

Київський національний університет імені Тараса Шевченка
факультет комп'ютерних наук та кібернетики

Випускна кваліфікаційна робота бакалавра

на тему:

“Розпізнавання раку молочної залози за допомогою
аналізу інтерфазних ядер bukalного епітелію”

Студента 4 курсу
кафедри обчислювальної математики
Чан Ха Ву

Науковий керівник
професор, доктор фізико-математичних наук
Клюшин Дмитро Анатолійович

Робота заслухана на засіданні кафедри обчислювальної математики та
рекомендована до захисту в ДЕК, протокол № 11 від 18 травня 2017 р.

Завідувач кафедри обчислювальної математики

проф. Ляшко С.І.

Київ 2017

Анотація

У зв'язку зі стрімким зростанням кількості випадків захворювань раком молочної залози, проблема неінвазивної діагностики цієї хвороби є дуже актуальною. Новітні дослідження в області онкології [1] виявили закономірність, яка полягає в тому, що існує зв'язок між змінами у інтерфазному ядрі bukalного епітелію та розвитком злюкісної або доброкісної пухлини (раком молочної залози та фіброаденоматозом відповідно). В роботі досліджується цей зв'язок за допомогою технік комп'ютерного зору, фрактальної геометрії та машинного навчання.

Ця робота складається з трьох частин. У першій частині пропонується алгоритм сегментації інтерфазних ядер bukalного епітелію на зображені. У другій частині використовується фрактальний аналіз інтерфазних ядер bukalного епітелію, щоб відрізнити здорових пацієнтів від хворих раком молочної залози або фіброаденоматозом. У третьій частині використовуються методи глибинного навчання для класифікації зразків на групи хворих раком молочної залози та хворих фіброаденоматозом.

Подяки

Хочу виразити щиру подяку своєму научному керівнику, професору та доктору ф-м наук Клюшину Дмитрові Анатолійовичу, хто надав мені можливість займатися цим неймовірно цікавим дослідженням.

Моїй матері, яка успішно боролася з раком молочної залози у 2011-му році та завжди мене підтримувала в навчанні.

Та Анні Разумовій, що допомогла мені в оформленні тексту цієї роботи.

Зміст

Анотація	i
Подяки	ii
1 Вступ	1
1.1 Метод дослідження	1
1.2 Опис даних	1
2 Попередня обробка даних	2
2.1 Сегментація ядер	2
2.1.1 Зменшення шуму. Медіанна фільтрація.	3
2.1.2 Відділення фону. Бінарізація Отсу.	4
2.1.3 Бінарні морфологічні операції над зображенням.	5
2.1.4 Трансформація дистанції	7
2.1.5 Алгоритм водорозділу	8
2.1.6 Кінцевий алгоритм сегментації	9
2.2 Нормалізація зображень	10
2.3 Фільтр на основі еліпсоїдів Петуніна	11
2.3.1 Еліпсоїди Петуніна	11
2.3.2 Фільтр зображень на основі еліпсоїдів Петуніна	12
2.3.3 Швидкий алгоритм на основі динамічного програмування	12
3 Оцінка якості даних	15
3.1 Алгоритм РСА для зменшення розмірності	15
3.2 Класифікатор SVM	15
3.3 Результати оцінки якості	15
4 Фрактальний аналіз	16
4.1 Фрактальна характеристика	16
4.1.1 Теоретичні відомості	16
4.1.2 Алгоритм box-counting	17
4.1.3 Лінійна регресія	17
4.2 Класифікація пацієнтів	18
4.2.1 Міра близькості та p -статистика	18
4.2.2 Метод к найближчих сусідів. Кросвалідація.	19
4.3 Результати та спостереження	19
5 Глибинне навчання	20
5.1 Нейронні мережі	20
5.1.1 Багатошаровий перцептрон	21
5.1.2 Функція активації ReLU	21
5.1.3 Згорткові нейронні мережі	21
5.1.4 Алгоритм RMSProp	22

5.1.5	Регуляризація нейронної мережі	23
5.1.6	Функція втрати. Перехрестна ентропія.	23
5.2	Класифікація за допомогою нейронних мереж	24
5.2.1	Прирошення даних	24
5.2.2	Проблема незбалансованості вибірки для тренування	24
5.2.3	Архітектура мережі	25
5.2.4	Тренування мереж та результати	26
5.2.5	Аналіз результатів	27
5.3	Підвищення якості розпізнавання	27
5.3.1	Класифікатор на основі нейронної мережі та критерії еквівалентності	27
5.3.2	Аналіз результатів	28
6	Висновок	30
A	Програмний пакет	31
	Bibliography	32

Розділ 1

Вступ

Незважаючи на стрімкий розвиток сучасної онкології, захворюваність на рак молочної залози як в Україні, так і в більшості розвинених країн світу продовжує зростати. Щороку у світі реєструють більше п'яти тисяч нових випадків захворювання на рак молочної залози, що становить понад 25% всіх ракових захворювань у жінок. Тому задача неінвазивної попередньої діагностики раку молочної залози є надзвичайно актуальною.

Обробка цифрових зображень давно стала складовою досліджень практично у всіх галузях науки. Робота зі зразками з будь-якої предметної області має на увазі не тільки витяг даних із зображень, але і класифікацію знімків, роботу зі складноструктураними зразками, з неочевидними закономірностями і особливостями, часто помітними лише фахівцям у цій галузі. У медицині можливість автоматично обробляти та розпізнавати великі набори знімків мікроскопа певної тематики може значно вплинути на хід досліджень, полегшити процес роботи із зображеннями та у результаті, наприклад, прискорити виявлення хвороби та постановку діагнозу, що допоможе підібрати своєчасне і необхідне лікування.

1.1 Метод дослідження

Як показано у попередніх роботах в області онкології [1], зложісні та доброкісні утворення певним чином впливають на інтерфазні ядра bucalного епітелію. Отже, ми сподіваємося що використовуючи методи комп’ютерного зору, фрактального аналізу, та машинного навчання, можна діагнозувати рак молочної залози та фіброаденоматоз з певною точністю.

1.2 Опис даних

Для дослідження було використано набір зображень інтерфазних ядер bucalного епітелію, отриманих від здорових пацієнтів, а також хворих раком молочної залози та фіброаденоматозом, які лікувались у інституті онкології.

Розділ 2

Попередня обробка даних

2.1 Сегментація ядер

Мікроскопічні зображення у чистому вигляді зазвичай неприйнятні для аналізу у зв'язку з наявністю дефектів, цифрових шумів, спричинених необхідністю підвищення світлоочутливості фотоматеріалів, а також сторонніх об'єктів. Отже, для подальшого аналізу вхідного набору зображень (або набору пацієнтів), необхідно зробити попередню обробку.

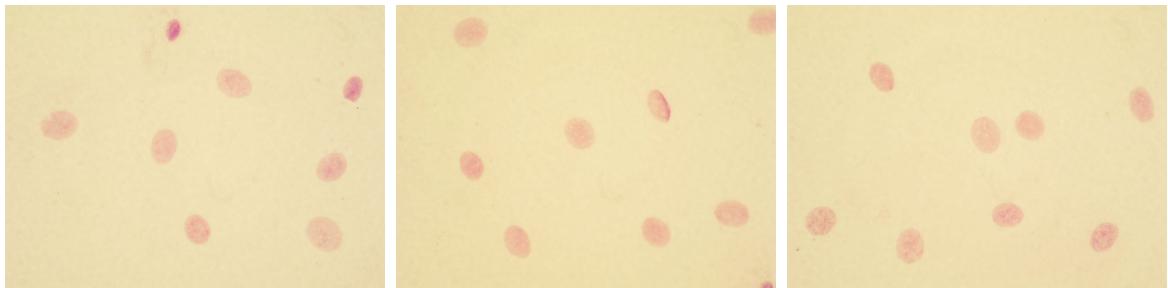
На (2.1) зображені приклади знімків інтерфазних ядер букального епітелію. З них треба виділити окрім клітини для подальшого аналізу. Задача обробки та сегментації зображень розбивається на наступні підзадачі:

1. Зменшення шуму
2. Відділення фону
3. Виділення окремого ядра клітини

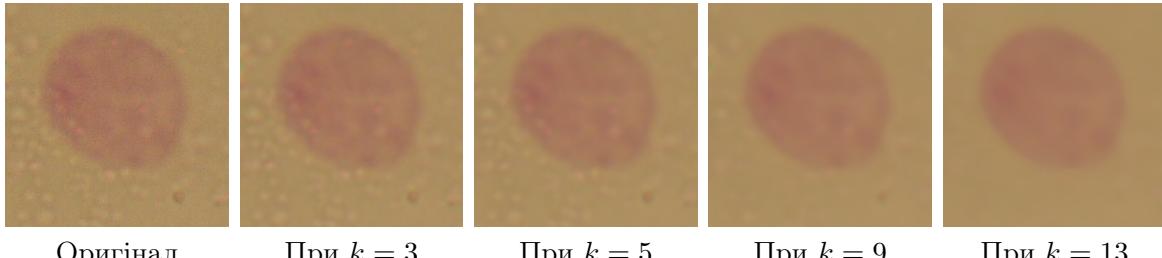
Перша та друга підзадачі стандартно розв'язуються за допомогою алгоритмів фільтрації зображень, а друга підзадача розв'язується за допомогою алгоритму бінаризації Оцу.

У загальному випадку, якщо на зображенні присутні різні об'єкти (у даному випадку – ядра різних типів клітин чи різні клітини), третя підзадача дуже складно розв'язується детермінованими алгоритмами, оскільки об'єкти можуть перекриватися, торкатися один одного та мати абсолютно різні форми, кольори та розміри. Так, наприклад, [2] використав нейронні мережі для класифікації та сегментації зображень клітин.

У нашому випадку, на досліджуваних зображеннях присутні тільки інтерфазні ядра букального епітелію. Отже, ми припускаємо, що радіуси цих ядер приблизно однакові та



Мал. 2.1: Зображення інтерфазних ядер букального епітелію, яких треба сегментувати.



Мал. 2.2: Дія медіанного фільтру на зображення.

можемо використовувати алгоритми на основі морфологічних операцій, трансформації відстані та алгоритму водоподілу, які будуть детально розглянуті далі.

2.1.1 Зменшення шуму. Медіанна фільтрація.

Для зменшення рівня шуму знімків мікроскопа будемо використовувати медіанний фільтр – один з видів цифрових фільтрів, запропонований [3], що широко застосовується в цифровій обробці сигналів та зображень. Його використання дає найкращі результати для збереження перепадів відтінків, різноманітних кордонів та локальних піків яскравості на спотворених імпульсним шумом зображеннях. Швидкий алгоритм медіанної фільтрації виглядає наступним чином:

Алгоритм 1: Медіанна фільтрація

Вхід : Зображення X розміром $m \times n$, радіус ядра (вікна фільтру) r
 Вихід: Зображення Y такого ж розміру

```

Ініціалізація гістограму ядра  $H$ ;
for  $i = 1$  to  $m$  do
    for  $j = 1$  to  $n$  do
        for  $k = -r$  to  $r$  do
            Вилучити  $X_{i+k,j-r-1}$  від  $H$ ;
            Додати  $X_{i+k,j+r}$  до  $H$ ;
         $Y_{i,j} \leftarrow$  медіана( $H$ );
    
```

На (Мал. 2.2) показано, як змінюється рівень шуму та чіткість зображення в залежності від розміру вікна фільтру $k = 2r + 1$. Можна побачити, що при $k = 3$, зображення мінімально позбавлене шуму, а текстура ядра майже повністю зберігається. При $k = 13$ дуже гарно виділяються контури ядра клітини та зображення зовсім позбавлене шуму, але майже повністю втрачено текстуру ядра.

Отже, на різних етапах аналізу можна використовувати медіанні фільтри з різними розмірами вікна фільтру. Наприклад, при сегментації зображень (Мал. 2.1) доречно використовувати вікна фільтру з великими розмірами, оскільки на цьому етапі нас більше цікавлять контури об'єктів, а при класифікації (здоровий, хворий раком або фіброаденоматозом) краще використовувати якомога менші вікна фільтру, оскільки нас цікавить текстура інтерфазного ядра клітини bukalного епітелію.

2.1.2 Відділення фону. Бінарізація Отсу.

Наступний етап обробки зображення – відділення фонових пікселів від пікселів об'єктів, наявних на зображенні (ядер клітин або сторонніх об'єктів). Для досягнення цієї мети підходить алгоритм бінаризації зображення Оцу [4].

Нехай пікселі зображення приймають значення у L рівнях сірого кольору, тобто значення з множини $\{1, 2, \dots, L\}$, кількість пікселів зі значенням i дорівнює n_i , а кількість пікселів у зображенні – N . Треба розділити множину пікселів на два класи C_0 та C_1 (фонові пікселі та пікселі об'єктів, або навпаки) за допомогою деякого значення порогу k . До C_0 віднесемо усі пікселі зі значеннями з $\{1, 2, \dots, k\}$, а до C_1 – усі пікселі зі значеннями з множини $\{k + 1, k + 2, \dots, L\}$. Нехай ω_0 та ω_1 – частота класів C_0 та C_1 відповідно.

$$\omega_0 = \omega_0(k) = \sum_{i=0}^k \frac{n_i}{N}, \quad \omega_1 = \omega_1(k) = \sum_{i=k+1}^L \frac{n_i}{N} = 1 - \omega_0$$

оді метод Оцу полягає у пошуку порогу k , який мінімізує дисперсію в середині класу, яка визначається як зважена сума дисперсій двох класів:

$$\sigma_\omega^2(k) = \sigma_0^2(k)\omega_0(k) + \sigma_1^2(k)\omega_1(k) \rightarrow \min$$

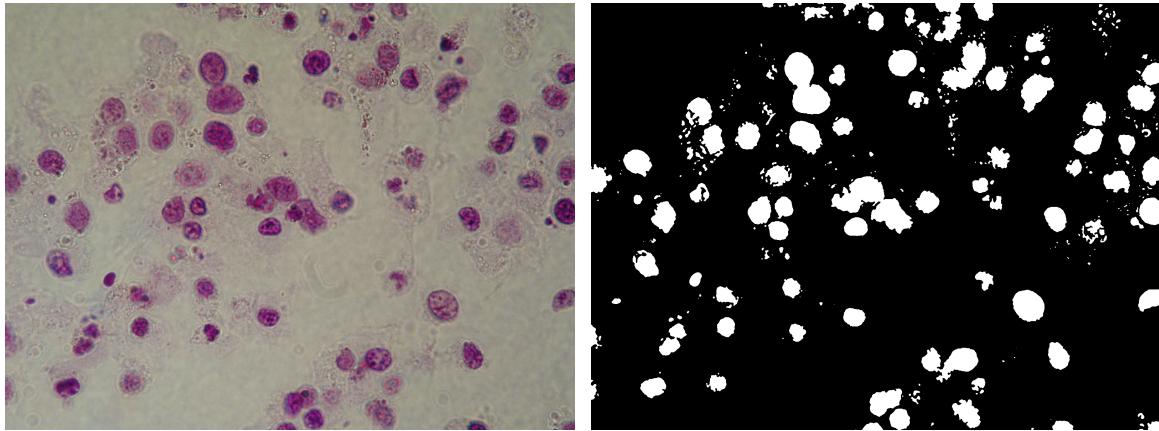
Де $\sigma_0^2(k)$ та $\sigma_1^2(k)$ – дисперсія класів C_0 та C_1 для порогу k відповідно. У своїй роботі Оцу також показав, що мінімізація дисперсії в середині класу рівносильна максимізації дисперсії між класами:

$$\sigma_b^2(k) = \omega_0(k)\omega_1(k)(\mu_0(k) - \mu_1(k))^2 \rightarrow \max$$

Де $\mu_0(k)$ та $\mu_1(k)$ – середнє арифметичне класів C_0 та C_1 при порозі k відповідно, тобто:

$$\begin{aligned} \mu_T &= \frac{\sum_{i=0}^L i \cdot n_i}{N} \\ \mu_0(k) &= \frac{\sum_{i=0}^{k-1} i \cdot n_i}{N \cdot \omega_0(k)}, \quad \mu_1(k) = \frac{\mu_T - \mu_0(k) \cdot \omega_0(k)}{\omega_1(k)}. \end{aligned}$$

Значення $\omega_0(k+1)$, $\omega_1(k+1)$, $\mu_0(k+1)$ та $\mu_1(k+1)$ достатньо легко виражаються через значення $\omega_0(k)$, $\omega_1(k)$, $\mu_0(k)$ та $\mu_1(k)$, що дозволяє швидко обчислити оптимальний поріг



Мал. 2.3: Результат бінарізації.

k. Отже, алгоритм бінаризації Оцу виглядає наступним чином:

Алгоритм 2: Бінарізація зображення Отсу

Вхід : Зображення X розміром $m \times n$

Вихід: Поріг k

обчислити гістограму n_i зображення та частоту $N(i) = \frac{n_i}{N}$ для кожного рівня інтенсивності зображення X ;

обчислити початкові значення $\omega_0(0)$, $\omega_1(0)$, $\mu_0(0)$ та $\mu_1(0)$;

поставити $\sigma_{\max}^2 \leftarrow 0$;

for $k = 1$ to L do

Оновити значення $\omega_0(k)$, $\omega_1(k)$, $\mu_0(k)$ та $\mu_1(k)$;

Обчислити $\sigma_b^2(k) = \omega_0(k)\omega_1(k)(\mu_0(k) - \mu_1(k))^2$;

if $\sigma_b^2(k) > \sigma_{\max}^2$ then

поставити $\sigma_{\max}^2 \leftarrow \sigma_b^2(k)$

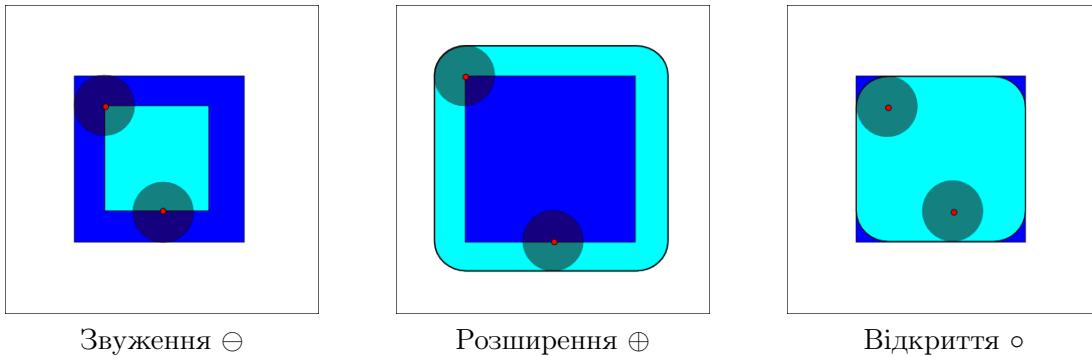
поставити $k \leftarrow \sigma_{\max}^2$;

На зображенні (Мал. 2.3) результат бінаризації фільтрованого медіанним фільтром зображення з розміром $k = 5$. Негладкі контури інтерфазних ядер клітин після бінаризації, а також неточності у деяких регіонах зображення нас не турбують на даному етапі, оскільки це можна виправити морфологічними операціями, які будуть детально розглянуті далі.

2.1.3 Бінарні морфологічні операції над зображенням.

У бінарній морфології, зображення розглядається як підмножина Евклідового простору \mathbb{R}^d або цілочисельної сітки \mathbb{Z}^d , де d – розмірність простору. Операції морфології застосовуються до зображень разом з заданим структурним елементом, або ядром, яке визначає природу операції. Для обробки зображень найчастіше використовують наступний структурний елемент:

$$B = \{(-1, -1), (-1, 0), (-1, 1), (0, -1), (0, 0), (0, 1), (1, -1), (1, 0), (1, 1)\}.$$



Мал. 2.4: Морфологічні операції з круглим структурним елементом над синім квадратом.

Двома основними морфологічними операціями є Erosion(звуження) і Dilation(розширення). Операція звуження зображення задається наступним чином:

$$A \ominus B = \left\{ z \in \mathbb{Z}^d \mid B_z \subseteq A \right\},$$

де A – зображення, B – структурний елемент, B_z – трансляція вектору z , тобто

$$B_z = \{b + z \mid b \in B\}, \quad \forall z \in \mathbb{Z}^d$$

Ідея цієї операції полягає в наступному. Ядро (структурний елемент) ковзає по зображеню. Піксель в бінаризованому зображені (1 або 0) буде вважатися 1, тільки якщо всі пікселі під ядром рівні 1, в іншому випадку він буде зруйнований (стане рівним нулю). Таким чином, в результаті операції всі пікселі, що знаходяться на краю, будуть відкинуті в залежності від розміру ядра. Таким чином, товщина або розмір об'єкта переднього плану зменшується. Цей метод може використовуватися для видалення невеликих шумів, роз'єднання двох пов'язаних об'єктів та іншого.

Операція розширення зображення задається наступним чином:

$$A \oplus B = \bigcup_{b \in B} A_b.$$

Ця операція протилежна еrozії. Тут піксель бінаризованого зображення дорівнює “1”, якщо принаймні один піксель під ядром дорівнює “1”. Таким чином, збільшується біла область зображення або збільшується розмір об'єкта переднього плану. Зазвичай у випадках, коли потрібно видалити шум навколо об'єкта, спочатку використовують операцію звуження, а потім операцію розширення. Розширення необхідне, оскільки еrozія, крім видалення білого шуму, скорочує наш об'єкт.

На основі основних операцій звуження (\ominus) та розширення (\oplus) задаються інші операції. Операція морфологічного відкриття (opening) задається наступним чином:

$$A \circ B = (A \ominus B) \oplus B.$$

або

$$A \circ B = \bigcup_{B_z \subseteq A} B_z.$$

У [5] більш детально розглянуто теоретичні та практичні аспекти операції математичної морфології.

Позначимо зображення, отримане у результаті розширення бінарного зображення за 2 ітерації у (2.5) як *zoneOfInterest*. Очевидно, що всі пікселі, які дорівнюють 0 у *sureBG* є пікселями фону. Отже, ядра треба шукати тільки серед тих точок *zoneOfInterest*, які дорівнюють 1 (позначені білим кольором). Позначимо також відкриття у (Мал. 2.5) як *threshOpening*.

2.1.4 Трансформація дистанції

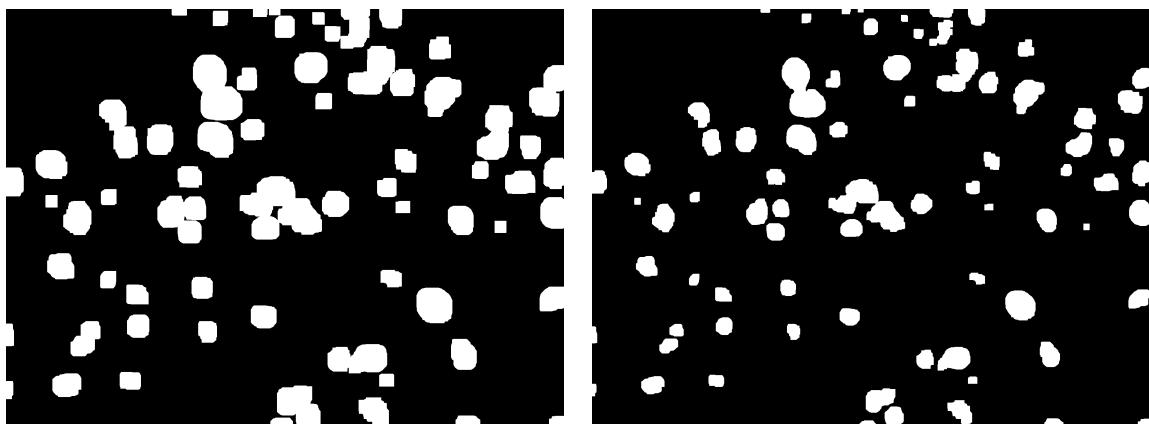
Трансформація дистанції (Distance Transformation) є також операцією математичної морфології. Кожному пікселю переднього плану (що належить E) бінарного зображення $E \subseteq \mathbb{Z}^d$ ставиться у відповідність відстань цієї точки до найближчого пікселя заднього плану (що не належить E). Для бінарних зображень найчастіше використовують Евклідову відстань, тобто для зображення $E \subseteq F = \{1, 2, \dots, height\} \times \{1, 2, \dots, width\}$, кожній точці $z = (z_x, z_y) \in E$ ставимо у відповідність:

$$DT(z) = \inf_{(x,y) \in F} (|x - z_x| + |y - z_y|)$$

Результат цієї операції для $E = threshOpening$ можна інтерпретувати як показник вірогідності, що піксель $z = (z_x, z_y)$ належить інтерфазному ядру клітини bukalного епітелію. Отже, можна задати емпіричне правило, за яким будемо приймати певні пікселі як “точно належать ядру клітини”:

$$DT(z) \geq \alpha \cdot \sup_{s \in E} DT(s) \Rightarrow "z \text{ точно належать ядру клітини}" . \quad (2.1)$$

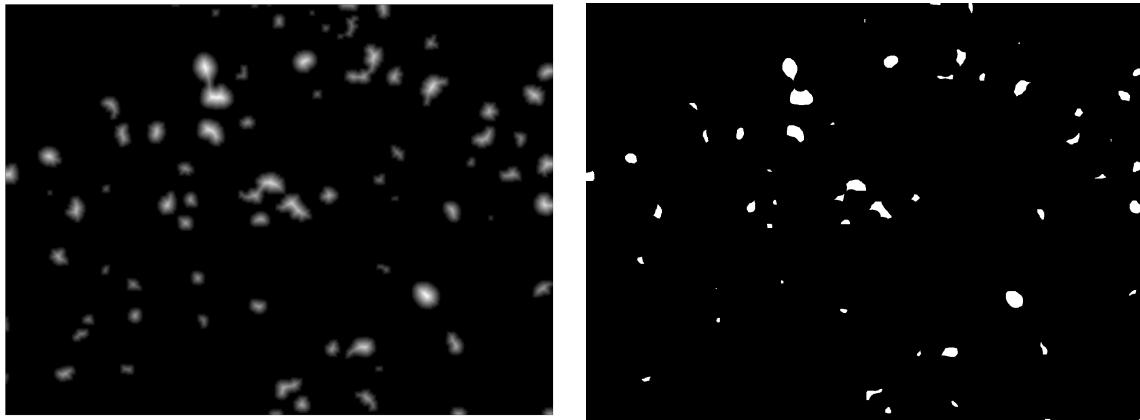
Позначимо множину, що задоволяє цю нерівність, як *sureFG*. Множник $\alpha = 0.475$ був обраний експериментальним чином, на досліджуваних даних. Тут також використовується припущення, що усі ядра клітин мають приблизно одинаковий розмір, а отже, центри цих клітин будуть знаходитись на приблизно однакому відстані від точок заднього плану. Для інших задач, дані яких зібрані за інших умов, слід обрати інший множник.



Розширене зображення за 2 ітерації

Відкриття зображення за 2 ітерації

Мал. 2.5: Морфологічні операції над бінарним зображенням (Мал. 2.3).



Трансформація дистанції

sureFG

Мал. 2.6: Трансформація дистанції.

У (2.6) показано результат трансформації відстані та обчислення *sureFG* над бінарним зображенням *threshOpening*, зображенім на (2.5). На *sureFG* ядра не дотикаються, завдяки обраному параметру α , тобто наш метод сегментації коректно обробляє ситуації, коли клітини торкаються чи частково перекриваються.

Використовуючи алгоритм що запропонованний у [6] який використовує так звану рекурсивну морфологію, можна обчислити трансформацію дистанції за два обходу. Для повного виявлення окремих клітин, використаємо алгоритм водоподілу.

2.1.5 Алгоритм водорозділу

Трансформація водоподілу (Watershed) розглядає одноканальне зображення як топографічну карту, де значення пікселя відображає висоту цієї точки.

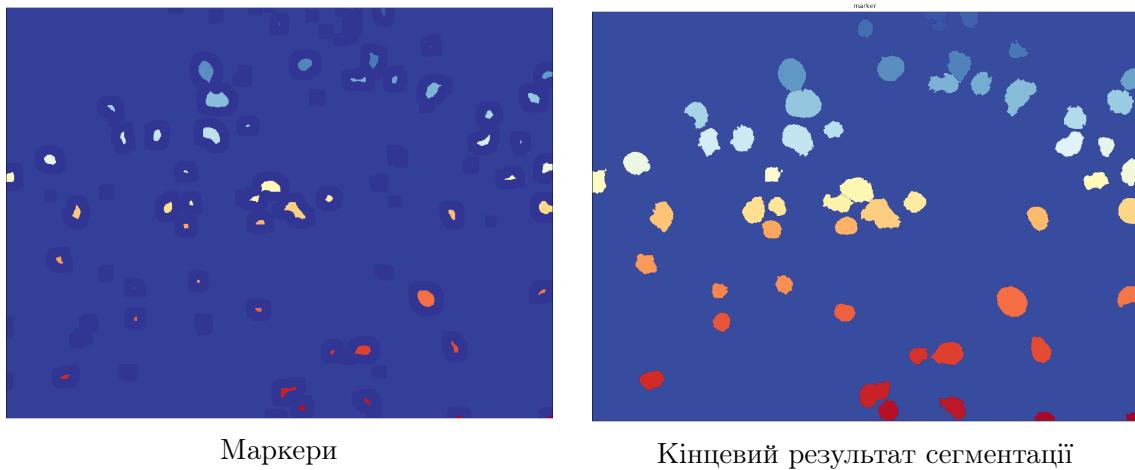
Починаючи з точок локальних мінімумів (або з заданих точок), кожна “яма” заповнюється рідиною різних кольорів. З підвищеннем рівня рідини, настає момент, коли дві рідини різного кольору починають змішуватися. Щоб уникнути цього, будуються границі між регіонами цих рідин. Границі будуються, доки рівень рідини не стане вищим за найвищу точку карти. Побудовані границі є результатами сегментації.

Використаємо цю ідею для сегментації ядер клітин. Кожну зв’язну компоненту *sureFG* пофарбуємо в певний “колір”. Назвемо таку сегментацію “набором маркерів”. Позначимо

$$\text{sureBG} = E \setminus \text{zoneOfInterest}$$

Зв’язну множину пікселів заднього плану *sureBG* також пофарбуємо в якийсь колір. Непофарбованою залишається така множина:

$$\text{unknownRegion} = \text{zoneOfInterest} \setminus \text{sureFG}$$



Мал. 2.7: Результат роботи алгоритму водорозділу (watershed).

Тоді алгоритм водоподілу з заданими маркерами наступний [7]:

Алгоритм 3: Алгоритм водоподілу Майера з заданими маркерами

Вхід : Зображення X розміром $m \times n$, набір маркерів

Вихід: Сегментація

0. Будемо вважати, що “затоплення” починається саме с пікселів маркерів.;
 1. Сусідні пікселі кожного маркеру додаємо в чергу з пріоритетом, де рівень пріоритету залежить від величини градієнту пікселя.;
 2. Розглядається піксель з найменшим рівнем пріоритету. Якщо сусідні пікселі вже пофарбовані та мають одинаковий колір, то фарбуємо цей піксель у колір сусідів. Усі сусіди, які ще не пофарбовані та ще не належать до черги з пріоритетом, додаються у чергу.;
 3. Повторюємо крок 2, поки черга не стане пустою.;
 4. Незафарбовані пікселі є границями водоподілу.;
-

2.1.6 Кінцевий алгоритм сегментації

Отже, кінцевий алгоритм сегментації інтерфазних ядер bukalного епітелію виглядає наступним чином:

Алгоритм 4: Сегментація ядер

Вхід : Зображення X розміром m

Вихід: Зображення Y такого ж розміру, де різні ядра зображені різними кольорами

$gray \leftarrow$ сіре зображення X (grayscale);

$blured \leftarrow$ медіанна фільтрація $gray$ з вікном радіусом $r = 2$;

$thresh \leftarrow$ бінаризація Оцу для $blured$;

$threshOpening \leftarrow$ морфологічне відкриття $thresh$ з ядром 9×9 за 2 ітерації;

$zoneOfInterest \leftarrow$ морфологічне розширення $opening$ з ядром 9×9 за 2 ітерації;

$distTransform \leftarrow$ трансформація відстані для $threshOpening$;

$sureFG \leftarrow z$, які задовільняють формулу (2.1);

$unknownRegion \leftarrow zoneOfInterest \setminus sureFG$;

$markers \leftarrow$ зв’язні компоненти $sureFG \cup sureBG$;

Виконати алгоритм водоподілу.

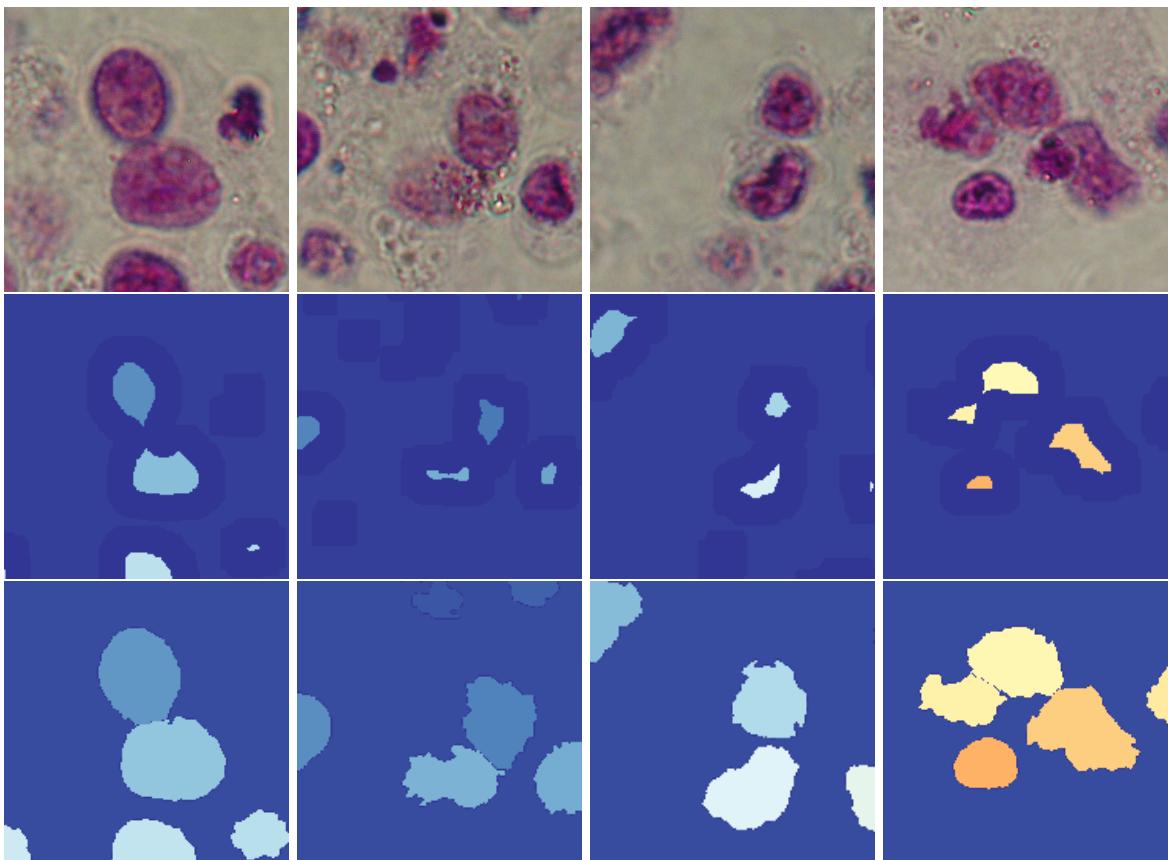
На (Мал. 2.8) бачимо, що наш алгоритм коректно сегментує також ті випадки, коли ядра клітин торкаються одне одного.

2.2 Нормалізація зображень

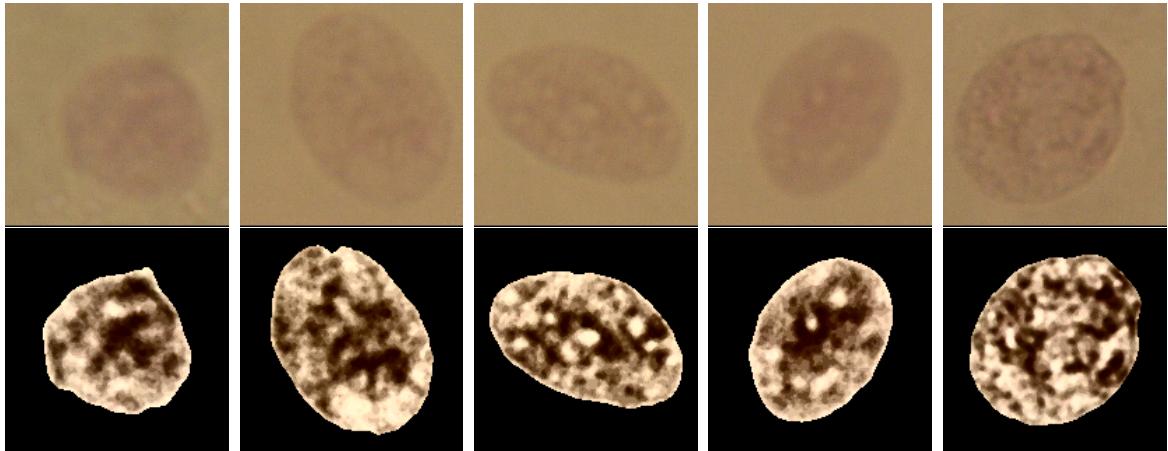
Для нейронної мережі важливо переконатися, що вхідні дані нормалізовані. Для випадку зображень інтерфазних ядер букального епітелію, зображення можуть бути отримані у трохи різних умовах. Це може вплинути на швидкість тренування нейронної мережі, а також на точність розпізнавання.

Задньому плану зображення присвоюємо значення 0, щоб воно не активувало нейрони у мережі. Виконуємо нормалізацію гістограми для зображення одного ядра клітини так, щоб пікселі ядра приймали значення від 0 до 255. Це також може надати прискорення тренуванню нейронної мережі, оскільки більший діапазон значень означає більшу різницю градієнтів у різних точках зображення.

Використовуючи попередні результати, видалимо задній план із зображення. Ядро клітини нормалізуємо за допомогою еквалізації гістограми, результат показано на (Мал. 2.9).



Мал. 2.8: Випадки, коли ядра торкаються чи частково перекриваються.



Мал. 2.9: Нормалізація зображень.

2.3 Фільтр на основі еліпсоїдів Петуніна

Згідно з роботою [8], фільтр зображень на основі використання еліпсоїдів Петуніна має багато переваг над медіанним фільтром. У випадку фільтрації зображень інтерфазних ядер клітин, цей метод краще зберігає форму текстури. Отже, його можна використовувати для попередньої обробки даних для фрактального аналізу.

2.3.1 Еліпсоїди Петуніна

Нехай маємо множину точок $M = \{x_1, x_2, \dots, x_n\}$, де $x_i \in \mathbb{R}^m$ – множина точок у m -вимірному просторі. Знайдемо точки x_l, x_k такі, щоб

$$\sup_{x_i, x_j \in M} \|x_i - x_j\| = \|x_l - x_k\|$$

тобто відрізок $L = x_l x_k$ – діаметр множини точок M . Повернемо і перенесемо систему координат так, щоб діаметр L належав $O_{x'_1}$, де x'_1, x'_2, \dots, x'_n – координати множини M у новій системі координат.

Побудуємо найменший прямокутний паралелепіпед у \mathbb{R}^m , який містив би елементи множини $M' = \{x'_1, x'_2, \dots, x'_n\}$. Стискаючи цей найменший прямокутний паралелепіпед, відобразимо його точки у гіперкуб. Знайдемо центр x_0 гіперкуба та обчислимо відстані r_1, r_2, \dots, r_n від нього до кожної точки.

Знайдемо найбільше $R = \max(r_1, r_2, \dots, r_n)$ та побудуємо гіперкулю навколо центру x_0 з радіусом R . Зробимо обернене перетворення розтягу, повороту та переносу, отримаємо еліпсоїд Петуніна у m -вимірному просторі для множини M . Система еліпсоїдів, отриманих з кіл радіусами r_i , де $i \in 1, \dots, n$ називається концентрованою системою еліпсоїдів Петуніна для множини M .

Для системи концентрованих еліпсоїдів також справжується властивості, що є наслідками роботи Хілла [9] та детально описані та доведені у [10].

2.3.2 Фільтр зображень на основі еліпсоїдів Петуніна

Нехай маємо на вході зображення, кожен піксель якого має d каналів. Зображення розглядається як відображення $M : \mathbb{Z}^k \rightarrow \mathbb{R}^d$, тобто кожній координаті присвоюється d -канальне значення. Для кожного пікселя зображення, розглядаємо ще декілька пікселів, координати яких знаходяться в радіусі r від цього пікселя. Побудуємо для множини значень цих пікселів концентровану систему еліпсоїдів Петуніна, та у вихідному зображені присвоїмо пікселю з тими ж координатами значення, яке відповідає внутрішньому (або центральному) еліпсоїду.

Детальний алгоритм виглядає наступним чином:

Алгоритм 5: Фільтр Петуніна

Вхід : Зображення X розміром $w \times h$, радіус ядра (вікна фільтру) r

Вихід: Зображення Y такого ж розміру

```

for  $i = 1$  to  $w$  do
    for  $j = 1$  to  $h$  do
         $D \leftarrow$  нескінчений відрізок;
        for  $k = -r$  to  $r$  do
            for  $l = -r$  to  $r$  do
                for  $o = -r$  to  $r$  do
                    for  $p = -r$  to  $r$  do
                        Якщо  $\|D\| < \|X_{i+k,j+l} - X_{i+o,j+p}\|$  то  $D \leftarrow X_{i+k,j+l} X_{i+o,j+p}$ 

```

Побудувати систему концентрованих еліпсоїдів Петуніна, маючи діаметр D ;
 $Y_{i,j} \leftarrow$ значення, що відповідає внутрішньому (або центральному) еліпсоїду;

2.3.3 Швидкий алгоритм на основі динамічного програмування

Незважаючи на переваги фільтру на основі еліпсоїдів Петуніна над медіанним фільтром (Алгоритм 5) має серйозний недолік. Так, для зображення розміром $w \times h$ та фільтру радіусом r , асимптотична оцінка швидкості цього алгоритму становить:

$$O(w \cdot h \cdot r^4).$$

Для порівняння, асимптотична оцінка швидкості швидкого алгоритму медіанної фільтрації (Алгоритм 1) або фільтра Гауса становить

$$O(\min\{w, h\} \cdot r^2 + w \cdot h \cdot r).$$

Отже, при великих розмірах ядра (наприклад, $r > 10$) та великому об'ємі зображень, фільтр зображень на основі еліпсоїдів Петуніна значно повільніше працює, ніж інші алгоритми фільтрації. Основна причина – повільний пошук діаметру.

Існує багато швидких алгоритмів пошуку діаметру множини точок у d -вимірному просторі. Так, наприклад, Har-Peled запропонував евристичний алгоритм знаходження діаметру [11]. У [12] запропоновано детермінований алгоритм, що дозволяє знайти діаметр за час $O(n \log n)$, де n – кількість точок множини. Недоліком цих алгоритмів при використанні для фільтру на основі еліпсоїдів Петуніна є те, що незважаючи на меншу асимптотику, ці алгоритми при кількості точок менших, ніж 1000, працюють не швидше, ніж простий

алгоритм (повний перебір) пошуку діаметру. Оскільки розміри фільтрів більших ніж 30 майже не мають практичного використання, ми не можемо використовувати ці алгоритми для прискорення процесу пошуку діаметру.

Можна помітити, що при “пересування” вікна фільтру, ми повторно виконуємо деякі обчислення. Якщо зможемо оптимальніше виконувати обчислення (без повторень), то зможемо швидше фільтрувати зображення. Автором цієї роботи пропонується алгоритм на основі динамічного програмування, що значно прискорює процес пошуку діаметру.

Опишемо структуру даних “двовимірна черга що зберігає діаметр”. Ця структура представляє собою двійка (H, p) , де p – координата центру вікна фільтру відносно зображення X , а H – таблиця розміром $(2r+1) \times (2r+1)$, кожен елемент якого зберігає діаметр множини точок, що відповідає пікселям, що потрапляють у фільтр нижче чи правіше, тобто:

$$H_{i,j} = \|D(\{X_{y,x} | x \in [(p_x + i), (p_x + r)]_Z, y \in [(p_y + j), (p_y + r)]_Z\})\|, \quad i, j \in -r \dots r \quad (2.2)$$

де D – діаметр множини, $X_{y,x}$ – d -вимірне значення пікселю з координатами (y, x) . Очевидно, що в $H_{-r,-r}$ зберігається діаметр множини точок значень пікселів, що потрапляють у вікно фільтру.

Опишемо алгоритм оновлення таблиці H при пересуванні фільтру праворуч. Алгоритм при пересуванні фільтру вниз задається аналогічним чином.

Алгоритм 6: Оновлення таблиці H при пересуванні вікна фільтру праворуч

Вхід : Зображення X , радіус ядра (вікна фільтру) r , двовимірна черга що зберігає діаметр (H, p) що відповідає вікну
Вихід: двовимірна черга що зберігає діаметр (H', p') що відповідає вікну після зміщення праворуч на 1 піксель

```

 $H'_{r,r} \leftarrow 0;$ 
for  $i = r - 1$  downto  $-r$  do
     $H'_{r,j} \leftarrow \max \left\{ \max \left\{ \|X_{p_y+i,p_x+r} - X_{p_y+k,p_x+r}\| \mid \forall k \in i \dots r \right\}, \quad H'_{i+1,r} \right\}$ 
for  $j = r - 1$  downto  $-r$  do
    for  $i = r$  downto  $-r$  do
         $H'_{i,j} \leftarrow \max \left\{ \max \left\{ \|X_{p_y+i,p_x+j} - X_{p_y+k,p_x+r}\| \mid \forall k \in i \dots r \right\}, \quad H'_{i+1,j}, \quad H'_{i,j+1} \right\}$ 

```

Також опишемо алгоритм для обчислення (H, p) , знаючи (H^{left}, p^{left}) та (H^{top}, p^{top}) , де $p^{left} = (p_y, p_x - 1)$, $p^{top} = (p_y - 1, p_x)$, тобто знаючи черги, що відповідають вікнам на піксель вище чи лівіше.

Алгоритм 7: Обчислення таблиці H знаючи відповідні таблиці ліворуч та вище

Вхід : Зображення X , радіус ядра (вікна фільтру) r , двовимірні черги що зберігає діаметр (H^{left}, p^{left}) та (H^{top}, p^{top}) що відповідає вікну лівіше та вище відповідно
Вихід: двовимірна черга що зберігає діаметр (H, p)

```

 $H_{r,r} \leftarrow 0;$ 
for  $i = r$  downto  $-r$  do
    for  $j = r$  downto  $-r$  do
         $H_{i,j} = \max \left\{ H_{i,j+1}^{left}, \quad H_{i+1,j}^{top}, \quad H_{i,j+1}, \quad H_{i+1,j}, \quad \|X_{p_y+r,p_x+r} - X_{p_y+i,p_x+j}\|, \right\}$ 

```

Слід зазначити, що в (Алгоритм 6) та (Алгоритм 7), після кожного кроку присвоєння значення діаметру, слід також у таблиці $D_{i,j}$ зберігати точки кінців діаметру. Опишемо кінцевий алгоритм швидкої фільтрації зображення на основі системи концентрованих еліпсоїдів Петуніна:

Алгоритм 8: Швидкий алгоритм фільтрації Петуніна

Вхід : Зображення X розміром $w \times h$, радіус ядра (вікна фільтру) r

Вихід: Зображення Y такого ж розміру

Обчислити $(H_0, p_0 = (1, 1))$ згідно за формулою (2.2), а також відповідний відрізок L_0 ;

for $i = 1$ to h do

За (Алгоритм 6) обчислити $(H^{i,1}, p^{i,1} = (i, 1))$ та відповідний відрізок $L^{i,1}$;
Побудувати систему концентрованих еліпсоїдів Петуніна, маючи діаметр $L^{i,1}$;
 $Y_{i,1} \leftarrow$ значення, що відповідає внутрішньому (або центральному) еліпсоїду;

for $j = 1$ to w do

За (Алгоритм 6) обчислити $(H^{1,j}, p^{1,j} = (1, j))$ та відповідний відрізок $L^{1,j}$;
Побудувати систему концентрованих еліпсоїдів Петуніна, маючи діаметр $L^{1,j}$;
 $Y_{1,j} \leftarrow$ значення, що відповідає внутрішньому (або центральному) еліпсоїду;

for $i = 1$ to h do

for $j = 1$ to w do

За (Алгоритм 7) обчислити $(H^{i,j}, p^{i,j})$ та відповідний відрізок $L^{i,j}$;
Побудувати систему концентрованих еліпсоїдів Петуніна, маючи діаметр $L^{i,j}$;
 $Y_{i,j} \leftarrow$ значення, що відповідає внутрішньому (або центральному) еліпсоїду;

Очевидно, що (Алгоритм 6) має асимптотичну оцінку $O(r^3)$, а (Алгоритм 7) має асимптотичну оцінку $O(r^2)$. Тоді асимптотична оцінка для (Алгоритм 8) дорівнює:

$$O(r^4 + r^3(w+h) + w \cdot h \cdot r^2)$$

Оскільки значення r значно менша ніж розмір зображення $w \times h$, то (Алгоритм 8) можна оцінити як

$$O(w \cdot h \cdot r^2)$$

Що дорівнює швидкості звичайного алгоритму медіанної фільтрації.

Роздiл 3

Оцiнка якостi даних

3.1 Алгоритм РСА для зменшення розмiрностi

3.2 Класифiкатор SVM

3.3 Результати оцiнки якостi

Розділ 4

Фрактальний аналіз

Наступним етапом цієї роботи є обчислення фрактальної розмірності для попередньо оброблених зображень. Таким чином, для кожного пацієнта буде побудовано набір фрактальних характеристик.

4.1 Фрактальна характеристика

4.1.1 Теоретичні відомості

Для визначення поняття фракталу треба спочатку навести визначення розмірності Хаусдорфа. Нехай Ω – обмежена множина у метричному просторі X .

Означення 1. Нехай $\epsilon > 0$. Не більш ніж зліченне сімейство $\{\omega_i\}_{i \in I}$ підмножин простору X будемо називати ϵ -покриттям множини Ω , якщо виконуються наступні умови:

$$(1) \quad \Omega \subset \bigcup_{i \in I} \omega_i$$

$$(2) \quad \forall i \in I : |\omega_i| < \epsilon$$

де $|\omega_i|$ – діаметр множини ω_i .

Нехай $\alpha > 0$. Нехай $\Theta = \{\omega_i\}_{i \in I}$ – покриття множини Ω . Визначимо наступну функцію, яка в деякому плані визначає розмір цього покриття:

$$F_\alpha(\Theta) := \sum_{i \in I} |\omega_i|^\alpha$$

Позначимо через $M_\alpha^\varepsilon(\Omega)$ “мінімальний розмір” ϵ -покриття множини Ω :

$$M_\alpha^\varepsilon(\Omega) := \inf(F_\alpha(\Theta))$$

де інфімум береться по всіх ϵ -покриттях множини Ω . Очевидно, що функція $M_\alpha^\varepsilon(\Omega)$ не спадає при зменшенні ϵ , оскільки при зменшенні ϵ ми також звужуємо множину можливих ϵ -покриттів. Отже, у неї є скінчenna або нескінчenna границя при $\varepsilon \rightarrow 0+$:

$$M_\alpha(\Omega) = \lim_{\varepsilon \rightarrow 0+} M_\alpha^\varepsilon(\Omega)$$

Означення 2. Величину $M_\alpha(\Omega)$ називають α -мірою Хаусдорфа множини Ω .

З властивостей α -міри Хаусдорфа слід визначити те, що $M_\alpha(\Omega)$ спадає по α . Більш того для будь-якої множини Ω існує критичне значення α_0 , таке, що:

1. $M_\alpha(\Omega) = 0$ для всіх $\alpha > \alpha_0$

2. $M_\alpha(\Omega) = +\infty$ для всіх $\alpha < \alpha_0$

Значення $M_{\alpha_0}(\Omega)$ може бути нульовим, скінченно додатнім або нескінченим.

Означення 3. Число α_0 називають розмірністю Хаусдорфа множини Ω .

Наведемо означення фракталу (або фрактальних множин), запропонований Мандельбротом [13]:

Означення 4. (Мандельброт) Множина називається фракталом, якщо її розмірність Хаусдорфа строго перевищує його топологічну розмірність.

4.1.2 Алгоритм box-counting

Як правило, фрактальні множини мають складну геометричну структуру, а також мають властивість самоподібності. Характеристика, що описує цю властивість, є фрактальною розмірністю [14].

При вимірюванні фрактальної розмірності різних природних і штучних об'єктів виникає ряд проблем, пов'язаних з тим, що існує кілька визначень фрактальної розмірності. Базовим поняттям є розмірність Хаусдорфа, але її обчислення часто виявляється досить непростою задачею. Тому на практиці частіше використовуються інші розмірності, що відносяться до так званого класу box-computing (або box-counting) [15]. Цей метод полягає у тому, що для довільного додатного δ обчислюється деяка функція $M_\delta(\Omega)$. Якщо $M_\delta(\Omega) \propto \delta^{-D}$, то множина Ω має фрактальну розмірність D . З пропорції випливає, що

$$\dim_B \Omega = D = \lim_{\delta \rightarrow 0} \frac{\log(M_\delta(\Omega))}{-\log(\delta)},$$

де значення $M_\delta(\Omega)$ дорівнює кількості S -вимірних кубів зі сторонами δ , необхідних для покриття множини Ω .

Означення 5. Число $\dim_B \Omega$ називають box-counting розмірністю множини Ω .

Зауважимо, що $\dim_B \Omega$ не завжди може існувати. Зв'язок між box-counting розмірністю $\dim_B \Omega$ та розмірністю Хаусдорфа $\dim_H \Omega$ виражається наступною теоремою:

Теорема 1. $\dim_B \Omega \leq \dim_H \Omega$

4.1.3 Лінійна регресія

Як зазначено у попередньому пункті, використання алгоритму обчислення фрактальної розмірності box-counting вимагає перевірки пропорції $M_\delta(\Omega) \propto \delta^{-D}$, а отже, і обчислення регресійної прямої. Кутовий коефіцієнт цієї прямої буде дорівнювати $\dim_B \Omega$. Пряму регресії будемо шукати методом найкращого середньоквадратичного наближення – треба знайти такий елемент Φ_0 , який мінімізує значення $\sum_{i=0}^n (\Phi_0(x_i) - y_i)^2$.

Нехай M_n – лінійна оболонка базису $\{\phi_0, \phi_1, \dots, \phi_n\}$, за яким ми шукаємо елемент найкращого середньоквадратичного наближення Φ_0 . Його існування для випадку гільбертового простору H випливає за теоремою. Нехай ми шукаємо функцію найкращого середньоквадратичного наближення для функції f , (\cdot, \cdot) – внутрішній добуток в даному гільбертовому просторі (у нашому випадку в якості такого простору логічно взяти простір неперервних функцій на області визначення функції f). Тоді за теоремою:

$$\forall i \in M_n (f - \Phi_0, \Phi) = 0.$$

Тоді

$$\forall i \in \{0, 1, \dots, n\} (f - \Phi_0, \phi_i) = 0.$$

Нехай

$$\Phi_0 = \sum_{i=0}^n c_i \phi_i,$$

тоді

$$\forall i \in \{0, 1, \dots, n\} \left(f - \sum_{j=0}^n c_j \phi_j, \phi_i \right) = 0,$$

$$\sum_{j=0}^n c_j (\phi_j, \phi_i) = (f, \phi_i), i \in \{0, 1, \dots, n\}.$$

Таким чином ми одержали систему лінійних алгебраїчних рівнянь з матрицею, що є матрицею Грамма $G = \|(\phi_i, \phi_j)\|_{i,j \in \{0,1,\dots,n\}}$. Отже, якщо $\det(G) = 0$, то можна знайти розв'язок даної системи (в нашому випадку використовується метод Гауса), а отже і шукану функцію Φ_0 .

Отже, регресійну пряму можна знайти, побудувавши елемент найкращого середньоквадратичного наближення Φ_0 , обравши базис $\{\phi_0, \phi_1\} = \{1, x\}$.

4.2 Класифікація пацієнтів

Слід помітити, що кожному пацієнту ставиться у відповідність набір (вибірка) фрактальних характеристик, розмір якої не є сталим. Отже, для порівняння набору характеристик різних пацієнтів, треба використовувати міру близькості між вибірками.

Оскільки міра близькості не є метрикою, вибір методів класифікації є дуже обмеженим. В цій роботі було використано найпростіший спосіб – метод k найближчих сусідів.

4.2.1 Міра близькості та p -статистика

Через H позначимо гіпотезу про рівноть неперервних функцій розподілу $F_G(u)$ і $F_{G'}(u)$ генеральних сукупностей G і G' відповідно. Нехай $x = (x_1, x_2, \dots, x_n) \in G$ і $x' = (x'_1, x'_2, \dots, x'_m) \in G'$, $x_{(1)} \leq \dots \leq x_{(n)}$, $x_{(1)} \leq \dots \leq x_{(m)}$ – порядкові статистики. Припустимо що $F_G(u) = F_{G'}(u)$. Позначимо через $A_{ij}^{(k)}$, $k = 1, 2, \dots, m$ випадкова подія, яка полягає у тому, що x' потрапляє в інтервал $(x_{(i)}, x_{(j)})$, тобто $A_{ij}^{(k)} = \{x'_k \in (x_{(i)}, x_{(j)})\}$. Як відомо, ймовірність цієї події обчислюється за формулою:

$$P(A_{ij}^{(k)}) = P\{x'_k \in (x_{(i)}, x_{(j)})\} = p_{ij}^{(n)} = \frac{j-i}{n+1} = \frac{q}{n+1}, \quad q = j-i.$$

Покладемо

$$p_{ij}^{(1)} = \frac{h_{ij}^{(n)} m + \frac{1}{2} g^2 - g \sqrt{h_{ij}^{(n)} (1-h)m + \frac{1}{4} g^2}}{m + g^2},$$

$$p_{ij}^{(2)} = \frac{h_{ij}^{(n)} m + \frac{1}{2} g^2 + g \sqrt{h_{ij}^{(n)} (1-h)m + \frac{1}{4} g^2}}{m + g^2},$$

де $h_{ij}^{(n)}$ – частота події $A_{ij}^{(k)}$ в m експериментів, величина $g = 3$.

Позначимо через N кількість усіх довірчих інтервалів $I_{ij}^{(n,m)} = (p_{ij}^{(1)}, p_{ij}^{(2)})$ (тобто $N = n(n-1)/2$) і L – кількість інтервалів $I_{ij}^{(n,m)}$, які містять ймовірності $p_{ij}^{(n)}$. Покладемо $h^{(n,m)} = \rho(F^*, F^{*'}) = \rho(x, x') = L/N$. Оскільки $h^{(n,m)}$ – частота випадкової події $B = \{p_{ij}^{(n)} \in I_{ij}^{(n,m)}\}$, ймовірність якої дорівнює $p(B) = 1 - \beta$, тоді, вважаючи що $h^{n,m} = h^{(n)}$, $m = N$ та $g = 3$, отримаємо довірчий інтервал $I^{(n,m)} = (p^{(1)}, p^{(2)})$ для ймовірності $p(B)$. Статистика $h^{(n)}$ називається p -статистикою. Вона також є мірою близькості $\rho(x, x')$ між вибірками x та x' [16].

4.2.2 Метод k найближчих сусідів. Кросвалідація.

Використовуючи Метод k найближчих сусідів для класифікації пацієнтів, необхідно зробити припущення коректності постулату про компактність:

Припущення 1. переважна більшість об'єктів, що належать до одного класу, є близчими один до одного, ніж до об'єктів іншого класу, і лежать в області з відносно простою межею.

Суть цього методу для класифікації елементу x полягає у тому, що з навчальної вибірки обираємо k найближчих (у розумінні міри близькості) до x елементів (характеристик пацієнтів) та відносимо x до домінантного класу з цих k елементів.

Для пошуку оптимального числа k , а також розміру навчальної вибірки R використовують метод кросвалідації. Суть цього методу заключається у тому, що для кожного значення k та R випадковим чином підбирається набір даних для навчання.

4.3 Результати та спостереження

Вхідний набір даних для дослідження кладається з знімків 6751 інтерфазних ядер букального епітелія, для кожного було зроблено 3 знімки мікроскопу: без фільтру, через жовтий фільтр та через пурпурного фільтру (отже всього 20253 фотографії), взятого з 130 пацієнтів, з них 68 хворих раком, 29 здорових та 33 хворих іншою хворобою. У таблиць (Мал. ?? - ??) приведено результати кросвалідації.

Отже, при $k = 12$ та $R = 0.8$ досягається найкраща чутливість, а при $k = 6$ та $R = 0.9$ досягається найкраща специфічність.

Розділ 5

Глибинне навчання

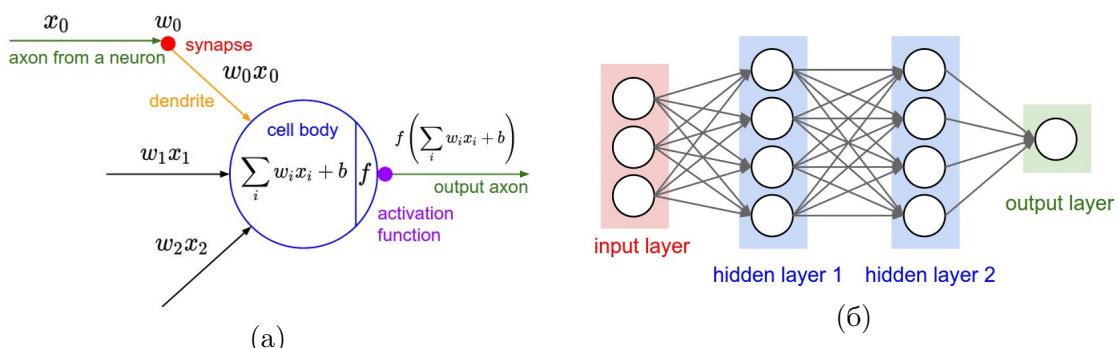
Як зазначено у попередньому розділі (Розділ 4), фрактальний аналіз текстури інтерфазного ядра букального епітелію і взагалі будь-який алгоритм класифікації, який використовує характеристики текстури та контуру ядра клітини, добре розпізнають клас здорових пацієнтів від пацієнтів, хворих раком молочної залози чи фіброаденоматозом, але майже не можуть розділяти останні два класи між собою. Це пояснюється тим, що ці хвороби викликають дуже схожі гормональні реакції, що призводить до схожих змін у ядрах ДНК клітин букального епітелію.

Цей розділ буде присвячено задачі класифікації тільки класів хворих раком та фіброаденоматозом. За допомогою технік глибинного машинного навчання, ми намагаємося аналізувати та використовувати локальні ознаки та особливості текстури поверхні ядра клітини, які не можна характеризувати одним числом.

5.1 Нейронні мережі

Штучна нейронна мережа – це математична модель, яка побудована за принципом функціонування біологічних нейронних мереж. Штучні нейронні мережі широко використовують для задач машинного навчання, розпізнавання образів, дискримінантного аналізу, кластерування тощо.

Нейронну мережу можна інтерпретувати як граф. У підручнику [17], нейронні мережі описуються як граф потоку сигналу (signal-flow graph), де кожен нейрон (або вершина у графі) є окремою одиницею обчислення. Схема одного штучного нейрону зображена на (Мал. 5.1 (а)).



Мал. 5.1: Структура одного нейрону (а) та схема багатошарових перцептронів як граф потоку сигналів (б).

Кожен нейрон мережі має справу лише з сигналами, які він періодично отримує, і сигналами, які він періодично надсилає іншим нейронам. і тим не менш, будучи з'єднаними в достатньо велику мережу з керованою взаємодією, такі прості нейрони разом здатні моделювати складні функції, та виявляти складні залежності між вхідними й вихідними даними, а також здійснювати узагальнення.

5.1.1 Багатошаровий перцептрон

Штучний багатошаровий перцептрон – це така архітектура нейронних мереж, у якій групи нейронів формують “шари”, в яких кожний нейрон з одного шару зв’язана з усіма нейронами сусіднього шару. Отже, сигнали на вихід одного шару нейронів є сигналами на вхід нейронів наступного шару. На (Мал. 5.1 (б)) зображено схему багатошарового перцептрону. Такі шари нейронів також називають “повністю зв’язаними” (Fully-connected layers).

5.1.2 Функція активації ReLU

У цій роботі було використано функцію зрізаних лінійних вузлів, або ReLU, для активації нейронів замість сигмоїдальних функцій чи гіперболічного тангенсу. Функція активації ReLU (Rectified Linear Unit), яка задається наступним чином:

$$f(x) = \max(0, x)$$

має багато переваг над сигмоїдальними функціями чи гіперболічним тангенсом при використанні в якості функції активації:

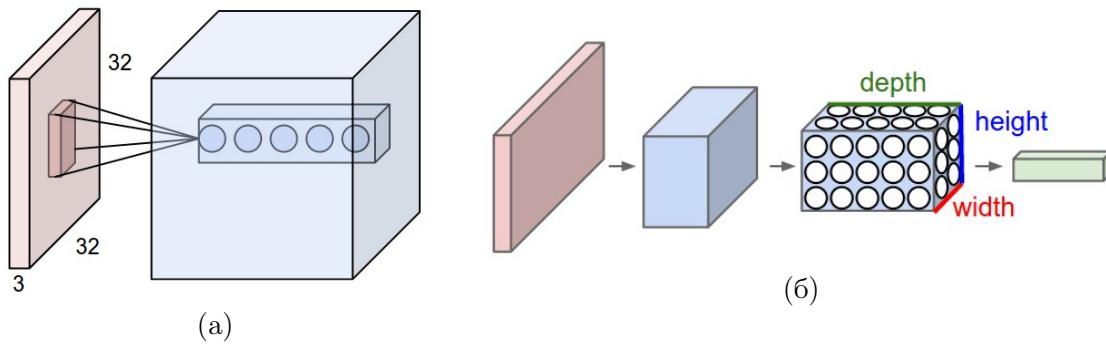
- В роботі [18], було доведено що така функція активації може прискорити збіжність спуску за градієнтом до 6 разів у порівнянні з сигмоїдальними функціями чи гіперболічними тангенсами.
- Сигмоїдальні функції та гіперболічні тангенси використовують складні операції, такі як зведення у степінь.

Детальний аналіз переваг та недоліків цієї функції, а також позбавлення недоліків за допомогою модифікованої версії “Leaky ReLU”, описаний в роботі [19].

5.1.3 Згорткові нейронні мережі

Згорткові нейронні мережі – це мережі, які містять згорткові шари нейронів, уперше запропонованій [20]. Згортковий шар нейронів містить фільтри – набір ваг. При подачі даних на вхід, вихідний сигнал обчислюється як згортка фільтру з локальним регіоном зображення. На (Мал. 5.2 (а)) показано схему першого згорткового шару, де кожен нейрон локально зв’язаний з пікселями зображення, та “бачить” усі канали. Кількість шарів нейронів у згортковому шару (зображені кружечками) дорівнює кількості фільтрів згорткового шару, тобто для різних шарів нейронів використовується різні фільтри.

Згорткові нейронні шари використовують для виявлення локальних особливостей зображення. Ці шари також дозволяють мережі бути більш стійкою до трансляції, тобто надають деяку інваріантність. Тому, будемо використовувати згорткові шари для побудови класифікатора між хворими раком та фіброаденоматозом.



Мал. 5.2: На (а) схематично зображене згортковий шар. На (б) зображенено схему згорткової мережі.

Разом із згортковими нейронними шарами також будемо використовувати підвибіркові шари (pooling layers), які значно зменшують обсяг вхідних даних на наступні шари мережі.

5.1.4 Алгоритм RMSProp

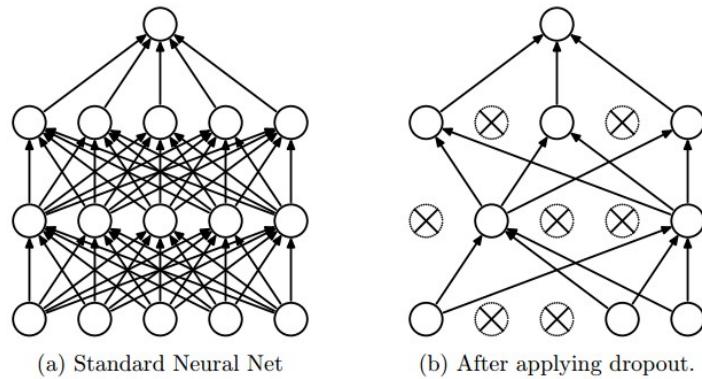
Після обчислення градієнтів за допомогою алгоритму зворотного поширення помилок (backpropagation), ці градієнти використовуються для оновлення параметрів (ваг) нейронної мережі. Зазначимо, що оптимізація алгоритму оновлення нейронних мереж є на момент написання цієї роботи дуже активною сфорою досліджень.

Алгоритм RMSProp є дуже ефективним, але ще не опублікований метод адаптивного оновлювання параметрів. Він налаштовує та регулює алгоритм Adagrad [21] таким чином щоб зменшити різкість монотонного зменшення швидкості навчання. На відміну від Adagrad, цей алгоритм використовує так звані “ковзаючі” середньоквадратичні значення градієнтів:

$$\begin{aligned} \text{cache} &= \text{decay_rate} \cdot \text{cache} + (1 - \text{decay_rate}) \cdot dx^2 \\ x &= x - \frac{\text{learning_rate} \cdot dx}{\sqrt{\text{cache}} + \text{eps}} \end{aligned}$$

де decay_rate це гіперпараметр, якому зазвичай присвоюють значення з множини $\{0.9, 0.99, 0.999\}$. Отже, метод RMSProp також модулює швидкість навчання кожного параметру з урахуванням величин його градієнтів, що має вигідний ефект вирівнювання. Але, на відміну від Adagrad, величини оновлень не зменшуються монотонно.

У цій роботі використовується саме алгоритм RMSProp. На практиці більш стандартним методом для оновлення параметрів є алгоритм Adam (adaptive moment estimation) [22], який дуже схожий на RMSProp та у деяких випадках є трохи ефективнішим ніж RMSProp, але цей алгоритм є більш громіздким та потребує трохи більше часу для обчислень. Повний алгоритм Adam також включає механізм корекції зміщень. Альтернативою RMSProp та Adam є також комбінація алгоритмів SGD та Nesterov Momentum.



Мал. 5.3: Візуалізація методу Dropout з ймовірністю $p = 0.5$.

5.1.5 Регуляризація нейронної мережі

Процес регуляризації нейронної мережі – це процес уникнення перенавчання (overfitting) мережі. існує кілька способів контролювання тренування нейронної мережі, щоб максимальним чином уникнути перенавчання.

У цій роботі використовується метод Dropout – простий та дуже ефективний метод регуляризації, запропонований [23]. При тренуванні, кожен нейрон є активним за ймовірністю p (що є гіперпараметром мережі).

На (Мал. 5.3) схематично зображенено принцип роботи методу Dropout. На практиці також часто використовують алгоритми $L1$ - та $L2$ - регуляризації нейронних мереж, разом з методом Dropout.

5.1.6 Функція втрати. Перехрестна ентропія.

Вибір функції втрати є одним з найважливіших кроків при побудові моделі класифікатора. Вдало обрана функція втрат може значно покращити результат та швидкість тренування класифікатора.

На практиці при тренуванні класифікаторів на основі нейронних мереж найчастіше використовують перехресну ентропію (categorical crossentropy):

$$L_i = -\log \left(\frac{e^{f_{y_i}}}{\sum_j e^{f_j}} \right)$$

яка для випадку задач класифікації дає кращі результати, ніж середньоквадратична функція втрат (mean square error) та функція помилки класифікації (classification error). У цій роботі використовується саме ця функція оцінювання втрат. Слід також зазначити, що при тренуванні моделі та при класифікації чи обчисленні можна використовувати різні оцінки втрат.



Мал. 5.4: Прирошення даних на прикладі одного зображення.

5.2 Класифікація за допомогою нейронних мереж

5.2.1 Прирошення даних

Як вказано у розділі 1, для дослідження маємо набір даних який містить 3403 зображень інтерфазних ядер букального епітелію, взятих з 68 пацієнтів хворих раком молочної залози, та 1741 зображені ядер взятих з 33 пацієнтів хворих фіброаденоматозом. Така кількість є дуже малою для тренування глибоких згорткових нейронних мереж. Отже, будемо штучно збільшувати кількість даних для тренування.

При подачі в нейронну мережу випадковим чином зробимо декілька з наступних перетворень над зображеннями:

- Перевернути зображення відносно горизонтальної осі
- Перевернути зображення відносно вертикальної осі
- Повернути зображення на випадковий кут від 0 до 2π , при чому пікселі, що до перетворення знаходилися за межами входу, присвоюються значенню сусідніх пікселів.

5.2.2 Проблема незбалансованості вибірки для тренування

Для тренування було обрано випадковим чином 80% від всього набору даних. Це становить 1392 зображень взятих з 33 пацієнтів хворих фіброаденоматозом та 2722 зображень з 68 пацієнтів хворих раком молочної залози. Отже, в набір даних для валідації входять інші зображення ядер букального епітелію тих самих пацієнтів. Маємо незбалансованість вибірки для тренування, де кількість даних з одного класу (з класу хворих раком) майже вдвічі перевищує кількість даних з іншого класу (хворих фіброаденоматозом).

Незбалансованість набору даних для тренування може значно негативно вплинути на навчання класифікатора. Слід зазначити, що набір даних для валідації теж є незбалансованим. Так, після тренування моделі на незбалансованій вибірці, автором цієї роботи було отримано класифікатор, точність розпізнавання якого на обох наборах даних для тренування і тестування становила більш ніж 65%. Подальший аналіз результатів виявив, що отримана модель має чутливість, близьку до 100% та специфічність, близьку до 0%, тобто модель у більшості випадків просто класифікувала зображення як “належить хвортого раком”. Отже, необхідно знайти спосіб боротьби з проблемою незбалансованості набору даних для тренування.

Існує багато методів для боротьби з цією проблемою. Найбільш поширеними на практиці є методи зміни об'єму даних (subsampling methods) та методи навчання з вагами (cost-sensitive methods), де різним класам присвоюється деяке значення “критичної” помилки. Було вирішено використовувати перший метод.

Отже, після зміни об'єму даних, кількість зображень для тренування становить 1392 з хворих фіброаденоматозом та 1632 з хворих раком молочної залози, а кількість зображень для валідації становить 349 зображень з хворих фіброаденоматозом та 1771 зображення з хворих раком молочної залози.

5.2.3 Архітектура мережі

Для зручного опису архітектури мереж, використаємо наступні позначення:

- Позначимо повністю зв'язані шари нейронів (fully-connected layers) як xFC , де x – кількість нейронів шару. Наприклад, якщо повністю зв'язаний шар має 128 нейронів, то позначаємо $128FC$.
- Згорткові нейронні шари (convolutional layers), де фільтри шару мають розміри $x \times x$, позначимо як $yCONVx$, де y – кількість фільтрів. Наприклад, $64CONV3$.
- Підвибіркові шари (max-pooling layers) позначимо як MPx , де x – розмір ядра фільтру. Наприклад, $MP2$.
- Вхідні (input) та вихідні (output) дані позначимо як I та O відповідно.
- Домовимось, що сигнал протікає з шару, що знаходиться лівше у записі, до шару, що знаходиться правіше. Сусідні шари нейронів з'єднуємо символом “ \rightarrow ”. Повна архітектура мережі так: $I \rightarrow 32CONV3 \rightarrow MP2 \rightarrow FC16 \rightarrow O$.

Автором було розглянуто багато різних архітектур нейронної мережі. Okрім простих відносно неглибоких архітектур згорткової нейронної мережі, які схожі на AlexNet, також були розглянуті мережі, подібні до VGG [24].

ARCH-0 $I \rightarrow 32CONV3 \rightarrow MP2 \rightarrow 32CONV3 \rightarrow MP2 \rightarrow 64CONV5 \rightarrow MP2 \rightarrow FC32 \rightarrow FC32 \rightarrow O$

ARCH-1 $I \rightarrow 32CONV3 \rightarrow MP2 \rightarrow 48CONV3 \rightarrow MP2 \rightarrow 64CONV3 \rightarrow MP2 \rightarrow 72CONV3 \rightarrow MP2 \rightarrow FC64 \rightarrow FC64 \rightarrow O$

ARCH-2 $I \rightarrow 64CONV5 \rightarrow 48CONV3 \rightarrow MP2 \rightarrow FC32 \rightarrow FC16 \rightarrow O$

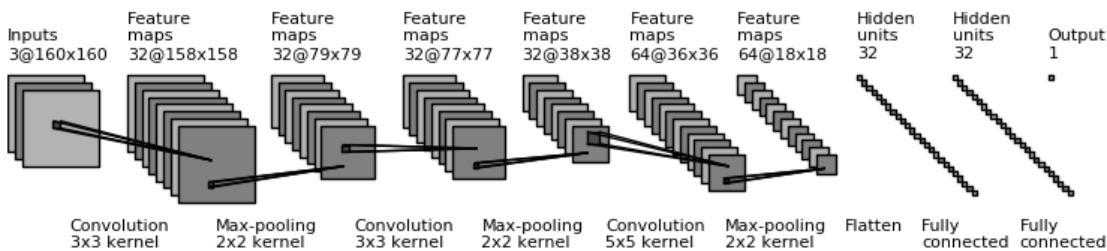
ARCH-3 $I \rightarrow 64CONV3 \rightarrow MP2 \rightarrow 64CONV3 \rightarrow MP2 \rightarrow 96CONV3 \rightarrow MP2 \rightarrow 96CONV3 \rightarrow MP2 \rightarrow FC48 \rightarrow FC48 \rightarrow O$

VGG-9 $I \rightarrow 32CONV3 \rightarrow 32CONV3 \rightarrow MP2 \rightarrow 48CONV3 \rightarrow 48CONV3 \rightarrow MP2 \rightarrow 64CONV5 \rightarrow 64CONV5 \rightarrow MP2 \rightarrow FC32 \rightarrow FC32 \rightarrow O$

VGG-11 $I \rightarrow 48CONV3 \rightarrow 48CONV3 \rightarrow MP2 \rightarrow 64CONV3 \rightarrow 64CONV3 \rightarrow MP2 \rightarrow 72CONV3 \rightarrow 72CONV3 \rightarrow MP2 \rightarrow 96CONV3 \rightarrow 96CONV3 \rightarrow MP2 \rightarrow FC96 \rightarrow FC96 \rightarrow O$

VGG-13 $I \rightarrow 32CONV3 \rightarrow 32CONV3 \rightarrow MP2 \rightarrow 64CONV3 \rightarrow 64CONV3 \rightarrow MP2 \rightarrow 128CONV3 \rightarrow 128CONV3 \rightarrow MP2 \rightarrow 256CONV3 \rightarrow 256CONV3 \rightarrow MP2 \rightarrow 256CONV3 \rightarrow 256CONV3 \rightarrow MP2 \rightarrow FC128 \rightarrow FC128 \rightarrow O$

Для кожного шару нейронів цих архітектур було використано ReLU в якості функції активації, а для кінцевого шару O з одного нейрону (якщо на виході 0 то зображення класифікується як належним до хворого раком молочної залози, а якщо 1 – то класифікується як належним до хворого фіброаденоматозом) було використано функцію активації Soft-Max. З метою регулювання, для кожного повністю зв'язаного шару було використано



Мал. 5.5: Архітектура найкращої мережі.

алгоритм Dropout. В якості алгоритму задання швидкості навчання (тобто оновлення параметрів) використовується RMSProp.

5.2.4 Тренування мереж та результати

Слід зазначити, що вибір архітектури обмежувався потужністю доступних автору обчислювальних ресурсів. У розпорядженні автора на момент дослідження був комп’ютер з процесором Intel Core i7, комп’ютер з процесором Intel Core i3, а також відеокарта Nvidia Geforce GT 520M. Отже, використання більш глибоких нейронних мереж, або мереж з більшою кількістю вільних параметрів, не було можливим.

Кожна з перерахованих вище мереж тренувалася по декілька (для ARCH-0 – декілька десятків) разів, кожен раз нейронна мережа ініціалізувалася випадковим набором ваг (параметрів) та навчалася протягом 300 епох. Кожна епоха – це цикл, в якому нейронна мережа навчається на даних, які по обсягу дорівнюють об’єму тренувальної вибірки. Неповторність входних даних для тренування гарантує алгоритм прирошення даних. Також були використані різні значення для Dropout, від 0.2 до 0.5.

Сумарно навчання вищеписаних мереж проводилося протягом трьох тижнів. Найбільш повільними для навчання є моделі VGG-9, VGG-11, та VGG-13. У середньому на навчання цих моделей протягом 300 епох пішло більш ніж 48 годин процесорного часу.

Найбільш вдалою є архітектура ARCH-0, яка схематично зображена на (Мал. 5.5). Нижче у таблиці приведені найкращі результати класифікації окремого зображення інтерфазного ядра клітини букального епітелію цих мереж.

Архітектура	Чутливість	Специфічність	Точність
ARCH-0	59.17%	67.33%	62.93%
ARCH-1	55.74%	57.59%	56.60%
ARCH-2	67.24%	53.58%	60.95%
ARCH-3	56.72%	58.45%	57.52%
VGG-9	55.25%	59.89%	57.39%
VGG-11	58.43%	58.45%	58.44%
VGG-13	56.23%	55.59%	55.94%

Найкращі результати за точністю дає мережа ARCH-0. Надалі будемо детальніше аналізувати її результати, та намагатися побудувати класифікатор з кращими результатами на основі цієї моделі.

Можна розглядати кожне зображення ядра клітини певного пацієнта не окремо, а у сукупності з іншими зображеннями клітин цього пацієнта, тобто класифікувати не належність окремого ядра букального епітелію до класу хвороби, а належність самого пацієнта до цих класів. Наприклад, якщо більшість зображень ядер клітин цього пацієнта класифікується належними до класу “хворі раком”, то віднесемо цього пацієнта до класу “хворий раком” (аналогічно з фіброаденоматозом). Тоді отримаємо наступні результати:

Архітектура	Чутливість	Специфічність	Точність
ARCH-0	58.82%	69.70%	62.38%

що незначно відрізняються від результатів, отриманих при класифікації окремих клітин.

5.2.5 Аналіз результатів

Помітимо, що мережа з відносно простою архітектурою ARCH-0 дає кращі результати ніж VGG-подібні мережі. Це можна пояснити тим, що VGG-подібні глибокі архітектури містять у декілька десятків разів більше вільних параметрів, а отже, потребують більше ніж 300 епох для навчання. При тренуванні моделі VGG-11, наприклад, спостерігалося дуже вільне спадання значення функції втрат. Це свідчить про те, що, можливо, треба було обирати інші гіперпараметри для мережі, такі як *learning_rate* для алгоритму RMSProp, або використовувати менший *batch_size* (обсяг даних при тренуванні, після якого обчислюється градієнт та виконується одне оновлення параметрів).

5.3 Підвищення якості розпізнавання

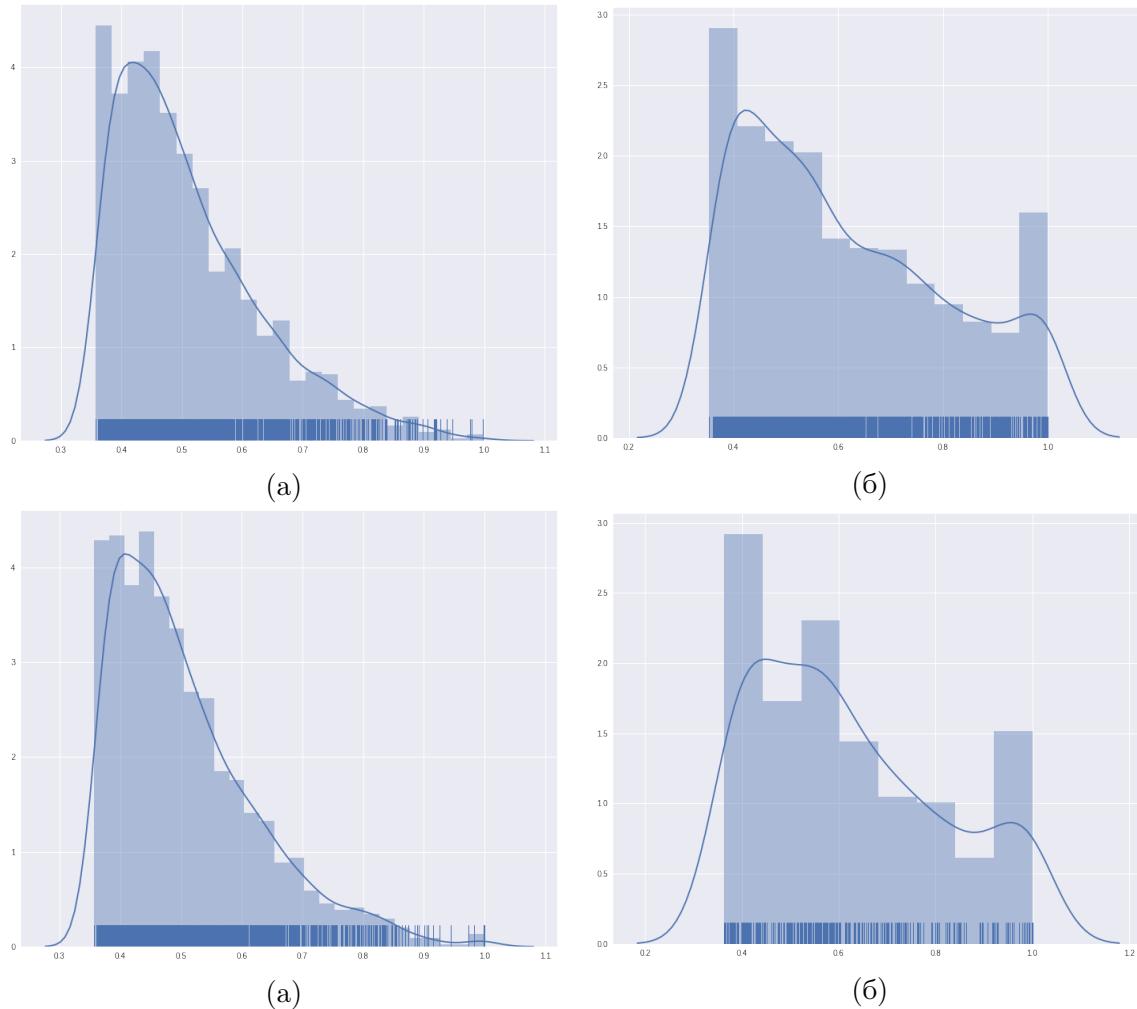
Надалі будемо будувати більш точний класифікатор на основі отриманої моделі нейронної мережі ARCH-0. Щоб переконатися у тому, що ця модель дійсно генералізується на невідомі дані (ті, які не було використано при тренуванні), розглянемо розподіл вихідних значень мережі для обох наборів: тренування та валідації, що показані на (Мал. 5.6). Можна побачити, що отримана мережа добре генералізується на вибірці для валідації.

Отже, припускаючи, що розподіл вихідних значень мережі для пацієнтів з однаковою хворобою (рак молочної залози чи фіброаденоматоз) буде також приблизно одинаковими, замість простого підходу, що описаний у попередньому підрозділі, можна розглядати нейронну мережу не як бінарний класифікатор, а як генератор ознак, тобто кожному зображенню ставити у відповідність деяке число (ознаку). Отже, до кожного пацієнта (набору зображень) можна застосовувати статистичні тести.

5.3.1 Класифікатор на основі нейронної мережі та критерії еквівалентності

В якості критерію еквівалентності генеральних сукупностей було використано критерій Колмогорова-Смірнова, що є одним з найефективніших непараметричних критеріїв для порівняння вибірок. Альтернативно можна використовувати *p*-статистику, що детально описана у попередньому розділі.

Нехай *can_trainX* – сукупність зображень інтерфазних ядер букального епітелію у вибірці для тренування, взятих з пацієнтів хворих раком, а *fib_trainX* – сукупність зображень ядер у вибірці для тренування взятих з пацієнтів хворих фіброаденоматозом.



Мал. 5.6: Структура одного нейрону (а) та схема багатошарових перцентронів як граф потоку сигналів (б).

При класифікації деякого пацієнта $P = \{X^1, X^2, \dots, X_n\}$ обчислюється:

$$\text{can_train} = \text{ARCH-0}(\text{can_train}X) = \{\text{ARCH-0}(X) \mid X \in \text{can_train}X\}$$

$$\text{fib_train} = \text{ARCH-0}(\text{fib_train}X) = \{\text{ARCH-0}(X) \mid X \in \text{fib_train}X\}$$

$$M = \text{ARCH-0}(P) = \{\text{ARCH-0}(X^i) \mid i \in 1 \dots n\}$$

де $\text{ARCH-0}(X)$ – значення виходу нейронної мережі після обробки зображення X . Потім, обчислюється p -значення статистики Колмогорова-Смірнова спочатку для can_train та M , а потім для fib_train та M . Ці значення використаємо як показників схожості цих вибірок.

5.3.2 Аналіз результатів

У таблиці показано результати класифікації пацієнтів як сукупність ядер клітин, отримані за підходом, що описаний вище.

Архітектура	Чутливість	Спеціфічність	Точність
ARCH-0	70.59%	63.64%	68.32%

Слід зазначити, що в нашому випадку у зв'язку зі зміною обсягу даних для навчання спостерігається значна незбалансованість класів у наборі даних для валідації, а саме: 1771 зображені ядер взятих з 68 пацієнтів хворих раком, та 349 - з хворих фіброаденоматозом. Отже, кожний пацієнт, хворий раком, має у середньому 26 зображень клітин ядер у наборі даних для валідації, та в той же час кожен пацієнт, що хворий фіброаденоматозом, має у середньому 10 зображень клітин ядер у вибірці для валідації. Саме тому спостерігається така різниця між чутливістю та специфічністю моделі після валідації. Можна висловити припущення, що маючи достатню кількість зображень для класифікації конкретного пацієнта, можна досягти чутливості та специфічності більші за 71%.

Перегрупуємо дані за пацієнтами, залишаючи мережу ARCH-0 тою самою. Отже, для тестування маємо окремих пацієнтів, для кожного з яких доступно по 50 клітин для валідації. Припускаючи, що розподіл вихідних значень мережі ARCH-0 одинаковий для всіх пацієнтів однакової категорії (з однаковою хворобою) та використовуючи метод, описаний вище, отримуємо наступні результати (чутливість, специфічність та точність у середньому після кросвалідації):

Архітектура	Чутливість	Специфічність	Точність
ARCH-0	71.32%	78.95%	73.81%

Отже, підхід з використанням згорткових нейронних мереж та статистичних методів має дуже перспективні результати для подальшого дослідження.

Роздiл 6

Висновок

Додаток А

Програмний пакет

Bibliography

- [1] N. Boroday et al. “Analysis of malignancy-associated DNA changes in interphase nuclei of buccal epithelium in persons with breast diseases”. In: Experimental Oncology 26.2 (2004), pp. 158–160.
- [2] Weidi Xie, J. Alison Noble, and Andrew Zisserman. “Microscopy Cell Counting with Fully Convolutional Regression Networks”. In: Department of Engineering Science, University of Oxford, UK (2015).
- [3] S. Perreault and P. Hébert. “Median filtering in constant time”. In: IEEE Transactions on Image Processing 16.9 (2007), 2389–2394.
- [4] Nobuyuki Otsu. “A threshold selection method from gray-level histograms”. In: Automatica 11 (285–296 1975), pp. 23–27.
- [5] Jean Serra. Image Analysis and Mathematical Morphology. Academic Press Inc., 1983.
- [6] A. Rosenfeld and J. Pfaltz. “Distance Functions in Digital Pictures”. In: Pattern Recognition 1 (1968), pp. 33–61.
- [7] R. Barnes, C. Lehman, and D. Mulla. “Priority-flood: An optimal depression-filling and watershed-labeling algorithm for digital elevation models”. In: Computers and Geosciences 62 (), pp. 117–127.
- [8] М. В. Присяжная. “Фільтрація зображеній з використанням еліпсів Петуна”. In: Журнал обчислювальної та прикладної математики 1.115 (2014).
- [9] B. M. Hill. “Posterior Distribution of Percentiles: Bayes’ Theorem for Sampling from a Population.” In: Journal of the American Statistical Association 63.322 (1968), p. 677.
- [10] С. И. Ляшко, Д. А. Ключин, and В. В. Алексеенко. “Многомерное ранжирование с помощью эллиптического пилинга”. In: Кибернетика и системный анализ 4 (2013).
- [11] Sariel Har-Peled. “A practical approach for computing the diameter of a point set”. In: Proceedings of the seventeenth annual symposium on Computational geometry (2001).
- [12] G. Malandain and J.-D. Boissonnat. In: International Journal of Computational Geometry and Applications 12.6 (2002), pp. 498–509.
- [13] Benoît Mandelbrot. The Fractal Geometry of Nature. W. H. Freeman and Co., 1982.
- [14] F. Richard Voss. “Fundamental algorithms for computer graphics”. In: Springer, 1985. Chap. Random fractal forgeries, pp. 805–835.
- [15] Ю. В. Шуплецов and Н. Б. Ампилова. “Алгоритм вычисления размерности Минковского для полутоновых изображений”. In: Известия Российского Государственного Педагогического Университета Им. АИ Герцена 165 (2014).
- [16] Д. А. Ключин and Ю. И. Петунин. “Непараметрический критерий эквивалентности генеральных совокупностей, основанный на мере близости между выборками”. In: Український математичний журнал 5 (2003), 147–163.
- [17] Simon Haykin. Neural Networks and Learning Machines. 3rd ed. Pearson Education, 2009.
- [18] A. Krizhevsky, I. Sutskever, and G. E. Hinton. “Imagenet classification with deep convolutional neural networks”. In: Advances in neural information processing systems (2012), 1097–1105.
- [19] He Kaiming et al. “Delving Deep into Rectifiers: Surpassing Human-Level Performance on ImageNet Classification”. In: (2015).

- [20] Yann LeCun et al. In: Proceedings of the IEEE 86.11 (1998), pp. 2278–2324.
- [21] J. Duchi, E. Hazan, and Y. Singer. “Adaptive Subgradient Methods for Online Learning and Stochastic Optimization”. In: 12 (2011), pp. 2121–2159.
- [22] P. Kingma Diederik and Ba Jimmy. “Adam: A Method for Stochastic Optimization”. In: International Conference for Learning Representations (2015).
- [23] N. Srivastava et al. “Dropout: a simple way to prevent neural networks from overfitting.” In: Journal of Machine Learning Research 15.1 (2014), 1929–1958.
- [24] K. Simonyan and A. Zisserman. “Very deep convolutional networks for large-scale image recognition”. In: International Conference for Learning Representations (2014).