



《大数据分析方法》

(2021年秋季学期)

翟祥
北京林业大学

E-mail: zhaixbh@126.com



万物充满智慧，一切皆可学习



机器学习之生成模型

万物充满智慧，一切皆可学习

知山知水 树木树人



生成模型

- 生成模型和判别模型
- 概率生成模型
- 性能对比
- 其他生成模型
- 补充知识

机器学习中的模型分类

机器学习（统计学习）当中，按照模型的产生方式，可以分为两大类，分别是：

1）判别模型（Discriminative Model）

目前，我们所见到的模型当中，判别模型居多，回归（感知机）、logistic回归、SVM、KNN、决策树、前馈神经网络、DNN中的CNN等。

2）生成模型（Generative Model）

以概率、随机过程和贝叶斯理论基础，主要有概率生成模型（logistic回归）、朴素贝叶斯分类、贝叶斯网络、高斯混合模型（GMM）隐马尔科夫模型（HMM）、GAN等。

相关概念

判别模型：由数据直接学习函数 $Y=f(X)$ 或者条件概率分布 $P(Y|X)$ 作为预测的模型，即判别模型。基本思想是有限样本条件下建立判别函数（模型），不考虑样本的概率特性特性，直接产生预测模型。

生成模型：由数据学习联合概率密度分布 $P(X,Y)$ ，然后求出条件概率分布 $P(Y|X)$ 作为预测的模型，即生成模型： $P(Y|X)= P(X,Y)/ P(X)$ 。基本思想是首先建立样本的联合概率密度模型 $P(X,Y)$ ，然后再得到后验概率 $P(Y|X)$ ，从而产生预测模型，然后进行预测。

一个典型示例

以一个分类模型为例---来自维基百科

我们有一组数据

y中0代表‘NO’，1代表‘YES’

Input	Output
x	y
0	No
0	No
1	No
1	Yes

Input	Output
x	y
0	0
0	0
1	0
1	1

当一种模型同时有两种形式的时候，生成模型可以导出判别模型，而判别模型不能反推出生成模型。因此，从随机数学角度来看生成模型更“原始”，利于衍生出更多的模型和算法。

判别模型只关注条件分布

$$P(y|x) \quad \text{其中} \left(\sum_y P(y|x) = 1 \right)$$

		y	
		0	1
x	0	1	0
	1	1/2	1/2

生成模型从关注联合分布 **开始!!!**

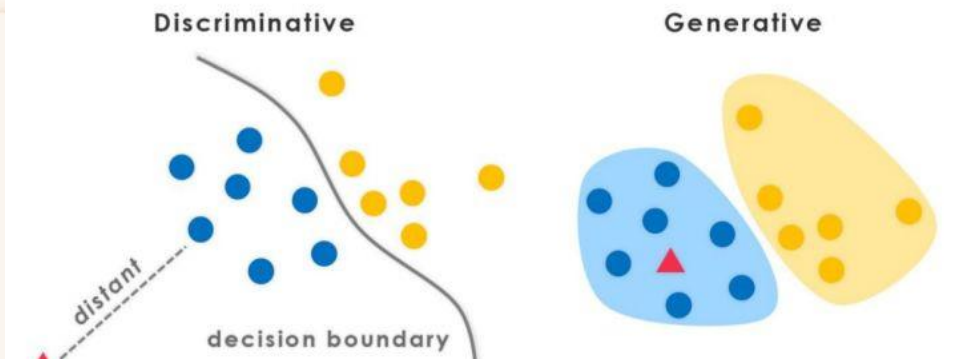
$$P(x, y) \quad \text{其中} \left(\sum_x \sum_y P(x, y) = 1 \right)$$

		y	
		0	1
x	0	1/2	0
	1	1/4	1/4

模型对比

对比	判别式模型	生成式模型
特点	寻找不同类别之间的最优分类面，反映异类数据之间的差异	以概率的角度表示数据的分布情况，能够反映同类数据本身的相似度
区别（输入 x ,输出 y ）	估计的是条件概率分布： $P(y x)$	估计的是联合概率分布 $P(x,y)$
联系	由判别式模型不能得到生成式模型	由生成式模型可以得到判别式模型（贝叶斯公式）
优势	(1)能清晰地分辨出多类或某一类与其他类之间的差异特征；（2）适用于较多类别的识别；（3）模型更简单	（1）研究单类问题比判别式模型更灵活；（2）模型可以通过增强学习得到；（3）能用于数据不完整的情况。
缺点	不能反映训练数据本身的特性；	学习和计算过程比较复杂
预测性能	较好（因为利用了训练数据的类别标识信息）	比判别模型要差？？？
常见模型举例	KNN，SVM，决策树，线性回归，逻辑回归，boosting，线性判别分析（LDA），感知机，传统神经网络	概率生成模型，朴素贝叶斯，隐马尔科夫模型，高斯混合模型，限制玻尔兹曼机

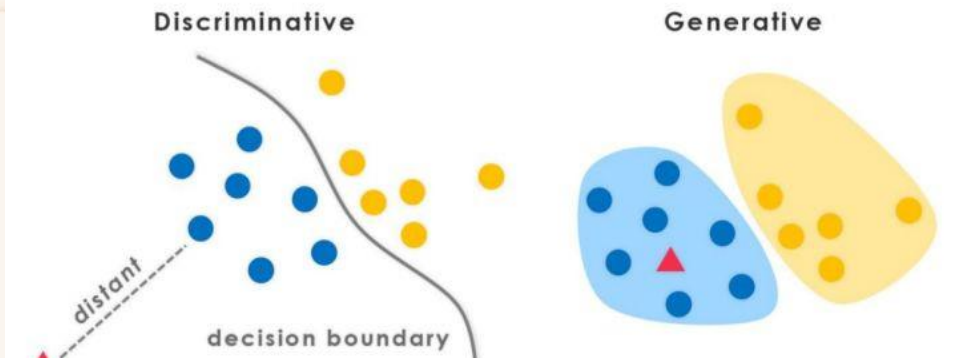
Discriminative vs. Generative



判别方法直接学习的是函数 $Y=f(X)$ 或者条件概率分布 $P(Y|X)$ 。不能反映训练数据本身的特性。但它寻找不同类别之间的最优分类面，反映的是异类数据之间的差异。直接面对预测，往往学习的准确率更高。由于直接学习 $P(Y|X)$ 或 $P(X)$ ，可以对数据进行各种程度上的抽象、定义特征并使用特征，因此可以简化学习问题。

对于判别式模型来说求得 $P(Y|X)$ ，对未知观测 X ，根据 $P(Y|X)$ 可以求得标记 Y ，即可以直接判别出来，如上图的左边所示，实际是就是直接得到了判别边界，所以传统的、耳熟能详的机器学习算法和模型都是判别式模型这些模型的特点都是输入属性 X 可以直接得到 Y ，对于二分类任务来说，实际得到一个score（概率），当score大于某个阈值（threshold）时则为正类，否则为负类。

Discriminative vs. Generative



首先建立样本的联合概率密度模型 $P(X,Y)$ ，然后再得到后验概率 $P(Y|X)$ ，再利用它进行预测。因此，生成模型可以导出判别模型，反之不成立。!!!

生成方法学习联合概率密度分布 $P(X,Y)$ ，所以就可以从概率的角度表示数据的分布情况，能够反映同类数据本身的相似度。但它不关心到底划分各类的那个分类边界在哪。生成方法可以还原出联合概率分布 $P(Y|X)$ ，而判别方法不能。生成方法的学习收敛速度更快，即当样本容量增加的时候，学到的模型可以更快的收敛于真实模型，当存在隐变量时，仍可以用生成方法学习。此时判别方法就不能用。

生成式模型求得 $P(X, Y)$ ，对于未知观测 X ，你要求出 X 与不同类之间的联合概率分布，然后大的类获胜，如上图右边所示，并没有什么边界存在，对于未知观测（红三角），分别求两个联合概率分布（有两个类），比较一下，取概率大的那个分布，也就是相应的类别或值。

如何区分两种建模方式（ approach ）---通俗来讲

生成算法尝试去找到底这个数据是怎么生成的（产生的），然后再对一个信号进行分类。基于你的生成假设，那么那个类别最有可能产生这个信号，这个信号就属于那个类别。判别模型不关心数据是怎么生成的，它只关心信号之间的差别，然后用差别来简单对给定的一个信号进行分类。

假如你的任务是识别一个语音属于哪种语言。例如对面一个人走过来，和你说了一句话，你需要识别出她说的到底是汉语、英语还是法语等。那么你可以有两种方法达到这个目的：

1、学习每一种语言，你花了大量精力把汉语、英语和法语等都学会了，我指的学会是你知道什么样的语音对应什么样的语言。然后再有人过来对你哄，你就可以知道他说的是什么语音。

2、不去学习每一种语言，你只学习这些语言之间的差别，然后再进行区分。意思是指我学会了汉语和英语等语言的发音是有差别的，我学会这种差别就好了。

第一种方法就是生成式，第二种方法是判别式。

两种方式的对比

生成模型

1) 生成给出的是**联合分布**，不仅能够由联合分布计算条件分布（反之则不行），还可以给出其他信息，比如可以使用 **$P(x)=\sum_i P(x|c_i)$** 来计算边缘分布 **$P(x)$** 。如果一个输入样本的边缘分布 **$P(x)$** 很小的话，那么可以认为学习出的这个模型可能不太适合对这个样本进行分类，分类效果可能会不好，这也是所谓的**outlier detection**。生成方法可以原出联合概率分布分布 **$P(X,Y)$** ，而判别方法不能。

2) 生成模型收敛速度比较快，即当样本数量较多时，生成模型能更快地收敛于真实模型。

3) 生成模型能够应付存在**隐变量**的情况，比如**混合高斯模型**就是含有隐变量的生成方法。此时判别方法就不能用。

1) 天下没有免费午餐，联合分布是能提供更多的信息，但也需要更多的样本和更多计算，尤其是为了更准确估计类别条件分布，需要增加样本的数目，而且类别条件概率的许多信息是我们做分类用不到，因而如果我们只需要做分类任务，就浪费了计算资源。

2) 另外，实践中多数情况下判别模型效果更好。

判别模型

1) 由于直接学习 **$P(Y|X)$** 或 **$f(X)$** ，可以对数据进行各种程度上的抽象、定义特征并使用特征，因此可以简化学习问题。与生成模型缺点对应，首先是节省计算资源，另外，需要的样本数量也少于生成模型。

2) 直接面对预测，准确率往往较生成模型高。

3) 判别方法直接学习的是决策函数 **$Y=f(X)$** 或者条件概率分布 **$P(Y|X)$** ，不能反映训练数据本身的特性。但它寻找不同类别之间的最优分类面，反映的是异类数据之间的差异。所以允许我们对输入进行抽象（比如降维、构造等），从而能够简化学习问题。

1) 不具备更加细致的附加能力，只能进行预测，虽然效果更好，但是容易过拟合。

2) 没有额外的诊断功能

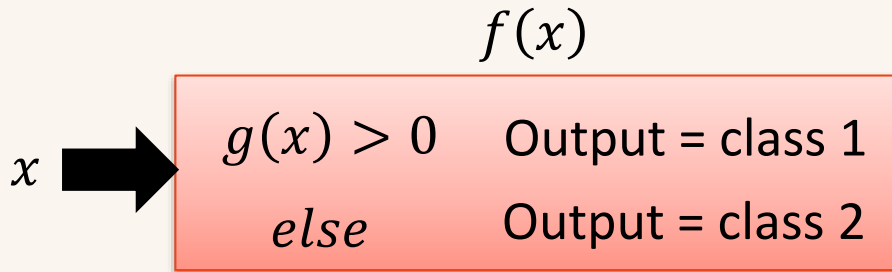


生成模型

- 生成模型和判别模型
- 概率生成模型
- 性能对比
- 其他生成模型
- 补充知识

概率生成模型(Probabilistic Generative Model)

- Function (Model):



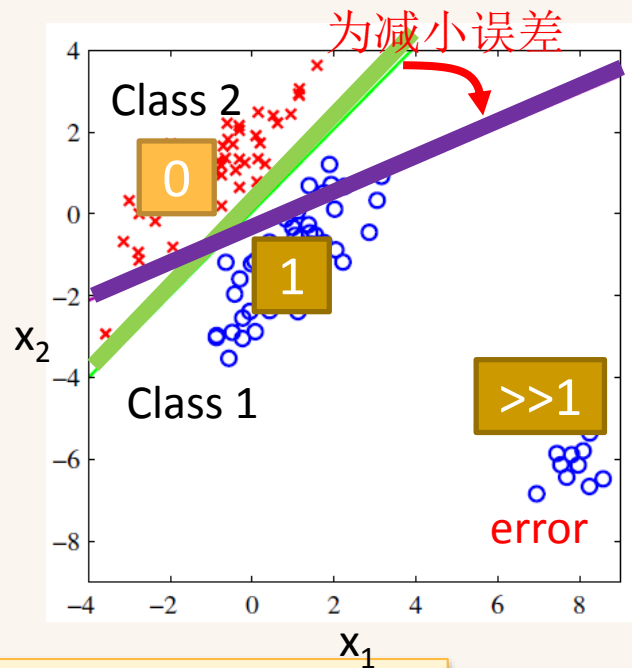
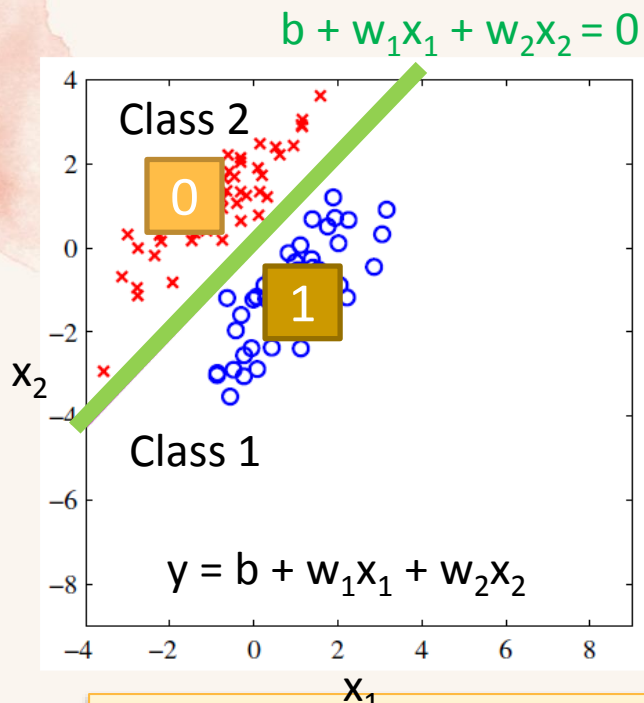
- 损失函数:

$$L(f) = \sum_n \delta(f(x^n) \neq \hat{y}^n)$$

The number of times f get incorrect results on training data.

- 获得最优模型:
 - Example: 感知器(Perceptron), SVM

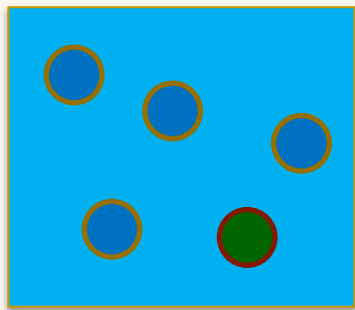
像回归那样进行分类



惩罚那些太正确的观测

两个盒子

Box 1

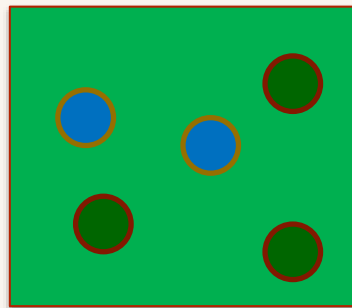


$$P(B_1) = 2/3$$

$$P(\text{Blue} | B_1) = 4/5$$

$$P(\text{Green} | B_1) = 1/5$$

Box 2



$$P(B_2) = 1/3$$

$$P(\text{Blue} | B_2) = 2/5$$

$$P(\text{Green} | B_2) = 3/5$$

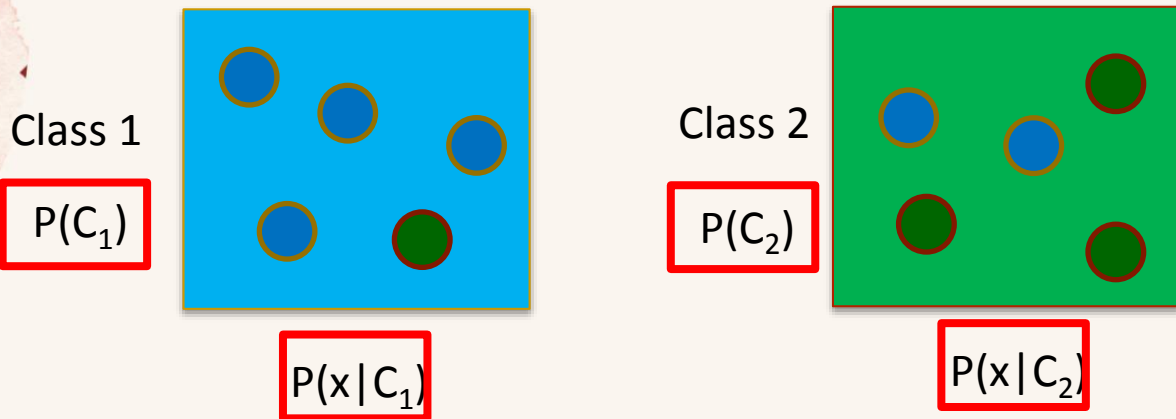
● 来自于其中一个盒子

那个呢?

$$P(B_1 | \text{Blue}) = \frac{P(\text{Blue} | B_1)P(B_1)}{P(\text{Blue} | B_1)P(B_1) + P(\text{Blue} | B_2)P(B_2)}$$

两个类别

从训练数据中估计概率。

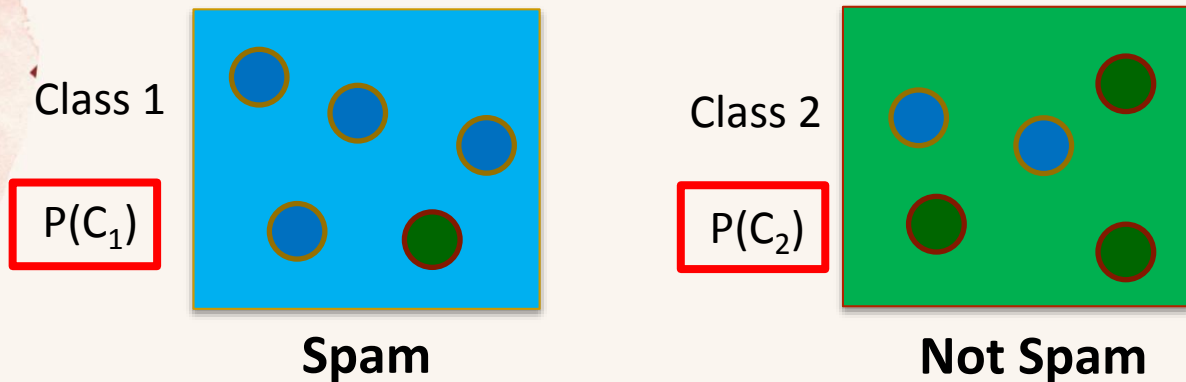


给一个 x , 他属于那个类别呢?

$$P(C_1|x) = \frac{P(x|C_1)P(C_1)}{P(x|C_1)P(C_1) + P(x|C_2)P(C_2)}$$

生成模型(Generative Model) $P(x) = P(x|C_1)P(C_1) + P(x|C_2)P(C_2)$

先验信息(Prior)



Spam和Not Spam 的数据中一部分作为训练集,
其他为测试集。

Training: 79 Spam, 61 Not Spam

$$P(C_1) = 79 / (79 + 61) = 0.56$$

$$P(C_2) = 61 / (79 + 61) = 0.44$$

来自类别的概率

$$P(x|C_1) = ?$$

P(

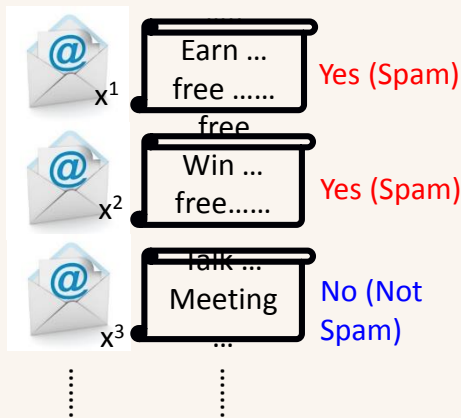


$$|\text{Spam}) = ?$$

邮件的各个特征用向量表示

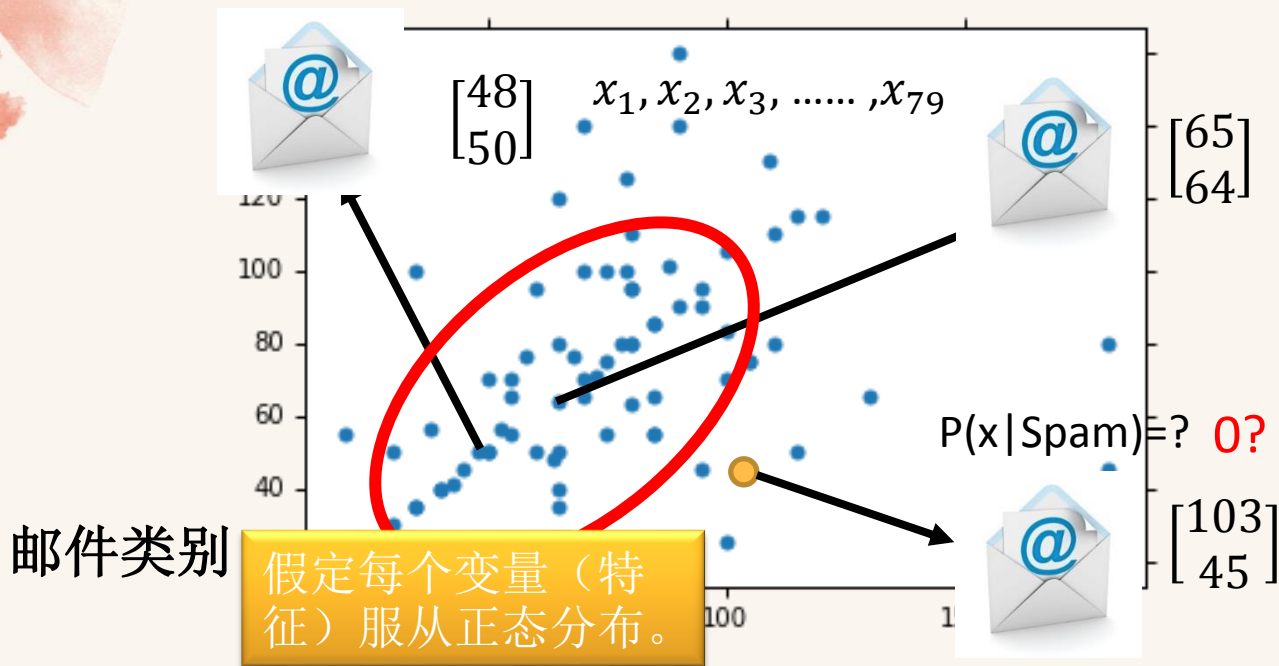


特征（变量）
feature



类别的概率- Feature

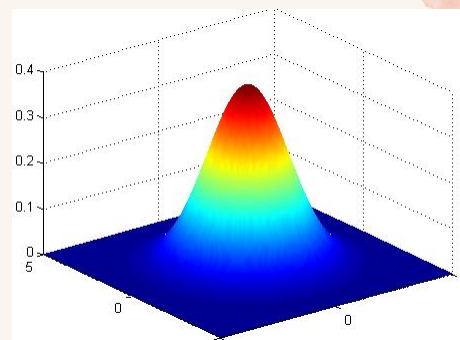
- 考虑两个特征



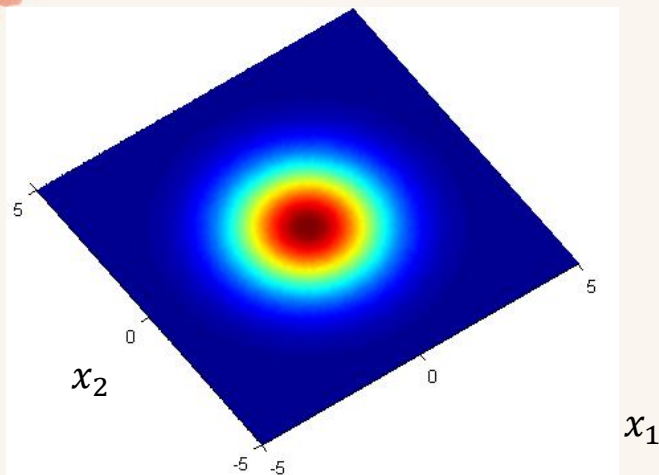
正态分布

$$f_{\mu, \Sigma}(x) = \frac{1}{(2\pi)^{D/2}} \frac{1}{|\Sigma|^{1/2}} \exp\left\{-\frac{1}{2}(x - \mu)^T \Sigma^{-1}(x - \mu)\right\}$$

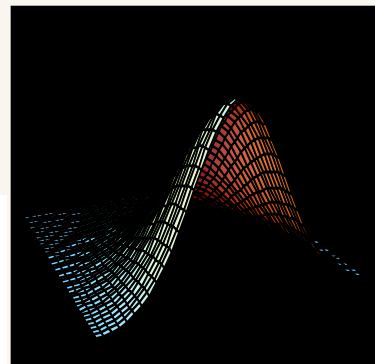
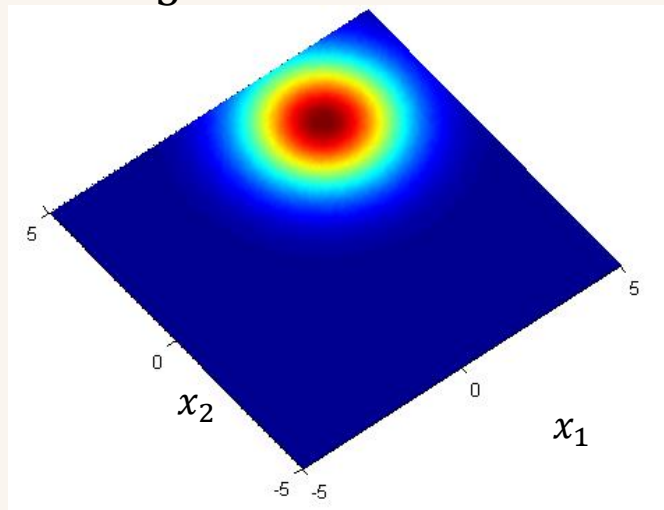
正态分布的形态由均值 μ 和协方差矩阵 Σ 决定。




$$\mu = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \quad \Sigma = \begin{bmatrix} 2 & 0 \\ 0 & 2 \end{bmatrix}$$



$$\mu = \begin{bmatrix} 2 \\ 3 \end{bmatrix} \quad \Sigma = \begin{bmatrix} 2 & 0 \\ 0 & 2 \end{bmatrix}$$

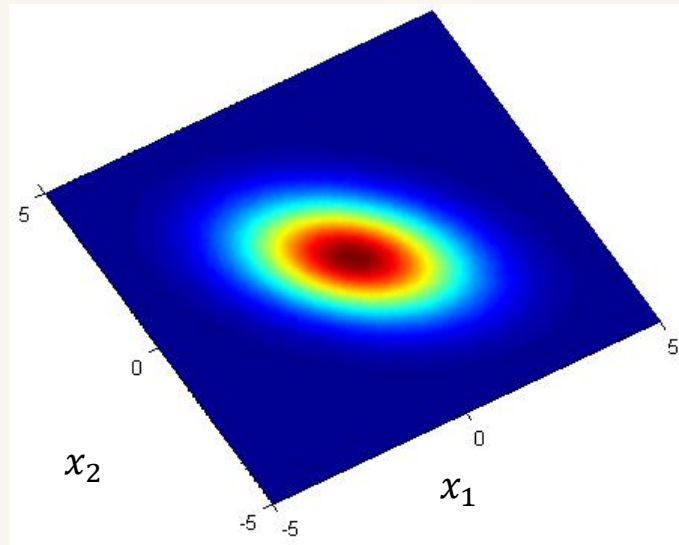
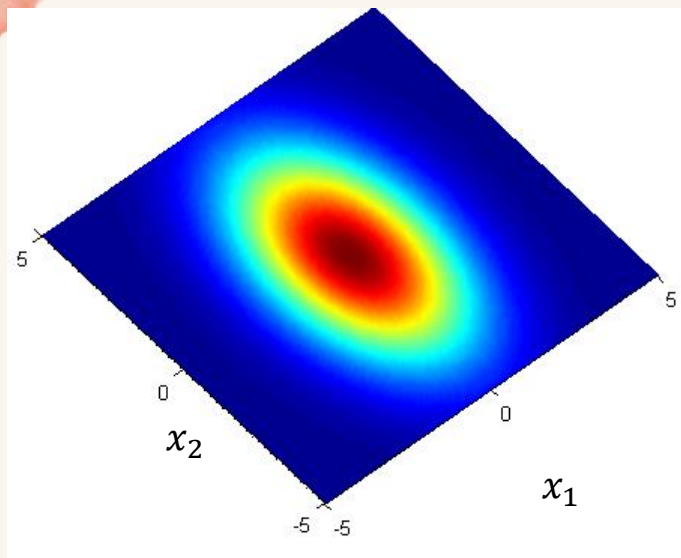


正态分布

$$f_{\mu, \Sigma}(x) = \frac{1}{(2\pi)^{D/2}} \frac{1}{|\Sigma|^{1/2}} \exp \left\{ -\frac{1}{2} (x - \mu)^T \Sigma^{-1} (x - \mu) \right\}$$

正态分布的形态由均值 μ 和协方差矩阵 Σ 决定。

$$\mu = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \quad \Sigma = \begin{bmatrix} 2 & 0 \\ 0 & 6 \end{bmatrix} \quad \mu = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \quad \Sigma = \begin{bmatrix} 2 & -1 \\ -1 & 2 \end{bmatrix}$$

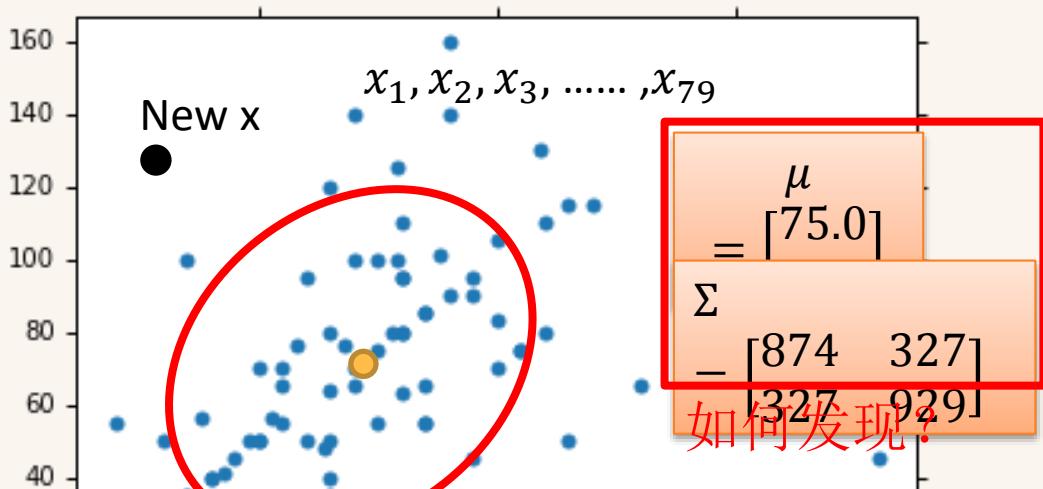


类别中计算概率

假定数据来自正态分布。

在新数据来之前先预测到它。

新观测的概率

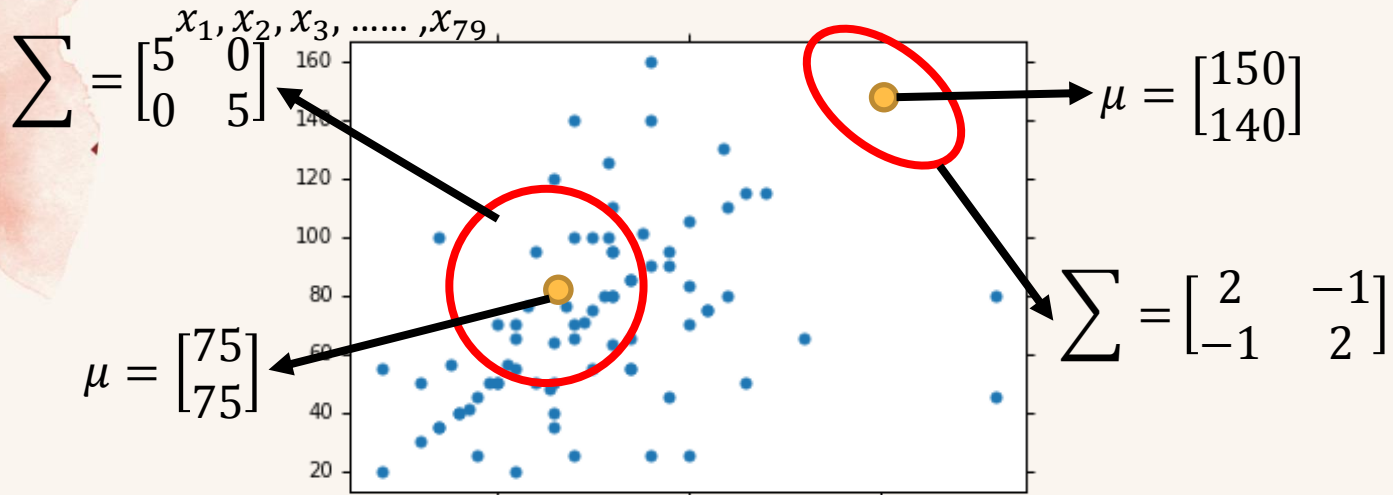


Spam

$$f_{\mu, \Sigma}(x) = \frac{1}{(2\pi)^{D/2}} \frac{1}{|\Sigma|^{1/2}} \exp \left\{ -\frac{1}{2} (x - \mu)^T \Sigma^{-1} (x - \mu) \right\}$$

极大似然 **Maximum Likelihood**

$$f_{\mu, \Sigma}(x) = \frac{1}{(2\pi)^{D/2}} \frac{1}{|\Sigma|^{1/2}} \exp \left\{ -\frac{1}{2} (x - \mu)^T \Sigma^{-1} (x - \mu) \right\}$$



任何均值 μ 和协方差矩阵 Σ 的正态分布都可以产生这些观测。

利用样本 $x_1, x_2, x_3, \dots, x_{79}$, 采用极大似然法来对正态分布的参数进行估计, μ 和 Σ

$$L(\mu, \Sigma) = f_{\mu, \Sigma}(x^1) f_{\mu, \Sigma}(x^2) f_{\mu, \Sigma}(x^3) \dots f_{\mu, \Sigma}(x^{79})$$

极大似然(Maximum Likelihood)

我们有spam类的观测79个: $x_1, x_2, x_3, \dots, x_{79}$

假定 $x_1, x_2, x_3, \dots, x_{79}$ 来自正态分布 (μ^*, Σ^*) , 采用极大似然估计

$$L(\mu, \Sigma) = f_{\mu, \Sigma}(x_1) f_{\mu, \Sigma}(x_2) f_{\mu, \Sigma}(x_3) \dots f_{\mu, \Sigma}(x_{79})$$

$$f_{\mu, \Sigma}(x) = \frac{1}{(2\pi)^{D/2}} \frac{1}{|\Sigma|^{1/2}} \exp \left\{ -\frac{1}{2} (x - \mu)^T \Sigma^{-1} (x - \mu) \right\}$$

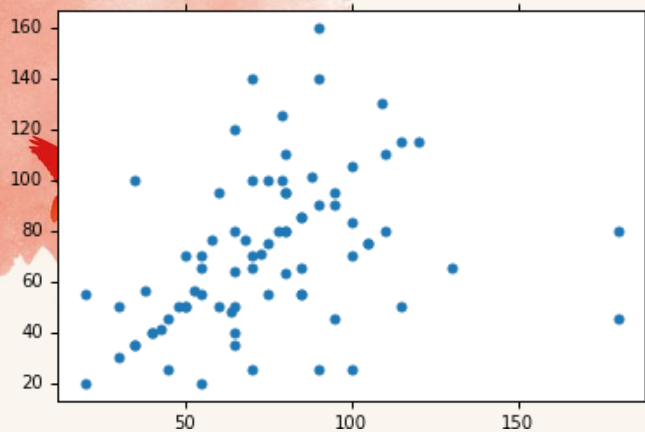
$$\mu^*, \Sigma^* = \arg \max_{\mu, \Sigma} L(\mu, \Sigma)$$

$$\mu^* = \frac{1}{79} \sum_{n=1}^{79} x_n$$

$$\Sigma^* = \frac{1}{79} \sum_{n=1}^{79} (x_n - \mu^*) (x_n - \mu^*)^T$$

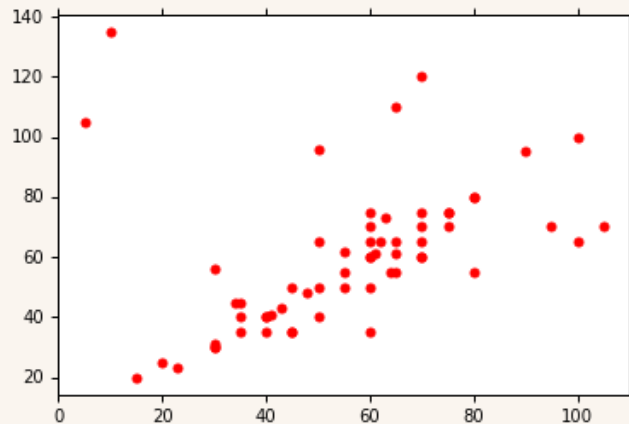
极大似然

Class 1: Spam



$$\mu_1 = \begin{bmatrix} 75.0 \\ 71.3 \end{bmatrix} \quad \Sigma_1 = \begin{bmatrix} 874 & 327 \\ 327 & 929 \end{bmatrix}$$

Class 2: Not Spam



$$\mu_2 = \begin{bmatrix} 55.6 \\ 59.8 \end{bmatrix} \quad \Sigma_2 = \begin{bmatrix} 847 & 422 \\ 422 & 685 \end{bmatrix}$$

通过计算概率做分类

$$f_{\mu_1, \Sigma_1}(x) = \frac{1}{(2\pi)^{D/2}} \frac{1}{|\Sigma_1|^{1/2}} \exp \left\{ -\frac{1}{2} (x - \mu_1)^T (\Sigma_1)^{-1} (x - \mu_1) \right\}$$

$$\mu_1 = \begin{bmatrix} 75.0 \\ 71.3 \end{bmatrix} \quad \Sigma_1 = \begin{bmatrix} 874 & 327 \\ 327 & 929 \end{bmatrix}$$

$P(C_1)$

$$= 79 / (79 + 61) \\ = 0.56$$

$$P(C_1|x) = \frac{P(x|C_1)P(C_1)}{P(x|C_1)P(C_1) + P(x|C_2)P(C_2)}$$

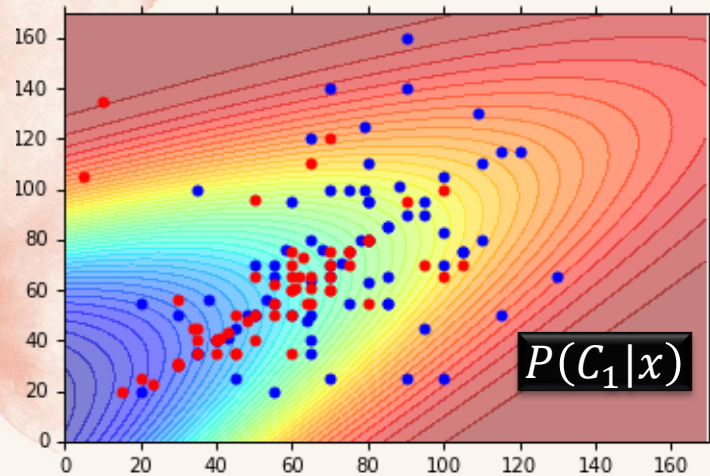
$$f_{\mu_2, \Sigma_2}(x) = \frac{1}{(2\pi)^{D/2}} \frac{1}{|\Sigma_2|^{1/2}} \exp \left\{ -\frac{1}{2} (x - \mu_2)^T (\Sigma_2)^{-1} (x - \mu_2) \right\}$$

$P(C_2)$

$$= 61 / (79 + 61) \\ = 0.44$$

$$\mu_2 = \begin{bmatrix} 55.6 \\ 59.8 \end{bmatrix} \quad \Sigma_2 = \begin{bmatrix} 847 & 422 \\ 422 & 685 \end{bmatrix}$$

If $P(C_1|x) > 0.5$  x 属于class 1 (Spam)



蓝点: C_1 (Spam), 红点: C_2 (Not Spam)

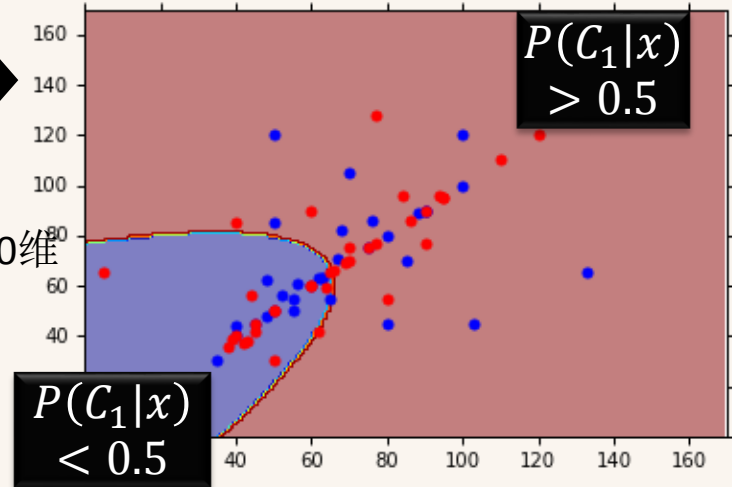
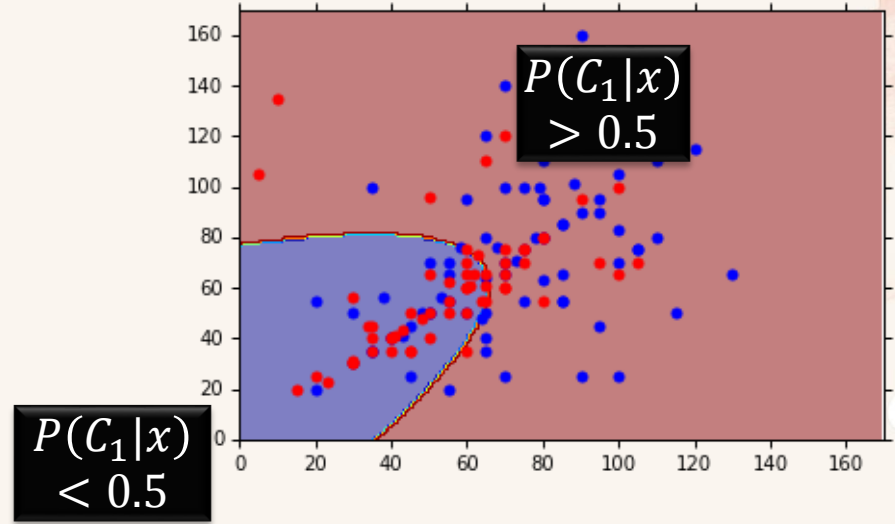
Testing data: 47% accuracy



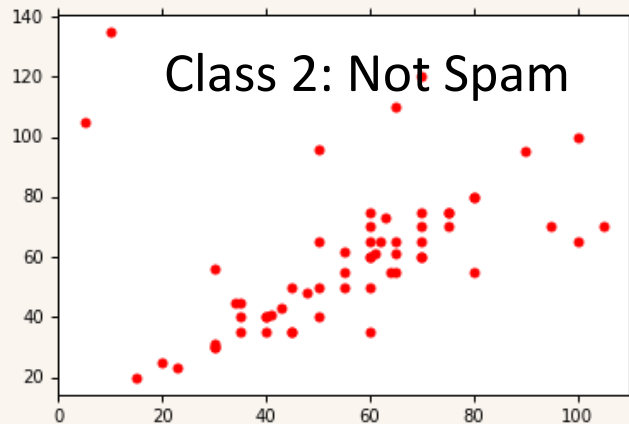
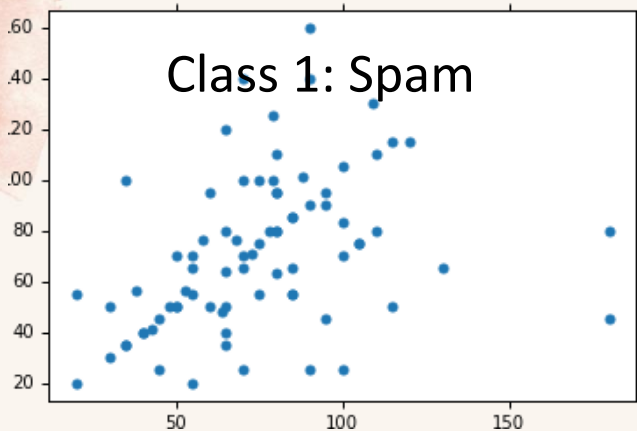
μ_1, μ_2 : 向量, 几个输入变量就是几维, 比如10维

Σ_1, Σ_2 : 10×10 矩阵

64% accuracy ...



模型修正(Modifying Model)



$$\mu_1 = \begin{bmatrix} 75.0 \\ 71.3 \end{bmatrix} \quad \Sigma_1 = \begin{bmatrix} 874 & 327 \\ 327 & 929 \end{bmatrix}$$

$$\mu_2 = \begin{bmatrix} 55.6 \\ 59.8 \end{bmatrix} \quad \Sigma_2 = \begin{bmatrix} 847 & 422 \\ 422 & 685 \end{bmatrix}$$

相等的 Σ
减少参数个数

模型修正(Modifying Model)

- 极大似然

“Spam” :

$x_1, x_2, x_3, \dots, x_{79}$

μ_1

“Not Spam” :

$x_{80}, x_{81}, x_{82}, \dots, x_{140}$

μ_2

Σ

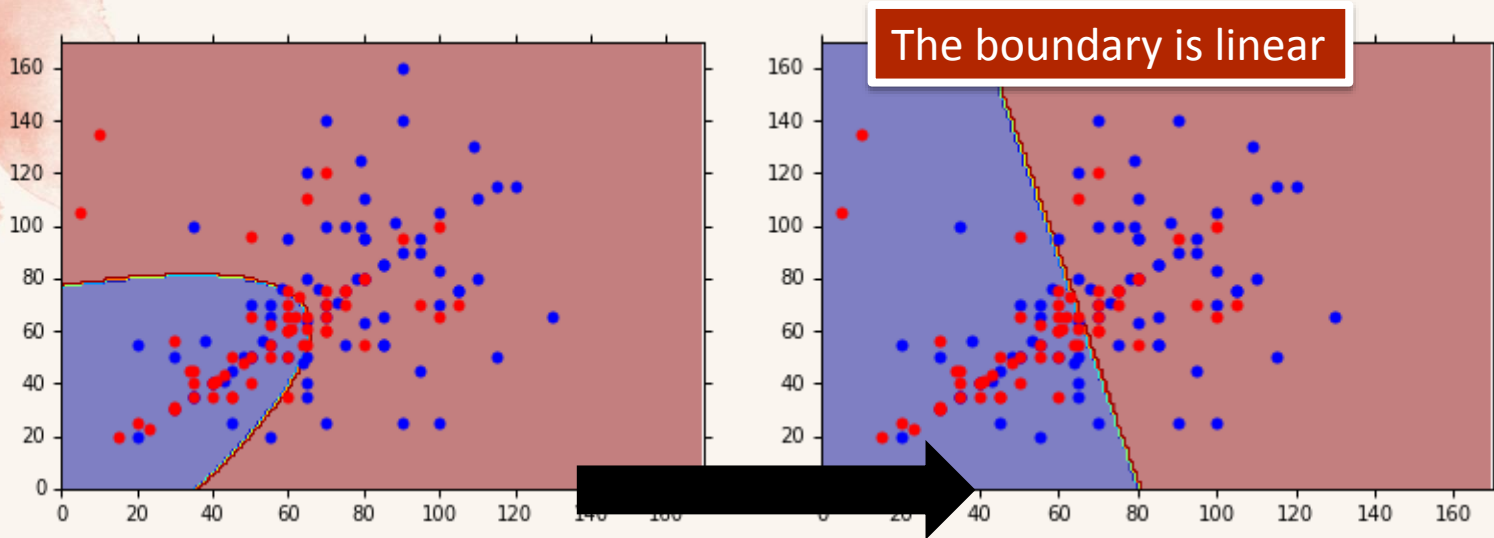
最小化似然函数 $L(\mu_1, \mu_2, \Sigma)$ 获得 μ_1, μ_2, Σ

$$L(\mu_1, \mu_2, \Sigma) = f_{\mu_1, \Sigma}(x_1) f_{\mu_1, \Sigma}(x_2) \cdots f_{\mu_1, \Sigma}(x_{79}) \\ \times f_{\mu_2, \Sigma}(x_{80}) f_{\mu_2, \Sigma}(x_{81}) \cdots f_{\mu_2, \Sigma}(x_{140})$$


如果 μ_1 和 μ_2 相同

$$\Sigma = \frac{79}{140} \Sigma_1 + \frac{61}{140} \Sigma_2$$

模型修正(Modifying Model)



相同的协方差矩阵

54% accuracy  73% accuracy

三部曲

- 函数集(模型):

x 

$$P(C_1|x) = \frac{P(x|C_1)P(C_1)}{P(x|C_1)P(C_1) + P(x|C_2)P(C_2)}$$

If $P(C_1|x) > 0.5$, output: class 1
否则, output: class 2

- 损失函数(Goodness of a function) :
 - 极大似然估计, μ , Σ 等
- 获得最优模型: 水到渠成

分布的选择

- 根据实际情况和你的经验，反复尝试，获得适合的分布

$$P(x|C_1) = P(x_1|C_1) P(x_2|C_1) \cdots P(x_k|C_1) \cdots$$

$$\begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_k \\ \vdots \\ x_K \end{bmatrix}$$

一维正态分布？

如果是离散数据，根据数据的具体情况，选择分布。

如果你的每个特征（变量）之间完全独立，那么模型就变成朴素贝叶斯分类(Naive Bayes Classifier)。

朴素贝叶斯分类(Naive Bayes Classifier)

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

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

```
>>> import numpy as np
>>> X = np.array([[-1, -1], [-2, -1], [-3, -2], [1, 1], [2, 1], [3, 2]])
>>> Y = np.array([1, 1, 1, 2, 2, 2])
>>> from sklearn.naive_bayes import GaussianNB
>>> clf = GaussianNB()
>>> clf.fit(X, Y)
GaussianNB()
>>> print(clf.predict([[-0.8, -1]]))
[1]
>>> clf_pf = GaussianNB()
>>> clf_pf.partial_fit(X, Y, np.unique(Y))
GaussianNB()
>>> print(clf_pf.predict([[-0.8, -1]]))
[1]
```

```
>>> import numpy as np
>>> rng = np.random.RandomState(1)
>>> X = rng.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()
>>> print(clf.predict(X[2:3]))
[3]
```

`sklearn.naive_bayes.ComplementNB(*, alpha=1.0, fit_prior=True, class_prior=None, norm=False)`

```
>>> import numpy as np
>>> rng = np.random.RandomState(1)
>>> X = rng.randint(5, size=(6, 100))
>>> y = np.array([1, 2, 3, 4, 5, 6])
>>> from sklearn.naive_bayes import ComplementNB
>>> clf = ComplementNB()
>>> clf.fit(X, y)
ComplementNB()
>>> print(clf.predict(X[2:3]))
[3]
```

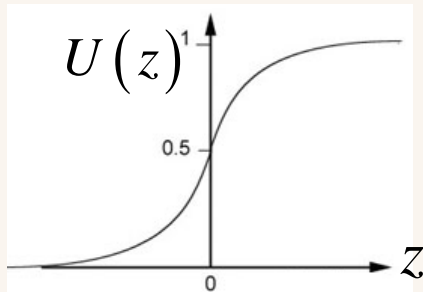
后验概率(Posterior Probability)

$$P(C_1|x) = \frac{P(x|C_1)P(C_1)}{P(x|C_1)P(C_1) + P(x|C_2)P(C_2)}$$

$$= \frac{1}{1 + \frac{P(x|C_2)P(C_2)}{P(x|C_1)P(C_1)}} = \frac{1}{1 + \exp(-z)} = U(z)$$

Sigmoid function

$$z = \ln \frac{P(x|C_1)P(C_1)}{P(x|C_2)P(C_2)}$$



后验概率(Posterior Probability)

$$P(C_1|x) = U(z) \quad \text{sigmoid} \quad z = \ln \frac{P(x|C_1)P(C_1)}{P(x|C_2)P(C_2)}$$

$$z = \ln \frac{P(x|C_1)}{P(x|C_2)} + \ln \frac{P(C_1)}{P(C_2)} \rightarrow \frac{\frac{N_1}{N_1 + N_2}}{\frac{N_2}{N_1 + N_2}} = \frac{N_1}{N_2}$$

$$P(x|C_1) = \frac{1}{(2\pi)^{D/2}} \frac{1}{|\Sigma_1|^{1/2}} \exp \left\{ -\frac{1}{2} (x - \mu_1)^T (\Sigma_1)^{-1} (x - \mu_1) \right\}$$

$$P(x|C_2) = \frac{1}{(2\pi)^{D/2}} \frac{1}{|\Sigma_2|^{1/2}} \exp \left\{ -\frac{1}{2} (x - \mu_2)^T (\Sigma_2)^{-1} (x - \mu_2) \right\}$$

$$z = \ln \frac{P(x|C_1)}{P(x|C_2)} + \ln \frac{P(C_1)}{P(C_2)} = \frac{N_1}{N_2}$$

$$P(x|C_1) = \frac{1}{(2\pi)^{D/2}} \frac{1}{|\Sigma_1|^{1/2}} \exp \left\{ -\frac{1}{2} (x - \mu_1)^T (\Sigma_1)^{-1} (x - \mu_1) \right\}$$

$$P(x|C_2) = \frac{1}{(2\pi)^{D/2}} \frac{1}{|\Sigma_2|^{1/2}} \exp \left\{ -\frac{1}{2} (x - \mu_2)^T (\Sigma_2)^{-1} (x - \mu_2) \right\}$$

$$\ln \frac{\cancel{\frac{1}{(2\pi)^{D/2}}} \frac{1}{|\Sigma_1|^{1/2}} \exp \left\{ -\frac{1}{2} (x - \mu_1)^T (\Sigma_1)^{-1} (x - \mu_1) \right\}}{\cancel{\frac{1}{(2\pi)^{D/2}}} \frac{1}{|\Sigma_2|^{1/2}} \exp \left\{ -\frac{1}{2} (x - \mu_2)^T (\Sigma_2)^{-1} (x - \mu_2) \right\}}$$

$$= \ln \frac{|\Sigma_2|^{1/2}}{|\Sigma_1|^{1/2}} \exp \left\{ -\frac{1}{2} [(x - \mu_1)^T (\Sigma_1)^{-1} (x - \mu_1) - (x - \mu_2)^T (\Sigma_2)^{-1} (x - \mu_2)] \right\}$$

$$= \ln \frac{|\Sigma_2|^{1/2}}{|\Sigma_1|^{1/2}} - \frac{1}{2} [(x - \mu_1)^T (\Sigma_1)^{-1} (x - \mu_1) - (x - \mu_2)^T (\Sigma_2)^{-1} (x - \mu_2)]$$

$$z = \ln \frac{P(x|C_1)}{P(x|C_2)} + \ln \frac{P(C_1)}{P(C_2)} = \frac{N_1}{N_2}$$

$$= \ln \frac{|\Sigma_2|^{1/2}}{|\Sigma_1|^{1/2}} - \frac{1}{2} [(x - \mu_1)^T (\Sigma_1)^{-1} (x - \mu_1) - (x - \mu_2)^T (\Sigma_2)^{-1} (x - \mu_2)]$$

$$(x - \mu_1)^T (\Sigma_1)^{-1} (x - \mu_1)$$

$$= x^T (\Sigma_1)^{-1} x - x^T (\Sigma_1)^{-1} \mu_1 - (\mu_1)^T (\Sigma_1)^{-1} x + (\mu_1)^T (\Sigma_1)^{-1} \mu_1$$

$$= x^T (\Sigma_1)^{-1} x - 2(\mu_1)^T (\Sigma_1)^{-1} x + (\mu_1)^T (\Sigma_1)^{-1} \mu_1$$

$$(x - \mu_2)^T (\Sigma_2)^{-1} (x - \mu_2)$$

$$= x^T (\Sigma_2)^{-1} x - 2(\mu_2)^T (\Sigma_2)^{-1} x + (\mu_2)^T (\Sigma_2)^{-1} \mu_2$$

$$z = \ln \frac{|\Sigma_2|^{1/2}}{|\Sigma_1|^{1/2}} - \frac{1}{2} x^T (\Sigma_1)^{-1} x + (\mu_1)^T (\Sigma_1)^{-1} x - \frac{1}{2} (\mu_1)^T (\Sigma_1)^{-1} \mu_1 \\ + \frac{1}{2} x^T (\Sigma_2)^{-1} x - (\mu_2)^T (\Sigma_2)^{-1} x + \frac{1}{2} (\mu_2)^T (\Sigma_2)^{-1} \mu_2 + \ln \frac{N_1}{N_2}$$

$$P(C_1|x) = U(z)$$

$$z = \ln \frac{|\Sigma_2|^{1/2}}{|\Sigma_1|^{1/2}} - \frac{1}{2} x^T (\Sigma_1)^{-1} x + (\mu_1)^T (\Sigma_1)^{-1} x - \frac{1}{2} (\mu_1)^T (\Sigma_1)^{-1} \mu_1 \\ + \frac{1}{2} x^T (\Sigma_2)^{-1} x - (\mu_2)^T (\Sigma_2)^{-1} x + \frac{1}{2} (\mu_2)^T (\Sigma_2)^{-1} \mu_2 + \ln \frac{N_1}{N_2}$$

$$\Sigma_1 = \Sigma_2 = \Sigma$$

$$z = \underbrace{(\mu_1 - \mu_2)^T \Sigma^{-1} x}_{\mathbf{w}^T} - \underbrace{\frac{1}{2} (\mu_2)^T \Sigma^{-1} \mu_1 + \frac{1}{2} (\mu_2)^T \Sigma^{-1} \mu_2}_{\mathbf{b}} + \ln \frac{N_1}{N_2}$$

$$P(C_1|x) = U(\mathbf{w} \cdot \mathbf{x} + b)$$

我们可以直接获得参数值—判别模型

在生成模型中，我们估计 $N_1, N_2, \mu_1, \mu_2, \Sigma$

然后就可以获得 \mathbf{w} 和 \mathbf{b}



生成模型

- 生成模型和判别模型
- 概率生成模型
- 性能对比
- 其他生成模型
- 补充知识

Discriminative v.s. Generative

$$P(C_1|x) = \sigma(w \cdot x + b)$$

直接估计 w 和 b

估计 μ_1, μ_2, Σ

$$w^T = (\mu_1 - \mu_2)^T \Sigma^{-1}$$

$$b = -\frac{1}{2}(\mu_1)^T(\Sigma_1)^{-1}\mu_1$$

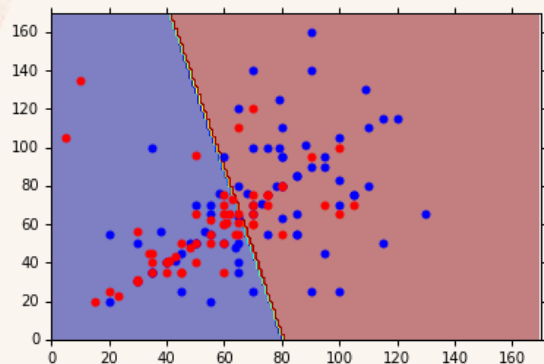
$$+\frac{1}{2}(\mu_2)^T(\Sigma_2)^{-1}\mu_2 + \ln \frac{N_1}{N_2}$$

同样的数据会获得相同的参数吗？

由于估计方式等技术问题，经常会出现同样的样本数据获得的最终模型不同。

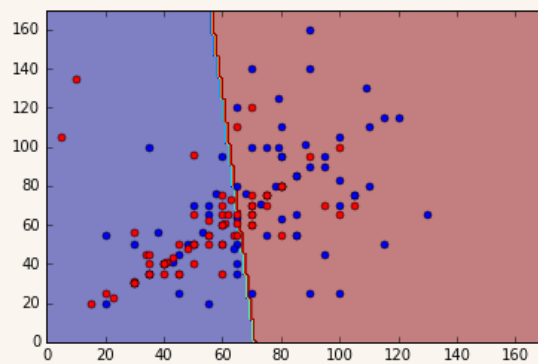
Generative v.s. Discriminative

Generative



73% accuracy

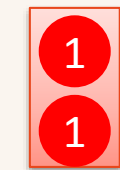
Discriminative



79% accuracy

Generative v.s. Discriminative

Training
Data



Class 1



Class 2

X
4



Class 2

X
4



Class 2

X
4

Testing
Data



Class 1?

Class 2?

朴素贝叶斯怎样?

$$P(x|C_i) = P(x_1|C_i)P(x_2|C_i)$$

Generative v.s. Discriminative

Training
Data



Class 1



Class 2

X
4



Class 2

X
4



Class 2

X
4

$$P(C_1) = \frac{1}{13}$$

$$P(x_1 = 1|C_1) = 1$$

$$P(x_2 = 1|C_1) = 1$$

$$P(C_2) = \frac{12}{13}$$

$$P(x_1 = 1|C_2) = \frac{1}{3}$$

$$P(x_2 = 1|C_2) = \frac{1}{3}$$

Training
Data



Class 1



Class 2

X
4



Class 2

X
4



Class 2

X
4

Testing
Data



$$P(C_1|x) = \frac{P(x|C_1)P(C_1)}{P(x|C_1)P(C_1) + P(x|C_2)P(C_2)}$$

<0.5

Diagram illustrating the calculation of the posterior probability $P(C_1|x)$ for the testing data (two red circles, both containing 1). The numerator represents the joint probability for Class 1, and the denominator represents the total joint probability for both classes.

Arrows indicate the components of the formula:

- 1×1 (top left) points to $P(x|C_1)$
- $\frac{1}{13}$ (top right) points to $P(C_1)$
- 1×1 (bottom left) points to $P(x|C_1)$
- $\frac{1}{13}$ (bottom middle) points to $P(C_1)$
- $\frac{1}{3} \times \frac{1}{3}$ (bottom right) points to $P(x|C_2)$
- $\frac{12}{13}$ (bottom far right) points to $P(C_2)$

$$P(C_1) = \frac{1}{13}$$

$$P(x_1 = 1|C_1) = 1$$

$$P(x_2 = 1|C_1) = 1$$

$$P(C_2) = \frac{12}{13}$$

$$P(x_1 = 1|C_2) = \frac{1}{3}$$

$$P(x_2 = 1|C_2) = \frac{1}{3}$$

Generative v.s. Discriminative

- 通常情况下，判别式模型的预测效果更好
- 生成模型的优势在于
 - 对于特定的分布来说
 - 对噪声更加稳健
 - 先验分布的获取更加灵活
 - 由于拟合过程更加精细，便于演化出其他的模型和算法，比如GAN（生成对抗网络）和其他深度生成模型。



机器学习基础（二）

- 生成模型和判别模型
- 概率生成模型
- 性能对比
- 其他生成模型
- 补充知识

拉普拉斯平滑

朴素贝叶斯模型可以在大部分情况下工作良好。但是该模型有一个缺点：对数据稀疏问题敏感。

比如在邮件分类中，假设NIPS这个词在词典中的位置为35000，然而NIPS这个词从来没有在训练数据中出现过，这是第一次出现NIPS，于是算概率时：

$$\begin{aligned}\phi_{35000|y=1} &= \frac{\sum_{i=1}^m 1\{x_{35000}^{(i)} = 1 \wedge y^{(i)} = 1\}}{\sum_{i=1}^m 1\{y^{(i)} = 1\}} = 0 \\ \phi_{35000|y=0} &= \frac{\sum_{i=1}^m 1\{x_{35000}^{(i)} = 1 \wedge y^{(i)} = 0\}}{\sum_{i=1}^m 1\{y^{(i)} = 0\}} = 0\end{aligned}$$

由于NIPS从未在垃圾邮件和正常邮件中出现过，所以结果只能是0了。于是最后的后验概率：

$$\begin{aligned}p(y=1|x) &= \frac{\prod_{i=1}^n p(x_i|y=1)p(y=1)}{\prod_{i=1}^n p(x_i|y=1)p(y=1) + \prod_{i=1}^n p(x_i|y=0)p(y=0)} \\ &= \frac{0}{0}.\end{aligned}$$

拉普拉斯平滑

对于这样的情况，我们可以采用拉普拉斯平滑，对于未出现的特征，我们赋予一个小的值而不是0。具体平滑方法为：

$$\phi_j = \frac{\sum_{i=1}^m 1\{z^{(i)} = j\}}{m}$$

变为

$$\phi_j = \frac{\sum_{i=1}^m 1\{z^{(i)} = j\} + 1}{m + k}$$

$$\begin{aligned}\phi_{j|y=1} &= \frac{\sum_{i=1}^m 1\{x_j^{(i)} = 1 \wedge y^{(i)} = 1\} + 1}{\sum_{i=1}^m 1\{y^{(i)} = 1\} + 2} \\ \phi_{j|y=0} &= \frac{\sum_{i=1}^m 1\{x_j^{(i)} = 1 \wedge y^{(i)} = 0\} + 1}{\sum_{i=1}^m 1\{y^{(i)} = 0\} + 2}\end{aligned}$$

A decorative red watercolor splash is located in the top-left corner of the slide. Within this splash, three small red birds are depicted in flight, moving towards the right.

END
(Q&A)