

1. (15 points) Short questions: Provide justification for all answers to receive full credit.

(a) Under what conditions does weak-sense stationarity imply strong-sense stationarity?

When the timeseries is a linear combination of gaussian distributions

(b) Is the auto-covariance at lag zero, $C_Y(k, 0)$, necessarily equal to the variance of the process Y if the process is not weak-sense stationary (WSS)?

No, if the process is not WSS, then $C_Y(k, 0)$ may also depend on k rather than solely $h \Rightarrow C_Y(k, 0) \neq \text{Var}(Y)$

(c) Describe the autocorrelation function (ACF) of a white noise process.

A white noise process is WSS, its ACF is independent with time k

(d) Consider the following process:

$$Y_k = Y_{k-1} + \epsilon_k, \quad \text{where } \epsilon_k \sim \mathcal{N}(0, \sigma_\epsilon^2) \text{ and } Y_0 = 0.$$

is not a WSS for all k

Is this process weak-sense stationary (WSS) for all k ? Justify your answer.

Mean: $E[Y_k] = E[Y_{k-1}] + E[\epsilon_k] \Rightarrow \text{Var: } \text{Var}[Y_k] = \text{Var}[Y_{k-1}] + \sigma_\epsilon^2$
 $= E[Y_{k-1}] + 0$
 $\Rightarrow E[Y_k] = 0 \text{ for } k \in \mathbb{Z} \checkmark$
 $\Rightarrow \text{Var}[Y_k] = k\sigma_\epsilon^2 \text{ depends on } k,$

(e) True or False: If a stochastic process X Granger-causes Y , then there is a true causation from X to Y .

True

(f) What issue arises when calculating the empirical autocorrelation function (ACF) of an AR(1) process with a linear trend added? The empirical ACF uses a single realization of the path integral to estimate the theoretical autocorrelation.

The empirical ACF may be drifted to infinity as k goes to infinity

(g) Consider a one-step-ahead forecast density that is symmetric around its mean μ , i.e., $f_{Y_{k+1}|\mathcal{F}_k}(\mu + y | \mathcal{F}_k) = f_{Y_{k+1}|\mathcal{F}_k}(\mu - y | \mathcal{F}_k)$ for all $y \in \mathbb{R}$. What is the optimal one-step-ahead predictor \hat{y}_{k+1}^* under the absolute error (AE) loss function? Justify your answer.

$\hat{y}_{k+1}^* \equiv f_{Y_{k+1}|\mathcal{F}_k}(\mu | \mathcal{F}_k)$, because after symmetry + AE means the peak is at $y = \mu$

(h) What is the Gibbs sampler used for in Bayesian estimation?

It's used to sample random process in gibbs method

- (i) What is the most appropriate loss function when the one-step-ahead forecast distribution is uniform between -1 and 1?

Squared Error Loss

- (j) What is one of the key advantages of the deep learning approach compared to traditional machine learning methods?

① More powerful, it can learn representation in data automatically thus having better accuracy in inference.

- (k) What is one significant disadvantage of the traditional machine learning approach compared to the model-based approach?

It may overfit on the training set data,
and ML ~~costs~~ ^{lost} more resource / data

- (l) What are the key differences between the model-based approach and the Bayesian model-based approach?

Bayesian model-based one involves random variable to handle estimation

- (m) What are the advantages of Transformers compared to vanilla RNNs? Vanilla RNNs do not include LSTMs or GRUs.

① Self-attention mechanism could obtain global info of data

② Multitask makes GPU-based multitasked processing easily

- (n) In the context of Maximum Likelihood Estimation (MLE), why do we maximize the likelihood (i.e., the probability density of the observed data given the parameters) rather than maximizing the probability of observing the data directly?

Likelihood is the summary of all probabilities given partial condition.

It could better indicate the global optimum, which may be local

- (o) Provide the definition of the conditional risk, $\mathcal{R}(\hat{y}_{k+1}; \mathcal{F}_k) = \mathbb{E}[\ell(Y_{k+1}, \hat{y}_{k+1}) | \mathcal{F}_k]$ in the form of an integral. in prob. ...

$$\begin{aligned} \mathcal{R}(\hat{y}_{k+1}; \mathcal{F}_k) &= \mathbb{E}[\ell(Y_{k+1}, \hat{y}_{k+1}) | \mathcal{F}_k] \\ &= \mathbb{E}\left[\prod_{k=1}^L P(y_k | \mathcal{F}_k)\right] \\ &= \int_{\mathcal{F}_k} P(y_k; \mathcal{F}_k) \end{aligned}$$

2. (15 points) Programming questions:

(a) What type of stochastic process does the following Python function generate? What do the parameters param1 and param2 represent in the process?

```
def generate_process(param1, param2, n):
    process = np.zeros(n)
    meas = np.random.normal(0, param2, n)
    for t in range(1, n):
        process[t] = param1 * process[t-1] + meas[t]
    return process
```

$$Y_k = \phi Y_{k-1} + \epsilon_k$$

Answer:

AR(1) :

param1 is ϕ for the AR(1)

param2 is σ^2 of white noise
variance

(b) Describe the type of stochastic process generated by the following Python function, and explain the role of the parameters in the process.

```
def generate_process(L=1000, C=2, A=None, AC=None):
    if A is None:
        A = np.array([[0.3, 0.5], [0.2, 0.1]])
    if AC is None:
        AC = np.eye(C)
    meas = np.random.multivariate_normal(np.zeros(C), AC, L)
    data = np.zeros((L, C))
    for t in range(1, L):
        data[t] = A @ data[t-1] + meas[t]
    df = pd.DataFrame(data, columns=[f'Y{i+1}' for i in range(C)])
    return df
```

L: length of ~~process~~ data

C: matrix dim

A: ~~State transition~~ A

State transition matrix

AC: ~~matrix~~ (C x C) gaussian matrix

Multivariate Random Process

Answer:

Vector Autoregression process with lag (1)
(2x2 matrix)

(c) The following Python code generates a time series based on a certain stochastic process. Identify the type of process generated by this function and explain the meaning of each parameter in the context of this process.

```
def generate_time_series(param1, param2, param3, n):
    process = np.zeros(n)
    meas_values = np.random.normal(0, param3, n)
    for t in range(1, n):
        process[t] = param1 * process[t-1] + param2 * meas[t-1] + meas[t]
    return process
```

ϵ_k is normal distributed

$\epsilon_k \sim N(0, \text{param}^3)$

Answer:

MAAR(1)

$$Y_k = \phi Y_{k-1} + \theta \epsilon_{k-1} + \epsilon_k$$

param1 is ϕ for random process

param2 is θ for ϵ_{k-1} term

param3 is σ^2 variance of white noise

3. (10 points) Below is a matrix representing the autocovariance matrix of a stationary multivariate random process with two components: Z^1 and Z^2 . Information about the autocorrelation functions and cross-correlation functions between the processes are shown for lag values less than 3 in the following table:

Autocorrelation: $R_{Z^1}(h)$		Cross-correlation: $R_{Z^1 Z^2}(h)$	
Lag	Value	Lag	Value
0	1.0	0	0.7
1	0.6	1	0.35
2	0.3	2	0.15
3	0.15	3	0.05
Cross-correlation: $R_{Z^2 Z^1}(h)$		Autocorrelation: $R_{Z^2}(h)$	
Lag	Value	Lag	Value
0	0.7	0	1.0
1	0.4	1	0.5
2	0.05	2	0.25
3	0.01	3	0.1

Assuming that the variance of both Z^1 and Z^2 is 2, answer the following questions:

- (a) Compute $\text{Cov}(Z_k^2, Z_{k+3}^2)$.

$$\begin{aligned}\text{Cov}(Z_k^2, Z_{k+3}^2) &= \sigma_{Z^2}^2 R_{Z^2}(h), \quad h=3 \\ &= 2^2 \times 0.1 = 0.4\end{aligned}$$

- (b) Compute $\text{Cov}(Z_k^1, Z_{k+3}^2)$.

$$\begin{aligned}\text{Cov}(Z_k^1, Z_{k+3}^2) &= \sigma_{Z^1} \sigma_{Z^2} R_{Z^1 Z^2}(h=3) \\ &= 2 \times 2 \times 0.05 \\ &= 0.2\end{aligned}$$

- (c) Compute $\text{Cov}(Z_k^2, Z_{k+3}^1)$.

$$\begin{aligned}\text{Cov}(Z_k^2, Z_{k+3}^1) &= \sigma_{Z^2} \sigma_{Z^1} R_{Z^2 Z^1}(h=3) \\ &= 2 \times 2 \times 0.01 \\ &= 0.04\end{aligned}$$

4. (15 points) Consider a fair coin. At each time step k , you flip the coin. If it lands heads, you assign the value $N_k = 1$; if it lands tails, you assign the value $N_k = -1$. Based on this random process \mathcal{N} , define the following stochastic process:

$$X_k = \alpha X_{k-1} + N_k, \quad \text{with } X_0 = 1 \text{ and } |\alpha| < 1,$$

where $\alpha \in \mathbb{R}$ is a constant.

- (a) Compute the theoretical mean $E[X_k]$ of the process for all $k \in \mathbb{N}$.

$$\begin{aligned} E[X_k] &= E[\alpha X_{k-1} + N_k] \\ &= E[\alpha X_{k-1}] + E[N_k] \\ &= \alpha E[X_{k-1}] + 0 \end{aligned}$$

$$\text{So } E[X_k] = 0 \text{ for } k \in \mathbb{N}$$

- (b) Compute the theoretical variance $\text{Var}(X_k)$ of the process for all $k \in \mathbb{N}$. *Hint: Feel free to leave your answer as a summation. Hint: The variance of the coin-flip process is 1.*

$$\begin{aligned} \text{Var}(X_k) &= \text{Var}(\alpha X_{k-1} + N_k) \\ &= \text{Var}(\alpha X_{k-1}) + \text{Var}(N_k) \\ &= \alpha^2 \text{Var}(X_{k-1}) + 1 \end{aligned}$$

as hint. and proved in HW2

let $V_k = \text{Var}(X_k)$, then $V_k = \alpha^2 V_{k-1} + 1$

$$\begin{aligned} V_{k-1} &= \alpha^2 V_{k-2} + 1 \Rightarrow V_k = \alpha^2 \left(\frac{1 - \alpha^{k-1}}{1 - \alpha} \right) + 1 \\ &\vdots \\ V_2 &= \alpha^2 V_1 + 1 \\ V_1 &= \alpha^2 V_0 + 1 = 1 \end{aligned}$$

- (c) For any time step k , construct a 100% confidence interval for X_k that is as narrow as possible. For example, the interval $(-\infty, \infty)$ is not a valid answer. *Hint: Feel free to define the interval using summations.*

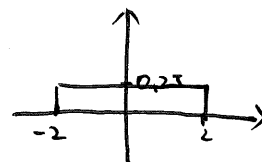
Within k step, we know

X_k can reach: highest as $1+k$, with prob = 0.5^k
lowest as $1-k$, with prob = 0.5^k

\Rightarrow ~~$X_k \in$~~ 100% conf interval is $[-k, k]$

5. (15 points) Consider the following model with uniform noise:

$$Y_k = \beta Y_{k-1} + \epsilon_k, \quad \text{with } \epsilon_k \stackrel{\text{iid}}{\sim} \text{Unif}(-2, 2) \text{ and } Y_0 = 2,$$



where β is the unknown parameter we want to estimate, and $\text{Unif}(a, b)$ is the uniform distribution for $a < b$, i.e., $X \sim \text{Unif}(a, b) \Rightarrow f_X(x) = \frac{1}{b-a}$ for $x \in [a, b]$; and 0 otherwise. In the following questions, you are required to perform one step of Bayesian parameter estimation assuming the following prior uniform distribution for β , i.e., $\beta \sim \text{Unif}(0, 2)$.

- (a) Compute the likelihood function
- $\mathcal{L}(\beta; y_0 = 2, y_1 = 0)$
- .

$\mathcal{L}(\beta; y_0 = 2, y_1 = 0)$ is where we only observe $Y_0 = 2$

$$= \cancel{P(Y_1 | Y_0 = 2)} =$$

$$\beta \sim \text{Unif}(0, 2)$$

$$\Rightarrow \beta \mid Y_1 \sim \text{Unif}(-2 + \beta, 2 + \beta)$$

$$\mathcal{L}(\beta; y_0 = 2, y_1 = 0) = \text{Unif}(-2, 4)$$

- (b) Compute the posterior distribution
- $f_{\beta|F_1}$
- given
- $F_1 = \{Y_0 = 2, Y_1 = 0\}$
- .

Based on (a), we now observe $Y_1 = 0$



$$f_{\beta|F_1} = \begin{cases} \frac{1}{2} & , 0 < \beta < 2 \\ 0 & , \text{otherwise} \end{cases}$$

- (c) Interpret the implication of
- $f_{\beta|F_1}$
- on the possible values of
- β
- .

$f_{\beta|F_1}$ means β is uniformly sample on $[0, 2]$,
with Y_0 & Y_1 fixed for observing

6. (15 points) Consider the following VAR(2) process:

$$y_k^1 = \frac{1}{2} y_{k-1}^1 + \frac{1}{3} y_{k-2}^1 + \eta_k^1$$

$$y_k^2 = \frac{1}{2} y_{k-1}^2 + \frac{1}{3} y_{k-1}^1 + \eta_k^2$$

$$\begin{bmatrix} y_k^1 \\ y_k^2 \end{bmatrix} = \begin{bmatrix} \frac{1}{2} & 0 \\ 0 & \frac{1}{2} \end{bmatrix} \begin{bmatrix} y_{k-1}^1 \\ y_{k-1}^2 \end{bmatrix} + \begin{bmatrix} 0 & \frac{1}{3} \\ \frac{1}{3} & 0 \end{bmatrix} \begin{bmatrix} y_{k-2}^1 \\ y_{k-2}^2 \end{bmatrix} + \begin{bmatrix} \eta_k^1 \\ \eta_k^2 \end{bmatrix},$$

where $n_k^1 \sim \mathcal{N}(0, 1)$, $n_k^2 \sim \mathcal{N}(1, 2)$, and initial conditions are given by $y_0^1 = 1$, $y_0^2 = 2$, $y_{-1}^1 = 2$, $y_{-1}^2 = 1$.

We would like to build an equivalent VAR(1) process of the form:

$$\mathbf{X}_k = A\mathbf{X}_{k-1} + \varepsilon_k.$$

- (a) What should be the dimension of the vector \mathbf{X}_k in the equivalent VAR(1) process?

~~2x2~~ 4×1

- (b) Explain how you would assign values to the entries of the vector \mathbf{X}_k .

$$\mathbf{X}_k = \begin{bmatrix} y_k^1 \\ y_{k-1}^1 \\ y_k^2 \\ y_{k-1}^2 \end{bmatrix} \text{ as } \mathbf{X}_{k-1} = \begin{bmatrix} y_{k-1}^1 \\ y_{k-2}^1 \\ y_{k-1}^2 \\ y_{k-2}^2 \end{bmatrix} \quad \varepsilon_k = \begin{bmatrix} \eta_k^1 \\ 0 \\ \eta_k^2 \\ 0 \end{bmatrix}$$

- (c) Provide an expression for the autoregressive matrix A in the VAR(1) process.

~~$$A = \begin{bmatrix} \frac{1}{2} & 0 & 0 & \frac{1}{3} \\ 0 & 0 & 0 & 0 \\ 0 & \frac{1}{3} & 0 & \frac{1}{2} \\ 0 & 0 & 0 & 0 \end{bmatrix}$$~~

$$A = \begin{bmatrix} \frac{1}{2} & 0 & 0 & \frac{1}{3} \\ 0 & 0 & 0 & 0 \\ 0 & \frac{1}{3} & 0 & \frac{1}{2} \\ 0 & 0 & 0 & 0 \end{bmatrix}$$

- (d) Provide the mean μ and covariance matrix Σ of the noise vector ε_k in the VAR(1) process.

$$\vec{\mu} = \begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \end{bmatrix} \quad \Sigma = \begin{bmatrix} 1 & 2 \\ 2 & 4 \end{bmatrix}$$

- (e) Provide the initial condition \mathbf{X}_0 for the VAR(1) process.

$$\mathbf{X}_0 = \begin{bmatrix} 1 \\ 2 \\ 2 \\ 1 \end{bmatrix}$$

7. (15 points) Consider a second-order Markov chain modeling the daily weather with two possible states, $\mathcal{S} = \{s_1 = \text{Hot}, s_2 = \text{Cold}\}$. To transform this process into a first-order Markov chain, we introduce a set of *compound states*, as defined as follows:

$$\mathcal{S}^2 = \{(\text{Hot}, \text{Hot}), (\text{Hot}, \text{Cold}), (\text{Cold}, \text{Hot}), (\text{Cold}, \text{Cold})\},$$

$$(H, H) \rightarrow (*, H)$$

↑
new

where the pairs of states are temporally ordered, i.e., the first entry in the compound state represents the weather the day before the second entry. In what follows, we analyze the resulting first-order Markov chain with these four compound states.

- (a) The transition matrix of the lifted first-order Markov chain is given by:

$$P = \begin{bmatrix} * & * & 0.5 & * \\ 0.5 & * & * & * \\ * & 0.5 & * & * \\ * & * & * & 0.5 \end{bmatrix} \quad \begin{bmatrix} 0.5 & 0 & 0.5 & 0 \\ 0.5 & 0 & 0 & 0.5 \\ 0 & 0.5 & 0.5 & 0 \\ 0 & 0 & 0 & 0.5 \end{bmatrix}$$

Complete the entries marked with *, ensuring the matrix remains stochastic. *Hints:* Recall that a stochastic matrix must have non-negative entries and each row must sum to 1. Additionally, certain entries are constrained to be zero due to the structure of the augmented Markov chain.

$$P = \begin{bmatrix} 0 & 0.5 & 0.5 & 0 \\ 0.5 & 0 & 0 & 0.5 \\ 0 & 0.5 & 0 & 0.5 \\ 0.5 & 0 & 0 & 0.5 \end{bmatrix}$$

in order to ensure

① no overlap (same) row

② each row sum to 1

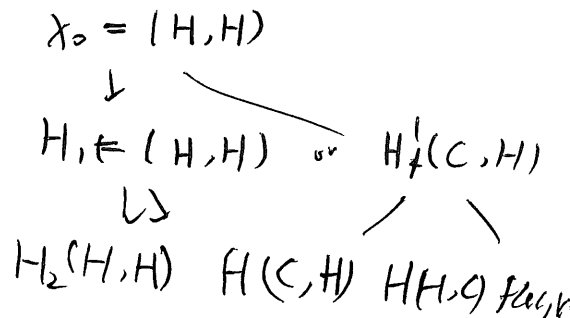
- (b) Given that the augmented Markov chain starts in state $X_0 = (\text{Hot}, \text{Hot})$, compute the probability that $X_2 = (\text{Cold}, \text{Hot})$ after two time steps. *Hint:* $\pi_k = (P^T)^k \pi_0$.

$$X_0 = (H, H) \rightarrow X_1 = (H, H) \rightarrow X_2 = (C, H) \quad \text{only 1 chance}$$

$$P = 0.5 \times 0.5 \\ = 0.25$$

- (c) Starting from $X_0 = (\text{Hot}, \text{Hot})$, compute the distribution vector π_2 , which represents the state probabilities at time step 2. *Hint: $\pi_k = (P^T)^k \pi_0$.*

$$\pi_2 =$$



- (d) In the original second-order Markov chain, what is the probability that a day is Hot when the previous two days are Hot?

$$(H, H, H) \rightarrow (H, H)$$

$$P_{\text{ans.}}(x_2 = H | x_1 = H, x_0 = H)$$

$$= P(x_2 = H) = 0.5$$

because markov chain depends on and only
on current state