

$$\hat{Y} = f(\hat{x})$$

$$Y = f(x) + \varepsilon$$

Pattern Recognition
and machine learning

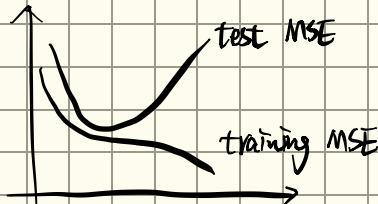
Inference need stronger assumptions.

$$\begin{aligned} E(Y - \hat{Y})^2 &= E((f(x) + \varepsilon_i) - (\hat{f}(x)))^2 \\ &= E((f(x) - \hat{f}(x))^2 + 2\varepsilon_i(f(x) - \hat{f}(x)) + \varepsilon_i^2) \\ &= E((f(x) - \hat{f}(x))^2) + \text{Var}(\varepsilon_i) \\ &\quad + 2E[\varepsilon_i(f(x) - \hat{f}(x))] \end{aligned}$$

parametric model: assume $f(x)$ is known

non-parametric: 不是没有参数而是没有预先假设的参数形式

$$MSE = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{f}(x_i))^2 \xrightarrow{\sum (y_i - g(x_i))^2 \rightarrow \int g''(t)^2 dt} g(t) \text{ 没有假设}$$



$$\sum_{i=1}^n (y_i - g(x_i))^2 \xrightarrow{\int g''(t)^2 dt} \text{non parametric}$$

根据不同输入下不同的 testing error (过拟合模型) \Rightarrow MSE is difficult \Leftarrow no test data
 \Rightarrow Cross Validation

Bias-Variance Trade off

$$\begin{aligned} &E_P[(f(x_0) - \hat{f}(x_0))^2] \xrightarrow{\text{如果增加 } Y, \text{ 增加 } \text{Var}(\varepsilon)} \text{样本只有1个数据点后, 不是 test 而是} \\ &= E[(f(x_0) - E_D\hat{f})^2 + E_D(\hat{f}) - f]^2 \xrightarrow{\text{没有 randomness 可将 } E_D \text{ 直接去掉} \Rightarrow \text{bias}} \\ &= E[(f - E_D\hat{f})^2] + E[(E_D\hat{f} - f)]^2 \\ &\quad + 2E_D[(f(x_0) - E_D\hat{f})(E_D\hat{f} - f)] = 0 \\ &= \text{bias}^2 + \text{Var}(E_D\hat{f} - f) + (E_D[(E_D\hat{f} - f)])^2 \\ &= \text{bias}^2 + \text{Var}(E_D\hat{f} - f) \end{aligned}$$

\rightarrow non training error or testing error.

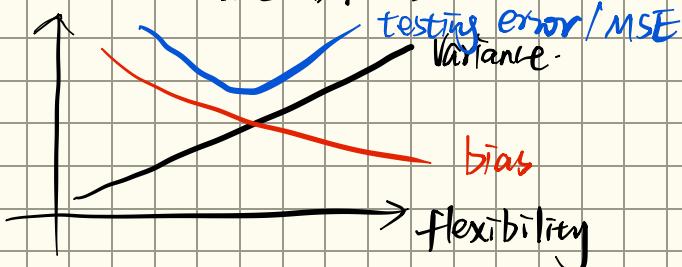
$$\frac{1}{n} \sum (y_i - \hat{f}(x_i))^2 \text{ empirical testing error.}$$

$$\begin{aligned} &E[(Y - \hat{f}(x))^2] \xrightarrow{Y \in \mathcal{X}^N} \text{expected testing error} \\ &\text{condition on the training data} \\ &\text{but } E_D = E_D E_{\text{CV}}[(Y - \hat{f}(x))^2 | D] = \text{bias}^2 + \text{variance} + \text{noise} \end{aligned}$$

unconditional one

in CV: un conditional

In AIC, BIC: conditional



overfit: small bias & large variance.

why ML flexible but small error?

Signal \gg noise overfit will not be a big issue.

ML 仍然会 overfit e.g. k-lines. noise very big.

事实上, 对于 bias 来说我们 不是 true fun. $f(x)$.
 \Rightarrow simulate.

interaction: effect of x_1 on Y depends on x_2 . (tree model)

$$Y = \beta_0 + \beta_1 x$$

$$\hat{\beta}_1 = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2} \quad \beta_0 = \bar{y} - \hat{\beta}_1 \bar{x}$$

$$= \max_i \left(\sum_{j=1}^n (y_j x_j^\top \beta - \log(1 + \exp(x_j^\top \beta))) \right) \quad \text{对 \beta 求导}$$

Another. ~~TFIDF~~ distribution.

$$P(y_i) = P_i^{y_i} (1-P_i)^{1-y_i}$$

$$P(D|\beta) = \prod_{i=1}^n P_i^{y_i} (1-P_i)^{1-y_i}$$

$$\log P(D|\beta) = \sum_{i=1}^n y_i \log P_i + (1-y_i) \log(1-P_i)$$

CROSS entropy
TEXT 和 PDI multiclass ~~normalize~~

$$P(y_i) = y_1 P_1 y_2 P_2 y_3 P_3$$

(\frac{1}{3} \text{ 的 } 10 \text{ 倍}) (\frac{1}{3}) non causal effect
B: joint effect and margin effect.

generative model: can generate data without anything given parameters

The logistic regression \hat{x}

discriminative model (generative model) ~~本应是 fit \mu \Sigma (f(x))~~

Linear: $X|Y \sim k \sim N(\mu_k, \Sigma_k)$ 利用 μ_k, Σ_k fit $P(Y|X)$, ~~本应是 max P~~

$$f(x) \log \frac{P(Y=1|x)}{P(Y=0|x)} = (\mu_k^\top - \mu_0^\top)x + b_0$$

TK, μ_k , how to estimate? linear: $S(\beta) = \mu_k^\top x - \frac{1}{2} \mu_k^\top \Sigma^{-1} \mu_k$ $L = \sum_{i=1}^n (y_i \beta^\top x_i - \log(1 + \exp(\beta^\top x_i)))$

$$\mu_k = \frac{1}{nk} \sum_{i=1}^n y_{ik} x_i \quad (\text{用小范围内数据})$$

$$\hat{\Sigma} = \frac{1}{nk} \sum_{i=1}^n \sum_{j=1}^k (x_{ij} - \mu_k)(x_{ij} - \mu_k)^\top \quad (\text{用所有数据})$$

$$\text{LDA 没有 numerical issue. } \exp > 0 \rightarrow x^\top \beta > 0 \rightarrow \text{numerical issue}$$

LR: neurons, optimal = ∞ . \Rightarrow not converge \Rightarrow logistic Reg 问题.
LDA 有数值问题.

sensitivity: TP specificity: TN.

AUC ROC Curve \uparrow 每个点对应 threshold
imbalance data: $FP = 1 - TN$ specificity

LDA (μ_k, Σ_k) ~~只限于单维~~
这个模型涉及 bias & variance trade off. 不取决于 assumption. 其他正确.

If enough data \Rightarrow variance is not important.

Validation set approach: AUC vs BIC . \uparrow 各自合适?
怎样选择数据?

Cross Validation $\left\{ \begin{array}{l} \text{leave one out: variance is large. no randomness} \\ \text{correlation很大, 导致 variance 很大, although bias is small} \end{array} \right. \quad \text{nearby small bias}$
 k -folders CV

CV 为什么 variance (因为 x_i 与 x_j uncorrelated), confidence interval \uparrow mean \downarrow confidence interval 不准确

用同一组选取 variable selection: training data: CV + selection \Rightarrow data snapping

但 PCA 没有用. \therefore PCA 不适用

naive classification: $X|Y \sim k \sim N(\mu_k, \Sigma_k)$ Σ_k is diagonal \Rightarrow independent

(LOOCV 很麻烦. 可以直接计算. Ridge 也能用) $[e_i] = \frac{e_i}{1 - e_i}$ if bias \propto $e_i / (1 - e_i)$
 $\Rightarrow g = Sg \Rightarrow$ similar works (Woodbury formula)

Variable selection: best subset.
forward.
backward.

$$AIC = -2 \log(\text{likelihood}) + 2g$$

$$BIC = -2 \log(\text{likelihood}) + 2 \log(n) d$$

for better prediction AIC > BIC

If enough data. BIC > AIC both use approximation, less time.

CV \uparrow BIC

最坏情况 all data 在 β 附近 ≈ 0 para

bootstrap: mean is different, variance structure is similar

$$\sum y_i x_i - \frac{\exp(x_i^\top \beta)}{1 + \exp(x_i^\top \beta)} \quad (\text{prior } n \times 1)$$

$$\sum x_i (y_i - p_i) = 0 = x^\top (y - p)$$

再求一阶导. $\frac{\partial}{\partial \beta} = \sum_i \frac{(x_i^\top \exp(-x_i^\top \beta)) (-x_i)}{(1 + \exp(-x_i^\top \beta))^2}$

$$H = \frac{\partial^2 L}{\partial \beta \partial \beta^\top} \quad \text{not compatible}$$

\Rightarrow Newton Method

$$\beta_{\text{new}} = \beta_{\text{old}} - H^{-1} \nabla L \quad | \quad \beta_{\text{old}} = -x^\top w$$

$$\min f(x). \quad g(x) = f'(x) = 0 \quad w = \text{diag} [P_i(1-P_i)]$$

$$= g(x_0) + g'(x_0)(x - x_0) = 0$$

$$x = -g'(x_0)^{-1} g(x_0) + x_0$$

$$\Rightarrow \text{应用. } L = \sum_{i=1}^n [y_i \beta^\top x_i - \log(1 + \exp(\beta^\top x_i))]$$

$$L = \sum_{i=1}^n (y_i - p_i) x_i = x^\top (y - p)$$

$$H = -x^\top w x$$

$$\beta_{\text{new}} = \beta_{\text{old}} - (x^\top w x)^{-1} x^\top (y - p)$$

reweighted LS. $\min \sum_{i=1}^n w_i (y_i - x_i^\top \beta)^2$

$$= (y - x\beta)^\top w (I - x^\top w)$$

$$\hat{\beta} = (x^\top w x)^{-1} x^\top w y$$

different: ① LDA assume normal error distribution

② LDA no numerical issue

every \rightarrow 1 ROC

difference: AUC. not assume model

loss fun: testing data comes

from distribution (eg. cross entropy)

sublinear: subgradient method

linear: gradient method + strongly convex

super linear: Newton's method

$\beta_{\text{new}} = \beta_{\text{old}} - H^{-1} \nabla L / \text{point}$

time-consuming

H 为对称 optimal s

so much quicker than gradient

best subset approach: ℓ_0 norm i forward: approximation for ℓ_0 norm.

Regularization: explicit: Ridge/Lasso

implicit: forward stepwise. variable selection.

/ gradient descent / XGBoost. / dropout \approx Ridge

Ridge: L_2 standardize is important. ① fit $\hat{\beta}$ ②入 $\hat{\beta}$ rescale

$$\mathbf{x}^T(\mathbf{x} + \lambda \mathbf{I})^{-1} \mathbf{y} = (\mathbf{x}^T \mathbf{x} + \lambda I)^{-1} \mathbf{x}^T \mathbf{y} \Leftarrow \mathbf{U}(\mathbf{D} + \lambda \mathbf{I})^{-1} \mathbf{W}^T \mathbf{x}^T \mathbf{y}$$

if # number of parameters

Lasso: L_1 variable selection. coordinate descent.

Convex + differentiable.

or convex + non-diff + separable.

soft shrinkage operator $\text{sign}(\hat{\beta})(1 - |\hat{\beta}|) +$

revised: ① inner product, track $y_i^T \Rightarrow \hat{\beta}^T \mathbf{r}_i$

② active set

③ warm start + active set

SNR小时. 效果优于其它

$$\hat{\beta}_{\lambda}^{\text{relax}}(x) = \mathbf{r} \cdot \hat{\beta}_n + (-\lambda) \hat{\beta}^{\text{LS}}$$

$$\mathbf{r}^T = \frac{1}{G^2} \sum_{j=1}^G \text{Cov}(y_j, \hat{y}_j)$$

$$S(\hat{\beta}, \lambda) = \arg \min \frac{1}{2} (\hat{\beta} - \beta)^T + \lambda \|\beta\|$$

$$\frac{\partial}{\partial \beta} = -(\hat{\beta} - \beta) + \lambda \text{sign}(\beta) = 0$$

$$\beta > 0. \quad \hat{\beta} + \beta + \lambda = 0$$

$$\hat{\beta} = (\hat{\beta} - \lambda) + \hat{\beta} > 0$$

$$\beta < 0. \quad -\hat{\beta} + \beta - \lambda = 0$$

$$\beta = \lambda - |\hat{\beta}|$$

$$\text{Lasso. } \frac{\partial P}{\partial \beta_j} = \sum_i (-x_{ij}(y_i - \hat{y}_i) + \lambda \text{sign}(\beta_j))$$

$$\Rightarrow \hat{\beta}_j = S\left(\sum_i x_{ij}(y_i - \hat{y}_i), \lambda\right).$$

standardize.

RF: if classification: voting
if regression: mean

out of Bagging error: it's about $\sqrt{\text{var}}$.

$$\sqrt{\frac{1}{n} \sum (y_i - \hat{f}_i)^2} \approx \text{LOOCV}$$

advantage: when training, already know LOOCV test error.

bagging: correlated.

Search P. paras

RF: $\frac{m}{P} = 0.1, 0.5, 1$. more search, off T.

search less.

e.g. best subset. most search

search m paras.

不是 bagging to. Signal/noise ratio

而是用 bagging / random to forward stepwise.

Variance ↑ bias ↑

approximation to

uncorrelated best subset.
to ↓ variance its variance is not
that small.

(1) 对所有数据领一个 w_i (2) (a) 以代数树. 对每棵树 $f_i(x)$

$$\text{loss function} = \sum_{i=1}^n w_i^{(m)} I(f_m(x_i) + y_i) \quad \text{损失函数为 } v$$

$$(b) \quad \hat{\varepsilon}_m = \frac{\sum_{i=1}^n w_i^{(m)} I(f_m(x_i) + y_i)}{\sum_{i=1}^n w_i}$$

$$\lambda_m = \log \left\{ \frac{1}{\hat{\varepsilon}_m} \right\}$$

$$w_i^{(m+1)} = w_i^{(m)} \exp \lambda_m I(f_m(x_i) + y_i)$$

(3) 最终预测: $F_m(x) = \text{sign}(\sum_{i=1}^m \lambda_i f_i(x))$ 只看前 m 次.

$$L(y, F) = \exp(-y \cdot F(x)) \quad F_m(x) = \frac{1}{2} \sum_{i=1}^m \lambda_i f_i(x)$$

在 N 步时构建的树. $F(x)$ 依赖于之前所有的 f_i , 显然不能这么写所以用归纳法. 假设已知 $\lambda_1, \dots, \lambda_{m-1}, F_1, \dots, F_{m-1}$.B. 由 F_m 得 λ_m 和 F_m .

$$L(y, F) \leq \exp(-y \cdot F_m(x))$$

$$\leq \exp(-y_i(F_{m-1}(x) + \frac{1}{2} \lambda_m f_m(x)))$$

$$\leq \exp(-y_i(F_{m-1}(x) + \frac{1}{2} y_i \lambda_m f_m(x))) \quad w_i^{(m)} = \exp(-y_i F_{m-1}(x_i))$$

$$\leq w_i^{(m)} \exp(-\frac{1}{2} y_i \lambda_m f_m(x)) \quad \text{即 } y_i f_m(x) =$$

$$\text{one fact: } \sum_{i \in I} w_i^{(m)} = \sum_{i \in I} w_i^{(m-1)} + \sum_{i \in M} w_i^{(m)}$$

$$L(y, F) = \sum_{i \in I} w_i^{(m)} \exp(-\frac{1}{2} \lambda_m) + \sum_{i \in M} w_i^{(m)} \exp(\frac{1}{2} \lambda_m)$$

$$= \exp(-\frac{1}{2} \lambda_m) \left(\sum_{i \in I} w_i^{(m-1)} - \sum_{i \in M} w_i^{(m-1)} \right) + \sum_{i \in M} w_i^{(m)} \exp(\frac{1}{2} \lambda_m)$$

$$= \exp(-\frac{1}{2} \lambda_m) \sum_{i \in I} w_i^{(m-1)} + \sum_{i \in M} w_i^{(m)} \left(\exp(\frac{1}{2} \lambda_m) - \exp(-\frac{1}{2} \lambda_m) \right)$$

令 $L(y, F)$ 为 0:

$$\begin{aligned} \text{对 } \lambda_m: & \sum_{i \in I} w_i^{(m)} \exp(-\frac{1}{2} \lambda_m) \cdot (-\frac{1}{2}) \\ & + \sum_{i \in M} w_i^{(m)} \left(\exp(\frac{1}{2} \lambda_m) + \frac{1}{2} \exp(-\frac{1}{2} \lambda_m) \right) = 0 \\ & - \sum_{i \in I} w_i^{(m)} + \sum_{i \in M} w_i^{(m)} \left(\exp(\lambda_m) + 1 \right) = 0 \end{aligned}$$

$$\lambda_m = \log \left(\frac{\sum_{i \in I} w_i^{(m)}}{\sum_{i \in M} w_i^{(m)}} \right)$$

$$\begin{aligned} \frac{1 - \hat{\varepsilon}_m}{\hat{\varepsilon}_m} &= 1 - \frac{\sum_{i=1}^n w_i^{(m)} I(f_m(x_i) + y_i)}{\sum_{i=1}^n w_i^{(m)}} \\ &> w_i^{(m)} \quad \sum_{i=1}^n w_i^{(m)} I(f_m(x_i) + y_i) \\ &\geq \sum_{i=1}^n w_i^{(m)} I(f_m(x_i) + y_i) \\ &\geq w_i^{(m)} \end{aligned}$$

(这个部分和 λ 无关, 不是 λ 的). 所以 y_i 和 f 会一起和 λ

$$\therefore \lambda_m = \log \frac{1 - \hat{\varepsilon}_m}{\hat{\varepsilon}_m}$$

这样证明 AdaBoost 其实就是 logistic regression.

$$\begin{aligned}
w_i^{(m)} &= \exp(-y_i f_m(s_i)) \\
w_i^{(m+1)} &= \exp(-y_i f_m(s_i)) \\
&= w_i^{(m)} \cdot \exp(-\frac{1}{n} \sum_{j=1}^n f_m(s_j) - \frac{1}{n} \sum_{j=1}^n y_j) \\
&= w_i^{(m)} \exp(-\frac{1}{n} \sum_{j=1}^n (f_m(s_j) + y_j)) \\
&= w_i^{(m)} \exp(-\frac{1}{n} \sum_{j=1}^n f_m(s_j) - \frac{1}{n} \sum_{j=1}^n y_j) \\
&= w_i^{(m)} \cdot \exp(-\frac{1}{n} \sum_{j=1}^n f_m(s_j) - \frac{1}{n} \sum_{j=1}^n y_j)
\end{aligned}$$

对于 m , 对所有的 n 来说, 都是 1 个常数, 所以忽略

$$\Rightarrow w_i^{(m+1)} = w_i^{(m)} \exp(-\frac{1}{n} \sum_{j=1}^n (f_m(s_j) + y_j))$$

为什么是 logistic?

$$E[y|f(s)] = \frac{1}{2} \int \exp(-y f(s)) p(y|f(s)) dy$$

$$\begin{aligned}
\text{对于 } y. \quad & \frac{d}{dy} E[y|f(s)] = \frac{1}{2} \int y \exp(-y f(s)) p(y|f(s)) dy \\
&= \int \exp(f(s)) p(y=f(s)) \exp(-f(s)) p(y=f(s)) dy = 0 \\
\Rightarrow \quad & f(s) = \frac{1}{2} \log \frac{p(y=1|s)}{p(y=0|s)} \rightarrow \text{logistic regression.}
\end{aligned}$$

Forward stagewise Additive Model

$$f(s) = \sum_{m=1}^M \beta_m b_m(s, y_m)$$

如果 b 是 树 那么 y_m 包含所有树的 parameter (包含从 m 分支). n estimators.

square loss:

$$L(y, f(s)) = (y - f(s))^2$$

$$\begin{aligned}
L(y_i, f_{m-1}(s) + \beta_m b_m(s_i, y_m)) &= (y_i - f_{m-1}(s) - \beta_m b_m(s_i, y_m))^2 \\
&= (y_i - \underbrace{\beta_m b_m(s_i, y_m)}_{\text{current residual}})^2
\end{aligned}$$

所以不用每次都对 β 求梯度.

$$\text{initial: } \hat{f}_0(s) = \arg \min_{f_0} \sum_{i=1}^n L(y_i, f_0)$$

Friedman's Gradient Boosting

iteration:

$$(1) r_i = y_i - \hat{f}_{m-1}(s_i) \quad | \quad \hat{f}_{m-1}(s_i) = \hat{f}_{m-1}(s_i)$$

相当之前的 r_i .

$$(2) \text{fit } g(s) \text{ 用 } r_i \text{ 来拟合 } r_i$$

即我们以新的 square loss 为 residual $\hat{r}_i = y_i - \hat{f}_0(s_i)$

$$(3) \hat{\beta} = \arg \min_{\beta} \sum_{i=1}^n L(y_i, \hat{f}_{m-1}(s_i) + \beta g(s_i))$$
即新的梯度的 loss function 为 L .

$$(4) \hat{f}_m(s) = \hat{f}_{m-1}(s) + \hat{\beta} g(s)$$

if L = squared loss function.

$$L(y_i, f) = \frac{1}{2} (y_i - f(s_i))^2$$

$$\frac{\partial L}{\partial f} = y_i - f(s_i) \quad | \quad f(s_i) = \hat{f}_{m-1}(s_i) + \hat{\beta} g(s_i)$$

还是 current residual

$$\beta^{\text{new}} = \beta^{\text{old}} + s(g) \quad | \quad \text{old} \cdot \text{梯度 descent}$$

$$\text{initial: } f_0(s), \text{ constant. } \hat{f}_0(s) = \arg \min_{f_0} \sum_{i=1}^n L(y_i, f_0)$$

Stochastic Gradient Boosting

$$(1) r_i = y_i - \hat{f}_{m-1}(s_i) \quad | \quad \hat{f}_{m-1}(s_i) = \hat{f}_{m-1}(s_i)$$

(2) 随机抽样 n 个 data = A

(3) 用 A 构建一个 regression tree. $A \rightarrow r_i$ (r_i 为 terminal nodes)

(4) 对每个 terminal node k 相应属于它 s_k . $\hat{f}_k(s_k) \rightarrow r_i$. $g(s_k) \geq f_k(s_k)$

$$\hat{\beta}_k = \arg \min_{\beta_k} \sum_{s \in A} L(y_i, \hat{f}_{m-1}(s_i) + \beta_k)$$

$$(5) \hat{f}_m(s) = \hat{f}_{m-1}(s) + \hat{\beta}_k g(s_k)$$

$\hat{\beta}_k \rightarrow g(s_k) \geq f_k(s_k)$

$$\text{loss: } -\frac{1}{n} \sum_{i=1}^n [y_i - f(s_i) - \log(\frac{1}{1 + \exp(f(s_i))})]$$

$$\text{initial: } f_0 = \log \frac{\sum y_i}{\sum (1 - y_i)}$$

P_i

$$\textcircled{1} \quad y_i = \frac{1}{1 + \exp(-f(x_i))} \quad \text{if } f(x_i) > 0 \quad \text{else } p_i$$

$$\textcircled{2} \quad \theta_{\min} = \arg \min \left[-g_{\min} - T(\theta_i; \theta) \right]^2$$

$$p_k = \frac{\exp(y_i p_i)}{\sum_i \exp(p_i)} \quad \text{why?}$$

$$H = \frac{\partial^2}{\partial p^2} = -\frac{2}{n} \sum_i p_i(1-p_i)$$

$$L - H'g = -\frac{\sum_i p_i(y_i - p_i)}{\sum_i p_i(1-p_i)}$$

α : learning rate. Shrinkage.

Lasso and Boosting.

$$\text{square loss} \quad \min_b \sum_i (y_i - b t(x_i))^2 = \min_b \sum_i (y_i^2 - 2x_i y_i b + x_i^2 b^2)$$

$$\frac{\partial^2}{\partial b^2} = 1 \rightarrow \max_j \left\{ \frac{x_i^2}{x_i t(x_i)} \right\}$$

$$x_i^T x_i (y_i - b t(x_i))$$

inner product with x_i and y_i .

(\times 轴) Lasso $\| \cdot \|_1$ regularization.

EM

$$\text{E-step: } E_z[\log P(z|x| \theta)]$$

$$\text{M-step: } Q(\theta, \theta^{old}). \quad E_{z|x, \theta^{old}} [\ln P(z|x| \theta)] = \sum_z \ln P(z|x| \theta) \cdot P(z|x, \theta^{old})$$

$$\text{and max } Q(\theta, \theta^{old})$$

$$\theta^{old} \leftarrow \theta^{new}$$

首先, the missing data 指的不是 nan, 而是不能直接观察到的数据。

因此需要找到 z , 我们设其为 indicator.

$$\text{假设 } z \text{ 有 } k \text{ 个值. } \quad P(z=k) = \pi_k$$

$$P(z_1|z=1) = N(\mu_1, \Sigma_1), \quad P(z_1|z=2) = N(\mu_2, \Sigma_2), \quad P(z_1|z=3) = N(\mu_3, \Sigma_3)$$

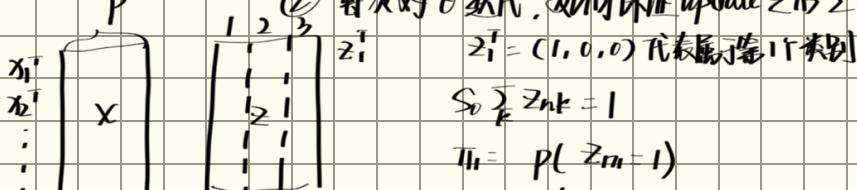
$$p(x) = \sum_z p(x|z) = \sum_z p(x|z) p(z) = \sum_k p(x|z_k, \mu_k, \Sigma_k) \cdot \pi_k.$$

$$p(z) = \prod_{k=1}^K \left(\sum_i p(z_k|z_i, \mu_k, \Sigma_k) \cdot \pi_k \right)$$

$$\log P(z) = \sum_{i=1}^n \log \left(\sum_k p(z_k|z_i, \mu_k, \Sigma_k) \cdot \pi_k \right) \quad \pi_k, \mu_k, \Sigma_k \text{ 是 } \theta^{old}.$$

Problem: ① 有 n 对 μ_k, Σ_k , π_k 手写很麻烦, 看数太多.

② 每次对 θ 迭代, 如何更新 $\pi_k, \sum_k \pi_k = 1$.



$$\sum_k \pi_k = 1$$

$$\pi_1 = p(z_m=1)$$

$$\pi_2 = p(z_m=1)$$

$$\pi_3 = p(z_m=1)$$

$$p(z_n) = \prod_{k=1}^3 \pi_k^{z_{nk}}$$

$$\text{So } p(z_n|z_n) = \prod_{k=1}^3 p(z_n|z_n, \mu_k, \Sigma_k)^{z_{nk}}$$

$$\log P(z_n|z_n) = \sum_{k=1}^3 z_{nk} \log p(z_n|z_n, \mu_k, \Sigma_k)$$

$$p(z|z, \theta) = \prod_{i=1}^n \prod_{k=1}^3 p(z_i|z_i, \mu_k, \Sigma_k)^{z_{ik}}$$

$$p(z_i|z, \theta) = \prod_{k=1}^3 \pi_k^{z_{ik}} p(z_i|z_i, \mu_k, \Sigma_k)^{z_{ik}}$$

$$= \prod_{k=1}^3 \pi_k^{z_{ik}} \left(p(z_i|z_i, \mu_k, \Sigma_k) \right)^{z_{ik}}$$

$$E_z[\log P(x_i \geq |z|)]$$

$$\begin{aligned} p(z|\theta) &= p(\theta|z)p(z) \\ &= \prod_{k=1}^K p(\pi_k|\theta_k)^{\pi_k} \cdot \prod_{k=1}^K \log \pi_k + \log p(\pi_k|\theta_k) \end{aligned}$$

$$\boxed{\sum_{k=1}^K \log p(x_i|z, \theta) \cdot p(z|\theta)}$$

$$= E_z \left[\sum_{i=1}^n \sum_{k=1}^K \pi_k (\log \pi_k + \log p(x_i|\mu_k, \Sigma_k)) \right]$$

old, $\pi_k \geq \pi_{k+1}$.

因此 E_z 时为常数

$$= E_z \left[\sum_{i=1}^n \sum_{k=1}^K \pi_k (\log \pi_k + \log p(x_i|\mu_k, \Sigma_k)) \right]$$

$$= E_z \left[\sum_{i=1}^n \sum_{k=1}^K p(z_{nk}=1|x_i, \theta^{old}) (\log \pi_k + \log p(x_i|\mu_k, \Sigma_k)) \right] \text{ 由 } E_z$$

可以理解为是 posterior mean

$$P(z_{nk}=1|x_i, \theta^{old}) = \frac{p(x_n|z_{nk}=1, \theta^{old})}{p(x_n|\theta^{old})} \text{ how to calculate}$$

$$= p(x_n|z_{nk}=1, \theta^{old}) / p(z_{nk}=1)$$

$$\stackrel{k}{\sum} p(x_n|z_{nk}=1, \theta^{old})$$

$$P(x_n|\mu_k \geq z_{nk}) \cdot \text{all old, can get}$$

$$\stackrel{j=1}{\sum} p(x_n|\mu_j, \Sigma_j) \pi_j \Rightarrow \text{第 } k \text{ 步 E-step}$$

$$\text{let } E[z_{nk}|x_i, \theta] = r(z_{nk})$$

$$\begin{aligned} \sum_{i=1}^n \sum_{k=1}^K r(z_{nk}) &= \sum_{i=1}^n \sum_{k=1}^K p(x_n|\mu_k, \Sigma_k) \pi_k \\ &\quad \left(\sum_{j=1}^J p(x_n|\mu_j, \Sigma_j) \pi_j \right) \end{aligned}$$

$$= \sum_{i=1}^n 1 = N.$$

$$\text{且 } \sum_{k=1}^K r(z_{nk}) = \sum_{k=1}^K E[z_{nk}|x_i, \theta] = 1$$

即对 x_i 每一行的期望为 1. 即 x_i 落在 K 类中的概率为 1. 合理.

然后考虑 M -step.

$$M\text{-step: } Q(\theta, \theta^{old}). E_{z|x_i, \theta^{old}} [\ln p(x_i \geq |z|)] = \sum_k \ln p(x_i \geq |z|) \cdot p(z_{nk}|\theta^{old})$$

and max $Q(\theta, \theta^{old})$

$$\theta^{old} < \theta^{new}$$

$$Q(\theta, \theta^{old}) = \sum_{i=1}^n \sum_{k=1}^K r(z_{nk}) (\log \pi_k + \log p(x_n|\mu_k, \Sigma_k)) \theta.$$

$$\max Q(\theta, \theta^{old}) \text{ with constraints } \sum_k \pi_k = 1$$

\Rightarrow lagrange

$$(1) \text{ 固定 } \pi_k \text{ 时. } \sum_{k=1}^K r(z_{nk}) \cdot \frac{1}{\pi_k} + \lambda = 0. \quad \lambda \neq 0.$$

固定 π_k .

$$\text{由 } \sum_{k=1}^K r(z_{nk}) = 1$$

$$\Rightarrow \sum_{k=1}^K \frac{1}{\pi_k} + \lambda = 0$$

$$\Rightarrow N + \sum_{k=1}^K \lambda \pi_k = 0. \quad \sum \pi_k = 1 \mid \lambda = N$$

$$\begin{aligned} \sum_{k=1}^K r(z_{nk}) - N = 0 \quad \pi_k = \frac{\sum_{k=1}^K r(z_{nk})}{N}. \end{aligned}$$

同样，对 μ_k 和 Σ_k 也使用同样的方法。但不难，观察

$$\sum_{n=1}^N \sum_{k=1}^K r(z_{nk}) (\underbrace{\log \pi_k}_{\text{fixed}} \cdot \log p(x_n | \mu_k, \Sigma_k))$$

$$\begin{aligned} \text{根据 } x_i \sim N(\mu, \Sigma) \Rightarrow \hat{\mu} = \frac{\sum x_i}{n} \quad \sum = \frac{\sum x_i}{n} \\ \Rightarrow \mu_k = \frac{\sum x_i r(z_{ik})}{\sum r(z_{ik})} \quad \Sigma_k = \frac{\sum r(z_{ik})(x_i - \mu_k)(x_i - \mu_k)^T}{\sum r(z_{ik})} \end{aligned}$$

$$\begin{aligned} \text{一个很重要的性质: } \ln p(x | \pi, \mu, \Sigma) &= \ln \left(\prod_{k=1}^K \pi_k \cdot p(x_n | \mu_k, \Sigma_k) \right) \quad \text{区间不相加了.} \\ &\uparrow \\ &= \sum_{i=1}^N \ln \sum_{k=1}^K \pi_k p(x_i | \mu_k, \Sigma_k) \end{aligned}$$

单个递增. \Rightarrow 收敛 ZM. stable

why?

$$\begin{aligned} \ln p(x | \pi, \mu, \Sigma) &= L + KL \\ L(q, \theta) &= \sum z q(z) \log p(x, z | \theta) = \sum z q(z) \log q(z) \\ KL(q || p) &= - \sum q(z) \log p(z | x, \theta) = \sum q(z) \log q(z) \end{aligned}$$

$$\begin{aligned} L + KL &= \sum q(z) (\log p(x, z | \theta) - \log(p(z | x, \theta))) \\ &= \sum q(z) \log p(x | \theta) \\ &= \log p(x | \theta) \quad \text{分解完成.} \end{aligned}$$

由下 $KL \geq 0$. Jensen's equality: $f(E(x)) \leq E[f(x)]$ for convex function

$$\begin{aligned} KL &= \sum q(z) \log p(z | x, \theta) + \sum q(z) \log \frac{p(z)}{q(z)} \geq 0 \\ f &= -\ln s \quad X = \frac{p(z)}{q(z)} \\ &= -\sum q(z) \log \frac{p(z)}{q(z)} = E[f(x)] \\ &\Rightarrow -\ln \cdot \left(\frac{p(z)}{q(z)} \cdot q(z) \right) = 0 \end{aligned}$$

$\therefore L \leq \log p(x | \theta)$ for lower bound.

at E step. $\ln p(x | \theta^{old}) = L(q, \theta^{old}) \Leftrightarrow q(z) = p(z | x, \theta^{old})$

at M step. $L(q, \theta^{old}) \leq L(q, \theta^{new}) = \ln p(x | \theta^{new})$

递增.