

از نظر دقت از بین خطی درش بهتر است و مقدار اولی این کم دارد حتی اگر بنا سنجی به عنوان مقدار اولی انتخاب کنیم، مقدار  
تصادفی هم

اولی این را انتخاب می‌دهیم باز هم این روش حقیقتاً مسو در جواب خوبی می‌دهد.  
تصادفی  
نمونه: تولید عدد از mixture of gaussian ما در اینجا یک مقیسه را می‌بینیم، مثلاً چند مقیسه را انجام دهیم!

(k-means-hour.pdf) Data mining 92/08/06 8 حباب

این کتاب داده‌ها را به عکس مایه‌ها درش می‌بیند به فارسی است می‌تواند تهیه استفاده کنید مایه Tan است =  
کتاب داده‌های کاربردی دکتر محمد صنیعی زاده دانشگاه تربیت مدرس.

این حباب موضوع در مورد k-means صحبت کنیم، در واقع در مورد روش مای clustering که صحبت کردیم، ما حبابی  
پس یک اشاره این کردیم به clustering بر روش Model Based. در Model Based اقدام کنیم: این روش به خوبی  
روش مسئله‌های نیست ولی اگر داده‌های ما فوضی زیاد بود و داده‌ها به هم آمیخته بودند با سیستم‌های تقاربی می‌تواند  
ندارد و تقریباً به یک جواب می‌رسد اگر چه مسئله‌های از norm اولیه استفاده کنیم و اقدام می‌کنیم با چهار روش  
که کردیم حل کنیم البته این به جنبه‌های آموزشی است نه پژوهشی، می‌دانیم یک مقدار واقعاً ما می‌توانیم با چهار روش توان  
خطی‌تر حاصل کنیم: به روش دقیق - به روش ریاضی - به روش EM الگوریتم و نهایتاً با استفاده از شبکه‌های

این ما چهار روش این کار را می‌کنیم Model Based clustering انجام می‌دهیم از حوصله دانش درس ما خارج است. ولی ما باید بتوانیم  
با نرم افزار این کار را انجام دهیم، یعنی ما دارا فاضی خواهد بود به سبب اولی که برنامه‌نویسی در Model Based clustering انجام  
برود، ممکن است که قادر درش clustering هستیم کار انجام به سبب اولی که در درس داده‌های ما فاضی خواهیم این قدر

ریز نوع و در clemantine می‌توانیم این کار (Model Based clustering) را انجام دهیم. ولی سبب‌های راستی را نتوانیم  
که با سبب انجام دهیم. امروز ما خواهیم نوبت دوم clustering بر داریم البته انواع clustering بر روش نیست و به روش  
نمونه این را ما بهترین می‌توانیم: اول = مسئله‌های، دوم = دقیق، سوم = Model Based، چهارم = آسانی

یعنی k-means که گویا اولی است ولی این روش هم اگر که داده‌های ما از توزیع نرمال تبعیت کنند و norm ای  
که در نظر بگیریم، norm اولیه‌ها باشد توان این k میانگین و یا اندک فاصله این که در این k میانگین در نظر بگیریم فاصله  
اولی‌ها باشد، آن موقع جواب این روش هم یکی می‌شود ولی خوبی اینها این است که سرعتش از Model Based  
هم بیشتر است، بنابراین من هم است که خیلی کمتر یا اگر واقعاً در جدول نوع clustering صحبت می‌کنیم با میانگین  
شروع می‌کنیم، پس امروز ما می‌خواهیم روش سوم را بررسی کنیم، روش سوم که با سبب است.



و  $k$  میانگین یکی از روش های پارتیشن بندی است. چطور می تواند ما را در کلاسترینگ با  $k$  میانگین؟! مسکنه ترین روش  
 کلاسترینگ  $k$  میانگین است. سریع ترین روش کلاسترینگ هم می باشد و متاسفانه کم وقت ترین روش هم می باشد که البته  
 کم وقت ترینش شبیهی دارد حالا بعد از آن که می بینیم چه دارد؟! هدف ما در clustering چیست؟!  $n$  تا نقطه  
 دارند در فضای  $d$  می بینیم آنها را تقسیم می کنیم به  $k$  تا، وقت کنید که اگر  $k$  معلوم باشد، تقسیم کنیم به چقدر می شود؟

classification پس بحث clustering نیست. ولی clustering برای  $k$  معلوم خیلی خیلی پیچیده است به نظر

همین ما  $k$  را معلوم می گیریم و برای  $k$  های مختلف این را انجام می دهیم، بعد تقسیم می کنیم به  $k$  اینی انتخاب کنیم! گفتیم  
 که در روش ما محدوداً  $k$  را بین  $1$  تا  $k$  در نظر می گیریم و بعد به مسئله، حالا از یک هم بحث می شروع می کنیم، از یک جایی شروع  
 می کنند، فقط تقسیم می کنند که کدام cluster بهتر است با آن خیر روش که در سطح جایی گفتیم. حالا اگر  $C(x)$  کوچکترین

کلاستر چقدر  $x$  باشد که آن موقع  $k$  مرکز می باشد  $\min_{x \in X} \max \|x - C(x)\|$ ، یعنی ما دنبال یک کلاستر می گیریم که  
 کلاستران  $C(x)$  است که  $x$  های  $C(x)$  با  $x$  عامل تا حد تا مرکز کمترین باشد، و فاصله بین خود ما با مرکز کمترین

باشد، چطور همین  $\min \max$  قرار داده، اگر می توانیم این کار را انجام بدهیم به این می گویند  $k$ -center. بله اینطوری

between و within را بگویم، اگر بخواهم  $k$ -median را حساب کنیم، وقت کنید median چه بود؟

میان است، میان تو ای بعدی تعریف زده ولی میانگین در این تعریف شده، یعنی اگر ما به فرض صد تا عدد داشته باشیم  
 تو یک بعدی، چطور میانگین را پیدا کنیم، ما هم اعداد را مرتب می کنیم، وسطی را به عنوان میانگین در نظر می گیریم، این می  
 میانگین اگر اعداد ما دو بعدی باشد، ما هم می توانیم میانگین را در ۲ بعد تعریف کنیم، پس می بینیم که در اینجا بر خود می گذاریم، این

معدی  $k$ -median اول median این شیت که ما در آمار خوانیم، بزرگ آن است ولی آن نیست، چون اول در

حالت غیر بعدی اصلاً تعریف نکرده.  $k$ -median را تعریف می کنیم، این یک تعریف جدید است و آن هم همین (median)

قبلی نیست.  $k$ -median را تعریف می کنیم به صورتی که  $x$  های متعلق به  $x$ ،  $\text{norm}(x - C(x))$ ؛

$$k\text{-median} = \min \sum_{x \in X} \|x - C(x)\|$$

در قبلی  $\text{norm}$  و آن هم ماکزیمم این  $\text{norm}$  است، یعنی برای تمام  $x$  های متعلق به  $x$ ،  $x - C(x)$  را حساب می کنند

بعد ماکزیمم این را که حساب می کنیم، می بینیم این را در نظر می گیریم ولی در اینجا ما می آیم سافت می کنیم،  $x - C(x)$  را حساب می کنیم،

$\text{norm}$  را حساب می کنیم بعد سافت می کنیم  $\text{norm}$  را می گیریم، بعد می بینیم آن را حساب می کنیم. هنوز آنقدریم که چطور این

کار را انجام بدهیم تقسیم. اول نقطه وسطی را چگونه پیدا می کنیم؟ به صورت Random چون می تواند نقطه وسطی باشد (۵)

[More on www.faruk.ba](#)[Article Top](#)

## Comments and Discussions

Add a Comment or Question

Search this forum

Go

[Profile popups](#)[Spacing](#)[Relaxed](#)[Noise](#)[Medium](#)[Layout](#)[Normal](#)[Per page](#)[25](#)[Update](#)

There are no messages in this forum

[Permalink](#) | [Advertise](#) | [Privacy](#) | [Mobile](#)  
 Web01 | 2.7.1310022.1 | Last Updated 13 Oct 2013

Layout: fixed / fluid

Article Copyright 2013 by Faruk Basi  
 Everything else Copyright © CodeProject, 1999-2013  
[Terms of Use](#)

توجه کنید در  $k$ -median ما  $\min \max$  را در نظر نمی گیریم، به جای این  $\min \max$  بگیریم می آید  $\min$  می بینیم  
 ما می بینیم راضی گیریم و هوایمان باشد که این اسمش  $k$ -median است و در چند پرسش سوال آن را تعریف کرد.

$k$ -Median Square:  $\min \sum_{x \in X} \|x - C(x)\|^2$  می توان ۲ صریحانه  $k$ -Median Square

قبلی می بینیم که پر زور بود و بعدی آن را دهی برت داشتیم خیلی حساس بود در این یکی حساس است نسبت به داده های  
 برت. و بعضی وقتها از این استفاده می کنند البته معمولاً قبلی بهتر است. حالا  $k$  میانیست چه طوری است؟!

مثلاً: ما میخواهیم این داده ها را به ۳ گروه تقسیم کنیم، به ما می گویند  $k=3$  است و باید به ۳ تا از داده ها را به  
 تعداد انتخاب کنید، این اولین نکته می باشد که می بینیم در  $k$  میانیست است، یعنی چگونه می خواهیم روش یک روشی باشد  
 می بینیم بتواند به این بگوید، تا جایی که می بینیم این به عدد داریم و انتخاب کنیم، گفتیم این روش دست می بینیم  
 به راه ترین راه این است که می بینیم به صورت تعداد انتخاب می کنیم، این دو تا ایراد دارد نوی  $k$  میانیست: اگر  
 این به عدد معلوم از این بار در یک خوشه انتخاب شوند،  $k$  میانیست به بدترین جواب ممکن می رسد مناسب و ایراد  
 بعدی این روش هم این است که هر دفعه ما این روش را اجرا کنیم، باید جواب می رسد یعنی ممکن است ما دو بار انجام  
 این کار را انجام دهیم و دو تا جواب مختلف می رسد اگر چه

صفحه 7 از 7  
 11/04/2013 01:04



در حالی که در سلسله مراتبی همچین چیزی امکان ندارد، نوی مدل پس امکان دارد، اما امکانی خیلی خیلی ضعیف است اما نوی K  
میائین احتمال می دهیم ضعیف نیست، پس عدد تصادفی بد مثل همان را حل کنید، روش را توضیح می دهم، که به یک میائین  
برای آن انجام دهیم، ولی دوتا ایراد قبل را دارد. ایراد اول قابل رفع است، یعنی ما میائینم اگر داده های را که انتخاب کردیم،  
فصل شان از یک حدی کمتر بود، یکم داده های تصادفی را دوباره انتخاب کنیم، این کار را ما میائیم اول بعضی برنام هان K  
میائین انجام می دهیم ولی در مورد این که اگر حجم در خوشه های تصادفی انتخاب شدند، ولی امکان داره باشد که جوابهای K  
میائین ما، بی نباشد، باز هم وجود دارد. اینه این چه فایده ای دارد این است که ما میائیم لا نوی روش هان پیچیده میائیم  
یک خوشه بندی داریم. K میائین بهترین روش خوشه بندی داریم است یعنی خوشی K میائین می و درون خوشه های بد  
ما. پس این نیست که K میائین روش بد باشد، بلکه آنها بد نیست، اگر که این مقدار اولیم را خوب انتخاب کنیم، K میائین  
خوب. روشی که K میائین حساب میائیم را انجام می دهند بر اساس حل نوی است، یعنی یک روش بنی ریاضی را داریم  
است. K میائین هم با دست راحت نیست انجامش دهیم. پس اولین قدم در نوی K میائین انتخاب K نقطه به تصادف  
است. بعد میائیم؟ بعد میائیم؟ فاصله ی تمام نقاط را تا آن سه مرکزی که انتخاب کردیم حساب میائیم. فاصله هر نقطه که  
هر مرکزی نزدیک تر بود (36) آن نقطه را جزو آن خوشه در نظر میگیریم، پس خیلی راحت است، تنها سوالی که پیش می آید این  
نقاطی که به تصادف انتخاب کردیم، مرکزی شدند؟ آن مرکزی شدند؟ اولی شده را ما اینجا اعمال کردیم!! بله، اما فعلاً مرکزی اولیم  
را پیدا کردیم این بهترین مرکزی نیست، باید بیاییم الان بهترین مرکزی را حساب کنیم (فعلاً اولیم می و در خوشه میائیم فاصله ها  
را حساب کنیم) حالا میائیم بهترین مرکزی را پیدا میائیم، هدف کنید که لازم نیست مرکزی از خود داده ها باشد، مثل میائین است  
فعلاً درباری آنس آنرا تنها یکبار خواهم تا س را بنویسم احتمال 3.5 توان ما برنده میائیم، یعنی نقطه را آیدیم حساب کردیم.  
3.5 در اعداد تا س نیست ولی به طور متوسط ما 3.5 توان برنده میائیم اگر درباری تا س شرکت کنید. اینجا طم حسی است  
الان ما دوتا داده داریم در خوشه ولی مرکز مشخص نبود، پس ما بیاییم مرکزی اصلی شان را حساب کنیم، پس مرکزی اصلی  
داده ها، اول اعداد تصادفی ما بود، ~~خوب این مرکزی~~ یک مرکزی بی روی هستند، پس میائیم دوباره مرکزی  
را حساب میائیم، حالا دوباره میائیم فاصله تمام نقاط را تا مرکزی حساب میائیم، فاصله هر نقطه تا مرکزی میائیم بود، آن  
را جزو همان خوشه میائیم، ممکن است یک داده که همان طور میائیم از خوشه دوم یک داده می و رفت تو خوشه اول و  
یک داده از خوشه سوم می و رفت تو خوشه دوم، چون فاصله را با مرکزی جوید حساب کردیم. (صفحه ۱۰۹) (۷)



حالا دوباره می آیم مرکز جدید را معیار حساب میکنیم، مرکز جدیدی که الان تشکیل شده، دوباره مرکز عوض میشود، حالا دوباره  
می آیم فاصله ی تمام نقاط را تا مرکز جدیدی که پیدا شده، حساب میکنیم تا ببینیم چه تغییری رخ می دهد، کنتورهای بعدی را رخ  
نمی دهد؟! این کار را تا کی ادامه می دهیم؟! تا وقتی که مراکز تغییر نکنند، این هم تدریس کلاسیک است. پس مراکز را  
تنها به صورت تعدادی انتخاب کردیم و مراکز بعدی را به صورت کلیل بدست می آوریم. حالا بگذاریم این الگوریتم را بتوانیم با  
کامپیوتر پیاده کنیم با توجه به فاصله ی که در اختیار داریم، برنامه نویسی ما چند تفاوت میشود. سوال؟! آمارش های دیگر هم داریم؟! سوال

بیشتر مای خواص به صورت unsupervised کار کنیم، اگر بخواهیم supervise کار کنیم، باید می توانیم همیشه کمر داده بداریم  
رسم کنیم، خودمان یک threshold (حد آستانه) تعیین می کنیم، وسط را می بریم به عنوان مراکز، بعد هم می توانیم این کار  
را انجام بدهیم. اسلایدهای (تیرجفت های تخصصی) داده انجام می دهیم و بدو هم نیست ببینیم، مثلاً گفته چقدر این الگوریتم سریع است؟  
گفته خیلی سریع است، یعنی بعد از iteration های خیلی کمی همگرا می شود و دقیقاً هم همین جور است و درست است  
image processing یک ایده توی ما iteration ها می خورد. از این بعد بیشتر به جنبه های ریاضی درگیر می شدیم که  
ما واقعاً اینجا کار می نداریم، پس ما تا همان اسلاید 12 را گفتیم می کند برآیند. در SPSS تنها فرض این الگوریتم نشان  
دارد و می شود مراحل نشان داده نمی شود.

ارباب و بیعت 8 clst  
میخواهیم در مورد clustering کمی بیشتر صحبت کنیم، دلیل این هم این است که ممکن است داده های که ما با آنها سر و کار داریم،  
صیغ داده های پیوسته نباشد و بعد برای clustering این جور داده ها نتوانیم از روش های معمول استفاده کنیم. اول  
بحث clustering را مطرح کرده، که ما اینجا بحث آشنایی داریم، پس کاربردهایش را گفته که اینم قبلاً گفتیم، چند مثال  
هم برآیند زده و cluster خوب را تعریف کرده، مثل چیزی خوب، و غیره می گویم که گاهی گاهی خوب است که داخل کلاس ها  
و آرایش ها کم باشد بین کلاس ها و آرایش ها زیاد باشد و کیفیت کلاس ها بد باشد هم به measure این داریم بقیه  
که هم در روش ما می باشد measure این توی روش ما داریم و روش های مختلف هم داریم. یک measure را  
میتوان در روش های مختلف به کاربرد. از این دو تا را خوب انتخاب کرده با هم مقایسه می کنیم جواب را خواهیم گرفت که البته اینجا  
داریم هم داریم گفته کیفیت clustering خوب هم اول clustering این خوب است که بتواند الگوهای پنهان را بران ما  
کشف کند و واقعاً در داده ها می بینیم است یعنی یک وقتها می بینیم عنوان سال استان ها می شود و اوقتی که cluster 8



میکنیم، مثلاً بر اساس خیلی جواب باشد، تهران با چ استانی مثلاً هم گروه میشود؟ با چنین شهر و عظیم بوشهر، در صورتی که  
به نظر ما کافی که در ضمن علوم، به شهر و ضمن شهر در امتحان هیچ شباهتی بین این دو شهر تهران نمی بینند ولی یک چیز  
که در این سه تا خیلی پررنگ است، در آمد زیاد این سه شهر است یعنی از تمام شهرهای کشور این سه شهر در آمدشان بیشتر است.  
چنانچه همین هم است که در روش clustering انجام می دهیم اینها را هم اندازه دیک خوشه براساس معیارهای مختلفی میجوید و در از این معیار

معیارهای فرقی، تفاوت، آمیزش و... و در آمد هم جزو این است ولی در آمد آنقدر کم است که نمیتوانیم معیارهای  
Dominate می کند (نکته: ما بیشتر خود را می بینیم). ما در روش clustering اینی که میخوانیم از این (همه باید این ۹ تا ویژگی <sup>منقسم 7</sup>)  
را بتواند توان خودش داشته باشد، یعنی اگر scale داده ها را عوض کنیم در کلاسهایشان تأثیری نداشته باشد، تفسیر رخ  
ندهد، Attribute های مختلفی اگر در بر داشته باشیم، بتواند کلاستریک کند، اگر کلاسهای ما اشکال مختلفی داشته  
باشند همچون به عنوان مثال آبراهه های ماه (عیب دیگر کامیابین)، چون در سری فراتر داده ~~صعب~~ صعبت میکند، آنرا خوشه  
های خاص صورت ببیند و دایره نباشد، کامیابین افصح، بارش آورد، کافی است به شکل کامیابین، این را را

دایره حساب میکند، خیلی از داده های دیگر را هم ماطل این داده ها میکند، پارامترهای ورودی اش نباید زیاد باشد، باید بتواند  
outliers یعنی اختساش و داده های پرت را بتواند جدا کند، کار کند، نسبت به داده های ورودی حساسیت زیادی  
نداشته باشد، با ابعاد بالا بتواند کار کند، اگر شرطی، معدر رسی توسط کاربرمان اعمال شده بتواند آنها را در نظر بگیرد و آخرش  
هم قابل تفسیر و استقاده باشد چون ممکن است یک cluster اینی که ما انجام دهیم، قابل تفسیر نباشد، ساختار داده های ما  
در clustering بصورت یک ماتریس بود، ما  $p$  تا تفسیر داریم،  $n$  تا عدد داریم و Dissimilarity matrix، آن را

هم همان طوری که در صفا ۹ می بینید درست میکنیم، چطور این quality این clustering مان را اندازه گیری میکنیم؟!  
خوب تفسیر با  $D$  نشان می دهیم که حاصل از این معادله را به ماتریس  $D$  می دهیم،  $D_{ij} = 0$  می شود یعنی فاصله از هم  
دستر ندارند. تنهایی که ما به توجع کنیم این است که این  $D$  به نوع داده ها همان باید تعریف مناسبی داشته باشد، یعنی  
آورد داده ها مثلاً با نیزی باشد یا Categorical باشد یعنی تعداد باشد یا ordinal باشد که این  $D$  اینی که تعریف  
میکنیم مقارنت است با  $D$  اینی که مثلاً در داده های پیوسته که به صورت آماری باشند، این چهار مورد دایره ای که ممکن

است ما با هاشی برخورد کنیم را اینجا گفته (صفا ۱۱) یعنی مقایسه شان کنیم یا حاصل این باشند - با نیزی باشند -  
⑨



$$S_f = \frac{1}{n} (|x_1 f - m_f| + |x_2 f - m_f| + \dots + |x_n f - m_f|)$$

$$S_f = \frac{1}{n} (|x_1 f - m_f| + |x_2 f - m_f| + \dots + |x_n f - m_f|)$$

$$z_{if} = \frac{x_{if} - \mu_f}{s_f}$$

$$m_f = \frac{1}{n} (x_1 f + x_2 f + \dots + x_n f)$$

- مقایسه با الگوریتم Categorical باشند یعنی  $a, b, c$  و  $d$  تعداد اند (صفحه 15)  
 این تعداد را ما مستقیم به دست آوریم و حالا استفاده کنیم تعریف میکنیم که

	1	0	sum
1	a	b	a+b
0	c	d	c+d
sum	a+c	b+d	p

$$d(i, j) = \frac{b+c}{a+b+c+d} \quad \text{symmetric binary variable} \quad \text{simple matching coefficient}$$

•  $d(i, j) = \frac{b+c}{a+b+c}$  is not invariant

این مثال زده که این مثال برای آن جدید است (مقدم 16) نیست که یک صفت *symmetric* است ولی به نظر می‌رسد

symmetric ست، نه نامتبر دارم! اینجا 7 نامتبر داریم مثلاً  $d(jack, mary)$  را هم فقط اطمینان حساب کنیم

برای کدام متغیر خواهیم حساب یعنی؟!  $d = 0.33$   $P$  فرض کنید ما به جای  $Y$  و  $P$  را بنویسیم  $\frac{1}{N}$

اینجا هم ما هم کار را انجام می دهیم.  $d = \frac{b+c}{a+b+c}$  و این را هم symetric بنامید.  $d = \frac{b+c}{a+b+c}$



چون مقدار اولی symmetric است، جنبش را آنجا هم بگذاریم کنار، با این مسئله را به جدول ۲ بعدی تبدیل می‌کنیم. وقت کش این جدول ۲ بعدی است ولی آن مقدار دارد، آنجا هم ۲ بعدی است ولی مقدارشان زیادتر است. ما یک قشر را می‌توانیم به عنوان سطر بگیریم، بقیه را به عنوان جدول ۲ بعدی درش بگذاریم و آنجا هم جدول

۲ بعدی درش بگذاریم که فراوانی‌های آن را بنویسیم، آن موقع می‌تواند این را حساب کند. (منظور این را به عنوان تمرین

بگذاریم برای حل بعدی) جدولی بعدی به آن داده داریم که اسامی ما را درباره بتواند در خودشان حساب کند.

آن قشر ما اسامی باشد، اول موقع فاصله را به این صورت حساب می‌کنیم: تعداد اونهایی که با  $P$  برابرند، تقسیم بر کل آن

$$m \text{ رتبه } \rightarrow \text{تعداد اونهایی که با } m \text{ برابرند قبل وقت که داریم احتمال} \quad \alpha(\vec{r}, \vec{r}_i) = \frac{p - m}{p}$$

حساب می‌کنیم، یک هم چنین حالتی دارد، پس آن قشر ما اسمی بود، وقت فرغ از هم برای قشر اسامی فاصله را حساب می‌کنیم، فاصله

یک تغییر می‌شود، احتمال دارد ~~آن قشر~~ آن قشر ما بتویسد باشد، رتبه می‌باشد، آن رتبه این باشد، دوباره برای معادله فاصله

اول یک تغییر قشر می‌دهیم، پس از روش‌های معمول برای فاصله حساب می‌کنیم و

$$r_i \neq \in \{1, \dots, m\} \quad z_i \neq = \frac{r_i - 1}{m - 1}$$

پس آن قشر ما تیرگی بود، ما می‌آیم بر حسب آنکه مرتبه‌ش می‌کنیم، تیرگی که هست یعنی مثلا  $a, b, c, d, e, f, g$

اینها کار را بهتر می‌کنند ولی ترتیب دارند، یعنی از نظر فاصله برای sort کردن اینها باید مرتب باشند، ما می‌توانیم با

قشرهای  $a$  و  $b$  با خوردن کار کنیم و به جدول بگذاریم  $n$ ، با rank این کار می‌کنیم، می‌توانیم Rank را این

است در حرف انگلیسی،  $p$  رنگش ۲ در حرف انگلیسی و ... علامت می‌آیم این اعداد را که به حرف نسبت داریم،

بر اساس این اعداد که به نسبت داریم، این  $z$  را درست می‌کنیم، حالا این  $z$  این که درست می‌کنیم برای

rank صفای ۱ بخش بر  $m - 1$  که  $m - 1$  بزرگترین rank می‌باشد. این  $z$  که درست آوریم،  $z$

این کار می‌آیم با روش معمول فاصله این را حساب می‌کنیم. آن قشرها ~~و~~ نسبتی باشد، اول موقع فاصله می‌توانیم با

یک قشر قشر با هم با هاش کار کنیم البته تا اول جایی که می‌توانیم سعی می‌کنیم قشر قشر استفاده نکنیم ولی خوبی

تفسیر قشر این است که مثلا  $g$  که  $g$  که بگیم، برای داده‌هایی که ضلع بزرگ هستند را ضلع کوچک می‌کنند، داده‌ها

که یک را کمتر کوچک می‌کنند، داده‌ها حالت نزول به آن کشیده وقت که داده‌ها نزول می‌کنند، یکسره های خوش بین ضلع

مقدار عمل می‌کنند. بعد از خوش می‌آیم را انجام داریم، معادله دوباره، یک تبدیل مقدار می‌زنیم روی داده‌ها ~~و~~ ۱۱



به داده های اولیه برنگردیم، در حالت معمولی سعی میکنیم ماهیت و وقت این کار را انجام ندهیم هر دو داده ما تبدیل نزنیم اگر  
 داده های دیگری داشته باشیم این کار را انجام می دهیم، ولی هرگز سعی نمیکنیم این کار را انجام می دهیم یعنی وقت  
 می خواهد خوش نبندی انجام دهیم یک تبدیل روی داده حاصلی زنده مثلاً می یاریم تو حاصلی صفر و یک که نتوانند  
 خطی راحت تر با هم کار کنند و هم در برنام نویسی ~~سخت تر~~ است و هم با  $overflow$  و  $underflow$  از این چیزها  
 حاصلی کمتر می بینیم. اگر فرضی که داریم فوعلی قاطع باشد یعنی چون  $multivariate$  هسته، یکی این مثلاً

عددی باشد، این حاصل این باشد ما این را معمولاً از این فرمول استفاده میکنیم (صفحه ۷۲)  

$$\alpha(i, j) = \frac{\sum_{f=1}^P \alpha(i, j, f)}{\sum_{f=1}^P 1}$$
 مثلاً ما همان قبلی ها را با  $\alpha$  بگیریم، فقط توی ذهنمان  
 باشد که  $\alpha$  را حاصل بردن هم می باشد در باره اینها را خوش نبندی کرد.

انواع روش های خوش نبندی را  $\alpha$  گفته که قبلاً در موردش صحبت کردم  $Partitioning$  Methods همان  $k$   
 میگویند است که می خواهیم طبق بعدی در موردش صحبت کنیم.  $Partitioning$  داده = یعنی یک  $k$  میگویند  
 با هر نرم افزار که بلدی باید یک  $k$  داده را خوش نبندی کنند.

TIME SERIES DATA

MINING - 8data series.ppt 92, 8, 13

صفحه ۹