# Machine Learning Methods for Neural Data Analysis

**Lecture 11: Unsupervised modeling**

Scott Linderman

# Announcements

- **Proposals** look great!

  - Feedback coming in the next day or two.

  - Next update due Mar 5 (next Friday).

# Agenda
## Finish decoding and start unsupervised modeling

- "Direct" decoders and structured prediction

- Unit III: Unsupervised models of neural and behavioral data

- Bayesian inference in latent variable models

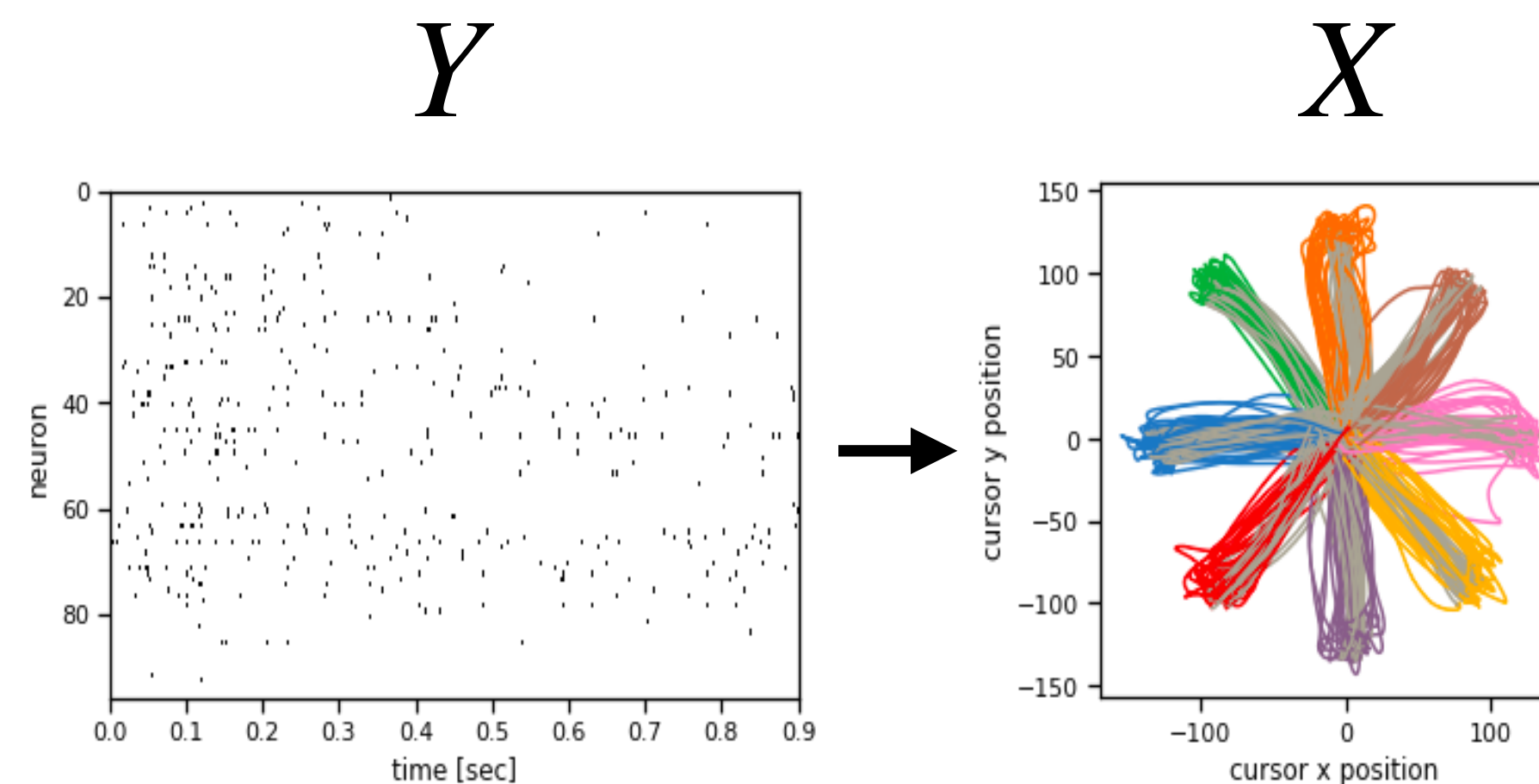# Decoding movement from neural spike trains
## A linear dynamical system (LDS) model

- One of the problems with the basic model (Lab 6 Part 1) is that it treated each time bin as independent.

- Instead, consider the following prior

$$p(X) = p(x_1) \prod_{t=2}^{T} p(x_t \mid x_{t-1})$$

$$= \mathcal{N}(x_1 \mid 0, Q) \prod_{t=2}^{T} \mathcal{N}(x_t \mid Ax_{t-1}, Q)$$

- Parameterized by **dynamics matrix** $A \in \mathbb{R}^{D \times D}$ and **dynamics covariance** $Q \in \mathbb{R}^{D \times D}$.

$Y$

$X$

# Decoding movement from neural spike trains
## The posterior distribution

The posterior is given by

$$p(X \mid Y) \propto \left[ \mathcal{N}(x_1 \mid 0, Q) \prod_{t=2}^{T} \mathcal{N}(x_t \mid Ax_{t-1}, Q) \right] \left[ \prod_{t=1}^{T} \mathcal{N}(y_t \mid Cx_t + d, R) \right]$$

$$= \exp \left\{ \text{"a big quadratic function of } X\text{"} \right\}$$

$$= \mathcal{N}(\text{vec}(X) \mid \mu, \Sigma)$$

Where $\Sigma = J^{-1}$ and $\mu = J^{-1}h$ with,

$$J = \begin{bmatrix} J_{11} & J_{21}^{\top} & & & \\ J_{21} & J_{22} & J_{32}^{\top} & & \\ & \ddots & \ddots & \ddots & \\ & & \ddots & \ddots & J_{T,T-1}^{\top} \\ & & & J_{T,T-1} & J_{TT} \end{bmatrix} \qquad h = \begin{bmatrix} h_1 \\ h_2 \\ \vdots \\ h_T \end{bmatrix}$$

- The blocks are given by $J_{tt} = Q^{-1} + A^{\top}Q^{-1}A + C^{\top}R^{-1}C$ (except for $J_{11}$ and $J_{TT}$), $J_{t,t-1} = -Q^{-1}A$, and $h_t = C^{\top}R^{-1}(y_t - d)$.

# Decoding movement from neural spike trains
## Poisson observations

- So far we've used a linear, Gaussian encoder for the spikes, even though they are counts!

- Suppose instead,

$$p(Y \mid X) = \prod_{t=1}^{T} \prod_{n=1}^{N} \text{Po}\left(y_{tn} \mid f(c_n^{\top} x_t + d_n)\right)$$

- The posterior is no longer Gaussian, but it's common to approximate it as one.

# Decoding movement from neural spike trains
## Laplace approximation

Approximate the posterior as

$$p(X \mid Y) \approx \mathcal{N}(\mu, \Sigma)$$

where

$$\mathcal{L}(X) = -\log p(X, Y)$$

$$\mu = \text{argmin}_X \, \mathcal{L}(X)$$

$$\Sigma = \left[ \nabla^2 \mathcal{L}(X) \Big|_{X=\mu} \right]^{-1}$$

For GLM encoders, the log joint is concave
and $\mu$ and $\Sigma$ can be found efficiently.

# Decoding movement from neural spike trains
## Laplace approximation under a Poisson GLM encoder

Derive the Hessian under the Poisson GLM encoder with exponential mean function $f(a) = e^a$:

$$\frac{\partial^2}{\partial x_t \partial x_t} \mathcal{L}(X) = J_{tt} - \sum_{n=1}^{N} \frac{\partial^2}{\partial x_t \partial x_t} \log \mathrm{Po}\left(y_{tn} \mid f(c_n^\top x_t + d_n)\right)$$

$$= J_{tt} - \sum_{n=1}^{N} \frac{\partial^2}{\partial x_t \partial x_t} \left[-f(c_n^\top x_t + d_n) + y_{tn} \log f(c_n^\top x_t + d_n)\right]$$

$$= J_{tt} + \sum_{n=1}^{N} \exp\{c_n^\top x_t + d_n\} c_n c_n^\top$$

# Decoding movement from neural spike trains
## Structured decoders

- If we're going to make a Gaussian approximation anyway, why not learn more flexible means and covariances?

- Recall the form of the LDS posterior,

$$J_{tt} = Q^{-1} + A^\top Q^{-1} A + C^\top R^{-1} C$$

$$J_{t,t-1} = -Q^{-1} A$$

$$h_t = C^\top R^{-1}(y_t - d)$$

- **Idea**: replace these with learned functions of $y_{1:T}$.

# Decoding movement from neural spike trains
## Structured decoders

For example,

$$p(X \mid Y) = \mathcal{N}(\text{vec}(X) \mid \mu, \Sigma)$$
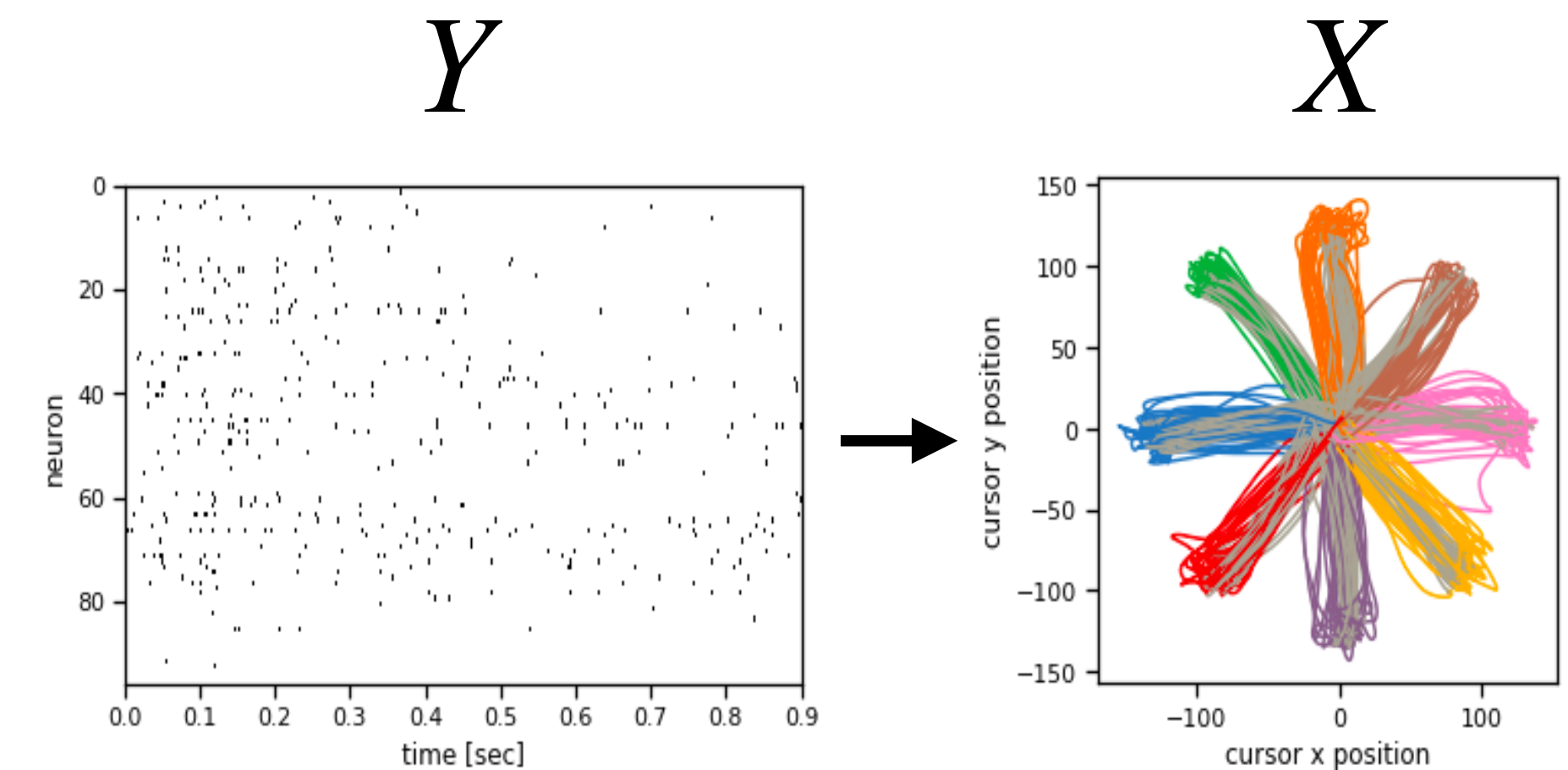
$$\mu = J(Y)^{-1}h(Y)$$

$$\Sigma = J(Y)^{-1}$$

Where

$$J_{tt} = Q^{-1} + A^\top Q^{-1}A + \textcolor{red}{f(y_{t-\Delta:t+\Delta})}$$
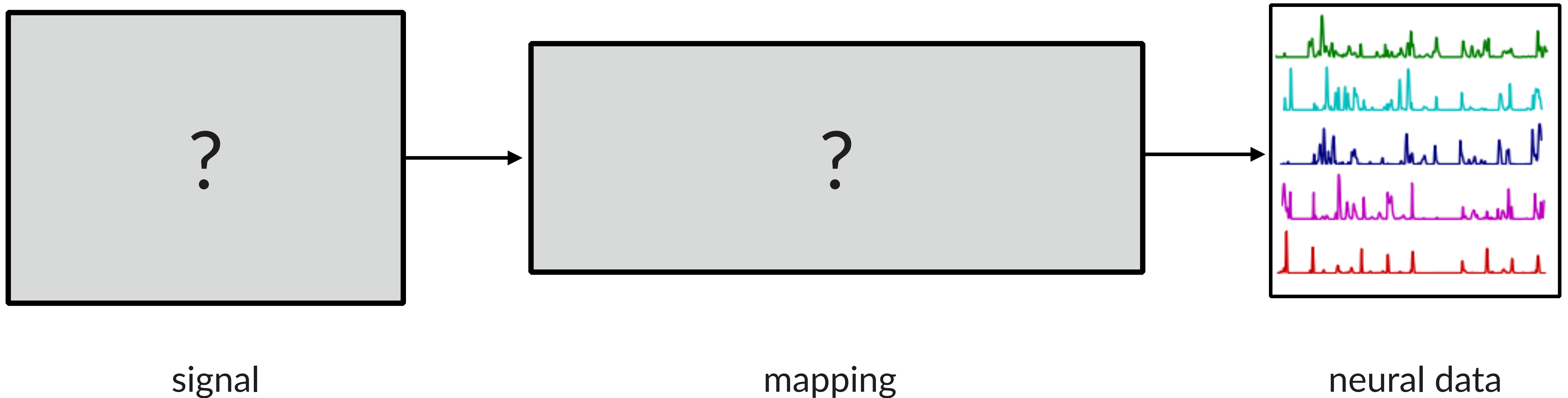
$$J_{t,t-1} = -Q^{-1}A$$

$$h_t = \textcolor{red}{g(y_{t-\Delta:t+\Delta})}$$

$Y$



$X$

# Unit III: Unsupervised modeling

# Data-driven modeling
## Searching for signals to explain neural activity



signal                          mapping                          neural data

# Data-driven modeling
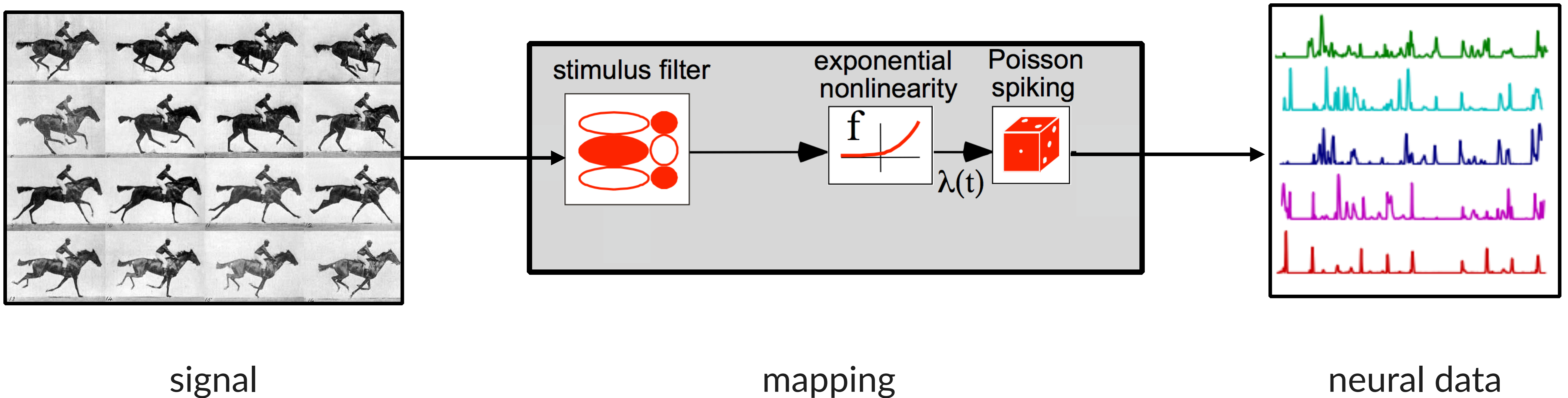## Searching for signals to explain neural activity



signal                              mapping                              neural data

Encoding models: given stimulus (covariates) and response, find mapping.

# Data-driven modeling
## Searching for signals to explain neural activity



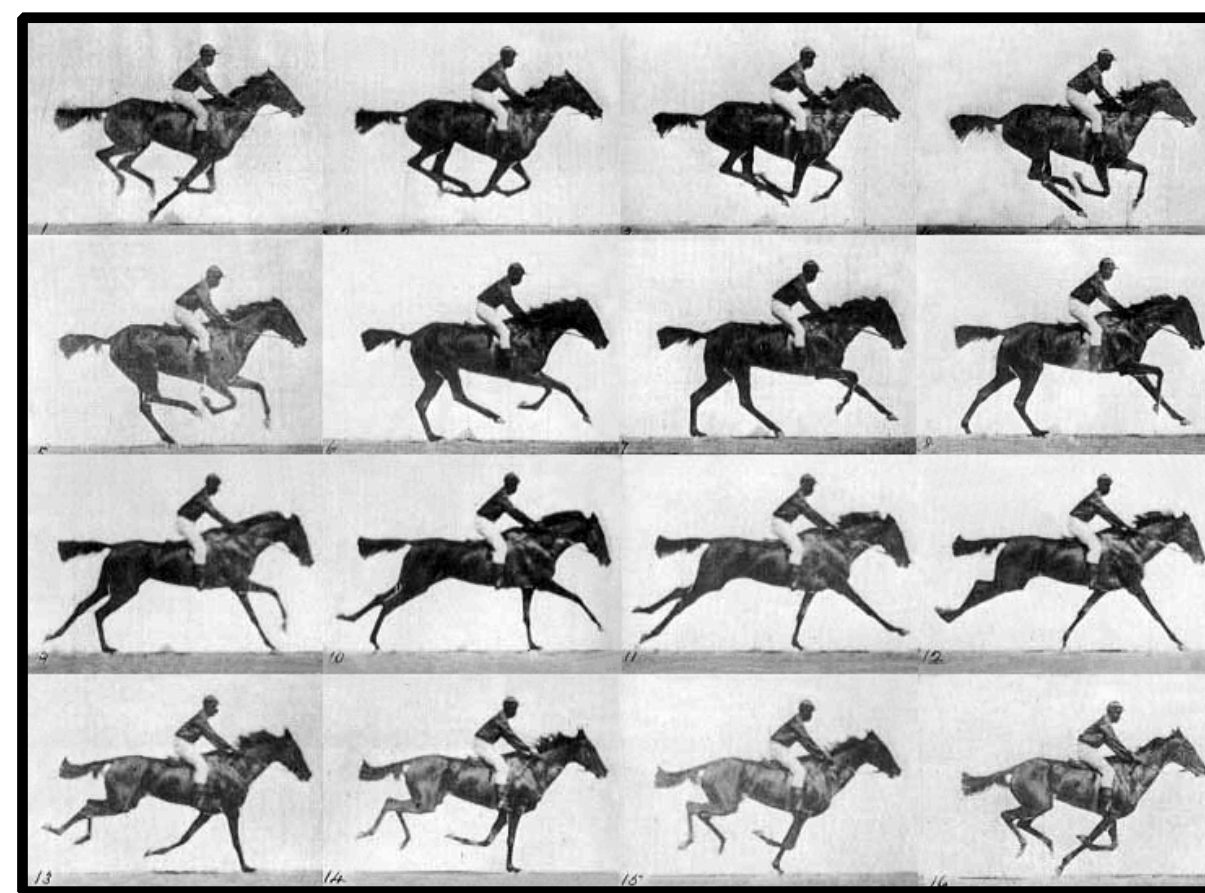signal                           mapping                          neural data

Recent examples: Musall et al (2018), Stringer et al (2018)
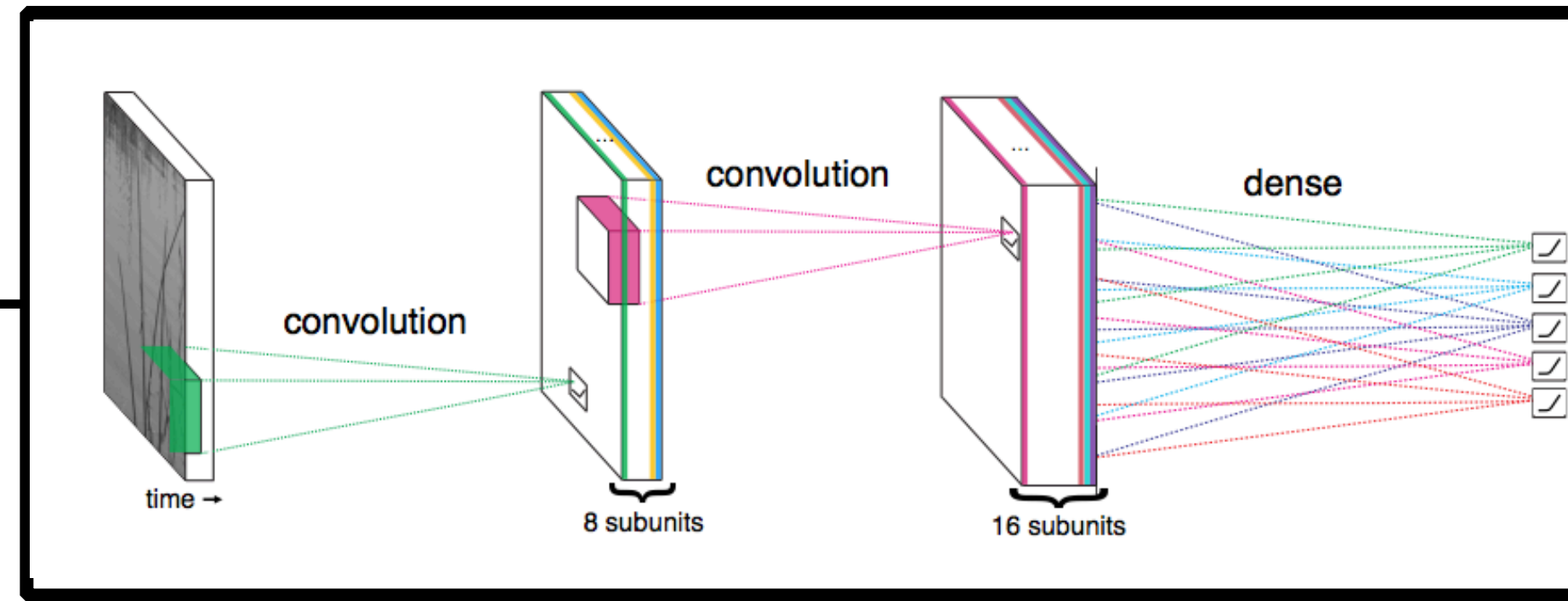
Paninski (2004)
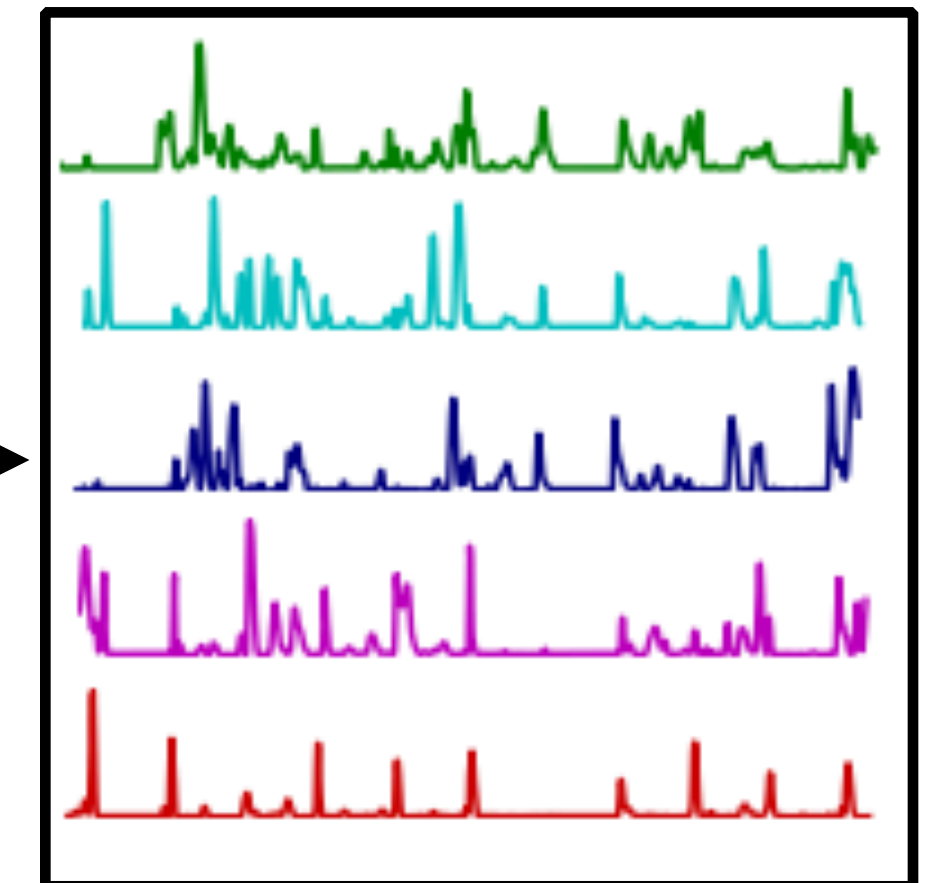Truccolo et al (2005)
Pillow et al (2008)

# Data-driven modeling
## Searching for signals to explain neural activity
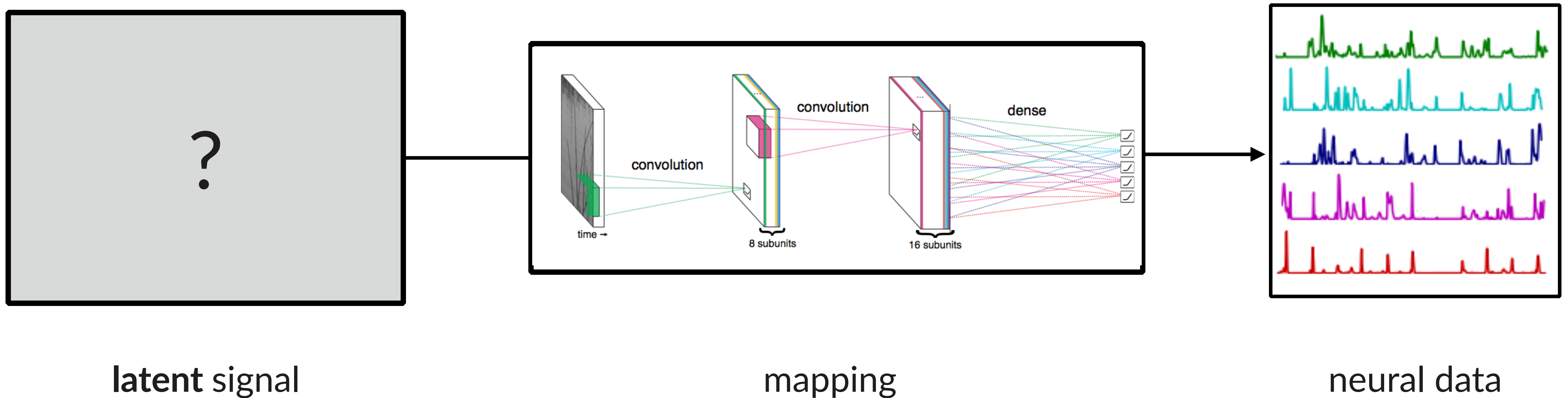


signal                    mapping                    neural data

Toward nonlinear and/or more biophysically plausible mappings.
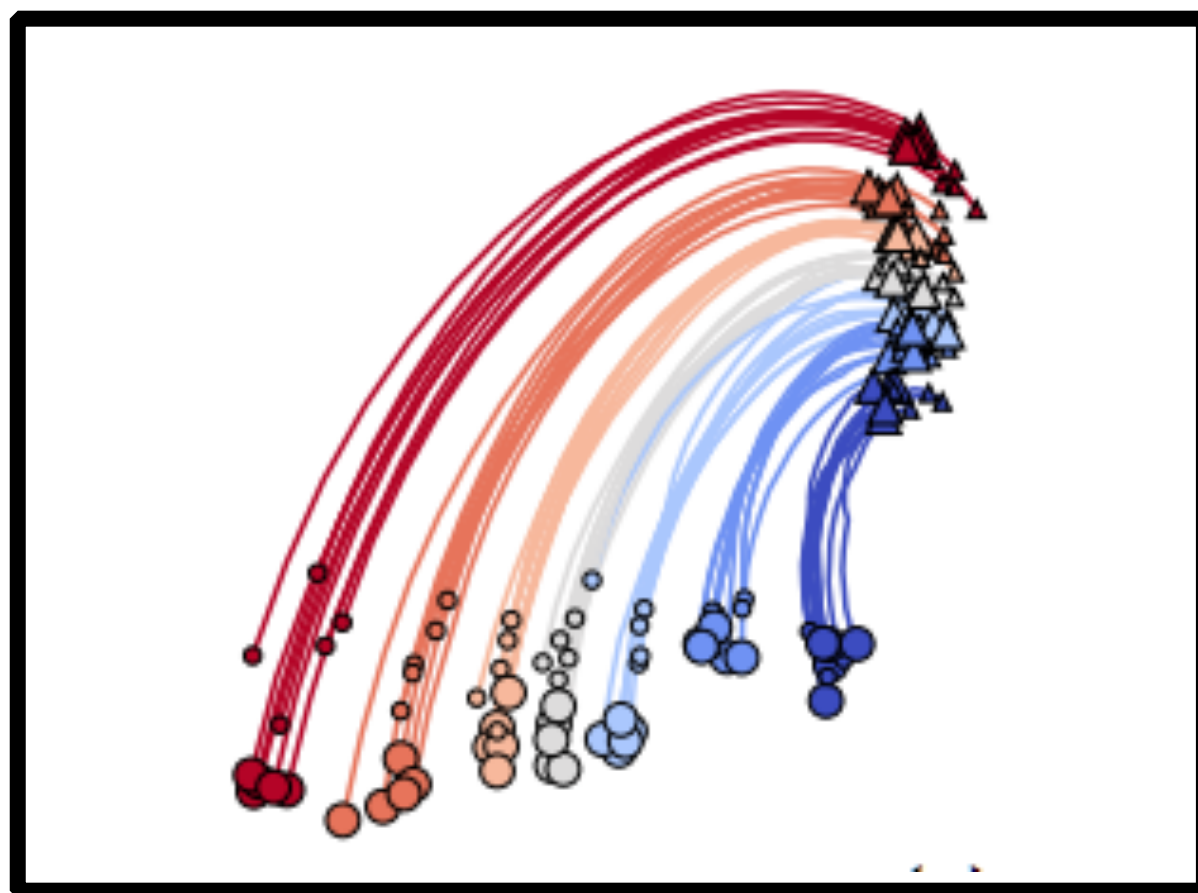
McIntosh et al (2017)

# Data-driven modeling
## Searching for signals to explain neural activity



latent signal                                  mapping                          neural data
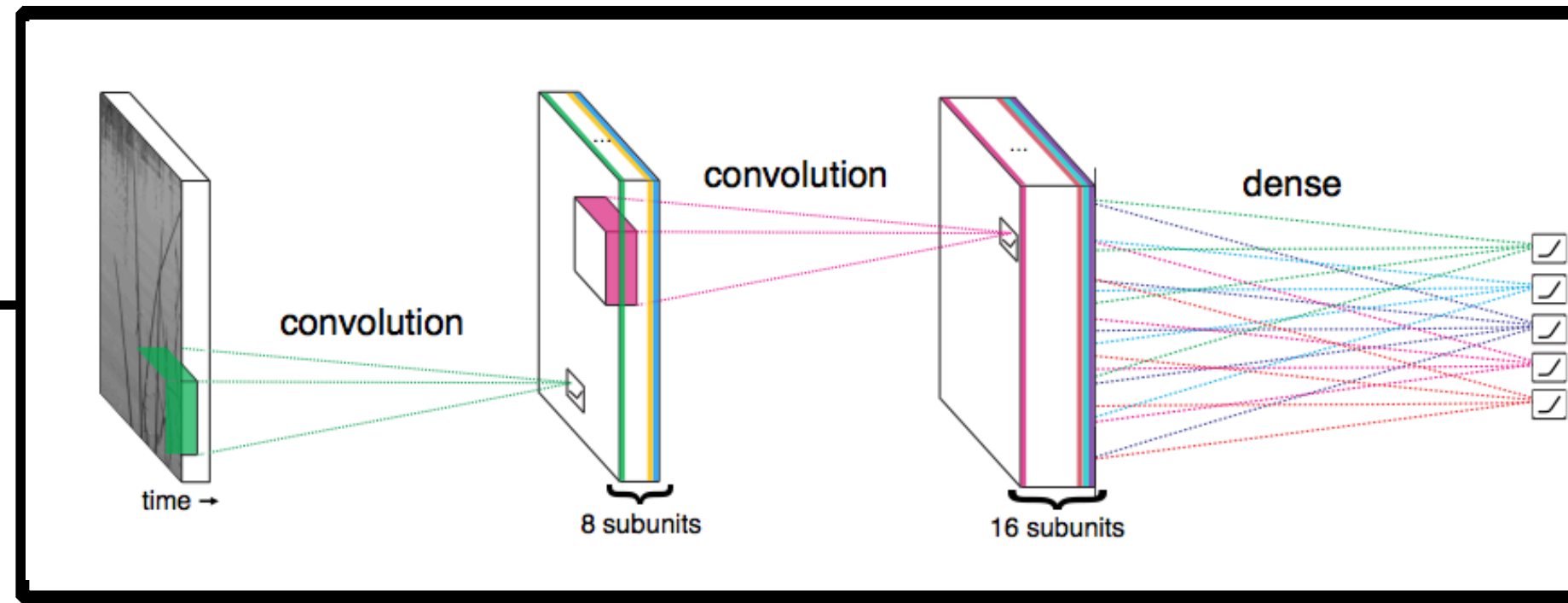
**Alternative:** try to infer latent signals from the data
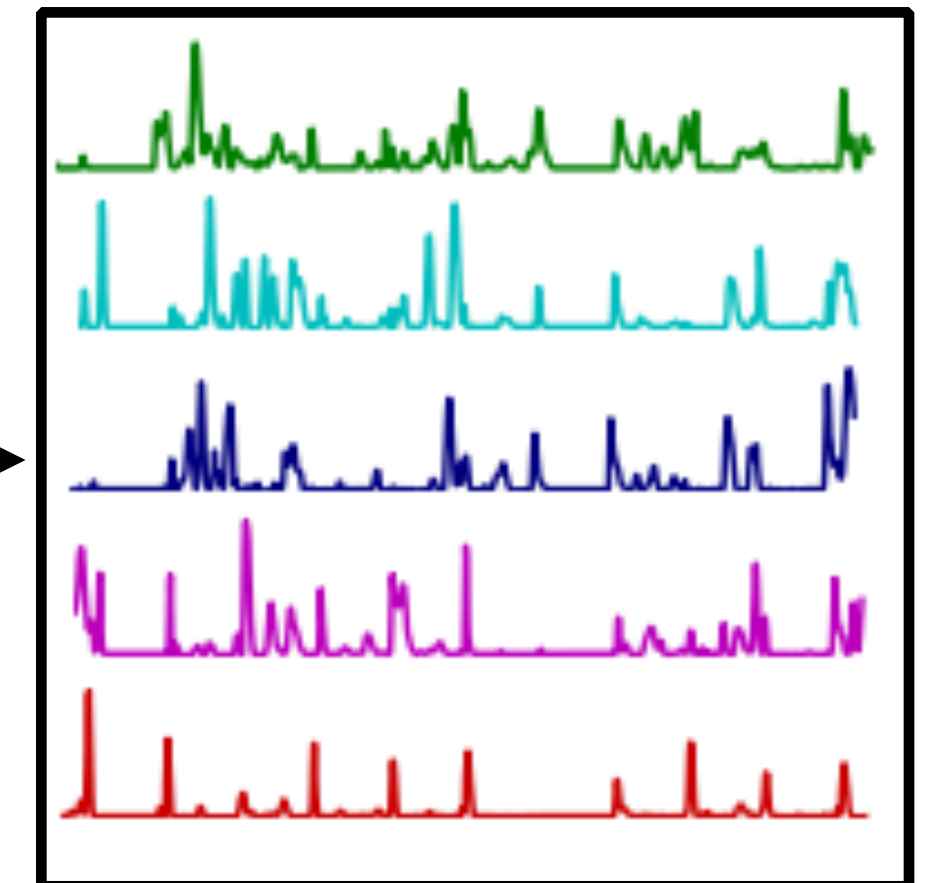
# Data-driven modeling
## Searching for signals to explain neural activity



latent signal        mapping        neural data

**Alternative:** try to infer latent signals from the data, *subject to constraints.*

# Latent variable modeling is all about constraints
## The five D's

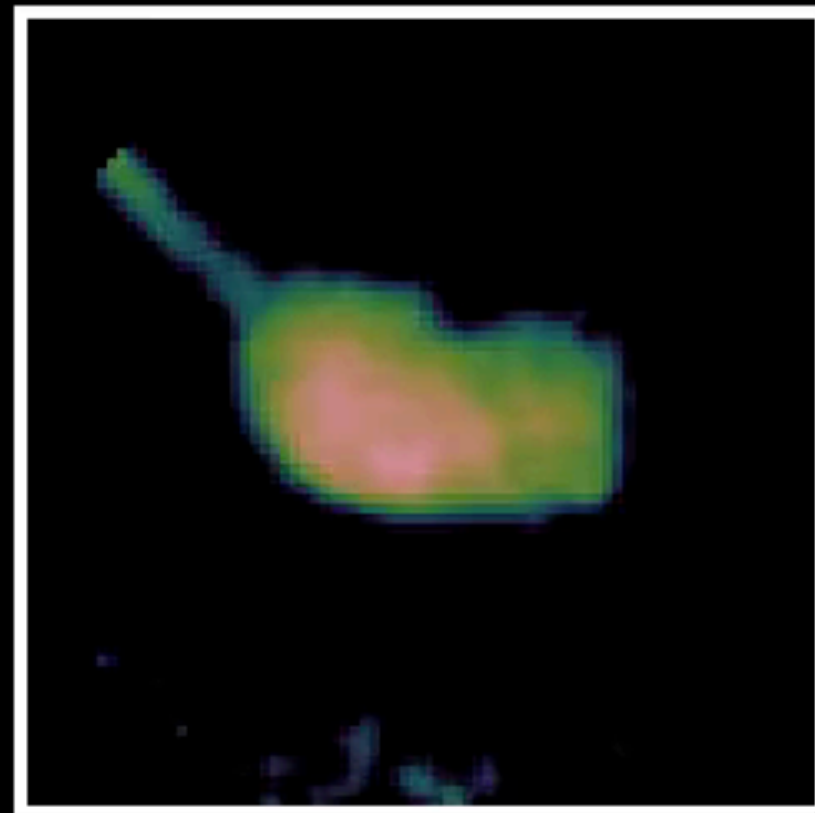- *Dimensionality*: how many latent clusters, factors, etc.?

- *Domain*: are the latent variables discrete, continuous, bounded, sparse, etc.?

- *Dynamics*: how do the latent variables change over time?

- *Dependencies*: how do the latent variables relate to the observed data?

- *Distribution*: do we have prior knowledge about the variables' probability?


- We've already seen some examples in Unit 1!

# Latent variable modeling is all about constraints

**Domain/Dependency/Distribution**

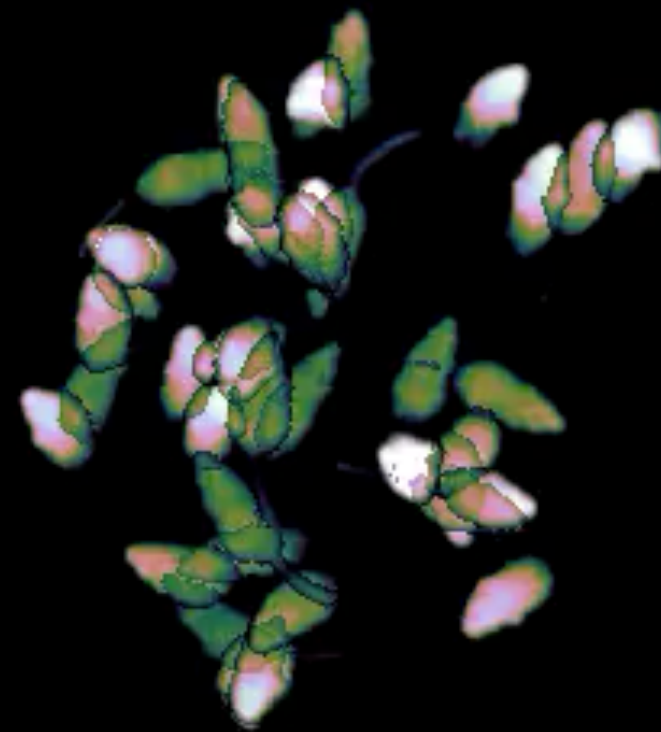| | Continuous<br>Linear<br>Gaussian | Discrete<br>(Gen.) Linear<br>Bernoulli/Poisson/etc. | Nonlinear Observation Models |
|---|---|---|---|
| **Discrete<br>Markovian<br>Categorical** | **HMM**<br>*Rabiner (1989)* | **HMM**<br>*Rabiner (1989)* | **Structured VAE**<br>*Johnson et al (2016)* |
| **Continuous<br>Linear<br>Gaussian** | **LDS**<br>*Kalman (1960)* | **Poisson LDS**<br>*Smith and Brown (2003), Paninski et al (2010)*<br>*Macke et al (2011)* | **Deep PfLDS**<br>*Archer et al (2015); Gao et al (2016)* |
| **Continuous<br>Nonlinear (parametric)<br>Gaussian** | **NLDS, e.g. Hodgkin-Huxley**<br>*Ahrens, Huys, Paninski (2006)*<br>*Huys and Paninski (2009)* | **NLDS, e.g. Hodgkin-Huxley**<br>*Meng, Kramer, Eden (2011)* | **GPSSM, DKF, LFADS, VIND**<br>*Frigola et al (2013), Krishnan et al (2015),*<br>*Sussillo et al (2016), Hernandez et al (2018)* |
| **Mixed<br>Switching Linear** | **SLDS**<br>*Ghahramani and Hinton (1996)*<br>*Murphy (1998)* | **Poisson SLDS**<br>*Petreska et al (2013)* | **Structured VAE**<br>*Johnson et al (2016)* |
| **Mixed<br>Recurrent Linear** | **recurrent/augmented SLDS**<br>*Barber (2006); Pachitariu et al (2014);*<br>*Linderman et al (2017); Nassar et al (2019)* | **rSLDS**<br>*Linderman et al (2017)*<br>*Nassar et al (2019)* | **Structured VAE**<br>*Johnson et al (2016)* |
| **Continuous<br>Nonlinear (smoothing)<br>Gaussian** | **GPFA**<br>*Yu, Cunningham, et al (2009)* | **vLGP**<br>*Zhao and Park (2017)* | **GPLVM**<br>*Lawerence (2005), Wu et al (2017)* |
| **Continuous<br>Nonlinear (nonparametric)<br>Gaussian** | **GPSSM, DKF, LFADS, VIND**<br>*Frigola et al (2013), Krishnan et al (2015),*<br>*Sussillo et al (2016), Hernandez et al (2018)* | **GPSSM, DKF, LFADS, VIND**<br>*Frigola et al (2013), Krishnan et al (2015),*<br>*Sussillo et al (2016), Hernandez et al (2018)* | **GPSSM, DKF, LFADS, VIND**<br>*Frigola et al (2013), Krishnan et al (2015),*<br>*Sussillo et al (2016), Hernandez et al (2018)* |

Dynamics /Domain

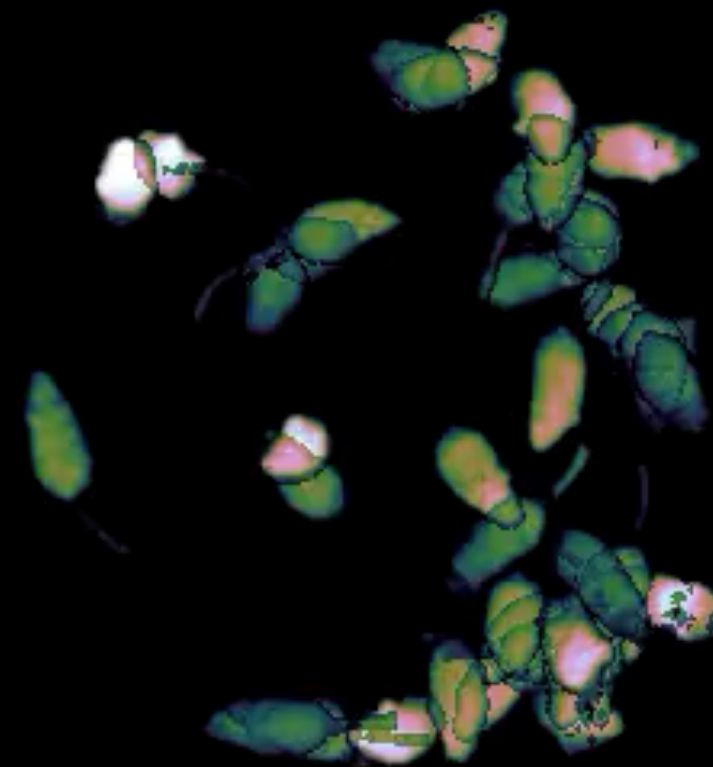# Motivating Example: summarizing videos with behavioral states



Frame 0

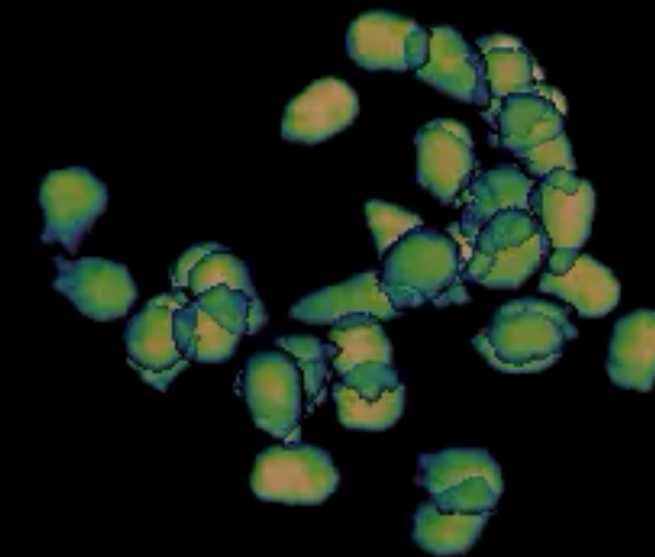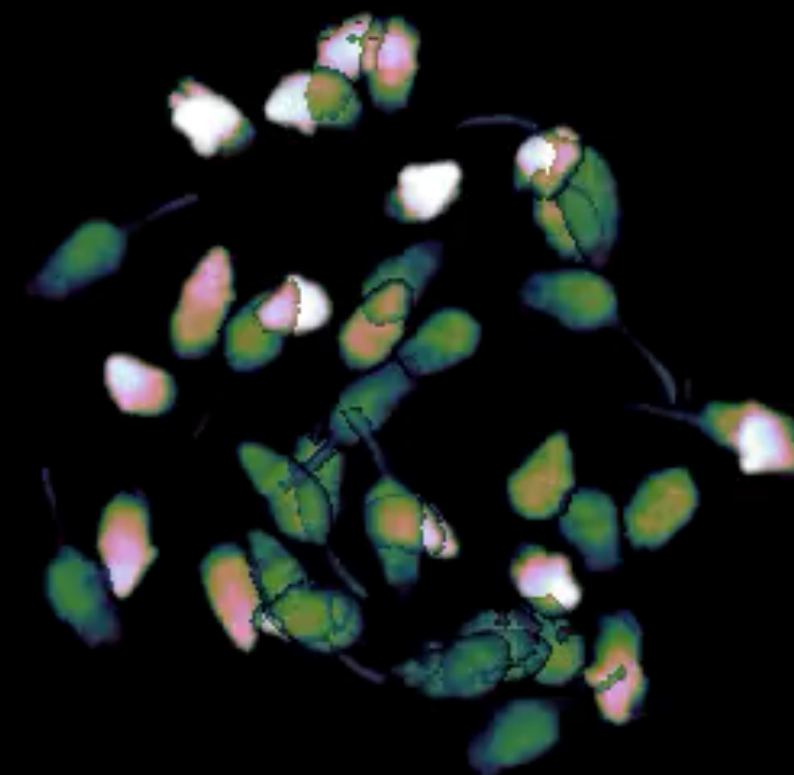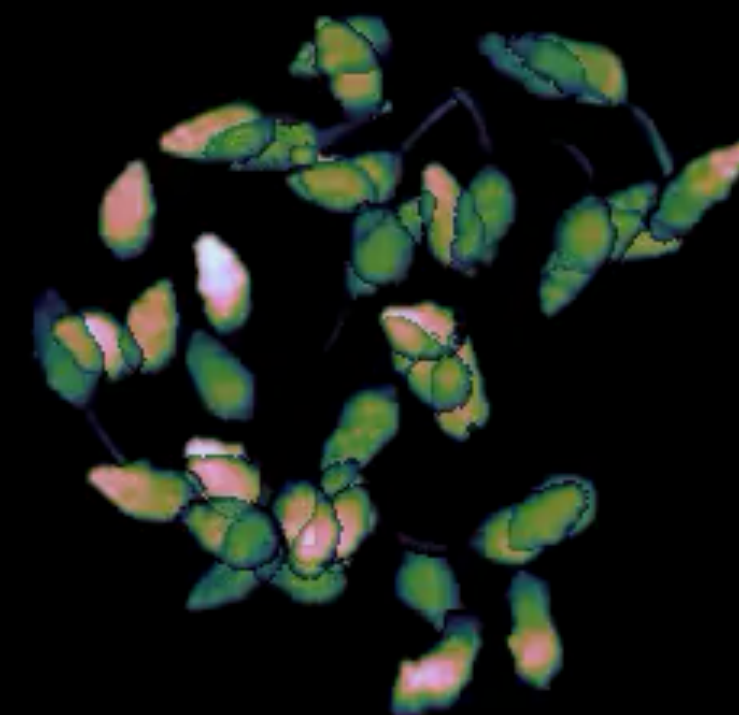# Motivating Example: summarizing videos with behavioral states



Rear down

Walk forward

Grooming

Scrunch

Rear up

Jump

*Wiltschko et al, 2015*

# Formulating as a probabilistic model

- **Variables:** Let,

    - $y_t \in \mathbb{R}^P$ denote the (vectorized) image at time $t$.

    - $z_t \in \{1, \ldots, K\}$ denote the discrete latent state (aka behavioral "syllable") at time $t$.

- **Model:** Assume each time frame is independent and,

$$z_t \sim \mathrm{Cat}(\pi)$$
$$y_t \mid z_t \sim \mathcal{N}(d_{z_t}, R_{z_t})$$

- **Parameters:** Let $\Theta = \pi, \{d_k, R_k\}_{k=1}^K$ denote the parameters,

    - $\pi \in \Delta_K$ is the prior probability of each state

    - $(d_k, R_k) \in \mathbb{R}^P \times \mathbb{R}^{P \times P}$ are the conditional mean and variance of images for discrete state $z_t = k$.

# Bayesian inference in latent variable models

# Bayesian inference in latent variable models
## MAP Estimation

- In Unit 1 we used ***maximum a posteriori* (MAP) estimation** to find,

$$z^{\star}, \Theta^{\star} = \arg\max_{z,\Theta} \log p(y, z, \Theta)$$

- This gave us a **point estimate** of the latent variables $z$ and parameters $\Theta$.

- Point estimates can lead to an **overly optimistic** view of the model.

- Specifically, MAP estimation found **the best assignment**, which may not reflect the **average performance** under the prior $p(z, \Theta)$.

# Bayesian inference in latent variable models
## Integrating over the latent variables

- A more **Bayesian approach** is to **integrate** over the latent variables.

- First, **learn** a point estimate of the parameters,

$$\Theta^{\star} = \arg\max_{\Theta} \log p(y, \Theta)$$

  where $p(y, \Theta) = \int p(y, z, \Theta)\, \mathrm{d}z = \mathbb{E}_{p(z,\Theta)}[p(y \mid z, \Theta)]$ is the **marginal likelihood.**

- Then, **infer** the posterior distribution over latent variables given observed data and parameters,

$$p(z \mid y, \Theta) = \frac{p(y \mid z, \Theta)\, p(z \mid \Theta)\, p(\Theta)}{p(y, \Theta)}$$

- (A "fully Bayesian" approach would integrate over both $z$ and $\Theta$.)

# Bayesian inference in latent variable models
## Maximizing the marginal likelihood

- How to learn the parameters?

- First idea: **gradient ascent**,

$$\nabla_\Theta \log p(y, \Theta) = \frac{\nabla_\Theta p(y, \Theta)}{p(y, \Theta)} = \frac{\int \nabla_\Theta p(y, z, \Theta)\, \mathrm{d}z}{\int p(y, z, \Theta)\, \mathrm{d}z}$$

- Sometimes, these integrals are available in **closed form**.

  - For example, when $z$ **is discrete** the integrals become sums.

- Can we do better?

# Bayesian inference in latent variable models

## Lower bound the marginal likelihood

- Next idea: lower bound the marginal likelihood with a more tractable form,

$$\log p(y, \Theta) = \log \int p(y, z, \Theta) \, \mathrm{d}z$$

# Bayesian inference in latent variable models
## Lower bound the marginal likelihood

- Next idea: lower bound the marginal likelihood with a more tractable form,

$$\log p(y, \Theta) = \log \int p(y, z, \Theta)\, \mathrm{d}z$$

$$= \log \int \frac{q(z)}{q(z)} p(y, z, \Theta)\, \mathrm{d}z \qquad \text{for any distribution } q(z)$$

$$= \log \mathbb{E}_{q(z)} \left[ \frac{p(y, z, \Theta)}{q(z)} \right]$$

$$\geq \mathbb{E}_{q(z)} \left[ \log p(y, z, \Theta) - \log q(z) \right] \qquad \text{by Jensen's inequality}$$

$$\triangleq \mathscr{L}[q, \Theta]$$

- $\mathscr{L}$ is called the **evidence lower bound** or the **ELBO** for short.

# Bayesian inference in latent variable models
## Coordinate ascent on the ELBO

- Update the parameters,

$$\Theta \leftarrow \arg\max_\Theta \mathscr{L}[q, \Theta] = \arg\max_\Theta \mathbb{E}_{q(z)}[\log p(y, z, \Theta)]$$

- Update the distribution on latent variables,

$$q \leftarrow \arg\max_q \mathscr{L}[q, \Theta]$$

$$= \arg\max_q \mathbb{E}_{q(z)}\left[\frac{\log p(y, z, \Theta)}{q(z)}\right]$$

$$= \arg\min_q \mathrm{KL}\left(q(z) \,\|\, p(z \mid y, \Theta)\right)$$

$$= p(z \mid y, \Theta)$$

# Bayesian inference in latent variable models
## The Expectation-Maximization (EM) algorithm

- **M-step**: Maximize the expected log probability

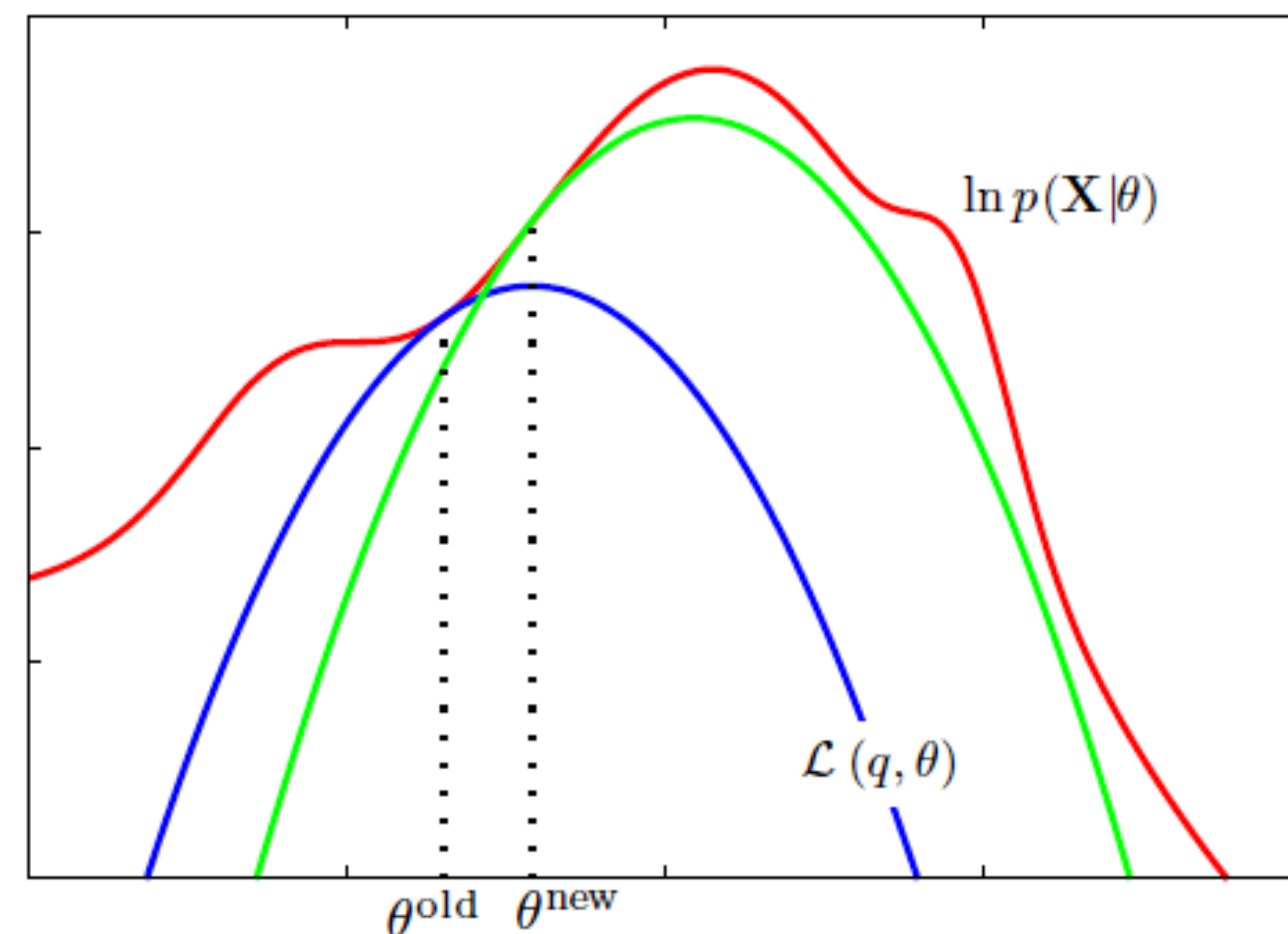$$\Theta \leftarrow = \arg\max_\Theta \mathbb{E}_{q(z)}[\log p(y, z, \Theta)]$$

- **E-step**: Update the posterior over latent variables

$$q \leftarrow p(z \mid y, \Theta)$$

- After each E-step, the **ELBO is tight**:

$$\mathscr{L}[p(z \mid y, \Theta), \Theta] = \mathbb{E}_{p(z|y,\Theta)} \left[ \log \frac{p(y, z, \Theta)}{p(z \mid y, \Theta)} \right]$$

$$= \mathbb{E}_{p(z|y,\Theta)} \left[ \log p(y, \Theta) \right]$$

$$= \log p(y, \Theta)$$

- EM converges to **local optima** of the marginal distribution.



Bishop (2006). Pattern Recognition and Machine Learning, Ch 9.4.

# Bayesian inference in latent variable models
## EM for the Gaussian mixture model

Recall the model,

$$z_t \sim \mathrm{Cat}(\pi)$$
$$y_t \mid z_t \sim \mathcal{N}(d_{z_t}, R_{z_t})$$

with parameters $\Theta = \pi, \{d_k, R_k\}_{k=1}^{K}$.

- **E-step**: Update the posterior over latent variables,

$$q(z_t = k) \leftarrow p(z_t = k \mid y_t, \Theta) \propto \frac{\pi_k \mathcal{N}(y_T \mid d_k, R_k)}{\sum_{j=1}^{K} \pi_j \mathcal{N}(y_t \mid d_j, R_j)}$$

- **M-step**: Update the parameters. Let $N_k = \sum_{t=1}^{T} q(z_t = k)$, then

$$\pi_k \leftarrow \frac{N_k}{T}$$

$$d_k \leftarrow \frac{1}{N_k} \sum_{t=1}^{T} q(z_t = k)\, y_t$$

$$R_k \leftarrow \frac{1}{N_k} \sum_{t=1}^{T} q(z_t = k)\, (y_t - d_k)(y_t - d_k)^{\top}$$

# Conclusion

- Unsupervised models, specifically **latent variable models,** seek simple underlying variables to explain neural or behavioral data.

- LVMs are all about **constraints.**

- MAP estimation, which we used in Unit 1, yielded a point estimate, but can be over-optimistic.

- EM maximizes the marginal likelihood by coordinate ascent on the **latent variable posterior** and the **parameters.**