# Linear Programming: Chapter 12
# Regression

Robert J. Vanderbei

October 17, 2007

Operations Research and Financial Engineering
Princeton University
Princeton, NJ 08544
http://www.princeton.edu/∼rvdb

# Outline

- Means and Medians

- Least Squares Regression

- Least Absolute Deviation (LAD) Regression

- LAD via LP

- Average Complexity of Parametric Self-Dual Simplex Method

# 1995 Adjusted Gross Incomes

Consider 1995 Adjusted Gross Incomes on Individual Tax Returns:

| $b_1$ | $b_2$ | $b_3$ | $\ldots$ | $b_{m-1}$ | $b_m$ |
|---|---|---|---|---|---|
| $25,462 | $45,110 | $15,505 | $\ldots$ | $33,265 | $75,420 |

Real summary data is shown on the next slide...

# Table 1.--1995, Individual Income Tax Returns

[All figures are estimates based on samples--money amounts are in thousands of dollars]

| Size of adjusted gross income | Number of returns | Adjusted gross income |
|---|---|---|
| | (1) | (2) |
| **All returns** | **118,218,327** | **4,189,353,615** |
| No adjusted gross income | 944,141 | -55,253,648 |
| $1 under $5,000 | 14,646,131 | 37,604,828 |
| $5,000 under $10,000 | 13,982,404 | 104,603,365 |
| $10,000 under $15,000 | 13,562,088 | 169,317,443 |
| $15,000 under $20,000 | 11,385,632 | 198,418,324 |
| $20,000 under $25,000 | 9,970,099 | 223,400,219 |
| $25,000 under $30,000 | 7,847,862 | 215,200,244 |
| $30,000 under $40,000 | 12,380,339 | 430,491,242 |
| $40,000 under $50,000 | 9,098,760 | 406,638,597 |
| $50,000 under $75,000 | 13,679,023 | 828,349,278 |
| $75,000 under $100,000 | 5,374,489 | 458,505,650 |
| $100,000 under $200,000 | 4,074,852 | 532,030,480 |
| $200,000 under $500,000 | 1,007,136 | 292,117,517 |
| $500,000 under $1,000,000 | 178,374 | 120,347,093 |
| $1,000,000 or more | 86,998 | 227,582,987 |
| **Taxable returns** | **89,252,989** | **4,007,580,441** |
| **Nontaxable returns** | **28,965,338** | **181,773,174** |

# Means and Medians

Median:

$$\hat{x} = b_{\frac{1+m}{2}} \approx \$22{,}500.$$

Mean:

$$\bar{x} = \frac{1}{m} \sum_{i=1}^{m} b_i = \$4{,}189{,}353{,}615{,}000/118{,}218{,}327 = \$35{,}437.$$

# Mean's Connection with Optimization

$$\bar{x} = \mathsf{argmin}_x \sum_{i=1}^{m} (x - b_i)^2 \,.$$

Proof:

$$
\begin{aligned}
f(x) &= \sum_{i=1}^{m} (x - b_i)^2 \\
f'(x) &= \sum_{i=1}^{m} 2\,(x - b_i) \\
f'(\bar{x}) &= 0 \qquad \Longrightarrow \qquad \bar{x} = \frac{1}{m} \sum_{i=1}^{m} b_i \\
\lim_{x \to \pm\infty} f(x) &= +\infty \qquad \Longrightarrow \qquad \bar{x} \text{ is a minimum}
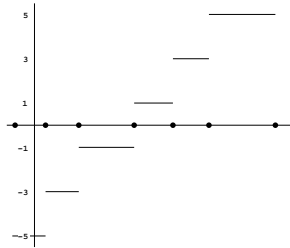\end{aligned}
$$

# Median's Connection with Optimization

$$\hat{x} = \mathsf{argmin}_x \sum_{i=1}^{m} |x - b_i|.$$

Proof:

$$f(x) = \sum_{i=1}^{m} |x - b_i|$$

$$f'(x) = \sum_{i=1}^{m} \mathsf{sgn}\,(x - b_i) \qquad \text{where } \mathsf{sgn}(x) = \begin{cases} 1 & x > 0 \\ 0 & x = 0 \\ -1 & x < 0 \end{cases}$$

$$= (\# \text{ of } b_i\text{'s smaller than } x) - (\# \text{ of } b_i\text{'s larger than } x).$$

If $m$ is odd:

# A Regression Model for Algorithm Efficiency

*Observed Data:*

$$
\begin{aligned}
t &= \#\text{ of iterations} \\
m &= \#\text{ of constraints} \\
n &= \#\text{ of variables}
\end{aligned}
$$

*Model:*

$$t \approx 2^{\alpha}(m+n)^{\beta}$$

*Linearization:* Take logs:

$$\log t = \alpha \log 2 + \beta \log(m+n) + \quad \epsilon$$

$\uparrow$
error

# Regression Model Continued

Solve several instances:

$$\begin{bmatrix} \log t_1 \\ \log t_2 \\ \vdots \\ \log t_k \end{bmatrix} = \begin{bmatrix} \log 2 & \log(m_1 + n_1) \\ \log 2 & \log(m_2 + n_2) \\ \vdots & \vdots \\ \log 2 & \log(m_k + n_k) \end{bmatrix} \begin{bmatrix} \alpha \\ \beta \end{bmatrix} + \begin{bmatrix} \epsilon_1 \\ \epsilon_2 \\ \vdots \\ \epsilon_k \end{bmatrix}$$

In matrix notation:

$$b = Ax + \epsilon$$

*Goal:* find $x$ that "minimizes" $\epsilon$.

# Least Squares Regression

Euclidean Distance: $\|x\|_2 = \left(\sum_i x_i^2\right)^{1/2}$

Least Squares Regression: $\bar{x} = \text{argmin}_x \|b - Ax\|_2^2$

Calculus:

$$f(x) = \|b - Ax\|_2^2 = \sum_i \left(b_i - \sum_j a_{ij}x_j\right)^2$$

$$\frac{\partial f}{\partial x_k}(\bar{x}) = \sum_i 2\left(b_i - \sum_j a_{ij}\bar{x}_j\right)(-a_{ik}) = 0, \qquad k = 1, 2, \ldots, n$$

Rearranging,

$$\sum_i a_{ik}b_i = \sum_i \sum_j a_{ik}a_{ij}\bar{x}_j, \qquad k = 1, 2, \ldots, n$$

In matrix notation,

$$A^T b = A^T A \bar{x}$$

Assuming $A^T A$ is invertible,

$$\bar{x} = \left(A^T A\right)^{-1} A^T b$$

# Least Absolute Deviation Regression

*Manhattan Distance:* $\|x\|_1 = \sum_i |x_i|$
*Least Absolute Deviation Regression:* $\hat{x} = \mathsf{argmin}_x \|b - Ax\|_1$
*Calculus:*

$$f(x) = \|b - Ax\|_1 = \sum_i \left| b_i - \sum_j a_{ij} x_j \right|$$

$$\frac{\partial f}{\partial x_k}(\hat{x}) = \sum_i \frac{b_i - \sum_j a_{ij} \hat{x}_j}{\left| b_i - \sum_j a_{ij} \hat{x}_j \right|}(-a_{ik}) = 0, \qquad k = 1, 2, \ldots, n$$

Rearranging,

$$\sum_i \frac{a_{ik} b_i}{\epsilon_i(\hat{x})} = \sum_i \sum_j \frac{a_{ik} a_{ij} \hat{x}_j}{\epsilon_i(\hat{x})}, \qquad k = 1, 2, \ldots, n$$

In matrix notation,

$$A^T E(\hat{x}) b = A^T E(\hat{x}) A \hat{x}, \qquad \text{where } E(\hat{x}) = \mathsf{Diag}(\epsilon(\hat{x}))^{-1}$$

Assuming $A^T E(\hat{x}) A$ is invertible,

$$\hat{x} = \left( A^T E(\hat{x}) A \right)^{-1} A^T E(\hat{x}) b$$

# Least Absolute Deviation Regression— Continued

An implicit equation.
Can be solved using *successive approximations*:

$$
\begin{aligned}
x^0 &= 0 \\
x^1 &= \left(A^T E(x^0)A\right)^{-1} A^T E(x^0)b \\
x^2 &= \left(A^T E(x^1)A\right)^{-1} A^T E(x^1)b \\
&\vdots \\
x^{k+1} &= \left(A^T E(x^k)A\right)^{-1} A^T E(x^k)b \\
&\vdots \\
\hat{x} &= \lim_{k\to\infty} x^k
\end{aligned}
$$

# Least Absolute Deviation Regression via Linear Programming

$$\min \sum_i \left| b_i - \sum_j a_{ij} x_j \right|$$

Equivalent Linear Program:

$$\min \sum_i t_i$$

$$-t_i \ \leq \ b_i - \sum_j a_{ij} x_j \ \leq \ t_i \qquad i = 1, 2, \ldots, m$$

# AMPL Model

```
param m;
param n;

set I := {1..m};
set J := {1..n};

param A {I,J};
param b {I};

var x{J};
var t{I};

minimize sum_dev:
      sum {i in I} t[i];

subject to lower_bound {i in I}:
    -t[i] <= b[i] - sum {j in J} A[i,j]*x[j];

subject to upper_bound {i in I}:
            b[i] - sum {j in J} A[i,j]*x[j] <= t[i];
```

# Parametric Self-Dual Simplex Method

Thought experiment:
- $\mu$ starts at $\infty$.

- In reducing $\mu$, there are $n + m$ barriers.

- At each iteration, one barrier is passed—the others move about randomly.

- To get $\mu$ to zero, we must on average pass half the barriers.

- Therefore, on average the algorithm should take $(m + n)/2$ iterations.

Least Squares Regression:

$$\begin{bmatrix} \bar{\alpha} \\ \bar{\beta} \end{bmatrix} = \begin{bmatrix} -1.03561 \\ 1.05152 \end{bmatrix} \qquad \Longrightarrow \qquad T \approx 0.488(m + n)^{1.052}$$

Least Absolute Deviation Regression:

$$\begin{bmatrix} \hat{\alpha} \\ \hat{\beta} \end{bmatrix} = \begin{bmatrix} -0.9508 \\ 1.0491 \end{bmatrix} \qquad \Longrightarrow \qquad T \approx 0.517(m + n)^{1.049}$$

# Parametric Self-Dual Simplex Method: Data

| Name | $m$ | $n$ | iters | Name | $m$ | $n$ | iters |
|------|-----|-----|-------|------|-----|-----|-------|
| 25fv47 | 777 | 1545 | 5089 | nesm | 646 | 2740 | 5829 |
| 80bau3b | 2021 | 9195 | 10514 | recipe | 74 | 136 | 80 |
| adlittle | 53 | 96 | 141 | sc105 | 104 | 103 | 92 |
| afiro | 25 | 32 | 16 | sc205 | 203 | 202 | 191 |
| agg2 | 481 | 301 | 204 | sc50a | 49 | 48 | 46 |
| agg3 | 481 | 301 | 193 | sc50b | 48 | 48 | 53 |
| bandm | 224 | 379 | 1139 | scagr25 | 347 | 499 | 1336 |
| beaconfd | 111 | 172 | 113 | scagr7 | 95 | 139 | 339 |
| blend | 72 | 83 | 117 | scfxm1 | 282 | 439 | 531 |
| bnl1 | 564 | 1113 | 2580 | scfxm2 | 564 | 878 | 1197 |
| bnl2 | 1874 | 3134 | 6381 | scfxm3 | 846 | 1317 | 1886 |
| boeing1 | 298 | 373 | 619 | scorpion | 292 | 331 | 411 |
| boeing2 | 125 | 143 | 168 | scrs8 | 447 | 1131 | 783 |
| bore3d | 138 | 188 | 227 | scsd1 | 77 | 760 | 172 |
| brandy | 123 | 205 | 585 | scsd6 | 147 | 1350 | 494 |
| czprob | 689 | 2770 | 2635 | scsd8 | 397 | 2750 | 1548 |
| d6cube | 403 | 6183 | 5883 | sctap1 | 284 | 480 | 643 |
| degen2 | 444 | 534 | 1421 | sctap2 | 1033 | 1880 | 1037 |
| degen3 | 1503 | 1818 | 6398 | sctap3 | 1408 | 2480 | 1339 |
| e226 | 162 | 260 | 598 | seba | 449 | 896 | 766 |

# Data Continued

| Name | $m$ | $n$ | iters | Name | $m$ | $n$ | iters |
|---|---|---|---|---|---|---|---|
| etamacro | 334 | 542 | 1580 | share1b | 107 | 217 | 404 |
| fffff800 | 476 | 817 | 1029 | share2b | 93 | 79 | 189 |
| finnis | 398 | 541 | 680 | shell | 487 | 1476 | 1155 |
| fit1d | 24 | 1026 | 925 | ship04l | 317 | 1915 | 597 |
| fit1p | 627 | 1677 | 15284 | ship04s | 241 | 1291 | 560 |
| forplan | 133 | 415 | 576 | ship08l | 520 | 3149 | 1091 |
| ganges | 1121 | 1493 | 2716 | ship08s | 326 | 1632 | 897 |
| greenbea | 1948 | 4131 | 21476 | ship12l | 687 | 4224 | 1654 |
| grow15 | 300 | 645 | 681 | ship12s | 417 | 1996 | 1360 |
| grow22 | 440 | 946 | 999 | sierra | 1212 | 2016 | 793 |
| grow7 | 140 | 301 | 322 | standata | 301 | 1038 | 74 |
| israel | 163 | 142 | 209 | standmps | 409 | 1038 | 295 |
| kb2 | 43 | 41 | 63 | stocfor1 | 98 | 100 | 81 |
| lotfi | 134 | 300 | 242 | stocfor2 | 2129 | 2015 | 2127 |
| maros | 680 | 1062 | 2998 | | | | |

# Parametric Self-Dual Simplex Method: Plot



A log–log plot of $T$ vs. $m+n$ and the $L^1$ and $L^2$ regression lines.