	結果が2択	[O • X]	結果は程度(複数)
	「合格不合格」	「あり・なし」	「数値」「点数」「金額」とか
どっちを 使う?			
使う?使わない?			t 検定
2択			
量は			

| 結果がバラバラだから、割合は無理!

結果の平均を比較してを検定

答え複数

Α	0)	13	
В	12	13	
C	10	0)	
D	10	7	
E	12	16	
F	11	7	
G	13	11	
Н	11	15	
	11	80	
J	11	11	
平均	11	11	

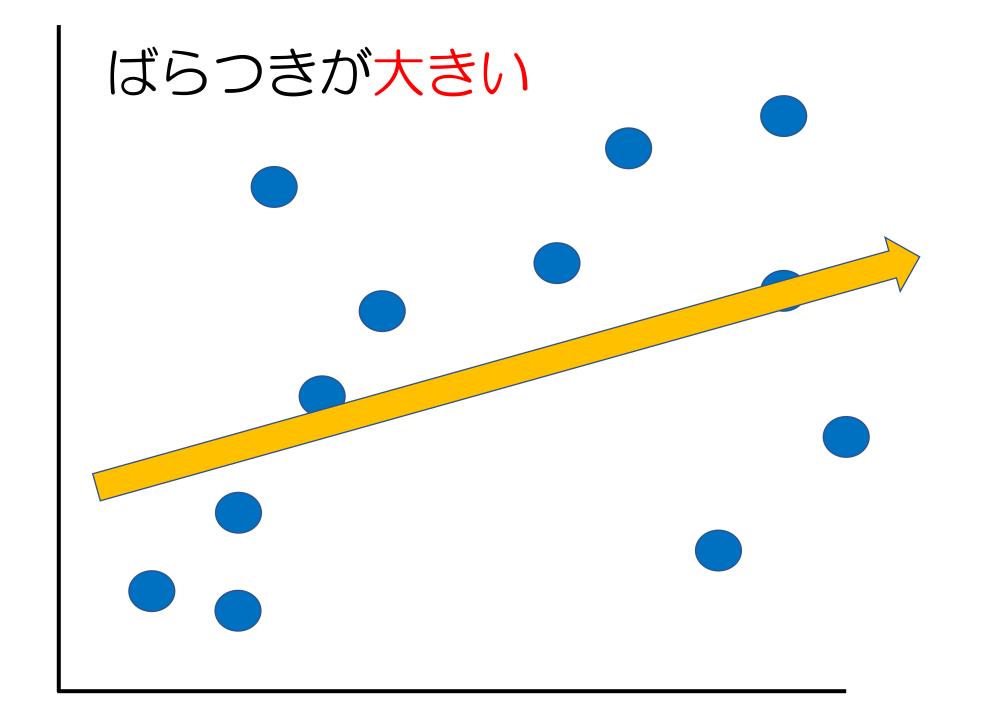
平均が同じでも

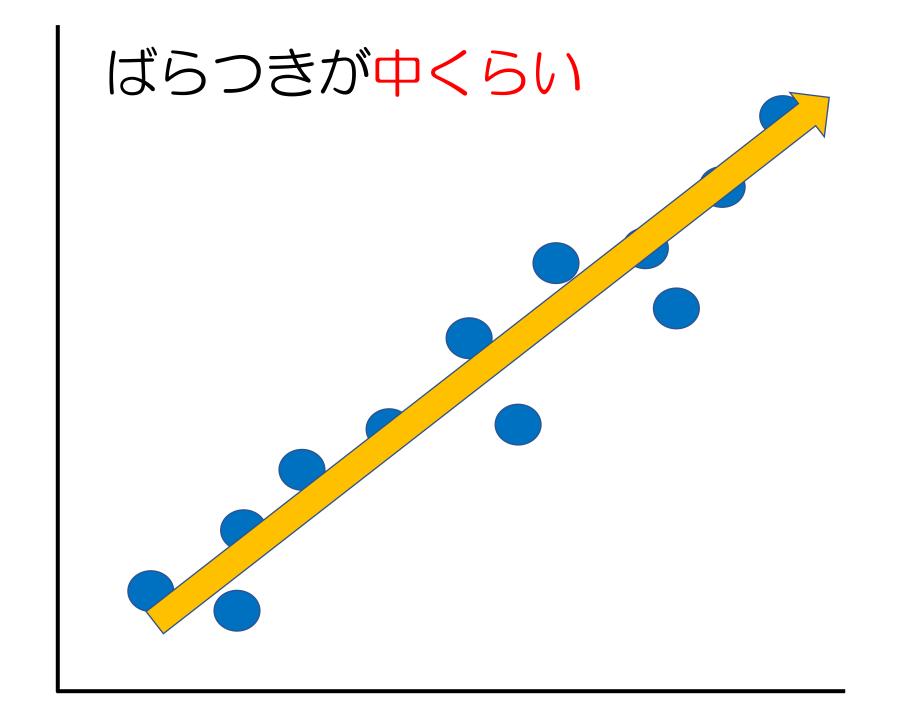
意味が違う!

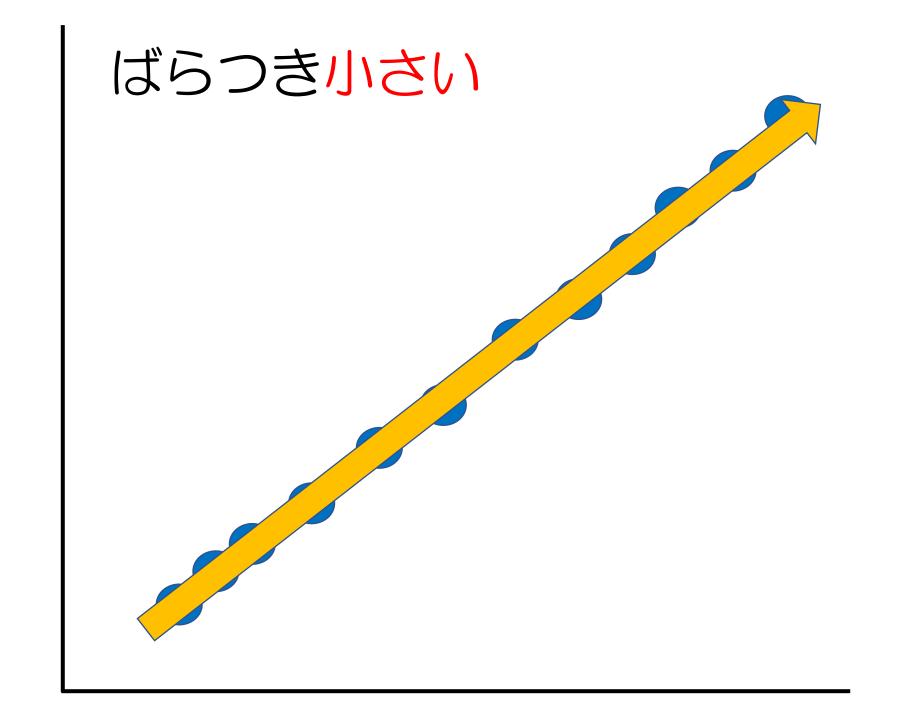
ばらつきは

左:小さい

右:大きい







必要な標本数を考えるとき重要なのが 「誤差」(ばらつき)の割合

「誤差が少ない」

⇒ 少ない数でもだいたい正確

「誤差が多い」

⇒ 少ない数では不正確



「効果量 d 」 介入によってどれだけ変化が起きたか? を表す値(通常-3 から 3 の範囲)

例えば、新薬を飲んだグループと、 飲まないグループを比較して <u>薬の効果がどの程度あったか</u> を数値で示したのが効果量



「効果量 d 」 介入によってどれだけ変化が起きたか? を表す値(通常-3 から 3 の範囲)

```
例えば、 d=0.1 効果がほとんどない (2つの間に差がない) d=1.0 効果がめっちゃあった(2つの間の差が大きい) d=10 ありえない結果 (データミスとか)
```

「効果量 d 」 介入によってどれだけ変化が起きたか? を表す値(通常-3 から 3 の範囲)

```
例えば、 d=0.1 効果がほとんどない (2つの間に差がない) d=1.0 効果がめっちゃあった(2つの間の差が大きい) d=10 ありえない結果 (データミスとか)
```

代表的な効果量の指標

「Cohen's d」「Hedges'g」「r: 相関係数」

Cohen's d

Hedges' g

効果量	記述的効果量	推定的効果量		
意味	平均値の差を標準偏差で割った値	平均値の差を偏分散の平方根で割った値		
特徴	サンプルの2つの差を記述した値	2つの母集団の差を推測する値		

必要な標本数を考えるとき重要なのが

「誤差」(ばらつき)の割合

とも言えるやんね!

「誤差が少ない」 効果が大きい!

⇒ 少ない数でもだいたい正確

「誤差が多い」

効果が小さい

⇒ 少ない数では不正確



#### 効果量 d とサンプルサイズの関係

効果量が大きい:信号が強い状態

小さなサンプルサイズ(少ない数)でも信号を捉えることができる

効果量が少ない:信号が弱い

大きなサンプルサイズじゃないと信号を捉えることができない

サンプルサイズが小さい

ただ、どっちも可能性!

効果量が少ないと検出できない、大きいと検出できる

サンプルサイズが大きい

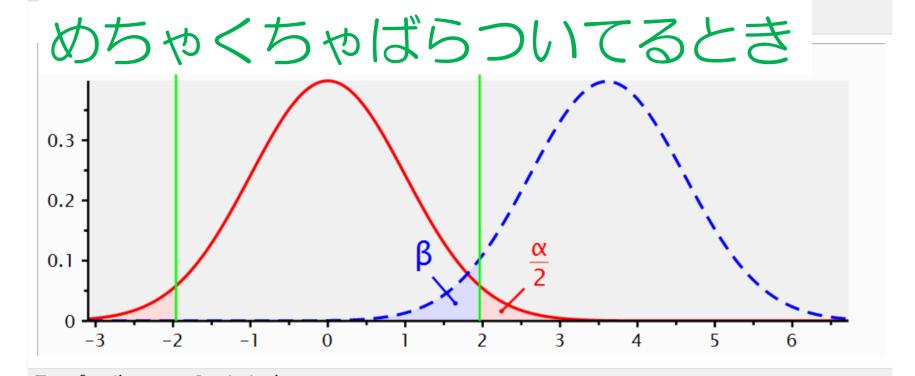
効果量が少なくても検出できる、大きいとより正確に検出できる

## 効果量 d とサンプルサイズの関係

計算はめんどくさいから、 統計ソフトに任せて

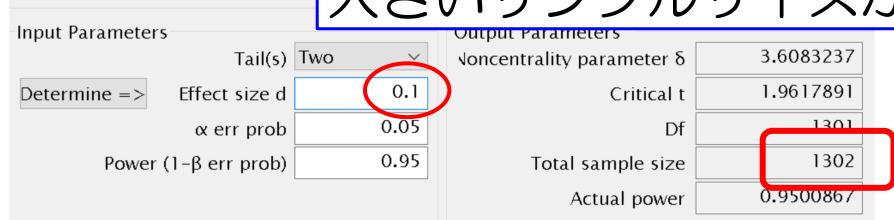


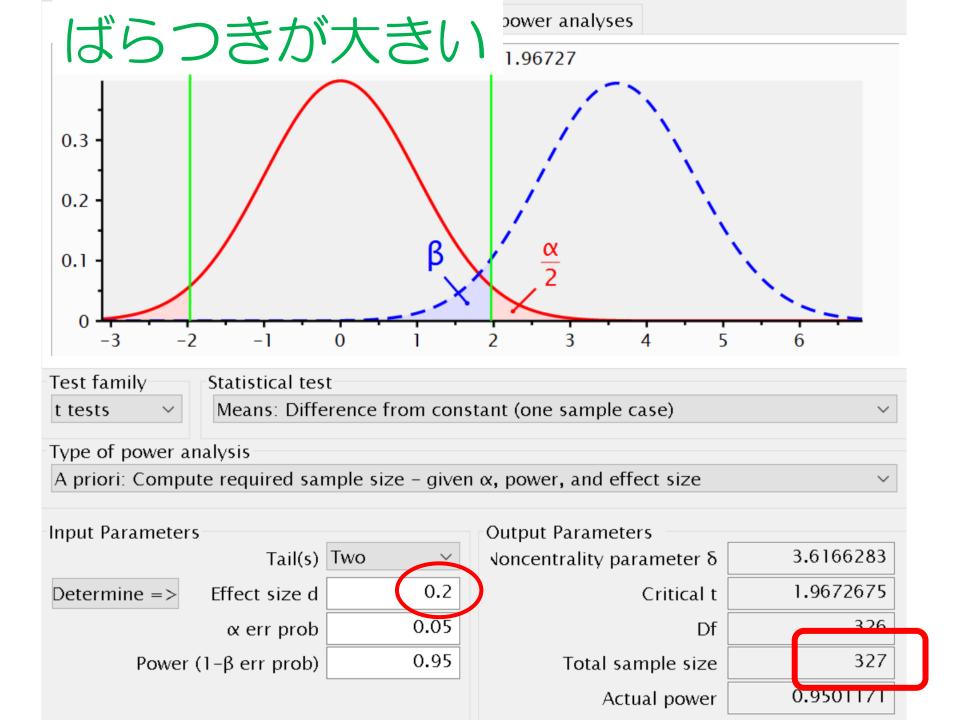


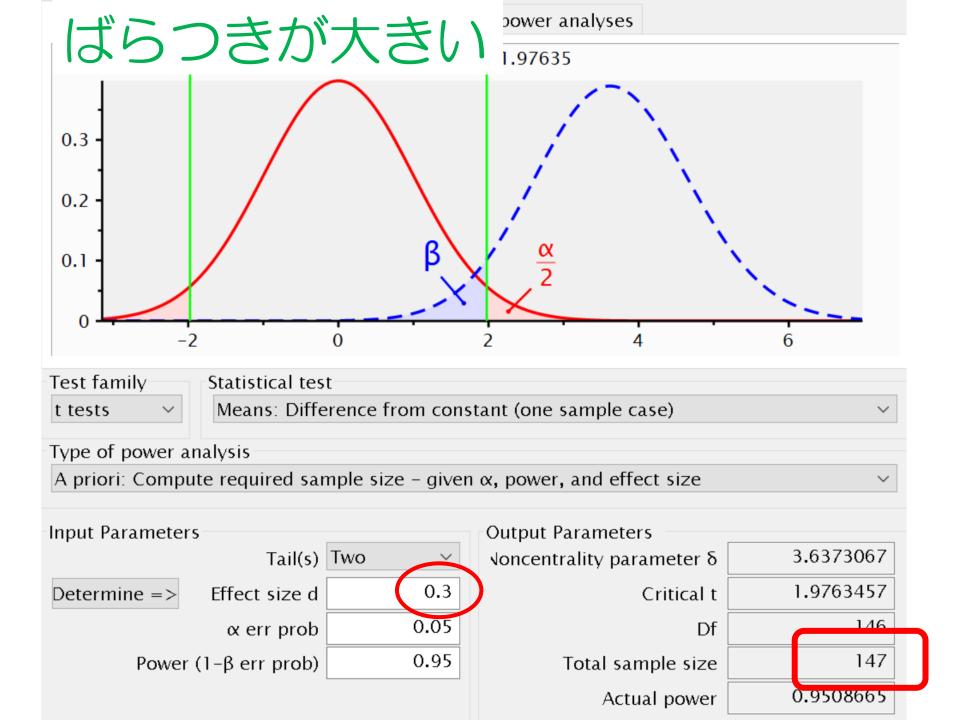


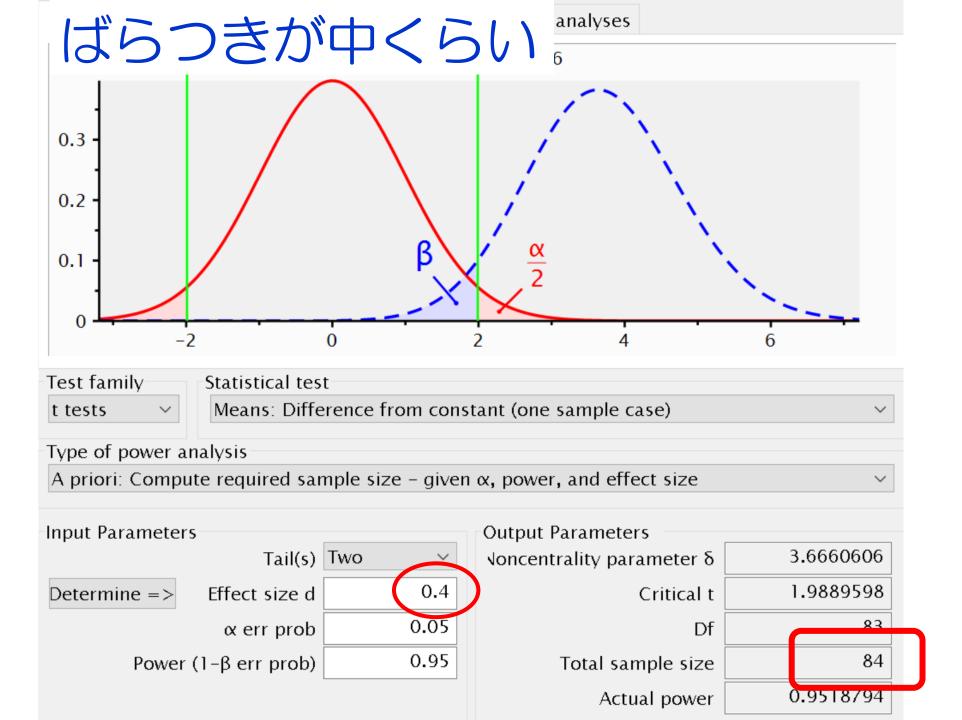
効果量 d が小さい ference from constant (one sample case)

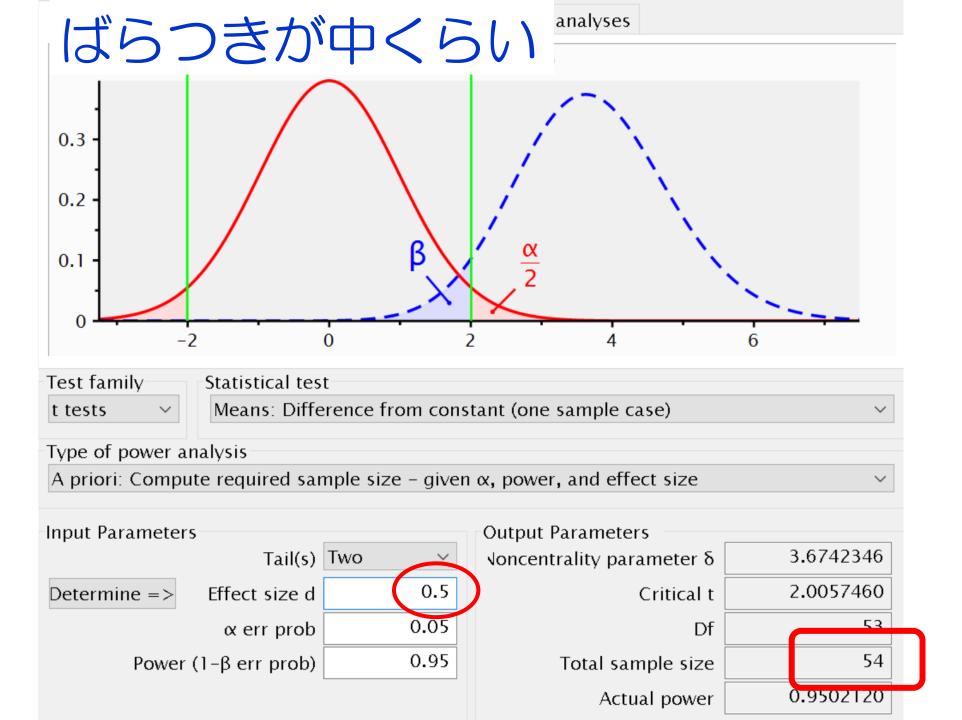
> A priori: Compute required samp ハサンプルサイズが必要! Output Parameters

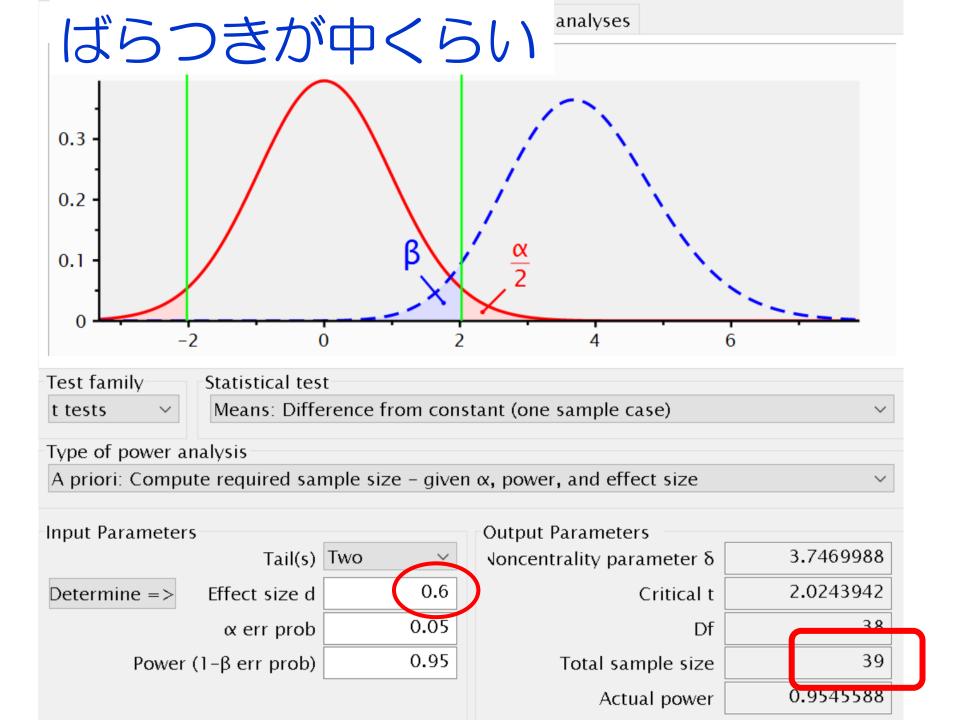


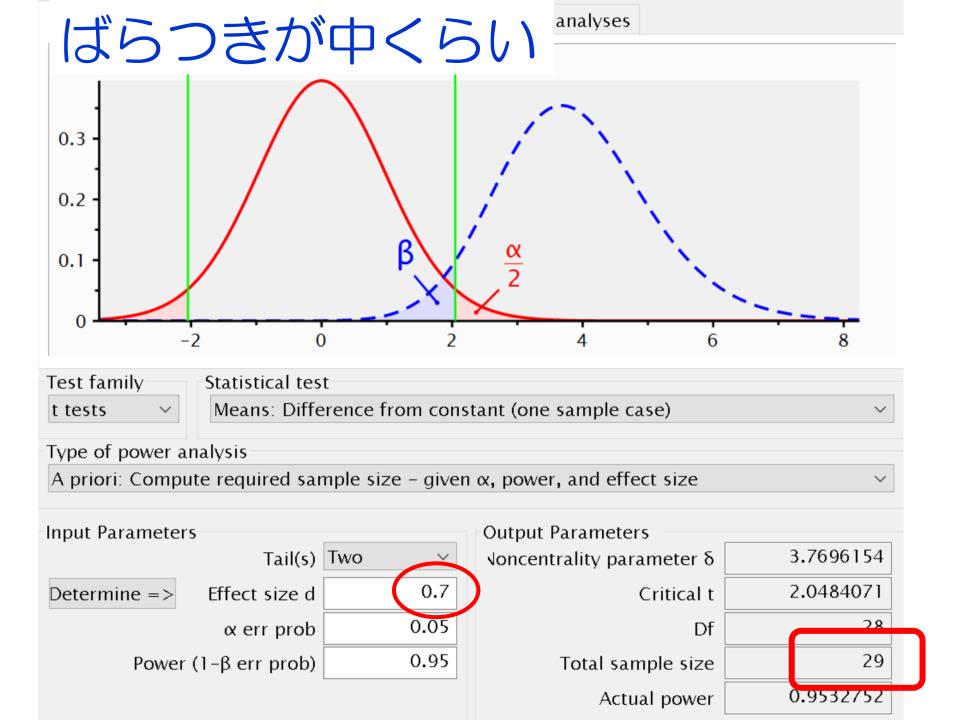


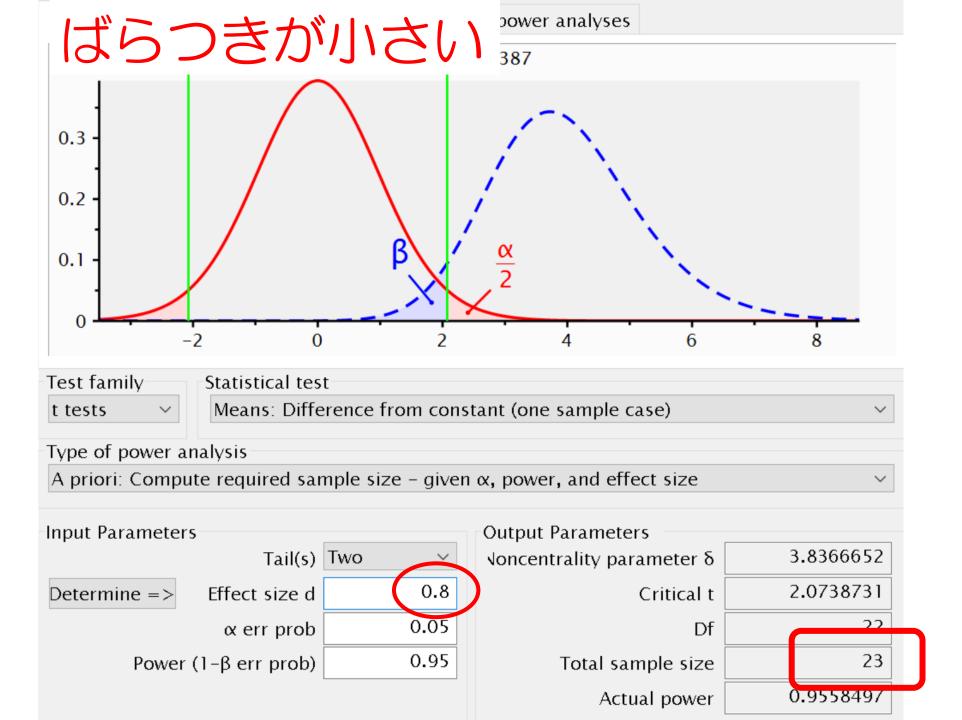


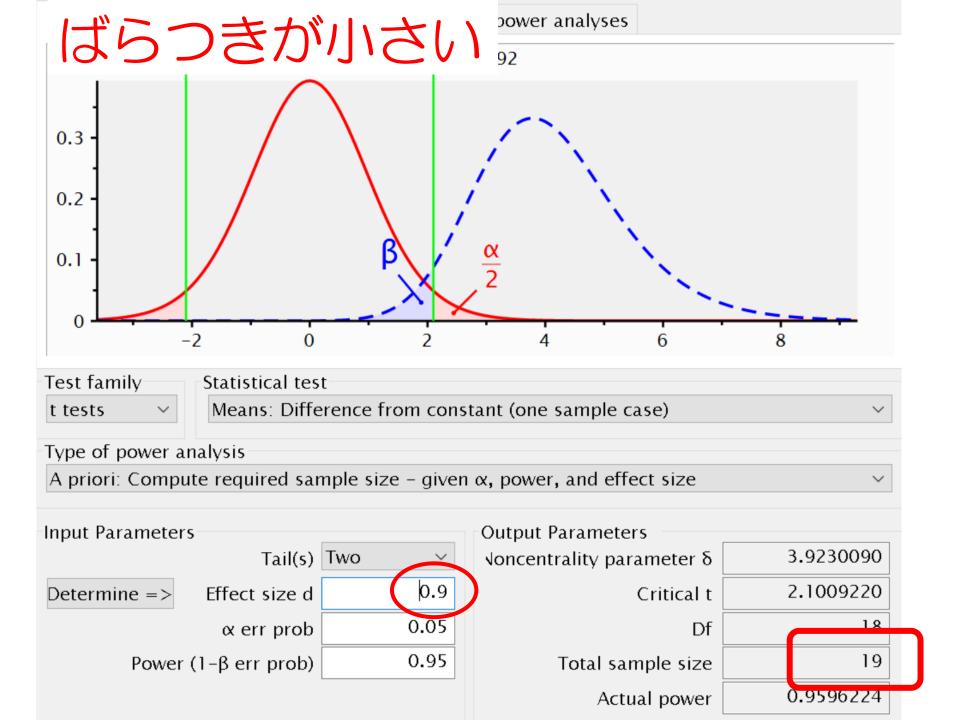


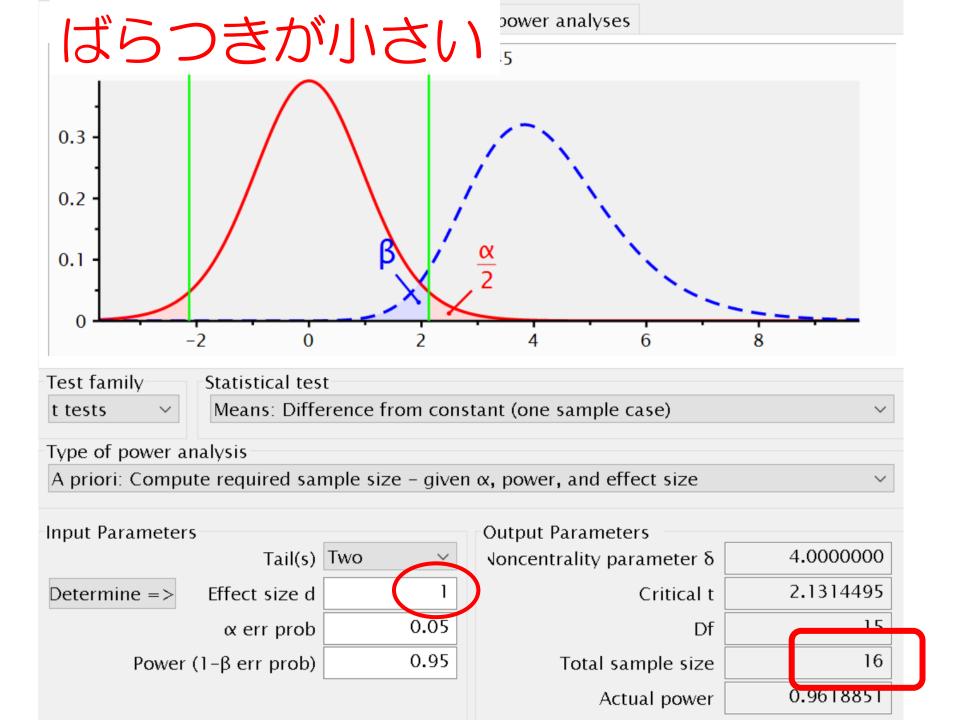


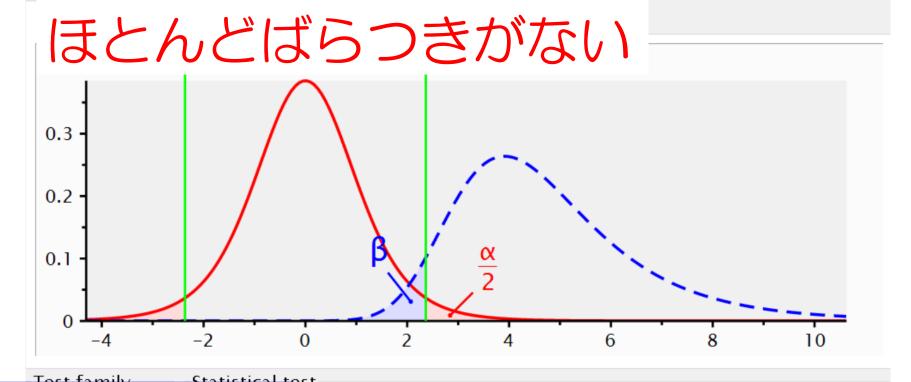










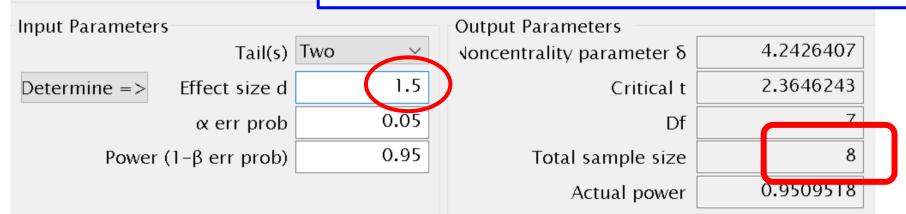


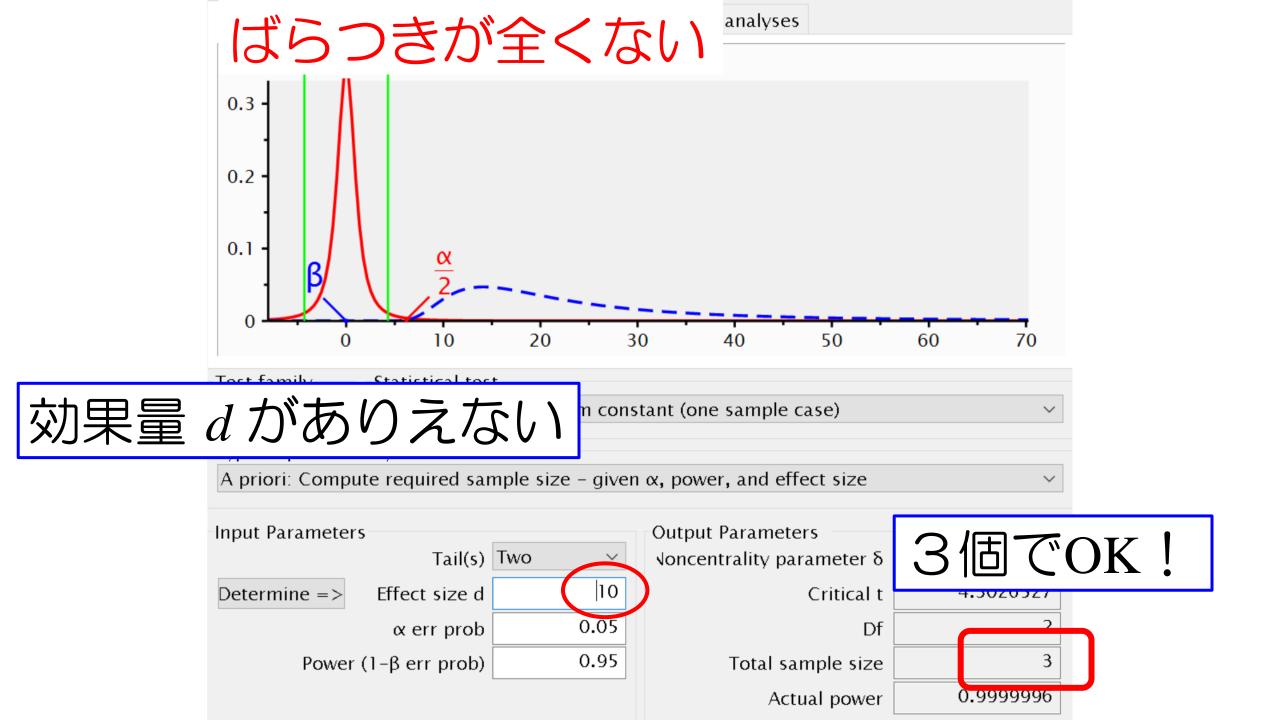
効果量 d が大きい

A priori: Compute required samp

ference from constant (one sample case)

## 小さいサンプルサイズでOK!





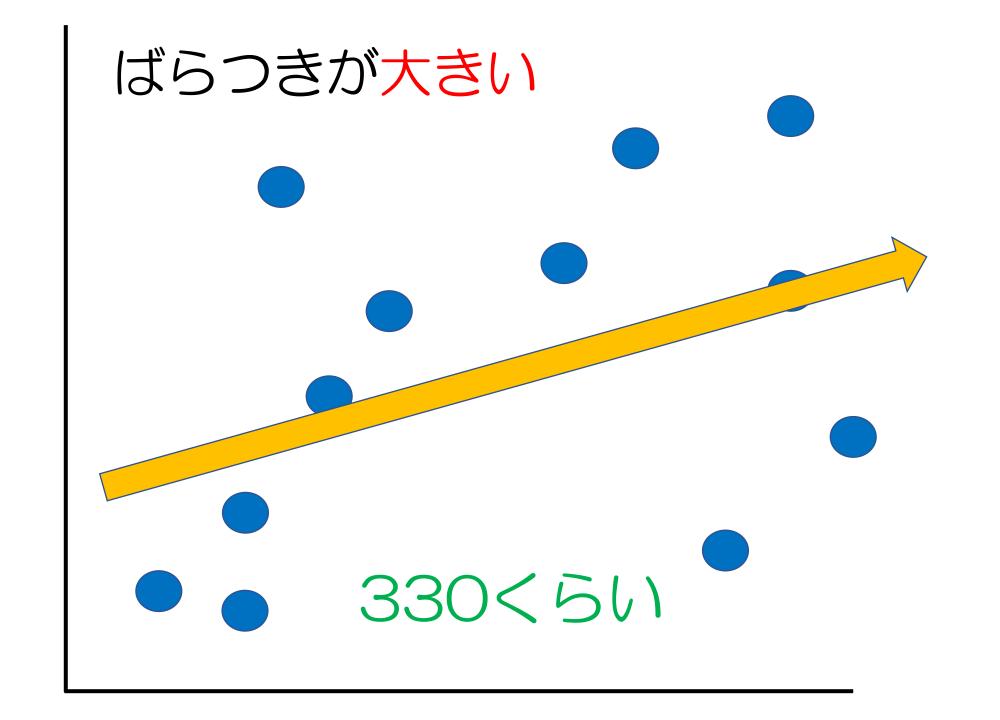
詳しい原理はわからんけど…

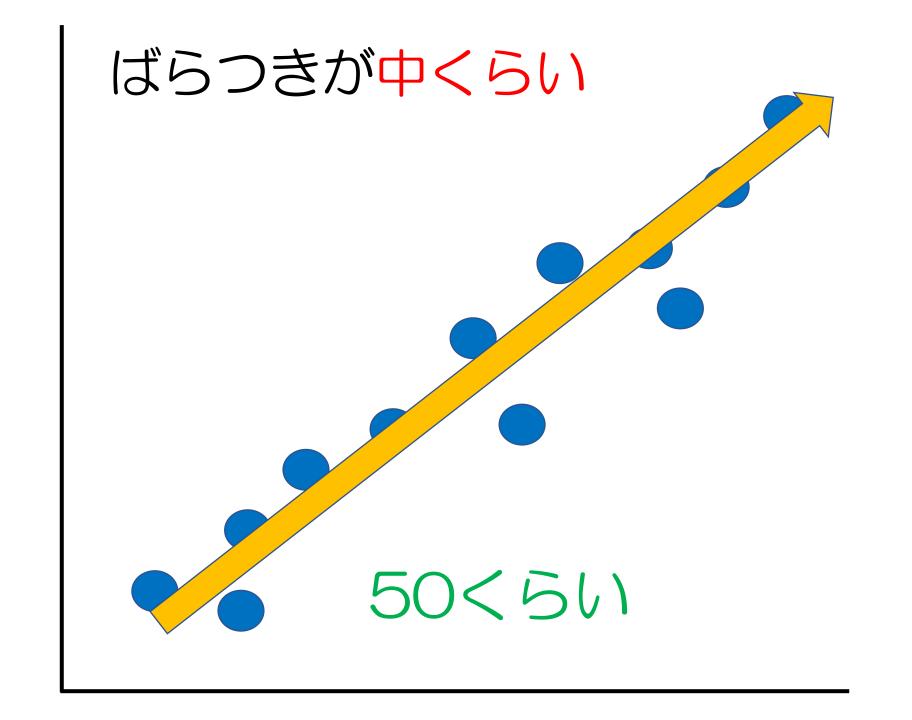
「ばらつき大」:330くらい

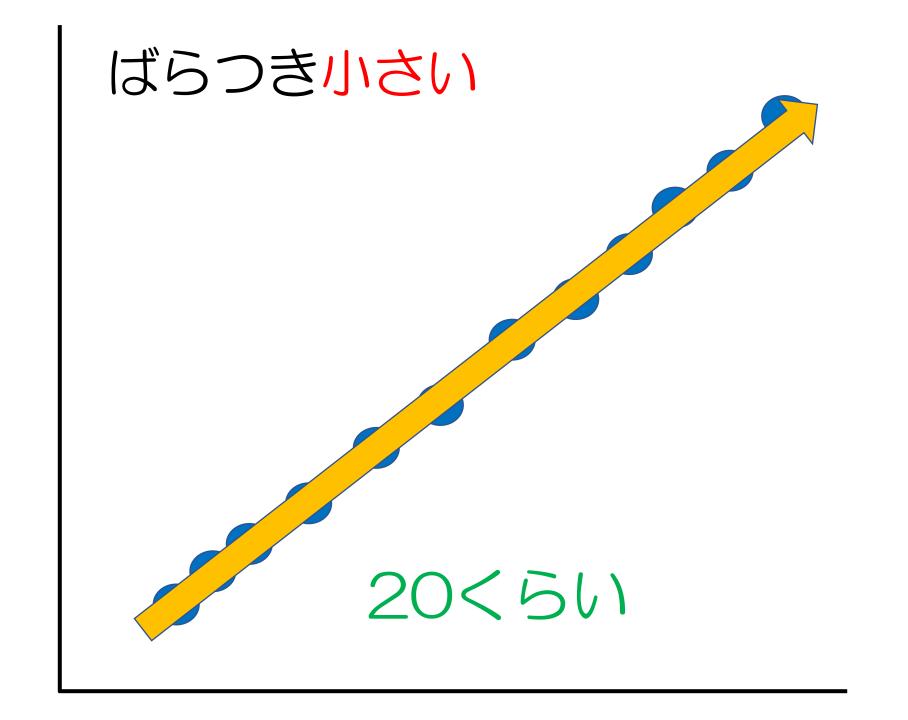
「ばらつき中」:50くらい

「ばらつき小」:20くらい

のところが大体の目安

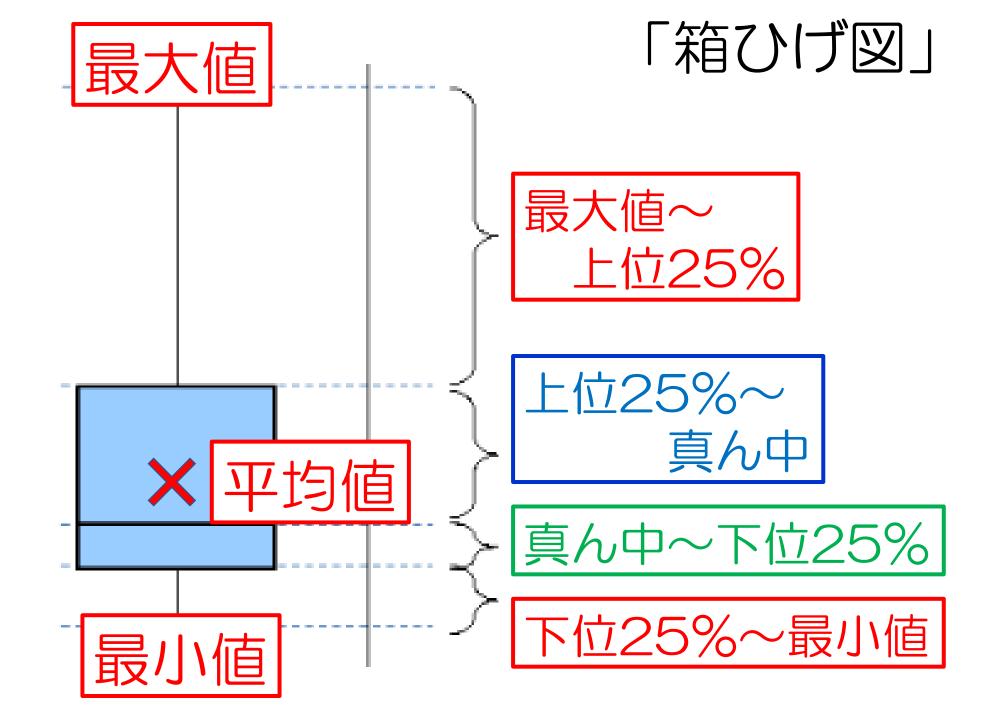




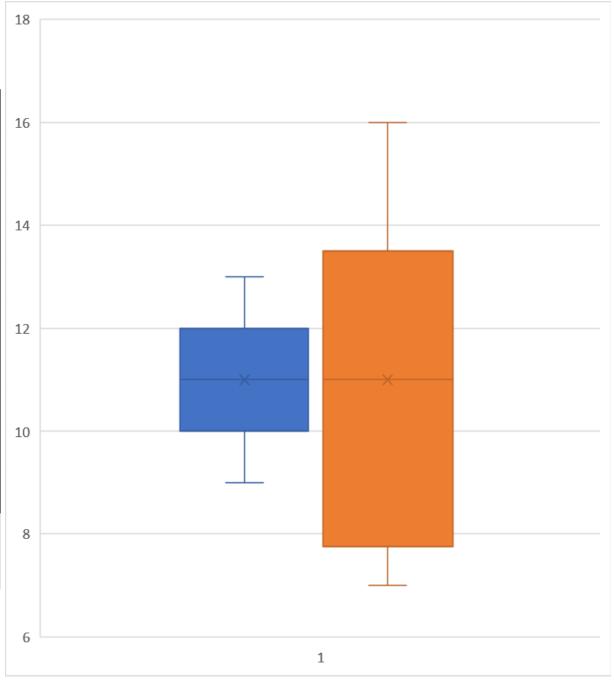


#### そのばらつきをどうやって表すのか?

## 「箱ひげ図」



A	9	13	
В	12	13	
B	10	9	
D	10	7	
E	12	16	
F	11	7	
G	13	11	
Н	11	15	
	11	8	
J	11	11	
平均	11	11	



# 仮説検定 (統計的検定)

1 カイ二乗検定

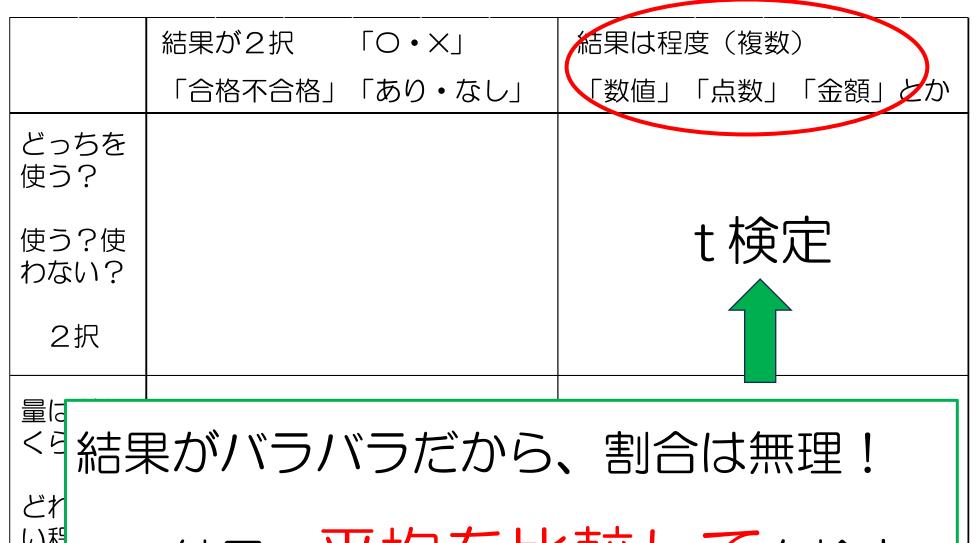


2 t 検定

3 回帰分析

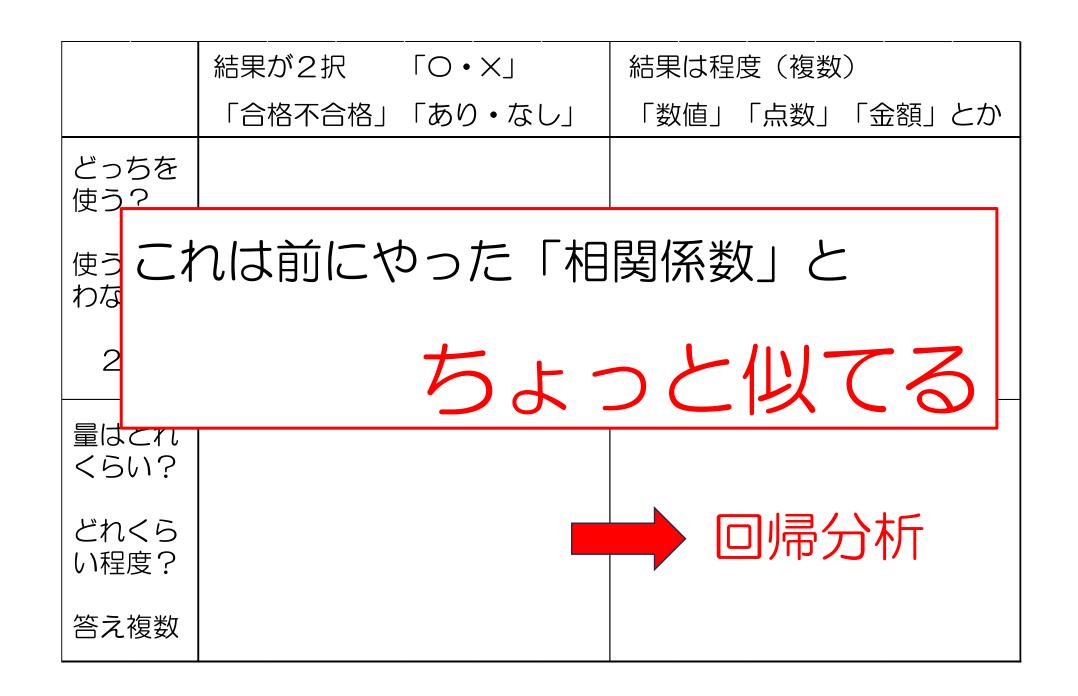
	結果が2択	[O • X]	結果は程度(複数)		
	「合格不合格」	「あり・なし」	「数値」	「点数」	「金額」とか
どっちを 使う?					
使う?使わない?	X二氢	t 検定			
2択					
量はどれくらい?					
どれくら い程度?	ロジステク	ィック回帰		回帰り	分析
答え複数					

	結果が2択	[O • X]	結果は程度(複数)			
	「合格不合格」	「あり・なし」	「数値」	「点数」	「金額」とか	
どっちを 使う?						
使う?使わない?	X二氢	<b>長検定</b>				
2択						
量はどれくらい?	赤と青のお皿で、食べるかどうかの					
どれくら い程度?			割包	うを核	<b>食定</b>	
答え複数						



結果の平均を比較してを検定

答え複数



結果が2択 「〇·×」

「合格不合格」「あり・なし」

結果は程度(複数)

「数値」「点数」「金額」とか

どっちを 使う?

使う?使わない?

2択

4

結果は変わる?



量はどれ くらい?

どれくらい程度?

答え複数

これによって



#### 「相関分析と回帰分析の違い」

相関係数がメイン。1対1のみ

XとYの関係性を見るだけ

XとYの関係式は出てこない

どっちが原因とかは考えない

相関分析

1対1でも、1対複数でも!

XからみたYの変化を見る

数式 Y = aX + b を立てて予測

Xが原因でYが結果として進める

回帰分析

## 回帰分析は意外と簡単!

いろいろ考えずにとりあえず

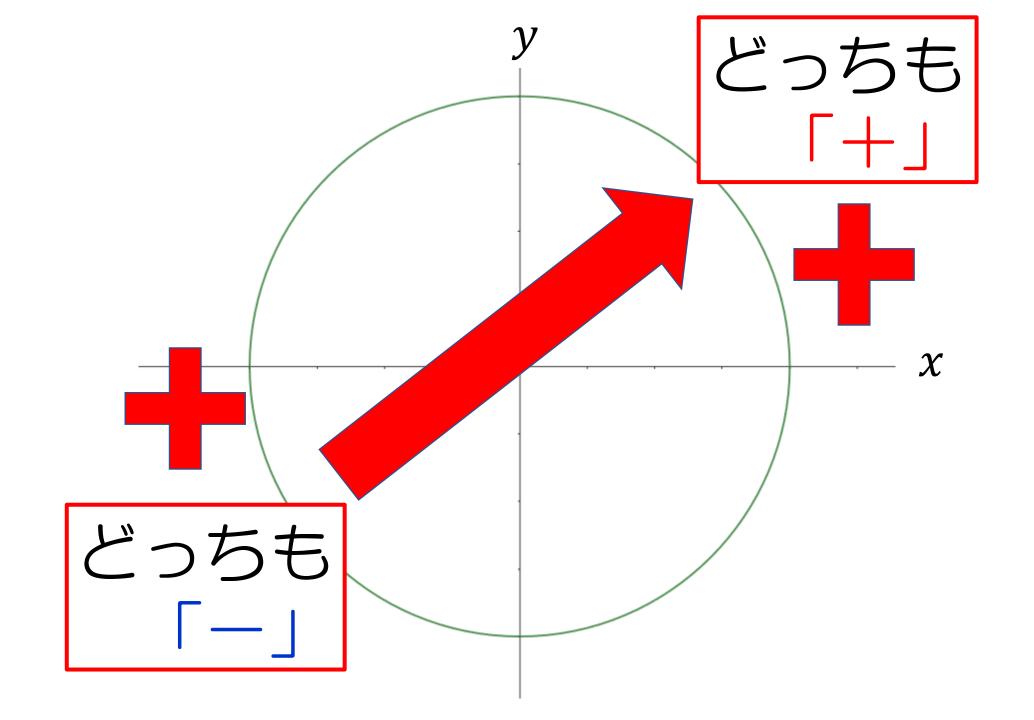
やってみよう!

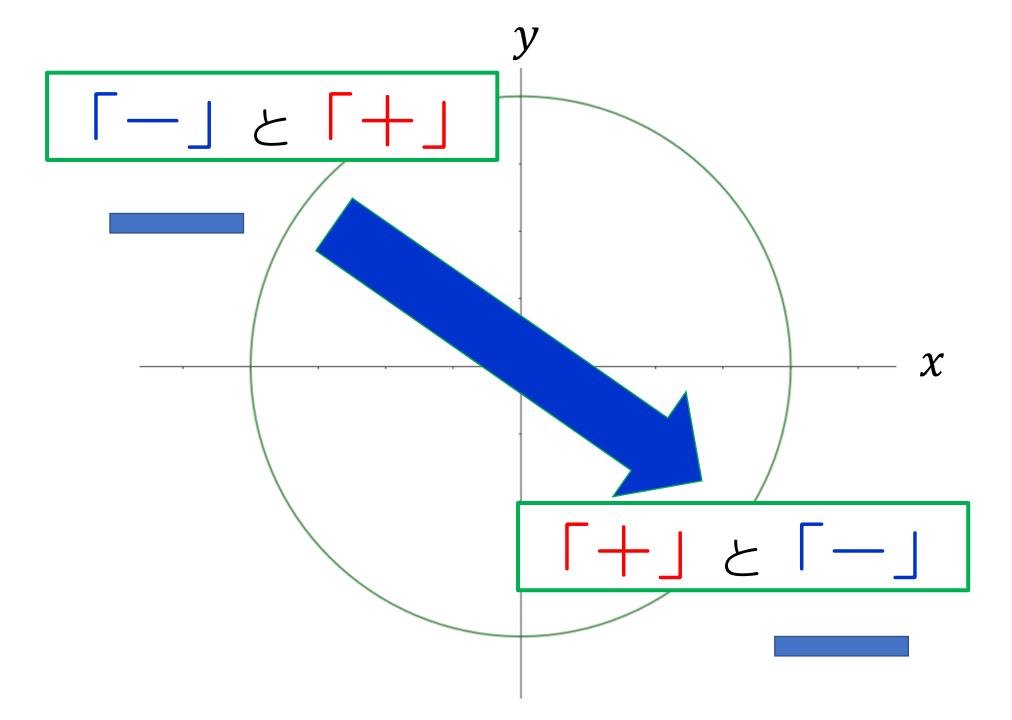
何となくわかってくると思う!



### 全部の空欄を埋めて

			平均との差	平均との差	身長の平均との差×
ID	身長	体重	身長	体重	体重の平均との差
1	163.1	62			
2	174.2	65.1			
3	175.1	67			
4	174.8	68.2			
5	169.4	67.2			
6	159.5	64.6			
7	152.0	62.5			
平均					
					(平均の差×平均の差)の平均





# これと同じようなのどっかで見た覚えない?



# 相関係数

2種類のデータの関係性の強さを

「一1から十1」

の間の値で表した数

「r」で表されることが多い

#### 相関の強さ

「Xが大きくなると

「Y」も大きくなる

$$0.7 \le r \le 1.0$$
 強い正の相関

$$0.4 \le r \le 0.7$$
 正の相関

$$0.2 \le r \le 0.4$$
 弱い正の相関

$$-0.2 \le r \le 0.2$$
 相関なし

$$-0.4 \le r \le -0.2$$
 弱い負の相関

$$-0.7 \le r \le -0.4$$
 負の相関

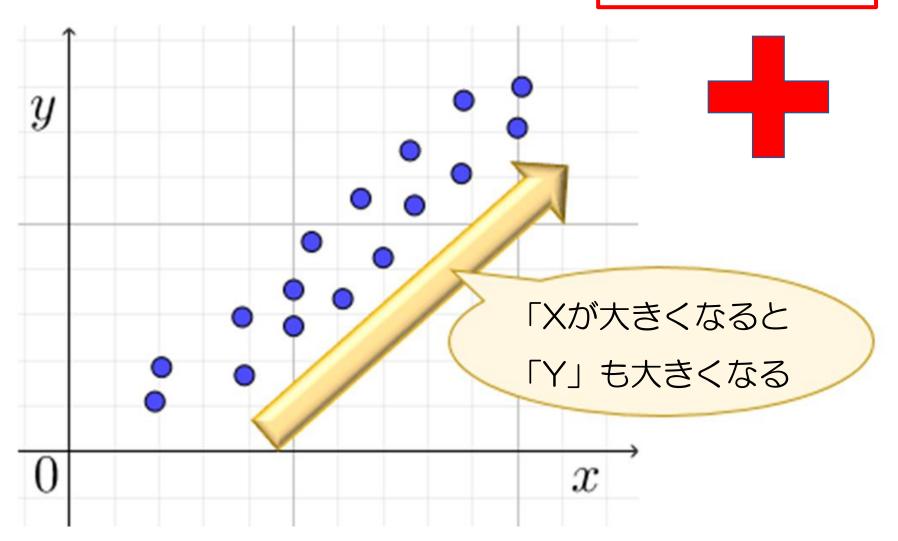
$$-1.0 \le r \le -0.7$$
 強い負の相関

「Xが大きくなると

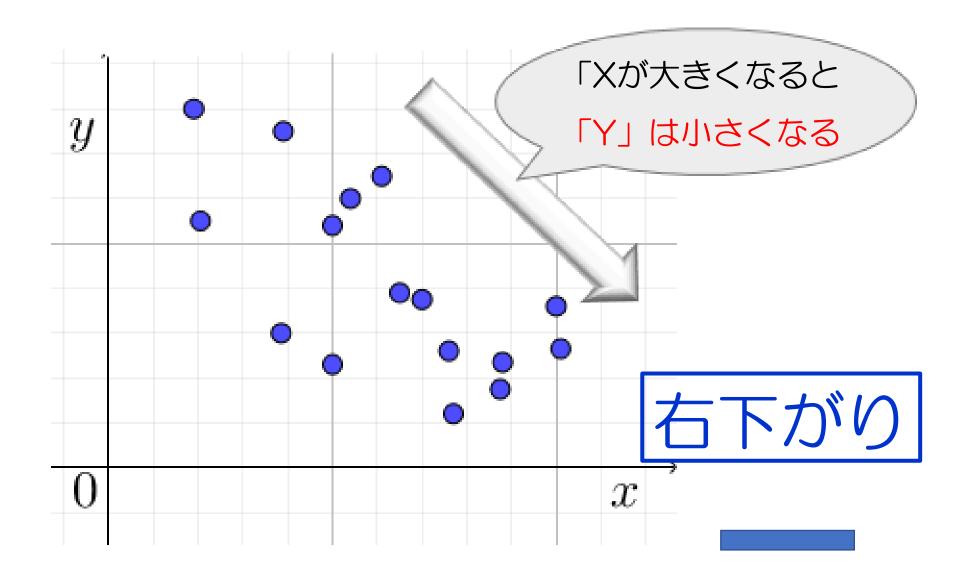
「Y」は小さくなる

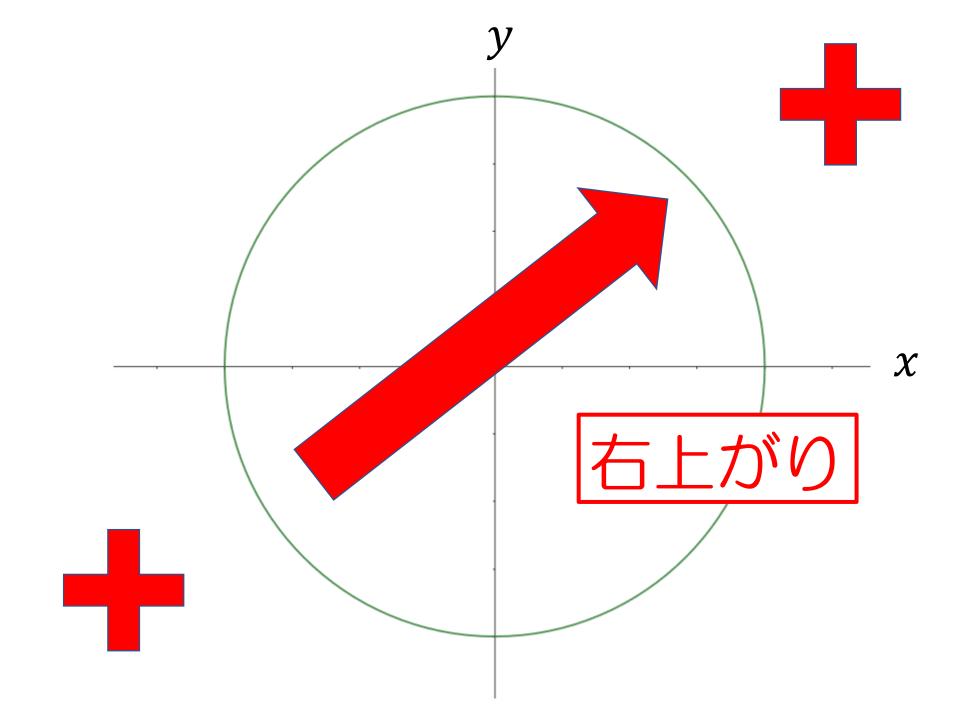
相関係数:0.94

# 右上がり

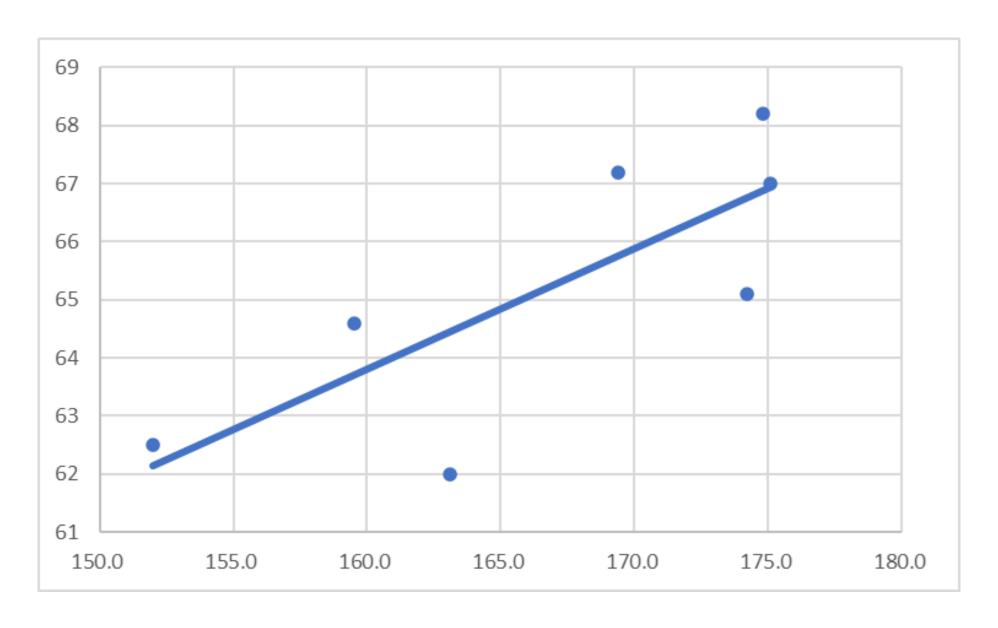


# 相関係数:一0.74





#### 散布図・近似式で見ても「右上がり」



共分散:「平均の差」×「平均の差」の事

今わかったのは

「右上がりの関係にある」

ってことだけ

### Excelだともっと簡単

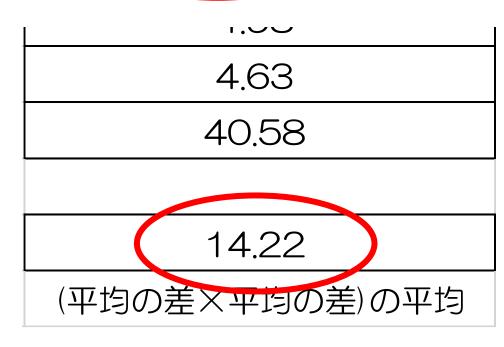


Excelの「データ」タブの中にある 「データ分析」 → 「共分散」

#### こんなの出るはず

	身長	体重
身長	68.63	
体重	14.22	4.876

# 一瞬!



同じようにやってみて!

シート「共分散2」

- 1 年齢と最高血圧
- 2 年齢と最低血圧
- 3 年齡と年収
- 4 最高血圧と最低血圧
- 5 最高血圧と年収
- 6 最低血圧と年収 の相関は?

#### こんなの出るはず

	年龄	最高血圧	最低血圧	年収
年齡	174.0			
最高血圧	121.1	113.1		
最低血圧	70.9	52.2	38.8	
年収	1491.4	796.8	604.7	17678.0

「共分散」は向きがわかるだけ…

どれくらい強い関係か 調べるために行うのが

「相関係数の有意差検定」

#### 前に「相関係数」を調べたときは

	強い正の相関	$0.7 \le r \le 1.0$
強さを	正の相関	$0.4 \le r \le 0.7$
出して	弱い正の相関	$0.2 \le r \le 0.4$
おおまかに	相関なし	$-0.2 \le r \le 0.2$
0000007510	弱い負の相関	$-0.4 \le r \le -0.2$
調べた	負の相関	$-0.7 \le r \le -0.4$
	強い負の相関	$-1.0 \le r \le -0.7$

## 相関係数の基準値

$$r_0 = \sqrt{\frac{4}{n+2}}$$

n:データ数

# シート「相関1」 AグループとBグループの相関係数

それぞれ求めてみて

# A身長体重身長1体重0.7781

В	身長	体重
身長	1	
体重	0.922	1

基準値は 
$$r_0 = \sqrt{\frac{4}{7+2}} = 0.667$$

Α	身長	体重	В	身長	体重
身長	1		身長	1	
体重	0.778	) 1	体重	0.922	1

「0.667」より大きいとき

統計的に「相関関係にある」と言える

シート「相関2」

平均気温とビール消費量の相関係数 を求めてみて

基準値は 
$$r_0 = \sqrt{\frac{4}{12+2}} = 0.535$$

	平均気温	ビール消費量
平均気温	1	
ビール消費量	0.510	1

「0.535」より小さいから

統計的に「相関関係にない」と言える

ここでちょっと考えてみる

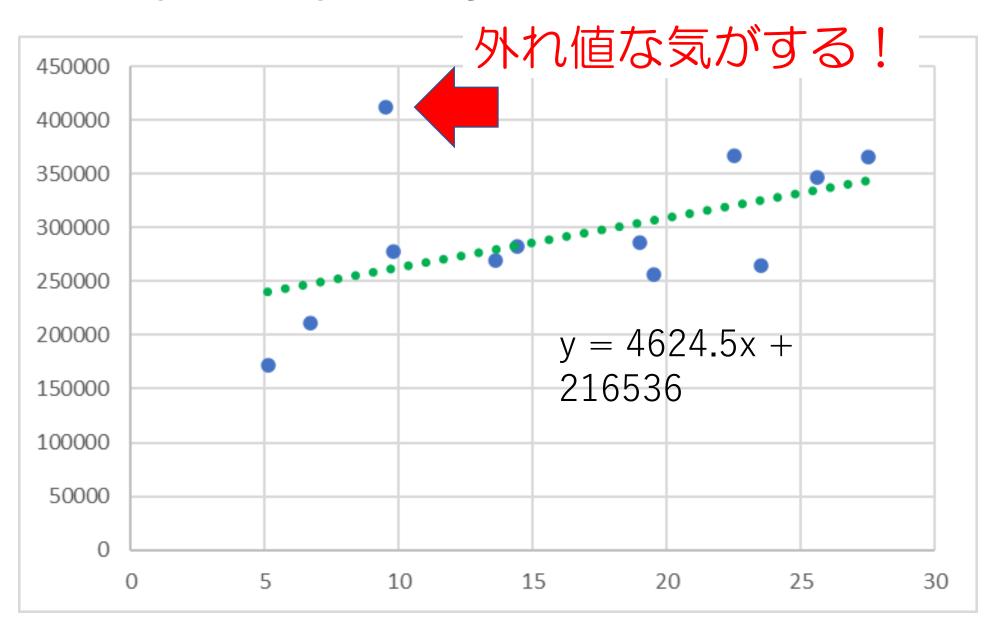
前やった時も「外れ値」があると

相関係数がむちゃくちゃになった。



「散布図」を使って確認してみる

#### こんなのになるはず



## 「外れ値」のデータを抜いて

もう一度相関係数を調べる

基準値は 
$$r_0 = \sqrt{\frac{4}{11+2}} = 0.555$$

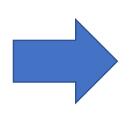
$$r_0 = \sqrt{\frac{4}{12+2}} = 0.555$$

	平均気温	ビール消費量
平均気温	1	
ビール消費量	0.829	1

「0.555」より大きいから

統計的に「相関関係にある」と言える

# 「外れ値」以外を調べると 相関関係にあった



「外れ値」になった理由が 推測できるかも…

「12月は忘年会などで気温に関係なく ビールの消費量が増えたのではないか」 と推測できる!

# 相関関係の調べ方はわかったけどいつになったら

## 「回帰分析」が始まるのか?

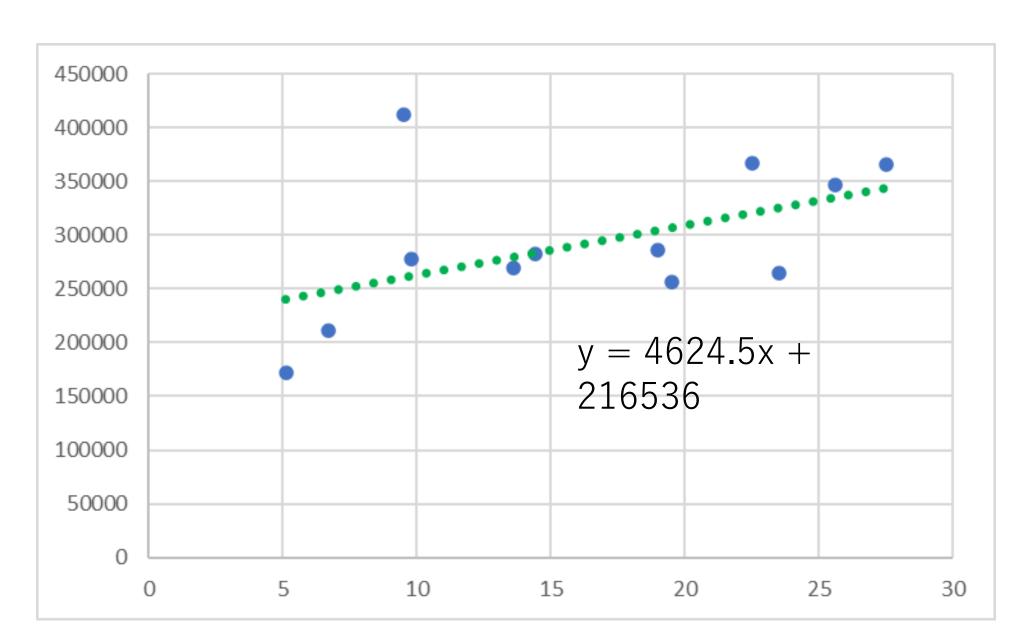
 2択

 量はどれ くらい?

 どれくらい程度?

 答え複数

## 実はもう「回帰分析」をやってる!



#### も一回さっきの

「平均気温とビール消費量」に戻って 「回帰分析」をやってみよう



#### 「回帰分析」

入力Y範囲:ビール消費量(11月まで)

入力X範囲:平均気温 (11月まで)

☑ ラバル( <u>L</u> ) ☑ 有意水準( <u>O</u> ) 95	定数に 0 を使用( <u>Z</u> ) %
出力オプション <ul> <li>一覧の出力先(<u>S</u>):</li> <li>新規ワークシート(<u>P</u>):</li> <li>新規ブック(<u>W</u>)</li> </ul>	ここをチェック
残差	<ul><li></li></ul>

#### こんなのになるはず

	帰統計		
重相関 R	0.829		
重決定 R2	0.688		
補正 R2	0.653	<u> </u>	
標準誤差	35696.891	さっきと	
観測数	11.000		,
分散分析表			
	自由度	変動	
回帰	自由度 1.000	変動 25271620194.279	
回帰残差			
	1.000	25271620194.279	
残差	1.000	25271620194.279 11468412424.267	
残差	1.000	25271620194.279 11468412424.267	
残差	1.000 9.000 10.000	25271620194.279 11468412424.267 36740032618.546	

#### 「相関分析と回帰分析の違い」

相関係数がメイン。1対1のみ

XとYの関係性を見るだけ

XとYの関係式は出てこない

どっちが原因とかは考えない

相関分析

1対1でも、1対複数でも!

XからみたYの変化を見る

数式 Y = aX + b を立てて予測

Xが原因でYが結果として進める

回帰分析

# 「相関分析と回帰分析の違い」

# 回帰分析

1対1でも、1対複数でも!

XからみたYの変化を見る

数式 Y = aX + b を立てて<u>予測</u>

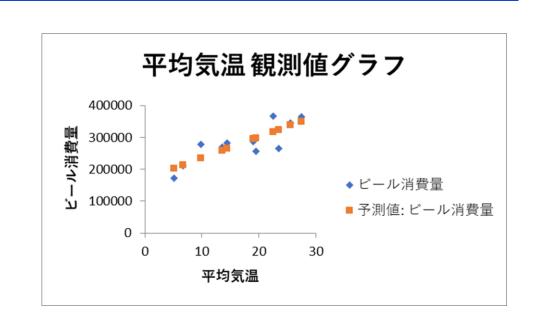
Xが原因でYが結果として進める

# 「回帰分析」 新たなデータXからデータY

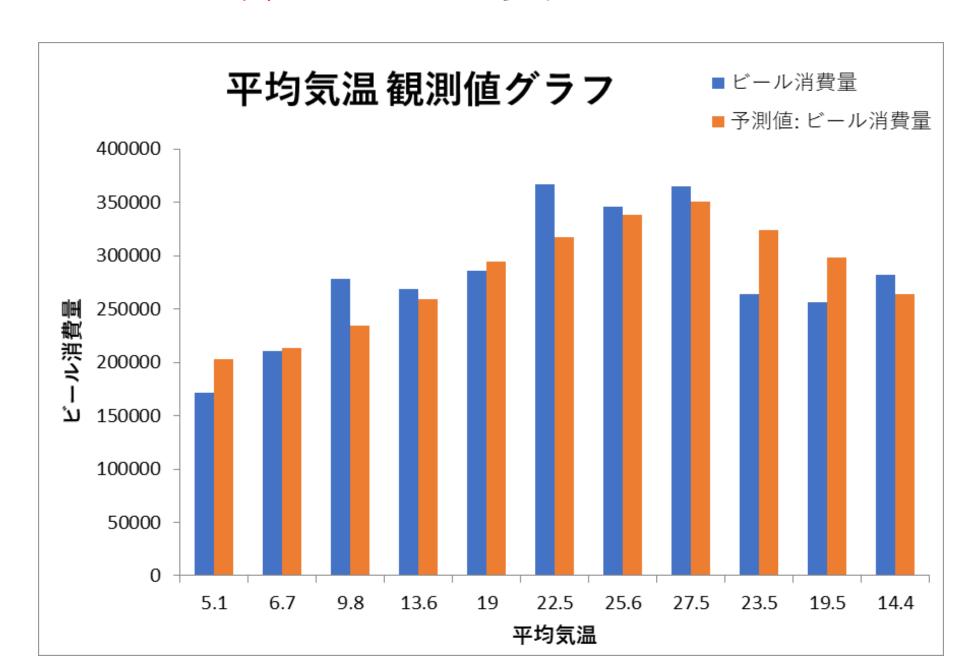
# を予測する

$$Y = aX + b$$

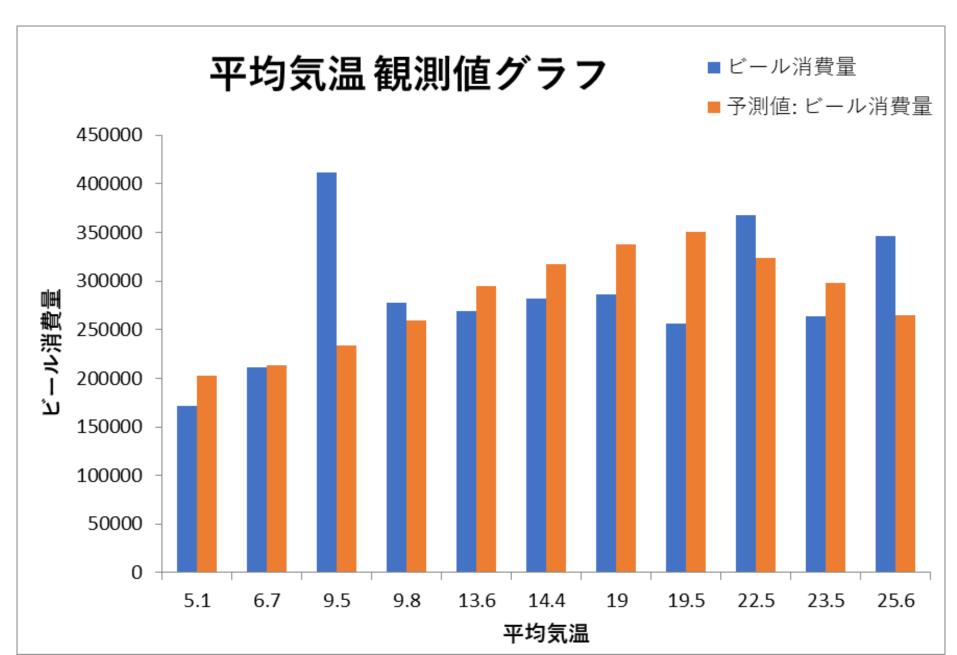
横のほうに出てる グラフ を見てみて



#### グラフを棒グラフに変えると



### 元データを昇順にすると完成!



# 「相関分析と回帰分析の違い」

# 回帰分析

1対1でも、1対複数でも!

XからみたYの変化を見る

数式 Y = aX + b を立てて予測

Xが原因でYが結果として進める

「回帰分析」 新たなデータXからデータY を予測する

$$Y = aX + b$$

次はこの式の意味を考える! さっきの表に戻って

### 次に見るのはココ

回帰統計		
重相関 R	0.829	
重決定 R2	0.688	
補正 R2	0.653	
標準誤差	35696.891	
観測数	11.000	
分散分析表		
	自由度	変動
回帰	自由度 1.000	変動 25271620194.279
回帰残差		2 3 2 7 2
	1.000	25271620194.279
残差	1.000	25271620194.279 11468412424.267
残差	1.000	25271620194.279 11468412424.267
残差	1.000 9.000 10.000	25271620194.279 11468412424.267 36740032618.546

Yを「ビール消費量」 Xを「平均気温」にしたから

これを代入すると…

$$Y = aX + b$$

# $Y = 6573 \times X + 169580$

Y:「ビール消費量」

X:「平均気温」

#### 予測できること

・平均気温が1°上がると

消費量が6573増える

平均気温が0°では消費量は169580

# 「相関分析と回帰分析の違い」

# 回帰分析

1対1でも、<u>1対複数でも!</u>

XからみたYの変化を見る

数式 Y = aX + b を立てて予測

Xが原因でYが結果として進める

さっきやったのは

「(単)回帰分析」(1対1)

問題設定は

「家賃」と「面積・築年数・時間」の関係を調べたい

これをメインに調べたい

#### 「データ分析」 ⇒ 「重回帰分析」

入力Y範囲:家賃

入力X範囲:面積 · 築年数 · 時間全部

(	☑ ラベル( <u>L</u> ) ☑ 有意水準( <u>O</u> )	□ 定数に 0 を使用( <u>Z</u> ) 95 %
	出力オプション <ul> <li>一覧の出力先(<u>S</u>):</li> <li>新規ワークシート(<u>P</u>):</li> <li>新規ブック(<u>W</u>)</li> </ul>	\$E\$2
	残差 □ 残差( <u>R</u> ) □ 標準化された残差( <u>T</u> )	<ul><li>             □ 残差ケラスの作成(<u>D</u>)         </li><li>             ☑ 観測値グラフの作成(<u>I</u>)         </li></ul>

	係数
切片	40900
面積	642.32
築年数	-213.3
駅までの時間	-507.3

#### これになった?

$$Y = 642X_1 - 213X_2$$
$$-507X_3 + 40899$$

Y :家賃  $X_1$ :面積

 $X_2$ : 築年数  $X_3$ : 時間

# $Y = 642X_1 - 213X_2$ $-507X_3 + 40899$

Y :家賃  $X_1$ :面積

 $X_2$ : 築年数  $X_3$ : 時間

家賃は、基本が40899円で

面積が1増えると642円上がり

築年数が1増えると213円安くなり

時間が1増えると507円安くなる

#### ここから先で大事な話

### 「パラメトリック」

#### 「ノンパラメトリック」

特徴	パラメトリック	ノンパラメトリック
データの分布	特定の分布を仮定 (通常は正規分布)	特定の分布を仮定しない
分析対象	平均、標準偏差など	順位
代表的な手法	t検定	マン・ホイットニーのU検定 ウィルコクソンの符号順位検定
メリット	小さな差も検出できる	適用範囲が広い

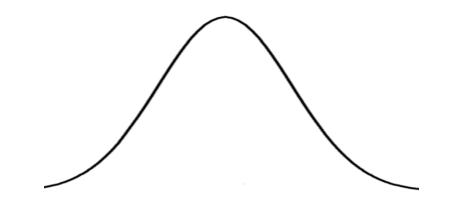
# 「パラメトリック」

母集団の分布(正規分布とか)

が事前にわかってるとき

この分布に従ってパラメータを

統計的に推測する(t検定とか)



## 「ノンパラメトリック」

#### 母集団の分布がわからないとき

分布がわからないとして推測 (どんな分布でもOK)

・母集団に正規性がなく、

サンプルサイズが小さいとき

• 極端な外れ値があるけど無視できないとき

- マンホイットニーのU検定、
- ウィルコクソンの符号順位検定

「使い分け」

まず、パラメトリックが無理かを考える



ダメなときはノンパラメトリックな手法を用いる

(パラメトリックな時に、 ノンパラメトリックを使っても いいけど検出力は不利になる) これで

「医療統計」終わり

テストはPCを使ってやります

授業でやった

記述統計・推測統計の中で

「提出」してないやつを重点的に!

