

# NTT DATA

NTT DATA Mathematical Systems Inc.

アジア特許情報研究会 設立10周年記念講演会  
テキストマイニングと深層学習を用いた  
英文技術文書の分析ソリューション

2018年11月30日  
株式会社 NTTデータ数理システム  
岩本 圭介

Trusted Global Innovator

NTT DATA Group

**NTT DATA**

## 1. NTTデータ数理システムのご紹介

- 会社紹介
- ソリューション紹介

## 2. 英文技術文書の分析

- テキストマイニングツール **Text Mining Studio**
- 英語解析における当社の取り組み と 分析例

## 3. データマイニング～深層学習の利用

- データマイニングツール **Visual Mining Studio**  
深層学習デザインツール **Deep Learner**
- 深層学習を用いての文書分類モデル構築

# NTTデータ数理システムのご紹介

- 会社名：株式会社NTTデータ数理システム
- 資本金：5,600万円（NTTデータ100%出資）
- 所在地：東京都新宿区信濃町35 信濃町煉瓦館1階
- 従業員数：約100名（80%が技術者）
- 沿革： 1982年4月 (株)数理システム設立  
2012年2月 NTTデータグループ入り  
2013年9月 (株)NTTデータ数理システムに社名変更



**数理学とコンピュータサイエンスを軸にして、  
社会のあらゆる分野に起こる問題解決のための  
ソリューションを提供する専門家集団です**

# NTTデータ数理システムのご紹介

## ■ 主な業務内容

- ・ ビジネス・アナリティクス領域における、パッケージソフトウェアの開発・販売
- ・ アプリケーション開発
- ・ 分析コンサルティング事業

## ■ 開発・分析対象領域

- ・ AI、機械学習、深層学習
- ・ 数理計画、最適化
- ・ 統計解析、マイニング
- ・ 知識データベース、言語処理、パターン認識
- ・ シミュレーションなどの科学計算



## NTTデータ数理 システムの強み

パッケージ製品・  
受託コンサルティング  
導入実績数2,000社以上

数理科学のプロ集団としての  
豊富な経験・取引実績

自社パッケージによる開発・チューニング

NTTデータグループとの  
事業連携の強化による  
包括的なサービス提供

# 主なパッケージ製品

- データマイニング Visual Mining Studio (需要予測・傾向分析・クラスタリング)  
Big Data Module (大規模データ分析・Hadoop)
- テキストマイニング Text Mining Studio (ポジネガ分析・特徴分析・話題分析)
- 深層学習 デザインツール Deep Learner (ディープラーニング)
- ベイジアンネットワーク BayoLinkS (医療・故障診断)
- 最適化 Numerical Optimizer(スケジューリング・組み合わせ最適化)
- シミュレーション S<sup>4</sup> Simulation System (離散イベント・連続系・エージェント)
- 統計解析 S-PLUS ※ (回帰・検定・多変量解析)
- 統計解析 Visual R Platform (Rユーザー向け分析プラットフォーム)
- CRM分析 CRM Insight (購買傾向、ターゲティング)
- 特許分析 Patent Mining eXpress (特許情報分析)
- 超高速シミュレーション Monaco
- 半導体TCADシミュレータ Paradise World II
- 金融工学 Fiopt (ポートフォリオ最適化・シナリオ発生) ※以外はすべて自社開発です

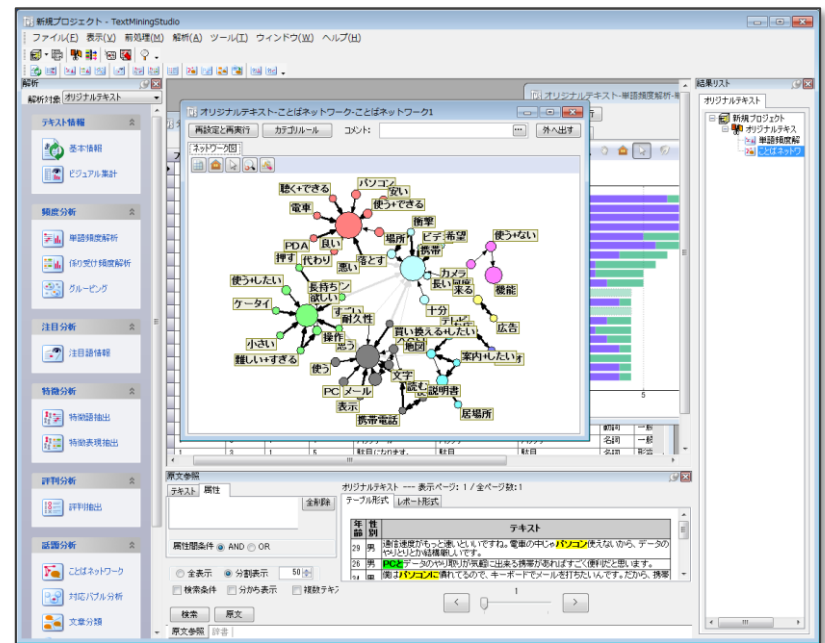
## 2. 英文技術文書の分析





## テキストデータから有益な情報を抽出するための テキストマイニングツール

- **誰にでも**高度なテキストマイニングを
  - マウスの操作のみで分析が可能
  - 豊富な分析機能 と 強力なグラフ機能
- より**自由度**の高い分析を
  - カテゴリ機能等 「意味」 に着目した分析
  - データマイニングツール **Visual Mining Studio** とのシームレスな関係
  - アドオンツール **英語アドオン** の導入で英文の分析も可能



Text Mining Studio (TMS)



文章「従来手法よりも検出精度を向上させることができた。」

見出し語	原形
従来	従来
の	の
手法	手法
より	より
も	も
検出	検出
精度	精度
を	を
向上	向上
さ	する
せる	せる
こと	こと
が	が
でき	できる
た	た
。	。

一次処理である  
形態素解析適用時の結果

## 自動連結機能

文節の単位をとらえ「**欲しい単語**」を適切に抽出する

- ・ 連語の自動判定
- ・ 文節内からキーワード相当部分のみを抽出

単語ID	見出し語	原形	品詞	係り先	態度表現
1	従来	従来	名詞	2	なし
2	手法	手法	名詞	4	なし
3	検出精度	検出精度	名詞	4	なし
4	向上	向上	名詞	-1	可能

## 態度表現機能

テキスト記述者の「**態度**」を適切に抽出する

- ・ 否定、要望、疑問、可能/不可能、容易/困難など  
9種類の態度表現ラベルを文節に自動的に付与

- この「**欲しい単語**」を適切に抽出する機能を、英語の自然言語解析エンジンにも組み込み

文章 "The system should be able to predict traffic jam."

単語ID ▾	見出し語 ▾	原形 ▾	品詞 ▾	係り先 ▾
1	The system	system	名詞	2
2	should be able	should able	形容詞	-1
3	to predict	predict	動詞	2
4	traffic jam	traffic jam	名詞	3
5	.	.	記号	2

## 自動連結機能

単語がばらばらにならず、**日本語の「文節」と類似した感覚で単語をまとめ上げる。**

その結果、ユーザーが辞書を整備することなく  
“traffic jam” といった連語が抽出可能、  
またより自然な係り受け関係を抽出できる

テキストマイニングの  
アウトプットとしては不要な  
**前置詞・冠詞などは  
自動で削除**され、現れない。

データクリーニングの手間が大幅に軽減

# 英語解析における当社の取り組み

- データから連語の候補を自動的に抽出する機能を備え、特別な技術分野依存の語彙にも対応可能

➤ 接続語群の「**連語らしさ**」を統計的指標により評価

登録支援

連語抽出設定

品詞設定

連語全般  
品詞を制限せず連語候補を抽出する  
例) at least, one or more

名詞系  
名詞として使われる連語候補を抽出する  
例) New York City  
City of New York

数字のみの単語を含む連語を抽出する

単語数

2 単語以上

4 単語以下

からなる連語を抽出する

連語抽出

単語検索(スペース区切り)

検索単語

AND検索  OR検索

検索 リセット

全選択 全解除

	選択	連語	品詞	指標値	頻度	単語数
7	<input type="checkbox"/>	control system	名詞一般	48.17	49	2
8	<input type="checkbox"/>	lane change	名詞一般	41.06	42	2
9	<input type="checkbox"/>	experimental result	名詞一般	40.00	41	2
10	<input type="checkbox"/>	advanced driver assistance sy	名詞一般	39.00	13	4
11	<input type="checkbox"/>	automatic operation	名詞一般	35.00	36	2
12	<input type="checkbox"/>	autonomous vehicle	名詞一般	34.11	35	2
13	<input type="checkbox"/>	driver assistance system	名詞一般	30.00	19	3
14	<input type="checkbox"/>	driver assistance	名詞一般	24.74	26	2
15	<input type="checkbox"/>	autonomous driving	名詞一般	24.00	25	2
16	<input type="checkbox"/>	human driver	名詞一般	24.00	25	2
17	<input type="checkbox"/>	assistance system	名詞一般	21.67	25	2
18	<input type="checkbox"/>	Advanced Driver Assistance	名詞同有名詞	21.00	7	4
19	<input type="checkbox"/>	real world	名詞一般	19.14	20	2
20	<input type="checkbox"/>	automatic driving	名詞一般	18.18	20	2

ファイル出力 全 8709 件

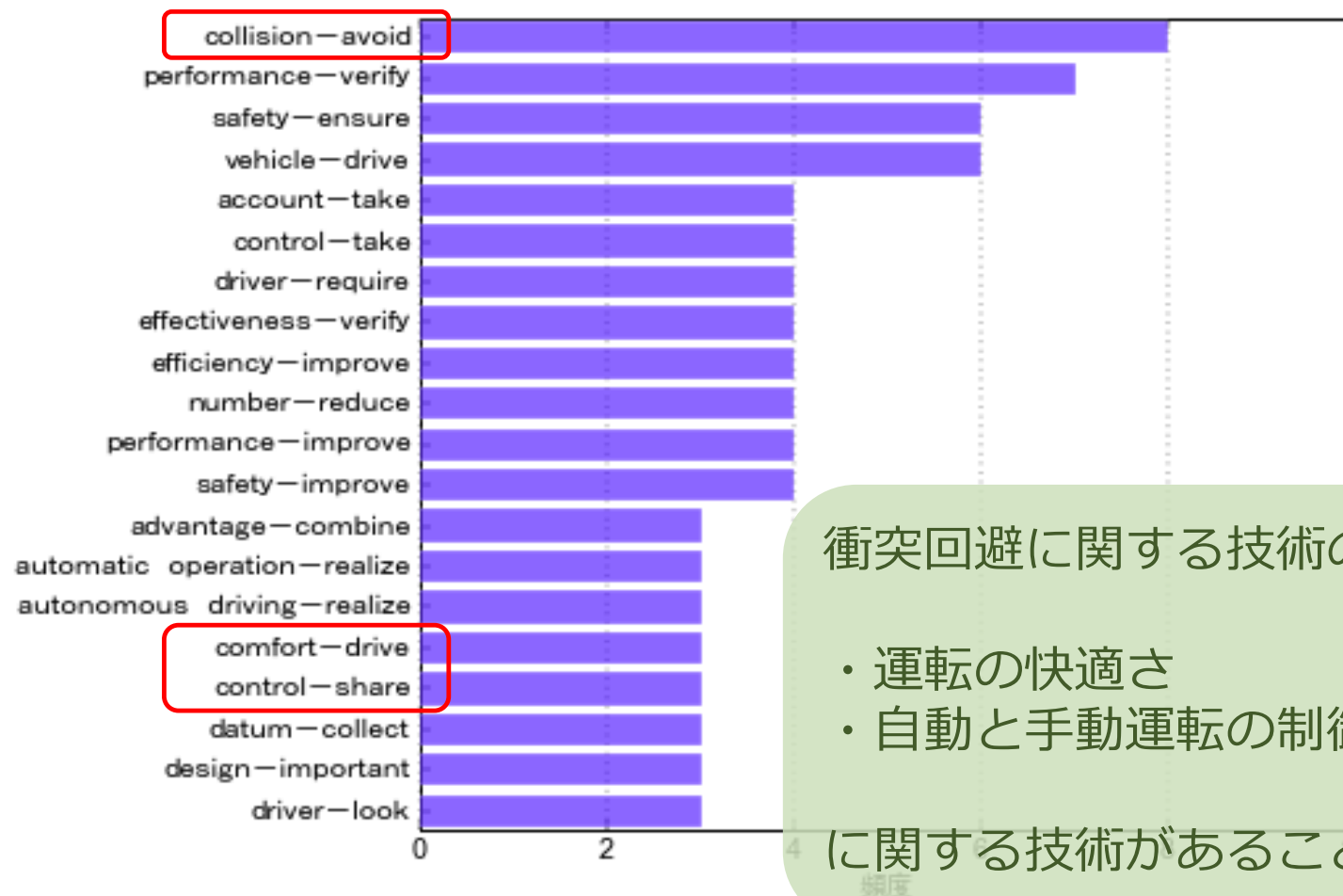
8709 件表示しました

自動抽出された連語群

辞書作成 閉じる

# 分析例：係り受け解析による技術傾向の把握

## ■ 対象データ：「自動運転」に関する論文の英語抄録部分 438件



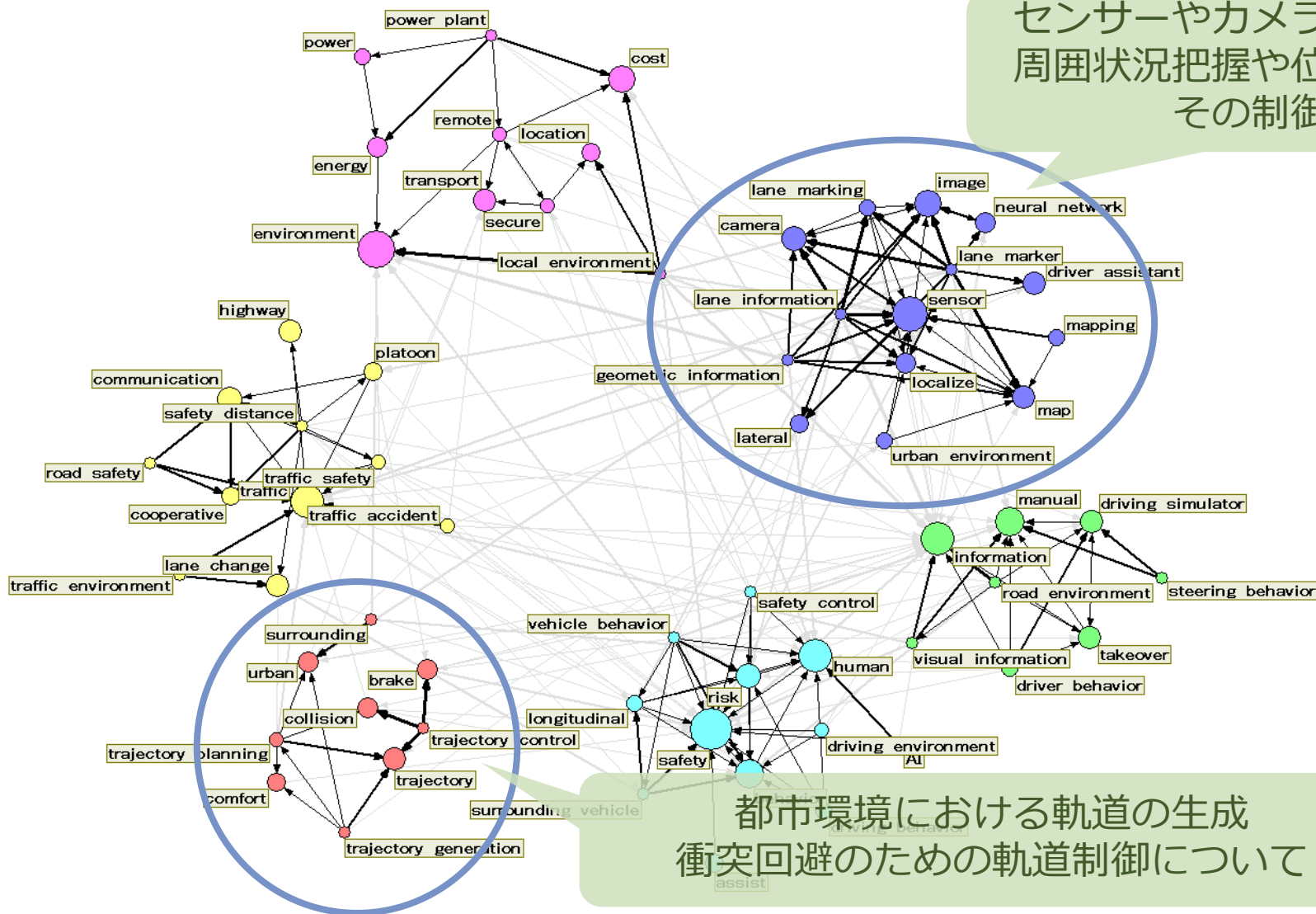
衝突回避に関する技術のほか

- 運転の快適さ
- 自動と手動運転の制御の割り当て

に関する技術があることがわかる

# 分析例：ネットワーク図による技術間の関係図示

センサーやカメラを用いた  
周囲状況把握や位置確認と  
その制御



都市環境における軌道の生成  
衝突回避のための軌道制御について



# 3. データマイニング～ 深層学習の利用

## 文書の教師あり分類

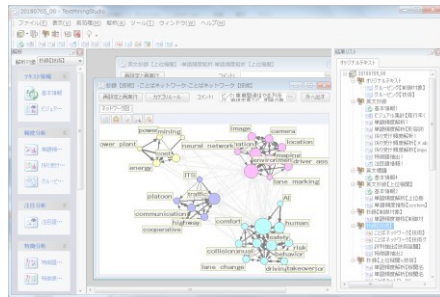
- 機械学習による  
分類モデル構築・ラベル付与

この論文は  
どういう分類？

文書

モデル

技術分類は  
〇〇です

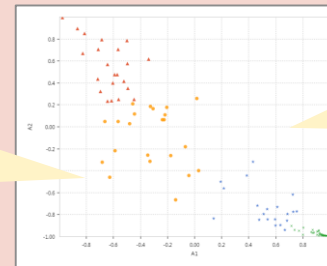


自然言語処理  
内容理解・把握  
ルールベース分類  
(グルーピング)

## 文書を数値情報化して利用

– クラスタリング・可視化

どんな文書が  
近い文書？



どんな分布、  
どんなかたまりを  
なしている？



# テキストマイニング×データマイニングでできること



自然言語処理  
内容理解・把握  
ルールベース分類  
(グルーピング)

## データマイニングツール Visual Mining Studio



※Deep Learner は、当社製品  
- Visual Mining Studio  
- Visual R Platform  
- Big Data Module  
- BayoLinkS  
のアドオンです。  
上記4製品のいずれかが必要です。

# テキストマイニング×データマイニング でできること



**Text Mining Studio**

自然言語処理  
内容理解・把握  
ルールベース分類  
(グルーピング)

	<b>Visual Mining Studio</b>	<b>Deep Learner</b>
機械学習による分類手法	<b>実用的な手法を多数搭載</b> 決定木 Support Vector Machine Neural Network (3層) ...等々	<b>より複雑な Neural Network モデル</b> 多層パーセプトロン <b>テキストの扱いに適したモデル</b> RNN, LSTM
テキストの数値情報化(ベクトル化)及び活用手法	<b>文書ベクトルの作成次元圧縮</b> 主成分分析 対応分析 ...等々 <b>クラスタリング</b> 階層型クラスタリング k-Means ...等々	<b>文書をベクトル化する Neural Network モデル</b> RNN + AutoEncoder
可視化手法	<b>共通分析基盤 Visual Analytics Platform</b> が備える強力な可視化機能を利用可能	

## 文書の教師あり分類

- 従来手法では 単語の有無 で文書の特徴付ける **Bag of Words** という考えを用いることが多いが 一般にこれは語順・文脈の情報が反映されない

区分	TCP-IP	ポート	インストール	停止	Win-98	XP	CD	カスタム	用紙サイズ	設定	...
インストール	1	1	1	1	0	0	0	0	0	0	
インストール	0	0	1	0	1	1	1	0	0	0	
設定	0	0	0	0	1	0	1	1	1	1	

- RNN では可変長の系列情報を自然に扱える
  - 含まれる単語数は各文書で異なるため、テキストは**可変長情報**
  - さらに、その中の**単語出現の順序も考慮**される

## 文書を数値情報化して利用

- 語順・文脈を考慮して文書を**数値（ベクトル）化**できる
  - ベクトル化することで、互いがどれだけ似ているかを定量的に評価できる

# Deep Learner が前提とするデータ

テキスト	分類
視線情報を用いた自動運転制御について	注目
プラント自動運転のための油圧自動制御方式	非注目
協調運転のための視覚情報検出技術	注目



行ID	単語ID	単語	分類
1	1	視線情報	注目
1	2	用いる	注目
1	...	...	注目
2	1	プラント	非注目
2	2	自動運転	非注目
2	...	...	非注目

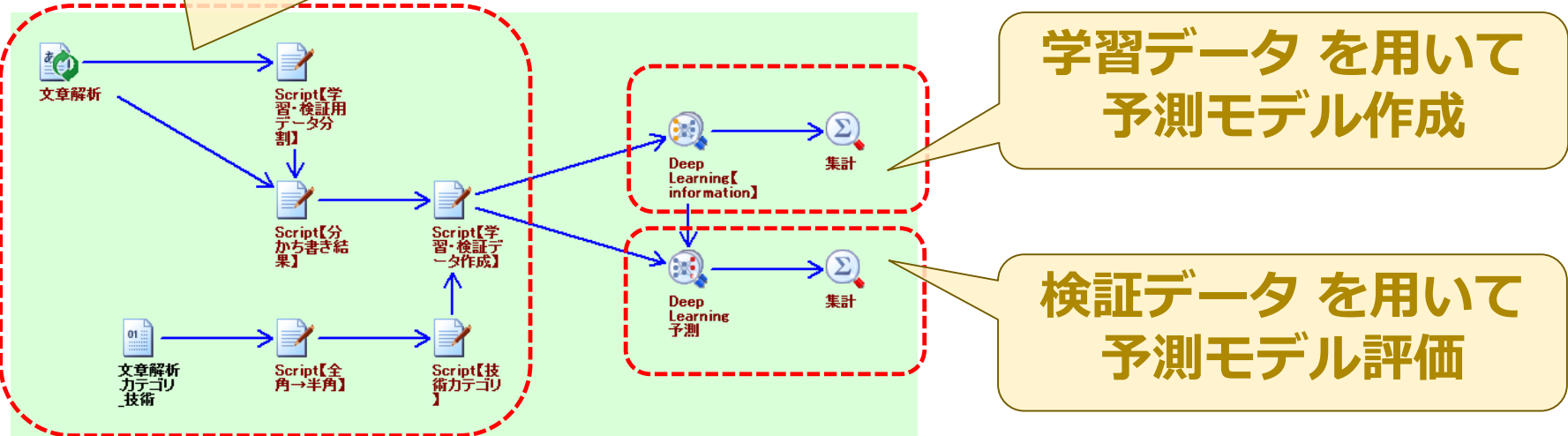
キー列：「1件のデータ」として扱う単位

説明変数列：予測の根拠とする情報

目的変数列：予測対象の情報

## 前処理

- TMS出力結果の分かち書き結果と属性情報をマージし Deep Learner への入力データを作成
- 学習データ と 検証データ を分割



# 業界初！シームレス連携

テキスト  
マイニング

データ  
マイニング

ベイジアンネッ  
トワーク

統計解析

深層学習

最適化

シミュレ  
ーション

CRM

Text  
Mining  
Studio

Visual  
Mining  
Studio

BayoLinkS

Visual  
R  
Platform

Deep  
Learner

Numerical  
Optimizer

S<sup>4</sup>  
Simulation  
System

CRM  
Insight

Patent  
Mining  
eXpress

Big Data  
Module

S-PLUS

数理学とコンピュータサイエンスによる問題解決環境

## Visual Analytics Platform

同一マシンに各種ツールをインストールすることで、  
VAP上ですべてのツールをシームレスに利用することが可能

## 導入前から導入後まで、完全サポート！

導入前 無料サポート

**体験セミナー**  
基本的な操作方法を  
実習を通して理解

**分析個別相談会**  
体験セミナー後、  
分析の進め方を  
技術コンサルタントと相談

**テスト使用**  
(30日間)  
安心の  
操作サポートつき  
(メール対応)

**分析個別  
コンサルティング**  
(1h)  
VMSの疑問点や、  
今後の分析について  
を相談

導入後 保守サポート

**スキルアップセミナー**  
一歩進んだ分析活用法を体得  
※開催有無は各製品に  
よって異なります

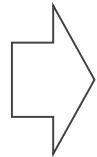
**分析個別相談会**  
スキルアップセミナー後  
日頃のお悩みを  
技術コンサルタントに相談

いつでも操作サポート (メール対応)

## ユーザー様のご予算・スケジュールにあわせた コンサルティング・教育支援の充実！

導入後 各種サポート

マンツーマンで  
一緒に分析を  
してほしい



**スポット  
コンサルティング**

複数人向けに  
分析スキルがつく  
講習をしてほしい



**教育支援**

自分で分析を進める  
上で、定期的なアド  
バイスや作業援助を  
してほしい



**定期個別サポート**

分析作業を任せて作  
業結果を出してほし  
い



**分析受託**



# NTT DATA

NTT DATA Mathematical Systems Inc.

Trusted Global Innovator

NTT DATA Group

NTT DATA