

The Poisson transform for unnormalised statistical models

Nicolas Chopin (ENSAE)
joint work with Simon Barthelmé (CNRS, Gipsa-LAB)

Part I

Unnormalised statistical models

Unnormalised statistical models

- ▶ “Unnormalised” statistical models: models with an intractable normalisation constant in the likelihood.
- ▶ Example: Ising model for binary vectors $\mathbf{y} \in \{0, 1\}^m$

$$p(\mathbf{y}|\mathbf{a}, \mathbf{Q}) \propto \exp(\mathbf{a}^t \mathbf{y} + \mathbf{y}^t \mathbf{Q} \mathbf{y})$$

- ▶ Very popular in Machine Learning, Computer Vision (deep learning), neuroscience.
- ▶ Creates computational difficulties (“*doubly intractable problems*” in Bayesian context).

Unnormalised *sequential* models

- ▶ Markov sequence $\mathbf{y}_0, \dots, \mathbf{y}_n$ where the transition kernel is defined up to a constant.
- ▶ Example: sequential Ising

$$p(\mathbf{y}_t | \mathbf{y}_{t-1}, \mathbf{a}, \mathbf{Q}, \mathbf{R}) \propto \exp(\mathbf{a}^t \mathbf{y}_t + \mathbf{y}_t \mathbf{Q} \mathbf{y}_t + \mathbf{y}_t \mathbf{R} \mathbf{y}_{t-1})$$

- ▶ Nastier than IID version: n normalisation constants missing.

Current strategies for inference

- ▶ Classical estimation: MCMC-MLE, contrastive divergence (Bengio and Delalleau, 2009), noise-contrastive divergence (Gutmann and Hyvärinen, 2012).
- ▶ Bayesian: exchange algorithm (Murray et al., 2006), ABC, russian roulette (Girolami et al., 2013).
- ▶ I do not know of methods for *sequential* unnormalised models.

Our contribution

- ▶ Poisson transform shows you can treat the missing normalisation constant as just another parameter. Gives you an alternative likelihood function.
- ▶ Applies to sequential problems as well.
- ▶ Noise-contrastive divergence is an approximation of the Poisson transform and we can now extend it to the sequential setting.
- ▶ Sequential estimation can be turned into a *semiparametric logistic regression* problem.

Part II

The Poisson transform

Poisson point processes

- ▶ Poisson processes are distributions over countable subsets of a domain Ω (e.g., $\Omega = \mathbb{R}$ for a temporal point process).
- ▶ Let S be a realisation from a PP. For all (measurable) $\mathcal{A} \subseteq \Omega$, the number of points of S in \mathcal{A} follows a Poisson distribution with parameter

$$\lambda_{\mathcal{A}} = \mathbb{E}(|S \cap \mathcal{A}|) = \int_{\mathcal{A}} \lambda(\mathbf{y}) \, d\mathbf{y}$$

where $\lambda(\mathbf{y})$ is the intensity function.

Poisson point processes (II)

Let's assume that $\int_{\Omega} \lambda(\mathbf{y}) \, d\mathbf{y} < \infty$, then

- ▶ The cardinal of S is Poisson, with parameter $\int_{\Omega} \lambda(\mathbf{y}) \, d\mathbf{y} < \infty$;
- ▶ conditional on $|S| = k$, the elements of S are IID with density

$$\propto \exp\{\lambda(\mathbf{y})\}.$$

Likelihood of a Poisson process

$$\log p(S|\lambda) = \sum_{\mathbf{y}_i \in S} \log \lambda(\mathbf{y}_i) - \int_{\Omega} \lambda(\mathbf{y}) d\mathbf{y}$$

The Poisson transform

- ▶ Generalisation of the Poisson-Multinomial transform (Baker, 1994)
- ▶ For estimation purposes, you can treat IID data in just about any space as coming from a Poisson process.
- ▶ New likelihood function: no loss of information, one extra latent parameter.

Theorem statement (I)

Data: $\mathbf{y}_1, \dots, \mathbf{y}_n \in \Omega$, density $p(\mathbf{y}|\boldsymbol{\theta}) \propto \exp\{f_{\boldsymbol{\theta}}(\mathbf{y})\}$, so log-likelihood is

$$\mathcal{L}(\boldsymbol{\theta}) = \sum_{i=1}^n f_{\boldsymbol{\theta}}(\mathbf{y}_i) - n \log \int_{\Omega} \exp\{f_{\boldsymbol{\theta}}(\mathbf{y})\} d\mathbf{y}.$$

Theorem statement (I)

Data: $\mathbf{y}_1, \dots, \mathbf{y}_n \in \Omega$, density $p(\mathbf{y}|\boldsymbol{\theta}) \propto \exp\{f_{\boldsymbol{\theta}}(\mathbf{y})\}$, so log-likelihood is

$$\mathcal{L}(\boldsymbol{\theta}) = \sum_{i=1}^n f_{\boldsymbol{\theta}}(\mathbf{y}_i) - n \log \int_{\Omega} \exp\{f_{\boldsymbol{\theta}}(\mathbf{y})\} d\mathbf{y}.$$

Poisson log-likelihood:

$$\mathcal{M}(\boldsymbol{\theta}, \nu) = \sum_{i=1}^n \{f_{\boldsymbol{\theta}}(\mathbf{y}_i) + \nu\} - n \int_{\Omega} \exp\{f_{\boldsymbol{\theta}}(\mathbf{y}) + \nu\} d\mathbf{y}$$

i.e. log-likelihood of a PP with intensity $\lambda(\mathbf{y}) = f_{\boldsymbol{\theta}}(\mathbf{y}) + \nu$.

Theorem statement (II)

Theorem

Let $\theta^* = \underset{\theta \in \Theta}{\operatorname{argmax}} \mathcal{L}(\theta)$ and $(\tilde{\theta}, \nu^*) = \underset{\theta \in \Theta, \nu \in \mathbb{R}}{\operatorname{argmax}} \mathcal{M}(\theta, \nu)$. Then

$\tilde{\theta} = \theta^*$ and $\nu^* = -\log(\int \exp\{f_{\theta^*}(\mathbf{y})\} d\mathbf{y})$.

In other words, the MLE can be computed by maximising $\mathcal{M}(\theta, \nu)$ in both variables. There is no loss of information. Also, asymptotic confidence intervals for θ are the same. The latent variable ν “estimates” the normalisation constant.

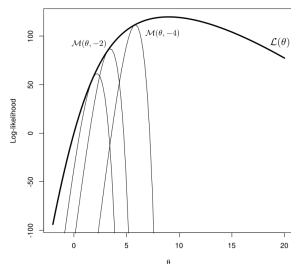
For fixed θ , maximise $\mathcal{M}(\theta, \nu)$ wrt ν leads to:

$$\nu^*(\theta) = -\log \left(\int \exp \{f_{\theta^*}(\mathbf{y})\} d\mathbf{y} \right)$$

and

$$\mathcal{M}(\theta, \nu^*(\theta)) = \mathcal{L}(\theta) - n.$$

Poisson vs. standard likelihood



Running example: truncated exponential distribution:

$$y \in [0, 1], p(y|\theta) \propto \exp(\theta y)$$

Extension to sequential models

The same logic can be applied to sequential models:

$$p_{\theta}(\mathbf{y}_t | \mathbf{y}_{t-1}) \propto \exp \{f_{\theta}(\mathbf{y}_t, \mathbf{y}_{t-1})\}$$

We will apply the Poisson transform to each conditional distribution.

Extension to sequential models

- ▶ Original log-likelihood of sequence:

$$\mathcal{L}(\theta) = \sum_{t=1}^n \left[f_{\theta}(\mathbf{y}_t; \mathbf{y}_{t-1}) - \log \left(\int_{\Omega} \exp \{ f_{\theta}(\mathbf{y}; \mathbf{y}_{t-1}) \} d\mathbf{y} \right) \right]$$

Extension to sequential models

- ▶ Original log-likelihood of sequence:

$$\mathcal{L}(\boldsymbol{\theta}) = \sum_{t=1}^n \left[f_{\boldsymbol{\theta}}(\mathbf{y}_t; \mathbf{y}_{t-1}) - \log \left(\int_{\Omega} \exp \{ f_{\boldsymbol{\theta}}(\mathbf{y}; \mathbf{y}_{t-1}) \} d\mathbf{y} \right) \right]$$

- ▶ Poisson-transformed log-likelihood:

$$\mathcal{M}(\boldsymbol{\theta}, \boldsymbol{\nu}) = \sum_{t=1}^n \{ f_{\boldsymbol{\theta}}(\mathbf{y}_t; \mathbf{y}_{t-1}) + \nu_{t-1} \} - \int_{\Omega} \sum_{t=1}^n \exp \{ f_{\boldsymbol{\theta}}(\mathbf{y}; \mathbf{y}_{t-1}) + \nu_{t-1} \} d\mathbf{y}$$

We have introduced one latent variable ν_t per observation. Sum of integrals becomes integral of a sum.

Extension to sequential models

Maximising the Poisson-transformed likelihood wrt $(\boldsymbol{\theta}, \boldsymbol{\nu})$, gives the MLE for $\boldsymbol{\theta}$, and

$$\nu_{t-1}^*(\boldsymbol{\theta}^*) = -\log \left(\int \exp \{f_{\boldsymbol{\theta}^*}(\mathbf{y}; \mathbf{y}_{t-1})\} d\mathbf{y} \right),$$

i.e. minus the log-marginalisation constant for the conditional

$$p(\mathbf{y}|\mathbf{y}_{t-1}, \boldsymbol{\theta}^*) \propto \exp \{f_{\boldsymbol{\theta}^*}(\mathbf{y}; \mathbf{y}_{t-1})\}.$$

From parametric to semi-parametric inference

The value of the latent variables at the mode are a function of \mathbf{y}_{t-1} :

$$\nu_{t-1}^*(\boldsymbol{\theta}^*) = -\log \left(\int \exp \{f_{\boldsymbol{\theta}^*}(\mathbf{y}; \mathbf{y}_{t-1})\} d\mathbf{y} \right) = \chi(\mathbf{y}_{t-1}).$$

If $\mathbf{y}_t, \mathbf{y}_{t'}$ are close, ν_t, ν_{t-1} should be close as well, i.e., $\chi(\mathbf{y})$ is (hopefully) smooth.

From parametric to semi-parametric inference

The value of the latent variables at the mode are a function of \mathbf{y}_{t-1} :

$$\nu_{t-1}^*(\boldsymbol{\theta}^*) = -\log \left(\int \exp \{f_{\boldsymbol{\theta}^*}(\mathbf{y}; \mathbf{y}_{t-1})\} d\mathbf{y} \right) = \chi(\mathbf{y}_{t-1}).$$

If $\mathbf{y}_t, \mathbf{y}_{t'}$ are close, ν_t, ν_{t-1} should be close as well, i.e., $\chi(\mathbf{y})$ is (hopefully) smooth.

⇒ Do inference over χ : e.g. if you have n points but χ is well captured by a spline basis with $k \ll n$ components, use spline basis instead. Poisson likelihood becomes:

$$\begin{aligned} \mathcal{M}(\boldsymbol{\theta}, \chi) &= \sum_{t=1}^n \{f_{\boldsymbol{\theta}}(\mathbf{y}_t; \mathbf{y}_{t-1}) + \chi(\mathbf{y}_{t-1})\} \\ &\quad - \int_{\Omega} \sum_{t=1}^n \exp \{f_{\boldsymbol{\theta}}(\mathbf{y}; \mathbf{y}_{t-1}) + \chi(\mathbf{y}_{t-1})\} d\mathbf{y} \end{aligned}$$

Using the Poisson transform in practice

Back to the IID case: Poisson-transformed likelihood still involves an intractable integral

$$\mathcal{M}(\boldsymbol{\theta}, \nu) = \sum_{i=1}^n \{f_{\boldsymbol{\theta}}(\mathbf{y}_i) + \nu\} - n \int_{\Omega} \exp \{f_{\boldsymbol{\theta}}(\mathbf{y}) + \nu\} d\mathbf{y}$$

which we need to approximate.

Using the Poisson transform in practice

Back to the IID case: Poisson-transformed likelihood still involves an intractable integral

$$\mathcal{M}(\boldsymbol{\theta}, \nu) = \sum_{i=1}^n \{f_{\boldsymbol{\theta}}(\mathbf{y}_i) + \nu\} - n \int_{\Omega} \exp \{f_{\boldsymbol{\theta}}(\mathbf{y}) + \nu\} d\mathbf{y}$$

which we need to approximate.

Several ways, but an interesting one is to go through logistic regression.

Stochastic gradient descent

Before we go to logistic regression, note that another approach would be to use Monte Carlo (importance sampling) to obtain an *unbiased* estimate of the gradient:

$$\frac{1}{n} \frac{\partial}{\partial \boldsymbol{\theta}} \mathcal{M}(\boldsymbol{\theta}, \nu) = \frac{1}{n} \sum_{i=1}^n \frac{\partial}{\partial \boldsymbol{\theta}} f(\mathbf{y}_i; \boldsymbol{\theta}) - \int_{\Omega} \frac{\partial}{\partial \boldsymbol{\theta}} f(\mathbf{y}; \boldsymbol{\theta}) \exp(f_{\boldsymbol{\theta}}(\mathbf{y}) + \nu) d\mathbf{y}$$

$$\frac{1}{n} \frac{\partial}{\partial \nu} \mathcal{M}(\boldsymbol{\theta}, \nu) = 1 - \int_{\Omega} \exp(f_{\boldsymbol{\theta}}(\mathbf{y}) + \nu) d\mathbf{y}$$

The we could use SGD (stochastic gradient descent) to maximise $\mathcal{M}(\boldsymbol{\theta}, \nu)$.

Part III

The logistic trick & noise-contrastive divergence

The logistic trick

- ▶ Idea: reduce an estimation problem to a classification problem.
- ▶ Several versions:
 - ▶ Logistic regression for density estimation: Hastie et al. (2003), intensity estimation: Baddeley et al. (2010).
 - ▶ Logistic regression for normalisation constants: Geyer (1994).
 - ▶ Logistic regression for estimation in unnormalised models: Gutmann and Hyvärinen (2012).
- ▶ The last one is called “noise-contrastive divergence” by the authors.

The logistic trick

We have n random points from distributions $p(y)$ and n points from $q(y)$. We note $z_i = 1$ if the i -th point is from p , $z_i = 0$ otherwise. Logistic regression models the log-odds ratio:

$$\eta(y) = \log \frac{p(z = 1|y)}{p(z = 0|y)}.$$

We have that:

$$\eta(y) = \log \frac{p(y)}{q(y)}$$

\Rightarrow provided $q(y)$ is known, we can first estimate η (doing some form of logistic regression), and then recover $p(y)$ from $\eta(y)$.

From the logistic trick to noise-contrastive divergence

If we have a *normalised* model $p_{\theta}(y)$ then we can run a logistic regression with the following model for the log-odds:

$$\eta(y; \theta) = \log p_{\theta}(y) - \log q(y).$$

From the logistic trick to noise-contrastive divergence

If we have a *normalised* model $p_{\theta}(y)$ then we can run a logistic regression with the following model for the log-odds:

$$\eta(y; \theta) = \log p_{\theta}(y) - \log q(y).$$

If the model is *unnormalised*, $p_{\theta}(y) \propto \exp\{f_{\theta}(y)\}$, we introduce an intercept in the logistic regression

$$\eta(y; \theta) = f_{\theta}(y) + \nu - \log q(y).$$

This is the *noise-contrastive divergence* (NCD) technique of Gutmann and Hyvärinen (2012).

Truncated exponential

Recall the truncated exponential model:

$$p(y|\theta) \propto \exp(\theta y)$$

We produce reference samples from $U(0, 1)$, so that the logistic model for NCD is just:

$$\eta(y; \theta) = \theta y + \nu$$

Fitting in R:

```
m <- glm(z~y+offset(logratio),data=df,family=binomial)
```

Summary

- ▶ Logistic trick: get a logistic classifier to discriminate true data from random reference data (from a known distribution). It implicitly learns a model for the true data
- ▶ NCD: in unnormalised models, introduce an intercept for the missing normalisation constant
- ▶ Our interpretation: NCD is an approximation of the Poisson-transformed likelihood

NCD approximates the Poisson transform

- ▶ In NCD, you can introduce as many reference points (points simulated from q) as you like.
- ▶ Parametrise the log-odds by

$$\eta(\mathbf{y}) = f_{\boldsymbol{\theta}}(\mathbf{y}) + \nu + \log \frac{n}{m} - \log q(\mathbf{y})$$

where m is the number of reference points.

- ▶ Theorem: as $m \rightarrow +\infty$, the logistic log-likelihood $\mathcal{R}^m(\boldsymbol{\theta}, \nu)$ tends to the Poisson log-likelihood $\mathcal{M}(\boldsymbol{\theta}, \nu)$ (pointwise).

NCD approximates the Poisson transform

To sum up: take your true n datapoints, add m random reference datapoints, and estimate the model

$$p_{\theta}(\mathbf{y}|\theta) \propto \exp \{f_{\theta}(\mathbf{y})\}$$

using a logistic regression with log-odds

$$\eta(\mathbf{y}) = f_{\theta}(\mathbf{y}) + \nu + \log \frac{n}{m} - \log q(\mathbf{y})$$

The intercept will be used to estimate the missing normalisation constant. The technique is effectively a practical way of approximating a Poisson-transformed likelihood.

NCD for sequential models

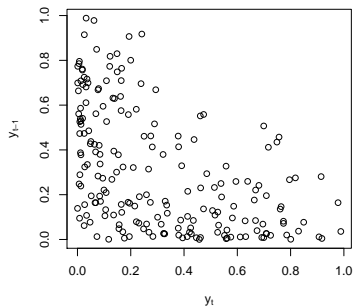
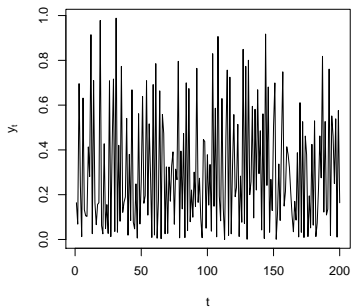
The relationship between the Poisson transform and NCD shows directly how to adapt NCD to sequential models: apply NCD to each conditional distribution (the transition kernels)

- ▶ Reference density $q(\mathbf{y})$ becomes a reference kernel $q(\mathbf{y}_t|\mathbf{y}_{t-1})$
- ▶ Include an intercept ν_t per conditional distribution $p(\mathbf{y}_t|\mathbf{y}_{t-1}, \boldsymbol{\theta})$

Truncated exponential, revisited

We turn our previous example into a Markov chain:

$$p(y_t | y_{t-1}, \theta) \propto \exp(\theta y_t y_{t-1})$$



Truncated exponential, revisited

Consider the NCD approximation for *fixed* y_{t-1} . The model for the log-odds will take the form:

$$\eta(y_t) = \theta y_t y_{t-1} + \nu_{t-1} + \log \frac{n}{m} - \log q(y_t | y_{t-1})$$

This leads to a linear logistic regression with $y_t y_{t-1}$ as a covariate.

Parametric vs. semi-parametric model

It is wasteful to fit a separate intercept per time-point. As in the semi-parametric version of the Poisson transform, we can use:

$$\eta(y_t) = \theta y_t y_{t-1} + \chi(y_{t-1}) + \log \frac{n}{n_r} - \log q(y_t | y_{t-1})$$

where $\chi(y_{t-1})$ will be fitted using splines.

In practice (I)

Positive examples are given by:

Value at time $t - 1$	Value at time t	Label
y_1	y_2	1
y_2	y_3	1
\vdots	\vdots	\vdots
y_{n-1}	y_n	1

While negative examples are given by:

Value at time $t - 1$	Value at time t	Label
y_1	r_2	0
y_2	r_3	0
\vdots	\vdots	\vdots
y_{n-1}	r_n	0

In practice (II)

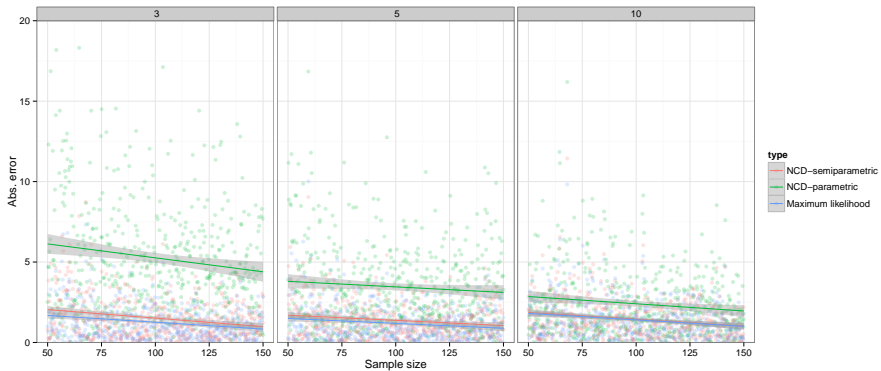
We can fit the (semi-parametric) model via:

```
m <- gam(label ~ I(y_t*y_tminusone)+s(y_tminusone), data=
```

The fully parametric model corresponds to:

```
m <- gam(label ~ I(y_t*y_tminusone)+as.factor(y_tminusone)
```


Parametric vs. semi-parametric model



Part IV

Application: LATKES

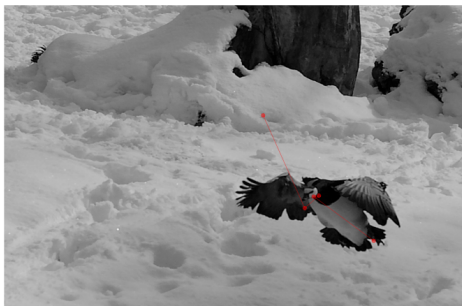


Figure : A sequence of eye movements extracted from the dataset of Kienzle et al. (2009). Fixation locations are in red and successive locations are linked by a straight line.

Eye movements recorded while 14 subjects were exploring a set of photographs (Fig. 1); each contributing between 600 and 2,000 datapoints.

LATKES: Log-Additive Transition Kernels. A class of spatial Markov chain models, with applications to eye movement data:

$$p(y_t | y_{t-1}, \dots, y_{t-k}) \propto \exp \left\{ \sum \beta_i v_i(y_t) + g(y_t, y_{t-1}, \dots, y_{t-k}) \right\}$$

where $y_1 \dots y_t$ are spatial locations (e.g. on a screen), $v_i(y)$ are spatial covariates, $g(\dots)$ is an interaction kernel.

Fitting LATKES using logistic regression

- ▶ Transition kernel only specified up to normalisation constant.
- ▶ Can use sequential version of NCD to turn the problem into (semiparametric) logistic regression.
- ▶ Standard packages can be used (mgcv, INLA).

Example

We fit the model:

$$p(y_t|y_{t-1}) \propto \exp \{ b(\|y_t\|) + r_{\text{dist}} (\|y_t - y_{t-1}\|) + r_{\text{ang}} (\angle (y_t - y_{t-1})) \}$$

where:

- ▶ $b(\|y_t\|)$ should reflect a centrality bias;
- ▶ $r_{\text{dist}} (\|y_t - y_{t-1}\|)$ should reflect the fact that successive fixations are close together;
- ▶ $r_{\text{ang}} (\angle (y_t - y_{t-1}))$ should reflect a tendency for making movements along the cardinal axes (vertical and horizontal).

Note on NCD implementation

- ▶ We fitted functions b , r_{dist} and r_{ang} (plus the log-normalising constant χ , as already explained) using smoothing splines. (Extension of NCD to smoothing splines is direct: simply add appropriate penalty to log-likelihood).
- ▶ We used R package *mgcv* (Wood, 2006).
- ▶ Reference points were sampled from an Uniform distribution (over the screen); 20 reference datapoints per datapoint.
- ▶ Requires one line of code of R, took about 5 minutes.

Results

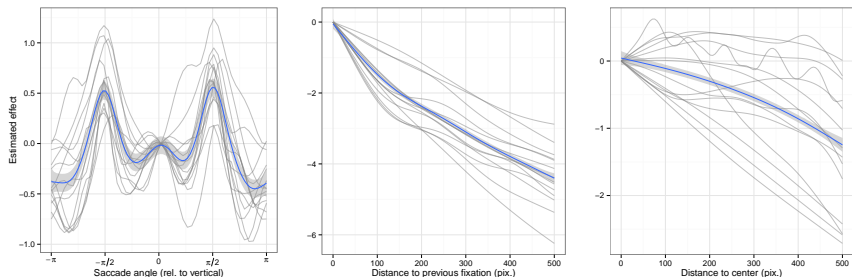


Figure : The different panels display the estimated effects of saccade angle (r_{ang}), distance to previous fixation (r_{dist}) and centrality bias (s). Individual subjects are in gray, and the group average is in blue.

Conclusion

- ▶ Poisson transform: you can treat any data as coming from a Poisson point process in the appropriate space, and infer the *intensity* rather than the density.
 - ▶ *It is OK to treat the normalisation constant as a free parameter!*
- ▶ NCD effectively approximates the Poisson transform via *logistic regression*
- ▶ Inference for unnormalised sequential models can be turned into semi-parametric logistic regression
 - ▶ True as well for unnormalised models with covariates
- ▶ See paper on arxiv (1406.2839) and soon in Statistics and Computing.

References

- Baddeley, A., Berman, M., Fisher, N. I., Hardegen, A., Milne, R. K., Schuhmacher, D., Shah, R., and Turner, R. (2010). Spatial logistic regression and change-of-support in poisson point processes. *Electronic Journal of Statistics*, 4(0):1151–1201.
- Baker, S. G. (1994). The Multinomial-Poisson transformation. *Journal of the Royal Statistical Society. Series D (The Statistician)*, 43(4):495–504.
- Bengio, Y. and Delalleau, O. (2009). Justifying and generalizing contrastive divergence. *Neural computation*, 21(6):1601–1621.
- Geyer, C. J. (1994). Estimating normalizing constants and reweighting mixtures in markov chain monte carlo. Technical Report 568, School of Statistics, University of Minnesota.
- Girolami, M., Lyne, A.-M., Strathmann, H., Simpson, D., and Atchade, Y. (2013). Playing russian roulette with intractable likelihoods. *arxiv 1306.4032*.
- Gutmann, M. U. and Hyvärinen, A. (2012). Noise-contrastive estimation of unnormalized statistical models, with applications to natural image statistics. *J. Mach. Learn. Res.*, 13(1):307–361.
- Hastie, T., Tibshirani, R., and Friedman, J. H. (2003). *The Elements of Statistical Learning*. Springer, corrected edition.
- Kienzle, W., Franz, M. O., Schölkopf, B., and Wichmann, F. A. (2009). Center-surround patterns emerge as optimal predictors for human saccade targets. *Journal of vision*, 9(5).
- Wood, S. (2006). *Generalized Additive Models: An Introduction with R*