

ارائه یک سیستم پیشنهادی به منظور کشف تقلب در متون فارسی با پیچیدگی ساده بر اساس تکنیک‌های پردازش تصویر و سیستم‌های هوشمند مصنوعی

محمد فیوضی^{۱*}، جواد حدادنیا^۲

^۱ بخش مهندسی پزشکی، دانشکده مهندسی برق و کامپیوتر، دانشگاه حکیم سبزواری، سبزوار، ایران. mohammad.fiuzy@yahoo.com

^۲ بخش مهندسی پزشکی، دانشکده مهندسی برق و کامپیوتر، دانشگاه حکیم سبزواری، سبزوار، ایران. haddadnia@sttu.ac.ir

چکیده

با توجه به گسترش روزافزون تقلب^۱ در متون مختلف از قبیل بانکداری، ارتباطات، تجارت و در سال‌های اخیر تکنیک‌های گوناگونی برای کشف تقلب ارائه شده است. تشخیص تقلب در متون به طور چشمگیری باعث کاهش این پدیده می‌شود. پردازش تصویر، دقیق، ارزان و موثر برای تشخیص این پدیده است. شاید بتوان از مهم ترین مشکلات روش‌های تشخیص برابری ظاهر اسناد و متون را دانست، که گاه با چشم غیر مسلح قابل تشخیص نیست. به همین خاطر امروزه تلاش می‌شود تا تشخیص نهایی به صورت هوشمند و ماشینی انجام گیرد. متأسفانه در این روش‌ها عدم نظر داشتن ویژگی‌های مناسب در آزمایشات، روش‌های ماشینی و هوشمند مساله‌ای مهم است حتی منجر به عدم تشخیص صحیح این پدیده می‌شود. در این تحقیق سعی شد تا با استفاده از ترکیب روش‌های هوشمند و پردازش تصویر به منظور افزایش کیفیت تصاویر متون فارسی، الگوریتم‌های طبقه بندی^۳ (FCM) به منظور جدا سازی قسمت‌هایی از متون مورد نظر، الگوریتم تحلیل مولفه‌های اصلی^۴ (PCA) بمنظور کاهش ویژگی‌ها و شبکه‌های عصبی مصنوعی^۵ (NN) برای مدل‌سازی، شناسایی ساختار و شناسایی پارامترهای متون دارای تقلب روشی نو در جهت تشخیص صحیح و دقیق این پدیده در متون فارسی پیشنهاد شده است. روش پیشنهادی بر روی بانک تهیه شده شامل ۷۶ نمونه (دستخط) (۴۵ نمونه متن بدون تقلب و ۳۱ نمونه متن دارای تقلب) پیاده شد، که در نهایت به دقت شناسایی " بیش از ۹۹٪" دست یافت. که در نوع خود با توجه به هوشمند، غیر تهاجمی، سریع و دقیق بودن می‌تواند به عنوان سیستم تشخیص کارآمدی به منظور کشف تقلب در متون فارسی با پیچیدگی ساده معرفی، بررسی و مورد استفاده قرار گیرد.

کلمات کلیدی

متون فارسی، تقلب، ویژگی، پردازش تصویر، هوش مصنوعی.

۱- مقدمه

می‌باشد. همه ساله تحقیقات گسترده‌ای برای ارایه سیستم‌های OCR^۷ با کارایی و سرعت بهتر انجام می‌شود. یکی از اساسی ترین مراحل در تحلیل اسناد و تشخیص متون، بازشناسی قلم است، که شامل تعیین مواردی چون زبان متن، اندازه قلم، استیل قلم و نوع قلم می‌تواند کارایی سیستم درک سند را به مراتب بالا ببرد چرا که با استفاده از اطلاعات قلم می‌توان پارامترهای سیستم را بطور خاصی تنظیم کرد [8]. در یک تحقیق، از ویژگی‌های گشتاوری و طبقه بندی کننده نیز برای بازشناسی حروف دستنویس فارسی استفاده شده است [9]. روش‌های یادگیری آماری زیادی در زمینه دسته بندی متن‌ها در سال‌های اخیر بکار برده شده است، که شامل مدل‌های برگشتی^۸، دسته بندی نزدیکترین همسایه^۹ (KNN)، شبکه‌های بیز^{۱۰}، درخت تصمیم گیری، الگوریتم‌های یادگیری بر اساس قانون، شبکه‌های عصبی مصنوعی^{۱۱} و تکنیک‌های یادگیری استنتاجی^{۱۲} می‌باشد [10, 11, 12]. ظهور و کشف پدیده تقلب در متون از مسائل و

بازشناسی حروف و ارقام دستنویس همواره یکی از موضوعات مورد علاقه برای تحقیق بوده است [1-4]. در آینده مستندات الکترونیکی به عنوان اصلی‌ترین ابزار ارتباطات نوشتاری مطرح خواهند شد و کتاب‌ها و مجلات کاغذی بخشی از تاریخ خواهند بود [5]. آنچه امروزه از اهمیت بسیار زیادی برخوردار گردیده، کمبود یا نبود اطلاعات نیست بلکه کمبود روش‌هایی در جهت استخراج و بهره‌برداری از اطلاعات در دسترس، به صورت مطلوب است. یک ویراستار انسانی، تنها به وسیله دنبال کردن دقیق همه منابع متنی می‌تواند متوجه وقوع یک مساله جدید شود و یا اسناد را دسته‌بندی کند. این روش کار برای حجم بالایی اطلاعات در سیستم‌های اطلاعاتی با پیچیدگی زیاد، نامناسب است [6]. روش‌های پویا به دلیل دسترسی به ویژگی‌های زمانی از دقت و حساسیت بالاتری در تصدیق افراد برخوردار هستند [7]. بازشناسی نوری حروف یکی از مباحث مهم بازشناسی الگو