

進捗報告

2022/05/11

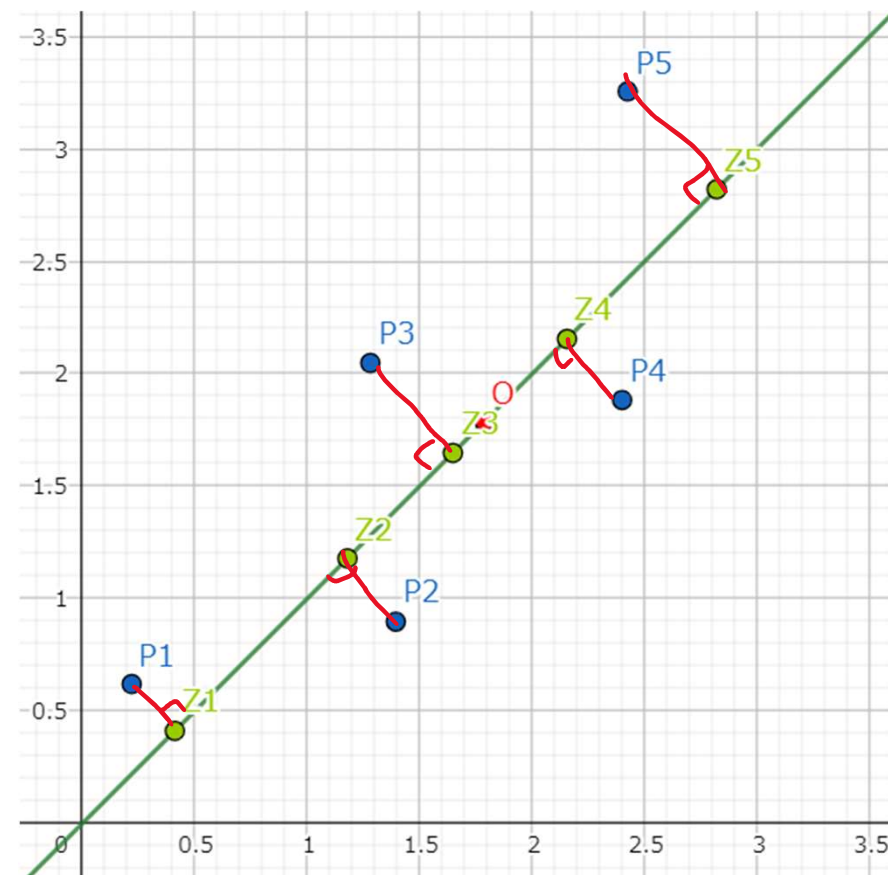
今回やったこと

- 主成分分析について

主成分分析：イメージ

- 点 $P_i (i = 1 \dots n)$ の平均点 O に対して
 $OZ_1^2 + OZ_2^2 + \dots + OZ_n^2$ が最大となる新しい軸（緑色）を探す

→新しい軸の分散を最大にする



主成分分析：データの準備

- データ $\{x_1, x_2, \dots, x_n\}$ に対してその平均を求める

$$\bar{x} = \frac{1}{N} \sum_{n=1}^N x_n$$

- また、平均が原点に来るように平行移動する

$$x'_n = x_n - \bar{x}$$

※以降、 x' は x と表記する

主成分分析：射影されたデータを求める

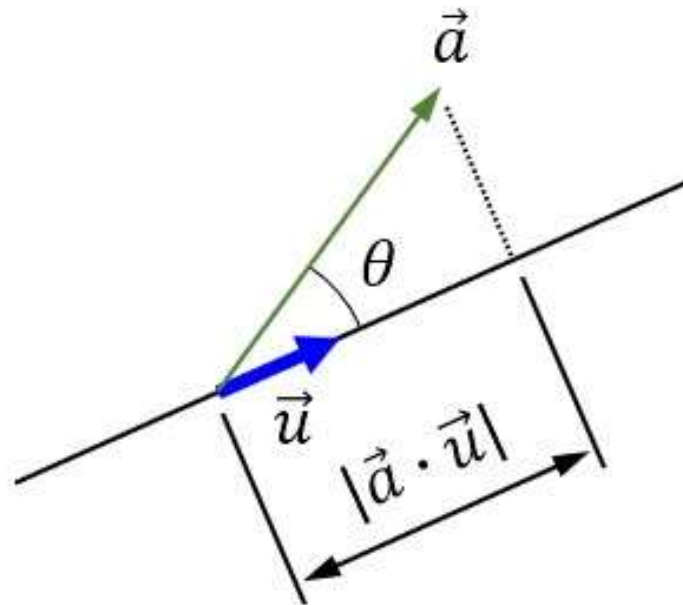
- データ $\{x_1, x_2, \dots, x_n\}$ に対して、単位ベクトル u を使って直交射影し、次の $\{z_1, z_2, \dots, z_n\}$ のスカラー値を得る

$$z_i = u^T x_i$$

- 単位ベクトルと任意ベクトルの内積を計算すると、単位ベクトル方向の成分（大きさ）が取得できる

主成分分析：射影されたデータを求める

- 例えば、 \vec{a} の単位ベクトル \vec{u} の方向の成分は、 \vec{a} と \vec{u} の内積の絶対値で求まる



主成分分析：射影されたデータを求める

■ \vec{a} の単位ベクトル \vec{u} の方向の成分は、 \vec{a} と \vec{u} の内積の絶対値で求まる

● \vec{a} と \vec{u} のなす角を θ とすると

$$L = |\vec{a}| \cos \theta$$

ここで、 \vec{a} と \vec{u} の内積は

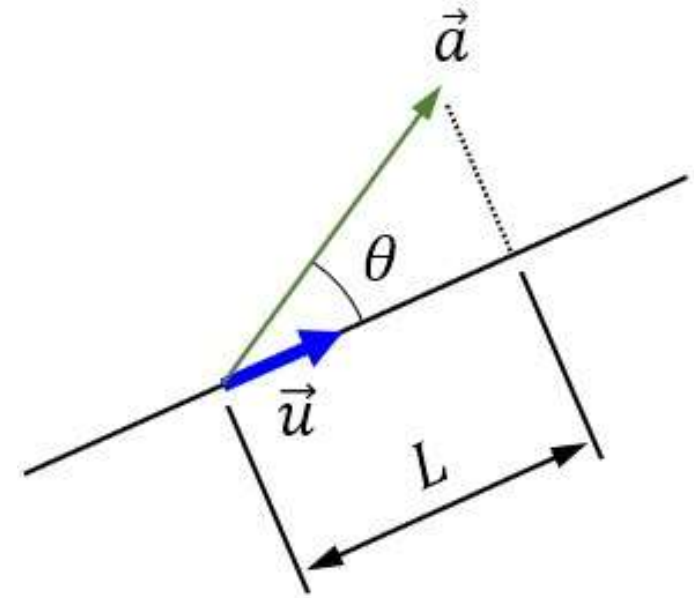
$$\vec{a} \cdot \vec{u} = |\vec{a}| |\vec{u}| \cos \theta$$

単位ベクトルの大きさは1なので

$$\vec{a} \cdot \vec{u} = |\vec{a}| \cos \theta$$

$$L = \vec{a} \cdot \vec{u}$$

$$z_i = \mathbf{u}^T \mathbf{x}_i$$



主成分分析：問題の定式化

- 射影されたデータの集合 $\{z_n\}$ の分散が最大化されるように u の値を決定したい。
つまり、分散を V_z とすると、

$$V_z = \frac{1}{N} \sum_{n=1}^N z_n^2$$

を最大化したいという問題になる

$$\text{分散} = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2$$

主成分分析：問題の定式化

$$V_z = \frac{1}{N} \sum_{n=1}^N z_n^2$$

$$\therefore z_n = \mathbf{u}^T \mathbf{x}_n$$

$$= \frac{1}{N} \sum_{n=1}^N (\mathbf{u}^T \mathbf{x}_n)^2$$

$$= \mathbf{u}^T \left(\frac{1}{N} \sum_{n=1}^N \mathbf{x}_n \mathbf{x}_n^T \right) \mathbf{u}$$

$$= \mathbf{u}^T \mathbf{V}_x \mathbf{u}$$

主成分分析：問題の定式化

$$= \mathbf{u}^T \mathbf{V}_x \mathbf{u}$$

$$\mathbf{V}_x = \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n \mathbf{x}_n^T \quad \Sigma = \begin{bmatrix} E[(X_1 - \mu_1)(X_1 - \mu_1)] & E[(X_1 - \mu_1)(X_2 - \mu_2)] & \cdots & E[(X_1 - \mu_1)(X_n - \mu_n)] \\ E[(X_2 - \mu_2)(X_1 - \mu_1)] & E[(X_2 - \mu_2)(X_2 - \mu_2)] & \cdots & E[(X_2 - \mu_2)(X_n - \mu_n)] \\ \vdots & \vdots & \ddots & \vdots \\ E[(X_n - \mu_n)(X_1 - \mu_1)] & E[(X_n - \mu_n)(X_2 - \mu_2)] & \cdots & E[(X_n - \mu_n)(X_n - \mu_n)] \end{bmatrix}$$

は分散共分散行列

分散…データのばらつきを表す $\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2$

共分散…二組の対応するデータの間の関係を表す $\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})$

分散共分散行列…n変量のデータから得られるn個の分散とn(n-1)個の共分散を
n次の正方行列にまとめたもの

主成分分析：問題を解く

■ 式 $V_Z = \mathbf{u}^T \mathbf{V}_X \mathbf{u}$ を \mathbf{u} に関して最大化する

- \mathbf{u} は単位ベクトルなので $\mathbf{u}^T \mathbf{u} = 1$
- ラグランジュの未定乗数 λ を導入して

$$L(\lambda, \mathbf{u}) = \mathbf{u}^T \mathbf{V}_X \mathbf{u} + \lambda(1 - \mathbf{u}^T \mathbf{u})$$

\mathbf{u} に関して微分して0と置くことで

$$2\mathbf{V}_X \mathbf{u} - 2\lambda \mathbf{u} = 0$$

$$\mathbf{V}_X \mathbf{u} = \lambda \mathbf{u}$$

→分散共分散行列 \mathbf{V}_X の固有値問題に帰着する

主成分分析：問題を解く

■ 式 $V_Z = \mathbf{u}^T \mathbf{V}_X \mathbf{u}$ を \mathbf{u} に関して最大化する

$$\mathbf{V}_X \mathbf{u} = \lambda \mathbf{u}$$

→ 分散共分散行列 \mathbf{V}_X の固有値問題に帰着する

両辺に \mathbf{u}^T をかける。また、 $\mathbf{u}^T \mathbf{u} = 1$ より、

$$\lambda = \mathbf{u}^T \mathbf{V}_X \mathbf{u}$$

これはもともと最大化したいと考えていた V_Z にほかならないので、求めている問題は固有値問題を解いたときの最大固有値が最大値となり、その最大値を満たす \mathbf{u} は最大固有値に属する固有ベクトルである