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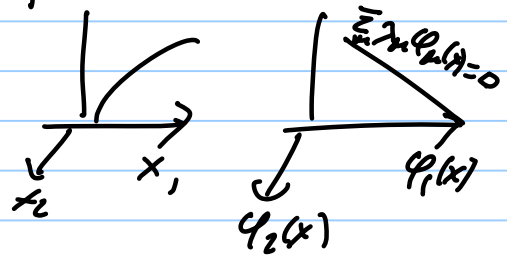
AdaBoost is a method for combining a number of weak classifiers to make a strong classifier.

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

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

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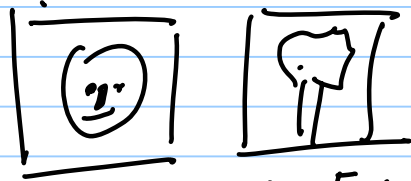
The motivation is that it is often possible to obtain weak classifier for classification tasks — i.e. a weak classifier that is effective 60% of the time. Want to build a strong classifier — effective 99% of the time — that is built by combining weak classifiers.

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).

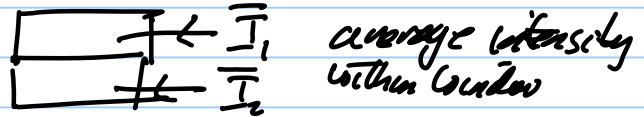

Face          Non Face.
              threshold.

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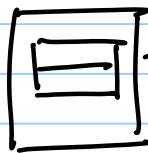contains a face, $y^t = -1$ if not.

Weak classifier:

$\bar{I}_1 \leftarrow$ average intensity within window
$\bar{I}_2$

Face: if $\quad \bar{I}_1(x) - \bar{I}_2(x) > T$

Intensity in forehead is
bigger than intensity in eye
region.

$\leftarrow$ window position

Forehead
Detector.

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

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

$\bar{I}_1$

$\bar{I}_2$

Symmetry
Detector
$\rightarrow$ 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)                    (Algorithm later)

## AdaBoost: Mathematical Description.

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

This is an upper bound of the error rate of the strong classifier $S(x) = \text{sign}\left(\overset{M}{\underset{\mu=1}{\sum}} \lambda_\mu \varphi_\mu(x)\right)$.

Error Rate: $E[\bar{\lambda}\bar{\lambda}_i] = \overset{N}{\underset{t=1}{\sum}} \langle 1 - I(S(x^t), y^t)\rangle$

where $\quad I(S(x^t), y^t) = 1,\ \ if\ \ S(x^t) = y^t \quad (correct)$ answer

$\qquad\qquad\qquad\qquad = 0,\ \ otherwise.$

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

compare each term in the summation $\overset{\tilde{N}}{\underset{t=1}{\sum}}$

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

$\qquad$ So $\quad y^t \overset{M}{\underset{\mu=1}{\sum}} \lambda_\mu \varphi_\mu(x^t) = A \geq 0 \quad$ (Defines A)

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

$\qquad$ Z Term $\quad e^{-y^t \sum_\mu \lambda_\mu \varphi_\mu(x^t)} = e^{-A} \geq 0 \quad \checkmark$

Case (ii) If $S(x^t) \neq y^t$, then $\quad y^t \overset{M}{\underset{\mu=1}{\sum}} \lambda_\mu \varphi_\mu(x^t) = -B < 0$

$\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad (B > 0)$

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

$\qquad$ Z term $\quad e^{-y^t(\sum_\mu \lambda_\mu \varphi_\mu(x^t))} = e^{B} > 1 \quad \checkmark$

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

## Goal: AdaBoost minimizer

(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 <u>coordinate descent</u>.

At time step $\ell$.

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

For each $i \rightarrow$ minimize $Z$ w.r.t. $\lambda_i$ with $\lambda_j^\ell$ fixed for $j \neq i$.

<u>Solve</u> $\dfrac{\partial Z}{\partial \lambda_i} = 0$ to solve for $\hat{\lambda}_i$ for each $i$.

Compute $Z[\lambda_1^\ell, \ldots \lambda_{i-1}^\ell, \hat{\lambda}_i, \lambda_{i+1}^\ell, \ldots \lambda^\ell]$ for each.

select $\hat{i} = \underset{i}{\text{ARGMIN}} \; Z[\lambda_1^\ell, \ldots \lambda_{i-1}^\ell, \hat{\lambda}_i, \lambda_{i+1}^\ell, \ldots \lambda_n^\ell]$

<u>set</u> $\lambda_j^{\ell+1} = \lambda_j^\ell, \; j \neq \hat{i}, \; \lambda_{\hat{i}}^{\ell+1} = \hat{\lambda}_{\hat{i}}$.

(6)    Intuition: at each time
            step $\ell$.

Calculate how much you
can decrease $Z$ by changing
only one of the $\{\lambda_i\}$.

    select the $\hat{\lambda_i}$ 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 minimum).

    why $Z$? why this algorithm?

Practical — we can compute $\frac{\partial Z}{\partial \lambda_i}$ and the
        minimum of $Z$ easily.

(7)                    AdaBoost Algorithm.

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

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

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

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

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

$$D_t^\ell = \frac{e^{-y^t \sum_{\mu=1}^{m} \lambda_\mu^\ell \varphi_\mu(x^t)}}{\sum_t e^{-y^t \sum_{\mu=1}^{M} \lambda_\mu^\ell \varphi_\mu(x^t)}} \qquad \sum_t D_t^\ell = 1$$
$$D_t^\ell > 0.$$

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

    i.e. $D_t^\ell$ is large if $y^t \sum_{\mu=1}^{m} \lambda_\mu^\ell \varphi_\mu(x^t) < 0$
        implies $\text{sign}(\sum_{\mu=1}^{M} \lambda_\mu^\ell \varphi_\mu(x^t)) \neq y^t$.
    $D_t^\ell$ is small if $y^t \sum_{\mu=1}^{M} \lambda_\mu^\ell \varphi_\mu(x^t) > 0$.
        implies $\text{sign}(\sum_{\mu=1}^{m} \lambda_\mu^\ell \varphi_\mu(x^t)) = y^t$

(8)

# Adaboost Algorithm.

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

At time step $\lambda_1^\ell, \lambda_2^\ell \ldots \lambda_N^\ell$

For each $i$, calculate $\Delta_i^\ell = \frac{1}{2} \log \left( \dfrac{\sum\limits_{t \in w_i^+} D_t^\ell}{\sum\limits_{t \in w_i^-} D_t^\ell} \right)$

(change in $\Delta_i^\ell$ due to solving $\frac{\partial Z}{\partial \lambda_i} = 0$, see later)

calculate $\sqrt{\overrightarrow{\sum\limits_{t \in w_i^+} D_t^\ell}} \cdot \sqrt{\overrightarrow{\sum\limits_{t \in w_i^-} D_t^\ell}}$

(change in $Z$ due to setting $\lambda_i$ to $\hat{\lambda}_i$, see later)

select $\hat{i} = \text{ARG MAX}_i \sqrt{\overrightarrow{\sum\limits_{t \in w_i^+} D_t^\ell}} \cdot \sqrt{\overrightarrow{\sum\limits_{t \in w_i^-} D_t^\ell}}$

set $\lambda_j^{\ell+1} = \lambda_j^\ell$, $j \neq i$

$\lambda_i^{\ell+1} = \lambda_i^\ell + \Delta_i^\ell$.

repeat, until convergence.

(9) Intuition for the $\sum_{t \in w_i^+} D_t^\ell$ and $\sum_{t \in w_i^-} D_t^\ell$ terms.

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

$\sum_{t \in w_i^+} D_t^{\ell=0}$ is the proportion of data that is correctly classified by $\varphi_i(x)$

$\sum_{t \in w_i^-} D_t^{\ell=0}$ is proportion incorrectly classified

. i.e. $\sum_{t \in w_i^-} D_t^{\ell=0}$ is the normalized error rate if we just use classifier $\varphi_i(x)$
— normalized by $\frac{1}{N}$. //

For $\ell \neq 0$, $\sum_{t \in w_i^+} D_t^\ell$ is data correctly classified by $\varphi_i(x)$ taking into account the previously selected classifiers (those for which $\lambda_j^\ell \neq 0$).

$\varphi_i(x)$ is a useless classifier if $\sum_{t \in w_i^+} D_t^\ell = \sum_{t \in w_i^-} D_t^\ell$
i.e. weighted error = $\frac{1}{2}$.

corresponds to $\Delta_i^\ell = 0$ (i.e. no change in weight)
$\lambda_i$

(10)

Term $\sqrt{\sum_{t \in w_i^+} D_t^\ell} \sqrt{\sum_{t \in w_i^-} D_t^\ell}$ is a

non-linear function of the weighted error rate
of $\varphi_i(x)$.

It can be rewritten as

$$\sqrt{\left(\sum_{t \in w_i^-} D_t^\ell\right)\left(1 - \sum_{t \in w_i^-} D_t^\ell\right)}$$

because $\sum_{t \in w_i^+} D_t^\ell + \sum_{t \in w_i^-} D_t^\ell = 1$

It's smallest values are if

$\sum_{t \in w_i^-} D_t^\ell = 0$ , $\varphi_i(x)$ has optimal weighted
classified

$\sum_{t \in w_i^-} D_t^\ell = 1$ , $\varphi_i(x)$ worst possible classify
$\rightarrow$ so implies $-\varphi_i(x)$ best possible class

its largest values are when

$\sum_{t \in w_i^-} D_t^\ell = \frac{1}{2}$ , i.e. when weighted error
is $\frac{1}{2}$, $\varphi_i(x)$ useless

(11)

### When does AdaBoost converge?

It stops when all weak classifiers are useless — i.e. when $\sum_{t \in w_i^-} D_t^\ell = \frac{1}{2}$, for all $i$.

In this case $\sqrt{\sum_{t \in w_i^+} D_\epsilon^\ell} \sqrt{\sum_{t \in w_i^-} D_t^\ell}$ takes its biggest (i.e. worst) value of $\frac{1}{2}$, for all $i$.

The weight update $\Delta_i^\ell = \frac{1}{2} \log \left\{ \frac{\sum_{t \in w_i^+} D_t^\ell}{\sum_{t \in w_i^-} D_t^\ell} \right\}$ is $0$ for all $q_i(x)$ (since $\log 1 = 0$).

In general, at time step $\ell$ select the classifier $q_i(x)$ with smallest "weighted error rate"

$$\sqrt{\sum_{t \in w_i^+} D_\epsilon^\ell} \sqrt{\sum_{t \in w_i^-} D_t^\ell} \qquad \Delta_i^\ell = \frac{1}{2} \log \left\{ \frac{\sum_{t \in w_i^+} D_t^\ell}{\sum_{t \in w_i^-} D_t^\ell} \right\}$$

Update $\lambda_i^\ell \Rightarrow \lambda_i^\ell + \Delta_i^\ell$

the smaller the weighted error rate — i.e. smaller $\sum_{t \in w_i^-} D_t^\ell$ or $\sum_{t \in w_i^+} D_t^\ell$ — then the bigger the change $\Delta_i^\ell$.

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

AdaBoost mathematics requires:

(i) efficient solutions of $\frac{\partial Z}{\partial \lambda_i} = 0$.

(ii) efficient computation of $Z$.

(i) $\frac{\partial Z}{\partial \lambda_i} = 0$

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

set.   $\lambda_i \to \lambda_i + \Delta_i$        solve for $\Delta_i$.

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

$$\sum_{t=1}^{N} \left(y^t \varphi_i(x^t)\right) D_t \, e^{-y^t \Delta_i \varphi_i(x^t)} = 0.$$

Recall $D_t$
$$= \frac{e^{-y^t \sum_{\mu=1}^{M} \lambda_\mu \varphi_\mu(x^t)}}{\sum_t e^{-y^t \sum_\mu \lambda_\mu \varphi_\mu(x^t)}}$$

Divide $\sum_{t=1}^{N} = \sum_{t \in \omega_i^+} + \sum_{t \in \omega_i^-}$

$$\sum_{t \in \omega_i^+} D_t \, e^{-\Delta_i} - \sum_{t \in \omega_i^-} D_t \, e^{\Delta_i} = 0.$$

$$e^{2\Delta_i} = \left(\sum_{t \in \omega_i^+} D_t\right) / \left(\sum_{t \in \omega_i^-} D_t\right)$$

$$\Delta_i = \frac{1}{2} \log\left\{\sum_{t \in \omega_i^+} D_t / \sum_{t \in \omega_i^-} D_t\right\}. \; /\!/$$

(13)  (2)  Computation of $Z$.

$$Z[\lambda_1, \ldots, \lambda_i, \lambda_i + \Delta_i, \ldots \lambda_N]$$

$$= \sum_{t=1}^{\tilde{N}} e^{-y^t \left\{ \sum_{\mu=1}^{\tilde{N}} \lambda_\mu \phi_\mu(x^t) + \Delta_i \phi_i(x^t) \right\}}$$

$$= K \sum_{t=1}^{\tilde{N}} D_t \, e^{-y^t \Delta_i \phi_i(x^t)}$$

where $K = \sum_{t=1}^{\tilde{N}} D_t e^{-y^t \Delta_i \phi_i(x^t)}$

is independent of $i$.

$$Z[\lambda_1, \ldots \lambda_i, \lambda_i + \Delta_y, \ldots \lambda_N]$$

$$= K \left\{ \sum_{t \in w_i^+} D_t \, e^{-\Delta_i} + \sum_{t \in w_i^-} D_t \, e^{\Delta_i} \right\}$$

$$= 2K \sqrt{\sum_{t \in w_i^+} D_t} \sqrt{\sum_{t \in w_i^-} D_t}$$

using $e^{\Delta_i}$ from previous page.

Hence, coordinate descent reduces to computing $\Delta_i = \frac{1}{2} \log \left( \dfrac{\sum_{t \in w_i^+} D_t}{\sum_{t \in w_i^-} D_t} \right)$   $\rightarrow$ how much to change $\lambda_i$

solving $\hat{i} = \underset{i}{\text{arg max}} \left\{ \sum_{t \in w_i^+} D_t \, e^{-\Delta_i} + \sum_{t \in w_i^-} D_t \, e^{\Delta_i} \right\}$

to find best $\lambda_i$ to change.

# Error to Avoid in AdaBoost.

Once a weak classifier $\varphi_i(x)$ has been selected, it can be selected again.

This should be obvious from the coordinate descent formulation — if you decide to update $\lambda_i$ at time step $\ell$, then you can also update $\lambda_i$ 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 $w$ 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 $\lambda_i$.

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 $\Rightarrow$ 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.

(8) Second Reason:

Adaboost converges to the posteriors.

$$p(\omega=1|x) = \frac{e^{\sum_{\mu=1}^{M}\lambda_\mu \varphi_\mu(x)}}{e^{\sum_{\mu=1}^{M}\lambda_\mu \varphi_\mu(x)} + e^{-\sum_{\mu=1}^{M}\lambda_\mu \varphi_\mu(x)}}$$

$$p(\omega=-1|x) = \frac{e^{-\sum_{\mu=1}^{M}\lambda_\mu \varphi_\mu(x)}}{e^{\sum_{\mu=1}^{M}\lambda_\mu \varphi_\mu(x)} + e^{-\sum_{\mu=1}^{M}\lambda_\mu \varphi_\mu(x)}}$$

$$\frac{\partial Z}{\partial \lambda_\nu} = \sum_{i=1}^{N} \omega_i \varphi_\nu(x_i) e^{-\omega_i \sum_{\mu=1}^{M}\lambda_\mu \varphi_\mu(x_i)} = 0$$

hence,

$$\sum_{i=1}^{N} \delta_{\omega_i,1} \varphi_\nu(x_i) e^{-\sum_{\mu=1}^{M}\lambda_\mu \varphi_\mu(x_i)} = \sum_{i=1}^{N} \delta_{\omega_i,-1} \varphi_\nu(x_i) e^{+\sum_{\mu=1}^{M}\lambda_\mu \varphi_\mu(x_i)}$$

$$\forall \nu.$$

$$\frac{1}{N}\sum_{i=1}^{N} \delta_{\omega_i,1}\, \delta_{x_i,x} \to p(\omega=1,x)$$

$$\frac{1}{N}\sum_{i=1}^{N} \delta_{\omega_i,-1}\, \delta_{x_i,x} \to p(\omega=-1|x) \cdot \mathbin{/\!/}$$

$$\frac{1}{N}\sum_{i=1}^{N} \delta_{\omega_i,1}\, f(x_i) = \sum_x \left( \frac{1}{N}\sum_{i=1}^{N} \delta_{\omega_i,1}\, \delta_{x_i,x}\, f(x) \right)$$

$$\frac{1}{N}\sum_{i=1}^{N} \delta_{\omega_i,1}\, f(x_i) \to \sum_x p(\omega=1,x)\, f(x) \qquad \forall \nu.$$

So,

$$\sum_x p(\omega=1,x)\, \varphi_\nu(x)\, e^{\sum_{\mu=1}^{M}\lambda_\mu \varphi_\mu(x)} = \sum_x p(\omega=-1,x)\, \varphi_\nu(x)\, e^{-\sum_{\mu=1}^{M}\lambda_\mu \varphi_\mu(x)}$$

$$p(\omega=1|x) = \frac{e^{+\sum_{\mu=1}^{M}\lambda_\mu \varphi_\mu(x)}}{Z} \quad , \quad p(\omega=-1|x) = \frac{e^{-\sum_{\mu=1}^{M}\lambda_\mu \varphi_\mu(x)}}{Z}$$