
Modeling trend in temperature volatility using generalized LASSO

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 words, words,

2 1 Introduction

3 Nonparametric variance estimation for spatio-temporal data.

4 1.1 Motivating applications

5 There is a considerable interest in determining if there is an increasing trend in the climate variability
6 [8, 10]. An increase in the temperature variability will increase the probability of extreme hot outliers.
7 It might be harder for the society to adapt to these extremes than to the gradual increase in the mean
8 temperature [10].

9 In this project, we consider the problem of detecting the trend in the temperature volatility. All
10 analyses are performed on a sub-set of the European Centre for Medium-Range Weather Forecasts
11 (ECMWF) ERA-40 dataset [26]. This dataset include the temperature measurements over a grid over
12 the earth from 1957 to 2002. [6, 18, 19, 25, 27]

13 Research on analyzing the trend in the volatility of spatio-temporal data is scarce. [8] studied the
14 change in the standard deviation (SD) of the surface temperature in the NASA Goddard Institute
15 for Space Studies gridded temperature data set. In their analysis, for each geographical position,
16 the mean of the temperature computed for the period 1951-1980 (called the base-period) at that
17 position, is subtracted from the corresponding time series. Each time series is then divided by the
18 standard deviation computed at each position and during the same time period. The distribution of
19 the resulting data is then plotted for different periods. These distributions represent the deviation
20 of the temperature for a specific period, from the mean in the base period, in units of the standard
21 deviation in that period. The results showed that these distributions are wider for the recent time
22 periods compared to 1951-1980. [10] took a similar approach in analysing the ERA-40 data set.
23 However, in addition to the aforementioned method, they computed the distribution of the SDs in
24 an alternative way: for each position and each time period, the deviation of the time-series at that
25 position from the mean in that time period at that position was computed, and then divided by the SD
26 of that position in the period before 1981. The results showed that there still is an increase in the SDs
27 from 1958-1970 to 1991-2001, but this is much less than what is obtained from the method used in
28 [8]. The authors also computed the time-evolving global SD from the de-trended time-series at each
29 position. The resulting curve suggested that the global SD has been stable.

30 These previous work (and other related research, e.g., [16]) have several shortcomings. First, no
31 statistical analysis has been performed to examine if the change in the SD is statistically significant.
32 Second, the methodologies for computing the SDs are rather arbitrary. The deviation of each time-
33 series in a given period, is computed from either the mean of a base-period (as in [8]), or from the

given period (as in [10, 16]). These deviations are then normalized using the SD of the base-period or the given period. No justification is provided for these choices. Third, the correlation between the observations is ignored. The observations in subsequent days and close geographical positions could be highly correlated. Without considering these correlations, any conclusion based on the averaged data could be flawed.

The main contribution of this work is to develop a new methodology for detecting the trend in the volatility of spatio-temporal data. In this methodology, the variance at each position and time, is considered as a hidden (un-observed) variable. The value of these hidden variables are then estimated by maximizing the likelihood of the observed data. We show that this formulation per se, is not appropriate for detecting the trend. To overcome this issue, we penalize the differences between the estimated variances of the observations which are temporally and/or spatially close to each other. This will result in an optimization problem called the *generalized LASSO problem* [21]. As we will see, the dimension of this optimization problem is very high and so the standard methods for solving the generalized LASSO cannot be applied directly. We investigate two methods for solving this optimization problem. In the first method, we adopt an optimization technique called alternative direction method of multipliers (ADMM) [4], to divide the total problem into several sub-problems of much lower dimension and show how the total problem can be solved by iteratively solving these sub-problems. The second method, called the *linearized ADMM algorithm* [14] solves the main problem by iteratively solving a linearized version of it. We will compare the benefits of each method. Also neuroscience.

1.2 Related work

Mention [7, 12]. Also, [22, 23]. Find Sharpnack Tibshirani. ARCH/GARCH. [13, 17, 28] [15]

1.3 Main contributions

2 ℓ_1 -trend filtering for estimating variance of a time-series

ℓ_1 -trend filtering was proposed by [11] as a method for estimating a smooth, time-varying trend. It is formulated as the optimization problem

$$\min_{\beta} \frac{1}{2} \sum_{t=1}^T (y_t - \beta_t)^2 + \lambda \sum_{t=1}^{T-2} |\beta_t - 2\beta_{t+1} + \beta_{t+2}|$$

or equivalently:

$$\min_{\beta} \frac{1}{2} \|y - \beta\|_2^2 + \lambda \|D\beta\|_1 \quad (1)$$

where y_t is an observed time-series, β is the smooth trend, D is a $T - 2 \times T$ matrix, and λ is a tuning parameter which balances fidelity to the data (small errors in the first term) with a desire for smoothness. With the penalty matrix D , the estimated β will be piecewise linear. [11] proposed a specialized primal-dual interior point (PDIP) algorithm for solving (1). From a statistical perspective, (1) is a constrained maximum likelihood problem with independent observations from a normal distribution with common variance, $y_t \sim N(\beta_t, \sigma^2)$, subject to a piecewise linear constraint on β .

2.1 Extension to variance estimation

Inspired by the ℓ_1 -trend filtering algorithm, we propose a non-parametric model for estimating the variance of a time-series. To this end, we assume that at each time step t , there is a hidden variable h_t such that conditioned on h_t the observations y_t are independent normal variables with zero mean and variance $\exp(h_t)$. The negative log-likelihood of the observed data in this model is $l(y|h) \propto -\sum_{t=1}^T h_t - y_t^2 e^{-h_t}$. Crucially, we assume that the hidden variables h_t vary smoothly. To impose this assumption, we estimate h_t by solving the penalized, negative log-likelihood:

$$\min_h -l(y|h) + \lambda \|Dh\|_1 \quad (2)$$

75 where D has the same structure as above.

76 **TODO:** explain the objective more. give the AR(1) example. Explain what you loses by this
 77 assumption (ACF, forecasting). Also explain that the covariace matrix is diagonal so it cannot capture
 78 the covariance structure. But in contrast to spatial stat literature, it does not make any assumption on
 79 estimated variances. Compare to Hallac et al and Lingren et al.

80 As with (1), one can solve (2) using the PDIP algorithm (as in, e.g., `cvxopt` [1]). First, we note that
 81 this is a generalized LASSO problem [21]. The dual of a generalized LASSO with the objective
 82 $f(x) + \lambda \|Dx\|_1$ is:

$$\begin{aligned} \min_{\nu} \quad & f^*(-D^\top \nu) \\ \text{s.t.} \quad & \|\nu\|_\infty \leq \lambda \end{aligned}$$

83 where $f^*(\cdot)$ is the Fenchel conjugate of f : $f^*(u) = \max_x u^\top x - f(x)$. It is simple to show that

$$f^*(u) = \sum_t (u_t - 1) \log \frac{y_t^2}{1 - u_t} + u_t - 1. \quad (3)$$

84 Writing

$$r_w(v, \mu_1, \mu_2) := \begin{bmatrix} \nabla f^*(-D^\top v) + D(v - \lambda \mathbf{1})^\top \mu_1 - D(v + \lambda \mathbf{1})^\top \mu_2 \\ -\mu_1(v - \lambda \mathbf{1}) + \mu_2(v + \lambda \mathbf{1}) - w^{-1} \mathbf{1} \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

85 for $w > 0$, where μ_1 and μ_2 are dual variables for the ℓ_∞ constraint in the dual problem. As $w \rightarrow \infty$,
 86 the solution to this nonlinear system reduces to the KKT conditions. Therefore, the PDIP method
 takes Newton steps to solve the system for a series of increasing values of w . 1 provides the details.

Algorithm 1 PDIP for ℓ_1 variance estimation

Require: $\lambda > 0, w > 0, \nu \leftarrow 0, \mu_1 \leftarrow 0, \mu_2 \leftarrow 0, J \in \mathbb{Z}^+, \{w_k\}$ ▷ Initialization
for $k = 1, 2, \dots$ **do** ▷ Central path
 for $j = 1, \dots, J$ **do** ▷ Newton updates
 Solve $A[\Delta \nu \ \Delta \mu_1 \ \Delta \mu_2]^\top = r_{w_k}^j$ to find the search direction.
 A is the Jacobian of r_w .
 Update $r^{j+1} \leftarrow r^j + [\Delta \nu \ \Delta \mu_1 \ \Delta \mu_2]^\top$.
 end for
end for
return $h = -D^\top \nu$ **TODO: Is this right??**

87

88 2.2 Extension to spatio-temporal data

89 The method in the previous section can be used to estimate the variance of a single time-series. In
 90 this section, we extend this method to the estimation of the variance of spatio-temporal data.

91 At a specific time t , the data is measured on a grid of points with n_r rows and n_c columns. Let
 92 y_{ijt} denote the value of the observation at time t on the i^{th} row and j^{th} column of the grid, and h_{ijt}
 93 denote the corresponding hidden variable. We seek to impose both temporal and spatial smoothness
 94 constraints on the hidden variables. Specifically, we seek a solution for h which is piecewise linear in
 95 time and piecewise constant in space¹. This can be achieved by solving the following optimization
 96 problem:

¹The assumption that the variance is spatially piecewise constant simplifies the computations. The justification for this assumption is that we are interested in examining the trend in the variance over time and not over space.

$$\begin{aligned}
& \min_h \sum_{i,j,t} h_{ijt} + g_{ijt}^2 e^{-h_{ijt}} \\
& + \lambda_t \sum_{i,j} \sum_{t=1}^{T-2} |h_{ijt} - 2h_{ij(t+1)} + h_{ij(t+2)}| \\
& + \lambda_s \sum_{t,j} \sum_{i=1}^{n_r-1} |h_{ijt} - h_{(i+1)jt}| \\
& + \lambda_s \sum_{t,i} \sum_{j=1}^{n_c-1} |h_{ijt} - h_{i(j+1)t}|
\end{aligned} \tag{4}$$

97 The first term in the objective is the negative log-likelihood (minus a constant term). The second
 98 term is the temporal penalty for the time-series at each location (i, j) . The third and fourth terms,
 99 penalize the difference between the estimated variance of two vertically and horizontally adjacent
 100 points, respectively. This penalty is a special case of the penalty used in the *trend filtering on graphs*
 101 [29] (where the difference between the estimated values of the signal at each two nodes connected
 102 with an edge is penalized.). The optimization problem (4) can be written in the matrix form. Let h be
 103 a vector whose first T entries are h_{11t} for $t = 1, \dots, T$, its next T entries are h_{21t} and so on. Then the
 104 optimization problem in the matrix form is as follows:

$$\min_h -l(y|h) + \Lambda^t \|D_{total} h\|_1 \tag{5}$$

105 Let, \mathbf{e}_n denote a vector of size n with all entries being equal to 1. We have: $\Lambda^t = (\lambda_t \mathbf{e}_{n_t}^t | \lambda_s \mathbf{e}_{n_s}^t)$,
 106 $D_{total}^t = [D_t^t | D_s^t]$, where n_t and n_s are the number of rows of D_t and D_s , respectively (see below).
 107 The matrix D_t is the following block-diagonal matrix and corresponds to the temporal penalty:

$$D_t = \begin{bmatrix} D & & \\ & \ddots & \\ & & D \end{bmatrix}$$

108 where D was defined in (??). The number of the diagonal blocks is equal to the grid size $n_r \times n_c$.
 109 Each row of the matrix D_s corresponds to one spatial constraint in (4). For example, the first T rows
 110 correspond to $|h_{11t} - h_{21t}|$ for $t = 1, \dots, T$, the next T rows correspond to $|h_{11t} - h_{12t}|$, and so on.

111 **TODO: Detailed discussion of the PDIP to here**

112 The dual of this problem is:

$$\begin{aligned}
& \min_{\nu} f^*(-D_{total}^t \nu) \\
& \text{s.t. } |\nu_i| \leq \Lambda_i
\end{aligned} \tag{6}$$

113 where f^* is given in (3).

114 **3 Optimization**

115 For a spatial grid of size $n_r \times n_c$ and for T time steps, we have $n_t = 3n_r n_c - T n_c - 2n_r n_c$ and
 116 $n_s = n_r n_c T$. For a grid over the united states and for weekly averaged data of 10 years we have
 117 $n_r = 32$, $n_c = 68$, $T = 521$ and so $n_t \approx 3.5 \times 10^6$ and $n_s \approx 1.0 \times 10^6$. Therefore, the size of the
 118 optimization problem (6) is about $n_t \approx 4.5 \times 10^6$. In each step of the PDIP algorithm, we need to
 119 solve a linear system of equations in the form $z = Ax$ where A is a square matrix of size $2(n_t + n_s)$
 120 (see [3] equation 11.54). Therefore, applying the PDIP directly for solving the optimization problem
 121 (6) is infeasible.

122 In the next section, we develop an ADMM algorithm for solving this problem. The idea is to cast the
 123 problem (6) as a so-called *consensus optimization problem* [4] and solve it by breaking the problem

124 into smaller sub-problems using an ADMM-based algorithm. Next, we first give a brief overview
 125 of the consensus optimization problem and then explain how the problem (5) can be solved in this
 126 framework.

127 3.1 Consensus optimization

128 A general form consensus optimization problem is in the form $\min_z \sum_i f_i(z)$, $z \in \mathbb{R}^n$. We can
 129 define a set of *local variables* $x_i \in \mathbb{R}^{n_i}$ such that $\sum_i f_i(z) = \sum_i f_i(x_i)$. We follow the notation of
 130 [4] closely. Let $k = \mathcal{G}(i, j)$ which means that the j^{th} entry of x_i is z_k (or $(x_i)_j = z_k$) and define
 131 $\tilde{z}_i \in \mathbb{R}^{n_i}$ by $(\tilde{z}_i)_j = (x_i)_j$. Then the original unconstrained optimization problem is equivalent to
 132 the following constrained optimization problem:

$$\begin{aligned} \min_{\{x_1, \dots, x_N\}} \quad & \sum_i f_i(x_i) \\ \text{s.t.} \quad & \tilde{z}_i = x_i \end{aligned}$$

133 Now, we can apply ADMM to the *augmented Lagrangian* of this problem. This results in the
 134 following ADMM updating steps at each iteration m :

$$\begin{aligned} x_i^{m+1} &:= \underset{x_i}{\operatorname{argmin}} \left(f_i(x_i) + (u_i^m)^t x_i + (\rho/2) \|x_i - \tilde{z}_i^m\|_2^2 \right) \\ z_k^{m+1} &:= (1/S_k) \sum_{\mathcal{G}(i,j)=k} (x_i^{m+1})_j \\ u_i^{m+1} &:= u_i^m + \rho(x_i^{m+1} - \tilde{z}_i^{m+1}) \end{aligned} \tag{7}$$

135 Here, S_k is the number of local variable entries that correspond to z_k , and u_i are the Lagrange
 136 multipliers.

137 To solve the optimization problem (5) or (4) using the method explained in the previous section we
 138 need to two address two questions: first, how to choose the local variables x_i , and second, how to
 139 solve the optimization problem for updating these variables (the first line of (7)).

140 Obviously, the global variable in the problem (4) is $z = h$. The global variable can be index by
 141 the spatial coordinate and the temporal step. In Figure 1, z is represented as a cube. We can
 142 decompose z into sub-cubes. An example of this decomposition is shown in the figure by white
 143 lines. It is easy to see that with this definition of x_i , the objective (5) decomposes as $\sum_i f_i(x_i)$ where
 144 $f_i(x_i) = -l(y_i|x_i) + \|\Lambda_i^t D_i x_i\|_1$.

145 By this definition of x_i , the update step for x_i is the following optimization problem: $x_i^{m+1} :=$
 146 $\underset{x_i}{\operatorname{argmin}} (f_i(x_i) + (u_i^m)^t x_i + (\rho/2) \|x_i - \tilde{z}_i^m\|_2^2)$.

147 **TODO: Redo the below and move to the section on PDIP above**

148 We solve this using the PDIP method. To this end we first need to compute the dual of the problem.
 149 It can be shown that the dual of this problem is:

$$\begin{aligned} \min_{\nu} \quad & f_i^*(-D_i^t \nu_i) \\ \text{s.t.} \quad & |(\nu_i)_j| \leq (\Lambda_i)_j \end{aligned} \tag{8}$$

150 where:

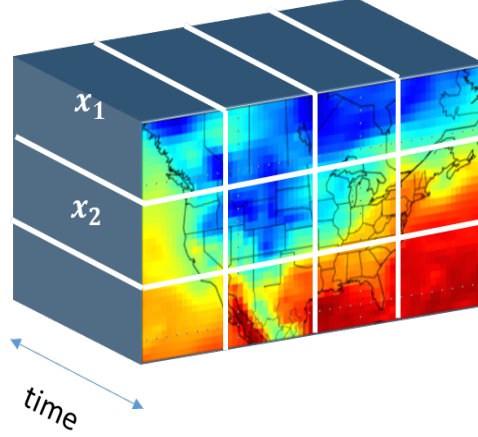


Figure 1: The cube represents the global variable z in space and time. The sub-cubes specified by the white lines are x_i .

Algorithm 2 ADMM for sparse estimation of variance of spatio-temporal data

Input: data y , mapping $\mathcal{G}(i, j)$, ρ , λ_t , λ_s

Initialization: $x_i^0 = z^0 = u_i^0 = \mathbf{0}$.

for $m = 1, 2, \dots$ **do**

for $i = 1$ **to** $N_{sub-cubes}$ **do**

 compute ν_i from (8)

 compute w_i from (9)

 set $x_i^m := w_i$

end for

 Compute z^m from (7)

 Compute u_i^m from (7)

end for

$$\begin{aligned}
 f_i^*(q) &= \sum_j (q_i)_j (w_i)_j - (w_i)_j - (y_i)_j^2 e^{-(w_i)_j} - \\
 &\quad (\rho/2) ((w_i)_j - (\alpha_i^m)_j)^2 \\
 \alpha_i^m &= \hat{z}_i^m - u_i^m \\
 (w_i)_j &= \mathcal{W} \left(\frac{(y_i)_j^2}{\rho} \exp \left[\frac{1 - q_j - \rho(\alpha_i^m)_j}{\rho} \right] \right) - \\
 &\quad \frac{1 - q_j - \rho(\alpha_i^m)_j}{\rho}
 \end{aligned} \tag{9}$$

151 In these equations, \mathcal{W} is the *Lambert function* [5]. To use the PDIP method, we also need to compute
 152 the gradient and Hessian of $f_i^*(-D_i^t \nu_i)$. This involves computing the derivatives of the Lambert
 153 functions. The primal solution x_i can be obtained from the dual solution ν_i in (8) by setting $x_i = w_i$
 154 where w_i is defined in the last equation of (9).

155 The complete ADMM algorithm for estimating the variances is represented in 2. All the computations
 156 in the three updating steps (7) can be performed in parallel. The number of rows and columns of
 157 the sub-cubes should be chosen so that the updating of x_i could be performed in one processor. We
 158 choose $3 \times 3 \times 521$ sub-cubes.

159 This algorithm divides the main optimization problem into several sub-problems and in each iteration
 160 solves these sub-problems and then performs a consensus step to make the solutions of these sub-

161 problems agree with each other. These sub-problems can be solved independently which makes this
 162 algorithm appealing for parallelization: in each iteration, each of these sub-problems can be sent into
 163 a separate computer and then all the solutions can be sent into a single computer which performs
 164 the consensus step. Therefore, in each iteration, the computation time will be equal to the time
 165 of solving each sub-problem plus the time of communicating the solutions to the master computer
 166 and then performing the consensus step. Since each sub-problem is small, with parallelization, the
 167 computation time in each iteration will be small. In addition, our experiments with several values of
 168 λ_t and λ_s showed that the algorithm converges in few hundreds iterations. Solving each sub-problem
 169 on a machine with four 3.20GHz Intel i5-3470 cores takes less than 3 seconds on average, and so
 170 for example if we assume that communication time is 10 seconds and the algorithm converges in
 171 300 iterations, with parallelization on $N_{sub-cubes}$ machines, the algorithm will converge in about 1
 172 hour. Assuming that we use $N_{sub-cubes}$ machines and that the convergence rate of the algorithm is
 173 independent of the grid size, this time will be independent of the grid size.

174 If we perform these computations on a single machine, the computation time grows linearly with
 175 $N_{sub-cubes}$. For example, for the data in a grid over the united states and using $3 \times 3 \times 521$ sub-cubes
 176 each iteration of the algorithm will take about 20 minutes on a single machine and so with 300
 177 iterations it will take several days to converge. Given that we need to compute the solution for several
 178 values of the parameters λ_t and λ_s , this computation time is not feasible.

179 Therefore, this algorithm is only useful if we can parallelize the computation over several machines.
 180 In the next section, we describe another algorithm which makes the computation feasible on a single
 181 machine.

182 3.2 Linearized ADMM

183 In this section, we describe *Linearized ADMM algorithm* [14] which, as we will see, makes the
 184 computation on a single machine feasible.

185 Consider the following optimization problem:

$$\min_x f(x) + g(Dx) \quad (10)$$

186 where $x \in \mathbb{R}^n$ and $D \in \mathbb{R}^{m \times n}$. Each iteration of the linearized ADMM algorithm for solving this
 187 problem has the following form:

$$\begin{aligned} x^{k+1} &:= \text{prox}_{\mu f}(x^k - (\mu/\rho)D^T(Dx^k - z^k + u^k)) \\ z^{k+1} &:= \text{prox}_{\rho g}(Dx^k + u^k) \\ u^{k+1} &:= u^k + Dx^{k+1} - z^{k+1} \end{aligned}$$

188 where $z, u \in \mathbb{R}^m$ and the *proximal operator* is defined as follows:

$$\text{prox}_{\alpha f}(u) = \min_x \alpha \cdot f(x) + \frac{1}{2} \|x - u\|^2$$

189 This algorithm belongs to the general class of *proximal algorithms* for solving convex optimization
 190 problems. For more details about these algorithms see [14].

191 The optimization problem (4) can be put into the form (10) as follows:

$$\begin{aligned} f(x) &:= \sum_k f_k(x_k) := \sum_k x_k + y_k^2 e^{-x_k} \\ g(z) &:= \sum_l g_l(z_l) := \sum_l \lambda_l |z_l| \\ z &= Dx \end{aligned}$$

192 where y_k is the k^{th} entry of the vector whose entries are y_{ijt} , and λ_l is the l^{th} entry of the vector
 193 $\Lambda^t = (\lambda_t \mathbf{e}_{n_t}^t | \lambda_s \mathbf{e}_{n_s}^t)$ (see Section 2.2).

194 To perform the steps in (11), we need to evaluate \mathbf{prox}_{μ_f} and $\mathbf{prox}_{\rho g}$. The proximal algorithms are
 195 feasible only if these proximal operators can be evaluated efficiently which, as we show next, is the
 196 case for our problem.

197 Let $(\mathbf{prox}_{\mu_f}(u))_k$ be the k^{th} entry of $\mathbf{prox}_{\mu_f}(u)$. From the *separable sum* property of the proximal
 198 operators we have (see [14], section 2.1):

199 **TODO: Make this a theorem**

$$\left(\mathbf{prox}_{\mu_f}(u) \right)_k = \mathbf{prox}_{\mu_{f_k}}(u_k)$$

200 Similarly, evaluating $\mathbf{prox}_{\rho g}$ reduces to evaluating the proximal operators of scalar functions. We
 201 have:

$$\mathbf{prox}_{\mu_{f_k}}(u_k) = \min_{x_k} \mu(x_k + y_k^2 e^{-x_k}) + \frac{1}{2}(x_k - u_k)^2$$

202 By setting the derivative with respect to x_k of the function to zero we obtain:

$$\mathbf{prox}_{\mu_{f_k}}(u_k) = \mathcal{W}\left(\frac{y_k^2}{\mu} \exp\left(\frac{1 - \mu u_k}{\mu}\right)\right) + \frac{1 - \mu u_k}{\mu}$$

203 where \mathcal{W} is the Lambert function.

204 Next we compute $\mathbf{prox}_{\mu_{g_l}}(u_l)$. The function $\rho\lambda_l|z_l| + 1/2(z_l - u_l)^2$ is not differentiable. However,
 205 at the optimal solution we have: $\rho \cdot \lambda_l \cdot \partial(|z_l|) = u_l - z_l$, where $\partial(|z_l|)$ is the sub-gradient of $|z_l|$.
 206 This results in the following solution:

$$\mathbf{prox}_{\rho g_l}(u_l) = \begin{cases} u_l - \rho\lambda_l & \text{if } u_l > \rho\lambda_l \\ 0 & \text{if } |u_l| \leq \rho\lambda_l \\ u_l + \rho\lambda_l & \text{if } u_l < -\rho\lambda_l \end{cases}$$

207 Therefore, both proximal operators in (11) can be evaluated easily and so we can use the linearized
 208 ADMM algorithm to solve the optimization problem (4).

209 4 Empirical evaluation

210 4.1 Simulations

211 To analyze the proposed optimization methods, we first fit the model on some synthetic spatio-
 212 temporal data. The observations at all time steps and all locations were generated from independent
 213 Gaussian random variables with zero mean. However, the variance of these random variables changes
 214 smoothly in time and space. Specifically, the observation at time step t and at location (r, c) on the
 215 grid is generated from a Gaussian $N(0, \sigma^2(t, r, c))$, where:

$$\sigma(t, r, c) = \sum_{s=1}^S W_s(t) \cdot \exp\left(\frac{(r - r_s)^2 + (c - c_s)^2}{2\sigma_s^2}\right)$$

$$W_s(t) = \alpha_s \cdot t + \exp(\sin(2\pi\omega_s t + \phi_s))$$

216 In words, the variance at each time and location is computed as the weighted sum of S bell-shaped
 217 functions where the weights are time-varying functions. Specifically, the weights consist of a linear
 218 trend $\alpha_s \cdot t$ and a periodic term $\beta_s \cdot \sin(2\pi\omega_s t + \phi_s)$. The bell-shaped functions impose the spatial
 219 smoothness, and the linear trend and the periodic terms enforce the temporal smoothness. We
 220 simulated the data on a 5 by 7 grid and for 780 time steps. We also used $S = 4$. The parameters of

Table 1: Parameters used to simulate data

s	r_s	c_s	σ_s	α_s	ω_s	ϕ_s
1	0	0	5	0.5	0.121	0
2	0	5	5	0.1	0.121	0
3	3	0	5	-0.5	0.121	$\pi/2$
4	3	5	5	-0.1	0.121	$\pi/2$

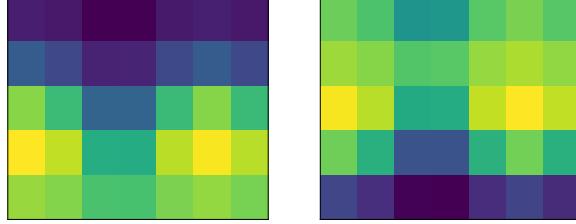


Figure 2: Variance function at $t = 25$ (left) and $t = 45$ (right)

the variance function are shown in Table 1. The value of the variance function for all locations on the grid and at $t = 25$ and $t = 45$ is shown in Figure 2. Also, the variance for all time steps and at the location (0,0) on the grid is shown in Figure 3. The linear trend and the period term can be seen clearly in this figure.

The left panel of Figure 3 shows the convergence of the two algorithms. Each iteration of the linearized algorithm takes about 0.01 seconds on average while each iteration of the consensus ADMM takes about 20 seconds.

We estimated the linearized ADMM for all combinations of values of λ_t and λ_s from the sets $\lambda_t \in \{0, 1, 5, 10, 50, 100\}$ and $\lambda_s \in \{0, 0.05, 0.1, 0.2, 0.3\}$. For each pair, we then compute the mean absolute error (MAE) between the estimated variance and the true variance at all locations and all time steps. For $\lambda_t = 5$ and $\lambda_s = 0.1$ MAE was minimized. The right panel of Figure 3 shows the true and the estimated standard deviation at location (0,0) using $\lambda_s = 0.1$ and $\lambda_t = 5$ (blue) and $\lambda_t = 100$ (green). As we can see, larger than optimal value of λ_t leads to estimated values which are “too smooth”.

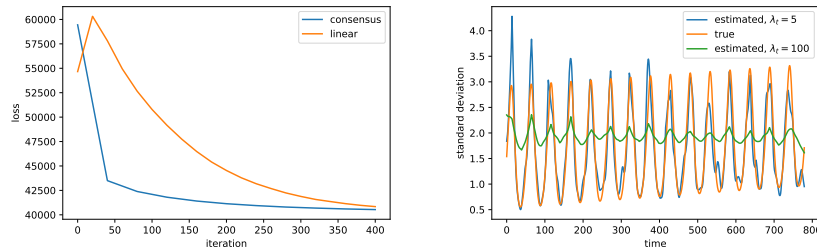


Figure 3: Left: Convergence speed of linearized and consensus ADMM. Right: The true (orange) and estimated standard deviation function at the location (0,0). The estimated values are obtained using linearized ADMM with $\lambda_s = 0.1$ and two values of λ_t : $\lambda_t = 5$ (blue) and $\lambda_t = 100$ (green).

4.2 Data analysis

As it was mentioned before, the algorithm proposed in Section 3.1 is appropriate only if we parallelize it over multiple machines and so we do not pursue it further here. All the results reported in this section are obtained using the linearize ADMM algorithm of Section 3.2. We applied this algorithm on a subset of the ERA-40 dataset. The data is the 2 meter temperature measured daily at 12 p.m

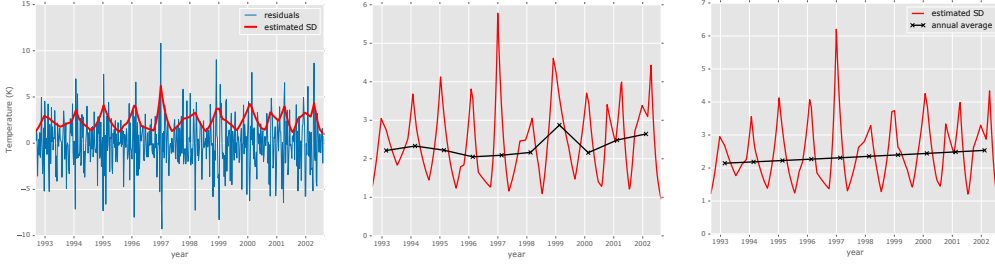


Figure 4: Left: The residuals of the time-series of Bloomington (averaged weekly) and the estimated SD obtained from the method of Section 2.1 (red). Middle: the estimated SDs (red) and their annual average (black) without imposing the long horizon penalty. Right: the same as middle panel but here the long horizon penalty is imposed. See the text for more details.

from August 31 of 1992 to 2002. To reduce the noise we first computed the weekly average of this data. To further reduce the size of the data, we will only analyze the data of the locations inside a rectangle extended from $(58^\circ, 226^\circ)$ to $(22^\circ, 302^\circ)$. This rectangle covers the united states. All the computations were performed on a machine with four 3.20GHz Intel i5-3470 cores.

Data Exploration The red curve in the left panel of Figure 4 shows the estimated SD (which is $\exp(h_t/2)$) of the residuals of the time-series of Bloomington. To reduce the number of time-steps in this figure and in the remainder of the paper we work on the weekly averaged of the data.

The curve of the estimated SD captures the periodic variations in the SD of the signal. Just by looking at this curve, it is hard to say if the SD is decreasing or increasing. Therefore, we compute the average of the estimated SD for each year. The estimated SD together with this annual average is shown in the middle panel of Figure 4. As it can be seen, the annual trend is not smooth. This is because in the optimization problem (2), the smoothness of the annual trend is not encouraged. To remedy this, we add the following long horizon penalty to (2):

$$\sum_{i=1}^{N_{year}-2} \left| \sum_{t=1}^{52} h_{t_1} - 2h_{t_2} + h_{t_3} \right| \quad (12)$$

where $t_1 = 52(i-1) + t$, $t_2 = 52i + t$ and $t_3 = 52(i+1) + t$. Also, N_{year} is the number of years over which we are performing our analysis (here $N_{year} = 10$). Since we are working on the weekly averaged data, each year corresponds to 52 observations. In the matrix form, the penalty (12) adds N_{year} rows to the matrix D . The estimated SDs using this penalty matrix is shown in the right panel of Figure 4. The annual average of the estimated SDs shows a linear trend with a positive slope.

This section is devoted to exploring some of the properties of the ERA-40 surface temperature data set. The goal here is to demonstrate some of the difficulties in modeling the trend in the temperature volatility and motivate our methodology for doing so.

The right panel of Figure 5 shows the time-series of the temperature of Bloomington, after removing the cyclic terms and de-trending using the method explained in the next section. The goal is to investigate the trend in the variance of this signal. This figure, reveals another issue toward this goal: the variance of this signal, shows cyclic behavior. Also, the cycles are not regular and their amplitude and frequency change. Even if one can describe the behavior of the variance of all the time-series using a single parametric model (for example a variant of the GARCH models [2]), it is not clear how the trend in the variance should be investigated in this framework. These observations motivate the need to develop a non-parametric framework for the problem at hand.

Convergence We used the following rule to determine when to stop the optimization: the optimization was stopped if the value of the loss did not improve by at least 0.1% in 1000 trials. As we can see, the algorithm converged in about 2000 iterations. This took about 11 minutes. Our experiments

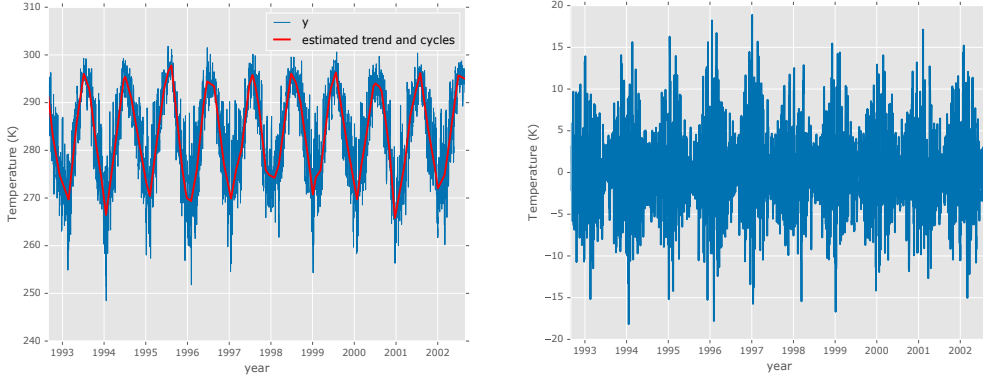


Figure 5: Left: Time-series of the temperature of Bloomington (blue) and the estimated trend and cycles obtained from the ℓ_1 -trend filtering (red). Right: the same time-series after removing the cyclic terms and de-trending using ℓ_1 -trend filtering.

showed that the convergence speed depends on the value of λ_t and λ_s . Also, if we use the solution obtained for smaller values of these parameters as the initial value for the larger values (*warm start*), the converges speed improves.

Model selection One common method for choosing the penalty parameters in the Lasso problems is to find the solution for a range of the values of these parameters and then choose the values which minimize a model selection criterion. However, such analyses needs the computation of the degrees of freedom (df). Several previous work have investigated the df in generalized lasso problems [9, 24, 30]. However, all these studies have considered the linear regression problem and, to the best of our knowledge, the problem of computing the df for generalized lasso with general objective function has not been considered yet.

Another approach is to choose the set of values which minimize an estimate of the expected prediction error obtained by k-fold cross-validation [20]. Although this method is applicable for our problem, it needs k times more computation.

In this paper, we use a heuristic method for choosing λ_t and λ_s : we compute the optimal solution for a range of values of these parameters and choose the values which minimize $\mathcal{L}(\lambda_t, \lambda_s) = -l(y|h) + \sum \|D_{total}h\|$. This objective is a compromise between the negative log likelihood ($-l(y|h)$) and the complexity of the solution ($\sum \|D_{total}h\|$). For smoother solutions the value of $\sum \|D_{total}h\|$ will be smaller but with the cost of larger $-l(y|h)$.

We computed the optimal solution for all the combinations of the following sets of values: $\lambda_t \in \{1, 5, 10, 20\}$, $\lambda_s \in \{0, .1, 1, 5, 10\}$. The best combination based on a held out set was $\lambda_t = 5$ and $\lambda_s = 1$. All the analyses in the next section are performed on the solution for these values.

Analysis of trend of temperature volatility The top row of Figure 6 shows the detrended data, the estimated standard deviation and the yearly average of these estimates for two cities in the US. The estimated SD captures the periodic behavior in the variance of the time-series. In addition, the number of linear segments changes adaptively in each time window depending on how fast the variance is changing. The yearly average of the estimated SD captures the trend in the temperature volatility. For example, we can see that in Bloomington, there is a small positive trend. To determine how the volatility has changed in each location, we subtract the average of the estimated variance in 1992 from the average in the following years and compute their sum. The value of this change in the variance in each location is depicted in the right panel of Figure 6. The left panel of this figure, shows the average estimated variance in each location.

It is interesting to note that the trend in volatility is almost zero over the oceans. The most positive trend can be observed in the south-east and the most negative trend has happened in the north-east.

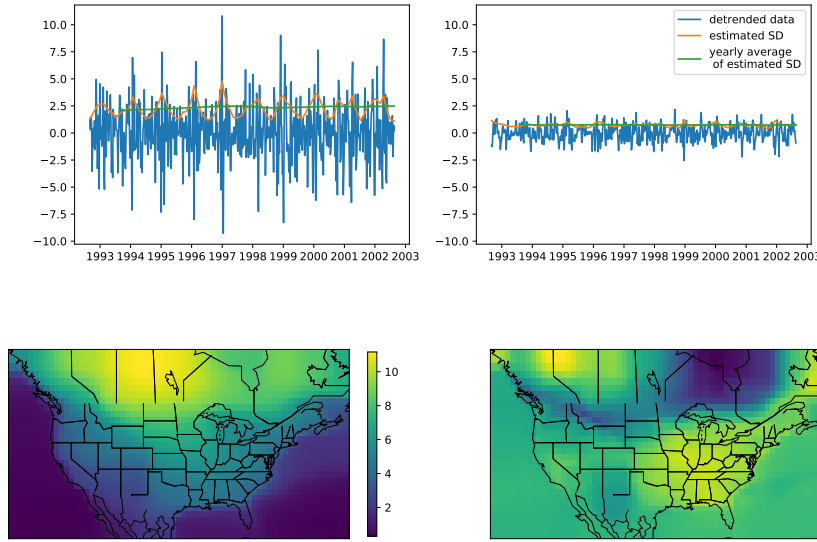


Figure 6: Top row: Detrended data and the estimated SD for Bloomington (left) and San Diego (right). Bottom: the average of the estimated variance over the US (left) and the change in the variance from 1992 to 2002 (right)

5 Discussion

In this paper, we proposed a new method for estimating the variance of spatio-temporal data. The main idea is to cast this problem as a constrained optimization problem where the constraints enforce smooth changes in the variance for neighboring points in time and space. In particular, the solution is piecewise linear in time and piecewise constant in space. The resulting optimization is in the form of a generalized LASSO problem with high-dimension, and so applying the PDIP method directly is infeasible. We therefore developed two ADMM-based algorithms to solve this problem: the consensus ADMM and linearized ADMM.

The consensus ADMM algorithm converges in few hundreds of iterations but each iteration takes much longer than the linearized ADMM algorithm. The appealing feature of the consensus ADMM algorithm is that if it is parallelized on enough number of machines the computation time per iteration remains constant as the problem size increases. The linearized ADMM algorithm, on the other hand converges in few thousands of iterations but each iteration is performed in split second. However, since the algorithm converges in many iterations it is not very appropriate for parallelization. The reason is that after each iteration the solution computed in each machine should be broadcast to the master machine and this operation takes some time which depends on the speed of the network connecting the slave machines to the master. A direction for future research would be to combine these two algorithms in the following way: the problem should be split into the sub-problems (as in the consensus ADMM) but each sub-problem can be solved using linearized ADMM.

We applied the linearized ADMM algorithm to the surface temperature data on a grid over the united states, for years 1992-2002. The results showed that in many locations the variance of the temperature has increased about 1 unit in 10 years.

The goal of this paper, however, is not to make any conclusions about the trend in the variance because we solved the problem only for a grid over the united states and for 10 years of the data. A thorough analysis, needs the full solution over the globe and for a longer time period. The goal of the paper, was to propose the idea of estimating the trend in variance of spatio-temporal signals using generalized lasso and to investigate the algorithms for solving the resulting optimization problem.

References

- [1] M. S. Andersen, J. Dahl, and L. Vandenberghe. CVXOPT: A Python package for convex optimization, version 1.1. 6. Available at cvxopt.org 54, 2013.
- [2] T. Bollerslev. Generalized autoregressive conditional heteroskedasticity. *Journal of Econometrics*, 31(3): 307–327, Apr. 1986. ISSN 0304-4076.
- [3] S. Boyd and L. Vandenberghe. *Convex optimization*. Cambridge university press, 2004.
- [4] S. Boyd, N. Parikh, E. Chu, B. Peleato, and J. Eckstein. Distributed Optimization and Statistical Learning via the Alternating Direction Method of Multipliers. *Foundations and Trends in Machine Learning*, 3(1): 1–122, 2011. ISSN 1935-8237.
- [5] R. M. Corless, G. H. Gonnet, D. E. G. Hare, D. J. Jeffrey, and D. E. Knuth. On the LambertW function. *Advances in Computational Mathematics*, 5(1):329–359, Dec. 1996.
- [6] E. M. Fischer, U. Beyerle, and R. Knutti. Robust spatially aggregated projections of climate extremes. *Nature Climate Change*, 3:1033–1038, 2013. URL <http://dx.doi.org/10.1038/nclimate2051>.
- [7] D. Hallac, Y. Park, S. Boyd, and J. Leskovec. Network inference via the time-varying graphical lasso. In *Proceedings of the 23rd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, KDD ’17, pages 205–213, New York, NY, USA, 2017. ACM. doi: 10.1145/3097983.3098037. URL <http://doi.acm.org/10.1145/3097983.3098037>.
- [8] J. Hansen, M. Sato, and R. Ruedy. Perception of climate change. *Proceedings of the National Academy of Sciences*, 109(37), Sept. 2012.
- [9] Q. Hu, P. Zeng, and L. Lin. The dual and degrees of freedom of linearly constrained generalized lasso. *Computational Statistics & Data Analysis*, 86:13–26, June 2015.
- [10] C. Huntingford, P. D. Jones, V. N. Livina, T. M. Lenton, and P. M. Cox. No increase in global temperature variability despite changing regional patterns. *Nature*, 500(7462):327–330, Aug. 2013. ISSN 0028-0836.
- [11] S. Kim, K. Koh, S. Boyd, and D. Gorinevsky. ℓ_1 Trend Filtering. *SIAM Review*, 51(2):339–360, May 2009. ISSN 0036-1445. doi: 10.1137/070690274. URL <http://epubs.siam.org/doi/abs/10.1137/070690274>.
- [12] S.-J. Kim, K. Koh, S. Boyd, and D. Gorinevsky. ℓ_1 trend filtering. *SIAM Review*, 51(2):339–360, 2009. doi: 10.1137/070690274. URL <https://doi.org/10.1137/070690274>.
- [13] K. Lin, J. L. Sharpnack, A. Rinaldo, and R. J. Tibshirani. A sharp error analysis for the fused lasso, with application to approximate changepoint screening. In I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, editors, *Advances in Neural Information Processing Systems 30*, pages 6884–6893. Curran Associates, Inc., 2017. URL <http://papers.nips.cc/paper/7264-a-sharp-error-analysis-for-the-fused-lasso-with-application-to-approximate-changepoint-screening.pdf>.
- [14] N. Parikh and S. Boyd. Proximal Algorithms. *Foundations and Trends® in Optimization*, 1(3):127–239, Jan. 2014.
- [15] A. Ramdas and R. J. Tibshirani. Fast and flexible admm algorithms for trend filtering. *Journal of Computational and Graphical Statistics*, 25(3):839–858, 2016.
- [16] A. Rhines and P. Huybers. Frequent summer temperature extremes reflect changes in the mean, not the variance. *Proceedings of the National Academy of Sciences*, 110(7):E546–E546, Feb. 2013.
- [17] V. Sadhanala, Y.-X. Wang, J. L. Sharpnack, and R. J. Tibshirani. Higher-order total variation classes on grids: Minimax theory and trend filtering methods. In I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, editors, *Advances in Neural Information Processing Systems 30*, pages 5800–5810. Curran Associates, Inc., 2017. URL <http://papers.nips.cc/paper/7162-higher-order-total-variation-classes-on-grids-minimax-theory-and-trend-filtering-methods.pdf>.
- [18] J. A. Screen. Arctic amplification decreases temperature variance in northern mid- to high-latitudes. *Nature Climate Change*, 4:577–582, 2014. URL <http://dx.doi.org/10.1038/nclimate2268>.
- [19] P. W. Staten, B. H. Kahn, M. M. Schreier, and A. K. Heidinger. Subpixel characterization of HIRS spectral radiances using cloud properties from AVHRR. *Journal of Atmospheric and Oceanic Technology*, 33(7): 1519–1538, 2016. doi: 10.1175/JTECH-D-15-0187.1.

- 383 [20] R. Tibshirani. Regression Shrinkage and Selection via the Lasso. *Journal of the Royal Statistical Society.*
384 *Series B (Methodological)*, 58(1):267–288, 1996.
- 385 [21] R. J. Tibshirani. *The Solution Path of the Generalized Lasso*. PhD Thesis, Stanford University, 2011.
- 386 [22] R. J. Tibshirani. Adaptive piecewise polynomial estimation via trend filtering. *Annals of Statistics*, 42:
387 285–323, 2014. URL <http://www.stat.cmu.edu/~ryantibs/papers/trendfilter.pdf>.
- 388 [23] R. J. Tibshirani and J. Taylor. The solution path of the generalized lasso. *Annals of Statistics*, 39(3):
389 1335–1371, 2011.
- 390 [24] R. J. Tibshirani and J. Taylor. Degrees of freedom in lasso problems. *The Annals of Statistics*, 40(2):
391 1198–1232, 2012.
- 392 [25] K. E. Trenberth, Y. Zhang, J. T. Fasullo, and S. Taguchi. Climate variability and relationships between
393 top-of-atmosphere radiation and temperatures on earth. *Journal of Geophysical Research: Atmospheres*,
394 120(9):3642–3659, 2014. doi: 10.1002/2014JD022887.
- 395 [26] S. M. Uppala, P. W. K  llberg, A. J. Simmons, U. Andrae, and e. al. The ERA-40 re-analysis. *Quarterly*
396 *Journal of the Royal Meteorological Society*, 131(612):2961–3012, Oct. 2005.
- 397 [27] D. A. Vasseur, J. P. DeLong, B. Gilbert, H. S. Greig, C. D. G. Harley, K. S. McCann, V. Savage, T. D.
398 Tunney, and M. I. O’Connor. Increased temperature variation poses a greater risk to species than climate
399 warming. *Proceedings of the Royal Society of London B: Biological Sciences*, 281(1779), 2014. doi:
400 10.1098/rspb.2013.2612.
- 401 [28] Y.-X. Wang, J. Sharpnack, A. J. Smola, and R. J. Tibshirani. Trend filtering on graphs. *Journal of Machine*
402 *Learning Research*, 17(105):1–41, 2016. URL <http://jmlr.org/papers/v17/15-147.html>.
- 403 [29] Y.-X. Wang, J. Sharpnack, A. J. Smola, and R. J. Tibshirani. Trend Filtering on Graphs. *Journal of Machine*
404 *Learning Research*, 17(105):1–41, 2016. URL <http://jmlr.org/papers/v17/15-147.html>.
- 405 [30] P. Zeng, Q. Hu, and X. Li. Geometry and Degrees of Freedom of Linearly Constrained Generalized Lasso.
406 *Scandinavian Journal of Statistics*, 44(4):989–1008, Nov. 2017. ISSN 0303-6898.