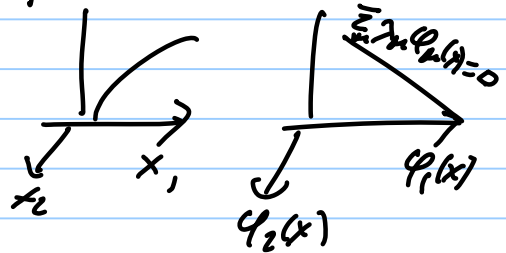# AdaBoost.

Note Title

AdaBoost is a method for combining a number of weak classifiers to make a strong classifier.

Input: set of weak classifiers $\{ \varphi_\mu(x) : \mu = 1 \text{ to } M \}$

labelled data $\{ (\underline{x}^i, y^i) : i = 1 \text{ to } N \}$

$$ y^i \in \{ -1, 1 \} $$

Output: strong classifier

$$ \text{sign} \left( \sum_{\mu=1}^{M} \lambda_\mu \, \varphi_\mu(x) \right) $$

$\{ \lambda_\mu \}$ weights / coefficients.

• The strong classifier is a plane in feature space $\{ \varphi_\mu(x) \}$.

In practice, most of the $\lambda_\mu = 0$.

The "selected" weak classifiers are those with $\lambda_\mu \neq 0$.

(2)    The task of AdaBoost is to
select weights $\{\lambda_\mu\}$ to make the strong
classifier as effective as possible.

   The motivation is that it is often
possible to obtain weak classifiers for
classification tasks — i.e. a weak classifier
that is effective 60% of the time. Want
to build a strong classifier — effective 99%
of the time — that is built by combining weak
classifiers.

   The combination is by weighted
summation        $\sum_\mu \lambda_\mu \rho_\mu (x)$

   Why linear weighted combination?
   Because we can use an efficient
algorithm — AdaBoost — to estimate them.

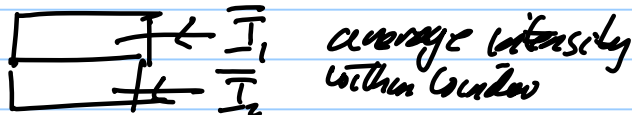# (3)   Example : Face Detection (Viola & Jones)

Face

Non Face. threshold

Training examples are a set of images $x^t$ labelled by $y^t = 1$ if image contains a face, $y^t = -1$ if not.
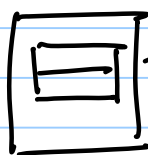
## Weak classifier :

$\bar{I}_1$ ← average intensity within window

$\bar{I}_2$

Face : if   $\bar{I}_1(x) - \bar{I}_2(x) > T$
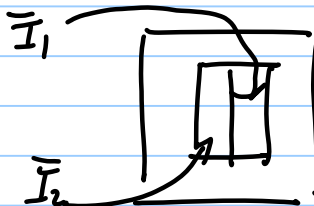
Intensity in forehead is bigger than intensity in eye region.

← window position

Forehead Detector.

Face : if   $\bar{I}_1 - \bar{I}_2 < \epsilon_1$

$\qquad\quad \bar{I}_2 - \bar{I}_1 < \epsilon_2$

$\bar{I}_1$

$\bar{I}_2$

Symmetry Detection

→ Faces symmetric.

where $\epsilon_1$ & $\epsilon_2$ are small constants.

Easy to get weak classifiers of this type — weak classifier is features $\bar{I}_1(x), \bar{I}_2(x) +$ threshold $T$ — but each weak classifier is only partially success. AdaBoost gives a way to select weak classifiers and combine them to make a strong classifier.

(4)  AdaBoost: Mathematical Descript.

Define  $Z[\lambda_1, \dots \lambda_M] = \sum_{t=1}^{N} e^{-y^i \sum_{\mu=1} \lambda_\mu \varphi_\mu(x^i)}$.

This a **convex** upper bound of the error rate of the strong classifier $S(x) = \text{sign}\left(\sum_{\mu=1}^{M} \lambda_\mu \varphi_\mu(x)\right)$.

← Identity function

Error Rate: $E[\lambda_i] = \sum_{i=1}^{N} \langle 1 - I(S(x_i), y_i)\rangle$  ← Empirical Risk.

whae  $I(S(x), y) = 1$, if $S(x) = y$  (correct)
answer

$= 0$, otherwise.

Claim :  $E[\lambda_1 \dots \lambda_M] \leq Z[\lambda_1 \dots \lambda_M]$

compare each term in the summcatoun  $\sum_{i=1}^{N}$

Case (i) If $S(x^i) = y^i$, then $\text{sign}\left(\sum_{\mu=1}^{M} \lambda_\mu \varphi_\mu(x^i)\right) = \text{sign } y^i$

So  $y^i \sum_{\mu=1}^{M} \lambda_\mu \varphi_\mu(x^i) = A \geq 0$  (Defines A)

Error Term  $\langle 1 - I(S(x^i), y^i)\rangle = 0$

$Z$ Term  $e^{-y^i \sum_\mu \lambda_\mu \varphi_\mu(x^i)} = e^{-A} \geq 0$  ✓

Case (ii) If $S(x^i) \neq y^i$, then  $y^i \sum_{\mu=1}^{M} \lambda_\mu \varphi_\mu(x^i) = -B < 0$
($B > 0$)

Error Term  $\langle 1 - I(S(x^i), y^i)\rangle = 1$

$Z$ term  $e^{-y^i(\sum_\mu \lambda_\mu \varphi_\mu(x_i))} = e^{B} \geq 1$  ✓

So  $Z$ term is bigger than error term in both cases.

## Goal: AdaBoost minimizes

(5) $Z[\lambda_1 \ldots \lambda_N]$. This will guarantee that the error rate is small (but not necessarily the minimum error rate).

Strategy to minimize $Z[\lambda_1 \ldots \lambda_N]$.

**Initialize** by $\lambda_1 = \ldots = \lambda_N = 0$.
(i.e. $H_0(x) = 0 \rightarrow$ no weak classifier selected).

**Minimize** $Z[\lambda_1 \ldots \lambda_N]$ by coordinate descent.

At time step $\ell$.

State $\lambda_1^\ell, \ldots \lambda_N^\ell$.

**For each** $i \rightarrow$ minimize $Z$ w.r.t. $\lambda_\mu$
with $\lambda_\nu^\ell$ fixed for $\nu \neq \mu$

**Solve** $\dfrac{\partial Z}{\partial \lambda_\mu} = 0$ to solve for $\hat{\lambda}_\mu$
for each $\mu$

Compute $Z[\lambda_1^\ell, \ldots \lambda_{\mu-1}^\ell, \hat{\lambda}_\mu, \lambda_{\mu+1}^\ell, \ldots \lambda_N^\ell]$ for each $\mu$
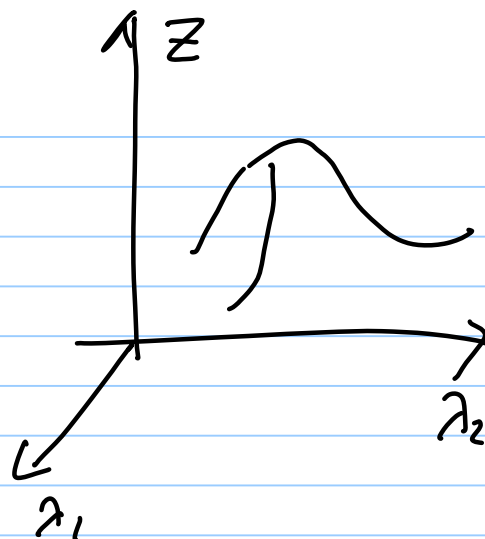
**select** $\hat{\mu} = \underset{\mu}{ARG\,MIN}\ Z[\lambda_1^\ell, \ldots \lambda_{\mu-1}^\ell, \lambda_\mu, \lambda_{\mu+1}^\ell, \ldots \lambda_N^\ell]$

**set** $\lambda_\nu^{\ell+1} = \lambda_\nu^\ell, \nu \neq \hat{\mu}, \lambda_{\hat{\mu}}^{\ell+1} = \hat{\lambda}_{\hat{\mu}}$.

(6)   Intuition: at each time
        ___step l.___

calculate how much you
can decrease $Z$ by changing
only one of the $\{\lambda_n\}$.
    select the $\lambda_{\hat\mu}$ which
maximally decreases $Z$,



   Each step of this algorithm is guaranteed
to decrease $Z$.
   So algorithm will converge to a minimum
of $Z$. ( But $Z$ is convex in $\lambda$, so the algorithm
converges to global min(mum).

   Why $Z$? Why this algorithm?

__Practical__ — we can compute $\frac{\partial Z}{\partial \lambda_\mu}$ and the
   minimum of $Z$ easily,

Why not steppest descent?   We want to keep
     most of the $\lambda_\mu = 0$ .   use as few non-zero
                                    $\lambda$'s as possible

(7)

# AdaBoost Algorithm.

Data $\{(x^i, y^i) : i = 1, .., N\}$

Set of weak classifiers. $\{\varphi_\mu(x), \mu = 1, ... M\}$.

For each weak classifier, divide the training data into two classes

$$W_\mu^+ = \{i : y \, \varphi_\mu(x^i) = 1\} \quad \varphi_\mu \text{ correct}$$
$$W_\mu^- = \{i : y \, \varphi_\mu(x^i) = -1\} \quad \varphi_\mu \text{ wrong.}$$
$$W_\mu^+ \cup W_\mu^- = \{1, ... N\}$$

At each time step $t$, define a set of "weights" for the training examples.

$$D_i^t = \frac{e^{-y^i \sum_{\mu=1}^{M} \lambda_\mu^t \varphi_\mu(x^i)}}{\sum_{i=1}^{N} e^{-y^i \sum_{\mu=1}^{M} \lambda_\mu^t \varphi_\mu(x^i)}} \qquad \sum_i D_i^t = 1$$
$$D_i^t \geq 0.$$

Gives bigger weights to data incorrectly classified by current "strong classifier"

i.e. $D_i^t$ is large if $y^i \sum_{\mu=1}^{M} \lambda_\mu^t \varphi_\mu(x^i) < 0$.

implies $\text{sign}\left(\sum_{\mu=1}^{M} \lambda_\mu^t \varphi_\mu(x^i)\right) \neq y^i.$

$D_i^t$ is small if $y^i \sum_{\mu=1}^{M} \lambda_\mu^t \varphi_\mu(x^i) > 0.$

implies $\text{sign}\left(\sum_{\mu=1}^{M} \lambda_\mu^t \varphi_\mu(x^i)\right) = y^i$

(8)

# Adaboost Algorithm.

Initialize $\quad \lambda_1 = \lambda_2 = \ldots = \lambda_N = 0$

At time step $\quad \lambda_1^\ell, \lambda_2^\ell \ldots \lambda_n^\ell$

For each $\mu$, calculate $\quad \Delta_\mu^t = \frac{1}{2} \log \left\{ \dfrac{\sum\limits_{i \in w_\mu^+} D_i^t}{\sum\limits_{i \in w_\mu^-} D_i^t} \right\}$

(change in $\Delta_\mu^t$ due to solving $\frac{\partial Z}{\partial \lambda_\mu} = 0$, see later)

calculate $\quad \sqrt{\sum\limits_{i \in w_\mu^+} D_i^t} \; \sqrt{\sum\limits_{i \in w_\mu^-} D_i^t}$

(change in $Z$ due to setting $\lambda_\mu$ to $\hat{\lambda}_\mu$, see later)

select $\quad \hat{\mu} = \underset{\mu}{ARG\;MIN} \; \sqrt{\sum\limits_{i \in w_\mu^+} D_i^t} \; \sqrt{\sum\limits_{i \in w_\mu^-} D_i^t}$

set $\quad \lambda_\nu^{t+1} = \lambda_\nu^t \quad, \quad \nu \neq \hat{\mu}$

$\qquad \lambda_{\hat{\mu}}^{t+1} = \lambda_{\hat{\mu}}^t + \Delta_{\hat{\mu}}^t .$

repeat, until convergence.

(9) Intuition for the $\sum_{i \in w_\mu^+} D_i^t$ and $\sum_{i \in w_\mu^-} D_i^t$ terms.

Initially, $D_i^{t=0} = \frac{1}{N}$    the data is equally weighted.

$\sum_{i \in w_\mu^+} D_i^{t=0}$ is the proportion of data that is correctly classified by $\varphi_\mu(.)$

$\sum_{i \in w_\mu^-} D_i^{t=0}$ is proportion incorrectly classified

i.e. $\sum_{i \in w_\mu^-} D_i^{t=0}$ is the normalized

error rate if we just use classifier $\varphi_\mu(.)$
— normalized by $\frac{1}{N}$. //

For $t \neq 0$, $\sum_{i \in w_\mu^+} D_i^t$ is data correctly classified by $\varphi_\mu(x)$ taking into account the previously selected classifiers (those for which $\lambda_\mu^t \neq 0$).

$\varphi_\mu(.)$ is a useless classifier if $\sum_{i \in w_\mu^+} D_i^t = \sum_{i \in w_\mu^-} D_i^t$
i.e. weighted error = $\frac{1}{2}$.

corresponds to $\Delta_\mu^t = 0$ (i.e. no change in weight)
$\lambda_\mu$

(10)

Term $\left|\sqrt{\sum\limits_{i \in w_\mu^+} D_i^t} \sqrt{\sum\limits_{i \in w_\nu^-} D_i^t}\right|$ is a

non-linear function of the weighted error rate

of $\varphi_\mu(.)$.

It can be rewritten as

$$\left|\sqrt{\left(\sum\limits_{i \in w_\mu^-} D_i^t\right)\left(1 - \sum\limits_{i \in w_\mu^+} D_i^t\right)}\right|$$

because $\sum\limits_{i \in w_\mu^+} D_i^t + \sum\limits_{i \in w_\mu^-} D_i^t = 1$

It's smallest values are if

$\sum\limits_{i \in w_\mu^-} D_i^t = 0$ , $\varphi_\mu$ has optimal weighted classified

$\sum\limits_{i \in w_\mu^-} D_i^t = 1$ , $\varphi_\mu$ worst classifier implies $-\varphi_\mu$ best classifier

its largest values are when

$\sum\limits_{i \in w_\mu^-} D_i^t = \frac{1}{2}$ , i.e. when weighted error is $\frac{1}{2}$, i.e. $\varphi_\mu()$ useless

## When does AdaBoost converge?

It stops when all weak classifiers are useless — ie when $\sum\limits_{i \in w_\mu^-} D_i^t = \frac{1}{2}$, for all $\mu$

In this case $\sqrt{\sum\limits_{i \in w_\mu^+} D_i^t} \; \sqrt{\sum\limits_{i \in w_\mu^-} D_i^t}$ takes its biggest (i.e. worst) value of $\frac{1}{2}$, for all $\mu$

The weight update $\Delta_\mu^t = \frac{1}{2} \log \left\{ \dfrac{\sum\limits_{i \in w_\mu^+} D_i^t}{\sum\limits_{i \in w_\mu^-} D_i^t} \right\}$ is $0$ for all $q_\mu(.)$ (since $\log 1 = 0$).

In general, at time step $t$ select the classifier $q_\mu(.)$ with smallest "weighted error rate" $\sqrt{\sum\limits_{i \in w_\mu^+} D_i^t} \; \sqrt{\sum\limits_{i \in w_\mu^-} D_i^t}$ $\qquad \Delta_{\hat{\mu}}^t = \frac{1}{2} \log \left\{ \dfrac{\sum\limits_{i \in w_\mu^+} D_i^t}{\sum\limits_{i \in w_\mu^-} D_i^t} \right\}$

Update $\lambda_{\hat{\mu}}^t \Rightarrow \lambda_{\hat{\mu}}^t + \Delta_{\hat{}}$

the smaller the weighted error rate — i.e. smaller $\sum\limits_{i \in w_\mu^-} D_i^t$ or $\sum\limits_{i \in w_\mu^+} D_i^t$ — then the bigger the change $\Delta_{\hat{\mu}}$.

(12) How does AdaBoost algorithm relate to AdaBoost mathematics?

AdaBoost mathematics requires:

(i) efficient solution of $\frac{\partial Z}{\partial \lambda_\mu} = 0$.

(ii) efficient computation of $Z$.

(1) $\frac{\partial Z}{\partial \lambda_\mu} = 0$

$$\frac{\partial Z}{\partial \lambda_\mu} = \sum_{i=1}^{N} \left\{-y^i \varphi_\mu(x^i)\right\} e^{-y^i \sum_{\mu=1}^{M} \lambda_\mu \varphi_\mu(x^i)}$$

set. $\lambda_\mu \to \lambda_\mu + \Delta_\mu$    solve for $\Delta_\mu$

$$\frac{\partial Z}{\partial \lambda_\mu} = 0 \Rightarrow \sum_{i=1}^{N} \left\{y^i \varphi_\mu(x^i)\right\} e^{-y^i \sum_{\mu=1}^{M} \lambda_\mu \varphi_\mu(x^i)} e^{-y^i \Delta_\mu \varphi_\mu(x^i)} = 0$$

$$\sum_{i=1}^{N} \left\{y^i \varphi_\mu(x^i)\right\} D_i \, e^{-y^i \Delta_\mu \varphi_\mu(x^i)} = 0.$$

Divide

$$\sum_{i=1}^{N} = \sum_{i \in \omega_\mu^+} + \sum_{i \in \omega_\mu^-}$$

$$\boxed{\begin{array}{l} \text{Recall } D_i \\ = \dfrac{e^{-y^i \sum_{\mu=1}^{M} \lambda_\mu \varphi_\mu(x^i)}}{\sum_i e^{-y^i \sum_\mu \lambda_\mu \varphi_\mu(x^i)}} \end{array}}$$

$$\sum_{i \in \omega^+} D_i \, e^{-\Delta_\mu} - \sum_{i \in \omega^-} D_i \, e^{\Delta_\mu} = 0.$$

$$e^{2\Delta_\mu} = \left(\sum_{i \in \omega_\mu^+} D_i\right) \bigg/ \left(\sum_{i \in \omega_\mu^-} D_i\right)$$

$$\Delta_\mu = \frac{1}{2} \log\left\{\sum_{i \in \omega_\mu^+} D_i \bigg/ \sum_{i \in \omega_\mu^-} D_i\right\} . \; /\!/$$

(2) Computation of $Z$.

$$Z[\lambda_1, \ldots, \lambda_\mu + \Delta_\mu, \lambda_{\mu+1}, \ldots, \lambda_N]$$

$$= \sum_{i=1}^{N} e^{-y^i \left\{ \sum_{\mu=1}^{m} \lambda_\mu \varphi_\mu(x^i) + \Delta_\mu \varphi_\mu(x^i) \right\}}$$

$$= K \sum_{i=1}^{N} D_i \, e^{-y^i \Delta_\mu \varphi_\mu(x^i)}$$

where $K = \sum_{i=1}^{N} D_i \, e^{-y^i \Delta_\mu \varphi_\mu(x^i)}$
is independent of $\mu$.

$$Z[\lambda_1, \ldots, \lambda_\mu + \Delta_\mu, \ldots, \lambda_N]$$

$$= K \left\{ \sum_{i \in W_\mu^+} D_i \, e^{-\Delta_\mu} + \sum_{i \in W_\mu^-} D_i \, e^{\Delta_\mu} \right\}$$

$$= 2K \sqrt{\sum_{i \in W_\mu^+} D_i} \sqrt{\sum_{i \in W_\mu^-} D_i}$$

using $e^{\Delta_\mu}$ from previous page.

**Hence**, coordinate descent reduces to
computing $\Delta_\mu = \frac{1}{2} \log \left( \dfrac{\sum_{i \in W_\mu^+} D_i}{\sum_{i \in W_\mu^-} D_i} \right)$  $\rightarrow$ how much to
change $\lambda_\mu$

Solving $\hat{\mu} = \underset{\mu}{\text{ARG MIN}} \left\{ \sum_{i \in W_\mu^+} D_i \, e^{-\Delta_\mu} + \sum_{i \in W_\mu^-} D_i \, e^{\Delta_\mu} \right\}$

to find best $\lambda_\mu$ to change.

# Error to Avoid in AdaBoost.

Once a weak classifier $\varphi_\mu(\cdot)$ has been selected, it can be selected again.

This should be obvious from the coordinate descent formulation — if you decide to update $\lambda_\mu$ at time step $t$, then you can also update $\lambda_\mu$ at a later time step.

## Probabilistic Interpretation.

It can be shown (Friedman, Hastie, Tibshirani) that AdaBoost relates to logistic regression.

$$P(y \mid x) = \frac{e^{y \sum_\mu \lambda_\mu \varphi_\mu(x)}}{e^{\sum_\mu \lambda_\mu \varphi_\mu(x)} + e^{-\sum_\mu \lambda_\mu \varphi_\mu(x)}}$$

This result is asymptotic — only true in the limit as the number of samples $N$ becomes infinitely large.

Note: standard sigmoid regression means that you specify a small number of features $\varphi_\mu(x)$ that are not necessarily binary-valued.

(15)

The main advantage of AdaBoost is that you can specify a large set of weak classifiers and the algorithm decides which weak classifiers to use – by assigning them non-zero $\lambda_\mu$.

Standard logistic regression only uses a small set of features (like weak classifiers).

SVM uses the kernel trick $K(\underline{x}, \underline{x}') = \underline{\varphi}(\underline{x}) \cdot \underline{\varphi}(\underline{x}')$ to simplify the dependence on $\underline{\varphi}(\underline{x})$, but doesn't say how to select $K(\cdot, \cdot)$ or $\underline{\varphi}(\cdot)$. //

Multilayer perceptron can be interpreted as selecting weak classifiers ⇒ but in a non-optimal manner.

Recent work, suggests that AdaBoost can be improved by making it more similar to logistic regression.

AdaBoost can be extended to multiclasses.