

# テキストマイニングを用いたコンサルティングサービスの支援手法 —対応分析とDEA判別分析による 不正予測—

1515050 山本聖也

# 発表の流れ

1. はじめに
2. 研究手法
3. 計算機実験
4. 結論

# 1. はじめに

- 中小企業向けコンサルティングの支援体制の需要が高まってきている.

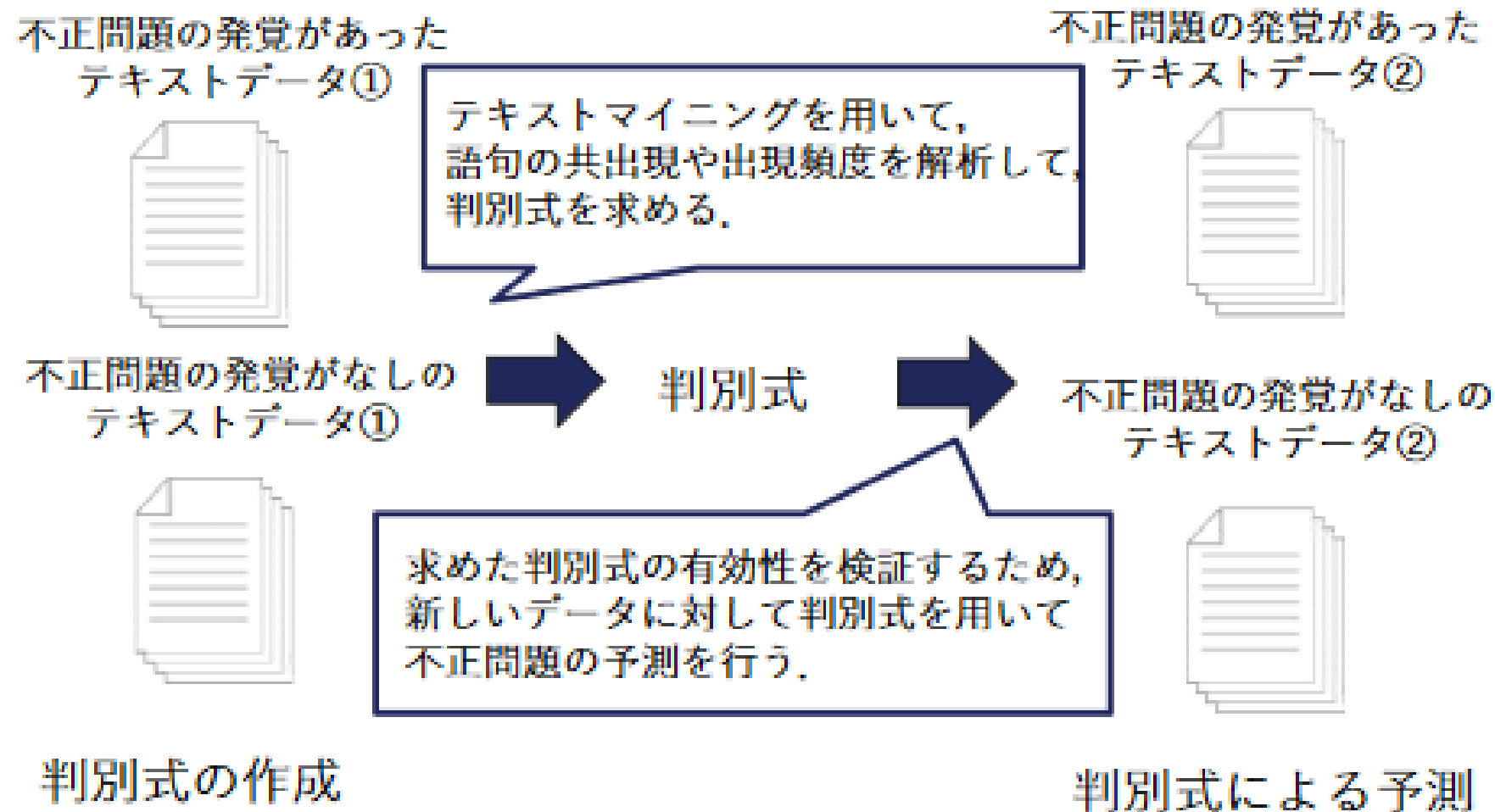


- 専門知識の欠如しているコンサルタントであっても専門性の高いサービスの提供を！

# どうやって？

- 顧客から受けた相談事項のコミュニケーションを記したテキストデータを解析する.
- これによりクライアント企業における不正問題発生の有無を判別する

## 2. 研究手法



I. テキストデータの分類

不正発覚の有無,  
タイミングで分類

II. 形態素解析

不要な語をノイズとし  
て省く

III. 対応分析

IV. 要因となる語の抽出

V. DEA判別分析

VI. 新しいデータを用いた判別率の検証

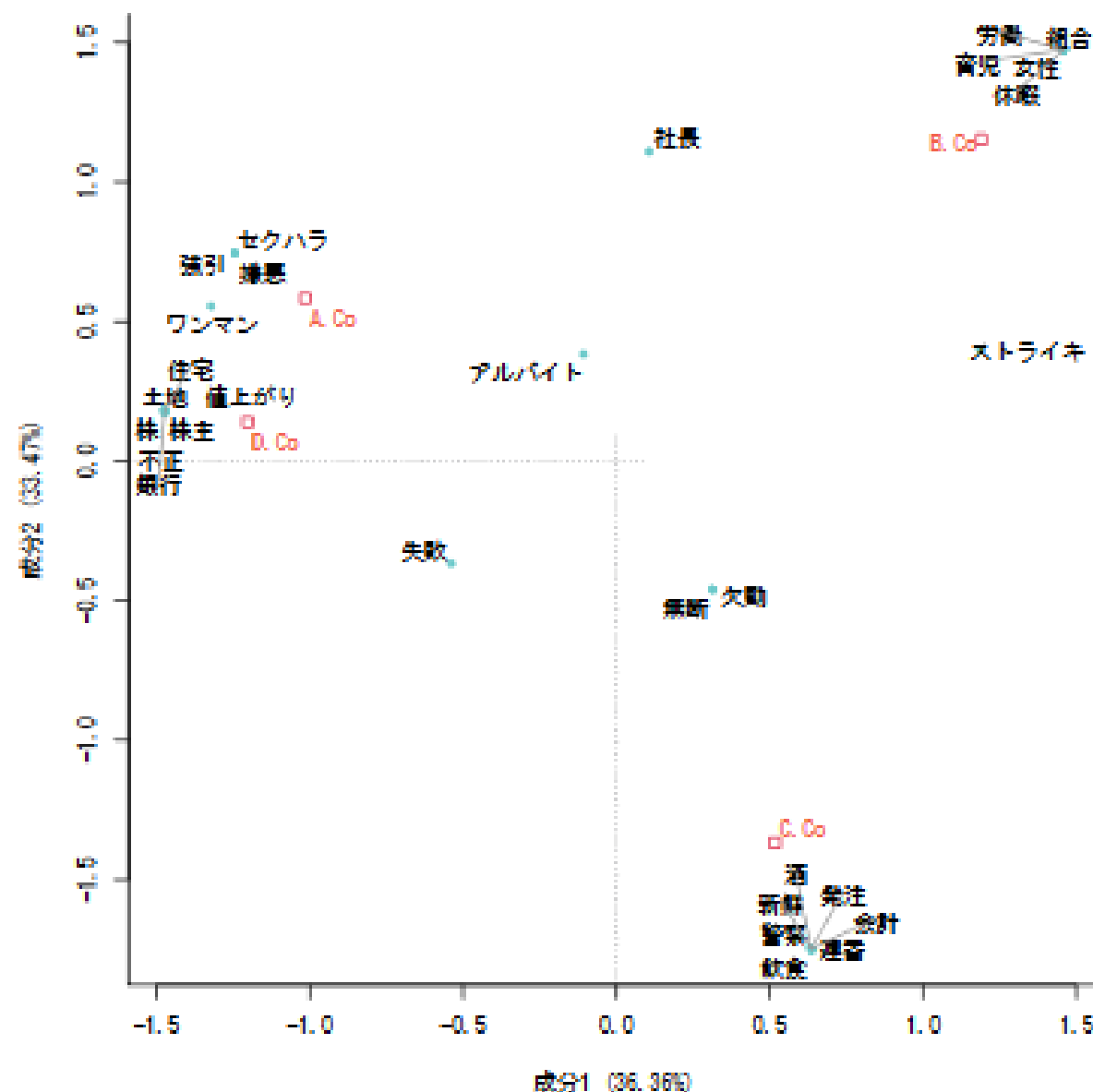
# • 対応分析

データ票の行や列に含まれる情報を小数の成分に圧縮する手法

行に語句，列に企業  
重みとなる変数をそれぞれ設定し，それらの  
変数の相関が最大となるように計算を行う

	語句	AAA	BBB	...	KKK
企業名		$a_1$	$a_2$	$\cdots$	$a_K$
A Co.	$b_1$	$t_{11}$	$t_{12}$	$\cdots$	$t_{1K}$
B Co.	$b_2$	$t_{21}$	$t_{22}$	$\cdots$	$t_{2K}$
$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\ddots$	$\vdots$
N Co.	$b_N$	$t_{N1}$	$t_{N2}$	$\cdots$	$t_{NK}$

各成分のサンプルスコア, カテゴリースコアの値を求め, 散布図に配置する. これにより語句の対応関係が可視化できる





# 要因となる語の抽出

- 対応分析の結果から全成分を考慮し、各グループ $G_1 G_2$ において単語 $i$ と原点の距離 $d_{iG}$ を式(1)で算出

$$d_{iG} = \sqrt{\sum_{j=1}^{D_G} \{(x_{ijG} * C_{jG})^2\}}$$

- 続いて(2), (3)より $d$ を更新
- 各グループの要因となる語をDEA分析の変数に用いる

$$\begin{array}{ll} \begin{array}{l} G_1 < G_2 \text{ のとき} \\ \left\{ \begin{array}{l} e_{iG_1} = d_{iG_1} + M - d_{iG_2} \\ e_{iG_2} = \infty \end{array} \right. \end{array} & (2) \end{array} \quad \begin{array}{ll} \begin{array}{l} G_2 \geq G_1 \text{ のとき} \\ \left\{ \begin{array}{l} e_{iG_2} = d_{iG_2} + M - d_{iG_1} \\ e_{iG_1} = \infty \end{array} \right. \end{array} & (3) \end{array}$$

# DEA判別法

- テキストデータを利用して不正予測を行う手法として採用
- 二段階で行われる判別法
- 一段階目

誤判別のデータとどちらに分類されるか **不明なデータを検出**

- 二段階目

一段階で検出されたデータに対して判別分析を行い分類することで **判別の制度を高める**

- Stage1は誤判別を最小化するようになっており判別境界に $\eta$ の幅をもたせることでグループ1、グループ2のどちらに半別されるか不明なデータを洗い出すことができる.

stage1

$$\begin{aligned}
 \min \quad & \sum_{j \in G_1} S_{1j}^+ + \sum_{j \in G_2} S_{2j}^- \\
 \text{s.t.} \quad & \sum_{i=1}^k \lambda_i z_{ij} + S_{1j}^+ - S_{1j}^- = d + \eta \quad (j \in G_1) \\
 & \sum_{i=1}^k (\lambda_i^+ - \lambda_i z_{ij} + S_{2j}^+ - S_{2j}^-) = d \quad (j \in G_2) \\
 & \sum_{i=1}^k |\lambda_i| = 1 \\
 & S_{1j}^+, S_{1j}^-, S_{2j}^+, S_{2j}^- \geq 0
 \end{aligned} \quad ($$

- Stage1で得られた最適解を $\lambda_i^*$ と $d^*$ としたとき, 要因 $i$ の出現回数 $z_{ij}$ は次の判断基準により5種類に分類される

- $C_1, C_2$ は正しく分類されたこと

- 誤判別の集合

$G_1 \cap R_2, G_2 \cap R_1$

$R_0$ はオーバーラップ領域にある

データ

$$R_1 = \left\{ j \in G \mid \sum_{i=1}^k \lambda_i^* z_{ij} \geq d^* + \eta \right\}$$

$$R_0 = \left\{ j \in G \mid d^* + 1 > \sum_{i=1}^k \lambda_i^* z_{ij} > d^* \right\}$$

$$R_2 = \left\{ j \in F \mid d^* \geq \sum_{i=1}^k \lambda_i^* z_{ij} \right\}$$

$$C_1 = \{ j \in R_1 \mid j \in G_1 \}$$

$$C_2 = \{ j \in R_2 \mid j \in G_2 \}$$

- stage2では誤判別となったデータとオーバーラップ領域に存在するデータへの対処を行う
- オーバーラップ領域に存在していたデータの判別が可能となっている

stage2

$$\begin{aligned}
& \min \quad \sum_{j \in G_1 \cap (R_0 \cup R_2)} S_{1j}^+ + \sum_{j \in G_2 \cap (R_0 \cup R_1)} S_{2j}^+ \\
& \text{s.t.} \quad \sum_{i=1}^k \lambda_i z_{ij} \geq d + \eta \quad (j \in C_1) \\
& \quad \sum_{i=1}^k \lambda_i z_{ij} + S_{1j}^+ - S_{1j}^- = c \quad (j \in G_1 \cap (R_0 \cup R_2)) \\
& \quad \sum_{i=1}^k \lambda_i z_{ij} + S_{2j}^+ - S_{2j}^- = c \quad (j \in G_2 \cap (R_0 \cup R_1)) \\
& \quad \sum_{i=1}^k \lambda_i a_{ij} \leq d \quad (j \in C_2) \\
& \quad \sum_{i=1}^k |\lambda_i| = 1 \\
& \quad d \leq c \leq d + \eta \\
& \quad S_{1j}^+, S_{1j}^-, S_{2j}^+, S_{2j}^- \geq 0
\end{aligned} \tag{10}$$

- Stage2の最適解より, 次の基準で判別される

$$\sum_{i=1}^k \lambda_i^* z_{ij} \geq c^*$$

$$\sum_{i=1}^k \lambda_i^* z_{ij} < c^*$$

- 上式のと看企業jはG1
- 下式のと看企業jはG2

# 3. 計算機実験

- 実験内容

クライアント企業の業種や地域の偏りによって判別式に及ぼす影響を確かめる

条件 1 業種に偏りを持たせる

1. 製造業のみ
2. 卸売業のみ

条件 2 地域に偏りを持たせる (a 地域のみ)

条件 3 ランダム (5 回)

条件 4 不正問題発覚ありのテキストデータに  
業種・地域をそろえる

# 実験結果

学習良

抽出条件	条件 1-1	条件 1-2	条件 1 平均	条件 2
学習 (%)	97.5	90	93.75	95
予測 $G_1$ (%)	35	75	60	65
予測 $G_2$ (%)	55	30	42.5	30

抽出条件	条件 3-1	条件 3-2	条件 3-3	条件 3-4	条件 3-5	条件 3 平均
学習 (%)	90	87.5	92.5	92.5	95	91.5
予測 $G_1$ (%)	70	55	65	60	55	61
予測 $G_2$ (%)	46	50	40	40	35	42.2

予測良



- 条件4の場合の結果

	学習	予測 ( $G_1$ )	予測 ( $G_2$ )	予測平均
判別率 (%)	97.37	80	40	60

- 1の条件より判別率が向上
- 不正問題発生の有無で判別するための判別式を作成できたため
- 1～3と同様にG1の判別率の方が高くなった

## 4. 結論

- 業種・地域によって文章に特徴があることが分かった．判別式を作成する際は業種や地域の差のノイズが判別分析に影響しないように，各グループの業種と地域の比率をそろえるべき
- 不正発覚ありのグループよりも不正問題発覚なしのグループの判別率が低下している．これは不正問題発覚なしのグループの中に不正問題発生企業が存在すること．