

# SASによるテキスト・アナリティクス入門

津田高治  
アナリティクス推進グループ

## Introduction of SAS Text Analytics

Takaharu Tsuda  
Analytics business development

## 要旨:

以下に要旨を記載(100文字以内)

テキスト・アナリティクスの分類問題(教師付き・教師なし学習)やトピック抽出・次元圧縮などのテーマに対するSASの方法論の紹介を行う。

ベイジアンネット, ベクター空間モデル, SVDなど、ビジネス現場での事例紹介と共にその方法・適用例・結果を議論する。

キーワード: 続けてキーワードを記載

テキスト・アナリティクス, ベイジアンネット, 正規表現, SVD, 活用事例

## SASの考える先進アナリティクス

### 先進アナリティクス

#### 時系列予測

時系列を分析し、将来の意思決定に貢献する洞察を導く

#### データマイニング

データから関連性を発見し将来に関する精確な予測を実現する

#### テキストマイニング

意味ある話題の発見、ソーシャルメディア、コールセンターから知見を発見し、さらに深い分析へ結びつける  
そのことでビジネスへの貢献をする

#### 最適化

ビジネスのダイナミクスを捉え、リソースの最適配置を提案する



アナリティクス基盤

# テキスト・アナリティクスとは

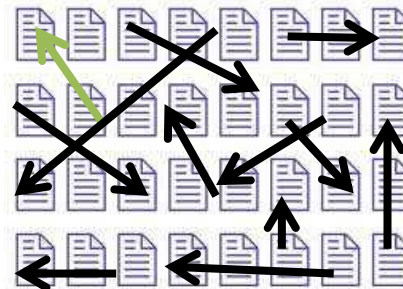


## テキスト・アナリティクスとは

大量の文書から



隠れたパターン（知見）  
を発見



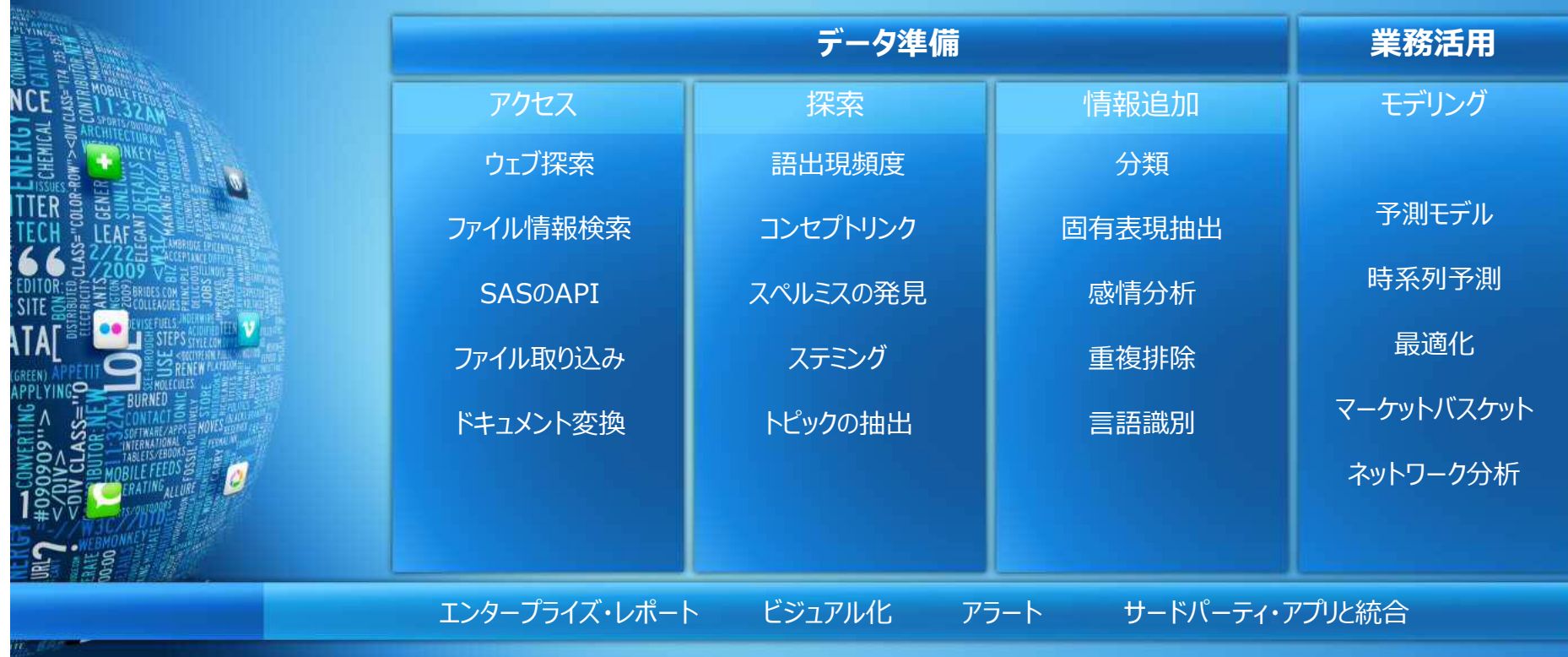
- 業務データと結びつける
- データマイニングする

ビジネス活用をする

でも知見で終わりではない

## SASのテキスト・アナリティクスとは

SAS  
テキストマイニング



## SASのテキスト・アナリティクスとは

収益拡大

マーケティング  
ROI最大化

人件費削減

作業品質  
の標準化

コスト削減

イメージアップ

活用



マーケティング

顧客の興味エリアを理解し、リコメンデーションする  
(協調フィルタリング)

コールセンター

クレーム/要望などの感情・緊急度、製品区分等でテキストを自動振分け

セキュリティ・ディフェンス

事件の関係者や場所・モノ(クルマ、武器)などの関係性を描く、それと自社との関係をアラート

製品品質・製品開発

品質業務の高度化をテキストに基づく予測分析で行う、評判を製品開発に活かす

分析



SAS® Text Miner

自動分類

事前定義に従って自動振り分け

多様な分類器

ベイズ推定、K近傍等

SAS® Contextual Analysis

自動トピック

書かれている内容に従って文書群を自動トピック抽出

文書レベルでの感情分析

固有表現・事実の抽出

人名・地名・時間・モノなどを抽出する

ドキュメントレベルで感情分析を行う

SAS® Sentiment Analysis

詳細なレベルの感情分析

一つの文書の中で好評・不評が混在する、最も詳細レベルで評判分析

多言語対応

世界の27言語以上で自然言語処理を行う

データ



SNS

コールセンター

フェイスブック

営業日報

ブログ

論文・技術文書

ツイッター

社内文書

ニュース

Eメール

多様なデータアクセス (各種データ形式対応、データソース接続、等)

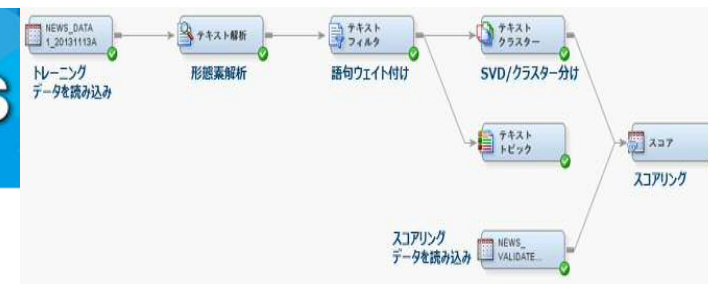
統合メタデータ管理



## 事例と基本技術



# 大規模統計処理に対する期待と現実、そしてSAS



## 導入事例A-声を発見する

### ①元の文章

行数  
||  
文書数

「私はカフェで紅茶を飲むのが毎日です。」

毎朝エクセルシオールカフェでコーヒーを飲んでいる

私は紅茶は飲まないのです

### ②形態素解析の結果（列数が多い）

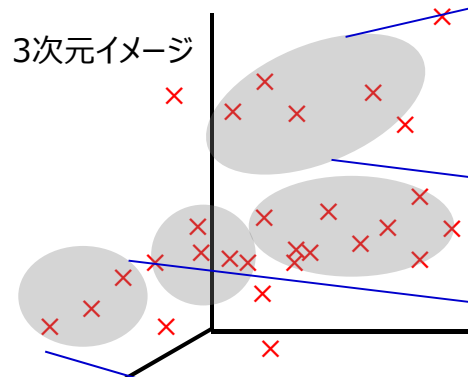
語句の一致										
	私	は	カ フ エ	で	紅 茶	を	飲 む	の が	毎 日	で す
一致 = 1	1	1	1	1	1	1	1	1	1	1
			1			1	1			
	1	1			1		1			1

### ③次元数を落とす

少ない列数（例.100）										
行数    文書数	1	0	1	0	0	1	1	0	..	
	0	0	0	1	1	0	0	0	..	
	0	0	1	0	0	1	1	0	..	
	0	1	0	0	1	0	0	0	..	

### ④n次元空間で近い文書をグループ化⇒「トピック」

低次元座標でクラスタリング



トピック＝携帯電話やスマートホン

- 1 よほど在庫が積み上がってるのか、ドコモ SC-04D Galaxy Nexus
- 2 私が持っている電話の中で、主に通話用に使っているのはウィルコム
- 3 スマートフォンじゃなくて、携帯電話をまだまだ使っている人って多い
- 4 うっかり買ってしまった iPhone 5s のおかげで余った iPhone 5 を、
- 5 先週やっと自宅のネット回線をWIMAXから固定回線に変更しました
- 6 『WIMAX 2+』のサービスが発表されたので、現在WIMAXサービスを
- 7 「わださん、このサービス使えば日本への電話が激安です！」短期

トピック＝B級グルメ

- 1 過去数回にわたって「こんな店ができますよ」的な紹介してきた板橋本
- 2 [image via Frederic Poirat's flickr ] 世界でも広く知られている日本料理
- 3 ベトナム料理にハマって、近所のラーメン屋について色々調べていた
- 4 昨日はムスメと地元パスタ屋ランチなどを食しラーメン。となれば今日
- 5 新潟の燕三条に行ったら何を食べたら良い！？という質問に、速攻で、し
- 6 気がつけば10月。そろそろ鍋の季節がやってまいりますね。恋人と、友
- 7 西台店に行った。今日のところは、特にメルマガなどで告知されてなかつ

### ⑤注目すべき記事を特定

おすすめ⇒メンズ108春もの柿落し

トピック＝渋谷ファッション

## 導入事例A-声を発見する～結果のイメージ

## クラスターのイメージ

- ・ トピック（クラスタリング）は一定の正確性がある
- ・ 特徴が顕著な記事に高い注目度が付く傾向がある
- ・ 情報量の多い記事に高い注目度が付く傾向がある

## 事件など

## パソコン 2013秋冬モデルについて

17060	VAIO 2013年秋冬モデルが発表されて、この夏は
524	VAIO 2013年秋冬モデルを発表！・ノートPCとして
3345	VAIO 2013年秋冬モデルが発表されて、この夏は
9757	VAIO 2013年秋冬モデルを発表！・ノートPCとして 妥協のない
11216	カスタマイズといっても実際にどれを選んだら良いのか迷います
10752	カスタマイズといっても実際にどれを選んだら良いのか迷います
13777	VAIO 2013年秋冬モデルが発表されて、この夏に登場した「VA
12223	VAIO 2013年秋冬モデル
13205	ますます高性能になった
16053	VAIO 2013年秋冬モデル

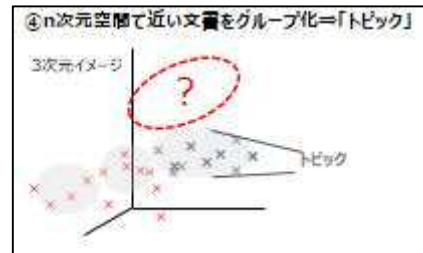
24	1	東京都三鷹市の女子高生殺害事件を受け、各メディアがこぞ
24	2	東京都三鷹市の女子高生殺害事件を受け、各メディアがこぞ
24	3	1:そーきそば 中 ★:2013/10/09(水) .97 ID???0東京都千代田
24	4	(前略)…三鷹殺害 ( )茨城・常陸太田市などで震度4 ( )かつ
24	5	(前略)…、東京・三鷹市内で起きた女子高生刺殺事件で、被
24	6	(前略)…三鷹の女子高生のストーカー事件のアニメオタク報
24	7	神戸拓光 (kobe_takumi55) on Twitter より 不適切ツイートで
24	8	東京都三鷹市の女子高生殺害事件を受け、各メディアがこぞ
24	9	千葉ロッテの選手を名乗る人物が2013年10月10日、三鷹市
24	10	8日夕方、東京・三鷹市の住宅街で、18歳の女子高校生が首

## グルメについて

8153	ベトコンラーメンにハマって、近所のラーメン屋について色々と調べていたら
24820	過去数回にわたって「こんな店ができますよ」的な紹介をしてきた板橋本町の
15490	[image via Frederic Poirot's flickr ] 世界でも広く知られている日本料理。し
37310	ベトコンラーメンにハマって、近所のラーメン屋について色々と調べていたら
6941	昨日はムスメと地元パスタ屋ランチなどを食しノーラーメン。となれば今日は
15693	新潟の燕三条に行ったら何を食べたら良い！？という質問に、速攻で、しか
13212	気がつけば10月。そろそろ鍋の季節がやってまいりますね。恋人と、友人や
7045	西台店に行った。今日のところは、特にメルマガなどで告知されてなかった
24817	前回同様にまだまだ書かなきゃならない店もありますが、気になるネタがあ
19082	新橋駅(鳥森口)から徒歩3分程にあるお店。今日の昼はガッツリ食べたい

## 導入事例A-声を発見する～変化を捉える

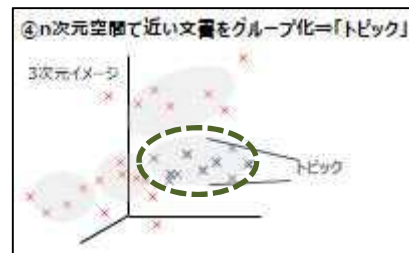
変化  
1



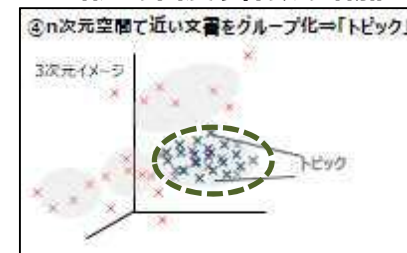
新しい話題の出現



変化  
2



話題の密度や件数の増加



※密度は注目度、件数は盛り上がりを指すと考えられる

## 導入事例B-声を分類する

まず、固有表現を抽出するロジックと結果

「異音」  
に関する固有表現

「場所や状況」  
に関する固有表現

Enter Names

☐ Use same name for both fields

Enter Display Name:  
音

Enter name for internal data files:  
(Latin=1 characters only)

Oto

OK Cancel

CLASSIFIER:音  
CLASSIFIER:異音  
CLASSIFIER:いやな音

Enter Names

☐ Use same name for both fields

Enter Display Name:  
場所や状況

Enter name for internal data files:  
(Latin=1 characters only)

Where\_When

OK

SAMPLE1.k2 - SAS Content Categorization Studio

File Edit View Build Project Category Concept Testing

SAMPLE1

- Japanese-UTF8
- Concepts
  - Top
    - クルマの部位名
    - 場所や状況
    - クルマの異音

CLASSIFIER:より  
CLASSIFIER:から  
CLASSIFIER:フロント  
CLASSIFIER:リア  
CLASSIFIER:とき  
CLASSIFIER:時に  
CLASSIFIER:時の  
CLASSIFIER:ときに  
CLASSIFIER:ときの

「クルマの異音」  
に関する事実抽出

Enter Names

☐ Use same name for both fields

Enter Display Name:  
クルマの異音

Enter name for internal data files:  
(Latin=1 characters only)

Bad\_sound

OK Cancel

PREDICATE\_RULE:(arg1,arg2):(SENT, (ORD, "\_arg1{Where\_When}", "\_arg2{Oto}"))

【テスト結果】  
テキストの中から  
「クルマの異音」  
に該当する  
事実の抽出が行われた



Test File: Go Stop

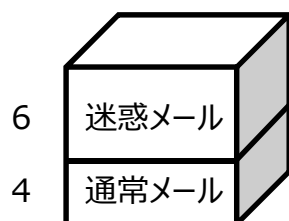
不具合が130件あった後のリコール発表ですから、隠れたかったけど、隠せなくなった。てところでしょうな 私は、本日ロービームのハーネス交換してもらいましたよ。その時に、リアのハブの異音のことも告げると、これまた無償で交換しますって言われました。来週の土曜日交換予定です。みなさんのフラット走行中に後輪のほうから、ゴーーーーって音しませんか？



## 導入事例C-声を分類する

## 1) 予め過去のデータで学習しルールを導いておく

過去のメール集



例. 「アイドル」や「無料」など迷惑メールに特徴的な語句、「統計」や「科学」など通常メールに特徴的な語句を選択(ツールで選択)

過去のデータ(メール)から以下の値を得る

(ツールで算出)

⇒通常メールと迷惑メールの比率は4:6

⇒キーワード別の

通常メールと迷惑メールそれぞれに含まれる確率

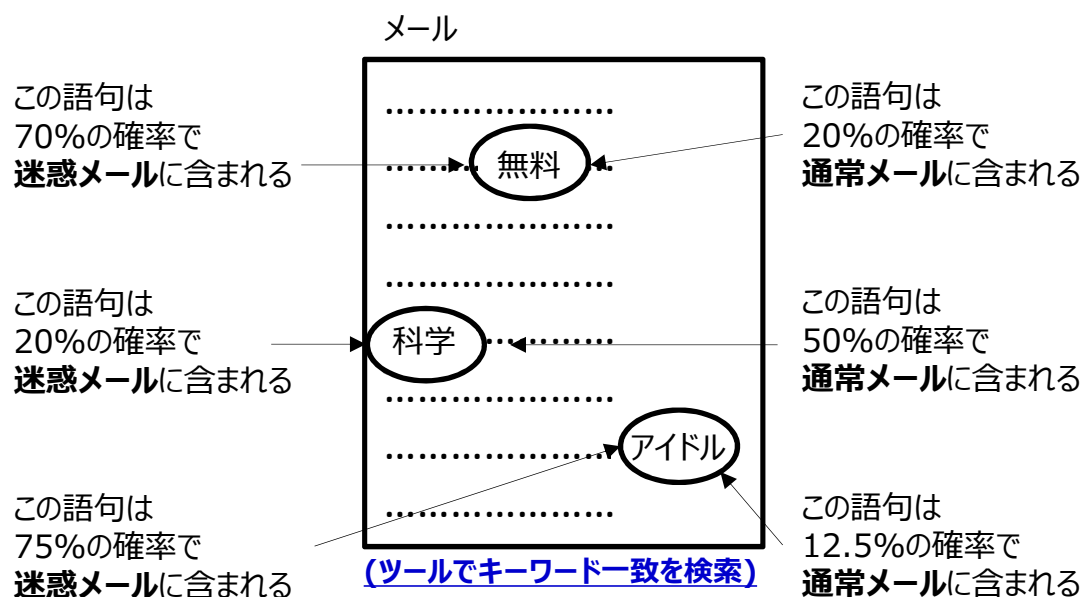
語句	通常メール	迷惑メール
統計	0.8(*1)	0.1
科学	0.5	0.2
無料	0.2	0.7
アイドル	0.125	0.75

(\*1)

通常メールが100通あれば、  
その中の80通に「統計」という語句が含まれる

## 2) ルールを新しいデータに当てはめ、結果を導く

通常メール : 迷惑メール = 4:6



(このメールが)  
通常メールである確率 : 迷惑メールである確率  

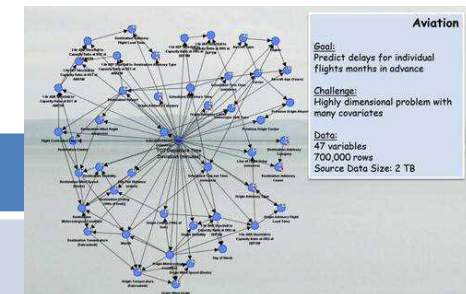
$$= (0.2 \times 0.4) \times (0.5 \times 0.4) \times (0.125 \times 0.4) : (0.7 \times 0.6) \times (0.2 \times 0.6) \times (0.75 \times 0.6)$$

$$= 0.0008 : 0.02268$$
よって迷惑メールである確率の方が高い (ツールで算出)

## 導入事例C-声を分類する～分類精度

下記の精度報告は特定データを用いた例示であり、一般的な精度保証をするものではありません

ベイジアンネットのイメージ図

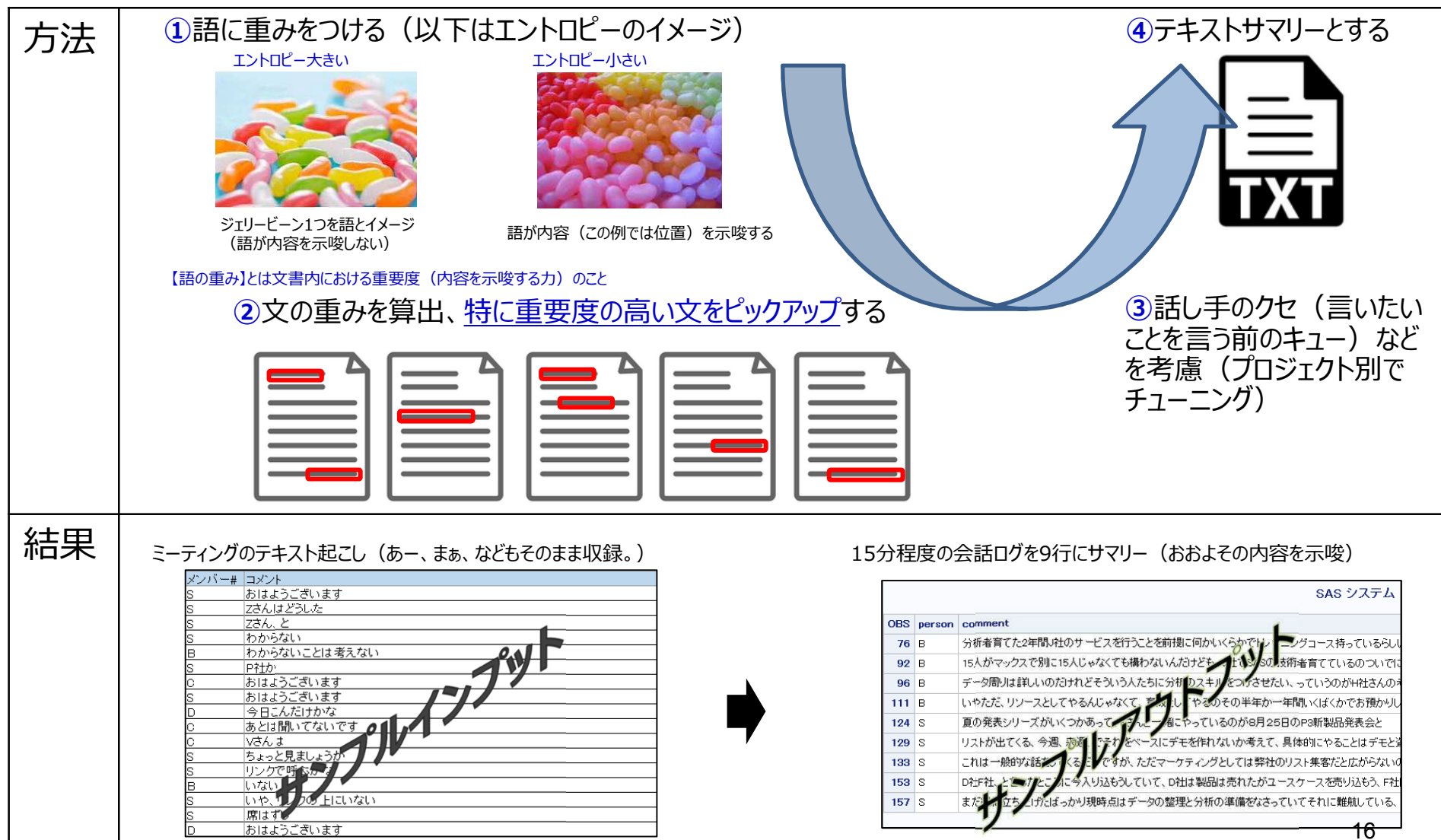


	カテゴリー判別	感情判別
分類精度	95%	91%
分類対象	ニュース、SNS	ニュース、SNS
活用した手法	<ul style="list-style-type: none"> <li>Naive Bayes</li> <li>Tree augmented network (TAN),</li> <li>Bayesian network-augmented naïve (BAN), その他</li> </ul>	左記に加え 係り受けを考慮した分析
分類精度の 算出手法・根拠	<ul style="list-style-type: none"> <li>教師データから分類ルールを機械学習で導く</li> <li>検証データを分類ルールに則って自動分類</li> <li>検証データの分類結果を検証</li> </ul>	
誤分類の原因	同一キーワードが複数の意味で使われることに対応できていない	皮肉表現やスラング、抑制表現に対応できていないため
誤分類の改善	ニュースとSNSの2種類のデータそれぞれにモデル構築するなど、表現の多様性に対応する	多様なニュアンスに対応できる教師データを用意するなど機械学習の学習材料を豊富にする

## ビジネス活用～アプローチ別

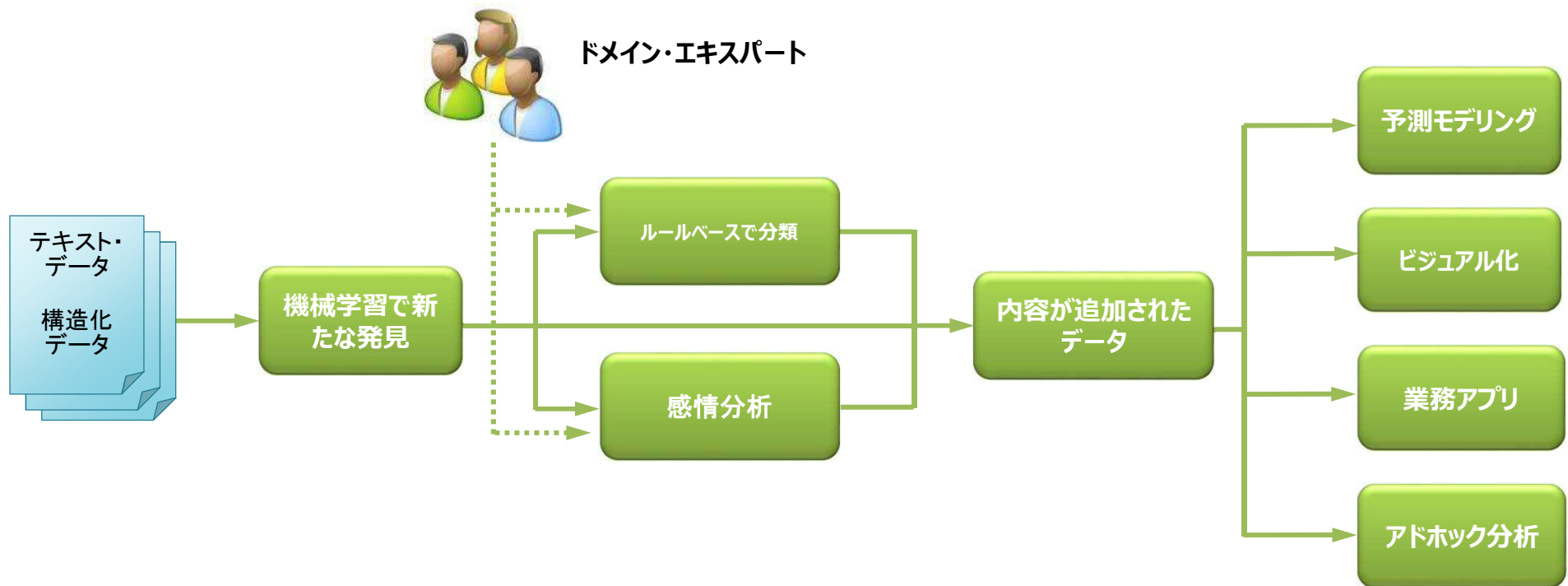
導入事例	技術	分析のタイプ	手法の例	目的	活用例
導入事例A	機械学習	教師なし	SVD	未知のトピックを発見	ネットの声を把握
導入事例B	ルールベース	教師あり	ベイジアンネット	サブジェクト・エキスパートが高い精度で分類ルール作成	製品品質など基幹の業務
導入事例C	機械学習	教師あり	正規表現	一定品質で標準的な分類	顧客の声からリコメンド

## 参考～テキストのサマリー

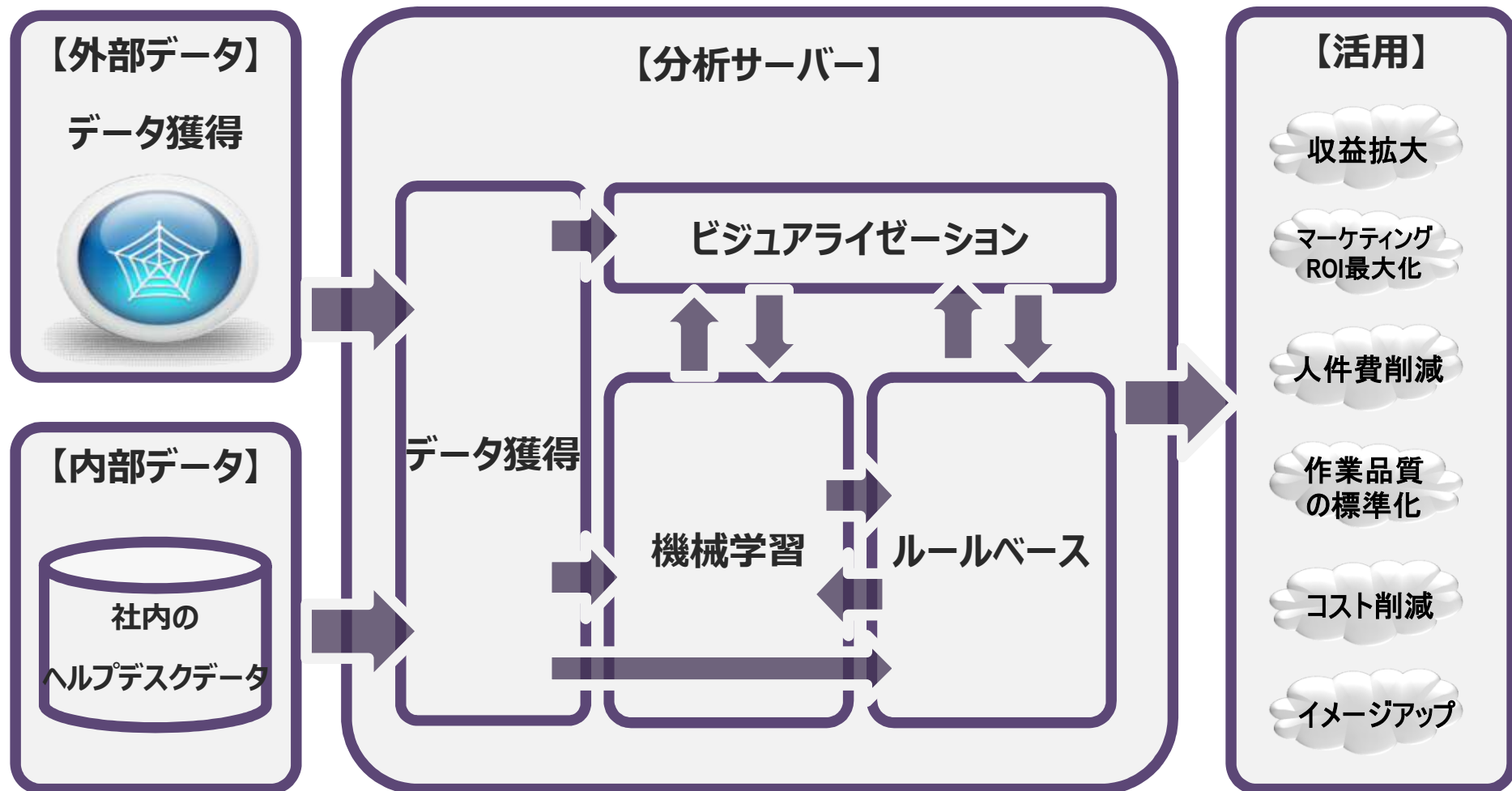




## テキスト・アナリティクスのプロセス



## システムの概念図



## テキスト・アナリティクスの発展

テキストを分析する

テキストから  
知見を得る

テキストから  
アクションする

アクションを精緻化する

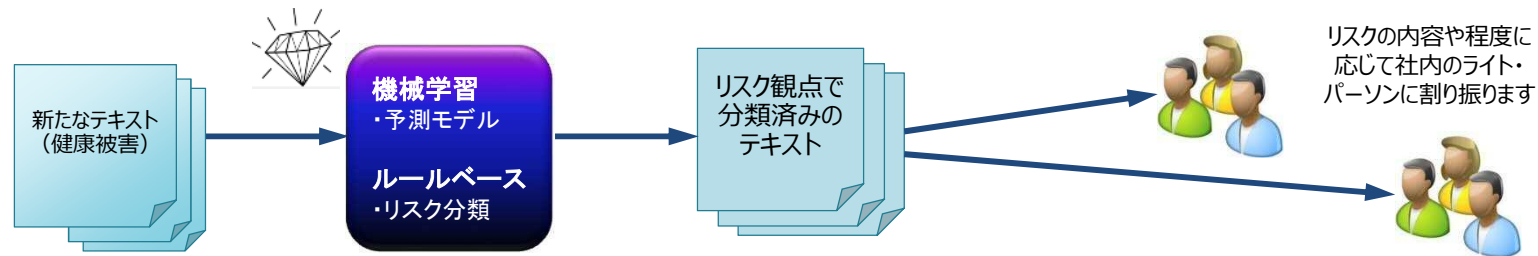
<ul style="list-style-type: none"> <li>・頻度</li> <li>・分類</li> <li>・インデックス付け</li> </ul> <p>例. 顧客を知るために</p> <ul style="list-style-type: none"> <li>• テキストを読む</li> <li>• 人海戦術</li> </ul>	<ul style="list-style-type: none"> <li>・自動トピック抽出</li> <li>・エンティティの抽出</li> <li>・要約</li> </ul> <p>例. 顧客を知るために</p> <ul style="list-style-type: none"> <li>• 興味エリア・マップ</li> <li>• 関係者のマップ</li> </ul>	<ul style="list-style-type: none"> <li>・ビジネスを予測する</li> <li>・リコメンドする</li> <li>・プロファイリングする</li> </ul> <p>例. 顧客を知るために</p> <ul style="list-style-type: none"> <li>• 未来のビジネスを知る</li> <li>• 個人の興味を知る</li> </ul>	<ul style="list-style-type: none"> <li>・パワフルな機械学習</li> <li>・精緻なルールベース</li> <li>・2つの統合と発展</li> </ul> <p>例. 顧客を知るために</p> <ul style="list-style-type: none"> <li>• 未知の声を知る</li> <li>• 既知のルールにする</li> <li>• リコメンド精度の向上</li> </ul>
--	--	---	---

## 事例紹介

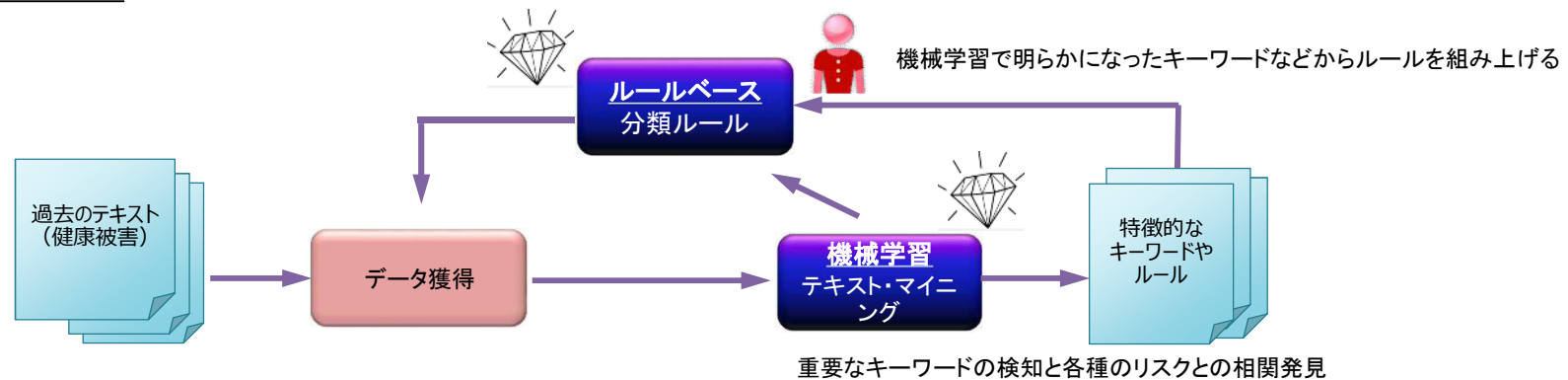


## 事例～製薬業界

**オンライン処理**（顧客の声などテキストが発生した際にリスクスコアを付与）



**学習処理**（機械学習、ルールベースでリスクとの関連付けを構築する）



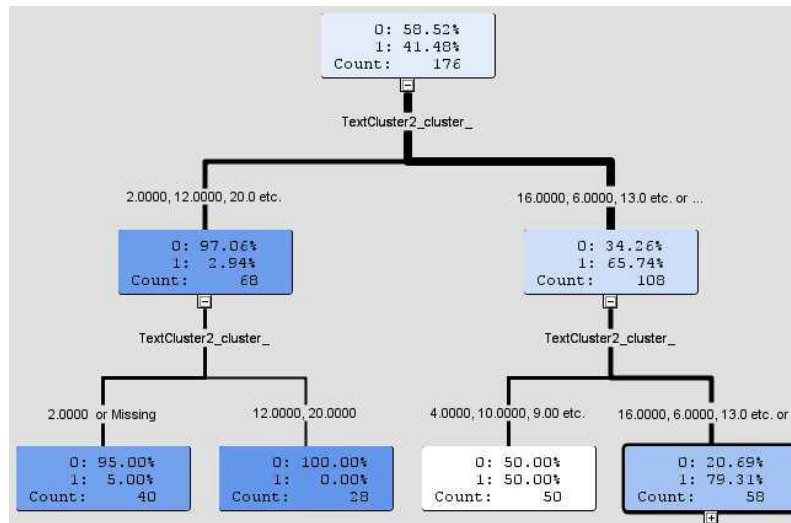
## 機械学習のイメージ

機械学習

**特定リスクの有無**をテキストの内容で説明する

## 活用

新たなテキストに対して、そのテキストの内容から**特定リスクの可能性**を算出する

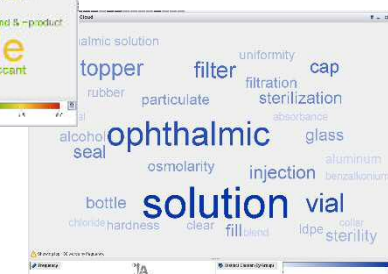


## リスク関連キーワードのイメージ

## ハイ・リスク



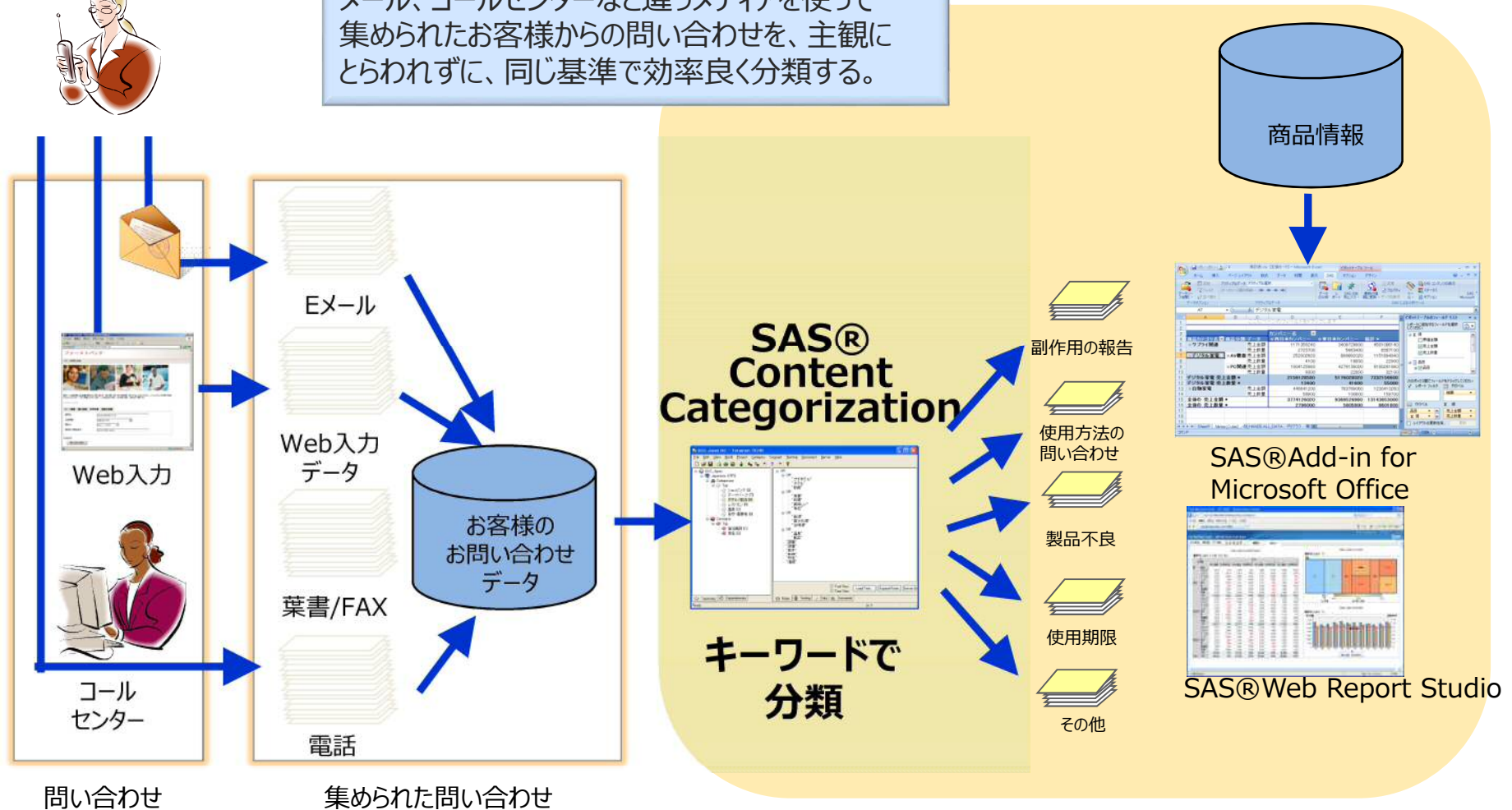
## ロー・リスク



## 事例～製薬業界

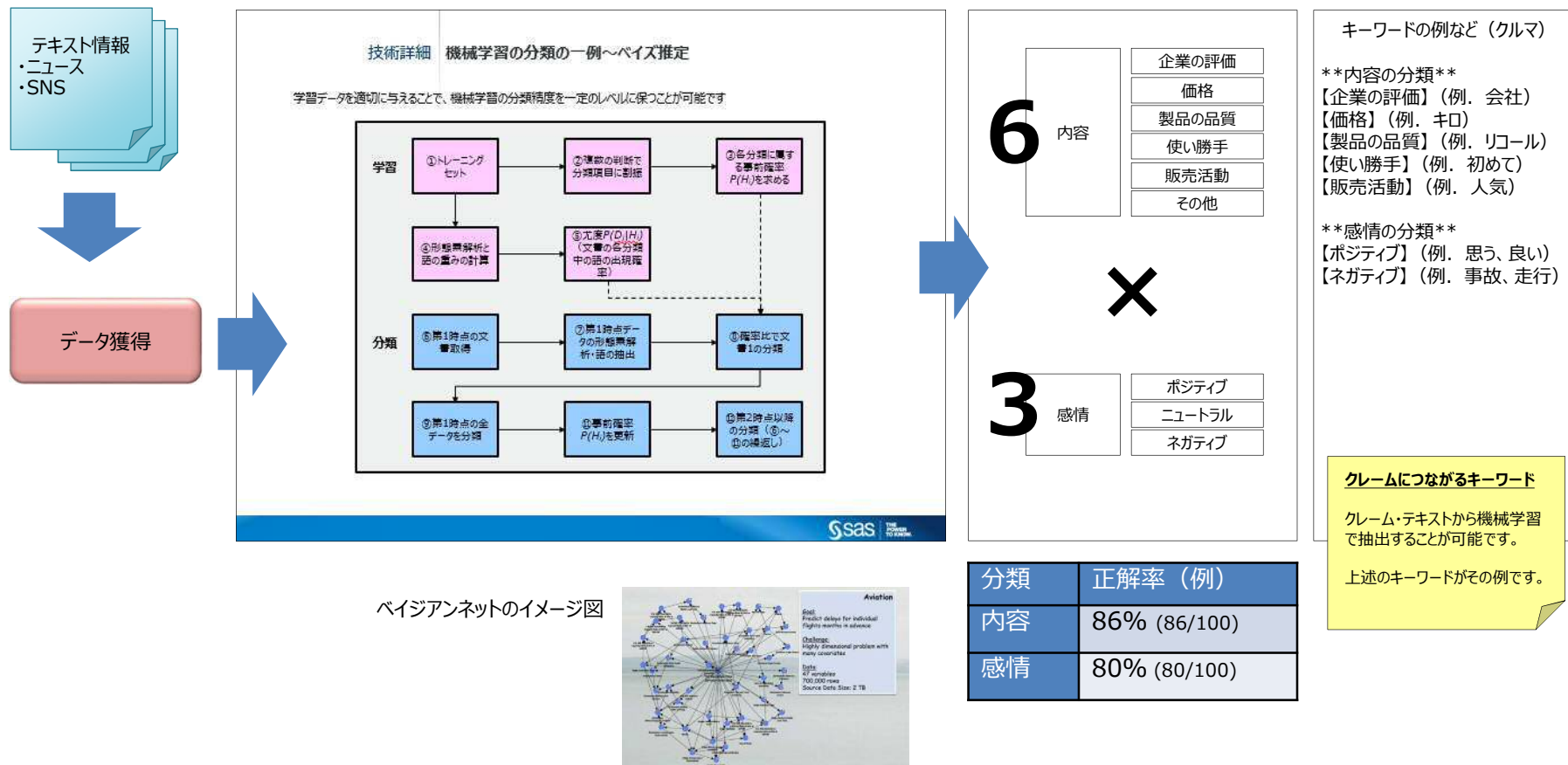


メール、コールセンターなど違うメディアを使って集められたお客様からの問い合わせを、主観にとらわれずに、同じ基準で効率良く分類する。



## 事例～製造業

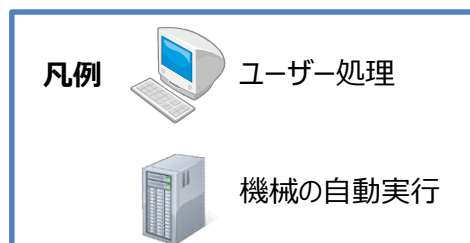
### 機械学習のイメージ (ベイジアンネット)





## 事例～製造業

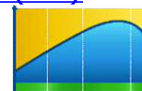
青字：既知の不具合事象に対する処理  
緑字：未知の不具合事象に対する処理



価格.com  
みんなの

**ECC:** Enterprise content categorization  
**EM:** Enterprise Miner  
**EG:** Enterprise Guide  
**VA:** Visual Analytics  
**TM:** Text Miner  
**MM:** Model Manager

③分析実行(EG)



「今を見る」：ワイブル分析

④予測実行(EM)



「将来を予測する」：ニューラルネット

アウトプット・イメージ



⑤レポート作成(VA) ⑥レポート閲覧(VA)

未知～新しい知見の発見

①ルール実行(ECC)

保証

エンジニア

コールセンター

SNS

ルールベース  
分類

既知

既知～分類完了

②書出実行(EG)

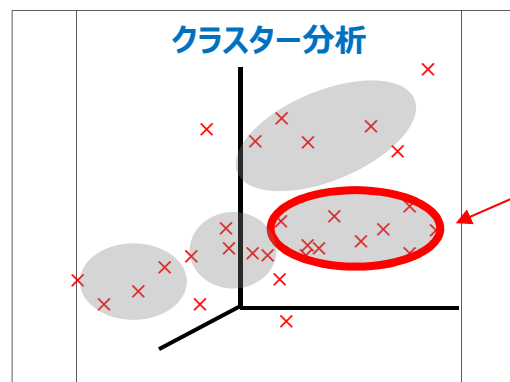
未知※



ルール  
追加

ルール

クラスター分析



⑨事象名付与  
⑩教師テキスト選定

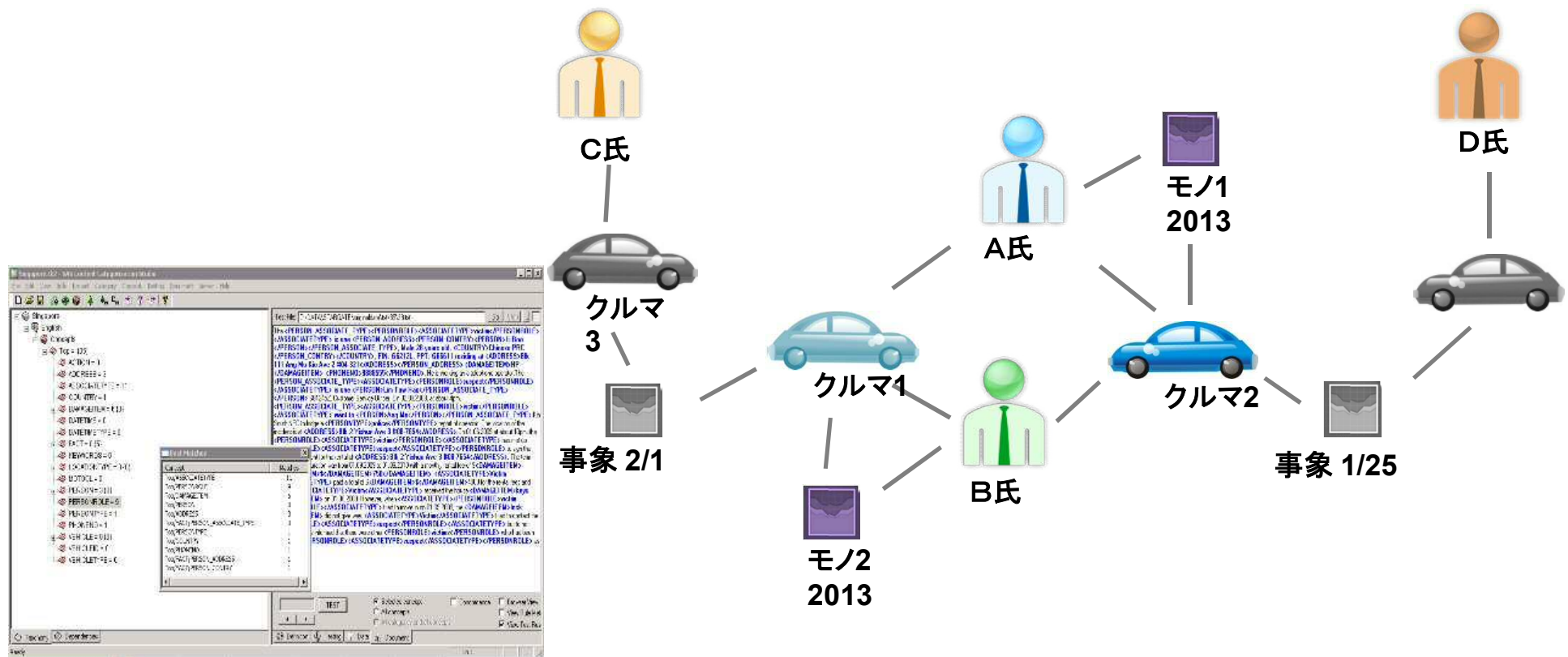
事象名付与  
教師テキスト選定  
(20-30テキスト)

## 事例～通信



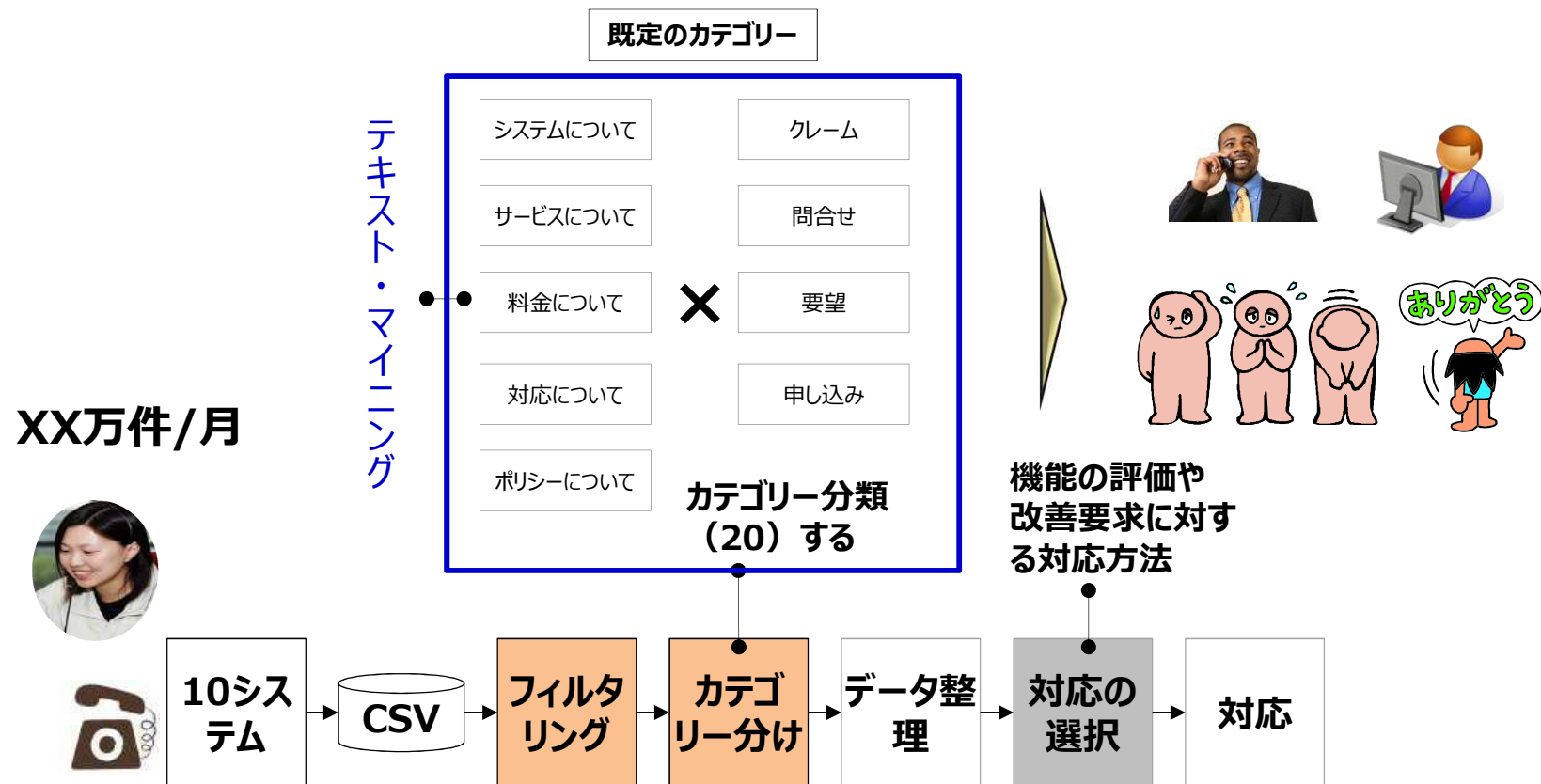
## 事例～公共

- ・ エンティティ- 人物 / 時間 / 場所 / 行動 / モノ（車両）
- ・ 例. A氏が2014年4月9日B氏所有のクルマに乗りXYZショッピング・モールで薬品を買った





## 事例～公共





## 事例まとめ

業種	会社名	業務課題	SASによる解決
製造業	A社 (国内)	テキストデータを用いた顧客へのリコメンドを実施していたが、データ量や海外展開などを考慮すると生産性を大きく向上する必要があった	SAS® Text MinerによりGUIベースのアナリティクスを可能にすることで生産性を大きく向上させることができた。将来の海外展開もより容易な点を評価。
製造業	B社 (国内)	新規顧客の顧客像を明確に描くことができず、マスに対する画一的な訴求しかできていなかった	SAS® Enterprise Miner とText Minerによりインターネット上の消費者の感想と受注データの相関分析をし、新規顧客のプロファイルを浮かび上げらせ、施策に結びつける
製造業	C社 (国内)	製品品質のさらなる向上のために社内データ・インターネットデータを活用したい	SAS® Enterprise Miner とText Miner, SAS® Enterprise Content categorizationにより製品品質と社内データ・SNSデータの相関を導き、よりプロアクティブな活動を容易にした
金融業	D社 (国内)	マーケティングでアナリティクスを活用しているが、さらに詳細・幅広い顧客の理解やマーケティングの施策を実施したい	顧客の取引履歴のテキストデータをEnterprise Content Categorization とText Minerで分析して顧客の興味やセグメントをより深く・広く理解、マーケティングなどの施策に活用
金融業	E社 (国内)	コールセンターのお客様との膨大な会話ログをテキスト化して、ビッグデータ活用につなげたい	テキスト化された音声データを大量データストアに格納した上で、顧客の要望や興味を分析し、営業的顧客対応やマーケティング等に活用する
パブリックセキュリティ	F 国家警察 (海外)	事件の関係者・場所・時間・もの（クルマ・武器・薬物等々）とそれらの関係性を取り調べ記録からマップとして抽出することに多大な時間とワークがかかっていた	SAS® Enterprise Content Categorizationにより取り調べの電子データから左記のエンティティ（関係者・場所・等々）を自動抽出し、関係者マップを描き出すことを可能に。また過去の類似事件の関係者一覧等。