

Hints and Solutions Sketches:

Q1.

i) Use formula for maximum likelihood estimation (in this case, just the relative frequency of heads). For Bayesian estimate: Beta prior (a,b) + Bernoulli likelihood (m heads, n tails) = Beta posterior (m+a, n+b). Repeat this step by step for every toss. Posterior after toss "i" becomes the prior for toss "i+1".

ii) Estimate the "lambda" parameter of Poisson distribution using maximum likelihood formula. This is the rate of bus arrivals per hour. So rate for 2 hours = 2*lambda. Calculate probability of 6 buses using 2*lambda as parameter.

iii) Use formula for mean (2D vector) and covariance matrix (2x2) for multivariate Gaussian.

Q2.

Prob. (Fair dice used twice | sum is 8) = (Prob.(Fair dice used twice)* Prob.(sum is 8|fair dice used twice))/Prob.(sum is 8)

Prob. (Fair dice used twice) = Prob.(coin gave two heads) = $0.7 \cdot 0.7 = 0.49$.

Prob. (loaded dice used twice) = Prob.(coin gave two tails) = $0.3 \cdot 0.3 = 0.09$.

Prob. (fair dice used once) = Prob.(coin gave 1 heads) = $2 \cdot 0.7 \cdot 0.3 = 0.42$.

Prob.(sum is 8) = Prob.(sum is 8, Fair dice used twice) + Prob.(sum is 8, loaded dice used twice) + Prob.(sum is 8, Fair dice used once)

Prob.(sum is 8) calculation: break up the cases, as 8 can be 2+6, 3+5, 4+4, 5+3, 6+2.

If fair dice used twice, all of these are possible with probability $(1/6) \cdot (1/6) = 1/36$ in each case. Total probability = $5/36$.

But if loaded dice used twice, then only 3+5 and 5+3 are possible. Probabilities: $1/3 \cdot 1/3 = 1/9$ in each case. Total probability = $2/9$.

Once fair dice used + once loaded dice used (prob = $2 \cdot 0.7 \cdot 0.3 = 0.42$): then also only 3+5 and 5+3 are possible. Probabilities: $1/3 \cdot 1/6 = 1/18$ in each case.

.....

Q3.

i) Expected value of a sample from each Gaussian distribution = mean vector of the Gaussian.

So expected value = $p \cdot M1 + (1-p) \cdot M2$

ii) For any observation X, its probability under the model = $p \cdot N(X, M1, C1) + (1-p) \cdot N(X, M2, C2)$ where $N(x, a, b)$ denotes pdf at x for a Gaussian distribution with mean a and covariance matrix b.

Choose p to maximize the above likelihood function.

Q4.

Mean stationarity: calculate the mean observation at each hour. Are they roughly equal?

Covariance stationarity: calculate the covariance between observations at each pair of time-points with equal difference, eg. 4:00-8:00, 8:00-12noon (4hour difference), 4:00-12noon, 8:00am-4:00pm (8 hour difference) etc. are they roughly equal?

Q5.

1st order autoregressive process using the same parameters: $X_4 = a \cdot X_0 + b$, $X_8 = a \cdot X_4 + b$, $X_{12} = a \cdot X_8 + b$, etc

1st order autoregressive process using different parameters: $X_4 = a_0 \cdot X_0 + b_4$, $X_8 = a_4 \cdot X_4 + b_8$, $X_{12} = a_8 \cdot X_8 + b_{12}$ etc

2nd order autoregressive process: $X_8 = c_0 \cdot X_0 + c_4 \cdot X_4 + d_8$ etc

Estimate the parameters by least-square regression. In which case is the error lowest?

Q6.

Consider a simple regression for each location separately, with other locations as predictors and their coefficients as defined. Calculate w and the bias (local component) by least square minimization.

Q7.

Let's say the observations at the two locations are $[x_1, x_2]$. As it is a Gaussian Process, this 2D vector follows a Gaussian distribution, with mean as $[0 \ 0]$, and covariance matrix calculated according to the covariance function. Now say $x_2=5$. Write the pdf of (x_1, x_2) , and choose x_1 so as to maximize this pdf.

Q8.

Directly apply Kriging formula as in the slides. Solve simultaneous equations.

Q9.

For all locations except s' (21,60), you know $m(s)$. So calculate $Y(s) = X(s) - m(s)$ for all the other locations. Treat the one unknown value as "y". Calculate the Gaussian distribution pdf, as they follow a Gaussian Process. Estimate "y" as to maximize the PDF just like Q7. Also estimate $Y(s')$ at the location s' for all 5 days in same way, i.e. maximizing Gaussian PDF. Estimate $m(s')$ as the mean difference between $X(s')$ (observed) and $Y(s')$ (estimated).

Q10.

At location s_1 , $a \cdot m(s_1)$ can be calculated. So $(1-a) \cdot n(t) + e(s_1, t)$ can also be calculated for each year. If we take the mean of these, we will get $(1-a) \cdot n(t)$, as mean of $e(s_1, t) = 0$. Thus, we get $n(t)$. Now consider every other location, use a and $n(t)$ to calculate $m(s)$.

Q11.

Climatology: at every location, calculate the mean over the 20 years. Anomaly: Observation – mean.

90th quantile: sort the 20 values at each location in ascending order, mean of the 18th and 19th values.

Return period: calculate the quantile of the mean value of each location, i.e. what fraction of observations at that location are more than the mean. Return period is its inverse.

Conditional probability calculation for temporal coherence: identify all positive anomalies and the gaps (no. of years) between successive positive anomalies. What fraction of these gaps is 1 year?

Conditional probability calculation for spatial coherence: at each location, identify all positive anomalies. In each case, see how many neighboring locations are also having positive anomaly on same year.

Q 12.

Suppose the matrix is approximately rank-1, so there is a 4-dim vector u and a 6-dim vector v such that $u \cdot v' = X$. For every known entry in X , write down the equation, eg. $u_1 \cdot v_1 = 20$, $u_1 \cdot v_2 = 16$ etc. Using these, try to estimate the values of X , i.e. $u_2 \cdot v_2 = ?$ $u_1 \cdot v_5 = ?$ etc

Try for rank-2 also, i.e. consider U as 4x2 matrix and V as 6x2 matrix. Again you will get many simultaneous equations.

Q13.

Suppose $X(s, t)$ is missing, denote its value by x . $p(x) = N(x, a, b)$, where a, b are Gaussian parameters

Also, suppose $X(s_1, t) = y_1$, $X(s_2, t) = y_2$ etc. So their differences with x are $(y_1 - x)$, $(y_2 - x)$ etc.

$p(y_1 - x) = N(y_1 - x, 0, k|s - s_1|)$, $p(y_2 - x) = N(y_2 - x, 0, k|s - s_2|)$ etc

The PDF of x then is given by $p(x) \cdot p(y_1 - x) \cdot p(y_2 - x) \cdot \dots$

Choose x to maximize this PDF.

Q14.

Return period of 10 years \Rightarrow probability of exceedance in any year is 0.1, i.e. 90th quantile.

You know the Gamma distribution parameters. Draw about 1000 samples from them, and find the 90th quantile. This is the threshold for anomaly.

Repeat the process for return period of 5 years \Rightarrow probability of exceedance in any year is 0.2, i.e. 80th quantile.

Q15.

At any location say s_2 , consider $\mu(s_2) = x$. Due to Gaussian Process assumption, $[x \ \mu(s_1)]$ follows multivariate Gaussian distribution with mean vector $[5 \ 5]$ and 2x2 covariance matrix C defined according to covariance function. Also, the GEV distribution PDF of the observations (X_1, X_2, \dots) at s_2 can be calculated in terms of x .

So PDF of $X = N([x \ \mu(s_1)], [5 \ 5], C) \cdot \text{GEV}(X_1, \mu(s_1), \sigma, \xi) \cdot \text{GEV}(X_2, \mu(s_1), \sigma, \xi) \cdot \dots$

Now choose x to maximize this PDF.