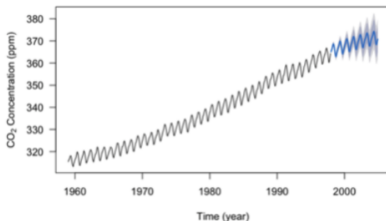
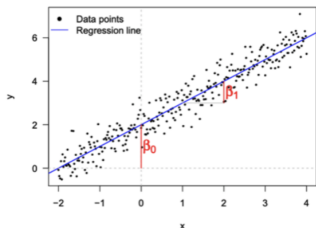


# Lecture 16

## Additional Topics in Regression and Time Series Analysis

*MATH 4070: Regression and Time-Series Analysis*



Whitney Huang  
Clemson University

- We have mainly focused on **linear regression** so far

**Model:**  $y = x\beta + \varepsilon, \quad \varepsilon \sim N(0, \sigma^2)$

**Data:**  $y$  (response vector);  $X$  (design matrix)

- $\hat{\beta} = (X^T X)^{-1} X^T y; \hat{y} = X\hat{\beta} = \underbrace{X(X^T X)^{-1} X^T}_{H: \text{"Hat" matrix}} y$
- $\hat{\beta} \sim N(\beta, \sigma^2(X^T X)^{-1})$

- **Non-parametric** regression modeling

**Model:**  $y = f(x) + \varepsilon \Rightarrow E[y|x] = f(x)$

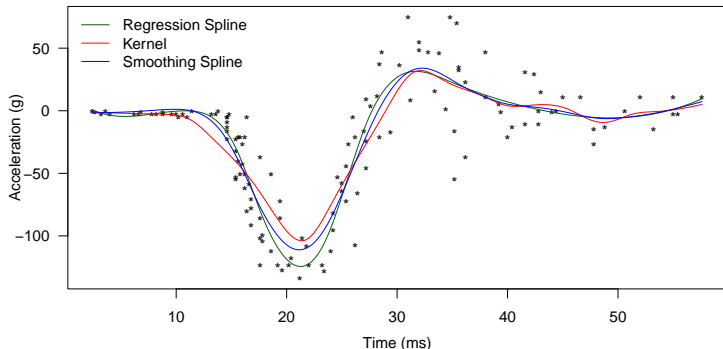
- The (smooth) function  $f(x)$  must be represented somehow
- The degree of smoothness of  $f(x)$  must be made controllable
- Some means for estimating the most appropriate degree of smoothness from data is required

# Examples of Nonparametric Regression Fits

**Regression Spline:** 10 degrees of freedom quantile knot

**Smoothing Spline:** the amount of smoothness is estimated from the data by GCV

**Kernel Regression:**  $K$ : Epanechnikov kernel and  $h = 5$



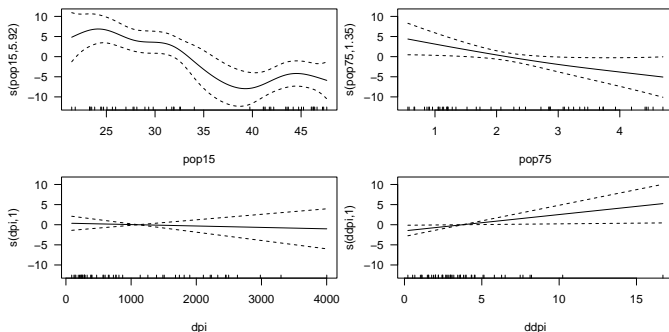
General non-parametric regression models

$$y = f(x_1, x_2, \dots, x_p) + \varepsilon$$

suffer from the “curse of dimensionality”

**Generalized Additive Models:**

$$y = \beta_0 + f_1(x_1) + f_2(x_2) + \dots + f_p(x_p) + \varepsilon$$



LASSO

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \cdots + \beta_{p-1} x_{p-1} + \varepsilon, \quad \varepsilon \sim N(0, \sigma^2)$$

$x_1, x_2, \dots, x_{p-1}$  are the predictors.

**Question:** What if we have too many predictors (i.e.,  $p$  is “large”)?

- Explanation can be difficult due to collinearity
- Can lead to overfitting by using too many predictors

Two methods, namely **Ridge regression** and **LASSO**, allow us to “shrink” the information contained in all the predictors into a more useful form

Ridge regression assumes that the regression coefficients (after normalization) should not be very large

- The ridge regression estimate chooses the  $\beta$  that minimizes:

$$\sum_{i=1}^n (y_i - \beta_0 - \sum_{j=1}^{p-1} \beta_j x_{ij})^2 + \lambda \sum_{j=1}^{p-1} \beta_j^2,$$

where  $\lambda \geq 0$  is a **tuning parameter** to be determined via cross-validation

- The ridge regression estimates:

$$\hat{\beta}_{\text{ridge}} = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y}$$

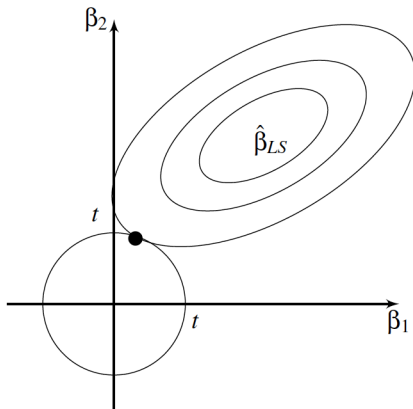
- Ridge regression is particularly effective when the model matrix is collinear

## Graphical Illustration of Ridge Regression

Estimation of ridge regression can also be solved by choosing  $\beta$  to minimize

$$\sum_{i=1}^n (y_i - \beta_0 - \sum_{j=1}^{p-1} \beta_j x_{ij})^2$$

subject to  $\sum_{j=1}^p \beta_j^2 \leq t^2$



Source: p. 175, Fig. 11.9 *Linear Models with R*, Faraway, 2014

# Least Absolute Shrinkage and Selection Operator (LASSO)

Tibshirani, 1996

LASSO assumes the effects are **sparse** in that the response can be explained by a small number of predictors with the rest having no effect

- LASSO choose  $\hat{\beta}$  to minimize:

$$\sum_{i=1}^n (y_i - \beta_0 - \sum_{j=1}^{p-1} \beta_j x_{ij})^2 + \lambda \sum_{j=1}^{p-1} |\beta_j|$$

- No explicit solution to this minimization problem
- The penalty term has the effect of forcing some of the coefficient estimates to be zero when the tuning parameter  $\lambda$  is “large”  $\Rightarrow$  performs **shrinkage** and **variable selection**

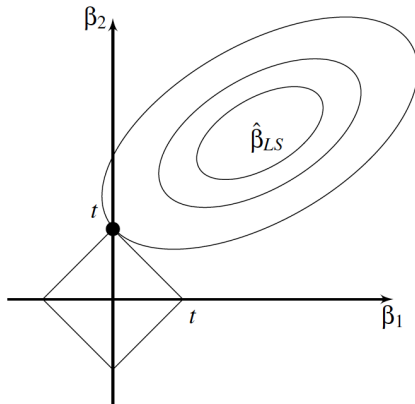


## Graphical Illustration of LASSO

Estimation of LASSO can also be solved by choosing  $\beta$  to minimize

$$\sum_{i=1}^n (y_i - \beta_0 - \sum_{j=1}^{p-1} \beta_j x_{ij})^2$$

subject to  $\sum_{j=1}^p |\beta_j| \leq t$



Source: p. 175, Fig. 11.9 *Linear Models with R*, Faraway, 2014

- **Gaussian Linear Model:**

$$y \sim N(\mu, \sigma^2), \quad \mu = \mathbf{x}^T \boldsymbol{\beta}$$

- **Bernoulli Linear Model:**

$$y \sim \text{Bernoulli}(\pi), \quad \log\left(\frac{\pi}{1 - \pi}\right) = \mathbf{x}^T \boldsymbol{\beta}$$

- **Poisson Linear Regression:**

$$y \sim \text{Poisson}(\lambda), \quad \log \lambda = \mathbf{x}^T \boldsymbol{\beta}$$

These models fall into the family of **generalized linear models** [Nelder and Wedderburn (1972); McCullagh and Nelder (1989)] with the **distributional assumptions** (normal, Bernoulli, Poisson) and the **link functions** (identity, logit, and log)

- Time domain methods [Box and Jenkins, 1970]:

- Regress present on past

**Example:**  $Y_t = \phi Y_{t-1} + Z_t$ ,  $|\phi| < 1$ ,  $\{Z_t\} \sim \text{WN}(0, \sigma^2)$

- Capture dynamics in terms of “velocity”, “acceleration”, etc

- Frequency domain methods [Priestley, 1981]:

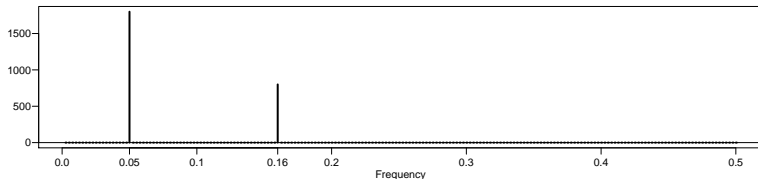
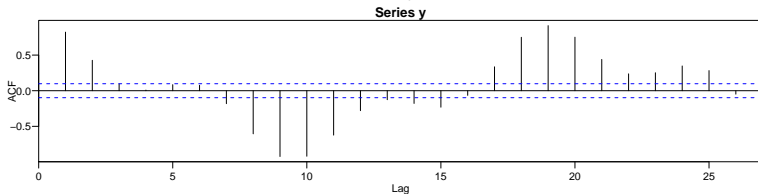
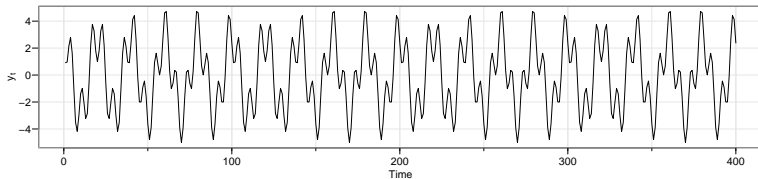
- Regress present on periodic sines and cosines

**Example:**  $Y_t = \alpha_0 + \sum_{j=1}^p [\alpha_{1j} \cos(2\pi\omega_j t) + \alpha_{2j} \sin(2\pi\omega_j t)]$

- Capture dynamics in terms of **resonant frequencies**

# Searching Hidden Periodicities

$$y_t = 3 \cos\left(2\pi\left(\frac{10}{200}\right)t\right) + 2 \cos\left(2\pi\left(\frac{32}{200}t + 0.3\right)\right)$$



LASSO

## Spectral density $\iff$ Covariance function

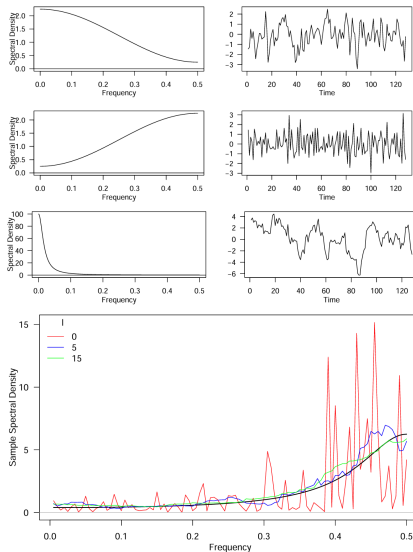
If  $\{Y_t\}$  has  $\sum_{h=-\infty}^{\infty} |\gamma(h)| < \infty$ ,  
then its **spectral density** is

$$f(\omega) = \sum_{h=-\infty}^{\infty} \gamma(h) e^{-2\pi i \omega h}$$

for  $-\infty < \omega < \infty$ . We have

$$\gamma(h) = \int_{-\frac{1}{2}}^{\frac{1}{2}} e^{2\pi i \omega h} f(\omega) d\omega$$

**Smoothing techniques**, like  
those in **nonparametric  
regression**, are needed to  
estimate  $f(\omega)$  well



Log-returns,  $r_t = \log\left(\frac{y_t}{y_{t-1}}\right)$ , are often modeled instead of daily stock prices,  $y_t$ , in financial time series analysis

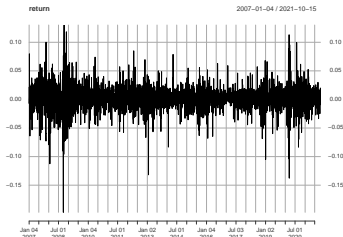
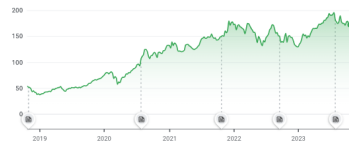
Apple Inc

\$171.95 +217.95% +117.87 5Y

Oct 25, 10:26:13 AM UTC-4 USD NASDAQ Disclaimer

1D 5D 1M 6M YTD 1Y 5Y MAX

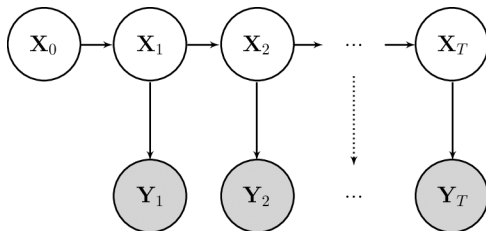
Key events



Generalized Autoregressive Conditional Heteroskedasticity (GARCH) is commonly used to model the dynamics of fluctuations in log-returns to capture volatility clustering.

$$r_t = \mu_t + a_t, \quad a_t = \sigma_t \epsilon_t, \quad \sigma_t^2 = \alpha_0 + \sum_{i=1}^m \alpha_i a_{t-i}^2 + \sum_{j=1}^s \beta_j \sigma_{t-j}^2$$

LASSO



State:  $\mathbf{X}_t = \mathbf{M}_t \mathbf{X}_{t-1} + \mathbf{V}_t$ ,  $\mathbf{V}_t \stackrel{i.i.d.}{\sim} \text{WN}(\mathbf{0}, \mathbf{Q}_t)$ ,  $t = 1, 2, \dots$

Observation:  $\mathbf{Y}_t = \mathbf{H}_t \mathbf{X}_t + \mathbf{W}_t$ ,  $\mathbf{W}_t \stackrel{i.i.d.}{\sim} \text{WN}(\mathbf{0}, \mathbf{R}_t)$ ,  $t = 1, 2, \dots$

- $\mathbf{X}_t \in \mathbb{R}^p$  and  $\mathbf{Y}_t \in \mathbb{R}^q$  are the **state vector** and the **observation vector** at time  $t$
- $\mathbf{M}_t$  is the  $p \times p$  **transition matrix**, and  $\mathbf{H}_t$  is the  $q \times p$  **observation matrix**
- $\mathbf{V}_t$  and  $\mathbf{W}_t$  are the state and observation noises

**Goal:** To estimate the underlying unobserved signal  $X_t$ , given the data  $\mathbf{Y}_{1:s} = \mathbf{y}_{1:s} = \{\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_s\}$ :

- When  $s < t$ , the problem is called **forecasting** or **prediction**
- When  $s = t$ , the problem is called **filtering**
- When  $s > t$ , the problem is called **smoothing**

In addition to these estimates, we would also want to measure their precision. The solution to these problems is accomplished via the **Kalman filter** and **Kalman smoother**



Assume the filtering distribution at time  $t - 1$  is

$$[\mathbf{X}_{t-1} | \mathbf{Y}_{1:t-1}] \sim N(\boldsymbol{\mu}_{t-1}^a, \Sigma_{t-1}^a)$$

- **Forecast Step:** Gives the forecast distribution at time  $t$ :

$$[\mathbf{X}_t | \mathbf{Y}_{1:t-1}] \sim N(\boldsymbol{\mu}_t^f, \Sigma_t^f),$$

where  $\boldsymbol{\mu}_t^f = M_t \boldsymbol{\mu}_{t-1}^a$ , and  $\Sigma_t^f = M_t \Sigma_{t-1}^a M_t^T + Q_t$ .

- **Update Step:** updates the forecast distribution using new data  $\mathbf{Y}_t$

$$[\mathbf{X}_t | \mathbf{Y}_{1:t}] \sim N(\boldsymbol{\mu}_t^a, \Sigma_t^a),$$

where  $\boldsymbol{\mu}_t^a = \boldsymbol{\mu}_t^f + K_t (\mathbf{Y}_t - H_t \boldsymbol{\mu}_t^f)$ , and  $\Sigma_t^a = (I - K_t H_t^T) \Sigma_t^f$ ,  
and

$$K_t = \Sigma_t^f H_t^T (H_t \Sigma_t^f H_t^T + R_t)^{-1}$$

is the **Kalman gain matrix**

All the methods presented for univariate time series also apply to multivariate processes

$$\{\mathbf{Y}_t \in \mathbb{R}^p\}$$

- The theory becomes more involved as we generalize to the cross-covariance:

$$\text{Cov}(\mathbf{Y}_s, \mathbf{Y}_t) = \mathbf{C}(s, t),$$

where  $\mathbf{C}(\cdot, \cdot)$  is the  $p \times p$  matrix-valued **cross-covariance function (CCVF)**

- Similarly, in the frequency domain approach, the **cross-spectrum** is given by:

$$f_{XY}(\omega) = \sum_{h=-\infty}^{\infty} \gamma_{XY}(h) e^{-2\pi i \omega h}$$

VAR( $p$ ) model:

$$\mathbf{Y}_t = \boldsymbol{\mu} + A_1 \mathbf{Y}_{t-1} + \cdots + A_p \mathbf{Y}_{t-p} + \mathbf{W}_t, \quad t = 0, 1, 2, \dots,$$

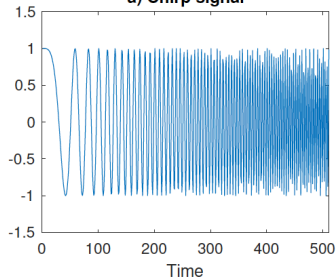
where

- $\mathbf{Y}_t = (Y_{1t}, \dots, Y_{pt})^T$  is a  $(p \times 1)$  random vector
- $A_i$  are  $(p \times p)$  coefficient matrices
- $\boldsymbol{\mu} = (\mu_1, \dots, \mu_p)^T$  is the intercept vector
- $\mathbf{W}_t = (W_{1t}, \dots, W_{pt})^T$  is a  $p$ -dimensional white noise, i.e.,  $E[\mathbf{W}_t] = \mathbf{0}$ ,  $E[\mathbf{W}_t \mathbf{W}_t^T] = \Sigma_{\mathbf{W}}$  and  $E[\mathbf{W}_s \mathbf{W}_t^T] = \mathbf{0}$  for  $s \neq t$ .

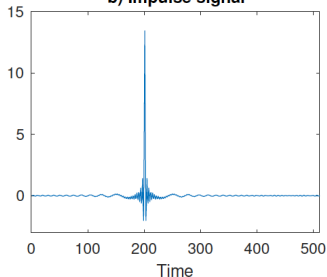
LASSO

# Time-Frequency Analysis: A Motivation Example

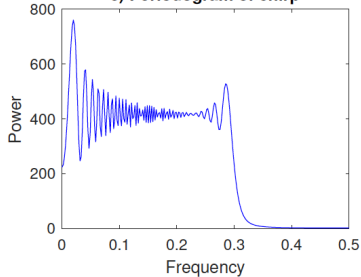
a) Chirp signal



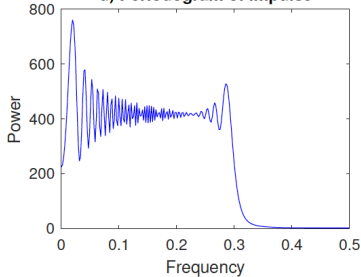
b) Impulse signal



c) Periodogram of chirp



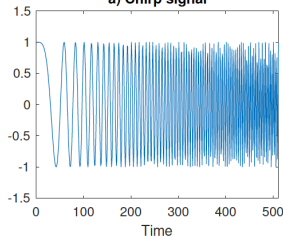
d) Periodogram of impulse



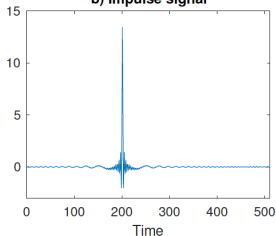
# Time-Frequency Analysis: Spectrogram

A **spectrogram** is a visual representation of the spectrum of frequencies of a signal as it **varies with time**

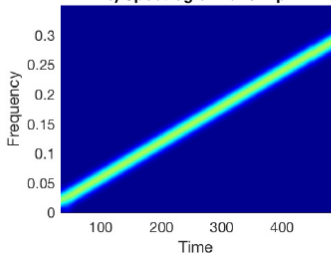
a) Chirp signal



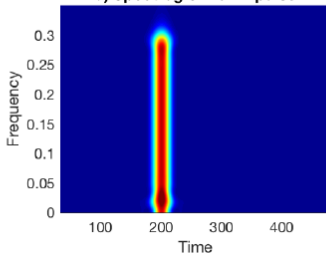
b) Impulse signal



c) Spectrogram of chirp



d) Spectrogram of impulse



Some selected references:

- Regression models for time series analysis, Kedem and Fokianos, 2002
- Handbook of discrete-valued time series, edited by Davis, Holan, Lund, Ravishanker, 2016
- Bayesian Dynamic Generalized Linear Models, Gamerman *et. al*, 2016
- Count Time Series: A Methodological Review, Davis *et. al.*, 2021

LASSO