DeepTransport: Learning Spatial-Temporal Dependency for Traffic Condition Forecasting

Xingyi Cheng^{†*}, Ruiqing Zhang[†], Jie Zhou, Wei Xu

Baidu Research - Institute of Deep Learning {chengxingyi,zhangruiqing01,zhoujie01,wei.xu}@baidu.com

Abstract

Predicting traffic conditions has been recently explored as a way to relieve traffic congestion. Several pioneering approaches have been proposed based on traffic observations of the target location as well as its adjacent regions, but they obtain somewhat limited accuracy due to lack of mining road topology. To address the effect attenuation problem, we propose to take account of the traffic of surrounding locations(wider than adjacent range). We propose an endto-end framework called DeepTransport, in which Convolutional Neural Networks (CNN) and Recurrent Neural Networks (RNN) are utilized to obtain spatial-temporal traffic information within a transport network topology. In addition, attention mechanism is introduced to align spatial and temporal information. Moreover, we constructed and released a real-world large traffic condition dataset with 5-minute resolution. Our experiments on this dataset demonstrate our method captures the complex relationship in temporal and spatial domain. It significantly outperforms traditional statistical methods and a state-of-the-art deep learning method.

Introduction

With the development of location-acquisition and wireless device, a vast amount of data with spatial transport networks and timestamps can be collected by mobile phone map app. The majority of map apps can tell users real-time traffic conditions, as shown in Figure 1. However, only the current traffic conditions are not enough for making effective route planing, a traffic system to predict future road condition may be more valuable.

In the past, there are mainly two approaches for traffic prediction: time-series analysis based on classical statistics and data-driven methods based on machine learning. Most former methods are univariate; they predict the traffic of a place at a certain time. The fundamental work was Auto Regressive Integrated Moving Average (ARIMA) (Ahmed and Cook 1979) and its variations (Pan, Demiryurek, and Shahabi 2012; Williams and Hoel 1999). Motivated by the fact (Williams 2001) that traffic evolution is a temporal-spatial phenomenon, multivariate methods with both temporal and spatial features was proposed. (Stathopoulos and Karlaftis 2003) developed a model that feeds on data



Figure 1: A real-time traffic network example from a commercial map app, the networks including many locations and the color(green, yellow, red, dark red) depth illustrated the condition of a location(a stretch of road).

from upstream detectors to improve the predictions of downstream locations. However, many statistics are needed in such methods. On the other hand, data-driven methods (Jeong et al. 2013; Vlahogianni, Karlaftis, and Golias 2005) fit a single model from vector-valued observations including historical scalar measurements with the trend, seasonal, cyclical, and calendar variations. For instance, (Deng et al. 2016) expressed traffic pattern by mapping road attributes to a latent space. However, the linear model here is limited in its ability to extract effective features.

Neural networks and deep learning have been demonstrated as a unified learning framework for feature extraction and data modeling. Since its applicability in this topic, significant progress has been made in related work. Firstly, both temporal and spatial dependencies between observations in time and space are complex and can be strongly nonlinear. While the statistics frequently fail when dealing with nonlinearity, neural networks are powerful to capture very complex relations (LeCun, Bengio, and Hinton 2015). Secondly, neural networks can be trained with raw data in an end-to-end manner. Apparently, hand-crafted engineered features that extract all information from data spread in time and space is laborious. Data-driven based neural networks extracts features without the need of statistical feature, e.g., mean or variance of all adjacent lo-

^{*}Xingyi Cheng is the corresponding author.

[†]main contribution

cations of the current location. The advantage of neural networks for traffic prediction has long been discovered by researchers. Some early work (Chang and Su 1995; Innamaa 2000) simply put observations into input layer, or take sequential feature into consideration (Dia 2001a) to capture temporal patterns in time-series. Until the last few years, some works of deep learning was applied. For instance, Deep Belief Networks (DBN) (Huang et al. 2014) and Stack Autoencoders (SAEs) (Lv et al. 2015). However, input data in these works are directly concatenated from different locations, which ignored the spatial relationship. In general, the existing methods either concerns with the time series or just a little use of the spatial information. Depending on traffic condition of a "narrow" spatial range will undoubtely degrades prediction accuracy. To achieve a better undestanding of spatial information, we propose to solve this problem by taking the intricate topological graph as a key feature in traffic condition forecasting, especially for long prediction horizon.

To any target location as the center of radiation, surrounding locations with same order form a "width" region, and regions with different order constitute a "depth" sequence. We propose a double sequential deep learning model to explore the traffic condition pattern. This model adopts a combination of convolutional neural networks (CNN) (LeCun, Bengio, and others 1995) and recurrent networks with long short-term memory (LSTM) units (Hochreiter and Schmidhuber 1997) to deal with spatial dependencies. CNN is responsible for maintaining the "width" structure, while LSTM for the "depth" structure. To depict the complicated spatial dependency, we utilize attention mechanism to demonstrate the relationships between time and space.

The main contribution of the paper is summarized as follows:

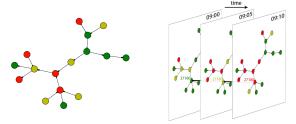
- We introduce a novel deep architecture to enable temporal and dynamical spatial modeling for traffic condition forecasting.
- We propose the necessity of aligning spatial and temporal information and introduce attention mechanism into the model to quantify their relationship. The obtained attention weight is helpful for daily traveling and path planning.
- Experiment results demonstrate that the proposed model significantly outperforms existing methods based on deep learning and time series forecasting methods.
- We also release a real large (millions) traffic dataset with topological networks and temporal traffic condition ¹.

Preliminary

In this section, we briefly revisit the traffic prediction problem and introduce notations in this work.

Common Notations and Definition

A traffic network can be represented in a graph in two ways. Either monitoring the traffic flow of crossings, take crossing



(a) A plain graph at a time point (b) A graph with time-series

Figure 2: Traffic condition. Five colors in this graph denote five states for visually displaying: green(1, fluency), yellow(2, slow), red(3, congestion) and dark red(4, extreme congestion). "27180" is the ID number of a location(road section).

as node and road as an edge of graph, or conversely, monitoring the condition of roads, take roads as nodes and crossings as connecting edges. The latter annotation is adopted in our work. Taking figure 2(a) as an example, each colored node corresponds to a stretch of road in a map app.

We consider a graph consists of weighted vertices and directed edges. Denote the graph as $G = \langle V, E \rangle$. V is the set of vertices and $E \subseteq \{(u,v)|u\in V,v\in V\}$ is the set of edges, where (u,v) is an ordered pair. A location(vertex) v at any time point t have five traffic condition states $c(v,t)\in\{0,1,2,3,4\}$, expressing not-released, fluency, slow, congestion, extreme congestion respectively. Figure 2(b) presents an example of road traffic at three time points in an area.

Observations: Each vertex in the graph is associated with a feature vector, which consists of two parts, time-varying O and time-invariant variables F. Time-varying variables that characterize the traffic network dynamically are traffic flow observation aggregated by a 5-minutes interval. Time-invariant variables are static features as natural properties which do not change with time s, such as the number of input and output degrees of a road, its length, limit speed and so forth.

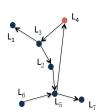
In particular, the time-varying and time-invariant variables are denoted as:

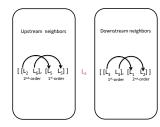
$$\mathbf{O}_{v,t} = \begin{bmatrix} c(v,t) \\ c(v,t-1) \\ \vdots \\ c(v,t-p) \end{bmatrix} \quad \mathbf{F}_v = \begin{bmatrix} f_{v,1} \\ f_{v,2} \\ \vdots \\ f_{v,k} \end{bmatrix}$$
(1)

where c(v,t) is traffic condition of vertex v at time t, p is the length of historical measurement. $f_{v,k}$ are time-invariant features.

Order Slot: In a path of the directed graph, the number of edges required to take from one vertex to another is called order. Vertices of the same order constitute an order slot. Directly linked vertices are termed first-order neighbors. Second-order spatial neighbors of a vertex are the first-order neighbors of its first-order neighbors and so forth. For any vertex in our directed graph, we define the incoming

¹https://github.com/cxysteven/MapBJ





- (a) The direct graph of a location.
- (b) Upstream flow and down-stream flow neighbor of L_4 with in order 2.

Figure 3: An example of directed graph and order slot notation in DeepTransport

traffic flow as its upstream flow and the outflow as its downstream flow. Take figure 3(a) as an example, L_4 is the target location to be predict. L_3 is the first-order downstream vertex of L_4 . L_1 , L_2 is the first order downstream set of L_3 and they constitute the second order slot of L_4 . Each vertex in the traffic flow that goes in one direction is affected by its upstream flow and downstream flow. The first and second order slots of L_4 is shown in Figure 3(b). Introducing the dimension of time series, any location $L_{v,t}$ is composed of two vectors, $O_{v,t}$ and F_v . Any order slot consists of some locations:

$$\mathbf{L}_{v,t} = \begin{bmatrix} \mathbf{O}_{v,t} \\ \mathbf{F}_v \end{bmatrix} \quad \mathbf{X}_{v,t}^j = \begin{bmatrix} \mathbf{L}_{u_1,t}^T \\ \mathbf{L}_{u_2,t}^T \\ \vdots \\ \mathbf{L}_{u_k,t}^T \end{bmatrix}$$
(2)

where location index u is one of the jth order neighbors of v.

Perceptive Radius: The maximum ordered number controls the perceptive scope of the target location. It is an important hyperparameter describing spatial information, we call it perceptive radius and denote it as r.

Problem Definition: According to the above notation, we define the problem as follows: Predict a sequence of traffic flow $\mathbf{L}_{v,t+h}$ for prediction horizon h given the historical observations of $\mathbf{L}_{v',t'}$, where $v' \in neighbor(v,r)$, $t' \in \{t-p,\cdots,t\}, r \in \{0,\cdots,R\}$ is perceptive radius and p is the length of historical measurement.

Model

As shown in Figure 4, our model consists of four parts: upstream flow observation(left), target location module(middle), downstream flow observation(right), and training cost module(top). In this section, we detail the work process of each module.

Spatial-temporal Relation Construction

Since the traffic condition of a road is strongly impacted by its upstream and downstream flow, we use a convolutional subnetwork and a recurrent subnetwork to maintain the road topology in the proposed model. Convolutional Layer CNN is used to extract temporal and "width" spatial information. As demonstrated in the example of figure 3, when feeding into our model, L_4 's first upstream neighbor L_5 should be copied twice, because there are two paths to L_4 , that are $[L_6, L_5]$ and $[L_2, L_5]$. With the exponential growth of paths, the model suffers from the high dimension and intensive computation. Therefore, we employ a convolution operation with multiple encoders and shared weights (LeCun, Bengio, and others 1995). To further reduce the parameter space while maintaining independence among vertices with the same order, we set the convolution stride to the convolution kernel window size, which is equal to the length of a vertex's observation representation.

The non-linear convolutional feature is obtained as follows:

$$\mathbf{e}_{up,q}^r = \sigma(\mathbf{W}_{up,q} * \mathbf{U}_{v,t} + \mathbf{b}_{up,q}), \tag{3}$$

$$\mathbf{e}_{down,q}^r = \sigma(\mathbf{W}_{down,q} * \mathbf{D}_{v,t} + \mathbf{b}_{down,q}), \quad (4)$$

where $\mathbf{U}_{v,t} = [\mathbf{X}_{v,t}^1, \cdots, \mathbf{X}_{v,t}^r]$ (only upstream neighbors) is denoted as upstream input matrix, while $\mathbf{D}_{v,t}$ is downstream input matrix. The $\mathbf{e}_{\cdot,q}^r$ is at rth order vector of upstream or downstream module where $q \in \{1,2...m\}$ and m is the number of feature map. We set $\mathbf{e}_{up}^r = [\mathbf{e}_{up,1}^r, \cdots, \mathbf{e}_{up,m}^r]$ and $\mathbf{e}_{up}^r \in \mathbb{R}^{l \times m}$, l is the number of observations in a slot. Similarly, we can get the \mathbf{e}_{down}^r . The weights \mathbf{W} and bias b composes parameters of CNN subnetworks. σ represents nonlinear activation, we empirically adopt the tanh function here.

Recurrent Layer RNN is utilized to represent each path that goes to the target location(upstream path) or go out from the target location(downstream path). The use of RNN have been investigated for traffic prediction for a long time, (Dia 2001b) used a Time-Lag RNN for short-term speed prediction(from 20 seconds to 15 minutes) and (Lint, Hooqendoorn, and Zuvlen 2002) adopted RNN to model state space dynamics for travel time prediction. In our proposed method, since the upstream flow from high-order to low-order, while the downstream flow is contrary, the output of the CNN layer in upstream module and downstream module are fed into RNN layer separately.

The structure of vehicle flow direction uses LSTM with "peephole" connections to encode a path as a sequential representation. In LSTM, the forget gate ${\bf f}$ controls memory cell ${\bf c}$ to erase, the input gate ${\bf i}$ helps to ingest new information, and the output gate ${\bf o}$ exposes the internal memory state outward. Specifically, given a rth slot matrix ${\bf e}_{down}^r \in \mathbb{R}^{l \times m}$, map it to a hidden representation ${\bf h}_{down}^r \in \mathbb{R}^{l \times d}$ with LSTM as follows:

$$\begin{bmatrix} \mathbf{\tilde{c}}^{r} \\ \mathbf{o}^{r} \\ \mathbf{i}^{r} \\ \mathbf{f}^{r} \end{bmatrix} = \begin{bmatrix} \tanh \\ \sigma \\ \sigma \\ \sigma \end{bmatrix} \left(\mathbf{W}_{p} \begin{bmatrix} \mathbf{e}^{r} \\ \mathbf{h}^{r-1} \end{bmatrix} + \mathbf{b}_{p} \right), \tag{5}$$

$$\mathbf{c}^r = \tilde{\mathbf{c}}^r \odot \mathbf{i}^r + \mathbf{c}^{r-1} \odot \mathbf{f}^r, \tag{6}$$

$$\mathbf{h}^r = [\mathbf{o}^r \odot \tanh(\mathbf{c}^r)]^T, \tag{7}$$

where $\mathbf{e}^r \in \mathbb{R}^{l \times m}$ is the input at the rth order step; $\mathbf{W}_p \in \mathbb{R}^{4d \times (m+d)}$ and $\mathbf{b}_p \in \mathbb{R}^{4d}$ are parameters of affine trans-

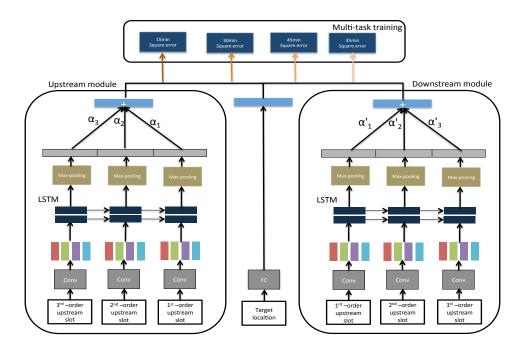


Figure 4: An example of the model architecture. There are three slots in upstream and downstream module respectively, each with input vertices length of two. The convolution operation has four sharing feature map. The middle block demonstrates that the target location is propagated by the fully-connected operation. A multi-task module that with four cost layers on the top block. Conv: Convolution; FC: Fully-connection.

formation; σ denotes the logistic sigmoid function and \odot denotes elementwise multiplication.

The update of upstream and downstream LSTM unit can be written precisely as follows:

$$\mathbf{h}_{down}^{r} = \mathbf{LSTM}(\mathbf{h}_{down}^{r-1}, \mathbf{e}_{down}^{r}, \theta_{p}). \tag{8}$$

$$\mathbf{h}_{up}^{r} = \mathbf{LSTM}(\mathbf{h}_{up}^{r+1}, \mathbf{e}_{up}^{r}, \theta_{p}). \tag{9}$$

The function $\mathbf{LSTM}(\cdot,\cdot,\cdot)$ is a shorthand for Eq. (5-7), in which θ_p represents all the parameters of \mathbf{LSTM} .

Slot Attention To get the representation of each order slot, max-pooling is performed on the output of LSTM. As \mathbf{h}^r represents the status sequence of the vertices in the corresponding order slot, we pool on each order slot to get r number of slot embeddings $\mathbf{S}_{up} = [\mathbf{s}_{up}^1, \cdots, \mathbf{s}_{up}^r]$ and $\mathbf{S}_{down} = [\mathbf{s}_{down}^1, \cdots, \mathbf{s}_{down}^r]$. Since different order slot have different effects on target prediction, we introduce attention mechanisms to align these embeddings. Given the target location hidden representation \mathbf{g} , we get the jth slot attention weights (Bahdanau, Cho, and Bengio 2014; Rocktschel et al. 2015) as follows:

$$\alpha_j = \frac{\exp a(\mathbf{g}, \mathbf{s}^j)}{\sum_{k=1}^r \exp a(\mathbf{g}, \mathbf{s}^k)}.$$
 (10)

We parametrize the model a as a Feedforward Neural Networks that is used to compute the relevance between target location and corresponding order slot. The weight α_j is normalized by a softmax function. To write it precisely, we let

ATTW(\mathbf{s}^{j}) is a shorthand for Eq.(10), we get the upstream and downstream hidden representation by weighting sum of these slots:

$$\mathbf{z}_{down} = \sum_{i=1}^{r} \mathbf{ATTW}(\mathbf{s}_{down}^{j}) \mathbf{s}_{down}^{j}.$$
 (11)

$$\mathbf{z}_{up} = \sum_{j=1}^{r} \mathbf{ATTW}(\mathbf{s}_{up}^{j}) \mathbf{s}_{up}^{j}.$$
 (12)

Lastly, we concatenate the \mathbf{z}_{up} , \mathbf{z}_{down} and target location's hidden representation \mathbf{g} and then sent them to cost layer.

Top Layers with Multi-task Learning

The choice of cost function on the top layer is tightly coupled with the choice of the output unit. We simply use square error to fit the future conditions of the target locations.

Multi-task learning is first introduced by (Huang et al. 2014) for traffic forecasting task. It is considered as soft constraints imposed on the parameters arising out of several tasks (Evgeniou and Pontil 2004). These Additional training examples put more pressure on the parameters of the model towards values that generalize well when part of a model is shared across tasks. Forecasting traffic future condition is a multi-task problem as time goes on and different time points correspond to different tasks. In DeepTransport model, in addition to the computation of the attention weights and affine transformations of the output layer, all other parameters are shared.

Experiments

Dataset

We adopt *snowball sampling method* (Biernacki and Waldorf 1981) to collect an urban areal dataset in Beijing from a commercial map app and named it "MapBJ". The dataset provides traffic condition in {fluency, slow, congestion, extreme congestion}. The dataset contains about 349 locations which are collected from March 2016 to June for every five minutes. We select the first two months data for training and the remaining half month for testing. Besides traffic topological graph and time-varying traffic condition, we also provide the limit speed of each road. Since the limit speed of different roads may be very distinct, and locations segmentations method regards this as an important reference index. We introduce a time-invariable feature called limit level and discretize it into four classes.

Evaluation

Evaluation is ranked based on *quadratic weighted Cohen's Kappa* (Ben-David 2008), a criterion for evaluating the performance of categorical sorting.

In our problem, quadratic weighted Cohen's Kappa is characterized by three 4×4 matrices: observed matrix \mathbf{O} , expected matrix \mathbf{E} and weight matrix \mathbf{w} . Given Rater \mathbf{A} (ground truth) and Rater \mathbf{B} (prediction), $\mathbf{O}_{i,j}$ denotes the number of records rating i in \mathbf{A} while rating j in \mathbf{B} , $\mathbf{E}_{i,j}$ indicates how many samples with label i is expected to be rated as j by \mathbf{B} and $\mathbf{w}_{i,j}$ is the weight of different rating,

$$\mathbf{w}_{i,j} = \frac{(i-j)^2}{(N-1)^2},\tag{13}$$

where N is the number of subjects, we have N=4 in our problem. From these three matrices, the quadratic weighted kappa is calculated as:

$$\kappa = 1 - \frac{\sum_{i,j} \mathbf{w}_{i,j} \mathbf{O}_{i,j}}{\sum_{i,j} \mathbf{w}_{i,j} \mathbf{E}_{i,j}}.$$
 (14)

This metric typically in the range of 0 (random agreement between raters) to 1 (complete agreement between raters).

Implementation Details

Since the condition value ranges in $\{1,2,3,4\}$, the multiclassification loss can be treated as the objective function. However, cost layer with softmax cross-entropy for does not take into account the magnitude of the rating. Thus, square error loss is applied as the training objective. But another disadvantage straightforward use linear regression is that the predicted value may be out of the range in $\{1,2,3,4\}$. However, we can avoid this problem by label projection as follows:

We have a statistical analysis on the state distribution of training data. Fluency occupies 88.2% of all records, fluency and slower occupies about 96.7%, fluency, slower and congestion occupies about 99.5%, the extreme congestion is very rare that it accounts for only 0.5%. Therefore, we rank the prediction result in ascending order and set the first 88.2% to fluency, 88.2%-96.7% to slower, 96.7%-99.5% to congestion, 99.5%-100% to extreme congestion.

We put the all the observation into 32 dimension continuous vectors. The training optimization is optimized by backpropagation using Adam (Kingma and Ba 2014). Parameters are initialized with uniformly distributed random variables and we use batch size 1100 for 11 CPU threads, with each thread processes 100 records. All models are trained until convergence. Besides, there are two important hyperparameters in our model, the length of historical measurement p and perceptive radius r that controls temporal and spatial magnitude respectively.

Choosing Hyperparamerters

We intuitively suppose that expanding perceptive radius would improve prediction accuracy, but also increase the amount of computation, so it is necessary to explore the correlation between the target location and its corresponding rth order neighbors.

Mutual Infomation(MI) measures the degree of correlation between two random variables. When MI is 0, it means the given two random variables are completely irrelevant. When MI reaches the maximum value, it equals to the entropy of one of them, and the uncertainty of the other variable can be eliminated. MI is defined as

$$MI(\mathbf{X}; \mathbf{Y}) = H(\mathbf{X}) - H(\mathbf{X}|\mathbf{Y})$$

$$= \sum_{x \in \mathbf{X}, y \in \mathbf{Y}} p(x, y) \log \frac{p(x, y)}{p(x)p(y)}, (15)$$

where $H(\mathbf{X})$ and $H(\mathbf{X}|\mathbf{Y})$ are marginal entropy and conditional entropy respectively. MI describes how much uncertainty is reduced.

With MI divided by the average of entropy of the given two variables, we get Normalized mutual information(NMI) in [0, 1]:

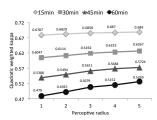
$$NMI(\mathbf{X}; \mathbf{Y}) = 2\frac{MI(\mathbf{X}, \mathbf{Y})}{H(\mathbf{X}) + H(\mathbf{Y})}.$$
 (16)

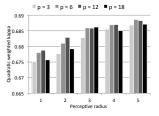
We calculated NMI between observation of each vertex and its rth neighbors over all time points. The NMI gradually decreases as the order increases, it values 0.116, 0.052, 0.038, 0.035, 0.034 for r in $\{1,2,3,4,5\}$ respectively and hardly change after r>5.

Therefore, we set the two hyperparameters as: $p \in \{3, 6, 12, 18\}$ (corresponding to 15, 30, 60, 90 minutes past measurements as 5-minutes record interval) and $r \in \{1, 2, 3, 4, 5\}$.

Effects of Hyperparameters

Figure 5 shows the averaged quadratic weighted kappa of corresponding prediction horizon. Figure 5(a) illustrates 1) closer prediction horizon always performs better; 2) As r increases, its impaction on the prediction also increases. This can be seen from the slope between r=1 and r=5, the slope at 60-min is greater than the same segment of 15-min. Figure 5(b) takes 60-min estimation as an example, indicating that the predictive effect is not monotonically increasing as the length of measurement p, and the same result can be obtained at other time points. This is because the increase in p brings an increase in the amount of parameter, which leads to overfitting.





- (a) Prediction with p = 12
- (b) 60-minute prediction

Figure 5: Averaged quadratic weighted kappa over the perceptive radius r and the length of historical measurement p on validation data. The left figure illustrates that as a function of perceptive radius r increase, the longer horizon prediction has more growth. The right figure shows that the optimal p should be chosen by observing the perceptive radius.

Comparison with Other Methods

We compared DeepTransport with four representative approaches: Random Walk(RW), Autoregressive Integrated Moving Average(ARIMA) and Stacked AutoEncoders(SAEs).

RW: In this baseline, the traffic condition at the next moment is estimated as a result of the random walk at the current moment condition that adds a white noise(a normal variable with zero mean and variance one).

ARIMA: It (Ahmed and Cook 1979) is a common statistical method for learning and predicting future values with time series data. We take a grid search over all admissible values of p, d and q which are less than p = 5, d = 2 and q = 5.

FNN: We also implemented Feed-forward Neural Networks (FNN), with a single hidden layer and an output layer with regression cost. Hidden layer has 32 neurons, and four output neurons refer to the prediction horizon. Hyperbolic tangent function and linear transfer function are used for activation function and output respectively.

SAEs: We also implemented SAEs (Lv et al. 2015), one of the most effective deep learning based methods for traffic condition forecasting. It concatenates observations of all locations to a large vector as inputs. SAEs also can be viewed as a pre-training version of FNN with large input vector that proposed by (Polson and Sokolov 2017). The stacked autoencoder is configured with four layers with [256, 256, 256, 256] hidden units for pre-train. After that, a multi-task linear regression model is trained on the top layer.

Besides, we also provides the result of DeepTransport with two configurations, with r=1, p=12 (DeepTransport-R1P12) and r=5, p=12 (DeepTransport-R5P12).

Table 1 shows the results of our model and other baselines on MapBJ. In summary, the models that use spatial information(SAEs, DeepTransport) significantly have higher performance than those that do not use(RW, ARIMA, FNN), especially in longer prediction horizon. On the other hand, SAEs is a the fully-connected form, meaning that it assumes that

Quadratic Weighted Kappa					
	0 11				
Model	15-min	30-min	45-min	60-min	Avg.
RW	0.5106	0.4474	0.3917	0.3427	0.4231
ARIMA	0.6716	0.5943	0.5389	0.4545	0.5648
FNN-P12	0.6729	0.596	0.5292	0.4689	0.5667
SAEs	0.6782	0.6157	0.5553	0.4919	0.5852
DeepTransport-R1P12	0.6787	0.6114	0.5494	0.4925	0.5841
DeepTransport-R5p12	0.6889	0.6267	0.5724	0.5259	0.6035

Table 1: Models performance comparison at various future time points.

any couple locations directly connect each other so that it neglects the topology structure of transport networks. On the contrary, DeepTransport considers traffic structure results into higher performance than these baselines, demonstrating that our proposed model has good generalization performance.



(a) Downstream attention weights (b) Upstream attention weights

Figure 6: Average attention weights alignments. It quantifies the spatial-temporal dependency relationships. The left figure is downstream alignments; it captures our intuition that as predict time increased, the attention weights shifts from low order slot to higher ones. The right figure is upstream alignments; the model pay more attention to lower orders because traffic flow in higher order is dispersed.

Slot Attention Weights

DeepTransport also can observe the influence of each slot on the target location by checking slot attention weights. Figure 6 illustrates the attention weights between prediction minutes and perceptive radius by averaging all target locations. For downstream order slots, as shown in figure 6(a), it can be seen that as predict time increased, the attention weights shifts from low order slot to higher ones. On the other side, figure 6(b) shows that upstream first order slot has more impact on target location for any future time. To capture this intuition, we utilized sandglass as a metaphor to depict the spatial-temporal dependencies of traffic flow. The flowing sand passes through the aperture of a sandglass just like traffic flow through the target location. For the downstream part, the sand is first to sink to the bottom, after a period, these accumulated sand will affect the aperture just like the cumulative congestion from the higher order to the lower order. Thus, when we predict the long-period condition of the target location, our model is more willing to refer to higher order current conditions. On the other hand, the upstream part is a little different. Higher order slots are no longer important references because traffic flow in higher order is dispersed. The target location may not be the only channel of upstream traffic flow. The nearest locations are that can directly affect the target location just like the sand gather to the aperture of the sandglass. So the future condition of target location put more attention on the lower order. Although higher order row receives less attention in the upstream module, there is still a gradual change as prediction minutes increase.

Case Study

For office workers, it might be more valuable to tell when traffic congestion comes and when the traffic condition will ease. We analyze the model performance over time in figure 7, which shows the Root Mean Square Error(RMSE) between ground truth and prediction result of RW, ARIMA, SAEs, DeepTransport-R5P12. It has two peak periods, during morning and evening rush hours. We summed up three points from this figure:

- During flat periods, especially in the early morning, there is almost no difference between models as almost all roads are fluency
- 2. Rush hours are usually used to test the effectiveness of models. When the prediction horizon is 15 minutes, DeepTransport has lower errors than other models, and the advantage of DeepTransport is more obvious when predicting the far point of time(60-minute prediction).
- After the traffic peak, it is helpful to tell when the traffic condition can be mitigated. The result just after traffic peaks shows that DeepTransport predicts better over these periods.

Related Works

There has been a long thread of statistical models based on solid mathematical foundations for traffic prediction. Such as ARIMA (Ahmed and Cook 1979) and its large variety (Kamarianakis and Vouton 2003; Kamarianakis and Prastacos 2005; Kamarianakis, Shen, and Wynter 2012) played a central role due to effectiveness and interpretability. However, the statistical methods rely on a set of constraining assumptions that may fail when dealing when complex and highly nonlinear data. (Karlaftis and Vlahogianni 2011) compare the difference and similarity between statistical methods versus neural networks in transportation research.

To our knowledge, the first deep learning approach to traffic prediction was published by (Huang et al. 2014), they used a hierarchical structure with a Deep Belief Networks(DBN) in the bottom and a (multi-task) regression layer on the top. Afterward, (Lv et al. 2015) used deep stacked autoencoders(SAEs) model for traffic prediction. A comparison (Tan et al. 2016) between SAEs and DNB for traffic flow prediction was investigated. More recently, (Polson and Sokolov 2017) concatenated all observations to a large vector as inputs and send them to Feed-forward Neural Networks(FNN) that predicted future traffic conditions at each location.

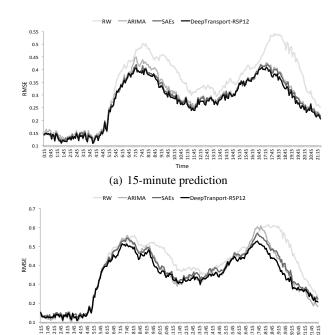


Figure 7: Model comparison with RMSE over time when prediction horizon equals 3 (15-minute) and 12 (60-minute)

(b) 60-minute prediction

On other spatial-temporal tasks, several recent deep learning works attempt to capture both time and space information. DeepST (Zhang et al. 2016) uses convolutional neural networks to predict citywide crowd flows. Meanwhile, ST-ResNet (Zhang, Zheng, and Qi 2016) uses the framework of the residual neural networks to forecast the surrounding crowds in each region through a city. These works partition a city into an $I \times J$ grid map based on the longitude and latitude (Lint, Hooqendoorn, and Zuvlen 2002) where a grid denotes a region. However, MapBJ provides the traffic networks in the form of traffic sections instead of longitude and latitude, and the road partition method should be considered the speed limit level rather than equally cut by road length. Due to the differences in data granularity, we do not follow these methods on traffic forecasting.

Conclusions

In this paper, we demonstrate the importance of using road temporal and spatial information in traffic condition forecasting. We proposed a novel deep learning model (Deep-Transport) to learn the spatial-temporal dependency. The model not only adopts two sequential models(CNN and RNN) to capture the spatial-temporal information but also take attention mechanism to quantify the spatial-temporal dependency relationships. We further released a real-world large traffic condition dataset including millions of recordings. Our experiment shows that DeepTransport significantly outperformed other previous statistical and deep learning methods for traffic forecasting.

References

- [Ahmed and Cook 1979] Ahmed, M. S., and Cook, A. R. 1979. *Analysis of freeway traffic time-series data by using Box-Jenkins techniques*. Number 722.
- [Bahdanau, Cho, and Bengio 2014] Bahdanau, D.; Cho, K.; and Bengio, Y. 2014. Neural machine translation by jointly learning to align and translate. *Computer Science*.
- [Ben-David 2008] Ben-David, A. 2008. Comparison of classification accuracy using cohen's weighted kappa. *Expert Systems with Applications* 34(2):825–832.
- [Biernacki and Waldorf 1981] Biernacki, P., and Waldorf, D. 1981. Snowball sampling: Problems and techniques of chain referral sampling. *Sociological methods & research* 10(2):141–163.
- [Chang and Su 1995] Chang, G.-L., and Su, C.-C. 1995. Predicting intersection queue with neural network models. *Transportation Research Part C: Emerging Technologies* 3(3):175–191.
- [Deng et al. 2016] Deng, D.; Shahabi, C.; Demiryurek, U.; Zhu, L.; Yu, R.; and Liu, Y. 2016. Latent space model for road networks to predict time-varying traffic. *arXiv preprint arXiv:1602.04301*.
- [Dia 2001a] Dia, H. 2001a. An object-oriented neural network approach to short-term traffic forecasting. *European Journal of Operational Research* 131(2):253–261.
- [Dia 2001b] Dia, H. 2001b. An object-oriented neural network approach to short-term traffic forecasting. *European Journal of Operational Research* 131(2):253–261.
- [Evgeniou and Pontil 2004] Evgeniou, T., and Pontil, M. 2004. Regularized multi-task learning. In *Proceedings of the tenth ACM SIGKDD international conference on Knowledge discovery and data mining*, 109–117. ACM.
- [Hochreiter and Schmidhuber 1997] Hochreiter, S., and Schmidhuber, J. 1997. Long short-term memory. *Neural computation* 9(8):1735–1780.
- [Huang et al. 2014] Huang, W.; Song, G.; Hong, H.; and Xie, K. 2014. Deep architecture for traffic flow prediction: Deep belief networks with multitask learning. *IEEE Transactions on Intelligent Transportation Systems* 15(5):2191–2201.
- [Innamaa 2000] Innamaa, S. 2000. Short-term prediction of traffic situation using mlp-neural networks. In *Proceedings* of the 7th world congress on intelligent transport systems, Turin, Italy, 6–9.
- [Jeong et al. 2013] Jeong, Y.-S.; Byon, Y.-J.; Castro-Neto, M. M.; and Easa, S. M. 2013. Supervised weighting-online learning algorithm for short-term traffic flow prediction. *IEEE Transactions on Intelligent Transportation Systems* 14(4):1700–1707.
- [Kamarianakis and Prastacos 2005] Kamarianakis, Y., and Prastacos, P. 2005. Space-time modeling of traffic flow. *Computers & Geosciences* 31(2):119–133.
- [Kamarianakis and Vouton 2003] Kamarianakis, Y., and Vouton, V. 2003. Forecasting traffic flow conditions in an urban network: Comparison of multivariate and univariate approaches. *Transportation Research Record* 1857(1):74–84.

- [Kamarianakis, Shen, and Wynter 2012] Kamarianakis, Y.; Shen, W.; and Wynter, L. 2012. Real-time road traffic forecasting using regime-switching space-time models and adaptive lasso. *Applied stochastic models in business and industry* 28(4):297–315.
- [Karlaftis and Vlahogianni 2011] Karlaftis, M. G., and Vlahogianni, E. I. 2011. Statistical methods versus neural networks in transportation research: Differences, similarities and some insights. *Transportation Research Part C: Emerging Technologies* 19(3):387–399.
- [Kingma and Ba 2014] Kingma, D., and Ba, J. 2014. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*.
- [LeCun, Bengio, and Hinton 2015] LeCun, Y.; Bengio, Y.; and Hinton, G. 2015. Deep learning. *Nature* 521(7553):436–444.
- [LeCun, Bengio, and others 1995] LeCun, Y.; Bengio, Y.; et al. 1995. Convolutional networks for images, speech, and time series. *The handbook of brain theory and neural networks* 3361(10):1995.
- [Lint, Hooqendoorn, and Zuvlen 2002] Lint, J. W. C. V.; Hooqendoorn, S. P.; and Zuvlen, H. J. V. 2002. Freeway travel time prediction with state-space neural networks: Modeling state-space dynamics with recurrent neural networks. *Transportation Research Record Journal of the Transportation Research Board* 1811(1):347369.
- [Lv et al. 2015] Lv, Y.; Duan, Y.; Kang, W.; Li, Z.; and Wang, F.-Y. 2015. Traffic flow prediction with big data: a deep learning approach. *IEEE Transactions on Intelligent Transportation Systems* 16(2):865–873.
- [Pan, Demiryurek, and Shahabi 2012] Pan, B.; Demiryurek, U.; and Shahabi, C. 2012. Utilizing real-world transportation data for accurate traffic prediction. In *Data Mining (ICDM)*, 2012 IEEE 12th International Conference on, 595–604. IEEE.
- [Polson and Sokolov 2017] Polson, N. G., and Sokolov, V. O. 2017. Deep learning for short-term traffic flow prediction. *Transportation Research Part C Emerging Technologies* 79:1–17.
- [Rocktschel et al. 2015] Rocktschel, T.; Grefenstette, E.; Hermann, K. M.; Koisk, T.; and Blunsom, P. 2015. Reasoning about entailment with neural attention.
- [Stathopoulos and Karlaftis 2003] Stathopoulos, A., and Karlaftis, M. G. 2003. A multivariate state space approach for urban traffic flow modeling and prediction. *Transportation Research Part C: Emerging Technologies* 11(2):121–135.
- [Tan et al. 2016] Tan, H.; Xuan, X.; Wu, Y.; Zhong, Z.; and Ran, B. 2016. A comparison of traffic flow prediction methods based on dbn. In *CICTP 2016*. 273–283.
- [Vlahogianni, Karlaftis, and Golias 2005] Vlahogianni, E. I.; Karlaftis, M. G.; and Golias, J. C. 2005. Optimized and meta-optimized neural networks for short-term traffic flow prediction: A genetic approach. *Transportation Research Part C: Emerging Technologies* 13(3):211–234.

- [Williams and Hoel 1999] Williams, B. M., and Hoel, L. A. 1999. Modeling and forecasting vehicular traffic flow as a seasonal stochastic time series process. Technical report.
- [Williams 2001] Williams, B. 2001. Multivariate vehicular traffic flow prediction: Evaluation of arimax modeling. *Transportation Research Record: Journal of the Transportation Research Board*.
- [Zhang et al. 2016] Zhang, J.; Zheng, Y.; Qi, D.; Li, R.; and Yi, X. 2016. Dnn-based prediction model for spatio-temporal data. In *Proceedings of the 24th ACM SIGSPA-TIAL International Conference on Advances in Geographic Information Systems*, 92. ACM.
- [Zhang, Zheng, and Qi 2016] Zhang, J.; Zheng, Y.; and Qi, D. 2016. Deep spatio-temporal residual networks for citywide crowd flows prediction. *arXiv preprint arXiv:1610.00081*.