

بِسْمِ اللَّهِ الرَّحْمَنِ الرَّحِيمِ



دانشگاه تبریز

دانشگاه تبریز

دانشکده‌ی ریاضی

گروه علوم کامپیوتر

پایان نامه

جهت دریافت درجه کارشناسی ارشد در رشته علوم کامپیوتر

عنوان

بررسی الگوریتم‌های خوشه‌بندی بارویکرد بهبود در وقت آنها

استاد راهنما:

دکتر محمد رضا فیضی درخشی

استاد مشاور:

میر محمد اتفاق

پژوهشگر:

فاطمه محمودلو

شهریور ۱۳۹۲

نام خانوادگی دانشجو: محمودلو	نام: فاطمه
عنوان پایان نامه: بررسی الگوریتم‌های خوشه‌بندی با رویکرد بهبود در دقت آنها	
استاد راهنما: دکتر محمدرضا فیضی درختی	
استاد مشاور: دکتر میر محمد اتفاق	
مقطع تحصیلی: کارشناسی ارشد	رشته: علوم کامپیوتر
دانشگاه: دانشگاه تبریز	دانشکده: علوم ریاضی
تاریخ فارغ التحصیلی: ۱۳۹۲/۰۶/۱۶	تعداد صفحه: ۸۲
کلید واژه‌ها: خوشه‌بندی، الگوریتم جنگل و الگوریتم‌های فراابتکاری	
<p>چکیده: خوشه‌بندی قرار دادن داده‌ها در گروه‌هایی است، که اعضای هر گروه از زاویه‌ی خاصی شبیه یکدیگرند، بطوریکه شباهت درون هر خوشه حداکثر و شباهت بین داده‌های درون خوشه‌های متفاوت، حداقل باشد خوشه‌بندی فرآیند یادگیری . بدون ناظر است که از قبل هیچ دانشی درباره برچسب داده‌ها ندارد. روش‌های زیادی برای خوشه‌بندی وجود دارد که می‌توان آنها را به روش‌های افرازبندی و سلسله مراتبی تقسیم کرد. در این پایان نامه از الگوریتم فراابتکاری جدیدی به نام الگوریتم جنگل که از طبیعت جنگل الهام گرفته شده، برای خوشه‌بندی استفاده شده است. در این الگوریتم، برای رهایی از بهینه‌های محلی، تغییراتی در قسمت دانه پراکنی محلی انجام شد و نتایج خوبی بدست آمد. برای ارزیابی، روش ارائه شده را روی دو سری مجموعه داده مورد آزمایش قرار دادیم. سری اول شامل داده‌های استاندارد و سری دوم شامل مجموعه داده‌های حقیقی استخراج شده جهت عیب یابی سیستم‌های دوار مکانیکی، است. همچنین از روش‌های شناخته شده‌ای نظیر روش‌های GA، PSO، ACO، CAS_C و K-means برای مقایسه نتایج بدست آمده استفاده شده است. این الگوریتم برای داده‌های استاندارد کمترین مقدار مجموع مجذور فاصله درون خوشه‌ها را بدست آورد. این مقدار در داده iris برابر ۹۶.۶۵۵۷ با انحراف معیار ۰.۰۰۱، در wine برابر ۱۶۲۹۲.۴۱۰۰ با انحراف معیار ۵.۳۴۷۴ و در داده‌های Glass برابر ۲۱۰.۵۳۴۰ با انحراف معیار ۱.۸۰۲۹ است که در مقایسه با دیگر روش‌ها، روش فوق نتیجه مطلوبی را تولید می‌کند. همچنین در این روش فاصله درون خوشه‌ها کاهش و فاصله بیرون خوشه‌ها افزایش پیدا کرده است. علاوه بر این روش فوق درصد خطای خوشه‌بندی پایین‌تری نسبت به دیگر روش‌ها دارد.</p>	

ماحصل آموختیم را تقدیم می‌کنم به آنان که مرا آسانی‌شان آرام بخش آلام زمینی ام است

به استوارترین تکیه‌گاه، دستان پر مهر پدرم

و به دلسوزترین نگاه، چشمان شورانگیز مادرم

که هرچه آموختیم در مکتب عشق شاست و هرچه بلو شتم قطره‌ای از دریای بی‌کران عشق تان را پاس گزار نیستم

باشد که حاصل تلاشیم نسیم کوزه غبار خشکیتان را بروداید

و تقدیم به خواهرم

به همسفر مهربان زندگی ام، شهربانوی نازنینم

که با هم آغاز کردیم، در کنار هم آموختیم و به امید هم به آینده چشم می‌دوزیم. قلمم لبریز از عشق به توست و خوشبختی ات منتهای آرزویم.

فهرست مطالب

شماره‌ی صفحه	عنوان
۱	فصل اول: مقدمه
۲	۱-۱ مقدمه
۲	۲-۱ شرح مسئله
۳	۳-۱ ساختار پایان نامه
۵	فصل دوم: بررسی مفاهیم و الگوهای اولیه
۶	۱-۲ مقدمه
۶	۲-۲ خوشه‌بندی
۶	۲-۲-۱ تعریف خوشه‌بندی
۸	۲-۲-۲ کاربردهای خوشه‌بندی
۹	۲-۲-۳ تفاوت خوشه‌بندی با طبقه‌بندی
۹	۲-۲-۴ هدف از خوشه‌بندی
۱۰	۲-۲-۵ مراحل انجام خوشه‌بندی
۱۲	۳-۲ روش‌های خوشه‌بندی
۱۳	۲-۱ روش‌های سلسله‌مراتبی
۱۶	۲-۳ روش‌های افرازبندی
۲۵	۴-۲ اعتبارسنجی خوشه‌ها
۲۶	۲-۴-۱ اعتبارسنجی خارجی خوشه‌ها
۲۷	۲-۴-۲ فاصله‌های درون خوشه‌ها و بیرون خوشه‌ها
۲۸	۲-۴-۳ اعتبارسنجی درونی خوشه‌ها
۱۷	۵-۲ الگوریتم‌های مورد استفاده در خوشه‌بندی
۱۷	۲-۵-۱ الگوریتم <i>k-means</i>
۱۸	۲-۵-۲ الگوریتم ژنتیک
۱۹	۲-۵-۳ الگوریتم <i>psa</i>
۲۰	۲-۵-۴ الگوریتم جستجوی گرانشی (<i>GSA</i>)
۲۲	۲-۵-۵ الگوریتم مورچگان
۲۴	۲-۵-۶ الگوریتم جنگل (<i>FA</i>)
۲۵	۷-۲ جمع‌بندی فصل
۳۰	فصل سوم: بررسی روش‌های موجود

ب

۱-۳ مقدمه ۳۱

۲-۳ روش‌های مرکز محور ۳۱

۱-۲-۳ بررسی روش‌های مرکز محور ۳۱

۳-۳ روش‌های فرابتکاری ۴۰

۱-۳-۳ بررسی روش‌های فرابتکاری ۴۰

۴-۳ جمع بندی فصل ۵۱

۵۳ فصل چهارم: خوشه‌بندی با استفاده از الگوریتم جنگل (FAC)

۱-۴ مقدمه ۵۴

۲-۴ خوشه‌بندی بر اساس الگوریتم جنگل ۵۴

۳-۴ تنظیم پارامترهای الگوریتم جنگل ۵۷

۴-۴ بررسی نتایج حاصل از اجرای الگوریتم جنگل و مقایسه آن با دیگر الگوریتم‌ها ۵۷

۱-۴-۴ معرفی داده‌های استفاده شده و نتایج شبیه‌سازی مربوط به آن ۵۸

۵-۴ جمع‌بندی فصل ۶۸

۶۹ فصل پنجم: نتیجه‌گیری و پیشنهادهای آینده

۱-۵ نتیجه‌گیری ۷۰

۲-۵ پیشنهادهای آینده ۷۱

۷۳ مراجع

فهرست اشکال

شماره ی صفحه	عنوان
۷	شکل ۱-۲: خوشه بندی نمونه های اولیه.....
۱۱	شکل ۲-۲: مراحل انجام خوشه بندی [۶].....
۱۳	شکل ۳-۲: مراحل انجام خوشه بندی سلسله مراتبی.....
۱۴	شکل ۴-۲: خوشه بندی single linkage.....
۱۴	شکل ۵-۲: خوشه بندی complete linkage.....
۱۵	شکل ۶-۲: خوشه بندی average linkage.....
۱۶	شکل ۷-۲: گام های الگوریتم تقسیم شونده و متراکم شونده.....
۱۹	شکل ۸-۲: روند کلی الگوریتم ژنتیک.....
۲۰	شکل ۹-۲: روند کلی الگوریتم PSO.....
۲۵	شکل ۱۰-۲: فلوچارت الگوریتم FA [۳].....
۳۶	شکل ۱-۳: الگوریتم ترکیبی GSO [۱].....
۳۷	شکل ۲-۳: عمل برش مربوط به ژنتیک GSOKHM [۱].....
۳۸	شکل ۳-۳: نمایش کروموزوم [۳].....
۴۵	شکل ۴-۳: نمایش ذره در الگوریتم GSA.....
۵۰	شکل ۵-۳: الگوریتم CAS_A [۴۴].....

فهرست جداول

عنوان	شماره ی صفحه
جدول ۱-۳: داده های مورد آزمایش [۴۰].....	۳۴
جدول ۲-۳: مقایسه الگوریتم K-meanse بهبود یافته با الگوریتم K-means [۴۰].....	۳۵
جدول ۱-۴: پارامترهای مربوط به الگوریتم FAC.....	۵۷
جدول ۲-۴: مراکز خوشه بدست آمده با اجرای الگوریتم FAC روی مجموعه داده Iris.....	۵۹
جدول ۳-۴: نتیجه الگوریتم‌های موجود روی داده‌های Iris.....	۶۰
جدول ۴-۴: مقادیر intra و inter و درصد خطای بدست آمده از الگوریتم‌ها روی داده‌های Iris.....	۶۰
جدول ۵-۴: مراکز خوشه بدست آمده با اجرای الگوریتم FAC روی مجموعه داده Wine.....	۶۲
جدول ۶-۴: نتیجه الگوریتم‌های موجود روی داده‌های Wine.....	۶۳
جدول ۷-۴: مقادیر intra و inter و درصد خطای بدست آمده از الگوریتم‌ها روی داده‌های Wine.....	۶۳
جدول ۸-۴: مراکز خوشه بدست آمده با اجرای الگوریتم FAC روی مجموعه داده Glass.....	۶۴
جدول ۹-۴: نتیجه الگوریتم‌های موجود روی داده‌های Glass.....	۶۵
جدول ۱۰-۴: مقادیر intra و inter و درصد خطای بدست آمده از الگوریتم‌ها روی داده‌های Glass.....	۶۶
جدول ۱۱-۴: مراکز خوشه بدست آمده با اجرای الگوریتم FAC روی مجموعه داده حقیقی.....	۶۷
جدول ۱۲-۴: نتیجه الگوریتم FAC روی داده‌های حقیقی.....	۶۷
جدول ۱۳-۴: مقادیر intra و inter و درصد خطای بدست آمده از الگوریتم FAC روی داده‌های حقیقی.....	۶۸

فصل اول:

مقدمہ

۱-۱ مقدمه

خوشه‌بندی^۱، دسته‌بندی بدون ناظر الگوها (مشاهدات، داده‌ها یا بردار ویژگی‌ها) درون گروه‌ها (خوشه‌ها) است [۱]. امروزه با توجه به رشد روز افزون داده‌ها، خوشه‌بندی یکی از ایده‌آل‌ترین مکانیزم‌ها برای ورود به دنیای عظیم داده‌ها است چرا که تشخیص ساختار داده‌ها را امکان پذیر می‌کند. دلیل اهمیت خوشه‌بندی این است که یک دید کلی و سریع از یک مجموعه بزرگی از اسناد را در اختیار کاربر قرار می‌دهد. می‌توان خوشه‌بندی را به این صورت تعریف کرد که خوشه‌بندی قرار دادن داده‌ها درون خوشه‌ها است بطوریکه داده‌هایی که در یک خوشه قرار دارند بیشترین شباهت را با یکدیگر و با داده‌های خوشه‌های دیگر کمترین شباهت را داشته باشند.

۲-۱ شرح مسئله

خوشه‌بندی روشی شناخته شده برای دریافت مفاهیم پنهان در متن داده‌ها است. بحث اصلی خوشه‌بندی، تقسیم داده‌ها به گروه‌هایی از عناصر است که شبیه یکدیگر باشند. هر کدام از این گروه‌ها که خوشه نامیده می‌شود، شامل عناصری است که با یکدیگر شباهت دارند و با عناصر سایر گروه‌ها متفاوتند.

خوشه‌بندی را می‌توان، یافتن ساختاری در مجموعه‌ای از داده‌ها دانست که دسته‌بندی نشده‌اند. به بیان دیگر قرار دادن داده‌ها در گروه‌هایی است که اعضای هر گروه از زاویه‌ی خاصی شبیه یکدیگرند، بطوریکه شباهت درون هر خوشه حداکثر و شباهت بین داده‌های درون خوشه‌های متفاوت حداقل باشد. برای انجام این کار روش‌های زیادی وجود دارد که می‌توان آنها را به دو دسته روش‌های افرازبندی و سلسله مراتبی تقسیم کرد. امروزه مسئله خوشه‌بندی به یک مسئله بهینه‌سازی تبدیل شده، که هدف آن پیدا کردن مراکز است که مجموع فاصله آنها با دیگر داده‌ها حداقل شود. از این‌رو می‌توان از الگوریتم‌های

¹ clustering

فراابتکاری^۱ برای حل آن استفاده کرد؛ به عبارتی هدف این الگوریتم‌ها پیدا کردن مراکز خوشه‌ها است بطوری که فاصله آنها از دیگر داده‌ها حداقل شود.

در این پایان نامه برای حل مسئله خوشه‌بندی، از الگوریتم فراابتکاری جدیدی به نام الگوریتم جنگل استفاده شده است. برای بررسی کارایی الگوریتم ارائه شده از دو سری مجموعه داده که اولی مجموعه داده‌های استاندارد و دیگری مجموعه داده‌های حقیقی استخراج شده برای عیب‌یابی سیستم‌های دوار مکانیکی استفاده شده است.

۱-۳ ساختار پایان نامه

در ادامه مبحث، در فصل دوم، مسائل مربوط به خوشه‌بندی و مفاهیم پایه‌ای لازم برای خوشه‌بندی توضیح داده شده است. در ابتدا در مورد مفهوم خوشه‌بندی و اینکه خوشه‌بندی چیست و چه کاربردهایی دارد، توضیح مختصری ارائه شده؛ سپس تفاوت آن با طبقه‌بندی ذکر و اهداف خوشه‌بندی و مراحل انجام آن شرح داده شده است. برای خوشه‌بندی روش‌های زیادی وجود دارد که می‌توان آنها را به دو دسته روش‌های افرازبندی و سلسله مراتبی تقسیم نمود؛ هر یک از این روش‌ها مزایا و معایبی دارند که در ادامه فصل دوم اشاره‌ی مختصری به آنها شده است. برای انجام خوشه‌بندی الگوریتم‌های زیادی ارائه شده است [۵] که در ادامه این فصل اشاره کوتاهی به برخی از این الگوریتم‌ها خواهیم داشت. بعد از اینکه خوشه‌بندی انجام شد معیارهایی مورد نیاز است تا بتوان خوشه‌بندی انجام شده را ارزیابی کرد، تعدادی از این معیارها در انتهای فصل دوم آورده شده است.

همانطور که ذکر شد یکی از روش‌های خوشه‌بندی، خوشه‌بندی افرازبندی است؛ روش‌هایی که به صورت افرازبندی خوشه‌بندی را انجام می‌دهند خود به دو دسته تقسیم می‌شود، روش‌های کلاسیک و روش‌های فراابتکاری. برای هر دو این روش‌ها الگوریتم‌های بسیاری وجود دارد که در فصل سوم II به مطالعه

¹ heu

و بررسی تعدادی از آنها خواهیم پرداخت.

یکی از الگوریتم‌های فرامکاشفه‌ای جدید الگوریتم جنگل می‌باشد؛ این الگوریتم نیز همانند دیگر الگوریتم‌های فرامکاشفه‌ای مبتنی بر جمعیت است. در ابتدای فصل چهارم نگاه اجمالی به الگوریتم جنگل داریم و سپس نحوه بکارگیری این الگوریتم را در مسئله خوشه‌بندی (FAC) مورد بررسی قرار دادیم. داده‌هایی که برای ارزیابی این الگوریتم استفاده شده، شامل دو سری داده است، که سری اول شامل داده‌های استاندارد و سری دوم شامل داده‌های حقیقی می‌باشد که برای عیب‌یابی سیستم‌های دوار مکانیکی تولید شده‌اند. در ادامه فصل چهارم این داده‌ها را معرفی کرده‌ایم و نتایج آنها را در جداول جداگانه آوردیم.

در نهایت در فصل پنجم یک نتیجه‌گیری کلی از روش ارائه شده خواهیم داشت و در نهایت

راه‌کارهایی را برای کارهای آتی ارائه دادیم.

فصل دوم:

بررسی معانیسم و الگوهای اولیه

۱-۲ مقدمه

امروزه با حجم عظیمی از داده‌ها روبرو هستیم که برای استفاده از آنها به ابزارهای کشف دانش نیاز داریم. داده‌کاوی به عنوان یک توانایی پیشرفته، در تحلیل داده و کشف دانش مورد استفاده قرار می‌گیرد. برای این کار روش‌های متعددی وجود دارد که هر یک از آنها برای اهداف خاصی مورد استفاده قرار می‌گیرند. یکی از مهمترین روش‌های داده‌کاوی، خوشه‌بندی است که کاربرد بسیاری در کشف دانش دارد.

۲-۲ خوشه‌بندی

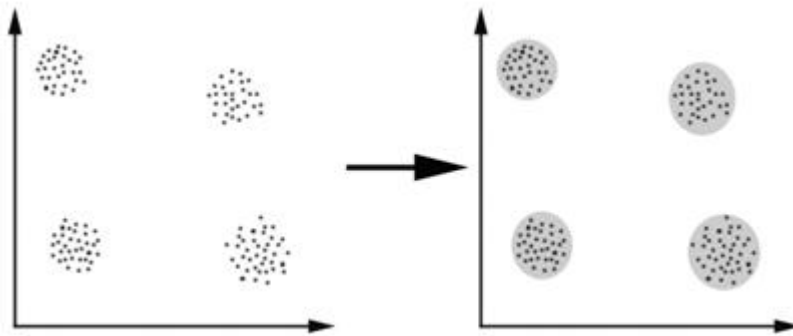
یکی از مراحل پایه‌ای درک و یادگیری داده‌ها، گروه‌بندی آن‌ها در دسته‌های ملموس است. در حالت پایه، سیستم‌های گروه‌بندی داده، اگر اطلاعاتی درباره گروه داده‌ها را در اختیار داشته باشند، با ناظر، در غیر این صورت بدون ناظر نامیده می‌شوند [۶]. خوشه‌بندی یکی از شاخه‌های یادگیری بدون ناظر می‌باشد و فرآیند خودکاری است که نمونه‌ها را به دسته‌هایی که اعضای آن مشابه یکدیگر هستند، تقسیم می‌کند؛ که به این دسته‌ها خوشه گفته می‌شود.

۱-۲-۲ تعریف خوشه‌بندی

خوشه‌بندی یک تکنیک دسته‌بندی بدون نظارت است که مجموعه‌ی داده‌ها که معمولاً بردارهایی در فضای چند بعدی هستند را براساس یک معیار شباهت یا عدم شباهت به تعداد مشخصی خوشه تقسیم می‌کند [۷]؛ بنابراین خوشه مجموعه‌ای از اشیاء است که در آن اشیاء با یکدیگر مشابه بوده و با اشیاء موجود در خوشه‌های دیگر غیر مشابه هستند. برای مشابه بودن می‌توان معیارهای مختلفی در نظر گرفت؛ مثلاً می‌توان معیار فاصله را برای خوشه‌بندی مورد استفاده قرار داد و اشیائی که به یکدیگر نزدیکتر هستند را به عنوان یک خوشه در نظر گرفت که به این نوع خوشه‌بندی، خوشه‌بندی مبتنی بر فاصله^۳ می‌گویند. به عنوان مثال در شکل ۱-۲ نمونه‌های ورودی در سمت چپ به چهار خوشه مشابه شکل سمت راست تقسیم

³Distance-based Clustering

می‌شوند. در این مثال هر یک از نمونه‌های ورودی به یکی از خوشه‌ها تعلق دارد و نمونه‌ای وجود ندارد که متعلق به بیش از یک خوشه باشد.



شکل ۱-۲: خوشه‌بندی نمونه‌های اولیه

ولی برخی شباهت‌ها را نمی‌توان به صورت عددی تخمین زد، لذا معیار شباهت دیگر، شباهت مفهومی نام دارد که به این نوع از خوشه‌بندی، خوشه‌بندی مفهومی گفته می‌شود. در خوشه‌بندی مفهومی، خوشه‌ها تنها گروهی از اشیاء با شباهت عددی نیستند، بلکه خوشه‌ها به عنوان گروهی از اشیاء که با یکدیگر یک مفهوم را نشان می‌دهند، هستند؛ در واقع تنها تعدادی خوشه تولید می‌شود که مفاهیم مرتبط را توصیف می‌کند. برای خوشه‌بندی مفهومی، ما به مجموعه‌ای از صفات برخی اشیاء (یک زبان توصیف برای مشخص کردن خوشه‌های چنین اشیایی) و یک معیار کیفیت خوشه‌بندی نیاز داریم. هدف، تقسیم‌بندی کردن اشیاء در خوشه‌ها به گونه‌ای است که معیار کیفیت ماکزیمم شود و در عین حال توصیفات عمومی از این خوشه‌ها را تعیین کند. توجه شود که در خوشه‌های مفهومی خصوصیات خوشه، با بررسی و دقت به فرآیند مشخص کردن خوشه‌ها، به وجود می‌آید. این یک تفاوت اصلی بین خوشه‌بندی مفهومی و خوشه‌بندی مبتنی بر فاصله است. در روش خوشه‌بندی مبتنی بر فاصله، خوشه‌ها مطابق با یک معیار شباهت مشخص می‌شوند. این معیار شباهت یک تابع است که فقط خصوصیات اشیاء را مقایسه می‌کند. در مقابل، در خوشه‌بندی مفهومی به شرح یا توصیف هم توجه می‌شود.

۲-۲-۲ کاربردهای خوشه‌بندی

خوشه‌بندی را می‌توان یکی از مهمترین زیر گروه‌های یادگیری بدون نظارت دانست که در موارد بسیاری کاربرد دارد [۸]. به عبارت دیگر خوشه‌بندی، نقش حیاتی در روش‌های طبقه‌بندی اطلاعات بازی می‌کند برخی از این نمونه‌ها عبارتند از:

◀ در زمینه مهندسی (مانند مهندسی برق، مهندسی کامپیوتر، یادگیری ماشین، هوش مصنوعی، تشخیص الگو و مهندسی مکانیک) کاربردهایی مانند:

- داده کاوی^۴: کشف اطلاعات و ساختار جدید از داده‌های موجود

- تشخیص گفتار^۵: در ساخت کتاب کد از بردارهای ویژگی، رد تقسیم کردن گفتار بر حسب گویندگان آن و یا فشرده سازی گفتار

- تقسیم بندی تصاویر^۶: تقسیم‌بندی تصاویر پزشکی و یا ماهواره‌ای

- وب (www): دسته‌بندی اسناد و یا دسته‌بندی مشتریان به سایت‌ها و ...

◀ علوم پزشکی (شاخه‌های ژنتیک، زیست شناسی، میکروب شناسی، فسیل شناسی، روان شناسی بالینی، آسیب شناسی) مانند:

- زیست شناسی: دسته‌بندی حیوانات و گیاهان از روی ویژگی‌های آنها

◀ علوم زمین شناسی (جغرافیا، زمین شناسی، نقشه برداری از زمین) به عنوان نمونه

- نقشه برداری شهری^۷: دسته‌بندی خانه‌ها براساس نوع و موقعیت جغرافیایی آنها

- مطالعات زلزله نگاری^۸: تشخیص مناطق حادثه خیز براساس مشاهدات قبلی

◀ علوم اجتماعی (جامعه شناسی، روان شناسی، تاریخ، آموزش و پرورش) که در این راستا به موارد

⁴ Data mining

⁵ Speech Recognition

⁶ Image Segmentation

⁷ City-Planning

⁸ Earthquake studies

زیر می‌توان اشاره کرد:

- کتابداری: دسته‌بندی کتابها

- بیمه: تشخیص افراد متقلب، تشخیص افرادی که بیمه موتور دارند و بیشترین میزان درخواست از بیمه را نیز در سال مشخصی داشته‌اند.

◀ اقتصاد (بازاریابی، تجارت) به عنوان مثال:

- در بازاریابی^۹: دسته‌بندی مشتری‌ها به دسته‌هایی بر حسب رفتارها و نیازهای آنها از طریق مجموعه زیادی از ویژگی‌ها و آخرین خریدهای آنها

۳-۲-۲ تفاوت خوشه‌بندی با طبقه‌بندی

برای تقسیم داده در گروه‌های مختلف از دو واژه طبقه‌بندی^{۱۰} و خوشه‌بندی استفاده می‌شود. در بسیاری از موارد این دو واژه به جای یکدیگر بکار برده می‌شود. اما به این نقطه باید توجه داشت که اگر چه در هر دو روش به دنبال آن هستیم که داده‌ها را براساس شباهت‌های آنها، در چندین گروه قرار دهیم، ولی طبقه‌بندی با خوشه‌بندی متفاوت است. در طبقه‌بندی نمونه‌های ورودی برچسب گذاری شده‌اند ولی در خوشه‌بندی، نمونه‌های ورودی دارای برچسب اولیه نمی‌باشند و در واقع با استفاده از روش‌های خوشه‌بندی، داده‌های مشابه، مشخص و بطور ضمنی برچسب‌گذاری می‌شوند. در واقع می‌توان قبل از عملیات طبقه‌بندی داده‌ها، یک خوشه‌بندی روی نمونه‌ها انجام داد و سپس مراکز خوشه‌های حاصل را محاسبه کرد و یک برچسب به مراکز خوشه‌ها نسبت و سپس عملیات طبقه‌بندی را برای نمونه‌های ورودی جدید انجام داد.

۴-۲-۲ هدف از خوشه‌بندی

هدف روش‌های خوشه‌بندی، گروه‌بندی مجموعه‌ای از داده‌های بدون برچسب است، به طوری که دو داده

^۹ Marketing

^{۱۰} classification

در یک خوشه تا حد امکان به هم شبیه، و در دو خوشه متفاوت، تا حد امکان از یکدیگر متمایز باشند [۹].
سوالی که مطرح می‌شود این است که چگونه یک خوشه‌بندی مفید داشته باشیم؟ این موضوع به هدف
نهایی خوشه‌بندی بستگی دارد. به عبارتی، کاربر باید ملاک را به گونه‌ای تعیین کند که بهترین نتیجه را
مطابق با نیازهایش بدست آورد.

خوشه‌بندی برای چهار هدف زیر می‌تواند مورد استفاده قرار بگیرد:

◀ یافتن ساختار زیربنایی^{۱۱} (پیدا کردن شهود راجع به داده‌ها، ایجاد فرضیه، تشخیص ناهنجاری^{۱۲} و

تعیین ویژگی‌های برجسته)

◀ گروه‌بندی طبیعی داده‌ها

◀ فشرده‌سازی به منظور سازمان دهی یا خلاصه‌سازی داده

◀ کشف گروه‌های طبیعی داده‌ها

۲-۲-۵ مراحل انجام خوشه‌بندی

یک روند کلی برای فرایند خوشه‌بندی شامل مراحل زیر است [۶]:

۱. نمایش الگو که معمولاً شامل انتخاب یا استخراج خصیصه می‌باشد.
۲. تعریف یک معیار ارزیابی شباهت با توجه به دامنه داده‌ها.
۳. فرایند خوشه‌بندی یا گروه‌بندی.
۴. خلاصه سازی داده‌ها در صورت نیاز.
۵. اعتبارسنجی سیستم.

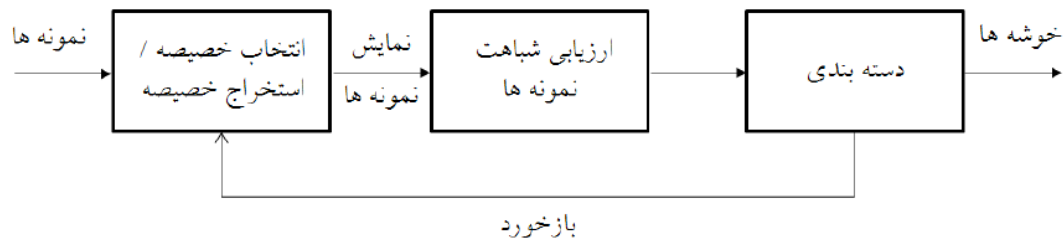
شکل ۲-۲ سه مرحله از این فرایند، که شامل بازخورد حاصل از نتایج خوشه‌بندی روی دو مرحله

¹³ underlying

¹⁴ Anomaly

ابتدایی است، را نشان می‌دهد.

نمایش الگوها،^{۱۳} به تعداد کلاس‌ها، تعداد نمونه‌های موجود و تعداد، نوع و مقیاس خصیصه‌های موجود در الگوریتم خوشه‌بندی اشاره می‌کند. برخی از این اطلاعات توسط کاربر کنترل نمی‌باشد.



شکل ۲-۲: مراحل انجام خوشه بندی [۶]

انتخاب خصیصه،^{۱۴} فرایند شناسایی یک زیر مجموعه از مؤثرترین خصیصه‌ها برای استفاده در خوشه‌بندی است، و استخراج خصیصه،^{۱۵} فرایند تغییر برخی خصیصه‌های موجود و تولید خصیصه‌های جدید می‌باشد. هر دوی این تکنیک‌ها به منظور دست یافتن به مجموعه مناسبی از خصیصه‌ها و افزایش کارایی خوشه‌بندی است.

مجاورت نمونه‌ها،^{۱۶} معمولا بوسیله یک تابع فاصله میان زوج الگوهای ورودی، اندازه‌گیری می‌شود. معیارهای گوناگونی برای اندازه‌گیری فاصله در حوزه‌های مختلف استفاده می‌شود [۱۰، ۱۱، ۱۲]. یک معیار اندازه‌گیری ساده مثل فاصله اقلیدسی معمولا برای نمایش عدم تشابه میان دو الگو بکار می‌رود [۱۳]، همچنین برای تشخیص شباهت‌های مفهومی میان الگوها می‌توان از معیارهای دیگری استفاده کرد.

مرحله گروه‌بندی به روش‌های گوناگونی انجام می‌گیرد. نتایج حاصل از خوشه‌بندی می‌تواند به شکل تقسیم‌بندی سخت (تقسیم داده‌ها درون گروه‌های مجزا) و یا تقسیم‌بندی فازی (هر نمونه با درجه‌های

¹³ Pattern Representatio

¹⁴ Feature Selection

¹⁵ Feature Extraction

¹⁶ Pattern Proximity

عضویت متفاوتی در گروه‌های مختلف قرار دارد) باشد. الگوریتم‌های خوشه‌بندی سلسله‌مراتبی از توابع شباهت برای ترکیب و تقسیم خوشه‌ها استفاده می‌کنند، تا زنجیره‌های تودرتویی از تقسیم‌بندی‌های گوناگون را ارائه دهند. الگوریتم‌های خوشه‌بندی افرازبندی تنها یک تقسیم‌بندی از نمونه‌ها را در اختیار ما می‌گذارند که در ادامه به بررسی هر یک از آنها خواهیم پرداخت.

خلاصه‌سازی داده‌ها،^{۱۷} فرایند استخراج یک تعریف مختصر و ساده از مجموعه داده‌های اصلی می‌باشد. در مبحث خوشه‌بندی، نمونه‌ای از خلاصه‌سازی داده را می‌توان یک توصیف کوتاه برای هر خوشه دانست که معمولاً آن را مدل پیش‌الگو^{۱۸} می‌نامند، برای مثال مرکز ثقل هر خوشه [۱۱] می‌تواند یک توصیف خلاصه‌ای از نمونه‌های درون آن خوشه باشد.

اعتبار سنجی، نتایج حاصل از الگوریتم‌های خوشه‌بندی را ارزیابی می‌کند. اعتبار سنجی قابل مشاهده بوده [۱۴] و برای تشخیص با معنی بودن خروجی مورد استفاده قرار می‌گیرد. یک ساختار خوشه‌بندی زمانی معتبر است که به صورت شانسی اتفاق نیفتد. زمانیکه از روش‌های آماری برای خوشه‌بندی استفاده می‌شود، اعتبارسنجی نیز با استفاده از روش‌های آماری و آزمایش فرضیات صورت می‌گیرد. دو نوع اعتبارسنجی برای مطالعه وجود دارد. اعتبارسنجی بیرونی،^{۱۹} ساختار بدست آمده برای داده‌ها را با ساختارهای قبل از خوشه‌بندی مقایسه می‌کند. اعتبارسنجی درونی،^{۲۰} سعی دارد تا مشخص کند آیا ساختار بدست آمده بطور ذاتی برای داده‌ها مناسب است یا نه.

۳-۲ روش‌های خوشه‌بندی

با توجه به کاربرد و تنوع مسایل، برای خوشه‌بندی روش‌های مختلفی وجود دارد [۵] که می‌توان آنها را به طور کلی به دو روش عمده دسته‌بندی کرد:

¹⁷ Data Abstraction

¹⁸ Prototype

¹⁹ External Examination

²⁰ Internal Examination