# *Elasticsearchを活用して、統計データから新たなインサイトを生成する取り組み*

Jesus Jackson, Chief Data Scientist

Daewoo Chong, Lead Data Scientist

*ブーズ・アレン・ハミルトン*

2016年10月13日

United States™
Census
Bureau

# アジェンダ

- はじめに
- 検索プロトタイプ
  - プロトタイプとは？
  - 統計データについて
- インデキシングストラテジー
  - 統計データのインデキシングにおける課題とは？
  - 候補となるストラテジー現行のストラテジー
- プロトタイプの機能
  - プロトタイプの主要機能
  - Elasticsearchの主要機能
- デモ
- ご質問をどうぞ！

# はじめに



### Jesus Jackson

- ブーズ・アレン・ハミルトンに7年+勤務
- 金融サービス業界向け連邦データサイエンスサービス部門責任者
- 自身について一言：先日スカイダイビングしました



### Daewoo Chong

- ブーズ・アレン・ハミルトン1.5年+勤務
- 機械学習マニア
- 自身について一言：オフの日はニューラルネットを作成してます

## 他のチームメンバー：

Julia Stevens
Raj Cheekatamarla
Ram Anusuri

## クライアントチームメンバー：

Richie Wang
David MacCormack
Zachary Whitman

elastic

2

# 検索プロトタイプ

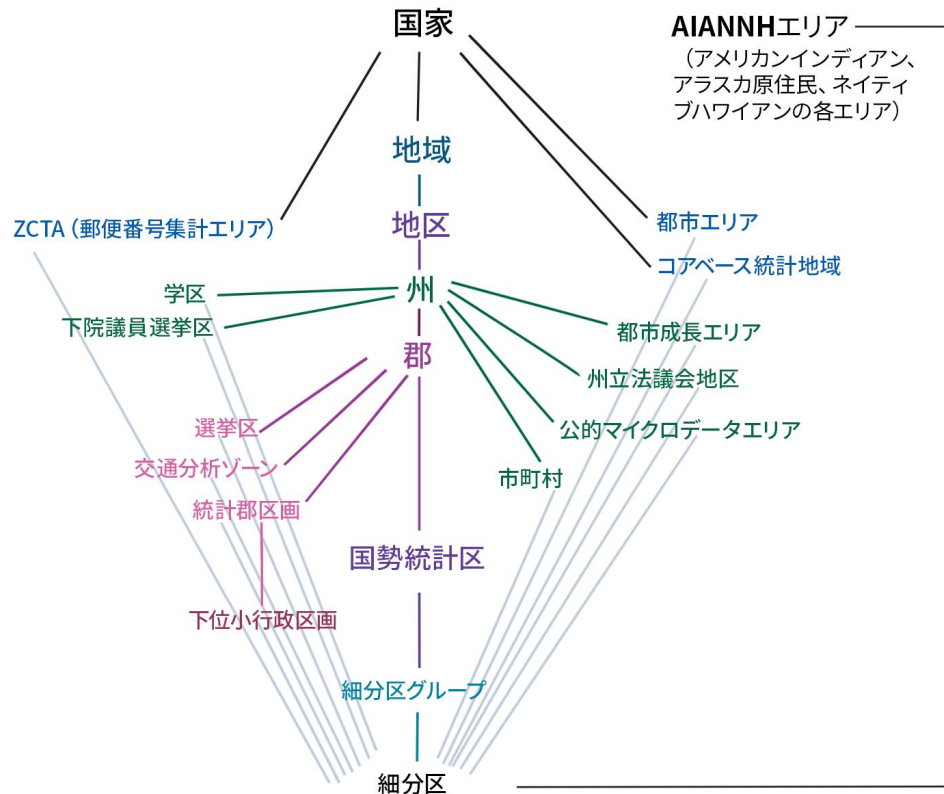# プロトタイプとは？

ユーザーがシンプルに、いつでも妥当なデータにアクセスできる機能を実現する

# 検索プロトタイプ

統計データについて
- 統計データは非常に情報量が多く、複雑。たとえば：
  - 調査/プログラムの種類：米国コミュニティ調査、国勢調査、事業主調査など
  - 地理的区分：地域、州、郡、学区、細分区など
  - トピック：年齢、コンピューターとインターネットの使用状況、貧困度、人種、性別など
  - 業界：建設、金融・保険、製造、卸売など

| | |
|---|---|
| Journey To Work (543) > | 00 - Total for all sectors |
| Labor Force Status (1505) | 11 - Agriculture, forestry, fishing and ... > |
| Language Spoken At Home (313) > | 21 - Mining, quarrying, and oil and gas ... > |
| Legal Form of Organization (4) | 22 - Utilities > |
| Living Quarters (2) | 23 - Construction > |
| Migration (296) > | 31-33 - Manufacturing > |
| Nativity And Foreign Born Population (8) > | 42 - Wholesale trade > |
| Population Total (57) | 44-45 - Retail trade > |

# 検索プロトタイプ：*地理階層*

# インデキシングストラテジー

# インデキシングストラテジー

統計データのインデキシングにおける課題とは？
- ブール検索モデル
- トピック、地理、業界、データセット、調査年を横断するフィルタリング
- デッドエンドなし
- 応答時間100ミリ秒以下

# インデキシングストラテジー

候補となるストラテジー現行のストラテジー

### 非正規化

クエリしやすいが
サイズが大きくコスト高

### ペアレント/チャイルド

小サイズでコストは
高くないが応答時間
が遅い

### 準正規化

小サイズで経済的、応
答時間も早いが、セル
値を持たない

elastic

プロトタイプの機能

# プロトタイプの機能

プロトタイプの主要機能
- デッドエンドなし
- 応答時間100ミリ秒以下
- ログのパース/インデックス/分析にElastic Stackを使用
- 認証とアクセスベースのロール制御に Shieldを使用

# プロトタイプの機能

プロトタイプにメリットをもたらした主要な Elasticsearchの機能
- クエリキャッシングで応答時間が最大 25%向上
- マッピング API
- 関連性モデル（例：TF/IDF、BM25）
- クエリ API
- 複数言語のサポート（例：JavaScript、Python、Java）

デモ

ご質問をどうぞ！

www.elastic.co