生成AIで迷惑SMS対策やってみた

KDDI株式会社 コア技術統括本部 オペレーション本部 プラットフォームエンジニアリング部 鹿野良太

2025/10/24



- 自己紹介



- 名前
 - ・鹿野良太
 - ■出身
 - ・東京都

- 経歴
 - 2024年度~
 - KDDI株式会社入社
 - SMS設備担当
 - →SMS設備の開発・保守運用
 - →生成AIを使用した迷惑SMS対策

■趣味

- ・美味しいものを食べる
- →高知で鰹のたたきを食べたい

今回の発表テーマ

- 生成AIを活用するに至った背景

<従来の課題>

・迷惑SMSの疑いのあるお客様からの申告データをもとに、SMS設備に適用する迷惑SMSブロック用のファイルを人力で作成していた

稼働 **国視で**正規SMSと迷惑SMSに分類し ファイル作成 専任で約1.5営業日

鮮度

分析〜設定までに<u>一週間</u> 迷惑SMSのトレンドを追えない

精度

<u>一部の</u>申告データのみで分析作業 を行うためブロック率が上がらない

く仮説>生成AIは申告SMSから正規のSMSと迷惑SMSを判別できるのではないか。もし可能であれば、課題の解決になる。

<検証>迷惑SMSの判定を指示するプロンプトを作成し、実際の迷惑SMSを入力

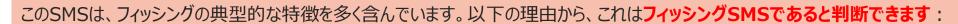
入力本文

゚ロンプト(生成スエスの指言

プロンプト(生成AIへの指示文)

迷惑SMSであるか判 定してください

生成AIの回答



- 1.緊急性を強調する言葉遣い:「緊急通知」という表現を使用しています。
- 2.銀行を装っている:三菱UFJ銀行の名前を使用し、正当性を装っています。
- 3.不自然なURL:短縮URLが使用されており、実際のリンク先が隠されています。これは典型的なフィッシング手法です。

<結果>生成AIは正規のSMSか迷惑SMSかおおむね判別できる。

しかし、商用適用にあたって、正規SMSを迷惑SMSと誤判定する(FP)をゼロにすることが必要である。



FP=0への取り組み ~SMS本文判定~

- FPゼロ対策実施内容
- ・無理に白か黒かの2値分類をさせない
- →わからないものはわからない(判定不可)と答えてもらう
- 各社HPの注意喚起をプロンプトに反映
 - ・業界ごとに判定プロンプトを作成し、各社HP上で「この本 文(URL)に注意」と、注意喚起されている内容をプロンプ トへ反映
 - →しかし、各社の公式HP内容を継続的に追従する必要がある

【迷惑メール・SMS】最近多い迷惑メール・詐欺メールの事例と特徴を 知りたい KDDIやauを装う迷惑SMS 【迷惑メール・SMS】KDDIやauを装う迷惑SMSの特徴を知りたい A KDDIやau、通信会社を装う迷惑SMSの事例です。 迷惑SMSに記載されたURLを開いたことで、不正アプリのダウンロード、偽のMy auへのログイン、KDDIで受け付けていな い電子マネー(ギフトカード)での料金支払いなどが求められる事例が報告されています。 KDDIやauが下記のようなSMSをお客さまに配信することは一切ありません。 身に覚えのないものや不審に感じた場合は、記載のURLなどは絶対に開かないようにご注意ください。

KDDI系に特化させたプロンプトの例

<task> <role>あなたは大手通信会社セキュリティ 各社HPのURLとHTML文 <instructions> をDBに保存し更新がある SMSの文章を見て、<KDDI、au、UOモ たびに**赤枠プロンプト** イル、POVOなどのサービス]</KDDI、au、し シングSMSであるか判定してください。 を自動再作成 以下のフィッシングSMSの特徴を参考に </instructions> <features> <feature># KDDI、au、UQモバイル、POVOなどのサービスにおける正規のSMSの特徴: - 公式アプリやあらかじめ登録したブックマークからアクセスしてログインする · 2段階認証時は内容をしっかり確認する - アプリのインストールは公式サイトから行う - 送信元が正規のものである(例: auto@connect.auone.jp) - 公式サイトのURLが正しく記載されている

- 正式な社名・サービス名が使用されている

- 日本語として自然な文章で書かれている

#KDDI、au、UQモバイル、POVOなどのサービスにおける迷惑SMSの特徴:

- auやKDDIを装ったメッセージとURLが記載されている

"訴訟最終通知"</feature>

<judgmentCriteria>

判定結果は分析しやすい形式で指定 <criteria>#判定結果:</cr

<options>

<option>フィッシングの場合: black</option>

<option>フィッシングでない場合: white</option>

<option>判定不可: gray/option>

</options>

</judgmentCriteria>

- FP=0への取り組み ~HTML判定~

■ 判定材料がSMS本文だけではFP=0が実現せず

HTMLの中に不審なリン クが含まれている

実例

URL先は正規の ホテルサイト

お荷物について

https://xxxxxxxx

xxxxxx.xx/xxxxxxxx/xxxxxx

xxxx/xxxx

HPに連絡先や問い 合わせ情報が載って いない

HTMLの構文が不自然



- ・URL先のHTML文も生成AIの判定材料、 特にSMS本文では読み取れない情報を判断材料 に白黒判定を行う
 - →本文とHTMLの2重で判定することで誤判定を防ぐ

HTML文 判定プロンプトの例

<role>あなたは大手通信会社セキュリティ部門のプロフェッショナルです。</role> <instructions>

URLのHTMLの中身を見て、迷惑SMSメッセージであるか判定してください。 以下の迷惑SMSメッセージのHTMLの特徴を参考にしてください。

</instructions>

<features>

- <feature>不審なリンクやURLが含まれている</feature>
- <feature>緊急性を強調する表現が使われている</feature>
- <feature>送信者の正当性が疑わしい(公式なものではない)</feature>
- <feature>個人情報の提供を要求する内容</feature>
- <feature>文法やスペルの誤りが多い</feature>
- <feature>一般的な挨拶や名前の欠如</feature>
- <feature>偽のロゴやデザインが使用されている</feature>
- <feature>連絡先情報や問い合わせ先の情報が明確に記載されていない</feature>
- <feature>html構文が不自然</feature>
- </features>

<judgmentCriteria>

<criteria>#判定結果:</criteria>

<options

<option>迷惑SMSの場合: black

<option>迷惑SMSでない場合: white</option>

<option>判定不可: gray/option>

</options>

</judgmentCriteria>

<output/nstructions>

<instruction>#出力方法:</instruction>

<format><ans>判定結果</ans>(出力では<ans>判定結果</ans>のみ出力してください)

</format>

</output/instructions>

<example>

<description>例: </description>

<sampleInput>入力: HTMLの中身</sampleInput>

<sampleOutput>出力:<ans>black</ans></sampleOutput>

</example>

prompt>

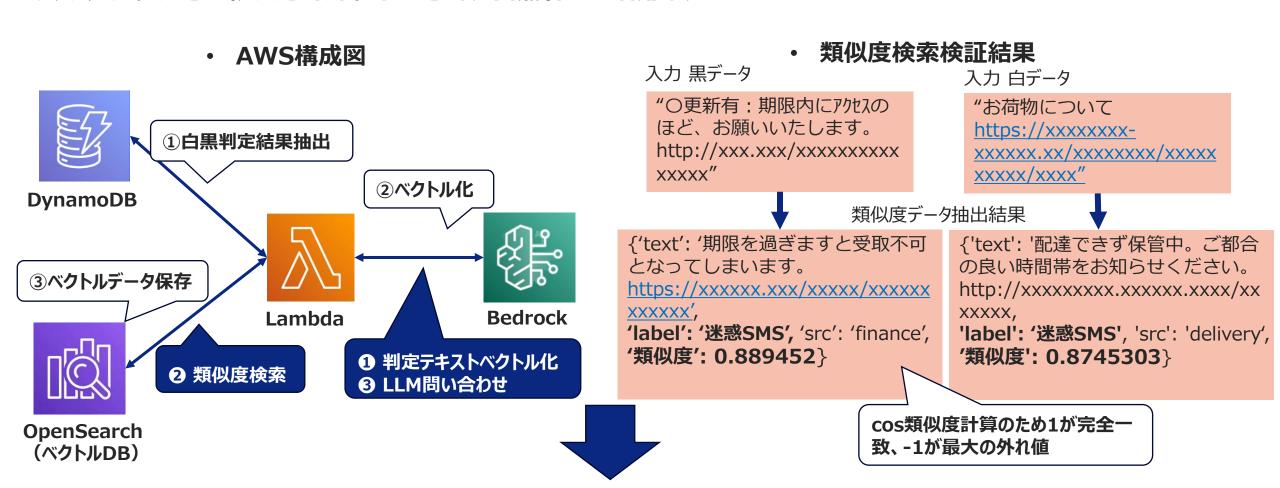
以下のURLのHTMLの中身を判定してください:

© KDDI CORPORATION

- FP=0への取り組み ~RAGの導入検討~

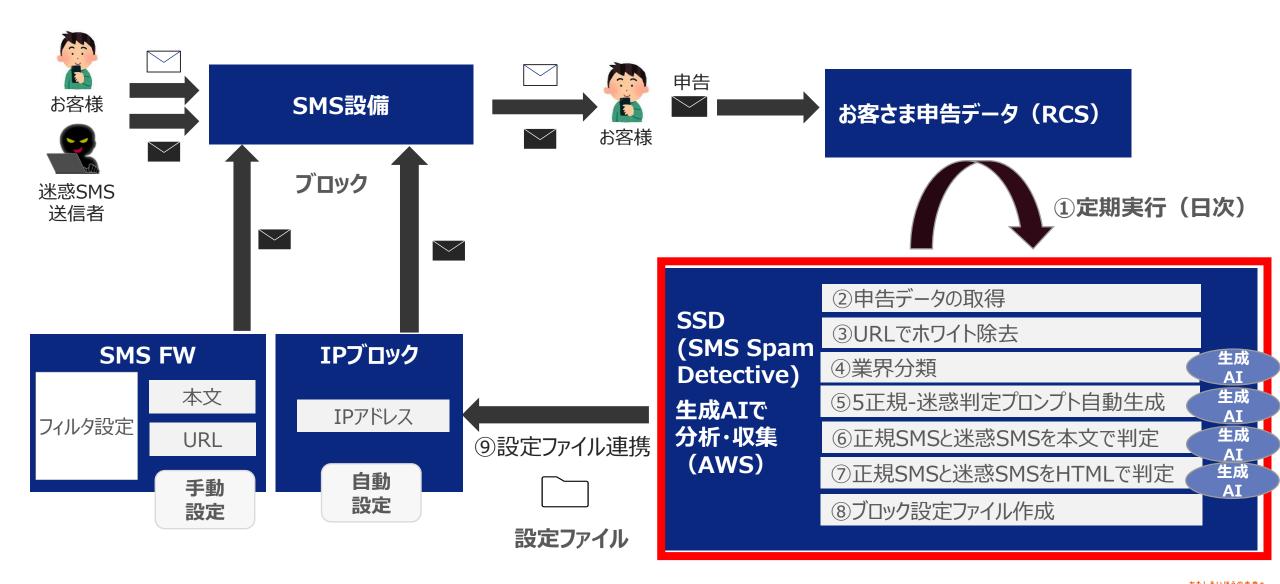
■ 精度を更に高めるために

過去に白黒判定した本文をベクトル化(数値化)しDBに保存し、今後判定するSMSも同様にベクトル化して、 過去のナレッジと類似度を計算することで、白黒判定に活用したい





- 構築したシステム(SSD)の概要



- 運用実績と課題

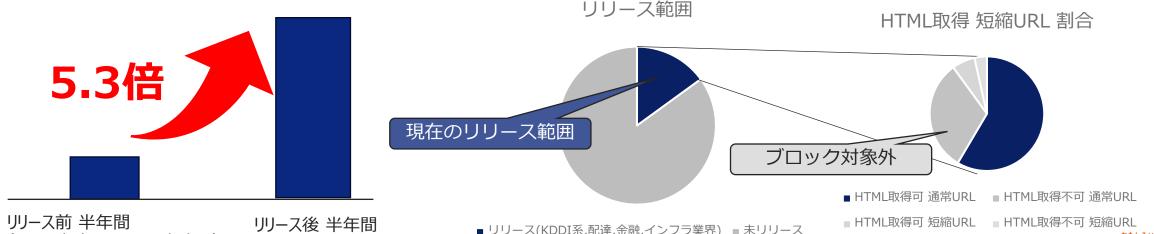
■ 運用実績(2025/3/25~2025/9/24)

- ・リリース範囲:業界分類でKDDI、配達、金融、インフラ業界に分類された全体15%の出力IPアドレスだけSMS設備と連携
- ・半年間 FP数: **0件**
- ・ IPブロック数 リリース前後の半年間の比較で**5.3倍**に

■ 現状の課題

- ・HTML判定をする際、リンク先サイトが削除されており、HTMLを取得出来ずに判定不可が35%。
- ・ 短縮URLから返されるIPアドレスをブロックすると、正規のSMSもブロックしてしまう可能性があるため、

短縮URLの約10%がブロックできない。



(2024/9/25~2025/3/24)

 $(2025/3/25\sim2025/9/24)$

■ リリース(KDDI系,配達,金融,インフラ業界) ■ 未リリース



社外秘B

- 今後の方針と活動を通じての学び

■ 今後の方針

- ・本文判定にRAGを取り入れる
- ・業界分類の残り85%リリースする

- ・IPブロックだけではなくURLブロックを取れ入れる
- ・迷惑メール対策にも適用する

■ 活動を通じて

- ・ 生成AIをシステム内で利用する際の難しさを感じた。
- →生成AI出力結果の確認は最終的に人が行うことになり、リリース後もどの程度まで確認を続けるべきか、 不明確になる
- ・生成AIはブラックボックスであり、出力結果はモデルの性能の良さに依存してしまうと考えていたが、
 - プロンプトやRAGなどの工夫により出力結果を変えられる





「つなぐチカラ」を進化させ、 誰もが思いを実現できる社会をつくる。

- KDDI VISION 2030

