

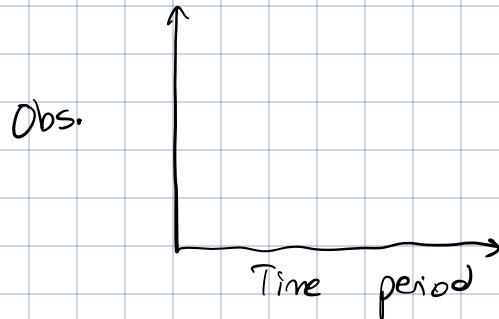
INTRODUCTION TO FORECASTING, CONTROL AND TIME SERIES

Examples

What makes this diff. than other statistics?

- ① Data is not indep.
- ② Data is not i.i.d. \leftarrow can vary across time
- ③ Abrupt changes to data generation process

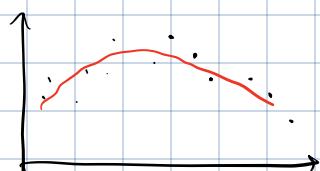
Time series plot:



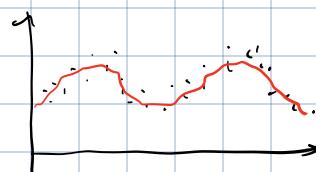
Level: local mean of observations



Trend: level varies w/ time



Seasonal effect: calendar-related effect which repeats in a given period



Change point: time where 1/more of following changes

- ① Data gen. process
- ② Method of data measurement
- ③ Obs. defn.

Additive-based Modelling

Modelling stationary data is easier than modelling changing data. We split data into stationary & non-stationary data.

$$x_t = m_t + s_t + r_t$$

↑ ↑ ↗
trend seasonality Randomness. Hopefully 0-mean & uncorr. (stationary)
deterministic func. of time

General process:

- ① Remove non-stationary elements: stabilize variance, remove seasonality & trend
 - ↳ Similar to L.R.: remove trend to produce random residuals
- ② Model stationary data
- ③ Forecast on stationary data & add back non-stationary

Time Series Models

Stochastic process: family of random variables $\{x_t, t \in T\}$ defined on a prob. space

↳ T : index set, could be real-valued func.

Ideally, we can find the joint distr. of all $\{x_t\}$

$$\text{F}(x_1, \dots, x_n) = P(x_1 \leq x_1, \dots, x_n \leq x_n)$$

We know that this joint distr. exists b/c of Kolmogorov's Existence Theorem:

Kolmogorov's Existence Theorem:

$F_T(\cdot)$ ac distribution func. of some stoch. process $\Leftrightarrow \forall n \in \{1, 2, \dots\} \quad \& \quad t = (t_1, \dots, t_n) \in T$ and $1 \leq i \leq n$:

$$\lim_{x_i \rightarrow \infty} F_t(x) = F_{t(i)}(x(i)) \quad \left. \begin{array}{l} \uparrow \quad \uparrow \\ (n-1)\text{-component vectors w/ } i^{\text{th}} \text{ comp. deleted} \end{array} \right\} \text{Like how joint} \rightarrow \text{marginal!}$$

Intuition: if we can't get CDF of $n-1$ R.Vs from n R.Vs, it's not a stoch. process

In reality, we don't need joint distr., we only need:

$$\begin{array}{ll} \textcircled{1} & E(X_t) \\ \textcircled{2} & E(X_t^2) \\ \textcircled{3} & E(X_t, X_t^*) \end{array} \quad \left\{ \begin{array}{l} \text{Var}(X_t) \\ \text{Cov}(X_t, X_t^*) \end{array} \right\}$$

Since we can't find joint distr., we don't know if $\forall i, j, X_i$ is indep. of X_j

\hookrightarrow Recall: X_i is indep. of $X_j \Leftrightarrow F(X_i, X_j) = F(X_i)F(X_j)$. We don't know any $F(\cdot)$!

Zero-Mean Time Series Models

If we follow our process listed above, our residuals ideally have $\mu=0$. Such models that assume this are called zero-mean time series models, which include:

① i.i.d. Noise

Defn: $\{X_i, i=1, \dots, n\}$ where $E(X_i) = 0 \quad \forall i$

Joint distr:

$$\begin{aligned} F(x_{i_1}, \dots, x_{i_n}) &= P(X_{i_1} \leq x_{i_1}, \dots, X_{i_n} \leq x_{i_n}) \\ &= \prod_{k=1}^n P(X_{i_k} \leq x_{i_k}) \end{aligned}$$

This implies:

$$P(X_{n+h} \leq x | X_1 \leq x_1, \dots, X_n \leq x_n) = P(X_{n+h} \leq x)$$

In other words, past does not help w/ predicting future

② Random Walk

Defn: $\{S_t, t=0, 1, \dots\}$ where $S_t = \sum_{i=1}^t X_i$, where $X_i \in \{X_t, t=1, 2, \dots\}$ of i.i.d noise. Usually, $S_0 = 0$ (or some S_0 w/ $E(S_0) = 0$)

Note:

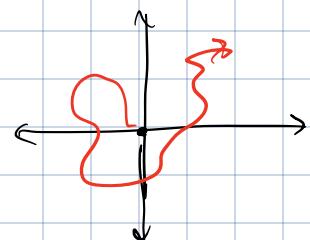
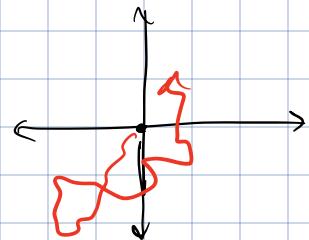
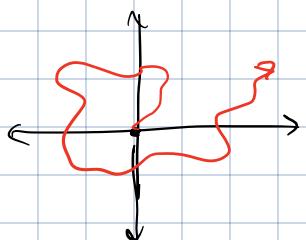
$$E(S_t) = E\left(\sum_{i=1}^t X_i\right) = \sum_{i=1}^t E(X_i) = 0 \quad \text{BUT not finite variance}$$

Ex:// Consider 2D discrete-time walk where $|dx| = |dy| = 1$

Sim. 1

Sim. 2

Sim. 3



③ White noise

Defn. $\{X_i, i=1, 2, \dots\}$ where each X_i is:

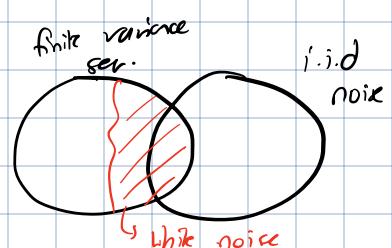
- i) Not corr. w/ another (doesn't say anything about indep.)
- ii) $E(X_i) = 0$
- iii) $\text{Var}(X_i) = \sigma^2 < \infty$

Notation: $X_i \sim WN(0, \sigma^2)$

Random variables are a sum of white noise as well

If i.i.d noise has finite variance \Rightarrow white noise.

If white noise is indep. \Rightarrow i.i.d noise



WHY ZERO-MEAN RESIDUALS \Rightarrow we will not be under/over predicting data

TREND, SEASONALITY, AND CLASSICAL DECOMPOSITION

Classical Decomposition

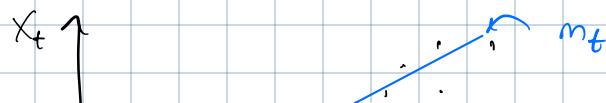
① Models w/ (non-periodic) trend

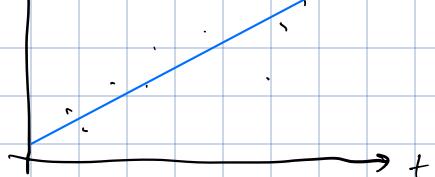
$$X_t = m_t + Y_t \leftarrow \text{random. var w/ } E(Y_t) = 0 \quad \forall t$$

↑
deterministic, slowly
changing

Note that $E(X_t) = m_t$

Ex:// $X_t = m_t + Y_t, m_t = 1 + 2t, Y_t \stackrel{\text{i.i.d.}}{\sim} N(0, 36)$





\Rightarrow Similar to linear regression

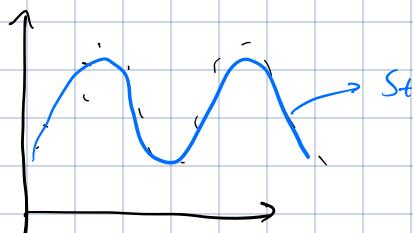
② Models w/ Seasonality (Periodic trends)

$$X_t = S_t + Y_t \quad \leftarrow \text{random. var w/ } E(Y_t) = 0 \quad \forall t$$

↑
 periodic deterministic
 function w/ period δ

Note that $E(X_t) = S_t$

Ex:// $X_t = S_t + Y_t$, $S_t = 5 - 3\cos(0.2\pi t)$, $Y_t \stackrel{\text{i.i.d.}}{\sim} N(0, 1)$.

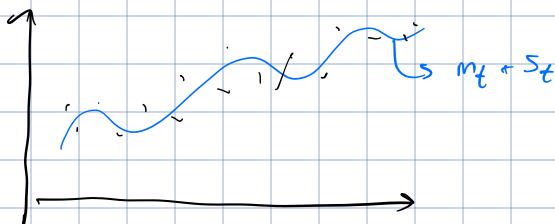


③ Models w/ both trend & seasonality

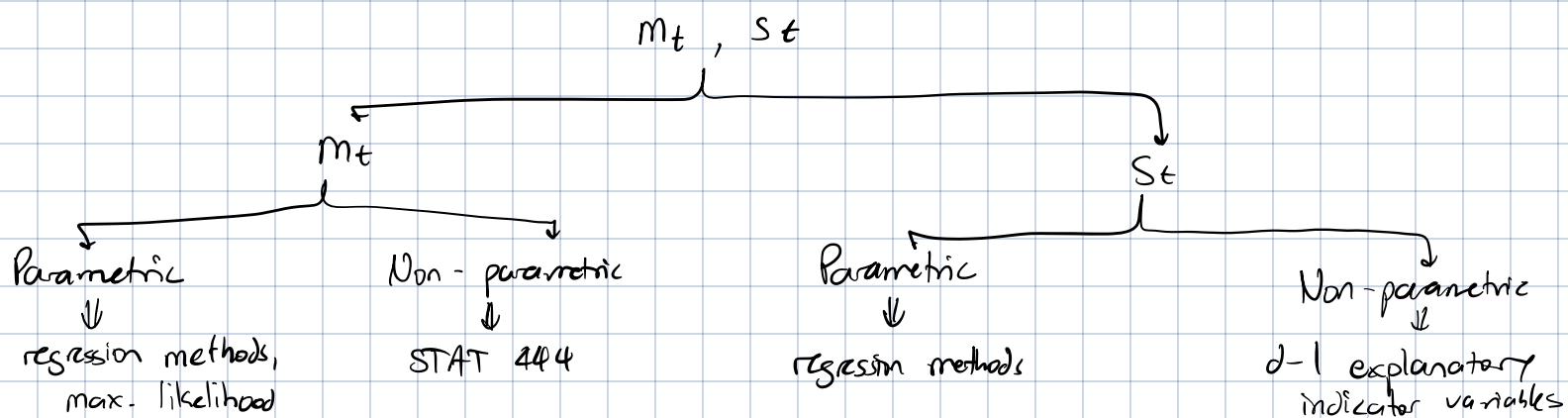
$$X_t = m_t + S_t + Y_t$$

$$E(X_t) = S_t + m_t$$

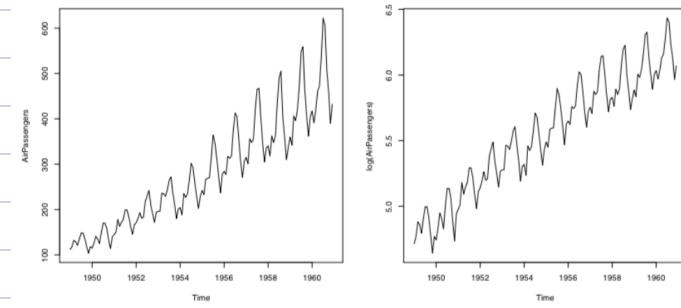
Ex:// $X_t = m_t + S_t + Y_t$, $m_t = 1 + 2t$, $S_t = 15\sin(0.3\pi t)$, $Y_t \stackrel{\text{i.i.d.}}{\sim} N(0, 7)$.



Estimating trend & seasonality via regression



Ex:// Fit models for the following pattern



If seems like there is a periodic, increasing trend. Log both sides to remove exponential behaviour when fitting.

$$\begin{aligned}\log(Y_t) &= m_t + S_t + \beta_t \\ &= \underbrace{\beta_0 + \beta_1 t + \beta_2 x_1 + \dots + \beta_{12} x_{12}}_{m_t} + \underbrace{\beta_t}_{S_t}\end{aligned}$$

S_t composed of 12 indicator variables b/c period is 12 months
 ↳ $S_t = \beta_0 + \beta_1 x_1 + \dots + \beta_{12} x_{12}$ ↑ Expected change from Jan. baseline

Exponentiate back:

$$Y_t = e^{m_t + S_t + \beta_t}$$

Evaluating model, we see problems (add)

Add a quadratic component to trend:

$$m_t = \beta_0 + \beta_1 t + \beta_2 t^2$$

Loss FUNCTIONS & FORECASTING

Introduction & Definition

Loss function: measures inconsistency b/w true & fitted values in model

↳ Notation: $L(Y, f(x))$ ↑ true values ↓ estimated values

↳ Conditions:

$$\textcircled{1} \quad L(Y, f(x)) > 0$$

$$\textcircled{2} \quad Y = f(x) \Rightarrow L(Y, f(x)) = 0$$

Risk: $E_{x,y} [L(Y, f(x))]$. If $X=x \Rightarrow E_{Y|x} [L(Y, f(x)) | X=x]$

Know how to take conditional expectations

Types of Loss Functions

① Quadratic loss / squared error

$$L_{SE}(Y, f(x)) = (Y - f(x))^2$$

Risk / expected value is mean squared error:

$$MSE(f) := E_{x,y} ((Y - f(x))^2)$$

Conditional expectation theorem:

X, Y are R.V., $E(Y) = \mu$, $\text{Var}(Y) < \infty \Rightarrow f$ that minimizes MSE

is $f(x) = E_{Y|x} [Y | X=x]$ (R.V. depends on x , NOT #)

↳ Proof outline:

① Show that constant c that minimizes $E_y [(Y - c)^2]$ is $c = \mu$
by solving following for c :

$$\frac{\partial}{\partial c} E_y [(Y - c)^2] = 0$$

② From ①, $\hat{c}(x)$ that minimizes $E_y [(Y - c(x))^2 | X=x]$ is $E[Y | X=x]$

③ We know $Z := z(x) \geq w := w(x) \Leftrightarrow \forall x, z(x) \geq w(x)$.
From ②, for any $s(\cdot)$:

$$E_y [(Y - s(x))^2 | X=x] \geq E[(Y - \hat{c}(x))^2 | X=x]$$

④ Take E_x on both sides & by law of iterated expectations:

$$E_{x,y} [(Y - s(x))^2] \geq E[(Y - \hat{c}(x))^2]$$

⑤ Conclude minimality

Ex:// In least squares regression

$$Y = X\beta + \varepsilon$$

Since $E(Y) = X\beta$, $\text{Var}(Y) < \infty$, MSE function is $E[Y|X] = X\beta$

② Absolute error loss

$$L_{\text{abs}}(Y, f(x)) = |Y - f(x)|$$

Function that minimizes $E(L_{\text{abs}}(Y, f(x)))$:

$$\hat{f}(x) = \text{Median}(Y|x)$$

Know proof

③ Zero-one loss

If Y is discrete:

$$L_{01}(Y, f(x)) = \begin{cases} 0 & \text{if } Y = f(x) \\ 1 & \text{o.w.} \end{cases}$$

Expected value:

$$E(L_{01}(Y, f(x))) = P(Y \neq f(x))$$

Mostly for classification

Function that minimizes risk

$$\hat{f}(x) = \max_{y \in S} P(y|x) \Rightarrow \text{Bayes classifier}$$

Mostly interested in quadratic loss \Rightarrow conditional expectations

REGRESSION REVIEW

Main assumption of linear model: If $X = (x_1, \dots, x_p)$, then:

$$f(x; \beta) = E[Y|X] = \beta_0 + \sum_{i=1}^p \beta_i x_i$$

In matrix notation:

$$Y = X\beta + \epsilon$$

↓ variable
 obs. 1 → $\begin{bmatrix} 1 & x_{11} & x_{12} & \dots & x_{1p} \end{bmatrix}$ $\begin{bmatrix} \beta_0 \\ \vdots \\ \beta_p \end{bmatrix}$ $\begin{bmatrix} \epsilon_1 \\ \vdots \\ \epsilon_n \end{bmatrix}$
 $(n \times (p+1))$ $(p+1) \times 1$ $(n \times 1)$

Parameter Estimation

Set loss function to residual sum of squares

$$RSS(\beta) := \sum_{i=1}^n \left(y_i - \underbrace{\beta_0 + \sum_{j=1}^p \beta_j x_{ij}}_{X\beta} \right)^2$$

To solve this, solve:

$$\hat{\beta} = \min_{\beta} \underbrace{((Y - X\beta)^T (Y - X\beta))}_{\epsilon}$$

Taking derivative:

$$\frac{\partial}{\partial \beta} RSS(\beta) = -2X^T(Y - X\beta) = 0$$

$$\hat{\beta} = (X^T X)^{-1} X^T Y$$

Use second derivative test \rightarrow positive definite \rightarrow solution is minimal

Hat matrix

Fitted values:

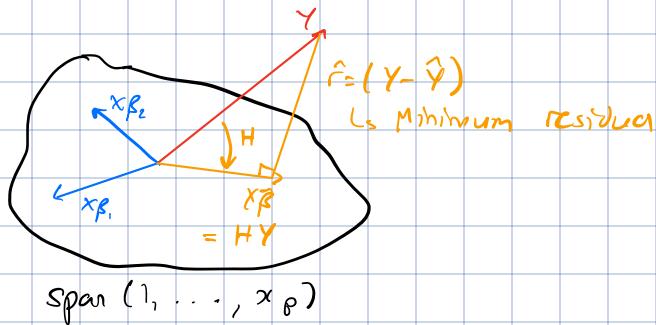
$$\begin{aligned}\hat{Y} &= X\hat{\beta} \\ \hat{Y} &= X(X^T X)^{-1} X^T Y \\ &\text{transforms } Y \rightarrow \hat{Y} \Rightarrow \text{hat matrix } (n \times n) (H)\end{aligned}$$

Properties:

① Idempotent

$$H^2 = H \rightarrow \text{orthogonal}$$

② Geometry:



How do we know \hat{Y} & r are orthogonal:

$$\hat{r} = Y - \hat{Y} = Y - HY = (I - H)Y$$

$$\hat{r} \cdot \hat{Y} = H(I - H)Y = HY - H^2Y = HY - HY = 0 \checkmark$$

③ Symmetric

$$H = H^T$$

④ Trace:

$$\text{trace } ((\cdot)) = \sum_{i=1}^n h_{ii} = p + 1$$

⑤ Prediction of new value is just transformation of \mathbf{Y}

$$\hat{Y}_i = \sum_{i=1}^n h(x_i) y_i$$

Statistical inference

Another assumption: $\varepsilon \stackrel{\text{iid}}{\sim} N(0, \sigma^2)$ \rightarrow only for inference

Distr. of $\tilde{\beta} = (\mathbf{x}^\top \mathbf{x})^{-1} \mathbf{x}^\top \tilde{\mathbf{y}}$:

$$\tilde{\beta} \sim MVN(\beta, \sigma^2 (\mathbf{x}^\top \mathbf{x})^{-1})$$

Estimating σ^2 :

$$\tilde{\sigma}^2 = \frac{SSE}{n - p - 1}$$

↳ Distr.:

$$(n - p - 1) \cdot \frac{\tilde{\sigma}^2}{\sigma^2} \sim \chi^2_{n - p - 1}$$

Possible tests:

A: $H_0: \beta_j = 0$

$$z_j = \frac{\hat{\beta}_j}{\sigma \sqrt{v_j}} \sim t_{n-p-1}$$

$(j+1)^{\text{th}}$ diagonal elem $(\mathbf{x}^\top \mathbf{x})^{-1}$

B: Test if smaller model w/ p_0 variates is better than one w/ p_1 variates ($p_0 < p_1$)

$$\frac{(RSS_0 - RSS_1) \div (p_1 - p_0)}{RSS_1 \div (n - p_1 - 1)} \sim F_{p_1 - p_0, n - p_1 - 1}$$

Further notes

① Sum of squares

$$\underbrace{\sum_{i=1}^n (y_i - \bar{y})^2}_{\text{Total sum of squares.}} = \underbrace{\sum_{i=1}^n (\hat{y}_i - \bar{y})^2}_{SS_{\text{Res}}} + \underbrace{\sum_{i=1}^n (y_i - \hat{y}_i)^2}_{SSE}$$

total variability = model variability + error

Balance \rightarrow bias-variance tradeoff (eg. decreasing SSE (variance) \rightarrow increase SS_{Reg} (bias))

② $\hat{\sigma}$:

residual standard error $\rightarrow \hat{\sigma}$

③ Prediction interval

If predict on x is \hat{y} w/ $[a, b]$ as 95% interval $\rightarrow P(a < Y_0 < b) = 0.95$
the value

Why is this interpretation diff than C.I. $\Rightarrow Y_0$ is R.V., param in C.I. is fixed

④ Interpretation:

$\hat{\beta}_j = a \Rightarrow$ expected diff. in y w/ 1 unit change in x_j w/ all other
variables the same yields a

BESTIAL DIAGNOSTICS

Assuming $\varepsilon_i \sim N(0, \sigma^2) \rightarrow$ residual checking for assumption correctness

① Normality

Graphical: plot residuals on Q-Q plot. If straight line \rightarrow approx. Normal

* Note: histogram \Leftrightarrow Q-Q plot*

Formal test: Shapiro-Wilk test (non-parametric)

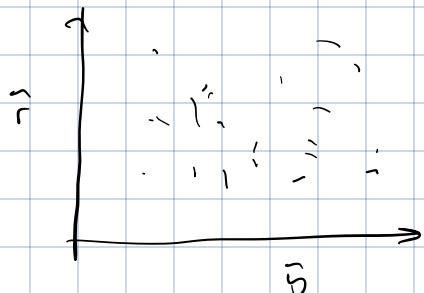
↳ H₀: $Y_1, \dots, Y_n \sim N(\dots, \dots)$

↳ Low p-val \rightarrow evidence of non-normality. If very high p-value

② Constant mean & variance

Graphical:

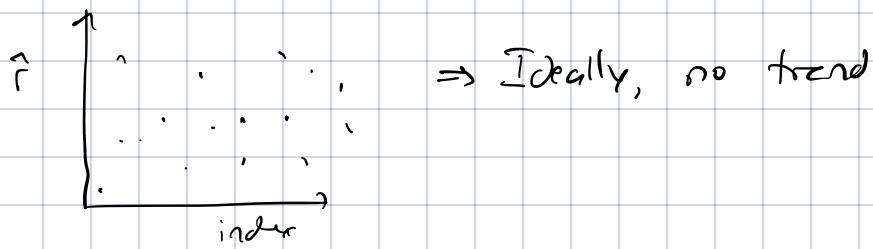
A: residual vs. fitted values



\Rightarrow Ideally, no trend

If funnel-shaped \Rightarrow variance not constant

B: Residuals vs. index/time (only if data order known)



Formal test: non-parametric Fligner-Killeen test for non-constant variance

$$\hookrightarrow H_0: \sigma_1^2 = \dots = \sigma_k^2$$

↳ Test divides data into k segments & checks variance for each segment

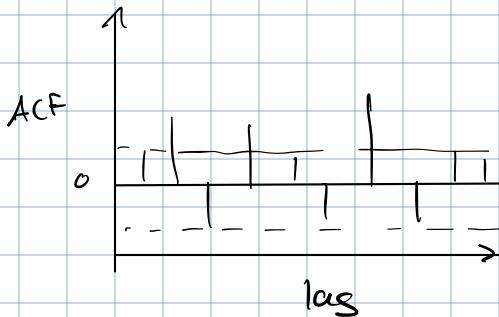
↳ Low-pvalue \rightarrow evidence of non-constant variance

③ Uncorrelatedness

If correlated \Rightarrow dependent. Independence \Rightarrow uncorrelated, uncorrelated $\xrightarrow{\text{joint normality}}$ indep.

Test to check if residuals are correlated \rightarrow not indep.

Residuals' sample ACF: 95% of spikes should be \in confidence bands



\Rightarrow Spikes should also be in random lags

If trend \Rightarrow dependence

Spikes at lags = 1, 2 \rightarrow concern

④ Randomness

A: Difference sign test:

$$1. S = \# \text{ of values s.t. } y_i - y_{i-1} > 0$$

$$2. \text{ For large } n, S \sim N\left(\mu_S = \frac{n-1}{2}, \sigma_S^2 = \frac{n-1}{12}\right). \text{ Thus}$$

$$\frac{S - \mu_S}{\sqrt{\sigma_S^2}} \sim N(0, 1)$$

$$3. H_0: \text{data is random} \rightarrow \text{run a test.}$$

Doesn't work if data has a lot of ~~seasonality~~

$y_i = y_{i-1} \Rightarrow$ exclude

False negative problem: test doesn't consider order of sign increase/decrease.



B: Runs signs test

1. Estimate median m

2. $n_1 := \text{num. of obs} > m$, $n_2 := \text{num. of obs. } < m$

3. $R := \# \text{ of consec. seq. of data where } \underline{\text{all}} < m$

$$4. \frac{R - \mu_R}{\sigma_R} \sim N(0, 1) \Rightarrow \mu_R = 1 + \frac{2n_1 n_2}{n_1 + n_2}$$

$$\sigma_R^2 = \frac{(\mu_R - 1)(\mu_R - 2)}{n_1 + n_2 - 1}$$

5. Run test w/ $H_0: \text{data is random}$

False negative still a problem \Rightarrow clustered sequences not weighted highly

Run many tests \rightarrow if any show problem, assumptions may be faulty.

PREDICTION WITH LINEAR REGRESSION

Forecasting w/ Multiple L.R. Models

Divide dataset into 2:

(1) Training set: data for model fitting

(2) Testing set: evaluate model perf.

Prediction interval theorem: (y_0, x_0) \notin training set. Define point forecast of y_0 at $x = x_0$ as $\hat{y} = x_0^T \hat{\beta}$

$$100(1-\alpha)\% \text{ prediction interval} = \hat{y} \pm c_{\alpha} \hat{\sigma} \sqrt{1 + x_0^T (x^T x)^{-1} x_0}$$

Std. error: estimate
std. dev. of estimator

$(1-\frac{\alpha}{2})$ quantile of t_{n-p-1}
 $\hat{\sigma}^2$: MSE of model

Proof sketch:

$$y_0 - x_0^T \hat{\beta} \sim N(0, \sigma^2 (1 + x_0^T (x^T x)^{-1} x_0))$$

$$\therefore T := \frac{y_0 - x_0^\top \hat{\beta}}{\hat{\sigma} \sqrt{1 + \dots}} \sim t_{n-p-1} \text{ if } \hat{\sigma} \text{ is estimate}$$

Thus:

$$P(-c_\alpha \leq T \leq c_\alpha) = 1 - \frac{\alpha}{2}$$

Remembering:

$$x_0^\top \hat{\beta} \pm c_\alpha \hat{\sigma} \sqrt{1 + \dots}$$

Q: What can we do to reduce the width of prediction interval?

Components

c_α ↓ → less confidence
in our interval

Inherent var.
in data → can't
control

$$\sqrt{1 + x_0^\top (X^\top X)^{-1} x_0}$$

$$(X^\top X)^{-1} = \frac{1}{\det(X^\top X)} \text{adj}(X^\top X)$$

Increased $\det(X^\top X) \rightarrow$
P.I. length ↓.

Multicollinearity →
 $\det(X^\top X) \downarrow$

Multicollinearity

Defn: two columns of X are highly correlated

Impacts:

① $\text{Var}(\tilde{\beta})$ is high

$$\text{Var}(\tilde{\beta}) = \sigma^2 (X^\top X)^{-1} = \sigma^2 \cdot \frac{\text{adj}(X^\top X)}{\det(X^\top X)}$$

Since multicollinearity \Rightarrow small $\det(\dots)$, $\text{Var}(\tilde{\beta})$ is high

$\hat{\beta}$ is not a stable estimate

② Prediction interval width

$$\text{P.I.} = \hat{\mu} \pm c_\alpha \hat{\sigma} \sqrt{1 + x_0^\top (X^\top X)^{-1} x_0}$$

↓
High value \rightarrow interval ↑

Identification

① Check corr. b/w explanatory variables

- ② $\text{SE}(\hat{\beta}_j)$ is high
- ③ Overall test of sig. does not match test of sig. of $\hat{\beta}_j$
- ④ $\hat{\beta}_j$ sign is not expected
- ⑤ Variance inflation factor

For some variate j , regress j on all other variates & extract R^2 .

$$\text{VIF}(\hat{\beta}_j) = \frac{1}{1 - R_{x_j|x_{-j}}^2}, j=1, \dots, p$$

↳ If $R^2 \uparrow$, VIF \uparrow

$$\text{VIF}(\hat{\beta}_j) \geq 5 \Rightarrow \text{collinearity}$$

Q: How to deal w/ this?

- ① Drop variables
- ② Regularization models

Interpolation & Extrapolation

Idea: find range of validity where prediction makes sense

$$h_{\max} := \max(H_{ii}) \Rightarrow \text{convex hull of data}$$

Assume predicting on point x .

$$x^T (x^T x)^{-1} x \leq h_{\max} \Rightarrow \text{interpolation}$$

Why is extrapolation bad? Assumptions may not hold for data outside of range

BIAIS-VARIANCE TRADEOFF

Bias-variance decomposition:

Assume we are estimating on $x = x_0$, $\hat{y}_0 = x_0' \hat{\beta}$:

$$\text{MSE}(x_0) = \underbrace{(\hat{y}_0 - \mathbb{E}(\hat{y}_0))^2}_{\text{Bias}(\hat{y}_0)^2} + \underbrace{\mathbb{E}((\hat{y}_0 - \mathbb{E}(\hat{y}_0))^2)}_{\text{Var}(\hat{y}_0)}$$

Bias: error of $f(x)$ (model) not belonging in actual model space

Variance: variability of params \leftarrow new data

Expected prediction error:

Let $Y = f(x) + \epsilon$, $\text{Var}(\epsilon) = \sigma_\epsilon^2$. We choose a model $\hat{f}(x)$. We fit value x_0 .

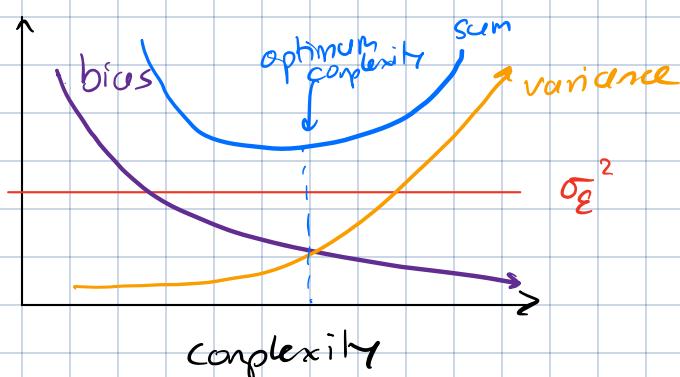
$$\begin{aligned}\text{Err}(x_0) &= E[(Y - \hat{f}(x_0))^2] \\ &= \sigma_\epsilon^2 + \underbrace{\text{Bias}^2}_{\text{Choosing } \hat{f}(x_0) \text{ instead } f(x)} + \text{Variance}\end{aligned}$$

Irreducible error: error inherent in data that we model (ϵ term)

Irreducible error not in control, others depend on model selection

The more complex the model, more variability but lower bias

" simple " , less " higher bias "



VARIABLE SELECTION

R^2 & Adjusted R^2

Defn of R^2 :

$$R^2 = \frac{SS_{reg}}{SS_{total}} = 1 - \frac{SSE}{SS_{Total}}$$

Issue: adding more variables will always increase $R^2 \rightarrow$ overfitting

Fix: adjusted R^2

$$\bar{R}^2 = 1 - (1 - R^2) \cdot \frac{n-1}{n-p-1}$$

- Penalizes large, complex models

Akaike's Information Criterion

Defn: $AIC = -2 \ell(\hat{\beta}) + 2N_p$ → # of model params

The lower the AIC, the better

This measures model fit, not pred.

Issue: $n \downarrow \rightarrow AIC$ doesn't work well

- Corrected AIC: $AIC_c = AIC + \frac{2N_p(N_p+1)}{n-N_p+1}$
→ 0 as $n \uparrow$

Bayesian Information Criteria

Defn: $BIC = -2 \ell(\hat{\beta}) + \log(n) \cdot N_p$ → Penalizes harsher than AIC

Same as AIC in interpretation & usage

Stepwise model selection

Either add more variables / subtract variables one by one → choose best AIC / BIC model

Average Prediction Squared Error (APSE)

Unlike others, this is a measure of predictive power

$$APSE = \frac{\sum_v (y - \hat{y})^2}{|v|} \quad (v \text{ is test set})$$

Smaller APSE → better forecasting power

k-fold cross validation

1st fold:

Test	Train	Train	Test
------	-------	-------	------

2nd fold:

Train	Test	Train	Test
-------	------	-------	------

Partition data into k disjoint & ~equal sets

After each fold:

$$sMSE_i = \frac{1}{n_k} \sum_{j \in T_i} (y_j - \hat{f}_i(x_j))^2$$

\hat{f}_i fitting on T_{-i}

This is just an example of a loss func, can be anything

Then, calculate:

$$CV(k) = \frac{1}{k} \sum_{i=1}^k sMSE_i$$

How to choose k ? Large k is low bias & high variance + more comp. time

Randomness in creating k -folds \Rightarrow won't get same ans.

Model selection techniques

① Subset selection

Create subsets of variables (2^p models) \rightarrow train \rightarrow evaluate on metrics.

Comp. expensive

② Shrinkage

Penalize loss function on # of variables (known as regularization)

Effect: some $\beta \rightarrow 0$, $\text{Var}(\beta) \downarrow$

③ Dimensionality reduction

Project p variables \rightarrow M -dimensional space ($M < p$) \Rightarrow dimensions incorporate all p variables in L.C.

REGULARIZATION METHODS

$$\tilde{\beta}_{\text{regulariz}} = \arg \min_{\beta} \left\{ (Y - X\beta)'(Y - X\beta) + \lambda \underbrace{\text{Pen}(\beta)}_{>0} \right\}$$

Pros:

- + multicollinearity ↓
- + $n > p$ not an issue
- + variable selection
- + $\text{Var}(\tilde{\beta}_{\text{regulariz}}) < \text{Var}(\tilde{\beta}_{\text{OLS}})$

Cons:

- + Bias: $E[\tilde{\beta}_{\text{reg.}}] \neq \beta$ can bias correct
- + Comp. expensive
- + Optimality of β not achieved

Ridge Regression

$$\text{Defn: } \text{Pen}(\beta) = \ell_2 = \sum_{j=1}^p \beta_j^2 \rightarrow \beta_0 \text{ not penalized}$$

Equivalent form:

$$(\beta_0, \hat{\beta}_{\text{ridge}}) := \arg \min_{\beta_0, \beta} \sum_{i=1}^n (y_i - \beta_0 - \sum_{j=1}^p x_{ij} \beta_j)^2, \text{ s.t. } \sum_{i=1}^p \beta_j \leq t \quad (t \leftarrow \lambda)$$

Solving minimization.

$$\text{RSS}(\lambda) = (y - \beta_0 \vec{1} - X\beta)'(y - \beta_0 \vec{1} - X\beta) + \lambda \beta' \beta$$

$$\hat{\beta}_{\text{ridge}} = (X'X + \lambda I)^{-1} X' Y, \quad \hat{\beta}_0 = \bar{y} - \frac{1}{n} \sum_{i=1}^n \hat{\beta}_j x_{ij}$$

↓
diagonal shifted
by λ

Since $(X'X + \lambda I)$ is invertible for some λ , multicollinearity not an issue

Why do we not penalize β_0 ? β_0 should absorb constant changes to $Y \rightarrow$ no change to $\hat{\beta}_{\text{ridge}}$ soln.

↳ Can drop β_0 : Replace $x_{ij} \rightarrow x_{ij} - \bar{x}_j, j = 1 \dots p \Rightarrow \hat{\beta}_0 = 5. y \rightarrow y - \bar{y} \rightarrow \beta_0 = 0$.

We need to standardize predictors s.t. all contrib. to penalty is equal.

$$\hookrightarrow \frac{x_{ij} - \bar{x}_j}{\sigma_x}$$

The larger $-\lambda$, the more shrinkage. Does not set $\beta \rightarrow 0 \rightarrow$ no variable selection.

Lasso Regression

$$\text{Defn: Pen}(\beta) = \ell_1 = \sum_{j=1}^p |\beta_j|$$

Minimization problem

$$(\hat{\beta}_0, \hat{\beta}_{\text{lasso}}) := \arg \min_{\beta_0, \beta} \left(\sum_{i=1}^n (y_i - \beta_0 - \sum_{j=1}^p x_{ij} \beta_j)^2 + \lambda \sum_{j=1}^p |\beta_j| \right)$$

Full form: least absolute selection and shrinkage operator (does variable sel.)

Compared to ridge:

① Harder to minimize b/c abs. val.

② No closed form soln for $\hat{\beta}_{\text{lasso}}$

④ Good if few predictors w/ outsized effects

⑤ Chooses collinear vars randomly \rightarrow not

③ Variable selection

snow for interpretation

$\lambda \uparrow \rightarrow$ more shrinkage

How to choose λ ? C.V.

Elastic net regression

$$\hat{\beta}_{\text{en}} := \underset{\beta}{\arg \min} \left\{ \sum_{i=1}^n (y_i - \beta_0 - \sum_{j=1}^p x_{ij} \beta_j)^2 + \lambda \sum_{j=1}^p [(1-\alpha) \beta_j^2 + \alpha | \beta_j |] \right\}, \quad \alpha \in [0, 1], \lambda > 0$$

ridge lasso

How to choose α & λ ? Grid search (coarse interval for α , fine for λ)

STATIONARY PROCESS

Ideal time series model is finite dimensional distributions of process $\{X_t, t \in T\}$

- Finite dimensional distn.: joint distn. of finite subset of time series

Stationarity: finite dimensional distn. unchanged under time shift

Strong stationarity

Defn: \forall finite subsets $\{t_1, \dots, t_n\} \subseteq T$ & $\forall h$, $F(X_{t_1}, \dots, X_{t_n}) = F(X_{t_1+h}, \dots, X_{t_n+h})$

Ex:// Show any seq. of IID R.V. is strictly stationary

$$\begin{aligned} F_{x_1, \dots, x_n}(x_1, \dots, x_n) &= \prod_{i=1}^n F_{x_i}(x_i) \\ &= \prod_{i=1}^n F_{x_{i+h}}(x_i) \\ &= F_{t_1+h, \dots, t_n+h}(x_1, \dots, x_n) \end{aligned} \quad \left. \begin{array}{l} x_i \stackrel{D}{=} x_j \end{array} \right.$$

Practically, strict stationarity is v. hard to prove & rarely true

Weak stationarity

Mean function: $\mu_x(t) = E[X_t]$, $t \in T$

Auto-covariance function (ACVF): $\gamma(r, s) = \text{Cov}(X_r, X_s)$ (assumes all $\text{Var}[X_t] < \infty$)
 $= E(X_r X_s) - E(X_r)E(X_s)$

Ex:// Derive mean & ACVF for $X_t = 2 + 3t + Z_t - 0.5Z_{t-1}$, $Z_j \stackrel{iid}{\sim} N(0, \sigma^2)$

Mean function:

$$\mu(t) = E[X_t] = 2 + 3t \quad \forall t = 1, 2, \dots$$

ACUF:

$$\begin{aligned}\gamma(r, s) &= \text{Cov}(2 + 3r + Z_r - 0.5Z_{r-1}, 2 + 3s + Z_s - 0.5Z_{s-1}) \\ &= \text{Cov}(Z_r - 0.5Z_{r-1}, Z_s - 0.5Z_{s-1}) \quad \begin{matrix} \text{Constants don't} \\ \text{matter} \end{matrix} \\ &= \text{Cov}(Z_r, Z_s) - 0.5 \text{Cov}(Z_r, Z_{s-1}) - 0.5 \text{Cov}(Z_{r-1}, Z_s) + 0.25 \text{Cov}(Z_{r-1}, Z_{s-1})\end{aligned}$$

Case analysis: For each term, check when time index is equal

$$\textcircled{1} \quad r = s$$

$$\gamma(r, s) = \sigma^2 + 0 + 0 + 0.25\sigma^2 = 1.25\sigma^2$$

$$\textcircled{2} \quad r = s-1$$

$$\gamma(r, s) = -0.5\sigma^2$$

$$\textcircled{3} \quad s = r-1$$

$$\gamma(r, s) = -0.5\sigma^2$$

Thus:

$$\gamma(r, s) = \begin{cases} 1.25\sigma^2 & |r-s| = 0 \\ -0.5\sigma^2 & |r-s| = 1 \\ 0 & |r-s| > 1 \end{cases}$$

Weak stationarity defn: $\{X_T\}$ is weakly stationary if:

$$\textcircled{1} \quad E(X_t^2) < \infty \quad \forall t \in T \quad (\text{or } \text{Var}(X_t) < \infty)$$

$$\textcircled{2} \quad E(X_t) = \mu \quad \forall t \in T \Rightarrow \text{indep. of } t$$

$$\textcircled{3} \quad \forall r, s, r+t, s+t \in T, \quad \gamma(r, s) = \gamma(r+t, s+t) \Rightarrow$$

$\gamma(r, s) \text{ & } \gamma(r+t, s+t)$
only function of $|r-s|$

If $\textcircled{3}$ is true, then ACVF at lag h is:

$$\gamma_x(h) := \gamma_x(h, 0) = \gamma_x(t, t+h) = \gamma_x(t+h, t)$$

Corollary: Strictly stationary process \Rightarrow weakly stationary, converse not true!

• Proof: Assume $\{X_t\}$ is strictly stationary w/ $E(X_t^2) < \infty$. Proving weak:

$$\textcircled{1} \quad E(X_t^2) < \infty \text{ by defn.}$$

$$\textcircled{2} \quad \text{If } \{X_t\} \text{ is strictly stationary, then } F_{X_1}(x) = \dots = F_{X_n}(x). \quad \forall x.$$

$$\Rightarrow F_{X_t}(x) = F_{X_0}(x) \quad \forall x, t$$

$$\Rightarrow E[X_t] = E[X_0]$$

$$\Rightarrow E[X_t] = \mu$$

$$\textcircled{3} \quad \gamma(h) = \text{Cov}(X_t, X_{t+h})$$

$$= \vdots$$

$$= \text{Cov}(X_1, X_{h+1})$$

) Prove. Show
 $X_t, X_{t+h} \stackrel{D}{=} X_1, X_{h+1}$

Auto-correlation function (ACF):

$$\rho_X(h) = \frac{\gamma_X(h)}{\gamma_X(0)} = \text{Cor}(X_{t+h}, X_t)$$

↳ $\text{Var}(X_t)$

$\left. \begin{array}{l} \text{ACVF can} \\ \text{be scaled up/down.} \end{array} \right\}$

Use b/c

Ex:// Is $\{Z_t\} \sim \text{WN}(0, \sigma^2)$ stationary?

$$\textcircled{1} \quad \text{Var}(Z_t) = \sigma^2 < \infty$$

$$\textcircled{2} \quad E(Z_t) = 0 \Rightarrow \text{not a func. of time}$$

$$\textcircled{3} \quad \gamma(h) = \begin{cases} 0 & h \neq 0 \\ \sigma^2 & h = 0 \end{cases} \Rightarrow \text{not a func. of time.}$$

Concl. Yup! $\{Z_t\}$ is stationary

Ex:// Is random walk stationary?

Recall that $\{S_t\}$ is random walk if $S_t = \sum_{i=1}^t X_i$, X_i is iid noise.

Consider 3rd cond.:

$$\gamma(t, t+h) = \text{Cov}(S_{t+h}, S_t)$$

$$= \text{Cov}\left(\sum_{i=1}^t X_i + \sum_{i=t+1}^{t+h} X_i, \sum_{i=1}^t X_i\right)$$

$$= \text{Cov}\left(\sum_{i=1}^t X_i, \sum_{i=1}^t X_i\right) + \text{Cov}\left(\sum_{i=t+1}^{t+h} X_i, \sum_{i=1}^t X_i\right)$$

$$= \sum_{i=1}^t \text{Cov}(X_i, X_i)$$

$$= t\sigma^2$$

O b/c no overlap

Since $\gamma(h) \propto t$, 3rd cond. is violated

Ex:// Show that $X_t = 5 + 4t + Z_t$, Z_t is white noise is not stationary.

From cond. 2 $\Rightarrow E[X_t] = 5 + 4t \rightarrow$ dependent on time.

Exercise: show cond. 1 & 3 are met

Ex:// Example of strict stationary time series which is not weakly stationary

Intuition: strictly stationary $\xrightarrow{\text{Var}(x_t) < \infty}$ weakly stationary

$\{X_t\}$ is iid. Cauchy. Strictly stationary b/c finite distn. or identical

Not weakly stationary b/c variance is infinite.

Ex:// Example of weak stationary process that is not strictly stationary

Intuition: we want $\{X_t\}$ whose distn. changes under time shift.

Let $\{X_t: t \geq 1\}$ be:

$$X_t = \begin{cases} Z_t, & t \text{ is even} \\ \frac{1}{\sqrt{2}}(Z_{t-1}^2 - 1), & t \text{ is odd} \end{cases}, Z_t \stackrel{iid}{\sim} N(0, \sigma^2)$$

Proving weak stationarity:

$$\textcircled{1} \quad \text{Var}(X_t) = \begin{cases} \sigma^2 & t \text{ even} \\ \frac{1}{2} \text{Var}(Z_{t-1}^2) & t \text{ odd} \end{cases} = \begin{cases} \sigma^2 & t \text{ even} \\ \frac{1}{2} \cdot 2 = 1 & t \text{ odd} \end{cases}$$

$\downarrow \sim X_1^2, \text{Var} = 2 \cdot \text{df.}$

$$\textcircled{2} \quad E[X_t] = \begin{cases} 0 & \text{E of } X_1^2 \text{ is 1} \\ \frac{1}{\sqrt{2}}(E(X_{t-1}^2) - 1) & \end{cases} = 0$$

$$\textcircled{3} \quad \gamma_x(h) = 0 \quad \forall h \neq 0 \quad \text{b/c no overlap}$$

To prove this is not strongly stationary, consider prob. func. of X_2 & X_3 .

$t=2$

$t=3$

$$P(X_2 \leq 0) = P(Z_2 \leq 0) = \frac{1}{2}$$

$$P(X_3 \leq 0) = P(\frac{1}{\sqrt{2}}(Z_2^2 - 1) \leq 0)$$

$$= P(Z_2^2 \leq 1)$$

$$\neq \frac{1}{2}$$

Univariate distn. are diff. under time shift \rightarrow not strictly stationary.

Ex:// $X_t = Z_t + \theta Z_{t-1}$, $Z_t \sim WN(0, \sigma^2)$. Show that $\{X_t\}$ is weakly stationary.

$$\begin{aligned} \textcircled{1} \quad \text{Var}(X_t) &= \text{Var}(Z_t + \theta Z_{t-1}) \\ &= \text{Var}(Z_t) + \theta^2 \text{Var}(Z_{t-1}) + 2\theta \text{Cov}(Z_t, Z_{t-1}) \xrightarrow{\text{O b/c no overlap}} \\ &= \sigma^2 + \theta^2 \sigma^2 \leq \infty \end{aligned}$$

$$\textcircled{2} \quad E(X_t) = E(Z_t) + \theta E[Z_{t-1}] = 0, \text{ doesn't depend on } t$$

$$\begin{aligned} \textcircled{3} \quad \text{Cov}(X_t, X_{t+h}) &= \text{Cov}(Z_t + \theta Z_{t-1}, Z_{t+h} + \theta Z_{t+h-1}) \\ &= \text{Cov}(Z_t, Z_{t+h}) + \theta \text{Cov}(Z_t, Z_{t+h-1}) \\ &\quad + \theta \text{Cov}(Z_{t-1}, Z_{t+h}) + \theta^2 \text{Cov}(Z_{t-1}, Z_{t+h-1}) \end{aligned}$$

Cov analysis:

$$\gamma(h) = \begin{cases} \sigma^2 + \theta^2 \sigma^2, & h=0 \\ \theta \sigma^2, & h=1 \\ \theta \sigma^2, & h=-1 \\ 0, & \text{o.v.} \end{cases} \Rightarrow \text{Doesn't depend on } t$$

①, ②, ③ $\Rightarrow \{X_t\}$ is weakly stationary.

Properties of ACVF & ACF

Assume $\gamma(r, s)$ is ACVF of stationary process $\{X_t\}$. Then:

$$\textcircled{1} \quad \forall r, \quad \gamma(h) := \gamma(r+h, r) \rightarrow r \text{ doesn't matter}$$

$$\textcircled{2} \quad \gamma(0) \geq 0 \rightarrow \text{Cov}(X_r, X_r) = \text{Var}(X_r) \geq 0$$

$$\textcircled{3} \quad |\gamma(h)| \leq \gamma(0) \rightarrow |\gamma(h)| \leq 1 \Rightarrow \frac{|\gamma(h)|}{\gamma(0)} \leq 1 \Rightarrow |\gamma(h)| \leq \gamma(0)$$

$$\textcircled{4} \quad \text{ACVF is even: } \gamma(h) = \gamma(-h) \Rightarrow \gamma(h) = \text{Cov}(X_t, X_{t+h}) = \text{Cov}(X_t, X_{t-h})$$

Q: How do we go about modelling stationary process?

We want to keep model simple yet useful.

Eg:// If we tried to model w MVN under stationarity, variance-covariance matrix $\in O(n^2)$ params. \rightarrow more data, more params. Too complicated.

Examples: MA(1), AR(1) \rightarrow will see later.

Estimating ACF

In real data, we don't have model to construct ACF; need to use actual data.

Sample ACVF:

$$\hat{\gamma}(h) = \frac{1}{n} \sum_{t=1}^{n-|h|} (x_{t+|h|} - \bar{x})(x_t - \bar{x})$$

$\frac{1}{\text{Cov}(x_{t+|h|}, x_t)}$

Sample ACF:

$$\hat{\rho}(h) = \frac{\hat{\gamma}(h)}{\hat{\gamma}(0)}$$

If $\{x_t\}$ is i.i.d. noise w/ finite variance

$$\hat{\rho}(h) \sim N(0, 1/n)$$

so,

$$\text{C.I.}_{95\%} = \hat{\rho}(h) \pm 1.96 \cdot \sqrt{1/n}$$

If data has trend, $|\hat{\rho}(h)|$ has slow decay.

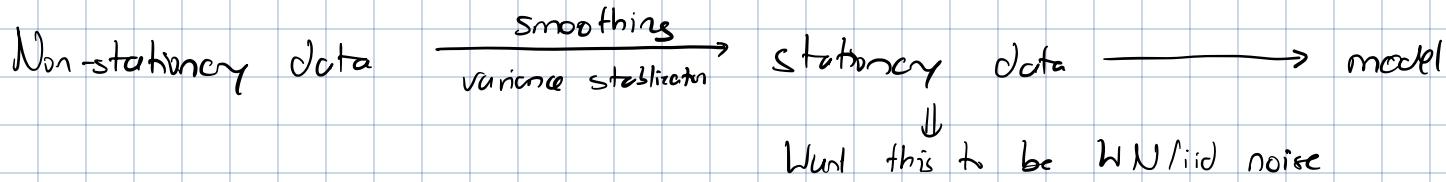
If data has periodicity, $|\hat{\rho}(h)|$ will also have periodicity w/ same period.

Note: it is fine for ACF plot to have 5% of values $\geq \text{C.I.}$, but must be in random order!

SMOOTHING METHODS

Defn: estimating patterns while leaving out noise

$\xrightarrow{\text{deterministic}}$



Trend estimation

Estimate m_t & $s_t \rightarrow$ analyze & model residuals

① Finite moving avg. filter:

Assume $X_t = m_t + Y_t$ (Y_t is random component), q is non-neg. #

$$W_t := \frac{1}{2q+1} \sum_{j=-q}^q X_{t-j}$$

$$= \frac{1}{2q+1} \sum_{j=-q}^q m_{t-j} + \frac{1}{2q+1} \sum_{j=-q}^q Y_{t-j}$$

$W_t \approx m_t$

Law of large #s: Sample mean \rightarrow pop. mean.
Avg. of $Y_{t-j} \rightarrow 0$.

Intuition: avg. of q points b/f & after each X_t datapoint

Effect of q : $q \uparrow \rightarrow$ bias \uparrow , variance \downarrow

Problem: can't use in forecasting b/c needs future data

↳ Can't use C.V. to find optimal q

② Exponential smoothing

Defn:

$$\begin{cases} \hat{m}_t = \alpha X_t + (1-\alpha) \hat{m}_{t-1} \\ \hat{m}_1 = x_1 \end{cases}$$

↓

$$\hat{m}_t = \sum_{j=0}^{t-2} \alpha(1-\alpha)^j X_{t-j} + (1-\alpha)^{t-1} x_1$$

Intuition: estim. trend is avg. of current data point & previous estim. trends exponentially weighted

Effect: $\alpha \uparrow \rightarrow$ bias \downarrow , variance \uparrow

Finding α : 1. $\hat{\alpha} = \min_{\alpha} (\sum (x_t - \hat{x}_t)^2)$

2. leave-one-out cross-validation

③ Polynomial regression

1. Make model assumption
 2. Fit model
 3. Model selection
- None of this requires $\epsilon \sim N(0, \sigma^2)$
- assumption!
- Try your best to make residuals have
constant mean of 0

Trend Elimination (Differencing)

Important note: Cannot help w/ non-constant variance, only removes trend & seasonality

Backwards shift operator: $B X_t := X_{t-1}$

$$B^k X_t := X_{t-k}$$

Differencing operator:

$$\nabla X_t = (1 - B) X_t = X_t - X_{t-1}$$

$$\nabla^k X_t = (1 - B)^k X_t \Rightarrow \text{Differencing } k \text{ times}$$

How many times should we difference (k ?)

- ↳ More you difference \rightarrow the more variability you introduce, model more complicated
- ↳ Only want to difference until trend/seasonality removed.

Ex:// $Y_t = 1 + t + Z_t$, $Z_t \sim WN(0, \sigma^2)$, $\sigma^2 < \infty$.

1) Is Y_t stationary

$$E[Y_t] = 1 + t \rightarrow \text{dependent on } t \rightarrow \text{non-stationary}$$

2) Is ∇Y_t stationary

$$W_t = \nabla Y_t = (1 + t + Z_t) - (1 + t - 1 + Z_{t-1})$$

$$= Z_t - Z_{t-1} + 1$$

$$l. \quad \text{Var}(W_t) = \text{Var}(Z_t) + \text{Var}(Z_{t-1}) - 2\text{Cov}(Z_t, Z_{t-1})$$

$$= \sigma^2 + \sigma^2$$

$$= 2\sigma^2 < \infty$$

σ b/c WN
(no corr.)

2. $E(U_t) = 1 \rightarrow$ not dependent on t

3. $\text{Cov}(U_t, U_{t+h}) = \text{Cov}(Z_t - Z_{t-1}, Z_{t+h} - Z_{t+h-1} + 1)$
⋮
⋮
Not a function of t

1, 2, 3 \rightarrow non-stationary

Theoretical result: if trend is polynomial of degree p , $\nabla^p X_t$ will eliminate trend

Can also eliminate seasonality by differencing in order of seasonal period.

↳ Ex:// If X_t has period of $d=6 \rightarrow X_t - X_{t-6}$ will eliminate seasonality
 $(1 - B^6) X_t$

↳ Notation: $\nabla^k X_t \rightarrow$ difference k times, lag of 1 $\Rightarrow (X_t - X_{t-1})^k$
 $\nabla_k X_t \rightarrow$ difference 1 time, lag of $k \Rightarrow (X_t - X_{t-k})$
 $= (1 - B^k) X_t$

Ex:// $X_t = \sin\left(\frac{\pi t}{10}\right) + Z_t, Z_t \sim WN(0, \sigma^2)$.

1. Is X_t stationary?

$E[X_t] = \sin\left(\frac{\pi t}{10}\right) \rightarrow$ dependent on t , non-stationary

2. $\nabla_{20} X_t$ stationary?

Aside: What is period of X_t ?

$$\sin\left(\frac{\pi(t+d)}{10}\right) = \sin\left(\frac{\pi t}{10}\right)$$

$$\sin\left(\frac{\pi t}{10} + \frac{\pi d}{10}\right) = \sin\left(\frac{\pi t}{10}\right)$$

$$\frac{\pi d}{10} = 2\pi$$

$$d = 20$$

$$\nabla_{20} X_t = \underbrace{\sin\left(\frac{\pi t}{10}\right) - \sin\left(\frac{\pi(t-20)}{10}\right)}_{\text{Equiv.}} + Z_t - Z_{t-20}$$

$$= Z_t - Z_{t-20}$$

⋮ Show that this is stationary

If you have both trend & seasonality, try regular & seasonal diff. separately, and only combine if neither works.

Decompose

We introduce a method to remove trend & seasonality \Rightarrow Only for EDA

Decompose function in R does following:

- ① Estimate trend: Apply moving avg. filter w/ $q = d/2$ (d is period of time series we model).

$$\hat{m}_t = \begin{cases} \frac{0.5x_{t-q} + x_{t-q+1} + \dots + x_{t+q-1} + 0.5x_{t+q}}{d}, & d = 2q \\ \frac{1}{2q+1} \sum_{i=-q}^q x_{t-i}, & d = 2q+1 \end{cases}, \quad \begin{array}{l} \text{Even period} \\ q \leq t \leq n-q \end{array} \quad \begin{array}{l} \text{odd period} \end{array}$$

Some period, so sum x_{t+i}

- ② Detrend data: $x_t - \hat{m}_t$

- ③ Estimate seasonality from detrended data

Assumptions for seasonality:

1. $S_t = S_{t+d} \Rightarrow d$ is period

2. $\sum_{j=1}^d S_j = 0$

Process:

1. Take avg. of every period value in detrended data $\Rightarrow \{w_k\}$

$E_S //$ Avg. for Jan., Feb., ..., Dec.

$k=1, \dots, d$

2. Centralize each to overall avg.

$$\hat{s}_k = w_k - \frac{\sum w_j}{d} \Rightarrow \begin{array}{l} \text{Guarantees assump.} \\ 2 \text{ is fulfilled} \end{array}$$

- ④ Model $x_t - \hat{m}_t - \hat{s}_k$ (residuals)

Special note: decompose will always estimate season., it's not always relevant

\hookrightarrow Signature seasonality not important: scale of seasonal plot is \approx scale of residuals.

Decomposition model can also be used for multiplicative models ($X_t = m_t s_t r_t$)

↳ Outcomes: avg. of seasonal is 1, not 0

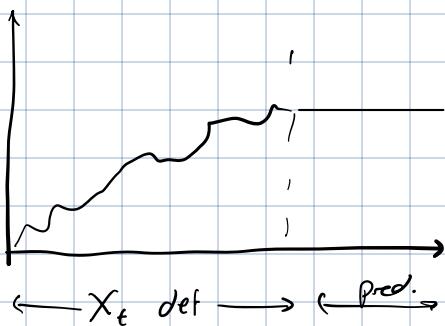
Another note: freq. in time series model may not be correct. Look @ ACF, det. correct period & change the series for proper decomposition.

Holt - Winters

Problems w/ exp. smoothing

- ① Cannot accommodate seasonality
- ② Prediction results in constant value predictor

Obj. = generalize exp. smoothing & solve this



Also (additive):

$$\begin{cases} \text{level} \rightarrow L_t = \alpha(X_t - T_{t-p}) + (1-\alpha)(L_{t-1} + T_{t-1}) \\ \text{trend} \rightarrow T_t = \beta(L_t - L_{t-1}) + (1-\beta)T_{t-1} \\ \text{seasonal} \rightarrow I_t = \gamma(X_t - L_t) + (1-\gamma)I_{t-p} \end{cases} \quad \text{Exp. Smoothing}$$

where L_0, T_0 & I_0 are defined.

↳ No need to memorize

Forecast of h periods ahead:

$$\hat{X}_{t+h} = L_t + hT_t + I_{\text{period of } t+h}$$

Also (multiplicative):

$$\begin{cases} L_t \\ T_t \\ I_t \end{cases} \Rightarrow \text{Forecast: } \hat{X}_{t+h} = (L_t + hT_t) I_{\text{period of } t+h}$$

In R: need starting values & $\{\alpha, \beta, \gamma\}$ → R does this automatically

Special cases:

① $\beta = \gamma = 0$ (regular exp. smoothing).

$$\hookrightarrow \hat{x}_{t+h} = l_t$$



② $\gamma = 0$ (no season, double exp. smoothing)

$$\hookrightarrow \hat{x}_{t+h} = l_t + hT_t$$



Try out both additive & multiplicative \rightarrow see which is better

IMPORTANT NOTE: No stat. Assumptions needed!

\hookrightarrow Needed only for prediction interval:

$$\hat{x}_{t+h} \pm z_{\alpha/2} \text{SE}(\hat{x}_{t+h} - \bar{x}_{t+h})$$

STATIONARITY & LINEAR PROCESSES

We will use stationarity as an assumption when developing models

Moving Average Process: MA(q)

Defn: $Z_t \sim WN(0, \sigma^2)$, θ_i are non-zero constants. Then:

$$X_t = Z_t + \theta_1 Z_{t-1} + \dots + \theta_q Z_{t-q}$$

\Rightarrow Not to same
as MA filter!

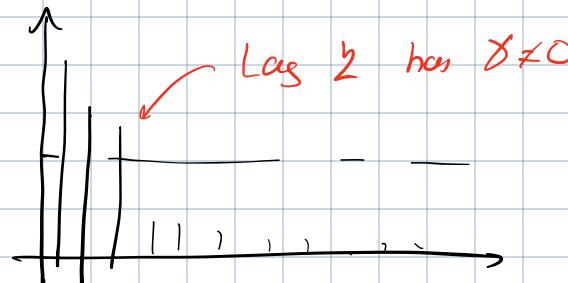
Q: How do you know what q should be?

A: Check ACF! q^{th} lag will be last non-zero lag

Ex://



ACF of MA(1)



ACF of MA(2)

Theorem: $MA(q)$ is stationary

$$\begin{aligned} \textcircled{1} \quad \text{Var}[X_t] &= \text{Var}(Z_t + \theta_1 Z_{t-1} + \dots + \theta_q Z_{t-q}) \\ &= \sigma^2 + \theta_1^2 \sigma^2 + \dots + \theta_q^2 \sigma^2 \\ &= \sigma^2 (1 + \theta_1^2 + \dots + \theta_q^2) < \infty \text{ b/c } \theta_i \text{ & } \sigma \text{ are finite} \end{aligned}$$

$$\begin{aligned} \textcircled{2} \quad E[X_t] &= E[Z_t + \theta_1 Z_{t-1} + \dots + \theta_q Z_{t-q}] \\ &= 0 \Rightarrow \text{indep. of } t \text{ & finite} \end{aligned}$$

$$\textcircled{3} \quad \text{Cov}(X_t, X_{t+h}) = \text{Cov}(\underbrace{Z_t + \theta_1 Z_{t-1} + \dots + \theta_q Z_{t-q}}_i, \underbrace{Z_{t+h} + \theta_1 Z_{t+h-1} + \dots + \theta_q Z_{t+q}}_i)$$

↳ Case analysis:

$$\left\{ \begin{array}{l} \text{A: } h=0 \Rightarrow i \text{ & ii perfectly match} \\ \therefore \text{Cov}(X_t, X_t) = \text{Var}(X_t) = \sigma^2 (1 + \theta_1^2 + \dots + \theta_q^2) \\ \text{Func. of } h, \text{ not } t \\ \text{B: } h \leq q \Rightarrow \text{overlap} \\ \therefore \text{Cov}(X_t, X_{t+h}) = \theta_0 \sigma^2 + \theta_1 \theta_{1+h} \sigma^2 + \dots + \theta_{q-h} \theta_q \sigma^2 \\ \text{C: } h > q \Rightarrow 0, \text{ no overlap} \end{array} \right.$$

\textcircled{1}, \textcircled{2}, \textcircled{3} \Rightarrow \{X_t\} \text{ is stationary}

AVCF of $MA(q)$:

$$\gamma(h) = \sigma^2 \sum_{j=0}^{q-|h|} \theta_j \theta_{j+|h|}$$

ACF of $MA(q)$:

$$p(h) = \frac{\sum_{j=0}^{q-|h|} \theta_j \theta_{j+|h|}}{\sum_{j=0}^q \theta_j^2}$$

Note that for $h > 0$,
 $p(h) = 0$

↓
Signature of $MA(q)$ proc.

q -dependence: Process $\{X_t\}$ is q -dep. if X_t & X_s are dep. $\Leftrightarrow |t-s| \leq q$.
and indep. if $|t-s| > q$

↳ pretty strict, need joint prob distns.

γ -correlated: Process $\{x_t\}$ is γ -corr. if

$$\gamma(h) = \begin{cases} \text{non-zero} & h=\alpha \\ 0 & h>\alpha \end{cases}$$

Ex:// Prove / disprove that iid sq. of R.V. is α -dependent.

True:

Case A: $|t-s|=0 \Rightarrow x_t = x_s$. This is trivially dependent.

Case B: $|t-s|>0 \Rightarrow$ b/c iid, $x_t \perp x_s$.

Ex:// Prove / disprove that white noise is α -corr. & α -dep.

α -corr:

Suppose $Z_t \sim WN(0, \sigma^2)$. At lag 0, there is corr. ($\rho=0^2$) b/c it's same variable

At lag $h>0$, $\text{Corr}(Z_t, Z_{t+h}) = 0$ by defn. of WN

α -dep.:

Uncorrelation of $Z_t \not\Rightarrow Z_t \perp Z_{t+h}$

\Rightarrow
jointly normal

Ex:// Show all $MA(q)$ processes are γ -corr.

Suppose X_t is $MA(q)$ process:

$$\gamma(h) = \begin{cases} \sigma^2(1+\theta_1^2 + \dots + \theta_q^2), & h=0 \\ \sigma^2 \sum_{j=0}^{q-h} \theta_j \theta_{j+h}, & h \leq q \\ 0 & h > q \end{cases}$$

At $h=\alpha$, $\gamma(h) \neq 0$. At $h>q$, $\gamma(h)=0 \Rightarrow \gamma$ -corr.

Proposition: γ -correlated w/ mean 0 $\Leftrightarrow MA(q)$ process & stationary

AR(P) Process

$$AR(1): X_t = \Phi X_{t-1} + Z_t, Z_t \sim WN(0, \sigma^2).$$

$$\Rightarrow (1 - \Phi B) X_t = Z_t$$

• $|\Phi| < 1 \Rightarrow X_t$ is stationary

• ACF of AR processes have exp. decay $\rightarrow 0$

Another way of representing AR(1) is via MA(∞)

• Defn: X_t is MA(∞) if

$$\textcircled{1} \quad X_t := \sum_{j=0}^{\infty} \psi_j Z_{t-j}, \quad Z_t \sim WN(0, \sigma^2)$$

$$\textcircled{2} \quad \sum_{j=0}^{\infty} |\psi_j| < \infty$$

• Wold decomposition theorem: any stationary process = MA(∞) process + indep. det. process

$$\Rightarrow X_t = \sum_{j=0}^{\infty} \psi_j Z_{t-j} + \mu_t$$

• Conditions:

$$1. \psi_0 = 1, \sum_{j=0}^{\infty} |\psi_j| < \infty$$

2. $\{\psi_j\}$, $\{\mu_t\}$ are unique

3. $\{\mu_t, t \in \mathbb{Z}\}$ is a linear combo of past & det.

$$4. \text{Cor}(\mu_t, \mu_s) = 0 \quad \forall t, s$$

• AR(1) \rightarrow MA(∞)

$$X_t = \Phi X_{t-1} + Z_t$$

$$= \Phi (\Phi X_{t-2} + Z_{t-1}) + Z_t$$

$$= \Phi^2 X_{t-2} + \Phi Z_{t-1} + Z_t$$

= ;

$$= \sum_{j=0}^{\infty} \Phi^j Z_{t-j} \Rightarrow$$

We also know converges

if $|\Phi| < 1$

$$\therefore \sum_{j=0}^{\infty} |\Phi^j| = \frac{1}{1 - |\Phi|} < \infty$$

AR(P) process defn:

$$\Phi(B) := 1 - \Phi_1 B - \Phi_2 B^2 - \dots - \Phi_p B^p$$

$$\hookrightarrow \Phi(B) X_t = Z_t \quad (\text{AR}(P) \text{ process})$$

$$\Rightarrow X_t = \Phi_1 X_{t-1} + \Phi_2 X_{t-2} + \dots + \Phi_p X_{t-p} + Z_t$$

- Signature: exp. decay behavior on ACF

- By Wold decomposition, we know:

$$\{X_t\} = \sum_{j=0}^{\infty} \psi_j Z_{t-j} \Rightarrow \text{exists } \psi_j \text{ is some func. of } \Phi;$$

- ACVF:

$$\gamma(h) = \sigma^2 \sum_{j=0}^{\infty} \psi_j \psi_{j+h}$$

$$\rho(h) = \frac{\sum_{j=0}^{\infty} \psi_j \psi_{j+h}}{\sum_{j=0}^{\infty} \psi_j^2}$$

} No φ s.t. $\rho(h) = 0 \forall h > 2$
 } Not φ -corr.

Linearity

Defn: $\{X_t\}$ is linear \Leftrightarrow

$$\textcircled{1} \quad X_t = \sum_{j=-\infty}^{\infty} \psi_j Z_{t-j}, \quad Z_t \sim WN(0, \sigma^2)$$

$$\textcircled{2} \quad \sum_{j=-\infty}^{\infty} |\psi_j| < \infty$$

To show a process is linear:

① Show it can be rep. as $\sum_{j=-\infty}^{\infty} \psi_j Z_{t-j}$, define ψ_j

② Show that ψ_j converges

Ex:// Show that AR(1) process is linear

$$\textcircled{1} \quad \text{Under } |\Phi| < 1, \quad X_t = \sum_{j=0}^{\infty} \Phi^j Z_{t-j}$$

$$\Rightarrow X_t = \sum_{j=-\infty}^{\infty} \psi_j Z_{t-j}, \quad \psi_j = \begin{cases} \Phi^j & \text{if } j \geq 0 \\ 0 & \text{o.w.} \end{cases}$$

$$\textcircled{1} \quad \sum_{j=-\infty}^{\infty} |\psi_j| = \sum_{j=-\infty}^0 + \underbrace{\sum_{j=0}^{\infty} |\psi_j|}_{<\infty} \Rightarrow \text{converges.}$$

Causality: $X_t = f(Z_s)$, $s < t$ (doesn't need future info)

- If X_t is linear, X_t is causal $\Leftrightarrow \psi_j = 0 \forall j < 0$
- Ex:// AR(1) & MA(2) are causal processes

Linear predictors

Linear processes are easy to work w/ \rightarrow best predictors will usually be linear

Let's create a best linear predictor for a Gaussian process

Gaussian process: $\{X_t\}$ where all finite dimensional distn. are MVN

Note:

$$\begin{pmatrix} X_1 \\ X_2 \end{pmatrix} \sim N \left[\begin{pmatrix} \mu_1 \\ \mu_2 \end{pmatrix}, \begin{pmatrix} \sigma_1^2 & \sigma_{12} \\ \sigma_{21} & \sigma_2^2 \end{pmatrix} \right]$$

$$\Rightarrow X_1 | X_2 = x_2 \sim N \left(\mu_1 + \frac{\sigma_{12}}{\sigma_2^2} (x_2 - \mu_2), \sigma_1^2 + \frac{\sigma_{12}^2}{\sigma_2^2} \right)$$

$$= N \left(\mu_1 + \frac{\sigma_{12}}{\sigma_2^2} (x_2 - \mu_2), \sigma_1^2 (1 - p^2) \right)$$

\downarrow If X_1 & X_2 are com.
 $p^2 \uparrow \Rightarrow \text{Var}(X_1 | X_2 = x_2) \downarrow$
so pred. is easy!

Want to forecast X_{n+h} ($h > 0$) using a func. of X_n , $m(X_n)$. $\{X_n\}$ is stat.

$$\text{MSE} = E[(X_{n+h} - m(X_n))^2]$$

What is best $m(X_n)$ given X_n ?

We already know from Loss Function, that

$$m(X_n) = E(X_{n+h} | X_n) \Rightarrow \text{minimizes MSE}$$

Finding $E(X_{n+h} | X_n)$:

$$\begin{aligned} \begin{pmatrix} X_{n+h} \\ X_n \end{pmatrix} &\sim N \left(\begin{pmatrix} \mu \\ \mu \end{pmatrix}, \begin{pmatrix} \sigma^2 & \gamma(h) \\ \gamma(h) & \sigma^2 \end{pmatrix} \right) \\ b/c \text{ stat.} & \\ \Rightarrow X_{n+h} | X_n = x_n &\sim N \left(\mu + \frac{\gamma(h)}{\gamma(0)} (x_n - \mu), \gamma(0) (1 - \varphi(h)^2) \right) \\ &= N \left(\mu + \varphi(h) (x_n - \mu), \gamma(0) (1 - \varphi(h)^2) \right) \end{aligned}$$

From
c-side

This is best predictor of X_{n+h} . Note that this is linear wrt X_n .
b/c Gaussian.

W/b if process is not Gaussian?

↳ Best predictor of X_{n+h} may not be linear, but we don't care! Focus on best linear predictor b/c easy

↳ For any process $\{X_n\}$, we derive $m(X_n) = aX_n + b$ via

$$\min_{a, b} E \left[(X_{n+h} - aX_n - b)^2 \right]$$

$$\frac{\partial M}{\partial a} = -2E[X_n(X_{n+h} - aX_n - b)] = 0 \quad \text{--- (1)}$$

$$\frac{\partial M}{\partial b} = -2E[X_{n+h} - aX_n - b] = 0$$

$$\begin{aligned} E[X_{n+h}] - a E[X_n] - b &= 0 \\ \mu - a\mu - b &= 0 \end{aligned} \quad \begin{array}{l} E[X_{n+h}] = E[X_n] = \mu \text{ b/c} \\ \text{stationary} \end{array}$$

$$\hat{b} = \mu(1 - \hat{a})$$

$$(1): E[X_n X_{n+h}] - \hat{a} E[X_n^2] - \hat{b} E[X_n] = 0$$

$$E[X_n X_{n+h}] - \hat{a} E[X_n^2] - \mu^2(1 - \hat{a}^2) = 0$$

$$E[X_n X_{n+h}] - \mu^2 - \hat{a}(E[X_n^2] - \mu^2) = 0$$

$$\underbrace{\gamma(h)}_{\text{---}} \quad \underbrace{\sigma^2}_{\text{---}}$$

$$\gamma(h) - \hat{a}\sigma^2 = 0$$

$$\hat{a} = \frac{\gamma(h)}{\sigma^2} = \varphi(h)$$

$$\hat{m}(X_n) = \varphi(h)X_n + \mu(1 - \varphi(h)) \Rightarrow \text{Best linear predictor for any stationary process}$$

Reminder: Best linear predictor may not be good for non-linear process predictions.

Q: Forecast using all historical data \rightarrow best linear predictor?

A: Theorem, best linear predictor of X_{n+h} given X_1, \dots, X_n is

$$\hat{X}_{n+h} = P_n X_{n+h} = \mu + \sum_{j=1}^n \alpha_j (X_{n+1-j} - \mu)$$

predicting X_{n+h}
w/ n history

$$\mu = E[X_t], \quad \alpha_0 = \mu(1 - \sum_{i=1}^n \alpha_i), \quad \Gamma_n \alpha_n = \gamma_n(h) \Rightarrow \alpha_n = \Gamma_n^{-1} \gamma_n(h)$$

$$\Gamma_n = \begin{bmatrix} \gamma(0) & \gamma(1) & \dots & \gamma(n-1) \\ \gamma(1) & \ddots & & \gamma(n-2) \\ \vdots & & \ddots & \gamma(0) \\ \gamma(n-1) & \gamma(n-2) & \dots & \gamma(1) \end{bmatrix}, \quad \gamma_n(h) = \begin{bmatrix} \gamma(h) \\ \gamma(h+1) \\ \vdots \\ \gamma(n+h-1) \end{bmatrix}$$

Properties:

① Only need $\vec{X}, \mu, \gamma(h)$

② MSE = $E(\|X_{n+h} - P_n X_{n+h}\|^2) = \gamma(0) - \alpha_n' \gamma_n(h)$

③ Unbiased: $E(X_{n+h} - P_n X_{n+h}) = 0$

④ $E([X_{n+h} - P_n X_{n+h}] X_j) = 0, j = 1, 2, \dots$

↳ Proof: $\text{Cov}(E_m, X_j) = E[E_m X_j] - E[\overline{E_m}] E[X_j]$

$$= E[\overline{E_m} X_j]$$

Error should be uncorr. w/ past values, so $E[\overline{E_m} X_j] = 0$

⑤ $P_n [\alpha_1 X_{n+h_1} + \alpha_2 X_{n+h_2} + \beta] = \alpha_1 P_n X_{n+h_1} + \alpha_2 P_n X_{n+h_2} + \beta$

↳ Predictions are linearly decomposable

⑥ $P_{\text{pred}}(\sum_{i=1}^n \alpha_i X_i + \beta | X_1, \dots, X_n) = \sum \alpha_i X_i + \beta$

↳ Predictor on fitted values will be fitted values themselves

⑦ $P_n X_{n+h} = \mu \Leftrightarrow \gamma_n(h) = \text{Cov}(X_{n+h}, \vec{X}) = 0$

Ex:// Given X_1, \dots, X_n , denote $P_n X_{n+1}$ for an AR(1) process & calculate MSE

① AR(1) process representation:

$$X_t = \Phi X_{t-1} + Z_t, \quad |\Phi| < 1, \quad Z_t \sim WN(0, \sigma^2)$$

② Solve for a_n :

$$a_n = \Gamma_n^{-1} \gamma_n(h)$$

$\vdots \in \mathbb{R}$

③ One-step prediction ($h=1$)

$$\Gamma_n a_n = \gamma_n(h)$$

$$\begin{aligned} & \left[\begin{array}{cccc} \gamma(0) & \Phi \gamma(0) & \dots & \Phi^{n-1} \gamma(0) \\ & \ddots & \ddots & \ddots \\ & & \ddots & \gamma(0) \end{array} \right] \begin{bmatrix} a_1 \\ \vdots \\ a_n \end{bmatrix} = \begin{bmatrix} \Phi \gamma(0) \\ \Phi^2 \gamma(0) \\ \vdots \\ \Phi^n \gamma(0) \end{bmatrix} \\ & \stackrel{\div \gamma(0)}{\downarrow} \left[\begin{array}{cccc} 1 & \Phi & \dots & \Phi^{n-1} \\ & \vdots & & \vdots \\ & & \ddots & \end{array} \right] \begin{bmatrix} a_1 \\ \vdots \\ a_n \end{bmatrix} = \begin{bmatrix} \Phi \\ \vdots \\ \Phi^n \end{bmatrix} \\ & \therefore \begin{bmatrix} a_1 \\ \vdots \\ a_n \end{bmatrix} = \begin{bmatrix} \Phi \\ 0 \\ \vdots \\ 0 \end{bmatrix} \end{aligned}$$

④ Write in form:

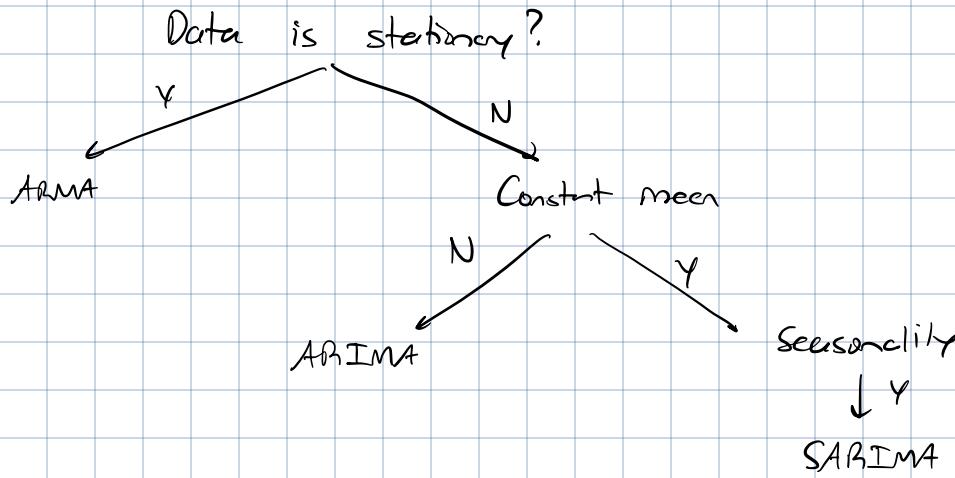
$$\begin{aligned} \hat{X}_{n+1} &= \mu + \sum_{i=1}^n a_i (X_{n+1-i} - \mu) \quad \mu = 0 \\ &= \sum_{i=1}^n a_i X_{n+1-i} \quad a_i = \begin{cases} \Phi, & i=1 \\ 0, & \text{o.w.} \end{cases} \\ &= \Phi X_n \end{aligned}$$

⑤ MSE:

$$\begin{aligned} \text{MSE} &= E[(\hat{X}_{n+1} - X_{n+1})^2] = E[(\Phi X_n - X_{n+1})^2] \quad \text{from AR(1) defn.} \\ &= E[\Phi X_n - \Phi X_n - Z_t^2] \\ &= E[Z_t^2] \quad E[X_t^2] = \text{Var}(X) + E[X]^2 \\ &= \sigma^2 \end{aligned}$$

BOX JENKINS MODELS

ARMA: autoregressive moving average process



ARMA (p, q)

Defn: $\{X_t\}$ is ARMA (p, q) if

$$1. \quad X_t - \underbrace{\mathbb{D}_1 X_{t-1} - \dots - \mathbb{D}_p X_{t-p}}_{\text{AR } (p) \text{- like}} = \underbrace{Z_t + \theta_1 Z_{t-1} + \dots + \theta_q Z_{t-q}}_{\text{MA } (q) \text{- like}}, \quad Z_t \sim WN$$

2. $\{X_t\}$ is stationary (constant mean is fine, we can adjust $\{X_t - \mu\}$)

$$3. \quad \begin{cases} \mathbb{D}(z) = 1 - \mathbb{D}_1 z - \dots - \mathbb{D}_p z^p \\ \theta(z) = 1 + \theta_1 z + \dots + \theta_q z^q \end{cases} \quad \left\{ \begin{array}{l} \text{no common factors.} \end{array} \right.$$

Alt. way of writing: $\mathbb{D}(B) X_t = \theta(B) Z_t$

$$\begin{cases} \mathbb{D}(B) = 1 - \mathbb{D}_1 B - \dots - \mathbb{D}_p B^p \\ \theta(B) = 1 + \theta_1 B + \dots + \theta_q B^q \end{cases}$$

Theorem: $\mathbb{D}(z)$ has no roots on unit circle (i.e. $\mathbb{D}(z) \neq 0 \quad \forall z \in \mathbb{C}: |z|=1$)
then unique stationary sdn. exists for X_t .

↳ Intuition: Solving for X_t :

$$X_t = \frac{1}{\mathbb{D}(B)} \theta(B) Z_t$$

From Laurent series, we know that if no roots are on unit circle, we can write:

$$\frac{1}{\mathbb{D}(B)} = \sum_{j=-\infty}^{\infty} x_j z^j \quad \text{for } 1 - \delta < |z| < 1 + \delta \text{ where } \delta \text{ is really small.}$$

Let $\frac{1}{\Phi(B)} = X(B)$

$$\begin{aligned} \therefore X_t &= X(B) \theta(B) Z_t \\ &= \sum_{j=-\infty}^{\infty} \psi_j B^j Z_t \\ &= \sum_{j=-\infty}^{\infty} \psi_j Z_{t-j} \quad \& \quad \sum_{j=-\infty}^{\infty} |\psi_j| < \infty \end{aligned}$$

We just showed that X_t can be written as a linear process if X_t is stationary!

Special cases for ARMA(p, q)

$$1) \text{ARMA}(0, q) \Rightarrow \Phi(B) = 1 \Rightarrow X_t = \theta(B) Z_t, \text{MA}(q)$$

$$2) \text{ARMA}(p, 0) \Rightarrow \theta(B) = 1 \Rightarrow \Phi(B) X_t = Z_t, \text{AR}(p)$$

$$3) \text{ARMA}(0, 0) \Rightarrow X_t = Z_t, \text{WN}$$

Causality & Invertibility

① Causal:

ARMA(p, q) causal \Rightarrow can be represented as

$$X_t = \sum_{j=0}^{\infty} \psi_j Z_{t-j} \quad \& \quad \sum_{j=0}^{\infty} |\psi_j| < \infty$$

more restrictive than stationary cond.

Equivalent to showing no roots of $\Phi(z)$ in/on unit circle

Basically asking if we can write X_t as an MA process

② Invertibility:

ARMA(p, q) invertible \Rightarrow can be represented as

$$Z_t = \sum_{j=0}^{\infty} \pi_j X_{t-j}, \quad \sum_{j=0}^{\infty} |\pi_j| < \infty$$

Equivalent to showing $\theta(z)$ has no roots in/on unit circle.

Basically asking if X_t can be written as AR process

$$\text{Ex:// } X_t - 0.5X_{t-1} = Z_t + 0.4Z_{t-1}, \quad Z_t \sim WN(0, \sigma^2).$$

a) Is X_t causal? If so, show causal soln.

① Check roots of $\Phi(z)$

$$\Phi(z) = 1 - 0.5z \rightarrow z = 2, \text{ so not on unit circle.}$$

This is causal

② Find ψ_j

Note that:

$$\frac{\Theta(z)}{\Phi(z)} = \psi(z) \Rightarrow \Theta(z) = \Phi(z)\psi(z)$$

Thus

$$1 + 0.4z = (1 - 0.5z)(\psi_0 + \psi_1 z + \psi_2 z^2 + \dots) \quad \forall z$$

Find pattern by equating coefficients of some power of z :

$$z^0: 1 = \psi_0$$

$$z^1: 0.4 = \psi_1 - 0.5\psi_0 \Rightarrow \psi_1 = 0.9$$

$$z^2: 0 = \psi_2 - 0.5\psi_1 \Rightarrow \psi_2 = 0.5 \cdot 0.9$$

$$z^3: 0 = \psi_3 - 0.5\psi_2 \Rightarrow \psi_3 = 0.5^2 \cdot 0.9$$

To prove general form of ψ_j , let's use induction

$$\text{Base case: } \psi_3 = 0.5^2 \cdot 0.9$$

$$\text{Assume: } \psi_k = 0.5^{k-1} \cdot 0.9$$

$$\text{Inductive step: prove } \psi_{k+1} = 0.5^k \cdot 0.9$$

$$\text{For } z^{k+1}: 0 = \psi_{k+1} - 0.5\psi_k$$

$$\begin{aligned} \psi_{k+1} &= 0.5\psi_k \\ &= 0.5^k \cdot 0.9 \quad \square \end{aligned}$$

Thus?

$$\psi_j = \begin{cases} 0.5^{j-1} \cdot 0.9, & j \geq 1 \\ 1, & j = 0 \end{cases}$$

③ Soln:

Causal soln:

$$X_t = \sum_{j=0}^{\infty} \psi_j Z_{t-j} = Z_t + 0.9 \sum_{j=1}^{\infty} 0.5^{j-1} Z_{t-j}$$

Also note that $\sum |\psi_j| < \infty$!

b) Is X_t invertible? If so show invertible solution

① Check if invertible:

$$\Theta(z) = 1 + 0.4z \Rightarrow z = -2.5 \text{ not on unit circle} \\ \text{so invertible}$$

② Find ψ_j

$$\Phi(z) = \Theta(z) \pi(z)$$

$$(-0.5z = (1 + 0.4z)(\pi_0 + \pi_1 z + \pi_2 z^2 + \dots))$$

$$1 = \pi_0$$

$$-0.5 = \pi_1 + 0.4\pi_0 \Rightarrow \pi_1 = -0.9$$

$$\cdot 0 = \pi_2 + 0.4\pi_1 \Rightarrow \pi_2 = 0.4(0.9)$$

⋮

$$\pi_j = \begin{cases} -0.9(-0.4)^{j-1}, & j \geq 1 \\ 1, & j=0 \end{cases}$$

③ Soln:

$$Z_t = \sum_{j=0}^{\infty} \pi_j X_{t-j} \Rightarrow Z_t = X_t - 0.9 \sum_{j=1}^{\infty} (-0.4)^{j-1} X_{t-j}$$

$$\text{Also, } \sum_{j=0}^{\infty} |\pi_j| < \infty$$

Ex:// $X_t = 0.7X_{t-1} - 0.1X_{t-2} + Z_t$ (AR(2) process).

c) Is this invertible?

Check $\Theta(z) = 1 \rightarrow$ no roots on unit circle

$$\Theta(z) \rightsquigarrow (-0.7z - 0.1z^2 = \pi_0 + \pi_1 z + \pi_2 z^2 + \dots)$$

$$\pi_0 = 1, \pi_1 = -0.7, \pi_2 = -0.1, \pi_3 = \pi_4 = \dots = 0$$

$$Z_t = \sum_{j=0}^{\infty} \pi_j X_{t-j}$$

b) Show MA(∞) esp.

To show X_t as MA, we need to get causal soln.

(Can show: $\psi_j = \begin{cases} 0.7\psi_{j-1} - 0.1\psi_{j-2}, & j \geq 2 \\ 0.7, & j=1 \\ 1, & j=0 \end{cases} \rightarrow \text{no closed form}$)

$$\Rightarrow X_t = \sum_{j=0}^{\infty} \psi_j Z_{t-j}$$

ACF of ARMA(p, q)

Consider a causal & stationary ARMA(p, q) process: $\Phi(B)X_t = \Theta(B)Z_t$

This will have an MA(∞) rep: $X_t = \sum_{j=0}^{\infty} \psi_j Z_{t-j}$

We have:

$$\gamma(h) = \sigma^2 \sum_{j=0}^{\infty} \psi_j \psi_{j+h}$$

Ex:// Derive ACVF & ACF of ARMA(1, 1).

$$X_t - \Phi X_{t-1} = Z_t + \theta Z_{t-1}$$

(1) Get MA(∞) if causal

$$X_t = \sum_{j=0}^{\infty} \psi_j Z_{t-j}, \quad \psi_j = \begin{cases} 1, & j=0 \\ \Phi^{j-1}(\theta + \Phi), & j \geq 1 \end{cases}$$

(2) $\gamma(0)$:

$$\begin{aligned} \gamma(0) &= \sigma^2 \sum \psi_j^2 \\ &= \sigma^2 \left(1 + \sum_{j=1}^{\infty} (\theta + \Phi)^2 (\Phi^{j-1})^2 \right) \\ &= \sigma^2 \left(1 + \frac{(\theta + \Phi)^2}{\Phi^2} \sum_{j=1}^{\infty} \Phi^{2j} \right) \\ &= \sigma^2 \left(1 + \frac{(\theta + \Phi)^2}{\Phi^2} \cdot \frac{\Phi^2}{1 - \Phi^2} \right) \\ &= \sigma^2 \left(1 + \frac{(\theta + \Phi)^2}{(1 - \Phi^2)} \right) \end{aligned}$$

(3) $\gamma(h)$:

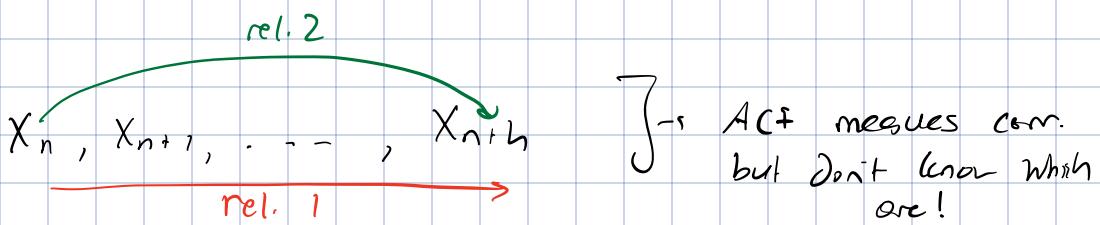
$$\begin{aligned} \gamma(h) &= \sigma^2 \left[\Phi^{h-1} (\theta + \Phi) + \sum_{j=1}^{\infty} \Phi^{j-1} (\theta + \Phi)^2 \Phi^{j+h-1} \right] \\ &= \Phi^{h-1} \cdot \sigma^2 \left[(\theta + \Phi) + (\theta + \Phi)^2 \cdot \frac{\Phi}{1 - \Phi^2} \right] \end{aligned}$$

PARTIAL AUTO CORRELATION FUNCTION

Idea: If X is related to Y but Z is a confounding variable, then we should regress X on Z & Y on Z . Correlation b/w residuals

is partial corr.

In forecasting:



- This is a problem w/ AR processes b/c AR is a func. of many X_i

PACF looks at direct corr. b/w X_n & X_{n+h} by conditioning/regressing on $X_{n+1} + \dots + X_{n+h-1}$.

Definition: PACF of process $\{X_t\}$ is defined as

$$\alpha(h) = \begin{cases} 1 & h=0 \\ \text{Cor}(X_{n+1}, X_n), |h|=1 \\ \rho_{X_{n+h}, X_n \cdot \{X_{n+1}, \dots, X_{n+h-1}\}}, |h| > 1 \end{cases}$$

*Cor b/w X_n & X_{n+h} condition
on $X_{n+1} \dots X_{n+h-1}$*

Note that we can write ρ_{X_{n+h}, X_n} in another way

$$\begin{aligned} \rho_{X_{n+h}, X_n \cdot \{X_{n+1}, \dots, X_{n+h-1}\}} &= \text{Cor}(X_n, X_{n+h} \mid X_{n+1}, \dots, X_{n+h-1}) \\ &= \text{Cor}(\textcircled{1}, \textcircled{2}) \\ &\quad \begin{array}{c} \text{---} \\ \boxed{\frac{X_n - \bar{X}_n}{\sqrt{\text{Var}(X_n)}} \quad \frac{X_{n+h} - \bar{X}_{n+h}}{\sqrt{\text{Var}(X_{n+h})}}} \end{array} \\ &\quad \text{Cor. of residuals!} \\ &\quad \begin{array}{c} \text{---} \\ \boxed{\frac{\text{Pred}(X_n \mid X_{n+1}, \dots, X_{n+h-1})}{\sqrt{\text{Var}(\text{Pred}(X_n \mid X_{n+1}, \dots, X_{n+h-1}))}}} \end{array} \end{aligned}$$

Ex:// PACF of AR(1) process

① AR(1) defn:

$$X_t = \varphi X_{t-1} + Z_t$$

② Find predictor

$$\text{For } AR(p) \Rightarrow \hat{X}_{n+1} = \sum_{j=1}^p \varphi_j X_{n+1-j}$$

$$\text{For } AR(1) \Rightarrow \hat{X}_{n+1} = \varphi X_n$$

③ Cases:

$$h=0 \Rightarrow 1$$

$$h=1 \Rightarrow \text{Cor}(X_{t+1}, X_t) = \mathbb{D}$$

$$\begin{aligned} h=2 &\Rightarrow \text{Cor}(X_t - \hat{X}_t, X_{t+2} - \hat{X}_{t+2}) \\ &= \text{Cor}(X_t - \hat{X}_t, X_{t+2} - \mathbb{D}X_{t+1}) \\ &= \text{Cor}(\underbrace{X_t - \hat{X}_t}_{\text{L.C. of } Z_i}, Z_{t+2}) \\ &= 0 \end{aligned}$$

$$\therefore \alpha(h) = \begin{cases} 1, & h=0 \\ \mathbb{D}, & h=1 \\ 0, & h \geq 2 \end{cases}$$

For ARMA processes:

$$\alpha(0) = 1$$

$$\alpha(h) = \mathbb{D}(hh) \rightarrow \text{lent comp. of } \Gamma_h^{-1} \gamma_h$$

Ex:// Calculate $\alpha(2)$ for MA(1)

① MA(1) process

$$X_t = Z_t + \theta Z_{t-1}$$

② $\gamma(h)$:

We know from before that:

$$\gamma(h) = \begin{cases} (1+\theta^2)\sigma^2, & h=0 \\ \theta\sigma^2, & h=1 \\ 0, & h \geq 1 \end{cases}$$

③ Calculate:

$$\alpha(0) = 1$$

$$\alpha(1) = \gamma(1) = \frac{\theta}{(1+\theta)^2}$$

$$\begin{aligned} \alpha(2) \Rightarrow \mathbb{D}_2 &= \Gamma_2^{-1} \gamma_2 = \begin{bmatrix} \gamma(0) & \gamma(1) \\ \gamma(1) & \gamma(0) \end{bmatrix}^{-1} \begin{bmatrix} \gamma(1) \\ \gamma(2) \end{bmatrix} \\ &= \frac{1}{(1+\theta^2)\sigma^4 - \theta^2\sigma^4} \begin{bmatrix} \theta((1+\theta^2)\sigma^4) \\ -\theta^2\sigma^4 \end{bmatrix} \end{aligned}$$

$$\alpha(2) = \mathbb{D}_{22} = \frac{-\theta^2}{(1+\theta^2) + \theta^4}$$

Recursive formula:

$$\alpha(h) = \hat{\mathbb{D}}_h = \frac{\gamma(h) - \sum_{j=1}^{h-1} \hat{\mathbb{D}}_{h-1,j} \cdot \rho(h-j)}{1 - \sum_{j=1}^{h-1} \hat{\mathbb{D}}_{h-1,j} \cdot \alpha(j)}$$

Ex: // Calculate $\alpha(2)$ fr MA(1) recursively

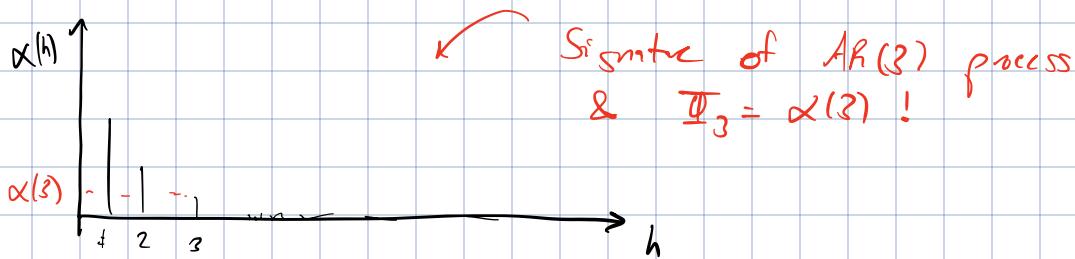
$$\begin{aligned}\alpha(2) &= \frac{\rho(2) - \Phi_{1,1} \cdot \rho(1)}{1 - \Phi_{1,1} \cdot \rho(1)} \\ &= \frac{\rho(2) - \rho(1)^2}{1 - \rho(1)^2} \\ &= \frac{-\theta^2}{1 + \theta^2 + \theta^4}\end{aligned}$$

} Much simpler! Use this

Why is PACF important? Theorem:

$\{X_t\}$ is a causal AR(p) process $\Leftrightarrow \alpha(p) \neq 0, \alpha(h) = 0 \quad \forall h > p$

$$\Leftrightarrow \alpha(p) = \Phi_p$$



General ACF & PACF shapes of ARMA processes:

	ACF	PACF
MA(q)	cuts off after lag q	exp. decay / damped sinusoid
AR(p)	exp. decay / damped sinusoid	cut off after lag p
ARMA	!!	exp. decay / damped sinusoid.

To propose models:

- ① Construct ACF & PACF plots
- ② Look @ cutoffs & decays. Mcp to table doe.

NOTE: PACF not useful if data is non-stationary

ARIMA & SARIMA MODELS

Up to this point:

Data \longrightarrow Remove trend & seasonality \longrightarrow Apply ARIMA on data

Can we do this all at once? Yup!

ARIMA: autoregressive integrated moving avg.

$\hookrightarrow \{X_t\}$ is ARIMA (p, d, q) if $Y_t := (1 - B)^d X_t$ is ARMA (p, q) process

↳

$$\underbrace{\Phi(B)}_{\Phi^*(B)} (1 - B)^d X_t = \Theta(B) Z_t$$

This type of model can only handle trend non-stationarity, not seasonality

Ex:// Process is $X_t = 0.8 X_{t-1} + 2t + Z_t$, $Z_t \stackrel{iid}{\sim} N(0, 35^2)$. Write process in ARIMA format.

① Rewrite:

$$(1 - 0.8B) X_t = Z_t + 2t \rightarrow \text{Like AR}(1)$$

② Diff. \rightarrow get rid of $2t$

$$\begin{aligned} \nabla X_t &= X_t - X_{t-1} \\ &= 0.8(X_{t-1} - X_{t-2}) + Z_t - Z_{t-1} + 2 \end{aligned}$$

Rewrite in terms of $Y_t := \nabla X_t$

$$Y_t = 0.8 Y_{t-1} + Z_t - Z_{t-1} + 2$$

Clearly, this process has a mean (+2). Finding mean.

$$(Y_t - \mu) - 0.8(Y_{t-1} - \mu) = Z_t - Z_{t-1}$$

$$\mathbb{E} \hookrightarrow -\mu + 0.8\mu = 2$$

$$\mu = 10$$

$$\therefore (Y_t - 10) - 0.8(Y_{t-1} - 10) = Z_t - Z_{t-1}$$

③ Model:

$$\text{ARIMA } (1, 1, 1)$$

Can we incorporate seasonality? Yes! \rightarrow SARIMA

\hookrightarrow Remember: ACF shows seasonality via periodicity! Corr. at season is NOT seasonality.

SARIMA defn: $\{X_t\} \in \text{SARIMA}(p, d, q) \times (P, D, Q)_s$ if.

$Y_t := \nabla^d \nabla_s^D X_t$
is causal ARMA process

\hookrightarrow How many times seasonal diff'g needs to happen.

$$\therefore \Phi(B) \Phi^*(B^s) Y_t = \Theta(B) \Theta^*(B^s) Z_t$$

\hookrightarrow Seasonal polynomials, only hor Z^a w/ a is multiple of s

Confusing! Just think of SARIMA as 2 ARMA model.

1. ARMA on regular lags (w/out seasonal lags)

2. ARMA on just seasonal lags

$\Phi(z) \wedge \Phi(z^s)$ cannot have roots on unit circle if causal.

If any params. missing $\rightarrow 0$

If $d = D = 0 \rightarrow X_t$ is stationary

How to find params of SARIMA model?

① Plot ACF to check for patterns

② Difference seasonal & regular until stationarity $\rightarrow d, D, s$

③ ACF & PACF on differenced data $\rightarrow P, Q, \theta, \phi$

④ Fit

NOTE: ACF of X_t^2 has trends $\rightarrow X_t$ has non-constant variance.

\hookrightarrow Why? $\text{Var}[X_t] = E[X_t^2] - E[X_t]^2$ $\text{O if trend removed}$

\hookrightarrow Trend if $E[X_t^2]$ has trend.

SARIMA output on R gives us 4 plots:

- ① Residual plot \Rightarrow should look like white noise
- ② Residual ACF \Rightarrow no corr. except first one
- ③ Normal Q-Q \Rightarrow normality for prediction intervals
- ④ p-val for Ljung-Box \Rightarrow $H_0: \rho(k) = 0 \forall k < h$
Ideally large p-values, no corr.

ESTIMATING ARMA(p, q) PARAMS

We have $p+q+1$ params to estimate:

- ① $\varnothing_1, \dots, \varnothing_p$
- ② $\theta_1, \dots, \theta_q$
- ③ σ

Two methods of estimating:

- ① Likelihood.

Assume $(x_1, \dots, x_n) \sim \text{MVN}$. Then:

$$\text{Lik}(\varnothing_1, \dots, \varnothing_p, \theta_1, \dots, \theta_q, \sigma^2) = \frac{1}{(2\pi)^n \prod_{i=1}^n \Gamma_i} \exp\left(-\frac{1}{2} \sum_{i=1}^n x_i^\top \Gamma_i^{-1} x_i\right)$$

Optimize Likelihood via params. This is comp. expensive!

- ② Yule-Walker Method (only $AB(p)$)

$$\begin{aligned}\hat{\varnothing} &= \hat{\Gamma}_p^{-1} \hat{\gamma}_p, \quad \hat{\sigma}^2 = \hat{\gamma}(0) - \hat{\varnothing}^\top \hat{\gamma}_p \\ &\doteq \hat{\gamma}(0) \left(\hat{\varnothing} = \left(\frac{\hat{\gamma}(k)}{\hat{\gamma}(0)} \right) \cdot \hat{\gamma}_p^\top, \quad \hat{\sigma}^2 = \hat{\gamma}(0) \left(1 - \hat{\varnothing}^\top \hat{\gamma}_p \right) \right)\end{aligned}$$

where $\hat{\gamma}_p = \begin{bmatrix} \hat{\gamma}(1) \\ \vdots \\ \hat{\gamma}(p) \end{bmatrix}, \quad \hat{\Gamma}_p^{-1} = \begin{bmatrix} \hat{\gamma}(0) & \hat{\gamma}(1) & \cdots & \hat{\gamma}(p-1) \\ \vdots & \ddots & \ddots & \vdots \\ \hat{\gamma}(1) & \hat{\gamma}(2) & \cdots & \hat{\gamma}(p) \end{bmatrix}^{-1}$

For large sample size:

$$\hat{\Phi} \sim MVN \left(\bar{\Phi}, \frac{\sigma^2}{n} \bar{P}_\Phi^{-1} \right)$$

FORECASTING IN ARMA MODELS

AR(p) forecasting

$$\hat{x}_{n+h} = \begin{cases} \sum_{j=1}^p \hat{\Phi}_j x_{n+h-j}, & h=1 \\ \underbrace{\sum_{j=1}^{h-1} \hat{\Phi}_j \hat{x}_{n+h-j}}_{\text{use pred. values}} + \underbrace{\sum_{j=h}^p \hat{\Phi}_j x_{n+h-j}}_{\text{history}}, & h=2 \dots p \\ \sum_{j=1}^h \hat{\Phi}_j \hat{x}_{n+h-j}, & h>p \end{cases}$$

All in history
history
All pred. values

h-step predictor is linear function of either historical data points or predicted values.

Ex:// Assume annual sales of store is modelled as AR(2) w/ $\hat{\Phi}_1 = 1$, $\hat{\Phi}_2 = -0.21$. Sales of 2021, 2020 & 2019 are 9, 11 & 10 million.

Forecast 2022:

$$\hat{x}_{2022} = \sum_{j=1}^p \hat{\Phi}_j \hat{x}_{2022-j} = \hat{\Phi}_1 \hat{x}_{2021} + \hat{\Phi}_2 \hat{x}_{2020}$$

$$\hat{x}_{2023} = \hat{\Phi}_1 \hat{x}_{2022} + \hat{\Phi}_2 \hat{x}_{2021}$$

$$\hat{x}_{2024} = \hat{\Phi}_1 \hat{x}_{2023} + \hat{\Phi}_2 \hat{x}_{2022}$$

MA(q) forecasting

$$\hat{x}_{n+h} = \begin{cases} \sum_{j=h}^q \theta_j z_{n+h-j}, & 1 \leq h \leq q \\ 0, & h > q \end{cases}$$

$$z_j = 0 \quad \forall j \leq 0, \quad z_n = x_n - \sum_{j=1}^q \theta_j z_{n-j}, \quad n \in 1, 2, \dots$$

Ex:// Let $X_t = Z_t + 0.5 Z_{t-1}$. $\alpha_1 = 0.3$, $\alpha_2 = -0.1$, $\alpha_3 = 0.1$. Forecast X_4 & X_5

① Classify X_t as MA

$X_t \sim MA(1)$ & invertible

② Let $Z_0 \dots Z_3$

$$Z_0 = 0$$

$$Z_1 = \alpha_1 - 0.5 Z_0 = 0.3$$

$$Z_2 = \alpha_2 - 0.5 Z_1 = -0.25$$

$$Z_3 = \alpha_3 - 0.5 Z_2 = 0.225$$

③ Forecast

$$\hat{X}_4 = \theta Z_3$$

$$= 0.5 \cdot 0.225$$

$$= 0.1125$$

$$\hat{X}_5 = \dots = 0$$

ARMA(p, q) forecast

We have $X_n = \Phi_1 X_{n-1} + \dots + \Phi_p X_{n-p} + Z_n + \Theta_1 Z_{n-1} + \dots + \Theta_q Z_{n-q}$

It is just a combination of AR(p) & MA(q) combined for forecast

$$\hat{X}_{n+h} = \begin{cases} \sum_{j=1}^p \Phi_j X_{n+h-j} + \sum_{i=1}^q \Theta_i Z_{n+h-i}, & h=1 \\ \sum_{j=1}^{h-1} \Phi_j \hat{X}_{n+h-j} + \sum_{j=h}^p \Phi_j X_{n+h-j} + \sum_{i=h}^q \Theta_i Z_{n+h-i}, & h>1 \end{cases}$$

No. 1:

$$\textcircled{1} \quad Z_n = \begin{cases} X_n - \sum_{j=1}^p \Phi_j X_{n-j} - \sum_{i=1}^q \Theta_i Z_{n-i}, & n>p \\ Z_n = 0, & n \leq p \end{cases}$$

$$\textcircled{2} \quad \Phi_j = 0, \quad \forall j > p$$

$$\textcircled{3} \quad \sum_{i=h}^q \Theta_i Z_{n+h-i} = 0, \quad h>q \rightarrow \text{obs.}$$

$$\textcircled{24} \quad \sum_{j=h}^p \Phi_j X_{n+h-j} = 0, \quad h > p \quad \rightarrow \text{obv.}$$