



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

گزارش پروژه درس داده کاوی محاسباتی
پروژه ۲

استخراج ویژگی مبتنی بر استقلال خطی و تحلیل تأثیر آن بر فرآیند
بهینه سازی توابع زیان

نگارش
مهندی حیدری

استاد درس
نام کامل استاد راهنما

تدریسیار
مهندس بهنام یوسفی مهر

آبان ۱۴۰۴

چکیده

هدف از این پژوهه، درک عمیق ارتباط میان کیفیت فضای ویژگی و عملکرد الگوریتم‌های بهینه‌سازی بود. ما به طور خاص، نأثیر «استقلال خطی» (Linear Independence) که از طریق روش استخراج ویژگی PCA به دست می‌آید را بر پایداری و سرعت همگرایی مدل‌های یادگیری ماشین بررسی کردیم.

آزمایش‌های ما بر روی سه مجموعه داده (رگرسیون، طبقه‌بندی و خوشه‌بندی) نشان داد که: ۱. پایداری: استفاده از PCA ضرایب ناپایدار و بسیار بزرگ مدل [LinearRegression](#) را که ناشی از همخطی (Collinearity) بود، به طور کامل برطرف و پایدار ساخت. ۲. سرعت همگرایی: فضای ویژگی متعامد حاصل از PCA، سرعت همگرایی بهینه‌ساز مبتنی بر گرادیان ([SGDRegressor](#)) را به شکلی چشمگیر (بیش از ۱۰ برابر سریع‌تر) افزایش داد. ۳. کیفیت: در مدل‌های مبتنی بر فاصله ([KNN](#) و [KMeans](#))، کاهش ابعاد با PCA منجر به خوشه‌های باکیفیت‌تر (Silhouette بالاتر) و پیش‌بینی سریع‌تر شد. ۴. مقاومت: مدل [RandomForest](#) (مبتنی بر درخت) مقاومت کامل خود را در برابر همخطی نشان داد و عملکرد آن روی داده‌های اصلی با ۳۰ ویژگی و داده‌های انتخاب‌شده (RFE) با ۱۰ ویژگی یکسان بود.

این نتایج به صورت تجربی ثابت می‌کند که طراحی فضای ویژگی مناسب، یک گام حیاتی برای موفقیت فرآیند بهینه‌سازی است.

صفحه

۵

فهرست مطالب

چکیده	۱
مقدمه و بیان مسئله	۴
مراحل انجام پروژه و آماده‌سازی داده‌ها	۵
نتایج آزمایش‌ها و تحلیل بهینه‌سازی (مرحله ۵)	۶
(Housing) بخش الف: رگرسیون	۶
(Iris) بخش ب: خوشه‌بندی	۶
(Cancer) بخش ج: طبقه‌بندی	۷
تحلیل نتایج (پاسخ به سوالات مرحله ۶)	۸
نتیجه‌گیری (مرحله ۷)	۹
پیوند مخزن پروژه	۱۰

مقدمه و بیان مسئله

در بسیاری از مسائل داده‌کاوی، داده‌های خام دارای ویژگی‌های متعدد و همبسته هستند. این پدیده که «همخطی» (Collinearity) نام دارد، مشکلات جدی در فرآیند مدل‌سازی ایجاد می‌کند؛ از جمله ناپایداری در ضرایب مدل‌های خطی و کند شدن شدید الگوریتم‌های بهینه‌سازی مبتنی بر مشتق.

برای حل این مشکل، دو رویکرد اصلی وجود دارد:

- **انتخاب ویژگی (Feature Selection)**: حذف ویژگی‌های کم‌اهمیت یا همبسته .(مانند RFE)

- استخراج ویژگی (Feature Extraction): تبدیل فضای ویژگی به یک فضای جدید و کمابعدتر که ویژگی‌های آن مستقل خطی باشند (مانند PCA).

در این پروژه، ما هر دو رویکرد را پیاده‌سازی کرده و تأثیر آن‌ها را بر پایداری (در مدل‌های تحلیلی)، سرعت همگرایی (در مدل‌های مبتنی بر گرادیان) و کیفیت (در مدل‌های مبتنی بر فاصله و درخت) به صورت تجربی تحلیل می‌کنیم.

مراحل انجام پروژه و آماده‌سازی داده‌ها

پروژه در ۷ مرحله طبق دستورالعمل^۱ انجام شد:

مرحله ۱: بارگیری داده‌ها

سه مجموعه داده کلاسیک برای سه وظیفه مختلف بارگیری شدند:

- طبقه‌بندی (Classification): Wisconsin Breast Cancer (569 نمونه، 30 ویژگی).

Boston Housing (جایگزین Regression): California Housing (20640 نمونه، 8 ویژگی).

- خوشه‌بندی (Clustering): UCI Iris (150 نمونه، 4 ویژگی).

مرحله ۲: تحلیل همخطی اولیه

با استفاده از نقشه حرارتی (Heatmap) ماتریس همبستگی، همخطی شدید در داده‌ها تایید شد:

- همخطی بسیار شدید بین ویژگی‌های مرتبط با شعاع، محیط و مساحت تومور (مانند mean perimeter و mean radius با همبستگی 0.998 مشاهده شد).

- همبستگی بالا بین Latitude و Longitude (0.92) و AveRooms و AveBedrms وجود داشت.

• Iris: همبستگی قوی بین petal length (0.96) و petal width (0.96) تایید شد.

مرحله ۳: استخراج ویژگی (PCA)

ما از PCA با هدف حفظ ۹۵٪ واریانس داده‌های اصلی استفاده کردیم. این کار منجر به کاهش ابعاد قابل توجهی شد:

• Cancer: از ۳۰ ویژگی به ۱۰ مؤلفه اصلی (PC) کاهش یافت.

• Housing: از ۸ ویژگی به ۶ مؤلفه اصلی کاهش یافت.

• Iris: برای اهداف خوشبندی و مصورسازی، به ۲ مؤلفه اصلی کاهش یافت (که ۹۵.۸٪ واریانس را پوشش داد).

مرحله ۴: انتخاب ویژگی (RFE)

برای ایجاد یک معیار مقایسه، از روش (RFE) با یک مدل پایه (مانند RandomForest برای طبقه‌بندی) استفاده کردیم تا داده‌ها را به همان ابعاد PCA کاهش دهیم (مثلاً ۱۰ ویژگی برای Cancer). این مرحله برای Iris (خوشبندی) به دلیل ماهیت بدون ناظر آن انجام نشد.

نتایج آزمایش‌ها و تحلیل بهینه‌سازی (مرحله ۵)

در این مرحله، مدل‌ها بر روی سه نسخه داده آموزش داده شدند.

بخش الف: رگرسیون (Housing)

۱. مدل تحلیلی (LinearRegression):

• پایداری ضرایب: همانطور که در جدول ۱ مشاهده می‌شود، ضرایب مدل اصلی به

دلیل هم خطی، بسیار بزرگ و ناپایدار هستند (مثلاً -0.89 برای Latitude). در مقابل، ضرایب مدل PCA کاملاً پایدار و معنادار هستند.

• خطای PCA (MSE): مدل PCA خطای کمی بالاتری (0.6789) نسبت به مدل اصلی (0.5416) داشت. این یک موازنۀ Trade-off (قابل انتظار است، زیرا PCA برای رسیدن به استقلال خطی، بخشی از اطلاعات (۵٪ واریانس) را حذف کرده است).

مدل	MSE آزمون	ضرایب ناپایدار (مثال)
Original	0.5416	Latitude: -0.896, Longitude: -0.866
PCA	0.6789	PC4: 0.742, PC6: 0.323 (همگی پایدار)
جدول ۱: مقایسه پایداری ضرایب رگرسیون خطی		

۲. مدل مبتنی بر مشتق (SGDRegressor)

- سرعت همگرایی (Iteration): جدول تحلیل حساسیت (Run 3) نشان داد که مدل PCA در ۹۲ تکرار و مدل اصلی در ۹۵ تکرار به همگرایی کامل ($tol=1e-6$) رسیدند.
- نمودار همگرایی (الزامی): نمودار ۱ پاسخ قاطع را نشان می‌دهد. مدل SGD روی داده‌های PCA (خط نارنجی) در کمتر از ۵ اپوک به خطای نهایی خود همگرا شد، در حالی که مدل روی داده‌های اصلی (خط آبی) به حدود ۴۰ تا ۵۰ اپوک برای همگرایی نیاز داشت.

نمودار ۱: مقایسه سرعت همگرایی $SGDRegressor$ (اپوک در برابر MSE)
بخش ب: خوشبندی (Iris)
:(KMeans) مدل:

همانطور که در جدول ۲ مشخص است، PCA هم کیفیت و هم سرعت خوشبندی را بهبود بخشدید:

- کیفیت: اینرسی (Inertia) کاهش (خوشههای فشردهتر) و شاخص Silhouette افزایش یافت (خوشههای جدا از هم).
- سرعت: خوشبندی روی ۲ مؤلفه PCA در تکرارهای کمتر (۴ در برابر ۶) و زمان کمتری انجام شد.

تکرار (Iter)	Silhouette بهتر	Inertia بهتر	مدل (کمتر)
6	0.4599	139.82	Original
4	0.5092	115.02	PCA
			جدول ۳: مقایسه عملکرد KMeans

بخش ج: طبقه‌بندی (Cancer)

۱. مدل مبتنی بر فاصله (KNN):

- دقت (Accuracy): دقت روی داده‌های اصلی (۳۰ ویژگی) و (۱۰ PCA ویژگی) یکسان (%98.25) بود.

- سرعت (Predict Time): کاهش ابعاد تأثیر مستقیم بر زمان پیش‌بینی داشت. مدل (۱۰ RFE ویژگی) بیش از ۲.۳ برابر سریع‌تر از مدل اصلی بود، که حساسیت KNN به ابعاد بالا را نشان می‌دهد.

۲. مدل مبتنی بر درخت (RandomForest):

- مقاومت: این مدل نتیجه‌ای کلیدی را نشان داد. دقت روی داده‌های اصلی (۳۰ ویژگی) و (۱۰ RFE ویژگی) کاملاً یکسان (95.61%) بود.
- در برابر RFE: دقت مدل PCA (۹۳.۸۶%) پایین‌تر بود.

مدل	داده‌ها (ویژگی)	دقت (Accuracy)	زمان پیش‌بینی (ثانیه)
KNN	Original ((30	0.9825	0.0084
	(PCA (10	0.9825	0.0052
	(RFE (10	0.9737	0.0036
RandomForest	Original ((30	0.9561	0.0159
	(PCA (10	0.9386	0.0141
	(RFE (10	0.9561	0.0144
جدول ۳: مقایسه عملکرد مدل‌های طبقه‌بندی			

تحلیل نتایج (پاسخ به سوالات مرحله ۶)

۱. تأثیر PCA بر پایداری ضرایب LinearRegression چه بود؟

PCA مشکل ناپایداری را کاملاً حل کرد. ضرایب بزرگ و متضاد (مانند -0.86 و -0.89) در داده‌های اصلی، که ناشی از همخطی بود، در مدل PCA به ضرایبی پایدار و قابل تفسیر تبدیل شدند.

۲. تأثیر استقلال خطی بر سرعت همگرایی SGDRegressor چه بود؟

تأثیر آن بسیار چشمگیر بود. همانطور که نمودار ۱ نشان می‌دهد، داده‌های متعامد PCA سرعت همگرایی را بیش از ۱۰ برابر افزایش دادند (همگرایی در کمتر از ۵ اپوک به جای ۵۰ اپوک).

۳. تأثیر کاهش ابعاد (PCA) بر KMeans و KNN چه بود؟ چرا حساس هستند؟

PCA عملکرد هر دو مدل را بهبود بخشد. برای KMeans، کیفیت خوشها (Silhouette) بالاتر رفت و سرعت افزایش یافت. برای KNN، سرعت پیش‌بینی به طور قابل توجهی (۲.۳ برابر) بهتر شد.

این مدل‌ها به ابعاد حساس هستند زیرا هر دو مبتنی بر «فاصله اقلیدسی» می‌باشند. در ابعاد بالا (نفرین ابعاد)، مفهوم فاصله کم‌رنگ‌تر شده و محاسبات کند می‌شود ۲.

۴. عملکرد RandomForest روی داده‌های اصلی (با همخطی) چگونه بود؟ آیا مقاوم است؟

عملکرد آن عالی و با داده‌های انتخاب‌شده (RFE) یکسان بود (دقت ۹۵.۶۱٪). این نشان می‌دهد RandomForest به دلیل مکانیسم انتخاب ویژگی تصادفی در هر گره، ذاتاً به همخطی مقاوم است و توانست ویژگی‌های اضافی را نادیده بگیرد ۳.

۵. در نهایت، آیا استخراج ویژگی (PCA) بهتر از انتخاب ویژگی (RFE) عمل کرد؟

پاسخ به مدل بستگی دارد:

- برای RandomForest: RFE بهتر بود (دقت ۹۵.۶۱٪ در برابر ۹۳.۸۶٪ برای PCA). مدل‌های درختی ویژگی‌های اصلی و قابل تفسیر را ترجیح می‌دهند.

- برای **SGD** و **PCA** و **KMeans**: بهتر بود، زیرا مستقیماً مشکل استقلال خطی و ابعاد بالا را که این مدل‌ها به آن حساس بودند، حل کرد.
- برای **RFE** و **KNN**: سریع‌تر بود، اما **PCA** دقیق‌تر بود.

نتیجه‌گیری (مرحله ۷)

این پژوهه به طور تجربی ارتباط مستقیم بین هندسه فضای ویژگی و فرآیند بهینه‌سازی را نشان داد. هم خطی (Collinearity) یک مشکل واقعی است که منجر به ناپایداری مدل‌های تحلیلی و کندی شدید مدل‌های مبتنی بر گرادیان می‌شود.

استخراج ویژگی با PCA یک راه حل قدرتمند برای مدل‌های «حساس» (مبتنی بر فاصله یا گرادیان) است که سرعت و کیفیت را به طور همزمان بهبود می‌بخشد. با این حال، PCA یک راه حل جادویی برای همه مدل‌ها نیست؛ مدل‌های « مقاوم » (مانند RandomForest) نه تنها به آن نیازی ندارند، بلکه ممکن است عملکرد ضعیفتری روی مؤلفه‌های انتزاعی PCA نشان دهند و انتخاب ویژگی (RFE) برای آن‌ها گزینه بهتری باشد.

درک این موازندها برای یک متخصص داده‌کاوی جهت انتخاب روش پیش‌پردازش مناسب بر اساس الگوریتم یادگیری نهایی، امری حیاتی است.

پیوند مخزن پژوهه

جهت بررسی کدها و فایل نوت‌بوک، می‌توانید به مخزن عمومی گیت‌هاب زیر مراجعه نمایید:

<https://github.com/MrHidr/Computational-Data-Mining-Homeworks/tree/main/HW2>
<https://colab.research.google.com/drive/1u6-wkW05y5PYALTlqv8WLYZCWIM96uAx?usp=sharing>

مراجع

این بخش، منابع استفاده شده برای مفاهیم نظری، الگوریتم‌ها، مجموعه داده‌ها و کتابخانه‌های نرم‌افزاری مورد استفاده در این پژوهه را فهرست می‌کند.

کتاب‌ها و مقالات مرجع (مفاهیم و الگوریتم‌ها)

- .**Breiman, L. (2001).** Random forests. *Machine learning*, 45(1), 5-32 .1
 (منبع اصلی برای معرفی الگوریتم RandomForestClassifier ○
- Cover, T., & Hart, P. (1967).** Nearest neighbor pattern classification. .2
.IEEE transactions on information theory, 13(1), 21-27
 (مقاله کلاسیک برای معرفی الگوریتم KNN ○
- Fisher, R. A. (1936).** The use of multiple measurements in taxonomic .3
 problems. *Annals of eugenics*, 7(2), 179-188
 (مقاله اصلی که مجموعه داده Iris را معرفی کرد ○
- Guyon, I., Weston, J., Barnhill, S., & Vapnik, V. (2002).** Gene selection .4
 for cancer classification using support vector machines. *Machine
 learning*, 46(1), 389-422
 (منبع معتبر برای روش Recursive Feature Elimination - RFE ○
- Hastie, T., Tibshirani, R., & Friedman, J. (2009).** *The Elements of Statistical Learning: Data Mining, Inference, and Prediction*. Springer .5
 .Science & Business Media
 (کتاب مرجع اصلی داده‌کاوی برای مفاهیم پایه مانند SVD، PCA ○
 (رگرسیون و همخطی)
- Hyvärinen, A. (1999).** Fast and robust fixed-point algorithms for .6
 independent component analysis. *IEEE transactions on neural
 networks*, 10(3), 626-634
 (مقاله اصلی برای الگوریتم FastICA که در Scikit-learn استفاده شد ○

- Lloyd, S. P. (1982).** Least squares quantization in PCM. *IEEE transactions on information theory*, 28(2), 129-137
- (منبع استاندارد برای الگوریتم خوشبندی KMeans - مبتنی بر M-E)
- Pace, R. K., & Barry, R. (1997).** Sparse spatial autoregressions. .8
- .*Statistics & Probability Letters*, 33(3), 291-297
- (منبع معرفی مجموعه داده Boston California Housing که جایگزین شد)
- Pearson, K. (1901).** LIII. On lines and planes of closest fit to systems .9 of points in space. *The London, Edinburgh, and Dublin Philosophical Magazine and Journal of Science*, 2(11), 559-572
- (مقاله اصلی PCA برای تحلیل مؤلفه‌های اصلی - foundational)
- Street, W. N., Wolberg, W. H., & Mangasarian, O. L. (1993).** Nuclear .10 feature extraction for breast tumor diagnosis. *In Biomedical Image Processing and Biomedical Visualization* (Vol. 1905, pp. 861-870).
- .SPIE
- (مقاله اصلی معرفی کننده مجموعه داده Wisconsin Breast Cancer)
- کتابخانه‌های نرم‌افزاری و ابزارها
- Hunter, J. D. (2007).** Matplotlib: A 2D graphics environment. *Computing in science & engineering*, 9(3), 90-95
- (منبع استاندارد برای کتابخانه مصورسازی Matplotlib)
- McKinney, W. (2010).** Data structures for statistical computing in .12 python. *In Proceedings of the 9th Python in Science Conference* (Vol. 445, pp. 51-56)
- (مقاله معرفی کتابخانه Pandas)

- Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., .13**
- Grisel, O., ... & Duchesnay, É. (2011).** Scikit-learn: Machine learning .in Python. *Journal of machine learning research*, 12(Oct), 2825-2830
- (مقاله رسمی و منبع اصلی برای ارجاع به کتابخانه Scikit-learn)
- Van der Walt, S., Colbert, S. C., & Varoquaux, G. (2011).** The NumPy .14 array: a structure for efficient numerical computation. *Computing in Science & Engineering*, 13(2), 22-30
- (مقاله رسمی برای ارجاع به کتابخانه NumPy)



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

... Department ...

MSc or PhD Thesis

Title of Thesis

**By
Name**

**Supervisor
Dr.**

**Advisor
Dr.**

Month & Year