# USING DATA TO FIND THE OPTIMAL MIX OF RETAIL LOCATIONS AND RESOURCES

**north**highland®

WORLDWIDE CONSULTING

Data and Analytics

# INTRODUCTION

Education

- BS CS, Georgia Tech 2009 – Theory and Machine Learning
- MS CS, Georgia Tech 2011 – Heavy Tail Network Analysis

Work

- Institute of Nuclear Power Operations (2010-16)
  - Build, deploy, maintain a model that predicted nuclear power station performance along 13 key functional areas
- North Highland (2016-)
  - ETL, BI, Advanced Analytics for Fortune 100 retailer

**north**highland® | Data and Analytics
WORLDWIDE CONSULTING

1. Data & Analytics outside academia

2. Case Study: Reassigning territories for district managers

3. Q&A

northhighland®
WORLDWIDE CONSULTING

Data and
Analytics

# WORKING WITH CLIENTS

- Problems are never stated formally

- "Interesting" problems can be few and far between

- But they can build your personal brand

Proprietary & Confidential

# REMAPPING TERRITORIES – PROBLEM DESCRIPTION

- Minimizing travel time for regional managers can reduce incurred travel costs and boost morale
- Aligning districts to strategic goals can help ensure a variety of goals:
  - A level playing field where top talent can be evaluated evenly
  - Specialized focus for individual district owners
  - No one regional leader becomes overburdened compared to the others

northhighland® | Data and Analytics
WORLDWIDE CONSULTING

# AVAILABLE DATA

- Store Metadata – geocoding, age, size, store annual sales category, etc.
- Sales Data – department, class, subclass, SKU grain data anywhere from monthly roll-ups to individual transactions
- Inventory Data
- Online Transactions
- Current Territory

# (ABBREVIATED) TOOLBOX OF TECHNIQUES

**Technique 1: k-means**

- Unsupervised Learning
- Identifies a number of means around the map and builds clusters with equal variance inside them
- Very much a black box-hard to specify, and requires a lot of tuning
- **Use if:** You want to explore your data, equal size isn't as important
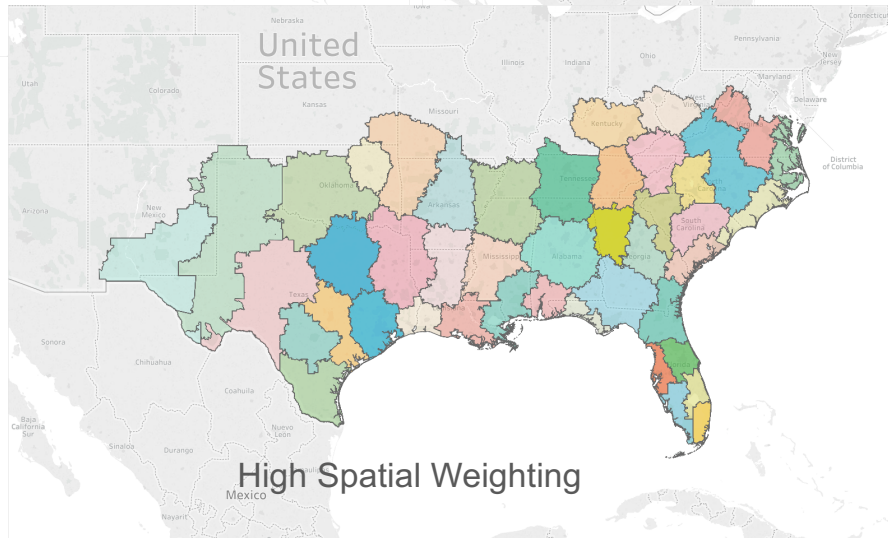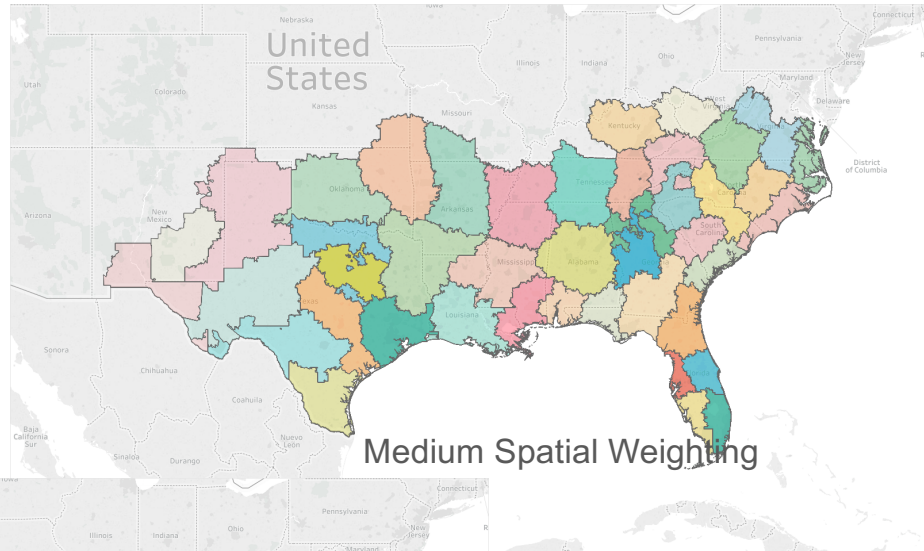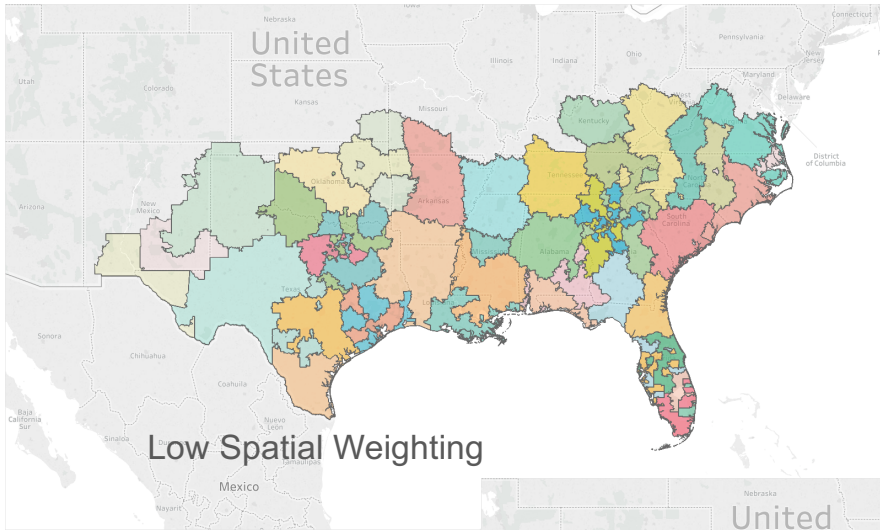
**Technique 2: Integer programming**

- Can specify exactly what you want, but rules are rigid
- Computationally impossible for large datasets—constraints have to be relaxed
- **Use if:** You have little data

**Technique 3: Network construction**

- Randomized (or can be non-random) algorithm to build out a network 'greedily'
- Easy to specify and tune parameters as you go
- **Use if:** Iteration is OK, exact solutions aren't required

northhighland | Data and Analytics
WORLDWIDE CONSULTING

# PULLING IT TOGETHER

- SQL

- Python
  - ○ Pandas
    - High performance data management/manipulation, SQL-like interface
  - ○ Numpy
    - N-dimensional arrays, math libraries
  - ○ Scikit-learn
    - Huge number of supervised and unsupervised ML algorithms prewritten
  - ○ Networkx
    - Network/Graph analysis library

- Brute force

northhighland® | Data and Analytics
WORLDWIDE CONSULTING

Low Spatial Weighting


Medium Spatial Weighting


High Spatial Weighting

northhighland®
WORLDWIDE CONSULTING

Data and
Analytics

# NETWORKX

**https://networkx.github.io/**

- Graph data structure with huge library of built-ins
  - Graph Operations
    - Edge/Node maintenance, weighting, node attributes, etc.
  - Graph Algorithms
    - Connectivity, Neighborhoods, k-core, max-flow, matching, bipartite, approximation algorithms, and on and on…
  - Linear algebra library that takes graph objects
    - Eigenvalue spectrums, laplacians, PageRank
  - Generators
    - Random graph generators (e.g. random normal, Erdős–Rényi, power law)
    - Canonical graphs (Karate club, Florentine families graph)
  - Visualization Tools

northhighland® | Data and Analytics
WORLDWIDE CONSULTING

# GREEDY ALGORITHM OVERVIEW

- Load data
- Using networkx, build an approximately-planar graph based on district mean locations
  - Find the norm of the district centers, pick n-closest
- Set parameters for "optimizer"
- Loop:
  - Pick manager with lowest score, assign them a random district that's a neighbor as long as constraints are met
    - If that manager has no districts, pick a random district to add.
  - Simulated annealing—jostle where districts are in an attempt to avoid local minima, cooling over time
- Once all districts are assigned, score districts and reshuffle them to minimize variance

northhighland®    Data and Analytics
WORLDWIDE CONSULTING

# LOAD DATA

```python
65 #Import our data into dataframes
66 storeList = pd.read_csv(appDir + 'RMM Territories/rmm_map.txt', sep='|', na_values='?')
67
68
69 storeClass = pd.read_csv(appDir + 'RMM Territories/Str_Class_Sales.txt', sep='\t', na_values='?')
70 #storeList = pd.concat([storeList, storeClass['District Nbr']], join='inner', axis=0)
71 storeList = storeList.set_index(['STR_NBR'])
72 storeClass = storeClass.set_index(['Str Nbr'])
73 allStore = pd.concat([storeList, storeClass], join='inner', axis=1)
```

# BUILD GRAPH

```
119 distrLocCenters = pd.concat([storeIndex['District'], distrList], join='inner', axis=1).groupby(['District']).mean()
120
121 distr_distances = np.zeros((distrLocCenters.shape[0],distrLocCenters.shape[0]))
122 for i in range(distrLocCenters.shape[0]):
123         for j in range(distrLocCenters.shape[0]):
124             if i!=j:
125                 distr_distances[i,j] = np.linalg.norm(distrLocCenters.values[i,:]-distrLocCenters.values[j,:])
```

```
140 rmmAddr = pd.read_csv(appDir + 'RMM Territories/RMM_Cities_Coord_V2.csv')
141
142 rmmAddr = rmmAddr.loc[rmmAddr['DIV_NBR'] == divAssign]
143 n_Distr = rmmAddr.shape[0]
144
145 m_distances = np.zeros((distrLocCenters.shape[0], rmmAddr.shape[0]))
146
147 for i in range(rmmAddr.shape[0]):
148     m_distances[:,i] += (np.linalg.norm(distrLocCenters[['LAT_NBR', 'LNG_NBR']].values - rmmAddr[['LAT', 'LNG']].values[i], axis=-1))
149
150
151
152 #Keep only n best neighbors
153 closestN = 4
154
155 for i in range(distrLocCenters.shape[0]):
156     a = distr_distances[i].argsort()[closestN+1:]
157     distr_distances[i, a] = 0
158
159 G2 = nx.from_numpy_matrix(distr_distances)
```

northhighland® | Data and
WORLDWIDE CONSULTING | Analytics

# PARAMETERS AND CONTROLS

```
161 nx.set_node_attributes(G2, 'RMM', -1)
162 storeDollar = pd.read_csv(appDir + 'MerchFinance/RegionalRealignment/strDlr.csv', sep=',', na_values='?')
163 storeDollar = storeDollar.set_index(['STR_NBR'])
164 storeIndex = pd.concat([storeIndex, storeDollar], join='inner', axis=1)
165 #Get District Dollar Values
166 distrDollar = storeIndex[['SLS_DLR','District']].groupby(['District']).mean()
167 rmmStoreCount = storeIndex[['SLS_DLR','District']].groupby(['District']).count()
168
169
170
171 #initialize the variables.
172 nx.set_node_attributes(G2, 'sls', None)
173 nx.set_node_attributes(G2, 'distNbr', 0)
174 nx.set_node_attributes(G2, 'move', 0)
175 nx.set_node_attributes(G2, 'strcnt', 0)
176
177 for i in range(len(G2.nodes())):
178     G2.node[i]['sls'] = np.round(distrDollar.values[i][0])
179     G2.node[i]['strcnt'] = rmmStoreCount.values[i][0]
180     G2.node[i]['distNbr'] = rmmStoreCount.index.values[i]
```
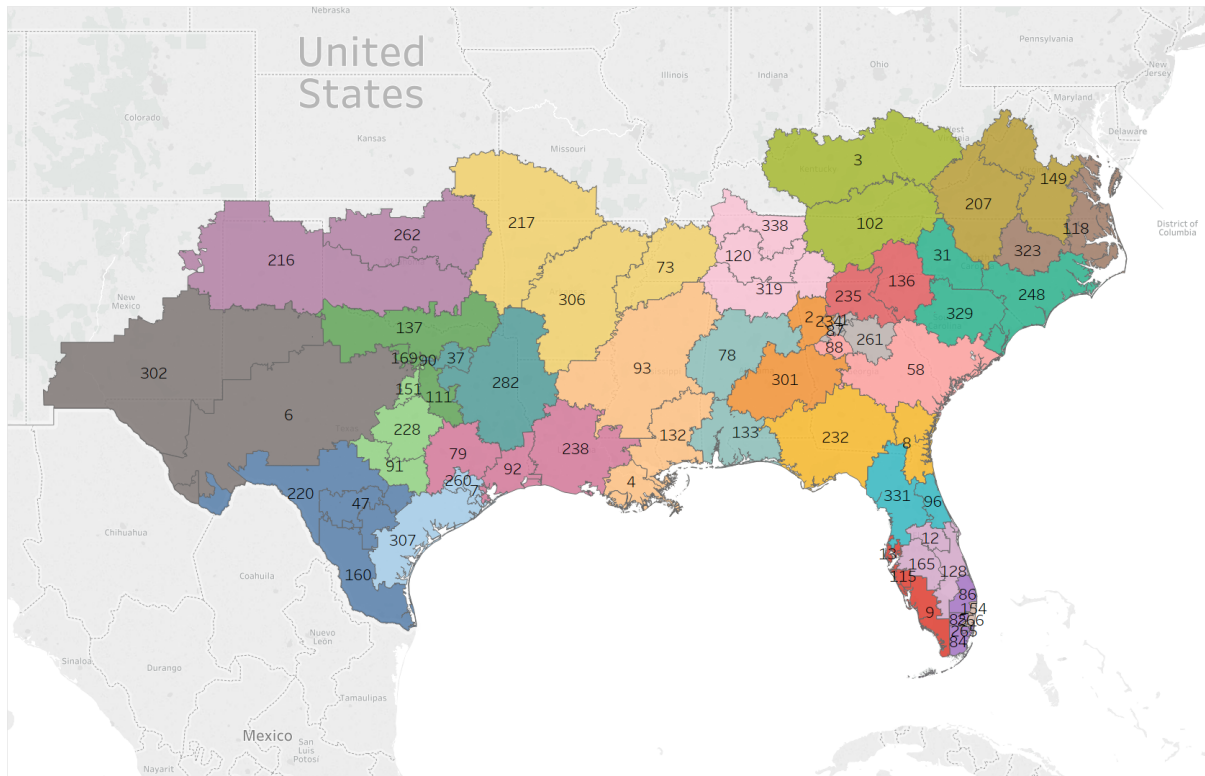
# ITERATE AND BE GREEDY

```
221    in_itns = 0
222    while ((rmmDistrInfo[rmm, 0] >= maxDistr or rmmDistrInfo[rmm,1]>maxDlr
223            or rmmDistrInfo[rmm,2]>maxStr)
224            and in_itns < 100):
225        in_itns += 1
226        possible = np.where(rmmDistrInfo[:, 2] <= rmmDistrInfo[:, 2].min()+5)[0]
227        rmm = random.choice(possible)
```

```
246    else:
247        di = distrs[np.random.randint(len(distrs))]
248        #choose a random neighbor of that rmm
249        ngh = list(G2[di])
250
251        minMove = 999
252        for i in ngh:
253            if G2.node[i]['move'] < minMove:
254                minMove = G2.node[i]['move']
255
256        minDist = 999
257        for i in ngh:
258            if G2.node[i]['move'] == minMove:
259                if distrAssignments[di] != distrAssignments[i]:
260                    if G2[di][i]['weight'] < minDist:
261                        minDist = G2[di][i]['weight']
262        mv = -1
263        for i in ngh:
264            if G2[di][i]['weight'] == minDist:
265                mv = i
266
267        dold = distrAssignments[mv]
```

- Pick a random manager from the ones that have approximately the lowest score
- Get a list of possible districts they could have, and randomly pick one of those
- Verify all the constraints (lots of IFs) are met
- Perform some simulated annealing along the way—some random chance to jostle districts from one manager to another adjacent manager occasionally to avoid local minima
- If all districts are assigned, still grab a local district if it improves your score more than it decreases your neighbor's score

northhighland® | Data and Analytics
WORLDWIDE CONSULTING

# RESULTS

Division 1 Map with "Optimal" RMM Hometown



| RMM | Distrs | Strs | RMM Hometown |
|-----|--------|------|--------------|
| 0 | 3 | 21 | 0 - San Antonio, TX |
| 1 | 3 | 20 | 1 - Houston, TX |
| 2 | 3 | 26 | 2 - Atlanta, GA |
| 3 | 3 | 25 | 3 - New Orleans, LA |
| 4 | 3 | 26 | 4 - Dallas, TX |
| 5 | 3 | 23 | 5 - Dallas, TX |
| 6 | 2 | 23 | 6 - Roanoke, VA |
| 7 | 3 | 27 | 7 - Memphis, TN |
| 8 | 3 | 20 | 8 - Dallas, TX |
| 9 | 2 | 20 | 9 - Birmingham, AL |
| 10 | 2 | 25 | 10 - Greenville, SC |
| 11 | 2 | 23 | 11 - Atlanta, GA |
| 12 | 2 | 19 | 12 - Midland, TX |
| 13 | 3 | 28 | 13 - Atlanta, GA |
| 14 | 3 | 27 | 14 - Houston, TX |
| 15 | 3 | 28 | 15 - Nashville, TN |
| 16 | 2 | 16 | 16 - Oklahoma City.. |
| 17 | 3 | 27 | 17 - Orlando, FL |
| 18 | 2 | 23 | 18 - Raleigh, NC |
| 19 | 3 | 20 | 19 - Miami, FL |
| 20 | 2 | 22 | 20 - Daytona Beac.. |
| 21 | 3 | 27 | 21 - Charlotte, NC |
| 22 | 2 | 25 | 22 - Knoxville, TN |
| 23 | 2 | 22 | 23 - Jacksonville, FL |
| 24 | 3 | 24 | 24 - Tampa, FL |
| 25 | 4 | 23 | 25 - Miami, FL |

northhighland® | Data and Analytics
WORLDWIDE CONSULTING

# WHY DO IT THIS WAY?

- Explainable
  - Client has minimal experience and trust of advanced analytics, a simple algorithm makes it easier to get buy-in
- Repeatable, with little variation
  - Similar but not identical results allow fine-tuning / re-running to smooth out client concerns
- Very easy to tweak in live sessions
  - Simple code, simple algorithms mean you can modify on-the-fly in response to questions
- In this case, all solutions are approximations
  - There's no right answer

# SOME OTHER PROJECTS

Are there natural clusters and needs of customers/employees?

Based on forecasted vs. actual sales, what stores are under-performing? Where should the next store be located?

Which customers/employees are likely to churn? Why?

Among elderly population, who is likely to need assisted living?

Which customers are likely to click/convert

**Key Business Questions**

Which patients are likely to be readmitted? Why?

How do we create robust tests of content customers are most likely to respond to?

Can we use predictive maintenance to minimize production impacts?

Who are most likely social influencers?

What is the next best action/offer for each customer?

**Advanced Analytics Toolkit**

Predictive/Explanatory modeling

Behavioral segmentation

Survey segmentation and projection

Forecasting

Pricing analytics

Design of Experiments (A/B and MVT)

Text/VOC analytics

Social influence propensity

**north**highland® | Data and Analytics
WORLDWIDE CONSULTING

# THANK YOU

www.northhighland.com

**CHARLIE MORN**

Sr. Data Analyst
North Highland
charlie.morn@northhighland.com

**northhighland**® | Data and Analytics
WORLDWIDE CONSULTING

## QUICK OIL CHANGE CHAIN

### PROBLEM

Our client has a large base of customers that are "oil-only" and have never used them for mechanical services (e.g., belts, brakes, hoses)

### SOLUTION

Develop a predictive model used to target customers most likely to convert so they can receive a differentiated experience on their next visit.
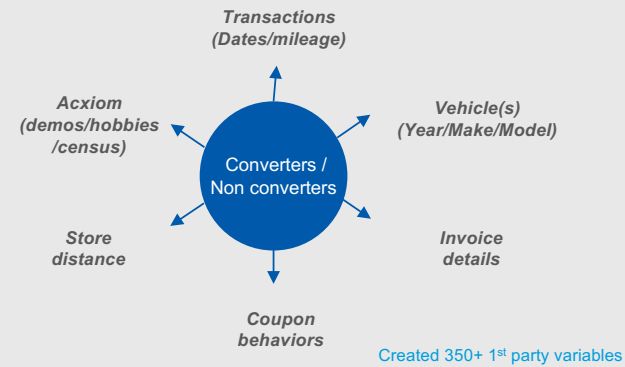
*Perform deep data-mining of prevailing customer behaviors to identify ones that tend to lead to conversion and just as important, ones that might turn off customers (e.g., "over-selling")*

*A sound byte from the modeling process is that air filter replacement recommendations tend to turn customers off and reduce their chance of mechanical conversion by 25%.*
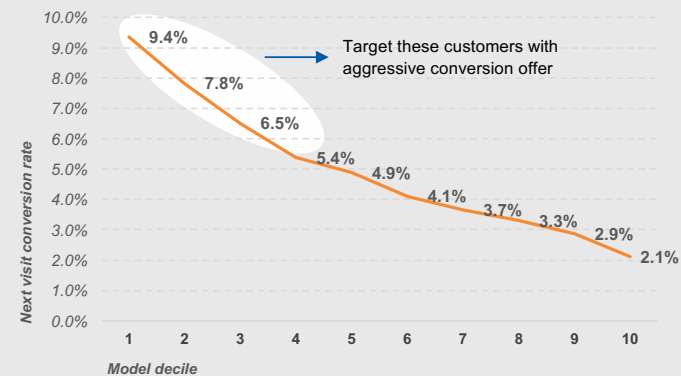
### RESULTS

Paid back initial investment at two month mark (based on net EBIT)

At three months (mid-October 2016), converted 1,377 customers for a total of $350k net NEW mechanical revenue.

---

*Transactions (Dates/mileage)*

*Acxiom (demos/hobbies /census)*

*Vehicle(s) (Year/Make/Model)*

Converters / Non converters

*Store distance*

*Invoice details*

*Coupon behaviors*

Created 350+ 1st party variables

---

**PREDICTIVE MODEL PERFORMANCE**

Next visit conversion rate

- 9.4%
- 7.8%
- 6.5%
- 5.4%
- 4.9%
- 4.1%
- 3.7%
- 3.3%
- 2.9%
- 2.1%

Target these customers with aggressive conversion offer

Model decile: 1 2 3 4 5 6 7 8 9 10

*Theory matches reality*
*Decile 1 – Most likely to convert >> highest next visit conversion (9.4%)*
*Decile 10 – Least likely to convert >> lowest next visit conversion (2.1%)*

northhighland® | Data and Analytics
WORLDWIDE CONSULTING