

# POLITECHNIKA KRAKOWSKA IM. TADEUSZA KOŚCIUSZKI

## KARTA PRZEDMIOTU

obowiązuje studentów rozpoczynających studia w roku akademickim 2023/2024

Wydział Informatyki i Telekomunikacji

Kierunek studiów: Matematyka

Profil: Ogólnoakademicki

Forma studiów: stacjonarne

Kod kierunku: M

Stopień studiów: II

Specjalności: Modelowanie matematyczne, Matematyka w finansach i ekonomii

### 1 INFORMACJE O PRZEDMIOCIE

NAZWA PRZEDMIOTU	Zaawansowana eksploracja dużych zbiorów danych
NAZWA PRZEDMIOTU W JĘZYKU ANGIELSKIM	Advanced exploration of big datasets
KOD PRZEDMIOTU	WiT M oIIS C11 23/24
KATEGORIA PRZEDMIOTU	Przedmioty kierunkowe
LICZBA PUNKTÓW ECTS	4.00
SEMESTRY	4

### 2 RODZAJ ZAJĘĆ, LICZBA GODZIN W PLANIE STUDIÓW

SEMESTR	WYKŁAD	ĆWICZENIA	LABORATORIUM	LABORATORIUM KOMPUTERO- WE	SEMINARIUM	PROJEKT
4	30	0	0	0	0	30

### 3 CELE PRZEDMIOTU

**Cel 1** Zapoznanie studentów z nowymi, ważniejszymi algorytmami i metodami stosowanymi obecnie do przechowywania, przetwarzania, analizy, modelowania i wizualizacji olbrzymich ilości danych.

**Cel 2** Zapoznanie studentów z ważniejszym oprogramowaniem stosowanym do przechowywania, przetwarzania, analizy, modelowania i wizualizacji olbrzymich ilości danych.

**Cel 3** Celem będzie nabycie umiejętności zespołowej (parami) realizacji projektów, pomagającej w utrwaleniu praktycznych umiejętności, aby na podstawie zdobytej wiedzy, w sposób staranny i terminowy realizować te projekty, być gotowym do rozwiązywania problemów ze wspomnianego zakresu wiedzy, oraz utrwalić umiejętność pisanie odpowiedniej dokumentacji zrealizowanego projektu, a także zwiększenia swych umiejętności klarownego przedstawienia i wytłumaczenia swoich ukończonych projektów wobec całej klasy, wraz z odpowiadaniem na zadawane im przez pozostałe grupy pytań dotyczących realizacji tych projektów.

## 4 WYMAGANIA WSTĘPNE W ZAKRESIE WIEDZY, UMIEJĘTNOŚCI I INNYCH KOMPETENCJI

1 Podstawowa wiedza ze statystyki i z rachunku prawdopodobieństwa.

2 Podstawowa wiedza z algebry liniowej.

## 5 EFEKTY KSZTAŁCENIA

**EK1 Wiedza** Student będzie potrafił wytłumaczyć działanie oraz zinterpretować wyniki ważniejszych algorytmów stosowanych przy przetwarzaniu, analizie i wizualizacji danych.

**EK2 Umiejętności** Student utrwali swe umiejętności zastosowania do przechowywania, przetwarzania, analizy i wizualizacji olbrzymich ilości danych ważniejszych algorytmów (tzw. funkcji) zawartych w wybranych bibliotekach (pakietach) języka R i środowiska zintegrowanego RStudio.

**EK3 Umiejętności** Student utrwali swe umiejętności samodzielnego programowania w języku R w środowisku RStudio oraz łączenia się internetowo poprzez RStudio zarówno z uniwersalnym silnikiem dla Big data (tj. ze SPARKiem), jak i z darmowymi serwerami RStudio, wspomagającymi przetwarzanie dużych zbiorów danych w Chmurze.

**EK4 Kompetencje społeczne** Student utrwali swe umiejętności wzajemnej współpracy w małych zespołach, w samokształceniu, oraz we właściwej komunikacji z innymi grupami i z nauczycielem.

## 6 TREŚCI PROGRAMOWE

PROJEKT		
LP	TEMATYKA ZAJĘĆ OPIS SZCZEGÓŁOWY BLOKÓW TEMATYCZNYCH	LICZBA GODZIN
P1	Treści programowe 1. Projekt pierwszy. Dokonanie analizy porównawczej wybranych (dla każdego zespołu innych) modeli predykcyjnych dostępnych w pakietach tidymodels oraz parsnip, w języku R. Każdy zespół będzie miał inny zbiór danych. Analiza porównawcza będzie opierać się o kluczowe czynniki, takie jak: wskazania wybranych metryk oceny modelu, np.: RMSE, MAPE, MASE, oraz MAE, najczęściej dążąc do minimalizacji wartości każdej z nich. Dane prognozowane będą porównywane z danymi rzeczywistymi, dostępnymi w zbiorze testującym. Zespoły przygotują sprawozdania z projektu nr 1 oraz przedstawia je wszystkim pozostałym grupom.	8
P2	Treści programowe 2. Projekt drugi. Rozbudowa poprzednio wykonanej analizy o zrealizowanie w niej kroku dostrajania (tuning) posiadanych przez dany model hiperparametrów, oraz porównanie otrzymanych nowych wartości czterech metryk. Zespoły przygotują sprawozdania z projektu nr 2 oraz przedstawia je wszystkim pozostałym grupom.	7

PROJEKT		
LP	TEMATYKA ZAJĘĆ OPIS SZCZEGÓŁOWY BLOKÓW TEMATYCZNYCH	LICZBA GODZIN
<b>P3</b>	Treści programowe 3. Projekt trzeci. Rozbudowa poprzednio wykonanych analiz o zrealizowanie w swej analizie kroku walidacji krzyżowej i wybranie modelu najbardziej dostrojonego o najbardziej optymalnych nowych wartościach czterech metryk. Zespoły przygotowują sprawozdania z projektu nr 3 oraz przedstawia je wszystkim pozostałym grupom.	7
<b>P4</b>	Treści programowe 4. Projekt czwarty. Modyfikacja poprzednio wykonanych kroków i zastąpienie ich funkcjami automatyzującymi cały wspomniany w poprzednich trzech analizach proces generacji modeli, ich dostrajania i walidacji krzyżowej, wykorzystując w tym celu dostępne w rodzinie pakietów tidymodels funkcje scalające poszczególne kroki wspomnianych trzech analiz. Zespoły przygotowują sprawozdania z projektu nr 4 oraz przedstawia je wszystkim pozostałym grupom.	8

WYKŁAD		
LP	TEMATYKA ZAJĘĆ OPIS SZCZEGÓŁOWY BLOKÓW TEMATYCZNYCH	LICZBA GODZIN
<b>W1</b>	Poznanie nowych elementów architektury i nowych funkcjonalności, dodanych do ekosystemu Apache SPARK i do jego podsystemów. Poznanie nowych funkcji dodanych do pakietu sparklyr (autorstwa twórców środowiska RStudio), pozwalającego na łączenie się ze SPARKiem oraz korzystanie z pakietów SPARKa w środowisku skalowalnym i rozproszonym, będąc w środowisku programistycznym RStudio.	4
<b>W2</b>	Treści programowe 2. Poznanie nowych rozproszonych algorytmów dodanych do bogatego asortymentu biblioteki spark.ml w SPARKu. Poznanie głównych możliwości i funkcjonalności i sposobu wykorzystania tych algorytmów w chmurze oraz możliwości współpracy tych algorytmów ze środowiskiem RStudio.	4
<b>W3</b>	Treści programowe 3. Poznanie nowych funkcjonalności rozszerzonego pakietu rsparkling, dostarczającego możliwości połączenia się z dużą liczbą najnowszych rozproszonych algorytmów uczenia maszynowego w środowisku H2O, mocno współpracującym ze SPARKiem. Z platformy H2O korzysta ponad 8000 największych firm na świecie. Wszystko to jest możliwe do sterowania, będąc w środowisku programistycznym RStudio.	8
<b>W4</b>	Treści programowe 4. Poznanie nowych pakietów służących do modelowania tekstu, eksploracji sieci społecznościowych (social networks) oraz nowych metod tworzenia dashboardów, współpracujących z rodziną pakietów tidyverse i tidymodels oraz zaznajomienie się z przykładami ich zastosowania.	4

WYKŁAD		
LP	TEMATYKA ZAJĘĆ OPIS SZCZEGÓŁOWY BLOKÓW TEMATYCZNYCH	LICZBA GODZIN
W5	Treści programowe 8. Poznanie nowych funkcjonalności w drugiej co do ważności (obok rodziny tidyverse) rodziny pakietów uczenia maszynowego "Tidymodels", której celem jest uproszczenie i zautomatyzowanie przetwarzania potokowego zadań w uczeniu maszynowym, oraz pomoc w zrównolegleniu całego ciągu procesów, w skład których wchodzi walidacja krzyżowa, wstępne przetwarzanie danych (dzięki pakietowi recipe, należącemu do tej rodziny) oraz dostrajanie (tuning) algorytmów uczenia maszynowego, polegające na przetwarzaniu olbrzymiej ilości wygenerowanych modeli dla różnych kombinacji wartości hiperparametrów tych algorytmów. Rodzina pakietów "Tidymodels" pozwala na wybranie pewnej liczby miar, służących do badania jakości wybranych modeli oraz na zautomatyzowanie procesu znalezienia najlepszego modelu o najlepszej jakości, gdy otrzymanie tego wyniku wymaga generacji olbrzymiej ilości modeli, a każdy z nich musi podlegać większości procesom wspomnianego przetwarzania potokowego.	10

## 7 NARZĘDZIA DYDAKTYCZNE

- N1 Wykłady (w przypadku realizacji zajęć w trybie zdalnym z wykorzystaniem stosownych narzędzi teleinformatycznych, np. MS TEAMS)
- N2 Praca w 2-3 osobowych grupkach (w przypadku realizacji zajęć w trybie zdalnym z wykorzystaniem stosownych narzędzi teleinformatycznych, np. MS TEAMS)
- N3 Prezentacje multimedialne (w przypadku realizacji zajęć w trybie zdalnym z wykorzystaniem stosownych narzędzi teleinformatycznych, np. MS TEAMS)
- N4 Konsultacje (w przypadku realizacji zajęć w trybie zdalnym z wykorzystaniem stosownych narzędzi teleinformatycznych, np. MS TEAMS)
- N5 Dyskusja (w przypadku realizacji zajęć w trybie zdalnym z wykorzystaniem stosownych narzędzi teleinformatycznych, np. MS TEAMS)

## 8 OBCIĄŻENIE PRACĄ STUDENTA

FORMA AKTYWNOŚCI	ŚREDNIA LICZBA GODZIN NA ZREALIZOWANIE AKTYWNOŚCI
<b>Godziny kontaktowe z nauczycielem akademickim, w tym:</b>	
Godziny wynikające z planu studiów	60
Konsultacje przedmiotowe	15
Egzaminy i zaliczenia w sesji	5
<b>Godziny bez udziału nauczyciela akademickiego wynikające z nakładu pracy studenta, w tym:</b>	
Przygotowanie się do zajęć, w tym studiowanie zalecanej literatury	10
Opracowanie wyników	10
Przygotowanie raportu, projektu, prezentacji, dyskusji	20
<b>SUMARYCZNA LICZBA GODZIN DLA PRZEDMIOTU WYNIKAJĄCA Z CAŁEGO NAKŁADU PRACY STUDENTA</b>	<b>120</b>
SUMARYCZNA LICZBA PUNKTÓW ECTS DLA PRZEDMIOTU	4.00

## 9 SPOSOBY OCENY

### OCENA FORMUJĄCA

**F1** Wysłanie każdego z czterech sprawozdań, dotyczących czterech projektów do prowadzącego zajęcia projektowe.

**F2** Prezentowanie każdego z czterech projektów przed pozostałymi zespołami w klasie.

### OCENA PODSUMOWUJĄCA

**P1** Średnia ważona ocen formujących

### WARUNKI ZALICZENIA PRZEDMIOTU

**W1** Uzyskanie odpowiedniej liczby punktów z czterech sprawozdań oraz z czterech prezentacji.

**W2** Spełnienie warunku obecności na obowiązkowych formach zajęć (dopuszczalna jedna nieobecność nieusprawiedliwiona na każdej z obowiązkowych form zajęć).

### OCENA AKTYWNOŚCI BEZ UDZIAŁU NAUCZYCIELA

**B1** Oddanie wszystkich zleconych do napisania sprawozdań wykonanych w środowisku programistycznym RStudio i Apache SPARK. Wszystkie sprawozdania będą punktowane.

**B2** Ocena za odpowiedzi ustne podczas zajęć.

**B3** Ocena za aktywność podczas wykonywania ćwiczeń praktycznych w klasie.

## KRYTERIA OCENY

EFEKT KSZTAŁCENIA 1	
NA OCENĘ 2.0	Student nie spełnia warunków określonych dla oceny 3.0
NA OCENĘ 3.0	Opanowanie zagadnień w stopniu powyżej 50%.
NA OCENĘ 3.5	Opanowanie zagadnień w stopniu powyżej 60%.
NA OCENĘ 4.0	Opanowanie zagadnień w stopniu powyżej 70%.
NA OCENĘ 4.5	Opanowanie zagadnień w stopniu powyżej 80%.
NA OCENĘ 5.0	Opanowanie zagadnień w stopniu powyżej 90%.
EFEKT KSZTAŁCENIA 2	
NA OCENĘ 2.0	Student nie spełnia warunków określonych dla oceny 3.0
NA OCENĘ 3.0	Opanowanie zagadnień w stopniu powyżej 50%.
NA OCENĘ 3.5	Opanowanie zagadnień w stopniu powyżej 60%.
NA OCENĘ 4.0	Opanowanie zagadnień w stopniu powyżej 70%.
NA OCENĘ 4.5	Opanowanie zagadnień w stopniu powyżej 80%.
NA OCENĘ 5.0	Opanowanie zagadnień w stopniu powyżej 90%.
EFEKT KSZTAŁCENIA 3	
NA OCENĘ 2.0	Student nie spełnia warunków określonych dla oceny 3.0
NA OCENĘ 3.0	Opanowanie zagadnień w stopniu powyżej 50%.
NA OCENĘ 3.5	Opanowanie zagadnień w stopniu powyżej 60%.
NA OCENĘ 4.0	Opanowanie zagadnień w stopniu powyżej 70%.
NA OCENĘ 4.5	Opanowanie zagadnień w stopniu powyżej 80%..
NA OCENĘ 5.0	Opanowanie zagadnień w stopniu powyżej 90%.
EFEKT KSZTAŁCENIA 4	
NA OCENĘ 2.0	Student nie spełnia warunków określonych dla oceny 3.0
NA OCENĘ 3.0	Opanowanie zagadnień w stopniu powyżej 50%.
NA OCENĘ 3.5	Opanowanie zagadnień w stopniu powyżej 60%.
NA OCENĘ 4.0	Opanowanie zagadnień w stopniu powyżej 70%.
NA OCENĘ 4.5	Opanowanie zagadnień w stopniu powyżej 80%.
NA OCENĘ 5.0	Opanowanie zagadnień w stopniu powyżej 90%.

## 10 MACIERZ REALIZACJI PRZEDMIOTU

EFEKT KSZTAŁCENIA	ODNIESIENIE DANEGO EFEKTU DO SZCZEGÓŁOWYCH EFEKTÓW ZDEFINIOWANYCH DLA PROGRAMU	CELE PRZEDMIOTU	TREŚCI PROGRAMOWE	NARZĘDZIA DYDAKTYCZNE	SPOSOBY OCENY
EK1	K_W01 K_W02 K_W04 K_W06 K_W07 K_W08 K_W10 K_W11 K_W12	Cel 1	P1 P2 P3 P4 W1 W2 W3 W4 W5	N1 N2 N3 N4 N5	F1 F2 P1
EK2	K_U02 K_U10 K_U11 K_U12 K_U13 K_U16 K_U20 K_U22	Cel 2	P1 P2 P3 P4 W1 W2 W3 W4 W5	N1 N2 N3 N4 N5	F1 F2 P1
EK3	K_U02 K_U10 K_U11 K_U12 K_U13 K_U15 K_U16 K_U20 K_U21 K_U22	Cel 2	P1 P2 P3 P4 W1 W2 W3 W4 W5	N1 N2 N3 N4 N5	F1 F2 P1
EK4	K_K01 K_K02 K_K03 K_K04 K_K05 K_K06 K_K07	Cel 3	P1 P2 P3 P4	N1 N2 N3 N4 N5	F1 F2 P1

## 11 WYKAZ LITERATURY

### LITERATURA PODSTAWOWA

- [1] Podstawowa dokumentacja wraz z bardzo licznymi przykładami: <https://spark.rstudio.com/guides/> ; Liczne przykłady: <https://stat545.com/index.html> ; <https://ggplot2-book.org/> ; <https://spark.apache.org/docs/latest/ml-guide.html> ; Tysiące przykładowych notebooków: <https://rpubs.com> ; <https://www.kaggle.com/code?language=R> ; <https://www.kaggle.com/>
- [2] Rozszerzenie książki ISLR (T. Hastie, R. Tibshirani) o rodzinę pakietów Tidymodels: <https://emilvitfeldt.github.io/ISLR-tidymodels-labs/index.html> ; Książki online: <https://bookdown.org/> ; <https://www.tnwr.org/> ; <https://smltar.com/> ; <https://www.tidytextmining.com/> ; <https://therinspark.com/> ; <http://r4ds.had.co.nz/>
- [3] Darmowy serwer online z zainstalowanym RStudio i z pakietami: <https://rdr.io/snippets/> ; Środowisko darmowe w Chmurze: <https://rstudio.cloud/> ; <https://databricks.com/try>
- [4] Dwie najważniejsze rodziny pakietów w RStudio: tidyverse i Tidymodels; tutorialy: <https://rpubs.com/cliex159/885971> ; <https://jhudatasience.org/tidyversecourse/model.html> ; <https://www.hackerearth.com/practice/machine-learning/machine-learning-algorithms/beginners-tutorial-on-xgboost-parameter-tuning-r/tutorial/> ; <https://www.startuptips.com/>



/https://www.startupengineer.io/\_repos/\_transfer/data\_science/ ; https://jhudatascience.org/tidyversecourse/model.ht  
; https://juliasilge.com/blog/xgboost-tune-volleyball/ ; https://www.alexbaeher.com/post/ml-tidymodels/  
; https://cran.rstudio.com/web/packages/workboots/vignettes/Getting-Started-with-workboots.html ; https://tanthiamhu  
ecosystem/ ; https://scienceofdata.org/2021/07/07/look-what-the-cat-dragged-in-catboost-with-tidymodels/  
; https://dniel.com/posts/tidymodels-intro/ ; https://www.r-bloggers.com/2020/05/using-xgboost-with-  
tidymodels/

## LITERATURA UZUPEŁNIAJĄCA

- [1] | Dokumentacja SPARK'a: <http://spark.apache.org/docs/latest/> ; Uczenie maszynowe <https://spark.apache.org/docs/latest/guide.html> ; <https://docs.databricks.com/> ; <https://sparkhub.databricks.com/resources/>
- [2] | Y. Zhao, "R and Data Mining: Examples and Case Studies", 2014, książka dostępna w WWW z licznymi materiałami, np. książka: <http://www.rdatamining.com/docs/introduction-to-data-mining-with-r> oraz przykłady w R: <http://www.rdatamining.com/examples> oraz wiele inn. wartosciowych plików: <http://www.rdatamining.com>

## 12 INFORMACJE O NAUCZYCIELACH AKADEMICKICH

### OSOBA ODPOWIEDZIALNA ZA KARTĘ

dr Barbara Borowik (kontakt: [bborowik@pk.edu.pl](mailto:bborowik@pk.edu.pl))

### OSOBY PROWADZĄCE PRZEDMIOT

- 1 dr Barbara Borowik (kontakt: [barbara.borowik@pk.edu.pl](mailto:barbara.borowik@pk.edu.pl))

## 13 ZATWIERDZENIE KARTY PRZEDMIOTU DO REALIZACJI

---

(miejscowość, data)

(odpowiedzialny za przedmiot)

(dziekan)

PRZYJMUJĘ DO REALIZACJI (data i podpisy osób prowadzących przedmiot)

.....