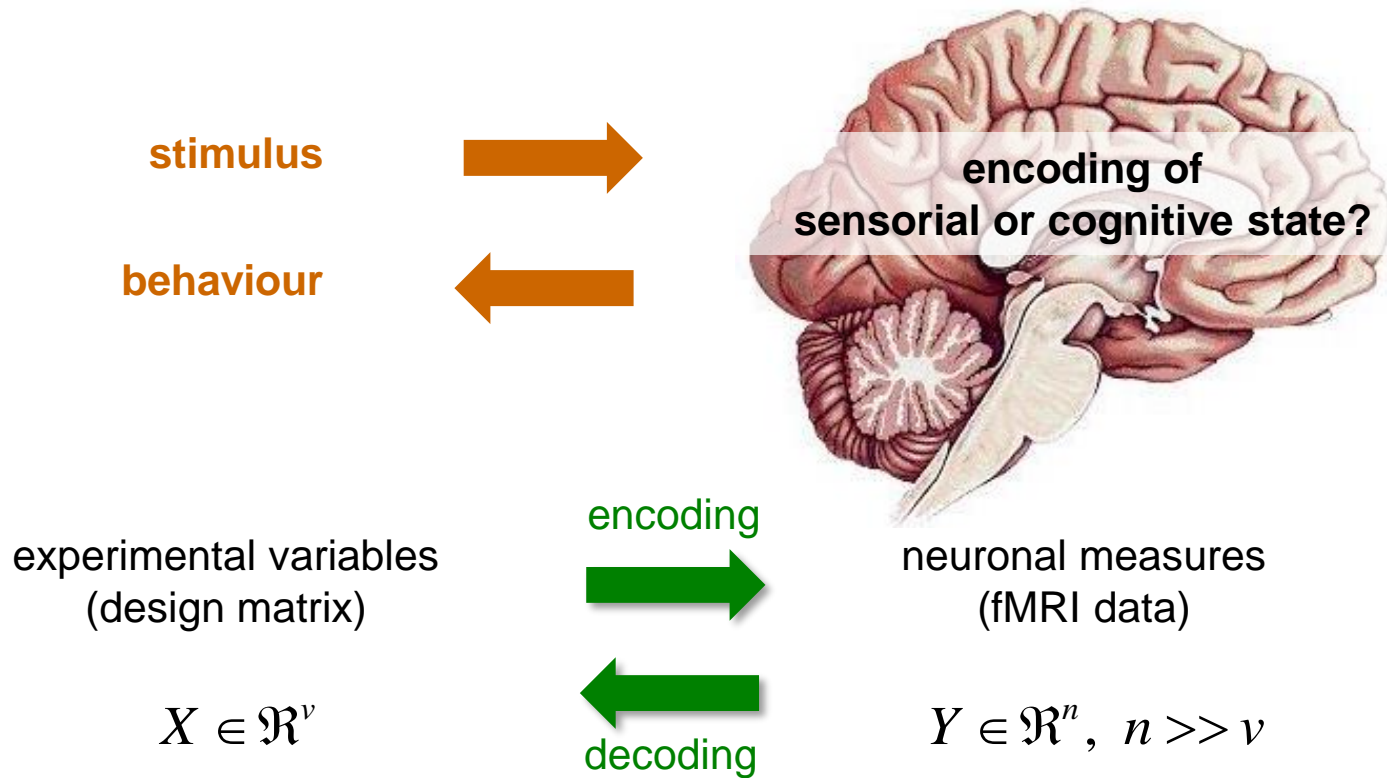


# **MultiVariate Bayesian (MVB) decoding of brain images**

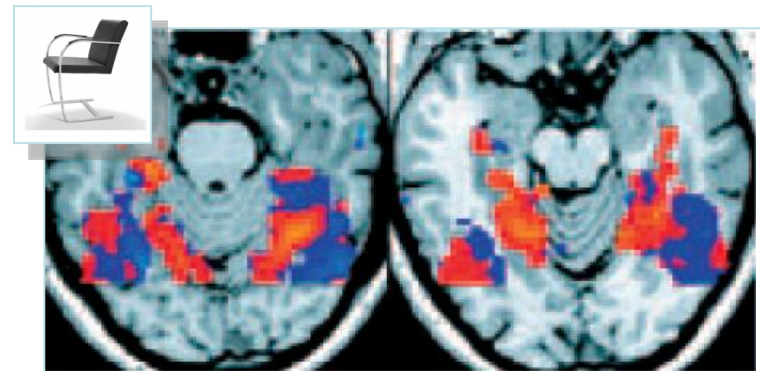
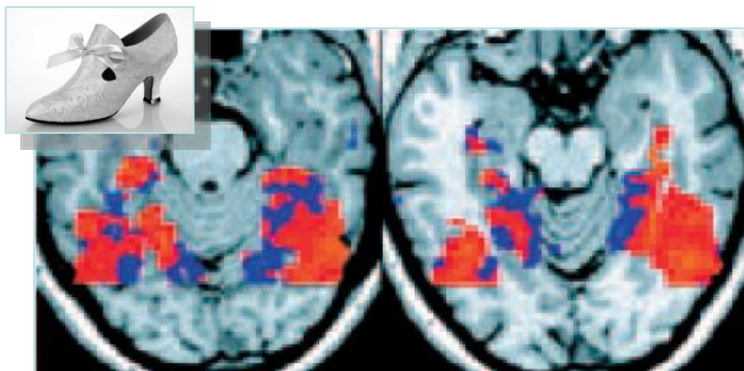
**Alexa Morcom**

**Edinburgh SPM course 2015**

**With thanks to J. Daunizeau, K. Brodersen for slides**



What if neuronal responses are distributed (over space)?



# Overview of the talk

## 1 Introduction

*1.1 Lexicon*

*1.2 “Decoding”: so what?*

*1.3 Multivariate: so what?*

*1.4 Preliminary statistical considerations*

## 2 Multivariate Bayesian decoding

*2.1 From classical encoding to Bayesian decoding*

*2.2 Hierarchical priors on patterns*

*2.3 Probabilistic inference*

## 3 Example

## 4 Summary

# Overview of the talk

## 1 Introduction

*1.1 Lexicon*

*1.2 “Decoding”: so what?*

*1.3 Multivariate: so what?*

*1.4 Preliminary statistical considerations*

## 2 Multivariate Bayesian decoding

*2.1 From classical encoding to Bayesian decoding*

*2.2 Hierarchical priors on patterns*

*2.3 Probabilistic inference*

## 3 Example

## 4 Summary

# Lexicon

the jargon to swallow

## 1 Encoding or decoding?

- An **encoding** model (or generative model) relates context (independent variable) to brain activity (dependent variable).
- A **decoding** model (or recognition model) relates brain activity (independent variable) to context (dependent variable).

$$X \rightarrow Y$$

$$Y \rightarrow X$$

## 2 Univariate or multivariate?

- In a **univariate** model, brain activity is the signal measured in one voxel.
- In a **multivariate** model, brain activity is the signal measured in many voxels (NB: *decoding*  $\rightarrow$  *ill-posed problem*).

$$Y \in \mathfrak{R}$$

$$Y \in \mathfrak{R}^n, \quad n \gg v$$

## 3 Regression or classification?

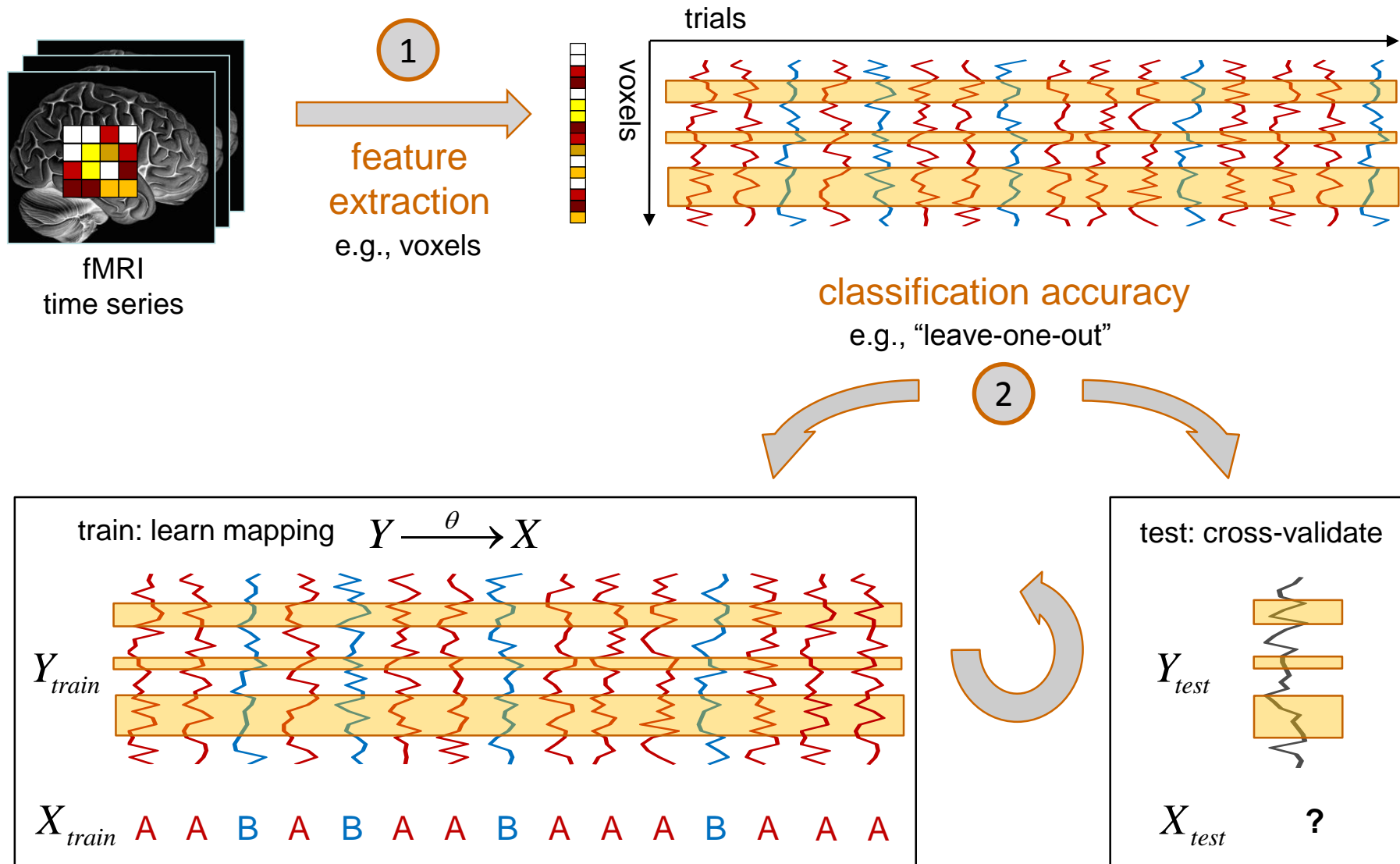
- In a **regression** model, the dependent variable is continuous.
- In a **classification** model, the dependent variable is categorical (typically binary).

$$X \in \mathfrak{R} \quad \text{or} \quad Y \in \mathfrak{R}^n$$

$$X \in \{-1, +1\}$$

# “Decoding”: so what?

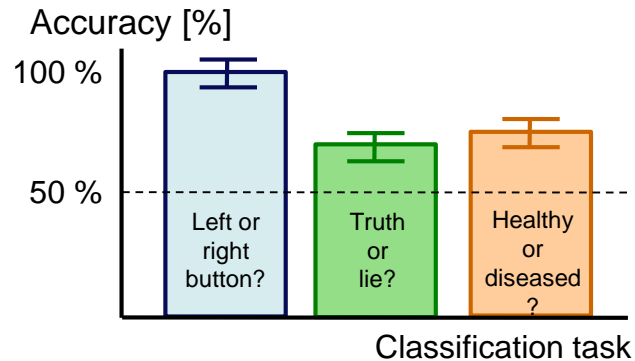
The seminal approach: classification



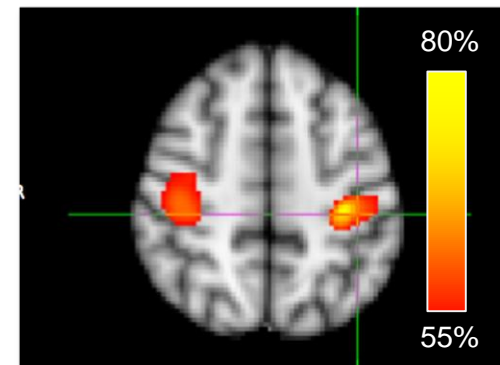
# “Decoding”: so what?

Reversing the X-Y mapping: target questions

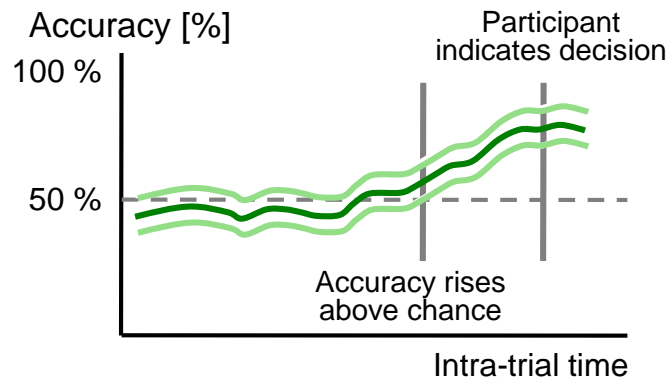
## (1) X-Y mapping overall reliability



## (2) X-Y mapping spatial deployment

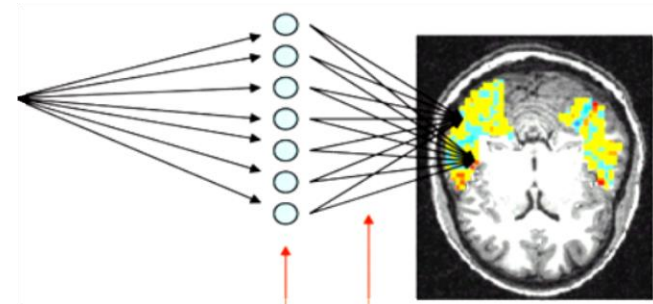


## (3) X-Y mapping temporal evolution



## (4) X-Y mapping: subtle issues

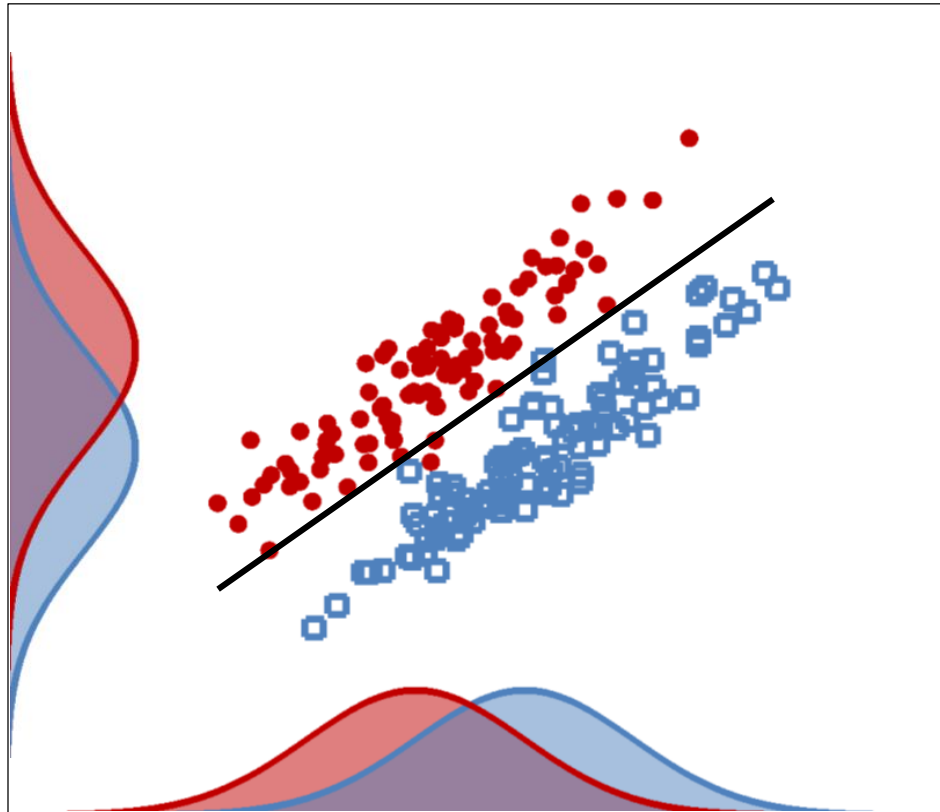
- functionally selective vs segregated representations
- degenerate (many-to-one) structure-function mappings



# Multivariate: so what?

Well, we might need it.

- Multivariate approaches can reveal information jointly encoded by several voxels. This is because multivariate distance measures can take into account correlations among voxels

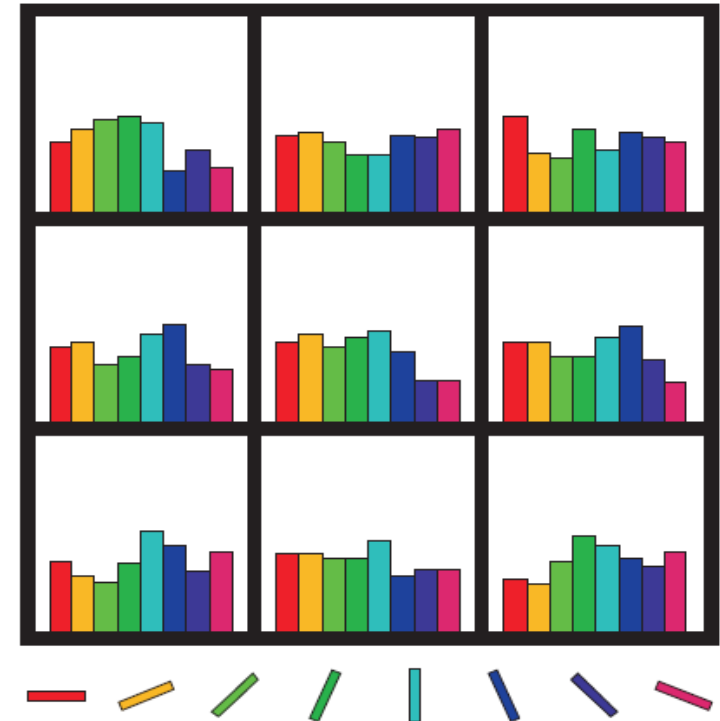
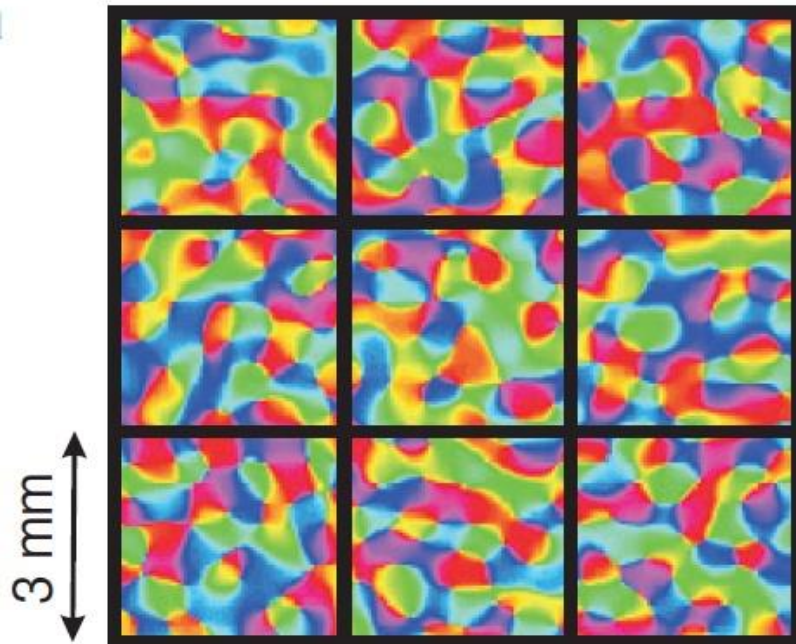




# Multivariate: so what?

Why we might need it: subvoxel processing.

- Multivariate approaches can exploit a sampling bias in voxelized images. Such subvoxel processing is unlikely to be detected by univariate methods.



Boynton 2005 *Nature Neuroscience*

# Preliminary statistical considerations

## lessons from the Neyman-Pearson lemma

- Do neuronal responses encode some sensorial or cognitive state of the subject?
- Null assumption: there is no dependency between Y and X

$$H_0 : p(Y|X) = p(Y)$$

- **Neyman-Pearson lemma**: the likelihood ratio (or Bayes factor)

$$\Lambda = \frac{p(Y|X)}{p(Y)} = \frac{p(X|Y)}{p(X)} \geq u$$

is the most powerful test of size  $\alpha$  to test the null

...choose threshold  $u$  such that  $P(\Lambda(Y) \geq u | H_0) = \alpha$ .

- So what? Well...

- 1 All we have to do is compare a model that links Y to X with a model that does not.
- 2 The link can be from X to Y or from Y to X. From the point of view of inferring a link exists, its direction is not important (but...).

# Preliminary statistical considerations

## prediction and inference

- Some confusion about the roles of prediction and inference may arise from the use of classification accuracy to infer a significant relationship between  $X$  and  $Y$ .

- This is because « cross-validation » relies on the predictive density:

$$p(X_{new} | Y_{new}, X, Y) = \int p(X_{new} | Y_{new}, \theta) p(\theta | X, Y) d\theta$$

where  $\theta$  are unknown parameters of the mapping  $Y \xrightarrow{\theta} X$   
to check the « generalization error » of the inferred mapping.

- Note:

- 1 The only situation that legitimately requires us to predict a new target is when we do not know it, e.g.:
  - brain-computer interface
  - automated diagnostic classification
- 2 When used in the context of experimental neuroscience, standard classifiers provide suboptimal inference on the mapping  $Y \rightarrow X$

# Preliminary statistical considerations

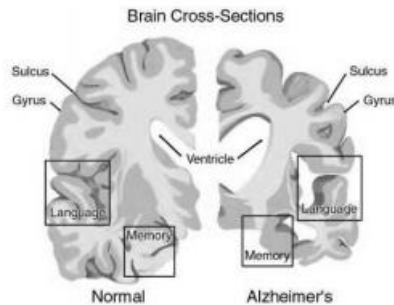
## prediction and inference

1

The goal of **prediction** is to find a highly accurate encoding or decoding function.



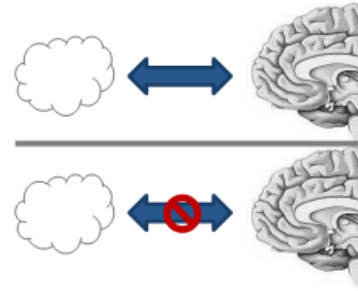
predicting a cognitive state using a brain-machine interface



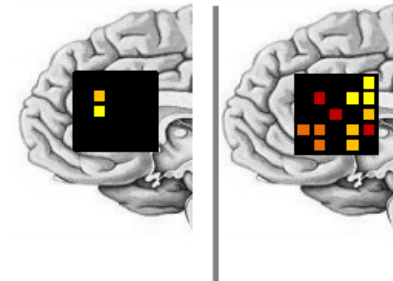
predicting a subject-specific diagnostic status

2

The goal of **inference** is to decide between competing hypotheses about structure-function mappings in the brain.



comparing a model that links distributed neuronal activity to a cognitive state with a model that does not \*



weighing the evidence for sparse coding vs. dense coding

predictive density

$$p(X_{new}|Y_{new}, X, Y) = \int p(X_{new}|Y_{new}, \theta) p(\theta|X, Y) d\theta$$

marginal likelihood

$$p(X|Y) = \int p(X|Y, \theta) p(\theta) d\theta$$

\* Although MVB alone does not provide this

# Overview of the talk

## 1 Introduction

*1.1 Lexicon*

*1.2 “Decoding”: so what?*

*1.3 Multivariate: so what?*

*1.4 Preliminary statistical considerations*

## 2 Multivariate Bayesian decoding

*2.1 From classical encoding to Bayesian decoding*

*2.2 Hierarchical priors on patterns*

*2.3 Probabilistic inference*

## 3 Example

## 4 Summary

# From classical encoding to Bayesian decoding

## MVB: inferring on the multivariate X-Y mapping

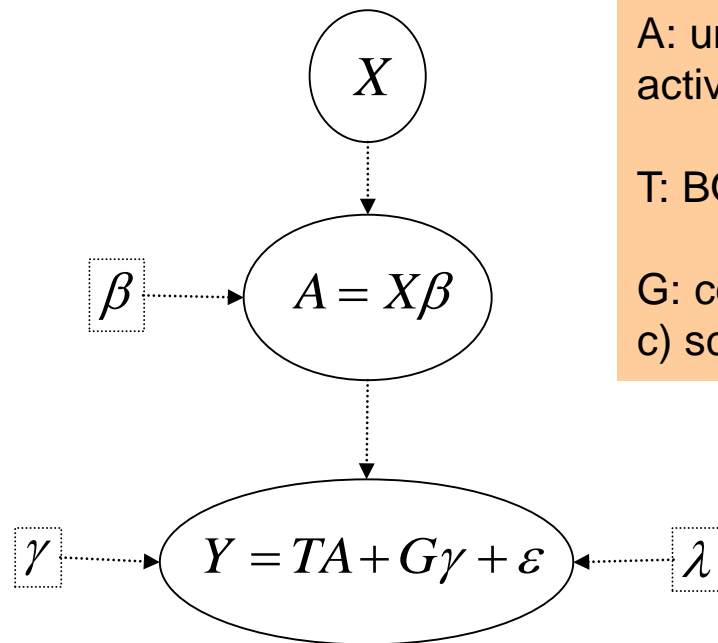
- Multivariate analyses in SPM are not implemented in terms of the classification schemes outlined in the previous section.
- Instead, SPM brings decoding into the conventional inference framework of hierarchical models and their inversion (c.f. Neyman-Pearson lemma).
- MVB can be used to address two questions:
  - **Overall significance of the X-Y mapping** (as with classical SPM or classifiers)  
... using probabilistic inference (model comparison, cross-validation)
  - **Inference on the form of the X-Y mapping** (no other alternative), e.g.
    - 1 Identify the spatial structure of the X-Y mapping (smooth, sparse, etc...)
    - 2 Tell whether the X-Y mapping is degenerate (many-to-one regions-to-function).
    - 3 Disambiguate between category-specific representations that are functionally selective (with overlap) and functionally segregated (without).

# From classical encoding to Bayesian decoding

reversing the standard GLM

## Encoding models

*X as a cause*



$$g(\theta) : X \rightarrow Y$$

$$Y = TX\beta + G\gamma + \varepsilon$$

X: scalar psychological target variable

Y: noisy measurements of A

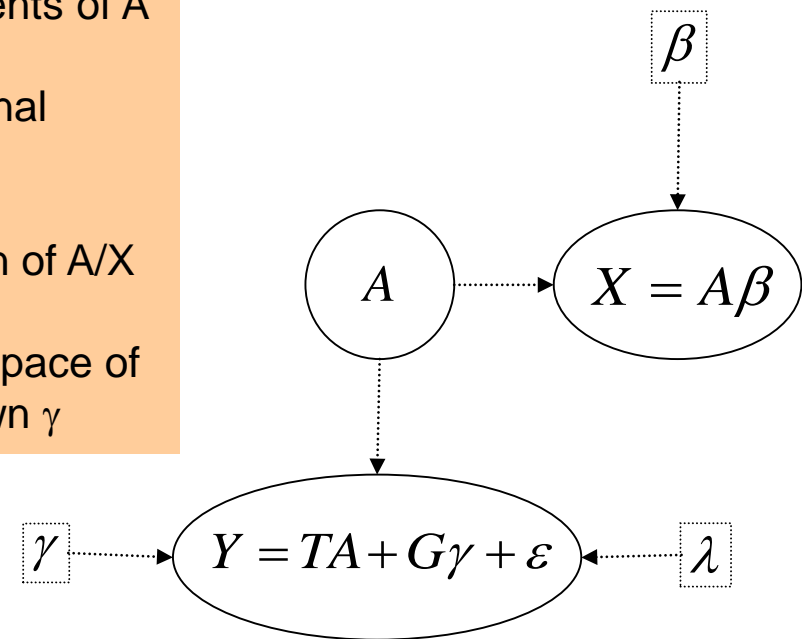
A: underlying neuronal activity in n voxels

T: BOLD convolution of A/X

G: confounds (null space of c) scaled by unknown  $\gamma$

## Decoding models

*X as a consequence*



$$g(\theta) : Y \rightarrow X$$

$$X = A\beta$$

$$TX = Y\beta - G\gamma\beta - \varepsilon\beta$$

# Hierarchical priors on patterns

spatial deployment of the X-Y mapping

- Decoding models are typically ill-posed: there is an infinite number of equally likely solutions. We therefore require constraints or priors to estimate the voxel weights  $\beta$ .
- MVB specifies several alternative coding hypotheses in terms of empirical spatial priors on voxel weights.

→ project onto spatial basis function set:

$$\beta = U\eta \quad \text{patterns}$$

$$\text{cov}(\beta) = U \text{cov}(\eta) U^T$$

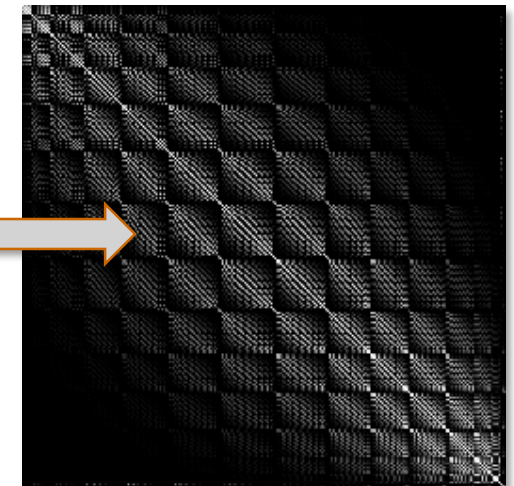
null:  $U = \emptyset$

sparse vectors:  $U = I$

smooth vectors:  $U(\vec{x}_i, \vec{x}_j) = \exp(-\frac{1}{2}(\vec{x}_i - \vec{x}_j)^2 \sigma^{-2})$

compact vectors:  $UDV^T = RY^T$  (compact: previously singular, now SVD of the smooth vectors)

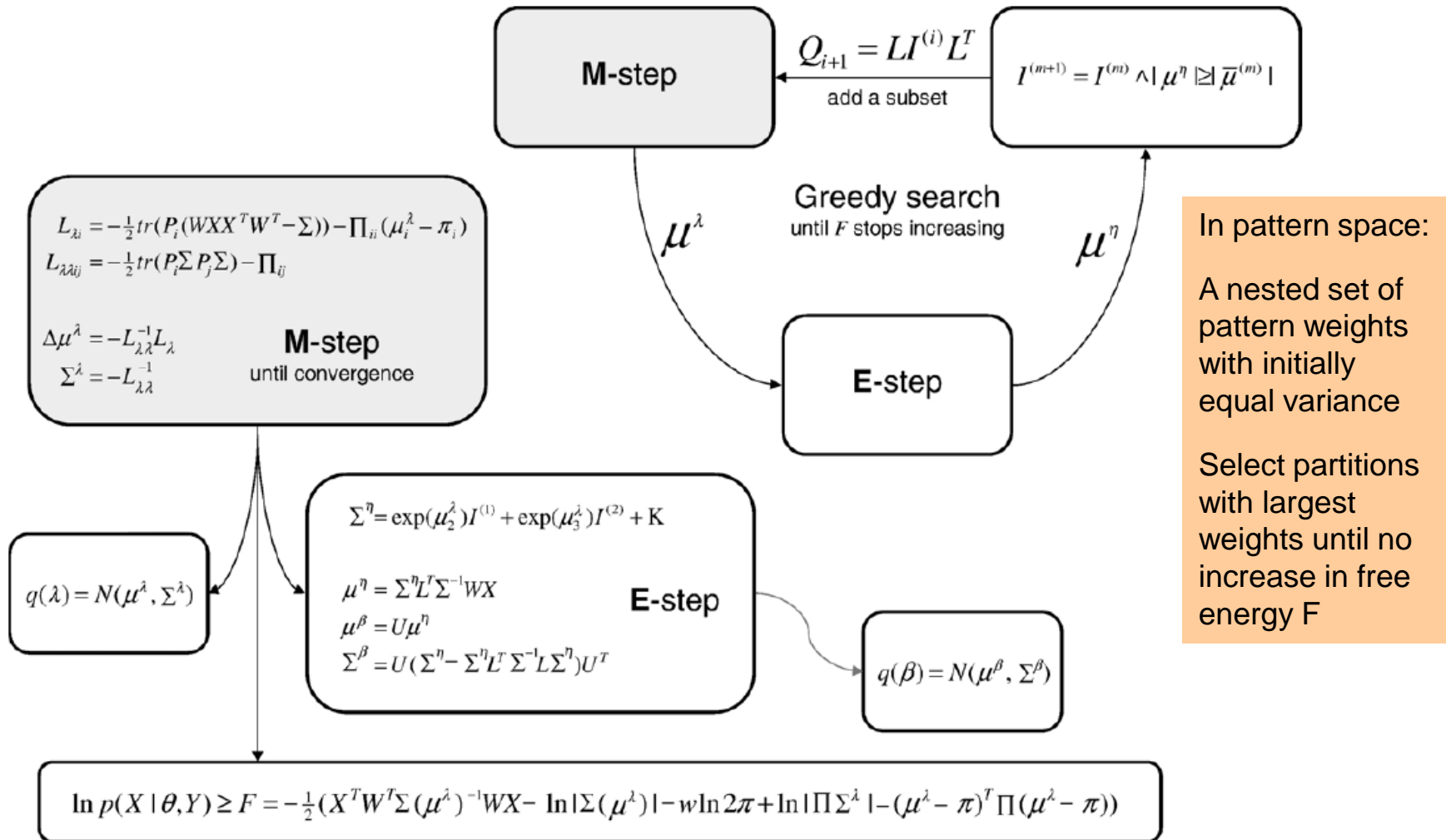
support vectors:  $U = RY^T$





# Hierarchical priors on patterns

Expectation-Maximization and the greedy search



Simplified EM algorithm (see Friston et al., 2008)

# Probabilistic inference

classical inference with cross-validation

- p-values from a standard leave-one-out scheme can't be used for inference (train and test data *are not* independent)
- Recall compact form for the decoding model:

$$WX = RY\beta + \zeta \quad \text{target variable}$$

$$W = RT \quad \text{weighting matrix: temporal convolution + confounds removal}$$

$$R = I - GG^{-} \quad \text{residual forming matrix: confounds removal}$$

- Use train/test k-fold data features that are linearly independent:

**train** (identify mapping)

$$\hat{\beta}_{(-k)} = \left\langle \beta \middle| Y_{(-k)} \right\rangle$$

$$Y_{(-k)} = R_{(-k)} Y$$

$$R_{(-k)} = \left( I - G_{(-k)} G_{(-k)}^{-} \right)$$

$$G_{(-k)} = \begin{bmatrix} G & I^{(k)} \end{bmatrix}$$

**test** (measure generalization error)

$$WX \stackrel{?}{=} \hat{X}_{(k)}$$

$$\hat{X}_{(k)} = R_{(k)} Y \hat{\beta}_{(-k)}$$

$$R_{(k)} = \left( I - G_{(k)} G_{(k)}^{-} \right)$$

$$G_{(k)} = \begin{bmatrix} G & I - I^{(k)} \end{bmatrix}$$

# Overview of the talk

## 1 Introduction

*1.1 Lexicon*

*1.2 “Decoding”: so what?*

*1.3 Multivariate: so what?*

*1.4 Preliminary statistical considerations*

## 2 Multivariate Bayesian decoding

*2.1 From classical encoding to Bayesian decoding*

*2.2 Hierarchical priors on patterns*

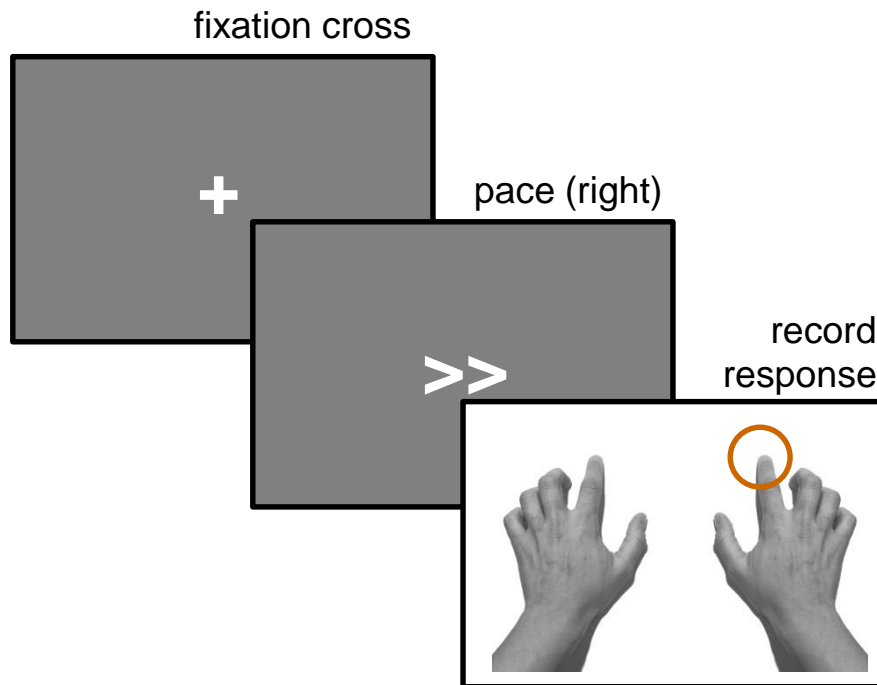
*2.3 Probabilistic inference*

## 3 Example

## 4 Summary

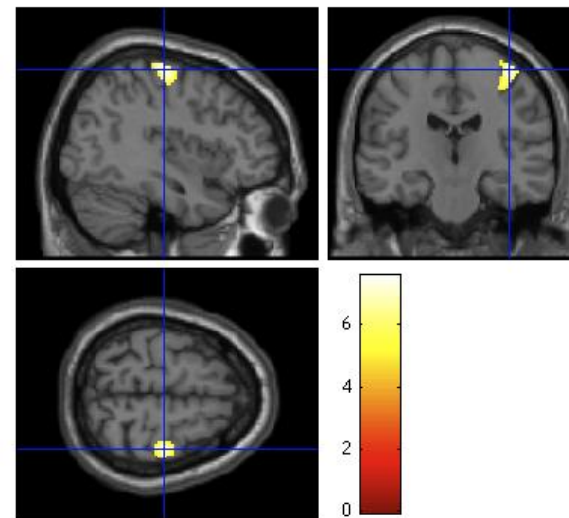
# Example

## finger tapping dataset

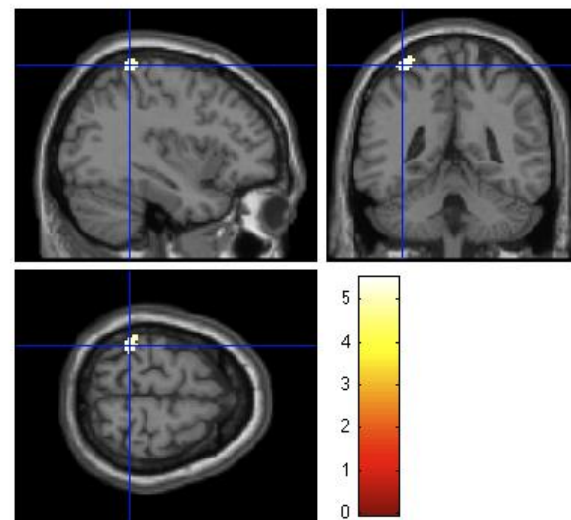


- 400 events (100 left, 100 right, 100 left & right, 100 null)
- average ITI = 2 sec
- block design (10 trials/block)
- TR = 1.3 sec

SPM{T} : left > right

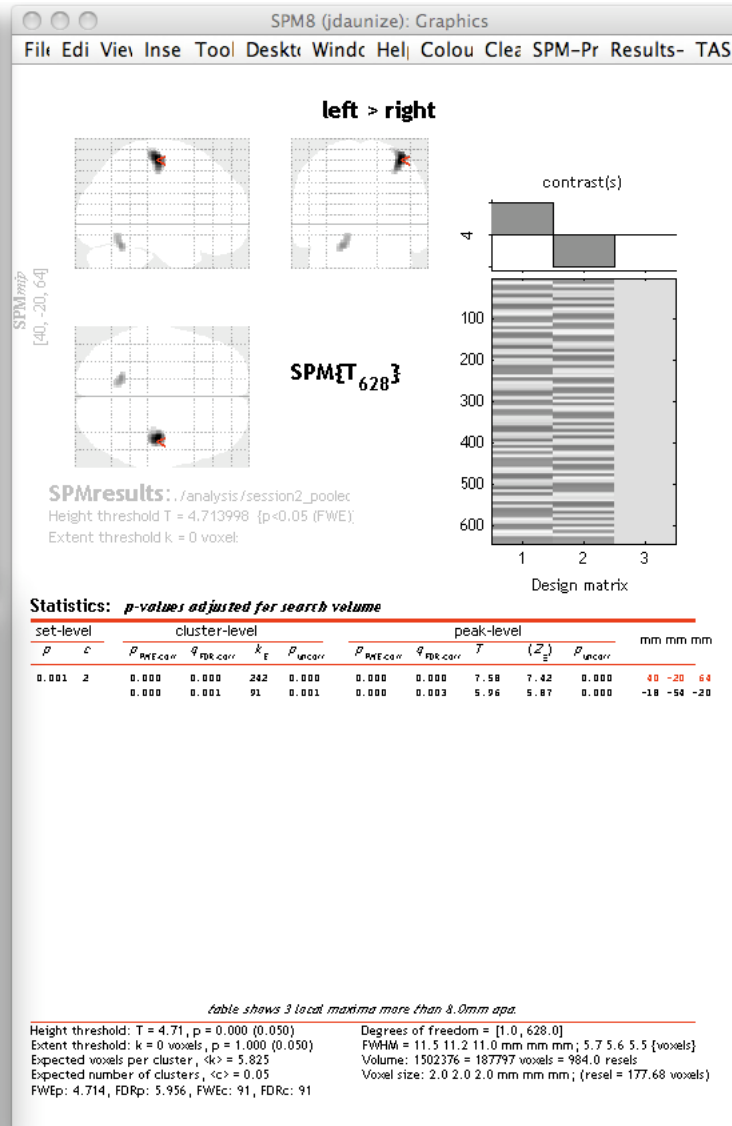
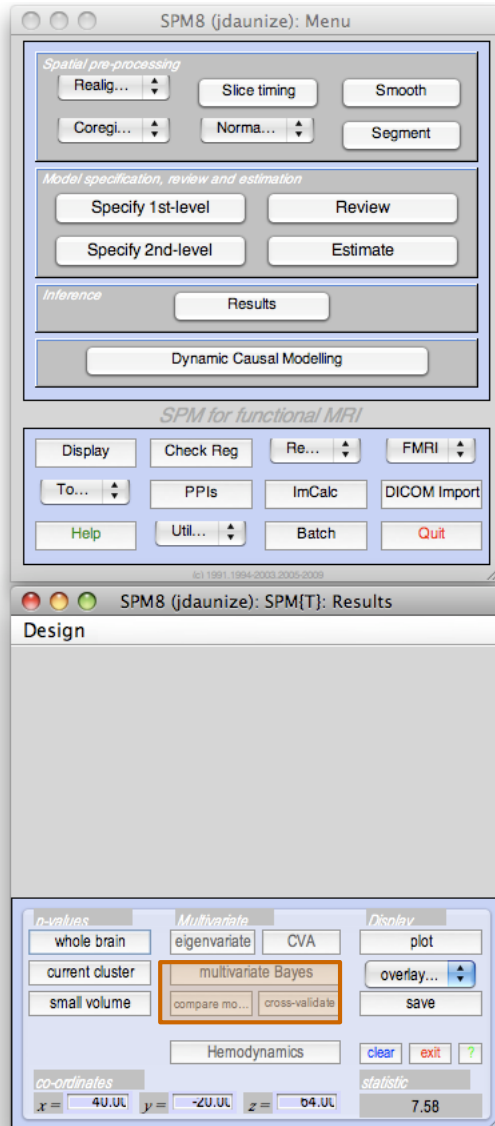


SPM{T} : right > left



# Example

## MVB in SPM: decoding within a search volume



- target:

$$TX = \mathbf{X}c$$

design  
matrix

contrast

- confounds:

$$G = \mathbf{X}(I - cc^{-1})$$

$$\Rightarrow Gc = 0$$

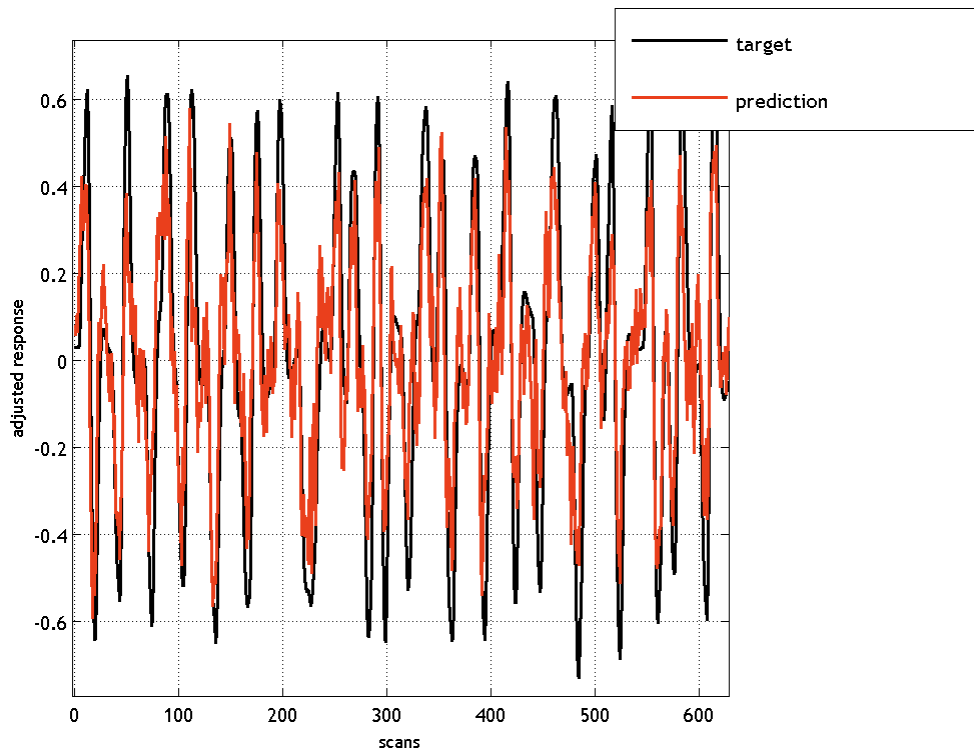
(confounds =  
null space of  
the contrast)

# Example

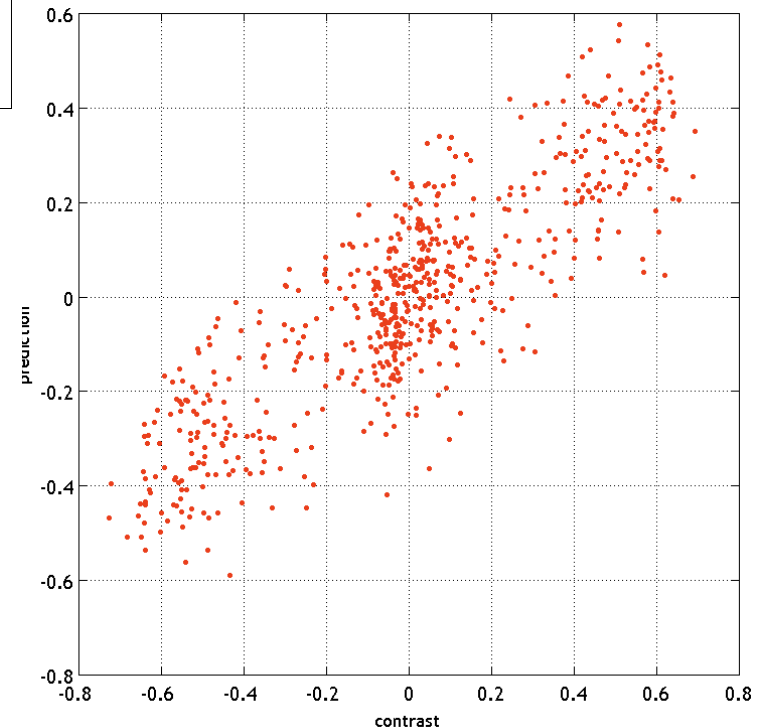
predicted responses from left & right motor cortices

- MVB-based predictions closely match the observed responses. But crucially, they don't perfectly match them. Perfect match would indicate overfitting.

target and (MVB) predicted responses over scans



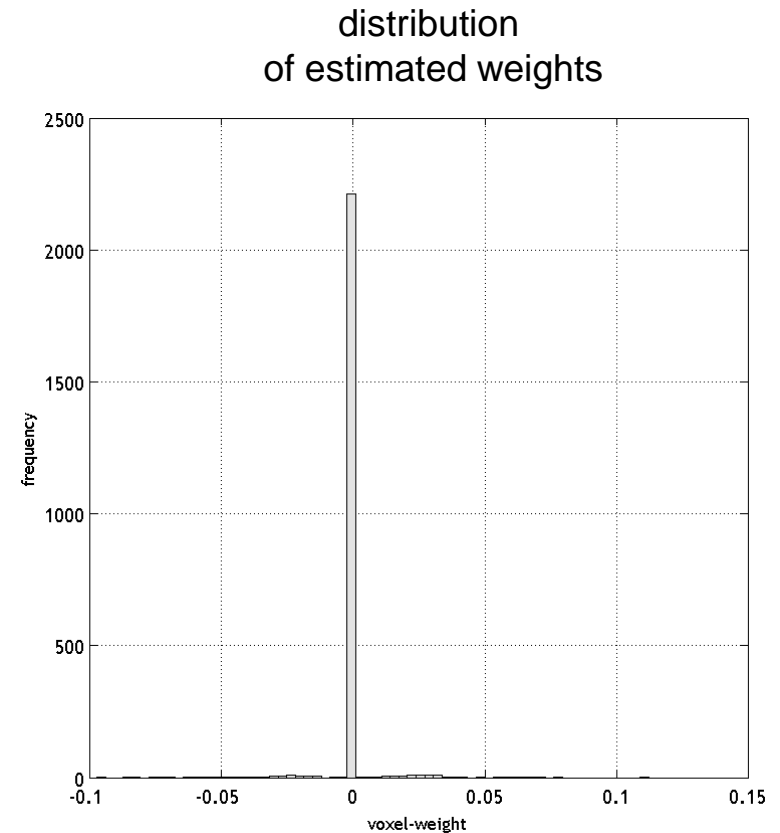
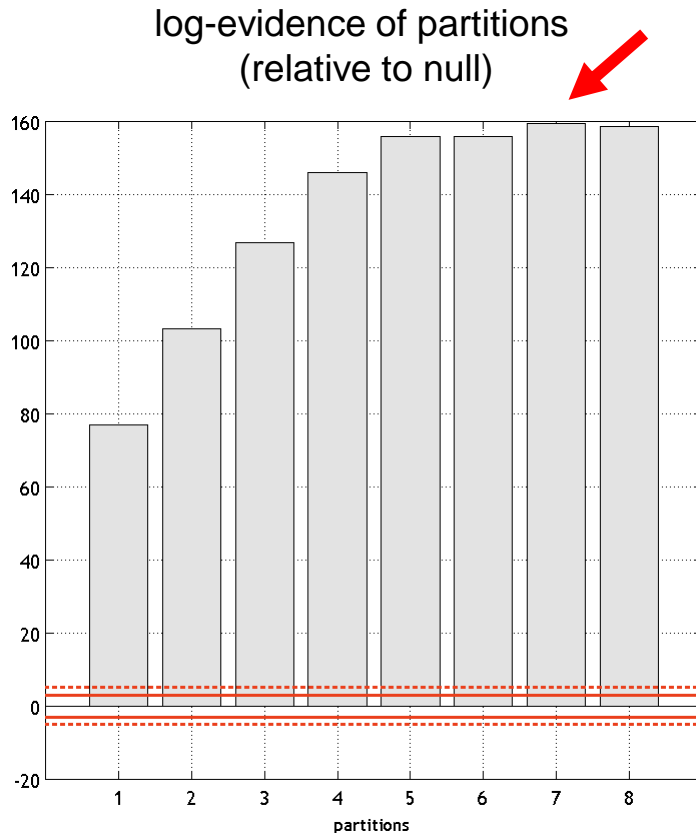
target versus (MVB) predicted responses



# Example

## pattern sparsity

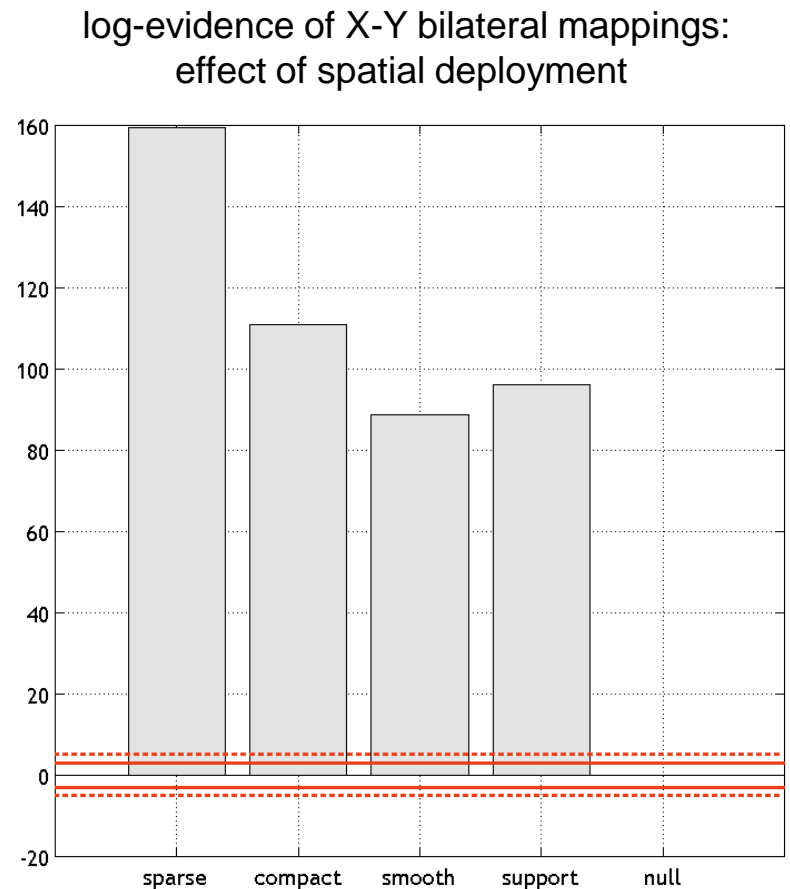
- The highest model evidence is achieved by a model that recruits 7 partitions. The weights attributed to each voxel in the sphere are sparse and bidirectional. This suggests sparse coding (in pattern space)



# Example

## model comparison illustration

- The best model corresponds to a sparse representation of motion ; as one would expect from functional segregation.
- Better evidence for spatially sparse than smooth (clustered) coding
- Sparse spatial model also better than compact (SVD on smooth model) or support models
- see also Morcom & Friston (2012), Neuroimage



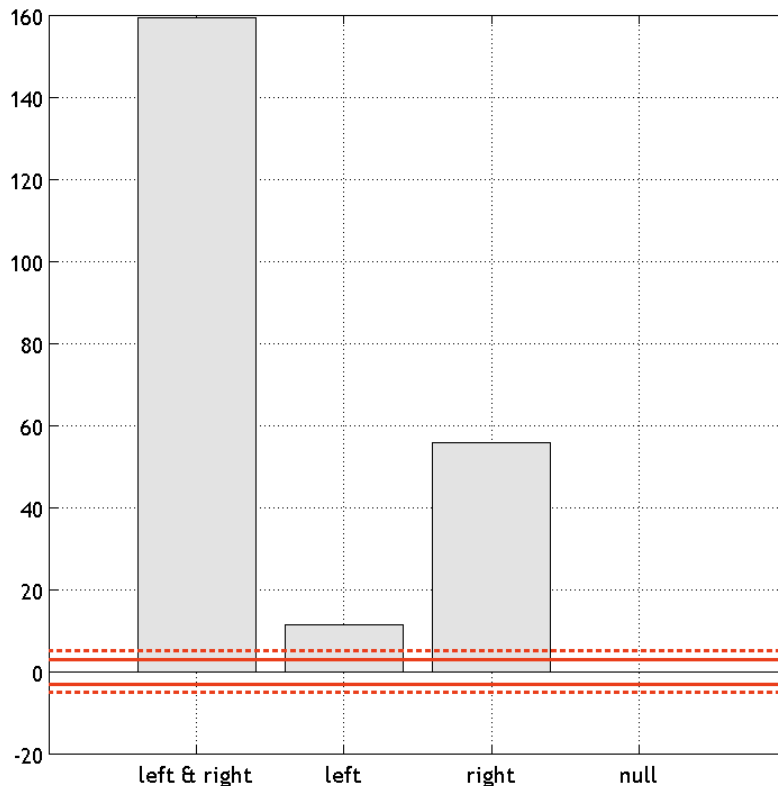


# Example

## model comparison illustration

- Model comparison between regions  
(given optimal – sparse – model of spatial deployment)

log-evidence of X-Y sparse mappings:  
effect of lateralization



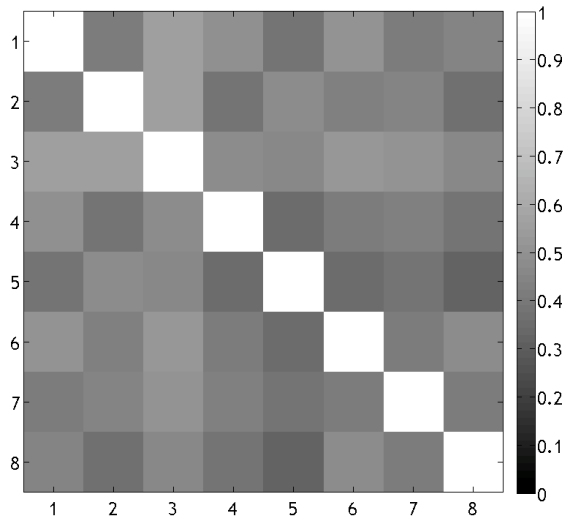
- The right-lateralized model is better than the left-lateralized model
- The bilateral model is better than the right-lateralized model
- Consistent with non-redundant (joint) coding of finger tapping

# Example

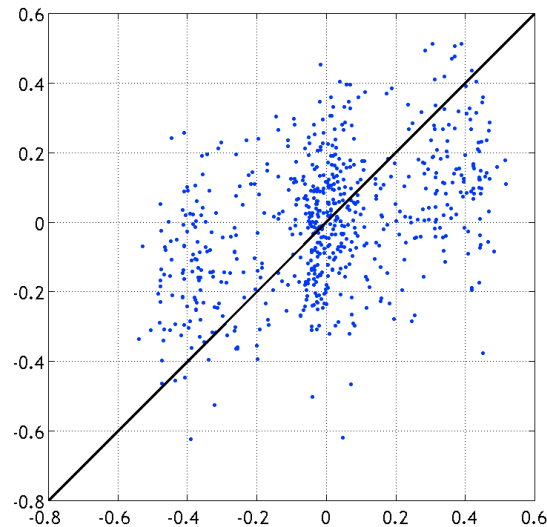
cross-validation : k-fold scheme

- $k = 8$
- p-value < 0.0001
- classification accuracy = 65.8%
- R-squared = 20.7%

absolute correlation among  
k-fold feature weights

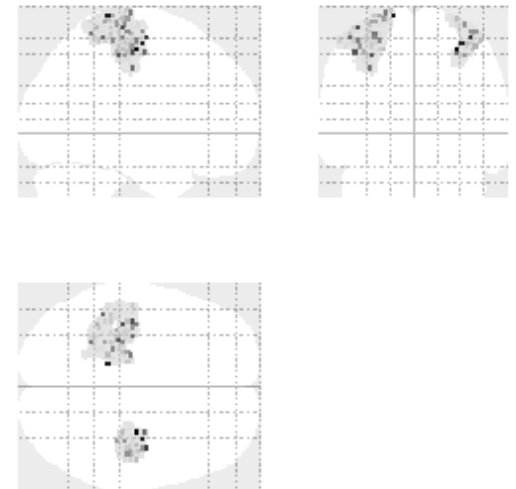


test predictions versus  
test k-fold features



maximum intensity projection:

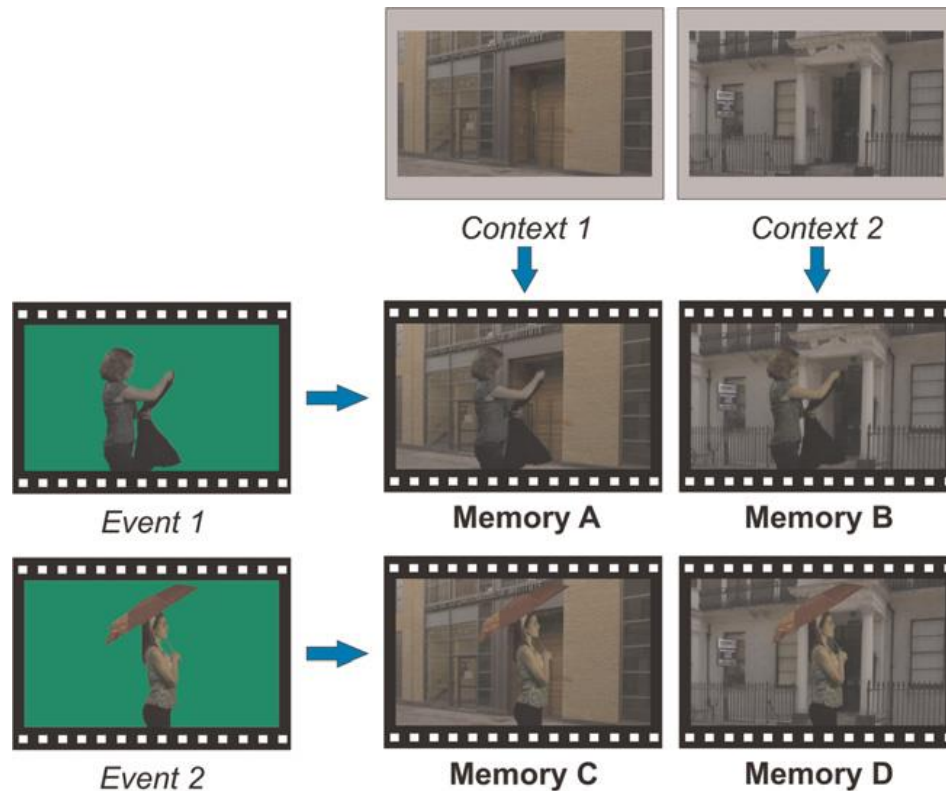
$$\prod_k P(|\beta| > 0 | Y_{(-k)})$$



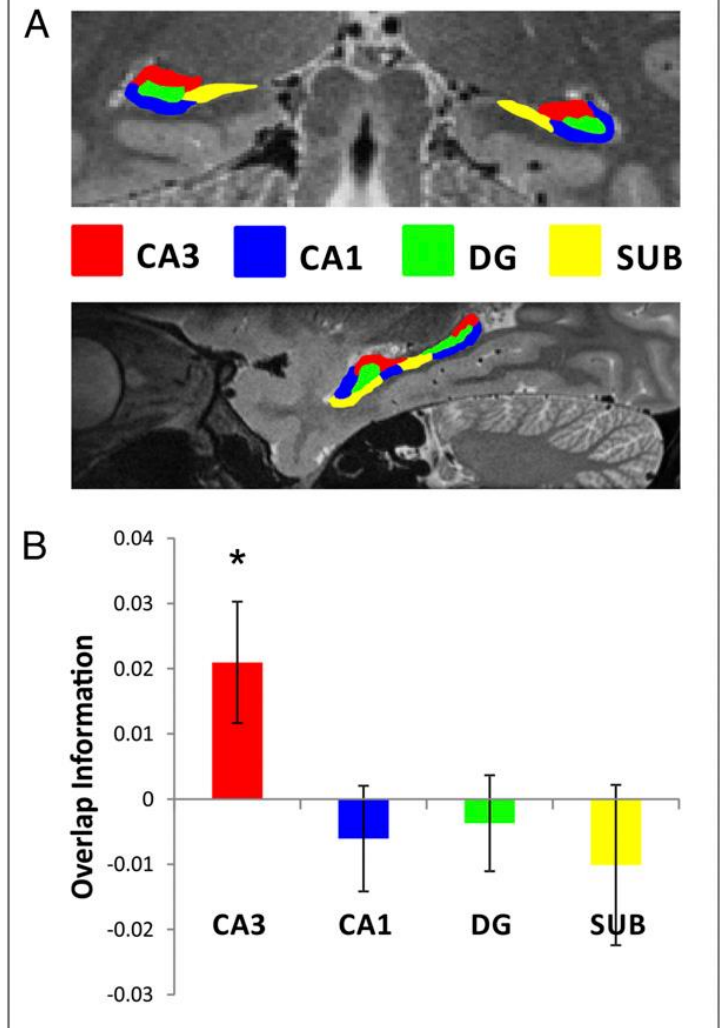
# Example

episodic memory, MTL

Distinct spatial patterns for overlapping memories only in hippocampus CA3, consistent w/ pattern separation



- Overlapping episodic memories
- Same events, different contexts (scenes)
- Use MVB to fit sparse patterns, then measure their correlation over voxels in MTL sub-regions
- **A form of representational similarity analysis**



Chadwick, Bonnici, Maguire (2014, PNAS)

# Overview of the talk

## 1 Introduction

*1.1 Lexicon*

*1.2 “Decoding”: so what?*

*1.3 Multivariate: so what?*

*1.4 Preliminary statistical considerations*

## 2 Multivariate Bayesian decoding

*2.1 From classical encoding to Bayesian decoding*

*2.2 Hierarchical priors on patterns*

*2.3 Probabilistic inference*

## 3 Example

## 4 Summary

# Summary

- 1 Inference on the form of the X-Y mapping rests on model comparison, using the marginal likelihood of competing models. The marginal likelihood derives from the specification of a generative model prescribing the form of the joint density over observations (X,Y) and model parameters ( $\theta$ ).
- 2 Multivariate models can map from experimental variables (X) to brain responses (Y) or from Y to X. In the latter case (i.e., decoding), identifying the mapping is an ill-posed problem, which is resolved with appropriate constraints or priors on model parameters. These constraints are part of the model and can be evaluated using model comparison.
- 3 Cross-validation is not necessary for decoding brain activity but generalization error is a proxy for testing whether the observed X-Y mapping is unlikely to have occurred by chance. This can be useful when the null distribution of the likelihood ratio (i.e. Bayes factor) is not evaluated easily.

## References

Chadwick, M. J., Bonnici, H. M. and Maguire, E. a. (2014) 'CA3 size predicts the precision of memory recall', *Proceedings of the National Academy of Sciences*, 111(29), pp. 10720–5. doi: 10.1073/pnas.1319641111.

Friston, K., Chu, C., Mouruo-Miranda, J., Hulme, O., Rees, G., Penny, W. and Ashburner, J. (2008) 'Bayesian decoding of brain images', *NeuroImage*, 39(1), pp. 181–205. doi: 10.1016/j.neuroimage.2007.08.013.

Friston, K., Mattout, J., Trujillo-Barreto, N., Ashburner, J. and Penny, W. (2007) 'Variational free energy and the Laplace approximation', *NeuroImage*, 34(1), pp. 220–234. doi: 10.1016/j.neuroimage.2006.08.035.

Morcom, A. M. and Friston, K. J. (2012) 'Decoding episodic memory in ageing: A Bayesian analysis of activity patterns predicting memory', *NeuroImage*, 59(2). doi: 10.1016/j.neuroimage.2011.08.071.