

در مدت فعالیت من ، بعنوان پژوهش گرا رشد و رئیس دپارتمان داده‌کاوی در شرکت بنز آلمان ، این دپارتمان در چند پروژه داده‌کاوی نیز شرکت داشت که توسط بازار مشترک اروپا پشتیبانی مالی می‌شد. CRISP-DM بنظر من مهم‌ترین آن‌ها بود. ایده این پروژه که در انجام آن علاوه بر دپارتمان داده‌کاوی شرکت بنز سه شرکت دیگر از جمله شرکت SPSS نیز سهیم بودند در سال ۱۹۹۶ میلادی شکل گرفت و انجام آن سه سال بطول انجامید. این شرکت که چند سال پیش توسط IBM خریداری شد هماهنگی پروژه را نیز بعهده داشت.

امروزه که بیش از چهارده سال از پایان این پروژه می‌گذرد، بدون اغراق می‌توان گفت که CRISP-DM مهم‌ترین استاندارد موجود جهت انجام پروژه‌های کاربردی داده‌کاوی در جهان می‌باشد. بسیاری از پروژه‌های کاربردی داده‌کاوی که در مدت حضور من در دپارتمان داده‌کاوی شرکت بنز انجام گرفت، با استفاده از این استاندارد بود. افزون بر این، تاجایی که من اطلاع دارم، در ایران نیز تاکنون این استاندارد در بسیاری از رساله‌های کارشناسی ارشد که در زمینه داده‌کاوی نوشته شده اند، مورد استفاده قرار گرفته است. همچنین مدت زیادی است که این استاندارد در ابزار داده‌کاوی SPSS Modeler نیز تلفیق گردیده و در دسترس کاربران این ابزار قرار دارد. برای من بعنوان شخصی که نقش کوچکی در بوجود آمدن این استاندارد داشته، بسیار مسرت بخش است که اکنون ترجمه توصیف این استاندارد به زبان فارسی، به همت گروه داده‌کاوی دایکه و بویژه با کوشش مستمر آقای محمد روزه، در اختیار پژوهش‌گران و کاربران فارسی زبان داده‌کاوی قرار می‌گیرد. ضمن ابراز قدردانی از زحمات ایشان برای تمام خوانندگان و بویژه برای دانشجویان عزیز آرزوی بهره‌گیری فراوان از این اثر را دارم.

غلامرضا نخعی زاده

استاد دانشگاه کارلسروهه آلمان

مدیر اسبق دپارتمان داده‌کاوی شرکت بنز

مقدمه مترجم

تصمیم به ترجمه و چاپ این استاندارد تقریباً همزمان با شروع فعالیت‌های ما در قالب دپارتمان داده‌کاوی در زمستان ۱۳۸۶ بوده است. کمتر از ۵ ماه بعد نسخه اولیه ترجمه با همکاری اعضای اولیه دپارتمان آماده شد. ولی نوپا بودن این دانش در سطح کشور و نداشتن تجربه اجرایی مناسب اعضای تیم در آن مقطع، باعث شد نسخه اولیه صرفاً ترجمه‌ای از متن اصلی بدون درک درستی از مفاهیم بکار رفته در آن باشد. بنابراین تصمیم به چاپ آن معلق شد تا در وقتی مناسب با کیفیتی شایسته محتوای پر اهمیت آن ارائه گردد.

طی سال‌های بعد متن اولیه بارها دستخوش تغییر شد و با توجه به تجربیات روزافزون که در قالب آموزش، مشاوره و اجرای پروژه‌های متعدد بدست آمده بود مورد بازنگری قرار می‌گرفت. ولی گذر زمان و مشغله‌های کاری باعث کم‌رنگ شدن تصمیم اولیه برای چاپ شده بود؛ تا اینکه برگزاری مستمر دوره‌های آموزشی «پیاده‌سازی فرآیند داده‌کاوی در ابزارهای تخصصی داده‌کاوی» از جمله (Clementine SPSS Modeler) در دو سال اخیر و مشاوره‌های متعددی که برای پیاده‌سازی و اجرای پروژه‌های داده‌کاوی به مدیران و کارشناسان شرکت‌ها، سازمان‌ها و همچنین دانشجویان مقاطع تکمیلی ارائه می‌شد، ضرورت چاپ این استاندارد را بیش از پیش نمایان ساخت.

با توجه به ماهیت استاندارد و همچنین طیف مخاطبان آن، سعی شده تا در کنار پایبندی به ساختار اصلی متن و همچنین نوع نگارش آن، متن ترجمه شده بصورت روان و شناور بوده تا کمترین میزان خستگی را در مطالعه آن ایجاد نماید.

برای رسیدن به این هدف، اساتید محترم و دوستان گرانقدری کمک نموده‌اند که صمیمانه از همه آن‌ها تشکر می‌کنم. جناب پروفسور نخعی زاده که راهنمایی‌ها و همراهی‌های ایشان مسلماً برای من فراموش نشدنیست. همکاران خوبم در هسته اولیه دپارتمان داده‌کاوی، خانم زهرا زجاجی و آقای عباسعلی ناطقی که در تهیه نسخه اولیه این ترجمه از اعضای اصلی بوده‌اند. همچنین دوستان با ارزش و عزیزم خانم زهرا ذوالقدر و آقایان احسان اژدری، خشایار مقدم، صادق طولابی فر، پویان رضانی و سایر دوستانی که در بازبینی متن، ویرایش و همچنین حمایت‌هایشان برای تسریع در اتمام این کار نقش بسزایی داشتند.

با تمام این صحبت‌ها معتقدم متن حاضر خالی از اشکال نیست و بهبود کیفیت آن مستلزم ارائه پیشنهادات و انتقادات خوانندگان عزیز می‌باشد.

محمد روزبه

مرداد ۱۳۹۲

فهرست مطالب

پیشگفتار

■ فصل اول: مقدمه

روش شناسی CRISP-DM

جداسازی سلسله مراتبی

مدل مرجع و راهنمای کاربر

نگاشت مدل های عمومی به مدل های اختصاصی

چهارچوب داده کاوی

نگاشت با چهارچوب ها

چگونگی نگاشت

توضیح فصل ها

مطالب

اهداف

■ فصل دوم: مدل مرجع CRISP-DM

۱ شناخت و درک کسب و کار

۱-۱ تعیین اهداف فعالیت تجاری

۲-۱ ارزیابی وضعیت

۳-۱ تعیین اهداف داده کاوی

۴-۱ ارائه طرح پروژه

۲ شناسایی و درک داده ها

۱-۲ جمع آوری اولیه داده ها

۲-۲ توصیف داده ها

۳-۲ کاوش داده ها

۴-۲ بررسی کیفیت داده ها

۳ آماده سازی داده ها

۱-۳ انتخاب داده ها

۲-۳ پاکسازی داده ها

۳-۳ ساخت داده ها

- ۴-۳ یکپارچه سازی داده ها
- ۴ مدل سازی
- ۱-۴ انتخاب تکنیک مدل سازی
- ۲-۴ ایجاد طرح آزمون
- ۳-۴ ساخت مدل
- ۴-۴ ارزیابی مدل
- ۵ ارزیابی
- ۱-۵ ارزیابی نتایج
- ۲-۵ مرور فرآیند
- ۳-۵ تعیین مراحل بعدی
- ۶ گسترش و استقرار
- ۱-۶ طراحی برای گسترش و استقرار
- ۲-۶ طراحی برای نظارت و نگهداری
- ۳-۶ تهیه گزارش نهایی
- ۴-۶ بازبینی پروژه

■ فصل سوم: راهنمای کاربر CRISP-DM

۱ شناخت و درک کسب و کار

۱-۱ تعیین اهداف فعالیت تجاری

۲-۱ ارزیابی وضعیت

۳-۱ تعیین اهداف داده کاوی

۴-۱ ارائه طرح پروژه

۲ شناسایی و درک داده ها

۱-۲ جمع آوری اولیه داده

۲-۲ توصیف داده ها

۳-۲ کاوش داده ها

۴-۲ بررسی کیفیت داده ها

۳ آماده سازی داده ها

۱-۳ انتخاب داده ها

۲-۳ پاکسازی داده ها

۳-۳ ساخت داده ها

۴-۳ یکپارچه سازی داده ها

- ۳-۵ فرمت داده‌ها
- ۴ مدلسازی
- ۴-۱ انتخاب تکنیک مدلسازی
- ۲-۴ ایجاد طرح آزمون
- ۴-۳ ساخت مدل
- ۴-۴ ارزیابی مدل
- ۵ ارزیابی
- ۵-۱ ارزیابی نتایج
- ۵-۲ مرور فرآیند
- ۵-۳ تعیین مراحل بعدی
- ۶ گسترش و استقرار
- ۶-۱ طراحی برای گسترش و استقرار
- ۶-۲ طراحی برای نظارت و نگهداری
- ۶-۳ تهیه گزارش نهایی
- ۶-۴ بازبینی پروژه

■ فصل چهارم: خروجی‌های CRISP-DM

- ۱ شناخت و درک کسب و کار
- ۲ شناسایی و درک داده‌ها
- ۳ آماده سازی داده‌ها
- ۴ مدلسازی
- ۶ گسترش و استقرار
- ۷ خلاصه وابستگی‌ها

■ فصل پنجم: پیوست

- ۱ واژه نامه و شرح اصطلاحات تخصصی
- ۲ انواع مسئله داده کاوی
- ۲-۱ توصیف و خلاصه سازی داده‌ها
- ۲-۲ خوشه بندی
- ۲-۳ توصیف مفاهیم
- ۲-۴ طبقه بندی
- ۲-۵ پیش بینی
- ۲-۶ تحلیل وابستگی

پیشگفتار

CRISP-DM محصول تلاش های سه شرکت پر سابقه بازار نوپا و ناپخته داده کاوی، در اواخر سال ۱۹۹۶ است. دایملر کرایسلر (که بعداً دایملر بنز شد) پیش از دیگران در فرآیند های صنعتی و امور تجاری اش از داده کاوی استفاده کرده بود. SPSS (یا ISL) از سال ۱۹۹۰ خدماتی بر اساس داده کاوی ارائه داده بود و در سال ۱۹۹۴، نخستین نرم افزار تجاری داده کاوی (کلمنتاین) را معرفی کرد. NCR نیز در راستای ارائه ارزش افزوده بیشتر برای مشتریان انبار داده ترادیتای خود، تیمی را متشکل از مشاوران داده کاوی و متخصصان صنعت، برای رفع نیازمندی های آنان، به کار گرفته بود.

گرایش های بازار آن زمان حاکی از گسترش روزافزون داده کاوی بود؛ گسترشی که هیجان برانگیز بود و البته نگران کننده. همه ما با همراهی هم دستاوردهای داده کاوی را توسعه داده بودیم اما آیا ما درست عمل کرده بودیم؟ آیا همه کسانی که برای اولین بار از داده کاوی استفاده می کردند می بایستی از طریق سعی و خطا (مانند ما!) آن را بیاموزند؟ و ما که ارائه دهنده آن بودیم چگونه باید مشتریان خود را به این باور می رساندیم که داده کاوی می تواند راهنمای امور تجاری آن ها باشد؟

سپس به این نتیجه رسیدیم یک مدل فرآیند استاندارد که به جایی اختصاص نداشته باشد و به صورت مجانی در دسترس باشد می تواند این موضوع را به ما و همه کاربران نشان دهد.

یک سال بعد کنسرسیومی تجاری تشکیل دادیم (با عنوان فرآیند استاندارد فرا صنعتی داده کاوی) که با تأمین بودجه خود از جانب کمیسیون اروپا همان اهداف اولیه ما را دنبال می کرد. از آنجا که نایبست CRISP-DM از نظر صنعت و ابزار و کاربرد متعلق به گروه خاصی می بود، دریافتیم که برای ادامه کار می بایست از گروهی به اندازه کاربران و سایر سرویس دهندگان (مانند ارائه دهندگان انبار داده ها و مدیران مشاوره که علاقمند به داده کاوی بودند) انرژی بگیریم. این کار را با ایجاد گروه علاقمندان ویژه CRISP-DM (SIG) انجام دادیم. به این صورت که از طریق تبلیغات رادیویی از علاقمندان دعوت کردیم در کارگاه یک روزه ای که در آمستردام برگزار شد شرکت کنند. ما می خواستیم ایده های خود را با آن ها در میان بگذاریم و از آن ها بخواهیم ضمن معرفی خود به بحث و تبادل نظر در مورد پیشبرد CRISP-DM بپردازند.

در روز برگزاری کارگاه اعضای گروه نگران بودند نکند کسی در کارگاه شرکت نکند؟ و اگر کسی شرکت کرد برای این آمده باشد که بگوید نیازی ضروری به یک فرآیند استاندارد احساس نمی‌کند. یا این که تلاش‌های ما آنقدر پراکنده بوده که هرگونه ایده استاندارد کردن مشتتی تخیلات غیر عملی باشد. نتیجه کارگاه فراتر از انتظار ما بود. آنچه بیش از همه جلب نظر می‌کرد سه مورد زیر بود:

- جمعیت حاضرین دو برابر انتظار اولیه ما بود.
- همگی متفق‌القول بر این باور بودند که صنعت امروز نیاز مبرم به فرآیند استاندارد دارد که بایستی هم-اکنون انجام پذیرد.
- هر بار که یکی از حضار دیدگاه خود را در مورد داده کاوی با توجه به تجربیات پروژه اش ارائه می‌کرد، علی‌رغم تفاوت‌های ظاهری به ویژه در تأیین حد و مرزها و مجموعه لغات مورد استفاده، وجوه قابل توجه مشترک در دیدگاه‌های آن‌ها وضوح بیشتری پیدا می‌کرد.
- با خاتمه کارگاه اطمینان کامل داشتیم که با تکیه بر امیدبخشی و نقدهای SIG خواهیم توانست یک مدل فرآیند استاندارد برای جامعه داده کاوی ارائه دهیم.

در طی دو سال و نیم پس از آن تلاش ما معطوف به توسعه و رفع اشکالات CRISP-DM بود. تلاش‌هایمان را روی پروژه‌های واقعی و عظیمی مانند داده کاوی برای مرسدس بنز و شریک بخش بیمه مان OHRA ادامه داده و طی آن روی یکپارچه سازی CRISP-DM با ابزارهای تجاری کار کردیم. SIG بسیار ارزشمند از آب درآمد بود و اعضای آن به بیش از ۲۰۰ نفر رسیده بود و کارگاه‌هایی در لندن، نیویورک و بروکسل برگزار کرده بود.

پس از پایان یافتن بخش تخصیص اعتبار پروژه در اواسط سال ۱۹۹۹ یک پیش نویس با کیفیت از مدل فرآیند تهیه کردیم و برای بهبود هر چه بیشتر آن تلاش کردیم که نتیجه آن تلاش‌ها CRISP-DM 1.0 بود. در حین انجام کار می‌دانستیم که فرآیند مدل CRISP-DM هنوز در حال طی کردن مراحل ساخت‌امی باشد و امتحان خود را تنها روی پروژه‌های انگشت شمار پس داده است. در سال گذشته دایملر کرایسلر فرصت پیدا کرد تا CRISP-DM را به طور گسترده به کار گیرد. همچنین گروه‌های خدماتی حرفه‌ای مانند SPSS و CRISP-DM، NCR را پذیرفته‌اند، و از آن برای حل مسائل تجاری و صنعتی مشتریان‌شان استفاده می‌کنند. در این مدت سرویس‌های خدماتی زیادی که CRISP-DM را پذیرفته‌اند با تحلیلگران مختلف مکرراً به آن به عنوان یک استاندارد غیر رسمی صنعت مراجعه کرده‌اند. و در کنار همه اینها افزایش آگاهی مشتریان از اهمیت CRISP-DM نکته ایست که به آسانی نمی‌توان از کنار آن گذشت (CRISP-DM هم اکنون مرتباً به مناقصه گذاشته می‌شود).

بر این باوریم که ابتکار عمل ما کاملاً به اثبات رسیده است و اگرچه گسترش‌ها و بهسازی‌هایی که بعدها صورت می‌گیرد مطلوب و البته اجتناب ناپذیر است اما نسخه CRISP-DM-1.0 آنقدر ارزشمند

هست که نشر و گسترش داده شود.

CRISP-DM نه از طریق نظریه پردازی و کارآکادمیک بر روی اصول فنی شکل گرفته است و نه به وسیله گروه نخبه ای که پشت درهای بسته در مورد آن تحقیق کرده باشند. این دو مسیر به توسعه روش هایی منجر شده است که پیش از این امتحان شده اند اما به ندرت به یک استاندارد مورد اقبال عموم که عملی و موفق باشد منجر شده اند. دلیل موفقیت CRISP-DM استوار بودن بر مبنای تجربه های عملی و واقعی مرسوم است که پروژه های داده کاوی با استفاده از آن ها انجام شده اند. از این رو ما بی نهایت مرهون کاربرانی هستیم که در طول پروژه ایده ها و تلاش های خود را در اختیار ما قرار داده اند.

کنسرسیوم CRISP-DM

آگوست ۲۰۰۰



فصل اول مقدمه

روش شناسی CRISP-DM

جداسازی سلسله مراتبی

روش داده کاوی CRISP-DM در شمای کلی با مدلی که شامل چهار سطح از وظایف است نشان داده می شود. این چهار سطح عبارتند از: فاز^۳، وظیفه عمومی^۴، وظیفه تخصصی^۵ و نمونه فرآیند^۶. در واقع این مدل نشان دهنده فرآیندی است که داده کاوی از طریق آن صورت می گیرد. در سطح نخست، فرآیند شامل چند فاز است؛ و هر فاز دربرگیرنده چندین وظیفه عمومی (سطح دوم)، سطح دوم به این دلیل عمومی نامیده می شود که باید شامل همه حالات ممکن داده کاوی باشد. دو ویژگی مهم هر سطح عمومی «کامل بودن» و «پایداری» است. «کامل بودن» موجب می شود این وظیفه دربردارنده همه قسمت های فرآیند داده کاوی و همه کاربردهای ممکن آن باشد؛ و «پایداری» قابلیت توسعه یافتن آن را در مواجهه با موارد پیش بینی نشده (مانند تکنیک های نوین مدلسازی) تضمین می کند.

سطح سوم شامل وظایف تخصصی ست. این وظایف نحوه انجام وظایف عمومی را در شرایط مختلف شرح می دهند. به عنوان مثال اگر وظیفه عمومی «پاکسازی» باشد، وظیفه تخصصی (که توضیح دهنده نحوه اجرای وظیفه عمومی در موقعیت های مختلف است) اجرای پاکسازی در مورد مقادیر عددی یا مقادیر گسسته را توضیح خواهد داد.

اجرا شدن وظایف و فازها به صورت مراحل مجزا و با ترتیب مشخص حالتی ایده آل است. در عمل ترتیب اجرای بسیاری از وظایف با حالت ایده آل تفاوت دارد و در موارد زیادی لازم است مجدداً به مراحل پیشین برگشته و عملیات خاصی را تکرار کنیم. با این توضیحات باید این نکته را هم خاطر نشان کنیم که مدل حاضر در صدد پوشش دادن همه مسیرهای موجود در فرآیند داده کاوی نیست، چون این امر مستلزم توصیف مدلی پیچیده خواهد بود.

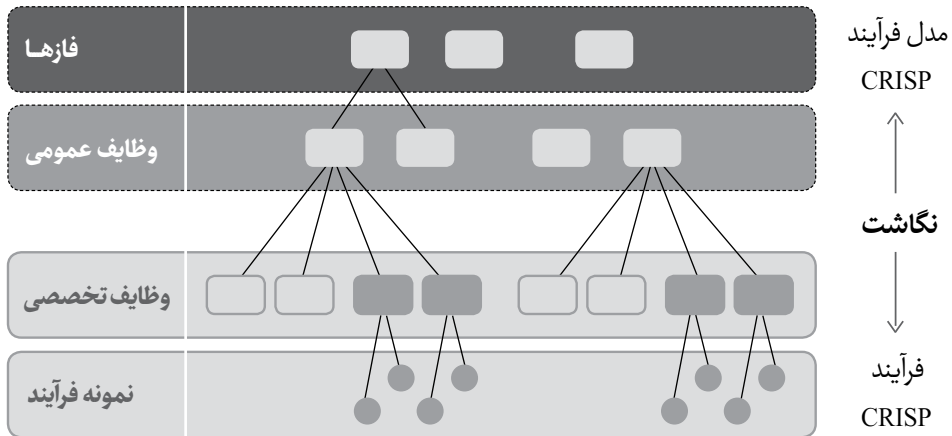
سطح چهارم نمونه فرآیند نام دارد و شامل ثبت دقیق عملیات، تصمیم ها و نتایج یک مسئله واقعی داده کاوی است. نمونه فرآیندها بر اساس وظایفی که در سطوح بالاتر تعریف شده اند سازمان دهی می شوند و شامل جزئیات عملیات یک مسئله واقعی هستند.

Phase ۳

Generic Task ۴

Specialized Task ۵

Process Instance ۶



شکل ۱: تقسیم بندی چهارسطحی روش شناسی CRISP-DM

مدل مرجع^۷ و راهنمای کاربر^۸

در روش شناسی^۹ CRISP-DM مدل مرجع و راهنمای کاربر با هم تمایز دارند. مدل مرجع شمایی کلی از فازها، وظایف و خروجی آن هاست و عملیاتی را که در یک پروژه داده کاوی صورت می گیرد توضیح می دهد. اما راهنمای کاربر نکات و راهنمایی های جزئی تری در مورد هر فاز و وظایف مربوط به آن در اختیار می گذارد و روش انجام پروژه داده کاوی را شرح می دهد. این کتاب مدل مرجع و راهنمای کاربر را به صورتی کلی پوشش می دهد.

نگاشت مدل های عمومی به مدل های اختصاصی

چهارچوب داده کاوی^{۱۰}

چهارچوب داده کاوی، نگاشت بین سطوح عمومی و اختصاصی را به عهده دارد. چهار بعد چهارچوب داده کاوی عبارتند از:

- **دامنه کاربردی**، که حوزه خاصی است که داده کاوی در آن رخ می دهد.
- **نوع مسئله داده کاوی**، که بیان کننده طبقه بندی خاصی از نوع موضوعات پروژه داده کاوی است.

Reference Model ۷

User Guide ۸

Methodology ۹

Data Mining Context ۱۰

- جنبه های تکنیکی، که در بر گیرنده نکات خاصی است که در یک پروژه داده کاوی صورت می گیرند.
 - ابزار و تکنیک ها، که مشخص کننده ابزارها و تکنیک های مورد استفاده در پروژه داده کاوی است.
- در جدول ۱ چهار بعد داده کاوی همراه با مثال شرح داده شده اند.
- هر چهارچوب خاص با نسبت دادن مقادیر واقعی به یک یا چند بعد از ابعاد مذکور به وجود می آید. برای مثال یک پروژه داده کاوی که مسئله آن از نوع پیش بینی است، یک چهارچوب داده کاوی خاص است. هر چه مقادیر مشخص شده بیشتری به ابعاد داده کاوی نسبت داده شود چهارچوب داده کاوی روشن تر و محکم تر خواهد بود.

چهارچوب داده کاوی				
ابعاد	دامنه کاربرد	نوع مسئله داده کاوی	جنبه های تکنیکی	ابزارها و تکنیک ها
مثال ها	مدل پاسخ	توصیف و خلاصه سازی	مقادیر گمشده	کلمنتاین
	پیش بینی رویگردانی مشتریان	خوشه بندی	مقادیر پرت	MineSet
	...	توصیف مفاهیم		درخت های تصمیم
		طبقه بندی		
		پیش بینی		
		تحلیل وابستگی ها		

جدول ۱: ابعاد چهارچوب داده کاوی همراه با مثال

نگاشت^{۱۱} با چهارچوب ها

در CRISP-DM نگاشت های عمومی و تخصصی از هم متمایزند:

نگاشت برای وضعیت کنونی:

اگر مدل کلی فرآیند فقط برای یک پروژه داده کاوی خاص در نظر گرفته شده باشد و هدف ما تطبیق وظایف عمومی و توصیفاتشان با خواسته ها و نیازهای آن پروژه خاص باشد، آنگاه پای یک انطباق یگانه در میان است که (به احتمال زیاد) فقط یک بار مورد استفاده قرار می گیرد.

نگاشت برای آینده:

هر گاه یک مدل را به صورت سیستماتیک و براساس یک چهارچوب از پیش تعریف شده تخصصی کنیم (به عبارت دیگر تجربیات یک پروژه را در راستای اختصاصی سازی مدل، تحلیل و جمع بندی کنیم) تا بعدها از آن در چهارچوب های دیگر استفاده کنیم، یک مدل ویژه CRISP-DM ایجاد کرده ایم. این که کدام انطباق مفید تر خواهد بود، به چهارچوب های خاص داده کاوی و نیازهای شرکت بستگی دارد.

چگونگی نگاشت

استراتژی پایه برای نگاشت مدل کلی به تخصصی برای هر دو نگاشت یکسان است:

- چهارچوب خاص خود را تحلیل کنید.
- جزئیات غیر قابل اعمال در چهارچوب مورد نظر را حذف کنید.
- جزئیات مورد نظر خود را اضافه کنید.
- بخش های کلی را براساس چهارچوب مورد نظر تنظیم کنید.
- در صورت امکان نام قسمت های کلی را با نام های تازه ای که القا کننده معانی گویاتری در چهارچوب اند، عوض کنید.

توضیح فصل ها

مطالب

مدل فرآیند CRISP-DM در پنج فصل مختلف تنظیم شده است.

- فصل اول مقدمه ای بر روش CRISP-DM است که راهنمایی های کلی برای نگاشت مدل عمومی فرآیند به مدل تخصصی فرآیند را فراهم آورده است.
- فصل دوم مدل مرجع CRISP-DM و اجزای آن از قبیل فازها، وظایف عمومی و خروجی ها را شرح می دهد.
- فصل سوم راهنمای کاربر CRISP-DM را توضیح می دهد و شامل توصیف کامل مباحث فازها، وظایف عمومی و خروجی هاست. در این فصل همچنین انجام یک پروژه داده کاوی به صورت دقیق شرح داده شده است.
- فصل چهارم درباره گزارش هایی ست که در حین انجام پروژه و پس از آن بایستی نوشته شوند و طرح کلی این گزارش ها را ارائه می دهد. همچنین جدول متقاطع خروجی ها و وظایف نمایش داده شده است.
- فصل پنجم پیوست است که شامل توصیف اصطلاحات و واژگان است و خصوصیات انواع مسائل داده کاوی را شرح می دهد.

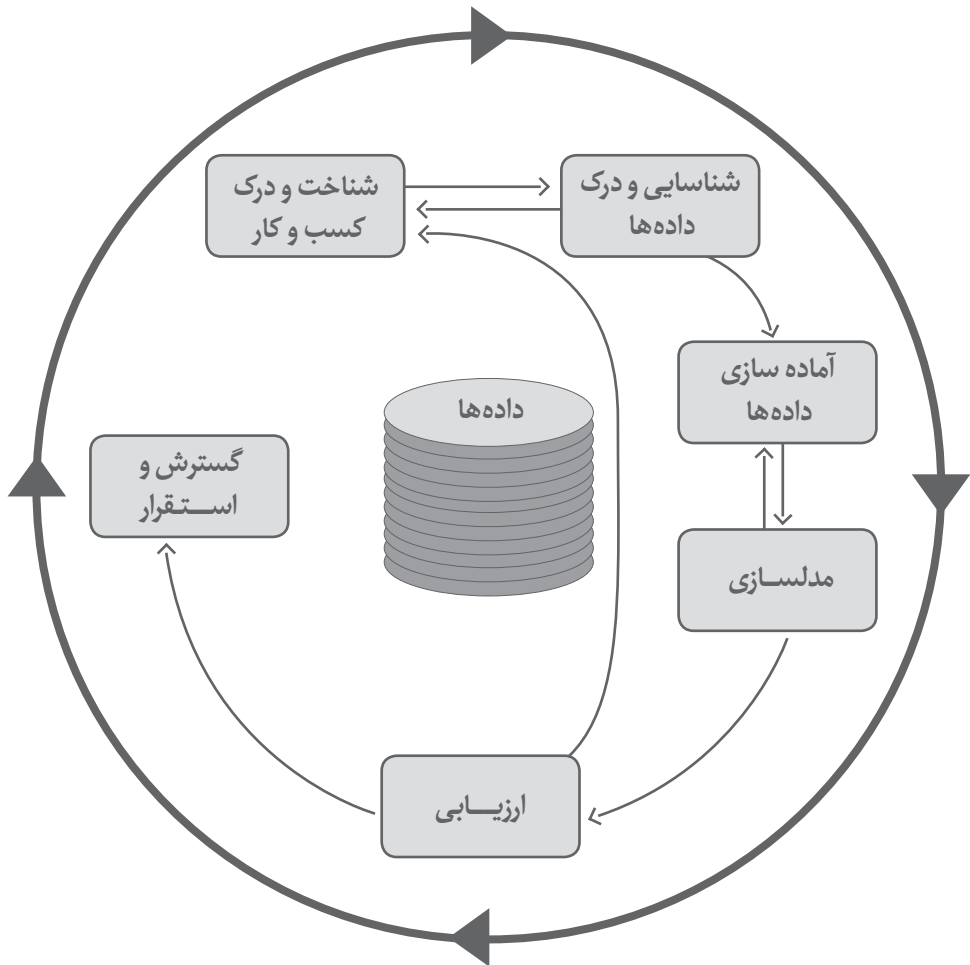
اهداف

- نکات زیر امکان استفاده هر چه بهتر از این کتاب را فراهم می کند:
- اگر نخستین باری است که مدل CRISP-DM را مطالعه می کنید، مطالعه کتاب را از فصل اول (مقدمه) شروع کنید تا درک مناسب تری از مفاهیم این روش و ارتباط بین این مفاهیم به دست آورید. در این صورت در پروژه های بعد از مطالعه مجدد آن معاف خواهید بود!
 - اگر دستیابی هر چه سریع تر به یک شمای کلی از مدل CRISP-DM مد نظر شماست، به فصل دوم (مدل مرجع CRISP-DM) مراجعه کنید. همچنین اگر قصد دارید به سرعت یک پروژه داده کاوی را شروع کنید برای آشنایی مقدماتی با راهنمای کاربر مطالعه این فصل را از دست ندهید.
 - اگر برای انجام پروژه داده کاوی به راهنمایی های دقیق و کامل نیازمندید به فصل سوم (راهنمای کاربر CRISP-DM) مراجعه کنید. این فصل مهم ترین و با ارزش ترین فصل این کتاب است (پیش از مطالعه این فصل بایستی دو فصل نخست را مطالعه کرده باشید).
 - اگر در حال نوشتن گزارش های داده کاوی هستید به فصل چهارم مراجعه کنید. در صورت تمایل به تهیه گزارش ها در حین انجام پروژه، از مطالب دو فصل ۳ و ۴ به صورت همزمان استفاده کنید.
 - پیوست پایان کتاب می تواند برای به دست آوردن اطلاعاتی اولیه در مورد داده کاوی مفید واقع شود. اگر هنوز تخصص چندانی در زمینه داده کاوی ندارید برای پیدا کردن معانی و مفاهیم واژه های تخصصی از این فصل استفاده کنید.



فصل دوم
مدل مرجع
CRISP-DM

مدل مرجع CRISP-DM (شکل ۲) شمایی اجمالی از مراحل مختلف پروژه‌های داده کاوی است. این مدل شامل فازهای یک پروژه، وظایف مربوط به آن‌ها و روابط بین این وظایف است. بنابراین تعیین دقیق همه روابط بین وظایف غیر ممکن می‌باشد. اساساً وجود ارتباط بین هر کدام از وظایف داده کاوی، به هدف، زمینه، علایق کاربر و اهمیت داده‌ها بستگی دارد.



شکل ۲: فازهای مدل مرجع CRISP-DM

مراحل مختلف یک پروژه داده کاوی شامل ۶ فاز است. این فازها در شکل ۲ نمایش داده شده اند. البته توالی اجرای فازها از پیش تعیین شده نیست و معمولاً لازم است به برخی فازها مجدداً مراجعه شود. لزوم انجام یا عدم انجام این عمل بستگی به نتیجه هر فاز دارد و اینکه کدام فاز یا کدام وظیفه خاص از یک فاز، تعیین کننده گام بعدی خواهد بود. پیکان های نمایش داده شده در شکل ۲ مهم ترین و پر کاربرد ترین ارتباطات بین فازها را نشان می دهد.

دایره بیرونی در شکل ۲ نشان دهنده طبیعت چرخشی داده کاوی است. در واقع داده کاوی چیزی بیش از یک راهکار تعمیر یافته برای حل مسئله نیست. مطالبی که در حین فرآیند داده کاوی و از طریق راهکارهای تعمیر یافته آموخته می شوند بیشتر تمرکز خود را به سؤالات کسب و کار معطوف می کنند. بر همین اساس است که بیشترین بهره برداری در فرآیندهای داده کاوی از تجربیات پیشین صورت می گیرد.

در بخش های زیر هر کدام از فازها را به طور خلاصه مرور می کنیم.

شناخت و درک کسب و کار^{۱۲}

نخستین فاز بر به دست آوردن درک مناسبی از اهداف و ملزومات پروژه از منظر کسب و کار تمرکز می کند و از آن برای تعریف یک مسئله داده کاوی استفاده می شود تا طرحی ابتدایی برای رسیدن به اهداف مسئله فراهم گردد.

شناسایی و درک داده ها^{۱۳}

فاز شناسایی و درک داده ها با جمع آوری مقدماتی داده ها آغاز شده و با فعالیت هایی به منظور آشنایی با داده ها برای تعیین کیفیت آن ها ادامه می یابد تا به بینش اولیه ای از داده ها یا فرضیه هایی در مورد زیرمجموعه های مستتر در داده ها بینجامد.

آماده سازی داده ها^{۱۴}

فاز آماده سازی داده ها شامل همه فعالیت هایی است که مجموعه داده های نهایی (داده هایی که در اختیار ابزار مدل سازی قرار می گیرند) را از داده های اولیه ایجاد می کند. معمولاً فاز آماده سازی داده ها به زمان بیشتری نسبت به سایر فازها نیاز دارد و البته هیچ راهکار از پیش تعیین شده ای برای وظایف این فاز وجود ندارد. وظایف این فاز عبارتند از: جدول بندی، انتخاب و گزینش رکورد و صفت، تبدیل و پاکسازی داده ها برای مورد استفاده قرار گرفتن به وسیله ابزارهای مدل سازی.

Business Understanding ۱۲

Data Understanding ۱۳

Data Preparation ۱۴

مدلسازی^{۱۵}

در این فاز تکنیک های مختلف داده کاوی انتخاب و استفاده می شوند و پارامترهای مدل برای رسیدن به مقادیر بهینه تنظیم می شوند. به عنوان مثال برای یک مسئله داده کاوی چندین نوع تکنیک وجود دارد. برخی تکنیک ها به اعمال تغییرات خاصی در مورد فرم داده ها نیازمندند و این موضوع ایجاب می کند مجدداً به فاز آماده سازی داده ها رجوع شود.

ارزیابی^{۱۶}

در این مرحله مدل ساخته شده از نظر تحلیل داده ها به کیفیت مطلوبی رسیده است. قبل از ادامه مسیر و گسترش و استقرار نهایی مدل، ارزیابی کاملی از مدل و بررسی مراحل ساخت آن ضروری است تا در سایه آن، اطمینان از دستیابی مناسب به اهداف فعالیت تجاری حاصل شود. یک نکته کلیدی در این مرحله شناسایی فعالیت های تجاری مهمی است که به صورت مناسبی بررسی نشده اند. در پایان این فاز باید شرایط مناسب برای تصمیم گیری در مورد قابل استفاده بودن نتایج داده کاوی فراهم شده باشد.

گسترش و استقرار^{۱۷}

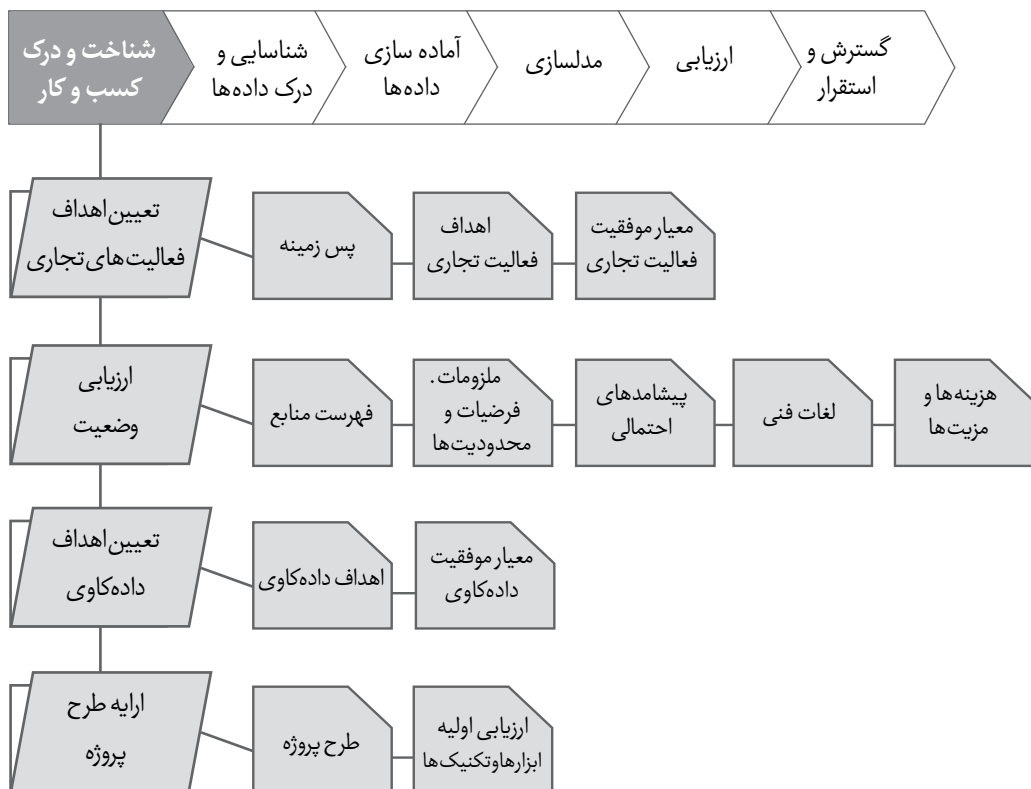
ساختن مدل پایان پروژه نیست. حتی اگر کسب اطلاعات بیشتر از داده ها هدف مدلسازی باشد، اطلاعات به دست آمده به سازماندهی و ارائه شدن نیازمندند به نحوی که مشتری بتواند از آن ها استفاده کند.

شکل ۳ شمایی کلی از فازها با وظایف عمومی (تیره) و خروجی ها را (رنگ خاکستری) نشان می دهد. در بخش های بعد، هر یک از وظایف و خروجی ها با جزئیات بیشتری توضیح داده شده اند.

<p>تعیین اهداف فعالیت تجاری (پس زمینه، اهداف فعالیت تجاری، معیار موفقیت فعالیت تجاری)</p> <p>ارزیابی وضعیت (فهرست منابع، ملزومات، فرضیات و محدودیت‌ها، پیشامدهای احتمالی، لغات فنی، هزینه‌ها و مزیت‌ها)</p> <p>تعیین اهداف داده کاوی (اهداف داده کاوی، معیار موفقیت داده کاوی)</p> <p>ارائه طرح پروژه (طرح پروژه، ارزیابی اولیه ابزارها و تکنیک‌ها)</p>	<p>شناخت و درک کسب و کار</p>
<p>جمع آوری اولیه داده‌ها (گزارش جمع آوری اولیه داده‌ها)</p> <p>توصیف داده‌ها (گزارش توصیف داده‌ها)</p> <p>کاوش داده‌ها (گزارش کاوش داده‌ها)</p> <p>بررسی کیفیت داده‌ها (گزارش بررسی کیفیت داده‌ها)</p>	<p>شناسایی و درک داده‌ها</p>
<p>انتخاب داده‌ها (دلیل انتخاب داده‌ها/عدم انتخاب داده‌ها)</p> <p>پاکسازی داده‌ها (گزارش پاکسازی داده‌ها)</p> <p>ساخت داده‌ها (صفت‌های اشتقاقی، تولید رکورد)</p> <p>یکپارچه سازی داده‌ها (ادغام داده‌ها)</p> <p>فرمت داده‌ها (فرمت بندی مجدد داده‌ها)</p> <p>مجموعه داده‌ها (توصیف مجموعه داده‌ها)</p>	<p>آماده سازی داده‌ها</p>
<p>انتخاب تکنیک مدلسازی (تکنیک‌های مدلسازی، فرضیات مدلسازی)</p> <p>ایجاد طرح آزمون (طرح آزمون)</p> <p>ساخت مدل (تنظیم پارامترها، مدل‌ها، توصیف مدل‌ها)</p> <p>ارزیابی مدل (ارزیابی مدل، بازیابی تنظیمات پارامترها)</p>	<p>مدلسازی</p>
<p>ارزیابی نتایج (ارزیابی نتایج داده کاوی با توجه به معیارهای موفقیت تجاری، مدل‌های تاییدشده)</p> <p>مرور فرآیند (مرور فرآیند)</p> <p>تعیین مراحل بعدی (فهرست راهکارهای ممکن، تصمیم گیری)</p>	<p>ارزیابی</p>
<p>طراحی برای گسترش و استقرار (طرح گسترش و استقرار)</p> <p>طراحی برای نظارت و نگهداری (طرح نظارت و نگهداری)</p> <p>تهیه گزارش نهایی (گزارش نهایی، ارائه نهایی)</p> <p>بازیابی پروژه (مستندسازی تجربیات)</p>	<p>گسترش و استقرار</p>

شکل ۳: وظایف عمومی به همراه نتایج در مدل مرجع CRISP-DM

۱. شناخت و درک کسب و کار



شکل ۴: شناخت و درک کسب و کار

۱-۱ تعیین اهداف فعالیت تجاری

وظیفه: تعیین اهداف فعالیت تجاری

یک تحلیلگر داده در نخستین گام باید بتواند با مد نظر قرار دادن چشم انداز فعالیت تجاری درک کاملی از خواسته های واقعی مشتری حاصل کند. مشتری ها غالباً چندین هدف موازی و محدود دارند که باید در توازی منطقی قرار گیرند. هدف تحلیلگر در آغاز آنست که عوامل مؤثر بر خروجی های پروژه را شناسایی کند. بی دقتی در این مرحله ممکن است ما را متحمل هزینه های زیادی کند که به خاطر اصرار ما در یافتن جوابی درست برای سؤالی نادرست ایجاد می شوند.

خروجی ها:

پس زمینه^{۱۸}

ثبات اطلاعات مربوط به فعالیت های تجاری شرکت در آغاز پروژه

اهداف فعالیت تجاری

با توجه به چشم انداز فعالیت تجاری هدف اصلی مشتری را شرح دهید. علاوه بر هدف اصلی فعالیت تجاری سؤالات دیگری وجود دارد که مشتری مایل است مورد بحث قرار دهد. برای مثال هدف اصلی یک فعالیت تجاری ممکن است حفظ مشتری از طریق پیش بینی زمانی که او به همکاری با شرکت رقیب تمایل پیدا می کند باشد. در این صورت سؤالات مربوطه می توانند اینگونه طرح شوند: «چگونه کانال اصلی ارتباط با مشتریان یک بانک (از قبیل ATM، حضور در شعبه یا از طریق اینترنت) اثرگذاری بیشتری خواهد داشت، چه آن ها بمانند یا بروند». یا اینکه «آیا با کاهش کارمزد سیستم ATM تعداد مشتریان ارزشمندی که ما را ترک خواهند کرد، به طور معناداری کاهش خواهد یافت؟»

معیار موفقیت فعالیت تجاری

از منظر فعالیت تجاری، معیاری برای موفقیت خروجی ها تعیین کنید. این معیار می تواند کاملاً ملموس و واقعی باشد، مانند «کاهش مشتریانی که شما را ترک میکنند و به رقیب می پیوندند»، و هم می تواند کلی و ذهنی باشد، مثلاً «کسب یک بینش مفید از روابط». در مورد دوم باید مشخص شود ارزیابی ذهنی از جانب چه کسی صورت می گیرد.

۲-۱ ارزیابی وضعیت^{۱۹}

وظیفه: ارزیابی وضعیت

این وظیفه جزئیات بیشتری از اطلاعات به دست آمده در مورد همه منابع، محدودیت ها و فرضیات است و همچنین شامل اطلاعات لازم از سایر عواملی که در تعیین هدف داده کاوی و طراحی پروژه باید مورد توجه قرار گیرند. در وظیفه قبلی هدف دستیابی هر چه سریع تر به صورت مسأله یک وضعیت بود. ولی هدف این وظیفه دقت هر چه بیشتر در جزئیات است.

خروجی ها:

فهرست منابع

این فهرست موارد زیر را در بر می گیرد:

پرسنل (متخصص کسب و کار، متخصص داده ها، پشتیبانی فنی، تیم داده کاوی)، داده ها (استخراج داده های ثبت شده، دسترسی به انبار داده های مؤثر یا داده های عملیاتی)، منابع محاسباتی (سخت افزارها) و نرم افزار (ابزارهای داده کاوی، دیگر نرم افزارهای مرتبط)

ملزومات، فرضیات و محدودیت ها

ملزومات پروژه عبارتند از: برنامه زمانی اجرا، قابل درک بودن، امنیت، کیفیت نتایج و همچنین رسیدگی حقوقی. به عنوان بخشی از خروجی مطمئن شوید اجازه استفاده از داده ها را دارید.

فهرستی از فرضیاتی که در حین پروژه ساخته می شوند تهیه کنید. این فهرست هم می تواند شامل فرضیاتی باشد که در حین داده کاوی رسیدگی می شوند هم شامل فرضیاتی غیر قابل بررسی در مورد فعالیت تجاری که در صورت مواجهه با آن ها پروژه متوقف می شود.

همچنین فهرستی در مورد محدودیت های پروژه تهیه کنید. این فهرست می تواند هم شامل محدودیت هایی در مورد دسترسی بودن منابع باشد هم شامل محدودیت های فنی. (مانند اندازه داده هایی که به صورت مؤثر در مدلسازی مورد استفاده قرار می گیرند)

پیشامدهای احتمالی

فهرستی از خطرات و حوادثی که ممکن است به تأخیر یا ناکارآمدی در روند پروژه منجر شوند تهیه کنید. همچنین فهرستی از اقدامات مناسب برای مواجهه با هر کدام از آن پیشامدها تهیه کنید.

لغات فنی

واژه نامه ای از لغات فنی مربوط به پروژه تهیه کنید. این واژه نامه می تواند شامل دو بخش باشد: (۱) واژه نامه ای در مورد کلمات تخصصی فعالیت تجاری که تشکیل دهنده بخشی از درک فعالیت تجاری است و برای انجام پروژه مفید خواهد بود.

۲) واژه نامه ای در مورد لغات تخصصی داده کاوی که با مثال هایی مرتبط با فعالیت تجاری توضیح داده شده است.

هزینه ها و مزیت ها

برای پروژه یک تجزیه و تحلیل هزینه-فایده صورت دهید که هزینه های پروژه را با سودآوری آن (در صورت موفقیت فعالیت تجاری) مقایسه کند. مقایسه تا حد ممکن باید اختصاصی باشد، مثلاً استفاده از پیش بینی های مالی در یک وضعیت تجاری خاص.

۳-۱ تعیین اهداف داده کاوی

وظیفه: تعیین اهداف داده کاوی

اهداف فعالیت تجاری، هدف پروژه را با اصطلاحات خاص فعالیت تجاری و اهداف داده کاوی، هدف را به زبان فنی و تخصصی داده کاوی توضیح می دهد. به طور مثال یک هدف فعالیت تجاری می تواند «افزایش فروش کاتالوگ به مشتریان باشد». هدف داده کاوی متناظر می تواند «پیش بینی میزان خرید یک محصول از جانب مشتری، با توجه به خریدهای او در بیش از سه سال گذشته، اطلاعات جمعیت شناختی (سن، درآمد، شهر و غیره) و قیمت آن محصول» باشد.

خروجی ها:

اهداف داده کاوی

خروجی های منتخب پروژه که قابلیت دستیابی به اهداف فعالیت تجاری را دارند، شرح دهید.

معیار موفقیت داده کاوی

برای یک خروجی موفق، با بیانی فنی و تخصصی معیاری تعریف کنید. مثلاً سطحی معین برای صحت پیش بینی یا برآورد گرایش به خرید با میزان صعود^{۲۰} معین. همچنین با داشتن یک معیار برای موفقیت تجاری ممکن است لازم باشد آن را با عباراتی کیفی توصیف کنیم تا قضاوت در مورد میزان موفقیت تسهیل شود.

۴-۱ ارائه طرح پروژه

وظیفه: ارائه طرح پروژه

طرح های منتخب را برای رسیدن به هدف داده کاوی و در نتیجه رسیدن به هدف فعالیت تجاری تبیین کنید. طرح باید مجموعه پیش بینی شده ای از مراحل قابل اجرا در حین توقف پروژه به انضمام انتخاب اولیه ابزارها و تکنیک ها باشد.

خروجی ها:

طرح پروژه

فهرستی از مراحل که در پروژه باید اجرا شود، به اضافه مدت زمان اجرای آن‌ها، منابع مورد نیاز، ورودی‌ها، خروجی‌ها و وابستگی‌ها تهیه کنید. در پیاده سازی فرآیند داده کاوی در موارد مورد نیاز، عملیات را به دفعات زیاد تکرار کنید. به عنوان مثال تکرار در فازهای مدلسازی و ارزیابی. به عنوان بخشی از طرح پروژه، تجزیه و تحلیل وابستگی‌های میان برنامه زمانی و ریسکها از اهمیت بسزایی برخوردار است. تأثیر نتایج تجزیه و تحلیل را در طرح پروژه دقیقاً مشخص کنید. همچنین اقدامات و پیشنهادات مقتضی در مواجهه با هر ریسک را تعیین کنید.

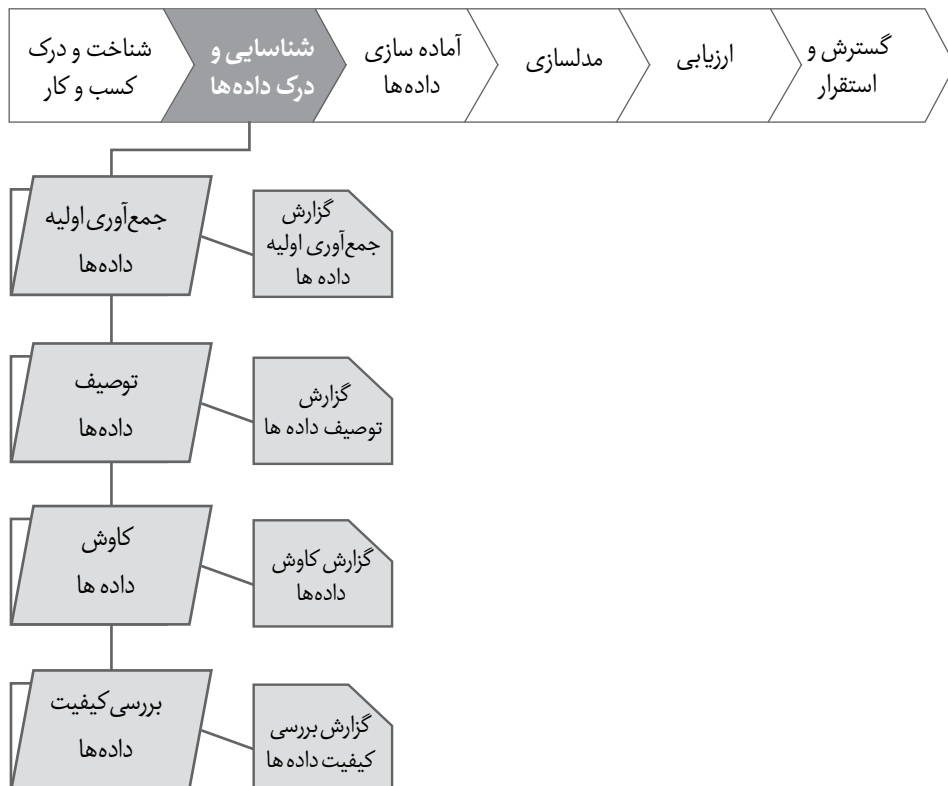
توجه: طرح پروژه شامل جزئیات طرح‌های هر فاز می‌باشد. به طور مثال در همین مرحله مشخص کنید قصد دارید از کدام استراتژی ارزیابی در فاز ارزیابی استفاده کنید.

طرح پروژه در حکم یک مستند پویاست که بررسی مجدد میزان پیشرفت و موفقیت آن در پایان هر فاز ضروریست و مطابق با پیشنهادات جدید باید بروز شود. همچنین زمان‌های بررسی مجدد را مشخص کنید چرا که این بررسی‌ها هم جزئی از طرح پروژه اند.

ارزیابی اولیه ابزارها و تکنیک‌ها

هم‌چنین در پایان فاز اول پروژه، باید یک ارزیابی اولیه از ابزارها و تکنیک‌ها انجام دهید. در این مرحله می‌بایستی ابزارهایی برای داده کاوی انتخاب کنید که کارایی لازم را در حین به کارگیری روشهای گوناگون برای مراحل مختلف داده کاوی داشته باشند. برای مثال، ارزیابی ابزارها و تکنیکها پیش از فرآیند، از اهمیت بسزایی برخوردار است چون انتخاب ابزارها و تکنیکها می‌تواند بر همه پروژه تأثیر گذار باشد.

۲. شناسایی و درک داده ها



شکل ۵: شناسایی و درک داده ها

۱-۲ جمع آوری اولیه داده ها

وظیفه: جمع آوری اولیه داده ها

دسترسی به داده‌ها در یک پروژه (یا دستیابی به داده‌ها)، در منابع پروژه لیست می‌شود. اگر برای فاز درک داده لازم باشد این جمع آوری اولیه شامل بارگذاری داده می‌شود. به عنوان مثال، در صورت استفاده از یک ابزار خاص برای درک داده، بارگذاری داده‌ها می‌تواند درک مناسبی از آن برای شما ایجاد نماید. این کار ممکن است منجر به مراحل اولیه آماده سازی داده‌ها شود.

تذکر: اگر چندین منبع داده به دست آورید آنگاه بایستی یکپارچه سازی را نیز یا در این مرحله یا در مرحله بعد صورت دهید.

خروجی: گزارش جمع آوری اولیه داده ها

فهرستی از مجموعه داده (یا مجموعه داده‌های) قابل دستیابی، بعلاوه موقعیت آن‌ها در پروژه، روش‌های دسترسی به آن‌ها در پروژه و مسائل و مشکلاتی که امکان رخ دادن دارند، تهیه کنید. مسائل پیش آمده و راه حل آن‌ها برای استفاده در دیگر مراحل همین پروژه (اگر باز هم تکرار شوند) یا پروژه‌های مشابه دیگر را یادداشت کنید.

۲-۲ توصیف داده ها

وظیفه: توصیف داده ها

ناخالصی و دیگر خواص ظاهری داده‌های به دست آمده را بررسی کنید و نتایج را گزارش دهید.

خروجی: گزارش توصیف داده ها

داده‌های به دست آمده را از نظر فرمت داده‌ها، مقدار داده‌ها (تعداد رکوردها و متغیرها در جدول)، ماهیت متغیرها و سایر ویژگی‌های ظاهری توصیف کنید. آیا داده‌های موجود نیازهای پروژه را برآورده می‌کنند؟

۳-۲ کاوش داده ها

وظیفه: کاوش داده ها

در این وظیفه پرسش‌های داده کاوی با استفاده از تحقیق، مجسم سازی و گزارش دهی طرح می‌شوند. این مرحله شامل توزیع فیلهای کلیدی (مانند متغیر هدف یک عمل پیش بینی)، ارتباط بین دو یا چند صفت، نتایج ترکیبات ساده، خواص معنی دار خرده جمعیت‌ها، و آنالیز آماری ساده است.

این تجزیه و تحلیل ها می تواند به طور مستقیم اهداف داده کاوی را مشخص کند و با استفاده از آن می توان توصیف داده ها و گزارش کیفیت ها را تصحیح کرد، و همچنین تبدیل ها و دیگر کارهای آماده سازی داده های مورد نیاز را برای تجزیه و تحلیل های ثانویه صورت داد.

خروجی ها: گزارش کاوش داده ها

نتایج این مرحله، شامل رابطه های کشف شده و فرضیات اولیه و تأثیر آن بر باقیمانده پروژه را بررسی کنید. در صورت لزوم از رسم نمودار که تعیین کننده خصوصیات داده هاست و موجد زیر مجموعه جالب توجهی از داده ها که امکان بررسی بیشتر را فراهم می کند استفاده کنید.

۲-۴ بررسی کیفیت داده ها

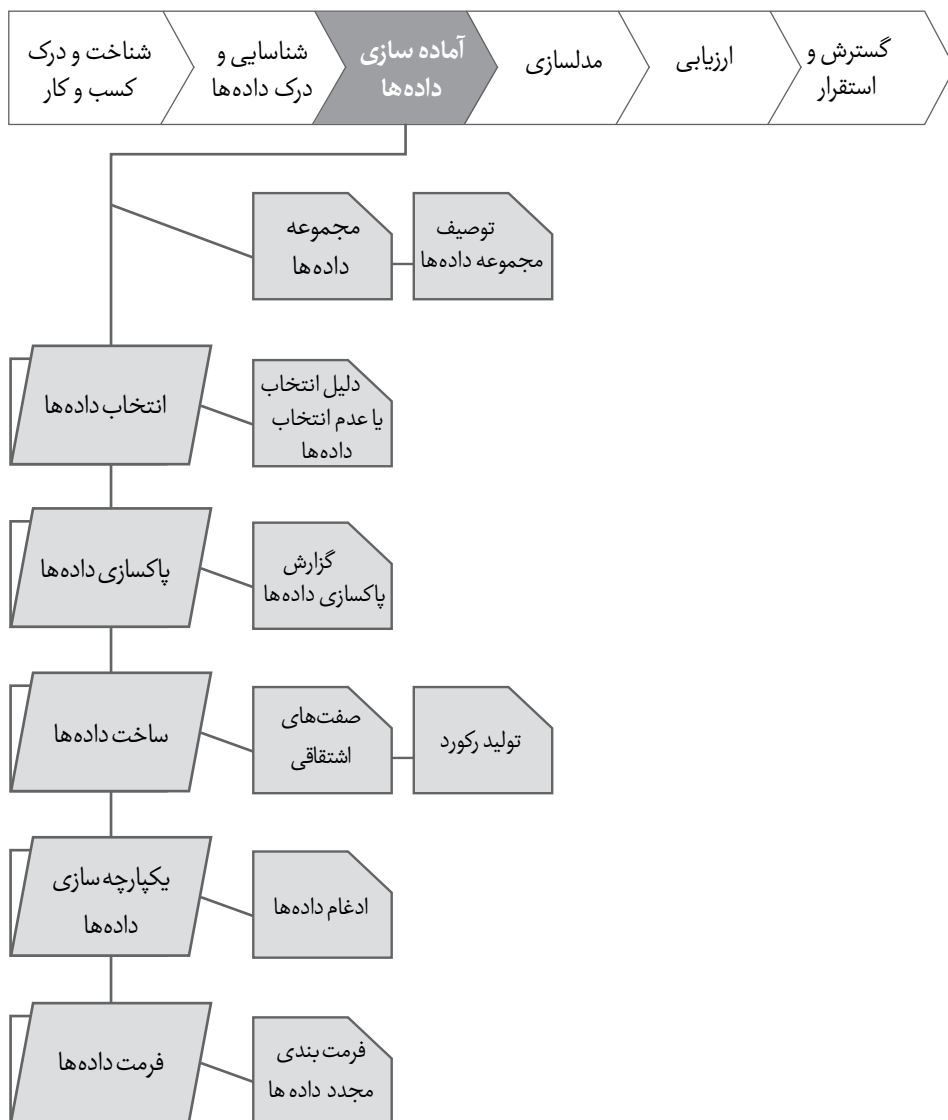
وظیفه: بررسی کیفیت داده ها

کیفیت داده ها را با سوالاتی مانند: "آیا داده ها کامل اند؟ (آیا تمام موارد لازم را می پوشانند؟) آیا داده ها صحیح اند یا خطا دارند؟ و اگر خطا دارند، خطای آن ها چقدر شایع است؟ آیا در داده ها مقادیر گمشده وجود دارد؟ اگر هست چگونه نمایش داده می شوند، کجا اتفاق افتاده اند و چقدر شایع اند؟" بررسی نمایید.

خروجی: گزارش کیفیت داده ها

لیستی از نتایج بازرسی کیفیت داده ها تهیه کنید. اگر مشکل کیفی وجود دارد، راه حل های ممکن را لیست کنید. این راه حل ها وابستگی زیادی به اطلاعات داده کاوی و تجاری دارند.

۳. آماده سازی داده ها



شکل ۶: آماده سازی داده ها

خروجی ها:**مجموعه داده ها**

مجموعه داده‌ها محصول فاز آماده سازی داده است، که برای مدل‌سازی یا تجزیه و تحلیل‌های عمده و اصلی پروژه استفاده خواهد شد.

توصیف مجموعه داده ها

مجموعه داده ای را که برای مدل‌سازی یا تجزیه و تحلیل‌های اصلی پروژه استفاده خواهند شد، تشریح کنید.

۱-۳ انتخاب داده ها**وظیفه: انتخاب داده ها**

داده‌هایی را که برای تجزیه و تحلیل استفاده می‌شوند انتخاب کنید. معیارها شامل مرتبط بودن با اهداف داده کاوی، محدودیت‌های فنی و کیفی از قبیل محدودیت در حجم داده‌ها یا نوع داده‌هاست. توجه داشته باشید که انتخاب داده‌ها شامل انتخاب صفت‌ها (ستون‌ها) و همچنین انتخاب رکورد‌ها (سطرها) در جدول است.

خروجی: دلیل انتخاب یا عدم انتخاب داده ها

فهرستی از داده‌هایی که تصمیم بر انتخاب یا عدم انتخاب آن‌ها گرفته‌اید تهیه کنید و دلایل تصمیم خود را نیز ذکر کنید.

۲-۳ پاکسازی داده ها**وظیفه: پاکسازی داده ها**

کیفیت داده‌ها را تا سطح مورد نیاز تکنیک‌های تجزیه و تحلیل انتخاب شده افزایش دهید. این کار ممکن است شامل انتخاب زیرمجموعه‌ای بی نقص از داده‌ها، اضافه کردن مقادیر قراردادی مناسب یا تکنیک‌های جاه طلبانه تری مانند برآورد داده‌های گمشده با استفاده از مدل‌سازی باشد.

خروجی: گزارش پاکسازی داده ها

چگونگی اجرای اقدامات و تصمیماتی که تصمیم به انجام آن‌ها گرفته بودید را شرح دهید (این اقدامات و تصمیمات در گزارش مشکلات کیفی داده‌ها در جریان وظیفه بازرسی کیفیت داده‌ها در فاز شناسایی و درک داده‌ها ثبت شده بودند). همچنین تبدیلات روی داده‌ها، به منظور پاکسازی و بررسی تأثیرات آن روی نتایج تجزیه و تحلیل بایستی صورت گیرد.

۳-۳ ساخت داده ها

وظیفه: ساخت داده ها

این وظیفه شامل عملیات آماده سازی برای ایجاد داده های مفید می باشد، از قبیل ساخت صفت های اشتقاقی^{۲۱}، ورود رکوردهای جدید دست نخورده یا تبدیل مقادیر صفت های موجود.

خروجی ها:

صفت های اشتقاقی

صفت های اشتقاقی صفت های جدیدی هستند که از یک یا تعداد بیشتری از صفت های موجود با رکوردهای یکسان ساخته می شوند. مثلاً مساحت یک صفت اشتقاقی است که از دو صفت طول و عرض به این صورت ساخته می شود: مساحت = طول × عرض

تولید رکورد

ایجاد رکوردهای جدید را به طور کامل تشریح کنید. مثال: ساخت رکورد برای مشتری هایی که در طی سال قبل خریدی نداشته اند. دلیلی برای وجود این نوع رکوردها در داده های اولیه نیست؛ اما برای مدلسازی، عمل صحیح آن است که مشتریان بدون خرید هم نمایش داده شوند.

۳-۴ یکپارچه سازی داده ها^{۲۲}

وظیفه: یکپارچه سازی داده ها

روش هایی وجود دارد که به وسیله آن ها اطلاعات از چندین جدول یا منابع اطلاعاتی مختلف ترکیب می شوند و رکوردها یا مقادیر جدید می سازند.

خروجی: ادغام داده ها

ادغام جداول یعنی دو یا چند جدول را که حاوی اطلاعات متفاوتی در مورد موضوع واحدی هستند به هم الحاق کنیم. به طور مثال یک خرده فروش زنجیره ای یک جدول با اطلاعاتی در مورد مشخصه های عمومی فروشگاه (از قبیل سطح بنا، پیاده رو مجاور)، جدولی با اطلاعات خلاصه شده فروش (از قبیل سود و درصد تغییر در فروش از یک سال قبل) و جدولی با اطلاعات جمعیت شناختی درباره مناطق مجاور دارد. در هر کدام از این جداول یک سطر برای هر فروشگاه در نظر گرفته شده است. می توان این جداول را با هم ادغام کرد و یک جدول جدید ساخت که در آن برای هر فروشگاه یک سطر در نظر گرفته شده است و اطلاعات هر سطر ترکیبی از اطلاعات جداول قبل است.

همچنین داده های ادغام شده، جمعیت داده ها را تحت پوشش قرار می دهد. تجمیع داده ها عبارت

است از عملیاتی که مقادیر جدید را به وسیله خلاصه کردن اطلاعات چندین رکورد یا جدول محاسبه می‌کند. برای مثال، تبدیل یک جدول خرید مشتریان که در آن برای هر خرید یک رکورد ثبت شده است، به جدول جدیدی که برای هر مشتری یک رکورد ثبت شده است و صفت‌های آن عبارتند از تعداد خرید، مقدار متوسط خرید، درصد سفارشات، هزینه کردن کارت اعتباری، درصد اقلام در حال گسترش و غیره.

وظیفه: فرمت داده‌ها

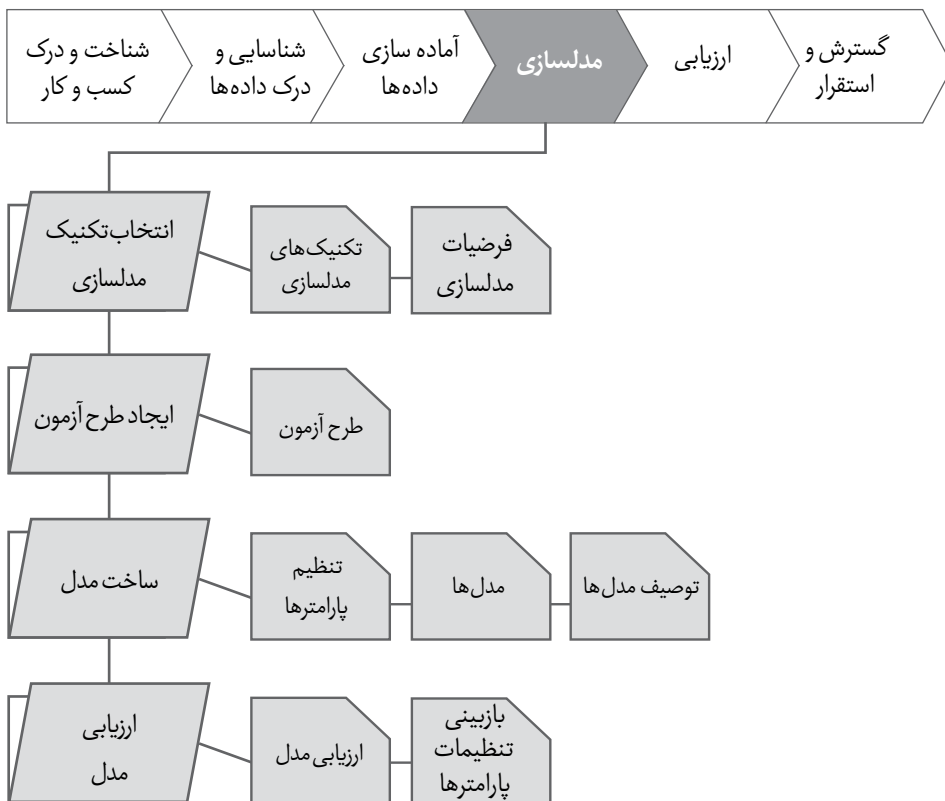
تغییر فرمت داده‌ها شامل اصلاحات نحوی است و در آن ارزش و مفهوم داده‌ها تغییر نمی‌کند. این کار در صورت لزوم برای ابزارهای مدلسازی استفاده می‌شود.

خروجی: فرمت بندی مجدد داده‌ها

بعضی نرم افزارها، شرایطی برای متغیرها دارند، مثلاً محتوای اولین متغیر، شناسه‌های یکتا برای هر سطر از داده‌ها باشد، یا محتوای آخرین متغیر می‌بایستی متغیر خروجی باشد که به وسیله مدل پیش بینی می‌شود. گاهی لازم است ترتیب رکوردها را در مجموعه داده‌ها تغییر دهیم. ممکن است ابزارهای مدلسازی به مرتب کردن رکوردها بر حسب مقادیر متغیرهای خروجی نیاز داشته باشند. معمولاً رکوردهای یک مجموعه داده، در ابتدا به صورت‌های خاصی مرتب شده‌اند، اما الگوریتم مدلسازی بر حسب لزوم، آن‌ها را به صورت کاملاً تصادفی مرتب می‌کند. برای مثال، هنگامی که از شبکه‌های عصبی استفاده می‌کنیم، به طور کلی بهتر است که رکوردها را با ترتیب تصادفی نمایش دهیم، اگر چه نرم افزارها، این کار را به صورت اتوماتیک و بدون دخالت مستقیم کاربر انجام می‌دهند.

به علاوه به طور کلی تغییرات نحوی دیگری نیز وجود دارد که برای برآورده کردن نیازهای نرم افزارهای مدلسازی خاصی به کار می‌روند. مثلاً حذف کاما از درون فایل‌های متنی که در آن‌ها از جدا کننده کاما استفاده شده باشد، یا حذف تمام مقادیری که بیش از ۳۲ کاراکتر دارند.

۴. مدل سازی



شکل ۷: مدل سازی

۴-۱ انتخاب تکنیک مدلسازی

وظیفه: انتخاب تکنیک مدلسازی

نخستین گام در مدلسازی، انتخاب تکنیکی است که می‌خواهیم مورد استفاده قرار دهیم. از آنجا که احتمالاً نرم افزار مورد استفاده را پیش از این در فاز شناسایی و درک کسب و کار انتخاب نموده‌اید، این وظیفه شامل چند تکنیک خاص مدلسازی ست؛ تکنیک‌هایی مانند: ساخت درخت تصمیم با استفاده از C4.5 یا ایجاد شبکه عصبی به روش پس انتشار^{۲۳}. اگر از چندین تکنیک استفاده شود برای هر یک از تکنیک‌ها این وظیفه را جداگانه مورد استفاده قرار دهید.

خروجی‌ها:

تکنیک‌های مدلسازی

تکنیک‌هایی را که برای مدلسازی از آن‌ها استفاده کرده‌اید یادداشت کنید.

فرضیات مدلسازی

در بسیاری از تکنیک‌های مدلسازی فرض‌هایی در مورد داده‌ها صورت می‌گیرد، مانند: همه صفت‌ها توزیع یکنواخت دارند، هیچ داده‌ای گم نشده است، صفت طبقه بندی کننده باید به صورت نمادی^{۲۴} باشد. دیگر فرض‌های صورت داده شده را نیز یادداشت کنید.

۴-۲ ایجاد طرح آزمون^{۲۵}

وظیفه: ایجاد طرح آزمون

پیش از ساخت یک مدل، باید روشی برای سنجش کیفیت و اعتبار مدل ایجاد کنیم. به عنوان مثال در وظایف مربوط به داده کاوی نظارتی مانند طبقه بندی، از نرخ خطا برای سنجش کیفیت پیش‌بینی‌ها در مدل‌های داده کاوی استفاده می‌شود. بنابراین معمولاً داده‌ها را به مجموعه‌های آموزشی و آزمایشی تقسیم می‌کنیم و پس از ایجاد مدل روی داده‌های آموزشی، کیفیت آن را روی مجموعه جدا شده آزمایشی مورد بررسی قرار می‌دهیم.

خروجی: طرح آزمون

طرح انتخاب شده برای آموزش، آزمون و ارزیابی را تشریح کنید. یکی از بخش‌های ابتدایی طرح، نحوه تقسیم بندی مجموعه داده به داده‌های آموزش و آزمایشی، و ارزیابی این مجموعه هاست.

۲۳ Back Propagation

۲۴ Symbolic

۲۵ Test Design

۳-۴ ساخت مدل

وظیفه: ساخت مدل

با استفاده از نرم افزارهای داده کاوی، برای مجموعه داده‌های آماده شده یک یا چند مدل بسازید.

خروجی‌ها:

تنظیم پارامترها

در هر نرم افزار مدل‌سازی، عموماً پارامترهایی وجود دارند که باید تنظیم شوند. فهرستی از پارامترها و مقادیر انتخاب شده برای آن‌ها همراه با دلیل انتخاب آن‌ها، تهیه کنید.

مدل‌ها

مدل‌های واقعی که به وسیله نرم افزارهای مدل‌سازی ساخته شده‌اند.

توصیف مدل‌ها

مدل حاصله را تشریح کنید. گزارشی درباره تفسیر مدل‌ها تهیه کرده و اشکالات مدل و مفهوم آن‌ها را یادداشت کنید.

۴-۴ ارزیابی مدل

وظیفه: ارزیابی مدل

تفسیر یک مهندس داده کاو از مدل، بر اساس میزان اطلاعاتش از حوزه کاربرد مدل، معیارهای موفقیت داده کاوی و تست‌های طراحی شده شکل می‌گیرد. این وظیفه با فاز ارزیابی (بخش بعد) در تداخل است. به دلیل آنکه مهندس داده کاو مرجع قضاوت در مورد میزان موفقیت مدل‌سازی است و بر مبنای اصول فنی در جستجوی تکنیک‌های بیشتر است، باید در تماس دائم با تحلیلگران تجاری و متخصصین آن حوزه خاص باشد تا بتواند نتایج داده کاوی را در زمینه تجاری توضیح دهد. در این قسمت فاز ارزیابی مدل و دیگر نتایج حاصل در ضمن پروژه به یک اندازه مورد بررسی قرار می‌گیرند.

یک مهندس داده کاو باید جدولی برای رتبه بندی مدل‌ها ایجاد کند. این رتبه بندی با توجه به معیارهای ارزیابی صورت می‌گیرد و حتی الامکان معیارهای موفقیت تجاری و نیز مد نظر قرار دادن موارد واقعی تجارت را در بر می‌گیرد. در بیشتر پروژه‌های داده کاوی، مهندس داده کاو از یک تکنیک چند بار استفاده می‌کند یا نتایج داده کاوی را با تکنیک‌ها جایگزین به دست می‌آورد. در این بخش او همه نتایج را با توجه به معیارهای موفقیت با هم مقایسه می‌کند.

خروجی‌ها:

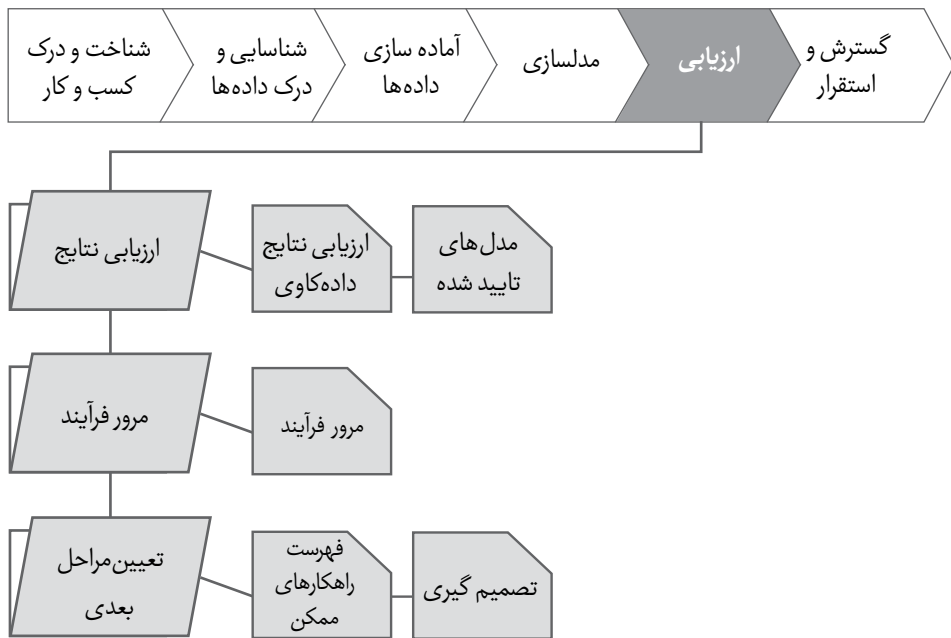
ارزیابی مدل

نتایج این بخش را خلاصه کنید. فهرستی از کیفیت مدل‌های تولید شده (مثلاً بر حسب دقت) تهیه کنید و آن را رتبه بندی کنید.

بازبینی تنظیمات پارامترها

با توجه به ارزیابی مدل، تنظیمات پارامترها بازبینی می‌شود و برای اجراهای آینده در بخش مدلسازی آن‌ها را تنظیم می‌کنیم. مدلسازی را تکرار می‌کنیم و مدل‌های تازه را ارزیابی می‌کنیم تا کاملاً مطمئن شویم مدل حاصل بهترین مدل است. تمام بازنگری‌ها و ارزیابی‌ها را یادداشت می‌کنیم.

۵. ارزیابی



شکل ۸: ارزیابی

۱-۵ ارزیابی نتایج

وظیفه: ارزیابی نتایج

در مراحل قبلی ارزیابی ویژگی هایی مانند صحت و جامعیت مدل ها مورد بررسی قرار گرفت. در این مرحله میزان تناسب مدل با اهداف فعالیت تجاری بررسی می شود و در صورت وجود نقص در مدل، پاره ای از دلایل تجاری آن تشخیص داده می شوند. از دیگر خصیصه های فاز ارزیابی، آزمودن مدل در موارد واقعی است؛ البته اگر این امکان از نظر هزینه و زمان میسر باشد.

هم چنین دیگر نتایج ناشی از مدلسازی را ارزیابی کنید. نتایج داده کاوی می تواند هم پوشش دهنده مدلهایی که مرتبط با هدف اصلی داده کاوی هستند باشد، هم دربرگیرنده یافته های دیگری که لزوماً ارتباط چندانی به هدف اصلی داده کاوی نداشته باشند؛ اگر چه ممکن است حاوی اطلاعاتی باشند که در ادامه مورد استفاده قرار گیرد.

خروجی ها:

ارزیابی نتایج داده کاوی با توجه به معیارهای موفقیت تجاری نتایج ارزیابی را با توجه به معیارهای موفقیت تجاری خلاصه کنید و قضاوت نهایی خود را در مورد تناسب پروژه با اهداف ابتدایی، به آن ضمیمه کنید.

مدل های تأیید شده

پس از ارزیابی مدل با توجه به معیارهای موفقیت تجاری، مدل های منطبق با معیارها، در لیست مدل های تأیید شده قرار می گیرند.

۲-۵ مرور فرآیند

وظیفه: مرور فرآیند

در این مرحله به مدل های رضایت بخشی دست یافته ایم که خواسته های تجاری مان را هم برآورده میکند. اکنون زمان مناسبی است که مروری کلی بر تعهدات داده کاوی داشته باشیم. آیا عامل ها یا ویژگیهای مهمی هست که مورد توجه قرار نگرفته باشد؟ در این مرحله از داده کاوی، فرآیند بازنگری از یک بازنگری کیفی آغاز می شود؛ آیا مدل را به درستی ساخته ایم؟ آیا فقط از متغیرهایی که مجاز بوده ایم استفاده کرده ایم و آیا آن ها برای تجزیه و تحلیل های بعدی هم در اختیار ما قرار خواهند داشت؟

خروجی: مرور فرآیند

فرآیند بازنگری را خلاصه کنید و پیشنهادهایی برای فعالیت های از قلم افتاده یا فعالیت هایی که مجدداً باید تکرار شوند، ارائه کنید.

۳-۵ تعیین مراحل بعدی

وظیفه: تعیین مراحل بعدی

با توجه به ارزیابی نتایج و بازبینی فرآیند باید در مورد نحوه ادامه پروژه تصمیم گیری شود. اینکه آیا پروژه در این مرحله تمام شده و گسترش و استقرار آن شروع گردد یا اینکه تکرارهای بعدی آغاز شده و یا اصلاً مقدمات انجام پروژه داده کاوی تازه ای فراهم شود. این بخش شامل بررسی منابع باقی مانده و تأثیرگذاری بودجه بر نحوه تصمیم گیری خواهد بود.

خروجی ها:

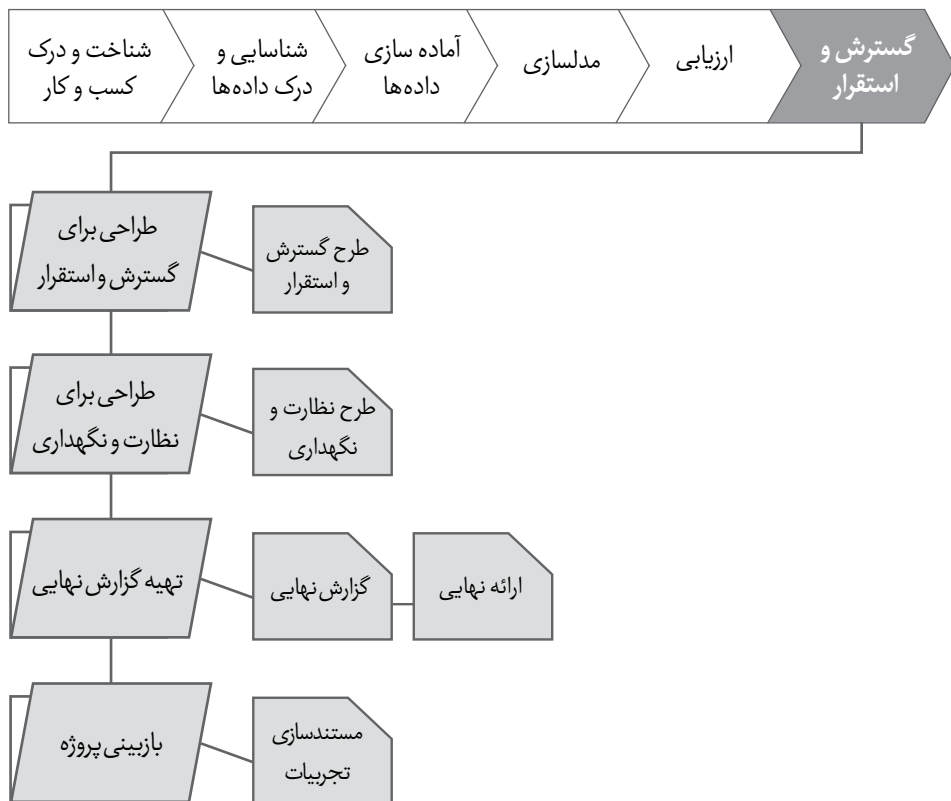
فهرست راهکارهای ممکن

فهرستی از راهکارهای ممکن بعدی، به همراه دلایلی در تأیید یا رد آن ها تهیه کنید.

تصمیم گیری

تصمیم خود در مورد نحوه پیگیری کار را به همراه دلایل آن شرح دهید.

۶. گسترش و استقرار



شکل ۹: گسترش و استقرار و استقرار

۱-۶ طراحی برای گسترش و استقرار

وظیفه: طرح گسترش و استقرار

این فاز بر اساس نتایج ارزیابی، یک استراتژی مناسب برای گسترش و استقرار نتایج تعیین می‌کند. اگر برای ساخت مدل روال عمومی مناسبی به دست آورده‌اید پرونده‌ای شامل مشروح آن برای توسعه‌های بعدی تهیه کنید.

خروجی: طرح گسترش و استقرار

استراتژی گسترش و استقرار را همراه با گام‌های لازم و چگونگی انجام آن‌ها به طور خلاصه بیان کنید.

۲-۶ طراحی برای نظارت و نگهداری^{۲۶}

وظیفه: طرح نظارت و نگهداری

اگر نتایج داده‌کاوی از اطلاعات روزانه تجاری و محیط آن حاصل شده باشد، نظارت و نگهداری آن بسیار حائز اهمیت است. یک آماده‌سازی با دقت در استراتژی نگهداری منجر به عدم استفاده نادرست در بلندمدت از نتایج داده‌کاوی خواهد شد. با هدف نظارت بر گسترش و استقرار نتایج داده‌کاوی، پروژه به یک طرح تفصیلی از فرآیند نظارت نیاز دارد. این طرح نیازمند برخی گسترش و استقرارهای خاص می‌باشد.

خروجی: طرح نظارت و نگهداری

استراتژی نظارت و نگهداری را همراه با گام‌های لازم و چگونگی انجام آن‌ها به طور خلاصه بیان کنید.

۳-۶ تهیه گزارش نهایی

وظیفه: تهیه گزارش نهایی

در پایان کار، سرپرستان پروژه و تیم آن‌ها یک گزارش نهایی تهیه می‌کنند. موضوع گزارش بستگی به طرح گسترش و استقرار دارد، این گزارش هم می‌تواند خلاصه‌ای از وضعیت پروژه و نتایج تجربی آن باشد هم می‌تواند گزارشی نهایی برای آشنایی با نتایج داده‌کاوی باشد.

خروجی‌ها:

گزارش نهایی

این آخرین گزارش تعهد شده در پروژه داده‌کاوی است. این گزارش شامل همه موارد ذکر شده قبلی و همچنین نتایج خلاصه‌سازی شده و سازماندهی شده می‌باشد.

ارائه نهایی

این امر در صورت برگزاری نشست در پایان پروژه برای معرفی شفاهی نتایج به مشتری انجام می شود.

۴-۶ بازبینی پروژه

وظیفه: بازبینی پروژه

این امر شامل ارزیابی صحت امور انجام شده و بهبود مواردی که نیازمند اصلاحند می باشد.

خروجی: مستندسازی تجربیات

خلاصه ای از تجربیات مهم کسب شده در طی پروژه تهیه کنید. به عنوان مثال ریسک ها، موارد گمراه کننده یا راهنمایی هایی برای انتخاب مناسب ترین تکنیک های داده کاوی در شرایط مشابه می تواند بخشی از این خلاصه سازی را به خود اختصاص دهد. در یک پروژه ایده آل، پرونده تجارب، شامل گزارش های شخصی همه اعضا از قسمت های مختلف پروژه می باشد.



فصل سوم
راهنمای کاربر
CRISP-DM

۱. شناخت و درک کسب و کار

۱-۱ تعیین اهداف فعالیت تجاری

وظیفه: تعیین اهداف فعالیت تجاری

یک تحلیلگر داده در نخستین گام باید بتواند با مد نظر قراردادن چشم انداز فعالیت تجاری درک کاملی از خواسته های واقعی مشتری حاصل کند. مشتری ها غالباً چندین هدف موازی و محدود دارند که باید در توافقی منطقی قرار گیرند. هدف تحلیلگر در آغاز آنست که عوامل مؤثر بر خروجی های پروژه را شناسایی کند. بیدقتی در این مرحله ممکن است ما را متحمل هزینه های زیادی کند که به خاطر اصرار ما در یافتن جوابی درست برای سؤالی نادرست ایجاد می شود.

خروجی: پس زمینه

جمع آوری و تطبیق اطلاعات اولیه درباره وضعیت کسب و کار شرکت در شروع پروژه انجام می پذیرد. این جزئیات علاوه بر کمک به شناخت دقیق اهداف تجاری مورد نظر، برای شناخت دقیقتر منابع انسانی و غیر انسانی که در خلال پروژه مورد استفاده قرار می گیرند نیز سودمند است.

فعالیت ها:

سازماندهی

- ایجاد چارت سازمانی برای بخشها و گروه های درگیر پروژه. چارت تنها باید شامل نام مدیران و مسئولیت هایشان باشد.
- شناسایی افراد کلیدی و نقش آن ها در فعالیت تجاری
- شناخت یک حامی داخلی (که هم حامی مالی باشد، هم یک کاربر اصلی و متخصص در حوزه مربوط)
- بررسی وجود یک هیأت هدایت کننده و شناخت اعضای آن
- شناخت واحد های تجاری که با پروژه داده کاوی دارای اثر متقابل اند. (بازاریابی، فروش، قسمت مالی)

محدوده مسئله

- شناخت محدوده مسئله (مانند بازاریابی، توجه به مشتری، توسعه تجارت و...)
- توضیح مسئله در حالت کلی
- بررسی شرایط فعلی پروژه (بررسی اینکه آیا واحدی که برای آن داده کاوی انجام می شود شناخت کافی از داده کاوی دارد یا اینکه داده کاوی به عنوان یک تکنولوژی کلیدی باید به آن ها شناسانده شود؟)
- تعیین پیش نیازهای پروژه (انگیزه داده کاوی چیست؟ آیا داده کاوی کاربردی در کسب و

کار کنونی دارد؟)

- معرفی داده کاوی در فعالیت تجاری در صورت لزوم
- شناخت گروه هدف برای نتایج پروژه (آیا قصد شما تهیه یک گزارش برای مدیریت عالی است یا راه اندازی سیستمی برای کاربرانی بی اطلاع)
- شناخت انتظارات و خواسته های کاربران.

راه حل های کنونی

- بیان تمام راه حل هایی که هم اکنون برای این مسئله مورد استفاده قرار می گیرند.
- بیان معایب و مزایای راه حل های کنونی و سطح مقبولیت آن ها نزد کاربران.

خروجی: اهداف فعالیت تجاری

با توجه به چشم انداز فعالیت تجاری هدف اصلی مشتری را شرح دهید. علاوه بر هدف اصلی فعالیت تجاری سؤالات دیگری وجود دارد که مشتری مایل است مورد بحث قرار دهد. برای مثال هدف اصلی یک فعالیت تجاری ممکن است حفظ مشتری از طریق پیش بینی زمانی که او به همکاری با شرکت رقیب تمایل پیدا می کند باشد. در حالی که هدف دوم پروژه ممکن است پی بردن به این سؤال باشد آیا کاهش حق الزحمه، فقط روی بخشی از مشتریان تأثیرگذار است یا روی همه آن ها.

فعالیت ها:

- توضیح غیر رسمی مسأله ای که باید با داده کاوی حل شود.
- آماده کردن سؤالات در چهارچوب فعالیت تجاری مورد بررسی با دقت هر چه بیشتر.
- آماده کردن سایر نیازمندی ها. (در تجارت هیچ کدام از مشتری ها نباید از دست داده شوند)
- تعیین منافع مورد انتظار در چهارچوب فعالیت تجاری.

توجه!

اهداف غیر قابل دستیابی را با اهدافی واقعی و قابل دسترس جایگزین کنید.

خروجی: معیار های موفقیت فعالیت تجاری

معیار های مفید بودن و موفقیت آمیز بودن نتایج یک پروژه را از دیدگاه تجاری شرح دهید. این معیار ها ممکن است کاملاً معلوم و سهل الوصول باشند، مانند کاهش مشتری های در گردش از یک سطح معین؛ و یا کلی و غیر قابل لمس باشند، مانند «ایجاد نگاهی نو در روابط». در مورد دوم باید مشخص شود قضاوت درباره موفقیت آمیز بودن پروژه از جانب چه کسی صورت خواهد گرفت.

فعالیت ها:

- تعیین معیارهای تجارت موفق. (مثلاً بهبود ده درصدی سرعت پاسخ دهی در تبادلات و افزایش بیست درصدی سرعت ثبت نام)
- تعیین ارزیابی کننده معیارها.

یاد آوری!

هر کدام از معیارهای موفقیت باید حداقل با یکی از اهداف تعیین شده فعالیت تجاری در ارتباط باشد.

پیشنهاد!

پیش از ارزیابی وضعیت (گام بعدی) می توانید تجربه های پیشین در مورد حل این مسأله که با CRISP-DM یا روش های دیگر انجام شده اند را ملاحظه کنید.

۲-۱ ارزیابی وضعیت**وظیفه: ارزیابی وضعیت**

این وظیفه جزئیات بیشتری از اطلاعات به دست آمده در مورد همه منابع، محدودیت ها و فرضیات است و همچنین سایر عواملی که در تعیین هدف داده کاوی و طراحی پروژه باید مورد توجه قرار گیرند.

خروجی: فهرست منابع

فهرستی از منابع موجود پروژه تهیه کنید. این فهرست شامل پرسنل (کارشناس داده، پشتیبان فنی، تیم داده کاوی)، داده (استخراج داده های ثبت شده، دسترسی به پایگاه داده های فعال یا غیر قابل استفاده)، منابع محاسباتی (سخت افزارها) و نرم افزار (ابزارهای داده کاوی و دیگر نرم افزارهای مربوط) خواهد بود.

فعالیت ها:**منابع سخت افزاری**

- شناخت سخت افزار پایه
- فراهم کردن سخت افزارهای پایه ای مورد نیاز برای پروژه داده کاوی
- چک کردن برنامه نگهداری سخت افزار در صورت تعارض با دیگر سخت افزارهای موجود برای پروژه
- شناسایی سخت افزارهای موجود برای ابزار داده کاوی مورد استفاده (در صورت شناخته شده بودن ابزار این مرحله)

منابع داده و دانش

- شناسایی منابع داده
- شناسایی انواع منابع داده (منابع آنلاین، منابع خاص، اسناد نوشتاری و...)
- شناسایی منابع دانش
- شناسایی انواع منابع دانش (منابع آنلاین، خاص، اسناد نوشتاری و...)
- بررسی ابزار و تکنیک های موجود
- تشریح زمینه های دانشی مربوطه (رسمی و غیررسمی)

منابع پرسنلی

- شناسایی حامیان پروژه (در صورت تمایز از حامیان داخلی مطرح شده در قسمت ۱-۱-۱)
- شناسایی مدیران سیستم، مدیر پایگاه داده و تیم پشتیبان فنی برای سؤالات بعدی.
- شناسایی تحلیلگر بازار، کارشناسان داده کاوی و آمار و بررسی وجود این افراد.
- بررسی وجود متخصص در زمینه های مرتبط به فازهای بعدی

یادآوری!

از یاد نبرید ممکن است در قسمت هایی از پروژه به کارمندان تیم پشتیبانی فنی احتیاج داشته باشید، مثلاً در مرحله تبدیل داده.

خروجی: نیازها، فرضیات و محدودیت ها^{۲۷}

فهرستی از همه نیازهای پروژه همراه با برنامه زمان بندی تکمیل آن ها، کیفیت و قابلیت درک نتایج، تضمین ها و همچنین پی آمد های قانونی ممکن تهیه کنید. به عنوان قسمتی از خروجی مطمئن شوید اجازه استفاده از داده ها را دارید.

فهرستی از فرضیات حاصل از پروژه تهیه کنید. این فهرست هم می تواند شامل فرضیاتی باشد که در حین داده کاوی رسیدگی می شوند یا فرضیاتی غیر قابل بررسی که بستگی به شرایط تجاری داشته باشد. اهمیت این مسئله زمانی بیشتر می شود که این فرضیات بر درستی نتایج تاثیرگذار باشد.

همچنین فهرستی از محدودیت های ایجاد شده به وسیله پروژه تهیه کنید. این محدودیت ها ممکن است شامل کمبود منابع برای انجام وظایف پروژه یا محدودیت های قانونی و اخلاقی برای استفاده از خود داده ها یا راه حل های داده کاوی باشد.

فعالیت ها:

نیازمندی ها

- نمایه سازی^{۲۸} گروه هدف
- ثبت نیازمندی ها با برنامه زمانی
- ثبت نیازمندی ها به صورتی که قابل درک، دقیق و قابل توسعه باشند و تکرار پذیری پروژه داده کاوی و نتیجه دهی مدل ها را امکان پذیر کند.
- ثبت نیازمندی ها با ذکر تضمین ها، موانع قانونی، موارد محرمانه، اطلاع رسانی ها و زمان بندی های در نظر گرفته شده برای پروژه.

فرضیات

- توضیح همه فرض ها (شامل موارد نامشخص) و شفاف سازی آن ها. (به عنوان مثال، با توجه به سؤالات تجاری، حداقل به چه تعداد مشتری با سن حدود ۵۰ سال نیاز داریم؟).
- تهیه فهرستی از فرضیات با توجه به وضعیت کیفی داده ها. (دقت، میزان دسترسی).
- تهیه فهرستی از فرضیات با توجه به عوامل خارجی. (نتایج اقتصادی، تولیدات رقابتی و پیشرفت های تکنیکی)
- تشریح فرضیاتی که منجر به تولید برآوردهایی گشته اند. (بطور مثال قیمت نرم افزار خاصی کمتر از ۱۰۰۰ دلار برآورد شده است)
- لیست کردن همه فرضیات با توجه به قابل فهم بودن و توضیح آن ها یا توضیح مدل. (مدل و نتایج آن از چه طریقی باید به مدیران و حامیان معرفی شود)

محدودیت ها

- چک کردن محدودیت های عمومی (پی آمدهای قانونی، بودجه، زمان بندی، و منابع)
- چک کردن دسترسی درست به منابع داده ها (محدردیت های دسترسی، رمزهای مورد نیاز)
- چک کردن قواعد فنی دسترسی به داده ها (سیستم های عملیاتی، سیستم مدیریت داده ها، فرمت داده ها) چک کردن دسترسی به اطلاعات مورد نیاز
- چک کردن محدودیت های بودجه (هزینه های ثابت، هزینه های جاری و...)

یادآوری!

لیست فرضیات شامل فرضیات مرتبط با شروع پروژه نیز می باشد، مانند نقطه شروع پروژه.

خروجی: ریسک‌ها و پیشامدها

ریسک‌ها یا پیشامدهایی که ممکن است رخ دهد، برنامه زمانبندی فشرده، هزینه‌ها و نتایج را لیست کنید. همچنین طرح مقتضی^{۲۹} خود را نیز مشخص کنید؛ برای کم کردن فشار و رهایی از پیامدهای منفی چه اقداماتی باید انجام داد؟

فعالیت‌ها:

شناسایی ریسک‌ها

- شناسایی ریسک‌های کسب و کار (مثال: رقیب با نتایج بهتر رشد بالاتری دارد)
- شناسایی ریسک‌های سازمانی (مثال: واحد سفارش دهنده پروژه قادر به تامین هزینه‌های پروژه نباشد)
- شناسایی ریسک‌های مالی (مثال: تامین مالی پروژه در فازهای مختلف آن به نتایج اولیه پروژه داده کاوی بستگی خواهد داشت)
- شناسایی ریسک‌های فنی
- شناسایی ریسک‌هایی که به داده‌ها و پایگاه داده‌ها بستگی دارد. (مثال: ضعف‌های کیفیتی و پوششی داده‌ها)

توسعه طرح‌های مقتضی

- شرایط ممکن برای رخ دادن هر ریسک را تعیین کنید.
- طرح‌های مقتضی را توسعه دهید.

خروجی: مجموعه لغات

واژه نامه ای مرتبط با مفاهیم بکار رفته در پروژه تهیه کنید. این واژه نامه حداقل شامل دو بخش زیر باید باشد:

- (۱) واژه نامه لغات جاری مربوطه که شامل مفاهیم تجاری موجود در پروژه است.
- (۲) واژه نامه لغات داده کاوی که از طریق مثالی مرتبط با مسأله تجاری، توضیح داده شده باشد.

فعالیت‌ها

- واژه نامه‌های قبلی را بررسی کنید. در صورت عدم وجود چنین واژه نامه‌ای برای آماده کردن پیش نویس یک واژه نامه تازه دست به کار شوید.
- برای درک بهتر لغات از متخصصان آن بخش کمک بگیرید.
- با اصطلاحات تجاری آشنا شوید.

خروجی: هزینه و فایده

یک تجزیه و تحلیل سود و زیان برای پروژه تهیه کنید. این تجزیه و تحلیل باید هزینه های پروژه را با پتانسیل های سود دهی آن (در صورت موفقیت) مقایسه کند.

پیشنهاد!

مقایسه باید حتی الامکان اختصاصی باشد. در این صورت شما قادر به بهبود فعالیت تجاری خود خواهید بود.

فعالیت ها

- هزینه جمع آوری داده را برآورد کنید.
- هزینه تهیه و اجرای یک راه حل را برآورد کنید.
- مزایای استقرار یک راه حل را شناسایی کنید. (افزایش رضایت مشتری، بازدهی سرمایه، و افزایش سود)
- هزینه عملیاتی شدن طرح را برآورد کنید.

توجه!

برآورد هزینه های پنهان از قبیل استخراج داده های تکراری و آماده سازی آن ها، تغییر در جریان کاری و زمان آموزش برای یادگیری را فراموش نکنید.

۳-۱ تعیین اهداف داده کاوی

وظیفه: تعیین اهداف داده کاوی

اهداف فعالیت تجاری، هدف پروژه را با اصطلاحات خاص فعالیت تجاری و اهداف داده کاوی، هدف را به زبان فنی و تخصصی داده کاوی توضیح می دهد. برای مثال یک هدف تجاری ممکن است «افزایش فروش به مشتری های موجود» باشد؛ در حالی که یک هدف داده کاوی ممکن است «پیش بینی تعداد جنسی که یک مشتری خرید خواهد کرد»، «برآورد درآمد یک سال خاص» یا «ارتباط آمار نفوس و قیمت یک کالا» باشد.

خروجی: اهداف داده کاوی

خروجی های قابل توسعه پروژه را که قابلیت دستیابی به اهداف پروژه دارند مشخص نمایید. به کار بستن این تکنیک در خروجی ها رایج است.

فعالیت ها

- مسئله تجاری را به اهداف داده کاوی تبدیل کنید. (بطور مثال یک کمپین بازاریابی نیاز به تقسیم

بندی مشتریان دارد تا درباره افرادی که از آن کمپین استقبال می نمایند تصمیم گیری نماید؛ ضمن آنکه اندازه بخش های تقسیم بندی شده باید مشخص باشد)

- نوع مسئله داده کاوی را تعیین کنید. (به عنوان مثال آیا مسئله طبقه بندی است یا پیش بینی یا خوشه بندی) برای مطالعه بیشتر در این زمینه به پیوست ۲ مراجعه کنید.

پیشنهاد!

در مواقع خاصی بهتر است تعریف دوباره ای از مسأله ارائه دهیم. به عنوان مثال، مدلسازی بقای محصول نسبت به بقای مشتری (زمانی که هدف گذاری بقای مشتری به علت زمان بر بودن مشاهده اثر آن بر نتایج تأثیر اندکی بر خروجی داشته باشد)، در اولویت قرار دارد.

خروجی: معیارهای موفقیت داده کاوی

با بیانی فنی معیارهایی برای موفقیت پروژه تعریف کنید؛ برای مثال سطح معینی برای دقت پیش بینی، یا درجه معینی از صعود^{۳۰} در نمودار تمایل به خرید. گاهی اوقات علاوه بر یک معیار تجارت موفق توضیح شرایط غیرعینی هم لازم است تا امکان قضاوت ذهنی فراهم شود.

فعالیت ها

- معیارهایی برای ارزشیابی مدل تعیین کنید. (بطور مثال: دقت مدل، کارایی و پیچیدگی).
- مبنایی برای معیارهای ارزشیابی معین کنید.
- معیاری برای ارزشیابی های ذهنی تعیین کنید. (بطور مثال: توانایی تفسیرپذیری مدل و قابلیت آن در درک مسائل بازاریابی)

توجه!

فراموش نکنید معیارهای موفقیت داده کاوی با معیارهای موفقیت تجارت متفاوت است. همچنین توصیه می شود طرح گسترش و استقرار پروژه از همان ابتدا تهیه شود.

۴-۱ ارائه طرح پروژه

وظیفه: ارائه طرح پروژه

طرح مورد نظر برای رسیدن به اهداف داده کاوی و به تبع آن رسیدن به اهداف تجاری را تشریح کنید

خروجی: طرح پروژه

فهرستی از مراحل مختلف پروژه همراه با مدت زمان لازم، منابع مورد نیاز، ورودی ها، خروجی ها و وابستگی های آن ها تهیه کنید. تا حد امکان، لزوم تکرارها در حین انجام پروژه را تصریح کنید. به عنوان مثال تکرار در فاز های مدل سازی و ارزیابی. قسمت مهمی از طرح پروژه، آنالیز زمانبندی پروژه و ریسک هاست. در طرح پروژه نتایج پروژه را همراه با پیشنهادی برای عملکرد مناسب در هنگام بروز ریسک صراحتاً مشخص کنید.

اگر چه این بخش تنها قسمتی است که در آن مستقیماً از عنوان طرح پروژه یاد می شود، با این وجود پروژه بایستی دائماً مورد بازبینی قرار گیرد و درباره آن مشورت شود. این کار باید در شروع هر قسمت یا شروع هر اقدامی صورت گیرد.

فعالیت ها

- قسمت ابتدایی طرح پروژه را مشخص کنید و توضیح دهید آیا این قسمت با درگیر کردن تمام کارکنان قابلیت اجرایی دارد؟
- همه اهداف انتخاب شده را وارد کنید و تکنیک ها و شیوه ای انتخاب کنید که مسئله را حل کند و با معیارهای تجاری موفق مطابقت داشته باشد.
- نیرو و منابع مورد نیاز را برای دستیابی به راه حل و توسعه آن برآورد کنید. (این امر می تواند برآورد مقیاس زمانی پروژه داده کاوی را برای دیگر متخصصان تسهیل کند. برای مثال، اغلب ۵۰ تا ۷۰ درصد از وقت و نیرو را برای مرحله آماده سازی داده ها، ۲۰ تا ۳۰ درصد را برای مرحله درک داده، ۱۰ تا ۲۰ درصد برای مدل سازی، ارزیابی و فهم تجاری و ۱۰ تا ۱۵ درصد را مرحله گسترش و استقرار دربر می گیرد).
- مراحل حساس را مشخص کنید.
- نقاط تصمیم گیری را معلوم کنید.
- نقاط تجدید نظر را معلوم کنید.
- تکرارهای عمده را مشخص کنید.

خروجی: ارزیابی اولیه از ابزارها و تکنیک ها

در پایان مرحله اول، فقط یک ارزشیابی اولیه از ابزار و تکنیک ها شکل گرفته است. در اینجا یک ابزار داده کاوی که برای روش های متنوع مراحل مختلف داده کاوی قابل استفاده باشد انتخاب کنید. ارزیابی ابزارها و تکنیک ها از همان لحظه انتخاب آن ها مهم است و می تواند در کل پروژه تأثیرگذار باشد.

فعالیت ها:

- فهرستی از معیارهای انتخاب ابزار و تکنیک ها تهیه کنید (یا از یک لیست قبلاً تهیه شده استفاده کنید)
- ابزارها و تکنیک های بالقوه را انتخاب کنید.
- تناسب تکنیک های موجود را ارزیابی کنید.
- تکنیک های اجرایی را با توجه به ارزیابی راه حلها بررسی و رتبه بندی کنید.

۲. شناسایی و درک داده ها

۱-۲ جمع آوری اولیه داده

وظیفه: جمع آوری اولیه داده

به دست آوردن داده ها در یک پروژه (یا دستیابی به داده ها)، در منابع پروژه لیست می شود. اگر برای فاز درک داده لازم باشد این جمع آوری اولیه شامل بارگذاری^{۳۱} داده می شود. به عنوان مثال، در صورت استفاده از یک ابزار خاص برای درک داده، بهتر است داده هایتان را در این ابزار بارگذاری کنید.

خروجی: گزارش جمع آوری اولیه داده ها

فهرستی از مجموعه داده (یا مجموعه داده های) قابل دستیابی در پروژه به همراه جزئیات مورد نیاز برای این کار تهیه کنید. این گزارش بایستی مشخص کننده متغیرهایی که دارای اهمیت بیشتری هستند باشد. به خاطر داشته باشید که ارزیابی کیفی مجموعه داده ها نباید بر اساس منابع جداگانه باشد. الحاق مجموعه داده ها به یکدیگر باعث کشف اشکالاتی خواهد شد که در منابع اصلی داده ها نیست و این امر ناشی از تفاوت بین منابع است.

فعالیت ها:

برنامه ریزی برای داده های مورد نیاز

- برنامه ریزی برای اطلاعاتی که مورد نیاز است. (بطور مثال: آیا صفت^{۳۲} های فعلی کافیت، یا اطلاعات اضافی دیگر هم مورد نیاز است؟)

Loading ۳۱

Attribute ۳۲

- بررسی میزان دسترس پذیری اطلاعات مورد نیاز.

معیارهای انتخاب

- مشخص کردن معیارهای انتخاب (کدام صفت ها برای اهداف تعیین شده داده کاوی مورد نیاز است؟ کدام صفت ها بی ارتباط شناخته شده اند؟ با توجه به تکنیک های انتخاب شده چند ویژگی (صفت) نیاز ما را برطرف می کند؟)
- انتخاب جداول و فایل های قابل توجه
- انتخاب داده ها از جداول و فایل ها
- بررسی اینکه تحلیل چه مدت زمانی از داده ها مورد نیاز است؟ (بطور مثال: داده های ۱۸ ماه موجود است اما با داده های ۱۲ ماه نیاز ما برطرف می شود.)

توجه!

هنگام جمع آوری داده ها از منابع مختلف، ممکن است در الحاق آن ها با مشکل مواجه شویم (ناسازگاری فرمت ها، بی اعتباری داده ها و ...)

الحاق داده ها

- اگر داده ها دارای متن هستند آیا برای مدل نیاز به تبدیل کد وجود دارد یا بایستی آن ها را گروه بندی نمود؟
 - چگونه می توان متغیرهای گم شده را به دست آورد؟
 - چگونه می توان داده ها را استخراج کرد؟
- پیشنهاد!
- برخی داده ها از منابع غیر الکترونیکی به دست آمده اند. مانند دفاتر ثبتی و علاوه بر آن گاهی اوقات داده ها باید پیش پردازش شوند. (سری های زمانی، میانگین وزنی و ...)

۲-۲ توصیف داده ها

وظیفه: توصیف داده ها

پس از بررسی های لازم، ناخالصی داده های به دست آمده را گزارش کنید.

خروجی: گزارش توصیف داده ها

داده های بدست آمده را توصیف کنید. این توصیف شامل فرمت داده ها، ابعاد داده ها (تعداد رکوردها و فیلدهای هر جدول) به همراه معرفی و شناسایی هر یک از فیلدها می باشد.

فعالیت‌ها:

- سنجش حجمی داده‌ها
- شناسایی داده‌ها و روشهای جمع‌آوری
- دستیابی به پایگاه داده
- استفاده از آنالیز آماری (در صورت نیاز)
- توضیح جداول و روابط حاکم بر آن‌ها
- بررسی حجم داده‌ها، تعداد ابعاد و نوع ترکیب‌ها
- و بالاخره اینکه آیا داده‌ها شامل متن هستند یا خیر؟
- نوع صفت‌ها و مقادیر آن‌ها
- بررسی صفت‌های موجود و میزان دسترسی به آن‌ها
- بررسی نوع صفت‌ها (عددی، اسمی، ترتیبی و ...)
- بررسی دامنه صفت‌ها
- تحلیل همبستگی صفت‌ها
- درک مفهوم و ماهیت هر صفت و مقادیر آن در حوزه کسب و کار
- محاسبات آماری پایه برای هر صفت (محاسبه توزیع داده‌ها، میانگین، ماکزیمم، مینیمم، انحراف معیار، واریانس، مد، چولگی و ...)
- آنالیزهای آماری پایه ای و بیان مفهوم نتایج آن در حوزه کسب و کار
- بررسی مناسبیت هر صفت با هدف خاص داده کاوی
- آیا ماهیت هر یک از صفت‌ها با هدف پروژه سازگار می‌باشند؟
- بررسی نظرات متخصصین مربوطه در مورد نقش هر یک از صفت‌ها
- آیا داده‌ها نیازمند متوازن سازی^{۳۳} می‌باشند؟ (بسته به مدل مورد استفاده)

کلیدها

- تجزیه و تحلیل روابط کلیدی
 - بررسی میزان اشتراک بین صفت‌های کلیدی در جدول
- ### بررسی فرضیات و اهداف
- در صورت نیاز، بروزرسانی لیست فرضیات انجام شود.

۲-۳ کاوش داده ها

وظیفه: کاوش داده ها

این وظیفه شامل ابزارهای پرس و جو، مجسم سازی روابط و گزارش دهی می باشد. این وظایف به واسطه توضیحاتی در ارتباط با داده ها، گزارش های کیفی، تبدیل ها و تکنیک های آماده سازی و تصحیح آن ها، نقش مهمی در نزدیک شدن به اهداف داده کاوی دارد.

خروجی: گزارش کاوش داده ها

نتایج این مرحله، شامل رابطه های کشف شده و فرضیات اولیه و تأثیر آن بر باقیمانده پروژه را بررسی کنید. در صورت لزوم از رسم نمودار که تعیین کننده خصوصیات داده هاست و مورد زیر مجموعه جالب توجهی از داده ها که امکان بررسی بیشتر را فراهم می کند استفاده کنید.

فعالیت ها:

کاوش داده ها

- تجزیه و تحلیل ویژگی های جالب صفت ها با جزئیات کامل (بطور مثال: آماره های مقدماتی، زیر گروه های قابل توجه)

- شناسایی خواص زیر گروه ها

شکل دهی فرضیات برای آنالیزهای بعدی

- ارزیابی اطلاعات و کاوش در گزارش تشریحی داده ها
- شکل دهی فرضیات اولیه و تعیین عملیات مورد نیاز
- (در صورت امکان) تبدیل فرضیات به یک هدف داده کاوی
- اهداف داده کاوی را به صورت دقیق و روشن تنظیم کنید. جستجوی کورکورانه لزوماً بی فایده نیست
- اما یک جستجوی هدایت شده به سمت اهداف داده کاوی ارجحیت دارد.
- اجرای آنالیزهای مقدماتی برای اعتبارسنجی فرضیات

۲-۴ بررسی کیفیت داده ها

وظیفه: بررسی کیفیت داده ها

کیفیت داده ها با را با سؤالات زیر بررسی کنید: آیا داده ها کامل هستند؟ (همه موارد مورد نیاز تحت پوشش قرار میگیرند؟) آیا آن ها صحیح اند یا خطا دارند؟ و اگر خطا دارند، خطای آن ها چقدر شایع است؟ آیا در داده ها مقادیر گمشده وجود دارد؟ اگر هست چگونه نمایش داده می شوند، کجا اتفاق افتاده اند و چقدر شایع اند؟

خروجی: گزارش کیفیت داده ها

فهرستی از نتایج بررسی کیفیت داده ها تهیه کنید؛ اگر مشکلی وجود داشت فهرستی هم برای راه حل های ممکن تهیه کنید.

فعالیت ها:

- مقادیر خاص و متمایز را شناسایی کرده و مفهوم آن ها را مشخص نمایید.
- **مرور کلیدها و متغیرها**
- بررسی شمولیت داده ها (آیا تمامی مقادیر ممکن در مجموعه داده ها وجود دارد؟)
- بررسی کلیدها
- آیا مفهوم هر یک از صفت ها و مقادیر آن ها همخوانی دارند؟
- شناسایی صفت های گمشده و خالی
- درک مفهوم داده های گمشده
- بررسی صفت هایی که مفهوم یکسانی دارند ولی دارای مقادیر مختلفی می باشند. (به طور مثال: رژیم غذایی و میزان چربی غذا)
- بررسی املائی مقادیر (ممکن است مقادیر یکی باشد اما حرف شروع آن ها یکسان نباشد، مثلاً با حروف کوچک شروع شده باشد یا حروف بزرگ)
- بررسی انحرافات و تصمیم گیری در مورد اینکه انحراف موجود، نشان دهنده اختلال است یا نشانه یک پدیده قابل توجه!
- بررسی قابل قبول بودن مقادیر فیلدها (به عنوان مثال: تمام فیلدها مقادیر همسان یا نزدیک به هم داشته باشند)

پیشنهاد!

پس از بررسی صفت ها، معقول بودن مقادیر آن ها را بررسی کنید. (مثلاً یک نوجوان با درآمد بالا).
با استفاده از نمودارها، هیستوگرام و... ناسازگاری ها را آشکار کنید.

کیفیت داده ها در فایل های Flat

- جداکننده^{۳۴} استفاده شده در فایل را بررسی کنید. آیا بین همه صفت ها از جداکننده یکسان استفاده شده است؟
- تعداد صفت ها را در هر رکورد بررسی کنید. آیا تناقضی وجود دارد؟

اختلال ۳° و ناسازگاری بین منابع

- میزان سازگاری و همچنین وجود اطلاعات تکراری بین منابع مختلف را بررسی کنید.
- طرحی برای مواجهه با اختلال ها ترتیب دهید.
- نوع اختلال و همچنین متغیرهای تأثیر گرفته از آن را پیدا کنید.

پیشنهاد!

گاهی لازم است داده‌هایی که رفتارهای مشخصی ندارند از سایر داده‌ها جدا شوند و مورد استفاده قرار نگیرند. (بطور مثال: برای بررسی رفتار وام گیرندگان، حذف اطلاعات کسانی که وامی دریافت نکرده اند و تمام کسانی که هنوز زمان سررسید وامشان نرسیده است، نوعی جداسازی داده‌های با رفتارهای نامشخص است)

قابل قبول بودن فرضیات و همچنین قابلیت آن‌ها برای به دست آوردن اطلاعات کافی را بررسی کنید.

۳. آماده سازی داده ها

خروجی: مجموعه داده

مجموعه داده‌ها محصول فاز آماده سازی داده است، که برای مدلسازی یا سایر تجزیه و تحلیل ها استفاده خواهد شد.

خروجی: توصیف مجموعه داده

مجموعه داده ای را که برای مدلسازی یا تجزیه و تحلیل های اصلی پروژه استفاده خواهند شد، تشریح کنید.

۱-۳ انتخاب داده ها

وظیفه: انتخاب داده ها

داده‌هایی را که برای تجزیه و تحلیل استفاده می‌شوند انتخاب کنید. معیارها شامل مرتبط بودن با اهداف داده کاوی، محدودیت های فنی و کیفی از قبیل محدودیت در حجم داده‌ها یا نوع داده‌هاست.

خروجی: دلیل انتخاب یا عدم انتخاب داده ها

فهرستی از داده‌های انتخاب شده یا حذف شده را به همراه دلیل انتخاب یا عدم انتخاب هر یک از آن‌ها تهیه کنید.

فعالیت ها:

- داده‌های مناسب بیشتری را جمع‌آوری کنید. (از سایر منابع - داخلی و خارجی)
- برای انتخاب فیله‌های مناسب، تست‌های همبستگی و معنا داری را انجام دهید.
- معیارهای انتخاب داده را، با توجه به تجربه مرحله کیفیت داده‌ها، بازبینی کنید (وظیفه ۲-۱ را ببینید؛ چون ممکن است بخواهید حذف یا اضافه‌ای در داده‌ها اعمال کنید).
- معیارهای انتخاب داده را با توجه به نتایج مدل‌سازی بازبینی کنید. (وظیفه ۲-۱ را ببینید؛ ممکن است ارزیابی مدل نشان دهنده نیاز به مجموعه دیگری از داده‌ها باشد).
- زیر مجموعه‌های مختلفی از داده‌ها انتخاب کنید. (بطور مثال: صفت‌های مختلف، داده‌های با شرایط خاص)
- استفاده از تکنیک‌های نمونه‌گیری را در نظر داشته باشید (بطور مثال: یک راه حل سریع ممکن است شامل جداسازی مجموعه داده‌ها به دو بخش آموزشی و آزمایشی می‌باشد، و یا در مواقعی که نرم‌افزار قادر به استفاده از تمامی داده‌ها نیست، می‌توان کاهش حجم داده‌های آزمایشی را مد نظر داشت. همچنین ممکن است نمونه‌گیری بصورت وزنی برای متغیرهایی با درجه اهمیت متفاوت انجام گردد و یا در داخل متغیرهای همسان، مقادیر دارای وزن‌های متفاوت باشند).
- دلایل حذف و اضافه را ثبت کنید.
- تکنیک‌های موجود نمونه‌گیری را بررسی کنید.

پیشنهاد!

بر اساس معیارهای انتخاب داده، مهم‌ترین متغیر یا متغیرها را شناسایی کنید و با قرار دادن آن متغیر به عنوان مبنا، سایر متغیرها را وزن دهی کنید.

۲-۳ پاکسازی داده‌ها

وظیفه: پاکسازی داده‌ها

کیفیت داده‌ها را تا سطح مورد نیاز تکنیک‌های تجزیه و تحلیل انتخاب شده افزایش دهید. این کار ممکن است شامل انتخاب زیرمجموعه‌ای بی‌نقص از داده‌ها، اضافه کردن مقادیر قراردادی مناسب یا تکنیک‌های جاه طلبانه تری مانند برآورد مقادیر گمشده با استفاده از مدل‌سازی باشد.

خروجی: گزارش پاکسازی داده‌ها

چگونگی اجرای اقدامات و تصمیماتی که تصمیم به انجام آن‌ها گرفته بودید را شرح دهید. (این اقدامات و تصمیمات در گزارش مشکلات کیفی داده‌ها در جریان وظیفه بازرسی کیفیت داده‌ها در فاز

شناسایی و درک داده‌ها ثبت شده بودند) گزارش تهیه شده بایستی به شاخص های اصلی کیفیت داده‌های مورد استفاده در پروژه داده کاوی اشاره نموده و تاثیرات احتمالی آن بر نتایج را بیان نماید.

فعالیت‌ها:

- عملکرد خود در مواجهه با انواع اختلال‌ها را بازبینی کنید.
- در مواجهه با اختلال‌ها یکی از این سه راهکار را برگزینید: تصحیح، حذف، یا نادیده گرفتن
- در خصوص مواجهه با مقادیر خاص و مبهم تصمیم‌گیری کنید. دامنه این مقادیر می‌تواند باعث نتایج عجیبی شود و باید به دقت آن‌ها را امتحان کرد. مقادیر خاص می‌تواند ناشی از نتایج بررسی‌هایی باشد که در آن بعضی از سؤالات یا پرسیده نشده باشند یا جوابی نگرفته باشند. این مسئله ممکن است ناشی از ثبت مقدار ۹۹ برای داده‌های ناشناخته باشد. برای مثال ۹۹ را برای وضعیت تأهل یا گرایش سیاسی ناشناخته در نظر بگیریم. همچنین مقادیر خاص می‌تواند برخاسته از داده‌های ناقص و بریده شده باشد، برای مثال «۰۰» برای افراد صد ساله یا برای همه ماشین‌هایی که بیش از ۱۰۰۰۰۰ کیلومتر مسافت طی نموده‌اند.
- معیارهای انتخاب داده را با توجه به تجربه مرحله پاک‌سازی داده‌ها، بازبینی کنید. (ممکن است مجموعه داده‌ای را اضافه یا حذف کنیم)

پیشنهاد!

گاهی اوقات برخی صفت‌ها ارتباط چندانی با اهداف داده کاوی ندارند. در این صورت اختلال در آن‌ها بی‌معنی است. اما در صورت تغییر شرایط، این اختلال‌ها نیز بایستی لحاظ شوند.

۳-۳ ساخت داده‌ها

وظیفه: ساخت داده‌ها

این وظیفه شامل عملیات آماده‌سازی برای ایجاد داده‌های مفید می‌باشد، از قبیل ساخت صفت‌های اشتقاقی، ورود رکورد‌های جدید دست‌نخورده یا تبدیل مقادیر صفت‌های موجود.

فعالیت‌ها:

- مکانیزم ترکیب‌های قابل‌حصول را با توجه به لیست ابزارهای^{۳۶} پیشنهاد شده در پروژه بررسی کنید.
- در مورد انجام ترکیب در داخل ابزار یا خارج از آن تصمیم‌گیری کنید. (کدام یک مؤثرتر است؟)
- معیارهای انتخاب داده را با توجه به تجربیات مرحله ترکیب داده‌ها، بازنگری کنید. (ممکن است بخواهیم مجموعه داده‌ای را اضافه یا حذف کنیم)

خروجی: صفت های اشتقاقی

صفت های اشتقاقی صفت های جدیدی هستند که از یک یا تعداد بیشتری از صفت های موجود با رکوردهای یکسان ساخته می شوند: مثلاً مساحت = طول × عرض .

چرا نیاز به ایجاد صفت های ترکیبی در طی پروژه داده کاوی وجود دارد؟ بایستی توجه نمود که برای ساخت مدل ضرورتی برای استفاده بدون تغییر از داده های خام اولیه وجود ندارد. دلایل ایجاد صفت های ترکیبی، از این قرارند:

- تجربه نشان داده است در بسیاری از مواقع، فاکتورهای بسیار مهمی وجود دارد که بایستی در ساخت مدل وجود داشته باشد، ولی در لیست داده های اولیه نیستند.
- الگوریتم های مورد استفاده انواع خاصی از داده ها را استفاده می کنند؛ برای مثال هنگام استفاده از رگرسیون خطی به داده های خاصی نیاز داریم که باید با ترکیب سایر داده ها ایجاد شوند.
- ممکن است خروجی مرحله مدل سازی نشان دهد فاکتورهای مهمی در نظر گرفته نشده است.

فعالیت ها:

صفت های اشتقاقی

- در مورد متغیرهایی که بایستی نرمالیزه شوند، تصمیم گیری کنید. (مثلاً وقتی از الگوریتم خوشه بندی برای سن و درآمد استفاده می شود، درآمد غالب می گردد)
- نکات لازم در مورد اهمیت ارتباط متغیرها با متغیرهای جدید اضافه شده، را ذکر کنید. (بطور مثال: وزن ها و وزن های استاندارد شده)
- چگونه می توان صفت های گمشده را بسازیم یا وارد کنیم؟ (از کدام نوع ترکیب استفاده کنیم: تجمیع، میانگین، استقرا)
- اضافه کردن صفت های جدید برای دستیابی به داده ها.

پیشنهاد!

پیش از اضافه کردن صفت ترکیبی، صفت هایی که پروسه مدل کردن را راحت تر می کنند، یا به الگوریتم ها کمک می کنند را مشخص کنید. شاید استفاده از صفت «درآمد سرانه» نسبت به استفاده از صفت «درآمد خانوار» بهتر یا راحت تر باشد.

از دیگر انواع صفت های ترکیبی، تبدیلات تک متغیری است، که معمولاً برای تأمین نیازهای ابزار انتخاب شده انجام می گیرد.

فعالیت ها:**تبدیلات تک متغیری**

- تعیین مراحل تبدیلات مورد نیاز با توجه به امکان پذیری آن ها. (بطور مثال: تغییر متغیرهای عددی به گسسته)
- انجام مراحل تعیین شده.

راهنمایی!

تبدیل ها می توانند شامل تغییر داده های عددی به مقادیر نمادین^{۳۷} باشند (مثلاً تبدیل سن به مقاطع سنی) و یا تبدیل فیلدهای نمادین (عالی، خوب، متوسط، بد) به مقادیر عددی. ابزار مدلسازی یا الگوریتم ها معمولاً به این موارد نیاز دارند.

خروجی: تولید رکورد

رکوردهای تولیدی رکوردهای کاملاً جدیدی هستند، که اطلاعات جدیدی را اضافه می کنند یا داده های جدیدی که به صورت دیگری قابل شناسایی نیستند با آن ها نمایش داده می شود. (برای مثال در داده های بخش بندی شده، تولید و استفاده از رکوردهای جدیدی که اشاره به نمونه های اولیه^{۳۸} هر بخش داشته باشد، در پردازش های آتی می تواند مفید باشد.)

فعالیت:

تکنیک های موجود را در صورت نیاز بررسی کنید. (مثلاً مکانیسم ساخت نمونه های اولیه برای هر بخش از داده های بخش بندی شده)

۳-۴ یکپارچه سازی داده ها**وظیفه: یکپارچه سازی داده ها**

روش هایی وجود دارد که به وسیله آن ها اطلاعات از چندین جدول یا منابع اطلاعاتی مختلف ترکیب می شوند و رکوردها یا مقادیر جدید می سازند.

خروجی: ادغام داده ها

ادغام جدول یعنی دو یا چند جدول را که حاوی اطلاعات متفاوتی در مورد موضوع واحدی هستند به هم الحاق کنیم. در این مرحله ممکن است لازم باشد رکورد جدیدی تولید کنیم. همچنین ممکن است

تولید مقادیر تجمیعی^{۳۹} نیز در دستور کار قرار گیرد.

تجمیع داده‌ها عبارت است از عملیاتی که مقادیر جدید را به وسیله خلاصه کردن اطلاعات چندین رکورد یا جدول محاسبه می‌کند.

فعالیت‌ها:

- در صورت نیاز امکان یکپارچه سازی و موثر بودن الحاق منابع ورودی را بررسی کنید.
- منابع داده ای مختلف را یکپارچه کنید.
- معیارهای انتخاب داده را، با توجه به مرحله یکپارچه سازی داده‌ها، بازبینی کنید. (ممکن است مجموعه داده ای را اضافه یا حذف کنیم)

پیشنهاد!

به خاطر داشته باشید برخی اطلاعات دارای فرمت غیرالکترونیکی می باشند.

۳-۵ فرمت داده‌ها

وظیفه: فرمت داده‌ها

تغییر فرمت داده‌ها شامل اصلاحات نحوی است و در آن ارزش و مفهوم داده‌ها تغییر نمی‌کند. این کار در صورت لزوم برای ابزارهای مدلسازی استفاده می‌شود.

خروجی: فرمت بندی مجدد داده‌ها

بعضی نرم افزارها، شرایطی برای متغیرها دارند، مثلا محتوای اولین فیلد، شناسه‌های یکتا برای هر سطر از داده‌ها باشد، یا محتوای آخرین فیلد می‌بایستی متغیر خروجی (هدف) باشد که به وسیله مدل پیش بینی می‌شود.

فعالیت‌ها:

بازآزایی فیلدها

بعضی نرم افزارها، شرایطی برای متغیرها دارند، مثلا محتوای اولین فیلد، شناسه‌های یکتا برای هر سطر از داده‌ها باشد، یا محتوای آخرین فیلد می‌بایستی متغیر خروجی (هدف) باشد که به وسیله مدل پیش بینی می‌شود.

بازآزایی رکوردها

بعضی مواقع تغییر در ترتیب رکوردها در مجموعه داده‌ها مهم می‌باشد. ممکن است ابزار مدلسازی به مرتب سازی رکوردها بر اساس مقادیر خروجی نیاز داشته باشد.

فرمت بندی مجدد مقادیر داخلی

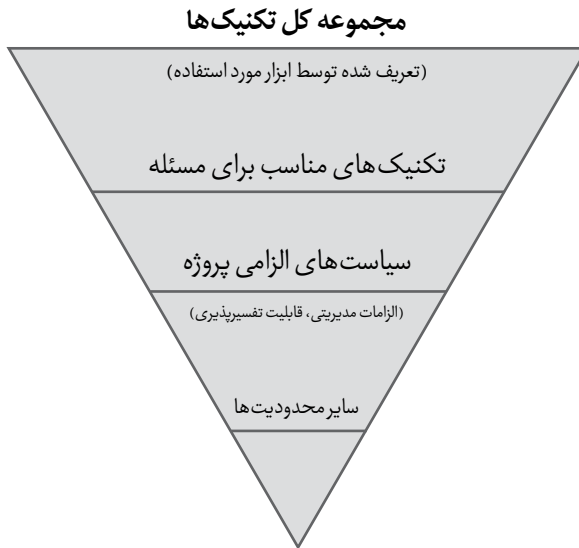
- واضح است که هر تغییر خاص، نیازمندی های ابزار مدل سازی خاصی را تأمین می کند.
- معیارهای انتخاب داده را، با توجه به مرحله پاک سازی داده ها بازبینی کنید. (ممکن است مجموعه داده ای را اضافه یا حذف کنیم)

۴. مدل سازی

۴-۱ انتخاب تکنیک مدل سازی

وظیفه: انتخاب تکنیک مدل سازی

نخستین گام در مدل سازی، انتخاب تکنیکی است که می خواهیم مورد استفاده قرار دهیم. اگر از چندین تکنیک استفاده شود برای هر یک از تکنیک ها این وظیفه را جداگانه مورد استفاده قرار دهید. فراموش نکنید همه تکنیک ها و ابزارها در همه قسمت ها مورد استفاده نیستند. برای هر مسئله خاص فقط تعداد محدودی از تکنیک ها مورد محاسبه قرار می گیرد. (ضمیمه ۲ به بررسی تکنیک های مناسب برای انواع خاص مسائل داده کاوی می پردازد) سیاست های^۴ موجود در پروژه (شامل الزامات مدیریتی، قابلیت تفسیرپذیری نتایج و ...) به همراه سایر محدودیت ها نیز با تنگتر نمودن این حلقه، امکان انتخاب درست را به تحلیلگر می دهند. این امکان وجود خواهد داشت تا فقط یک ابزار یا تکنیک برای حل مسئله در اختیار باشد و یا حتی ابزار قابل دسترس دارای بهترین نتیجه از نقطه نظر تکنیک های مورد نیاز نباشد.



شکل ۱۰: روند انتخاب تکنیک‌های مناسب

خروجی: تکنیک‌های مدل‌سازی

مدل‌هایی که تاکنون مورد استفاده قرار گرفته‌اند را ثبت کنید.

فعالیت:

تکنیک‌هایی را که برای مدل‌سازی از آن‌ها استفاده کرده‌اید یادداشت کنید.

خروجی: فرضیات مدل‌سازی

در بسیاری از تکنیک‌های مدل‌سازی فرض‌هایی در مورد داده‌ها صورت می‌گیرد.

فعالیت‌ها:

- فرضیاتی که برای ساخت مدل در خصوص داده‌ها در نظر گرفته شده‌است را تعریف نمایید. (کیفیت، فرمت و توزیع داده‌ها)
- فرضیات مورد استفاده را با فرض‌هایی که در گزارش شرح داده‌ها آمده، مقایسه کنید.
- از برقراری فرضیات مورد استفاده در داده‌ها اطمینان حاصل کنید و در صورت لزوم به مرحله آماده‌سازی داده‌ها بازگردید.

۴-۲ ایجاد طرح آزمون

وظیفه: ایجاد طرح آزمون

پیش از ساختن یک مدل، باید روشی برای سنجش کیفیت و اعتبار مدل ایجاد کنیم. به عنوان

مثال در وظایف مربوط به داده کاوی نظارتی، مانند طبقه بندی، از نرخ خطا برای سنجش کیفیت پیش بینی ها در مدل های داده کاوی استفاده می شود. بنابراین معمولاً داده ها را به مجموعه های آموزشی و آزمایشی تقسیم می کنیم و پس از ایجاد مدل روی داده های آموزشی، کیفیت آن را روی مجموعه جدا شده آزمایشی مورد بررسی قرار می دهیم.

خروجی: طرح آزمون

طرح انتخاب شده برای آموزش، آزمون و ارزیابی را تشریح کنید. یکی از بخش های ابتدایی طرح، نحوه تقسیم بندی مجموعه داده به داده های آموزش و آزمایشی، و ارزیابی این مجموعه هاست.

فعالیت ها:

- طرح های آزمون موجود را به صورت جداگانه برای هر هدف داده کاوی بررسی کنید.
- برای مراحل مورد نیاز تصمیم گیری کنید. (تعداد تکرارها، تعداد فیلدها و ...)
- تغییرات مورد نیاز داده ها را برای آزموده شدن مهیا کنید.

۳-۴ ساخت مدل

وظیفه: ساخت مدل

با استفاده از نرم افزارهای داده کاوی، برای مجموعه داده های آماده شده یک یا چند مدل بسازید.

خروجی: تنظیم پارامترها

در هر نرم افزار مدل سازی، عموماً پارامترهایی وجود دارند که باید تنظیم شوند. فهرستی از پارامترها و مقادیر انتخاب شده برای آن ها همراه با دلیل انتخاب صورت گرفته، تهیه کنید.

فعالیت ها:

- پارامترهای اولیه را تنظیم کنید.
- دلایل انتخاب مقادیر را مستند کنید.

خروجی: مدل ها

با استفاده از نرم افزارهای داده کاوی، برای مجموعه داده های آماده شده یک یا چند مدل بسازید.

فعالیت ها:

- تکنیک های انتخاب شده را روی داده های ورودی اجرا کنید و مدل بسازید.
- نتایج داده کاوی را پردازش کنید. (بطور مثال: ویرایش قوانین و نمایش درختی)

خروجی: توصیف مدل

مدل حاصله را تشریح کنید. گزارشی درباره تفسیر مدل ها تهیه کنید و اشکالات مدل و مفهوم آن ها را یادداشت کنید.

فعالیت ها:

- ویژگی های مدل موجود را شرح دهید، چون ممکن است در آینده مورد استفاده قرار گیرد.
- تنظیمات مربوط به هریک از پارامترهای مدل را ثبت کنید.
- جزئیات مدل و هریک از ویژگی خاص آن را به صورت دقیق تشریح کنید.
- فهرستی برای مدل های قاعده مند تهیه کنید که شامل قواعد مدل و ارزیابی هریک از قوانین براساس میزان پوشش و صحت باشد.
- برای مدل های بی قاعده، هر گونه اطلاعات تکنیکی در مورد مدل (مانند توپولوژی شبکه عصبی) و دیگر ویژگی های ایجاد شده به وسیله فرآیند مدلسازی (مانند صحت و حساسیت) را فهرست کنید.
- رفتارهای مدل و تفسیر آن ها را شرح دهید.
- نتیجه مربوط به الگوها را بیان کنید (در صورت وجود)، گاهی اوقات مدل بدون انجام پروسه ارزیابی، ویژگی های مهمی در مورد داده ها بیان می کند. (بطور مثال: خروجی مدل، با یکی از ورودی ها برابر یا همسان است)

۴-۴ ارزیابی مدل**وظیفه: ارزیابی مدل**

ارزیابی به معنای بررسی تناسب مدل با معیارهای موفقیت داده کاوی است. این بررسی یک ارزیابی کاملاً تکنیکی است که بر مبنای خروجی های وظیفه مدلسازی صورت می گیرد.

خروجی: ارزیابی مدل

نتایج این بخش را خلاصه کنید. فهرستی از کیفیت مدل های تولید شده (مثلاً بر حسب دقت) تهیه کنید و آن را رتبه بندی کنید.

فعالیت ها:

- نتایج را با توجه به معیارهای ارزیابی بررسی کنید.
- با یک استراتژی مناسب نتایج را آزمایش کنید. (بطور مثال: ایجاد مجموعه داده آموزشی و آزمایشی، ارزیابی متقاطع^{۴۱}، بوت استراپ^{۴۲} و ...)
- نتایج ارزیابی و تفسیر آن ها را مقایسه کنید.
- با توجه به معیارهای ارزیابی و موفقیت نتایج را رتبه بندی کنید.
- بهترین مدل را انتخاب کنید.

- نتایج را از نگاه تجاری تفسیر کنید. (تا حدی که در این مرحله ممکن است)
- به کمک متخصصین حوزه کسب و کار و داده ها، توضیحاتی به مدل ضمیمه کنید.
- معقول بودن مدل را بررسی کنید.
- هم جهت بودن آن با اهداف داده کاوی را بررسی کنید.
- بدیع و مفید بودن اطلاعات اکتشافی حاصل از مدل را با توجه به دانش موجود بررسی نمائید.
- قابلیت اطمینان مدل را بررسی کنید.
- قابلیت گسترش و استقرار هریک از نتایج را تجزیه و تحلیل کنید.
- در صورت وجود توضیحات شفاهی در مورد ایجاد مدل (ارزیابی قواعد، منطقی بودن آن ها، عملی بودن آن ها، تعداد آن ها، و...) آن ها را بیان کنید.
- نتایج را بطور کلی ارزیابی کنید.
- در خصوص چرایی تاثیر مثبت یا منفی یک تکنیک یا پارامتر بر نتایج مدل، بینش مناسبی ایجاد نمائید.

پیشنهاد!

برای تعیین میزان خوب بودن قدرت مدل ها در پیش بینی، ایجاد "جداول صعود"^۳ یا "جداول بهره"^۴ می تواند کمک کننده باشد.

خروجی: بازیابی پارامترها

با توجه به ارزیابی مدل، تنظیمات پارامترها بازیابی می شود و برای اجراهای آینده در بخش مدل سازی آن ها را تنظیم می کنیم. مدل سازی را تکرار میکنیم و مدل های تازه را ارزیابی می کنیم تا کاملاً مطمئن شویم مدل حاصل بهترین مدل است.

فعالیت:

پارامترها را تنظیم کنید تا به مدل بهتری دست پیدا کنید.

۵. ارزیابی

در مراحل قبلی ارزیابی ویژگی هایی مانند صحت و جامعیت مدل ها مورد بررسی قرار گرفت. در این مرحله میزان تناسب مدل با اهداف فعالیت تجاری بررسی می شود و در صورت وجود نقص در مدل، پاره

ای از دلایل تجاری آن تشخیص داده می‌شوند. در این مرحله نتایج کسب شده با معیارهای ارزیابی مشخص شده در ابتدای پروژه مقایسه می‌شوند.

یک راه مناسب برای تعیین مجموع خروجی های یک پروژه داده کاوی استفاده از تساوی زیر است:
نتایج = مدل ها + یافته ها

این تساوی نشان می‌دهد که مدل ها، همه خروجی های یک پروژه داده کاوی نیستند (اگرچه مدل بخش مهمی از آن است)، بلکه یافته ها نیز بخش مهمی از خروجی اند و نشان دهنده مسائل مهمی در مسیر نیل به اهداف داده کاوی می‌باشند. یافته ها همچنین به سؤالات تازه و برخی اثرات جانبی (مانند مسئله کیفیت داده‌ها که در داده کاوی تحت پوشش قرار نمی‌گیرد) اشاره می‌کنند.
تذکر: اگرچه مدل مستقیماً با مسائل تجاری در رابطه است، اما یافته ها لزوماً مرتبط به مسأله یا هدف خاصی نیستند و اهمیت آن‌ها در مراحل ابتدایی پروژه است.

۵-۱ ارزیابی نتایج

وظیفه: ارزیابی نتایج

این مرحله میزان تطابق مدل با اهداف تجاری را بررسی نموده و در صورت مشاهده نقص در نتایج مدل سعی در یافتن دلایل تجاری برای آن می‌باشد. از دیگر اختیارات فاز ارزیابی، آزمودن مدل در موارد کاربردی و واقعی است؛ البته اگر محدودیت های زمان و بودجه اجازه این امر را بدهد. علاوه بر این، فاز ارزیابی سایر نتایج تولید شده به وسیله داده کاوی را نیز مورد سنجش قرار می‌دهد. نتایج داده کاوی شامل مدل های مرتبط با اهداف اصلی کسب و کار و سایر یافته ها می‌باشد. بعضی از آن‌ها می‌تواند در ارتباط با اهداف تجاری بوده و برخی نیز شامل اطلاعات و راهنمایی هایی باشند که در ادامه سودمند واقع شوند.

خروجی: ارزیابی نتایج داده کاوی با توجه به معیارهای موفقیت تجاری

نتایج ارزیابی را با توجه به معیارهای موفقیت تجاری خلاصه کنید و قضاوت نهایی خود را در مورد تناسب پروژه با اهداف ابتدایی، به آن ضمیمه کنید.

فعالیت ها:

- سعی کنید نتایج بدست آمده را فهمیده و درک نمائید.
- نتایج را بر حسب کاربردشان تفسیر کنید.
- اثر آن‌ها را بر اهداف داده کاوی بررسی نمائید.
- بدیع بودن و سودمندی نتایج داده کاوی را با توجه به پایگاه دانش موجود بررسی کنید.
- نتایج را با توجه به معیارهای موفقیت تجاری ارزیابی کنید. (آیا پروژه به اهداف اصلی کسب و کار

دست یافته است؟)

- نتایج ارزیابی و تفاسیر آن را با یکدیگر مقایسه کنید.
- نتایج را با توجه به معیارهای موفقیت تجاری رتبه بندی کنید.
- تأثیر نتایج را بر هدف کاربردی اولیه، بررسی کنید.
- آیا نشانی از اهداف تجاری جدید در این پروژه یا پروژه‌های بعدی هست؟
- نتایج را برای پروژه‌های بعدی داده کاوی شرح دهید.

خروجی: مدل های تأیید شده

بعد از ارزیابی مدل با توجه به معیارهای موفقیت تجاری، مدل های منطبق با معیارها، در لیست مدل های تأیید شده قرار می گیرند.

۲-۵ مرور فرآیند

وظیفه: مرور فرآیند

در این مرحله از کار به مدل های رضایت بخشی دست یافته ایم که خواسته های تجاری مان را هم برآورده میکند. اکنون زمان مناسبی است که مروری کلی بر تعهدات داده کاوی داشته باشیم، تا اگر عامل مهمی از نظر دور مانده است، لحاظ شود. فرآیند بازنگری در این مرحله از داده کاوی شامل بازنگری کیفی است.

خروجی: مرور فرآیند

فرآیند بازنگری را خلاصه کنید و پیشنهادهایی برای فعالیت های از قلم افتاده یا فعالیت هایی که مجدداً باید تکرار شوند، ارائه کنید.

فعالیت ها:

- خلاصه ای از فرآیند داده کاوی صورت گرفته ارائه کنید.
- فرآیند داده کاوی را مورد تجزیه و تحلیل قرار دهید.
- برای هر مرحله فرآیند به این سوالات پاسخ دهید:
- آیا لزومی برای انجام آن فعالیت در آن مرحله خاص هست؟
- آیا آن فعالیت به صورت بهینه اجرا شده است؟
- فرآیند از چه طریقی می تواند بهبود پیدا کند؟
- نواقص موجود را شناسایی کنید.
- مراحل گمراه کننده را شناسایی کنید.
- راهکارهای جایگزین و مراحل غیر قابل انتظار را شناسایی کنید.

- با توجه به معیارهای موفقیت تجاری نتایج داده کاوی را بازبینی کنید.

۳-۵ تعیین مراحل بعدی

وظیفه: تعیین مراحل بعدی

با توجه به ارزیابی نتایج و بازبینی فرآیند باید در مورد نحوه ادامه پروژه تصمیم گیری کنیم. این که آیا پروژه را در این مرحله تمام کنیم و گسترش و استقرار آن را شروع کنیم یا اینکه تکرارهای بعدی را آغاز کنیم و یا اصلاً مقدمات انجام پروژه داده کاوی تازه ای را فراهم کنیم.

خروجی: فهرست راهکارهای ممکن

فهرستی از راهکارهای ممکن، به همراه دلایلی در تأیید یا رد آن‌ها تهیه کنید.

فعالیت‌ها:

- قابلیت‌های هر نتیجه را برای مورد گسترش قرار گرفتن بررسی کنید.
- قابلیت‌های فرآیند حاضر، برای بهبود یافتن را برآورد کنید.
- منابع باقیمانده را بررسی کنید تا امکان تکرارهای بیشتر فرآیند مشخص شود. (همچنین می‌توانید امکان دستیابی به منابع بیشتر را نیز بررسی کنید).
- یک مسیر جایگزین پیشنهاد کنید.
- اصلاحات لازم را در طرح فرآیند ایجاد کنید.

خروجی: تصمیم گیری

تصمیم خود در مورد نحوه پیگیری کار را به همراه دلایل آن شرح دهید.

فعالیت‌ها:

- راهکارهای ممکن را رتبه بندی کنید.
- یکی از راهکارهای ممکن را انتخاب کنید.
- دلایل انتخاب خود را یادداشت کنید.

۶. گسترش و استقرار

۱-۶ طراحی برای گسترش و استقرار

وظیفه: طرح گسترش و استقرار

این فاز بر اساس نتایج ارزیابی، یک استراتژی مناسب برای گسترش و استقرار نتایج تعیین می‌کند.

خروجی: طرح گسترش و استقرار

استراتژی گسترش و استقرار را همراه با گام های لازم و چگونگی انجام آن ها به طور خلاصه بیان کنید.

فعالیت ها:

- نتایج قابل گسترش و استقرار را خلاصه کنید.
- طرح های دیگری برای گسترش و استقرار ارائه دهید و آن ها را ارزیابی کنید.
- برای هر یک از اطلاعات و دانش بدست آمده از نتایج تصمیم گیری کنید.
- چگونگی انتشار اطلاعات و دانش بدست آمده برای کاربران را تعیین نمایید.
- تعیین کنید چگونه می توان بر استفاده از نتایج نظارت کرد و مزایای آن را محاسبه نمود؟ (در چه مکان هایی قابل اجراست؟)
- درباره نتایج نرم افزاری و مدل های قابل گسترش و استقرار تصمیم گیری کنید.
- چگونه یک مدل یا خروجی نرم افزاری از طریق یک نظام سازمانی قابل گسترش و استقرار است؟
- در یک نظام سازمانی چگونه می توان بر استفاده از نتایج نظارت کرد و مزایای آن را محاسبه کرد؟ (در چه مکان هایی قابل اجراست؟)
- مسائل و مشکلاتی که می توانند هنگام گسترش و استقرار نتایج داده کاوی رخ دهند تعیین کنید.

۶-۲ طراحی برای نظارت و نگهداری

وظیفه: طرح نظارت و نگهداری

اگر نتایج داده کاوی از اطلاعات روزانه تجاری و محیط آن حاصل شده باشد، نظارت و نگهداری آن اهمیت دارد. یک آماده سازی با دقت در استراتژی نگهداری منجر به عدم استفاده نادرست در بلند مدت از نتایج داده کاوی خواهد شد. با هدف نظارت بر گسترش و استقرار نتایج داده کاوی، پروژه به یک طرح تفصیلی از فرآیند نظارت نیاز دارد. این طرح به برخی گسترش و استقرار های خاص نیاز دارد.

خروجی: طرح نظارت و نگهداری

استراتژی نظارت و نگهداری را همراه با گام های لازم و چگونگی انجام آن ها بطور خلاصه بیان کنید.

فعالیت ها:

- حالات پویا را بررسی کنید. (چه چیزهایی در حیطه کاری تغییر می کند؟)
- چگونه می توان بر صحت نتایج نظارت کرد؟
- در چه مواقعی نباید از خروجی های داده کاوی و یا سایر مدل های بدست آمده استفاده کرد؟ معیارهایی برای این موضوع تعریف کنید. (بطور مثال: میزان درستی، مرز صحت قابل قبول، داده های جدید، تغییر در دامنه کاربردهای مدل و...) همچنین مشخص نمایید در صورتیکه از یک زمانی قادر به استفاده

- مناسب از مدل نباشیم چه اتفاقی خواهد افتاد؟ (بروزرسانی مدل، شروع پروژه داده کاوی جدید و...)
- آیا هدف تجاری مرتبط با مدل، با گذشت زمان تغییر می‌کند؟ مسئله اولیه ای را که مدل به قصد حل آن ساخته شد را به صورت کامل ثبت کنید.
 - طرح نظارت و نگهداری را گسترش دهید.

۳-۶ تهیه گزارش نهایی

وظیفه: تهیه گزارش نهایی

در پایان کار، سرپرستان پروژه و تیم آن‌ها یک گزارش نهایی تهیه می‌کنند. موضوع گزارش بستگی به طرح گسترش و استقرار دارد، این گزارش هم می‌تواند خلاصه ای از وضعیت پروژه و نتایج تجربی آن باشد هم می‌تواند گزارشی نهایی برای آشنایی با نتایج داده کاوی باشد.

خروجی: گزارش نهایی

در پایان پروژه حداقل یک گزارش نهایی وجود دارد که شامل همه جزئیات خواهد بود. در این گزارش باید شرحی از نتایج به دست آمده، شرح فرآیند و هزینه های صرف شده، هرگونه انحراف از طرح اصلی، شرح طرح اجرایی و پیشنهاداتی برای پروژه‌های بعدی گنجانده شود. محتوای جزئیات واقعی پروژه تا حد زیادی بستگی به مخاطبان گزارش ویژه دارد.

فعالیت‌ها:

- مشخص کنید چه نوع گزارشی لازم است. (معرفی با اسلاید، خلاصه مدیریت، جزئیات یافته ها، توضیح مدل و...)
- بررسی کنید چقدر به اهداف اولیه داده کاوی پرداخته شده است؟
- گروه هدف را برای گزارش معین کنید.
- ساختار و محتویات گزارش را نشان دهید.
- یافته هایی را که گزارش شامل آن‌ها خواهد بود، انتخاب کنید.
- یک گزارش بنویسید.

خروجی: ارائه نهایی

یک گزارش نهایی ممکن است نیاز به ارائه نهایی به مدیران حامی داشته باشد تا در جریان خلاصه پروژه قرار گیرند. یک ارائه معمولی شامل زیر مجموعه ای از اطلاعات موجود در گزارش است، اما به طرق مختلفی ساختار بندی می‌شود.

فعالیت ها:

- برای ارائه نهایی، گروه هدف خود را انتخاب کنید. (آیا گزارش نهایی تاکنون به آن ها رسیده است؟)
- بخش هایی از گزارش نهایی را که ارائه نهایی شامل آن ها می شود انتخاب کنید.

۴-۶ بازبینی پروژه**وظیفه: بازبینی پروژه**

این امر شامل ارزیابی صحت امور انجام شده و بهبود مواردی که نیازمند اصلاحند می باشد.

خروجی: مستندسازی تجربیات

خلاصه ای از تجربیات مهم انجام شده در طی پروژه تهیه کنید. به عنوان مثال ریسک ها، موارد گمراه کننده یا راهنمایی هایی برای انتخاب مناسب ترین تکنیک های داده کاوی در شرایط مشابه می تواند بخشی از این خلاصه سازی را به خود اختصاص دهد. در یک پروژه ایده آل، پرونده تجارب، شامل گزارش های شخصی همه اعضا از قسمت های مختلف پروژه می باشد.

فعالیت ها:

- با همه افرادی که نقشی در پروژه داشته اند مصاحبه کنید و از آن ها در مورد تجربیاتشان در طول پروژه بپرسید.
- در صورت مصاحبه با کاربران اصلی که پروژه برای آن ها انجام شده است سؤالات زیر را از آن ها بپرسید: آیا راضی هستند؟ چه کارهایی می توانست به صورت بهتری انجام شود؟ آیا نیاز به خدمات بیشتر احساس می شود؟
- خلاصه ای از بازخوردها تهیه کنید و تجربیات کسب شده را مستند کنید.
- فرآیندهای خاص داده کاوی را یادداشت کنید. (چگونه نتایج و تجربیات حاصل از بکارگیری مدل می توانند در فرآیند بازخورد داشته باشند؟)
- چکیده تجربیات مفید خود را برای پروژه های بعدی یادداشت کنید.



فصل چهارم
خروجی‌های
CRISP-DM

این فصل شامل توضیح مختصری در مورد اهداف و محتوای گزارش دهی است. گزارش هایی که در این فصل توضیح داده می شوند برای انتقال نتایج یک فاز به افرادی که درگیر مراحل آن فاز و یا حتی کل پروژه نبوده اند، استفاده می شود. این گزارش ها لزوماً همان خروجی هایی که در مدل مرجع و راهنمای کاربر توضیح داده شده اند نیستند. هدف خروجی ها عموماً مستند سازی نتایج در حین انجام پروژه است.

۱. شناخت و درک کسب و کار

نتایج این فاز را می توان در یک گزارش خلاصه کرد. برای شروع قسمت های زیر پیشنهاد می شوند.

پیش زمینه

پیش زمینه یک نمای کلی از پروژه فراهم می کند که زمینه های فعالیت پروژه را مشخص می کند و نشان می دهد چه مشکلاتی برای انجام پروژه شناخته شده، و چرا داده کاوی به عنوان راه حل این مشکلات انتخاب شده است.

اهداف کسب و کار و معیارهای موفقیت

اهداف کسب و کار، اهداف انجام پروژه را به زبان تجاری بیان می کند. برای هر هدف باید معیارهای موفقیت کسب و کار تعریف شوند که این معیارها، شاخص های روشنی برای تشخیص موفقیت یا عدم موفقیت یک پروژه اند.

موجودی منابع

هدف این بخش تشخیص پرسنل، منابع داده ای، تسهیلات تکنیکی و سایر منابع مفید برای انجام پروژه است.

نیازها، مفروضات و محدودیت ها

این خروجی، نیازهای کلی پروژه در مورد نحوه اجرای آن، نتایج پروژه، فرضیاتی که در مورد طبیعت مسئله صورت می گیرد، داده های مورد استفاده، و محدودیت های اعمال شده بر پروژه را لیست می کند.

ریسک ها و احتمالات

این خروجی مشکلاتی را که احتمال دارد در حین پروژه رخ دهند و عملیاتی را که در شرایط مختلف برای به حداقل رسانی آثار مخرب این مشکلات می توان انجام داد، شرح می دهد.

لغات تخصصی

این بخش امکان یک آشنایی نسبی با موضوع پروژه را در اختیار افراد ناآشنا با پروژه قرار می دهد.

هزینه و فایده

این بخش هزینه های پروژه و سود تجاری پیش بینی شده (در صورت موفقیت پروژه) را شرح

می‌دهد. مزایای دیگری نیز که کمتر محسوس اند (مانند رضایت مشتری) باید در نظر گرفته شوند.

اهداف داده کاوی و معیارهای موفقیت

اهداف داده کاوی بیانگر نتایجی از پروژه است که اهداف کسب و کار را برآورده می‌کنند. همچنین در این بخش راه حل‌ها و رهیافت‌های محتمل و معیارهای موفقیت برای رسیدن به نتایج مذکور، با توجه به ادبیات داده کاوی لیست می‌شوند.

طرح و برنامه پروژه

این قسمت شامل مراحل مختلف اجرای پروژه همراه با زمان و منابع لازم، ورودی‌ها، خروجی‌ها و وابستگی هاست. همچنین چگونگی تکرارهای اصلی فرآیند داده کاوی باید به صورت واضح بیان شود؛ مثلاً تکرار فازهای مدلسازی و ارزیابی.

ارزیابی مقدماتی ابزار و تکنیک‌ها

این قسمت دید اولیه‌ای از ابزارها و تکنیک‌هایی که برای استفاده مناسب به نظر می‌رسند و چگونگی استفاده از آن‌ها را به دست می‌دهد. این خروجی، ابزارها و تکنیک‌های در دسترس و مورد نیاز را لیست کرده و آن‌ها را با نیازمندی‌ها هماهنگ می‌کند.

۲. شناسایی و درک داده‌ها

خروجی‌های فاز درک داده در گزارش‌های مختلفی که باید در حین انجام پروژه و وظایف مربوطه تهیه شوند، مستندسازی می‌شود. این گزارش‌ها داده‌هایی را که در زمان انجام فرآیند داده کاوی کاوش می‌شوند توصیف می‌کند. برای گزارش نهایی خلاصه‌ای از قسمت‌های مرتبط و مهم کافی خواهد بود.

گزارش جمع‌آوری اولیه داده‌ها

این گزارش چگونگی تشخیص و استخراج انواع پایگاه داده را بیان می‌کند. مباحثی که باید ذکر شوند عبارتند از:

- پس زمینه داده‌ها
- لیست پایگاه داده‌ها و حوزه و دامنه‌ای که هر یک شامل می‌شوند.
- روش استخراج یا دسترسی به داده‌ها در هر پایگاه داده
- مشکلاتی که در مرحله استخراج یا دسترسی داده‌ها پیش آمده است.

گزارش توصیف داده‌ها

در این بخش هر یک از مجموعه داده‌های به دست آمده توصیف می‌شوند. مباحثی که باید مطرح شوند عبارتند از:

- توصیف جزئیات هر یک از منابع داده ای
- لیست جداول و دیگر موضوعات مرتبط با هر پایگاه داده
- توصیف هر فیلد شامل واحدها، سیستم کدگذاری مورد استفاده و ...

گزارش کاوش داده ها

- این گزارش، جستجوی داده و نتایج آن را توضیح می دهد.
- مباحثی که باید مطرح شوند عبارتند از:
 - پس زمینه، شامل اهداف کلان کاوش داده ها.
 - برای هر نوع از کاوش داده موارد زیر مشخص گردد:
 - قواعد و الگوهای قابل انتظار
 - روش های تشخیص
 - قواعد و الگوهای یافت شده (قابل انتظار و غیر قابل انتظار)
 - هر موضوع غیر قابل انتظار دیگر
 - نتایج تبدیل داده، پاکسازی و سایر مراحل پیش پردازش
 - نتایج مربوط به اهداف داده کاوی و دورنماهای کسب و کار
 - خلاصه کلیه نتایج بدست آمده

گزارش کیفیت داده ها

- این گزارش کلیت و صحت داده ها را توصیف می کند.
- مباحثی که باید ذکر شوند عبارتند از:
 - پس زمینه، شامل کیفیت مورد انتظار داده ها
 - برای هر مجموعه داده:
 - راهکار مورد استفاده برای ارزیابی کیفیت داده ها
 - نتایج ارزیابی کیفیت داده ها
 - خلاصه ای از نتایج کیفیت داده ها

۳. آماده سازی داده ها

گزارش های مرحله آماده سازی داده، بیشتر بر مراحل پیش پردازش داده به منظور آماده ساختن آن ها برای ساخت مدل های داده کاوی، تمرکز می کنند.

گزارش توصیف مجموعه داده

این گزارش شامل توصیف یک مجموعه داده (پس از پیش پردازش) و فرآیند تولید آن می‌شود. مباحثی که باید مطرح شوند عبارتند از:

- پس زمینه، شامل اهداف کلان و طرح کلی پیش پردازش
- ارائه توضیحات لازم در خصوص انتخاب یا حذف مجموعه داده‌ها. برای هر مجموعه داده:
- توصیف مراحل پیش پردازش، شامل عملیات لازم برای اطمینان از دارا بودن کیفیت مطلوب داده‌ها.
- توصیف مفصل مجموعه داده‌های حاصل، بصورت جدول به جدول و فیلد به فیلد.
- ارائه توضیحات لازم در خصوص انتخاب یا حذف صفات و ویژگی‌ها.
- ارائه هر یافته‌ای که در طول پیش پردازش بدست آمده و کلیه نکاتی که برای پیشرفت کار مناسب به نظر می‌رسد.
- نتیجه گیری و خلاصه سازی.

۴. مدل‌سازی

خروجی‌های تولید شده در طول فاز مدل‌سازی می‌توانند در قالب یک گزارش جمع آوری شوند. بخش‌های زیر برای این گزارش پیشنهاد می‌گردد:

مفروضات مدل‌سازی

این قسمت تمام مفروضاتی را که درباره داده‌ها در نظر گرفته شده، یا به صورت ضمنی در تکنیک‌های مدل‌سازی وارد شده‌اند، به وضوح تعریف می‌کند.

طرح آزمون

این قسمت نحوه ایجاد مدل و چگونگی مورد ارزیابی قرار گرفتن آن‌ها را توضیح می‌دهد. مباحثی که باید مطرح شوند عبارتند از:

- پس زمینه، شرح مدل‌های استفاده شده و ارتباط آن‌ها با اهداف داده کاوی به صورت مختصر.

برای هریک از وظایف مدلسازی:

- شرح مفصل نوع مدل و داده‌های آموزشی استفاده شده
- توضیح چگونگی آزمون و ارزیابی مدل
- توصیف داده‌های مورد نیاز برای آزمون
- برنامه تولید داده‌های آزمون در صورت وجود
- در صورت وجود، ارائه طرح چگونگی تولید مجموعه داده آزمون
- توصیف هرگونه طرح ارزیابی مدل‌ها توسط متخصصین حوزه کسب و کار
- خلاصه طرح آزمون

توصیف مدل

این گزارش مدل‌های نهایی و خلاصه‌ای از فرآیند ایجاد آن‌ها را بیان می‌کند. مباحثی که باید در این گزارش مطرح شوند عبارتند از:

- خلاصه مدل‌های تولید شده

برای هر مدل:

- نوع مدل و رابطه آن با اهداف داده کاوی
- تنظیمات بکاررفته برای تولید مدل
- توصیف جزئیات مدل و همه ویژگی‌های خاص آن

برای مثال:

- برای مدل‌های قانون مند، همه قواعد شکل گرفته، همچنین ارزیابی‌های انجام شده برای هر قانون و مدل کلی از لحاظ دقت و پوشش
- برای مدل‌های بی‌قاعده، همه اطلاعات تکنیکی در مورد مدل (از قبیل توپولوژی شبکه‌های عصبی و...) و نیز لیست توصیفات رفتاری.
- توصیف رفتار مدل و تفسیر آن.
- ثبت نتایج مربوط به الگوهای موجود در داده‌ها (در صورت وجود)؛ گاهی اوقات مدل واقعیات مهمی را در مورد داده‌ها آشکار می‌کند، بدون آنکه ارزیابی جداگانه‌ای بر روی داده‌ها انجام شود. (مثلاً خروجی عیناً با یکی از ورودی‌ها برابر شده است)

- خلاصه‌ای از نتایج

ارزیابی مدل‌ها

این قسمت نتایج آزمایش مدل را بر اساس طرح آزمون شرح می‌دهد.

مباحثی که باید ذکر شوند عبارتند از:

- چکیده ای از مراحل فرآیند و نتایج ارزیابی و نیز هر تخلفی از طرح آزمون برای هر مدل:

- ارزیابی مفصل مدل شامل اندازه‌گیری‌ها (مثل دقت مدل) و تفسیر نتایج و رفتار مدل
- ثبت هر گونه نظری در مورد مدل که متخصصین داده و متخصصین حوزه تجاری مربوط اظهار کرده باشند.
- خلاصه ارزیابی مدل
- ملاحظات و نکاتی در مورد اینکه چرا یک تکنیک مدلسازی خاص به همراه حالت خاصی از تنظیم پارامترها به نتایج خوب یا بدی منجر شده است.
- خلاصه ارزیابی مجموعه کامل مدل‌ها

۵. ارزیابی

ارزیابی نتایج داده کاوی با توجه به معیارهای موفقیت کسب و کار.

این گزارش اهداف داده کاوی را با اهداف کسب و کار و معیارهای موفقیت آن مقایسه می‌کند.

مباحثی که باید مطرح شوند عبارتند از:

- مروری بر اهداف کسب و کار و معیارهای موفقیت آن (که ممکن است در طول دوره زمانی داده کاوی و یا پس از حصول نتایج تغییر کرده باشند)

برای هر معیار موفقیت کسب و کار

- مقایسه دقیق معیارهای موفقیت و نتایج به دست آمده داده کاوی
- نتیجه‌گیری در مورد میزان دستیابی به موفقیت و میزان مناسب بودن فرآیند داده کاوی
- مروری بر موفقیت پروژه؛
- آیا پروژه به اهداف اصلی کسب و کار دست یافته است؟
- آیا خواسته‌ها و اهداف جدیدی در کسب و کار به وجود آمده‌اند که در ادامه پروژه یا در پروژه‌های بعدی باید مورد توجه قرار گیرند؟
- نتیجه‌گیری در خصوص پروژه‌های داده کاوی بعدی.

مرور فرآیند

این قسمت کارایی پروژه را ارزیابی می‌کند و هر عاملی را که بایستی در صورت تکرار پروژه در نظر گرفته شود مد نظر قرار می‌دهد.

لیست فعالیت های ممکن

این قسمت پیشنهاداتی در مورد مراحل بعدی پروژه در اختیار می‌گذارد.

۶. گسترش و استقرار

طرح گسترش و استقرار

این قسمت نحوه به کارگیری نتایج داده کاوی را مشخص می‌کند.

مباحثی که باید ذکر شوند عبارتند از:

- خلاصه نتایج قابل گسترش و استقرار (این خلاصه از گزارش مراحل بعدی به دست می‌آید)
- شرح طرح گسترش و استقرار

طرح نظارت و نگهداری

این طرح چگونگی نگهداری نتایج گسترش و استقرار یافته را مشخص می‌کند.

مباحثی که باید مطرح شوند عبارتند از:

- خلاصه نتایج گسترش و استقرار و شناسایی نتایجی که ممکن است به بروز رسانی نیاز داشته باشند. (با ذکر دلیل)

برای هر نتیجه گسترش و استقرار یافته:

- شرح چگونگی آغاز بروز رسانی (بروز رسانی عادی، بروز رسانی پس از یک رخداد یا نظارت اجرایی)

- شرح چگونگی انجام بروز رسانی

- خلاصه نتایج فرآیند بروز رسانی

گزارش نهایی

گزارش نهایی برای خلاصه سازی پروژه و نتایج آن استفاده می‌شود.

مندرجات این گزارش عبارتند از:

- خلاصه شناخت و درک کسب و کار؛ پس زمینه، اهداف و معیارهای موفقیت
- خلاصه فرآیند داده کاوی
- خلاصه نتایج داده کاوی


- خلاصه نتایج ارزیابی
- خلاصه طرح‌های گسترش و استقرار و نگهداری
- تجزیه و تحلیل هزینه و فایده
- نتیجه‌گیری برای کسب و کار
- نتیجه‌گیری برای داده‌کاوی‌های آینده

۷. خلاصه وابستگی‌ها

جدول زیر ورودی‌های اصلی را برای هر یک از محتواهای خروجی خلاصه می‌کند. این به این معنا نیست که تنها ورودی‌های لیست شده باید مد نظر قرار گیرند. به عنوان مثال اهداف کسب و کار باید درمورد همه خروجی‌ها لحاظ شوند. ضمن آن که مستندات خروجی باید پوشش دهنده نکات خاصی باشند که از جانب ورودی‌ها تصریح شده‌اند.

ارتباط نزدیک با	اشاره به	محتوای خروجی	فازها
		پس‌زمینه	شناخت و درک
لغات فنی	پس‌زمینه	اهداف فعالیت تجاری	کسب و کار
	اهداف فعالیت تجاری	معیار موفقیت فعالیت تجاری	
		فهرست منابع	
	اهداف فعالیت تجاری	ملزومات فرضیات و محدودیت‌ها	
	اهداف فعالیت تجاری : معیار موفقیت فعالیت تجاری	پسامدهای احتمالی	
اهداف فعالیت تجاری	پس‌زمینه	لغات فنی	
طرح پروژه	اهداف فعالیت تجاری	هزینه‌ها و مزیت‌ها	
	اهداف فعالیت تجاری ملزومات فرضیات و محدودیت‌ها	اهداف داده‌کاوی	
هزینه‌ها و مزیت‌ها	اهداف ، فعالیت تجاری ، فهرست منابع ، ملزومات ، فرضیات و محدودیت‌ها ، پیشامدهای احتمالی	طرح پروژه	

شنایایی و درک داده ها	گزارش جمع آوری اولیه داده ها	اهداف فعالیت تجاری، فهرست منابع، اهداف داده کاوی	
	گزارش توصیف داده ها	اهداف فعالیت تجاری، گزارش جمع آوری اولیه داده ها	گزارش کمیت داده ها
	گزارش کیفیت داده ها	اهداف فعالیت تجاری، گزارش جمع آوری اولیه داده ها	
	گزارش کاوش داده ها	اهداف فعالیت تجاری، گزارش جمع آوری اولیه داده ها	
آماده سازی داده ها	مجموعه داده ها و توصیف آن	اهداف فعالیت تجاری، اهداف داده کاوی، گزارش توصیف داده ها، گزارش کمیت داده ها، گزارش کاوش داده ها	
مدلسازی	طرح آزمون	اهداف داده کاوی، معیار موفقیت داده کاوی	
ارزیابی	مدل ها	اهداف داده کاوی	تنظیم پارامترها
	تنظیم پارامترها	اهداف داده کاوی	
	توصیف مدل ها	مدل ها، تنظیم پارامترها، طرح آزمون	
	ارزیابی	معیار موفقیت، داده کاوی، طرح آزمون مدل ها	
	مراحل بعدی	طرح پروژه ارزیابی نتایج بر اساس معیار موفقیت فعالیت تجاری	
گسترش و استقرار	طرح گسترش و استقرار	اهداف فعالیت، تجاری، ملزومات، فرضیات و محدودیت ها	طرح نگهداری
	گزارش و ارائه نهایی	اهداف فعالیت تجاری، لغات فنی، ارزیابی نتایج بر اساس معیار موفقیت فعالیت تجاری	طرح گسترش و استمرار
	مستند سازی تجربیات	طرح پروژه بازبینی فرآیند	

A decorative graphic consisting of numerous thin, parallel lines that fan out from the bottom-left corner towards the center of the page, creating a sense of depth and movement.

فصل پنجم پیوست

۱. واژه نامه و شرح اصطلاحات تخصصی

فعالیت

بخشی از وظیفه در راهنمای کاربر که اعمال لازم برای انجام وظایف را شرح می دهد.

روش شناسی CRISP-DM

عبارتی جامع برای همه مفاهیم شکل گرفته در CRISP-DM.

چهارچوب داده کاوی^{۴۵}

مجموعه ای از محدودیت ها و مفروضات. (مانند نوع مسئله، ابزارها و تکنیک ها و دامنه کاربرد)

نوع مسئله داده کاوی^{۴۶}

طبقه بندی ویژه مسائل داده کاوی (مانند توصیف و خلاصه سازی داده ها، خوشه بندی، توصیف

مفاهیم، طبقه بندی، پیش بینی و تحلیل وابستگی)

وظیفه عمومی^{۴۷}

وظیفه ای که در همه پروژه های داده کاوی وجود دارد.

مدل

مجموعه اعمالی که با اجرای آن ها روی مجموعه داده، توانایی پیش بینی متغیر هدف میسر می شود.

خروجی

نتایج محسوسی که از انجام یک وظیفه حاصل می شوند.

فاز^{۴۸}

عبارتی است برای بخش سطح بالایی از مدل فرآیند CRISP-DM که شامل وظایفی خاص می باشد.

Data Mining Context ۴۵

Data Mining Problem Type ۴۶

Generic Task ۴۷

Phase ۴۸

نمونه فرآیند^{۴۹}

پروژه ای خاص که در قالب مدل فرآیند توصیف می شود.

مدل فرآیند^{۵۰}

ساختار پروژه داده کاوی را تعریف می کند و راهنمایی هایی برای اجرای آن فراهم می نماید و شامل مدل مرجع و راهنمای کاربر است.

مدل مرجع

تقسیم بندی پروژه داده کاوی به فازها، وظایف و خروجی ها.

وظیفه خصوصی سازی^{۵۱}

وظیفه ای که مفروضات خاصی را در چهارچوب خاصی از داده کاوی ایجاد می کند.

وظیفه

قسمتی از فاز که مجموعه ای از فعالیت ها را برای تولید یک یا چند خروجی مشخص می کند.

راهنمای کاربر

راهنمای تخصصی در مورد نحوه انجام پروژه داده کاوی.

۲. انواع مسئله داده کاوی

یک پروژه داده کاوی معمولاً با ترکیبی از انواع مسائل داده کاوی سروکار دارد که از حل شدن مجموع آن ها، مسئله کسب و کار حل می شود.

۱-۲ توصیف و خلاصه سازی داده ها

هدف از توصیف و خلاصه سازی داده ها، توصیف دقیق خصوصیات داده به صورت ابتدایی و جمع بندی شده است. این توصیف دیدی کلی از ساختار داده در اختیار کاربر قرار می دهد. گاهی اوقات توصیف و خلاصه سازی داده می تواند به تنهایی هدف پروژه داده کاوی باشد. به عنوان مثال ممکن است یک خرده فروش علاقمند به تفکیک حجم معاملات مالی بازار بر اساس اقلام خاصی باشد. در این مورد تغییرات و تفاوت ها نسبت به دوره قبل می تواند خلاصه و مورد توجه قرار

Process Instance ۴۹

Process Model ۵۰

Specialized Task ۵۱

گیرد. این نوع مسائل در پایین ترین سطح از مقیاس مسائل داده کاوی قرار می گیرند. اگر چه تقریباً در تمام مسائل داده کاوی، توصیف و خلاصه سازی داده یک هدف میانی در فرآیند است، اما کاربرد اغلب در زمان شروع داده کاوی نه از هدف دقیق تحلیل مطلع است نه از مقدار دقیق داده ها. تحلیل اولیه داده ها می تواند به درک طبیعت داده ها کمک کند و مفروضات بالقوه ای برای اطلاعات ناآشکار پیدا کند. توصیف ساده آماری و تکنیک های مجسم سازی، بینش اولیه ای نسبت به داده ها ایجاد می کند. برای مثال توزیع سن مشتریان و محل زندگی آنان نکات مفیدی در مورد مشتریانی که برای استراتژی های آینده بازار باید مورد توجه بیشتری قرار گیرند در اختیار می نهد. توصیف و خلاصه سازی داده عموماً همراه با سایر مسائل داده کاوی انجام می شود. برای مثال توصیف داده ممکن است به شکل گیری فرضیاتی در مورد بخش های جالب توجهی از داده ها منجر شود. توصیف و خلاصه سازی داده ها زمانی مفید است که بخش های مختلف شناسایی و تعریف شوند. توصیه می شود توصیف و خلاصه سازی پیش از هر نوع مسئله داده کاوی دیگر انجام شود. این امر باعث شده است در این کتاب، بخش توصیف و خلاصه سازی داده به عنوان یکی از وظایف فاز درک داده مطرح شود.

خلاصه سازی از طرفی نقش مهمی در تهیه و نمایش نتایج نهایی پروژه خواهد داشت. خروجی سایر مسائل داده کاوی (مانند توصیف مفاهیم و مدل های پیش بینی کننده) نیز ممکن است نوعی خلاصه سازی داده تلقی شوند که در سطح مفهومی بالاتری ارائه شده اند. بسیاری از سیستم های گزارش گیری، بسته های آماری و سیستم های EIS^{ot} و OLAP نیز می توانند عملیات توصیف و خلاصه سازی داده را انجام دهند. اما هیچ روشی برای انجام مدلسازی های پیشرفته فراهم نمی کنند. بنابراین اگر توصیف و خلاصه سازی به تنهایی به عنوان نوع مسئله تلقی شود و هیچ نیازی به مدلسازی نباشد این ابزار برای انجام مسئله داده کاوی مناسب خواهد بود.

۲-۲ خوشه بندی

خوشه بندی نوعی مسئله داده کاوی است که داده ها را به خوشه ها یا زیر گروه های با معنی تقسیم می کند. تمام اعضای یک زیر گروه خصوصیات مشترکی دارند. برای مثال در تحلیل سبد خرید می توان خوشه های سبد را بر اساس کالاهایی که در آن وجود دارد تعریف کرد. (یعنی خصوصیت مشترک اعضای زیر گروه یا خوشه، نوع کالای خریداری شده باشد) خوشه بندی می تواند به صورت دستی یا نیمه- خودکار انجام شود. تحلیلگر می تواند زیر گروه های

خاصی را مرتبط با سوالات کسب و کار فرض کند که در نتیجه دانسته های قبلی او و یا نتایج مرحله توصیف و خلاصه سازی داده است. با این حال تکنیک های خوشه بندی خودکار نیز وجود دارند که می توانند ساختارهایی را در داده ها تشخیص دهند که قبلاً نامشخص و پنهان بوده اند و خوشه بندی را به این صورت انجام دهند.

خوشه بندی می تواند به خودی خود یکی از انواع مسائل داده کاوی باشد که در این صورت تشخیص خوشه ها هدف اصلی داده کاوی خواهد بود. برای مثال ممکن است برای ارسال آگهی های بیمه پرستاری، بخواهیم آدرس ها و کدهای پستی تمام افرادی را که از لحاظ سن و درآمد بالاتر از میانگین جامعه هستند به دست آوریم.

با این وجود معمولاً خوشه بندی یک مرحله موقت در جهت حل سایر مسائل داده کاوی است و در این موارد هدف ممکن است تقسیم مجموعه داده ها به قسمت های قابل مدیریت و کنترل پذیر بوده و یا یافتن زیرمجموعه های همگن تری از داده ها برای ساده نمودن تحلیل باشد. معمولاً در مجموعه داده های بزرگ تأثیر انواع داده های مختلف بر روی هم باعث می شود تا در نهایت الگوهای موجود و مطلوب در داده ها مبهم و غیر قابل تشخیص شوند. بنابراین خوشه بندی مناسب می تواند به آسانتر شدن مراحل دیگر داده کاوی کمک کند. برای مثال تحلیل وابستگی های بین کالاهای موجود در میلیونها سبد خرید بسیار مشکل است در حالی که تشخیص وابستگی ها در خوشه های مطلوب سبد خرید بسیار آسانتر و منطقی تر خواهد بود. تعیین خوشه های مرتبط با سبدهای گران قیمت، سبدهای لوازم دم دستی و راحتی و یا سبدهای مربوط به روز یا یک زمان خاص.

تذکر: گاهی میان مفهوم واژه خوشه بندی و سایر مفاهیم داده کاوی مانند طبقه بندی، تداخل به وجود می آید. این تداخل به این دلیل است که برخی افراد از این کلمه برای اشاره به ایجاد کلاس ها و زیر گروه ها استفاده می کنند. منظور برخی دیگر ایجاد مدل هایی برای پیش بینی کلاس های از پیش شناخته شده، روی داده های جدید است. اما در این کتاب واژه طبقه بندی را به مفهوم دوم اختصاص داده و کلمه خوشه بندی را برای اشاره به مفهوم اول استفاده می کنیم. با این وجود تکنیک های طبقه بندی می توانند برای فهمیدن توصیف خوشه های یافت شده استفاده شوند.

تکنیک های مناسب عبارتند از:

- تکنیک های خوشه بندی
- شبکه های عصبی
- مجسم سازی

مثال:

یک شرکت اتومبیل سازی معمولاً اطلاعات مربوط به مشتریان را در رابطه با ویژگی های اجتماعی و اقتصادی آن ها جمع آوری می کند. خصوصیات مانند میزان درآمد، جنسیت، سن، شغل و... با استفاده از تحلیل خوشه ای، شرکت می تواند مشتری های خود را به زیر گروه های قابل فهم تر و معنی دارتر تقسیم کند و ساختار هر گروه را به صورت جداگانه تحلیل کند و استراتژی های بازاریابی خاصی برای هر یک از زیرگروه ها اعمال کند.

۳-۲ توصیف مفاهیم

هدف از توصیف مفاهیم، ارائه توصیفی قابل فهم از مفاهیم یا کلاس هاست. هدف، ساخت مدل های کامل با دقت پیش بینی بالا نیست بلکه هدف کسب بینش و درک بهتر است. برای مثال ممکن است شرکتی بخواهد در مورد مشتری های وفادار و مشتری های گذرای خود بیشتر بداند و بر اساس توصیف مفاهیم «وفادار» و «گذرا» نتیجه بگیرد چه کارهایی می تواند انجام دهد تا مشتریان وفادار، بمانند، یا مشتریان گذرا، وفادار شوند.

توصیف مفاهیم هم با خوشه بندی و هم با طبقه بندی رابطه نزدیکی دارد. خوشه بندی ممکن است بدون این که هیچ توصیف قابل فهمی در اختیار قرار دهد به صورت بندی اشیاء متعلق به یک کلاس منجر شود. اصولاً خوشه بندی قبل از توصیف مفاهیم انجام می شود اما برخی تکنیک ها مانند خوشه بندی مفهومی، خوشه بندی و توصیف مفاهیم را به طور همزمان انجام می دهند.

توصیف مفاهیم برای اهداف طبقه بندی نیز می تواند استفاده شود. از سوی دیگر برخی از تکنیک های طبقه بندی مدل های قابل فهمی ایجاد می کنند که می تواند به عنوان توصیف مفاهیم تلقی شود. تفاوت مهم این دو این است که طبقه بندی باید کامل باشد زیرا نتایج آن باید روی تمام اعضای جامعه مورد مطالعه اعمال شود. اما از سوی دیگر، توصیف مفاهیم نیازی به کامل بودن ندارد و فقط کفایت قسمت های مهم مفاهیم و یا کلاس ها را شرح دهد. در مثال بالا توصیف مفهوم مشتریان وفادار، با توجه به ویژگی های مشتریانی که وفاداری آن ها واضح است، کافی خواهد بود.

تکنیک های مناسب عبارتند از:

- روش های استنتاج قانونی
- خوشه بندی مفهومی

مثال:

یک شرکت می تواند با استفاده از داده های مربوط به خریداران اتومبیل های جدید و تکنیک استنتاج قانون، قوانینی ایجاد کند که مشتریان وفادار و گذرایش را توصیف کند. قوانین زیر نمونه هایی از قوانین تدوین شده شرکت مذکورند:

- اگر جنس: مرد ، سن < ۵۱ ← مشتری: وفادار
 اگر جنس: زن ، سن < ۲۱ ← مشتری: وفادار
 اگر شغل: مدیر ، سن > ۵۱ ← مشتری: گذرا
 اگر وضعیت تأهل: مجرد ، سن > ۵۱ ← مشتری: گذرا

۴-۲ طبقه بندی

طبقه بندی برای فرض استوار است که مجموعه ای از اشیاء به واسطه چندین ویژگی از یکدیگر متمایز می‌شوند و متعلق به کلاس‌های متفاوتی هستند. برچسب کلاس یک مقدار گسسته (اسمی) است و برای هر شیء مقدار مشخصی است. هدف از ساخت مدل‌های طبقه بندی کننده آن است که مدل بتواند به اشیائی که قبلاً دیده و برچسب گذاری نشده اند، برچسب صحیحی نسبت دهد. مدل‌های طبقه بندی کننده غالباً برای مدلسازی پیشبینانه استفاده می‌شوند.

برچسب کلاس می‌تواند به وسیله کاربر یا با استفاده از خوشه بندی تعیین گردد. طبقه بندی یکی از مهم‌ترین انواع مسائل داده کاوی است که گستره کاربرد وسیعی دارد. بسیاری از مسائل داده کاوی قابل تبدیل به مسائل طبقه بندی هستند. به عنوان مثال، تعیین اعتبار سعی می‌کند تا ریسک اعتباردهی به یک مشتری جدید را ارزیابی کند. این موضوع می‌تواند به یک مسئله طبقه بندی تبدیل شود به این صورت که دو کلاس برای مشتریان در نظر گرفته می‌شود. مشتریان قابل اعتماد و غیر قابل اعتماد. یک مدل طبقه بندی کننده می‌تواند با توجه به اطلاعات موجود مشتری و رفتار اعتباری او ساخته شود. سپس مدل ساخته شده می‌تواند برای طبقه بندی مشتریان جدید به یکی از دو کلاس فوق استفاده شود و بر اساس آن، مشتری مورد نظر پذیرفته یا رد گردد.

طبقه بندی تقریباً با همه انواع مسائل داده کاوی در ارتباط است. مسائل پیش بینی با گسسته سازی مقادیر پیوسته فیلد هدف به مسائل طبقه بندی تبدیل می‌شوند. «گسسته سازی» امکان تبدیل طیف‌های پیوسته به بازه‌های گسسته را فراهم می‌کند. این مقادیر گسسته سپس به جای مقادیر عددی دقیق، به عنوان برچسب کلاس‌ها استفاده می‌شود و مسئله را به یک مسئله طبقه بندی تبدیل می‌کند. برخی از تکنیک‌های طبقه بندی، توصیفی از مفاهیم و یا کلاس‌های قابل فهم ایجاد می‌کنند. طبقه بندی با تحلیل وابستگی‌ها نیز مرتبط است زیرا مدل‌های طبقه بندی کننده معمولاً وابستگی‌های بین ویژگی‌ها را استخراج می‌کنند.

خوشه بندی می‌تواند برای به دست آوردن برچسب اولیه کلاس‌ها مورد استفاده قرار گیرد و یا مجموعه داده‌ها را طوری محدود می‌کند که مدل طبقه بندی کننده بتواند به طور بهتر و کارتری ساخته شود. تجزیه و تحلیل انحرافات قبل از ساخت مدل طبقه بندی کننده، بسیار مفید است زیرا وجود

انحرافات و مقادیر پرت ممکن است باعث شود که الگوهایی که در مجموعه داده وجود دارد مبهم و غیر قابل تشخیص شوند و در نتیجه به دست آوردن یک مدل طبقه بندی کننده کارا، ممکن نباشد. از سوی دیگر یک مدل طبقه بندی کننده می تواند برای تشخیص داده های پرت، انحرافات و سایر مشکلات موجود در داده ها نیز استفاده شود.

تکنیک های مناسب عبارتند از:

- تحلیل ممیزی
- روش استنتاج قانون
- یادگیری درخت تصمیم
- شبکه های عصبی
- k -KNN - نزدیک ترین همسایه
- استدلال مبتنی بر موارد
- الگوریتم ژنتیک

مثال:

بانک ها عموماً اطلاعاتی در زمینه رفتار اعتباری متقاضیان وام خود دارند. اگر این اطلاعات مالی را با سایر اطلاعات شخصی مانند جنس، سن و درآمد ترکیب کنیم، آنگاه امکان توسعه سیستم برای طبقه بندی مشتریان جدید به کلاس های قابل اعتماد و غیر قابل اعتماد به وجود خواهد آمد. (این طبقه بندی مشخص کننده میزان ریسک اعتباری در پذیرفتن یک مشتری خواهد بود)

۲-۵ پیش بینی

مسئله مهم دیگری که کاربرد وسیعی دارد پیش بینی است. پیش بینی شباهت زیادی با طبقه بندی دارد. تنها تفاوت آن در این است که در پیش بینی، متغیر هدف (کلاس) یک ویژگی کیفی گسسته نیست، بلکه مقداری پیوسته است. هدف پیش بینی یافتن مقادیر عددی برای نمونه های مشاهده نشده است. این مسئله گاهی رگرسیون نامیده می شود و در صورتی که پیش بینی با داده های سری زمانی سروکار داشته باشد، پیش گویی^{۵۳} نامیده می شود.

تکنیک های مناسب عبارتند از:

- تحلیل رگرسیون

- درخت های رگرسیون
- شبکه های عصبی
- K-NN (k- نزدیک ترین همسایه)
- روش های Box-Jenkins
- الگوریتم ژنتیک

مثال:

سود سالانه یک شرکت بین المللی با ویژگی هایی مانند تبلیغات، نرخ تبدیل، نرخ تورم و... مرتبط است. در این صورت شرکت با در دست داشتن این مقادیر (یا تخمین قابل اعتمادی از آن ها) برای سال بعد می تواند سود مورد انتظار خود را در سال آینده پیش بینی کند.

۶-۲ تحلیل وابستگی

تحلیل وابستگی عبارت است از یافتن مدلی که وابستگی بین داده ها یا رخداد ها را شرح دهد. یکی از کاربرد های تحلیل وابستگی، پیش بینی مقدار یک عنصر داده ای، با در اختیار داشتن سایر عناصر است. اگر چه وابستگی ها در مدل سازی پیشبینانه هم قابل استفاده اند اما آن ها اغلب برای توصیف روابط در بین داده ها استفاده می شوند. وابستگی ها می توانند اکید یا احتمالی باشند.

«روابط انجمنی» حالت خاصی از وابستگی است که اخیراً بسیار رایج شده است. روابط انجمنی، رابطه های موجود بین داده ها را شرح می دهند. (یعنی داده ها یا رخدادهایی که مکرراً با هم رخ می دهند را کشف می کنند). یکی از کاربردهای رایج روابط انجمنی، تحلیل سبدهای خرید است. یک مثال رایج برای روابط انجمنی از این قرار است: «در ۳۰ درصد از کل خریدها نوشیدنی و بادام زمینی با هم خریداری می شوند».

الگوریتم هایی که روابط انجمنی را تشخیص می دهند بسیار سریع هستند و تعداد زیادی از این نوع رابطه ها ایجاد می کنند. اما نکته مهم، یافتن مطلوب ترین آن ها است.

تحلیل وابستگی رابطه نزدیکی با پیش بینی و طبقه بندی دارد. مثلاً برای ساخت مدل های پیشبینانه، تلویحاً از تحلیل وابستگی ها استفاده می شود. همچنین این تحلیل وقتی که ارتباطات خاصی را برجسته می سازد با مسئله توصیف مفاهیم نیز در ارتباط است.

در کاربردها تحلیل وابستگی معمولاً همراه با خوشه بندی انجام می شود. در مجموعه داده های بزرگ به علت تاثیر متقابل داده ها بر روی هم وابستگی ها به ندرت مشخص و برجسته می شوند. در این مورد پیشنهاد می شود که تحلیل وابستگی به جای کل مجموعه داده روی خوشه های همگن انجام شود.

«تحلیل توالی» حالت خاصی از تحلیل وابستگی می باشد که تمرکز آن بر روی ترتیب زمانی رخدادها هست. بطور مثال در تحلیل سبد خرید، روابط انجمنی، وابستگی بین اقلامی که در یک زمان باهم خرید می شوند را مشخص می نماید. اما تحلیل توالی، به توصیف الگوی خرید یک مشتری خاص یا گروهی از مشتریان در طول زمان می پردازد.

تکنیک های مناسب عبارتند از:

- تحلیل همبستگی
- تحلیل رگرسیون
- قوانین انجمنی
- شبکه های بیزین
- برنامه نویسی منطقی استنتاجی
- تکنیک های مجسم سازی

مثال ۱:

با استفاده از تحلیل رگرسیونی، یک تحلیل گر کسب و کار در می یابد که وابستگی های مهمی بین فروش کل یک محصول، قیمت و میزان تبلیغات آن وجود دارد. به محض به دست آوردن این اطلاعات او می تواند با تغییر دادن قیمت یا تبلیغات، فروش خود را به میزان مناسبی برساند.

مثال ۲:

یک شرکت سازنده اتومبیل با اعمال الگوریتم های قوانین انجمنی روی داده های مربوط به لوازم جانبی اتومبیل، دریافته است که در ۹۵ درصد از موارد، در صورت سفارش سیستم صوتی پیشرفته، جعبه دنده اتوماتیک نیز سفارش داده می شود. بر اساس دریافت این وابستگی، شرکت سازنده اتومبیل تصمیم می گیرد که این محصولات را با هم و به صورت یک مجموعه ارائه دهد و در نتیجه هزینه های مربوط را به نحو قابل توجهی کاهش دهد.



CRISP-DM

Step-by-Step Data Mining Guide

«در مدت فعالیت من، بعنوان پژوهشگر ارشد و رئیس دپارتمان داده‌کاوی در شرکت بنز آلمان، این دپارتمان در چند پروژه داده‌کاوی نیز شرکت داشت که توسط بازار مشترک اروپا پشتیبانی مالی می‌شد. CRISP-DM بنظر من مهمترین آنها بود. امروزه که بیش از چهارده سال از پایان این پروژه میگذرد، بدون انحراف می‌توان گفت که CRISP-DM مهمترین استاندارد موجود جهت انجام پروژه‌های کاربردی داده‌کاوی در جهان می‌باشد.»*

غلامرضا نخعی زاده - استاد دانشگاه کلسروهه آلمان و مدیر سابق دپارتمان داده‌کاوی شرکت بنز

