

Inżynieria Mechatroniczna

Podstawy Sztucznej Inteligencji i Ucznia Głębokiego:  
**3: Uczenie (maszynowe) na podstawie danych**

Ziemowit Dworakowski  
 AGH w Krakowie

1

---

---

---

---

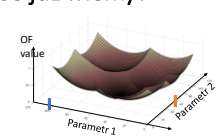
---

---

---

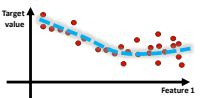
---

**Co już wiemy?** SD



OF value

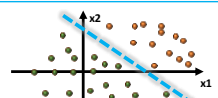
Jak dotąd, wiemy jak znaleźć współrzędne minimum globalnego funkcji celu



Target value

Feature 1

Wiemy czym jest regresja i jak zaprojektować liniowe i nieliniowe modele regresyjne



x2

x1

Wiemy czym jest klasyfikacja i potrafimy zaproponować liniowy model klasyfikacyjny

*Dzisiaj nieco uogólnimy, spojrzymy na te problemy od strony przestrzeni cech i poznamy kilka nowych narzędzi i pojęć które przydadzą się po drodze*

2

---

---

---

---

---

---

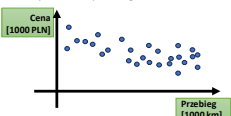
---

---

**Przestrzeń parametrów a przestrzeń cech** SD

Jeśli jakąś wartość mierzymy (lub z góry otrzymujemy) – jest to **cecha** (ang. **feature**). Jeśli mamy nad nią kontrolę i aktywnie poszukujemy jej konkretnej wartości - jest to **parametr** (ang. **parameter**).

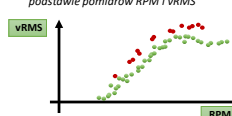
**Regresja:**  
 „Podaj spodziewaną cenę samochodu na podstawie przebiegu”



Cena [10000 PLN]

Przebieg [10000 km]

**Klasyfikacja:**  
 „Podaj spodziewany stan turbiny wiatrowej na podstawie pomiarów RPM i vRMS”



vRMS

RPM

3

---

---

---

---

---

---

---

---

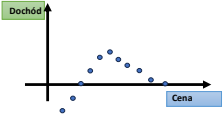
### Przestrzeń parametrów a przestrzeń cech SD

Jeśli jakąś wartość **mierzmy** (lub z góry otrzymujemy) – jest to **cecha** (ang. **feature**). Jeśli mamy nad nią kontrolę i aktywnie poszukujemy jej konkretnej wartości – jest to **parametr** (ang. **parameter**).

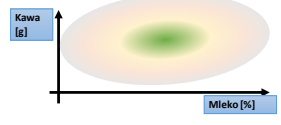
Dana zmienna może być parametrem lub cechą – w zależności od kontekstu

Optymalizację wykonujemy w przestrzeni parametrów. Regresję i klasyfikację w przestrzeni cech

**Optymalizacja:**  
Znajdź taką cenę za samochód, by zmaksymalizować dochód.



**Optymalizacja:**  
„Podaj dla jakich proporcji kawy i mleka otrzymamy najsmaczniejsze latte”



4

---

---

---

---

---

---

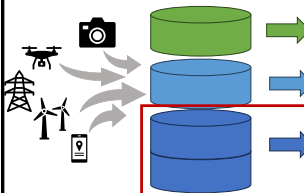
---

---

---

---

### Organizacja danych SD



**Te podzbiory muszą być rozłączne!**

- Na podstawie podzbioru **testowego** oceniamy spodziewaną docelową skuteczność metod
- Na podstawie podzbioru **walidacyjnego** konfigurujemy modele i sprawdzamy czy wzorce są ogólne
- Na podstawie podzbioru **treningowego** uczymy się rozpoznawania wzorców w danych

**Dzisiaj skupimy się wyłącznie na tym początkowym kroku, finalną konfigurację i weryfikację zostawiając na kolejne okazje...**

5

---

---

---

---

---

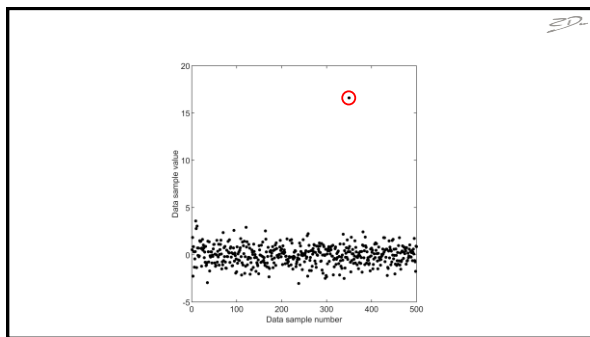
---

---

---

---

---



6

---

---

---

---

---

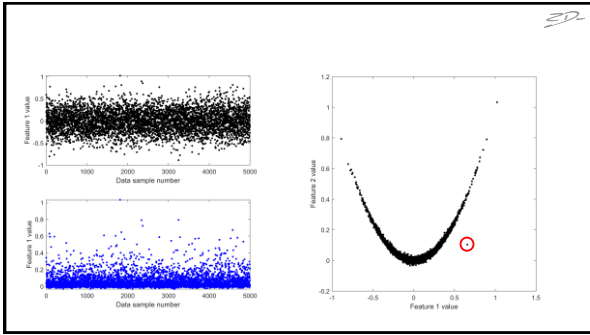
---

---

---

---

---



7

---

---

---

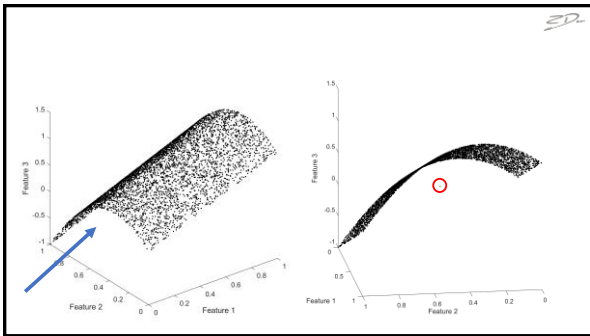
---

---

---

---

---



8

---

---

---

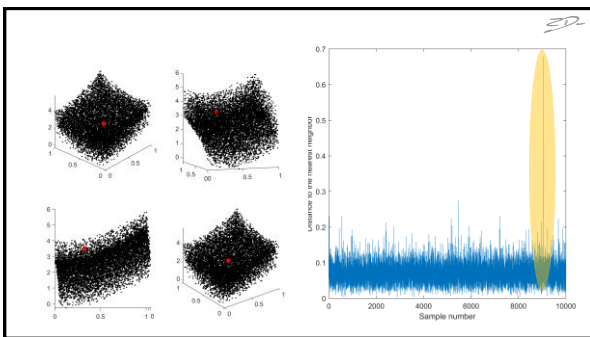
---

---

---

---

---



9

---

---

---

---

---

---

---

---

## Wizualizacje danych mogą być mylące...

Więc – chcemy zaprojektować takie metody reprezentowania danych, by być w stanie automatycznie **rozpoznawać wzorce** w wielowymiarowej przestrzeni cech



- Zbudować reguły określające gdzie można się spodziewać obecności danych
- Zorganizować dane poprzez rozpoznanie zależności między nimi i ew. przynależności do podkategorii

10

---

---

---

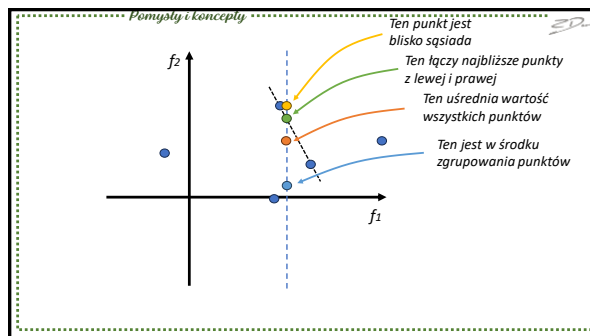
---

---

---

---

---



11

---

---

---

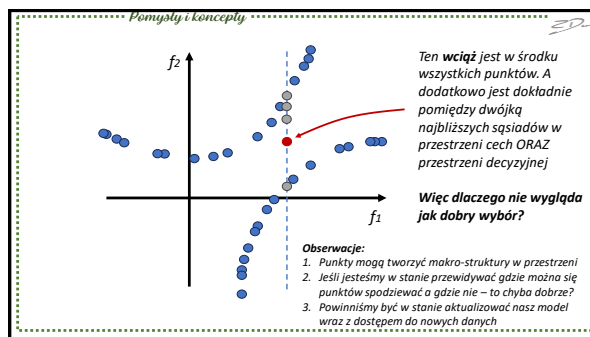
---

---

---

---

---



12

---

---

---

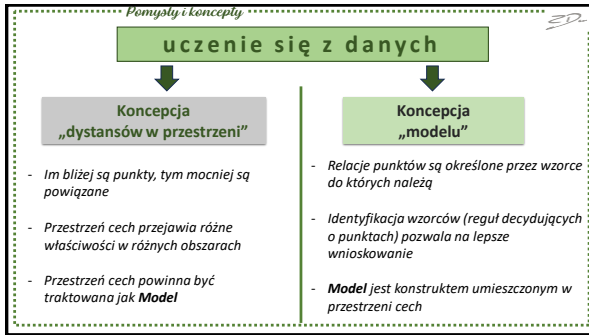
---

---

---

---

---



13

---

---

---

---

---

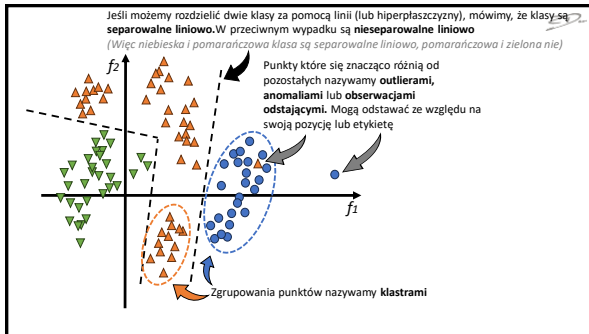
---

---

---

---

---



14

---

---

---

---

---

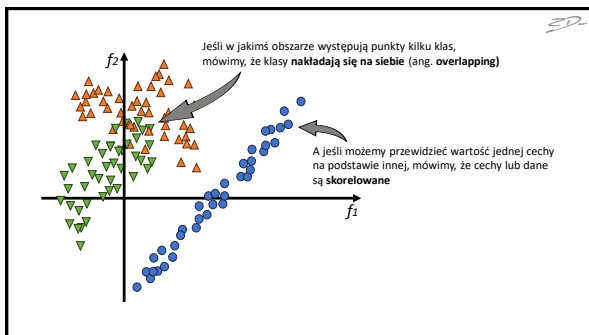
---

---

---

---

---



15

---

---

---

---

---

---

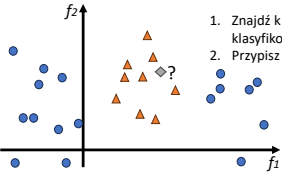
---

---

---

---

**Metoda k-najbliższych sąsiadów (k Nearest Neighbor classifier, kNN)** SD



1. Znajdź k (nieparzyste) najbliższych sąsiadów klasyfikowanego punktu
2. Przypisz punktowi najpopularniejszą etykietę

Jak metoda zareaguje na anomalie?

↓

k powinno być  $\geq 3$

+ Intuicyjna  
+ Prosta w implementacji  
+ Prosta w konfiguracji

- Wymaga dużo pamięci  
- Wymaga dużo mocy obliczeniowej (powolna)  
- Złe skaluje się w wysokowymiarową przestrzeń cech

16

---

---

---

---

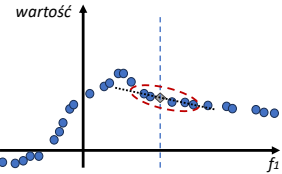
---

---

---

---

**Regresja z zastosowaniem metody kNN** SD



... znamy już coś bardzo podobnego!  
(Prawie dokładnie tak samo działa estymator lokalnie liniowy)

17

---

---

---

---

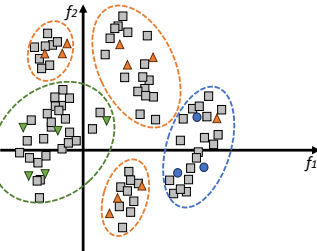
---

---

---

---

**Rozwiązania bazujące na klasteryzacji** SD



1. Zaczynamy od nieoetykietowanych danych (to metoda **nienadzorowana**)
2. Znajdźmy klastry danych
3. „Pokolorujmy” klastry na podstawie wybranych oetykietowanych punktów (np. wybierając większośćową etykietę)

18

---

---

---

---

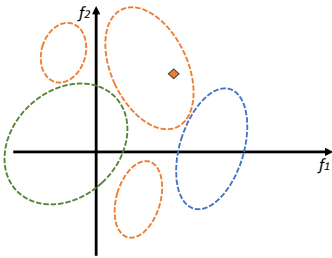
---

---

---

---

## Rozwiązania bazujące na klasteryzacji



1. Zaczynamy od nieoetykietowanych danych (to metoda **nienadzorowana**)
2. **Znajdźmy klastry danych**
3. „Pokolorujemy” klastry na podstawie wybranych oetykietowanych punktów (np. wybierając **większościową etykietę**)
4. Dla nowych danych sprawdzimy do których klastrów należą punkty i klasyfikujemy je odpowiednio

19

---

---

---

---

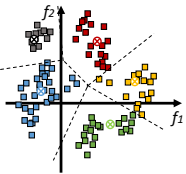
---

---

---

---

## Algorytm centroidów (K-means clustering, KMC)



- + Intuicyjny
- + Jedna z najprostszych metod klasteryzacji
- **Musimy wiedzieć, ile klastrów jest w danych**
- **Granice klastrów nie mają dobrego uzasadnienia**
- **Metoda źle skaluje się w przestrzeń wielowymiarową**

- Założmy ilość klastrów, wybierzmy ich centroidy losowo
- ↓
- Przypiszmy punkty do klastrów bazując na ich dystansie do centroidu
- ↓
- Przeliczmy położenia centroidów na podstawie przypisanych do nich punktów
- ↓
- Powtarzajmy, dopóki położenia centroidów się zmieniają

20

---

---

---

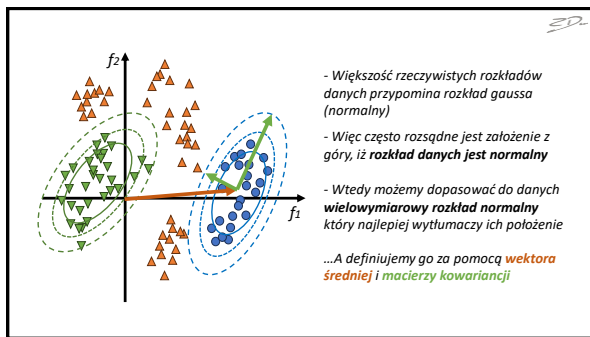
---

---

---

---

---



- Większość rzeczywistych rozkładów danych przypomina rozkład gaussa (normalny)
- Więc często rozsądne jest założenie z góry, iż **rozkład danych jest normalny**
- Wtedy możemy dopasować do danych **wielowymiarowy rozkład normalny** który najlepiej wytłumaczy ich położenie
- ...A definiujemy go za pomocą **wektora średniej** i **macierzy kowariancji**

21

---

---

---

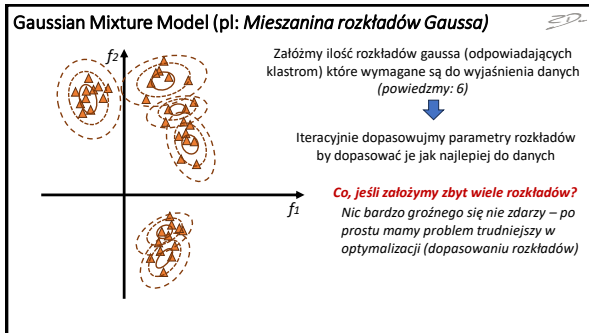
---

---

---

---

---



22

---

---

---

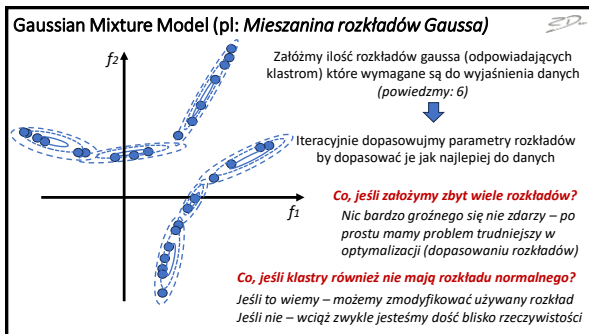
---

---

---

---

---



23

---

---

---

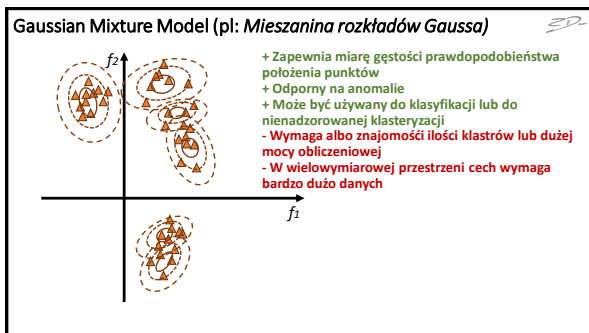
---

---

---

---

---



24

---

---

---

---

---

---

---

---

## Pytania do powtórki:

SD

1. Przedstaw jak dane eksperymentalne dzielone są na podzbiory, nazwij te podzbiory i wyjaśnij do czego są używane
2. Wyjaśnij separowalność liniową i jej brak, anomalie, klastry i korelację danych. Zaprezentuj graficzne ilustracje w.w. pojęć
3. Wyjaśnij jak działa metoda kNN w klasyfikacji (podaj też wady i zalety)
4. Wyjaśnij na czym polega klasteryzacja danych, wyjaśnij metodę centroidów z jej wadami i zaletami
5. Wyjaśnij jak działa metoda mieszaniny rozkładów Gaussa – z graficznym przykładem Gaussianów dopasowanych do klastrów danych. Przedstaw wady i zalety tej metody

---

---

---

---

---

---

---

---