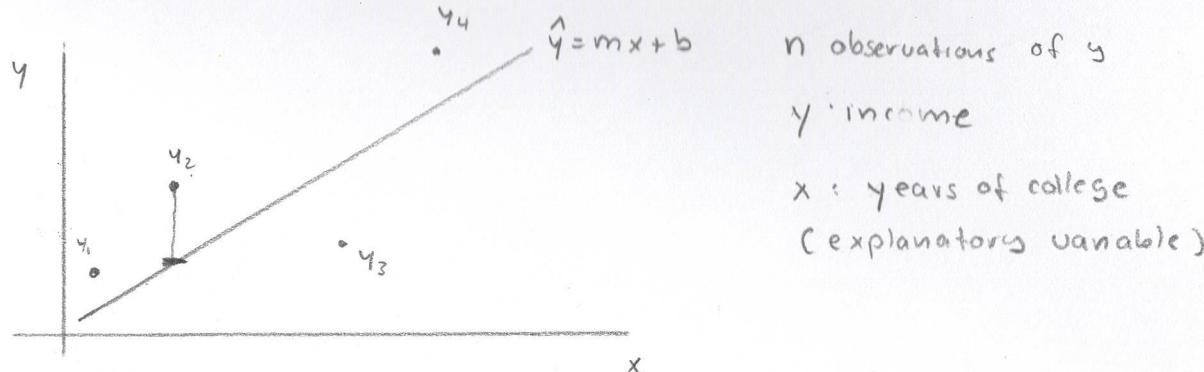


Linear Regression



- Q: (1) degree of correlation?
 (2) best (linear) fit?

SSE: sum squared error

$$\sum (\hat{y}_i - y_i)^2 = \|\vec{\hat{y}} - \vec{y}\|_{\text{norm}}^2$$

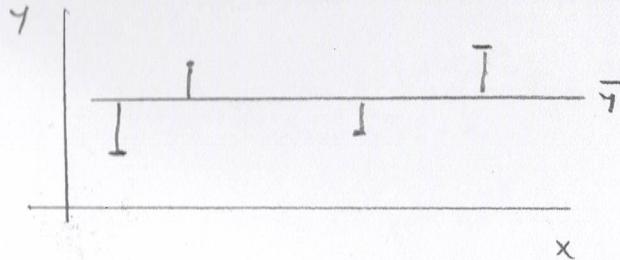
$$\|\vec{x}\| = \sqrt{x_1^2 + x_2^2 + \dots + x_n^2}$$

$\text{SSE} = 0 \Rightarrow \hat{y}_i = y_i \text{ for all } i \Rightarrow \text{perfect fit}$

SST: sum squared total "error if using mean as predictor"

$$\sum (\bar{y}_i - \bar{y})^2$$

mean of all observations



$SST = 0 \Rightarrow$ all values equal mean

$$R^2 = 1 - \frac{SSE}{SST}$$

correlation coefficient

$SSE = 0 \Rightarrow$ perfect fit $R^2 = 1$

$SSE = SST \Rightarrow$ no correlation $R^2 = 0$

"no linear model is better than mean"

What is the best linear model?

$$\hat{y} = mx + b \quad \text{Predictor model}$$

$$\left. \begin{array}{l} \hat{y}_1 = mx_1 + b \\ \hat{y}_2 = mx_2 + b \\ \vdots \\ \hat{y}_n = mx_n + b \end{array} \right\} \quad \left[\begin{array}{c} \hat{y}_1 \\ \hat{y}_2 \\ \vdots \\ \hat{y}_n \end{array} \right] = \left[\begin{array}{cc|c} 1 & x_1 & b \\ 1 & x_2 & m \\ \vdots & \vdots & \\ 1 & x_n & \end{array} \right]$$

$$\hat{y} = A\beta$$

"best fit" \Rightarrow minimize SSE

$$\text{minimize } \|\hat{y} - y\|^2$$

$$\text{minimize } \|A\beta - y\|^2 \Rightarrow \text{find } \beta \text{ that minimizes SSE}$$

steps $\frac{\partial \text{SSE}}{\partial \beta} = 0 \Rightarrow \text{solve for } \beta$

$$SSE = \|A\beta - y\|^2$$

$$= (A\beta - y)^T (A\beta - y)$$

$$= (\beta^T A^T - y^T) (A\beta - y)$$

$$= \beta^T A^T A\beta - y^T A\beta - \underbrace{\beta^T A^T y}_{(A\beta)^T y} + \underbrace{y^T y}_{y^T A\beta^T} \text{ transpose}$$

$$= \beta^T A^T A\beta - 2y^T A\beta + y^T y$$

$$\frac{\partial SSE}{\partial \beta} = 2$$

derivative rules

$$\frac{\partial \vec{y}}{\partial \vec{x}} \stackrel{\text{def}}{=} \begin{bmatrix} \frac{\partial y_1}{\partial x_1} & \cdots & \frac{\partial y_n}{\partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial y_1}{\partial x_m} & \cdots & \frac{\partial y_n}{\partial x_m} \end{bmatrix}$$

useful identities

$$\|x^2\| = \sum x_i^2 = x^T x$$

$$\|x\| = \sqrt{\sum x_i^2}$$

$$\begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix} [x_1 \dots x_n] = x_1^2 + \dots + x_n^2$$

$$(A^T)^T = A$$

$$(AB)^T = B^T A$$

$$(A+B)^T = A^T + B^T$$

$$SSE = \beta^T A^T A \beta - 2 y^T A \beta + y^T y$$

$$\frac{\partial SSE}{\partial \beta} = 2 A^T A \beta - 2 A^T y + 0 = 0$$

$$A^T A \beta = A^T y$$

$$(A^T A)^{-1} (A^T A) \beta = (A^T A)^{-1} A^T y$$

$$\boxed{\beta = (A^T A)^{-1} A^T y}$$

note: $A^T A$ is always invertible if column vectors are linearly independent

Common matrix derivatives

$$\frac{\partial \vec{y}}{\partial \vec{x}} = \frac{\partial \vec{y}}{\partial \vec{x}}$$

$$A \vec{x} \quad | \quad A$$

$$\vec{x} A \quad | \quad A^T$$

$$\vec{x}^T \vec{x} \quad | \quad 2x$$

$$\vec{x}^T A^T A \vec{x} \quad | \quad 2A^T A x$$

$$\vec{y}^T A \vec{x} \quad | \quad A^T y$$

|

|

Linear regression in practice

X	Y
1	-0.2
2	2.2
3	3.5
.	.
.	.
.	.
8	12

n explanatory
variable

Depressed
variable

Column of 1 Column of x_i Vector of y

\downarrow \swarrow

$$A = \begin{bmatrix} 1 & -0.2 \\ 1 & 2.2 \\ \vdots & \vdots \\ 1 & 8 \end{bmatrix} \quad y = \begin{bmatrix} -0.2 \\ 2.2 \\ \vdots \\ 12 \end{bmatrix}$$

$$\beta = \text{inv}(A^T A) A^T y = \begin{bmatrix} m \\ b \end{bmatrix}$$

↑

defined
in any language

may need some
Linear Algebra package

$$SSE = \text{norm}(Ax - y)^2$$

$$SST = \text{norm}(y - \text{mean}(y))^2$$

$$R^2 = 1 - \frac{SSE}{SST}$$

may need Statistics package

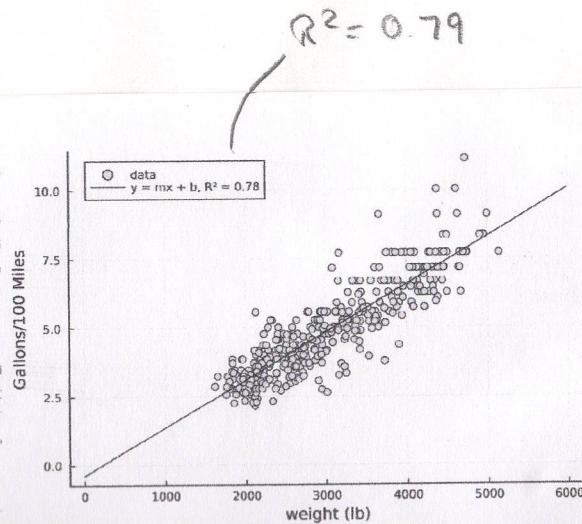
Multilinear Regression

base case $y = \frac{1}{mpg}$ (Gallons per 100 miles)

$x = \text{weight}$

$$y = mx + b$$

	mpg	cylinders	displacement	hp	weight	acceleration	year	mc
1	18.0	8	307.0	130.0	3604.0	12.0	70	"chevrolet che
2	15.0	8	350.0	165.0	3693.0	11.5	70	"buick skylar
3	18.0	8	318.0	150.0	3436.0	11.0	70	"plymouth sat
4	16.0	8	304.0	150.0	3433.0	12.0	70	"amc rebel ss
5	17.0	8	302.0	140.0	3449.0	10.5	70	"ford torino"
6	15.0	8	429.0	198.0	4341.0	10.0	70	"ford galaxie
7	14.0	8	454.0	220.0	4354.0	9.0	70	"chevrolet imp
8	14.0	8	440.0	215.0	4312.0	8.5	70	"plymouth fur
9	14.0	8	455.0	225.0	4429.0	10.0	70	"Pontiac catal
10	15.0	8	390.0	190.0	3850.0	8.5	70	"amc ambassad
... more								
392	31.0	4	119.0	82.0	2720.0	19.4	82	"chevy s-10"



"Cars" Dataset

Multivariable model

$$y = c_1 x_1 + c_2 x_2 + c_3 x_3 + c_4 \leftarrow \text{"intercept"}$$

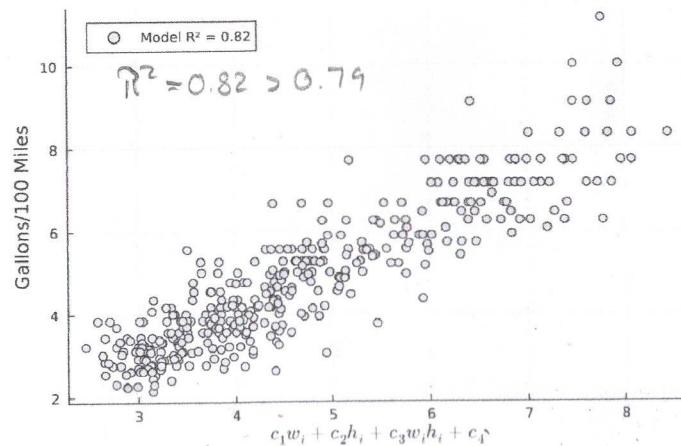
↑ weight ↑
 horse power weight × horse power

$$\hat{Y} = \hat{A} \hat{\beta}$$

$$\begin{bmatrix} Y_1 \\ \vdots \\ Y_n \end{bmatrix} = \begin{bmatrix} X_{11} & X_{21} & X_{31} & \dots \\ \vdots & \vdots & \vdots & \vdots \\ X_{1n} & X_{2n} & X_{3n} & \dots \end{bmatrix} \begin{bmatrix} c_1 \\ c_2 \\ c_3 \\ c_4 \end{bmatrix}$$

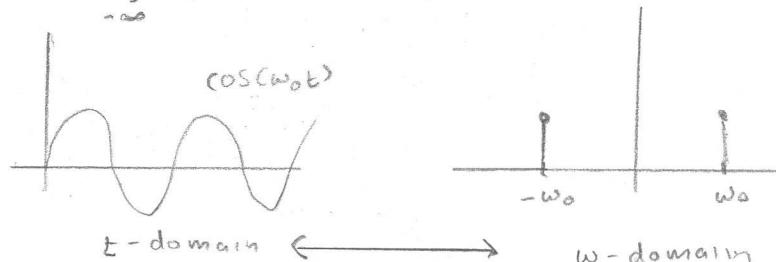
$$\hat{\beta} = (\hat{A}^T \hat{A})^{-1} \hat{A}^T \hat{y}$$

Adding variables increases explanatory power



Integral Equations (Domain Transformations)

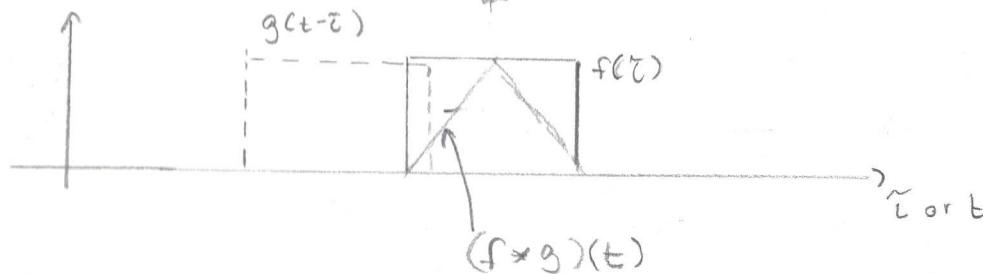
Fourier Transform $F(\omega) = \int_{-\infty}^{\infty} \exp(-i\omega t) f(t) dt$



Laplace Transform $F(s) = \int_{-\infty}^{\infty} \exp(-st) f(t) dt$ s is $a+bi$ (complex number)

t-domain \longleftrightarrow s = $a + bi$

Convolution $(f * g)(t) = \int_{-\infty}^{\infty} g(t-\tau) f(\tau) d\tau$



General Integral Equation

$$y(x) = \int_a^x k(x,t) f(t) dt$$

limits \rightarrow ↑
can differ from
 $\pm \infty$

Kernel function

$k(x,t) = g(x-t) \Rightarrow$ convolution

$k(x,t) = \exp(-xt) \Rightarrow$ Laplace

$k(x,t) = \exp(-ixt) \Rightarrow$ Fourier

⋮

t -domain \longleftrightarrow x -domain

Applications / Problem

- Measurement is in the x -domain \Rightarrow quantity of interest is in t -domain
 - surface temperature evolution
 - size distribution function
 - Optical measurements / remote sensing
 - FTIR
 - ⋮
- Inversion \Rightarrow find $f(t)$ from discrete $y(x)$

Example Inversion Problem

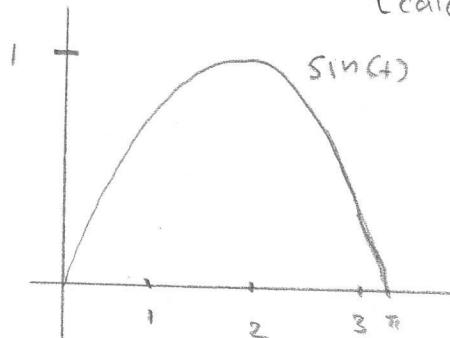
$$y(x) = \int_a^b k(x, t) u(t) dt \quad \text{"Fredholm Integral Eq."}$$

$$k(x, t) = \exp(x \cdot \cos(t)) \quad u(t) = \sin(t) \quad y(x) = \frac{2 \sinh(x)}{x}$$

$$a=0, b=\pi \quad t \in [0, \pi] \quad \text{and} \quad x \in [0, \pi/2]$$

test : $\int_0^\pi \exp(x \cos(t)) \sin(t) dt = \frac{2 \sinh}{x}$

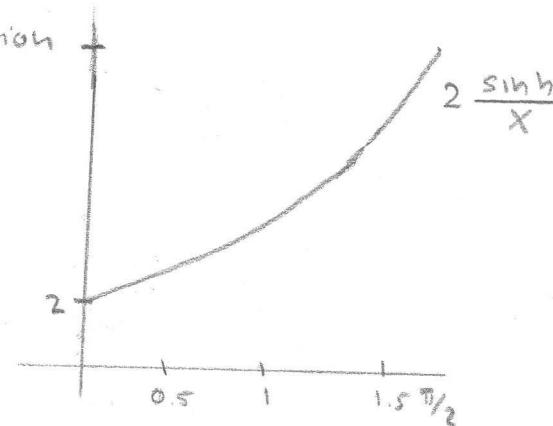
→ Verify using integration
(calculus)



\leftarrow t -domain



x -domain



Discretized Integral Equation

$$\int_a^b k(x, t) u(t) dt = y(x)$$

$$A \vec{u} = \vec{y}$$

$$\vec{u} = [u_1, \dots, u_n]$$

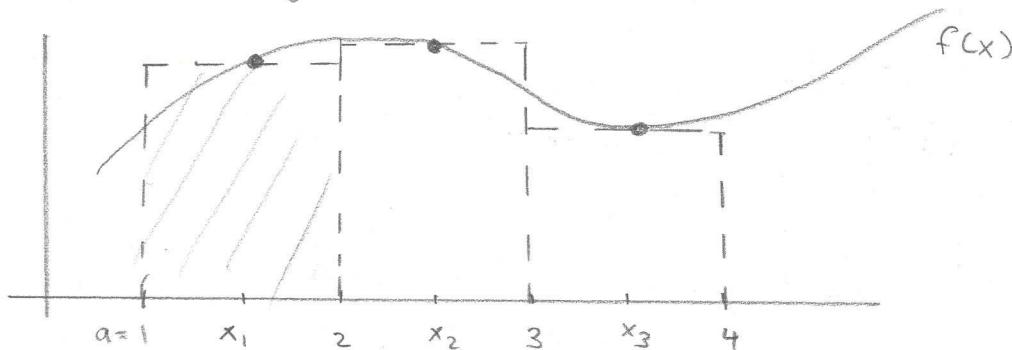
$$\vec{y} = [y_1, \dots, y_m]$$

$$A \in \mathbb{R}^{m \times n}$$

Question: how to find A ? (next page)

$$\sum_{i=1}^n w_i k(x_j, t_i) u(t_i) \approx y_j$$

Numerical Integration



$n=3$: number of data points

$a=1$, $b=4$: limits of integration

$$w = \frac{b-a}{n} \quad \text{bin width (here } w=1\text{)}$$
$$\int_a^b f(x) dx = \sum_{i=1}^n f(x_i) w \quad \left. \right\} \text{quadrature algorithm}$$

$$x_i = (i - \frac{1}{2}) w + a$$

$$x_1 = (1 - 0.5)1 + 1 = 1.5$$

$$x_2 = (2 - 0.5)1 + 1 = 2.5$$

$$x_3 = (3 - 0.5)1 + 1 = 3.5$$

$$\int_a^b k(x, t) u(t) dt = y(x) \quad t \in [a, b] \quad x \in [c, d]$$

u is discretized into n points $i = 1 \dots n$ $w_n = \frac{b-a}{n}$

y is discretized into m points $j = 1 \dots m$ $w_m = \frac{d-c}{m}$

$$y_1 = \int_a^b k(x_1, t) u(t) dt$$

$$y_1 = k(x_1, t_1) u(t_1) w_n + k(x_1, t_2) u(t_2) w_n + \dots + k(x_1, t_n) u(t_n) w_n$$

$$y_m = k(x_m, t_1) u(t_1) w_n + k(x_m, t_2) u(t_2) w_n + \dots + k(x_m, t_n) u(t_n) w_n$$

$$\begin{bmatrix} y_1 \\ \vdots \\ y_m \end{bmatrix} = \begin{bmatrix} k(x_1, t_1) w_n & \dots & k(x_1, t_n) w_n \\ \vdots & \ddots & \vdots \\ k(x_m, t_1) w_n & \dots & k(x_m, t_n) w_n \end{bmatrix} \begin{bmatrix} u_1 \\ \vdots \\ u_n \end{bmatrix}$$

$a_{j,i} = w_n k(x_j, t_i)$

$$A = \mathbb{R}^{m \times n}$$

$$y = A u$$

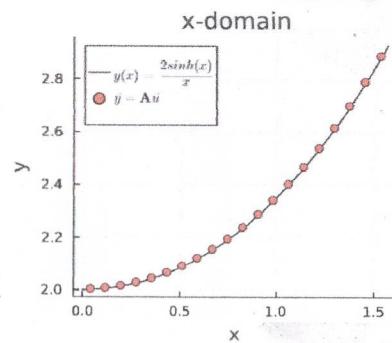
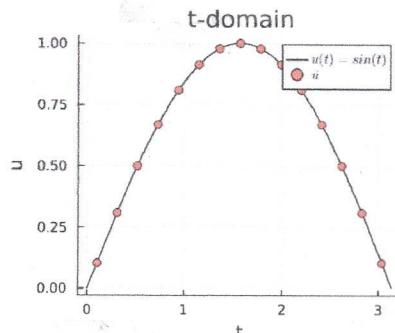
$$\mathbb{R}^m$$

$$u = \mathbb{R}^n$$

Example Discretization

→ n, m can be independently varied

→ over-determined and under-determined transforms



$n = 15$

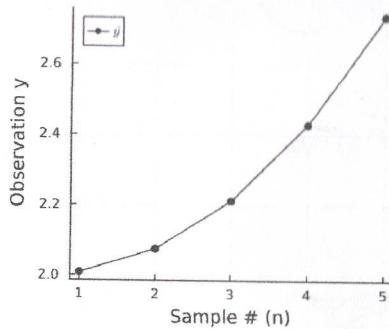
$m = 20$

Example Inversion

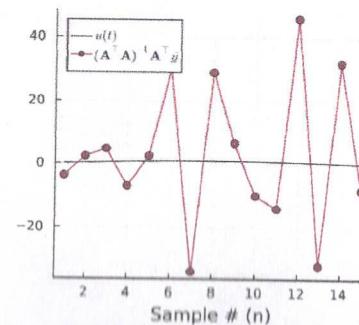
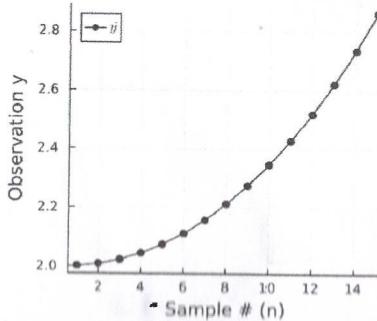
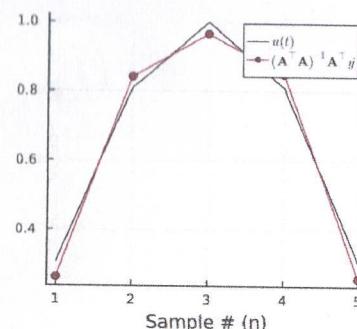
$$y = Au(t)$$

SSE: minimize $\|Au - y\|^2$
 ↑
 find \hat{u}
 that is best

Observation in
X-domain



Prediction in
t-domain



works well for small n

fails for large n

why?

III posed Problems

- Lack of a unique solution
- Sensitivity to initial conditions (noise) \rightarrow weather prediction
- nonexistence of a solution

$$A = \begin{bmatrix} 1 & 2 \\ 0 & 0 \end{bmatrix} \quad \det A = 0 \quad \left. \begin{array}{l} \text{- underdetermined} \\ \text{- no unique solution} \\ A^{-1} \text{ not defined} \end{array} \right\}$$
$$A = \begin{bmatrix} 1 & 2 \\ 2 & 4 \end{bmatrix} \quad \det A = 0$$

$\text{rank}(A) = \text{maximum number of linearly independent columns}$

here $\text{rank}(A) = 1$

Full rank = 2 (max possible)

A is "rank deficient"

Moore-Penrose Inverse

$$A^+ = (A^T A)^{-1} A^T$$

- (1) for $m=n$ $A^+ = A^{-1}$ if A is full rank
- (2) $m > n$:
 - over determined system (e.g. linear regression case)
 - more eq. than unknowns
 - minimizes $\|A u - y\|^2$
- (3) $n > m$:
 - under determined system
 - fewer eq. than unknowns
 - infinite solutions
- if A is not full rank A^+ must be computed using SVD
- use `pinv(A)` in most languages

Singular Value Decomposition

Fall SVD

$$A = U \Sigma V^T$$

singular values

$$M = \begin{pmatrix} 0 & - & - & - & - & - & - \\ 0 & 0 & - & - & - & - & - \\ 0 & 0 & 0 & - & - & - & - \\ 0 & 0 & 0 & 0 & - & - & - \\ 0 & 0 & 0 & 0 & 0 & - & - \\ 0 & 0 & 0 & 0 & 0 & 0 & - \end{pmatrix}$$

$$\sigma_1 > \sigma_2 \dots > \sigma_p$$

- can find using numerical methods
 - $O(n^2)$ \Rightarrow expensive
 - use $\text{svd}(A)$ from linear algebra

Moore - Penrose Inverse

$$A^+ = V \Sigma^+ U^\top$$

$$\Sigma^+ = \begin{bmatrix} \frac{1}{\sqrt{6}} & 0 & - & - & - & - & 0 \\ 0 & \ddots & & & & & \\ 0 & & \ddots & & & & \\ \vdots & & & \ddots & & & \\ \vdots & & & & \ddots & & \\ 0 & 0 & - & - & - & - & \frac{1}{\sqrt{6}} & 0 \end{bmatrix}$$

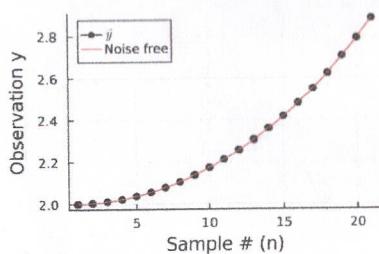
- numerical method to find A^+
- if $\sigma_i = 0$ then a zero is placed
- if A is rank deficient one or more $\sigma_i = 0$

Inversion Revisited

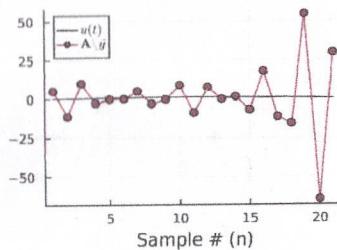
$$\frac{1}{10^{-15}} = 10^{15} \Rightarrow \text{large entries in } \Sigma^+$$

"infected by numerical or measurement noise."

x-domain

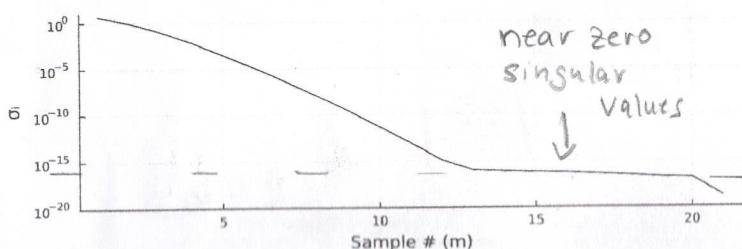


t-domain



$$u = (A^T A)^{-1} A^T y$$

$$u = A^+ y$$



near zero
singular
values

↑ log scale
↓

— machine precision
for 64 bit floating point
 $\approx 10^{-16}$

No new information

Regularized Inverse (Noise Filtering)

$$u_R = \underset{\text{Ordinary regression}}{\text{minimize}} \left\{ \|Au - y\|^2 + \lambda^2 \|L(u - u_0)\|^2 \right\}$$

filtermatrix
 initial guess
 regularization parameter

$$u_R = (A^T A + \lambda^2 L^T L)^{-1} (A^T y + \lambda^2 L^T L u_0)$$

$$u_0 = A^+ y$$

$$\lim_{\lambda \rightarrow \infty} u_R = u_0$$

Starting point for $L = \overline{I}$

\overline{I}
identity matrix

$u_0 = 0$
(no initial guess)

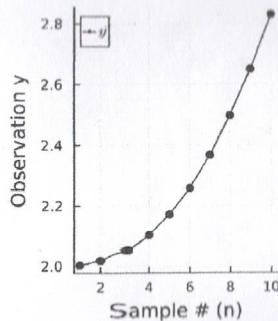
then

$$u_R = (A^T A + \lambda^2 \overline{I})^{-1} A^T y \quad \text{"Ridge Regression"}$$

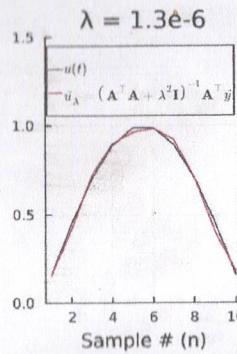
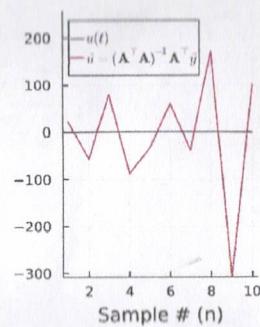
Inversion Revisited

- need to try range of λ
- $\lambda = 0 \Rightarrow$ bad solution
- λ large \Rightarrow approach zero \Rightarrow bad solution
- a λ exists to find good inversion
 \Rightarrow how to find λ ?

x-domain



t-domain



$$u_L = (A^T A + \lambda^2 I)^{-1} A^T y$$

$\lambda = 0$

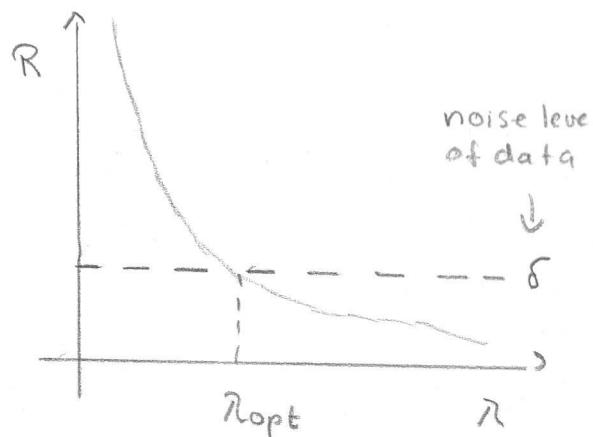
Methods to find λ

$$u_\lambda = \text{minimize} \left\{ \underbrace{\|Au - y\|^2}_{R = \text{Residual Norm}} + \lambda^2 \underbrace{\|L(u - u_0)\|^2}_{S = \text{Solution Norm}} \right\}$$

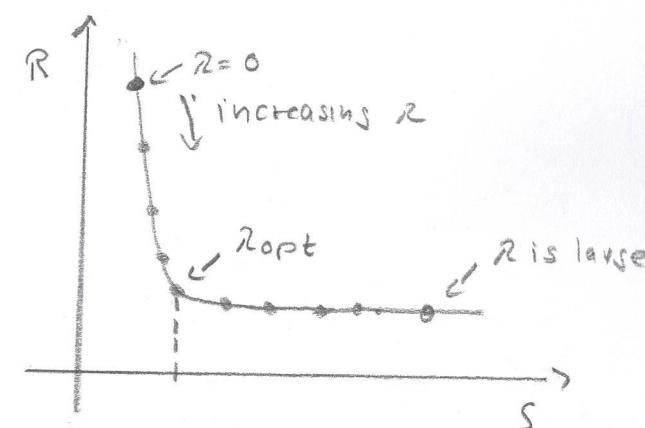
$R = \text{Residual Norm}$

$S = \text{Solution Norm}$

Morozev



L-curve



- λ_{opt} is "corner" at L-curve
- find visually or algorithm