

## 1. Motywacja

Aby białka pełniły właściwie swoje funkcje, muszą znajdować się w odpowiednim regionie komórki. „Adres” miejsca przeznaczenia białka znajduje się w obrębie jego sekwencji, w postaci tak zwanego sygnału kierującego. Do takich sygnałów należą peptydy sygnałowe, które są odpowiedzialne za kierowanie białek do systemu błon wewnętrznych komórki i w konsekwencji za ich wydzielanie na zewnątrz. Białka wyposażone w peptydy sygnałowe stanowią bardzo duży udział w różnorodności wszystkich białek w komórce i odgrywają istotną rolę w metabolizmie, tworzeniu struktur tkankowych, odporności i regulacji funkcji organizmu. Ponadto, przez wiele białek przez system błon wewnętrznych jest niezbędne do przyjmowania przez nie właściwej struktury i różnorodnych modyfikacji, jak glikozylacja (dodawanie reszt cukrowych) i fosforylacja. Peptydy sygnałowe występują na początku białka i mają długość zwykle 20-30 reszt aminokwasowych. Powszechny model typowych peptydów sygnałowych zakłada, że składają się z trzech części o różnym składzie aminokwasowym z czego szczególnie istotny jest region z dużą ilością aminokwasów hydrofobowych. Na końcu znajduje się miejsce odcinania peptydu od dojrzałego białka po wykonaniu przez niego funkcji.

Dotychczasowe programy rozpoznające peptydy sygnałowe były głównie oparte na sekwencjach najczęściej badanych gatunków jak drożdże, rośliny i ssaki. Jednakże nasze analizy dużej liczby sekwencji reprezentujących wiele różnorodnych grup organizmów gatunków wskazują na dużą zmienność peptydów sygnałowych. W chwili obecnej nie ma uniwersalnego modelu peptydu sygnałowego. Popularne programy przewidujące peptydy sygnałowe są z reguły oparte na metodach uczenia maszynowego, które uniemożliwiają badaczowi zapoznanie się z cechami peptydów służąc do ich rozpoznawania. Ponadto rozważania tego rodzaju nie są w stanie oddać pełnego znaczenia peptydów sygnałowych i nie rozpoznają nietypowych przypadków. Nie dostarczają one również odpowiedzi na pytania dotyczące budowy i zmienności międzygatunkowej peptydów sygnałowych, co utrudnia stosowanie ich do projektowania nowych peptydów.

## 2. Cel pracy

Dlatego celem naszego projektu jest stworzenie modelu statystycznego opisu tego peptydu sygnałowego, który będzie uwzględniał wiedzę o ich organizacji, składzie aminokwasowym oraz zmienności między różnymi grupami taksonomicznymi i typami peptydów. Składający się z reguł decyzyjnych model będzie mógł być również stosowany do analizy i wykrywania peptydów sygnałowych oraz do ich projektowania ich sztucznych odpowiedników.

## 3. Metoda badawcza

Model peptydu sygnałowego proponowany w tym projekcie będzie oparty na tzw. ukrytych modelach semi-Markowa (HSMMS). W najprostszym ukrytym modelu Markowa obliczy na podstawie składu aminokwasowego badanej sekwencji do jakiego regionu peptydu sygnałowego ona należy. Dużą zaletą tej metody jest możliwość uwzględnienia wiedzy biologicznej. Informacje o układzie regionów, ich długości oraz właściwości fizykochemiczne tworzących je aminokwasów dodane do modelu wnoszą istotny wkład do tworzenia reguł decyzyjnych. Stworzony w ten sposób model jest bardzo elastyczny i skutecznie opisuje również nietypowe peptydy sygnałowe, które posiadają różne domeny funkcjonalne albo pochodzą z białek specyficznych grup taksonomicznych takich jak np. pasożyty. Dzięki tym właściwościom model oparty na HSMMS nie musi być szczególnie douczony z powodu ciągłego wzrostu liczby sekwencji z peptydami sygnałowymi w bazach danych. Nowa metoda kodowania aminokwasów oparta na podobieństwach ich cech fizykochemicznych umożliwia wydobywanie z sekwencji peptydu sygnałowego najbardziej istotnych informacji. Dzięki temu możliwe jest tworzenie modeli peptydu sygnałowego dla nielicznych i specyficznych grup białek.

Stworzony model zostanie sprawdzony czy skutecznie wykrywa peptydy sygnałowe. Aby dokonać rzetelnego porównania modeli planujemy określić obszary stosowalności badanych algorytmów i stworzyć meta-predyktora peptydów sygnałowych wykorzystujący wiedzę z wielu dostępnych programów.

## 3. Wstępne wyniki

Pierwszy ze stworzonych modeli nazwany roboczo signalHsmm wykazał najwyższą wartość współczynnika  $AUC=0,98$  w porównaniu z innymi programami (ang. Area Under the Curve, oznacza w najprostszym ukrytym wyznacznik jako skuteczności przewidywania w skali od 0 do 1, gdzie 0 to całkowicie niewłaściwa predykcja, 0,5 losowa, a 1 idealna). Ponadto model ten efektywnie wykrywał peptydy sygnałowe nawet po trenowaniu na małych zbiorach sekwencji. Przykładowo, wersja programu wyuczona na skrajnie małym zbiorze danych z 1987 roku (336 sekwencji) zidentyfikowała tak samo skutecznie peptydy jak wersja wyuczona na o wiele większym zbiorze danych (2311 sekwencji).

Skutecznie stworzonego modelu sprawdzono dla peptydów sygnałowych z różnymi słabo poznanymi i reprezentowanymi w bazach danych grupami taksonomicznymi. Przykładowo, dla peptydów pasożytniczych zarodźca malarii i spokrewnionych rodzajów,  $AUC$  wynosi 0.92, istotnie lepiej niż w przypadku popularnych programów (0.84).

## 4. Wpływ rezultatów

Ostateczny model będzie zaimplementowany w formie programu komputerowego i stanie się użyteczny zarówno dla biologów eksperymentalnych, jak i bioinformatyków chcących analizować peptydy sygnałowe. Dodatkowo, opracowany model może zostać użyty do przyspieszenia procesu projektowania sztucznych peptydów sygnałowych w kosztownych i czasochłonnych badaniach laboratoryjnych. Nasz model peptydu sygnałowego będzie również opublikowany jako narzędzie do wykrywania peptydów sygnałowych.

Białka wyposażone w peptydy sygnałowe pełnią istotne funkcje w komórce, a zaburzenie ich głównej roli związanej z importem białek do siateczki wewnętrznej jest przyczyną wielu chorób. Proponowany algorytm wykrywania skutecznie takie peptydy sygnałowe pozwoli znacznie zmniejszyć koszty i przyspieszyć opracowywanie odpowiednich leków zorientowanych na te peptydy. Przykładem mogą być białka sekrecyjne pasożytów takich jak zarodźca malarii, które nie są dostatecznie rozpoznawane przez istniejące oprogramowanie.