Machine Learning Methods for Neural Data Analysis Spike Sorting by Deconvolution

Scott Linderman

STATS 220/320 (NBIO220, CS339N).



Spike Sorting by Deconvolution

Improving upon the simple model

- Our simple model was a good warm-up, but we made some strong assumptions:
 - Namely, that we could extract windows for each spike and assign them to one neuron.
 - All spikes from the same neuron share the same mean,
- With lots of recording channels, there are likely to be overlapping spikes, and they could have different amplitudes.
- Next, we'll extend the simple model with a more realistic one using **convolutional matrix** factorization.
- The resulting model will be very similar to **Kilosort** [Pachitariu et al., 2023]



20 ms

10,000ft view

- Idea: each time a neuron spikes, it adds a scaled copy of its template to the measured voltage.
- Formally, we model the data as a sum of convolutions of templates and amplitudes for each neuron, plus noise.



ConvolutionIn one dimension

• **Convolution** is an operation that takes in a signal a(t) and a filter w(t) and outputs

$$y(t) = [a \circledast w](t) = \int a(t-\tau)w(\tau) \,\mathrm{d}\tau.$$

- In **discrete time** this becomes, $y_t = [a \circledast w]_t = \sum_{d=-\infty}^{\infty} a_{t-d}w_d$
- **Causal** filters are constrained so that $w_d = 0$ for d < 0. Then y_t is only influenced by $a_{1:t}$.
- Our filters will also have **bounded support** so that $w_d = 0$ for $d \ge D$. Then y_t is only influenced by $a_{t-D+1:t}$.
- In our case, the signal is the time series of spike amplitudes, and the filter is the waveform template.
 Every time there's a spike, we plop down a scaled template.





Convolution With multiple output channels

 We need to convolve the amplitude signal with multiple filters in parallel, one for each channel of the voltage recording.

$$\mathbf{y}_{t} = \begin{pmatrix} [a \circledast w_{1}]_{t} \\ \vdots \\ [a \circledast w_{N}]_{t} \end{pmatrix} = \begin{pmatrix} \sum_{d=1}^{D} a_{t-d} w_{1,d} \\ \vdots \\ \sum_{d=1}^{D} a_{t-d} w_{N,d} \end{pmatrix}$$

• (I'm going to index d from $1, \ldots, D$ because the notation is a bit simpler.)



Convolution With multiple input & output channels

 Finally, we need to sum convolutions of multiple input signals, one for each neuron in the model.

 $\mathbf{X} = \mathbf{A} \circledast \mathbf{W}$

by which we mean

$$x_{n,t} = \sum_{k=1}^{K} \sum_{d=1}^{D} a_{k,t-d} w_{k,n,d}$$



Cross-Correlation In one dimension

- Cross-correlation essentially goes in reverse.
- In signal processing, the cross-correlation is a sliding dot product of data y(t) and template w(t), which produces a new function $[y \star w](t)$.
- For discrete time, real-valued inputs,

$$[y \star w]_t = \sum_{d=-\infty}^{\infty} y_{t+d} w_d.$$

• With a change of variables, we see that crosscorrelation is equivalent to convolution with a timereversed filter \widehat{w} :

$$[y \star w]_t = \sum_{d=\infty}^{-\infty} y_{t-d} w_{-d} = [y \circledast \overleftarrow{w}]_t.$$

 (Note: the definition of cross-correlation is not unique. This definition is consistent with np.correlate but opposite of Wikipedia.)



Cross-CorrelationWith multiple channels

As before, we can extend this definition to handle multiple channels

$$[\mathbf{Y} \star \mathbf{W}]_t = \sum_{n=1}^N \sum_{d=-\infty}^\infty y_{n,t+d} w_{n,d}.$$

- The cross-correlation measures the similarity of the data and the template at each point in time.
- The **auto-correlation** is the crosscorrelation of a signal with itself.



Cross-Correlation With multiple input & output channels

As before, we can extend this definition to handle multiple channels

$$[Y \star W]_{t} = \begin{pmatrix} [Y \star W_{1}]_{t} \\ \vdots \\ [Y \star W_{K}]_{t} \end{pmatrix}$$
$$= \begin{pmatrix} \sum_{n=1}^{N} \sum_{d=-\infty}^{\infty} y_{n,t+d} W_{1,n,d} \\ \vdots \\ \sum_{n=1}^{N} \sum_{d=-\infty}^{\infty} y_{n,t+d} W_{K,n,d} \end{pmatrix}$$



Convolution and Cross-Correlation in Pytorch

- PyTorch (and other deep learning libraries) have fast, GPU-backed implementations of convolutions.
- What they call convolution is actually cross-correlation!
- But remember, we can always get convolution by cross-correlating with the flipped filter.
- For discrete time signals, you have to play with **padding** to handle **edge effects**.
- By default, these functions operate on minibatches of inputs, so you need to add an extra leading dimension to your signal.
- There are lots of other options to read about (strides, dilations, groups), but we won't use them this week.

torch.nn.functional.conv1d(*input*, *weight*, *bias=None*, *stride=1*, *padding=0*, *dilation=1*, *groups=1*) → Tensor

Applies a 1D convolution over an input signal composed of several input planes.

This operator supports TensorFloat32.

See Conv1d for details and output shape.

• NOTE

In some circumstances when using the CUDA backend with CuDNN, this operator may select a nondeterministic algorithm to increase performance. If this is undesirable, you can try to make the operation deterministic (potentially at a performance cost) by setting torch.backends.cudnn.deterministic = True. Please see the notes on Reproducibility for background.

Parameters

- **input** input tensor of shape (minibatch, in_channels, *iW*)
- weight filters of shape (out_channels, $\frac{\text{in_channels}}{\text{groups}}, kW$)
- **bias** optional bias of shape (out_channels). Default: None
- stride the stride of the convolving kernel. Can be a single number or a one-element tuple (sW,). Default: 1
- padding implicit paddings on both sides of the input. Can be a single number or a one-element tuple (padW,). Default: 0
- dilation the spacing between kernel elements. Can be a single number or a oneelement tuple (dW,). Default: 1
- groups split input into groups, in_channels should be divisible by the number of groups. Default: 1

Spike sorting by deconvolution

Probabilistic Model Likelihood

• Assume each spike is a noisy, scaled version of the template of the neuron that generated it.

$$p(\mathbf{X} \mid \mathbf{A}, \mathbf{W}) = \prod_{t=1}^{T} \mathcal{N}\left(\mathbf{x}_{t} \mid \sum_{k=1}^{K} [\mathbf{a}_{k} \circledast \mathbf{W}]\right)$$





Probabilistic Model Prior on spike amplitudes

 Assume the spike amplitudes are drawn from an exponential distribution.

$$a_{k,t} \sim \operatorname{Exp}(\lambda)$$

- This simple prior will lead to sparse amplitudes, but it does not encode any dependencies between time steps.
- Ideally, we would also like to prohibit a neuron from producing two spikes within D samples of each other.



https://en.wikipedia.org/wiki/Exponential_distribution

Probabilistic Model Scale invariance via Frobenius norm constraint

- For the amplitudes to be meaningful, we need to constrain the scale of the templates.
- Assume the matrix $\mathbf{W}_k \in \mathbb{R}^{N \times D}$ has unit **Frobenius norm** $\|\mathbf{W}_k\|_{\mathrm{F}} = 1$.

Probabilistic Model Aside: The Frobenius norm and the SVD

• The Frobenius norm of a matrix is induced by the Frobenius inner product $\langle \mathbf{A}, \mathbf{B} \rangle_F = \text{Tr}(\mathbf{A}^\top \mathbf{B})$. With this definition,

$$\|\mathbf{W}\|_F^2 = \langle \mathbf{W}, \mathbf{W} \rangle_F = \mathrm{Tr}(\mathbf{W}^\top \mathbf{W}).$$

• It is equal to the ℓ_2 norm of the flattened matrix,

$$\|\mathbf{W}\|_{F}^{2} = \sum_{n=1}^{N} \sum_{d=1}^{D} w_{n,d}^{2} = \operatorname{vec}(\mathbf{W})^{\mathsf{T}}\operatorname{vec}(\mathbf{W}) = \|\mathbf{w}\|_{F}^{2}$$

• We can also write it in terms of the singular values of \mathbf{W} ,

$$\|\mathbf{W}\|_{\mathrm{F}} = \mathrm{Tr}(\mathbf{V}\mathbf{S}\mathbf{U}^{\mathsf{T}}\mathbf{U}\mathbf{S}\mathbf{V}^{\mathsf{T}}) = \mathrm{Tr}(\mathbf{S}^{2}) = \|\mathbf{W}\|_{\mathrm{F}}$$

• where $W = USV^{\top}$ with S = diag(s) is the singular value decomposition (SVD).





 $\|\mathbf{S}\|_2$,

Probabilistic Model Scale invariance via Frobenius norm constraint

- For the amplitudes to be meaningful, we need to constrain the scale of the templates.
- Assume the matrix $\mathbf{W}_k \in \mathbb{R}^{N \times D}$ has unit Frobenius norm $\|\mathbf{W}_k\|_{\mathrm{F}} = 1$.
- This is equivalent to constraining the **singular** values to be normalized $\|\mathbf{s}_k\|_2 = 1$.



Probabilistic Model Low-rank constraint

- This view suggests a further assumption: constraint the **rank** of the templates as well.
- If we constrain it to be rank 1 (i.e., only one nonzero singular value), then

$$\mathbf{W}_k = \mathbf{u}_k \mathbf{v}_k^{\mathsf{T}}$$

where $\mathbf{u}_k \in \mathbb{S}_{N-1}$ is the **spatial footprint** and $\mathbf{v}_k \in \mathbb{S}_{D-1}$ is the **temporal profile**,

and $\mathbb{S}_{N-1} = \{ \mathbf{u} : \mathbf{u} \in \mathbb{R}^N \text{ and } \|\mathbf{u}\|_2 = 1 \}$ is the set of unit vectors in \mathbb{R}^N .



MAP estimation

Maximum a posteriori estimation **Coordinate ascent**

- Initialize templates W and set A = 0.
- Iterate until convergence:
 - For neuron $k = 1, \dots, K$:
 - a. Optimize **amplitudes** a_k for neuron k.

b. Optimize **templates** \mathbf{W}_k for neuron k.

fixed.]

- [In each case, maximize log joint probability wrt one variable, holding others

Maximum a posteriori estimation Deriving the log likelihood

The likelihood is,



Maximum a posteriori estimation **Deriving the log likelihood**

Taking the log of both sides and expanding the Gaussian density yields,



Maximum a posteriori estimation **Deriving the log likelihood**

$\log p(\mathbf{X} \mid \mathbf{A}, \mathbf{W}) = -$

where $\mathbf{R}_k \in \mathbb{R}^{N \times T}$ is the residual for neuron k, defined as

$$\mathbf{R}_k = \mathbf{X} -$$

As a function of \mathbf{a}_k and \mathbf{W}_k , fixing the other neurons' amplitudes and templates,

$$-\frac{1}{2\sigma^2} \|\mathbf{R}_k - \mathbf{a}_k \circledast \mathbf{W}_k\|_{\mathrm{F}}^2$$

$$\sum_{\substack{j \neq k}} [\mathbf{a}_j \circledast \mathbf{W}_j]$$

As a function of \mathbf{a}_k , the log *joint* probability is the log *likelihood* plus the log *prior* probability,

 $\log p(\mathbf{X}, \mathbf{W}, \mathbf{A}) =$ $= \log p(\mathbf{X} \mid \mathbf{A}, \mathbf{W}) + \log p(\mathbf{a}_k; \lambda)$ $= -\frac{1}{2\sigma^2} \|\mathbf{R}_k - \mathbf{a}_k \circledast \mathbf{W}_k\|_{\mathrm{F}}^2 + \log p(\mathbf{a}_k) + c$



Now, expand the square in the Frobenius norm,

$$\log p(\mathbf{X}, \mathbf{W}, \mathbf{A}) = -\frac{1}{2\sigma^2} \|\mathbf{R}_k - \mathbf{a}_k \otimes \mathbf{W}_k\|_F^2 + \log p(\mathbf{a}_k) + c$$
$$= -\frac{1}{2\sigma^2} \|\mathbf{a}_k \otimes \mathbf{W}_k\|_F^2 + \frac{1}{\sigma^2} \langle \mathbf{R}_k, \mathbf{a}_K \otimes \mathbf{W}_k \rangle_F + \log p(\mathbf{a}_k) + c$$

 $\mathscr{L}_2(\mathbf{a}_k)$

 $\mathscr{L}_1(\mathbf{a}_k)$

Further expanding the quadratic term,

$$\begin{aligned} \mathscr{L}_{2}(\mathbf{a}_{k}) &= -\frac{1}{2\sigma^{2}} \|\mathbf{a}_{k} \circledast \mathbf{W}_{k}\|_{F}^{2} \\ &= -\frac{1}{2\sigma^{2}} \sum_{t=1}^{T} \sum_{n=1}^{N} \left(\sum_{d=1}^{D} a_{k,t-d}^{2} w_{k,n,d}^{2} + 2 \sum_{d=1}^{D} a_{k,t-d}^{2} w_{k,n,d}^{2} + 2 \sum_{d=1}^{D} a_{k,t-d}^{2} w_{k,n,d}^{2} + 2 \sum_{d=1}^{D} a_{k,t-d}^{2} \|\mathbf{W}_{k}\|_{F}^{2} \\ &\approx -\frac{1}{2\sigma^{2}} \sum_{t=1}^{T} a_{k,t}^{2} \|\mathbf{W}_{k}\|_{F}^{2} \end{aligned}$$

with equality when nonzero entries (i.e. "spikes") of \mathbf{a}_k are separated by at least D samples.

 $\sum_{d=1}^{D} \sum_{d'=1}^{d-1} a_{k,t-d} a_{k,t-d'} w_{k,n,d} w_{k,n,d'} \bigg)$

Now take the linear term...

$$\begin{aligned} \mathscr{L}_{1}(\mathbf{a}_{k}) &= \frac{1}{\sigma^{2}} \langle \mathbf{R}_{k}, \mathbf{a}_{k} \circledast \mathbf{W}_{k} \rangle \\ &= \frac{1}{\sigma^{2}} \sum_{t=1}^{T} \sum_{n=1}^{N} \sum_{n=1}^{N} r_{k,n,t} [\mathbf{a}_{k} \circledast \mathbf{w}_{k,n}]_{t} \\ &= \frac{1}{\sigma^{2}} \sum_{t=1}^{T} \sum_{n=1}^{N} \sum_{d=1}^{D} a_{k,t-d} r_{k,n,t} w_{k,n,d} \\ &= \frac{1}{\sigma^{2}} \sum_{t=1}^{T} a_{k,t} \sum_{n=1}^{N} \sum_{d=1}^{D} r_{k,n,t+d} w_{k,n,d} \\ &= \frac{1}{\sigma^{2}} \sum_{t=1}^{T} a_{k,t} [\mathbf{R}_{k} \star \mathbf{W}_{k}]_{t} \end{aligned}$$

where $[\mathbf{R} \star \mathbf{W}_k]_t$ is the cross-correlation of the residual and the template for neuron k.



Putting it all together

$$\mathscr{L}(\mathbf{a}_k) = \sum_{t=1}^T \left[-\frac{1}{2\sigma^2} a_{k,t}^2 + \frac{1}{\sigma^2} \mu_{k,t} a_{k,t} - \lambda a_{k,t} \right] + c,$$

which separates into a sum of quadratic objective functions for each time t.

Maximum a posteriori estimation **Completing the square and solving for the optimal amplitudes**

 The maximum of this quadratic objective, subject to non-negativity constraints, is obtained at

$$a_{k,t} = \max\left\{0, \mu_{k,t} - \sigma^2 \lambda\right\}$$

- However, we also want spikes to be wellseparated; i.e. $a_{k,t} > 0 \implies a_{k,t+d} = 0$ for d = 1, ..., D.
- We'll enforce this with a **simple heuristic**: use find peaks to select local maxima of this "score" signal.





As a function of \mathbf{W}_k $\log p(\mathbf{X}, \mathbf{A}, \mathbf{W}) = -\frac{1}{2\sigma^2} \sum_{i=1}^{I} \langle a_{k,t} \mathbf{R}_{k,t}, \mathbf{W}_k \rangle + c'$ $= -\frac{1}{2\sigma^2} \left\langle \sum_{t=1}^{T} a_{k,t} \mathbf{R}_{k,t}, \mathbf{W}_k \right\rangle + c'$ where $\mathbf{R}_{k,t} = \begin{bmatrix} r_{k,1,t} & \cdots & r_{k,1,t+D} \\ \vdots & & \vdots \\ r_{k,N,t} & \cdots & r_{k,N,t+D} \end{bmatrix}$

is a slice of the residual matrix (R[k,:,t:t+D] in code).

We want to maximize this log joint probability over the space of low-rank, unit-norm matrices,

$$\mathbf{W}_{k}^{\star} = \arg \max_{\mathbf{W}_{k} \in \mathbb{S}_{R}^{N,D}} \left\langle \sum_{t=1}^{T} a_{k,t} \mathbf{R}_{k,t}, \mathbf{W}_{k} \right\rangle$$

The solution is to set the waveform matrix "proportional to" the weighted sum of residual matrices by taking its SVD and renormalizing the singular values.

$$\mathbf{W}_{k}^{\star} = \sum_{r=1}^{R} \bar{s}_{r} \mathbf{u}_{r} \mathbf{v}_{r}^{\mathsf{T}} \quad \text{where} \quad \bar{s}_{r} = \frac{s_{r}}{\sqrt{\sum_{r'=1}^{R} \frac{s_{r'}}{\sqrt{\sum_{r'=1}^{R} \frac{s_{r'}}}{\sqrt{\sum_{r'=1}^{R} \frac{s_{r'}}{\sqrt{\sum_{r'=1}^{R} \frac{s_{r'}}{\sqrt{\sum_{r'}} \frac{s_{r'}}{\sqrt{\sum_{r'}} \frac{s_{r'}}{\sqrt{\sum_{r'}}$$



singular vectors

$$=1 \frac{s_{r'}^2}{r}$$

e

More efficient computation Leveraging the low-rank templates

We can compute the "scores" for amplitude updates more efficiently by leveraging the low-rank templates,

$$[\mathbf{R}_{k} \star \mathbf{W}_{k}]_{t} = \sum_{n=1}^{N} \sum_{d=1}^{D} r_{k,n,t+d} \mathbf{W}_{k,n,d}$$
$$= \sum_{d=1}^{D} \mathbf{r}_{k,:,t+d}^{\top} \mathbf{W}_{k,:,d}$$
$$= \sum_{d=1}^{D} \mathbf{r}_{k,:,t+d}^{\top} \mathbf{U}_{k} \mathbf{S}_{k} \mathbf{v}_{k,:,d}$$
$$= \sum_{d=1}^{D} (\mathbf{U}_{k}^{\top} \mathbf{r}_{k,:,t+d})^{\top} [\mathbf{S}_{k} \mathbf{V}_{k}^{\top}]_{:,d}$$
$$= [(\mathbf{U}_{k}^{\top} \mathbf{R}_{k}) \star (\mathbf{S}_{k} \mathbf{V}_{k}^{\top})]_{t}$$

In other words, we cross-correlate the projected residual.



►	

Conclusion

- We developed a basic spike sorting model that was good for building intuition, but not very practical.
- We developed a new model for the voltage in terms of a superposition of templates convolved with spike amplitudes for each neuron.
 - Along the way, we learned about convolution and cross-correlation.
- We derived a **coordinate ascent algorithm** for *maximum a posteriori* (MAP) inference.
- Next time: you'll implement the algorithm in lab! You'll learn a bit of PyTorch for implementing the convolutions and cross-correlations, then test it out on the GPU.

Further reading

- Simple Spike Sorting and Spike Sorting by Deconvolution course notes.
- Convolution and cross-correlation:
 - <u>convnets.html</u>)
 - generated/torch.nn.functional.conv1d.html
- Spike sorting:
 - sorting problem with Kilosort." bioRxiv (2023).
 - The model we presented is a slightly modified version of *Kilosort*

Chapter 9 of The Deep Learning Book (deeplearningbook.org/contents/

Start reading up on PyTorch convolutions! <u>https://pytorch.org/docs/stable/</u>

Pachitariu, Marius, Shashwat Sridhar, and Carsen Stringer. "Solving the spike