TSSL Lab 2 - Structural model, Kalman filtering and EM We will continue to work with the Global Mean Sea Level (GMSL) data that we got acquainted with in lab 1. The data is taken from https://climate.nasa.gov/vital-signs/sea-level/ and is available on LISAM in the file sealevel.csv. In this lab we will analyse this data using a structural time series model. We will first set up a model and implement a Kalman filter to infer the latet states of the model, as well doing long-term prediction. We will then implement a disturbance smoother and an expectation maximization algorithm to tune the parameters of the model. We load a few packages that are useful for solving this lab assignment. import pandas # Loading data / handling data frames import numpy as np import math import matplotlib.pyplot as plt plt.rcParams["figure.figsize"] = (12,8) # Increase default size of plots 2.1 Setting up a structural state space model We start by loading and plotting data to reming ourselves what it looks like. data=pandas.read_csv('sealevel.csv', header=0) y = data['GMSL'].values u = data['Year'].values ndata = len(y)plt.plot(u,y) plt.xlabel('Year') plt.ylabel('GMSL deviation (mm)') plt.show() 1995 60 40 GMSL deviation (mm) 1995 2000 2005 2010 2015 2020 Year len(y) Out[4]: 1047 In this lab we will use a structural time series model to analys this data set. Specifically, we assume that the data $\{y_t\}_{t\geq 1}$ is generated by (1) $y_t = \mu_t + \gamma_t + \varepsilon_t$ where μ_t is a trend component, γ_t is a seasonal component, and ε_t is an observation noise. The model is expressed using a state space representation, $lpha_{t+1} = Tlpha_t + R\eta_t, \qquad \eta_t \sim N(0,Q),$ (2) $y_t = Zlpha_t + arepsilon_t, \qquad arepsilon_t \sim N(0,\sigma_arepsilon^2).$ (3)**Q0:** Let $d = \dim(\alpha_t)$ denote the state dimension and $d_\eta = \dim(\eta_t)$ denote the dimension of the state noise. Then, what are the dimensions of the matrices T, R, and Z of the state space model? A: dim(T) = d * d $\dim(R) = d * d_n$ dim(Z)=1*d**Q1:** Create the state space matrices $T_{[\mu]}$, $R_{[\mu]}$, and $Z_{[\mu]}$ corresponding to the trend component μ_t . We assume a local linear trend (that is, of order k=2). Hint: Use **2-dimensional** numpy . ndarray s of the correct sizes to represent all the matrices. $T_{mu} = np.array([[2,-1],[1,0]])$ $R_{mu} = np.array([[1],[0]])$ $Z_{mu} = np.array([[1,0]])$ # T_mu.shape # R_mu.shape # Z_mu.shape **Q2:** There is a yearly seasonal pattern present in the data. What should we set the periodicity s of the seasonal component to, to capture this pattern? Hint: Count the average number of observations per (whole) year and round to the closest integer. In [6]: u1=np.around(u) np.unique(u1) (unique, counts) = np.unique(u1, return_counts=True) f = np.asarray((unique, counts)).T print(f) [[1993. 18.] Ī1994. 37.] [1995. 37.] [1996. 37.] [1997. 37.] [1998. 37.] [1999. 36.] [2000. 37.] [2001. 37.] [2002. 37.] [2003. 37.] [2004. 37.] [2005. 37.] [2006. 36.] [2007. [2008. 37.] [2009. 37.] Γ2010. 37.] [2011. 36.] [2012. 37.] [2013. 37.] [2014. 37.] [2015. 37.] [2016. 37.] [2017. 37.] [2018. 36.] [2019. 37.] [2020. 37.] [2021. **Q3:** What is the state dimension of a seasonal component with periodicity s? That is, how many states are needed in the corresponding state space representation? A: The s is determined by periodicity, where the seasonal pattern is constant over time, the seasonal values for year 1 to s can be modelled by the constants. So here the value s=37 is constantly repeating over the year. s=37 d=int(s)+1**Q4:** Create the state space matrices $T_{[\gamma]}$, $R_{[\gamma]}$, and $Z_{[\gamma]}$ corresponding to the seasonal component γ_t . Hint: Use **2-dimensional** numpy . ndarray s of the correct sizes to represent all the matrices. first_row = np.repeat(-1,36,axis=0) remain_row = np.hstack((np.identity(35),np.zeros(35).reshape(-1,1))) remain_row Out[8]: array([[1., 0., 0., ..., 0., 0., 0.], $[0., 1., 0., \ldots, 0., 0., 0.]$ $[0., 0., 1., \ldots, 0., 0., 0.]$ $[0., 0., 0., \dots, 0., 0., 0.]$ $[0., 0., 0., \ldots, 1., 0., 0.],$ $[0., 0., 0., \ldots, 0., 1., 0.]])$ In [9]: T_gamma = np.vstack((first_row,remain_row)) $R_{gamma} = np.append(1, np.zeros(35)).reshape(-1,1)$ $Z_{gamma} = np.append(1, np.zeros(35)).reshape(1, -1)$ Q5: Using the matrices that you have constructed above, create the state space matrices for the complete structural time series model. Print out the shapes of the resulting system matrices and check that they correspond to what you expect (cf Q0). Hint: Use scipy.linalg.block_diag and numpy.concatenate. In [10]: from scipy.linalg import block_diag T = block_diag(T_mu, T_gamma) R = block_diag(R_mu, R_gamma) Z = np.concatenate((Z_mu, Z_gamma), axis=1) print("Shape of T is {} ".format(T.shape)) In [11]: print("Shape of R is {}".format(R.shape)) print("shape of Z is {}".format(Z.shape)) Shape of T is (38, 38) Shape of R is (38, 2) shape of Z is (1, 38) We also need to specify the variances of the process noise η_t and measurement noise ε_t . Below, we will estimate (two of) these variances from data, but for now we set them arbitrarily to get an initial model to work with. # Some arbitrary noise values for now $sigma_trend = 0.01$ sigma_seas = 1 $sigma_eps = 1$ Q = np.array([[sigma_trend**2, 0.], [0., sigma_seas**2]]) # Process noise covariance matrix In [13]: Q.shape Out[13]: (2, 2) Finally, to complete the model we need to specify the distribution of the initial state. This encodes our a priori belief about the actual values of the trend and seasonality, i.e., before observing any data. **Q6:** Set up the mean vector of the initial state $a_1 = \mathbb{E}[\alpha_1]$ such that: • The trend component starts at the first observation, $\mathbb{E}[\mu_1] = y_1$, • The slope of the trend is a *priori* zero in expectation, $\mathbb{E}[\mu_1 - \mu_0] = 0$, • The initial mean of all states related to the seasonal component are zero. Also, create an initial state covariance matrix $P_1 = \text{Cov}(\alpha_1)$ as an identity matrix of the correct dimension, multiplied with a large value (say, 100) to represent our uncertainty about the initial state. a1 = np.zeros(d)P1 = np.identity(d)*100In [15]: a1[:2]=y[0] a1=a1.reshape(-1,1)We have now defined all the matrices etc. that make up the structural state space model. For convenience, we can create an object of the class LGSS available in the module tssltools_lab2 as a container for these quantities. In [16]: from tssltools_lab2 import LGSS $model = LGSS(T, R, Q, Z, sigma_eps**2, a1, P1)$ help(model.get_params) Help on method get_params in module tssltools_lab2: get_params() method of tssltools_lab2.LGSS instance Return all model parameters. T, R, Q, Z, H, a1, P1 = $model.get_params()$ 2.2 Kalman filtering for the structural model Now we have the data and a model available. Next, we will turn our attention to the inference problem, which is a central task when analysing time series data using the state space framework. State inference is the problem of estimating the unknown (latent) state variables given the data. For the time being we assume that the *model parameters* are completely specified, according to above, and only consider how to estimate the states using the Kalman filter. In the questions below we will treat the first n=800 time steps as training data and the remaining m=197 observations as validation data. In [17]: n = 800m = ndata-n # ndata- number of observation in y is 1047 # From the give data the value of m is 247 but we use m = 197 #np.isnan(y) In [18]: In [19]: $\#alpha_pred = np.zeros((d, 1, n))$ #np.zeros((4,2,5))In [20]: #np.matmul(np.matmul(R,Q),R.T) Q7: Complete the Kalman filter implementation below. The function should be able to handle missing observations, which are encoded as "not a number", i.e. y[t] = np.nan for certain time steps t. Hint: The Kalman filter involves a lot of matrix-matrix and matrix-vector multiplications. It turns out to be convient to store sequences of vectors (such as the predicted and filtered state estimates) as (d,1,n) arrays, instead of (d, n) or (n, d) arrays. In this way the matrix multiplications will result in 2d-arrays of the correct shapes without having to use a lot of explicit reshape. However, clearly, this is just a matter of coding style preferences! from tssltools_lab2 import kfs_res def kalman_filter(y, model: LGSS): """Kalman filter for LGSS model with one-dimensional observation. :param y: (n,) array of observations. May contain nan, which encodes missing observations. :param model: LGSS object with the model specification. :return kfs_res: Container class with member variables, alpha $_pred: (d,1,n)$ array of predicted state means. P_{pred} : (d,d,n) array of predicted state covariances. alpha_filt: (d,1,n) array of filtered state means. P_filt: (d,d,n) array of filtered state covariances. y_pred: (n,) array of means of $p(y_t | y_{1:t-1})$ F_pred: (n,) array of variances of $p(y_t | y_{1:t-1})$ n = len(y)d = model.d # State dimension $alpha_pred = np.zeros((d, 1, n))$ $P_{pred} = np.zeros((d, d, n))$ $alpha_filt = np.zeros((d, 1, n))$ $P_{filt} = np.zeros((d, d, n))$ y_pred = np.zeros(n) $F_pred = np.zeros(n)$ T, R, Q, Z, H, a1, P1 = model.get_params() # Get all model parameters (for brevity) for t in range(n): # Time update (predict) **if** t**==**0: $alpha_pred[:,:,t] = a1$ $P_pred[:,:,t] = P1$ else: alpha_pred[:,:,t] = T@alpha_filt[:,:,t-1] $P_pred[:,:,t] = T@P_filt[:,:,t-1]@T.T + R@Q@R.T # @ numpy matrix multiplication$ # Compute prediction of current output y_pred[t] = Z@alpha_pred[:,:,t] $F_pred[t] = Z@P_pred[:,:,t]@Z.T+H$ # Measurement update (correct) if np.isnan(y[t]): alpha_filt[:,:,t] = alpha_pred[:,:,t] P_filt[:,:,t] = P_pred[:,:,t] else: # Kalman Gain $k_t = (P_pred[:,:,t]@Z.T@(np.array([F_pred[t]]))**-1).reshape(-1,1)$ # update filter $alpha_filt[:,:,t] = alpha_pred[:,:,t] + np.array([k_t*(y[t]-y_pred[t])]).reshape(-1,1)$ $P_{filt}[:,:,t] = (np.identity(d)-k_t@Z)@P_pred[:,:,t]$ kf = kfs_res(alpha_pred, P_pred, alpha_filt, P_filt, y_pred, F_pred) return kf kalman_filter(y,model) In [22]: Out[22]: <tssltools_lab2.kfs_res at 0x2749187dc40> **Q8:** Use the Kalman filter to infer the states of the structural time series applied to the sealevel data. Run the filter on the training data (i.e., first n=700 time steps), followed by a long-range prediction of y_t for the remaining time points. Generate a plot which shows: 1. The data $y_{1:n+m}$, 2. The one-step predictions $\hat{y}_{t|t-1} \pm 1$ standard deviation for the training data, i.e., $t=1,\ldots,n$, 3. The long-range predictions $\hat{y}_{t|n}\pm 1$ standard deviation for the validation data, i.e., $t=n+1,\ldots,n+m$, 4. A vertical line indicating the switch between training and validation data, using plt.axvline(x=u[n]). Hint: It is enough to call the kalman_filter function once. Make use of the missing data functionality! $train_data = y[:n]$ In [23]: validation_data = np.repeat(np.nan,m) train = np.concatenate((train_data, validation_data)) K = kalman_filter(train, model) $y_hat = K.y_pred$ In [24]: $sd1 = y_hat+(K.F_pred)**(1/2)$ $sd2 = y_hat-(K.F_pred)**(1/2)$ plt.plot(u, y, label="Orginal values", color="black") In [25]: plt.plot(u, y_hat, label="one step Predicted values",color="red") plt.plot(u, sd1, linestyle=':',color="blue",label="Upper CI") plt.plot(u, sd2, linestyle=':',color="blue",label="Lower CI") plt.plot(u[801:997], y_hat[801:997], label="long range prediction", color="green") plt.axvline(x=u[n]) plt.xlabel('Year') plt.ylabel('GMSL deviation (mm)') plt.legend() plt.show() Orginal values 100 one step Predicted values ····· Upper CI Lower Cl 75 50 ation (mm) 25 GMSL devi -50 -75 -1001995 2000 2010 2015 2020 2005 **Q9:** Based on the output of the Kalman filter, compute the training data log-likelihood $\log p(y_{1:n})$. $F_t = K.F_pred[:n]$ $y_v = (y[:n]-y_hat[:n])$ $ml = y_v.T*F_t**-1*y_v$ $logllhood = (-1/2)*sum(np.log(F_t)+ml)$ logllhood -2107.31179640246 Out[26]: print("Log-likelihood for the training data is {}".format(logllhood)) In [27]: Log-likelihood for the training data is -2107.31179640246 2.3 Identifying the noise variances using the EM algorithm So far we have used fixed model parameters when running the filter. In this section we will see how the model parameters can be learnt from data using the EM algorithm. Specifically, we will try to learn the variance of the state noise affecting the seasonal component as well as the variance of the observation noise, $\theta = (\sigma_{\gamma}^2, \sigma_{\varepsilon}^2).$ (4)For brevity, the variance of the trend component σ_{μ}^2 is fixed to the value $\sigma_{\mu}^2=0.01^2$ as above. (See Appendix A below for an explanation.) Recall that we consider $y_{1:n}$ as the training data, i.e., we will estimate θ using only the first n=800 observations. Q10: Which optimization problem is it that the EM algorithm is designed to solve? Complete the line below! A: $\hat{ heta} = rg \max_{ heta} E[\log P_{ heta}(lpha_{1:n}, y_{1:n}) | y_{1:n}, \tilde{ heta})]$ **Q11:** Write down the updating equations on closed form for the M-step in the EM algorithm. Hint: Look at Exercise Session 2 $\sigma_arepsilon^2 = rac{1}{n} \sum_{t=1}^n [(\hat{arepsilon}_{t|n}^2 + Var(arepsilon_t \mid y_{1:n})]$ $Q_{t+1} = rac{1}{n} \sum_{t=1}^n [tr(\hat{\eta}_{t\mid n}\hat{\eta}_{t\mid n}^T + Var(\eta \mid y_{1:n})]$ To implement the EM algorithm we need to solve a smoothing problem. The Kalman filter that we implemented above is based only on a forward propagation of information. The smoother complements the forward filter with a backward pass to compute refined state estimates. Specifically, the smoothed state estimates comprise the mean and covariances of $p(lpha_t \mid y_{1:n}), \qquad t = 1, \dots, n$ (5)Furthermore, when implementing the EM algorithm it is convenient to work with the (closely related) smoothed estimates of the disturbances, i.e., the state and measurement noise, $egin{aligned} p(\eta_t \mid y_{1:n}), & & t = 1, \ldots, n-1 \ p(arepsilon_t \mid y_{1:n}), & & t = 1, \ldots, n \end{aligned}$ (6)(7)An implementation of a state and disturbance smoother is available in the tssltools_lab2 module. You may use this when implementing the EM algorithm below. from tssltools_lab2 import kalman_smoother help(kalman_smoother) Help on function kalman_smoother in module tssltools_lab2: kalman_smoother(y, model: tssltools_lab2.LGSS, kf: tssltools_lab2.kfs_res) Kalman (state and disturbance) smoother for LGSS model with one-dimensional observation. :param y: (n,) array of observations. May contain nan, which encodes missing observations. :param model: LGSS object with the model specification. :parma kf: kfs_res object with result from a Kalman filter foward pass. return kfs_res: Container class. The original Kalman filter result is augmented with the following member variables, alpha_sm: (d,1,n) array of smoothed state means. V: (d,d,n) array of smoothed state covariances. eps_hat: (n,) array of smoothed means of observation disturbances. eps_var: (n,) array of smoothed variances of observation disturbances. eta_hat: (deta,1,n) array of smoothed means of state disturbances. eta_cov: (deta,deta,n) array of smoothed covariances of state disturbances. # An implementation of disturbance smoother Ks=kalman_smoother(y=train, model=model, kf=K) #Ks.eta_cov[:,:,1] In [31]: #uon=Ks.eta_hat.T #Ks.eta_cov.shape **Q12:** Implement an EM algorithm by completing the code below. Run the algorithm for 100 iterations and plot the traces of the parameter estimates, i.e., the values θ_r , for $r=0,\ldots,100$. Note: When running the Kalman filter as part of the EM loop you should only filter the training data (i.e. excluding the prediction for validation data). def maxim_fn(s): In [32]: $sigma_eps = (1/n)*sum(s.eps_hat**2 + s.eps_var)$ $q = sum([(s.eta_hat[:,:,t]@s.eta_hat.T[t,:,:] + s.eta_cov[:,:,t])$ for t in range(n)])/nreturn[np.sqrt(sigma_eps), np.sqrt(q[1,1])] maxim_fn(Ks) Out[32]: [1.7931350962845962, 1.1969272540918143] $num_iter = 100$ In [33]: $sigma_t = 0.01$ theta_hat = {} theta_hat $[0] = [1, sigma_seas]$ In [34]: for r in range(num_iter): # E-step $Q_n = np.array([[sigma_t**2, 0.], [0., theta_hat[r][1]**2]])$ $m = LGSS(T, R, Q_n, Z, theta_hat[r][0]**2, a1, P1)$ k = kalman_filter(train_data,m) ks= kalman_smoother(train_data, m, k) theta_hat[r+1] = $maxim_fn(ks)$ In [35]: # converged values of sigma_epsilon and variance of the seasonal trend list(theta_hat.items())[95:] Out[35]: [(95, [2.7383364389614258, 0.19631707507634844]), (96, [2.738399546499734, 0.1949934245331376]), (97, [2.7384610369701123, 0.19369286043512043]), (98, [2.7385209665568935, 0.19241476382189718]), (99, [2.7385793890588106, 0.19115853778671746]), (100, [2.7386363560064653, 0.18992360650126286])] # Converting the dictionary to dataframe theta_h=pandas.DataFrame.from_dict(theta_hat) theta_h=theta_h.T plt.plot(theta_h[0],label="Estimates of sigma_epsilon") plt.plot(theta_h[1],label="Estimates of sigma_seasonal") plt.legend() plt.show() 2.5 2.0 1.5 1.0 0.5 Estimates of sigma_epsilon Estimates of sigma_seasonal 100 theta_h.iloc[100][0] Out[38]: 2.7386363560064653 2.4 Further analysing the data We will now fix the model according to the final output from the EM algorithm and further analyse the data using this model. **Q13:** Rerun the Kalman filter to compute a *long range prediction for the validation data points*, analogously to **Q8** (you can copy-paste code from that question). That is, generate a plot which shows: 1. The data $y_{1:n+m}$, 2. The one-step predictions $\hat{y}_{t|t-1}\pm 1$ standard deviation for the training data, i.e., $t=1,\dots,n$, 3. The long-range predictions $\hat{y}_{t|n}\pm 1$ standard deviation for the validation data, i.e., $t=n+1,\ldots,n+m$, 4. A vertical line indicating the switch between training and validation data, using plt.axvline(x=u[n]). Furthermore, compute the training data log-likelihood $\log p(y_{1:n})$ using the estimated model (cf. **Q9**). # optimal Q is the array of new sigma_trend and sigma_seasonal &&&& new In [39]: $Optimal_Q = np.array([[sigma_t**2, 0.], [0., theta_hat[99][1]**2]])$ $o_m = LGSS(T, R, Optimal_Q, Z, theta_hat[99][0]**2, a1, P1)$ ok = kalman_filter(train,o_m) In [40]: y_h=ok.y_pred $sd11 = y_h+(ok.F_pred)**(1/2)$ $sd22 = y_h-(ok.F_pred)**(1/2)$ plt.plot(u,y,label="Orginal values",color="black") In [41]: plt.plot(u,y_h,label="Predicted values",color="red") plt.plot(u,sd11,linestyle=':',color="blue",label="Upper CI") plt.plot(u, sd22, linestyle=':', color="blue", label="Lower CI") plt.axvline(u[n]) plt.xlabel('Year') plt.ylabel('GMSL deviation (mm)') plt.legend() plt.show() Orginal values Predicted values 75 ····· Upper CI Lower CI 50 25 GMSL deviation (mm) 0 -50-75-1001995 2000 2020 2005 2010 2015 $F_t = ok.F_pred[:n]$ $y_v = (y[:n]-y_h[:n])$ $ml = y_v.T*F_t**-1*y_v$ $logllhood = (-1/2)*sum(np.log(F_t)+ml)$ print("Log-likelihood for the training data is {}".format(logllhood)) Log-likelihood for the training data is -1370.4945075287667 Note that we can view the model for the data y_t as being comprised of an underlying "signal", $s_t = \mu_t + \gamma_t$ plus observation noise ε_t (8) $y_t = s_t + \varepsilon_t$ We can obtain refined, *smoothed*, estimates of this signal by conditioning on all the training data $y_{1:n}$. **Q14:** Run a Kalman smoother to compute smoothed estimates of the signal, $\mathbb{E}[s_t|y_{1:n}]$, conditionally on all the *training data*. Then, similarly to above, plot the following: 1. The data $y_{1:n+m}$, 2. The smoothed estimates $\mathbb{E}[s_t|y_{1:n}]\pm 1$ standard deviation for the training data, i.e., $t=1,\ldots,n$, 3. The predictions $\mathbb{E}[s_t|y_{1:n}]\pm 1$ standard deviation for the validation data, i.e., $t=n+1,\ldots,n+m$, 4. A vertical line indicating the switch between training and validation data, using plt.axvline(x=u[n]). *Hint:* Express s_t in terms of α_t . Based on this expression, compute the smoothed mean and variance of s_t based on the smoothed mean and covariance of α_t . ks2 = kalman_smoother(train,o_m,ok) In [43]: In [44]: y_t=np.zeros(ndata) for i in range(ndata): $y_t[i]=o_m.Z@ks2.alpha_sm[:,:,i] + ks2.eps_hat[i]$ # alpha_sm: (d,1,n) array of smoothed state means. # eps_hat: (n,) array of smoothed means of observation disturbances plt.plot(u, y, color="black", label="data") In [45]: plt.plot(u,y_t,color="blue",label="fitted values of smoothened observatiton") plt.plot(u, y_t + (ks2.eps_var)**(1/2), linestyle=":", color= "red", label="CI to smooth estimates") plt.plot(u, y_t - (ks2.eps_var)**(1/2), linestyle=":", color= "red") plt.axvline(x=u[n]) plt.xlabel('Year') plt.ylabel('GMSL deviation (mm)') plt.legend() plt.show() 795 40 GMSL deviation (mm) 1995 2000 2015 2020 2005 2010 Q15: Explain, using a few sentences, the qualitative differences (or similarities) between the Kalman filter predictions plotted in Q13 and the smoothed signal estimates plotted in Q14 for, 1. Training data points, $t \leq n$ 2. Validation data points, t > nA: However, both the kalman filter and kalman smoother predicted the data exactly. There is small difference in log-likelihood values of smoother and filter. Likelihood of smoother values is higher than the filter on the training $data(y_{1:n})$. The smoother add the backward pass on the top of filtering. Where the confidence interval are also found to be closer to the smoothened estimates. We can shed additional light on the properties of the process under study by further decomposing the signal into its trend and seasonal components. **Q16:** Using the results of the state smoother, compute and plot the *smoothed* estimates of the two signal components, i.e.: 1. Trend: $\hat{\mu}_{t|n} = \mathbb{E}[\mu_t|y_{1:n}]$ for $t=1,\ldots,n$ 2. Seasonal: $\hat{\gamma}_{t|n} = \mathbb{E}[\gamma_t|y_{1:n}]$ for $t=1,\ldots,n$ (You don't have to include confidence intervals here if don't want to, for brevity.) $E_mu = np.zeros(ndata)$ In [46]: $E_{gamma} = np.zeros(ndata)$ for i in range(ndata): $E_mu[i] = Z_mu@ks2.alpha_sm[:2,:,i]$ E_gamma[i] = Z_gamma@ks2.alpha_sm[2:,:,i] plt.plot(u,y,label="data") plt.plot(u, E_mu, label="Trend") plt.plot(u,E_gamma,label="Seasonal",color="red") plt.xlabel('Year') plt.ylabel('GMSL deviation (mm)') plt.legend() plt.show() Trend Seasonal 40 GMSL deviation (mm) 1995 2015 2020 2005 2010 2.5 Missing data We conclude this section by illustrating one of the key merits of the state space approach to time series analysis, namely the simplicity of handling missing data. To this end we will assume that a chunk of observations in the middle of the training data is missing. **Q17:** Let the values y_t for $300 < t \le 400$ be missing. Modify the data and rerun the Kalman filter and smoother. Plot, 1. The Kalman filter predictions, analogously to **Q8** 2. The Kalman smoother predictions, analogously to Q13 Comment on the qualitative differences between the filter and smoother estimates and explain what you see (in a couple of sentences). $data_p = np.concatenate((y[:300], np.repeat(np.nan, 100), y[400:n]))$ In [49]: kh = kalman_filter(data_p, model) ks=kalman_smoother(data_p, model, kh) ks_y_pred = np.zeros(len(data_p)) for i in range(len(data_p)): ks_y_pred[i] = model.Z@ks.alpha_sm[:,:,i] plt.plot(u[:len(data_p)], data_p, label="Orginal values", color="black") plt.plot(u[:len(data_p)],kh.y_pred,label="Kalman Filter prediction",color="red") plt.plot(u[:len(data_p)], ks_y_pred, label="Kalman smoother prediction", color="green") plt.axvline(u[300]) plt.axvline(u[400]) plt.legend() plt.show() Orginal values Kalman Filter prediction 30 Kalman smoother prediction 20 10 -102015 I notice that the filter and smoother prediction are exactly same on the validation data points. From the likelihood calculation, I observe there is a slight increase in smoother. where the current time point predicted with the help of previous time points with upto current value, so the Kalman filter is satisfying the otimal condition of the observing the null values. In addition, smoothing is used to achieve the exact result of filter by iterating the backpass through the data, but smoothing does not give more information. Appendix A. Why didn't we learn the trend noise variance as well? In the assignment above we have fixed σ_{μ} to a small value. Conceptually it would have been straightforward to learn also this parameter with the EM algorithm. However, unfortunately, the maximum likelihood estimate of σ_{μ} often ends up being too large to result in accurate long term predictions. The reason for this issue is that the structural model (9) $y_t = \mu_t + \gamma_t + \varepsilon_t$ is not a perfect description of reality. As a consequence, when learning the parameters the mismatch between the model and the data is compensated for by increasing the noise variances. This results in a trend component which does not only capture the long term trends of the data, but also seemingly random variations due to a model misspecification, possibly resulting in poor long range predictions. Kitagawa (Introduction to Time Series Modeling, CRC Press, 2010, Section 12.3) discusses this issue and proposes two solutions. The first is a simple and pragmatic one: simply fix σ_u^2 to a value smaller than the maximum likelihood estimate. This is the approach we have taken in this assignment. The issue is of course that in practice it is hard to know what value to pick, which boild down to manual trial and error (or, if you are lucky, the designer of the lab assignment will tell you which value to use!). The second, more principled, solution proposed by Kitagawa is to augment the model with a stationary AR component as well. That is, we model $y_t = \mu_t + \gamma_t + \nu_t + \varepsilon_t$ (10)where $\nu_t \sim \mathsf{AR}(p)$. By doing so, the stationary AR component can compensate for the discrepancies between the original structural model and the "true data generating process". It is straightforward to include this new component in the state space representation (how?) and to run the Kalman filter and smoother on the resulting model. Indeed, this is one of the beauties with working with the state space representation of time series data! However, the M-step of the EM algorithm becomes a bit more involved if we want to use the method to estimate also the AR coefficients of the ν -component, which is beyond the scope of this lab assignment.