

贝叶斯决策理论

● 贝叶斯公式：

P
(

ω

i

|
x
¯
)
=

p
(
x
¯
|

ω

i

)
P
(

ω

i

)

∑
p
(
x
¯
|

ω

j

)
P
(

ω

j

)

● 联合概率分布和条件概率分布：

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

● 最小错误率方法：

$$\arg \max_i P(\omega_i|\vec{x})=\arg \max_i P(\vec{x}|\omega_i)P(\omega_i)$$

● 最小错误率分析：

P
(
e
)
=

P
(

ω

2

)

P

2

(
e
)
+
P
(

ω

1

)

P

1

(
e
)
,

P

2

(
e
)

是把2分类错误成1的概率

● 多类正确率分析：

P
(
c
)
=
∑

P

i

(
c
)
P
(

ω

j

)

 对应最高的曲线加先验的面积

● False positive: 将负样本判成正样本。False negative: 将正样本判成负样本

● ROC曲线：横轴 False Postive，纵轴 True Positive

● 最小风险贝叶斯决策：

λ
(

α

i

,

ω

j

)

是在*ω_j*状态上采用*α_i*决策的代价。（选对了可以认为是0）。决策arg min ∑_j λ(α_i, ω_j) p(x̄|ω_j) p(ω_j)

● 限定一类错误率，求第二类错误率最小值：拉格朗日待定系数法 γ = P₁(e) + λ(P₂(e) − ε₀) = 1 − P₁(c) + λP₂(e) − λε₀。对λ求偏导，一维情况得到决策面满足λ =

p
(
t
|

ω

1

)

p
(
t
|

ω

2

)

，决策面限定错误率达到指标上限。（图解：看残余面积）。

● 最小最大决策： 风险已知， 先验概率未知。

R
=

∫

R

1

λ

11

p
(

ω

1

)
P
(
x
¯
|

ω

1

)
+

λ

12

p
(

ω

2

)
P
(
x
¯
|

ω

2

)
+

∫

R

2

λ

21

p
(

ω

1

)
P
(
x
¯
|

ω

1

)
+

λ

22

p
(

ω

2

)
P
(
x
¯
|

ω

2

)
=
a
+
p
(

ω

1

)
b

风险不可确定，以b = 0作为最坏情况的最优解（防止P(ω₂)产生变化

● 多类决策：可以进行逐一比较。或者算最大的P(ω_i|x̄) > P(ω_j|x̄) ⇒ p(x̄|ω_i)p(ω_i) > p(x̄|ω_j)p(ω_j) ⇒ p(x̄|ω_i)p(x̄|ω_j) > p(ω_j)p(ω_i) ⇒ ln p(x̄|ω_i) + ln p(ω_i) > ln p(x̄|ω_j) + ln p(ω_j)（对数似然）

● 决策面以及两类情况的决策面：

g
(
x
¯
)
{
>
0
⇒
ω

1

<
0
⇒
ω

2

}

g
(
x
¯
)
=
0

是决策面

朴素贝叶斯决策和正态分布决策

● 朴素贝叶斯：认为所有变量独立：P(x̄|ω) = ∏_i p(x_i|ω)

● 正态分布基本性质：

P
(
x
¯
)
=

1

(
2
π

)

d

/

2

|
Σ

|

1

/

2

exp
⁡
(
−

1

2

(
x
¯
−
μ
)

T

Σ

−
1

(
x
¯
−
μ
)

)

● 性质：等密度点是超椭圆面(x̄−μ)^TΣ^{−1}(x̄−μ) = const, 马氏距离，不相关性E₁₂ = E₁E₂ == 独立性P₁₂ = P₁P₂（仅在正态分布下）

● 边缘分布和条件分布，线性变化正态性。g̃ = A x̄ → N(Aμ, AΣA^{A^T})

● Σ_i = σ²I（协方差相等均对角），分类面为圆形（先验不相等）或中位（先验相等）

● Σ_i = Σ（协方差相等），分类面为椭圆形（先验不相等）或直线（先验相等，相对中位线发生偏转）

● 其他可能性：二次曲线。计算错误率：公式，近似计算上界，实验估计

概率密度函数的估计参数和非参数估计

● 参数空间，点估计，区间估计，统计量等概念。MLE方法max ∑_i ln(p(x̄_i|θ))（有的时候不需要求导，

● 正态分布的MLE估计方法：

μ
=

1

N

∑

i

x
¯

i

,
Σ
=

1

N

∑
(
x
¯
−
μ
)
(
x
¯
−
μ

)

T

● 可识别性：参数空间内每两个不同参数均有不同的概率密度函数，连续的往往是可以识别的，离散的往往不能识别

● 贝叶斯参数估计：将参数认为成随机变量， 并有先验。P(θ|x̄) =

P
(
x
¯
|
θ
)
P
(
θ
)

p
(
x
¯
)

● 正态分布MAP：

x
¯
∼
N
(
μ
,

σ

2

)
,
p
(
μ
)
∼
N
(

μ

0

,

σ

0

2

)

(
均
值
有
不
确
定
性
)

P
(
μ
|
x
¯
)
=

1

p
(
x
¯
)
∏

1

√
2
π
σ

exp
⁡
(
−

1

2

x
¯

k

−

μ

2

σ

2

)

1

√
2
π

σ

0

exp
⁡
(
−

1

2

m
¯
u
−

μ

0

σ

0

2

)

∼
N
(

μ

N

,

σ

N

)
,

μ

N

=

N

σ

0

2

+

σ

2

m
N

+

σ

2

N

σ

0

2

+

σ

2

μ

0

.

结果仍然是正态分布。N = 0 →仅能靠先验，没有实验信息。N → ∞实验信息足够多，先验信息没有用。σ₀ → 0先验太强，实验信息被忽视。σ₀ → ∞先验信息太弱，被忽视。当先验可靠时，利用更多的信息。目的是max P(x̄|θ)p(θ)。同等情况下MLE更简单

● 最小二乘等价于对认为误差是正态分布的MLE方法，此方法上的加参数先验为正态分布的MAP方法就是L-2正则化

● Parzen窗法：窗口大小的选择，正态分布Parzen窗

● 三类误差：贝叶斯误差：特征一旦选定就无法改变的误差。模型误差：模型不标准或者错误。估计误差：参数估计的时候产生的误差。维数问题：普通x100个样本一维。依靠PCA，独立性等方式降低维数灾难

● 过拟合： MLE的普遍特性， 考虑增加样本或者引入MAP方案，通过简单模型（参数化模型， 对角矩阵， 共享参数）降低参数

● 错误率估计：先验未知ε =

k
N

,
E
(
k
)
=
N
ϵ
,
V
a
r
(
k
)
=
N
ϵ
(
1
−
ϵ
)
,
V
a
r
(
ê
)
=

ϵ
(
1
−
ϵ
)

N

.

先验已知：

E
(
ê
)
=
ϵ
,
V
a
r
(
ê
)
=

1

N

∑
P
(

ω

i

)

ϵ

i

(
1
−

ϵ

i

)
.

已知先验方差更低。交叉验证和留一法等

EM, GMM 算法

● GMM算法：

P
(
X
|
Θ
)
=

∑

α

i

p

i

(
X
|

θ

i

)
,
∑

α

i

=
1
,

p

i

(
X
|

θ

i

)
∼
N
(

μ

i

,

σ

i

)
log
⁡
M
L
E

：

∑

i

ln
∑

j

N
(

x

i

|

μ

j

,

Σ

j

)
P
(

ω

j

)

● 已知概率是多类高斯混合而成的，每一类的比例（先验）未知，每一类高斯的参数未知

$$P(\omega_k|x_i,\mu_k,\Sigma_k)=\frac{N(x_i|\mu_k,\Sigma_k)P(\omega_k)}{\sum_jN(x_i|\mu_j,\Sigma_j)P(\omega_j)}$$

$$\hat{P}(\omega_k|x_i,\mu_k,\Sigma_k)=\frac{1}{N}\sum_iP(\omega_k|x_i,\mu_k,\Sigma_k)$$

$$\text{对}\mu_k\text{求偏导为零得到}\hat{\mu}_k=\frac{\sum_i\frac{P(\omega_k|x_i,\mu_k,\Sigma_k)x_i}{\sum_iP(\omega_k|x_i,\mu_k,\Sigma_k)}}{\sum_i\frac{P(\omega_k|x_i,\mu_k,\Sigma_k)}{\sum_iP(\omega_k|x_i,\mu_k,\Sigma_k)}}$$

$$\text{对}\Sigma_k:\hat{\Sigma}_k=\frac{\sum_i\frac{P(\omega_k|x_i,\mu_k,\Sigma_k)(x_i-\mu_k)(x_i-\mu_k)^T}{\sum_iP(\omega_k|x_i,\mu_k,\Sigma_k)}}{\sum_i\frac{P(\omega_k|x_i,\mu_k,\Sigma_k)}{\sum_iP(\omega_k|x_i,\mu_k,\Sigma_k)}}$$

● 解决策略：E步利用上一步的东西计算P(ω_k|x_i, μ_k, Σ_k)，M步计算参数迭代

● EM 算法：解决数据缺失或隐变量的问题X显变量，Y隐变量。P(X,Y|θ) = p(Y|X, θ)p(X|θ)

● L(θ) = ln p(X|θ) = ln(∑_Y p(X,Y|θ))，设Y ~ q(Y)

线性判别函数

● 线性判别，广义线性判别，增广向量等

● 感知准则函数：假设样本线性可分，

{

a

T

y

i

>
0
对
一
切

y

i

∈

ω

1

a

T

y

i

<
0
对
一
切

y

i

∈

ω

2

}

采用y_i' =

{

−

y

i

得到
∀

i

,

a

T

y

i

′

>
0
,
错
分
样
本
集
合
为

Y

k

}

● 目标函数J_P(a) = ∑_{y∈Y^k} −a^Ty ≥ 0，希望优化到0。梯度下降得到ak + 1 = a_k + ρ ∑_{y∈Y^k} y

● 1-bit SGD, ρ = 1 可以达到最小值（可以硬证明）

SVM

● 优化问题：min

1

2

w

T

w
,
subject
to

d

i

(

w

T

x

i

+
b
)
>
1
,
凸
优
化
问
题
加
拉
格
朗
日
乘
子
转
对
偶
方
法
优
化

● J =

1

2

w

T

w
−
∑

α

i

[

d

i

(

w

T

x

i

+
b
)
−
1
]
,

α

i

>
0

 转优化J(W̄, ᾱ).

● 鞍点是最值点：J(W̄', α) ≤ J(W̄', α') ≤ J(W̄, α')，左侧得到∑(α'_i − α)g_i(W̄') ≤ 0 ⇒ ∑g_i(w̄') > 0说明是可行点，取α_i = 0 得到α'_ig_i(w̄') = 0

● 右侧：f(w̄') ≤ f(w̄) − α'g(w̄) ≤ f(w̄)得到最优点（强对偶条件）

● 对偶求鞍点：w = ∑ α_id_ix_i, ∑ alpha_id_i = 0

● 对偶方法：

max

Q
(
α
)
=
∑

α

i

−

1

2

∑
∑

α

i

α

j

d

i

d

j

H

i
j

,

H

i
j

=

x

i

T

x

j

,

之后可以更换为核函数。

● 线性不可分松弛：min

1

2

w

T

w
+
C
∑
ξ
.

松弛d_i(w^Tx_i + b) > 1 − ξ_i，C越大，要求性能越好，间隔越小，C越小要求间隔越大。

● 非线性核函数方法，多项式K(x, y) = (xy + 1)^p，RBF：exp(−||x − y||²/2σ²)等、

● 采用L-1正则化和L-2正则化的函数的SVM， L-1稀疏学习，L-2防过拟合

● 产生式模型（优点：可以提供大量信息量）和判别式模型（优点：简单）

神经网络

● 损失函数：TSSE：

1

2

(

d

i

−

y

i

)

2

,

可以进一步补充为对整个batch的求导

● BP算法，矩阵求导，反向传播，链式法则等

● CNN, RNN, LSTM的结构（有空补上）。序列长度：一个sample的长度（sample之间没有关系）

决策树

● 交叉熵：E = − ∑ p_i ln p_i

● Gini 不纯度：1 − ∑ P(w_n)²

● 误差不纯度：1 − max_j P(ω_j)

● 信息熵的增量：l − p(l)l_l − p(r)l_r（左右两支加权计算）

近邻法和距离

● 最近邻法错误率在12倍贝叶斯错误率中甸

● 压缩近邻法：首先用近邻法测试所有测试样本，如果测错了放进评测集，通过这种方式扩充评测集

● 距离度量：1. s-Minkowski: D(x, y) = [∑ |x_i − y_i|^s]^{1/s}. s = 2为欧几里得距离。Chebyshev 距离：max_i |x_i − y_i|，马氏距离(x − y)^TQ(x − y)

● 距离的正定性，三角形，Holder不等式等（不会）

● KL 散度作为概率PDF之间距离，切距离；做流形变换之后点到直线的距离（单边切距离，双边切距离）

特征提取和特征选择PCA, KL变换等线性方法

● 目的：防止维数灾难，方便可视化和理解

● Fisher 方法： 找一个方向使得投影结束后类间方差最大，类间方差最小。类间方差(m₁ − m₂)(m₁ − m₂)^T. 类内方差∑_j ∑_i(x_{ij} − m_j)(x_{ij} − m_j)^T

● 投影之后方差变化做比J(w) =

w

T

S

b

w

w

T

S

w

w

最大转拉格朗日乘子：L = w^TS_bw − λ(w^TS_ww − c)（控制一个类内方差不变）得到w是S_w^{−1}S_b的最大特征值对应的特征向量。

● Fisher投影之后的分类问题：按照两类重心的重点、或者按照加权重心或者

m
1
+

m

2

2

+

ln
⁡
P
(

ω

1

)
−
ln
⁡
P
(

ω

2

)

N

1

+

N

2

−
2

先验。也可以采用Bayes决策手段

● 其他的Fisher手段，局部方法，非线性方法等，选取其他判别方案等（如矩阵迹）

● 多维度Fisher方法：选取最大的几个向量

● PCA：最小化剩余方差min E(x − x̂)^T(x − x̂) = E[∑_{i= d+1} c²]

● PCA：找S（对称正交矩阵）的从大到小的特征值，对应的特征向量组成变换矩阵。

特征选择方案

● 向前向后方法，顺序前进法，顺序后退法（找收益最大的和损失最小的）

● Relief方案，遗传算法等，存在问题：特征选择过学习

MDS, LLE, ISOMap, 遗传算法

● MDS方案：给定两点之间的距离，找p维的x使得d²_{r_s} = (x_r − x_s)^T(x_r − x_s) = x_r^Tx_r + x_s^Tx_s − 2x_r^Tx_s可以构造B = XX^T, n × n, 进行p维PCA提出 B = VΛV^T, X = V₁Λ^{1/2}, 展开成P维空间

● ISOMAP 方案：先通过kNN或者ϵ圆构建一个整个图的支撑树，之后采用Dijkstra算法找到图中任意两点的最近距离构建距离矩阵，之后使用MDS方法确定坐标

● LLE（局部线性化）：先通过kNN或者ϵ圆找到每一个点的近邻然然后把每一个面元铺平

● GA：初始化种群，交叉和编码方式等

Ensemble

● Adaboost：采用一组弱学习器进行学习。方案：1. 初始化w₁(i) =

1

N

.

 2. 归一化：p_l(i) =

w

l

(
i
)

∑

i

w

l

(
i
)

 3. ϵ_l = ∑_{m is} p_l(i) 4. a_l =

1

2

ln
⁡

1
−

ϵ

l

ϵ

l

}

 5. 调整：w_{l+1}(i) = w_l(i) exp(±a_l) 正：错分样本加权重。负：正确样本减权重。5. 训练结束之后按a_l对L个分类器加权投票。

● 随机森林：随机抽取若干特征构造树，按照准确率进行投票，投票悖论等

聚类分析，非监督学习方法

● GMM聚类：对样本进行GMM算法，每一个成分分成一类

● kmeans方案：对分割进行调整规则：

N

j

N

j
+
1

||
y
−

m

j

||

2

<

N

k

N

k
−
1

||
y
−

m

k

||

2

将y从k类移到j类，复杂度O(n)

● 分类个数的选取：选择肘点（二阶差分最大（一阶差分变化最大））

● 核函数方法，选择核函数代替欧几里得距离

● 多级聚类方案：将两个最近/最远/平均距离最近的两类每次聚类，直到只剩下一类。形成聚类树。O(n³)复杂度

● 谱聚类：首先构造相似距离矩阵D（越远的距离越大，可以使用RBF核或者其他ϵ圆或者kNN方案。将行相加构造对角矩阵D。L = D − W, L_{r_{rw}} = D^{−1}L构造非归一化和归一化的拉普拉斯矩阵。计算两个矩阵的前k个向量组成n行k列的矩阵，对此矩阵进行聚类。复杂度O(n²)

各种情况下的EM算法

一般情况下的EM算法

- X 是显变量, Z 是隐变量 (一般来说是离散的)。联合分布 $p(X, Z|\theta)$ 确定。目标是最大化 $p(X|\theta) = \sum_Z p(X, Z|\theta)$
- 强行引入 Z 的分布 $q(Z)$, 立即得到 $\ln(p(X|\theta)) = \mathcal{L}(q, \theta) + KL(q\|p)$
- $\mathcal{L}(q, \theta) = \sum_Z q(Z) \ln \frac{p(X, Z|\theta)}{q(Z)}$
- $KL(q\|p) = -\sum_Z q(Z) \ln \frac{p(Z|X, \theta)}{q(Z)}$
- 以上根据 $\ln P(X, Z|\theta) = \ln P(Z|X, \theta) + \ln P(X|\theta)$, $KL(q\|p)$ 是KL散度, 当且仅当 $q(Z) = p(Z|X, \theta)$ (估计完全正确) 时等于0否则大于零, 注意 $\sum_Z q(z) \ln p(X|\theta) = \ln(p(X|\theta))$ ($\ln p(X|\theta)$ 与 Z 无关
- $\mathcal{L}(q, \theta) \leq \ln(p(X|\theta))$, 得到优化下界
- E-step: 使用 q 最大化 $\mathcal{L}(q, \theta)$, 由于 $\ln(p(X|\theta))$ 不变, 故只能KL散度为0, $q(Z) = p(Z|X, \theta)$ 时 \mathcal{L} 达到最大化。
- M-step: 使用 θ 优化 \mathcal{L} , 这将使得 $p(Z|X, \theta)$ 导致非零的KL散度, 给E-step留空间。

- 当总共有 n 个样本时, $P(Z|X, \theta) = \frac{P(X, Z|\theta)}{\sum_Z P(X, Z|\theta)} = \frac{\prod p(x_n, z_n|\theta)}{\sum_Z \prod p(x_n, z_n|\theta)} = \prod p(z_n|x_n, \theta)$ 。其中用到了交换求和求积顺序, 同时需要条件概率。

基于上述解释的EM算法

- E-step: 计算 $p(Z|X, \theta)$
- M-step: 计算 $\theta = \arg \max_{\theta} Q(\theta, \theta^{old}) = \sum_Z p(Z|X, \theta^{old}) \ln p(X, Z|\theta)$ 。 这里忽略了上面的分母因为和 θ 无关

EM算法用于有先验的贝叶斯参数估计, GEM算法

- $\ln P(\theta|X) = \ln P(X|\theta) + \ln P(\theta) - \ln P(X)$ 进行优化
- GEM 算法: 对M-step进行梯度下降而不是一次到最优

GMM算法的EM表示

- E-step: $\gamma(z_{nk}) = \frac{\pi_k \mathcal{N}(x_n|\mu_k, \Sigma_k)}{\sum_j \pi_j \mathcal{N}(x_n|\mu_j, \Sigma_j)}$
- M-step: $\begin{cases} \mu_k = \frac{1}{N_k} \sum \gamma(z_{nk}) x_n \\ \Sigma_k = \frac{1}{N_k} \sum \gamma(z_{nk}) (x_n - \mu_k)(x_n - \mu_k)^T \end{cases}$
- 对数似然: $\sum \ln \sum \pi_k \mathcal{N}(x_n|\mu_k, \Sigma_k)$
- 隐变量是二值的 Z, Z_{nk} 表征第 n 个样本是否在 k 的聚类中。 $P(X, Z|\theta) = \prod_n \prod_k \pi_k^{z_{nk}} \mathcal{N}(x_n|\theta)^{z_{nk}}$
- $p(Z|X, \theta) = const \prod_n \prod_k \pi_k^{z_{nk}} \mathcal{N}(x_n|\theta)^{z_{nk}}$ 具有 X 在第 k 个聚类中的概率, 可以用条件公式贝叶斯等
- M-step 必然得到GMM算法的结果

使用琴生不等式推出的EM算法

- $\ln P(x|\theta) = \ln \sum_Z P(x, z|\theta) = \ln \sum_Z q(z) \frac{p(x, z|\theta)}{q(z)} \geq \sum_Z q(z) \ln p(x, z|\theta) - \sum_Z q(z) \ln q(z) = F(q, \theta)$
- E-step 仍然是对 $F(q, \theta)$ 进行优化, 得到结论仍是 $q(Z) = p(Z|X, \theta)$, 原因是此时处理之后 $F(q, \theta) = \mathcal{L}(\theta)$ 。我们优化的是 q
- M-step 优化 $Q(\theta, \theta^{old}) = \sum_Z p(Z|X, \theta^{old}) \ln p(x, z|\theta)$ 和前面的相同

GMM对应一个batch的样本

- $P(Z|X, \theta) = \prod p(z_n|x_n, \theta)$, 混合模型下 $p(x, z|\theta) = \sum \alpha p_{z_i} p_{z_i}(x|\theta)$ 。注意这里的表示方式和上面的二值表示不同, 这里表示的是下标参数, 因此 \sum_Z 需要展开成那个很复杂的式子。设 $y_i = k$ 表示的是 i 样本在 k 类中, 直接得到M步为 $\sum_{y_1=1}^M \cdots \sum_{y_1=1}^M \sum_{i=1}^N \ln(\alpha_{z_i} p_{z_i}(x|\theta_{z_i}))$ 。
- 这样做没有上面那么做好, 二值的可以表示的更简洁。求解过程中注意区分 $\theta, \theta^{old}, p(z|x, \theta_{old})$ 中不含我们需要的东西