

Metody eksploracji danych

Laboratorium 1

Weka + Python + regresja

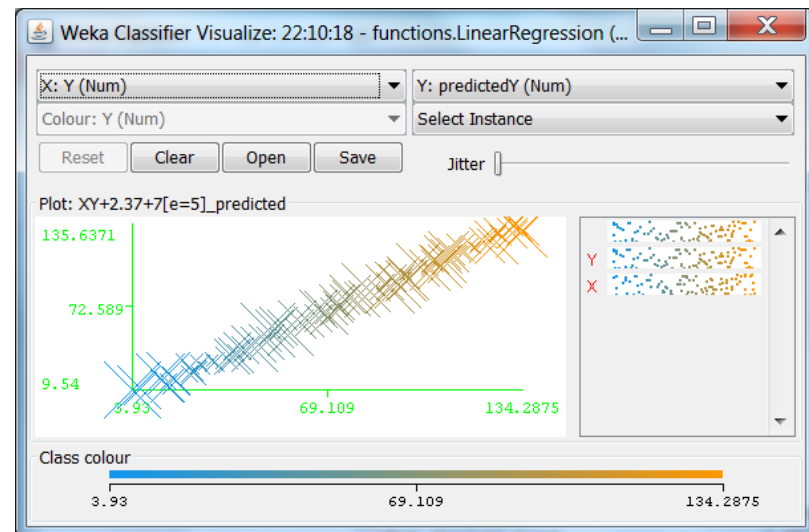
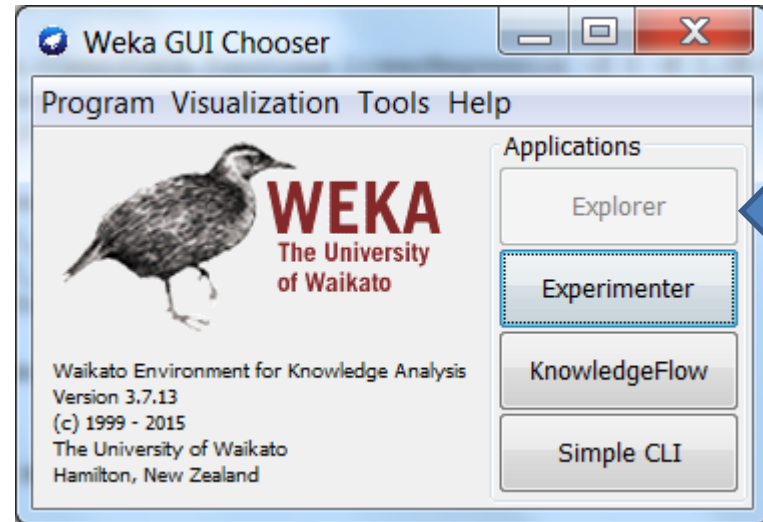
Metody eksploracji danych

- Zasoby
 - Weka (gdzieś na dysku)
 - Środowisko dla języka Python (Spyder, Jupyter, gdzieś na dysku)
 - Zbiory danych (dostępne na stronie http://home.agh.edu.pl/~pszwed/wiki/doku.php?id=metody_eksploracji_danych)
- Cel
 - Poznanie możliwości Weki
 - Wykonanie klika przykładów regresji
 - Alternatywne środowisko i biblioteki języka Python
 - Typowy workflow podczas analizy danych

Część 1 – Nawigacja po GUI Weka + nieco bibliotek Pythona

1.1 W oknie startowym Weka uruchom opcję Explorer

- Załaduj plik **xy-001.arff** w zakładce **Preprocess**
- Możesz wyświetlić dane w tabeli **Edit**
- Wyświetl zawartość danych w zakładce **Visualize**
- W zakładce **Classify** wybierz **Choose:** functions -> LinearRegression
- Naciśnij Start. Powinno się pojawić równanie prostej
- Niestety za pomocą Weka nie można wyświetlić prostej, co najwyżej zakres błędów klasyfikacji.
- Kliknij prawym klawiszem na result list i wybierz jedną z opcji wizualizacji.



1.2. Wyświetlamy dane i krzywą

Korzystamy z bibliotek w języku Python

1. Używamy IDE Spyder lub Jupyter. Można też przygotować plik tekstowy i uruchomić z poziomu konsoli (mając nadzieję, że biblioteki są zainstalowane). Czasem też działa <https://try.jupyter.org/>
2. Kopiujemy fragment kodu ze strony: http://home.agh.edu.pl/~pszwed/wiki/doku.php?id=metody_eksploracji_danych
3. Wpisujemy dane w osobnych liniach. `"""` to specjalny zapis dla tekstu wprowadzanego w wielu wierszach.
4. `x, y = np.loadtxt(inp, delimiter=',', usecols=(0, 1), unpack=True, skiprows=6)` – ładuje tekst z kolumn 0 i 1 do tablic x,y
5. `plt.scatter(x,y,s=80, marker='+')` – rysuje punkty
6. `fx=np.linspace(-10,60,100)` 100 równomiernie rozłożonych punktów w zakresie [-10,60]
7. `fy=2.3702*fx+6.1973` - wartości funkcji regresji (to są operacje na wektorach)
8. `ftrue=2.37*fx+7` – „prawdziwa” funkcja użyta do generacji danych
9. `plt.plot(fx,fy,linewidth=2,color='r')` – wykres liniowy
10. `plt.xlim(-10,60)` – zakres wyświetlanych zmiennych

Oczekiwany wynik

- Wykres może być osadzony na stronie HTML w Jupyter: usuń komentarze z pierwszej linii: `%matplotlib notebook`
- Wykres można zapisać w jednym z formatów (bitmapowych/wektorowych)

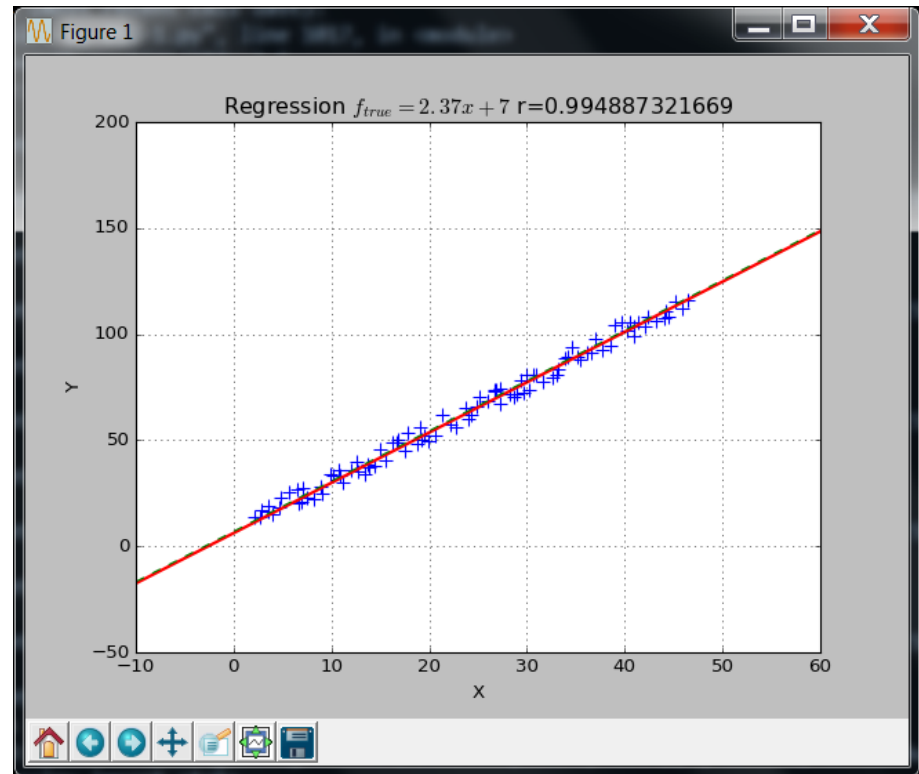
Wartość r która jest wyświetlana dodatkowo?

Jest to współczynnik korelacji.

Przeczytaj tekst i obejrzyj przykłady:

https://en.wikipedia.org/wiki/Correlation_and_dependence

Jeśli $r \approx 0$, zazwyczaj równanie regresji będzie miało postać:
 $y = 0x + c$



Inne pliki

Analogicznie załaduj pliki, wyznacz równanie regresji, oceń błędy i wyświetl przebieg krzywej regresji dla:

1.3 Pliku xy-002.arff

1.4 Pliku xy-003.arff

1.5 Pliku xy-005.arff

1.6 Pliku xy-006.arff

1.7 Pliku xy-007.arff

1.8 Pliku xy-008.arff [Skomentuj przebieg krzywej regresji]

1.8 Pliku xy-009.arff

1.9 Pliku xy-010.arff

Zajrzyj do 1.10. Może będzie bardziej efektywnie?

1.10 Realizacja za pomocą bibliotek Python

Najprostsza forma ładowania pliku arff

```
inp = "xy-001.arff",  
x, y = np.loadtxt(inp, delimiter=',', usecols=(0, 1), unpack=True, skiprows=6)
```

Liczba wierszy do przeskoczenia zależy od pliku

Tworzenie macierzy cech

```
features=x.reshape(x.size,1)
```

Musi być to macierz dwuwymiarowa $m \times 1$ (w ogólnym przypadku $m \times n$)

Dopasowanie krzywej

```
regr = linear_model.LinearRegression()  
regr.fit(features, y)
```

Parametry modelu (regr.coef_ jest tablicą)

```
print('Coefficients: ', regr.coef_, ' Intercept: ',regr.intercept_)
```

...

```
fy= regr.coef_[0]* fx + regr.intercept_
```

Metody regresji

Metoda najmniejszych kwadratów

```
regr = linear_model.LinearRegression()  
regr.fit(features, y)
```

Regresja grzbietowa: metoda gradientu prostego z regularyzacją współczynników, norma L^2

```
regr = linear_model.Ridge(alpha = .5)  
regr.fit(features, y)
```

Lasso: metoda gradientu prostego z regularyzacją współczynników, norma L^1

```
regr = linear_model.Lasso(alpha=0.1)  
regr.fit(features, y)
```