

## چکیده:

در حال حاضر تجزیه و تحلیل داده های حجیم نقش مهمی در اکثر زمینه های علمی و پژوهشی ایفا می کند. با وجود این که در پژوهش های اخیر، صورت مساله بعضا متداول بوده و قبلا در کار های پیشین معرفی شده، اما پدیده عظیم داده<sup>۱</sup> چالش های جدیدی را مطرح می کند که منجر به باز نگری در کارهای قبلی خواهد شد. یک نمونه از کاربرد های آن تحلیل سری های زمانی<sup>۲</sup> به صورت بلادرنگ در حوزه مراقبت های درمانی است. مواقعی که داده های ورودی طولانی باشند، رویکرد های سنتی برای ردیابی الگو<sup>۳</sup> های پر تکرار، وقت گیر و ناکارآمد هستند. در آن سال ها که این مدل ها پیشنهاد میشد، صورت مساله تک بعدی<sup>۴</sup> و تعداد متغیر ها کم بود یا محققین با داده های کم حجم و کوتاه مدت سر و کار داشتند. برخی از ویژگی های این نوع مسایل که استفاده از روش های پردازش عظیم داده را تشویق می کند، عبارت اند از: قابلیت شکسته شدن داده های ورودی به تکه های کوچکتر، پراکندگی<sup>۵</sup> داده ها، ... امروزه مدل های بهبود یافته و ابزار های پردازش عظیم داده باید قادر باشند بر این چالش ها غلبه کنند.

به منظور تشخیص بیماری، حسگر های متعددی در نقاط مختلف بدن مریض تعبیه و داده ها جمع آوری می شود. مجموع آن ها جدولی از داده های چند بعدی<sup>۶</sup> تولید می کند. هر بُعد یک جنبه<sup>۷</sup> از مشخصات پزشکی فرد بیمار را بازگو می کند. متخصص با مقایسه ابعاد مختلف از این داده ها تلاش می کند ارتباطی بین آن ها پیدا کند. مثلا بیمار با دیدن عکسی خاص، پلک چشم راستش می زند و هم زمان ضربان قلب اش تندتر می شود. متخصص با کنار هم قرار دادن سری های زمانی حاصل از این چند حسگر و جست و جوی موتیف<sup>۸</sup> های مشابه در آن ها، بهتر می تواند برای درمان شخص تصمیم بگیرد. معمولا برای بررسی مجموعه داده های حاصل از این حسگر ها از مدل های طبقه بندی با نظارت<sup>۹</sup> استفاده می شود. اگر این مدل ها در مرحله یادگیری، فقط از نتایج خروجی مربوط به یک بیمار به عنوان مجموعه داده های آموزشی<sup>۱۰</sup> استفاده کند، مدل مربوطه بیش از حد بیش برآزش<sup>۱۱</sup> خواهد شد. این در حالی است که هر چه قدر مورد های بیشتری را معاینه کنیم و همین طور ابعاد متنوع تری را زیر نظر بگیریم، مدل پیشنهادی ما دقیق تر قادر به تشخیص و طبقه بندی نوع بیماری خواهد بود.

در این پژوهش، با استفاده از پنجره های لغزان<sup>۱۲</sup> زیر دنباله های یک سری زمانی از نوع ولگشت<sup>۱۳</sup> را از هم جدا و هر زیر دنباله را با بکار گیری مدل SAX<sup>۱۴</sup> به موتیف های کوچکتر شکسته و برچسب گذاری می کنیم. سپس با استفاده از ماتریس تصادم<sup>۱۵</sup> و داده کاوی در رشته<sup>۱۶</sup> ای از این حروف تلاش می کنیم تا الگویی های مشابه که در جاهای مختلف سری تکرار شده اند پیدا کنیم. از آن جایی که این زیر دنباله های توزیع شده را می توان به قسمت های کوچکتر تقسیم کرد، از روش نگاشت-کاهش<sup>۱۷</sup> در زیر ساخت Hadoop<sup>۱۸</sup> استفاده می شود. داده های مورد نظر برای آزمایش از سیگنال های ECG<sup>۱۹</sup> واقعی در مجموعه داده UCR<sup>۲۰</sup> دریافت شده است.

<sup>1</sup> Big Data

<sup>2</sup> Time Series

<sup>3</sup> Pattern

<sup>4</sup> One Dimensional

<sup>5</sup> Distributed Data

<sup>6</sup> Multi-Dimensional

<sup>7</sup> Feature

<sup>8</sup> Motif

<sup>9</sup> Supervised Classification

<sup>10</sup> Learning

<sup>11</sup> Overfitting

<sup>12</sup> Sliding Window

<sup>13</sup> Random Walk

<sup>14</sup> Symbolic Aggregate approXimation

<sup>15</sup> Collision Matrix

<sup>16</sup> String

<sup>17</sup> Map-Reduced

<sup>18</sup> [hadoop.apache.org/](http://hadoop.apache.org/)

<sup>19</sup> Electrocardiography

<sup>20</sup> [http://www.cs.ucr.edu/~eamonn/time\\_series\\_data/](http://www.cs.ucr.edu/~eamonn/time_series_data/)