# Probabilistic models for neural data:
# From single neurons to population dynamics

## NEUROBIO 316QC

Jan Drugowitsch

jan_drugowitsch@hms.harvard.edu

**Session 1**: Course overview & Bayesian recap

# Couse outline

**Aim:** Understand modular structure of probabilistic models

> Framework for thinking about models

> Reveals assumptions

> Supports changing/refining models

This course is *only a starting point!*

**Structure**

> This/next week: recap of/intro to Bayesian statistics

> Future sessions:

>> discussion of paper uses techniques of previous sessions

>> introduction of new techniques

> Between sessions:

>> Reading (Perusall) & preparing discussion (Google Slides)

>> exercises & brief write-up

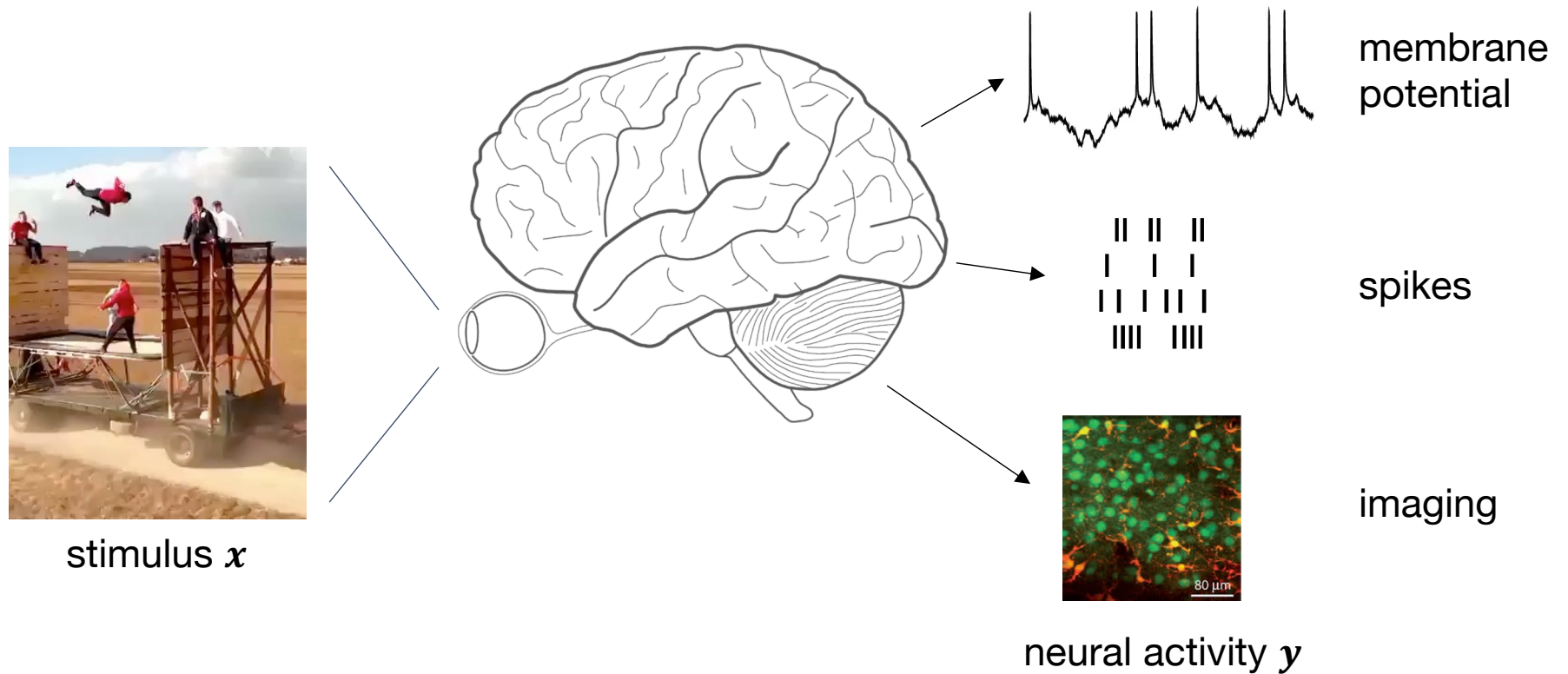>> completing quiz by noon of day of session

# What to expect …

You should **not** expect to

    become an expert in Bayesian modeling;

    understand all the details of the discussed papers;

    design and implement new models from ground up.


Ideally, you would learn to

    understand structure of different Bayesian generative models for neural data,
    the associated graphical models and assumptions;

    read up on new models and understand how they relate to existing models;

    know what you would need to learn (and where you could find information)
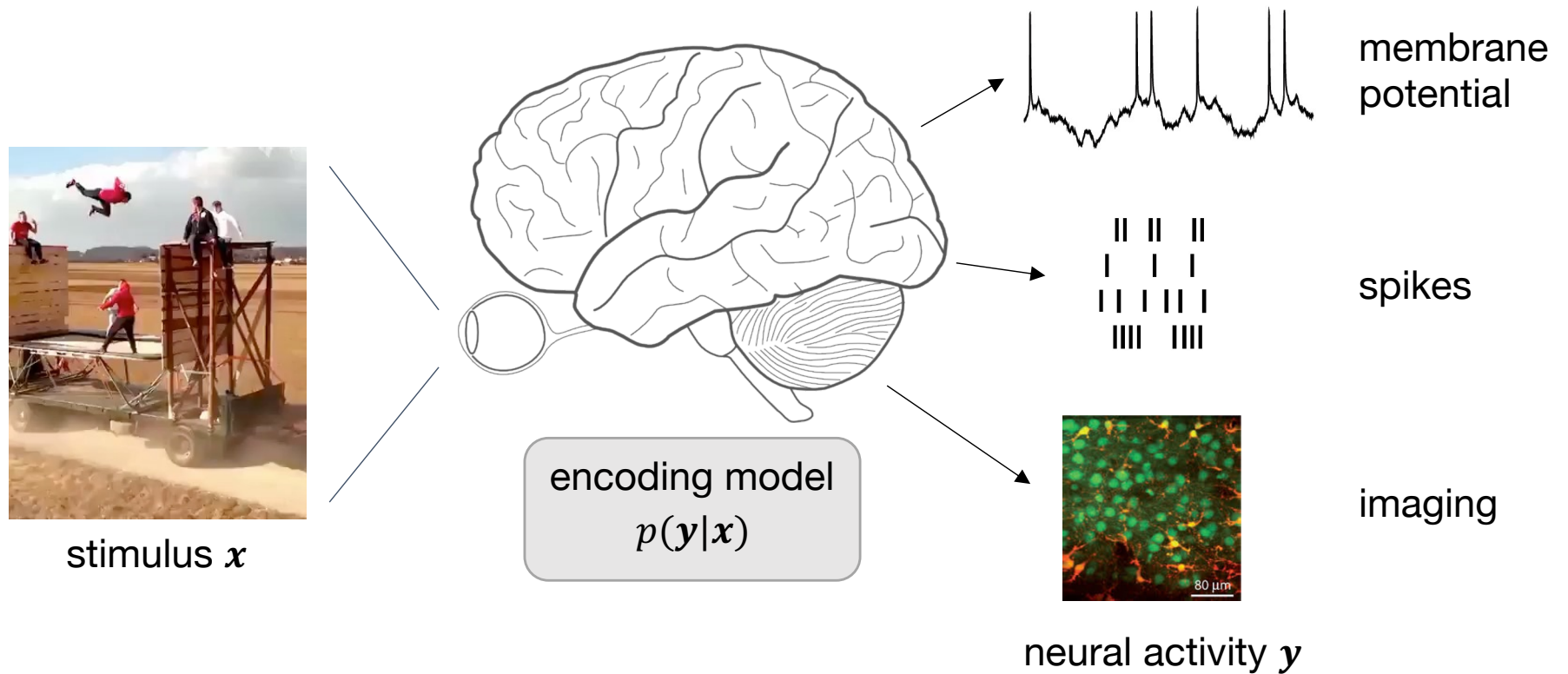    to design and implement new models.

# Why build models?



stimulus $x$

membrane potential

spikes

imaging

neural activity $y$

How are stimuli $x$ encoded in neural activity $y$?

What aspects of neural activity carry information?

What can we/the brain say about the stimulus given neural activity?

(after Jonathan Pillow)

# Why build models?



stimulus $x$

encoding model
$p(y|x)$
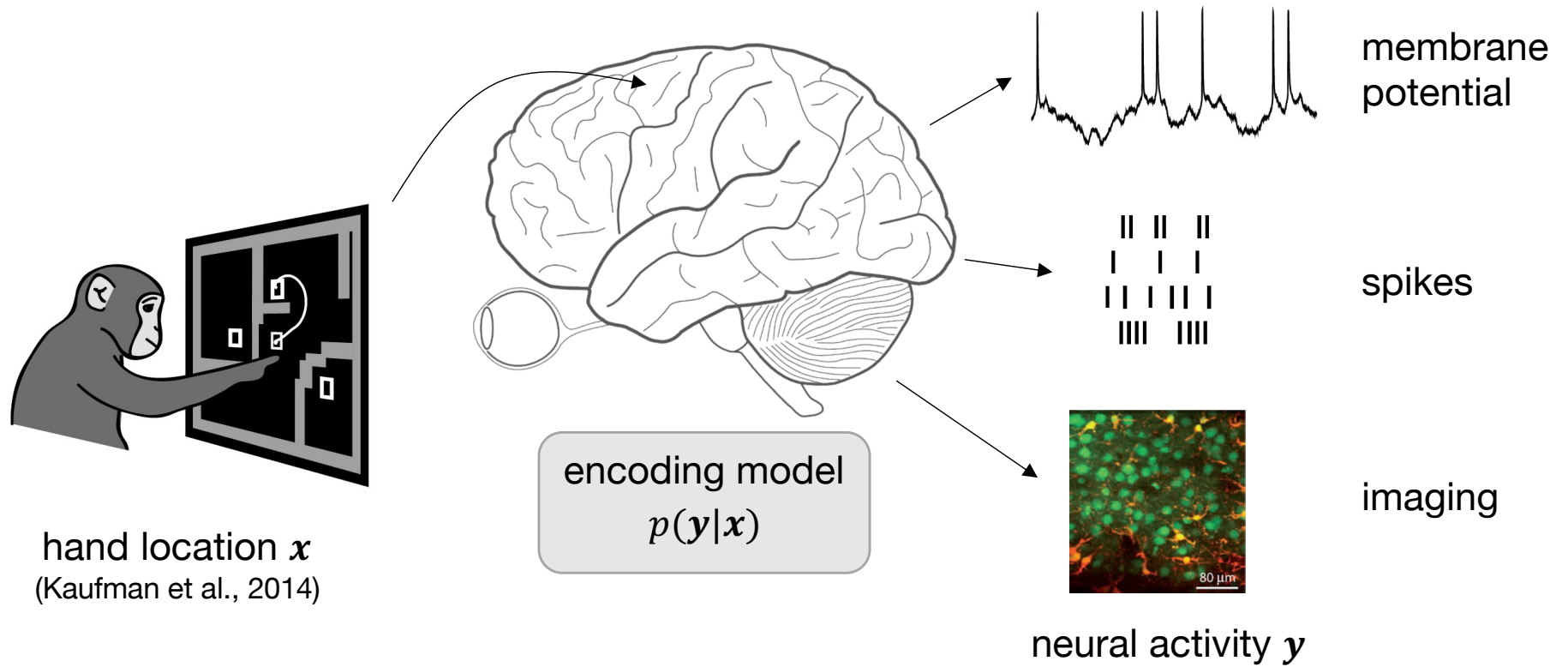
membrane potential

spikes

imaging

neural activity $y$

Build flexible statistical encoding model $p(y|x)$

Quantify information carried in neural responses

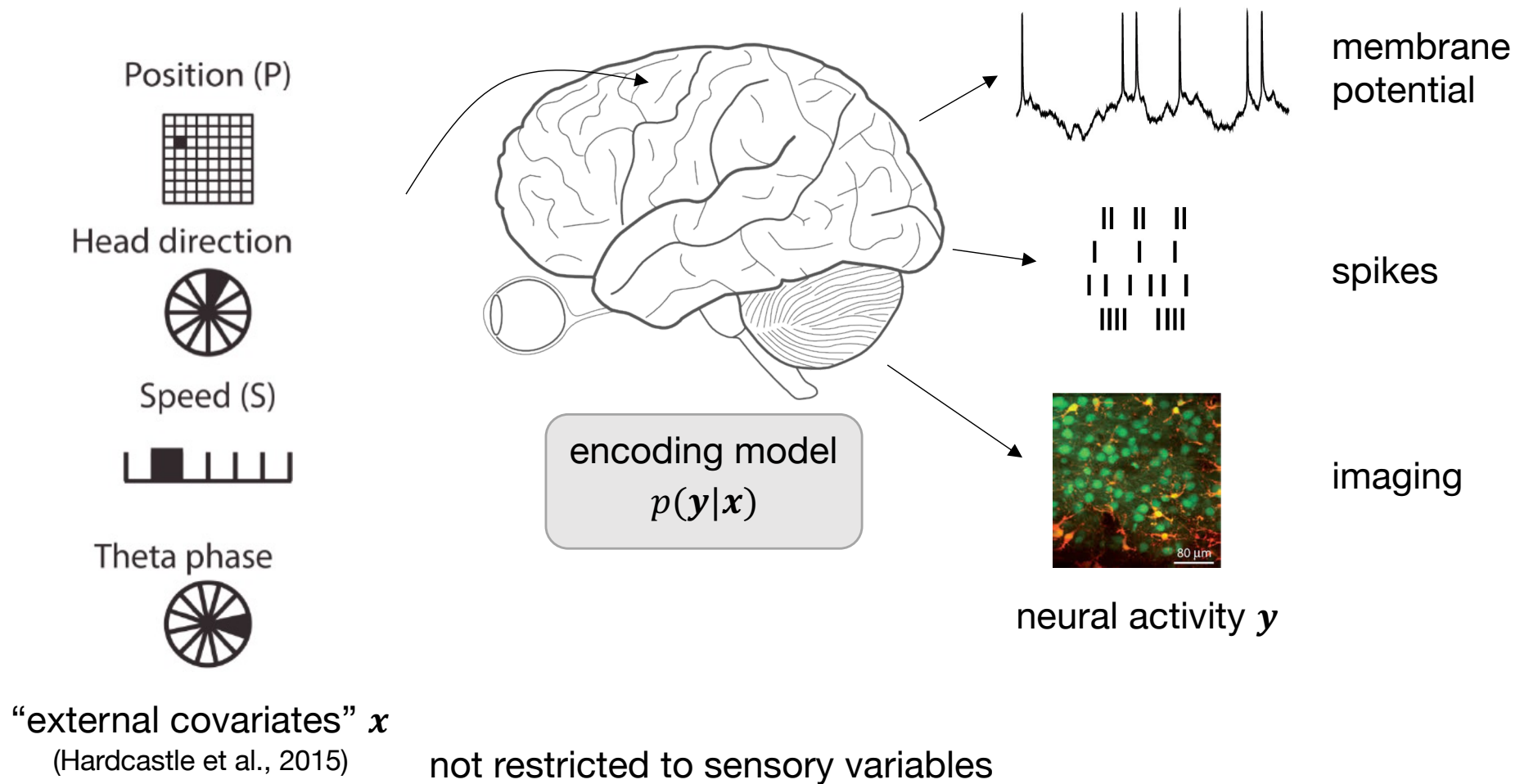Invert encoding model for decoding, $p(x|y)$

(after Jonathan Pillow)
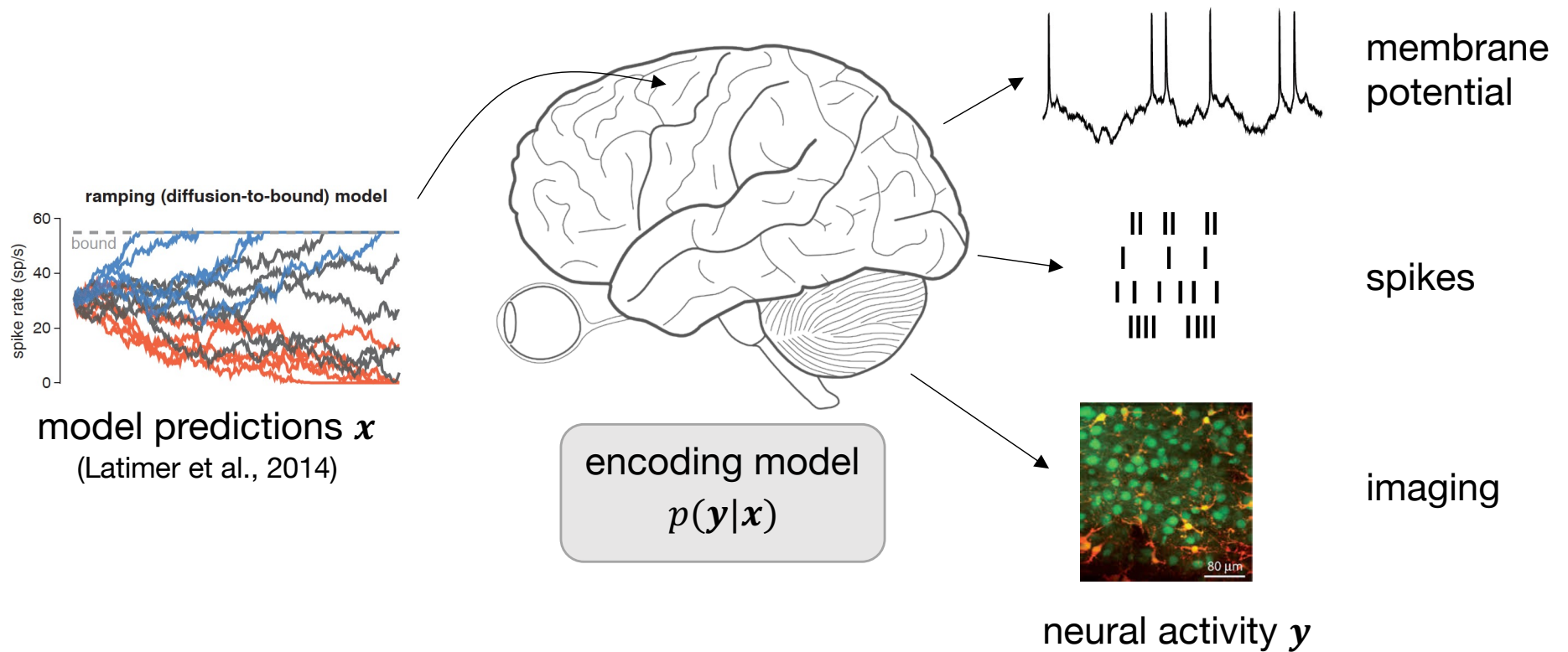
# Why build models?



membrane potential

spikes

imaging

encoding model
$p(\boldsymbol{y}|\boldsymbol{x})$

hand location $\boldsymbol{x}$

(Kaufman et al., 2014)

neural activity $\boldsymbol{y}$

not restricted to sensory variables

(after Jonathan Pillow)

# Why build models?

Position (P)

Head direction

Speed (S)

Theta phase

"external covariates" $\boldsymbol{x}$
(Hardcastle et al., 2015)

encoding model
$p(\boldsymbol{y}|\boldsymbol{x})$

membrane potential

spikes

imaging

neural activity $\boldsymbol{y}$

not restricted to sensory variables

(after Jonathan Pillow)

# Why build models?



ramping (diffusion-to-bound) model

model predictions $x$
(Latimer et al., 2014)

encoding model
$p(y|x)$

membrane
potential

spikes

imaging

neural activity $y$

not restricted to sensory variables

(after Jonathan Pillow)

# Why build models?



latent variables $\mathbf{z}$

membrane potential

spikes

imaging

latent encoding models
$p(\mathbf{y}|\mathbf{z})p(\mathbf{z})$

neural activity $\mathbf{y}$

80 µm

capture hidden structure underlying neural activity

(after Jonathan Pillow)

# Why build models?

latent dynamics

latent variables $z_t$

latent dynamic encoding models
$p(\boldsymbol{y}_t|\boldsymbol{z}_t)p(\boldsymbol{z}_t|\boldsymbol{z}_{t-1})$

membrane potential

spikes

imaging

80 µm

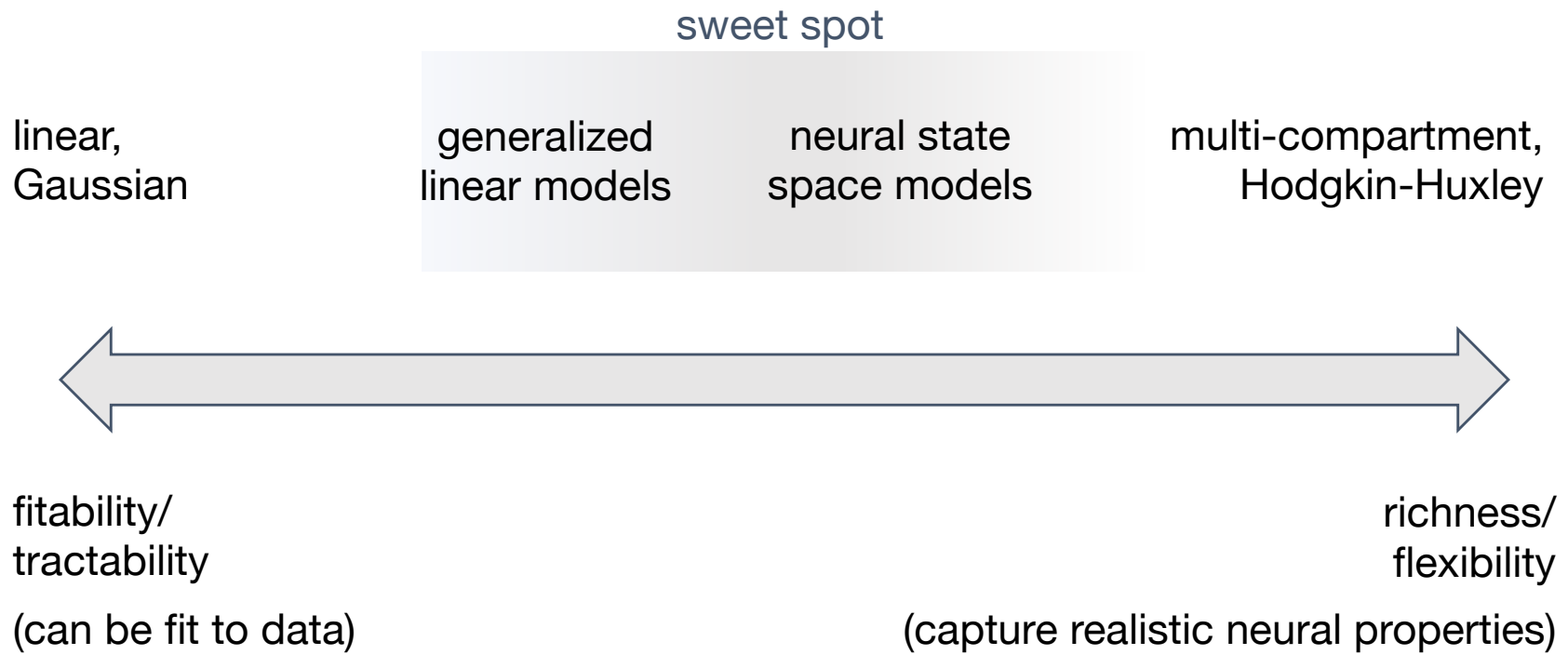neural activity $\boldsymbol{y}$

capture hidden dynamics underlying neural activity
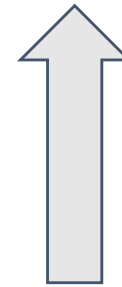
The same can be done for behavior (but not in this course)

(after Jonathan Pillow)

# Model desiderata

sweet spot

| linear,<br>Gaussian | generalized<br>linear models | neural state<br>space models | multi-compartment,<br>Hodgkin-Huxley |

fitability/<br>tractability

(can be fit to data)

richness/<br>flexibility

(capture realistic neural properties)

(after Jonathan Pillow)

# Descriptive statistical models

normative theories
(e.g., efficient coding)

"Why does the code
take this form?"

$p(\boldsymbol{y}|\boldsymbol{x})$

**descriptive statistical models**

"What is the code?"

anatomy, biophysics

"How is it implemented?"

(after Jonathan Pillow)

# Sessions

**Session 1** (today): Bayesian recap

<span style="color:red">Exercise: Bayesian histogram tuning curve fits</span>

**Session 2**: linear models
Topics: linear-Gaussian models, priors as regularizers

**Session 3**: generalized linear models #1
Topics: LNP neurons, single-neuron GLMs, IF neurons

**Session 4**: generalized linear models #2
Topics: GLMs for neural populations, decoding with GLMs

<span style="color:red">Exercise: GLMs</span>

encoding &
decoding

$$p(\boldsymbol{y}|\boldsymbol{x})$$
&
$$p(\boldsymbol{x}|\boldsymbol{y})$$

**Session 5**: Dimensionality reduction
Topics: PCA, probabilistic PCA & Factor Analysis, TCA

latent encoding
$$p(\boldsymbol{x}|\boldsymbol{z})p(\boldsymbol{z})$$

**Session 6**: State space models #1
Topics: Laplace approx., Expectation Maximization, Variational Bayes

**Session 7**: State space models #2
Topics: Gaussian processes

<span style="color:red">Exercise: GPFA? (TBD)</span>

latent dynamic
encoding

$$p(\boldsymbol{x}_t|\boldsymbol{z}_t)p(\boldsymbol{z}_t|\boldsymbol{z}_{t-1})$$

**Session 8**: State space models #3
Topics: Artificial neural networks

**Session 9**: Paper discussion & wrap-up

**Session 1**: Bayesian recap

# Overview

**Probabilities and probabilistic models**

Simple stimulus → response models

Rules of probabilities

Parametric models and their graphical representation

Independent and identically distributed data

**Inference with probabilistic models**

Maximum likelihood estimates

Bayesian inference and its components, generative model inversion

Maximum a-posteriori estimates

Conjugacy and tractability

**Model comparison**

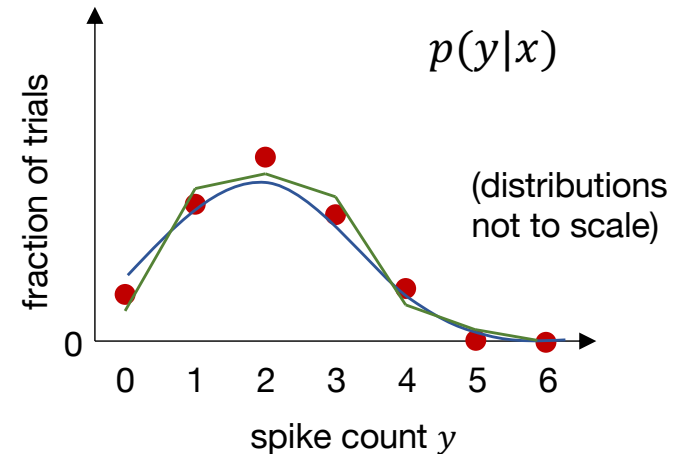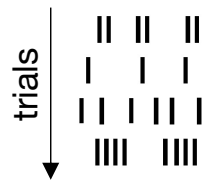Bayesian decision theory

Posterior predictive checks

Bayesian model comparison

# Overview

**Probabilities and probabilistic models**

Simple stimulus → response models

Rules of probabilities

Parametric models and their graphical representation

Independent and identically distributed data

**Inference with probabilistic models**

Maximum likelihood estimates

Bayesian inference and its components, generative model inversion

Maximum a-posteriori estimates

Conjugacy and tractability

**Model comparison**

Bayesian decision theory

Posterior predictive checks

Bayesian model comparison

# Simple stimulus → response models

stimulus $x$    response $y$

trials

$p(y|x)$

fraction of trials

(distributions not to scale)

0

0  1  2  3  4  5  6

spike count $y$

**Directly measure $p(y|x)$?**

Either $x$ or $y$ might be too large/continuous

Cannot extrapolate beyond seen data

**Instead use (parametric) models**, for example

|  | $p(y|x)$ | parameters |
|---|---|---|
| Poisson | $p(y|x) = \text{Pois}(y|\lambda(x))$ | rate $\lambda(x)$ |
| Gaussian | $p(y|x) = \text{N}(y|\mu(x), \sigma^2(x))$ | mean $\mu(x)$, variance $\sigma^2(x)$ |

# Fundamentals of probabilities

**Probability distributions are *functions* that return probabilities**

$p(X = x)$ (short: $p(x)$) returns probability that random variable $X$ takes value $X = x$

**Probability *mass* differs from probability *density***

Discrete $x$ (e.g., spike count, $x \in \{0,1,2, \dots\}$)   $p(x)$ returns *probability mass* ($p(x) \in [0,1]$)

Continuous $x$ (e.g., $\Delta F / F$, $x \in [0, \infty]$)   $p(x)$ returns *probability density* ($p(x) \in [0, \infty]$)



(Bishop, 2006)

mass and density are related

$$p(a \leq x \leq b) = \int_a^b p(x)\mathrm{d}x$$

mass          density

**Probabilities sum to one**

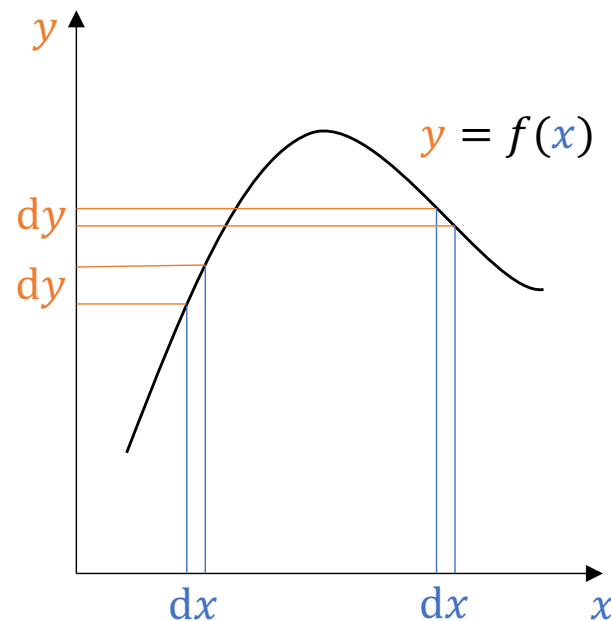Discrete $x$: $\sum_x p(x) = 1$          Continuous $x$: $\int p(x)\mathrm{d}x = 1$

**Probabilities can be defined across multiple random variables**

$p(X = x, Y = y)$ (short: $p(x, y)$) returns *joint probability* that $X = x$ and $Y = y$

# Transformations of random variables

We know $p_x(x)$ and $y = f(x)$

What is $p_y(y)$?



Matching probability mass $\qquad p_x(x)dx = p_y(y)dy \qquad$ s.t. $\qquad p_y(y) = p_x(x)\dfrac{dx}{dy}$

...but ignore the sign of the derivative $\qquad p_x(x) = p_y(y)\left|\dfrac{dy}{dx}\right|$

$$p_y(y) = p_x(X = f^{-1}(y))\left|\frac{dx}{dy}\right| = p_x(X = f^{-1}(y))\left|\frac{1}{f'(x)}\right|$$

For vector-valued $x$ and $y$: $\left|\dfrac{dx}{dy}\right|$ becomes determinant of Jacobian

# Rules of probabilities

**Sum rule** (also called *marginalization*)    $p(y) = \sum_x p(x, y)$        (still a function!)

**Product rule**                $p(x, y) = \underbrace{p(y|x)}p(x)$

conditional probability that $Y = y$ given that $X = x$

$$p(y|x) = \frac{p(x, y)}{p(x)}$$

**Independence**                $p(x, y) = p(x)p(y)$

(what does sum/product rule simplify to?)

**Preview: results in Bayes rule**

$$p(x|y) = \frac{p(y|x)p(x)}{p(y)} = \frac{p(y|x)p(x)}{\sum_x p(y|x)p(x)}$$

# Parametric distributions

**Gaussian distribution**

probability density function (pdf)

$$p(x; \mu, \sigma^2) = \mathrm{N}(x|\mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

$x-$dependent

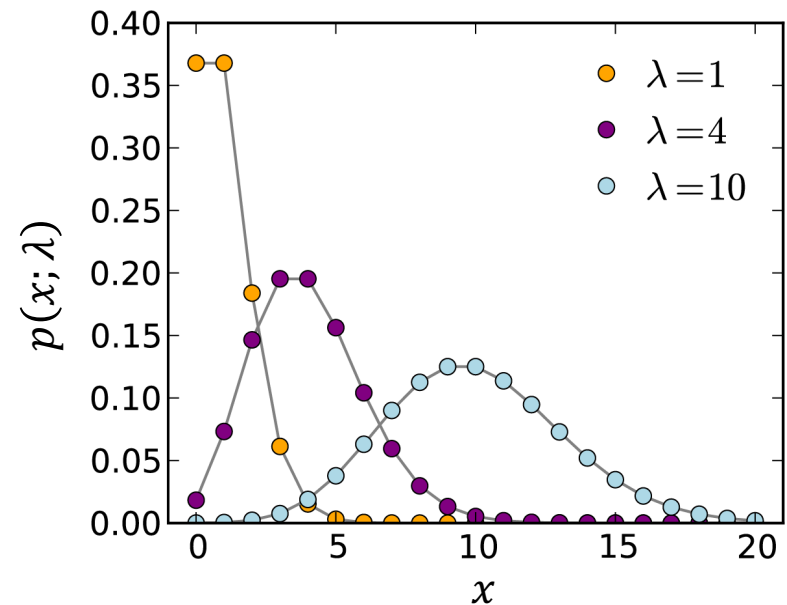mean $\quad \mathrm{E}[x] = \mu$

variance $\quad \mathrm{var}[x] = \sigma^2$



**Poisson distribution**

probability mass function (pmf)

$$p(x; \lambda) = \mathrm{Pois}(x|\lambda) = e^{-\lambda}\frac{\lambda^x}{x!}$$

$x-$dependent

mean $\quad \mathrm{E}[x] = \lambda$

variance $\quad \mathrm{var}[x] = \lambda$

# i.i.d. data & (directed) graphical models

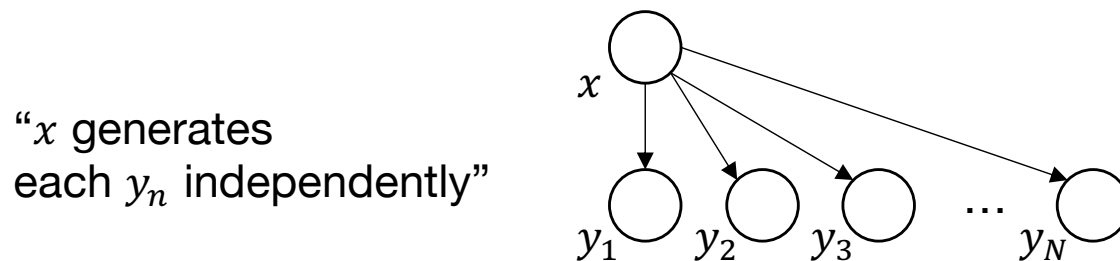***Independent*** and ***identically*** distributed (i.i.d.) data
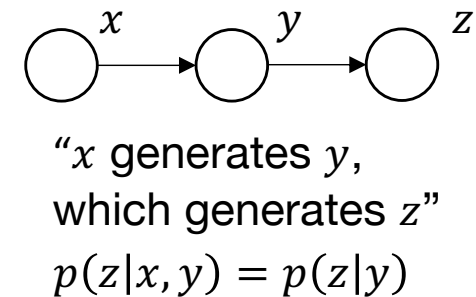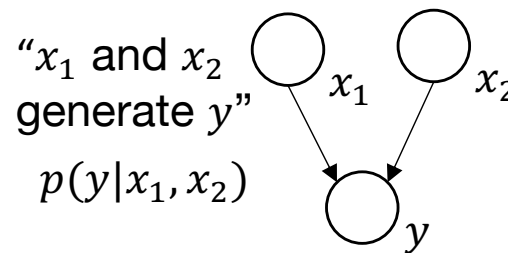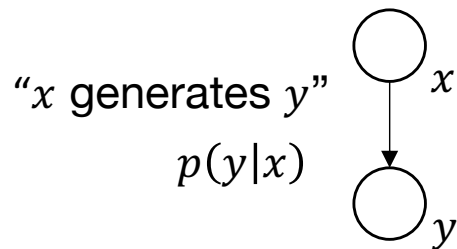
e.g., data $y_{1:N} = y_1, y_2, \ldots, y_N$ from set of trials with same stimulus $x$

*independent*, conditional on stimulus $x$

$$p(y_{1:N}|x) = \prod_{n=1}^{N} p(y_n|x) = \prod_{n=1}^{N} \mathrm{Pois}(y_n|\lambda(x))$$

assume *identical* Poisson "emissions"

**(directed) graphical models**

"$x$ generates $y$"
$p(y|x)$

"$x_1$ and $x_2$ generate $y$"
$p(y|x_1, x_2)$

"$x$ generates $y$, which generates $z$"
$p(z|x, y) = p(z|y)$

"$x$ generates each $y_n$ independently"

$=$

"plate"

# Overview

**Probabilities and probabilistic models**

Simple stimulus → response models

Rules of probabilities

Parametric models and their graphical representation

Independent and identically distributed data

**Inference with probabilistic models**

Maximum likelihood estimates

Bayesian inference and its components, generative model inversion

Maximum a-posteriori estimates

Conjugacy and tractability

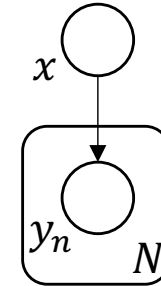**Model comparison**

Bayesian decision theory

Posterior predictive checks

Bayesian model comparison

# maximum likelihood estimates

Generative model of data $y_{1:N}$ in response to (fixed) stimulus $x$

$$p(y_{1:N}|x) = \prod_{n=1}^{N} p(y_n|x) = \prod_{n=1}^{N} \text{Pois}(y_n|\lambda(x))$$

Assume: single rate $\lambda$ for fixed stimulus $x$

**Maximum likelihood: what is the rate $\lambda$ that makes the observed data most likely?**

(Which model parameters make the observed data most likely?)

$$\hat{\lambda}_{\text{ML}} = \text{argmax}_\lambda p(y_{1:N}|\lambda) = \text{argmax}_\lambda \log p(y_{1:N}|\lambda) = \text{argmax}_\lambda \sum_{n=1}^{N} \log \text{Pois}(y_n|\lambda)$$

...

$$= \frac{1}{N} \sum_{n=1}^{N} y_n \qquad \text{the average spike count!}$$

**Pros** Consistent (converges to true $\lambda$)
Efficient (asymptotically no better estimator)
...

**Cons** Noisy for little data
No estimate of uncertainty
...

# Proportionality, $\propto$

$$\overbrace{p(x)}^{\text{normalized}} \underset{x}{\propto} \overbrace{f(x)}^{\text{unnormalized}} \longleftrightarrow p(x) = \frac{1}{Z_p} f(x)$$

"proportional in $x$ to"

independent of $x$

**This works because probability distributions sum/integrate to one!**

$$\int p(x)\,dx = 1 \qquad \frac{1}{Z_p}\int f(x)\,dx = 1 \qquad Z_p = \int f(x)\,dx \qquad p(x) = \frac{1}{\int f(x)\,dx} f(x)$$

**Examples**

$$N(x; \mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}} \underset{x}{\propto} e^{-\frac{(x-\mu)^2}{2\sigma^2}} \qquad\qquad \text{Pois}(x; \lambda) = e^{-\lambda}\frac{\lambda^x}{x!} \underset{\lambda}{\propto} \lambda^x e^{-\lambda}$$
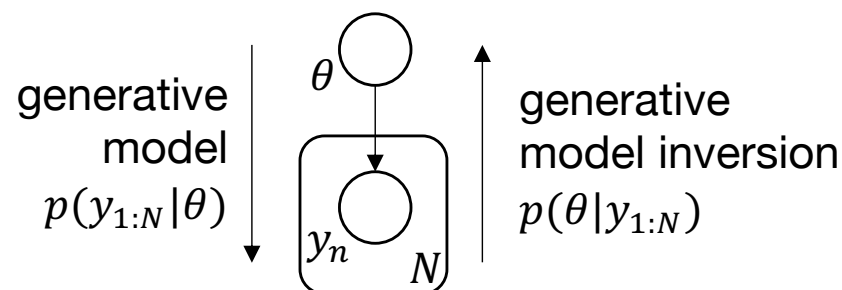
ML inference

$$\text{argmax}_\lambda \sum_{n=1}^{N} \log \text{Pois}(y_n | \lambda) = \text{argmax}_\lambda \sum_{n=1}^{N} \left( \log(\lambda^{y_n} e^{-\lambda}) + \log\frac{1}{y_n!} \right)$$

# Bayesian inference

**Most likely model parameters → *p*(model parameters | data)**

Assume (again): single rate $\lambda$ (model parameter $\theta = \lambda$) for fixed stimulus $x$

$$\underbrace{p(\theta|y_{1:N})}_{\text{posterior}} = \frac{\overbrace{p(y_{1:N}|\theta)}^{\text{likelihood}}\overbrace{p(\theta)}^{\text{prior}}}{\underbrace{p(y_{1:N})}_{\text{marginal likelihood}}} \propto_{\theta} p(y_{1:N}|\theta)p(\theta)$$

generative
model
$p(y_{1:N}|\theta)$

$\theta$

$y_n$

$N$

generative
model inversion
$p(\theta|y_{1:N})$

| | | |
|---|---|---|
| *prior* | $p(\theta)$ | belief about model parameter value(s) before observing data |
| *likelihood* | $p(y_{1:N}|\theta)$ | likelihood of data given model parameter value(s) (function of $\theta$) |
| *posterior* | $p(\theta|y_{1:N})$ | belief about model parameter value(s) after observing data |

| | | |
|---|---|---|
| *marginal likelihood* | $p(y_{1:N}) = \int p(y_{1:N}|\theta)p(\theta)\mathrm{d}\theta$ | likelihood of data under model a.k.a. *model evidence* |

**Beliefs vs. probabilities**

| | |
|---|---|
| Probability | relative frequency of $x$ across "trials" |
| Belief | belief that $x$ is true value within "trial" |

# Maximum a-posteriori inference

Including prior information → *regularizes* the estimate

$$p(\theta|y_{1:N}) \propto_\theta p(y_{1:N}|\theta)p(\theta) \qquad\qquad \log p(\theta|y_{1:N}) = \log p(y_{1:N}|\theta) + \log p(\theta) + \text{const.}$$

MAP estimate:  $\hat{\theta}_{\text{MAP}} = \text{argmax}_\theta(\log p(y_{1:N}|\theta) + \log p(\theta))$

Compare to ML estimate:  $\hat{\theta}_{\text{ML}} = \text{argmax}_\theta \log p(y_{1:N}|\theta)$  (assumes $p(\theta) \propto 1$)
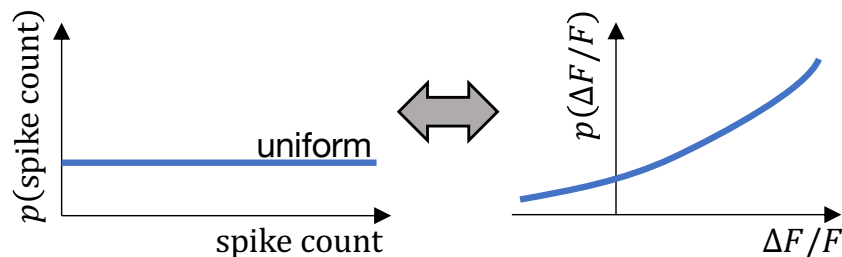
**Benefits of MAP estimates**

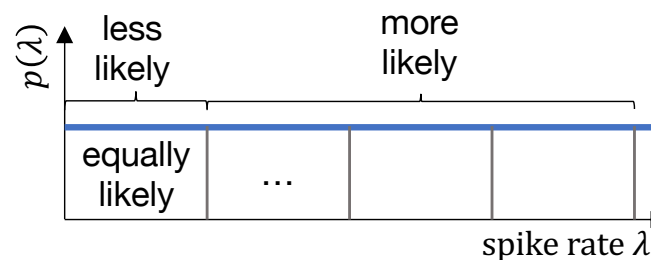Little/uninformative data → minimize impact of noise

Includes prior information

**The myth of "uninformative" priors** (i.e., ML estimates also make assumptions)

e.g., $\Delta F/F = \log(\text{spike count})$       uniform priors, e.g. on spike rate $\lambda$?

# Full posteriors & conjugate priors

**Sequential inference with i.i.d. data**

$$p(\theta|y_1, y_2) \propto p(y_1, y_2|\theta)p(\theta) = p(y_2|\theta)\underbrace{p(y_1|\theta)p(\theta)}$$

$$\propto p(\theta|y_1)$$

Split into two steps

$$p(\theta|y_1) \propto p(y_1|\theta)p(\theta)$$

$$p(\theta|y_1, y_2) \propto p(y_2|\theta)p(\theta|y_1)$$

**Challenge**: "prior" $p(\theta|y_{1:n})$ and "posterior" $p(\theta|y_{1:n+1})$ should have same distribution
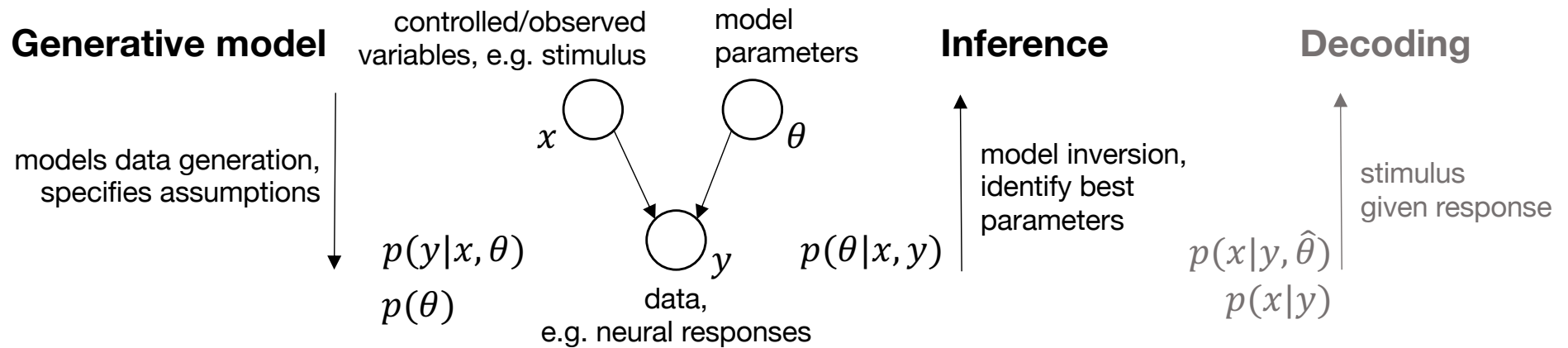
**Solution**: use priors that are *conjugate* to likelihood

| **Examples** | likelihood | parameter(s) | conjugate prior |
|---|---|---|---|
| | $N(y_n|\mu, \sigma^2)$ | $\mu$ | Gaussian |
| | $N(y_n|\mu, \sigma^2)$ | $\mu, \sigma^2$ | Normal-inverse-gamma (NIG) |
| | $\text{Pois}(y_n|\lambda)$ | $\lambda$ | Gamma |
| | | | (see "Conjugate prior" on Wikipedia) |

| **Pros** | mathematical tractability | **Cons** | inflexible |
|---|---|---|---|
| | interpretable parameters | | reflect undesired assumptions |

# Inference summary

**Generative model**   controlled/observed variables, e.g. stimulus   model parameters   **Inference**   **Decoding**

$x$   $\theta$

models data generation, specifies assumptions

model inversion, identify best parameters

stimulus given response

$p(y|x,\theta)$   $y$   $p(\theta|x,y)$   $p(x|y,\hat{\theta})$

$p(\theta)$   data, e.g. neural responses   $p(x|y)$

**Methods of inference**

Full Bayesian   find posterior, (often) requires approximations, or conjugacy

MAP estimates   find most likely parameter posterior, $\hat{\theta}_{\mathrm{MAP}} = \mathrm{argmax}_\theta\, p(\theta|x,y)$

ML estimates   find most likely data likelihood, $\hat{\theta} = \mathrm{argmax}_\theta\, p(y|x,\theta)$
(beware implicit prior assumptions)

# Overview

**Probabilities and probabilistic models**

    Simple stimulus → response models

    Rules of probabilities

    Parametric models and their graphical representation

    Independent and identically distributed data

**Inference with probabilistic models**

    Maximum likelihood estimates

    Bayesian inference and its components, generative model inversion

    Maximum a-posteriori estimates

    Conjugacy and tractability

**Model comparison**

    Bayesian decision theory

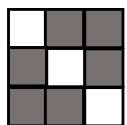    Posterior predictive checks

    Bayesian model comparison

# Bayesian decision theory

Tuning posterior probability distributions into decisions/estimates

$$\underbrace{\hat{\theta}_L}_{} = \mathrm{argmin}_{\hat{\theta}} \int \underbrace{L(\hat{\theta}, \theta)}_{} p(\theta | x, y) \mathrm{d}\theta = \mathrm{argmin}_{\hat{\theta}} \mathrm{E}\left[ L(\hat{\theta}, \theta) | x, y \right]$$

best estimate under loss $L$     loss for choosing $\hat{\theta}$ when $\theta$ is correct

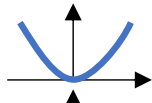**Decision problems** (e.g., $\theta$ is nominal, i.e., unordered and discrete)

    0-1 loss      $L(\hat{\theta}, \theta) = \begin{cases} 0, & \hat{\theta} = \theta, \\ 1, & \hat{\theta} \neq \theta. \end{cases}$     pick most likely $\theta$
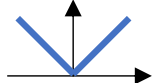
**Estimation problems** (e.g., $\theta$ is ordinal or continuous)

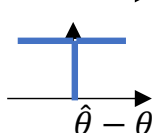squared loss   $L(\hat{\theta}, \theta) = (\hat{\theta} - \theta)^2$      $\hat{\theta}_L = \mathrm{E}[\theta | x, y]$

absolute loss   $L(\hat{\theta}, \theta) = |\hat{\theta} - \theta|$      $\hat{\theta}_L = \mathrm{median\ of\ p}(\theta | x, y)$

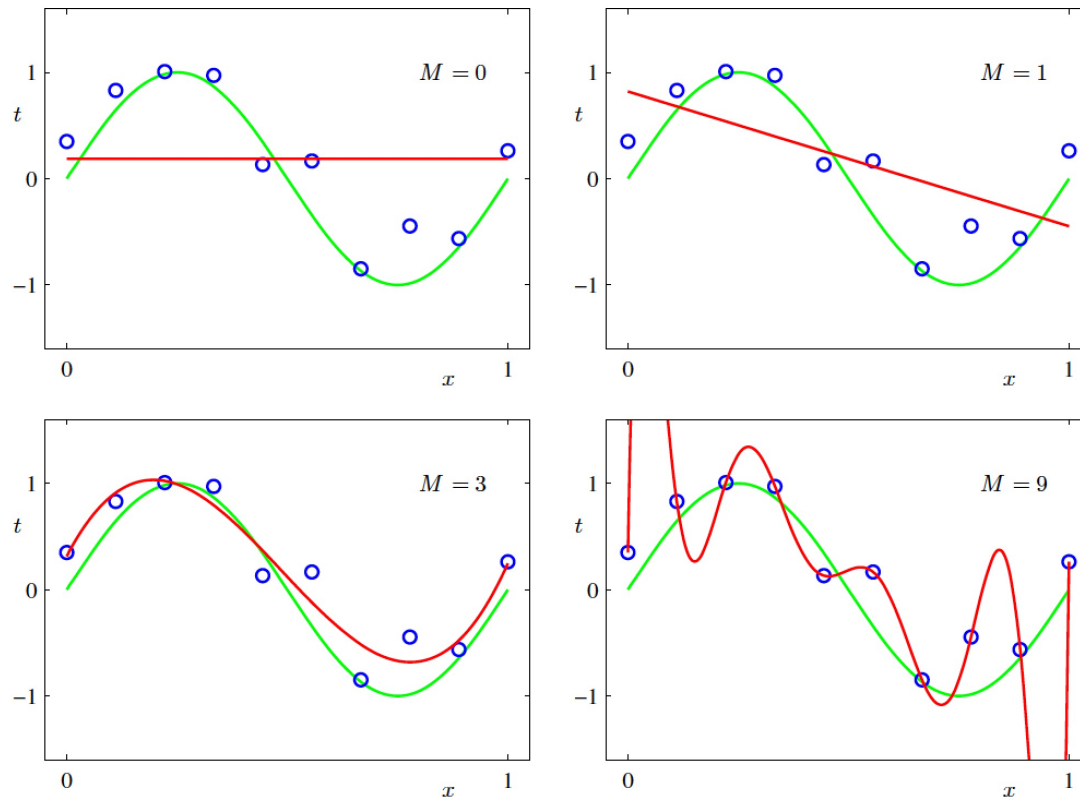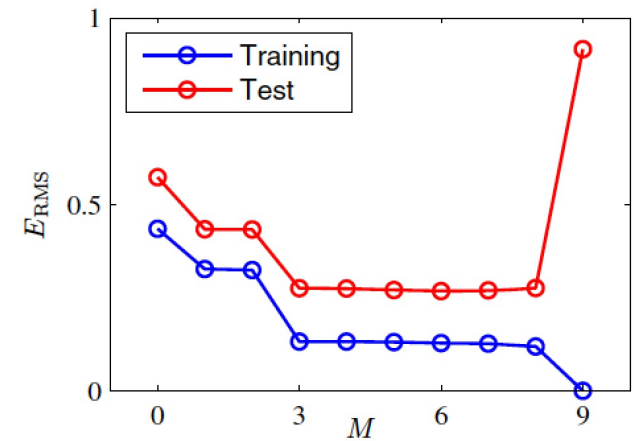notch loss    $L(\hat{\theta}, \theta) = \lim\limits_{c \to 0} \begin{cases} 0, & |\hat{\theta} - \theta| < c, \\ 1, & \text{otherwise.} \end{cases}$   $\hat{\theta}_L = \hat{\theta}_{\mathrm{MAP}} = \mathrm{mode\ of\ } p(\theta | x, y)$

$\hat{\theta} - \theta$

**Why care?** The chosen estimator determines the assumed loss function

# Comparing models



(Bishop, 2006)

too simple ⟷ too complex

unfitting ⟷ overfitting

**Posterior predictive checks** — assess model performance on hold-out dataset (e.g., cross-validation)

**Bayesian model comparison** — use model evidence (marginal likelihood) to reward high data likelihood while penalizing model complexity
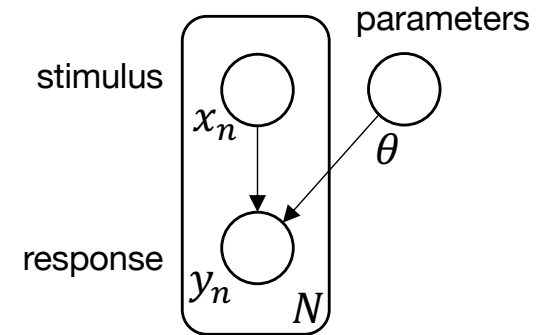
# Posterior predictive checks

Estimate prediction quality by comparing model predictions to hold-out data

**Posterior predictive distribution**

Training data $x_{1:N}, y_{1:N}$; test set instance $\tilde{x}, \tilde{y}$

$$p(\tilde{y}|\tilde{x}, x_{1:N}, y_{1:N}) = \int \overbrace{p(\tilde{y}|\theta, \tilde{x})}^{\substack{\text{likelihood of} \\ \text{test set instance}}} \underbrace{p(\theta|x_{1:N}, y_{1:N})}_{\text{training set posterior}}\mathrm{d}\theta \overset{\substack{\text{ignoring posterior} \\ \text{parameter uncertainty}}}{\approx} p(\tilde{y}|\hat{\theta}, \tilde{x})$$

stimulus $x_n$  parameters $\theta$

response $y_n$  $N$

**Assessing prediction quality** (here across training instances $\tilde{x}_{1:M}, \tilde{y}_{1:M}$)

Choose loss function → measure average test set loss/error

e.g., absolute loss:     $\dfrac{1}{M}\displaystyle\sum_{m=1}^{M}|\tilde{y}_m - \underbrace{\mathrm{median}(\tilde{y}|\tilde{x}_m)}_{\substack{\text{median estimate} \\ \text{compatible with absolute loss}}}|$

... or assess hold-out data log-likelihood  $\log p(\tilde{y}_{1:M}|\tilde{x}_{1:M}) = \displaystyle\sum_{m=1}^{M} \log p(\tilde{y}_m|\tilde{x}_m)$

(requires comparable likelihoods across models)

# Bayesian model comparison

Marginal likelihood (a.k.a. *model evidence*) captures model fit and complexity

$$\underbrace{p(y_{1:N}|x_{1:N}, M_j)}_{\text{model evidence for model } M_j} = \int p(y_{1:N}|x_{1:N}, \theta, M_j)\underbrace{p(\theta|M_j)}_{\text{parameter prior for model } M_j}\,\mathrm{d}\theta$$

conjugacy might make integral tractable

Compare $M_1$ and $M_2$ by (log-)Bayes' factor

Bayes' rule,
assuming uniform model prior, $p(M_j) \propto 1$

$$\log\frac{p(M_1|x_{1:N}, y_{1:N})}{p(M_2|x_{1:N}, y_{1:N})} = \log\frac{p(y_{1:N}|x_{1:N}, M_1)}{p(y_{1:N}|x_{1:N}, M_2)} \overset{?}{\lessgtr} 0$$

See "Bayes factor" on Wikipedia
for guideline values of significant differences

**Comparison to posterior predictive checks**

**Pros**  does not require hold-out data

(sometimes) computationally cheaper

(usually) more sensitive to model details

**Cons**  spurious results for bad models

sensitive to choice of prior

(often) hard/impossible to compute

In general, posterior predictive checks are the safer choice!

# Topics that we won't discuss

**Intractability**

Bayesian inference is in most cases intractable, need to be approximated

Approximations: variational Bayes, Markov Chain Monte Carlo, etc.


**Calibration**

Bayesian inference is sensitive to model misspecification

How to ensure that posterior beliefs correspond to variability across datasets?


**Combining Bayesian inference and deep learning**

Deep neural networks are 'just another function approximator'


**Normative computational models**

Bayesian inference as a model for how brain processes uncertain information

Generates predictions for neural dynamics through encoding/decoding models


…

# Overview

**Probabilities and probabilistic models**

    Simple stimulus → response models

    Rules of probabilities

    Parametric models and their graphical representation

    Independent and identically distributed data

**Inference with probabilistic models**

    Maximum likelihood estimates

    Bayesian inference and its components, generative model inversion

    Maximum a-posteriori estimates
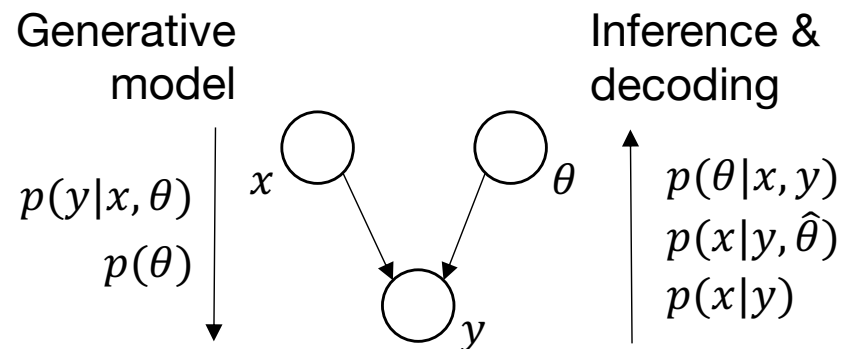
    Conjugacy and tractability

**Model comparison**

    Bayesian decision theory

    Posterior predictive checks

    Bayesian model comparison

# Summary

Generative
model

$p(y|x,\theta)$
$p(\theta)$

$x$        $\theta$

$y$

Inference &
decoding

$p(\theta|x,y)$
$p(x|y,\hat{\theta})$
$p(x|y)$

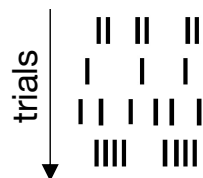| | |
|---:|:---|
| **Handle uncertainty/noise** | Probabilities (and associated rules) |
| **Model structure / independence** | Graphical models |
| **Inference** | Approximate (ML, MAP), or full posteriors |
| **Assess/compare models** | Posterior predictive checks & Bayesian model comp. |

# Exercise

Simple stimulus-response models $p(y|x, \theta)$

stimulus $x$        response $y$



6 neurons
16 different drift directions $x$
40 trials/neuron and drift direction
bin drift directions into {1,2,4,8,16} bins
assume same response model per bin

Poisson likelihood $\text{Pois}(y|\lambda(x))$, fit rate $\lambda(x)$ per stimulus bin
Gaussian likelihood $N(y|\mu(x), \sigma^2(x))$, fit mean $\mu(x)$ and variance $\sigma^2(x)$ per stimulus bin

Compare ML and MAP estimates, impact of prior

Compare models across different bin sizes, and different prior variances

Deliverable: brief write-up
See session notes for instructions

# Until next week

Complete exercise and write-up

Read statistical methods sections (see notes for Session 2)

**Next session**

Discussing the exercise (~15min)

Theory of Gaussians & linear models (remaining time)