

Machine-Learning for Brain Signal Analysis

Vincent Guigue
vincent.guigue@lip6.fr

September 9th 2016

SMART Summer School

○ Which signals ?

[Non-invasive technologies]

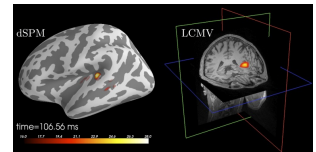
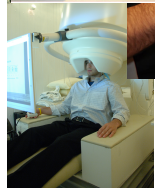
- EEG
- MEG
- fMRI

○ Real-life issues ?

- Medical diagnose
- Brain understanding
 - Source localisation
 - Brain reading

○ Machine-Learning issues ?

- Classification
- Regression
- + Transfer
- + Specific framework : 0-shot learning



Raw data

- ⇒ Spatio-(temporal) data, sensor networks
- ⇒ Personalized signal

Non-invasive technologies

- fMRI
- MEG



Vincent Guigue

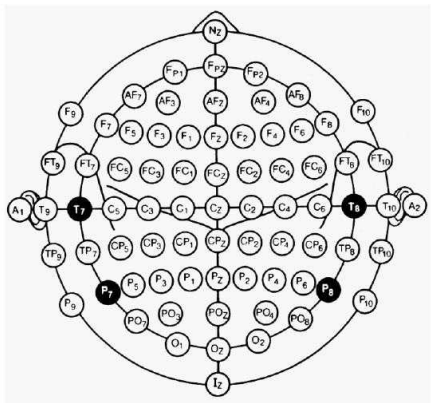
high noise level

Issues & machine learning approaches

General problem		ML techniques	Specific settings
Signal classification	P300 BCI	Signal (pre-)Processing Classifier (SVM, Ridge, LASSO) Riemannian Geometry	Transfer learning
	Seizure detection	Convolutional network (deep learning)	
	Brain Reading	Neural network Latent representation	Transfer learning 0-shot learning
Source localization		Regression	Inverse problem

X-EG Datasets

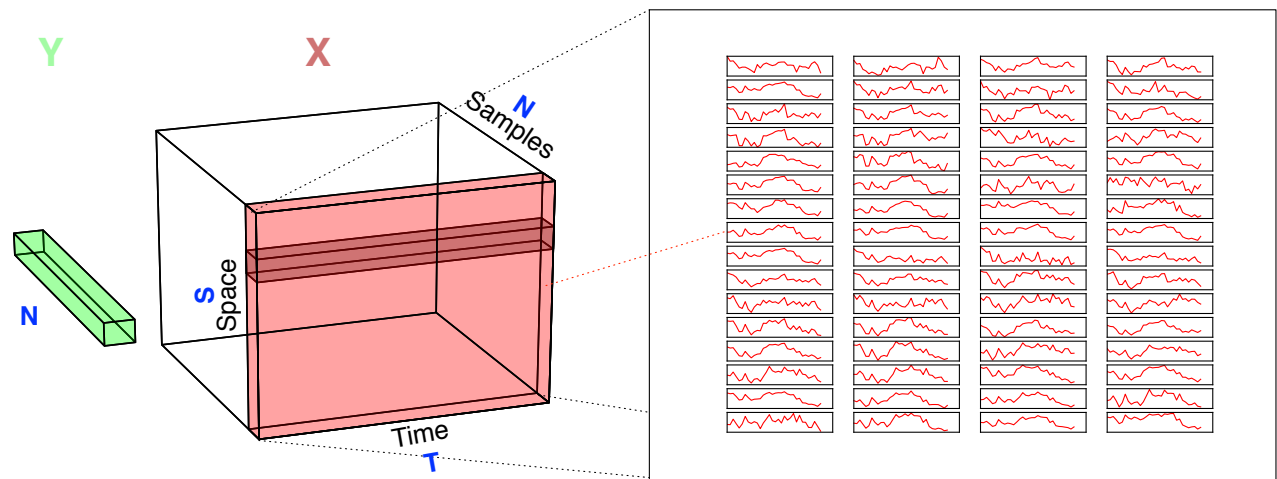
- **Spatial information** : sensors are placed according to standard patterns, e.g. :



EEG : 14 (epoc), 64 (usually), 118...
 MEG : > 300, 2 kind of sensors

- **Temporal Information** : usual sampling $30Hz < f < 1000Hz$

Dataset & notations

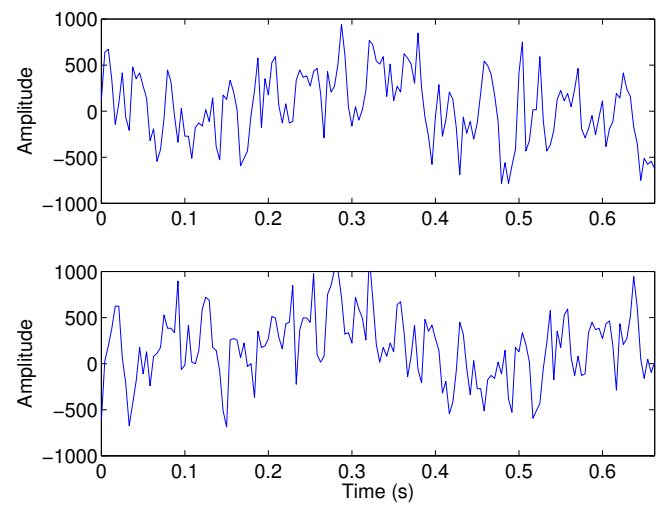


- **N** samples can be divided in **U** users
- Each user can be splitted in **Ns** sessions

Difficulties

P300 exemples :

- Data from the BCI Competition 2003 provided by the Wadsworth Institute
- EEG acquisition : 64 Channels scalp sampled at 240 Hz
- Single user and 3 acquisition sessions spelling (5,6 and 8 words)



Positive & negative samples

Unbalanced dataset & (very) **high noise** level !

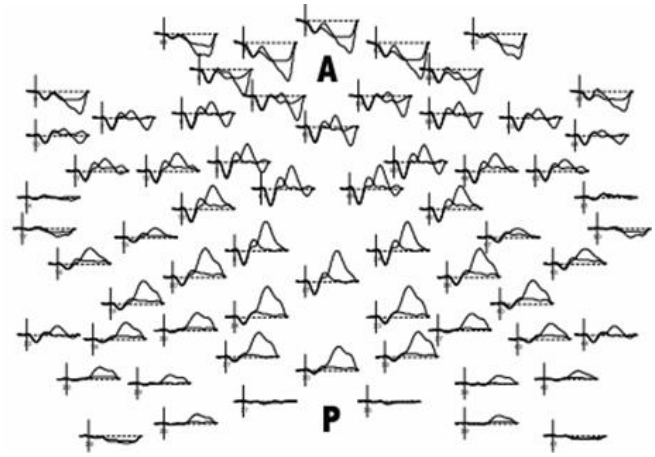
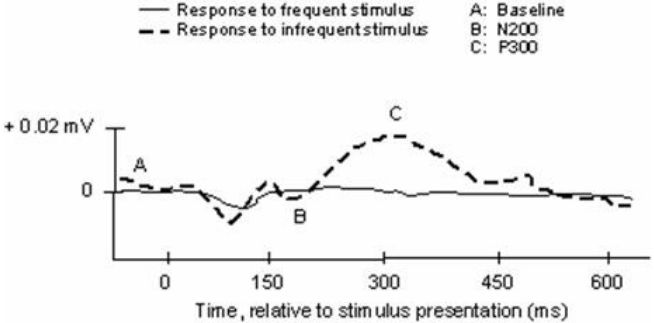
⇒ ML techniques are not able to tackle efficiently raw data (yet)

Sample aggregation

Credit : Patel and Azzam, 2005

Event Related Potential (mean of signals) : Over the scalp

Over one channel



- A powerful tool to understand...
- ... Harder to classify single sample.

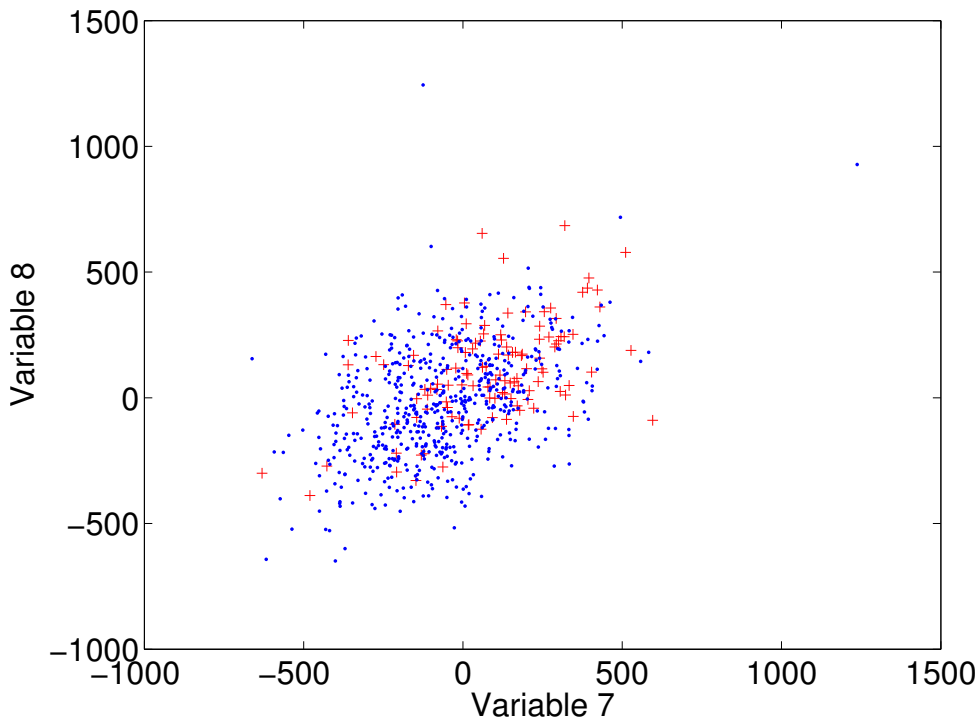


Patel and Azzam, 2005

Characterization of N200 and P300 : Selected Studies of the Event-Related Potential

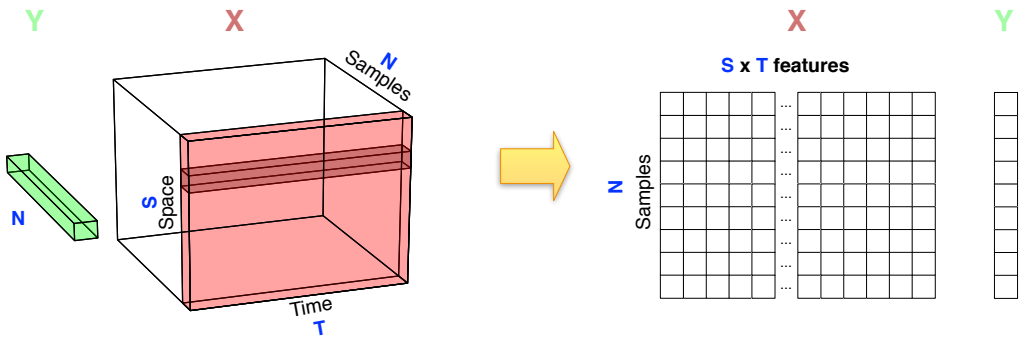
Filtering + sample aggregation

The problem remain difficult :



Plot of variable 7 vs variable 8 (≈ 300 ms)

Spatial aggregation : Concatenation

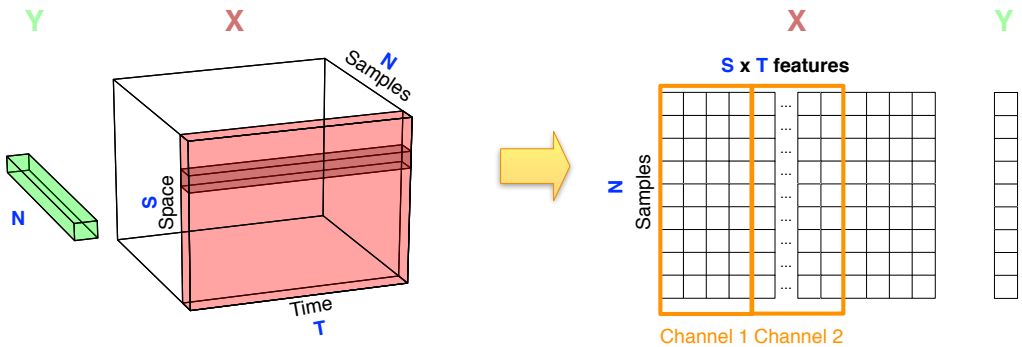


- Linear classifier :

$$f(\mathbf{x}_i) = \sum_j w_j x_{ij} \approx y_i$$

- No satisfactory performances

Spatial aggregation : Concatenation



- Linear classifier :

$$f(\mathbf{x}_i) = \sum_j w_j x_{ij} \approx y_i$$

- No satisfactory performances
- \Rightarrow (bloc) feature selection : finding which channel are important... Or not. = eliminating bloc of \mathbf{w}

Recursive Channel Elimination

A simple (& costly) approach :

```

Initialization : RANKED =  $\emptyset$  ; CHANNEL =  $[1, \dots, d]$  ;
while CHANNEL is not empty do
  for  $i$  in CHANNEL do
    Remove temporarily channel  $i$  in CHANNEL;
    Learn a linear SVM with the remaining channel;
    Compute ranking criterion  $Crit^{-i}$ ;
  end
  RANKCHAN =  $\arg \min_i Crit^{-i}$  ;
  RANKED = [ RANKCHAN RANKED ] ;
  CHANNEL = CHANNEL / RANKCHAN ;
end

```

Algorithm 1: Variable ranking with backwards algorithm

Recursive Channel Elimination

A simple (& costly) approach :

```

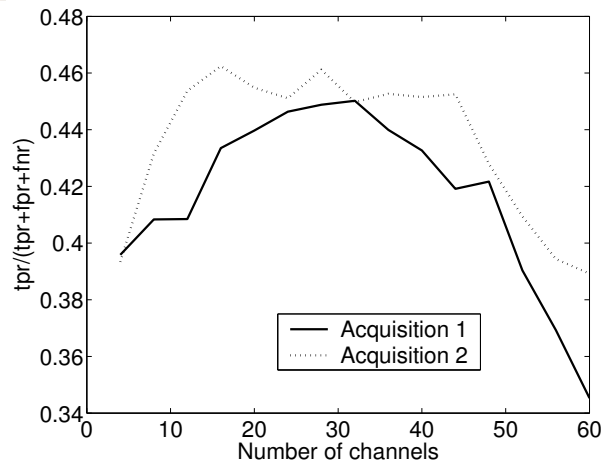
Initialization : RANKED =  $\emptyset$  ; CHANNEL =  $[1, \dots, d]$  ;
while CHANNEL is not empty do
  for  $i$  in CHANNEL do
    Remove temporarily channel  $i$  in CHANNEL;
    Learn a linear SVM with the remaining channel;
    Compute ranking criterion  $Crit^{-(i)}$ ;
  end
  RANKCHAN =  $\arg \min_i Crit^{-(i)}$  ;
  RANKED = [ RANKCHAN RANKED ] ;
  CHANNEL = CHANNEL / RANKCHAN ;
end

```

Algorithm 2: Variable ranking with backwards algorithm

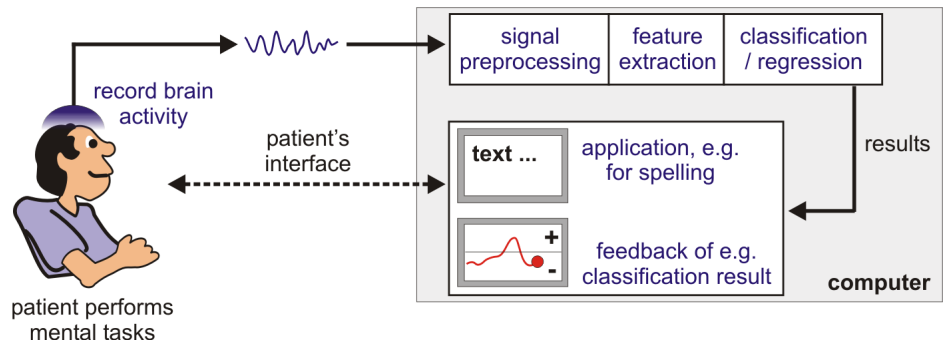
Feature Selection Results

- learning with 2 different sets lead to very different results
- best number of channels varies between 10 and 30
- performance varies between 0.35 and 0.46



Sessions	10 Top Ranked Channels									
1	9	15	18	36	40	55	56	59	63	64
2	18	39	53	55	56	58	59	60	61	64
3	9	18	40	48	53	55	56	58	61	64
4	10	18	33	42	46	55	56	58	60	64
5	16	22	39	40	50	56	57	60	61	62

Processing chain



Crédit : M. Tangermann

Step 2 : which classifier ?

Several alternatives (even for linear classifier)

Classical(& robust) linear classifier :

$$f(\mathbf{x}_i) = \sum_j w_j x_{ij} \approx y_i$$

- Logistic Regression (max likelihood)

$$\mathbf{w}^* = \arg \max_{\mathbf{w}} \prod_i P(s_{\mathbf{w}}(\mathbf{x}_i) = 1 | \mathbf{x}_i)^{y_i} \times [1 - P(s_{\mathbf{w}}(\mathbf{x}_i) = 1 | \mathbf{x}_i)]^{1 - y_i}$$

- SVM (L1 cost, L2 regularization)
- LASSO (L2 cost, L1 regularization)
- Ridge regression (L2 cost, L2 regularization)

$$\mathbf{w}^* = \arg \max_{\mathbf{w}} \sum_i \Delta(f_{\mathbf{w}}(\mathbf{x}_i), y_i) + \lambda \Omega(\mathbf{w})$$

No impact in our chain... But many opportunities in other contexts.

Merging classifiers

Each classifier is trained on a word (sessions contain resp. 5, 6 and 8 words).

How to recognize a character from the 15 sequences ?

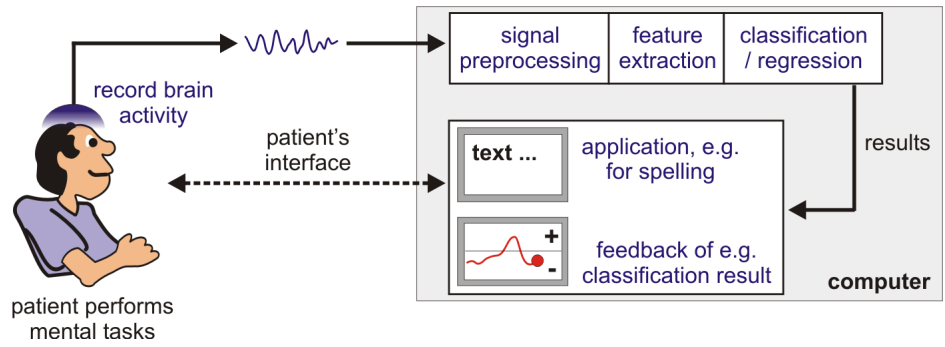
- Let x_i be post-stimulus signal associated to the illumination of a row or a column
- Each classifier scores bx_i through $f_k(bx_i)$
- Update the overall score of the given row/column

$$S_{rc} = S_{rc} + \sum_k f_k(bx_i)$$

- After all the sequences, select the character which corresponds to the highest row and columns scores.

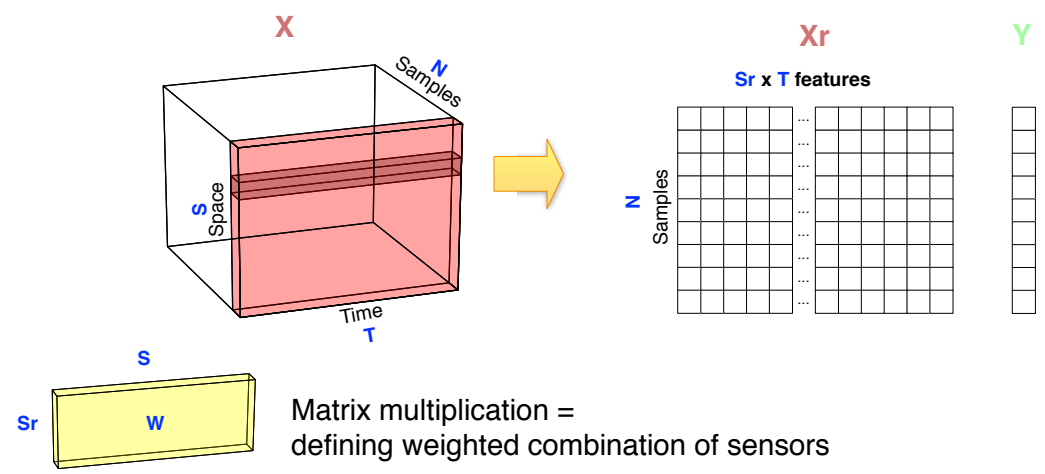
Algorithms	Nb. of sequences							
	1	2	3	4	5	6	7	10
10 preselected channels and single SVM	14	6	6	0	1	0	0	0
all channels and single SVM	14	10	9	5	5	5	1	0
10 preselected channels and Ens. SVM	13	8	3	1	2	0	0	0
all channels and Ens. SVM	7	4	3	0	0	0	0	0
4 relevant channels and Ens. SVM	8	7	4	0	1	0	0	0
10 relevant channels and Ens. SVM	8	5	5	1	0	1	0	0
26 relevant channels and Ens. SVM	4	2	0	0	0	0	0	0
30 relevant channels and Ens. SVM	5	3	0	0	0	0	0	0
optimal relevant channels and Ens. SVM	4	2	1	0	0	0	0	0

TABLE: Errors wrt the nb of illumination sequences



Crédit : M. Tangermann

Which alternatives ?
Can we merge pre-processing & training steps ?
 (≈) New issues in ML techniques for EEG analysis



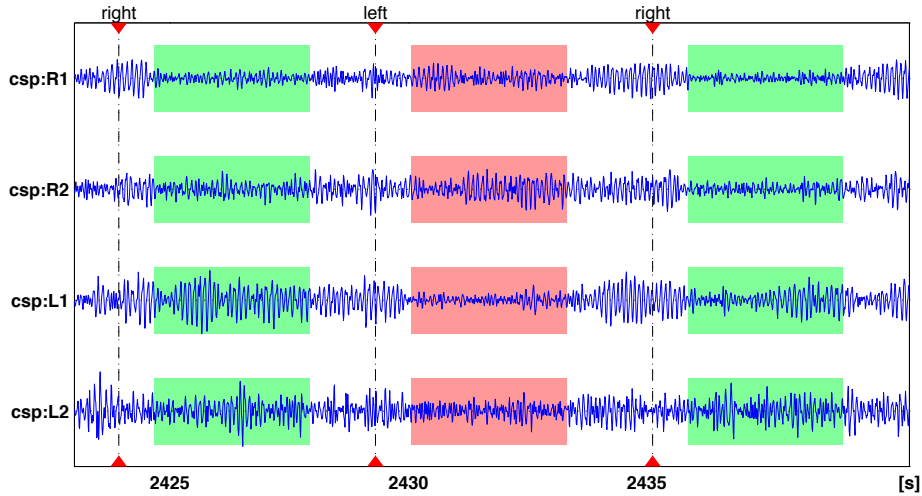
- Orthogonal sensor combinations maximizing the variance (\approx PCA in sensor space)
- Combining sensor = noise reduction

 ZJ Koles, MS Lazar, SZ Zhou, 1990

Spatial patterns underlying population differences in the background EEG

CSP : one of the key of Robust EEG Single-Trial Analysis

Exemples of use in motor cortex imagery



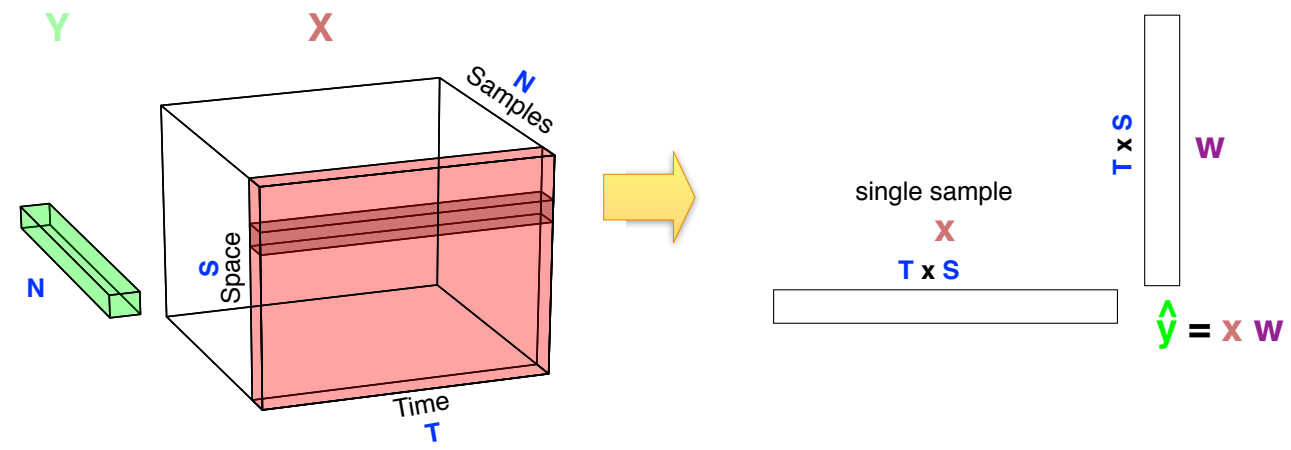
Left vs right hand move mapped to 4 aggregated channels.

[Blankertz et al., 2008](#)

Optimizing Spatial Filters for Robust EEG Single-Trial Analysis

Bi-linear SVM

- Using a linear classifier = losing structure information

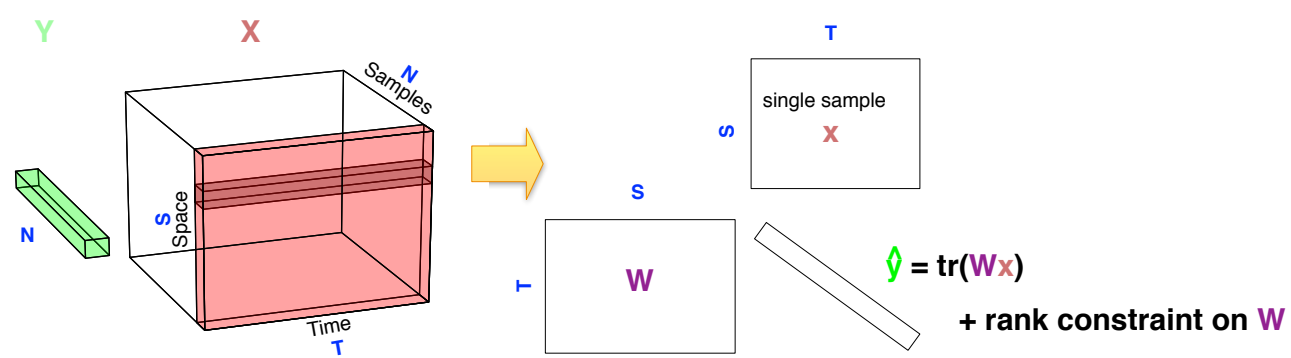



How can we impose structural constraints on w ??

Pirsiavash et al., NIPS 2009
 Bilinear classifiers for visual recognition

Bi-linear SVM

- Using a linear classifier = losing structure information
- bilinear classifiers \Rightarrow Modeling variable dependencies on 2 axis (time/space)

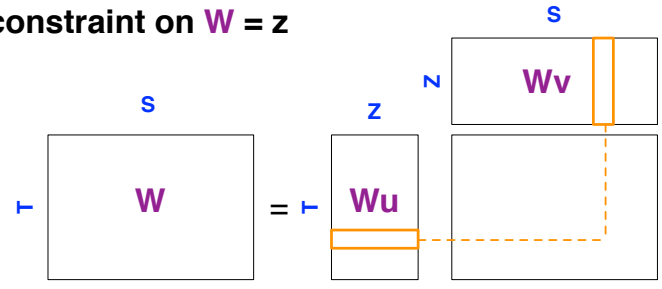



 Pirsiavash et al., NIPS 2009
 Bilinear classifiers for visual recognition

Bi-linear SVM

- Using a linear classifier = losing structure information
- bilinear classifiers \Rightarrow Modeling variable dependencies on 2 axis (time/space)

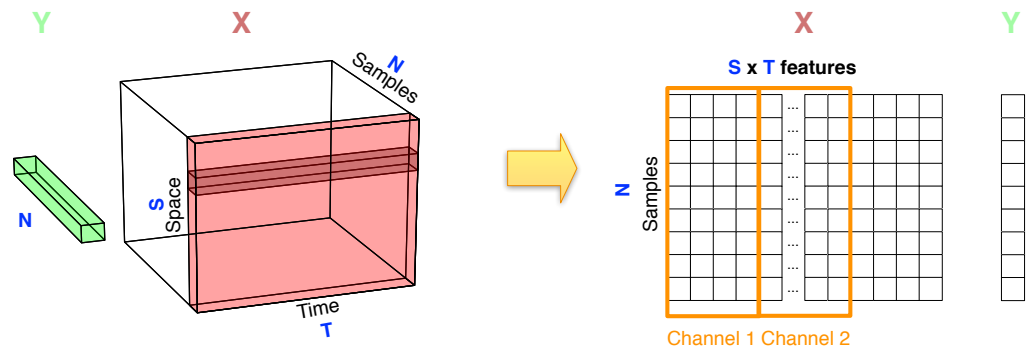
Rank constraint on $W = z$



\Rightarrow Structural consistency in the way of building W

Pirsiavash et al., NIPS 2009
 Bilinear classifiers for visual recognition

Regularization as a selection procedure with linear classifiers



$$f(\mathbf{x}_i) = \sum_j w_j x_{ij} \approx y_i$$

General training formulation :

$$\mathbf{W}^* = \arg \max_{\mathbf{W}} \sum_i \Delta(f_{\mathbf{W}}(\mathbf{x}_i), y_i) + \lambda \Omega(\mathbf{W}), \quad \mathbf{W}^* \in \mathbb{R}^{S \times T}$$

Regularization as a selection procedure with linear classifiers (2)

$$\mathbf{W}^* = \arg \max_{\mathbf{W}} \sum_i \Delta(f_{\mathbf{W}}(\mathbf{x}_i), y_i) + \lambda \Omega(\mathbf{W}), \quad \mathbf{W}^* \in \mathbb{R}^{S \times T}$$

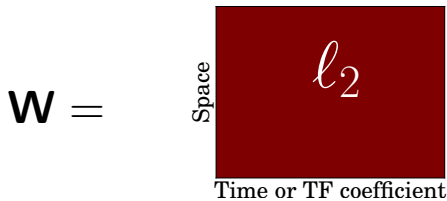
Regularization as a selection procedure with linear classifiers (2)

$$\mathbf{W}^* = \arg \max_{\mathbf{W}} \sum_i \Delta(f_{\mathbf{W}}(\mathbf{x}_i), y_i) + \lambda \Omega(\mathbf{W}), \quad \mathbf{W}^* \in \mathbb{R}^{S \times T}$$

- **L2 regularization** : $\Omega(\mathbf{W}) = \sum_{j,k} w_{jk}^2$

Associated update in a gradient descent procedure :

$$w_{jk} \leftarrow w_{jk} - 2\epsilon w_{jk} \Leftrightarrow w_{jk} \leftarrow w_{jk}(1 - 2\epsilon)$$



[credit Gramfort]

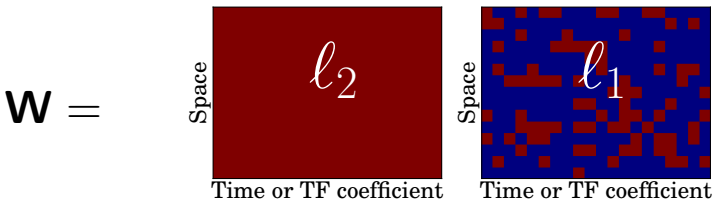
Regularization as a selection procedure with linear classifiers (2)

$$\mathbf{W}^* = \arg \max_{\mathbf{W}} \sum_i \Delta(f_{\mathbf{W}}(\mathbf{x}_i), y_i) + \lambda \Omega(\mathbf{W}), \quad \mathbf{W}^* \in \mathbb{R}^{S \times T}$$

- **L1 regularization** : $\Omega(\mathbf{W}) = \sum_{j,k} |w_{jk}|$

Associated update in a gradient descent procedure = **soft-thresholding** :

$$w_{jk} \leftarrow \begin{cases} w_{jk} - \epsilon \text{ sign}(w_{jk}) & \text{if } |w_{jk}| > \epsilon \\ 0 & \text{else} \end{cases}$$



[credit Gramfort]

Regularization as a selection procedure with linear classifiers (2)

$$\mathbf{W}^* = \arg \max_{\mathbf{W}} \sum_i \Delta(f_{\mathbf{W}}(\mathbf{x}_i), y_i) + \lambda \Omega(\mathbf{W}), \quad \mathbf{W}^* \in \mathbb{R}^{S \times T}$$

Elastic net variant combines L1 and L2 for more stability

- Sparseness of L1,
- Robustness of L2



Zou, Hastie, 2005

Regularization and variable selection via the elastic net

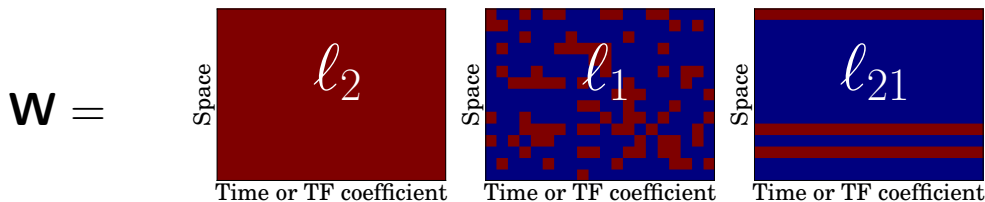
Regularization as a selection procedure with linear classifiers (2)

$$\mathbf{W}^* = \arg \max_{\mathbf{W}} \sum_i \Delta(f_{\mathbf{W}}(\mathbf{x}_i), y_i) + \lambda \Omega(\mathbf{W}), \quad \mathbf{W}^* \in \mathbb{R}^{S \times T}$$

- **L21 regularization** : $\Omega(\mathbf{W}) = \sum_j \sqrt{\sum_k w_{jk}^2} = \sum_j \|\mathbf{w}_j\|$

Sparsity at the sensor level Gradient descent update :

$$w_{jk} \leftarrow \begin{cases} w_{jk} \left(1 - \frac{\epsilon}{\|\mathbf{w}_j\|}\right) & \text{if } \|\mathbf{w}_j\| > \epsilon \\ 0 & \text{else} \end{cases}$$



[credit Gramfort]



G. Obozinski, B. Taskar, and M. I. Jordan, 2006

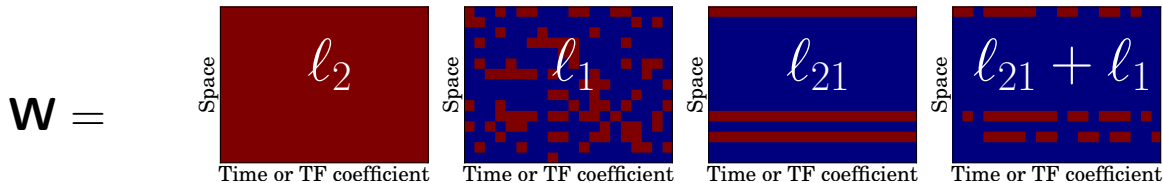
Multi-task feature selection

Regularization as a selection procedure with linear classifiers (2)

$$\mathbf{W}^* = \arg \max_{\mathbf{W}} \sum_i \Delta(f_{\mathbf{W}}(\mathbf{x}_i), y_i) + \lambda \Omega(\mathbf{W}), \quad \mathbf{W}^* \in \mathbb{R}^{S \times T}$$

- o **L21 regularization + L1 :**

Playing with advanced (and dedicated formulation)



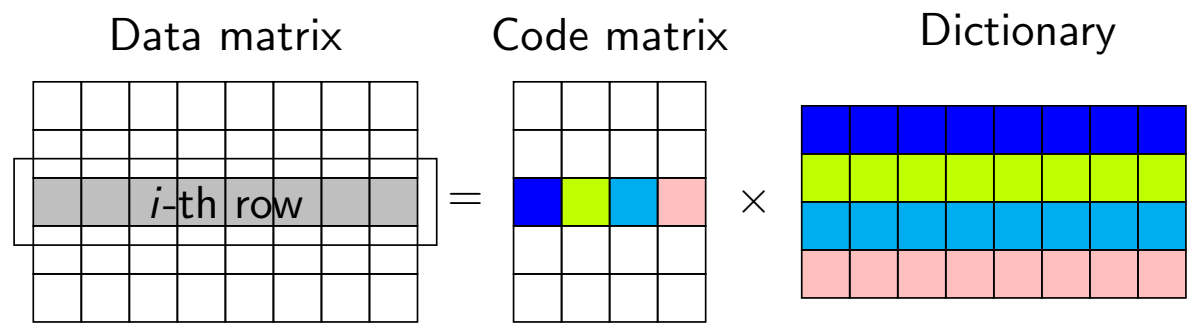
[credit Gramfort]



Gramfort et al., 2013

Time-Frequency Mixed-Norm Estimates : Sparse M/EEG imaging with non-stationary source activations

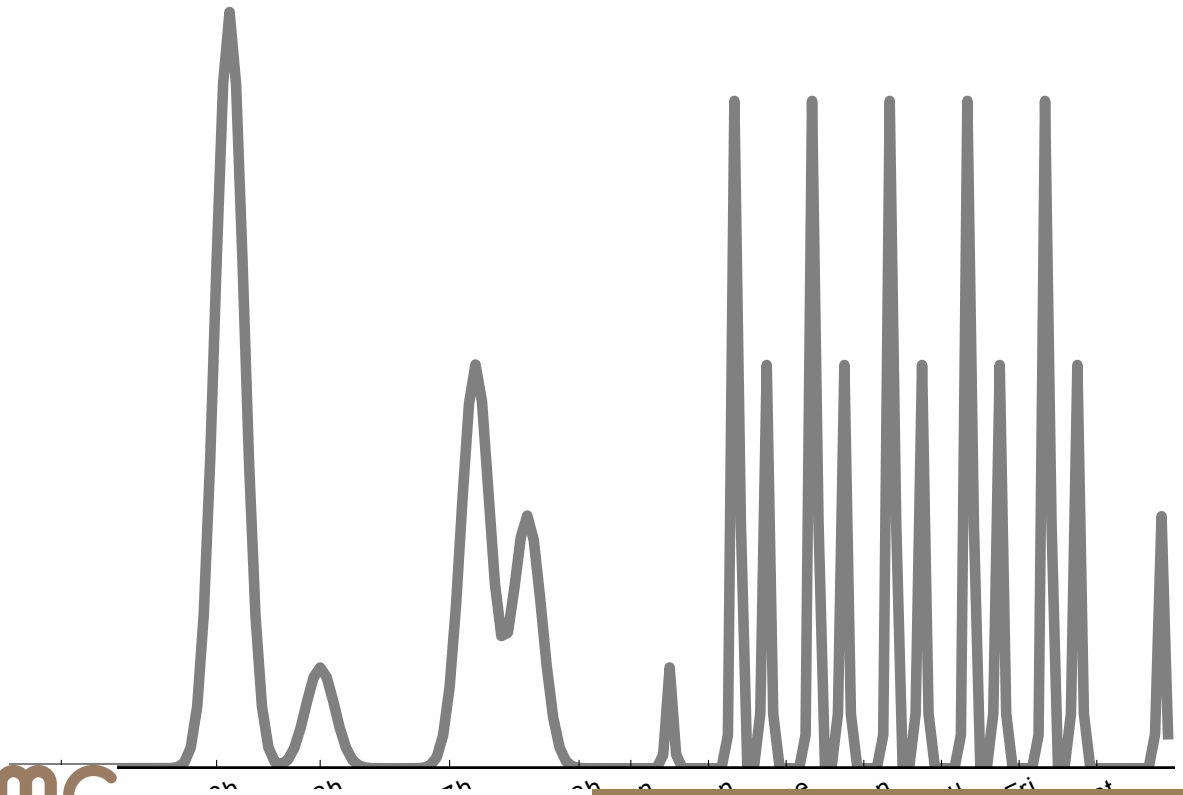
Raw signal are very difficult to handle...
 ... Let learn a new space where the problem is easy to solve !



- Variations SVD [Golub 96].
 - Non-negative matrix factorization [Lee 2000]
 - Sparseness [Hoyer 2002]
- Learning criterion = reconstruction error
- Easy constraint design (to adapt to specific problems)
- Efficient solvers

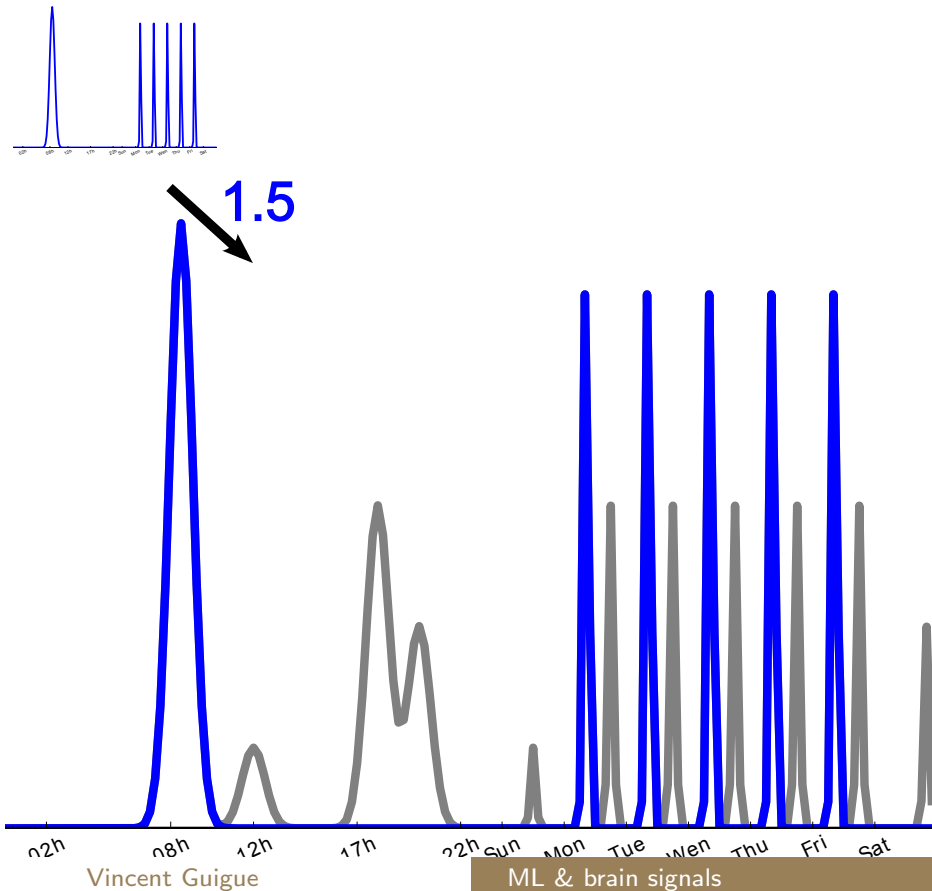
Representation learning / dictionary learning

With an exemple (far away from EEG...)



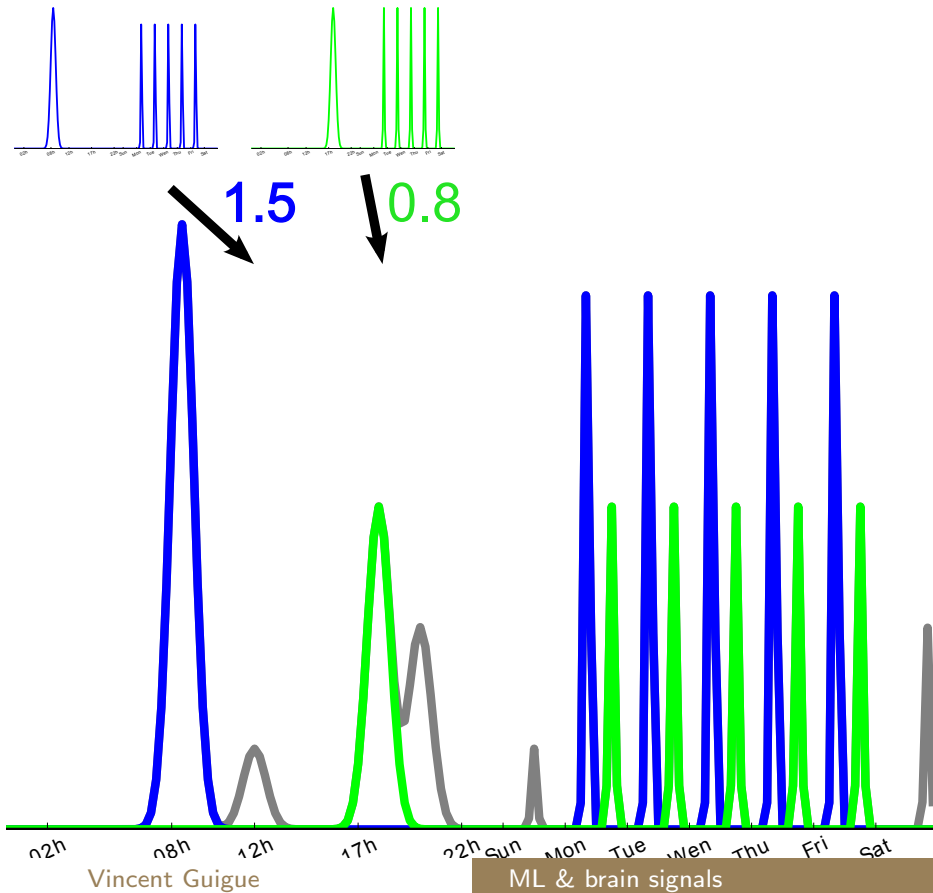
Representation learning / dictionary learning

With an exemple (far away from EEG...)

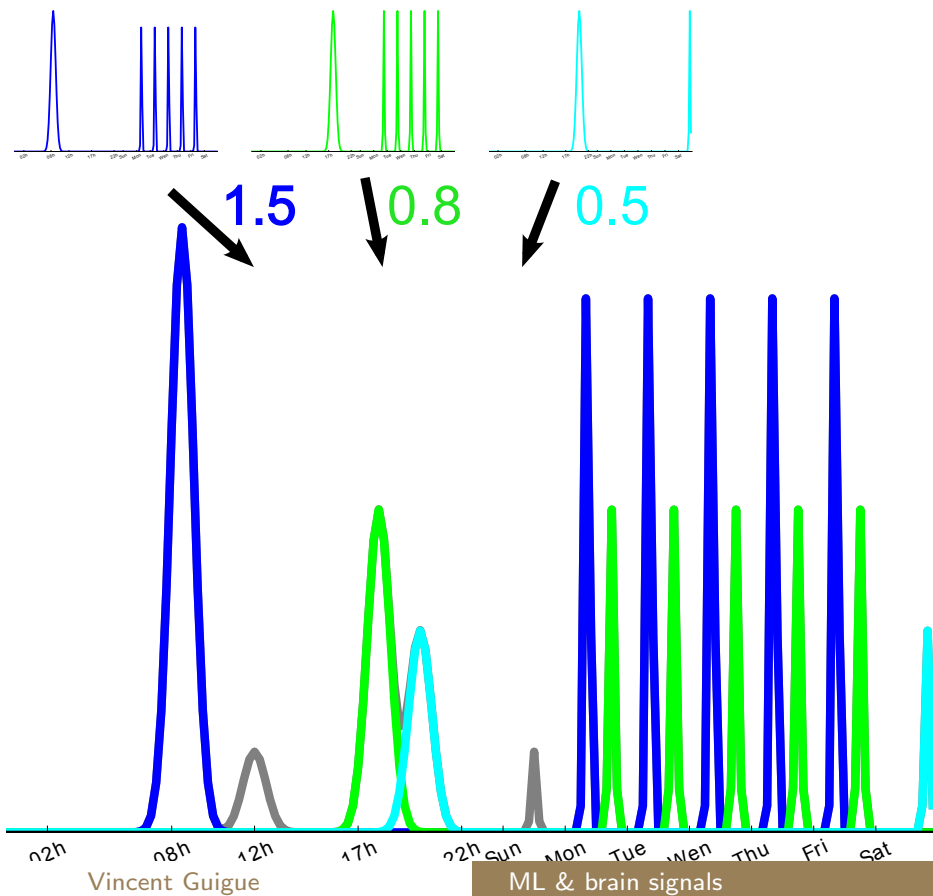


Representation learning / dictionary learning

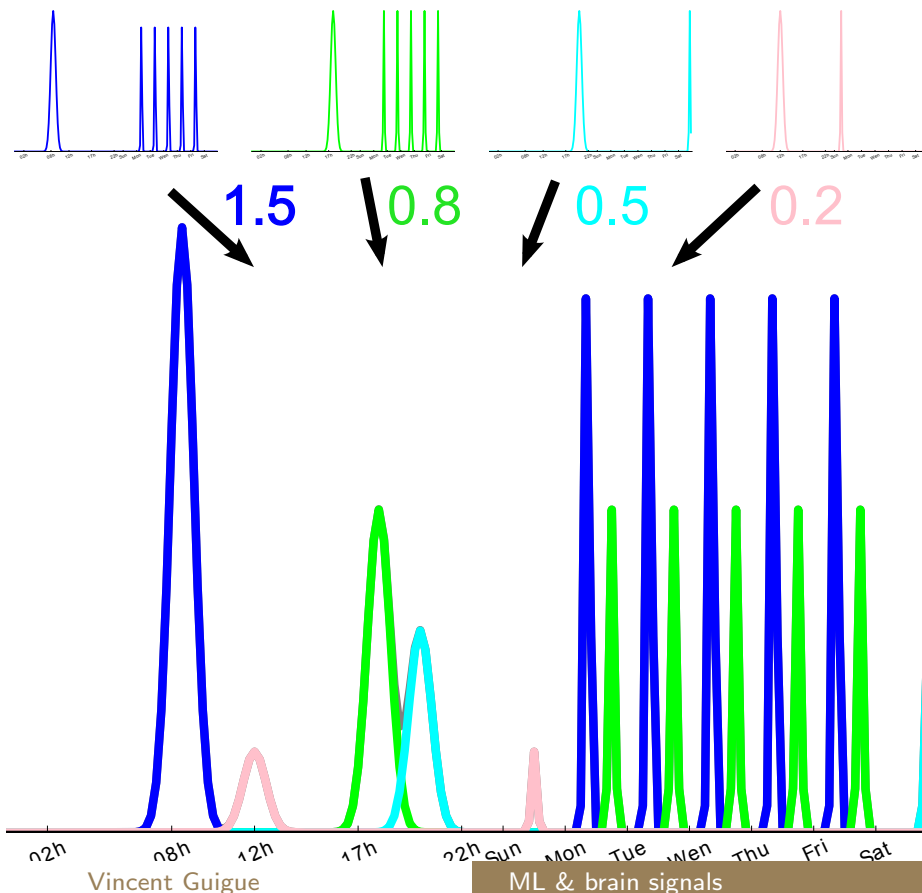
With an exemple (far away from EEG...)



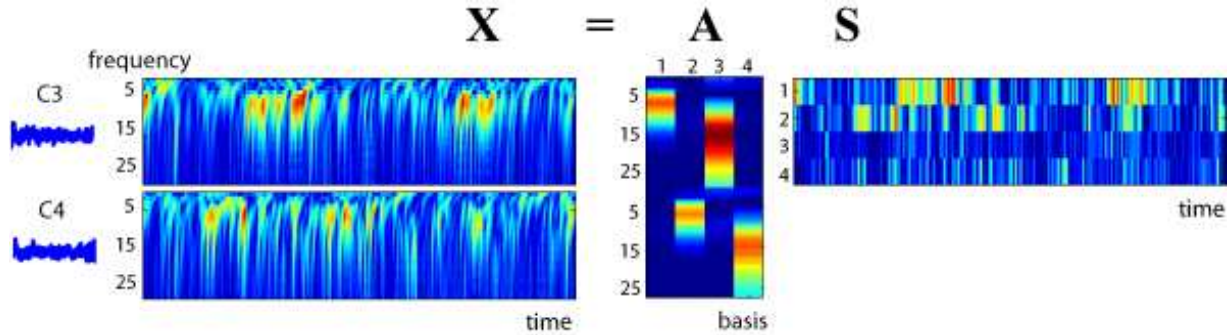
With an exemple (far away from EEG...)



With an exemple (far away from EEG...)



Extracting common pattern in time-frequency representation of EEG :



- Adding extra-constraints
- Gain when classifying A instead of X on BCI Challenge III (motor imagery)

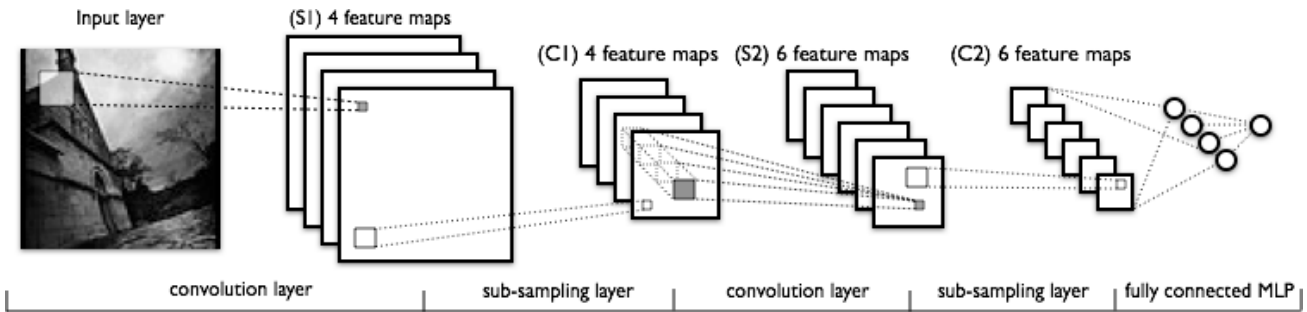

[Lee and Choi, AISTATS 2009](#)
 Group Nonnegative Matrix Factorization for EEG Classification

Deep learning & EEG

Neural networks opportunities for EEG

- Extracting auto-learned features
- Modeling invariances (both time/space)

General CNN architecture



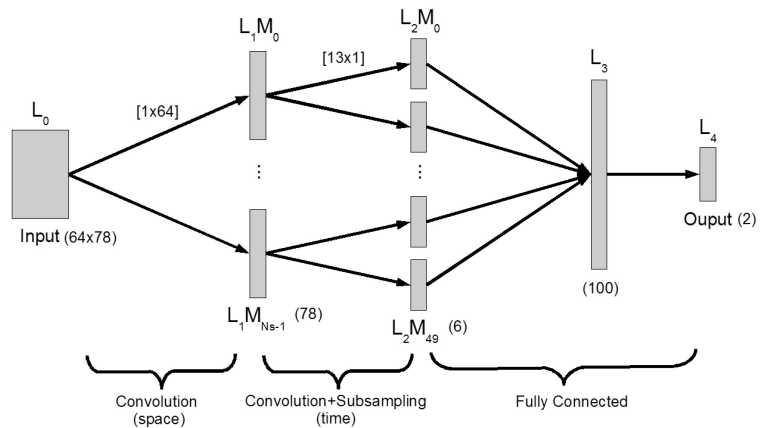
- **Very efficient** on many problems... But **not so robust** to noise
- **Easy** to understand... But **hard** to implement

Deep learning & EEG

Neural networks opportunities for EEG

- Extracting auto-learned features
- Modeling invariances (both time/space)

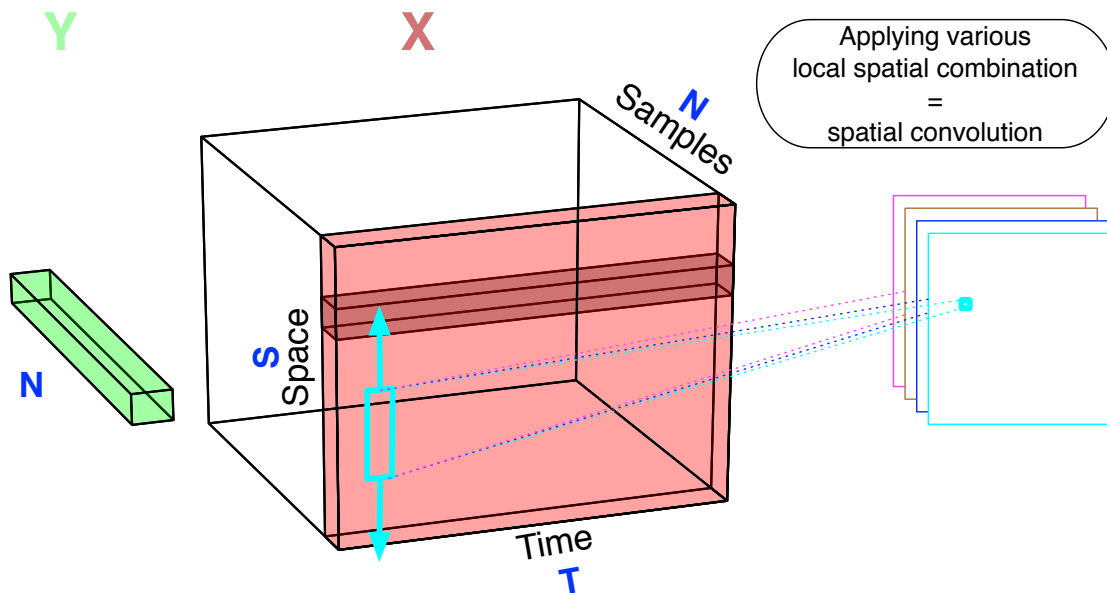
Cecotti Architecture dedicated to P300 :



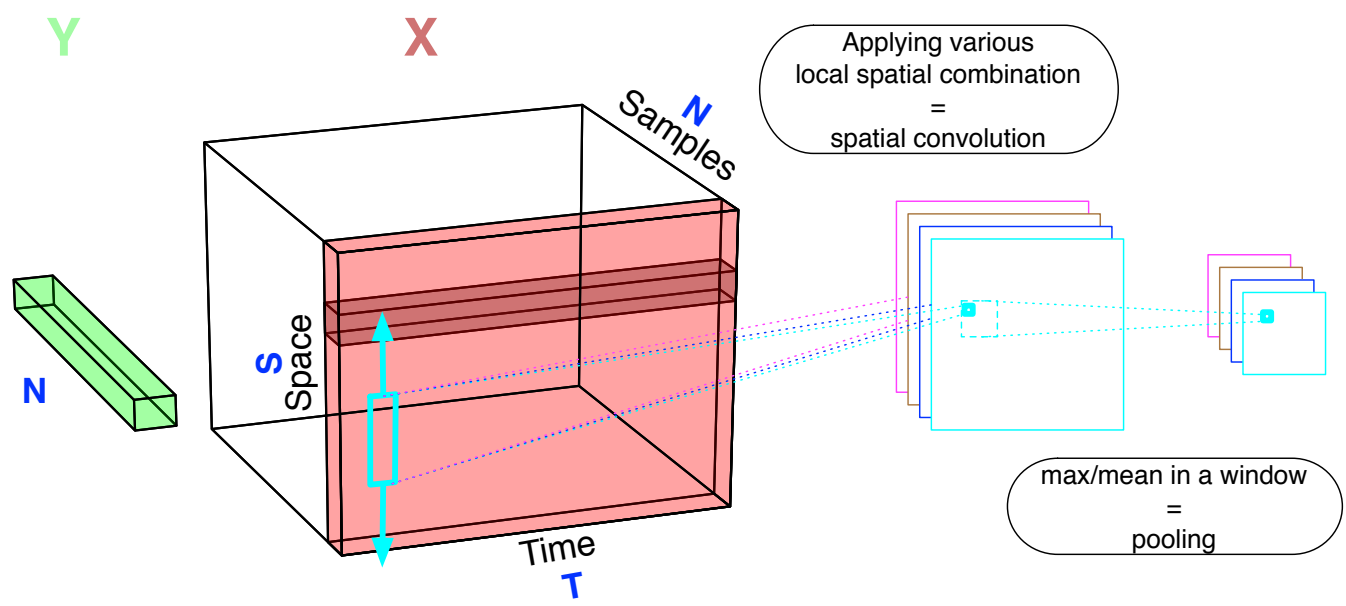
Cecotti and Gräser, PAMI 2011

Convolutional Neural Networks for P300 Detection with Application to Brain-Computer Interfaces

Detail of the architecture :

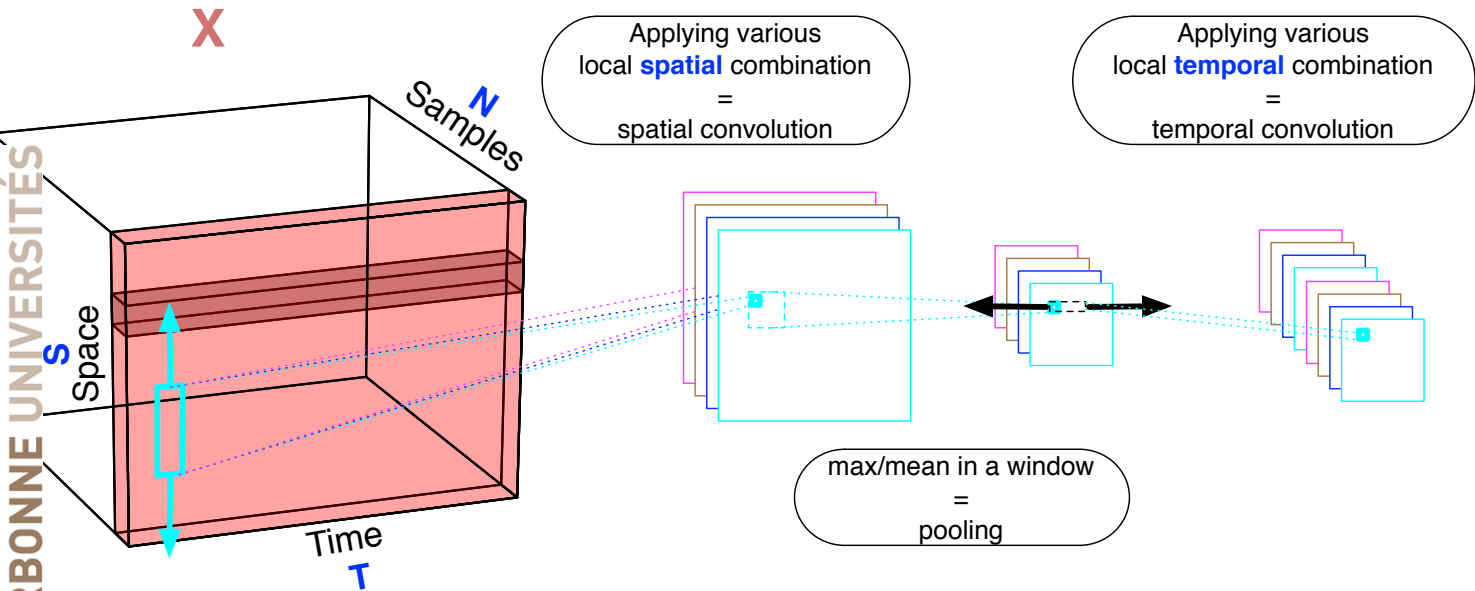


Detail of the architecture :



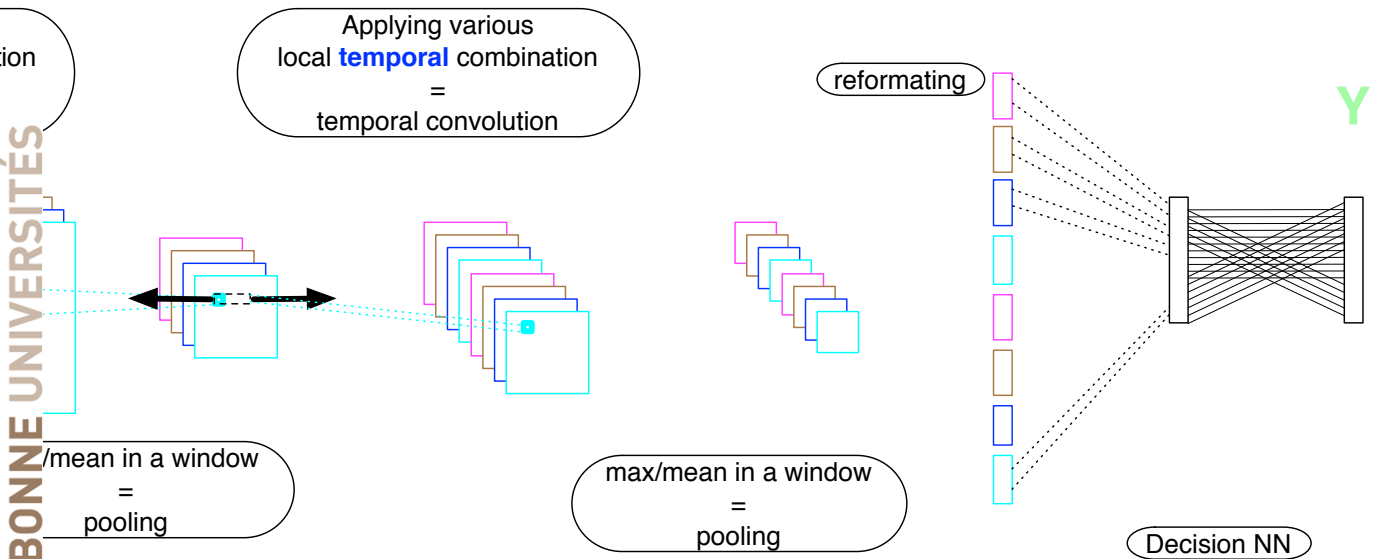
Deep learning & EEG

Detail of the architecture :

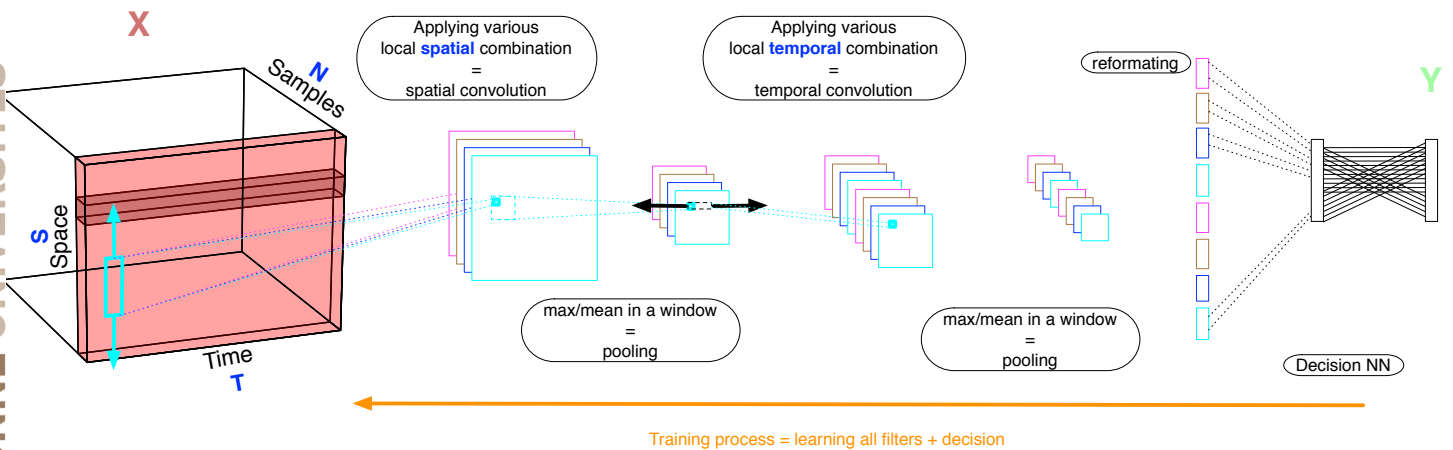


Deep learning & EEG

Detail of the architecture :



Detail of the architecture :



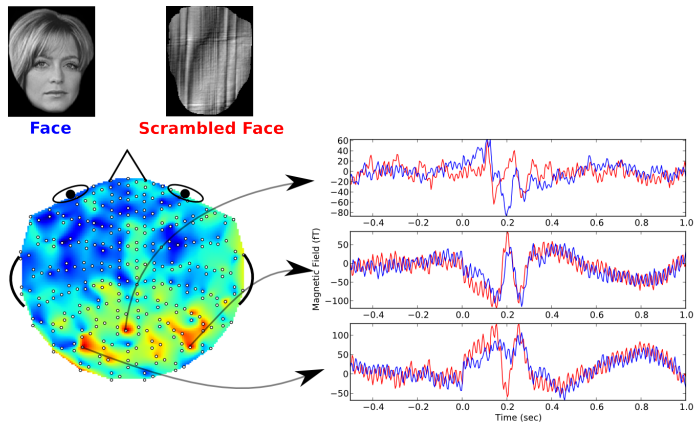
Transfer Learning

Our aim :

- ① Training models on existing EEG dataset
- ② Testing algorithms on new subjects

⇒ classical algorithms fail !

Kaggle MEG 2014 : *train* = 16 subjects ; *test* = 6 **different** subjects



Transfer Learning : which solutions ?

- Learning many classifier adapted to various topologies + aggregation/vote [easy]
- Extracting subject invariant features
- Aligning data/classifiers from one patient to another
 - Iterative Procrustean alignment + classifier in a *universal* space

Solving : $\min_T \|VT - L\|$ with :

$V \in \mathbb{R}^{k \times n}$ data to align

$L \in \mathbb{R}^{k \times n}$ well known reference

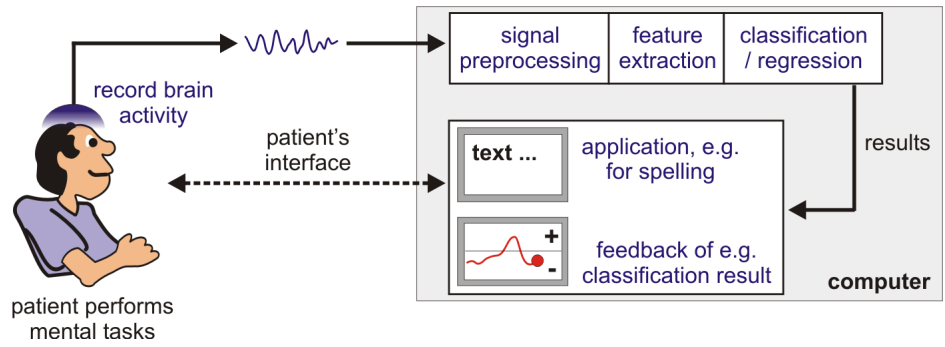
$T \in \mathbb{R}^{n \times n}$ transfer matrix



[Haxby et al., 2011](#)

A Common, High-Dimensional Model of the Representational Space in Human Ventral Temporal Cortex

Alexandre Barachant



Crédit : M. Tangermann

A real breakthrough for EEG classification...
... And transfer !
 Winner of Kaggle competition MEG 2014, EEG 2015

Transfer scheme

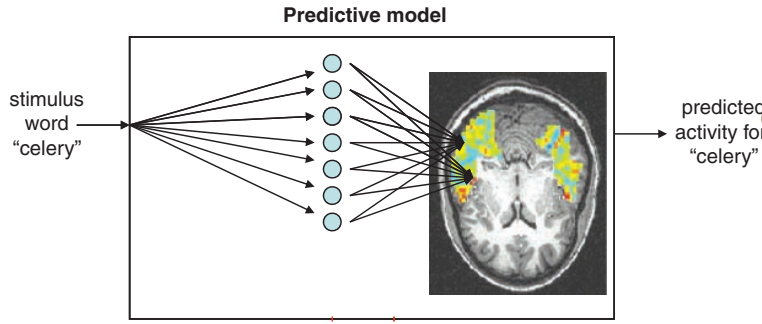
- ① Computing $P_{u,c}^*$ for all users & classes
- ② Using tangent space features :

$$x \Rightarrow \left[\phi_{P_{u,c}^*}(\Sigma) \right] \in \mathbb{R}^{S' \times S' \times U \times C}$$

- ③ LASSO linear classifier in the new space

- 1 Introduction
- 2 Signal classification for BCI applications
 - Old school processing chain
 - Opportunities in ML for EEG
 - Riemannian Geometry
- 3 Brain Reading
- 4 Source localization

Predicting brain activity



- 1 How to represent stimuli ?
Transformation ϕ
- 2 How to map ϕ to the fMRI voxel activations ?
Multi-dimensional regression :

$$\text{word } w \Rightarrow \phi(w) \in \mathbb{R}^Z$$

$$\tilde{\mathbf{y}} \in \mathbb{R}^V = \phi(w)R, \quad R \in \mathbb{R}^{Z \times V}$$

$$R^* \in \mathbb{R}^{Z \times V} =$$

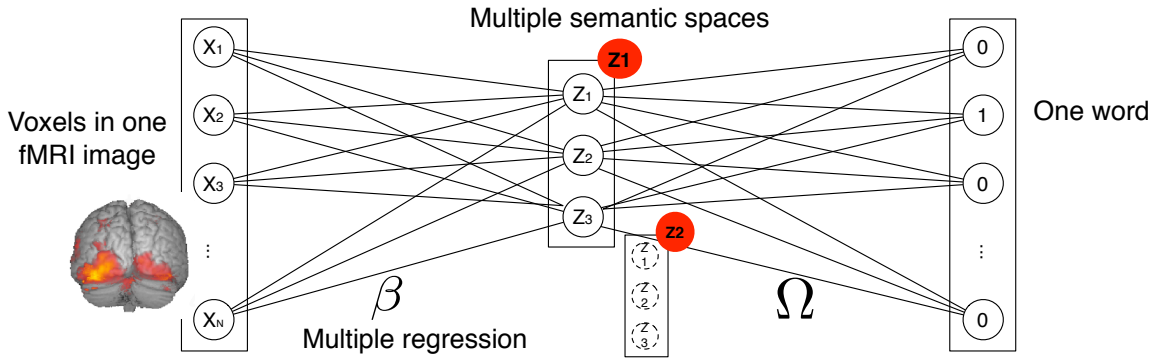
$$\arg \min_R \sum_i (\phi(w)R - \mathbf{y})^2$$

Mitchell et al., Sciences 2008

Predicting Human Brain Activity Associated with the Meanings of Nouns

0-shot learning

An original framework :
 Are we able to find a label that we didn't see in the training step ?

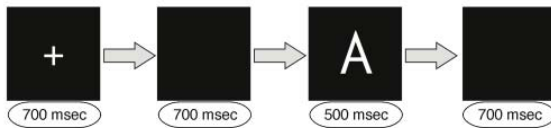


- Ω is a semantic (learned or manually designed)
- Corpus = 60 words ; 58 for training, 2 for testing
- > 80% accuracy (several semantics & bloc regularization)

 [Pipanmaekaporn et al., 2015](#)
 Designing Semantic Feature Spaces for Brain-Reading

Can we tackle brain-reading in EEG ?

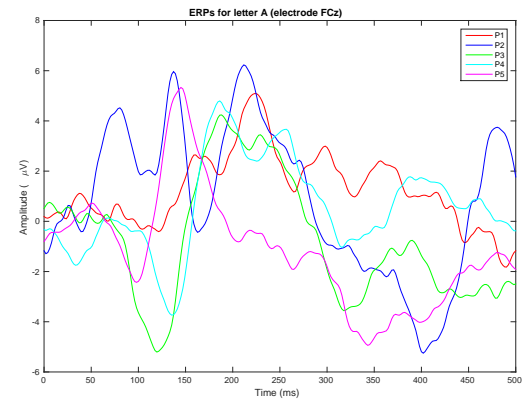
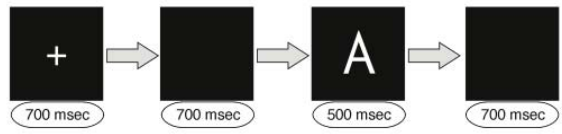
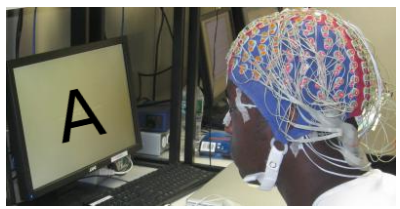
- J. Grainger (Marseille) built some datasets



Can we tackle brain-reading in EEG ?

- J. Grainger (Marseille) built some datasets

High variability :





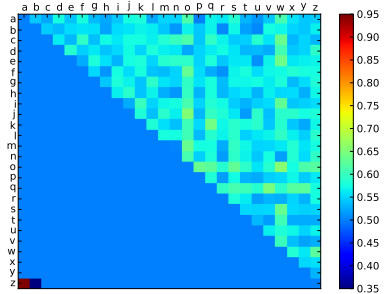
Preliminary results

- SVM \approx Ridge
- Binary classification of couples of letters
 - 325 experiments
 - Baseline (random) = 50%

Ridge Regression

56.46%

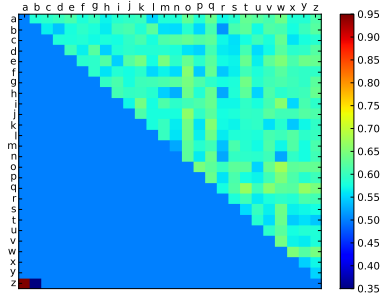
perf moy= 0.5646
min=0.5016 (a,p) max=0.6591 (o,w)



Global

59.58%

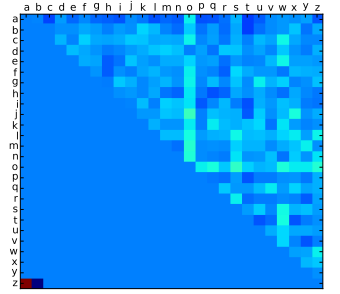
perf moy= 0.5958
min=0.5134 (b,r) max=0.6775 (o,w)



Per participant

52.11%

perf moy= 0.5211
min=0.4623 (a,t) max=0.6022 (j,o)



Transfer

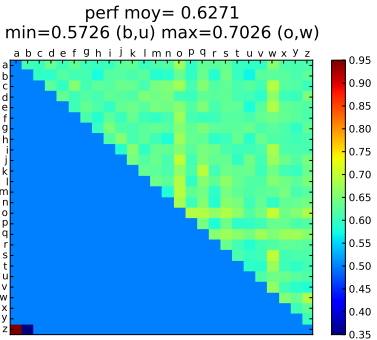


Preliminary results

- SVM \approx Ridge
- Binary classification of couples of letters
 - 325 experiments
 - Baseline (random) = 50%

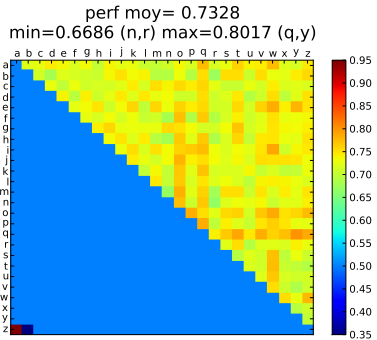
BE-C + Ridge Regression

62.71%



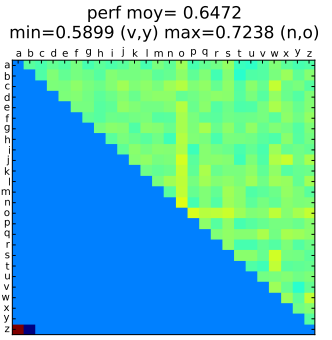
Global

73.28%



Per participant

64.72%



Transfer

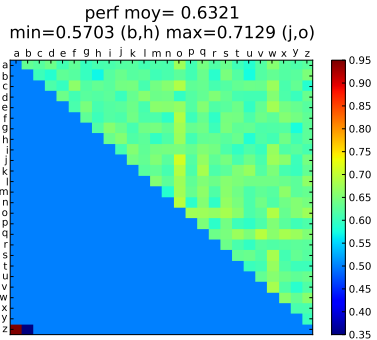


Preliminary results

- SVM \approx Ridge
- Binary classification of couples of letters
 - 325 experiments
 - Baseline (random) = 50%

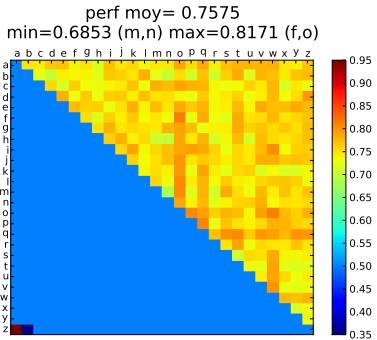
Lasso + Ridge Regression

63.21%



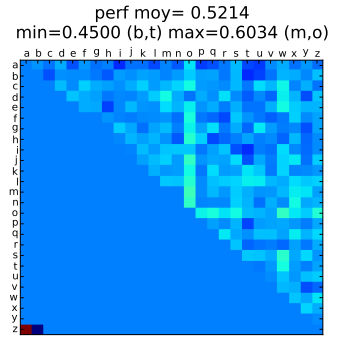
Global

75.75%



Per participant

52.14%

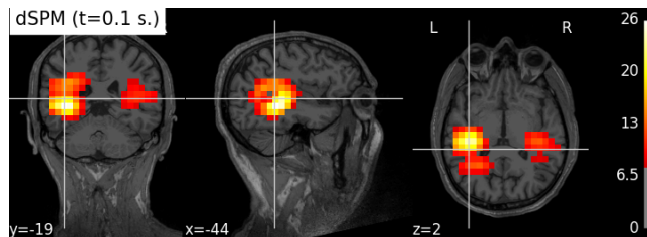


Transfer

- 1 Introduction
- 2 Signal classification for BCI applications
 - Old school processing chain
 - Opportunities in ML for EEG
 - Riemannian Geometry
- 3 Brain Reading
- 4 Source localization

Source Localization

Credit : A. Gramfort



Formulation :

$$\tilde{\mathbf{X}}^{star} = \arg \min_{\tilde{\mathbf{X}}} \|\mathbf{M} - \mathbf{G}\tilde{\mathbf{X}}\|_F$$

Major problem : noise level (!)

A. Gramfort's proposals :

- Using a time-frequency representation
- Exploiting mix-norm regularizations

Conclusion

- Many beautiful problem (from both real life & ML point of view)
- Many existing dataset (BCI, fMRI...)
- Many existing tools (sklearn, mne...)

Let's decode the brain !