

Adaptive basis function models

Jiali Lin

Virginia Tech

January 15, 2017

Outline

Introduction

Classification and regression trees (CART)

Boosting

Ensemble learning

Introduction

- ▶ Consider kernel methods to create non-linear models for regression and classification.
- ▶ The prediction takes the form $f(\mathbf{x}) = \mathbf{w}^T \phi(\mathbf{x})$, where $\phi(\mathbf{x})$ is kernel.
- ▶ Coming up with a good kernel function is quite difficult.
- ▶ **Goal:** dispense with kernels altogether, and try to learn useful features $\phi(\mathbf{x})$ directly from the input data.
- ▶ **Adaptive basis function model (ABM)** has the form

$$f(\mathbf{x}) = w_0 + \sum_{m=1}^M w_m \phi_m(\mathbf{x})$$

Outline

Introduction

Classification and regression trees (CART)

Boosting

Ensemble learning

Basics

Consider the tree in Figure.

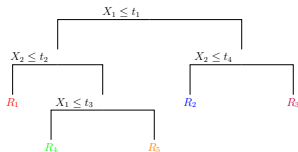


Figure: The first node asks if x_1 is less than some threshold t_1 . If yes, we then ask if x_2 is less than some other threshold t_2 . If yes, we are in the bottom left quadrant of space, R_1 . If no, we ask if x_1 is less than t_3 . And so on. The result of these axis parallel splits is to partition 2d space into 5 regions.

We can write the model in the following form

$$f(\mathbf{x}) = E[y|\mathbf{x}] = \sum_{m=1}^M w_m \mathbb{I}(\mathbf{x} \in R_m) = \sum_{m=1}^M w_m \phi(\mathbf{x}; \mathbf{v}_m)$$

where R_m is the m 'th region, w_m is the mean response in this region, and \mathbf{v}_m encodes the choice of variable to split on.

Growing a tree

- ▶ Finding the optimal partitioning of the data is NP-complete.
- ▶ Methods to compute a locally optimal MLE: **CART**, **C4.5** and **ID3**.
- ▶ The split function chooses the best feature, and the best value for that feature

$$(j^*, t^*) = \min_{j \in \{1, \dots, D\}} \min_{t \in \mathcal{T}_j} \text{cost}(\mathbf{x}_i, y_i : x_{ij} \leq t) + \text{cost}(\mathbf{x}_i, y_i : x_{ij} > t)$$

where \mathcal{T}_j is set of possible thresholds for feature j .

- ▶ Split function checks if a node is worth splitting can use several stopping heuristics.
- ▶ Cost measure depends on whether our goal is regression or classification.

Cost

► **Regression cost**

$$\text{cost}(\mathcal{D}) = \sum_{i \in \mathcal{D}} (y_i - \bar{y})^2$$

where $\bar{y} = \frac{1}{|\mathcal{D}|} \sum_{i \in \mathcal{D}} y_i$.

- **Classification cost:** First, we fit a multinoulli model to the data in the leaf satisfying the test $X_j < t$ by estimating

$$\hat{\pi}_c = \frac{1}{|D|} \sum_{i \in D} \mathbb{I}(y_i = c)$$

Several ways to measure the quality of a split:

- Misclassification rate.
- Entropy, or deviance.
- Information gain.
- Gini index.

Example

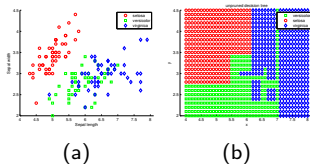


Figure: (a) Iris data. We only show the first two features, sepal length and sepal width, and ignore petal length and petal width. (b) Decision boundaries induced by the decision tree in Figure 2(a).

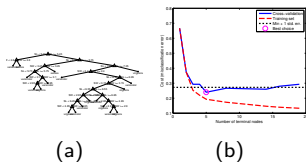


Figure: (a) Unpruned decision tree for Iris data. (b) Plot of misclassification error rate vs depth of tree. Figure generated by dtreeDemoIris.

Pruning a tree

- ▶ Motivation: prevent overfitting.
- ▶ Idea: grow a “full” tree, and then to perform **pruning**.
- ▶ How far to prune back? Evaluate the cross-validated error on each subtree.
- ▶ Biggest **cons**: unstable! outlier can leads to high variance trees.

Random forests

- ▶ Reduce the variance of an estimate: average together many estimates (say M trees)

$$f(\mathbf{x}) = \sum_{m=1}^M \frac{1}{M} f_m(\mathbf{x})$$

where f_m is the m 'th tree (This is called **bagging**).

- ▶ Problems: can result in highly correlated predictors.
- ▶ **Random forests** decorrelates the base learners by learning trees based on a randomly chosen subset of input variables and datasets.

Outline

Introduction

Classification and regression trees (CART)

Boosting

Ensemble learning

Introduction

- ▶ **Boosting:** a greedy algorithm for fitting ABMs, where the ϕ_m are generated by an algorithm called a **weak learner** or a **base learner**.
- ▶ Idea: applying the weak learner sequentially to weighted versions of the data.
- ▶ More weight is given to examples that were misclassified by earlier rounds.
- ▶ **Goal:** solve the optimization problem

$$\min_f \sum_{i=1}^N L(y_i, f(\mathbf{x}_i))$$

Name	Loss	Derivative	Minimizers f^*	Algorithm
Squared error	$\frac{1}{2}(y_i - f(\mathbf{x}_i))^2$	$y_i - f(\mathbf{x}_i)$	$E[Y \mathbf{x}_i]$	L2Boosting
Absolute error	$ y_i - f(\mathbf{x}_i) $	$\text{sign}(y_i - f(\mathbf{x}_i))$	$\text{median}(y \mathbf{x}_i)$	Gradient boosting
Exponential loss	$\exp(-\tilde{y}_i f(\mathbf{x}_i))$	$-\tilde{y}_i \exp(-\tilde{y}_i f(\mathbf{x}_i))$	$\frac{1}{2} \log \frac{\pi_i}{1-\pi_i}$	AdaBoost
Logloss	$\log(1 + e^{-\tilde{y}_i f(\mathbf{x}_i)})$	$y_i - \pi_i$	$\frac{1}{2} \log \frac{\pi_i}{1-\pi_i}$	LogitBoost

Table: Some commonly used loss functions. Assume $\tilde{y}_i \in \{-1, +1\}$, $y_i \in \{0, 1\}$ and $\pi_i = \text{sigm}(2f(\mathbf{x}_i))$. For regression problems, we assume $y_i \in \mathbb{R}$.

Generalized additive models

- ▶ Finding the optimal f may be hard, we shall tackle it sequentially.
- ▶ We initialise by defining

$$f_0(\mathbf{x}) = \min_{\gamma} \sum_{i=1}^N L(y_i, f(\mathbf{x}_i; \gamma))$$

- ▶ Then at iteration m , we compute

$$(\beta_m, \gamma_m) = \min_{\beta, \gamma} \sum_{i=1}^N L(y_i, f_{m-1}(\mathbf{x}_i) + \beta \phi(\mathbf{x}_i; \gamma))$$

- ▶ Then, we set

$$f_m(\mathbf{x}) = f_{m-1}(\mathbf{x}) + \beta_m \phi(\mathbf{x}; \gamma_m)$$

- ▶ **Key:** do not go back and adjust earlier parameters. This is called **forward stagewise additive modeling**.
- ▶ In practice, we do “partial updates”

$$f_m(\mathbf{x}) = f_{m-1}(\mathbf{x}) + \nu \beta_m \phi(\mathbf{x}; \gamma_m)$$

where $0 < \nu \leq 1$ is a step-size parameter (called **shrinkage**).

L2boosting

- Suppose we used squared error loss. Then at step m the loss

$$L(y_i, f_{m-1}(\mathbf{x}_i) + \beta\phi(\mathbf{x}_i; \gamma)) = (r_{im} - \phi(\mathbf{x}_i; \gamma))^2$$

where $r_{im} = y_i - f_{m-1}(\mathbf{x}_i)$ is the current residual, and we have set $\beta = 1$.

- We can find the new basis function by using the weak learner to predict r_m .

AdaBoost

- Consider a binary classification problem with exponential loss. At step m we have to minimize

$$L_m(\phi) = \sum_{i=1}^N \exp[-\tilde{y}_i(f_{m-1}(\mathbf{x}) + \beta\phi(\mathbf{x}_i))] = \sum_{i=1}^N w_{i,m} \exp(-\beta\tilde{y}_i\phi(\mathbf{x}_i))$$

where $w_{i,m} = \exp(-\tilde{y}_i(f_{m-1}(\mathbf{x})))$ is a weight applied to datacase i .

- Rewrite this objective as follows

$$\begin{aligned} L_m &= e^{-\beta} \sum_{\tilde{y}_i = \phi(\mathbf{x}_i)} w_{i,m} + e^{\beta} \sum_{\tilde{y}_i \neq \phi(\mathbf{x}_i)} w_{i,m} \\ &= (e^{\beta} - e^{-\beta}) \sum_{i=1}^N w_{i,m} \mathbb{I}(\tilde{y}_i \neq \phi(\mathbf{x}_i)) + e^{-\beta} \sum_{i=1}^N w_{i,m} \end{aligned}$$

AdaBoost (cont'd)

- Consequently the optimal function to add is

$$\phi_m = \operatorname{argmin}_{\phi} w_{i,m} \mathbb{I}(\tilde{y}_i \neq \phi(\mathbf{x}_i))$$

- After finding ϕ_m , substitute ϕ_m into L_m

$$\beta_m = \frac{1}{2} \log \frac{1 - \text{err}_m}{\text{err}_m}$$

where

$$\text{err}_m = \frac{\sum_{i=1}^N w_i (\tilde{y}_i \neq \phi_m(\mathbf{x}_i))}{\sum_{i=1}^N w_{i,m}}$$

- The overall updates becomes

$$f_m(\mathbf{x}) = f_{m-1}(\mathbf{x}) + \beta_m \phi_m(\mathbf{x})$$

AdaBoost (cont'd)

The weights at the next iteration become

$$\begin{aligned}w_{i,m+1} &= w_{i,m} e^{-\beta_m \tilde{y}_i \phi_m(\mathbf{x}_i)} \\&= w_{i,m} e^{\beta_m (2\mathbb{I}(\tilde{y}_i \neq \phi_m(\mathbf{x}_i)) - 1)} \\&= w_{i,m} e^{2\beta_m \mathbb{I}(\tilde{y}_i \neq \phi_m(\mathbf{x}_i))} e^{-\beta_m}\end{aligned}$$

where we exploited that $-\tilde{y}_i \phi_m(\mathbf{x}_i) = -1$ if $\tilde{y}_i = \phi_m(\mathbf{x}_i)$ and $-\tilde{y}_i \phi_m(\mathbf{x}_i) = +1$ otherwise. Since $e^{-\beta_m}$ will cancel out in the normalization step, we can drop it.

given $w_i = 1/N$.

for $m = 1 : M$ **do**

1. Fit a classifier $\phi_m(\mathbf{x})$ to the training set using weights w .
2. Compute $\text{err}_m = \frac{\sum_{i=1}^N w_i \mathbb{I}(\tilde{y}_i \neq \phi_m(\mathbf{x}_i))}{\sum_{i=1}^N w_{i,m}}$.
3. Compute $\alpha_m = \log \frac{1 - \text{err}_m}{\text{err}_m}$.
4. Set $w_i \leftarrow w_i \exp[\alpha_m \mathbb{I}(\tilde{y}_i \neq \phi_m(\mathbf{x}_i))]$.

return $f(\mathbf{x}) = \text{sgn} \left[\sum_{m=1}^M \alpha_m \phi_m(\mathbf{x}) \right]$.

Why does boosting work so well?

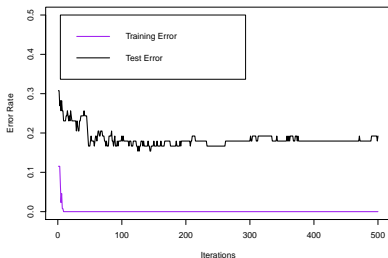


Figure: Illustration of AdaBoost. AdaBoost is very slow to overfit. Figure generated by AdaboostDemo.

- ▶ Can be viewed as “ ℓ_1 regularization”: use ℓ_1 regularization to select a subset, $\phi(x) = [\phi_1(x), \dots, \phi_K(x)]$.
- ▶ It can be proved that AdaBoost maximizes the margin on the training data.

LogitBoost

- ▶ Exponential loss puts a lot of weight on misclassified examples.
- ▶ Therefore, exponential loss is sensitive to outliers (misclassified examples).
- ▶ Cannot recover probability estimates from $f(\mathbf{x})$, $\exp(-\tilde{y}_i f(\mathbf{x}_i))$ is not the logarithm of any pmf for binary variables $\tilde{y}_i \in \{-1, +1\}$.

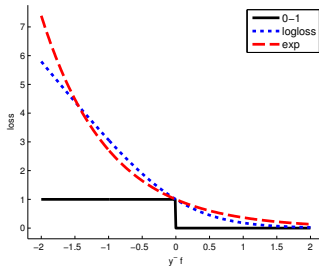


Figure: Illustration of various loss functions for binary classification. Figure generated by `hingeLossPlot`.

LogitBoost (cont'd)

- ▶ **Logloss** only punishes mistakes linearly.
- ▶ It means we will be able to extract probabilities from the final learned function, by

$$p(y = 1|\mathbf{x}) = \frac{e^{f(\mathbf{x})}}{e^{-f(\mathbf{x})} + e^{f(\mathbf{x})}} = \frac{1}{e^{-2f(\mathbf{x})} + 1}$$

- ▶ The goal is to minimize the expected log-loss, given by

$$L_m(\phi) = \sum_{i=1}^N \log [1 + \exp(-2\tilde{y}_i(f_{m-1}(\mathbf{x}) + \phi(\mathbf{x}_i)))]$$

given $w_i = 1/N, \pi_i = 1/2$.

for $m = 1 : M$ **do**

1. Compute the working response $z_i = \frac{y^* - \pi}{\pi_i(1 - \pi_i)}$.
2. Compute the weights $w_i = \pi(1 - \pi)$.
3. Compute $\phi_m = \operatorname{argmin}_{\phi} \sum_{i=1}^N w_i (z_i - \phi(\mathbf{x}_i))^2$.
4. Update $f(\mathbf{x}) \leftarrow f(\mathbf{x}) + \frac{1}{2}\phi_m(\mathbf{x})$.
5. Compute $\pi = 1/(e^{-2f(\mathbf{x}_i)} + 1)$.

return $f(\mathbf{x}) = \operatorname{sgn} \left[\sum_{m=1}^M \phi_m(\mathbf{x}) \right]$.

Boosting as functional gradient descent

- **Gradient boosting**: imagine minimizing

$$\hat{\mathbf{f}} = \min_{\mathbf{f}} L(\mathbf{f})$$

where $\mathbf{f} = (f(\mathbf{x}_1), \dots, f(\mathbf{x}_N))$ are the “parameters”.

- At step m let \mathbf{g}_m be the gradient of $L(\mathbf{f})$ evaluated at $\mathbf{f} = \mathbf{f}_{m-1}$:

$$g_{im} = \left[\frac{\partial L(y_i, f(\mathbf{x}_i))}{\partial f(\mathbf{x}_i)} \right]_{\mathbf{f}=\mathbf{f}_{m-1}}$$

- Then make the update

$$\mathbf{f}_m = \mathbf{f}_{m-1} - \rho_m \mathbf{g}_m$$

where ρ_m is the step length, chosen by

$$\rho_m = \min_{\rho} L(\mathbf{f}_{m-1} - \rho \mathbf{g}_m)$$

Boosting as functional gradient descent (cont'd)

- ▶ We only optimize f , but do not learn a function that can generalize.
- ▶ Modify the algorithm by fitting a weak learner to approximate the negative gradient signal

$$\gamma_m = \min_{\gamma} \sum_{i=1}^N (-g_{im} - \phi(\mathbf{x}_i, \gamma))^2$$

Initialize $f_0(\mathbf{x}) = \min_{\gamma} \sum_{i=1}^N L(y_i, f(\mathbf{x}_i))$.

for $m = 1 : M$ **do**

1. Compute the gradient residual using $r_{im} = -[\frac{\partial L(y_i, f(\mathbf{x}_i))}{\partial f(\mathbf{x}_i)}]_{f(\mathbf{x}_i)=f_{m-1}(\mathbf{x}_i)}$.
2. Compute the weights $w_i = \pi(1 - \pi)$.
3. Use the weak learner to compute γ_m which minimizes $\sum_{i=1}^N (r_{im} - \phi(\mathbf{x}_i, \gamma))^2$.
4. Update $f_m(\mathbf{x}) = f_{m-1}(\mathbf{x}) + \nu \phi(\mathbf{x}; \gamma_m)$.

return $f(\mathbf{x}) = f_M(\mathbf{x})$.

Outline

Introduction

Classification and regression trees (CART)

Boosting

Ensemble learning

Stacking

- ▶ **Ensemble learning:** learning a weighted combination of base models

$$f(y|\mathbf{x}, \boldsymbol{\pi}) = \sum_{m \in M} w_m f_m(y|\mathbf{x})$$

- ▶ Estimate the weights

$$\hat{\mathbf{w}} = \min_{\mathbf{w}} \sum_{i=1}^N L(y_i, \sum_{m=1}^M w_m f_m(\mathbf{x}))$$

- ▶ However, this will result in overfitting, with w_m being large for the most complex model.
- ▶ **Solution:** use the LOOCV estimate

$$\hat{\mathbf{w}} = \min_{\mathbf{w}} \sum_{i=1}^N L(y_i, \sum_{m=1}^M w_m f_m^{-1}(\mathbf{x})).$$

where $f_m^{-1}(\mathbf{x})$ is the predictor obtained by training on data excluding (\mathbf{x}_i, y_i)

- ▶ This is known as **stacking**, which stands for “stacked generalization”.