

# TemperatureImputations

Maxime Rischard

August 31, 2017

## Contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
1.1	Illustration of bias induces by measurement hour . . . . .	2
1.2	Proposed solution . . . . .	3
<b>2</b>	<b>First Spatiotemporal Model</b>	<b>5</b>
2.1	Fitting the spatiotemporal model . . . . .	5
<b>3</b>	<b>Predictions using nearby data</b>	<b>6</b>
<b>4</b>	<b>Imputations</b>	<b>7</b>
4.1	Imputing by Conditioning on Extrema . . . . .	7
4.2	Illustration of Hamiltonian Monte Carlo with Smoothmax Approximation . . . . .	9
4.3	Smoothmax Temperature Model . . . . .	12
<b>5</b>	<b>Model diagnostics</b>	<b>14</b>
5.1	Variogram . . . . .	14
5.2	Error and expected error . . . . .	14
<b>6</b>	<b>Improving model</b>	<b>15</b>
<b>7</b>	<b>Analysis</b>	<b>17</b>
<b>8</b>	<b>Inference on measurement hour</b>	<b>17</b>
<b>9</b>	<b>Appendices</b>	<b>19</b>
9.1	Stan programs for illustration of smoothmax . . . . .	19
9.1.1	Without smoothmax Approximation . . . . .	19
9.1.2	With smoothmax Approximation . . . . .	20
<b>10</b>	<b>Stan model for temperature imputations</b>	<b>20</b>

## 1 Introduction

Long, high-quality records of temperature provide an important basis for our understanding of climate variability and change. Historically, there has been a focus on monthly-average temperature records, which are sufficient for certain analyses, such as quantifying long-term changes in temperature. As our knowledge of climate change expands, however, there is increasing interest in understanding changes in temperature on shorter timescales, with a particular focus on extreme events. To do so, it is necessary to utilize higher-resolution temperature data.

Recent work has led to the development of the Global Historical Climatology Network-Daily (GHCND) database (Menne et al., 2012), which contains, among other variables, daily maximum and minimum temperatures from weather stations around the globe. The database draws from a range of different sources, and the data within it undergoes basic quality control to remove erroneous values.

The current quality control methodology, however, does not account for so-called ‘inhomogeneities’. Inhomogeneities result from changes in measurement practices that impact the recorded temperatures. For temperature, known inhomogeneities include (a) changes in the time of observation, (b) changes in the thermometer technology, (c) station relocation, and (d) changes in land use around a station (Menne et al., 2009). While these inhomogeneities have a small effect on, e.g., the estimation of global mean temperature, they can have a large effect on estimation of temperature variability and change at a more local scale.

There is a large body of work focused on homogenizing monthly-average temperatures (e.g., Karl et al., 1986; Easterling et al., 1996; Peterson et al., 1998; Ducré-Robitaille et al., 2003; Menne and Williams Jr, 2009; Vincent et al., 2012), resulting in widely available, large-scale homogenized monthly temperature datasets. Homogenization typically proceeds through identifying non-climatic ‘breakpoints’ in a given time series through comparison with neighboring stations. Once a breakpoint is identified, the measurements recorded after the breakpoint are adjusted in some way to reduce or remove the inhomogeneity. Most applications of these methods, however, focus on adjusting the mean state of the data rather than the shape of the distribution (see Della-Marta and Wanner, 2006, and references therein). While this may be sufficient for monthly data, it is known that certain changes in measurement practices affect different percentiles of daily temperature in different ways. To address this issue, some homogenization methods have also employed percentile matching techniques, wherein the adjustment to a timeseries after a breakpoint is a function of percentile (Della-Marta and Wanner, 2006; Trewin, 2013).

Here, we focus primarily on addressing the time of observation bias, as well as its time trend, because of its known impact on the distribution of daily maximum and minimum temperature ( $T_x$  and  $T_n$ ) measurements. The bias exists because  $T_x$  and  $T_n$  are often recorded by an observer who visits a weather station every 24 hours, and notes the maximum and minimum temperatures measured by the thermometer over the previous 24 hours. Ideally, the observer would visit the station at midnight, and the highest and lowest temperatures over the past 24 hours would typically be representative of the high and low during the prior day. For convenience, however, most observers record data at a daytime hour instead. As can be seen in Fig. XX, measurements recorded in the early morning may not properly register the low of the night before if it was unusually warm. Similarly, measurements recorded in the late afternoon may not properly register the high of the prior day if it was usually cool. In both cases, this will lead to a reduction in the variance of  $T_x$  and  $T_n$  distributions, but the effect will be greater at low (high) percentiles for  $T_x$  ( $T_n$ ).

If the time of observation remained constant over time, the bias would still exist, but it would not be linked to spurious trends in the data. However, there have been known (and likely unknown) changes in the time of observation. In the United States, for example, observers were instructed to switch from recording data in the afternoon to recording data in the morning beginning in the 1950s. This change has led to an apparent decrease in both  $T_x$  and  $T_n$  over time (Menne et al., 2009).

The goal of our approach is to infer the true  $T_n$  and  $T_x$  values throughout the data records, thereby correcting both the variance biases and the spurious trends. This stands in contrast to previous work, which has focused primarily on addressing spurious trends. We approach the problem as a missing data problem, wherein we are trying to recover the values of  $T_x$  and  $T_n$  that may have been overwritten due to measurement practices. Furthermore, by employing a Gaussian process framework and nearby stations with hourly data, we are able to simulate multiple realizations of temperature timeseries at each station, thereby providing estimates of uncertainty.

## 1.1 Illustration of bias induces by measurement hour

(Baker, 1975).

We illustrate the measurement bias in daily maxima and minima with ten days of hourly temperature measurements from the Waterloo Municipal Airport station in Iowa. Ideally,  $T_x$  measurements should capture the peak of each diurnal cycle, and  $T_n$  its trough. In Figure 1, those ideal measurements are indicated by the red and blue triangles respectively. The actual measurements are obtained by dividing the data into 24 hour measurement windows, and extracting the minimum and maximum. For each window, we plot

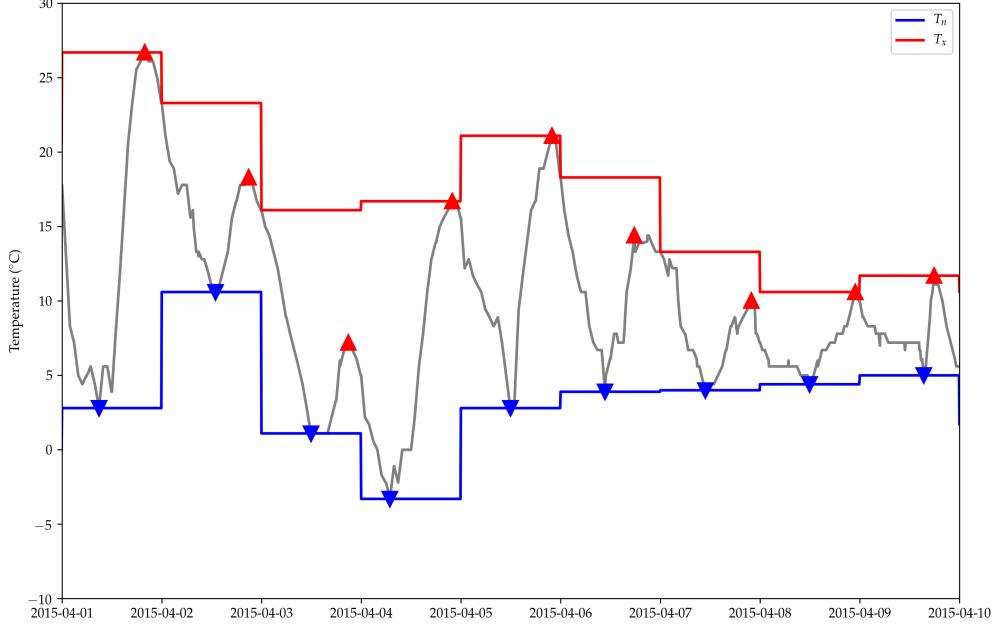


Figure 1:

these extrema with a red and blue horizontal line.

On most days, the ideal measurement and the actual measurement coincide: the triangle is on that day's line. But there are also several misses. The most blatant example occurs on April 3rd, where the peak of the diurnal cycle is 7.2°C and occurs at 21:00 UTC. However, because the previous day was much warmer, the day's  $T_x$  record of 16.1°C is reached immediately after the previous day's measurement. The measured  $T_x$  therefore overestimates the diurnal cycle's peak by 8.9°C.

This subtle bias in the daily records can in turn bias long-term summary statistics that are of climatological interest. A measure as simple as the average daily maximum temperature for an entire year (2015) increases by over 1°C if the measurements are made at the warmest time of day 21:00 UTC rather than 14:00 UTC (see Figure 2). Conversely, the average  $T_n$  is colder by over 1°C if  $T_n$  is measured at 10:00 UTC (the coldest time of day on average) rather than 17:00 UTC.

A climatologist studying weather variability might be interested in summary statistics such as the average absolute change in the daily temperature maxima and minima from one day to the next. The answer to that question too depends on the time of day at which the temperatures are recorded. Collecting the measurements at the hottest time of day means that the peaks on a warm day gets recorded twice, erasing the diurnal peaks of the following colder day, and hence the variability gets underestimated. We can see this in Figure 3), where the respective variability estimates drop if the maxima get measured at the warmest time, or if the minima get measured at the coldest time.

## 1.2 Proposed solution

We have seen that the daily maxima and minima do not faithfully record each diurnal cycle's peak and trough. The peaks on a relatively cold day can get overwritten by temperatures at either end of the measurement window that properly belong to the previous or the next diurnal cycle. Troughs on relatively warm days can be similarly overwritten. Our goal is to undo this damage, and recover estimates of summary statistics, such as the average daily maximum temperature, that do not suffer from the consequent bias. We need to address the erasure of information caused by the measurement mechanism, and therefore view this as a missing data problem.

Taking the missing data perspective, we seek to impute the hourly temperatures that have been replaced by a maximum and minimum over a 24 hour period. To do so, we combine information from two sources:

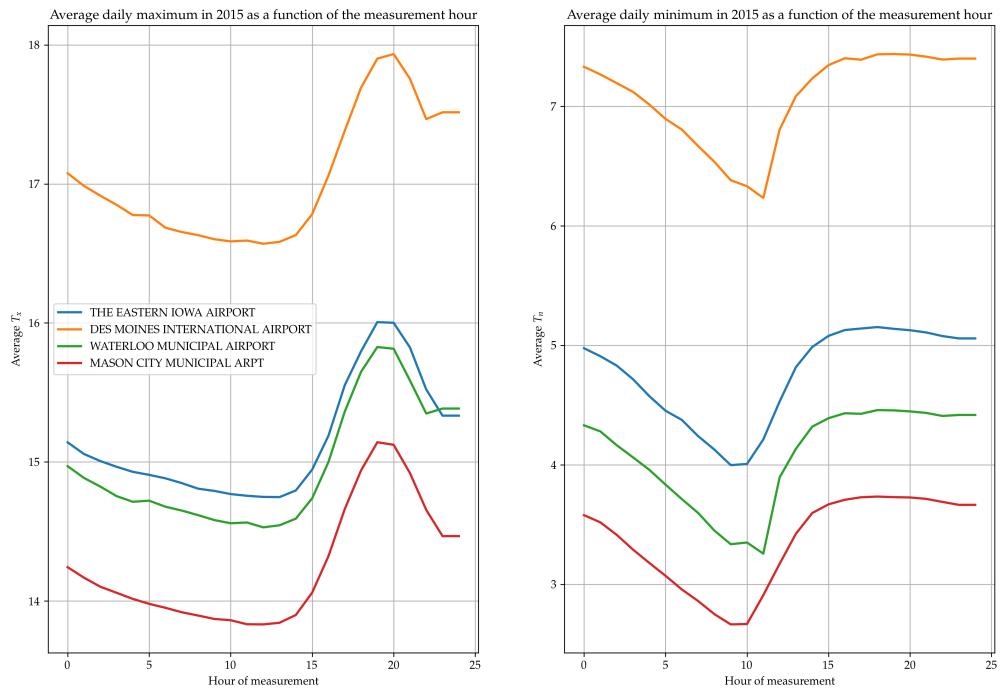


Figure 2:

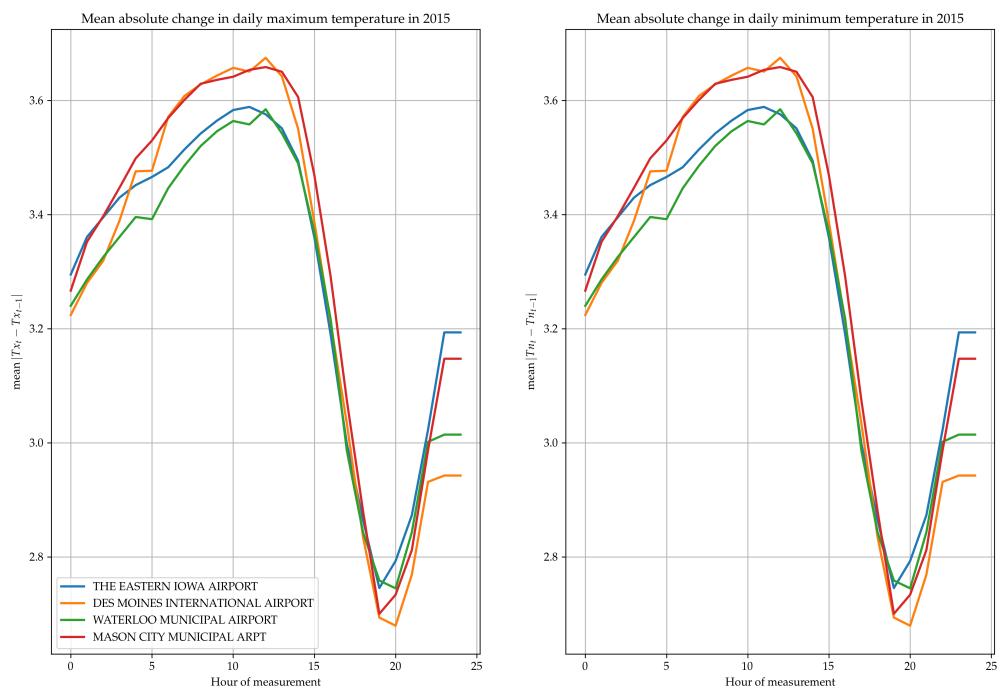


Figure 3:

the recorded daily temperature extremes at the station of interest, and also hourly temperatures recorded at nearby meteorological stations. These hourly measurements are considered less reliable by climatologists, as they aren't as carefully documented, calibrated, and situated. The meteorological stations are often in locations (like airports) where human activity will affect temperatures. Therefore, summary statistics extracted directly from those measurements would not be directly usable for climatology, as they could suffer from systematic bias. However, even if miscalibrated, the meteorological data do contain valuable information about the hourly changes in temperatures on any given day. We therefore use them to inform the shape of the imputed temperature time-series at our location of interest, while we use the recorded temperature extrema to calibrate and constrain them.

## 2 First Spatiotemporal Model

To model measured temperatures at various locations and times, we use a spatio-temporal Gaussian process model. In its simplest form, we believe that temperatures from stations that are near each other are more correlated than distant stations, and that those correlations should also decay in time. In the spatial statistics literature, squared exponential covariance functions are commonly used to model correlations decaying as a function of distance. Ignoring the time dimension, we would model the simultaneous temperatures throughout a region as a Gaussian process, with the covariance of two locations  $\mathbf{x}$  and  $\mathbf{x}'$

$$\text{cov} (T(\mathbf{x}), T(\mathbf{x}') | t) = k_{\text{space}}(\mathbf{x}, \mathbf{x}') = \sigma_{\text{GP}}^2 \exp \left( -\frac{(\mathbf{x} - \mathbf{x}')^\top (\mathbf{x} - \mathbf{x}')}{2\ell_x^2} \right). \quad (1)$$

Similarly, ignoring the spatial dimension, the time series of temperatures at a single location can be modeled as a Gaussian process with covariance between two moments  $t$  and  $t'$

$$\text{cov} (T(t), T(t') | \mathbf{x}) = k_{\text{time}}(t, t') = \sigma_{\text{GP}}^2 \exp \left( -\frac{(t - t')^2}{2\ell_t^2} \right). \quad (2)$$

We then combine the spatial and temporal model by multiplying the covariance functions

$$k_{\text{st}}(\mathbf{x}, \mathbf{x}', t, t') = k_{\text{time}}(t, t') \cdot k_{\text{space}}(\mathbf{x}, \mathbf{x}'). \quad (3)$$

This gives us the covariance of the Gaussian process underlying the full spatio-temporal model of temperatures. To complete the model specification, we add a mean temperature for each station  $\mu_{\text{station}[i]}$ , and iid measurement noise  $\epsilon_i$ .

$$T_i = \mu_{\text{station}[i]} + f(\mathbf{x}_i, t_i) + \epsilon_i \quad (4)$$

$$f(\mathbf{x}_i, t_i) \sim \mathcal{GP}(0, k_{\text{st}}(\mathbf{x}, \mathbf{x}', t, t')) \quad (5)$$

$$\epsilon_i \stackrel{\text{iid}}{\sim} \mathcal{N}(0, \sigma_\epsilon^2) \quad (6)$$

$$(7)$$

### 2.1 Fitting the spatiotemporal model

Software is readily available in many programming languages for fitting Gaussian process models, including inference on the covariance parameters. We chose to use the julia `GaussianProcesses.jl` package to fit the above spatiotemporal model to the hourly temperatures at four Iowa weather stations. However, the Iowa data includes 47,864 measurements, which is computationally infeasible to fit directly with a single Gaussian process. While approximation techniques exist to fit such large datasets, we chose the less efficient but simpler approach of dividing the data into 10-day chunks, modeled as independent Gaussian processes with shared hyperparameters. To simplify our implementation, we replaced the  $\mu_{\text{station}[i]}$  terms by a spatial squared exponential component

$$k_\mu(\mathbf{x}, \mathbf{x}') = \sigma_\mu^2 \exp \left( -\frac{(\mathbf{x} - \mathbf{x}')^\top (\mathbf{x} - \mathbf{x}')}{2\ell_\mu^2} \right) \quad (8)$$

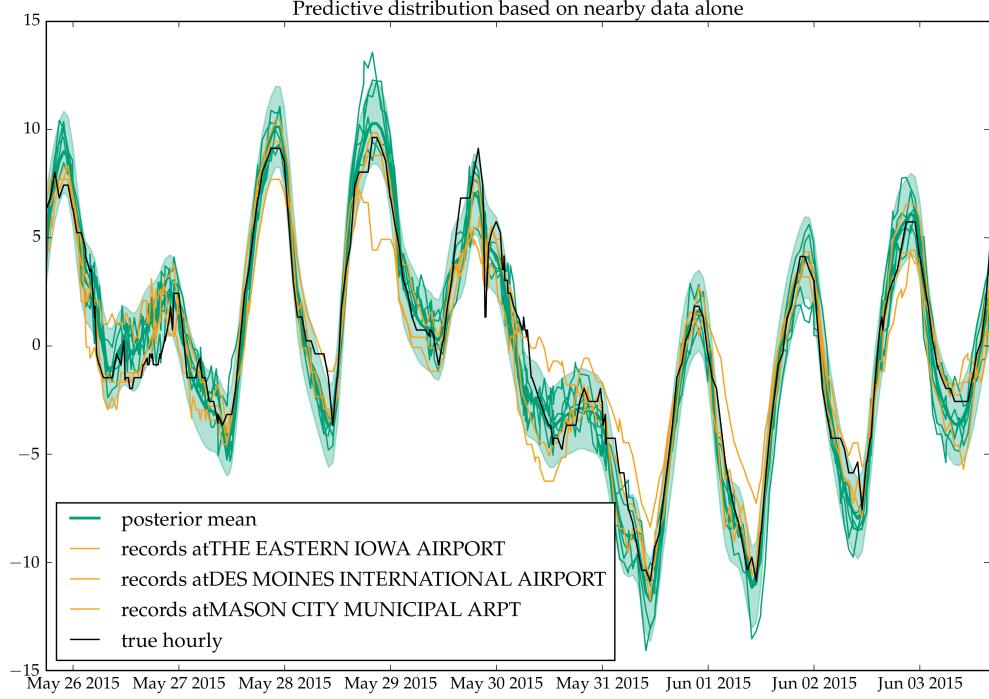


Figure 4: Predictive distribution using only nearby data and the simple product of square exponentials model. The orange lines are the measurements at nearby stations that are being used to inform the predictions. The black line is the true temperatures that have been withheld from the model, while the green line is the posterior mean of the predictions. The credible range (in green) is twice the standard deviations extracted from the diagonal entries of the posterior covariance matrix.

with large variance  $\sigma_\mu^2$  and low lengthscale  $\ell_\mu$  added to the covariance function so that the spatio-temporal kernel becomes

$$k_{st}(\mathbf{x}, \mathbf{x}', t, t') = k_{time}(t, t') \cdot k_{space}(\mathbf{x}, \mathbf{x}') + k_\mu(\mathbf{x}, \mathbf{x}'). \quad (9)$$

The model therefore has 4 free parameters:  $\sigma_{GP}$ ,  $\ell_t$ ,  $\ell_x$  and  $\sigma_\epsilon$ . We optimize the marginal likelihood of the Iowa data as a function of these three parameters

$$\hat{\sigma}_{GP}, \hat{\ell}_t, \hat{\ell}_x, \hat{\sigma}_\epsilon = \arg \max_{\sigma_{GP}, \ell_t, \ell_x, \sigma_\epsilon} \{ \mathbb{P}(Y | \sigma_{GP}, \ell_t, \ell_x, \sigma_\epsilon) \}, \quad (10)$$

and obtain  $\hat{\sigma}_{GP} = 3.73^\circ \text{C}$ ,  $\hat{\ell}_t = 2.7 \text{ hours}$ ,  $\hat{\ell}_x = 176.4 \text{ km}$  and  $\hat{\sigma}_\epsilon = 0.44^\circ \text{C}$ .

### 3 Predictions using nearby data

Once we have a spatio-temporal Gaussian process model with optimized covariance parameters, we can use it to generate predictions at the station where we aim to generate imputations based on nearby measurements. Gaussian processes make this a closed-form procedure. We'll denote the temperatures we wish to impute as  $T_{\text{miss}}$  at times  $t_{\text{miss}}$  and location  $\mathbf{x}_{\text{miss}}$  and those observed at nearby stations as  $T_{\text{nearby}}$ , at times  $t_{\text{nearby}}$  and locations  $\mathbf{x}_{\text{nearby}}$ . Under the spatio-temporal model,  $T_{\text{miss}}$  and  $T_{\text{nearby}}$  are jointly multivariate normal, with mean zero and covariance given by  $k_{st}(\mathbf{x}, \mathbf{x}', t, t')$ . Standard results for conditioning within multivariate normals then yields

$$\begin{aligned}
T_{\text{miss}} \mid T_{\text{nearby}} &\sim \mathcal{N}(\mu_{\text{miss}|\text{nearby}}, \Sigma_{\text{miss}|\text{nearby}}), \\
\mu_{\text{miss}|\text{nearby}} &= \mathbb{E}(T_{\text{miss}} \mid T_{\text{nearby}}) \\
&= \text{cov}(T_{\text{miss}}, T_{\text{nearby}}) \text{cov}(T_{\text{nearby}}, T_{\text{nearby}})^{-1} T_{\text{nearby}}, \\
\Sigma_{\text{miss}|\text{nearby}} &= \text{var}(T_{\text{miss}} \mid T_{\text{nearby}}) \\
&= \text{cov}(T_{\text{miss}}, T_{\text{miss}}) - \text{cov}(T_{\text{miss}}, T_{\text{nearby}}) \text{cov}(T_{\text{nearby}}, T_{\text{nearby}})^{-1} \text{cov}(T_{\text{nearby}}, T_{\text{miss}}).
\end{aligned} \tag{11}$$

All covariance matrices can be obtained by plugging into  $k_{st}$ . For example, the  $ij$ th entry of  $\text{cov}(T_{\text{miss}}, T_{\text{nearby}})$  is given by  $k_{st}(x_{\text{miss}}, X_{\text{nearby}}[j], t_{\text{miss}}[i], t_{\text{nearby}}[j])$ , where  $X_{\text{nearby}}[j]$  gives the spatial covariates of the  $j$ th observation, and  $t_{\text{nearby}}[j]$  its time.

In Figure 4, we show an example of predictions obtained from this spatio-temporal model. We withheld measurements from the Waterloo Municipal Airport, and then used data from three nearby stations between May 25, 2015 and June 3, 2015 to predict the Waterloo temperatures during the same time window. This allows us to assess the quality of the predictions on this example.

## 4 Implications

### 4.1 Imputing by Conditioning on Extrema

- introduce general problem
- rejection sampling: doesn't work
- MCMC approach: sample from a distribution
  - spell out model
  - sharp condition: add some tolerance
  - still a difficult distribution to sample from
    - \* 3D diagram of constraint (see Constrained distribution in OneNote)
    - \* MCMC algorithm must navigate corners, and travel down edges efficiently
  - HMC and Stan are particularly good at sampling from complicated distributions (citation?)
- HMC uses the gradient to navigate complicated distributions
  - but the maximum function is ill-behaved (expand)
  - replace it with the smoothmax approximation (plot)
- Demonstration
  - without smoothmax and with smoothmax
  - Stan code in appendix

Our aim isn't simply to predict temperatures at a location with no measurements, but rather to impute hourly temperatures at a location with accurate measurements of the daily temperature extrema. This is an instance of a more general statistical problem: if a random  $p$ -vector  $\{X_i : i = 1, \dots, p\}$  has a known distribution  $F_X$ , and its maximum  $X_{\max} \equiv \max_i \{X_i\}$  and minimum  $X_{\min} \equiv \min_i \{X_i\}$  are measured, how does one draw samples from  $F_{X|X_{\max}, X_{\min}}$ , the distribution of  $X$  conditional on  $X_{\max}$  and  $X_{\min}$ ? Conditional draws from  $F_{X|X_{\max}, X_{\min}}$  need to respect three constraints: one component of  $X$  must be equal to  $X_{\min}$ , another to  $X_{\max}$ , and all other components must lie between  $X_{\min}$  and  $X_{\max}$ .

Conceptually, we could implement a valid imputation algorithm by drawing random samples  $F_X$  (in our application, this is the posterior predictive multivariate normal distribution  $T_{\text{miss}} \mid T_{\text{nearby}}$  obtained from nearby measurements), and only keeping the samples that satisfy the three constraints. Unfortunately, if  $F_X$  is a continuous distribution, the probability of a random draw from  $F_X$  exactly satisfying such sharp constraints is zero. One could envision adding some tolerance, so that samples with minimum and maximum within  $\epsilon$  of  $X_{\max}$  and  $X_{\min}$  are retained, but as  $p$  grows, the rejection probability will rapidly go to 1, thus requiring huge sample sizes. Ultimately, this rejection sampling strategy is therefore bound to fail.

Markov Chain Monte Carlo (MCMC) techniques can also be used to draw samples from arbitrary distributions with densities known up to a constant. The density of  $F_{X|X_{\max}, X_{\min}}$  is obtained up to a constant multiplier through a simple application of Bayes' theorem. It is proportional to the prior density of  $F_X$  multiplied by indicators ensuring that the extrema are respected.

$$\begin{aligned} \mathbb{P}(X | X_{\max}, X_{\min}) &\propto \mathbb{P}(X) \mathbb{P}(X_{\max}, X_{\min} | X), \\ &\propto \mathbb{P}(X) \mathbb{I}\left(\max_i \{X_i\} = X_{\max}\right) \mathbb{I}\left(\min_i \{X_i\} = X_{\min}\right). \end{aligned} \quad (12)$$

However, this distribution is zero everywhere in  $\mathbb{R}^p$ , except in a ( $p-2$ ) dimensional subspace where the min and max constraints are met. This doomed the rejection sampler, and will also prevent any unmodified MCMC algorithm from converging onto  $F_{X|X_{\max}, X_{\min}}$ . We therefore approximate the constraint by replacing the likelihood term  $\mathbb{P}(X_{\max}, X_{\min} | X)$  with two narrow independent normal distributions around the minimum and maximum of  $X$ . This "softens" the conditional distribution,

$$\mathbb{P}(X | X_{\max}, X_{\min}) \propto \mathbb{P}(X) \mathcal{N}\left(X_{\max} | \max_i \{X_i\}, \epsilon^2\right) \mathcal{N}\left(X_{\min} | \min_i \{X_i\}, \epsilon^2\right), \quad (13)$$

where  $\mathcal{N}(x | \mu, \sigma^2)$  is the density of a normal distribution with mean  $\mu$  and variance  $\sigma^2$  evaluated at  $x$ . For small  $\epsilon$ , this is a very tolerable approximation which enables the use of MCMC techniques.

This approximation to  $F_{X|X_{\max}, X_{\min}}$  remains a difficult distribution to sample from. We illustrate the constraint in a 3-dimensional setting in Figure ???. The MCMC algorithm must travel efficiently along the edges of the allowed subspace, and navigate corners when the index of the extremum components change. Hamiltonian Monte Carlo (HMC) has shown a remarkable ability to navigate complicated distributions, including distributions where the typical set has "pinch points" of strong curvature (Betancourt, 2017), similar to the "corners" in  $F_{X|X_{\max}, X_{\min}}$ . We therefore used HMC as implemented by the Stan probabilistic programming language (Carpenter et al., 2017) to obtain draws from  $F_{X|X_{\max}, X_{\min}}$ .

HMC's efficient sampling relies on gradient information in order to move towards regions of high probability. The normal likelihood (13) softened the extrema constraints, but the maximum and minimum functions also remove information from the gradient. The partial derivative of the log-likelihood of the maximum term with respect to  $X_i$  is proportional to

$$\frac{\partial \log \mathcal{N}(X_{\max} | \max_i \{X_i\}, \epsilon^2)}{\partial X_i} \propto (X_{\max} - X_i) \mathbb{I}\left\{\arg \max_j (X_j) = i\right\}, \quad (14)$$

where  $\arg \max$  is the function that returns the index of the maximum component. In other words, the gradient pulls the maximum of the current sample towards  $X_{\max}$ , and ignores all other components. This makes it difficult for HMC to efficiently explore scenarios where other components are the maximum.

We can assist the HMC algorithm with another approximation. We replace the max and min functions with the smoothmax and smoothmin functions, which take real inputs  $x_1, \dots, x_p$  and a sharpness parameter  $k$  and return

$$\begin{aligned} \text{smoothmax}(x_1, \dots, x_p; k) &= \frac{1}{k} \log \left( \sum_{i=1}^p e^{kx_i} \right) \\ \text{smoothmin}(x_1, \dots, x_p; k) &= -\text{smoothmax}(-x_1, \dots, -x_p; k) \end{aligned} \quad (15)$$

As  $k \rightarrow \infty$ , smoothmax becomes the maximum, and smoothmin becomes the minimum. When smoothmax replaces max and smoothmin replaces min, there is a small price in precision due to the approximation, but there is a huge computational benefit: the gradient is now informative for all components of  $X$ :

$$\frac{\partial \log \mathcal{N}(X_{\max} | \text{smoothmax}(X_{1:p}; k), \epsilon^2)}{\partial X_i} \propto (X_{\max} - \text{smoothmax}(X_{1:p}; k)) \frac{e^{kx_i}}{\sum_{j=1}^p e^{kx_j}}. \quad (16)$$

These modifications make HMC a viable algorithm to efficiently draw samples from the constrained posterior. Setting  $k$  and  $\sigma_e$  is a compromise between exactness and efficiency; we found  $k = 10$  and  $\sigma_e = 0.1$  to perform well for this paper's application.

## 4.2 Illustration of Hamiltonian Monte Carlo with Smoothmax Approximation

We verify the ability of the Hamiltonian Monte Carlo algorithm with the smoothmax approximation to obtain draws from  $F_{X|X_{\max}, X_{\min}}$  in a simplified setting where the distribution function of  $F_{X|X_{\max}, X_{\min}}$  can be derived analytically and also computed easily. In our application,  $F_X$  is the posterior predictive multivariate normal distribution  $T_{\text{miss}} | T_{\text{nearby}}$  obtained from nearby measurements, with mean and marginal variance evolving smoothly from one prediction to the next. To retain a resemblance to this, we specify a random vector  $X$  with each component  $X_i$  normally distributed, and with sinusoidal means and variances, but without any correlations between them, so as to avoid a combinatorial explosion when obtaining the distribution function analytically:

$$\begin{aligned} X_i &\sim \mathcal{N}(\mu_i, \sigma_i), \\ X_i &\perp X_j \forall i \neq j, \\ X_{\max} &= \max_i \{X_i\}, \\ X_{\min} &= \min_i \{X_i\}, \\ \mu_i &= 10 + \sin(2\pi i/50), \\ \sigma_i &= 0.1 + \cos^2(2\pi i/50), \\ i &= 1, 2, \dots, 100. \end{aligned} \tag{17}$$

The unconstrained distribution of  $X_i$  is illustrated in Figure 5(a). In this example, we aim to sample from the distribution of  $X_i$  subject to the observation that  $X_{\max} = 12.5$  and  $X_{\min} = 8.8$ .

We first derive and compute  $F_{X|X_{\max}, X_{\min}}$  for this example. We denote by  $f_i(\cdot)$  and  $F_i(\cdot)$  the prior probability distribution function and cumulative distribution function of  $X_i$ , i.e. the normal PDF and CDF with means and variances given by (17). Let  $\mathbb{P}_{ij}$  be the probability that  $X_i$  is the minimum of  $X$ , and  $X_j$  is its maximum. We also define  $\mathbb{P}_{i\bullet} = \sum_{j=1}^{100} \mathbb{P}_{ij}$ , the probability that  $X_i$  is the minimum, and  $\mathbb{P}_{\bullet j} = \sum_{i=1}^{100} \mathbb{P}_{ij}$ , the probability that  $X_j$  is the maximum. The cumulative distribution function of  $X_i$  is then given by

$$\mathbb{P}(X_i \leq x | X_{\max}, X_{\min}) = \begin{cases} 0 & \text{when } x < X_{\min}, \\ \mathbb{P}_{i\bullet} + (1 - \mathbb{P}_{i\bullet} - \mathbb{P}_{\bullet i}) \left[ \frac{F_i(x) - F_i(X_{\min})}{F_i(X_{\max}) - F_i(X_{\min})} \right] & \text{when } X_{\min} \leq x < X_{\max}, \\ 1 & \text{when } x \geq X_{\max}. \end{cases} \tag{18}$$

Meanwhile,  $\mathbb{P}_{ij}$  is proportional to

$$f_i(X_{\min}) f_j(X_{\max}) \prod_{k \neq i, j}^{100} (F_k(X_{\max}) - F_k(X_{\min})), \tag{19}$$

which we compute for all  $i, j$  and renormalize to obtain the  $100 \times 100$  matrix of probabilities. We sum over its rows and columns to obtain  $\mathbb{P}_{\bullet j}$  and  $\mathbb{P}_{i\bullet}$ . While this algorithm has cubic complexity in the dimensionality  $p$  of  $X$ , for  $p = 100$  computers only take seconds to compute the entries of  $\mathbb{P}$  and evaluate  $\mathbb{P}(X_i \leq x | X_{\max}, X_{\min})$  over a range of  $x$ . Figure 5(b) shows the analytical quantiles of  $F_{X|X_{\max}, X_{\min}}$ . Roughly speaking, we see that the prior distribution  $F_X$  is stretched to fit between  $X_{\min}$  and  $X_{\max}$ .

To obtain samples from  $F_{X|X_{\max}, X_{\min}}$ , we use the implementation of Hamiltonian Monte Carlo provided by the probabilistic programming language Stan. In Stan, the user specifies a probabilistic data-generating process for the observed data, based on parameters and latent variables with accompanying priors. Stan then compiles this model into a custom C++ program that implements posterior sampling using HMC. We implement two Stan models to draw from  $F_{X|X_{\max}, X_{\min}}$ . The Stan model code for both are available in

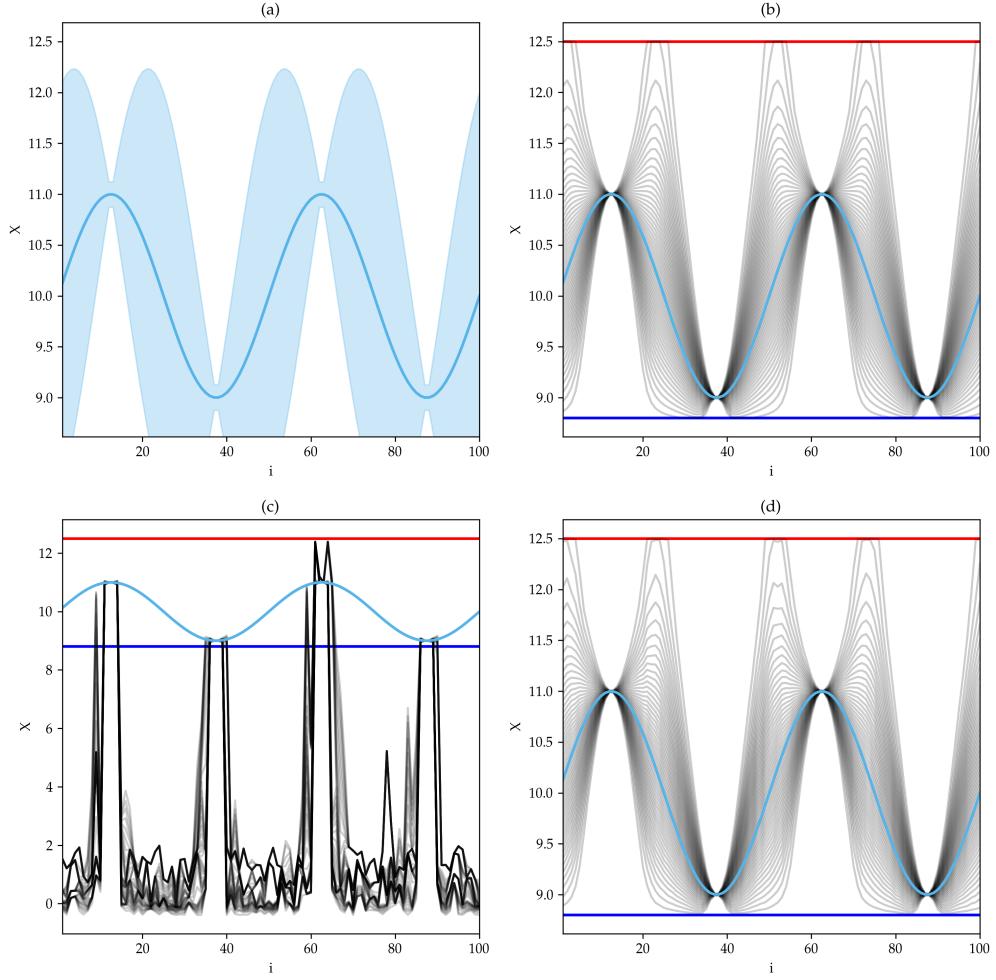


Figure 5: (a) Prior distribution of  $X_i$  displayed as mean function with  $2\sigma$  envelope; (b) Quantiles of the analytically derived posterior  $F_{X|X_{\max}, X_{\min}}$  conditioned on  $X_{\min}$  and  $X_{\max}$ , with prior  $\mu_i$  shown in blue; (c) Quantiles of the samples drawn from  $F_{X|X_{\max}, X_{\min}}$  using Stan without the smoothmax approximation, with prior  $\mu_i$  shown in blue; (d) Quantiles of the samples drawn from  $F_{X|X_{\max}, X_{\min}}$  using Stan with the smoothmax approximation, with prior  $\mu_i$  shown in blue.

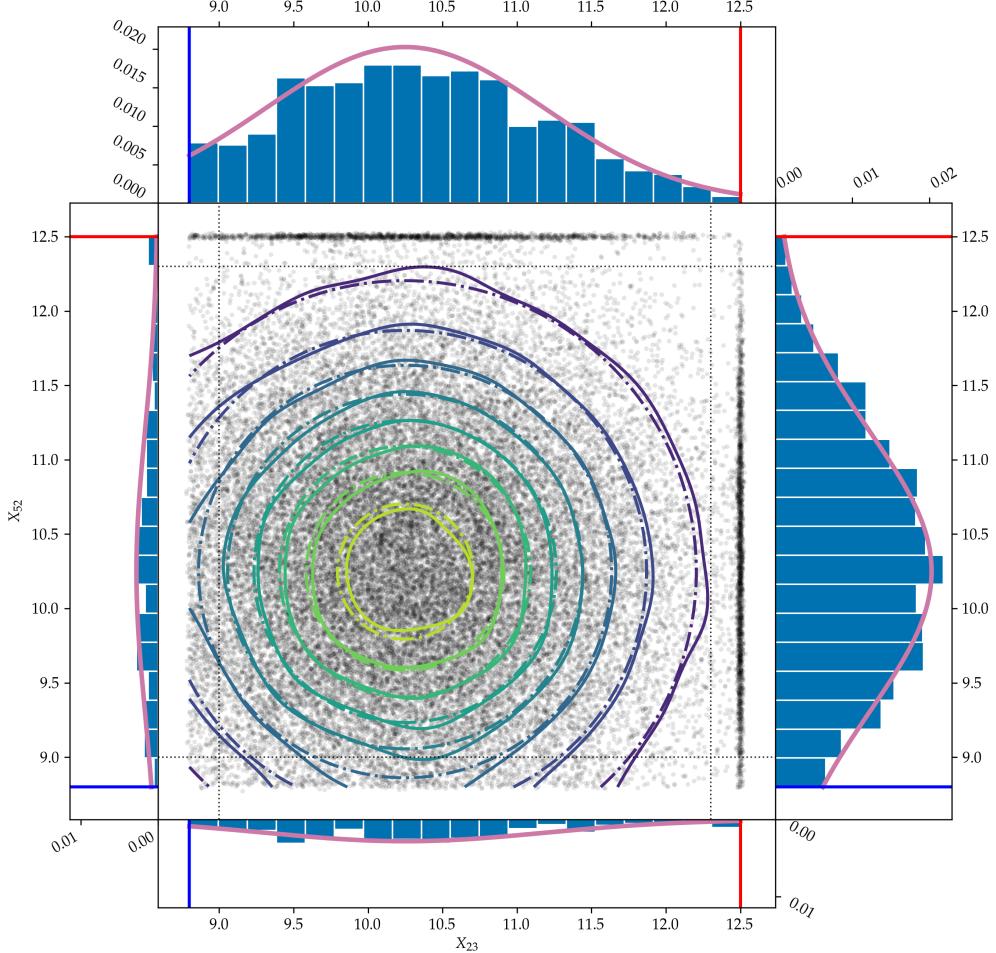


Figure 6: Comparison of the joint distribution of  $X_{23}$  and  $X_{52}$  obtained analytically and from Stan samples.

appendix 9.1. The first model implements equation (13), with the narrow normal likelihood term around the maximum and minimum, while the second model includes the smoothmax and smoothmin approximations to the maximum and minimum functions. For each Stan model, we obtain 4 HMC chains each with 10,000 warm-up samples followed by 10,000 samples. The quantiles of the samples obtained without the smoothmax approximation are shown in Figure 5(c). By default, Stan initializes each  $X_i$  uniformly at random between -2 and 2, and for most variables, the algorithm remains stuck near the initial values. Furthermore, most samples do not conform to the constraints imposed by the observed  $X_{\min}$  and  $X_{\max}$  values, which further invalidates these results. However, once we replace the maximum function with the smoothmax function, with quantiles shown in Figure 5(d), Stan is able to draw samples that respect the observed extrema. Furthermore, a visual comparison of the analytical quantiles in Figure 5(a) and the Stan sample quantiles in Figure 5(d) confirms that this sampling algorithm delivers a close approximation of the marginal distribution of each variable  $X_i$  in  $F_{X|X_{\max},X_{\min}}$ .

We may also wish to verify that Stan samples correctly from the *joint* distribution of any combination of variables. We do this visually for a pair of variables,  $X_{23}$  and  $X_{52}$ , with results shown in Figure 6. In that figure, the central scatterplot shows the 40,000 Stan samples obtained using the smoothmax approximation. Superimposed on the scatterplot are a contour plot (with dash-dotted lines) of the probability distribution function of the analytical conditional distribution  $F_{X|X_{\max},X_{\min}}$  of  $X_{23}$  and  $X_{52}$  when neither  $X_{23}$  nor  $X_{52}$  is one of the extrema, multiplied by the probability of that being the case (which is  $1 - (\mathbb{P}_{23\bullet} + \mathbb{P}_{52\bullet} + \mathbb{P}_{\bullet23} + \mathbb{P}_{\bullet52}) + (\mathbb{P}_{23,52} + \mathbb{P}_{23,52})$ ). This can be compared to the contour plot (solid lines) of the same probability distribution function obtained through a kernel density estimator of the subset of Stan

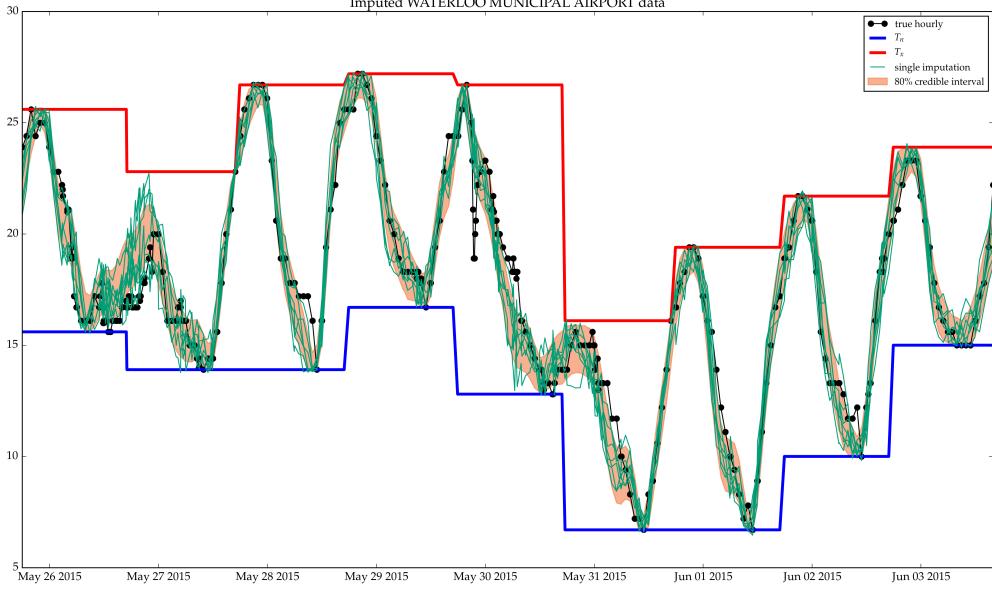


Figure 7: Imputations at Waterloo Airport from May 25, 2015 to June 3, 2015

samples where neither  $X_{23}$  and  $X_{52}$  is the minimum or maximum, using a normal kernel with bandwidth 0.2, and multiplied by the proportion of samples where that is the case. The kernel estimates are divided by the integrated probability mass of the kernel that is inside of the boundaries, in order to reduce boundary effects. The thin black dotted line are one kernel bandwidth away from the  $X_{\min}/X_{\max}$  boundaries. Outside the thin black dotted line, the kernel density estimates are less trustworthy. The four histograms around the scatter plot are of the Stan samples where one of the variables is an extremum, weighted so as to integrate to the fraction of samples that satisfy that condition. For example, the top histogram is of  $X_{23}$  for samples where  $X_{52}$  is the maximum, and integrates to the fraction of samples where that is the case. The super-imposed pink line is that of a truncated normal probability distribution function multiplied by the probability of the satisfied condition. For example, the pink line over the top histogram integrates to  $\mathbb{P}_{\bullet 52}$ . Lastly, the blue and red lines are  $X_{\min}$  and  $X_{\max}$  respectively. There is a close match between the contours of the analytical joint distribution function (dash-dotted lines) and of the kernel density estimate of the Stan samples. Each of the four histogram of samples where  $X_{23}$  or  $X_{52}$  occupies the minimum or maximum position matches the corresponding analytical distribution function. This visual comparison of the sample and analytical distributions should reassure us that Stan is yielding a good approximation of a sample drawn from the true  $F_{X|X_{\max}, X_{\min}}$  in this example. We did not examine the behavior of the sampling algorithm for the joint distribution of more than two variables due to the difficulty of visualizing such a distribution, but we see no reason to suspect that the algorithm suffers from pathological behavior that does not manifest itself in these univariate and bivariate inspections.

### 4.3 Smoothmax Temperature Model

Armed with the smoothmax approximation implemented in Stan, we finally return to the problem of imputing hourly temperature measurements. A small leap of faith is needed to accept that the success of the strategy that we implemented and tested in a toy example in the previous two sections will extend to this application. There are three important differences between the toy example and the temperature model. Firstly,  $F_{X|X_{\max}, X_{\min}}$  is now a multivariate normal distribution with strong correlations obtained as the posterior of a Gaussian process in section. Secondly, instead of a single minimum and maximum, we observe extrema for every 24 hour period. Thirdly, we allow for the mean temperature to be different at different locations, and so the imputed temperatures are shifted by an additional parameter  $\mu_{\text{miss}}$ , to which we will attach a vague prior. To summarize, the probabilistic model that we wish to draw posterior imputations of

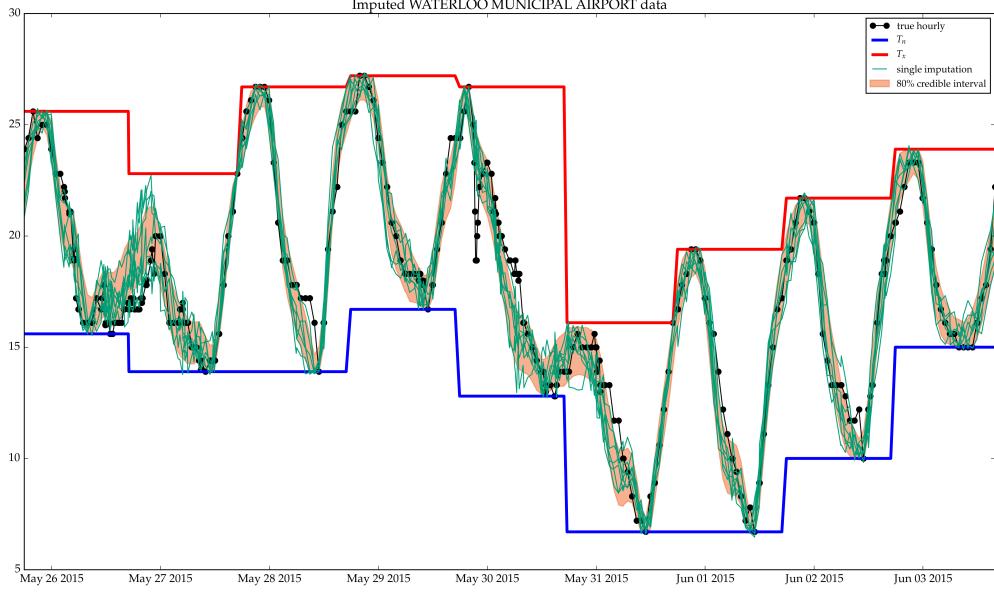


Figure 8: Imputations obtained using a product of squared exponential kernels.

$T_{\text{miss}}$  from is given in (20).

$$\begin{aligned}
 \mu_{\text{miss}} &\sim \mathcal{N}(0, 100) && \text{(vague prior on mean temperature)} \\
 f_{\text{miss}} &\sim \mathcal{N}(\mu_{\text{miss|nearby}}, \Sigma_{\text{miss|nearby}}) && \text{(posterior from } T_{\text{nearby}} \text{ becomes prior)} \\
 T_{\text{miss}} &= \mu_{\text{miss}} + f_{\text{miss}} \\
 T_x[\text{day}] &= \max_{i \in \{i\}_{\text{day}}} \{T_{\text{miss},i}\} && \text{(observe maximum in 24hr window)} \quad (20) \\
 T_n[\text{day}] &= \min_{i \in \{i\}_{\text{day}}} \{T_{\text{miss},i}\} && \text{(observe minimum in 24hr window)} \\
 \{i\}_{\text{day}} &= \left\{ i : \text{day} - 1 + \frac{\text{hour}}{24} < t_{\text{miss},i} \leq \text{day} + \frac{\text{hour}}{24} \right\} && \text{(indices of times in the 24hr window)}
 \end{aligned}$$

To sample from this model, we modify it with the smoothmax approximation to the maximum, and a normal likelihood. The resulting model is shown in (21), and the corresponding Stan code is in Appendix 10.

$$\begin{aligned}
 \mu_{\text{miss}} &\sim \mathcal{N}(0, 100) \\
 f_{\text{miss}} &\sim \mathcal{N}(\mu_{\text{miss|nearby}}, \Sigma_{\text{miss|nearby}}) \\
 T_{\text{miss}} &= \mu_{\text{miss}} + f_{\text{miss}} \\
 T_x[\text{day}] &\sim \mathcal{N} \left( \text{smoothmax}_{i \in \{i\}_{\text{day}}} \{T_{\text{miss},i}; k = 10\}, 0.1^2 \right) \\
 T_n[\text{day}] &\sim \mathcal{N} \left( \text{smoothmin}_{i \in \{i\}_{\text{day}}} \{T_{\text{miss},i}; k = 10\}, 0.1^2 \right)
 \end{aligned} \quad (21)$$

Example imputations from this procedure are shown in Figure 8. From May 25, 2015 to June 3, 2015, hourly temperatures are imputed at Waterloo Airport, using the hourly temperature measurements from nearby stations to inform the course of the temperatures, and using the daily minima and maxima “measurements” to constrain the imputed temperatures, and to infer the mean. Because we actually have hourly

data for Waterloo, yet only fed our algorithm a reduction of this data to daily extremes, we can also plot the hidden temperatures, and see how faithfully the imputations reproduce them. We see that the imputations indeed track the true measurements very closely. The error bars satisfactorily narrow and widen in accordance to the amount of information available at each moments. On May 27th, we can see that the imputations capture the fact that the  $T_x$  record *could* have been set early in the measurement window, but more likely at its very end.

## 5 Model diagnostics

### 5.1 Variogram

We can visually inspect our model by plotting temporal and spatial semi-variograms. The semi-variogram of a stationary spatio-temporal function  $Y(x, t)$  is a function of the spatial lag  $h$  and the temporal lag  $r$

$$\gamma(h, r) = \frac{1}{2} \mathbb{E} [(Y(x, t) - Y(x + h, t + r))^2] = \text{var}(Y(x, t)) - \text{cov}((Y(x, t)), Y(x + h, t + r)). \quad (22)$$

For a Gaussian Process model, with a stationary kernel  $k(h, r) = k(x, x + h, t, t + r)$  this can be expressed in terms of the observation noise  $\sigma_e$  and kernel function  $k(\cdot, \cdot)$ , as

$$\gamma(h, r) = \sigma_e^2 + k(0, 0) - k(h, r). \quad (23)$$

From the data, the semi-variogram can also be estimated empirically, by averaging the square differences of any two observations that are separated by  $h$  in space, and  $r$  in time (or, in practice, within half a bin width of  $h$  and  $r$ ). By comparing the empirical variogram to the variogram of our fitted GP model, we obtain a visual diagnosis of the model.

In our Iowa example, there are only four possible locations. For each location, we plot the empirical temporal variogram  $\hat{\gamma}(0, r)$ . For any pair of stations separated by  $h$  (fixed), we can also plot  $\hat{\gamma}(h, r)$ . We then overlay the model's semi-variogram obtained through equation (23), resulting in Figure 9.

We notice that the variogram of the simple SExSE model tracks the empirical variogram well at short lags, but fails to capture the diurnal cycle, and the fit degrades at long lag. We attempt to improve the model in section 6.

### 5.2 Error and expected error

The variogram gives us a visual diagnostic of the overall model fit. To quantify the model's predictive ability in the Iowa example, we compare the posterior mean temperature to the withheld truth, and obtain the empirical mean squared error as

$$\text{MSE}(\text{err} | T_{\text{nearby}}, T_x, T_n) = \frac{1}{N} \sum_{i=1}^N [\mathbb{E}(T_{\text{miss},i} | T_{\text{nearby}}, T_x, T_n) - T_{\text{miss},i}]^2. \quad (24)$$

This equation is for the final predictions obtained using nearby hourly temperatures and local daily maxima and minima. A similar diagnostic can be computed for the intermediary predictions, which exclude the local  $T_x$  and  $T_n$  information. At that stage, we are not concerned with any overall bias in the predicted temperatures, so we instead compute the sample variance of the errors as

$$\text{var}(\text{err} | T_{\text{nearby}}) = \text{var}_i \{ \mathbb{E}(T_{\text{miss},i} | T_{\text{nearby}}) - T_{\text{miss},i} \}. \quad (25)$$

Model	Log Likelihood	Var(err)	$\mathbb{E}(\text{Var}(\text{err}))$	MSE(err)	$\mathbb{E}(\text{MSE}(\text{err}))$
SE x SE	-55,614	1.589	0.875	1.104	0.614
SExSE + diurnal	-54,472	1.633	0.974	1.137	0.697
Sum of products, fixed variance	-48,589	4.991	8.791		

Model	Log Likelihood	Var(err)	$\mathbb{E}(\text{Var}(err))$	MSE(err)	$\mathbb{E}(\text{MSE}(err))$
SoP, fixed temporal, free var	-47,082	1.314	2.321	1.150	0.897
SoP, completely free	-46,184	1.423	1.765	1.152	0.950
SoP, simpler	-45,945	1.319	1.190	1.069	0.823

For our purposes, it isn't sufficient for the spatio-temporal model to yield good predictions; we also require a good estimate of its own accuracy. We estimate the expected MSE and predictive variance by sampling  $K$  random draws  $T_{\text{miss}}^k$  from the posterior distribution, again conditioned firstly on just  $T_{\text{nearby}}$  after fitting the spatio-temporal Gaussian process model, and then additionally on  $T_{\text{nearby}}, T_x$  and  $T_n$  after incorporating the local data using Stan. The draws are obtained from the posterior multivariate normal distribution in the first case, and the MCMC samples obtained through Stan in the second case. We then evaluate the variance or MSE between the samples and the posterior mean as

$$\begin{aligned}\mathbb{E}(\text{var}(\text{err} | T_{\text{nearby}})) &\approx \frac{1}{K} \sum_{k=1}^K \text{var}_i \left\{ T_{\text{miss},i}^{(k)} - \mathbb{E}(T_{\text{miss},i} | T_{\text{nearby}}) \right\} \\ \mathbb{E}(\text{MSE}(\text{err} | T_{\text{nearby}}, T_x, T_n)) &\approx \frac{1}{K} \sum_{k=1}^K \text{MSE}_i \left\{ T_{\text{miss},i}^{(k)} - \mathbb{E}(T_{\text{miss},i} | T_{\text{nearby}}, T_x, T_n) \right\}\end{aligned}\quad (26)$$

When evaluating models, we want the errors to be small, and so the empirical error variance and MSE to be low. A well-calibrated model should also have the expected error variances  $\mathbb{E}(\text{var}(\text{err} | \cdot))$  close to their empirical values.

These diagnostics for our first spatio-temporal model, the product of squared exponentials, are found in the first row of Table X. The empirical error variance using only nearby measurements is already fairly low, with typical errors of order  $\sqrt{1.589} = 1.26^\circ \text{C}$ . Incorporating the local measurements reduces it further to  $\sqrt{1.104} = 1.05^\circ \text{C}$ . However, the model is overly optimistic, and the expected errors underestimate the true errors.

## 6 Improving model

In this section, we develop more sophisticated Gaussian process models than the simple product of squared exponential kernels. We then assess whether these models improve the variogram and the predictive diagnostic measures that we developed in the previous sections.

The most salient feature of the empirical variogram that isn't captured by the SExSE model is the oscillation with a 24-hour period. It is intuitively obvious that the diurnal cycle induces this periodic covariance, and that our model should be improved by incorporating this feature. Gaussian process models allow for periodic components of the covariance, for example the periodic squared exponential kernel, which we will use with a 24-hour period

$$k_{24}(t, t') = \sigma_{24}^2 \exp \left[ -\frac{2}{\ell_{24}^2} \sin^2 \left( \pi \frac{t - t'}{24 \text{ hrs}} \right) \right]. \quad (27)$$

We modify the spatiotemporal model by adding this diurnal component to it, with its own spatial decay kernel  $k_{\text{space}24}$  (with the same specification as  $k_{\text{space}}$  in (1)).

$$k_{\text{SESE}_24}(\mathbf{x}, \mathbf{x}', t, t') = k_{\text{time}}(t, t') \cdot k_{\text{space}}(\mathbf{x}, \mathbf{x}') + k_{24}(t, t') \cdot k_{\text{space}24}(\mathbf{x}, \mathbf{x}') + k_\mu(\mathbf{x}, \mathbf{x}'). \quad (28)$$

We also develop a more complex model, which breaks up  $k_{\text{time}}$  into short-term, medium-term and long-term correlation components, each with their own spatial decay.

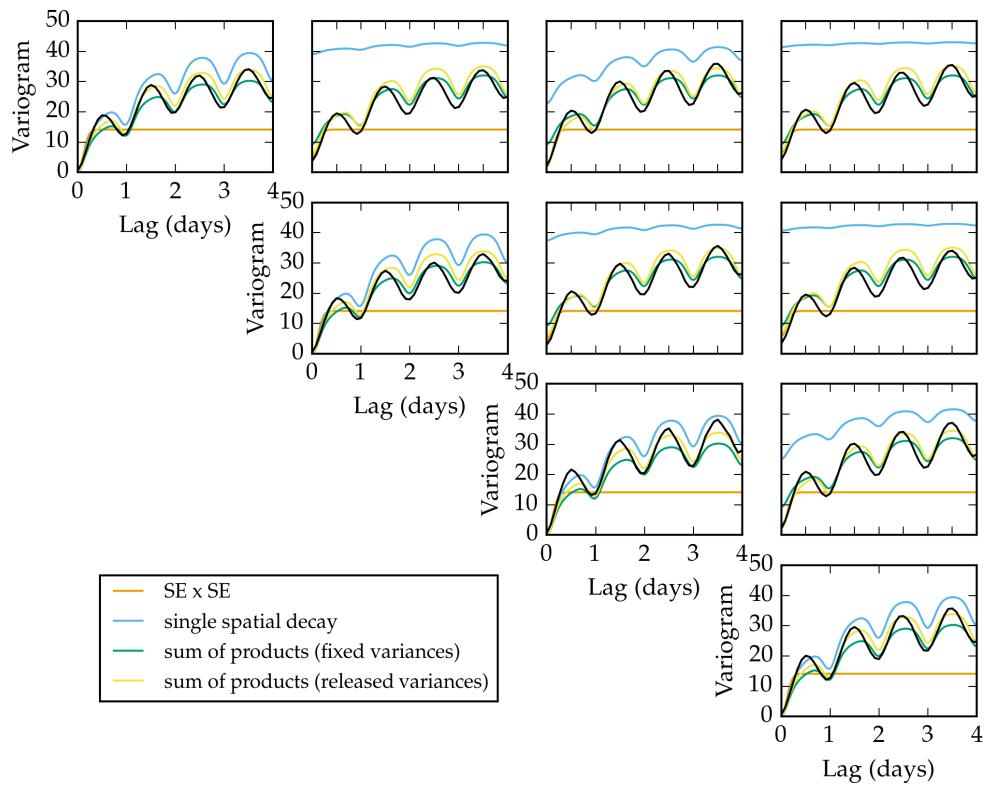


Figure 9: Semi-variogram

$$\begin{aligned}
k_{\text{sumprod}}(\mathbf{x}, \mathbf{x}', t, t') &= k_{\text{time}1}(t, t') \cdot k_{\text{space}1}(\mathbf{x}, \mathbf{x}') && (\text{short-term variation}) \\
&+ k_{\text{time}2}(t, t') \cdot k_{\text{space}2}(\mathbf{x}, \mathbf{x}') && (\text{medium-term variation}) \\
&+ k_{\text{time}3}(t, t') \cdot k_{\text{space}3}(\mathbf{x}, \mathbf{x}') && (\text{long-term variation}) \\
&+ k_{24}(t, t') \cdot k_{\text{space}24}(\mathbf{x}, \mathbf{x}') && (\text{diurnal cycle}) \\
&+ k_{\mu}(\mathbf{x}, \mathbf{x}') && (\text{station mean})
\end{aligned} \tag{29}$$

Each of  $k_{\text{time}1}$ ,  $k_{\text{time}2}$ , and  $k_{\text{time}3}$ , is a rational quadratic kernel

$$k_{\text{RQ}}(t, t') = \sigma^2 \left( 1 + \frac{(t - t')^2}{2\alpha\ell^2} \right)^{-\alpha} \tag{30}$$

which is accompanied by its spatial decay kernel, specified as a squared exponential covariance. This more complicated kernel therefore has  $3 \times 3 \times 2 + 2 \times 2 = 22$  free parameters, in addition to the noise parameter  $\sigma_e^2$ .

We now have three competing Gaussian process models, with covariance functions  $k_{\text{SEXSE}}$ ,  $k_{\text{SESE}_2 4}$ , and  $k_{\text{sumprod}}$  respectively. We can compare them in three ways. Firstly, the marginal log-likelihood is the quantity maximized by the parameter fitting procedure in (10). The maximized log-likelihood can be found in the second column of Table XX, and we see that the more complex models indeed yield a much higher log-likelihood, promising a better model fit which should yield better predictions. Secondly, we compare the variance of the error in the predicted temperatures specified in (25) when withholding all the data from a test station. Averaged over all of 2015, this is given in the third column, and shows more mixed results. The diurnal model  $k_{\text{SESE}_2 4}$  performs worse than the simple  $k_{\text{SEXSE}}$  model, and  $k_{\text{sumprod}}$  only yields a small improvement. Thirdly, we can reintroduce the daily minima and maxima from the withheld station, and compare the mean squared error specified in (24) for predictions at the test station. Results in the fifth column show even more modest improvements for the more complex models.

We interpret these results as a reminder that predictions using Gaussian process are sensitive to model specification when extrapolating, but fairly insensitive to the model when interpolating [cite?]. Since our imputations interpolate the temperatures from nearby stations, further aided by the constraints imposed by the daily  $T_n$  and  $T_x$  measurements, the choice of model does not have a large impact on the performance of our procedure. This insensitivity can be seen as reassuring, as it (to an extent) reduces our need to worry about the incorrectness of our model.

## 7 Analysis

- show imputations on interesting days
- show imputations can capture two possible explanations for a measurement
- discuss possibility of inferring measurement time

## 8 Inference on measurement hour

Our analysis so far has focused on the case where the hour of measurement `hour` is known in advance. This is an unrealistic assumption in practice, and so inference on `hour` is a desirable feature. It is conceptually straightforward to modify (21) with a uniform prior on `hour`. However, because we obtain our imputations in ten-day windows, in most windows precise information about `hour` will not be available, as moving the measurement time one hour earlier or later rarely affects the measured  $T_n$  and  $T_x$ . Furthermore, `hour` affects which observations are attributed to each day's measurements. This effect is discontinuous (observations suddenly jump from one day to the next) and non-differentiable, and so Hamiltonian Monte Carlo becomes unviable. This issue is similar to that caused by the non-differentiability of the minimum and maximum functions. We therefore do not consider the introduction of a uniform prior on `hour` in Stan to be feasible.

Our procedure allows us to obtain imputation samples of  $T_{\text{miss}}$  conditional on  $T_{\text{nearby}}$ ,  $T_n$ ,  $T_x$  and `hour`. If we do so for  $\text{hour} = 1, 2, \dots, 24$ , is there information available in these samples to infer `hour`? We will

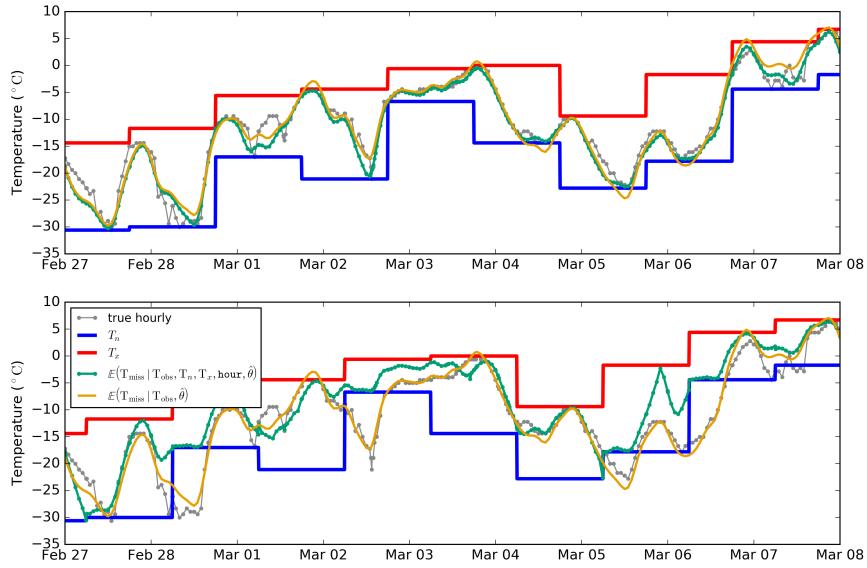


Figure 10: A sample window showing constrained and unconstrained imputations assuming (top) the correct measurement hour (17 UTC), and (bottom) the wrong measurement hour (5 UTC). Assuming the wrong measurement time drives the constrained mean imputation away from the unconstrained mean imputation.

examine sample imputations to answer this question. Figure 10 shows mean imputation for temperatures over nine days starting on February 27, 2015. The orange line is the mean using only nearby temperatures (shifted by a constant to match the true temperatures), while the green line is additionally conditional on  $T_n$  and  $T_x$ ; the true temperatures are shown in grey. The top plot shows the imputation under the correct daily measurement time (17 UTC), while the bottom plot is under an incorrect measurement time (5 UTC). The first unsurprising observation is that assuming an incorrect measurement time can lead to wildly inaccurate imputations. But we then also notice that assuming the wrong time also causes the mean constrained imputation to depart further from the unconstrained imputation (that is, the green and orange lines are further apart). This can be interpreted as an indication of an incompatibility between  $T_{\text{nearby}}$  and the daily extremes, caused by assuming the wrong hour. To quantify this discrepancy, we propose to calculate the probability of the mean constrained imputation under the unconstrained posterior given by (11):

$$\begin{aligned} \mu(\text{hour}) &\equiv \mathbb{E}(T_{\text{miss}} | T_{\text{nearby}}, T_n, T_x, \text{hour}) \quad (\text{the mean imputed temperature}), \\ \delta_{\text{hour}} &\equiv \mathbb{P}(T_{\text{miss}} = \mu(\text{hour}) | T_{\text{nearby}}). \end{aligned} \tag{31}$$

Our intuition is that  $\delta_{\text{hour}}$  will drop sharply when the wrong hour is assumed, and we may be able to infer the true hour by maximizing  $\delta_{\text{hour}}$ .

Fortunately, our discrepancy measure  $\delta_{\text{hour}}$  also admits a Bayesian interpretation: it is proportional to the marginal likelihood of hour under (admittedly fanciful) approximating assumptions. Ideally, we would evaluate the marginal likelihood  $\mathbb{P}(T_n, T_x | T_{\text{nearby}}, \text{hour})$ , and then appeal to Bayes theorem to obtain a posterior on hour

$$\mathbb{P}(\text{hour} | T_n, T_x, T_{\text{nearby}}) \propto \mathbb{P}(T_n, T_x | T_{\text{nearby}}, \text{hour}) \mathbb{P}(\text{hour}). \tag{32}$$

However, marginal likelihoods are notoriously difficult to estimate from posterior samples [cite? Raftery

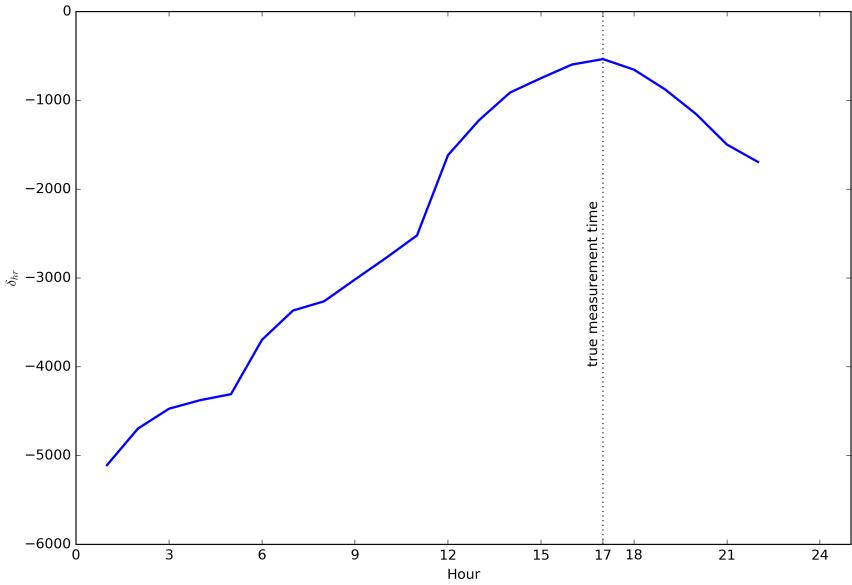


Figure 11: Discrepancy measure for imputations of temperatures at Waterloo Municipal Airport assuming measurement hours  $\text{hour} = 1, 2, \dots, 24$ . The true hour of measurement is 17, and obtains the highest  $\delta_{\text{hour}}$ .

1994?]. The marginal likelihood is the normalizing constant for the posterior (??) of  $T_{\text{miss}}$ , and therefore for any  $T_{\text{miss}}$

$$\mathbb{P}(T_n, T_x | T_{\text{nearby}}, \text{hour}) = \frac{\mathbb{P}(T_n, T_x | T_{\text{miss}}, \overline{T_{\text{nearby}}, \text{hour}}) \mathbb{P}(T_{\text{miss}} | T_{\text{nearby}}, \text{hour})}{\mathbb{P}(T_{\text{miss}} | T_n, T_x, T_{\text{nearby}}, \text{hour})}. \quad (33)$$

The first term in the numerator is either one or zero, as discussed near equation (??). If we assume the constraint is satisfied, pick  $T_{\text{miss}} = \mu(\text{hour})$ , and assume that the posterior density evaluated at its mean does not depend heavily on the time of measurement, we obtain that the marginal likelihood is proportional to  $\delta_{\text{hour}}$ . Both assumptions are fanciful: the posterior mean generally will violate the constraint imposed by  $T_n$  and  $T_x$ , and therefore the likelihood  $\mathbb{P}(T_n, T_x | T_{\text{miss}})$  should in fact be zero. Furthermore, there is no reason to think the posterior density at the posterior mean does not depend on  $\text{hour}$ , but we might reasonably hope that the wrongness of this assumption does not overwhelm the signal contained in  $\delta_{\text{hour}}$ . This reasoning at least confirms that  $\delta_{\text{hour}}$  captures information about the likelihood of  $\text{hour}$ , and that once renormalized it can be loosely interpreted as a posterior probability under a uniform prior.

## 9 Appendices

### 9.1 Stan programs for illustration of smoothmax

#### 9.1.1 Without **smoothmax** Approximation

```
data {
    int<lower=0> N; // number of observations
    real Xmax;
    real Xmin;
    vector[N] mu_i;
    real<lower=0> sigma_i[N];
```

```

    }
parameters {
    vector[N] X_i; // latent variables
}
model {
    X_i ~ normal(mu_i, sigma_i);
    Xmax ~ normal(max(X_i), 0.01);
    Xmin ~ normal(min(X_i), 0.01);
}

```

### 9.1.2 With **smoothmax** Approximation

```

functions {
    real smoothmax(vector x, real k, real maxkx) {
        return (maxkx+log(sum(exp(k*x - maxkx))))/k;
    }
    real smoothmin(vector x, real k, real minkx) {
        return -smoothmax(-x, k, -minkx);
    }
}
data {
    int<lower=0> N; // number of observations
    real Xmax;
    real Xmin;
    real mu_i[N];
    real<lower=0> sigma_i[N];
    real<lower=0> k;
}
parameters {
    vector[N] X_i; // latent variables
}
transformed parameters {
    real Xsmoothmax;
    real Xsmoothmin;
    Xsmoothmax = smoothmax(X_i, k, k*Xmax);
    Xsmoothmin = smoothmin(X_i, k, k*Xmin);
}
model {
    X_i ~ normal(mu_i, sigma_i);
    Xmax ~ normal(Xsmoothmax, 0.01);
    Xmin ~ normal(Xsmoothmin, 0.01);
}

```

## 10 Stan model for temperature imputations

```

functions {
    real smoothmax(vector x, real k, real maxkx) {
        return (maxkx+log(sum(exp(k*x - maxkx))))/k;
    }
    real smoothmin(vector x, real k, real minkx) {
        return -smoothmax(-x, k, -minkx);
    }
}

```

```

data {
    // Tn Tx data
    int<lower=1> N_TxTn; //
    vector[N_TxTn] Tx;
    vector[N_TxTn] Tn;

    // imputation points (for which we have )
    int<lower=1> Nimpt;
    int<lower=1,upper=N_TxTn> day_impute[Nimpt];
    // number of hours recorded within each day
    int<lower=1> impt_times_p_day[N_TxTn];

    // prior
    vector[Nimpt] predicted_mean;
    matrix[Nimpt,Nimpt] predicted_cov;
    matrix[Nimpt,Nimpt] predicted_cov_chol;

    // control soft max hardness
    real<lower=0> k_smoothmax;
}
parameters {
    vector[Nimpt] w_uncorr;
    real mu;
}
transformed parameters {
    vector[Nimpt] temp_impt;
    real Tsmoothmax[N_TxTn];
    real Tsmoothmin[N_TxTn];
    temp_impt = mu + predicted_mean + predicted_cov_chol*w_uncorr;
{
    int istart;
    istart = 1;
    for (i in 1:N_TxTn) {
        int ntimes;
        ntimes = impt_times_p_day[i];
        Tsoftmin[i] = smoothmin(segment(temp_impt,istart,ntimes),
                                k_smoothmax,
                                k_smoothmax*Tn[i]);
        Tsoftmax[i] = smoothmax(segment(temp_impt,istart,ntimes),
                                k_smoothmax,
                                k_smoothmax*Tx[i]);
        istart = istart + ntimes;
    }
}
model {
    w_uncorr ~ normal(0,1);
    mu ~ normal(0, 100.0);
    Tn ~ normal(Tsmoothmin, 0.1);
    Tx ~ normal(Tsmoothmax, 0.1);
}

```

## References

- Baker, D. G., 1975: Effect of observation time on mean temperature estimation. *Journal of Applied Meteorology*, **14** (4), 471–476.
- Betancourt, M., 2017: A conceptual introduction to hamiltonian monte carlo. *arXiv preprint arXiv:1701.02434*.
- Carpenter, B., and Coauthors, 2017: Stan: A probabilistic programming language. *Journal of Statistical Software, Articles*, **76** (1), 1–32, doi:10.18637/jss.v076.i01, URL <https://www.jstatsoft.org/v076/i01>.
- Della-Marta, P., and H. Wanner, 2006: A method of homogenizing the extremes and mean of daily temperature measurements. *Journal of Climate*, **19** (17), 4179–4197.
- Ducré-Robitaille, J.-F., L. A. Vincent, and G. Boulet, 2003: Comparison of techniques for detection of discontinuities in temperature series. *International Journal of Climatology*, **23** (9), 1087–1101.
- Easterling, D. R., T. C. Peterson, and T. R. Karl, 1996: On the development and use of homogenized climate datasets. *Journal of climate*, **9** (6), 1429–1434.
- Karl, T. R., C. N. Williams Jr, P. J. Young, and W. M. Wendland, 1986: A model to estimate the time of observation bias associated with monthly mean maximum, minimum and mean temperatures for the United States. *Journal of Climate and Applied Meteorology*, **25** (2), 145–160.
- Menne, M. J., I. Durre, R. S. Vose, B. E. Gleason, and T. G. Houston, 2012: An overview of the Global Historical Climatology Network-Daily database. *Journal of Atmospheric and Oceanic Technology*, **29** (7), 897–910.
- Menne, M. J., and C. N. Williams Jr, 2009: Homogenization of temperature series via pairwise comparisons. *Journal of Climate*, **22** (7), 1700–1717.
- Menne, M. J., C. N. Williams Jr, and R. S. Vose, 2009: The US Historical Climatology Network monthly temperature data, version 2. *Bulletin of the American Meteorological Society*, **90** (7), 993–1007.
- Peterson, T. C., and Coauthors, 1998: Homogeneity adjustments of in situ atmospheric climate data: a review. *International Journal of Climatology*, **18** (13), 1493–1517.
- Trewin, B., 2013: A daily homogenized temperature data set for Australia. *International Journal of Climatology*, **33** (6), 1510–1529.
- Vincent, L. A., X. L. Wang, E. J. Milewska, H. Wan, F. Yang, and V. Swail, 2012: A second generation of homogenized Canadian monthly surface air temperature for climate trend analysis. *Journal of Geophysical Research: Atmospheres*, **117** (D18).