

```

1 from IPython.display import Image
2
3 PATH = "https://raw.githubusercontent.com/SIIIKOR/sad2/main/projekt2/"
4 IMG_PATH = PATH + "/BN.png"
5 TASK_1_IMG_PATH = PATH + "/TASK_1.jpg"
6 TASK_9_IMG_PATH = PATH + "/TASK_9.jpg"
7 Image(url=IMG_PATH)

```

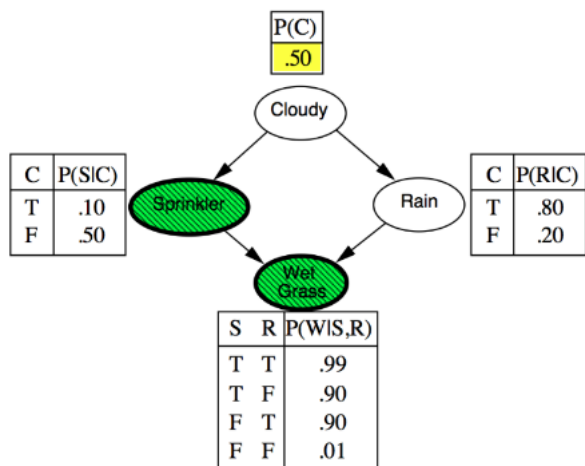


Figure 1: The *Rain Bayesian network*

```

1 from random import choices, random
2 import numpy as np
3 import matplotlib.pyplot as plt
4 import statsmodels.api as sm

```

## TASK 1

```

1 Image(url=TASK_1_IMG_PATH, width=3200//3, height=1809//3)

```

$$P(C=T | R=T, S=T, W=T) = P(C=T | R=T, S=T) =$$

$$P(C=T | R=F, S=T, W=T) = P(C=T | R=F, S=T) = \frac{P(C=T) P(R=F | C=T) P(S=T | C=T)}{P(C=T) P(R=F | C=T) P(S=T | C=T) + P(C=F) P(R=F | C=F) P(S=T | C=F)}$$

$$P(R=T | C=T, S=T, W=T) = \frac{P(R=T | C=T) P(W=T | R=T, S=T)}{P(R=T | C=T) P(W=T | R=T, S=T) + P(R=F | C=T) P(W=T | R=F, S=T)}$$

$$P(R=T | C=F, S=T, W=T) = \frac{P(R=T | C=F) P(W=T | R=T, S=T)}{P(R=T | C=F) P(W=T | R=T, S=T) + P(R=F | C=F) P(W=T | R=F, S=T)}$$

```

1 P_ONE = 0.444
2 P_TWO = 0.048
3 P_THREE = 0.815

```

```

# P(C=T | R=T, S=T, W=T)
# P(C=T | R=F, S=T, W=T)
# P(R=T | C=T, S=T, W=T)

```

```

4 P_FOUR = 0.216 # P(R=T | C=F, S=T, W=T)
5
6 CLOUDY_DIST = (
7     #C=FALSE    C=TRUE
8     (1-P_TWO,   P_TWO),    #R=FALSE
9     (1-P_ONE,   P_ONE)     #R=TRUE
10 )
11 RAIN_DIST = (
12     #R=FALSE    R=TRUE
13     (1-P_FOUR,  P_FOUR),    #C=FALSE
14     (1-P_THREE, P_THREE)    #C=TRUE
15 )
16
17 DIST = (CLOUDY_DIST, RAIN_DIST)
18
19 CHOICES = (0, 1)

```

## ▼ TASK 2

```

1 def gibbs_sample(n):
2     samples = np.empty((n, 2), dtype=int)
3     samples[0, 0] = 0 # cloudy
4     samples[0, 1] = 1 # rain
5
6     for i in range(1, n): # for each time step excluding
7         idx_var = int(random() > 0.5) # Randomly select index of the
8         idx_other_var = 1 - idx_var # Index of the other variable
9
10        x_other_var_prev = samples[i-1, idx_other_var] # Value of the other variable
11        x_var = choices(CHOICES, weights=DIST[idx_var][idx_other_var])[0] # Sample chosen variable value
12        x_other_var = choices(CHOICES, weights=DIST[idx_other_var][x_var])[0] # Sample other variable value
13
14        samples[i, idx_var] = x_var
15        samples[i, idx_other_var] = x_other_var
16
17    return samples

```

## ▼ TASK 3

```

1 samples = gibbs_sample(100)
2
3 print("P(R=T|S=T,W=T) =", samples[:, 1].mean())

```

$P(R=T|S=T,W=T) = 0.32$

## ▼ TASK 4

```

1 samples_one = gibbs_sample(50_000)
2 samples_two = gibbs_sample(50_000)

```

## ▼ TASK 5

```

1 def get_zero_count(ones_count, idx):
2     return (idx+1) - ones_count
3
4 def full_time_relative_counts(samples):
5     counts = np.empty((len(samples), 2), dtype=int) #cloudy rain one counts
6     counts[0, 0] = 0 # initial values
7     counts[0, 1] = 0
8
9     for i, (x_cloudy, x_rain) in enumerate(samples):
10        x_cloudy_freq = counts[i-1, 0] + x_cloudy if i != 0 else x_cloudy
11        x_rain_freq = counts[i-1, 1] + x_rain if i != 0 else x_rain
12        counts[i, 0] = x_cloudy_freq
13        counts[i, 1] = x_rain_freq
14
15    return counts
16

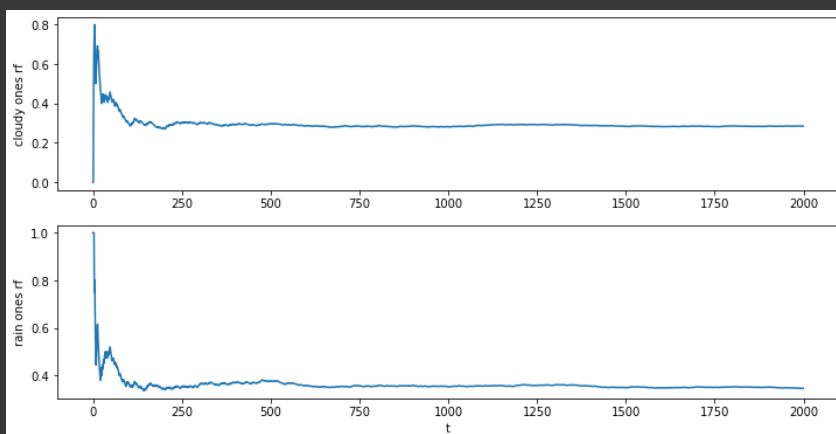
```

```

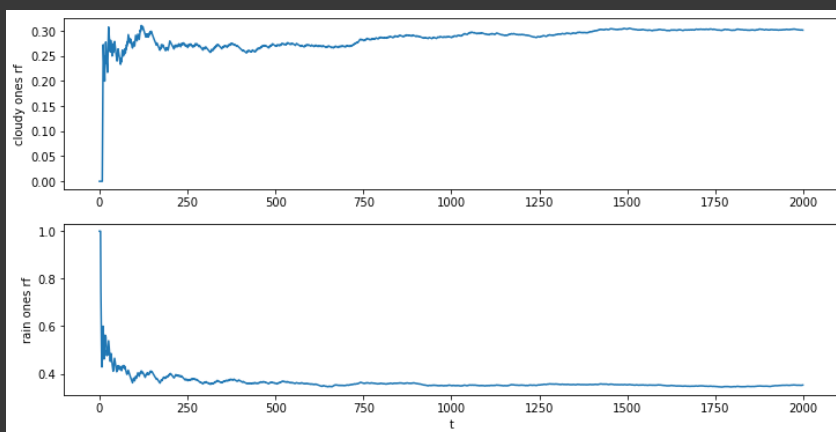
17 def counts_to_rf(samples):
18     rf = np.empty((len(samples), 2), dtype=float)
19     for i, val in enumerate(samples):
20         rf[i] = val / i if i != 0 else val
21     return rf
22
23 def plot_freqs(samples):
24     fig, axs = plt.subplots(2, figsize=(12, 6))
25     x = np.arange(len(samples))
26
27     axs[0].plot(x, samples[:, 0])
28     axs[0].set_ylabel("cloudy ones rf")
29     axs[1].plot(x, samples[:, 1])
30     axs[1].set_xlabel('t')
31     axs[1].set_ylabel("rain ones rf")
32
33     fig.show()
34
35 samples_one_rf = counts_to_rf(full_time_relative_counts(samples_one))
36 samples_two_rf = counts_to_rf(full_time_relative_counts(samples_two))

```

```
1 plot_freqs(samples_one_rf[:2000])
```



```
1 plot_freqs(samples_two_rf[:2000])
```



burn in time (takie najmniejsze  $t$ , że wykres zdaje się zbiegać) dla niektórych generacji może być tak małe jak  $t=100$  a dla innych może wymagać nawet  $t=400$ .

Mimo tego, że może wymagać tylko 400, czasem zdarza się tak, że po pewnym czasie wartości zaczynają w miarę liniowo dość powoli spadać lub wzrastać. W związku z tym decyduje się wziąć większe  $t$ , takie, po którym nie zauważy tych wzrostów lub spadków.

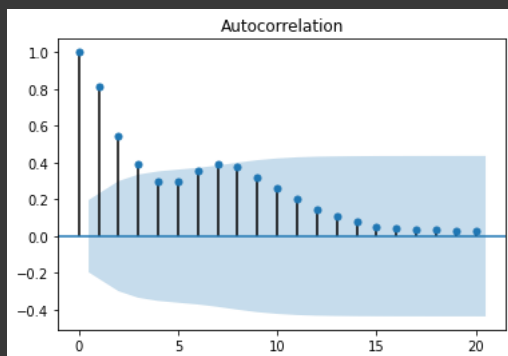
W związku z tym wydaje mi się, że w miarę bezpieczny burn in time to  $t = 2000$

## TASK 6

Zredukowałem ilość sampli do pokazania na wykresie, inaczej z wykresu nie jestem w stanie nic odczytać.

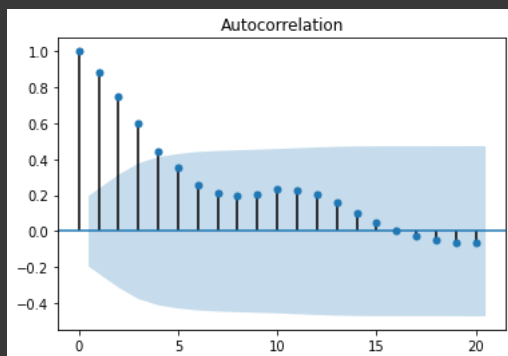
## Cloudy

```
1 sm.graphics.tsa.plot_acf(samples_one_rf[:100, 0])
2 plt.show()
```



## Rain

```
1 sm.graphics.tsa.plot_acf(samples_one_rf[:100, 1])
2 plt.show()
```



Możemy zauważyć, że korelacja na początku jest dodatnia. W przypadku cloudy spada do zera po interval=10 natomiast w przypadku rain jest zerowa tak szybko jak już po interval=4 lub interval=5. Dalej dla większych wartości oscyluje wokół zera.

Wybrałbym interval=10 tak aby dla obu było ok.

## TASK 7

```
1 def gibbs_sample_optimised(n, burn_in=2000, interval=10):
2     samples = np.empty((burn_in + interval*n, 2), dtype=int)           # we will have to make interval
3     samples[0, 0] = 0                                                  # cloudy
4     samples[0, 1] = 1                                                  # rain
5
6     for i in range(1, burn_in + n*interval):                          # for each time step excluding
7         idx_var = int(random() > 0.5)                                # Randomly select index of the
8         idx_other_var = 1 - idx_var                                    # Index of the other variable
9
10        x_other_var_prev = samples[i-1, idx_other_var]                # Value of the other variable
11        x_var = choices(CHOICES, weights=DISTS[idx_var][idx_other_var])[0] # Sample chosen variable value
12        x_other_var = choices(CHOICES, weights=DISTS[idx_other_var][x_var])[0] # Sample other variable value
13
14        samples[i, idx_var] = x_var
15        samples[i, idx_other_var] = x_other_var
16
17    return samples[burn_in::interval]                                   # take only interval'th element
```

## TASK 8

```
1 samples = gibbs_sample_optimised(100)
2
3 print("P(R=T|S=T,W=T) =", samples[:, 1].mean())
```

$$P(R=T|S=T, W=T) = 0.33$$

## TEST

```
1 var = 1
2 t = 100
3 T = 1000
4 base_version_samples = np.array([gibbs_sample(t)[: , var].mean() for _ in range(T)])
5 improved_version_samples = np.array([gibbs_sample_optimised(t)[: , var].mean() for _ in range(T)])
```

```
1 print("base gibbs mean: ", base_version_samples.mean())
2 print("improved gibbs mean: ", improved_version_samples.mean())
3 print()
4 print("base gibbs variance: ", base_version_samples.var())
5 print("improved gibbs variance: ", improved_version_samples.var())
```

```
base gibbs mean: 0.35996
improved gibbs mean: 0.34945
```

```
base gibbs variance: 0.0022059984
improved gibbs variance: 0.0022645975
```

W obu przypadkach mean jest prawie taki sam (powinien być identyczny). Natomiast wersja z burn in i spłaszczaniem ma mniejszą wariancję, czyli ogólnie jej wyniki są bliżej rzeczywistej wartości niż w przypadku wersji bazowej.

Pokazuje to, że mój burn in i interval poprawiają jakość samplowania, ale nie jakoś szczególnie bardzo.

## TASK 9

```
1 Image(url=TASK_9_IMG_PATH, width=3200//3, height=1809//3)
```

$$P(R=T|S=T, W=T) = \frac{P(R=T, S=T, W=T)}{P(S=T, W=T)} \approx \frac{0.178}{0.556}$$

$$P(R=T, S=T, W=T) = P(R=T|C=T) P(S=T|C=T) P(W=T|C=T) \\ P(R=T|C=F) P(S=T|C=F) P(W=T|C=F)$$

$$\approx 0.178$$

$$P(S=T, W=T) = P(R=F|C=T) P(S=T|C=T) P(W=T|C=T) \\ P(R=F|C=F) P(S=T|C=F) P(W=T|C=F) \\ P(R=T, S=T, W=T)$$

$$\approx 0.556$$

Otrzymywane przeze mnie wyniki nie są identyczne, ale są bardzo zbliżone do 0.32

## BONUS

Bazując na <https://www.imperial.ac.uk/media/imperial-college/research-centres-and-groups/astrophysics/public/icic/data-analysis-workshop/2018/Convergence-Tests.pdf?fbclid=IwAR2SkO2-oWiUOiNJK5CbuSZXJs9YJP4qmxuYBxrZ7jY3gHXs9vLSHBfmH08>

```
1 def gelman_rubin_test(x):
2     chain amount = x.shape[0]
```

