Fitting behavioural models

Quentin Huys Dr. med. MBBS PhD

Max Planck UCL Centre for Computational Psychiatry and Ageing and Camden and Islington NHS Foundation Trust

Ringberg 26/9/18









Think of it as four separate two-armed bandit tasks

Analysing behaviour

- Standard approach:
 - Decide which feature of the data you care about
 - Run descriptive statistical tests, e.g. ANOVA

- Many strengths
- Weakness
 - Piecemeal, not holistic / global
 - Descriptive, not generative
 - No internal variables

Analysing behaviour

- Standard approach:
 - Decide which feature of the data you care about
 - Run descriptive statistical tests, e.g. ANOVA



- Many strengths
- Weakness
 - Piecemeal, not holistic / global
 - Descriptive, not generative
 - No internal variables

Models

Holistic

• Aim to model the process by which the data came about in its "entirety"

Generative

• They can be run on the task to generate data as if a subject had done the task

Inference process

- Capture the inference process subjects have to make to perform the task.
- Do this in sufficient detail to replicate the data.

Parameters

- replace test statistics
- their meaning is explicit in the model
- their contribution to the data is assessed in a holistic manner

How to fit a model and believe the

Model building The first step is to build a series of models. Each contains an internal process by which different choice options are valued, and a link function which describes how preferences turn into observed decisions. At least two models should be built: a model M0 of 'no interest' that performs the task, but without involving the process of interest, and a model M1 that does contain the process of interest.

Validation on surrogate data

- 1. **Data generation**: Run each model on the experiment from which data will be examined. Do the generated data look reasonable?
- 2. **Surrogate model fitting**: Fit each model to the data generated from it. Are the true parameters readily recovered? Are some parameters not identifiable?
- 3. **Surrogate model comparison**: Does the model comparison procedure correctly identify the data generated by each model?

Real data analysis

- 1. **Real model fitting**: Fit each model to the real data.
- 2. **Real model validation**: Run each model with the fitted parameters on the exact experimental instance presented to that particular subject. Are the key features of the real data captured reasonably?
- 3. **Real model comparison**: choose the least complex model that best accounts for the data.
- 4. **Parameter examination**: only at this point should the parameters of the model be examined, and only the parameters of the most parsimonious model should be ascribed meaning.

Huys 2017

Actions

- Q values "the process" $Q_t(a_t, s_t) = Q_{t-1}(a_t, s_t) + \epsilon(r_t - Q_{t-1}(a_t, s_t))$
- Probabilities "link function"

$$p(a_t|s_t, h_t, \beta) = p(a_t|\mathcal{Q}(a_t, s_t), \beta)$$
$$= \frac{e^{\beta \mathcal{Q}(a_t, s_t)}}{\sum_{a'} e^{\beta \mathcal{Q}(a', s_t)}}$$

• Features:

$$p(a_t|s_t) \propto \mathcal{Q}(a_t, s_t)$$
$$0 \le p(a) \le 1$$

- Inks learning process and observations
 - choices, RTs, or any other data
 - link function in GLMs
 - many other forms

Link functions



- irreducible noise $p(a|s) = \frac{1-g}{2} + g \frac{e^{\beta Q(a,s)}}{\sum_{a'} e^{\beta Q(a',s)}}$
- critical sanity check 1: reasonable link function?





other link functions for other observations

Fitting models I

Maximum likelihood (ML) parameters

$$\hat{\theta} = \operatorname*{argmax}_{\theta} \mathcal{L}(\theta)$$

where the likelihood of all choices is:

$$\mathcal{L}(\theta) = \log p(\{a_t\}_{t=1}^T | \{s_t\}_{t=1}^T, \{r_t\}_{t=1}^T, \underbrace{\theta}_{\beta, \epsilon})$$

$$= \log p(\{a_t\}_{t=1}^T | \{\mathcal{Q}(s_t, a_t; \epsilon)\}_{t=1}^T, \beta)$$

$$= \log \prod_{t=1}^T p(a_t | \mathcal{Q}(s_t, a_t; \epsilon), \beta)$$

$$= \sum_{t=1}^T \log p(a_t | \mathcal{Q}(s_t, a_t; \epsilon), \beta)$$

Fitting models II

- No closed form
- Use your favourite method
 - gradients
 - fminunc / fmincon...
- Gradients for RW model

$$\begin{aligned} \frac{d\mathcal{L}(\theta)}{d\theta} &= \frac{d}{d\theta} \sum_{t} \log p(a_t | \mathcal{Q}_t(a_t, s_t; \epsilon), \beta) \\ &= \sum_{t} \frac{d}{d\theta} \beta \mathcal{Q}_t(a_t, s_t; \epsilon) - \sum_{a'} p(a' | \mathcal{Q}_t(a', s_t; \epsilon), \beta) \frac{d}{d\theta} \beta \mathcal{Q}_t(a', s_t; \epsilon) \\ \frac{d\mathcal{Q}_t(a_t, s_t; \epsilon)}{d\epsilon} &= (1 - \epsilon) \frac{d\mathcal{Q}_{t-1}(a_t, s_t; \epsilon)}{d\epsilon} + (r_t - \mathcal{Q}_{t-1}(a_t, s_t; \epsilon)) \end{aligned}$$

Little tricks

Transform your variables

$$\beta = e^{\beta'}$$

$$\Rightarrow \beta' = \log(\beta)$$

$$\epsilon = \frac{1}{1 + e^{-\epsilon'}}$$

$$\Rightarrow \epsilon' = \log\left(\frac{\epsilon}{1 - \epsilon}\right)$$

$$\frac{d\log \mathcal{L}(\theta')}{d\theta'}$$

Avoid over/underflow

$$y(a) = \beta Q(a)$$

$$y_m = \max_a y(a)$$

$$p = \frac{e^{y(a)}}{\sum_b e^{y(b)}} = \frac{e^{y(a) - y_m}}{\sum_b e^{y(b) - y_m}}$$

ML characteristics



ML characteristics

- ML is asymptotically consistent, but variance high
 - 10-armed bandit, infer beta and epsilon



ML characteristics

- ML is asymptotically consistent, but variance high
 - 10-armed bandit, infer beta and epsilon



Priors



Priors



Priors



Model fitting

Quentin Huys, UCL

Maximum a posteriori estimate

$$\mathcal{P}(\theta) = p(\theta|a_{1...T}) = \frac{p(a_{1...T}|\theta)p(\theta)}{\int d\theta p(\theta|a_{1...T})p(\theta)}$$
$$\log \mathcal{P}(\theta) = \sum_{t=1}^{T} \log p(a_t|\theta) + \log p(\theta) + const.$$
$$\frac{\log \mathcal{P}(\theta)}{d\alpha} = \frac{\log \mathcal{L}(\theta)}{d\alpha} + \frac{d p(\theta)}{d\theta}$$

If likelihood is strong, prior will have little effect

mainly has influence on poorly constrained parameters

ŨЙ

• if a parameter is strongly constrained to be outside the typical range of the prior, then it will win over the prior

Maximum a posteriori estimate



200 trials, I stimulus, I0 actions, learning rate = .05, beta=2 m_{beta}=0, m_{eps}=-3, n=1

Maximum a posteriori estimate



200 trials, I stimulus, I0 actions, learning rate = .05, beta=2 $m_{beta}=0$, $m_{eps}=-3$, n=1

What prior parameters should I use?







- Fixed effect
 - conflates within- and between- subject variability



- Fixed effect
 - conflates within- and between- subject variability
- Average behaviour
 - disregards between-subject variability
 - need to adapt model



- Fixed effect
 - conflates within- and between- subject variability
- Average behaviour
 - disregards between-subject variability
 - need to adapt model
- Summary statistic
 - treat parameters as random variable, one for each subject
 - overestimates group variance as ML estimates noisy







- Fixed effect
 - conflates within- and between- subject variability
- Average behaviour
 - disregards between-subject variability
 - need to adapt model
- Summary statistic
 - treat parameters as random variable, one for each subject
 - overestimates group variance as ML estimates noisy
- Random effects
 - prior mean = group mean







- Fixed effect
 - conflates within- and between- subject variability
- Average behaviour
 - disregards between-subject variability
 - need to adapt model
- Summary statistic
 - treat parameters as random variable, one for each subject
 - overestimates group variance as ML estimates noisy
- Random effects
 - prior mean = group mean

$$p(\mathcal{A}_i|\mu_{\theta},\sigma_{\theta}) = \int d\theta_i \, p(\mathcal{A}_i|\theta_i) \, p(\theta_i|\mu_{\theta},\sigma_{\theta})$$







- Fixed effect
 - conflates within- and between- subject variability
- Average behaviour
 - disregards between-subject variability
 - need to adapt model
- Summary statistic
 - treat parameters as random variable, one for each subject
 - overestimates group variance as ML estimates noisy
- Random effects
 - prior mean = group mean

$$p(\mathcal{A}_i | \mu_{\theta}, \sigma_{\theta}) = \int d\theta_i \, p(\mathcal{A}_i | \theta_i) \, p(\theta_i | \underbrace{\mu_{\theta}, \sigma_{\theta}})$$









- Fixed effect
 - conflates within- and between- subject variability
- Average behaviour
 - disregards between-subject variability
 - need to adapt model
- Summary statistic
 - treat parameters as random variable, one for each subject
 - overestimates group variance as ML estimates noisy
- Random effects
 - prior mean = group mean

$$p(\mathcal{A}_i | \mu_{\theta}, \sigma_{\theta}) = \int d\theta_i \, p(\mathcal{A}_i | \theta_i) \, p(\theta_i | \underbrace{\mu_{\theta}, \sigma_{\theta}})$$

Random effects

Random effects

- See subjects as drawn from group
- Fixed models



Random effects

- See subjects as drawn from group
- Fixed models
 - all the same: fixed effect wrt model




- See subjects as drawn from group
- Fixed models
 - all the same: fixed effect wrt model
 - parametrically nested





- See subjects as drawn from group
- Fixed models
 - all the same: fixed effect wrt model
 - parametrically nested

$$\mathcal{Q}(a,s) = \omega_1 \mathcal{Q}^1(a,s) + \omega_2 \mathcal{Q}^2(a,s)$$

 assumes within-subject mixture, rather than a group mixture of perfect types





- See subjects as drawn from group
- Fixed models
 - all the same: fixed effect wrt model
 - parametrically nested

- assumes within-subject mixture, rather than a group mixture of perfect types
- w/in subject model comparison?





- See subjects as drawn from group
- Fixed models
 - all the same: fixed effect wrt model
 - parametrically nested

$$\mathcal{Q}(a,s) = \omega_1 \mathcal{Q}^1(a,s) + \omega_2 \mathcal{Q}^2(a,s)$$

- assumes within-subject mixture, rather than a group mixture of perfect types
- w/in subject model comparison?
- HMM





- See subjects as drawn from group
- Fixed models
 - all the same: fixed effect wrt model
 - parametrically nested

- assumes within-subject mixture, rather than a group mixture of perfect types
- w/in subject model comparison?
- HMM
 - switch between models over trials





- See subjects as drawn from group
- Fixed models
 - all the same: fixed effect wrt model
 - parametrically nested

- assumes within-subject mixture, rather than a group mixture of perfect types
- w/in subject model comparison?
- HMM
 - switch between models over trials
- Random effects in models





- See subjects as drawn from group
- Fixed models
 - all the same: fixed effect wrt model
 - parametrically nested

- assumes within-subject mixture, rather than a group mixture of perfect types
- w/in subject model comparison?
- HMM
 - switch between models over trials
- Random effects in models
 - Bayesian model averaging





- See subjects as drawn from group
- Fixed models
 - all the same: fixed effect wrt model
 - parametrically nested

- assumes within-subject mixture, rather than a group mixture of perfect types
- w/in subject model comparison?
- HMM
 - switch between models over trials
- Random effects in models
 - Bayesian model averaging
 - parameter interpretation?







Estimating the hyperparameters

- Effectively we now want to do gradient ascent on: $\frac{d}{d\zeta}p(\mathcal{A}|\zeta)$
- But this contains an integral over individual parameters:

$$p(\mathcal{A}|\zeta) = \int d\theta p(\mathcal{A}|\theta) \, p(\theta|\zeta)$$

• So we need to: $\hat{\zeta} = \operatorname*{argmax}_{\zeta} p(\mathcal{A}|\zeta)$ = $\operatorname*{argmax}_{\zeta} \int d\theta p(\mathcal{A}|\theta) p(\theta|\zeta)$

Inference

$$\hat{\zeta} = \operatorname{argmax}_{\zeta} p(\mathcal{A}|\zeta)$$
$$= \operatorname{argmax}_{\zeta} \int d\theta p(\mathcal{A}|\theta) p(\theta|\zeta)$$

- analytical rare
- brute force for simple problems
- Expectation Maximisation approximate, easy
- Variational Bayes approximate, often hard
- Sampling / MCMC slow, easy

Expectation Maximisation

$$\begin{split} \log p(\mathcal{A}|\zeta) &= \log \int d\theta \, p(\mathcal{A}, \theta|\zeta) \\ &= \log \int d\theta \, q(\theta) \frac{p(\mathcal{A}, \theta|\zeta)}{q(\theta)} \\ &\geq \int d\theta \, q(\theta) \log \frac{p(\mathcal{A}, \theta|\zeta)}{q(\theta)} \end{split} \text{Jensen's inequality} \\ k^{\text{th}} \; \text{E step:} \; q^{(k+1)}(\theta) \; \leftarrow \; p(\theta|\mathcal{A}, \zeta^{(k)}) \\ k^{\text{th}} \; \text{M step:} \; \zeta^{(k+1)} \; \leftarrow \; \underset{\zeta}{\operatorname{argmax}} \int d\theta \, q(\theta) \log p(\mathcal{A}, \theta|\zeta) \end{split}$$

- Iterate between
 - Estimating MAP parameters given prior parameters
 - Estimating prior parameters from MAP parameters









EM with Laplace approximation

• E step:
$$q^{(k+1)}(\theta) \leftarrow p(\theta|\mathcal{A}, \zeta^{(k)})$$

- only need sufficient statistics to perform M step
- Approximate $p(\theta|\mathcal{A}, \zeta^{(k)}) \sim \mathcal{N}(\mathbf{m}_k, \mathbf{S}_k)$
- and hence:

E step:
$$q_k(\theta) = \mathcal{N}(\mathbf{m}_k, \mathbf{S}_k)$$

 $\mathbf{m}_k \leftarrow \operatorname*{argmax}_{\theta} p(\mathbf{a}_k | \theta) p(\theta | \zeta^{(i)})$
 $\mathbf{S}_k^{-1} \leftarrow \frac{\partial^2 p(\mathbf{a}^k | \theta) p(\theta | \zeta^{(i)})}{\partial \theta^2} \Big|_{\theta = \mathbf{m}_k}$
matlab: [m,L,,,S]=fminunc(...)

Just what we had before: MAP inference given some prior parameters

EM with Laplace approximation



And now iterate until convergence

Parameter recovery



Correlations



Are parameters ok for correlations?

- Individual subject parameter estimates NO LONGER INDEPENDENT!
 - Change group -> change parameter estimates
- compare different params
 - if different priors
- correlations, t-tests
 - within same prior ok
 - more power than ML



So far

- infer individual parameters
- apply standard tests

- So far
 - infer individual parameters
 - apply standard tests
- Alternative
 - View as variation across group
 - Specific more powerful?



- So far
 - infer individual parameters
 - apply standard tests
- Alternative
 - View as variation across group
 - Specific more powerful?

$$\mu_{\theta}^{i} = \mu_{\theta}^{\text{Group}} + \beta \psi_{i}$$



- So far
 - infer individual parameters
 - apply standard tests
- Alternative
 - View as variation across group
 - Specific more powerful?





Group-level regressor



Group-level regressor



Group error bars

- EM standardly does not provide error bars
- Shifted samples

$$p(\mathcal{A}|\zeta) \approx \sum_{i} p(\mathcal{A}|\theta_{i}); \quad \theta_{i} \sim p(\theta|\zeta)$$
 (1)

$$\frac{\partial}{\partial \mu} p(\mathcal{A}|\zeta) \approx \frac{1}{\delta} \left[\sum_{i} p(\mathcal{A}|\theta_{i} + \delta) - \sum_{i} p(\mathcal{A}|\theta_{i}) \right]$$
(2)

• Oakes 1999

- analytical description of gradients
- tricky, but combined with forward differentiation it is automatic (julialang.org Pkg ForwardDiff)

GLM error bars

Shifted samples RW 0.0 0.0 0.0 0.1 0.2 0.4 0.6 0.8 1 p value

GLM error bars



0 0.2 0.4 0.6 0.8 1 p value

GLM error bars

Shifted samples

Forward differentiation

RW



Hierarchical / random effects models

Advantages

- Accurate group-level mean and variance
- Outliers due to weak likelihood are regularised
- Strong outliers are not
- Useful for model selection

Disadvantages

- Individual estimates θ_i depend on other data, i.e. $\partial A_{j \neq i}$ and therefore need to be careful in interpreting these as summary statistics
- More involved; less transparent
- Psychiatry
 - Groups often not well defined, covariates better
- ► fMRI
 - Shrink variance of ML estimates fixed effects better still?

How does it do?



Overfitting



Model comparison

- A fit by itself is not meaningful
- Generative test
 - qualitative
- Comparisons
 - vs random
 - vs other model -> test specific hypotheses and isolate particular effects in a generative setting

Model comparison

Averaged over its parameter settings, how well does the model fit the data?

$$p(\mathcal{A}|\mathcal{M}) = \int d\theta \, p(\mathcal{A}|\theta) \, p(\theta|\mathcal{M})$$

Model comparison: Bayes factors

$$BF = \frac{p(\mathcal{A}|\mathcal{M}_1)}{p(\mathcal{A}|\mathcal{M}_2)}$$

- Problem:
 - integral rarely solvable
 - approximation: Laplace, sampling, variational...
Why integrals? The God Almighty test





Why integrals? The God Almighty test





Why integrals? The God Almighty test



 $\frac{1}{N} \left(\mathbf{p}(\mathbf{X}|\boldsymbol{\theta}_1) + p(X|\boldsymbol{\theta}_2) + \cdots \right)$

These two factors fight it out Model complexity vs model fit

Group-level BIC

$$\begin{split} \log p(\mathcal{A}|\mathcal{M}) &= \int d\boldsymbol{\zeta} \, p(\mathcal{A}|\boldsymbol{\zeta}) \, p(\boldsymbol{\zeta}|\mathcal{M}) \\ &\approx -\frac{1}{2} \mathsf{BIC}_{\mathsf{int}} \\ &= \log \hat{p}(\mathcal{A}|\hat{\boldsymbol{\zeta}}^{ML}) - \frac{1}{2} |\mathcal{M}| \log(|\mathcal{A}|) \end{split}$$

- Very simple
 - 1) EM to estimate group prior mean & variance
 - simply done using fminunc, which provides Hessians
 - 2) Sample from estimated priors
 - 3) Average

How does it do?







How does it do?



Group Model selection

Integrate out your parameters

Model comparison: overfitting?





Model comparison: overfitting?



Posterior distribution on models

Generative model for models



Bayesian model selection - equations

- Write down joint distribution of generative model
- Variational approximations lead to set of very simple update equations
 - start with flat prior over model probabilities

 $\alpha = \alpha_0$

• then update

$$u_{k}^{i} = \left(\int d\theta_{i} p(\mathcal{A}_{i}, \theta_{i} | \mathcal{M}_{k}) \right) \exp \left(\Psi(\alpha_{k}) - \Psi\left(\sum_{k} \alpha_{k}\right) \right)$$

$$\alpha_{k} \leftarrow \alpha_{0,k} + \sum_{i} \frac{u_{k}^{i}}{\sum_{k} u_{k}^{i}}$$

emfit

- www.quentinhuys.com/pub/emfit
- six model classes
 - basic RW to learn
 - Affective Go/Nogo (Guitart et al., 2012)
 - Probabilistic Reward (Pizzagalli et al., 2005)
 - Twostep (Daw et al., 2011)
 - Effort (Gold et al., 2013)
 - Pruning (Huys et al., 2012)
- ► ML, MAP and EM-MAP
- GLM (no implicit differentiation, so only for small models!)
- many models for each task standardly fit all
- model comparison
- data generation and plotting of standard sanity checks

emfit

Run example

- generate sample dataset
- fit all models
- do model comparison, output as figure & latex file
- plot sanity / generative checks
- batchRunEMfit('mTASKNAME')

Fit dataset

- check data format in mTASKNAME/dataformat.txt or look at example dataset generated above in fitResults/ Data.mat
- batchRunEMfit('mTASKNAME',Data.mat)

How to fit a model and believe the results

Model building The first step is to build a series of models. Each contains an internal process by which different choice options are valued, and a link function which describes how preferences turn into observed decisions. At least two models should be built: a model M0 of 'no interest' that performs the task, but without involving the process of interest, and a model M1 that does contain the process of interest.

Validation on surrogate data

- 1. **Data generation**: Run each model on the experiment from which data will be examined. Do the generated data look reasonable?
- 2. **Surrogate model fitting**: Fit each model to the data generated from it. Are the true parameters readily recovered? Are some parameters not identifiable?
- 3. **Surrogate model comparison**: Does the model comparison procedure correctly identify the data generated by each model?

Real data analysis

- 1. **Real model fitting**: Fit each model to the real data.
- 2. **Real model validation**: Run each model with the fitted parameters on the exact experimental instance presented to that particular subject. Are the key features of the real data captured reasonably?
- 3. **Real model comparison**: choose the least complex model that best accounts for the data.
- 4. **Parameter examination**: only at this point should the parameters of the model be examined, and only the parameters of the most parsimonious model should be ascribed meaning.

Huys 2017

Behavioural data modelling

Are no panacea

- statistics about specific aspects of decision machinery
- only account for part of the variance
- Model needs to match experiment
 - ensure subjects actually do the task the way you wrote it in the model
 - model comparison
- Model = Quantitative hypothesis
 - strong test
 - need to compare **models**, not **parameters**

Thanks

- Nathaniel Daw
- Peter Dayan
- Daniel Schad

