Data Modeling and Least Squares Fitting

COS 323

### Last time: Linear Systems Part 2

- Iterative refinement
- Fixed-point stationary methods
  - Formulations for root-finding and linear systems
  - Iterative refinement as a stationary method
  - Iterative methods for large systems:
    - Jacobi, Gauss-Seidel, Successive Over-Relaxation
- Sherman-Morrison: Rank-1 updating
- Conjugate gradient (formulating a linear system as an optimization problem)
- Representing sparse systems

### Outline

- What is data modeling and why do it?
- Why choose a model that minimizes sum of squared error?
- How to formulate and compute the optimal leastsquares linear model?
- Illustrating least-squares with special cases: constant, line
- Weighted least squares
- How to judge the quality of a model?

## Data Modeling

- Given: data points, functional form, find constants in function
- Example: given (x<sub>i</sub>, y<sub>i</sub>), find line through them;
   i.e., find a and b in y = ax+b



## Data Modeling

- You might do this because you actually care about those numbers...
  - Example: measure position of falling object, fit parabola



 $p = -1/_2 gt^2$ 

 $\Rightarrow$  Estimate g from fit

### Data Modeling

 ... or because some aspect of behavior is unknown and you want to ignore it







Note: 1972 to 2006. Sample size: 41,795. Each circle represents an income range of \$2,000 (e.g., \$10,001 to \$12,000), in 2006\$. Its diameter is proportional to the number of people in that range.

Source: My calculations from General Social Survey data.

#### Historical context



### Which model is best?



#### Best-fit lines under different metrics



# Least Squares

- Nearly universal formulation of fitting: minimize squares of differences between data and function
  - Example: for fitting a line, minimize

$$\chi^2 = \sum_i \left( y_i - (ax_i + b) \right)^2$$

with respect to a and b

- Finds one unique best-fit model for a dataset

### Linear Least Squares

- (Also called "Ordinary least squares"
- General pattern:

$$y_i = a f(\vec{x}_i) + b g(\vec{x}_i) + c h(\vec{x}_i) + \cdots$$
  
Given  $(\vec{x}_i, y_i)$ , solve for  $a, b, c, \dots$ 

 Dependence on unknowns (a, b, c...) is linear, but f, g, etc. might not be!

## Linear Least Squares Examples

- General form:  $y_i = a f(\vec{x}_i) + b g(\vec{x}_i) + c h(\vec{x}_i) + \cdots$ Given  $(\vec{x}_i, y_i)$ , solve for  $a, b, c, \dots$
- Linear regression:

$$y_{i} = a * x_{i} + b$$

Multiple linear regression: x has many dimensions

• **Polynomial** regression:

$$y_i = a * x_i^2 + b * x_i + c$$

How do we compute the model parameters?

Solving Linear Least Squares Problem (one simple approach)

• Take partial derivatives:

$$\chi^{2} = \sum_{i} (y_{i} - a f(x_{i}) - b g(x_{i}) - \cdots)^{2}$$

$$\frac{\partial}{\partial a} = \sum_{i} -2f(x_i)(y_i - a f(x_i) - b g(x_i) - \cdots) = 0$$
$$a \sum_{i} f(x_i)f(x_i) + b \sum_{i} f(x_i)g(x_i) + \cdots = \sum_{i} f(x_i)y_i$$

$$\frac{\partial}{\partial b} = \sum_{i} -2g(x_i)(y_i - a f(x_i) - b g(x_i) - \cdots) = 0$$
$$a \sum_{i} g(x_i)f(x_i) + b \sum_{i} g(x_i)g(x_i) + \cdots = \sum_{i} g(x_i)y_i$$

## Solving Linear Least Squares Problem

• For convenience, rewrite as matrix:

$$\begin{bmatrix} \sum_{i} f(x_{i})f(x_{i}) & \sum_{i} f(x_{i})g(x_{i}) & \cdots \\ \sum_{i} g(x_{i})f(x_{i}) & \sum_{i} g(x_{i})g(x_{i}) & \cdots \\ \vdots & \vdots & & \end{bmatrix} \begin{bmatrix} a \\ b \\ \vdots \end{bmatrix} = \begin{bmatrix} \sum_{i} f(x_{i})y_{i} \\ \sum_{i} g(x_{i})y_{i} \\ \vdots \end{bmatrix}$$

• Factor:

$$\sum_{i} \begin{bmatrix} f(x_i) \\ g(x_i) \\ \vdots \end{bmatrix} \begin{bmatrix} f(x_i) \\ g(x_i) \\ \vdots \end{bmatrix}^{\mathrm{T}} \begin{bmatrix} a \\ b \\ \vdots \end{bmatrix} = \sum_{i} y_i \begin{bmatrix} f(x_i) \\ g(x_i) \\ \vdots \end{bmatrix}$$

Alternative perspective: [Approximate] linear system

 There's a different derivation of this: overconstrained linear system

$$\mathbf{A}x = b$$
$$\begin{pmatrix} x \\ x \end{pmatrix} = \begin{pmatrix} b \\ b \end{pmatrix}$$

Notation change:

- A is now basis functions computed on observations (f(x<sub>i</sub>), g(x<sub>i</sub>), ...)
  x is now model parameters (a, b, c...)
  b is now "y"
- A has n rows and m<n columns: more equations than unknowns

Geometric Interpretation for Over-determined System

• Find the x that comes "closest" to satisfying Ax=b– i.e., minimize b-Ax b r=b-Ax $\theta$  y=Ax

### Geometric Interpretation

- Interpretation: find x that comes "closest" to satisfying Ax=b
  - i.e., minimize b-Ax
  - i.e., minimize  $\parallel$  b–Ax  $\parallel$



Equivalently, find x such that r is orthogonal to span(A)

$$0 = \mathbf{A}^{\mathrm{T}}\mathbf{r} = \mathbf{A}^{\mathrm{T}}(\mathbf{b} - \mathbf{A}\mathbf{x})$$
$$\mathbf{A}^{\mathrm{T}}\mathbf{A}\mathbf{x} = \mathbf{A}^{\mathrm{T}}\mathbf{b}$$

## Forming the equation

- What are A and b?
  - Row i of A is basis functions computed on x<sub>i</sub>
  - Row i of b is y<sub>i</sub>

$$\mathbf{A} = \begin{bmatrix} f(x_1) & g(x_1) & \cdots \\ f(x_2) & g(x_2) & \cdots \\ \vdots & \vdots & \end{bmatrix}, \quad b = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ \vdots \end{bmatrix}$$
$$\mathbf{A}^{\mathsf{T}} \mathbf{A} = \begin{bmatrix} \sum_i f(x_i) f(x_i) & \sum_i f(x_i) g(x_i) & \cdots \\ \sum_i g(x_i) f(x_i) & \sum_i g(x_i) g(x_i) & \cdots \\ \vdots & \vdots & \vdots & \end{bmatrix}, \quad \mathbf{A}^{\mathsf{T}} b = \begin{bmatrix} \sum_i y_i f(x_i) \\ \sum_i y_i g(x_i) \\ \vdots \\ \vdots & \vdots & \end{bmatrix}$$

Minimizing Sum of Squares = Finding Closest Ax in span(A)

Starting from goal of minimizing sum of squares

Starting from goal of finding Ax in span(A) closest to b outside span(A)

$$\sum_{i} \begin{bmatrix} f(x_i) \\ g(x_i) \\ \vdots \end{bmatrix} \begin{bmatrix} f(x_i) \\ g(x_i) \\ \vdots \end{bmatrix}^{\mathrm{T}} \begin{bmatrix} a \\ b \\ \vdots \end{bmatrix} = \sum_{i} y_i \begin{bmatrix} f(x_i) \\ g(x_i) \\ \vdots \end{bmatrix}$$

Great, but how do we solve it?

## 1: Normal Equations

Pseudocode:

for each  $x_i, y_i$ compute  $f(x_i), g(x_i),$  etc. store in column i of A store  $y_i$  in b compute  $A^TA, A^Tb$  **solve A^TAx=A^Tb**   $\sum_i \begin{bmatrix} f(x_i) \\ g(x_i) \\ \vdots \end{bmatrix} \begin{bmatrix} f(x_i) \\ g(x_i) \\ \vdots \end{bmatrix}^T$  $\begin{bmatrix} a \\ b \\ \vdots \end{bmatrix} = \sum_i y_i \begin{bmatrix} f(x_i) \\ g(x_i) \\ \vdots \end{bmatrix}$ 

 These can be inefficient, since A typically much larger than A<sup>T</sup>A and A<sup>T</sup>b

#### 2: More efficient normal equations

for each x<sub>i</sub>,y<sub>i</sub> compute f(x<sub>i</sub>), g(x<sub>i</sub>), etc. accumulate outer product in U accumulate product with y<sub>i</sub> in v solve Ux=v

### 3: Using the pseudoinverse

$$\min (b - \mathbf{A}x)^{\mathrm{T}} (b - \mathbf{A}x)$$
$$\nabla ((b - \mathbf{A}x)^{\mathrm{T}} (b - \mathbf{A}x)) = -2\mathbf{A}^{\mathrm{T}} (b - \mathbf{A}x) = \vec{0}$$
$$\mathbf{A}^{\mathrm{T}} \mathbf{A}x = \mathbf{A}^{\mathrm{T}} b$$

for each  $x_i, y_i$ compute  $f(x_i), g(x_i),$  etc. store in row i of A store  $y_i$  in b compute  $x = (A^TA)^{-1} A^Tb$ 

(A<sup>T</sup>A)<sup>-1</sup> A<sup>T</sup> is known as "pseudoinverse" of A

## The Problem with Normal Equations

- Involves solving A<sup>T</sup>Ax=A<sup>T</sup>b
- This can be inaccurate
  - Independent of solution method

- Remember: 
$$\frac{\|\Delta x\|}{\|x\|} \le cond(A) \frac{\|\Delta A\|}{\|A\|}$$
$$- \operatorname{cond}(A^{\mathsf{T}}A) = [\operatorname{cond}(A)]^2$$

- Next week: computing pseudoinverse
  - More expensive, but more accurate
  - Also allows diagnosing insufficient data

Special Cases

### Special Case: Constant

Let's try to model a function of the form

y = a

#### Special Case: Constant

Let's try to model a function of the form
 y = a

$$y_i = a f(\vec{x}_i) + b g(\vec{x}_i) + c h(\vec{x}_i) + \cdots$$

In this case, f(x<sub>i</sub>)=1 and we are solving

$$\sum_{i} [1] [a] = \sum_{i} [y_{i}]$$
$$\therefore a = \frac{\sum_{i} y_{i}}{n}$$

# Special Case: Line

• Fit to y=a+bx

• 
$$f(x_i)=1$$
,  $g(x_i)=x$ . So, solve:  

$$\sum_{i} \begin{bmatrix} 1 \\ x_i \end{bmatrix} \begin{bmatrix} 1 & x_i \end{bmatrix} \begin{bmatrix} a \\ b \end{bmatrix} = \sum_{i} y_i \begin{bmatrix} 1 \\ x_i \end{bmatrix}$$

$$(\mathbf{A}^{\mathrm{T}}\mathbf{A})^{-1} = \begin{bmatrix} n & \Sigma x_i \\ \Sigma x_i & \Sigma x_i^2 \end{bmatrix}^{-1} = \frac{\begin{bmatrix} \Sigma x_i^2 & -\Sigma x_i \\ -\Sigma x_i & n \end{bmatrix}}{n\Sigma x_i^2 - (\Sigma x_i)^2}, \quad \mathbf{A}^{\mathrm{T}}b = \begin{bmatrix} \Sigma y_i \\ \Sigma x_i y_i \end{bmatrix}$$

$$a = \frac{\Sigma x_i^2 \Sigma y_i - \Sigma x_i \Sigma x_i y_i}{n\Sigma x_i^2 - (\Sigma x_i)^2}, \quad b = \frac{n\Sigma x_i y_i - \Sigma x_i \Sigma y_i}{n\Sigma x_i^2 - (\Sigma x_i)^2}$$

Variant: Weighted Least Squares

# Weighted Least Squares

- Common case: the (x<sub>i</sub>,y<sub>i</sub>) have different uncertainties associated with them
- Want to give more weight to measurements of which you are more certain
- Weighted least squares minimization  $\min \chi^2 = \sum_i w_i (y_i - f(x_i))^2$
- If "uncertainty" (stdev) is  $\sigma$ , best to take  $W_i = \frac{1}{\sigma_i^2}$

### Weighted Least Squares

• Define weight matrix W as



Then solve weighted least squares via

$$\mathbf{A}^{\mathrm{T}}\mathbf{W}\mathbf{A} x = \mathbf{A}^{\mathrm{T}}\mathbf{W} b$$

Understanding Error and Uncertainty

# Error Estimates from Linear Least Squares

- For many applications, finding model is useless without estimate of its accuracy
- Residual is b Ax
- Can compute  $\chi^2 = (b Ax) \cdot (b Ax)$
- How do we tell whether answer is good?
  - Lots of measurements
  - $-\chi^2$  is small
  - $\chi^2$  increases quickly with perturbations to x ( $\rightarrow$  standard variance of estimate is small)
  - R<sup>2</sup> ("coefficient of determination") is near 1

## Error Estimates from Linear Least Squares

- C=(A<sup>T</sup>A)<sup>-1</sup> is called *covariance* of the data
- The "standard variance" in our estimate of x is  $2 \chi^2$

$$\sigma^2 = \frac{\chi}{n-m} \mathbf{C}$$

- This is a matrix:
  - Diagonal entries give variance of estimates of components of x: e.g., var(a<sub>0</sub>)
  - Off-diagonal entries explain mutual dependence:
     e.g., cov(a<sub>0</sub>, a<sub>1</sub>)
- n–m is (# of samples) minus (# of degrees of freedom in the fit): consult a statistician...



#### Coefficient of Determination

$$R^2 \equiv 1 - \frac{\chi^2}{\sum_i y_i - \overline{y}}$$

R<sup>2</sup> : Proportion of observed variability that is explained by the model

e.g.,  $R^2 = 0.7$  means 70% variability explained For linear regression,  $R^2$  is Pearson's correlation.



## Keep in mind...

- In general, uncertainty in estimated parameters goes down slowly: like 1/sqrt(# samples)
- Formulas for special cases (like fitting a line) are messy: simpler to think of A<sup>T</sup>Ax=A<sup>T</sup>b form
- Normal equations method often not numerically stable: orthogonal decomposition methods used instead
- Linear least squares is not always the most appropriate modeling technique...

#### Next time

- Non-linear models
  - Including logistic regression
- Dealing with outliers and bad data
- Practical considerations
  - Is least squares an appropriate method for my data?
- Examples with Excel and Matlab