

A BAYESIAN MODEL FOR DATA FLOW: BIKEMi

Andrea De Gobbis, Lorenzo Ghilotti, Giorgio Meretti

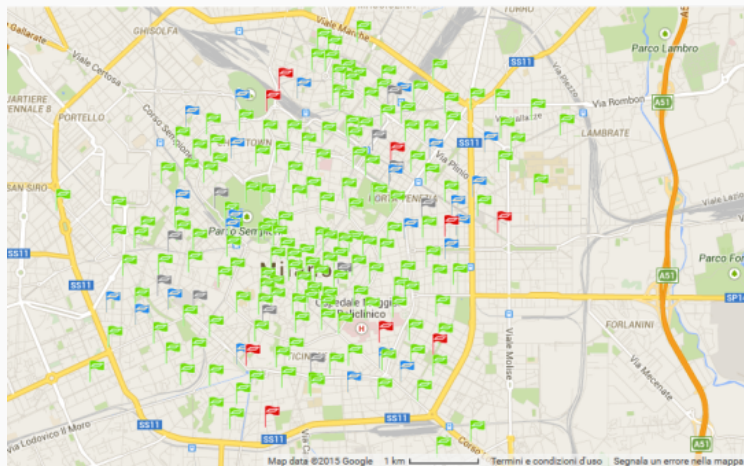
January 8, 2020

Politecnico di Milano

What we are doing

The BikeMi stations net in Milan

Andrea De Gobbis,
Lorenzo Ghilotti,
Giorgio Meretti



Two prospective of the problem

Andrea De Gobbis,
Lorenzo Ghilotti,
Giorgio Meretti

We followed two distinct paths:

Global model: the total volume of bikes travels in a specific day Y_t without considering the graph structure. This results in a single time series.

Network model: dividing in the different nodes and analysing the flow of bikes in the net. The dimensionality is much higher.

Day by day **Poisson**:

$$\left\{ \begin{array}{l} Y_t \sim \text{Po}(Y_t | \lambda_t) \\ \lambda_t = \exp\{\alpha + \boldsymbol{\beta} \cdot \mathbf{x}_t\} \\ \alpha \sim \mathcal{N}(0, \sigma_\alpha^2) \\ \beta_i \stackrel{iid}{\sim} \mathcal{N}(0, \sigma_\beta^2) \end{array} \right. \quad (1)$$

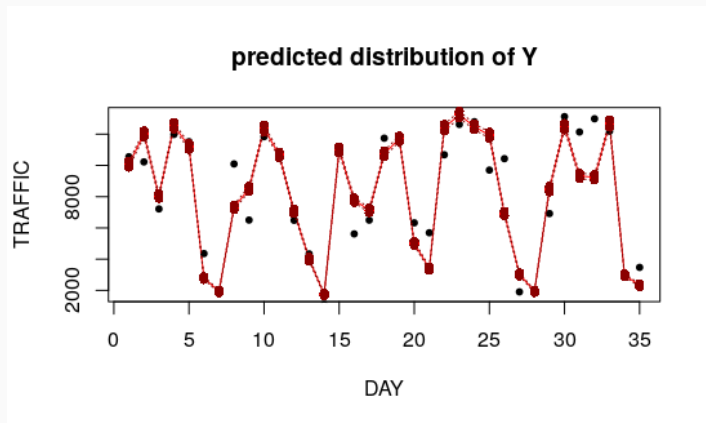
With **covariates** \mathbf{x}_t :

- Y_{t-1} volume on the previous day
- Y_{t-7} volume on the same weekday of the previous week
- W_t dummy for weekday / weekend
- R_t, R_{t-1} dummies for rain in the current and previous day
- T_t mean temperature for the day
- S_t, M_t dummies for Saturday and Monday

Predictive distribution of the poisson model

Andrea De Gobbis,
Lorenzo Ghilotti,
Giorgio Meretti

Only 3 of 35 in the 90% credible interval



$$\begin{cases} \mathbf{X}(t) = \mathbf{f}(\mathbf{X}(1 : (t - 1))) + \boldsymbol{\epsilon}_1(t) \\ \mathbf{Y}(t) = \mathbf{g}(\mathbf{X}(t)) + \boldsymbol{\epsilon}_2(t) \end{cases} \quad (2)$$

with suitable initial conditions and priors

Elementary model with precisions

$$\left\{ \begin{array}{l} Y_t = \mu_t + \gamma_t + \frac{1}{\sqrt{\tau_\epsilon}} \tilde{\epsilon}_t \\ \mu_t = \mu_{t-1} + \delta_{t-1} + \frac{1}{\sqrt{\tau_\eta}} \tilde{\eta}_t \\ \delta_t = \delta_{t-1} + \frac{1}{\sqrt{\tau_v}} \tilde{v}_t \\ \gamma_t = \sum_{i=1}^{S-1} \gamma_{t+i-S} + \frac{1}{\sqrt{\tau_w}} \tilde{w}_t \\ \tilde{\epsilon}_t, \tilde{\eta}_t, \tilde{v}_t, \tilde{w}_t, \stackrel{iid}{\sim} \mathcal{N}(0, 1) \end{array} \right. \quad (3)$$

Elementary model priors

$$\left\{ \begin{array}{l} \mu_0 \sim \mathcal{N}(m, \tau_m) \\ \delta_0 \sim \mathcal{N}(d, \tau_d) \\ \gamma_{0:(2-S)} \sim \mathcal{N}_{S-1}(\mathbf{g}, \tau_g \mathbf{I}) \\ \tau_* \sim \text{Gamma}(a_*, b_*), \text{ with } * = \{\epsilon, \eta, v, w\} \\ \{\tilde{\epsilon}_t, \tilde{\eta}_t, \tilde{v}_t, \tilde{w}_t, \mu_0, \delta_0, \gamma_{0:(-S+2)}, \tau_\epsilon, \tau_\eta, \tau_v, \tau_w\} \\ \text{independent.} \end{array} \right. \quad (4)$$

Complete model with precisions

$$\left\{ \begin{array}{l} Y_t = \mu_t + \gamma_t + \rho_t + \beta^T \mathbf{z}_t + \frac{1}{\sqrt{\tau_\epsilon}} \tilde{\epsilon}_t \\ \mu_t = \mu_{t-1} + \delta_{t-1} + \frac{1}{\sqrt{\tau_\eta}} \tilde{\eta}_t \\ \delta_t = \delta_{t-1} + \frac{1}{\sqrt{\tau_v}} \tilde{v}_t \\ \gamma_t = \sum_{i=1}^{S-1} \gamma_{t+i-S} + \frac{1}{\sqrt{\tau_w}} \tilde{w}_t \\ \rho_t = \alpha \rho_{t-1} + \frac{1}{\sqrt{\tau_u}} \tilde{u}_t \\ \tilde{\epsilon}_t, \tilde{\eta}_t, \tilde{v}_t, \tilde{w}_t, \tilde{u}_t \stackrel{iid}{\sim} \mathcal{N}(0, 1) \end{array} \right. \quad (5)$$

Complete model priors

$$\left\{ \begin{array}{l} \mu_0 \sim \mathcal{N}(m, \tau_m) \\ \delta_0 \sim \mathcal{N}(d, \tau_d) \\ \gamma_{0:(2-S)} \sim \mathcal{N}_{S-1}(\mathbf{g}, \tau_g \mathbf{I}) \\ \rho_0 \sim \mathcal{N}(r, \tau_r) \\ \beta \sim \mathcal{N}_p(\mathbf{0}, \tau_b \mathbf{I}) \\ \alpha \sim \mathcal{N}(a, \tau_a) \\ \tau_* \sim \text{Gamma}(a_*, b_*), \text{ with } * = \{\epsilon, \eta, v, w, u\} \\ \{\tilde{\epsilon}_t, \tilde{\eta}_t, \tilde{u}_t, \tilde{w}_t, \tilde{v}_t, \mu_0, \delta_0, \gamma_{0:(-S+2)}, \rho_0, \beta, \alpha, \tau_\epsilon, \tau_\eta, \tau_v, \tau_w, \tau_u\} \\ \text{independent.} \end{array} \right. \quad (6)$$

Robust model with precisions

Andrea De Gobbis,
Lorenzo Ghilotti,
Giorgio Meretti

$$\left\{ \begin{array}{l} Y_t = \mu_t + \gamma_t + \rho_t + \beta^T \mathbf{z}_t + \frac{1}{\sqrt{\tau_\epsilon}} \tilde{\epsilon}_t \\ \mu_t = \text{avgpred}(\mu_t) + \text{avgpred}(\delta_t) + \frac{1}{\sqrt{\tau_\eta}} \tilde{\eta}_t \\ \delta_t = \text{avgpred}(\delta_t) + \frac{1}{\sqrt{\tau_v}} \tilde{v}_t \\ \gamma_t = \sum_{i=1}^{S-1} \gamma_{t+i-S} + \frac{1}{\sqrt{\tau_w}} \tilde{w}_t \\ \rho_t = \alpha \rho_{t-1} + \frac{1}{\sqrt{\tau_u}} \tilde{u}_t \\ \tilde{\epsilon}_t, \tilde{\eta}_t, \tilde{v}_t, \tilde{w}_t, \tilde{u}_t \stackrel{iid}{\sim} \mathcal{N}(0, 1) \end{array} \right. \quad (7)$$

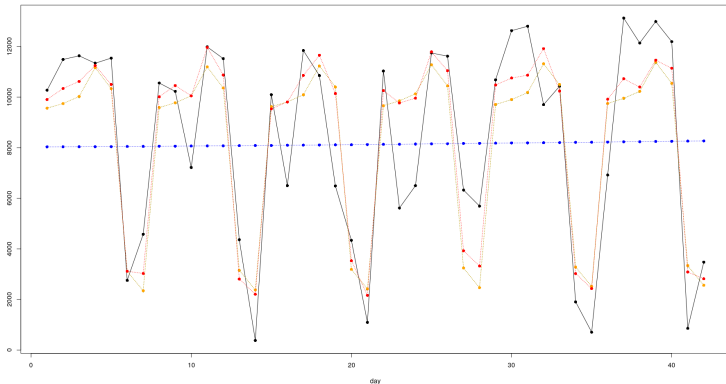
where $\text{avgpred}(\phi_t) = \frac{1}{2}\phi_{t-1} + \frac{1}{3}\phi_{t-S} + \frac{1}{6}\phi_{t-2S}$
and same priors as before

Time series model: BSTS

Instability of the errors, τ_ϵ explodes, bad autocorrelation

Andrea De Gobbis,
Lorenzo Ghilotti,
Giorgio Meretti

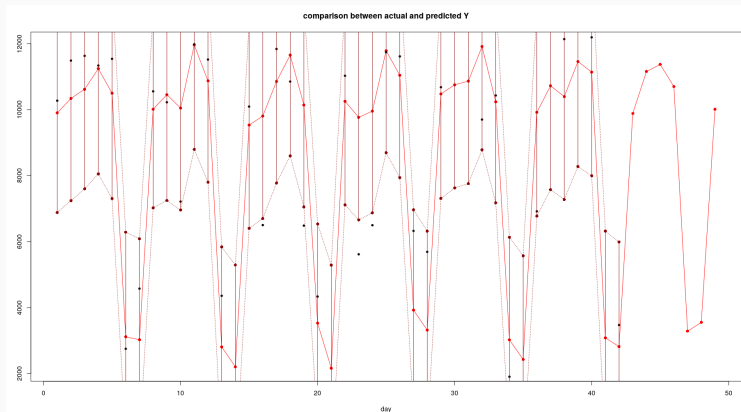
comparison between actual and predicted traffic



Time series model: BSTS

Instability of the errors, τ_ϵ explodes, bad autocorrelation

Andrea De Gobbis,
Lorenzo Ghilotti,
Giorgio Meretti

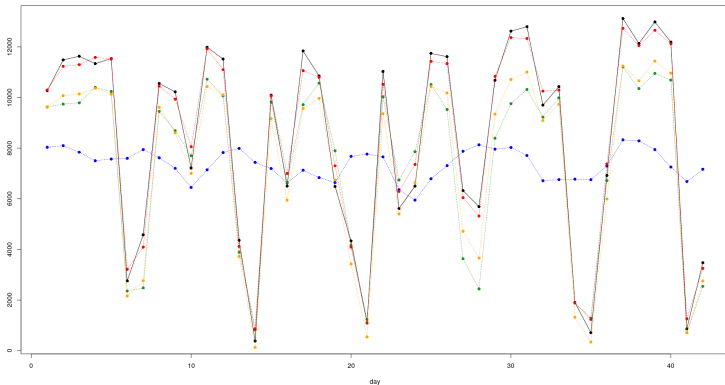


Time series model: BSTS

Blocking the maximum variance, τ_* under control,
partial mixing of state variables

Andrea De Gobbis,
Lorenzo Ghilotti,
Giorgio Meretti

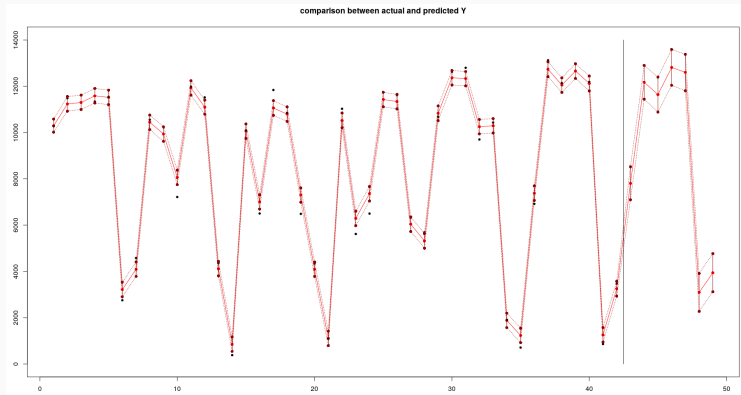
comparison between actual and predicted traffic



Time series model: BSTS

Smaller induced variability, **attempt of prediction with real weather.**

Andrea De Gobbis,
Lorenzo Ghilotti,
Giorgio Meretti



Comparison truncated-variance robust BSTS vs Poisson.

Observations inside **credibility intervals**:

- Poisson 3/35
- BSTS 24/42

WAIC:

- Poisson -0.0171
- BSTS -0.0160

Robust model for time zones with precisions

Andrea De Gobbis,
Lorenzo Ghilotti,
Giorgio Meretti

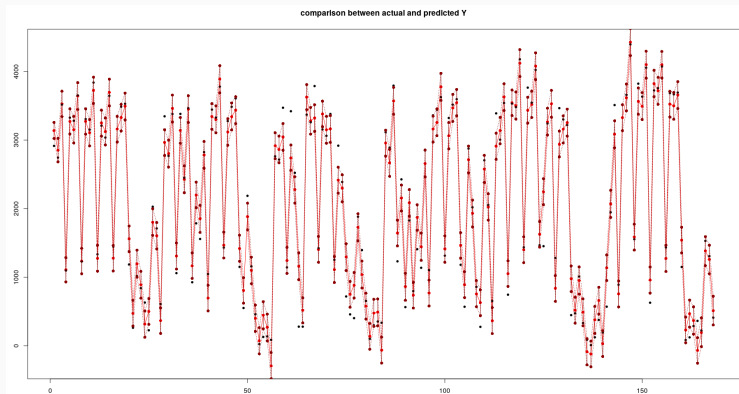
$$\left\{ \begin{array}{l} Y_{th} = \mu_t + \gamma_t + \rho_t + \chi_k + \beta^T \mathbf{z}_t + \frac{1}{\sqrt{\tau_\epsilon}} \tilde{\epsilon}_{th} \\ \mu_t = \text{avgpred}(\mu_t) + \text{avgpred}(\delta_t) + \frac{1}{\sqrt{\tau_\eta}} \tilde{\eta}_t \\ \delta_t = \text{avgpred}(\delta_t) + \frac{1}{\sqrt{\tau_v}} \tilde{v}_t \\ \gamma_t = \sum_{i=1}^{S-1} \gamma_{t+i-S} + \frac{1}{\sqrt{\tau_w}} \tilde{w}_t \\ \rho_t = \alpha \rho_{t-1} + \frac{1}{\sqrt{\tau_u}} \tilde{u}_t \\ \chi_k = \sum_{i=1}^{F-1} \chi_{k+i-F} + \xi \delta_{t(k)=6,7} + \frac{1}{\sqrt{\tau_\zeta}} \tilde{\zeta} \\ \tilde{\epsilon}_t, \tilde{\eta}_t, \tilde{v}_t, \tilde{w}_t, \tilde{u}_t, \tilde{\zeta} \stackrel{iid}{\sim} \mathcal{N}(0, 1) \end{array} \right. \quad (8)$$

where $\text{avgpred}(\phi_t) = \frac{1}{2}\phi_{t-1} + \frac{1}{3}\phi_{t-S} + \frac{1}{6}\phi_{t-2S}$, $h = \text{mod}(k, 4) + 1$ and **priors of the same class as before**

Time series model: BSTS

Robust truncated variance model for time zones

Andrea De Gobbis,
Lorenzo Ghilotti,
Giorgio Meretti



Two prospective of the problem

Andrea De Gobbis,
Lorenzo Ghilotti,
Giorgio Meretti

We followed two distinct paths:

Global model: the total volume of bikes travels in a specific day Y_t without considering the graph structure. This results in a single time series.

Network model: dividing in the different nodes and analysing the flow of bikes in the net. The dimensionality is much higher.

For every (i, j) edge of the graph and $t \in 1 : 42$ we have $Y_{ij}(t)$ the number of travels from node i to j at day t .

Problem:

more than 4 million variables \Rightarrow
Computationally untreatable

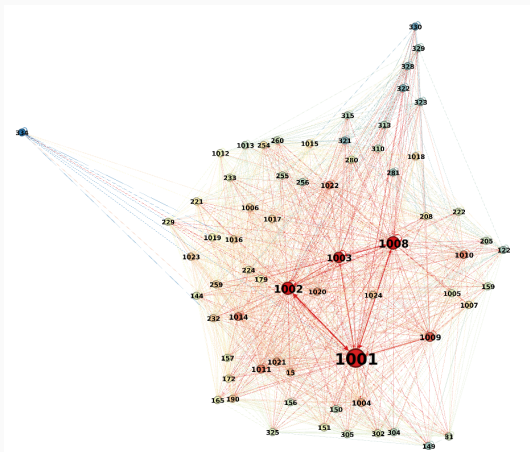
Solutions:

- clusterization through DBSCAN
- simplification of the variables

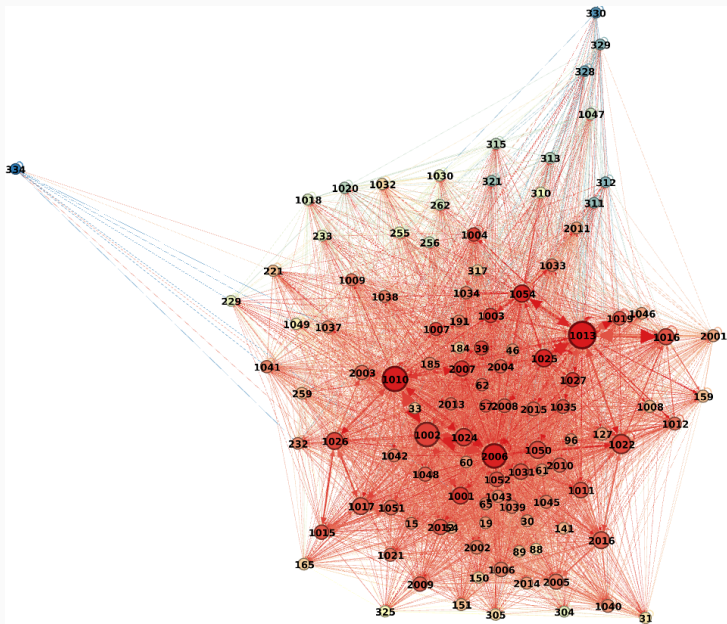
Preprocessing clusterization with DBSCAN

Algorithm to divide the nodes into initial clusters. We introduced a modified version with a weight to break up the **bigger clusters**, minimizing the **autorings** presence.

Andrea De Gobbis,
Lorenzo Ghilotti,
Giorgio Meretti



From 334 nodes to 140.



A BAYESIAN MODEL
FOR DATA FLOW:
BIKEMI

Andrea De Gobbis,
Lorenzo Ghilotti,
Giorgio Meretti

Modeling the flux

Andrea De Gobbis,
Lorenzo Ghilotti,
Giorgio Meretti

To further simplify the analysis we focus on the bikes arriving and departing from each node $N_i^{IN}(\Delta t)$ and $N_i^{OUT}(\Delta t)$ in the time interval Δt .

We are interested in the smallest Δt as possible but this would increase the number of variables \Rightarrow consider them as **functional data**.

$$V_i(t) = \lim_{\Delta t \rightarrow 0} \frac{N_i^{IN}(\Delta t) - N_i^{OUT}(\Delta t)}{\Delta t}$$

$$\Phi_i(t) = \int_0^t V_i(u) \, du$$

We can analyse when new bikes should be brought to which station.

Net time-Interval flow

