

# Zaawansowana wizualizacja danych

Jarosław Jasiewicz, dr hab. prof. UAM

[jarekj@amu.edu.pl](mailto:jarekj@amu.edu.pl)

## Tworzenie grafiki

Warstwa definicji

■ Opis co zrobić z danymi w jakimś w miarę zrozumiałym języku (np. python)

Warstwa przetwarzania

■ Tłumaczy język zrozumiały na obiekty graficzne (linie, poligony i ich zbiory)

Warstwa rysowania

■ Nakłada obiekty graficzne na jakieś tło

## Programowanie proceduralne

wszystkie parametry wykresu muszą być przekazane do procedury w jednym poleceniu

```
scatter(dane[0], dane[1], c="#aa1234", s=12, marker="o")
```

Po przekazaniu parametrów do procedury wykonywany jest wykres

## programowanie obiektowe

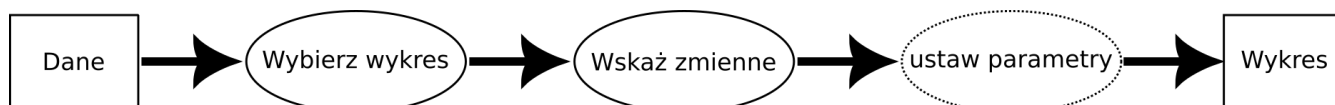
Tworzony jest obiekt a następnie modyfikowany jest jego stan wewnętrzny. Po ustaleniu parametrów, obiekt przekazywany jest do dalszego przetwarzania

```
fig, ax = plt.subplot() # tworzenie obiektów figura i axes  
ax.scatter(x[0], x[1]) # wywołanie metody scatter  
ax.draw() # wywołanie funkcji rysującej (niejawne)
```

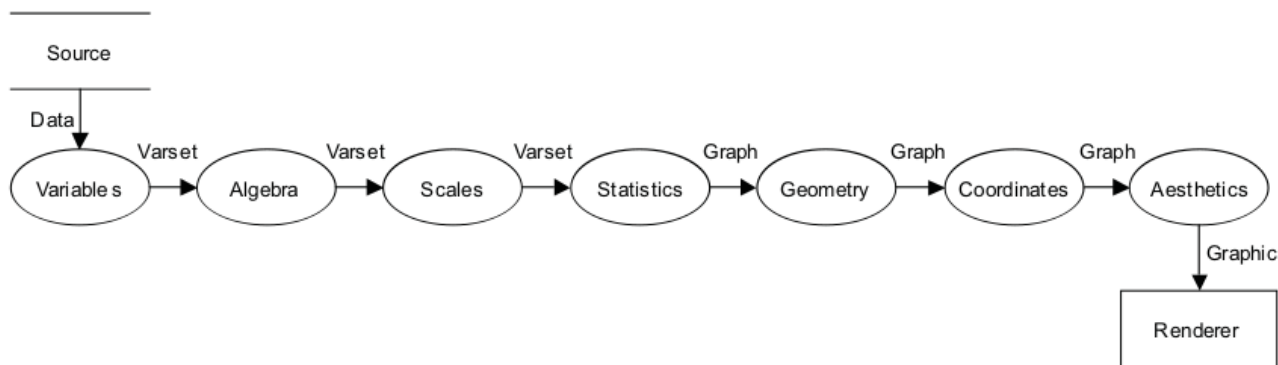
W przypadku ggplot jest to tylko imitowanie przepływu danych. R nie stosuje notacji *dot*.

## Definiowanie wykresu

Z punktu widzenia użytkownika



Z punktu widzenia programu



## Gramatyka grafiki

Niezależnie od implementacji istnieje konieczność przetłumaczenia opisu wykresu wykonanego w ramach języka programowania na język opisu geometrii, która zostanie użyta do narysowania wykresu na urządzeniu (np. PDF). Za tłumaczenie języka opisu danych na język opisu obiektów graficznych odpowiada gramatyka grafiki

## Leland Wilkinson *Grammar of Graphics*

1. **DATA** : zbiór operacji wyodrębniającej zmienne z danych
2. **TRANS** : transformacje zmiennych (grupowanie, sortowanie)
3. **SCALE** : transformacje skali (log, logit, power)
4. **COORD** : wybór układu odniesienia (polar, merkator, inne),
5. **ELEMENT** : geometria (punkt, linia) i estetyki (color),
6. **GUIDE** : elementy wykresu (axes, legends, etc.).

## Elementy graficzne wykresów

- punkty współrzędna X i Y
- łamane - uporządkowana lista współrzędnych x,y
- poligon
  - parametryczny: prostokąt, elipsa, linia
  - nieregularny
- macierz
- pole wektorowe

## Zmienne wizualne

Zmienne posługujące się kontekstem wizualnym do przekazywania treści. Zmienne mogą być planarne (położenie) lub retinalne (pozostałe)

- **pozycja**: zmiana położenia
- **rozmiar**: zmiana długości, powierzchni lub gęstości
- **kształt**: nieskończenie wiele, ale niw wszystkie rozróżnialne
- **jasność (value)**: gradient od najaśniejszego do najciemniejszego

- **kolor (walor)**
- **tekstura**: różnicowanie wzorców
- **orientacja**: różnicowanie wyrównania obiektów

## Rola elementów graficznych

### Punkty

Reprezentują położenie, nie mają długości ani obszaru. Rozmiar, kolor, znak, jasność są atrybutami niezależnymi od położenia i reprezentują niezależne zmienne

### Linie

Reprezentuje zjawiska o mierzalnej długości, ale nie ma powierzchni. Grubość linii podobnie jak powierzchnia punktu jest atrybutem.

### Powierzchnie

Reprezentuje zjawiska o mierzalnej powierzchni. Obszar może zmieniać położenie ale nie może zmieniać kształtu ani orientacji

## Zmienne wizualne a geometria

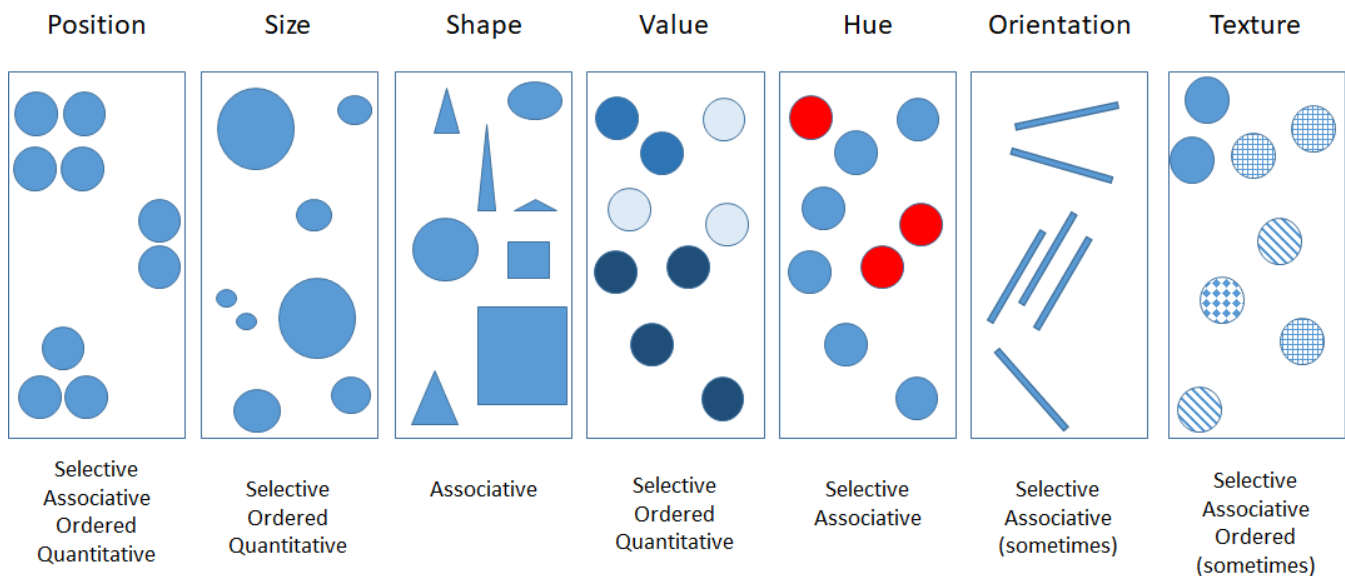
		symbols		
		point	line	area
differences in	size			
	value			
	grain			
	colour			
	orientation			
	shape			

## Rola zmiennych wizualnych

- **selekcja**: czy zmiana wartości oddziela od pozostałych obiektów

- **łączenie**: czy wartości zmiennych tworzą grupy i heirarchie
- **ilość**: czy forma zmiennej przekazuje informacje o ilości
- **porządek**: czy wartości zmiennej wywołują efekt porządkowania
- **skala zmian**: ile zmian jest wymaganych aby wartość zmiennej pokazała nieciągłość

## Rola poszczególnych zmiennych



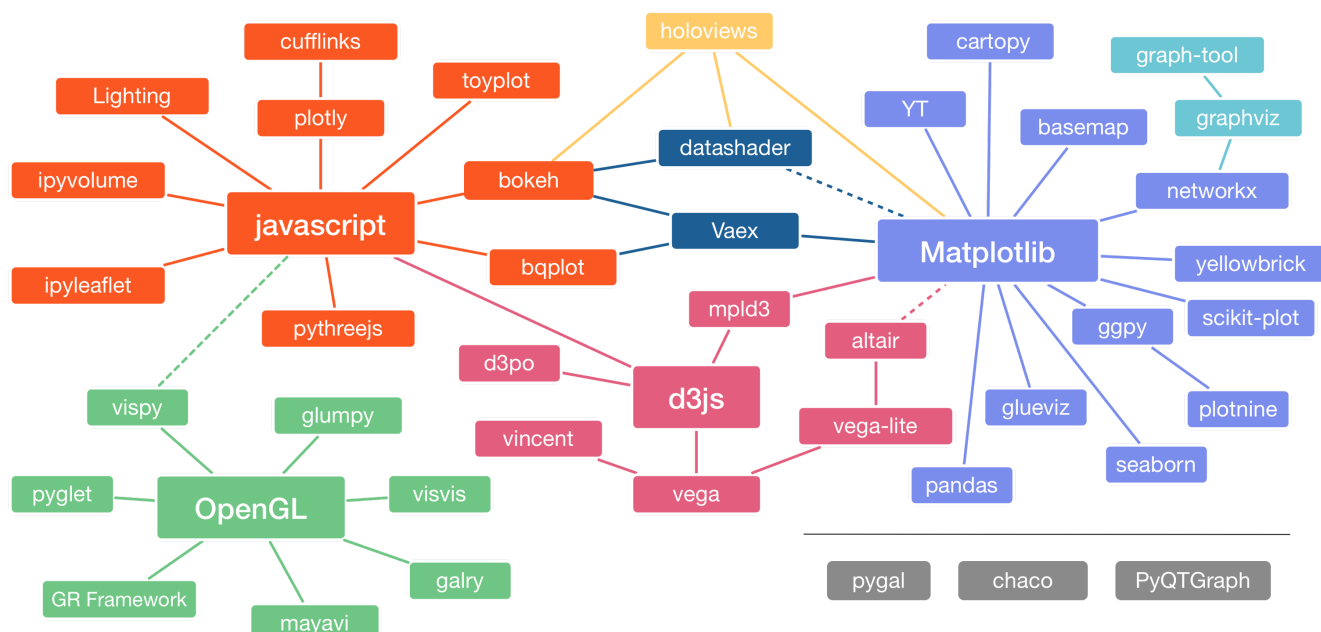
## Tworzenie wizualizacji

1. Wymaga zrozumienia relacji informacji i zdanań
2. Utworzenie rzutowania pomiędzy danymi a zmienną wizualną (jaka zmienna reprezentuje jaką informację)
3. Prezentacja wizualizacji
4. Weryfikacja użyteczności

## Weryfikacja użyteczności

1. Czy ma być drukowana bw/kolor?
2. Czy ma komunikować ogólną informację czy szczegóły
3. Z jakiej odległości ma być oglądana
4. Czy oświetlenie może wpływać na odbiór
5. Rodzaj odbiorcy

## System graficzny Pythona



- Zalety: dostępny wszędzie gdzie dostępne Numpy, liczne backendy, rozbudowany, wysoce konfigurowalny, liczne specjalistyczne nakładki
- Wady: trudny w obsłudze, nieprzystępny, wiele przestarzałych rozwiązań, nieprzystępna dokumentacja

- Zalety: interaktywność, duży wybór wykresów i możliwości konfiguracji, bardzo dobra dokumentacja
- Wady: ograniczony do rozwiązań przeglądarkowych, nie dostępny w rozwiązaniach osadzonych i innych backendach, wersja w Python to tylko wrapper.

- Zalety: Oparty o gramatykę wizualizacji i JSON
- Wady: takie same jak w przypadku **Plotly** + trudny w obsłudze i mało popularny

# Holoviews

---

Używa wyjścia matplotlib (wizualizacje statyczne) i Bokkeh/Plotly (interaktywne)

- Zalety: zalety obu systemów
- Wady: wady i ograniczenia obu systemów

## Oparty na SVG: PyGAL

---

Opisać

## Alternatywne API matplotlib

---

Matplotlib dostarcza przede wszystkim tworzywo do budowania zaawansowanych systemów. Alternatywne API generujące obiekty Matplotlib

Plotnine

■ Nakładka imitująca składnię *ggplot*.

Seaborn

■ Zaawansowany system tworzenia wykresów używający składni proceduralnej. (Będzie na zajęciach)

Seaborn.object

■ Eksperymentalny system wprowadzony w wersji 0.12 używający notacji funkcyjnej i gramatyki grafiki (rozwojowy)

## Wybrane dedykowane specjalistyczne nakładki

---

Pandas i XArray

■ Tabelaryczny przetwarzania danych. Posiada funkcję *plot()*, pozwalającą na tworzenie ograniczonych wykresów

Yellowbrick

■ System analizy i wizualizacji oparty na sci-kit learn. Przede wszystkim wykresy używane w data science

Cartopy

■ System przeznaczony do wizualizacji kartograficznych

---

NetworkX

■ System przeznaczony do analizy i wizualizacji sieci i grafów

SciPy

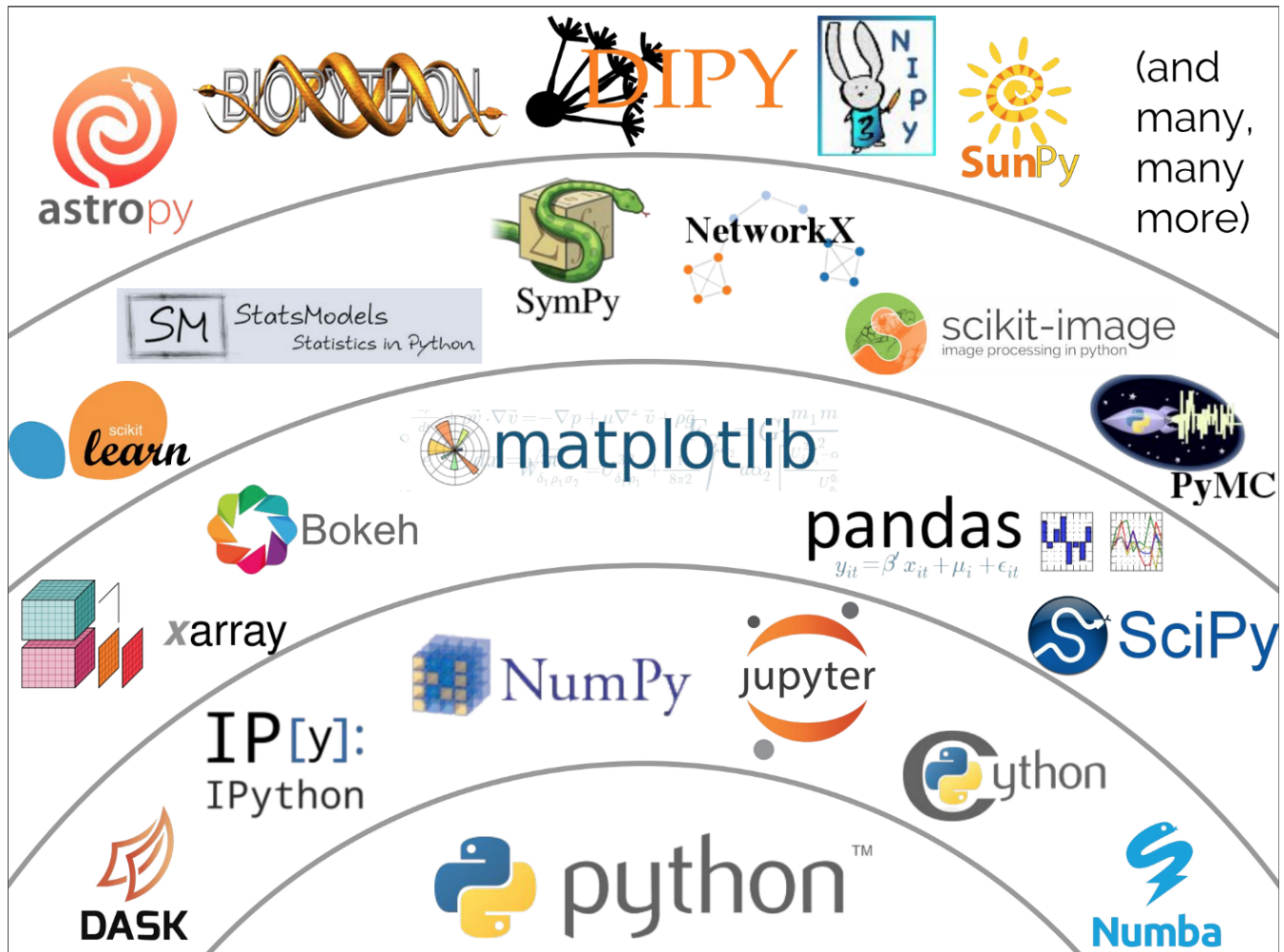
■ Dendrogramy

Statsmodel

Alternatywne wykresy podstawowe, wybrane złożone wykresy np. Treemap, mozaikowe dedykowane

Biblioteki dedykowane do pojedynczych wykresów, np *squarify*

## Ekosystem naukowy Pythona



## Warstwy przetwarzania danych w matplotlib

# **Scripting Layer**

*matplotlib.pyplot*

# **Artist Layer**

*matplotlib.artist*

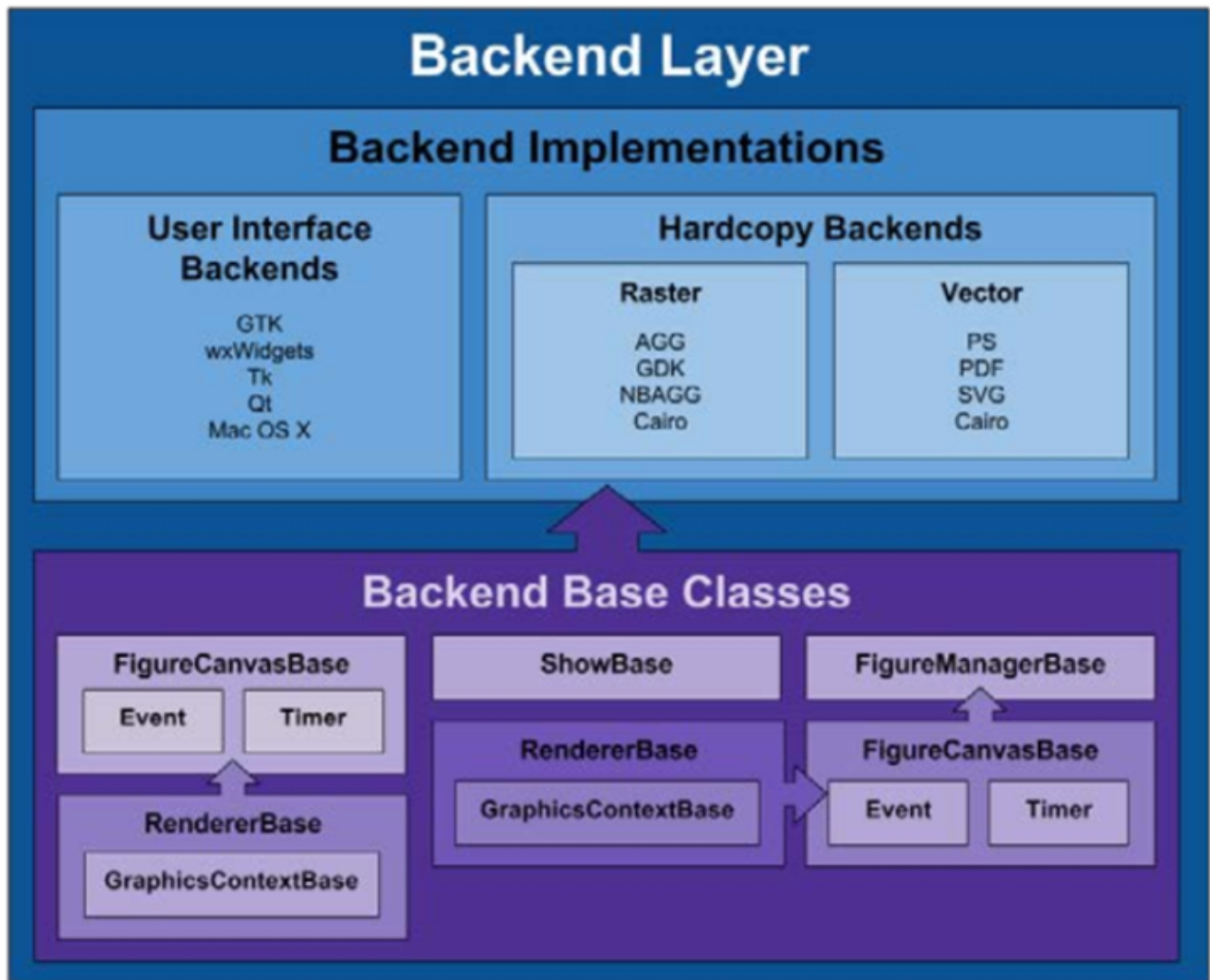
# **Backend Layer**

*matplotlib.backend\_bases*

## Backend Layer

---





FigureCanvas:

- Definiuje obszar rysowania (kartka)

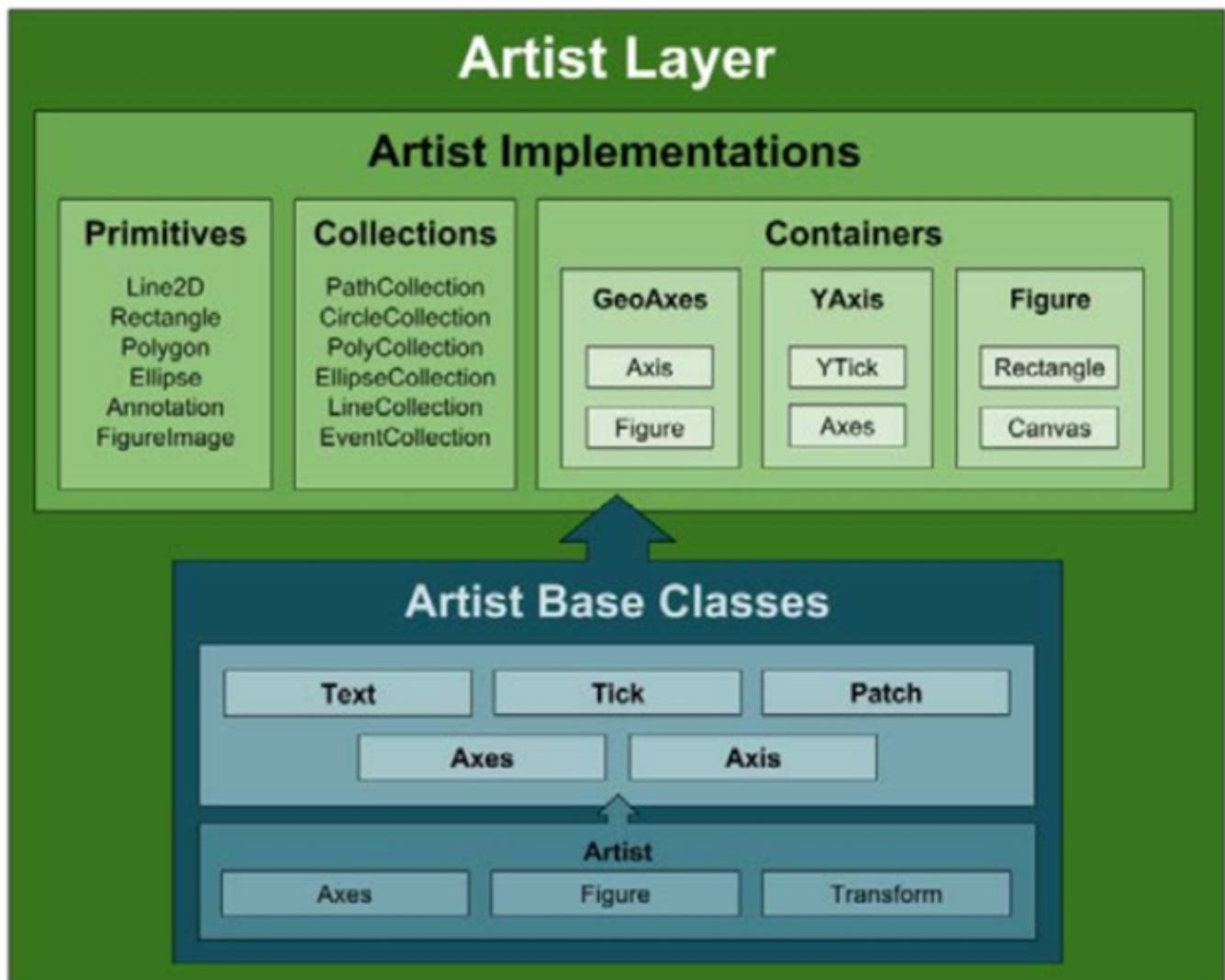
Renderer:

- Narzędzie rysujące (ołówek)

Event:

- Reakcja na działania użytkownika (kiedy rysować)

## Artist layer



Obsługuje zamianę opisu wykresów na rysowane elementy

Prymitywy

■ linie, prostokąty, elipsy, teksty obrazy

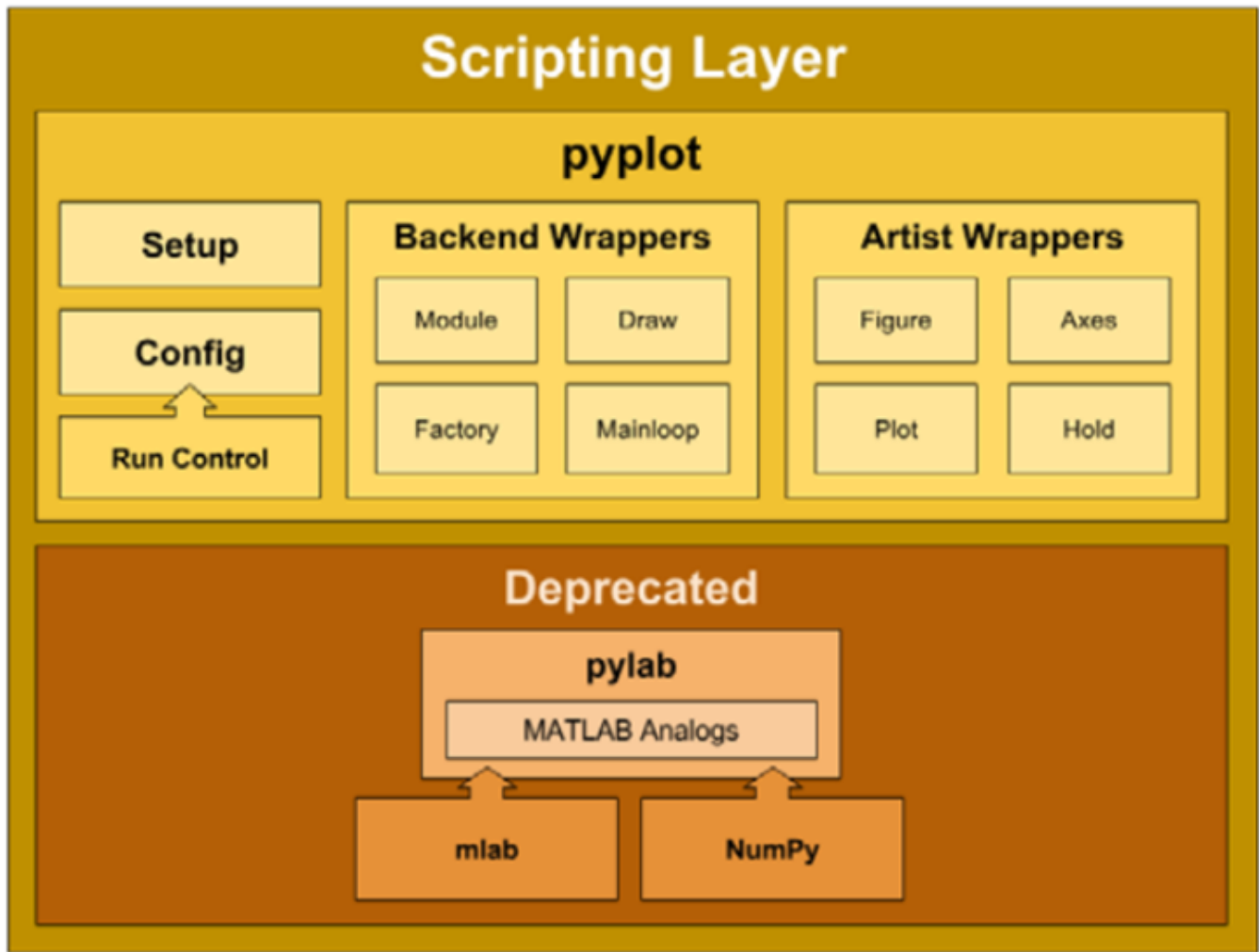
Kolekcje

■ złożone zarówno z prymitywów jak i innych kolekcji

Typy złożone

■ Figura, wykres (axes), osie (axis), znaczniki (ticks)

## Scripting layer



Język opisu figure

Manager figurey

▮ rozmiar, ilość i układ wykresów, itp

Przetwarzanie danych

▮ funkcje rysujące zwracające Artist

Wybór backendu

## Sposób definiowania wykresu:

```
fig, ax = plt.subplot() # manager figurey
ax.scatter(x[0],x[1]) # skrypt przetwarzania danych
fig.savefig("wykres.pdf") # wybór backendu (pdf)
```

## Implementacja

**matplotlib** jest zaimplementowany w ramach paradygmatu obiektowego (parogramowanie imperatywne), ale jego uproszczona wersja (pylab) imituje działanie proceduralne.

Każdy parametr w wykresie musi być określony. Jeżeli nie jest (z reguły *nie* jest) określony przez użytkownika, pobierana jest wartość. Domyślna. Każdy system graficzny posiada jakieś źródło wartości domyślnych. W matplotlib jest to **rcParams** (Run-time configuration). a następnie wykonywane jest rendering.

■ Żaden parametr nie może pozostać nieokreślony

## Sposób definiowania parametrów domyślnych w matplotlib

---

- **stała** - np. kolor tła lub osi
- **zależność funkcyjna** - np. wielkość czcionki zależy od rozmiaru wykresu
- **wartość cykliczna (cycler)** - np. kolor kolejnego obiektu jest pobierany z listy cyklicznej
- **wartość losowa** - pobierana losowo, najczęściej z predefiniowanej listy (nie występują w matplotlib, ale są powszechne np w Excelu)
- **wyliczana** - system graficzny próbuje znaleźć optymalną wartość parametru na podstawie analizy danych (z reguły zawodzi)

## Taksonomia wykresów - wizualizacja wartości

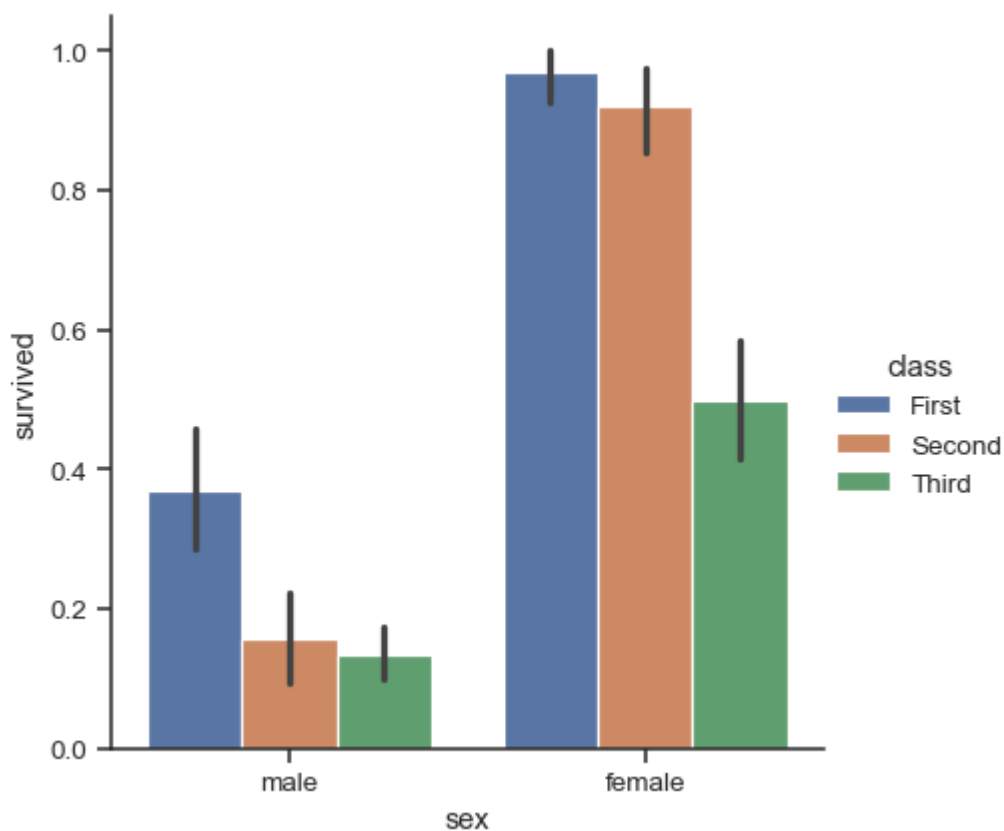
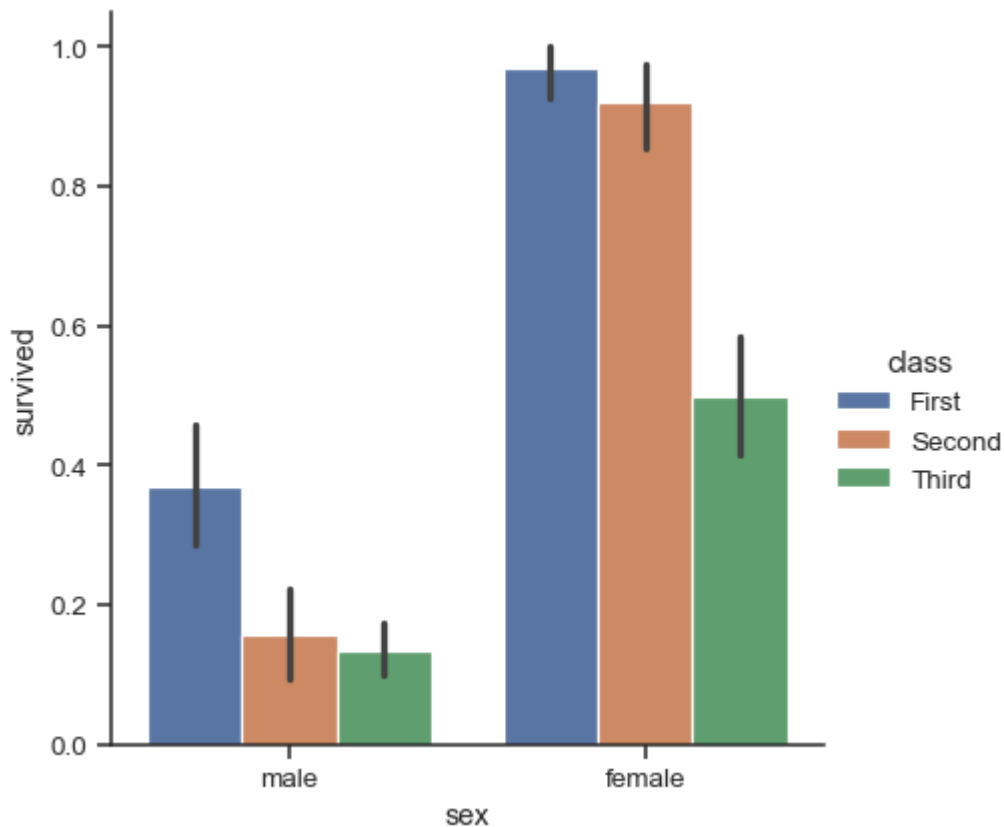
---

1. Wartości kategoryzowane
2. Wartości ciągłe (1D i 2D)
3. Wartości kategoryzowane i ciągłe
4. Relacje między wartościami
5. Wykresy wielozmienne
6. Wykresy grupujące

## Wartości kategoryzowane

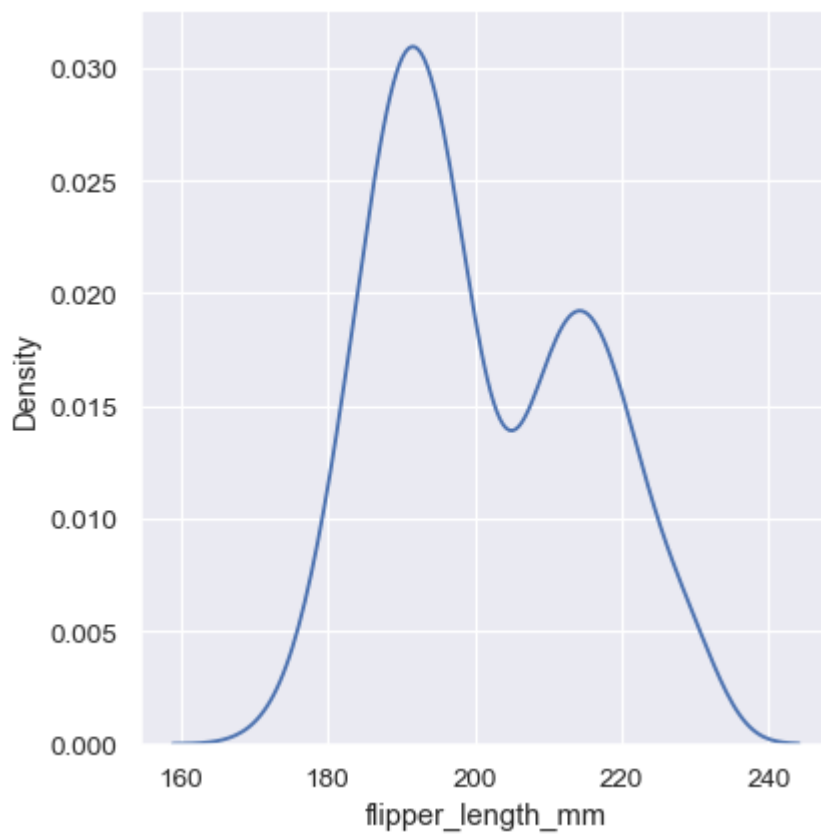
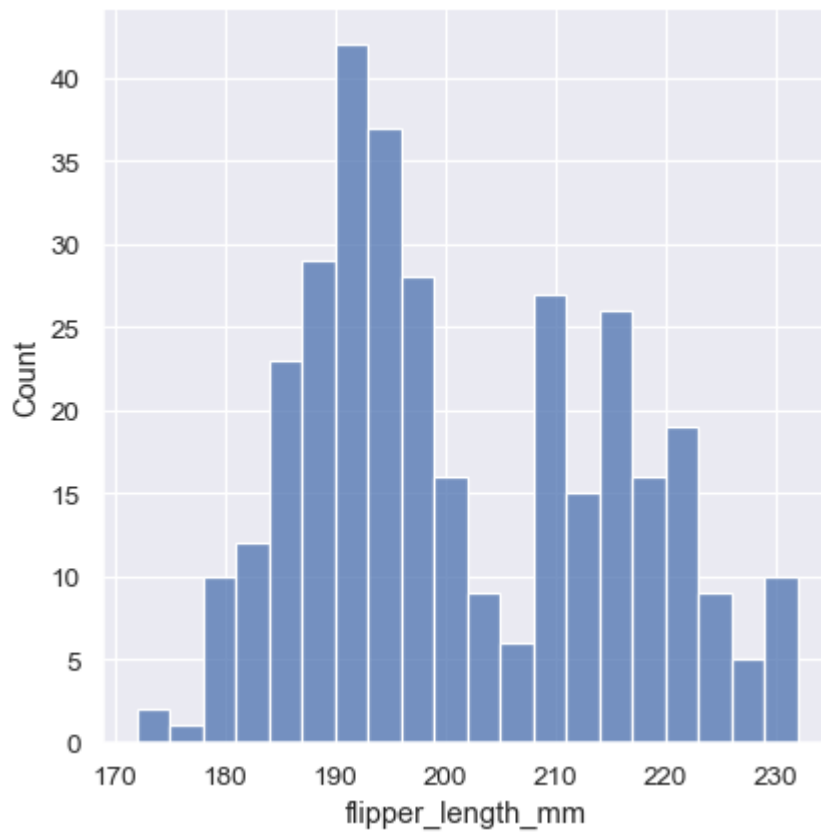
---

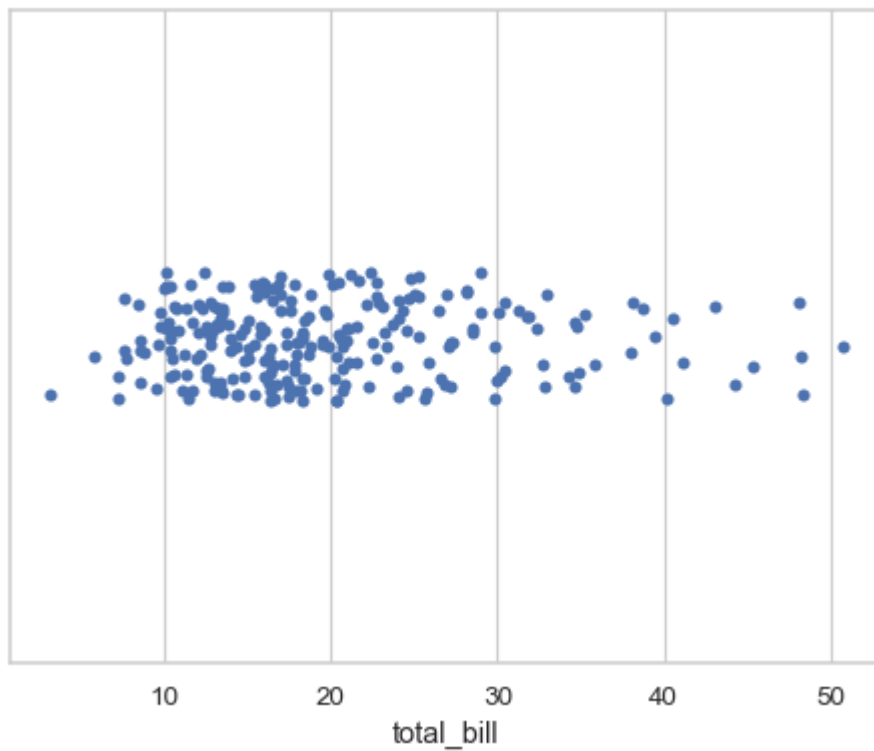
Wykresy mogą pokazywać tylko licznosc kategorii



## Wartości ciągłe

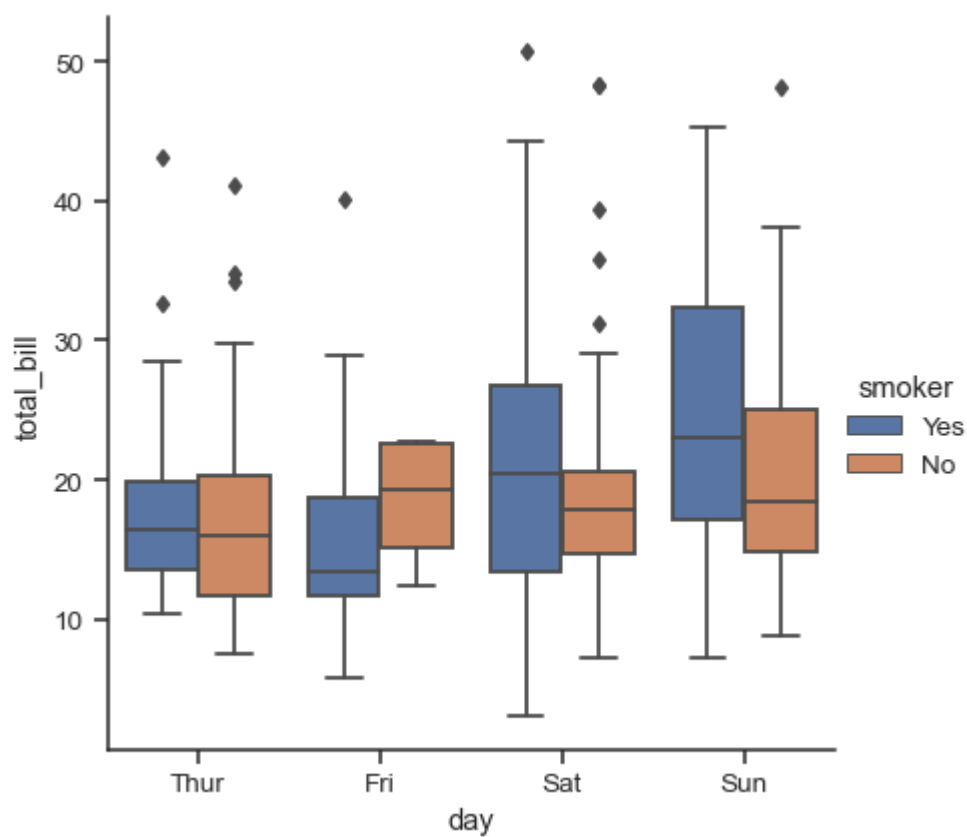
Wartości ciągłe prezentowane są w formach ilustrujących rozkłady statystyczne

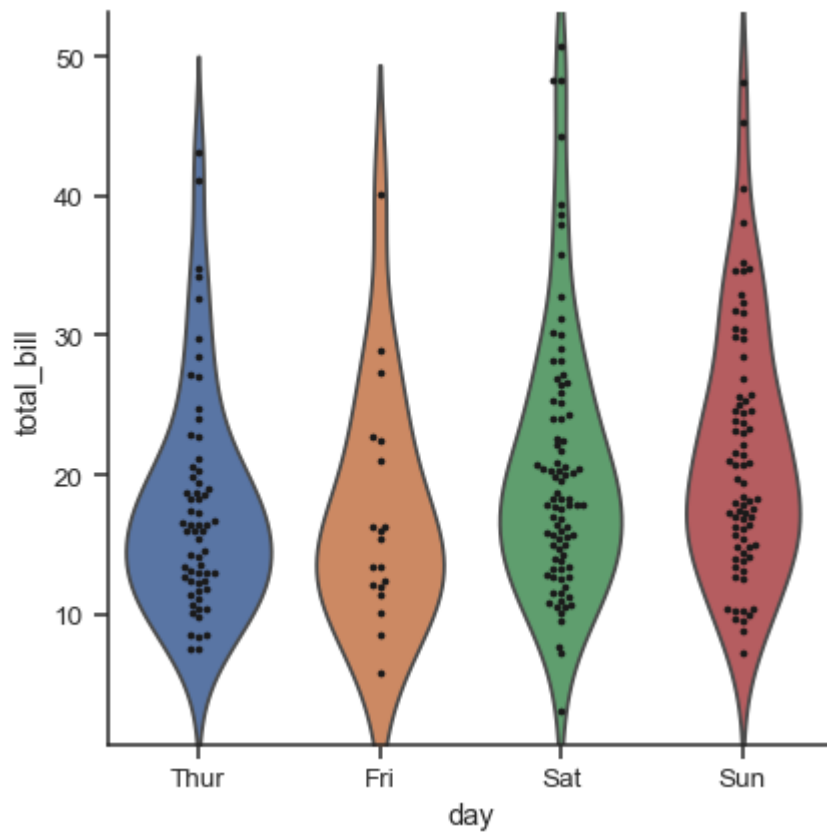




## Zmienne ciągłe w obrębie kategorii

Stosuje się wykresy dla rozkładu podzielone na kategorie

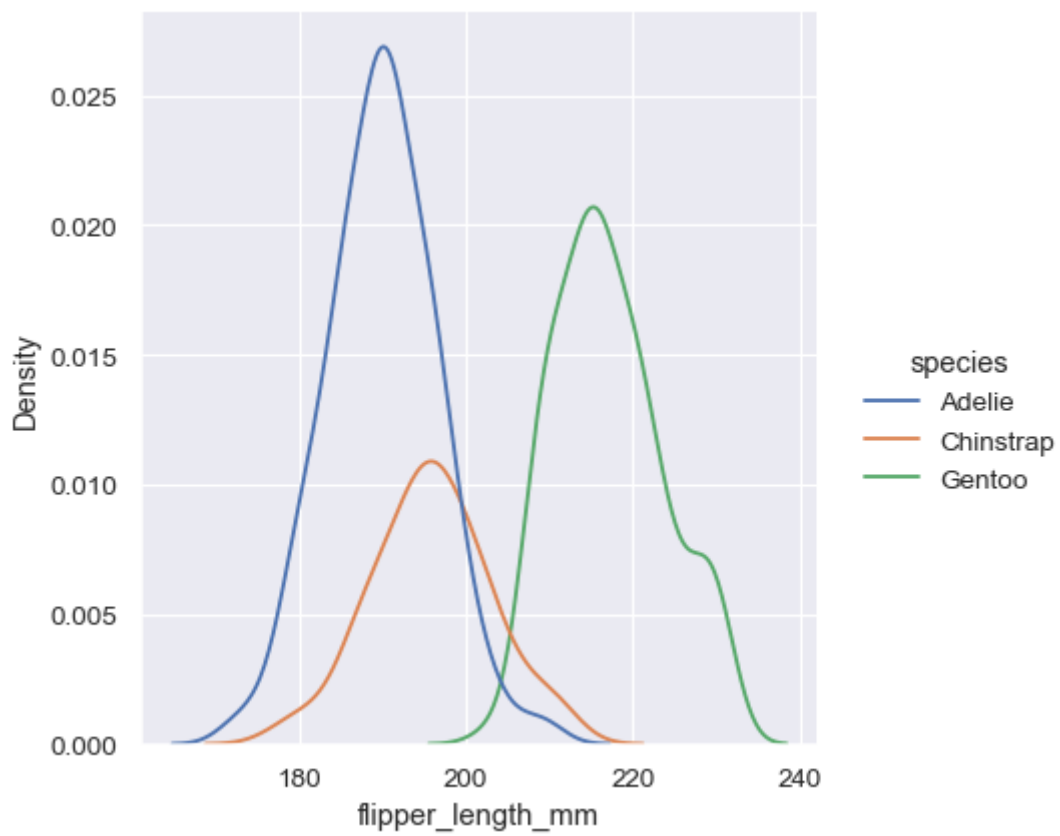
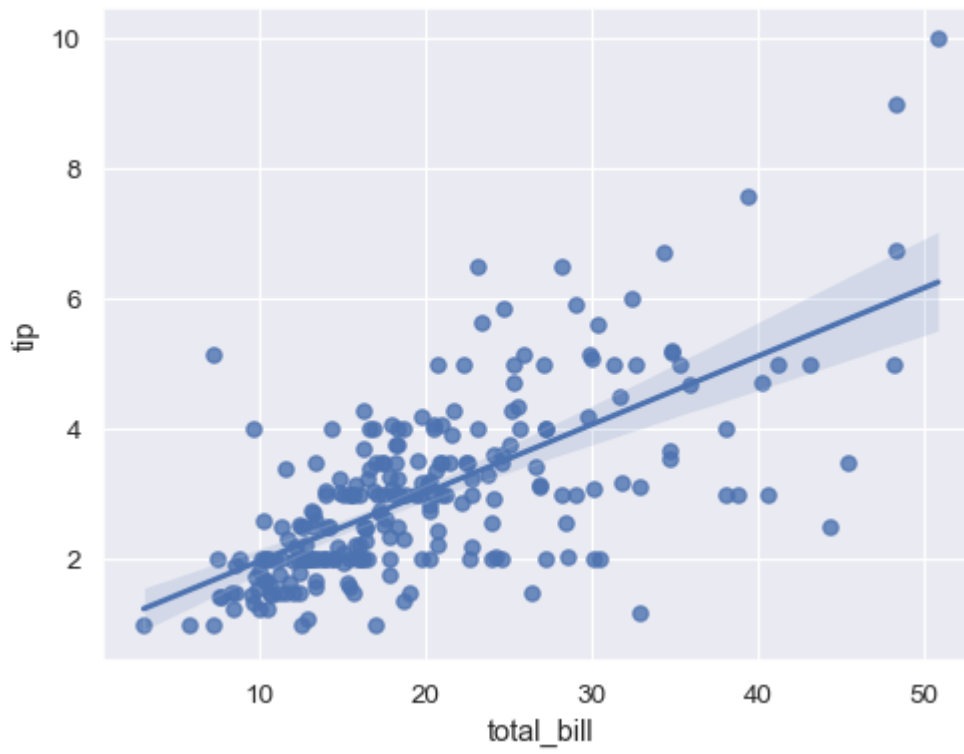


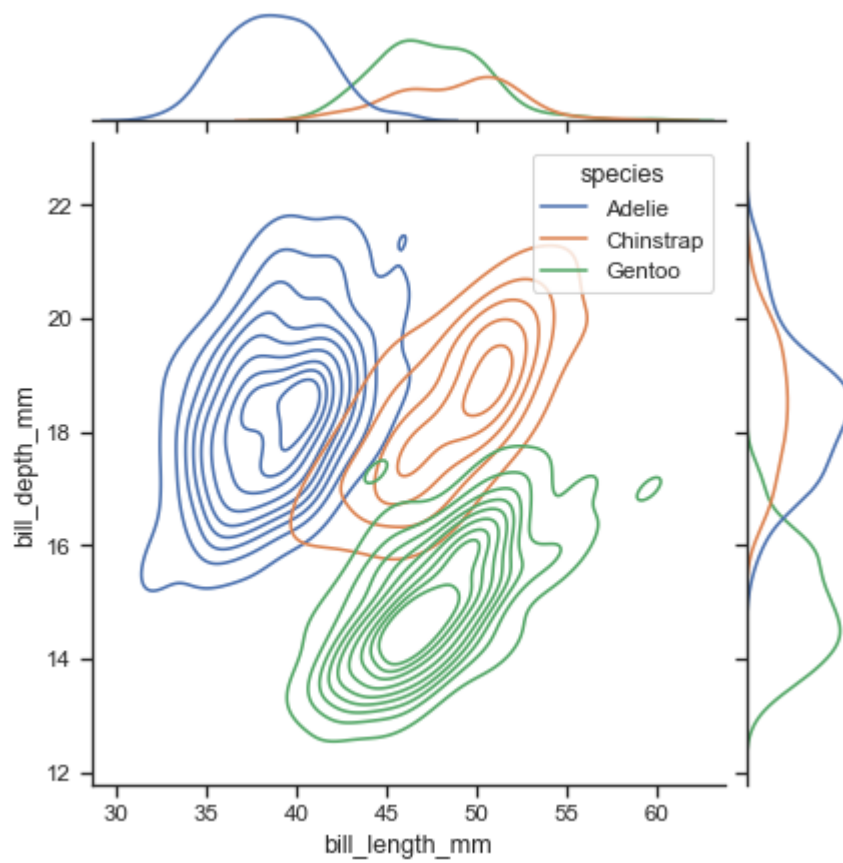


## Relacje między wartościami

---

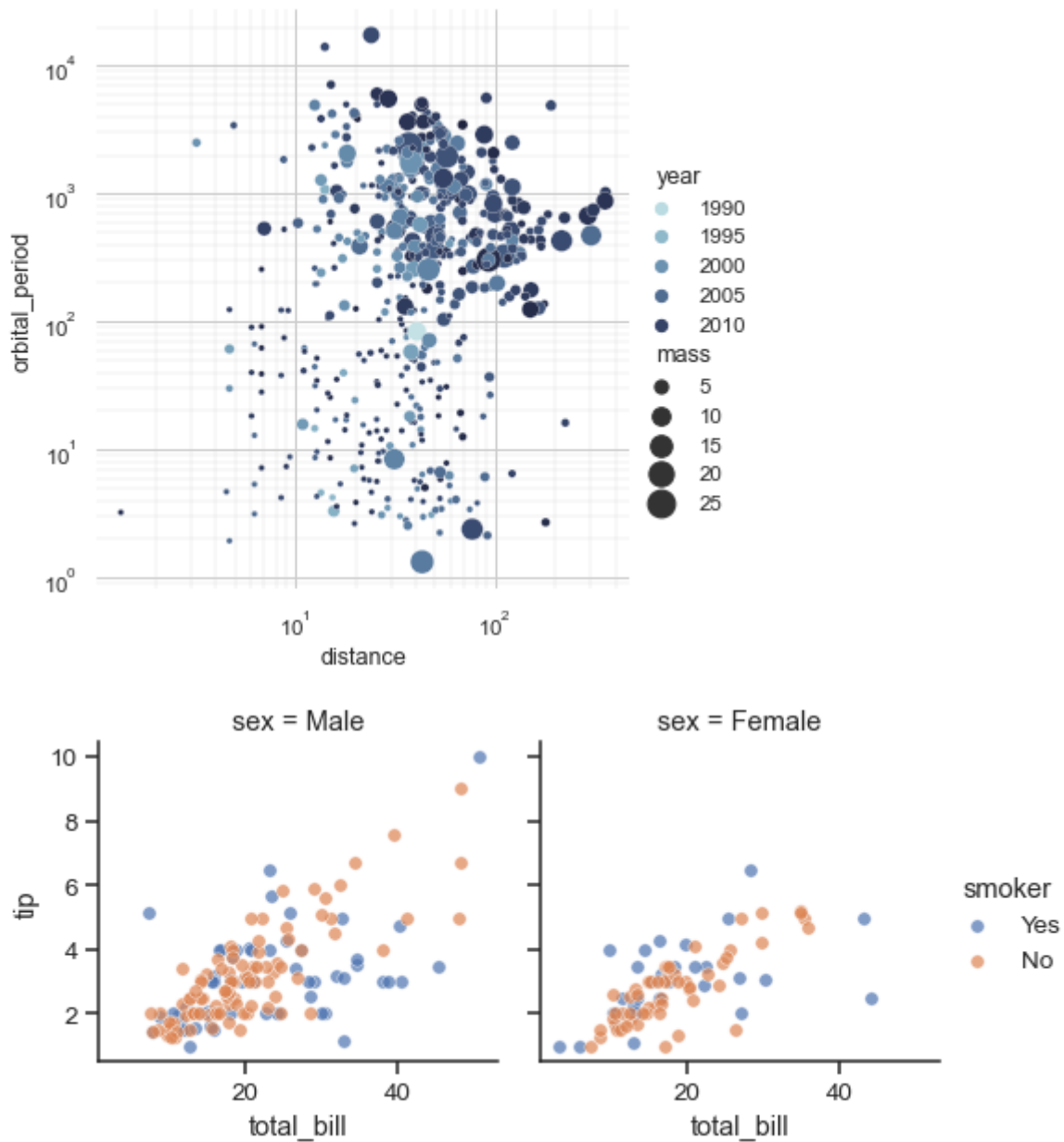






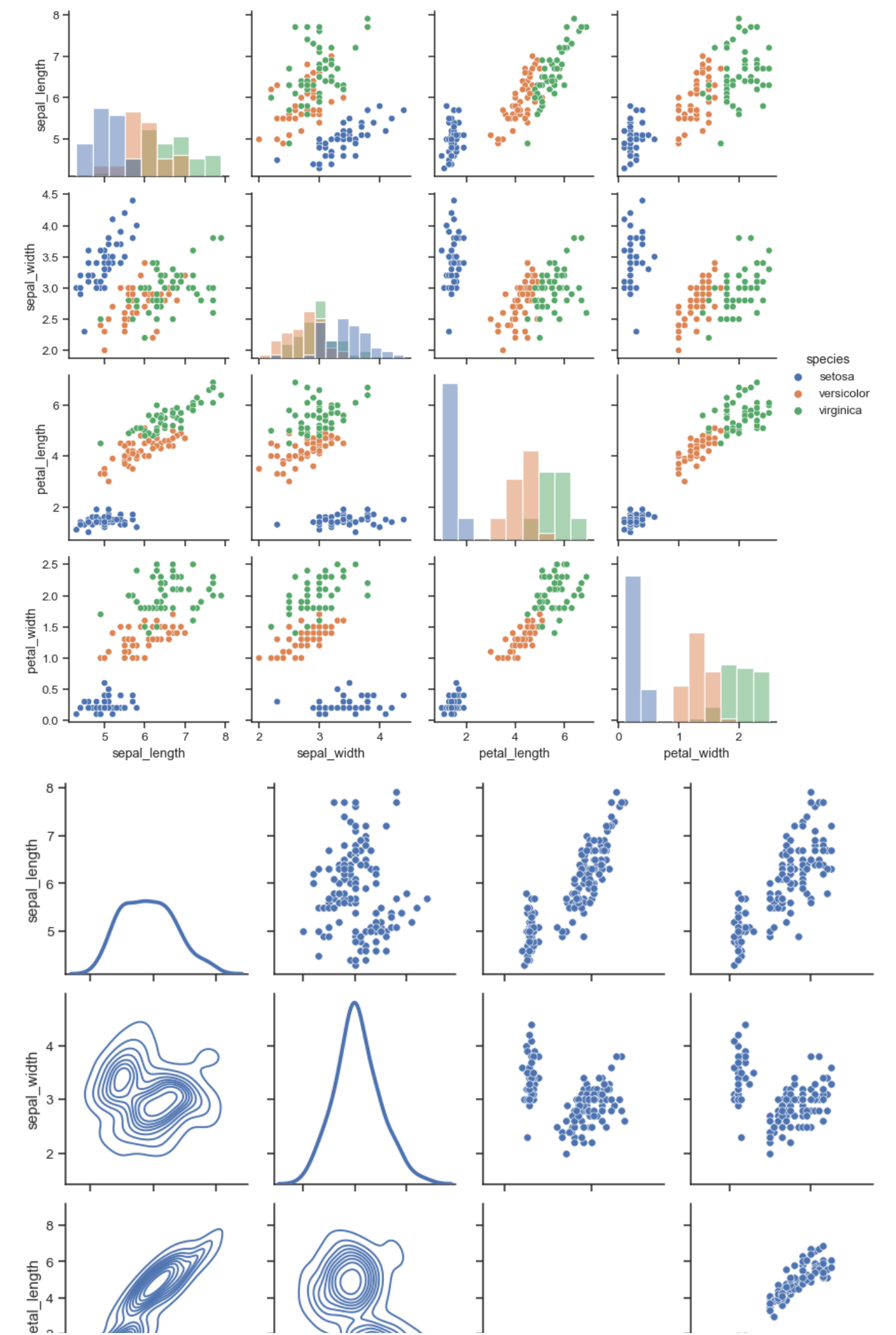
Wiele zmiennych na wykresie

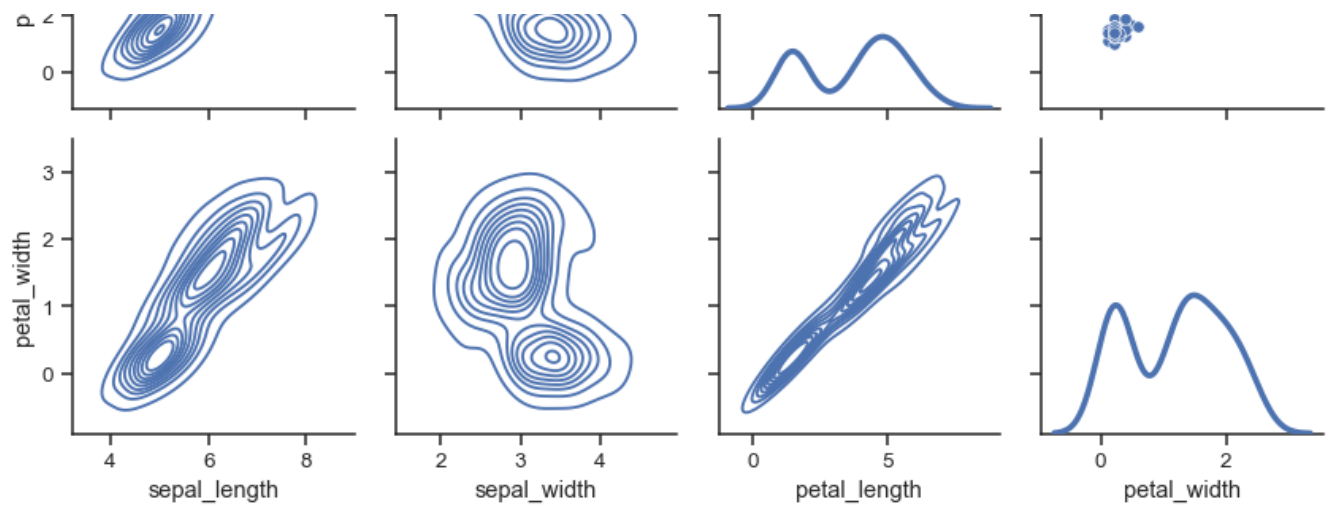
---



## Wykresy wielozmienne (pairplots)

---





## Wykresy złożone

---

