

第9讲 Naïve Bayes

- 概率基础回顾
- 极大似然估计
- Naïve Bayes

Fundamentals of Machine Learning_WANGBIANQIN

概率基础回顾

- 联合概率(joint probability)

表示 A 事件和 B 事件同时发生的概率, $P(A \cap B)$

- 边际概率(marginal probability)

在 A 和 B 的样本空间中, 只看 A 或 B 的概率, 称之边际概率

- 条件概率(conditional probability)


在发生 A 的条件下, 发生 B 的概率, 称为 $P(B|A)$

$$P(B | A) = \frac{P(A \cap B)}{P(A)}, P(A | B) = \frac{P(A \cap B)}{P(B)}$$

概率基础回顾

□ 贝叶斯定理 (Bayes theorem)

$$P(B | A) = \frac{P(A \cap B)}{P(A)}, P(A | B) = \frac{P(A \cap B)}{P(B)}$$


$$P(B | A) = P(A | B) \frac{P(B)}{P(A)}$$

贝叶斯定理是关于随机事件A和B的条件概率和边缘概率的一则定理

概率基础回顾

□ 全概率(Total probability)

假设 $\{B_n: n=1, 2, 3, \dots\}$ 是一个概率空间的有限或可数无限的分割，且每个集合 B_n 是一个可测集合，则对任意事件 A 有全概率公式：

$$\Pr(A) = \sum_n \Pr(A \cap B_n) = \sum_n \Pr(A | B_n) \Pr(B_n)$$

概率基础回顾

□ 乘法法则(Multiplicative rule)

$$P(A \cap B) = P(B) \times P(A|B) = P(A) \times P(B|A)$$

□ 独立事件(Independent events)

设事件 A 和事件 B 满足以下条件：

$$P(A \cap B) = P(A) \times P(B)$$

$$\text{或： } P(A) > 0, P(B|A) = P(B)$$

$$P(B) > 0, P(A|B) = P(A)$$

则称 A 与 B 为独立事件

概率基础回顾

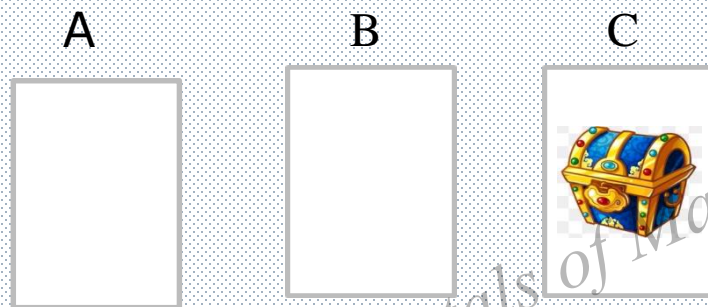
◆ 示例：概率计算

	赞成(B1)	反对(B2)	合计
男性(A1)	40	120	160
女性(A2)	10	30	40
合计	50	150	200

- 联合概率: $P(\text{男性}, \text{赞成}) = ?$
- 条件概率: $P(\text{赞成}|\text{男性}) = ?$
- 边际概率: $P(\text{赞成}) = P(B1) = ?$

概率基础回顾

◆ 示例--换门游戏



嘉宾选定一扇门
未开启

主持人开启
一扇空门

询问嘉宾是否换门？

换门得宝？
不换得宝？

换门
得宝

不换门
得宝

换门
不得宝

不换门
不得宝

极大似然估计

□ 极大似然估计(MLE)

- 输入：模型(全部或者部分未知参数)和样本数据集
- 输出：模型的未知参数
- 处理：根据样本值选择模型参数，使得产生给定样本的可能性(概率)最大

极大似然估计

□ **似然函数(likelihood function)**：一种关于统计模型参数的函数

离散型数据概率分布：假定一个关于参数 θ 、具有离散型概率分布 P 的随机变量 X ，则在给定 X 的输出 x_i 时，参数 θ 的似然函数：

$$L(\theta | X) = P(X | \theta) = P(X; \theta) = P_{\theta}(X) = \prod_{i=1}^N P_{\theta}(X = x_i), X = [x_1, x_2, \dots, x_N]$$

连续型数据概率分布：假定一个关于参数 θ 、具有连续概率密度函数 f 的随机变量 X ，则在给定 X 的输出 x_i 时，参数 θ 的似然函数

$$L(\theta | X) = f_{\theta}(X) = \prod_{i=1}^N f_{\theta}(X = x_i), X = [x_1, x_2, \dots, x_N]$$

$$\text{MLE: } \hat{\theta} = \arg \max_{\theta} L(\theta | x)$$

极大似然估计

□ MLE求解步骤

- ① 写出似然函数；
- ② 对似然函数取对数、整理；
- ③ 求导数；
- ④ 解似然方程

极大似然估计

- ◆ 示例 设总体 $X \sim N(\mu, \sigma^2)$, 样本 (x_1, \dots, x_n) , 其中 μ 和 σ^2 未知, 试求 μ 和 σ^2 的最大似然估计量

$$L = \prod_{i=1}^n \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x_i - \mu)^2}{2\sigma^2}} = \left(\frac{1}{\sqrt{2\pi}\sigma}\right)^n e^{-\sum_{i=1}^n \frac{(x_i - \mu)^2}{2\sigma^2}},$$

$$\ln L = -\frac{n}{2} \ln(2\pi) - \frac{n}{2} \ln \sigma^2 - \frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu)^2,$$

$$\text{令} \begin{cases} \frac{\partial \ln L}{\partial \mu} = \frac{1}{\sigma^2} \sum_{i=1}^n (x_i - \mu) = 0 \\ \frac{\partial \ln L}{\partial \sigma^2} = -\frac{n}{2\sigma^2} + \frac{1}{2(\sigma^2)^2} \sum_{i=1}^n (x_i - \mu)^2 = 0 \end{cases}$$

极大似然估计

令
$$\begin{cases} \frac{\partial \ln L}{\partial \mu} = \frac{1}{\sigma^2} \sum (x_i - \mu) = 0 \\ \frac{\partial \ln L}{\partial \sigma^2} = -\frac{n}{2\sigma^2} + \frac{1}{2(\sigma^2)^2} \sum (x_i - \mu)^2 = 0 \end{cases}$$

解得

$$\hat{\mu} = \frac{1}{n} \sum_{i=1}^n x_i = \bar{X}, \quad \hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{X})^2$$

即为 μ 和 σ^2 的最大似然估计量.

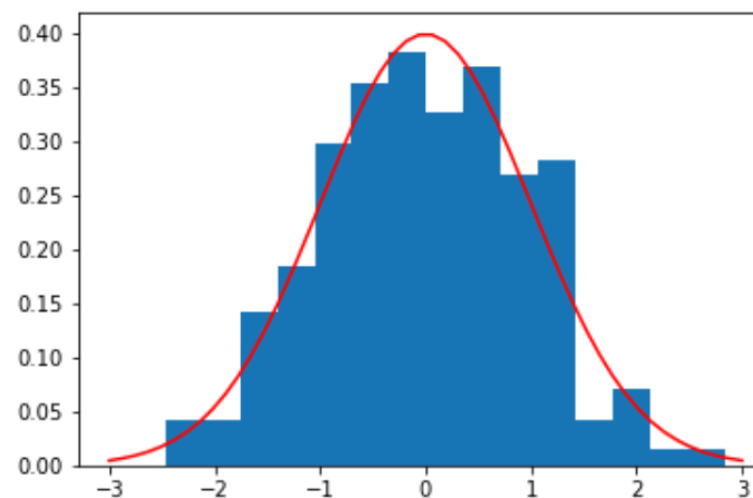
极大似然估计

◆ 示例：MLE

```
from scipy.stats import norm
import matplotlib.pyplot as plt
%matplotlib inline
import numpy as np

x_norm = norm.rvs(size=200)
x_mean, x_std = norm.fit(x_norm)
print('mean, ', x_mean)
print('x_std, ', x_std)
plt.hist(x_norm, normed=True, bins=15)
x = np.linspace(-3, 3, 50)
plt.plot(x, norm.pdf(x), 'r-')
```

```
mean, -0.009373023041395184
x_std, 0.9759251326915617
[<matplotlib.lines.Line2D at 0x1065c7b8>]
```



Naïve Bayes

□ 朴素贝叶斯分类器(Naïve Bayes classifier) :

每个数据样本包含 N 个特征 : $X = \{x_1, x_2, \dots, x_N\}$

属于 m 个不同类别之一 : $C = \{c_1, c_2, \dots, c_m\}$

- X 属于类别的概率:

$$h^*(x) = \arg \max_{c \in y} P(c | X) = \arg \max_{c \in y} \frac{P(X, c)}{P(X)} = \frac{P(X | c)P(c)}{P(X)}$$

似然概率

- 类条件独立假设(在给定类别情况下, 所有特征相互独立)

$$P(X | c) = \prod_{i=1}^N P(x_i | c)$$



贝叶斯(Thomas Bayes, 1702-1761)
英国数学家

Naïve Bayes

示例:

判断
是否
购买
计算
机

age	income	student	credit_rating	buys_computer
youth	high	no	fair	no
youth	high	no	excellent	no
middle_aged	high	no	fair	yes
senior	medium	no	fair	yes
senior	low	yes	fair	yes
senior	low	yes	excellent	no
middle_aged	low	yes	excellent	yes
youth	medium	no	fair	no
youth	low	yes	fair	yes
senior	medium	yes	fair	yes
youth	medium	yes	excellent	yes
middle_aged	medium	no	excellent	yes
middle_aged	high	yes	fair	yes
senior	medium	no	excellent	no

Naïve Bayes

◆ 示例：判断是否购买计算机

$X = (\text{age}=\text{youth}, \text{income}=\text{medium}, \text{student}=\text{yes}, \text{credit_rating}=\text{fair})$

buy computer?

- $P(\text{buy}=\text{yes}|X) > P(\text{buy}=\text{no}|X) ?$
- $P(\text{buy}=\text{yes}|X) = P(X|\text{buy}=\text{yes}) * P(\text{buy}=\text{yes}) = 0.028$
- $P(\text{buy}=\text{no}|X) = P(X|\text{buy}=\text{no}) * P(\text{buy}=\text{no}) = 0.007$
- $P(\text{buy}=\text{yes}) = 9/14 = 0.643, P(\text{buy}=\text{no}) = 5/14 = 0.357$
- $P(\text{age}=\text{youth}|\text{buy}=\text{yes}) = 2/9 = 0.222$
- $P(X|\text{buy}=\text{yes}) = P(\text{age}=\text{youth}|\text{buy}=\text{yes}) * P(\text{income}=\text{medium}|\text{buy}=\text{yes}) * P(\text{student}=\text{yes}|\text{buy}=\text{yes}) * P(\text{credit}=\text{fair}|\text{buy}=\text{yes}) = 0.044$
- $P(X|\text{buy}=\text{no}) = 0.019$

Naïve Bayes

- ◆ 示例：假设在文本分类中，有3个类：c1、c2、c3，在指定的1000个训练样本中，某个词语x在各个类别中观测计数分别为0，990，10，则x的概率为
 $0, 0.99, 0.01$
- 零概率问题：在计算实例的概率时，如果某个量x，在训练集中没有出现过，会导致整个实例的概率结果为0。

Naïve Bayes

□ 拉普拉斯平滑(Laplace Smoothing) :

$$P(x_i = x_{i^*} | c) = \frac{N_{Ci^*} + \lambda}{N_c + \alpha \times \lambda}$$

实际中，经常加 $\lambda(1 \geq \lambda \geq 0)$ 代替简单加1，如果对 N 个计数都加上 λ ，这时分母加上 $a * \lambda$ ， a 为类别数

当 $\lambda=1$ 时，称作Laplace平滑，当 $0 < \lambda < 1$ 时，称作Lidstone平滑， $\lambda=0$ 时，不做平滑

Naïve Bayes

- ◆ 示例：假设在文本分类中，有3个类：C1、C2、C3，在指定的1000个训练样本中，某个词语x在各个类别中观测计数分别为0，990，10，则x的概率为

0，0.99，0.01

拉普拉斯平滑：

$1/1003 = 0.001$ ， $991/1003 = 0.988$ ， $11/1003 = 0.011$

Naïve Bayes

◆ 示例：根据如下训练数据集建模判断性别

性别	身高(英尺)	体重(磅)	脚的尺寸(英寸)
男	6	180	12
男	5.92	190	11
男	5.58	170	12
男	5.92	165	10
女	5	100	6
女	5.5	150	8
女	5.42	130	7
女	5.75	150	9

Naïve Bayes

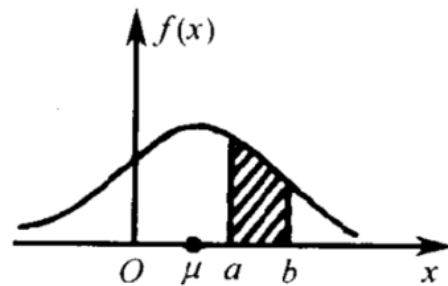
- 概率密度函数 (probability density function , PDF)
 $f(x)$ 应该满足下述两个条件：

$$f(x) \geq 0$$

$$\int_{-\infty}^{+\infty} f(x)dx = 1$$

随机变量X在a与b之间的概率

$$P(a < X < b) = \int_a^b f(x)dx$$



(b)

Naïve Bayes

- 假定连续值特征的概率服从高斯分布（正态分布）

高斯分布的概率密度函数

$$f(x, \mu, \sigma) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}}, \quad -\infty < x < \infty$$

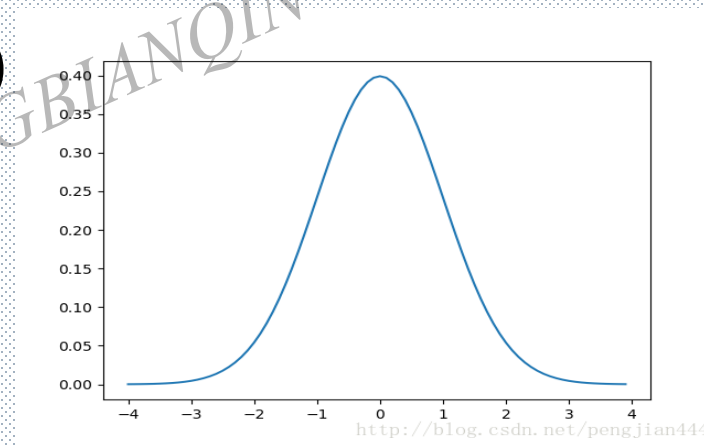
用来反映各特征值的相对可能性

$$P(X | c) = f(x, \mu, \sigma)$$

估算训练集中不同类别的不同特征值的均值、方差

$$\hat{\mu} = \frac{1}{n} \sum_{i=1}^n x_i = \bar{X},$$

$$\hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{X})^2$$



Naïve Bayes

◆ 示例：判断性别

性别	均值(身高)	方差(身高)	均值(体重)	方差(体重)	均值(脚的尺寸)	方差(脚的尺寸)
男性	5.855	3.5033e-02	176.25	1.2292e+02	11.25	9.1667e-01
女性	5.4175	9.7225e-02	132.5	5.5833e+02	7.5	1.6667e+00

Naïve Bayes

◆ 示例：判断性别

预测：某人身高6英尺，体重130磅，脚的尺寸为8英尺，
请推断某人是男性还是女性？

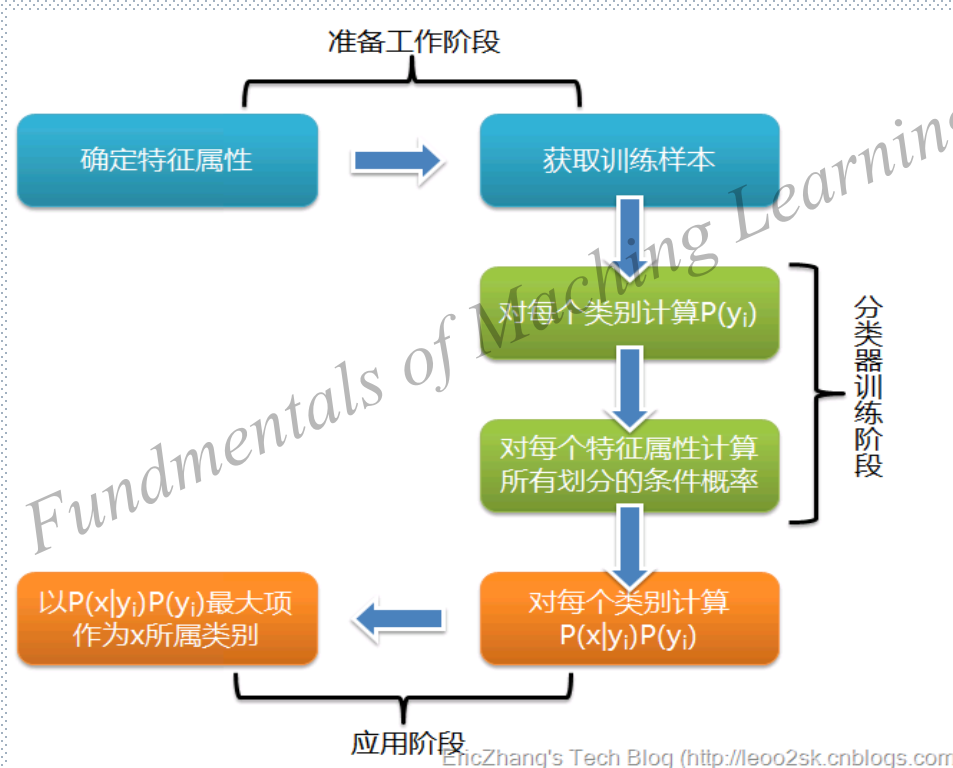
- $P(\text{male}) = P(\text{female}) = 0.5$
- $P(\text{male}/X) = P(\text{male})p(\text{height}|\text{male})p(\text{weight}|\text{male})p(\text{footsize}|\text{male})$
- $P(\text{female}/X) = P(\text{female})p(\text{height}|\text{female})p(\text{weight}|\text{female})p(\text{footsize}|\text{female})$

$$P(\text{male}/X) = 6.1984\text{e}^{-09}$$

$$P(\text{female}/X) = 5.3778\text{e}^{-04}$$

Naïve Bayes

□ Naïve Bayes流程



EricZhang's Tech Blog (<http://leo2sk.cnblogs.com>)

sklearn.naive_bayes

□ sklearn.naive_bayes

https://scikit-learn.org/stable/modules/classes.html#module-sklearn.naive_bayes

naive_bayes.BernoulliNB

naive_bayes.MultinomialNB

naive_bayes.GaussianNB

sklearn.naive_bayes

- **BernoulliNB** : 适合伯努利分布的数据集，即适合类特征为二元值的情况

伯努利分布(Bernoulli distribution)：又称两点分布或0-1分布
 n 次伯努利试验，每次试验的结果只有两种

$$f(x) = p^x(1-p)^{1-x} = \begin{cases} p & \text{if } x = 1 \\ 1-p & \text{if } x = 0 \end{cases}$$

二项分布(Binomial distribution)：

n 重伯努利试验成功次数的离散概率分布

$$P\{X = k\} = C_n^k p^k (1-p)^{n-k}, k = 0, 1, 2, \dots, n$$

sklearn.naive_bayes

❑ **class** sklearn.naive_bayes.**BernoulliNB**(* , alpha=1.0, binarize=0.0, fit_prior=True, class_prior=None)

__init__(* , alpha=1.0, binarize=0.0, fit_prior=True, class_prior=None

- Parameters: [alpha](#), fit_prior, class_prior
- Attributes: class_count_, class_log_prior_, classes_, [feature_count_](#), [feature_log_prob_](#), n_features_
- Methods: fit(X, y[, sample_weight]) , predict(X), [predict_proba\(X\)](#), [predict_log_proba\(X\)](#), score(X, y[, sample_weight])

sklearn.naive_bayes

◆ 示例：BernoulliNB

```
import numpy as np

# X包含4个特征: 刮风, 闷热, 多云
X = np.array([[0, 1, 0],
               [1, 1, 1],
               [0, 1, 1],
               [0, 0, 0],
               [0, 1, 1],
               [0, 1, 0],
               [1, 0, 0]])

# y为样本标签, 0代表没下雨, 1代表下雨
y = np.array([0, 1, 1, 0, 1, 0, 0])
```

```
# 构建BernoulliNB分类器
from sklearn.naive_bayes import BernoulliNB
clf = BernoulliNB()
clf.fit(X, y)

# 预测
Next_Day = [[0, 0, 1]]
pre = clf.predict(Next_Day)
if pre == [1]:
    print("下雨")
else:
    print("无雨")
```

```
clf.predict_proba(Next_Day)
```

```
array([[ 0.13848881,  0.86151119]])
```

sklearn.naive_bayes

- **MultinomialNB** : 适合类别特征多项式分布(Multinomial Distribution)的数据集, 二项式分布的推广
 n 次试验, 每次结果有 $m (> 2)$ 个, 且 m 个结果发生的概率互斥且之和为1, 则发生其中一个结果 X 的概率
- ◆ 例如, 扔骰子是典型的多项式分布
若有6点的骰子重复扔 n 次
试问有 k 次都是点数6朝上的概率

$$P\{X = k\} = C_n^k p_6^k (1 - p_6)^{n-k}, k = 0, 1, 2, \dots, n$$

sklearn.naive_bayes

❑ **class** sklearn.naive_bayes.MultinomialNB(*, alpha=1.0, fit_prior=True, class_prior=None)

__init__(*, alpha=1.0, fit_prior=True, class_prior=None)

- Parameters: alpha, fit_prior, class_prior
- Attributes: class_count_, class_log_prior_, classes_, feature_count_, feature_log_prob_, n_features_
coef_, intercept_
- Methods: fit(X, y[, sample_weight]) , predict(X), predict_proba(X), predict_log_proba(X), score(X, y[, sample_weight])

sklearn.naive_bayes

◆ 示例 : MultinomialNB

```
import numpy as np
X = np.random.randint(5, size=(6, 100))
y = np.array([1, 2, 3, 4, 5, 6])
from sklearn.naive_bayes import MultinomialNB
clf = MultinomialNB()
clf.fit(X, y)
MultinomialNB(alpha=1.0, class_prior=None, fit_prior=True)
print(clf.predict(X[2:3]))

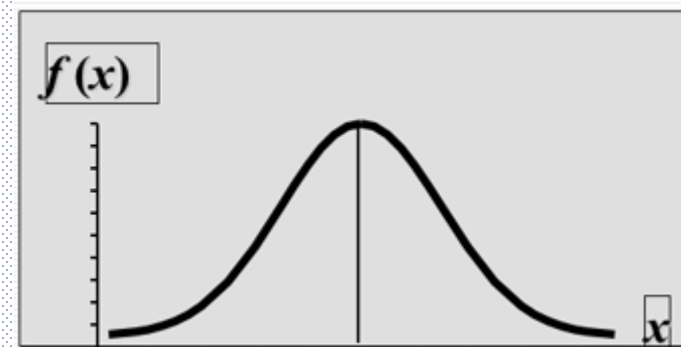
[3]
```


sklearn.naive_bayes

□ **GaussianNB** : 适合高斯分布的数据集，即适合连续特征值

- 高斯分布：也称正态分布(**normal distribution**)
- 连续值特征，概率密度函数：

$$f(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{1}{2\sigma^2}(x-\mu)^2}, \quad -\infty < x < +\infty$$



sklearn.naive_bayes

❑ **class** sklearn.naive_bayes.GaussianNB(*,
priors=None, var_smoothing=1e-09)

- Parameters: [prior\(\)](#), var_smoothing
- Attributes: class_count_, classes_,
[epsilon_](#), [sigma_](#), [theta_](#)
- Methods: fit(X, y[, sample_weight]), predict(X),
predict_proba(X), predict_log_proba(X),
score(X, y[, sample_weight])

sklearn.naive_bayes

◆ 示例：GaussianNB

```
# 导入数据集生成工具
from sklearn.datasets import make_blobs
# 生成500个样本数据, 5个类别
X, y = make_blobs(n_samples=500, centers=5, random=
```

```
print(X)
```

```
[[ -4.43344765e+00  -9.14511574e+00]
 [ -5.06998128e+00  -9.75464122e+00]
 [  6.54464509e+00   8.99873511e-01]
 [  3.25023324e-01   1.50633915e-01]
 [ -1.51028157e+00  -1.10581275e+00]
 [ -8.90489310e+00  -1.10427432e+01]
 [  9.28383472e-02  -2.00771121e-02]
 [ -6.21720086e+00  -1.11227678e+01]
 [  7.63027116e+00   8.69797933e+00]
 [  7.92430026e+00   1.04511206e-01]
 -
```

```
print(y)
```

```
[2 2 1 4 3 2 4 2 0 1 2 0 2 0 0 2 3 0 2 3 0 4 1 3 4 1 4 4 2 3 3 1 0 0 4 1 2
 0 2 4 4 3 1 0 3 0 0 3 3 2 2 4 4 4 2 0 3 1 0 0 4 0 4 3 2 0 0 3 2 0 0 0 3 2
 1 1 3 2 0 3 1 3 1 3 3 2 3 3 1 4 0 2 4 3 3 4 1 2 4 2 3 2 2 1 3 1 2 0 0 4 1
 4 4 0 2 1 0 4 0 1 0 0 0 3 0 2 0 3 1 3 0 1 2 0 3 4 1 1 2 1 2 1 2 2 3 2 0 4
 1 0 2 0 1 1 4 1 2 2 2 4 3 3 0 2 2 4 4 4 0 1 4 1 0 0 2 0 3 4 3 0 0 1 3 4 4
 0 2 4 4 3 0 1 1 1 4 2 0 1 4 2 2 1 3 1 3 4 3 3 2 3 3 2 3 4 3 3 1 1 3 0 4 4
 4 1 3 0 4 4 0 2 4 1 0 2 4 0 1 3 4 0 4 1 3 3 1 1 3 3 1 4 3 2 3 0 2 2 1 1 1
 2 2 0 1 2 4 3 0 2 3 0 2 0 4 4 1 0 0 4 1 1 1 4 3 0 1 2 1 1 0 3 3 2 2 1 4 3
 4 4 2 0 1 4 3 0 3 0 2 3 2 1 4 0 1 3 2 2 3 2 1 2 1 1 4 4 1 3 4 2 4 4 0 1 3
 0 1 1 2 1 3 4 0 2 2 3 1 4 2 3 1 1 4 0 3 4 0 2 3 2 0 1 1 3 3 0 2 2 1 2 0 0
 2 4 0 1 2 2 1 0 4 3 1 1 2 2 1 1 2 1 0 0 3 2 2 4 4 2 2 3 1 0 4 0 4 3 3 0 3
 1 0 1 4 4 4 1 3 0 2 2 4 0 3 2 2 1 4 4 4 0 4 3 2 4 3 3 1 3 0 0 2 0 0 3 3 2
 2 3 1 4 0 4 3 2 4 0 3 3 3 0 0 4 3 1 1 4 3 3 3 2 2 1 4 2 4 0 1 0 4 4 0 4 4
 3 3 4 2 4 1 4 0 1 0 1 4 4 4 1 2 1 1 3]
```

sklearn.naive_bayes

◆ 示例：GaussianNB

```
# 导入数据集生成工具
from sklearn.datasets import make_blobs
# 生成500个样本数据, 5个类别
X, y = make_blobs(n_samples=500, centers=5, random_state=8)

# 导入数据集拆分工具
from sklearn.model_selection import train_test_split

# 数据拆分成训练集和测试集
X_train, X_test, y_train, y_test = train_test_split(X, y, random_state=8)

# 构建GaussianNB分类器
from sklearn.naive_bayes import GaussianNB
gnb = GaussianNB()
gnb.fit(X_train, y_train)
print('模型得分: {:.3f}'.format(gnb.score(X_test, y_test)))

模型得分: 0.968
```

sklearn.naive_bayes

□ sklearn的分类模型可信度评估

- **.predict()**直接获得唯一的预测结果
- **.predict_proba()**获取属于类别的概率
- **.predict_log_proba()**获取属于类别的概率的对数

sklearn.naive_bayes

◆ 示例：.predict()直接获得唯一预测结果

```
# predict

# 导入GaussianNB模型
from sklearn.naive_bayes import GaussianNB
gnb = GaussianNB()

# 拆分数据集
from sklearn.model_selection import train_test_split
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=0)

# 训练模型
gnb.fit(X_train, y_train)

# 获取模型预测值
predict = gnb.predict(X_test)
print('预测值形态: {}'.format(predict.shape))
print(predict[:5])
```

```
预测值形态: (50,)
[0 1 1 0 1]
```

```
# 坐标轴范围
x_min, x_max = X[:, 0].min() - 0.5, X[:, 0].max() + 0.5
y_min, y_max = X[:, 1].min() - 0.5, X[:, 1].max() + 0.5

xx, yy = np.meshgrid(np.arange(x_min, x_max, 0.02),
                     np.arange(y_min, y_max, 0.02))

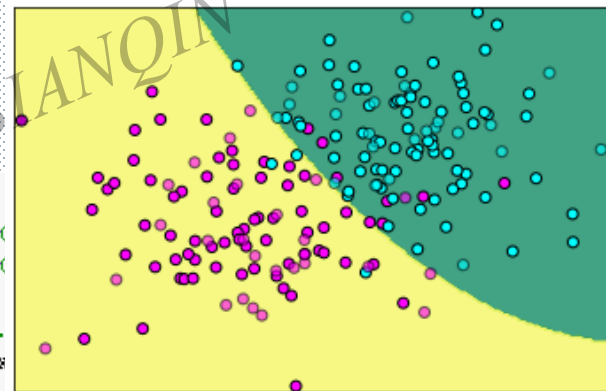
# 用不同色彩表示不同分类
Z = gnb.predict(np.c_[xx.ravel(), yy.ravel()])
Z = Z.reshape(xx.shape)

# 绘制等高线
plt.contourf(xx, yy, Z, cmap=plt.cm.summer, alpha=0.8)

# 绘制散点图
plt.scatter(X_train[:, 0], X_train[:, 1], c=y_train, cmap=plt.cm.cool, edgecolor='k')
plt.scatter(X_test[:, 0], X_test[:, 1], c=y_test, cmap=plt.cm.cool, edgecolor='k', alpha=0.6)

# 坐标轴范围
plt.xlim(xx.min(), xx.max())
plt.ylim(yy.min(), yy.max())

# 坐标轴单位
plt.xticks(())
plt.yticks(())
```



sklearn.naive_bayes

□ sklearn的一些分类模型中的预测概率

.predict_proba()

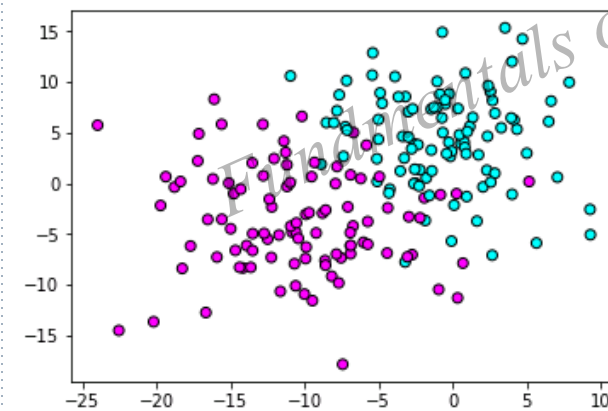
- 计算分类时每个样本属于不同类别概率
- 对于二元分类，它的形状是(**n_samples**, **2**)
- 对于多元分类，它的形状是(**n_samples**, **n_classes**)

sklearn.naive_bayes

◆ 示例：predict_proba()方法

```
import matplotlib.pyplot as plt
%matplotlib inline
import numpy as np
# 导入数据集生成工具
from sklearn.datasets import make_blobs
X, y = make_blobs(n_samples=200, random_state=1, centers=[[-10, 0], [10, 0]], cluster_std=10)
plt.scatter(X[:, 0], X[:, 1], c=y, cmap=plt.cm.cool, s=50)
```

<matplotlib.collections.PathCollection at 0x1091a9...



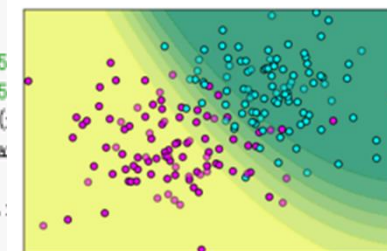
```
# 导入GaussianNB模型
from sklearn.naive_bayes import GaussianNB
gnb = GaussianNB()
# 拆分数据集
from sklearn.model_selection import train_test_split
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3, random_state=1)
# 训练模型
gnb.fit(X_train, y_train)
# 获取模型预测概率
predict_proba = gnb.predict_proba(X_test)
print('预测准确率形态: {}'.format(predict_proba[:5]))
```

```
预测准确率形态: (50, 2)
[[0.98849996 0.01150004]
 [0.0495985  0.9504015 ]
 [0.01648034 0.98351966]
 [0.8168274  0.1831726 ]
 [0.00282471 0.99717529]]
```

```
# predict_proba可简化
# 坐标轴范围
x_min, x_max = X[:, 0].min() - 0.5, X[:, 0].max() + 0.5
y_min, y_max = X[:, 1].min() - 0.5, X[:, 1].max() + 0.5
xx, yy = np.meshgrid(np.arange(x_min, x_max, 0.05), np.arange(y_min, y_max, 0.05))
Z = gnb.predict_proba(np.c_[xx.ravel(), yy.ravel()])
Z = Z.reshape(xx.shape)
```

```
# 绘制等高线
plt.contourf(xx, yy, Z, cmap=plt.cm.summer, alpha=0.8)
# 绘制散点图
plt.scatter(X_train[:, 0], X_train[:, 1], c=y_train, cmap=plt.cm.cool, edgecolor='k')
plt.scatter(X_test[:, 0], X_test[:, 1], c=y_test, cmap=plt.cm.cool, edgecolor='k', alpha=0.6)
# 坐标轴范围
plt.xlim(xx.min(), xx.max())
plt.ylim(yy.min(), yy.max())
# 坐标轴单位
plt.xticks(())
plt.yticks(())
```

([], <a list of 0 Text yticklabel objects>)



Bayes Classifier

- Naïve Bayes分类假定类条件独立，即给定样本的类标号，属性的值相互条件独立，也因此称为“朴素”。
- 实际中变量之间可能存在依赖。贝叶斯信念网络描述联合条件概率分布，它允许在变量的子集间定义类条件独立性，它提供一种因果关系的图形。
- 贝叶斯信念网络可用于分类，它是图形模型。相比Naïve Bayes，它能够处理属性子集间有依赖关系的分类。