1400/09/11

استدلال بیزین روشی بر پایه احتمالا برای استنتاج کردن است.

اساس کلی بیز این است که برای هر کمیتی یک توزیع احتمال وجود دارد که با مشاهده یک داده جدید و استدلال در مورد توزیع احتمال آن می توان تصمیمات بهینه ای در خصوص آن اتخاذ کرد.

اهمیت یادگیری بیزین:

- کاربرد در دسته بندی متن
- ارائه راه حل های مفید در روش های یادگیری
- کارایی بهتر بیز در مسائل یکسان نسبت به روش های درخت تصمیم و شبکه های عصبی

- دانش موجود درباره ی موضوع را با تعدادی احتمال ذخیره می کنیم.
- مقادیر کیفی دانش موجود را به صورت توزیع احتمال، فرضیات استقلال و ... مدل می کنیم.
 - مدلی که ایجاد می کنیم دارای پارامترهای ناشناخته خواهد بود.
- برای هر یک از مقادیر ناشناخته، تزیع احتمال پیشین در نظر گرفته ایم که بازگو کننده ی این احتمال داشتن را در هر یک از مقادیر بدون مشاهده نمونه های داده ایجاد می کند.
 - جمع آوری داده ها را انجام می دهیم.
 - با مشاهده داده های مختلف مقدار توزیع احتمال پسین را محاسبه می کنیم.
 - با احتمال پسین در مورد نمونه های جدید تصمیم گیری های خودمان را انجام می دهیم.

در بیز یادگیری به صورت تدریجی است.

هر نمونه داده باعث افزایش یا کاهش احتمال درست بودن فرضیه ی ما می شود.

احتمالی که برای یک فرضیه در نظر می گیریم با دانش قبلی و اطلاعات موجودی که در آن نمونه مشاهده شده است تصمیم گیری می شود.

پس احتمال قبلی برای هر فرضیه موجود می باشد.

فرضیه های حاصل از روش های بیزین قادر به پیش بینی احتمالی هستند.

داده های جدید را می توان با ترکیب وزنی چندین فرضیه دسته بندی نمود.

D به طور کلی در فضای فرضیه یا H با کمک تئوری بیز به دنبال بهترین فرضیه ای خواهیم بود که در مورد داده های آموزشی یا H محدق پیدا کند.

یک راه تعیین بهترین فرضیه گشتن به دنبال محتمل ترین فرضیه می باشد که با داشتن داده ی آموزشی D و احتمال قبلی در مورد فرضیه های مختلف می توان انتظار داشت.

فضای فرضیه H

مجموعه داده آموزشي D

احتمال پیش بینی که فرضیه h قبل از مشاهده مثال آموزشی D را داشته است. (P(h

احتمال پیشین مشاهده داده های آموزشی (P(D

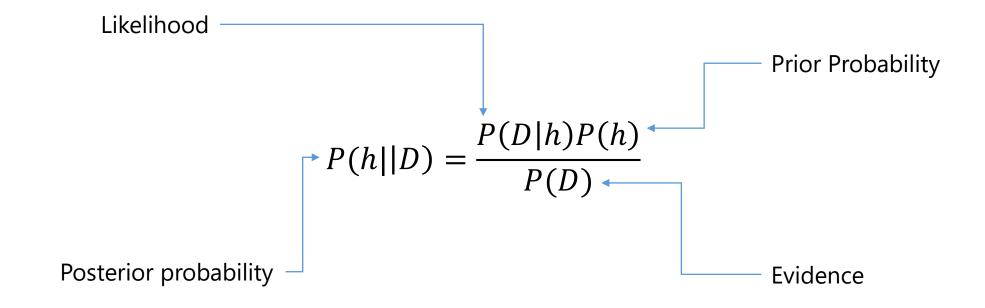
P(D|h) صادق باشد D به فرض اینکه فرضیه D صادق باشد

انتظار نهایی ما به دست آمدن (P(h|D است.

مفهوم P(h|D) -> احتمال اینکه با داده ی آموزشی D فرضیه D صادق باشد. (احتمال پسین)

احتمال پیشین همیشه مستقل از D است ولی احتمال پسین متاثر از D است.

یادگیری بیزین بر مبنای نظریه احتمال بیزین است.



نحوه تعیین محتمل ترین فرضیه (فرضیه با حداکثر پسین MAP) از بین مجموعه فرضیات ممکن H:

$$h_{MAP} \equiv \arg\max_{h \in H} P(h|D)$$

$$= \arg\max_{h \in H} \frac{P(D|h)P(h)}{P(D)}$$

$$= \arg\max_{h \in H} P(D|h)P(h)$$
h خستقل از h

در مواقعی که هیچ اطلاعی در مورد (P(h وجود نداشته باشد، می توان فرض کرد که تمام فرضیه های H دارای احتمال پیشین یکسانی هستند.

در اینصورت، برای محاسبه فرضیه با حداکثر احتمال، می توان فقط مقدار P(D|h) را در نظر گرفت. این مقدار، میزان احتمال Maximum داده D نسبت به فرض D نامیده می شود و هر فرضیه ای که این مقدار را بیشینه کند، فرضیه ی Likelihood یا D نامیده می شود:

$$h_{ML} \equiv \arg \max_{h \in H} P(h|D)$$

مثال: تشخیص سرطان

حالات بيمار:

• بیمار دارای بیماری سرطان است.

• بيمار سالم است.

P(cancer) = 0.008

P(+|cancer)=0.98

 $P(+|\sim cancer)=0.03$

 $P(\sim cancer) = 0.992$

P(-|cancer) = 0.02

 $P(-|\sim cancer) = 0.97$

مثال: تشخیص سرطان

یک بیمار جدید آمده و نیاز به تصمیم گیری داریم:

احتمال داشتن سرطان

P(cancer|+)=P(+|cancer|) P(cancer) / P(+) = (0.98)(0.008)/P(+)=0.0078/P(+)

احتمال نداشتن سرطان

 $P(\sim cancer | +) = P(+ | \sim cancer) P(\sim cancer) / P(+) = (0.03)(0.992)/P(+) = 0.0298/P(+)$

مثال: تشخیص سرطان

 $P(cancer|+)+P(\sim cancer|+)=1$

0.0078/P(+)+0.0298/P(+)=1

P(+)=0.0078+0.0298=0.0376

مثال: تشخیص سرطان

احتمال ابتلای بیمار به سرطان:

احتمال نداشتن سرطان در بیمار:

$$P(cancer|+)=0.0078/P(+)=0.21$$

 $P(\sim cancer +) = 0.0298/P(+) = 0.79$

مفهوم Brute-force Map Learning

مى توان با استفاده از نظريه بيزين، الگوريتمى براى يادگيرى مفهوم كه بتواند فرضيه با بيشترين احتمال را بدست آورد داشته باشيم و آن را Brue-force Map Learning Algorithm بناميم.

برای هر فرضیه h موجود در H مقدار احتمال پسین را حساب می کنیم.

فرضیه h_{map} را که بیشترین احتمال پسین را دارد مشخص می کنیم.

Bayes Optimal Classifier

دسته بندی بهینه بیزین

Brue-force Map Learning Algorithm : محتمل ترین فرضیه برای مجموعه داده آموزشی داده شده چیست؟

• محتمل ترین دسته بندی یک نمونه مشاهده شده چیست؟

در عمل محتمل ترین دسته بندی برای یک نمونه جدید از ترکیب پیش بینی تمامی فرضیه ها بدست می آید. مقدار پیش بینی هر فرضیه در احتمال پسین آن ضرب شده و حاصل آنها با هم ترکیب می شود.

مثال:

ورودی ها:

h1
$$P(h_1|D)=0.4$$

h2
$$P(h_2|D)=0.3$$

h3
$$P(h_3|D)=0.3$$

$$P(h_1) = + P(h_2) = - P(h_3) = -$$

احتمال مثبت بودن

احتمال منفى بودن

دسته بندی بهینه بیزین:

$$P(h_1|D) = 0.4 P(-|h_1) = 0 P(+|h_1) = 1$$

$$P(h_2|D) = 0.3 P(-|h_2) = 1 P(+|h_2) = 0$$

$$P(h_3|D) = 0.3 P(-|h_3) = 1 P(+|h_3) = 0$$

لذا:

$$\sum i P(+|h_i)P(h_i|D) = 0.4 \ and \ \sum i P(-|h_i)P(h_i|D) = 0.6$$

با توجه به این نمونه های به صورت منفی دسته بندی می شود؛ دقت کنید استفاده از این روش برای فضاهای فرضیه های بزرگ، غیرعملی است.

Naïve Bayes Classifier

کارایی بالا نسبت به شبکه های عصبی و درخت تصمیم

کاربرد:

- نمونه X توسط ترکیب عطفی ویژگی های قابل توصیف باشد.
 - این ویژگی ها به صورت شرطی مستقل از یکدیگر باشند.
- تابع هدف f(x) بتواند هر مقداری را از مجموعه محدوده V داشته باشد.
 - مجموعه مثال های آموزشی زیاد باشد.

Naïve Bayes Classifier

تابع هدف:

$$f:X->V$$

 $X=(a_1, ..., a_n)$

صورت مساله : برای یک نمونه مشاهده شده، مقدار تابع هدف یا به عبارت دیگر دسته بندی آن را مشخص کنید.

Naïve Bayes Classifier

حل مسئله روش بیزین برای حل مساله محتمل ترین مقدار هدف V_{MAP} محاسبه می شود:

 $v_{MAP} = \arg \max P(v_j | a_1, ..., a_n)$

در این رابطه با استفاده از نظریه بیزین به صورت زیر نوشته می شود:

$$v_{MAP} = \arg\max_{v_j \in V} \frac{P(a_1, ..., a_n | v_j) P(v_j)}{P(a_1, ..., a_n)} = \arg\max_{v_j \in V} P(a_1, ..., a_n | v_j) P(v_j)$$

Naïve Bayes Classifier

 $v_{MAP} = \arg \max P(v_j | a_1, ..., a_n)$

تعداد دفعاتی که v_j در مثالهای آموزشی مشاهده شده را با $P(v_j)$ می شناسیم و اینکه مصاحبه ی $P(a_1,, a_n | v_j)$ در دست به جز اینکه مجموعه داده آموزشی بسیار بزرگی در دست باشد.

روش یادگیری نایو بیز بر پایه ی فرض ساده (Naive) عمل می کند: مقادیر ویژگی ها به صورت شرطی مستقل هستند.

برای یک مقدار هدف احتمال مشاهده ی ترکیب عطفی $(a_1, ..., a_n)$ برابر است با حاصلضرب احتمال تک تک ویژگی ها در اینصورت رابطه ی زیر بدست می آید:

$$V_{NB} = \arg \max_{v_j \in V} P(v_j) \prod_{i=1}^{n} P(a_i | v_j)$$

Naïve Bayes Classifier

در Naïve Bayes مقادیر مختلف $P(v_j)$ و $P(v_j)$ با استفاده از تعداد دفعات تکرار تخمین زده می شود.

مجموعه این تخمین ها فرضیه ای را با استفاده از رابطه ی زیر تخمین می زند.

$$V_{NB} = \arg \max_{v_j \in V} P(v_j) \prod_{i=1}^{n} P(a_i | v_j)$$

در این روش هیچگونه عمل جستجوی آشکاری وجود ندارد.

Day	Outlook	Temperature	Humidity	Wind	Play Tennis
1	Sunny	Hot	High	weak	No
2	Sunny	Hot	High	Strong	No
3	Overcast	Hot	High	Weak	Yes
4	Rain	Mild	High	Weak	Yes
5	Rain	Cool	Normal	Weak	Yes
6	Rain	Cool	Normal	Strong	No
7	Overcast	Cool	Normal	String	Yes
8	Sunny	Mild	High	Weak	No
9	Sunny	Cool	Normal	Weak	Yes
10	Rain	Mild	Normal	Weak	Yes
11	Sunny	Mild	Normal	String	Yes
12	Overcast	Mild	High	Strong	Yes
13	Overcast	Hot	Normal	Weak	Yes
14	Rain	Mild	High	strong	no

Outlook	Temperature	Humidity	Wind	Play Tennis
Sunny	Hot	High	False	N
Sunny	Hot	High	True	N
Overcast	Hot	High	false	Р
Rain	Mild	High	False	Р
Rain	Cool	Normal	False	Р
Rain	Cool	Normal	true	N
Overcast	Cool	Normal	True	Р
Sunny	Mild	High	False	N
Sunny	Cool	Normal	False	Р
Rain	Mild	Normal	False	Р
Sunny	Mild	Normal	True	Р
Overcast	Mild	High	True	Р
Overcast	Hot	Normal	False	Р
Rain	Mild	High	true	N

P(p) = 9/14	
P(n) = 5/14	

Outlook	Temperature	Humidity	Wind	Play Tennis
Sunny	Hot	High	False	N
Sunny	Hot	High	True	N
Overcast	Hot	High	false	Р
Rain	Mild	High	False	Р
Rain	Cool	Normal	False	Р
Rain	Cool	Normal	true	N
Overcast	Cool	Normal	True	Р
Sunny	Mild	High	False	N
Sunny	Cool	Normal	False	Р
Rain	Mild	Normal	False	Р
Sunny	Mild	Normal	True	Р
Overcast	Mild	High	True	Р
Overcast	Hot	Normal	False	Р
Rain	Mild	High	true	N

Outlook::

P(sunny | p)=2/9P(sunny | n)=3/5

P(overcast | p)=4/9P(overcast | n)=0

P(rain | p)=3/9 P(rain | n)=2/5

Outlook	Temperature	Humidity	Wind	Play Tennis
Sunny	Hot	High	False	N
Sunny	Hot	High	True	N
Overcast	Hot	High	false	Р
Rain	Mild	High	False	Р
Rain	Cool	Normal	False	Р
Rain	Cool	Normal	true	N
Overcast	Cool	Normal	True	Р
Sunny	Mild	High	False	N
Sunny	Cool	Normal	False	Р
Rain	Mild	Normal	False	Р
Sunny	Mild	Normal	True	Р
Overcast	Mild	High	True	Р
Overcast	Hot	Normal	False	Р
Rain	Mild	High	true	N

temperature::

P(hot | p) = 2/9

P(hot | n) = 2/5

 $P(\text{mild} \mid p)=4/9$ $P(\text{mild} \mid n)=2/5$

P(cool | p)=3/9P(cool | n)=1/5

Outlook	Temperature	Humidity	Wind	Play Tennis
Sunny	Hot	High	False	N
Sunny	Hot	High	True	N
Overcast	Hot	High	false	Р
Rain	Mild	High	False	Р
Rain	Cool	Normal	False	Р
Rain	Cool	Normal	true	N
Overcast	Cool	Normal	True	Р
Sunny	Mild	High	False	N
Sunny	Cool	Normal	False	Р
Rain	Mild	Normal	False	Р
Sunny	Mild	Normal	True	Р
Overcast	Mild	High	True	Р
Overcast	Hot	Normal	False	Р
Rain	Mild	High	true	N

humidity::

P(high | p)=3/9P(high | n)=4/5

P(normal | p)=6/9P(normal | n)=2/5

Outlook	Temperature	Humidity	Wind	Play Tennis
Sunny	Hot	High	False	N
Sunny	Hot	High	True	N
Overcast	Hot	High	false	Р
Rain	Mild	High	False	Р
Rain	Cool	Normal	False	Р
Rain	Cool	Normal	true	N
Overcast	Cool	Normal	True	Р
Sunny	Mild	High	False	N
Sunny	Cool	Normal	False	Р
Rain	Mild	Normal	False	Р
Sunny	Mild	Normal	True	Р
Overcast	Mild	High	True	Р
Overcast	Hot	Normal	False	Р
Rain	Mild	High	true	N

windy::

P(true | p) = 3/9

P(true $\mid n$)=3/5

P(false | p)=6/9

P(false | n = 2/5

Naïve Bayes Classifier

X: (outlook=sunny, temperature=cool, Humidity=high, wind=strong)

$$V_{NB} = \arg \max_{v_k \in [yes,no]} P(v_k) \prod_{i=1}^{n} P(a_i | v_k)$$

$$= \arg \max_{v_k \in [yes,no]} P(v_k) P(outlook = sunny | v_k) P(temp$$

Naïve Bayes Classifier

مقدار احتمالات بر اساس نسبت تعداد مشاهده شده یک مقدار به تعداد کل حالات ممکن محاسبه شد: n_c / n این مقدار می تواند منجر به این مقدار می تواند منجر به نتایج غلط شود.

اگر برای دسته v_j مقدار a_i هرگز مشاهده نشود، مقدار $v_j = 0$ شده و در نتیجه کل حاصلضرب صفر خواهد شد. برای جلوگیری از این مشکل از روش m-estimate استفاده می شود.

$$\frac{n_c + mp}{n + m}$$

 $m-estimate\ of\ probability$

- و n_{c} همان مقادیر قبلی است. N_{c}
 - m تعداد مثالهای مجازی است.
- مقدار p = 1 / k هم معمولا p = 1 / k می باشد و به صورت یکنواخت است که k تعداد مقادیر ممکن برای ویژگی هاست.