# Linear regression with ROC

Anton Antonov
MathematicaForPrediction project at GitHub
MathematicaVsR project at GitHub
October 2016

## Introduction

This document demonstrates how to do in Mathematica linear regression (easily using the built-in function `LinearModelFit`) and to tune the binary classification with the derived model through the so called Receiver Operating Characteristic (ROC) framework, [5].

The data used in this document is from [1] and it has been analyzed in more detail in [2]. In this document we only show to how to ingest and do very basic analysis of that data before proceeding with the linear regression model and its tuning. The package [4] provides the needed ROC functionalities.

### Used packages

These commands load the packages [3,4]:

```
In[1]:= Import[
    "https://raw.githubusercontent.com/antononcube/MathematicaForPrediction/master/
      MathematicaForPredictionUtilities.m"]
    Import[
    "https://raw.githubusercontent.com/antononcube/MathematicaForPrediction/master/
      ROCFunctions.m"]
```

## Reading data

The code below imports the data.

```
In[3]:= lines = Import["~/Datasets/adult/adult.data"];
    lines = Select[lines, Length[#] > 3 &];
    Dimensions[lines]

Out[5]= {32 561, 15}
```

```
In[6]:= linesTest = Import["~/Datasets/adult/adult.test"];
    linesTest = Select[linesTest, Length[#] > 3 &];
    Dimensions[linesTest]

Out[8]= {16 281, 15}
```

```
In[9]:= columnNames = StringSplit[
        "age,workclass,fnlwgt,education,education-num,marital-status,occupation,
          relationship,race,sex,capital-gain,capital-loss,hours-per-week,native-
          country", ","]
```

```
Out[9]= {age, workclass, fnlwgt, education, education-num, marital-status, occupation,
        relationship, race, sex, capital-gain, capital-loss, hours-per-week, native-country}
```

```
In[10]:= AppendTo[columnNames, "income"]
```

```
Out[10]= {age, workclass, fnlwgt, education, education-num,
        marital-status, occupation, relationship, race, sex, capital-gain,
        capital-loss, hours-per-week, native-country, income}
```

```
In[11]:= aColumnNames = AssociationThread[columnNames → Range[Length[columnNames]]];
```

In[12]:= `Magnify[#, 0.6] &@@GridTableForm[lines[[1 ;; 12]],`
`TableHeadings → Map[Style[#, Blue, FontFamily → "Times"] &, columnNames]]`

Out[12]=

| # | age | workclass | fnlwgt | education | education–num | marital–status | occupation | relationship | race | sex | capital–gain | capital–loss | hours–per–week | native–country | incor |
|---|-----|-----------|--------|-----------|--------------|----------------|------------|--------------|------|-----|--------------|--------------|----------------|----------------|-------|
| 1 | 39 | State–gov | 77 516 | Bachelors | 13 | Never-married | Adm-clerical | Not-in–family | White | Male | 2174 | 0 | 40 | United–States | <= |
| 2 | 50 | Self–emp–not–inc | 83 311 | Bachelors | 13 | Married–civ–spouse | Exec-managerial | Husband | White | Male | 0 | 0 | 13 | United–States | <= |
| 3 | 38 | Private | 215 646 | HS–grad | 9 | Divorced | Handlers-cleaners | Not-in–family | White | Male | 0 | 0 | 40 | United–States | <= |
| 4 | 53 | Private | 234 721 | 11th | 7 | Married–civ–spouse | Handlers–cleaners | Husband | Black | Male | 0 | 0 | 40 | United–States | <= |
| 5 | 28 | Private | 338 409 | Bachelors | 13 | Married–civ–spouse | Prof-specialty | Wife | Black | Female | 0 | 0 | 40 | Cuba | <= |
| 6 | 37 | Private | 284 582 | Masters | 14 | Married–civ–spouse | Exec-managerial | Wife | White | Female | 0 | 0 | 40 | United–States | <= |
| 7 | 49 | Private | 160 187 | 9th | 5 | Married–spouse-absent | Other-service | Not-in–family | Black | Female | 0 | 0 | 16 | Jamaica | <= |
| 8 | 52 | Self–emp–not–inc | 209 642 | HS–grad | 9 | Married–civ–spouse | Exec-managerial | Husband | White | Male | 0 | 0 | 45 | United–States | >5 |
| 9 | 31 | Private | 45 781 | Masters | 14 | Never-married | Prof-specialty | Not-in–family | White | Female | 14 084 | 0 | 50 | United–States | >5 |
| 10 | 42 | Private | 159 449 | Bachelors | 13 | Married–civ–spouse | Exec-managerial | Husband | White | Male | 5178 | 0 | 40 | United–States | >5 |
| 11 | 37 | Private | 280 464 | Some-college | 10 | Married–civ–spouse | Exec-managerial | Husband | Black | Male | 0 | 0 | 80 | United–States | >5 |
| 12 | 30 | State–gov | 141 297 | Bachelors | 13 | Married–civ–spouse | Prof-specialty | Husband | Asian-Pac-Islander | Male | 0 | 0 | 40 | India | >5 |

In[13]:=
```
Magnify[#, 0.6] &@GridTableForm[linesTest[[1 ;; 12]],
    TableHeadings → Map[Style[#, Blue, FontFamily → "Times"] &, columnNames]]
```

Out[13]=

| # | age | workclass | fnlwgt | education | education–num | marital–status | occupation | relationship | race | sex | capital–gain | capital–loss | hours–per–week | native–country | incom |
|---|-----|-----------|--------|-----------|---------------|----------------|------------|--------------|------|-----|--------------|--------------|----------------|----------------|-------|
| 1 | 25 | Private | 226 802 | 11th | 7 | Never-married | Machine-op-inspct | Own-child | Black | Male | 0 | 0 | 40 | United-States | <=5 |
| 2 | 38 | Private | 89 814 | HS-grad | 9 | Married-civ-spouse | Farming-fishing | Husband | White | Male | 0 | 0 | 50 | United-States | <=5 |
| 3 | 28 | Local-gov | 336 951 | Assoc-acdm | 12 | Married-civ-spouse | Protective-serv | Husband | White | Male | 0 | 0 | 40 | United-States | >50 |
| 4 | 44 | Private | 160 323 | Some-college | 10 | Married-civ-spouse | Machine-op-inspct | Husband | Black | Male | 7688 | 0 | 40 | United-States | >50 |
| 5 | 18 | ? | 103 497 | Some-college | 10 | Never-married | ? | Own-child | White | Female | 0 | 0 | 30 | United-States | <=5 |
| 6 | 34 | Private | 198 693 | 10th | 6 | Never-married | Other-service | Not-in-family | White | Male | 0 | 0 | 30 | United-States | <=5 |
| 7 | 29 | ? | 227 026 | HS-grad | 9 | Never-married | ? | Unmarried | Black | Male | 0 | 0 | 40 | United-States | <=5 |
| 8 | 63 | Self-emp-not-inc | 104 626 | Prof-school | 15 | Married-civ-spouse | Prof-specialty | Husband | White | Male | 3103 | 0 | 32 | United-States | >50 |
| 9 | 24 | Private | 369 667 | Some-college | 10 | Never-married | Other-service | Unmarried | White | Female | 0 | 0 | 40 | United-States | <=5 |
| 10 | 55 | Private | 104 996 | 7th-8th | 4 | Married-civ-spouse | Craft-repair | Husband | White | Male | 0 | 0 | 10 | United-States | <=5 |
| 11 | 65 | Private | 184 454 | HS-grad | 9 | Married-civ-spouse | Machine-op-inspct | Husband | White | Male | 6418 | 0 | 40 | United-States | >50 |
| 12 | 36 | Federal-gov | 212 465 | Bachelors | 13 | Married-civ-spouse | Adm-clerical | Husband | White | Male | 0 | 0 | 40 | United-States | <=5 |

# Assignment of training and testing data

As usual in classification and regression problems we work with two data sets: a training data set and a testing data set. Here we split the original training set into two sets training set and tuning set. The tuning set is going to be used to find a good value of a tuning parameter through ROC.

### Training data

In[14]:= 
```
data = lines;
data[All, -1] = Map[StringTrim, data[All, -1]];
```

In[16]:= 
```
trainingInds = RandomSample[Range[Length[data]], Ceiling[Length[data] * 0.75]];
tuningInds = Complement[Range[Length[data]], trainingInds];
trainingData = data[trainingInds];
tuningData = data[tuningInds];
```

### Testing data

In[20]:= 
```
testData = linesTest;
testData[All, -1] = Map[StringDrop[StringTrim[#], -1] &, testData[All, -1]];
```

# Some preliminary data analysis

Before doing regression it is a good idea to do some preliminary analysis of the data. For that we are going to use functions defined in the package [3].

In[22]:= 
```
Magnify[#, 0.7] &@@Grid[ArrayReshape[RecordsSummary[data, columnNames], {3, 5}],
  Dividers → All, Alignment → {Left, Top}]
```

Out[22]=

| 1 age | | 2 workclass | | 3 fnlwgt | | 4 education | | 5 education-num | |
|---|---|---|---|---|---|---|---|---|---|
| Min | 17 | Private | 22 696 | Min | 12 285 | HS-grad | 10 501 | Min | 1 |
| 1st Qu | 28 | | 2541 | 1st Qu | 471 297 / 4 | Some-college | 7291 | 1st Qu | 9 |
| Median | 37 | Self-emp-not-inc | | Median | 178 356 | Bachelors | 5355 | Median | 10 |
| Mean | 1 256 257 / 32 561 | | | Mean | 6 179 373 392 / 32 561 | Masters | 1723 | Mean | 328 237 / 32 561 |
| 3rd Qu | 48 | Local-gov | 2093 | | | Assoc-voc | 1382 | 3rd Qu | 12 |
| Max | 90 | ? | 1836 | 3rd Qu | $\frac{474\,109}{2}$ | 11th | 1175 | Max | 16 |
| | | State-gov | 1298 | Max | 1 484 705 | (Other) | 5134 | | |
| | | Self-emp-inc | 1116 | | | | | | |
| | | (Other) | 981 | | | | | | |

| 6 marital-status | | 7 occupation | | 8 relationship | | 9 race | | 10 sex | |
|---|---|---|---|---|---|---|---|---|---|
| Married-civ-spouse | 14 976 | Prof-specialty | 4140 | Husband | 13 193 | White | 27 816 | Male | 21 790 |
| Never-married | 10 683 | Craft-repair | 4099 | Not-in-family | 8305 | Black | 3124 | Female | 10 771 |
| Divorced | 4443 | Exec-managerial | 4066 | Own-child | 5068 | Asian-Pac-Islander | 1039 | | |
| Separated | 1025 | Adm-clerical | 3770 | Unmarried | 3446 | | | | |
| Widowed | 993 | Sales | 3650 | Wife | 1568 | Amer-Indian-Eskimo | 311 | | |
| Married-spouse-absent | 418 | Other-service | 3295 | Other-relative | 981 | | | | |
| Married-AF-spouse | 23 | (Other) | 9541 | | | Other | 271 | | |

| 11 capital-gain | | 12 capital-loss | | 13 hours-per-week | | 14 native-country | | 15 income | |
|---|---|---|---|---|---|---|---|---|---|
| 1st Qu | 0 | 1st Qu | 0 | Min | 1 | United-States | 29 170 | <=50K | 24 720 |
| 3rd Qu | 0 | 3rd Qu | 0 | 1st Qu | 40 | Mexico | 643 | >50K | 7841 |
| Median | 0 | Median | 0 | Median | 40 | ? | 583 | | |
| Min | 0 | Min | 0 | Mean | 1 316 684 / 32 561 | Philippines | 198 | | |
| Mean | $\frac{35\,089\,324}{32\,561}$ | Mean | $\frac{2\,842\,700}{32\,561}$ | 3rd Qu | 45 | Germany | 137 | | |
| Max | 99 999 | Max | 4356 | Max | 99 | Canada | 121 | | |
| | | | | | | (Other) | 1709 | | |

In[23]:=
```
Magnify[#, 0.7] &@Grid[ArrayReshape[RecordsSummary[testData, columnNames], {3, 5}],
    Dividers → All, Alignment → {Left, Top}]
```

Out[23]=

| 1 age | 2 workclass | 3 fnlwgt | 4 education | 5 education−num |
|---|---|---|---|---|
| Min       17<br>1st Qu  28<br>Median  37<br>Mean    210 391 / 5427<br>3rd Qu  48<br>Max       90 | Private          11 210<br>                   1321<br><br>  Self-emp-not-<br>  inc<br>Local-gov        1043<br>?                 963<br>State-gov         683<br>Self-emp-inc      579<br>(Other)           482 | Min      13 492<br>1st Qu  466 865 / 4<br>Median  177 831<br>Mean     3 084 202 270 /<br>            16 281<br>3rd Qu  238 384<br>Max      1 490 400 | HS-grad         5283<br>Some-college  3587<br>Bachelors       2670<br>Masters          934<br>Assoc-voc        679<br>11th             637<br>(Other)         2491 | Min       1<br>1st Qu  9<br>Median  10<br>Mean     163 997 / 16 281<br>3rd Qu  12<br>Max      16 |

| 6 marital−status | 7 occupation | 8 relationship | 9 race | 10 sex |
|---|---|---|---|---|
| Married-civ-       7403<br>  spouse<br>Never-married    5434<br>Divorced           2190<br>Widowed            525<br>Separated          505<br>                    210<br><br>  Married-spouse<br>  -absent<br>Married-AF-          14<br>  spouse | Prof-specialty   2032<br>Exec-managerial 2020<br>Craft-repair      2013<br>Sales             1854<br>Adm-clerical      1841<br>Other-service     1628<br>(Other)           4893 | Husband          6523<br>Not-in-family   4278<br>Own-child        2513<br>Unmarried        1679<br>Wife              763<br>Other-relative  525 | White           13 946<br>Black            1561<br>Asian-Pac-       480<br>  Islander<br>Amer-Indian-     159<br>  Eskimo<br>Other            135 | Male     10 860<br>Female 5421 |

| 11 capital−gain | 12 capital−loss | 13 hours−per−week | 14 native−country | 15 income |
|---|---|---|---|---|
| 1st Qu  0<br>3rd Qu  0<br>Median  0<br>Min       0<br>Mean     $\frac{5\,871\,499}{5427}$<br>Max      99 999 | 1st Qu  0<br>3rd Qu  0<br>Median  0<br>Min       0<br>Mean     $\frac{1\,431\,088}{16\,281}$<br>Max      3770 | Min       1<br>1st Qu  40<br>Median  40<br>Mean     $\frac{657\,626}{16\,281}$<br>3rd Qu  45<br>Max       99 | United-States  14 662<br>Mexico           308<br>?                274<br>Philippines       97<br>Puerto-Rico       70<br>Germany           69<br>(Other)          801 | <=50K  12 435<br>>50K    3846 |

Looking at the column "income" we can see that for both datasets the people who earn more than $50000 is ≈ 25% of all people. We will consider ">50" to be the more important class label for the classifiers built below.
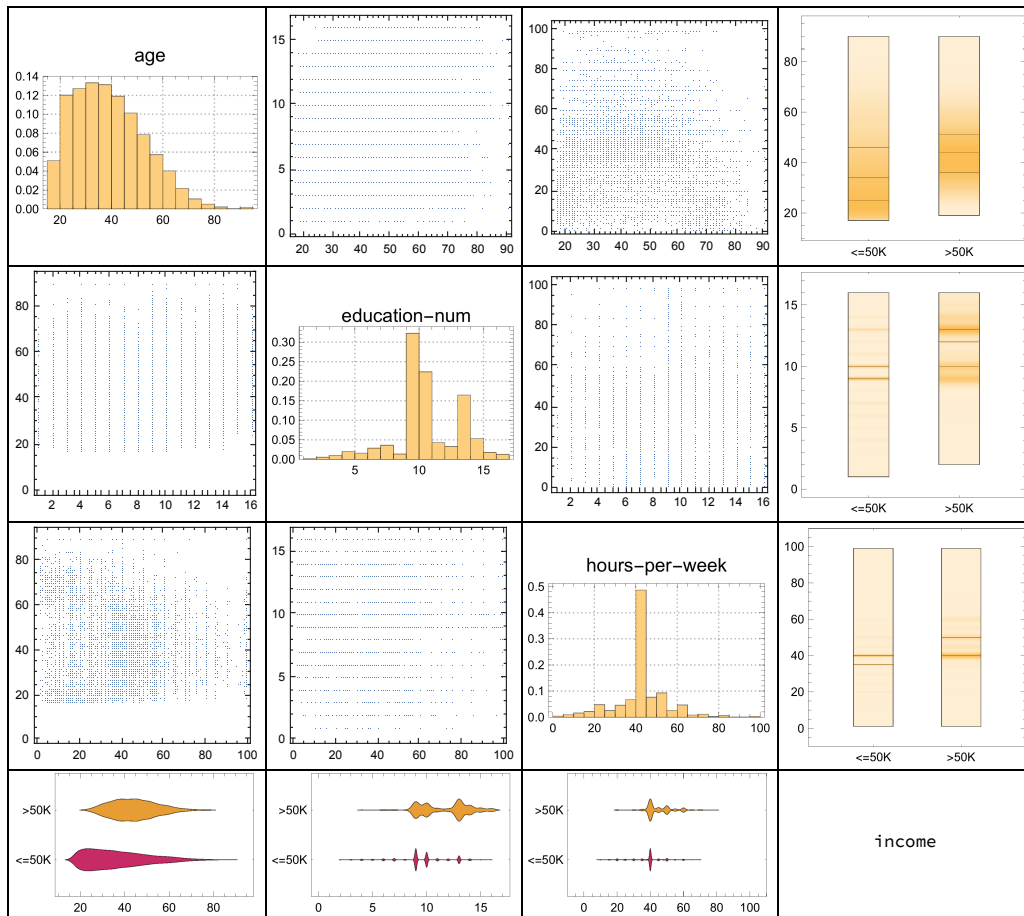
For simplicity of the exposition below we are going to use only the columns "age", "education-num", "hours-per-week", "income".

In[24]:=
```
columnNamesExplantoryVars = {"age", "education-num", "hours-per-week"};
columnNameResponseVar = "income";
columnNamesForAnalysis =
    Append[columnNamesExplantoryVars, columnNameResponseVar];
```

Here is variable dependence grid for those variables:

In[27]:= `Magnify[#, 0.7] &@VariableDependenceGrid[`
   `data[All, aColumnNames /@ columnNamesForAnalysis], columnNamesForAnalysis]`

Out[27]=



We can see from the last row of the plot above that the variables "age", "education-num", "hours-per-week" can explain "income" at least to a degree. We see that higher values of "age", "education-num", "hours-per-week" are associated closer with ">50K". For more detailed analysis see [2].

# LinearModelFit

`LinearModelFit` has several signatures. Doing Linear regression over the data we have is most convenient with the signature `LinearModelFit[{m,v}]` (using a design matrix *m* and a response vector *v*.)

As mentioned above in order to keep the exposition simple we do the regression with the three numerical columns "age", "education-num", and "hours-per-week". With the replacement rules `{"<=50K"→0,">50K"→1}` we convert the data column "income" into a vector of 0's and 1's. The result of `LinearModelFit` is a function based on the training set of data.

In[28]:= `lfmFunc = LinearModelFit[`
    `{trainingData[All, aColumnNames /@ {"age", "education-num", "hours-per-week"}],`
    `trainingData[All, aColumnNames["income"]] /. {"<=50K" → 0, ">50K" → 1}}]`

Out[28]= `FittedModel` | `0.00187859 #1 + 0.0167414 #2 + 0.000599763 #3` |

We use the model from the training data on the test data:

In[29]:= `tuningLables = tuningData[All, aColumnNames["income"]] /. {"<=50K" → 0, ">50K" → 1};`

Next we are going to evaluate what is the classification success of the derived model.

In[30]:= `modelValues = lfmFunc @@@`
    `tuningData[All, aColumnNames /@ {"age", "education-num", "hours-per-week"}];`

In[31]:= `modelValues[1 ;; 12]`

Out[31]= `{0.246049, 0.294229, 0.218629, 0.342148, 0.275292, 0.326165,`
    `0.217739, 0.315059, 0.266714, 0.218629, 0.16063, 0.256483}`

Although the response vector given to `LinearModelFit` is of 0's and 1's the regression model values are reals within a smaller than [0, 1] range.

In[32]:= `RecordsSummary[modelValues]`

Out[32]= 
```
1 column 1
Min      0.0689454
1st Qu   0.229504
Median   0.263319
Mean     0.264793
3rd Qu   0.299186
Max      0.463377
```

Here is a histogram of values from the regression model:

In[33]:= `Histogram[modelValues, Automatic, "Probability",`
    `PlotLabel → "Regression model values\nover the test data",`
    `AxesLabel → {"FittedModel values", "Probability"}]`

Out[33]=


We pick a threshold above which the model values are considered to be 1's (and hence ">=50K").

In[34]:= `With[{θ = 0.3}, modelLabels = Map[If[# > θ, 1, 0] &, modelValues]];`

In[35]:= `modelLabels[[1 ;; 12]]`

Out[35]= `{0, 0, 0, 1, 0, 1, 0, 1, 0, 0, 0, 0}`

Here is a table that shows classification success of the regression model with chosen threshold:

In[36]:= `labelsROC = SortBy[Tally[Transpose[{tuningLables, modelLabels}]], First];`
`labelsROC = Flatten /@`
`    MapThread[Append, {labelsROC, labelsROC[[All, 2]] / Total[labelsROC[[All, 2]]] // N}];`
`TableForm[labelsROC, TableHeadings →`
`   {{"true negative", "false positive", "false negative", "true positive"},`
`    {"test labels", "model labels", "freq", "%"}}]`

Out[38]//TableForm=

|  | test labels | model labels | freq | % |
|---|---|---|---|---|
| true negative | 0 | 0 | 5218 | 0.641032 |
| false positive | 0 | 1 | 953 | 0.117076 |
| false negative | 1 | 0 | 927 | 0.113882 |
| true positive | 1 | 1 | 1042 | 0.12801 |

We want to determine the threshold that gives the best classification success. What is "best" can be viewed and determined in several ways. We are going to use the so called Receiver Operating Characteristic (ROC); see [5].

(The table above is similar to the confusion matrix produced by *Mathematica*'s function `Classify` made available through `ClassifierMeasurements`. See the documentation for these function.)

# LinearModelFit with ROC

In this section we take a more systematic approach of determining the best threshold to be used to separate the regression model values.

We are going to call **positive** the income values ">50K" and **negative** the income values "<=50K". Again see [2]. As we mentioned above, we will consider ">50" to be the more important class label for the classifiers built below.

For the ROC functionalities are employed through the package [4].

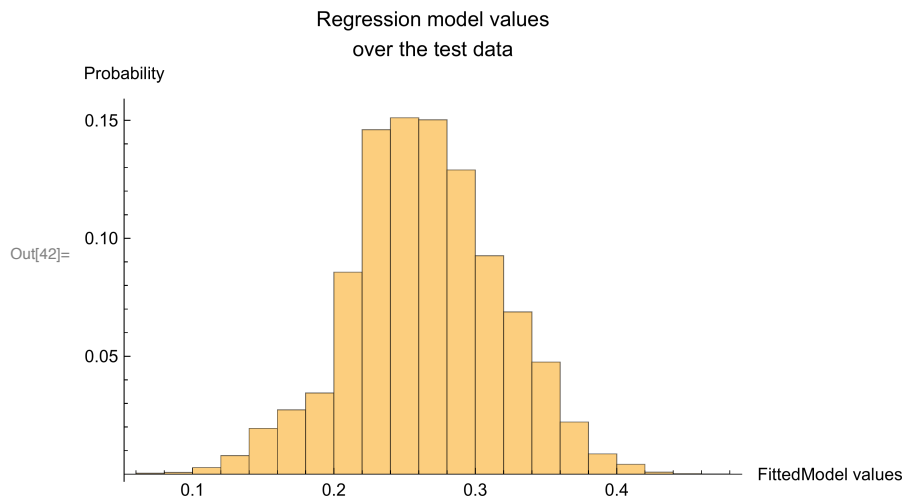## Linear regression classification definitions

```
In[39]:= Clear[ModelLabelsByThreshold]
    ModelLabelsByThreshold[modelValues_?VectorQ, θ_?NumberQ] :=
      Map[If[# > θ, 1, 0] &, modelValues];
    ModelLabelsByThreshold[lfmFunc_FittedModel, testData_?MatrixQ,
       aColumnNames_Association, θ_?NumberQ] :=
      Block[{t, testLables, modelValues, modelLabels},
       testLables = testData[[All, aColumnNames["income"]]] /. {"<=50K" → 0, ">50K" → 1};
       modelValues = lfmFunc @@@
         testData[[All, aColumnNames /@ {"age", "education-num", "hours-per-week"}]];
       ModelLabelsByThreshold[modelValues, testLables, θ]
       ];
```

## Computations to find the best threshold

Looking at the histogram of classification values:

```
In[42]:= Histogram[modelValues, Automatic, "Probability",
     PlotLabel → "Regression model values\nover the test data",
     AxesLabel → {"FittedModel values", "Probability"}]
```

Out[42]=



we decide to use the following a range of thresholds:

```
In[43]:= thRange = Range[0.15, 0.42, 0.01];
```

Compute the model values

In[44]:= ```
AbsoluteTiming[
  tuningLabels =
    testData〚All, aColumnNames[columnNameResponseVar]〛 /. {"<=50K" → 0, ">50K" → 1};
  modelValues = lfmFunc @@@ testData〚All, aColumnNames /@ columnNamesExplantoryVars〛];
 ]
```

Out[44]= {1.7642, Null}

Compute the ROC data for each threshold of the range of thresholds, and convert the ROC data to associations convenient for the computation of ROC functions:

In[45]:= ```
aROCs = Table[ToROCAssociation[{1, 0}, tuningLabels,
     ModelLabelsByThreshold[modelValues, θ]], {θ, thRange}];
```

Here is an example of evaluating some of the ROC functions over one of the elements of aROCs:

In[46]:= ```
Through[ROCFunctions[{"PPV", "NPV", "ACC"}][aROCs〚3〛]]
```

Out[46]= $\left\{ \dfrac{1917}{7796}, \dfrac{677}{689}, \dfrac{4511}{16\,281} \right\}$

Plot of the functions PPV, NPV, TPR, SPC, ACC:

In[47]:= ```
ListLinePlot[Map[Transpose[{thRange, #}] &, Transpose[
    Map[Through[ROCFunctions[{"PPV", "NPV", "TPR", "ACC", "SPC"}][#]] &, aROCs]]],
  Frame → True, FrameLabel → Map[Style[#, Larger] &, {"threshold, θ", "rate"}],
  PlotLegends → Map[# <> ", " <> (ROCFunctions["FunctionInterpretations"][#]) &,
    {"PPV", "NPV", "TPR", "ACC", "SPC"}], GridLines → Automatic]
```

Out[47]=



Using the plot above we can select a threshold for separating the model values into 0's and 1's ("<=50K" and ">50K" respectively). We can see that when the classification accuracy increases the positive predictive value decreases. A good threshold value is 0.3 because TPR and ACC have satisfactory, large enough values.
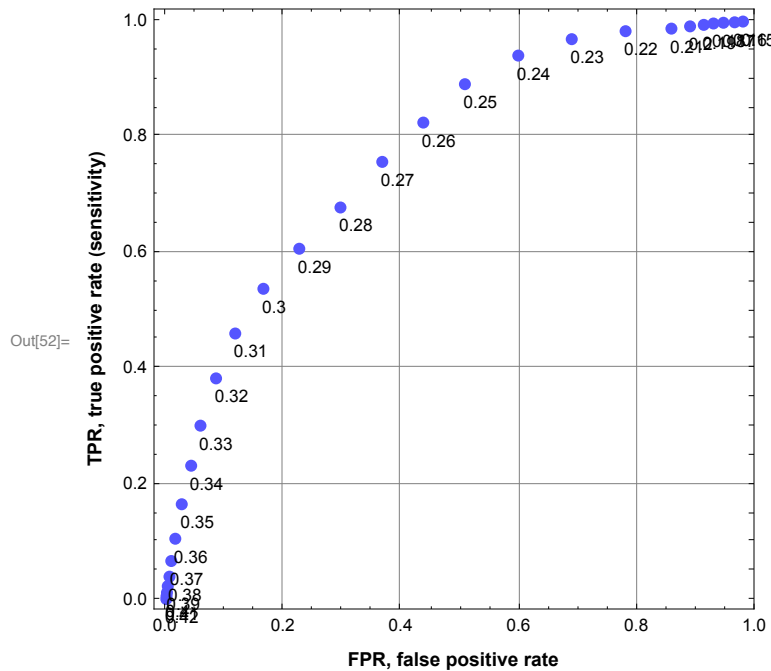
Let us find the intersection point of the curves PPV and TPR. Examining the plot above we can come up with the initial condition for *x*.

In[48]:= `Clear[θ]`

`ppvFunc = Interpolation[Transpose@{thRange, ROCFunctions["PPV"] /@ aROCs}];`

`tprFunc = Interpolation[Transpose@{thRange, ROCFunctions["TPR"] /@ aROCs}];`

`sol = FindRoot[ppvFunc[θ] - tprFunc[θ] == 0, {θ, 0.25}]`

Out[51]= $\{\theta \to 0.302858\}$

Here is a different plot used typically in ROC analysis:

In[52]:= `ROCPlot[thRange, aROCs, "PlotJoined" → False,`

`  "ROCPointCallouts" → True, "ROCPointTooltips" → True, GridLines → Automatic]`

Out[52]=



## Accuracy over the test data

We split the original training data into two parts for training and tuning. Using the found threshold, let us use evaluate the classification process over the test data.

In[53]:= `modelValues = lfmFunc @@@ testData〚All, aColumnNames /@ columnNamesExplantoryVars〛;`

`modelLabels = Map[If[# > (θ /. sol), ">50K", "<=50K"] &, modelValues];`

Using the actual labels and the predicted labels (`modelLabels`) let us find the contingency values with `ToROCAssociation`:

In[55]:= `clRes = ToROCAssociation[{">50K", "<=50K"},`

`  testData〚All, aColumnNames[columnNameResponseVar]〛, modelLabels]`

Out[55]= ⟨| TruePositive → 1983, FalsePositive → 1882,

TrueNegative → 10 553, FalseNegative → 1863 |⟩

Using the above result let us compute all of the ROC functions:

```
In[56]:= GridTableForm[
    Transpose[{ROCFunctions["FunctionNames"], N@Through[ROCFunctions[][clRes]]}],
    TableHeadings → {"name", "value"}]
```

Out[56]=

| # | name | value |
|---|------|-------|
| 1 | TPR | 0.515601 |
| 2 | SPC | 0.848653 |
| 3 | PPV | 0.513066 |
| 4 | NPV | 0.849952 |
| 5 | FPR | 0.151347 |
| 6 | FDR | 0.486934 |
| 7 | FNR | 0.484399 |
| 8 | ACC | 0.769977 |

## Using a built-in classifier

If we use one of the built-in classifiers we can see that we got better results for the important class label ">50K" through the ROC analysis using an inferior base classifier (linear regression).

```
In[57]:= AbsoluteTiming[cf = Classify[data[[All, aColumnNames /@ columnNamesExplantoryVars]] →
        data[[All, aColumnNames[columnNameResponseVar]]], Method → "NeuralNetwork"];
    cfPredictedLabels = cf /@ testData[[All, aColumnNames /@ columnNamesExplantoryVars]];
    ]
```

Out[57]= {24.3918, Null}

```
In[58]:= cfCLRes = ToROCAssociation[{">50K", "<=50K"},
    testData[[All, aColumnNames[columnNameResponseVar]]], cfPredictedLabels]
GridTableForm[Transpose[{ROCFunctions["FunctionNames"],
    N@Through[ROCFunctions[][cfCLRes]]}], TableHeadings → {"name", "value"}]
```

Out[58]= ⟨|TruePositive → 1497, FalsePositive → 874,
    TrueNegative → 11561, FalseNegative → 2349|⟩

Out[59]=

| # | name | value |
|---|------|-------|
| 1 | TPR | 0.389236 |
| 2 | SPC | 0.929715 |
| 3 | PPV | 0.631379 |
| 4 | NPV | 0.831129 |
| 5 | FPR | 0.0702855 |
| 6 | FDR | 0.368621 |
| 7 | FNR | 0.610764 |
| 8 | ACC | 0.802039 |

# References

[1] Bache, K. & Lichman, M. (2013). UCI Machine Learning Repository. Irvine, CA: University of California, School of Information and Computer Science. Census Income Data Set, URL: http://archive.ics.uci.edu/ml/datasets/Census+Income .

[2] Anton Antonov, "Classification and association rules for census income data", (2014), Mathematica-ForPrediction at WordPress.com , URL: https://mathematicaforprediction.wordpress.com/2014/03/30/-classification-and-association-rules-for-census-income-data/ .

[3] Anton Antonov, MathematicaForPrediction utilities, (2014), source code MathematicaForPrediction at GitHub*,* package MathematicaForPredictionUtilities.m.

[4] Anton Antonov, Receiver operating characteristic functions Mathematica package, (2016), source code MathematicaForPrediction at GitHub*,* package ROCFunctions.m .

[5] Wikipedia entry, Receiver operating characteristic. URL: http://en.wikipedia.org/wiki/Receiver_operating_characteristic .