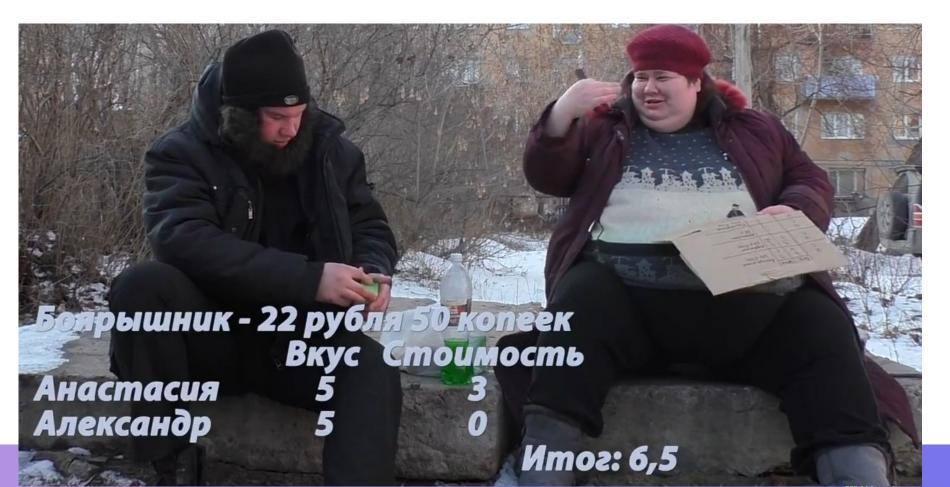
Наивный Байес и вероятностный подход к классификации

Лекция 6



Пример задачи регрессии (омская версия)



Ликбез: вероятность и формула Байеса





Вероятность события

```
Пусть А – событие. Например,
А={монета выпадет орлом}
А={докладчик помрёт во время лекции}
А={на игральном кубике выпадет 6}
Для событий можно найти их вероятность Pr(A)
следующим образом:
Пусть A1,A2,...,An – все элементарные исходы
эксперимента, тогда
Pr(A)={число исходов при которых наблюдается A}/n
```





Пример 1

```
Эксперимент: кидаем кубик.
Элементарные исходы: «выпадет 1»,..., «выпадет 6»
Для краткости:
                    1,2,3,4,5,6
Найдите вероятности следующих событий:
А={выпадет четное число}
А={выпадет число, делящиеся на 3}
А={выпадет делитель числа 6}
А={выпадет отрицательное число}
А={выпадет число меньше 10}
```





Пример 2

```
Эксперимент: заводим двоих детей.
Элементарные исходы (М – мальчик, Д- девочка):
                  ММ,МД,ДМ,ДД
Найдите вероятности следующих событий:
А={дети будут разного пола}
А={будет хотя бы одна девочка}
А={старшим ребенком будет мальчик}
А={число девочек не равно числу мальчиков}
```



Условные вероятности

- Pr(A|B) (читается как «вероятность А при условии В») Условная вероятность считается по правилу:
- 1. Мысленно представьте, что событие В уже произошло.
- 2.Удалите элементарные исходы, которые НЕ удовлетворяют событию В.
- 3. Найдите вероятность события А для нового множества элементарных исходов, используя формулу с предыдущих слайдов.





Пример 1*

Эксперимент: кидаем кубик.

Элементарные исходы:

1,2,3,4,5,6

Найдите вероятность Pr(A|B) для:

- 1) $A = \{ выпадет четное число \}$, $B = \{ выпадет число \}$ делящиеся на $3 \}$
- 2) $A = \{ выпадет 6 \}, B = \{ выпадет делитель числа 6 \}$
- 3) $A = \{ выпадет 1 \}$, $B = \{ выпадет четное число \}$



Пример 1*

Эксперимент: кидаем кубик.

Элементарные исходы:

1,2,3,4,5,6

Найдите вероятность Pr(A|B) для:

- 1) $A = \{ выпадет четное число \}$, $B = \{ выпадет число \}$ делящиеся на $3 \}$
- 2) $A = \{ выпадет 6 \}$, $B = \{ выпадет делитель числа 6 \}$
- 3) $A = \{выпадет 1\}$, $B = \{выпадет четное число\}$

А теперь подсчитайте Pr(B|A) для пп.1-3





Пример 2*

```
Эксперимент: заводим двоих детей.
Элементарные исходы (М – мальчик, Д- девочка):
                   ММ,МД,ДМ,ДД
Найдите вероятность Pr(A|B) для:
1) А={дети будут разного пола} В={среди детей есть
девочка}
2) A = \{дети будут разного пола\} B = \{старший ребенок
- девочка}
```



Пример 2*

```
Эксперимент: заводим двоих детей.
Элементарные исходы (М – мальчик, Д- девочка):
                   ММ,МД,ДМ,ДД
Найдите вероятность Pr(A|B) для:
1) А={дети будут разного пола} В={среди детей есть
девочка}
2) A = \{дети будут разного пола\} B = \{старший ребенок
- девочка}
```

А теперь подсчитайте Pr(B|A) для пп.1-2





Задачка более близкая к проблеме классификации

В выборке 60% женщин и 40% мужчин. Известно, что среди них курит 10% женщин и 70% мужчин. Про человека NN известно, что он курит. С какой вероятностью NN является женщиной (мужчиной)?

Если вы решите эту задачу, то вы фактически осуществите классификацию NN (предскажете значение его пола).





Решение

```
В выборке 60% женщин и 40% мужчин. Известно, что среди них курит 10% женщин и 70% мужчин. Про человека NN известно, что он курит. С какой вероятностью NN является женщиной (мужчиной)? Обозначим: М={человек является мужчиной}
Ж={человек является женшиной}
```

 $\mathcal{K}=\{$ человек является женщиной $\}$ $\mathcal{K}=\{$ человек курит $\}$ Из условия можно найти Pr(M), $Pr(\mathcal{K}|M)$, $Pr(\mathcal{K}|M)$.





Решение

Имеем:

Pr(M)=0.4, Pr(Ж)=0.6, Pr(K|M)=0.7, Pr(K|Ж)=0.1 Из этих чисел можно даже найти Pr(K): Pr(K)=Pr(M)*Pr(K|M)+Pr(Ж)*Pr(K|Ж)=0.4*0.7+0.6*0.1=0.34=34%

Но нас интересуют вероятности несколько других событий (каких?):





Решение

Имеем:

Pr(M)=0.4, Pr(K)=0.6, Pr(K|M)=0.7, Pr(K|K)=0.1 Из этих чисел можно даже найти Pr(K): Pr(K)=Pr(M)*Pr(K|M)+Pr(K)*Pr(K|K)=0.4*0.7+0.6*0.1=0.34=34%

Но нас интересуют вероятности несколько других событий (каких?):

Pr(M|K)=?Pr(X|K)=?

А как их найти?





Есть формула Байеса!

$$Pr(B|A)=Pr(A|B)*Pr(B)/Pr(A)$$

Для нашей задачи получаем: Pr(M|K)=Pr(K|M)*Pr(M)/Pr(K)=0.7*0.4/0.34=28/34 Pr(X|K)=Pr(K|X)*Pr(X)/Pr(K)=0.1*0.6/0.34=6/34

- фактически мы получили оценки принадлежности «классу мужчин» и «классу женщин»





Что делать с вероятностью при классификации?

На прошлом слайде были получены вероятности принадлежности объекта «классу мужчин» и «классу женщин» (28/34 и 8/34). Что с ними делать, чтобы получить точный (а не вероятностный) ответ?

- 1.Выбрать класс с наибольшей вероятностью (в данном примере выбрать «мужчин»).
 - 2. С распределением 28/34 и 8/34 сгенерировать случайную величину. Ее значение и будет окончательным ответом.





Что делать с вероятностью при классификации?

2. Установить пороговое значение п для вероятности. Если вероятность принадлежности «классу мужчин» больше п, то объект классифицируется как мужчина. Иначе – как женщина.

Такой подход оправдан, когда ущерб от ошибки классификации различен для первого и второго класса (см. пример на след. слайде).





Террорист или нет?

Допустим мы оцениваем не вероятности принадлежности классам «мужчин» и «женщин», а к классам «террористов» и «честных людей». Допустим, что для человека А были получены вероятностные оценки на принадлежности к классу террористов 6/34 и к классу «честных людей» 28/34.

Его «честность» в 4.5 раза выше его терроризма и в соответствии с методами пп1-2 он (скорей всего) будет окончательно отнесен к классу «честных».





Террорист или нет?

Но тут над понимать, что цена ошибки «террориста приняли за честного» гораздо выше цены ошибки «честного приняли за террориста». А для объекта А принадлежность террористам не такая уж и маленькая (6/34=18%) – заведомо выше доли террористов в населении страны.

Здесь оправдано установить маленький порог п (например, п=10%) и людей, у которых вероятность принадлежности «классу террористов» больше п, отправлять на дополнительную проверку.





Наивный Байес





А как получить вероятности классов, когда признаков >1?

В задаче про курение нецелевой признак был один («курит или нет»). А если нецелевых признаков больше?

Пусть

	Курит	Любит кошек	Пол (Y)
Α	0	1	1
В	0	1	0
С	1	0	1
D	1	0	0
Е	1	0	1

	Курит	Любит кошек	Пол (Y)
F	0	0	?



Вычисления для объекта F

```
Если применять формулу Байеса для объекта F, то
нужно будет знать вероятности:
Pr(K'*J' | M)
Pr(K'*Л' | Ж)
где
K = \{ K \neq K \}, K' = \{ H \in K \neq K \}, 
\Pi = \{ \text{любит кошек} \}, \Pi' = \{ \text{не любит кошек} \},
*= союз «и».
А как вычислить эти вероятности?
```



Теорема о произведении вероятности

Pr(A*B)=Pr(A)*Pr(B), если события A и B не зависят друг от друга.

Решая нашу задачу, сделаем допущение, что курение НЕ ЗАВИСИТ от любви к кошкам и наоборот (а в жизни это действительно так?). Тогда мы получаем: $Pr(K'*\Lambda' \mid M) = Pr(K' \mid M) * Pr(\Lambda' \mid M),$ $Pr(K'*\Lambda' \mid X) = Pr(K' \mid X) * Pr(\Lambda' \mid X),$ а вероятности из правых частей равенств можно найти по таблице.





Из таблицы получаем вероятности:

$$Pr(X) = 2/5$$
 $Pr(M) = 3/5$

$$Pr(K|X)=1/2$$
 $Pr(K'|X)=1/2$

$$Pr(K|M)=2/3$$
 $Pr(K'|M)=1/3$

$$Pr(\Pi|\mathcal{K})=1/2$$
 $Pr(\Pi'|\mathcal{K})=1/2$ $Pr(\Pi|M)=1/3$ $Pr(\Pi'|M)=2/3$

	Курит	Любит кошек	Пол (Y)
Α	0	1	1
В	0	1	0
С	1	0	1
D	1	0	0
Ε	1	0	1





По формуле Байеса для объекта F получаем вероятности принадлежности классам:

$$Pr(X|K'\Pi') = Pr(K'\Pi'|X)*Pr(X)/Pr(K'\Pi')$$

 $Pr(M|K'\Pi') = Pr(K'\Pi'|M)*Pr(M)/Pr(K'\Pi')$

Используем предположение о независимости курения от любви к кошкам:

$$Pr(\mathcal{K}|\mathcal{K}'\Pi')=Pr(\mathcal{K}'|\mathcal{K})*Pr(\Pi'|\mathcal{K})*Pr(\mathcal{K}'\Pi')$$

 $Pr(\mathcal{M}|\mathcal{K}'\Pi')=Pr(\mathcal{K}'|\mathcal{M})*Pr(\Pi'|\mathcal{M})*Pr(\mathcal{M})/Pr(\mathcal{K}'\Pi')$





Подставляем числа:

$$Pr(X|K'\Pi')= (1/2* 1/3* 2/5)/Pr(K'\Pi')=(1/15)/Pr(K'\Pi')$$

 $Pr(M|K'\Pi')= (2/3* 1/3* 3/5)/Pr(K'\Pi')=(2/15)/Pr(K'\Pi')$

Величину Pr(K'Л') можно не вычислять (она нам не нужна), а воспользоваться равенством $Pr(\mathcal{K}|K'Л') + Pr(M|K'Л') = 1$

И получить

$$Pr(X|K'J')=1/3 Pr(M|K'J')=2/3$$





Таким образом, принадлежность классу мужчин в 2 раза выше принадлежности классу женщин.

$$Pr(X|K'J')=1/3 Pr(M|K'J')=2/3$$

Вот таким макаром и происходит классификация с помощью формулы Байеса.

Аналогично обрабатываются таблицы с бОльшим числом признаков.





Допущение о независимости признаков

Работа классификатора существенно использует предположение о независимости признаков друг от друга.

Но независимость признаков (когда коэфф. корр=0) не часто наблюдается на практике.

Поэтому наш классификатор называют наивным байесовским классификатором (НБК).

Когда между признаками существует сильная зависимость НБК может сильно ошибаться.





Работа НБК при зависимых признаках

В этой таблице признаки коррелируют (более того: они совпадают).

Проведем расчет принадлежности классам объекта F 1)без использования второго признака,

2)с двумя признаками.

Курит	Любит	Пол (Y)
	кошек	
0	0	?





Любит

Пол (Y)

Курит

Используем только один признак

Имеем: Pr(Ж)=2/5, Pr(K'|Ж)=1/2, Pr(M)=3/5, Pr(K'|M)=2/3

Pr(K'|K') = Pr(K'|K') * Pr(K') / Pr(K') = (1/5) / Pr(K')Pr(M|K') = Pr(K'|M) * Pr(M) / Pr(K') = (2/5) / Pr(K')

Из равенства Pr(X|K')+Pr(M|K')=1 находим Pr(X|K')=1/3, Pr(M|K')=2/3

	•		
	Курит	Любит кошек	Пол (Y)
A	0	0	1
В	0	0	0
С	1	1	1
D	1	1	0
Ε	1	1	1



Используем оба признака

Имеем: $Pr(\mathcal{K})=2/5$, $Pr(\mathcal{K}'|\mathcal{K})=1/2$, $Pr(\mathcal{I}'|\mathcal{K})=1/2$, Pr(M)=3/5, Pr(K'|M)=2/3, Pr(J'|M)=2/3

$$Pr(X|K'\Pi')=Pr(K'|X)*Pr(\Pi'|X)*Pr(X)/Pr(K'\Pi')=(1/10)/Pr(K'\Pi')$$

 $Pr(M|K'\Pi')=Pr(K'|M)*Pr(\Pi'|M)*Pr(M)/Pr(K'\Pi')=(4/15)/Pr(K'\Pi')$

Из равенства Pr(X|K'J')+Pr(M|K'J')=1находим Pr(X|K')=3/11, Pr(M|K')=8/11

	(1/ ± 0	///	
	Курит	Любит кошек	Пол (Y)
Α	0	0	1
В	0	0	0
С	1	1	1
D	1	1	0
г	4	4	4

Ответ-то изменился!!!

... хотя никакой принципиально новой информации мы в таблицу не добавляли.

Вот к чему приводит наличие зависимостей между признаками!!!

Зависимые признаки нужно удалять (или как-то модифицировать) – об этом в теме «Отбор признаков»





Плюсы НБК:

- 1.Простота и быстрая работа. Используется в системах массового обслуживания. Например, при фильтрации спама.
- 2.В некоторых случаях НБК хорошо работает даже для коррелированных признаков (так что стоит рискнуть)))). Например, при фильтрации спама признаки, как правило, не являются независимыми этот эффект и рассмотрим на след. слайдах.





Фильтрация спама (общая постановка)

Задача: есть текст письма, определить является ли оно спамом (класс 1) или нет (класс 0).

Какие признаки есть у письма?

Самый простой способ: для каждого слова из словаря завести свой бинарный признак (он будет равен 1, если это слово встречается в рассматриваемом письме).

Получится примерно как на след. слайде...





Помогите Даше найти зависимые признаки

	В	
	С	
2	D	
	Е	
	F	
7		

		лекарство	запой	геморрой	заработать	миллион	Путин	Спам?
	Α	1	1	1	0	0	0	1
	В	1	1	0	0	0	0	1
	С	1	0	1	0	0	0	1
3	D	0	0	0	1	1	0	1
		1	0	0	0	1	0	1
	F	0	0	0	0	1	1	0

Вообще говоря...

... там много зависимых признаков (к тому же объем выборки небольшой).

Но НБК не обязательно на них начнет косячить (проверено на практике)!

Имейте в виду: в задаче про спам очень важно правильно установить порог п классификации.

Напомню смысл порога п: если

Pr(Спам | текст письма)> п,

то письмо относим к классу «Спам».

В этой задаче порог должен быть большим (а почему?).





Использованная литература

1. Т.Сегаран «Программируем коллективный разум» (фильтрация спама)



