



مقاله پژوهشی

## تعیین واحدهای ژئومکانیکی با استفاده از تجزیه و تحلیل الگوریتم‌های بدون نظارت یادگیری ماشین مبتنی بر نگاره‌های چاه‌های نفتی

حمید قالیباف محمدآبادی<sup>۱</sup>؛ ناصر حافظی مقدس<sup>۲\*</sup>؛ غلامرضا لشکری پور<sup>۳</sup>؛ رئوف غلامی<sup>۴</sup>؛ حسین طالبی<sup>۵</sup>

۱. دانشجوی دکترای تخصصی؛ گروه زمین‌شناسی مهندسی، دانشکده علوم، دانشگاه فردوسی مشهد، مشهد، ایران

۲. استاد؛ گروه زمین‌شناسی مهندسی، دانشکده علوم، دانشگاه فردوسی مشهد، مشهد، ایران

۳. استاد؛ گروه زمین‌شناسی مهندسی، دانشکده علوم، دانشگاه فردوسی مشهد، مشهد، ایران

۴. استاد؛ دانشگاه منابع انرژی نروژ

۵. استاد؛ شرکت مناطق نفت جنوب

دریافت مقاله: ۱۴۰۰/۱۲/۲۳ پذیرش مقاله: ۱۴۰۱/۰۱/۱۹

شناسه دیجیتال (DOI): 10.22107/jpg.2022.329417.1158

واژگان کلیدی	چکیده
الگوریتم‌های یادگیری ماشین، یادگیری بدون نظارت، یادگیری نظارت‌شده، مدل $k$ - میانگین، مدل آمیخته گوسی،	در این تحقیق از روش یادگیری بدون نظارت جهت تعیین واحدهای ژئومکانیکی در یکی از چاه‌های نفتی جنوب ایران با استفاده از لاگ‌های داده‌های چاه‌نگاری شامل نگاره گاما طبیعی ( $SGR$ )، نگاره گاما اصلاح‌شده ( $CGR$ )، چگالی ( $RHOB$ )، تخلخل نوترونی ( $NPHI$ )، زمان موج برشی ( $DTSM$ ) و زمان موج طولی ( $DTCO$ ) استفاده شده است. برنامه‌نویسی مورد نیاز در محیط پایتون انجام گرفته است. در این راستا ابتدا بعد از پردازش داده‌های چاه‌نگاری از دو الگوریتم محبوب قدرتمند نظارت‌شده یادگیری ماشین ایکس جی بوست ( $XGBoost$ ) و شبکه عصبی پرسپترون چندلایه ( $Multi-Layer Perceptron Neural Network$ ) جهت بازیابی داده‌های گمشده استفاده گردید. سپس از روش‌های بدون نظارت یادگیری ماشین شامل مدل $k$ - میانگین ( $K-Means$ Clustering)، الگوریتم خوشه‌بندی سلسله مراتبی ( $HAC$ )، الگوریتم خوشه‌بندی $DBSCAN$ مبتنی بر غلظت و مدل آمیخته گوسی ( $Gaussian Mixture Modelling$ ) جهت تعیین واحدهای ژئومکانیکی مخزنی پرفشار، آهک‌های نازک‌لایه و غیر مخزنی مسئله‌دار استفاده شد. در این روش‌ها الگوریتم‌ها خود الگوهای زیرسطحی را با استفاده از داده‌ها شناسایی می‌کنند که ممکن است به راحتی در طول کاوش داده قابل مشاهده نباشند. معیار ارزیابی دقت روش دقت در شناسایی آهک‌های نازک‌لایه، سازندهای غیر مخزنی مسئله‌دار و افق‌های پرفشار سازندهای مورد مطالعه در نظر گرفته شد. نتایج مطالعات نشان داد که از بین روش‌های مورد مطالعه روش $GMM$ به جای اینکه بر اساس فاصله باشد، مبتنی بر توزیع است و از مرزهای خوشه/تصمیم بیضی استفاده می‌کند؛ بنابراین، منجر به طبقه‌بندی نرم‌تری می‌شود. علاوه بر این، به خاطر قرار دادن الگوهای احتمالاتی مختلف برای شناسایی واحدهای ژئومکانیکی، روشی بهتر جهت تعیین واحدهای مخزنی پرفشار ایلام، سروک و آهک‌های نازک‌لایه می‌باشد.

## ۱. پیشگفتار

در این مقاله، یادگیری نظارت‌شده (*Supervised Machine Learning*) رایج‌ترین و کاربردی‌ترین وظایف یادگیری ماشین است و طراحی آن بر اساس یادگیری از مدل با استفاده از داده‌های ورودی است که به خروجی دقیق نگاشت شده‌اند. در روش دیگر انجام مدل‌سازی با استفاده از یادگیری بدون نظارت (*Unsupervised Machine Learning*) است، جایی که به الگوریتم‌ها خود الگوهای زیرزمینی را با استفاده از داده‌ها شناسایی می‌کنند که ممکن است به راحتی در طول کاوش داده قابل مشاهده نباشند. یادگیری بدون نظارت یک تکنیک یادگیری ماشینی (*Machine Learning/ML*) است که نیازی به نظارت بر مدل‌ها توسط کاربران را ندارد. در یادگیری ماشینی بدون نظارت، از یک الگوریتم یادگیری برای کشف الگوهای ناشناخته در مجموعه داده‌های بدون علامت‌گذاری استفاده می‌کنیم. این برخلاف یادگیری ماشینی نظارت‌شده است که از داده‌های برچسب‌گذاری شده توسط انسان استفاده می‌کند. الگوریتم‌های یادگیری بدون نظارت از داده‌های بدون ساختار استفاده می‌کنند که بر اساس شباهت‌ها و الگوها گروه‌بندی شده‌اند. یادگیری بدون نظارت می‌تواند داده‌های پیچیده را برای ایجاد ویژگی‌های کمتر مرتبط تجزیه و تحلیل کند. سپس می‌توان با حذف این ویژگی‌ها با تأثیرات ناچیز بر بینش‌های ارزشمند، مدل را ساده کرد (Hall, 2016; Bestagini et al., 2017).

در سال‌های اخیر استفاده از لاگ‌های چاه‌نگاری جهت طبقه‌بندی رخساره‌ها با استفاده از روش‌های یادگیری ماشین تحقیقات قابل توجهی را به خود دیده است (Hall, 2016). Bestagini et al., 2017) دانشمندان بسیاری با استفاده از الگوریتم‌های محاسباتی جهت تخمین پارامترهای ناشناخته مانند رخساره‌ها ارائه داده‌اند (Ashraf et al., 2019; Bestagini et al., 2017)، گسل‌ها و شکستگی‌ها (Ashraf et al., 2019; Wu et al., 2020a; et al., 2020a)، ساخت یک مدل پتروفیزیکی برای ارزیابی مخزن ماسه‌سنگی (Ali et al., 2020). در مطالعه دیگری با استفاده از لاگ‌های تصویری (*Image log*) و اطلاعات آزمایشگاهی از هر سازند و استفاده از ترکیب الگوریتم‌های نظارت‌شده و بدون نظارت در شناسایی و پیش‌بینی رخساره‌های مختلف به کار گرفته شد (Marco et al., 2021). محمد و همکاران با استفاده از این الگوریتم‌ها لاگ‌های صوتی برشی را برای عمق‌هایی که فاقد این اطلاعات بودند پیش‌بینی نمودند (Muhammad et al., 2010). تیاگو و همکاران با بررسی روش‌های طبقه‌بندی در یادگیری ماشین بهترین روش جهت شناخت لایه‌های سنگ‌شناسی، روش جنگل تصادفی (*Random forest method*) پیشنهاد داده‌اند.

(Tiago et al., 2020) سانگ و همکاران با استفاده از روش خوشه‌بندی *K-means* جهت تعیین رخساره‌های لرزه‌ای پیشنهاد داده‌اند (Song et al., 2021). دان هام و همکاران جهت بهبود تعیین رخساره‌ها با توجه به کمبود داده‌های آموزشی جهت تعیین متغیر وابسته (رخساره) با استفاده از فرا پارامترها با استفاده از یک مدل مخلوط گوسی (*Gaussian mixture models*) مطرح کرده‌اند (Dunham et al., 2020). فنگ و همکاران با استفاده از میدان‌های تصادفی مارکوف (*MRFs*) براساس مدل مخلوط گوسی جهت سنگ‌شناسی مخازن نفتی با توجه به تغییرات افقی و عمودی مخزن ارائه داده‌اند (Fang et al., 2018). با توجه به اینکه در تعیین رخساره با استفاده از روش‌های نظارت‌شده چنانچه برچسب‌گذاری زون‌ها به صورت خیلی دقیق نباشد مدل ارائه‌شده کارایی لازم را ندارد. مارکو ایپولیتو و همکاران با ترکیب روش‌های نظارت‌شده با بدون نظارت جهت بهبود پیش‌بینی رخساره‌ها معرفی کرده‌اند (Ippolito et al., 2021). از روش‌های بدون نظارت در تهیه نقشه‌های پهنه‌بندی با استفاده از پارامترهای ژئومکانیکی جهت بیشترین محدوده تزیق سیمان برای آب‌بندی پرده آب‌بند سد سررود با استفاده از مدل سلسله مراتبی (*Analytic Hierarchy process, AHP*) به کار گرفته شد (Ghalibaf. et al., 2020).

با توجه به اینکه در اکثر مخازن نفتی شناخت دقیقی از واحدهای ژئومکانیکی مخزن از لحاظ مکانیکی و فیزیکی وجود ندارد جهت شناسایی واحدها نیاز است از هر لایه چندین آزمایش سه محوری سنگ در شرایط مخزن از لحاظ فشار و دمای بالا انجام گیرد. حتی در صورت انجام این آزمایش‌ها بازهم سنگ مورد مطالعه بکر است و فاقد قسمت‌های ازدست‌رفته و خورد شده است. علاوه بر این، انجام آزمایش به تعداد زیاد غیرممکن می‌باشد؛ بنابراین، استفاده از روش‌های نظارت‌شده اگر به طور دقیق جهت تعیین رخساره‌ها برچسب‌گذاری نشده باشد. الگوریتم به کار گرفته شده و پیشنهادی جهت پیش‌بینی زون‌های ژئومکانیکی همراه با خطا بیش از حد می‌باشد. با توجه به اینکه داده‌های چاه‌نگاری به صورت پیوسته در هر سانتی‌متر از چاه با دقت بالا پارامترهای ژئومکانیکی را اندازه‌گیری کرده است. بهتر است که از روش‌های بدون نظارت استفاده گردد تا بدون انجام آزمایش‌های پرهزینه و محدود، با استفاده از پارامترهای ژئومکانیکی به الگوریتم‌ها اجازه دهیم تا لایه‌ها و سازندهای مخزنی و غیر مخزنی را خود با توجه به شباهت‌ها از لحاظ خصوصیات ژئومکانیکی و پتروفیزیکی در کنار هم قرار دهد. بعد از این مرحله که لایه‌های مخزنی و غیر مخزنی مسئله‌دار مشخص و تأیید شد. می‌توان جهت پیش‌بینی

رخساره‌های میدان از روش‌های نظارت‌شده استفاده کرد. در این مرحله ابتدا بعد از به‌کارگیری الگوریتم موردنظر با استفاده از ماتریس درهم ریختی می‌توان پارامترهای صحت (Accuracy score)، دقت (Precision)، حساسیت یا نرخ مثبت مدل با مقدار واقعی ((Recall or Sensitivity و F1-score را به دست آورد. چنانچه نتایج این شاخص‌ها بالا باشد می‌توان بهترین الگوریتم را جهت پیش‌بینی رخساره‌ها معرفی کرد. با توجه به اینکه استفاده از هریک از روش‌های نظارت‌شده جهت ارائه الگوریتم موردنیاز، یک تحقیق جداگانه نیاز دارد، در این تحقیق ارائه نشده است.

## ۱.۲ خوشه‌بندی k - میانگین

در این مطالعه، در مرحله اول با استفاده از روش‌های یادگیری بدون نظارت ماشین مدل ژئومکانیکی (MEM) سازندهای مخزنی و غیر مخزنی مسئله‌دار (با سنگ‌شناسی پیچیده شامل انیدریت، گچ، مارن و لایه‌های نازک سنگ‌آهک) که مسئول بسیاری از مسائل پایداری چاه و زمان ازدست‌رفته غیر تولیدی (NPT) است شناسایی شدند. بعد از پردازش این روش‌ها بهترین روش جهت تعیین زون‌های غیر مخزنی مسئله‌دار (آهک‌های نازک‌لایه) و سازندهای مخزنی پرفشار ایلام و سروک که یکی از اهداف اصلی این مطالعه می‌باشد معرفی شد. شناسایی آهک‌های نازک‌لایه به این علت مهم است که در یک حوضه/مخزن بسته لایه‌های گچی در طول زمان به انیدریت تبدیل می‌شود، به علت قرارگیری این آهک‌های نازک‌لایه پر درز و شکاف در سازندهای تبخیری و ماسه‌سنگی آسماری در طی این فرایند آب آزادشده را جذب می‌کنند. در نتیجه، در زمان حفاری فشار بیش از حد انتظار باعث ریزش دیواره‌های چاه در این زون‌ها اتفاق خواهد افتاد. دانشمندان بسیاری شامل (Kolawole et al., 2016) روی کرنات‌ها، (Aslannezhad et al., 2016) با بررسی پایداری چاه بر روی ماسه‌سنگ‌ها و همچنین در سال‌های اخیر پایداری دیوار چاه بر روی شیل‌ها صورت پذیرفته است (Wang et al., 2014).

$$\arg \min \sum_{i=0}^k \sum_{X \in S_i} \|X - \mu_i\|^2 = \arg \min \sum_{i=1}^k |S_i| \text{Var} S_i \quad (1)$$

که در آن  $\mu_i$  میانگین نقاط در  $S_i$  است. این معادل است با به حداقل رساندن دو به دو مربع انحراف از نقاط در همان خوشه که مطابق با رابطه (2) زیر است.

$$\sum_{\text{Cluster } C_i} \sum_{\text{Dimension } d} \sum_{x, y \in C_i} (x_d - y_d)^2 \quad (2)$$

چون کل واریانس ثابت است، از قانون واریانس کلی می‌توان نتیجه گرفت که این معادله برابر است با بیشینه کردن مربع انحرافات بین نقاط خوشه‌های مختلف (Kriegel, 2016). رایج‌ترین الگوریتم k- میانگین یا الگوریتم لوید با استفاده از یک تکرارشونده پالایش کار می‌کند (MacKay, 2003). مراحل الگوریتم به‌صورت زیر می‌باشد (شکل ۱).

ابتدا  $k$  میانگین یعنی  $(\mu_1^{(0)}, \mu_2^{(0)}, \dots, \mu_k^{(0)})$  را که نماینده خوشه‌ها هستند، به‌صورت تصادفی مقاردهی می‌کنیم.

سپس، این دو مرحله پایین را به‌تناوب چندین بار اجرا می‌کنیم تا میانگین‌ها به یک ثابت کافی برسند و یا مجموع واریانس‌های

در این مطالعه، در مرحله اول با استفاده از روش‌های یادگیری بدون نظارت ماشین مدل ژئومکانیکی (MEM) سازندهای مخزنی و غیر مخزنی مسئله‌دار (با سنگ‌شناسی پیچیده شامل انیدریت، گچ، مارن و لایه‌های نازک سنگ‌آهک) که مسئول بسیاری از مسائل پایداری چاه و زمان ازدست‌رفته غیر تولیدی (NPT) است شناسایی شدند. بعد از پردازش این روش‌ها بهترین روش جهت تعیین زون‌های غیر مخزنی مسئله‌دار (آهک‌های نازک‌لایه) و سازندهای مخزنی پرفشار ایلام و سروک که یکی از اهداف اصلی این مطالعه می‌باشد معرفی شد. شناسایی آهک‌های نازک‌لایه به این علت مهم است که در یک حوضه/مخزن بسته لایه‌های گچی در طول زمان به انیدریت تبدیل می‌شود، به علت قرارگیری این آهک‌های نازک‌لایه پر درز و شکاف در سازندهای تبخیری و ماسه‌سنگی آسماری در طی این فرایند آب آزادشده را جذب می‌کنند. در نتیجه، در زمان حفاری فشار بیش از حد انتظار باعث ریزش دیواره‌های چاه در این زون‌ها اتفاق خواهد افتاد. دانشمندان بسیاری شامل (Kolawole et al., 2016) روی کرنات‌ها، (Aslannezhad et al., 2016) با بررسی پایداری چاه بر روی ماسه‌سنگ‌ها و همچنین در سال‌های اخیر پایداری دیوار چاه بر روی شیل‌ها صورت پذیرفته است (Wang et al., 2014).

در مطالعه حاضر جهت تعیین واحدهای ژئومکانیکی سازندهای آسماری، جهرم، پابده، گورپی، ایلام، سروک و کژدمی از چهار روش خوشه‌بندی شامل: خوشه‌بندی k - میانگین (K Means Clustering)، مدل سلسله مراتبی (HAC)، مدل ترکیبی گوسی (Gaussian Mixture Modelling) و همین‌طور استفاده از خوشه‌بندی فضایی مبتنی بر چگالی در کاربردهای دارای نویز (Density-Based Spatial Clustering of Applications with Noise /DBSCAN) استفاده شده و با مقایسه نتایج مدل بهینه معرفی شده است.

## ۲. مبانی تئوری تحقیق

خوشه‌بندی داده‌ها شکل رایج تجزیه و تحلیل داده‌های اکتشافی

(۴) بروز می‌کنیم (Kriegel, 2016)

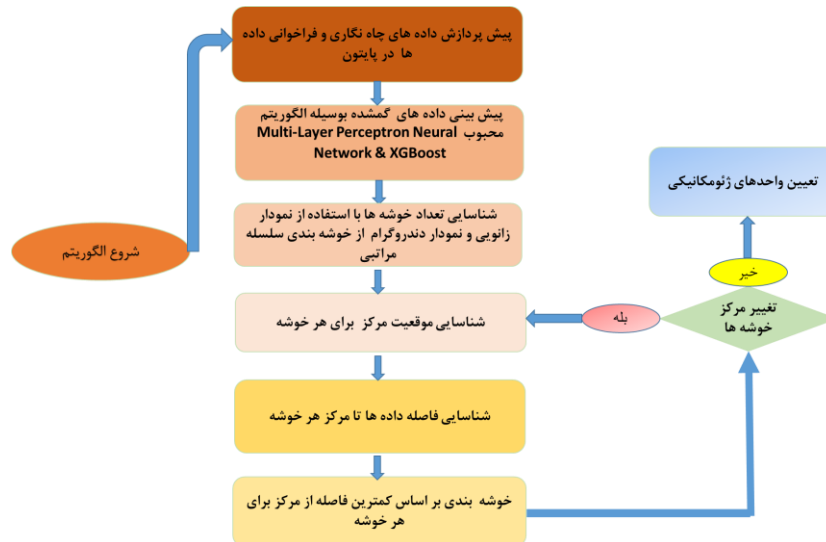
$$\mu_i^{(t+1)} = \frac{1}{|S_i^{(t)}|} \sum_{x_j \in S_i^{(t)}} x_j \quad (4)$$

در نهایت میانگین‌های مرحله آخر (در زمان  $T$ ) یعنی  $(\mu_1^{(T)}, \mu_2^{(T)}, \dots, \mu_k^{(T)})$  خوشه‌ها را نمایندگی خواهند کرد (Kriegel et al., 2016).

خوشه‌ها تغییر چندانی نکنند. از میانگین‌ها  $k$  خوشه می‌سازیم، خوشه  $i$  ام در زمان  $t$  تمام داده‌هایی هستند که از لحاظ اقلیدسی کمترین فاصله را با میانگین  $\mu_1^{(t)}$  یعنی میانگین  $i$  ام در زمان  $t$  دارند زبان ریاضی خوشه  $i$  ام در زمان  $t$  برابر رابطه (۳) خواهد بود:

$$S_i^{(t)} = \{x_p : \left\| x_p - \mu_i^{(t)} \right\|^2 \leq \left\| x_p - \mu_j^{(t)} \right\|^2 \forall j, 1 \leq j \leq k\} \quad (3)$$

حال میانگین‌ها را بر اساس این خوشه‌های جدید مطابق فرمول



شکل ۱. مراحل الگوریتم  $k$  - میانگین جهت تعیین واحدهای ژئومکانیکی

ترکیب شده و خوشه‌های سطح بالاتر را ایجاد می‌کنند. این شیوه خوشه‌بندی سلسله مراتبی به روش تجمیعی (Agglomerative) معروف است. در الگوریتم HAC پیچیدگی زمانی برابر با  $O(n^3)$  و فضای مورد نیاز حافظه نیز برابر با  $O(n^2)$  است. بنابراین با افزایش حجم داده‌ها، سرعت و فضای حافظه برای اجرای عملیات خوشه‌بندی به شدت افزایش می‌یابد.

برای اجرای محاسبات مربوط به این روش خوشه‌بندی، به دو معیار فاصله (شبهت) احتیاج داریم. ۱- میزان فاصله بین زوج مشاهدات. ۲- میزان فاصله بین خوشه‌ها. در حالت اول توابع فاصله برای داده‌های کمی یا کیفی قابل استفاده هستند؛ بنابراین اگر داده‌ها کمی باشند برای مثال می‌توان از فاصله اقلیدسی یا فاصله منهن استفاده کرد. همچنین برای داده‌های کیفی نیز می‌توان از میزان انطباق ساده (Simple Matching) و یا فاصله همینگ (Hamming Distance) برای داده‌های کمک گرفت. معمولاً قبل از شروع مراحل خوشه‌بندی سلسله مراتبی تجمیعی، از یک ماتریس فاصله (Distance Matrix) یا ماتریس شبهت (Similarity Matrix) برای تسریع در محاسبات کمک گرفته می‌شود. این ماتریس نشان می‌دهد که فاصله بین هر زوج از مشاهدات چقدر است. البته نوع تابعی که

به‌طور کلی خوشه‌بندی کی- میانگین یک روش خوشه‌بندی سخت است که در آن یک نقطه داده یا متعلق به یک خوشه است یا نیست. همچنین با اعمال یک دایره (ابر کره در مجموعه داده‌های چندبعدی) روی داده‌ها، خوشه‌بندی را انجام می‌دهد. از نظر محاسباتی بسیار کارآمد است و می‌توان به راحتی پیاده‌سازی کرد.

## ۲.۲ خوشه‌بندی سلسله مراتبی (HAC)

یکی از روش‌های یادگیری ماشین (Machine Learning) که به آموزش بدون نظارت (Unsupervised Learning) شهرت دارد، تحلیل خوشه‌بندی (Clustering Analysis) است. در این روش، برعکس خوشه بندی  $k$  - میانگین، هر مشاهده ممکن است در بیش از یک خوشه قرار گیرد زیرا بر اساس سطوح مختلف فاصله، خوشه‌ها تشکیل می‌شود. بنابراین هر خوشه ممکن است زیرمجموعه خوشه دیگر در سطحی از فاصله قرار گیرد. خوشه‌بندی سلسله مراتبی با روش تجمیعی (Hierarchical Agglomerative Clustering, or HAC) اگر دیدگاه به این نمودار از پایین به بالا باشد (Bottom-Up)، برحسب ارتفاع نمودار (Height) در سطح پایینی، خوشه‌ها زیرمجموعه خوشه‌های سطح بالاتر هستند در نتیجه به نظر می‌رسد که خوشه‌های زیرین با یکدیگر

توزیع است.

$$p(x) = \sum_{k=1}^K \pi_k N(x|\mu_k, \Sigma_k), \sum_{k=1}^K \pi_k = 1 \quad (5)$$

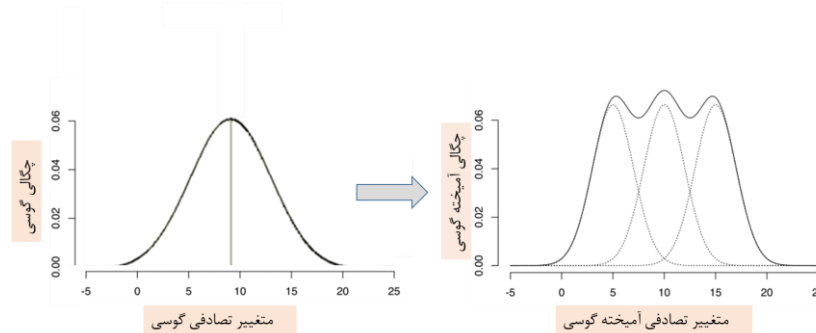
فرض کنید تعداد خوشه‌ها برابر با  $k$  باشد. در اینجا پارامترهای گاوسی (نرمال) نیز با  $\theta$  نشان داده می‌شوند. متغیر پنهان که به صورت یک بردار، برچسب هر یک از مشاهدات را نشان می‌دهد با  $Z$  مشخص شده است. در این صورت مدل آمیخته گاوسی یک‌بعدی برای این خوشه‌های به صورت رابطه ۶ نوشته می‌شوند.

$$p(X|\theta) = \sum_Z p(X, Z|\theta) \quad (6)$$

در نتیجه هدف از اجرای خوشه‌بندی بر مبنای مدل، بیشینه‌سازی این تابع درست‌نمایی است. مقدار  $Z$  به دو صورت دو دویی یا باینری است به این معنی که اگر مشاهده  $m$  به خوشه  $k$  ام تعلق داشته باشد مقدار این متغیر برابر با ۱ و در غیر این صورت برابر با صفر خواهد بود. البته از آنجایی که بیشینه سازی لگاریتم آن مشابه است و لگاریتم گیری، فرم تابع احتمال آمیخته گاوسی (نرمال) را ساده‌تر می‌کند از رابطه ۷ برای اجرای الگوریتم خوشه‌بندی بر مبنای مدل استفاده خواهیم کرد.

$$Ln[p(X|\theta)] = Ln\left\{\sum_Z p(X, Z|\theta)\right\} \quad (7)$$

در نمودار شکل ۲، یک نمونه از توزیع آمیخته گاوسی (ترکیب سه متغیر تصادفی نرمال) نشان داده می‌شود.



شکل ۲. نمایش گرافیکی مدل گاوسی و مدل‌های مخلوط گاوسی (Li et al., 2011)

متوسط گیری، هدف برآورد توزیع متغیر پنهان به شرط وزن  $(\pi)$  میانگین‌ها ( $Means$ ) و ماتریس کوواریانس  $(\Sigma)$  توزیع آمیخته نرمال است. بردار پارامترها را در اینجا با نماد  $\theta$  نشان می‌دهیم. برای این کار ابتدا یک حدس اولیه برای این پارامترها زده می‌شود. سپس گام‌ها به ترتیب برداشته می‌شوند. حدس‌های اولیه را می‌توان به وسیله الگوریتم  $k$ -means به دست آورد. به این معنی که برای خوشه‌های حاصل از الگوریتم  $k$ -means، میانگین، ماتریس کوواریانس و وزن‌ها محاسبه می‌شود. منظور از وزن، درصدی (احتمال) از داده‌ها

باید به وسیله آن فاصله اندازه‌گیری شود، در مقدارهای موجود در این ماتریس تأثیرگذار است.

در خوشه‌بندی سلسله مراتبی تجمیعی یا  $HAC$ ، با توجه به مقدارهای این ماتریس، مشاهدات یا خوشه‌هایی که دارای کمترین فاصله (بیشترین شباهت) هستند با هم ادغام می‌شوند و خوشه جدیدی می‌سازند. در مرحله بعد باز هم فاصله بین مشاهدات و یا خوشه‌های جدید، توسط ماتریس فاصله که به روزرسانی شده، محاسبه و کار ادغام ادامه پیدا می‌کند تا تنها یک خوشه باقی بماند (Hastie et al., 2009).

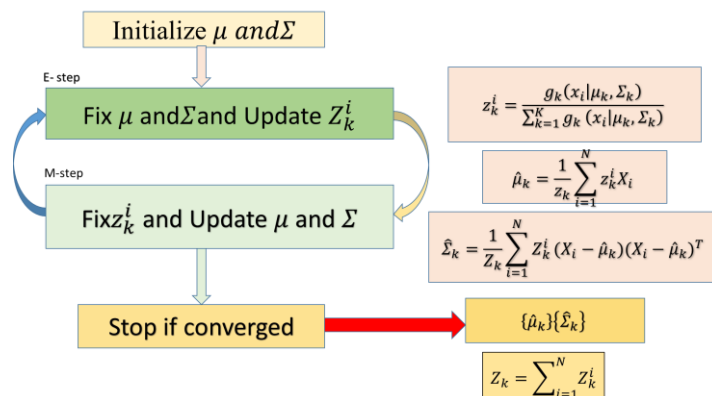
### ۳.۲ مدل‌های مخلوط گاوسی (GMM)

در این تکنیک خوشه‌بندی پیشرفته از ترکیبی از توزیع‌های گاوسی برای مدل‌سازی یک مجموعه داده استفاده می‌شود. این مدل‌های مخلوط احتمالی هستند و برای تولید نمونه‌های داده از مدل‌های خوشه‌بندی  $GMM$  استفاده می‌شود. در این مدل‌ها، هر نقطه داده عضوی از همه خوشه‌های مجموعه داده است، اما با درجات مختلف عضویت. احتمال عضویت در یک خوشه خاص بین ۰ و ۱ است. در مدل‌های مخلوط گاوسی، اطلاعات کلیدی شامل مراکز گاوسی پنهان و کوواریانس داده‌ها است. این باعث می‌شود که شبیه به  $K$ -means خوشه‌بندی شود. یک مدل آمیخته گاوسی در رابطه ۵ نشان داده می‌شود. در اینجا منظور از  $\pi_k$  وزن توزیع  $k$  ام یا فراوانی مشاهدات از آن

وجود پارامترهای زیاد در این تابع (مانند میانگین، واریانس و کوواریانس بین خوشه‌ها) امکان بیشینه‌سازی آن را با روش‌های تحلیلی نمی‌دهد؛ بنابراین در چنین مواقعی استفاده از تکنیک‌های عددی مانند الگوریتم  $EM$  کارساز است. در این تحقیق به منظور برآورد پارامترهای مدل آمیخته گاوسی از الگوریتم  $EM$  استفاده شد. این الگوریتم دارای دو بخش یا گام است. گام متوسط‌گیری ( $Expectation$ ) و گام بیشینه‌سازی ( $Maximization$ ) که به طور خلاصه مراحل این الگوریتم در نمودار شکل ۳ نشان داده می‌شود. در گام اول ( $E$ -step) یا گام

نکند)، ادامه خواهد داشت. به این ترتیب الگوریتم *EM* علاوه بر برآورد پارامترهای توزیع آمیخته گاوسی، برچسبها یا مقدار پنهان را نیز مشخص می کند.

است که در یک خوشه قرار دارند. در گام دوم (*M-step*) با استفاده از متغیرهای پنهان، سعی می شود تابع درست نمایی نسبت به پارامترهای  $\theta$  بیشینه شود. مراحل گام *E* و گام *M* تا زمانی که الگوریتم همگرا شود (مقدار تابع درست نمایی تغییر



شکل ۳. مراحل الگوریتم EM برای مدل آمیخته گاوسی (GMM)

(*quality*): نمای چندبعدی، دقت (درست یا غلط، دقیق یا نه)، کامل بودن (ثبت نشده، غیرقابل دسترس)، سازگاری (*Consistency*) (برخی اصلاح شده اما برخی نه، نامشخص)، به موقع بودن (به روز رسانی به موقع)، باورپذیری (داده ها چقدر قابل اعتماد هستند؟)، تفسیرپذیری (به چه راحتی می توان داده ها را درک کرد).

- جایگزینی داده های گمشده: مقادیر از دست رفته با استفاده از دو الگوریتم محبوب قدرتمند نظارت شده یادگیری ماشین ایکس جی بوست (*XGBoost*) و شبکه عصبی پرسپترون چندلایه *Multi-Layer Perceptron Neural Network* تعیین شد. این الگوریتم یک روش عمیق یادگیری ماشین نظارت شده است. استفاده از این روش برای پیش بینی طیف وسیعی از علوم رو به گسترش می باشد. این روش زمانی که شناخت خوبی از پیش بینی مدل وجود دارد استفاده می گردد. به طور مثال در مطالعات ژئومکانیک ما به مدل تغییرات سنگ شناسی در عمق های مختلف را می دهیم تا بهترین الگوریتم را آنالیز و پردازش کند. معمولاً جهت پر کردن داده های گمشده از روش های میانگین گیری یا استفاده از داده هایی بالا و پایین داده گمشده (بالا و پایین واحد ژئومکانیکی) استفاده می گردد. ارائه رابطه خطی نیز به دلیل محاسبه مقدار میانگین متغیر وابسته همیشه همراه با خطا بوده و مدل مناسبی ارائه نمی دهد. اما الگوریتم های نظارت شده

به طور کلی، روش *GMM* اجازه می دهد تا نقاط، داده را خوشه بندی کنند با این تفاوت که واریانس داده ها را محاسبه می کند، منجر به طبقه بندی نرم تری می شود و به جای اینکه بر اساس فاصله باشد، مبتنی بر توزیع است. همچنین، نقطه داده ای که طبقه بندی می شود، احتمال بودن در دو گروه خوشه بندی را دارد. در حالی که اگر خوشه های داده دایره ای باشند، خوشه بندی *K-Means* عالی عمل می کند، اما در موقعیت های پتروفیزیکی و زمین شناسی داده ها به ندرت الگوهای دایره ای خوبی را تشکیل می دهند. بنابراین مدل سازی *GMM* از مرزهای خوشه/تصمیم بیضی شکل استفاده می کند و بنابراین انعطاف پذیرتر است (*Jinghua et al., 2021*). برای این منظور از داده های چاه نگاری یکی از مخازن نفتی واقع در جنوب ایران جهت تعیین زون های ژئومکانیکی استفاده شده است.

### ۳. پیش پردازش داده ها

در مطالعه حاضر از داده های چاه نگاری یکی از میدان نفتی جنوب ایران شامل نگاره گاما طبیعی (*SGR*)، نگاره گاما اصلاح شده (*CGR*)، چگالی (*RHOB*)، تخلخل نوترونی (*NPHI*)، زمان موج برشی (*DTSM*) و زمان موج طولی (*DTCO*) که مستقیماً در تعیین واحدهای ژئومکانیکی تأثیر دارند استفاده شده است. ابتدا فایل داده های رقومی (*Las*) این داده ها در پایتون فراخوانی شد و یک دیتا فریم (*Data Frame*) ساخته شد بعد از این مرحله تمام مراحل پیش پردازش (*Data Preprocessing*) جهت ساخت مدل شامل موارد زیر به درستی صورت پذیرفت:

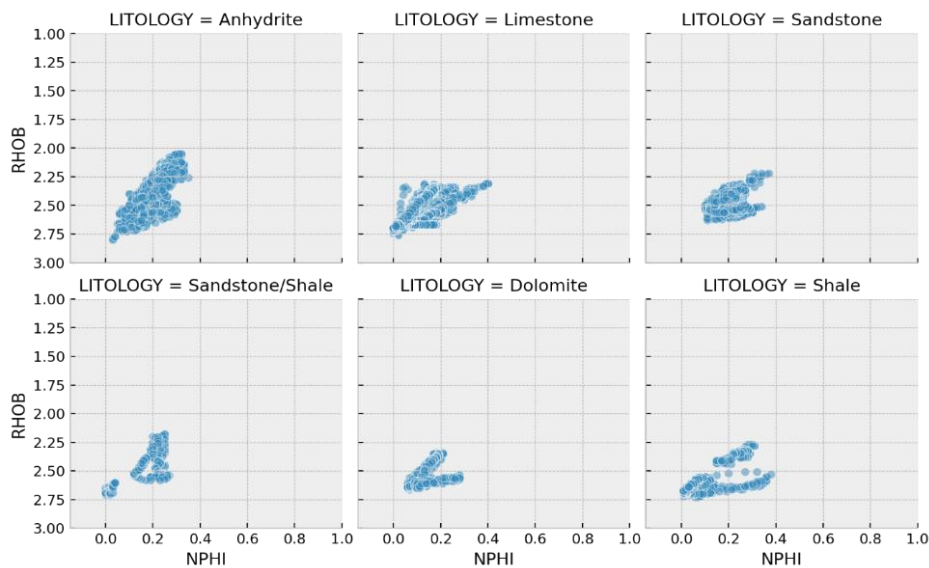
- اندازه گیری کیفیت داده ها *Measures for data*

بخشی از کاهش داده‌ها اما با اهمیت ویژه، به‌ویژه برای داده‌های عددی

بعد از این مرحله با توجه به اینکه در مطالعات زمین‌شناسی شناخت کاملی از نوع سازندهای مخازن وجود دارد، در این مدل‌سازی از یک ستون رخساره‌ای تحت عنوان لیتولوژی جهت شناسایی زون‌های ژئومکانیکی حاصل از مدل استفاده گردیده است. برای انجام این کار در محیط پایتون از روش *Seaborn-FacetGrid* برای رسم نمودارهای پراکنده چگالی-نوترون (*scatterplots*) برای هر سنگ‌شناسی که در ستون رخساره‌ای مشخص شده استفاده گردید. *FacetGrid* برای ایجاد یک ساختار زیرسطحی برای نمودار استفاده می‌شود. در این مثال، هنگامی که 6 ستون وجود دارد، داده‌ها تقسیم‌بندی می‌شوند. سپس می‌توانیم یک کراس پلات نوترونی چگالی را روی آن *FacetGrid* ترسیم کنیم. در این تحقیق با توجه به توالی رسوبی و چین‌شناسی، به ترتیب از بالا به پایین بودن سازندهای گچساران، آسماری، جهرم، پابده، گوری، ایلام، سروک و کژدمی از پنج رخساره قابل پیش‌بینی شامل انیدریت، کلسیت، دولومیت، ماسه‌سنگ، شیل، ماسه‌سنگ/ شیل استفاده شده است. (شکل ۴).

داده‌های گمشده را با در نظر گرفتن تمام داده‌های ورودی به مدل در نظر می‌گیرد. به‌طور مثال در عمقی از لاگ که ما داده تخلخل را نداریم این داده گمشده با بررسی تمام متغیرها از جمله سرعت موج طولی، موج برشی، چگالی و غیره پر می‌شود. در نتیجه این داده پر شده بهترین حالت ممکن برای نبود داده گمشده ارائه می‌دهد. قابل ذکر است داده‌های گمشده با  $Accuracy=99\%$  به‌دست‌آمده است.

- پاک‌سازی داده‌ها: در این بخش داده‌های نویز دار صاف شد ((*smooth noisy data*)، موارد پرت شناسایی یا حذف شد و ناسازگاری‌ها و خطاها را برطرف شد.
- یکپارچه‌سازی داده‌ها ((*Data integration*): ادغام چندین پایگاه داده، یکپارچگی داده یا فایل، تبدیل داده‌ها، عادی‌سازی و تجمیع
- کاهش داده‌ها: نمایش کاهش‌یافته در حجم را به دست می‌آورد اما نتایج تحلیلی یکسان یا مشابهی را تولید می‌کند.
- گسسته‌سازی داده‌ها ((*Data discretization*):



شکل ۴. پراکندگی داده‌های نوترون-چگالی به تفکیک خصوصیات سنگ‌شناسی در ستون چاه مورد مطالعه

خوشه‌ها اشتباه انتخاب شده باشد، الگوریتم‌ها ممکن است عملکرد خوبی نداشته باشند یا رفع آن‌ها بیشتر طول بکشد. از آنجایی که برخلاف عملیات تحلیل رده‌بندی (*Classification*)، خوشه‌بندی (*Clustering*) یک فرآیند

#### ۴. خوشه‌بندی - بدون نظارت

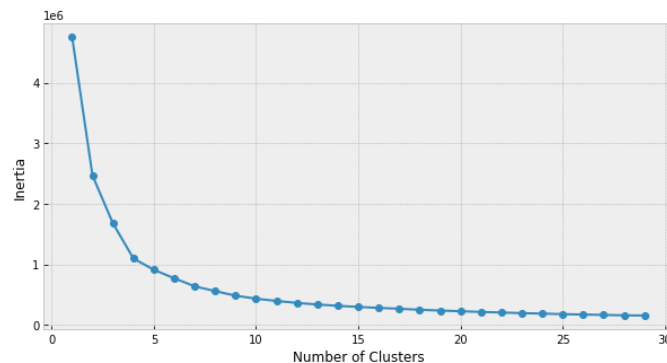
##### ۱,۴ یافتن تعداد بهینه خوشه‌ها

برای اطمینان از اینکه مدل‌های *HAC*، *K-Means* و *Gaussian Mixture Modeling* به‌طور مؤثر کار می‌کنند، باید تعداد اولیه خوشه‌ها را برای آن‌ها فراهم کنیم. اگر تعداد

با روش زانویی استفاده گردیده است. انتخاب تعداد خوشه‌ها در این روش اساساً با مدل  $k$ -means متفاوت است. در این خوشه‌بندی بجای انتخاب تعداد خوشه‌ها و آغاز با مرکز وارهای تصادفی، کار بدین صورت است که هر نقطه در مجموع داده یک خوشه است. سپس، دو نقطه نزدیک‌تر به هم پیدا و با یکدیگر در یک خوشه ادغام می‌شوند. پس از آن، دو نقطه نزدیک به هم بعدی یافت و با هم در یک خوشه قرار می‌گیرند. این فرآیند تا هنگامی انجام می‌شود که تنها یک خوشه بسیار بزرگ باقی بماند. در طول فرآیند بیان‌شده، چیزی ساخته می‌شود که به آن دندروگرام ( $Dendrogram$ ) گفته می‌شود. در نمودار شکل ۶ زیر تعداد خوشه‌ها با تقاطع خط افقی با محور عمودی در هر جای نمودار بسته به نظر متخصص می‌توان انتخاب کرد. در اینجا عدد ۵ در نظر گرفته شده است. می‌توان مشابه روش زانویی این عدد را ۵ تا ۱۰ انتخاب کرد.

بدون نظارت ( $Unsupervised$ ) است، بررسی صحت نتیجه خوشه‌بندی به راحتی امکان‌پذیر نیست. بنابراین احتیاج به معیارهای مناسب، هم برای بررسی کارایی یک روش خوشه‌بندی در بازیابی خوشه‌ها و هم برای مقایسه عملکرد روش‌های مختلف خوشه‌بندی، ضروری به نظر می‌رسد. در این میان دو گونه معیار ارزیابی نتایج خوشه‌بندی وجود دارد. معیارهای ارزیابی درونی ( $Internal Criteria Index$ ) و معیارهای ارزیابی بیرونی ( $External Criteria index$ ). از روش زانویی ( $elbow$ ) برای تعیین تعداد صحیح خوشه‌ها در این دیتاست استفاده شده است. در این روش مقادیر افزایشی  $k$  بر روی محور افقی و مجموع خط‌هایی که در هنگام استفاده از  $k$  میانگین رخ داده بر روی محور عمودی ترسیم می‌شود (شکل ۵).

در این تحقیق از روش خوشه‌بندی سلسله مراتبی ( $Agglomerative hierarchical Clustering$ ) جهت مقایسه

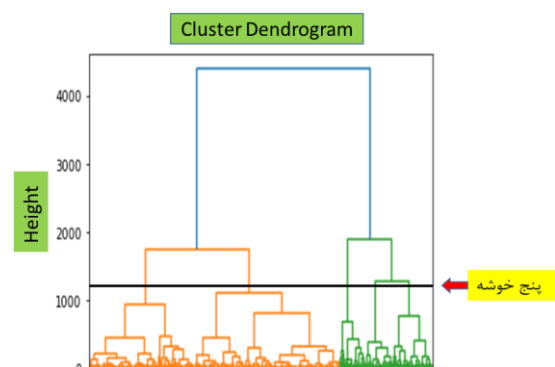


شکل ۵. نمودار زانوایی جهت تعیین تعداد بهینه خوشه‌ها، در این روش مقادیر افزایشی  $k$  بر روی محور افقی و مجموع خط‌هایی که در هنگام استفاده از  $k$  میانگین رخ داده بر روی محور عمودی ترسیم می‌شود در این تحقیق تعداد خوشه‌ها ۵ انتخاب شده است.

می‌یابد. هیچ شکاف مشخصی در این مجموعه داده وجود ندارد، با این حال، می‌توانیم ببینیم که شیب از حدود ۵ خوشه به بعد کم می‌شود. در نتیجه، مقدار تابع هزینه ثابت می‌شود. انتخاب این مقدار به مفسر بستگی دارد و می‌تواند از ۵ تا ۱۰ متغیر باشد. بنابراین برای این مثال ما ۵ را به عنوان تعداد بهینه خوشه‌ها در نظر گرفته شده است.

### ۵. ارزیابی نتایج

اکنون که خوشه‌ها با استفاده از روش‌های  $HAC$ ,  $KMeans$ ,  $DBSCAN$  و  $GMM$  محاسبه شده‌اند، می‌توانیم نتایج را رسم و تحلیل کنیم. لازم به ذکر است که این روش‌ها بدون نظارت هستند و از داده‌های برجسب‌گذاری شده برای آموزش استفاده نمی‌کنند لذا جهت مقایسه نتایج با نمودارهای لیتولوژی چاه تطبیق داده شده است. در نمودار (شکل ۷) رخساره‌های سنگی ( $Lithology$ ) و نتایج  $HAC$ ,  $DBSCAN$ ,  $KMeans$  و  $GMM$  جهت مقایسه نشان داده می‌شود. همچنین در شکل ۸

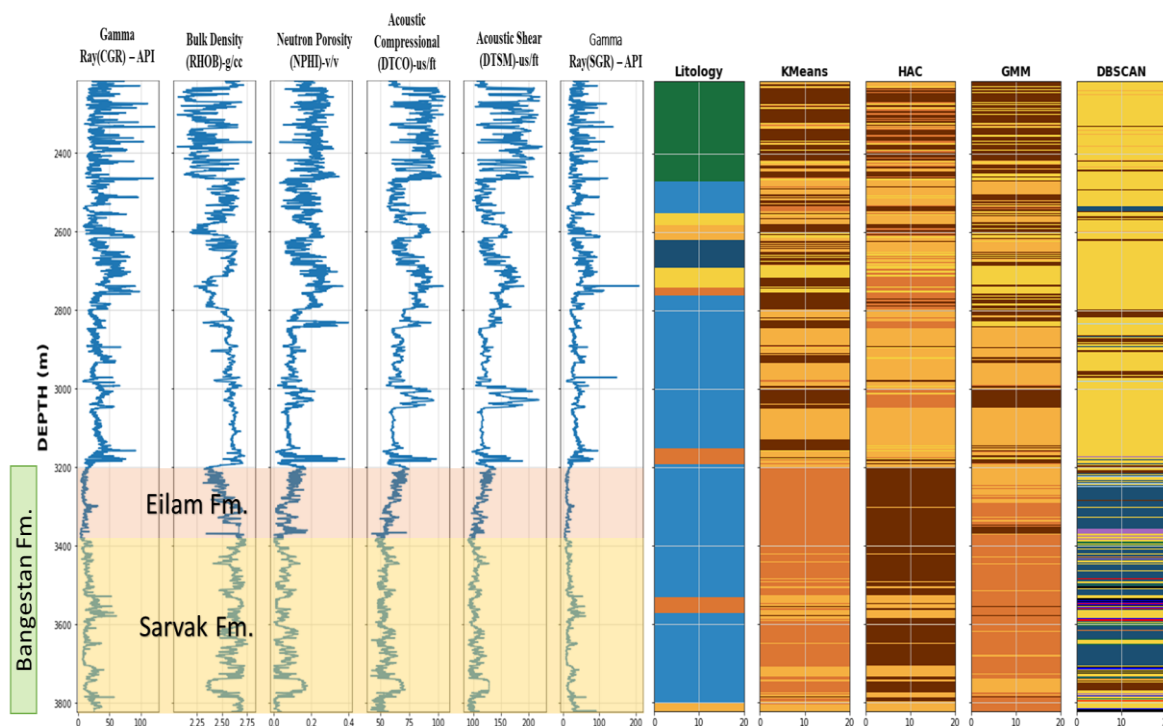


شکل ۶. نمودار دندروگرام جهت تعیین تعداد بهینه خوشه‌ها که در اینجا یک خط افقی از هر نقطه نمودار عبور داده می‌شود. تعداد بهینه خوشه‌ها از ۵ تا ۱۰ عدد می‌توان انتخاب کرد. در اینجا عدد پنج انتخاب شده است.

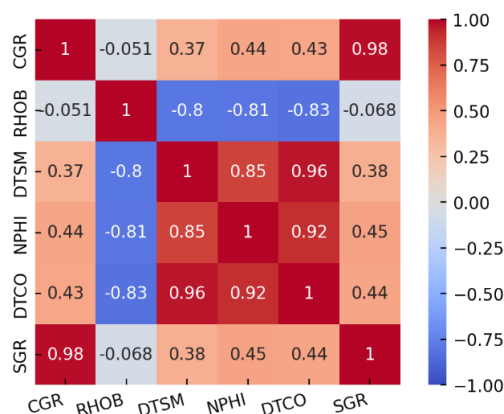
در نمودار بالا می‌بینیم که تابع هزینه (مجموع مجذور فواصل تا نزدیک‌ترین مرکز خوشه) با افزایش تعداد خوشه‌ها کاهش



با ماتریس در هم‌ریختی (*Confusion matrix*) ارتباط مثبت و منفی داده‌های چاه‌نگاری نشان داده می‌شود که در این شکل زمان عبوری موج فشاری و برشی با تخلخل ارتباط مستقیم دارد و با چگالی که قابل انتظار است ارتباط منفی را نشان می‌دهد.



شکل ۷. اولین چیزی که باید به آن توجه داشت این است که رنگ‌ها لزوماً به این معنی نیستند که داده‌ها از یک گروه در هر روش هستند. در این شکل مقایسه روش‌های ماشین لرنینگ شامل *k-means*، *HAC*، *GMM*، *DBSCAN* و یک ستون تحت عنوان لیتولوژی جهت شناسایی سازندهای مدل که از عمق ۳۲۰۰ تا ۳۸۰۰ سازند بنگستان می‌باشد

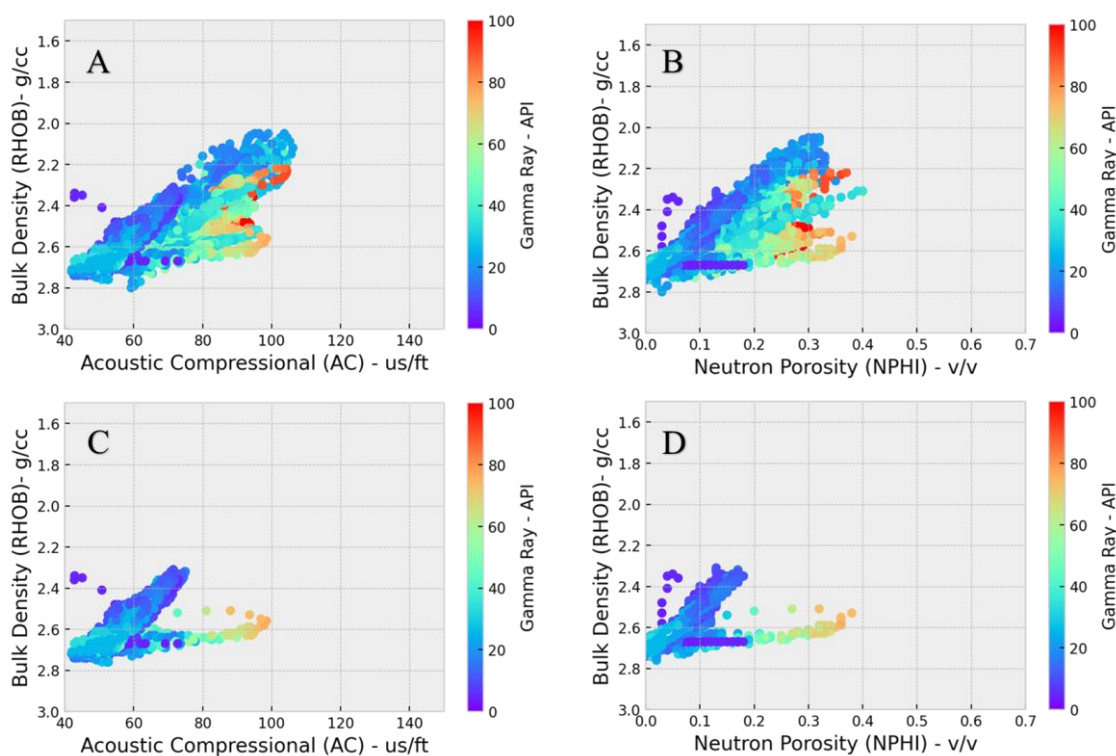


شکل ۸. در این ماتریس در هم‌ریختی (*Confusion matrix*) ارتباط مثبت و منفی داده‌های چاه‌نگاری نشان داده می‌شود

در نمودار شکل ۷، اولین چیزی که باید به آن توجه داشت این است که رنگ‌ها لزوماً به این معنی نیستند که داده‌ها از یک گروه در هر روش هستند. در این شکل مقایسه روش‌های ماشین لرنینگ شامل *k-means*، *HAC*، *GMM*، *DBSCAN* و یک ستون تحت عنوان لیتولوژی جهت شناسایی سازندهای مدل که از عمق ۳۲۰۰ تا ۳۸۰۰ سازند بنگستان می‌باشد. این یک سازند آهکی - دولومیتی با رخساره‌های پکستون و گریستون دارای تخلخل حفره‌ای و شکافی تعیین شده است که جهت مقایسه با مدل با رنگ آبی کدگذاری شده است. همان‌طور که مشخص است در روش خوشه‌بندی *k* - میانگین

دومی حداقل نقاط موجود در یک خوشه است که به آن *MinPoints* می‌گویند که برای نتایج بهتر باید این اعداد با تکرارهای مختلف عوض گردد تا بهترین نتیجه را به کاربر بدهد. مشاهده می‌گردد که زمان موج  $p$  با تخلخل رابطه مستقیم و با چگالی رابطه معکوس دارد که تناسب خوبی را نشان می‌دهد. اما این پارامترها با لاگ گاما رابطه معناداری را نشان نمی‌دهد. به‌طور مثال در شیل‌ها تخلخل بالا/نفوذپذیری پایین است. در نتیجه لاگ گاما بالا ولی در آهک‌های سازند بنگستان مقدار گاما پایین و تخلخل بالا است. در نتیجه، سنگ‌شناسی بیشتر از تخلخل و چگالی تبعیت می‌کند. جهت شناخت بیشتر مخزن بنگستان با استفاده از نمودار نوترون-چگالی و کندی موج فشاری-چگالی این قسمت از مخزن را از بقیه مخازن جدا شده است. همان‌طور که مشخص است در این نمودار مقدار گاما کم می‌باشد که نشان‌دهنده مخزن آهکی - دولومیتی می‌باشد (شکل ۹).

و خوشه‌بندی سلسله مراتبی پیش‌بینی سازند بنگستان به‌صورت یک سازند تقریباً مجزا به‌درستی نشان داده می‌شود با این وجود تفکیک بهتری را نشان نمی‌دهد. اما در روش *GMM* رخساره بنگستان را به‌صورت مشخص‌تر که در قسمت بالا و پایین سازند بنگستان با رنگ زرد متمایل به قهوه‌ای روشن شامل لایه‌های آهکی - دولومیتی دارای تخلخل بالا که سازندهای مخزنی ایلام و سروک را مشخص کرده است و رنگ قهوه‌ای روشن لایه‌های آهکی - دولومیتی با تخلخل کمتر (وکستون) و یک رخساره با رنگ قهوه‌ای تیره تفکیک کرده است. همچنین شناسایی لایه‌های ماسه‌سنگی سازند آسماری و میان لایه‌های آهک‌های نازک‌لایه در این سازند که مسئول ناپایداری چاه است به‌درستی شناسایی شد. روش *DBSCAN* جهت مقایسه با این روش‌ها انتخاب گردیده است. اما نتایج آن دل‌چسب نیست علت این است این روش برای داده‌های پراکنده دارای نویز بسیار خوب عمل می‌کند در الگوریتم *DBSCAN* دو پارامتر وجود دارد. یکی از آن‌ها شعاع است که به آن *Epsilon* نیز می‌گویند و

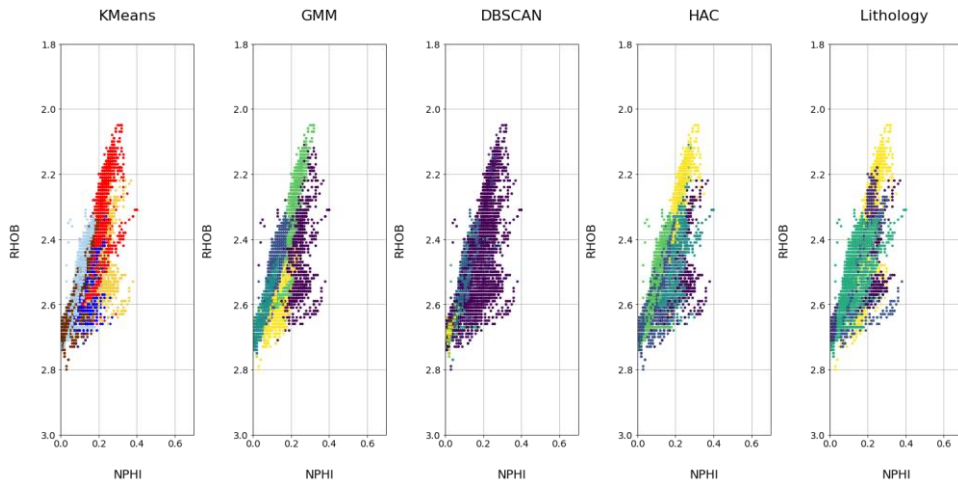


شکل ۹. نمودار A و B مربوط به کل ستون چاه می‌باشد. اما در نمودار C و D جهت شناخت بیشتر مخزن بنگستان از عمق ۳۲۰۰ تا ۳۸۰۰ جدا شده است

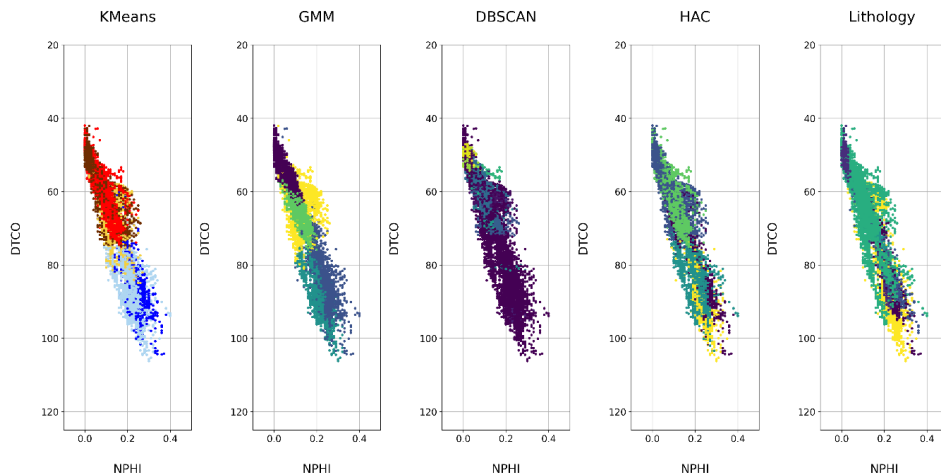
(شکل ۱۰)، ارتباط تخلخل با زمان موج  $P$  (شکل ۱۱)، ارتباط گاما با چگالی (شکل ۱۲) و ارتباط تخلخل-نوترونی با گاما نشان داده می‌شود (شکل ۱۳).

### ۵.۱ نمودارهای پراکنده (Scatter plots / Cross plots)

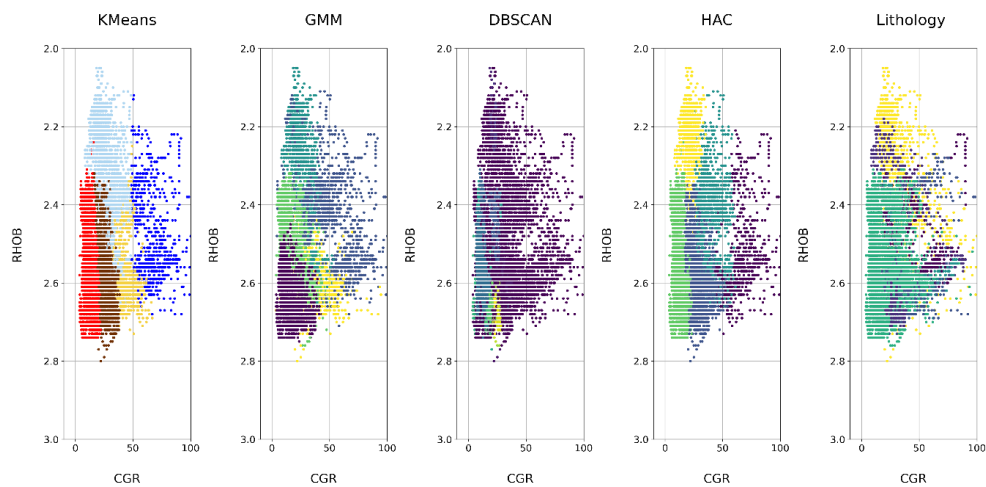
راه دیگر برای مشاهده عملکرد خوشه‌بندی از طریق نمودارهای پراکنده (*Scatter plots / Cross plots*) است. این کار را با استفاده از نمودارهای پراکنده متقاطع چگالی - تخلخل نوترونی



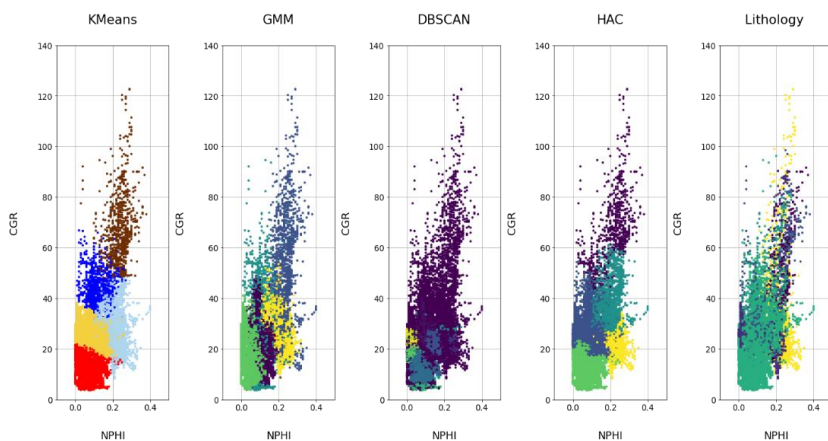
شکل ۱۰. مقایسه روش‌های ماشین لرنینگ به‌طور مثال در روش KMeans و HAC در گوشه پایین سمت چپ جایی که چگالی بالا و نوترون پایین است اختلاط دیده می‌شود ولی در روش GMM با تفکیک بهتر و مطابقت بیشتری نشان داده می‌شود.



شکل ۱۱. از بین روش‌های فوق روش خوشه‌بندی سلسله مراتبی،  $k$ - میانگین و مدل آمیخته گوسی بهتر از بقیه عمل کرده است



شکل ۱۲. از بین روش‌های فوق روش خوشه‌بندی سلسله مراتبی،  $k$ - میانگین و مدل آمیخته گوسی بهتر از بقیه عمل کرده است. با این وجود در روش GMM تفکیک واحدهای ژئومکانیکی و شناسایی مبتنی بر وارپانس می‌باشد.



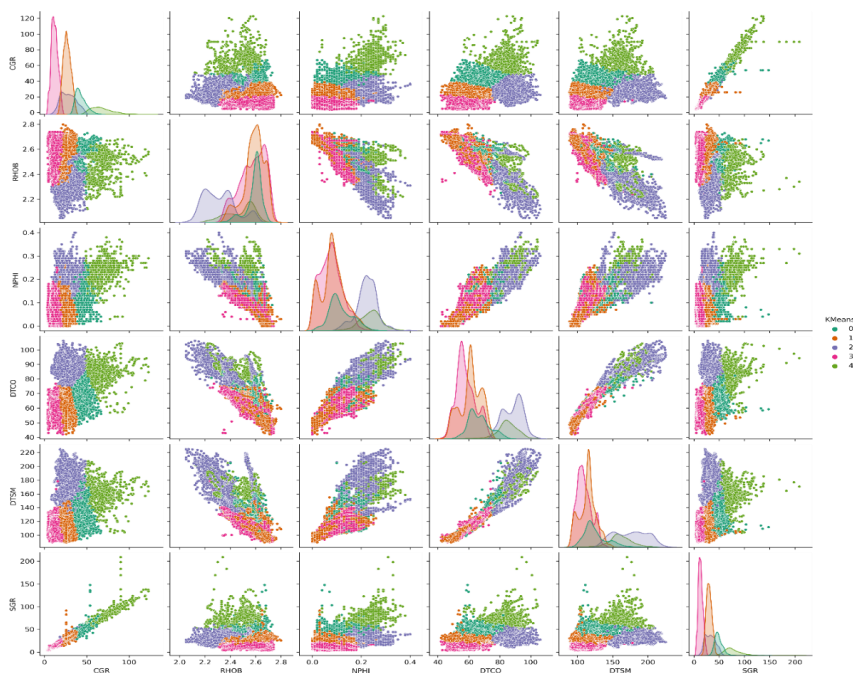
شکل ۱۳. از بین روش‌های فوق روش خوشه‌بندی سلسله مراتبی،  $k$ - میانگین و مدل آمیخته گوسی بهتر از بقیه عمل کرده است. مزیت بهتر روش GMM انتخاب چندین الگو احتمالاتی مختلف در زمان خروجی از پایتون می‌باشد.

خوشه‌بندی  $K$ -Means عالی عمل می‌کند. بنابراین بهتر است از روش مدل‌های مخلوط گاوسی (GMM) که در تشخیص لیتولوژی تفکیک بهتر و مطابقت بیشتری با سازند زمین‌شناسی دارد. استفاده شود. در نتیجه با توجه مطالب گفته‌شده، شکل ۱۷ جهت شناسایی لایه زمین‌شناسی ارائه شده است. در این شکل لایه‌های زمین به پنج واحد تقسیم‌بندی شده است. با توجه به بررسی‌های زمین‌شناسی  $Facies1$  با رنگ زرد معادل پکستون‌های با تخلخل بالا مربوط به آهک و دولومیت‌های سازند آسماری، جهرم، ایلام و سروک می‌باشد. لایه‌های مخزنی پرفشار ایلام و سروک بعد از بررسی و اعتبار سنجی با تحقیقات میدانی به‌درستی تشخیص داده شده است،  $Facies2$  با رنگ سبز کم‌رنگ مربوط به شیل و ماسه‌سنگ‌های گورپی می‌باشد،  $Facies3$  با رنگ سبز تیره با توجه به درصد کم رس مربوط به میان لایه‌های آهک نازک‌لایه سازند آسماری و گرینستون‌های محیط سدی سازند سروک می‌باشد  $Facies4$  مربوط به وکستون‌های سازند سروک و  $Facies5$  مربوط به ماسه‌سنگ‌های ناپایدار سازند آسماری می‌باشد. همان‌طور که مشخص است میان لایه‌های آهک نازک‌لایه دارای درز و شکاف که مسئول اصلی ناپایداری دیواره چاه می‌باشد در این واحد به‌درستی تشخیص داده شده است علت ریزش دیواره چاه (ناپایداری) در لاگ کالیپر به‌درستی در سازند آسماری دیده می‌شود و همین‌طور لایه‌های مخزنی پرفشار ایلام و سروک به‌درستی نشان داده شده است. در نتیجه با توجه به این تحقیق چهار روش یادگیری ماشین به‌کار گرفته شد که با وجود عملکرد خوب روش خوشه‌بندی سلسله مراتبی، مدل کی میانگین و مدل آمیخته گوسی، روش مدل مخلوط گوسی جهت پیش‌بینی دقیق‌تر زون‌های ژئومکانیکی پیشنهاد داده می‌شود. جهت شناخت بیشتر این زون‌ها درصد و فراوانی این زون‌ها در شکل ۱۸ نشان داده می‌شود.

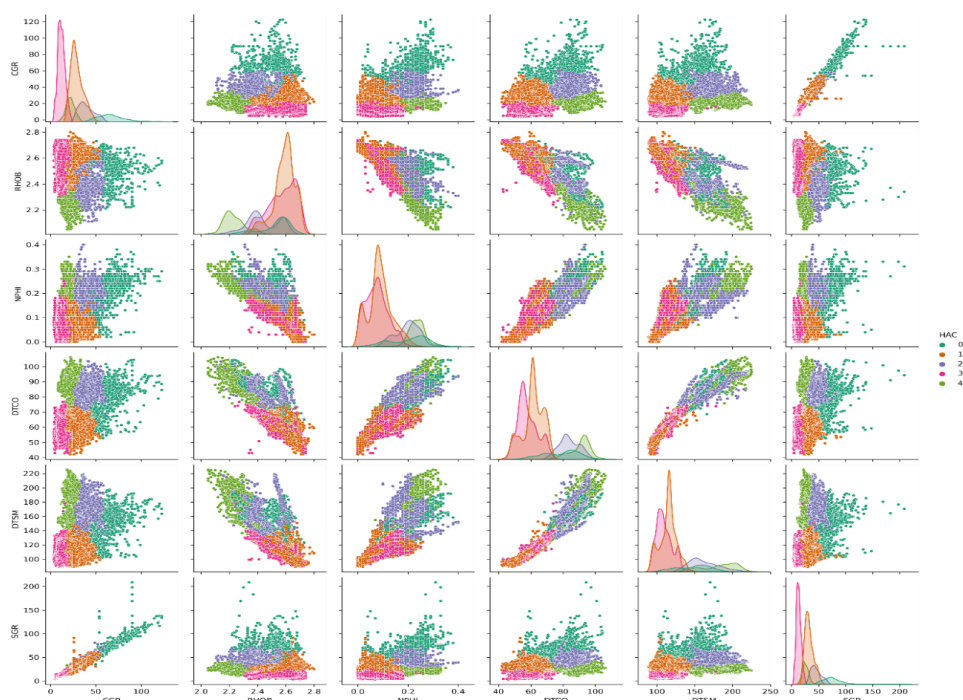
## ۲.۵ نمودارهای جفتی (Pair plot)

در این نمودارها از شش منحنی ورودی برای مدل‌سازی استفاده شده است تا تفاوت خوشه‌ها نسبت به هم و ارتباط هر خوشه با پارامترهای به‌کار گرفته شده مشخص‌تر باشد. این راه برای انجام این کار استفاده از نمودارهای زوجی می‌باشد. این نمودار روابط بین داده‌های مجموعه هر داده را در یک شبکه نمایش می‌دهد. همان‌طور که مشخص است این یک راه سریع و آسان برای شناسایی و تجسم داده‌ها می‌باشد. در امتداد مورب توزیع داده‌های تقسیم‌شده توسط خوشه نیز رسم شده است. در نمودار  $k$ - میانگین شکل ۱۴ و خوشه‌بندی سلسله مراتبی شکل ۱۵ زون‌های با رنگ قهوه‌ای کم‌رنگ، صورتی و آبی به ترتیب بیشترین حجم سازندها را به خود اختصاص داده‌اند و اختلاط خوشه‌ها در بعضی از این نمودارها به چشم می‌خورد. اما در شکل ۱۶ نمودار GMM تفکیک خوشه خیلی مشخص‌تر و بهتر نمایان شده است. ملاحظه می‌شود که مدل GMM بهبودهایی را در تعریف خوشه‌ها به‌خصوص در نمودار DTC در مقابل RHOB ارائه می‌دهد. بنابراین با توجه به این نمودارها و تحقیقات آزمایشگاهی سازند بنگستان از سه نوع لیتولوژی آهکی- دولومیتی با تخلخل بالا، متوسط و پایین تقسیم‌بندی می‌شود که در مقایسه با رخساره‌های رسوب‌شناسی به ترتیب شامل پکستون- گرینستون دارای تخلخل حفره‌ای و شکافی، گرینستون دارای تخلخل شکافی و وکستون- پکستون دارای تخلخل کمتر معادل است.

الگوریتم‌های Hierarchical K-Means Clustering و Gaussian Mixture Modeling و Agglomerative Clustering سه روش کاربردی در تحلیل خوشه‌بندی بدون نظارت یادگیری ماشین هستند اما در واقعیت‌های پتروفیزیکی، زمین‌شناسی و ژئومکانیکی داده‌ها به‌ندرت الگوهای دایره‌ای خوبی را تشکیل می‌دهند. درحالی‌که اگر خوشه‌های داده دایره‌ای باشند،

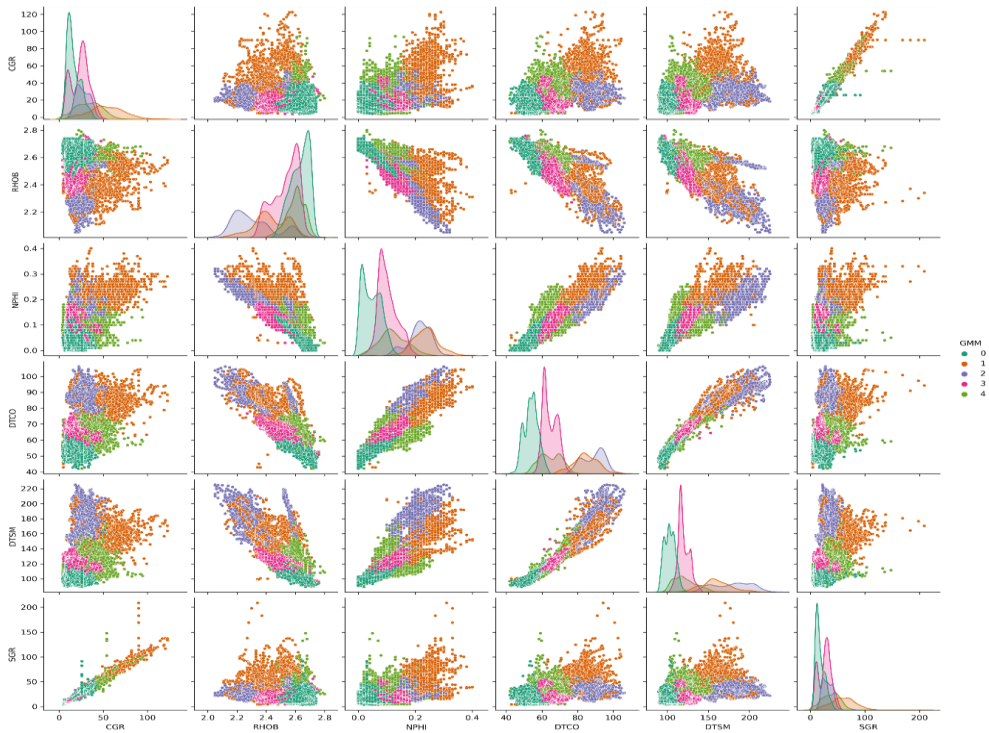


شکل ۱۴. با توجه به نمودار کی میانگین تفکیک خوشه‌ها را خوب نشان داده است. اما در بعضی نمودارها مثل DTCCO- RHOB اختلاط خوشه‌ها دیده می‌شود.

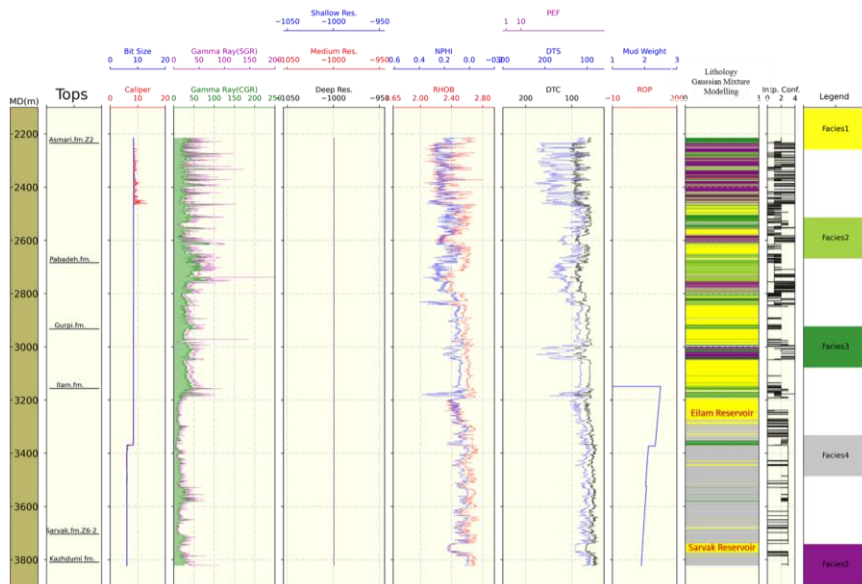


شکل ۱۵. با توجه به نمودار خوشه‌بندی سلسله مراتبی (HAC) تفکیک خوشه‌ها خوب نشان داده شده است.

تعیین واحدهای ژئومکانیکی با استفاده از ...

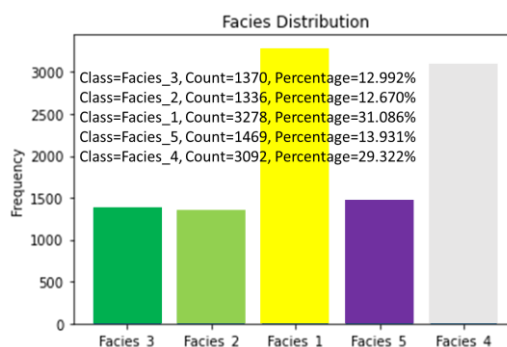


شکل ۱۶. در این روش داده‌ها بر اساس چندین احتمال مختلف خوشه‌بندی می‌شوند و نسبت به روش‌های دیگری انعطاف‌پذیر می‌باشد



شکل ۱۷. استفاده از روش (GMM) در تعیین واحدهای ژئومکانیکی در یکی از چاه‌های جنوب ایران. در این تصویر میان لایه‌های آهک نازک لایه دارای درز و شکاف با رنگ سبز و خاکستری که مسئول ناپایداری دیواره چاه در سازند آسماری (Facies 5) می‌باشد. همین‌طور سازندهای مخزنی پرفشار ایلام و سروک با رنگ زرد به‌درستی شناسایی شدند.

کی میانگین ((*K-Means Clustering*), الگوریتم خوشه بندی *DBSCAN* مبتنی بر غلظت و مدل آمیخته گوسی (*Gaussian Mixture Modelling*) جهت تعیین واحدهای ژئومکانیکی غیر مخزنی مسئله دار، آهک های نازک لایه و مخزنی پرفشار استفاده شد. در اولین قدم با استفاده از یک روش بهینه سازی و با استفاده از نمودار زانویی و همین طور استفاده از نمودار دندروگرام (*Dendrogram*) از روش خوشه بندی سلسله مراتبی (*Agglomerative hierarchical Clustering*) بهترین تعداد خوشه ها تعیین شد که تعداد بهینه خوشه ها عدد پنج انتخاب گردید. مقایسه نتایج نشان داد از آنجایی که در روش خوشه بندی *HAC* و *K-Means* جهت پیش بینی واحدهای ژئومکانیکی از خوشه های کروی استفاده می کنند، جهت مدل سازی مناسب تر است از الگوریتم *GMM* که از الگو مرزهای خوشه/تضمیم بیضی شکل استفاده می کند، انعطاف پذیرتر است و سازگاری بیشتری با داده های مطالعه حاضر نشان می دهد. بنابراین مدل *GMM* به عنوان روش مناسب تر جهت تعیین واحدهای غیر مخزنی مسئله دار، آهک های نازک لایه در سازند آسماری و سازندهای مخزنی پرفشار ایلام و سروک معرفی شد. همین طور با این مدل بخش های پرفشار سازندهای مخزنی ایلام و سروک با دقت خوب شناسایی شدند.



شکل ۱۸. در این شکل درصد و فراوانی واحدهای ژئومکانیکی نشان داده می شود.

#### ۶. نتیجه گیری

این مقاله روش انجام تحلیل خوشه های بدون نظارت با استفاده از چهار الگوریتم *K-Means*, *HAC*, *DBSCAN* و *Gaussian Mixture Modeling* پوشش داده شده است. در مرحله اول، بعد از پردازش داده های چاهنگاری از دو الگوریتم محبوب قدرتمند نظارت شده عمیق یادگیری ماشین ایکس جی بوست (*XGBoost*) و شبکه عصبی پرسپترون چندلایه *Perceptron Neural Multi-Layer Network* جهت بازیابی داده های گمشده استفاده گردید. سپس از روش های بدون نظارت یادگیری ماشین شامل مدل خوشه بندی سلسله مراتبی ((*HAC*), مدل

#### ۷. مراجع:

- Ali, M., Ma, H., Pan, H., Ashraf, U., Jiang, R., (2020). Building a rock physics model for the formation evaluation of the Lower Goru sand reservoir of the Southern Indus Basin in Pakistan. *J. Petrol. Sci. Eng.* <https://doi.org/10.1016/j.petrol.2020.107461>
- Ashraf, U., Zhu, P., Yasin, Q., Anees, A., Imraz, M., Mangi, H.N., Shakeel, S., (2019). Classification of reservoir facies using well log and 3D seismic attributes for prospect evaluation and field development: a case study of Sawan gas field, Pakistan. *J. Petrol. Sci. Eng.* 175, 338–351. <https://doi.org/10.1016/j.petrol.2018.12.060>
- Ashraf, U., Zhang, H., Anees, A., Mangi, H.N., Ali, M., Ullah, Z., Zhang, X., (2020a). Application of unconventional seismic attributes and unsupervised machine learning for the identification of fault and fracture network. *Appl. Sci.* <https://doi.org/10.3390/app10113864>
- Ashraf, U., Zhang, H., Anees, A., Ali, M., Zhang, X., Abbasi, S.S., Mangi, H. N., (2020b). Controls on Reservoir Heterogeneity of a Shallow-Marine Reservoir in Sawan Gas Field, SE Pakistan: Implications for Reservoir Quality Prediction Using Acoustic Impedance Inversion. *Water* 12 (11), 2972.

<https://doi.org/10.3390/w12112972>. In this issue

Bestagini, P., Lipari, V., Tubaro, S., aug, (2017). A machine learning approach to facies classification using well logs. In: SEG Technical Program Expanded Abstracts 2017. Society of Exploration Geophysicists, pp. 2137–2142. <https://doi.org/10.1190/segam2017-17729805.1>.

Forgy, E. W., (1965). "Cluster analysis of multivariate data: efficiency versus interpretability of classifications". *Biometrics*.

Jinghua, W., Jianmin J., (2021). Unsupervised deep clustering via adaptive GMM modeling and optimization. *Neurocomputing* 433 (2021) 199-211

Hall, B., oct, (2016). Facies classification using machine learning. *Lead. Edge* 35 (10), 906–909. <https://doi.org/10.1190/tle35100906.1>.

Lppolito, M., Ferguson, J., Jenson, F., (2021): Improving facies prediction by combining supervised and unsupervised learning methods. *Journal of Petroleum Science and Engineering*, Volume 200, May 2021, 108300. <https://doi.org/10.1016/j.petrol.2020.108300>

Ghalibaf, H., Ghafoori, M., Lashkaripoor, G.R., Hafezi Moghaddas, N., (2020). Preparation of Zoning Maps for Cut-off Wall Using the Geotechnical parameters and Analytic Hierarchal Process (AHP). Case study: Sarroud Dam "Journal of Dam and Hydroelectric PowerPlant 7th Year / No. 26 / December 2020

Kriegel, H.P., Schubert, E., Zimek, A., (2016). The (black) art of runtime evaluation: Are we comparing algorithms or implementations. *Knowledge and Information Systems*. 52: 341–378. doi:10.1007/s10115-016-1004-2. ISSN 0219-1377

Marco, I., John, F., Fred, J., (2021). Improving facies prediction by combining supervised and unsupervised learning methods: *Journal of Petroleum Science and Engineering* 200 (2021) 108300

Muhammad Ali, A., Ren Jiang, B., Huolin Ma, A., Heping Pan, A., Khizar Abbas, C., Umar Ashraf, D., (2021). Machine learning - A novel approach of well logs similarity based on synchronization measures to predict shear sonic logs: *Journal of Petroleum Science and Engineering* 203 (2021) 108602

MacQueen, J. B. (1967). Some Methods for classification and Analysis of Multivariate Observations. *Proceedings of 5th Berkeley Symposium on Mathematical Statistics and Probability*. 1. University of California Press. pp. 281–297. MR 0214227. Zbl 0214.46201. Retrieved 2009-04-07.

MacKay, D., (2003). "Chapter 20. An Example Inference Task: Clustering" (PDF). *Information Theory, Inference and Learning Algorithms*. Cambridge University Press. pp. 284–292. ISBN 0-521-64298-1. MR 2012999.

Song, C., Li, L., Li, K., (2021). Robust K-means algorithm with weighted window for seismic facies analysis: *Journal of Geophysics and Engineering* (2021) 18, 618–626. <https://doi.org/10.1093/jge/gxab039>

Steinhaus, H. (1957). "Sur la division des corps matériels en parties". *Bull. Acad. Polon. Sci. (French)*. 4 (12): 801–804. MR 0090073. Zbl 0079.16403.

Thiago Santi, B., Marcelo Kehl, D., Tiago, J., Girelli, F., Chemale, J., (2020) Evaluation of machine learning methods for lithology classification using geophysical data: *Computers and Geosciences* 139(2020)104475

W. Dunham, M., Malcolm, A., Welford, J. K., (2020). Improved well log classification using semisupervised Gaussian mixture models and a new hyper-parameter selection strategy: *Computers and Geosciences* Volume 140, July 2020, 104501. <https://doi.org/10.1016/j.cageo.2020.104501>