

'MATH+ECON+CODE' MASTERCLASS ON MATCHING MODELS, OPTIMAL TRANSPORT AND APPLICATIONS

Alfred Galichon (NYU)

Spring 2018

Day 1, January 15 2018: Linear programming
Block 1. Linear programming duality

- ▶ Schedule: Mon 1/15 – Sat 1/20, 2018, 8am-12 noon and 1pm-3pm.
- ▶ Location: Courant Institute (Warren Weaver Hall, 251 Mercer) WWH102.
- ▶ Office hours: by appointment (my email: ag133@nyu.edu).
- ▶ Course webpage: <http://alfredgalichon.com/matheconcode/>
- ▶ Students (including auditors) need to register on Albert (MATH-GA 2840.002 or ECON-GA 3002.015)
- ▶ Text (optional): Galichon (2016). *Optimal Transport Methods in Economics*, Princeton.

- ▶ A. Galichon: professor of economics and of mathematics at NYU
- ▶ Keith O'hara: NYU econ grad student (econometrics; R and C++)
- ▶ Yifei Sun: NYU math grad student (machine learning; R and Python).

- ▶ Jointly offered in the Math and Econ programs. Self-contained for both audiences
- ▶ Teaching format: 18 “blocks”; each block = 50 minutes of theory + 1 hour of coding
 - ▶ coding most often based on an empirical application related to the theory just seen
 - ▶ students are expected to write their own code; we’ll ensure that it is operational at the end of each block
- ▶ Programming: our demos will be done in R and the support will be in R only, but you are welcome to use the language of your choice e.g. Matlab, C++, Python, Julia... Solvers used will be Gurobi (for LP) and NLOPT (for nonlinear optimization), so make sure your language of choice has a convenient interface to these.
- ▶ Questions?

- ▶ This course is focused on models of demand, matching models, and optimal transport methods, with various applications pertaining to labor markets, economics of marriage, industrial organization, matching platforms, networks, and international trade, from the crossed perspectives of theory, empirics and computation.
- ▶ It will introduce tools from economic theory, mathematics, econometrics and computing, on a needs basis, without any particular prerequisite other than the equivalent of a first year graduate sequence in econ or in applied math.

► Part I: Tools

- Monday 1/15: linear programming
- Tuesday 1/16: optimal transport toolbox
- Wednesday 1/17: quantile methods

► Part II: Models

- Thursday 1/18: static and dynamic multinomial choice
- Friday 1/19: statistical estimation of models of matching with transfers
- Saturday 1/20: more general models of matching

- ▶ (<https://www.r-project.org/>), Rstudio (<https://www.rstudio.com/>), and, for Windows users only, Rtools (<https://cran.r-project.org/bin/windows/Rtools/>).
- ▶ Jupyter (<http://jupyter.org/>), the R kernel for Jupyter (follow instructions from <https://www.datacamp.com/community/blog/jupyter-notebook-r>)
- ▶ gurobi (www.gurobi.com; you'll need to obtain an academic license, which is free if you're located on an educational domain).
- ▶ Please also register with a free account on Github

On a sheet of paper, please indicate:

1. Your name and email
2. Your program, department and institution
3. Whether you are a registered student (RS) or an approved auditor (AA)
4. Whether you are taking this course for credit
5. What you are looking for in this course (2 sentences max).
6. From a scale of 1 to 5, whether you are less familiar (1) or more familiar (5) with the concepts below:

- | | |
|--------------------------|------------------------------------|
| 1. Linear programming | 2. Legendre-Fenchel transforms |
| 3. Min-cost flow problem | 4. Becker's theory of marriage |
| 5. Backward induction | 6. Compensating wage differentials |
| 7. Envelope theorem | 8. Complementary slackness |
| 9. Walrasian equilibrium | 10. Hotelling's spatial model |
| 11. Logit choice model | 12. Gradient descent |
| 13. Newton descent | 14. Coordinate descent |

Note: even if you answer mostly “1”, don’t worry! These concepts will be explained in due course.

- ▶ Choice between:
 - ▶ Either a take-home exam (24 hours) available e.g. from Saturday Jan 27 noon, until Sunday Jan 28 noon (**to be confirmed with the class**),
 - ▶ Or a short paper (12 pages or more), to be discussed with the instructor. The paper will bear some connections, in a broad sense, with the topics of the course. Many papers are considered acceptable: original research paper, survey paper, report on numerical experiments, replication of existing empirical results. . . are all acceptable. The requirement is to be innovative on a theoretical, empirical, or computational level. This work should be submitted before June 30, 2018.
- ▶ The course will be assessed on a pass/fail basis.
- ▶ Approved auditing students should morally commit to take part in the assessment too.

Block 1: Linear programming

- ▶ Linear programming duality
- ▶ Economic interpretation of the dual
- ▶ Numerical computation

- ▶ [OTME], App. B
- ▶ Stigler (1945), The cost of subsistence. *Journal of Farm Economics*.
- ▶ Dantzig (1990), The diet problem. *Interface*.
- ▶ Complements:
 - ▶ Gale (1960), *The theory of linear economic models*.
 - ▶ Vohra (2011), *Mechanism Design: A Linear Programming Approach*.
- ▶ www.gurobi.com
- ▶ www.gnu.org/software/glpk/

Section 1

THEORY

MOTIVATION: THE DIET PROBLEM

- ▶ During World War II, engineers in US Army were wondering how to feed their personnel at minimal cost, leading to what is now called the “optimal diet problem”.
 - ▶ Nutritionists have identified a number of vital nutrients (calories, protein, calcium, iron, etc.) that matter for a person's health, and have determined the minimum daily intake of each nutrient
 - ▶ For each basic food (pasta, butter, bread, etc), nutritionists have characterized the intake in each of the various nutrients
 - ▶ Each food has a unit cost, and the problem is to find the optimal diet = combination of foods that meet the minimal intake in each of the nutrients and achieves minimal cost
- ▶ The problem was taken on by G. Stigler, who published a paper about it in 1945, giving a first heuristic solution, exhibiting a diet that costs \$39.93 per year in 1939 dollars. Later (in 1947) it was one of the first application of G.B. Dantzig's method (the simplex algorithm), which provided the exact solution (\$39.67). It then took 120 man-day to perform this operation. At the end of this block, the computer will perform it for us in a fraction of second.
- ▶ However, don't try this diet at home! Dantzig did so and almost died from it...

ALFRED GALICHON

- 'MATH+ECON+CODE': MATCHING MODELS, OPTIMAL TRANSPORT AND APPLICATIONS

Nevertheless standards of dietary adequacy have been established, perhaps prematurely and certainly very tentatively. The "allowances" (a term used to indicate their preliminary nature) of the National Research Council embody what is presumably the 1948 consensus of the experts; they are given in Table 1. Other minerals and vitamins are believed to be supplied in adequate quantities if these nutrients are secured from natural foods. The requirements are not of lessor in the preparation of food. These standards are met by the minimum cost diets derived subsequently.

* National Research Council, *Recommended Dietary Allowances*, Reprint and Circular Series No. 113, January, 1943.

Nutritive Values of Foods

The minimum cost of an adequate diet is obviously governed by the nutritive values and costs of the foods eligible for inclusion. The very restricted list of foods considered in this study is discussed in Section 8 and the foods are listed in Tables A and B. It may be mentioned here that only natural foods are included; vitamin pills are excluded because they do not contain all of the nutrients (known and unknown) which are necessary to good health.¹

This content downloaded from 128.122.11.192 on Wed, 27 Dec 2016 22:19:50 UTC
All use subject to <http://about.jstor.org/terms>

306

GEORGE J. J.

This content downloaded from 128.122.11.192 on Wed, 27 Dec 2016 22:19:50 UTC
All use subject to <http://about.jstor.org/terms>

- ▶ Problem setup:
 - ▶ Assume there are nutrients $i \in \{1, \dots, m\}$ (calories, protein, calcium, iron, etc.) that matter for a person's health, in such way that the minimum daily intake of nutrient i should be d_i .
 - ▶ Nutrients do not come as standalone elements, but are combined into various foods. Each unit of food $j \in \{1, \dots, n\}$ yields a quantity N_{ij} of nutrient $i \in \{1, \dots, m\}$. The dollar cost of food j is c_j .
- ▶ The problem is to find the diet that achieves the minimal intake of each nutrient at a cheapest price. If $q \in \mathbb{R}^n$ is a vector such that $q_j \geq 0$ is the quantity of food j purchased, the quantity of nutrient i ingested is $\sum_{j=1}^n N_{ij}q_j$, and the cost of the diet is $\sum_{j=1}^n q_j c_j$. The optimal diet is therefore given by

$$\begin{aligned} \min_{q \geq 0} \quad & c^\top q \\ \text{s.t.} \quad & Nq \geq d. \end{aligned} \tag{1}$$

- Let $c \in \mathbb{R}^n$, $d \in \mathbb{R}^m$, A be a $m \times n$ matrix, and consider the following problem

$$\begin{aligned} V_P &= \max_{x \in \mathbb{R}_+^n} c^\top x \\ \text{s.t. } Ax &= d \end{aligned} \tag{2}$$

This problem is a *linear programming problem*, as the objective function, namely $x \rightarrow c^\top x$ is linear, and as the constraint, namely $x \in \mathbb{R}_+^n$ and $Ax = d$ are also linear (or more accurately, affine). Problem (2) is called *primal program*, for reasons to be explained soon. The set of x 's that satisfy the constraint are called *feasible solutions*; the set of solutions of problem (2) are called *optimal solutions*.

- Remarks:
 - The previous diet problem can be reformulate into this problem – why?
 - A problem does not necessarily have a feasible solution (e.g. if $A = 0$ and $d \neq 0$), in which case (by convention) $V_P = -\infty$.
 - The whole space may be solution (e.g. if $A = 0$ and $d = 0$), in which case $V_P = +\infty$.

There is a powerful tool called duality which provides much insight into the analysis of problem (2). The idea is to rewrite the problem as

$$V_P = \max_{x \in \mathbb{R}_+^n} \left\{ c^\top x + L_P(d - Ax) \right\}$$

where $L_P(z)$ is a penalty function whose value is zero if the constraint is met, that is if $z = 0$, and $-\infty$ if it is not, namely if $z \neq 0$. The simplest choice of such penalty function is given by $L_P(z) = \min_{y \in \mathbb{R}^m} \{z^\top y\}$. One has

$$V_P = \max_{x \in \mathbb{R}_+^n} \min_{y \in \mathbb{R}^m} \left\{ c^\top x + (d - Ax)^\top y \right\}.$$

However, the minimax inequality $\max_x \min_y \leq \min_y \max_x$ always holds, thus

$$\begin{aligned} V_P &\leq \min_{y \in \mathbb{R}^m} \max_{x \in \mathbb{R}_+^n} \left\{ c^\top x + (d - Ax)^\top y \right\} = \min_{y \in \mathbb{R}^m} \max_{x \in \mathbb{R}_+^n} \left\{ x^\top (c - A^\top y) + d^\top y \right\} \\ &\leq \min_{y \in \mathbb{R}^m} \left\{ d^\top y + L_D (c - A^\top y) \right\} =: V_D \end{aligned}$$

where $L_D(z) = \max_{x \in \mathbb{R}_+^n} \{x^\top z\}$ is equal to 0 if $z \in \mathbb{R}_-^n$, and to $+\infty$ if not. Therefore, the value V_D is expressed by the *dual program*

$$\begin{aligned} V_D &= \min_{y \in \mathbb{R}^m} d^\top y, \\ \text{s.t. } &A^\top y \geq c \end{aligned} \tag{3}$$

and the weak duality inequality $V_P \leq V_D$ holds. It turns out that as soon as either the primal or dual program has an optimal solution, then both programs have an optimal solution and the values of the two programs coincide, so the weak duality becomes an equality $V_P = V_D$ called strong duality. Further, if $x^* \in \mathbb{R}_+^n$ is an optimal primal solution, and $y^* \in \mathbb{R}^m$ is an optimal dual solution, then complementary slackness holds, that is $x_i^* > 0$ implies $(A^\top y^*)_i = c_i$.

We summarize these results into the following statement.

Theorem. In the setting described above:

(i) The weak duality inequality holds:

$$V_P \leq V_D.$$

(ii) As soon as the primal or the dual program have an optimal solution, then both programs have an optimal solution, and strong duality holds:

$$V_P = V_D.$$

(iii) If $x^* \in \mathbb{R}_+^n$ is an optimal primal solution, and $y^* \in \mathbb{R}^m$ is an optimal dual solution, then complementary slackness holds:

$$x_i^* > 0 \text{ implies } \left(A^\top y^* \right)_i = c_i.$$

- Recall the optimal diet problem

$$\begin{aligned} \min_{q \geq 0} c^\top q \\ \text{s.t. } Nq \geq d. \end{aligned}$$

which has minimax formulation $\min_{q \geq 0} \max_{\pi \geq 0} c^\top q + d^\top \pi - q^\top N^\top \pi$, so the dual is

$$\begin{aligned} \max_{\pi \geq 0} d^\top \pi \\ \text{s.t. } N^\top \pi \leq c \end{aligned}$$

- Interpretation: imagine that there is a new firm called Nutrient Shoppe, who sells raw nutrients. Let π_i be the price of nutrient i . The cost of the diet is $d^\top \pi$. Consumer purchase raw nutrients and can generate “synthetic” foods. The cost of the synthetic version of food j is $\sum_{i=1}^m N_{ij} \pi_i = (N^\top \pi)_j$. The constraint thus means that each “synthetic” food is more affordable than its natural counterpart.

- ▶ The duality means that it is possible to price the nutrients so that the synthetic foods are cheaper than the natural ones, in such a way that the price of the synthetic diet equals the price of the natural diet.
- ▶ Complementary slackness yields:
 - ▶ $q_j > 0$ implies $(N^T \pi)_j = c_j$; that is, if natural food j is actually purchased, then the prices of its synthetic and natural versions coincide
 - ▶ $\pi_i > 0$ implies $(Nq)_i = d_i$; that is, if nutrient i has a positive price, then the natural diet has the “just right” amount.

Section 2

CODING

- ▶ Install R (<https://www.r-project.org/>), a free statistical environment
- ▶ Install Rstudio (<https://www.rstudio.com/>), a free IDE for R
- ▶ [Windows users only] Install Rtools (<https://cran.r-project.org/bin/windows/Rtools/>), a toolchain for Windows needed to install certain R packages
- ▶ Install Gurobi (www.gurobi.com) – a state-of-the-art commercial lp solver
 - ▶ Install the program
 - ▶ Obtain a license for Gurobi
 - ▶ Install R package for Gurobi
- ▶ Install Rglpk (from Rstudio) – an R interface to the GLPK open-source lp solver

- ▶ Today we shall retrieve Stigler's diet data and compute the optimal diet in order to compare with Stigler's computations
- ▶ We shall do so from R, using in turn Gurobi and GLPK.

```
library(gurobi)
thepath = getwd()
filename="/StiglerData1939.txt"
thedata =
as.matrix(read.csv(paste0(thepath,filename),sep="\t",
header=T))
nbCommodities=length(which(thedata[,1] != "" ))-1
names = thedata[1:nbCommodities,1]
themat = matrix( as.numeric(thedata[,3:13]), ncol = 11)
themat[is.na(themat)] = 0
```

```
N = t(themat[1:nbCommodities,3:11])
d = themat[(nbCommodities+1),3:11]
c = rep(1,nbCommodities)
result = gurobi (
  list(A=N,obj=c,modelsense="min",rhs=d,sense=">"),
  params=list(OutputFlag=0) )
q_yearly = result$x * 365 # convert into yearly cost
pi = result$pi
cost_daily = result$objval
```

Remark: by default, Gurobi assumes the constraint $x \geq 0$. To remove this constraint, include `ub = -Inf` in the list passed to Gurobi.

```
toKeep = which(q_yearly !=0 )  
foods = q_yearly[toKeep]  
names(foods) = names[toKeep]  
print(foods)  
print(paste0("Total cost (optimal)= ", sum(q_yearly*c) ))
```

```
toKeepStigler = c(1,15,46,52,69)
foods_stigler = c(13.33, 3.84,4.11,1.85,16.80)
names(foods_stigler) = names[toKeepStigler]
print(foods_stigler)
print(paste0("Total cost (Stigler)= ",
sum(foods_stigler*c[toKeepStigler])) )
```

```
library(Rglpk)
resGlpk = Rglpk_solve_LP(obj=c, mat=N,
  dir=rep(">",length(d)), rhs=d, bounds = NULL, max = FALSE,
  control = list())
print(resGlpk$optimum*365)
```