# DS Roundtable: An Introduction to Linear Programming

**7/21/2020**

**Andrew Wheeler, PhD**

**DSML Team**

**andrew.wheeler@hms.com**

# Overview of the Problem

- Models give a predicted probability or a predicted value

- Predictions do not intrinsically translate to *what* you should do with the information

- Linear programming is how to make decisions based on those model outputs given costs/benefits and constraints

# Example 1

- Predicted probability *of rain* is 5%

  - would you bring your umbrella?

- Predicted probability *of used car not working* is 5%

  - Would you buy that car?

- What is the difference? **Costs/Benefits** of false negative

# Example 2

- Utility of getting fast food for dinner

  - Benefits – fast, cheap

  - Costs – unhealthy, not as good as home cooked meals

- Costs/Benefits don't change, so why do you get fast food sometimes and not others?

  - Kids have event after school, more likely to get fast food

- What is the difference? **Constraints**

# Cost/Benefits & Constraints

- Better Example: When to audit a claim to determine overpayment, based on predictive model probability

- Costs: time it takes to review claim

- Potential Benefits: Probability of overpayment * Magnitude of Overpayment

- Constraints: We want the finding rate to be above a particular level and to not overburden particular providers

# Simple Example (no constraints)

Select 2 cases to audit:

| | Probability of Overpayment | Estimate of Overpayment Amount | Expected Return |
|---|---|---|---|
| Claim 1 | 0.05 | $10,000 | ? |
| Claim 2 | 0.10 | $3,000 | ? |
| Claim 3 | 0.50 | $500 | ? |

# Simple Example (no constraints)

Select 2 cases to audit:

|  | Probability of Overpayment | Estimate of Overpayment Amount | Expected Return |
|---|---:|---:|---:|
| Claim 1 | 0.05 | $10,000 | $500 |
| Claim 2 | 0.10 | $3,000 | $300 |
| Claim 3 | 0.50 | $500 | $250 |

You would select Claims 1 and 2

- total expected return = (500 + 300)

- Estimated overall finding rate is (0.05 + 0.10)/2 = 0.075

# Simple Example (constrain finding rate above 20%)

Select 2 cases to audit, but need to keep finding rate above 20%

|  | Probability of Overpayment | Estimate of Overpayment Amount | Expected Return |
| --- | --- | --- | --- |
| Claim 1 | 0.05 | $10,000 | $500 |
| Claim 2 | 0.10 | $3,000 | $300 |
| Claim 3 | 0.50 | $500 | $250 |

You would select Claims 1 and 3

- total expected return = (500 + 250)

- Estimated overall finding rate is (0.05 + 0.50)/2 = 0.275

# Using python *pulp* library to specify models

Making Example Data

Setting Up Model & Solving

```python
1 '''
2 Linear programming examples
3 Andy Wheeler, andrew.wheeler@hms.com
4 '''
5
6 #This is the library I like to use
7 #Many exist though
8 import pulp
9
10 #These are only necessary for simulating data
11 #You can pass in lists to pulp
12 import numpy as np
13 import pandas as pd
14
15 ##############################################
16 #Simple Example
17
18 #Creating data
19 prob = [0.05, 0.10, 0.50]
20 over = [10000, 3000, 500]
21 exp_ret = [p*o for p,o in zip(prob,over)]
22 case_index = list(range(len(prob)))
23 tot_audit = 2 #total number of claims to select
24 hit_rate = 0.2 #finding rate constraint
25
```

```python
25
26 #This is the model
27 P = pulp.LpProblem("Choosing Cases to Select", pulp.LpMaximize)
28
29 #These are the binary decision variables
30 D = pulp.LpVariable.dicts("Decision Variable", [i for i in case_index],
31                           lowBound=0, upBound=1, cat=pulp.LpInteger)
32
33 #Objective Function
34 P += pulp.lpSum( D[i]*exp_ret[i] for i in case_index)
35
36 #Constraint on total number of claims selected
37 P += pulp.lpSum( D[i] for i in case_index ) == tot_audit
38
39 #Constraint on the overall hit rate
40 P += pulp.lpSum( D[i]*prob[i] for i in case_index ) >= hit_rate*tot_audit
41
42 #Solve the problem
43 P.solve()
44
```

# Simulating Data to Look Like Claims

## Making Example Data

```python
63  ###############################################
64  #More complicated example data
65  #Has example providers as well
66
67  np.random.seed(10)
68  n = 20000 #total number of cases
69  underpay_est = np.random.lognormal(mean=7.6,sigma=0.5,size=n)
70  #about 25% overall finding rate
71  prob_over = np.random.beta(2.5,10,size=n)
72  exp_return = prob_over*underpay_est
73
74  #providers have a differential probability of being selected
75  prov = list('ABCDE')
76  prov_n = len(prov)
77  prov_index = list(range(1,prov_n+1))
78  prov_tot = sum(prov_index)
79  prov_prob = [i/prov_tot for i in prov_index]
80  prov_claims = np.random.choice(a=prov, size=n, replace=True, p=prov_prob)
81
82  sim_dat = pd.DataFrame(zip(underpay_est, prob_over,
83                        exp_return, prov_claims),
84                    columns=['underpay_est', 'prob_over',
85                             'exp_return', 'prov_claims'])
86
87  ###############################################
```

## Function with Provider Constraints

```python
93  def selection_model(er,prob,provider,data,cases_const,finding_const,provider_const):
94      #Preparing simpler lists
95      er_list = list(data[er])
96      min_const_list = list( data[prob] - finding_const )
97      prov_list = list(data[provider])
98      index_list = list(range(len(er_list)))
99      #getting the locations for each provider in the data
100     all_prov = set(prov_list)
101     prov_loc = {}
102     for p in all_prov:
103         prov_loc[p] = list( np.where(data[provider] == p)[0])
104     #Now setting up the model
105     Sel_Mod = pulp.LpProblem("Selection Model", pulp.LpMaximize)
106     Dec_Vars = pulp.LpVariable.dicts("Selected Cases",
107                     [i for i in sim_dat.index],
108                     lowBound=0, upBound=1,
109                     cat=pulp.LpInteger)
110     #Objective Function
111     Sel_Mod += pulp.lpSum( Dec_Vars[i]*er_list[i] for i in index_list)
112     #Constraint on total number of claims selected, has to be equal or fewer
113     Sel_Mod += pulp.lpSum( Dec_Vars[i] for i in index_list ) <= cases_const
114     #Constraint on finding rate, taking into account total cases selected
115     Sel_Mod += pulp.lpSum( Dec_Vars[i]*min_const_list[i] for i in index_list ) >= 0
116     #Provider constraints
117     for p in all_prov:
118         Sel_Mod += pulp.lpSum( Dec_Vars[i] for i in prov_loc[p] ) <= provider_const
119     #Solve the Problem
120     Sel_Mod.solve()
121     #Get the decision variables
122     dec_list = [Dec_Vars[i].varValue for i in index_list]
123     return dec_list
```

# Provider Constraints

No Provider Constraints

```
In [11]: sim_dat['Selected1'] = selection_model(er='exp_return',prob='prob_over',provider='prov_claims',
    ...:                                   data=sim_dat,cases_const=1000,finding_const=0.3,
    ...:                                   provider_const=1000)
    ...:
    ...:
    ...: print( sim_dat.loc[ sim_dat['Selected1'] == 1, 'prov_claims'].value_counts() )
    ...:
E    350
D    269
C    199
B    114
A     68
Name: prov_claims, dtype: int64

In [12]:
```

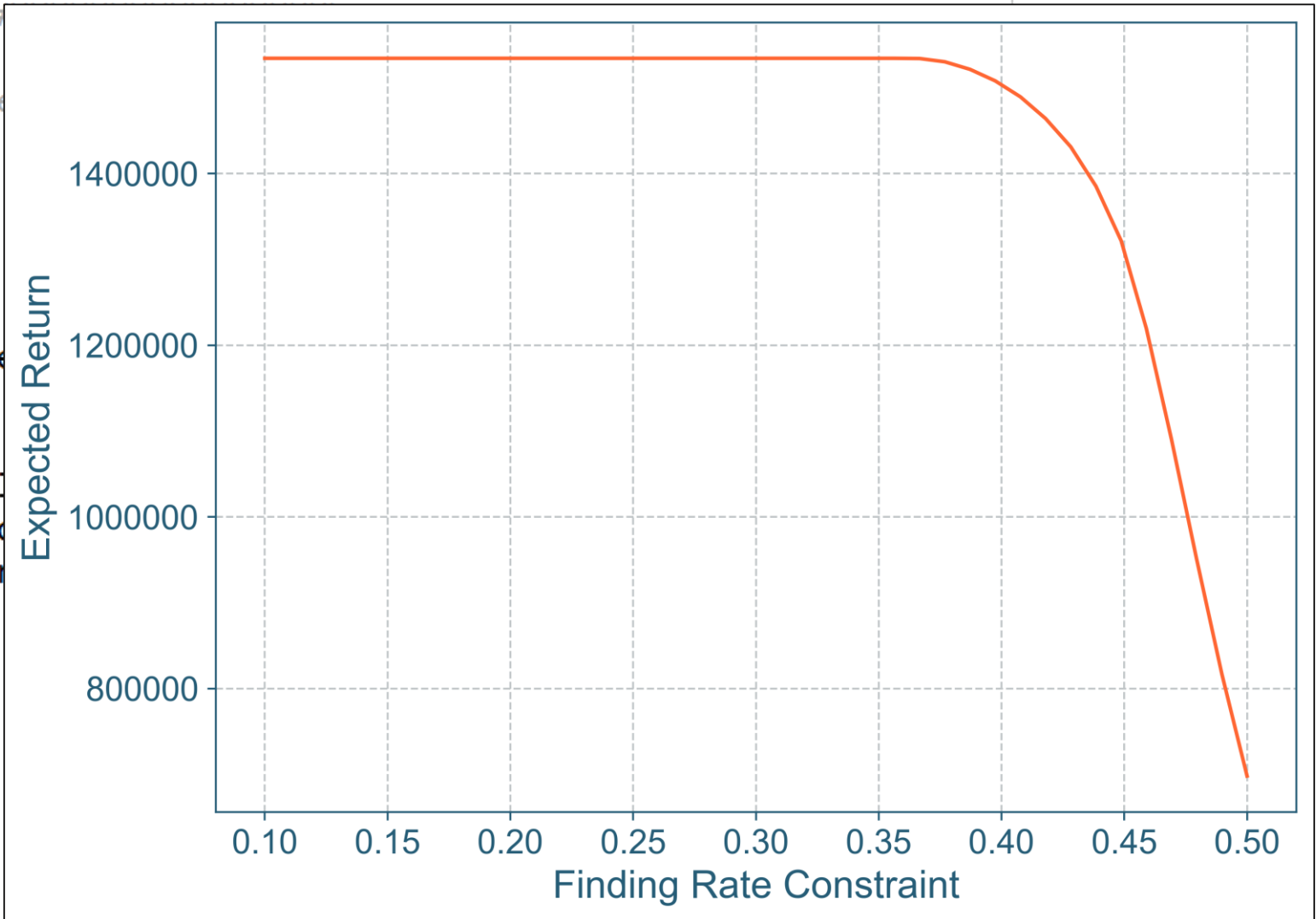Provider Constraints <= 250

```
In [12]: sim_dat['Selected2'] = selection_model(er='exp_return',prob='prob_over',provider='prov_claims',
    ...:                                   data=sim_dat,cases_const=1000,finding_const=0.3,
    ...:                                   provider_const=250)
    ...:
    ...: print( sim_dat.loc[ sim_dat['Selected2'] == 1, 'prov_claims'].value_counts() )
    ...: ################################################
E    250
C    250
D    250
B    160
A     90
Name: prov_claims, dtype: int64
```

# The Finding Rate Frontier

```python
142 ##############################################
143 #We can see how changing the finding rate constraint
144 #Effects our estimates of the expected return
145
146 hit_rate = np.linspace(0.10, 0.50, 40)
147 total_return = []
148
149 for h in hit_rate:
150     sel_list = selection_model(er='exp_return',prob='prob_over',provider='prov_claims',
151                                 data=sim_dat,cases_const=1000,finding_const=h,
152                                 provider_const=250)
153     sel_np = np.asarray(sel_list)
154     est_ret = (sel_np * sim_dat['exp_return']).sum()
155     total_return.append(est_ret)
```

# The Finding Rate Frontier

```
142 ##############################
143 #We can see how changing the
144 #Effects our estimates of the
145
146 hit_rate = np.linspace(0.10,
147 total_return = []
148
149 for h in hit_rate:
150     sel_list = selection_mode
151
152
153     sel_np = np.asarray(sel_
154     est_ret = (sel_np * sim_
155     total_return.append(est_
```

# Links to Resources

- Python code giving example use, https://github.com/hmsholdings/data-science-utils/tree/master/education/Advanced_DataScience/IntroLinearProgramming

- I have multiple blog posts/papers illustrating different linear programming problems:

    - An intro to linear programming for criminologists

    - Optimal treatment selection with network spillovers

    - Optimizing police patrol areas with workload constraints

    - Racial Equity constraints in predictive policing applications


- The blog Yet Another Math Programming Consultant gives many different examples of linear programming applications

# Questions?

# DS Roundtable: An Introduction to Linear Programming

**7/21/2020**

**Andrew Wheeler, PhD**

**DSML Team**

**andrew.wheeler@hms.com**