**Q2:** $H = 100.000$ terms

$$t = \frac{\bar{X} - \mu}{\sqrt{\frac{\Delta^2}{H}}}$$

a) "adolf" apare de 150 ori

"hitler" apare de 200 ori

"adolf hitler" apare de 175 ori

$H_0$: "adolf" și "hitler" apar împreună întâmplător (independent) =)

$$P(\text{"adolf hitler"}) = P(\text{"adolf"}) * P(\text{"hitler"})$$

$$= \frac{150}{100000} \cdot \frac{200}{100.000} = \frac{30}{10.000.000} \simeq 3 \cdot 10^{-6}$$

expected mean: $\mu = 3 \cdot 10^{-6}$

observed variation: $\sigma^2 = p(1-p) \approx p = \Delta^2 = \bar{X}(1-\bar{X}) \simeq \bar{X} = 175 \cdot 10^{-5}$

observed mean: $\bar{X} = \frac{175}{100000} = 175 \cdot 10^{-5}$

$$t = \frac{175 \cdot 10^{-5} - 3 \cdot 10^{-6}}{\sqrt{\frac{175 \cdot 10^{-5}}{10^5}}} = \frac{174,7 \cdot 10^{-5}}{\sqrt{\frac{175}{10^{10}}}} = \frac{174,7}{10^5} \cdot \frac{10^5}{13,228} \simeq 13,21$$

Deoarece $t = 13,21 > 2,573$ (critical value) =) Ipoteza $H_0$ este falsă =) avem de-a face cu o colocație.

b) "hitler" apare de 200 ori

"industrial" apare de 700 ori

"hitler industrial" apare de 4 ori

$H_0$: "hitler" și "industrial" apar împreună independent =)

$$P(\text{"hitler industrial"}) = P(\text{"hitler"}) \cdot P(\text{"industrial"})$$

$$= \frac{200}{100000} \cdot \frac{700}{100.000} = \frac{14}{1.000.000} = 14 \cdot 10^{-6}$$

①

$$\mu = 14 \cdot 10^{-6}$$
$$\Delta^2 = \overline{X}(1-\overline{X}) \approx \overline{X} = 4 \cdot 10^{-5}$$
$$\overline{X} = 4 \cdot 10^{-5}$$

$$\Rightarrow t = \frac{4 \cdot 10^{-5} - 14 \cdot 10^{-6}}{\sqrt{\frac{4 \cdot 10^{-5}}{10^5}}} = \frac{(4 - 1,4) \cdot 10^{-5}}{\sqrt{\frac{4}{10^{10}}}} = \frac{2,6}{10^5} \cdot \frac{10^5}{2} = 1,3$$

Deoarece $t = 1,3 < 2,573$ (the critical value) $\Rightarrow$ Ipoteza $H_0$ este adevărată $\Rightarrow$
" hitler industrial" nu este o colocație.

c) " hitler" $\rightarrow$ 200 apariții
" revolution" $\rightarrow$ 900 apariții
" hitler revolution" $\rightarrow$ 14 apariții

$H_0$: " hitler" și " revolution" apar împreună întâmplător $\Rightarrow$
$P(\text{" hitler revolution"}) = P(\text{" hitler"}) \cdot P(\text{" revolution"})$
$$= \frac{200}{100000} \cdot \frac{900}{100000} = 18 \cdot 10^{-6}$$

$$\mu = 18 \cdot 10^{-6}$$
$$\overline{X} = 14 \cdot 10^{-5}$$
$$\Delta^2 = \overline{X}(1-\overline{X}) \simeq \overline{X} = 14 \cdot 10^{-5}$$

$$\Rightarrow t = \frac{14 \cdot 10^{-5} - 18 \cdot 10^{-6}}{\sqrt{\frac{14 \cdot 10^5}{10^5}}} = \frac{12,2}{10^5} \cdot \frac{10^5}{3,741} = 3,261$$

Deoarece $t = 3,261 > 2,573$ (the critical value) $\Rightarrow$ Ipoteza este falsă $\Rightarrow$
" hitler revolution" este o colocație.

d) " revolution" $\rightarrow$ 900 apariții
" hitler" $\rightarrow$ 200 apariții
" revolution hitler" $\rightarrow$ 25 apariții

$H_0$: "revolution" și "hitler" apar împreună în mod independent $\Rightarrow$

$P(\text{"revolution hitler"}) = P(\text{"revolution"}) \cdot P(\text{"hitler"})$

$$= \frac{900}{100.000} \cdot \frac{200}{100.000} = 18 \cdot 10^{-6}$$

$\mu = 18 \cdot 10^{-6}$

$\bar{X} = 25 \cdot 10^{-5}$

$\Delta^2 \simeq \bar{X} = 25 \cdot 10^{-5}$

$\Rightarrow t = \dfrac{25 \cdot 10^{-5} - 18 \cdot 10^{-6}}{\sqrt{\dfrac{25 \cdot 10^{-5}}{10^5}}} = \dfrac{23,2}{10^5} \cdot \dfrac{10^5}{5} = 4,64$

Deoarece $t = 4,64 > 2,573$ (critical value) $\Rightarrow$ Ip. $H_0$ este falsă $\Rightarrow$
"revolution hitler" este o colocație.

e) "industrial" $\rightarrow$ 700 apariții

"revolution" $\rightarrow$ 900 apariții

"industrial revolution" $\rightarrow$ 250 apariții

$H_0$: "industrial" și "revolution" apar împreună întâmplător $\Rightarrow$

$P(\text{"industrial revolution"}) = P(\text{"industrial"}) \cdot P(\text{"revolution"})$

$$= \frac{700}{100.000} \cdot \frac{900}{100.000} = 63 \cdot 10^{-6}$$

$\mu = 63 \cdot 10^{-6}$

$\bar{X} = 250 \cdot 10^{-5}$

$\Delta^2 \simeq \bar{X} = 250 \cdot 10^{-5}$

$\Rightarrow t = \dfrac{250 \cdot 10^{-5} - 63 \cdot 10^{-6}}{\sqrt{\dfrac{250 \cdot 10^{-5}}{10^5}}} = \dfrac{243,7}{10^5} \cdot \dfrac{10^5}{15,811} = 15,413$

Deoarece $t = 15,413 > 2,573$ (the critical value) $\Rightarrow$ Ip. $H_0$ este falsă $\Rightarrow$
"industrial revolution" este colocație.

③

## Q4:

| | $w_1 =$ garden | $w_1 \neq$ garden | | $w_1 =$ watch | $w_1 \neq$ watch |
|---|---|---|---|---|---|
| $w_2 =$ soil | 15 | 50 | $w_2 =$ dog | 20 | 50 |
| $w_2 \neq$ soil | 200 | 400 | $w_2 \neq$ dog | 200 | 1000 |

a) Chi-squared for "garden soil"

$$X_1^2 = \frac{(15+50+200+400)(15\cdot400 - 200\cdot50)^2}{(15+200)(50+400)(15+50)(200+400)}$$

$$= \frac{665\cdot(6000 - 10000)^2}{215\cdot450\cdot65\cdot600}$$

$$= \frac{665\cdot16.000.000}{215\cdot450\cdot65\cdot600} = \frac{10.640.000}{3.773.250} = 2,819$$

b) Chi-squared for "watch dog"

$$X_2^2 = \frac{(20+50+200+1000)(20\cdot1000 - 200\cdot50)^2}{(20+200)\cdot(50+1000)\cdot(20+50)(200+1000)} =$$

$$= \frac{1270\cdot(10^4)^2}{220\cdot1050\cdot70\cdot1200} =$$

$$= \frac{1270\cdot100000000}{220\cdot1050\cdot70\cdot1200} = \frac{1270000}{194.040} = 6,545$$

c) degree of freedom: $r=c=2 \Rightarrow (r-1)(c-1)=1$ } $\Rightarrow$ critical value $=3,84$
$p=0,05$                             (valoare luată din tabelul de la $X^2$)

Deoarece $X_1^2 = 2,819 < 3,84$ (the critical value) $\Rightarrow$ "garden soil" nu este o colocație,
iar pt. că $X_2^2 > X_1^2 \Rightarrow$ "watch dog" are șanse mai mari decât "garden soil" să fie o colocație.

Deoarece $X_2^2 = 6,545 > 3,84$ (critical value) $\Rightarrow$ "watch dog" chiar este o colocație.

④