

METODY ANALIZY OBRAZU – ANALIZA OBRAZU MAMMOGRAFICZNEGO NA PODSTAWIE CECH WYZNACZONYCH Z TEKSTURY

Jagoda Lazarek

Politechnika Łódzka, Wydział Fizyki Technicznej, Informatyki i Matematyki Stosowanej, Instytut Informatyki

Streszczenie. W artykule przedstawiono analizę możliwości zastosowania cech wyznaczonych z tekstury do klasyfikacji wykrytych, na obrazie mammograficznym, obszarów zainteresowania – jako obszarów niezmienionych lub zmienionych chorobowo. Cechy tekstury wyznaczono na podstawie histogramu, macierzy gradientu, macierzy długości pasm oraz macierzy zdarzeń. Klasyfikację przeprowadzono z wykorzystaniem klasyfikatora k-NN. W wyniku przeprowadzonych eksperymentów poprawnie rozpoznano wszystkie zmienione chorobowo próbki.

Słowa kluczowe: mammografia, mbrazowanie medyczne, mnaliza tekstury, mklasyfikacja obrazów

IMAGE ANALYSIS METHODS – ANALYSIS OF MAMMOGRAPHIC IMAGE BASED ON TEXTURAL FEATURES

Abstract. This paper presents an analysis of the possibility of using textural features for mammographic images classification. Textural features are calculated base on histogram, gradient matrix, run-length matrix, co-occurrence matrix. Classification is based on k-NN classifier, the regions of interest can be classified as normal or abnormal. Results of some experiments are presented. All of abnormal regions were classified correctly.

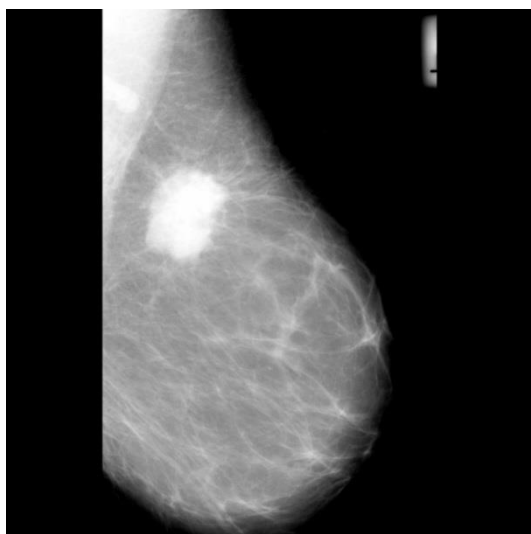
Keywords: mammography, medical diagnostic imaging, mmage texture analysis, mmage classification

Wstęp

Potrzeba tworzenia i udoskonalania metod przetwarzania i rozpoznawania różnego rodzaju obrazów wynika z ich szerokiego zastosowania w życiu codziennym, począwszy od medycyny, poprzez przemysł aż do rozrywki. Celem metod wspomaganych komputerowo jest semantyczna analiza obrazów, umożliwiająca automatyzację lub wspomaganie procesów decyzyjnych dokonywanych na ich podstawie. Zróżnicowanie otrzymywanych obrazów wynika ze sposobu i warunków ich wykonywania oraz rodzaju obiektów, które się na nich znajdują. Powoduje to konieczność dostosowania odpowiednich metod ich przetwarzania i analizy, w tym metod sztucznej inteligencji [9].

1. Zastosowania w medycynie - mammografia

Obrazy w medycynie oraz ich analiza stanowią podstawę w diagnostyce wielu schorzeń, m. in. złamań, nowotworów czy innych nieprawidłowości. Niektóre z badań, takie jak mammografia wykonywane są profilaktycznie, zwiększając tym samym szansę na wykrycie choroby we wczesnym stadium i wyleczenie. Ze względu na przesiewowy charakter badania ilość powstających do przeanalizowania danych jest bardzo duża [3, 6-7].

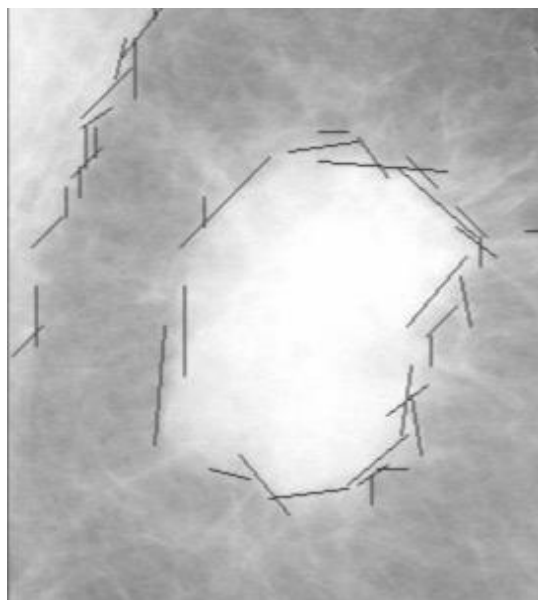


Rys. 1. Obraz nr 184 z bazy MIAS [11], centrum zmiany chorobowej znajduje się w punkcie o współrzędnych (352, 624) i obejmuje obszar ograniczony okręgiem o promieniu 114

Obrazy mammograficzne uzyskiwane są w wyniku prześwietlenia piersi promieniami Rentgena. Otrzymywany obraz jest obrazem w skali szarości. Metody komputerowej analizy obrazów mammograficznych wspomagają specjalistów w procesie diagnozy, dzięki czemu wpływają na zmniejszenie ryzyka pominięcia zmian chorobowych.

2. Wykrywanie zmian chorobowych

Wykrywanie zmian chorobowych jest skomplikowane ze względu na ich charakterystykę i niewyraźną granicę między nimi a pozostałym, niezmienionym obszarem [1]. Propozycje rozwiązania przedstawiono w artykule „Methods of Pattern Detection in Mammographic Images” [4].



Rys. 2. Obszar wykryty za pomocą metody opisanej w [4], który następnie będzie klasyfikowany jako niezmieniony lub zmieniony chorobowo

Metoda pozwala na wykrycie obszarów o wyróżniającej się w porównaniu do otoczenia jasności, w tym obszarów jednolitych oraz pasm. Obszary wykryte tą metodą, w następnym etapie powinny być klasyfikowane jako podejrzane lub niezmienione. W niniejszym artykule przedstawiono analizę możliwości wykorzystania cech statycznych obliczanych na podstawie tekstury do wspomnianej klasyfikacji.

3. Analiza tekstur

Tekstury reprezentowane są w formie dwuwymiarowej tablicy pikseli. Dla obrazów w skali szarości, każdy piksel reprezentuje jeden z poziomów szarości należących do zdefiniowanego przedziału. W artykule analizowane są obrazy, które mają 256 poziomów szarości, zatem poziom szarości g_L (ang. *grey level*) należy do przedziału od 0 do 255.

Tekstury mogą być wykorzystywane do detekcji obszarów zainteresowania (ang. *ROI – Region Of Interest*) znajdujących się na obrazach [2]. Ich analiza dokonywana jest między innymi w oparciu o cechy statystyczne obliczane na podstawie histogramu (ang. *histogram*), macierzy gradientu (ang. *gradient matrix*), macierzy długości pasm (ang. *run-length matrix*) oraz macierzy zdarzeń (ang. *co-occurrence matrix*) – zaproponowane przez Haralick'a [2].

4. Histogram

Histogram zawiera informacje o liczbie pikseli obrazu reprezentujących każdy z możliwych poziomów szarości, a po znormalizowaniu prawdopodobieństwo wystąpienia pikseli o danym poziomie szarości. Cechy wyznaczone na podstawie znormalizowanego histogramu:

- średnia (ang. *mean*):

$$\mu = \sum_{g_L=0}^{N_{g_L}-1} g_L * Hn_{g_L} \quad (1)$$

- wariancja (ang. *variance*):

$$\sigma^2 = \sum_{g_L=0}^{N_{g_L}-1} (g_L - \mu)^2 * Hn_{g_L} \quad (2)$$

- skośność (ang. *skewness*):

$$\mu_3 = \sigma^{-3} * \sum_{g_L=0}^{N_{g_L}-1} (g_L - \mu)^3 * Hn_{g_L} \quad (3)$$

- kurtozja (ang. *kurtosis*):

$$\mu_4 = \sigma^{-4} * \sum_{g_L=0}^{N_{g_L}-1} (g_L - \mu)^4 * Hn_{g_L} - 3 \quad (4)$$

gdzie: g_L – poziom szarości należący do zbioru $\langle 0, N_{g_L} - 1 \rangle$, Hn_{g_L} – wartość znormalizowanego histogramu dla g_L poziomu szarości.

5. Macierz gradientu

Macierz gradientu umożliwia analizę zmian poziomów szarości pikseli sąsiadujących z badanym pikselem, bądź znajdującym się w innej, zdefiniowanej odległości (w przeprowadzonych badaniach brano pod uwagę piksele sąsiednie). Gradient w punkcie o współrzędnych (i, j) należących do obszaru zainteresowania ROI, dla odległości d od badanego pikseli, wyrażony jest następującą zależnością:

$$G_{i,j} = \sqrt{(x_{i+d,j} - x_{i-d,j})^2 + (x_{i,j+d} - x_{i,j-d})^2} \quad (5)$$

Cechy wyznaczone na podstawie macierzy gradientu:

- średnia gradientu (ang. *gradient mean*):

$$G\mu = \frac{1}{M} \sum_{i,j \in ROI} G_{i,j} \quad (6)$$

- wariancja gradientu (ang. *gradient variance*):

$$(G\mu)^2 = \frac{1}{M} \sum_{i,j \in ROI} (G_{i,j} - G\mu)^2 \quad (7)$$

- skośność gradientu (ang. *gradient skewness*):

$$G\mu_3 = \frac{1}{M * (G\sigma)^3} \sum_{i,j \in ROI} (G_{i,j} - G\mu)^3 \quad (8)$$

- kurtozja gradientu (ang. *gradient kurtosis*):

$$G\mu_4 = \frac{1}{M * (G\sigma)^4} \sum_{i,j \in ROI} (G_{i,j} - G\mu)^4 - 3 \quad (9)$$

gdzie: M – liczba pikseli należących do ROI.

6. Macierz długości pasm

Macierz długości pasm przedstawia zależność między długością pasm występujących na obrazie a ich poziomem

szarości i zawiera informacje o ich liczbie. Obliczenia wykonywane są dla 4 macierzy wyznaczonych dla kierunków – 0° , 45° , 90° i 135° . Cechy wyznaczone na podstawie macierzy długości pasm:

- odwrotny moment uwydatnienia krótkich pasm (ang. *short run emphasis inverse moment*):

$$SRE = (\sum_{g_L=0}^{N_{g_L}-1} \sum_{l=1}^{N_r} \frac{R_{g_L,l}}{l^2}) / c \quad (10)$$

- moment uwydatnienia długich pasm (ang. *long run emphasis moment*):

$$LRE = (\sum_{g_L=0}^{N_{g_L}-1} \sum_{l=1}^{N_r} l^2 * R_{g_L,l}) / c \quad (11)$$

- niejednorodność poziomu szarości (ang. *grey level nonuniformity*):

$$GLN = (\sum_{g_L=0}^{N_{g_L}-1} (\sum_{l=1}^{N_r} R_{g_L,l})^2) / c \quad (12)$$

- niejednorodność pasm (ang. *run length nonuniformity*):

$$RLN = (\sum_{l=1}^{N_r} (\sum_{g_L=0}^{N_{g_L}-1} R_{g_L,l})^2) / c \quad (13)$$

- część obrazu w pasmach (ang. *fraction of image in runs*):

$$F = c / (\sum_{g_L=0}^{N_{g_L}-1} \sum_{l=1}^{N_r} l * R_{g_L,l}) \quad (14)$$

gdzie: g_L – poziom szarości należący do zbioru $\langle 0, N_{g_L} - 1 \rangle$, l – długość pasma należąca do zbioru $\langle 1, N_r \rangle$, $R_{g_L,l}$ – wartość z macierzy długości pasm wyrażająca liczbę wystąpień pasm o długości l i odcieniu szarości g_L , c – współczynnik zdefiniowany następująco:

$$c = \sum_{g_L=0}^{N_{g_L}-1} \sum_{l=1}^{N_r} R_{g_L,l} \quad (15)$$

7. Macierz zdarzeń

Macierz zdarzeń, inaczej histogram drugiego rzędu, wyraża ilościową zależność między odcieniami szarości pary pikseli znajdujących się w zdefiniowanej odległości od siebie. Obliczenia wykonywane są dla 4 macierzy wyznaczonych dla kierunków – 0° , 45° , 90° i 135° , dla odległości 1. Cechy wyznaczone na podstawie macierzy zdarzeń:

- drugi moment kątowy (ang. *angular second moment*):

$$ASM = \sum_{a,b} (C_{a,b})^2 \quad (16)$$

- kontrast (ang. *contrast*):

$$CON = \sum_{a,b} ((a - b)^2 * C_{a,b}) \quad (17)$$

- korelacja (ang. *correlation*):

$$CORR = - \sum_{a,b} \frac{(a-\mu)*(b-\mu)}{\sigma^2} * C_{a,b} \quad (18)$$

- wariancja (ang. *variance*):

$$V = \sum_{a,b} (a - \mu)^2 * C_{a,b} \quad (19)$$

- jednorodność (ang. *homogeneity*):

$$HOM = \sum_{a,b} \frac{C_{a,b}}{1 + (a-b)^2} \quad (20)$$

- sumaryczna średnia (ang. *sum average*):

$$SA = \sum_{k=2}^{2*L_g} k * P_{x+y}(k) \quad (21)$$

- sumaryczna wariancja (ang. *sum variance*):

$$SV = \sum_{k=2}^{2*L_g} (k - SA)^2 * P_{x+y}(k) \quad (22)$$

- sumaryczna entropia (ang. *sum entropy*):

$$SE = - \sum_{k=2}^{2*L_g} P_{x+y}(k) * \log(P_{x+y}(k)) \quad (23)$$

- entropia (ang. *entropy*):

$$E = - \sum_{a,b} C_{a,b} * \log(C_{a,b}) \quad (24)$$

- wariancja różnicowa (ang. *difference variance*):

$$DV = \sum_{k=0}^{L_g-1} k^2 * P_{x-y}(k) \quad (25)$$

- entropia różnicowa (ang. *difference entropy*):

$$DE = - \sum_{k=0}^{L_g-1} P_{x-y}(k) * \log(P_{x-y}(k)) \quad (26)$$

8. Wykorzystanie tekstur w analizie obrazów mammograficznych

W celu przeprowadzenia eksperymentów przygotowano zbiór próbek z obszarów niezmiennych oraz zmienionych chorobowo. Poniżej przedstawiono przykładowe próbki, wszystkie pochodzą z obrazów mammograficznych zawartych w bazie MIAS [11].



Rys. 3. Przykłady próbek z obszarów niezmiennych chorobowo (z obrazów odpowiednio nr 15, 148, 181, 184 z bazy MIAS [11] – od lewej strony)



Rys. 4. Przykłady próbek z obszarów zmienionych chorobowo (z obrazów odpowiednio nr 148, 181, 184, 193 z bazy MIAS [11] – od lewej strony)

9. Klasyfikacja

Na potrzeby etapu klasyfikacji wybrano jeden z najprostszych klasyfikatorów – klasyfikator k -najbliższych sąsiadów (k -NN, ang. *k-nearest neighbours*). Aby możliwe było wykorzystanie klasyfikatora do wskazania klasy, do której należy testowa próbka, konieczne jest jego wytrenowanie. W tym celu z przygotowanego zestawu 130 próbek o wymiarach 30x30 pikseli, wybrano zbiór treningowy składający się ze 100 próbek (w tym 55 próbek z obszarów niezmiennych oraz 45 próbek z obszarów zmienionych chorobowo) oraz zbiór testowy składający się z 30 próbek. Należy podkreślić, iż żadna z próbek ze zbioru testowego nie pochodziła z tych samych mammogramów co próbki tworzące zbiór treningowy.

Zbiór treningowy zbudowany został ze 100 wektorów cech, z których każdy reprezentował 72 cechy wyekstrahowane z odpowiadającej mu próbki. Każdemu wektorowi przypisano klasę, do której należy – w tym przypadku jedną z dwóch: zmieniony chorobowo, niezmienny chorobowo. Przypisanie do odpowiedniej klasy dokonano w oparciu o dane zawarte w bazie MIAS – zmiana chorobowa opisana jest poprzez podanie jej centrum oraz promienia okręgu, który ją ogranicza.

Kolejnym krokiem jest zaklasyfikowanie każdej próbki ze zbioru testowego do jednej z dwóch możliwych klas. W tym celu podobnie jak na etapie uczenia klasyfikatora, dla każdej z próbek utworzono reprezentujący ją wektor 72 cech, jednak odpowiadająca mu klasa wykorzystywana jest jedynie w późniejszym etapie oceny skuteczności klasyfikacji.

W celu dokonania klasyfikacji wektora w_T obliczana jest jego odległość od wszystkich wektorów ze zbioru treningowego, następnie wybieranych jest k -najbliższych mu wektorów. Wektorowi testowemu w_T przypisywana jest klasa, do której należy większość spośród k -najbliższych wektorów [8].

Za poprawną klasyfikację uznawane jest takie wskazanie klasy przez klasyfikator, które odpowiada klasie określonej dla danego wektora w zbiorze testowym.

10. Testy

Testy przeprowadzono dla zbioru składającego się z 30 próbek, w tym 15 próbek z obszarów niezmiennych oraz 15 próbek z obszarów zmienionych chorobowo. W wyniku klasyfikacji (zbiór wszystkich cech) z wykorzystaniem wytrenowanego wcześniej klasyfikatora k -NN (dla 5 najbliższych sąsiadów) poprawnie rozpoznano 26 próbek, w tym:

- TP (ang. *true positive*) – 15,
- TN (ang. *true negative*) – 11,
- FP (ang. *false positive*) – 4,

- FN (ang. *false negative*) – 0.

Szczegółowe wyniki klasyfikacji przedstawiają poniższe miary:

- dokładność (ang. *accuracy*):

$$AC = \frac{TP+TN}{TP+TN+FP+FN} \quad (27)$$

- precyzja (ang. *precision*):

$$P = \frac{TP}{TP+FP} \quad (28)$$

- czułość (ang. *sensitivity, recall*):

$$S = \frac{TP}{TP+FN} \quad (29)$$

Poniżej przedstawiono podsumowanie wyników klasyfikacji:

Tabela 1. Wyniki klasyfikacji dla całego zbioru cech

Liczba cech	TP	FP	TN	FN	Dokładność	Precyzja	Czułość
72	15	4	11	0	0,87	0,79	1

Otrzymane wyniki eksperymentu – czułość na poziomie 1,0 oznacza, iż żadna ze zmienionych chorobowo próbek nie została błędnie sklasyfikowana, co jest istotne w rozpatrywanym zagadnieniu. Wartości dokładności i precyzji, odpowiednio 0,87 i 0,79 wynikają z niepoprawnego zaklasyfikowania 4 niezmiennych chorobowo próbek jako zmienione chorobowo. Ze względu na specyfikę przedmiotu badań są dużo mniej znaczące w porównaniu z pominięciem obszaru zmienionego chorobowo. Powoduje to wskazanie większej liczby rejonów zainteresowania, które poddawane są ocenie lekarza, co nie ma negatywnego wpływu na stawianą diagnozę.

Redukcja liczby cech pozwala na uzyskanie wyników przedstawionych poniżej:

Tabela 2. Wyniki klasyfikacji po redukcji cech

Liczba cech	TP	FP	TN	FN	Dokładność	Precyzja	Czułość
65, 60, 55, 50	14	4	11	1	0,83	0,78	0,93
45	14	5	10	1	0,8	0,74	0,93
40, 35, 30	13	5	10	2	0,77	0,72	0,87

Zmniejszenie liczby cech wpływa na pogorszenie wyników klasyfikacji. Błędy klasyfikacji mogą wynikać również ze zbyt małego zbioru treningowego.

11. Dyskusja wyników i podsumowanie

Wyniki przeprowadzonego eksperymentu wskazują na dużą skuteczność klasyfikacji próbek z obszarów zmienionych chorobowo, dokonanej na podstawie cech statystycznych obliczanych z tekstury, a zatem możliwość wykorzystania wspomnianych cech do klasyfikacji obszarów zainteresowania, wykrytych we wcześniejszym etapie przetwarzania, jako obszarów niezmiennych lub zmienionych chorobowo.

W literaturze opublikowano wyniki badań, których autorzy zastosowali różne kombinacje cech do klasyfikacji zmian chorobowych na obrazach mammograficznych, m.in. różne cechy teksturowe, cechy opisujące kształt itd. Przykładowo autorzy pracy [5] obliczali cechy na podstawie histogramu oraz kilku macierzy zdarzeń (dla odległości od 5 do 30 między pikselami), uzyskując poprawną klasyfikację dla 84% przypadków. W artykule [10] zaprezentowano klasyfikację w oparciu o cechy wyznaczone z macierzy zdarzeń i cechy Tamura – podana przez autorów uzyskana średnia precyzja wynosi 69%.

Porównanie otrzymanych wyników ze znanymi z literatury nie jest całkowicie miarodajne ze względu na to, iż testy nie zostały przeprowadzone na tym samym zbiorze danych oraz nie zawsze podane są miary TP, TN, FP, FN pozwalające na dokładniejszą ocenę rezultatów.

W dalszym etapie badań wyznaczone zostaną dodatkowe cechy, które mogą zwiększyć skuteczność klasyfikacji. Kolejno przeprowadzone zostaną testy dla wszystkich obszarów znalezionych z wykorzystaniem metody opisanej w artykule [4] dla obrazów pochodzących z bazy MIAS [11].

Literatura

- [1] Guliato D., Rangayyan R., Carnielli W., Zuo J., Desautels J.: Segmentation of breast tumors in mammograms by fuzzy region growing, Engineering in Medicine and Biology Society, Proceedings of the 20th Annual International Conference of the IEEE, 1998, vol. 2, pp. 1002-1005.
- [2] Haralick R.M., Shanmugam K., Dinstein I.: Textural Features for Image Classification, IEEE Transactions on Systems, Man and Cybernetics, Vol. SMC-3, 1973, No. 6, pp. 610-621.
- [3] Huo Z., Giger M., Vyborny C., Metz C.: Breast cancer: Effectiveness of computer-aided diagnosis – observer study with independent database of mammograms1, Radiology, 2002, vol. 224, no. 2, pp. 560-568.
- [4] Lazarek J., Szczepaniak P.S., Tomczyk A.: Method of Pattern Detection in Mammographic Images, Intelligent Systems in Technical and Medical Diagnosis, Eds. Józef Korbicz, Marek Kowal. Springer, 2014, pp. 235-245.
- [5] Lyra M., Lyra S., Kostakis B., Drosos S., Georgosopoulus C., Skouroliakou K.: Digital mammography texture analysis by computer assisted image processing, IEEE International Workshop on Imaging Systems and Techniques – IST 2008 Chania, Greece, September 10–12, 2008.
- [6] Rangayyan R. M., Ayres F. J., Desautels J. L.: A review of computer-aided diagnosis of breast cancer: Toward the detection of subtle signs, Journal of the Franklin Institute, vol. 344, no. 34, 2007 pp. 312-348.
- [7] Sampat M., Markey M., Bovik A.: Computer-aided detection and diagnosis in mammography, Handbook of Image and Video Processing, vol. 2, 2005, pp. 1195-1217.

- [8] Szczepaniak P.S.: Obliczenia inteligentne, szybkie przekształcenia i klasyfikatory, Akademia Oficyna Wydawnicza EXIT, Warszawa, 2004.
- [9] Szczepaniak P.S., Tadeusiewicz R.: The role of artificial intelligence, knowledge and wisdom in automatic image understanding, Journal of Applied Computer Science – JACS, 18, 2010, No.1, pp. 75-85.
- [10] Zhou J., Feng Ch., Liu X., Tang J.: A Texture Features based Medical Image Retrieval System for Breast Cancer, 2012 7th International Conference on Computing and Convergence Technology (ICCT), IEEE, pp. 1010 – 1015.
- [11] „Mias dataset”. <http://peipa.essex.ac.uk/info/mias.html>. MIAS dataset.

Mgr inż. Jagoda Lazarek

e-mail: jagoda.lazarek@p.lodz.pl



Mgr inż. Jagoda Lazarek jest doktorantką w Instytucie Informatyki (Wydział Fizyki Technicznej, Informatyki i Matematyki Stosowanej), w którym zatrudniona jest także na stanowisku asystenta.

Jej zainteresowania naukowe związane są z metodami przetwarzania i rozpoznawania obrazów oraz metodami sztucznej inteligencji. Do jej ważnych osiągnięć należą między innymi wyróżnienia srebrnym i brązowym medalem rozwiązania „System do wykrywania znaków drogowych” (J. Lazarek, P.S. Szczepaniak), prezentowanego podczas międzynarodowych targów wynalazków i innowacji „IENA” oraz „IWIS”. W 2012 roku została finalistką międzynarodowego konkursu Great Minds organizowanego przez IBM i odbyła staż w laboratorium badawczym IBM w Izraelu (Hajfa).

otrzymano/received: 27.06.2013

przyjęto do druku/accepted: 04.09.2013



Dwusemestralne studia podyplomowe obejmują moduły kształcenia:

- | | |
|--|--|
| <ol style="list-style-type: none"> 1 Podstawy administrowania systemem Windows 2 Podstawy administrowania systemem Linux 3 Sieci LAN 4 Sieci IP 5 Routing w sieciach IP 6 Sieci TCP/IP 7 Usługi w sieciach TCP/IP 8 Wykorzystanie Windows w sieciach LAN 9 Wykorzystanie Linux w sieciach LAN 10 Zarządzanie systemami informatycznymi 11 Sieci telekomunikacyjne | <ol style="list-style-type: none"> 12 Serwery usług sieciowych 13 Bezpieczeństwo systemów komputerowych 14 Diagnostyka sieci komputerowych 15 Multimedia strumieniowe 16 Systemy telefonii VoIP 17 Sieci dostępne 18 Kodowanie informacji multimedialnych 19 Elementy sieci optoelektronicznych 20 Światłowody i kable światłowodowe 21 Podstawy normalizacji 22 Seminarium dyplomowe |
|--|--|

Telefon:

(0-81) 53 84 309 – sekretariat Instytutu Elektroniki i Technik Informatycznych

(0-81) 53 84 317 – opiekun ds. organizacyjnych, dr inż. Piotr Kisała

Faks:

(0-81) 53 84 312

Poczta elektroniczna:

ask@politechnika.lublin.pl

Adres pocztowy:

ul. Nadbystrzycka 38A, 20-618 Lublin