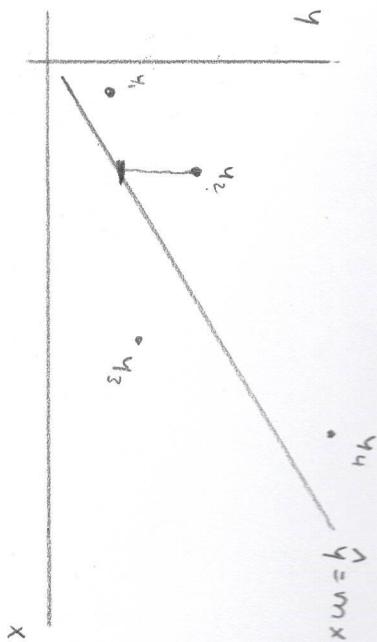


Linear Regression



n observations of y

y : income

x : years of college
(explanatory variable)

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

SSE: sum squared error

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

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

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

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

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

\bar{y} mean of all observations

$$SST = 0 \Rightarrow \text{all values equal mean}$$



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

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?

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

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

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

"best fit" \Rightarrow minimize SSE

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

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

steps

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

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

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

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

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

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

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

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

derivative rules

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

useful identities

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

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

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

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

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

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

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

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

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

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

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

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

Common matrix derivations

$$\vec{y} \quad | \quad \frac{\partial \vec{y}}{\partial \vec{x}}$$

$$A \vec{x}$$

$$x^T A$$

$$x^T x$$

$$A^T$$

$$2x$$

$$x^T A^T A x$$

$$2A^T A x$$

$$y^T A x$$

$$A^T y$$

Linear regression in practice

| X | y |
|-----|------|
| 1 | -0.2 |
| 2 | 2.2 |
| 3 | 3.5 |
| . | . |
| 7 | 7.5 |
| 8 | 12 |

$$A = \begin{matrix} n \times 2 \\ \begin{bmatrix} 1 & -0.2 \\ 1 & 2.2 \\ \vdots & \vdots \\ 1 & 7.5 \\ 1 & 12 \end{bmatrix} \end{matrix}$$

Column 1 Column of x_i

$$y = \begin{bmatrix} -0.2 \\ 2.2 \\ \vdots \\ 7.5 \\ 12 \end{bmatrix}$$

Regressed
variable

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

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

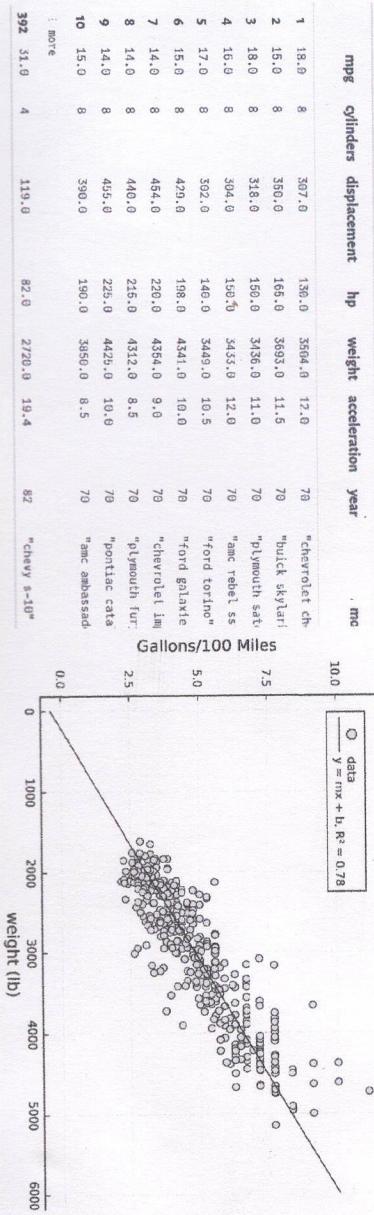
base case $y = \frac{1}{mpg}$

$x = \text{weight}$

$$y = mx + b$$

(Gallons per 100 miles)

$$R^2 = 0.79$$



Multi-variable model

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

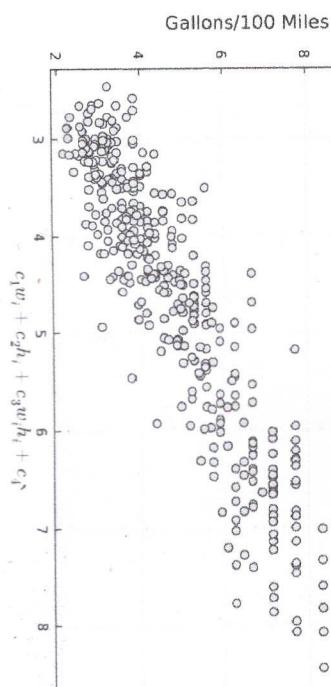
weight
horsepower
weight * horsepower

$$\hat{\beta} = A^{-1} B$$

$$\begin{bmatrix} Y \\ X_1 \\ X_2 \\ X_3 \end{bmatrix} = \begin{bmatrix} X_{11} & X_{21} & X_{31} \\ \vdots & \vdots & \vdots \\ X_{1n} & X_{2n} & X_{3n} \end{bmatrix} \begin{bmatrix} c_1 \\ c_2 \\ c_3 \\ c_4 \end{bmatrix}$$

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

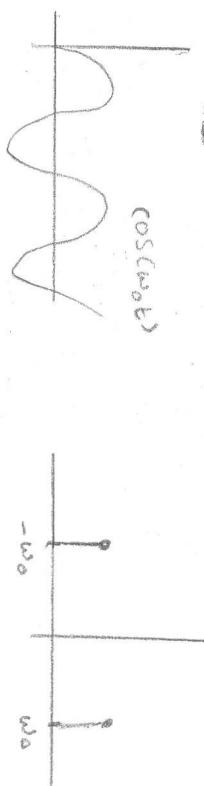
$$R^2 = 0.82 > 0.79$$



Adding variables increases explanatory power

Integral Equations (Domain Transformations)

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



t - domain \longleftrightarrow ω - domain

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

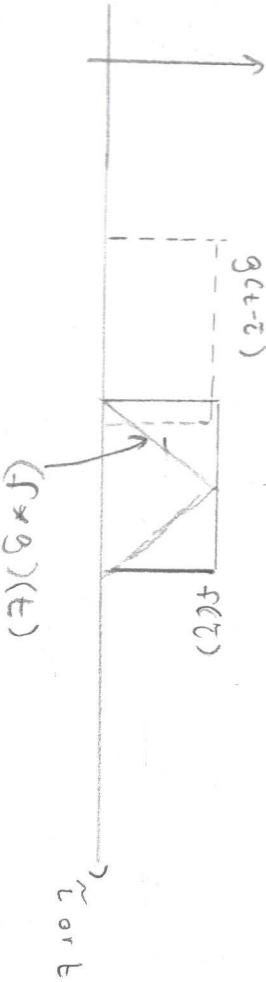
t - domain \longleftrightarrow s - domain

Convolution

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

$g(t-\tau)$

$f(\tau)$



L or b

General Integral Equation

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

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
- FTIR

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

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

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

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

Example Inversion Problem

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

"Fredholm Integral Eq."

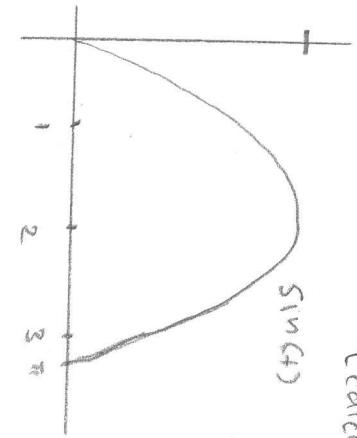
$$K(x,t) = \exp(x \cdot \cos(t)) \quad u(t) = \sin(t) \quad y(x) = \frac{2 \sinh(x)}{x}$$

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

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

→ Verify using integration
Calculus

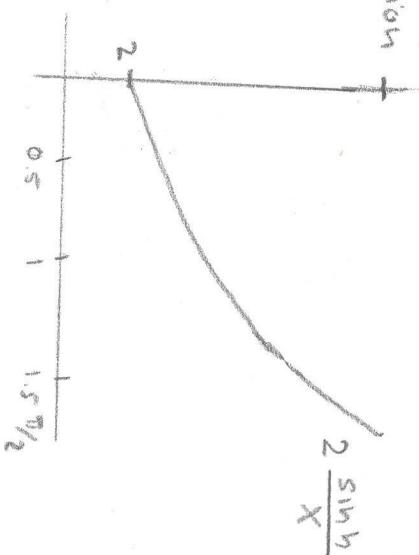
$$1 + \sin(x)$$



t - domain



x - domain



Discretized Integral Equation

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

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

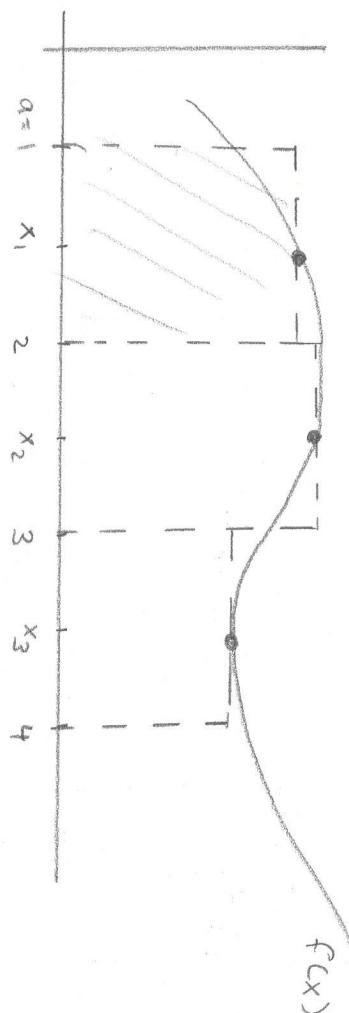
$$\begin{aligned}\vec{u} &= [u_1, \dots, u_n] \\ \vec{y} &= [y_1, \dots, y_m]\end{aligned}$$

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

Question: how to find A ? (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

$$w = \frac{b-a}{n} \quad \text{bin width} \quad (\text{where } n = 1)$$

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

} 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 $i = 1, \dots, n$

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

y is discretized into m points $j = 1, \dots, m$

$$w_m = \frac{d-c}{m}$$

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

$$y_1 = K(x_1, t_1) u(t_1) w_1 + K(x_1, t_2) u(t_2) w_2 + \dots + K(x_1, t_n) u(t_n) w_n$$

$$y_m = K(x_m, t_1) u(t_1) w_1 + K(x_m, t_2) u(t_2) w_2 + \dots + K(x_m, t_n) u(t_n) w_n$$

$$\begin{bmatrix} y_1 \\ \vdots \\ y_m \end{bmatrix} = \begin{bmatrix} K(x_1, t_1) w_1 & \dots & K(x_1, t_n) w_n \\ \vdots & \ddots & \vdots \\ K(x_m, t_1) w_1 & \dots & K(x_m, t_n) w_n \end{bmatrix} \begin{bmatrix} w_1 \\ \vdots \\ w_n \end{bmatrix}$$

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

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

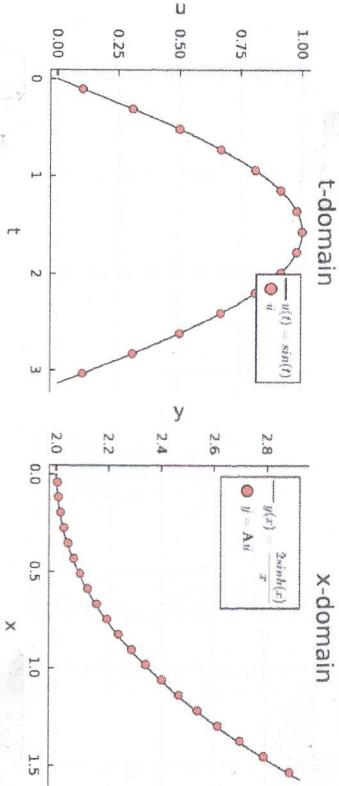
$$u = R^n$$

$$y = A u$$

$$R^m$$

Example Discretization

→ n, m can be independently varied
→ over determined and under determined transforms



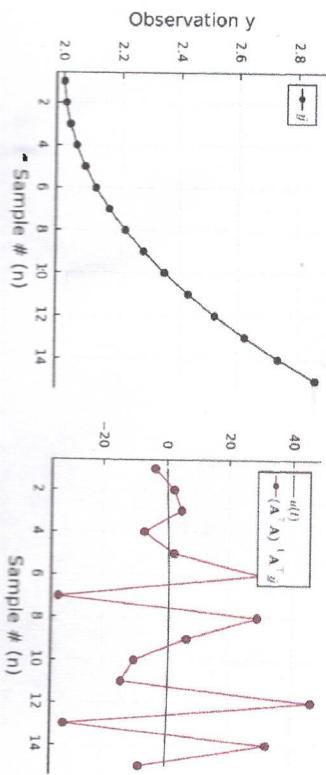
$n = 15$

$m = 20$

Example Inversion

Observation in X-domain

Prediction in t-domain



$$y = A u(t)$$

$$\text{SSE} : \underset{\uparrow}{\text{minimize}} \|Au - y\|^2$$

find \hat{u}
that is best

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

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

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

$\text{rank}(CA)$ = 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

Full SVD

Singular values

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

\downarrow
 $\mathbb{R}^{m \times n}$
 \uparrow
 $\mathbb{R}^{m \times m}$
 \swarrow
 $\mathbb{R}^{n \times n}$

left singular vectors
right singular vectors

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

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

$\rho = \min\{\sigma_i\}$

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

Moore - Penrose Inverse

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

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

- 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

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

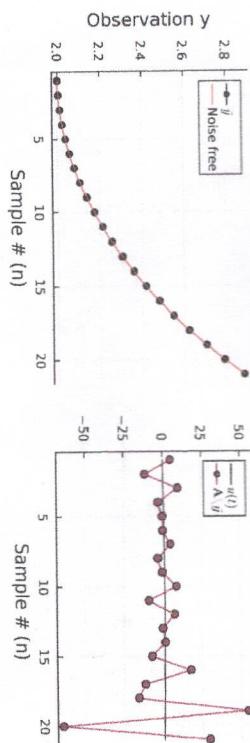
x -domain

t -domain

✓ "infected by
numerical or
measurement noise."

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

$$u = A^+ y$$



↓
near zero
singular
values

↑
log scale

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

↓
no new information

Regularized inverse (noise filtering)

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

↑ regularization parameter

Ordinary regression

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

$$u_0 = A^+ y$$

Starting point for $L = \frac{1}{n} I$

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

Identifying matrix

$$u_0 = 0$$

(no initial guess)

then

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

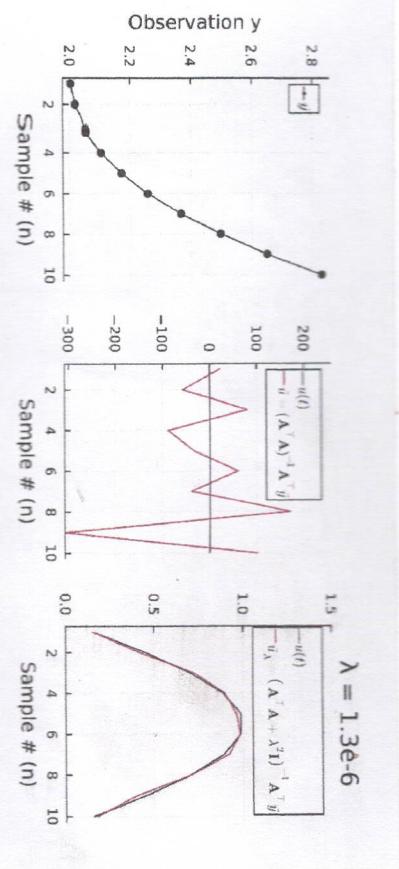
"Ridge Regression"

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



$$\lambda = 1.3e-6$$

$$u_\lambda = (\mathbf{A}^\top \mathbf{A} + \lambda^2 \mathbf{I})^{-1} \mathbf{A}^\top y$$

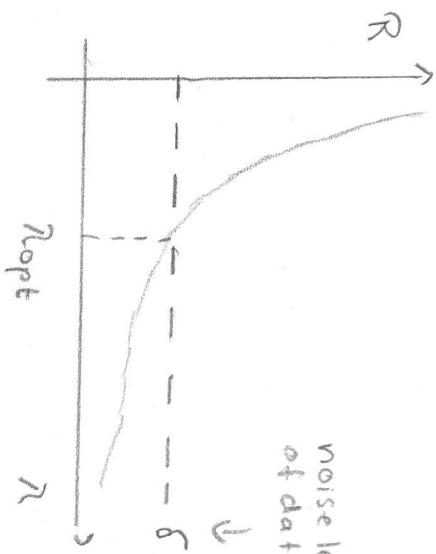
Methods to find λ

$$u_\lambda = \text{minimize } \{ \|Au - y\|^2 + \lambda^2 \|L(u - u_0)\|^2 \}$$

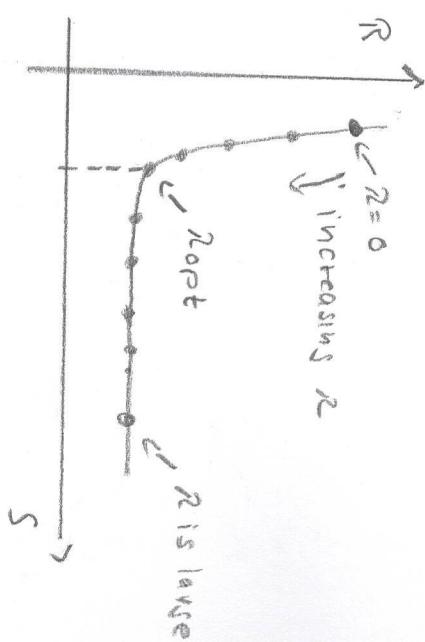
R = Residual Norm

S = Solution Norm

Morozov



L-curve



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