

Ensembles II

Jiayu Zhou

¹Department of Computer Science and Engineering
Michigan State University
East Lansing, MI USA

March 17, 2016

Slides from "A Gentle Introduction to Gradient Boosting." by Cheng Li.

Introduction

Gradient Boosting

- A powerful machine learning algorithm
- Regression/Classification/Ranking.
- Won Track 1 of the Yahoo Learning to Rank Challenge

What is Gradient Boosting

- **Gradient Boosting = Gradient Descent + Boosting**
- AdaBoost

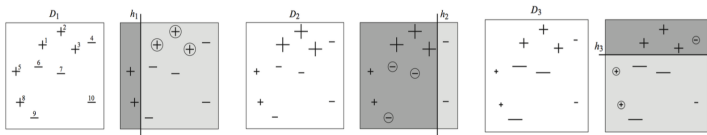


Figure : AdaBoost. Source: Figure 1.1 of [Schapire and Freund, 2012]

What is Gradient Boosting

- **Gradient Boosting = Gradient Descent + Boosting**
- AdaBoost

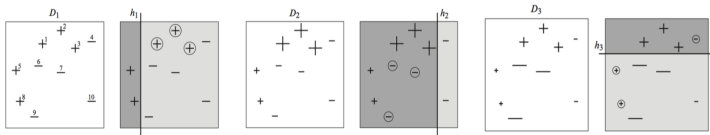


Figure : AdaBoost. Source: Figure 1.1 of [Schapire and Freund, 2012]

- Fit an additive model (ensemble) $\sum_t \rho_t h_t(x)$ in a forward stage-wise manner.
- In each stage, introduce a weak learner to compensate the shortcomings of existing weak learners.
- In Adaboost, “shortcomings” are identified by high-weight data points.

What is Gradient Boosting

- Gradient Boosting = Gradient Descent + Boosting
- AdaBoost

$$H(x) = \sum_t \rho_t h_t(x)$$

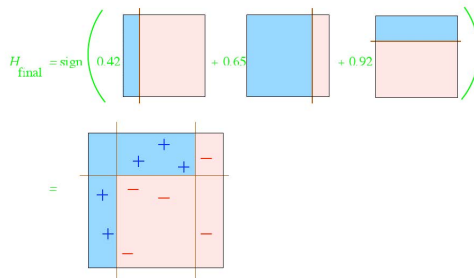


Figure : AdaBoost. Source: Figure 1.1 of [Schapire and Freund, 2012]

What is Gradient Boosting

- **Gradient Boosting = Gradient Descent + Boosting**
- Gradient Boosting
 - Fit an additive model (ensemble) $\sum_t \rho_t h_t(x)$ in a forward stage-wise manner.
 - In each stage, introduce a weak learner to compensate the shortcomings of existing weak learners.
 - In Gradient Boosting, shortcomings are identified by gradients.
 - In Adaboost, shortcomings are identified by high-weight data points.
 - Both high-weight data points and gradients tell us how to improve our model.

Regression

Gradient Boosting for Regression

Lets play a game...

You are given $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$, and the task is to fit a model $F(x)$ to minimize square loss.

Suppose your friend wants to help you and gives you a model F . You check his model and find the model is good but not perfect. There are some mistakes: $F(x_1) = 0.8$, while $y_1 = 0.9$, and $F(x_2) = 1.4$ while $y_2 = 1.3$... How can you improve this model?

Rule of the game:

- You are not allowed to remove anything from F or change any parameter in F .
- You can add an additional model (regression) h to F , so the new prediction will be $F(x) + h(x)$.

Gradient Boosting for Regression

Simple solution: You wish to improve the model such that

$$F(x_1) + h(x_1) = y_1$$

$$F(x_2) + h(x_2) = y_2$$

• • •

$$F(x_n) + h(x_n) = y_n.$$

Gradient Boosting for Regression

Simple solution: Or, equivalently, you wish

$$h(x_1) = y_1 - F(x_1)$$

$$h(x_2) = y_2 - F(x_2)$$

• • •

$$h(x_n) = y_n - F(x_n)$$

Can any regression model achieve goal perfectly?

Gradient Boosting for Regression

Simple solution: Or, equivalently, you wish

$$h(x_1) = y_1 - F(x_1)$$

$$h(x_2) = y_2 - F(x_2)$$

• • •

$$h(x_n) = y_n - F(x_n)$$

Can any regression model achieve goal perfectly?

Maybe not ...

Simple solution: Or, equivalently, you wish

$$h(x_n) = y_n - F(x_n)$$

Maybe not ...

But some regression models might be able to do this approximately.

Gradient Boosting for Regression

Simple solution: Or, equivalently, you wish

$$h(x_1) = y_1 - F(x_1)$$

$$h(x_2) = y_2 - F(x_2)$$

• • •

$$h(x_n) = y_n - F(x_n)$$

Can any regression model achieve goal perfectly?

Maybe not ...

But some regression models might be able to do this approximately.

How?

Gradient Boosting for Regression

Simple solution: Or, equivalently, you wish

$$h(x_1) = y_1 - F(x_1)$$

$$h(x_2) = y_2 - F(x_2)$$

• • •

$$h(x_n) = y_n - F(x_n)$$

Can any regression model achieve goal perfectly?

Maybe not ...

But some regression models might be able to do this approximately.

How?

Just fit a regression model h to data

$$(x_1, y_1 - F(x_1)), (x_2, y_2 - F(x_2)), \dots, (x_n, y_n - F(x_n))$$

Gradient Boosting for Regression

Simple solution: Or, equivalently, you wish

$$h(x_1) = y_1 - F(x_1)$$

$$h(x_2) = y_2 - F(x_2)$$

• • •

$$h(x_n) = y_n - F(x_n)$$

Can any regression model achieve goal perfectly?

Maybe not ...

But some regression models might be able to do this approximately.

How?

Just fit a regression model h to data

$$(x_1, y_1 - F(x_1)), (x_2, y_2 - F(x_2)), \dots, (x_n, y_n - F(x_n))$$

Congratulations, you get a better model!

Simple solution:

The role of h is to compensate the shortcoming of existing model F .

Simple solution:

The role of h is to compensate the shortcoming of existing model F . If the new model $F + h$ is still not satisfactory,

Simple solution:

The role of h is to compensate the shortcoming of existing model F . If the new model $F + h$ is still not satisfactory, we can add another regression model...

Simple solution:

The role of h is to compensate the shortcoming of existing model F .

If the new model $F + h$ is still not satisfactory, we can add another regression model...

We are improving the predictions of training data, is the procedure also useful for test data?

Gradient Boosting for Regression

Simple solution:

$y_i - F(x_i)$ are called residuals. These are the parts that existing model F cannot do well.

The role of h is to compensate the shortcoming of existing model F .

If the new model $F + h$ is still not satisfactory, we can add another regression model...

We are improving the predictions of training data, is the procedure also useful for test data?

Yes! Because we are building a model, and the model can be applied to test data as well.

Gradient Boosting for Regression

Simple solution:

$y_i - F(x_i)$ are called residuals. These are the parts that existing model F cannot do well.

The role of h is to compensate the shortcoming of existing model F .

If the new model $F + h$ is still not satisfactory, we can add another regression model...

We are improving the predictions of training data, is the procedure also useful for test data?

Yes! Because we are building a model, and the model can be applied to test data as well.

How is this related to gradient descent?

How is this related to gradient descent?

We want to minimize $J = \sum_i L(y_i, F(x_i))$ by adjusting

Notice that $F(x_1), F(x_2), \dots, F(x_n)$ are just some numbers. We can treat $F(x_i)$ as parameters and take derivatives.

Gradient Boosting for Regression

How is this related to gradient descent?

Loss function $L(y, F(x)) = (y - F(x))^2/2$

We want to minimize $J = \sum_i L(y_i, F(x_i))$ by adjusting $F(x_1), F(x_2), \dots, F(x_n)$.

Notice that $F(x_1), F(x_2), \dots, F(x_n)$ are just some numbers. We can treat $F(x_i)$ as parameters and take derivatives.

$$\frac{\partial \sum_i L(y_i, F(x_i))}{\partial F(x_i)} = \frac{\partial L(y_i, F(x_i))}{F(x_i)} = F(x_i) - y_i$$

So we can interpret residuals as negative gradients:

$$-\frac{\sum_i L(y_i, F(x_i))}{\partial F(x_i)} = y_i - F(x_i)$$

Gradient Boosting for Regression

How is this related to gradient descent?

$$F(x_i) = F(x_i) + h(x_i)$$

$$F(x_i) = F(x_i) + y_i - F(x_i)$$

$$F(x_i) = F(x_i) - 1 \frac{\partial \sum_i L(y_i, F(x_i))}{\partial F(x_i)}$$

$$\theta_i = \theta_i - \rho \frac{\partial J}{\partial \theta_i}$$

How is this related to gradient descent?

```

residual = negative gradient
fit h to residual = fit h to negative gradient
update F based on residual = update F based on negative gradient

```

How is this related to gradient descent?

```

    residual = negative gradient
    fit h to residual = fit h to negative gradient
    update F based on residual = update F based on negative gradient

```

So we are actually updating our model using **gradient descent**!

How is this related to gradient descent?

residual = negative gradient

fit h to residual = fit h to negative gradient

update F based on residual = update F based on negative gradient

So we are actually updating our model using **gradient descent**!

It turns out that the concept of **gradients** is more general and useful than the concept of **residuals**. So from now on, let's stick with gradients.

The reason will be explained later.

Loss Functions for Regression Problem

Why do we need to consider other loss functions? Isn't square loss good enough?

Loss Functions for Regression Problem

- $$L(y, F) = |y - F|$$

Gradient Boosting for Regression

Regression with Absolute Loss

Negative (sub)gradient

$$-g(x_i) = -\frac{\partial \sum_i L(y_i, F(x_i))}{\partial F(x_i)} = \text{sign}(y_i - F(x_i))$$

- ① start with an initial model, say, $F(x) = \frac{\sum_{i=1}^n y_i}{n}$
- ② iterate until converge:
 - calculate negative gradients $-g(x_i)$
 - fit a regression model h to negative gradients $-g(x_i)$
 - $F = F + \rho h$, where $\rho = 1$

Gradient Boosting for Regression

Regression with Huber Loss

Negative (sub)gradient

$$\begin{aligned} -g(x_i) &= -\frac{\partial \sum_i L(y_i, F(x_i))}{\partial F(x_i)} \\ &= \begin{cases} y_i - F(x_i) & |y_i - F(x_i)| \leq \delta \\ \delta \text{sign}(y_i - F(x_i)) & |y_i - F(x_i)| > \delta \end{cases} \end{aligned}$$

- 1 start with an initial model, say, $F(x) = \frac{\sum_{i=1}^n y_i}{n}$
- 2 iterate until converge:
 - calculate negative gradients $-g(x_i)$
 - fit a regression model h to negative gradients $-g(x_i)$
 - $F = F + \rho h$, where $\rho = 1$

Gradient Boosting for Regression

Regression with loss function L : general procedure

Given any (sub)differentiable loss function L

- 1 start with an initial model, say, $F(x) = \frac{\sum_{i=1}^n y_i}{n}$
- 2 iterate until converge:
 - calculate negative gradients $-g(x_i) = -\frac{\partial L(y_i, F(x_i))}{\partial F(x_i)}$
 - fit a regression model h to negative gradients $-g(x_i)$
 - $F = F + \rho h$, where $\rho = 1$

Gradient Boosting for Regression

Regression with loss function L : general procedure

Given any (sub)differentiable loss function L

- 1 start with an initial model, say, $F(x) = \frac{\sum_{i=1}^n y_i}{n}$
- 2 iterate until converge:
 - calculate negative gradients $-g(x_i) = -\frac{\partial L(y_i, F(x_i))}{\partial F(x_i)}$
 - fit a regression model h to negative gradients $-g(x_i)$
 - $F = F + \rho h$, where $\rho = 1$

In general,

negative gradients \neq residuals

We should follow negative gradients rather than residuals. Why?

Gradient Boosting for Regression

Negative Gradient vs Residual: An Example

Huber loss

$$L(y, F) = \begin{cases} \frac{1}{2}(y - F)^2 & |y - F| \leq \delta \\ \delta(|y - F| - \delta/2) & |y - F| > \delta \end{cases}$$

Update by Negative Gradient:

$$h(x_i) = -g(x_i) = \begin{cases} y_i - F(x_i) & |y_i - F(x_i)| \leq \delta \\ \delta \text{sign}(y_i - F(x_i)) & |y_i - F(x_i)| > \delta \end{cases}$$

Update by Residual:

$$h(x_i) = y_i - F(x_i)$$

Difference: negative gradient pays less attention to outliers.

Gradient Boosting for Regression

Summary of the Section

- Fit an additive model $F = \sum_t \rho_t h_t$ in a forward stage-wise manner.
- In each stage, introduce a new regression tree h to compensate the shortcomings of existing model.
- The “shortcomings” are identified by negative gradients.
- For any loss function, we can derive a gradient boosting algorithm.
- Absolute loss and Huber loss are more robust to outliers than square loss.

Things not covered

How to choose a proper learning rate for each gradient boosting algorithm. See [Friedman, 2001]

Classification

Gradient Boosting for Classification

Problem

Recognize the given hand written capital letter.

- Multi-class classification
- 26 classes. A,B,C,...,Z



Data Set

- <http://archive.ics.uci.edu/ml/datasets/Letter+Recognition>
- 20000 data points, 16 features (how do we extract features in this case?)

- $$Y_A(x_5) = 0, Y_B(x_5) = 0, \dots, Y_G(x_5) = 1, \dots, Y_Z(x_5) = 0$$

Loss Function for each data point

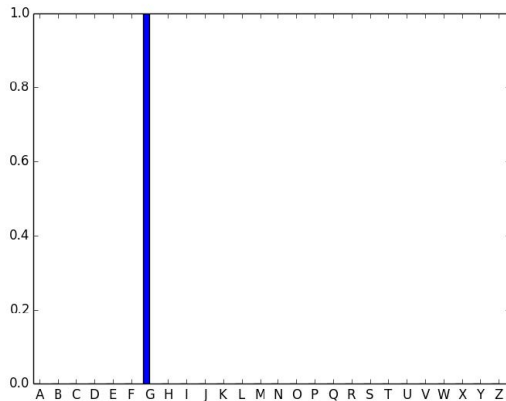


Figure: true probability distribution

Loss Function for each data point

- ➊ turn the label y_i into a (true) probability distribution $Y_c(x_i)$
For example $y_5 = G$, then

$$Y_A(x_5) = 0, Y_B(x_5) = 0, \dots, Y_G(x_5) = 1, \dots, Y_Z(x_5) = 0$$

- 2 calculate the predicted probability distribution $P_c(x_i)$ based on the current model F_A, F_B, \dots, F_Z .

$$P_A(x_5) = 0.03, P_B(x_5) = 0.05, \dots, P_G(x_5) = 0.3, \dots, P_Z(x_5) = 0.05$$

Loss Function for each data point

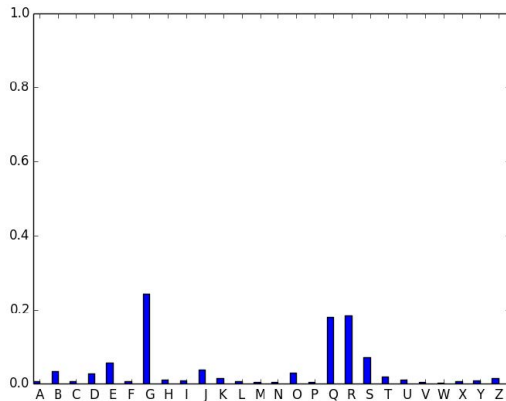
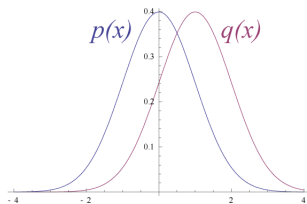


Figure: predicted probability distribution based on current model

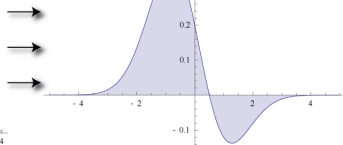
Kullback-Leibler (KL) divergence

- Measures of the difference between two probability distributions P (true distribution) and Q (model distribution)

$$D_{\text{KL}}(P||Q) = \sum_x P(x) \log \frac{P(x)}{Q(x)}$$
$$= - \sum_x p(x) \log q(x) + \sum_x p(x) \log p(x) = H(P, Q) - H(P)$$



Original Gaussian PDF's



KL Area to be Integrated

Kullback-Leibler Divergence

Example For a random variable $X = \{0, 1\}$ assume two distributions $P(x)$ and $Q(x)$ with $P(0) = 1 - r$, $P(1) = r$ and $Q(0) = 1 - s$, $Q(1) = s$:

$$D(P||Q) = (1 - r) \log \frac{1 - r}{1 - s} + r \log \frac{r}{s}$$

$$D(Q||P) = (1 - s) \log \frac{1 - s}{1 - r} + s \log \frac{s}{r}$$

If $r = s$

Kullback-Leibler Divergence

Example For a random variable $X = \{0, 1\}$ assume two distributions $P(x)$ and $Q(x)$ with $P(0) = 1 - r$, $P(1) = r$ and $Q(0) = 1 - s$, $Q(1) = s$:

$$D(P||Q) = (1 - r) \log \frac{1 - r}{1 - s} + r \log \frac{r}{s}$$

$$D(Q||P) = (1 - s) \log \frac{1 - s}{1 - r} + s \log \frac{s}{r}$$

If $r = s$ then $D(P||Q) = D(Q||P) = 0$. If $r = \frac{1}{2}$ and $s = \frac{1}{4}$:

Kullback-Leibler Divergence

Example For a random variable $X = \{0, 1\}$ assume two distributions $P(x)$ and $Q(x)$ with $P(0) = 1 - r$, $P(1) = r$ and $Q(0) = 1 - s$, $Q(1) = s$:

$$D(P||Q) = (1 - r) \log \frac{1 - r}{1 - s} + r \log \frac{r}{s}$$

$$D(Q||P) = (1 - s) \log \frac{1 - s}{1 - r} + s \log \frac{s}{r}$$

If $r = s$ then $D(P||Q) = D(Q||P) = 0$. If $r = \frac{1}{2}$ and $s = \frac{1}{4}$:

$$D(P||Q) = \frac{1}{2} \log \frac{1/2}{3/4} + \frac{1}{2} \log \frac{1/2}{1/4} = 0.2075$$

$$D(Q||P) = \frac{3}{4} \log \frac{3/4}{1/2} + \frac{1}{4} \log \frac{1/4}{1/2} = 0.1887$$

Properties of the Kullback-Leibler Divergence

- $D(P||Q) \geq 0$
- $D(P||Q) = 0$ iff $P(x) = Q(x)$ for all $x \in X$;
- Typically $D(P||Q) \neq D(Q||P)$

Multi-Class Classification via KL Divergence

Goal

- minimize the total loss (KL-divergence)
- for each data point, we wish the predicted probability distribution to match the true probability distribution as closely as possible

Minimizing KL Divergence

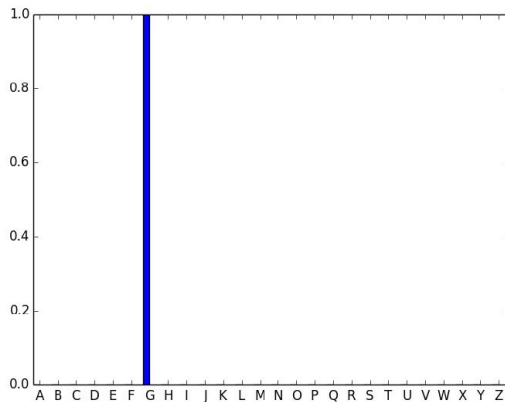


Figure: true probability distribution

Minimizing KL Divergence

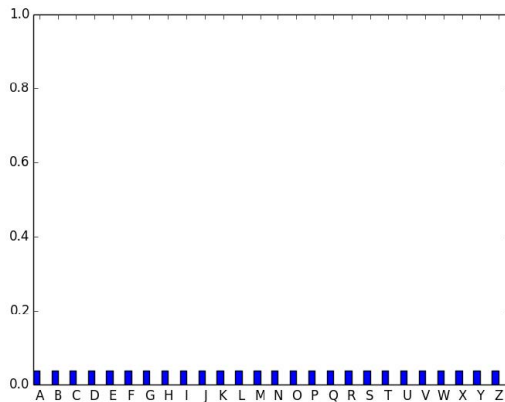


Figure: predicted probability distribution at round 0

Minimizing KL Divergence

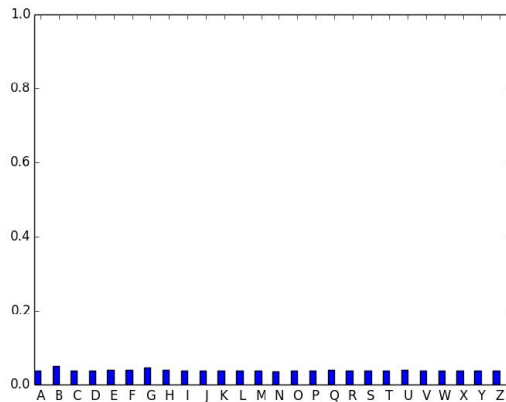


Figure: predicted probability distribution at round 1

Minimizing KL Divergence

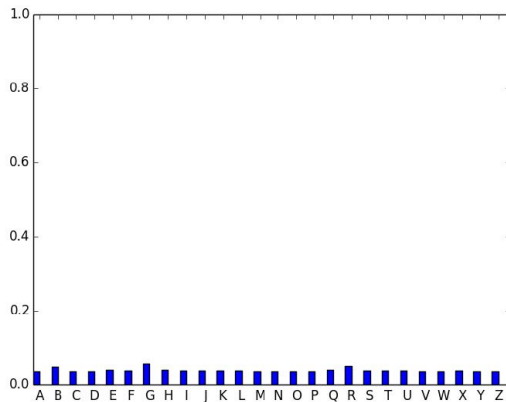


Figure: predicted probability distribution at round 2

Minimizing KL Divergence

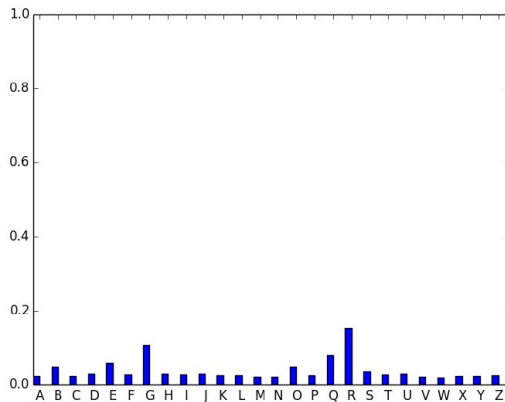


Figure: predicted probability distribution at round 10

Minimizing KL Divergence

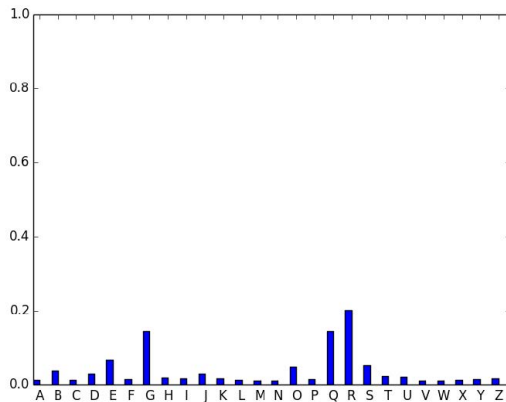


Figure: predicted probability distribution at round 20

Minimizing KL Divergence

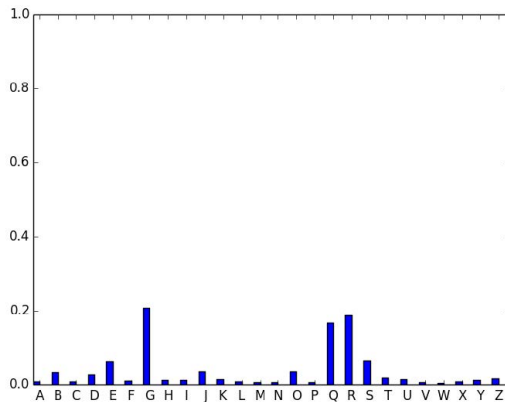


Figure: predicted probability distribution at round 30

Minimizing KL Divergence

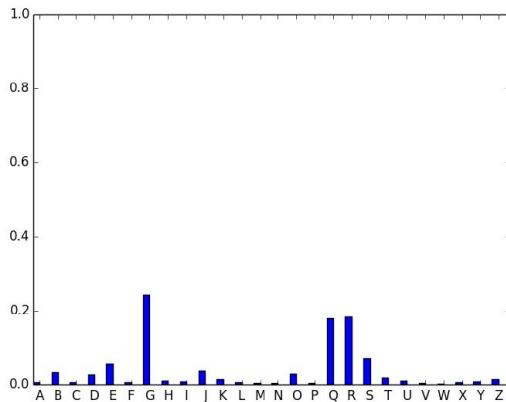


Figure: predicted probability distribution at round 40

Minimizing KL Divergence

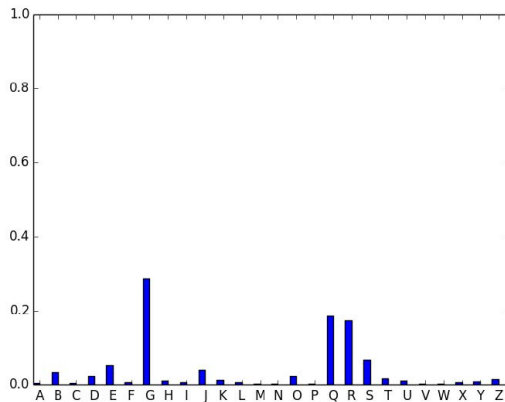


Figure: predicted probability distribution at round 50

Minimizing KL Divergence

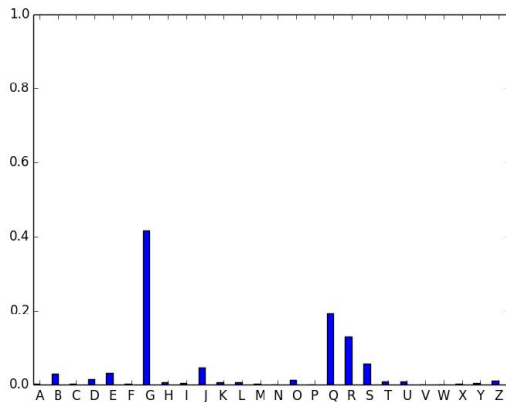


Figure: predicted probability distribution at round 100

The Classification Problem

Goal

- minimize the total loss (KL-divergence)
- for each data point, we wish the predicted probability distribution to match the true probability distribution as closely as possible
- we achieve this goal by adjusting our models F_A, F_B, \dots, F_Z .

Gradient Boosting for Regression: Review

Regression with loss function L : general procedure

Give any differentiable loss function L

start with an initial model F

iterate until converge:

- calculate negative gradients $-g(x_i) = -\frac{\partial L(y_i, F(x_i))}{\partial F(x_i)}$
- fit a regression model h to negative gradients $-g(x_i)$
- $F = F + \rho h$

Gradient Boosting for Classification

Difference between classification and regression

- F_A, F_B, \dots, F_Z **vs.** F
- a matrix of parameters to optimize **vs.** a column of parameters to optimize

$F_A(x_1)$	$F_B(x_1)$	\dots	$F_Z(x_1)$
$F_A(x_2)$	$F_B(x_2)$	\dots	$F_Z(x_2)$
\dots	\dots	\dots	\dots
$F_A(x_d)$	$F_B(x_d)$	\dots	$F_Z(x_d)$

- a matrix of gradients **vs.** a column of gradients

$\frac{\partial L}{\partial F_A(x_1)}$	$\frac{\partial L}{\partial F_B(x_1)}$	\dots	$\frac{\partial L}{\partial F_Z(x_1)}$
$\frac{\partial L}{\partial F_A(x_2)}$	$\frac{\partial L}{\partial F_B(x_2)}$	\dots	$\frac{\partial L}{\partial F_Z(x_2)}$
\dots	\dots	\dots	\dots
$\frac{\partial L}{\partial F_A(x_d)}$	$\frac{\partial L}{\partial F_B(x_d)}$	\dots	$\frac{\partial L}{\partial F_Z(x_d)}$

Gradient Boosting for Classification

start with an initial models $F_A, F_B, F_C, \dots, F_Z$

iterate until converge:

- calculate negative gradients for class A : $-g_A(x_i) = -\frac{\partial L(y_i, F(x_i))}{\partial F_A(x_i)}$
- calculate negative gradients for class B : $-g_B(x_i) = -\frac{\partial L(y_i, F(x_i))}{\partial F_B(x_i)}$
- ...
- calculate negative gradients for class Z : $-g_Z(x_i) = -\frac{\partial L(y_i, F(x_i))}{\partial F_Z(x_i)}$
- fit a regression model h_A to negative gradients $-g_A(x_i)$
- fit a regression model h_B to negative gradients $-g_B(x_i)$
- ...
- fit a regression model h_Z to negative gradients $-g_Z(x_i)$
- $F_A = F_A + \rho_A h_A$
- $F_B = F_B + \rho_B h_B$
- ...
- $F_Z = F_Z + \rho_Z h_Z$

Gradient Boosting for Classification

start with an initial models $F_A, F_B, F_C, \dots, F_Z$

iterate until converge:

- calculate negative gradients for class A : $-g_A(x_i) = Y_A(x_i) - P_A(x_i)$
- calculate negative gradients for class B :
 $-g_B(x_i) = Y_B(x_i) - P_B(x_i)$
- ...
- calculate negative gradients for class Z : $-g_Z(x_i) = Y_Z(x_i) - P_Z(x_i)$
- fit a regression model h_A to negative gradients $-g_A(x_i)$
- fit a regression model h_B to negative gradients $-g_B(x_i)$
- ...
- fit a regression model h_Z to negative gradients $-g_Z(x_i)$
- $F_A = F_A + \rho_A h_A$
- $F_B = F_B + \rho_B h_B$
- ...
- $F_Z = F_Z + \rho_Z h_Z$

Gradient Boosting for Classification

round 0

i	y	Y _A	Y _B	Y _C	Y _D	Y _E	Y _F	Y _G	Y _H	Y _I	Y _J	Y _K	Y _L	Y _M	Y _N	Y _O	Y _P	Y _Q	Y _R	Y _S	Y _T	Y _U	Y _V	Y _W	Y _X	Y _Y	Y _Z
1	T	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0
2	I	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
3	D	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
4	N	0	0	0	0	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0
5	G	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0

i	y	F _A	F _B	F _C	F _D	F _E	F _F	F _G	F _H	F _I	F _J	F _K	F _L	F _M	F _N	F _O	F _P	F _Q	F _R	F _S	F _T	F _U	F _V	F _W	F _X	F _Y	F _Z
1	T	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
2	I	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
3	D	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
4	N	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
5	G	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0

i	y	P _A	P _B	P _C	P _D	P _E	P _F	P _G	P _H	P _I	P _J	P _K	P _L	P _M	P _N	P _O	P _P	P _Q	P _R	P _S	P _T	P _U	P _V	P _W	P _X	P _Y	P _Z
1	T	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04
2	I	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04
3	D	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04
4	N	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04
5	G	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04

i	y	Y _A - P _A	Y _B - P _B	Y _C - P _C	Y _D - P _D	Y _E - P _E	Y _F - P _F	Y _G - P _G	Y _H - P _H	Y _I - P _I	Y _J - P _J	Y _K - P _K	Y _L - P _L	Y _M - P _M	Y _N - P _N	Y _O - P _O	Y _P - P _P	Y _Q - P _Q	Y _R - P _R	Y _S - P _S	Y _T - P _T	Y _U - P _U	Y _V - P _V	Y _W - P _W	Y _X - P _X	Y _Y - P _Y	Y _Z - P _Z
1	T	-0.04	-0.04	-0.04	-0.04	-0.04	-0.04	-0.04	-0.04	-0.04	-0.04	-0.04	-0.04	-0.04	-0.04	-0.04	-0.04	-0.04	-0.04	-0.04	-0.04	0.96	-0.04	-0.04	-0.04	-0.04	-0.04
2	I	-0.04	-0.04	-0.04	-0.04	-0.04	-0.04	-0.04	-0.04	0.96	-0.04	-0.04	-0.04	-0.04	-0.04	-0.04	-0.04	-0.04	-0.04	-0.04	-0.04	-0.04	-0.04	-0.04	-0.04	-0.04	-0.04
3	D	-0.04	-0.04	-0.04	0.96	-0.04	-0.04	-0.04	-0.04	-0.04	-0.04	-0.04	-0.04	-0.04	-0.04	-0.04	-0.04	-0.04	-0.04	-0.04	-0.04	-0.04	-0.04	-0.04	-0.04	-0.04	-0.04
4	N	-0.04	-0.04	-0.04	-0.04	-0.04	-0.04	-0.04	-0.04	-0.04	-0.04	-0.04	-0.04	-0.04	0.96	-0.04	-0.04	-0.04	-0.04	-0.04	-0.04	-0.04	-0.04	-0.04	-0.04	-0.04	-0.04
5	G	-0.04	-0.04	-0.04	-0.04	-0.04	-0.04	0.96	-0.04	-0.04	-0.04	-0.04	-0.04	-0.04	-0.04	-0.04	-0.04	-0.04	-0.04	-0.04	-0.04	-0.04	-0.04	-0.04	-0.04	-0.04	-0.04

Gradient Boosting for Classification

$$h_A(x) = \begin{cases} 0.98 & \text{feature 10 of } x \leq 2.0 \\ -0.07 & \text{feature 10 of } x > 2.0 \end{cases}$$

$$h_B(x) = \begin{cases} -0.07 & \text{feature 15 of } x \leq 8.0 \\ 0.22 & \text{feature 15 of } x > 8.0 \end{cases}$$

...

$$h_Z(x) = \begin{cases} -0.07 & \text{feature 8 of } x \leq 8.0 \\ 0.82 & \text{feature 8 of } x > 8.0 \end{cases}$$

$$F_A = F_A + \rho_A h_A$$

$$F_B = F_B + \rho_B h_B$$

...

$$F_Z = F_Z + \rho_Z h_Z$$

Gradient Boosting for Classification

round 1

i	y	Y _A	Y _B	Y _C	Y _D	Y _E	Y _F	Y _G	Y _H	Y _I	Y _J	Y _K	Y _L	Y _M	Y _N	Y _O	Y _P	Y _Q	Y _R	Y _S	Y _T	Y _U	Y _V	Y _W	Y _X	Y _Y	Y _Z
1	T	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0
2	I	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
3	D	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
4	N	0	0	0	0	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0
5	G	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0

i	y	F _A	F _B	F _C	F _D	F _E	F _F	F _G	F _H	F _I	F _J	F _K	F _L	F _M	F _N	F _O	F _P	F _Q	F _R	F _S	F _T	F _U	F _V	F _W	F _X	F _Y	F _Z
1	T	-0.08	-0.07	-0.06	-0.07	-0.02	-0.02	-0.08	-0.02	-0.03	-0.03	-0.06	-0.04	-0.08	-0.08	-0.07	-0.07	-0.02	-0.04	-0.04	0.59	-0.01	-0.07	-0.07	-0.05	-0.06	-0.07
2	I	-0.08	0.23	-0.06	-0.07	-0.02	-0.02	0.16	-0.02	-0.03	-0.03	-0.06	-0.04	-0.08	-0.08	-0.07	-0.07	-0.02	-0.04	-0.04	-0.07	-0.01	-0.07	-0.07	-0.05	-0.06	-0.07
3	D	-0.08	0.23	-0.06	-0.07	-0.02	-0.02	-0.08	-0.02	-0.03	-0.03	-0.06	-0.04	-0.08	-0.08	-0.07	-0.07	-0.02	-0.04	-0.04	-0.07	-0.01	-0.07	-0.07	-0.05	-0.06	-0.07
4	N	-0.08	-0.07	-0.06	-0.07	-0.02	-0.02	0.16	-0.02	-0.03	-0.03	0.26	-0.04	-0.08	0.3	-0.07	-0.07	-0.02	-0.04	-0.04	-0.07	-0.01	-0.07	-0.07	-0.05	-0.06	-0.07
5	G	-0.08	0.23	-0.06	-0.07	-0.02	-0.02	0.16	-0.02	-0.03	-0.03	-0.06	-0.04	-0.08	-0.08	-0.07	-0.07	-0.02	-0.04	-0.04	-0.07	-0.01	-0.07	-0.07	-0.05	-0.06	-0.07

i	y	P _A	P _B	P _C	P _D	P _E	P _F	P _G	P _H	P _I	P _J	P _K	P _L	P _M	P _N	P _O	P _P	P _Q	P _R	P _S	P _T	P _U	P _V	P _W	P _X	P _Y	P _Z
1	T	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.07	0.04	0.04	0.04	0.04	0.04	0.04
2	I	0.04	0.05	0.04	0.04	0.04	0.04	0.05	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04
3	D	0.04	0.05	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04
4	N	0.04	0.04	0.04	0.04	0.04	0.04	0.05	0.04	0.04	0.04	0.05	0.04	0.04	0.05	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04
5	G	0.04	0.05	0.04	0.04	0.04	0.04	0.05	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04

i	y	Y _A - P _A	Y _B - P _B	Y _C - P _C	Y _D - P _D	Y _E - P _E	Y _F - P _F	Y _G - P _G	Y _H - P _H	Y _I - P _I	Y _J - P _J	Y _K - P _K	Y _L - P _L	Y _M - P _M	Y _N - P _N	Y _O - P _O	Y _P - P _P	Y _Q - P _Q	Y _R - P _R	Y _S - P _S	Y _T - P _T	Y _U - P _U	Y _V - P _V	Y _W - P _W	Y _X - P _X	Y _Y - P _Y	Y _Z - P _Z
1	T	-0.04	-0.04	-0.04	-0.04	-0.04	-0.04	-0.04	-0.04	-0.04	-0.04	-0.04	-0.04	-0.04	-0.04	-0.04	-0.04	-0.04	-0.04	-0.04	0.93	-0.04	-0.04	-0.04	-0.04	-0.04	-0.04
2	I	-0.04	-0.05	-0.04	-0.04	-0.04	-0.04	-0.04	-0.04	0.96	-0.04	-0.04	-0.04	-0.04	-0.04	-0.04	-0.04	-0.04	-0.04	-0.04	-0.04	-0.04	-0.04	-0.04	-0.04	-0.04	-0.04
3	D	-0.04	-0.05	-0.04	0.96	-0.04	-0.04	-0.04	-0.04	-0.04	-0.04	-0.04	-0.04	-0.04	-0.04	-0.04	-0.04	-0.04	-0.04	-0.04	-0.04	-0.04	-0.04	-0.04	-0.04	-0.04	-0.04
4	N	-0.04	-0.04	-0.04	-0.04	-0.04	-0.04	-0.05	-0.04	-0.04	-0.04	-0.05	-0.04	-0.04	0.95	-0.04	-0.04	-0.04	-0.04	-0.04	-0.04	-0.04	-0.04	-0.04	-0.04	-0.04	-0.04
5	G	-0.04	-0.05	-0.04	-0.04	-0.04	-0.04	0.95	-0.04	-0.04	-0.04	-0.04	-0.04	-0.04	-0.04	-0.04	-0.04	-0.04	-0.04	-0.04	-0.04	-0.04	-0.04	-0.04	-0.04	-0.04	-0.04

Gradient Boosting for Classification

$$h_A(x) = \begin{cases} 0.37 & \text{feature 10 of } x \leq 2.0 \\ -0.07 & \text{feature 10 of } x > 2.0 \end{cases}$$

$$h_B(x) = \begin{cases} -0.07 & \text{feature 14 of } x \leq 5.0 \\ 0.22 & \text{feature 14 of } x > 5.0 \end{cases}$$

...

$$h_Z(x) = \begin{cases} -0.07 & \text{feature 8 of } x \leq 8.0 \\ 0.35 & \text{feature 8 of } x > 8.0 \end{cases}$$

$$F_A = F_A + \rho_A h_A$$

$$F_B = F_B + \rho_B h_B$$

...

$$F_Z = F_Z + \rho_Z h_Z$$

Gradient Boosting for Classification

round 2

i	y	Y _A	Y _B	Y _C	Y _D	Y _E	Y _F	Y _G	Y _H	Y _I	Y _J	Y _K	Y _L	Y _M	Y _N	Y _O	Y _P	Y _Q	Y _R	Y _S	Y _T	Y _U	Y _V	Y _W	Y _X	Y _Y	Y _Z
1	T	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0
2	I	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
3	D	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
4	N	0	0	0	0	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0
5	G	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0

i	y	F _A	F _B	F _C	F _D	F _E	F _F	F _G	F _H	F _I	F _J	F _K	F _L	F _M	F _N	F _O	F _P	F _Q	F _R	F _S	F _T	F _U	F _V	F _W	F _X	F _Y	F _Z
1	T	-0.15	-0.14	-0.12	-0.14	-0.03	0.28	-0.14	-0.04	1.49	-0.07	-0.11	-0.08	-0.14	-0.17	-0.13	-0.13	-0.04	-0.11	-0.07	1.05	0.19	0.25	-0.16	-0.09	0.33	-0.14
2	I	-0.15	0.16	-0.12	-0.14	-0.03	-0.08	0.33	-0.04	-0.07	-0.07	-0.11	-0.08	-0.14	-0.17	-0.13	-0.13	-0.04	-0.11	-0.07	-0.11	-0.07	-0.15	-0.16	-0.09	-0.13	-0.14
3	D	-0.15	0.16	-0.12	-0.14	-0.03	-0.08	0.1	-0.04	-0.07	-0.07	-0.11	-0.08	-0.14	-0.17	-0.13	-0.13	-0.04	0.19	-0.07	-0.11	-0.07	-0.15	-0.16	-0.09	-0.13	-0.14
4	N	-0.15	-0.14	-0.12	-0.14	-0.03	-0.08	0.1	-0.04	-0.07	-0.07	0.46	-0.08	-0.14	0.5	-0.13	-0.13	-0.04	-0.11	-0.07	-0.11	-0.07	-0.15	0.25	-0.09	-0.13	-0.14
5	G	-0.15	0.16	-0.12	-0.14	-0.03	-0.08	0.33	-0.04	-0.07	-0.07	-0.11	-0.08	-0.14	-0.17	-0.13	-0.13	-0.04	0.19	-0.07	-0.11	-0.07	-0.15	-0.16	-0.09	-0.13	-0.14

i	y	P _A	P _B	P _C	P _D	P _E	P _F	P _G	P _H	P _I	P _J	P _K	P _L	P _M	P _N	P _O	P _P	P _Q	P _R	P _S	P _T	P _U	P _V	P _W	P _X	P _Y	P _Z
1	T	0.03	0.03	0.03	0.03	0.03	0.04	0.03	0.03	0.15	0.03	0.03	0.03	0.03	0.03	0.03	0.03	0.03	0.03	0.03	0.09	0.04	0.04	0.03	0.03	0.05	0.03
2	I	0.04	0.05	0.04	0.04	0.04	0.04	0.06	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04
3	D	0.04	0.05	0.04	0.04	0.04	0.04	0.05	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.04	0.05	0.04	0.04	0.04	0.04	0.04	0.04	0.04
4	N	0.03	0.03	0.03	0.03	0.04	0.04	0.04	0.04	0.04	0.04	0.06	0.04	0.03	0.06	0.03	0.03	0.04	0.04	0.04	0.04	0.04	0.03	0.05	0.04	0.03	0.03
5	G	0.03	0.05	0.04	0.04	0.04	0.04	0.06	0.04	0.04	0.04	0.04	0.04	0.03	0.04	0.03	0.04	0.04	0.04	0.05	0.04	0.04	0.04	0.03	0.04	0.04	0.04

i	y	Y _A - P _A	Y _B - P _B	Y _C - P _C	Y _D - P _D	Y _E - P _E	Y _F - P _F	Y _G - P _G	Y _H - P _H	Y _I - P _I	Y _J - P _J	Y _K - P _K	Y _L - P _L	Y _M - P _M	Y _N - P _N	Y _O - P _O	Y _P - P _P	Y _Q - P _Q	Y _R - P _R	Y _S - P _S	Y _T - P _T	Y _U - P _U	Y _V - P _V	Y _W - P _W	Y _X - P _X	Y _Y - P _Y	Y _Z - P _Z
1	T	-0.03	-0.03	-0.03	-0.03	-0.03	-0.04	-0.03	-0.03	-0.15	-0.03	-0.03	-0.03	-0.03	-0.03	-0.03	-0.03	-0.03	-0.03	-0.03	0.91	-0.04	-0.04	-0.03	-0.03	-0.05	-0.03
2	I	-0.04	-0.05	-0.04	-0.04	-0.04	-0.04	-0.06	-0.04	0.96	-0.04	-0.04	-0.04	-0.04	-0.04	-0.04	-0.04	-0.04	-0.04	-0.04	-0.04	-0.04	-0.04	-0.04	-0.04	-0.04	-0.04
3	D	-0.04	-0.05	-0.04	0.96	-0.04	-0.04	-0.05	-0.04	-0.04	-0.04	-0.04	-0.04	-0.04	-0.04	-0.04	-0.04	-0.04	-0.04	-0.05	-0.04	-0.04	-0.04	-0.04	-0.04	-0.04	-0.04
4	N	-0.03	-0.03	-0.03	-0.03	-0.04	-0.04	-0.04	-0.04	-0.04	-0.04	-0.06	-0.04	-0.03	0.94	-0.03	-0.03	-0.03	-0.04	-0.04	-0.04	-0.04	-0.03	-0.05	-0.04	-0.03	-0.03
5	G	-0.03	-0.05	-0.04	-0.04	-0.04	-0.04	0.94	-0.04	-0.04	-0.04	-0.04	-0.04	-0.04	-0.03	-0.03	-0.04	-0.04	-0.05	-0.04	-0.04	-0.04	-0.04	-0.03	-0.04	-0.04	-0.04

Gradient Boosting for Classification

round 100

i	y	Y _A	Y _B	Y _C	Y _D	Y _E	Y _F	Y _G	Y _H	Y _I	Y _J	Y _K	Y _L	Y _M	Y _N	Y _O	Y _P	Y _Q	Y _R	Y _S	Y _T	Y _U	Y _V	Y _W	Y _X	Y _Y	Y _Z
1	T	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0
2	I	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
3	D	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
4	N	0	0	0	0	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0
5	G	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0

i	y	F _A	F _B	F _C	F _D	F _E	F _F	F _G	F _H	F _I	F _J	F _K	F _L	F _M	F _N	F _O	F _P	F _Q	F _R	F _S	F _T	F _U	F _V	F _W	F _X	F _Y	F _Z
1	T	-3.26	-2.7	-2.2	-2.22	-2.48	-0.31	-2.77	-1.19	2.77	0.1	-1.49	-1.02	-1.64	-0.8	-2.4	-3.57	-0.9	-2.45	-0.2	4.61	0.5	-0.71	-1.21	-0.24	0.49	-1.66
2	I	-1.64	-1.09	-2.29	-1.8	0.45	-0.43	2.14	-1.56	1.19	1.09	-1.5	-0.5	-3.64	-3.98	-0.39	-2.3	1.42	-0.59	0.27	-2.88	-1.96	-1.67	-4.38	-2.06	-2.95	-1.76
3	D	-2.45	0.18	-3.01	0.18	-2.79	-1.7	-2.21	0.43	-1.12	0.32	0.67	-2.16	-2.91	-2.76	-1.92	-3.04	-1.47	-0.48	-1.48	-1.25	-2.25	-3.23	-4.38	0.17	-2.95	-2.65
4	N	-3.95	-3.38	-0.22	-0.94	-1.33	-1.38	-1.22	-0.12	-2.33	-3.13	0.58	-0.65	-0.25	2.96	-2.84	-1.82	0.19	0.55	-1.22	-1.25	0.45	-1.8	0.11	-0.69	-1.6	-3.78
5	G	-3.14	-0.04	-2.37	-0.78	0.02	-2.68	2.6	-1.48	-1.93	0.42	-1.44	-1.45	-3.36	-3.98	-0.94	-3.42	1.84	1.44	0.62	-1.25	-1.33	-4.41	-4.71	-2.62	-2.15	-1.09

i	y	P _A	P _B	P _C	P _D	P _E	P _F	P _G	P _H	P _I	P _J	P _K	P _L	P _M	P _N	P _O	P _P	P _Q	P _R	P _S	P _T	P _U	P _V	P _W	P _X	P _Y	P _Z
1	T	0	0	0	0	0.01	0	0	0.13	0.01	0	0	0	0	0	0	0	0	0	0.01	0.79	0.01	0	0	0.01	0.01	0
2	I	0.01	0.01	0	0.01	0.06	0.02	0.32	0.01	0.12	0.11	0.01	0.02	0	0	0.03	0	0.16	0.02	0.05	0	0.01	0.01	0	0	0	0.01
3	D	0.01	0.11	0	0.11	0.01	0.02	0.01	0.14	0.03	0.12	0.17	0.01	0	0.01	0.01	0	0.02	0.05	0.02	0.03	0.01	0	0	0.11	0	0.01
4	N	0	0	0.02	0.01	0.01	0.01	0.01	0.03	0	0	0.05	0.02	0.02	0.59	0	0	0.04	0.05	0.01	0.01	0.05	0.01	0.03	0.02	0.01	0
5	G	0	0.03	0	0.01	0.03	0	0.42	0.01	0	0.05	0.01	0.01	0	0	0.01	0	0.19	0.13	0.06	0.01	0.01	0	0	0	0	0.01

i	y	Y _A - P _A	Y _B - P _B	Y _C - P _C	Y _D - P _D	Y _E - P _E	Y _F - P _F	Y _G - P _G	Y _H - P _H	Y _I - P _I	Y _J - P _J	Y _K - P _K	Y _L - P _L	Y _M - P _M	Y _N - P _N	Y _O - P _O	Y _P - P _P	Y _Q - P _Q	Y _R - P _R	Y _S - P _S	Y _T - P _T	Y _U - P _U	Y _V - P _V	Y _W - P _W	Y _X - P _X	Y _Y - P _Y	Y _Z - P _Z
1	T	-0	-0	-0	-0	-0	-0.01	-0	-0	-0.13	-0.01	-0	-0	-0	-0	-0	-0	-0	-0	-0	0.21	-0.01	-0	-0	-0.01	-0.01	-0
2	I	-0.01	-0.01	-0	-0.01	-0.06	-0.02	-0.32	-0.01	0.88	-0.11	-0.01	-0.02	-0	-0.03	-0	-0.03	-0	-0.16	-0.02	-0.05	-0	-0.01	-0.01	-0	-0	-0.01
3	D	-0.01	-0.11	-0	0.89	-0.01	-0.02	-0.01	-0.14	-0.03	-0.12	-0.17	-0.01	-0	-0.01	-0.01	-0	-0.02	-0.05	-0.02	-0.03	-0.01	-0	-0	-0.11	-0	-0.01
4	N	-0	-0	-0.02	-0.01	-0.01	-0.01	-0.01	-0.03	-0	-0	-0.05	-0.02	-0.02	0.41	-0	-0	-0.04	-0.05	-0.01	-0.01	-0.05	-0.01	-0.03	-0.02	-0.01	-0
5	G	-0	-0.03	-0	-0.01	-0.03	-0	0.58	-0.01	-0	-0.05	-0.01	-0.01	-0	-0	-0.01	-0	-0.19	-0.13	-0.06	-0.01	-0.01	-0	-0	-0	-0	-0.01