

## Program kursu *Data mining dla niezaawansowanych*

### Informacje wstępne

Liczba godzin	30
Język wykładowy	Polski
Punkty ECTS dla kształcenia na odległość	3
Wymagania wstępne	<ul style="list-style-type: none"><li>- Znajomość matematyki na poziomie szkoły średniej.</li><li>- Podstawowa umiejętność posługiwania się komputerem z systemem Windows.</li><li>- Samodyscyplina, umiejętność gospodarowania własnym czasem.</li><li>- Umiejętność czytania ze zrozumieniem tekstów napisanych językiem formalnym.</li><li>- Otwartość na nowe wyzwania.</li><li>- Świadomość konieczności samodoskonalenia się.</li><li>- Odwaga w zadawaniu pytań oraz wypowiedzaniu się na forum publicznym.</li></ul>
Skrócony opis	Duże zbiory danych pojawiają się aktualnie w każdym obszarze, w którym gromadzone są informacje. Oznacza to konieczność wykształcenia powszechniej sprawności w pracy z takimi danymi. Celem przedmiotu jest zapoznanie uczestników zajęć z najważniejszymi algorytmami data-miningu oraz wykształcenie umiejętności analizy danych z wykorzystaniem programu PS IMAGO PRO (IBM SPSS Statistics).
Całkowity nakład pracy	<ul style="list-style-type: none"><li>- Zajęcia prowadzone synchronicznie – 0 h,</li><li>- zapoznanie się z materiałami dostępnymi na platformie Moodle – 30 h,</li><li>- praca własna: studiowanie literatury, wykonywanie ćwiczeń, rozwiązywanie zadań, konsultacje z prowadzącymi zajęcia – 50 h,</li><li>- rozwiązywanie testów i zadań zaliczeniowych 10 h.</li></ul> Razem: 90 h (3 pkt. ECTS)

### Efekty uczenia się osiągnięte przez uczestnictwo w kursie

Wiedza	<p>W1. Rozumie potrzebę pozyskiwania wiedzy z danych. Zna podstawowe problemy eksploracji danych.</p> <p>W2. Zna wybrane algorytmy eksploracji danych i wie, które z nich stosują się do określonego typu zagadnień z tego zakresu.</p> <p>W3. Ma wiedzę na temat dostępnego oprogramowania służącego do eksploracji danych.</p>
Umiejętności	<p>U1. Potrafi znaleźć potrzebne dane w zbiorach danych ogólnie dostępnych, umie pobrać dane i poddać je analizie.</p> <p>U2. Umie zaproponować odpowiednie algorytmy eksploracji danych do konkretnego zagadnienia, w tym klasyfikacji, grupowania, szacowania i budowania</p>

	<p>reguł, oraz wyselekcjonować z ich użyciem najlepszy model.</p> <p>U3. Umie posługiwać się w stopniu podstawowym programem do eksploracji danych PS IMAGO PRO (IBM SPSS Statistics).</p> <p>U4. Potrafi przygotować raport z wynikami swoich analiz.</p>
Kompetencje społeczne	<p>K1. Potrafi sformułować problem eksploracji danych w sposób zrozumiały zarówno dla osób, z którymi współpracuje w tym obszarze, jak i ekspertów analityków.</p> <p>K2. Potrafi poddać krytycznej ocenie pozyskane dane.</p> <p>K3. Potrafi czerpać wiedzę z danych i na tej podstawie formułować propozycje rozwiązania sytuacji problemowych.</p>

#### Metody i techniki kształcenia stosowane przy realizacji kursu

Dydaktyczne	<ul style="list-style-type: none"> <li>- Metody odnoszące się do autentycznych lub fikcyjnych sytuacji (zbiory danych rzeczywistych i przykłady ich analiz),</li> <li>- metody wymiany i dyskusji (fora dyskusyjne),</li> <li>- metody ewaluacyjne (testy i zadania podlegające ocenie).</li> </ul>
Dydaktyczne eksponujące	Pokaz (filmy instruktażowe demonstrujące działanie programu).
Dydaktyczne podające	Wykład informacyjny (nagrania krótkich wykładów dotyczących teorii).
Dydaktyczne poszukujące	<ul style="list-style-type: none"> <li>- Ćwiczeniowa (testy i zadania zlecone przez prowadzących),</li> <li>- studium przypadku.</li> </ul>

#### Program przedmiotu

Szczegółowy program przedmiotu	<p>W czasie kursu uczestnik zapozna się z programem PS IMAGO PRO oraz jego możliwościami. Realizowane będą następujące zagadnienia:</p> <ol style="list-style-type: none"> <li>1. <b>Otwarte źródła informacji:</b> big data i data mining, obowiązek ustawowy eksploracji danych, zadania eksploracji danych, jawne źródła danych (2 godziny) - celem zajęć jest przedstawienie założeń i możliwości eksploracji danych oraz uświadomienie słuchaczy co do dostępności danych m.in. w Internecie.</li> <li>2. <b>Podstawy pracy w PS IMAGO PRO:</b> tabele danych, poziom pomiaru zmiennych, specyfikacja właściwości zmiennych w PS IMAGO PRO (2 godziny) - celem zajęć jest przedstawienie podstawowych możliwości okna danych programu PS IMAGO PRO oraz zwrócenie uwagi na konieczność sprawdzenia jakości danych i specyfikacji zmiennych.</li> <li>3. <b>Statystyka opisowa:</b> miary tendencji centralnej, miary rozproszenia, wykres skrzynka z wąsami, histogramy, standaryzacja i normalizacja zmiennych, zasada trzech sigma, rozkład normalny (2 godziny) - celem zajęć jest</li> </ol>
--------------------------------	---

	<p>przypomnienie podstawowych statystyk i wykresów używanych w analizie statystycznej zmiennych jakościowych i ilościowych.</p> <p>4. <b>Eksploracyjna analiza danych:</b> opis danych jedno- i wielowymiarowych, współczynniki korelacji Pearsona i Spearmana, wizualizacja danych (2 godziny) - celem zajęć jest omówienie najważniejszych typów wizualizacji danych oraz przedstawienie możliwości programu PS IMAGO PRO w zakresie tworzenia wizualizacji danych.</p> <p>5. <b>Klasyfikacja metodą k najbliższych sąsiadów:</b> przygotowanie zmiennych do analizy, zasada działania algorytmu, algorytm k najbliższych sąsiadów w PS IMAGO PRO (2 godziny) - celem zajęć jest omówienie algorytmu k najbliższych sąsiadów oraz jego implementacji w programie PS IMAGO PRO.</p> <p>6. <b>Ocena jakości modelu:</b> miary jakości klasyfikacji i szacowania (2 godziny) - celem zajęć jest omówienie procesu oceny modelu, w tym podziału na zbiory uczący i testowy, oraz najważniejszych miar jakości klasyfikacji (trafność, czułość, swoistość) i szacowania (średni błąd bezwzględny i pierwiastek z błędu średniokwadratowego).</p> <p>7. <b>Drzewa klasyfikacyjno-regresyjne CRT:</b> budowa i interpretacja drzew klasyfikacyjno-regresyjnych, zapobieganie przeuczeniu, poprawa jakości drzew, budowa drzew z wykorzystaniem programu PS IMAGO PRO, predykcja z użyciem zbudowanych modeli (2 godziny) - celem zajęć jest omówienie modelu drzew klasyfikacyjno-regresyjnych oraz wykształcenie praktycznej umiejętności tworzenia takich modeli w programie PS IMAGO PRO i stosowania ich do klasyfikacji i szacowania.</p> <p>8. <b>Sieci neuronowe:</b> budowa i uczenie modelu perceptronu wielowarstwowego (2 godziny) - celem zajęć jest zapoznanie słuchaczy z tematyką sieci neuronowych na przykładzie perceptronu wielowarstwowego, omówienie budowy tego typu sieci oraz tworzenie gotowych modeli w PS IMAGO PRO.</p> <p>9. <b>Prosta regresja liniowa:</b> model prostej regresji liniowej, interpretacja jego parametrów oraz ocena jakości (2 godziny) - celem zajęć jest wprowadzenie modeli liniowych na przykładzie prostej regresji liniowej.</p> <p>10. <b>Wieloraka regresja liniowa:</b> model regresji wielorakiej, ocena jakości modelu oraz metody wyboru zmiennych (2 godziny) - celem zajęć jest rozszerzenie modelu prostej regresji liniowej na przypadek z wieloma predyktorami oraz omówienie dostępnych w programie PS IMAGO PRO metod wyboru zmiennych do modelu.</p> <p>11. <b>Grupowanie metodą k-średnich:</b> cel i zastosowania grupowania, algorytm k średnich, jego własności, wady i zalety (2 godziny) - celem zajęć jest zapoznanie słuchaczy z algorytmem k średnich, możliwością wykorzystania programu PS IMAGO PRO do wyznaczenia grup z wykorzystaniem tego algorytmu oraz wykształcenie umiejętności interpretacji uzyskanych klastrów.</p> <p>12. <b>Hierarchiczna analiza skupień:</b> miary</p>
--	--

	<p>podobieństwa/niepodobieństwa dla zmiennych ilościowych i binarnych, metody określania odległości grup, tworzenie i interpretacja dendrogramów (2 godziny) - celem zajęć jest wprowadzenie metody grupowania hierarchicznego oraz wykształcenie umiejętności odczytywania dendrogramów.</p> <p>13. <b>Dwustopniowa analiza skupień:</b> opis algorytmu dwustopniowej analizy skupień oraz interpretacja otrzymanych grup (2 godziny) - celem zajęć jest wykształcenie umiejętności znajdowania grup z wykorzystaniem implementacji algorytmu dwustopniowej analizy skupień w programie PS IMAGO PRO, a także interpretacji i ewaluacji otrzymanych klastrów.</p> <p>14. <b>Analiza koszykowa:</b> typy i miary jakości reguł asocjacyjnych, algorytm A priori (2 godziny) - celem zajęć jest wprowadzenie do zagadnienia budowy reguł asocjacyjnych na przykładach z zakresu analizy koszykowej, demonstracja działania algorytmu A priori oraz ocena jakości reguł uzyskanych z wykorzystaniem tego algorytmu.</p> <p>15. <b>Studium przypadku</b> (2 godziny) - celem zajęć jest utrwalenie zdobytej przez słuchaczy wiedzy poprzez wykonanie w programie PS IMAGO PRO zadania wymagającego wykorzystania poznanych wcześniej algorytmów.</p> <p>Każdy temat będzie zawierał krótki wstęp teoretyczny, materiał instruktażowy, jak rozwiązać dany problem za pomocą PS IMAGO PRO, przykładowe zbiory danych, ćwiczenia do wykonania i zadania do rozwiązania.</p>
Literatura tradycyjna	<ol style="list-style-type: none"> <li>1. Bedyńska S., Książek M. (red.): <i>Statystyczny drogowkaz 1. Praktyczne wprowadzenie do wnioskowania statystycznego</i>. Warszawa, Wydawnictwo Akademickie Sedno, 2012.</li> <li>2. Górniak J., Wachnicki J.: <i>Pierwsze kroki w analizie danych</i>. Kraków, SPSS Polska, 2004.</li> <li>3. Daniel T. Larose: <i>Odkrywanie wiedzy z danych</i>, Wydawnictwo Naukowe PWN, Warszawa, 2013.</li> <li>4. Daniel T. Larose: <i>Metody i modele eksploracji danych</i>, Wydawnictwo Naukowe PWN, Warszawa, 2012.</li> <li>5. Tadeusz Morzy: <i>Eksploracja danych. Metody i algorytmy</i>. Wydawnictwo Naukowe PWN, Warszawa, 2013.</li> <li>6. Bartosz Saramak: <i>Wykorzystanie otwartych źródeł informacji w działalności wywiadowczej: historia, praktyka, perspektywy</i>. WDiNP UW, Warszawa, 2015 (dostępne on-line: <a href="https://wnpism.uw.edu.pl/wp-content/uploads/2019/08/Wykorzystanie-otwartych-zrodel.pdf">https://wnpism.uw.edu.pl/wp-content/uploads/2019/08/Wykorzystanie-otwartych-zrodel.pdf</a>, dostęp z dn. 22.10.2020).</li> <li>7. Sobczyk M.: <i>Statystyka, Wydanie piąte uzupełnione</i>. Warszawa, Wydawnictwo Naukowe PWN, 2020, str. 12-16.</li> <li>8. Stephane Tuffery: <i>Data Mining and Statistics for Decision Making</i>. Wiley, 2011.</li> <li>9. Xindong Wu, Vipin Kumar: <i>The Top Ten Algorithms in Data Mining</i>. Chapman &amp; Holl/CRC, 2009.</li> </ol>

Wykorzystywane e-materiały	<ul style="list-style-type: none"> <li>- Biecek P.: <i>Odkrywać! Ujawniać! Objaśniać! Zbiór esejów o sztuce prezentowania danych.</i> <a href="http://www.biecek.pl/Eseje/">http://www.biecek.pl/Eseje/</a> (dostęp z dn. 23.10.2020).</li> <li>- <i>Polska Statystyka Publiczna</i> <a href="https://bip.stat.gov.pl/dzialalnosc-statystyki-publicznej/polska-statystyka-publiczna/">https://bip.stat.gov.pl/dzialalnosc-statystyki-publicznej/polska-statystyka-publiczna/</a> (dostęp z dn. 21.10.2020).</li> <li>- Strona internetowa <i>Discovering Statistics</i> z tutorialami, filmami i zbiorami danych <a href="https://www.discoveringstatistics.com/">https://www.discoveringstatistics.com/</a> (dostęp z dn. 24.08.2020).</li> <li>- Strona internetowa IBM Knowledge Center z dokumentacją produktu IBM SPSS Statistics <a href="https://www.ibm.com/support/knowledgecenter/pl/SSLVMB26.0.0/statistics_kc_ddita/spss/product_landing.html">https://www.ibm.com/support/knowledgecenter/pl/SSLVMB26.0.0/statistics_kc_ddita/spss/product_landing.html</a> (dostęp z dn. 25.08.2020).</li> <li>- Predictive Solutions: <i>Samouczek.</i> <a href="http://samouczek.predictivesolutions.pl/">http://samouczek.predictivesolutions.pl/</a> (dostęp z dn. 12.09.2020).</li> <li>- Predictive Solutions: <i>Blog.</i> <a href="https://predictivesolutions.pl/2,blog">https://predictivesolutions.pl/2,blog</a> (dostęp z dn. 29.10.2020).</li> <li>- StatSoft: <i>Internetowy Podręcznik Statystyki.</i> <a href="https://www.statsoft.pl/textbook/stathome.html">https://www.statsoft.pl/textbook/stathome.html</a> (dostęp z dn. 25.10.2020).</li> <li>- Strona internetowa <i>Towards Data Science:</i> <a href="https://towardsdatascience.com/">https://towardsdatascience.com/</a> (dostęp z dn. 30.10.2020).</li> <li>- Portal <i>Sztuczna Inteligencja</i> <a href="https://www.sztucznainteligenca.org.pl/">https://www.sztucznainteligenca.org.pl/</a> (dostęp z dn. 29.10.2020).</li> <li>- Narzędzia doboru kolorów na wykresach: <ul style="list-style-type: none"> <li>- <a href="http://paletton.com">http://paletton.com</a> (dostęp z dn. 23.10.2020),</li> <li>- <a href="https://colorbrewer2.org/">https://colorbrewer2.org/</a> (dostęp z dn. 23.10.2020).</li> </ul> </li> <li>- Przemysław Biecek: Przewodnik po pakiecie R 4.0 <a href="http://pbiecek.github.io/Przewodnik/">http://pbiecek.github.io/Przewodnik/</a> (dostęp z dn. 30.10.2020).</li> <li>- Strona domowa projektu R <a href="https://www.r-project.org/">https://www.r-project.org/</a> (dostęp z dn. 27.10.2020).</li> </ul>
----------------------------	--

#### Weryfikacja efektów uczenia się i ewaluacja

Kryteria oceniania	Ocenie podlegają aktywności zamieszczone w kolejnych modułach, w tym: testy zawierające pytania teoretyczne oraz pytania odnoszące się do ćwiczeń praktycznych, krótkie zadania do wykonania w programie oraz końcowe zadanie zaliczeniowe w formie studium przypadku.
Sposób przeprowadzania oceny	Ocena jest przeprowadzana zdalnie. Kursant otrzymuje punkty za wszystkie wymienione aktywności, przy czym przy każdej z nich podana jest informacja o maksymalnej możliwej do zdobycie liczbie punktów. Wszystkie aktywności powinny być realizowane zgodnie z harmonogramem kursu. Prowadzący zajęcia może obniżyć ocenę za aktywności wykonane po terminie.

	Do zaliczenia kursu wymagane jest uzyskanie co najmniej połowy punktów możliwych do zdobycia.
Ewaluacja zajęć	<ul style="list-style-type: none"><li>– Fora zamieszczone przy każdym module i służące do zadawania pytań pozwolą na bieżącą ewaluację zajęć oraz reagowanie w sytuacjach wymagających korekty materiałów.</li><li>– Po każdym module zostanie zamieszczona krótka ankieta dotycząca poziomu trudności zajęć oraz stopnia ich zrozumienia.</li><li>– Kurs będzie kończył się ankieta pozwalająca na wszechstronną ocenę zajęć: poziomu merytorycznego, jakości materiałów, stopnia zadowolenia słuchacza z kursu.</li></ul>