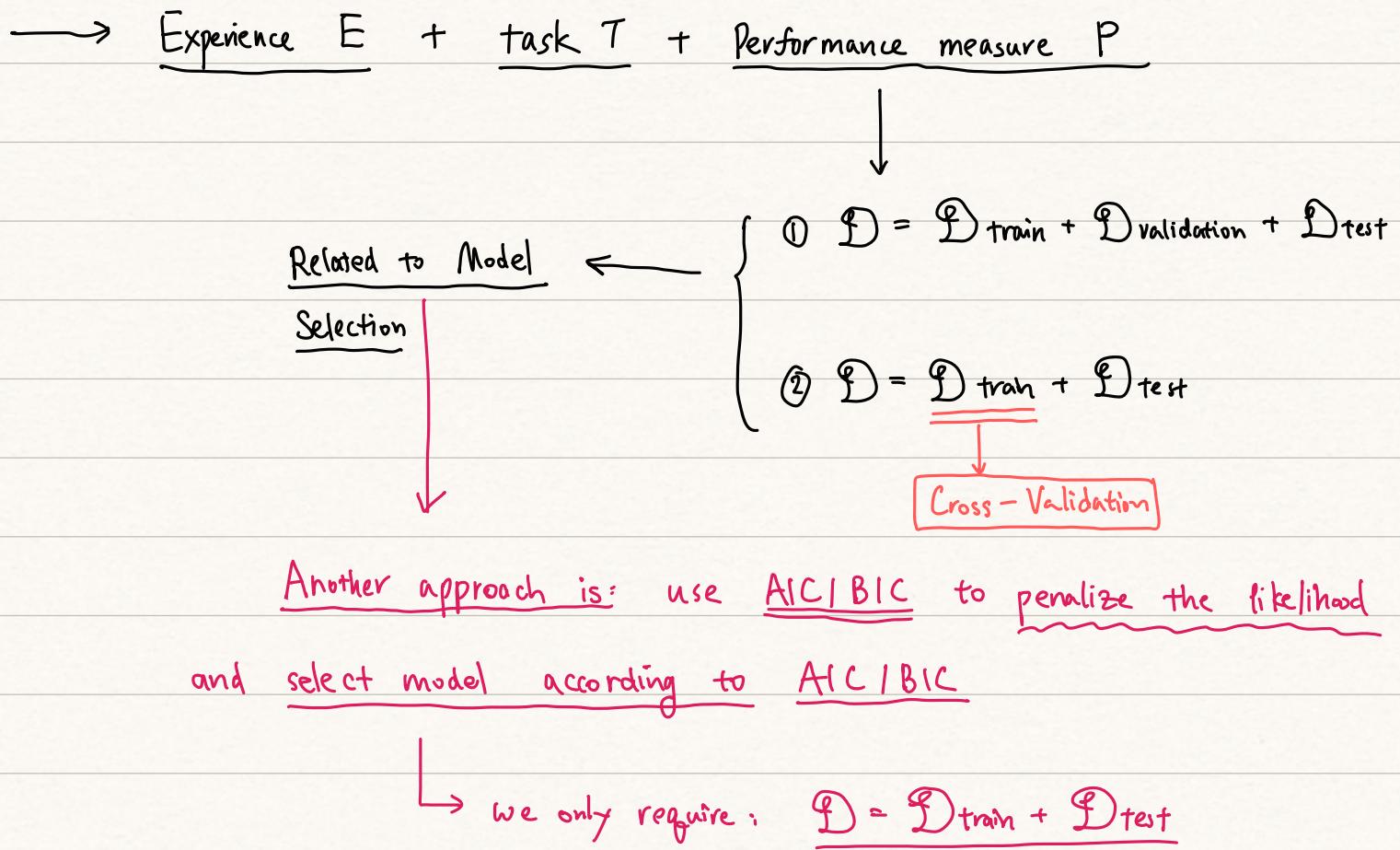


Recap 5105

1. Learning Framework



2. Supervised Learning



① Assumption: (strong form) deterministic

y_i is determined by x_i according to $y_i = f^*(x_i)$

$$f^* \leftrightarrow \text{oracle}$$

↑

$y_i = f^*(x_i)$

(weak form) stochastic

$$(x, y) \sim P^*$$

$$P^* \leftrightarrow \text{'Oracle'}$$



y_i is a sample from $P^*(\cdot | x_i) \Leftrightarrow \underline{y \sim P^*(\cdot | x)}$

② Population Risk Minimization Versus Empirical Risk Minimization

[ERM] → attainable

$$\hat{f} \leftarrow \min_{f \in \mathcal{H}} \hat{R}_n(f, f^*) = \frac{1}{n} \sum_{i=1}^n L(f(x_i), f^*(x_i)) \quad x_i \sim \mu$$

[PRM] → unattainable

$$\tilde{f} \leftarrow \min_{f \in \mathcal{H}} R(f, f^*) = E_{x \sim \mu} [L(f(x), f^*(x))]$$

Decomposition

$$\begin{aligned} & R(\hat{f}, f^*) - R(\tilde{f}, f^*) \\ &= \underline{R(\hat{f}, f^*) - \hat{R}_n(\hat{f}, f^*)} \rightarrow \text{uniform bound} \\ &\quad + \underline{\hat{R}_n(\hat{f}, f^*) - \hat{R}_n(\tilde{f}, f^*)} \rightarrow \text{optimization} \\ &\quad + \underline{\hat{R}_n(\tilde{f}, f^*) - R(\tilde{f}, f^*)} \rightarrow \text{point-wise bound} \end{aligned}$$

③ Three Paradigms of Supervised Learning

$$\begin{array}{c} \text{Approximation} \quad + \quad \text{Optimization} \quad + \quad \text{Generalization} \\ \hline f^* \leftrightarrow \hat{f} \qquad \qquad \hat{f} \qquad \qquad \hat{f} \leftrightarrow f^* \end{array}$$

3. Linear Models

① Linear Regression

1. Model: $\mathcal{H} = \{f: f(x) = w^\top x + b\}$

2. Optimization: $\min_w \frac{1}{N} \sum_{i=1}^N \|y_i - w^\top x_i\|_2^2$

→ ℓ_2 -loss

$$\Rightarrow \hat{\omega} = (\underline{X^T X})^{-1} X^T \underline{y}$$

Remarks: Actually there are other choices for regression problem

One example is Huber Loss $L(y, y') := \begin{cases} \frac{1}{2}(y - y')^2, & |y - y'| \leq \delta \\ \delta|y - y'| - \frac{1}{2}\delta^2, & \text{o/w} \end{cases}$

② Linear Basis Model \rightarrow feature maps

1. Model: $\mathcal{H}_M = \{ f: f(x) = \omega^T \phi(x) \}$

$\phi: x \in \mathbb{R}^n \mapsto \phi(x) \in \mathbb{R}^M$ \rightarrow pre-defined function

2. Optimization: $\min_{\omega} \frac{1}{2N} \sum_{i=1}^N \|y_i - \omega^T \phi(x_i)\|_2^2$.

$$\Rightarrow \hat{\omega} = (\underline{\Phi^T \Phi})^{-1} \underline{\Phi^T y}$$

where $\underline{\Phi} := \begin{pmatrix} \phi(x_1)^T \\ \vdots \\ \phi(x_N)^T \end{pmatrix} \in \mathbb{R}^{N \times m}$

3. When $\underline{\Phi}$ is singular matrix $\Leftrightarrow \det(\underline{\Phi^T \Phi}) = 0$

Regularization

$$\min_{\omega} \frac{1}{2N} \sum_{i=1}^N \|y_i - \omega^T \phi(x_i)\|_2^2 + \lambda C(\omega)$$

choices of $C(\omega)$ $\left\{ \begin{array}{l} l_1 - \text{regularization} \quad C(\omega) = \|\omega\|_1 \\ l_2 - \text{regularization} \quad C(\omega) = \|\omega\|_2 \end{array} \right.$

③ Linear Basis Model + Classification

1. Architecture

$$x \in \mathbb{R}^d \xrightarrow{\phi} \phi(x) \in \mathbb{R}^M \xrightarrow{W} W\phi(x) \in \mathbb{R}^K$$

2. Model

$$\mathcal{H} = \left\{ f : f(x) = g\left(\sum_{i=0}^{M-1} w_i \cdot \phi_i(x)\right), w_i \in \mathbb{R}^K \right\}$$

$$g(\cdot) \rightarrow \boxed{\text{soft max}} \rightarrow g_k(z) := \frac{\exp(2k)}{\sum_i \exp(2i)}$$

3. Optimization

$$\min_{W \in \mathbb{R}^{M \times K}} \frac{1}{N} \sum_{i=1}^N L(g(W^\top \phi(x_i)), y_i)$$

Here, for classification task, we often apply cross-entropy loss

$$\rightarrow \boxed{L(y, y') = - \sum_{k=1}^K y'_k \log y_k} \rightarrow \text{non-symmetric}$$

Remarks: Actually, this comes from the KL-divergence,

which can be viewed as the distance between 2 probability dist.

$$KL(p || q) = \mathbb{E}_{x \sim p} \left[\log \frac{p(x)}{q(x)} \right]$$

$$= \int_X \log \frac{p(x)}{q(x)} \cdot p(x) dx$$

$$= \underbrace{\int_X p(x) \log p(x) dx}_{\text{Entropy of } p(x)} - \underbrace{\int_X p(x) \log q(x) dx}_{\text{cross-entropy term}}$$

4. Kernel Tricks

①. Kernel Ridge Regression

1. Model $\mathcal{H} = \{ f: f(x) = w^\top \phi(x) \quad w \in \mathbb{R}^m \}$

2. Linear Basis Model + L_2 -reg

$$\min_w \frac{1}{2N} \left(\sum_{i=1}^N \|y_i - w^\top \phi(x_i)\|_2^2 + \lambda \|w\|_2^2 \right)$$

$$\Rightarrow \hat{w} = \underbrace{(\Phi^\top \Phi + \lambda I_N)^{-1}}_{\sim} \Phi^\top y$$

Extension: Observation: $\hat{w} = (\Phi^\top \Phi + \lambda I_N)^{-1} \Phi^\top y$

$$= \Phi^\top (\Phi \Phi^\top + \lambda I_N)^{-1} y$$

Therefore, $\hat{f}(x) = \phi(x)^\top \hat{w}$

$$= \phi(x)^\top \underbrace{\Phi^\top}_{\Phi^\top = (\phi(x_1), \dots, \phi(x_N))} \underbrace{(\Phi \Phi^\top + \lambda I_N)^{-1}}_{\alpha = (\alpha_1; \dots; \alpha_N)} y$$

$\Phi^\top = (\phi(x_1), \dots, \phi(x_N)) \alpha = \begin{pmatrix} \alpha_1 \\ \vdots \\ \alpha_N \end{pmatrix}$

$\Phi \Phi^\top = \begin{pmatrix} \phi(x_1) \\ \vdots \\ \phi(x_N) \end{pmatrix} (\phi(x_1)^\top, \dots, \phi(x_N)^\top)$

Gram Matrix

$$= \sum_{i=1}^N \alpha_i \langle \phi(x_i), \phi(x) \rangle$$

$$= \sum_{i=1}^N \alpha_i k(x_i, x)$$

3. Kernel Trick

idea: feature maps can be computed expensively



define kernel functions directly

Issue: what is a Valid Kernel Function?

→ Answer: Mercer's Theorem (SPD Kernel)

$K: \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$ is a valid kernel

↔ for $\forall n \in \mathbb{N}^+$, $\{x_1, \dots, x_n\} \quad x_i \in \mathbb{R}^d$

$(G)_{ij} = K(x_i, x_j)$ is symmetric + PSD

→ $\begin{cases} G_{ij} = G_{ji} \\ \forall c \in \mathbb{R}^n, \quad c^T G c \geq 0 \end{cases}$

4. Re-formulation

a) choose a valid kernel $k(x, y)$

b) construct the Gram Matrix $G_{ij} = k(x_i, x_j)$

c) $\hat{f}(x) = \sum_{i=1}^n [(G + \lambda I_n)^{-1} y]_i k(x, x_i)$

5. Property of Kernel Function

→ if $k_1(\cdot, \cdot)$, $k_2(\cdot, \cdot)$ is valid kernel

then a) $k(x, y) = \lambda k_1(x, y)$

b) $k(x, y) = k_1(x, y) + k_2(x, y)$

$$c) K(x,y) = g(x) k_1(x,y) g(y)$$

$$d) K(x,y) = \lim_{n \rightarrow \infty} k_n(x,y)$$

$$e) K(x,y) = k_1(x,y) k_2(x,y)$$

all above is valid kernel

From these 5 properties, it is convenient to deduce that:

$$K(x,y) = \exp \left\{ -\frac{1}{2S^2} \|x-y\|_2^2 \right\}$$

5. SVM → we only consider the linear separable case
(hard margin)

① Primal formulation

a) linearly separable

$\mathcal{D} = \{(x_i, y_i)\}_{i=1}^N$ is linearly separable

$$\Leftrightarrow \exists w, b, \text{ s.t. } \underbrace{y_i(w^\top x_i + b)}_{> 0} > 0$$

margin

$$\begin{cases} \gamma_i^g(w, b) = \frac{|w^\top x_i + b|}{\|w\|} \\ \gamma^g(w, b) = \min_{i \in [N]} \gamma_i^g(w, b) \end{cases}$$

b) formulation

$$\max_{w,b} \gamma^g(w, b)$$

s.t. 'Hyperplane' can separate those 2 classes

$$\iff \begin{cases} \max_{w,b} \gamma^g(w,b) \\ \text{s.t. } f(x) = w^T x + b \text{ can separate those 2 classes} \end{cases}$$

$$\iff \begin{cases} \max_{w,b} \min_{i \in [N]} \gamma_i^g(w,b) = \frac{|w^T x_i + b|}{\|w\|} \\ \text{s.t. } y_i(w^T x_i + b) > 0 \end{cases}$$

$$\iff \begin{cases} \max_{w,b,\gamma} \frac{1}{\|w\|} \cdot \gamma \\ \text{s.t. } y_i(w^T x_i + b) > 0 \end{cases}$$

$$\min_{i \in [N]} y_i(w^T x_i + b) = \gamma$$

$$\iff \begin{cases} \max_{w,b,\gamma} \frac{1}{\|w\|} \cdot \gamma \\ \text{s.t. } y_i(w^T x_i + b) \geq \gamma \end{cases}$$

↗

$$\begin{cases} \max_{w,b} \frac{1}{\|w\|} \\ \text{s.t. } y_i(w^T x_i + b) \geq 1 \end{cases}$$

$$\iff \begin{cases} \min_{w,b} \frac{1}{2} w^T w \\ \text{s.t. } y_i(w^T x_i + b) \geq 1 \quad i \in [N] \end{cases} \quad (*)$$

② KKT Condition & Duality

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

the Dual Form is $\begin{cases} \max_{\lambda} \min_{w,b} L(w, b; \lambda) \\ \text{s.t. } \lambda \geq 0 \end{cases} \quad (**)$

$Q(\lambda) := \min_{w,b} L(w, b; \lambda)$

→ Since (*) is a convex program,

then we have :

$(\hat{w}, \hat{b}; \hat{\lambda})$ is the optimal solution with respect to Primal / Dual

$$\Leftrightarrow \left\{ \begin{array}{l} \nabla_{w,b} L(\hat{w}, \hat{b}; \hat{\lambda}) = 0 \rightarrow \text{Stationary} \\ \hat{\lambda}_i \geq 0 \\ 1 - y_i (\hat{w}^T x_i + \hat{b}) \leq 0 \end{array} \right] \rightarrow \text{Feasibility}$$

$$\hat{\lambda}_i (1 - y_i (\hat{w}^T x_i + b)) = 0 \rightarrow \text{Complementary Slackness}$$

③ Dual Form of Hard Margin SVM

$$\left\{ \begin{array}{l} \nabla_w L(w, b; \lambda) = w - \sum_{i=1}^N \lambda_i y_i x_i \\ \nabla_b L(w, b; \lambda) = - \sum_{i=1}^N \lambda_i y_i \end{array} \right.$$

$$\Rightarrow \left\{ \begin{array}{l} \hat{w} = \sum_{i=1}^N \hat{\lambda}_i y_i x_i \\ \sum_{i=1}^N \hat{\lambda}_i y_i = 0 \end{array} \right.$$

Dual Form

$$\left\{ \begin{array}{l} \max_{\lambda} \min_{w,b} L(w, b; \lambda) \\ \text{s.t. } \lambda \geq 0 \end{array} \right.$$

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

$$\Leftrightarrow \left\{ \begin{array}{l} \max_{\lambda} \frac{1}{2} \left\| \sum_{i=1}^N \lambda_i y_i x_i \right\|_2^2 + \sum_{i=1}^N \lambda_i - w^T \sum_{i=1}^N \lambda_i y_i x_i \\ \text{s.t. } \lambda \geq 0 \\ \sum_{i=1}^N \lambda_i y_i = 0 \end{array} \right.$$

$$\longleftrightarrow \begin{cases} \max_{\lambda} & -\frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \lambda_i \lambda_j y_i y_j x_i^T x_j + \sum_{i=1}^N \lambda_i \\ \text{s.t.} & \lambda_i \geq 0 \quad i \in [N] \\ & \sum_{i=1}^N \lambda_i y_i = 0 \end{cases}$$

when achieving $(\hat{\lambda}_1, \dots, \hat{\lambda}_N)$, then

$$\begin{cases} \hat{w} = \sum_{i=1}^N \hat{\lambda}_i y_i x_i \\ \text{for } \alpha \in SV, \quad y_\alpha (\hat{w}^T x_\alpha + \hat{b}) = 1 \\ \Rightarrow \hat{w}^T x_\alpha + \hat{b} = y_\alpha \\ \Rightarrow \hat{b} = y_\alpha - \hat{w}^T x_\alpha \\ = y_\alpha - \sum_{i=1}^N \hat{\lambda}_i y_i x_i^T x_\alpha \end{cases}$$

$$\begin{aligned} \text{to make prediction, } \hat{f}(x) &= \operatorname{sgn}(\hat{w}^T x + \hat{b}) \\ &= \operatorname{sgn}\left(\sum_{i=1}^N \hat{\lambda}_i y_i \langle x_i, x \rangle + \hat{b}\right) \end{aligned}$$

④ Kernel Form

Recap: the dual formulation is

$$\begin{cases} \max_{\lambda} & -\frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \lambda_i \lambda_j y_i y_j x_i^T x_j + \sum_{i=1}^N \lambda_i \\ \text{s.t.} & \lambda_i \geq 0 \quad i \in [N] \\ & \sum_{i=1}^N \lambda_i y_i = 0 \end{cases}$$

$$\longrightarrow \phi: x \in \mathbb{R}^d \longmapsto \phi(x) \in \mathbb{R}^D$$



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

$$= k(x_i, x_j)$$

\Rightarrow Kernel formulation :

$$\left\{ \begin{array}{l} \max_{\lambda} -\frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \lambda_i \lambda_j y_i y_j k(x_i, x_j) + \sum_{i=1}^N \lambda_i \\ \text{s.t. } \lambda_i \geq 0 \quad i \in [N] \\ \sum_{i=1}^N \lambda_i y_i = 0 \end{array} \right.$$

$$\hat{w} = \sum_{i=1}^N \lambda_i y_i \phi(x_i)$$

$$\Rightarrow \hat{f}(x) = \text{sgn} \left(\sum_{i=1}^N \lambda_i y_i k(x_i, x) + \hat{b} \right)$$

where \hat{b} is determined by $\hat{b} = y_0 - \sum_{i=1}^N \lambda_i y_i k(x_i, x_0)$ $x \in SV$

6. Decision Tree

① Model :

$$\mathcal{H} = \{ f : f(x) = \sum_{j=1}^J a_j \mathbb{1}_{R_j}(x), \{R_j\}_{j=1}^J \text{ is partition of } \mathcal{X} \}$$

② Optimization :

$$\left\{ \begin{array}{l} \text{classification} \quad \min_{R_j} \sum_{j=1}^J \left\{ \sum_{k=1}^K -p(k|j) \log p(k|j) \right\} \\ \text{regression:} \quad \min_{a_j, R_j} \frac{1}{2} \sum_{i=1}^N \left(\sum_{j=1}^J a_j \mathbb{1}_{R_j}(x_i) - y_i \right)^2 \end{array} \right.$$

Issue → minimizing all partition of \mathcal{X} ⇒ Difficult / NP-Hard

\Rightarrow Strategy : greedy algorithm (Recursive Binary Search)

③ Ad. and Dis-ad.

Advantages: 1. Visualization

2. Implicit feature selection through the reduction
of error / impurity

3. Non-linear relationships

4. Robust to data types

Disadvantages: 1. prone to over-fitting

2. sensitive to { data imbalancing
data variation}

3. Greedy Algo \rightarrow sub-optimal

7. Ensemble {
Bagging
Boosting

a) Bagging (Bootstrap Aggregating)

\hookrightarrow { regression \rightarrow average
classification \rightarrow majority vote

Algo

$$1. \mathcal{D} = \{(X_i, Y_i)\}_{i=1}^N$$

2. for $j=1, 2, \dots, B$

Sample \mathcal{D}_j from \mathcal{D} with replacement

(use ECDF to approximate unknown CDF)

train f_j from \mathcal{D}_j

$$3. \bar{f}(x) = \frac{1}{B} \sum_{b=1}^B f_b(x)$$

Idea (Naive calculation)

assume $f_b(x) = f^*(x) + \varepsilon_b(x) \quad b=1, 2, \dots, B$

where $\varepsilon_b(x) \rightarrow \begin{cases} \mathbb{E}[\varepsilon_b(x)] = 0 \\ \text{Var}[\varepsilon_b(x)] = \sigma^2 \end{cases}$

then analyze the power of $\bar{f}(x) = \frac{1}{B} \sum_{b=1}^B f_b(x)$

$$\textcircled{1} \text{ MSE}(f_i) \rightarrow \mathbb{E}[(f_i(x) - f^*(x))^2]$$

$$= \mathbb{E}[(f_i(x) - \mathbb{E}[f_i(x)])^2]$$

$$= \text{var}[f_i(x)]$$

$$= \sigma^2$$

$$\textcircled{2} \text{ MSE}(\bar{f}) \rightarrow \mathbb{E}[(\bar{f}(x) - f^*(x))^2]$$

$$= \mathbb{E}\left[\left(\frac{1}{B} \sum_{b=1}^B f_b(x) - f^*(x)\right)^2\right]$$

$$= \mathbb{E}\left[\left(\frac{1}{B} \sum_{b=1}^B \varepsilon_b(x)\right)^2\right]$$

$$= \frac{1}{B^2} \left[\sum_{b=1}^B \varepsilon_b^2(x) + \sum_{i \neq j} \varepsilon_i(x) \varepsilon_j(x) \right]$$

if uncorrelated

$$\frac{1}{B} \cdot \sigma^2$$

Remarks: In real application, this significant shrink of variance is not

apparent. That is because, through sub-sampling, 2 things will happen:

1) # of sub-sampling large : highly correlated ;

2) # of sub-sampling small : σ^2 is large .

b) Boosting (AdaBoost) \rightarrow sequential model

↓

idea re-weight the data

correct \rightarrow less weight
incorrect \rightarrow larger weight

Formulation

1. Model (additive forward function)

$$F_m(x) = \sum_{i=1}^m d_i f_i(x)$$

$$\begin{aligned} &= \sum_{i=1}^{m-1} d_i f_i(x) + \alpha_m f_m(x) = \\ &= \underbrace{F_{m-1}(x)}_{\sim} + \alpha_m f_m(x) \end{aligned}$$

2. Loss function (exp loss)

$$L(f(x), y) = \exp \{-y f(x)\}$$

$$\Rightarrow L(F_m(x), y) = \exp \{-y F_m(x)\}$$

$$= \exp \{-y (F_{m-1}(x) + \alpha_m f_m(x))\}$$

$$= \exp \{-y F_{m-1}(x) - y \alpha_m f_m(x)\}$$

$$= \exp \{-y F_{m-1}(x)\} \exp \{-y \alpha_m f_m(x)\}$$

$$\text{Empirical Risk} = \frac{1}{N} \sum_{i=1}^N \exp \{-y_i F_{m-1}(x_i)\} \exp \{-y_i \alpha_m f_m(x_i)\}$$

(idea)

3. Greedy Algo involved \rightarrow Fix $F_{m-1}(x)$, find (α_m, f_m)

Define $w_i^{(m)} = \exp \{-y_i F_{m-1}(x_i)\}$

then Empirical Risk = $\frac{1}{N} \sum_{i=1}^N w_i^{(m)} \exp \{-y_i \alpha_m f_m(x_i)\}$

$$= \frac{1}{N} \left(\sum_{i \in C} w_i^{(m)} \exp \{-\alpha_m\} + \sum_{i \notin C} w_i^{(m)} \exp \{\alpha_m\} \right)$$

$$= \frac{1}{N} \sum_i w_i^{(m)} \exp \{-\alpha_m\} \cdot \mathbb{1}\{y_i = f_m(x_i)\}$$

$$\begin{aligned}
& + \frac{1}{N} \sum_i w_i^{(m)} \exp \{ \alpha_m \} \mathbb{1} \{ y_i \neq f_m(x_i) \} \\
& = \underbrace{\frac{1}{N} \exp(-\alpha_m) \sum_{i=1}^N w_i^{(m)}}_{+ \frac{1}{N} (\exp(\alpha_m) - \exp(-\alpha_m)) \sum_{i=1}^N w_i^{(m)} \mathbb{1} \{ y_i \neq f_m(x_i) \}}
\end{aligned}$$

4. Summary

$(\hat{f}_m, \hat{\alpha}_m) = \underset{f_m, \alpha_m}{\operatorname{argmin}} \text{ Empirical Risk on Exp Loss}$

$$= \underset{f_m, \alpha_m}{\operatorname{argmin}} \exp(-\alpha_m) \sum_{i=1}^N w_i^{(m)} + [\exp(\alpha_m) - \exp(-\alpha_m)] \sum_{i=1}^N w_i^{(m)} \mathbb{1} \{ y_i \neq f_m(x_i) \}$$

$$\Rightarrow \textcircled{1} \quad \hat{f}_m = \underset{f_m}{\operatorname{argmin}} \sum_{i=1}^N w_i^{(m)} \mathbb{1} \{ y_i \neq f_m(x_i) \}$$

$$\textcircled{2} \quad \hat{\alpha}_m = \underset{\alpha_m}{\operatorname{argmin}} \exp \{ \alpha_m \} \sum_{i \in I^c} w_i^{(m)} + \exp \{ -\alpha_m \} \sum_{i \in I^c} w_i^{(m)}$$

$$= \frac{1}{2} \log \left(\frac{1 - \hat{\varepsilon}_m}{\hat{\varepsilon}_m} \right)$$

$$\boxed{\hat{\varepsilon}_m := \frac{\sum_{i=1}^N w_i^{(m)} \mathbb{1} \{ y_i \neq \hat{f}_m(x_i) \}}{\sum_{i=1}^N w_i^{(m)}}}$$

8. Neural Network

$$\left\{ \begin{array}{l} \text{MLP} \rightarrow \text{LECS} \\ \text{CNN} \rightarrow \text{LEC6} \\ \text{RNN} \rightarrow \text{LEC7} \end{array} \right.$$

Summary : As for supervised learning task, we mainly introduce

a) Linear Model

Linear Regression

Linear Basis Model

regression
classification

Regularization Technique

b) Kernel Trick

Actually Linear Basis Model can also achieve!

Ridge Regression (with ℓ_2 -regularization) + kernel form

→ exact solution can be determined by $K(x,y) = \langle \phi(x), \phi(y) \rangle$
exact predictive function

SVM + kernel

→ 1. optimization problem can be formulated by kernel function, instead of feature map;
2. predictive function can be expressed by kernel function;

c) SVM

Hard-margin SVM

PRIMAL
DUAL

Soft-margin SVM

PRIMAL
DUAL

Kernel SVM

PRIMAL
DUAL

d) Decision Tree

Regression

Classification

Trick: Greedy algorithm involved → recursive binary search

e) Ensemble

1. Bagging → Bootstrap the Dataset



Aggregating through average

2. Boosting → AdaBoost

→ Re-weight + training + voted average

f) Neural Network

{

- FCNN
- RNN
- CNN