

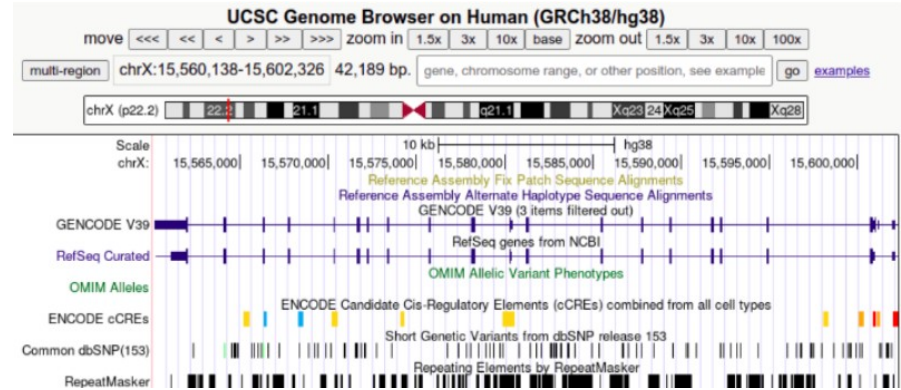


Modelovanie genomických anotácií pomocou diskretných distribúcií fázového typu

Hana Derková, školiteľ: Mgr. Askar Gafurov, PhD.

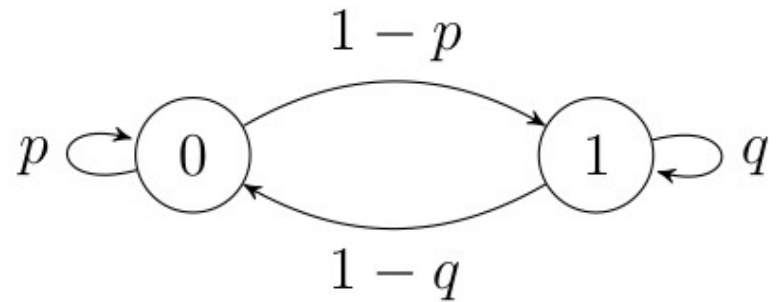
Genomická anotácia

- Kóduje významné časti genómu
- Súbor intervalov
- Významný prekryv -> biologická súvislosť



Miera štatistickej významnosti

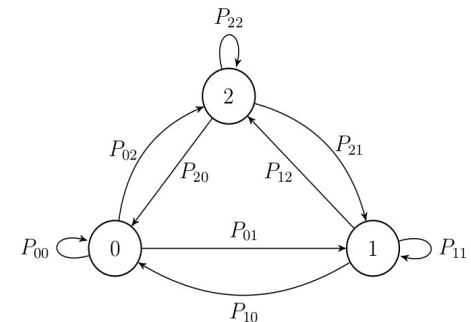
- Dve anotácie *reference* a *query*
- Nulová hypotéza : *query* je generovaná dvojstavovým Markovovským reťazcom
- P-hodnota vypočítaná v $O(m^2 + n)$ - nezávislá od dĺžky genómu
- Problém : distribúcia dĺžok je geometrická



Absorpčný Markovovský reťazec a diskretná distribúcia fázového typu

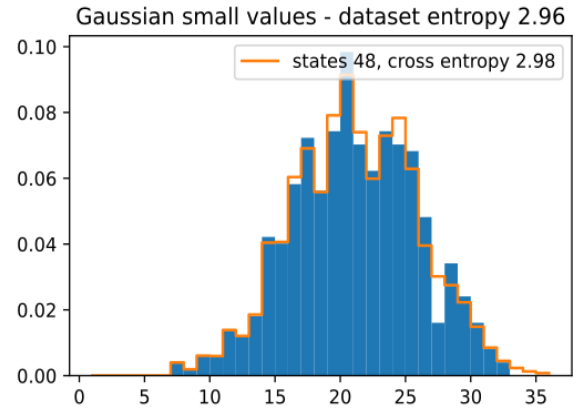
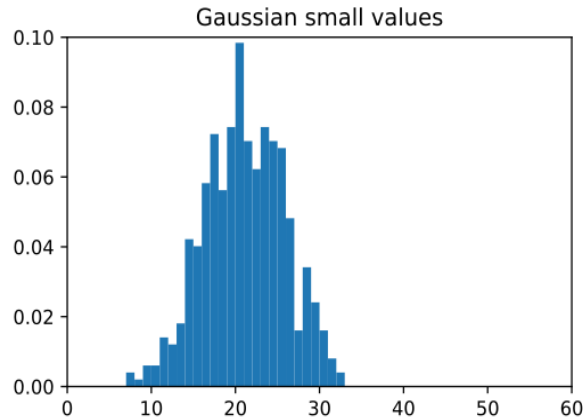
- Markovovský reťazec - stavové automaty
- Absorbčný Markovovský reťazec - každý stav vie dosiahnuť absorbčný stav
- ČAS ABSORPCIE = # krokov pred dosiahnutím absorbčného stavu
 - Tvorí triedu **diskrétnych distribúcií fázového typu**

$$P = \begin{bmatrix} P_{00} & P_{01} & P_{02} \\ P_{10} & P_{11} & P_{12} \\ P_{20} & P_{21} & P_{22} \end{bmatrix}$$



Úloha

- Pre danú distribúciu dĺžok intervalov / medzier, nájsť hodnoty pravdepodobnosti prechodov tak, aby sa výsledná distribúcia čo najviac podobala distribúcii dĺžok



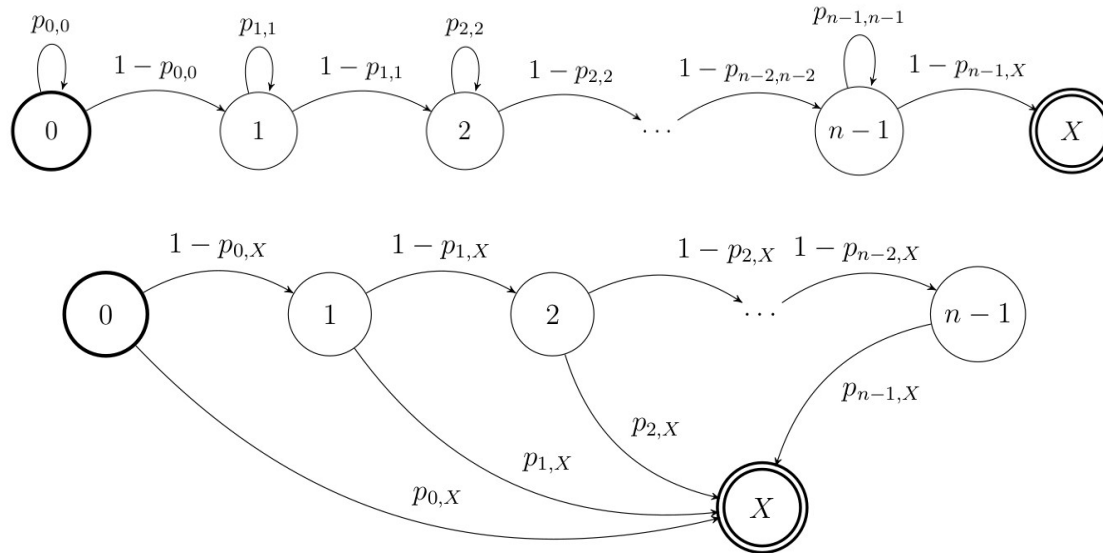


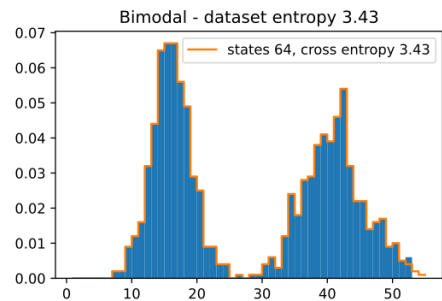
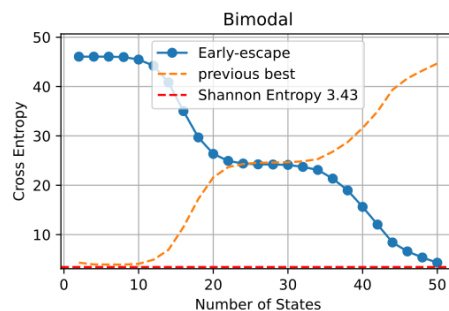
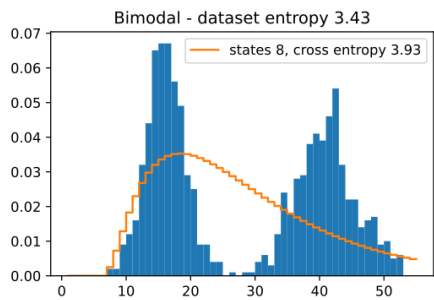
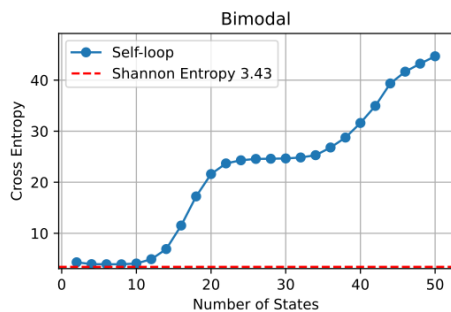
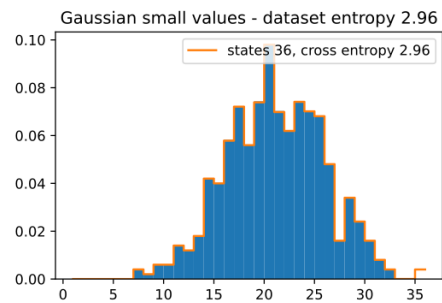
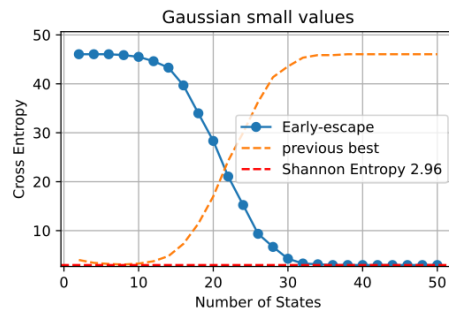
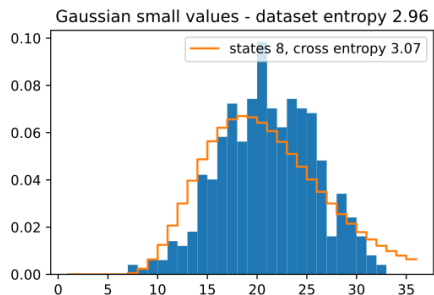
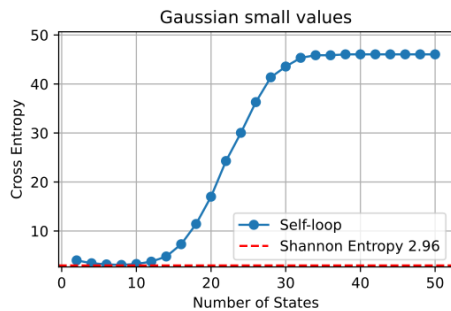
Stratová funkcia a parametre

- Minimalizovať negative log likelihood
- Cross entropia - diskretná fázová distribúcia a distribúcia dát
- Shanonova entropia datasetu
- Scipy - optimizer minimize
- Parametre - pravdepodobnosti prechodov
- Ďalšie parametre - počet stavov Markovovského reťazca, architektúra Markovovského reťazca

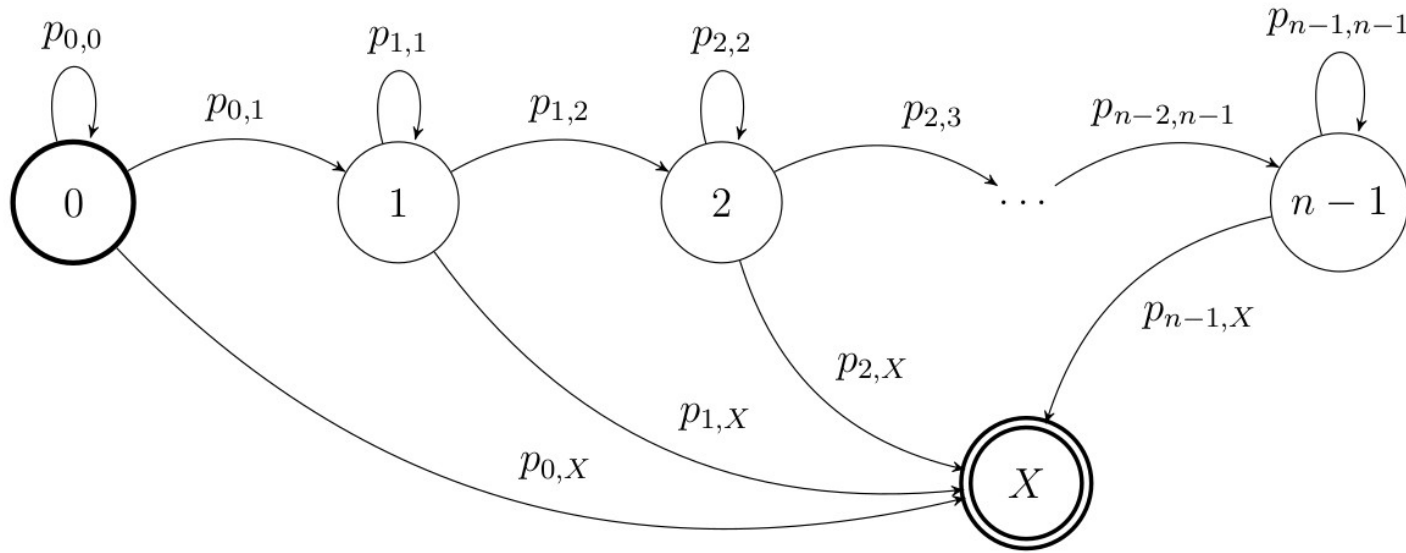
$$-\sum_{d \in data} \ln P(d|\theta)$$

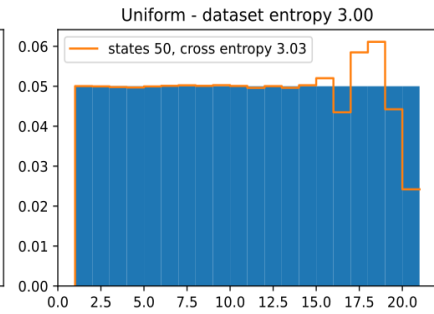
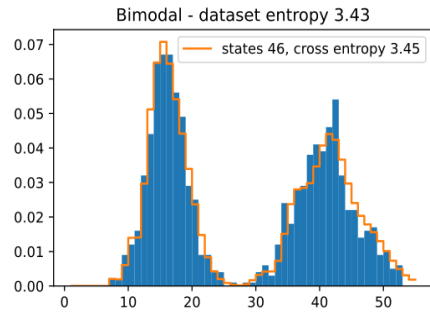
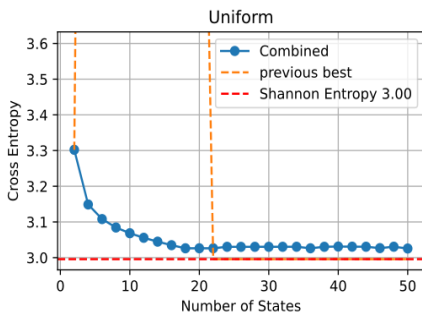
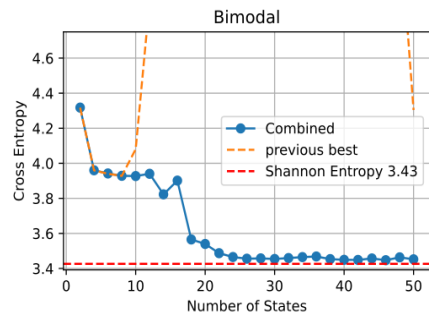
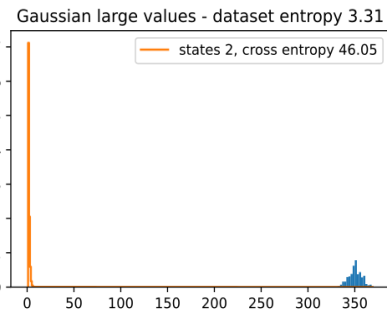
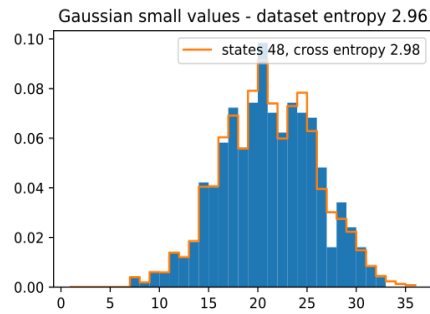
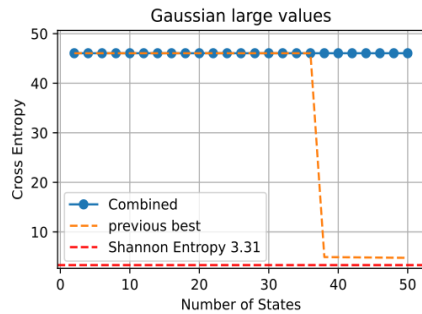
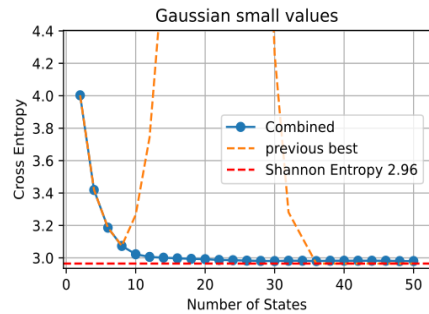
Self-loop a early-escape architektúra





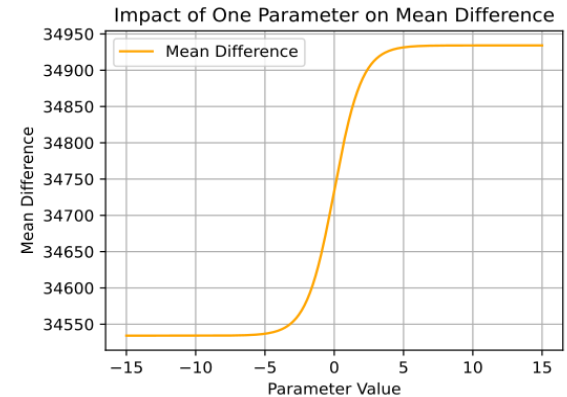
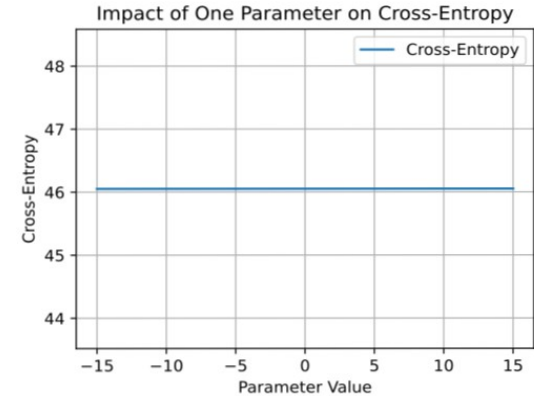
Combined architektúra



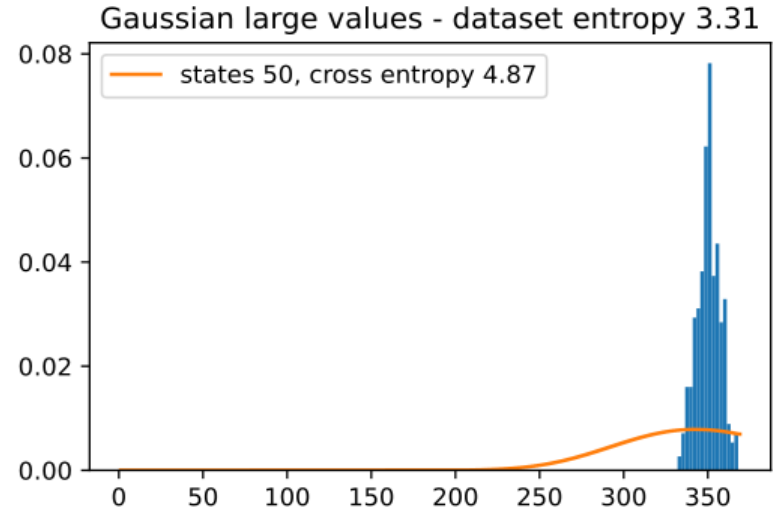
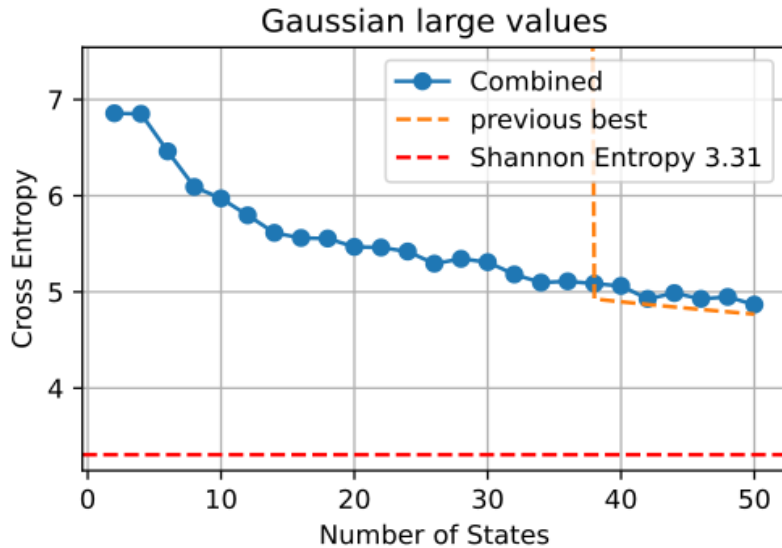


Technika - Mean shifting

- Nenatrérovaná architektúra pre veľké hodnoty
- Zastavovacia podmienka : projected gradient is too small
- Posunieme masu PHD bližšie ku strednej hodnote distribúcii dát
- Cieľová funkcia je rozdiel stredných hodnôt
- Následne pokračujeme v tréningu s výslednou inicializáciou



Efekt mean shiftingu





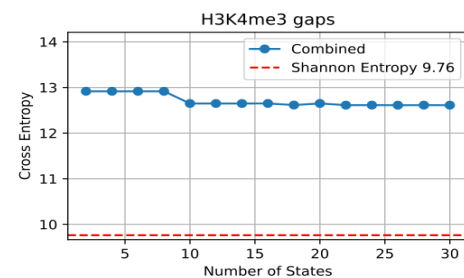
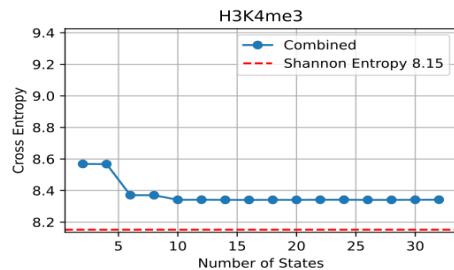
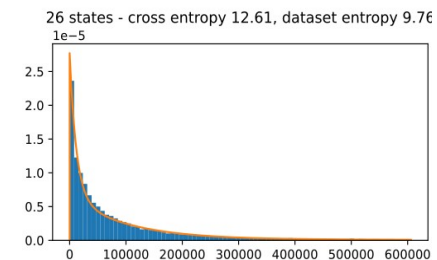
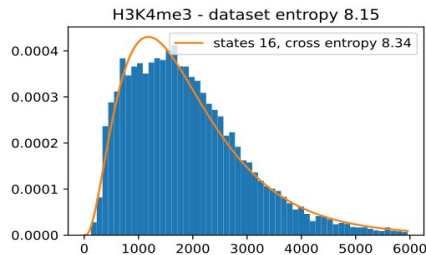
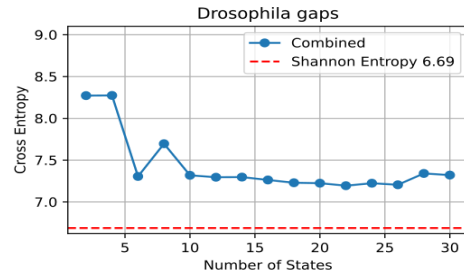
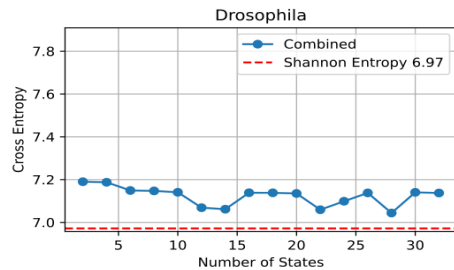
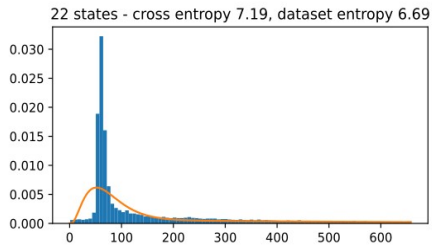
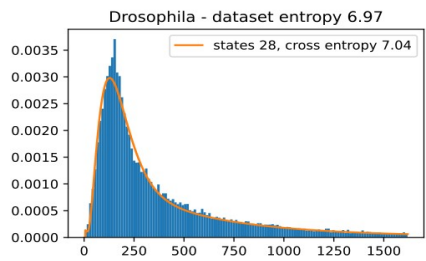
Modelovanie reálnych dát

- Combined architektúra + mean shifting
- Subsampling technika - 2000 vzoriek
- Hirt, Gains inc, H3K4me3, drosophila melanogaster exóny

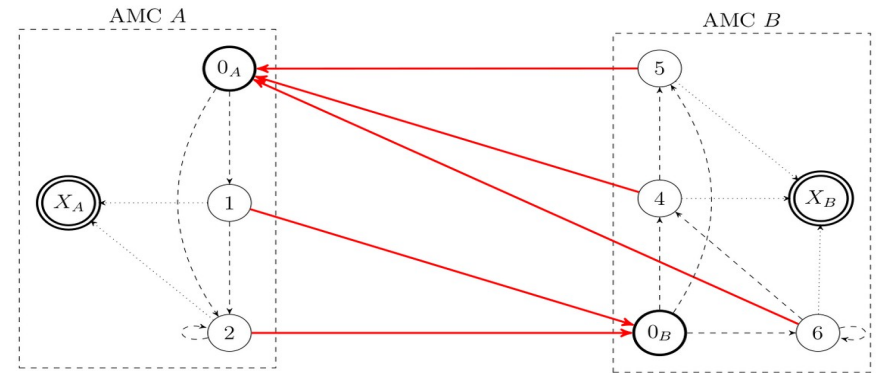
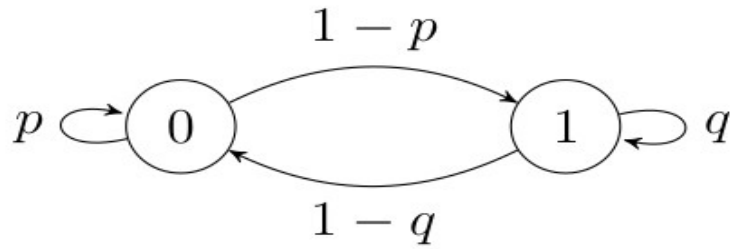
Dataset	Mean Gap Size	Range of Gap Lengths	Total Samples
Drosophila exons	1587	1-710734	63948
Gains inc	952046	4-30821993	3132
H3K4me3	156902	134-30542789	19333
Hirt	16296782	92740-87181600	116

Dataset	Mean	Range of interval lengths	Number of intervals
Drosophila exons	488	1-28074	188467
Gains inc	35581	73-3002725	3132
H3K4me3	1937	124-22590	19358
Hirt	1966196	6-24019400	247

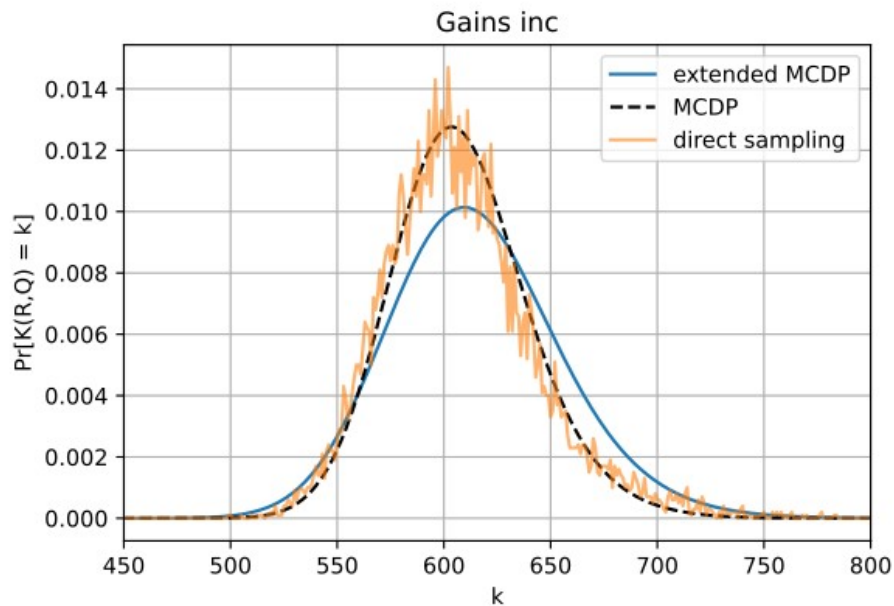
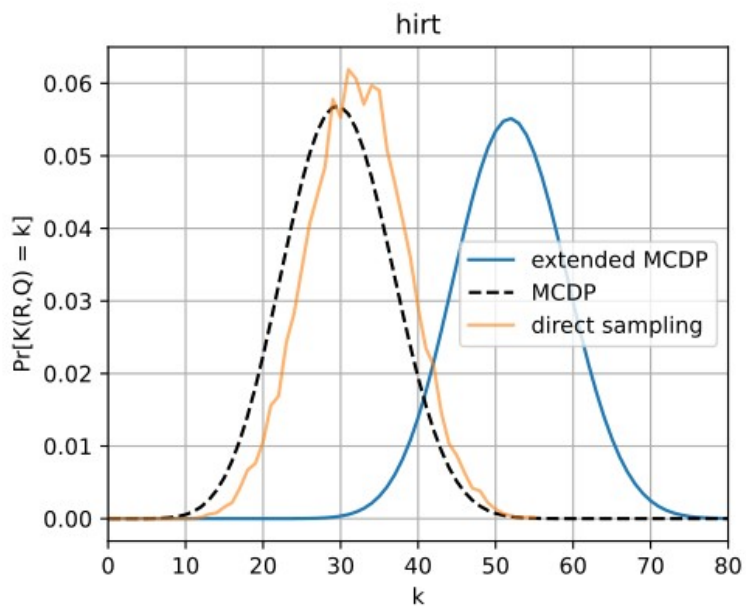
Reálne dáta výsledky



Extended MCDP



Výsledky





Ďakujem za pozornosť



Pri práci s reálnymi dátami vzorkujete iba podmnožinu dát. Ako závisí zložitosť výpočtu cross entropie vo vašom kóde od počtu stavov S , počtu dátových bodov N , počtu unikátnych dátových bodov U a maximálnej dĺžky intervalu v dátach L ?

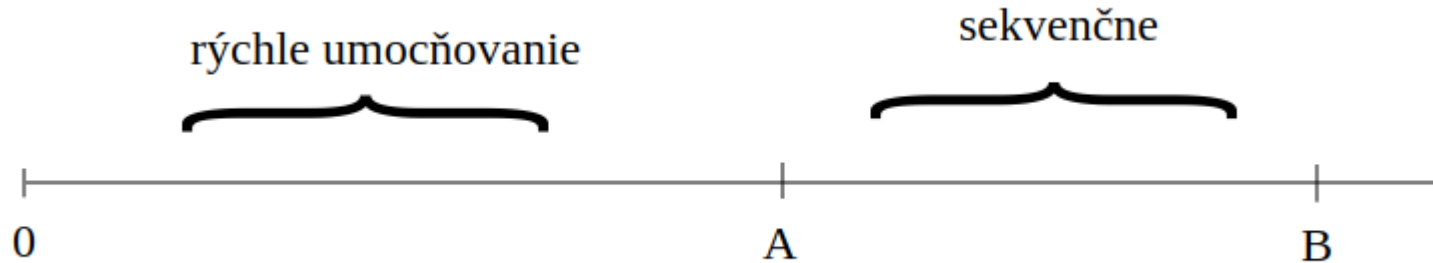
$$-\sum_{d \in \text{data}} \ln P(d|\theta) \quad O(S^m \cdot U \cdot \log L + N)$$



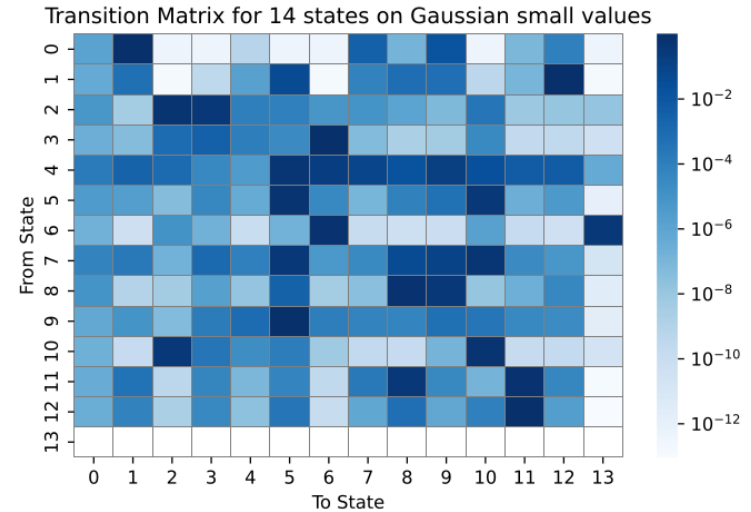
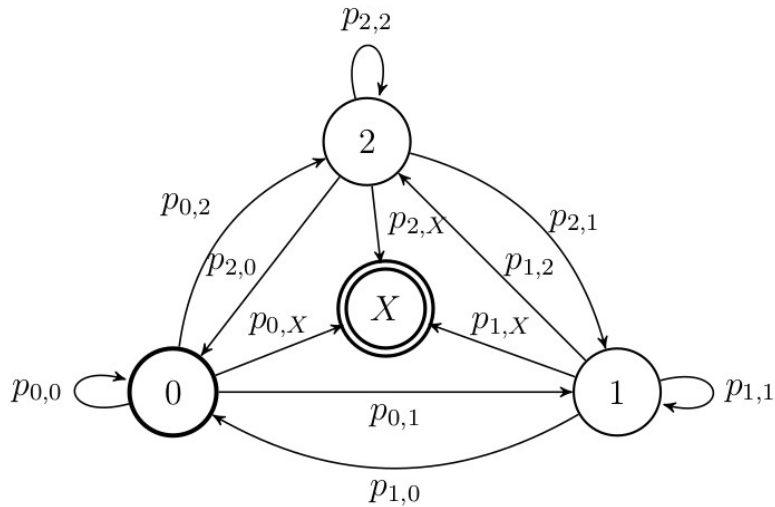
Ak by ste v predspracovaní zlúčili rovnaké hodnoty v dátach do jednej skupiny, vedeli by ste výpočet cross-entropie upraviť tak, aby nezávisel od N ale len od U?

$$O(S^m \cdot U \cdot \log L)$$

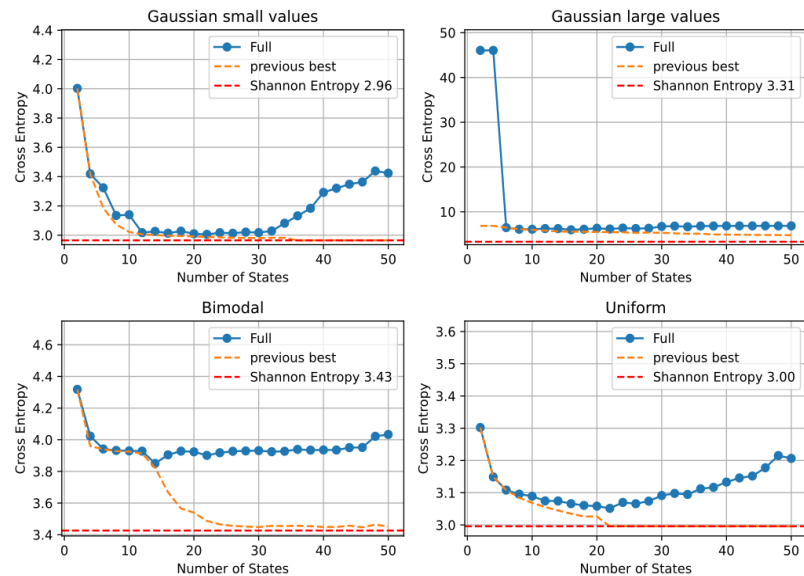
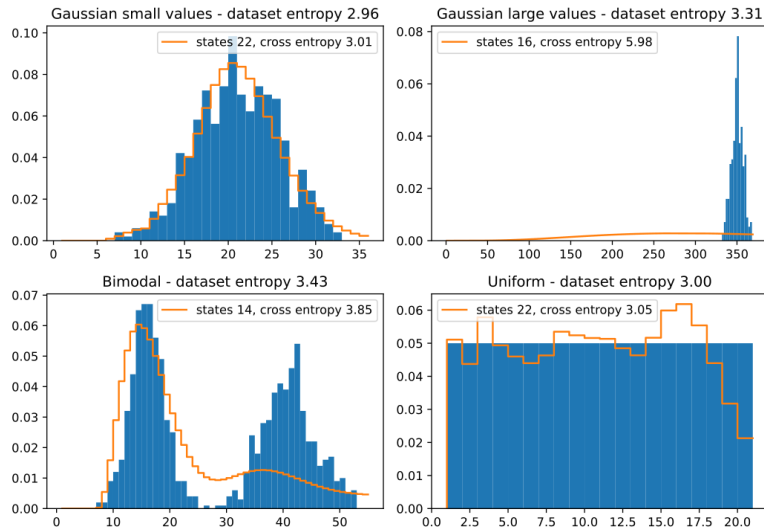
Ako by ste upravili výpočet, aby vedel rýchlo spracovať prípad, keď vaše dáta pochádzajú z intervalu $[A,B]$, pričom väčšina celých čísel z tohto intervalu sa v dátach nachádza aspoň raz?



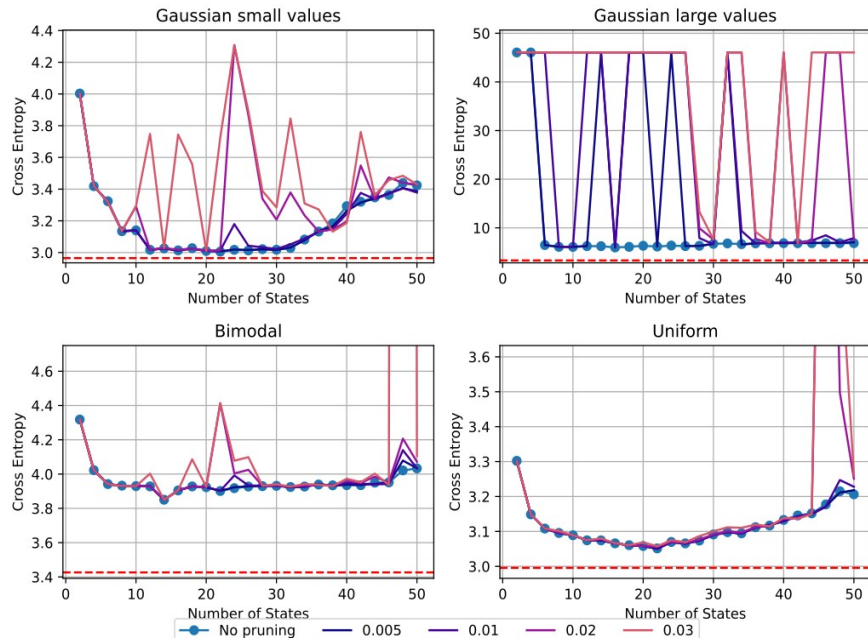
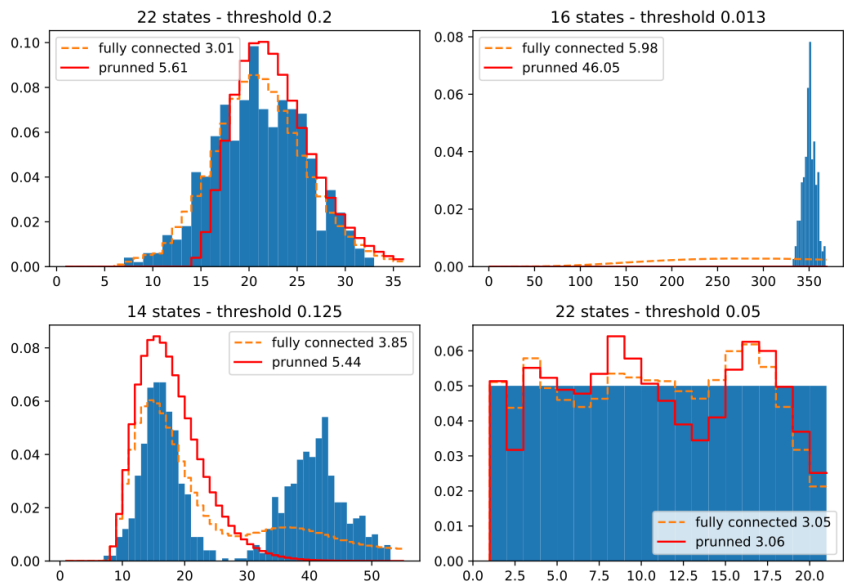
Fully connected architektúra a pruning



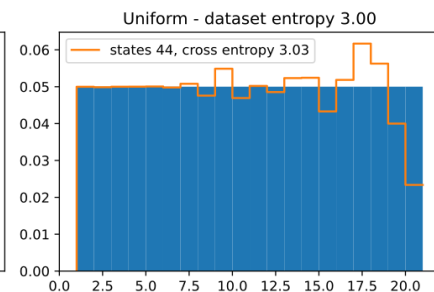
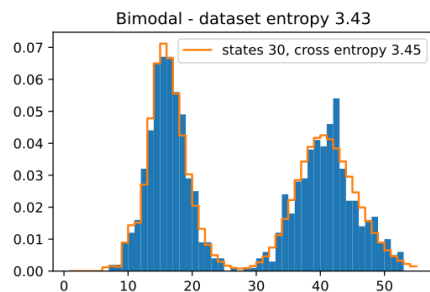
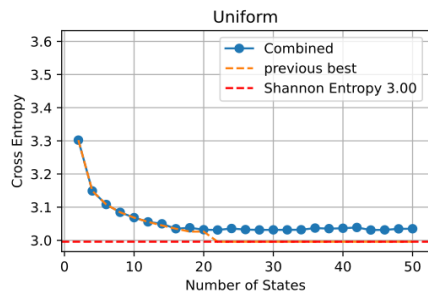
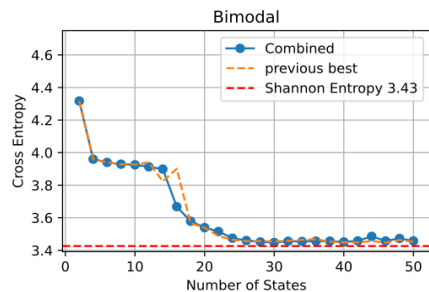
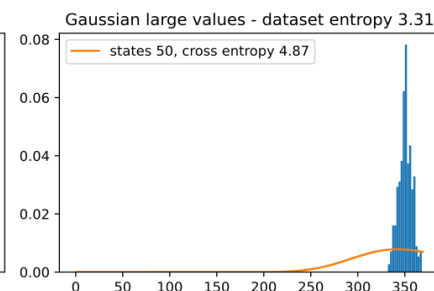
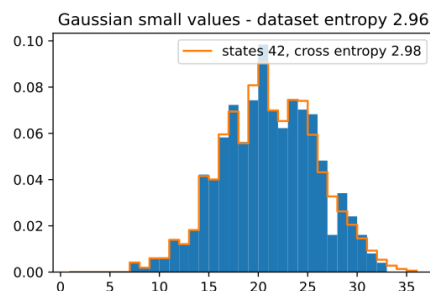
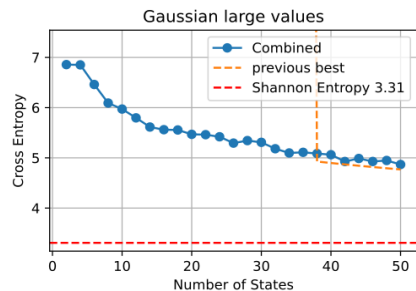
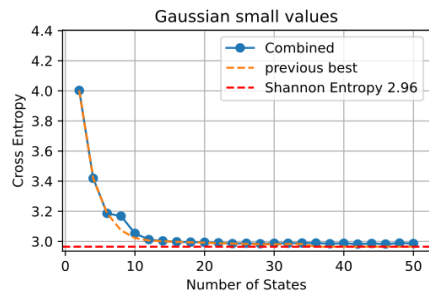
Výsledky pre fully-connected architektúru



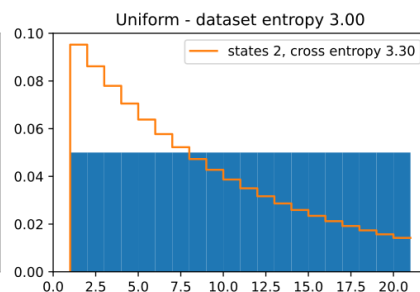
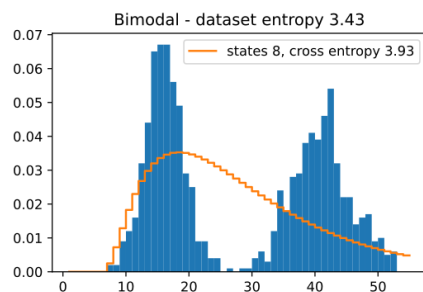
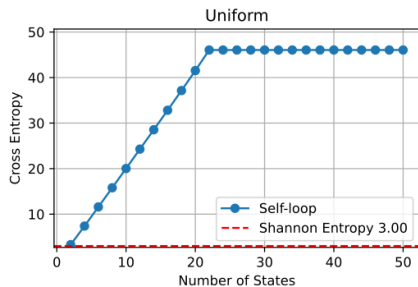
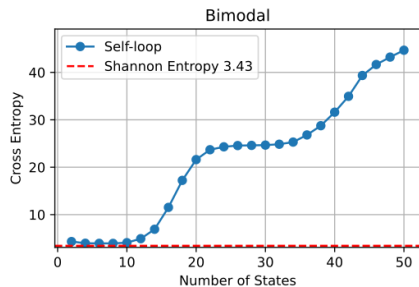
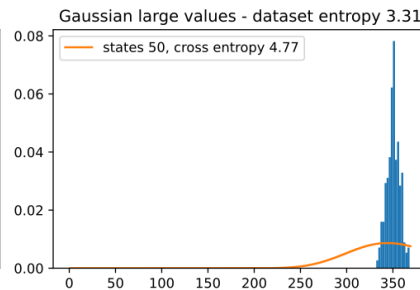
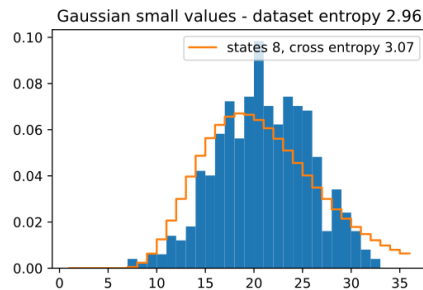
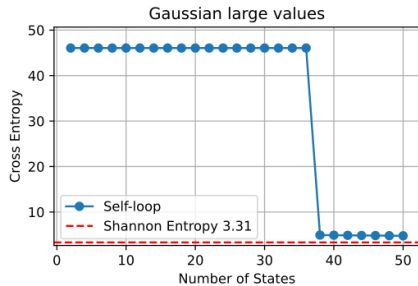
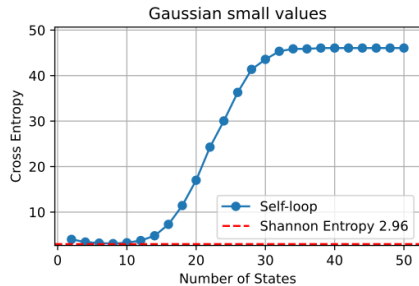
Pruning efekt



Všetky výsledky pre mean shifting



Všetky výsledky pre self-loop architektúru



Všetky výsledky pre early-escape architektúru

