① Term Project.

$\begin{cases} A & \text{Experimental} \\ B & \text{Theoretical} \end{cases}$   5-pages

① Summary

② $\begin{cases} \text{advocate} \\ \text{critic} \end{cases}$

---

This time : Boosting (Adaboost)

→ NUS Prof. (Generalization)

Recap: $\ell_1$ - regular. (Lasso / Compressed Sensing)

① Mathematical

② Graph Proof

$$\min_{\theta} \quad \frac{1}{2n} \sum_{t} (y_t - \theta^T x_t)^2 + \lambda \|\theta\|_1 .$$
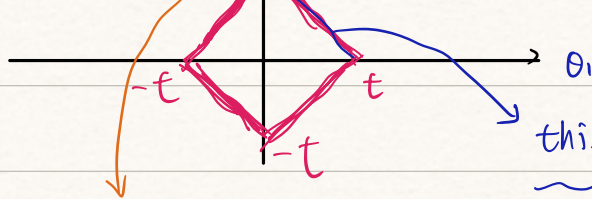
$\updownarrow$ Equivalent.

$$\min_{\theta} \quad \frac{1}{2n} \|y - X\theta\|^2 \quad \text{s.t} \quad \|\theta\|_1 \leq t \quad \underline{\text{for some } t > 0}$$

(LASSO)
Figure :



$\tilde{\theta} = (X^T X)^{-1} X^T y$

Sparse Solution

this probability is very small!

(Lebesgue Measure)

# Boosting / Adaboost  (Feature Selection)
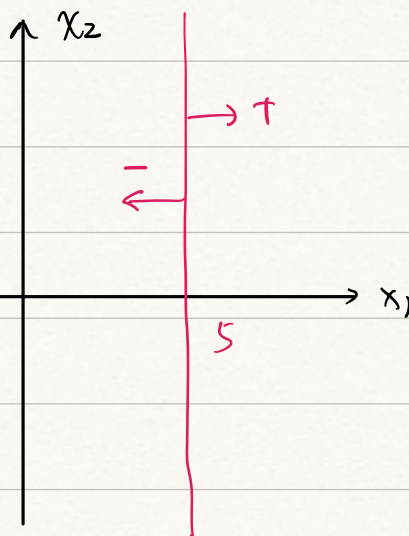
## Decision stump (Axis-aligned linear separators)

$$h(x; \underline{\theta}) = \text{sign}(s(x_k - \theta_0)), \quad \underline{x} \in \mathbb{R}^d \quad k \in \{1, \dots, d\} \quad s \in \{\pm 1\}$$
$$\theta_0 \in \mathbb{R}$$

[Example]

$$\begin{cases} s = +1 \\ k = 1 \Rightarrow \text{Focus on } x_1 \\ \theta_0 = 5 \end{cases}$$
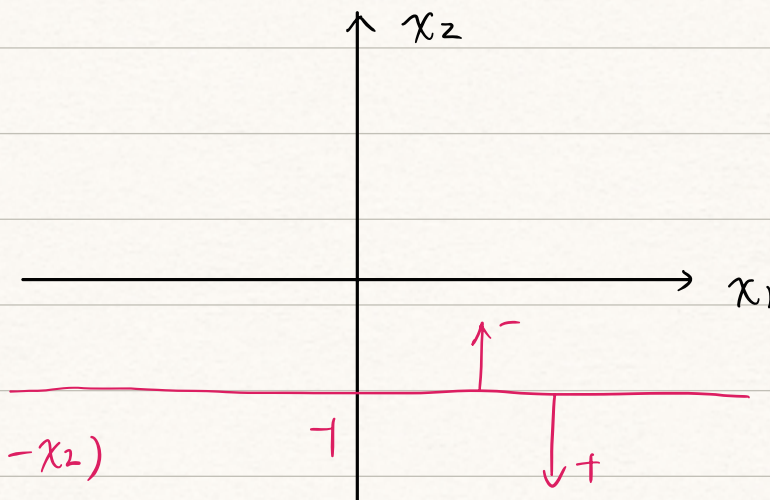
$$h(\underline{x}, \underline{\theta}) = \text{sign}(x_1 - 5)$$



[Example]

$$\begin{cases} s = - \\ k = 2 \\ \theta_0 = -1 \end{cases}$$
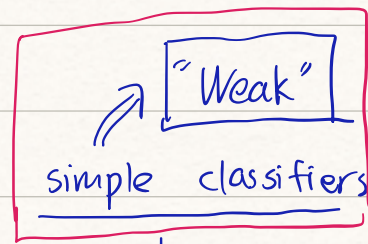
$$h(x; \theta) = \text{sign}(-1 - x_2)$$

Pretty Simple Classifier!

↓↓

Base Learner



Q: Can we "combine" several such simple classifiers to
Form a "Strong" classifier?

"Decision stumps"

Eg: $\phi(\underline{x}; \underline{\theta}) = [\, h(x_i; \underline{\theta_i}) ; \, i=1,...,m \,]^T \in \{\pm 1\}^m$  ⟹ Consider $m$ decision stumps

↓↓

$\underline{\theta_i} = \{s_i, k_i, \theta_{0i}\}$

⟹ Difficult to get (learn) parameters

Run a linear classifier based on $\phi(x, \underline{\theta})$

Eg: ✓ Collect the output of all decision stumps into a SINGLE
Classifier.

⟶ { Tractable
    Generalization Well }

Convex combin.

↓↓

Ensemble ⟶ $h_m(\underline{x}) = \sum_{j=1}^{m} \alpha_j \, h(\underline{x}; \underline{\theta_j})$ , $\alpha_j \geqslant 0$ & $\sum \alpha_j = 1$.

e.g. if $m=2$, and $h(\cdot; \theta_2)$ is more "reliable",
then we may choose $\alpha_2 = 0.8$, $\alpha_1 = 0.2$.

Rmk: Each stump has $\alpha_j$ ⟶ "votes".
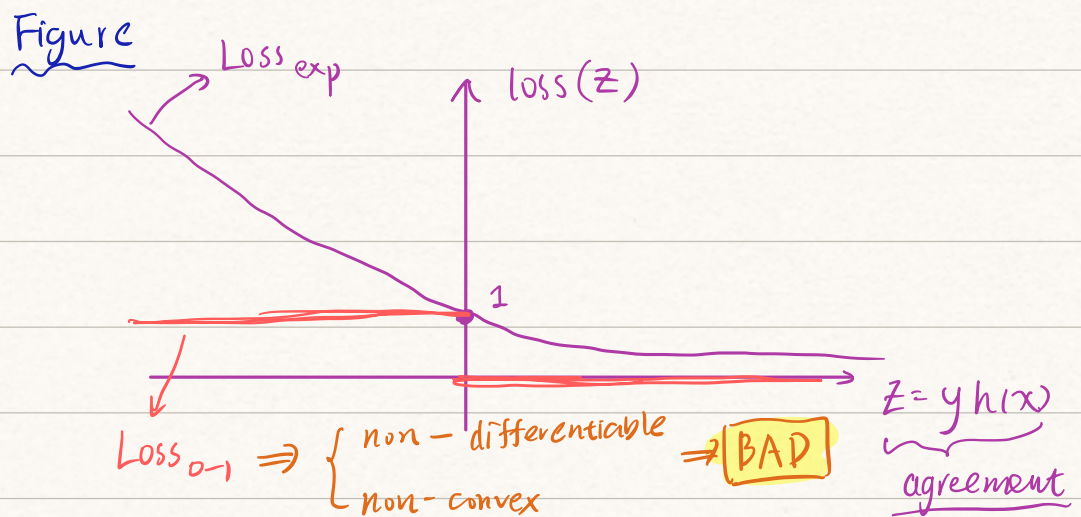
Ensemble $h_m(x)$ classifies a data $\underline{x}'$ according to
the votes $\alpha_j$ of each stump $h(\cdot; \underline{\theta_j})$

Our Aim: Learn $\{(\underline{\theta_j}, \alpha_j)\}_{j=1}^{m}$ where $\underline{\theta_j} = \{s_j, k_j, \theta_{0j}\}$

Weak Learner $\xRightarrow{\text{Ensemble}}$ Strong Learner

$\downarrow$

Stump

---

Exponential Loss $\rightarrow \text{Loss}_{exp}(y, h(x)) = \exp(-y\,h(x)) = \exp(-z)$

$\downarrow$ label      $\searrow$ prediction

Figure



$\text{Loss}_{0-1} \Rightarrow \begin{cases} \text{non - differentiable} \\ \text{non - convex} \end{cases} \Rightarrow \boxed{BAD}$

$z = y\,h(x)$ agreement

Rmk: $\text{Loss}_{0-1}(z) \leq \text{Loss}_{exp}(z)$

$\Downarrow$

The exp. loss is a SURROGATE to the

0-1 loss. $\searrow \begin{cases} \text{convex} \\ \text{differentiable} \end{cases}$

Aim: Learn parameters $\underline{\theta_j}$ and votes $\alpha_j$ GREEDILY.

$\Downarrow$ greedy alg.

Assume that we learn $(m-1)$-th weak classifier & their

votes, i.e., $\{(\theta_j, \alpha_j)\}_{j=1}^{m-1}$. These are $\boxed{\text{fixed}}$ for learning

the m-th term $\boxed{\alpha_m \cdot h(\cdot\ ; \theta_m)}$

Consider the ensemble:

$$h_m(x) = \sum_{j=1}^{m} \alpha_j\, h(x; \theta_j)$$

$$= \left[\sum_{j=1}^{m-1} \alpha_j\, h(x; \theta_j)\right] + \alpha_m\, h(x; \theta_m)$$

$$= h_{m-1}(x) + \alpha_m\, h(x; \theta_m)$$

Compute the exp. loss on a given dataset $\mathcal{D}_n = \{(x_t, y_t)\}_{t=1}^{n}$

$$J(\alpha_m, \theta_m) = \sum_{t=1}^{n} \text{loss}_{\exp}(y_t, h_m(x_t))$$

→ ensemble learner

↓ vote    ↓ stump

Our aim is to minimize this!

$$= \sum_{t=1}^{n} \exp(-y_t\, h_m(x_t))$$

→ ensemble at m-th iter

$$= \sum_{t=1}^{n} \exp(-y_t\, h_{m-1}(x_t) - y_t\, \alpha_m\, h(x_t; \theta_m))$$

→ ensemble at (m-1)-th iter

→ Loss → t & m-1

$$= \sum_{t=1}^{n} W_{m-1}(t)\, \exp(-y_t\, \alpha_m\, h(x_t; \theta_m))$$

$$W_{m-1}(t) := \exp(-y_t\, h_{m-1}(x_t))$$

"weights" Associated to data point $(x_t, y_t)$ after (m-1)-th iteration

Note: $W_m(t) = W_{m-1}(t) * \exp(-y_t\, \alpha_m\, h(x_t; \theta_m))$

# Ada Boost Alg.

⇒ Input : $\mathcal{D}_n = \{(x_t, y_t)\}_{t=1}^n$

Loss $= \exp(-yz)$

$z = \sum_{n} d_n g_n(x)$

— Initialize weight : $\underline{W_0(t) = \frac{1}{n}}$ For all $t$.

→ we have $(a_i, \hat{\theta}_i)$ $i = 1, 2, \ldots, m-1$
(stump)

— At Boosting stage $m$ , find a base learner $h(\cdot ; \hat{\theta}_m)$

that minimize :

$$\hat{\theta}_m = \underset{\theta_m}{\arg\min} - \sum_{t=1}^n \widetilde{W}_{m-1}(t) \, y_t \, h(x_t ; \theta_m)$$

↳ weighted trng loss.

↓

Learn the Stump.

$$\widetilde{W}_{m-1}(t) = \frac{W_{m-1}(t)}{\sum_t W_{m-1}(t)}$$

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

→ normalized weight

$\hat{\varepsilon}_m =$

— Choose vote $\hat{\alpha}_m \in \mathbb{R}$ using a formula (comes from Greedy)

— Update the weights $\widetilde{W}_m(t) = \dfrac{\widetilde{W}_{m-1}(t)}{Z_m} \exp\left(-y_t \, h(x_t, \hat{\theta}_m) \, \hat{\alpha}_m\right)$

$$= \bigcirc \begin{cases} \exp(\alpha_m), \text{ inequal} \\ \exp(-\alpha_m), \text{ equal} \end{cases}$$

---

Note: ①

$$\boxed{- \sum_{t=1}^n \widetilde{W}_{m-1}(t) \, y_t \, h(x_t ; \theta_m)}$$

⇒ Motivation

Claim: $-y_t \, h(x_t ; \theta_m) = 2 \cdot \mathbb{1}\{y_t \neq h(x_t ; \theta_m)\} - 1$

pf: ① $y_t = h(x_t; \theta_m) \Rightarrow$ LHS $= -1$  RHS $= -1$

② $y_t \neq h(x_t; \theta_m) \Rightarrow$ LHS $= +1$  RHS $= +1$

$$- \sum_{t=1}^{n} \widetilde{W}_{m-1}(t) \, y_t \, h(x_t; \underline{\theta_m}) \quad \begin{cases} y_t = h \Rightarrow 1 \\ y_t \neq h \Rightarrow -1 \end{cases}$$

$$= \sum_{t=1}^{n} \widehat{W}_{m-1}(t) \left[ 2 \cdot \mathbb{1}\{y_t \neq h(x_t; \theta_m)\} - 1 \right]$$

$$= \sum_{t=1}^{\hat{n}} 2 \, \widetilde{W}_{m-1}(t) \, \mathbb{1}\{y_t \neq h(x_t; \theta_m)\} - 1 \quad (\text{since } \sum_{1}^{n} \widetilde{W}_{m-1}(t) = 1)$$

$$:= 2\hat{\varepsilon}_m - 1.$$

$$\hat{\varepsilon}_m = \sum_{t=1}^{n} \widetilde{W}_{m-1}(t) \, \mathbb{1}\{y_t \neq h(x_t, \hat{\theta}_m)\}$$

weighted training $= \sum\limits_{t:\, \text{misclassified}} \widetilde{W}_{m-1}(t)$

loss

Rmk: Since $\widetilde{W}_0(t) = \frac{1}{n}$, $t = 1, 2, \dots, n$. $\varepsilon_1 = $ trg error $= \sum\limits_{t:\, y_t \neq h(x_t, \theta_1)} \frac{1}{n}$

② Update Rule of Weights  (when we attain $\hat{\theta}_m$)

$$\widetilde{W}_m(t) = \frac{\widetilde{W}_{m-1}(t)}{Z_m} \exp\left( \overbrace{-y_t \, h(x_t; \underline{\theta_m})}^{\in \{\pm 1\}} \, \hat{\alpha}_m \right)$$

$$= \frac{\widetilde{W}_{m-1}(t)}{Z_m} \times \begin{cases} e^{-\hat{\alpha}_m}, & \text{if } y_t = h(x_t, \hat{\theta}_m) \\ e^{+\hat{\alpha}_m}, & \text{if } y_t \neq h(x_t, \hat{\theta}_m) \end{cases} \quad \rightarrow \text{increase weight}$$

$\Rightarrow$ So, if $\hat{\alpha}_m > 0$ (usual case), trng examples that are misclassified by $h(\cdot; \hat{\theta}_m)$ are given higher weights in the next boosting stage.

Final Classifier $\rightarrow$ $h_m(x) = \sum_{j=1}^{m} \hat{\alpha}_j \, h(x; \hat{\theta}_m)$

Q: How to choose the votes $\hat{\alpha}_j$

$$\hat{\alpha}_m = \frac{1}{2} \ln\left(\frac{1 - \hat{\varepsilon}_m}{\hat{\varepsilon}_m}\right) \qquad \hat{\varepsilon}_m = \sum_{t=1}^{n} \hat{W}_{m-1}(t) \, \mathbb{1}\{y_t \neq h(x_t; \hat{\theta}_m)\}$$

$\hat{\varepsilon}_m$: weighted training error when we consider the optimized stump $h(\cdot; \hat{\theta}_m)$ $\Rightarrow$ m-th iteration (update)

Rmk: ① $\hat{\varepsilon}_m = 0$ $\Rightarrow$ $\hat{\alpha}_m = +\infty$ $\longleftrightarrow$ put all weights on
$\qquad\qquad\qquad \updownarrow \qquad\qquad\qquad\qquad$ classifier $h(\cdot; \hat{\theta}_m)$

$h(\cdot; \hat{\theta}_m)$ perfectly classifies all training samples

② $\hat{\varepsilon}_m = \frac{1}{2}$ $\Rightarrow$ $\hat{\alpha}_m = 0$ $\leftarrow$ No weights on
$\qquad\qquad \updownarrow \qquad\qquad\qquad\qquad\qquad$ classifier $h(\cdot; \hat{\theta}_m)$

$h(\cdot; \hat{\theta}_m)$ is 'maximally confused'

How does AdaBoost Perform?

$\downarrow$

Guarantee by a theorem:

[THM]: If each base classifier $h(\cdot ; \hat{\theta_j})$ is slightly better than random guessing, i.e., $\hat{\varepsilon_j} < \frac{1}{2}$

$\downarrow$

$$\sum \widetilde{W_{j-1}}(t)$$

$t$: classified wrongly on $h(\cdot ; \hat{\theta_j})$

the training error decrease to $0$ exponentially,

i.e., $\frac{1}{n} \sum_{t=1}^{n} \mathbb{1}\{y_t h_m(X_t) \leq 0\} \leq \exp\left(-2\sum_{j=1}^{m}\left(\frac{1}{2} - \hat{\varepsilon_j}\right)^2\right)$

$\downarrow$

In particular, if there exists $\gamma$ s.t.

$\frac{1}{2} - \hat{\varepsilon_j} \geq \gamma$

then $\frac{1}{n}\sum \mathbb{1}\{y_t h_m(X_t) \leq 0\} \leq \exp(-2m\gamma^2)$

$\downarrow$

training error (0-1)