

دانشگاه صنعتی امیرکبیر (پلی تکنیک تهران) دانشکده ریاضی و علوم کامپیوتر

پایاننامه کارشناسیارشد گرایش سیستمهای کامپیوتری

کاهش بعد دادههای بزرگ مقیاس با استفاده از نگاشت تصادفی

نگارش سیامک دهبد

استاد راهنما دکتر عادل محمدپور

> استاد مشاور دکتر هادی زارع

> > دی ۱۳۹۷

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

در این صفحه فرم دفاع یا تایید و تصویب پایان نامه موسوم به فرم کمیته دفاع- موجود در پرونده آموزشی- را قرار دهید.

نكات 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)
- رنگ جلد پایان نامه/رساله چاپی کارشناسی، کارشناسی ارشد و دکترا باید به ترتیب مشکی، طوسی و سفید رنگ باشد.
 - چاپ و صحافی پایان نامه/رساله بصورت پشت و رو(دورو) بلامانع است و انجام آن توصیه می شود.



تاریخ: دی ۱۳۹۷

دانشگاه صنعتی امیر کبیر (یلی تکنیک تهران)

اینجانب سیامک دهبد متعهد می شوم 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 پایاننامه 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 پایاننامه 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 كتبي دانشگاه صنعتي امير كبير ممنوع است. نقل مطالب با ذكر مآخذ بلامانع است.

سیامک دهبد

امضا



با تشکر از استاد گرامی دکتر محمدپور بابت همراهی و صبر ایشان

سامک دسد دی ۱۳۹۷

چکیده

روش تصویر تصادفی برای کاهش بعد دادههای بزرگ مقیاس مزایای متعددی نسبت به روشهای دیگر کاهش بعد دارد. در این پایاننامه این روش برای دادههای بزرگ مقیاس با دیگر روشهای کاهش بعد مقایسه شده است. همچنین توانایی این روش برای دادههای با توزیع پایدار غیر نرمال با دیگر روشهای کاهش بعد مقایسه شده است.

واژههای کلیدی:

کاهش بعد، تصویر تصادفی، توزیع پایدار، دادههای بزرگ مقیاس

فحه	صف	فهرست مطالب	عنوان
١			۱ مقدمه
٢			۱–۱ مقدمه
۲			١-١-١ مقدما
٣			۲ مرور ادبیات
۴			۲-۱ مقدمه
۵			۲-۲ دادههای حجی
۵		ای حجیم وب	۲-۲-۱ دادهه
۶			
۶			
۶			۵-۲ کاهش بعد .
۶		ِی مراجع	
٧			منابع و مراجع
٨			پيوست
٩		انگلیسی	واژەنامەي فارسى بە
11		به فا رسی	واژەنامەي انگلیسي ب

فهرست اشكال

فهرست اشكال

صفحه

شكل

صفحه	•	فهرست جداول	جدول
۵		ت برای کلمات با بازخورد بالا و کلمات با بازخورد نادر	۱-۲ تعداد بازدید صفحا

فهرست نمادها

نماد مفهوم \mathbb{R}^n n فضای اقلیدسی با بعد n کرہ یکه n بعدی \mathbb{S}^n M جمینهm-بعدی M^m M وی هموار روی M $\mathfrak{X}(M)$ (M,g) مجموعه میدانهای برداری هموار یکه روی $\mathfrak{X}^{\prime}(M)$ M مجموعه p-فرمیهای روی خمینه $\Omega^p(M)$ اپراتور ریچی Qتانسور انحنای ریمان \mathcal{R} تانسور ریچی ricمشتق لي L۲-فرم اساسی خمینه تماسی Φ التصاق لوى-چويتاي ∇ لاپلاسين ناهموار Δ عملگر خودالحاق صوری القا شده از التصاق لوی-چویتای ∇^* متر ساساكى g_s التصاق لوی-چویتای وابسته به متر ساساکی ∇ عملگر لاپلاس-بلترامی روی p-فرمها Δ

٥

فصل اول مقدمه

۱–۱ مقدمه

۱–۱–۱ مقدمه

فصل دوم مرور ادبیات

۱-۲ مقدمه

عمومیت پیدا کردن دادههای حجیم مانند دادههای حجیم تحت وب و جریانهای داده بزرگ در کاربردهای جدید، موجب به وجود آمدن فرصتهای و چالشهایی برای مهندسین و دانشمندان شده است. [T] برای مثال، زمانی که موجب به وجود آمدن فرصتهای و چالشهایی برای مهندسین و دانشمندان شده است. $\mathbf{A} \mathbf{A}^T$ سخت می شود. ماتریس داده $\mathbf{A} \mathbf{A}^T$ ابعادی در حد وب داشته باشد، عملیات سادهای مانند محاسبه $\mathbf{A} \mathbf{A}^T$ سخت می شود. برای ارائه و نگهداری دادههای حجیم در حافظهای کوچک و برای استخراج اطلاعات آماری اصلی از مجوعهای از برای بیانی محدود، روشهای گوناگونی نمونهبرداری توسعه یافته است. به طور کلی روش تصویر تصادفی پایدار \mathbf{A} برای دادههای با دم سنگین خیلی خوب کار می کند.

 $\mathbf{R} \in \mathbb{R}^{D \times k}$ روش تصویر تصادفی پایدار، ماتریس دادههای اولیه $\mathbf{A} \in \mathbb{R}^{n \times D}$ را در ماتریس تصادفی \mathbf{R} علی $\mathbf{A} \in \mathbb{R}^{n \times k}$ است. معمولا درایههای ماتریس تصادفی $\mathbf{B} = \mathbf{A} \mathbf{R} \in \mathbb{R}^{n \times k}$ می ضرب می کند و نتیجه ماتریس مقارن انتخاب می شوند ($^\circ$ حرص حرص ورد مخصه های می از یک توزیع $^\circ$ بایدار متوارن انتخاب می شوند و تصادفی پایدار توسط لم $^\mathsf{T}$ برجسته شده را در $^\mathsf{T}$ برجسته شده است. لم $^\mathsf{T}$ بیان می دارد که کافی است $^\mathsf{T}$ و باید و باید و به دویی با نرم $^\mathsf{T}$ در ابتوان با است. لم $^\mathsf{T}$ بیان می دارد که کافی است $^\mathsf{T}$ و باید و باید و به دویی با نرم $^\mathsf{T}$ در ابتوان با نرم $^\mathsf{T}$ و باید و باید و به دویی با نرم $^\mathsf{T}$ در ابتوان با است. روش تصویر تصادفی پایدار به یک مسئله تخمین آماری کاهش می باید برای تخمین پارامتر مقیاس برای است. روش تصویر تصادفی پایدار به یک مسئله تخمین آماری کاهش می باید برای تخمین پارامتر مقیاس برای یک توضیع پایدار $^\mathsf{T}$ متقارن. این مسئله از این جهت مورد توجه قرار می گیرد زیرا ما به دنبال برآوردی می گردیم که هم از نظر آماری درست باشند و هم از نظر محاسباتی مقرون به صرفه. برآوردگرهای مختلفی را مطالعه و مقایسه کردیم. شامل میانگین حسابی، میانگین هندسی، میانگین هارمونیک، تقسیم توانی $^\mathsf{T}$ و برآوردگر حداکثر مقایسی،

در این پایاننامه ما به بررسی موارد خاصی از تصویر تصادفی پایدار می پردازیم. برای نرم l_1 ارتقایی را با استفاده از اطلاعات حاشیهای پیشنهاد می کنیم. همچنین برای حالت l_1 می توان ماتریس تصویر گر را از یک توزیع زیر گوسی بسیار کوچکتر به جای توزیع نرمال انتخاب کرد. با در نظر گرفتن محدودیتهای معقولی می توان از یک توزیع خاص زیر گوسی استفاده کرد. این توزیع شامل [-1,0,0] با احتمالات $\{\frac{1}{s},1-\frac{1}{s},\frac{1}{s}\}$ با مقادیر بسیار بزرگی برای s (به عبارتی، تصویر تصادف خیلی گسسته s) می تواند به خوبی تصویر تصادفی نرمال عمل کند. برای حالت نرم s به عبارتی دیگر تصویر تصادفی کوچی انجام تخمین کاری نسبتا جذاب است. برای مثال، محاسبه برآوردگر بیشینه درستنمایی MLE در این حالت از لحاظ محاسباتی ممکن است. و یک توزیع معکوس گاوسی s برای مدل سازی دقیق توزیع s این شده است.

روش تصویر تصادفی از پراکندگی دادهها استفادهای نمی کند. در حالی که دادههای بزرگ مقیاس معمولاً بسیار پراکنده هستند. از روش تصویر تصادفی می توان برای حل مسائل بزرگ مقیاس در علوم و مهندسی در موتورهای جستجو و سیستمهای اخذ داده، پایگاههای داده، سیستمهای جریان داده جدید، جبر خطی عددی و بسیاری از کارهای یادگیری ماشین و داده کاوی که شامل محاسبه حجیم فاصلهها است، استفاده کرد.

¹Stable Random Projection

²Independent and identically distributed random variables

³Johnson-Lindenstrauss

⁴fractional power

 $^{^5 \}mathrm{sub}\text{-}\mathrm{Gaussian}$

⁶Very sparse random projections

⁷Cauchy random projections

⁸inverse Gaussian

۲-۲ دادههای حجیم

عبارات زیر از سایت Information Week نقل قول شدهاند $^{\circ}$:

- مقدار دادهای که توسط کسب و کارها ذخیره میشود تقریبا هر ۱۲ تا ۱۸ ماه دو برابر میشود.
- پایگاه دادهها بیشتر هم زمان شدهاند. فروشگاههای زنجیرهای Wall-Marat دادههای فروش را هر ساعت به روز می کند.
- اضافه شدن یک میلیون خط داده اجازه جستجوهای پیچیده تری را می دهد. شرکت EBay به کارمندان اجازه می دهد برای بدست آوردن در کی عمیق تر در خصوص رفتار مشتریان در میان داده های حراج در بازه های زمانی کوتاه جستجو کنند.
- بزرگترین پایگاه دادهها توسط، مرکز شتابدهنده خطی استاندارد، مرکز تحقیقات ناسا، آژانس امنیت ملی و
 ... در ابعادی در محدوده ی پتابایت (هزار ترابایت ۱۰۱۵ بایت)، اداره می شوند.

پدیده نو ظهور مجموعه داده ای حجیم، چالشهای محاسباتی در بسیاری کاربردهای علمی و تجاری به وجود آورده است. شامل اخترفیزیک، بیوتکنولوزی، جمعیت شناسی ۱۰۰ ، مالی، سیستمهای اطلاعات جغرافیایی، دولت، دارو، ارتباطات از راه دور، محیط زیست و اینترنت.

۲-۲-۱ دادههای حجیم وب

وب چقدر بزرگ است؟ جدول ۲-۱ نشان دهنده تعداد بازدید صفحات در موتورهای جستجوی امروزی است. به طور تخمینی حدود $D = 10^{10}$ صفحه وب را می توان بر اساس بازدید دو واژه ی بسیار پر کاربرد « $D = 10^{10}$ » تخمین زد. جدول ۲-۱ همچنین نشان می دهد که حتی کلماتی که به ندرت کاربرد دارند هم تعداد زیادی بازدید دارند.

Query	Google	Bing
A	25,270,000,000	175,000,000
The	25,270,000,000	101,000,000
Kalevala	7,440,000	939,000
Griseofulvin	1,163,000	332,000
Saccade	1,030,000	388,000

جدول ۲-۱: تعداد بازدید صفحات برای کلمات با بازخورد بالا و کلمات با بازخورد نادر

کلماتی با بازخورد معمولی چه میزان بازدید دارند؟ برای جواب این سوال ما به طور تصادفی ۱۵ صفحه از لغتنامهی آموزشی انتخاب میکنیم. [۲] (لغتنامهای با ۵۷،۰۰۰ کلمه) و اولین کلمه در هر صفحه را مد نظر قرار میدهیم. میانهی آماری بر اساس جستجوگر گوگل ۱۰ میلیون صفحه برای کلمه است.

زبان انگلیسی چند کلمه دارد؟ در اینجا عبارتی را از AskOxford.com نقل قول می کنیم:

« این بیان میدارد که حداقل یک چهارم میلیون واژه ی انگلیسی مستقل وجود دارد. به جز افعال صرفی و کلمات فنی و ناحیهای که توسط OED " تحت پوشش قرار نمی گیرند یا کلماتی که هنوز به لغتنامههای

 $^{^9} http://www.informationweek.com/news/showArticle.jhtml? article ID = 175801775$

¹⁰demographics

¹¹Oxford English Dictionary

منتشر شده اضافه نشدهاند. در صورتی که این موارد هم در نظر گرفته شوند تعداد لغات در حدود سه چهارم میلیون لغت خواهد بود »

بنابراین اگر یک ماتریس «عبارت به سند» $\mathbf{A} \in \mathbb{R}^{n \times D}$ در نظر بگیریم. در ابعاد وب این ماتریس در ابعاد بنابراین اگر یک ماتریس i می در سند i و i در سند i در سند i در اینجا عدد (i,j) در i تعداد ظهور واژه i در سند i در انشان می دهد.

کارکردن با ماترسی در این ابعاد بزرگ چالش برانگیز است. برای مثال، شاخص LSI ۱۱ [۱]

- ۲-۳ کاهش بعد
- ۲-۲ کاهش بعد
- ۲–۵ کاهش بعد
- ۲-۵-۲ بارگیری مراجع

¹²latent semantic indexing

منابع و مراجع

- [1] Deerwester, Scott, Dumais, Susan T, Furnas, George W, Landauer, Thomas K, and Harshman, Richard. Indexing by latent semantic analysis. *Journal of the American society for information science*, 41(6):391–407, 1990.
- [2] Hornby, Albert Sydney, editor. *Oxford Advanced Learner's Dictionary of Current English*. Oxford University Press, Oxford, UK, fourth edition, 1989.
- [3] Li, Ping. Stable random projections and conditional random sampling, two sampling techniques for modern massive datasets. Stanford, 2007.

پيوست

موضوعات مرتبط با متن گزارش پایان نامه که در یکی از گروههای زیر قرار می گیرد، در بخش پیوستها آورده شوند:

```
۱. اثبات های ریاضی یا عملیات ریاضی طولانی.
```

۲. داده و اطلاعات نمونه (های) مورد مطالعه (Case Study) چنانچه طولانی باشد.

۳. نتایج کارهای دیگران چنانچه نیاز به تفصیل باشد.

۴. مجموعه تعاریف متغیرها و پارامترها، چنانچه طولانی بوده و در متن به انجام نرسیده باشد.

کد میپل

```
with(DifferentialGeometry):
with(Tensor):
DGsetup([x, y, z], M)
frame name: M
a := evalDG(D_x)
D_x
b := evalDG(-2 y z D_x+2 x D_y/z^3-D_z/z^2)
```

واژهنامهی فارسی به انگلیسی

خودریختی Automorphism	Ĩ
s	اسکالر
Degree	ب
j	
ریز پر دازنده microprocessor	بالابر
j	پ
زيرمدولزيرمدول	پایا
س	ت
Character	تناظر
ص	ث
صادقانه Faithful	ثابتساز Stabilizer
ض	τ
ضرب داخلی	Permutation
ط	€
طوقه	چند جملهای Polynomial
ظ	τ
ظرفیت	حاصل ضرب دکارتی Cartesian product
3	έ

عدم مجاورت Nonadjacency
ف
فضای برداری Vector space
ى
كاملاً تحويل پذير Complete reducibility
گ
گراف
۴
ماتریس جایگشتی Permutation matrix
ن
Disconnected
9
وارون پذیر
c
همبند Connected
ى
يال Edge

واژهنامهی انگلیسی به فارسی

A	بالابر Lift
خودریختی Automorphism	M
В	مدول Module
Bijection	N
C	نگاشت طبیعی
گروه دوری	O
D	یک به یک
Degree	P
${f E}$	گروه جایگشتی Permutation group
يال	Q
F	_
تابع Function	گراف خارجقسمتی Quotient graph
G	R
گروه	تحویل پذیر Reducible
Н	S
همریختیهمریختی	Sequence
I	T
Jnvariant	سرشت بدیهی Trivial character
L	U

واژهنامهی انگلیسی به فارسی

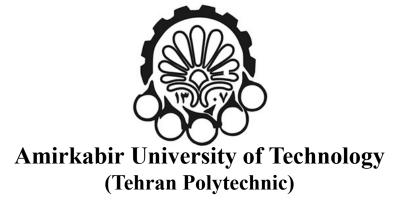
V

Abstract

This page is accurate translation from Persian abstract into English.

Key Words:

Write a 3 to 5 KeyWords is essential. Example: AUT, M.Sc., Ph. D,..



Department of ...

MSc Thesis

Title of Thesis

By

Name Surname

Supervisor

Dr.

Advisor

Dr.

Month & Year