

Lecture 7

1. Generalized additive models (GAM)
2. Tree-based methods
3. Boosting
4. Boosting trees

1 Generalized additive models (GAM)

In the previous lectures we talked about the generalized linear models (GLM), where we assumed that target variable y has a distribution from exponential family and the prediction is an expected value μ of this distribution. The link function g^{-1} has been introduced such that

$$g^{-1}(\mu) = f(\mu) = \theta^T x$$

(for convenience, we denoted $f = g^{-1}$). Recall that:

- if $f(\mu) = \mu$ and $y \sim N(\mu, \sigma^2)$, then the resulted model is a linear regression;
- if $f(\mu) = \ln\left(\frac{\mu}{1-\mu}\right)$ (or, $\mu = \frac{1}{1+e^{-\eta}}$) and $y \sim \text{Ber}(\mu)$, then the resulted model is a logistic regression.

The generalized additive models can be introduced as a modification of GLM, we assume that

$$f(\mu) = \alpha + f_1(x_1) + f_2(x_2) + \dots + f_n(x_n),$$

where $f_j, j = 1, \dots, n$, are nonlinear differentiable functions, α is a constant (notice that this is nonparametric model). Consider two cases: GAM for regression and GAM for classification.

1.1 GAM for regression

We assume that $f(\mu) = \mu$ and $y \sim N(\mu, \sigma^2)$, in other words,

$$y|x; \alpha; f_j \sim N(\alpha + f_1(x_1) + f_2(x_2) + \dots + f_n(x_n), \sigma^2).$$

We find the best value for α using maximum likelihood estimation. The log-likelihood has a form

$$l(\alpha, f_1, \dots, f_n) = \ln \prod_{i=1}^m \frac{1}{\sqrt{2\pi}\sigma} \exp \left(-\frac{(y^{(i)} - \alpha - \sum_{j=1}^n f_j(x_j^{(i)}))^2}{2\sigma^2} \right).$$

The derivative with respect to α is

$$l'_\alpha = \frac{1}{\sigma^2} \sum_{i=1}^m \left(y^{(i)} - \alpha - \sum_{j=1}^n f_j(x_j^{(i)}) \right) = 0.$$

Because of the term $\sum_{j=1}^n f_j(x_j^{(i)})$, α cannot be found explicitly. Assuming the additional condition

$$\sum_{i=1}^m f_j(x_j^{(i)}) = 0 \text{ for any } j$$

(mean of f_j along the data is zero), we will find

$$\alpha = \frac{1}{m} \sum_{i=1}^m y^{(i)}$$

(average of all targets in the dataset).

The next step will be to find the estimations for f_j . If we have already found the estimations for $f_j, j \neq 1$ we will find the best estimation for f_1 (general case can be considered by analogy):

$$l(\alpha, f_j) = \ln \prod_{i=1}^m \frac{1}{\sqrt{2\pi}\sigma} \exp \left(-\frac{(y^{(i)} - \alpha - f_1(x_1^{(i)}) - \sum_{j=2}^n f_j(x_j^{(i)}))^2}{2\sigma^2} \right)$$

Denote

$$z^{(i)} = y^{(i)} - \alpha - \sum_{j=2}^n f_j(x_j^{(i)}),$$

then the log-likelihood equals to

$$l(\alpha, f_j) = \ln \prod_{i=1}^m \frac{1}{\sqrt{2\pi}\sigma} \exp \left(-\frac{(z^{(i)} - f_1(x_1^{(i)}))^2}{2\sigma^2} \right).$$

We reformulate the problem: find the best estimation for f_1 with condition

$$z|x; f_1 \sim N(f_1(x_1), \sigma^2).$$

The advantage of this formulation is that the model with only one feature should be trained. It gives a rise to the **backfitting algorithm**.

Algorithm 1 Backfitting algorithm for GAM regression

```

1: set initial values  $\alpha = \frac{1}{m} \sum_{i=1}^m y^{(i)}$ ,  $f_j = 0$  for all  $j = 1, \dots, n$ 
2: repeat
3:   for  $j = 1$  to  $n$  do
4:     evaluate working targets  $z^{(i)} = y^{(i)} - \alpha - \sum_{k=1, k \neq j}^n f_k(x_k^{(i)})$ 
5:     train model with feature  $x_j$  and target  $z$  to estimate  $f_j$ 
6: until convergence
7: return  $\alpha, f_j$ 

```

For the line 5 of this algorithm we should choose a single variable model. Simple nonparametric approaches could be used (for example, weighted linear regression or cubic splines).

1.2 GAM for classification

For the classification we assume $f(\mu) = \ln \left(\frac{\mu}{1 - \mu} \right)$ and Bernoulli distribution for the target y , in other words:

$$y|x; \alpha; f_j \sim \text{Ber}(g(\alpha + f_1(x_1) + f_2(x_2) + \dots + f_n(x_n))),$$

1.2 GAM for classification GENERALIZED ADDITIVE MODELS (GAM)

where $g(z) = \frac{1}{1 + e^{-z}}$ is a sigmoid function. Denote

$$\eta = \alpha + f_1(x_1) + f_2(x_2) + \dots + f_n(x_n),$$

then the log-likelihood can be written as

$$l(\eta) = \sum_{i=1}^m y^{(i)} \ln \mu^{(i)} + (1 - y^{(i)}) \ln(1 - \mu^{(i)}).$$

To find the initial value for α we assume that $f_j = 0$ for all $j = 1, \dots, n$:

$$l'_\eta = \sum_{i=1}^m y^{(i)}(1 - \mu^{(i)}) - (1 - y^{(i)})\mu^{(i)} = 0 \Rightarrow \alpha = \ln \left(\frac{\mu}{1 - \mu} \right),$$

where

$$\mu = \frac{1}{m} \sum_{i=1}^m y^{(i)}.$$

To find f_j we will use the following procedure. First, we use Newton's method with respect to η :

$$\eta^{new} = \eta^{old} - \frac{l'_\eta(\eta^{old})}{l''_\eta(\eta^{old})} \text{ (for each training example).}$$

But

$$\begin{aligned} l'(\eta) &= y(1 - \mu) - (1 - y)\mu = y - \mu, \\ l''(\eta) &= -\mu(1 - \mu) \end{aligned}$$

(remember that $\mu = g(\eta)$). Then the updating formula for η is

$$\eta^{new} := \eta^{old} + \frac{y - \mu^{old}}{\mu^{old}(1 - \mu^{old})}.$$

The second step will be to use backfitting algorithm for GAM regression (previous section) with target values η^{new} .

The important difference between linear case and logistic case is that in the linear case we assumed that variance σ^2 for all training examples is constant. In the logistic case the variance for the Bernoulli distribution is calculated as $\mu(1 - \mu)$. It means that we should normalize variances for the training examples to the same value; it can be easily done by using weights $\mu(1 - \mu)$ in backfitting algorithm (**weighted backfitting algorithm**):

Algorithm 2 GAM for classification with weighted backfitting algorithm

```

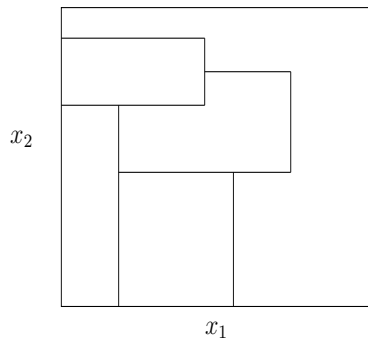
1: set initial values  $\mu = \frac{1}{m} \sum_{i=1}^m y^{(i)}$ ,  $\alpha = \ln \left( \frac{\mu}{1-\mu} \right)$ ,  $f_j = 0$ ,  $j = 1, \dots, n$ 
2: repeat
3:   evaluate  $\eta^{(i)} = \alpha + \sum_{j=1}^n f_j(x^{(i)})$  and  $\mu^{(i)} = \frac{1}{1 + e^{-\eta^{(i)}}}$ 
4:   make the Newton's step  $\eta^{(i)} := \eta^{(i)} + \frac{y^{(i)} - \mu^{(i)}}{\mu^{(i)}(1 - \mu^{(i)})}$ 
5:   evaluate weights  $w^{(i)} = \mu^{(i)}(1 - \mu^{(i)})$ 
6:   evaluate new value for  $\alpha = \frac{1}{m} \sum_{i=1}^m \eta^{(i)}$ 
7:   for  $j = 1$  to  $n$  do
8:     evaluate working targets  $z^{(i)} = \eta^{(i)} - \alpha - \sum_{k=1, k \neq j}^n f_k(x_k^{(i)})$ 
9:     train model with feature  $x_j$ , target  $z$  and weights  $w$  to estimate  $f_j$ 
10: until convergence
11: return  $\alpha$ ,  $f_j$ 

```

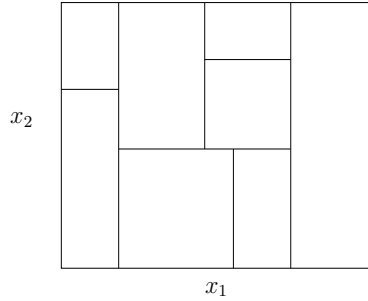
As before for the step 9 of the algorithm we can choose some nonparametric model, for example, cubic spline fitting or weighted linear regression.

2 Tree-based methods

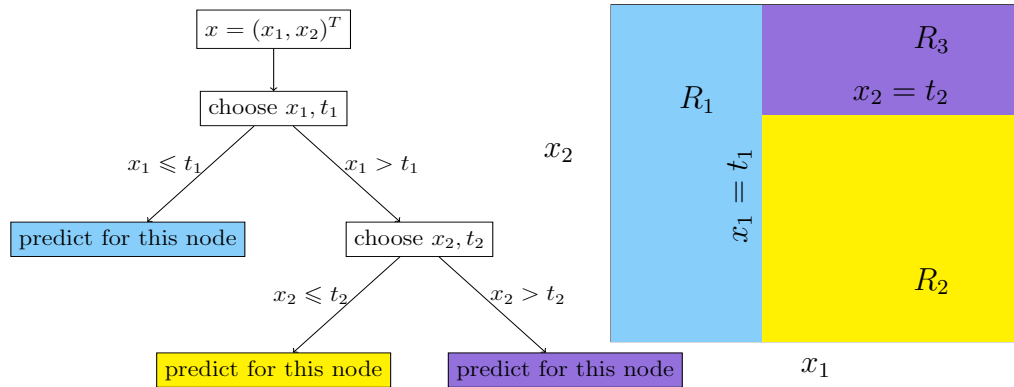
Another nonparametric approach utilizes the structure called decision tree. In many cases we should split the feature space in small regions and predict the target variable y for each region separately. The example of splitting is shown on the following picture:



Unfortunately, it is difficult to describe such splitting analytically. We consider simpler type of partition - binary partitions:



We can describe the last splitting using **decision trees**. The methods that utilize the decision trees are called tree-based methods. Consider the problem with two features $x = (x_1, x_2)^T$, the simplest decision tree with according binary partition could look like this:



Usually, the prediction is constant for each **terminal node** or **leaves** (coloured nodes on the picture). The advantage of the decision tree structure is that the hypothesis function $h_{c,t}(x)$ can be expressed as a linear combination of indicator functions. For the previous example, the hypothesis can be written as

$$h_{c,t}(x) = c_1 \cdot \mathbb{1}\{x_1 \leq t_1\} + c_2 \cdot \mathbb{1}\{x_1 > t_1, x_2 \leq t_2\} + c_3 \cdot \mathbb{1}\{x_1 > t_1, x_2 > t_2\}.$$

In this model we have 5 parameters: c_1, c_2, c_3 are predictions for the terminal nodes, t_1, t_2 are splitting points.

In general case if we have n features $x = (x_1, \dots, x_n)^T$, the hypothesis

function can be written as

$$h_{c,t}(x) = \sum_{k=1}^K c_k \mathbb{1}\{x \in R_k\}.$$

Here K is a number of terminal nodes and R_k is the region of the feature space that corresponds to the terminal node k . The main advantage of this model is interpretability.

2.1 Regression trees

We consider in details how to build the decision trees for regression problems. Assuming that the prediction c_k for each region R_k is constant it is easy to prove that the best prediction would be the average of $y^{(i)}$ for the training examples in this region.

Indeed, we could use the probabilistic approach in this case. Consider the simplest situation when we have one feature and one target. If the variable x is splitted by the number c and μ_1 is a constant prediction for $x < c$, μ_2 is a constant prediction for $x \geq c$, then

$$y|x, c, \mu_1, \mu_2 \sim N((\mu_2 - \mu_1)\mathbb{1}\{x \geq c\} + \mu_1, \sigma^2).$$

The log-likelihood

$$l(\mu_1, \mu_2, c) = m \ln \frac{1}{\sqrt{2\pi}\sigma} - \sum_{i=1}^m \frac{(y^{(i)} - (\mu_2 - \mu_1)\mathbb{1}\{x^{(i)} \geq c\} - \mu_1)^2}{2\sigma^2}.$$

Derivatives with respect to μ_1 and μ_2 give

$$\begin{aligned} \mu_1 &= \frac{\sum_{i=1}^m y^{(i)} \mathbb{1}\{x^{(i)} < c\}}{\sum_{i=1}^m \mathbb{1}\{x^{(i)} < c\}}, \\ \mu_2 &= \frac{\sum_{i=1}^m y^{(i)} \mathbb{1}\{x^{(i)} \geq c\}}{\sum_{i=1}^m \mathbb{1}\{x^{(i)} \geq c\}}. \end{aligned}$$

Differentiation with respect to c becomes more complicated, because the indicator function $\mathbb{1}$ is discontinuous. Another problem appears as in majority of

cases there are more than one feature, and we should build log-likelihood with respect to all features. It becomes computationally infeasible to optimize the error function for the hypothesis $h_{c,t}(x) = \sum_{k=1}^r c_k \mathbb{1}\{x \in R_k\}$ explicitly.

To resolve this problem we introduce the measure that will identify the best variable and the best split

$$Q_k(T) = \frac{1}{m_k} \sum_{i \in R_k} (y^{(i)} - \mu_k)^2,$$

where m_k is a number of training examples in the node R_k , and apply greedy algorithm that optimizes the tree step by step (not the hypothesis in general).

Algorithm 3 Greedy optimization algorithm for decision trees

```

1: initialize root node  $R_1 = \mathbb{R}^n$ , where  $n$  is a number of features
2: initialize a list of terminal nodes  $R = \{R_1\}$ 
3: repeat
4:   for each region  $R_k$  from  $R$  do
5:     for  $j = 1$  to  $n$  do
6:       for all splitting points  $c$  do
7:         define
             $R_{kl} = \{x^{(i)} \in R_k \mid x_j^{(i)} < c\},$ 
             $R_{kr} = \{x^{(i)} \in R_k \mid x_j^{(i)} \geq c\}$ 
8:         evaluate  $\mu_1 = \frac{1}{m_{kl}} \sum_{R_{kl}} y^{(i)}, \mu_2 = \frac{1}{m_{kr}} \sum_{R_{kr}} y^{(i)}$ 
9:         evaluate  $\varepsilon = m_{kl} Q_{kl}(T) + m_{kr} Q_{kr}(T)$ 
10:        choose  $j, c = \arg \min_{j,c} \varepsilon$ 
11:        remove  $R_k$  from  $R$ ; add  $R_{kl}, R_{kr}$  to  $R$ 
12: until convergence
13: return list of terminal nodes  $R$ 
```

This algorithm when run till the end gives a huge decision tree with small number of training examples in the terminal nodes. Obviously, that the resulted error on the training set will be equal to zero. Remember that we called such situation overfitting (excellent predictions on the training set and very bad prediction on the test set). To avoid such situation we can introduce additional stopping criteria:

- minimal size of terminal node: for example, we can require to have at least 10 training examples in each terminal node
- minimal change of error: for example, we can require not to split the node in case if the error does not decrease more than 0.01

Then the above algorithm should be modified by adding one of these conditions (or both) before the lines 10 and 11.

Usually, the decision tree is constructed using two steps. First, we build the maximal decision tree and second, we start pruning the tree by removing terminal nodes one by one. The procedure of pruning is stopped when the following function reaches the minimal value:

$$C_\alpha(T) = \sum_{k=1}^K m_k Q_k(T) + \alpha K.$$

This function creates some trade-off between the tree size and its goodness of fit to the data. If α is zero then K increases till the errors $Q_k = 0$ for all k . But if $\alpha > 0$, then big value of K implies the big value for the second term in $C_\alpha(T)$, for example, if there exists big α such that the minimal value of C_α is obtained if $K = 1$. Usually, α is defined using cross-validation procedure.

2.2 Classification trees

The main difference between regression and classification decision trees is the way to define the quality of split. If we have several classes $1, 2, \dots, D$, then it does not make sense to calculate the mean squared error $Q_k(T)$ like we did for the regression trees. Consider the regions R_1, R_2, \dots, R_K , obtained by the decision tree. We could evaluate the quantities

$$p_{kd} = \frac{1}{m_k} \sum_{i \in R_k} \mathbb{1}\{y^{(i)} = d\} \text{ (percentage of class } d \text{ in the region } R_k),$$

where as before m_k is a number of training examples in the region R_k , and $d \in \{1, 2, \dots, D\}$. There are several criteria for quality of split can be used:

- **Misclassification error:** denote $d(k) = \arg \max_d p_{kd}$ (the majority class in the region R_k), then

$$Q_k(T) = \frac{1}{m_k} \sum_{i \in R_k} \mathbb{1}\{y^{(i)} \neq d(k)\} = 1 - p_{k d(k)}.$$

- **Gini index:**

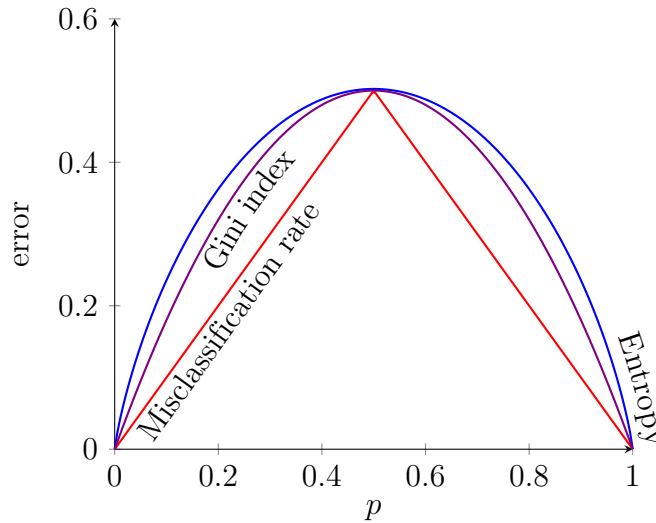
$$Q_k(T) = \sum_{d=1}^D p_{kd}(1 - p_{kd}).$$

- **Cross-entropy (or deviance):**

$$Q_k(T) = - \sum_{d=1}^D p_{kd} \ln p_{kd}.$$

For example, if we have two classes only ($D = 2$) and $p_{k1} = p$, then

- Misclassification error: $\min(p, 1 - p)$
- Gini index: $2p(1 - p)$
- Cross-entropy: $-p \ln p - (1 - p) \ln(1 - p)$



To build the classification decision tree we should choose the error $Q_k(T)$ and run the Algorithm 3.

3 Boosting

3.1 Exponential loss

The natural generalization for the additive models can be obtained by

$$\eta = \sum_{k=1}^K \beta_k b_k(x, \theta_k). \quad (1)$$

For classification problems as before we find the probability by applying the sigmoid function:

$$\mu = g\left(\sum_{k=1}^K \beta_k b_k(x, \theta_k)\right) = g(\eta),$$

where

$$g(z) = \frac{1}{1 + e^{-2z}}$$

(the coefficient 2 is introduced for convenience).

As before we assume

$$y \mid \beta_k; b_k; \theta_k \sim \text{Ber}(\mu),$$

and log-likelihood is expressed as:

$$l(\beta_k, b_k, \theta_k) = \sum_{i=1}^m y^{(i)} \ln \mu^{(i)} + (1 - y^{(i)}) \ln(1 - \mu^{(i)}).$$

Consider the function

$$h(\mu) = y \ln \mu + (1 - y) \ln(1 - \mu).$$

We introduce the transformed target by $z = 2y - 1 \Leftrightarrow y = \frac{z+1}{2}$ (with this notation if $y \in \{0, 1\}$, then $z \in \{-1, 1\}$), then

$$\begin{aligned} h(\mu) &= \frac{z+1}{2} \ln \frac{1}{1 + e^{-2\eta}} + \left(1 - \frac{z+1}{2}\right) \ln \left(1 - \frac{1}{1 + e^{-2\eta}}\right) = \\ &= \frac{1}{2} (-z \ln(1 + e^{-2\eta}) - \ln(1 + e^{-2\eta}) - 2\eta + 2z\eta - \ln(1 + e^{-2\eta}) + z \ln(1 + e^{-2\eta})) = \\ &= -(\ln(1 + e^{-2\eta}) + \eta(1 - z)) = -\ln((1 + e^{-2\eta})e^{\eta(1-z)}) = -\ln(1 + e^{-2z\eta}). \end{aligned}$$

In the last transition we use the fact that $z \in \{-1, 1\}$:

$$\begin{aligned} \text{if } z = 1, \text{ then } \ln((1 + e^{-2\eta})e^{\eta(1-z)}) &= \ln(1 + e^{-2\eta}) = \ln(1 + e^{-2z\eta}), \\ \text{if } z = -1, \text{ then } \ln((1 + e^{-2\eta})e^{\eta(1-z)}) &= \ln(e^{2\eta} + 1) = \ln(1 + e^{-2z\eta}). \end{aligned}$$

With the obtained expression for the likelihood we can formulate the maximum likelihood optimization problem in the simpler way:

$$\arg \max_{\beta_k, b_k, \theta_k} l(\beta_k, b_k, \theta_k) = \arg \min_{\beta_k, b_k, \theta_k} e^{-z\eta},$$

where η is defined with (1) and $z = 2y - 1$ is transformed target. The function

$$L(z, \eta) = \sum_{i=1}^m \exp(-z^{(i)}\eta^{(i)})$$

is called an **exponential loss**.

3.2 Adaboost

Boosting is the example of algorithm ensembling. The idea of ensembling is to combine different algorithms to increase the prediction accuracy. The simplest example of ensembling is to take the output from two algorithms (for example, SVM and neural net) and take the average of two predictions. Boosting is a very powerful algorithm that helps to combine “weak” models (the models with accuracy just a little bit higher than random guess).

Consider the equation (1) where we replace basis functions $b_k(x, \theta)$ by some classifiers $G_k(x) \in \{-1, 1\}$:

$$\eta = \sum_{k=1}^K \beta_k G_k(x). \quad (2)$$

We will use stagewise approach to find weights β_k and functions G_k in (2), the idea is similar to the backfitting algorithm for generalized additive models. Assuming that we have built the models $G_k(x)$ with weights β_k , $k = 1, \dots, K$, we try to find the next weight β and classifier $G(x)$ using exponential loss function:

$$\begin{aligned} \arg \min_{\beta, G} L(z, \eta) &= \arg \min_{\beta, G} \sum_{i=1}^m \exp(-z^{(i)}\eta^{(i)}) = \\ &= \arg \min_{\beta, G} \sum_{i=1}^m \exp\left(-z^{(i)} \left(\sum_{k=1}^K \beta_k G_k(x^{(i)}) + \beta G(x^{(i)})\right)\right) = \\ &= \arg \min_{\beta, G} \sum_{i=1}^m w^{(i)} \exp(-z^{(i)}\beta G(x^{(i)})), \end{aligned}$$

where

$$w^{(i)} = \exp \left(-z^{(i)} \sum_{k=1}^K \beta_k G_k(x^{(i)}) \right) \quad (3)$$

are constants during the algorithm step and can be considered as weights.

If the classifier G predicts the training example i correctly, then

$$z^{(i)} G(x^{(i)}) = 1,$$

otherwise,

$$z^{(i)} G(x^{(i)}) = -1.$$

Hence,

$$\begin{aligned} & \arg \min_{\beta, G} \left(\sum_{i=1}^m w^{(i)} \exp(-z^{(i)} \beta G(x^{(i)})) \right) = \\ & = \arg \min_{\beta, G} \left(\sum_{i: G(x^{(i)})=z^{(i)}} w^{(i)} e^{-\beta} + \sum_{i: G(x^{(i)}) \neq z^{(i)}} w^{(i)} e^{\beta} \right) = \\ & = \arg \min_{\beta, G} \left((e^{\beta} - e^{-\beta}) \cdot \sum_{i=1}^m w^{(i)} \mathbb{1}\{z^{(i)} \neq G(x^{(i)})\} + e^{-\beta} \cdot \sum_{i=1}^m w^{(i)} \right). \end{aligned}$$

This formula gives the following results:

- if β is fixed, then $G = \arg \min_G \sum_{i=1}^m w^{(i)} \mathbb{1}\{z^{(i)} \neq G(x^{(i)})\}$. In other words, our classifier G should minimize the sum of weights for wrongly predicted training examples.
- if G is fixed, then calculating the derivative with respect to β and assigning it to zero gives:

$$\beta = \frac{1}{2} \ln \frac{1-r}{r}, \text{ where } r = \frac{\sum_{i=1}^m w^{(i)} \mathbb{1}\{z^{(i)} \neq G(x^{(i)})\}}{\sum_{i=1}^m w^{(i)}}. \quad (4)$$

This result has a perfect common sense: if $r \rightarrow \frac{1}{2}$ (random classifier), then $\beta \rightarrow 0$; if $r \rightarrow 0$ (perfect classifier), then $\beta \rightarrow \infty$. As we can see the algorithm will give higher weight β to more accurate classifier G .

We can also notice that weights $w^{(i)}$ can be updated from previous step with the formula (3). Indeed,

$$\begin{aligned} w_{new}^{(i)} &= \exp \left(-z^{(i)} \sum_{k=1}^{K-1} \beta_k G_k(x^{(i)}) - z^{(i)} \beta_K G_K(x^{(i)}) \right) = \\ &= w_{old}^{(i)} \cdot \exp \left(-z^{(i)} \beta_K G_K(x^{(i)}) \right), \end{aligned}$$

but $-z^{(i)} G_K(x^{(i)}) = 2 \cdot \mathbb{1}\{z^{(i)} \neq G_K(x^{(i)})\} - 1$, thus

$$w_{new}^{(i)} = w_{old}^{(i)} \cdot e^{\alpha_K \mathbb{1}\{z^{(i)} \neq G_K(x^{(i)})\}} \cdot e^{-\beta_K},$$

where $\alpha_K = 2\beta_K$ and finally, we can remove the term $e^{-\beta_K}$, because it multiplies all weights by the same factor:

$$w_{new}^{(i)} = w_{old}^{(i)} \cdot e^{\alpha_K \mathbb{1}\{z^{(i)} \neq G_K(x^{(i)})\}} \quad (5)$$

If we start fitting with the constant weights for all training examples (for example, $w^{(i)} = \frac{1}{m}$ for all $i = 1, \dots, m$, then after each iteration k we multiply weight for the training example i by the factor

$$e^{\alpha_k \mathbb{1}\{z^{(i)} \neq G_k(x^{(i)})\}} = \begin{cases} 1, & \text{if } G_k \text{ classifies the example } i \text{ correctly} \\ e^{\alpha_k}, & \text{otherwise.} \end{cases}$$

Now we can formulate the first boosting algorithm (Adaboost.M1):

Algorithm 4 Adaboost.M1

- 1: Initialize weights $w^{(i)} = \frac{1}{m}$, $i = 1, 2, \dots, m$
 - 2: **for** $k = 1$ **to** K **do**
 - 3: Fit a classifier $G_k(x)$ to the training data with weights $w^{(i)}$
 - 4: Compute $r = \frac{\sum_{i=1}^m w^{(i)} \mathbb{1}\{z^{(i)} \neq G(x^{(i)})\}}{\sum_{i=1}^m w^{(i)}}$ $\triangleright (4)$
 - 5: Compute $\alpha_k = \ln \frac{1-r}{r}$ (weight for the classifier G_k)
 - 6: Set $w^{(i)} := w^{(i)} \cdot \exp(\alpha_k \cdot \mathbb{1}\{z^{(i)} \neq G_k(x^{(i)})\})$, $i = 1, \dots, m$ $\triangleright (5)$
 - 7: **return** $G(x) = \text{sign} \left(\sum_{k=1}^K \alpha_k G_k(x) \right)$
-

4 Boosting trees

4.1 Gradient boosting

In this section we combine the ideas of backfitting algorithm and decision tree. First, we introduce some loss function $L(f)$, for different problems we can choose squared loss function (for regression)

$$L(f) = \frac{1}{2m} \sum_{i=1}^m (y^{(i)} - f(x^{(i)}))^2,$$

or cross-entropy loss (for classification)

$$L(f) = -\frac{1}{m} \sum_{i=1}^m (y^{(i)} \log f(x^{(i)}) + (1 - y^{(i)}) \log(1 - f(x^{(i)}))) ,$$

or any other convenient loss function. The comparison of the Algorithms 1 and 2 can give the following observations:

- In the Algorithm 2 we try to find minimum of the loss function as a function of f . Replacing the Newton's method by Gradient Descent gives rise to the updating step:

$$f^{new} = f^{old} - \alpha \left. \frac{dL}{df} \right|_{f=f^{old}}, \quad (6)$$

where f^{old} is a vector of current values of f on all training examples, $\left. \frac{dL}{df} \right|_{f=f^{old}}$ is a derivative of the loss function $L(f)$ with respect to f calculated at all current values f^{old} , α is a learning rate.

- In the Algorithm 1 we find the new value for f as

$$f^{new} = f^{old} + g, \quad (7)$$

where g is fitted using some building block algorithm (for example, univariate cubic splines) with previous residues as as working target.

The equations (6) and (7) give a new idea: find the function g by fitting the building block algorithm to the working target $r = -\frac{dL}{df}\bigg|_{f=f^{old}}$ and find the coefficient α as

$$\alpha = \arg \min_{\alpha} L(f^{old} + \alpha g).$$

This idea is called **gradient boosting**.

4.2 Gradient tree boosting

Gradient boosting can be easily implemented when the building block algorithm is a classification or decision tree.

Algorithm 5 Gradient Tree Boosting

- 1: Initialize $f_0(x) = \arg \min_{\mu} \sum_{i=1}^m L(y^{(i)}, \mu)$
- 2: **for** $k = 1$ **to** K **do**
- 3: Compute working target $r_k^{(i)} = -\left(\frac{dL}{df}\right)\bigg|_{f=f_{k-1}(x^{(i)})}$ for all $i = 1, \dots, m$
- 4: Fit a regression tree to the targets $r_k^{(i)}$ with terminal nodes R_{kj} , $j = 1, \dots, J_k$
- 5: Compute

$$\gamma_{kj} = \arg \min_{\gamma} \sum_{x^{(i)} \in R_{kj}} L(y^{(i)}, f_{k-1}(x^{(i)}) + \gamma)$$

for all $j = 1, \dots, J_k$

- 6: Update $f_k(x) = f_{k-1}(x) + \sum_{j=1}^{J_k} \gamma_{kj} \mathbb{1}\{x \in R_{kj}\}$

- 7: **return** $f_K(x)$
-

In the line 5 of the algorithm we should find the best value for γ in each terminal node (for squared loss function it will be just an average of working target).

The classification algorithm is similar except that in the line 4 we should fit a classification tree and also we should repeat lines 3-6 for each class.

Exercise. Write down the pseudocode for the Gradient Tree Boosting for the classification problem with

- a) two classes;
- b) D classes.

References

- [1] T. Hastie and R. Tibshirani. “Generalized additive models”, *in Statistical Science, Vol.1, No.3*, 1998, pp. 297-318