

Московский Государственный Университет имени М.В. Ломоносова  
Факультет Вычислительной Математики и Кибернетики  
Кафедра Математических Методов Прогнозирования

# Краткий обзор ключевых возможностей пакета “gbm” системы R

Фонарев Александр  
317 группа

Апрель 2012

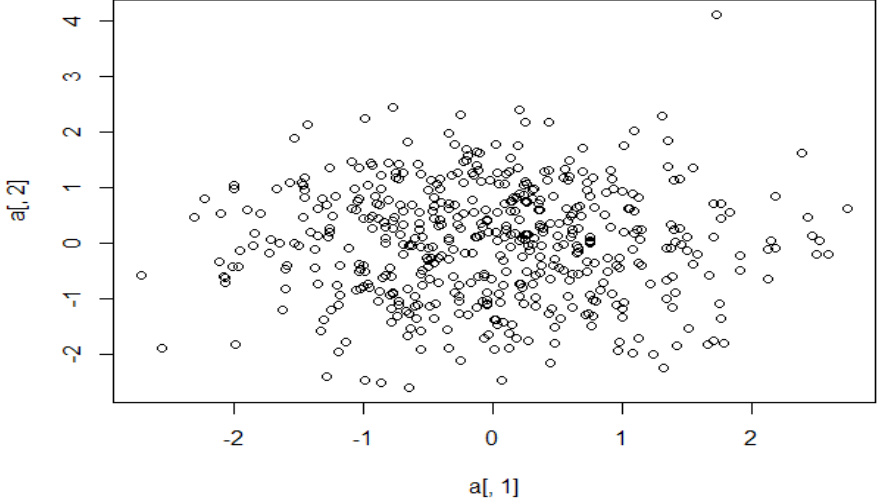
## Общее описание пакета “gbm” и проделанной работы

Полное название данного пакета – «Generalized Boosted Regression Models». Он включает в себя реализацию различных алгоритмов бустинга для решения задачи восстановления регрессии. Возможно использование различных функций потерь и других параметров.

Мною будут рассмотрены лишь несколько примеров, демонстрирующие основные возможности данного пакета. Далее будет рассмотрен только один вариант с построением композиции решающих деревьев. Для более детального изучения рекомендуется использовать официальное руководство к пакету (см. раздел «Использованные источники»).

## Генерация данных

Для начала сгенерируем модельные данные, на которых мы будем изучать пакет. Для этого создадим 500 точек из нормального распределения со стандартными параметрами и сохраним их координаты в переменную `a`. Далее для каждой точки посчитаем длину ее радиус-вектора и сохраним это значение в переменную `b`. Будем пробовать восстановить зависимость длины радиус-вектора от координат точки.

Код	Результат
<pre>a = matrix(rnorm(1000), 500, 2) plot(a[,1],a[,2]) b = sqrt(a[,1]^2 + a[,2]^2)</pre>	 <p>The figure shows a scatter plot of 500 data points. The x-axis is labeled 'a[, 1]' and ranges from approximately -3 to 3, with major ticks at -2, -1, 0, 1, and 2. The y-axis is labeled 'a[, 2]' and ranges from approximately -2 to 4, with major ticks at -2, -1, 0, 1, 2, 3, and 4. The points are distributed in a roughly circular cloud centered around the origin (0,0), representing a bivariate normal distribution with standard deviations of 1.</p>

## Создание модели

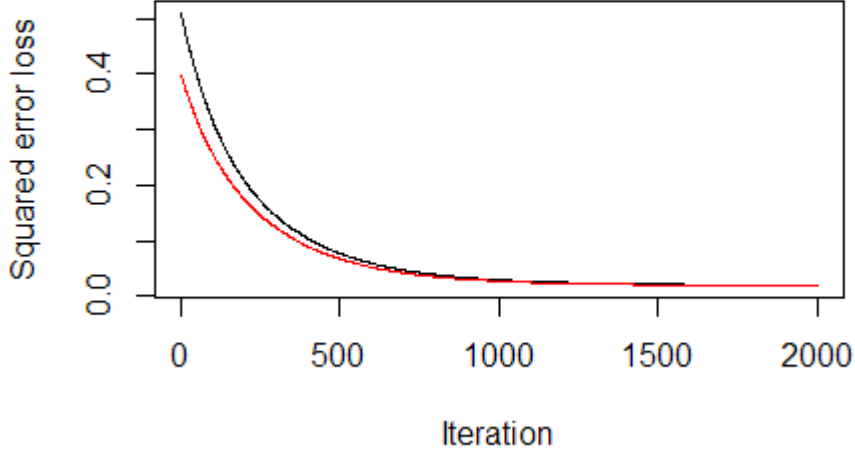
Пакет `gbm` предполагает обучение модели и сохранение ее в некоторый объект. Создадим объект-модель с названием `ourmodel`:

Код	Результат																																																																																																																																																											
<pre>ourmodel &lt;- gbm(  # формула, показывающая, что ищется зависимость поля Y от полей X1 и X2 в соответствующем фрейме Y ~ X1 + X2,  # создание фрейма из матрицы с данными data = data.frame(X1 = a[,1], X2 = a[,2], Y = b),  # выбор функции потерь distribution = "gaussian",  # максимальное количество используемых решающих деревьев n.trees = 2000,  # другие параметры алгоритма shrinkage = 0.005, interaction.depth=3,  # параметры валидации bag.fraction = 0.5, train.fraction = 0.5, n.minobsinnode = 10, cv.folds = 5,  # параметр, показывающий, надо ли сохранять копию данных keep.data = TRUE,  # параметр, показывающий, нужна ли печать в консоль результатов процесса настройки модели verbose=TRUE )</pre>	<p>Сначала идет вывод результатов для каждого из 5 этапов кросс-валидации отдельно. А затем общие результаты:</p> <table border="1"><thead><tr><th>Iter</th><th>TrainDeviance</th><th>ValidDeviance</th><th>StepSize</th><th>Improve</th></tr></thead><tbody><tr><td>1</td><td>0.5102</td><td>0.3978</td><td>0.0050</td><td>0.0021</td></tr><tr><td>2</td><td>0.5074</td><td>0.3960</td><td>0.0050</td><td>0.0025</td></tr><tr><td>3</td><td>0.5048</td><td>0.3942</td><td>0.0050</td><td>0.0022</td></tr><tr><td>4</td><td>0.5022</td><td>0.3925</td><td>0.0050</td><td>0.0022</td></tr><tr><td>5</td><td>0.4996</td><td>0.3906</td><td>0.0050</td><td>0.0026</td></tr><tr><td>6</td><td>0.4971</td><td>0.3889</td><td>0.0050</td><td>0.0024</td></tr><tr><td>7</td><td>0.4941</td><td>0.3868</td><td>0.0050</td><td>0.0025</td></tr><tr><td>8</td><td>0.4913</td><td>0.3848</td><td>0.0050</td><td>0.0023</td></tr><tr><td>9</td><td>0.4888</td><td>0.3829</td><td>0.0050</td><td>0.0025</td></tr><tr><td>10</td><td>0.4864</td><td>0.3814</td><td>0.0050</td><td>0.0022</td></tr><tr><td>100</td><td>0.3158</td><td>0.2569</td><td>0.0050</td><td>0.0013</td></tr><tr><td>200</td><td>0.2078</td><td>0.1742</td><td>0.0050</td><td>0.0007</td></tr><tr><td>300</td><td>0.1441</td><td>0.1233</td><td>0.0050</td><td>0.0005</td></tr><tr><td>400</td><td>0.1039</td><td>0.0893</td><td>0.0050</td><td>0.0002</td></tr><tr><td>500</td><td>0.0770</td><td>0.0672</td><td>0.0050</td><td>0.0002</td></tr><tr><td>600</td><td>0.0593</td><td>0.0522</td><td>0.0050</td><td>0.0001</td></tr><tr><td>700</td><td>0.0472</td><td>0.0422</td><td>0.0050</td><td>0.0001</td></tr><tr><td>800</td><td>0.0390</td><td>0.0351</td><td>0.0050</td><td>0.0000</td></tr><tr><td>900</td><td>0.0334</td><td>0.0304</td><td>0.0050</td><td>0.0000</td></tr><tr><td>1000</td><td>0.0296</td><td>0.0273</td><td>0.0050</td><td>0.0000</td></tr><tr><td>1100</td><td>0.0269</td><td>0.0248</td><td>0.0050</td><td>0.0000</td></tr><tr><td>1200</td><td>0.0251</td><td>0.0232</td><td>0.0050</td><td>-0.0000</td></tr><tr><td>1300</td><td>0.0236</td><td>0.0221</td><td>0.0050</td><td>-0.0000</td></tr><tr><td>1400</td><td>0.0226</td><td>0.0214</td><td>0.0050</td><td>0.0000</td></tr><tr><td>1500</td><td>0.0218</td><td>0.0207</td><td>0.0050</td><td>-0.0000</td></tr><tr><td>1600</td><td>0.0210</td><td>0.0203</td><td>0.0050</td><td>-0.0000</td></tr><tr><td>1700</td><td>0.0204</td><td>0.0199</td><td>0.0050</td><td>-0.0000</td></tr><tr><td>1800</td><td>0.0199</td><td>0.0197</td><td>0.0050</td><td>-0.0000</td></tr><tr><td>1900</td><td>0.0194</td><td>0.0194</td><td>0.0050</td><td>-0.0000</td></tr><tr><td>2000</td><td>0.0190</td><td>0.0191</td><td>0.0050</td><td>-0.0000</td></tr></tbody></table>	Iter	TrainDeviance	ValidDeviance	StepSize	Improve	1	0.5102	0.3978	0.0050	0.0021	2	0.5074	0.3960	0.0050	0.0025	3	0.5048	0.3942	0.0050	0.0022	4	0.5022	0.3925	0.0050	0.0022	5	0.4996	0.3906	0.0050	0.0026	6	0.4971	0.3889	0.0050	0.0024	7	0.4941	0.3868	0.0050	0.0025	8	0.4913	0.3848	0.0050	0.0023	9	0.4888	0.3829	0.0050	0.0025	10	0.4864	0.3814	0.0050	0.0022	100	0.3158	0.2569	0.0050	0.0013	200	0.2078	0.1742	0.0050	0.0007	300	0.1441	0.1233	0.0050	0.0005	400	0.1039	0.0893	0.0050	0.0002	500	0.0770	0.0672	0.0050	0.0002	600	0.0593	0.0522	0.0050	0.0001	700	0.0472	0.0422	0.0050	0.0001	800	0.0390	0.0351	0.0050	0.0000	900	0.0334	0.0304	0.0050	0.0000	1000	0.0296	0.0273	0.0050	0.0000	1100	0.0269	0.0248	0.0050	0.0000	1200	0.0251	0.0232	0.0050	-0.0000	1300	0.0236	0.0221	0.0050	-0.0000	1400	0.0226	0.0214	0.0050	0.0000	1500	0.0218	0.0207	0.0050	-0.0000	1600	0.0210	0.0203	0.0050	-0.0000	1700	0.0204	0.0199	0.0050	-0.0000	1800	0.0199	0.0197	0.0050	-0.0000	1900	0.0194	0.0194	0.0050	-0.0000	2000	0.0190	0.0191	0.0050	-0.0000
Iter	TrainDeviance	ValidDeviance	StepSize	Improve																																																																																																																																																								
1	0.5102	0.3978	0.0050	0.0021																																																																																																																																																								
2	0.5074	0.3960	0.0050	0.0025																																																																																																																																																								
3	0.5048	0.3942	0.0050	0.0022																																																																																																																																																								
4	0.5022	0.3925	0.0050	0.0022																																																																																																																																																								
5	0.4996	0.3906	0.0050	0.0026																																																																																																																																																								
6	0.4971	0.3889	0.0050	0.0024																																																																																																																																																								
7	0.4941	0.3868	0.0050	0.0025																																																																																																																																																								
8	0.4913	0.3848	0.0050	0.0023																																																																																																																																																								
9	0.4888	0.3829	0.0050	0.0025																																																																																																																																																								
10	0.4864	0.3814	0.0050	0.0022																																																																																																																																																								
100	0.3158	0.2569	0.0050	0.0013																																																																																																																																																								
200	0.2078	0.1742	0.0050	0.0007																																																																																																																																																								
300	0.1441	0.1233	0.0050	0.0005																																																																																																																																																								
400	0.1039	0.0893	0.0050	0.0002																																																																																																																																																								
500	0.0770	0.0672	0.0050	0.0002																																																																																																																																																								
600	0.0593	0.0522	0.0050	0.0001																																																																																																																																																								
700	0.0472	0.0422	0.0050	0.0001																																																																																																																																																								
800	0.0390	0.0351	0.0050	0.0000																																																																																																																																																								
900	0.0334	0.0304	0.0050	0.0000																																																																																																																																																								
1000	0.0296	0.0273	0.0050	0.0000																																																																																																																																																								
1100	0.0269	0.0248	0.0050	0.0000																																																																																																																																																								
1200	0.0251	0.0232	0.0050	-0.0000																																																																																																																																																								
1300	0.0236	0.0221	0.0050	-0.0000																																																																																																																																																								
1400	0.0226	0.0214	0.0050	0.0000																																																																																																																																																								
1500	0.0218	0.0207	0.0050	-0.0000																																																																																																																																																								
1600	0.0210	0.0203	0.0050	-0.0000																																																																																																																																																								
1700	0.0204	0.0199	0.0050	-0.0000																																																																																																																																																								
1800	0.0199	0.0197	0.0050	-0.0000																																																																																																																																																								
1900	0.0194	0.0194	0.0050	-0.0000																																																																																																																																																								
2000	0.0190	0.0191	0.0050	-0.0000																																																																																																																																																								

Это происходит с помощью вызова функции `gbm()`. Она содержит множество параметров, некоторые из которых даже опущены в этом примере. В качестве вывода можно видеть результат на каждой из 5 итераций кроссвалидации и его зависимость от количества деревьев, а также общий результат по обучению модели. Видно, что при увеличении числа деревьев начиная примерно с 1200 результат перестает улучшаться.

## Оптимальное количество алгоритмов

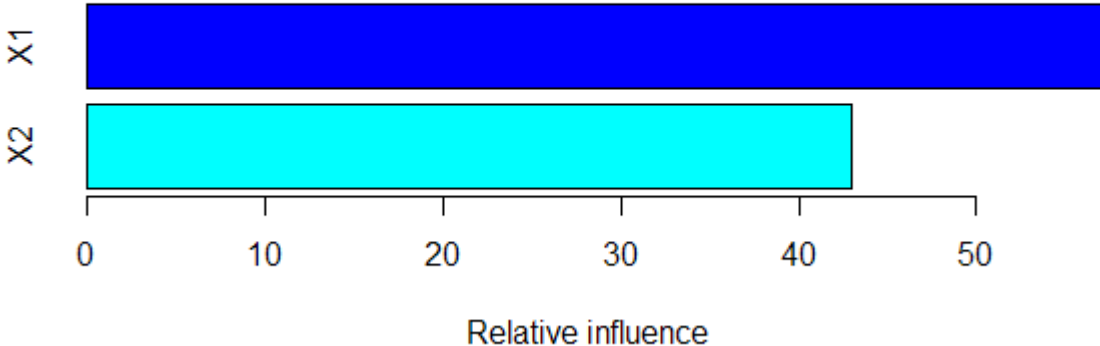
Посмотрим по-внимательнее на оптимальное количество деревьев после настройки модели:

Код	Результат
<pre data-bbox="91 544 842 678"># оптимальное число деревьев при оценке «out-of-bag» # также возможны варианты «test» и «cv»  best.iter &lt;- gbm.perf(ourmodel,method="OOB") print(best.iter)</pre>	 <p data-bbox="1131 951 1193 978">1263</p>

Красным показана ошибка на контроле в зависимости от количества деревьев, черным показана ошибка на обучении. Далее выводится оптимальное число деревьев в соответствии с данным методом проверки обучающей способности: 1263.

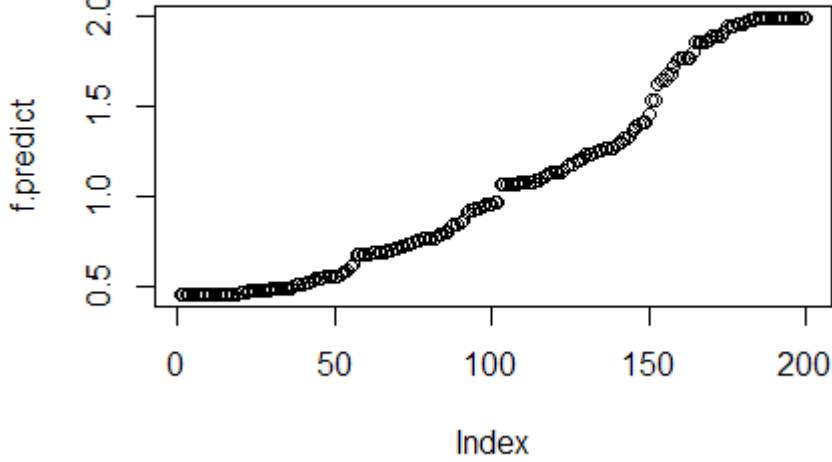
## Информация о каждом дереве в отдельности

Также в пакете присутствуют функции для просмотра информации и анализа деревьев и наборов деревьев по-отдельности:

Код	Результат																																																																																																			
<pre># также можно посмотреть на «влияние» каждого из двух признаков на предсказанный ответ при фиксированном числе деревьев  summary(ourmodel, n.trees = 1543)  # можно вывести фрейм с информацией о внутреннем строении дерева  print(pretty.gbm.tree(ourmodel, 1543))</pre>	 <p style="text-align: center;">Relative influence</p> <pre>var rel.inf 1 X1 57.06188 2 X2 42.93812</pre> <table border="1" data-bbox="786 962 2101 1257"> <thead> <tr> <th>SplitVar</th> <th>SplitCode</th> <th>Pred</th> <th>LeftNode</th> <th>RightNode</th> <th>MissingNode</th> <th>ErrorReduction</th> <th>Weight</th> <th>Prediction</th> </tr> </thead> <tbody> <tr><td>0</td><td>1</td><td>-1.741156e+00</td><td>1</td><td>2</td><td>9</td><td>0.04995455</td><td>125</td><td>3.364079e-05</td></tr> <tr><td>1</td><td>-1</td><td>3.726031e-04</td><td>-1</td><td>-1</td><td>-1</td><td>0.00000000</td><td>10</td><td>3.726031e-04</td></tr> <tr><td>2</td><td>1</td><td>5.936047e-01</td><td>3</td><td>7</td><td>8</td><td>0.06429663</td><td>115</td><td>4.165803e-06</td></tr> <tr><td>3</td><td>0</td><td>-2.090986e-01</td><td>4</td><td>5</td><td>6</td><td>0.06494628</td><td>82</td><td>-7.083488e-05</td></tr> <tr><td>4</td><td>-1</td><td>1.096515e-04</td><td>-1</td><td>-1</td><td>-1</td><td>0.00000000</td><td>31</td><td>1.096515e-04</td></tr> <tr><td>5</td><td>-1</td><td>-1.805423e-04</td><td>-1</td><td>-1</td><td>-1</td><td>0.00000000</td><td>51</td><td>-1.805423e-04</td></tr> <tr><td>6</td><td>-1</td><td>-7.083488e-05</td><td>-1</td><td>-1</td><td>-1</td><td>0.00000000</td><td>82</td><td>-7.083488e-05</td></tr> <tr><td>7</td><td>-1</td><td>1.905311e-04</td><td>-1</td><td>-1</td><td>-1</td><td>0.00000000</td><td>33</td><td>1.905311e-04</td></tr> <tr><td>8</td><td>-1</td><td>4.165803e-06</td><td>-1</td><td>-1</td><td>-1</td><td>0.00000000</td><td>115</td><td>4.165803e-06</td></tr> <tr><td>9</td><td>-1</td><td>3.364079e-05</td><td>-1</td><td>-1</td><td>-1</td><td>0.00000000</td><td>125</td><td>3.364079e-05</td></tr> </tbody> </table>	SplitVar	SplitCode	Pred	LeftNode	RightNode	MissingNode	ErrorReduction	Weight	Prediction	0	1	-1.741156e+00	1	2	9	0.04995455	125	3.364079e-05	1	-1	3.726031e-04	-1	-1	-1	0.00000000	10	3.726031e-04	2	1	5.936047e-01	3	7	8	0.06429663	115	4.165803e-06	3	0	-2.090986e-01	4	5	6	0.06494628	82	-7.083488e-05	4	-1	1.096515e-04	-1	-1	-1	0.00000000	31	1.096515e-04	5	-1	-1.805423e-04	-1	-1	-1	0.00000000	51	-1.805423e-04	6	-1	-7.083488e-05	-1	-1	-1	0.00000000	82	-7.083488e-05	7	-1	1.905311e-04	-1	-1	-1	0.00000000	33	1.905311e-04	8	-1	4.165803e-06	-1	-1	-1	0.00000000	115	4.165803e-06	9	-1	3.364079e-05	-1	-1	-1	0.00000000	125	3.364079e-05
SplitVar	SplitCode	Pred	LeftNode	RightNode	MissingNode	ErrorReduction	Weight	Prediction																																																																																												
0	1	-1.741156e+00	1	2	9	0.04995455	125	3.364079e-05																																																																																												
1	-1	3.726031e-04	-1	-1	-1	0.00000000	10	3.726031e-04																																																																																												
2	1	5.936047e-01	3	7	8	0.06429663	115	4.165803e-06																																																																																												
3	0	-2.090986e-01	4	5	6	0.06494628	82	-7.083488e-05																																																																																												
4	-1	1.096515e-04	-1	-1	-1	0.00000000	31	1.096515e-04																																																																																												
5	-1	-1.805423e-04	-1	-1	-1	0.00000000	51	-1.805423e-04																																																																																												
6	-1	-7.083488e-05	-1	-1	-1	0.00000000	82	-7.083488e-05																																																																																												
7	-1	1.905311e-04	-1	-1	-1	0.00000000	33	1.905311e-04																																																																																												
8	-1	4.165803e-06	-1	-1	-1	0.00000000	115	4.165803e-06																																																																																												
9	-1	3.364079e-05	-1	-1	-1	0.00000000	125	3.364079e-05																																																																																												

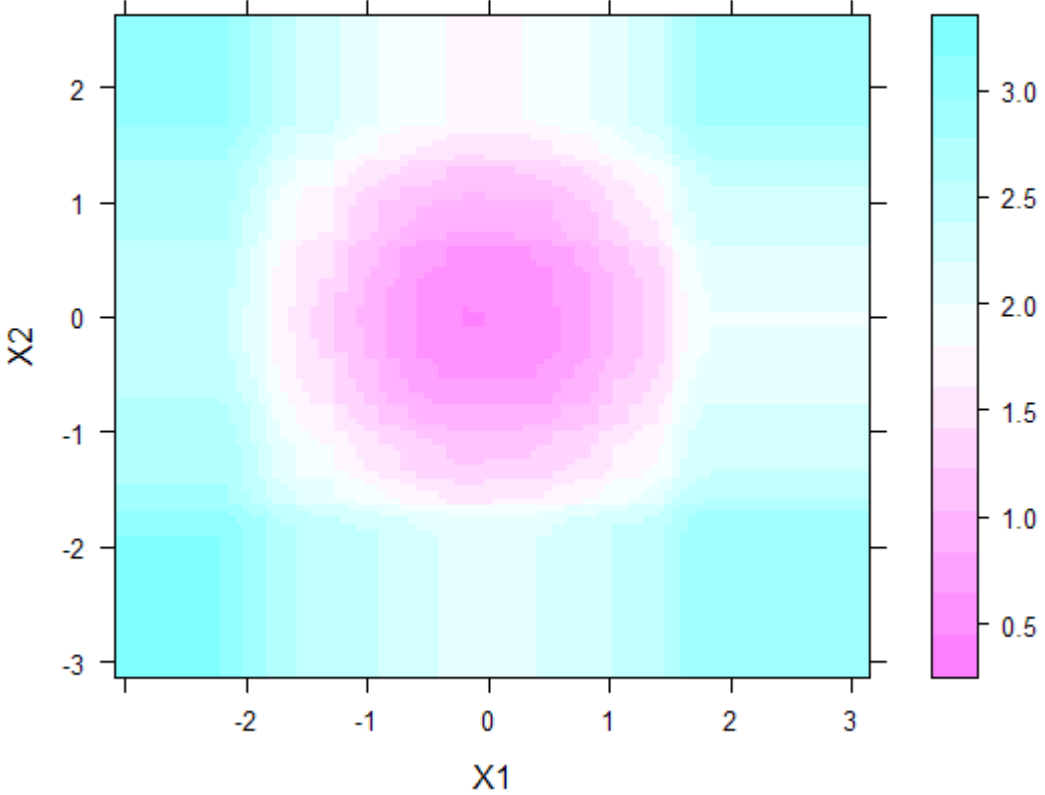
## Использование настроенной модели для предсказания

Сгенерируем набор точек на оси абсцисс с шагом 0.01 (их радиус-вектор просто равен абсциссе) и попробуем предсказать радиус-вектора для точек этого набора:

Код	Результат
<pre># генерируем данные a2 = cbind(1:200 / 100, 0)  # предсказываем, используя модель и оптимальное количество алгоритмов f.predict &lt;- predict.gbm(ourmodel, data.frame(X1 = a2[,1], X2 = a2[,2], Y = a2[,1]), best.iter)  # изобразим наш предсказанный результат plot(f.predict)</pre>	

Идеальный результат – линейная зависимость на графике. Видно, что нашему алгоритму ее найти не удалось. Однако все же он уже нашел достаточно неплохую закономерность в данных.

## Визуализация карты восстановленной регрессии

Код	Результат
<pre data-bbox="91 564 987 671"># визуализируем то, как восстановилась регрессия # в качестве координат используем признаки с первого по второй plot.gbm(ourmodel, 1:2, best.iter)</pre>	 <p>The figure is a heatmap representing the fitted regression surface. The horizontal axis is labeled 'X1' and has major tick marks at -2, -1, 0, 1, 2, and 3. The vertical axis is labeled 'X2' and has major tick marks at -3, -2, -1, 0, 1, and 2. The color scale on the right side of the plot ranges from 0.5 (dark pink) at the bottom to 3.0 (cyan) at the top, with intermediate values at 1.0, 1.5, 2.0, and 2.5. The heatmap shows a central region of high values (pink) centered around X1=0 and X2=0, which transitions through lighter colors to cyan as it moves away from the center. The plot is enclosed in a black frame.</p>



## Использованные источники

- Официальный сайт пакета – <http://cran.gis-lab.info/web/packages/gbm/index.html>
- Руководство к пакету “gbm” – <http://cran.gis-lab.info/web/packages/gbm/gbm.pdf>
- Методическое пособие А.Г. Дьяконова по системе R – <http://alexanderdyakonov.narod.ru/upR.pdf>
- Курс лекций К.В. Воронцова – [http://www.machinelearning.ru/wiki/index.php?title=Машинное\\_обучение\\_%28курс\\_лекций%2C\\_К.В.Воронцов%29](http://www.machinelearning.ru/wiki/index.php?title=Машинное_обучение_%28курс_лекций%2C_К.В.Воронцов%29)