

Spis treści

1	Rozkłady używane w statystyce	1
1.1	Rozkład normalny	1
1.2	Rozkład χ^2	2
1.3	Rozkład t-Studenta	2
1.4	Rozkład F-Snedecora	2
2	Metoda estymacji punktowej	3
3	Metoda Monte Carlo	3
4	Bootstrapping	3
5	Testowanie hipotez	3
6	Testy	4
6.1	Test t-Studenta	4
6.2	ANOVA	4
6.3	Test Shapiro-Wilka	4
6.4	F test	4
6.5	Testy statystycznej różności	4
7	Reference	4
7.1	Korelacja	4
7.2	Regresja	5
7.3	Szybka analiza danych	5
7.4	Wizualizacja danych	5
7.5	Podział danych na k grup	5

1 Rozkłady używane w statystyce

1.1 Rozkład normalny

$$X \sim N(\mu, \sigma^2)$$

gdzie μ to wartość oczekiwana, a σ^2 to wariancja. Rozkład normalny jest rozkładem ciągłym, który jest symetryczny względem średniej. Wartość oczekiwana i mediana są równe.

Listing 1: gęstość w punkcie x

```
1 dnorm(x, mean = 0, sd = 1)
```

Listing 2: dystrybuenta

```
1 pnorm(q, mean = 0, sd = 1)
```

Listing 3: kwantyl p-tego percentyla

```
1 qnorm(p, mean = 0, sd = 1)
```

Listing 4: n losowych zmiennych z rozkładu normalnego

```
1 rnorm(n, mean = 0, sd = 1)
```

1.2 Rozkład χ^2

Jeśli X_1, \dots, X_n są niezależne i $X_{1..n} \sim N(0, 1)$ to:

$$Z = X_1^2 + X_2^2 + \dots + X_n^2 \sim \chi^2(n)$$

$$\mathbb{E}Z = n \quad \text{Var}(Z) = 2n$$

Jeśli zmienne niezależne $X_1, \dots, X_n \sim N(m_n, 1)$ to zmienna losowa też ma rozkład χ^2 , lecz z parametrem niecentralności $m = \sqrt{m_1^2 + \dots + m_n^2}$. Wtedy $\mathbb{E}Z = n + m$ $\text{Var}(Z) = 2(n + 2m)$.

Listing 5: gęstość w punkcie x

```
1 dchisq(x, df = 1)
```

Listing 6: dystrybuenta

```
1 pchisq(q, df = 1)
```

Listing 7: kwantyl p-tego percentyla

```
1 qchisq(p, df = 1)
```

Listing 8: n losowych zmiennych z rozkładu chi-kwadrat

```
1 rchisq(n, df = 1)
```

1.3 Rozkład t-Studenta

Jeżeli $Z \sim N(0, 1)$, $X \sim \chi^2(k)$ to wtedy zmienna:

$$Y = \frac{Z}{\sqrt{\frac{X}{k}}} \sim t(k)$$

Listing 9: gęstość w punkcie x

```
1 dt(x, df = 1)
```

Listing 10: dystrybuenta

```
1 pt(q, df = 1)
```

Listing 11: kwantyl p-tego percentyla

```
1 qt(p, df = 1)
```

Listing 12: n losowych zmiennych z rozkładu t-Studenta

```
1 rt(n, df = 1)
```

1.4 Rozkład F-Snedecora

Jeżeli $X \sim \chi^2(k_1)$ i $Y \sim \chi^2(k_2)$ to wtedy zmienna:

$$Z = \frac{X/k_1}{Y/k_2} \sim F(k_1, k_2)$$

$$\mathbb{E}Z = \frac{k_2}{k_2 - 2} \quad \text{Var}(Z) = \frac{2k_2^2(k_1 + k_2 - 2)}{k_1(k_2 - 2)^2(k_2 - 4)}$$

Listing 13: gęstość w punkcie x

```
1 df(x, df1 = 1, df2 = 1)
```

Listing 14: dystrybuenta

```
1 pf(q, df1 = 1, df2 = 1)
```

Listing 15: kwantyl p-tego percentyla

```
1 qf(p, df1 = 1, df2 = 1)
```

Listing 16: n losowych zmiennych z rozkładu F

```
1 rf(n, df1 = 1, df2 = 1)
```

2 Metoda estymacji punktowej

Jak dopasować rozkład i parametry do danych?

1. Wybieramy n próbek z danych ($X_1 \dots X_n$)
2. Patrzymy na histogram i oceniamy vibe
3. Dla parametrów wybieramy estymator, i przy pomocy estymatorów obliczamy parametry rozkładu
4. Sprawdzamy, czy rozkład pasuje do danych

Każdy estymator ma swój zakres ufności, z reguły określany przy pomocy wzoru. Mając obliczony zakres ufności, możemy określić poziom ufności, czyli błąd estymacji. Jeśli realna wartość parametru leży poza przedziałem ufności, to sugeruje wadę w estymacji. Poziom ufności to procent prób, w których przedział ufności zawiera prawdziwą wartość parametru.

Listing 17: estymuj parametry rozkładu normalnego

```
1 enorm(x, method = "mle", ci=TRUE, ci.type="two-sided", conf.level = 0.95)
```

3 Metoda Monte Carlo

Metoda Monte Carlo to metoda numeryczna, która polega na symulacji losowych próbek z rozkładu i obliczeniu wartości funkcji na podstawie tych próbek. W estymacji, na przykład, wystarczająco duża liczba próbek danych, pozwala stworzyć wykres wartości estymatora, który przybliży dystrybucję estymatora, co może pozwolić na lepsze oszacowanie wartości parametru.

Każda metoda, w której wykorzystujemy losowe próbki do obliczenia wartości funkcji, to metoda Monte Carlo.

4 Bootstrapping

Mając próbkę danych, możemy stworzyć wiele próbek z tej samej populacji, z reguły poprzez losowanie z zwracaniem. Dla wystarczająco dużej próbki początkowej w ten sposób możemy stworzyć wiele próbek, które będą miały podobny rozkład do oryginalnej próbki. Próba stworzona w ten sposób nazywana jest próbą bootstrapową.

5 Testowanie hipotez

Hipoteza statystyczne, to przypuszczenie dotyczące danych. Do weryfikacji hipotez korzystamy z testów. Dla konkretnego testu wyznacza się poziom istotności α , który określa prawdopodobieństwo odrzucenia hipotezy zerowej. Wielkość $1 - \beta$ dla układu hipotez prostych nazywa się mocą testu hipotezy zerowej wobec (prostej) hipotezy alternatywnej. Testy również są parametryzowane przez c czyli wartość krytyczną testu. W obecnie używanych implementacjach komputerowych wartości krytyczne zastępowane są tzw. p-wartościami (p-value), według pomysłu Ronalda Fishera. Jest to prawdopodobieństwo wylosowania próby takiej lub bardziej skrajnej, jak zaobserwowana przy założeniu, że hipoteza zerowa jest prawdziwa. Inaczej mówiąc, jest to prawdopodobieństwo, że zależność, jaką otrzymaliśmy w próbie z populacji mogła wystąpić przypadkowo, wskutek losowej zmienności, chociaż w populacji nie występuje.

6 Testy

Konkretne algorytmy mające na celu weryfikację hipotez statystycznych.

6.1 Test t-Studenta

Test t-Studenta jest testem statystycznym, który służy do porównania średnich dwóch grup, aby sprawdzić, czy są one statystycznie różne.

6.2 ANOVA

ANOVA (analiza wariancji) jest testem statystycznym, który służy do porównania średnich więcej niż dwóch grup, aby sprawdzić, czy są one statystycznie różne. ANOVA jest rozszerzeniem testu t-Studenta, który służy do porównania średnich dwóch grup.

6.3 Test Shapiro-Wilka

Test Shapiro-Wilka jest testem statystycznym, który służy do sprawdzenia, czy dane pochodzą z rozkładu normalnego.

6.4 F test

F test jest testem statystycznym, który służy do porównania wariancji dwóch grup, aby sprawdzić, czy są one statystycznie różne.

6.5 Testy statystycznej różności

Wszystkie te testy służą do porównania średnich więcej niż dwóch grup, aby sprawdzić, czy są one statystycznie różne. Różnią się one tym jak bardzo są dokładne i konserwatywne.

- HSD Tukeya
- LSD Fishera
- Test Studenta-Newmana-Keulsa
- Test Scheffego
- Test Duncana
- Test Dunnetta dla porównania z kontrolą

7 Reference

Kilka powszechnych celów oraz jak je osiągnąć w R.

7.1 Korelacja

Zupełnie jak $cov(x, y)$.

Listing 18: oblicz korelację

```
1 cor(x, y, method = c("pearson", "kendall", "spearman"))
```

7.2 Regresja

Klasyczna regresja liniowa, czyli przyporządkowanie danym prostej, która najlepiej je opisuje.

Listing 19: regresja liniowa

```
1 lm(y ~ x, data = data.frame(x, y))
```

Listing 20: regresja wielomianowa

```
1 lm(y ~ poly(x, degree = 2), data = data.frame(x, y))
```

Listing 21: różnice pomiędzy danymi a ich regresją

```
1 residuals(lm(y ~ x, data = data.frame(x, y)))
```

7.3 Szybka analiza danych

Zwróci nam podsumowanie dot. danych, ich strukturę, pierwsze i ostatnie wiersze oraz ich wymiary.

Listing 22: szybka analiza danych

```
1 summary(data)
2 str(data)
3 head(data)
4 tail(data)
5 dim(data)
```

7.4 Wizualizacja danych

Listing 23: wykresy

```
1 plot(x, y)
2 hist(x)
3 boxplot(x)
4 barplot(x)
5 pairs(data)
```

7.5 Podział danych na k grup

Przyporządkuje dane do k grup, w których każda grupa ma jak najbliżej siebie elementy.

Listing 24: podział danych na k grup

```
1 kmeans(data, centers = k)
```