

NEU50 I: Learning & Memory



Lecture 3: action learning
trial by trial model fitting
model comparison

1

Act I: where were we?
(action selection)

2

Temporal Difference (TD) learning



The problem: optimal prediction of future reinforcement

The algorithm:

$$V_t = E[r_{t+1}] + V_{t+1}$$

$$V_t^{T+1} = V_t^T + \eta \left(r_{t+1}^T + V_{t+1}^T - V_t^T \right)$$

(note: t indexes time
within a trial,
 T indexes trials)

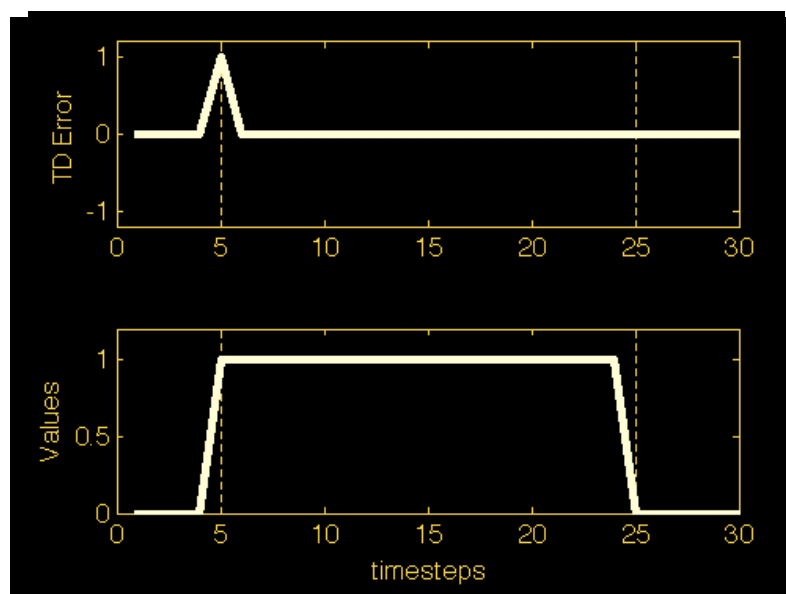
temporal difference prediction error $\delta(t+1)$

compare to: $V^{T+1} = V^T + \eta (r^T - V^T)$

Sutton & Barto 1983, 1990

3

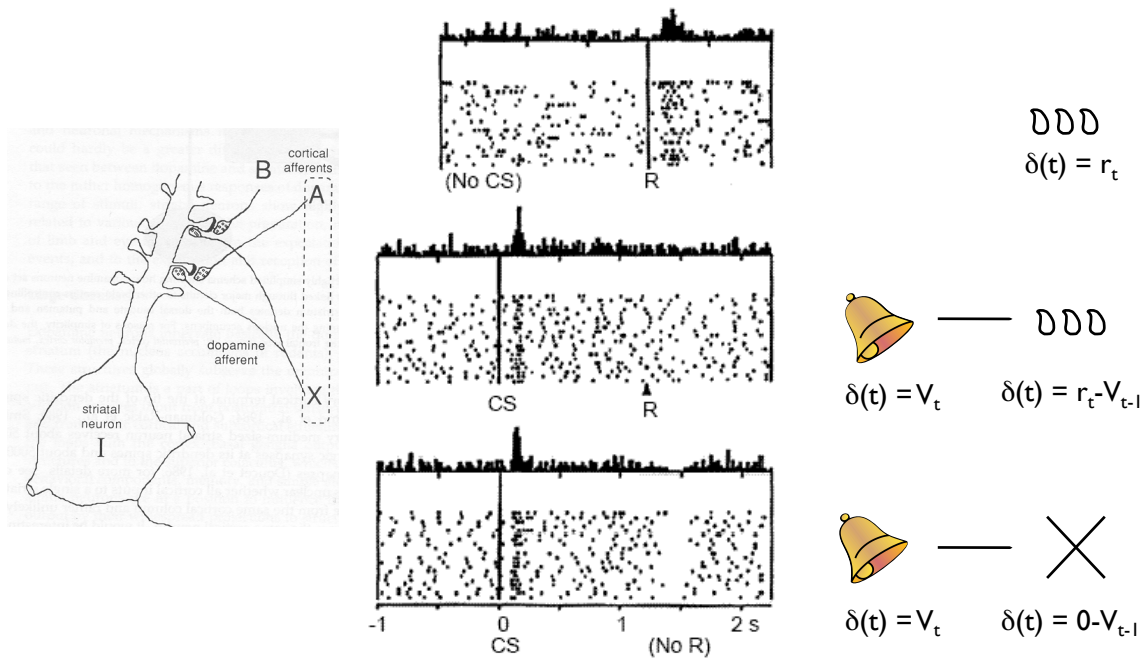
simulation



what would happen with partial reinforcement?
what would happen in second order conditioning?

4

implemented through dopamine

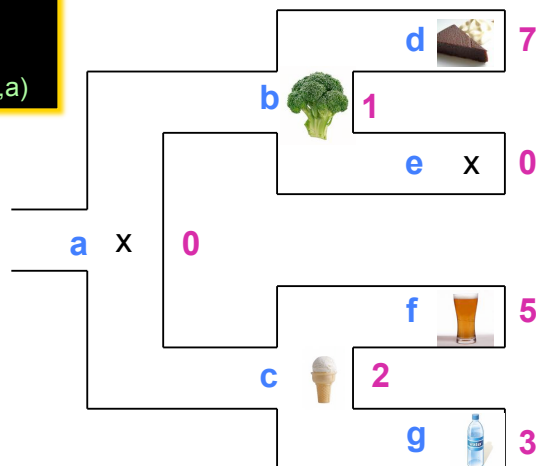


Schultz et al, 1997

5

MDPs

- States S
- Actions $\pi_{s,a} = P(a|S)$
- Transitions $T_{s \rightarrow s'} = P(S'|S, a)$
- Rewards $R_{s \rightarrow s'} = P(R|S, S', a)$



- The idea: given the current situation, history does not matter
- $P(S_{t+1}|S_t, a_t) = P(S_{t+1}|S_1, S_2, \dots, S_t, a_1, a_2, \dots, a_t)$
- $P(r_t|S_t, a_t) = P(r_t|S_1, S_2, \dots, S_t, a_1, a_2, \dots, a_t)$

6

Stylized task: described fully by S,A,R,T

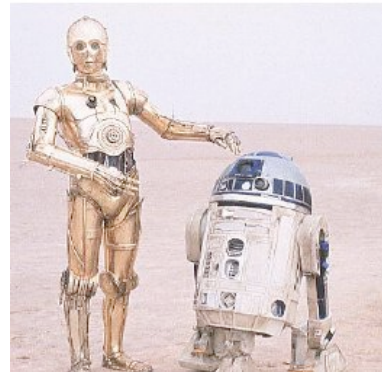
World: "You are in state 34. Your immediate reward is 3. You have 2 actions"

Robot: "I'll take action 1"

World: "You are in state 77. Your immediate reward is -7. You have 3 actions"

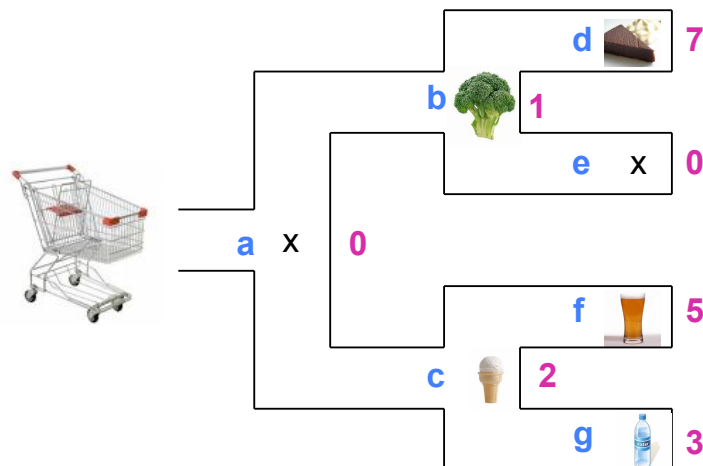
Robot: "I'll take action 3"

The task description requires no memory
(doesn't mean that the decision maker does not use memory to solve the task!)



7

what can we compute here?



state values: $V(S) = E[\text{sum of future rewards}|S]$

actually: $V^\pi(S) = E[\text{sum of future rewards}|\pi, S]$

8

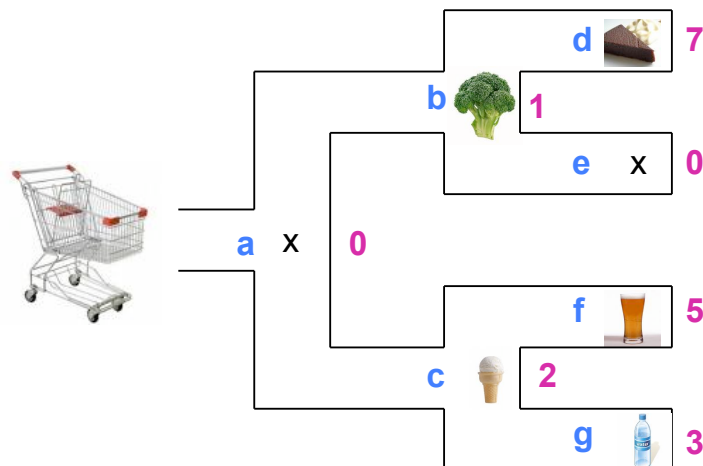
Key RL idea #1: Bellman's glorious equation

$$V^{\pi}(S) = \sum_a \pi_{s,a} \sum_{s'} T^a_{s \rightarrow s'} [R^a_{s \rightarrow s'} + V^{\pi}(S')]$$

In a Markov decision process state values are recursive

9

but there's more: computing the value of actions



(policy dependent) State-Action values:

$$Q^{\pi}(\text{action}|\text{state}) = E[\text{sum of future rewards}|S,a,\pi]$$

- $Q(\text{left}|a) = ?$ $Q(\text{right}|a) = ?$
- which action is better?

10

Key RL idea #1 (again): Bellman's glorious equation

$$Q(S,a) = \sum_{s'} T^a_{s \rightarrow s'} [R^a_{s \rightarrow s'} + V(s')]$$

But.. what if we don't know T, R?

11

model-free learning: sampling

World: "You are in state 34. Your immediate reward is 3. You have 2 actions"

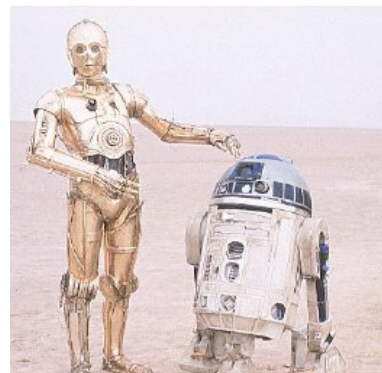
Robot: "I'll take action 1"

World: "You are in state 77. Your immediate reward is -7. You have 3 actions"

Robot: "I'll take action 3"

Take actions according to policy.

Treat experienced rewards and transitions as *samples*



12

Key RL idea #2: Model free learning

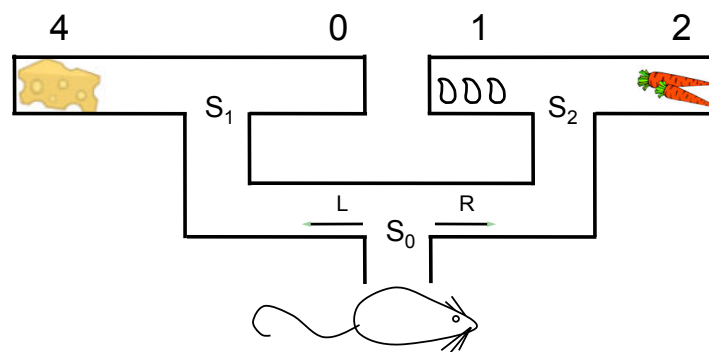
$$V^{\pi}(S) = \sum_a \pi_{S,a} \sum_{S'} T^a_{S \rightarrow S'} [R^a_{S \rightarrow S'} + V^{\pi}(S')]$$

1. choose initial values $V_0(S)$
2. at time point t and state S_t behave according to π
3. observe S_{t+1} and $r(S_{t+1})$
4. compute prediction error $r(S_{t+1}) + V(S_{t+1}) - V(S_t)$
5. update $V(S_t)$ according to prediction error

learning of long-term values can be done
using local information only

13

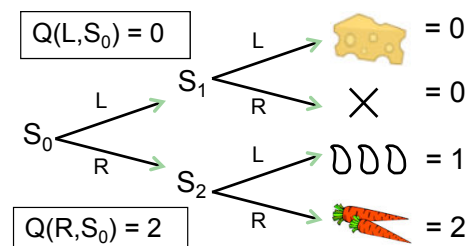
strategy I: “model-based” RL



learn model of task through experience (= cognitive map)

compute Q values by “looking ahead” in the map

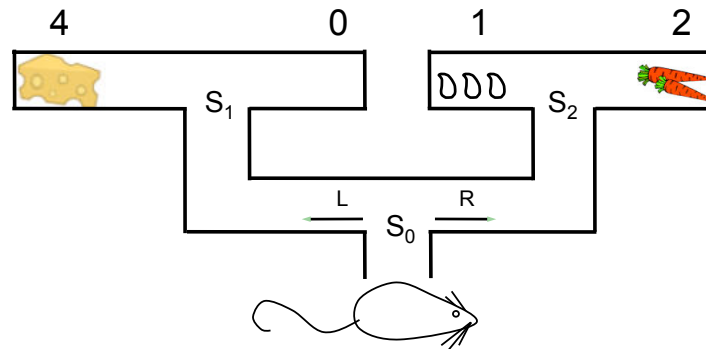
computationally costly, but also flexible
(immediately sensitive to change)



Daw et al (2005)

14

strategy II: “model-free” RL



- Shortcut: store long-term values
 - then simply retrieve them to choose action
- Can learn these from experience
 - without building or searching a model
 - incrementally through prediction errors
 - dopamine dependent SARSA/Q-learning or Actor/Critic

Stored:

$$Q(S_0, L) = 4$$

$$Q(S_0, R) = 2$$

$$Q(S_1, L) = 4$$

$$Q(S_1, R) = 0$$

$$Q(S_2, L) = 1$$

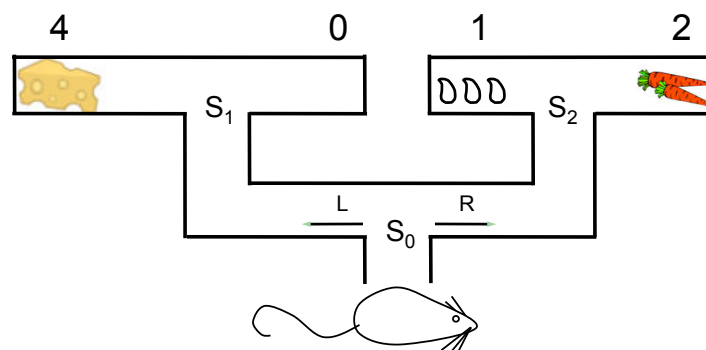
$$Q(S_2, R) = 2$$

Daw et al (2005)

15

15

strategy II: “model-free” RL



- choosing actions is easy so behavior is quick, reflexive (S-R)
- but needs a lot of experience to learn
- and inflexible, need relearning to adapt to any change (habitual)

Stored:

$$Q(S_0, L) = 4$$

$$Q(S_0, R) = 2$$

$$Q(S_1, L) = 4$$

$$Q(S_1, R) = 0$$

$$Q(S_2, L) = 1$$

$$Q(S_2, R) = 2$$

Daw et al (2005)

16

16

summary so far

Instrumental learning: an instance of learning optimal control

MDPs: class of stylized tasks

In a Markov process long term values can be defined that

- are self consistent (recursively defined)
- can be learned incrementally (dynamic programming)
- can be learned from experience even without a world model

These values are helpful because they can help us improve the policy!

17

Act II:
What does my *model* tell me
about my *data*?

18

bandit tasks



overall approach:

- learn values for options (how is this problem simpler than those we've been talking about?)
- choose the best option

suppose we ran this experiment on a person:

- what are the data?
- what do our models predict?
- what can we conclude/infer from the data?

our models are basically detailed hypotheses about behavior and about the brain... we can test these hypotheses!

19

Writing down a full model of learning

what do we know?

what can we measure?

what do we not know?

20

Estimating model parameters

why estimate parameters?

1. May measure quantities of interest (learning rates in different populations, how variance in the task affects learning rate etc.)
2. have to use these to generate hidden variables of interest (eg. prediction errors) in order to look for these in the brain

how to estimate parameters?

we want: $P(\alpha, \beta \mid D, M)$

wwBd?

$$P(\alpha, \beta \mid D, M) \propto P(D \mid \alpha, \beta, M) \quad \leftarrow \text{This we know!}$$

$$P(D \mid \alpha, \beta, M) = \prod P(c_t \mid \alpha, \beta, M)$$

21

Estimating model parameters

$$P(\alpha, \beta \mid D, M) \propto P(D \mid \alpha, \beta, M)$$

...this is a probability distribution

we often consider a point estimate: the maximal likelihood point

$$\operatorname{argmax}_{\alpha, \beta} P(D \mid \alpha, \beta, M)$$

(equivalently, can maximize the log likelihood)

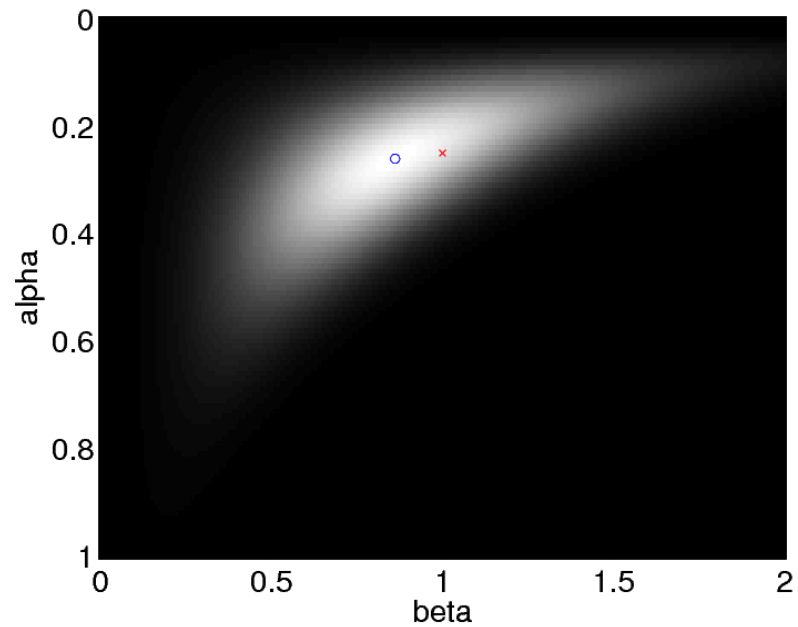
22

Estimating model parameters

$$P(D \mid \alpha, \beta, M) = \prod P(c_t \mid \alpha, \beta, M)$$

pragmatics:
fmincon, fminunc, fminsearch
minimize -LL
local minima: multiple starting points

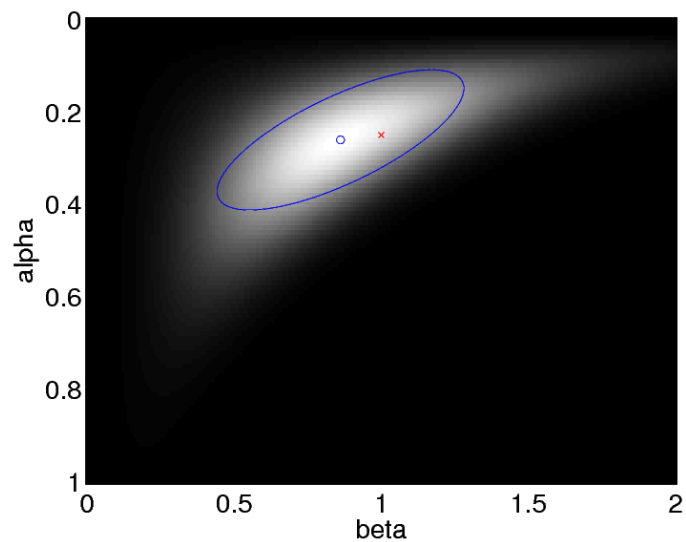
simulated 1000 trials
 $\alpha=0.25, \beta=1$
recovered:
 $\alpha=0.26, \beta=0.86$



23

Error bars on estimates

intuition: how fast likelihood is changing as parameters change

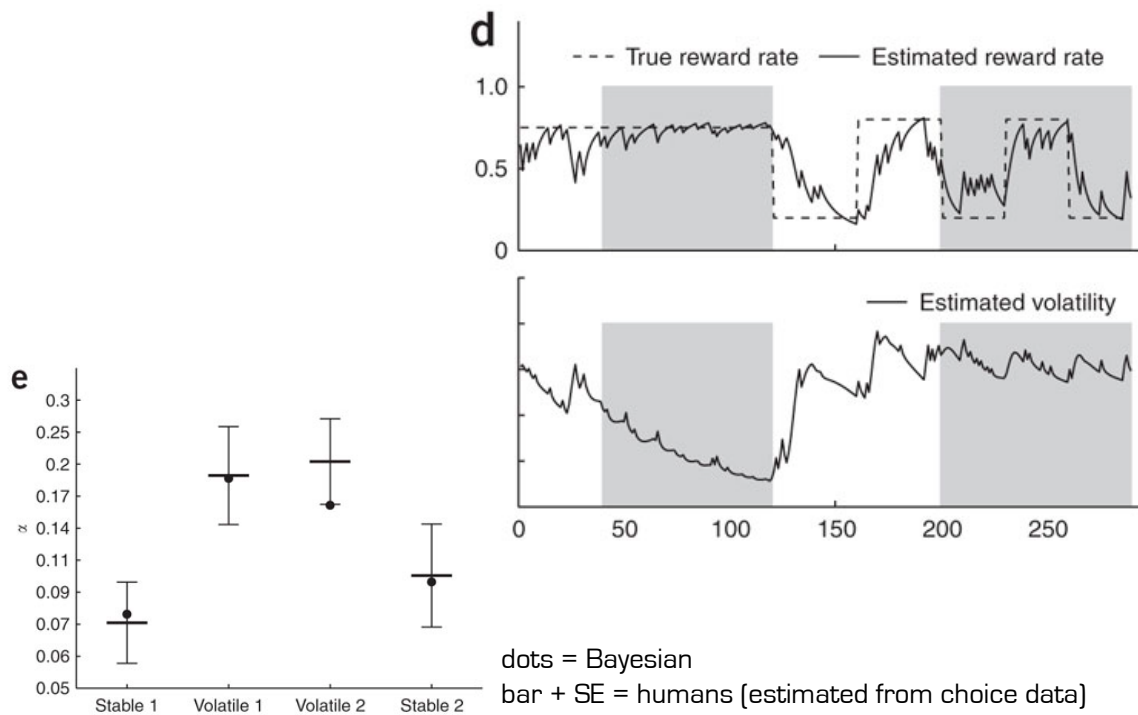


pragmatics: inverse Hessian (2nd derivative matrix; fmincon gives you this) of -LL
estimates parameter covariance matrix
-error bars along the diagonal (sqrt), covariation off the diagonal
-can also look at variation in fits across subjects (we won't go into detail on this)

24

Example: do learning rates adapt to volatility?

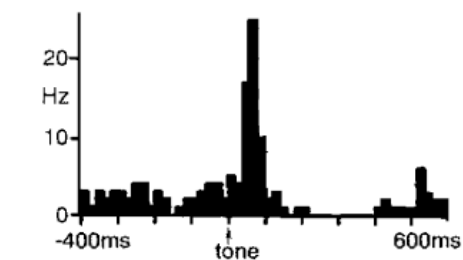
Behrens et al. (2007)



25

Example: novelty bonuses

Wittmann et al (2008)



Horvitz et al. (1997)
dopamine neurons

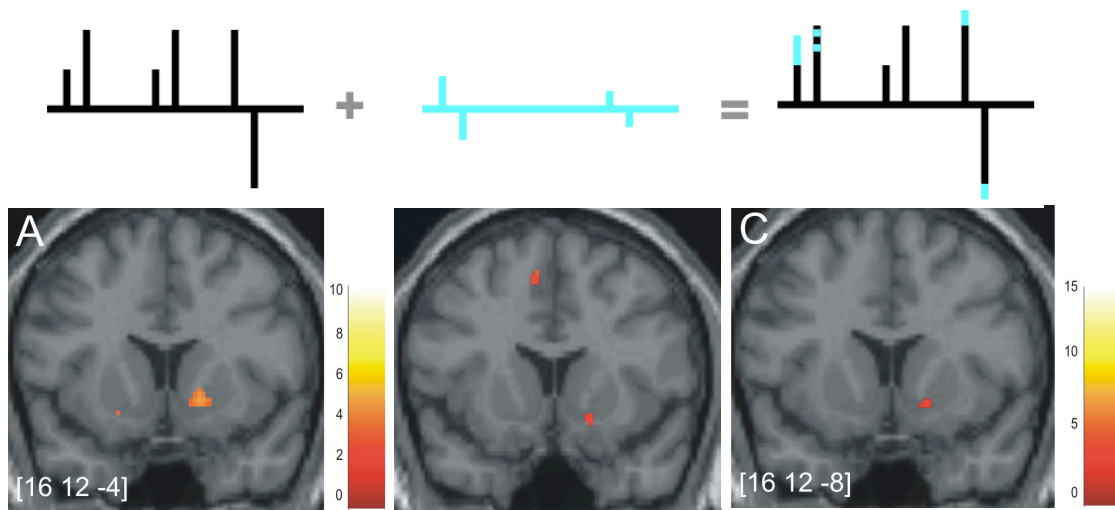


Fit initial value separately for novel
and preexposed images
initial value (preexposed) = £0.37
initial value (novel) = £0.41 (!!)

26

Example: novelty bonuses

Wittmann et al (2008)



27

Summary so far

- Learning models are detailed hypotheses about trial-by-trial overt and covert variables
- these can be tested against the brain
- help understand what different brain regions/networks are doing/computing
- a whole host of interesting results so far, but many questions still unanswered (relatively new method!)
- the models help us learn about the brain... can we also use the brain (or behavior) to learn about the models??

28

Act III:

What does my *data* tell me about my *model*?

29

Which model is best? Model comparison

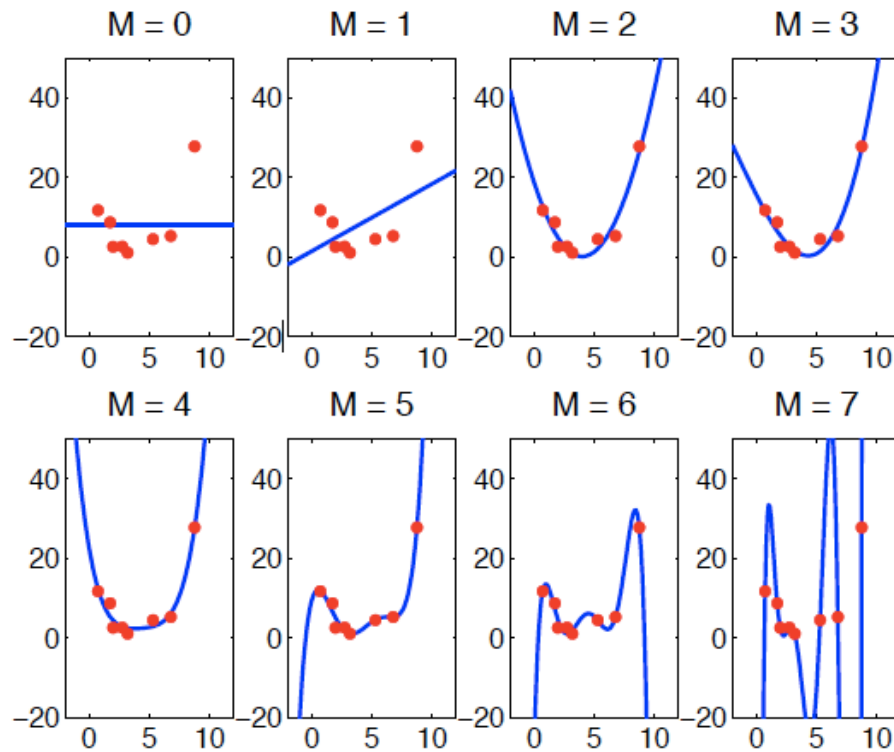
- $P(\text{Model} \mid \text{Data}) = ?$

- comparing two models:
$$\frac{P(M_1 \mid D)}{P(M_2 \mid D)} = \frac{P(D \mid M_1) \cdot P(M_1)}{P(D \mid M_2) \cdot P(M_2)}$$

Bayes factor

30

Which model is best? Model comparison



31

Which model is best? Model comparison

- $P(\text{Model} \mid \text{Data}) = ?$

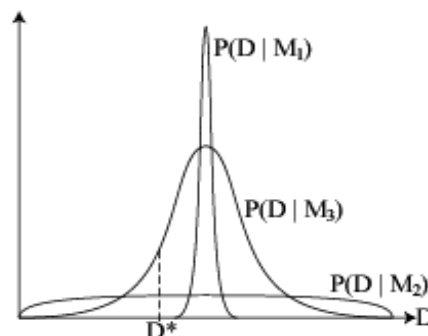
- comparing two models:

$$\frac{P(M_1 \mid D)}{P(M_2 \mid D)} = \frac{P(D \mid M_1) \cdot P(M_1)}{P(D \mid M_2) \cdot P(M_2)}$$

automatic Occam's razor:
simple models tend to make
precise predictions

can put in preference for
simple models here, but
don't need to...

- "Pluralitas non est ponenda sine necessitate"
Plurality should not be posited without necessity – William of Ockham (1349)
- we should go for the simplest model that explains the data



32

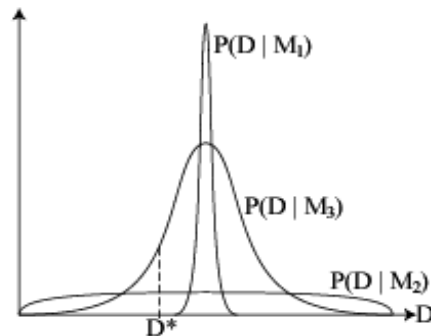
Which model is best? Model comparison

$$\frac{P(M_1|D)}{P(M_2|D)} = \frac{P(D|M_1) \cdot P(M_1)}{P(D|M_2) \cdot P(M_2)}$$

assuming uniform prior over models all we care about is $P(D | M)$

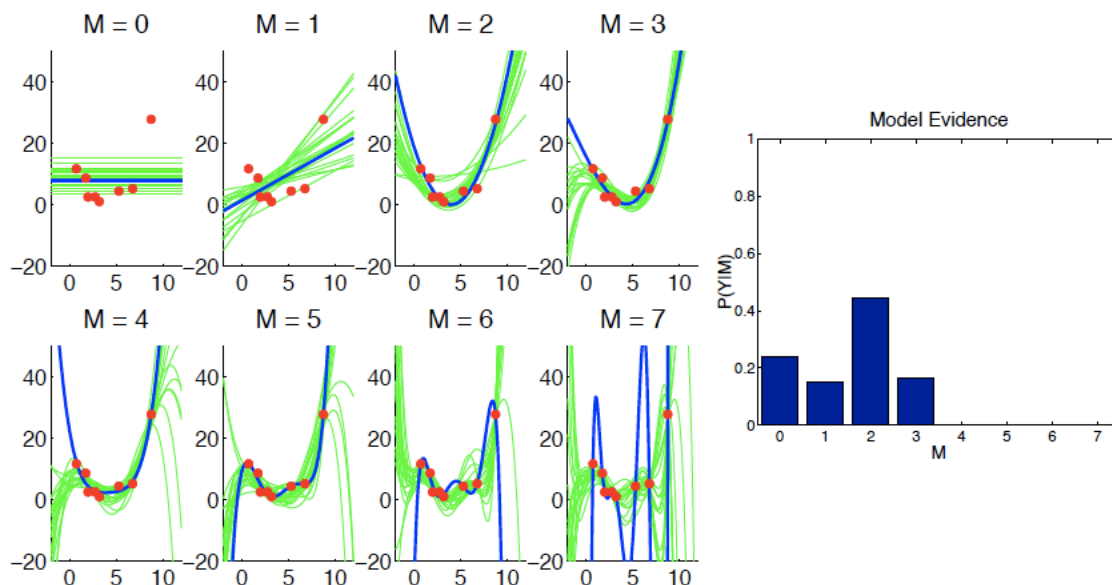
$$P(D|M) = \int d\theta P(D|M, \theta) \cdot P(\theta)$$

Bayesian evidence for model M (marginal likelihood)



33

Bayesian model comparison: Occam's Razor at work



34

Computing $P(D | M)$

$$P(D|M) = \int d\theta P(D|M, \theta) \cdot P(\theta)$$

- Integrating over all settings of the parameters is too hard...
- Approximate solutions:
 - sample posterior at many places to approximate integral and compute Bayes factor directly
 - Laplace approximation: make Gaussian approximation around MAP parameter estimate

35

Laplace approximation

$$P(D|M) = \int d\theta P(D|M, \theta) \cdot P(\theta)$$

- For large amounts of data (compared to # of parms d) the posterior is approximately Gaussian around the maximum a posteriori (MAP) estimate $\hat{\theta}$

$$P(\theta|D, M) \approx 2\pi^{-\frac{d}{2}} |A|^{\frac{1}{2}} \exp \left\{ -\frac{1}{2} (\theta - \hat{\theta})^T A (\theta - \hat{\theta}) \right\}$$

- and we also know that

$$P(D|M) = \frac{P(\theta, D|M)}{P(\theta|D, M)} = \frac{P(\theta|M)P(D|\theta, M)}{P(\theta|D, M)}$$

- so we can compute around the MAP estimate:

$$\ln P(D|M) \approx \ln P(\hat{\theta}|M) + \ln P(D|\hat{\theta}, M) + \frac{d}{2} \ln(2\pi) - \frac{1}{2} \ln |A|$$

- where $-A$ is the Hessian matrix of $\ln P(\theta | D, M)$

$$A_{kl} = -\frac{\partial^2}{\partial \theta_{mk} \partial \theta_{ml}} \ln P(\theta|D, M)|_{\hat{\theta}}$$

36

BIC approximation

$$\ln P(D|M) \approx \underbrace{\ln P(\hat{\theta}|M)}_{\text{prior on } \theta} + \underbrace{\ln P(D|\hat{\theta}, M)}_{\text{data log likelihood}} + \underbrace{\frac{d}{2} \ln(2\pi)}_{\text{easy}} - \underbrace{\frac{1}{2} \ln |A|}_{\text{easy}}$$

- In the limit of LOTS of data ($N \rightarrow \infty$) A grows as NA_0 (for fixed A_0) so $\ln |A| = \ln |NA_0| = \ln N^d |A_0| = d \ln N + \ln |A_0|$.
- Retaining only terms that grow with N , we can approximate further:

$$\ln P(D|M) \approx \ln P(D|\hat{\theta}, M) - \frac{d}{2} \ln(N)$$

(so, for each model we compute the log likelihood for the ML parameters and then add to that a penalty that depends on d (# of parameters), and then we compare the results between the models)

- Advantages: easy to compute; can use ML rather than MAP estimate
- Disadvantage: hard to determine d (only identifiable parameters) and N (only samples used to fit parameters; what if not same for diff parms?)

37

Non-Bayesian alternatives

- Likelihood ratio test: for nested models (one is a special case of the other; compares hypothesis H_1 to one where some parameters are fixed, H_0).
 - Statistical test on the likelihood differences: compare 2* difference in log likelihood (ML) to χ^2 statistic with $df = \#$ additional parameters
- AIC (Akaike's information criterion, 1974): measures goodness of a model based on bias and variance of the estimated model and measures of entropy.
 - Not statistical test, only ranks models.
 - Penalize log likelihood (ML) by adding # of parameters
- Fit models on training set and validate fit on hold-out set.
 - Problem: often hard to find two i.i.d. sets in a learning setting

38

Summary so far...

- Learning models are detailed hypotheses about trial-by-trial overt and covert variables
- trial-by-trial model fitting lets us test these hypotheses
- ...and compare alternatives
- special premium on detailed model fitting when considering learning data: non-stationary, can't use traditional averaging techniques
- a lot of leverage to pinpoint the neural correlates of learning and decision making in the brain

39

Additional reading

- Daw (2011) - Trial by trial data analysis using computational models
- Hare et al. (2008) - Dissociating the role of the orbitofrontal cortex and the striatum in the computation of goal values and prediction errors
- Niv (2009) - Reinforcement learning in the brain

40