

Boosting

Professor Ameet Talwalkar

Slide Credit: Professor Fei Sha

Outline

- 1 Administration
- 2 Review of last lecture
- 3 Boosting

Midterm

- Next Tuesday in class from 8am - 9:50am
- Completely closed-book (no notes allowed)
- Will include roughly 10 short answer questions and 2 long questions
 - ▶ Short questions should take no more than 5 minutes on average
 - ▶ Long questions should take no more than 20 minutes each
- Covers all material through (and including) this Thursday's lecture
 - ▶ Goal is to test conceptual understanding of the course material
 - ▶ Suggestion: carefully review lecture notes and problem sets
- Office hours / Section
 - ▶ me: Today 10:00am-11:00am (Boelter 4531F)
 - ▶ Amogh: Monday 11:30a - 12:30p; Friday 2:30p - 3:30p (Boelter 2432)
 - ▶ Nikos: Monday 3:30p - 4:30p; Wednesday 2:00p - 3:00p (Boelter 2432)
 - ▶ Section: Friday 12:00pm - 1:50pm (Kinsey Teaching Pavilion 1200B)

Other Announcements

- Nikos will be lecturing on Thursday about Neural Networks

Outline

- 1 Administration
- 2 Review of last lecture
- 3 Boosting

Support vector machines (SVM)

Primal Formulation

$$\begin{aligned} \min_{\mathbf{w}, b, \xi} \quad & \frac{1}{2} \|\mathbf{w}\|_2^2 + C \sum_n \xi_n \\ \text{s.t.} \quad & y_n [\mathbf{w}^T \phi(\mathbf{x}_n) + b] \geq 1 - \xi_n \quad \text{and} \quad \xi_n \geq 0, \quad \forall n \end{aligned}$$

Two equivalent interpretations

Support vector machines (SVM)

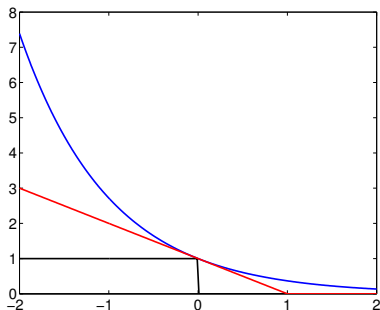
Primal Formulation

$$\begin{aligned} \min_{\mathbf{w}, b, \xi} \quad & \frac{1}{2} \|\mathbf{w}\|_2^2 + C \sum_n \xi_n \\ \text{s.t.} \quad & y_n [\mathbf{w}^T \phi(\mathbf{x}_n) + b] \geq 1 - \xi_n \quad \text{and} \quad \xi_n \geq 0, \quad \forall n \end{aligned}$$

Two equivalent interpretations

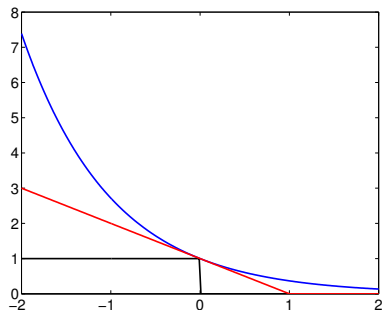
- Geometric: Maximizing (soft) margin
- Optimization: Minimize hinge loss with L2 regularization

Hinge Loss



- Upper-bound for 0/1 loss function (black line)
- Convex surrogate to 0/1 loss

Hinge Loss



- Upper-bound for 0/1 loss function (black line)
- Convex surrogate to 0/1 loss, though others exist as well
 - ▶ Hinge loss less sensitive to outliers than exponential (or logistic) loss
 - ▶ Logistic loss has a natural probabilistic interpretation
 - ▶ We can optimize exponential loss efficiently in a greedy manner (Adaboost)

Constrained Optimization

Primal Formulation

$$\begin{aligned} \min_{\mathbf{w}, b, \xi} \quad & \frac{1}{2} \|\mathbf{w}\|_2^2 + C \sum_n \xi_n \\ \text{s.t.} \quad & y_n [\mathbf{w}^T \phi(\mathbf{x}_n) + b] \geq 1 - \xi_n \quad \text{and} \quad \xi_n \geq 0, \quad \forall n \end{aligned}$$

- When working with constrained optimization problems with inequality constraints, we can write down primal and dual problems

Constrained Optimization

Primal Formulation

$$\begin{aligned} \min_{\mathbf{w}, b, \xi} \quad & \frac{1}{2} \|\mathbf{w}\|_2^2 + C \sum_n \xi_n \\ \text{s.t.} \quad & y_n [\mathbf{w}^T \phi(\mathbf{x}_n) + b] \geq 1 - \xi_n \quad \text{and} \quad \xi_n \geq 0, \quad \forall n \end{aligned}$$

- When working with constrained optimization problems with inequality constraints, we can write down primal and dual problems
- The dual solution is always a lower bound on the primal solution (weak duality)
- The duality gap equals 0 under certain conditions (strong duality), and in such cases we can either solve the primal or dual problem
- Strong duality holds for the SVM problem, and in particular the KKT conditions are necessary and sufficient for the optimal solution

Dual formulation of SVM and Kernel SVMs

$$\begin{aligned} \max_{\alpha} \quad & \sum_n \alpha_n - \frac{1}{2} \sum_{m,n} y_m y_n \alpha_m \alpha_n \phi(\mathbf{x}_m)^T \phi(\mathbf{x}_n) \\ \text{s.t.} \quad & 0 \leq \alpha_n \leq C, \quad \forall n \\ & \sum_n \alpha_n y_n = 0 \end{aligned}$$

- Dual problem is also a convex quadratic programming

Dual formulation of SVM and Kernel SVMs

$$\begin{aligned} \max_{\alpha} \quad & \sum_n \alpha_n - \frac{1}{2} \sum_{m,n} y_m y_n \alpha_m \alpha_n \phi(\mathbf{x}_m)^T \phi(\mathbf{x}_n) \\ \text{s.t.} \quad & 0 \leq \alpha_n \leq C, \quad \forall n \\ & \sum_n \alpha_n y_n = 0 \end{aligned}$$

- Dual problem is also a convex quadratic programming involving N dual variables α_n
- Kernel SVM:

Dual formulation of SVM and Kernel SVMs

$$\begin{aligned} \max_{\alpha} \quad & \sum_n \alpha_n - \frac{1}{2} \sum_{m,n} y_m y_n \alpha_m \alpha_n \phi(\mathbf{x}_m)^T \phi(\mathbf{x}_n) \\ \text{s.t.} \quad & 0 \leq \alpha_n \leq C, \quad \forall n \\ & \sum_n \alpha_n y_n = 0 \end{aligned}$$

- Dual problem is also a convex quadratic programming involving N dual variables α_n
- Kernel SVM: replace inner products $\phi(\mathbf{x}_m)^T \phi(\mathbf{x}_n)$ with a kernel function, $k(\mathbf{x}_m, \mathbf{x}_n)$

Recovering primal solution using dual variables

Why do we care?

Recovering primal solution using dual variables

Why do we care?

- Using solely a kernel function, we can solve the dual optimization problem and make predictions at test time!
- Prediction only depends on support vectors, i.e., points with $\alpha_n > 0$!

Recovering primal solution using dual variables

Why do we care?

- Using solely a kernel function, we can solve the dual optimization problem and make predictions at test time!
- Prediction only depends on support vectors, i.e., points with $\alpha_n > 0$!

Weights

$$\mathbf{w} = \sum_n y_n \alpha_n \phi(\mathbf{x}_n) \leftarrow \text{Linear combination of the input features}$$

Offset

$$b = [y_n - \mathbf{w}^T \phi(\mathbf{x}_n)] = [y_n - \sum_m y_m \alpha_m k(\mathbf{x}_m, \mathbf{x}_n)], \quad \text{for any } C > \alpha_n > 0$$

Prediction on a test point \mathbf{x}

$$h(\mathbf{x}) = \text{SIGN}(\mathbf{w}^T \phi(\mathbf{x}) + b) = \text{SIGN}\left(\sum_n y_n \alpha_n k(\mathbf{x}_n, \mathbf{x}) + b\right)$$

Outline

- 1 Administration
- 2 Review of last lecture
- 3 **Boosting**
 - AdaBoost
 - Derivation of AdaBoost

Boosting

High-level idea: combine a lot of classifiers

- Sequentially construct / identify these classifiers one at a time
- Use *weak* classifiers to arrive at complex decision boundaries (*strong* classifiers)

Our plan

- Describe AdaBoost algorithm
- Derive the algorithm

Adaboost Algorithm

- Given: N samples $\{\mathbf{x}_n, y_n\}$, where $y_n \in \{+1, -1\}$, and some way of constructing weak (or base) classifiers

Adaboost Algorithm

- Given: N samples $\{\mathbf{x}_n, y_n\}$, where $y_n \in \{+1, -1\}$, and some way of constructing weak (or base) classifiers
- Initialize weights $w_1(n) = \frac{1}{N}$ for every training sample

Adaboost Algorithm

- Given: N samples $\{\mathbf{x}_n, y_n\}$, where $y_n \in \{+1, -1\}$, and some way of constructing weak (or base) classifiers
- Initialize weights $w_1(n) = \frac{1}{N}$ for every training sample
- For $t = 1$ to T
 - ① Train a weak classifier $h_t(\mathbf{x})$ using current weights $w_t(n)$, by minimizing the weighted classification error

$$\epsilon_t = \sum_n w_t(n) \mathbb{I}[y_n \neq h_t(\mathbf{x}_n)]$$

Adaboost Algorithm

- Given: N samples $\{\mathbf{x}_n, y_n\}$, where $y_n \in \{+1, -1\}$, and some way of constructing weak (or base) classifiers
- Initialize weights $w_1(n) = \frac{1}{N}$ for every training sample
- For $t = 1$ to T
 - 1 Train a weak classifier $h_t(\mathbf{x})$ using current weights $w_t(n)$, by minimizing the weighted classification error

$$\epsilon_t = \sum_n w_t(n) \mathbb{I}[y_n \neq h_t(\mathbf{x}_n)]$$

- 2 Compute contribution for this classifier

Adaboost Algorithm

- Given: N samples $\{\mathbf{x}_n, y_n\}$, where $y_n \in \{+1, -1\}$, and some way of constructing weak (or base) classifiers
- Initialize weights $w_1(n) = \frac{1}{N}$ for every training sample
- For $t = 1$ to T
 - 1 Train a weak classifier $h_t(\mathbf{x})$ using current weights $w_t(n)$, by minimizing the weighted classification error

$$\epsilon_t = \sum_n w_t(n) \mathbb{I}[y_n \neq h_t(\mathbf{x}_n)]$$

- 2 Compute contribution for this classifier: $\beta_t = \frac{1}{2} \log \frac{1-\epsilon_t}{\epsilon_t}$

Adaboost Algorithm

- Given: N samples $\{\mathbf{x}_n, y_n\}$, where $y_n \in \{+1, -1\}$, and some way of constructing weak (or base) classifiers
- Initialize weights $w_1(n) = \frac{1}{N}$ for every training sample
- For $t = 1$ to T
 - 1 Train a weak classifier $h_t(\mathbf{x})$ using current weights $w_t(n)$, by minimizing the weighted classification error

$$\epsilon_t = \sum_n w_t(n) \mathbb{I}[y_n \neq h_t(\mathbf{x}_n)]$$

- 2 Compute contribution for this classifier: $\beta_t = \frac{1}{2} \log \frac{1-\epsilon_t}{\epsilon_t}$
- 3 Update weights on training points

Adaboost Algorithm

- Given: N samples $\{\mathbf{x}_n, y_n\}$, where $y_n \in \{+1, -1\}$, and some way of constructing weak (or base) classifiers
- Initialize weights $w_1(n) = \frac{1}{N}$ for every training sample
- For $t = 1$ to T
 - 1 Train a weak classifier $h_t(\mathbf{x})$ using current weights $w_t(n)$, by minimizing the weighted classification error

$$\epsilon_t = \sum_n w_t(n) \mathbb{I}[y_n \neq h_t(\mathbf{x}_n)]$$

- 2 Compute contribution for this classifier: $\beta_t = \frac{1}{2} \log \frac{1-\epsilon_t}{\epsilon_t}$
- 3 Update weights on training points

$$w_{t+1}(n) \propto w_t(n) e^{-\beta_t y_n h_t(\mathbf{x}_n)}$$

and normalize them such that $\sum_n w_{t+1}(n) = 1$

Adaboost Algorithm

- Given: N samples $\{\mathbf{x}_n, y_n\}$, where $y_n \in \{+1, -1\}$, and some way of constructing weak (or base) classifiers
- Initialize weights $w_1(n) = \frac{1}{N}$ for every training sample
- For $t = 1$ to T
 - 1 Train a weak classifier $h_t(\mathbf{x})$ using current weights $w_t(n)$, by minimizing the weighted classification error

$$\epsilon_t = \sum_n w_t(n) \mathbb{I}[y_n \neq h_t(\mathbf{x}_n)]$$

- 2 Compute contribution for this classifier: $\beta_t = \frac{1}{2} \log \frac{1-\epsilon_t}{\epsilon_t}$
- 3 Update weights on training points

$$w_{t+1}(n) \propto w_t(n) e^{-\beta_t y_n h_t(\mathbf{x}_n)}$$

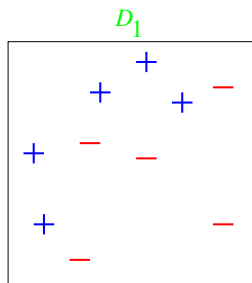
and normalize them such that $\sum_n w_{t+1}(n) = 1$

- Output the final classifier

$$h[\mathbf{x}] = \text{sign} \left[\sum_{t=1}^T \beta_t h_t(\mathbf{x}) \right]$$

Example

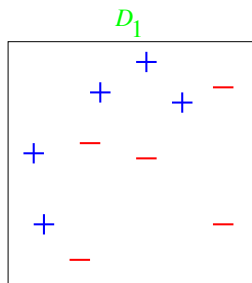
10 data points and 2 features



- The data points are clearly not linear separable
- In the beginning, all data points have equal weights (the size of the data markers “+” or “-”)

Example

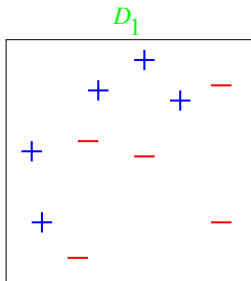
10 data points and 2 features



- The data points are clearly not linear separable
- In the beginning, all data points have equal weights (the size of the data markers “+” or “-”)
- Base classifier $h(\cdot)$: either horizontal or vertical lines

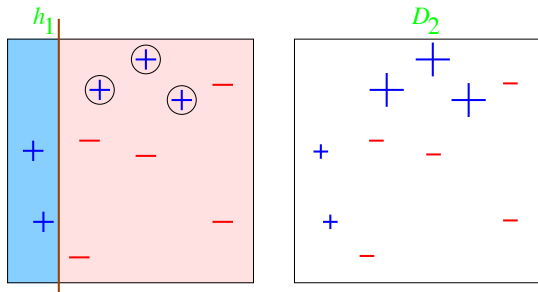
Example

10 data points and 2 features

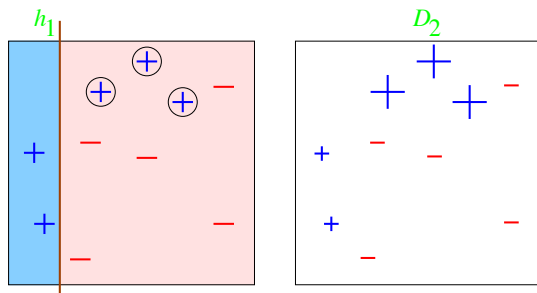


- The data points are clearly not linear separable
- In the beginning, all data points have equal weights (the size of the data markers “+” or “-”)
- Base classifier $h(\cdot)$: either horizontal or vertical lines
 - ▶ These 'decision stumps' are just trees with a single internal node, i.e., they classifying data based on a single attribute

Round 1: $t = 1$

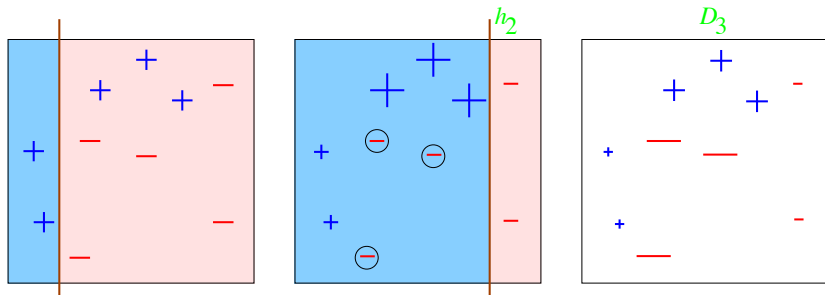


Round 1: $t = 1$

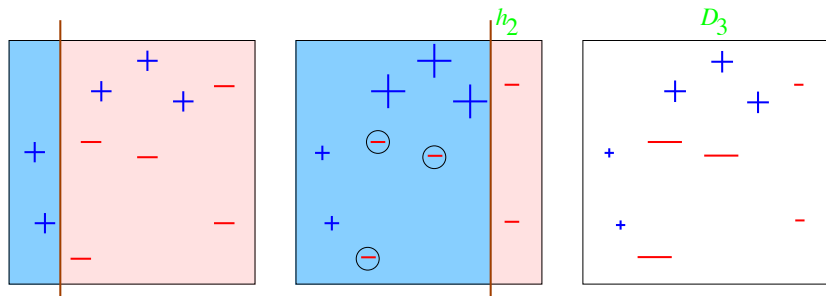


- 3 misclassified (with circles): $\epsilon_1 = 0.3 \rightarrow \beta_1 = 0.42$.
- Weights recomputed; the 3 misclassified data points receive larger weights

Round 2: $t = 2$

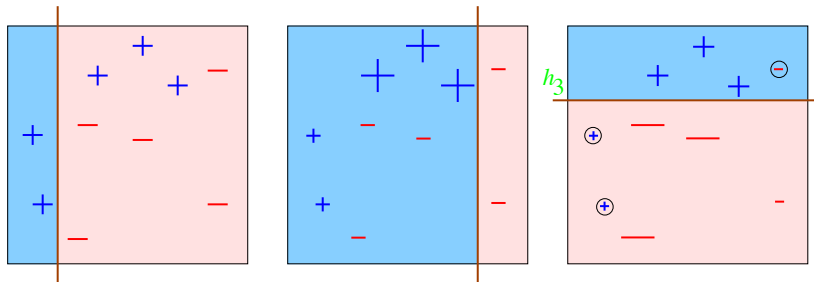


Round 2: $t = 2$

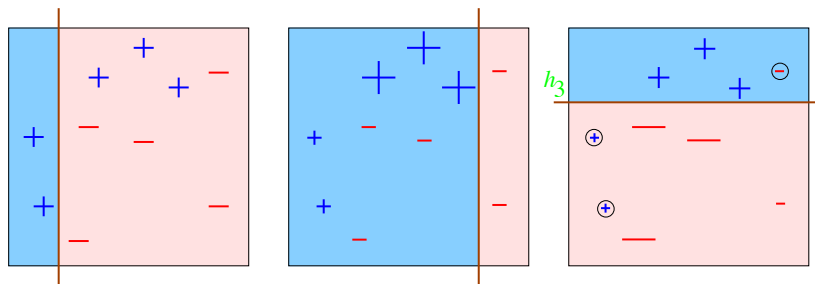


- 3 misclassified (with circles): $\epsilon_2 = 0.21 \rightarrow \beta_2 = 0.65$.
Note that $\epsilon_2 \neq 0.3$ as those 3 data points have weights less than $1/10$
- 3 misclassified data points get larger weights
- Data points classified correctly in both rounds have small weights

Round 3: $t = 3$

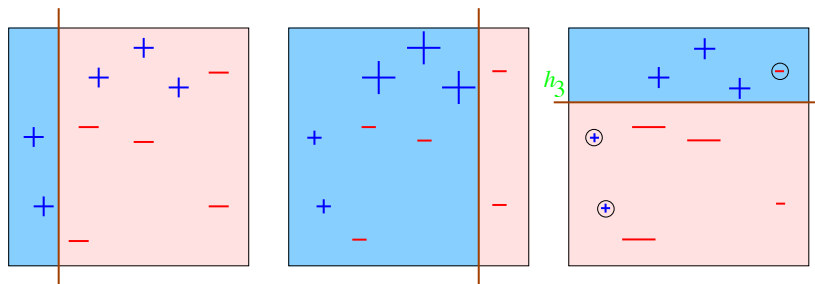


Round 3: $t = 3$



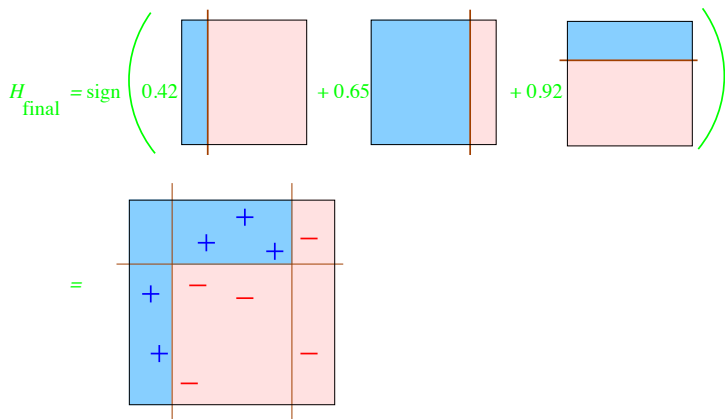
- 3 misclassified (with circles): $\epsilon_3 = 0.14 \rightarrow \beta_3 = 0.92$.
- Previously correctly classified data points are now misclassified, hence our error is low; what's the intuition?

Round 3: $t = 3$



- 3 misclassified (with circles): $\epsilon_3 = 0.14 \rightarrow \beta_3 = 0.92$.
- Previously correctly classified data points are now misclassified, hence our error is low; what's the intuition?
 - ▶ Since they have been consistently classified correctly, this round's mistake will hopefully not have a huge impact on the overall prediction

Final classifier: combining 3 classifiers



- All data points are now classified correctly!

Why AdaBoost works?

It minimizes a loss function related to classification error.

Classification loss

- Suppose we want to have a classifier

$$h(\mathbf{x}) = \text{sign}[f(\mathbf{x})] = \begin{cases} 1 & \text{if } f(\mathbf{x}) > 0 \\ -1 & \text{if } f(\mathbf{x}) < 0 \end{cases}$$

- Our loss function is thus

$$\ell(h(\mathbf{x}), y) = \begin{cases} 0 & \text{if } yf(\mathbf{x}) > 0 \\ 1 & \text{if } yf(\mathbf{x}) < 0 \end{cases}$$

Namely, the function $f(\mathbf{x})$ and the target label y should have the same sign to avoid a loss of 1.

Surrogate loss

0 – 1 loss function $\ell(h(\mathbf{x}), y)$ is non-convex and difficult to optimize.
We can instead use a surrogate loss – what are examples?

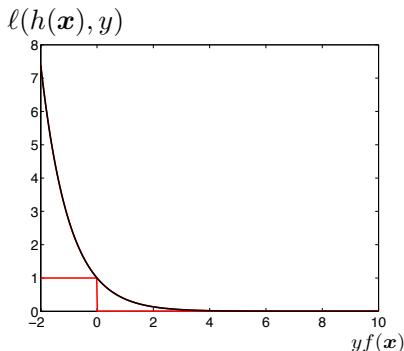
Surrogate loss

0 – 1 loss function $\ell(h(\mathbf{x}), y)$ is non-convex and difficult to optimize. We can instead use a surrogate loss – what are examples?

Exponential Loss

$$\ell^{\text{EXP}}(h(\mathbf{x}), y) = e^{-yf(\mathbf{x})}$$

$\ell^{\text{EXP}}(h(\mathbf{x}), y)$ is easier to handle numerically as it is differentiable



Choosing the t -th classifier

Suppose we have built a classifier $f_{t-1}(\mathbf{x})$, and we want to improve it by adding a weak learner $h_t(\mathbf{x})$

$$f(\mathbf{x}) = f_{t-1}(\mathbf{x}) + \beta_t h_t(\mathbf{x})$$

How can we choose optimally the new classifier $h_t(\mathbf{x})$ and the combination coefficient β_t ?

Adaboost greedily *minimizes the exponential loss function*.

$$(h_t^*(\mathbf{x}), \beta_t^*) = \arg \min_{(h_t(\mathbf{x}), \beta_t)} \sum_n e^{-y_n f(\mathbf{x}_n)}$$

Choosing the t -th classifier

Suppose we have built a classifier $f_{t-1}(\mathbf{x})$, and we want to improve it by adding a weak learner $h_t(\mathbf{x})$

$$f(\mathbf{x}) = f_{t-1}(\mathbf{x}) + \beta_t h_t(\mathbf{x})$$

How can we choose optimally the new classifier $h_t(\mathbf{x})$ and the combination coefficient β_t ?

Adaboost greedily *minimizes the exponential loss function*.

$$\begin{aligned}(h_t^*(\mathbf{x}), \beta_t^*) &= \arg \min_{(h_t(\mathbf{x}), \beta_t)} \sum_n e^{-y_n f(\mathbf{x}_n)} \\ &= \arg \min_{(h_t(\mathbf{x}), \beta_t)} \sum_n e^{-y_n [f_{t-1}(\mathbf{x}_n) + \beta_t h_t(\mathbf{x}_n)]}\end{aligned}$$

Choosing the t -th classifier

Suppose we have built a classifier $f_{t-1}(\mathbf{x})$, and we want to improve it by adding a weak learner $h_t(\mathbf{x})$

$$f(\mathbf{x}) = f_{t-1}(\mathbf{x}) + \beta_t h_t(\mathbf{x})$$

How can we choose optimally the new classifier $h_t(\mathbf{x})$ and the combination coefficient β_t ?

Adaboost greedily *minimizes the exponential loss function*.

$$\begin{aligned}(h_t^*(\mathbf{x}), \beta_t^*) &= \arg \min_{(h_t(\mathbf{x}), \beta_t)} \sum_n e^{-y_n f(\mathbf{x}_n)} \\ &= \arg \min_{(h_t(\mathbf{x}), \beta_t)} \sum_n e^{-y_n [f_{t-1}(\mathbf{x}_n) + \beta_t h_t(\mathbf{x}_n)]} \\ &= \arg \min_{(h_t(\mathbf{x}), \beta_t)} \sum_n w_t(n) e^{-y_n \beta_t h_t(\mathbf{x}_n)}\end{aligned}$$

where we have used $w_t(n)$ as a shorthand for $e^{-y_n f_{t-1}(\mathbf{x}_n)}$

The new classifier

We can decompose the *weighted* loss function into two parts

$$\begin{aligned} & \sum_n w_t(n) e^{-y_n \beta_t h_t(\mathbf{x}_n)} \\ &= \sum_n w_t(n) e^{\beta_t} \mathbb{I}[y_n \neq h_t(\mathbf{x}_n)] + \sum_n w_t(n) e^{-\beta_t} \mathbb{I}[y_n = h_t(\mathbf{x}_n)] \end{aligned}$$

The new classifier

We can decompose the *weighted* loss function into two parts

$$\begin{aligned} & \sum_n w_t(n) e^{-y_n \beta_t h_t(\mathbf{x}_n)} \\ &= \sum_n w_t(n) e^{\beta_t} \mathbb{I}[y_n \neq h_t(\mathbf{x}_n)] + \sum_n w_t(n) e^{-\beta_t} \mathbb{I}[y_n = h_t(\mathbf{x}_n)] \\ &= \sum_n w_t(n) e^{\beta_t} \mathbb{I}[y_n \neq h_t(\mathbf{x}_n)] + \sum_n w_t(n) e^{-\beta_t} (1 - \mathbb{I}[y_n \neq h_t(\mathbf{x}_n)]) \end{aligned}$$

The new classifier

We can decompose the *weighted* loss function into two parts

$$\begin{aligned} & \sum_n w_t(n) e^{-y_n \beta_t h_t(\mathbf{x}_n)} \\ &= \sum_n w_t(n) e^{\beta_t} \mathbb{I}[y_n \neq h_t(\mathbf{x}_n)] + \sum_n w_t(n) e^{-\beta_t} \mathbb{I}[y_n = h_t(\mathbf{x}_n)] \\ &= \sum_n w_t(n) e^{\beta_t} \mathbb{I}[y_n \neq h_t(\mathbf{x}_n)] + \sum_n w_t(n) e^{-\beta_t} (1 - \mathbb{I}[y_n \neq h_t(\mathbf{x}_n)]) \\ &= (e^{\beta_t} - e^{-\beta_t}) \sum_n w_t(n) \mathbb{I}[y_n \neq h_t(\mathbf{x}_n)] + e^{-\beta_t} \sum_n w_t(n) \end{aligned}$$

We have used the following properties to derive the above

- $y_n h_t(\mathbf{x}_n)$ is either 1 or -1 as $h_t(\mathbf{x}_n)$ is the output of a binary classifier
- The indicator function $\mathbb{I}[y_n = h_t(\mathbf{x}_n)]$ is either 0 or 1, so it equals $1 - \mathbb{I}[y_n \neq h_t(\mathbf{x}_n)]$

Finding the optimal weak learner

Summary

$$\begin{aligned}(h_t^*(\mathbf{x}), \beta_t^*) &= \arg \min_{(h_t(\mathbf{x}), \beta_t)} \sum_n w_t(n) e^{-y_n \beta_t h_t(\mathbf{x}_n)} \\ &= \arg \min_{(h_t(\mathbf{x}), \beta_t)} (e^{\beta_t} - e^{-\beta_t}) \sum_n w_t(n) \mathbb{I}[y_n \neq h_t(\mathbf{x}_n)] \\ &\quad + e^{-\beta_t} \sum_n w_t(n)\end{aligned}$$

What term(s) must we optimize to choose $h_t(\mathbf{x}_n)$?

Finding the optimal weak learner

Summary

$$\begin{aligned}(h_t^*(\mathbf{x}), \beta_t^*) &= \arg \min_{(h_t(\mathbf{x}), \beta_t)} \sum_n w_t(n) e^{-y_n \beta_t h_t(\mathbf{x}_n)} \\ &= \arg \min_{(h_t(\mathbf{x}), \beta_t)} (e^{\beta_t} - e^{-\beta_t}) \sum_n w_t(n) \mathbb{I}[y_n \neq h_t(\mathbf{x}_n)] \\ &\quad + e^{-\beta_t} \sum_n w_t(n)\end{aligned}$$

What term(s) must we optimize to choose $h_t(\mathbf{x}_n)$?

$$h_t^*(\mathbf{x}) = \arg \min_{h_t(\mathbf{x})} \epsilon_t = \sum_n w_t(n) \mathbb{I}[y_n \neq h_t(\mathbf{x}_n)]$$

Finding the optimal weak learner

Summary

$$\begin{aligned}(h_t^*(\mathbf{x}), \beta_t^*) &= \arg \min_{(h_t(\mathbf{x}), \beta_t)} \sum_n w_t(n) e^{-y_n \beta_t h_t(\mathbf{x}_n)} \\ &= \arg \min_{(h_t(\mathbf{x}), \beta_t)} (e^{\beta_t} - e^{-\beta_t}) \sum_n w_t(n) \mathbb{I}[y_n \neq h_t(\mathbf{x}_n)] \\ &\quad + e^{-\beta_t} \sum_n w_t(n)\end{aligned}$$

What term(s) must we optimize to choose $h_t(\mathbf{x}_n)$?

$$h_t^*(\mathbf{x}) = \arg \min_{h_t(\mathbf{x})} \epsilon_t = \sum_n w_t(n) \mathbb{I}[y_n \neq h_t(\mathbf{x}_n)]$$

Minimize weighted classification error as noted in step 1 of Adaboost!

How to choose β_t ?

Summary

$$\begin{aligned}(h_t^*(\mathbf{x}), \beta_t^*) &= \arg \min_{(h_t(\mathbf{x}), \beta_t)} \sum_n w_t(n) e^{-y_n \beta_t h_t(\mathbf{x}_n)} \\ &= \arg \min_{(h_t(\mathbf{x}), \beta_t)} (e^{\beta_t} - e^{-\beta_t}) \sum_n w_t(n) \mathbb{I}[y_n \neq h_t(\mathbf{x}_n)] \\ &\quad + e^{-\beta_t} \sum_n w_t(n)\end{aligned}$$

What term(s) must we optimize?

How to choose β_t ?

Summary

$$\begin{aligned}(h_t^*(\mathbf{x}), \beta_t^*) &= \arg \min_{(h_t(\mathbf{x}), \beta_t)} \sum_n w_t(n) e^{-y_n \beta_t h_t(\mathbf{x}_n)} \\ &= \arg \min_{(h_t(\mathbf{x}), \beta_t)} (e^{\beta_t} - e^{-\beta_t}) \sum_n w_t(n) \mathbb{I}[y_n \neq h_t(\mathbf{x}_n)] \\ &\quad + e^{-\beta_t} \sum_n w_t(n)\end{aligned}$$

What term(s) must we optimize?

We need to minimize the entire objective function with respect to β_t !

How to choose β_t ?

Summary

$$\begin{aligned}(h_t^*(\mathbf{x}), \beta_t^*) &= \arg \min_{(h_t(\mathbf{x}), \beta_t)} \sum_n w_t(n) e^{-y_n \beta_t h_t(\mathbf{x}_n)} \\ &= \arg \min_{(h_t(\mathbf{x}), \beta_t)} (e^{\beta_t} - e^{-\beta_t}) \sum_n w_t(n) \mathbb{I}[y_n \neq h_t(\mathbf{x}_n)] \\ &\quad + e^{-\beta_t} \sum_n w_t(n)\end{aligned}$$

What term(s) must we optimize?

We need to minimize the entire objective function with respect to β_t !

We can do this by taking derivative with respect to β_t , setting to zero, and solving for β_t . After some calculation and using $\sum_n w_t(n) = 1$, we find:

$$\beta_t^* = \frac{1}{2} \log \frac{1 - \epsilon_t}{\epsilon_t}$$

which is precisely step 2 of Adaboost! (*Exercise – verify the solution*)

Updating the weights

Once we find the optimal weak learner we can update our classifier:

$$f(\mathbf{x}) = f_{t-1}(\mathbf{x}) + \beta_t^* h_t^*(\mathbf{x})$$

Updating the weights

Once we find the optimal weak learner we can update our classifier:

$$f(\mathbf{x}) = f_{t-1}(\mathbf{x}) + \beta_t^* h_t^*(\mathbf{x})$$

We then need to compute the weights for the above classifier as:

$$w_{t+1}(n) = e^{-y_n f(\mathbf{x}_n)} = e^{-y_n [f_{t-1}(\mathbf{x}_n) + \beta_t^* h_t^*(\mathbf{x}_n)]}$$

Updating the weights

Once we find the optimal weak learner we can update our classifier:

$$f(\mathbf{x}) = f_{t-1}(\mathbf{x}) + \beta_t^* h_t^*(\mathbf{x})$$

We then need to compute the weights for the above classifier as:

$$\begin{aligned} w_{t+1}(n) &= e^{-y_n f(\mathbf{x}_n)} = e^{-y_n [f_{t-1}(\mathbf{x}) + \beta_t^* h_t^*(\mathbf{x}_n)]} \\ &= w_t(n) e^{-y_n \beta_t^* h_t^*(\mathbf{x}_n)} \end{aligned}$$

Updating the weights

Once we find the optimal weak learner we can update our classifier:

$$f(\mathbf{x}) = f_{t-1}(\mathbf{x}) + \beta_t^* h_t^*(\mathbf{x})$$

We then need to compute the weights for the above classifier as:

$$\begin{aligned} w_{t+1}(n) &= e^{-y_n f(\mathbf{x}_n)} = e^{-y_n [f_{t-1}(\mathbf{x}_n) + \beta_t^* h_t^*(\mathbf{x}_n)]} \\ &= w_t(n) e^{-y_n \beta_t^* h_t^*(\mathbf{x}_n)} = \begin{cases} w_t(n) e^{\beta_t^*} & \text{if } y_n \neq h_t^*(\mathbf{x}_n) \\ w_t(n) e^{-\beta_t^*} & \text{if } y_n = h_t^*(\mathbf{x}_n) \end{cases} \end{aligned}$$

Intuition Misclassified data points will get their weights increased, while correctly classified data points will get their weight decreased

Meta-Algorithm

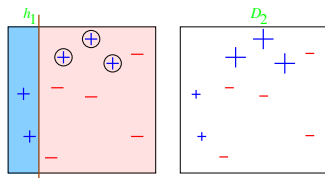
Note that the AdaBoost algorithm itself never specifies how we would get $h_t^*(\mathbf{x})$ as long as it minimizes the weighted classification error

$$\epsilon_t = \sum_n w_t(n) \mathbb{I}[y_n \neq h_t^*(\mathbf{x}_n)]$$

In this aspect, the AdaBoost algorithm is a meta-algorithm and can be used with any type of classifier

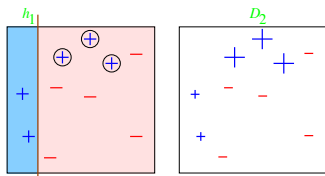
E.g., Decision Stumps

How do we choose the decision stump classifier given the weights at the second round of the following distribution?



E.g., Decision Stumps

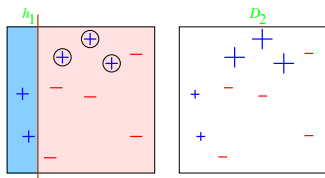
How do we choose the decision stump classifier given the weights at the second round of the following distribution?



We can simply enumerate all possible ways of putting vertical and horizontal lines to separate the data points into two classes and find the one with the smallest weighted classification error! Runtime?

E.g., Decision Stumps

How do we choose the decision stump classifier given the weights at the second round of the following distribution?



We can simply enumerate all possible ways of putting vertical and horizontal lines to separate the data points into two classes and find the one with the smallest weighted classification error! Runtime?

- Presort data by each feature in $O(dN \log N)$ time
- Evaluate $N + 1$ thresholds for each feature at each round in $O(dN)$ time
- In total $O(dN \log N + dNT)$ time – this efficiency is an attractive quality of boosting!

Interpreting boosting as learning nonlinear basis

Two-stage process

- Get $\text{SIGN}[h_1(\mathbf{x})], \text{SIGN}[h_2(\mathbf{x})], \dots,$
- Combine into a linear classification model

$$y = \text{SIGN} \left\{ \sum_t \beta_t \text{SIGN}[h_t(\mathbf{x})] \right\} = \boldsymbol{\beta}^\top \boldsymbol{\phi}(\mathbf{x})$$

Interpreting boosting as learning nonlinear basis

Two-stage process

- Get $\text{SIGN}[h_1(\mathbf{x})], \text{SIGN}[h_2(\mathbf{x})], \dots$,
- Combine into a linear classification model

$$y = \text{SIGN} \left\{ \sum_t \beta_t \text{SIGN}[h_t(\mathbf{x})] \right\} = \boldsymbol{\beta}^\top \boldsymbol{\phi}(\mathbf{x})$$

In other words, each stage learns a nonlinear basis $\phi_t(\mathbf{x}) = \text{SIGN}[h_t(\mathbf{x})]$

- This is an alternative way to introduce non-linearity aside from kernel methods
- We could also try to learn the basis functions and the classifier at the same time, as we'll talk about with neural networks next class