

Fast spike train inference from calcium imaging

Joshua T. Vogelstein, other fuckers, Liam Paninski

December 24, 2008

Abstract

1 Introduction

2 Methods

2.1 Model

We assume the following discrete-time state-space model relating spikes, n_t , baseline subtracted intracellular calcium concentration, denoted by C_t , and fluorescence measurements, F_t : (1) spikes follow Poisson statistics with rate $\lambda\Delta$, (2) C_t decays exponentially with decay γ to baseline ν , but jumps by ρ after each spike, and (3) fluorescent observations are linear functions of C_t with additive Gaussian noise with variance σ^2 . Together, these assumptions imply the following model:

$$F_t = \alpha C_t + \beta + \sigma \varepsilon_t, \quad \varepsilon_t \sim \mathcal{N}(0, 1) \quad (1)$$

$$C_t = \gamma C_{t-1} + \nu + \rho n_t, \quad n_t \sim \text{Poisson}(n_t; \lambda\Delta) \quad (2)$$

where α and β set the baseline and offset respectively. To enforce identifiability, we will let $\alpha = 1$ and $\beta = 0$, without loss of generality, as explained below.

2.2 Inference

Given such a model, our goal is to find the maximum *a posteriori* (MAP) spike train, i.e., the most likely spike train, $\mathbf{n} = [n_1, \dots, n_T]$, where T is the final time step, given the fluorescence measurements, $\mathbf{F} = [F_1, \dots, F_T]$. Thus, the objective function that we'd like to solve may be written as:

$$\begin{aligned}
\hat{\mathbf{n}} &= \underset{n_t \in \mathbb{N}_0 \forall t}{\operatorname{argmin}} P(\mathbf{n}|\mathbf{F}) = \underset{n_t \in \mathbb{N}_0 \forall t}{\operatorname{argmin}} P(\mathbf{F}|\mathbf{n})P(\mathbf{n}) \\
&= \underset{n_t \in \mathbb{N}_0 \forall t}{\operatorname{argmax}} \prod_{t=1}^T P(F_t|C_t)P(n_t) = \underset{n_t \in \mathbb{N}_0 \forall t}{\operatorname{argmax}} \sum_{t=1}^T (\log P(F_t|C_t) + \log P(n_t)) \\
&= \underset{n_t \in \mathbb{N}_0 \forall t}{\operatorname{argmin}} \sum_{t=1}^T \left(\frac{1}{2\sigma^2} (F_t - C_t)^2 - n_t \log \lambda \Delta + \log n_t! \right), \tag{3}
\end{aligned}$$

where $\mathbb{N}_0 = \{0, 1, 2, \dots\}$, and C_t is implicitly a function of n_t . Unfortunately, the computational complexity of solving Eq. (3) scales exponentially with T , i.e., is $O(e^T)$, making Eq. (3) an NP-hard problem. Thus, instead of an exact solution, we propose to make an approximation that reduces the complexity to be *polynomial* in T . In particular, we relax the assumption that we must have an integer number of spikes at any time step, by approximating the Poisson distribution with an exponential. Note that this is a common approximation technique in the machine learning literature [Hastie et al., 2001], as it is the closest convex relaxation to its non-convex counterpart. The constraint on n_t in Eq. (3) is therefore relaxed from $n_t \in \mathbb{N}_0$ to $n_t \geq 0$:

$$\hat{\mathbf{n}} \approx \underset{n_t \geq 0 \forall t}{\operatorname{argmin}} \sum_{t=1}^T \left(\frac{1}{2\sigma^2} (F_t - C_t)^2 + n_t \lambda \Delta \right). \tag{4}$$

While this convex relaxation makes the problem tractable, the “sharp” threshold imposed by the nonnegativity constraint prohibits the use of standard gradient ascent techniques [Boyd and Vandenberghe, 2004]. We therefore take an “interior-point” (or “barrier”) approach, in which we drop the sharp threshold, and add a barrier term, which must approach $-\infty$ as n_t approaches zero (e.g., $-\log n_t$) [Boyd and Vandenberghe, 2004]. By iteratively reducing the weight of the barrier term, $z > 0$, we are guaranteed to converge to the correct solution [Boyd and Vandenberghe, 2004], $\hat{\mathbf{n}}$. Thus, our goal is to efficiently solve:

$$\hat{\mathbf{n}}_z = \underset{n_t \geq 0 \forall t}{\operatorname{argmin}} \sum_{t=1}^T \left(\frac{1}{2\sigma^2} (F_t - C_t)^2 + n_t \lambda \Delta - z \log(n_t) \right), \tag{5}$$

which is concave and may therefore be solved exactly and efficiently using a gradient ascent technique. To see how, we first note that spikes and calcium are related to one another via a simple linear transformation, namely, $n_t = f(C_t, C_{t-1}) = (C_t - \gamma C_{t-1} - \nu)/\rho$, or, in matrix notation:

$$\mathbf{MC} - \mathbf{b} = \begin{bmatrix} 1/\rho & 0 & 0 & \dots & \dots \\ 1/\rho & -\gamma/\rho & 0 & \dots & \dots \\ 0 & 1/\rho & -\gamma/\rho & 0 & \dots \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ 0 & 0 & 0 & 1/\rho & -\gamma/\rho \end{bmatrix} \begin{bmatrix} C_1 \\ C_2 \\ \vdots \\ \vdots \\ C_T \end{bmatrix} - \begin{bmatrix} \nu/\rho \\ \nu/\rho \\ \vdots \\ \vdots \\ \nu/\rho \end{bmatrix} = \begin{bmatrix} n_1 \\ n_2 \\ \vdots \\ \vdots \\ n_T \end{bmatrix} = \mathbf{n} \tag{6}$$

where $\mathbf{M} \in \mathbb{R}^{T \times T}$ is a bidiagonal matrix, and \mathbf{b} , $\mathbf{C} = [C_1, \dots, C_T]$, and \mathbf{n} are T dimensional column vectors. Given Eq. (6), we may rewrite Eq. (5) in terms of \mathbf{C} :

$$\begin{aligned} \hat{\mathbf{C}}_z &= \underset{(C_t - \gamma C_{t-1} - \nu)/\rho \geq 0 \forall t}{\operatorname{argmin}} \sum_{t=1}^T \left(\frac{1}{2\sigma^2} (F_t - C_t)^2 + (C_t - \gamma C_{t-1} - \nu) \lambda \Delta / \rho - z \log[(C_t - \gamma C_{t-1} - \nu) / \rho] \right) \\ &= \underset{\mathbf{MC} - \mathbf{b} \geq \mathbf{0}}{\operatorname{argmin}} \frac{1}{2\sigma^2} \|\mathbf{F} - \mathbf{C}\|^2 + \lambda \Delta (\mathbf{MC} - \mathbf{b})' \mathbf{1} - z \log(\mathbf{MC} - \mathbf{b})' \mathbf{1}, \end{aligned} \quad (7)$$

where $\mathbf{MC} - \mathbf{b} \geq \mathbf{0}$ indicates that every element of $\mathbf{MC} - \mathbf{b}$ is greater than or equal to zero, $'$ indicates transpose, and $\mathbf{1}$ is a T dimensional column vectors. As (7) is concave (it is identical to Eq. (5) with different notation), we may use any descent technique to find $\hat{\mathbf{C}}_z$. We elect to use the Newton-Raphson approach:

$$\hat{\mathbf{C}}_z \leftarrow \hat{\mathbf{C}}_z + s\mathbf{d} \quad (8a)$$

$$\mathbf{H}\mathbf{d} = \mathbf{g} \quad (8b)$$

$$\mathbf{g} = -\frac{1}{\sigma^2} (\mathbf{F} - \hat{\mathbf{C}}_z) + \lambda \Delta \mathbf{M}' \mathbf{1} - z \mathbf{M}' (\mathbf{M} \hat{\mathbf{C}}_z - \mathbf{b})^{-1} \quad (8c)$$

$$\mathbf{H} = \frac{1}{\sigma^2} \mathbf{I} + 2z \mathbf{M}' (\mathbf{M} \hat{\mathbf{C}}_z - \mathbf{b})^{-2} \mathbf{M} \quad (8d)$$

where s is the step size, \mathbf{d} is the step direction, and \mathbf{g} and \mathbf{H} are the gradient (first derivative) and Hessian (second derivative) of the likelihood, respectively (the likelihood is the argument to be minimized in Eq. (7)). Note that we use “backtracking line searches”, meaning that for each iteration, we find the maximal s that is between 0 and 1 and decreases the likelihood.

Typically, implementing Newton-Raphson requires inverting the Hessian, i.e., $\mathbf{d} = \mathbf{H}^{-1} \mathbf{g}$, a computation consuming $O(T^3)$ time. Instead, because \mathbf{M} is bidiagonal, the Hessian is tridiagonal, so the solution may be found in $O(T)$ time via standard banded Gaussian elimination techniques (which can be implemented efficiently in Matlab using $\mathbf{H} \setminus \mathbf{g}$). The resulting fast algorithm for solving the optimization problem in (4) is the main result of this paper, i.e. we may approximately infer the optimal solution for models characterized by equations such as (1) and (2) in linear time, whereas an exact solution would require exponential time, and a naïve approximate solution would require polynomial time.

2.3 Learning

In the above, we assumed that the parameters governing our model, $\boldsymbol{\theta} = \{\alpha, \beta, \sigma, \gamma, \nu, \rho, \lambda\} \in \boldsymbol{\Theta}$, were known. In general, however, these parameters may be estimated from the data. To find the maximum likelihood estimator for the parameters, $\hat{\boldsymbol{\theta}}$, we must integrate over the unknown variable, \mathbf{n} . However, integrating over all possible spike trains is typically approximated by Monte Carlo approaches, which is relatively slow. Thus, one often approximates:

$$\hat{\boldsymbol{\theta}} = \underset{\boldsymbol{\theta} \in \boldsymbol{\Theta}}{\operatorname{argmax}} \int P(\mathbf{F} | \mathbf{n}, \boldsymbol{\theta}) P(\mathbf{n} | \boldsymbol{\theta}) d\mathbf{n} \approx \underset{\boldsymbol{\theta} \in \boldsymbol{\Theta}}{\operatorname{argmax}} P(\mathbf{F} | \hat{\mathbf{n}}, \boldsymbol{\theta}) P(\hat{\mathbf{n}} | \boldsymbol{\theta}) \quad (9)$$

where $\hat{\mathbf{n}} = [\hat{n}_1, \dots, \hat{n}_T]$ is estimate of $\hat{\mathbf{n}}$ from the inference algorithm described above ($\hat{\mathbf{C}}$ is defined similarly). The approximation in (14) is good whenever the likelihood is very peaky, meaning that most of the mass is around the MAP sequence.¹ In particular, for state-space models, one often approximates the integral in (14) with the Viterbi path, i.e., the MAP path of the hidden states [Rabiner, 1989]. Finding the Viterbi path is often far easier than solving the integral in (14), as in the case here. Due to the state space nature of the above model (Eqs (1) and (2)), the optimization in Eq (9) simplifies significantly. More specifically, the likelihood term, $P(\mathbf{F}|\hat{\mathbf{n}}, \boldsymbol{\theta})$ may be written as:

$$\begin{aligned} P(\mathbf{F}|\hat{\mathbf{n}}, \boldsymbol{\theta}) &= \prod_{t=1}^T P(F_t|n_t, \boldsymbol{\theta}) = \prod_{t=1}^T \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left\{ -\frac{1}{2\sigma^2} (F_t - \alpha\hat{C}_t - \beta)^2 \right\} \\ &= \prod_{t=1}^T \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left\{ -\frac{1}{2\sigma^2} [F_t - \alpha(\gamma\hat{C}_{t-1} + \nu + \rho n_t) - \beta]^2 \right\} \end{aligned} \quad (10)$$

which follows from Eq (2) because \hat{C}_t is deterministic given \hat{n}_t . Considering for a moment, only the quadratic term in Eq (10), the terms can be reorganized:

$$F_t - \alpha(\gamma\hat{C}_{t-1} + \nu + \rho\hat{n}_t) - \beta = F_t - (\alpha\gamma)\hat{C}_{t-1} - (\alpha\nu + \beta) - (\alpha\rho)\hat{n}_t \quad (11)$$

which emphasizes the over-parameterization in our model. Thus, without loss of generality, we can assume that $\alpha = 1$ and $\beta = 0$. Taking logs in Eq (10) then yields:

$$\begin{aligned} \log P(\mathbf{F}|\hat{\mathbf{n}}, \boldsymbol{\theta}) &= -\frac{T}{2} \log(2\pi\sigma^2) - \frac{1}{2\sigma^2} \sum_{t=1}^T (F_t - \gamma\hat{C}_{t-1} - \nu - \rho\hat{n}_t)^2 \\ &= -\frac{T}{2} \log(2\pi\sigma^2) - \frac{1}{2\sigma^2} \left\| \mathbf{F} - \hat{\mathbf{X}}\boldsymbol{\eta} \right\|^2 \end{aligned} \quad (12)$$

where $\boldsymbol{\eta} = [\gamma, \nu, \rho]$, $\hat{\mathbf{X}} = [\hat{\mathbf{X}}_1, \dots, \hat{\mathbf{X}}_T]'$, and $\hat{\mathbf{X}}_t = [\hat{C}_{t-1}, 1, \hat{n}_t]$. Similarly, taking the log of the prior term, $P(\hat{\mathbf{n}}|\boldsymbol{\theta})$ yields:

$$\log P(\hat{\mathbf{n}}|\boldsymbol{\theta}) = \sum_{t=1}^T \log(\lambda\Delta) - \lambda\Delta\hat{n}_t = T(\lambda\Delta) - \lambda\Delta\hat{\mathbf{n}}'\mathbf{1} \quad (13)$$

Plugging Eqs (12) and (13) back into Eq (9), and noting that log is a monotonic function, we have:

$$\hat{\boldsymbol{\theta}} = \underset{\boldsymbol{\theta} \in \Theta}{\operatorname{argmax}} -\frac{T}{2} \log(2\pi\sigma^2) - \frac{1}{2\sigma^2} \left\| \mathbf{F} - \hat{\mathbf{X}}\boldsymbol{\eta} \right\|_2^2 + T \log(\lambda\Delta) - \lambda\Delta\hat{\mathbf{n}}'\mathbf{1} \quad (14)$$

Estimating the parameters then separates into three log-concave problems. In particular, solving for $\hat{\boldsymbol{\eta}}$ is simply:

¹The approximation in (14) may be considered a first-order Laplace approximation

$$\hat{\eta} = \underset{\eta, 0 < \eta_1 < 1}{\operatorname{argmin}} \frac{1}{2} \left\| \mathbf{F} - \widehat{\mathbf{X}} \eta \right\|_2^2 = \underset{\eta, 0 < \eta_1 < 1}{\operatorname{argmin}} \frac{1}{2} \eta' \mathbf{Q} \eta - \mathbf{L}' \eta \quad (15)$$

where $\mathbf{Q} = \widehat{\mathbf{X}}' \widehat{\mathbf{X}}$, and $\mathbf{L} = \widehat{\mathbf{X}}' \mathbf{F}$, and the constraint, $0 < \eta_1 < 1$, ensures that the time constant is both non-negative and finite. Eq (15) may be solved using standard quadratic optimization tools, such as Matlab's `quadprog`.

Unfortunately, the recursive nature of our algorithm makes solving for $\hat{\eta}$ in this manner unidentifiable. In particular, because both $\widehat{\mathbf{X}}$ and η can scale and shift \mathbf{C} , alternating between inferring $\widehat{\mathbf{X}}$ and $\hat{\eta}$ does not converge (data not shown). However, this is not particularly problematic, because the fluorescence measurements are in arbitrary units. This unit arbitrariness suggests that \mathbf{F} may be scaled and shift without loss of generality. The prior term and non-negativity constraint on \mathbf{n} , however, imposes small costs associated with scaling and shifting \mathbf{n} and \mathbf{C} . We therefore estimate η from the raw fluorescence trace. More specifically, to estimate γ , we manually find a fluorescence sequence that seem to be decaying, estimate the time constant τ , and let $\gamma = 1 - \Delta/\tau$. For ν , we manually find a fluorescence sequence that seems to be near baseline, label it $[F_s, \dots, F_t]$, and let $\nu = \frac{1}{t-s} \sum_{u=s}^t F_u$. Finally, we estimate ρ by letting it be equal to what we manually determine it be, by eye. In practice, we have found that by first scaling and shifting \mathbf{F} to be between 0 and 1, subsequent traces from the same or different cells have very similar values for η , and therefore, these parameters need not be tweaked much. Furthermore, the effective signal-to-noise ratio (SNR) of the inferred spike train does not depend strongly on these parameters, in agreement with previous results [Yaksi and Friedrich, 2006]. The other parameters, σ and λ , can be solved for analytically:

$$\hat{\sigma}^2 = \frac{1}{T} \left\| \mathbf{F} - \widehat{\mathbf{X}} \hat{\eta} \right\|_2^2 \quad (16)$$

$$\hat{\lambda} = \frac{1}{T \Delta} \hat{\mathbf{n}}' \mathbf{1} \quad (17)$$

3 Results

3.1 Main Result

To evaluate the efficacy of our fast filter, we compare it with the optimal linear (i.e., Wiener) filter [Wiener, 1949]. The top three panels of Fig. 1 depict a typical dataset simulated according to our model, Eqs. (1) and (2). Beneath the simulation, we show both the Wiener filter output (fourth panel) and our fast filter output (bottom panel). For both filters, we provided only the fluorescence observations depicted in the top panel, and η . From this data, we estimate the remaining parameters, and infer the hidden spike train. Several differences between these two approaches should be apparent. First, the effective signal-to-noise ratio (SNR) of our fast filter improves upon the optimal linear filter. Second, while the Wiener filter induces a “ringing” effect (where the inferred signal oscillates above and below zero), our fast filter completely eliminates this effect. These two improvements are common observations upon imposing a non-negative constraint when appropriate [Shumway and Stoffer, 2006]. Importantly, both our implementation of the Wiener filter

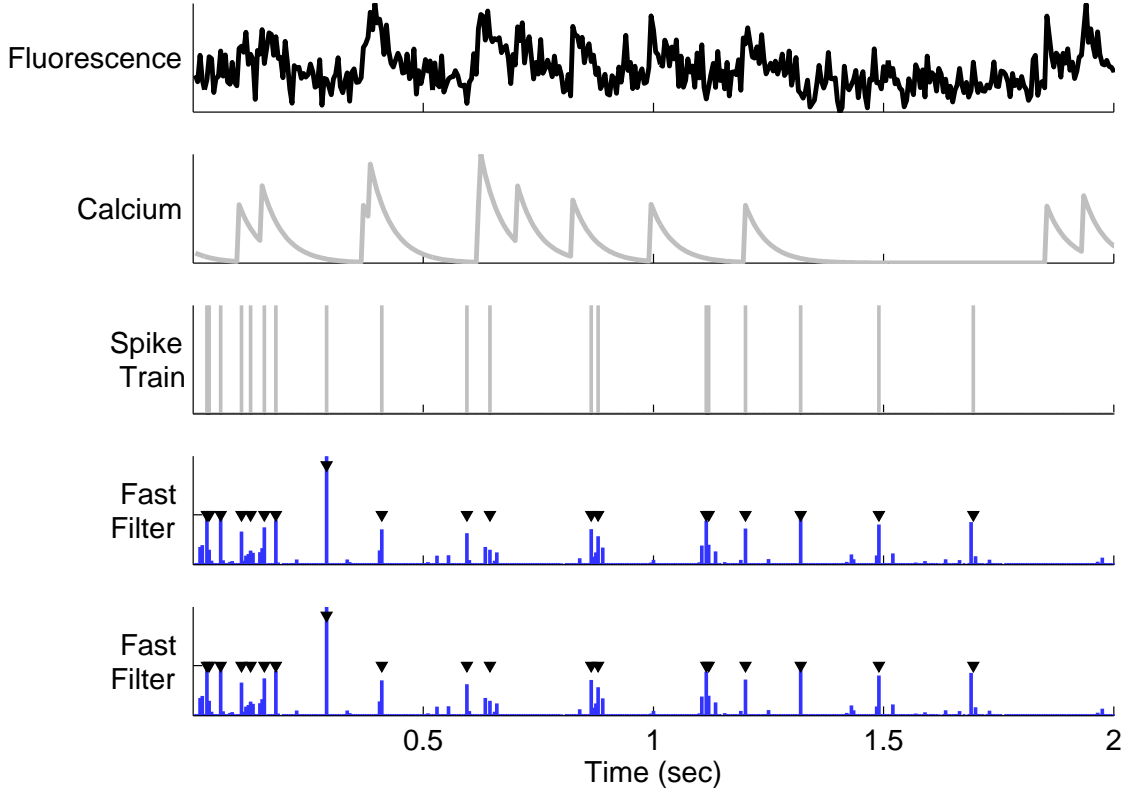


Figure 1: Simulation demonstrating our neuron model and inference. Note that the effective signal-to-noise ratio (SNR) of our fast filter improves on the Wiener filter, largely by eliminating the ringing effect present in the Wiener filter output. For both filters, only the fluorescence was provided to the algorithm, so both the parameters and the inference were performed using this minimal amount of data. Top panel: Simulated fluorescence. Second panel: Simulated intracellular calcium concentration. Third panel: Simulated spike train. Fourth panel: Wiener filter (blue for positive inference, red for negative), superimposed on simulated spike train (black triangles). Bottom panel: Fast filter (same conventions as fourth panel). Parameters: $\gamma = 0.94$, $\nu = 0$, $\rho = 1$, $\sigma = 0.3$, $\lambda = 8$, $\Delta = 0.005$ msec.

and our fast filter require only $O(T)$ time, whereas the naïve implementation of the Wiener filter requires $O(T \log(T))$, and the naïve implementation of a non-negative filter requires $O(T^3)$. Thus, when our approximation in Eq (4) is good, our filter outperforms the Wiener filter, and imposes approximately the same computational burden.

3.2 Fast Spiking Neuron

When the observed neuron is spiking quickly, the Poisson distribution may be well approximated by a Gaussian distribution, suggesting that the Wiener filter may be optimal in this regime. But the exponential approximation of a Poisson is also very accurate in the fast spiking regime. To compare these two strategies, we simulated a fast spiking neuron, where the expected number of spikes per bin exceeds 10 (Fig. 2). In this scenario, both the Wiener filter and our fast filter perform approximately equally well. Both filters infer peaks in the firing rate that are obscured

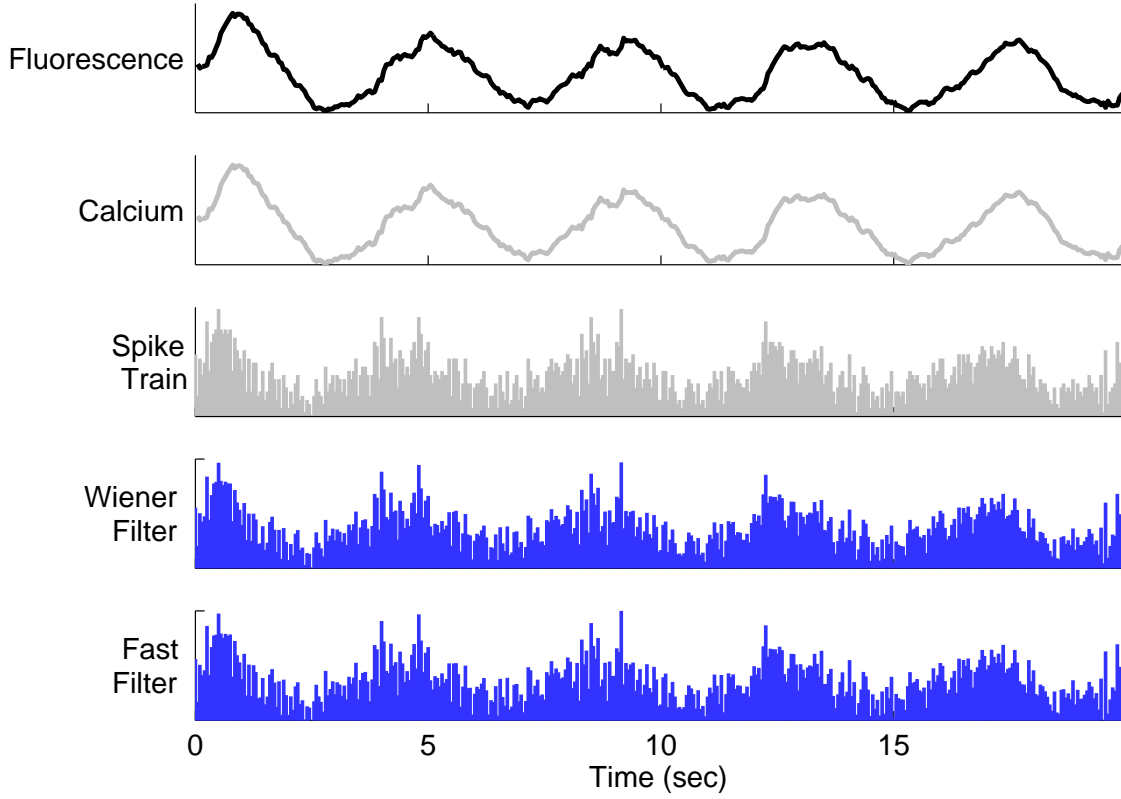


Figure 2: Comparing the Wiener filter and our fast filter for a fast spiking neuron. In this scenario, the two filters perform approximately equally well. Note that both recover fast fluctuations in the firing rate that are smoothed out by the calcium dynamics. Conventions as in Fig. 1. Actual and inferred spike trains were normalized similarly, to ease comparisons. Parameters: $\gamma = 0.94$, $\nu = 0$, $\rho = 1$, $\sigma = 0.3$, $\lambda = 100$ (modulated by a sinusoid), $\Delta = 0.05$ msec.

by the low-pass filter properties of the calcium dynamics. Thus, it seems from this analysis that regardless of the firing rate of the observable neuron, (1) filtering the signal may provide valuable information, and (2) our fast filter performs at least as well as the Wiener filter, without requiring more computational time.

3.3 Spatial Filtering

In the above, we assumed a one-dimensional (1-D) observation, F . The raw data, however, is actually a series of multidimensional images. To get the 1-D time-series, two pre-processing steps are required. First, for each neuron, one must define a region-of-interest (ROI), essentially assigning sets of pixels to neurons. This yields a vector time-series for each neuron, $\vec{F}_t \in \mathbb{R}^m$. Second, one must project the m -D observations into a 1-D time-series. Typically, this step is performed by averaging all the pixels within the ROI. In theory, one can improve on this uniform averaging by estimating the optimal spatial filter. Here, we provide details on how to incorporate this second step (projecting the m -D observations into 1-D) into our fast filter. First, we replace Eq (3) with:

$$\vec{F}_t = \alpha C_t + \beta + \Sigma \varepsilon_t, \quad \varepsilon_t \sim \mathcal{N}(\mathbf{0}, \mathbf{I}) \quad (18)$$

where \vec{F}_t , α , and β are all m dimensional column vectors, Σ is the covariance matrix, and ε_t is an m -dimensional standard normal random vector. α now represents the optimal spatial filter, and β represents the baseline for each pixel. Our goal then is to estimate α and β , to provide maximal information about the underlying spike train. We would expect this optimal spatial filtering improve the effective SNR after filtering in many situations. For example, if certain pixels are anti-correlated with others, a uniform spatial filter would average out those differences, whereas the optimal filter would take advantage of that information.

To estimate α and β , we first plug Eqs (18) and (2) into (9), to obtain:

$$\begin{aligned} \vec{\theta} = \underset{\vec{\theta} \in \vec{\Theta}}{\operatorname{argmax}} & -\frac{T}{2} \log(2\pi|\Sigma|) + T \log(\lambda\Delta) - \lambda\Delta \mathbf{n}'\mathbf{1} \\ & - \frac{1}{2} \sum_{t=1}^T (\vec{F}_t - \alpha(\gamma\hat{C}_{t-1} - \nu - \rho\hat{n}_t) - \beta)' \Sigma^{-1} (\vec{F}_t - \alpha(\gamma\hat{C}_{t-1} - \nu - \rho\hat{n}_t) - \beta) \end{aligned} \quad (19)$$

where $|\cdot|$ indicates the determinant. By pre-whitening the observation matrix, $\vec{F} = [\vec{F}_1, \vec{F}_2, \dots, \vec{F}_T]$, we can estimate the covariance matrix by the identity, i.e., $\Sigma = \mathbf{I}$. Without loss of generality, we may assume that $\nu = 0$ and $\rho = 1$, similar to our assumption of $\alpha = 1$ and $\beta = 0$ above. This leaves γ , α , and β . Assuming that γ is known (or estimated from the raw fluorescence signal, as we did above), we have:

$$\hat{\eta}_v = \underset{\|\eta_v\| < 1}{\operatorname{argmax}} \left\| \vec{F} - \eta_v \mathbf{X}_v \right\|^2 \quad (20)$$

where $\eta_v = [\alpha, \beta]$, and $\mathbf{X}_v = [\gamma\mathbf{C} + \mathbf{n}, \mathbf{1}]'$. Note that we have constrained the norm of η_v , which is necessary because otherwise there is a free scale between \mathbf{X}_v and η_v . Problems of the form as in Eq (20) are known as “trust region subproblems”. Fortunately, many algorithms for solving such a problem are available [Fortin, 2000].

3.4 Incorporating a nonlinear observation model

3.5 Optimal thresholding

3.6 Slow dynamics

4 Discussion

Acknowledgments Support for JTV was provided by NIDCD DC00109. LP is supported by an NSF CAREER award, by an Alfred P. Sloan Research Fellowship, and the McKnight Scholar Award. BOW was supported by NDS grant F30 NS051964. The authors would like to thank A. Packer for helpful discussions.

References

- [Boyd and Vandenberghe, 2004] Boyd, S. and Vandenberghe, L. (2004). *Convex Optimization*. Oxford University Press.
- [Fortin, 2000] Fortin, C. (2000). *A Survey of the Trust Region Subproblem within a Semidefinite Framework*. PhD thesis, University of Waterloo.
- [Hastie et al., 2001] Hastie, T., Tibshirani, R., and Friedman, J. (2001). *The Elements of Statistical Learning: Data Mining, Inference, and Prediction*. Springer.
- [Ikegaya et al., 2004] Ikegaya, Y., Aaron, G., Cossart, R., Aronov, D., Lampl, I., Ferster, D., and Yuste, R. (2004). Synfire chains and cortical songs: temporal modules of cortical activity. *Science*, 304(5670):559–564.
- [MacLean et al., 2005] MacLean, J. N., Watson, B. O., Aaron, G. B., and Yuste, R. (2005). Internal dynamics determine the cortical response to thalamic stimulation. *Neuron*, 48(5):811–823.
- [Niell and Smith, 2005] Niell, C. M. and Smith, S. J. (2005). Functional imaging reveals rapid development of visual response properties in the zebrafish tectum. *Neuron*, 45(6):941–951.
- [Ohki et al., 2005] Ohki, K., Chung, S., Ch’ng, Y. H., Kara, P., and Reid, R. C. (2005). Functional imaging with cellular resolution reveals precise micro-architecture in visual cortex. *Nature*, 433(7026):597–603.
- [Rabiner, 1989] Rabiner, L. R. (1989). A Tutorial on Hidden Markov Models and Selected Applications in Speech Recognition. *Proceedings of the IEEE*, 72(2):257–286.
- [Sato et al., 2007] Sato, T. R., Gray, N. W., Mainen, Z. F., and Svoboda, K. (2007). The functional microarchitecture of the mouse barrel cortex. *PLoS Biol*, 5(7):e189.
- [Shumway and Stoffer, 2006] Shumway, R. and Stoffer, D. (2006). *Time Series Analysis and Its Applications*. Springer, 2nd edition.
- [Vogelstein et al., 2008] Vogelstein, J., Watson, B., AM, P., Yuste, R., B, J., and Paninski, L. (2008). Spike inference from calcium imaging using sequential monte carlo methods. *Biophysical Journal*.
- [Wiener, 1949] Wiener, N. (1949). *Extrapolation, interpolation, and smoothing of stationary time series: with engineering applications*. MIT Press, Cambridge, Mass.
- [Yaksi and Friedrich, 2006] Yaksi, E. and Friedrich, R. W. (2006). Reconstruction of firing rate changes across neuronal populations by temporally deconvolved Ca^{2+} imaging. *Nat Methods*, 3(5):377–383.

A Wiener Filter

Sections 2.2 and 2.3 outline one approach to solving Eq. (3), by approximating the Poisson distribution with an exponential distribution, and imposing a non-negative constraint on the inferred $\hat{\mathbf{n}}$. Perhaps a more straightforward approach would be to approximate the Poisson distribution with a Gaussian distribution. In fact, as rate increases above about 10 spikes/sec, a Poisson distribution with rate $\lambda\Delta$ is well approximated by a Gaussian with mean and variance $\lambda\Delta$. Given such an approximation, instead of Eq. (4), we would obtain:

$$\hat{\mathbf{n}}_w \approx \underset{n_t \in \mathbb{R}, \forall t}{\operatorname{argmin}} \sum_{t=1}^T \left(\frac{1}{2\sigma^2} (F_t - C_t)^2 + \frac{1}{2\lambda\Delta} (n_t - \lambda\Delta)^2 \right) \quad (21)$$

As above, we can rewrite Eq. (21) in matrix notation in terms of \mathbf{C} :

$$\hat{\mathbf{C}}_w = \underset{C_t \in \mathbb{R}, \forall t}{\operatorname{argmin}} \frac{1}{2\sigma^2} \|\mathbf{F} - \mathbf{C}\|^2 + \frac{1}{2\lambda\Delta} \|\mathbf{M}\mathbf{C} - \lambda\Delta\mathbf{1}\|^2 \quad (22)$$

which is quadratic in \mathbf{C} , and may therefore be solved analytically using quadratic programming, $\hat{\mathbf{C}}_w = \hat{\mathbf{C}}_0 + \mathbf{d}_w$, where $\hat{\mathbf{C}}_0$ is the initial guess and $\mathbf{d}_w = \mathbf{H}_w \setminus \mathbf{g}_w$, where

$$\mathbf{g}_w = \frac{1}{\sigma^2} (\mathbf{C}'_0 - \mathbf{F}) + \frac{1}{\lambda\Delta} ((\mathbf{M}\hat{\mathbf{C}}_0)' \mathbf{M} - \lambda\Delta \mathbf{M}' \mathbf{1}) \quad (23)$$

$$\mathbf{H}_w = \frac{1}{\sigma^2} \mathbf{I} + \frac{1}{\lambda\Delta} \mathbf{M}' \mathbf{M} \quad (24)$$

Note that this solution is the optimal linear solution, under the assumption that spikes follow a Gaussian distribution, and is often referred to as the Wiener filter, regression with a smoothing prior, or ridge regression. To estimate the parameters for the Wiener filter, we take the same approach as above:

$$\hat{\boldsymbol{\theta}}_w \approx \underset{\boldsymbol{\theta}_w}{\operatorname{argmax}} P(\mathbf{F} | \hat{\mathbf{n}}_w, \boldsymbol{\theta}_w) P(\hat{\mathbf{n}}_w | \boldsymbol{\theta}_w) \quad (25a)$$

$$= \underset{\boldsymbol{\theta}_w}{\operatorname{argmax}} -\frac{T}{2} \log(4\pi^2 \sigma^2 \lambda\Delta) - \frac{1}{2\sigma^2} \|\mathbf{Y}_w + \boldsymbol{\eta}_w \mathbf{X}_w\|^2 - \frac{1}{2\lambda\Delta} \|\hat{\mathbf{n}}_w - \lambda\Delta\mathbf{1}\|^2 \quad (25b)$$

where \mathbf{Y}_w , $\boldsymbol{\eta}_w$, and \mathbf{X}_w are defined as their subscriptless counterparts in Eq. (9).

B old

Generalization of the optimal nonnegative filter The above filter design assumed a very simple model relating spikes to $[\text{Ca}^{2+}]$, and $[\text{Ca}^{2+}]$ to fluorescence. Both (1) and (2) may be generalized in a straightforward manner. In fact, our model may be thought of as a special case of the more general state-space framework:

Figure 3: Simulation demonstrating our neuron model and inference. Note that our filter accurately infers the unobserved spike train, given only the noisy fluorescence measurements. Top panel: Simulated fluorescence. Second panel: Simulated intracellular calcium concentration. Third panel: Simulated spike train. Fourth panel: Output of optimal nonnegative filter (green) superimposed on simulated spike train (black). Parameters: $\lambda = 10$, $\Delta = 0.025$ msec, $\tau = 0.5$ sec, $\sigma = 1$.

$$\mathbf{C}_t = \mathbf{D}\mathbf{C}_{t-1} + \mathbf{1}n_t \quad (26a)$$

$$\mathbf{F}_t = \mathbf{A}\mathbf{C}_t + \mathbf{b} + \boldsymbol{\varepsilon}_t \quad (26b)$$

which we can solve in linear time for any log-concave distribution of n_t and $\boldsymbol{\varepsilon}_t$. For our particular scenario of interest, i.e., that of inferring spike trains from calcium sensitive fluorescence observations, these generalizations facilitate considering a model with much richer dynamics. For instance, a more accurate model relating spikes to $[\text{Ca}^{2+}]$ would consider myriad calcium extrusion mechanisms, buffering, etc. In such a situation, we would represent $[\text{Ca}^{2+}]$ as a vector, $\mathbf{C}_t \in \mathbb{R}^{N \times 1}$, where each element of \mathbf{C}_t would correspond to a different spike dependent calcium process. The differential operator, $\mathbf{D} \in \mathbb{R}^{N \times N}$, accounts for the relative strength of each of these mechanisms, and their time constants. The spikes, n_t , are multiplied by a column vector of ones, $\mathbf{1} \in \mathbb{R}^{N \times 1}$, because the magnitude of the effect of a spike on each element in \mathbf{C}_t is scaled by ρ .² It should be clear that as in (2), spikes are related to calcium by a simple linear transformation. In the multi-dimensional scenario, however, \mathbf{D} — a matrix — would replace a — a scalar — so \mathbf{M} becomes block-bidiagonal, making the Hessian block tridiagonal. Nonetheless, Gaussian elimination may still solve $\mathbf{H}\mathbf{C} = \mathbf{g}$ in $O(T)$, so our filter may be applied to this scenario as well.

Another natural extension of this work would be to let \mathbf{F}_t also be multidimensional. In (1), although we assumed that we had access to a monodimensional fluorescent magnitude at each time, the raw data is actually a multidimensional movie. These images, however, are necessarily blurred by the point-spread-function of the camera. Thus, we could replace (1) with (26b), where \mathbf{A} is the point-spread function of the camera, and \mathbf{b} sets the relative baseline intensity per pixel. Furthermore, the noise, $\boldsymbol{\varepsilon}_t$, could have any log-concave distribution, and these results would still hold. Given these generalizations, this filter may be applied to a large class of problems.

Projection Pursuit Regression Projection pursuit regression (PPR) is another technique one could use to infer a spike train from fluorescence measurements. PPR is different from the optimal nonnegative filter in a few ways. First, it solves a related, but slightly different problem: instead of finding the MAP estimate of the spike train, PPR finds the maximum likelihood estimate, i.e., the spike train that makes the fluorescence measurements most likely. Second, PPR constrains the solution to have only nonnegative integers. Therefore, $\mathbf{n}_{PPR} = \arg\max_{\mathbf{n}_t \in \mathbb{N}_0} p(\mathbf{F}|\mathbf{n})$, where \mathbb{N}_0 is the set of nonnegative integers.. Third, because of this constraint, PPR requires an additional parameter, w , that sets the size of the calcium transient caused by a single action potential. This forces us to substitute (2) with $C_t = aC_{t-1} + wn_t$. Given these differences, to find \mathbf{n}_{PPR} , PPR proceeds iteratively, adding a spike with each iteration, as long as doing so reduces the residual square error (i.e., the sum of the squared difference between the inferred \mathbf{C} and \mathbf{F}). One obtains the

²For the same reasons that there is no parameter scaling the spikes in the monodimensional case

Figure 4: Comparison of various linear filters for a simulated Poisson neuron spiking with a rate of 1 Hz (left panels) and 200 Hz (right panels). The main result is that the two filters proposed outperform the optimal linear (i.e., Wiener) filter. Top panels: Fluorescence measurements. Second panels: Number of spikes per frame. Third panels: Optimal linear (i.e., Wiener) filter output given the above fluorescence signal. The Wiener filter does not impose a nonnegativity constraint, thus the inferred spike train can either be positive (green) or negative (red). Fourth panels: Same as third panels, but for the optimal nonnegative filter. Note the absence of negative spikes. Fifth panels: Same as third panels, but using the fast Projection Pursuit Regression (fPPR) filter. Parameters: $\Delta = 0.025$ sec, $\tau = 0.5$ sec, $\lambda = 1/(\text{firing rate})$, left $\sigma = 0.4$, right $\sigma = 1$.

time of the next inferred spike by finding the maximum of the convolution of the current residual square error with the calcium kernel, $we^{-t\Delta/\tau}$. In general, this procedure takes $O(T \log T)$ per iteration, as the convolution is computed in the Fourier domain. However, the recursive state-space representation in (26) implies that this convolution requires just $O(T)$ time here. Henceforth, we therefore refer to this approach as fast PPR (or fPPR). This approach may be generalized to the multidimensional cases, as in the previous filter. However, unlike the previous filter, the likelihood function is not concave (recall that PPR is a greedy optimization method), so we are no longer guaranteed to converge to the optimal solution.

C Results

To evaluate the efficacy of our filters, we compare their results to that of the optimal linear (i.e., Wiener) filter [Wiener, 1949]. The Wiener filter differs in construction from our filters in a few ways. First, it imposes no constraint on the spike train. Second, it is optimal upon assuming that the prior distribution of n_t is Gaussian. In our model, spiking was Poisson, (2), and the nonnegativity of n_t in the sparse-spiking regime makes the Gaussian model assumed by the Wiener filter inaccurate (Fig. 2, left). However, when spike rates are fast — e.g., on average, several spikes per image frame — a Poisson distribution is better approximated by a Gaussian distribution. Furthermore, at high firing rates, the mean of the Gaussian would be relatively far from zero, and the variance proportional, so the probability of sampling a negative number would be relatively small, obviating the need for the nonnegativity constraint. Thus, one might expect the Wiener filter to perform as well as our nonnegative filter (and fPPR) when the observed neuron has a high firing rate (and relatively low imaging rate). Furthermore, given that the calcium kernel is exponential, the Wiener filter, like the filters we developed here, only requires $O(T)$ time, as opposed to the typical $O(T \log T)$.

Fig. 4 depicts a comparison between the filters developed above and the Wiener filter for two different scenarios: a slow firing rate simulation (left panels) and a fast firing rate simulation (right panels). When action potentials are sparse, the two filters we propose above outperform the Wiener filter. Moreover, at high firing rates, all three filters perform approximately equally well.

While in simulations, all the above algorithms perform reasonably well, the true test is how well they perform given data from live cells. We simultaneously recorded both electrophysiologically and imaged with an epifluorescent microscope, from a pyramidal neuron in a somatosensory cortical slice, as described in [MacLean et al., 2005]. Fig. 5 shows an example fluorescence time-series

Figure 5: Inferring a spike train given fluorescence measurements from a live *in vitro* neuron. Simultaneous recording of a saturating fluorescence signal (black line in top panel), and its associated spike train (black impulses in all panels). Note how the three filters handle saturation differently. The optimal nonnegative filter’s signal-to-noise ratio degrades as saturation increases, but does not suffer a catastrophic failure (green in second panel). fPPR suffers more than the optimal nonlinear particle filter when the signal saturates, due to the hard threshold (green in third panel). The optimal nonlinear particle filter [Vogelstein et al., 2008], which explicitly incorporates saturation, correctly identifies each spike time (green in bottom panel).

in which the neuron spiked with a relatively low rate, but the calcium accumulated nonetheless, leading to fluorescence saturation (top panel). In practice, this kind of strong saturation is rarely observed (personal communication, anonymous), so this example is designed to test the limits of performance of our filters. In fact, the optimal nonnegative filter accurately infers every spike, even when the fluorescence is strongly saturating, and the signal-to-noise ratio is very poor (second panel). On the other hand, fPPR, which has a sharp threshold for including an additional spike, performs relatively poorly (third panel), demonstrating the dependency of this method on a good model fit. For comparison purposes, we also show the performance of an optimal nonlinear particle filter [Vogelstein et al., 2008], which specifically incorporates a nonlinear saturation function. While the optimal nonlinear particle filter performs better in scenarios such as this one, the computational burden is increased by approximately 100-fold relative to the fast nonnegative filter. Thus these two filters serve complementary purposes: the fast nonnegative filter is better geared to rapid online analysis of large-scale multi-neuronal data, whereas the nonlinear particle filter is better suited for offline refinement of the results from the fast nonnegative filter.

D Conclusions

We show here that for certain nonnegative deconvolution problems, we can derive an algorithm that is both optimal and efficient. More specifically, our algorithm may be applied to any model with a nonnegative signal that is linearly filtered by a matrix linear ordinary differential equation. We apply this approach to the problem of inferring the most likely spike train given noisy calcium sensitive fluorescence observations (c.f. Fig. 3), and demonstrate, in simulations, that the optimal nonnegative filter outperforms the optimal linear (i.e., Wiener) filter in both slow and fast firing rate regimes (c.f. Fig. 4). Furthermore, when applied to data from a live cell, the optimal nonnegative filter outperforms a fast projection pursuit regression filter, which constrains the inferred spike train to be nonnegative integers (c.f. Fig. 5). On the other hand, the nonnegative filter is based on a linear observation model, and therefore suffers a loss of precision in the presence of strong saturation effects, in contrast to the optimal nonlinear particle filter (c.f. Fig. 5).

The implications of these results are severalfold. First, it seems as if there is no reason to use the Wiener filter for scenarios in which our algorithm may apply. Second, as our filter is so efficient, it may be used for many real-time processing applications. Specifically, upon simultaneously imaging a population of neurons [Ikegaya et al., 2004, Niell and Smith, 2005, Ohki et al., 2005, Yaksi and Friedrich, 2006, Sato et al., 2007], our filter may be applied essentially online. This could greatly expedite the tuning of important experimental parameters — such as laser intensity

— to optimize signal-to-noise ratio for inferring spikes. Third, the parameters estimated from this filter may be used to initialize the parameters of the optimal nonlinear particle filter, which may then be used offline, to further refine the spike train inference. Future work will consider multidimensional models for this application, incorporating both more sophisticated calcium models, and spatial filtering for extracting the fluorescence signal, obviating the need for additional algorithms for image segmentation.

E other

We have found empirically that despite the approximation in (4), upon initializing the parameters with a reasonable guess, the algorithm tends to converge to parameters that are reasonably close to the true parameters, resulting in an accurate spike reconstruction. Importantly, estimating these parameters typically requires only a very short sequence of observations, e.g., several seconds of data including about 5–10 spikes (not shown).