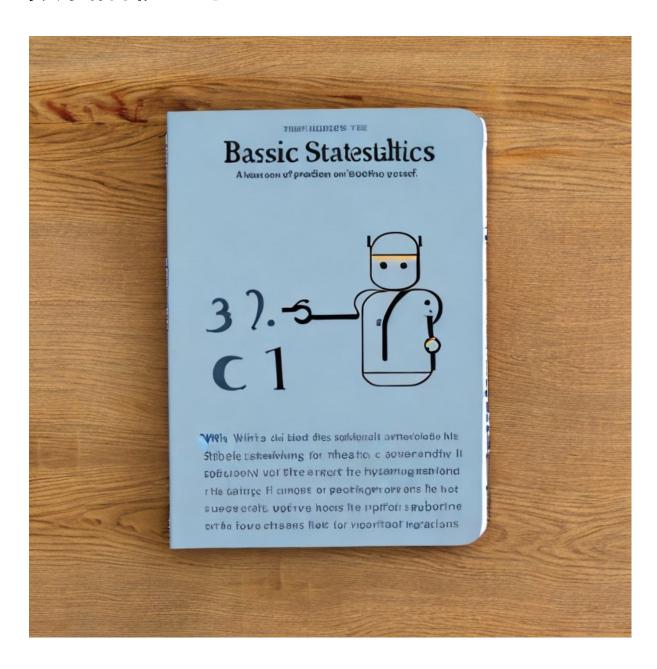
# ChatGPTからの贈り物: 初学者のための 頻度論統計学



# 序文

統計学は、現代社会において必要不可欠な学問分野の一つであり、様々な分野で広く利用されています。統計学は、データの集計、分析、解釈を行うことで、現象の本質を明らかにすることができます。そのため、ビジネス、医療、政治、科学、社会学などの分野で、統計学は不可欠な役割を果たしています。

本書は、統計学の基礎である頻度論に焦点を当て、平均、分散、不偏分散、分布、さらには相関 や回帰などの概念について詳しく説明します。また、これらの概念を用いて、データの中から有用 な情報を抽出するための手法についても解説します。

本書は、統計学初学者にとって最適な入門書として設計されています。統計学に不慣れな方でも、この本を通じて基本的な概念や手法を理解することができます。また、本書を通じて、統計学の重要性や実用的な側面についても理解することができます。

本書を読んでいただくことで、統計学の基礎概念についての深い理解を得ることができます。さらに、データの解釈や分析の手法を習得し、より高度な統計学の知識を獲得することができます。 本書を通じて、統計学に興味を持ち、その知識を応用することができるようになることを願っています。

# 第1章:統計学の基本的な概念

### 平均 (mean)

データの集合の平均値は、全てのデータを合計してデータ数で割った値です。例えば、データセット  $\{1, 2, 3, 4, 5\}$  の平均値は (1 + 2 + 3 + 4 + 5) / 5 = 3 となります。

### 中央値

中央値は、データの中央に位置する値のことです。つまり、データを小さい順に並べたときに、ちょうど真ん中に位置する値が中央値となります。中央値は、データの外れ値や異常値の影響を受けにくく、データの分布を表す上で平均よりも適切な代表値となります。

例えば、あるクラスのテストの点数が次のようになっている場合、

60, 70, 80, 90, 100

これを小さい順に並べると、

60, 70, 80, 90, 100

中央に位置する値は80なので、このデータの中央値は80となります。

### 最頻値:

最頻値は、データの中で最も頻度が高い値のことです。つまり、データを集計したときに、最も多く出現する値が最頻値となります。最頻値は、データの分布を表す上で有用な情報を提供することがあります。

例えば、あるクラスのテストの点数が次のようになっている場合、

60, 70, 80, 90, 90, 100

このデータを集計すると、90が最も頻度が高い値であるため、このデータの最頻値は90となります。

ただし、最頻値は、データの中で最も頻度が高い値が複数ある場合もあります。また、全てのデータが異なる値である場合は、最頻値は存在しません。

## 分散 (variance)

データの集合の分散は、各データと平均値の差の2乗の和をデータ数で割った値です。データが平均値からどの程度ばらついているかを表す指標であり、大きい値ほどデータがばらついていることを示します。例えば、データセット {1, 2, 3, 4, 5} の分散は ((1-3)^2 + (2-3)^2 + (3-3)^2 + (4-3)^2 + (5-3)^2) / 5 = 2 となります。

# 不偏分散 (unbiased variance)

不偏分散は、分散のバイアスを修正したものであり、標本分散とも呼ばれます。標本から推定された分散が、母集団分散と同じように偏らないようにするために用いられます。不偏分散は、分母に(n-1)を用いて計算されます。例えば、データセット {1, 2, 3, 4, 5} の不偏分散は ((1-3)^2 + (2-3)^2 + (3-3)^2 + (4-3)^2 + (5-3)^2) / (5-1) = 2.5 となります。

# 分布 (distribution)

分布は、データの値が取りうる範囲や出現頻度を表すグラフや表として表されるものであり、特定のデータセットがどのように分散しているかを視覚化するのに役立ちます。例えば、正規分布は、中心付近で頻繁に発生するが、中心から離れるほどデータが現れる頻度が低くなる分布です。他にも、一様分布、指数分布、ポアソン分布などがあります。

平均と最頻値が一致する分布の例としては、例えば以下のような場合が考えられます。

例1:テストの点数

100人の生徒が受けたテストの点数が次のようになっているとします。

60, 70, 80, 90, 100

この場合、平均は80であり、最頻値も80となります。つまり、この分布では平均と最頻値が一致しています。

### 例2:家賃

同じマンション内に住む人々の家賃が次のようになっているとします。

10万円、10万円、10万円、12万円、15万円

この場合、平均は11.4万円であり、最頻値も10万円となります。つまり、この分布でも平均と最頻値が一致しています。

一方、平均と最頻値が一致しない分布の例としては、以下のような場合が考えられます。

### 例1:月収

会社員の月収が次のようになっているとします。

20万円、25万円、30万円、40万円、50万円、100万円

この場合、平均は41.67万円であり、最頻値は30万円となります。つまり、この分布では平均と最 頻値が一致していません。

#### 例2:身長

あるグループの人々の身長が次のようになっているとします。

150cm, 155cm, 160cm, 170cm, 180cm, 190cm

この場合、平均は167.5cmであり、最頻値は存在しません。つまり、この分布でも平均と最頻値が一致していません。

# 第2章:単変量分析(相関と回帰)

相関と回帰は、単変量分析の代表的な手法の一つであり、類似点と相違点があります。

### 【類似点】

- ・相関と回帰は、変数間の関係を調べるために使用される手法である。
- ・相関と回帰は、どちらも2つ以上の変数を扱うため、多変量分析の前段階として用いられることがある。
- ・相関係数と回帰係数は、同じ単位で表されるため、比較が可能である。
- ・相関と回帰は、どちらも一方の変数から他方の変数を予測することができる。

### 【相違点】

- ・相関は、2つの変数の関係を調べる手法であり、回帰は、1つの変数が他の変数にどのような影響を与えるかを調べる手法である。
- ・相関は、変数間の関係を表す数値である相関係数を求める手法であり、回帰は、変数間の関係を表す回帰式を求める手法である。
- ・相関は、変数間の関係を表す数値を求めるため、どちらの変数が因果関係を示すかを判断することはできない。一方、回帰は、因果関係を仮定した上で回帰式を求めるため、因果関係を示すことができることがある。
- ・相関は、どちらの変数も独立変数として扱われるため、変数の関係にある種の均衡が保たれることがある。一方、回帰は、独立変数と従属変数があるため、変数の関係に均衡が保たれない場合がある。

以上を踏まえた上で、それぞれの定義と特徴と例を見てみましょう。

### 相関

相関とは、2つの変数の間にどの程度の関係があるかを表す指標であり、多くの場合、2つの変数が共に変化する傾向がある場合に高い相関が得られます。以下に相関の特徴を書き出します。

#### 【特徴】

- ・相関は、2つの変数の間に線形関係がある場合に用いることができる。
- ・相関係数は、-1から1の範囲の値をとり、絶対値が1に近いほど強い相関があるとされる。
- ・相関は、因果関係を示すものではなく、単に2つの変数の間にどの程度の関係があるかを示す 指標である。
- ・相関は、外れ値に弱いため、外れ値が存在する場合は注意が必要である。
- ・相関は、2つの変数のスケールに依存しないため、異なる尺度の変数でも比較が可能である。
- ・相関は、単変量分析のための手法であるため、複数の変数を同時に考慮する場合は、多変量解析が必要である。

相関は、統計解析において広く用いられる基本的な手法の一つであり、2つの変数の関係を理解する上で重要な役割を果たしています。以下に相関を用いた分析の例を2つ挙げます。

### 1. 年齢と犯罪率の相関分析

犯罪率が高い地域は、一般的に若年層が多いと言われています。そこで、ある地域の年齢構成と犯罪率との相関を分析することで、その地域の治安状況を把握することができます。例えば、20歳未満の人口割合が高い地域は、犯罪率が高い傾向があることが相関分析によって示されることがあります。

2. スマートフォン利用時間と睡眠時間の相関分析

近年、スマートフォンの普及に伴い、長時間のスマートフォン利用が睡眠不足につながるという報告がされています。そこで、スマートフォン利用時間と睡眠時間との相関を調べることで、スマートフォン利用が睡眠に与える影響を評価することができます。例えば、スマートフォン利用時間が長い人ほど、睡眠時間が短い傾向があることが相関分析によって示されることがあります。

これらの相関分析は、社会科学や医学などの分野でよく用いられます。

### 回帰

回帰とは、1つの変数(従属変数)を他の1つまたは複数の変数(独立変数)によって予測するための手法であり、2つ以上の変数の間の関係を分析するために用いられます。以下に回帰の特徴を書き出します。

### 【特徴】

- ・回帰は、線形回帰や非線形回帰など、様々な種類がある。
- ・回帰は、独立変数が従属変数に与える影響を定量的に評価することができる。
- ・回帰によって求められた予測値は、実測値に近づけるように最適化される。
- ・回帰には、最小二乗法や最尤法などの手法がある。
- ・回帰は、単回帰分析(1つの独立変数と1つの従属変数の関係を分析する)と重回帰分析(複数の独立変数と1つの従属変数の関係を分析する)がある。
- ・回帰は、因果関係を示すものではなく、単に予測値を求めるための手法である。

回帰分析は、多くの場面で用いられる基本的な手法の一つであり、データ分析や予測モデルの 構築において重要な役割を果たしています。以下に回帰を用いた分析の例を2つ挙げます。

### 1. 体重と身長の関係の回帰分析

身長がある程度伸びると、体重も増える傾向があります。そこで、身長と体重の関係を回帰分析することで、身長が分かれば体重を予測することができます。例えば、ある人の身長が180cmの場合、回帰分析によって予測される体重が80kgであることが分かります。

### 2. マーケティングの回帰分析

ある商品の売上が、価格や広告費などの要因によってどのように影響されるかを分析することがあります。これは、マーケティングの分野で広く用いられている手法です。例えば、ある商品の売上が増えるには、価格を下げるよりも広告費を増やす方が効果的であることが回帰分析によって示されることがあります。

これらの回帰分析は、ビジネスや経済学、健康管理などの分野でよく用いられます。

### コラム

Ken McAlinnとKota Matsuiと狸の置物の物語は、統計学の発展を物語る重要なエピソードとして、Twitter上で話題になりました。

物語の舞台は、英国の統計学者であるKen McAlinnが所有していた狸の置物です。ある日、McAlinnは自宅で置物を見ていたところ、その狸の置物が大量生産されていることに気づきました。そして、この狸の置物がどのように生産されているのか、そしてどのような特徴を持っているのかについて分析することにしました。

その分析に取り組むため、McAlinnは日本の統計学者であるKota Matsuiに協力を依頼しました。Matsuiは、日本国内で狸の置物を調査し、その特徴を分析することに成功しました。

しかし、その後に問題が発生しました。なぜなら、McAlinnが所有している置物と、Matsuiが調査した置物には、若干の違いがあったからです。具体的には、McAlinnの置物はより色彩が鮮やかで、より細かいディテールを持っていたのです。

この違いに対して、McAlinnは「これは優れた職人によって作られた置物だからだ」と主張しました。一方、Matsuiは「この違いは単なる生産工程の違いであり、統計的に有意な違いではない」と反論しました。

このように、McAlinnとMatsuiの間で狸の置物の分析について議論が交わされました。そして、この議論がTwitter上で話題になり、多くの人々が参加する大きな議論に発展しました。

この議論を通じて、統計学の発展について多くの人々が関心を持つようになりました。また、この議論が統計学の重要性を広く認知させるきっかけとなりました。