

# 情報管理

## 第3回：時系列データの可視化と相関

# 今回の講義内容

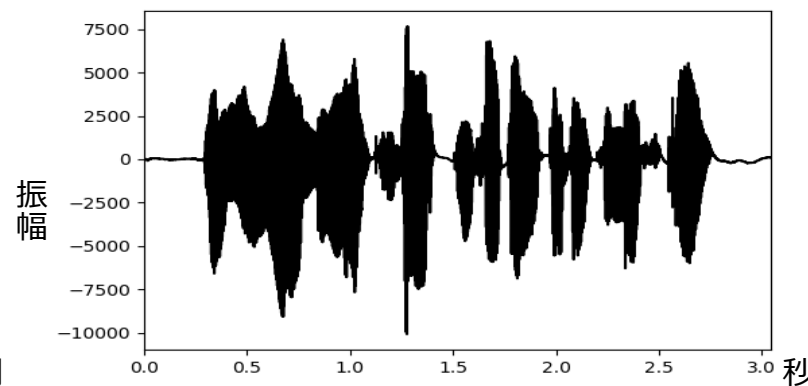
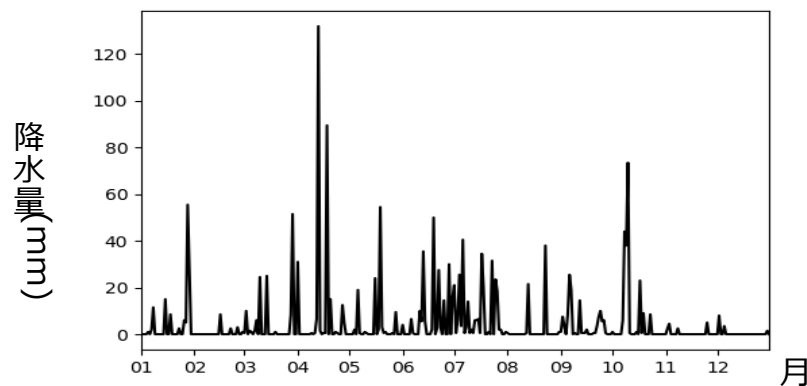
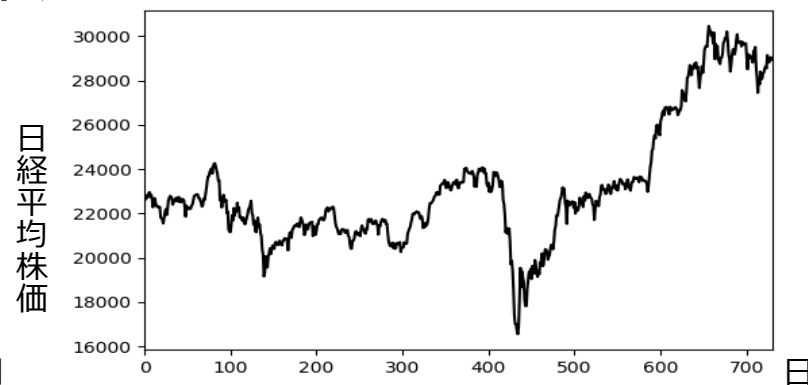
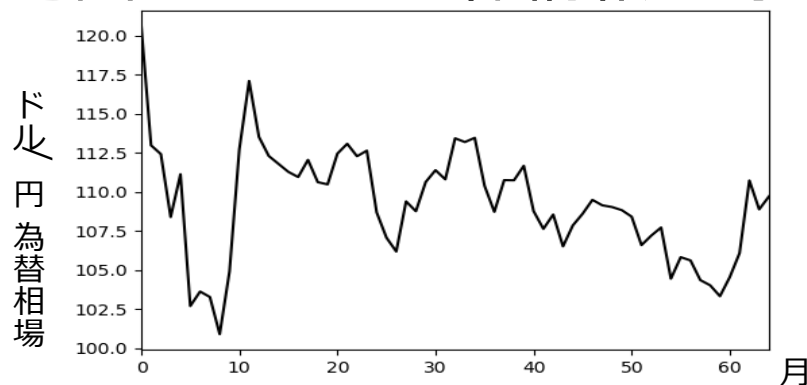
前回までは python の使い方説明が中心でしたが、今回からは本格的にデータの解析方法について勉強していきます。

今回は以下の内容を学んでもらいます。

- 時系列データの可視化
  - 長期的な変化と短期的な変化の可視化：移動平均と階差
  - 変化傾向の数値化：直線近似と曲線近似
- データ間の関係性の分析
  - 相関係数
  - 相互相関関数

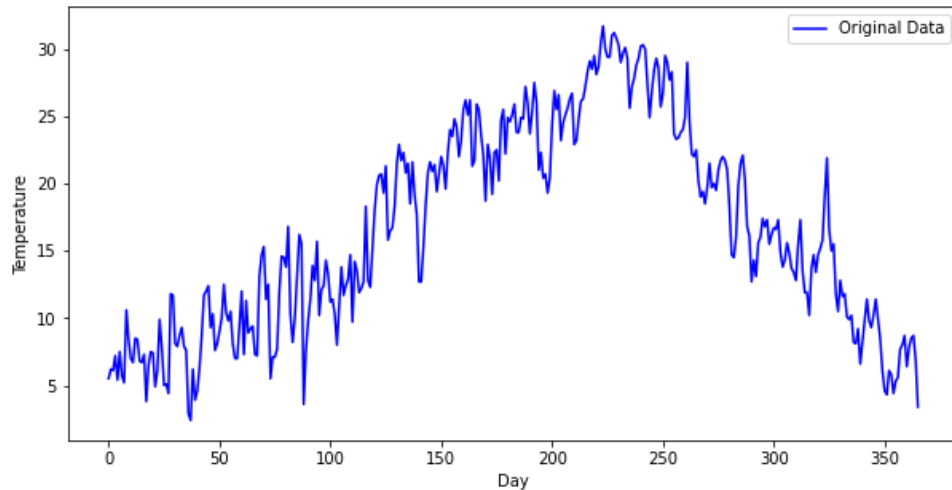
# 時系列データとは？

時間毎に変化する現象を一定間隔で記録したものを**時系列データ**と呼びます。  
前回扱った気温データは、1日間隔で記録した時系列データとなります。  
為替相場、株価のような経済データや気温・降水量のような気象データ、  
音声や心電図のような生体情報も時系列データです。



# 時系列データの可視化

時系列データをそのままプロットするだけでも分かる情報はありますが、様々な処理を行うことで、特定の情報を際立たせて見やすくできます。

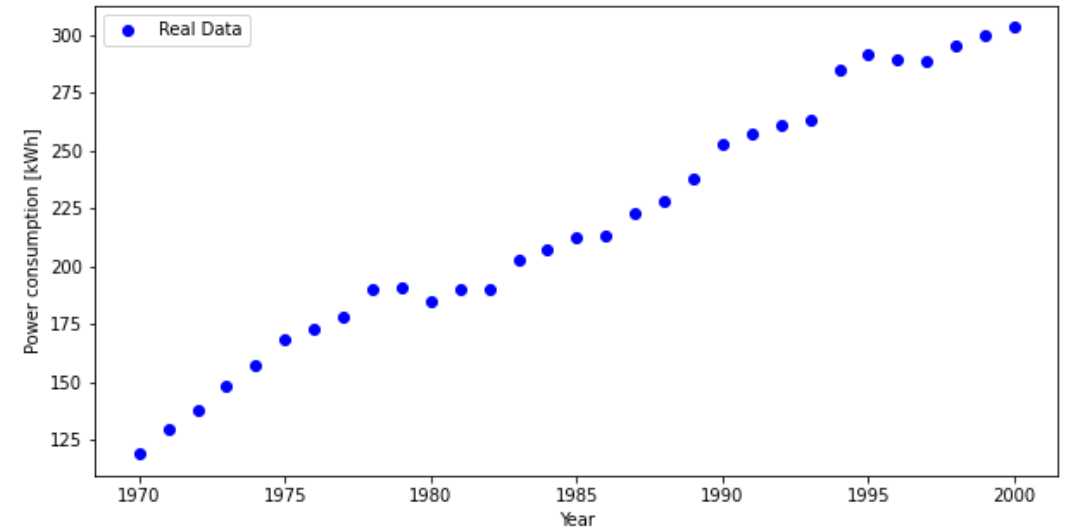


大まかにはどのように変化している？

→ 長期的変動の可視化

どれくらい小刻みに変化している？

→ 短期的変動の可視化

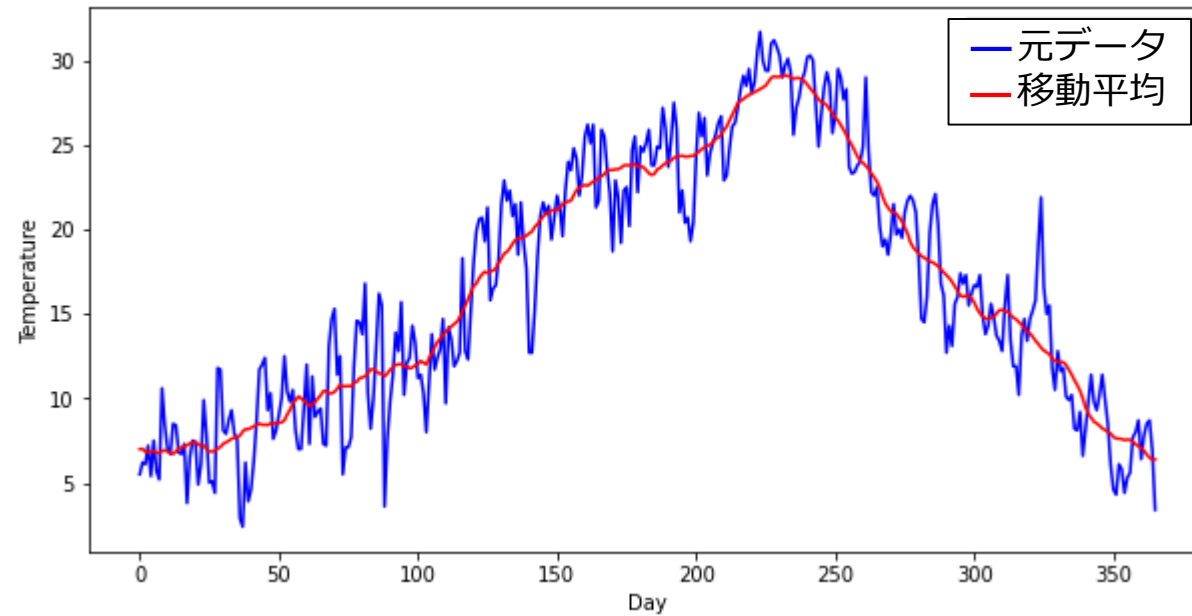


どれくらいの割合で増加している？

→ 傾きの数値化

# 時系列データの可視化方法 1-1：移動平均

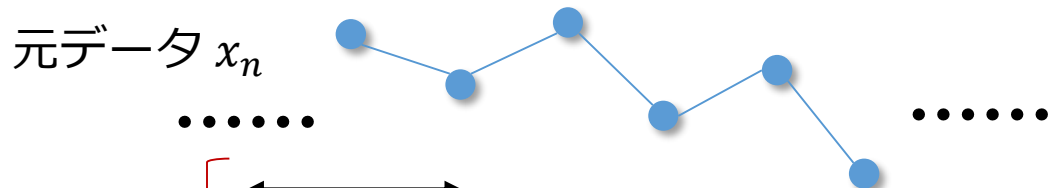
細かい変動成分を抑える（**平滑化**）ことで、大まかな変化傾向（＝長期的変動）を見やすくする方法です。



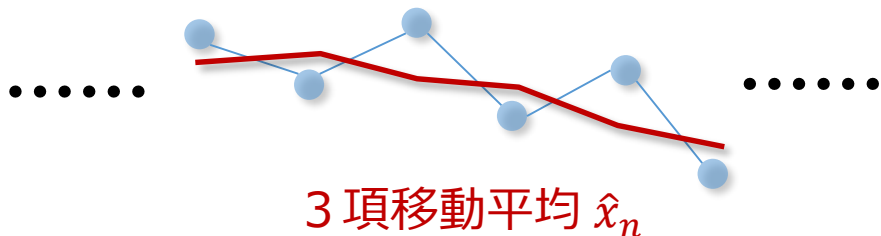
# 時系列データの可視化方法 1-1：移動平均

ある時刻に対して，その前後いくつかの時刻の値を用いて，平均を計算します。

前後 1 時刻 ( $k=1$ ) を用いる場合



それぞれの  
区間で  
平均値を  
求める



$$\hat{x}_n = \frac{1}{2k+1} (x_{n-k} + \cdots + x_n + \cdots + x_{n+k})$$

現在の時刻  $n$  に対して，  
前後  $k$  時刻の値を用いて  
平均を計算する。

★ 前後 2 方向に  $k$  項！

前後  $k$  時刻を用いた場合，  
 $K = 2k+1$  個のデータで平均計算することになる。  
→  $k$  項移動平均と呼ぶ。

# 時系列データの可視化方法 1-1：移動平均

移動平均は計算方法によって、さらに以下の3種類に分けられます。

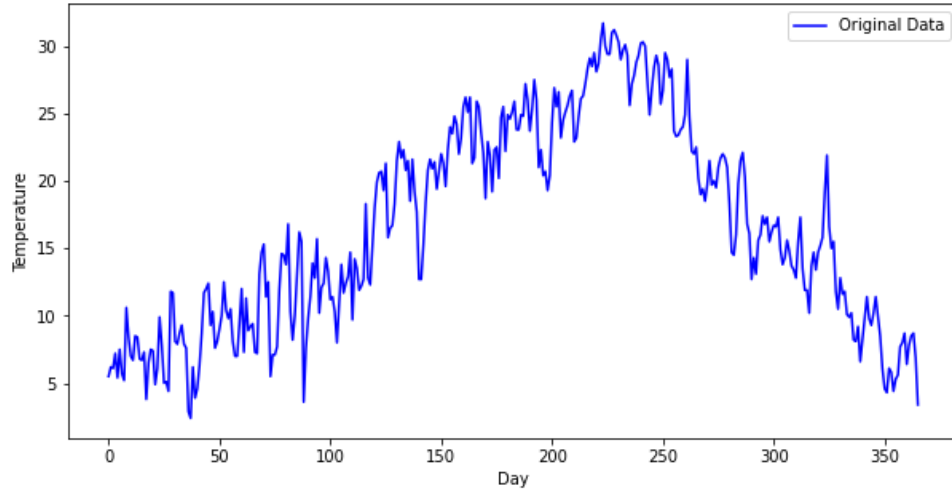
- 前後の時刻を用いる場合：中央移動平均  
前ページで紹介した方法。既に収集完了した区間のデータに対してよく用いられる。
- 後ろの時刻のみを用いる場合：前方移動平均
- 前の時刻のみを用いた場合：後方移動平均  
後ろの時刻が存在しない＝現在進行形で収集・分析をしている場合においてよく用いられる（例：COVID-19感染者の分析）。

# 時系列データの可視化方法 1-1：階差

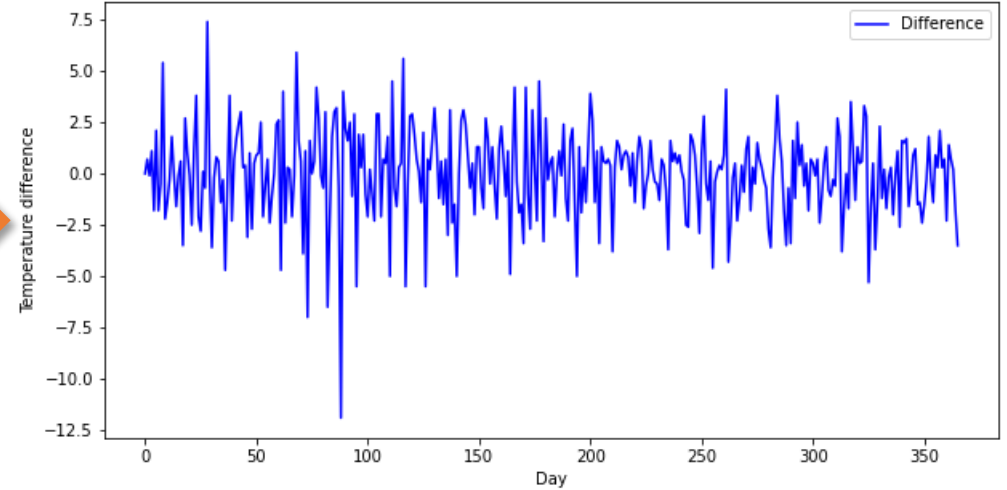
現在の時刻から一つ前の時刻の値との差を求めることで，細かい変化（＝短期的変動）を見やすくする方法です。

$$\Delta x_n = x_n - x_{n-1}$$

元データ  $x_n$



階差  $\Delta x_n$



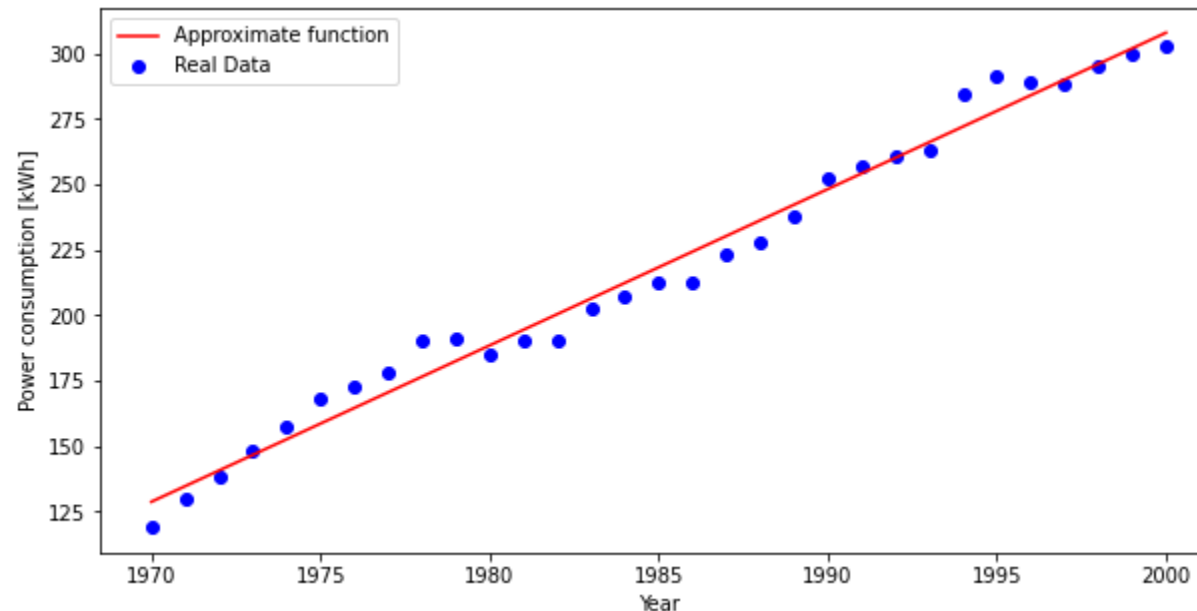


# 気温データに対して移動平均と階差を計算しよう

03\_01\_moving\_average.ipynb を動かして、気温データに対する移動平均と階差を計算し、どのようなことが分かるかを見てみましょう。

# 時系列データの可視化方法2-1：直線近似

時系列データを直線（線形関数： $y = ax + b$ ）で近似してみましょう。  
直線近似することで、時系列データの増減傾向や全体的な大きさが傾き  $a$  と切片  $b$  によって数値化できます。

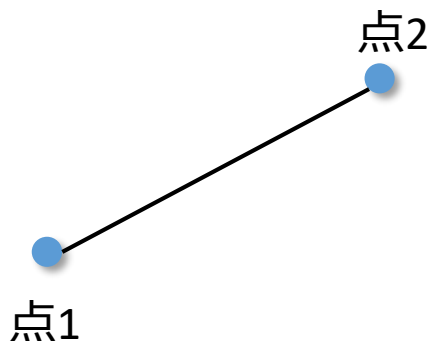


# 時系列データの可視化方法2-1：直線近似

中学・高校数学では、「点1と点2を通る直線を求めよ」という問題で、連立方程式を解くことで直線を求めていたと思います。

しかし、点が沢山ある場合、全てを通る直線を引くことは通常不可能です。そこで、できるだけ**全ての点に近い場所を通る直線を引く**ことを考えます。

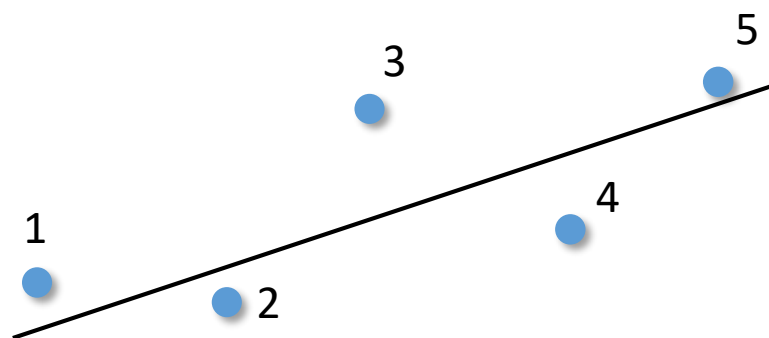
「2点を**通る**直線を求めよ」



$$\begin{cases} y_1 = ax_1 + b \\ y_2 = ax_2 + b \end{cases}$$

連立方程式を解けば良い

「全ての点の**近くを通る**直線を求めよ」



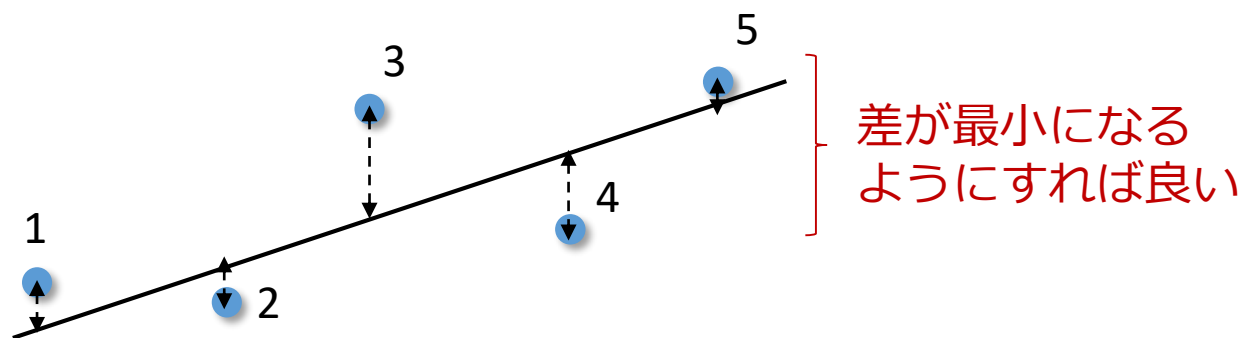
$$\begin{cases} y_1 = ax_1 + b \\ y_2 = ax_2 + b \\ y_3 = ax_3 + b \\ y_4 = ax_4 + b \\ y_5 = ax_5 + b \end{cases}$$

**全てを満たす  
a, b は存在しない！**

# 時系列データの可視化方法2-1：直線近似

できるだけ全ての点に近い場所を通る直線を引く

→直線  $ax + b$  が計算した  $y$  と、実際のデータの  $y$  との誤差が小さければよい。



つまり、以下の式を満たす  $a, b$  を求めればよい。

$$\sum_{\substack{n \\ \text{データ数で} \\ \text{総和}}} \underbrace{(y_n - (ax_n + b))}_{\substack{\text{実際の}y \\ \text{直線が計算した}y}}^2 \rightarrow \min$$

誤差の二乗を最小化することから、これを **最小二乗法** と呼びます。

# 時系列データの可視化方法2-1：直線近似

最小二乗法を用いて，傾き  $a$  と切片  $b$  を求めます。

$$L = \sum_n (y_n - (ax_n + b))^2 = \sum_n (y_n^2 + a^2 x_n^2 + b^2 - 2ax_n y_n + 2abx_n - 2by_n) \rightarrow \min$$

$L$  は  $a$  に対する下に凸の2次関数であり，また  $b$  に対しても下に凸の2次関数である。  
よって，最小値を求めるには  $a, b$  で偏微分し， $=0$  を解けばよい。

$L$  を  $a$  に対して偏微分し， $=0$  とおく

$$\begin{aligned} \frac{\partial L}{\partial a} &= 2a \sum_n x_n^2 - 2 \sum_n x_n y_n + 2b \sum_n x_n = 0 \\ &\rightarrow a \sum_n x_n^2 - \sum_n x_n y_n + b \sum_n x_n = 0 \quad \dots\dots \textcircled{1} \end{aligned}$$

$L$  を  $b$  に対して偏微分し， $=0$  とおく

$$\begin{aligned} \frac{\partial L}{\partial b} &= 2Nb + 2a \sum_n x_n - 2 \sum_n y_n = 0 \\ &\rightarrow a \sum_n x_n + Nb - \sum_n y_n = 0 \quad \dots\dots \textcircled{2} \end{aligned}$$

# 時系列データの可視化方法2-1：直線近似

$$a \sum_n x_n^2 - \sum_n x_n y_n + b \sum_n x_n = 0 \quad \dots\dots \textcircled{1}$$

$$a \sum_n x_n + Nb - \sum_n y_n = 0 \quad \dots\dots \textcircled{2}$$

上記の連立方程式を解くと、以下の結果が得られます（導出は省略）

## 直線近似の傾きと切片

$x_n$ :  $n$ 番目のデータの  $x$  軸の値  
 $y_n$ :  $n$ 番目のデータの  $y$  軸の値  
 $N$ : データの総数

$$a = \frac{N \sum_n x_n y_n - \sum_n x_n \sum_n y_n}{N \sum_n x_n^2 - (\sum_n x_n)^2}$$

$$b = \frac{\sum_n x_n^2 \sum_n y_n - \sum_n x_n y_n \sum_n y_n}{N \sum_n x_n^2 - (\sum_n x_n)^2}$$

# 時系列データの可視化方法2-1：直線近似

Q. なぜ誤差の二乗を最小化するのか？

A. いくつか理由があります。

- 誤差を二乗すればかならずゼロ以上の値になるので，誤差の正負を気にしなくて良い。
- 最適化関数が下に凸の二次関数になるため，偏微分して  $=0$  を求める方法が取れるため，使い勝手が良い。

誤差の二乗のほかに，誤差の絶対値を取る方法などもあります。

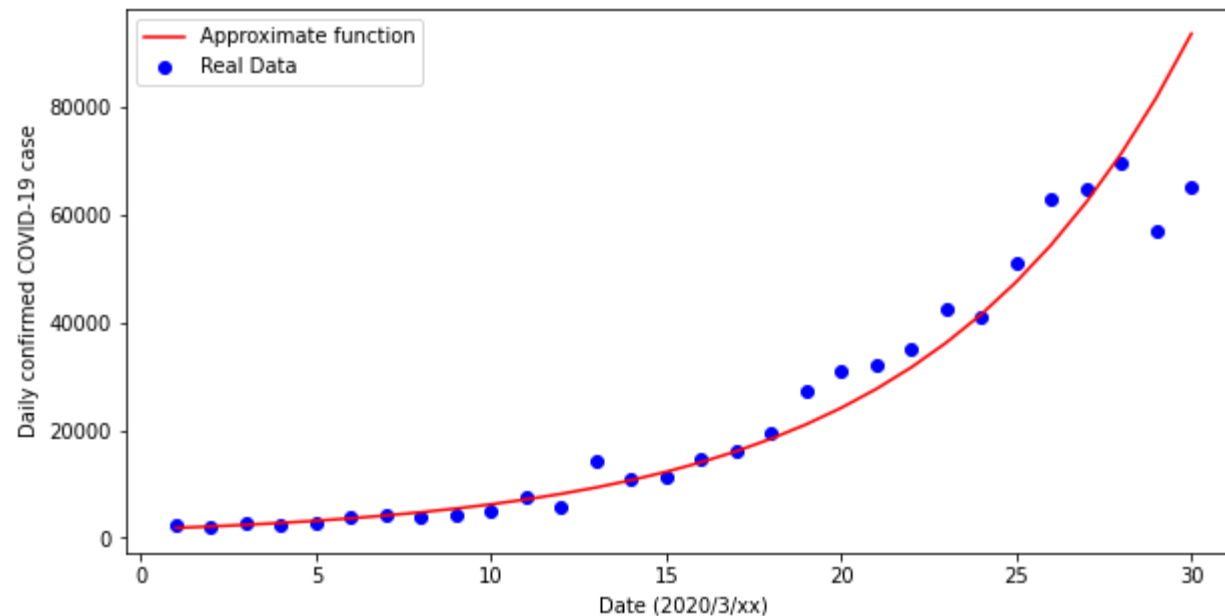
最小二乗法は今後も出てきますので，理解しておいてください。

# 時系列データの可視化方法2-1：曲線近似

時系列データを曲線で近似する場合、いくつか種類が考えられます。

(例：2次関数, 3次関数, ..., n次関数, 対数関数など)

ここでは、指数関数  $y = be^{ax}$  を使用します。





# 時系列データの可視化方法2-1：曲線近似

指数関数で近似する場合、データ全体に対数を計算することで、直線近似と同じ問題にすることができます。

$$y = be^{ax}$$

両辺の対数を取る

$$\begin{aligned}\log_e y &= \log_e be^{ax} \\ &= \log_e b + \log_e e^{ax} \\ &= \log_e b + ax \\ &= B + ax\end{aligned}$$

つまり、データ  $y$  に対して対数を取ったデータ  $\log_e y$  に対して、 $B + ax$  直線近似を行う。その後、 $b = e^B$  を計算すれば、 $a, b$  が求まる。

# 直線近似を動かしてみよう

03\_02\_approximate\_function.ipynb を動かして、直線近似の様子をみてみましょう。

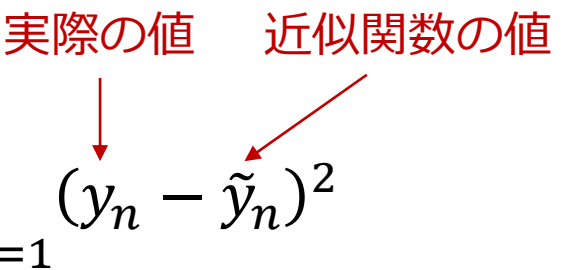
# [補足] どの程度うまく近似できているか？

近似関数の値と実際の値との誤差を計算することで、近似がどの程度うまくできているかを評価できます。

平均二乗誤差 (Mean square error: MSE)

$$\text{MSE} = \frac{1}{N} \sum_{n=1}^N (y_n - \tilde{y}_n)^2$$

実際の値      近似関数の値



最小二乗法において最小化している式そのもの。

平方根平均二乗誤差 (Root mean square error: RMSE)

$$\text{RMSE} = \sqrt{\text{MSE}}$$

平均二乗誤差の平方根を取ることで、目的変数と同じ単位で誤差が測れる。

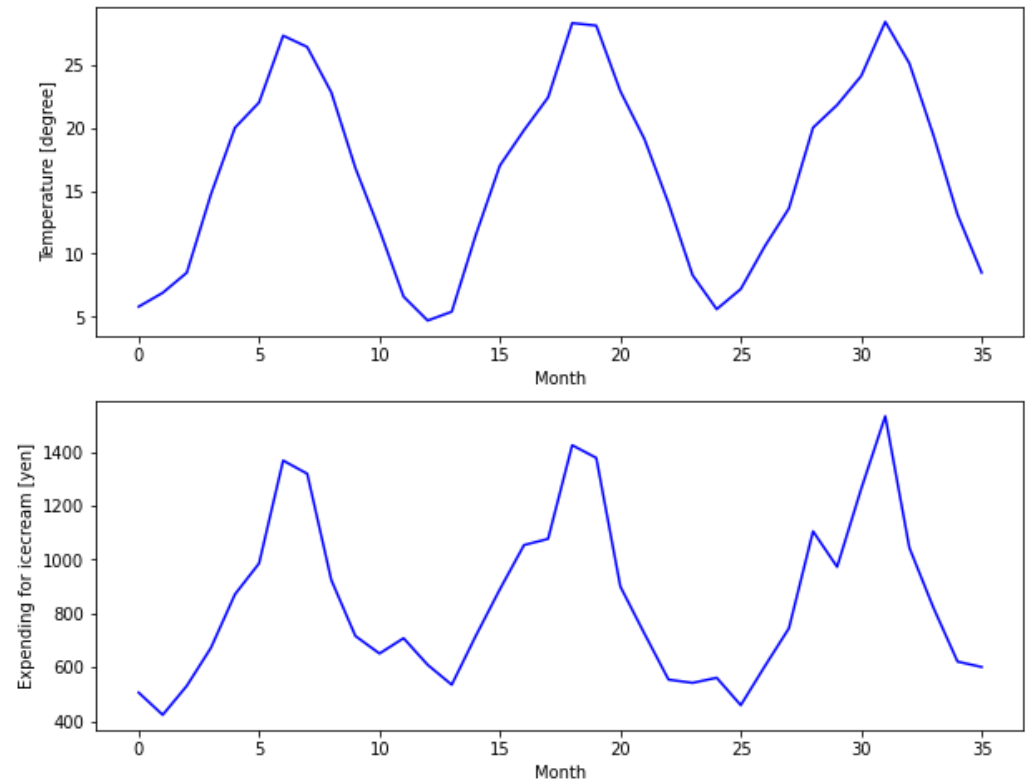
# データ間の関係性 1：相関係数

アイスクリームは夏の暑い時期に良く売れ、逆に寒い冬の時期にはあまり売れません。

このことから、アイスクリームの売り上げと気温には強い関係があると想像できます。

一方が増えると他方も増え、  
一方が減ると他方も減るというように、  
増減の仕方が似ているデータ同士を  
「**相関**がある」と呼びます。

3年間の気温の変化(上)とアイスの売上(下)



# データ間の関係性 1：相関係数

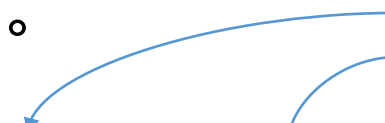
相関の強さを数値化したものとして、**相関係数**があります。

相関係数 R は以下の式で計算されます。

$$R = \frac{\sum_n (x_n - \bar{x})(y_n - \bar{y})}{\sqrt{\sum_n (x_n - \bar{x})^2} \sqrt{\sum_n (y_n - \bar{y})^2}}$$

$x$  の標準偏差                       $y$  の標準偏差

$x, y$  それぞれの平均値



R は -1 から 1 までの値を取り、R の値によって以下のように呼びます。

- R が正の時：正の相関
- R が負の時：負の相関
- R が 0 の時：無相関

# 相関係数の直感的説明

$$R = \frac{\sum_n (x_n - \bar{x})(y_n - \bar{y})}{\sqrt{\sum_n (x_n - \bar{x})^2} \sqrt{\sum_n (y_n - \bar{y})^2}}$$

$x_n$  と  $y_n$  が両方**増加**  $\rightarrow (x_n - \bar{x}), (y_n - \bar{y})$  どちらも**正**  $\rightarrow (x_n - \bar{x})(y_n - \bar{y})$  は**正**

$x_n$  と  $y_n$  が両方**減少**  $\rightarrow (x_n - \bar{x}), (y_n - \bar{y})$  どちらも**負**  $\rightarrow (x_n - \bar{x})(y_n - \bar{y})$  は**正**

**➡**  $x$  と  $y$  が**同じ増減傾向**をしていれば,  $\sum_n (x_n - \bar{x})(y_n - \bar{y})$  は**大きい正の値**になる。

一方,

$x_n$  と  $y_n$  の一方が**増加**, 他方が**減少**  $\rightarrow (x_n - \bar{x}), (y_n - \bar{y})$  の一方は**正**, 他方は**負**  $\rightarrow (x_n - \bar{x})(y_n - \bar{y})$  は**負**

**➡**  $x$  と  $y$  が**逆の増減傾向**をしていれば,  $\sum_n (x_n - \bar{x})(y_n - \bar{y})$  は**大きい負の値**になる。

また,

$x$  と  $y$  が**同じ増減**, **逆の増減**が入り混じっている場合, つまり**統一性が無い**場合,

**➡**  $\sum_n (x_n - \bar{x})(y_n - \bar{y})$  は**0**になる。

なお, 分母の標準偏差は,  $R$  を -1 から 1 のレンジに入れる (正規化) ための項である。

# データ間の関係性 2：相互相関関数

時系列データ同士で相関を調べる場合、**お互いの増減関係に時間差が存在している場合があります。**

例：広告を出した数と売上の関係

通常、広告の効果は遅れてやってくるため、広告数の増減と比べて売上は遅れて増減するはず。

**時間差を考慮しながら相関関係を調べる方法**として、**相互相関関数**があります。

# データ間の関係性 2：相互相関関数

相互相関関数は以下の式で計算されます。

$$R(m) = \sum_n x_n \underbrace{y_{n+m}}_{y_n \text{ を時刻 } m \text{ ずらしている}}$$

相関係数

$$R = \frac{\sum_n (x_n - \bar{x})(y_n - \bar{y})}{\sqrt{\sum_n (x_n - \bar{x})^2} \sqrt{\sum_n (y_n - \bar{y})^2}}$$

相関係数との違いは

- 分母が無い，また分子に対しても平均で引く処理が無い（標準化しない）
- $R$  が  $m$  の関数である（ $m$  の数だけ  $R$  を計算している）。
- $y$  が時刻  $m$  だけずれている。

つまり，一方のデータ(ここでは  $y$ ) の時刻を  $m$  ずらしつつ，  
相関(ただし標準化なし)を計算している。



# [補足] 正規化について

単位（スケール）が異なるデータを同一に扱くと，想定通りの分析ができない場合があります。

	気温 [°C]	アイス支出額 [円]
1	5.8	506
2	6.9	423
3	8.5	531

スケールを合わせるため，データ毎に平均が0，標準偏差が1になるように正規化します。（これを**標準化**と呼びます。）

$$x_{norm} = \frac{x - (x \text{の平均})}{(x \text{の標準偏差})}$$

相関係数の式は，標準化したデータ同士の内積に相当します。

相互相関関数を計算する際にも事前に標準化しておくと，相関係数のように扱うことができます。

（03\_03\_correlation.ipynb でも標準化しています。）

相関係数

$$R = \frac{\sum_n (x_n - \bar{x})(y_n - \bar{y})}{\sqrt{\sum_n (x_n - \bar{x})^2} \sqrt{\sum_n (y_n - \bar{y})^2}}$$

# 相関係数と相互相関関数を求めてみよう

03\_03\_correlation.ipynb を動かして，相関係数と相互相関関数を求めてみましょう。

# おわりに

今回は、主に時系列データの可視化方法と、相関について解説しました。

特に最小二乗法は、以降でもパラメータ推定方法として出てきますので、しっかり理解しておきましょう。

# レポート課題

(ゴールデンウィークを挟むので 2 題出します。)

課題 1 : 直線近似と曲線近似

課題 2 : 音源方向推定

前回のレポートと同じく, `convert_report.ipynb` を使って html ファイルに変換してから提出してください。

※課題 1, 2 のそれぞれ個別にBEEFの提出リンクを貼っています。

レポート提出期限 : 5/10(火) AM10:30

# レポート課題 1 : 直線近似と曲線近似

ファイル一式に入っている「covid19.csv」を読み込んでください。

このデータは、2020年3月1日から2020年3月30日までの、全世界のCOVID-19の新規感染者数を日ごとに記録したデータです。

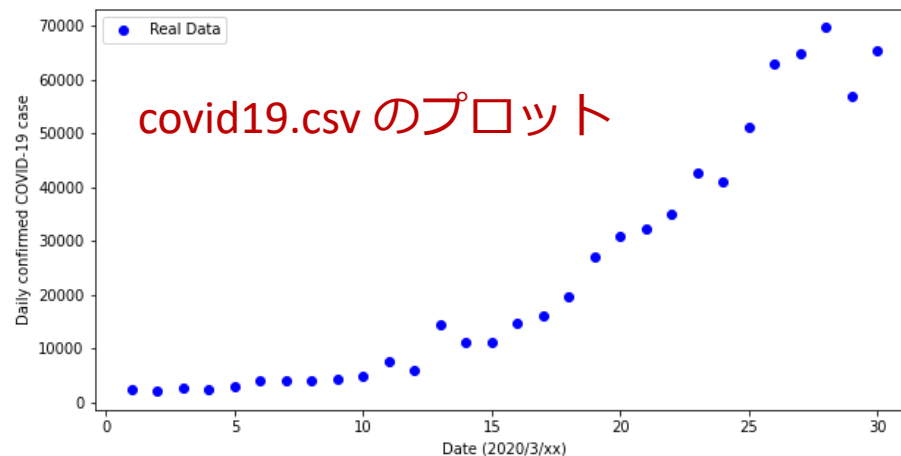
横軸( $x$ ) = 日, 縦軸( $y$ ) = 新規感染者数としてプロットしたデータに対して, 一次関数  $y = ax + b$  および指数関数  $y = be^{ax}$  を用いてそれぞれ近似し, その結果を比較・考察しなさい。

ヒント :

対数  $y = \log(x)$  および指数  $y = e^x$  はそれぞれ numpy の **log**, **exp関数** で計算できます。

```
import numpy as np
```

```
log_x = np.log(x)    log(x)を計算  
exp_x = np.exp(x)     $e^x$ を計算
```



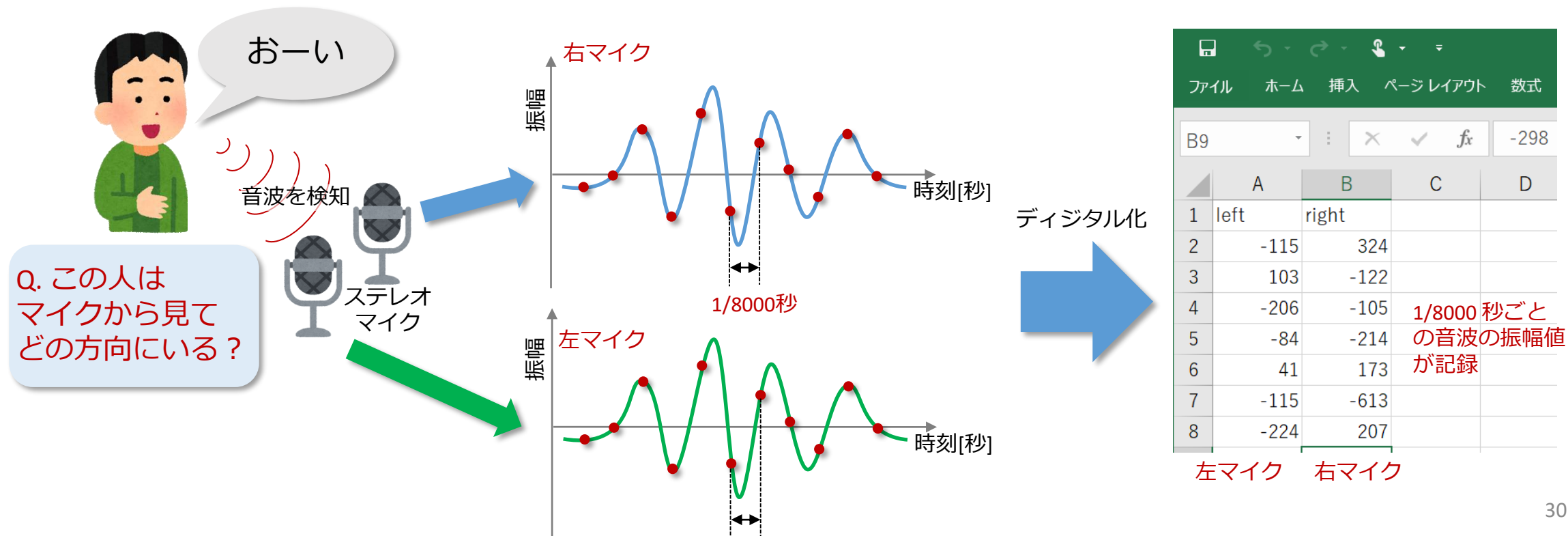
# レポート課題 2：音源方向推定

report03\_input.csv というデータを読み込んでください。

このデータは、左右ステレオのマイクで収録された音声信号です。

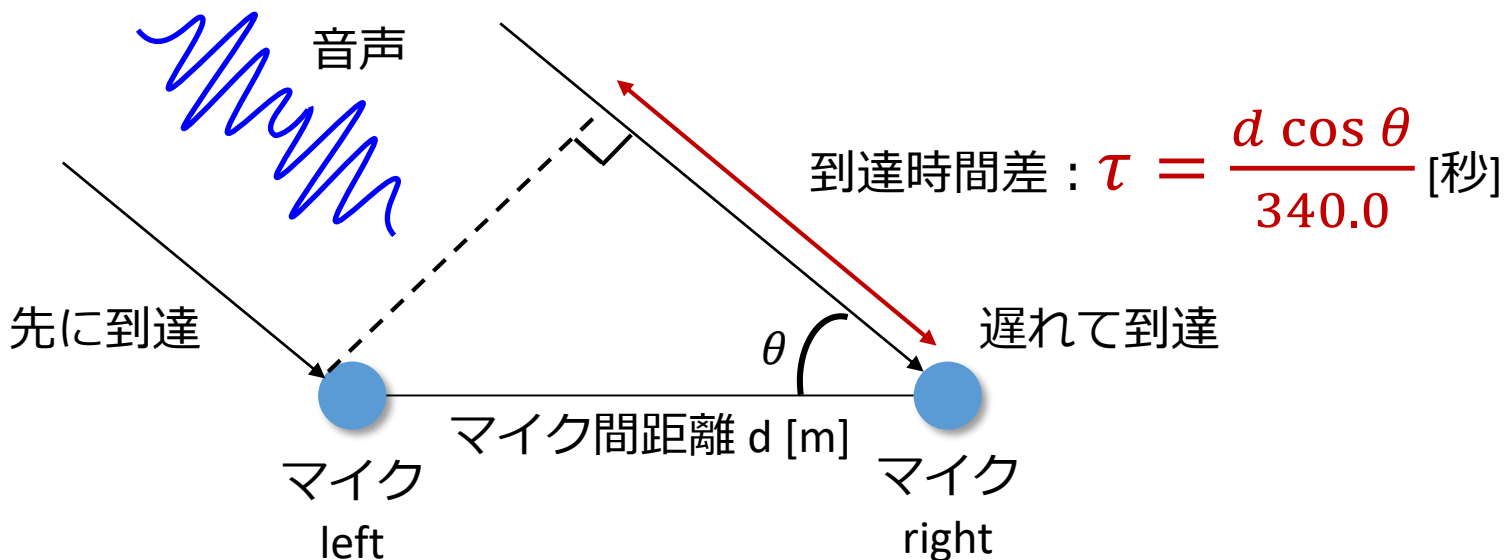
左右それぞれのマイクにおける音の振幅値を、**1/8000秒間隔**で記録しています。

次以降のページを読んで、**この音がどの方向からやって来たかを推定してください。**



# レポート課題 2 : 音源方向推定

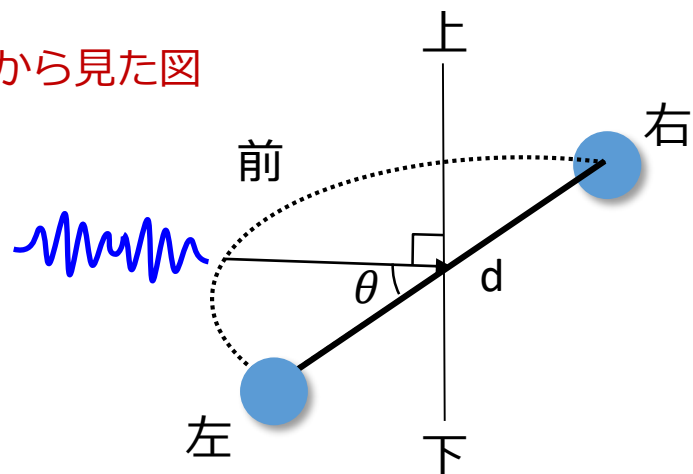
- 2個のマイクで音声を収録したとき、音声マイクに到達する時間差が生じる。
- マイク間距離を  $d$  [m], 音源方向を  $\theta$ , 音速を  $340.0$  [m/秒] としたとき、マイク間の音の到達時間差は  $\tau = d \cos \theta / 340.0$  [秒] である。
- つまり、マイク間の時間差  $\tau$  [秒] が分かれば、音源方向は以下の式で得られる。
$$\theta = \cos^{-1} \frac{340.0 \tau}{d}$$
- このデータのマイク間距離は  $d = 0.3$  [m] である。
- 音声データの時刻は 1 目盛りにつき  $1/8000$  秒である。



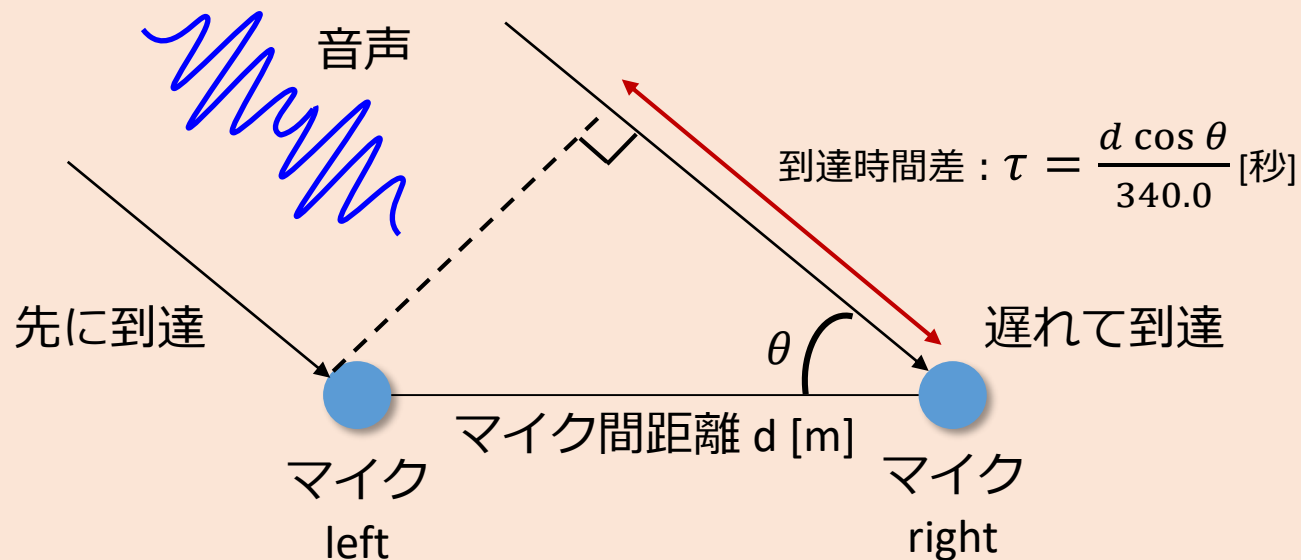
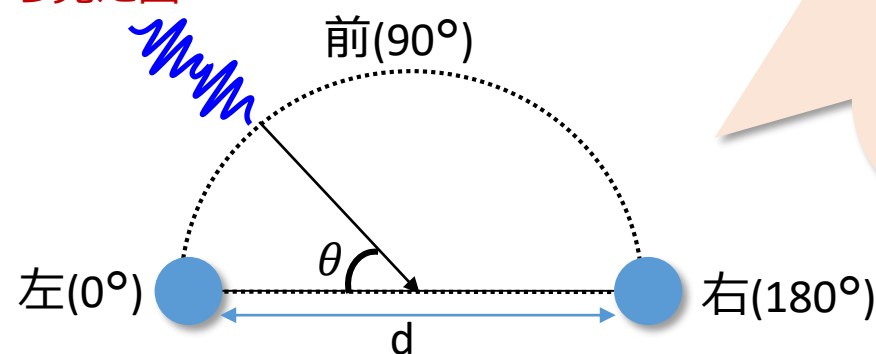
# レポート課題 2 : 注意点およびヒント

- 音源方向は下図のように定義する。マイクの真左から到来した場合は  $0^\circ$  , 真右から到来した場合は  $180^\circ$  である。
- 推定する角度  $\theta$  は「水平角」(前後左右)。高さ方向(仰角)は考慮しなくてよい。  
(便宜上, マイクと同じ高さから到来しているとする。)

斜め上から見た図



真上から見た図





# レポート課題 2 : 注意点およびヒント

- 求めた時間差の単位が [秒] に換算されているか注意すること。  
データの時刻の目盛りは 1/8000 秒である。
- $\cos^{-1}(A)$  の計算は, numpy の関数 `np.arccos(A)` で求められる。  
ただし戻り値の単位は radian である。
- radian から degree への変換は, `degree = np.rad2deg(radian)` で行える。