

Algorytmiczna Analiza Danych

Wykład 2

Uczenie nadzorowane

2025-10-09

1. Zbiór danych z etykietami

$$D = \{(\bar{x}^{(1)}, y^{(1)}), \dots, (\bar{x}^{(n)}, y^{(n)})\} \quad (1)$$

$\bar{x}^{(i)} = (x_1^{(i)}, \dots, x_p^{(i)})$ - cechy, zmienne wyjaśniające, predyktory
 $y^{(i)}$ - etykieta, zmienna wyjaśniająca

Problem: wyznaczyć ukrytą relację (ang. underlying relation) f

$$y^{(i)} = f(\bar{x}^{(i)}) + \varepsilon^{(i)} \quad (2)$$

$\varepsilon^{(i)}$ - szum

1.1. Etykieta

Regresja - zmienna ciągła

Klasyfikacja - zmienna dyskretna

2. Statistical learning vs Machine Learning

Statistical learning	Machine learning
<ul style="list-style-type: none">• proste modele z mocnym wsparciem teoretycznym, np. regresja liniowa,• łatwa interpretacja modelu, modele parametryczne	<ul style="list-style-type: none">• złożone modele, weryfikowane empirycznie, np. <i>KNN</i>, <i>DT</i>• zazwyczaj nieparametryczne lub z dużą ilością parametrów
<ul style="list-style-type: none">• łatwe wnioskowanie,• nie wymaga mocy obliczeniowej	<ul style="list-style-type: none">• duża moc predykcyjna• nie wymaga założeń o danych
<ul style="list-style-type: none">• słaba moc predykcyjna,• wymaga silnych założeń o danych,• wymaga aby szum (ε) miał rozkład normalny	<ul style="list-style-type: none">• trudne w interpretacji (black box),• wymaga dużej mocy obliczeniowej

\leadsto wnioskowanie (ang. inference)

\leadsto predykcja, przewidywanie

3. Przygotowanie danych

1. Brakujące dane

- usuwanie obserwacji,
 - zastępowanie:
 - średnią
 - medianą
 - dominantą dla danych dyskretnych
 - losowanie zgodne z rozkładem
 - dodatkowa etykieta „unknown”
2. Usuwanie podejrzanych obserwacji
 - outlier (nietypowy),
 - high leverage points,
 - histogram etykiet
 3. Inżynieria cech
 - dostosować dane, tak by predykcja była łatwiejsza
 4. Podział danych
 - część treningowa $\approx 75\%$,
 - część walidacyjna $\approx 15\%$,
 - część testowa $\approx 10\%$,
 - typowo podział jest losowy

4. Błąd średnio kwadratowy (ang. MSE - Mean Square Error)

$$\begin{aligned}\text{MSE} &= \mathbb{E}\left(\left(\hat{f}(\bar{x}^{(0)}) - y^{(0)}\right)^2\right) \\ &= \text{Bias}\left(\hat{f}(\bar{x}^{(0)})\right)^2 + \text{Var}\left(\hat{f}^2(\bar{x}^{(0)})\right) + \text{Var}(\varepsilon)\end{aligned}\tag{3}$$

$$\text{Bias}(\hat{f}(\bar{x}^{(0)})) = \mathbb{E}\left(f(\bar{x}^{(0)}) - \hat{f}(\bar{x}^{(0)})\right)\tag{4}$$