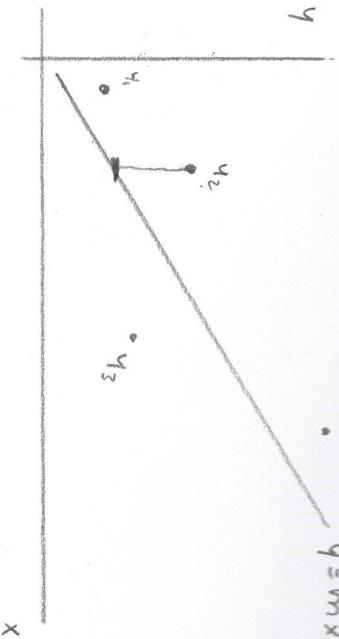


Linear Regression

$$y \quad \hat{y} = mx + b$$

y : income
 x : years of college (explanatory variable)



Q: (1) degree of correlation?

(2) best (linear) fit?

SSE: sum squared error

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

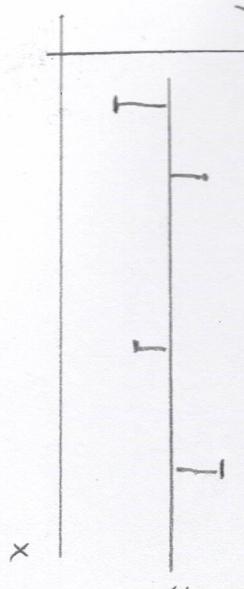
$SSE = 0 \Rightarrow \hat{y}_i = y_i \text{ for all } i \Rightarrow \text{perfect fit}$

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

$$\sum_{i=1}^n (\hat{y}_i - \bar{y})^2$$

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?

$$\hat{y} = mx + b$$

Predictor model

$$\begin{aligned}\hat{y}_1 &= m x_1 + b \\ \hat{y}_2 &= m x_2 + b \\ &\vdots \\ \hat{y}_n &= m x_n + b\end{aligned}\quad \left[\begin{array}{c} \hat{y}_1 \\ \hat{y}_2 \\ \vdots \\ \hat{y}_n \end{array} \right] = \left[\begin{array}{c} x_1 \\ x_2 \\ \vdots \\ x_n \end{array} \right] \left[\begin{array}{c} b \\ m \end{array} \right]$$

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

"best fit" \Rightarrow minimize SSE

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

minimize $\|A\beta - y\|^2 \Rightarrow$ find β that minimizes SSE

steps

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

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

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

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

$$= \beta^T A^T A\beta - y^T A\beta - \underbrace{\beta^T A^T y + y^T y}$$

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

$y^T A\beta$) transpose

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

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

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

derivative rules

$$\frac{\partial y}{\partial x} = \begin{bmatrix} \frac{\partial y_1}{\partial x_1} & \dots & \frac{\partial y_n}{\partial x_n} \\ \frac{\partial y_1}{\partial x_m} & \dots & \frac{\partial y_n}{\partial x_m} \end{bmatrix}$$

useful identities

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

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

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

$$\frac{\partial SSE}{\partial \beta_3} = 2$$

derivative rules

Common matrix derivations

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

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

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

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

$$\begin{array}{c|c} A^T x & A \\ \hline x^T A & A^T \\ \hline x^T x & 2x \\ \hline x^T A^T A x & 2 A^T A x \\ \hline y^T A x & A^T y \end{array}$$

Note: $A^T A$ is always invertible if column

vectors are linearly independent

Linear regression in practice

X	y
1	-0.2
2	2.2
3	3.5
4	4.8
5	5.1
6	6.4
7	7.7
8	8.0
12	12

Column of
vector of y

\downarrow

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

$y = \begin{bmatrix} -0.2 \\ 2.2 \\ \vdots \\ 12 \end{bmatrix}$

Column of
vector of x_i

\uparrow

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

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

$$\text{base case } y = \frac{1}{\text{mpg}}$$

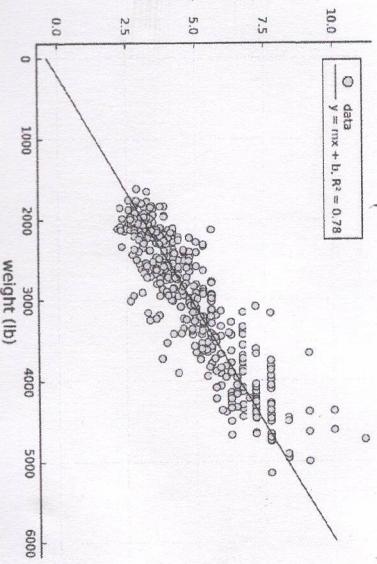
$x = \text{weight}$

$$y = mx + b$$

(Gallons per 100 miles)

$$R^2 = 0.79$$

mpg	cylinders	displacement	hp	weight	acceleration	year	mc
1	18.0	8	367.0	130.0	3664.0	12.0	70
2	15.0	8	350.0	165.0	3693.0	11.5	70
3	18.0	8	318.0	150.0	3435.0	11.0	70
4	16.0	8	394.0	158.0	3435.0	12.0	70
5	17.0	8	392.0	140.0	3443.0	10.5	70
6	15.0	8	429.0	198.0	4341.0	10.0	70
7	14.0	8	454.0	220.0	4356.0	9.0	70
8	14.0	8	440.0	215.0	4312.0	8.5	70
9	14.0	8	415.0	225.0	4428.0	10.0	70
10	15.0	8	390.0	190.0	3850.0	8.5	70
⋮ More							
392	31.0	4	119.0	82.0	2720.0	19.4	82
							"chevy s-10"



"Cars" Dataset

Multivariable model

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

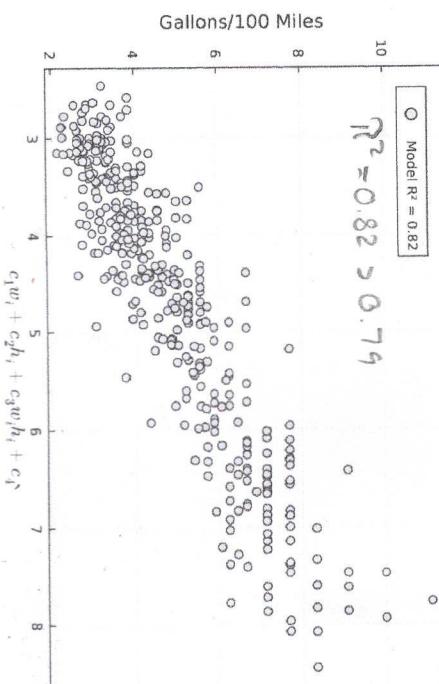
↑
weight
horse power

$$\text{weight} \times \text{horse power}$$

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

$$\begin{bmatrix} \hat{y}_1 \\ \vdots \\ \hat{y}_n \end{bmatrix} = \begin{bmatrix} x_{11} & x_{21} & x_{31} & 1 \\ \vdots & \vdots & \vdots & \vdots \\ x_{1n} & x_{2n} & x_{3n} & 1 \end{bmatrix} \begin{bmatrix} c_1 \\ c_2 \\ c_3 \\ c_4 \end{bmatrix}$$

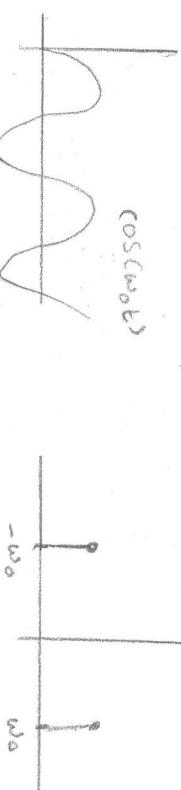
Adding variables increases explanatory power



Integral Equations (Domain Transformations)

Fourier Transform

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



\longleftrightarrow t - domain \longleftrightarrow ω - domain

Laplace transform

$$F(s) = \int_{-\infty}^{\infty} \exp(-st) f(t) dt$$

s is a + bi (complex number)

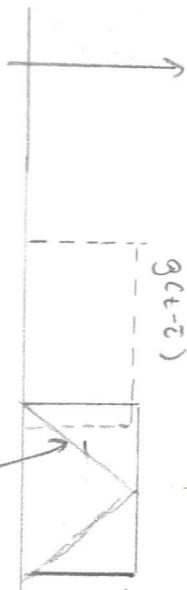
\longleftrightarrow t - domain \longleftrightarrow s - domain

Convolution

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

(2) $f * g$

f



$\text{or } g$

(7) $(g * f)$

General Integral Equation

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

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

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

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

↑
Kernel function

limits
can differ from
 $\pm \infty$

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
 - FIR
 - ⋮
- Inversion \Rightarrow find $f(\cdot)$ from discrete $y(x)$

Example Inversion Problem

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

"Fredholm Integral Eq."

$$K(x,t) = \exp(x \cdot \cos(t))$$

$$u(t) = \sin(t)$$

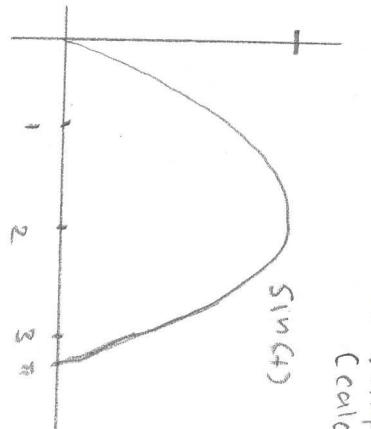
$$y(x) = \frac{2 \sinh(x)}{x}$$

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

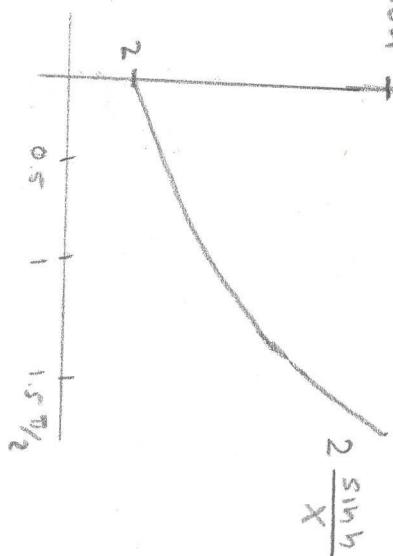
$$\text{test : } \int_0^{\pi} \exp(x \cos(t)) \sin(t) dt = \frac{2 \sinh}{x}$$

\rightarrow Verify using integration
(Calculus)

$$1 + \sin(x)$$



t - domain



x - domain



Discretized Integral Equation

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

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

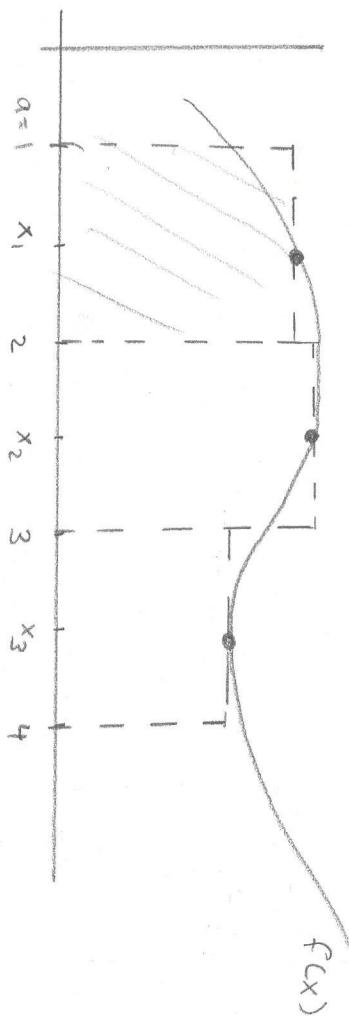
$$\begin{aligned}\vec{u} &= [u_1, \dots, u_n] \\ \vec{y} &= [y_1, \dots, y_m]\end{aligned}$$

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

Question: how to find $A \vec{u} = \vec{y}$ (next page)

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

Numerical Integration



$n = 3$: Number of data points

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

$$w = \frac{b-a}{n}$$

bin width (here $w = 1$)

$$\int_a^b f(x) dx = \sum_{i=1}^n f(x_i) w$$

} 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 $j = 1 \dots m$ $w_m = \frac{d-c}{m}$

$$y_i = \int_a^b K(x_i, t) u(t) dt$$

$$y_i = K(x_{i1}, t_1) u(t_1) w_n + K(x_{i1}, t_2) u(t_2) w_n + \dots + K(x_{i1}, t_n) u(t_n) w_n$$

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

$$= \begin{bmatrix} y_1 \\ \vdots \\ y_m \end{bmatrix} = \begin{bmatrix} K(x_{i1}, t_1) w_n & \dots & K(x_{i1}, t_n) w_n \\ \vdots & \ddots & \vdots \\ K(x_m, t_1) w_n & \dots & K(x_m, t_n) w_n \end{bmatrix} \begin{bmatrix} u_1 \\ \vdots \\ u_n \end{bmatrix}$$

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

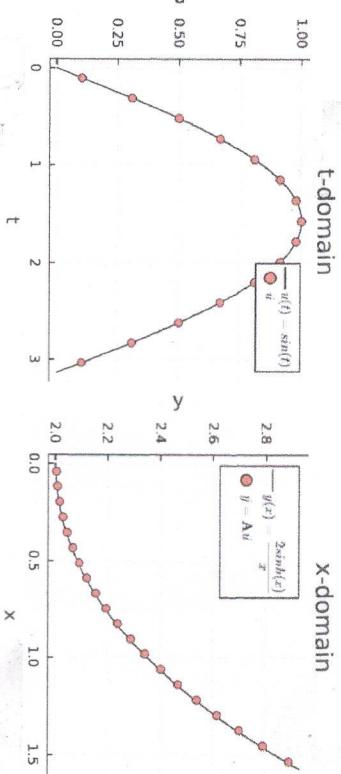
$$A = R^{m \times n}$$

$$y = A u$$

$$R^m$$

Example Discretization

$\rightarrow n_m$ can be independently varied
 \rightarrow over-determined and under-determined transforms



$$n = 15$$

$$m = 20$$

Example Inversion

$$y = A u$$

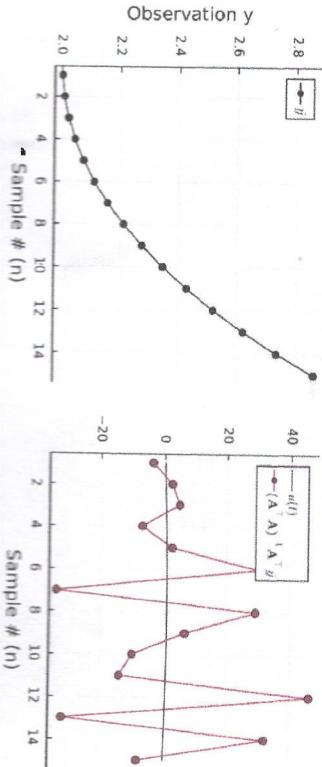
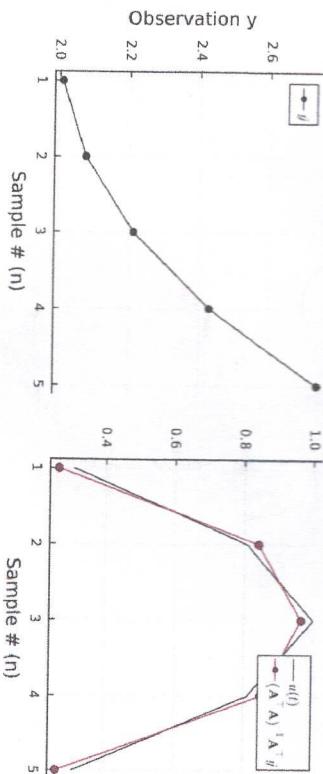
Observation in X-domain

Prediction in t-domain

$\hat{u} = \min_{\hat{u}} \|Au - y\|^2$
↑
find \hat{u}
that is best

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

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

$$A = \begin{bmatrix} 1 & 2 \\ 0 & 0 \end{bmatrix}$$

$$\det A = 0$$

} - underdetermined
- no unique solution

$$A = \begin{bmatrix} 1 & 2 \\ 2 & 4 \end{bmatrix}$$

$$\det A = 0$$

} A^{-1} not defined

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

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 $\|A\bar{u} - y\|^2$
- (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 $\text{pinv}(A)$ in most languages

Singular Value Decomposition

Full SVD

Singular values

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

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

$\mathbb{R}^{m \times m}$

$\mathbb{R}^{n \times n}$

left singular vectors

right singular vectors

$\Sigma =$

$$\Sigma = \begin{bmatrix} \sigma_1 & 0 & 0 & \dots & 0 \\ 0 & \ddots & & & 0 \\ 0 & & \ddots & & 0 \\ \vdots & & & \ddots & 0 \\ 0 & & & & \sigma_p \end{bmatrix}$$

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

$$p = \min(m, n)$$

- can find using numerical methods

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

Moore - Penrose Inverse

$$A^+ = \sqrt{\Sigma^+} U^\top$$

$$\Sigma^+ = \begin{bmatrix} \frac{1}{\sigma_1} & & & & & & & \\ & 0 & - & - & - & - & 0 & \\ & 0 & & & & & & \\ & & 1 & & & & & \\ & & & 1 & & & & \\ & & & & 1 & & & \\ & & & & & 1 & & \\ & & & & & & 1 & \\ 0 & 0 & - & - & - & - & \frac{1}{\sigma_p} & \\ & & & & & & & 0 \end{bmatrix}$$

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

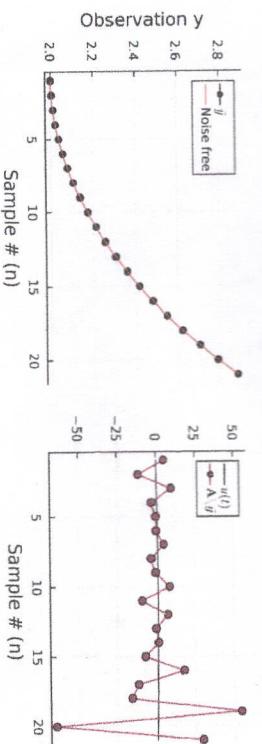
Inversion Revisited

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

"infected by
numerical or
measurement noise."

x -domain

t -domain



$$u = \underbrace{(A^T A)^{-1} A^T}_{\text{near zero singular values}} y$$

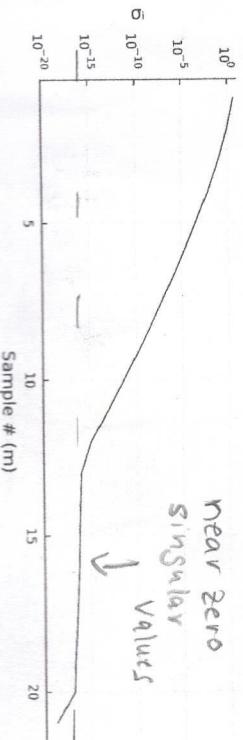
$$u = A^+ y$$

↑
log scale
↓

Machine precision

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

No new information



Near zero
singular
values

Regularized Inverse (Noise Filtering)

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

Ordinary
Regression

Starting point for $L = \frac{1}{\lambda}$
identity
matrix

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

$$u_0 = A^+ y$$

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

then

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

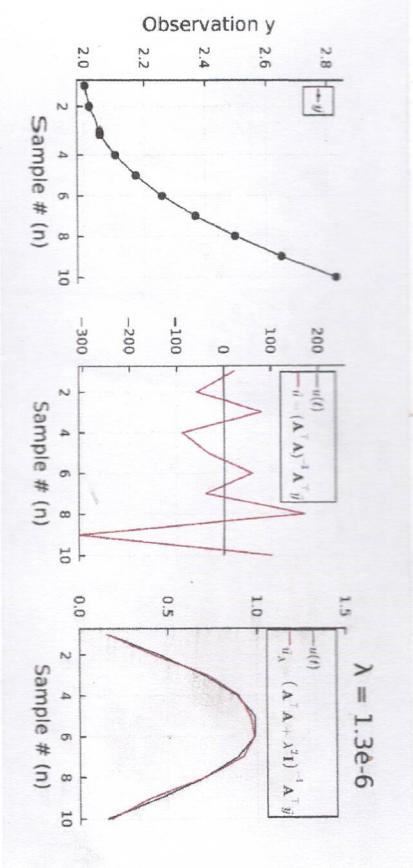
$u_0 = 0$
(no initial guess)

Inversion Revised

- 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



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

$$\lambda = 0$$

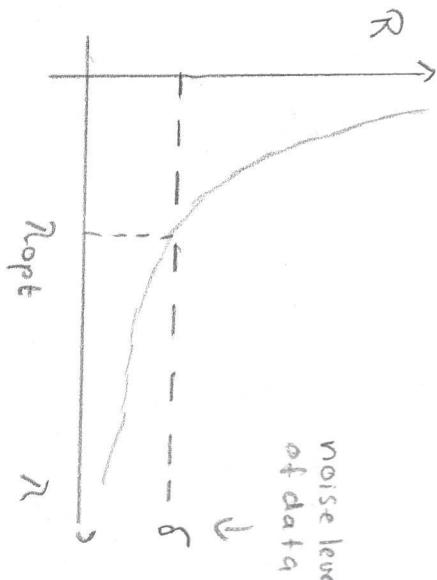
Methods to find λ

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

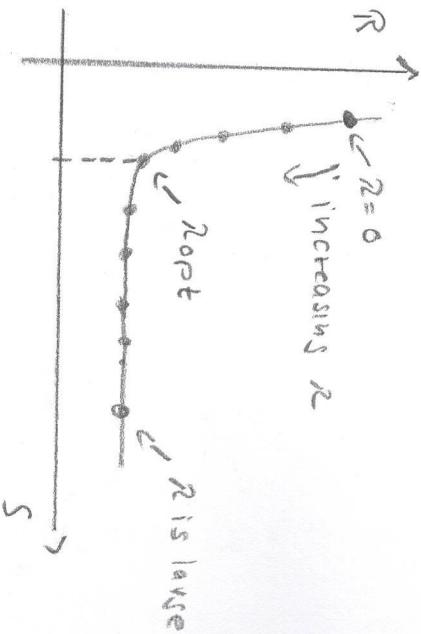
R = Residual Norm

S = Solution Norm

Morozov



L-curve



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