

چکیده:

پاکسازی داده ها از خطاها و نویز یکی از بخش های اصلی فرایند نگهداری داده محسوب می گردد. پاکسازی به توالی عملیات های انجام شده با هدف بهبود کیفیت کلی مجموعه های داده ها اتلاق می شود. در واقع پاکسازی داده کاهش و حذف نسخه های کپی در مجموعه داده می باشد. مسئله ای در برنامه های کاربردی پایگاه داده رخ می دهد و زمانی بدتر می شود که داده های منابع بخواهند با هم تلفیق شوند. از این رو ، پاکسازی داده به عنوان بخش اصلی فرایند یکپارچگی داده شناخته شده است.

هنوز مسائل و چالش های زیادی از جمله تصحیح خطا و حل تعارض، نگهداری داده های پاکسازی شده و پاکسازی داده در محیط های مجازی یکپارچه در زمینه پاکسازی داده وجود دارد. برای حل چالش های پیش رو پاکسازی داده توسعه یک چارچوب مناسب نیز احساس می شود.

از این رو در این تحقیق یک روش آماده سازی و پاکسازی داده مبتنی بر خوشه بندی ارایه شده است که در آن ابتدا داده های خام را به دلیل اینکه به سرعت به پاک سازی و نرمال سازی های آتی داده کمک میکند، خوشه بندی نموده سپس عملیات پاک سازی داده ها و در مرحله آخر عمل تبدیل داده انجام می شود. در این تحقیق، آزمایش ها در بدترین حالت میزان نویز، در نظر گرفته شده و با این وضعیت آزمایش پیش برده شده، تا بهترین نتیجه برای روش پیشنهادی بدست آید. به طور کلی روش پیشنهادی یک چارچوب پاکسازی داده جهت بهبود صحت روش های داده کاوی بر روی داده های پاکسازی شده را نشان می دهد و مزیت اصلی روش ارایه شده این است که نیاز به بروزرسانی در محیط های پویا را نداشته و با قراردادن هر داده جدید در دسته مناسب همیشه منبع را به روز نگه می دارد. در روش پیشنهادی از مجموعه داده تشخیص حروف با 20000 داده استفاده شده است که بعد از اعمال گام های روش پیشنهادی بر روی مجموعه داده، دو مجموعه پاکسازی شده با روش پیشنهادی مورد مقایسه قرار گرفته و طبق نتایج بدست آمده، روش پیشنهادی الگوریتم طبقه بندی Kstar را 25 درصد بهبود داده است.

واژه های کلیدی: پاکسازی، خوشه بندی، داده، نویز و خطا