

Modelowanie i analiza sieci złożonych

IV. Metryki sieci

Grzegorz Siudem

Politechnika Warszawska



**Politechnika
Warszawska**

Unia Europejska
Europejski Fundusz Społeczny



Zadanie 10 pn.

„Przygotowanie i uruchomienie nowego kierunku studiów na studiach II stopnia
- Inżynieria i Analiza Danych (IAD)”

realizowane jest w ramach projektu
„NERW PW. Nauka – Edukacja – Rozwój – Współpraca”
współfinansowanego ze środków Unii Europejskiej
w ramach Europejskiego Funduszu Społecznego

Przed zajęciami

Co to jest odległość w grafie

$$d(i, j) = ?$$

Poznane już metryki sieci

- średni stopień wierzchołka

$$\langle k \rangle = \frac{1}{N} \sum_{i=1}^N k_i,$$

- średnia droga w grafie

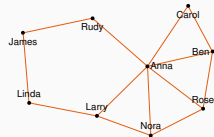
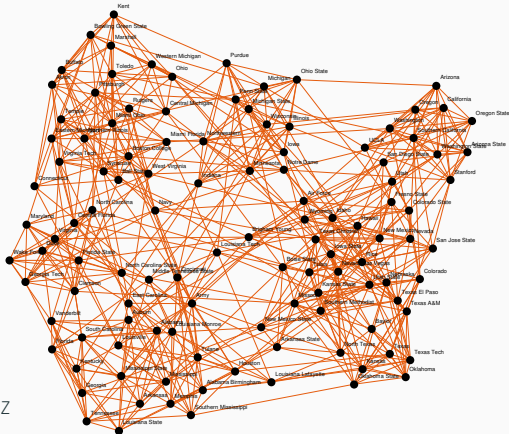
$$\ell = \frac{1}{N(N-1)} \sum_{i \neq j} d(i, j).$$

Lemat o uścisku dłoni

Wykład

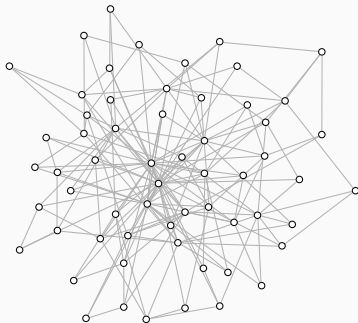
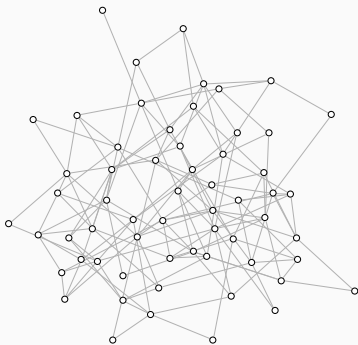
Jakie cechy sieci możemy zmierzyć?

- jak duża jest sieć?



Jakie cechy sieci możemy zmierzyć?

- jak duża jest sieć?
- jak gęsta jest sieć?



Jakie cechy sieci możemy zmierzyć?

- jak duża jest sieć?
- jak gęsta jest sieć?
- jaka jest topologia połączeń?

?

Uwaga!

W sieciologii topologia ma nieco inne znaczenie niż w matematyce!

Jakie znamy potoczne metryki/miary sieci?

Naturalne/naiwne metryki:

- liczba wierzchołków N (rozmiar),

Jakie znamy potoczne metryki/miary sieci?

Naturalne/naiwne metryki:

- liczba wierzchołków N (rozmiar),
- liczba krawędzi E (rozmiar, gęstość),

Jakie znamy potoczne metryki/miary sieci?

Naturalne/naiwne metryki:

- liczba wierzchołków N (rozmiar),
- liczba krawędzi E (rozmiar, gęstość),
- dlaczego $\langle k \rangle = 2E/N$? (gęstość)

Jakie znamy potoczne metryki/miary sieci?

Naturalne/naiwne metryki:

- liczba wierzchołków N (rozmiar),
- liczba krawędzi E (rozmiar, gęstość),
- dlaczego $\langle k \rangle = 2E/N$? (gęstość)
- największy stopień wężła (celebrytki/celebryci?),

Jakie znamy potoczne metryki/miary sieci?

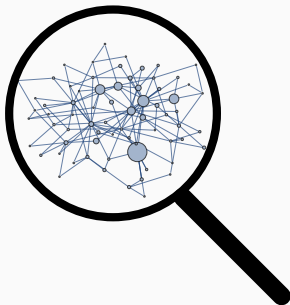
Naturalne/naiwne metryki:

- liczba wierzchołków N (rozmiar),
- liczba krawędzi E (rozmiar, gęstość),
- dlaczego $\langle k \rangle = 2E/N$? (gęstość)
- największy stopień wężła (celebrytki/celebryci?),
- inne?

Jakie znamy potoczne metryki/miary sieci?

Naturalne/naiwne metryki:

- liczba wierzchołków N (rozmiar),
- liczba krawędzi E (rozmiar, gęstość),
- dlaczego $\langle k \rangle = 2E/N$? (gęstość)
- największy stopień wężła (celebrytki/celebryci?),
- inne?



Na zajęciach mówiliśmy już o

- rozkładzie stopni wierzchołków $\mathcal{P}(k)$,

Na zajęciach mówiliśmy już o

- rozkładzie stopni wierzchołków $\mathcal{P}(k)$,
- w szczególności o rozkładach potęgowych z parametrem α

$$\mathcal{P}(k) \propto k^{-\alpha},$$

- jak jednak estymować α ?

Na zajęciach mówiliśmy już o

- rozkładzie stopni wierzchołków $\mathcal{P}(k)$,
- w szczególności o rozkładach potęgowych z parametrem α

$$\mathcal{P}(k) \propto k^{-\alpha},$$

- jak jednak estymować α ?

Kto odrobił pracę domową?

- M.E.J. Newman, *Power laws, Pareto distributions and Zipf's law*, Contemporary Physics, **46**, 323-351 (2005).
i/lub
- rozdziały 3.1-3.3 w A. Fronczak, P. Fronczak, *Świat sieci złożonych* PWN (2009).

Na zajęciach mówiliśmy już o

- rozkładzie stopni wierzchołków $\mathcal{P}(k)$,
- w szczególności o rozkładach potęgowych z parametrem α

$$\mathcal{P}(k) \propto k^{-\alpha},$$

- jak jednak estymować α ?

Kto odrobił pracę domową?

- M.E.J. Newman, *Power laws, Pareto distributions and Zipf's law*, Contemporary Physics, **46**, 323-351 (2005).
i/lub
- rozdziały 3.1-3.3 w A. Fronczak, P. Fronczak, *Świat sieci złożonych* PWN (2009).

Więcej na ten temat w części projektowej.

Miary korelacji węzłów

Znamy już:

Korelacje dwuwęzłowe $\mathcal{R}(k_i, k_j)$ czyli prawdopodobieństwo, że losowo wybrana krawędź łączy węzły o stopniach k_i i k_j

$$\mathcal{R}(k_i, k_j) = \frac{P(k_i, k_j)}{P_u(k_i, k_j)},$$

a P_u odpowiada sieci nieskorelowanej o tym samym rozkładzie.
Niestety nie jest to zbyt praktyczne narzędzie.

Miary korelacji węzłów

Znamy już:

Korelacje dwuwęzłowe $\mathcal{R}(k_i, k_j)$ czyli prawdopodobieństwo, że losowo wybrana krawędź łączy węzły o stopniach k_i i k_j

$$\mathcal{R}(k_i, k_j) = \frac{P(k_i, k_j)}{P_u(k_i, k_j)},$$

a P_u odpowiada sieci nieskorelowanej o tym samym rozkładzie. Niestety nie jest to zbyt praktyczne narzędzie.

Słyszeliśmy także o:

Prawdopodobieństwie warunkowym

$$\mathcal{P}(k_i|k_j) = \frac{\mathcal{P}(k_i, k_j)}{k_j \mathcal{P}(k_j) / \langle k \rangle}$$

Czy to się *dobrze* estymuje? Niestety nie...

Jak obniżyć korelacje?

Losowe przełączanie węzłów:



Zachowuje stopnie wierzchołków!

Czemu to robimy?

- żeby pozbyć się niechcianych korelacji,
- żeby określić ich znaczenie dla danej sieci,
- żeby zyskać model referencyjny o tym samym rozkładzie.

Wprowadźmy:

Średni stopień najbliższego węzła (dla węzła o stopniu k_i)

$$\langle k \rangle_{nn}(k_i) = \frac{1}{k_i} \sum_{j=1}^N a_{ij} k_j = \sum_{\ell} \ell \mathcal{P}(\ell | k_i).$$

Pytamy o zależność średniego stopnia najbliższego sąsiada $\langle k \rangle_{nn}$ od stopnia wierzchołka k_i

Miary korelacji węzłów cd.

Wprowadźmy:

Średni stopień najbliższego węzła (dla węzła o stopniu k_i)

$$\langle k \rangle_{nn}(k_i) = \frac{1}{k_i} \sum_{j=1}^N a_{ij} k_j = \sum_{\ell} \ell \mathcal{P}(\ell | k_i).$$

Pytamy o zależność średniego stopnia najbliższego sąsiada $\langle k \rangle_{nn}$ od stopnia wierzchołka k_i

Co mierzy ta miara?

- jeśli $\langle k \rangle_{nn}(k_i)$ jest funkcją rosnącą wówczas sieć jest asortatywna.

Wprowadźmy:

Średni stopień najbliższego węzła (dla węzła o stopniu k_i)

$$\langle k \rangle_{nn}(k_i) = \frac{1}{k_i} \sum_{j=1}^N a_{ij} k_j = \sum_{\ell} \ell \mathcal{P}(\ell | k_i).$$

Pytamy o zależność średniego stopnia najbliższego sąsiada $\langle k \rangle_{nn}$ od stopnia wierzchołka k_i

Co mierzy ta miara?

- jeśli $\langle k \rangle_{nn}(k_i)$ jest funkcją rosnącą wówczas sieć jest asortatywna.
- jeśli $\langle k \rangle_{nn}(k_i)$ jest funkcją malejącą wówczas sieć jest dysasortatywna.

Wprowadźmy:

Średni stopień najbliższego węzła (dla węzła o stopniu k_i)

$$\langle k \rangle_{nn}(k_i) = \frac{1}{k_i} \sum_{j=1}^N a_{ij} k_j = \sum_{\ell} \ell \mathcal{P}(\ell | k_i).$$

Pytamy o zależność średniego stopnia najbliższego sąsiada $\langle k \rangle_{nn}$ od stopnia wierzchołka k_i

Co mierzy ta miara?

- jeśli $\langle k \rangle_{nn}(k_i)$ jest funkcją rosnącą wówczas sieć jest asortatywna.
- jeśli $\langle k \rangle_{nn}(k_i)$ jest funkcją malejącą wówczas sieć jest dysasortatywna.
- jeśli jest stała, sieć jest nieskorelowana.

Wprowadźmy:

Średni stopień najbliższego węzła (dla węzła o stopniu k_i)

$$\langle k \rangle_{nn}(k_i) = \frac{1}{k_i} \sum_{j=1}^N a_{ij} k_j = \sum_{\ell} \ell \mathcal{P}(\ell | k_i).$$

Pytamy o zależność średniego stopnia najbliższego sąsiada $\langle k \rangle_{nn}$ od stopnia wierzchołka k_i

Co mierzy ta miara?

- jeśli $\langle k \rangle_{nn}(k_i)$ jest funkcją rosnącą wówczas sieć jest asortatywna.
- jeśli $\langle k \rangle_{nn}(k_i)$ jest funkcją malejącą wówczas sieć jest dysasortatywna.
- jeśli jest stała, sieć jest nieskorelowana.
- a co jeśli zależność jest niemonotoniczna?

Miary korelacji węzłów cd.

W praktyce wszystkie poznane miary były zbyt złożony...

A zatem pozostaje nam liczyć współczynnik korelacji

$$r = \frac{\sum_{jk} jk(e_{jk} - q_j q_k)}{\sigma_q^2}$$

co przy notacji

- e_{jk} – łączny rozkład pozostałych stopni.
- rozkład pozostałych stopni $q_k = \sum_j e_{jk}$, ale z drugiej strony $q_k = \frac{(k+1)p_{k+1}}{\sum_{j \geq 1} j p_j}$

prowadzi do

$$r = \frac{\frac{1}{M} \sum_i k_i j_i - \left[\frac{1}{2M} \sum_i (j_i + k_i) \right]^2}{\frac{1}{2M} \sum_i (j_i^2 + k_i^2) - \left[\frac{1}{2M} \sum_i (j_i + k_i) \right]^2},$$

MASZ gdzie $i = 1, 2, \dots, M$ numeruje krawędzie, a j_i i k_i to stopnie wierzchołków przyłączonych do i -tej krawędzi.

Zjawisko homofilii



Proszę pana, ja jestem
umysł ścisły. Mnie się
podobają melodie, które
już raz słyszałem.
Po prostu. No...
To... Poprzez...
No reminiscencję.
No jakże może podobać
mi się piosenka, którą
pierwszy raz slysze?

”

NAJLEPSZE CYTATY Z POLSKICH KOMEDII

Źródło: gazeta.pl

Model P-O-X Heidera w skrócie:

Model P-O-X Heidera w skrócie:

- Przyjaciel mojego przyjaciela jest moim przyjacielem.

Model P-O-X Heidera w skrócie:

- Przyjaciół mojego przyjaciela jest moim przyjacielem.
- Przyjaciół mojego wroga jest moim wrogiem.

Model P-O-X Heidera w skrócie:

- Przyjaciół mojego przyjaciela jest moim przyjacielem.
- Przyjaciół mojego wroga jest moim wrogiem.
- Wrog mojego przyjaciela jest moim wrogiem.

Model P-O-X Heidera w skrócie:

- Przyjaciel mojego przyjaciela jest moim przyjacielem.
- Przyjaciel mojego wroga jest moim wrogiem.
- Wrog mojego przyjaciela jest moim wrogiem.
- Wróg mojego wroga jest moim przyjacielem.

Model P-O-X Heidera w skrócie:

- Przyjaciel mojego przyjaciela jest moim przyjacielem.
- Przyjaciel mojego wroga jest moim wrogiem.
- Wrog mojego przyjaciela jest moim wrogiem.
- Wróg mojego wroga jest moim przyjacielem.

Dotyczy to jednak skierowanych sieci społecznych...

Uprośćmy nasze rozważania do sieci nieskierowanych.

Definicja

Współczynnik gronowania wierzchołka to stosunek liczby E_i istniejących krawędzi pomiędzy sąsiadami wierzchołka do wszystkich możliwych krawędzi pomiędzy tymi sąsiadami

$$C_i = \frac{2E_i}{k_i(k_i - 1)}.$$

Współczynnik grafu to średnia po wszystkich wierzchołkach

$$C = \langle C_i \rangle.$$

Definicja

Współczynnik gronowania wierzchołka to stosunek liczby E_i istniejących krawędzi pomiędzy sąsiadami wierzchołka do wszystkich możliwych krawędzi pomiędzy tymi sąsiadami

$$C_i = \frac{2E_i}{k_i(k_i - 1)}.$$

Współczynnik grafu to średnia po wszystkich wierzchołkach

$$C = \langle C_i \rangle.$$

Alternatywna definicja współczynnika gronowania:

$$C_{\Delta} = \frac{3 \times \text{liczba trójkątów w sieci}}{\text{liczba dróg o długości 2 w sieci}}.$$

Zliczamy *motywy w sieci*

przeważnie porównuje się *Z-score* z ансамblem sieci odkorelowanych

$$Z = \frac{p - \langle p \rangle}{\sigma}.$$



Jak zmierzyć jak mały jest świat sieci?

Średnia odległość

$$\ell = \frac{1}{N(N-1)} \sum_{i \neq j} d(i, j)$$

Wydajność

$$\mathcal{E} = \frac{1}{N(N-1)} \sum_{i \neq j} [d(i, j)]^{-1}.$$

Jak zmierzyć jak mały jest świat sieci?

Średnia odległość

$$\ell = \frac{1}{N(N-1)} \sum_{i \neq j} d(i, j)$$

Wydajność

$$\mathcal{E} = \frac{1}{N(N-1)} \sum_{i \neq j} [d(i, j)]^{-1}.$$

Pytanie.

Czym różnią się te dwie metryki? Która jest *lepsz*a?

Który węzeł jest najważniejszy w sieci?

Poszukujemy najważniejszych *stacji przesiadkowych*.

Oznaczenia:

- δ_{jk} oznacza liczbę najkrótszych dróg łączących węzły j oraz k ,
- $\delta_{jk}^{(i)}$ oznacza liczbę najkrótszych dróg łączących węzły j oraz k przechodzących przez węzeł i .

Definicja

$$B_i = \frac{2}{(N-1)(N-2)} \sum_k \sum_{j>k} \frac{\delta_{jk}^{(i)}}{\delta_{jk}}.$$

A jeśli zapytamy o najważniejszą krawędź?

Poszukujemy najważniejszych *linii*.

Oznaczenia:

- $\delta_{jk}^{(e)}$ oznacza liczbę najkrótszych dróg łączących węzły j oraz k przechodzących przez krawędź e .

Definicja

$$B_i = \frac{2}{N(N-1)} \sum_k \sum_{j>k} \frac{\delta_{jk}^{(e)}}{\delta_{jk}}.$$

To jeden z głównych celów *sieciologii*

- cała sieć jest zbyt bogatą strukturą, żeby mówić o niej bez uproszczeń,

To jeden z głównych celów *sieciologii*

- cała sieć jest zbyt bogatą strukturą, żeby mówić o niej bez uproszczeń,
- często różne osoby interesują zupełnie inne cechy sieci,

To jeden z głównych celów *sieciologii*

- cała sieć jest zbyt bogatą strukturą, żeby mówić o niej bez uproszczeń,
- często różne osoby interesują zupełnie inne cechy sieci,
- często pewne szczególne miary potrzebne są do opisu pewnych szczególnych typów sieci...

- Indeks Hirscha – w sieci cytowań,

- Indeks Hirscha – w sieci cytowań,
- Liczba Erdős'a – w sieci cytowań,

- Indeks Hirscha – w sieci cytowań,
- Liczba Erdős'a – w sieci cytowań,
- Liczba Bacona – w sieci aktorów,

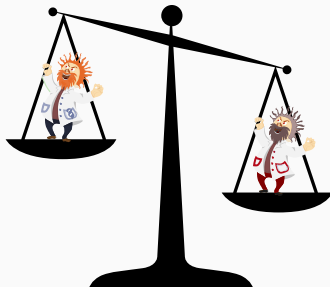
- Indeks Hirscha – w sieci cytowań,
- Liczba Erdős'a – w sieci cytowań,
- Liczba Bacona – w sieci aktorów,
- PageRank – w sieci www,

- Indeks Hirscha – w sieci cytowań,
- Liczba Erdős'a – w sieci cytowań,
- Liczba Bacona – w sieci aktorów,
- PageRank – w sieci www,
- próg epidemii – w epidemiologii,

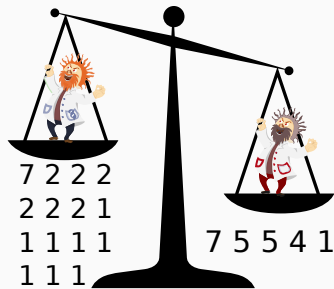
- Indeks Hirscha – w sieci cytowań,
- Liczba Erdős’a – w sieci cytowań,
- Liczba Bacona – w sieci aktorów,
- PageRank – w sieci www,
- próg epidemii – w epidemiologii,
- podatność na ataki/awarie – w inżynierii,

- Indeks Hirscha – w sieci cytowań,
- Liczba Erdős’a – w sieci cytowań,
- Liczba Bacona – w sieci aktorów,
- PageRank – w sieci www,
- próg epidemii – w epidemiologii,
- podatność na ataki/awarie – w inżynierii,
- detekcja społeczności (klastrow),

- Indeks Hirscha – w sieci cytowań,
- Liczba Erdős’a – w sieci cytowań,
- Liczba Bacona – w sieci aktorów,
- PageRank – w sieci www,
- próg epidemii – w epidemiologii,
- podatność na ataki/awarie – w inżynierii,
- detekcja społeczności (klastrow),
- wiele, wiele innych...



Jak zmierzyć sukces naukowy?



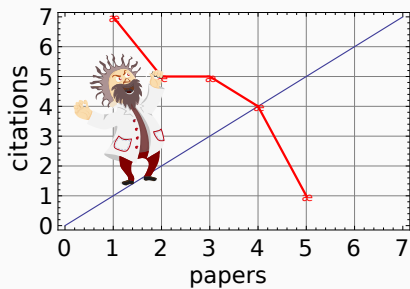
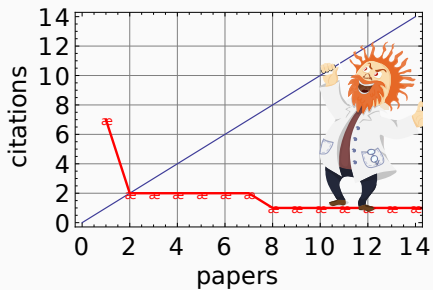
Policzmy cytowania (stopnie wierzchołków).

J.E. Hirsch, PNAS **102**, (2005).

$$h\text{-index} = \max \{h = 1, \dots, n : X_{(n-h+1)} \geq h\},$$

$$X_{(1)} \leq X_{(2)} \leq \dots \leq X_{(n)}.$$

Indeks Hirscha





wikipedia

Paul Erdős 1913-1996

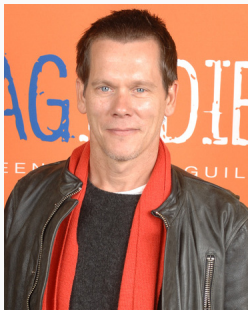
- węgierski matematyk
- na następnych zajęciach poznamy grafy Erdős-Rényiego.



wikipedia

Definicja

- Paul Erdős ma liczbę Erdősą równą 0.
- Liczbę Erdősą każdego innego naukowca określa się jako minimum z $e_i + 1$ gdzie e_i to liczba Erdősą jego/jej współautorów.



wikipedia

Kevin Bacon ur. 1958

- amerykański aktor, reżyser i producent filmowy,

Odpowiednik liczby Erősa w sieci aktorów

Przykładowe wartości:

- Elvis Preasley: 2,
- Ronald Reagan: 2,
- Andrzej Grabowski: 3,
- Andrzej Lepper: 3,
- Zdzisław Maklakiewicz: 3,
- Jan Himilsbach: 3,

Suma liczby Erdős i liczby Bacona:

- Steven Strogatz $E = 3 B = 1 \Rightarrow EB = 4$,
- Richard Feynman $E = 3 B = 3 \Rightarrow EB = 6$,
- Stephen Hawking $E = 4 B = 2 \Rightarrow EB = 6$,
- Natalie Portman $E = 5 B = 2 \Rightarrow EB = 7$,
- Colin Firth $E = 6 B = 1 \Rightarrow EB = 7$,
- Kristen Stewart $E = 5 B = 2 \Rightarrow EB = 7$,
- Mayim Bialik $E = 5 B = 2 \Rightarrow EB = 7$.

Suma liczby Erdős i liczby Bacona:

- Steven Strogatz $E = 3 B = 1 \Rightarrow EB = 4$,
- Richard Feynman $E = 3 B = 3 \Rightarrow EB = 6$,
- Stephen Hawking $E = 4 B = 2 \Rightarrow EB = 6$,
- Natalie Portman $E = 5 B = 2 \Rightarrow EB = 7$,
- Colin Firth $E = 6 B = 1 \Rightarrow EB = 7$,
- Kristen Stewart $E = 5 B = 2 \Rightarrow EB = 7$,
- Mayim Bialik $E = 5 B = 2 \Rightarrow EB = 7$.

Ciekawostka

Proszę poczytać o liczbie Erdős-Bacona-Black Sabbath...

PageRank

- metoda wyboru najistotniejszych stron www,
- zajmiemy się nią na 11. zajęciach.

PageRank

- metoda wyboru najistotniejszych stron www,
- zajmiemy się nią na 11. zajęciach.

Próg epidemii

- minimalna liczba zarażonych osób w sieci społecznej, która skutkuje wybuchem epidemii,
- zajmiemy się nim na 12. zajęciach.

PageRank

- metoda wyboru najistotniejszych stron www,
- zajmiemy się nią na 11. zajęciach.

Próg epidemii

- minimalna liczba zarażonych osób w sieci społecznej, która skutkuje wybuchem epidemii,
- zajmiemy się nim na 12. zajęciach.

Odporność na przypadkowe uszkodzenia i celowe ataki

- zajmiemy się nimi na 7. zajęciach.

PageRank

- metoda wyboru najistotniejszych stron www,
- zajmiemy się nią na 11. zajęciach.

Próg epidemii

- minimalna liczba zarażonych osób w sieci społecznej, która skutkuje wybuchem epidemii,
- zajmiemy się nim na 12. zajęciach.

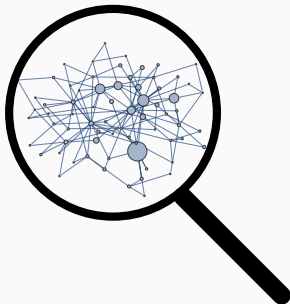
Odporność na przypadkowe uszkodzenia i celowe ataki

- zajmiemy się nimi na 7. zajęciach.

Detekcja społeczności (klastrow)

- zajmiemy się nią na 8. zajęciach.

Podsumowanie





**Politechnika
Warszawska**

Unia Europejska
Europejski Fundusz Społeczny



Zadanie 10 pn.

„Przygotowanie i uruchomienie nowego kierunku studiów na studiach II stopnia
- Inżynieria i Analiza Danych (IAD)”

realizowane jest w ramach projektu
„NERW PW. Nauka – Edukacja – Rozwój – Współpraca”
współfinansowanego ze środków Unii Europejskiej
w ramach Europejskiego Funduszu Społecznego

Dziękuję za uwagę!