





Original Article

## A novel technique based on principal component analysis and multi-layer perceptron with genetic algorithm optimization for diagnosis of lung cancer

Ali Sharifi, PhD Student in Analytical Chemistry, Department of Chemistry, Lorestan University, Khorramabad, Iran  
ID Kamal Alizadeh, Associate Professor, Department of Chemistry, Faculty of Basic Science, Lorestan University, Khorramabad, Iran (\*Corresponding author). alizadeh.k@lu.ac.ir

### Abstract

**Background:** Lung cancer was known as primary cancers. Early detection of lung cancer reduces the length of treatment and spends a great deal of cost on the survival and survival of the individual. In recent years, the use of computer techniques in the use of data mining and intelligent algorithms has accelerated the early diagnosis of this cancer. The purpose of this paper is to evaluate the role of the new method based on Principal Component Analysis and Multi-Layer Perceptron with Genetic Algorithm optimization for Diagnosis of Lung Cancer.

**Methods:** In this study, the lung cancer dataset used was taken from the UCI machine learning database, including 32 patient records with 57 features. After performing its preprocessing steps, in the process of extraction of features and reduction of data dimensions, the main data of lung cancer were reduced to 17 characteristics using a basic component analysis. Then, in the classification step, these characteristics were reduced to multilayer perceptron by optimizing the genetic algorithm and the sensitivity and specificity of the model were studied according to the accuracy, sensitivity and Specificity. All analysis and synthesis were performed using the software of MATLAB and SPSS.

**Results:** For the proposed model, the results of the simulations were the mean of classification accuracy, sensitivity and specificity, respectively, 98.86, 98 and 99.16%.

**Conclusion:** The results on real data indicate that the proposed system is very effective in the diagnosis of lung cancer and can be used for clinical applications.

**Conflicts of interest:** None

**Funding:** Lorestan University

### Keywords

Lung cancer,  
Principal component analysis,  
Artificial neural network,  
Multilayer perceptron,  
Genetic algorithm

Received: 13/07/2019

Accepted: 30/11/2019

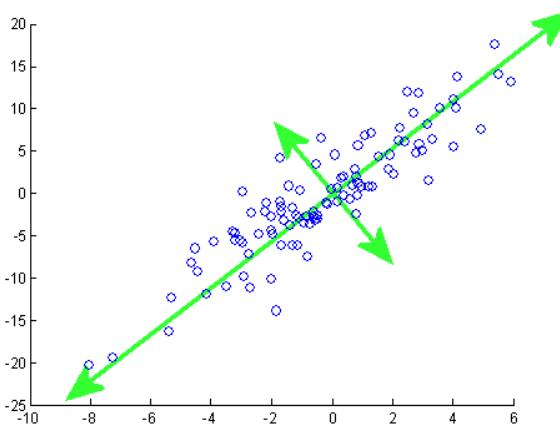
### Cite this article as:

Sharifi A, Alizadeh K. A novel technique based on principal component analysis and multi-layer perceptron with genetic algorithm optimization for diagnosis of lung cancer. Razi J Med Sci. 2019;26(10):48-56.

\*This work is published under CC BY-NC-SA 3.0 licence.







شکل ۲- تحلیل مؤلفه اساسی (محورهای جدید در جهت پرتوگیرمترین نقاط قرار دارند) (۹)

استفاده قرار گیرد. توابع فعالیت مختلفی به فراخور اسلوب مسئله در نرون‌ها مورد استفاده قرار می‌گیرد (۱۳). شبکه عصبی مصنوعی مورد استفاده در این مقاله، از سه لایه‌ی ورودی، خروجی و پردازش تشکیل می‌شود. هر لایه شامل گروهی از سلول‌های عصبی (نuron) است که عموماً با کلیه‌ی نuron‌های لایه‌های دیگر در ارتباط هستند، مگر این که کاربر ارتباط بین نuron‌ها را محدود کند؛ ولی نuron‌های هر لایه با سایر نuron‌های همان لایه، ارتباطی ندارند (۱۴). قبل از به کار بردن مدل شبکه عصبی، وزن‌ها و اریب‌های ارتباط دهنده نuron‌های شبکه تعیین می‌شوند. به همین منظور تمام داده‌ها برای تدوین ساختار مدل به سه گروه تقسیم می‌شوند (۱۵). اولین گروه داده‌ها، به عنوان داده‌های آموزش برای تعیین وزن‌ها و اریب‌های شبکه به کار می‌رود. دومین گروه از داده‌ها، که داده‌های اعتباری نامیده می‌شوند، برای ارزیابی نتایج مرحله آموزش و تصمیم‌گیری درخصوص توقف آموزش شبکه استفاده می‌شوند. تعیین دقت مدل، و یا به عبارتی آزمون مدل، با استفاده از سومین گروه داده‌ها، یعنی داده‌های آزمون که در تدوین مدل استفاده نشده‌اند، انجام می‌شود. در شکل ۳ ساختار کلی شبکه عصبی مصنوعی مبتنی بر ساختار پرسپترون چند لایه نشان داده شده است.

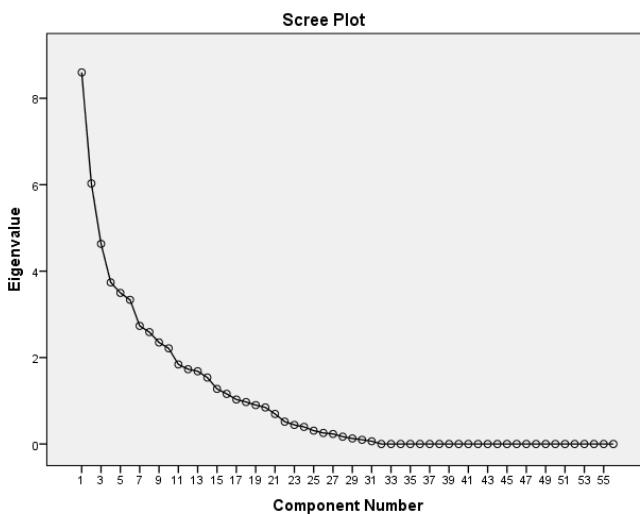
الگوریتم ژنتیک (GA: Genetic Algorithm) یک روش برای حل مسائل بهینه‌سازی می‌باشد که اساس آن بر انتخاب، بقاء و تکامل در محیط‌های طبیعی استوار است. روش الگوریتم ژنتیک را می‌توان برای انواع

کوواریانس داده‌ها، بیشترین مقدار ویژه بدست آمده دارای بالاترین اهمیت بردار ویژه در تحلیل مؤلفه‌های اساسی است.

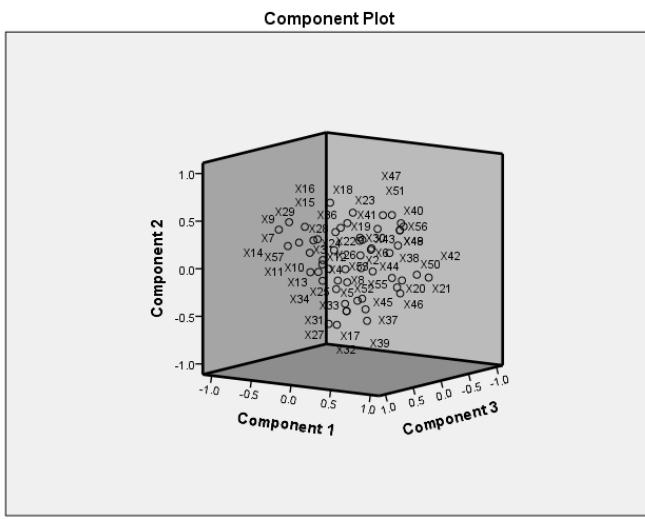
با استفاده از تحلیل مؤلفه اساسی در این مقاله به استخراج متغیرهای مهم (به شکل مؤلفه) از مجموعه متغیرهای موجود در مجموعه داده می‌پردازیم تا به ثبت اطلاعات بیشتر با تعداد کمتری از متغیرها کمک نماییم. بدین شکل، بصری‌سازی داده‌ها نیز معنادارتر می‌شود. تحلیل مؤلفه اساسی هنگامی که با داده‌های دارای سه یا تعداد بیشتری بعد سروکار داشته باشید، کاربرد پذیرتر است (۱۱) (شکل ۲). این مفهوم برای داده‌های دو بعدی برای این پژوهش نشان داده شده است.

پردازش داده و طبقه‌بندی: شبکه عصبی مصنوعی (ANN: Artificial Neural Network) برای پردازش اطلاعات می‌باشند که با تقلید از شبکه‌های عصبی بیولوژیکی مثل مغز انسان ساخته شده‌اند. عنصر کلیدی این الگو ساختار جدید سیستم پردازش اطلاعات آن می‌باشد و از تعداد زیادی عناصر (نuron) با ارتباطات قوی داخلی که هماهنگ با هم برای حل مسائل مخصوص کار می‌کنند تشکیل شده‌اند. یکی از شبکه‌های عصبی پرکاربرد شبکه عصبی پرسپترون چند لایه (MLP: Multi-Layer Perceptron) با روش یادگیری پس‌انتشار است که در صورت انتخاب صحیح ساختار داخلی، قادر است هر نوع سیستم غیر خطی را مدل‌بندی نماید. یک شبکه پرسپترون چند لایه می‌تواند به سادگی با تعریف اوزان و توابع مناسب مورد





شکل ۴- نمودار اسکری گراف برای تعیین تعداد عامل



شکل ۵- نمودار سه بعدی پراکنش متغیرها نسبت به سه عامل اول استخراج شده

معماری‌های مختلف شبکه عصبی برای حصول بهترین کلاسیندی موردن بررسی قرار گرفت و دقت، حساسیت و صحت به دست آمده با استفاده از بهینه‌سازی الگوریتم ژنتیک برای بهترین معماری به ترتیب برابر ۹۸/۶۵، ۹۸/۱۶ درصد بدست آمد. همان‌طور که در جدول ۱ مشاهده می‌کنید انتخاب دقیق بهترین معماری نقش عمده‌ای در افزایش دقت، حساسیت و صحت تشخیص شبکه دارد.

### بحث و نتیجه‌گیری

در این مطالعه سعی شده است برای بهبود تشخیص سرطان ریه از روش جدید بر پایه تحلیل مؤلفه اساسی جهت استخراج ویژگی و کاهش ابعاد داده با حفظ

تصادفی (Randomization) ردیف‌ها روی صفحه گسترده انجام گردید. پارامترهای شبکه عصبی، شامل متغیرهای ورودی و تعداد نمونه‌ها در لایه پنهان، با الگوریتم ژنتیک بهینه‌یابی شدند که در مرحله آزمون شبکه عصبی، نتایج برآورد شده توسط مدل، یا داده‌های خروجی بهترین برآذش را با مؤلفه‌های اصلی داشته باشند. در این مرحله از ۷۰ درصد داده‌های پیش‌پردازش شده توسط PCA که شامل هفده مؤلفه اول بود، جهت آموزش شبکه عصبی مصنوعی با استفاده از الگوریتم ژنتیک استفاده شده است. در مرحله بعد ۳۰ درصد داده‌های پیش‌پردازش شده توسط الگوریتم PCA به صورت بردار به شبکه عصبی مصنوعی پیاده‌سازی شده در نرم‌افزار اعمال گردید. همچنین



