

Final Review

Decision Tree

$\Delta H \& I$

$$H(Y) = -\sum P(Y_i) \log P(Y_i)$$

$$H(Y|X) = \sum P(X_i) H(Y|X_i)$$

$$I(Y; X) = H(Y) - H(Y|X) = H(X) - H(X|Y)$$

x_1	x_2	y
1	0	0
1	0	0
1	0	1
1	0	1
1	1	1
1	1	1
1	1	1

$$I(Y; X_1) = H(Y) - H(Y|X_1)$$

$$H(Y) = \frac{1}{4} \log \frac{1}{4} + \frac{3}{4} \log \frac{3}{4} = -\frac{1}{2} + \frac{3}{4} \log \frac{3}{4}$$

$$H(Y|X_1) = 1 \cdot H(Y|X_1=1) = 1 \cdot \log 2 = 0$$

$$H(Y|X_1) = \frac{1}{2} H(Y|X_1=0) + \frac{1}{2} H(Y|X_1=1)$$

$$= \frac{1}{2} \left(\frac{1}{2} \log \frac{1}{2} + \frac{1}{2} \log \frac{1}{2} \right) + \frac{1}{2} \cdot (1 \cdot \log 1)$$

$$= -\frac{1}{2}$$

$$\therefore I(Y; X_2) = \frac{3}{4} \log \frac{3}{4} > I(Y; X_1) = -\frac{1}{2} + \frac{3}{4} \log \frac{3}{4}$$

决策树分割：每次分割 $I(Y; X_i)$ 最大 $\geq X_i$

5. [20 points] [Decision Tree] A dataset is given below. Now we want to discover the relationship between the features and the target variable by using a Decision Tree.

Outlook (X_1)	Temperature (X_2)	Humidity (X_3)	Play Tennis? (Y)
sunny	hot	high	no
overcast	hot	high	yes
rain	mild	high	yes
rain	cool	normal	yes
sunny	mild	high	no
sunny	mild	normal	yes
rain	mild	normal	yes
overcast	hot	normal	yes

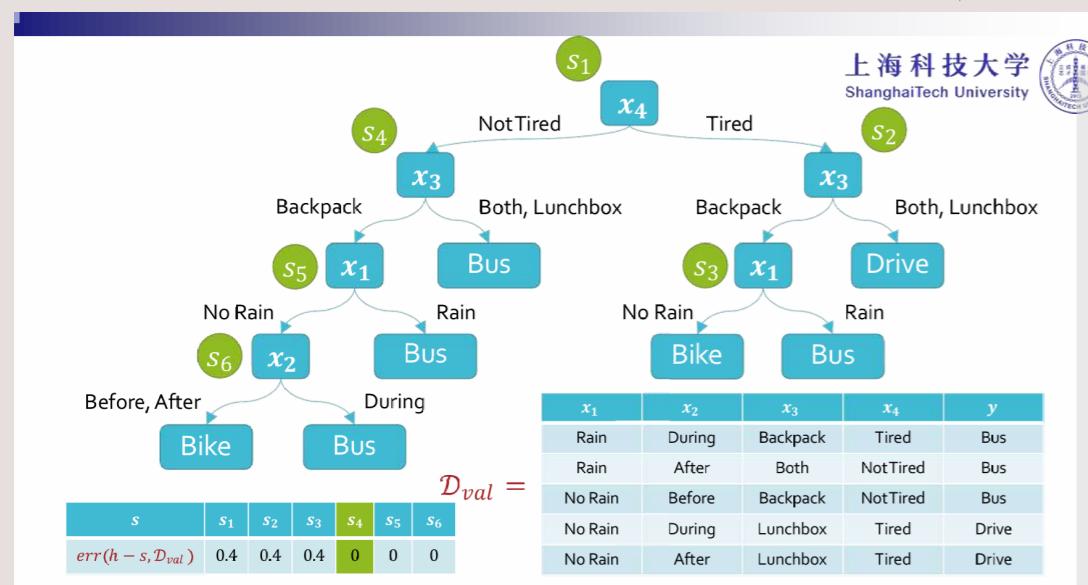
- (a) Using the dataset above, calculate the mutual information for each feature (X_1, X_2, X_3) to determine the root node for a Decision Tree trained on the above data.
- What is $I(Y; X_1)$? [3 points]
 - What is $I(Y; X_2)$? [3 points]
 - What is $I(Y; X_3)$? [3 points]
 - What feature should be split on at the root node? [1 points]
- (b) Calculate what the next split should be. [5 points]
- (c) Draw the resulting tree. [5 points]

x_1	x_2	\hat{y}	y	Mistake?
-1	2	+	-	Yes
1	0	+	+	No
1	1	-	+	Yes
-1	0	-	-	No
-1	-2	+	-	Yes
1	-1	+	+	No

(初值 $\vec{w} = [0]$, 又 $b = 0$)

剪枝：所加树叶节点编号 S_i

$err[S_i]$ 表示 S_i 子树缩为一点后训练误差



KNN

训练集：已有连接（有分类标签）

新测试点：寻找最近 k 个训练点 \Rightarrow 投票决定标签

$$N \rightarrow +\infty:$$

$$error(h) < 2 \times Bayes \text{ Error rate}$$

平局处理：投票 / 算数 / 加权 / 按距离函数

Perceptron

参数： $\vec{w}, b \xrightarrow{\text{bias}}$

$$\text{决策}: \text{sign}(\vec{w}^T \vec{x} + b) = \text{sign}([\vec{w}^T \vec{x}]^T [1]) = h^{(+)}$$

$h^{(+)} = \gamma^{(+)} : \text{参数不重}$

$$h^{(+)} \neq \gamma^{(+)}: \begin{cases} \vec{w}^+ = \gamma^{(+)} \vec{x}^{(+)}, \vec{w}^+ = [\vec{w}^T \vec{x}]^T [1] \\ b^+ = \gamma^{(+)} \end{cases}$$

x_1	x_2	\hat{y}	y	Mistake?
-1	2	+	-	Yes
1	0	+	+	No
1	1	-	+	Yes
-1	0	-	-	No
-1	-2	+	-	Yes
1	-1	+	+	No

ΔMargin : 在分界线上最大 \downarrow 一边距离

最小值 $\rightarrow \text{Margin} = \gamma \downarrow \vec{w}^T \vec{x} + b = 0$

算另一边距离：

边界上一点为 \vec{x} , 对面一点为 \vec{x}'

$$\text{距离} = |(\vec{x}' - \vec{x})^T \frac{\vec{w}}{\|\vec{w}\|}| \rightarrow \text{单位法向量}$$

$$= |\vec{x}^T \vec{w} - \vec{x}'^T \vec{w}| \frac{1}{\|\vec{w}\|} = |\vec{w}^T \vec{x}' + b| \frac{1}{\|\vec{w}\|}$$

△ Mistake Bound

$$R \cdot \gamma \Rightarrow \text{犯错次数} = \left(\frac{R}{\gamma}\right)^2$$

Def: $\vec{\theta}^{(k)}$ 在 $\vec{x}^{(i)}$, $\gamma^{(i)}$ 处犯错:

$$\gamma^{(i)} \vec{x}^{(i)\top} \vec{\theta}^{(k)} \leq 0$$

$$\Rightarrow \vec{\theta}^{(k+1)} = \vec{\theta}^{(k)} + \gamma^{(i)} \vec{x}^{(i)}$$

\therefore 边界 = γ

$\therefore \exists \vec{w}$, $\|\vec{w}\|=1$, 使得 \vec{w} 垂直于 $\vec{\theta}^{(k)}$:

$$\gamma^{(i)} \vec{x}^{(i)\top} \vec{w} \geq \gamma \rightarrow \begin{array}{c} \text{图示: } \\ \text{一个圆, 圆心在 } \vec{w} \text{ 上, } \vec{w} \text{ 垂直于 } \vec{\theta}^{(k)} \\ \text{圆的半径为 } \gamma^{(i)}, \text{ 圆与 } \vec{x}^{(i)} \text{ 相交} \end{array}$$

$$\vec{\theta}^{(k+1)\top} \vec{w} = \vec{\theta}^{(k)\top} \vec{w} + \gamma^{(i)} \vec{x}^{(i)\top} \vec{w}$$

$$\geq \vec{\theta}^{(k)\top} \vec{w} + \gamma$$

$$\therefore \vec{\theta}^{(k+1)\top} \vec{w} \geq k\gamma \quad \textcircled{1}$$

$$\therefore \|\vec{x}^{(i)}\| \leq R$$

$$\begin{aligned} \text{算模长: } & \|\vec{\theta}^{(k+1)}\|^2 = (\vec{\theta}^{(k)} + \gamma^{(i)} \vec{x}^{(i)})^2 \\ &= \vec{\theta}^{(k)\top} \vec{\theta}^{(k)} + \vec{x}^{(i)\top} \vec{x}^{(i)} + 2(\gamma^{(i)} \vec{x}^{(i)\top} \vec{\theta}^{(k)}) \\ &\leq \|\vec{\theta}^{(k)}\|^2 + R^2 \quad \leq 0 \end{aligned}$$

$$\Rightarrow \|\vec{\theta}^{(k+1)}\|^2 \leq kR^2 \Rightarrow \|\vec{\theta}^{(k+1)}\| \leq \sqrt{k}R \quad \textcircled{2}$$

$$\textcircled{1} \textcircled{2} \Rightarrow \sqrt{k}R \geq \|\vec{\theta}^{(k+1)}\| \geq \vec{\theta}^{(k+1)\top} \vec{w} \geq k\gamma$$

$$\Rightarrow k = \left(\frac{R}{\gamma}\right)^2$$

Kernel

$x^{(i)} \rightarrow \phi(x^{(i)}) \Rightarrow x^{(i)}$ 各特征以基函数映射

$$\phi(x_1) \phi(x_2) = K(x_1, x_2)$$

SVM

画出: 找出分类面 \vec{w} :

$$\text{满足约束 } \left\{ \begin{array}{l} \|\vec{w}\|_2 = 1 \\ f(\vec{x}_i, y_i), \forall y_i(\vec{w}^\top \vec{x}_i + b) \geq \gamma \end{array} \right.$$

边界计算

\rightarrow 情况下, 使 γ 最大

$$\Rightarrow \vec{w} = \arg \max_{\vec{w}} \gamma$$

$$\text{For: } \left\{ \begin{array}{l} \|\vec{w}\| = 1 \\ y_i(\vec{w}^\top \vec{x}_i + b) \geq \gamma \quad (\forall i) \end{array} \right.$$

此时可除去 \vec{w} , 只保留 γ 来间接
调整 γ 变化情况和 \vec{w} 原值

$$\left. \begin{array}{l} \text{可凸对偶变化: } \vec{w}' = \frac{\vec{w}}{\gamma}, \text{ 如} \\ \|\vec{w}'\| \text{ 变化只取决于 } \gamma \text{ 变化} \end{array} \right)$$

变动: $\|\vec{w}\|$ 固定 γ 变化 $\Rightarrow \gamma$ 固定 $\|\vec{w}\|$ 变化 \rightarrow \vec{w} 变化量 \rightarrow \vec{w} 原值

$$\left. \begin{array}{l} \text{约束: } y_i \left(\frac{\vec{w}^\top \vec{x}_i + b}{\gamma} \right) \geq 1 \Rightarrow y_i (\vec{w}'^\top \vec{x}_i + b) \geq 1 \end{array} \right)$$

$\therefore \gamma$ 最大 $\Rightarrow \vec{w}'$ 最小

问题变为:

$$\vec{w}' = \arg \min_{\vec{w}'} \frac{1}{2} \|\vec{w}'\|^2$$

$$\text{使得 } \forall i: y_i (\vec{w}'^\top \vec{x}_i + b) \geq 1$$

$$\text{即: } \min \frac{1}{2} \|\vec{w}'\|^2$$

$$\text{使得 } \forall i: y_i (\vec{w}^\top \vec{x}_i + b) \geq 1$$

$$\rightarrow -y_i (\vec{w}^\top \vec{x}_i + b) + 1 \leq 0$$

观察:

$$\left. \begin{array}{l} \|\vec{w}\|=1 \quad \angle \vec{w} = \theta_1 \text{ 时, 有 } \gamma_{(1)} \\ \angle \vec{w} = \theta_2 \text{ 时, 有 } \gamma_{(2)} \end{array} \right)$$

\Rightarrow \vec{w} 变化时, $\|\vec{w}\|$ 可用
 \vec{w} 于约束 γ , 于约束
于 \vec{w} 变化情况!

$$\text{目标: } \min \max \sum (w, b, \alpha)$$

(利用特征二利)

$$\text{Lagrange 对偶: } \max \min \sum L(w, b, \alpha)$$

$$\mathcal{L} = \frac{1}{2} \|\vec{w}\|^2 + \sum_{i=1}^m \alpha_i [y_i (\vec{w}^\top \vec{x}_i + b) + 1]$$

$$\nabla_w \mathcal{L}(\vec{w}, b, \vec{\alpha}) = \vec{w} + \sum_{i=1}^m \alpha_i (-y_i \vec{x}_i) = 0$$

$$\Rightarrow \vec{w} = \sum_{i=1}^m \alpha_i y_i \vec{x}_i$$

$$\nabla_b \mathcal{L}(\vec{w}, b, \vec{\alpha}) = \sum_{i=1}^m (\alpha_i y_i) = 0 \Rightarrow \sum_{i=1}^m \alpha_i y_i = 0$$

$$\text{四行: } \mathcal{L} = \frac{1}{2} \|\vec{w}\|^2 - \sum_{i=1}^m \alpha_i y_i (\vec{w}^\top \vec{x}_i + b) + \sum_{i=1}^m \alpha_i$$

$$= \frac{1}{2} \sum_{i,j=1}^m \alpha_i \alpha_j y_i y_j \vec{x}_i^\top \vec{x}_j - \sum_{i,j=1}^m (\alpha_i y_i \vec{x}_i)^\top (\alpha_j y_j \vec{x}_j) + \alpha_i y_i b + \sum_{i=1}^m \alpha_i$$

$$= -\frac{1}{2} \sum_{i,j=1}^m \alpha_i \alpha_j y_i y_j \vec{x}_i^\top \vec{x}_j - b \sum_{i=1}^m \alpha_i y_i + \sum_{i=1}^m \alpha_i$$

$$= \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i,j=1}^m \alpha_i \alpha_j y_i y_j \vec{x}_i^\top \vec{x}_j$$

$\hookrightarrow \mathcal{L}$ 的变量为 $\vec{\alpha}$ ($\alpha_i \geq 0$)

对偶问题:

$$\text{找 } \arg \max_{\vec{\alpha}} \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i,j=1}^m \alpha_i \alpha_j y_i y_j \vec{x}_i^\top \vec{x}_j$$

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

\rightarrow 可优化 \checkmark

找出 \vec{w} 后: $\vec{w} = \sum_{i=1}^m \alpha_i y_i \vec{x}_i$ (b 也可得出)

由 KKT 条件: 若 \vec{x}_i 在边界上, α_i 才 $\neq 0$

否则 $\alpha_i = 0$ \hookrightarrow 支持向量

\therefore 对偶问题:

$$\arg \min_{\vec{\alpha}} \frac{1}{2} \sum_{i,j=1}^m \alpha_i \alpha_j y_i y_j \vec{x}_i^\top \vec{x}_j - \sum_{i=1}^m \alpha_i$$

$$\text{s.t. } \begin{cases} \forall i \quad 0 \leq \alpha_i \leq C \quad (\text{对 } \alpha_i - \text{不放}) \\ \forall i \quad \sum_{i=1}^m \alpha_i y_i = 0 \end{cases}$$

Linear Regression

$$J(\vec{w}, b) = \frac{1}{N} \sum_{i=1}^N |y^{(i)} - (\vec{w}^\top \vec{x}^{(i)} + b)|^2$$

$$\text{优化: } \vec{\theta} = \arg \min_{\vec{\theta}} J(\vec{\theta}) \quad (\vec{\theta} = [b \vec{w}])$$

$$\text{梯度: } = \arg \min_{\vec{\theta}} \frac{1}{N} \sum_{i=1}^N (y^{(i)} - \sum_{j=1}^m \theta_j x_j^{(i)})^2 \rightarrow \text{梯度}$$

$$\vec{\theta} = \eta \frac{\partial J(\vec{\theta})}{\partial \vec{\theta}} \Rightarrow \theta_i = \frac{\partial J(\vec{\theta})}{\partial \theta_i}$$

$$= \frac{1}{N} \sum_{i=1}^N (\vec{\theta}^\top \vec{x}^{(i)} - y^{(i)}) \cdot x_j^{(i)}$$

\hookrightarrow 批量梯度

$$\text{随机梯度: } J(\vec{\theta}) = \frac{1}{N} \sum_{i=1}^N J^{(i)}(\vec{\theta})$$

$$J^{(i)}(\vec{\theta}) = \frac{1}{2} (\vec{\theta}^\top \vec{x}^{(i)} - y^{(i)})^2$$

$$\frac{\partial J^{(i)}(\vec{\theta})}{\partial \theta_j} = (\vec{\theta}^\top \vec{x}^{(i)} - y^{(i)}) \cdot \vec{x}_j^{(i)}$$

$$\text{闭式解: } \nabla_{\vec{\theta}} J(\vec{\theta}) = \frac{1}{N} \frac{\partial}{\partial \vec{\theta}} \sum_{i=1}^N \frac{1}{2} (\vec{\theta}^\top \vec{x}^{(i)} - y^{(i)})^2$$

$$= \frac{1}{2N} \frac{\partial}{\partial \vec{\theta}} (\vec{x} \vec{\theta}^\top - \vec{y})^\top (\vec{x} \vec{\theta}^\top - \vec{y})$$

$$\hookrightarrow (\vec{\theta}^\top \vec{x}^\top \vec{x} \vec{\theta} - 2 \vec{x}^\top \vec{x} \vec{\theta} + \vec{y}^\top \vec{y})$$

$$= \frac{1}{2N} (2 \vec{x}^\top \vec{x} \vec{\theta} - 2 \vec{x}^\top \vec{y}) = 0$$

$$\Rightarrow \vec{\theta} = (\vec{x}^\top \vec{x})^{-1} \vec{x}^\top \vec{y}$$

\hookrightarrow 常可逆

Δ 核心: 边界分类

$$\text{核二: } \arg \min_{\vec{w}, \vec{\xi}} \frac{1}{2} \|\vec{w}\|^2 + C \sum_i \xi_i$$

$$\text{s.t. } \begin{cases} \forall i \quad y_i \vec{w}^\top \vec{x}_i \geq 1 - \xi_i \quad \text{允许误差} \\ \xi_i \geq 0 \end{cases}$$

$$\mathcal{L} = \frac{1}{2} \|\vec{w}\|^2 + C \sum_i \xi_i - \sum_i [y_i (\vec{w}^\top \vec{x}_i + b) - 1 + \xi_i] - \sum_i \beta_i \xi_i$$

\hookrightarrow 参数: $\vec{w}, b, \vec{\xi}$

\hookrightarrow 成分: $C - \alpha_i - \beta_i = 0$

$$\Rightarrow 0 \leq \alpha_i \leq C$$

MLE & MAP (Lec 9)

演绎法 2

Logistic Regression

$$\sigma(z) = \frac{1}{1+e^{-z}} = \frac{1}{1+e^{-\theta^T \vec{x}}} = h_{\theta}(\vec{x})$$

$$\frac{\partial \sigma(z)}{\partial z} = \sigma(z)(1-\sigma(z))$$

$$P(Y=1 | \vec{x}; \theta) = h_{\theta}(\vec{x})$$

$$P(Y=0 | \vec{x}; \theta) = 1 - h_{\theta}(\vec{x})$$

$$\text{似然: } L(\theta) = \prod_{i=1}^n (h_{\theta}(\vec{x}^{(i)})^{y^{(i)}}) (1 - h_{\theta}(\vec{x}^{(i)}))^{(1-y^{(i)})}$$

$$l(\theta) = \sum_{i=1}^n \left(y^{(i)} \log h_{\theta}(\vec{x}^{(i)}) + (1-y^{(i)}) \log (1 - h_{\theta}(\vec{x}^{(i)})) \right)$$

$$= \sum_{i=1}^n y^{(i)} \log \frac{1}{1+e^{-\theta^T \vec{x}}} + (1-y^{(i)}) \log \frac{e^{-\theta^T \vec{x}}}{1+e^{-\theta^T \vec{x}}} \rightarrow \frac{1}{1+e^{-\theta^T \vec{x}}}$$

$$= \sum_{i=1}^n y^{(i)} (\theta^T \vec{x} - \log(1+e^{\theta^T \vec{x}})) + (1-y^{(i)}) (-\log(1+e^{\theta^T \vec{x}}))$$

$$J(\theta) = \frac{1}{n} l(\theta)$$

优化: 梯度下降.

$$\nabla_{\theta} J(\theta) = -\frac{1}{n} \sum_{i=1}^n \left(y^{(i)} \vec{x}^{(i)} - \frac{e^{\theta^T \vec{x}^{(i)}}}{1+e^{\theta^T \vec{x}^{(i)}}} \vec{x}^{(i)} \right) \rightarrow P(Y=1 | \vec{x}^{(i)}; \theta)$$

$$= -\frac{1}{n} \sum_{i=1}^n (y^{(i)} - g(\theta^T \vec{x}^{(i)})) \vec{x}^{(i)}$$

$$\Rightarrow J(\theta) = -\gamma \nabla_{\theta} J(\theta)$$

$$\text{若用 SGD: } \nabla_{\theta} J^{(i)}(\theta) = -\frac{1}{n} (y^{(i)} - g(\theta^T \vec{x}^{(i)})) \vec{x}^{(i)}$$

$$\text{Mini SGD: } \nabla_{\theta} J^{(i)}(\theta) = -\frac{1}{n} \sum_{s=1}^S (y^{(s)} - g(\theta^T \vec{x}^{(s)})) \vec{x}^{(s)}$$

Regularization

$$J(\theta) \rightarrow J(\theta) + \lambda R(\theta)$$

$$R(\theta) = \|\theta\|_2^2 \text{ or } \|\theta\|_1$$

$$\sum \theta_m^2 \quad \sum |\theta_m| \rightarrow \text{梯度下降法}$$

Neural Network

△ 每层结构: $a_j^{(l)} \rightarrow z_j^{(l)}$ (线性 & 激活)

$$a_j^{(l)} = \sum_{i=0}^D w_{j,i}^{(l)} z_i^{(l-1)}, z_j^{(l)} = f(a_j^{(l)})$$

↳ 权重 = 偏置 + (1)

$$\xrightarrow{\text{矩阵形式: }} \vec{a}^{(l)} = W^{(l)} \vec{z}^{(l-1)}, \vec{z}^{(l)} = [1, f(\vec{a}^{(l)})]^T$$

$$\begin{bmatrix} W^{(l)} \\ \vdots \end{bmatrix} \begin{bmatrix} \vec{z}^{(l-1)} \\ 1 \end{bmatrix} = \vec{a}^{(l)}$$

第 j 行第 i 列 \rightarrow 行向量即为该连接后 j 权重的值
↳ 行数 = 下一层向量维数

若多分类: 最后一层 Output 层先发到 Softmax

反向传播

$$L(W) = \frac{1}{N} \sum L_i + R(W)$$

$$L_i(x_i, y_i, W) = -\log P(Y=y_i | X=x_i)$$

↳ 分类
即: 预测值与真实值之差

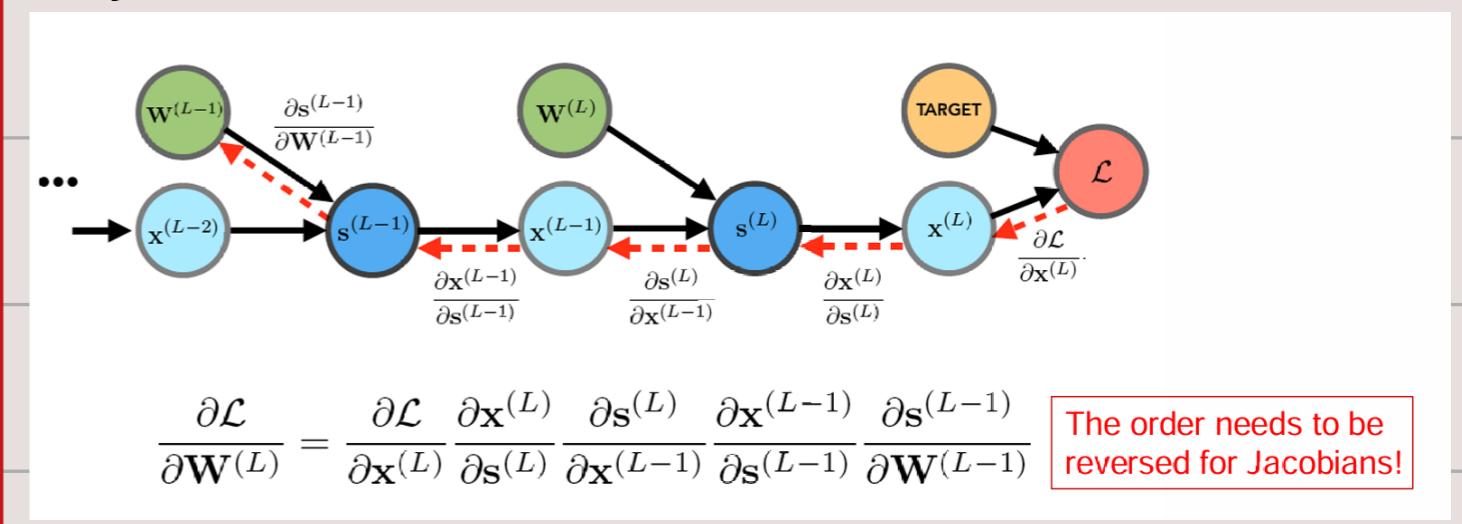
回溯: 最小二乘 or 交叉熵

$$W^{(l+1)} = W^{(l+1)} - \eta \frac{\partial E}{\partial W^{(l+1)}}$$

↳ 根据计算图逆向传播梯度

每层: 激发梯度 \rightarrow 对应梯度

顶层: 对改名为对 $W^{(L)}$



Computing the loss:

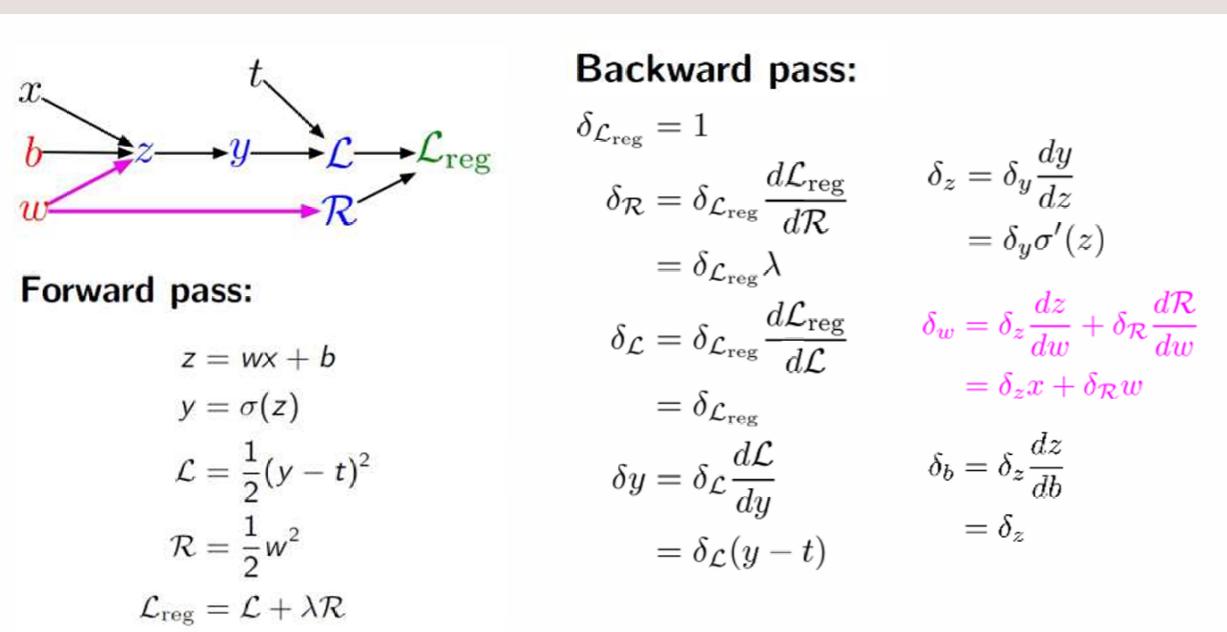
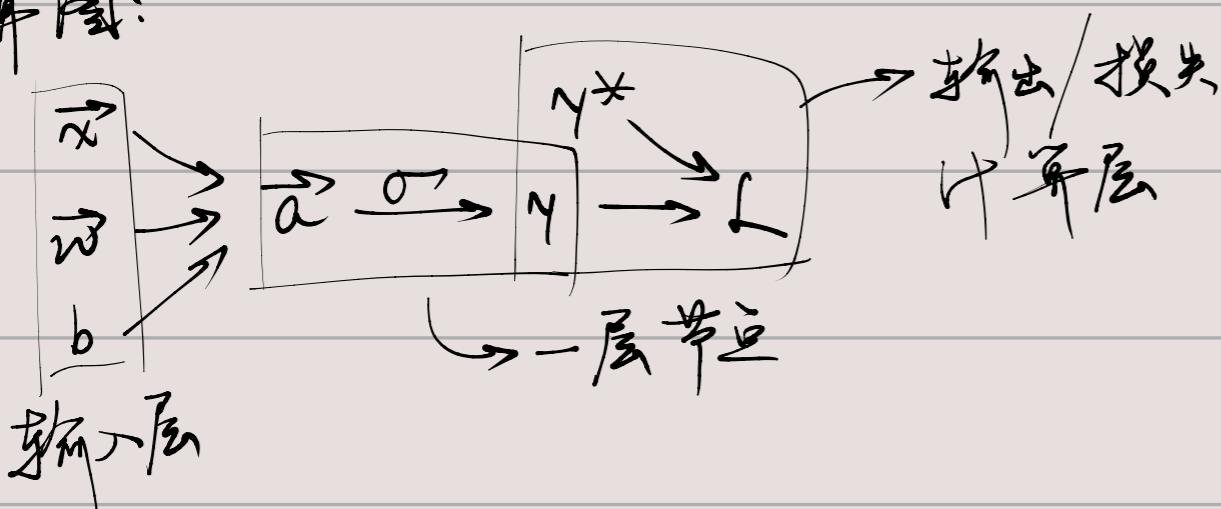
$$\begin{aligned}s &= wx + b \\y &= \sigma(s) \\L &= \frac{1}{2}(y - t)^2\end{aligned}$$

Computing the derivatives:

$$\begin{aligned}\frac{dL}{dy} &= y - t \\ \frac{dL}{ds} &= \frac{dL}{dy} \sigma'(s) \\ \frac{dL}{dw} &= \frac{dL}{ds} x \\ \frac{dL}{db} &= \frac{dL}{ds}\end{aligned}$$

$\vec{x} \rightarrow \vec{a} \rightarrow \vec{y} \rightarrow L$

计算图:



Recommender Systems

等价法 2

Boosting / AdaBoosting

初始数据: 权重 $D_t(i) = \frac{1}{N}$ (等重)

数据点 $(\vec{x}_i, y_i) \dots (\vec{x}_n, y_n), y_i \in \{-1, 1\}$

△ 遍历步骤:

- * 根据当前数据训练新弱分类器 h_t 一个
强分类器
权重和
- * 计算 h_t 错误 ϵ_t : $\epsilon_t = \sum_{i=1}^N \mathbb{1}\{h_t(\vec{x}_i) \neq y_i\} \cdot D_t(i)$
- * 计算该分类器权重 α_t : $\alpha_t = \frac{1}{2} \ln \frac{1 - \epsilon_t}{\epsilon_t}$
- ⇒ ϵ_t 高: 分类不好 $\Rightarrow \alpha_t$ 低
 ϵ_t 低: 分类效果好 $\Rightarrow \alpha_t$ 高

* 更新数据权重:

如果分类对 $S_i = \exp(-\alpha_t \cdot y_i \cdot h_t(\vec{x}_i))$ 下一次训练下降

Scoring 给分值: $S_i = \exp(-\alpha_t \cdot y_i \cdot h_t(\vec{x}_i))$ 如果 $S_i = \exp(\alpha_t)$ 下一次更关注

Scoring: $D_t(i) \cdot S_i$

$$\text{归一化为分布: } D_{t+1}(i) = \frac{D_t(i) \cdot S_i}{\sum_t D_t(i) \cdot S_i} \rightarrow = \sum_t D_t(i) \cdot S_i$$

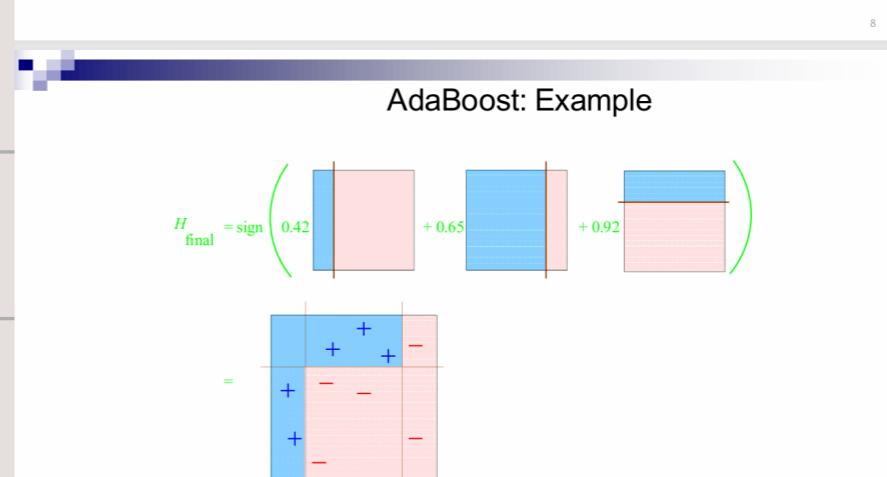
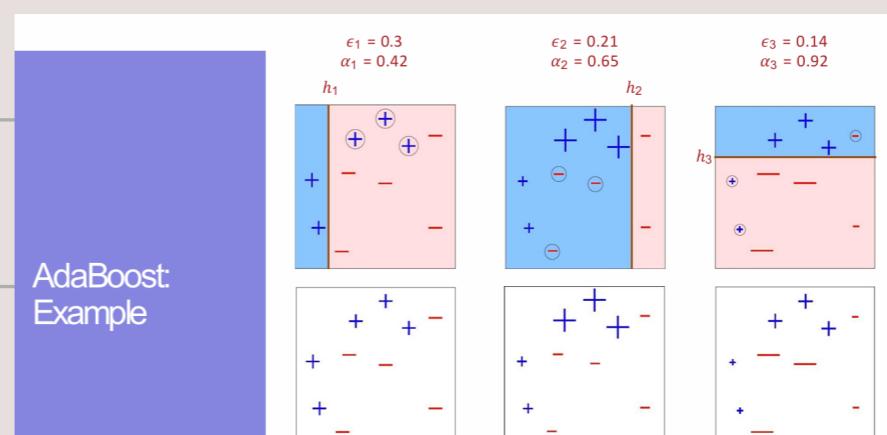
Matrix Calculus

等价法 2

CNN

等价法 2

$$H(x) = \text{sign} \left(\sum_{t=1}^T \alpha_t h_t(\vec{x}) \right) \quad (\text{分类得分})$$



AdaBoosting 可最小化指数损失函数:

等价版 2

Bagging

Sample Bagging

每组可放回抽样若干数据点，每组 m.

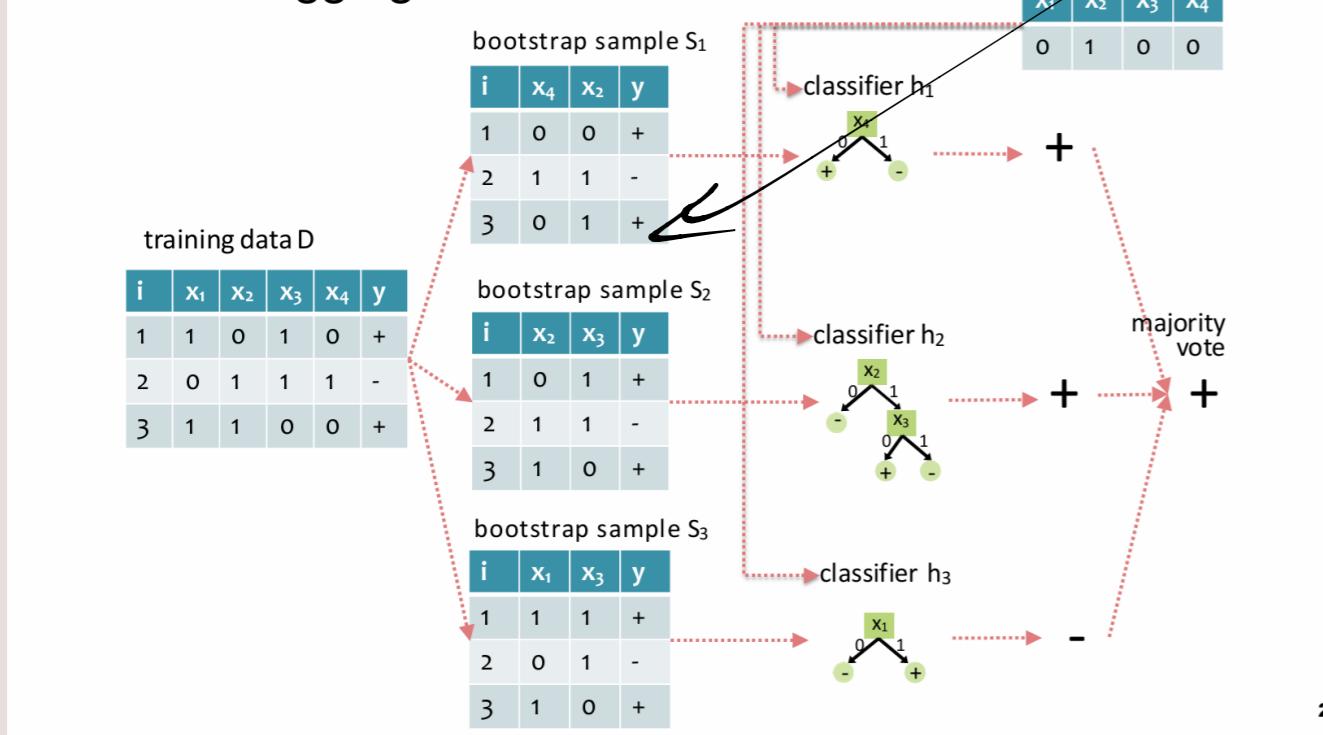
以训练模型。最终结果为模型投票(分类)或取平均(回归)产生。

Feature Bagging

同样 with Replacement

将“放样数据点”改为“放样特征”

Feature Bagging



Random Forests

Sample Bagging + Feature Bagging

子数据集构建 → 分裂时只考虑部分特征

⇒ 随机森林 + 分裂特征随机化

⇒ 训练出若干棵决策树，投票或取平均

(参数为 B, P)

随机森林选择：等价版 2

EM

△ Jensen 不等式

凸函数 \cup : $E(f(X)) \geq f(EX)$

凹函数 \cap (\log): $E(f(X)) \leq f(EX)$

△ 效率: $\rightarrow \rightarrow \cdots$ 无偏差。

固执: 这又～某模型 (大前提)

找出最佳参数 $\vec{\theta}$, 使 $l(\vec{\theta}) = \sum_{i=1}^m \log p(\vec{x}_i; \vec{z}_i; \vec{\theta})$

设 $z^{(i)}$ 为 $x^{(i)}$ 的潜在变量 (如 GMM 中, $z^{(i)}$ 是子 \Rightarrow 值域为 $\{1, 2, 3, \dots, k\}$)
 $x^{(i)} \sim$ 哪个高斯, $\therefore z^{(i)} \sim$ 多项式分布)

$$\Rightarrow l(\vec{\theta}) = \sum_{i=1}^m \log \sum_{z_i} p(\vec{x}_i; z_i; \vec{\theta})$$

即: 从潜在变量 z 为中, 找出最优 $\vec{\theta}$

△ 优化指导

对于 $l(\vec{\theta}) = \sum_{i=1}^m \log \sum_{z_i} p(\vec{x}_i; z_i; \vec{\theta})$:

$$l(\vec{\theta}) = \sum_{i=1}^m \log \sum_{z_i} Q_i(z_i) \frac{p(\vec{x}_i; z_i; \vec{\theta})}{Q_i(z_i)}$$

$$\text{令随机变量 } X = \frac{p(\vec{x}_i; z_i; \vec{\theta})}{Q_i(z_i)} = F(z_i), f(x) = \log x$$

$$\text{则 } l(\vec{\theta}) = \sum_{i=1}^m f(EX) \geq \sum_{i=1}^m E(f(X)) \quad (\text{Jensen 不等式})$$

$$= \sum_{i=1}^m \sum_{z_i} Q_i(z_i) \log \frac{p(\vec{x}_i; z_i; \vec{\theta})}{Q_i(z_i)} \quad (1)$$

意义说明: (1) 式为: 固定 $\vec{\theta}$ 时, $l(\vec{\theta})$ 为一个下界函

数. 该函数的变量为 $Q_i(z_i)$. 即: 取不同 $= Q_i(z_i)$

(即 z 分布), 有不同之下界值.

⇒ 期望: 找到一个 $Q_i(z_i)$, 使得下界如 - 确保

使的最紧. ⇒ 即: 该随机变量 X , Jensen

不等式都取等.

X 之构造为: $X = c$ (常数)

$$\Rightarrow \frac{p(\vec{x}_i, z_i; \vec{\theta})}{Q_i(z_i)} = c$$

求解 $Q_i(z_i)$:

$$Q_i(z_i) \propto p(\vec{x}_i, z_i; \vec{\theta})$$

$$\begin{aligned} Q_i(z_i) &\propto Q_i(z_i) = \frac{p(\vec{x}_i, z_i; \vec{\theta})}{\sum_z p(\vec{x}_i, z; \vec{\theta})} \rightarrow \text{所有点的和} \\ &= \frac{p(\vec{x}_i, z_i; \vec{\theta})}{p(\vec{x}_i; \vec{\theta})} \end{aligned}$$

$$= p(z_i | \vec{x}_i; \vec{\theta})$$

\Rightarrow 此为 E 步第二: 在固定 $\vec{\theta}$ 的情况下, 找到最合适潜在辅助变量.

数学本质为利用 Jensen 不等式右端函数

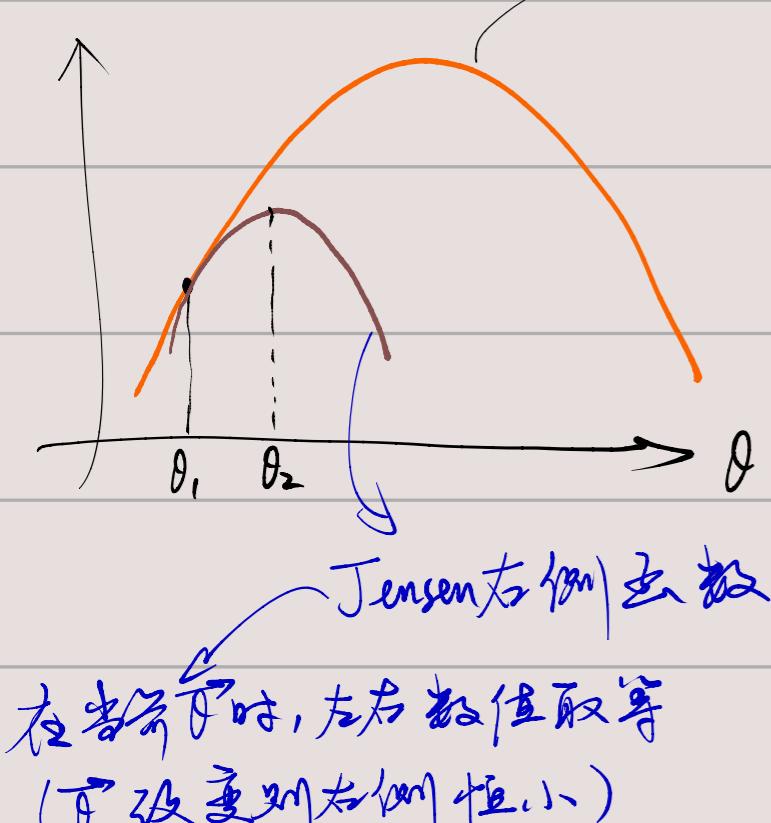
(右端为 $l(\vec{\theta})$) 找到 $l(\vec{\theta})$ 此时最大二下界, 即

最好一优化结果 ($= l(\vec{\theta})$)

M 步: $\vec{\theta} = \arg \max \sum_i \sum_z Q_i(z_i) \log \frac{p(\vec{x}_i, z_i; \vec{\theta})}{Q_i(z_i)}$ 找出 Jensen 右侧取最大值时 $= \vec{\theta}$.

图解:

E 步 1



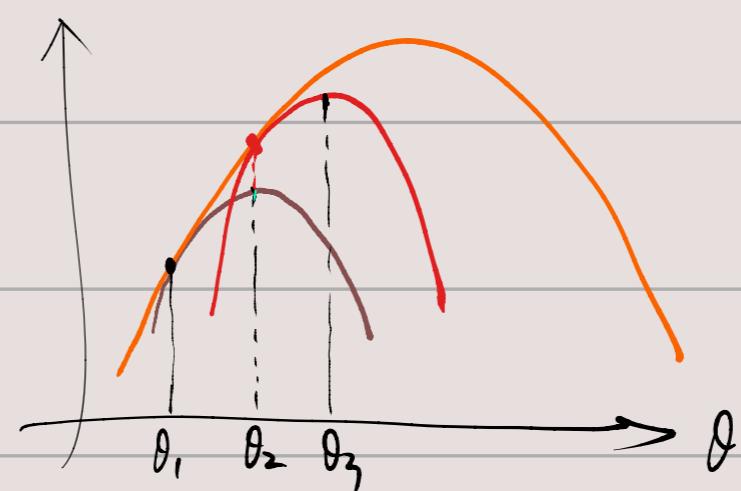
E 步: 在当前 $\vec{\theta}$ 下, 找到 Jensen 右侧函数, 使其线性化

即 \approx Jensen 右侧函数

M 步: 沿 Jensen 右侧函

数优化, 找到其取最
优时对应 $\vec{\theta}_2$

E 步 2:



E 步: 在当前 $\vec{\theta}_2$ 情况下, 找到最优

Jensen 右侧函数, 使其线性化后合

Jensen 右侧函数 $= l(\vec{\theta})$

M 步: 优化 Jensen 右侧函数, 迭代新 $\vec{\theta}$.

其余步骤依次类推

GMM 估计 $\vec{\theta}$ (参数)

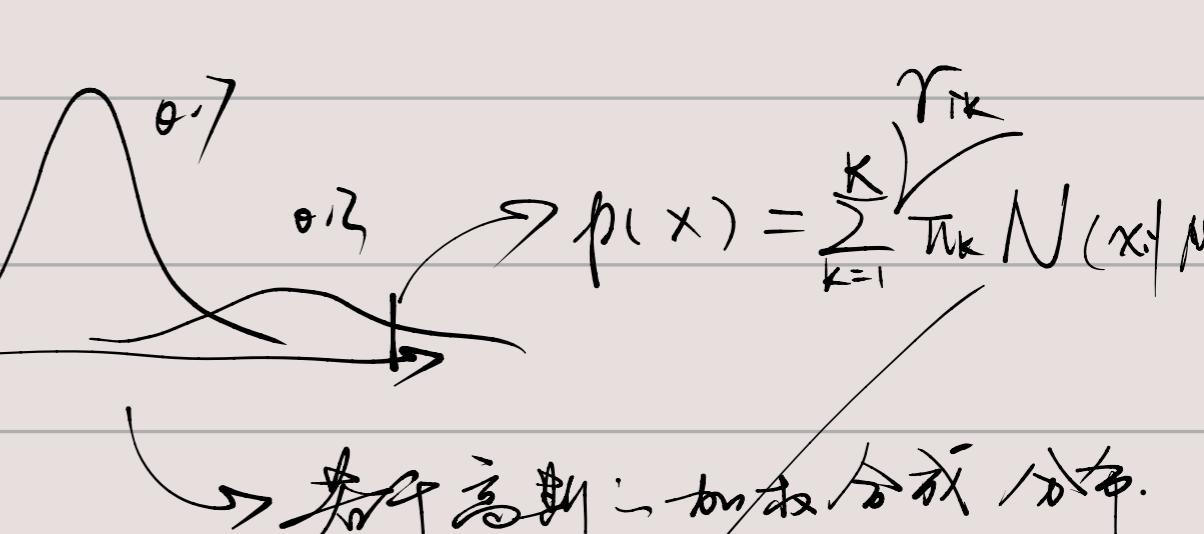
Problem: 已有观测值

π, μ, Σ

学习 k 个高斯分布 (即: GMM), 在此

分布中, 每一点 x 属于哪个高

斯) 推论: 使观测到该点概率



$$\sum_{x_i=k} p(x_i=k) N(x_i | \mu_k, \Sigma_k)$$

$E M \rightarrow \mu_k, \Sigma_k, \pi_k \rightarrow$ 板書
 \hookrightarrow 極端？

$$p(z_i=k | x_i) = \gamma_{ik}$$

$$E: \gamma_{ik} = p(z^{(i)}=k | x^{(i)}; \phi, \mu, \Sigma)$$

\propto (Rao-Blackwell)

$$\begin{aligned} M: l &= \log \prod_{i=1}^N \sum_{k=1}^K \gamma_{ik} (\pi_k \cdot N(x_i | \mu_k, \Sigma_k)) \\ &= \sum \sum \gamma_{ik} [\log \pi_k + \log N(x_i | \mu_k, \Sigma_k)] \\ &\quad \hookrightarrow f(\pi_k, \mu_k, \Sigma_k^{-1}) \\ \left\{ \begin{array}{l} \frac{\partial l}{\partial \pi_k} = 0 \\ \frac{\partial l}{\partial \mu_k} = 0 \\ \frac{\partial l}{\partial \Sigma_k^{-1}} = 0 \end{array} \right. \end{aligned}$$

$$\Rightarrow \pi_k = \frac{1}{n} \sum \gamma_{ik}$$

$$\mu_k = \frac{\sum \gamma_{ik} x_i}{\sum \gamma_{ik}}$$

$$\Sigma_k = \frac{\sum \gamma_{ik} (x_i - \mu_k)(x_i - \mu_k)^T}{\sum \gamma_{ik}}$$