

## 【データの偏りを見抜く】

### 見えていない情報に気づく

データ分析にとりかかる前に、ちょっと考えてみましょう。もし、これから分析するために集めた「データ」そのものに「偏り」があったとしたら…分析をしてもうまくいきません。このような集めたデータの偏りは（ ）と呼ばれ、いろいろなものが存在します。これらは隠れてしまいやすいもので、見抜く力が必要です。

### 例1. アメリカの戦闘機を強くした統計学

第2次世界大戦では、飛行機がとても活躍しましたが、戦闘機は撃墜されれば機体が助かることはまずありません。かと言って装甲を厚くすれば重くなって速度が遅くなって不利になってしまいます。そこで、厚くする場所と薄くする場所を選ばなければならないわけです。そのために、帰還した戦闘機の弾痕調査を実施して分析しました。

【考えてみよう1】 図1のような弾痕パターンが見られたとき、他より装甲を厚くすべき場所を斜線で図2に示してみましょう。

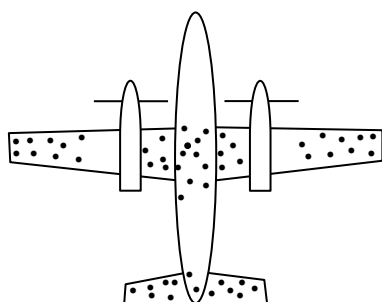


図1

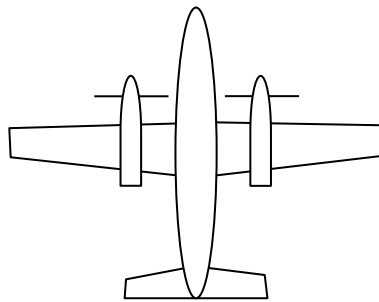


図2

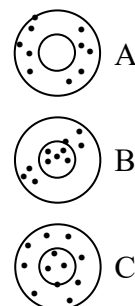


図3

- 考察① 動く丸い標的に銃を乱射したときの弾痕の予想は、のはず  
→考察② 飛行機も標的と同じであるとすると、被弾するはず  
→考察③ 弾痕の 場所に被弾したら、帰還できなかつたのではないか  
→結論 装甲を厚くすべき場所は、弾痕の 部分である。

●（ ）だけのデータを集めたため、（ ）のデータが全く集まっていない状態になっている。このような偏りを（ ）という。

【回答課題5】 次のデータ分析結果について、偏りを見抜いて説明しよう。(Classroom 質問)

「動物病院への調査データによると、ネコは5階以下の高さから落ちた場合よりも、6階以上の高さから落ちた場合の方が回復する割合が高い。終端速度に達して落下姿勢が安定したと考えられる。」

## 例2. ドイツを愛するハズレないタコ

2010年に一匹のタコが有名になったそうです。FIFA ワールドカップ 2010 南アフリカ大会で、ドイツ代表の試合8つの勝敗をすべての的中させたというのです。餌と一緒に置いた2つの国旗のどちらへ行くかで決めたそうです。

[考えてみよう2] 8つの試合の勝敗をすべての的中させる確率はどれくらいなのか、計算して分数で求めてみましょう。(引き分けは考えません。)

- これがとても珍しいことと言ってよいのか…というのは前提条件によるところがある。ワールドカップの勝敗は世界中の多くの人に関心を寄せ、同じように動物に勝敗を占わせていたのが( 256 ) 匹以上いたとすると、少なくとも1匹以上は確率的に的中するはずである。このとき、多くの( ) **存在がある**、ならば確率的に珍しいとは言えなくなる。これも一つの**バイアス**と言える。

## 例3. 食い違った2つの大統領選挙予想

1936年のアメリカ大統領選挙においては、民主党ルーズベルト大統領と共和党ランドン候補が争っていました。このとき、2つの世論調査の結果も対立しました。

①「リテラリー・ダイジェスト誌」…過去5回の大統領選挙結果を当てており、サンプル数は200万人。ランドン候補が57%の支持で有利と発表。

②「ギャラップ(当時はアメリカ世論研究所)」…当時は新興の調査会社であり、サンプル数は3000人。ルーズベルト大統領が54%の支持で有利と発表。

結果はルーズベルト大統領の圧勝となり、世間を驚かせました。

[考えてみよう3] 当時は①の調査結果が支持されていました。そのように信頼性が高いと思われていた理由を2点、簡潔に「~ので、~と思われていたから」という文で説明してみましょう。

- ・ と思われていたから
- ・ と思われていたから

- サンプル選びの違い…①と②ではサンプルの選び方に違いがあった。

①自社の雑誌読者、電話や自動車を所有する人を対象に得たサンプル

②都市の男女、農村の男女、中間層の男女…等の**層に細分化**し有権者数に合わせたサンプル

①のサンプルは結果として富裕層を中心に選んだことになり、それまでは母集団とあまり違わなかったものが、経済状況等の社会情勢の変化によって当てはまらなくなったと考えられる。②のようなサンプルの集め方は( ) という。

- サンプル選びの方法…人が選ぶと頼みやすさ等で偏りがでることから、コンピュータを用いた( ) **法**もしくは( ) **法**と呼ばれる方法が利用されている。その一つが無作為に電話番号を生成して電話して調査を依頼する **RDD (Random Digit Dialing) 法**であるが、これも間違いのない方法とは言い切れない。

## 【データの偏りを見抜く】

### 見えていない情報に気づく

データ分析にとりかかる前に、ちょっと考えてみましょう。もし、これから分析するために集めた「データ」そのものに「偏り」があったとしたら…分析をしてもうまくいきません。このような集めたデータの偏りは（ **バイアス** ）と呼ばれ、いろいろなものが存在します。これらは隠れてしまいやすいもので、見抜く力が必要です。

### 例1. アメリカの戦闘機を強くした統計学

第2次世界大戦では、飛行機がとても活躍しましたが、戦闘機は撃墜されれば機体が助かることはまずありません。かと言って装甲を厚くすれば重くなって速度が遅くなって不利になってしまいます。そこで、厚くする場所と薄くする場所を選ばなければならないわけです。そのために、帰還した戦闘機の弾痕調査を実施して分析しました。

【考えてみよう1】 図1のような弾痕パターンが見られたとき、他より装甲を厚くすべき場所を斜線で図2に示してみましょう。

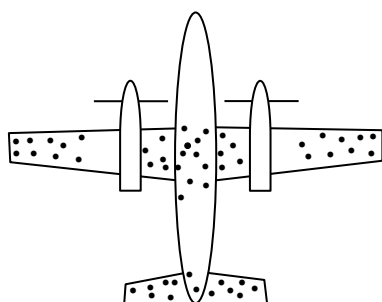


図1

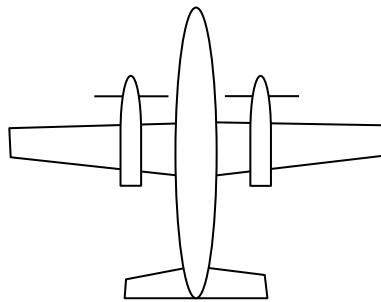


図2

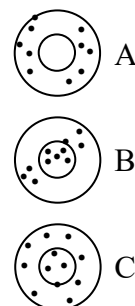


図3

- 考察① 動く丸い標的に銃を乱射したときの弾痕の予想は、 **図3-C** のはず
- 考察② 飛行機も標的と同じであるとすると、 **まんべんなく** 被弾するはず
- 考察③ 弾痕の **ない** 場所に被弾したら、帰還できなかったのではないか
- 結論 装甲を厚くすべき場所は、弾痕の **あまりない** 部分である。

- （ **生存者** ）だけのデータを集めたため、（ **死者** ）のデータが全く集まっていない状態になっている。このような偏りを（ **生存者バイアス** ）という。

【回答課題5】 次のデータ分析結果について、偏りを見抜いて説明しよう。（Classroom 質問）

「動物病院への調査データによると、ネコは5階以下の高さから落ちた場合よりも、6階以上の高さから落ちた場合の方が回復する割合が高い。終端速度に達して落下姿勢が安定したと考えられる。」

**6階以上の高さから落ちたネコには、病院に運ばれることなく亡くなったものが多くいたことが予想されるので、生存者バイアスが生じていると考えられる。POS データも同様。**

## 例2. ドイツを愛するハズレないタコ

2010年に一匹のタコが有名になったそうです。FIFA ワールドカップ 2010 南アフリカ大会で、ドイツ代表の試合8つの勝敗をすべての的中させたというのです。餌と一緒に置いた2つの国旗のどちらへ行くかで決めたそうです。

[考えてみよう2] 8つの試合の勝敗をすべての的中させる確率はどれくらいなのか、計算して分数で求めてみましょう。(引き分けは考えません。)

勝つか負けるかのどちらかとすると、的中の確率は $\frac{1}{2}$ ですから、8つ連続で $\left(\frac{1}{2}\right)^8 = \frac{1}{2^8} = \frac{1}{256}$ つまり、256回に1回的中する確率になります。

- これがとても珍しいことと言ってよいのか…というのは前提条件によるところがある。ワールドカップの勝敗は世界中の多くの人に関心を寄せ、同じように動物に勝敗を占わせていたのが(256)匹以上いたとすると、少なくとも1匹以上は確率的に的中するはずである。このとき、多くの(黙って隠れている)存在がある、ならば確率的に珍しいとは言えなくなる。これも一つのバイアスと言える。

## 例3. 食い違った2つの大統領選挙予想

1936年のアメリカ大統領選挙においては、民主党ルーズベルト大統領と共和党ランドン候補が争っていました。このとき、2つの世論調査の結果も対立しました。

①「リテラリー・ダイジェスト誌」…過去5回の大統領選挙結果を当てており、サンプル数は200万人。ランドン候補が57%の支持で有利と発表。

②「ギャラップ(当時はアメリカ世論研究所)」…当時は新興の調査会社であり、サンプル数は3000人。ルーズベルト大統領が54%の支持で有利と発表。

結果はルーズベルト大統領の圧勝となり、世間を驚かせました。

[考えてみよう3] 当時は①の調査結果が支持されていました。そのように信頼性が高いと思われていた理由を2点、簡潔に「～ので、～と思われていたから」という文で説明してみましょう。

- ・ 過去5回の結果を当てた実績があるので、今回も当たると思われていたから
- ・ サンプル数が一方の約700倍なので、より正確な結果だと思われていたから

- サンプル選びの違い…①と②ではサンプルの選び方に違いがあった。

①自社の雑誌読者、電話や自動車を所有する人を対象に得たサンプル

②都市の男女、農村の男女、中間層の男女…等の層に細分化し有権者数に合わせたサンプル

①のサンプルは結果として富裕層を中心に選んだことになり、それまでは母集団とあまり違わなかったものが、経済状況等の社会情勢の変化によって当てはまらなくなったと考えられる。②のようなサンプルの集め方は(層別サンプリング)という。

- サンプル選びの方法…人が選ぶと頼みやすさ等で偏りがでることから、コンピュータを用いた(ランダムサンプリング)法もしくは(無作為抽出)法と呼ばれる方法が利用されている。その一つが無作為に電話番号を生成して電話して調査を依頼するRDD(Random Digit Dialing)法であるが、これも間違いのない方法とはい切れぬ。