

Supplementary Material:

Prediction of Spatial Point Processes:

Regularized Method with Out-of-Sample Guarantees

Spatial basis $\phi(r)$

For space \mathcal{X} divided into R regions with each region \mathcal{X}_r denoted by its region index r , the spatial basis vector evaluated at r is an $R \times 1$ vector given by

$$\phi(r) = \text{col}\{\phi_1(r), \dots, \phi_R(r)\}.$$

Here $\phi(r)$ is the cubic b-spline in space with two parameters: center and support. The k^{th} component i.e $\phi_k(r)$ has its center at the center of region \mathcal{X}_k and peak value when evaluated at k . The support of $\phi_k(r)$ determines its value at the neighbouring regions and hence allows to control the local structure in intensity in our model. For details on cubic b-spline see [2].

Dual Norm $\tilde{f}(\cdot)$

Let $f(\theta) = \|\mathbf{w} \odot \theta\|_1$. By definition of dual norm,

$$\tilde{f}(\mathbf{g}) = \sup_{\theta: f(\theta) \leq 1} \mathbf{g}^\top \theta$$

The condition $f(\theta) \leq 1$ implies

$$\sum_{k=1}^R |w_k| |\theta_k| \leq 1, \quad \min_{k=1, \dots, R} |w_k| \sum_{k=1}^R |\theta_k| \leq 1, \quad \|\theta\|_1 \leq w_o^{-1},$$

where $w_o = \min_{k=1, \dots, R} |w_k|$. Moreover,

$$\mathbf{g}^\top \theta = \sum_{k=1}^R g_k \theta_k \leq \sum_{k=1}^R |g_k| |\theta_k| \leq \|\mathbf{g}\|_\infty \|\theta\|_1$$

Combining this with $\|\theta\|_1 \leq w_o^{-1}$ we get

$$\tilde{f}(\mathbf{g}) = \frac{\|\mathbf{g}\|_\infty}{w_o}$$

Hoeffding's inequality for z_k

We show that z_k is bounded in $[-Y, Y]$ and hence we can make use of Hoeffding's inequality to get eq. (14).

The gradient of eq. (10) evaluated at $\hat{\theta}$ is

$$\mathbf{g} = \left[\mathbb{E}_n[z_1], \dots, \mathbb{E}_n[z_R] \right]^\top,$$

where $z_k = (y - \mathbb{E}_{y|r}[y])\phi_k(r)$. Given that the maximum number of counts is bounded i.e. $\max y \leq Y$, we have

$$\begin{aligned}\max z_k &= \max \{(y - \mathbb{E}_{y|r}[y])\phi_k(r)\} = \max\{(y - \mathbb{E}_{y|r}[y])\} \max\{\phi_k(r)\} = Y, \\ \min z_k &= \min \{(y - \mathbb{E}_{y|r}[y])\phi_k(r)\} = \min\{(y - \mathbb{E}_{y|r}[y])\} \max\{\phi_k(r)\} = -Y,\end{aligned}$$

for all $k = 1, \dots, R$. Here $\max \phi_k(r) = 1$.

Union bound and DeMorgan's Law

Given that $\mathbb{E}[z_k] = 0$, from eq. (14) we get

$$\Pr(|\mathbb{E}_n[z_k]| \leq \epsilon) \geq 1 - 2 \exp \left[- \frac{n\epsilon^2}{2Y^2} \right].$$

Moreover,

$$\Pr \left(\max_{k=1, \dots, R} |\mathbb{E}_n[z_k]| \leq \epsilon \right) = \Pr \left(\bigcap_{k=1}^R |\mathbb{E}_n[z_k]| \leq \epsilon \right).$$

By DeMorgan's law,

$$\Pr \left(\bigcap_{k=1}^R |\mathbb{E}_n[z_k]| \leq \epsilon \right) = \Pr \left(\bigcup_{k=1}^R |\mathbb{E}_n[z_k]| \geq \epsilon \right)'.$$

By union bound,

$$\Pr \left(\bigcup_{k=1}^R |\mathbb{E}_n[z_k]| \geq \epsilon \right) \leq \sum_{i=1}^R \Pr \left(|\mathbb{E}_n[z_k]| \geq \epsilon \right) = 2R \exp \left[- \frac{n\epsilon^2}{2Y^2} \right],$$

which implies that

$$\Pr \left(\bigcap_{k=1}^R |\mathbb{E}_n[z_k]| \leq \epsilon \right) \geq 1 - 2R \exp \left[- \frac{n\epsilon^2}{2Y^2} \right].$$

Eq. (15) follows from above.

Minimization Algorithm

Here we derive the majorization-minimization (MM) algorithm that is used to solve eq. (7). Let $V(\boldsymbol{\theta}) = -n^{-1} \ln p_{\boldsymbol{\theta}}(\mathbf{y}|\mathbf{r})$ and $f(\boldsymbol{\theta}) = \|\mathbf{w} \odot \boldsymbol{\theta}\|_1$. For the Poisson model class considered in the paper,

$$V(\boldsymbol{\theta}) = n^{-1} \left(\sum_{i=1}^n \mathbb{E}_{\boldsymbol{\theta}}[y_i|r_i] - y_i \ln(\mathbb{E}_{\boldsymbol{\theta}}[y_i|r_i]) + \ln(y_i!) \right),$$

where $\mathbb{E}_{\boldsymbol{\theta}}[y_i|r_i] = \exp(\phi(r_i)^\top \boldsymbol{\theta})$. $V(\boldsymbol{\theta})$ is convex in $\boldsymbol{\theta}$ since

$$\partial_{\boldsymbol{\theta}}^2 V(\boldsymbol{\theta}) = n^{-1} \boldsymbol{\Phi} \mathbf{D} \boldsymbol{\Phi}^\top \geq 0.$$

Here,

$$\mathbf{\Phi} = [\phi(r_1), \dots, \phi(r_R)] \text{ and } \mathbf{D} = \text{diag}(\mathbf{h}(\boldsymbol{\theta}))$$

are an $R \times R$ basis and diagonal matrices respectively and

$$\mathbf{h}(\boldsymbol{\theta}) = \text{col}\{\mathbb{E}_{\boldsymbol{\theta}}[y_1|r_1], \dots, \mathbb{E}_{\boldsymbol{\theta}}[y_R|r_R]\}.$$

By convexity of $V(\boldsymbol{\theta})$, given an initial estimate $\tilde{\boldsymbol{\theta}}$, the objective in eq. (7) can be upper bounded as

$$V(\boldsymbol{\theta}) + n^{-\gamma} f(\boldsymbol{\theta}) \leq Q(\boldsymbol{\theta}; \tilde{\boldsymbol{\theta}}) + n^{-\gamma} f(\boldsymbol{\theta}),$$

where $Q(\boldsymbol{\theta}; \tilde{\boldsymbol{\theta}})$ is a quadratic majorization function (see [1], ch. 5) of $V(\boldsymbol{\theta})$ given by

$$Q(\boldsymbol{\theta}; \tilde{\boldsymbol{\theta}}) = V(\tilde{\boldsymbol{\theta}}) + \mathbf{v}^\top (\boldsymbol{\theta} - \tilde{\boldsymbol{\theta}}) + \frac{1}{2} (\boldsymbol{\theta} - \tilde{\boldsymbol{\theta}})^\top \mathbf{H} (\boldsymbol{\theta} - \tilde{\boldsymbol{\theta}}).$$

Here $\mathbf{v} = \partial_{\boldsymbol{\theta}} V(\boldsymbol{\theta})|_{\boldsymbol{\theta} = \tilde{\boldsymbol{\theta}}} = n^{-1} \mathbf{\Phi}(\mathbf{h}(\hat{\boldsymbol{\theta}}) - \mathbf{y})$ and $\mathbf{H} = \max_{\boldsymbol{\theta}} \{\partial_{\boldsymbol{\theta}}^2 V(\boldsymbol{\theta})\}$. The diagonal elements of \mathbf{D} represent the average number of counts in different regions. Given that the counts in any region are bounded i.e. $y \leq Y$, $\mathbf{H} \leq n^{-1} Y \mathbf{\Phi} \mathbf{\Phi}^\top$ therefore we have

$$V(\boldsymbol{\theta}) + n^{-\gamma} f(\boldsymbol{\theta}) \leq V(\hat{\boldsymbol{\theta}}) + n^{-1} (\mathbf{h}(\hat{\boldsymbol{\theta}}) - \mathbf{y})^\top \mathbf{\Phi}^\top (\boldsymbol{\theta} - \tilde{\boldsymbol{\theta}}) + \frac{Y}{2n} \|\mathbf{\Phi}^\top (\boldsymbol{\theta} - \tilde{\boldsymbol{\theta}})\|_2^2 + n^{-\gamma} f(\boldsymbol{\theta}). \quad (\text{SM1})$$

Therefore starting from an initial estimate $\tilde{\boldsymbol{\theta}}$, one can minimize the right hand side of (SM1) to obtain $\check{\boldsymbol{\theta}}$ then update $\tilde{\boldsymbol{\theta}} = \check{\boldsymbol{\theta}}$ and repeat until convergence to get final solution of eq. (7) $\hat{\boldsymbol{\theta}} = \check{\boldsymbol{\theta}}$. The pseudocode is given in algorithm (2).

Furthermore, the right hand side of (SM1) can be transformed into a weighted lasso regression problem and hence can be efficiently solved using coordinate descent algorithm [1]. Letting $\mathbf{q}(\tilde{\boldsymbol{\theta}}) = \mathbf{\Phi}^\top \tilde{\boldsymbol{\theta}} + Y(\mathbf{y} - \mathbf{h}(\tilde{\boldsymbol{\theta}}))$, the right hand side of (SM1) can be rewritten as

$$Y(2n)^{-1} (\mathbf{q}(\tilde{\boldsymbol{\theta}}) - \mathbf{\Phi}^\top \boldsymbol{\theta})^\top (\mathbf{q}(\tilde{\boldsymbol{\theta}}) - \mathbf{\Phi}^\top \boldsymbol{\theta}) + n^{-\gamma} f(\boldsymbol{\theta}) + K(\tilde{\boldsymbol{\theta}}),$$

where the first two terms form a weighted lasso regression problem in $\boldsymbol{\theta}$ and the last term $K(\tilde{\boldsymbol{\theta}}) = V(\tilde{\boldsymbol{\theta}}) - \mathbf{q}(\tilde{\boldsymbol{\theta}})^\top \tilde{\boldsymbol{\theta}}$ is independent of $\boldsymbol{\theta}$ and does not affect the minimization problem. This conclude the derivation of the MM algorithm.

References

- [1] R. Tibshirani, M. Wainwright, and T. Hastie. *Statistical learning with sparsity: the lasso and generalizations*. Chapman and Hall/CRC, 2015.
- [2] L. Wasserman. *All of nonparametric statistics*. Springer Science & Business Media, 2006.