

Full Bayesian Significance Testing for Neural Networks in Traffic Forecasting: Appendix

Zehua Liu¹, Jingyuan Wang^{1,2,3*}, Zimeng Li¹ and Yue He⁴

¹School of Computer Science and Engineering, Beihang University, Beijing, China

²School of Economics and Management, Beihang University, Beijing, China

³Key Laboratory of Data Intelligence and Management (Beihang University),

Ministry of Industry and Information Technology, Beijing, China

⁴Department of Computer Science and Technology, Tsinghua University, Beijing, China

{liuzehua, jywang, zimengli}@buaa.edu.cn, heyuethu@mail.tsinghua.edu.cn

1 Code

<https://github.com/liuzh-buaa/ST-nFBST>.

2 Theoretical Derivation of Optimization with Uncertainties

In this section, we provide a detailed theoretical derivation for the optimization of a Bayesian Neural Network (BNN) in the context of heteroscedastic aleatoric uncertainty and epistemic uncertainty. We demonstrate the optimization function \mathcal{L} through variational inference.

2.1 Optimization function \mathcal{L}

In the task of Bayesian Spatial-Temporal Modeling, we use a BNN, whose parameters θ follow a distribution rather than deterministic values, to represent the underlying function $f_0 : \mathbb{R}^{\tau_1 \times |\mathcal{V}|} \rightarrow \mathbb{R}^{\tau_2 \times |\mathcal{V}|}$. Given a dataset of n samples $\mathcal{D} = \{(\mathbf{X}_1, \mathbf{Y}_1), \dots, (\mathbf{X}_n, \mathbf{Y}_n)\}$, For a pair of observations $(\mathbf{X}_i, \mathbf{Y}_i)$, the regression process can be modeled as corrupted with Gaussian random noise:

$$\mathbf{Y}_i \sim \mathcal{N}(f_0(\mathbf{X}_i), \sigma_0^2(\mathbf{X}_i)). \quad (1)$$

Moreover, we assume the observation noise can vary with inputs \mathbf{X}_i at different time points, namely heteroscedastic. For convenience, we will abbreviate $\sigma_0(\mathbf{X}_i)$ as σ_0 if not causing any confusion.

Before optimization, a prior distribution is assigned to model parameters θ as an initial belief $\pi(\theta)$ according to experience. This belief is gradually adjusted to fit data \mathcal{D} by using the Bayesian rule. The final belief is presented as the posterior distribution

$$P(\theta|\mathcal{D}) = \frac{\pi(\theta)P(\mathcal{D}|\theta)}{P(\mathcal{D})} = \frac{\pi(\theta) \prod_{i=1}^n P(\mathbf{Y}_i|\mathbf{X}_i, \theta)}{\int_{\Theta} \prod_{i=1}^n P(\mathbf{Y}_i|\mathbf{X}_i, \theta) d\theta}. \quad (2)$$

However, it is intractable to solve the integral in Eq (2). Through variational inference, we use a tractable distribution to approximate the real but intractable posterior distribution. Formally, variational family $Q = \{q_{\vartheta} : \vartheta \in \Gamma\}$ is a pre-defined family of tractable distributions on model parameter space Θ , where ϑ is the parameter of variational distribution

and Γ is the range of ϑ . The optimal variational distribution q_{ϑ^*} is chosen from Q such that

$$\vartheta^* = \arg \min_{\vartheta \in \Gamma} \text{KL}(q_{\vartheta}(\theta) \| P(\theta|\mathcal{D})), \quad (3)$$

where KL divergence describes the “distance” between two distributions. According to its definition, the optimization objective function can be formulated as

$$\mathcal{L} = \text{KL}(q_{\vartheta}(\theta) \| P(\theta|\mathcal{D})) = \int_{\Theta} q_{\vartheta}(\theta) \log \frac{q_{\vartheta}(\theta)}{P(\theta|\mathcal{D})} d\theta. \quad (4)$$

From the Bayesian rule Eq (2), Eq (4) can be simplified to

$$\begin{aligned} \mathcal{L} &= \int_{\Theta} q_{\vartheta}(\theta) \log \frac{q_{\vartheta}(\theta)}{P(\theta|\mathcal{D})} d\theta \\ &= \int_{\Theta} q_{\vartheta}(\theta) \log \left(q_{\vartheta}(\theta) \frac{P(\mathcal{D})}{P(\mathcal{D}|\theta)\pi(\theta)} \right) d\theta \\ &= - \int_{\Theta} q_{\vartheta}(\theta) \log P(\mathcal{D}|\theta) d\theta + \int_{\Theta} q_{\vartheta}(\theta) \log \frac{q_{\vartheta}(\theta)}{\pi(\theta)} d\theta \\ &\quad + \int_{\Theta} q_{\vartheta}(\theta) \log p(\mathcal{D}) d\theta \\ &= -\mathbb{E} [\log P(\mathcal{D}|\theta)] + \text{KL}(q_{\vartheta}(\theta) \| \pi(\theta)) + \log P(\mathcal{D}). \end{aligned} \quad (5)$$

The third term $\log P(\mathcal{D})$ is a constant once the data \mathcal{D} is determined. In the following, we will present how to calculate the first two terms, respectively.

*Corresponding author

42 **2.2 The first term of \mathcal{L}**

43 Under the assumption of Eq (1), the log-likelihood of \mathcal{D} is
44 calculated as follows:

$$\begin{aligned} \log P(\mathcal{D}|\theta) &= \log \prod_{i=1}^n P(\mathbf{Y}_i|\mathbf{X}_i;\theta) \\ &= \sum_{i=1}^n \log P(\mathbf{Y}_i|\mathbf{X}_i;\theta) \\ &= \sum_{i=1}^n \log \frac{1}{\sqrt{2\pi}\sigma_0} \exp \left[-\frac{(\mathbf{Y}_i - f(\mathbf{X}_i;\theta))^2}{2\sigma_0^2} \right] \\ &= -\sum_{i=1}^n \log \sqrt{2\pi}\sigma_0(\mathbf{X}_i) \\ &\quad - \sum_{i=1}^n \frac{1}{2\sigma_0^2(\mathbf{X}_i)} (\mathbf{Y}_i - f(\mathbf{X}_i;\theta))^2. \end{aligned} \quad (6)$$

45 The second term of Eq (6) is similar to the Mean Squared
46 Error (MSE) in the regression task within a scaling factor.

47 Using Monte Carlo integration, the first term of \mathcal{L} can be
48 calculated as follows:

$$-\mathbb{E}[\log P(\mathcal{D}|\theta)] \approx -\frac{1}{m} \sum_{j=1}^m [\log P(\mathcal{D}|\theta_j)], \quad (7)$$

49 where θ_j is obtained by sampling m times based on the variational
50 distribution $q_\vartheta(\theta)$.

51 **2.3 The second term of \mathcal{L}**

52 In our experiment, we adopt popular diagonal Gaussian dis-
53 tributions as the prior distribution $\pi(\theta)$ and the variational
54 distribution $q_\vartheta(\theta)$ of parameters θ . We first consider θ as a
55 one-dimensional variable, that is:

$$\pi(\theta) \sim \mathcal{N}(\mu_\pi, \sigma_\pi^2) \quad (8)$$

56 and correspondingly

$$q_\vartheta(\theta) \sim \mathcal{N}(\mu, \sigma^2). \quad (9)$$

57 Then we get

$$\begin{aligned} \text{KL}(q_\vartheta(\theta)\|\pi(\theta)) &= \int q_\vartheta(\theta) \log \frac{q_\vartheta(\theta)}{\pi(\theta)} d\theta \\ &= \int q_\vartheta(\theta) \log q_\vartheta(\theta) d\theta \\ &\quad - \int q_\vartheta(\theta) \log \pi(\theta) d\theta \\ &= -\frac{1}{2} \left(1 + \log(2\pi\sigma^2) \right) \\ &\quad + \frac{1}{2} \left(\log 2\pi\sigma_\pi^2 + \frac{\sigma^2 + (\mu - \mu_\pi)^2}{\sigma_\pi^2} \right) \\ &= \log \frac{\sigma_\pi}{\sigma} + \frac{\sigma^2 + (\mu - \mu_\pi)^2}{2\sigma_\pi^2} - \frac{1}{2}. \end{aligned} \quad (10)$$

Then, we consider θ follows a n -dimensional diagonal Gaussian distribution. That is,

$$\pi(\theta) \sim \mathcal{N}(\mu_\pi, \Sigma_\pi), \quad (11)$$

where

$$\mu_\pi = \begin{bmatrix} \mu_{\pi 1} \\ \mu_{\pi 2} \\ \vdots \\ \mu_{\pi n} \end{bmatrix}, \Sigma_\pi = \begin{bmatrix} \sigma_{\pi 1}^2 & 0 & \cdots & 0 \\ 0 & \sigma_{\pi 2}^2 & \cdots & 0 \\ \vdots & \cdots & \cdots & \vdots \\ 0 & 0 & \cdots & \sigma_{\pi n}^2 \end{bmatrix}. \quad (12)$$

Correspondingly,

$$q_\vartheta(\theta) \sim \mathcal{N}(\mu, \Sigma), \quad (13)$$

where

$$\mu = \begin{bmatrix} \mu_1 \\ \mu_2 \\ \vdots \\ \mu_n \end{bmatrix}, \Sigma = \begin{bmatrix} \sigma_1^2 & 0 & \cdots & 0 \\ 0 & \sigma_2^2 & \cdots & 0 \\ \vdots & \cdots & \cdots & \vdots \\ 0 & 0 & \cdots & \sigma_n^2 \end{bmatrix}. \quad (14)$$

Further, $q_\vartheta(\theta)$ can be denotes as

$$q_\vartheta(\theta) = \frac{1}{\sqrt{(2\pi)^n |\Sigma|}} \exp[-\frac{1}{2}(\theta - \mu)^T \Sigma^{-1}(\theta - \mu)], \quad (15)$$

and $\pi(\theta)$ can be denotes as

$$\pi(\theta) = \frac{1}{\sqrt{(2\pi)^n |\Sigma_\pi|}} \exp[-\frac{1}{2}(\theta - \mu_\pi)^T \Sigma_\pi^{-1}(\theta - \mu_\pi)]. \quad (16)$$

Then, we can get

$$\begin{aligned} \text{KL}(q_\vartheta(\theta)\|\pi(\theta)) &= \int q_\vartheta(\theta) \log \frac{q_\vartheta(\theta)}{\pi(\theta)} d\theta \\ &= \int q_\vartheta(\theta) \log q_\vartheta(\theta) d\theta \\ &\quad - \int q_\vartheta(\theta) \log \pi(\theta) d\theta \\ &= \log \frac{1}{\sqrt{(2\pi)^n |\Sigma|}} - \log \frac{1}{\sqrt{(2\pi)^n |\Sigma_\pi|}} \\ &\quad - \frac{1}{2} \int q_\vartheta(\theta)(\theta - \mu)^T \Sigma^{-1}(\theta - \mu) d\theta \\ &\quad + \frac{1}{2} \int q_\vartheta(\theta)(\theta - \mu_\pi)^T \Sigma_\pi^{-1}(\theta - \mu_\pi) d\theta \\ &= \frac{1}{2} \log \frac{|\Sigma_\pi|}{|\Sigma|} - \frac{1}{2} \mathbb{E}(\theta - \mu)^T \Sigma^{-1}(\theta - \mu) \\ &\quad + \frac{1}{2} \mathbb{E}(\theta - \mu_\pi)^T \Sigma_\pi^{-1}(\theta - \mu_\pi) \end{aligned} \quad (17)$$

66 According to Theorem 1, we can simplify Eq (17) as

$$\begin{aligned}
\text{KL}(q_\theta(\theta)\|\pi(\theta)) &= \frac{1}{2} \log \frac{|\Sigma_\pi|}{|\Sigma|} - \frac{1}{2} \text{tr}(\Sigma^{-1}\Sigma) + \frac{1}{2} \text{tr}(\Sigma_\pi^{-1}\Sigma) \\
&\quad - (\mu - \mu)^T \Sigma^{-1}(\mu - \mu) \\
&\quad - (\mu - \mu_\pi)^T \Sigma_\pi^{-1}(\mu - \mu_\pi) \\
&= \frac{1}{2} \log \frac{|\Sigma_\pi|}{|\Sigma|} - \frac{1}{2}n + \frac{1}{2} \text{tr}(\Sigma_\pi^{-1}\Sigma) \\
&\quad + \frac{1}{2}(\mu - \mu_\pi)^T \Sigma_\pi^{-1}(\mu - \mu_\pi) \\
&= \sum_{i=1}^n \left(\log \frac{\sigma_{\pi i}}{\sigma_i} + \frac{\sigma_i^2 + (\mu_i - \mu_{\pi i})^2}{2\sigma_{\pi i}^2} - \frac{1}{2} \right)
\end{aligned} \tag{18}$$

67 The result is consistent with Eq (10).

68 In conclusion, we have analyzed the meanings of different
69 terms in the objective function Eq (5). The first term Eq (7)
70 is related to data. The second term Eq (18) is only related to
71 $\vartheta = (\mu, \sigma)$ further θ like a regularization term. The third term
72 $\log P(\mathcal{D})$ is a constant once the data \mathcal{D} is determined.

73 **Theorem 1.** Given $\mathbf{x} \sim \mathcal{N}(\mu, \Sigma)$, where \mathbf{x} is a n dimensional
74 vector, that is $\mathbf{x} \in \mathbb{R}^{n \times 1}$. If we assume A is a $n \times n$ matrix,
75 we have

$$\mathbb{E}(\mathbf{x}^T A \mathbf{x}) = \text{tr}(A\Sigma) + \mu^T A \mu. \tag{19}$$

76 *Proof.* The trace of a matrix (denoted as “tr”) is defined as
77 the sum of the principle diagonal elements of a matrix, which
78 has the following properties:

$$\begin{aligned}
\text{tr}(M_1 M_2) &= \text{tr}(M_2 M_1) \\
\Rightarrow \text{tr}(M_1 M_2 M_3) &= \text{tr}(M_2 M_3 M_1) = \text{tr}(M_3 M_1 M_2),
\end{aligned} \tag{20}$$

79 where M_1, M_2, M_3 are three compatible matrices. Therefore,
80 we have

$$\mathbf{x}^T A \mathbf{x} = \mathbb{E}(\mathbf{x}^T A \mathbf{x}) = \text{tr}(A \mathbf{x} \mathbf{x}^T). \tag{21}$$

81 According to the properties of n -dimensional diagonal Gaussian distribution, we have

$$\begin{aligned}
\Sigma &= \mathbb{E}((\mathbf{x} - \mu)(\mathbf{x} - \mu)^T) \\
&= \mathbb{E}((\mathbf{x} - \mu)(\mathbf{x}^T - \mu^T)) \\
&= \mathbb{E}(\mathbf{x}\mathbf{x}^T - \mathbf{x}\mu^T - \mu\mathbf{x}^T + \mu\mu^T) \\
&= \mathbb{E}(\mathbf{x}\mathbf{x}^T) - \mathbb{E}(\mathbf{x})\mu^T - \mu\mathbb{E}(\mathbf{x}^T) + \mu\mu^T \\
&= \mathbb{E}(\mathbf{x}\mathbf{x}^T) - \mu\mu^T.
\end{aligned} \tag{22}$$

83 This is equivalent to

$$\mathbb{E}(\mathbf{x}\mathbf{x}^T) = \Sigma + \mu\mu^T. \tag{23}$$

84 Therefore,

$$\begin{aligned}
\mathbb{E}(\mathbf{x}^T A \mathbf{x}) &= \mathbb{E}(\text{tr}(A \mathbf{x} \mathbf{x}^T)) = \text{tr}(\mathbb{E}(A \mathbf{x} \mathbf{x}^T)) \\
&= \text{tr}(A \mathbb{E}(\mathbf{x} \mathbf{x}^T)) \\
&= \text{tr}(A(\Sigma + \mu\mu^T)) \\
&= \text{tr}(A\Sigma) + \text{tr}(A\mu\mu^T).
\end{aligned} \tag{24}$$

Then, based on Eq (21),

$$\mathbb{E}(\mathbf{x}^T A \mathbf{x}) = \text{tr}(A\Sigma) + \mu^T A \mu. \tag{25}$$

□ 86

3 Experimental Setup

In this section, we introduce a detailed experimental setup, including hardware settings and algorithm parameters. All the experiments are conducted on an Ubuntu machine equipped with two Intel(R) Xeon(R) Silver 4214R CPUs @ 2.40GHz with 12 physical cores, and four NVIDIA GeForce RTX 3090 GPUs, armed with 24GB of GDDR6X memory. We train the BNN of ST-nFBST using the PyTorch package (v1.7.1). The descriptions and settings of related hyper-parameters are listed as follows:

- lr. The initial learning rate of Adam optimizer. The default value is 0.005.
 - lr_step_size. The period of learning rate decay. The default value is [5, 20, 40, 70].
 - lr_decay_ratio. The multiplicative factor of learning rate decay. The default value is 0.5.
 - τ . The threshold to control the precision of $\hat{\sigma}_0$ obtained by the Sigma-Decoder. That is, $\hat{\sigma}_0 = \max(\hat{\sigma}_0, \tau)$. The default value is 0.01.
 - μ_π . The mean of the prior distribution of θ in Eq (11). The default value is 0, which is a vector composed of 0.
 - Σ_π . The (co)variance matrix of the prior distribution of θ in Eq (11). The default value is an identity matrix I .
 - max_diffusion_step. The max diffusion step defined in the diffusion convolution operation. The default value is 2.
- Besides, as a core component of ST-nFBST, we need to sample through the Bayesian neural network and then approximate the posterior probability density by Kernel Density Estimation (KDE). The sample size we set is fifty. We use the Gaussian kernel function and choose the best bandwidth from {0.01, 0.05, 0.1, 1.0} by GridSearchCV for each sampling through 5-fold cross-validation. The bandwidth here acts as a smoothing parameter, controlling the trade-off between bias and variance in the result. A large bandwidth leads to a very smooth (that is, high-bias) density distribution while a small bandwidth leads to an unsmooth (that is, high-variance) density distribution.

4 More Experimental Results on the PEMS-BAY Dataset

Traffic forecasting on the METR-LA (Los Angeles, which is known for its complicated traffic conditions) dataset is more challenging than that on the PEMS-BAY dataset. Thus we use METR-LA as the default dataset for analysis in the main paper, and more results on the PEMS-BAY dataset are provided here. The conclusions obtained from two datasets are consistent.

Figure 1 illustrates the trends of MAE and uncertainties under different prediction horizons. Sensor 3 and 6 are central sensors, while sensor 212 belongs to marginal sensors. We have the following observations:

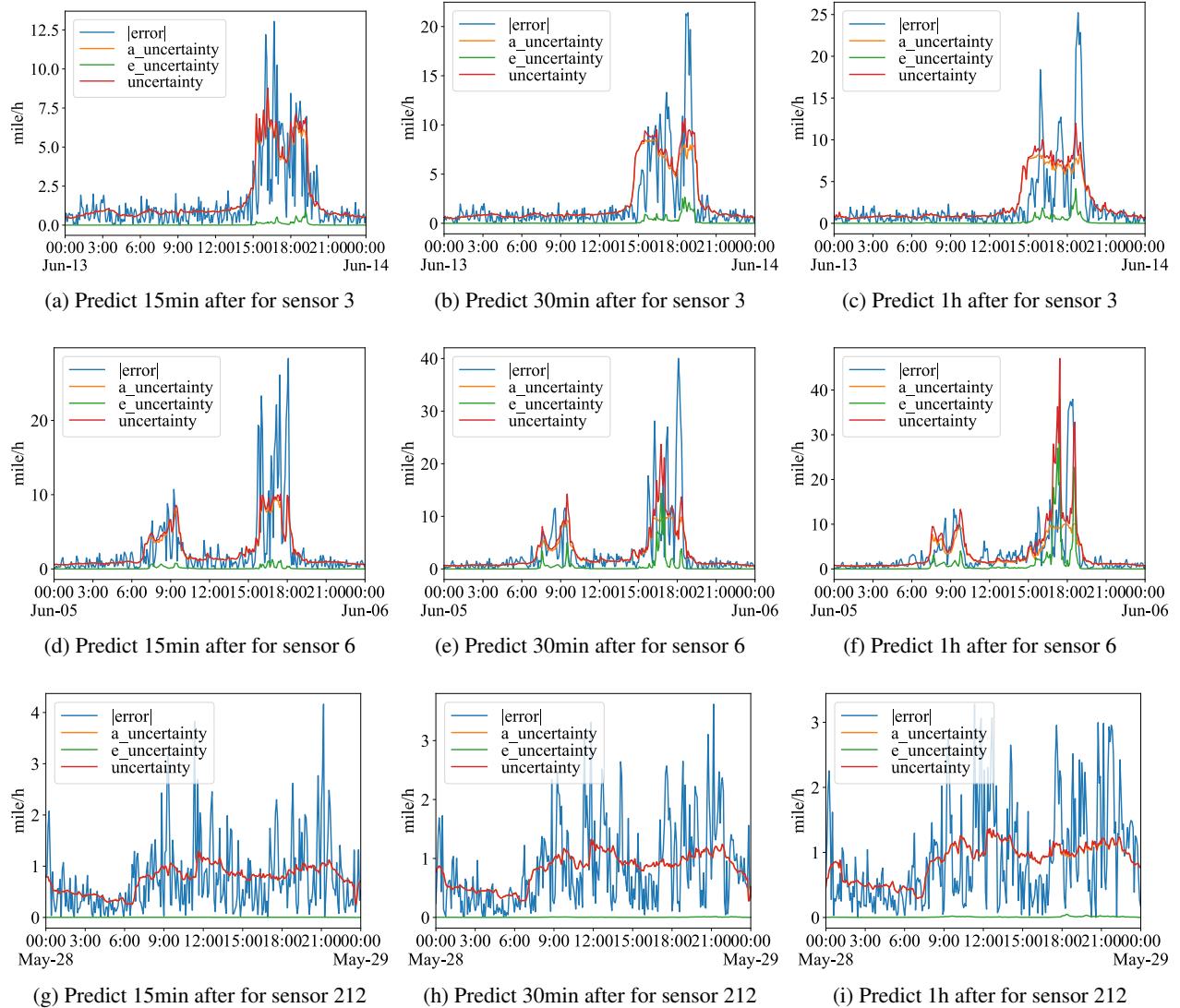


Figure 1: Comparison between uncertainties and $|\text{error}|$ for different forecasting horizons on the PEMS-BAY dataset.

- 137 • The fluctuation of MAE closely aligns with the variations
138 in uncertainty across diverse forecasting horizons. This
139 implies that our method accurately captures uncertainties.
140 Considering that we thoroughly leverage both uncertainties
141 in the optimization process by incorporating them into
142 the loss function, this further explains why our approach
143 outperforms other state-of-the-art models.
- 144 • In most cases, the aleatoric uncertainty is much larger than
145 the epistemic uncertainty, accounting for the majority of
146 uncertainty. This implies that the traffic uncertainty is
147 mainly data-related. Thanks to the powerful representation
148 capability of neural networks, we can train models effec-
149 tively to reduce epistemic uncertainty.
- 150 • There are some cases where error and uncertainties man-
151 ifest inconsistently, specifically, with low error but high
152 uncertainty. Moreover, epistemic uncertainty is high in

such cases. This occurs when changes in road conditions
153 at intersections, such as traffic congestion or unexpected
154 accidents, result in a sharp decrease in the speed of other
155 sensors. This implies that uncertainty, compared to error,
156 is more effective in monitoring changes in road conditions.
157 A more detailed case study is presented in Section 5.

- 158
- 159 • As the forecasting horizon expands, both error and uncer-
160 tainties of central sensors increase correspondingly, while
161 those of marginal sensors remain relatively stable. How-
162 ever, the predictions of marginal sensors fluctuate signifi-
163 cantly at different times. This aligns with the definition of
164 two types of sensors. The marginal sensors exhibit fewer
165 impacts from other sensors while utilizing less information
166 from other sensors. However, the central sensors are more
167 susceptible to the influence exerted by other sensors, ren-
168 dering predictions more intricate.

169 **5 Case Study**

170 In this section, we provide a case study to show the conclu-
171 sion obtained by significance testing spatially about the influ-
172 ence of intersections. Furthermore, we analyze the reasons
173 why error is inconsistent with uncertainty in some cases.

174 Figure 2,3,4,5 show the testing results for different sensors
175 under various error and uncertainty levels. They all represent
176 the testing results for 15 minutes forecasting horizon, repre-
177 senting short-term forecasting. Figure 2 and 4 correctly iden-
178 tify that sensors near intersections are insignificant, hence the
179 low error and uncertainty. On the other hand, the high error
180 and uncertainty in Figure 3 and 5 is attributed to their deter-
181 mination that surrounding sensors are insignificant, but the
182 sensors at intersections are significant. By comparing them,
183 it can be observed that short-term traffic forecasting is more
184 concerned with the road conditions where their own sensors
185 are located, but do not pay attention to the situation at inter-
186 sections. The conclusion is verified by the real dataset. For
187 example, the true speed in Figure 5c is 58.63 miles/h, while
188 the prediction is 19.97 miles/h. The prediction is significantly
189 lower than the true speed. This discrepancy arises from the
190 fact that, within a very short time frame in the recent past,
191 although the true speeds at the intersections are indeed slow,
192 the speeds of sensors on the same road as sensor 149 remain
193 high. However, the testing results indicate that information
194 from the intersection is utilized at this time, leading to high
195 errors and increased uncertainty.

196 Subsequently, we conduct the significance testing for the
197 same sensor at the same moment as Figure 3c, with the fore-
198 casting horizon changed to 1h, representing long-term fore-
199 casting. As depicted in Figure 6, the error significantly de-
200 creased, but uncertainty remained high. By observing the
201 speeds of different sensors, we found that those of sensors on
202 the same road as the interest sensor are high, while the speeds
203 at intersections are considerably slow. We conclude that the
204 changes at intersections take time to propagate to other road
205 segments. Therefore, short-term forecasting should not fo-
206 cus on them, but they have a significant impact on long-term
207 forecasting. Due to a sharp transition from high to low pre-
208 diction speed and the testing of a significant impact at the in-
209 tersection, the uncertainty of predictions in Figure 6 remains
210 high. This also confirms that uncertainty, compared to error,
211 is more effective in monitoring changes in road conditions.

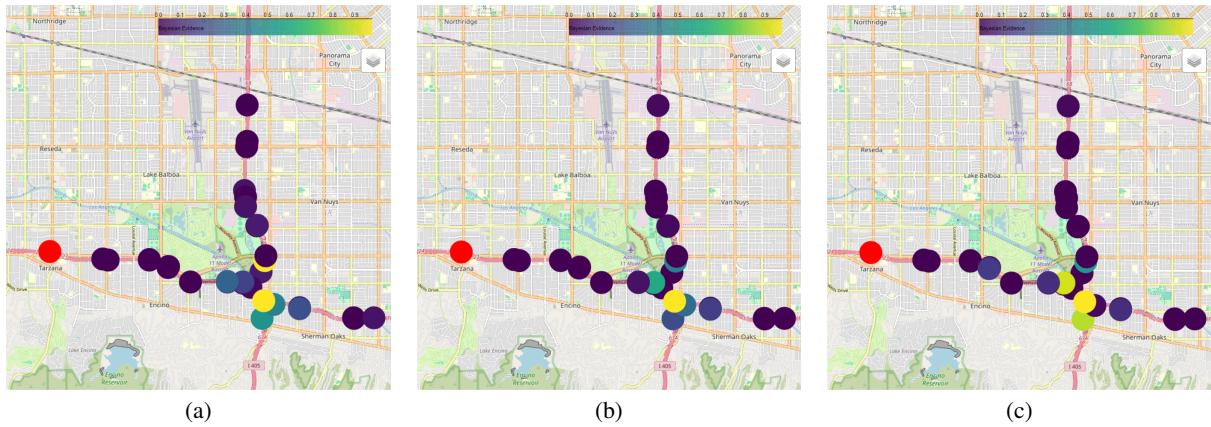


Figure 2: Spatial testing results for sensor 181 (highlighted in red) under low uncertainty and low error at different times.

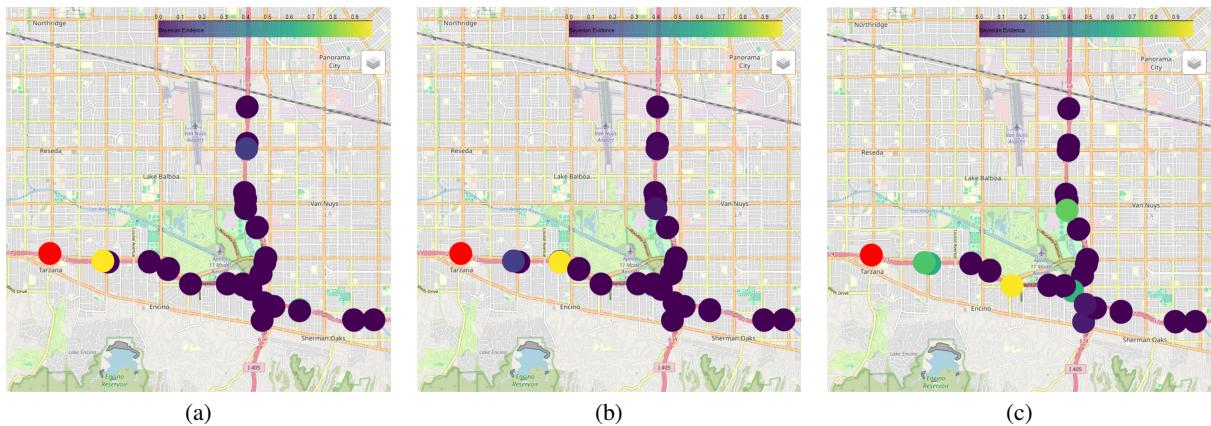


Figure 3: Spatial testing results for sensor 181 (highlighted in red) under high uncertainty and high error at different times.

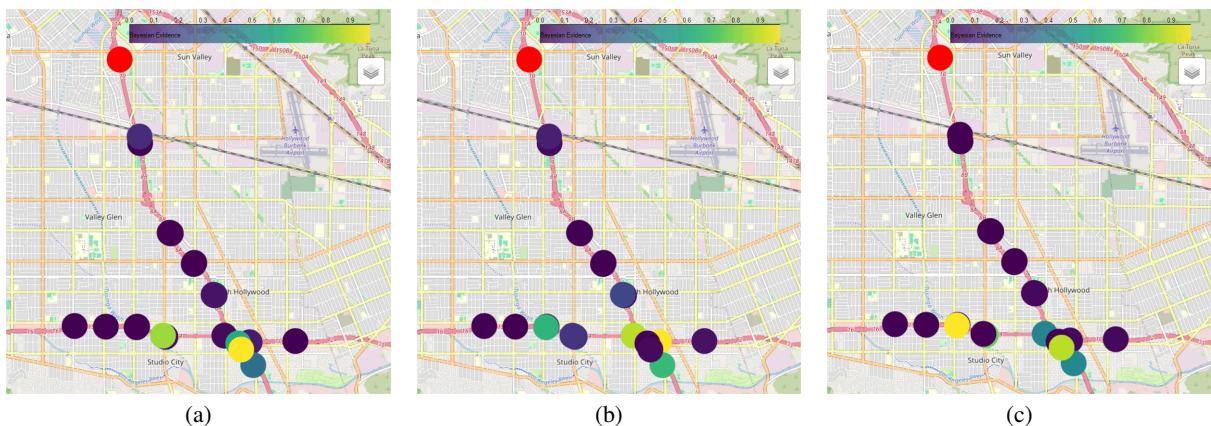


Figure 4: Spatial testing results for sensor 149 (highlighted in red) under low uncertainty and low error at different times.

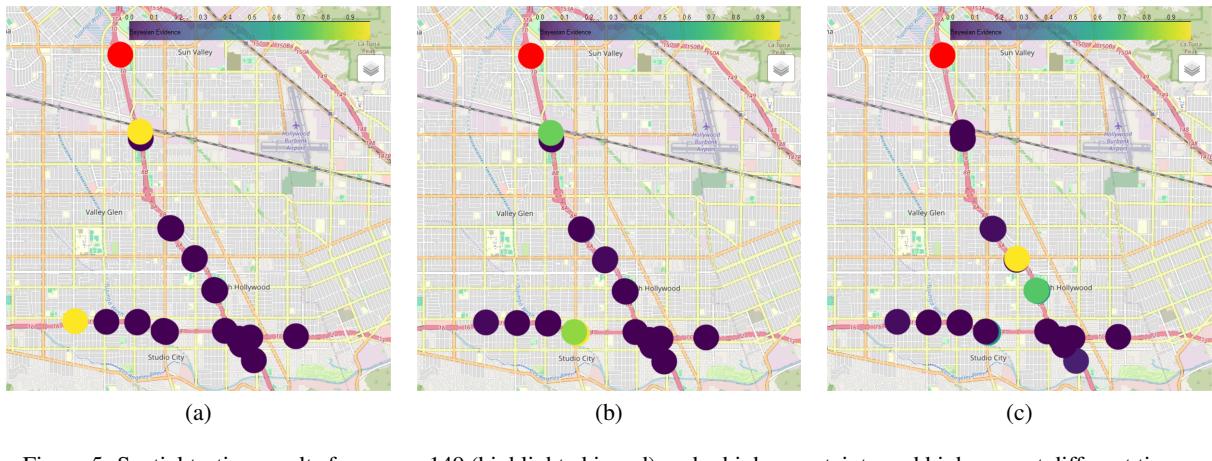


Figure 5: Spatial testing results for sensor 149 (highlighted in red) under high uncertainty and high error at different times.

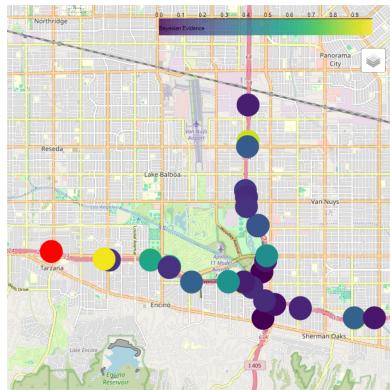


Figure 6: Spatial testing results for sensor 181 under low error but high uncertainty level for 1h forecasting horizon.