

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

NEUROBIO 316QC

Jan Drugowitsch

jan_drugowitsch@hms.harvard.edu

Session 3: Generalized linear models #1

Today

Q&A about previous session

Paper discussion (~1h)

Introducing generalized linear models (remaining time)

Overview

Bernoulli distribution, the exponential family & conjugate priors

Bernoulli distribution as spike generation model

Conjugate priors for the Bernoulli distribution

The exponential family of probability distributions

Linear classification models

Generative vs. discriminative models

Logistic regression: sigmoidal Bernoulli probability

Inhomogeneous Poisson process: sequence of Bernoulli events

Generalized linear models and canonical activation function

Overview

Bernoulli distribution, the exponential family & conjugate priors

Bernoulli distribution as spike generation model

Conjugate priors for the Bernoulli distribution

The exponential family of probability distributions

Linear classification models

Generative vs. discriminative models

Logistic regression: sigmoidal Bernoulli probability

Inhomogeneous Poisson process: sequence of Bernoulli events

Generalized linear models and canonical activation function

The Bernoulli distribution

To spike ($x = 1$) or not to spike ($x = 0$)?

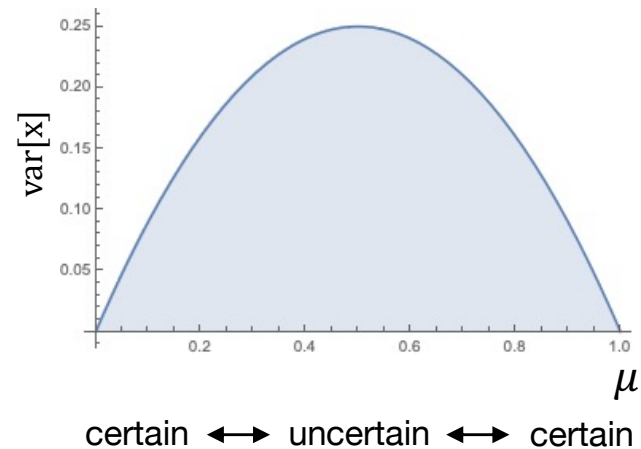
$$\left. \begin{array}{l} p(x = 1|\mu) = \mu \\ p(x = 0|\mu) = 1 - \mu \end{array} \right\} p(x|\mu) = \mu^x (1 - \mu)^{1-x}$$



Moments

$$E[x] = \mu$$

$$\text{var}[x] = \mu(1 - \mu)$$



Conjugate priors

Posterior has same functional form (e.g., Gaussian) as prior

$$\begin{array}{c} \text{same functional form} \\ \swarrow \quad \searrow \\ p(\theta|x) \propto p(x|\theta)p(\theta) \\ \uparrow \\ \text{depends on likelihood} \end{array}$$

Example: Bernoulli distribution

$$\log p(\mu|x) = \log p(x|\mu) + \log p(\mu) + \text{const.}$$

$$= x \log \mu + (1 - x) \log(1 - \mu) + \log p(\mu) + \text{const.}$$

$$= x \log \mu + (1 - x) \log(1 - \mu) + \underbrace{(\alpha - 1) \log \mu + (\beta - 1) \log(1 - \mu)} + \text{const.}$$

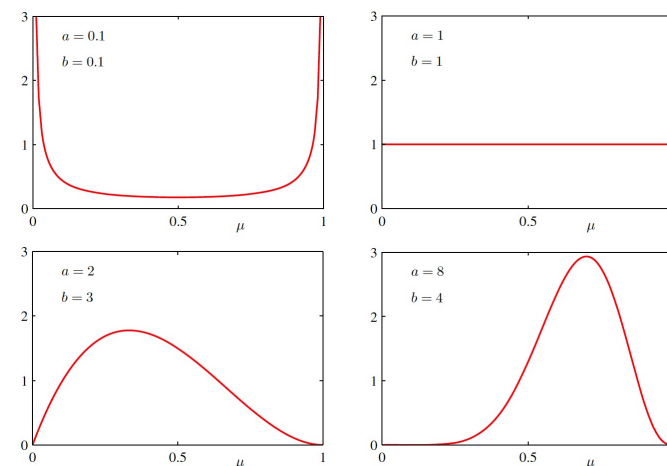
$$\text{Beta distribution prior } B(\mu|\alpha, \beta) = \frac{\Gamma(\alpha+\beta)}{\Gamma(\alpha)\Gamma(\beta)} \mu^{\alpha-1} (1 - \mu)^{\beta-1}$$

$$= \underbrace{(x + \alpha - 1) \log \mu + (1 - x + \beta - 1) \log(1 - \mu)} + \text{const.}$$

$$\text{Beta distribution posterior } B(\mu|\alpha + x, \beta + 1 - x)$$

For multiple observations, x_1, x_2, \dots :

$$\tilde{\alpha} = \alpha + \sum_n x_n \quad \tilde{\beta} = \beta + \sum_n (1 - x_n)$$



The exponential family of distributions

$$p(\mathbf{x}|\boldsymbol{\eta}) = \underbrace{h(\mathbf{x})}_{\text{base measure}} \underbrace{g(\boldsymbol{\eta})}_{\text{normalizer}} \exp(\underbrace{\boldsymbol{\eta}^T \mathbf{u}(\mathbf{x})}_{\text{natural parameters} \rightarrow \text{sufficient statistics}})$$

Base measure: density $p(\mathbf{x}|\boldsymbol{\eta}) = h(\mathbf{x})$ for $\boldsymbol{\eta} = \mathbf{0}$

Normalizer: ensures $\int p(\mathbf{x}|\boldsymbol{\eta}) d\mathbf{x} = 1$

Sufficient statistics: consider i.i.d. data x_1, \dots, x_N

$$p(\mathbf{x}_{1:N}|\boldsymbol{\eta}) \propto_{\boldsymbol{\eta}} g(\boldsymbol{\eta})^N \exp\left(\boldsymbol{\eta}^T \underbrace{\sum_{n=1}^N \mathbf{u}(x_n)}_{\text{summarizes what } \mathbf{x}_{1:N} \text{ tells us about } \boldsymbol{\eta}}\right)$$

summarizes what $\mathbf{x}_{1:N}$ tells us about $\boldsymbol{\eta}$

Fisher-Darmois-Koopman-Pitman theorem:

(finite) sufficient statistics only exist for exponential family

Conjugate prior: follows directly from form of exponential family

Instances: Bernoulli, Beta, Gaussian, Poisson, Laplace, ...

Mixture models are not members of the exponential family

Interim summary: Bernoulli distribution & exponential family

Bernoulli distribution: models binary events (e.g., spike/no spike in small time window)

Conjugate priors: posteriors have same functional form

Including information from likelihood: updating distribution parameters

Beta distribution is conjugate prior for Bernoulli probability

Exponential family: large family of probability distributions

Only family that can efficiently summarize information from data (i.e., finite suff. stats)

Members: Bernoulli Gaussian, Poisson, Exponential, Laplace, ...

Overview

Bernoulli distribution, the exponential family & conjugate priors

Bernoulli distribution as spike generation model

Conjugate priors for the Bernoulli distribution

The exponential family of probability distributions

Linear classification models

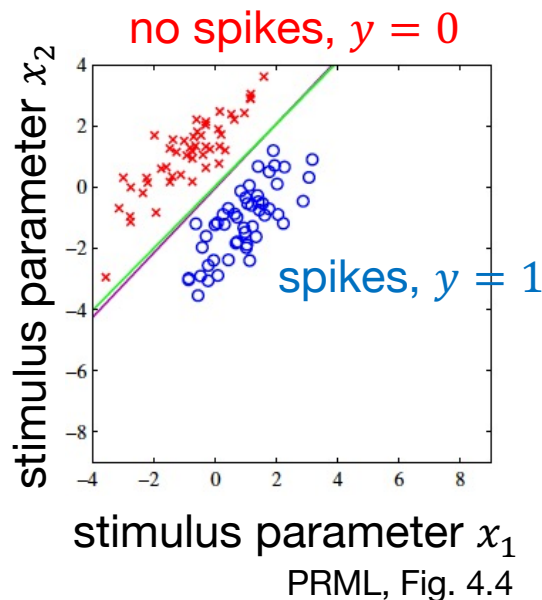
Generative vs. discriminative models

Logistic regression: sigmoidal Bernoulli probability

Inhomogeneous Poisson process: sequence of Bernoulli events

Generalized linear models and canonical activation function

Linear models for classification



Standard linear model?

$$y = \mathbf{w}^T \mathbf{x} + \eta$$

Better: non-linear linear model, s.t. $y \in \{0,1\}$

$$y = 1 \text{ if } \mathbf{w}^T \mathbf{x} + \eta > \theta, y = 0 \text{ otherwise}$$

leads to

$$p(y = 1 | \mathbf{x}, \mathbf{w}) = f(\mathbf{w}^T \mathbf{x})$$

Generative model: specify $p(\mathbf{x}|y, \mathbf{w})$, i.e., distribution over stimuli to generate spike

$$p(y|\mathbf{x}, \mathbf{w}) \propto p(\mathbf{x}|y, \mathbf{w})p(y)$$

provides full model of stimulus distribution (many parameters, might be hard to learn)

Discriminative model: specify $p(y|\mathbf{x}, \mathbf{w})$, i.e., spike probabilities given stimulus

Does not a-priori model stimulus distribution (fewer parameters, might be easier to learn)

They are related by Bayes' rule! (one implicitly determines part of the other)

Discriminative linear models

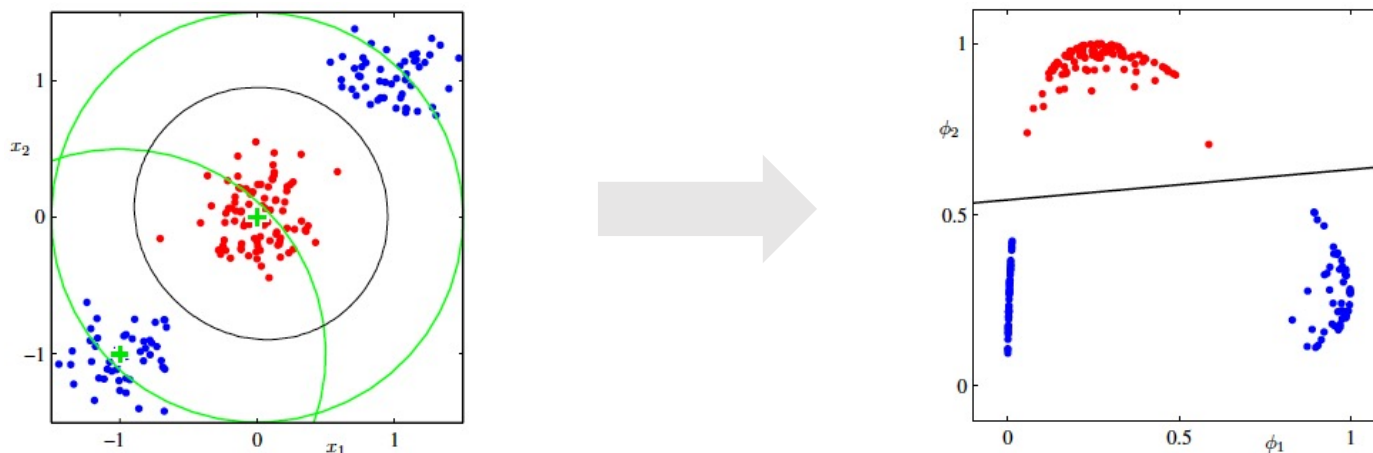
Monotonically increasing activation/link function $f(\cdot)$: $y = 1$ if $\underbrace{\mathbf{w}^T \mathbf{x} + \eta}_{\text{(noisy) linear threshold in } x} > \tilde{\theta}$

(Non-linear) **linear in parameters**, not ‘inputs’

as for linear models, can use non-linear functions $\Phi(x)$ of input

$$p(y = 1 | \mathbf{x}, \mathbf{w}) = f(\mathbf{w}^T \Phi(\mathbf{x}))$$

Non-linear functions can make (linear) non-discriminable data discriminable



Logistic regression

Assumes log-odds linear in \mathbf{x} (or $\Phi(\mathbf{x})$)

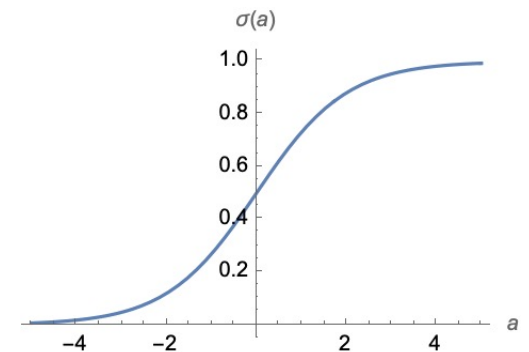
log-posterior ratio

log-likelihood ratio (LLR)

$$\log \frac{p(y=1|\mathbf{x}, \mathbf{w})}{p(y=0|\mathbf{x}, \mathbf{w})} = \log \frac{p(\mathbf{x}|y=1, \mathbf{w})}{p(\mathbf{x}|y=0, \mathbf{w})} + \log \frac{p(y=1)}{p(y=0)} \equiv a(\mathbf{x}, \mathbf{w})$$

leads to

$$p(y=1|\mathbf{x}, \mathbf{w}) = \underbrace{\frac{1}{1 + e^{-a(\mathbf{x}, \mathbf{w})}}}_{\text{logistic sigmoid}} \equiv \sigma(a(\mathbf{x}, \mathbf{w}))$$



Assuming **Gaussian likelihoods**, $p(\mathbf{x}|y) = N(\mathbf{x}|\boldsymbol{\mu}_y, \boldsymbol{\Sigma})$ $a(\mathbf{x}, \mathbf{w}) = \underbrace{-\frac{1}{2}(\boldsymbol{\mu}_0 - \boldsymbol{\mu}_1)^T \boldsymbol{\Sigma}^{-1} \mathbf{x}}_{\mathbf{w}^T} + \text{const.}$

$$p(y=1|\mathbf{x}, \mathbf{w}) = \frac{1}{1 + e^{-\mathbf{w}^T \mathbf{x}}}$$

In general: compatible with any generative model for which $\text{LLR} \sim \mathbf{w}^T \mathbf{x} + \text{const.}$

Moving to spike trains

Logistic regression: for individual spikes in small time bins (spike or no spike?)

Assume spike train $y_{1:N}$ with i.i.d. Bernoulli probabilities, $p(y_n | \mathbf{x}_n, \mathbf{w}) = \lambda(\mathbf{x}_n, \mathbf{w}) \delta t$

(very)small time bin

$$\begin{aligned} \log p(y_{1:N} | \mathbf{x}_{1:N}) &= \sum_n y_n \log \lambda(\mathbf{x}_n, \mathbf{w}) \delta t + \sum_n (1 - y_n) \log(1 - \lambda(\mathbf{x}_n, \mathbf{w}) \delta t) \quad \downarrow \log(1 + a\delta t) \approx a\delta t \\ &\approx \sum_n y_n \log \lambda(\mathbf{x}_n, \mathbf{w}) \delta t - \sum_n (1 - y_n) \lambda(\mathbf{x}_n, \mathbf{w}) \delta t \\ &= \sum_n y_n (\log \lambda(\mathbf{x}_n, \mathbf{w}) \delta t + \lambda(\mathbf{x}_n, \mathbf{w}) \delta t) - \sum_n \lambda(\mathbf{x}_n, \mathbf{w}) \delta t \quad \downarrow |\log \lambda \delta t| \gg |\lambda \delta t| \\ &\approx \sum_n y_n \log \lambda(\mathbf{x}_n, \mathbf{w}) \delta t - \sum_n \lambda(\mathbf{x}_n, \mathbf{w}) \delta t - \sum_n \log y_n! \quad \downarrow \log y_n! = 0 \\ &= \log \underbrace{\prod_n \text{Pois}(y_n | \lambda(\mathbf{x}_n, \mathbf{w}) \delta t)} \end{aligned}$$

time-discretized inhomogeneous Poisson process

Properties of the (in)homogeneous Poisson process

Derivable from i.i.d. Bernoulli probabilities:

spiking only depends on instantaneous rate λ , otherwise independent across time

More generally: spike counts of any two non-overlapping time periods are independent

Spike count within any time interval is distributed by a Poisson distribution

Within any small time-interval δt with instantaneous rate λ_t :

Spike count N satisfies $E[N] = \lambda_t \delta t$ and $\text{var}[N] = \lambda_t \delta t$; Fano factor = $\frac{\text{var}[N]}{E[N]} = 1$

Special case: homogeneous Poisson process ($\lambda_t = \text{const.}$)

Within any time interval Δ : spike count N has distribution $N \sim \text{Pois}(\lambda\Delta)$

Expected number of spikes in that interval are $E[N] = \lambda\Delta$

Inter-spike interval have exponential distribution

The exponential family and canonical activation functions

Linear regression: Gaussian noise + identity activation function, $f(x) = x$

$$\log p(y_{1:N} | \mathbf{x}_{1:N}, \mathbf{w}) \propto - \sum_n (\mathbf{w}^T \Phi(\mathbf{x}_n) - y_n)^2$$

gradient for single n :

$$\nabla_{\mathbf{w}} \log p(y_n | \mathbf{x}_n, \mathbf{w}) \propto -(\underbrace{\mathbf{w}^T \Phi(\mathbf{x}_n)}_{\text{estimate}} - \underbrace{y_n}_{\text{target}}) \Phi(\mathbf{x}_n)$$

Logistic regression: Bernoulli noise + sigmoidal activation function, $f(x) = \sigma(x)$

$$\log p(y_{1:N} | \mathbf{x}_{1:N}, \mathbf{w}) = \sum_n y_n \log \sigma(\mathbf{w}^T \Phi(\mathbf{x}_n)) + (1 - y_n) \log (1 - \sigma(\mathbf{w}^T \Phi(\mathbf{x}_n)))$$

gradient for single n :

$$\nabla_{\mathbf{w}} \log p(y_n | \mathbf{x}_n, \mathbf{w}) = -(\underbrace{\sigma(\mathbf{w}^T \Phi(\mathbf{x}_n))}_{\text{estimate}} - \underbrace{y_n}_{\text{target}}) \Phi(\mathbf{x}_n)$$

Poisson regression: Poisson noise + exponential activation function, $f(x) = e^x$

Same gradient form: supports the same Iterated Recursive Least Squares (IRLS) algorithm

Generally: likelihood $p(y_n | \mathbf{x}_n, \mathbf{w})$ determines *canonical* activation function (RPML 4.3.6)

Interim summary

Generative and discriminative linear models are related by Bayes' rule

Discriminative logistic regression assumes generative log-likelihood ratio (LLR) linear in \mathbf{w}

Inhomogeneous Poisson processes = Bernoulli spikes with probability rate $\times \delta t$

Assume spikes independent across time (conditional on rate)

Exponential family distributions have canonical activation function that support IRLS

Poisson regression has exponential activation function

Summary

Bernoulli distribution: simple model for instantaneous spiking

Assuming i.i.d. Bernoulli spikes in small time windows δt :

Inhomogeneous Poisson process, instantaneous rate = spike probability / δt

Exponential family: large family of probability distributions with convenient properties

- conjugate prior: easily parameter updates
- ability to summarize data in finite sufficient statistics for parameter inference

Linear models for classification

Generative vs. discriminative models

Logistic regression: discriminative model assuming LLR linear in $\Phi(x)$

Generalized linear models

Likelihood function & canonical activation function

Includes linear regression, logistic regression, Poisson regression, etc.

Until next week

Read paper and prepare presentation (see notes for Session 4)

Read statistical methods sections (see notes for Session 4)

Next session

Q&A for previous session (~15min)

Paper discussions (~1h)

Generalized linear models, graphical models & state space models (~30min)

