

$$y = f + \varepsilon, \quad \varepsilon \text{ noise, zero mean, variance } \sigma^2$$

$$\begin{aligned} \text{均方根差: } \text{MSE} &= E(y - \hat{f}(x))^2 \\ \text{方差: } \text{Var}(\hat{f}(x)) &= E[\hat{f}(x)^2] - E[\hat{f}(x)]^2 \\ \text{偏差: } \text{Bias}(\hat{f}(x)) &= E(\hat{f}(x)) - f(x) \end{aligned}$$

$$\begin{aligned} \text{Var}(x) &= E(x^2) - (E(x))^2 \\ \text{Var}(\bar{x}) &= \frac{1}{m} \left(\sum_{i=1}^m (x_i - \bar{x})^2 \right) = \frac{1}{m} \sum_{i=1}^m (x_i^2 + \bar{x}^2 - 2x_i \bar{x}) \\ &= \frac{\sum_{i=1}^m x_i^2}{m} + \frac{\sum_{i=1}^m \bar{x}^2}{m} - \frac{\sum_{i=1}^m 2x_i \bar{x}}{m} \\ &= E(x^2) + \bar{x}^2 - 2\bar{x} E(x) \\ &= E(x^2) + (E(x))^2 - 2(E(x))^2 \\ &= E(x^2) - (E(x))^2 \end{aligned}$$

$$\begin{aligned} E[(y - \hat{f})^2] &= \text{MSE} \\ &= E[(f + \varepsilon - \hat{f})^2] \\ &= E[f^2 + \varepsilon^2 - 2f\hat{f} + E(\hat{f})^2 - E(f)\hat{f}] \\ &= E[(\hat{f} - E(\hat{f}))^2] + E[(E(\hat{f}) - \hat{f})^2] \\ &\quad + E(\varepsilon^2) + 2E[(f - E(\hat{f})) \cdot (E(\hat{f}) - \hat{f})] \\ &\quad + 2E[(f - E(\hat{f}))\varepsilon] + 2E[(E(\hat{f}) - \hat{f}) \cdot \varepsilon] \\ &= (\hat{f} - E(\hat{f}))^2 + E(\varepsilon^2) + E[(E(\hat{f}) - \hat{f})^2] \\ &= (\hat{f} - E(\hat{f}))^2 + \text{Var}(\varepsilon) + \text{Var}(\hat{f}) \\ &= \text{Bias}(\hat{f})^2 + \sigma^2 + \text{Var}(\hat{f}) \end{aligned}$$



梯度下降：

$$\bar{J}(w, b) = \frac{1}{m} \sum_{i=1}^m J(y^i, \hat{y}^i)$$

find w and b to minimize $\bar{J}(w, b)$

$$w := w - \alpha \frac{d \bar{J}(w, b)}{d w}$$

$$b := b - \alpha \frac{d \bar{J}(w, b)}{d b}$$

2.5 导数 $e^x dx = e^x$

2.6 计算图

Final output
d Var.

$$\begin{aligned} a &= 5 \\ b &= 3 \end{aligned} \rightarrow v = bc \rightarrow \bar{J} = 37$$

$$y = x^{-1}$$



2.9 logistic 回归中的梯度下降
从最后向前一顶依次求导，得训练对
输入的导数.

$$\frac{\partial \bar{J}}{\partial w} = \frac{(1+e^{-z})^2 \cdot 2e^{-z}}{(1+e^{-z})^2} \cdot \frac{1}{1+e^z}$$

$$\begin{array}{c} x_1 \\ w \\ x_2 \\ w \\ \vdots \\ b \end{array} \Rightarrow z = w_0 x_0 + w_1 x_1 + \dots + b \Rightarrow a = g(z) \Rightarrow L(a, y)$$

$$da = \frac{d L(a, y)}{d a} = \frac{d (y \log a + (1-y) \log (1-a))}{d a}$$

反函数的导数
等于直接函数导数
的倒数.

$$= -\frac{y}{a} + \frac{1-y}{1-a}$$

$$\frac{d w_i^{(0)} d L(a, y)}{d w_i} = \frac{d L(a, y)}{d a} \cdot \frac{d a}{d z} \cdot \frac{d z}{d w_i}$$

$$= \left(-\frac{y}{a} + \frac{1-y}{1-a} \right) \cdot a (1-a) \cdot x_i$$

2.11 向量化 vectorization.

$$X = \begin{bmatrix} x^{(1)} \\ x^{(2)} \\ \vdots \\ x^{(m)} \end{bmatrix} \quad Y = \begin{bmatrix} y^{(1)} \\ y^{(2)} \\ \vdots \\ y^{(m)} \end{bmatrix}$$

$$Z = \begin{bmatrix} z^{(1)} \\ \vdots \\ z^{(m)} \end{bmatrix} = Z^{(m)} = W^\top X + [b \ b \ \dots \ b]$$

$$Z = \text{np.dot}(W^\top, X) + b$$

$$A = \begin{bmatrix} a^{(1)} \\ \vdots \\ a^{(m)} \end{bmatrix} \quad Y = \begin{bmatrix} y^{(1)} \\ \vdots \\ y^{(m)} \end{bmatrix}$$

$$dZ = A - Y$$

2.15 Broadcasting

use `np.random.randn(1, 5)`
but not `np.random.randn(5)`

Can use reshape to set the dim of the matrix.

3.6 線性函數

sigmoid function:

$$\alpha = \frac{1}{1 + e^{-x}}$$

$$\alpha = \tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}} \quad (\text{better})$$

ReLU: $\alpha = \max(0, z)$

2.18 logistic cost function



$$\text{if } y=1 : \quad P(y|x) = \hat{y}$$

$$\text{if } y=0 : \quad P(y|x) = 1 - \hat{y}$$

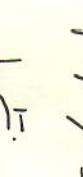
$$P(y|x) = \hat{y}^y (1 - \hat{y})^{1-y}$$

$$\max(-\log(P(y|x))) = \log(1 - \hat{y})$$

$$= y \log \hat{y} + (1-y) \log(1-\hat{y})$$

$$= -L(y, \hat{y})$$

$$\min(-\log(P(y|x))) = -\underbrace{(y \log \hat{y} + (1-y) \log(1-\hat{y}))}_{= L(y, \hat{y})}$$



3.10



Vazyme

$$\frac{dL}{d\bar{z}^{[2]}} = \frac{dL}{d\bar{a}^{[2]}} \cdot \frac{d\bar{a}^{[2]}}{d\bar{z}^{[2]}} = \bar{a}^{[2]} - y$$

$$\frac{dL}{d\bar{w}^{[2]}} = \frac{dL}{d\bar{z}^{[2]}} \cdot \frac{d\bar{z}^{[2]}}{d\bar{w}^{[2]}} = (\bar{a}^{[2]} - y) \cdot a^{[1]T}$$

$$\frac{dL}{db} = \frac{dL}{d\bar{z}^{[2]}} \cdot 1 = d\bar{z}^{[2]}$$

$$\begin{aligned} \frac{dL}{d\bar{z}^{[2]}} &= \frac{dL}{d\bar{z}^{[2]}} \cdot \frac{d\bar{z}^{[2]}}{d\bar{z}^{[2]}} \cdot \frac{d\bar{a}^{[2]}}{d\bar{a}^{[2]}} \cdot \frac{d\bar{a}^{[2]}}{d\bar{z}^{[2]}} \\ &= \frac{dL}{d\bar{z}^{[2]}} \cdot W^{[2]} \\ &= W^{[2]T} d\bar{z}^{[2]}, \quad g^{[2]}(z^{[2]}) \\ &= W^{[2]T} d\bar{z}^{[2]} \cdot X^T \\ \frac{dL}{db} &= d\bar{z}^{[2]} X^T \end{aligned}$$

2019.09.24.



Vazyme

为什么无无偏估计的方差是 $\frac{1}{n-1}$?

$$\bar{b}^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2$$

$$E(XY) = EX \cdot EY$$

$$Var(\bar{X}) = \frac{1}{n} Var(X) ?$$

$$\bar{b}(\bar{X})^2 = \frac{1}{n} \bar{b}^2(X)$$

二阶中心距:

$$\begin{aligned} E(X^2) &= Var(X) + (E(X))^2 \\ E(X^2) &= Var(X) + (E(X^2) - (E(X))^2) \\ E(X^2) &= Var(X) + (E(X^2) - (E(X))^2) = \bar{b}^2 + \bar{m}^2 \end{aligned}$$

无偏估计: 对变量 θ 的估计是 $\hat{\theta}$, 如果 $E(\hat{\theta}) = E(\theta)$, 则称 $\hat{\theta}$ 为 θ 的无偏估计.

$$E\left(\sum_{i=1}^n (x_i - \bar{x})^2\right)$$

$$= \sum_{i=1}^n E(x_i^2) - n \cdot E(\bar{x}^2)$$

$$= \sum_{i=1}^n (\bar{b}^2 + \bar{m}^2) - n (\frac{1}{n} \bar{b}^2 + \bar{m}^2)$$

$$= (n-1) \bar{b}^2$$

单样本

2019/10/02 BP (Back Propagation) Vazyme

输出层 第 j 个输出神经元的输出

$$\beta_j = \sum_{h=1}^d w_{jh} b_h$$

第 h 个隐藏神经元的输入
 $\alpha_h = \sum_{i=1}^d w_{hi} x_i$

输入层 (x_k, y_k)

输出: $\gamma_k = (\gamma_{1k}, \gamma_{2k}, \dots, \gamma_{dk})$

$$\gamma_j^k = f(\beta_j - \theta_j) \text{ (sigmoid)}$$

$$\beta_j = \sum_{h=1}^d w_{jh} b_h$$

$$\alpha_h = \sum_{i=1}^d w_{hi} x_i$$

网络在 (x_k, y_k) 的均方误差:

$$E_k = \frac{1}{2} \sum_{j=1}^L (\gamma_j^k - y_j^k)^2$$

【从这网络可以表示映射关系的前提下，计算均方误差来优化参数】

需确定的参数有:

$$\alpha \times q + q \times l + q + \text{输出层阈值}$$

输入
输出

$$\frac{\partial E_k}{\partial w_{kj}} = \frac{\partial E_k}{\partial \gamma_j^k} \cdot \frac{\partial \gamma_j^k}{\partial \beta_j} \cdot \frac{\partial \beta_j}{\partial w_{kj}}$$

$$\frac{\partial \beta_j}{\partial w_{kj}} = b_k$$

$$\frac{\partial E_k}{\partial \gamma_j^k} \cdot \frac{\partial \gamma_j^k}{\partial \beta_j} \quad \text{由于 } \gamma_j^k = f(\beta_j - \theta_j) \text{ 是 sigmoid 函数}$$

所以 $\frac{\partial \gamma_j^k}{\partial \beta_j} = (\gamma_j^k)' = f'(\beta_j) \cdot (1 - \gamma_j^k)$

$$\beta_j = -\frac{\partial E_k}{\partial \gamma_j^k} \cdot \frac{\partial \gamma_j^k}{\partial \beta_j} = -(\gamma_j^k - y_j^k) \cdot \gamma_j^k \cdot (1 - \gamma_j^k)$$

$$(\text{学习率}) \Delta w_{kj} = \eta \gamma_j^k b_k$$

$$\Delta \theta_j = -\eta \gamma_j^k (\text{目标}) \Delta r_k = -\eta \epsilon_k (\text{阈值})$$

$$(\text{修正}) \Delta w_{ih} = -\frac{\partial E_k}{\partial b_h} \cdot \frac{\partial b_h}{\partial \alpha_h} \cdot \frac{\partial \alpha_h}{\partial w_{ih}} \cdot \eta$$

$$= -\sum_{j=1}^L \frac{\partial E_k}{\partial \beta_j} \cdot \frac{\partial \beta_j}{\partial b_h} \cdot f'(\alpha_h - r_h) \cdot \eta$$

$$= -\sum_{j=1}^L \eta_j w_{kj} \cdot b_k (1 - b_k) \sum_{h=1}^d q_j w_{jh}$$

误差反传播对
权重求导
偏导数

$$\Delta w_{ij} = -\eta \frac{\partial E_k}{\partial w_{ij}}$$

使梯度下降: $\eta \leftarrow \eta + \Delta \eta$
更新参数: $\Delta w_{ij} = -\eta \frac{\partial E_k}{\partial w_{ij}}$
(单个训练，样本更新)

累积、BP

将单个 E_k 误差相加, 则对 w_k 未偏导也是各项的偏导相加.

$$w_k \leftarrow w_k + \frac{1}{n} \Delta w_k$$

梯度与偏差的平衡解: 使用正则化减小方差, 避免过拟合 (Overfitting).

$$E = \frac{1}{m} \sum_{k=1}^m E_k + (1-\lambda) \sum_i w_i^2$$

(连接权和
阈值, 也就进
也可用 n^2 , 精度、参数)

BP算去多到局部最小, 为了得到比某局部

最小更优的解, 可以:

① 以不同的参数值初始化多个神经网络

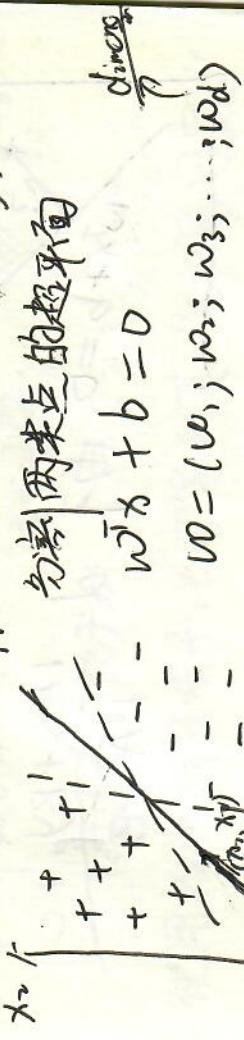
② 模拟退火 以一定的概率接受比当前解更差的解 (保证每个最小步数空间, 限速迭代次数)

③ 随机梯度下降, 即每在局部最小点, 它计算的梯度仍可能不为零.

增加隐层数目提高拟合能力.

2019.10.04

I 支持向量机 (Support vector machine).



当 $w^T x + b > 0$, 点在 $w^T x + b$ 的上方
当 $w^T x + b < 0$, 点在 $w^T x + b$ 的下方.

w 为超平面 $w^T x + b = 0$ 的法向量
由于向量 x 可能有解为 m 与 n . $x = m + n$.

$$w^T (m + n) + b = 0$$

$$w^T m + w^T n + b = 0$$

由于 $w^T n$ 垂直, $w^T n = 0$

$$w^T m + b = 0, w^T m = -b$$

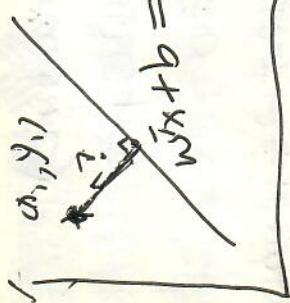
样本空间任意点 x 到超平面 (w, b) 距离:

$$r = \frac{|w^T x + b|}{\|w\|}$$

$$\begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_d \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} + b$$

L_1 范数是指向量各元素之和
然后求平方根

L_2 范数是指向量各元素平方和
之和



$$w^T x + b = 0 \text{ 由于 } d = \frac{|w^T x_i + b|}{\sqrt{w^2 + b^2}}$$

$$d = \frac{|w^T x_i + b|}{\|w\|} = r.$$

$$\begin{cases} w^T x_i + b > 0, & y_i = 1 \\ w^T x_i + b \leq 0, & y_i = -1 \end{cases}$$

两个异类支撑向量到超平面的距离之和：

$$Y = \frac{2}{\|w\|}, \text{ 称为间隔.}$$

这里先固定正样本点和负样本点，超平面心 $x + b = 0$ 的最小距离，再求最优的 w, b 。这里的问题是，我们不一定能“得到” 1 和 -1。这样样的最小距离，那我们可以通过非常小 $e.g. 10^{-10}$ ，只要大于 0，即可，要 Y 最大，

$$\begin{cases} \max \frac{2}{\|w\|} \quad (\text{距离大}) \\ y_i(w^T x_i + b) \geq 1 \quad (\text{分类正确}) \end{cases} \Rightarrow \begin{cases} \min \|w\|^2 \\ y_i(w^T x_i + b) \geq 1 \end{cases}$$

对子， $\begin{cases} \min_{w, b} \|w\|^2 \\ y_i(w^T x_i + b) \geq 1 \end{cases}$ (凸二次规划问题)

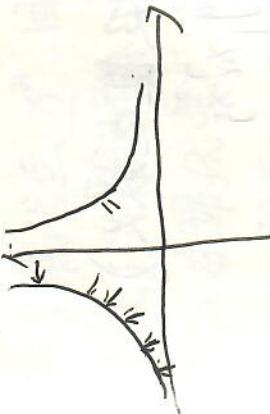
$$y_i(w^T x_i + b) \geq 1, \quad i = 1, 2, 3, \dots, m$$

使用拉格朗日乘子法可得到其“对偶问题”

什么是拉格朗日乘子法？

$$g(x, y) = x^2, \text{ 梯度向量 } \nabla g = \begin{pmatrix} \frac{\partial g}{\partial x} \\ \frac{\partial g}{\partial y} \end{pmatrix} = \begin{pmatrix} 2x \\ 0 \end{pmatrix}$$

梯度向量是等高线的法向量，梯度向量与等高线的切线垂直。



$$\nabla f = \lambda \nabla g \quad (\text{梯度向量平行})$$

$$\begin{cases} x^2 y = 3 \\ x^2 + y^2 \end{cases} \quad (\text{等高线，确定函数})$$

$$\begin{cases} g(x, y) = x^2 \\ f(x, y) = x^2 + y^2 \end{cases} \quad (\text{到原点最近})$$

定理：凸函数在 g 衍生下的极值，这种问题可以表示为：

$$\begin{array}{ll} \min & \int f(x_i) = x_i^2 \\ \text{s.t.} & \sum g_i(x_i) = 0 \end{array}$$



对每条约束加拉格朗日乘子
 $\lambda_i \geq 0$.

$$L(w, b, \lambda) = \frac{1}{2} \|w\|^2 + \sum_{i=1}^m \lambda_i (1 - y_i(w^T x_i + b))$$

$$\partial = (\partial_1, \partial_2, \dots, \partial_m)$$

对 w 和 b 求偏导为零

$$\begin{cases} \partial w = w - \sum_{i=1}^m \lambda_i y_i x_i = 0 \\ \partial b = \sum_{i=1}^m \lambda_i y_i = 0 \end{cases}$$

代入 $L(w, b, \lambda)$ 中: (极值点最大)

$$\max_w \sum_{i=1}^m \lambda_i - \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m \lambda_i \lambda_j y_i y_j x_i^T x_j$$

$$\text{s.t. } \sum_{i=1}^m \lambda_i y_i = 0$$

$$\lambda_i \geq 0$$

$$\begin{cases} y_i f(x_i) - 1 = y_i(w^T x_i + b) - 1 \geq 0 \\ \lambda_i (y_i f(x_i) - 1) = 0 \quad (\lambda_i = 0 \text{ 或 } y_i f(x_i) - 1 = 0) \end{cases}$$

解出 w 后, 求出 $w^T x + b$ 即可得到模型

$$f(x) = w^T x + b$$

$$= \sum_{i=1}^m \lambda_i y_i x_i^T x + b$$



对约束曲面上的任意点 x , 沿点的梯度 $\nabla g(x)$
 正交于约束曲面
 对于最优点 x^* , 固标函数在该点的梯度 $\nabla f(x^*)$
 正交于约束曲面.

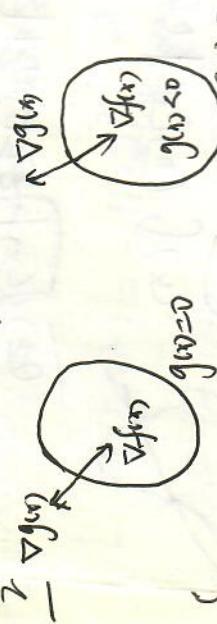
因此在最优点 $\nabla f(x^*)$ 的方向处处相同或
 相反, 即存在 $\nabla f(x^*) + \lambda \nabla g(x^*) = 0$

入为拉格朗日乘子: 于是定义拉格朗日函数
 $L(x, \lambda) = f(x) + \lambda g(x)$

其系数为 0 时的点 (最高点), 为极值与

$$KKT \text{ 条件 } L(x, \mu) = f(x) + \lambda g(x)$$

Vazyme



在等式 $g(x) = 0$ 或在不等式约束 $g(x) \leq 0$
下最小化目标函数 $f(x)$ 。
 $g(x) < 0$:

$g(x) < 0$ 时, 对 $f(x)$ 求极值相当于闭区间求极值, 最
值点即为极值点, 令 $\lambda = 0$, 直接对 $f(x)$ 求梯度即可
得到极值.

$g(x) = 0$,
极值点在边界取得, $f(x)$ 的梯度与 $g(x)$ 的梯度
相反, 从而存在常数 $\lambda > 0$, 使得 $\nabla f(x^*) + \lambda \nabla g(x^*) = 0$

$$\begin{cases} g(x) = 0 \\ \lambda > 0 \\ \lambda g(x) = 0 \end{cases}$$

KKT 条件是目标函数在约束条件下取得极值的充要
条件. 目标函数在约束条件下取得极值时对应
的 x , λ 必须满足 KKT 条件.

Vazyme

对偶问题:

$$\text{由 } \min \frac{1}{2} \|w\|^2 \Rightarrow \max L(v, b, \alpha)$$

对偶问题推导:

$$L(x, \mu) = f(x) + \sum_k \mu_k g_k(x) \quad \mu_k \geq 0, g_k(x) \leq 0.$$

$$\mu_k \geq 0, \quad g_k(x) \leq 0 \Rightarrow \mu_k g_k(x) \leq 0$$

$$\begin{aligned} \max_{\mu} L(x, \mu) &= f(x) \quad (1) \\ \min_x f(x) &= \min_x \max_{\mu} L(x, \mu) \quad (2) \\ \max_{\mu} \min_x L(x, \mu) &= \max_{\mu} \left[\min_x f(x) + \min_x \mu g(x) \right] \\ &= \max_{\mu} \min_x f(x) + \max_{\mu} \min_x \mu g(x) \end{aligned}$$

$$\min_x \mu g(x) = \begin{cases} 0, & \mu = 0 \text{ or } g(x) = 0 \\ -\infty, & \mu > 0 \text{ and } g(x) < 0 \end{cases}$$

\therefore 对 $\min \mu g(x)$ 取 max 时最大只能取 0, 此时

$$\mu = 0 \text{ 或 } g(x) = 0 \quad (3)$$

$$\begin{aligned} \max_{\mu} \min_x L(x, \mu) &= 0 \Rightarrow \mu = 0 \text{ or } g(x) = 0 \\ \max_{\mu} \min_x L(x, \mu) &= \min_x f(x) + 0 = \min_x f(x) \end{aligned}$$

你 $\max_{\mu} \min_x L(x, \mu)$ 为 $\min_{\mu} \max_x L(x, \mu)$
对偶问题的对偶问题



SVM 算法之所以高效，是由子在 Vazyme

固定其他参数后，仅优化两个参数的过程能做到非常高效

$$\alpha_i y_i + \alpha_j y_j = C, \alpha_i \geq 0, \alpha_j \geq 0$$

其 $K \times K$

$$\frac{\partial L(x, \mu)}{\partial x} = 0.$$

回到 SVM:

$$\left\{ \begin{array}{l} \alpha_i (y_i f(x_i) - 1) = 0 \\ y_i f(x_i) - 1 = 0 \end{array} \right.$$

当 $\alpha_i = 0$ ，该样本将不在式 $f(x) = \sum_{j=1}^m \alpha_j y_j x_i^T x + b$ 的求和中出现，也不会对 $f(x)$ 有任何影响。

当 $\alpha_i > 0$ ，则必有 $y_i f(x_i) = 1$

SMO (Sequential Minimal Optimization)

① 选取一对需要更新的变量 α_i 和 α_j

② 固定 α_i 和 α_j 以外的参数，求解

$$\max_{\alpha} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j x_i^T x_j$$

$$\frac{\partial L(x, \mu)}{\partial x} = 0$$

$$\min_x \max_{\mu} L(x, \mu) = \max_{\mu} \min_x L(x, \mu) = \min_x f(x) = f(x^*)$$

$$\left\{ \begin{array}{l} \nabla_k g_k(x) = 0 \end{array} \right.$$

$$C = - \sum_{k \neq i, j} \alpha_k y_k$$

使 $\sum_{i=1}^m \alpha_i y_i = 0$ 。用 $\alpha_i y_i + \alpha_j y_j = C$ 消去 α_j 可求出更新后 α_i 和 α_j 。

对任意支持向量 (x_s, y_s) ，都有 $y_s f(x_s) = 1$

$$y_s (\sum_{i \neq s} \alpha_i y_i x_i^T x_s + b) = 1$$

用所有支持向量求解的平均值：

$$b = \frac{1}{|S|} \sum_{s \in S} \left(\frac{1}{y_s} - \sum_{i \in S} \alpha_i y_i x_i^T x_s \right)$$

III 核函数: Kernel function



解决非线性可分。

将样本从原始空间映射到更高维空间，使其线性可分问题。

如果原始空间是有限维，那么一定存在一个高维特征空间使样本可分。

令 $\phi(x)$ 表示 x 映射后的特征向量， $\phi(x)$ 在特征空间中划为超平面所对应的模型可表示为：

$$\begin{aligned} f(x) &= w^T \phi(x) + b. \\ \min_{w, b} \frac{1}{2} \|w\|^2 \\ \text{s.t. } y_i(w^T \phi(x_i) + b) \geq 1, \quad i &= 1, 2, \dots, m. \end{aligned}$$

对偶问题： $\max_{\alpha} \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m \alpha_i \alpha_j y_i y_j \phi(x_i)^T \phi(x_j)$

$$\begin{cases} \text{s.t. } \sum_{i=1}^m \alpha_i y_i = 0 \\ (\alpha_i \geq 0) \end{cases}$$

$\phi(x_i)^T \phi(x_j)$ 是 x_i 与 x_j 映射到特征空间之后的内积

直接计算 $\phi(x_i)^T \phi(x_j)$ 通常是很困难的

Vazyume

$$K(x_i, x_j) = \langle \phi(x_i), \phi(x_j) \rangle = \phi(x_i)^T \phi(x_j)$$

对偶式重写为：

$$\max_{\alpha} \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m \alpha_i \alpha_j y_i y_j K(x_i, x_j)$$

$$\begin{aligned} \text{s.t. } \sum_{i=1}^m \alpha_i y_i &= 0 \\ (\alpha_i \geq 0) \end{aligned}$$

求解后即可得到：

$$f(x) = w^T \phi(x) + b$$

$$\begin{aligned} &= \sum_{i=1}^m \alpha_i y_i \phi(x_i)^T \phi(x) + b \\ &= \sum_{i=1}^m \alpha_i y_i K(x, x_i) + b. \end{aligned}$$

定理： $K(\cdot, \cdot)$ 是定义在 $\mathcal{X} \times \mathcal{X}$ 上的对称函数，当且仅当对于任意数据集 $D = \{x_1, x_2, \dots, x_m\}$ “核矩阵” K 总是半正定：

$$K = \begin{bmatrix} K(x_1, x_1) & \dots & K(x_1, x_m) \\ \vdots & \ddots & \vdots \\ K(x_m, x_1) & \dots & K(x_m, x_m) \end{bmatrix}$$



由于不知道特征映射的形式，
也不知道原本特征的样例，所
以我们并不知道经过怎样的映射，也就是选择怎样的
核函数，我们并不知道怎样的核函数是合适的。

III 软间隔

允许某些样本不满足约束： $y_i(\mathbf{w}^T \mathbf{x}_i + b) \geq 1$

优化目标： $\min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^m \max(0, 1 - y_i(\mathbf{w}^T \mathbf{x}_i + b) - 1)$

当 C 很大，所有样本都应满足 $y_i(\mathbf{w}^T \mathbf{x}_i + b) \geq 1$

当 C 取有限值，允许一些样本不满足约束

$$h_0(z) = \begin{cases} 1, & z \geq 0 \\ 0, & z < 0 \end{cases}$$

hinge损失： $l_{hinge}(z) = \max(0, 1 - z)$

指数损失： $l_{exp}(z) = \exp(-z)$

对数损失： $l_{log}(z) = \log(1 + \exp(-z))$

采用 hinge 损失： $\min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^m \max(0, 1 - y_i(\mathbf{w}^T \mathbf{x}_i + b))$

s.t. $\epsilon_i \geq 1 - y_i(\mathbf{w}^T \mathbf{x}_i + b)$



使用拉格朗日函数：

$$f(\mathbf{w}) = \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^m \epsilon_i$$

$$g_i(\mathbf{w}) = 1 - y_i(\mathbf{w}^T \mathbf{x}_i + b) - \epsilon_i \leq 0$$

$$g_i(\mathbf{w}) = \epsilon_i \geq 0 \quad -\epsilon_i \leq 0$$

$$\begin{aligned} L(\mathbf{w}, b, \mathbf{w}, \epsilon_i, \mu) &= f(\mathbf{w}) + \partial_i g_i(\mathbf{w}) + \mu_i \epsilon_i \\ &= \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^m \epsilon_i + \sum_{i=1}^m \partial_i (1 - y_i(\mathbf{w}^T \mathbf{x}_i + b) - \epsilon_i) - \sum_{i=1}^m \mu_i \epsilon_i \end{aligned}$$

其中 $\partial_i \geq 0, \mu_i \geq 0$ 是拉格朗日乘子

令 $L(\mathbf{w}, b, \mathbf{w}, \epsilon_i, \mu)$ 对 $\mathbf{w}, b, \epsilon_i$ 的偏导为 0

$$\mathbf{w} = \sum_{i=1}^m \partial_i y_i \mathbf{x}_i$$

$$0 = \sum_{i=1}^m \partial_i y_i$$

$$C = \partial_i + \mu_i$$

$$\begin{aligned} \text{对偶问题, 代入后:} \\ \max_{\mathbf{w}} \sum_{i=1}^m \sum_{j=1}^m \alpha_i \alpha_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j \\ \text{s.t. } \epsilon_i \geq 1 - y_i(\mathbf{w}^T \mathbf{x}_i + b) \end{aligned}$$

KKT 条件要求

$$\begin{cases} \alpha_i \geq 0, \quad \alpha_i \leq 1 \\ y_i f(x_i) - 1 + \varepsilon_i \geq 0 \\ \alpha_i(y_i f(x_i) - 1 + \varepsilon_i) = 0 \\ \varepsilon_i \geq 0, \quad \mu_i \varepsilon_i = 0 \end{cases}$$

IV

核方法

学习得的函数总能表示为 $f(x) = \sum_{i=1}^m \alpha_i y_i K(x, x_i) + b$
表示为核函数 $K(x, x_i)$ 的线性组合。

任何一个核函数都隐式地定义了一个称为“再生核
希尔伯特空间”

令 H 为核函数 K 对应的再生核希尔伯特空间, $\|h\|_H$
表示 H 空间中关于 h 的范数。
对任意单调增函数 $\Omega: [0, \infty] \mapsto R$
任意非负损火函数 $l: R^m \mapsto [0, \infty]$

优化问题: $\min_{h \in H} F(h) = \Omega(\|h\|_H) + l(h(x_1), h(x_2), \dots, h(x_m))$

解得可写为: $h^*(x) = \sum_{i=1}^m \alpha_i K(x, x_i)$ ①



核线性判别分析 (RLDA)



某种映射 $\phi: X \mapsto F$, 将样本映射到特征空间 F ,
然后在 F 中执行线性判别分析, 得到:

$$h(x) = \underline{w^T \phi(x)} \quad ②$$

$$\max_w f(w) = \frac{w^T S_b^\phi w}{w^T S_w^\phi w} \quad \begin{array}{l} \text{between: 美间散度矩阵} \\ \text{within: 美内散度矩阵} \end{array}$$

令 X 表示第 $i \in \{0, 1\}$ 类样本的集合, 样本数 m_i
第 i 类样本在特征空间 F 中的均值为:

$$\mu_i^\phi = \frac{1}{m_i} \sum_{x \in X_i} \phi(x)$$

两/内散度矩阵:

$$S_b^\phi = (\mu_1^\phi - \mu_0^\phi)(\mu_1^\phi - \mu_0^\phi)^T$$

$$S_w^\phi = \sum_{i=1}^2 \sum_{x \in X_i} (\phi(x) - \mu_i^\phi)(\phi(x) - \mu_i^\phi)^T$$

$$K(x, x_i) = \phi(x_i)^T \phi(x) \quad ③$$

$$w = \sum_{i=1}^2 \alpha_i \phi(x_i) \quad ① ② ③$$

令 $K \in \mathbb{R}^{m \times m}$ 为核函数对应的核矩阵. $K_{ij} = K(x_i, x_j)$

$l_i \in \{0, 1\}^m$, 当 $x_i \in \mathcal{X}_1, l_j = 1$ (分类正确)
否则 $l_j = 0$ (分类错误)

$l_0 = (1, 1, 0, 0, 0)$ 属于 0, 为 1.
 $l_1 = (0, 0, 1, 1, 1)$ 属于 1, 为 1.

$$\hat{\mu}_0 = \frac{1}{m_0} K l_0,$$

$$\hat{\mu}_1 = \frac{1}{m_1} K l_1,$$

$$M = (\hat{\mu}_0 - \hat{\mu}_1)(\hat{\mu}_0 - \hat{\mu}_1)^T,$$
$$N = K K^T - \sum_{i=0}^1 m_i \hat{\mu}_i \hat{\mu}_i^T,$$

$$\text{于是: } \max_w J(w) = \frac{w^T S_b^{\phi} w}{w^T S_w^{\phi} w}$$
$$= \max_{\alpha} J(\alpha)$$
$$= \frac{\partial^T M \alpha}{2^T N \alpha}$$

2019.10.05

vazyme

Vazyme

I. 贝叶斯分类器.

在样本 x 上的“条件风险”

$$R(c_i|x) = \sum_{j=1}^N \lambda_{ij} P(c_j|x)$$

的损失

$h: x \mapsto$ 最小化总体风险

$$R(h) = E_x[R(h(x)|x)]$$

$$h^*(x) = \arg \min_{c \in \mathcal{Y}} R(c|x) \quad \text{对每个样本最小化风险.}$$

h^* 为最优贝叶斯分类器, $R(h^*)$ 为贝叶斯风险.

若最小化分类错误, 则误判损失 λ :

$$\lambda_{ij} = \begin{cases} 0, & \text{if } i=j \\ 1, & \text{otherwise} \end{cases}$$

条件风险: $R(c|x) = 1 - P(c|x)$

$$h^*(x) = \arg \max_{c \in \mathcal{Y}} P(c|x)$$

对每个样本 x , 做后验概率 $P(c|x)$ 最大

生成模型必然考虑：

$$P(c|x) = \frac{P(x, c)}{P(x)}$$

$$P(c|x) = \frac{P(c)}{P(x)}$$



II

朴素贝叶斯分类器

利用“属性条件独立性假设”：对已知类别，假设所有属性相互独立，每个属性独立地对分类结果发生影响。

$$\begin{aligned} P(c|x) &= \frac{P(c) P(x|c)}{P(x)} \\ &= \frac{P(c)}{P(x)} \prod_{i=1}^d P(x_i|c) \quad x_1 = \text{青绿, 大, 酸} \\ &= \frac{P(c)}{P(x)} \prod_{i=1}^d P(x_i|c) \quad x_2 = \text{青绿} \\ &= P(x_1, x_2, x_3|c) \quad x_3 = \text{西瓜} \end{aligned}$$

$P(c)$ ：样本空间中各类样本所占的比值

极大似然估计：(源自频率主义学派，根据数据样本来估计概率分布参数)

$$P(D_c|\theta_c) = \prod_{x \in D_c} P(x|\theta_c) \quad (\text{使用 } \theta_c \text{ 对应的 } D_c \text{ 中的所有样本的频率})$$

在所有 θ_c 的可能取值中，找到一个使 $P(D_c|\theta_c)$ 最大的值。

$$\mathcal{L}(\theta_c) = \log P(D_c|\theta_c) = \sum_{x \in D_c} \log P(x|\theta_c)$$

$$P(x|c) \sim N(\mu_c, \sigma_c^2)$$

$$P(c) = \frac{|D_c|}{|D|}$$

D_c, x_i 表示 D_c 中在第 i 个属性上取值为 x_i 的样本组成的集合。

$$P(x_i|c) = \frac{|D_c, x_i|}{|D_c|}$$

为了避免其他属性携带的信息， Vazyme
被训练集中未出现的属性“抹去”，在估计梯度
值时通常要进行“平滑”，常用“拉普拉斯修正”
或者利用训练集中可能的类别数。

$$\hat{P}_{C|C} = \frac{|D_C|}{|D|}$$

$$\hat{P}_{C|C} = \frac{|D_C| + 1}{|D| + N}$$

$$P(x_i|C) = \frac{|D_{C,x_i}|}{|D_C|} \rightarrow \hat{P}_{C|x_i|C} = \frac{|D_{C,x_i}| + 1}{|D_C| + N_i}$$

半朴素贝叶斯分类器
对“属性条件独立性假设”进行一定程度的放松

III.

EM 算法

又表示隐变量集

θ 表示模型参数

x 表示已观测变量

$$LL(\theta | X, Z) = \ln P_C(X, Z | \theta)$$

$$LL(\theta | X) = \ln P_C(X | \theta) = \ln \sum_z P_C(X, Z | \theta)$$

EM 算法用来估计参数隐变量，迭代式的方法。

以初状态 θ^0 为起点：

E. ① 基于 θ^t 推断隐变量 Z 的期望，记为 Z^t
M. ② 基于 Z^t 和 X 对参数 θ 做极大似然估计，记为 θ^{t+1}
直到 θ^t 和 Z^t 稳定

基于 θ^t 计算隐变量 Z 的概率分布 $P(Z | X, \theta^t)$
E. ① 以当前 θ^t 推断 $P(Z | X, \theta^t)$ ，并计算对数似然：

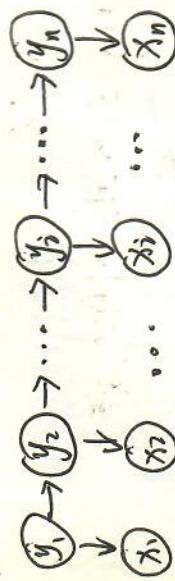
$$LL(\theta | X, Z) = \sum_z Q(\theta | \theta^t) LL(\theta | X, Z)$$

2019.10.07.

概率图模型

Vazyme

隐马尔可夫模型的图结构



所有变量的联合概率分布为：

$$P(x_1, y_1, \dots, x_n, y_n) = P(y_1) P(x_1 | y_1) \prod_{i=2}^n P(y_i | y_{i-1}) P(x_i | y_i)$$

确定一个隐马尔可夫模型还需：

① 状态转移概率： $A = [a_{ij}]_{NM}$

$$a_{ij} = P(y_t = s_j | y_{t-1} = s_i), \quad 1 \leq i, j \leq N.$$

② 观测概率： $B = [b_{ij}]_{NM}$

$$b_{ij} = P(y_t = y_j | x_t = s_i), \quad 1 \leq i \leq N, 1 \leq j \leq M$$

$M \leq N$, 因为有些状态并不能转换为任意其他状态。

③ 初始状态概率： $\pi_i = \underline{\pi_i} = \pi_{s_i}, \quad 1 \leq i \leq N$

$$\pi_i = P(y_1 = s_i), \quad 1 \leq i \leq N.$$

HMM 的问题

Vazyme

① 给定模型 $\lambda = [A, B, \pi]$, 如果

计算 $x = \{x_1, x_2, \dots, x_n\}$ 的概率 $P(x | \lambda)$

② 给定模型 $\lambda = [A, B, \pi]$ 和观测序列

$y = \{y_1, y_2, \dots, y_n\}$, 找出与观测序列 y 最匹配的状态序列 $x = \{x_1, x_2, \dots, x_n\}$

③ 给定观测序列 $y = \{y_1, y_2, \dots, y_n\}$, 调整模型参数 $\lambda = [A, B, \pi]$, 使 $P(x | \lambda)$ 最大。

马尔可夫随机场

M 个变量 $X = \{x_1, x_2, \dots, x_n\}$

所有结构构成的集合为 C

与团 $Q \in C$ 对应的变量集合记为 X_Q

联合概率 $P(x)$:

$$P(x) = \frac{1}{Z} \prod_{Q \in C} \psi_Q(x_Q)$$

ψ_Q 为团 Q 对应的势函数。

$$Z = \sum_x \prod_{Q \in C} \psi_Q(x_Q)$$

若 团 Q 等于 Q^* 极大团。

$$P(X_A, X_B | X_C) = P(X_A | X_C) P(X_B | X_C)$$

X_A 和 X_B 在给定 X_C 时条件独立。

局部马尔可夫性： $X_v \perp X_{V \setminus n^*(v)} \mid X_n(v)$

给定结点 v 的邻近结点 $X_n(v)$, 则 v 与除给定 v 及 $X_n(v)$ 之外的所有结点 $X_{V \setminus n^*(v)}$ 相互独立。

对马尔可夫性：给定所有其他变量, 两个非邻接变量条件独立。

条件随机场 2019. 10. 14. Vazyme
判别式模型

MCMC: (Markov Chain Monte Carlo)

计算 $f(x)$ 在 $P(x)$ 上的期望

$$P(A) = \int_A p(x) dx$$

$$P(f) = E_p[f(x)] = \int_x f(x) p(x) dx.$$

$$\hat{P}(f) = \frac{1}{N} \sum_{i=1}^N f(x_i).$$

22/9. 11. 15.

为什么梯度下降更新参数有效。

$$\Delta n_{\text{eff}} = -n \frac{\partial E_k}{\partial \omega_{\text{in}}}$$

$$f(x) = \int_0^x F(t) dt$$

$$y_2 = f(x_2)(x - x_0)^0$$

$$f(x) = \underbrace{f'(x_0)(x-x_0)}^0 + f(x_0)(x-x_0)^1$$

$$\text{是但是, } f(x) = \sum_{i=0}^n f^{(i)}(x_0)(x-x_0)^i$$

$$f'(x) = f'(x_0)$$

$$f''(x) = 2f''(x_0) \quad ??$$

对 $f^{(m)}(x_0)$ $(x-x_0)^m$ 余项看:

$$m! f_m(x) = m! f_{m+1}(x) \text{ for } x = 0.$$

—m阶领导——：m阶领导以后都是0, m阶领导之前 $(1-x)$ 都是0.

Վաշինգտոն

项目组的沟通与协调

$$y = f(x)$$

$$= w_{k+1} - \eta \frac{\partial E_k}{\partial w_k}$$

η：学习率

$$\frac{\partial E_k}{\partial w_k} : \text{梯度}$$

$$f(x)_{\text{Taylor}} = \sum_{n=0}^{\infty} \frac{f^{(n)}(a)}{n!} \times (x-a)^n \quad (7)$$

$$= f(a) + \frac{f'(a)}{1!}(x-a) + \frac{f^{(2)}(a)}{2!}(x-a)^2 + \dots + \frac{f^{(n)}(a)}{n!}(x-a)^n + R_n(x)$$

$$f(x)_{\text{Taylor}} = \sum_{i=0}^n \frac{f^{(i)}(x_0)}{i!} (x-x_0)^i, \quad x \in (x_0 - \Delta x, x_0 + \Delta x)$$

卷之三

用一阶泰勒展开式证明梯度下降：

$$f(\theta) = f(\theta_0) + f'(\theta_0) \cdot (\theta - \theta_0)$$

设 θ_0 是原点, θ 是向下一步后所处的点。

$$\min: f(\theta) - f(\theta_0) = (\theta - \theta_0) \cdot f'(\theta_0)$$

用向量 \vec{v} 表示 $\theta - \theta_0$: $\vec{v} = \vec{\theta} - \theta_0$

$f(\theta) - f(\theta_0) = \vec{v} \cdot f'(\theta_0) = \|\vec{v}\| \cdot \|f'(\theta_0)\| \cdot \cos \alpha$
要使 $f(\theta) - f(\theta_0)$ 最小, $\cos \alpha = -1$, \vec{v} 与 $f'(\theta_0)$ 方向相反。

$$\vec{v} = -\eta \cdot f'(\theta_0)$$

$$\text{因为: } \vec{v} = \theta - \theta_0, \theta - \theta_0 = -\eta \cdot f'(\theta_0)$$

$$\theta = \theta_0 - \eta \cdot f'(\theta_0)$$

因为要满足一阶泰勒展开式, 所以
 $\theta \rightarrow \theta_0$, $\eta \cdot f'(\theta_0)$ 要小, η 要小, 学习率要小。

生信中的多元统计

Vazyme

矩阵 A. $\begin{bmatrix} m \\ n \end{bmatrix}$: 对 $\begin{bmatrix} m \\ n \end{bmatrix}$ 的线性变换.

$$\begin{bmatrix} m \\ n \end{bmatrix} = m \begin{bmatrix} 1 \\ 0 \end{bmatrix} + n \begin{bmatrix} 0 \\ 1 \end{bmatrix}$$

矩阵 A. 次矩阵 : 对多个变量(样本) 同时做相同的线性变换.
(多变量回归)

矩阵的行列式 $\det(A)$ 表示矩阵的面积、体积等.

正交矩阵 : 方阵 $AA^T = I$

正交向量 : $a^T b = a_1 b_1 + a_2 b_2 + \dots + a_n b_n = 0$

如果某方阵 A 为正交方阵

$$A = \begin{bmatrix} a_{11} & a_{12} & a_{13} \\ a_{21} & a_{22} & a_{23} \\ a_{31} & a_{32} & a_{33} \end{bmatrix} \quad \text{向量 } a \text{ 与 } b \text{ 正交, 且 } \|a\| = \sqrt{a_{11}^2 + a_{12}^2 + a_{13}^2} \text{ 为单位向量}$$

①相似矩阵.

$A = PBP^{-1}$, A 与 B 互为相似矩阵, 有相同的特征值, 那么 A 与

阵可以被对角化 $|A| \neq 0$

②若 $A \in \mathbb{R}^{n \times n}$ 有 n 个线性无关特征向量. $|A| \neq 0$

$A = P \Lambda P^{-1}$, P 为对称矩阵, Λ 为对称矩阵

③若 A 为对称矩阵, 有 n 个线性无关特征向量. $|A| \neq 0$

$A = T \Lambda T^{-1}$, T 为对称矩阵, 特征向量矩阵

$$= (t_1, t_2, \dots, t_p) \begin{pmatrix} \lambda_1 & & \\ & \ddots & \\ & & \lambda_p \end{pmatrix} \begin{pmatrix} t_1 \\ t_2 \\ \vdots \\ t_p \end{pmatrix} = \sum_{i=1}^p \lambda_i t_i t_i^T$$

(特征值分解)

概率论与数理统计.

$$\text{pdf: } \text{cdf: } f(x) = F(x) \quad \int_{-\infty}^x f(x) dx = F(x)$$

$$F(x) = P(X \leq x) \quad X \sim F(x)$$

$$\text{离散型: } F(x) = \sum_{i \leq x} p_i$$

$$\text{连续型: } F(x) = \int_{-\infty}^x f(x) dx.$$

pdf: 概率密度函数 $f(x) = f(x)$

cdf: 累积分布函数 $F(x)$

- 均匀分布 $f(x) = C F(x) = C x$
- 二项分布 $C_n p^k (1-p)^{n-k} \sum_{k=0}^n C_n p^k (1-p)^{n-k} X \sim B(n, p)$
- 负二项分布: 連續不断且独立地重复进行一个参数为 p 的伯努利试验, 记 X 为第 r 次“成功”所离试验次数. k 次试验中, 第 $k-1$ 次成功 $r-1$ 次.

$$P(X=k) = P \cdot P(B(k-1, p) = r-1)$$

$$= P C_{k-1}^{r-1} p^{r-1} q^{k-1-(r-1)} = C_{k-1}^{r-1} p^{r-1} q^{k-r}.$$

$X \sim NB(n, p)$, 限制了第 k 次为成功的 r 次.

$$q = 1-p, \quad k = r, r+1, \dots$$

$$④ \text{泊松分布: } p_k = P(X=k) = e^{-\lambda} \frac{\lambda^k}{k!}, \quad k=0, 1, 2, \dots$$

$$0! = 1 \quad \sum_{k=0}^{\infty} p_k = \sum_{k=0}^{\infty} e^{-\lambda} \frac{\lambda^k}{k!} = e^{-\lambda} \frac{\lambda^0}{0!} = e^{-\lambda} \frac{1}{0!} = 1$$

泊松分布与二项分布的关系:

$$X \sim B(n, p), Y \sim P(np) \quad X \sim np$$

p 很小, n 很大.

$$\frac{B_k(n, p)}{B_{k-1}(n, p)} = \frac{C_n^k p^k (1-p)^{n-k}}{C_n^{k-1} p^{k-1} (1-p)^{n-k+1}} = \frac{\frac{p}{n} \frac{(n-1)p}{n-1} \frac{(n-2)p}{n-2} \cdots \frac{p}{2} \frac{p}{1}}{(n-1)! (n-k)!} \frac{(n-k+1)p}{1-p}$$

$$B_0(n, p) = (1-p)^n = \left(1 - \frac{p}{n}\right)^n \approx e^{-p}$$

$$\begin{aligned} &= \frac{np - kp + p}{k(1-p)} = \frac{(n-k)p + p}{k(1-p)} \\ &= \left(1 + \left(-\frac{p}{n}\right)\right)^{(-\frac{p}{n}) \cdot (n)} \approx \frac{1}{k} \\ &\approx e^{-p} \end{aligned}$$

几何分布 (逆向概率).

几何分布: $(X \geq 1)$ 为首次成功.

$$P(X=k) = (1-p)^{k-1} p \quad 0 < p < 1$$

$$X \sim Ge(p)$$

指数分布: $(X \geq 0)$ 逆线.

$$X \sim \text{Exp}(\lambda), \lambda > 0 \quad F(x) = P(X \leq x) = \int_0^x \lambda e^{-\lambda x} dx$$

$$f(x) = \begin{cases} 0 & x \leq 0 \\ \lambda e^{-\lambda x} & x > 0 \end{cases}$$

正态分布: $X \sim N(\mu, \sigma^2)$

$$f(x) = \frac{1}{\sqrt{2\pi} \cdot \sigma} \cdot e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

$\mu=0, \sigma^2=1, X \sim N(0, 1)$ 标准正态分布.

$X \sim N(\mu, \sigma^2), \frac{X-\mu}{\sigma} \sim N(0, 1)$.

若 X 为分析时的标准化可能就是转换为
标准正态分布.

χ^2 分布: $X \sim N(0, 1)$, X^2 为一个自由度的 χ^2 分布.

伯努利分布

$$P(X=k) = \binom{n}{k} p^k q^{n-k}$$

$$P(X=k) = e^{-p} \frac{p^k}{k!}$$

$$p = \frac{1}{n}$$

二项分布

$$P(X=k) = \binom{n}{k} p^k q^{n-k}$$

$$P(X=k) = \binom{n}{k} p^k q^{n-k}$$

$$p = \frac{1}{n}$$

$$p = \frac{1}{n}$$

随机变量 X 期望期望

$$E(g(X)) = \int_{-\infty}^{+\infty} g(x) f(x) dx$$

随机变量 X 的离散程度，
随机变量 X 与随机变量 Y 的相关程度。

$$\begin{aligned} \text{Var}(X) &= E((X - E(X))^2) \\ &= E(X^2 - 2X E(X) + E(X)^2) \\ &= E(X^2) - 2E(X)E(X) + E(X)^2 \\ &= E(X^2) - E(X)^2 \end{aligned}$$

$$\begin{aligned} X \text{ 与 } Y \text{ 独立: } \text{Var}(X+Y) &= \text{Var}(X) + \text{Var}(Y) \\ X \text{ 与 } Y \text{ 不独立: } \text{Var}(X+Y) &= \text{Var}(X) + \text{Var}(Y) + \text{Cov}(X, Y) \end{aligned}$$

1. 不相关可以不独立

Vazyme



协方差

$$\begin{aligned} \text{Cov}(X, Y) &= E[(X - E(X))(Y - E(Y))] \\ &= E[XY - XE(Y) - YE(X) + E(X)E(Y)] \\ &= E(XY) - E(X)E(Y) - E(Y)E(X) + E(X)E(Y) \\ &= E(XY) - E(X)E(Y) \\ \text{Cov}(X, Y) &= \frac{\text{Cov}(X, Y)}{\sqrt{\text{Var}(X) \cdot \text{Var}(Y)}} = \frac{\text{Cov}(X, Y)}{\sqrt{6X \cdot 6Y}} = R \end{aligned}$$

马尔科夫不等式：

Vazyme

$$\text{Cov}(X, Y) = E[(X - E(X))(Y - E(Y))]$$

$$\begin{aligned} &= E[XY - XE(Y) - YE(X) + E(X)E(Y)] \\ &= E(XY) - E(X)E(Y) - E(Y)E(X) + E(X)E(Y) \\ &= E(XY) - E(X)E(Y) \end{aligned}$$

$$\begin{aligned} \text{Cov}(X, Y) &= \frac{\text{Cov}(X, Y)}{\sqrt{\text{Var}(X) \cdot \text{Var}(Y)}} = \frac{\text{Cov}(X, Y)}{\sqrt{6X \cdot 6Y}} = R \end{aligned}$$

$$\text{Cov}(X, Y) = \frac{\text{Cov}(X, Y)}{\sqrt{\text{Var}(X) \cdot \text{Var}(Y)}} = \frac{\text{Cov}(X, Y)}{\sqrt{6X \cdot 6Y}} = R$$



፩፻፲፭

假設檢驗：

第一类错误： α ：拒绝真
第二类错误： β ：接受假

原假设没什。一般用不好的结果,为样假设H,用好的

PCA 理解 (通俗解)

对A矩阵做PCA,线性转换

$PA = C$ 转换矩阵为 P (P为过渡矩阵, 因为投票)
均为 1-1-1

卷之三

$$\Delta = \text{Cov}(C, C) = \frac{C^T C}{n-1} = \text{Var}(C)$$

后方的胸椎小方差，不同后胸椎小方差为0。

$$\mathbf{D} = \text{Var}(\mathbf{C}) = \text{Var}(\mathbf{PA}) = \mathbf{P} \text{Var}(\mathbf{A}) \mathbf{P}^T$$

$$= \mathbf{P} \cancel{\mathbf{A}} \cdot \mathbf{A}^T \cancel{\mathbf{P}}^T$$

$$= \mathbf{P} \mathbf{B} \mathbf{P}^T$$

$$\mathbf{B} = \mathbf{A} \mathbf{A}^T$$

$$D = P B P^{-1} = P B P^{-1}$$

$$P^T D P = B \quad \begin{matrix} \triangleright \text{(特征值矩阵)} \\ \triangleright \text{(特征向量矩阵)} \end{matrix}$$

(这里根据定义 P 正交, D 对角, 则 P 为 B 的特征向量矩阵, D 为 B 的特征值矩阵).

R 漢語 gerl geerz P
sample
spraker

$$\Rightarrow A^T = A_{\text{rep}}^T$$

Apxn.

$$n = \sqrt{1 + \sqrt{1 + \sqrt{1 + \dots}}}$$

PCA之前可以挑选高变异的 gene.
挑选方差, C1 高的基因.
500, 1000, 3000 个基因.



Vazyme

Vazyme

用考夫斯基距离

$$d(x, y) = \left[\sum_{i=1}^p |x_i - y_i|^p \right]^{\frac{1}{p}}$$

当 $p=1$ 时, 为绝对值距离 (曼哈顿距离, 街区距离)

$$d(x, y) = \sum_{i=1}^p |x_i - y_i|$$

当 $p=2$ 时, 为欧氏距离

当 $p=\infty$ 时, 为切比雪夫距离

$$d(x, y) = \max_{1 \leq i \leq p} |x_i - y_i|$$

聚类

层次聚类

hierarchical clustering method.
聚类法, 按距离, 每次设置 cutoff, 最后聚成一类.

分层法, 由少到多



(1) cutoff 设置?

- (2) 计算距离?
- (3) 选择参考点

兰氏距离 (Lance and Williams distance)

$$d(x, y) = \sum_{i=1}^p \frac{|x_i - y_i|}{x_i + y_i}$$

所有数据为正, $d(x, y) = \sum_{i=1}^p \frac{|x_i - y_i|}{x_i + y_i}$
该距离与各度量的单位无关, 且适用于高度偏斜或含异常值的数据.

度量 x, y 之间的距离满足:

- (1) 非负性: $d(x, y) \geq 0$
- (2) 对称性: $d(x, y) = d(y, x)$
- (3) 三角不等式: $d(x, y) \leq d(x, z) + d(z, y)$

相似系数、相关系数 \rightarrow

$$\cos \theta = \frac{x^T y}{\|x\| \|y\|} = \frac{\text{Cor}(x, y)}{\sqrt{V(x)} \sqrt{V(y)}}$$

Pearson 相关系数

相似系数与距离之间的关系

距离 \Rightarrow 相似系数

$$C_{ij} = \frac{1}{1+dist_{ij}}$$

相似系数 \Rightarrow 距离 相似系数矩阵 (C_{ij}) 为:

$$dist_{ij} = \sqrt{1 - C_{ij}}$$

$dist$ (距离矩阵) \rightarrow $hclust$ (层次聚类)
($hclust$ 也是层次聚类, 聚类 gene module).



④ 指引 $Y_i = \begin{bmatrix} Y_{1i} \\ Y_{2i} \\ \vdots \\ Y_{pi} \end{bmatrix}$, 全 $B = Y^T Y = \begin{bmatrix} Y_1 \\ Y_2 \\ \vdots \\ Y_m \end{bmatrix} \begin{bmatrix} Y_1 & Y_2 & \cdots & Y_m \end{bmatrix}$
中 \sim B (B 为矩阵)

③ $dist_{ij} = d(x_i, x_j) = \|\vec{Y}_i - \vec{Y}_j\|$ (欧式距离)

$$dist_{ij} = \|Y_{1i}\|^2 + \|Y_{2i}\|^2 + \cdots + \|Y_{pi}\|^2$$

$$dist_{ij} = b_{ii} + b_{jj} - 2b_{ij}$$

$$b_{ij} = \frac{1}{2} (b_{ii} + b_{jj} - dist_{ij})$$

⑥ $\frac{m}{2} b_{ij} = \sum_{i=1}^m \vec{Y}_{i \cdot} \cdot \vec{Y}_j = \sum_{i=1}^m (Y_{1i} Y_{1j} + Y_{2i} Y_{2j} + \cdots + Y_{pi} Y_{pj})$

MDS (多维标度法)

核心目的: 在低维空间中寻找一个矩阵 Z ,
使得 Z 构成的欧式距离矩阵 B 尽可能和原来矩阵保持
一致。此种多维标度法也称作主坐标分析。

① 样本矩阵 $X = \begin{bmatrix} x_{11} & \cdots & x_{1m} \\ \vdots & \ddots & \vdots \\ x_{n1} & \cdots & x_{nm} \end{bmatrix} = \begin{bmatrix} \vec{x}_1 & \vec{x}_2 & \cdots & \vec{x}_m \end{bmatrix}$ n sample, m gene.

$$② D = \begin{bmatrix} d(x_1, x_1) & \cdots & d(x_1, x_m) \\ \vdots & \ddots & \vdots \\ d(x_m, x_1) & \cdots & d(x_m, x_m) \end{bmatrix}$$

距离矩阵.

③ $\vec{x}_i = \begin{bmatrix} x_{1i} \\ x_{2i} \\ \vdots \\ x_{pi} \end{bmatrix} \Rightarrow \vec{Y}_i = \begin{bmatrix} Y_{1i} \\ Y_{2i} \\ \vdots \\ Y_{pi} \end{bmatrix}$ 它新的距离矩阵 B ?

$$\begin{aligned} &= y_{11} y_{1j} + y_{21} y_{2j} + \cdots + y_{p1} y_{pj} + \\ &= y_{12} y_{1j} + y_{22} y_{2j} + \cdots + y_{p2} y_{pj} + \\ &= y_{1m} y_{1j} + y_{2m} y_{2j} + \cdots + y_{pm} y_{pj} \\ &= y_{1j} \frac{\sum_{i=1}^m y_{1i}}{m} + y_{2j} \frac{\sum_{i=1}^m y_{2i}}{m} + \cdots + y_{pj} \frac{\sum_{i=1}^m y_{pi}}{m} \\ &= 0 \end{aligned}$$

⑦ $\sum_{i=1}^m dist_{ij}^2 = \sum_{i=1}^m (b_{ii} + b_{jj} - 2b_{ij}) = \sum_{i=1}^m b_{ii} + \sum_{i=1}^m b_{jj} - 2 \sum_{i=1}^m b_{ij}$

$$= \text{tr}(B) + \text{tr}(B) + m b_{jj} + 0$$

$$\sum_{i=1}^m dist_{ij}^2 = \text{tr}(B) + m b_{jj}$$

$$\begin{cases} \text{tr}(B) = \frac{1}{2m} \sum_{i=1}^m \sum_{j=1}^m \text{dist}_{ij}^2 \\ \sum_{j=1}^m \text{dist}_{ij}^2 = \text{tr}(B) + m b_{ij} \\ \sum_{j=1}^m \text{dist}_{ij}^2 = \text{tr}(B) + m b_{ji} \end{cases}$$

$$b_{ij} = \frac{1}{2} \{ (b_{ji} + b_{ij}) - \text{dist}_{ij} \}$$

对于: $B = Y^T Y$, B 是对称矩阵, 正交对角化

$$\begin{aligned} B &= Q \Lambda Q^T = Y^T Y \\ Y &= Q \Lambda^{\frac{1}{2}} \end{aligned}$$

主要思想推导:

$$\begin{cases} \text{dist}_{ij} = d(x_i, x_j) = \|y_i - y_j\| \\ B_{ij} = b_{ij} = d(y_i, y_j) = \|y_i - y_j\| \end{cases}$$

(方便计算).

注意: ① 当 dist 不矩阵本身为欧氏距离 $MDS = PCA$
 ② MDS 强调对不同种类 dist 都可转化为欧氏距离



t-SNE

① 高维空间概率分布 P ,

② 低维空间概率分布 Q ,

③ 使 P 与 Q 尽可能相同, 通过迭代.

最小信息熵

$$H = -\sum_{i=1}^n P_i \log(P_i)$$

motif

$$\text{bits} \sqrt{A \quad C \quad G \quad T} \quad \text{通过信息熵计算.}$$

KL 故度: 测量两个分布是否相似

$$D_{KL}(P||Q) = \sum_i P(i) \ln \frac{P(i)}{Q(i)}$$

当两个概率分布完全相同时, $D_{KL} = 0$, 在 R 中 $D_{KL} > 0$.

t-SNE 的数学过程.

1. 他通过欧式距离离散点的分布.

$$P_{j|i} = \frac{\exp(-\|x_i - x_j\|^2 / (2\sigma_i^2))}{\sum_{k \neq i} \exp(-\|x_i - x_k\|^2 / (2\sigma_i^2))} \quad (\text{以 } x_i \text{ 为中心})$$

$$Q_{j|i} = \exp(-\|y_i - y_j\|^2)$$

2. 希望高维和低维之间的分布尽可能相同。

$$C = \sum_i KL(P_i || Q_i) = \sum_i \sum_j P_{ji} \log \frac{P_{ji}}{Q_{ji}} \quad (\text{最小化})$$

3. 初始化时需要设置参数困惑度 (perplexity), SNE 会通过二分搜索的方法计算最佳的初始方差。

一般困惑度选择 5-50。

4. 对目标函数求梯度

$$\frac{\partial C}{\partial y_j} = 2 \sum_i (P_{ji} \ln -q_{(j|i)} + p_{(j|i)} - q_{(j|i)}) (y_j - y_i)$$

$$y^{(t+1)} = y^{(t-1)} + \eta \frac{\partial C}{\partial y_j} + \Delta(t) (y^{(t-1)} - y^{(t-2)})$$

5. 退代

$y^{(t+1)} = y^{(t-1)} + \eta \frac{2C}{2\eta} + \Delta(t) (y^{(t-1)} - y^{(t-2)})$

从 $y^{(t+1)} = y^{(t-1)} + \eta \frac{2C}{2\eta} + \Delta(t) (y^{(t-1)} - y^{(t-2)})$ 可以看出 $y^{(t+1)}$ 与 $y^{(t-1)}$ 之间的关系。

6. SNE

$P_{ji} = \frac{P_{ji} + P_{ji}}{2n}$

$$q_{(j|i)} = \frac{(1 + \|y_i - y_j\|^2)^{-1}}{\sum_{k \neq i} (1 + \|y_i - y_k\|^2)^{-1}} \quad (\text{梯度})$$

受初值影响。 (RTSNE)

PCA 理解 (from Cross Validated)



Vazyme

$X (n \times p)$

$C = X^T X$

$C = U L V^T$ (对角化, 对称矩阵正交对角化)

$X = U S V^T$ (奇异值分解)

$$C = X^T X = U S^2 V^T U S V^T = U S^2 V^T V U^T = U L V^T$$

$\lambda_{ij} = S_i^2$

坐标 (坐标): $X V = U S V^T V = U S$

① $X = U S V^T$, V 的列向量是特征向量, 是投影方向。

② $X V$ 或 $U S$ 是投影后的坐标。

③ 特征矩阵 C 的特征值 $\lambda_{ij} = S_i^2$

④ 标准化的 (除以特征值的) 行数是 U 的列向量, 从 $U S$ 为坐标可以看出。

⑤ $V \Sigma$ 是 loading (系数矩阵与主成分的相关系数)

⑥ 以上的条件是 行为样本, 列为变量, 否则 U, V 变换。

⑦ PCA on correlation matrix X , X 的列需要中心化 + 标准化。

⑧ $X_k = U_k S_k V_k^T$ (U 前 K 列, S 前 $K \times K$, V 前 K 行)。

X_k 包含前 K 个 PCs (主成分) 的压缩数据。

⑨ 若 $n > p$, U_{nn} , V_{pp} , V_{pp} 的最后 $n-p$ 列 (次要), 对应的 S 的行为 0。

多元线性回归

偏最小二乘回归 \star 变量多时效果较好。
岭回归

一元线性回归 $R^2 = r^2$

多元线性回归最小二乘：

$b = (X^T X)^{-1} X^T y$, 如果 X 有共线向量 b 无解。
如果 X 中有变量强相关, b 估计不稳。

判断是否存在共线性：

tolerance = $1 - R_j^2$ (第 j 个变量).

VIF (variance inflation factor)

$$VIF_j = \frac{1}{1 - R_j^2}$$

Kappa: $K = \sqrt{\frac{\lambda_{\max}}{\lambda_{\min}}}$ λ_i 为 $X^T X$ 的特征值.

($Kappa > 15$) 或 $VIF > 25$ 有强共线性

三、线性回归 / 2 跛声

$$f(w) = w^T x \quad \begin{array}{l} \text{数据满足正态分布} \\ \text{最小二乘都是最小二乘} \end{array}$$

$$\text{Loss function: } L(w) = \sum_{i=1}^N \|w^T x_i - y_i\|^2$$

$$\frac{dL(w)}{dw} \Rightarrow \hat{w} = (x^T x)^{-1} x^T y$$

$$X \in \mathbb{R}^{N \times P}, \quad N \text{ 个样本}, \quad x_i \in \mathbb{R}^P \quad N \rightarrow P$$

- 过拟合 \Rightarrow  ① 加数据
 ② 特征选择 / 特征提取 (PCA)
 ③ 正则化

正则化框架

$$\arg \min [L(w) + P(w)]$$

↓
loss
penalty.

 loss 可以将参数取 0
(lasso / ridge?)

$$L_1: \text{Lasso}, \quad P(w) = \|w\|_1 \quad L_2: \text{Ridge}, \quad P(w) = \|w\|^2$$

推导 L_2 的 \hat{w} :

$$\begin{aligned} \min_w \sum_{i=1}^N & \|w^T x_i - y_i\|^2 + \lambda w^T w \\ (w^T x_i - y_i, w^T x_i - y_i, \dots, w^T x_n - y_n) & \begin{pmatrix} w^T x_1 - y_1 \\ w^T x_2 - y_2 \\ \vdots \\ w^T x_n - y_n \end{pmatrix} = \hat{w} (x_1, x_2, \dots, x_n) - (y_1, y_2, \dots, y_n) \\ & = (x^T x^T - y^T) (x w - y) \end{aligned}$$

$$\begin{aligned} & \hat{w} = \underbrace{(x^T x^T - y^T)^{-1}}_{X^T X + \lambda I} \underbrace{(x^T x^T - y^T) y + \lambda w^T w}_{X^T X y + \lambda w^T w} \\ & = X^T X y + \lambda w^T w \end{aligned}$$



Vazyme

$$\hat{w} = \arg \min J(w)$$

$$\frac{\partial J(w)}{\partial w} = 2(X^T X + \lambda I)w - 2X^T y = 0$$

$$\hat{w} = (X^T X + \lambda I)^{-1} X^T y$$

惩罚项

最小二乘估计 \Rightarrow 极大似然估计
MLE (Maximum Likelihood Estimation)

LSE \Leftrightarrow MAP (最大后验)
Regularized LSE \Leftrightarrow MAP (noise 和 prior 也是 CID)

(IV) 线性分类
硬输出: 线性分类机 (输出结果)
软输出: 判别: logistic regression $P(X)$
(输出概率)

$$P(X) = \frac{P(X|Y)P(Y)}{P(X)}$$

恒等
max $P(X|Y) \Rightarrow \max(P(X|Y)P(Y))$
必然
后验

$$\bar{z} = \sum_{i=1}^N z_i = \frac{1}{N} \sum_{i=1}^N w^T x_i$$

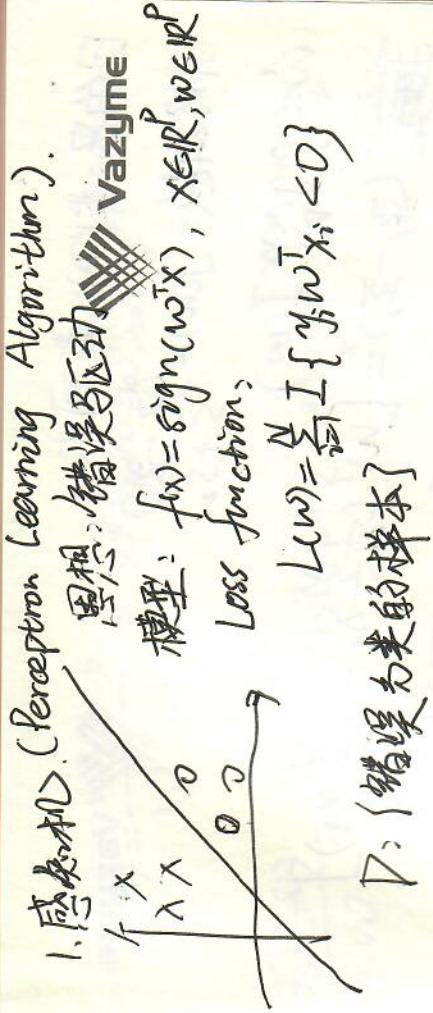
$$S_1 = \frac{1}{N} \sum_{i=1}^N (z_i - \bar{z}) (w^T x_i - \bar{z})^T$$

$$C_1: \bar{z}_1 = \frac{1}{N} \sum_{i=1}^N w^T x_i$$

$$C_2: \bar{z}_2 = \frac{1}{N} \sum_{i=1}^N w^T x_i$$

$$S_0 = \frac{1}{N} \sum_{i=1}^N (w^T x_i - \bar{z}_2)^T$$

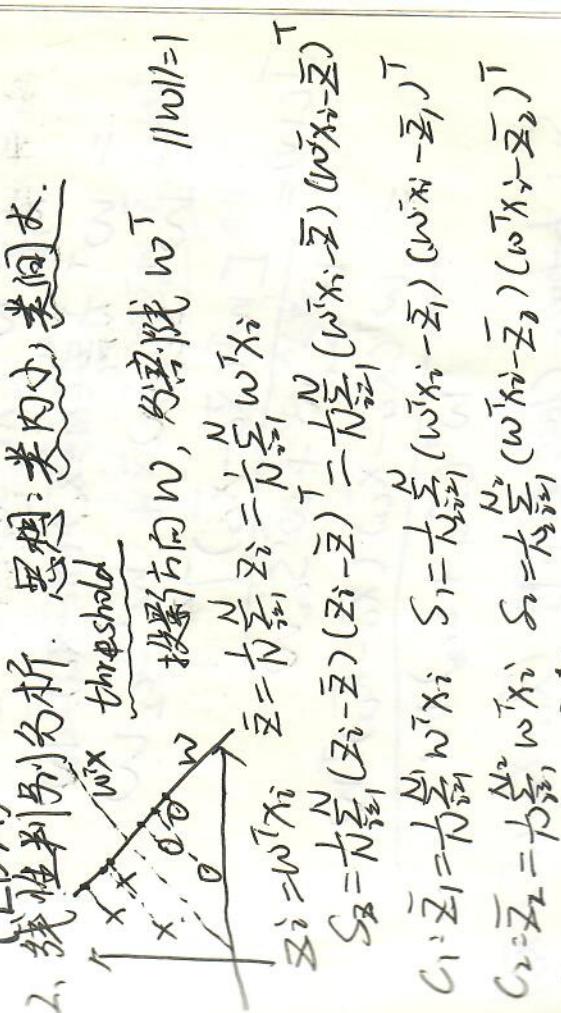
类间: $(\bar{z}_1 - \bar{z}_2)^2 \uparrow$



$$L(w) = \sum_{i \in E} -y_i w^T x_i$$

$$\frac{\partial L(w)}{\partial w} = \sum_{i \in E} -y_i x_i$$

SGD: $w^{(t+1)} = w^{(t)} - \frac{\partial L(w)}{\partial w}$
 $= w^{(t)} + \underbrace{y_i x_i}_{\text{权重更新}} \quad (\text{权重更新})$



目的是：类内小，类间大。

$$\text{目标函数: } J(w) = \frac{(\bar{x}_1 - \bar{x}_2)^2}{S_1 + S_2}$$

$$J(w) = \arg \max_w J(w)$$

$$J(w) = (\bar{x}_1 - \bar{x}_2)^2 = \left[w^T \left(\frac{1}{N} \sum_{i=1}^N x_i - \frac{1}{N} \sum_{i=1}^N x_i \right) \right]^2$$

$$= \left[w^T (\bar{x}_{c1} - \bar{x}_{c2}) \right]^2$$

$$S_1 + S_2 = \frac{1}{N} \sum_{i=1}^N (w^T x_i - \bar{x}_1) (w^T x_i - \bar{x}_1)^T +$$

$$\frac{1}{N} \sum_{i=1}^N (w^T x_i - \bar{x}_2) (w^T x_i - \bar{x}_2)^T$$

$$= \frac{1}{N} \sum_{i=1}^N (w^T x_i - \frac{1}{N} \sum_{j=1}^N w^T x_j) (w^T x_i - \frac{1}{N} \sum_{j=1}^N w^T x_j)^T + \dots$$

$$= w^T \frac{1}{N} \sum_{i=1}^N (x_i - \bar{x}_{c1}) (x_i - \bar{x}_{c1})^T w$$

$$= w^T S_{c1} w + w^T S_{c2} w.$$

$$J(w) = \frac{\left[w^T S_{c1} w + w^T S_{c2} w \right]}{w^T (S_{c1} + S_{c2}) w}$$

$$= \frac{w^T (\bar{x}_{c1} - \bar{x}_{c2}) (\bar{x}_{c1} - \bar{x}_{c2})^T w}{w^T (S_{c1} + S_{c2}) w}$$

S_b : between-class 距离差 $(\bar{x}_{c1} - \bar{x}_{c2})(\bar{x}_{c1} - \bar{x}_{c2})^T$



$$J(w) = \frac{w^T S_b w}{w^T S_w w} = w^T S_b w (w^T S_w w)^{-1} \cdot w$$

$$\frac{\partial J(w)}{\partial w} = 2 S_b w (w^T S_w w)^{-1} + w^T S_b w \cdot (-1) \cdot (w^T S_w w)^{-2} \cdot 2w$$

$$= 0$$

$$S_b w (w^T S_w w)^{-1} = w^T S_b w (w^T S_w w)^{-2} \cdot w$$

$$S_b w (w^T S_w w) = \underbrace{w^T S_b w}_{\in \mathbb{R}} S_w$$

$$(w) = \underbrace{w^T S_w w}_{w^T S_b w} S_w^{-1} \cdot S_b \cdot w$$

$$w = P x$$

$$w^T = I x P$$

$$S_w = P x P$$

$$(Euler) \quad \& \quad S_w^{-1} \underbrace{S_b \cdot w}_{(Euler)}$$

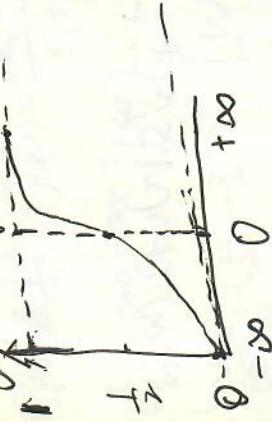
$$(w^T x_{c1} - \bar{x}_{c1})(w^T x_{c1} - \bar{x}_{c1})^T w$$

$$\alpha S_w^{-1} (\bar{x}_{c1} - \bar{x}_{c2})$$

Logistic Regression. (线性模型)



Vazyme



Gaussian Discriminant Analysis

Vazyme

Gaussian Discriminant Analysis

Vazyme

$P(y|x) \propto P(x|y) P(y)$

$\hat{y} = \arg \max_{y \in \{0, 1\}} P(y|x)$

$$P_1 = P(y=1|x) = \sigma(w^T x) = \frac{1}{1+e^{-w^T x}}, \quad y=1$$

$$P_0 = P(y=0|x) = 1 - P(y=1|x) = \frac{e^{-w^T x}}{1+e^{-w^T x}}, \quad y=0$$

$$P(y|x) = P_1^y P_0^{1-y}$$

MLE: $\hat{w} = \arg \max_w \log P(y|x)$

Maximum Likelihood Estimate. $= \arg \max_w \log \prod_{i=1}^N P(y_i|x_i)$

Maximum Likelihood Estimate. $= \arg \max_w \sum_{i=1}^N \log P(y_i|x_i)$

$= \arg \max_w \sum_{i=1}^N \frac{\log P_1^y_i}{P_0} + (1-y_i) \frac{\log P_0}{P_1}$

MLE \Rightarrow loss function
min. max

$$x|y=1 \sim N(\mu_1, \Sigma)$$

$$x|y=0 \sim N(\mu_0, \Sigma)$$

$$y \sim \text{Bernoulli}(\phi) \Rightarrow \frac{y}{P} \cdot \frac{1}{1-P}, y=0$$

log-likelihood:

$$L(\theta) = \log \prod_{i=1}^N P(x_i, y_i) = \sum_{i=1}^N \log (P(x_i|y_i) P(y_i))$$

$$= \sum_{i=1}^N \left[\log P(x_i|y_i) + \log P(y_i) \right]$$

$$= \sum_{i=1}^N \left[\log N(\mu_1, \Sigma)^{y_i} + \log N(\mu_0, \Sigma)^{1-y_i} + \log \phi^{y_i} (1-\phi)^{1-y_i} \right]$$

Naive Bayes Classifier.



最简单的有向概率图

$\hat{y} = \arg \max P(y|x)$

$$= \arg \max_{y=0,1} \frac{P(x,y)}{P(x)} = \arg \max_{y=0,1} P(x,y)$$

$$= \arg \max_y P(y) \cdot P(x|y)$$

条件独立性假设： $P(x|y) = \prod_{j=1}^p P(x_j|y)$
 \times 离散 $\rightarrow x_j$ un Categorical Dist
 $\left\{ \times$ 连续 $\rightarrow x_j \sim N(\mu_j, \sigma_j^2) \right\}$

$$\begin{aligned} \text{Sample Mean: } \bar{x}_{\text{pri}} &= \frac{1}{N} \sum_{i=1}^N x_i = \frac{1}{N} (x_1 - \bar{x}_1) \begin{pmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{pmatrix} = \frac{1}{N} x^T \bar{1}_N \\ \text{Sample Covariance: } S &= \frac{1}{N} \sum_{i=1}^N (x_i - \bar{x})(x_i - \bar{x})^T \\ &= \frac{1}{N} \sum_{i=1}^N (x_i - \bar{x})(x_i - \bar{x})^T = \frac{1}{N} x^T H x \\ &= \frac{1}{N} \sum_{i=1}^N (x_i - \frac{1}{N} x^T \bar{1}_N)(x_i - \frac{1}{N} x^T \bar{1}_N)^T \\ &= \frac{1}{N} \underbrace{(x_1 - \bar{x}, x_2 - \bar{x}, \dots, x_N - \bar{x})}_{(x_1, x_2, \dots, x_N - \bar{x})^T} \begin{pmatrix} (x_1 - \bar{x})^T \\ (x_2 - \bar{x})^T \\ \vdots \\ (x_N - \bar{x})^T \end{pmatrix} \end{aligned}$$

H_N centering matrix

$$\begin{aligned} H^T \cdot H &= (\bar{1}_N - \frac{1}{N} x^T \bar{1}_N \bar{1}_N^T) (\bar{1}_N - \frac{1}{N} x^T \bar{1}_N \bar{1}_N^T) \\ &= x^T - \bar{x} \begin{pmatrix} 1 & 1 & \dots & 1 \end{pmatrix} = x^T - \bar{x} \bar{1}_N^T = x^T - \frac{1}{N} x^T \bar{1}_N \bar{1}_N^T \\ &= \frac{1}{N} x^T (\bar{1}_N - \frac{1}{N} x^T \bar{1}_N \bar{1}_N^T) \cdot (\bar{1}_N - \frac{1}{N} x^T \bar{1}_N \bar{1}_N^T)^T x \\ &= \frac{1}{N} x^T H \cdot H^T x \Rightarrow \underbrace{\frac{1}{N} x^T H x}_H \end{aligned}$$

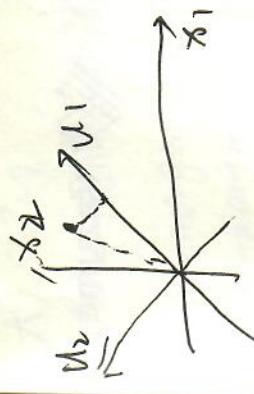
$$H^T = H$$

$$H^T = \frac{1}{N} \sum_{i=1}^N (x_i - \bar{x})(x_i - \bar{x})^T$$

一个中心：原始特征空间的重构
相关 \rightarrow 无关

两个基本点：最大投影方差
最小重构距离

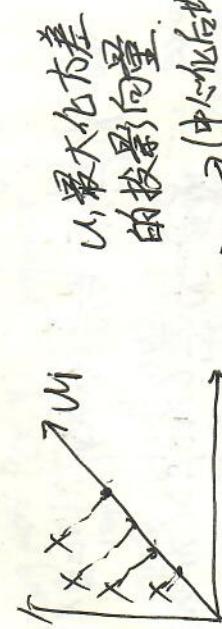
Vazyme



解释 $\sum_{i=1}^N (x_i^T u_1)^2$ u_1
 x_i 为第 i 个样本
特征向量

Vazyme

$$\begin{aligned}
 & \text{J} = \frac{1}{N} \sum_{i=1}^N ((x_i - \bar{x})^T u_1)^2 \rightarrow \text{(中心化后方差已归为0)} \\
 & \text{J} = \frac{1}{N} \sum_{i=1}^N u_1^T (x_i - \bar{x}) (x_i - \bar{x})^T u_1 \\
 & = \frac{1}{N} \sum_{i=1}^N u_1^T S u_1 \\
 & \text{s.t. } u_1^T u_1 = 1
 \end{aligned}$$



U₁ 最大化方差
的投影向量

$$\begin{aligned}
 & \text{J} = \frac{1}{N} \sum_{i=1}^N \|x_i - \bar{x}\|^2 \quad \text{(前 Q 维相同)} \\
 & = \frac{1}{N} \sum_{i=1}^N \left\| \frac{1}{\|u_1\|} (x_i^T u_1) u_1 \right\|^2 \\
 & = \frac{1}{N} \sum_{i=1}^N \frac{1}{\|u_1\|^2} \left[(x_i^T u_1) u_1 \right]^2
 \end{aligned}$$

特征
向量

$$\begin{aligned}
 & \triangleq \frac{1}{N} \sum_{i=1}^N \frac{1}{\|u_1\|^2} ((x_i - \bar{x})^T u_1)^2 \\
 & = \frac{1}{N} \sum_{i=1}^N \frac{1}{\|u_1\|^2} ((x_i - \bar{x})^T u_1)^2 \\
 & = \underbrace{\frac{1}{N} \sum_{i=1}^N ((x_i - \bar{x})^T u_1)^2}_{\text{特征
向量
的方差}} \\
 & = \frac{1}{\|u_1\|^2} \sum_{i=1}^N ((x_i - \bar{x})^T u_1)^2
 \end{aligned}$$

$$\begin{aligned}
 & u_1^T S u_1 \\
 & = \frac{1}{\|u_1\|^2} u_1^T S u_1
 \end{aligned}$$

$$\begin{aligned}
 & \text{s.t. } u_1^T u_1 = 1
 \end{aligned}$$

$$\begin{aligned}
 & u_k = \arg \min_{k \geq 1} \frac{1}{\|u_k\|^2} u_k^T S u_k
 \end{aligned}$$

$$\begin{aligned}
 & \text{s.t. } u_k^T u_k = 1
 \end{aligned}$$

特征的特征值最小化
= 留下的特征
值最大

$\frac{\partial L}{\partial u_1} = 2 S u_1 + \lambda (2 u_1) = 0$

$$\begin{aligned}
 & S u_1 = \lambda u_1 \\
 & \lambda \text{ 为 } S \text{ 的 eigen-vector}
 \end{aligned}$$

λ 为 S 的 eigen-value

SVD for PCA.

$$H^T X = U \Sigma V^T \quad \begin{cases} U^T U = I \\ V^T V = V V^T = I \end{cases}$$

特征向量

$$* \underbrace{S = X^T H X}_{(P \times P)} = X^T H^T H X = V \Sigma U^T V \Sigma U^T = V \Sigma^2 V^T$$

S : 特征分解, 得到方向(主轴), 然后 $(H \cdot V) \rightarrow$ 坐标

T: 特征分解, 直接得到坐标

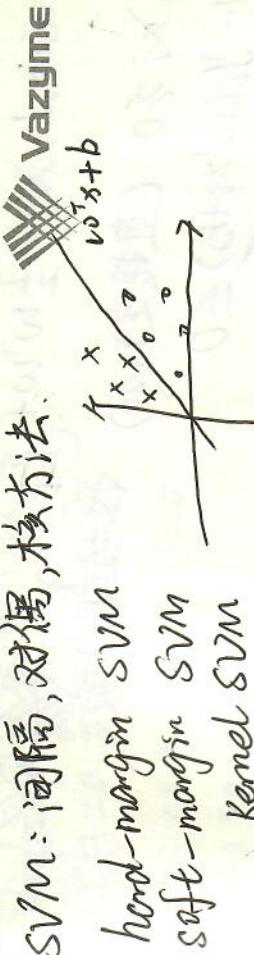
$$H X \cdot V = V \Sigma V^T V = V \Sigma \rightarrow$$

factor

$$X = (x_1, x_2, \dots, x_n)^T = \begin{pmatrix} x_1^T \\ x_2^T \\ \vdots \\ x_n^T \end{pmatrix} = \underbrace{\begin{pmatrix} x_1^T & \dots & x_n^T \end{pmatrix}}_{\text{Sample}} = \underbrace{\begin{pmatrix} x_1^T & \dots & x_n^T \end{pmatrix}}_{\text{特征矩阵}}$$

$H X = \underbrace{U \Sigma}_{\text{坐标矩阵}} \underbrace{V^T}_{\text{坐标方向}}$

(4). SVM



SVM: 间隔, 对偶, 核方法.

$$\begin{array}{l} \text{hard-margin SVM} \\ \text{soft-margin SVM} \\ \text{Kernel SVM} \end{array}$$

点到直线的距离公式:

$$\begin{array}{l} \min_{w, b, x_i} \|w\| \\ \text{s.t. } y_i (w^T x_i + b) \geq 1 \quad \text{for } i=1, 2, \dots, N \end{array}$$

最大间隔分类器

$$\begin{array}{l} \max_{w, b} \text{margin}(w, b) \\ \text{s.t. } y_i (w^T x_i + b) \geq 1 \quad \text{for } i=1, 2, \dots, N \end{array}$$

$$H X = U \Sigma V^T \quad \begin{array}{l} \max_{w, b} \min_{x_i} y_i (w^T x_i + b) \\ \text{s.t. } y_i (w^T x_i + b) \geq 1 \quad \text{for } i=1, 2, \dots, N \end{array}$$

$$\begin{array}{l} \max_{w, b} \min_{x_i} y_i (w^T x_i + b) \\ \text{s.t. } y_i (w^T x_i + b) \geq 1 \quad \text{for } i=1, 2, \dots, N \end{array}$$

$$\Rightarrow \begin{cases} \min_{w, b} \frac{1}{2} \|w\|^2 \\ \text{s.t. } y_i (w^T x_i + b) \geq 1 \quad \text{for } i=1, 2, \dots, N. \end{cases}$$

(对偶优化)
问题

找到对偶问题 dual problem.

$$L(w, b, \lambda) = \frac{1}{2} w^T w + \sum_{i=1}^N \lambda_i (1 - y_i (w^T x_i + b))$$

拉格朗日项

$\star \quad x_i \geq 0 \quad (\text{直接进} \lambda_i \geq 0)$

$1 - y_i (w^T x_i + b) \leq 0$

$\max_w L(w, b, \lambda) = \frac{1}{2} w^T w + 0$

$$\min_{w, b} \frac{1}{2} w^T w = \min_{w, b} \max_{\lambda} L(w, b, \lambda)$$

对偶问题

$$\begin{cases} \text{s.t. } \lambda_i \geq 0 \\ \text{max}_{\lambda} \min_{w, b} [L(w, b, \lambda)] \end{cases}$$

$\Delta \leq 0$

$\text{s.t. } \lambda_i \geq 0$

KKT 条件

$$\begin{cases} \frac{\partial L}{\partial w} = 0 \Rightarrow \frac{\partial L}{\partial b} = 0 \Rightarrow \frac{\partial L}{\partial \lambda} = 0 \\ \lambda_i (1 - y_i (w^T x_i + b)) = 0 \\ \lambda_i \geq 0 \\ 1 - y_i (w^T x_i + b) \leq 0 \end{cases} \rightarrow (\text{只有 } w^T x_i + b = 1 \text{ 且 } \lambda_i \neq 0)$$

$$w^* = \sum_{i=0}^N \lambda_i y_i x_i$$

$$\exists (x_k, y_k), w^T x + b = 1 \Rightarrow 1 - y_k (w^T x_k + b) = 0$$

$$y_k (w^T x_k + b) = 1$$

$$y_k^2 (w^T x_k + b) = y_k \quad (y_k^2 = 1)$$

$$b^* = y_k - w^T x_k = y_k - \sum_{i=0}^N \lambda_i y_i x_i^T x_k$$

$$w^* = \sum_{i=0}^N \lambda_i y_i x_i \quad (\text{data 线性组合})$$

$$b^* = y_k - \sum_{i=0}^N \lambda_i y_i x_i^T x_k$$

定理: $\min \max f(x) \geq \max \min f(x) \quad (\text{弱对偶})$
 $\min \max f(x) = \max \min f(x) \quad (\text{强对偶})$
 这里的条件满足强对偶。

$$f(x) = \text{sign}(w^T x + b^*)$$

$$\frac{\partial L}{\partial b} = 0 \Rightarrow w = \sum_{i=1}^N \lambda_i y_i x_i. \text{ 代入}$$

$$\min_{w, b} L(w, b, \lambda) = -\frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \lambda_i \lambda_j y_i y_j x_i^T x_j + \sum_{i=1}^N \lambda_i$$

Soft margin SVM.



Vazyme

$$\min \frac{1}{2} w^T w + \text{loss}$$

loss 用距离表示

$$\begin{cases} \text{如果 } y_i(w^T x_i + b) \geq 1, & \text{loss} = 0 \\ \text{如果 } y_i(w^T x_i + b) < 1, & \text{loss} = 1 - y_i(w^T x_i + b) \end{cases}$$

$$\text{loss} = \max \{0, 1 - y_i(w^T x_i + b)\}$$

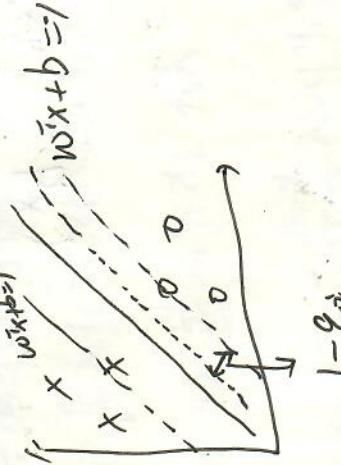
$$\min_{w, b} \frac{1}{2} w^T w + C \sum_{i=1}^N \max \{0, 1 - y_i(w^T x_i + b)\}$$

$$\text{s.t. } y_i(w^T x_i + b) \geq 1 - \varepsilon_i$$

$$\min_{w, b} \frac{1}{2} w^T w + C \sum_{i=1}^N \varepsilon_i$$

$$\begin{cases} \text{s.t. } y_i(w^T x_i + b) \geq 1 - \varepsilon_i \\ \varepsilon_i \geq 0 \end{cases}$$

$$\varepsilon_i = \max(0, 1 - y_i(w^T x_i + b))$$



约束优化问题 (原问题) I

$$\min_{w, b} \text{loss}$$

$$\begin{cases} \text{s.t. } \min_i(x) \leq 0 & i=1 \dots m \\ \eta_j(x) = 0 & j=1 \dots n \end{cases}$$

拉格朗日函数 (原问题拉格朗日化) II

$$L(x, \lambda, \eta) = f(x) + \sum_{i=1}^m \lambda_i \min_i(x) + \sum_{j=1}^n \eta_j \eta_j$$

$$\begin{cases} \min_{x, \lambda, \eta} L(x, \lambda, \eta) \\ \text{s.t. } \lambda_i \geq 0 \end{cases}$$

① 与 ② 等价

如果 x 通过约束 $\min_i(x) > 0$, $\max_i(x) \rightarrow \infty$

如果 x 不通过 $\min_i(x) \leq 0$, $\max_i(x) \neq +\infty$

$$\min_{x} \max_{\lambda} \min_{\eta} L = \min_{x} \{ \max_{\lambda} \{ \max_{\eta} L, +\infty \} \} = \min_{x} L$$

对偶问题: III

$$\begin{cases} \max_{\lambda} \min_{x} L(x, \lambda, \eta) \\ \text{s.t. } \lambda_i \geq 0 \end{cases}$$

对偶问题: III \leq II

$$\begin{cases} \max_{\lambda} \min_{x} L(x, \lambda, \eta) \\ \text{s.t. } \lambda_i \geq 0 \end{cases} \leq \min_{x} \max_{\lambda} L(x, \lambda, \eta)$$

$$\min_{x} L(x, \lambda, \eta) \leq L(x, \lambda, \eta) \leq \max_{\lambda} \min_{x} L(x, \lambda, \eta)$$

对偶性的几何解释



D : 定义域

$$\begin{cases} \min f(x) \\ x \in D \end{cases} \quad D = \text{dom } f \cap \text{dom } m_i \\ \text{s.t. } m_i(x) \leq 0$$

$$d(x, \lambda) = f(x) + \lambda m_i(x)$$

$$p^* = \min_{\lambda} \max_x d(x, \lambda) \quad (\text{原问题最优解})$$

$$\text{集 } G: G = \{ (m_i(x), f(x)) \mid x \in D \}$$

$$= \{ (u, t) \mid x \in D \}$$

$$p^* = \inf_{u \leq 0} \{ t \mid (u, t) \in G \} \quad \text{inf: 下确界}$$

$$d^* = \max_{\lambda} \min_x L(x, \lambda) = \max_{\lambda} \min_x \frac{\min_x (t + \lambda u)}{g(\lambda)} \\ = \max_{\lambda} \frac{g(\lambda)}{g(\lambda)}$$

$$g(\lambda) = \inf_t \{ t + \lambda u \mid (u, t) \in G \}$$

$$\rightarrow G(\text{非凸集}) \quad d^* < p^*$$

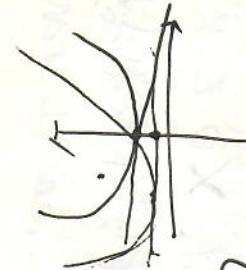
$$p^* = \inf_{u \leq 0} G(u)$$

Slater Condition

$$\exists x \in \text{relint } D \quad (\text{relative interior}) \\ \text{s.t. } \forall i=1, \dots, m \quad m_i(x) < 0$$

对于大多数RB优化, slater成立

$$\begin{array}{c} \text{X} \\ \text{f} \\ \text{d} \\ \text{G} \end{array}$$



对于大多数RB优化, slater成立

$$\begin{array}{c} \text{X} \\ \text{f} \\ \text{d} \\ \text{G} \end{array}$$

$$\text{KKT condition: } \begin{cases} \min_i f_i(x^*) \leq 0 \\ \lambda_i^* \geq 0 \\ \lambda_i^* m_i(x^*) = 0 \end{cases} \quad d^* = \max_{\lambda} \min_i g_i(\lambda, \eta^*) = g(\lambda^*, \eta^*)$$

$$= \min_{\lambda} L(x^*, \lambda^*, \eta^*)$$

$$\leq L(x^*, \lambda^*, \eta^*) \\ = f(x^*) + \frac{1}{2} \sum_{i=1}^m \lambda_i^* m_i(x^*) + \sum_{i=1}^m \lambda_i^* \eta_i^* \\ \leq f(x^*)$$

$$= p^*$$

$$p^* = \inf_{u \leq 0} G(u)$$



Kernel Method. 核方法 (七)



非线性性带来高维转换 (从模型角度)

又保留表示带来的内积 (从优化角度)

PLA (Perceptron Learning Algorithm) \rightarrow 多层感知机
(神经网络)

\rightarrow 深度学习

Cover Theorem: 高维比低维更容易线性可分.

通常说的核函数是正定核函数.

$\forall x, z \in X, K(x, z) = \langle \phi(x), \phi(z) \rangle$ 内积

$K: X \times X \rightarrow \mathbb{R}$ 核函数

正定核: $K: X \times X \rightarrow \mathbb{R}$

$\forall x, z \in X, K(x, z) \geq 0$

如果 $\phi: X \rightarrow \mathbb{R}^n, \phi \in \mathcal{H}$

s.t. $K(x, z) = \langle \phi(x), \phi(z) \rangle, K(x, z)$ 为正定核函数

正定核: $K: X \times X \rightarrow \mathbb{R}, \forall x, z \in X, K(x, z) \geq 0$

如果 $K(x, z)$ 满足如下性质:

① 对称性 ② 正定性

那么称 $K(x, z)$ 为正定核函数.

③ 对称性 $\Leftrightarrow K(x, z) = K(z, x)$

要证: $K(x, z) = \langle \phi(x), \phi(z) \rangle \Leftrightarrow$ Gram matrix 正定.

Hilbert Space:

完备的, 可能是无限维的, 被赋予内积的线性空间.

对极限是封闭的

$\{K_n\} \lim_{n \rightarrow \infty} K_n = K \in \mathcal{H}$

对称性

正定性

线性

$\langle f, g \rangle = \langle g, f \rangle$

$\langle f, f \rangle \geq 0$

$\langle rf_1 + sf_2, g \rangle = r \langle f_1, g \rangle + s \langle f_2, g \rangle$

必要性:

$K(x, z) = \langle \phi(x), \phi(z) \rangle \Rightarrow$ Gram matrix 正定.

对称: $K(x, z) = \langle \phi(x), \phi(z) \rangle = \langle \phi(z), \phi(x) \rangle = K(z, x)$

正定: $\begin{pmatrix} K(x_1, x_1) & \cdots & K(x_1, x_N) \\ \vdots & \ddots & \vdots \\ K(x_N, x_1) & \cdots & K(x_N, x_N) \end{pmatrix}$

$\Rightarrow K \in \mathbb{R}^{N \times N}, K \geq 0$

即 $\forall x \in \mathbb{R}^N, x^T K x \geq 0$

$K = (x_1, x_2, \dots, x_N) \begin{pmatrix} K_{11} & K_{12} & \cdots & K_{1N} \\ K_{21} & \ddots & & \vdots \\ \vdots & & \ddots & \vdots \\ K_{N1} & K_{N2} & \cdots & K_{NN} \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_N \end{pmatrix}$

$$\sum_{i=1}^N \sum_{j=1}^N \partial_i \partial_j K_{ij}$$

$$= \sum_{i=1}^N \sum_{j=1}^N \partial_i \partial_j \langle \phi(x_i), \phi(x_j) \rangle$$

$$= \sum_{i=1}^N \sum_{j=1}^N \partial_i \partial_j \phi(x_i)^T \phi(x_j)$$

$$= \sum_{i=1}^N \partial_i \phi(x_i)^T \cdot \sum_{j=1}^N \partial_j \phi(x_j)$$

$$= \left[\sum_{i=1}^N \partial_i \phi(x_i) \right]^T \cdot \sum_{j=1}^N \partial_j \phi(x_j)$$

$$= \left\langle \sum_{i=1}^N \partial_i \phi(x_i), \sum_{j=1}^N \partial_j \phi(x_j) \right\rangle$$

$$= \left\| \sum_{i=1}^N \partial_i \phi(x_i) \right\|^2 \geq 0.$$



vazyme

(11) 指数族分布.

二项分布

负二项分布

泊松分布

几何分布

指数分布

Beta

Dirichlet

Gamma

Cassian (高斯分布/正态分布)

指数族分布标准型势: $P(x|N) = h(x) \exp(\eta^T \phi(x) - A(\eta))$

N: 参数向量, $x \in \mathbb{R}^p$

A(N) log partition function.

$$A'(\eta) = \log [A(\eta)]$$

对数. 配分函数.

充分统计量: 能够表示分布的统计量 (均值, 方差, ...)

$\phi(x)$ / 后验 似然 先验

共轭: $P(z|x) \propto P(x|z) P(z)$

当似然为指数族分布, 先验的分布形式与后验相同

计算上的方便.

线性模型: $\begin{cases} \text{线性组合 } w^T x \\ \text{link function } \rightarrow \text{“激活函数”} \end{cases}$

极大似然估计与参数估计量.



vazyme

高维随机变量 $P(x_1, x_2, \dots, x_p)$ 边缘概率: $P(x_i)$

$$\begin{aligned} \eta_{MLE} &= \arg \max \log P(\mathbf{D} | \boldsymbol{\eta}) \\ &= \arg \max \sum_{i=1}^N (\eta^T \phi(x_i) - A(\eta)) \end{aligned}$$

$$\begin{aligned} \frac{\partial}{\partial \eta} \left(\sum_{i=1}^N (\eta^T \phi(x_i) - A(\eta)) \right) &= \frac{1}{N} \sum_{i=1}^N \phi'(x_i) \\ &= \sum_{i=1}^N \phi(x_i) - A'(\eta) = 0 \end{aligned}$$

最大熵:

$$\begin{aligned} \text{信息量: } -\log P &= \int -p(x) \cdot \log p(x) dx \\ \text{熵: } H(p) &= \int -p(x) \cdot \log p(x) \\ &= -\frac{1}{N} \sum_i p(x_i) \cdot \log p(x_i) \\ \max H(p) &= \max \left(-\sum_{i=1}^K p_i \log p_i \right) \Rightarrow \begin{cases} \min \sum_{i=1}^K p_i \log p_i \\ \text{s.t. } \sum_{i=1}^K p_i = 1 \end{cases} \\ L(p, \lambda) &= \sum_{i=1}^K (p_i \log p_i) + \lambda \left(1 - \sum_{i=1}^K p_i \right) \end{aligned}$$

$$\begin{aligned} p_i &= \exp(\lambda^{-1}) && \begin{cases} \text{参数分布已知:} \\ f(x) \text{是关于 } \lambda \text{ 的函数} \end{cases} \\ \therefore \hat{p}_1 &= \hat{p}_2 \dots = \hat{p}_K && E_p[f_{\lambda}] = \Delta(\lambda) \\ \therefore p(x) &\text{均匀分布} && \begin{cases} \min \sum_i p_i \log p_i \\ \text{s.t. } \sum_i p_i = 1 \end{cases} \end{aligned}$$

(九) 概率图.



Vazyme

条件概率: $P(x_j | x_i)$

$$\begin{aligned} \text{Sum Rule: } P(x_i) &= \int P(x_1, x_2) dx_2 \\ \text{Product Rule: } P(x_1, x_2) &= P(x_1) P(x_2 | x_1) = P(x_2) P(x_1 | x_2) \\ \text{Chain Rule: } P(x_1, x_2, \dots, x_p) &= \prod_{i=1}^p P(x_i | x_1, x_2, \dots, x_{i-1}) \\ \text{Bayesian Rule: } P(x_i | x_i) &= \frac{P(x_i, x_i)}{P(x_i)} = \frac{P(x_i, x_i)}{\int p(x_i, x_0) dx_0} = \frac{P(x_i)}{\int p(x_i, x_0) dx_0} \end{aligned}$$

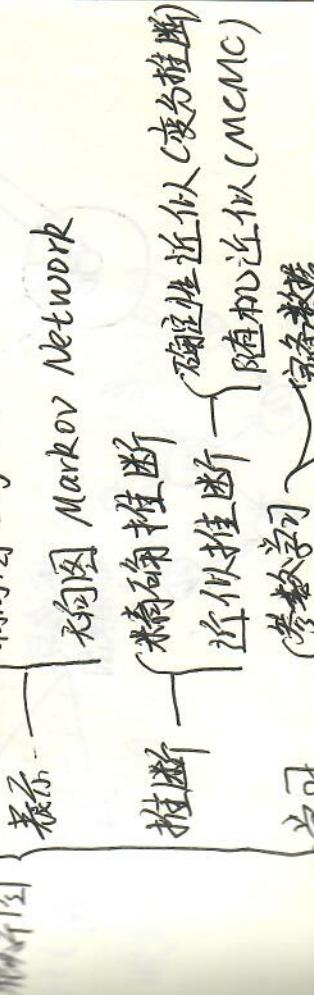
计算复杂度

$$\begin{aligned} \downarrow & \text{相互独立} \Rightarrow P(x_1, x_2, \dots, x_p) = \prod_{i=1}^p P(x_i) \\ \hookrightarrow & \text{Naive Bayes: } P(x_i | y) = \frac{P(y)}{\prod_{j \neq i} P(x_j | y)} \end{aligned}$$

Markov Property: $x_i \perp x_{i+1} | x_1, \dots, x_{i-1}$ (x_i 与 x_{i+1} 条件独立)

独立性 $x_A \perp x_B | x_C$

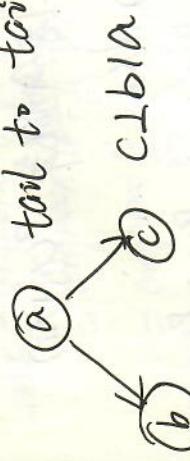
Bayesian Network



Bayesian Network

因子分解 $P(x_1, x_2, \dots, x_p) = \prod_{i=1}^p P(x_i | \text{父节点})$
 若子节点 x_i 的父节点为 x_1, x_2, \dots, x_{i-1}

tail to tail

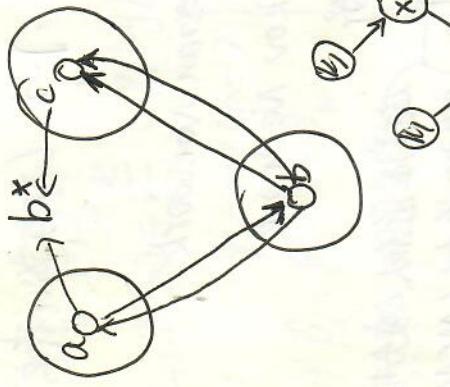


head to tail
 $a \rightarrow b \rightarrow c$ arc $a \rightarrow b$

head to head
 $a \rightarrow b$ arc $a \rightarrow b$

若 C 未被观察, $A \perp B$
 若 C 被观察, $A \perp B$ 不独立.

D-separation



$$P(x_i | x_{-i}) = \frac{P(x_i, x_{-i})}{P(x_{-i})}$$

$$= \frac{P(x_i)}{\int_{x_{-i}} P(x_i, x_{-i}) dx_{-i}}$$

挑出与 x_i 相关的

EM (Expectation Maximum) 期望最大.



Vazyme

$$\text{MLE: } P(x|\theta) \quad \theta_{\text{MLE}} = \underset{\theta}{\operatorname{arg\max}} \log P(x|\theta)$$

$$\theta^{t+1} = \underset{\theta}{\operatorname{arg\max}} \int_{\mathcal{Z}} \log P(x, z|\theta) \cdot P(z|x, \theta^{(t)}) dz$$

$$E_{\mathcal{Z}|x, \theta^{(t)}} [\log P(x, z|\theta)]$$

$$\theta^t \rightarrow \theta^{t+1}$$

$$\log P(x|\theta^t) \leq \log P(x|\theta^{t+1})$$

$$\log P(x|\theta) = \log P(x, z|\theta) - \log P(z|x, \theta)$$

$$\textcircled{1}$$

$$\textcircled{2} = \log P(x|\theta) \int_{\mathcal{Z}} P(z|x, \theta^{(t)}) dz = \log P(x|\theta)$$

$$\textcircled{3} = \int_{\mathcal{Z}} P(z|x, \theta^t) \cdot \log P(x, z|\theta) dz - \int_{\mathcal{Z}} P(z|x, \theta^t) \cdot \log P(z|\theta) dz$$

$$\textcircled{4} \quad \theta^{(t+1)}, \quad \theta^{(t)}$$

$$\textcircled{5}$$

$$H(\theta, \theta^t)$$

$$H(\theta^{t+1}, \theta^t) \leq H(\theta^t, \theta^t)$$

$$H(\theta^{t+1}, \theta^t) - H(\theta^t, \theta^t) \leq 0$$

$$\text{收敛性: } \log P(x|\theta^{t+1}) \geq \log P(x|\theta^t)$$

$$\log P(x|\theta) = \log P(x, z|\theta) - \log P(z|\theta)$$

$$P(x, z) = P(z|x) P(x)$$

$$\log P(x) = \log \frac{P(x, z)}{P(z|x)} = \log P(x, z) - \log P(z|x)$$

$$\begin{aligned} \log P(x|\theta) &= \log P(x, z|\theta) - \log P(z|\theta) \\ &= \log \frac{P(x, z|\theta)}{q(z)} - \log \frac{P(z|x, \theta)}{q(z)} \quad q(z) \neq 0 \\ &\quad \text{左边} = \int_Z q(z) \cdot \log P(x|\theta) dZ = \log P(x|\theta) \cdot \int_Z q(z) dZ = \log P(x|\theta) \end{aligned}$$

$$\text{右边} = \int_Z q(z) \log \frac{P(x, z|\theta)}{q(z)} dZ = \underbrace{\int_Z q(z) \log \frac{P(z|x, \theta)}{q(z)} dZ}_{\text{为拉格朗日量.}}$$

$$\begin{aligned} \text{ELBO} &= \text{Evidence (lower Bound)} \quad KL(q(z) \parallel P(z|x, \theta)) \\ \log P(x|\theta) &= \text{ELBO} + KL(q \parallel P) \quad (\text{KL}(q \parallel P) \geq 0) \\ \log P(x|\theta) &\uparrow \quad \text{ELBO} \uparrow \end{aligned}$$

$$\begin{aligned} \theta &= \arg \max \text{ELBO} = \arg \max_{\theta} \int_Z q(z) \log \frac{P(x, z|\theta)}{q(z)} dZ \\ &= \arg \max_{\theta} \int_Z P(z|x, \theta^t) \log \frac{P(x, z|\theta)}{P(z|x, \theta^t)} dZ \\ &= \arg \max_{\theta} \int_Z P(z|x, \theta^t) \log \frac{P(x, z|\theta)}{P(z|x, \theta^t)} dZ \end{aligned}$$

$$\hat{q} = P(z|x, \theta^t)$$

$$\text{E-step: } P(z|x, \theta^t) \rightarrow E_{z|x, \theta^t} [\log P(x, z|\theta)]$$

$$\text{M-step: } \theta^{t+1} = \arg \max_{\theta} E_{z|x, \theta^t} [\log P(x, z|\theta)]$$

$$\log P(x|\theta) = \log \int_Z P(x, z|\theta) dz = \log \int_Z \frac{P(x, z|\theta)}{q(z)} q(z) dz$$

$$= \log E_{q(z)} \left[\frac{P(x, z|\theta)}{q(z)} \right]$$

$$\geq E_{q(z)} \left[\log \frac{P(x, z|\theta)}{q(z)} \right]$$

$$\text{ELBO} = \underbrace{\log E_{q(z)} \left[\frac{P(x, z|\theta)}{q(z)} \right]}_{\text{ELBO} \cdot \underbrace{P(x|\theta)}_{\text{似然量.}}}$$

$$\text{EM Algorithm: } \theta = \arg \max_{\theta} P(x|\theta) = \arg \max_{\theta} \log P(x|\theta)$$

$$\log P(x|\theta) = \text{ELBO} + KL(q \parallel P)$$

$$\geq \underbrace{\text{ELBO}}_{L(\theta, \theta)} \quad \uparrow$$

$$L(\theta) \text{ 可能取不到 } P(x|\theta)$$

$$\theta \text{ 固定, } \log P(x|\theta) \text{ 固定, } \text{ELBO} \uparrow, KL(q \parallel P) \downarrow$$

$$\text{相当于求 } \arg \min_{\theta} KL(q \parallel P) = \arg \max_{\theta} L(\theta, \theta)$$

$$\text{固定 } \hat{\theta}, \theta = \arg \max_{\theta} L(\hat{\theta}, \theta)$$

$$\begin{aligned}
 L(a, \theta) &= E_Q[\log P(x, z) - \log Q] \\
 &= E_Q[\log P(x, z)] - \underbrace{E_Q[\log Q]}_{H[\theta]}
 \end{aligned}$$

Monte Carlo
Markov Chain and Monte Carlo
(+ =) 变分推断
Variational Inference (变分推断).

版书阅读化问题
(1. 回归: $f(w) = w^T x$ (模型)
loss function: $L(w) = \sum_{i=1}^N \|w^T x_i - y_i\|$
 $\hat{w} = \arg \min_w L(w)$ (参数)
解法: $\frac{\partial L(w)}{\partial w} = 0 \Rightarrow w^* = (X^T X)^{-1} X^T y$
(2) 梯度下降: GD

Inference
SVM
EM

精确推断
近似推断
(准确性降低)
精确推断
近似推断
(准确性降低)
VI
↓
MC, MC
MH, Gibbs

贝叶斯推断
对新样本求 $P(z|X)$.
 $P(z|X) = \int P(z|z) P(z|X)$.

X : observed data

Z : latent variable + parameter.

度数推断为了逼近后验分布 $P(\theta|X)$.

(X, Z) , complete data.

$$\log P(x) = \log P(x, z) - \log P(z|x)$$

$$= \log \frac{P(x, z)}{q(z)} - \log \frac{P(z|x)}{q(z)}$$

$$\int_Z q(z|p(x)) dz = \log P(x) = \int_Z \log \frac{P(x, z)}{q(z)} - \int_Z \log \frac{P(z|x)}{q(z)}$$

$$= \int_Z q(z) \log \frac{P(x, z)}{q(z)} - \int_Z q(z) \log \frac{P(z|x)}{q(z)}$$

KL $Q||P$

$$= L(q) - \underbrace{KL(Q||P)}_{\geq 0}$$

要使 $q(z) \approx P(z|x)$

$$\arg \max_{q(z)} L(q)$$

$$q(z) = \arg \max_{q(z)} L(q) \Rightarrow q(z) \approx P(z|x)$$

构成 M 个相互独立的组.



$$q(z) = \prod_{i=1}^M q_i(z_i)$$

$$\begin{aligned} L(q) &= \int_Z q(z) \log \frac{P(x, z)}{q(z)} \\ &= \int_Z q(z) \log P(x, z) - \int_Z q(z) \log q(z) \\ &\quad \text{①} \quad \text{②} \\ ① &= \int_Z^M q_i(z_i) \log (P(x, z_i) d z_i \dots d z_m \\ &= \int_{Z_j} q_j(z_j) \left[\prod_{i \neq j}^M q_i(z_i) \log P(x, z_i) d z_i \dots d z_m \right] d z_j \\ &\quad \text{③} \quad \text{④} \\ &= \int_{Z_j} q_j(z_j) \cdot \prod_{i \neq j}^M q_i(z_i) d z_j \\ &\quad \text{⑤} \\ ② &= \int_{Z_j} q_j(z_j) \cdot \underbrace{\prod_{i \neq j}^M q_i(z_i)}_{(z_i \neq z_j)} d z_j \\ &\quad \text{⑥} \end{aligned}$$

$$\int_Z q(z) \log q(z) d z = \sum_{j=1}^M \int_{Z_j} q_j(z_j) \log q_j(z_j) d z_j$$

$$\begin{aligned} &= \int_{Z_j} q_j(z_j) \log q_j(z_j) d z_j + C \\ ① - ② &= \int_{Z_j} q_j(z_j) \cdot \log \frac{P(x, z_j)}{q_j(z_j)} d z_j \end{aligned}$$

(十三) MCMC (蒙特卡洛法, 为了近似后验)

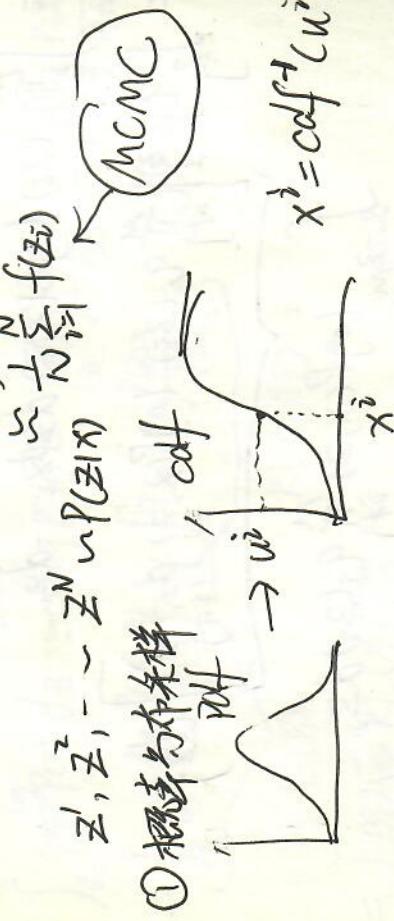
Vazyme

推断后验 $P(\theta|x)$
近似推断 $\begin{cases} \text{确定性} \\ \text{近似} \end{cases} \rightarrow VI$

随机性 $\rightarrow MCMC$

Monte Carlo method: 基于采样的随机近似法。
 \rightarrow observed data

$$P(z|x) \rightarrow E_{z|x}[f(z)] = \int p(z|x) \cdot f(z) dz$$



① 根据样本
p(z|x)



$$x^i = \text{off}^{-1}(u^i)$$

② Rejection Sampling
 $q(z)$ proposed distribution.



$\forall z_i, Mq(z_i) \geq p(z_i)$
取样, 附录部分接收, \exists : 接收率 $\alpha = \frac{p(z^*)}{Mq(z^*)}, 0 \leq \alpha \leq 1$

③ Importance Sampling

Vazyme

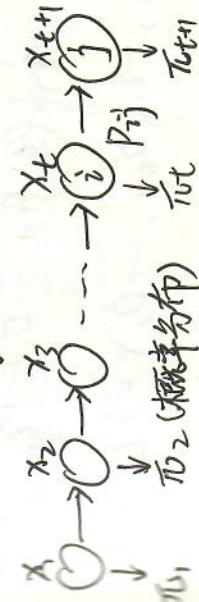
$$\begin{aligned} E_{p(z)}[f(z)] &= \int_z p(z) f(z) dz = \int \frac{p(z)}{q(z)} \cdot q(z) f(z) dz \\ &= \int f(z) \cdot \frac{p(z)}{q(z)} \cdot q(z) dz \\ &\approx \frac{1}{N} \sum_{i=1}^N f(z^i) \frac{p(z^i)}{q(z^i)} \quad z_i \sim q(z) \quad i=1, 2, \dots, N. \end{aligned}$$

weight

- $P_t(z|x)$ Markov Chain $P(x_{t+1} = x | x_1, x_2, \dots, x_t)$

$P \rightarrow$ 转移矩阵 $[P_{ij}]$

$P_{ij} = P(x_{t+1} = j | x_t = i)$



$$\pi_{tn}(x^*) = \int \pi_t(x) \cdot P(x \rightarrow x^*) dx$$

$\pi_t = [\pi_t^{(1)}, \pi_t^{(2)}, \dots, \pi_t^{(N)}]$

$$\sum_{i=1}^N \pi_t(i) = 1$$

维分布: $\pi_t^{(i)} \sim q(z^i)$

(+3) HMM 隐马尔可夫模型



$\lambda = (\pi, A, B)$ 状态转移
初态 prob dist
观测 prob dist

观测量: $o_1, o_2, \dots, o_t \rightarrow V = \{v_1, v_2, \dots, v_m\}$ (观测量的值域)
隐状态变量: $i_1, i_2, \dots, i_t \rightarrow Q = \{q_1, q_2, \dots, q_n\}$ (隐状态的值域)

$A = [a_{ij}]$, $a_{ij} = P(i_{t+1} = q_j | i_t = q_i)$ 状态转移矩阵.

$B = [b_j^{(k)}]$
 $b_j^{(k)} = P(o_t = v_k | i_t = q_j)$ 发射矩阵, 状态
发射矩阵

两个假设: ① 状态 Markov 假设. ② 观察独立.

$$P(i_{t+1} | i_t, i_{t+1}, \dots, i_1, o_t, o_{t-1}, \dots, o_1) = P(i_{t+1} | i_t)$$

$$P(o_t | i_t, i_{t-1}, \dots, i_1, o_{t-1}, \dots, o_1) = P(o_t | i_t)$$

三个问题:

① Evaluation $\rightarrow P(\lambda)$ 前向后向

② Learning $\rightarrow EM$

③ Decoding $\lambda = \arg \max P(\lambda | O)$

预测 $\rightarrow P(i_{t+1} | o_1, o_2, \dots, o_t)$

Evaluation 问题: Given λ , 求 $P(O)$

$$P(\lambda) = \sum_I P(I, \lambda) = \sum_I P(O, I, \lambda) \cdot P(I | \lambda)$$

$$P(I | \lambda) = P(i_1, i_2, \dots, i_t | \lambda) = P(i_1 | \lambda, i_2, \dots, i_{t-1}, \lambda) \cdot P(i_2, i_3, \dots, i_t, \lambda)$$

$$= \pi(a_{i1}) \prod_{t=2}^T a_{it-1, it} = \underbrace{P(i_1 | \lambda)}_{a_{i1, it}}$$

$$P(O | I, \lambda) = \prod_{t=1}^T b_{it}(o_t)$$

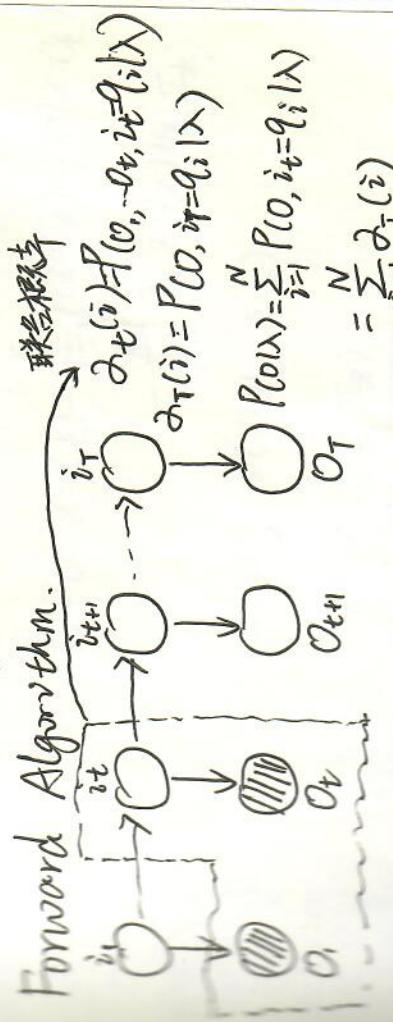
$$P(O | \lambda) = \sum_I \pi(a_{i1}) \prod_{t=2}^T a_{it-1, it} \prod_{t=1}^T b_{it}(o_t)$$

$$= \sum_I \prod_{i=1}^T \pi(a_{ii}) \prod_{t=2}^T a_{it-1, it} \prod_{t=1}^T b_{it}(o_t)$$

$$= \sum_I \prod_{i=1}^T \pi(a_{ii}) \prod_{t=2}^T a_{it-1, it} \prod_{t=1}^T b_{it}(o_t)$$

$$= \sum_I P(O, i_t = q_i | \lambda)$$

$$= \sum_i P(O, i_t = q_i | \lambda)$$



$$= \sum_i \frac{1}{N}$$

$$\alpha_{t+1}(j) = P(c_0, \dots, c_t, 0_{t+1}, \dots, 0_{t+j-1}) \quad \text{Vazyme}$$

$$= \sum_{j=1}^N P(o_j, \sim, \theta_t, \theta_{t+1}, \gamma_{t+1} = \theta_j) / \pi_j$$

$$= \sum_{t=1}^T P(O_{t+1}) \theta_t - \theta^* = \theta^* - \theta_t = \theta_t$$

$$= \sum_{i=1}^N P(i_{t+1} | \text{obs}) \cdot P(0, \dots, 0_t, i_t = q_i, i_{t+1} = q_{i+1} | \text{obs})$$

$$= \sum_{j=1}^N P(O_{t+1} | \theta_t) \cdot P(q_j | O_1, \dots, O_t, \theta_t, \lambda)$$

$$= \prod_{j=1}^J P(a_{t+1} | \underbrace{a_j}_{\text{fixed}}) \cdot P(i_{t+1} = a_j | i_t = a_j, \lambda) \cdot \mathcal{D}_t(i_j)$$

$$= \prod_{j=1}^J b_j(a_{t+1}) \cdot \alpha_j \cdot \mathcal{D}_t(i_j)$$

$$V(k) = \frac{1}{k} \left(\frac{p_1}{p_2} - 1 \right)$$

$$P(x) = \sum_z P(x, z) = \sum_{k=1}^K P(x, z=k) = \sum_{k=1}^K P(z=k) \cdot P(x|z=k)$$

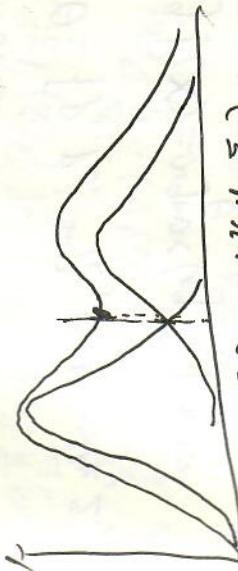
$$= \sum_{k=1}^K p_k \cdot N(x | \mu_k, \Sigma_k)$$

$$P(x, z) = P(z) \cdot P(x|z) = P_z \cdot N(x|\mu_z, \Sigma_z)$$

$$P(\mathcal{Z}|X) = \frac{P(X, \mathcal{Z})}{P(X)} = \frac{P_{\mathcal{Z}} \cdot N(X|\mu_{\mathcal{Z}}, \Sigma_{\mathcal{Z}})}{\sum_{\mathcal{Z}'} P_{\mathcal{Z}'} N(X|\mu_{\mathcal{Z}'}, \Sigma_{\mathcal{Z}'})}$$

GMM: Gaussian Mixture Model (高斯混合模型)
从几何角度看: 加权平均 \rightarrow (多个高斯分布叠加而成)


从几何角度来看：加权平均 \rightarrow (多个高斯分布叠加而成)



$$\mathcal{N}(\mu_1, \Sigma_1) \quad \mathcal{N}(\mu_2, \Sigma_2)$$

$$P(x) = \sum_{k=1}^K \partial_k N(\mu_k, \Sigma_k), \quad \sum_{k=1}^K \partial_k = 1$$

从混合模型角度来看：（生成模型）
 $R=1$ 生成N个样本。

X: observed variable
Latent variable. (对应的样本属于哪一个高斯分布)

卷之三

1. 传统随机变量

卷之三

$$B = \frac{B_0}{2} \left(1 - \frac{2}{\pi} \operatorname{erf} \left(\frac{x}{\sqrt{2}} \right) \right)$$

$$v = k_2 - k_1$$

$$P(X) = \sum_{i=1}^k P(X=i) = \sum_{i=1}^k P(X=i) \cdot P(X=i)$$

$$\frac{1}{2} \left(\sum_{i=1}^n a_i^2 - \sum_{i=1}^n b_i^2 \right) = \frac{1}{2} \left(\sum_{i=1}^n (a_i - b_i)^2 \right)$$

$$= K \mathcal{P} \mathcal{V}(\mathbf{z}) \mathcal{W} =$$

$$= \sum_{k=1}^{\infty} \frac{1}{k} \cdot N(k) \left(\frac{1}{k}, \frac{1}{k} \right)$$

$$P(X, Y) = P(Z) \cdot P(X|Z) = P_{\text{prior}} \cdot P(X|Y \leq)$$

$$(\omega_1, \omega_2, \omega_3) = (\omega_1, -\omega_2, \omega_3)$$

$$\beta(x, z) \beta_2(x) \beta_1(x) \beta_0(x) \leq 1$$

X : Observed data $\rightarrow X = (X_1, X_2, \dots, X_N)$

(X, Z) : complete data

θ : parameter $\rightarrow \theta = \{P_1, P_2, \dots, P_K, \mu_1, \mu_2, \dots, \mu_K, \Sigma_1, \Sigma_2, \dots, \Sigma_K\}$

$$\theta_{MLE} = \arg \max_{\theta} \log P(X) = \arg \max_{\theta} \log \prod_{i=1}^N P(X_i, Z_i)$$

$$\begin{aligned} &= \arg \max_{\theta} \sum_{i=1}^N \log P(X_i, Z_i) \\ &= \arg \max_{\theta} \sum_{i=1}^N \log \underbrace{\sum_{k=1}^K P_k N(X_i | \mu_k, \Sigma_k)}_{\text{由加高斯分布}} \end{aligned}$$

无法直接得到解析解，所以用EM算法求近似解。

使用EM算法。

$$\theta^{t+1} = \arg \max_{\theta} \underbrace{E_{Z|X, \theta^t} [\log P(X, Z | \theta)]}_{Q(\theta, \theta^t)}$$

$$= \sum_{i=1}^N \log P(X_i, Z_i | \theta)$$

$$= \sum_{i=1}^N \log P(X_i, Z_i | \theta) \cdot P(Z_i | X_i, \theta^t)$$

$$= \sum_{i=1}^N \log P(X_i, Z_i | \theta)$$

$$Q(\theta, \theta^t) = \int_Z \log P(X, Z | \theta) P(Z | X, \theta^t) dZ$$

Vaizyme

$$= \sum_{i=1}^N \log \prod_{j=1}^N P(X_{ij}, Z_{ij} | \theta) \prod_{j=1}^N P(Z_{ij} | X_{ij}, \theta^t)$$

$$= \sum_{i=1}^N \sum_{j=1}^N \log P(X_{ij}, Z_{ij} | \theta) \cdot \prod_{j=1}^N P(Z_{ij} | X_{ij}, \theta^t)$$

2: 第一个样本属于各个分布的概率

$$= \sum_{i=1, Z_{i1} \sim Z_{iN}} \left[\log P(X_i, Z_i | \theta) + \log P(X_i, Z_2 | \theta) + \dots + \log P(X_i, Z_N | \theta) \right] \prod_{j=1}^N P(Z_{ij} | X_{ij}, \theta^t)$$

$$\begin{aligned} &= \sum_{i=1, Z_{i1} \sim Z_{iN}} \log P(X_i, Z_i | \theta) \prod_{j=1}^N P(Z_{ij} | X_{ij}, \theta^t) \\ &= \sum_{i=1, Z_{i1} \sim Z_{iN}} \underbrace{\log P(X_i, Z_i | \theta)}_{\text{由加高斯分布}} \underbrace{\prod_{j=1}^N P(Z_{ij} | X_{ij}, \theta^t)}_{\text{由加高斯分布}} \end{aligned}$$

$$= \sum_{i=1}^N \log P(X_i, Z_i | \theta) \cdot P(Z_i | X_i, \theta^t) + \dots + \sum_{i=1}^N \log P(X_i, Z_i | \theta) \cdot P(Z_N | X_N, \theta^t)$$

$$= \sum_{i=1}^N \log P(X_i, Z_i | \theta)$$

$$= \sum_{i=1}^N \log P(X_i, Z_i | \theta)$$

E-step:

$$\begin{aligned}
 Q(\theta, \theta^t) &= \int_Z \log P(x, z | \theta) P(z | x, \theta^t) dz \\
 &= \sum_{i=1}^N \sum_{z_i} \log p_{z_i} \cdot N(x_i | \mu_{z_i}, \Sigma_{z_i}) \frac{p_{z_i} \cdot N(x_i | \mu_{z_i}^t, \Sigma_{z_i}^t)}{\sum_{z_i} p_{z_i}^t \cdot N(x_i | \mu_{z_i}^t, \Sigma_{z_i}^t)} \\
 &= \sum_{j=1}^N \sum_{z_i} \log [p_{z_i} \cdot N(x_i | \mu_{z_i}, \Sigma_{z_i})] P(z_i | x_i, \theta^t) \\
 &= \sum_{z_i} \sum_{j=1}^N \log [p_{z_i} \cdot N(x_i | \mu_{z_i}, \Sigma_{z_i})] P(z_i | x_i, \theta^t) \\
 &= \sum_{k=1}^K \sum_{j=1}^N \log [p_k \cdot N(x_i | \mu_k, \Sigma_k)] \cdot P(z_i = c_k | x_i, \theta^t) \\
 &= \sum_{k=1}^K \sum_{j=1}^N [\log p_k + \log N(x_i | \mu_k, \Sigma_k)] \cdot P(z_i = c_k | x_i, \theta^t)
 \end{aligned}$$

M-step:

$$\begin{aligned}
 p^{t+1} &= (p_1^{t+1}, p_2^{t+1}, \dots, p_K^{t+1}) \\
 \max_p & \sum_{k=1}^K \sum_{j=1}^N \log p_k \cdot P(z_i = c_k | x_i, \theta^t) \\
 \text{s.t.} \quad & \sum_{k=1}^K p_k = 1
 \end{aligned}$$

$$L(p, \lambda) = \sum_{k=1}^K \sum_{j=1}^N \log p_k \cdot P(z_i = c_k | x_i, \theta^t) + \lambda \left(\sum_{k=1}^K p_k - 1 \right)$$

Vazyme

$$\sum_{i=1}^N P(z_i = c_k | x_i, \theta^t) + p_k \cdot \lambda = 0$$

Vazyme

$$\begin{aligned}
 \sum_{i=1}^N \sum_{k=1}^K P(z_i = c_k | x_i, \theta^t) + \sum_{k=1}^K p_k \cdot \lambda &= 0 \\
 \sum_{i=1}^N \sum_{k=1}^K P(z_i = c_k | x_i, \theta^t) &= -\lambda \\
 \sum_{i=1}^N P(z_i = c_k | x_i, \theta^t) &= -\lambda \\
 \lambda &= -N \\
 P(z_i = c_k | x_i, \theta^t) &= -N
 \end{aligned}$$

$$\theta = \{p_1, p_2, \dots, p_K, \mu_1, \mu_2, \dots, \mu_K, \Sigma_1, \Sigma_2, \dots, \Sigma_K\}$$

(三十一) 生成模型 (关注样本分布本身) $P(x)$ ~~Vazyme~~

监督) 概率模型 $P(x|y)$: ~~离散~~ \rightarrow ~~连续~~ \rightarrow ~~最大熵~~
生成模型

非概率模型: PLA, SVM, LDA, KNN, NN

Tree Model

非监督) 概率模型: 生成模型

非概率模型: PCA, K-means, Auto-encoder

生成模型:

Naive Bayes 假设: $P(x|y) = \prod_{i=1}^N P(x_i|y)$

Mixture model: GMM Kalman Particle

Time-series Model: HMM | KF | PF

非参数贝叶斯模型: GP / DP

(三十二) 可以离散或连续 ~~Vazyme~~

监督) 离散 \rightarrow 形状 \rightarrow ~~参数~~ \rightarrow ~~无向~~ \rightarrow ~~无向~~ \rightarrow ~~无向~~
生成模型 \rightarrow ~~无向~~ \rightarrow ~~无向~~ \rightarrow ~~无向~~ \rightarrow ~~无向~~

种类: $\left\{ \begin{array}{l} \text{Discrete vs. Continuous} \\ \text{Directed Model vs. Undirected} \\ \text{Latent Variable vs. Fully-observed} \\ \text{shallow vs. Deep (多层)} \\ \text{Sparse vs. Dense (连接)} \end{array} \right\}$

参数: $\left\{ \begin{array}{l} \text{parameteric vs. non-parameteric} \\ \text{Implicit Density vs. Explicit Density} \\ \text{GAN} \end{array} \right\}$

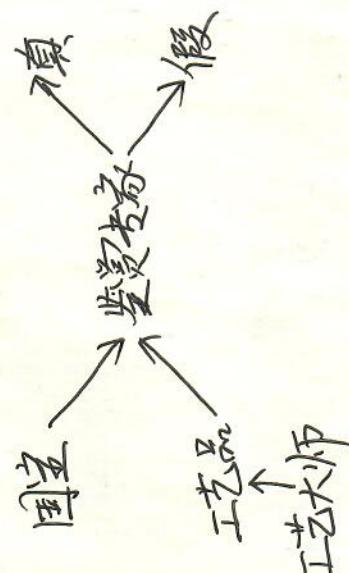
维. MLE: tractable vs. intractable

参数: MLE: maximum likelihood - based
vs. likelihood - free Model.

(三十一) Generative Adversarial Network. 生成对抗网络.



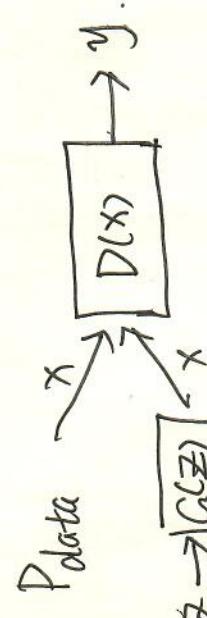
Vazyme



目标: $D_{\text{data}}(x) = \{x_i\}_{i=1}^N$
 $D_{\text{gen}}^* = D_g(x, \theta_g) : \text{generator} (D_g(z) + G(z, \theta_g))$
 生成数据: $y|x$ discriminator

$$\begin{array}{c|cc} y|x & 1 & 0 \\ \hline P & D(x) & 1 - D(x) \end{array}$$

目标: ① 生成数据鉴别能力强 (手段)
 ② 工艺大师(生成器)造能力强 (目标)



$$\min_{G} \max_{D} E_{x \sim p_{\text{data}}} [\log D(x)] + E_{x \sim p_g} [1 - \log D(x)]$$

$$\text{记 } V(D, G) = E_{x \sim p_{\text{data}}} [\log D(x)] + E_{x \sim p_g} [1 - \log D(x)]$$

① for fixed G , 求 $\max_{D} V(D, G)$

$$\max_{D} V(D, G) = \int_{\text{data}} \log D(x) dx + \int_{\text{gen}} \log(1 - \log D(x)) dx$$

$$= \int [P_d \log D(x) + P_g (1 - \log D(x))] dx$$

$$\frac{\partial}{\partial D} = \int \frac{\partial}{\partial D} [P_d \log D(x) + P_g (1 - \log D(x))] dx$$

$$= \int \left[P_d \frac{1}{D} + P_g \cdot \frac{-1}{1-D} \right] dx = 0$$

$$\Rightarrow D_g^* = \frac{P_d}{P_d + P_g}$$

② 将 D_g^* 代入 $V(D, G)$

$$\min_{G} \max_{D} V(D, G) = \min_{G} V(D_g^*, G)$$



(= + =) 谱聚类 Spectral Clustering

$$\min_{G} \mathcal{J}(D_G^*, C_1)$$

$$= \min_G \mathbb{E}_{P_d} \left[\log \frac{P_d}{P_d + P_g} \right] + \mathbb{E}_{P_g} \left[\log \frac{P_g}{P_d + P_g} \right]$$

$$= \min_G \mathbb{E}_{P_d} \left[\log \left(\frac{P_d}{(P_d + P_g)/2} \cdot \frac{1}{2} \right) \right] +$$

$$\mathbb{E}_{P_g} \left[\log \left(\frac{P_g}{(P_d + P_g)/2} \cdot \frac{1}{2} \right) \right]$$

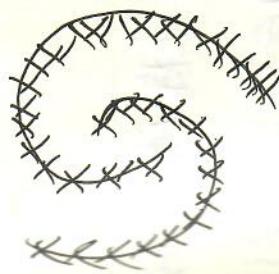
$$= \min_G KL \left(P_d \parallel \frac{P_d + P_g}{2} \right) + KL \left(P_g \parallel \frac{P_d + P_g}{2} \right) - \log 4$$

$\geq \log 4$
 当 $P_d = \frac{P_d + P_g}{2} = P_g$ 时, 等号成立

$$P_g^* = P_d, D_g^* = \frac{P_d}{P_d + P_g} = \frac{1}{2}$$



\Rightarrow compactness:
 K-means Clustering



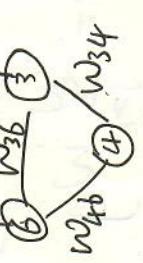
\Rightarrow Connectivity:
 spectral clustering.

Graph-based 谱聚类的无向图

$$G = \{V, E\}$$

$$V = \{1, 2, \dots, N\} \text{ (节点集)}$$

一个节点一个样本

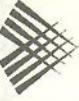


$$W = [w_{ij}] \quad 1 \leq i, j \leq N$$

(相似度矩阵)
 (相似度集)

$$(W) \text{ Similarity matrix (affinity matrix)}$$

$$W_{ij} = \exp \left(- \frac{\|x_i - x_j\|^2}{2\sigma^2} \right), (i, j) \in E$$

定理: $A \subset V, B \subset V, A \cap B = \emptyset$ 

$$w(A, B) = \sum_{i \in A, j \in B} w_{ij}$$

$w(A, B)$ 是有相连的节点的和, 可以看做是
将集多切割以后的损失

假如有 K 个类别

$$\text{cut}(V) = \text{cut}(A_1, A_2, \dots, A_K) \quad (\text{将 } V \text{ 分成 } K \text{ 个子集})$$

$$V = \bigcup_{k=1}^K A_k \quad A_i \cap A_j = \emptyset, \forall i, j \in \{1, 2, \dots, K\}$$

$$\text{cut}(V) = \sum_{k=1}^K w(A_k, \bar{A}_k) = \sum_{k=1}^K w(A_k, V) - w(A_k, A_k)$$

$$\text{目标: } \min \text{Cut}(V) \quad \{A_k\}_{k=1}^K$$

$$\text{Cut} = \sum_{k=1}^K \frac{w(A_k, \bar{A}_k)}{\Delta} \quad \Delta = \text{degree}(A_k)$$

$$w_{ij} = \sum_{j \geq 1} w_{ij} \quad d_i = \sum_{j \in A_k} d_i$$

$$\text{Cut} = \sum_{k=1}^K \frac{w(A_k, \bar{A}_k)}{\sum_{i \in A_k} d_i}$$

Model: 

$$\{\hat{A}_k\}_{k=1}^K = \arg \min \text{Cut}(V)$$

indicator vector:

$$\{y_i \in \{0, 1\}^K \mid y_{ij} = 1 \iff \text{第 } i \text{ 个样本属于 } j \text{ 个类别}\}$$

$$\sum_{j=1}^K y_{ij} = 1 \quad \begin{cases} 1 \leq i \leq N \\ 1 \leq j \leq K \end{cases}$$

$$y = (y_1, y_2, \dots, y_N)^T \quad N \times K$$

$$y = \arg \min_y \underbrace{\sum_{i=1}^N w(A_k, \bar{A}_k)}_{\text{Cut}(V)} \quad \underbrace{\sum_{k=1}^K \frac{w(A_k, \bar{A}_k)}{\sum_{i \in A_k} d_i}}$$

$$= \text{tr} \begin{pmatrix} w(A, \bar{A}) \\ \sum_{i \in A} d_i \end{pmatrix} \begin{pmatrix} 0 \\ \frac{w(A_k, \bar{A}_k)}{\sum_{i \in A_k} d_i} \end{pmatrix}$$

$$= \left(\begin{pmatrix} w(A, \bar{A}) \\ \sum_{i \in A} d_i \end{pmatrix} \cdot \begin{pmatrix} \sum_{i \in A} d_i \\ w(A_k, \bar{A}_k) \end{pmatrix} \right)^{-1}$$

已知 α, γ, β 求 θ, ρ .

$$\tilde{Y}^T = \underbrace{(y_1, y_2, \dots, y_N)}_{K \times K} \left(\underbrace{\begin{pmatrix} y_1^T \\ y_2^T \\ \vdots \\ y_N^T \end{pmatrix}}_{N \times 1} \right) = \sum_{j=1}^N y_j y_j^T$$

$$= \begin{pmatrix} N_1 & & \\ & \ddots & \\ & & N_K \end{pmatrix}_{K \times K} = \begin{pmatrix} \sum_{i \in A_1} 1 & & \\ & \ddots & \\ & & \sum_{i \in A_K} 1 \end{pmatrix}_{K \times K}$$

N_k : 在 N 个样本中, 属于 K 的样本个数.

$$\sum_{k=1}^K N_k = N, \quad N_k = |A_k| = \sum_{j \in A_k} 1$$

$$\sum_{i=1}^N y_i y_i^T d_i = \sum_{i=1}^N y_i d_i y_i^T$$

$$P = \left(\sum_{i \in A_1} d_i, \sum_{i \in A_2} d_i \right)$$

$$= Y^1 \cdot D \cdot Y$$

$$D = \begin{pmatrix} d_1 & & \\ & \ddots & \\ & & d_N \end{pmatrix} = \text{diag}(w_0, 1_N)$$

$$0 = \nabla_{\bar{D}}^T \bar{y} - \bar{y}_m$$

$$\begin{aligned}
 \bar{Y}_i^T \bar{W} \bar{Y} &= (\bar{y}_1 - \bar{y}_N) \begin{pmatrix} w_{11} & \dots & w_{1N} \\ \vdots & \ddots & \vdots \\ w_{N1} & \dots & w_{NN} \end{pmatrix} \begin{pmatrix} \bar{y}_1^T \\ \vdots \\ \bar{y}_N^T \end{pmatrix} \\
 &= \left(\sum_{i=1}^N y_i w_{i1} - \sum_{j=1}^N y_j w_{jN} \right) \begin{pmatrix} \bar{y}_1^T \\ \vdots \\ \bar{y}_N^T \end{pmatrix}
 \end{aligned}$$

Vazyme

$$D = \left(\begin{array}{cc} w(A_1, \bar{A}_1) & w(A_k, \bar{A}_k) \\ \sum_{i \neq k} d_i & \sum_{i \neq k} \sum_{j \neq i} w_{ij} \end{array} \right)$$

$$w(A_1, \bar{A}_1) = w(A_1, v) - w(A_1, A_1)$$

$$D = \left(\begin{array}{cc} w(A_1, A_1) & - \left(\begin{array}{c} \sum_{i \neq k} d_i \\ \sum_{j \neq k} d_j \end{array} \right) \\ - \left(\begin{array}{c} \sum_{i \neq k} d_i \\ \sum_{j \neq k} d_j \end{array} \right) & w(A_k, A_k) \end{array} \right)$$

Vazyme

$$D = \left(\begin{array}{cc} w(A_1, \bar{A}_1) & w(A_k, \bar{A}_k) \\ \sum_{i \neq k} d_i & \sum_{i \neq k} \sum_{j \neq i} w_{ij} \end{array} \right)$$

$$w(A_1, \bar{A}_1) = w(A_1, v) - w(A_1, A_1)$$

$$D = \left(\begin{array}{cc} w(A_1, A_1) & - \left(\begin{array}{c} \sum_{i \neq k} d_i \\ \sum_{j \neq k} d_j \end{array} \right) \\ - \left(\begin{array}{c} \sum_{i \neq k} d_i \\ \sum_{j \neq k} d_j \end{array} \right) & w(A_k, A_k) \end{array} \right)$$

$$N_{\text{cut}}(v) = \text{tr}(O \cdot P) \\ = \text{tr}(O' \cdot P)$$



(二十一) VAE
Variational Auto-encoder.



Latent Variable Model

$$\hat{z} \rightarrow \hat{x}$$

GMN: RT Gaussian Dist: $z \sim \text{Categorical Dist}$
VAE: infinite Gaussian Dist: $z \sim N(0, I)$

$x|z \sim N(\mu_\theta(z), \Sigma_\theta(z))$

$$p_{\text{in}} = \int_z p(x, z) dz = \int_z p(z) \cdot p(x|z) dz$$

intractable

$$p_{\theta}(z|x) = \frac{p(z) \cdot p(x|z)}{p(x)}$$

intractable

$$p(z) = N(0, I)$$

$p_\theta(x|z) = N(\mu_\theta(z), \Sigma_\theta(z))$
 $p_\theta(z|x)$ is intractable.

$$p_\theta(x) = \int_z p_\theta(x|z) p_\theta(z) dz$$

$$\log p_{\theta}(x) = ELBO + KL(p_\theta(z|x) || p_0(z|x))$$

$$\text{E-step: } \hat{z} = p_\theta(z|x) \quad KL=0, E=ELBO \\ \text{M-step: } \theta = \arg \max E_{\text{LBO}}$$

$$\langle \hat{\theta}, \hat{\phi} \rangle = \arg \min_{\theta, \phi} \text{KL}(q_{\phi}(z|x) \| p_{\theta}(z|x))$$

($z + z$) Normalizing Flow 標榜  Vazyme

$= \arg \max \text{ELBO}$

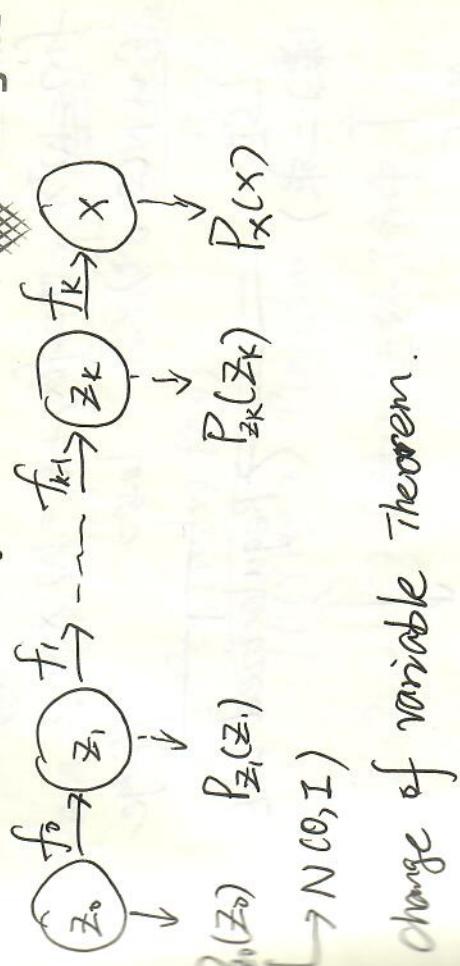
$$= \arg \max \mathbb{E}_{q_{\phi}(z|x)} [\log p_{\theta}(x, z)] - H[q_{\phi}]$$

$$= \arg \max \mathbb{E}_{q_{\phi}(z|x)} [\log p_{\theta}(x|z) + \log p_{\phi}(z)] - H[\phi]$$

$$= \arg \max \mathbb{E}_{q_{\phi}(z|x)} [\log p_{\theta}(x|z)] + \mathbb{E}_{q_{\phi}(z|x)} \log p_{\theta}(z)$$

$$= \arg \max \mathbb{E}_{q_{\phi}(z|x)} [\log p_{\theta}(x|z)] - \text{KL}(q_{\phi}(z|x) \| p_{\theta}(z))$$

$\hat{\phi}$ 痘最大



change of variable theorem.

$$\begin{aligned} x &= f(z), \quad z, x \in \mathbb{R}^P \\ z &\sim p_z(z), \quad x \sim p_x(x) \end{aligned}$$

f is continuous, invertible.

$$\int_z p_z(z) dz = 1 = \int_x p_x(x) dx \Rightarrow \left| \frac{dp_z(z)}{dz} dz \right| = \left| \frac{dp_x(x)}{dx} dx \right|$$

$$p_x(x) = \left| \frac{dp_x(x)}{dx} \right| \cdot p_z(z)$$

$x = f(z)$, f is invertible

$$p_x(x) = \left| \frac{\partial f^{-1}(x)}{\partial x} \right| \cdot p_z(z) \frac{\partial f(x)}{\partial x}$$

Jacobian Matrix

$$p_x(x) = \left| \det \left(\frac{\partial f^{-1}(x)}{\partial x} \right) \right| \cdot p_z(z) \det \left(\frac{\partial f(x)}{\partial x} \right)$$

$$p_x(x) = \left| \det \left(\frac{\partial f^{-1}(x)}{\partial x} \right) \right| \cdot p_z(z) \left| \frac{\partial f(x)}{\partial x} \right|$$

(+ 7c) Ridge 算法

$$f(x) = w^T x, \quad y = f(x) + \varepsilon = w^T x + \varepsilon$$

$\varepsilon \sim N(0, \sigma^2)$

lasso
ridge

LSE
(最小二乘)

MLE
noise is Gaussian

MAP
noise is Gaussian

$$w_{MLE} = \arg \max_w P(w | \text{Data})$$

$$P(w) = N(w, \Sigma_p)$$

$$P(w | \text{Data}) \propto P(w | \text{Data}) P(w)$$

$$w_{MAP} = \arg \max_w P(w | \text{Data})$$

$$= \arg \max_w P(\text{Data} | w) P(w)$$

$$\propto \prod_{i=1}^N N(y_i | w^T x_i, \sigma^2) \cdot N(w, \Sigma_p)$$

Inference:

$$P(w | \text{Data}) = P(w | X, Y) = \frac{P(w, Y | X)}{P(Y | X)} = \frac{\frac{P(w, Y | X)}{P(w)}}{\frac{P(Y | X) P(w)}{P(w)}}$$

$$= \frac{P(Y | w, X) P(w)}{\int P(Y | w, X) P(w) dw}$$

$$P(w) = P(w | X) \text{ (若然分布了参数)}$$

$$P(Y | w, X) = \prod_{i=1}^N P(y_i | w, X_i) = \prod_{i=1}^N N(y_i | w^T x_i, \sigma^2)$$

$$P(w | X, \sigma^2) = N(w | \bar{w}, \Sigma_p)$$

$$P(w | \text{Data}) \propto P(Y | w, X) \cdot P(w)$$

Gaussian

$$\propto \prod_{i=1}^N N(y_i | w^T x_i, \sigma^2) \cdot N(w, \Sigma_p)$$

$$\propto \prod_{i=1}^N N(w | \mu, \Sigma_p)$$



Vazyme

softmax Classifier Vectorized. Vazyme

$$L_i = -\log\left(\frac{e^{f_{y_i}}}{\sum_j e^{f_j}}\right) = -f_{y_i} + \log\sum_j e^{f_j}$$

$$f_j = x_i w_j$$

$$\frac{dL_i}{dw_j} = \frac{1}{\sum_j e^{f_j}} \cdot x_i \cdot e^{x_i w_j} \cdot \frac{e^{f_j}}{\sum_j e^{f_j}} \cdot x_i$$

$$\frac{dL_i}{dw_j} = -x_i + \frac{e^{f_{y_i}}}{\sum_j e^{f_j}} \cdot x_i = \left(\frac{e^{f_{y_i}}}{\sum_j e^{f_j}} - 1\right) x_i$$

$$\frac{dL_i}{dw} = \left[\frac{\frac{dL_i}{dw_1}}{\sum_j e^{f_j}} \cdot \frac{\frac{dL_i}{dw_2}}{\sum_j e^{f_j}} \cdot \dots \frac{\frac{dL_i}{dw_n}}{\sum_j e^{f_j}} \right]$$

$$= \left[\frac{x_i w_1}{\sum_j e^{f_j}} \cdot x_i - \frac{e^{x_i w_1}}{\sum_j e^{f_j}} \cdot x_i \right] \cdot \left[\frac{x_i w_2}{\sum_j e^{f_j}} \cdot x_i - \frac{e^{x_i w_2}}{\sum_j e^{f_j}} \cdot x_i \right] \cdot \dots \cdot \left[\frac{x_i w_n}{\sum_j e^{f_j}} \cdot x_i - \frac{e^{x_i w_n}}{\sum_j e^{f_j}} \cdot x_i \right]$$

$$Vazyme$$

$$= X \left[\begin{array}{c} \frac{e^{x_i w_1}}{\sum_j e^{f_j}} \cdot x_i - \frac{e^{x_i w_1}}{\sum_j e^{f_j}} \cdot x_i \\ \vdots \\ \frac{e^{x_i w_n}}{\sum_j e^{f_j}} \cdot x_i - \frac{e^{x_i w_n}}{\sum_j e^{f_j}} \cdot x_i \end{array} \right]$$

$$= [x_1, x_2, \dots, x_i] \cdot \left[\begin{array}{c} \frac{e^{f_j}}{\sum_j e^{f_j}} \cdot x_i \\ \vdots \\ \frac{e^{f_j}}{\sum_j e^{f_j}} \cdot x_i \end{array} \right]$$

$$\frac{dL_i}{dw} = \left[\frac{\frac{dL_i}{dw_1}}{\sum_j e^{f_j}} \cdot x_i - \frac{dL_i}{dw_1} \right]$$

$$= \left[\frac{e^{x_i w_1}}{\sum_j e^{f_j}} \cdot x_i - \frac{dL_i}{dw_1} \right]$$

$$= \left[\frac{e^{x_i w_1}}{\sum_j e^{f_j}} \cdot x_i - \frac{e^{x_i w_1}}{\sum_j e^{f_j}} \cdot x_i \right] \cdot \left[\frac{e^{x_i w_2}}{\sum_j e^{f_j}} \cdot x_i - \frac{e^{x_i w_2}}{\sum_j e^{f_j}} \cdot x_i \right] \cdot \dots \cdot \left[\frac{e^{x_i w_n}}{\sum_j e^{f_j}} \cdot x_i - \frac{e^{x_i w_n}}{\sum_j e^{f_j}} \cdot x_i \right]$$