# Deep Learning Extraction for Counterparty Risk Signals from a Corpus of Millions of Documents

# Deep Learning Extraction for Counterparty Risk Signals from a Corpus of Millions of Documents

**Abstract**

China has been experiencing rapid growth over the last decade due to economically friendly reforms and a growing skilled and young population. With this increasing growth, China's interconnectedness with the global economy has increased significantly. In parallel to this economic evolution, technology has experienced rapid acceleration, which has enabled firms and governments to track and record vast amounts of data. The side effect of this unstructured big data growth is that datasets may be polluted, meaning information can be conflicting, missing, and/or unreliable. This creates a gap in the ability to provide transparency to the exposed firms importing from China: both timely early warning signals and wide coverage of small- and medium-sized enterprises (SMEs). We have been able to address this problem for our end-users by using deep learning* to extract information value and opinion from a public corpus to create the needed transparency.

**Moody Hadi**
Group Manager – Innovation & Product Research
Risk Services
S&P Global Market Intelligence

*Deep learning (also known as deep structured learning or hierarchical learning) is part of a broader family of machine learning methods based on learning data representations, as opposed to task-specific algorithms. Learning can be supervised, semi-supervised, or unsupervised.

**S&P Global**

# Outline

- Overview of S&P Global

- Commercial Problem & Solution Overview

- Process Explained
  - Preprocessing
  - Training Set
  - Tokenization
  - Vectorization
  - Neural Networks
  - Benchmarking
  - Model Explanation

- Conclusion

**S&P Global**

# Overview of S&P Global

- Four major business lines:
  - S&P Global Ratings
  - S&P Global Indices
  - S&P Global Market Intelligence
  - S&P Global Platts

- Within S&P Global Market Intelligence
  - S&P Global Risk Services focuses on the varied form of risk management applications including:
    - Credit Risk
    - Market Risk
    - Operational Risk

**S&P Global**

# Commercial Problem & Solution Overview

# Current credit risk models primarily rely on financial statements

Financials are not always available or timely, so there is a need for alternative risk indicators

- Most firms, whether financial or non-financial corporates, have exposures to SMEs worldwide

- Current data sources include financial statements and trade payments

- These sources are not always available or timely

- The challenges are compounded in China, given language barriers and data availability

- There is a need for alternative timely risk signals

We are seeking useful, trusty, and timely first-party source information from companies with a broad reach

# A common problem identified by S&P Global Market Intelligence is that our clients need to review large bodies of unstructured data to find a domain-specific signal

- Problem:
  - Clients need to manage a significant amount of information to extract useful signals
  - The process today is manual, inefficient, and dependent on human expertise
  - While not specific to China, China presents its own idiosyncratic challenges

- Solution:
  - We started with Chinese mandatory company announcements:
    - In line with our strategic focus on China
    - Deals with language complexity
  - We created a robust technology solution to automate the processing, extraction, and opinion classification of documents
  - Our end-to-end workflow processes PDFs and extracts sentiment specifically around a company's financial health
  - The results are signals that clients can use to help sort through massive amounts of data, set their own notifications, and identify insights where they should focus their time
  - Domain specificity is handled by multiple modules in S&P Global Market Intelligence's offering within the same tech stack conditional on the domains being targeted

- The process is not intended to be a credit model:
  - First phase is to extract the opinion that a company discloses about the business within a financial health context
  - Second phase is to combine that opinion with credit scoring from their financial statements

**S&P Global**

# Systematically identifying real-time sentiment from announcements has many advantages

- Efficiency: A firm monitoring only 20 companies would need to sift through over 2,000 announcements and 33,000 pages*

- Effectiveness: One individual reading 130 pages a day will see their performance reduce over time

- Timeliness: Identify signals more frequently in order to analyze information faster (as compared to financial statements that are only available once a year)

*This data is in the public domain.

# CAESAR* is an opinion extraction system focused on the Chinese Market. CAESAR = Chinese Automated Extraction for Sentiment Analysis Reporting

- China is a unique market:
  - Documents are difficult to handle
  - Getting reliable information on the Chinese market is a unique and non-trivial problem
  - Difficult language for non-native speakers

- CAESAR handles the above by:
  - Optical Character Recognition (OCR) engine has image enhancements to handle noisy documents
  - Focused on a corpus of the most reliable documents in China, those directly from the companies' public announcements
  - Language difficulty is handled by a context-sensitive system:
    - OCR extraction on the Chinese corpus
    - Sentence and word detection on Chinese financial corpus, where words and sentence beginnings and endings are language dependent
    - The conversion of tokens via word2vec is Chinese financial health specific

- Reusable assets include:
  - A series of systems that cover the gamut of Natural Language Processing (NLP) from BagOfWords to Neural Networks
  - Context sensitive vectorization
  - A modular extraction system, where components can be replaced to handle other languages and contexts
  - Sentiment attribution to create efficiencies for the analyst, as well as transparency and self-consistency** of the extraction
  - Simplifies localization to the end-users
  - Internal know-how and experience

**S&P Global**

# Detect and reflect: From document to sentiment

Our first phase is to transform unstructured text into sentiment signals relating to financial health of companies

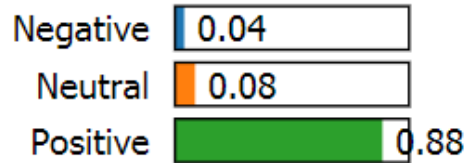Large body of PDFs

Structured data

Sentiment

# We are identifying relevant pieces of text and its embedded opinion
## Context Sensitive Extraction, Attribution, and Actionable Notifications via Sentiment Analysis

True class is: Positive
Prediction probabilities

| | |
|---|---|
| Negative | 0.04 |
| Neutral | 0.08 |
| Positive | 0.88 |

**Text with highlighted words**

证券代码：Θ 证券简称：力尊信通 公告编号
主办券商：中航证券 北京力尊信通科技股份有限公司关于获得高新技术企业证书的公告
本公司及董事会全体成员保证公告内容的真实、准确和完整，没有虚假记载、误导性陈述或者重大遗漏，并对其内容的真实性、准确性和完整性承担 个别及连带法律责任
北京力尊信通科技股份有限公司（以下简称"公司"）于Θ年Θ月Θ日收到由北京市科学技术委员会、北京市财政局、北京市国家税务局、北京市地方税务局联合颁发的《高新技术企业证书》（证书编号：αΘ），发证时间：Θ 年Θ月Θ日，有效期三年
根据相关规定，公司自通过高新技术企业认定并向主管税务机关办理完减免税手续后三年内，可享受国家关于高新技术企业的相关优惠政策，即按Θ税率缴纳企业所得税

Positive Connotation – Green
Negative Connotation – Blue

**Google Translate:**

Stock code: 证券 Securities abbreviation: Lizun Xintong Bulletin number

Sponsored Broker: AVIC Securities Beijing Lizun Xintong Technology Co., Ltd. Announcement on Obtaining High-tech Enterprise Certificate

The company and all members of the board of directors guarantee the truthfulness, accuracy and completeness of the contents of the announcement, and there are no false records, misleading statements or major omissions, and bear individual and joint legal responsibility for the truthfulness, accuracy and completeness of the contents.

Beijing Lizun Xintong Technology Co., Ltd. (hereinafter referred to as "the company") was jointly awarded by the Beijing Municipal Science and Technology Commission, the Beijing Municipal Finance Bureau, the Beijing Municipal State Taxation Bureau and the Beijing Municipal Local Taxation Bureau on the following day. "High-tech Enterprise Certificate" (Certificate No.: αΘ), issuance time: Θ Year Θ Month Day, valid for three years

According to the relevant regulations, the company can enjoy the relevant preferential policies of the state on high-tech enterprises within three years after passing the high-tech enterprise certification and handling the tax reduction and exemption procedures with the competent tax authorities, that is, paying the enterprise income tax at the tax rate.

*Note: the translation is provided for context only. The process deals with the native language announcements.*
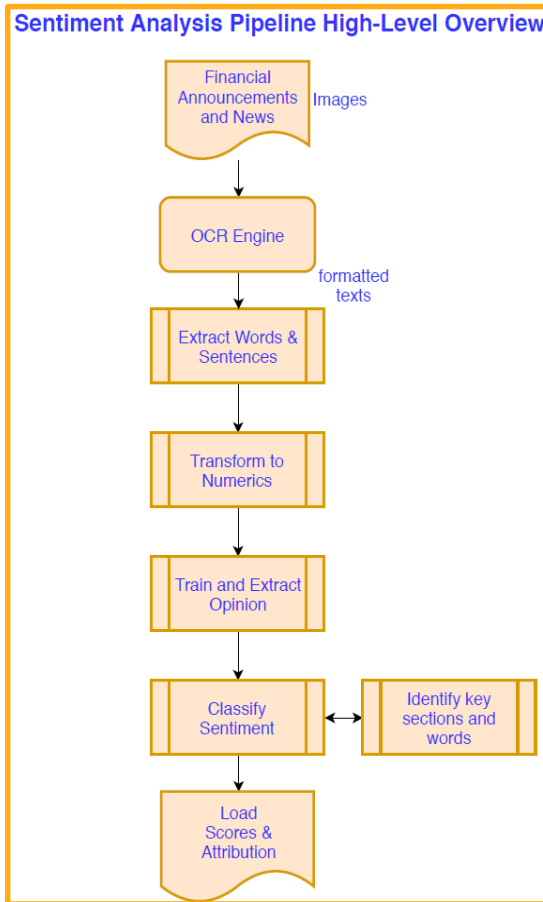
S&P Global

# Data set

We selected a broad set of company announcements

- There is a large body of public disclosures by Chinese companies

- These disclosures are mandatory and factual

- They cover:
  - Firmographics: e.g., board and management resolutions
  - Legal changes: e.g., legal opinions, litigations
  - Financial information: e.g., annual reports, shareholder meetings

- These factual and timely disclosures can be used to extract valuable signals about the financial health of a company

**S&P Global**

# Process

The end-to-end process accesses, parses, processes, and analyzes announcements to extract each article's embedded opinion on the financial health of a company – sentiment classification

**Sentiment Analysis Pipeline High-Level Overview**

```
Financial
Announcements  Images
and News
      |
      v
  OCR Engine
      |          formatted
      v          texts
Extract Words &
  Sentences
      |
      v
 Transform to
  Numerics
      |
      v
Train and Extract
   Opinion
      |
      v
  Classify     <-->   Identify key
  Sentiment           sections and
      |               words
      v
  Load
Scores &
Attribution
```

- Automated collection of announcement documents

- Transform unstructured to structured text via OCR and document analysis

- Pre-process the structured body of text to define sentences, paragraphs, and sections

- Vectorization model

- Model training and validation

- Sentiment classification model

- Sentiment indicators, visualization

**S&P Global**

# Summary of Technical Considerations for NLP techniques

## Technical considerations

| Approach | Advantage | Disadvantage |
|---|---|---|
| Bag of Words | Deterministic<br>Tractable<br>Transparent (Simple Attribution) | Dictionary Maintenance<br>Hard to capture sequencing<br>High Bias<br>Context insensitive |
| Regularized Bag of Words | More targeted to extract veracity<br>Deterministic<br>Tractable<br>Transparent (Simple Attribution)<br>Can control context sensitivity | Less Transparent<br>Requires Recalibration/Higher Maintenance<br>Target Variables need to be defined |
| Neural Networks | Captures Sequencing<br>Context sensitive<br>Low Bias (highly adaptable)<br>Reinforcement allows for automated re-training | Low Transparency (black box)<br>Training is complex (computationally… eg : H/W, Algo) |

# Summary of Business Considerations for NLP techniques

## Business considerations

| Use Case | Advantage | Disadvantage |
|---|---|---|
| Public opinions on companies (e.g. forums, surveys, Twitter….) | Public data (no privacy considerations) | Large body of prior work |
| | Easy to understand the output | Relatively new to extract veracity |
| | Attribution | Public opion data is suspect |
| | Complements existing scoring methods | |
| Public disclosures directly from "informed audience" (e.g. company, analyst research …) | More reliable data - easier to extract information value | To extract veracity more complex approaches are required |
| | Provides wider coverage for companies that do not disclose financials or generally have unreliable information | Target variable needs to be defined (e.g.. credit ranking …) |
| | Public data (no privacy considerations) | Non-deterministic approaches are complex computationally |
| | Easy to understand the output | |
| | Attribution | |
| | Complements existing scoring methods | |
| | Creates new methods for classification given a target | |

S&P Global

# Process Explained

S&P Global

# Pre-processing: A PDF to enhanced images

Image enhancements include removing borders, background colors, and adjusting brightness and contrasts of the text.

Example of a good quality source PDF file (PDF that is constructed from a text document)

*Image enhancement routine*

Example of a medium quality source PDF file (scanned pages that have borders, shadows, and blurred text)

*Image enhancement routine*

Example of a poor quality source PDF file (low quality scanned pages, with noisy background and text)

*Image enhancement routine*

# Pre-processing: Enhanced images to 'structured' text

Enhanced images are converted to a table with coordinates of each token (page, block, paragraph, and line number, in addition to the position and width and height of each token)

| level | page_num | block_num | par_num | line_num | word_num | left | top | width | height | conf | text |
|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 2480 | 3507 | -1 | |
| 2 | 1 | 1 | 0 | 0 | 0 | 1656 | 185 | 445 | 49 | -1 | |
| 3 | 1 | 1 | 1 | 0 | 0 | 1656 | 185 | 445 | 49 | -1 | |
| 4 | 1 | 1 | 1 | 1 | 0 | 1656 | 185 | 445 | 49 | -1 | |
| 5 | 1 | 1 | 1 | 1 | 1 | 1656 | 185 | 123 | 48 | 93 | 公告 |
| 5 | 1 | 1 | 1 | 1 | 2 | 1779 | 185 | 89 | 49 | 93 | 编号 |
| 5 | 1 | 1 | 1 | 1 | 3 | 0 | 0 | 2480 | 3507 | 71 | ： |
| 5 | 1 | 1 | 1 | 1 | 4 | 1905 | 193 | 196 | 36 | 71 | 2018-001 |
| 2 | 1 | 2 | 0 | 0 | 0 | 377 | 307 | 1703 | 419 | -1 | |
| 3 | 1 | 2 | 1 | 0 | 0 | 377 | 307 | 1703 | 419 | -1 | |
| 4 | 1 | 2 | 1 | 1 | 0 | 377 | 307 | 1703 | 50 | -1 | |
| 5 | 1 | 2 | 1 | 1 | 1 | 377 | 307 | 96 | 49 | 93 | 证 |
| 5 | 1 | 2 | 1 | 1 | 2 | 0 | 0 | 2480 | 3507 | 79 | 券 |
| 5 | 1 | 2 | 1 | 1 | 3 | 479 | 308 | 94 | 48 | 92 | 代码 |
| 5 | 1 | 2 | 1 | 1 | 4 | 582 | 331 | 9 | 21 | 88 | ： |
| 5 | 1 | 2 | 1 | 1 | 5 | 627 | 315 | 145 | 36 | 91 | 839771 |
| 5 | 1 | 2 | 1 | 1 | 6 | 979 | 310 | 47 | 42 | 93 | 证 |
| 5 | 1 | 2 | 1 | 1 | 7 | 1032 | 307 | 44 | 49 | 92 | 券 |
| 5 | 1 | 2 | 1 | 1 | 8 | 1083 | 308 | 91 | 48 | 82 | 简 |
| 5 | 1 | 2 | 1 | 1 | 9 | 1185 | 331 | 9 | 21 | 91 | 称 |
| 5 | 1 | 2 | 1 | 1 | 10 | 0 | 0 | 2480 | 3507 | 34 | ： |
| 5 | 1 | 2 | 1 | 1 | 11 | 1231 | 308 | 45 | 47 | 26 | 麒 |
| 5 | 1 | 2 | 1 | 1 | 12 | 1281 | 308 | 46 | 47 | 72 | 腊 |

S&P Global

# Pre-processing: Structured text to text without tables and headers/footers

Tables should be removed from the documents. Using the constructed structured text, we are able to extract all tables from documents with circa 100% accuracy.

Additionally, all headers and footers (like company logo, title of document, legal disclaimer, etc.) that are being repeated across all pages and have no information value have been extracted from the documents.

S&P Global

# Pre-processing: Cleaned documents to 'sectioned' documents

Section 1:



Section 2:



Section n:



| Document Type | Value | Why? |
|---|---|---|
| Long documents (>=100 pages) | Dividing the document into chapters allows for identification of meaningful sections. | • Provides transparency and self-consistency to the end-user on the reasoning of why a document has been assigned a certain classification.<br>• Allows the end-user to preform actions with the notifications. They can "skim-read" large documents, especially when they have a vast body to review. |
| Short documents | Identifying the key words and phrases that had the highest influence on the assigned classification allows for easy to understand reasoning of the classification. | • Provides transparency and self-consistency to the end user on the reasoning of why a document has been assigned a certain classification.<br>• Allows for easy interpretation of the classification, similar to the deterministic approach (Bag of Words). |

S&P Global

# Training set (end state)



Sampling Population:
1.3 million announcements – 2016 to now

Categorically 'without sentiment' (registration, name change, etc.):
240,000

Adjusted Population:
1.1 million

50,000* representative sample

10,000** sample being labeled three times, by three different reviewers randomly chosen from the reviewers pool

40,000 sample being labeled once, by a reviewer randomly chosen from the reviewers pool

Reviewer

Reviewer

Reviewer

Reviewer

\*  First training set is 5,500 articles: we are adding 50,000 to the sample (in progress)

\*\* 1,000 of the original 5,500 articles were also shared with a second group

| Labels | |
|---|---|
| - Positive | - Without Sentiment |
| - Negative | - Illegible |
| - Neutral | |

# Building the end state training set

The following steps have been taken to build the training set:

1.  Interview and select qualified native Chinese financial professionals, with experience in risk management and assessment of Chinese companies.

2.  Identify categories of announcements that are 'without sentiment' by nature (e.g. legal registration, name change, etc.) based on reviewers feedback.

3.  Sample 50,000 announcements from the population, after removing items identified in Step 2. Samples are randomly drawn from year 2016 to present, while maintaining the overall distribution of announcements across categories.

4.  From the 50,000 sampled announcements, randomly choose 10,000 of them as a 'hold-out'* sample. Assign each hold-out sample to three random reviewers, in order to control quality the human reviews.

5.  Send sectioned documents to reviewers and receive back labels for each document under the following categories: positive, negative, neutral, without sentiment, and illegible.

*The hold-out sample is a sample of data not used in fitting a model; used to assess performance.
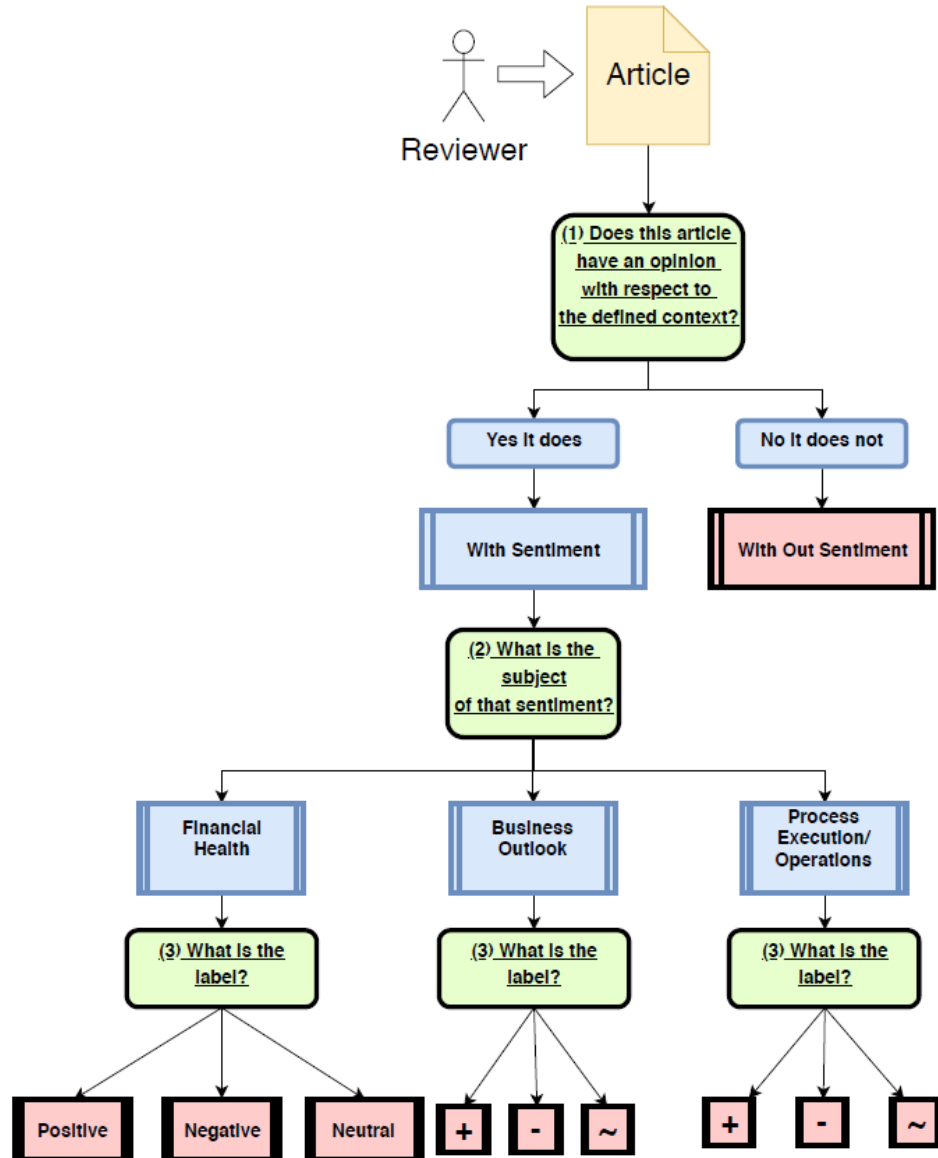
**S&P Global**

# Definitions for training set design and labeling

Objective:

This product extracts and then labels the company's explicit message in its own disclosures on activities that relate to the following context (a group of subjects):

- *Financial Health*:
    - This is the company's opinion on the subject of financial risk and health. In its opinion (if it exists), is there financial risk on liquidity i.e. the company's ability to pay future obligations?
    - This covers connotation on its capital structure, revenues, assets, liabilities.
    - Meritable observations on its financials, for example has there been fraud…?

- *Business Outlook*:
    - This is the company's opinion (if it exists) on the subject of its business model.
    - Business model activities are strategic and tactical activities, such as new products, partnerships/alliances.

- *Process Execution*:
    - This is the company's opinion (if it exists) on the subject of its business process execution.
    - These are operational activities that have no impact on the strategic or tactical activities.

# Reviewer labeling process flow



- Simple examples of *extracting* a label:
  - "Our operations are strong and we see good growth forward …" →Positive
  - "Our operations are proceeding as expected …"→Neutral
  - "Our operations had weakness and the market is softening …"→Negative
  - "Our operations are provided in this statement and we certify …" →Without Sentiment

- For short articles, does the overall connotation in the article, while following the labeling process flow, have a label?

- For long articles, we section them and ask the reviewers to follow the above label process flow to define a label to the section.

# Tokenization: Sentence and word tokenization

Input to the models are sequences of sentences and words:

- Sentence Tokenization: From the text, identify the beginning and end of sentences.

- Word Tokenization: Break each sentence into sequences of words.

深圳市嘉兰图设计股份有限公司关于子公司破产申请已被法院受理的公告。本公司及董事会全体成员保证公告内容的真实、准确和完整，没有虚假记载、误导性陈述或者重大遗漏，并对其内容的真实性、准确性和完整性承担个别及连带法律责任。

[
["深圳市", "嘉兰图", "设计", "股份", "有限公司", "关于", "子公司", "破产", "申请", "已", "被", "法院", "受理", "的", "公告"],

["本", "公司", "及", "董事会", "全体成员" "保证", "公告", "内容", "的", "真实", "准确", "和", "完整", "没有", "虚假", "记载", "误导性", "陈述", "或者", "重大", "遗漏", "并", "对", "其", "内容", "的", "真实性", "准确性", "和", "完整性", "承担", "个别", "及", "连带", "法律责任"]
]

**S&P Global**

# Vector representation of words

Once we have converted our text into sequences of sentences and words, we need to turn these sequences into numerical vectors. There are two main methods for this:

1. One-hot encoding

   This representation works when vocabulary is small and so not useful, especially for the Chinese language.

```
'The mouse ran up the clock' = [
    [0, 1, 0, 0, 0, 0, 0],
    [0, 0, 1, 0, 0, 0, 0],
    [0, 0, 0, 1, 0, 0, 0],
    [0, 0, 0, 0, 1, 0, 0],
    [0, 1, 0, 0, 0, 0, 0],
    [0, 0, 0, 0, 0, 1, 0]
]
```

2. Word embeddings

   Model learns to turn word index sequences into vectors during the training process, such that each word index gets mapped to a dense vector of real values representing that word's location in semantic space.

S&P Global

# Neural Networks



Neural Networks are nothing but a series of linear transformations from inputs to outputs.

For this system, we have assessed alternative methodologies and chose a framework that has two main components:

- Convolutional Neural Networks (CNN): They learn local patterns of the input sequence. However, they have no memory. Each input shown to them is processed independently, with no state kept in between inputs.

- Long Short-Term Memory (LSTM) networks: They iterate over the input sequence, and 'remember' what they have seen so far.

Our analysis shows that combining CNN and LSTM networks gives the system the ability to learn complexities of the Chinese language.

# Neural Networks - limitations

Neural Networks are revolutionizing our day-to-day life. However, they have two main drawbacks:

1. Lack of transparency:
   - Utilized method for overcoming it:
     - Assess the effect of changing the model's input on its output. Effectively, utilize sensitivity analysis to understand how the model works and its main drivers

2. Overfitting the training set:
   - Utilized methods for overcoming it:
     - Maintaining a hold-out validation set
     - Drop outs: randomly *dropping out* (setting to zero) a number of output features of the layer during training
     - K-fold cross validation (split data into K partitions of equal size. For each partition $i$, train a model on the remaining $K-1$ partitions, and evaluate it on partition $i$)

# Benchmarking the approach

Benchmarking a complex workflow and deep learning is challenging – it is context and language sensitive.
We are focusing on the process and training set.

Approaches we leveraged

| Approach | Result |
|---|---|
| Bag of Words on English translation | ~40% accuracy |
| FaceBook FastText vectorization (native Chinese – non financial context) | ~ 30% accuracy |
| Strong performing Deep Learning model | 60%+ is considered a strong model |

- Chosen architecture, compared to alternative architectures, shows robust out-of-sample performance. Additionally, results of conducted sensitivity analysis shows model is self consistent via Local Interpretable Model-agnostic Explanations (LIME).
- Performance will be monitored through capturing users' feedback within the User Interface, and quarterly performance monitoring against a challenger architecture (Bag of Words approach).
- Re-training is deterministic in case of significant divergence in model performance against the challenger systems and based on clients' feedback, or an introduction of a new announcement category by data provider not seen during the training.

**S&P Global**

# UI: Transparency via LIME to show the Extraction Algorithm is *Self Consistent*

True class is: Negative

Prediction probabilities

| | |
|---|---|
| Negative | 0.89 |
| Neutral | 0.09 |
| Positive | 0.02 |

**Text with highlighted words**

证券代码：⊖ 证券简称：中艺股份　公告编号
南京中艺建筑设计院股份有限公司关于收到扬州市城乡建设局行政处罚决定书的公告
本公司及董事会全体成员保证公告内容的真实、准确和完整，没有虚假记 载、误导性陈述或者重大遗漏，并对其内容的真实性、准确性和完整性承担 个别及连带法律责任
南京中艺建筑设计院股份有限公司（以下简称"公司"）近日收到扬州市城乡建设局开具的《行政处罚决定书》（扬 建罚
扬州市城乡建设局城市建设监察支队调查发现公司未 按工程建设强制性标准（条文）对承揽的扬子江世贸商业综合体工程进行设计，该行为违反了《建设工程质量管理条例》第十九条第一款
扬州市城乡建设局依据《建设工程质量管理条例》第六 十三条第一款第（四）项以及第七十三条对当事人进行处罚
鉴于该工程尚未开工，公司积极配合行政机关查处违法行为并整改完成，主动消除违法行为危害后果，根据《中华人民
公告编号：⊖ 共和国行政处罚法》第二十七条第一款的规定，给予公司减轻行政处罚
综上，扬州市城乡建设局作出对中艺公司处⊖ 元罚款、对直接负责的主管人员王刚处⊖ 元罚款、对其他直接责任人员祁岭处⊖ 元罚款的行政处罚
公司将根据《行政处罚决定书》的规定，在缴款期限内 足额缴纳人民币伍万元罚款，该笔罚款占营业收入比重较小，情节较轻
此外，该工程尚未开工，我公司积极配合行政机关查处违法行为并整改完成，主动消除了危害后果，本次事件不会对公司经营产生重大影响
公司管理层高度重视本次事件：由总经理召开紧急会议，进行全公司通告，对相关责任人进行严肃地批评和处罚，督促全公司认真吸取本次行政处罚的教训，积极整改；责令公司员工对公司承接的所有项目进行全面排查，避免此类问题再次发生；总结本次事项的教训，制定预防措施，特别是进一步加强对此类外地工程设计项目成果质量的严格把控，定期对公司员工进行工程建设强制性标准及相关法律法规的培训
扬州市城乡建设局行政处罚决定书》（扬建罚

Positive Connotation – Green
Negative Connotation – Blue

**Google Translate:**

Stock code: 证券 Securities abbreviation: Zhongyi Shares Announcement No.

Announcement of Nanjing Zhongyi Architectural Design Institute Co., Ltd. on Receiving the Administrative Penalty Decision of Yangzhou Urban and Rural Construction Bureau

[ ...]

Nanjing Zhongyi Architectural Design Institute Co., Ltd. (hereinafter referred to as "the company") recently received the "Administrative Punishment Decision" issued by Yangzhou Urban and Rural Construction Bureau.

[...]

Yangzhou Urban and Rural Construction Bureau will impose penalties on the parties in accordance with Article 63, paragraph 1 (4), and Article 73 of the Regulations on Quality Management of Construction Projects.

[...]

Announcement No.: 规定 The provisions of the first paragraph of Article 27 of the Administrative Punishment Law of the Republic give the company less administrative penalties

In summary, the Yangzhou Urban and Rural Construction Bureau made a fine for the Chinese company, a fine for the directly responsible person, Wang Gang, and a fine for other directly responsible personnel.

According to the "Administrative Punishment Decision", the company will pay a fine of RMB 10,000 in full within the payment period. The fine accounts for a small proportion of the operating income, and the plot is lighter.

In addition, the project has not yet started, and our company actively cooperates with the administrative organs to investigate and deal with illegal acts and complete the rectification, and actively eliminates the harmful consequences. This incident will not have a major impact on the company's operations.

The management of the company attaches great importance to this incident: the general meeting convened an emergency meeting, issued a company-wide notice, seriously criticized and punished the relevant responsible persons, urged the whole company to seriously absorb the lessons of this administrative punishment, and actively rectified; ordered the employees of the company Conduct a comprehensive investigation of all the projects undertaken by the company to avoid such problems;, and develop preventive measures, especially to further strengthen summarize the lessons learned from this matter the strict control of the quality of such field engineering design projects, and regularly to the employees of the company. Training in mandatory construction standards and related laws and regulations Yangzhou Urban and Rural Construction Bureau Administrative Punishment Decision

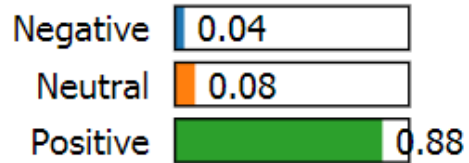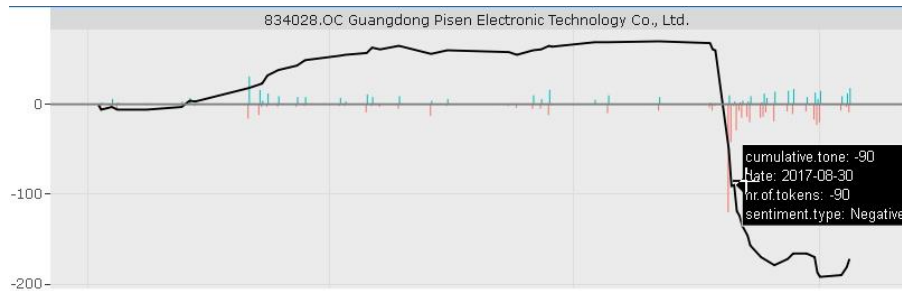*Note: the translation is provided for context only. The process deals with the native language announcements.*

# We are identifying relevant pieces of text and its embedded opinion
## Context Sensitive Extraction, Attribution, and Actionable Notifications via Sentiment Analysis

True class is: Positive

Prediction probabilities

| | |
|---|---|
| Negative | 0.04 |
| Neutral | 0.08 |
| Positive | 0.88 |

**Text with highlighted words**

证券代码：Θ 证券简称：力尊信通 公告编号
主办券商：中航证券 北京力尊信通科技股份有限公司关于获得高新技术企业证书的公告
本公司及董事会全体成员保证公告内容的真实、准确和完整，没有虚假记载、误导性陈述或者重大遗漏，并对其内容的真实性、准确性和完整性承担 个别及连带法律责任
北京力尊信通科技股份有限公司（以下简称"公司"）于Θ年Θ月Θ日收到由北京市科学技术委员会、北京市财政局、北京市国家税务局、北京市地方税务局联合颁发的《高新技术企业证书》（证书编号：αΘ），发证时间：Θ 年Θ 月Θ 日，有效期三年
根据相关规定，公司自通过高新技术企业认定并向主管税务机关办理完减免税手续后三年内，可享受国家关于高新技术企业的相关优惠政策，即按Θ税率缴纳企业所得税

Positive Connotation – Green
Negative Connotation – Blue

**Google Translate:**

Stock code: 证券 Securities abbreviation: Lizun Xintong Bulletin number

Sponsored Broker: AVIC Securities Beijing Lizun Xintong Technology Co., Ltd. Announcement on Obtaining High-tech Enterprise Certificate

The company and all members of the board of directors guarantee the truthfulness, accuracy and completeness of the contents of the announcement, and there are no false records, misleading statements or major omissions, and bear individual and joint legal responsibility for the truthfulness, accuracy and completeness of the contents.

Beijing Lizun Xintong Technology Co., Ltd. (hereinafter referred to as "the company") was jointly awarded by the Beijing Municipal Science and Technology Commission, the Beijing Municipal Finance Bureau, the Beijing Municipal State Taxation Bureau and the Beijing Municipal Local Taxation Bureau on the following day. "High-tech Enterprise Certificate" (Certificate No.: αΘ), issuance time: Θ Year Θ Month Day, valid for three years

According to the relevant regulations, the company can enjoy the relevant preferential policies of the state on high-tech enterprises within three years after passing the high-tech enterprise certification and handling the tax reduction and exemption procedures with the competent tax authorities, that is, paying the enterprise income tax at the tax rate.

*Note:* *the translation is provided for context only. The process deals with the native language announcements.*

S&P Global

# UI: We have implemented a simple demo app that aggregates company-level sentiment and attributes the score to key words

## Screenshots of actual indicative output

### Time series of company-level aggregated sentiment

834028.OC Guangdong Pisen Electronic Technology Co., Ltd.

cumulative.tone: -90
Hate: 2017-08-30
nr.of.tokens: -90
sentiment.type: Negative

### Word cloud: key drivers of positive/negative sentiment

### Highlighting of key sentiment driver words in original text

Prediction probabilities

| | |
|---|---|
| Negative | 0.89 |
| Neutral | 0.09 |
| Positive | 0.02 |

Positive Connotation – Green
Neutral Connotation – Orange
Negative Connotation – Blue

**Text with highlighted words**

证券代码：⊙ 证券简称：中艺股份　公告编号
南京中艺建筑设计院股份有限公司关于收到扬州市城乡建设局行政处罚决定书的公告
本公司及董事会全体成员保证公告内容的真实、准确和完整，没有虚假记 载、误导
性陈述或者重大遗漏，并对其内容的真实性、准确性和完整性承担 个别及连带法律
责任
南京中艺建筑设计院股份有限公司（以下简称"公司"）近日收到扬州市城乡建设局开
具的《行政处罚决定书》（扬 建罚
扬州市城乡建设局城市建设监察支队调查发现公司未 按工程建设强制性标准（条文）
对承揽的扬子江世贸商业综合体工程进行设计，该行为违反了《建设工程质量管理条
例》第十九条第一款
扬州市城乡建设局依据《建设工程质量管理条例》第六 十三条第一款第（四）项以及
第七十三条对当事人进行处罚
鉴于该工程尚未开工，公司积极配合行政机关查处违法行为并整改完成，主动消除违
法行为危害后果，根据《中华人民
公告编号：⊙ 共和国行政处罚法》第二十七条第一款的规定，给予公司减轻行政处罚
综上，扬州市城乡建设局作出对中艺公司处⊙ 元罚款、对直接负责的主管人员王刚处
⊙ 元罚款、对其他直接责任人员祁岭处⊙ 元罚款的行政处罚
公司将根据《行政处罚决定书》的规定，在缴款期限内 足额缴纳人民币伍万元罚款，
该笔罚款占营业收入比重较小，情节较轻
此外，该工程尚未开工，我公司积极配合行政机关查处违法行为并整改完成，主动消
除了危害后果，本次事件不会对公司经营产生重大影响
公司管理层高度重视本次事件：由总经理召开紧急会议，进行全公司通告，对相关责
任人进行严肃地批评和处罚，督促全公司认真吸取本次行政处罚的教训，积极整改；
责令公司员工对公司承接的所有项目进行全面排查，避免此类问题再次发生；总结本
次事项的教训，制定预防措施，特别是进一步加强对此类外地工程设计项目成果质量
的严格把控，定期对公司员工进行工程建设强制性标准及相关法律法规的培训
扬州市城乡建设局行政处罚决定书》（扬建罚

S&P Global

30

# Further details on overcoming model transparency

- LIME: An algorithm that can explain the predictions of any classifier or regressor in a faithful way, by approximating it locally with an interpretable model.

- Allows us to analyze the effect of removing a token within a corpus to see how it will impact the overall sentiment classification.

- We can rank these perturbations by article and across the entire corpus via Gini Importance to identify the most significant tokens (words) that have an impact on a certain polarity (positive, negative)

- The paper below explains the approach in detail:

- "Why Should I Trust You?" Explaining the Predictions of Any Classifier
  - https://arxiv.org/abs/1602.04938
  - https://arxiv.org/pdf/1602.04938.pdf
  - University of Washington:
    - Marco Tulio Ribeiro
    - Sameer Singh
    - Carlos Guestrin

**S&P Global**

# Conclusion

- Hope we have demonstrated how we have converted unstructured data to reliable information using a Deep Learning Pipeline.

▪ We are actively using cutting-edge technologies and techniques to improve our clients' day-to-day workflow by building on top of S&P Global Market Intelligence's models and data to solve complex business problems for our clients.

▪ Thank You!

**Moody Hadi**

Group Manager – Innovation & Product Research

Risk Services

S&P Global Market Intelligence

**S&P Global**

**S&P Global**