

# Ground truth inference of binary classifier accuracies - the independent classifiers case

The accuracy of binary classifiers can be estimated exactly without knowing the true labels for the classifier decisions. Three independent classifiers are sufficient. This notebook demonstrates the above claims by construction.

The method proposed here is algebraic as opposed to the maximum-likelihood methods commonly seen in ground truth inference algorithms.

---

## The Polynomial System - Algebra of Error for Independent Binary Classifiers

The following set of quartic polynomials define algebraic system that must be solved to estimate the accuracy of three classifiers as well as the true prevalence of the labels. There are  $8 = 2^3$  polynomials, one for each of the possible voting pattern of three classifiers doing binary label classification.

$$\begin{aligned} \text{groundTruthPolynomials} = \{ & -f_{\alpha,\alpha,\alpha} + P_{\alpha} P_{1,\alpha,\alpha} P_{2,\alpha,\alpha} P_{3,\alpha,\alpha} + P_{\beta} P_{1,\alpha,\beta} P_{2,\alpha,\beta} P_{3,\alpha,\beta}, \\ & -f_{\alpha,\alpha,\beta} + P_{\alpha} P_{1,\alpha,\alpha} P_{2,\alpha,\alpha} P_{3,\beta,\alpha} + P_{\beta} P_{1,\alpha,\beta} P_{2,\alpha,\beta} P_{3,\beta,\beta}, \\ & -f_{\alpha,\beta,\alpha} + P_{\alpha} P_{1,\alpha,\alpha} P_{2,\beta,\alpha} P_{3,\alpha,\alpha} + P_{\beta} P_{1,\alpha,\beta} P_{2,\beta,\beta} P_{3,\alpha,\beta}, \\ & -f_{\alpha,\beta,\beta} + P_{\alpha} P_{1,\alpha,\alpha} P_{2,\beta,\alpha} P_{3,\beta,\alpha} + P_{\beta} P_{1,\alpha,\beta} P_{2,\beta,\beta} P_{3,\beta,\beta}, \\ & -f_{\beta,\alpha,\alpha} + P_{\alpha} P_{1,\beta,\alpha} P_{2,\alpha,\alpha} P_{3,\alpha,\alpha} + P_{\beta} P_{1,\beta,\beta} P_{2,\alpha,\beta} P_{3,\alpha,\beta}, \\ & -f_{\beta,\alpha,\beta} + P_{\alpha} P_{1,\beta,\alpha} P_{2,\alpha,\alpha} P_{3,\beta,\alpha} + P_{\beta} P_{1,\beta,\beta} P_{2,\alpha,\beta} P_{3,\beta,\beta}, \\ & -f_{\beta,\beta,\alpha} + P_{\alpha} P_{1,\beta,\alpha} P_{2,\beta,\alpha} P_{3,\alpha,\alpha} + P_{\beta} P_{1,\beta,\beta} P_{2,\beta,\beta} P_{3,\alpha,\beta}, \\ & -f_{\beta,\beta,\beta} + P_{\alpha} P_{1,\beta,\alpha} P_{2,\beta,\alpha} P_{3,\beta,\alpha} + P_{\beta} P_{1,\beta,\beta} P_{2,\beta,\beta} P_{3,\beta,\beta} \} \\ \{ & -f_{\alpha,\alpha,\alpha} + P_{\alpha} P_{1,\alpha,\alpha} P_{2,\alpha,\alpha} P_{3,\alpha,\alpha} + P_{\beta} P_{1,\alpha,\beta} P_{2,\alpha,\beta} P_{3,\alpha,\beta}, \\ & -f_{\alpha,\alpha,\beta} + P_{\alpha} P_{1,\alpha,\alpha} P_{2,\alpha,\alpha} P_{3,\beta,\alpha} + P_{\beta} P_{1,\alpha,\beta} P_{2,\alpha,\beta} P_{3,\beta,\beta}, \\ & -f_{\alpha,\beta,\alpha} + P_{\alpha} P_{1,\alpha,\alpha} P_{2,\beta,\alpha} P_{3,\alpha,\alpha} + P_{\beta} P_{1,\alpha,\beta} P_{2,\beta,\beta} P_{3,\alpha,\beta}, \\ & -f_{\alpha,\beta,\beta} + P_{\alpha} P_{1,\alpha,\alpha} P_{2,\beta,\alpha} P_{3,\beta,\alpha} + P_{\beta} P_{1,\alpha,\beta} P_{2,\beta,\beta} P_{3,\beta,\beta}, \\ & -f_{\beta,\alpha,\alpha} + P_{\alpha} P_{1,\beta,\alpha} P_{2,\alpha,\alpha} P_{3,\alpha,\alpha} + P_{\beta} P_{1,\beta,\beta} P_{2,\alpha,\beta} P_{3,\alpha,\beta}, \\ & -f_{\beta,\alpha,\beta} + P_{\alpha} P_{1,\beta,\alpha} P_{2,\alpha,\alpha} P_{3,\beta,\alpha} + P_{\beta} P_{1,\beta,\beta} P_{2,\alpha,\beta} P_{3,\beta,\beta}, \\ & -f_{\beta,\beta,\alpha} + P_{\alpha} P_{1,\beta,\alpha} P_{2,\beta,\alpha} P_{3,\alpha,\alpha} + P_{\beta} P_{1,\beta,\beta} P_{2,\beta,\beta} P_{3,\alpha,\beta}, \\ & -f_{\beta,\beta,\beta} + P_{\alpha} P_{1,\beta,\alpha} P_{2,\beta,\alpha} P_{3,\beta,\alpha} + P_{\beta} P_{1,\beta,\beta} P_{2,\beta,\beta} P_{3,\beta,\beta} \} \end{aligned}$$

## The Ground Truth and the Inference Problem

In classification problems, the ground truth is the correct label for each data point in a dataset. The statistics of ground truth solved by the above polynomial system are:

- the prevalence of the labels.
- the accuracies of the classifiers for each label.

Note that these are sample statistics. The method is not trying to infer anything about the process that created the data over which the binary classifiers have been run. There is some data, you ran binary classifiers over it and now you want to ask statistical questions about the data that would be trivial to answer IF you had the true labels for each of the data points in your dataset. But you don't have the ground truth so you need techniques like the one explained here.

### Synthetic/Fake/Manufactured Ground Truth

The purpose of this notebook is to demonstrate the mathematics of ground truth inference with polynomials. So to keep it simple, we will be manufacturing our own ground truth. The following commands create the ground truth for this problem - a list of labels for the true label of each data point in the dataset.

```
desiredALabelPrevalence = 23 / 100
datasetSize = 10 000
labelGroundTruth =
  Table[If[RandomReal[] < desiredALabelPrevalence,  $\alpha$ ,  $\beta$ ], {datasetSize}];
  23
  100
  10 000
```

Let's check that the ground truth makes sense

```
RandomSample[labelGroundTruth, 10]
```

```
{ $\alpha$ ,  $\beta$ ,  $\alpha$ ,  $\beta$ ,  $\beta$ ,  $\beta$ ,  $\beta$ ,  $\alpha$ ,  $\beta$ ,  $\beta$ }
```

```
labelGroundTruth // Tally
```

```
{{ $\beta$ , 7675}, { $\alpha$ , 2325}}
```

Now we construct the synthetic data for the classifier label decisions. Again, we are using synthetic data to create a simple example that focuses on the math of the algorithm. You can substitute numbers from real independent classifiers if you have them.

1. The first thing we need to specify is the accuracy of each of the classifiers. This accuracy is one of the statistics we seek to estimate with the polynomial system that is constructed from the classifier decisions alone. I'll arbitrarily pick them uniformly from the same range. I need two accuracies, one for each label since we are doing binary classification.

```

classifierAccuracies = Table[
  {(* The  $\alpha$  accuracy *)  $\alpha \rightarrow \text{RandomReal}[\{0.7, 0.9\}]$ ,
    (* The  $\beta$  accuracy *)  $\beta \rightarrow \text{RandomReal}[\{0.7, 0.9\}]$ } // Association, {3}]
{<| $\alpha \rightarrow 0.724074$ ,  $\beta \rightarrow 0.762341$ |>,
  <| $\alpha \rightarrow 0.873416$ ,  $\beta \rightarrow 0.777599$ |>, <| $\alpha \rightarrow 0.758298$ ,  $\beta \rightarrow 0.757377$ |> }

```

2. Now we need a function that takes a classifier's accuracies, a list true labels, and produces a sample of its classification decisions.

- OtherLabel is a convenience function to create incorrect classification decisions.
- SyntheticClassification is the actual function that produces a sample of what a classifier with the specified accuracies would produce on a dataset.

```
OtherLabel[ $\alpha$ ] =  $\beta$ 
```

```
OtherLabel[ $\beta$ ] =  $\alpha$ 
```

```
 $\beta$ 
```

```
 $\alpha$ 
```

```

SyntheticClassification[classifierAccuracies_Association, dataset_List] :=
  dataset //
  Map[
    If[
      (* Throw a die, if below accuracy,
        pick the true label, if not, pick the other label *)
      RandomReal[] < classifierAccuracies[#, #, OtherLabel[#]] &,
      #] &

```

3. The classifier decisions are next. Three lists, one for each classifier, as long as the dataset.

```

classifierDecisions =
  Map[SyntheticClassification[#, labelGroundTruth] &, classifierAccuracies];

```

Let's make sure things are okay

```

classifierDecisions[[1]] // {RandomSample[#, 10], # // Length} &
classifierDecisions[[2]] // {RandomSample[#, 10], # // Length} &
classifierDecisions[[3]] // {RandomSample[#, 10], # // Length} &
{{ $\alpha$ ,  $\beta$ ,  $\beta$ ,  $\beta$ ,  $\beta$ ,  $\alpha$ ,  $\alpha$ ,  $\beta$ ,  $\alpha$ ,  $\beta$ }, 10 000}
{{ $\beta$ ,  $\beta$ ,  $\beta$ ,  $\beta$ ,  $\beta$ ,  $\beta$ ,  $\alpha$ ,  $\beta$ ,  $\alpha$ ,  $\alpha$ }, 10 000}
{{ $\beta$ ,  $\alpha$ ,  $\beta$ ,  $\beta$ ,  $\beta$ ,  $\alpha$ ,  $\beta$ ,  $\beta$ ,  $\beta$ ,  $\beta$ }, 10 000}

```

Let's also make sure that we are getting noisy classification of the true labels

```

Map[Tally, classifierDecisions]
{{{ $\beta$ , 6385}, { $\alpha$ , 3615}}, {{ $\beta$ , 6287}, { $\alpha$ , 3713}}, {{ $\alpha$ , 3653}, { $\beta$ , 6347}}}

```

4. Now we can construct the left side input of the ground truth polynomials for independent binary

classifiers - the frequency of the voting patters

```
votingPatternCounts = classifierDecisions // Transpose // Tally
{{{β, β, α}, 1181}, {{β, β, β}, 3433}, {{α, α, β}, 695}, {{α, β, α}, 526},
 {{β, α, β}, 1072}, {{α, β, β}, 1147}, {{α, α, α}, 1247}, {{β, α, α}, 699}}
```

Let's sort them and pretty print them to see the pattern better

```
votingPatternCounts // Sort // Column
{{α, α, α}, 1247}
{{α, α, β}, 695}
{{α, β, α}, 526}
{{α, β, β}, 1147}
{{β, α, α}, 699}
{{β, α, β}, 1072}
{{β, β, α}, 1181}
{{β, β, β}, 3433}
```

Everything looks okay with our synthetic data. Note that votingPatternCounts no longer contains any information about the true label for any of the data points in labelGroundTruth. Hence, any algorithm that only uses votingPatternCounts - the observed frequency of classifier voting patterns, would be carrying out ground truth inference. The goal of the next section is to estimate classifierAccuracies and the prevalence of the labels using only votingPatternCounts.

---

## Exact Solution for Three Independent Binary Classifiers

The above polynomial can be solved exactly for the unknown ground truth statistics given the tally of voting patters by the three classifiers.

The solution for the prevalence is a quadratic equation for the prevalence of, say, label A and complicated coefficients based on the 8 voting pattern frequencies.

This means that there are two solutions for a label prevalence -  $x\%$  or  $(1-x)\%$ . Likewise for the accuracies of the classifiers, although there is a symmetry between the accuracies for both labels.

This illustrates, in the case of classification tasks, the difference between accuracy and precision.

Curating a very small subset of the target dataset immediately establishes the correct solution. Accuracy is cheap, precision is expensive.

In general, these polynomial systems are solved using ideas from algebraic geometry such as a Groebner basis and Buchberger's algorithm.

### Choosing an independent variables set

There are two labels so we have two prevalences that are unknown. But, by definition, they must sum up to 100% in the dataset. So one variable is enough. I arbitrarily picked  $P_\alpha$ .

Similar thing happens with the accuracies of the classifiers. There are two labels and they can give two answers for each label so each classifier has 4 variables. And similar to the case of the prevalence, there are two equations the 4 variables have to obey. I chose representing everything in terms of the probabil-

ity of being correct given label.

If you put together the two paragraphs above, you can figure out that  $N$  independent classifiers would require  $2*N + 1$  variables to be solved for. So three classifiers should yield a polynomial system that can be reduced to 7 unknown variables ( $2*3 + 1$ ).

I've encapsulated the arbitrary choice of what those 7 variables need to be in the following rules

```
reduceToIndependentVariablesRules = {
  (* The prevalence *)
   $P_\beta \rightarrow (1 - P_\alpha)$ ,
  (* The alpha label rules for each of the three classifiers *)
  (*Sequence@@Flatten@Table[ $P_{i,\beta,\alpha} \rightarrow (1 - P_{i,\alpha,\alpha})$ , {i,3}])*
   $P_{i,\beta,\alpha} \rightarrow (1 - P_{i,\alpha,\alpha})$ ,
  (* The beta label rules for each of the three classifiers *)
   $P_{i,\alpha,\beta} \rightarrow (1 - P_{i,\beta,\beta})$ 
}
{ $P_\beta \rightarrow 1 - P_\alpha$ ,  $P_{i,\beta,\alpha} \rightarrow 1 - P_{i,\alpha,\alpha}$ ,  $P_{i,\alpha,\beta} \rightarrow 1 - P_{i,\beta,\beta}$ }
```

And let's transform our polynomials into using an independent set of variables

```
groundTruthPolynomials /. reduceToIndependentVariablesRules
{- fα,α,α + Pα P1,α,α P2,α,α P3,α,α + (1 - Pα) (1 - P1,β,β) (1 - P2,β,β) (1 - P3,β,β),
- fα,α,β + Pα P1,α,α P2,α,α (1 - P3,α,α) + (1 - Pα) (1 - P1,β,β) (1 - P2,β,β) P3,β,β,
- fα,β,α + Pα P1,α,α (1 - P2,α,α) P3,α,α + (1 - Pα) (1 - P1,β,β) P2,β,β (1 - P3,β,β),
- fα,β,β + Pα P1,α,α (1 - P2,α,α) (1 - P3,α,α) + (1 - Pα) (1 - P1,β,β) P2,β,β P3,β,β,
- fβ,α,α + Pα (1 - P1,α,α) P2,α,α P3,α,α + (1 - Pα) P1,β,β (1 - P2,β,β) (1 - P3,β,β),
- fβ,α,β + Pα (1 - P1,α,α) P2,α,α (1 - P3,α,α) + (1 - Pα) P1,β,β (1 - P2,β,β) P3,β,β,
- fβ,β,α + Pα (1 - P1,α,α) (1 - P2,α,α) P3,α,α + (1 - Pα) P1,β,β P2,β,β (1 - P3,β,β),
- fβ,β,β + Pα (1 - P1,α,α) (1 - P2,α,α) (1 - P3,α,α) + (1 - Pα) P1,β,β P2,β,β P3,β,β}
```

And let's check that we only have seven non-frequency variables

```
groundTruthPolynomials /. reduceToIndependentVariablesRules // Variables /@ # & //
  Flatten // DeleteDuplicates // SortBy[#, First] &
{fα,α,α, fα,α,β, fα,β,α, fα,β,β, fβ,α,α, fβ,α,β,
fβ,β,α, fβ,β,β, Pα, P1,α,α, P1,β,β, P2,α,α, P2,β,β, P3,α,α, P3,β,β}
```

## Preparing the observed voting frequencies rules

Recall that we have the voting pattern counts already

```
votingPatternCounts // Sort // Column
```

```
{ {α, α, α}, 1247 }
{ {α, α, β}, 695 }
{ {α, β, α}, 526 }
{ {α, β, β}, 1147 }
{ {β, α, α}, 699 }
{ {β, α, β}, 1072 }
{ {β, β, α}, 1181 }
{ {β, β, β}, 3433 }
```

Let's turn them into rules so we can substitute their value into the ground truth polynomials. We need to normalize them so we need the total count first. We know it is 10K, but let's check it by computation so we can verify everything is okay

```
totalVotes = votingPatternCounts // Last /@# & // Total
```

```
totalVotes == datasetSize
```

```
10 000
```

```
True
```

```
frequencyRules = votingPatternCounts //
```

```
Map[(Subscript[f, Sequence @@ #[[1]]] → #[[2]] / datasetSize) &, #] &
```

```
{ fβ,β,α →  $\frac{1181}{10\,000}$ , fβ,β,β →  $\frac{3433}{10\,000}$ , fα,α,β →  $\frac{139}{2000}$ , fα,β,α →  $\frac{263}{5000}$ ,  
fβ,α,β →  $\frac{67}{625}$ , fα,β,β →  $\frac{1147}{10\,000}$ , fα,α,α →  $\frac{1247}{10\,000}$ , fβ,α,α →  $\frac{699}{10\,000}$  }
```

## The final polynomial set

```
finalPolynomialSet =
```

```
(groundTruthPolynomials /. reduceToIndependentVariablesRules) /. frequencyRules
```

```
{ -  $\frac{1247}{10\,000}$  + Pα P1,α,α P2,α,α P3,α,α + (1 - Pα) (1 - P1,β,β) (1 - P2,β,β) (1 - P3,β,β),  
-  $\frac{139}{2000}$  + Pα P1,α,α P2,α,α (1 - P3,α,α) + (1 - Pα) (1 - P1,β,β) (1 - P2,β,β) P3,β,β,  
-  $\frac{263}{5000}$  + Pα P1,α,α (1 - P2,α,α) P3,α,α + (1 - Pα) (1 - P1,β,β) P2,β,β (1 - P3,β,β),  
-  $\frac{1147}{10\,000}$  + Pα P1,α,α (1 - P2,α,α) (1 - P3,α,α) + (1 - Pα) (1 - P1,β,β) P2,β,β P3,β,β,  
-  $\frac{699}{10\,000}$  + Pα (1 - P1,α,α) P2,α,α P3,α,α + (1 - Pα) P1,β,β (1 - P2,β,β) (1 - P3,β,β),  
-  $\frac{67}{625}$  + Pα (1 - P1,α,α) P2,α,α (1 - P3,α,α) + (1 - Pα) P1,β,β (1 - P2,β,β) P3,β,β,  
-  $\frac{1181}{10\,000}$  + Pα (1 - P1,α,α) (1 - P2,α,α) P3,α,α + (1 - Pα) P1,β,β P2,β,β (1 - P3,β,β),  
-  $\frac{3433}{10\,000}$  + Pα (1 - P1,α,α) (1 - P2,α,α) (1 - P3,α,α) + (1 - Pα) P1,β,β P2,β,β P3,β,β }
```

Let's check that only the seven unknown statistics of the ground truth remain in the final polynomial system.

```
statisticsOfGroundTruth =
  finalPolynomialSet // Variables /@# & // Flatten // DeleteDuplicates
{P $\alpha$ , P $_{1,\alpha,\alpha}$ , P $_{2,\alpha,\alpha}$ , P $_{3,\alpha,\alpha}$ , P $_{1,\beta,\beta}$ , P $_{2,\beta,\beta}$ , P $_{3,\beta,\beta}$ }
```

## The “magical” solution

```
Solve[Map[# == 0 &, finalPolynomialSet], statisticsOfGroundTruth] // N
{{P $\alpha$  → 0.237026, P $_{1,\alpha,\alpha}$  → 0.746384, P $_{2,\alpha,\alpha}$  → 0.872898,
  P $_{3,\alpha,\alpha}$  → 0.743696, P $_{1,\beta,\beta}$  → 0.758068, P $_{2,\beta,\beta}$  → 0.784526, P $_{3,\beta,\beta}$  → 0.752253},
 {P $\alpha$  → 0.762974, P $_{1,\alpha,\alpha}$  → 0.241932, P $_{2,\alpha,\alpha}$  → 0.215474, P $_{3,\alpha,\alpha}$  → 0.247747,
  P $_{1,\beta,\beta}$  → 0.253616, P $_{2,\beta,\beta}$  → 0.127102, P $_{3,\beta,\beta}$  → 0.256304}}
```

Recall what the “real” ground truth in this manufactured case was produced by a process that consisting of throwing loaded dice for being correct on the given label. For example, classifier 1 had a 72.4% chance of getting the  $\alpha$  label right.

```
classifierAccuracies
{{<| $\alpha$  → 0.724074,  $\beta$  → 0.762341|>,
  <| $\alpha$  → 0.873416,  $\beta$  → 0.777599|>, <| $\alpha$  → 0.758298,  $\beta$  → 0.757377|>}}
```

But note that classifierAccuracies are the dice we threw to produce the dataset. These need not be the same as the sample statistics or close to it! For example, if you threw the dice only once, the sample statistic for the accuracy would be 0 or 1. So we should really be comparing our answer to the sample statistic - ground truth inference is always about the sample statistics as befits a non-parametric method!

Let's calculate the sample accuracies of classifier 1 to see how far off the 0.746 answer is from it. We align our ground truth with the classifiers decisions

```
Length@classifierDecisions[[1]]
10 000

Length@labelGroundTruth
10 000

Transpose[{labelGroundTruth, classifierDecisions[[1]]}] // GroupBy[#, First] & //
  Map[Last, #, {2}] & // Tally /@# &
<| $\beta$  → {{ $\beta$ , 5796}, { $\alpha$ , 1879}},  $\alpha$  → {{ $\beta$ , 589}, { $\alpha$ , 1736}}|>
```

So the sample accuracy for the alpha data points for classifier 1 is

```
1736 / (1736 + 589) // N  
0.746667
```

Ha ha! Cool. This simple experiment should clarify what it means to say that ground truth inference algorithms measures statistics of the ground truth. And ground truth is always referring to the sample at hand in your study. The method estimated the sample statistic of classifier 1 on the alpha label - 0.746667 as 0.746384. That estimate was not trying to infer a statistic of the process, the accuracy for classifier 1 that I used to produce the synthetic data, that 0.724074 value that is stored in classifierAccuracies.

This approach was patented by Data Engines Corporation in 2010. Countless other patents/papers/investigations are possible once the reader realizes that the idea presented here - polynomial methods for statistics of ground truth - is widely applicable to other statistics and machine learning tasks. One example of this generalization is in the paper in this repository measuring the accuracy of an online ID service.