

$$\text{sign}(w_t h_t(x))$$

The output of Adaboost is a composition of half space over the predictions of the T weak hypotheses.

Assume that the weak hypotheses belong to the base hypothesis class, B .

The hypothesis produced by Adaboost belongs to a hypothesis class of linear half spaces defined over predictions by T hypotheses $h_t, t \in [1, T]$ where $h_t \in B$.
 Base Hypothesis class \nearrow stages

We denote the hypothesis class of Adaboost as $L(B, T)$.
 Linear Half Space \nwarrow

$$L(B, T) = \left\{ x \mapsto \text{sign}\left(\sum_{t=1}^T w_t h_t(x)\right) : \underline{w} \in \mathbb{R}^T, h_t \in B \right\}$$

To learn a hypothesis $h \in L(B, T)$, the AdaBoost learner needs to learn $\underline{w} \in \mathbb{R}^T$ and $h_t \in B, t \in [1, T]$.

During the inference phase, the learned h will classify a given test instance x by constructing $\varphi(x) = (\underline{h_1(x)}, \underline{h_2(x)}, \dots, \underline{h_T(x)})$ and then applying the half space defined by \underline{w} on $\varphi(x)$.

✓ VC dimension of Adaboost Hypothesis Class $L(B, T)$ is bounded by T times the VC dimension of the base hypothesis class B .

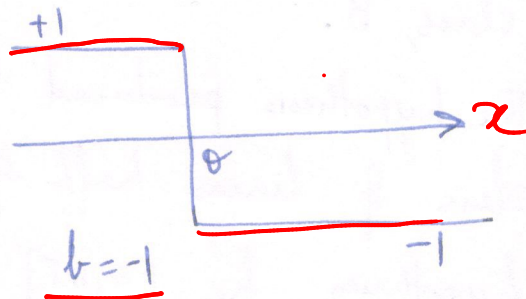
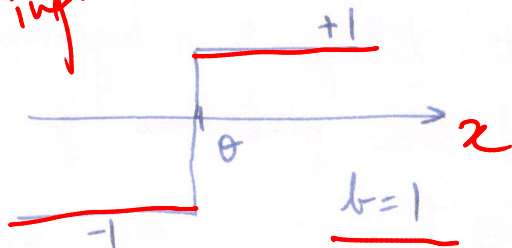
The empirical risk of Adaboost decreases with T .

Consider a base hypothesis class of Decision Stumps.



$$H_{DS} = \left\{ x \mapsto \text{sign}(x - \theta) \cdot b : \begin{array}{l} \theta \in \mathbb{R} \\ b \in \{\pm 1\} \end{array} \right\}$$

scalar input

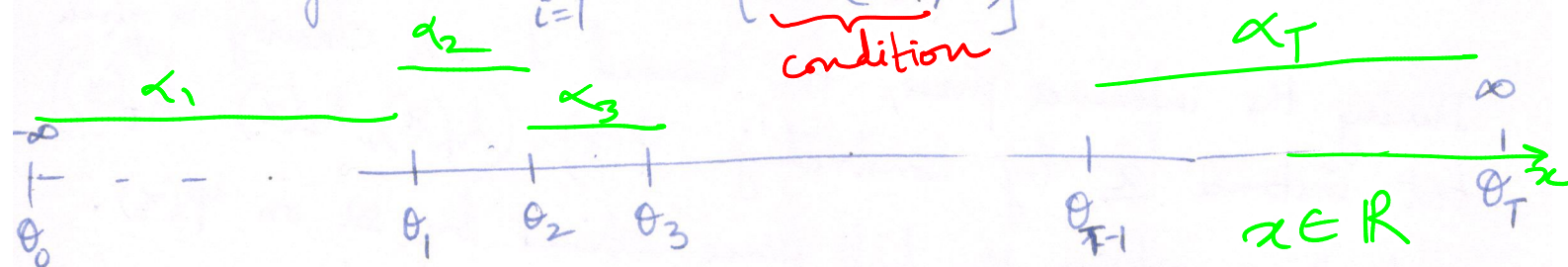


We show that Adaboost can use $B \equiv H_{DS}$ and construct a rather complex class of piecewise constant functions.

A piecewise constant function over T steps can be defined as

$$g_T(x) = \sum_{i=1}^T \alpha_i \mathbb{1}_{\{x \in (\theta_{i-1}, \theta_i]\}} \quad \forall i \quad \alpha_i \in \{\pm 1\}$$

indicator
condition



There are T pieces in $g_T(x)$.

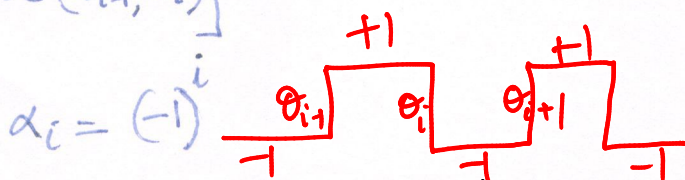
The class of half spaces over T decision stumps yields all the piece-wise constant classifiers with at most T pieces.

We demonstrate this by taking a specific example that

$$g_T(x) = \sum_{i=1}^T (-1)^i \mathbb{1}_{\{x \in (\theta_{i-1}, \theta_i]\}}$$

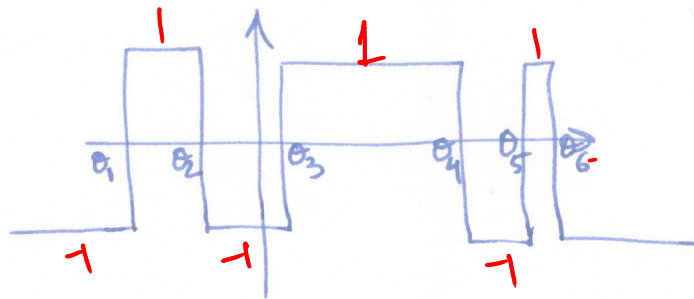
$$\alpha_i = (-1)^i$$

i.e.



Such a function $g_T(x) = \sum_{i=1}^T (-1)^i \mathbb{1}_{[x \in (\theta_{i-1}, \theta_i]}$

would look like a square wave



A given x would belong to only one interval.

The hypothesis produced by the Adaboost classifier

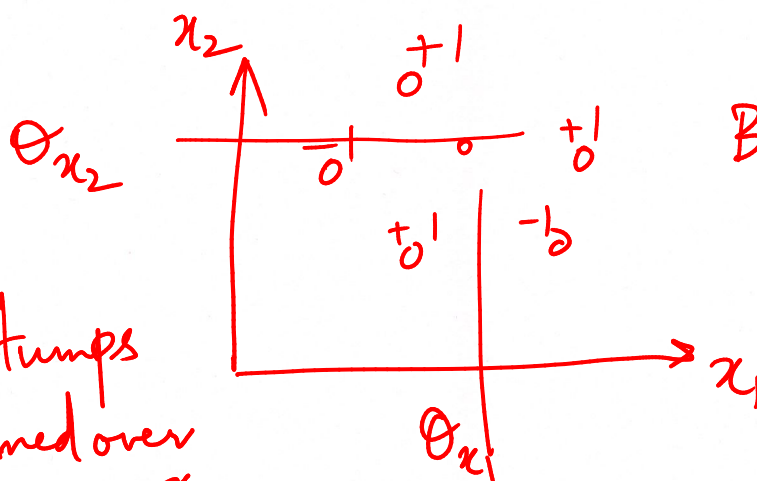
$$h(x) = \text{sign} \left(\sum_{t=1}^T w_t \text{sign}(x - \theta_{t+1}) \right)$$

Try it

with $w_1 = -0.5$
 $w_t = (-1)^t$

This will give an output same as $g_T(x)$

will behave exactly in the same way as $g_T(x)$



$$B = H_{DS}$$

$$\underline{x} \in \mathbb{R}^2$$

Decision stumps
 Can be defined over
 x_1 or x_2

supervisory input

We have seen construction of supervised machines when the training examples ^S were in the form of $\{(x, y)\}$

We used learning rules such as ERM, Hard margin, soft margin, etc. ^{target}

$x \rightarrow y$

Now we consider training data where examples are represented using features x but there is no target label.

Unsupervised learning.

We would like to learn the structure / distribution of the features $P(x)$. ^{density function} $P(x)$. ^{continuous valued features}

$P(x)$ is a data generating distribution.

We can generate observations x by sampling from $P(x)$
 $x \sim P(x)$.

Generative model can generate observations x

Now there are two questions.

^{parameterized}

1) What model would best describe the data generating distribution?

$x \in \mathbb{R}^d$ ^{d-dim space}

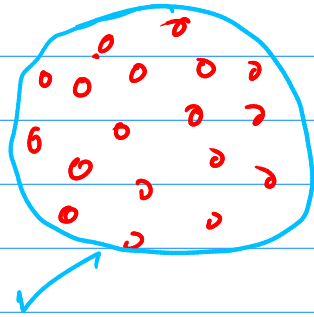
2) How do we estimate the parameters of the data generating distribution? ^{model chosen for}

✓ We use maximum likelihood estimation (MLE) for estimating the model parameters.

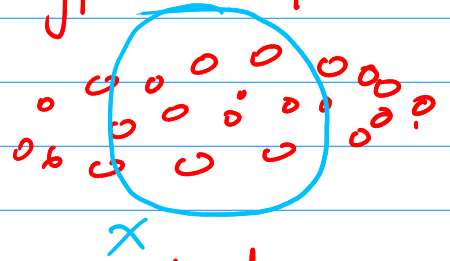
✓ We formulate a likelihood function which gives the likelihood (probability) of observing all the examples x_i in a training set S ^{generating}

hypersphere.

Multivariate
Gaussian with
 $C = I$
Covariance
matrix



hyper ellipse



manifold in a d-dim space

Choose

Decide the model

Shape of the distribution
→ parameterized model

MLE fit the model
(estimation of parameters)

↓
parameters are
calculated.