EXTRACTING FEATURES FROM SMRI

John Ashburner

Wellcome Trust Centre for Neuroimaging, UCL Institute of Neurology, 12 Queen Square, London WC1N 3BG, UK.

イロト イポト イヨト イヨト

э

MEASURING GENERALISATION ACCURACY INCORPORATING PRIOR KNOWLEDGE "DATA-DRIVEN FEATURE SELECTION"

FEATURE ENGINEERING

First-timers are often surprised by how little time in a machine learning project is spent actually doing machine learning. But it makes sense if you consider how time-consuming it is to gather data, integrate it, clean it and pre-process it, and how much trial and error can go into feature design. Also, machine learning is not a one-shot process of building a data set and running a learner, but rather an iterative process of running the learner, analyzing the results, modifying the data and/or the learner, and repeating. Learning is often the quickest part of this, but that's because we've already mastered it pretty well! Feature engineering is more difficult because it's domain-specific, while learners can be largely general-purpose. However, there is no sharp frontier between the two, and this is another reason the most useful learners are those that facilitate incorporating knowledge.

Domingos, Pedro. "A few useful things to know about machine learning." Communications of the ACM 55, no. 10 (2012): 78-87.

イロト イポト イヨト イヨト

MEASURING GENERALISATION ACCURACY INCORPORATING PRIOR KNOWLEDGE "DATA-DRIVEN FEATURE SELECTION"

ACCURACY

- Proportion of guesses that are correct.
- Assessed by cross-validation.
- A very simple measure of generalisation.
- Very noisy.

If 90% of subjects are controls and 10% are patients, then guessing that everyone is a control will give 90% accuracy.

Other measures (sensitivity, specificity, etc) will be discussed later.

< ロ > < 同 > < 回 > < 回 > < 回 >

MEASURING GENERALISATION ACCURACY INCORPORATING PRIOR KNOWLEDGE "DATA-DRIVEN FEATURE SELECTION"

AREA UNDER THE CURCE (AUC)

Area under the Receiver Operating Characteristic (ROC) curve. Assessed by cross-validation.





イロト イポト イヨト イヨト

INTRODUCTION

MEASURING GENERALISATION ACCURACY

HARD V PROBABILISTIC CLASSIFICATION



Probabilistic classification

(ロ) (部) (目) (日) (日)

MEASURING GENERALISATION ACCURACY INCORPORATING PRIOR KNOWLEDGE "DATA-DRIVEN FEATURE SELECTION"

TARGET INFORMATION

Using cross-validation with binary classification, the number of *bits* of information obtained for each subject is:

$$I = rac{1}{N} \sum_{n=1}^{N} \left(t_n \log_2 p_n + (1-t_n) \log_2 (1-p_n)
ight) \ - \left(\overline{t} \log_2 \overline{t^*} + (1-\overline{t} \log_2 (1-\overline{t^*}))
ight)$$

where t_n is the label of the *n*th test subject (0 or 1)

 p_n is the predicted probability for the *n*th test subject $\bar{t^*}$ is the average of the labels of the training data.

A similar scheme may be used for regression, where information is given in *nats* (used log_e , rather than log_2).

< 日 > < 同 > < 回 > < 回 > < 回 >

MEASURING GENERALISATION ACCURACY INCORPORATING PRIOR KNOWLEDGE "DATA-DRIVEN FEATURE SELECTION"

Log Marginal Likelihood

Bayesian methods give a measure known as log marginal likelihood.

$$P(\mathbf{y}|\mathbf{X}) = \int_{\mathbf{w}} P(\mathbf{y}|\mathbf{X}, \mathbf{w}) p(\mathbf{w}) d\mathbf{w}$$

- An established Bayesian model selection approach (see papers by David MacKay and others).
- Does not involve cross-validation.
- Not trusted by some machine learning people.

< 日 > < 同 > < 回 > < 回 > < 回 >

Measuring Generalisation Accuracy Incorporating Prior Knowledge "Data-driven Feature Selection"

NO FREE DUCKLINGS

No Free Lunch theorem says that

learning is impossible without prior

knowledge.

http://en.wikipedia.org/wiki/No_free_lunch_in_search_and_ optimization

Ugly Duckling theorem says that things are all equivalently similar to each other without prior knowledge. http://en.wikipedia.org/wiki/Ugly_duckling_theorem



Ryan Ebert from Portland, US (Flickr) [CC BY 2.0], via Wikimedia Commons. https://creativecommons.org/licenses/by/2.0/

イロト イポト イヨト イヨト

What prior knowledge do we have about the variability among people that can be measured using MRI? How do we use this knowledge?

INCORPORATING PRIOR KNOWLEDGE INTO KERNELS

Linear kernel matrices are often computed from the raw features:

$$\mathbf{K} = \mathbf{X}\mathbf{X}^{T}$$

A simple spatial feature selection may be considered as the following, where Σ_0 is a (scaled) diagonal matrix of ones and zeros:

$$\mathbf{K} = \mathbf{X} \mathbf{\Sigma}_0 \mathbf{X}^T$$

 Σ_0 may be more complicated, for example encoding spatial smoothing, high-pass filtering or any number of other things.

- 4 同 6 4 日 6 4 日 6

WEIGHTING SUSPECTED REGIONS MORE HEAVILY

- The best way would be to augment the training data with data from previous studies.
- Lack of data-sharing means this is generally not possible, so we need to extract information from publications.
- The neuroimaging literature is mostly blobs.
- These give pointers about how best to weight the data $(\Sigma_0 = diag(\mathbf{s}), s_i \in \mathbb{R}^+).$

イロト 人間ト イヨト イヨト

Measuring Generalisation Accuracy Incorporating Prior Knowledge "Data-driven Feature Selection"

WEIGHTING SUSPECTED REGIONS MORE HEAVILY



Chu et al. "Does feature selection improve classification accuracy? Impact of sample size and feature selection on classification using anatomical magnetic resonance images". NeuroImage 60:59–70 (2012).

イロト イポト イヨト イヨト

MEASURING GENERALISATION ACCURACY INCORPORATING PRIOR KNOWLEDGE "DATA-DRIVEN FEATURE SELECTION"

Smoothing

If we know that higher frequency signal is more likely to be noise.

 $\mathbf{K} = \mathbf{X} \mathbf{\Sigma}_0 \mathbf{X}^T$

 Σ_0 no longer diagional.



(ロ) (部) (目) (日) (日)

"DATA-DRIVEN FEATURE SELECTION"

Two main approaches:

• Non-embedded feature selection, where approaches such as t- or F-tests, or *recursive feature elimination* are used to switch off certain features. Not very principled, but can save computation time.

We should only do feature selection if there is a cost associated with measuring features or predicting with many features. Note: Radford Neal won the NIPS feature selection competition using Bayesian methods that used 100% of the features.

— Zoubin Ghahramani

イロト 人間ト イヨト イヨト

• Embedded feature selection, where features are weighted differently as part of the machine learning model. Works best when features are of different types so need different weighting (*a priori*).

DIMENSIONALITY \neq NUMBER OF VOXELS

- Lots of effort on data-driven feature selection methods.
 - Involves estimating

$$\boldsymbol{\Sigma}_0 = diag(\mathbf{s}), s_i \in \{0, w\}, \text{ where } w \in \mathcal{R}^+.$$

- Lots of parameters needed to achieve this.
- Many papers claim excellent results.
- Little evidence to suggest that most voxel-based feature selection methods help.
 - Little or no increase in predictive accuracy.
 - Commonly perceived as being more "interpretable".

< 日 > < 同 > < 回 > < 回 > < 回 >

Measuring Generalisation Accuracy Incorporating Prior Knowledge "Data-driven Feature Selection"

"DATA-DRIVEN FEATURE SELECTION"

"In our evaluation, two methods included a feature selection step: Voxel-STAND and Voxel-COMPARE. Overall, these methods did not perform substantially better than simpler ones... ... A more robust way to decrease the dimensionality of the features way would be to use more prior knowledge of the disease."

Cuingnet et al. "Automatic classification of patients with Alzheimer's disease from structural MRI: A comparison of ten methods using the ADNI database". NeuroImage 56(2):766–781 (2011).

- 4 同 6 4 日 6 4 日 6

Measuring Generalisation Accuracy Incorporating Prior Knowledge "Data-driven Feature Selection"

"DATA-DRIVEN FEATURE SELECTION"



Chu et al. "Does feature selection improve classification accuracy? Impact of sample size and feature selection on classification using anatomical magnetic resonance images". NeuroImage 60:59–70 (2012).

(日) (周) (王) (王)

3

Measuring Generalisation Accuracy Incorporating Prior Knowledge "Data-driven Feature Selection"

REMOVING NONLINEARITIES

Instead of using nonlinear pattern recognition methods, we can...

- Capture nonlinearities by appropriate preprocessing.
 - Accurate nonlinear registration can remove much of the nonlinearity.
- Allows nonlinear effects to be modelled by a linear classifier.
- Gives more interpretable characterisations of differences.
- May lead to more accurate predictions particularly with smaller amounts of training data.

- 4 同 6 4 日 6 4 日 6

Measuring Generalisation Accuracy Incorporating Prior Knowledge "Data-driven Feature Selection"

Removing nonlinearities

Simulated images



Principal components



A suitable model would reduce this variability to two dimensions.

Aligned Tissue Maps Deformation Features Scalar Momentum

RAW PIXEL VALUES

Raw pixel data could be another option. Data needs to be "spatially normalised" (and possibly skull-stripped). Results may not generalise well to data from other scanners.



(ロ) (部) (目) (日) (日)

Aligned Tissue Maps Deformation Features Scalar Momentum

REGION VOLUMES

Label propagation or other methods can be used to subdivide brain into regions.



(ロ) (部) (目) (日) (日)

Aligned Tissue Maps Deformation Features Scalar Momentum

OTHER FEATURES

Other features include:

- Cortical thickness.
- Shape features.
- PCA/ICA weights.
- Lesion maps.

etc



<ロ> (日) (日) (日) (日) (日)

ALIGNED TISSUE MAPS DEFORMATION FEATURES SCALAR MOMENTUM

SPM12 PROCESSING

Tissue class segmentation



Alignment with Shoot



ALIGNED TISSUE MAPS DEFORMATION FEATURES SCALAR MOMENTUM

"Unmodulated" GM, WM & BG



Pattern recognition run using: GM alone; WM alone; BG alone; GM + WM; GM + WM + BG.

Aligned Tissue Maps Deformation Features Scalar Momentum

"Modulated" GM, WM & BG



Pattern recognition run using: GM alone; WM alone; BG alone; GM + WM; GM + WM + BG.

JOHN ASHBURNER ANATOMICAL FEATURE REPRESENTATION

Aligned Tissue Maps Deformation Features Scalar Momentum

JACOBIAN DETERMINANTS



Encodes relative volumes before and after warping.

ヘロト 人間 ト 人 ヨト 人 ヨトー

John Ashburner

ANATOMICAL FEATURE REPRESENTATION

Aligned Tissue Maps Deformation Features Scalar Momentum

LOGARITHMS OF JACOBIAN DETERMINANTS



There are sometimes simple logarithmic relationships among volumes.



Zhang and Sejnowski. "A universal scaling law between gray matter and white matter of cerebral cortex." Proceedings of the National Academy of Sciences 97(10):5621–5626 (2000).

<ロ> (日) (日) (日) (日) (日)

John Ashburner

ANATOMICAL FEATURE REPRESENTATION

Aligned Tissue Maps Deformation Features Scalar Momentum

DIVERGENCES OF VELOCITY FIELDS



Very similar to logarithms of Jacobians. Not easy to explain.

<ロ> (日) (日) (日) (日) (日)

FEATURE TYPES

Scalar Momentum

Scalar Momentum

 $\mathbf{a} = |D\phi|(\boldsymbol{\mu} - \mathbf{c}(\phi))$



Scalar momentum





Vector momentun

Mean image gradients

Deformation







JOHN ASHBURNER

ANATOMICAL FEATURE REPRESENTATION

イロト イヨト イヨト イヨト

3

Aligned Tissue Maps Deformation Features Scalar Momentum

Scalar Momentum

$$\mathbf{a} = |D\phi|(\boldsymbol{\mu} - \mathbf{c}(\phi))|$$



SPM12 GUI for scalar

momentum.

		Batch Editor	-		×	
E	le Edit View SPM Bi	id0				
n	📽 🖬 🕨					
	Module List	Current Module: Generate Scalar Momenta				
	Generate Scalar Mon	 Help on: Generate Scalar Momenta 			•	
		TemplateRE/mprage/Templa	ste 4.r	1Î		
		Images 1	10 110	~		
		. Images 1	48 file	ŝ		
		Deformation fields 1	48 file	s		
		Smoothing	48 file 10mr	s n		
					•	
	•					
	Constants Contra University					
	Generate spatially smoothed ""scalar momenta" in a form suitable for using with pattern					
	recognition. In principle, a Gaussian Process model can be used to determine the optimal					
	(positive) inear combination of kernel matrices. The idea would be to combined kernel — matrix derived from these, with a kernel derived from the velocity-fields. Such a combined kernel should then encode more relevant information than the individual kernels alone.					
	The scalar momentur	fields that are generated contain a number of volumes equ	al to t	he	•	

<ロ> (日) (日) (日) (日) (日)

ANATOMICAL FEATURE REPRESENTATION

IXI DATASET ABIDE I DATASET COBRE DATASET

IXI: DATASET

580 T1w brain MRI from IXI (Information eXtraction from Images) dataset. http://www. brain-development.org/ Data from three different hospitals in London:

- Hammersmith Hospital using a Philips 3T system
- Guy's Hospital using a Philips 1.5T system
- Institute of Psychiatry using a GE 1.5T system

10-fold cross-validation.



ANATOMICAL FEATURE REPRESENTATION

IXI DATASET ABIDE I DATASET COBRE DATASET



IXI DATASET ABIDE I DATASET COBRE DATASET



IXI DATASET ABIDE I DATASET COBRE DATASET



IXI DATASET ABIDE I DATASET COBRE DATASET



IXI DATASET ABIDE I DATASET COBRE DATASET



IXI DATASET ABIDE I DATASET COBRE DATASET


IXI DATASET ABIDE I DATASET COBRE DATASET



IXI DATASET ABIDE I DATASET COBRE DATASET



IXI DATASET ABIDE I DATASET COBRE DATASET



IXI DATASET ABIDE I DATASET COBRE DATASET



IXI DATASET ABIDE I DATASET COBRE DATASET



IXI DATASET ABIDE I DATASET COBRE DATASET



IXI DATASET ABIDE I DATASET COBRE DATASET



IXI DATASET ABIDE I DATASET COBRE DATASET



IXI DATASET ABIDE I DATASET COBRE DATASET

ABIDE: DATASET

The Autism Brain Imaging Data Exchange initiative. http://fcon_1000.projects.nitrc.org/indi/abide/.

T1w brain MRI from 1,102 subjects.

- 531 with Autism Spectrum Disorder (Gender ratio: 64:467).
- 571 controls (Gender ratio: 99:472).

Data from 17 international sites.

The 20 greatest outliers were excluded.

5-fold cross-validation.

IXI DATASET ABIDE I DATASET COBRE DATASET



IXI DATASET ABIDE I DATASET COBRE DATASET



IXI DATASET ABIDE I DATASET COBRE DATASET



IXI DATASET ABIDE I DATASET COBRE DATASET



IXI DATASET ABIDE I DATASET COBRE DATASET



IXI DATASET ABIDE I DATASET COBRE DATASET



INTRODUCTION Feature Types Data Conclusions

IXI DATASET ABIDE I DATASET COBRE DATASET



IXI DATASET ABIDE I DATASET COBRE DATASET



IXI DATASET ABIDE I DATASET COBRE DATASET



IXI DATASET ABIDE I DATASET COBRE DATASET



IXI DATASET ABIDE I DATASET COBRE DATASET



IXI DATASET ABIDE I DATASET COBRE DATASET



IXI DATASET ABIDE I DATASET COBRE DATASET

COBRE: DATASET

Centre for Biomedical Research Excellence http:

//fcon_1000.projects.nitrc.org/indi/retro/cobre.html.

T1w brain MRI from 146 subjects.

- 72 with schizophrenia (14 male : 58 female).
- 74 controls (23 male : 51 female).

All from a single scanner.

5-fold cross-validation, repeated 10 times.

- 4 同 6 4 日 6 4 日 6

IXI DATASET ABIDE I DATASET COBRE DATASET



IXI DATASET ABIDE I DATASET COBRE DATASET



IXI DATASET ABIDE I DATASET COBRE DATASET



IXI DATASET ABIDE I DATASET COBRE DATASET



IXI DATASET ABIDE I DATASET COBRE DATASET



IXI DATASET ABIDE I DATASET COBRE DATASET



IXI DATASET ABIDE I DATASET COBRE DATASET



IXI DATASET ABIDE I DATASET COBRE DATASET



IXI DATASET ABIDE I DATASET COBRE DATASET



IXI DATASET ABIDE I DATASET COBRE DATASET



IXI DATASET ABIDE I DATASET COBRE DATASET

COBRE: SCHIZ. V. CONTROL (SVC)



IXI DATASET ABIDE I DATASET COBRE DATASET

COBRE: SCHIZ. V. CONTROL (SVC)



IXI DATASET ABIDE I DATASET COBRE DATASET

COBRE: SCHIZ. V. CONTROL (GPC)



IXI DATASET ABIDE I DATASET COBRE DATASET

COBRE: SCHIZ. V. CONTROL (GPC)


IXI DATASET ABIDE I DATASET COBRE DATASET

COBRE: SCHIZ. V. CONTROL (GPC)



IXI DATASET ABIDE I DATASET COBRE DATASET

COBRE: SCHIZ. V. CONTROL (GPC)



OVERALL SCORES



John Ashburner

ANATOMICAL FEATURE REPRESENTATION

OVERALL SCORES



John Ashburner

ANATOMICAL FEATURE REPRESENTATION

CONCLUSIONS

- No feature set was best in all situations (no free lunch).
- Scalar momentum appears to be a useful feature set, although its effectiveness was not statistically significantly better than other methods that also considered the BG class.
- Jacobian-scaled warped GM alone, or with WM, is surprisingly poor.
- Amount of spatial smoothing makes a difference, with the best results from smoothing of about 12mm FWHM.
- Further dependencies on the details of the registration still need exploring.

Thanks to help from Gemma Monté. Much of this work has been submitted to NeuroImage.

(4 回) (4 回) (4 回)