



ترموگراف: ترکیب روش‌های متن کاوی و گراف کاوی به منظور پیش‌بینی اصناف در شبکه‌ی پرداخت

ارائه‌دهنده: بهناز پورابراهیم

9th Annual Conference on
Electronic Banking &
Payment Systems



شرکت ملی انفورماتیک



بانک مرکزی جمهوری اسلامی ایران



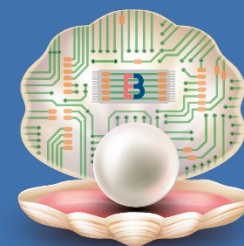
پژوهشکده پولی و بانکی
بانک مرکزی جمهوری اسلامی ایران



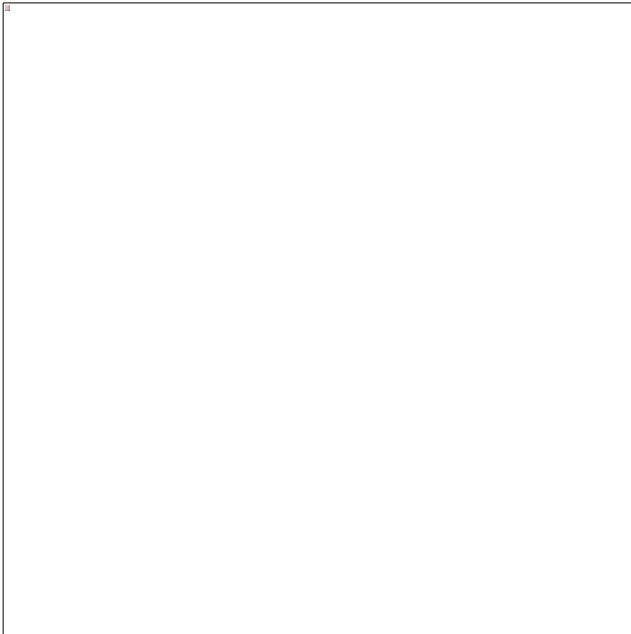
نهمین همایش سالانه
بانکداری الکترونیک
و نظام‌های پرداخت

(ارزش آفرینی دیجیتالی)

تهران، مرکز همایش‌های بین‌المللی برج میلاد - ۱ و ۲ اسفند ۱۴۰۱



مقدمه



۱ شبکه پرداخت الکترونیک

۲ پذیرنده و کد صنف مرتبط
با فعالیت آن

ناهنجاری صنف



چالش‌ها

وجود ناهنجاری صنف و
عدم تشخیص آن برای
پذیرندگان



هدف

تشخیص ناهنجاری
اصناف



تعریف

عدم ثبت کد صنف
مرتبط با فعالیت
پذیرندگان

هدف



۱ تجزیه و تحلیل نام پذیرندگان

۲ تشخیص نوع فعالیت صنفی هر پذیرنده

۳ پیش‌بینی صنف درست جایگزین

۴ تکنیک‌های متن‌کاوی و گراف‌کاوی



مجموعه داده

➤ مجموعه S شامل صنف‌های ۱ تا N هستند.

$$S = \{s_1, s_2, \dots, s_N\}$$

➤ هر صنف s_i شامل k پذیرنده (p) است.

$$s_i = \{p_1, p_2, \dots, p_k\}$$

➤ مجموعه داده‌ی ما مجموعه نام فروشگاه‌ی و شامل هشت میلیون پذیرنده است.

$$Dataset = \{p_1, p_2, \dots, p_{8000000}\}$$



مدل پیشنهادی: ترموگراف

تحلیل متن کاوی نام پذیرندگان و محاسبه‌ی شباهت معنایی بین اصناف

ساخت گراف میان اصناف با توجه به شباهت معنایی اصناف

پیش‌بینی کد صنف پذیرندگان با رویکردهای گراف کاوی و خوشه‌بندی

قدم اول

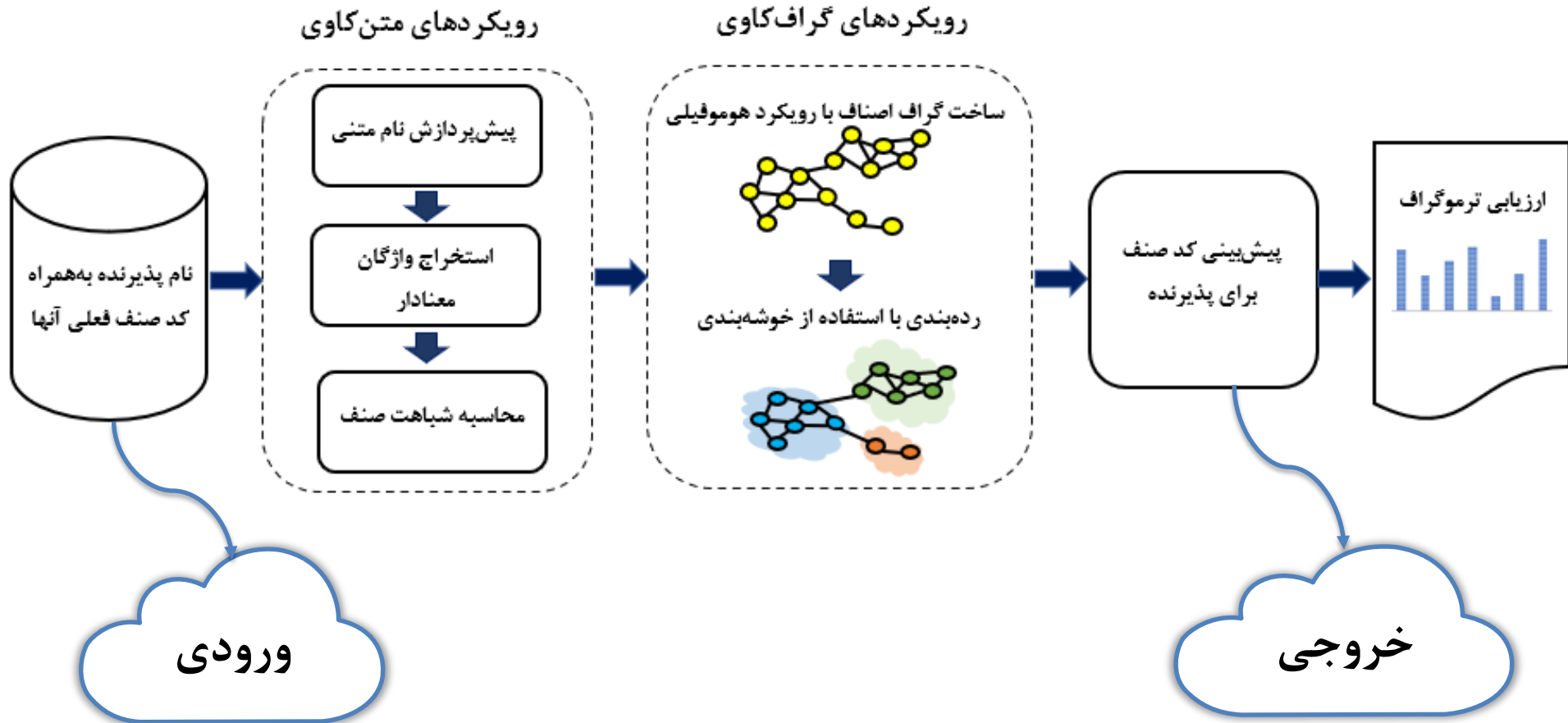
قدم دوم

قدم سوم

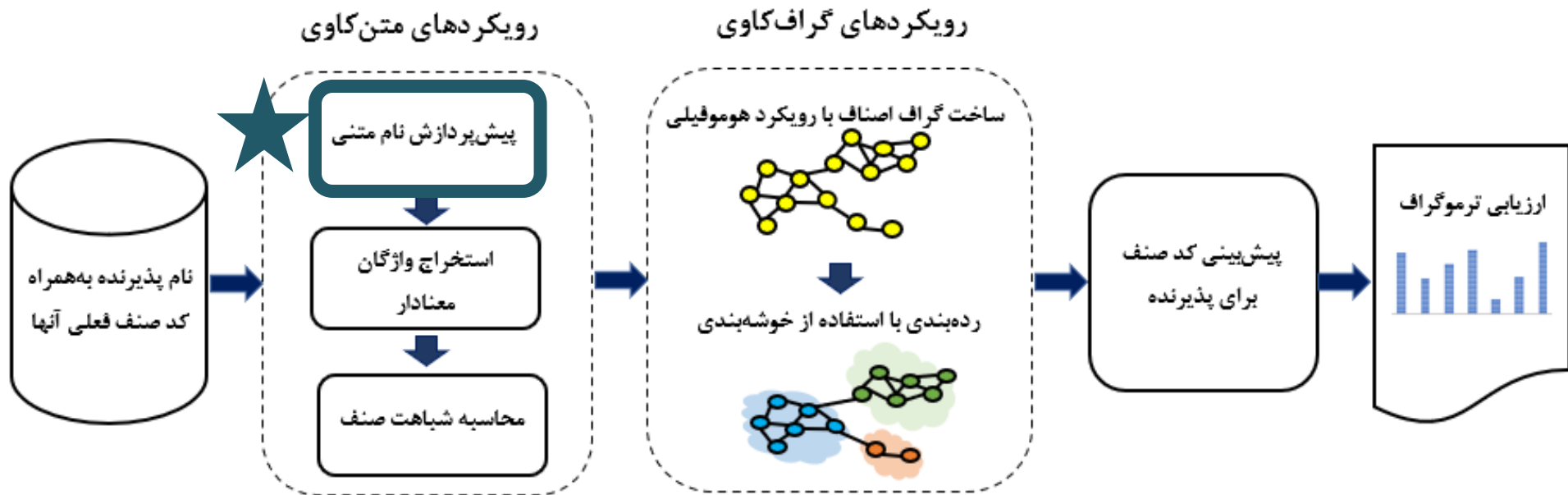
نوآوری‌ها

الف) استفاده از داده‌های متنی نام پذیرندگان
ب) رویکرد ترکیبی جدید متن کاوی و گراف کاوی

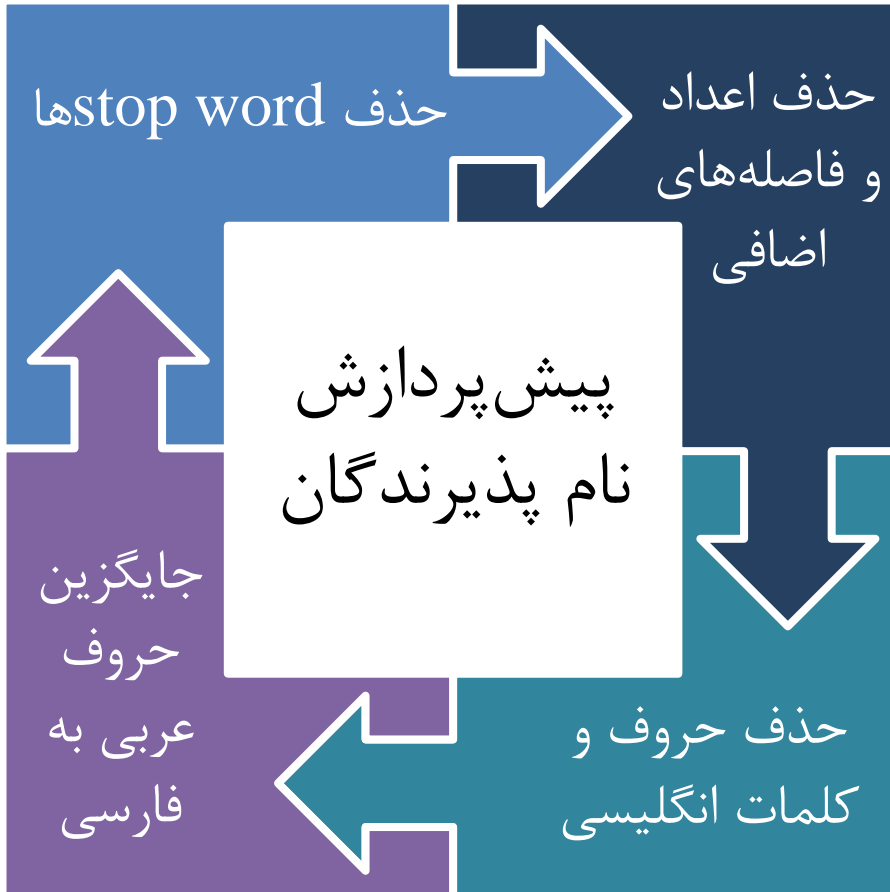
معماری کلی ترموگراف



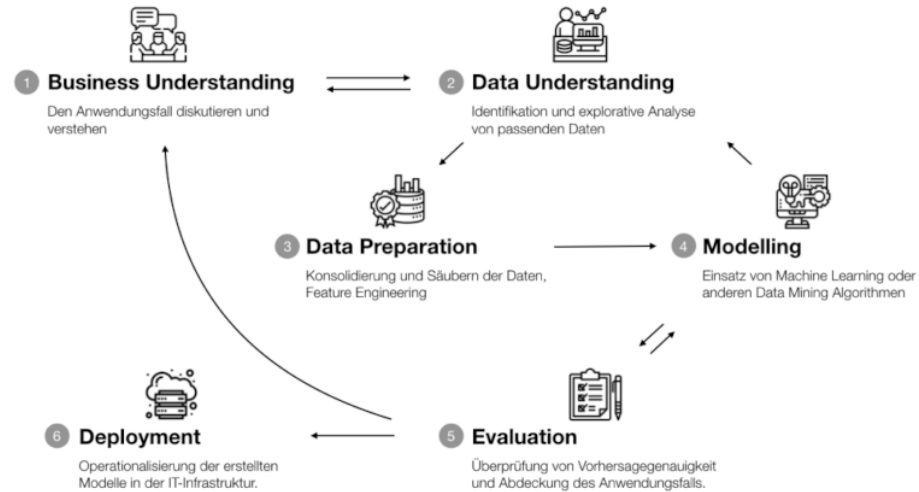
معماری کلی ترموگراف



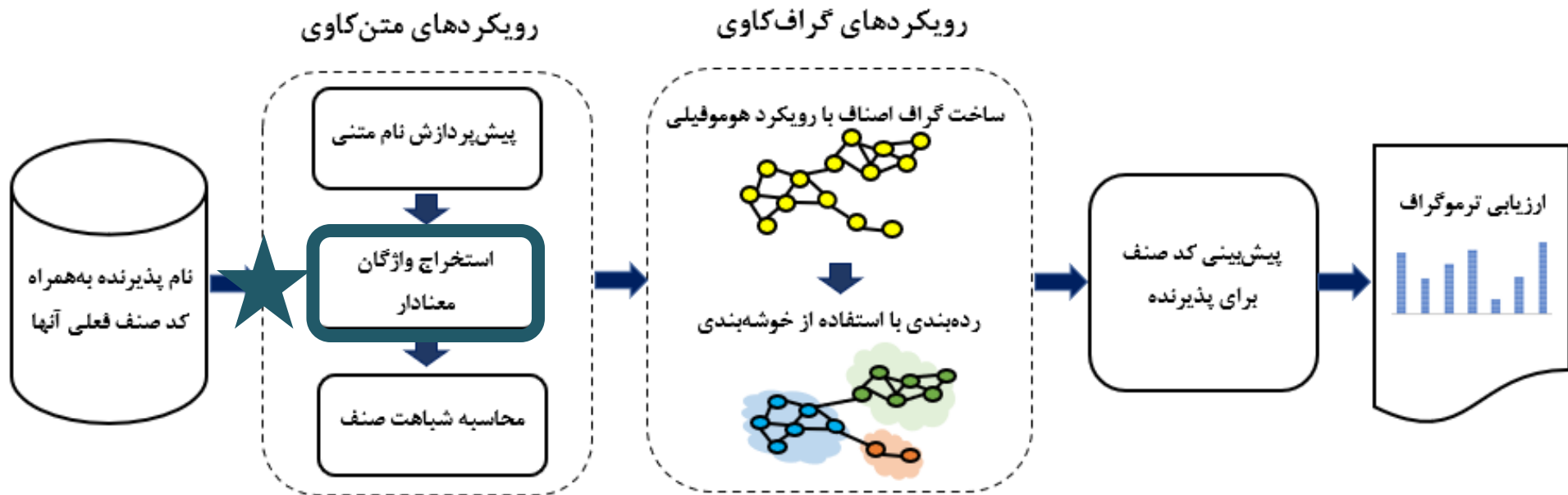
پیش پردازش داده‌های متنی



CRISP DM



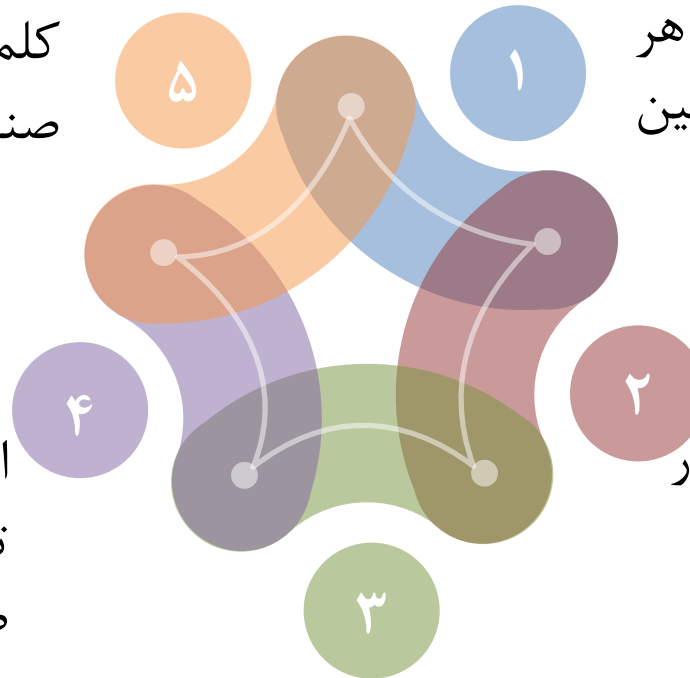
معماری کلی ترموگراف



استخراج واژگان معنادار

کلمات معنادار و مرتبط با
صنف را استخراج می‌کنیم.

نام پذیرنده‌های موجود در هر
صنف براساس فاصله‌ی بین
کلمات استخراج می‌کنیم.



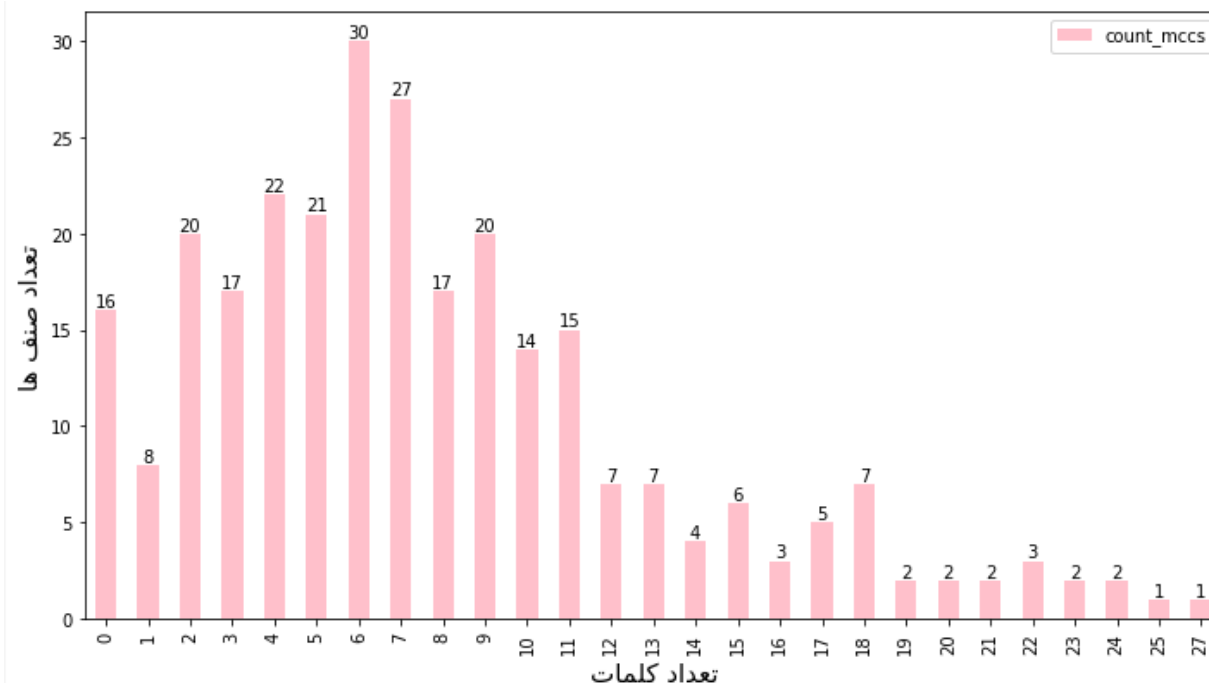
تعداد تکرار هر کلمه در
صنف نیز محاسبه می‌شود.

۱۰۰ کلمه‌ی پرتکرار هر
صنف به‌عنوان کلمات کاندید
صنف استخراج می‌شوند.

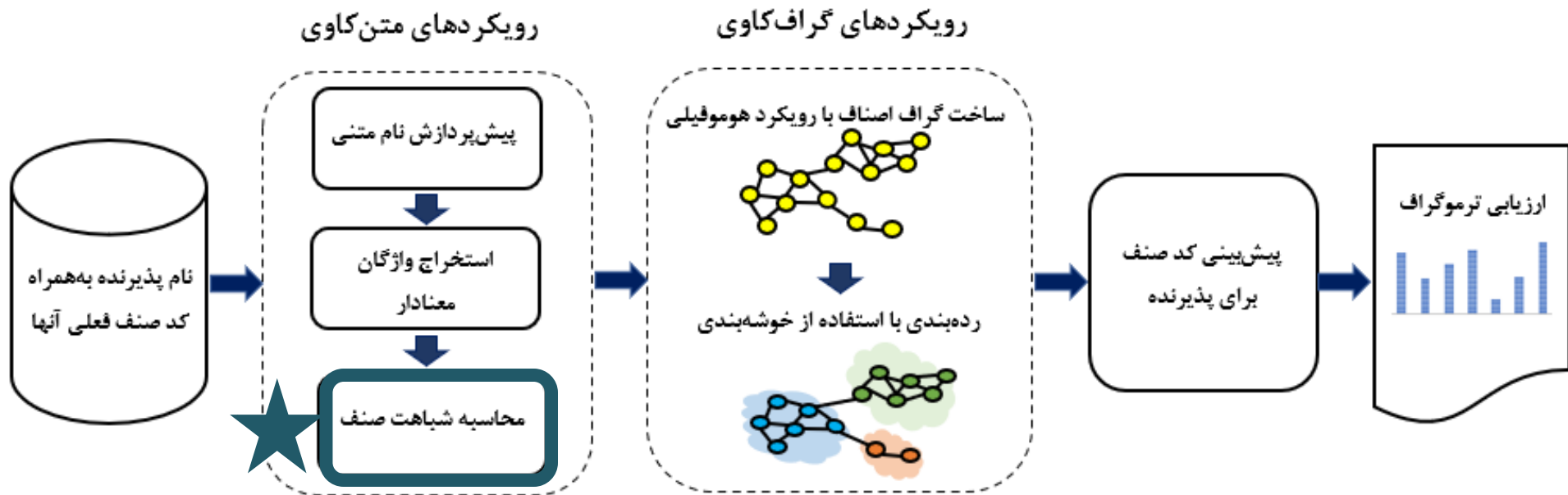
از بین ۱۰۰ کلمه‌ی کاندید،
تعدادی کلمه مرتبط معنادار با
صنف و تعدادی کلمه نامرتب با
صنف وجود دارد.

تنوع و پراکندگی نام پذیرندگان

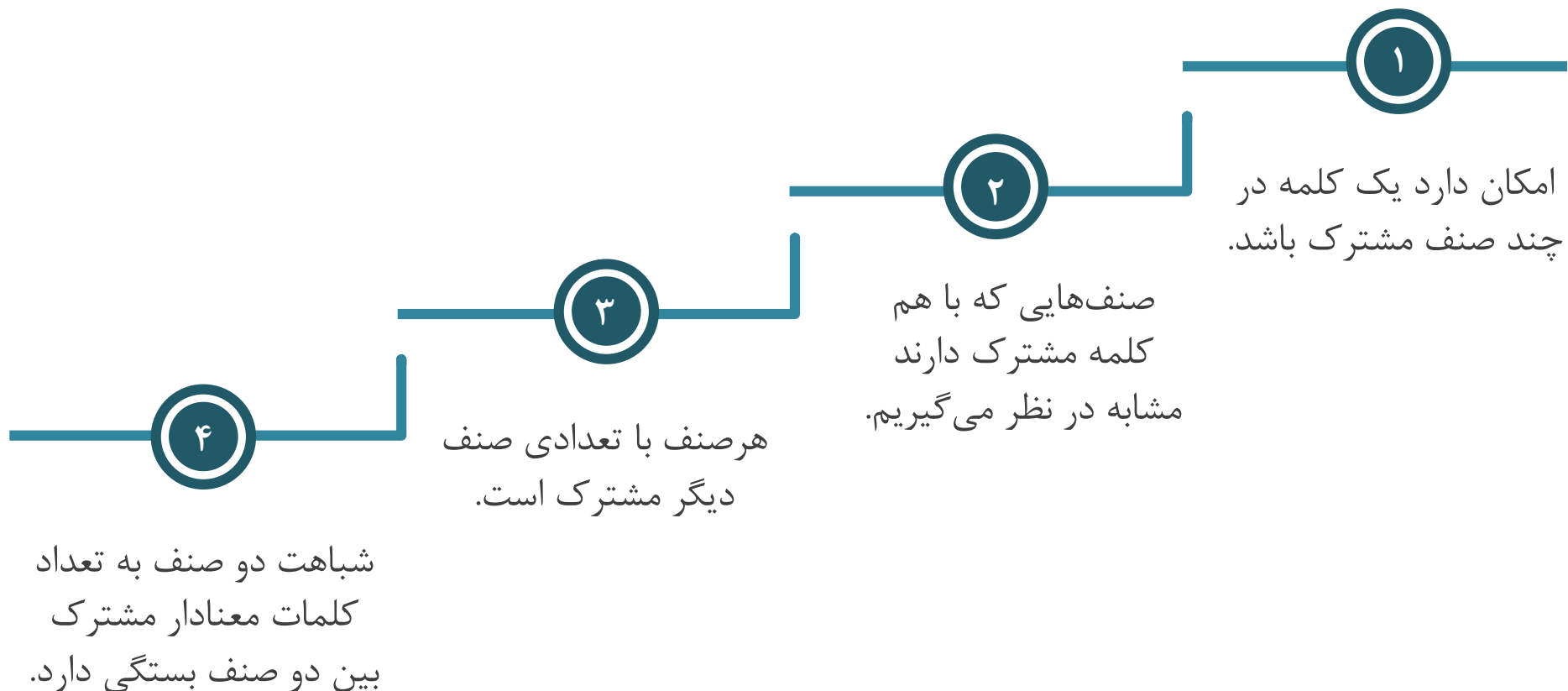
- برای بعضی از صنوف تعداد کلمات معنادار زیادی انتخاب شده است که دلیل آن تنوع نام پذیرنده‌ها است.
- برای بعضی از صنوف تعداد کلمات معنادار کمی انتخاب شده است که دلیل آن پراکندگی نام پذیرنده‌ها است.



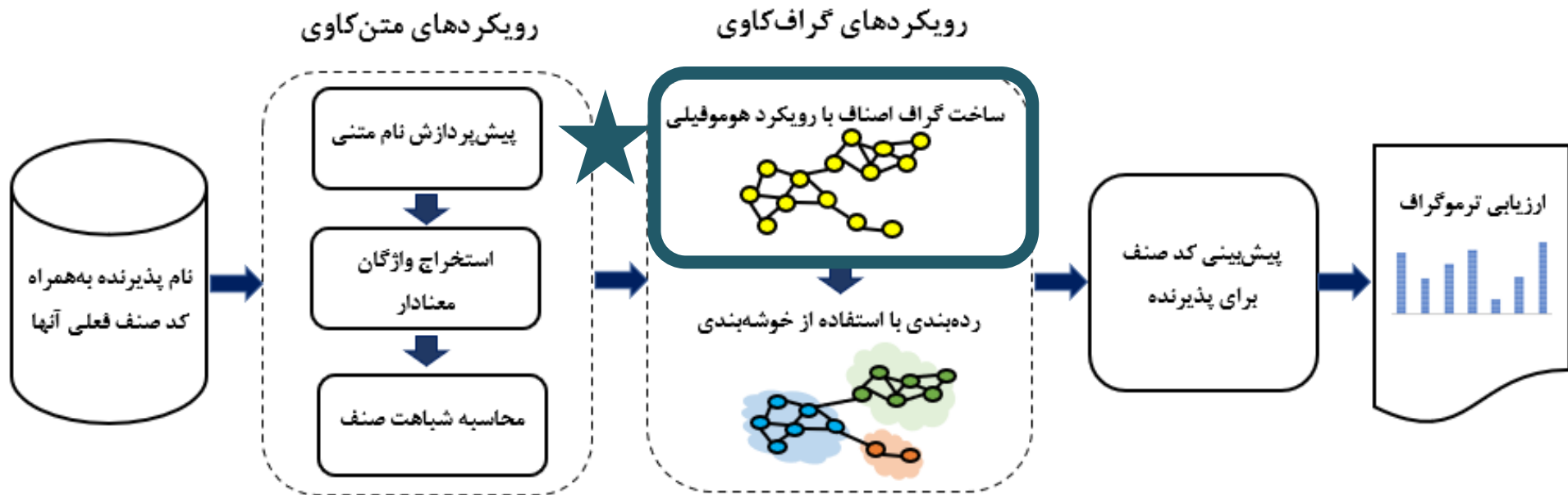
معماری کلی ترموگراف



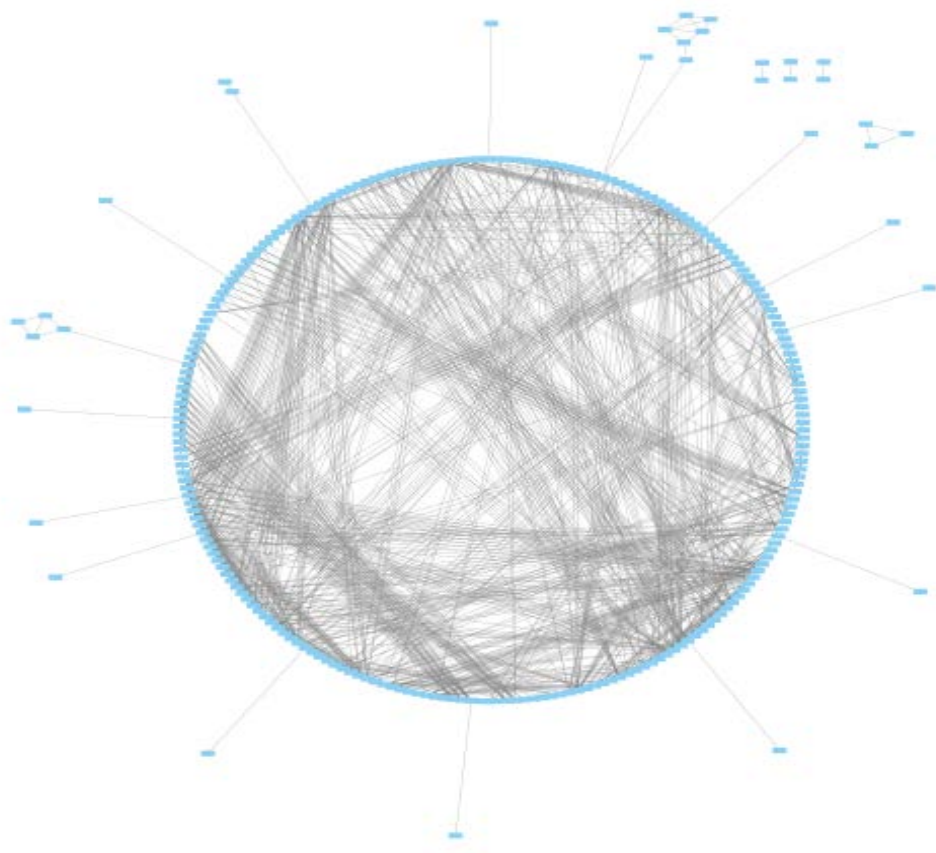
محاسبه‌ی شباهت اصناف



معماری کلی ترموگراف



تحلیل گرافی اصناف

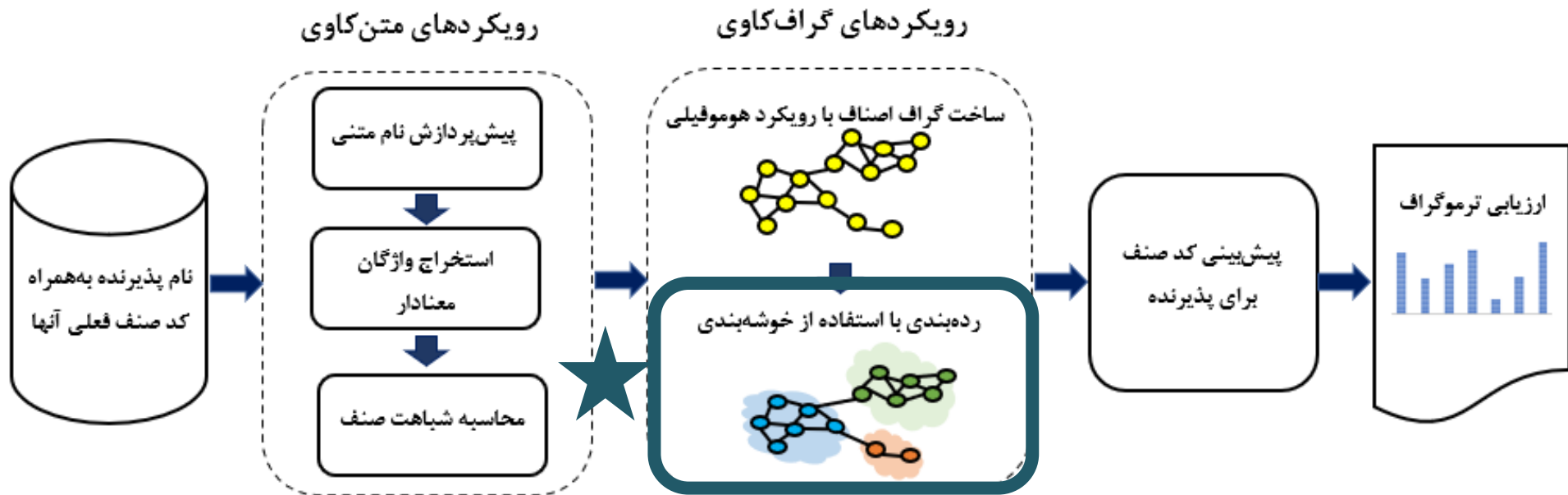


❖ گره: صنف

❖ یال: کلمه معنادار مشترک

➤ تحلیل گرافی اصناف: اگر بین هر دو صنف حداقل یک کلمه‌ی مشترک وجود داشته باشد، دو صنف با یک یال به هم متصل می‌شوند.

معماری کلی ترموگراف



خوشه‌بندی اصناف

چالش تحلیل ریزدانه

عدم وجود برچسب

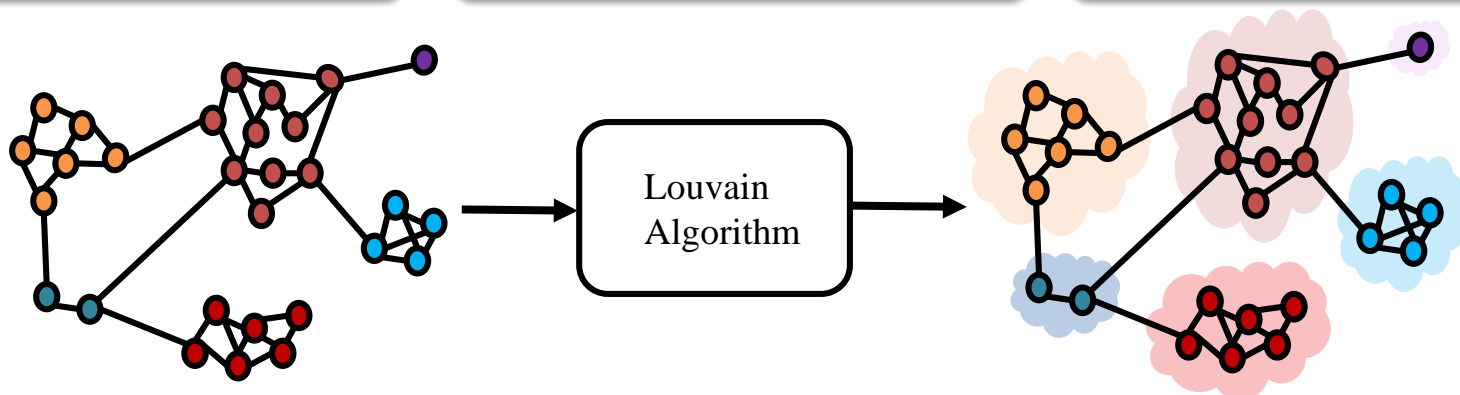
نیاز به تحلیل درشت دانه

خوشه بندی (تحلیل بدون ناظر)

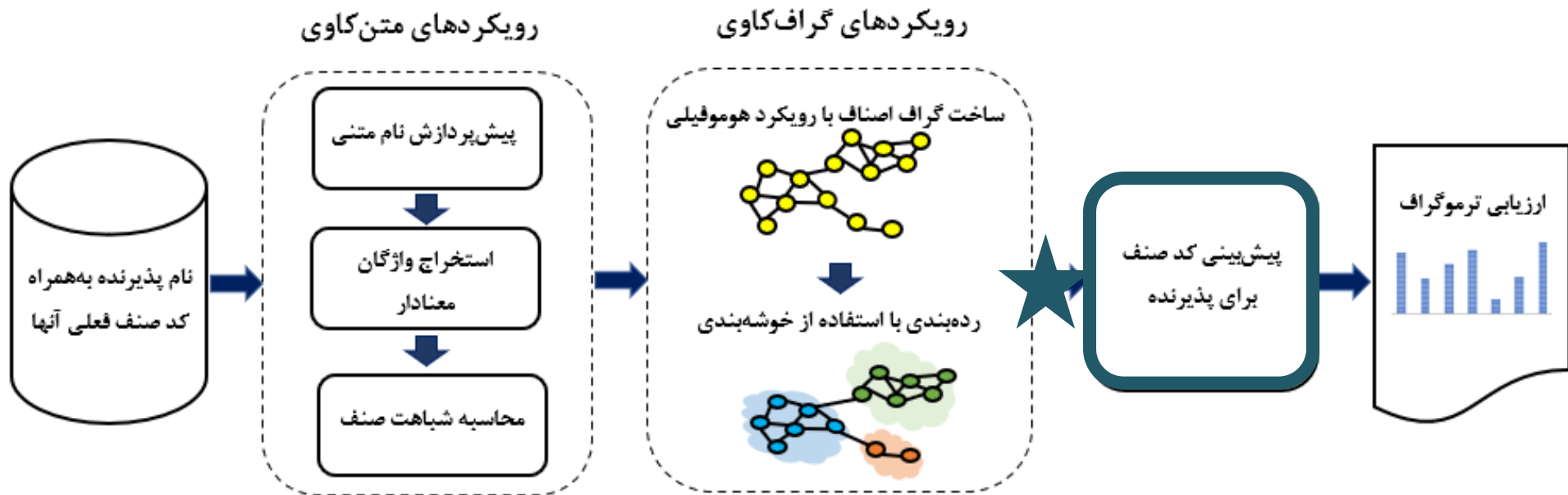
خوشه بندی گرافی Louvain

شباهت بین خوشه ای پایین

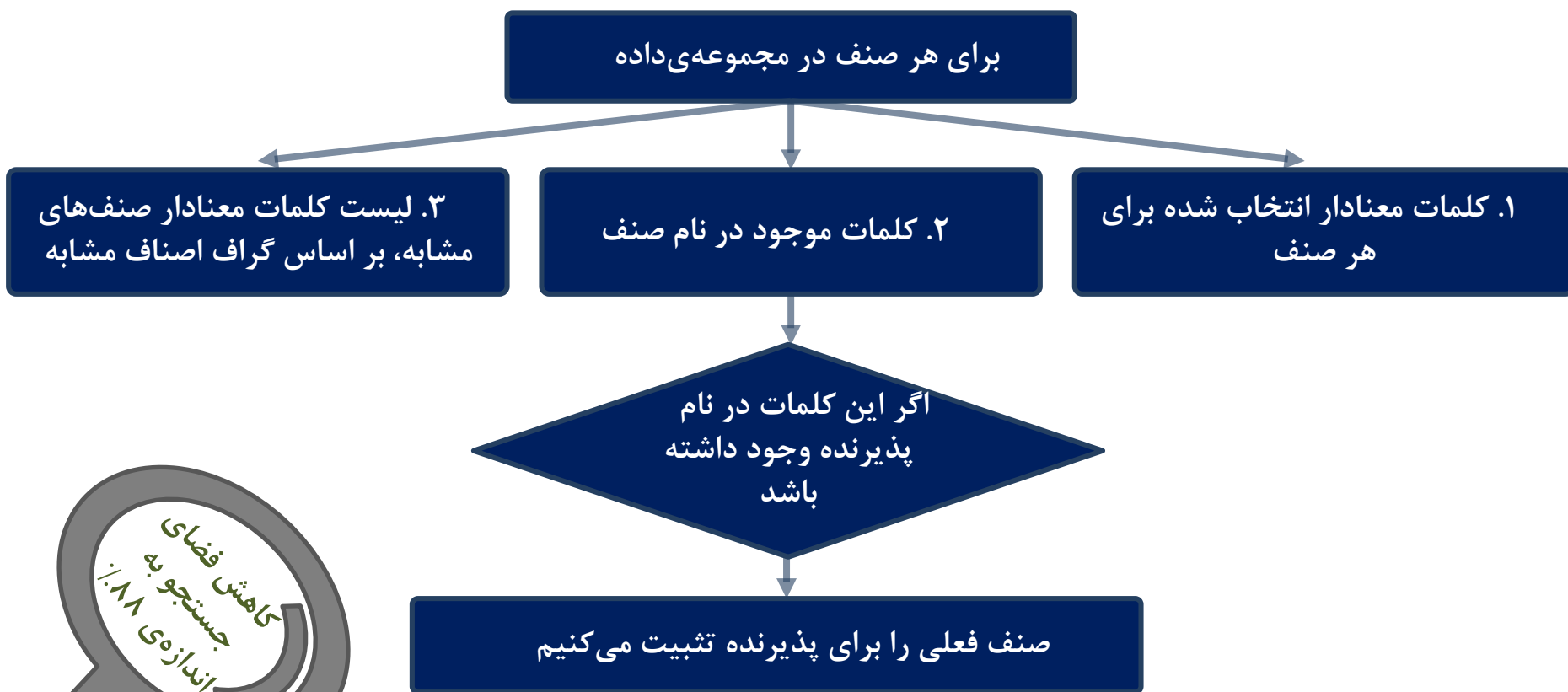
شباهت درون خوشه ای بالا



معماری کلی ترموگراف

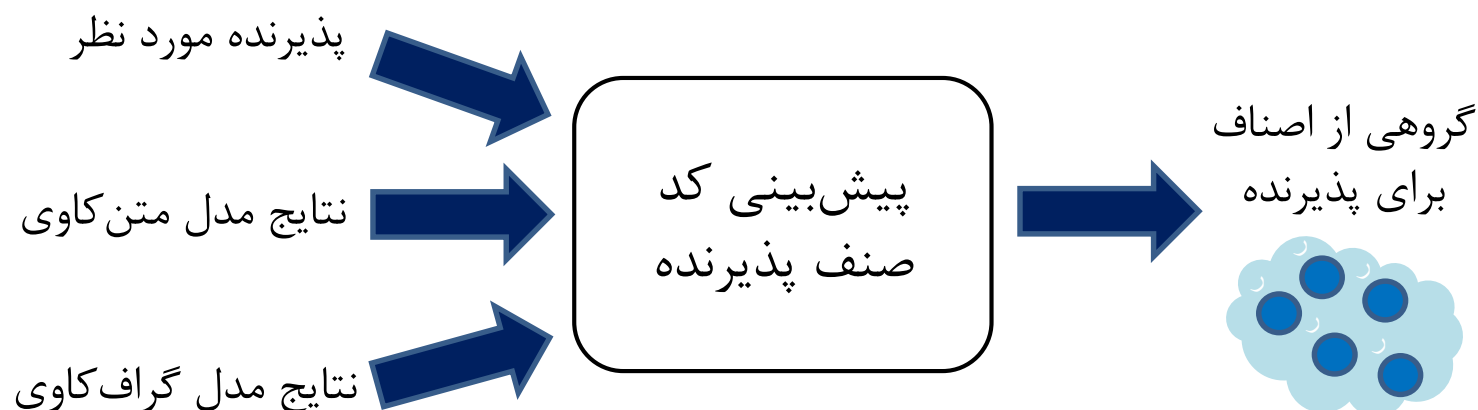


فیلترهای اعمال شده روی نام پذیرنده

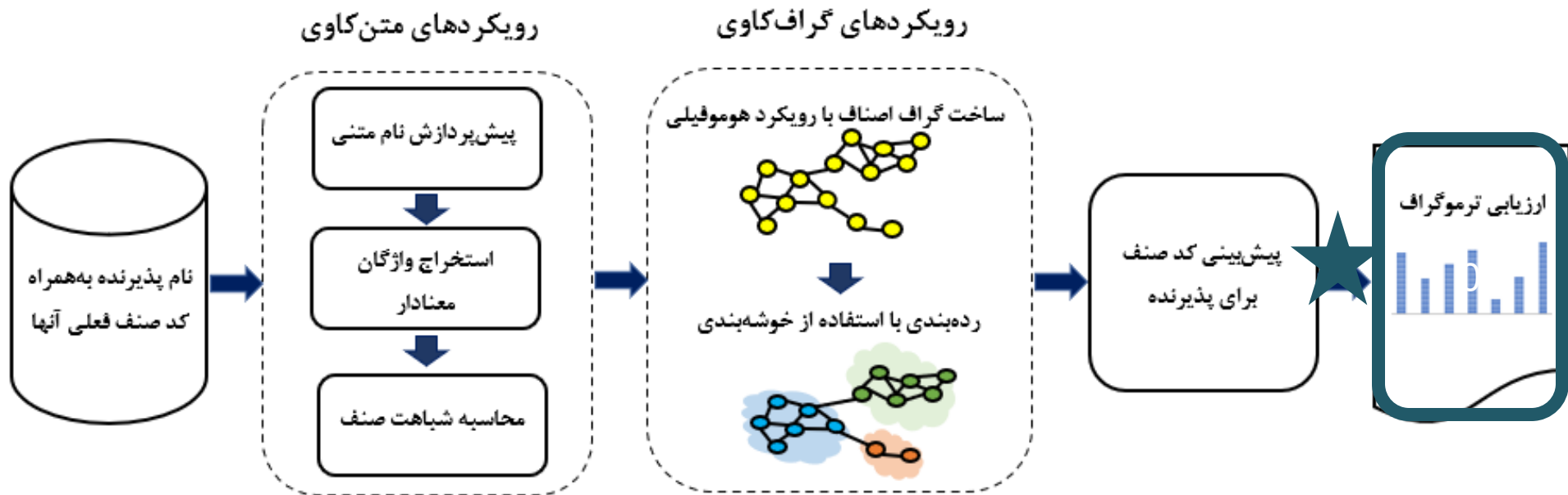


پیش‌بینی کد صنف پذیرنده

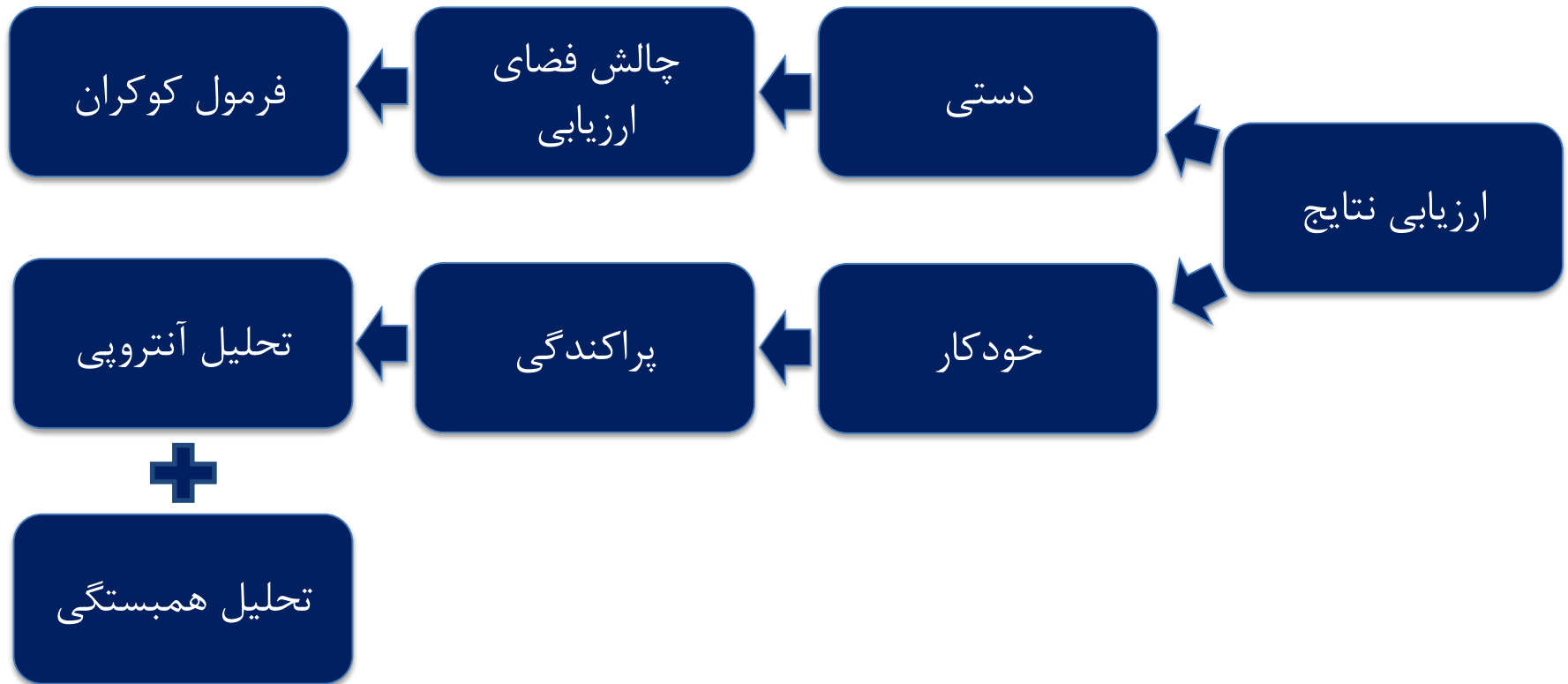
- کلمات معنادار هر صنف را به لیست تبدیل می‌کنیم.
- اگر از این کلمات در نام پذیرنده‌ای وجود داشته باشد، در نتیجه صنف فعلی پذیرنده اشتباه ثبت شده.
- صنف مورد نظر و خوشه‌ای که صنف در آن قرار دارد برای آن پیش‌بینی می‌شود.



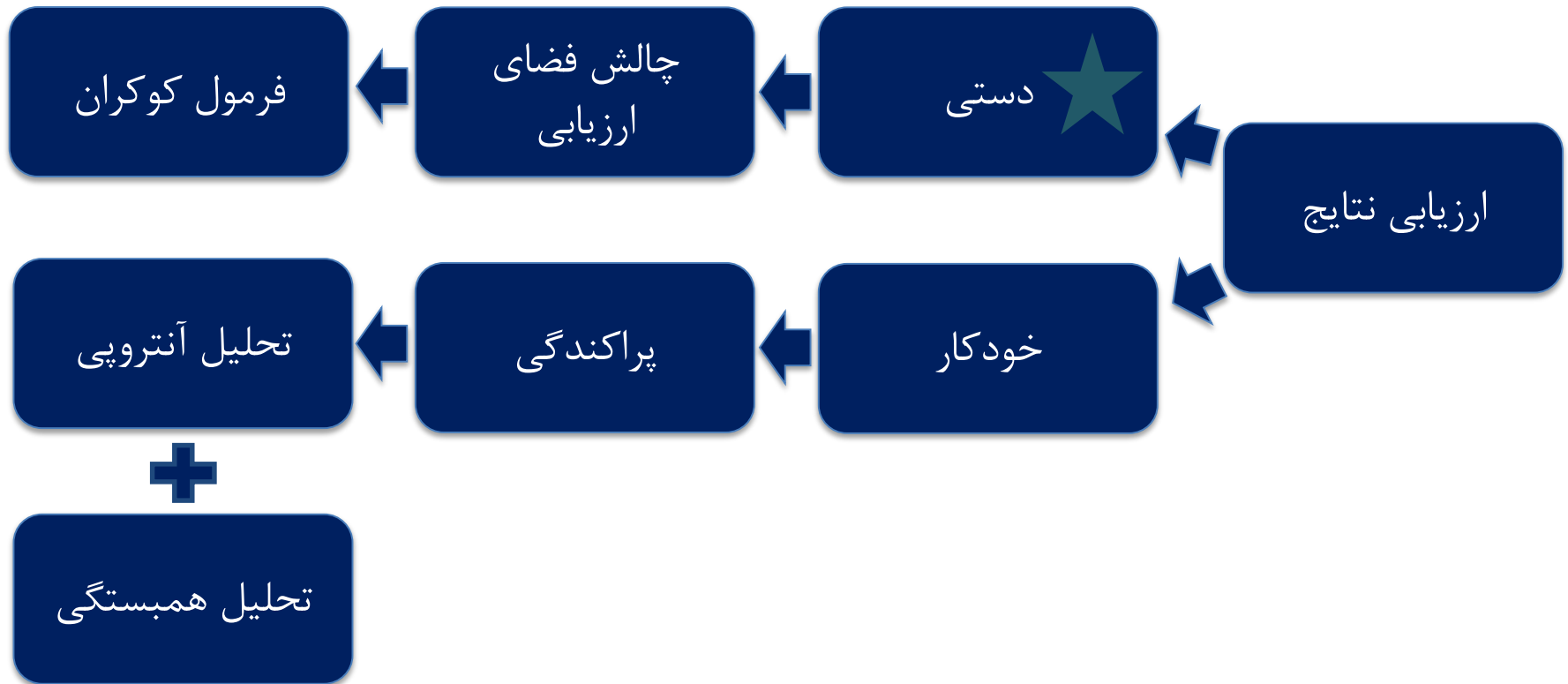
معماری کلی ترموگراف



ارزیابی نتایج

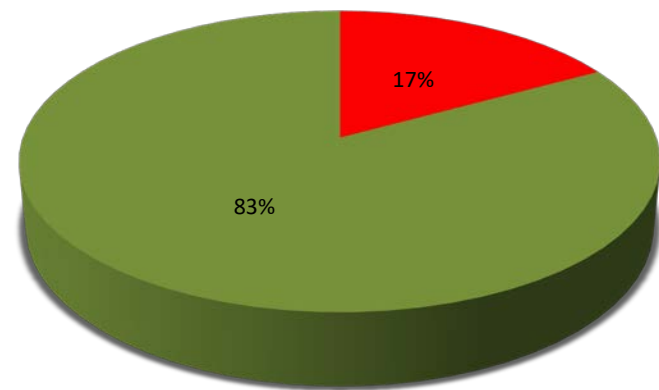
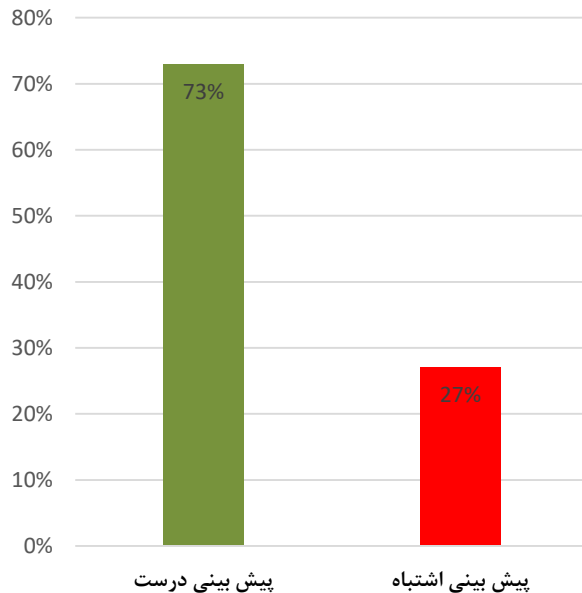


ارزیابی نتایج



ارزیابی دستی نتایج

- جامعه آماری نتایج نهایی تشخیص و پیش‌بینی به دست آمده بزرگ است و نمی‌توان همه‌ی پذیرنده‌ها را بررسی دستی کرد.
- به همین خاطر از فرمول کوکران برای بدست آوردن حجم نمونه استفاده می‌کنیم.



- صنف فعلی فروشگاه درست ثبت شده است
- صنف فعلی فروشگاه اشتباه ثبت شده است

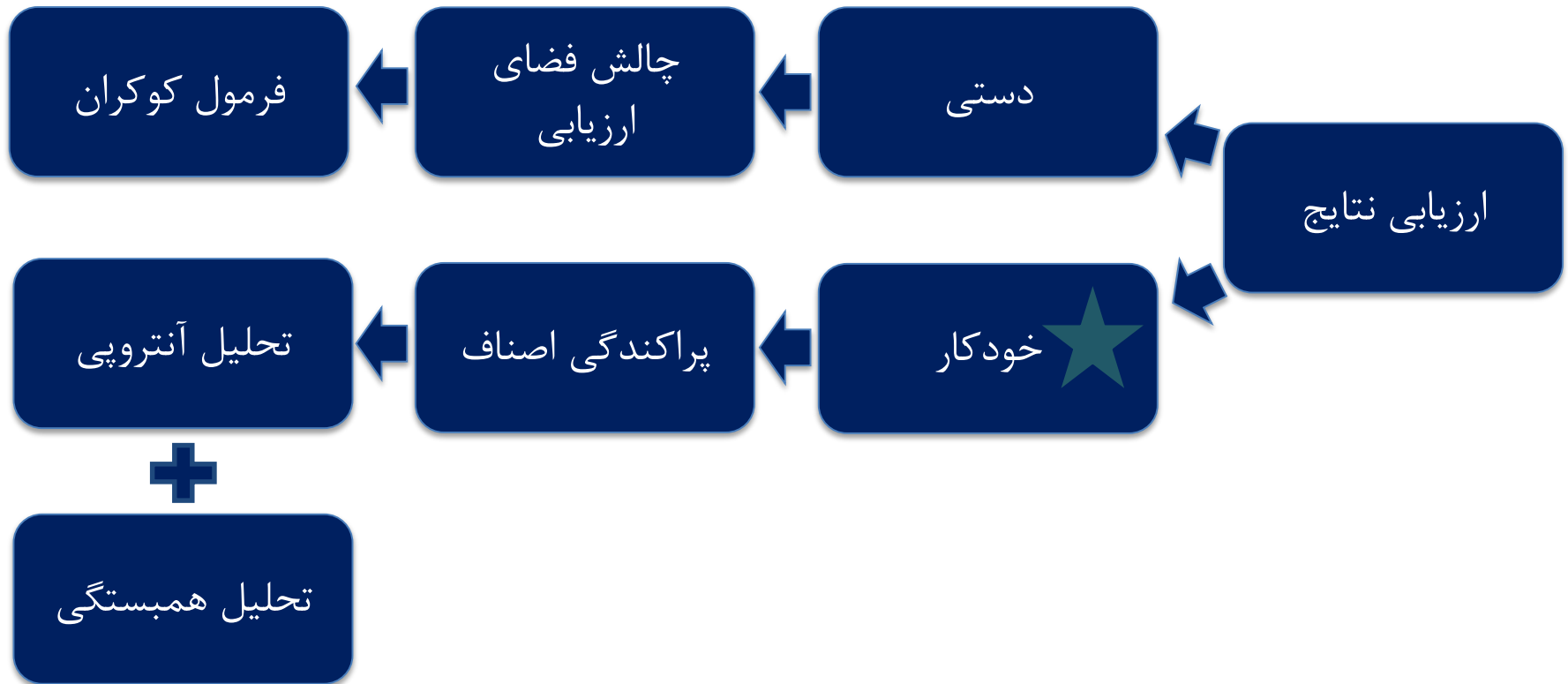


پیش‌بینی صنف جایگزین

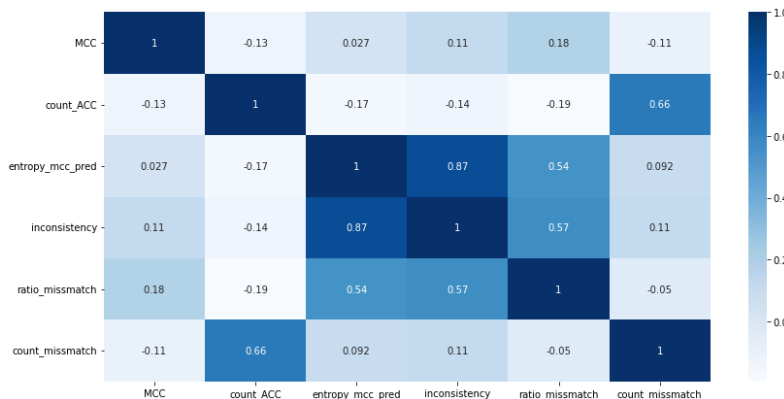
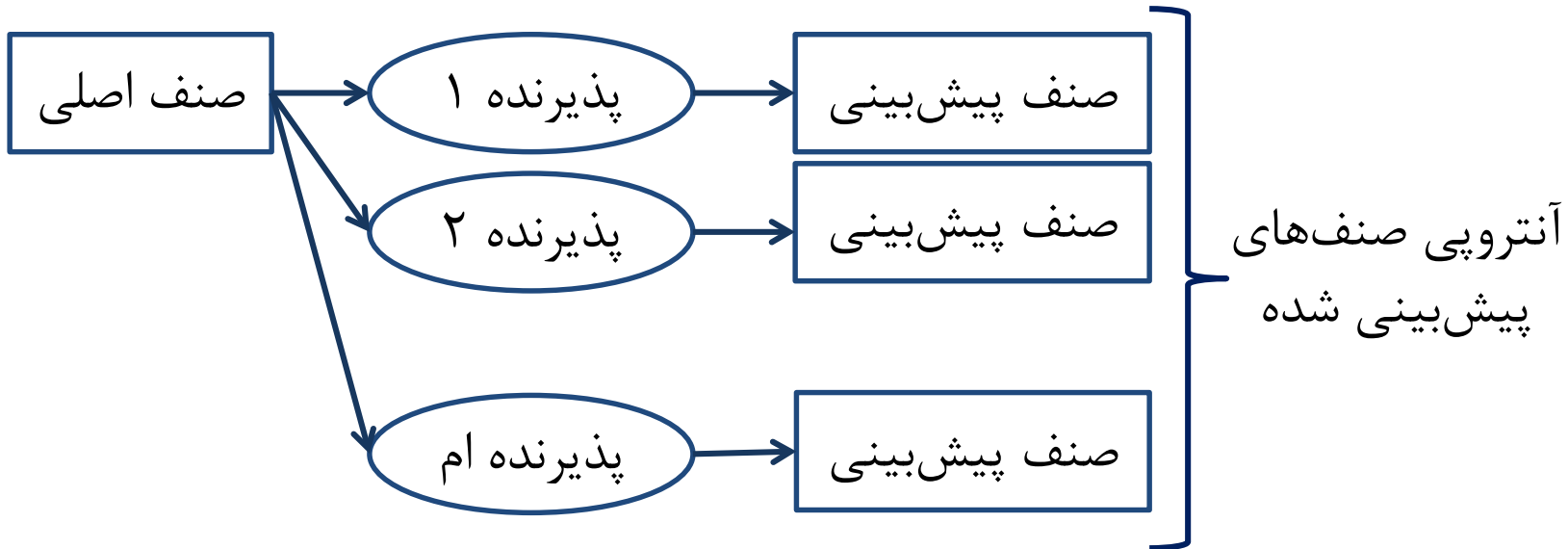
تشخیص ناهنجاری صنف



ارزیابی نتایج

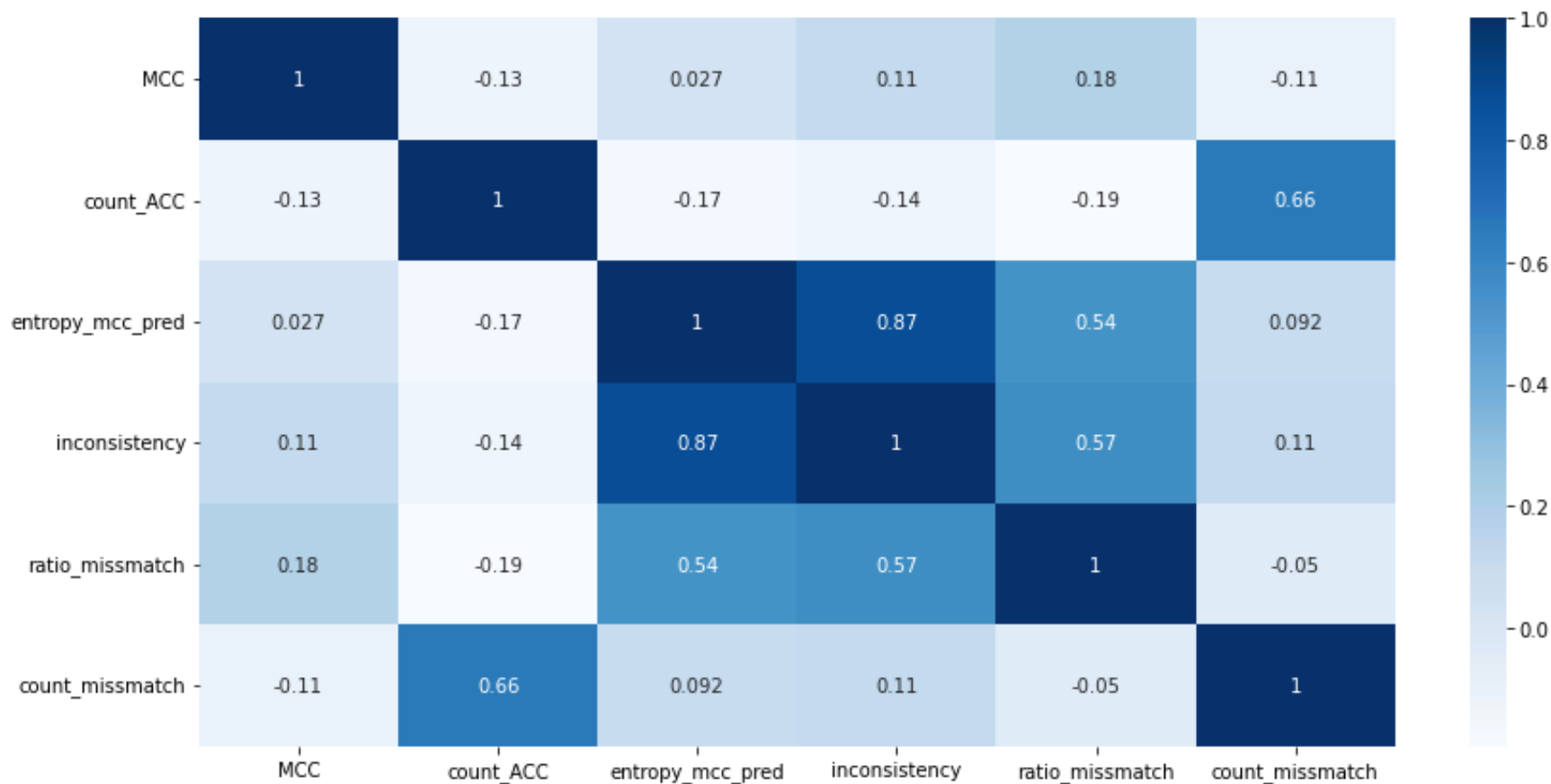


ارزیابی خودکار نتایج

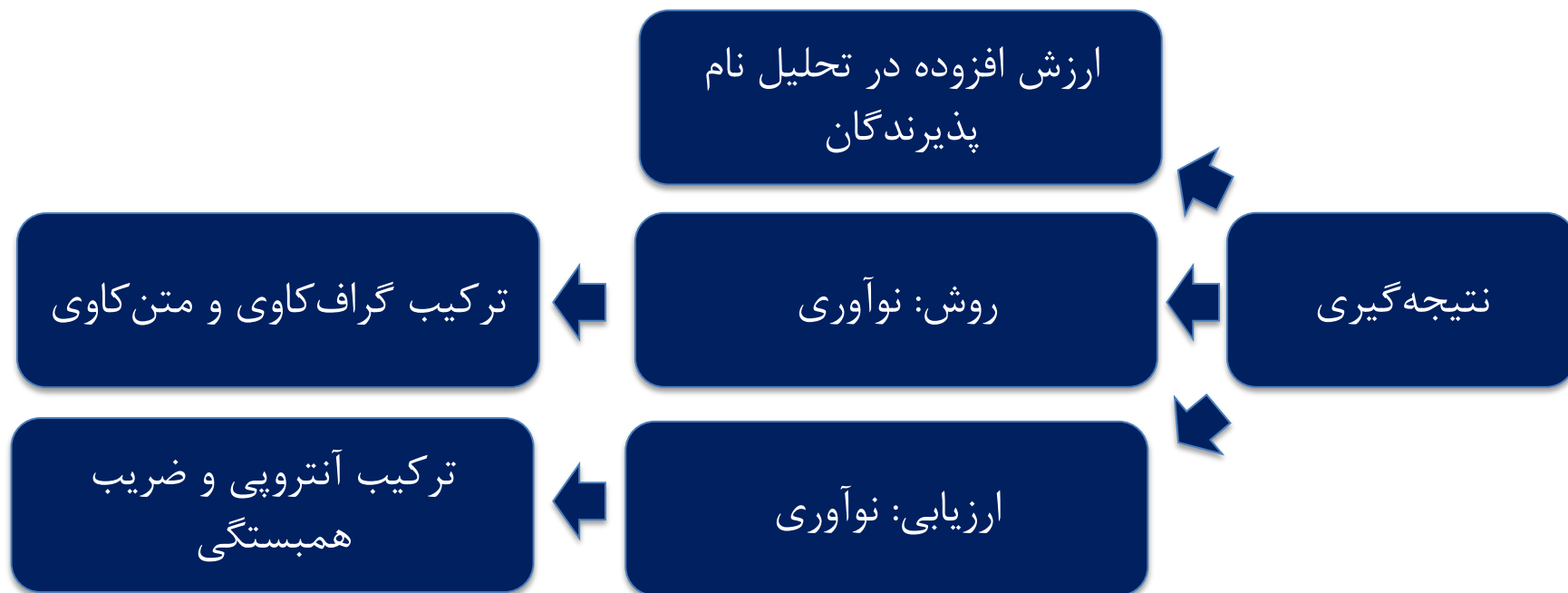


میزان همبستگی بین
آنتروپی صنف‌ها و عدم
تطابق پذیرندگان صنف‌ها

میزان همبستگی



نتیجه گیری



با تشکر از وقت و توجه شما

بهناز پورا براهیم