

تشخیص موضوع یکی از مسائل حوزه‌ی پردازش زبان طبیعی است که در سال‌های اخیر مورد توجه قرار گرفته و پژوهش‌های گوناگونی بر روی آن انجام شده است. در این مسئله هدف اصلی خوشه‌بندی اسناد متنی در دسته‌های مختلف است به نحوی که اسناد هر خوشه از نظر موضوعی به هم نزدیک باشند. راه حل‌های مختلفی برای این مسئله ارائه شده است که بخش قابل توجهی از آن‌ها از الگوریتم‌های خوشه‌بندی مانند K-means برای حل این مسئله استفاده می‌کنند. علاوه بر روش‌های مبتنی بر خوشه‌بندی اسناد در برخی از پژوهش‌ها از روش‌های مدل‌سازی موضوعی برای حل این مسئله استفاده شده است. در این پژوهش دو روش تشخیص موضوع ارائه می‌شود که مبتنی بر الگوریتم خوشه‌بندی K-means عمل می‌کنند و با ارائه‌ی روش‌هایی برای تعیین مراکز اولیه مناسب برای این الگوریتم کیفیت آن را بهبود می‌دهند. با توجه به حساسیت قابل توجه الگوریتم K-means به انتخاب مراکز اولیه، در این پژوهش ابتدا این میزان از حساسیت را بصورت عملی در مسئله تشخیص موضوع نشان می‌دهیم و سپس دو روش برای انتخاب هوشمندانه مراکز اولیه در الگوریتم K-means ارائه می‌شود که باعث بهبود کیفیت در مسئله‌ی تشخیص موضوع می‌شوند. در روش پیشنهادی اول برای انتخاب مراکز اولیه ابتدا گراف اخبار را ساخته و بر روی آن الگوریتم DivRank را که یک روش گام‌برداری تصادفی تقویتی است اجرا می‌کنیم تا با استفاده از آن مراکز اولیه مناسب را شناسایی کنیم. در روش پیشنهادی دوم نیز یک الگوریتم ارائه شده است که با بهره‌گیری از مدل‌سازی موضوعی (LDA (Latent Dirichlet Allocation اقدام به انتخاب هوشمندانه‌ی مراکز اولیه می‌کند و در نهایت با داشتن مراکز اولیه مناسب اسناد را خوشه‌بندی می‌کنیم. در روش‌های ارائه شده برای محاسبه فاصله اسناد از توزیع موضوع حاصل از LDA آن‌ها استفاده شده است. در نهایت به انجام آزمایش بر روی سه مجموعه دادگان مختلف پرداختیم، آزمایش‌ها نشان می‌دهند که استفاده از روش‌های ارائه شده در این پژوهش باعث بهبود چشم‌گیر نتایج نسبت به روش LDA در دو مجموعه از سه مجموعه دادگان می‌شود. همچنین نتایج نشان می‌دهند که مراکز اولیه حاصل از روش‌های پیشنهادی برای انتخاب مراکز اولیه در مقایسه با روش‌های انتخاب مراکز اولیه تصادفی و روش K-mean++ در دو مجموعه دادگان همیشه مناسب‌تر بوده و احتمال بهتر بودن مراکز اولیه انتخابی در مجموعه دادگان دیگر مورد آزمایش بیش از ۷۰ درصد است.