciu, potom vyrobíme zmutovanú sekvenciu tak, že s určitou pravdepodobnosťou nahradíme bázu zo základnej sekvencie inou bázou. Pravdepodobnosť závisí aj od toho, či je na danej pozícii gén v oboch sekvenciách, v jednej, alebo v žiadnej. Na rozhodovanie máme pre každú možnosť jednu falošnú mincu a podľa toho, ktorá z možností nastala, hodíme si príslušnou mincou.

Deléciu simulujeme opäť pomocou podobnej Markovovskej reťaze, pretože počas evolúcie majú tendenciu vypadávať súvislé úseky. Stav nula bude, že nemažeme a stav jedna, že mažeme, teda nahradzujeme danú bázu znakom '-'.

Využitie

Simulátor je prvá vec, ktorú sme implementovali a slúžil na väčšinu našich experimentov. Simulované dáta majú totiž oproti biologickým niekoľko výhod. Prvou je,

že poznáme správnu odpoveď, a teda môžeme objektívne zhodnotiť úspešnosť našich modelov. Druhou je možnosť zvoliť si parametre, ktoré menia závislosť mutácií na anotáciách a testovať správanie modelov, ak zvýšime túto závislosť. Ďalšou výhodou je, že si môžeme vygenerovať ľubovoľne veľa ľubovoľne dlhých sekvencií.

5.3.2 Trénovanie modelov

Trénovanie modelov má na starosti trieda *ModelTraining* a je urobené pomocou metódy maximálnej vierohodnosti (kapitola 1.6.3). Trénovanie prechodových pravdepodobností je pre všetky modely spoločné. Trénovanie emisných pravdepodobností je v každom stave spravené zvlášť, pretože stavy emitujú rôzne symboly a jednotlivé stavy majú rôzne spôsoby tréovania. Trieda *ModelTraining* najskôr načíta model z ktorého zistí stavy. Potom načíta trénovacie sekvencie a oannotuje ich stavmi. Následne vypočíta prechodové pravdepodobnosti. Spočíta všetky výskyty nasledujúcich stavov vo všetkých sekvenciách pre každý stav zvlášť. To vydelí celkovým počtom výskytov daného stavu. Potom oannotované sekvencie predloží jednotlivým stavom, aby vypočítali emisné pravdepodobnosti.

Máme štyri základné typy stavov, ktoré sme si popísali v časti 5.2.3 a pre každý máme triedu pre Match a Inzert verziu:

- SimpleMatchState, SimpleIndelState
- ClassifierState, ClassifierIndelState
- ClassifierAnnotationState, ClassifierAnnotationIndelState
- ContinuousClassifierAnnotationState, ContinuousClassifierAnnotationIndelState

V SimpleMatchState, SimpleIndelState algoritmus priamo spočíta výskyty príslušných dvojíc báz (resp. báz) a vydelí celkovým počtom výskytu stavu. V ClassifierState, ClassifierIndelState sa emisie netrénujú, takže tabuľka bude prázdna. V ClassifierAnnotationState, ClassifierAnnotationIndelState sa pre každý stav a pre každú dvojicu báz (resp. bázu) spočítajú hodnoty z klasifikátora a tie sa následne aproximujú pomocou gausiánov. Na toto využijeme triedu *gaussian_kde* z modulu *scipy.stats*. Táto trieda robí odhad distribučnej funkcie z množiny hodnôt. Šírku pásma (parameter,

ktorý určuje hladkosť výslednej funkcie) sme nechali na predvolenej hodnote – „scott”. Tá nastavuje šírku pásma podľa vzorca:

$$n^{-1/(d+4)},$$

kde n je počet vstupov a d je počet dimenzií (v našom prípade $d = 1$). Viac detailov je v [com13] alebo [Wik14b]. Následne túto funkciu diskretizujeme na k hodnôt. To spravíme pomocou určitého integrálu, pomocou metódy `integrate_box(a, b)`, ktorá spočíta určitý integrál od a po b . Keďže `gaussian_kde` je definovaný na intervale $\langle -\infty, \infty \rangle$ a

$$\int_{-\infty}^{\infty} \text{gaussian_kde}(x) dx = 1,$$

orezaním definičného oboru na $\langle 0, 1 \rangle$ sa môže stať, že

$$\int_0^1 \text{gaussian_kde}(x) dx < 1.$$

V tom prípade výslednú hodnotu ešte musíme normalizovať vydelením hodnotou

$$\mu = \int_0^1 \text{gaussian_kde}(x) dx.$$

Ako sme si ukázali v 4.2 výslednú pravdepodobnosť dopočítame vynásobením diskretizovanej pravdepodobnosti daného výstupu klasifikátora v prípade daných báz a pravdepodobnosti daných báz. V `ContinuousClassifierAnnotationState`, `ContinuousClassifierAnnotationIndelState` postupujeme rovnako, ako v predchádzajúcom prípade, ale distribúciu už nediskretizujeme. Namiesto toho si ju uložíme do súboru ako g spolu s číslom $\lambda = p/\mu$, kde p je pravdepodobnosť daných báz a μ je normalizačná konštanta – tá istá ako v predchádzajúcom prípade. Hustotu emisie potom počas zarovňavania vypočítame tak, že zoberieme parametre pre danú dvojicu báz (resp. bázu) – g a λ – a pre výstup klasifikátora c bude hustota $f(c) = \lambda g(c)$.

5.4 Použitie

Hlavný program sa spúšťa príkazom

```
PYTHONPATH=. python bin/Realigner.py vstup výstup
```

s nasledujúcimi parametrami

- `-mathtype`: float/LogNum – označuje, či sa má počítať s logaritmovanými hodnotami alebo s reálnymi
- `-algorithm`: viterbi – v našich modeloch je podporovaný len viterbi
- `-beam_width`: číslo – šírka okolo lúča vstupného zarovnania, ktorá sa má prehľadávať vo viterbiho algoritme
- `-model`: ClassificationHMM.js/OracleHMM.js/SimpleHMM2.js/vlastný – model, ktorý sa má použiť (naše modely sa nachádzajú v data/models)
- `-sequence_regexp`: sekvenciaX sekvenciaY regulárne výrazy na vyfiltrovanie dvoch sekvencií zo vstupného súboru zarovnania
- `-annotation_model`: súbor.js – súbor s konfiguráciou anotácií pre jednotlivé sekvencie

Pozičné parametre

- 1: vstup – súbor so vstupným zarovnaním
- 2: výstup – súbor kam sa má uložiť výstupné zarovnanie

Napríklad:

```
PYTHONPATH=. python bin/Realigner.py\  
    --mathType LogNum --algorithm viterbi --beam_width 30\  
    --model data/models/ClassificationHMM.js\  
    --sequence_regexp ^sequence1$ ^sequence2$\  
    --annotation_model annotations.js\  
    input_alignment.fa output_alignment.fa
```

Predtým, ako prvýkrát spustíme zarovnávač, je nutné model najskôr natrénovať. To spravíme spustením programu

```
python model_training.py model trénovacie_dáta,
```

kde prvý parameter je model, ktorý chceme trénovať a druhý je adresár s trénovacími sekvenciami. Ak ešte nemáme natrénovaný klasifikátor, tak sa natrénuje automaticky z dát ktoré sa nachádzajú v priečinku data/sequences/train_sequences. Napríklad

```
python model_training.py\
    data/models/OracleHMM.js data/sequences/model_training/bio,
```

Klasifikátor sa potom uloží do `data/clf` odkiaľ sa bude automaticky načítavať aj pri trénovaní modelov aj pri zarovnávaní. Pokiaľ treba klasifikátor natrénovať odznovu, je nutné ho z `data/clf` vymazať.

5.4.1 Konfigurácia a Formáty súborov

Konfigurácia našich modelov sa nachádza v dvoch súboroch:

- `classifier_alignment/constants.py` – tu sa dá zmeniť veľkosť okna a či sú povolené anotácie
- `classifier_alignment/config.py` – tu sa dá zmeniť klasifikátor, datapreparer a či sa má použiť ten istý klasifikátor aj na Match aj na Inzert stav alebo nie. V tomto súbore sú preddefinované klasifikátory a datapreparer-y takže stačí zmeniť index označujúci, ktorý z nich sa má použiť.

Súbory zarovnania sú vo formáte FASTA³, súbory anotácií sú vo formáte BED⁴.

Súbory modelov sú vo formáte JSON⁵, nasledujúceho tvaru:

```
1 {
2   "model": {
3     "states": [
4       {
5         "name": "Match",
6         "durations": [[[1, 1], 1.0]],
7         "startprob": 0.33,
8         "endprob": 1.0,
9         "emission": [],
10        "onechar": "M",
11        "__name__": "SimpleMatchState"
12      },
13      {
```

³http://en.wikipedia.org/wiki/FASTA_format

⁴<http://genome.ucsc.edu/FAQ/FAQformat.html#format1>

⁵JavaScript object notation


```

14     "name": "InsertX",
15     "durations": [[[1, 0], 1.0]],
16     "startprob": 0.33,
17     "endprob": 1.0,
18     "emission": [],
19     "onechar": "X",
20     "__name__": "SimpleIndelState"
21 },
22 {
23     "name": "InsertY",
24     "durations": [[[0, 1], 1.0]],
25     "startprob": 0.33,
26     "endprob": 1.0,
27     "emission": [],
28     "onechar": "Y",
29     "__name__": "SimpleIndelState"
30 }
31 ],
32 "__name__": "GeneralizedPairHMM",
33 "transitions": []
34 }
35 }

```

Emisné a prechodové tabuľky sa vyplnia pri trénovaní modelu. Zaujímavé sú pre nás len položky `__name__` v jednotlivých stavoch. Tam treba vložiť meno príslušnej triedy pre stav.

Súbor pre konfiguráciu anotácií je JSON nasledujúceho tvaru:

```

1 {
2     "__name__": "Annotations",
3     "annotations": ["gene"],
4     "sequences": [
5         {
6             "name": "sequence1", "annotations": [
7                 {
8                     "id": "gene",
9                     "file": "data/sequences/simulated/
                        simulated_alignment_sequence1_gene.bed"

```

```
10         "offset": 47
11     }
12 ]
13 },
14 {
15     "name": "sequence2", "annotations": [
16     {
17         "id": "gene",
18         "file": "data/sequences/simulated/
19             simulated_alignment_sequence2_gene.bed"
20         "offset": 47
21     }
22 ],
23 },
24 {
25     "name": "sequence3", "annotations": [
26     {
27         "id": "gene",
28         "file": "data/sequences/simulated/
29             simulated_alignment_sequence3_gene.bed"
30         "offset": 47
31     }
32 ],
33 }
```

Do **annotations** treba dať zoznam anotácií ktoré sa majú použiť a do **sequences** treba dať zoznam sekvencií – meno musí byť rovnaké ako vo FASTA súbore. Pre každú sekvenciu treba vyplniť v **annotations** pre každú anotáciu (id musí byť rovnaké ako v zozname anotácií) zdrojový súbor, odkiaľ sa má daná anotácia načítať. Navyše je možné špecifikovať **offset** v prípade, že sekvencia, ktorú chceme zarovnať, nezačína v genóme na pozícii 0 a anotácie sú číslované podľa pozícií v genóme.

5.5 Rozšírenie programu

Program je možné ľahko rozšíriť rôznymi spôsobmi. Je možné pridať:

- novú triedu pre stav
- nový DataPreparer
- nový klasifikátor

Nová trieda pre stav musí dediť od GeneralizedPairState, alebo niektorého z jej potomkov. Najdôležitejšiou metódou je `emission`, ktorá dostane dve pozície, dve sekvencie a dve hodnoty, ktoré označujú koľko symbolov z ktorej sekvencie treba vygenerovať, jej výstupom je pravdepodobnosť, že sa vygenerujú dané znaky z daných sekvencií od daných pozícií. Ak chceme túto triedu použiť v modeli, musíme ju ešte zaregistrovať do triedy *HMMLoader* pomocou metódy `addFunction`.

Nový DataPreparer musí dediť od triedy DataPreparer, alebo niektorej z jej potomkov. Najdôležitejšou je metóda `prepare_data`, ktorá dostane obe sekvencie, anotácie a aj dve pozície v sekvenciách. Výstupom musí byť pole číselných atribútov. Preťažením metódy `prepare_training_data` sa dá zmeniť generovanie trénovacej množiny.

Nový klasifikátor musí implementovať nasledujúce metódy:

- `fit(X, y)` – dostane pole trénovacích vektorov a pole cieľových hodnôt, z ktorých natrénuje svoje parametre
- `predict_proba(X)` – dostane pole vstupných vektorov a výstupom je dvojrozmerné pole – pre každý vstupný vektor a každú triedu pravdepodobnosť, že vektor patrí do danej triedy

Navyše musí byť serializovateľný pomocou modulu *pickle*.

Záver

Literatúra

- [BC] Leo Breiman and Adele Cutler. Random forests. [Online; accessed 14-Jan-2013].
- [BPSS11] Alvis Brazma, Helen Parkinson, Thomas Schlitt, and Mohammadreza Shojatalab. EBI Research - Functional Genomics - Introduction To Biology. 2011. [Online; accessed 26-Oct-2012].
- [Bre01] L. Breiman. Random forests. *Machine learning*, 45(1):5–32, 2001.
- [BV11] Broňa Brejová and Tomáš Vinař. *Metódy v bioinformatike [Methods in Bioinformatics]*. Knižničné a edičné centrum FMFI UK, 2011. Lecture notes.
- [com13] SciPy community. *SciPy Reference Guid*, 2013.
<http://docs.scipy.org/doc/scipy/scipy-ref.pdf>.
- [DEKM98] R. Durbin, S.R. Eddy, A. Krogh, and G. Mitchison. *Biological Sequence Analysis: Probabilistic Models of Proteins and Nucleic Acids*. Cambridge University Press, 1998.
- [DGB06] Chuong Do, Samuel Gross, and Serafim Batzoglou. Contralign: Discriminative training for protein sequence alignment. In Alberto Apostolico, Concettina Guerra, Sorin Istrail, Pavel Pevzner, and Michael Waterman, editors, *Research in Computational Molecular Biology*, volume 3909 of *Lecture Notes in Computer Science*, pages 160–174. Springer Berlin / Heidelberg, 2006. 10.1007/11732990_15.
- [Edg04] Robert C Edgar. Muscle: multiple sequence alignment with high accuracy and high throughput. *Nucleic acids research*, 32(5):1792–1797, 2004.

- [GS96] D. Gusfield and P. Stelling. [28] parametric and inverse-parametric sequence alignment with xparal. *Methods in enzymology*, 266:481–494, 1996.
- [HHL89] XD Huang, Hsiao-Wuen Hon, and Kai-Fu Lee. Multiple codebook semi-continuous hidden markov models for speaker-independent continuous speech recognition. 1989.
- [Hir75] D. S. Hirschberg. A linear space algorithm for computing maximal common subsequences. *Commun. ACM*, 18(6):341–343, June 1975.
- [KA90] S Karlin and S F Altschul. Methods for assessing the statistical significance of molecular sequence features by using general scoring schemes. *Proceedings of the National Academy of Sciences*, 87(6):2264–2268, 1990.
- [KBD⁺03] Donna Karolchik, Robert Baertsch, Mark Diekhans, Terrence S. Furey, A Hinrichs, YT Lu, Krishna M. Roskin, M Schwartz, Charles W. Sugnet, Daryl J Thomas, et al. The ucsc genome browser database. *Nucleic acids research*, 31(1):51–54, 2003.
- [MB06] Alexander Yu. Mitrophanov and Mark Borodovsky. Statistical significance in biological sequence analysis. *Briefings in Bioinformatics*, pages 2–24, 2006.
- [NJ02] Andrew Y Ng and Michael I Jordan. On discriminative vs. generative classifiers: A comparison of logistic regression and naive bayes. *Advances in neural information processing systems*, 2:841–848, 2002.
- [NVB13] Michal Nánási, Tomáš Vinař, and Broňa Brejová. Probabilistic approaches to alignment with tandem repeats. In *Algorithms in Bioinformatics*, pages 287–299. Springer, 2013.
- [NW70] Saul B. Needleman and Christian D. Wunsch. A general method applicable to the search for similarities in the amino acid sequence of two proteins. *Journal of Molecular Biology*, 48(3):443 – 453, 1970.
- [Sri] Sargur N. Srihari. Machine Learning: Generative and Discriminative Models.

- <http://www.cedar.buffalo.edu/~srihari/CSE574/Discriminative-Generative.pdf>. [Online; accessed 14-Jan-2013].
- [Sut07] Ivan Sutóris. Tvorba klasifikačných stromov pri procese data miningu, 2007.
- [SW81] T.F. Smith and M.S. Waterman. Identification of common molecular subsequences. *Journal of Molecular Biology*, 147(1):195 – 197, 1981.
- [Wik14a] Wikipedia. Deoxyribonukleová kyselina.
http://sk.wikipedia.org/wiki/Deoxyribonukleov%C3%A1_kyselina, 2014. [Online; accessed 3-May-2014].
- [Wik14b] Wikipedia. Kernel density estimation.
http://en.wikipedia.org/wiki/Kernel_density_estimation, 2014. [Online; accessed 1-May-2014].
- [Wik14c] Wikipedia. Maximum likelihood.
http://en.wikipedia.org/wiki/Maximum_likelihood, 2014. [Online; accessed 17-Apr-2014].
- [Wik14d] Wikipedia. Maximum likelihood.
http://en.wikipedia.org/wiki/ID3_algorithm, 2014. [Online; accessed 26-Apr-2014].
- [Wik14e] Wikipedia. Random forest.
http://en.wikipedia.org/wiki/Random_forest, 2014. [Online; accessed 26-Apr-2014].
- [YJEP07] Chun-Nam Yu, Thorsten Joachims, Ron Elber, and Jaroslaw Pillardy. Support vector training of protein alignment models. In Terry Speed and Haiyan Huang, editors, *Research in Computational Molecular Biology*, volume 4453 of *Lecture Notes in Computer Science*, pages 253–267. Springer Berlin / Heidelberg, 2007. 10.1007/978-3-540-71681-5_18.