

1. サイコロの目の分布: ヒストグラム

前回、サイコロの目をランダムに出す方法を学んだ。[0,1)の一樣乱数を発生させる関数、rand()を使えば、`=int(rand()*6)+1`で1~6までの整数をランダムに出せることがわかった。

さて、正しく作られたサイコロならば、どの目が出る確率も1/6である(場合の数が6通りで、どれも等確率なので)。そこで、実際に1/6の確率になっているか、図示してみよう。

まず手始めに、サイコロを6回ふり、そのヒストグラムを書いてみよう。

1. A1~A6の6つのセルに、上と同じく`=int(rand()*6)+1`と書こう(一つだけ書いて、あとはコピーすればよい)。6回ふれば、原理的にはすべての目が1回ずつ出る可能性はあるが、たった6回では偏りが生じているのがわかるだろう。
2. 次に、ヒストグラムの横軸を明示しておく。つまり、サイコロの目の値、1~6を書いておく。今回は、C1~C6に1~6の数字を書いておこう(場所はどこでもよい)。
3. 「分析ツール」を使ってヒストグラムを作る。「データ」タブの右端にある「分析」をクリックする
※もし分析ツールがなければ、「Office ボタン」を押し、下に出てくる「Excel のオプション」をクリック、「アドイン」を選び、下にある「管理」から「Excel アドイン」を選び「設定」ボタンを押す。出てきたパネルの中の「分析ツール」にチェックを入れ「OK」を押すと使えるようになる。
4. 「ヒストグラム」を選ぶ。「入力範囲」は、サイコロの目の範囲、すなわちA1~A6である。「データ区間」はさきほど入力したヒストグラムのデータ範囲、C1~C6である。「出力オプション」ではどこにヒストグラムを出力するか指定できる。ここではE1を指定しておこう。これでOKをクリックすると、出力オプションで指定したセルにヒストグラムの数値データが出力される(OKするたびに新しい乱数に変わってしまうので、最初に出した乱数の分布を覚えておき、正しいかチェックしよう)。
5. グラフウィザードを使い、今度は棒グラフで示してみよう。「項目軸ラベルに使用」で「データ区間」のセルが選ばれるようにしておくこと。正しく図がかけたか、数値と照らし合わせてみよう。
6. F9を押してサイコロを振り直し、上と同様の操作を繰り返してみよう。

6回ふっただけでは、「確率が1/6」かどうかなんて確かめられないことがわかったと思う。1/6の割合で目が出るためには、もっと沢山ふらなければならぬ。そこで、今度は**60回ふってみて、そのヒストグラムを書いてみよう。**
課題:ヒストグラムの図を印刷して提出する。

2. ズレの度合い

次に、確率1/6からのズレを調べてみよう。

60回サイコロをふったら、それぞれの目が出る「期待値」は10回ずつである。これは、(全試行回数)×(その事象が実現する確率)であらわされる。そこで、それぞれの目が出た回数と、「10」とのズレを見る。

1の目が出た回数はE2に出力されていると思う(違う人は、以下の説明も適切に読み替えること)。そこで、その数値の隣(F2)に、`=E2-60*(1/6)`と書いてみよう(=E2-10でも同じことである)。これを他のすべての目についても行う(セルの右下をドラッグして、下にコピーすればよい)。これが、期待値からのズレである。

ズレの平均はいくらだろうか? F8セルに、F2~F7セルの値の平均を書いてみよう(ヒント:average関数)。結果はいくらになっただろうか? 0になったはずだ。なぜ0になるかはよく考えてみよう。

さて、これでは、どれくらいズレるのかの見通しがつかぬ。そこで、ズレを自乗してすべて正の値にしよう。さらに隣のG2セルに、`=F2^2`と書いてみよう。これで、F2セルの値の自乗が書き込まれたはずだ。これを他のすべての目に付いても実行し、さらにF8セルにF2~F7までの平均値を計算してみよう。

最後に、その平均値(ズレの自乗の平均値)の平方根をとってみよう。これが、**確率分布から期待される平均値に対する典型的なゆらぎの大きさ**となる。これを**標準偏差**と呼ぶ(自乗したものは**分散**と呼ばれる)。

60回ふった場合は、各目が出る期待値が10なので、標準偏差はおおよそ $\sqrt{10} \approx 3$ 程度になるはずだ(なぜ $\sqrt{10}$ で導けるのかは来週学ぶ)。

3. 標準偏差の試行回数依存性

さて、直感的には、試行回数を増やすほど、まんべんなくそれぞれの目が出現し、ズレの度合いは相対的に小さくなっていくと予想されるだろう。実際そうになっているか調べてみる。

まず、60回のときの相対的な「誤差」を求めておこう。これは(標準偏差) / (期待値)で見積もられる。なんとかサイコロを振りなおしてみ、大体どれくらいの値になるか調べてみよう。

次に、600回ふってみよう。期待値は100になる。この場合の、誤差を何回か求めてみよう。

結論を先に述べると、だいたい $1/\sqrt{N}$ 程度になるはずだ。つまり、期待値が10ならおよそ0.33, 100なら0.1、もし10000ならば、0.01となる。パーセントで表すと、それぞれ33%, 10%, 1%となる。従って、誤差1%で物事を言うならば、サンプルの大きさは10000程度ないといけないことがわかる。

4. 練習問題:血液型の分布

日本人の場合、血液型分布はおよそ A:O:B:AB=4:3:2:1 になっていることが知られている(これは民族によってかなり異なる)。では、10人の集団があったら、それぞれの血液型の人が4人、3人、2人、1人となるであろうか? そうはならないだろうということは、今までの結果からわかるだろう。では、少人数集団の場合、だいたいどれくらいのズレが生じるだろうか?

そこで、全体では上のような血液型分布になっている集団の中から、10人をランダムにピックアップしたとしよう。今回は、ランダムではあるが、期待値が一樣ではなく、血液型によって違いがあることに注意する。

まず、エクセルでは文字列よりも数値データが扱いやすいので、Aを1、Bを2、Oを3、ABを4と番号をつけておこう。次に、rand関数で生成される[0,1)一様乱数を用いて、期待値としては4:2:3:1になるようにうまく値を変換する。もっとも単純な方法は、出てきた乱数が0.4未満ならA(1)、そうではなくて0.6未満ならB(2)、さらにそうではなくて0.9未満ならO(3)、さらにそうでなければAB(4)とする。これならば、上記の期待値を実現できるであろう。

1. まず、**10人分の乱数を発生**させよう。A1セルに=rand()とやって、[0,1)一様乱数を発生させる。次に、それを10人分になるように下にコピーしよう。
2. 次に血液型に変換するため、数値を読み取り、条件によってどの値(1,2,3,4)を出力するかを判別しなければならない。そのために使えるのが**IF文**である。

IF構文は次のようになっている。=IF(A,a,b)、この意味は、「Aが真ならaを、偽ならbを実行せよ」ということである。試しにB1セルに次のように書いてみよう。=IF(A1<0.4,1,0)、すると、A1の値が0.4未満ならB1には1が表示され、0.4以上ならば0が表示されるはずである。いま、0は「A以外」に対応する。

「A以外」を具体的にしていこう。IF文の中にIF文を入れることも可能である。そこで、以下のようになれば、4つの血液型に分けられることになる、=IF(A1<0.4,1,IF(A1<0.6,2,IF(A1<0.9,3,4)))。これによって、A1の値が0.4未満なら1,0.4以上0.6未満なら2,0.6以上0.9未満なら3、0.9以上なら4が出力されることになる。この1,2,3,4がA,B,O,ABに対応していると読めばよい。

3. では、これらのヒストグラムを作ってみよう。やり方は3と同様であるのでやってみてほしい。ただし、図にする際は、横軸が数字ではわかりづらいので、どこかに(例えばD1~D4、どこでもよい)A,B,O,ABと書いておき、それを横軸に使用するようにしよう。

これも何度も乱数を振りなおして図を書いてみて欲しい。10人の集団では、血液型分布に偏りの出る場合がかなりあることがわかるだろう。

4. 次に、もっとサンプルを増やして、**100人でやってみよう**。するとどうなるか。10人のときよりも4:2:3:1に近づき、偏りの目立つ場合が少なくなっていることがわかると思う。

この場合も、典型的には \sqrt{N} 人程度の「誤差」がある。相対誤差にすれば、 $\sqrt{N}/N=1/\sqrt{N}$ となるので、サンプルが大きいくほがより期待値に近づくのである。

逆に言えば、少人数のサンプルから全体を推測することは極めて危険であることがわかる。「歴代首相にはO型が多いからO型はリーダー向き」のような主張の無意味さを読み取って欲しい¹。また、高々数百人のサンプルで得た視聴率では0.1ポイントの違いは完全に誤差の範囲で比較する意味がないこともわかるだろう。

課題:100人の集団に対するヒストグラムの図を印刷して提出する。

5. 課題(再掲)

1. サイコロを60回ふったときのヒストグラムのグラフ
2. ランダムに選んだ100人の血液型分布(ヒストグラム)のグラフ

これらを印刷して提出せよ(番号・名前を忘れずに)。

1 1万人以上のサンプルで慎重に調べた結果、血液型と性格の間には強い相関(血液型から性格を当てる、あるいはその逆が可能なくらいの強さ)はないということが既に明らかになっている。つまり、仮になんらかの相関があったとしても、それは1%程度のごくわずかな違いしかない、ということであり、日常生活のレベルでは血液型と性格は無関係であると言ってよい。