# Math 基础：矩阵求导

$x, y$: scalar / m-vec; $x$: scalar / n-vec.

$$\frac{\partial y}{\partial x} = \left(\frac{\partial y}{\partial x_i}\right), \quad \frac{\partial y}{\partial x} = \left(\frac{\partial y_i}{\partial x_j}\right), \quad \left(\frac{\partial y}{\partial x}\right)_{ij} = \frac{\partial y_i}{\partial x_j}$$

$$\frac{\partial y}{\partial x} = \begin{bmatrix} \frac{\partial y}{\partial x_1} \\ \vdots \\ \frac{\partial y}{\partial x_n} \end{bmatrix}, \quad \frac{\partial y}{\partial x} = \begin{bmatrix} \frac{\partial y_1}{\partial x_1} & \cdots & \frac{\partial y_m}{\partial x_1} \\ \vdots & & \vdots \\ \frac{\partial y_1}{\partial x_n} & \cdots & \frac{\partial y_m}{\partial x_n} \end{bmatrix}$$

① $\frac{\partial x^T A}{\partial x} = \frac{\partial x^T \cdot a}{\partial x} = a$ **必重要计算！**

② $\frac{\partial x^T A x}{\partial x} = (A + A^T) x$

Lagrange: 对于优化问题：

$$\min f_0(x), \quad s.t.: \quad f_i(x) \le 0, i=1,\cdots,m; \quad h_i(x)=0, i=1,\cdots,n$$

则：$L(x, \lambda, \nu) = f_0(x) + \sum_{i=1}^{m}\lambda_i f_i(x) + \sum_{j=1}^{n}\nu_j h_j(x)$

其中 $\lambda$ 和 $\nu$ 是拉格朗日乘子，$\lambda_i \ge 0$, $\nu_i$ 无约束

**KKT条件**：
$$\begin{cases} f_i(x) \le 0, i=1,\cdots,m \\ h_i(x)=0, i=1,\cdots,n \end{cases}, \quad \lambda_i \ge 0$$
且 $\lambda_i f_i(x) = 0, \; i=1,\cdots,m, \quad \nabla_x L(x,\lambda,\nu)=0$

**奇异值分解**：$A = U\Sigma V^T$, $U \in \mathbb{R}^{m\times n}$, $UU^T = I$;
$V \in \mathbb{R}^{n\times n}$, $V^TV = I$; $\Sigma \in \mathbb{R}^{m\times n}$, $(\Sigma)_{ii} = \sigma_i$, 且 $\sigma_i$
是排序负数且满足：$\sigma_1 \ge \sigma_2 \cdots \ge 0$

$A^TAx = \lambda x$, $|A^TA - \lambda I| = 0$
$\sigma_i = \sqrt{\lambda_i}$, 而 $V$ 中每行列向量是 $A^TA$ 的 eigen vector
$V$ 中每个列向量为 $AA^T$ 中 eigenvector

若 Data Linearly Unseperable, 则 Network 且有 overfit

**GD**: $x^{k+1} \leftarrow x^k - d_k \nabla f(x^k)$
**MLE**: $\hat\theta = \underset{\theta}{argmax}\ P(D|\theta)$
**MAP**: $\hat\theta = \underset{\theta}{argmax}\ P(\theta|D) \propto \underset{\theta}{argmax}\ P(\theta) \cdot P(D|\theta)$
**GD**: $\log p(x) = \int q_{(z)} \log \frac{p(x,z)}{q_{(z)}} dz + \int KL(q_{(z)} \| p(z|x))]$
$\log p(x) = \int q_{(z)} \log \frac{p(x,z)}{q_{(z)}} dz + KL(q_{(z)} \| p(z|x))]$ → **ELBO** $\ge 0$

---

$$H(X) = -\sum_{x\in X} P(x)\log P(x) = E[\log \frac{1}{P(x)}]$$
$$H(Y|X) = -\sum_{x\in X} P(x)\sum_{y\in Y} P(y|x)\log P(y|x)$$
$$= \sum_{x\in X}\sum_{y\in Y} P(x,y)\log \frac{1}{P(y|x)}$$
$$I(X;Y) = \sum_{x,y} P(x,y)\log \frac{P(x,y)}{P(x)P(y)}$$

① $H(X) = ①+②$   $H(Y) = ②+③$
② ③ $H(X|Y) = ①$  $H(Y|X) = ③$  $I(X;Y) = ②$
$H(X,Y) = ①+②+③$
$H(X|Y) \cdot ① \; H(Y|X) \cdot ③ \; I(X;Y) \cdot ②$

$H(X)-H(X|Y)$ or: $H(Y)-H(Y|X)$

**决策树**：(离散) 归纳地取出属性 (metric: I: Gini: Gini error)
以有可解释性，易 overfit (可用 pruning 剪枝), 对数据敏感
连值敏感, inductive bias

**KNN**: Lazy learning $O(\nu)$; Predict: $O(MN)$
对 noise, K值, 度量 metric 敏感, 对数据敏感, 易 overfit; 易用于分类&回归; 表现在理论观与取
量关系; 易 overfit; $\dim$ 样本数

**Perceptron**: $proj_b a = \frac{|a^Tb|}{\|b\|_2} \cdot \frac{b}{\|b\|_2}$
初: $w=0\; b=0$, $\hat{y} = sign(w^Tx+b) = \begin{cases} -1, & else \\ +1, & w^Tx+b \ge 0 \end{cases}$
$w \leftarrow w + y^{(t)} x^{(t)}$  $b \leftarrow b + y^{(t)}$ if wrong $\hat y \ne y$
$\Rightarrow \theta = [b], \; \theta \leftarrow \theta + y^{(t)} x^{(t)}$ if $\hat y \ne y$

在上述中, $w$ 可视为 $\hat x$ 的线性组合 $\Rightarrow$ 有 overfitting
mistake $\le (R/\gamma)^2$, $\gamma$ 为 margin, $\|x^{(t)}\| \le R$
$\theta^*$ 为 $x$ 表达平面函数且 $\|\theta^*\|=1$, $|\theta^* x^{(t)}| \ge \gamma$

**SVM**: Formulation, input: $S = \{(x_1,y_1),\cdots,(x_m,y_m)\}$
Find $\underset{w,\hat u}{argmin}\ \|w\|^2 + C\sum_{i=1}^{m}\hat\xi_i$, $s.t.$:

Dual: $\underset{\alpha}{argmin}\ \frac{1}{2}\sum_{i,j}\hat\xi_i \cdot y_i y_j a_i a_j x_i \cdot x_j - \sum_i a_i$
$s.t., \quad 0 \le a_i \le C_i, \; \sum y_i a_i = 0$
points $(x_i, y_i)$ whose $a_i \ne 0$: support vector
or: $y_i (w x_i) = 1$: final: $w = \sum_i a_i y_i x_i$

$\forall i, \quad y_i(w x_i) \ge 1 - \hat\xi_i, \; \hat\xi_i \ge 0$

---

**Kernel**: 将样本映射至更高维 $h$ 维 solve 非线性可分
**问题**, 但也将组在高维中有存在, 可用 Kernel
数据在低维组中的操作得到！ 计算量大！即

$$K(x_i, x_j) = \langle \phi(x_i), \phi(x_j) \rangle = \phi(x_i)^T \phi(x_j)$$
$K_{ij} = K(x_i, x_j)$  ① 对称 $\Delta: \forall x \ne 0$, $x^T K x \ge 0$
* $K_{ij} = K(x_i, x_j)$, 若 $k_1, k_2$ 核函数, 则 $r_1 k_1 + r_2 k_2$ 也是, $g(x) K(x_a, x_b) g$
-tanh$(\beta x_i^T x_j + \theta)^d$, $\exp\left(-\frac{\|x-x'\|^2}{2\sigma^2}\right)$; $\exp\left(-\frac{\|x-x'\|}{\sigma}\right)$;
$\beta x_i^T x_j + \theta) \to$ sigmoid  gauss  laplace

**Linear Regression**: $\underset{w,b}{argmin} \sum_{i=1}^{m}(y_i - w x_i - b)^2$
闭式解: $w = \frac{\sum_{i=1}^m x_i^2 - \frac{1}{m}(\sum_{i=1}^m x_i)^2}{\cdots}$, $b = \frac{1}{m}\sum_{i=1}^m(y_i - w x_i)$
$w^* = \underset{w}{argmin}\ [(y - Xw)^T(y-Xw)], w^* = (X^TX)^{-1}X^Ty$

**Logistic Regression** → 逻辑回归, 但依然为线性!

$$y = \frac{1}{1+e^{-(w^Tx+b)}}, \quad \ln\frac{P(y=1|x)}{P(y=0|x)} = w^Tx + b$$
$P(y_i | x_i; \beta) = y_i \cdot P_1(x_i; \beta) + (1-y_i) P_0(x_i; \beta)$.
$P_1(x_i; \beta) = \frac{e^{\beta^Tx_i}}{1+e^{\beta^Tx_i}}$; $P_0(x_i; \beta) = \frac{1}{1+e^{\beta^Tx_i}}$
$\mathcal{L}(\beta) = \sum_{i=1}^m (-y_i \beta^T x_i + \ln(1+e^{\beta^Tx}))$
$\nabla_\beta \mathcal{L}(\beta) = \sum_{i=1}^m x_i (P(Y=1|x_i, \beta) - y_i)$ (有时 $x_i$)

**必** $\theta_k \leftarrow \theta_k + \lambda (h_\theta(x^{(i)}) - y^{(i)}) x_k^{(i)}$

**GD&BP**: GD: $\theta \leftarrow \theta - \gamma \nabla_\theta J(\theta)$,
**SGD**: $\theta \leftarrow \theta - \gamma \nabla_\theta J^{(i)}(\theta)$  ($i$ are shuffled)
GD $\gamma$ 沿动个 (曲线); SGD 代 GD 易曲折, 朝
Minibatch: 取 $K$ 个样本子集作 GD. $S$ 介于
**BP**: Basic: $f_L(w_L, f_{L-1}(\cdots f_1(w_1, x))\cdots)$, 对每导数
$\varphi(z) = \varphi(z)(1-\varphi(z))$

$w^L$ → target

$$\frac{\partial L}{\partial w^L} = \sum_s s^L \xleftarrow x^L \quad \text{Loss:} f_0(u)^2 \sx^2$$

$$\to x^{L-1} \xleftarrow{} x^L$$

$$\frac{\partial L}{\partial w^L} = \frac{\partial L}{\partial x^L} \cdot \frac{\partial x^L}{\partial s^L} \cdot \frac{\partial s^L}{\partial w^L} , \quad \text{链式法则}$$

**Recom Sys & Matrix Factorization**

Unconstrained: $R \in \mathbb{R}^{m\times n}$, $R$'s rand $k \ll \min(m,n)$.

$U \in \mathbb{R}^{m\times k}$, $V \in \mathbb{R}^{n\times k}$, s.t., $R = UV^T$

$f(U,V) = \frac{1}{2}\|R - UV^T\|_F^2 = \frac{1}{2}\sum_{i=1}^{m}\sum_{j=1}^{n}(R_{ij} - U_{i\cdot} \cdot V_{j\cdot}^T)^2$

where $R_{ij}$ is known. 引入 $\frac{\lambda}{2}(\|U_{i}\|^2 + \|V_{j}\|^2)$ 正则化

$\nabla U_j: J_{ij}(U,V) = -E_{ij}V_{\cdot,i} + \lambda U_j$.

$\nabla V_{\cdot,i}: J_{ij}(U,V) = -E_{ij}U_{\cdot,i} + \lambda V_j$. (SGD)

$U_{j\cdot} \leftarrow U_{j\cdot} - \eta \nabla U_j; \quad V_{\cdot,i} \leftarrow V_{\cdot,i} - \eta (\nabla V_{\cdot,i})^T$

Constrained: $U, V = \arg\min \frac{1}{2}\|R - UV^T\|_F^2$ s.t.

$U_{ij} \geq 0$, $V_{ij} \geq 0$

Or: Fix $V^T$, 解 $U$; Fix $U$, 解 $V^T$, 循环

若 $R$ 为 SVD $R = Q\Sigma P^T$, 则可 truncate $Q \to k$.

$U \cong Q_k \Sigma_k$, $V \cong P_k$, 并作低秩近似

$\triangle$: $k = \log_2 d$; $k=1$:随机; $k=d$:穷举搜索 learn

**Bagging**: for $t = 1, \dots, T$ do

  for $s = 1, \dots, S$ do

    $i \sim \text{Uniform}(1, \dots, N)$

    $S_t = \{(x^{(is)}, y^{(is)})\}_{s=1}^{S}$

  $h_t = \text{train}(S_t)$

return $h(x) = \text{aggregate}^*(h_1, \dots, h_T)$

$*$:分类集成,为 majority vote;回归时,average

**Random Forest**: 抽 $n$ 个 sample(数 $\leq$ 重复低)构成每个

feature Bagging:随机从$S_t$ feature每个样本抽

未注意 $S_t$, $i^{\text{th}}$; 每个 $h_t$, $D_t$重新低样本$R$,

$\triangle$: Bootstrapping: 抽似样本(有放回)$\leq 36.8\%$没采样到

表现不错, 且 收敛性与 Bagging 相似 $n$ $\lim_{m\to\infty}(1-\frac{1}{m})^m = \frac{1}{e}$

**K-means**: Initialize $c = \{c_1, \dots, c_k\}$.

Repeat:   for $i$ in $\{1, \dots, N\}$

  $z^{(i)} \leftarrow \arg\min_j \|x^{(i)} - c_j\|_2^2$

  for $j$ in $\{1, \dots, k\}$

    $c_j \leftarrow \arg\min \sum_{i: z^{(i)} = j}\|x^{(i)} - c_j\|_2^2$

Converge   $\Rightarrow c_j = \frac{1}{|C_j|}\sum_{i \in C_j} x$

$*$: 可随机, 也可 **Furthest Point Heuristic** c,也可:

  先任选 $c_1$, Pick $c_j$ from 分布:

  $P(C_j) = \pi^{(i)} \propto \min_{j' < j}\|x^{(i)} - c_{j'}\|^2$, $x^{(i)}$ are not picked

**EM**: E: $Q(\theta|\theta^t) = \sum_z P(z|x, \theta^t)\log P(x, z|\theta)$

M: $\theta^{t+1} = \arg\max_\theta [Q(\theta|\theta^t) + \log P(\theta)]$ 有:MAP  不有:MLE

**GMM**: 假设 $x_i$ 生成: $k$ 个高斯 $h$ $\pi_k$ 概率 随机

$\to$ 考虑 $x_i$ 生成数据: Invariance; 输出不变; 不移动

$\Rightarrow$ E step: $\gamma_{ik} = P(z_i = k | x_i, \theta)$

$$\text{例 }\log p(x, z|\theta) = \sum_i \log \pi_{z_i} + \log N(x_i | \mu_{z_i}, \Sigma_{z_i})$$

$$\gamma_{ik} = P(z_i = k) = \frac{P(z_i = k)\pi_{z_i}}{P(x_i|\theta^{old})} = \frac{\pi_k^{old} N(x_i|\mu_k, \Sigma_k^{old})}{\sum_{j=1}^{k}\pi_j^{old} N(x_i|\mu_j^{old}, \Sigma_j^{old})}$$

$= \sum_{j=1}^{k}\pi_j N(x_i|\mu_k, \Sigma_k)$

M: $Q(\theta, \theta^{old}) = \sum_i \sum_{k=1}^{K}\gamma_{ik}[\log\pi_k + \log N(x_i|\mu_k, \Sigma_k)] + \lambda(1 - \sum_{k=1}^{k}\pi_k)$

且: $\pi_k^{new} = \frac{\sum_{i=1}^{N}\gamma_{ik}}{N}$, $\mathcal{L} = Q(\theta, \theta^{old})$

$$\mu_k^{new} = \frac{\sum_{i=1}^{N}\gamma_{ik} x_i}{\sum_{i=1}^{N}\gamma_{ik}}$$

$$\Sigma_k^{new} = \frac{\sum_{i=1}^{N}\gamma_{ik}(x_i - \mu_k^{new})(x_i - \mu_k^{new})^T}{\sum_{i=1}^{N}\gamma_{ik}}$$

**PAC**: 主成分分解:

① $\bar{x}_i \leftarrow x_i - \frac{1}{m}\sum_{i=1}^{m}x_i; \quad X \odot XX^T$ 计算方差

② 取前大 $d'$ eigenvalue 对应

的 eigenvector, 组成 $\omega^* \in \mathbb{R}^{d\times d'}$.

则 $\omega^{*T} x$ 后 进入 低维

**Adaboost** 6-Boosting: $h_t(x)$有低权/长 (**red**)

for $t = 1, \dots, T$ do $h_t \leftarrow T$ 上学习器个数

choose $h_t$ (一般:weak) (and train)

  算 error rate $\varepsilon \in (\alpha \sim D_t)$

  if $\varepsilon > 0.5$, break (太差!!)

  $\alpha_t = \frac{1}{2}\ln(\frac{1-\varepsilon_t}{\varepsilon_t})$

$$D_{t+1(\alpha)} = \frac{D_{t(\alpha)}}{Z_t^*} \times \begin{cases} \exp(-\alpha_t), \text{if 预测对} \\ \exp(\alpha_t), \text{if 错} \end{cases}$$

$$= \frac{D_{t(\alpha)}}{Z_t}\exp(-\alpha_t h_t(x)f_t(x))$$

end for

$h(x) = \text{sign}(\sum_t \alpha_t h_t(x)) \Rightarrow$ 只能作一维

**Boosting**: Goal: $\hat{h}(x) = \text{sign}(\sum_t \alpha_t h_t(x))$

$$Z_t = \sum_{s} e^{\alpha_t \varepsilon_t} + e^{-\alpha_t}(1-\varepsilon_t) , \quad \alpha_t = \frac{1}{2}\ln\frac{1-\varepsilon_t}{\varepsilon_t}$$

$= 2\sqrt{\varepsilon_t(1-\varepsilon_t)}$

$\triangle$: 对无法提升学校样本的基准法。可重采样

**CNN**: 注重于 local connectivity; parameter

sharing; pooling/subsampling hidden units

老依不: kernel 数 过少 参为 $i$; 考察 $(K, K, Cin)$ 中 **没有 bias**

Input: $(H_{in}, W_{in}, Cin)$, 输出 $(H_{out}, W_{out}, Cout)$.

stride $S$, 则 输出 $(H_{out}, W_{out}, Cout)$:

$$\begin{cases} H_{out} = \lfloor\frac{H_{in} + 2P - K}{S}\rfloor + 1, & \text{(Para 才有 bias = } C_{out}) \\ C_{out} = N \\ W_{out} = \lfloor\frac{W_{in} + 2P - K}{S}\rfloor + 1 \end{cases}$$

每个 channel 的形状都是相同的而非单独的

Pooling layers: e.g. Max pooling (**no learnable参数**)

CNN两个关键性质: Invariance; 输出不变 $!$ 不移动

希望 $x$ 平移后经过卷积 → 池化 (e.g. Conv + Pooling $\Rightarrow$ 变)

→ 考察 (e.g. Conv unit 移)

SVM推导：约束 $w^Tx+b=-1, (w)^Tx+b=1$，则构造问题

$\gamma = \frac{1}{\|w\|}$, $d = \frac{2}{\|w\|}$, 则构造问题

$\min_{w,b} \frac{1}{2}\|w\|^2$, s.t. $y_i(w^Tx_i+b)\geq 1$

$L(w,b,\alpha) = \frac{1}{2}\|w\|^2 + \sum_{i=1}^m \alpha_i(1-y_i(w^Tx_i+b))$

$\alpha = (\alpha_1; \alpha_2; \cdots; \alpha_m)$, 则KKT条件：

① primal: $y_i(w^Tx_i+b)\geq 1$   ② dual: $\alpha_i \geq 0$

③ Complementary: $\alpha_i(y_i(w^Tx_i+b)-1)=0$

④ Stationary: $\nabla_w L = 0$

对 $w$ 偏导: $w = \sum_{i=1}^m \alpha_i y_i x_i$, 代回:

$b$: $\sum_{i=1}^m \alpha_i y_i = 0$

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

s.t. $\sum_{i=1}^m \alpha_i y_i = 0$, $\alpha_i \geq 0$, 解出 $\alpha$ 后:

$f(x) = w^Tx+b = \sum_{i=1}^m \alpha_i y_i x_i^T x + b$

Why biased? $E[\hat\sigma^2_{MLE}] = E[\frac{1}{n}\sum(x_i-\hat\mu_{MLE})^2]$

则 $E[\hat\sigma^2_{MLE}] = \frac{n-1}{n}\sigma^2$, bias $= -\frac{\sigma^2}{n}$

Probably Approximate Correct learning:

$P(|\theta - \theta^*|\geq\epsilon)\leq 2e^{-2n\epsilon^2}\leq\delta,\ n\geq\frac{\ln(2/\delta)}{2\epsilon^2}$

K-mean++: achieves a O(logk) approximation to optimal clustering; 这种思路是对K法的改进

ELBO: $\log p(X) = \log\frac{p(X,z)}{q(z)}\cdot\frac{q(z)}{p(z|X)}\cdot p(X)$

$= [\log\frac{p(X,z)}{q(z)}]\cdot\frac{q(z)}{p(z|X)}\cdot p(X)$

$= \int q(z)\log\frac{p(X,z)}{q(z)}dz + \int q(z)\log\frac{q(z)}{p(z|X)}dz$

---

$= \int q(z)\log\frac{p(X,z)}{q(z)}dz + KL(q(z)\|p(z|X))$

ELBO $= \int q(z)\log p(X,z)dz - \int q(z)\log q(z)dz$  ← ELBO 先验 ≥ 0

$= E_{q(z)}[\log p(X,z)] + H(q(z))$

Perceptron mistake bound: $(R/\gamma)^2$

Proof: 设 $\theta^*$ 为最大间隔下的参数, 且 $\|\theta^*\|=1$

则 $\theta^{(k+1)}$, $\theta^{(k+1)} = \theta^{(k)} + y^{(i)}x^{(i)}$

且 $\forall i$, $y^{(i)}\cdot(\theta^*\cdot x^{(i)})\geq\gamma$

$\theta^{(k+1)}\cdot\theta^* \geq \theta^{(k)}\cdot\theta^* + \gamma$, 又 $\theta^{(0)}$ 正常初始化

$\Rightarrow \|\theta^{(k+1)}\|\geq k\gamma$

同时 $\|\theta^{(k+1)}\|_2^2 = \|\theta^{(k)} + y^{(i)}x^{(i)}\|_2^2$

$= \|\theta^{(k)}\|_2^2 + y^{(i)2}\|x^{(i)}\|_2^2 + 2y^{(i)}(\theta^{(k)}\cdot x^{(i)})$  $<0$, 分母点的:

$\leq \|\theta^{(k)}\|_2^2 + \|x^{(i)}\|_2^2$

$\leq \|\theta^{(k)}\|_2^2 \leq kR^2$, $k\leq\frac{R^2}{\gamma^2}$

$\therefore k\gamma^2 \leq \|\theta^{(k+1)}\|_2^2 \leq kR^2$, $k\leq\frac{R^2}{\gamma^2}$

Logistics: $p(y_i|x_i; w,b) = y_i p_i(x_i; \beta) + (1-y_i)p_0(x_i;\beta)$

$\mathcal{L}(\beta) = \sum_{i=1}^m (-y_i\beta^Tx_i + \ln(1+e^{\beta^Tx_i}))$

$\frac{\partial\mathcal{L}(\beta)}{\partial\beta} = -\sum_{i=1}^m x_i(y_i - p_1(\hat x_i;\beta))$

Adaboost: $\epsilon = \frac{1}{n}\sum_{i=1}^n \mathbb{I}(y_i\neq H_{final}(x_i))$

$\leq \frac{1}{n}\sum_{i=1}^n \exp(-y_i(\sum_{t=1}^T \alpha_t h_t(x_i)))$

$D_{T+1(i)} = \frac{D_{T(i)}}{Z_T}\exp(-\alpha_T y_i h_T(x_i))$

$D_{T(i)} = \frac{D_{T-1(i)}}{Z_{T-1}}\exp(-\alpha_{T-1}y_i h_{T-1}(x_i))$

$\vdots$

$D_{1(i)} = 1/n$

$\therefore \frac{1}{n}\prod_{t=1}^T Z_t\exp(-y_i\sum_{t=1}^T \alpha_t h_t(x_i)) = \frac{\sum_{i=1}^n}{\prod_{t=1}^T D_{T+1}(i)}$

---

$\therefore \epsilon \leq \prod_{t=1}^T Z_t$, 又 $Z_t = e^{\alpha t}\epsilon_t + e^{-\alpha t}(1-\epsilon_t)$

$= \epsilon_t e^{\alpha t} + (1-\epsilon_t)e^{-\alpha t}$

$\frac{\partial Z_t}{\partial\alpha_t} = \epsilon_t e^{\alpha t} - (1-\epsilon_t)e^{-\alpha t} \Rightarrow \alpha_t = \frac{1}{2}\log\frac{1-\epsilon_t}{\epsilon_t}$

$\frac{\partial^2 Z_t}{\partial\alpha_t^2} = \epsilon_t e^{\alpha t} + (1-\epsilon_t)e^{-\alpha t}\Big|_{\alpha_t=\frac{1}{2}\log\frac{1-\epsilon_t}{\epsilon_t}} = 2\sqrt{\epsilon_t(1-\epsilon_t)} > 0$

$H(x) = \sum_t \alpha_t h_t(x)$

求 $\frac{\partial\ell}{\partial H(x)}$: $\ell = \exp(-H(x))p(f(x)=1|x) + e^{H(x)}p(f(x)=-1|x)$

$\frac{\partial\ell}{\partial H(x)} = -e^{-H(x)}p(f(x)=1|x) + e^{H(x)}p(f(x)=-1|x) = 0$

$\Rightarrow H(x) = \frac{1}{2}\ln\frac{p(f(x)=1|x)}{p(f(x)=-1|x)}$