# Chap 9 Linear Predictors.

$\boxed{L_d} := \{ h_{w,b} \mid w \in \mathbb{R}^d, b \in \mathbb{R} \}$, where

$$\boxed{h_{w,b}(x) := \langle w, x \rangle + b} = \left( \sum_{i=1}^{d} w_i x_i \right) + b,$$

(bias)

is called the class of $\boxed{\text{affine functions}}$.

$\boxed{\text{Rmk}:}$

The different hypothesis classes of linear predictors are compositions of a function

$$\phi : \mathbb{R} \to \mathcal{Y} \quad \text{on } L_d.$$

ex

① binary classification : $\phi(x) := \begin{cases} 1 & , x > 0 \\ 0 & , x = 0 \\ -1 & , x < 0 \end{cases}$  ← this is the $\boxed{\text{sign fun.}}$

② regression : $\phi(x) = x$, the identity fun.

$\boxed{\text{Rmk}:}$

$\boxed{w'} = (b, w_1, \cdots, w_d) \in \mathbb{R}^{d+1}$, $\boxed{x'} := (1, x_1, \cdots, x_d) \in \mathbb{R}^{d+1}$.

Then

$$h_{w,b}(x) = \langle w, x \rangle + b = \langle w', x' \rangle, \text{ called}$$

the $\boxed{\text{homogeneous representation}}$ of $h_{w,b}$,

simplifying the representation.

## 9.1 Halfspaces.

$\boxed{HS_d} = \text{sign} \circ L_d := \{ x \mapsto \text{sign}(h_{w,b}(x)) : h_{w,b} \in L_d \}$

called the class of $\boxed{\text{halfspace hypothesis}}$.

$\underline{\text{Rmk}}:$

In this setting, the realizable case is often called $\boxed{\text{separable}}$. ($\leftrightarrow \boxed{\text{nonseparable}}$).

$\underline{\text{Note}}:$

Implementing ERM in the nonseparable case w.r.t. 0-1 loss is known to be computationally hard. As a substitute, some people use $\boxed{\text{surrogate loss functions}}$ instead. (as in the logistic regression)

## 9.1.1 Linear Programming for the Class of Halfsp.

Linear program (LP) :

$$\boxed{\begin{array}{l} \max\limits_{w \in \mathbb{R}^d} \langle u, w \rangle \\ \text{subject to } Aw \geq V \end{array}}, \text{ where } \begin{array}{l} A : m \times d \\ V \in \mathbb{R}^m \text{ are given} \\ u \in \mathbb{R}^d \quad \text{and} \\ w \in \mathbb{R}^d \text{ is to be determined} \end{array}$$

$\boxed{\text{Prop}}$

In the realizable case, the ERM problem for halfspaces can be expressed as a linear program.

$\langle Pf \rangle$

$\underline{\text{Goal}}$ : Find $w$ s.t. $\text{sign}(\langle w, x_i \rangle) = y_i, \forall i$.

$(\Leftarrow) \ y_i \langle w, x_i \rangle > 0, \forall i.$

$\because$ separable $\therefore \exists w^*$ s.t. $y_i \langle w^*, x_i \rangle > 0$.

Let $\gamma := \min\limits_{i} y_i \langle w^*, x_i \rangle$ so that $\gamma \leq y_i \langle w^*, x_i \rangle, \forall i.$

$\Rightarrow y_i \langle \frac{w^*}{\gamma}, x_i \rangle \geq 1, \forall i.$

i.e. $y_i \langle \overline{w}, x_i \rangle \geq 1, \forall i, \text{ w/ } \overline{w} = \frac{w^*}{\gamma}.$  ✳

Note that $y_i x_i = y_i \begin{pmatrix} x_{i1} \\ \vdots \\ x_{id} \end{pmatrix}$

Thus, ✳ is $y_i \cdot (x_{i1} \cdots x_{id}) \cdot \overline{w} \geq 1, \forall i.$

Let

$$A = \begin{pmatrix} y_1 x_{11} & y_1 x_{12} & \cdots & y_1 x_{1d} \\ \vdots & \vdots & & \vdots \\ y_d x_{d1} & y_d x_{d2} & \cdots & y_d x_{dd} \end{pmatrix} \text{ and } V = \begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix},$$

✳ becomes $A \overline{w} \geq V.$

Setting an arbitrary $u$, say $u = 0$, finding $\overline{w}$ becomes an LP problem. ⌗

## 9.1.2 Perceptron for Halfspaces.

A different approach to implement ERM is the $\boxed{\text{Perceptron algorithm of Rosenblatt}}$.

$$\boxed{\begin{array}{l} \underline{\text{Batch Perceptron}} \\ \text{Input : A training set } (x_1, y_1), \cdots, (x_m, y_m). \\ \text{Initialize : } w^{(1)} = (0, \cdots, 0) \\ \quad \text{for } t = 1, 2, \cdots \\ \qquad \text{if } (\exists i \text{ s.t. } y_i \langle w^{(t)}, x_i \rangle \leq 0) : \text{do} \\ \qquad \quad w^{(t+1)} = w^{(t)} + y_i x_i \\ \qquad \text{else} \\ \qquad \quad \text{output } w^{(t)}. \end{array}}$$

## Idea:

The update of the Perceptron guides the sol'n to be _more correct_ on the $i$th example.:

$$y_i \langle w^{(t+1)}, x_i \rangle = y_i \langle w^{(t)} + y_i x_i, x_i \rangle$$

$$= y_i \langle w^{(t)}, x_i \rangle + \|x_i\|^2.$$

### Thm 9.1 (Perceptron under separable case).

$(x_1, y_1), \cdots, (x_m, y_m)$ : separable.

$B := \min \{ \|w\| : y_i \langle w, x_i \rangle \geq 1, \forall i \in [m] \}$.

$R := \max_i \|x_i\|$.

Then the Perceptron algorithm stops after at most $(RB)^2$ interations.

(when it stops, we have $y_i \langle w^{(t)}, x_i \rangle > 0$, $\forall i \in [m]$).

### 9.1.3 The VC dim of Half spaces.

### Thm 9.2 (See P122 for proof).

The VC dimension of the class of homogeneous half sp. in $\mathbb{R}^d$ is $d$.

### Thm 9.3 (See P122 for proof).

The VC dimension of the class of non homogeneous half sp. in $\mathbb{R}^d$ is $d+1$.

### 9.2 Linear Regression.

$\mathcal{H}_{reg} = L_d = \{ x \mapsto \langle w, x \rangle + b : w \in \mathbb{R}^d, b \in \mathbb{R} \}$.

One common way is to use the squared-loss fun. In this case, the empirical risk is called the **Mean Squared Error**

Another option: absolute value loss funct'n.

Note : This can be solved by LP. (Exercise 1).

Rmk:

Regression problems __cannot__ be analyzed using VC dimension. Rigorous means will be introduced later to analyze regression problems.
in the book.

### 9.2.1 Least Squares.

Least squares is the algorithm for solving ERM problems for the hypothesis class of linear regression predictors w.r.t. squared loss.

### Goal:

Solve $\arg\min_w L_S(h_w) = \arg\min_w \frac{1}{m} \sum_{i=1}^{m} (\langle w, x_i \rangle - y_i)^2$.

Taking gradient, we have

$$\nabla_w L_S(h_w) = (Aw - b) \cdot \frac{2}{m}, \text{ where}$$

$$A = \left( \sum_{i=1}^{m} x_i x_i^T \right) \text{ and } b = \sum_{i=1}^{m} y_i x_i$$

$$= \begin{pmatrix} x_1 \cdots x_m \end{pmatrix} \begin{pmatrix} x_1 \cdots x_m \end{pmatrix}^T \quad = \begin{pmatrix} x_1 \cdots x_m \end{pmatrix} \begin{pmatrix} y_1 \\ \vdots \\ y_m \end{pmatrix}.$$

Setting the gradient to 0, we are solving

$$Aw = b.$$

(Note: $L_S(h_w)$ is convex, as a fun. of $w$, and thus every local min. is a global min.)

### Prop

W/ $A$ and $b$ as above, the equation $Aw = b$ always has a sol'n.

⟨Pf⟩

∵ $A$ is real symmetric

∴ We can write it as the eigen decomp. :

$$A = VDV^T, \text{ where } V : \text{orthogonal and } D : \text{diag.}$$

Define $D^\dagger$ as a $d \times d$ diagonal matrix w/

$$D_{ii}^\dagger = \begin{cases} 1/D_{ii} & \text{if } D_{ii} \neq 0 \\ 0 & \text{o.w.} \end{cases}$$

Define

$$A^\dagger = VD^\dagger V^T, \text{ and } \hat{w} = A^\dagger b$$

(to be claimed as a sol'n to $Aw = b$.)

Then $A\hat{w} = VDV^T VD^\dagger V^T b = VDD^\dagger V^T b$

$$= \sum_{i : D_{ii} \neq 0} \langle b, v_i \rangle v_i, \text{ where } V = \begin{pmatrix} v_1 \cdots v_d \end{pmatrix}.$$

Thus, $A\hat{w}$ is the orthogonal proj. of $b$ onto $R(A)$.

(since $v_i$ w/ $D_{ii} \neq 0$ are the e-vec. w/ nonzero e-val.)

$$X := \begin{pmatrix} x_1 \cdots x_m \end{pmatrix}. \text{ (so } A = XX^T).$$

Note $\text{rank}(A) = \text{rank}(XX^T) = d - \text{nullity}(XX^T)$

$$= d - \text{nullity}(X^T) = \text{rank}(X^T) = \text{rank}(X) \text{ and}$$

$R(A) \subseteq R(X)$. Thus, $R(A) = R(X)$.

$\therefore b \in R(X)$ $\therefore A\hat{w} = b$. #

## 9.3 Logistic Regression.

In logistic regression, we learn a family of functions $h: \mathbb{R}^d \to [0,1]$. Mainly used for classification task.

We compose $L_d$ w/ a sigmoid function (meaning "S-shaped" function).

In logistic regression, we use

$$\phi_{sig}(z) = \frac{1}{1+e^{-z}}, \text{ called the logistic function.}$$

$$\mathcal{H}_{sig} := \phi_{sig} \circ L_d = \{x \mapsto \phi_{sig}(\langle w, x \rangle) : w \in \mathbb{R}^d\}.$$

Note:

① When $\langle w, x \rangle$ is very positive, $\phi_{sig}(\langle w, x \rangle) \approx 1$.

② " " " " " negative, " " $\approx 0$.

③ " " " " " close to 0, " " $\approx \frac{1}{2}$.

Thus, $\mathcal{H}_{sig}$ is solving classification problem via a probabilistic approach.

The loss function used by logistic regression is

$$\ell(h_w(x,y)) = \log(1 + e^{-y\langle w, x \rangle}).$$

Thus, given a training set $S = (x_1, y_1), \cdots, (x_m, y_m)$, we are solving (if we use ERM paradigm).

$$\underset{w \in \mathbb{R}^d}{\arg\min} \frac{1}{m} \sum_{i=1}^{m} \log(1 + e^{-y_i \langle w, x_i \rangle}).$$

Rmk:

① $\because \log(1 + e^{-y\langle w, x \rangle})$, the loss functn, is convex
$\therefore$ ERM is easy to implement. (via gradient descent)

② Implementing ERM on logistic regression is equivalent to finding Maximum Likelihood Estimator (to be introduced in Chap 24).

---

# Chap 10 Boosting.

Q: Can an efficient weak learner be "boosted" into an efficient strong learner?

Ada Boost (= Adaptive Boosting) will be introduced here.

Rmk:
successfully
Ada Boost has been used to detect faces in images.

## 10.1 Weak Learnability.

Def $\mathcal{H}$: a hypothesis class.

(1) A learning algorithm $A$ is a $\gamma$-weak-learner for $\mathcal{H}$ if

$\exists m_{\mathcal{H}}: (0,1) \to \mathbb{N}$ s.t. $\forall \delta \in (0,1)$, dist. $\mathcal{D}$ over $\mathcal{X}$, and labelling $f: \mathcal{X} \to \{\pm 1\}$,
if realizability holds w.r.t. $\mathcal{H}, \mathcal{D}, f$, then

$$L_{(\mathcal{D}, f)}(h) \leq \frac{1}{2} - \gamma \text{ w/ prob.} \geq 1-\delta \text{ over}$$

where $h$ is the ← returned hypothesis by the algorithm.

$m \geq m_{\mathcal{H}}(\delta)$ iid samples generated by $\mathcal{D}$ and labelled by $f$,

It seems that $h \notin \mathcal{H}$ is also allowed.

(2) $\mathcal{H}$ is $\gamma$-weak-learnable if it has a $\gamma$-weak-learner.

Rmk:

① $L_{(\mathcal{D}, f)}(h) \leq \frac{1}{2} - \gamma$ means the hypothesis returned is only a little (quantified by $\gamma$) better then random guess.

② The potential advantage of weak learning is that maybe there is an efficient algorithm for it. This motivates the opening question in this Chap.

Note: (Possible approach to construct weak learner).
For given $\mathcal{H}$, choose a simple $B$ and consider $ERM_B$ instead of $ERM_{\mathcal{H}}$.

Example:

$\mathcal{X} := \mathbb{R}$

$\mathcal{H} := \{3\text{-piece classifiers}\} = \{h_{\theta_1, \theta_2, b} : \theta_1, \theta_2 \in \mathbb{R}, \theta_1 < \theta_2, b \in \{\pm 1\}\}$

where $h_{\theta_1, \theta_2, b}(x) = \begin{cases} +b & \text{if } x < \theta_1 \text{ or } x > \theta_2 \\ -b & \text{if } \theta_1 \leq x \leq \theta_2. \end{cases}$

$B := \{\text{decision stumps}\} = \{x \mapsto \text{sign}(x-\theta)\cdot b : \begin{smallmatrix}\theta \in \mathbb{R}\\ b \in \{\pm 1\}\end{smallmatrix}\}$

**Prop**

$\text{ERM}_B$ is a $\gamma$-weak-learner for $\mathcal{H}$, for $\gamma = 1/12$.

⟨Pf⟩

Not hard. See P133, upper part. ⌗

**10.1.1 Efficient Implementation of ERM for Decision Stumps.**

↖ __SKIPPED__

**10.2 AdaBoost.**

AdaBoost is an algorithm having access to a weak learner and finds a hypothesis w/ low empirical risk. The following is the pseudo code.

__Ada Boost__ (for binary classification).

(Input)

  training set : $S = (x_1, y_1), \cdots, (x_m, y_m)$.
  weak learner : WL
  number of rounds : $T$.

(Initialize): $D^{(1)} = (\frac{1}{m}, \cdots, \frac{1}{m})$. ( a dist. on $x_1, \cdots, x_m$ )

for $t = 1, \cdots, T$ :

  invoke weak learner : $h_t = WL(D^{(t)}, S)$.
  compute : $\varepsilon_t = \sum_{i=1}^{m} D_i^{(t)} \cdot 1_{[y_i \neq h_t(x_i)]}$
  let $w_t = \frac{1}{2} \log(\frac{1}{\varepsilon_t} - 1)$ (the smaller $\varepsilon_t$, the larger $w_t$).
  update the dist. $D_i^{(t+1)} = \dfrac{D_i^{(t)} \exp(-w_t y_i h_t(x_i))}{\sum_{j=1}^{m} D_j^{(t)} \exp(-w_t y_j h_t(x_j))}$ for $i = 1, \cdots, m$.

(Output) :

  the hypothesis $h_S(x) = \text{sign}(\sum_{t=1}^{T} w_t h_t(x))$.

(Rmk) :

① The smaller $\varepsilon_t$ is, the larger $w_t$ is.
  i.e. if $h_t$ is more correct, then it has more weight.

② If $x_i$ is wrongly classified, then
  $-w_t y_i h_t(x_i) = w_t > 0$. $\Rightarrow D_i^{(t+1)}$ exceptionally large.
  $\Rightarrow x_i$ is taken care of more at the next round.
  On the other hand, if $x_i$ : correctly classified,
  then $D_i^{(t+1)}$ : exceptionally small. $\Rightarrow x_i$ : less taken care of next round.

---

Other than the above intuitions, we also have the following theoretical guarantee :

**Thm 10.2**

  $S$ : training set.

Suppose at each iteration of AdaBoost, the weak learner returns a hypothesis w/ $\varepsilon_t \leq \frac{1}{2} - \gamma$. Then

$$L_S(h_S) = \frac{1}{m} \sum_{i=1}^{m} 1_{[h_S(x_i) \neq y_i]} \leq e^{-2\gamma^2 T}.$$

(Rmk) :

By weak learnability, $\varepsilon_t \leq \frac{1}{2} - \gamma$ may fail w/ prob. at most $\delta$.
By union bd, the assumption in Thm 10.2 holds w/ prob. $\geq 1 - T\delta$.

**10.3 Linear Comb. of Base Hypotheses.**

$B$ : a hypothesis class ($B$ stands for base).
$T \in \mathbb{N}$.
Define

$$\boxed{L(B,T)} = \{x \mapsto \text{sign}(\sum_{t=1}^{T} w_t h_t(x)) : \begin{smallmatrix}w \in \mathbb{R}^T\\ h_t \in B, \forall t\end{smallmatrix}\}$$

In this sec, we estimate the VCdim of $L(B,T)$ (bound) in terms of ① the VCdim of $B$ and ② $T$.

**10.3.1 VC dim of $L(B,T)$.**

**Lemma 10.3**

$B$ : a base hypothesis class
Assume both $T$ and $VCdim(B) \geq 3$. Then

$$VCdim(L(B,T)) \leq T \cdot (VCdim(B)+1) \cdot (3\log(T(VCdim(B)+1))+2)$$

(Rmk) :

Thm 10.2 allows us to use AdaBoost to reduce empirical risk while Lemma 10.3 guarantees the true risk not far from empirical risk.

**10.4 AdaBoost for Face Detection.**

see P140-141 or the paper by Viola and Jones for details.

# Chap 11  Model Selection and Validation.

## 11.1 Model Selection using SRM.

Given $H_1, H_2, H_3, \cdots$, a countable seq. of hypothesis classes.

Assume each $H_d$ enjoys the UC property w/

$$m_{H_d}^{UC}(\varepsilon, \delta) \leq \frac{g(d) \log(1/\delta)}{\varepsilon^2}, \text{ where}$$

$g: \mathbb{N} \to \mathbb{R}$ : some monotone increasing fun.

Rmk:

① For binary classification (see Thm 6.8), take

$g(d) = C_2(D+1)$, where $D = VC\dim(H_d)$.

② For AdaBoost, see Thm 10.3.

Taking $w(d) = \frac{6}{\pi^2} \cdot \frac{1}{d^2}$, Thm 7.4 implies

$$L_{\mathcal{D}}(h) \leq L_S(h) + \sqrt{\frac{g(d) \cdot (\log(1/\delta) + 2\log(d) + \log(\pi^2/6))}{m}}$$

Recall that SRM will search for "$d$" and "$h \in H_d$" that minimizes the RHS. In this sense, SRM is doing model selection.

Issue:

In many practical situations, the bd on the RHS is pessimistic and not useable.!!

## 11.2 Validation.

### 11.2.1 Hold-out Set.

Let $V = (x_1, y_1), \cdots, (x_{m_v}, y_{m_v})$ be $m_v$ fresh examples sampled w.r.t. $\mathcal{D}$. (v: validation).

We have:

Thm 11.1

$h \in H$, some predictor.

Assume the loss fun. is in $[0,1]$.

Then, $\forall \delta \in (0,1)$,

$$|L_V(h) - L_{\mathcal{D}}(h)| \leq \sqrt{\frac{\log(2/\delta)}{2m_v}}$$

w/ prob. $\geq 1-\delta$ over the choices of a validation set $V$ of size $m_v$.

⟨Pf⟩

Recall the Hoeffding's ineq.:

Lemma 4.5 (Hoeffding's ineq.)

$\Theta_1, \cdots, \Theta_m$ : iid w/ $E[\Theta_i] = \mu$ and $P[a \leq \Theta_i \leq b] = 1$.

Then, $\forall \varepsilon > 0$,

$$P\left[|\frac{1}{m}\sum_{i=1}^{m} \Theta_i - \mu| > \varepsilon\right] \leq 2 \cdot \exp(-2m\varepsilon^2/(b-a)^2).$$

Take $\Theta_i = \ell(h, (x_i, y_i))$, where $(x_i, y_i) \sim \mathcal{D}$.

Then $\mu = L_{\mathcal{D}}(h)$ and $\frac{1}{m}\sum_{i=1}^{m_v} \Theta_i = L_V(h)$.

Solving $2\exp(-2m_v \varepsilon^2) = \delta$ for $\varepsilon$, the result follows. #

Rmk:

① Bound in Thm 11.1 does not depend on the algorithm to construct $h$ and is often tighter than that given by SRM.

② Price we pay is that we need additional fresh examples other than the training set.

In practice, $V$ is constructed by holding out part of the whole examples. Thus, the validation set is also called a hold-out set.

### 11.2.2 Validation for Model Selection.

Validation can be used for model selection as follows:

Assume $h_i$ is the returned hypothesis from $H_i$, $i = 1, \cdots, r$. Denote $H = \{h_1, \cdots, h_r\}$.

Then we have:

Thm 11.2

$H = \{h_1, \cdots, h_r\}$, an arbitrary set of predictors.

Assume the loss fun. is in $[0,1]$.

$V$: a validation set of size $m_v$, sampled indep. of $H$.

Then, w/ prob. $\geq 1-\delta$ over the choice of $V$,

$$|L_{\mathcal{D}}(h) - L_V(h)| \leq \sqrt{\frac{\log(2|H|/\delta)}{2m_v}}, \forall h \in H.$$

⟨Pf⟩

Use Hoeffding and union bd. #

Rmk: (Gist)

Too many models may result in overfitting.