


2019年度秋学期 統計学 第7回  
データの関係を知る(2)—回帰分析

浅野 晃  
関西大学総合情報学部



回帰分析とは🤔

## 回帰分析とは

多変量データがあるとき  
ある変量の変化を他の変量の変化で  
[説明] する方法

説明? 🤔

## 回帰分析とは

緯度と気温のデータを例にとると

### 相関分析

「緯度が上がると、気温が下がる」という  
傾向がはっきりしている

### 回帰分析

「緯度が上がるから気温が下がる」と考える  
緯度が1度上がると、気温が○℃下がる

## 回帰分析とは

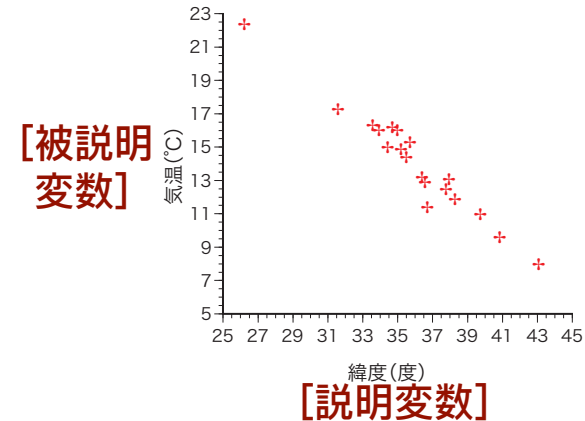
緯度が上がるから気温が下がると考える  
緯度が1度上がると、気温が○℃下がる

各都市の気温の違いは、緯度によって決まっているという [モデル] を考える

統計学では、  
気温 (のばらつき) は、緯度によって [説明] されるという

## 説明変数・被説明変数

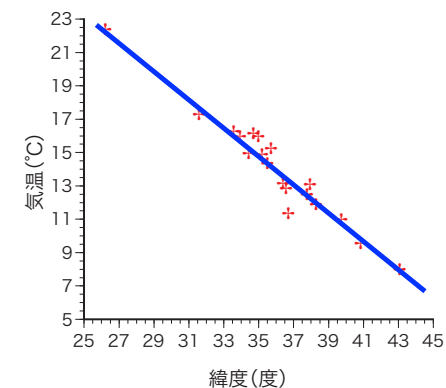
気温は緯度によって説明される  
(というモデル)



線形単回帰 🤔

## 線形単回帰

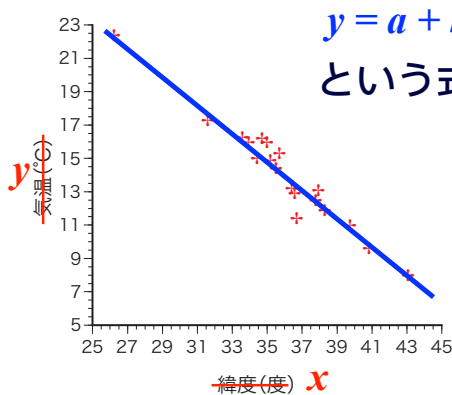
気温は緯度によって説明される  
どう説明される？



散布図上で直線の関係がある、と考える

# 線形単回帰

散布図上で直線の関係がある

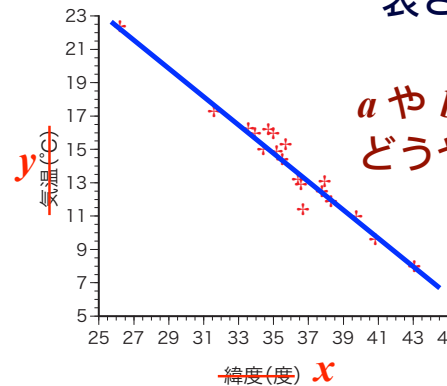


$y = a + bx$   
という式で表される関係

[線形単回帰]  
という

# 線形単回帰

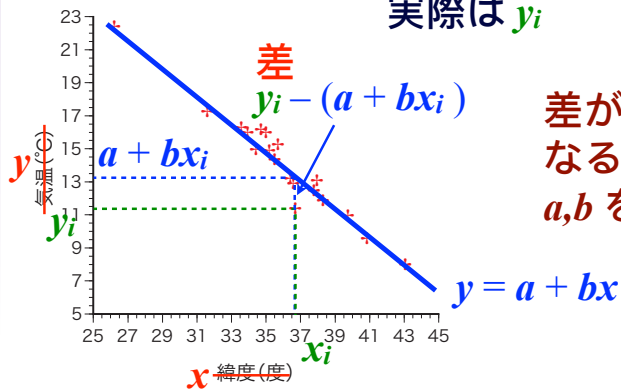
$y = a + bx$  という式で  
表される関係



$a$  や  $b$  (パラメータ) は  
どうやって求める?

# パラメータの決定

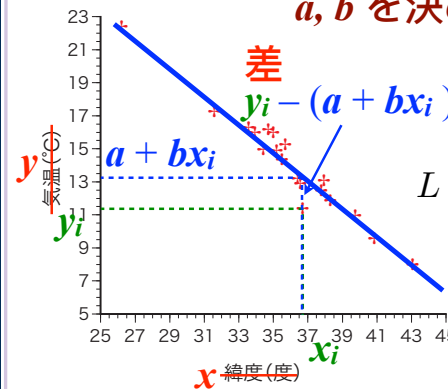
$x = x_i$  のとき  
モデルによれば  $a + bx_i$   
実際は  $y_i$



差が最小に  
なるように  
 $a, b$  を決める

# パラメータの決定

すべての  $x_i$  について、  
差の合計が最小になるように  
 $a, b$  を決める



$$L = \sum_{i=1}^n \{y_i - (a + bx_i)\}^2$$

が最小になる  
 $a, b$  を求める

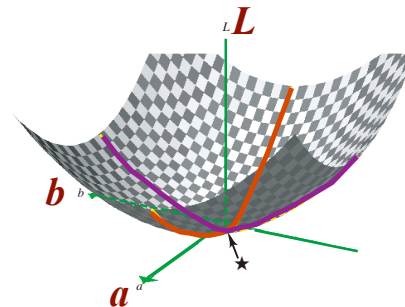
## Lが最小になるa,bを求める

- 偏微分による方法 (付録 1)
- 「2次関数の最大・最小」による方法 (付録 2)

## 「偏微分」による方法

$$L = \sum_{i=1}^n \{y_i - (a + bx_i)\}^2$$

が最小になる  $a, b$  を求める  
 $a, b$  の2次関数

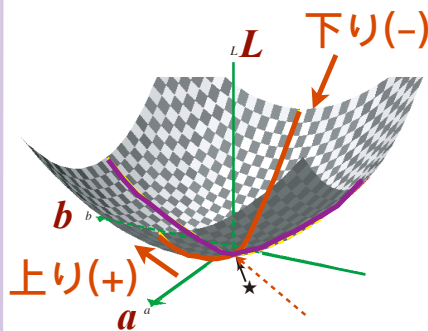


$a$ だけの関数  
と考える微分

$b$ だけの関数  
と考える微分

微分? 🤖

## 微分?



$a$ だけの関数  
と考える微分

微分は、傾きを  
求める計算

底では微分=0

$b$ についても同じ、  
底では微分=0

これらから  
 $a, b$ を求める

## 計算はともかく結論は

- 偏微分による方法 (付録 1)
- 「2次関数の最大・最小」による方法 (付録 2)

$$b = \frac{\sigma_{xy}}{\sigma_x^2}$$

$\sigma_{xy}$  ←  $x, y$  の共分散  
 $\sigma_x^2$  ←  $x$  の分散

$$a = \bar{y} - b\bar{x}$$

$\bar{y}$  ←  $y$  の平均  
 $\bar{x}$  ←  $x$  の平均

## 最小二乗法

$$b = \frac{\sigma_{xy}}{\sigma_x^2}$$

$$L = \sum_{i=1}^n \{y_i - (a + bx_i)\}^2$$

を最小にしたので  
[最小二乗法]

$$a = \bar{y} - b\bar{x}$$

$$y = a + bx$$

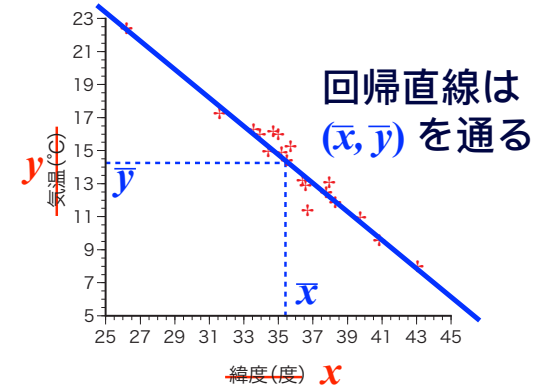
[回帰方程式] あるいは  
[回帰直線]

[回帰係数]

## ところで

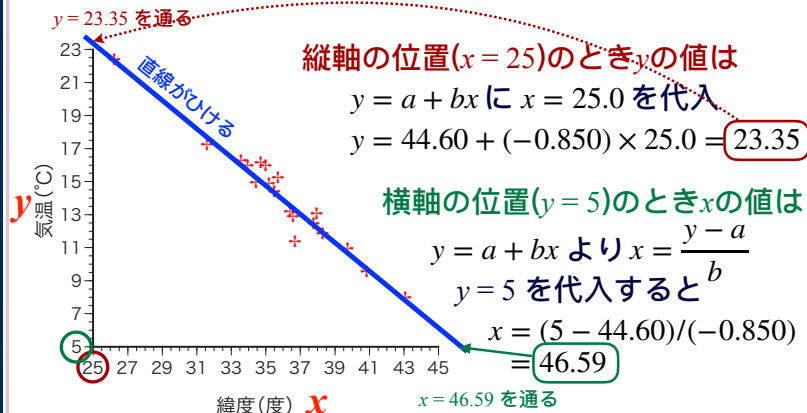
$$y = a + bx \quad \text{から} \quad y - \bar{y} = b(x - \bar{x})$$

$$a = \bar{y} - b\bar{x}$$



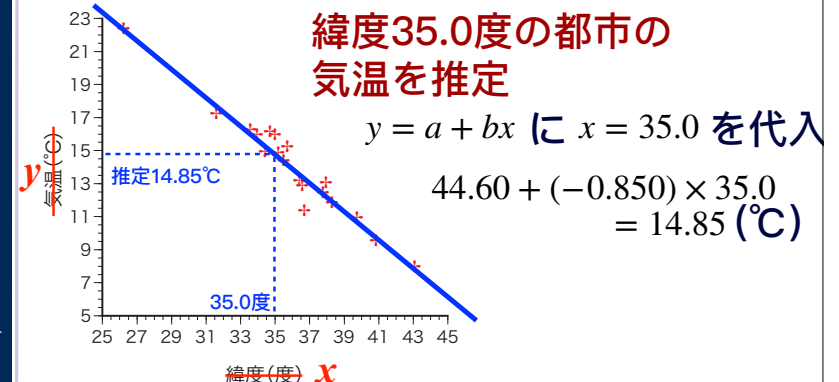
## 緯度と気温の例では

緯度を  $x$  として回帰直線  $y = a + bx$  を求める  
気温を  $y$  として  
→  $b = -0.850$ ,  $a = 44.60$



## 緯度と気温の例では

緯度を  $x$  として回帰直線  $y = a + bx$  を求める  
気温を  $y$  として  
→  $b = -0.850$ ,  $a = 44.60$

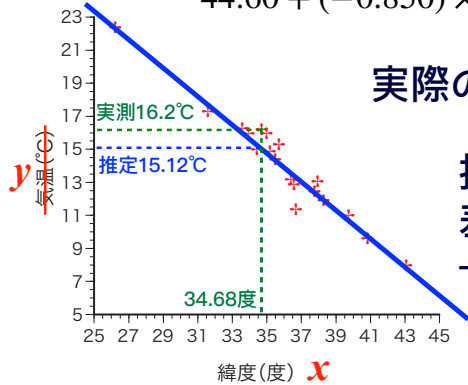


## 緯度と気温の例では

表にある大阪市(緯度34.68度)の気温を推定

$y = a + bx$  に  $x = 34.68$  を代入

$$44.60 + (-0.850) \times 34.68 = 15.12 (^{\circ}\text{C})$$



実際の気温は16.2°C

推定値と実測値に  
差がある  
→次の話へ

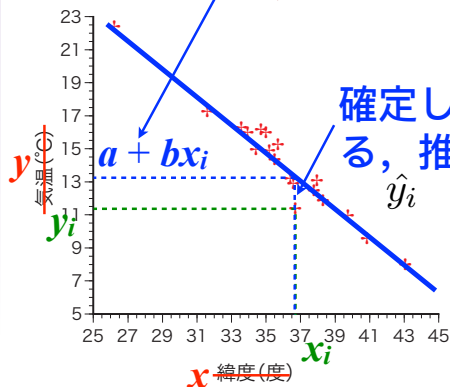
決定係数 🤔

## 残差

$a, b$  が求められて、回帰直線が確定

$x_i$  に対する、回帰直線による  $y$  の推定値

$$\hat{y}_i = a + bx_i$$



確定しても残っている、推定値と実際の差

[残差] という  
 $d_i$

## 残差と決定係数

回帰方程式を使って  $y_i$  を予測したときの、予測によって表現できなかった部分

残差について (付録3)

$$\sum d_i^2 = (1 - r_{xy}^2) \sum (y_i - \bar{y})^2$$

残差                      相関   決定  
                                 係数   係数

## 決定係数の意味

$$\sum d_i^2 = (1 - r_{xy}^2) \sum (y_i - \bar{y})^2 \text{ より}$$

残差の2乗の平均

$$1 - r_{xy}^2 = \frac{\sum d_i^2 / n}{\sum (y_i - \bar{y})^2 / n}$$

決定  
係数

yの偏差の2乗の平均  
= yの分散

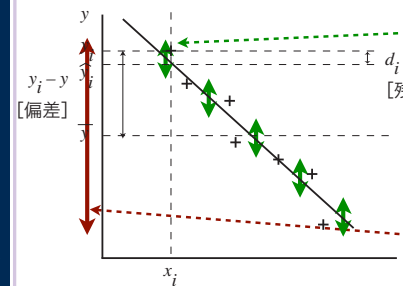
## 決定係数の意味

$$1 - r_{xy}^2 = \frac{\sum d_i^2 / n}{\sum (y_i - \bar{y})^2 / n}$$

決定  
係数

残差の2乗の平均

yの偏差の2乗の平均  
(yの分散)



回帰直線から見ると  
ばらつきは  
こんなに減った

もともとyはこんなに  
ばらついていたが、

## 決定係数の意味

$$1 - r_{xy}^2 = \frac{\sum d_i^2 / n}{\sum (y_i - \bar{y})^2 / n}$$

決定  
係数

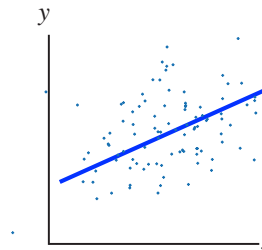
回帰直線からの  
ばらつき

yのもともとの  
ばらつき

決定係数 = 回帰直線によるばらつきの  
減少の割合

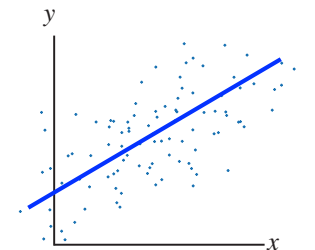
= 回帰直線によって、  
ばらつきの何%が説明できたか

## 「中くらいの相関」とは



相関係数0.5  
決定係数0.25

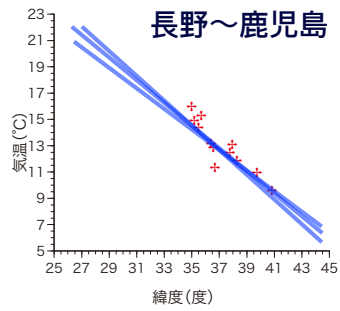
回帰直線では  
ばらつきの25%  
しか説明できない



相関係数0.7  
決定係数0.49

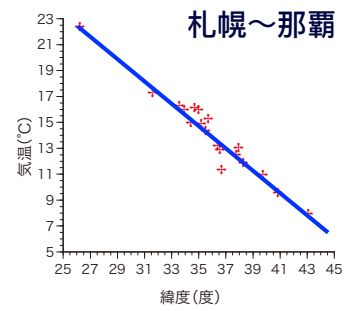
こちらが、  
中くらいの相関関係

## 前回の演習問題の例



決定係数0.712

平均付近に密集して  
いると不安定



決定係数0.949

平均から離れた  
個体があると  
安定する

## 困った例

(教室でのみ資料を呈示します)

(2014年10月29日 日経新聞)