

روش‌های یادگیری ماشین در پردازش زبان طبیعی

دسته‌بندی متون

هادی ویسی

h.veisi@ut.ac.ir

دانشگاه تهران - دانشکده علوم و فنون نوین



فهرست

○ دسته بندی متون

- کاربردها

○ بیز ساده (Naïve Bayes)

- آموزش: تخمین پارامترهای مدل
- مثال
- ویژگی‌ها (ارتباط با مدل زبانی)
- نکات کاربردی

دسته‌بندی متون: کاربردها ...

تشخیص عنوان (موضوع)

Türkçe | Español | Русский | اردو | ۲۰ آبان ۱۳۹۴

خبرگزاری جمهوری اسلامی

صفحه اصلی | سیاسی | اقتصادی | اجتماعی | فرهنگی | علمی | ورزشی | بین الملل | استان ها | پژوهش | حوادث | عکس | فیلم

آزمایشگاه خبری | هجوم نظامیان آل خلیفه به عزاداران عاشورایی - 52 دقیقه پیش

کد خبر: ۸۱۸۱۰۴۹۸ (۴۹۴۵۴۶۰) | تاریخ خبر: ۲۰/۰۸/۱۳۹۴ | ساعت: ۱۴:۴۲

نسخه چاپی | ارسال به دوستان

🏠 📄 📧 📧 📧 📧

کتاب نذری برای دومین بار در اهواز توزیع شد

اهواز- ایرنا- با همت فعالان اجتماعی و به منظور سوق دادن مردم به کتاب و کتابخوانی برنامه «نذر کتاب» برای دومین بار در شب های محرم در اهواز برگزار شد.

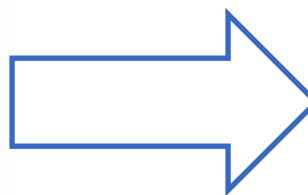


رئیس هیات مدیره کانون آموزش شهروندی برای توسعه (کاشت) در گفت وگو با ایرنا اظهار کرد: برنامه نذر کتاب در شب نهم ماه محرم برای دومین بار در طول دهه اول محرم در اهواز برگزار شد.

فهمیه صیادی افزود: این کانون دومین برنامه خود را در ساعت ۲۱ منطقه حصیر آباد اهواز برگزار کرد. به گفته وی، اولین برنامه نذر کتاب در منطقه کوی نفت اهواز در شب هفتم محرم انجام شد.

وی هدف اولیه این کانون فرهنگی را ترویج نذر فرهنگی عنوان کرد و گفت: ترویج نذر فرهنگی برای این است که به جامعه آموزش داده شود که به جای غذای نذری می توانند کتاب یا کالاهای فرهنگی دیگر را جایگزین کنند.

صیادی با بیان اینکه این کانون پیش از این برنامه اهدای کتاب در روستاهای شادگان و ایذه را اجرا کرده است، گفت: نذر کتاب در محرم و در هیات های عزاداری برای نخستین بار انجام می شود. وی اظهار کرد: از علاقه مندان به مشارکت در این برنامه ۲۴ میلیون ریال جمع آوری شد که بخشی از آن در برنامه های نذر کتاب توزیع شده و مابقی در برنامه های آتی این کانون توزیع خواهد شد.



- سیاسی
- مذهبی
- علمی
- هنری
- ورزشی
- ...



دسته بندی متون: کاربردها ...

○ عقیده کاوی (دسته بندی نظرات کاربران)

● موافق / مخالف

- | | | |
|------------|---|---|
| + 1 | <p>امید دلیر (1393/1/24) :</p> <p>فوق العادس</p> |  ● |
| + 0 | <p>محمود ابراهیمی جاویدی (1393/1/22) :</p> <p>من این لب تاب 5 ماه دارم و از همه نظر عالی متل: کیفیت صفحه نمایش با کیفیت لولاش برای لمسی بودنش و...</p> |  ● |
| + 1 | <p>مرتضی اکبری (1393/1/17) :</p> <p>یکی از مشکلات از نظر من کارت گرافیکشه. منظورم باس کارت گرافیکه!</p> |  ● |
| + 0 | <p>هادی کریم (1392/12/22) :</p> <p>سلام..اصلا تو خریدنش شک نکنید.من خریدم خیلی عالی.خیلی بیشتر از عکسش خوشگله.واقعا زیباست.اصلا لولایی دیده نمیشه واستحکامش خیلی بیشتره.به قیمتش می ارزه.صدای فنش هم اصلا اذیت نمیکه.جون میده برای بازی</p> |  ● |
| + 5 | <p>ابوذر هومن (1392/12/17) :</p> <p>با سلام
بعد از 4ماه کار مداوم
از هر نظر عالی
قدرت خوبی داره.بسیار سریع
تاج خیلی خوبی داره.کیفیت واقعا عالی.صفحه عجیبی داره که اصلا خسته نمیشن موقع کار باهاش
وزنش هم که عالی.به همراه همیشگی میشه براتون
مرسی</p> |  ● |
| + 0 | <p>هادی کریم (1392/12/1) :</p> <p>سلام..اصلا تو خریدنش شک نکنید.من خریدم خیلی عالی.خیلی بیشتر از عکسش خوشگله.واقعا زیباست.اصلا لولایی دیده نمیشه واستحکامش خیلی بیشتره.به قیمتش می ارزه.صدای فنش هم اصلا اذیت نمیکه.جون میده برای بازی</p> |  ● |



دسته بندی متون: کاربردها ...

○ تشخیص هرزنامه‌ها

- ایمیل دریافتی هرزنامه است یا نه؟

Fri 10/16/2015 6:29 PM

maair@host.maail.ir on behalf of <info@websit.ir> فروشگاه فایل رایگان بسازید

درآمد میلیونی از فروش فایل

To veisi@ce.sharif.edu

فایل برای فروش دارید؟

فقط کافیست در سایت عضو شوید و بدون داشتن هیچگونه دانش فنی هر فایلی که فکر میکنید فروش خوبی دارد آپلود کنید و بفروش برسانید
گاهی دیده شده هر فایل در روز چندین بار بفروش برسد

فضا و امکانات سایت نامحدود است و میتوانید بدون صرف هزینه های زیاد برای راه اندازی یک سایت فروش فایل به رایگان با سرعت بالا صاحب یک سایت با پنل مدیریت کاملا فارسی شوید

تمام هزینه های سایت رایگان است و در صورت فروش هر فایل مبلغ اندکی از سود شما برداشت می شود

فقط یک قدم تا درآمد عالی چندین میلیون در ماه فاصله دارید



دسته بندی متون: کاربردها ...

○ تشخیص نویسنده/جنسیت نویسنده

• نویسنده متن کیست؟ زن است یا مرد؟

چراغهای امامزاده، از دور در تاریکی؛ مثل چراغ خانه ای بود که تو را میخواهد. گرم، روشن و منتظر. سرم را به ضریح چسباندم. سلام آقا. دوستش دارم از بین این همه آدم ، فقط اون! شاید بچه گیام فقط برای ظاهرش بود ،اما روزی که به خاطر من ، دعوا کرد،دیدم جوونمرده. مثل قهرمونای قصه. وقتی منو سر مزار دوستش برد و گریه کرد، دیدم مهربونه.دل داره؛ و پاکه، مگه آدم چند بار میتونه دلشو هدیه بده؟ من هیچوقت روم نشده از خدا چیزی بخوام.اما این بار میخوام ! عمر در برابر عمر! از من نگیرش خدا! چیزی ندارم بت بدم،جز عشقی که خودت تو قلبم گذاشتی. پیرزن بخش زنانه گفت: تو اتاقشه. اما گفته کسی رو نمیبینه! گفتم:بگو دخترت اومده ! پشت در اتاقش بودم.از اتاقهای کوچک اجاره ای

چهارشنبه آخر ماه پیش وقتی از مسیر دور تر می رفت، سر یک کوچه ترمز کرد نگاهی به این طرف و آن طرف انداخت، بعد گفت: «بخشید الان برمی گردم» و از ماشین پیاده شد. دوباره کمی این طرف و آن طرف را نگاه کرد، یک کوچه را تا نیمه رفت و برگشت بعد سوار شد و رفتیم. به دست هایش نگاه کردم، فرمان را آنقدر محکم گرفته بود که ترسیدم از جا کنده شود، اما لرزش دست هایش پیدا بود، پرسیدم «حالتون خوبه؟» گفت «نه.» نگاهش کردم و بعد برایم تعریف کرد. چهل و شش سال پیش عاشق دختر جوانی می شود. چهارشنبه آخر یک ماه دختر جوان به او می گوید خانواده اش اجازه نمی دهند با او ازدواج کند. راننده از دختر جوان می خواهد لاقبل ماهی یک بار او را از دور ببیند. دختر جوان قول می دهد تا آخر عمر چهارشنبه آخر هر ماه سر این کوچه بیاید. چهل و شش سال دختر جوان چهارشنبه آخر هر ماه سر کوچه آمده، راننده او را از دور دیده و رفته است. از راننده پرسیدم «دختر جوان ازدواج کرد؟» نمی دانست. پرسیدم «آدرسشو دارین؟» نداشت. در این چهل و شش سال با او حتی یک کلمه هم حرف نزده بود فقط چهارشنبه های آخر هر ماه دختر جوان را دیده بود و رفته بود. راننده گفت «چهل و شش سال چهارشنبه آخر هر ماه اومد، ولی



دسته‌بندی متون: کاربردها

تشخیص زبان متن

• فارسی

• کردی

• عربی

زبان‌شناسی علمی است که به مطالعه و بررسی روشمند زبان می‌پردازد. در واقع، زبان‌شناسی می‌کوشد تا به پرسش‌هایی بنیادین همچون «زبان چیست؟»، «زبان چگونه عمل می‌کند و از چه ساخت‌هایی تشکیل شده‌است؟»، «انسان‌ها چگونه با یکدیگر ارتباط برقرار می‌کنند؟»، «زبان آدمی با سامانه ارتباطی دیگر جانوران چه تفاوتی دارد؟»، «کودک چگونه سخن گفتن می‌آموزد؟»، «زبان بشر چگونه تکامل یافته‌است؟»، «زبان‌ها چه فراتنی با یکدیگر دارند؟»، «ویژگی‌های مشترک زبان‌های جهان کدامند؟»، «انسان چگونه می‌نویسد و از چه راهی زبان ناوشتاری را واکاوی

زمانه‌وانی ، زمانناسی یا (به ینگیسی Linguistics) بریتیه له زانستی لئکولینه‌وهی ورد و همه‌لایه‌هی زمان. به وانایه‌کی تر زمانه‌وانی به دواى دؤزینه‌وهی نهم پرسیارانه‌دایه: «زمان چیه؟»، «چون کار ده‌کا و له چی ینکھانوهه؟»، «مروفه‌کان چون په‌یوه‌ندی به یه‌که‌وه ده‌که‌ن؟»، «زمانی مروف له چیدا له که‌ل سیسته‌مه ینوه‌ندیه‌کانی گیانداران حیاوازه؟»، «مندال چون زمان ده‌گری؟»، «زمانی مروف چون گه‌شهی سه‌ندوهه؟»، «خزمایه‌تی ینوان زمانه‌کان چونه؟»، «تابه‌نمه‌نده هاوه‌شه‌کانی زمانه‌کانی دنیا کامانه‌ن؟»، «مروف چون ده‌نوسن و چلون زمانی نوسن لئک

اللغويات أو اللسانيات أو الألسن(هامش) هي العلم الذي يهتم بدراسة اللغات الإنسانية ودراسة خصائصها وتراكيبها ودرجات التشابه والتباين فيما بينها. أما اللغوي فهو الشخص الذي يقوم بهذه الدراسة. ظهرت في القرن ١٩م وهي متعلقة بدراسة اللغة. جاءت بفكرة رئيسة مع العالم دي سوسير فمع علمنة الثورة الصناعية اراد علمنة اللغة أيضا في كتابه /محاضرات في اللغويات العامة/فاللغة عنده تحمل هويات من القيم الدين.المحيط.الثقافة.الفكر الفلسفي.

لسانيات (Linguistics)ایک ایسا مضمون ہے جس میں انسانی زبانوں کا ، زبانوں کی موجودہ صورت کا اور زبانوں میں وقت کے ساتھ ساتھ ہونے والی تبدیلیوں کا مطالعہ کیا جاتا ہے۔ اس علم میں مختلف زبانوں کی آپس میں مشابہت کے بارے میں مطالعہ کے ساتھ ساتھ اس چیز کا بھی مطالعہ کیا جاتا ہے کہ زبانوں کا اس دنیا کی دیگر چیزوں کے ساتھ کیا تعلق ہے۔ گویا لسانیات در اصل وہ علم ہے جس میں صرف انسانی زبان پر بحث کی جاتی ہے اور اس کے علاوہ کسی دوسری نظام کا مطالعہ نہیں کیا جاتا۔ اردومیں اسم "لسان" کے ساتھ "یات" بطور لاحقہ نسبت لگانے سے لسانیات کی اصطلاح بنتی ہی ۔ یہ اصطلاح اردو میں بطور اسم استعمال ہوتی ہے اور علم اللسان یا علم زبان

• اردو

• پشتو

انسانان د زمانې په لمن کې پېښې زيروي او مورخين د هغو پېښو په پرله پسې ډولو تاريخ ليکنه کوي . تر ټولو ښه تاريخ ، هغو مورخينو ليکلی چې پېښې يې په هغه ډول چې واقع شوېدي ، د زمانې په قيد کې راټولې کړي او د پوهيدو وړ يې گرځولې دي . تاريخ ليکنه د هر هيواد په ملت جوړونه کې ستر رول لري ، هغومره چې د ټولنې سياست پوهان ، واکمنان ، ټولواکان د ټولنې په شتون ، وده او پراختيا او يا هم بدحالی کې ونډه لري د مورخينو ونډه هم په هماغه پيمانه ده . تاريخ د ټولنې له پاره هغه هنداره ده چې هر لوستونکی کولای شي چې د هغې ټولنې خيره ، قد



دسته بندی متون ...

○ مساله

• ورودی

○ یک سند d

○ تعداد مشخص (و ثابت) دسته $C = \{c_1, c_2, \dots, c_J\}$

• خروجی

○ پیش بینی یکی از دسته‌ها $c \in C$ برای سند d

○ راه حل

• استفاده از تعدادی قانون برای تعیین هر دسته بر اساس ویژگی‌های مشخصی از متن

○ استخراج مجموعه قوانین مناسب زمان بر، پرهزینه و گاهی مشکل است

• استفاده از یادگیری ماشین

○ یادگیری قوانین از روی داده‌ها



دسته بندی متون ...

○ استفاده از یادگیری ماشین (با ناظر) ...

● فاز آموزش

○ تعداد مشخص (و ثابت) دسته $C = \{c_1, c_2, \dots, c_J\}$

○ ورودی: داده آموزش = تعداد M نمونه سند دارای برچسب (دسته سندها مشخص است)

○ $(d_1, c_1), \dots, (d_m, c_m)$

○ خروجی: مدل‌های دسته‌ها/تابع جداکننده دسته‌ها

● فاز استفاده (آزمون)

○ پیش‌بینی یکی از دسته‌ها $c \in C$ برای سند d

○ ورودی: یک سند d با دسته نامشخص + مدل‌های دسته‌ها/تابع جداکننده دسته‌ها

○ خروجی: دسته $c \in C$ پیش‌بینی شده برای سند d



دسته بندی متون ...

○ روش‌های یادگیری ماشین (با ناظر)

- بیز ساده (Naïve Bayes)
- رگرسیون (Regression)
- نزدیک‌ترین همسایه (k-Nearest Neighbors)
- شبکه عصبی مصنوعی (Artificial Neural Network)
- ماشین بردار پشتیبان (SVM: Support Vector Machines)
- مدل مخفی مارکوف (HMM: Hidden Markov Model)
- ...



دسته بندی متون: پیز ساده ...

○ استفاده از قانون بیز

$$P(\text{Class} | \text{Doc}) = \frac{P(\text{Doc} | \text{Class})P(\text{Class})}{P(\text{Doc})}$$

• برای سند Doc و دسته Class داریم

○ دسته‌ای انتخاب می شود که احتمال $p(\text{Class} | \text{Doc})$ بالاتری داشته باشد

• تخمین بیشینه احتمال پسین (MAP: maximum a posteriori)

$$c_{MAP} = \underset{c \in C}{\operatorname{argmax}} P(c | d) = \underset{c \in C}{\operatorname{argmax}} \frac{P(d | c)P(c)}{P(d)} = \underset{c \in C}{\operatorname{argmax}} P(d | c)P(c)$$

عدم تاثیر مقدار مخرج بر تصمیم گیری

• با فرض نمایش سند d به صورت بردار $[x_1, x_2, \dots, x_n]$

$$c_{MAP} = \underset{c \in C}{\operatorname{argmax}} P(d | c)P(c) = \underset{c \in C}{\operatorname{argmax}} P(x_1, x_2, \dots, x_n | c)P(c)$$



دسته بندی متون: بیز ساده ...

○ بیز ساده (Naïve Bayes) ...

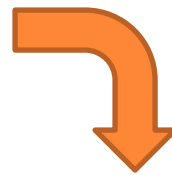
• فرض اصلی: مستقل بودن ویژگی‌ها از همدیگر

○ استقلال شرطی: ویژگی به شرط دسته $P(x_i | c_j)$

$$P(x_1, \dots, x_n | c) = P(x_1 | c) \times P(x_2 | c) \times P(x_3 | c) \times \dots \times P(x_n | c)$$

• بنابراین

$$c_{MAP} = \operatorname{argmax}_{c \in C} P(x_1, x_2, \dots, x_n | c) P(c)$$



$$c_{NB} = \operatorname{argmax}_{c_j \in C} P(c_j) \prod_{i=1}^n P(x_i | c_j)$$



دسته بندی متون: پیز ساده ...

بیز ساده (Naïve Bayes)

• فرض کیسه لغات: محل کلمات در متن مهم نیست

1

$Y(\text{I love this movie! It's sweet, but with satirical humor. The dialogue is great and the adventure scenes are fun... It manages to be whimsical and romantic while laughing at the conventions of the fairy tale genre. I would recommend it to just about anyone. I've seen it several times, and I'm always happy to see it again whenever I have a friend who hasn't seen it yet.}) = C$

thumbs up / thumbs down

2

$Y(\text{I love this movie! It's sweet, but with satirical humor. The dialogue is great and the adventure scenes are fun... It manages to be whimsical and romantic while laughing at the conventions of the fairy tale genre. I would recommend it to just about anyone. I've seen it several times, and I'm always happy to see it again whenever I have a friend who hasn't seen it yet.}) = C$

thumbs up / thumbs down

3

$Y(\text{x love xxxxxxxxxxxxxxxxxxxx sweet
xxxxxxx satirical xxxxxxxxxxxx
xxxxxxxxxxxx great xxxxxxxx
xxxxxxxxxxxxxxxxxxxxxxxx fun xxxx
xxxxxxxxxxxxxxxxx whimsical xxxx
romantic xxxx laughing
xxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxx
xxxxxxxxxxxxxxxxx recommend xxxxx
xxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxx
xx several xxxxxxxxxxxxxxxxxxxx
xxxxx happy xxxxxxxxxxxx again
xxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxx
xxxxxxxxxxxxxxxxxxxxxxxx}) = C$

thumbs up / thumbs down

4

$Y(\text{great 2
love 2
recommend 1
laugh 1
happy 1
... ..}) = C$

thumbs up / thumbs down



دسته بندی متون: پیز ساده (آموزش) ...

آموزش: تخمین پارامترهای مدل

- $P(c_j)$: احتمال دسته زام
- $P(x_i | c_j)$: احتمال ویژگی (کلمه نام) به شرط دسته زام

- تخمین شباهت بیشینه (MLE) برای احتمال دسته زام $P(c_j)$
 - تعداد مستندات که دسته زام هستند به تعداد کل مستندات
 - در بیشتر کاربردها از توزیع یکنواخت (برابر بودن احتمال برای همه دسته ها)

$$\hat{P}(c_j) = \frac{DocCount(c = c_j)}{N_{Doc}}$$

- تخمین شباهت بیشینه (MLE) برای احتمال ویژگی (کلمه نام) به شرط دسته زام $P(x_i | c_j)$
 - تعداد تکرار کلمه w_i در کل مستندات که دسته زام هستند، به تعداد کل کلمات موجود در مستندات دسته زام

$$\hat{P}(x_i = w_i | c_j) = \frac{Count(w_i, c_j)}{\sum_{w \in V} Count(w, c_j)}$$

- در عمل: تبدیل کلیه فایل‌های هر دسته به یک فایل برای محاسبه ساده‌تر



دسته بندی متون: پیز ساده (آموزش) ...

○ مشکل

- طبق قانون پیز ساده، اگر یکی از احتمال‌های $P(x_i | c_j)$ صفر باشد، کل احتمال برای آن دسته صفر خواهد بود (مستقل از اینکه مابقی $P(x_i | c_j)$ ها چه مقداری دارند)
- در داده‌های آموزش ویژگی (کلمه) x_i در دسته c_j رخ نداده است

$$c_{NB} = \operatorname{argmax}_{c_j \in C} P(c_j) \prod_{i=1}^n P(x_i | c_j)$$

• مثال: عقیده کاوی (موافق/مخالف)

- صفر شدن احتمال برای کلمه «خفن» در دسته «موافق» به دلیل عدم وقوع این کلمه در داده آموزش

$$\hat{p}(\text{موافق} | \text{خفن}) = \frac{\text{Count}(\text{"خفن", موافق})}{\sum_{w \in V} \text{Count}(w, \text{موافق})} = 0$$

- باعث صفر شدن احتمال انتخاب دسته «موافق» در متن (نظر) زیر می‌شود: «طراحی زیبا و کیفیت صفحه خفن به همراه باطری قوی. این گوشی محبوب منه»

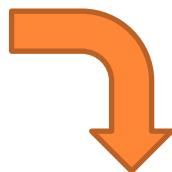


دسته بندی متون: پیز ساده (آموزش) ...

○ رفع مشکل صفر شدن احتمال‌ها

- اضافه کردن مقداری ثابت (مثلا ۱) به تعداد همه کلمات واژگان در محاسبه احتمال

$$\hat{P}(w_i | c_j) = \frac{\text{count}(w_i, c_j)}{\sum_{w \in V} (\text{count}(w, c_j))}$$



$$\hat{P}(w_i | c_j) = \frac{\text{Count}(w_i, c_j) + 1}{\sum_{w \in V} (\text{Count}(w, c_j) + 1)} = \frac{\text{Count}(w_i, c_j) + 1}{\left(\sum_{w \in V} \text{Count}(w, c_j) \right) + |V|}$$

اندازه واژگان = تعداد کلمات

$$\hat{P}(w_i | c_j) = \frac{\text{Count}(w_i, c_j) + \alpha}{\left(\sum_{w \in V} \text{Count}(w, c_j) \right) + \alpha |V|}$$

- می توان به جای ۱ مقدار ثابتی مانند α به تعداد اضافه کرد



دسته بندی متون: پیز ساده (الگوریتم آموزش) ...

○ گام ۱: در نظر گرفتن (استخراج) واژگان V و تعداد C دسته

• اندازه واژگان = تعداد لغات = $|V|$

○ گام ۲: محاسبه $P(c_j)$ ها و $P(w_i | c_j)$

• ۱-۲ برای هر دسته c_j

$$\hat{p}(c_j) = \frac{|Docs_j|}{|Total\ No\ of\ Docs|}$$

○ ۱-۱-۲ محاسبه احتمال $P(c_j)$

○ تقسیم تعداد سندهای دسته c_j به کل سندها

• $docs_j \leftarrow$ all docs with class = c_j

○ ۲-۱-۲ ایجاد یک فایل حاوی تمام متون آموزش دسته c_j : $Text_j \leftarrow$ single doc containing all $docs_j$

○ ۲-۱-۳ برای هر کلمه w_i در واژگان V

○ ۲-۱-۳-۱ قرار بده $n_j \leftarrow$ # of occurrences of w_i in $Text_j$ و $n_j \leftarrow$ # of all words in $Text_j$

• تعداد رخداد کلمه w_i در متون دسته c_j

○ ۲-۱-۳-۲ محاسبه احتمال $P(w_i | c_j)$

$$\hat{p}(w_i | c_j) = \frac{n_{ij} + \alpha}{n_j + \alpha |V|}$$

• تقسیم فراوانی کلمه w_i در دسته c_j بر کل کلمات این دسته



دسته بندی متون: پیز ساده (مثال) ...

○ هدف: دسته بندی متون به دو دسته «ایران» و «آمریکا»

○ داده آموزش و آزمون

• تعداد واژه‌ها = $|V| = 10$

آموزش		
دسته	متن	شماره سند
ایران	تهران شهر بزرگ ایران	۱
ایران	ایران بزرگ آمریکا	۲
ایران	تهران فارسی زبان ایران	۳
آمریکا	واشنگتن شهر آمریکا	۴
آمریکا	زبان نیویورک انگلیسی	۵
آزمون		
؟	زبان فارسی انگلیسی تهران	۱

○ آموزش ...

• احتمال دسته‌ها

تعداد سندهای «ایران»

$$p(c = \text{ایران}) = \frac{3}{5}$$

تعداد کل سندها

$$p(c = \text{آمریکا}) = \frac{2}{5}$$



دسته بندی متون: پیز ساده (مثال) ...

آموزش

$$p(w_i|c_j) = \frac{n_{ij} + 1}{n_j + |V|}$$

• احتمال کلمات در دسته‌ها

تعداد کلمات دسته «ایران» = ۱۱

تعداد کلمات دسته «آمریکا» = ۶

آموزش		
دسته	متن	شماره سند
ایران	تهران شهر بزرگ ایران	۱
ایران	ایران بزرگ آمریکا	۲
ایران	تهران فارسی زبان ایران	۳
آمریکا	واشنگتن شهر آمریکا	۴
آمریکا	زبان نیویورک انگلیسی	۵

کلمه	P(ایران کلمه)	P(آمریکا کلمه)
تهران	$P(\text{ایران} \text{تهران}) = (2+1)/(11+10) = 3/21$	$P(\text{آمریکا} \text{تهران}) = (0+1)/(6+10) = 1/16$
شهر	$P(\text{ایران} \text{شهر}) = (1+1)/(11+10) = 2/21$	$P(\text{آمریکا} \text{شهر}) = (1+1)/(6+10) = 2/16$
بزرگ	$P(\text{ایران} \text{بزرگ}) = (2+1)/(11+10) = 3/21$	$P(\text{آمریکا} \text{بزرگ}) = (0+1)/(6+10) = 1/16$
ایران	$P(\text{ایران} \text{ایران}) = (3+1)/(11+10) = 4/21$	$P(\text{آمریکا} \text{ایران}) = (0+1)/(6+10) = 1/16$
آمریکا	$P(\text{ایران} \text{آمریکا}) = (1+1)/(11+10) = 2/21$	$P(\text{آمریکا} \text{آمریکا}) = (1+1)/(6+10) = 2/16$
فارسی	$P(\text{ایران} \text{فارسی}) = (1+1)/(11+10) = 2/21$	$P(\text{آمریکا} \text{فارسی}) = (0+1)/(6+10) = 1/16$
زبان	$P(\text{ایران} \text{زبان}) = (1+1)/(11+10) = 2/21$	$P(\text{آمریکا} \text{زبان}) = (1+1)/(6+10) = 2/16$
واشنگتن	$P(\text{ایران} \text{واشنگتن}) = (0+1)/(11+10) = 1/21$	$P(\text{آمریکا} \text{واشنگتن}) = (1+1)/(6+10) = 2/16$
نیویورک	$P(\text{ایران} \text{نیویورک}) = (0+1)/(11+10) = 1/21$	$P(\text{آمریکا} \text{نیویورک}) = (1+1)/(6+10) = 2/16$
انگلیسی	$P(\text{ایران} \text{انگلیسی}) = (0+1)/(11+10) = 1/21$	$P(\text{آمریکا} \text{انگلیسی}) = (1+1)/(6+10) = 2/16$



دسته بندی متون: پیز ساده (مثال) ...

$$c_{NB} = \operatorname{argmax}_{c_j \in C} P(c_j) \prod_{i=1}^n P(x_i | c_j)$$

○ آزمون

• محاسبه $p(\text{class} | \text{doc})$ برای $\text{doc} = \text{«زبان فارسی انگلیسی تهران»}$

• محاسبه برای $\text{class} = \text{ایران}$

$$\begin{aligned} p(\text{ایران} | \text{تهران}) &= p(\text{ایران} | \text{انگلیسی}) \cdot p(\text{ایران} | \text{فارسی}) \cdot p(\text{ایران} | \text{زبان}) \cdot p(\text{ایران}) \\ &= \frac{3}{5} \cdot \frac{2}{21} \cdot \frac{2}{21} \cdot \frac{1}{21} \cdot \frac{3}{21} = 0.0004 \end{aligned}$$

• محاسبه برای $\text{class} = \text{آمریکا}$

$$\begin{aligned} p(\text{آمریکا} | \text{تهران}) &= p(\text{آمریکا} | \text{انگلیسی}) \cdot p(\text{آمریکا} | \text{فارسی}) \cdot p(\text{آمریکا} | \text{زبان}) \cdot p(\text{آمریکا}) \\ &= \frac{2}{5} \cdot \frac{2}{16} \cdot \frac{1}{16} \cdot \frac{2}{16} \cdot \frac{1}{16} = 0.0002 \end{aligned}$$

• بنابراین داریم $p(\text{ایران} | \text{doc}) > p(\text{آمریکا} | \text{doc})$

○ سند متعلق به دسته «ایران» است



دسته بندی متون: پیز ساده ...

○ ویژگی‌ها

- خیلی سریع است
- عدم نیاز به حافظه زیاد
- مقاوم به ویژگی‌های غیرمرتبط (به علت لغو اثر این ویژگی‌ها توسط همدیگر)
- کارایی خوب در کاربردهایی که ویژگی‌هایی با اهمیت مشابه دارند
- در صورت برقراری شرط استقلال، این روش تبدیل به دسته‌بند **بهینه** خواهد شد

- درارای کارایی خوب برای دسته‌بندی متن
- دارای ارتباط مفهومی با مدل زبانی

- هنوز نیاز به روشی با کارایی بالاتر
 - روش‌های دیگر یادگیری ماشین



دسته بندی متون: بیز ساده (ارتباط با مدل زبانی) ...

○ ویژگی‌های مورد استفاده بیز ساده در دسته بندی متون = کلمات

- امکان استفاده از هر نوع ویژگی‌ای در این روش
- اگر از همه کلمات متن در بیز ساده استفاده کنیم، آنگاه دسته بندی مشابه محاسبه احتمال در مدل زبانی است

○ مدل زبانی یک تایی (Unigram)

- محاسبه احتمال برای هر کلمه $P(w_i)$
- محاسبه احتمال برای هر جمله با ضرب احتمال تک تک کلمات در همدیگر

$$P(W) = P(w_1 w_2 \cdots w_m) \cong \prod_{i=1}^m P(w_i | w_{i-1} \cdots w_1) \cong \prod_{i=1}^m P(w_i) = P(w_1) P(w_2) P(w_3) \cdots P(w_m)$$

- در صورتی که مقدار احتمال هر کلمه $P(w_i)$ برای یک دسته محاسبه شود = $P(w_i | c_j)$

$$P(\text{Sentence} | c_j) = P(W | c_j) = \prod_{i=1}^m P(w_i | c_j)$$



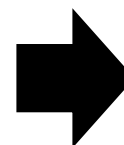
دسته بندی متون: پیز ساده (ارتباط با مدل زبانی) ...

مثال

محاسبه احتمال هر کلمه در هر دسته

Model Positive		Model Negative	
0.1	I	0.2	I
0.1	love	0.001	love
0.01	this	0.01	this
0.05	fun	0.005	fun
0.1	film	0.1	film

<u>I</u>	<u>love</u>	<u>this</u>	<u>fun</u>	<u>film</u>
0.1	0.1	0.01	0.05	0.1
0.2	0.001	0.01	0.005	0.1



$$P(s|pos) > P(s|neg)$$



دسته بندی متون: پیز ساده (نکات کاربردی) ...

○ رفع مشکل Underflow

- ضرب تعداد زیادی مقدار احتمال (با مقادیر کوچک) = صفر شدن

$$c_{NB} = \operatorname{argmax}_{c_j \in C} P(c_j) \prod_{i=1}^n P(x_i | c_j)$$



$$c_{NB} = \operatorname{argmax}_{c_j \in C} \log P(c_j) + \sum_{i \in \text{positions}} \log P(x_i | c_j)$$

- استفاده از لگاریتم برای رفع مشکل

- دسته‌ای با لگاریتم احتمال بیشتر (در صورت عدم نرمال کردن)، دسته برنده است

○ دارای کارایی خوب با داده کم



دسته بندی متون: پیز ساده

○ مثال (فارسی)

- تعداد ۳۱۹۳ سند
- ۸ موضوع (عنوان): ادبی، مذهبی، اقتصادی، هنری، پزشکی، تاریخی، سیاسی و ورزشی
- تعداد کلمه‌ها (ویژگی‌ها): ۲۹۰۰ کلمه
- کارایی روش‌های مختلف (بر حسب درصد)

DT	50-NN	Naïve Bayes	SVM	ANN	HMM	
۵۹	۶۴	۶۳	۶۹	۷۴	۸۳	میانگین Recall
۵۵	۷۰	۶۰	۷۳	۷۳	۷۶	میانگین Precision
۵۷	۶۷	۶۱	۷۱	۷۳	۷۹	میانگین معیار F

