



دانشگاه صنعتی امیرکبیر
(پلی‌تکنیک تهران)
دانشکده ریاضی و علوم کامپیوتر

گزارش پروژه درس داده کاوی محاسباتی پروژه 2

استخراج ویژگی مبتنی بر استقلال خطی و تحلیل تاثیر آن بر فرایند بهینه
سازی توابع زیان

نگارش
هومن ذوالفقاری

استاد راهنما
دکتر مهدی قطعی

تدریس یار
مهندس بهنام یوسفی مهر

آبان 1404

صفحه فرم ارزیابی و تصویب پایان نامه- فرم تأیید اعضای کمیته دفاع

در این صفحه (هر سه مقطع تحصیلی) باید تصویر فرم ارزیابی یا تأیید و تصویب پایان نامه/رساله موسوم به فرم کمیته دفاع برای مقاطع کارشناسی ارشد و دکتری و تصویر فرم تصویب برای مقطع کارشناسی، موجود در **پرونده آموزشی** را قرار دهند.

نکات مهم:

- ✓ نگارش پایان نامه/رساله باید به **زبان فارسی** و بر اساس آخرین نسخه دستورالعمل و راهنمای تدوین پایان نامه های دانشگاه صنعتی امیرکبیر باشد. (دستورالعمل و راهنمای حاضر)؛
- ✓ تحویل پایان نامه به زبان انگلیسی، برای دانشجویان بین الملل با شرایط دستورالعمل حاضر بلامانع است و داشتن صفحه عنوان فارسی به همراه چکیده مبسوط فارسی، 30 صفحه برای پایان نامه کارشناسی ارشد و 50 صفحه برای رساله دکتری در ابتدای آن الزامی است؛
- ✓ دریافت پایان نامه کارشناسی، کارشناسی ارشد و رساله دکتری، **بصورت نسخه دیجیتال** مطابق راهنمای وبسایت و دستورالعمل حاضر می باشد؛
- ✓ در صورتی که يك عنوان پایان نامه دارای **دو نویسنده** است، فقط یکبار فایل و فرم اطلاعات آن با ذکر هر دو نویسنده بارگذاری و تکمیل گردد؛
- ✓ با توجه به اینکه در ورود 2016 یا بالاتر، احتمال تغییر ترتیب ذکر زیر فصل ها وجود دارد لطفا در انتها به شماره دهی زیر فصل ها توجه نمایید که بصورت صحیح باشد. از راست به چپ: شماره فصل-زیرفصل 1-زیرفصل 2-زیرفصل 3 و



به نام خدا

تاریخ:

تعهدنامه اصالت اثر

اینجانب هومن ذوالفقاری متعهد می‌شوم که مطالب مندرج در این پایان نامه حاصل کار پژوهشی اینجانب تحت نظارت و راهنمایی اساتید دانشگاه صنعتی امیرکبیر بوده و به دستاوردهای دیگران که در این پژوهش از آنها استفاده شده است مطابق مقررات و روال متعارف ارجاع و در فهرست منابع و مآخذ ذکر گردیده است. این پایان نامه قبلاً برای احراز هیچ مدرک هم‌سطح یا بالاتر ارائه نگردیده است.

در صورت اثبات تخلف در هر زمان، مدرک تحصیلی صادر شده توسط دانشگاه از درجه اعتبار ساقط بوده و دانشگاه حق پیگیری قانونی خواهد داشت.

کلیه نتایج و حقوق حاصل از این پایان نامه متعلق به دانشگاه صنعتی امیرکبیر می‌باشد. هرگونه استفاده از نتایج علمی و عملی، واگذاری اطلاعات به دیگران یا چاپ و تکثیر، نسخه‌برداری، ترجمه و اقتباس از این پایان نامه بدون موافقت کتبی دانشگاه صنعتی امیرکبیر ممنوع است.

نقل مطالب با ذکر مآخذ بلامانع است.

در صفحه تعهدنامه اصالت اثر، در قسمت بالا سمت چپ، تاریخ دفاع خود را جایگزین تاریخ نوشته شده کنید.

همچنین در صفحه تعهدنامه اصالت اثر، در خط اول، نام و نام خانوادگی خود را به صورت کامل با نام و نام خانوادگی نمونه، جایگزین کنید. در انتهای متن تعهد، در قسمت امضا نیز باید نام و نام خانوادگی کامل خود را وارد نمایید.

هومن ذوالفقاری

امضا

چکیده

چکیده باید جامع و بیان‌کننده خلاصه‌ای از اقدامات انجام‌شده باشد. در قسمت چکیده، چکیده پایان‌نامه خود را که حداکثر می‌تواند شامل 250 کلمه باشد، بنویسید. در آخر چکیده و در قسمت واژگان کلیدی، کلمات کلیدی خود را وارد کنید. کلمات کلیدی بین 3 تا 5 کلمه می‌تواند باشد که طبق فرمت باید با ویرگول از هم جدا شوند.

واژه‌های کلیدی:

داده کاوی محاسباتی، کاهش ابعاد، انتخاب فیچر

چکیده ا

فصل اول مقدمه (دستور العمل) 1

فصل دوم مشخصات یک پایان نامه و گزارش علمی 9

2-1- برخورداری از غنای علمی Error! Bookmark not defined.

2-2- ارجاع به موقع و صحیح به منابع دیگر Error! Bookmark not defined.

2-3- ساده نویسی Error! Bookmark not defined.

2-3-1- وحدت موضوع Error! Bookmark not defined.

2-3-2- اختصار Error! Bookmark not defined.

2-3-3- رعایت نکات دستوری و نشانه گذاری Error! Bookmark not defined.

2-3-4- توجه به معلومات ذهنی مخاطب Error! Bookmark not defined.

2-3-5- رعایت مراحل اصولی نگارش Error! Bookmark not defined.

فصل سوم نگارش صحیح 12

3-1- فارسی نویسی Error! Bookmark not defined.

3-2- رعایت املاي صحیح فارسی Error! Bookmark not defined.

3-3- رعایت قواعد نشانه گذاری Error! Bookmark not defined.

3-3-1- ویرگول و نقطه Error! Bookmark not defined.

3-3-2- دو نقطه Error! Bookmark not defined.

3-3-3- گیومه Error! Bookmark not defined.

3-3-4- نشانه پرششی Error! Bookmark not defined.

3-3-5- خط تیره Error! Bookmark not defined.

3-3-6- پرانتز Error! Bookmark not defined.

فصل چهارم سبک ها و قلم ها 16

4-1- قلم های فارسی Error! Bookmark not defined.

4-2- قلم های انگلیسی Error! Bookmark not defined.

4-3- فرمول ها (روابط ریاضی) Error! Bookmark not defined.

4-4- فاصله های افقی و عمودی Error! Bookmark not defined.

4-4-1- فاصله کلی از چهار طرف کاغذ Error! Bookmark not defined.

4-4-2- فاصله خطها Error! Bookmark not defined.

4-4-3- فاصله های تفکیک کننده Error! Bookmark not defined.

4-5- فواصل بین کلمات Error! Bookmark not defined.

4-6- جدانشتن کلمات بدون گذاشتن فاصله بین آنها Error! Bookmark not defined.

4-7- فهرست گزارش، فهرست شکل ها و فهرست جداول Error! Bookmark not defined.

4-8- سربرگ و تهبرگ (Header and Footer) Error! Bookmark not defined.

4-9- جداول، منحنی ها، شکل ها Error! Bookmark not defined.

4-10- ارجاع به جداول، شکل ها، روابط، مراجع و بخش ها Error! Bookmark not defined.

فصل پنجم بررسی ساختار پایان نامه 20

5-1- بررسی سرفصل ها Error! Bookmark not defined.

5-2- بررسی ساختار کلی Error! Bookmark not defined.

5-3- بررسی مفهومی Error! Bookmark not defined.

5-4- مطالعه مفهومی و جمله بندی Error! Bookmark not defined.

5-5- تنظیم بندها	Error! Bookmark not defined.
5-6- بررسی قواعد نگارشی	Error! Bookmark not defined.
5-7- بررسی روابط	Error! Bookmark not defined.
5-8- بررسی شکل‌ها	Error! Bookmark not defined.
5-8-1- بررسی کیفیت شکل و تطابق عنوان آن	Error! Bookmark not defined.
5-8-2- بررسی تطابق روابط، برنامه و شکل	Error! Bookmark not defined.
5-9- بررسی جداول	Error! Bookmark not defined.
5-9-1- بررسی کیفیت جدول و تطابق عنوان آن	Error! Bookmark not defined.
5-9-2- بررسی تطابق روابط، برنامه و جدول	Error! Bookmark not defined.
5-10- به‌روزرسانی مراجع	Error! Bookmark not defined.
5-11- صفحه‌بندی	Error! Bookmark not defined.
5-12- سربرگ و ته‌برگ‌ها	Error! Bookmark not defined.
فصل ششم نتیجه‌گیری	33
منابع و مراجع	36
پیوست‌ها	Error! Bookmark not defined.
Abstract	37

صفحه

فهرست اشکال

شکل 4-1- فرایند کواکستروژن Error! Bookmark not defined.

صفحه

فهرست جداول

جدول 4-1 - قلم‌هاي فارسي.....	Error! Bookmark not defined.
جدول 4-2 - قلم‌هاي انگليسي.....	Error! Bookmark not defined.
جدول 4-3 - قلم و سبك فرمول‌ها.....	Error! Bookmark not defined.
جدول 4-4 - اندازه فرمول‌ها.....	Error! Bookmark not defined.
جدول 4-5 - عنوان جدول.....	Error! Bookmark not defined.

فهرست علائم

علائم لاتین

ارتفاع	h
طول موج توربولانس	L
پریود توربولانس	T
سرعت تعادل وسیله پرنده	U_0
مولفه سرعت تندباد در راستای محور طولی دستگاه مختصات بدنی نسبت به اینرسی	u_g^B

علائم یونانی

چگالی طیفی قدرت توربولانس	$\Phi(\omega)$
شدت توربولانس	σ
بسامد توربولانس	ω
بسامد فاصله‌ای	Ω

بالانویس‌ها

دستگاه مختصات بدنی	B
--------------------	-----

زیرنویس‌ها

تندباد (گاست)	g
---------------	-----

فصل اول

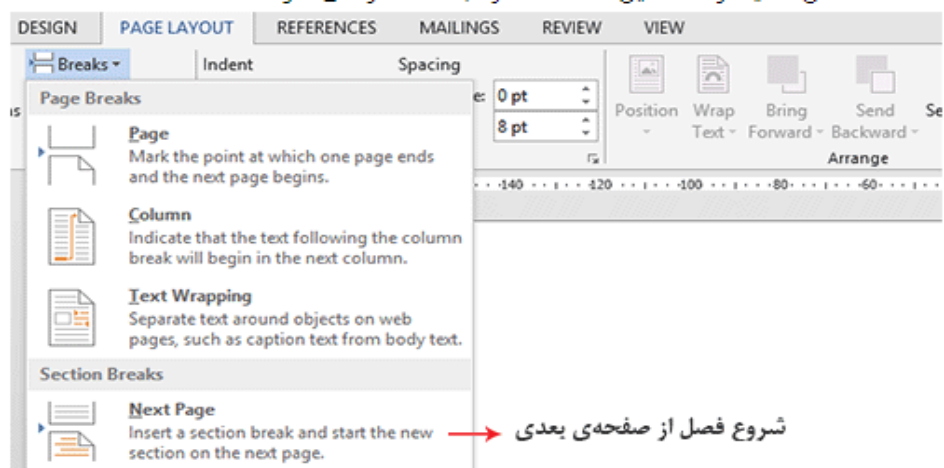
مقدمه (دستور العمل)

مقدمه

1- پایان نامه نمونه برای 5 فصل طراحی شده است، چنانچه تعداد فصل های پایان نامه شما، کمتر از پنج فصل است، فصول اضافه را پاک کنید.

2- اگر تعداد فصل ها، بیشتر از پنج فصل باشد، برای اضافه کردن یک فصل جدید، باید قسمت (Section) جدیدی ایجاد کنید. برای ایجاد یک قسمت جدید با تنظیمات متفاوت نسبت به قسمت قبل در یک فایل ورد، کفایت مکان نما را در جایی که باید قسمت جدید آغاز شود (مثلا پس از اتمام فصل 5) قرار دهید. سپس از تب Page Layout و از گروه Page Setup، روی دکمه بازشونده Breaks کلیک کنید. همان طور که می بینید، چهار نوع Section Break وجود دارد:

- Next Page : فصل جدید از صفحه بعد شروع می شود.
- Continuous : بدون شکست صفحه، فصل جدید در ادامه ی فصل قبلی شروع می شود.
- Even Page : فصل جدید از نخستین صفحه ی زوج بعدی آغاز می شود.
- Odd Page : فصل جدید از نخستین صفحه ی فرد بعدی آغاز می شود.



شکل 1: نحوه ساخت قسمت جدید

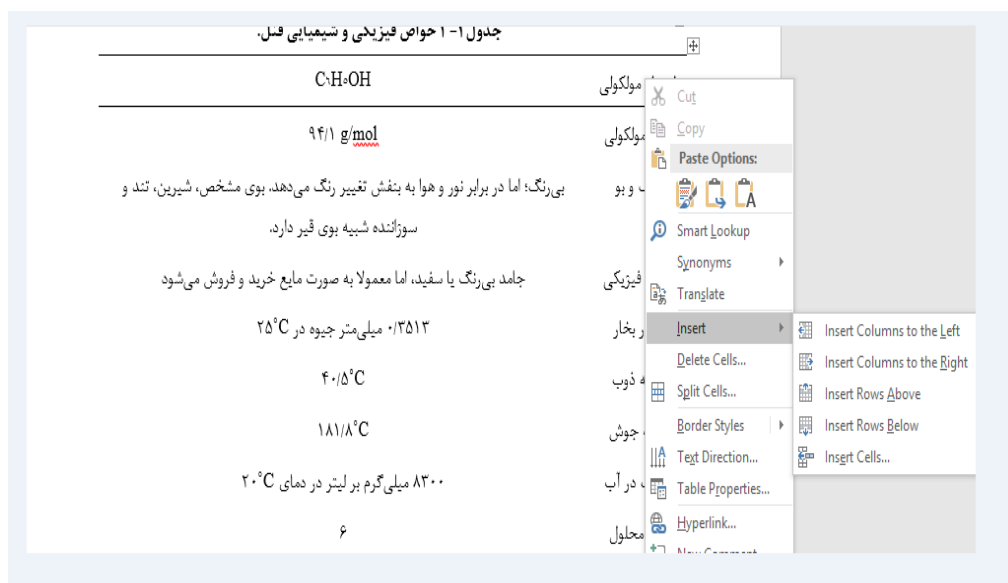
مطابق شکل (1) گزینه Next Page را انتخاب کنید. در ادامه برای تنظیم سربرگ فصل جدید، در قسمت سربرگ صفحه دوبار کلیک کنید تا قابلیت ویرایش آن فعال شود و بخش Header and Footer Tools به تب ها اضافه شود. در تب Design و در گروه Navigation، دکمه Link to

previous را پیدا کنید. همان طور که می بینید، این دکمه به طور پیش فرض روشن است. روی آن کلیک کنید تا دکمه خاموش شود و ارتباط این فصل با فصل پیشین قطع شود. اکنون اگر سربرگ قسمت فعلی را ویرایش کنید، سربرگ قسمت قبلی تغییر نخواهد کرد. برای اینکار باید گزینه Link to Previous را غیرفعال کرد.

برای سرفصل ها مطابق جدول « جدول 4-1 قلم های فارسی » اقدام کنید و یکی از سرفصل مورد نظر را در این قالب کپی کنید و بعد جایگزین (paste) کنید، چنانچه بعد از جایگزین کردن شماره سرفصل به هر دلیلی دچار بهم ریختگی شد از راست به چپ، عدد چهارم را تغییر داده و سپس متن آن را ویرایش کنید.

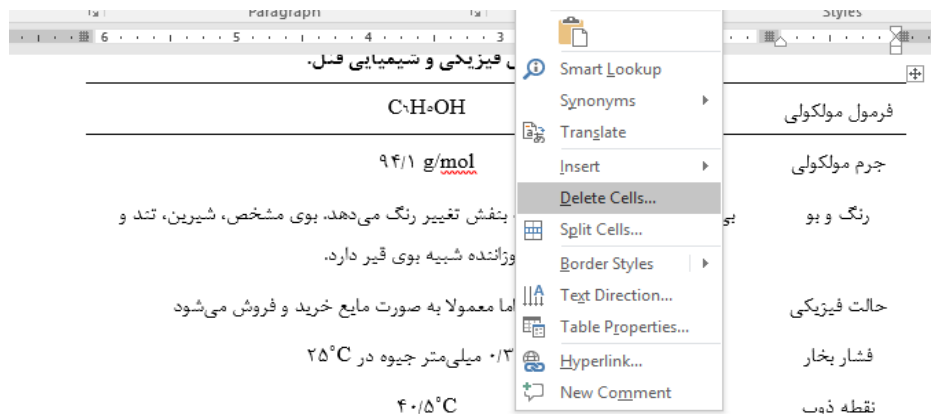
- در هر قسمتی از پایان نامه که نیاز به جدول دارید، مراحل زیر را انجام دهید:

- 1- یکی از جدول های موجود در متن نمونه همراه با عنوان آن به طور کامل انتخاب کنید.
- 2- با کلید ترکیبی $ctrl+c$ آن جدول را کپی کنید.
- 3- سپس به محل مورد نظر رفته و در آنجا با کلید ترکیبی $ctrl+v$ جدول را جایگزین (paste) کنید.
- 4- داده های مورد نظر خود را در جدول وارد کنید.
- 5- اگر در جدول انتخابی نیاز به اضافه کردن سطر یا ستون است، بر روی جدول کلیک راست کرده و از گزینه Insert، عملیات مورد نظر را انتخاب کنید (شکل 2).



شکل 2: اضافه کردن سطر یا ستون در جدول

6- اگر تعداد سطر و ستون جدول شما، کمتر از جدول انتخابی است، در این حالت نیاز به حذف تعدادی سطر یا ستون دارید. برای اینکار بر روی جدول کلیک راست کرده و با انتخاب گزینه Delete cells...، پنجره زیر باز می‌شود (شکل 3)، برای حذف سطر اضافی گزینه سوم و برای حذف ستون اضافی گزینه چهارم را انتخاب کنید.



شکل 3: حذف سطر یا ستون در جدول

اگر به هر دلیلی دچار بهم ریختگی شدید مطابق « جدول 4-2 قلم‌های فارسی» از سبک (Style) مورد اشاره استفاده و **در آخر هم برای تهیه فهرست جداول، از استایل TABLE TITLE فهرست گیری کنید.**

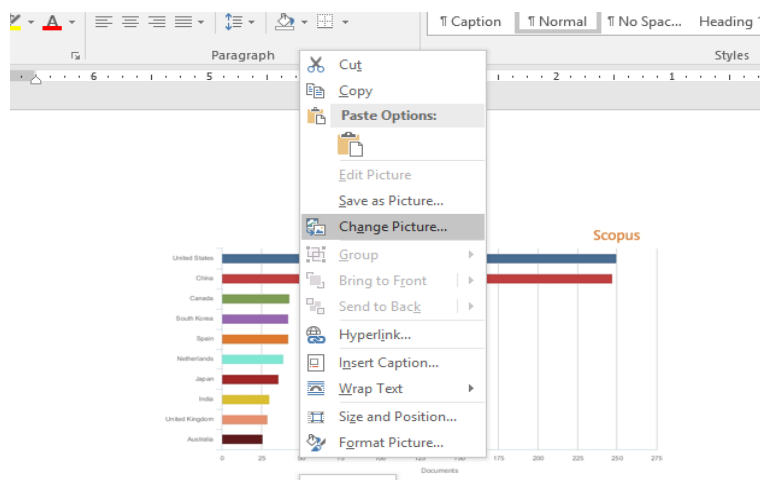
- برای قرار دادن عکس در پایان نامه، باید مراحل زیر را به ترتیب انجام دهید:

1- یکی از عکس‌های موجود در متن پایان‌نامه نمونه همراه با عنوان آن به طور کامل انتخاب کنید.

2- با کلید ترکیبی $\text{ctrl}+\text{c}$ آن عکس را کپی کنید.

3- سپس به محل مورد نظر رفته و در آنجا با کلید ترکیبی $\text{ctrl}+\text{v}$ عکس را جایگزین کنید.

4- روی عکس کلیک راست کرده مطابق شکل 4، گزینه change picture... انتخاب کنید.



شکل 4: تغییر عکس در پایان نامه

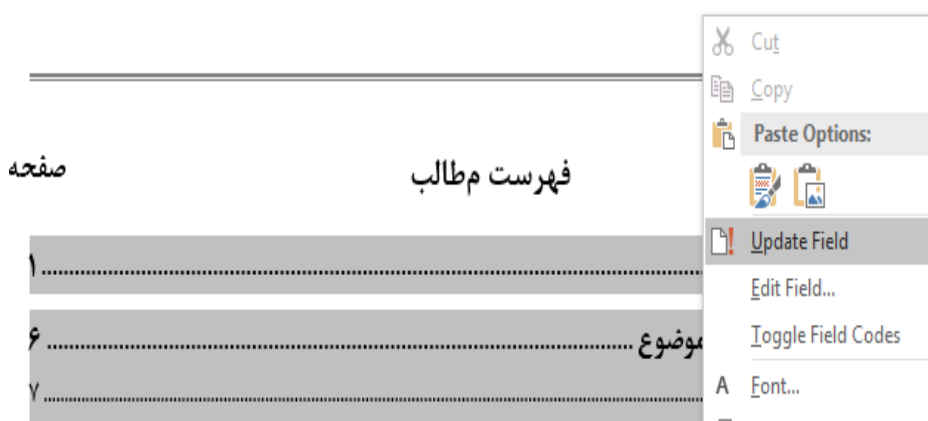
5- با انتخاب گزینه بالا، صفحه زیر باز می‌شود، از قسمت جستجو (Browse)، عکس مورد نظر را انتخاب کرده و جایگزین عکس فعلی می‌کنید.

6- متن مورد نظر خود را جایگزین متن عکس انتخابی کنید. البته نباید در سایز و فونت آن تغییری دهید.

اگر به هر دلیلی دچار بهم ریختگی شدید مطابق « جدول 3-4 قلم‌های فارسی » از سبک مورد اشاره استفاده و **در آخر هم برای تهیه فهرست جداول ، از سبک PIC TITLE فهرست‌گیری کنید.**

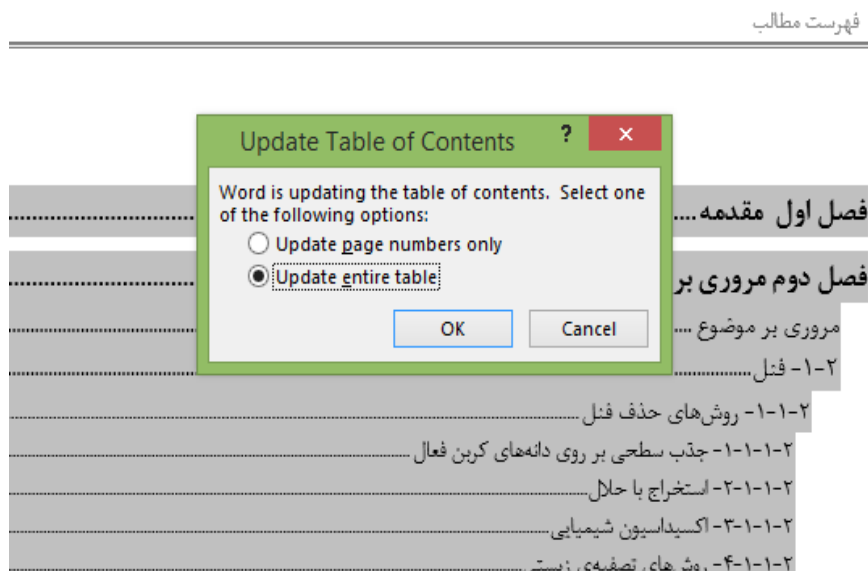
بعد از وارد و جایگزین کردن کامل مطالب خود در پایان نامه نمونه، به قسمت فهرست مطالب برگردید.

18- برای بروزرسانی فهرست‌ها (مطالب، اشکال و جداول)، بر روی جدول فهرست مورد نظر کلیک راست کرده و گزینه Update Field را انتخاب کنید. (شکل 5).



شک 5: انتخاب گزینه بروزرسانی فهرست

19- بعد از انتخاب این گزینه، پنجره زیر باز می‌شود (شکل 6)، برای بروزرسانی کامل فهرست گزینه دوم را انتخاب کنید. فهرست جدید برای پایان‌نامه شما ساخته می‌شود.

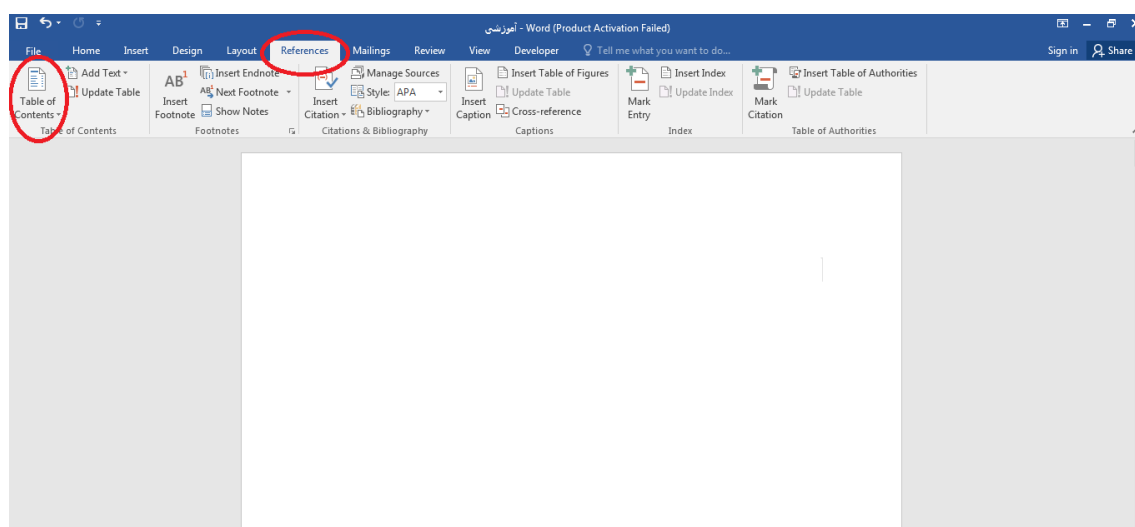


شکل 6: بروز رسانی کامل فهرست

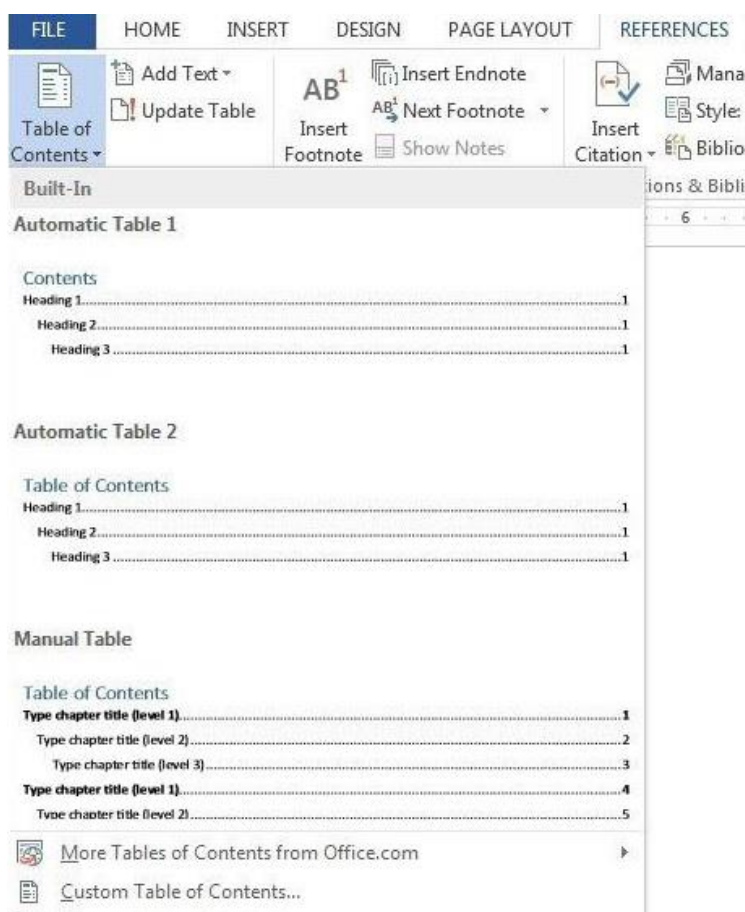
در آخر هم در خصوص فهرست‌گیری در صورتی که با جایگزین کردن فهرست دچار مشکل شدید از **HEADING1** و **HEADING2** و **HEADING3** و **HEADING4** مطابق دستورالعمل زیر فهرست‌گیری نمایید.

نحوه نوشتن فهرست مطالب

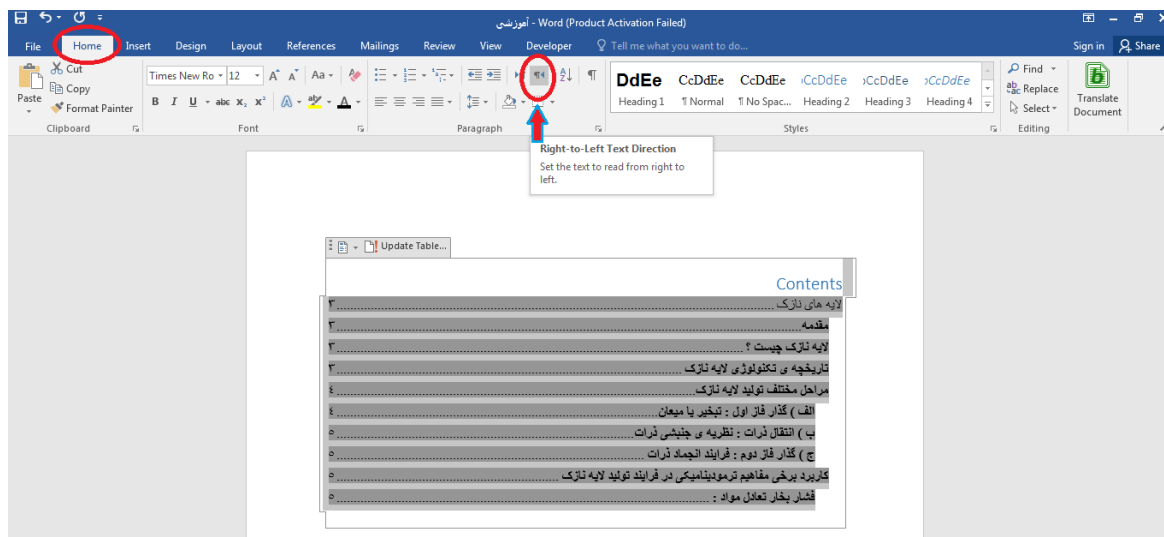
از سربرگ‌های موجود در صفحه‌ی ورد خود وارد سربرگ References شوید و بر روی گزینه‌ی Table of Contents کلیک کنید.



با کلیک بر روی Table of Contents نمونه حالت‌های پیش فرض فهرست‌بندی در ورد به صورت الگو برایتان به نمایش در می‌آید. که با انتخاب هر کدام از آنها فهرست شما به همان شکل به فایل وردتان اضافه می‌شود.



باید توجه داشت که فهرست ایجاد شده برای زبان فارسی مناسب نیست، کافی است تا کل فهرست را انتخاب کرده و سپس در تب Home بر روی گزینه‌ی Right-to-Left کلیک کنید تا فهرست از راست به چپ قرار گیرد.



حالا به ویرایش فهرست می‌رسیم، می‌توانیم عبارت Contents را پاک کرده و به جای آن عبارت “فهرست” را قرار دهیم و با استفاده از قسمت Font نوع فونت و اندازه‌ی فونت متن فهرست و عناوین را به دلخواه تنظیم کنیم.

فصل دوم

داده ها

داده ها

1-2- داده های کلاس بندی

مجموعه داده سرطان پستان ویسکانسین (تشخیصی):

ویژگی های مجموعه داده:

- تعداد نمونه ها: ۵۶۹
 - تعداد ویژگی ها: ۳۰ ویژگی عددی پیش بین و یک ویژگی کلاسی
- توضیحات ویژگی ها:

- شعاع (میانگین فاصله از مرکز تا نقاط روی محیط)
- بافت (انحراف معیار مقادیر خاکستری)
- محیط
- مساحت
- یکنواختی (تغییرات محلی در طول شعاع)
- فشردگی (محیط² / مساحت - ۱/۰)
- فرورفتگی (شدت بخش های فرورفته در مرز)
- نقاط فرورفته (تعداد بخش های فرورفته در مرز)
- تقارن
- بعد فراکتال ("تقریب خط ساحلی" - ۱)

کلاس ها:

- WDBC-Malignant بدخیم)
- WDBC-Benign خوش خیم)

توزیع کلاس ها:

- ۲۱۲ مورد بدخیم
- ۳۵۷ مورد خوش خیم

2-2- داده های رگرسیون

	(20640, 8)	(20640,)						
	MedInc	HouseAge	AveRooms	AveBedrms	Population	AveOccup	Latitude	\
0	8.3252	41.0	6.984127	1.023810	322.0	2.555556	37.88	
1	8.3014	21.0	6.238137	0.971880	2401.0	2.109842	37.86	
2	7.2574	52.0	8.288136	1.073446	496.0	2.802260	37.85	
3	5.6431	52.0	5.817352	1.073059	558.0	2.547945	37.85	
4	3.8462	52.0	6.281853	1.081081	565.0	2.181467	37.85	
	Longitude							
0	-122.23							
1	-122.22							
2	-122.24							
3	-122.25							
4	-122.25							

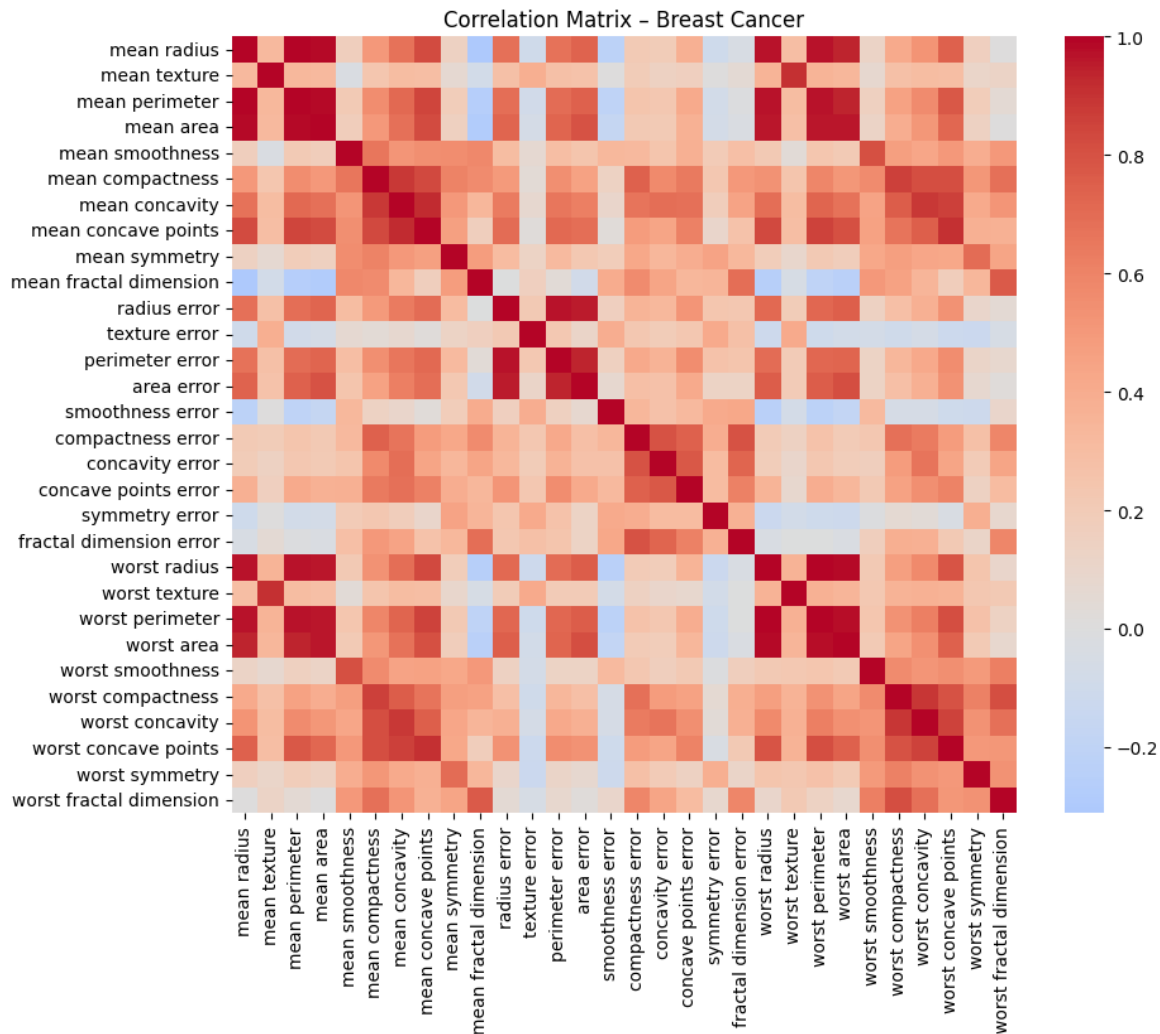
2-3- داده های خوشه بندی

(150, 4)	(150,)
['setosa' 'versicolor' 'virginica']	

فصل سوم بررسی کورلیشن

بررسی کورلیشن

3-1- داده های طبقه بندی



خلاصه ویژگی‌های با همبستگی بالا در مجموعه داده سرطان پستان (ضریب همبستگی > 0.75)

- ویژگی‌های شعاع، محیط و مساحت (در حالت میانگین و بدترین مقدار) تقریباً همپوشانی کامل دارند؛ ضریب همبستگی بین آن‌ها تا ۱.۰۰ می‌رسد.
- نقاط فرورفته (concave points)، فرورفتگی (concavity) و فشردگی (compactness) نیز به‌طور قوی با یکدیگر مرتبطاند (تا حدود ۰.۹۲).
- نسخه‌های میانگین و بدترین هر ویژگی (مثل بافت، یکنواختی، بعد فراکتال و...) معمولاً همبستگی بالایی دارند (بین ۰.۸ تا ۰.۹).

- خطاهای (error) مربوط به شعاع، محیط و مساحت نیز بسیار همبسته‌اند (حدود ۰.۹۵-۰.۹۷).

- در کل، مجموعه داده شامل چند خوشه بزرگ از ویژگی‌های همبسته است که اطلاعات مشابهی ارائه می‌دهند:

1. خوشه اندازه (radius-perimeter-area)

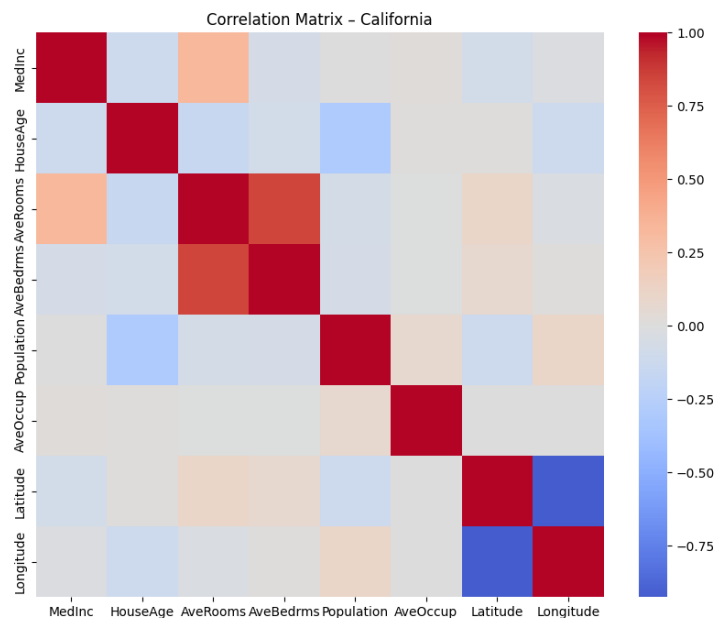
2. خوشه شکل (compactness-concavity-concave points)

3. خوشه خطاها (radius error-perimeter error-area error)

4. خوشه نسخه‌های «میانگین» و «بدترین» هر ویژگی

نتیجه: بسیاری از ویژگی‌ها در این داده همبستگی بسیار بالا دارند، بنابراین در روش‌های کاهش بُعد (مثل PCA یا انتخاب ویژگی، می‌توان برخی از آن‌ها را حذف کرد بدون از دست دادن اطلاعات مهم.

2-3- داده های رگرسیون

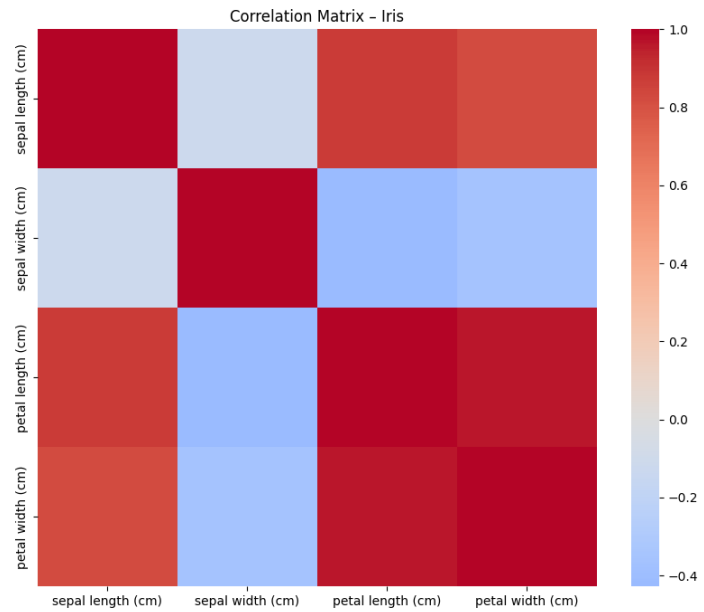


California – Highly Correlated Features (>0.75):

AveRooms ↔ AveBedrms: 0.85

Latitude ↔ Longitude: 0.92

3-3-1 داده های خوشه بندی



Iris – Highly Correlated Features (>0.75):

sepal length (cm) \leftrightarrow petal length (cm): 0.87

sepal length (cm) \leftrightarrow petal width (cm): 0.82

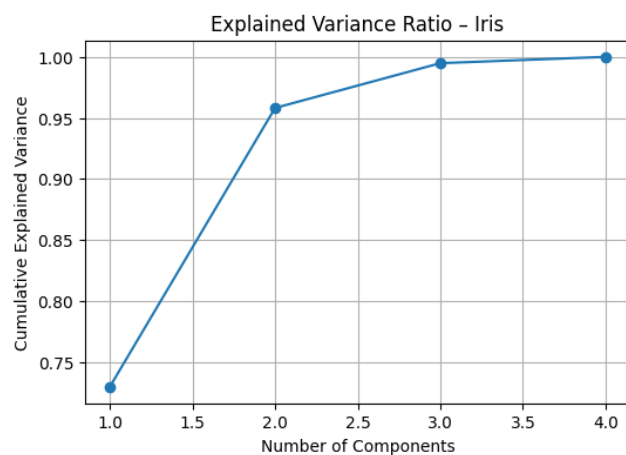
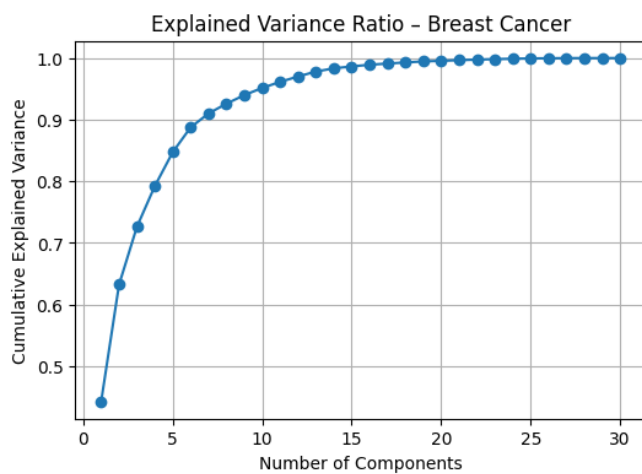
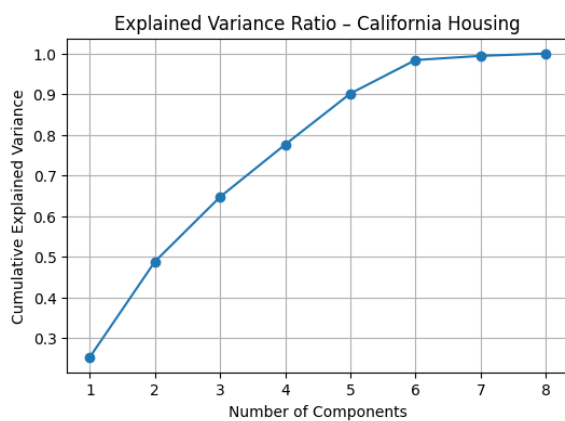
petal length (cm) \leftrightarrow petal width (cm): 0.96.

فصل چهارم كاهش بعد و انتخاب فيچر

کاهش بعد و انتخاب فیچر

در تعریف سبک‌های مختلف این دستورالعمل از قلم‌های Times New Roman و B Nazanin استفاده شده‌است که خصوصیات کامل آنها در بخش‌های بعدی تشریح می‌گردد.

4-1- کاهش بعد



• خلاصه نتایج تحلیل مؤلفه‌ها و خطاها در سه مجموعه داده:

- **Breast Cancer**
- تعداد مؤلفه‌ها برای ۹۵٪ واریانس: ۱۰
- میانگین خطای بازسازی (PCA MSE): ۰/۰۴۸۴

- میانگین قلم در مطلق همبستگی ویژگی ها:

$$PCA = 0.100 \mid ICA = 0.100 \mid SVD = 0.100$$

- توضیح: این مجموعه داده پیچیده تر است و برای پوشش ۹۵٪ واریانس به ۱۰ مؤلفه نیاز دارد.

California Housing

- تعداد مؤلفه ها برای ۹۵٪ واریانس: ۶

$$PCA \text{ MSE: } 0.0159$$

- Mean|Corr|:

$$PCA = 0.167 \mid ICA = 0.167 \mid SVD = 0.167$$

- توضیح: ویژگی ها همبستگی نسبتاً بالاتری دارند اما با مؤلفه های کمتر می توان واریانس را توضیح داد.

Iris

- تعداد مؤلفه ها برای ۹۵٪ واریانس: ۲

$$PCA \text{ MSE: } 0.0419$$

- Mean|Corr|:

$$PCA = 0.500 \mid ICA = 0.500 \mid SVD = 0.500$$

- توضیح: داده آیریس کم بعد است و تنها با دو مؤلفه می توان تقریباً تمام واریانس را پوشش داد. همبستگی ویژگی ها در این مجموعه بالاتر است.

جمع بندی کلی:

- مجموعه Breast Cancer بیشترین تعداد مؤلفه ها را برای حفظ واریانس نیاز دارد.
- در California Housing همبستگی ویژگی ها متوسط است.
- مجموعه Iris با ابعاد کمتر و همبستگی بالا ساده تر فشرده می شود.

4-2- انتخاب فیچر

سرطان پستان (دسته بندی)

SelectKBest: میانگین شعاع، محیط، مساحت، فرورفتگی، نقاط فرورفته، و همچنین بدترین شعاع، محیط، مساحت، فرورفتگی و نقاط فرورفته.

RFE: نقاط فرورفته میانگین، خطای شعاع، خطای مساحت، خطای فشردگی، و همچنین بدترین شعاع، بافت، محیط، مساحت، فروفتگی و نقاط فروفته. توضیح: هر دو روش بیشتر بر ویژگی‌های مربوط به شکل و اندازه توده‌ها تمرکز دارند.

مسکن کالیفرنیا (رگرسیون)

SelectKBest: درآمد متوسط، سن ساختمان، میانگین تعداد اتاق‌ها، میانگین تعداد اتاق‌خواب‌ها، عرض جغرافیایی. RFE: درآمد متوسط، میانگین تعداد اتاق‌ها، میانگین تعداد اتاق‌خواب‌ها، عرض و طول جغرافیایی. توضیح: هر دو روش نشان می‌دهند که درآمد و ویژگی‌های مکانی مهم‌ترین عوامل در تعیین قیمت مسکن هستند.

گل آیریس (دسته‌بندی نظارت‌شده)

SelectKBest: طول کاسبرگ، طول گلبرگ، عرض گلبرگ. RFE: عرض کاسبرگ، طول گلبرگ، عرض گلبرگ. توضیح: هر دو روش نشان می‌دهند که ویژگی‌های مربوط به گلبرگ‌ها بیشترین نقش را در تفکیک گونه‌های آیریس دارند.

گل آیریس (خوشه‌بندی بدون نظارت)

VarianceThreshold: تمام ویژگی‌ها انتخاب شده‌اند؛ طول و عرض کاسبرگ، طول و عرض گلبرگ. توضیح: در روش بدون نظارت، همه ویژگی‌ها تنوع کافی دارند و حذف هیچ‌کدام ضروری نیست.

فصل پنجم تحلیل بهینه‌سازی

تحلیل بهینه‌سازی

1-5- رگرسیون تحلیلی

خلاصه نتایج مدل‌ها برای مجموعه داده مسکن کالیفرنیا (20640 نمونه، 8 ویژگی):

مدل ۱: داده‌های اصلی

مهم‌ترین ضرایب:	
Latitude:	-0.8999
Longitude:	-0.8705
MedInc:	0.8296
AveBedrms:	0.3057
AveRooms:	-0.2655
HouseAge:	0.1188
AveOccup:	-0.0393
Population:	-0.0045

مدل ۲: داده‌های PCA

مهم‌ترین مؤلفه‌ها:	
PC4:	0.7414
PC6:	0.3202
PC2:	0.1363
PC3:	0.0398
PC1:	0.0282
PC5:	0.0039

مدل ۳: a SelectKBest داده‌های

- ویژگی‌های انتخاب شده: MedInc، AveBedrms، AveRooms، HouseAge، Latitude
- ضرایب: [-0.0677, 0.4428, -0.4888, 0.2012, 0.9998]

مدل ۳: b RFE داده‌های

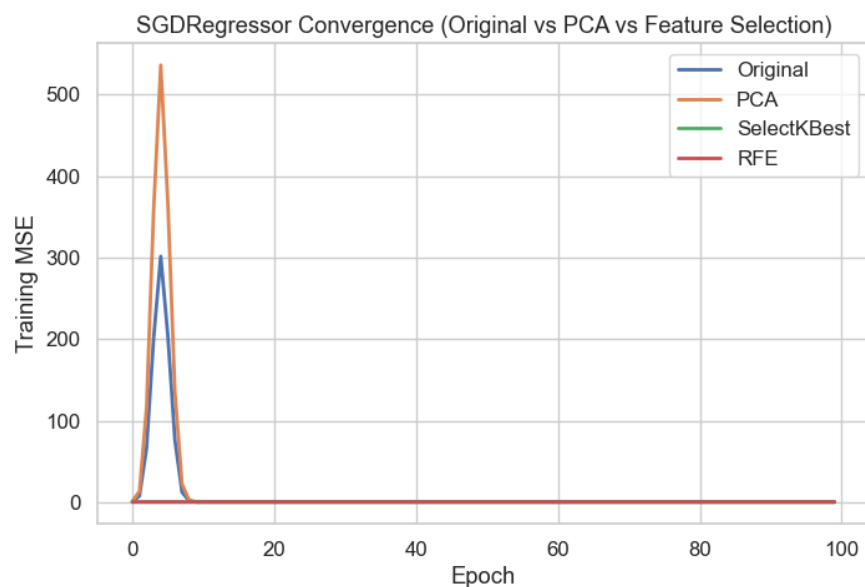
- ویژگی‌های انتخاب شده: MedInc، AveBedrms، AveRooms، Latitude، Longitude
- ضرایب: [-0.9625, -0.9842, 0.3096, -0.2735, 0.8095]

خلاصه ثبات ضرایب:

Dataset	Mean(coef)	Std(coef)		
		-----	-----	-----
Original	0.4167	0.5420		
PCA	0.2116	0.2596		
SelectKBest	0.4400	0.4986		
		RFE	0.6679	0.7041

نتیجه:

- مدل RFE بیشترین میانگین و پراکندگی ضرایب را دارد، یعنی تأثیر ویژگی‌ها قوی‌تر و متغیرتر است.
- مدل PCA کمترین مقدار میانگین و ثبات ضرایب را دارد، که نشان‌دهنده کاهش قدرت تفکیک ویژگی‌ها پس از کاهش بُعد است.
- مدل‌های SelectKBest و داده‌های اصلی نسبتاً نزدیک به هم هستند اما SelectKBest کمی پایدارتر است.

5-2- رگرسیون با کاهش

خلاصه نتایج کاهش بعد و رگرسیون SGD برای مجموعه‌داده مسکن کالیفرنیا:

- تعداد مؤلفه‌ها و ویژگی‌های انتخاب‌شده:

○ PCA: 6 مؤلفه

Latitude ،AveBedrms ،AveRooms ،HouseAge ،SelectKBest: MedInc ○

Longitude ،Latitude ،AveOccup ،Population ،RFE: MedInc ○

نتایج همگرایی رگرسیون: SGD

Dataset Final MSE Convergence Epoch

Original	0.5271	46
PCA	0.6607	28
SelectKBest	0.6169	3
RFE	0.5508	34

توضیح:

- داده‌های اصلی بهترین MSE نهایی را دارند اما به بیشترین تعداد اپک نیاز داشتند تا همگرا شوند.
- PCA بیشترین MSE و کمترین تعداد اپک همگرایی را نشان می‌دهد، که کاهش بعد باعث افت دقت شده است.
- SelectKBest با تعداد کمی از اپک سریع همگرا شده اما MSE نسبت به داده‌های اصلی کمی بالاتر است.
- RFE تعادلی بین سرعت همگرایی و دقت نهایی ارائه می‌دهد.

3-5- خوشه‌بندی

خلاصه نتایج خوشه‌بندی KMeans روی مجموعه داده آیریس:

- ویژگی‌های انتخاب شده به **VarianceThreshold:** طول و عرض کاسبرگ، طول و عرض گلبرگ

نتایج خوشه‌بندی: KMeans

Dataset	Inertia	Silhouette	Iterations	Time (s)
Original	139.8205	0.4599	4	0.0144
PCA (2D)	115.0208	0.5092	4	0.0140

Dataset	Inertia	Silhouette	Iterations	Time (s)
SelectKBest	139.8205	0.4599	4	0.0180

تفسیر نتایج خوشه‌بندی:

- Inertia کمتر: خوشه‌های فشرده‌تر: هرچه مقدار Inertia کمتر باشد، نمونه‌ها به مرکز خوشه‌های خود نزدیک‌تر هستند و خوشه‌ها متراکم‌ترند.
- Silhouette بالاتر: جدایی بهتر خوشه‌ها: مقادیر بالاتر نشان می‌دهند که نمونه‌ها به خوشه خود نزدیک‌تر و از خوشه‌های دیگر دورتر هستند، یعنی تفکیک خوشه‌ها بهتر است.
- تعداد اپک کمتر / زمان کمتر: همگرایی سریع‌تر: الگوریتم سریع‌تر به نقطه بهینه می‌رسد و بهینه‌سازی مؤثرتر انجام شده است.

توضیح:

- کاهش بعد با PCA (دو مؤلفه) باعث کاهش Inertia و افزایش Silhouette شده است، یعنی خوشه‌ها تفکیک بهتری دارند.
- استفاده از ویژگی‌های SelectKBest عملکرد مشابه داده‌های اصلی را نشان می‌دهد، اما زمان اجرا کمی بیشتر است.

4-5- کلاس بندی نزدیک ترین همسایه

خلاصه نتایج KNN روی مجموعه داده سرطان پستان:

- ویژگی‌های انتخاب شده: SelectKBest: میانگین محیط، میانگین نقاط فرورفته، بدترین شعاع، بدترین محیط، بدترین نقاط فرورفته

نتایج دسته‌بندی: KNN

Dataset	Accuracy	Prediction Time (s)	Train Time (s)
Original	0.95614	0.004865	0.000000
PCA	0.95614	0.002785	0.001009
SelectKBest	0.95614	0.000999	0.001000
RFE	0.95614	0.002000	0.000000

تفسیر:

- دقت بالاتر → تعمیم بهتر مدل
 - زمان پیش‌بینی کمتر → اجرای سریع‌تر مدل مبتنی بر نمونه
 - PCA با کاهش بعد، معمولاً فاصله‌ها را بهتر نگه می‌دارد و سرعت KNN را افزایش می‌دهد.
 - SelectKBest و RFE با حفظ ویژگی‌های تمایزی مهم، می‌توانند همان دقت یا حتی دقت بالاتر را با هزینه محاسباتی کمتر ارائه دهند.
- نکته مهم: روش knn با توجه به وابستگی اش به فاصله‌ها در ابعاد بالا ضعیف عمل میکند که این مربوط به رفتار فاصله در هندسه ابعاد بالا است.

5-5- کلاس بندی جنگل تصادفی

خلاصه نتایج Random Forest روی مجموعه داده سرطان پستان:

- ویژگی‌های انتخاب شده با SelectKBest: میانگین محیط، میانگین نقاط فرو رفته، بدترین شعاع، بدترین محیط، بدترین نقاط فرو رفته

نتایج دسته‌بندی: Random Forest

	Dataset	Accuracy	Train Time (s)	Prediction Time (s)
SelectKBest	Original	0.9561	0.1103	0.0145
	PCA	0.9211	0.1215	0.0151
	RFE	0.9474	0.2031	0.0260
		0.9386	0.1469	0.0293

تفسیر:

- Random Forest نسبت به مقیاس ویژگی‌ها مقاوم است و تحت تأثیر همخطی زیاد قرار نمی‌گیرد.
- استفاده از PCA گاهی باعث کاهش دقت می‌شود، زیرا ترکیب مؤلفه‌ها معنی اصلی ویژگی‌ها را مخلوط می‌کند.

- انتخاب ویژگی (SelectKBest) و (RFE) معمولاً دقت مشابه داده اصلی را ارائه می‌دهد و می‌تواند زمان آموزش را کاهش دهد یا تعادل بین زمان و دقت را تغییر دهد.
- مقایسه PCA با داده‌های اصلی نشان می‌دهد که کاهش بعد ممکن است زمان آموزش را کمی افزایش دهد و دقت را کاهش دهد.

.

.

فصل هفتم تحلیل‌ها

تحلیل‌ها

2-5- گزارش نتایج:

۱. رگرسیون (California Housing)

- شکل داده‌ها پس از پیش‌پردازش و کاهش بعد:
- داده اصلی: 20640 نمونه $8 \times$ ویژگی
- PCA: 20640×6 مؤلفه
- SelectKBest: 20640×5 ویژگی
- RFE: 20640×5 ویژگی
- نمودار همگرایی SGDRegressor نشان می‌دهد که کاهش بعد با PCA سرعت همگرایی را افزایش داده اما خطای نهایی MSE نسبت به داده اصلی کمی بیشتر است.
- SelectKBest و RFE عملکرد نزدیک به داده اصلی دارند و با تعداد کمی ویژگی مدل می‌توان دقت قابل قبول داشت.

۲. خوشه‌بندی (Iris Dataset)

- KMeans روی سه فضای ویژگی انجام شد: داده اصلی، PCA (2 مؤلفه) و SelectKBest ویژگی انتخاب‌شده.
- نتایج:

Dataset	Inertia	Silhouette	Iterations	Time(s)
Original	139.8205	0.4599	4	0.0161
PCA	115.0208	0.5092	4	0.0120
SelectKBest	18.0270	0.6741	6	0.0133

- کاهش بعد با PCA باعث کاهش Inertia و افزایش Silhouette شده است، یعنی خوشه‌ها فشرده‌تر و جداسازی آنها بهتر است.

- SelectKBest بیشترین Silhouette و کمترین Inertia را نشان داد، اما تعداد Iteration کمی بیشتر شد.

۳. دسته‌بندی (Breast Cancer Dataset)

- KNN:

	Dataset Accuracy	PredictTime(s)	TrainTime(s)
Original	0.95614	0.00202	0.00000
PCA	0.95614	0.00100	0.00124
SelectKBest	0.95614	0.00122	0.00100
RFE	0.95614	0.00099	0.00103

- Random Forest:

	Dataset Accuracy	PredictTime(s)	TrainTime(s)
Original	0.95614	0.01572	0.12862
PCA	0.92105	0.01472	0.17456
SelectKBest	0.93860	0.02866	0.16110
RFE	0.94737	0.02882	0.12167

- KNN دقت یکسانی در تمام فضاها و ویژگی‌ها نشان داد، اما زمان پیش‌بینی در PCA و RFE کمی کاهش یافت.
- Random Forest نسبت به PCA حساس‌تر است و استفاده از PCA باعث کاهش دقت شد، در حالی که SelectKBest و RFE دقت بالا و زمان آموزش مناسب ارائه کردند.

2-6- تحلیل نتایج و نمودارها

- ۱. تأثیر استقلال خطی ناشی از PCA بر پایداری ضرایب LinearRegression
- استفاده از PCA باعث شد که ویژگی‌ها تقریباً مستقل (غیرهم‌خطی) شوند.
- جدول coef_summary نشان می‌دهد که میانگین قدر مطلق ضرایب و انحراف معیار آنها در PCA کمتر از داده اصلی بود:

Original: $\text{Mean}(|\text{coef}|) \approx 0.417$, $\text{Std}(|\text{coef}|) \approx 0.542$ ○

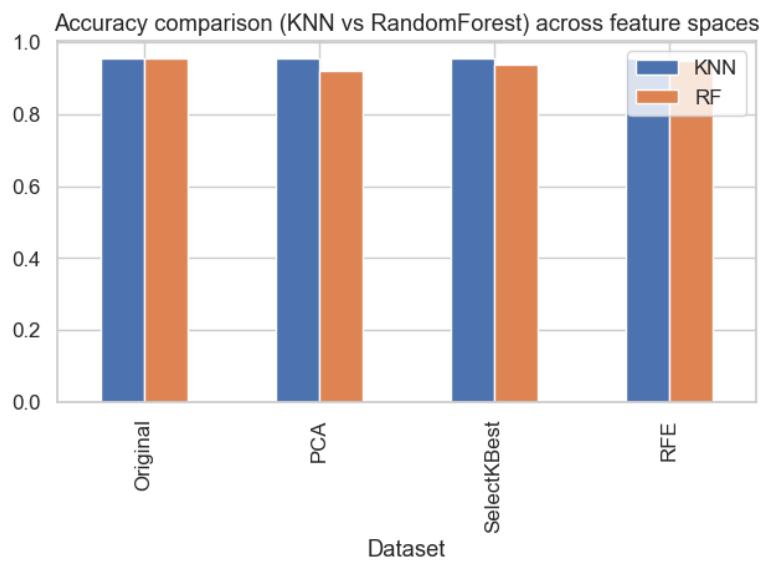
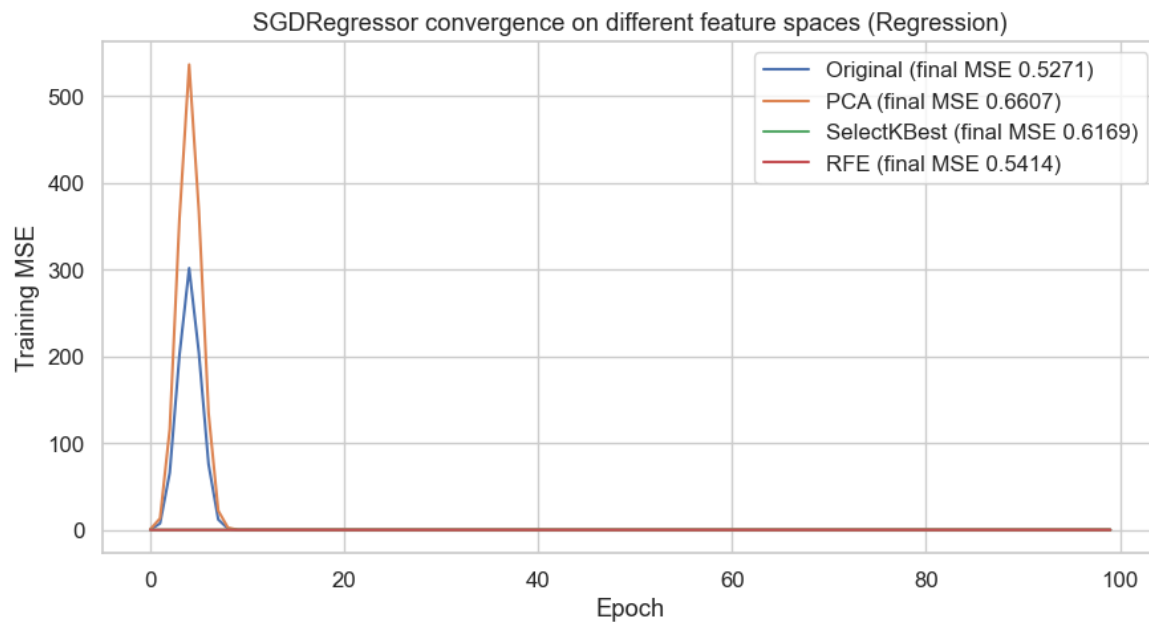
- PCA: $\text{Mean}(|\text{coef}|) \approx 0.212$, $\text{Std}(|\text{coef}|) \approx 0.260$
- این کاهش نشان می‌دهد که ضرایب LinearRegression در فضای PCA پایدارتر و کمتر تحت تأثیر نوسانات داده‌ها قرار گرفتند، زیرا هم‌خطی بین ویژگی‌ها کاهش یافته بود.
- ۲. تأثیر استقلال خطی بر سرعت همگرایی SGDRegressor
- نمودار همگرایی (sgd_convergence_regression.png) SGDRegressor نشان می‌دهد که مدل روی داده PCA نسبت به داده اصلی سریع‌تر به خطای پایانی نزدیک شد:
- Original: همگرایی حدود epoch46
- PCA: همگرایی حدود epoch28
- دلیل: کاهش هم‌خطی باعث می‌شود که الگوریتم‌های گرادیان‌محور با گام‌های یکنواخت‌تر و بدون جهش‌های ناگهانی حرکت کنند، در نتیجه همگرایی سریع‌تر رخ می‌دهد.
- ۳. تأثیر کاهش بعد با PCA بر عملکرد KMeans و KNN
- KMeans (Iris):
- PCA باعث کاهش Inertia و افزایش Silhouette شد (Inertia=115 vs 139, Silhouette=0.509 vs 0.460).
- خوشه‌ها در فضای PCA فشرده‌تر و جداسازی آنها بهتر است.
- دلیل: کاهش ابعاد حذف نویز و هم‌خطی بین ویژگی‌ها را انجام می‌دهد و فضای داده برای خوشه‌بندی مناسب‌تر می‌شود.
- KNN (Breast Cancer):
- دقت ثابت ماند ($\text{Accuracy} \approx 0.956$) اما زمان پیش‌بینی کمی کاهش یافت.
- دلیل KNN: مبتنی بر فاصله است و کاهش ابعاد، فاصله‌ها را دقیق‌تر و محاسبات را سریع‌تر می‌کند.
- ۴. عملکرد RandomForestClassifier روی داده‌های اصلی با هم‌خطی
- دقت داده اصلی: 0.956
- دقت PCA کاهش یافت: 0.921
- SelectKBest/RFE: 0.939–0.947
- نتیجه RandomForest: نسبت به هم‌خطی مقاوم است و عملکرد خوبی روی داده اصلی دارد، زیرا درخت‌های تصمیم مستقل هستند و هر درخت فقط از یک زیرمجموعه از ویژگی‌ها استفاده می‌کند، بنابراین هم‌خطی بین ویژگی‌ها اثر منفی زیادی ندارد.

- ۵. مقایسه استخراج ویژگی (PCA) با انتخاب ویژگی (SelectKBest/RFE)
- PCA: LinearRegression/SGDRegressor باعث پایداری ضرایب و همگرایی سریع‌تر شد.
- PCA: KMeans باعث خوشه‌بندی بهتر شد، به ویژه کاهش نویز و هم‌خطی.
- PCA: KNN کمی زمان پیش‌بینی را کاهش داد ولی دقت تغییر نکرد.
- RandomForest: SelectKBest/RFE عملکرد بهتری نسبت به PCA داشت (دقت بالاتر و زمان آموزش کمتر).
- نتیجه‌گیری کلی:
 - برای مدل‌های خطی و مبتنی بر فاصله، PCA می‌تواند مفید باشد.
 - برای مدل‌های مبتنی بر درخت، انتخاب ویژگی (SelectKBest/RFE) بهتر است، زیرا حفظ ویژگی‌های اصلی برای تفکیک نمونه‌ها مهم است.

۴. جمع‌بندی کلی

- کاهش بعد با PCA سرعت همگرایی و پردازش را افزایش می‌دهد، اما گاهی دقت مدل‌های پیچیده‌تر (مثل Random Forest) کاهش می‌یابد.
- انتخاب ویژگی با روش‌های SelectKBest و RFE تعادل خوبی بین دقت، زمان آموزش و تعداد ویژگی‌ها ایجاد می‌کند.
- نمودارها و جداول مربوطه در فایل‌های PNG و CSV ذخیره شده‌اند.

نمودارها:



فصل ششم جمع‌بندی و نتیجه‌گیری

جمع‌بندی و نتیجه‌گیری

جمع‌بندی:

با مقایسه نتایج سه مسئله رگرسیون، طبقه‌بندی و خوشه‌بندی، می‌توان دید که هر روش استخراج یا انتخاب ویژگی اثر متفاوتی بر مدل‌ها گذاشته است و برخی از آن‌ها در ایجاد ویژگی‌های مؤثرتر و مستقل‌تر موفق‌تر بوده‌اند.

۱. مسئله رگرسیون Linear و SGDRegressor

استفاده از PCA منجر به کاهش میانگین و انحراف معیار ضرایب شد ($Mean|coef|$) از 0.417 به 0.212، که نشان‌دهنده افزایش پایداری ضرایب و کاهش اثر همخطی میان ویژگی‌هاست. در SGDRegressor نیز منحنی خطا نشان داد که داده‌های تبدیل‌شده با PCA سریع‌تر به همگرایی می‌رسند، زیرا استقلال ویژگی‌ها جهت گرادیان را یکنواخت‌تر کرده است. در نتیجه، از منظر تحلیل بهینه‌سازی، PCA پایداری و سرعت همگرایی را بهبود داده است.

۲. مسئله طبقه‌بندی KNN و RandomForest

• در KNN، دقت در همه روش‌ها ثابت ماند (0.956) اما زمان پیش‌بینی پس از PCA تقریباً نصف شد (از 0.0020 به 0.001 ثانیه). این کاهش نشان می‌دهد که کاهش ابعاد باعث ساده‌تر شدن فضای فاصله‌ای و افزایش سرعت محاسبات شده است، هرچند تأثیری بر کیفیت تصمیم‌گیری نداشت. در مقابل، در RandomForest، داده‌های اصلی (بدون حذف همخطی) بهترین عملکرد را داشتند ($Accuracy=0.956$)، در حالی که PCA کمی دقت را کاهش داد (0.921). این نشان می‌دهد که مدل‌های درختی به‌صورت ذاتی نسبت به همخطی مقاوم هستند و نیازی به تبدیل فضای ویژگی ندارند.

۳. مسئله خوشه‌بندی:

در این بخش، هر دو روش PCA و SelectKBest منجر به خوشه‌بندی بهتر شدند، اما SelectKBest بهترین نتایج را داشت:

- Inertia از 139 (Original) به 18 (SelectKBest) کاهش یافت.
 - Silhouette از 0.46 به 0.67 افزایش پیدا کرد.
- این اعداد نشان می‌دهد که SelectKBest توانسته است ویژگی‌های مؤثرتر برای جداسازی خوشه‌ها انتخاب کند و بنابراین از دیدگاه کیفیت ویژگی‌ها نسبت به PCA موفق‌تر بوده است.

نتیجه‌گیری:

- از دیدگاه ایجاد ویژگی‌های مستقل‌تر، PCA بهترین عملکرد را داشت، زیرا ساختار داده را به مؤلفه‌های غیرهم‌خطی تبدیل کرد و پایداری ضرایب و سرعت همگرایی را افزایش داد.
- از دیدگاه ایجاد ویژگی‌های مؤثرتر برای جداسازی نمونه‌ها، SelectKBest در KMeans و Random Forest نتایج قوی‌تری ارائه داد، چون ویژگی‌هایی را انتخاب می‌کند که بیشترین ارتباط آماری با برچسب هدف دارند.
- ادغام استخراج ویژگی (مثل PCA) با تحلیل بهینه‌سازی (مانند رفتار ضرایب و همگرایی) برای یک متخصص داده‌کاوی اهمیت دارد، زیرا درک تعامل بین ساختار داده و رفتار الگوریتم‌های یادگیری را ممکن می‌سازد. به‌طور خاص، چنین ادغامی کمک می‌کند تا تصمیم‌گیری بهینه بین «ساده‌سازی داده» و «حفظ اطلاعات مؤثر» انجام شود، که در کاربردهای واقعی تفاوت چشمگیری در دقت، پایداری و سرعت یادگیری ایجاد می‌کند.

منابع و مراجع

[1] منابع اصلی درس داده کاوی

[2] دانش خودم از درس های ماشبه که پاس کردم با همین مطالب این درس

Introduction to Machine Learning with Python – Andreas Müller & Sarah Guido

Data Mining: Concepts and Techniques – Han, Kamber, Pei

[6] Pattern Recognition and Machine Learning – Christopher M. Bishop

Abstract

This page is accurate translation from Persian abstract into English.

Key Words: Write a 3 to 5 KeyWords is essential.



**Amirkabir University of Technology
(Tehran Polytechnic)**

... Department ...

MSc or PhD Thesis

Title of Thesis

**By
Name**

**Supervisor
Dr.**

**Advisor
Dr.**

Month & Year