

Detection and Classification of Immature Leukocytes for Diagnosis of Acute Myeloid Leukemia

Satvik Dasariraju

BACKGROUND: ACUTE MYELOID LEUKEMIA

- Acute Myeloid Leukemia (AML) is the deadliest type of leukemia, accounting for 11,000 deaths annually in the US [1,2]
- AML progresses quickly, and is fatal within months or weeks if not treated [2]
- Leukemia is characterized by overproduction of immature leukocytes, which are only found in peripheral blood under pathological conditions [3,4]

BACKGROUND: DIAGNOSIS

- Current method for diagnosis: microscopic examination and complete blood count to classify leukocytes [5, 6]
- Initial assessment of each blood smear takes **over 3 minutes** [7]
- Manual examination is inefficient and inaccurate (30-40% error rate) [3, 8]

Table 1. Interval Between Patient Help-seeking and Diagnosis of AML. All images (figures and tables) are created by the student researcher.

Country	Median Interval Between Help-seeking and Diagnosis of AML
United Kingdom	10 days [9]
Italy	14 days [10]
Nicaragua	29 days [10]

LITERATURE REVIEW

Table 2. Summary of previous studies on AML classification. All images (figures and tables) are created by the student researcher.

Year	Classification	Classifier	Performance	Limitation(s) and Obstacle(s)
2016	4 subtypes	Support Vector Machine	87% accuracy	<ul style="list-style-type: none"> • Small data set [8]
2017	3 leukocyte types	Support Vector Machine	80% accuracy	<ul style="list-style-type: none"> • Low accuracy [11]
2017	2 subtypes	k-Nearest Neighbor	67% accuracy	<ul style="list-style-type: none"> • Low accuracy • Extracted 3 features [12]
2019	3 immature leukocyte types	Random Forest	90% accuracy	<ul style="list-style-type: none"> • Low sensitivity • Small data set [13]
2019	Classification and detection of leukocytes in AML and healthy patients	Convolutional Neural Network	Less than 65% precision for most cell types in multi class prediction	<ul style="list-style-type: none"> • Low classification performance despite accurate detection • Imbalanced data set [3]

OBJECTIVES

1) To develop an algorithm, capable of accurate detection and classification of 4 types of immature leukocytes in AML cells

2) To calculate and identify the most important features for classification of leukocytes

MATERIALS

- Images were obtained from a publicly available data set in The Cancer Imaging Archive [14, 15]
- 1,070 images (mature and immature) were used for detection of immature leukocytes
- 613 images were used for classification of immature leukocytes into 4 types
- The project was coded with the python programming language and open source libraries [16-21]

METHODS: SEGMENTATION

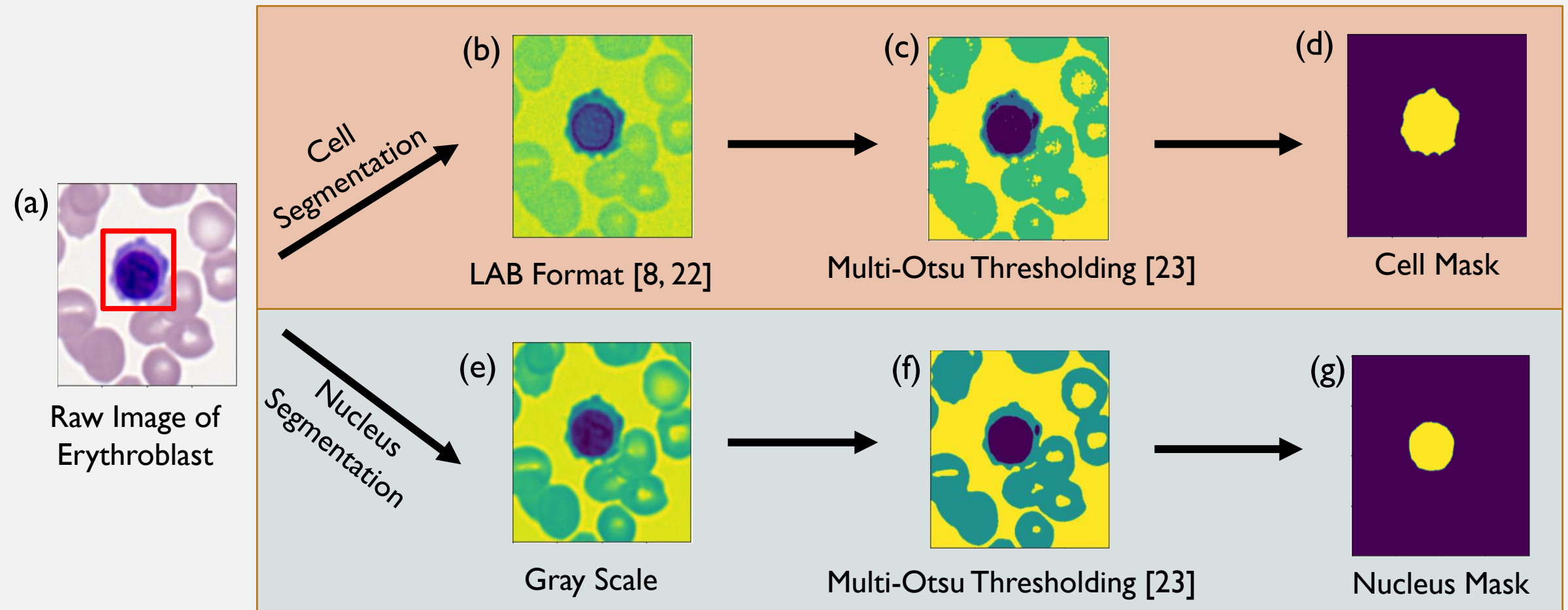


Figure 1. Flowchart of Segmentation Process. (a) The original image. (b) Conversion to LAB color space to differentiate cytoplasm from background cells. (c) Multi-Otsu thresholding to isolate the leukocyte. (d) Result: Binary mask of cell. (e) Conversion to grayscale to differentiate nucleus. (f) Multi-Otsu thresholding to isolate dark components of image. (g) Result: Binary mask of nucleus.

METHODS: FEATURE EXTRACTION

- From each image, 16 features were extracted and inputted to a features matrix [8, 22]
- **2 new proposed features** in this study from LAB color space:
 - Average of B Channel Intensity of Nucleus in LAB Color Space
 - Standard Deviation of B Channel Intensity of Nucleus in LAB Color Space

METHODS: MODEL TRAINING

- Random Forest classifier was chosen for imbalanced data and feature importance [25-27]
- Binary classification (immature, atypical vs. mature, typical) data was split 80-20 for training and testing
- Data for immature leukocyte classification was split 70-30 for training and testing
- Optimization was performed to improve classification

RESULTS & DISCUSSION: DETECTION OF IMMATURE LEUKOCYTES

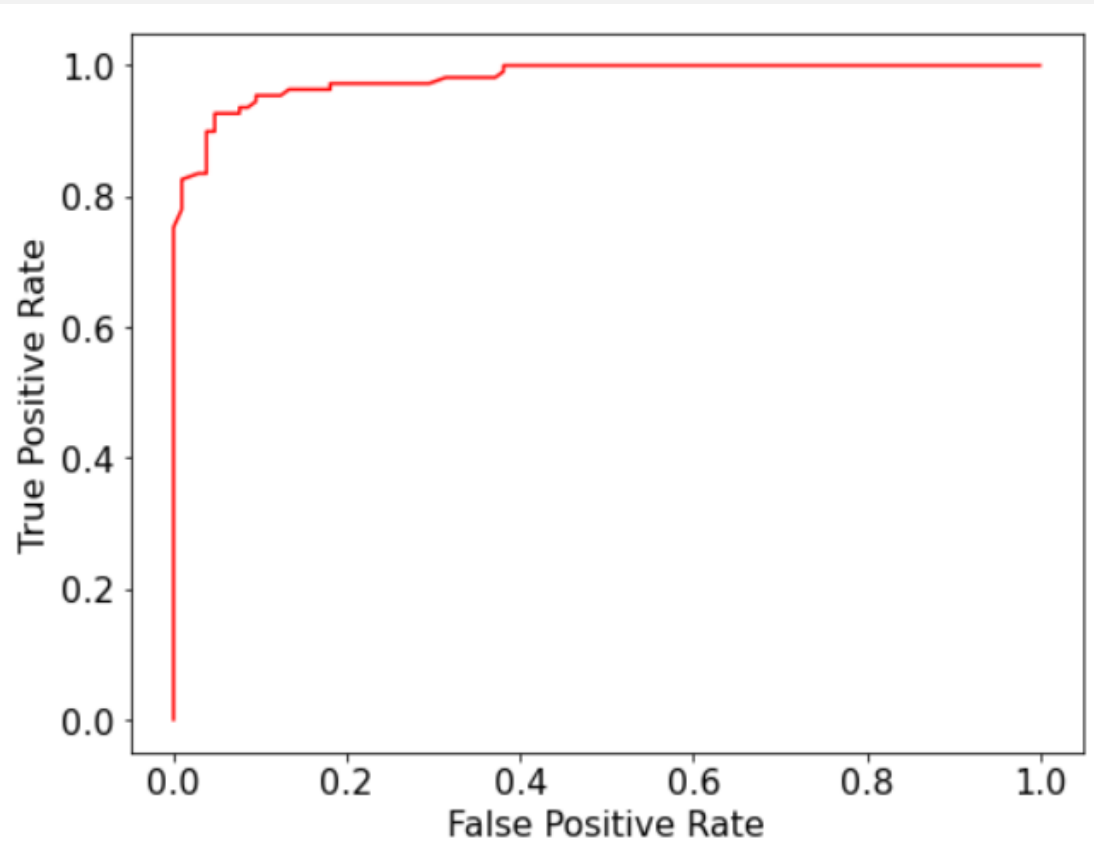


Figure 2. Receiver Operating Characteristic Curve for Binary Classification (Immature or Mature). Graph created with Matplotlib.

Table 3. Performance Metrics for Binary Classification. Results are on par with the current state of art [3].

Performance Metric	Score on Testing Set
Accuracy	92.99%
Precision	91.23%
Recall (Sensitivity)	95.41%
Specificity	90.48%
Area Under Curve of Receiver Operating Characteristic	0.9803

RESULTS & DISCUSSION: CLASSIFICATION OF IMMATURE LEUKOCYTES

Table 4. Performance Metrics for Immature Leukocyte Classification. Optimized model was optimized with weighted precision as the metric. Results are superior to previous studies on classification of AML leukocytes.

Model Type	Performance Metric	Score on Erythroblast Class	Score on Monoblast Class	Score on Promyelocyte Class	Score on Myeloblast Class
Initial Random Forest	Precision	100.00%	87.50%	62.50%	96.75%
	Recall (Sensitivity)	91.30%	100.00%	83.33%	94.44%
	Overall Accuracy	93.45%			
Optimized Random Forest	Precision	100.00%	77.78%	69.23%	97.56%
	Recall (Sensitivity)	91.30%	100.00%	75.00%	96.77%
	Overall Accuracy	93.45%			

RESULTS & DISCUSSION: CONFUSION MATRICES FOR CLASSIFICATION

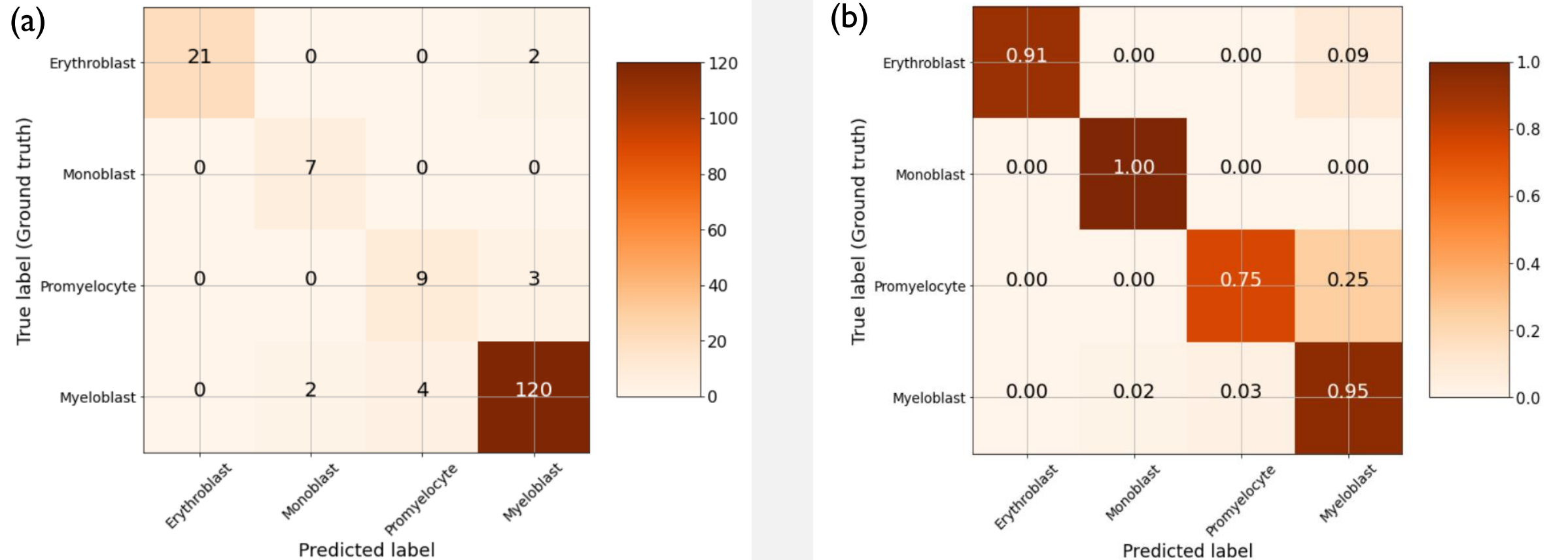


Figure 3. Multi Class Confusion Matrices for Optimized Model. (a) Absolute Confusion Matrix: Numbers refer to the number of images in a class classified with a given label. (b) Normalized Confusion Matrix: Numbers refer to the proportion of images in a class classified with a given label. Images created with scikit learn and matplotlib.

RESULTS & DISCUSSION: MOST IMPORTANT FEATURES

Table 5. Most Important Features for Random Forest Classifier. (a) Most Important Features for Detection of Immature Leukocytes. (b) Most Important Features for Classification of Immature Leukocytes.

(a) Detection

Feature	Gini Importance [27]
Nucleus to Cytoplasm Area Ratio [28]	0.2801
Area to Perimeter Ratio	0.1076
Nucleus Minor Axis Length [29]	0.0829
Nucleus Major Axis Length [29]	0.0803
Area [8]	0.0627

(b) Classification

Feature	Gini Importance [27]
Average Nucleus Color Intensity in B Channel	0.2532
Standard Deviation of Nucleus Color Intensity in B Channel	0.1853
Nucleus to Cytoplasm Area Ratio [28]	0.1765
Standard Deviation of Cytoplasm Color Intensity in B Channel [22]	0.0618
Average Cytoplasm Color Intensity in B Channel [22]	0.0571

CONCLUSIONS

- This project explored a vital, yet less researched computer diagnosis task [4, 30]
- An algorithm capable of accurate detection and precise classification of immature leukocytes was developed
- Nucleus to cytoplasm area ratio was established as an important morphological feature for detection and classification of immature leukocytes
- 2 new nucleus color features were displayed to be significant for classification

APPLICATIONS

- **The model can be used as an efficient support tool for pathologists to detect and classify immature leukocytes for the diagnosis of AML due to its efficiency and accuracy [31]**
- The features calculated to be the most important in this study can be used in future research
- The proposed new features (cytoplasm color intensity) can be used to elevate the performance of future models for leukocyte classification

FUTURE INVESTIGATIONS

- Improve the classification between similar cell types (eg. promyelocytes and myeloblasts)
- Developing algorithms to remove overlapping cells from blood smear images
- Establish additional morphological features for leukocyte classification [22]
- Develop systems that can be completely integrated into the current diagnosis methodology [32, 33]

ACKNOWLEDGEMENTS

This project was completed with the guidance and mentorship of Marc Huo (Stanford University) and Dr. Serena McCalla (Jericho High School) at iResearch Institute.



REFERENCES

- [1] “Acute Myeloid Leukemia - Cancer Stat Facts,” *SEER*. <https://seer.cancer.gov/statfacts/html/amyl.html> (accessed Jul. 22, 2020).
- [2] C. C. Kumar, “Genetic abnormalities and challenges in the treatment of acute myeloid leukemia,” *Genes Cancer*, vol. 2, no. 2, pp. 95–107, Feb. 2011, doi: 10.1177/1947601911408076.
- [3] C. Matek, S. Schwarz, K. Spiekermann, and C. Marr, “Human-level recognition of blast cells in acute myeloid leukaemia with convolutional neural networks,” *Nature Machine Intelligence*, vol. 1, no. 11, pp. 538–544, Nov. 2019, doi: 10.1038/s42256-019-0101-9.
- [4] S. Shafique and S. Tehsin, “Computer-Aided Diagnosis of Acute Lymphoblastic Leukaemia,” *Comput Math Methods Med*, vol. 2018, Feb. 2018, doi: 10.1155/2018/6125289.
- [5] “Blood Smear: MedlinePlus Medical Test.” <https://medlineplus.gov/lab-tests/blood-smear/> (accessed Jul. 22, 2020).
- [6] “Leukemia.” <https://www.hematology.org/443/education/patients/blood-cancers/leukemia> (accessed Jul. 22, 2020).
- [7] A. Adewoyin and B. Nwogoh, “PERIPHERAL BLOOD FILM - A REVIEW,” *Ann Ib Postgrad Med*, vol. 12, no. 2, pp. 71–79, Dec. 2014.
- [8] F. Kazemi, T. A. Najafabadi, and B. N. Araabi, “Automatic Recognition of Acute Myelogenous Leukemia in Blood Microscopic Images Using K-means Clustering and Support Vector Machine,” *J Med Signals Sens*, vol. 6, no. 3, pp. 183–193, 2016.
- [9] D. A. Howell *et al.*, “Time-to-diagnosis and symptoms of myeloma, lymphomas and leukaemias: a report from the Haematological Malignancy Research Network,” *BMC Blood Disorders*, vol. 13, no. 1, p. 9, Oct. 2013, doi: 10.1186/2052-1839-13-9.
- [10] C. De Angelis *et al.*, “The experience in nicaragua: childhood leukemia in low income countries—the main cause of late diagnosis may be ‘medical delay,’” *International Journal of Pediatrics*, Feb. 12, 2012. <https://www.hindawi.com/journals/ijpedi/2012/129707/> (accessed Jul. 28, 2020).
- [11] L. Bigorra, A. Merino, S. Alférez, and J. Rodellar, “Feature Analysis and Automatic Identification of Leukemic Lineage Blast Cells and Reactive Lymphoid Cells from Peripheral Blood Cell Images,” *J. Clin. Lab. Anal.*, vol. 31, no. 2, Mar. 2017, doi: 10.1002/jcla.22024.
- [12] E. S. Wiharto, S. Palgunadi, Y. R. Putra, and E. Suryani, “Cells identification of acute myeloid leukemia AML M0 and AML M1 using K-nearest neighbour based on morphological images,” in *2017 International Conference on Data and Software Engineering (ICoDSE)*, Nov. 2017, pp. 1–6, doi: 10.1109/ICoDSE.2017.8285851.
- [13] W. Wiharto, E. Suryani, and Y. R. Putra, “Classification of blast cell type on acute myeloid leukemia (AML) based on image morphology of white blood cells,” *TELKOMNIKA*, vol. 17, no. 2, p. 645, Aug. 2018, doi: 10.12928/telkomnika.v17i2.8666.
- [14] C. Matek, S. Schwarz, C. Marr, and K. Spiekermann, “A Single-cell Morphological Dataset of Leukocytes from AML Patients and Non-malignant Controls [Data Set],” in *The Cancer Imaging Archive*. doi: 10.7937/tcia.2019.36f5o9ld.
- [15] K. Clark *et al.*, “The Cancer Imaging Archive (TCIA): Maintaining and Operating a Public Information Repository,” *J Digit Imaging*, vol. 26, no. 6, pp. 1045–1057, Dec. 2013, doi: 10.1007/s10278-013-9622-7.
- [16] “The Python Language Reference — Python 3.8.5 documentation.” <https://docs.python.org/3/reference/> (accessed Jul. 23, 2020).

REFERENCES

- [17] S. van der Walt, S. C. Colbert, and G. Varoquaux, "The NumPy Array: A Structure for Efficient Numerical Computation," *Comput. Sci. Eng.*, vol. 13, no. 2, pp. 22–30, Mar. 2011, doi: 10.1109/MCSE.2011.37.
- [18] J. D. Hunter, "Matplotlib: A 2D Graphics Environment," *Comput. Sci. Eng.*, vol. 9, no. 3, pp. 90–95, 2007, doi: 10.1109/MCSE.2007.55.
- [19] W. McKinney, "Data Structures for Statistical Computing in Python," Austin, Texas, 2010, pp. 56–61, doi: 10.25080/Majora-92bf1922-00a.
- [20] S. van der Walt *et al.*, "scikit-image: image processing in Python," *PeerJ*, vol. 2, p. e453, Jun. 2014, doi: 10.7717/peerj.453.
- [21] R. Garreta and G. Moncecchi, *Learning scikit-learn: machine learning in Python : experience the benefits of machine learning techniques by applying them to real-world problems using Python and the open source scikit-learn library*. 2013.
- [22] N. Ghane, A. Yard, A. Talebi, and P. Nematollahy, "Segmentation of White Blood Cells From Microscopic Images Using a Novel Combination of K-Means Clustering and Modified Watershed Algorithm," *J Med Signals Sens*, vol. 7, no. 2, pp. 92–101, 2017.
- [23] P. Liao, T. Chen, and P. Chung, "A fast algorithm for multilevel thresholding," *Journal of Information Science and Engineering*, vol. 17, pp. 713–727, 2001.
- [24] J. Prinyakupt and C. Pluempitwiriyawej, "Segmentation of white blood cells and comparison of cell morphology by linear and naïve Bayes classifiers," *Biomed Eng Online*, vol. 14, Jun. 2015, doi: 10.1186/s12938-015-0037-1.
- [25] A. Parmar, R. Katariya, and V. Patel, "A Review on Random Forest: An Ensemble Classifier," in *International Conference on Intelligent Data Communication Technologies and Internet of Things (ICICI) 2018*, Cham, 2019, pp. 758–763, doi: 10.1007/978-3-030-03146-6_86.
- [26] M. Khalilia, S. Chakraborty, and M. Popescu, "Predicting disease risks from highly imbalanced data using random forest," *BMC Med Inform Decis Mak*, vol. 11, p. 51, Jul. 2011, doi: 10.1186/1472-6947-11-51.
- [27] G. Louppe, L. Wehenkel, A. Sutura, and P. Geurts, "Understanding variable importances in forests of randomized trees," in *Proceedings of the 26th International Conference on Neural Information Processing Systems - Volume 1*, Lake Tahoe, Nevada, Dec. 2013, pp. 431–439, Accessed: Jul. 23, 2020. [Online].
- [28] A. Mathur, A. S. Tripathi, and M. Kuse, "Scalable system for classification of white blood cells from Leishman stained blood stain images," *Journal of Pathology Informatics*, vol. 4, no. 2, p. 15, Jan. 2013, doi: 10.4103/2153-3539.109883.
- [29] K. Fukuma, V. B. S. Prasath, H. Kawanaka, B. J. Aronow, and H. Takase, "A Study on Nuclei Segmentation, Feature Extraction and Disease Stage Classification for Human Brain Histopathological Images," *Procedia Computer Science*, vol. 96, pp. 1202–1210, 2016, doi: 10.1016/j.procs.2016.08.164.
- [30] S. Shafique and S. Tehsin, "Acute Lymphoblastic Leukemia Detection and Classification of Its Subtypes Using Pretrained Deep Convolutional Neural Networks," *Technol. Cancer Res. Treat.*, vol. 17, p. 1533033818802789, 01 2018, doi: 10.1177/1533033818802789.
- [31] M. J. Feldman, E. P. Hoffer, G. O. Barnett, R. J. Kim, K. T. Famiglietti, and H. C. Chueh, "Impact of a Computer-Based Diagnostic Decision Support Tool on the Differential Diagnoses of Medicine Residents," *J Grad Med Educ*, vol. 4, no. 2, pp. 227–231, Jun. 2012, doi: 10.4300/JGME-D-11-00180.1.
- [32] K. Doi, "Computer-Aided Diagnosis in Medical Imaging: Historical Review, Current Status and Future Potential," *Comput Med Imaging Graph*, vol. 31, no. 4–5, pp. 198–211, 2007, doi: 10.1016/j.compmedimag.2007.02.002.
- [33] J. Shiraishi, Q. Li, D. Appelbaum, and K. Doi, "Computer-aided diagnosis and artificial intelligence in clinical imaging," *Semin Nucl Med*, vol. 41, no. 6, pp. 449–462, Nov. 2011, doi: 10.1053/j.semnuclmed.2011.06.004.