

دانشکده مهندسی کامپیوتر

بررسی الگوریتمهای یادگیری ماشین برای جریانهای متنی

گزارش سمینار کارشناسی ارشد در رشته مهندسی کامپیوتر گرایش هوش مصنوعی و رباتیک

وحيد خرازي

استاد راهنما

دکتر بهروز مینایی بیدگلی

آبان ۹۵



پدر و مادر عزیزم که پیغمبر نگاهشان همیشه مرا چشم و چراغ خواهد بود.

قدرداني

سپاس خداوندگار حکیم را که با لطف بی کران خود، آدمی را زیور عقل آراست.

در آغاز وظیفه خود می دانم از زحمات بی دریغ استاد راهنمای خود، جناب آقای دکتر مینایی، صمیمانه تشکر و قدردانی کنم که قطعاً بدون راهنمایی های ارزنده ایشان، این مجموعه به انجام نمی رسید.

وحید خرازی آبان ۹۵

چکیده

این پایاننامه، به بحث در مورد نوشتن پروژه، پایاننامه و رساله با استفاده از کلاس IUST-Thesis میپردازد. حروف چینی پروژه کارشناسی، پایاننامه یا رساله یکی از موارد پرکاربرد استفاده از زیپرشین است. زیپرشین بسته ای است که به همت آقای وفا خلیقی آماده شده است و امکان حروف چینی فارسی در ۱ΔΤΕΧ و برای فارسی زبانان فراهم کرده است. از جمله مزایای لاتک آن است که در صورت وجود یک کلاس آماده برای حروف چینی یک سند خاص مانند یک پایاننامه، کاربر بدون درگیری با جزییات حروف چینی و صفحه آرایی می تواند سند خود را آماده نماید.

شاید با قالبهای لاتکی که برخی از مجلات برای مقالات خود عرضه میکنند مواجه شده باشید. اگر نظیر این کار در دانشگاههای مختلف برای اسناد متنوع آنها مانند پایاننامهها آماده شود، دانشجویان به جای وقت گذاشتن روی صفحهآرایی مطالب خود، روی محتوای متن خود تمرکز خواهند نمود. به علاوه با آشنایی با لاتک خواهند توانست از امکانات بسیار این نرمافزار جهت نمایش بهتر دستآوردهای خود استفاده کنند. به همین خاطر، یک کلاس با نام IUST-Thesis برای حروف چینی پروژهها، پایاننامهها و رسالههای دانشگاه علم و صنعت ایران با استفاده از نرمافزار زیپرشین، آماده شده است. این فایل به گونهای طراحی شده است که کلیات خواستههای مورد نیاز مدیریت تحصیلات تکمیلی دانشگاه علم و صنعت ایران را برآورده میکند و نیز، حروف چینی بسیاری از قسمتهای آن، به طور خود کار انجام می شود.

واژگان کلیدی: یادگیری ماشین، داده های جریانی، متن کاوی، پردازش متن

فهرست مطالب

رست تصاویر
رست جداول
رست الگوريتمها
رست علائم اختصاری
وست آ: مدیریت مراجع در لاتک ذ
آ_۱ مدیریت مراجع با BibT _E X
آ_۱_۱ سبکهای فعلی قابل استفاده در زیپرشین ر
آ-۱-۲ نحوه استفاده از سبکهای فارسی
وست ب: جدول، نمودار و الگوريتم در لاتک ش
ب_۱ مدلهای حرکت دوبعدی
ب_۲ ماتریس
ب_٣ الگوريتم با دستورات فارسي
ب_۴ الگوريتم با دستورات لاتين
ب_۵ نمودار
ب_۶ تصویر
ژهنامه فارس <i>ی</i> به انگلیسی

فهرست مطالب

واژهنامه انگلیسی به فارسی

فهرست تصاوير

س	•	•						•	•	•	•	•				as	a- 1	fa	ک	بب	ا س	، ب	جى	رو-	خ	نمونه	١	_Ĩ
ط																									,	دو شي	١_	ب

فهرست جداول

فهرست الگوريتمها

ض	•					•					ئرافي.	<i>ن</i> وموً	ں ہ	ريس	, مات	فمين	تخ	ای	I بر	D LT	نم آ	ورية	الگو	۱_	ب
ض										اف	هو موگر	بس	مات	ين	نخم	ای ن	[بر	RA	NS	SA(نہ 5	و ريت	الگو	۲_	ب

فهرست علائم اختصارى

a (m/s	s^2).	 	 	 	 			٠.		 •	 	 							 			٠ ر	ثر	راند	گر	ب	اب	ئىڌ	*
F(N))	 	 	 	 	 								 			 				 						9	,_	٠

پیوست آ

مديريت مراجع در لاتک

در بخش ؟؟ اشاره شد که با دستور bibitem میتوان یک مرجع را تعریف نمود و با فرمان ارجاع ادامه به آن ارجاع داد. این روش برای تعداد مراجع زیاد و تغییرات آنها مناسب نیست. در ادامه به صورت مختصر توضیحی در خصوص برنامه BibTeX که همراه با توزیعهای معروف تِک عرضه میشود و نحوه استفاده از آن در زیپرشین خواهیم داشت.

آ_۱ مدیریت مراجع با BibT_EX

یکی از روشهای قدرتمند و انعطافپذیر برای نوشتن مراجع مقالات و مدیریت مراجع در لاتک، استفاده از BibTeX است. روش کار با BibTeX به این صورت است که مجموعهی همهی مراجعی را که در پروژه/پایاننامه/رساله استفاده کرده یا خواهیم کرد، در پروندهی جداگانهای نوشته و به آن فایل در سند خودمان به صورت مناسب لینک می دهیم. کنفرانسها یا مجلههای گوناگون برای نوشتن مراجع، قالبها یا قراردادهای متفاوتی دارند که به آنها استیلهای مراجع گفته می شود. در این حالت به کمک استیلهای کنید. بیشتر توانست تنها با تغییر یک پارامتر در پروندهی ورودی خود، مراجع را مطابق قالب موردنظر تنظیم کنید. بیشتر مجلات و کنفرانسهای معتبر یک پرونده ی سبک (BibTeX Style) با پسوند bst در وبگاه خود می گذارند

به جز نوشتن مقالات این سبکها کمک بسیار خوبی برای تهیهی مستندات علمی همچون پایاننامههاست

که فرد می تواند هر قسمت از کارش را که نوشت مراجع مربوطه را به بانک مراجع خود اضافه نماید. با داشتن چنین بانکی از مراجع، وی خواهد توانست به راحتی یک یا چند ارجاع به مراجع و یا یک یا چند بخش را حذف یا اضافه نماید؛ مراجع به صورت خودکار مرتب شده و فقط مراجع ارجاع داده شده در قسمت کتابنامه خواهند آمد. قالب مراجع به صورت یکدست مطابق سبک داده شده بوده و نیازی نیست که کاربر درگیر قالب دهی به مراجع باشد. در این جا مجموعه سبکهای بسته Persian-bib که برای زی پرشین آماده شده اند به صورت مختصر معرفی شده و روش کار با آنها گفته می شود. برای اطلاع بیشتر به راهنمای بسته Persian-bib مراجعه فرمایید.

آ ـ ۱ ـ ۱ سبکهای فعلی قابل استفاده در زیپرشین

در حال حاضر فایلهای سبک زیر برای استفاده در زیپرشین آماده شدهاند:

unsrt-fa.bst این سبک متناظر با unsrt.bst میباشد. مراجع به ترتیب ارجاع در متن ظاهر میشوند.

plain-fa.bst این سبک متناظر با plain.bst میباشد. مراجع بر اساس نامخانوادگی نویسندگان، به ترتیب صعودی مرتب میشوند. همچنین ابتدا مراجع فارسی و سپس مراجع انگلیسی خواهند آمد.

acm-fa.bst این سبک متناظر با acm.bst میباشد. شبیه plain-fa.bst است. قالب مراجع کمی متفاوت است. اسامی نویسندگان انگلیسی با حروف بزرگ انگلیسی نمایش داده میشوند. (مراجع مرتب میشوند)

ieeetr-fa.bst این سبک متناظر با ieeetr.bst میباشد. (مراجع مرتب نمی شوند)

plainnat-fa.bst این سبک متناظر با plainnat.bst میباشد. نیاز به بستهٔ plainnat.bst دارد. (مراجع مرتب می شوند)

chicago-fa.bst این سبک متناظر با chicago.bst میباشد. نیاز به بستهٔ chicago.bst این سبک متناظر با میشوند)

asa-fa.bst این سبک متناظر با asa.bst میباشد. نیاز به بستهٔ datbib دارد. (مراجع مرتب میشوند)

با استفاده از استیلهای فوق می توانید به انواع مختلفی از مراجع فارسی و لاتین ارجاع دهید. به عنوان نمونه مرجع [؟] یک نمونه مجله فارسی است. مرجع [؟] یک نمونه مقاله مجله فارسی است. مرجع [؟] یک نمونه مقاله کنفرانس فارسی و مرجع [؟] یک نمونه کتاب فارسی با ذکر مترجمان و ویراستاران فارسی است. مرجع [؟] یک نمونه پروژه کارشناسی ارشد انگلیسی و [؟] هم یک نمونه متفرقه می باشند.

مراجع [؟، ؟] نمونه کتاب و مقاله انگلیسی هستند. استیل مورد استفاده در این پروژه/پایاننامه/رساله asa-fa است که خروجی سبک asa-fa در شکل میتوانید مشاهده کنید. نمونه خروجی سبک asa-fa در شکل آ_۱ آمده است.

آ_۱_۲ نحوه استفاده از سبکهای فارسی

برای استفاده از بیبتک باید مراجع خود را در یک فایل با پسوند bib ذخیره نمایید. یک فایل bib در واقع یک پایگاه داده از مراجع شماست که هر مرجع در آن به عنوان یک رکورد از این پایگاه داده با قالبی خاص ذخیره می شود. به هر رکورد یک مدخل گفته می شود. یک نمونه مدخل برای معرفی کتاب Digital Image در ادامه آمده است:

در مثال فوق، BOOK شخصه ی شروع یک مدخل مربوط به یک کتاب و BOOK برچسبی است که به این مرجع منتسب شده است. این برچسب بایستی یکتا باشد. برای آنکه فرد به راحتی بتواند برچسب مراجع خود را به خاطر بسپارد و حتی الامکان برچسبها متفاوت با هم باشند معمولاً از قوانین خاصی به این منظور استفاده می شود. یک قانون می تواند فامیل نویسنده ی اول + دورقم سال نشر + اولین کلمه ی عنوان اثر باشد. به AUTHOR و . . . و ADDRESS فیلدهای این مدخل گفته می شود؛ که هر یک با مقادیر مربوط به مرجع مقدار گرفته اند. ترتیب فیلدها مهم نیست.

انواع متنوعی از مدخلها برای اقسام مختلف مراجع همچون کتاب، مقالهی کنفرانس و مقالهی ژورنال وجود دارد که برخی فیلدهای آنها با هم متفاوت است. نام فیلدها بیانگر نوع اطلاعات آن میباشد. مثالهای ذکر شده در فایل MyReferences.bib کمک خوبی به شما خواهد بود. با استفاده از سبکهای فارسی آماده

¹Bibliography Database

^YEntry

شده، محتویات هر فیلد میتواند به فارسی نوشته شود، ترتیب مراجع و نحوه ی چینش فیلدهای هر مرجع را سبک مورد استفاده مشخص خواهد کرد.

نکته: بدون اعمال تنظیمات موردنیاز BibT_EX در TeXWorks، مراجع فارسی در استیلهایی که مراجع را به صورت مرتب شده چاپ میکنند، ترتیب کاملاً درستی نخواهند داشت. برای توضیحات بیشتر [؟] را ببینید یا به سایت پارسیلاتک مراجعه فرمایید. تنظیمات موردنیاز در TeXMaker اصلاح شده اعمال شدهاند.

برای درج مراجع خود لازم نیست نگران موارد فوق باشید. در فایل MyReferences.bib که همراه با این پروژه/پایاننامه/رساله هست، موارد مختلفی درج شده است و کافیست مراجع خود را جایگزین موارد مندرج در آن نمایید.

پس از قرار دادن مراجع خود، یک بار XeLaTeX را روی سند خود اجرا نمایید، سپس bibtex و پس از آن دوبار XeLaTeX را. در TeXMaker کلید F11 و در TeXWorks هم گزینهی BibTeX را موی سند شما اجرا میکنند.

برای بسیاری از مقالات لاتین حتی لازم نیست که مدخل مربوط به آنرا خودتان بنویسید. با جستجوی نام مقاله + کلمه bibtex در اینترنت سایتهای بسیاری همچون ACM و ScienceDirect را خواهید یافت که مدخل bibtex مربوط به مقاله شما را دارند و کافیست آنرا به انتهای فایل MyReferences اضافه کنید.

از هر یک از سبکهای Persian-bib میتوانید استفاده کنید، البته اگر از سه استیل آخر استفاده میکنید و مایلید که مراجع شما شماره بخورند باید بسته natbib را با گزینه numbers فراخوانی نمایید.

نمونه خروجی با استیل فارسی asa-fa برای BibTeX در زیپرشین

محمود امين طوسي

مرجع امیدعلی (۱۳۸۷) یک نمونه پروژه دکترا و مرجع واحدی (۱۳۸۷) یک نمونه مقاله مجله فارسی است. مرجع امینطوسی و دیگران (۱۳۸۷) یک نمونه مقاله کنفرانس فارسی و مرجع استالینگ (۱۳۸۰) یک نمونه کتاب فارسی با ذکر مترجمان و ویراستاران فارسی است. مرجع خلیقی (۱۳۸۷) یک نمونه پروژه کارشناسی ارشد انگلیسی و خلیقی (۱۳۸۷) هم یک نمونه متفرقه می باشند.

مرجع گنزالس و وودس (۲۰۰۶) یک نمونه کتاب لاتین است که از آنجا که دارای فیلد مرجع گنزالس و وودس (۲۰۰۶) یک نمونه کتاب لاتین است که از آنجا که دارای فیلد authorfa است، نام نویسندگان آن در استیلهای plainnat-fa ، asa-fa به فارسی نام دیده می شود. مرجع Kanade and Baker (۲۰۰۲) مقاله انگلیسی است که معادل فارسی نام نویسندگان آن ذکر نشده بوده است.

مراجع

استالینگ، ویلیام (۱۳۸۰)، اصول طراحی و ویژگیهای داخلی سیستم های عامل. ترجمهی صدیقی مشکنانی، محسن و پدرام، حسین، (ویراستار)برنجکوب، محمود، اصفهان: نشر شیخ بهایی، ویرایش سوم.

امیدعلی، مهدی (۱۳۸۲)، "تابع هیلبرت،" پایاننامه دکترا، دانشکده ریاضی، دانشگاه امیرکبیر.

امین طوسی، محمود، مزینی، ناصر، و فتحی، محمود (۱۳۸۷)، "افزایش وضوح ناحیه ای،" در چهاردهمین کنفرانس ملی سالانه انجمن کامپیوتر ایران، دانشگاه امیرکبیر، تهران، ایران، صفحات ۱۰۱-۸-۱.

خلیقی، وفا (۱۳۸۷)، "زیپرشین (XaPersian): بسته فارسی برای حروفچینی در ŁTEX2e،" HTTP://BITBUCKET.ORG/VAFA/XEPERSIAN.

واحدى، مصطفى (١٣٨٧)، "موضوعي جديد در هندسه محاسباتي،" مجله فارسي نمونه، ١، ٢٢-٣٠.

Baker, S. and Kanade, T. (2002), "Limits on Super-Resolution and How to Break Them," *IEEE Trans. Pattern Anal. Mach. Intell.*, 24, 1167–1183.

Gonzalez, R. C. and Woods, R. E. (2006), *Digital Image Processing*, Upper Saddle River, NJ, USA: Prentice-Hall, Inc., 3rd ed. .

Khalighi, V. (2007), "Category Theory," Master's thesis, Sydny Univ.

شكل آ_ ١: نمونه خروجي با سبك asa-fa

پیوست ب

جدول، نمودار و الگوریتم در لاتک

در این بخش نمونه مثالهایی از جدول، نمودار و الگوریتم در لاتک را خواهیم دید.

ب_۱ مدلهای حرکت دوبعدی

بسیاری از اوقات حرکت بین دو تصویر از یک صحنه با یکی از مدلهای پارامتری ذکر شده در جدول (ب۱) قابل مدل نمودن می باشد.

ب_۲ ماتریس

شناخته شده ترین روش تخمین ماتریس هوموگرافی الگوریتم تبدیل خطی مستقیم (DLT) است. فرض کنید چهار زوج نقطهٔ متناظر در دو تصویر در دست هستند، $\mathbf{x}_i' = H\mathbf{x}_i$ و تبدیل با رابطهٔ $\mathbf{x}_i' = H\mathbf{x}_i$ نشان داده می شود که در آن:

$$\mathbf{x}_i' = (x_i', y_i', w_i')^\top$$

^¹Direct Linear Transform

مدلهای تبدیل.	جدول ب_١:
---------------	-----------

توضيح	تبديل مختصات	درجه آزادی	نام مدل
انتقال دوبعدي	$x' = x + t_x$ $y' = y + t_y$	۲	انتقالى
انتقالی+دوران	$x' = x\cos\theta - y\sin\theta + t_x$ $y' = x\sin\theta + y\cos\theta + t_y$	٣	اقليدسى
اقليدسى+تغييرمقياس	$x' = sxcos\theta - sysin\theta + t_x$ $y' = sxsin\theta + sycos\theta + t_y$	۴	مشابهت
مشابهت+اریبشدگی	$x' = a_{11}x + a_{17}y + t_x$ $y' = a_{11}x + a_{17}y + t_y$	۶	آفين
آفین+keystone+chirping	$x' = (m_1 x + m_1 y + m_2)/D$ $y' = (m_2 x + m_2 y + m_2)/D$ $D = m_1 x + m_2 y + 1$	٨	پروجکتيو
حرکت آزاد	$x' = x + v_x(x, y)$ $y' = y + v_y(x, y)$	∞	شارنوري

$$H = \left[egin{array}{cccc} h_{1} & h_{7} & h_{7} \ h_{7} & h_{5} & h_{9} \ h_{7} & h_{A} & h_{9} \end{array}
ight]$$

رابطه زیر را برای الگوریتم (ب-۱) لازم دارم.

$$\begin{bmatrix} \cdot^{\top} & -w_i' \mathbf{x}_i^{\top} & y_i' \mathbf{x}_i^{\top} \\ w_i' \mathbf{x}_i & \cdot^{\top} & -x_i' \mathbf{x}_i^{\top} \\ -y_i' \mathbf{x}_i^{\top} & x_i' \mathbf{x}_i^{\top} & \cdot^{\top} \end{bmatrix} \begin{pmatrix} \mathbf{h}^{\mathsf{t}} \\ \mathbf{h}^{\mathsf{r}} \\ \mathbf{h}^{\mathsf{r}} \end{pmatrix} = \boldsymbol{\cdot} \tag{1----}$$

ب-۳ الگوریتم با دستورات فارسی

با مفروضات فوق، الگوریتم DLT به صورت نشان داده شده در الگوریتم (ب۱) خواهد بود.

الگوريتم بــ الگوريتم DLT براي تخمين ماتريس هوموگرافي.

 $\mathbf{x}_i \leftrightarrow \mathbf{x}_i'$ ورودی: ۴ $\mathbf{x}_i \leftrightarrow \mathbf{x}_i'$ ورودی: اوج نقطهٔ متناظر در دو تصویر

 $\mathbf{x}_i' = H\mathbf{x}_i$ ماتریس هوموگرافی H به نحویکه: $\mathbf{x}_i' = H\mathbf{x}_i$.

۱: برای هر زوج نقطهٔ متناظر $\mathbf{x}_i \leftrightarrow \mathbf{x}_i'$ ماتریس \mathbf{A}_i را با استفاده از رابطهٔ ب \mathbf{x}_i محاسبه کنید.

۲: ماتریسهای ۹ ستونی \mathbf{A}_i را در قالب یک ماتریس \mathbf{A} ۹ ستونی ترکیب کنید.

۳: تجزیهٔ مقادیر منفرد (SVD) ماتریس ${\bf A}$ را بدست آورید. بردار واحد متناظر با کمترین مقدار منفرد جواب ${\bf h}$

۴: ماتریس هوموگرافی H با تغییر شکل h حاصل خواهد شد.

الگوریتم ب_۲ الگوریتم RANSAC برای تخمین ماتریس هوموگرافی.

Require: $n \geq 4$ putative correspondences, number of estimations, N, distance threshold T_{dist} . **Ensure:** Set of inliers and Homography matrix H.

1: for k = 1 to N do

2: Randomly choose 4 correspondence,

3: Check whether these points are colinear, if so, redo the above step

4: Compute the homography H_{curr} by DLT algorithm from the 4 points pairs,

5: ...

6: end for

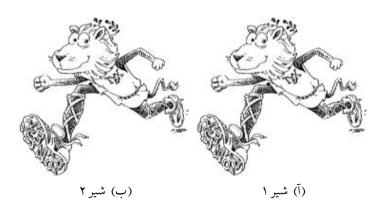
7: Refinement: re-estimate H from all the inliers using the DLT algorithm.

ب-۴ الگوریتم با دستورات لاتین

الگوريتم ب_٢ يك الگوريتم با دستورات لاتين است.

ب_۵ نمودار

لاتک بسته هایی با قابلیت های زیاد برای رسم انواع مختلف نمودارها دارد. مانند بسته های Tikz و PSTricks. توضیح اینها فراتر از این پیوست کوچک است. مثالهایی از رسم نمودار را در مجموعه پارسی لاتک خواهید یافت. توصیه می کنم که حتماً مثالهایی از برخی از آنها را ببینید. راهنمای همه آنها در تک لایو هست. نمونه مثالهایی از بسته Tikz را می توانید در /http://www.texample.net/tikz/examples ببینید.



شكل ب_١: دو شير

ب_ع تصوير

نمونه تصاویری در بخش قبل دیدیم. دو تصویر شیر کنار هم را هم در شکل ب_۱ مشاهده میکنید.

واژهنامه فارسی به انگلیسی

احتماليtrobabilistic
رزیابی
اندازه
پایدارtably
توپولوژی ضعیف
دامنه تو انی
فضای تابع
دامنه معناّتيي
قطعهبرنامه
مجموعه جزئاً مرتب كامل جهتدار
مر ت

واژهنامه انگلیسی به فارسی

Dcpo	مجموعه جزئاً مرتب كامل جهتدار .
Function Space	
Measure	اندازه
Ordered	مرتب
Powerdomain	دامنەتوانى
Probabilistic	احتمالي
Program Fragment	قطعەبرنامەقطعەبرنامە
Semantic Domain	دامنه معنایی
Stably	پايدار
Valuation	ارزیابی
Weak Topology	تو په له ژې ضعیف

Abstract:

-abstract

Keywords: Machine Learning, Data Stream, Text Mining, Text Processing



Iran University of Science and Technology Computer Engineering Department

A Survey on Learning Algorithms for Text Streams

A Thesis Submitted in Partial Fulfillment of the Requirement for the Degree of Master of Science in Computer Engineering

By:

Vahid Kharazi

Supervisor:

Dr. Behroz Minaei Bidgoli

Octobr 2016