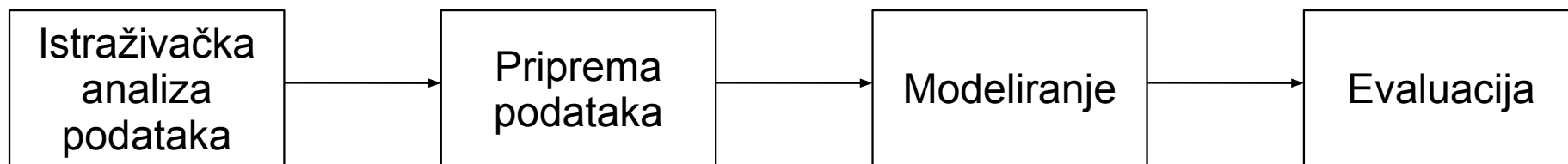# Hadoop implementacija algoritama za klasifikaciju podataka

Miloš Veljković 1174

1

# Algoritmi za klasifikaciju

- KNN (K-nearest neighbors)

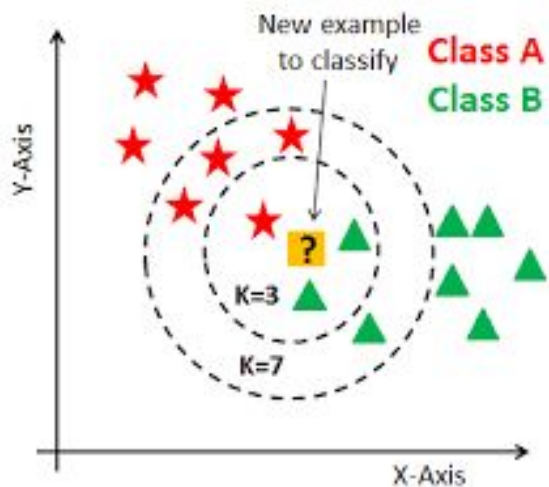- Naive-Bayes

```
┌─────────────┐     ┌─────────────┐     ┌─────────────┐     ┌─────────────┐
│ Istraživačka│     │             │     │             │     │             │
│   analiza   │ ──▶ │  Priprema   │ ──▶ │ Modeliranje │ ──▶ │  Evaluacija │
│  podataka   │     │  podataka   │     │             │     │             │
└─────────────┘     └─────────────┘     └─────────────┘     └─────────────┘
```

Faze klasifikacije

**2**

# Koji problem rešavamo?

- Određivanje da li osoba ima srčanu bolest ili ne !?

| | age | sex | cp | trestbps | chol | fbs | restecg | thalach | exang | oldpeak | slope | ca | thal | target |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 2 | 63 | 1 | 3 | 145 | 233 | 1 | 0 | 150 | 0 | 2.3 | 0 | 0 | 1 | 1 |
| 3 | 37 | 1 | 2 | 130 | 250 | 0 | 1 | 187 | 0 | 3.5 | 0 | 0 | 2 | 1 |
| 4 | 41 | 0 | 1 | 130 | 204 | 0 | 0 | 172 | 0 | 1.4 | 2 | 0 | 2 | 1 |
| 5 | 56 | 1 | 1 | 120 | 236 | 0 | 1 | 178 | 0 | 0.8 | 2 | 0 | 2 | 1 |
| 6 | 57 | 0 | 0 | 120 | 354 | 0 | 1 | 163 | 1 | 0.6 | 2 | 0 | 2 | 1 |
| 7 | 57 | 1 | 0 | 140 | 192 | 0 | 1 | 148 | 0 | 0.4 | 1 | 0 | 1 | 1 |
| 8 | 56 | 0 | 1 | 140 | 294 | 0 | 0 | 153 | 0 | 1.3 | 1 | 0 | 2 | 1 |
| 9 | 44 | 1 | 1 | 120 | 263 | 0 | 1 | 173 | 0 | 0 | 2 | 0 | 3 | 1 |
| 10 | 52 | 1 | 2 | 172 | 199 | 1 | 1 | 162 | 0 | 0.5 | 2 | 0 | 3 | 1 |
| 11 | 57 | 1 | 2 | 150 | 168 | 0 | 1 | 174 | 0 | 1.6 | 2 | 0 | 2 | 1 |
| 12 | 54 | 1 | 0 | 140 | 239 | 0 | 1 | 160 | 0 | 1.2 | 2 | 0 | 2 | 1 |
| 13 | 48 | 0 | 2 | 130 | 275 | 0 | 1 | 139 | 0 | 0.2 | 2 | 0 | 2 | 1 |
| 14 | 49 | 1 | 1 | 130 | 266 | 0 | 1 | 171 | 0 | 0.6 | 2 | 0 | 2 | 1 |
| 15 | 64 | 1 | 3 | 110 | 211 | 0 | 0 | 144 | 1 | 1.8 | 1 | 0 | 2 | 1 |
| 16 | 58 | 0 | 3 | 150 | 283 | 1 | 0 | 162 | 0 | 1 | 2 | 0 | 2 | 1 |
| 17 | 50 | 0 | 2 | 120 | 219 | 0 | 1 | 158 | 0 | 1.6 | 1 | 0 | 2 | 1 |
| 18 | 58 | 0 | 2 | 120 | 340 | 0 | 1 | 172 | 0 | 0 | 2 | 0 | 2 | 1 |

# KNN (K-nearest neighbors) ?



TEST

| 3 | 2 | ? |

TRAINING

| 1 | 2 | 1 |
| 3 | 6 | 0 |
| 2 | 4 | 1 |
| 4 | 5 | 0 |

New example to classify — Class A / Class B

$$d(x, y) = \sqrt{\sum_{i=1}^{n}(x_i - y_i)^2}$$

| 2 | 1 |
| 2.23 | 1 |
| 3,16 | 0 |
| 4 | 0 |

K=3

4

# Naive-Bayes ?

- Bazira se na verovatnoći i statistici

- Kreće od pretpostavke da su svi podaci podjednako važni

- Predviđen za rad sa diskretnim podacima

- Kontinualni podaci podrazumevaju primenu Gausove normalne raspodele

| age | sex | cp | trestbps | chol | fbs | restecg | thalach | exang | oldpeak | slope | ca | thal | target |
|-----|-----|----|----------|------|-----|---------|---------|-------|---------|-------|----|------|--------|
| 41 | 1 | 1 | 120 | 157 | 0 | 1 | 182 | 0 | 0 | 2 | 0 | 2 | 1 |
| 38 | 1 | 2 | 138 | 175 | 0 | 1 | 173 | 0 | 0 | 2 | 4 | 2 | 1 |
| 38 | 1 | 2 | 138 | 175 | 0 | 1 | 173 | 0 | 0 | 2 | 4 | 2 | 1 |
| 67 | 1 | 0 | 160 | 286 | 0 | 0 | 108 | 1 | 1.5 | 1 | 3 | 2 | 0 |
| 67 | 1 | 0 | 120 | 229 | 0 | 0 | 129 | 1 | 2.6 | 1 | 2 | 3 | 1 |
| 62 | 0 | 0 | 140 | 268 | 0 | 0 | 160 | 0 | 3.6 | 0 | 2 | 2 | 0 |
| 63 | 1 | 0 | 130 | 254 | 0 | 0 | 147 | 0 | 1.4 | 1 | 1 | 3 | 0 |
| 53 | 1 | 0 | 140 | 203 | 1 | 0 | 155 | 1 | 3.1 | 0 | 0 | 3 | 0 |
| 56 | 1 | 2 | 130 | 256 | 1 | 0 | 142 | 1 | 0.6 | 1 | 1 | 1 | 0 |
| 48 | 1 | 1 | 110 | 229 | 0 | 1 | 168 | 0 | 1 | 0 | 0 | 3 | 0 |

| | | target | |
|----|---|---|---|
| | | 1 | 0 |
| cp | 0 | 1 | 4 |
| | 1 | 1 | 1 |
| | 2 | 2 | 1 |

1/4          4/6

**5**

# Kontinualni podaci ?
# Gausova normalna raspodela !

| age | sex | cp | trestbps | chol | fbs | restecg | thalach | exang | oldpeak | slope | ca | thal | target |
|-----|-----|----|----------|------|-----|---------|---------|-------|---------|-------|----|------|--------|
| 41 | 1 | 1 | 120 | 157 | 0 | 1 | 182 | 0 | 0 | 2 | 0 | 2 | 1 |
| 38 | 1 | 2 | 138 | 175 | 0 | 1 | 173 | 0 | 0 | 2 | 4 | 2 | 1 |
| 38 | 1 | 2 | 138 | 175 | 0 | 1 | 173 | 0 | 0 | 2 | 4 | 2 | 1 |
| 67 | 1 | 0 | 160 | 286 | 0 | 0 | 108 | 1 | 1.5 | 1 | 3 | 2 | 0 |
| 67 | 1 | 0 | 120 | 229 | 0 | 0 | 129 | 1 | 2.6 | 1 | 2 | 3 | 1 |
| 62 | 0 | 0 | 140 | 268 | 0 | 0 | 160 | 0 | 3.6 | 0 | 2 | 2 | 0 |
| 63 | 1 | 0 | 130 | 254 | 0 | 0 | 147 | 0 | 1.4 | 1 | 1 | 3 | 0 |
| 53 | 1 | 0 | 140 | 203 | 1 | 0 | 155 | 1 | 3.1 | 0 | 0 | 3 | 0 |
| 56 | 1 | 2 | 130 | 256 | 1 | 0 | 142 | 1 | 0.6 | 1 | 1 | 1 | 0 |
| 48 | 1 | 1 | 110 | 229 | 0 | 1 | 168 | 0 | 1 | 0 | 0 | 3 | 0 |

$$\mu = \frac{1}{n}\sum_{i=1}^{n} x_i \qquad \text{Mean}$$

$$\sigma = \left[\frac{1}{n-1}\sum_{i=1}^{n}(x_i-\mu)^2\right]^{0.5} \qquad \text{Standard deviation}$$

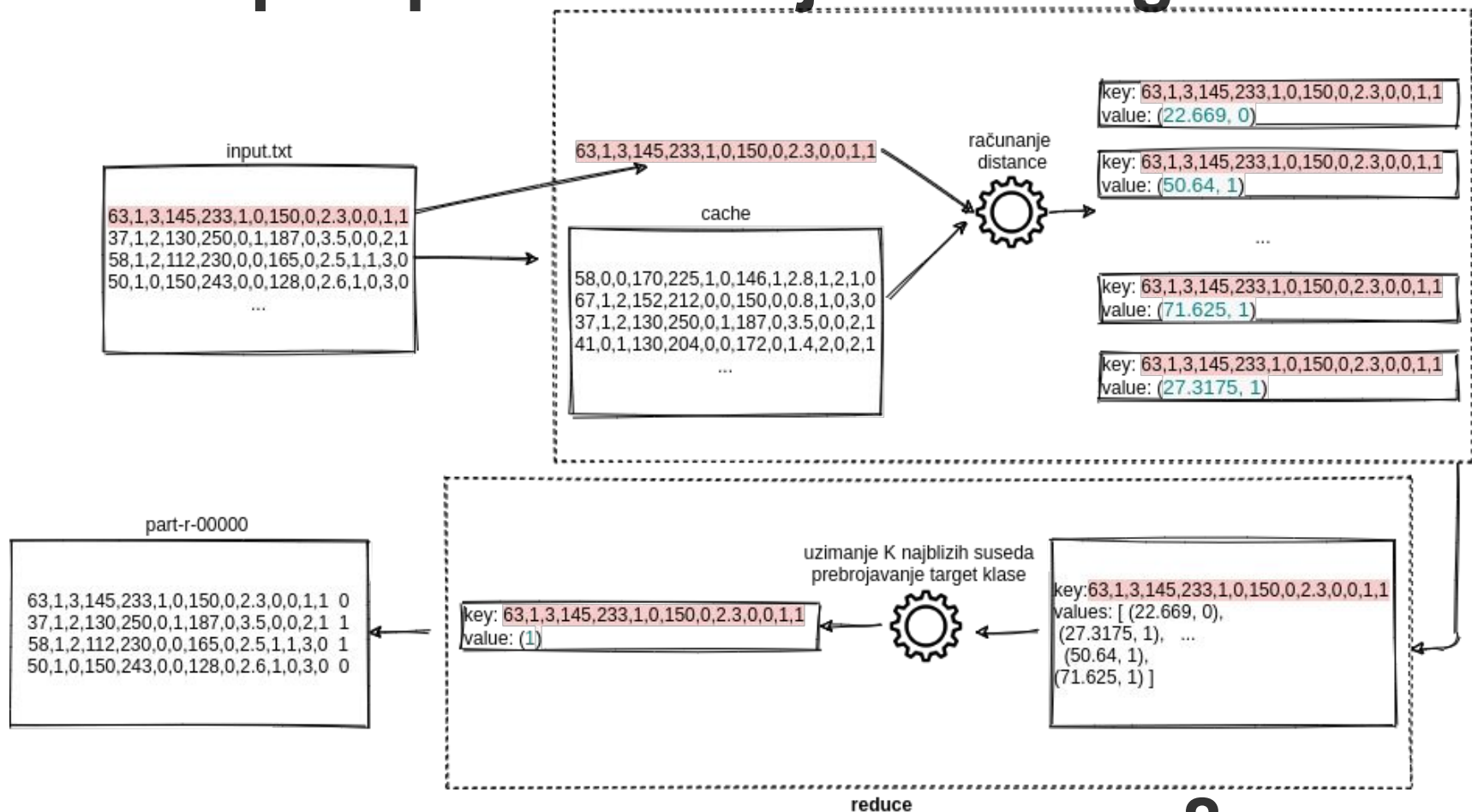$$f(x) = \frac{1}{\sqrt{2\pi}\sigma}e^{-\frac{(x-\mu)^2}{2\sigma^2}} \qquad \text{Normal distribution}$$

**6**

# Kako odrediti pripadnos klasi kod Naive-Bayes?

TEST

| age | sex | cp | trestbps | chol | fbs | restecg | thalach | exang | oldpeak | slope | ca | thal | target |
|-----|-----|----|----------|------|-----|---------|---------|-------|---------|-------|----|------|--------|
| 41 | 1 | 1 | 120 | 157 | 0 | 1 | 182 | 0 | 0 | 2 | 0 | 2 | |
| 38 | 1 | 2 | 138 | 175 | 0 | 1 | 173 | 0 | 0 | 2 | 4 | 2 | |
| 38 | 1 | 2 | 138 | 175 | 0 | 1 | 173 | 0 | 0 | 2 | 4 | 2 | |
| 67 | 1 | 0 | 160 | 286 | 0 | 0 | 108 | 1 | 1.5 | 1 | 3 | 2 | |
| 67 | 1 | 0 | 120 | 229 | 0 | 0 | 129 | 1 | 2.6 | 1 | 2 | 3 | |
| 62 | 0 | 0 | 140 | 268 | 0 | 0 | 160 | 0 | 3.6 | 0 | 2 | 2 | ? |
| 63 | 1 | 0 | 130 | 254 | 0 | 0 | 147 | 0 | 1.4 | 1 | 1 | 3 | |
| 53 | 1 | 0 | 140 | 203 | 1 | 0 | 155 | 1 | 3.1 | 0 | 0 | 3 | |
| 56 | 1 | 2 | 130 | 256 | 1 | 0 | 142 | 1 | 0.6 | 1 | 1 | 1 | |
| 48 | 1 | 1 | 110 | 229 | 0 | 1 | 168 | 0 | 1 | 0 | 0 | 3 | |

$$\mu = \frac{1}{n}\sum_{i=1}^{n} x_i \quad \text{Mean}$$

$$\sigma = \left[\frac{1}{n-1}\sum_{i=1}^{n}(x_i - \mu)^2\right]^{0.5} \quad \text{Standard deviation}$$

$$f(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}} \quad \text{Normal distribution}$$

| | | target | |
|---|---|---|---|
| | | 1 | 0 |
| cp | 0 | 1 | 4 |
| | 1 | 1 | 1 |
| | 2 | 2 | 1 |

1/4        1/6

P(1)

P(0)

P(1)>P(0) - target=1

P(0)>P(1) - target=0

7

# Hadoop implementacija KNN algoritma



**8**

# Mapper

```java
9    public class MapClass extends Mapper<LongWritable, Text, Text, DistanceTarget>{
10
11       private RecordsArray trainingSet;
12
13       @Override
14       protected void setup(Context context) throws IOException, InterruptedException {
15           trainingSet = new RecordsArray();
16           trainingSet.populate(new File("./KNN/cache/TrainingRecords.txt"));
17       }

19    @Override
20     public void map(LongWritable key, Text value, Context context) throws IOException, InterruptedException {
21
22         String line = value.toString();
23         String[] data = line.split(",");
24         int age = Integer.parseInt(data[0]);
25         int sex = Integer.parseInt(data[1]);
26         int cp = Integer.parseInt(data[2]);
27         int trestbps = Integer.parseInt(data[3]);
28         int chol = Integer.parseInt(data[4]);
29         int fbs = Integer.parseInt(data[5]);
30         int restecg = Integer.parseInt(data[6]);
31         int thalach = Integer.parseInt(data[7]);
32         int exang = Integer.parseInt(data[8]);
33         double oldpeak = Double.parseDouble(data[9]);
34         int slope = Integer.parseInt(data[10]);
35         int ca = Integer.parseInt(data[11]);
36         int thal = Integer.parseInt(data[12]);
37         int target = Integer.parseInt(data[13]);
38
39         Record testRecord = new Record( age,  sex,  cp,  trestbps,
40                 chol,  fbs,  restecg,  thalach,  exang,
41                 oldpeak, slope,  ca,  thal,  target);
42
43         for(Record r: trainingSet.records){
44             double distance = r.calculateEuclideanDistance(testRecord);
45             DistanceTarget dt = new DistanceTarget(new DoubleWritable(distance), new IntWritable(r.target));
46             Text outputKey = new Text();
47             outputKey.set(line);
48             context.write(outputKey, dt);
49         }
50     }
51 }
```

# Reducer

```java
8   public class ReduceClass extends Reducer<Text, DistanceTarget, Text, IntWritable>{
9
10
11      protected void reduce(Text key, Iterable<DistanceTarget> values, Context context)
12              throws IOException, InterruptedException {
13          int k = 5;
14          int class1=0, class0=0;
15          TreeMap<Double, Integer> currKnnMap = new TreeMap<Double, Integer>();
16          for(DistanceTarget val: values){
17              int target = val.getTarget().get();
18              double distance = val.getDistance().get();
19              currKnnMap.put(distance, target);
20          }
21          for(int i=0; i<k; i++) {
22              int target = currKnnMap.pollFirstEntry().getValue();
23              if(target == 1){ //has value
24                  class1++;
25              }else {
26                  class0++;
27              }
28          }
29
30          if(class1 > class0){
31              context.write(key, new IntWritable(1));
32          }else {
33              context.write(key, new IntWritable(0));
34          }
35
36      }
37  }
```
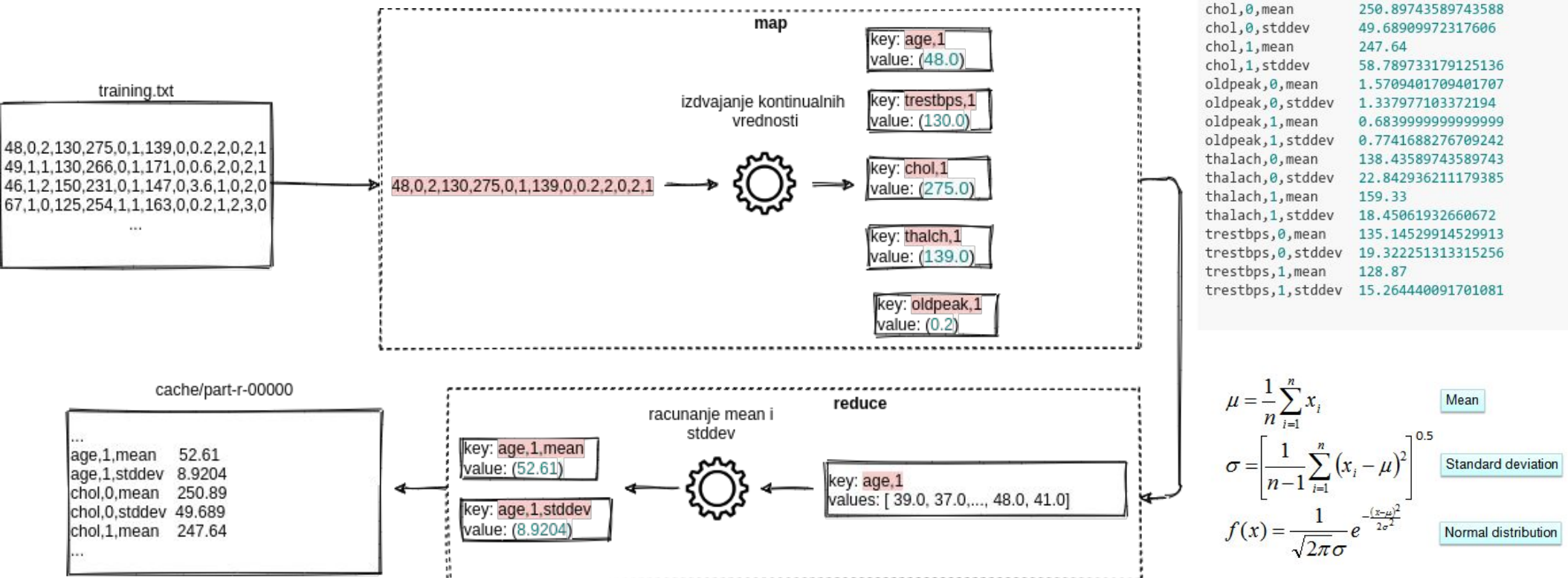
**10**

# Hadoop implementacija Naive-Bayes algoritma

- 3 faze

  - Faza 1: Prepocesiranje-Računanje standardne devijacije i srednje vrednosti

  - Faza 2: Kreiranja frekvencione tabele za diskretne podatke

  - Faza 3: Određivanje target klase

**11**

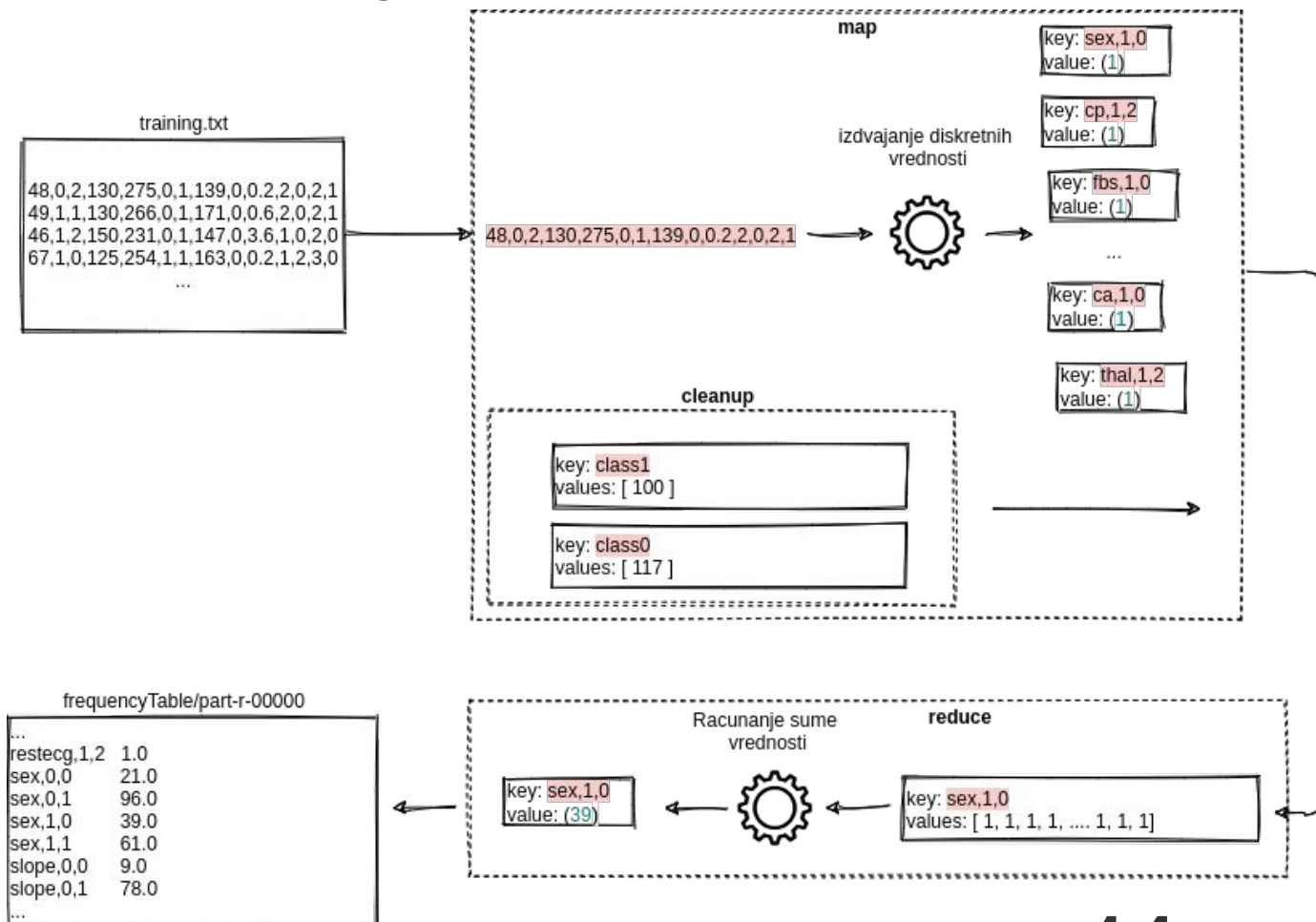# Faza 1: Računanje standardne devijacije i srednje vrednosti

# Faza 1: Hadoop implementacija

```java
 9  public class MapClass extends Mapper<LongWritable, Text, Text, DoubleWritable>{
10
11      @Override
12       public void map(LongWritable key, Text value, Context context) throws IOException, InterruptedException {
13
14          String line = value.toString();
15          String[] data = line.split(",");
16          for (int i = 0; i < 10; i++) {
17              if(i == 0){
18                  contextWrite("age,"+data[13].trim(), Double.parseDouble(data[0]), context);
19              } else if (i == 3) {
20                  contextWrite("trestbps,"+data[13].trim(), Double.parseDouble(data[3]), context);
21              } else if (i == 4) {
22                  contextWrite("chol,"+data[13].trim(), Double.parseDouble(data[4]), context);
23              } else if (i == 7) {
24                  contextWrite("thalach,"+data[13].trim(), Double.parseDouble(data[7]), context);
25              } else if (i == 9) {
26                  contextWrite("oldpeak,"+data[13].trim(), Double.parseDouble(data[9]), context);
27              }
28          }
29      }
30
31      public void contextWrite(String key, Double value, Context context){
32          try {
33              Text outputKey = new Text();
34              outputKey.set(key);
35              context.write(outputKey, new DoubleWritable(Double.valueOf(value)));
36          } catch (InterruptedException e) {
37              e.printStackTrace();
38          } catch (IOException e) {
39              e.printStackTrace();
40          }
41
42      }
43  }
```

```java
 7  public class ReduceClass extends Reducer<Text, DoubleWritable, Text, DoubleWritable>{
 8
 9      protected void reduce(Text key, Iterable<DoubleWritable> values, Context context)
10          throws IOException, InterruptedException {
11          double mean = 0, stddev = 0, sum = 0, stddevSum = 0;
12          ArrayList<Double> curValues = new ArrayList<>();
13          int counter = 0;
14          for(DoubleWritable val: values){
15              double value = val.get();
16              curValues.add(value);
17              sum += value;
18              counter++;
19          }
20          mean = sum / counter ;
21          for(Double val: curValues) {
22              double value = val.doubleValue();
23              double squareAddtions = Math.pow((value - mean), 2);
24              stddevSum += squareAddtions;
25          }
26          stddev = Math.pow( stddevSum/(counter-1), 0.5);
27          Text outputKeyMean = new Text(key+",mean");
28          Text outputKeyStddev = new Text(key+",stddev");
29          context.write(outputKeyMean, new DoubleWritable(mean));
30          context.write(outputKeyStddev, new DoubleWritable(stddev));
```

**13**

# Faza 2: Kreiranje frekvencione tabele



**14**

# Faza 2: Hadoop implementacija

Part1

Part2

```java
7   public class MapClassFT extends Mapper<LongWritable, Text, Text, IntWritable>{
8
9       private int class1 = 0, class0 = 0;
10
11      @Override
12      public void map(LongWritable key, Text value, Context context) throws IOException, Inte
13
14          String line = value.toString();
15          String[] data = line.split(",");
16          int target = Integer.parseInt(data[13]);
17          if(target == 1) {
18              class1++;
19          } else {
20              class0++;
21          }
22          for (int i = 0; i < 14; i++) {
23              if(i == 1){
24                  contextWrite("sex,"+data[13].trim()+","+data[1].trim(), 1, context);
25              } else if (i == 2) {
26                  contextWrite("cp,"+data[13].trim()+","+data[2].trim(), 1, context);
27              } else if (i == 5) {
28                  contextWrite("fbs,"+data[13].trim()+","+data[5].trim(), 1, context);
29              } else if (i == 6) {
30                  contextWrite("restecg,"+data[13].trim()+","+data[6].trim(), 1, context);
31              } else if (i == 8) {
32                  contextWrite("exang,"+data[13].trim()+","+data[8].trim(), 1, context);
33              } else if (i == 10) {
34                  contextWrite("slope,"+data[13].trim()+","+data[10].trim(), 1, context);
35              } else if (i == 11) {
36                  contextWrite("ca,"+data[13].trim()+","+data[11].trim(), 1, context);
37              } else if (i == 12) {
38                  contextWrite("thal,"+data[13].trim()+","+data[12].trim(), 1, context);
39              }
40          }
```
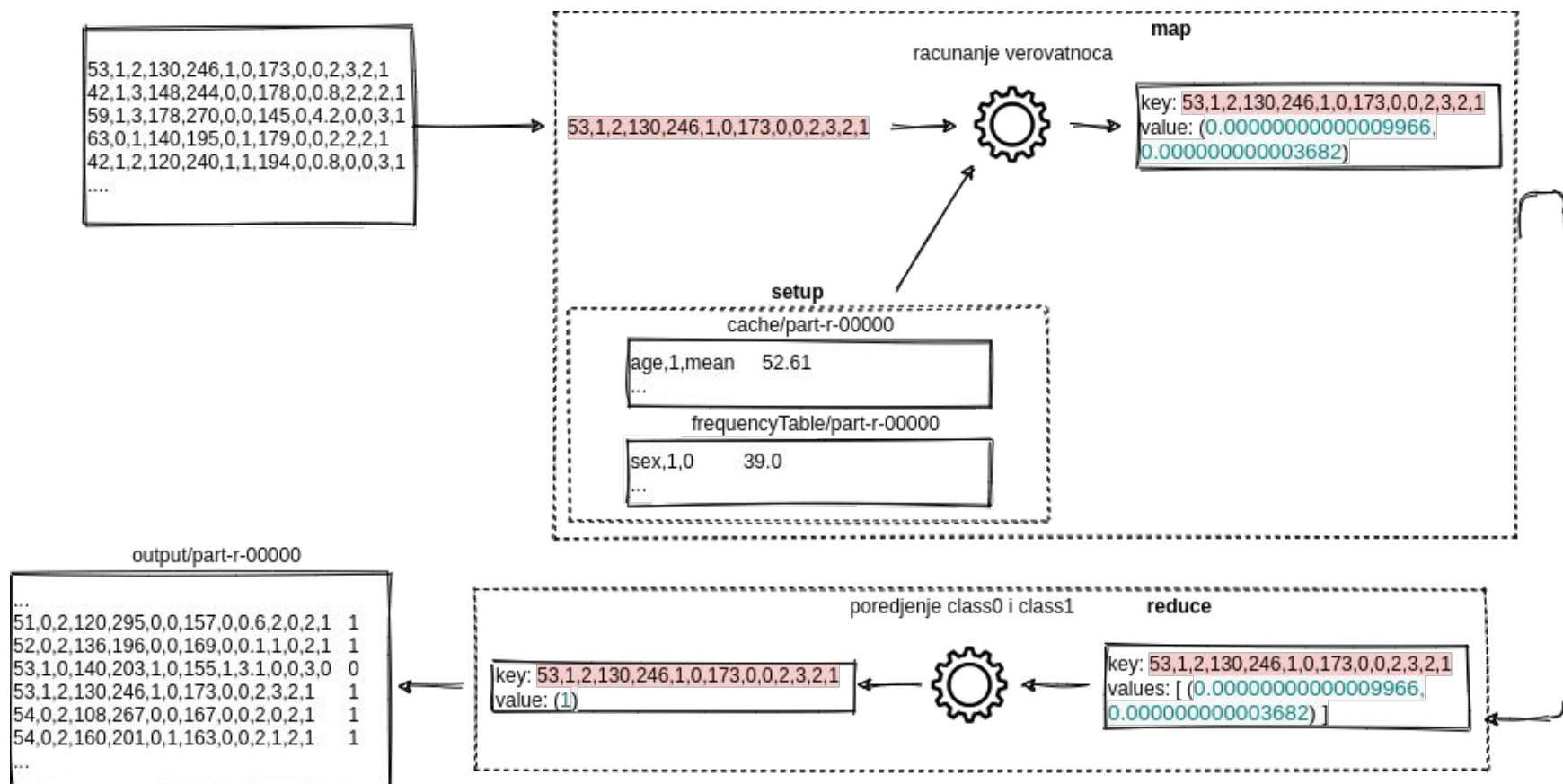
```java
37              } else if (i == 12) {
38                  contextWrite("thal,"+data[13].trim()+","+data[12].trim(), 1, context);
39              }
40          }
41      }
42
43      public void contextWrite(String key, Integer value, Context context){
44          try {
45              Text outputKey = new Text();
46              outputKey.set(key);
47              context.write(outputKey, new IntWritable(value));
48          } catch (InterruptedException e) {
49              e.printStackTrace();
50          } catch (IOException e) {
51              e.printStackTrace();
52          }
53      }
54
55      @Override
56      protected void cleanup(Mapper.Context context) throws IOException, InterruptedException
57          contextWrite("class1", class1, context);
58          contextWrite("class0", class0, context);
59      }
60  }
```

**15**

# Faza 2: Reducer

```java
7   public class ReduceClassFT extends Reducer<Text, IntWritable, Text, DoubleWritable> {
8
9       protected void reduce(Text key, Iterable<IntWritable> values, Context context)
10              throws IOException, InterruptedException {
11          if(key.toString().equals("class1") || key.toString().equals("class1") ){
12              int value = 0;
13              for(IntWritable val: values){
14                  value = val.get();
15              }
16              context.write(key, new DoubleWritable(value));
17          }else {
18              int sum = 0;
19              for(IntWritable val: values){
20                  int value = val.get();
21                  sum += value;
22              }
23              context.write(key, new DoubleWritable(sum));
24          }
25      }
26  }
```

# Faza 3: Određivanje target klase

# Faza 3: Hadoop implementacija

## Part1

```java
public class MapClassFIT extends Mapper<LongWritable, Text, Text, Estimation>{

    private HashMap<String, Double> preProcessingTable, frequencyTable;

    @Override
    protected void setup(Mapper.Context context) throws IOException, InterruptedException {
        BufferedReader cacheReader = new BufferedReader(new FileReader("./NB/cache/part-r-00000"));
        BufferedReader ftReader = new BufferedReader(new FileReader("./NB/frequencyTable/part-r-00000"));
        preProcessingTable = new HashMap<>();
        frequencyTable = new HashMap<>();
        String line;
        try {
            while((line = cacheReader.readLine()) != null) {
                String[] data = line.split("\t");
                preProcessingTable.put(data[0], Double.parseDouble(data[1]));
            }

            while((line = ftReader.readLine()) != null) {
                String[] data = line.split("\t");
                frequencyTable.put(data[0], Double.parseDouble(data[1]));
            }
        }catch (NumberFormatException e){

        } finally {
            if (cacheReader != null || ftReader != null) {
                try {
                    cacheReader.close();
                    ftReader.close();
                } catch (IOException e) {
                }
            }
        }
    }
}
```

```java
46      @Override
47      public void map(LongWritable key, Text value, Context context) throws IOException, InterruptedException {
48
49          String line = value.toString();
50          String[] data = line.split(",");
51          double age1,age0,sex1,sex0,cp1,cp0,trestbps1,trestbps0,chol1,chol0,fbs1,fbs0,restecg1,restecg0,thalach1,thalach0,exang1,exang0,oldpeak1,oldpeak
52          age1=age0=sex1=sex0=cp1=cp0=trestbps1=trestbps0=chol1=chol0=fbs1=fbs0=restecg1=restecg0=thalach1=thalach0=exang1=exang0=oldpeak1=oldpeak0=slope
53          for (int i = 0; i < 14; i++) {
54              if( i == 0){
55                  age1 = normalDistribution(Double.parseDouble(data[0]), preProcessingTable.get("age,1,mean"), preProcessingTable.get("age,1,stddev"));
56                  age0 = normalDistribution(Double.parseDouble(data[0]), preProcessingTable.get("age,0,mean"), preProcessingTable.get("age,0,stddev"));
57              } else if(i == 1){
58                  sex1 = frequencyTable.get("sex,1,"+data[1].trim())/frequencyTable.get("class1");
59                  sex0 = frequencyTable.get("sex,0,"+data[1].trim())/frequencyTable.get("class0");
60              } else if (i == 2) {
61                  cp1 = frequencyTable.get("cp,1,"+data[2].trim())/frequencyTable.get("class1");
62                  cp0 = frequencyTable.get("cp,0,"+data[2].trim())/frequencyTable.get("class0");
63              } else if (i == 3) {
64                  trestbps1 = normalDistribution(Double.parseDouble(data[3]), preProcessingTable.get("trestbps,1,mean"), preProcessingTable.get("trestbps
65                  trestbps0 = normalDistribution(Double.parseDouble(data[3]), preProcessingTable.get("trestbps,0,mean"), preProcessingTable.get("trestbps

                                            ...

89                  ca0 = frequencyTable.get("ca,0,"+data[11].trim())/frequencyTable.get("class0");
90              } else if (i == 12) {
91                  thal1 = frequencyTable.get("thal,1,"+data[12].trim())/frequencyTable.get("class1");
92                  thal0 = frequencyTable.get("thal,0,"+data[12].trim())/frequencyTable.get("class0");
93              }
94          }
95          DecimalFormat df = new DecimalFormat("#.#####");
96          double total1 = frequencyTable.get("class1") / (frequencyTable.get("class1")+frequencyTable.get("class0"));
97          double total0 = frequencyTable.get("class0") / (frequencyTable.get("class1")+frequencyTable.get("class0"));
98          double class1 = age1*sex1*cp1*trestbps1*chol1*fbs1*restecg1*thalach1*exang1*oldpeak1*slope1*ca1*thal1*total1;
99          double class0 = age0*sex0*cp0*trestbps0*chol0*fbs0*restecg0*thalach0*exang0*oldpeak0*slope0*ca0*thal0*total0;
100         Estimation e = new Estimation(new DoubleWritable(class1), new DoubleWritable(class0));
101         context.write(value, e);
102     }
103
104     public double normalDistribution(Double value, Double mean, Double stddev){
105         return (1/(Math.sqrt(2*Math.PI)*stddev))*(Math.pow(Math.E,(-(Math.pow(value-mean, 2))/(2*Math.pow(stddev, 2)))));
106     }
```

# Faza 3: Reducer

```java
6  public class ReduceClassFIT extends Reducer<Text, Estimation, Text, IntWritable> {
7
8      protected void reduce(Text key, Iterable<Estimation> values, Context context)
9              throws IOException, InterruptedException {
10         double class1, class0;
11         class0 = class1 = 0;
12         for(Estimation val: values){
13             class1 = val.getClass1().get();
14             class0 = val.getClass0().get();
15         }
16         if(class1 > class0) {
17             context.write(key, new IntWritable(1));
18         } else {
19             context.write(key, new IntWritable(0));
20         }
21     }
22 }
```

**20**

# Hvala na pažnji!