



Lecture 8: Optimization-based Approaches

Introduction to Modern Brain-Computer Interface Design

Christian A. Kothe
SCCN, UCSD



Outline

1. Introduction
2. Going Beyond CSP
3. Large-Scale Machine Learning
4. Application to the Spectral Model
5. Application to ERPs
6. Learning ERP and Oscillatory Weights Simultaneously
7. Practical Remarks





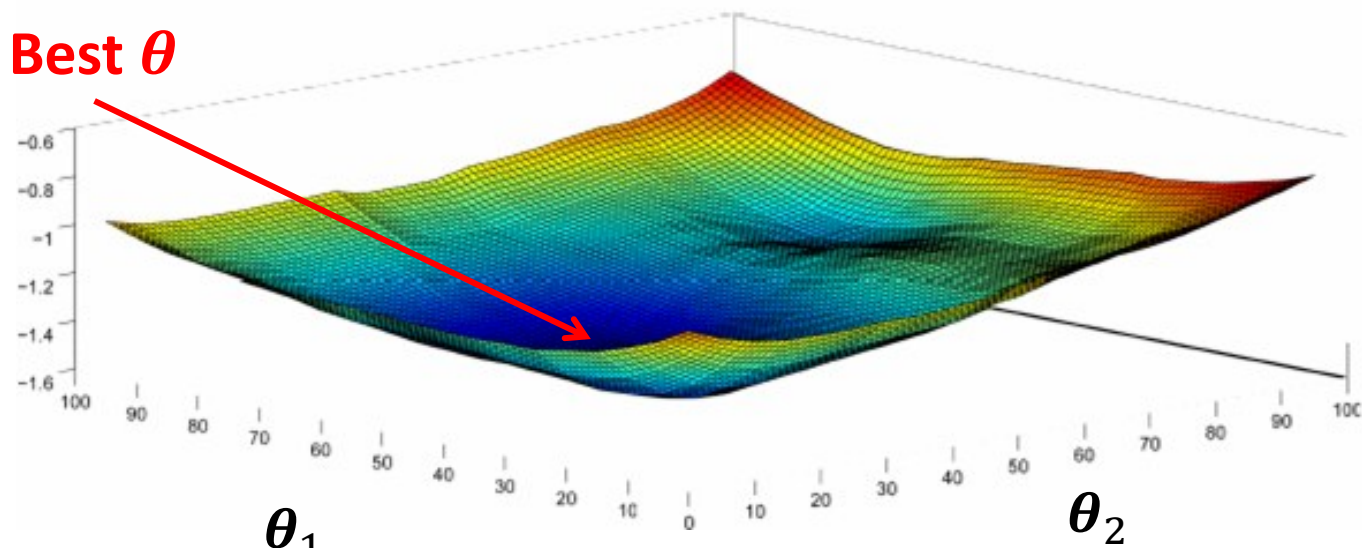
8.1 Introduction

Beyond CSP

- A unified, globally optimal solution to spatial filter estimation has recently been proposed (as an alternative to CSP+LDA)
- The method learns in a single step both the spatial filters and the relative weights for the filtered variance
- This is an *optimization*-based approach

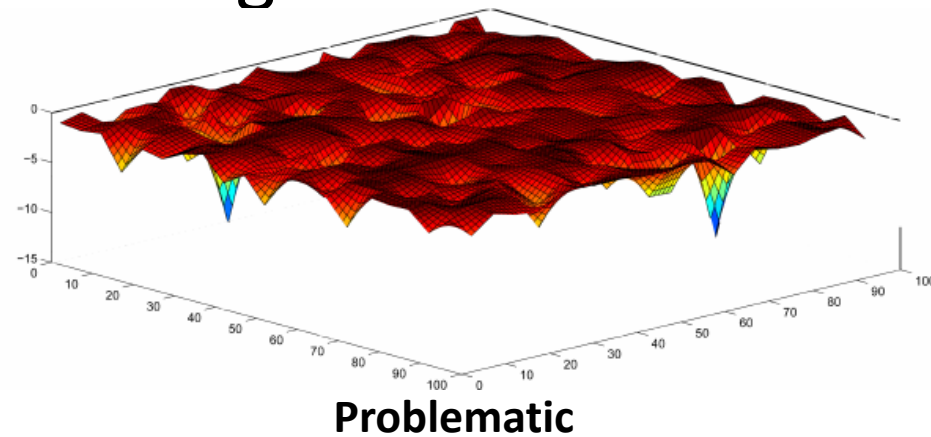
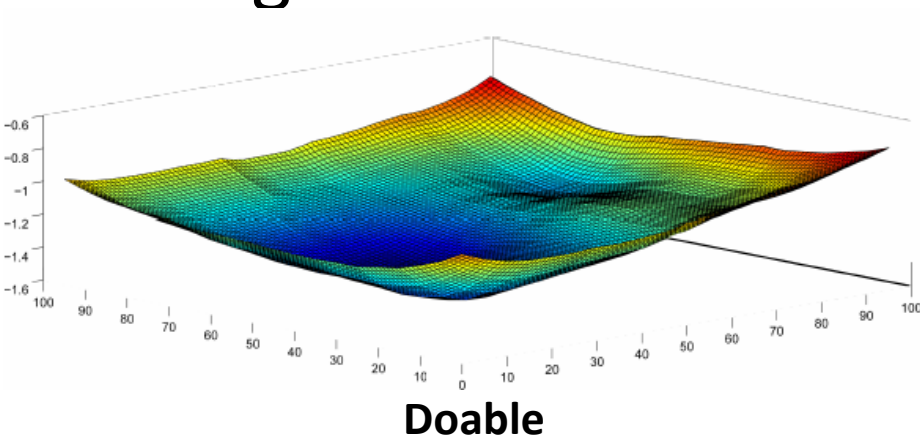
Optimization

- Broad field, concerned with finding assignments to parameters that minimize a *cost function f*
- Two branches: *Local optimization* and *Global Optimization*



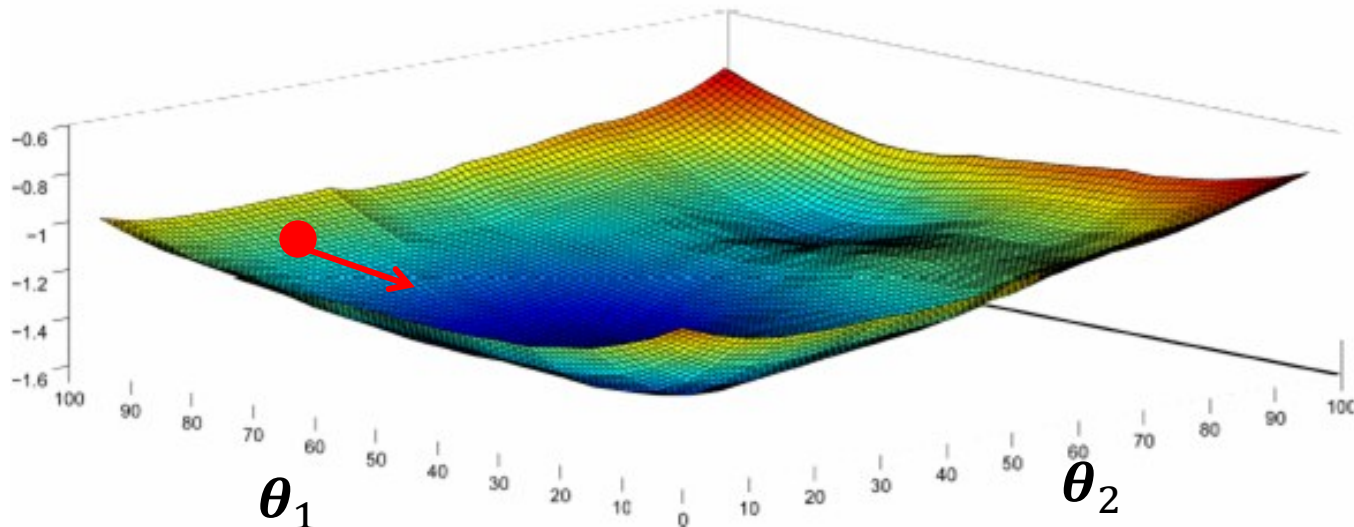
Global Optimization

- Aims to find the global optimum of a function, even if it has multiple local optima
- Can be approximate (e.g., Simulated Annealing) or exact (e.g., Branch-and-Bound)
- **Problem:** Can be *extremely slow*, especially on high-dimensional and pathological data



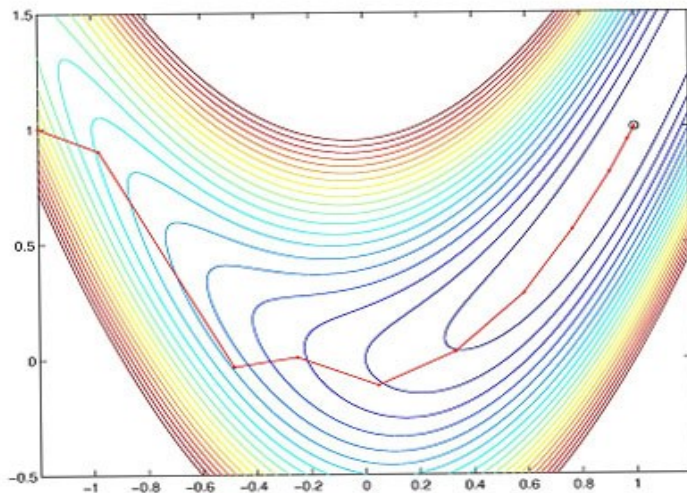
Local Optimization

- Improving an initial guess of a parameter by incremental updates
- **Method 1, Gradient Descent:** walk into the direction of steepest descent, requires 1st derivative g of cost function

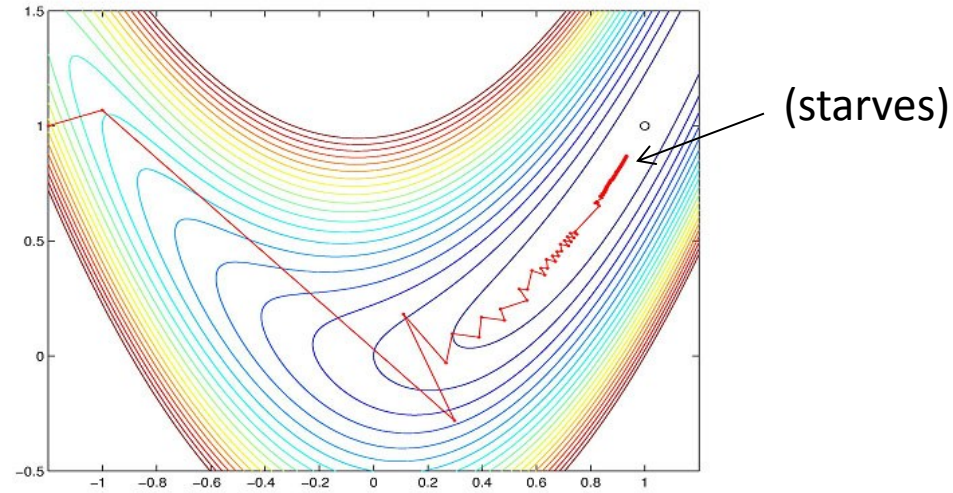


Local Optimization

- **Method 2, Newton Method:** account for local curvature, requires 2nd derivative \mathbf{B} of function (Hessian matrix)
- Can be very efficient if \mathbf{g}/\mathbf{B} is easy to compute (or easy to approximate otherwise)



Newton



Gradient Descent



Local Optimization

- **Method 3, Quasi-Newton:** approximates B^{-1} incrementally from gradients gathered along the way (no second derivatives necessary)
- **Practical variant:** L-BFGS (Limited-Memory Broyden-Fletcher-Goldfarb-Shanno) – one of the best known “off-the-shelf” second-order optimizers, *very easy* to use

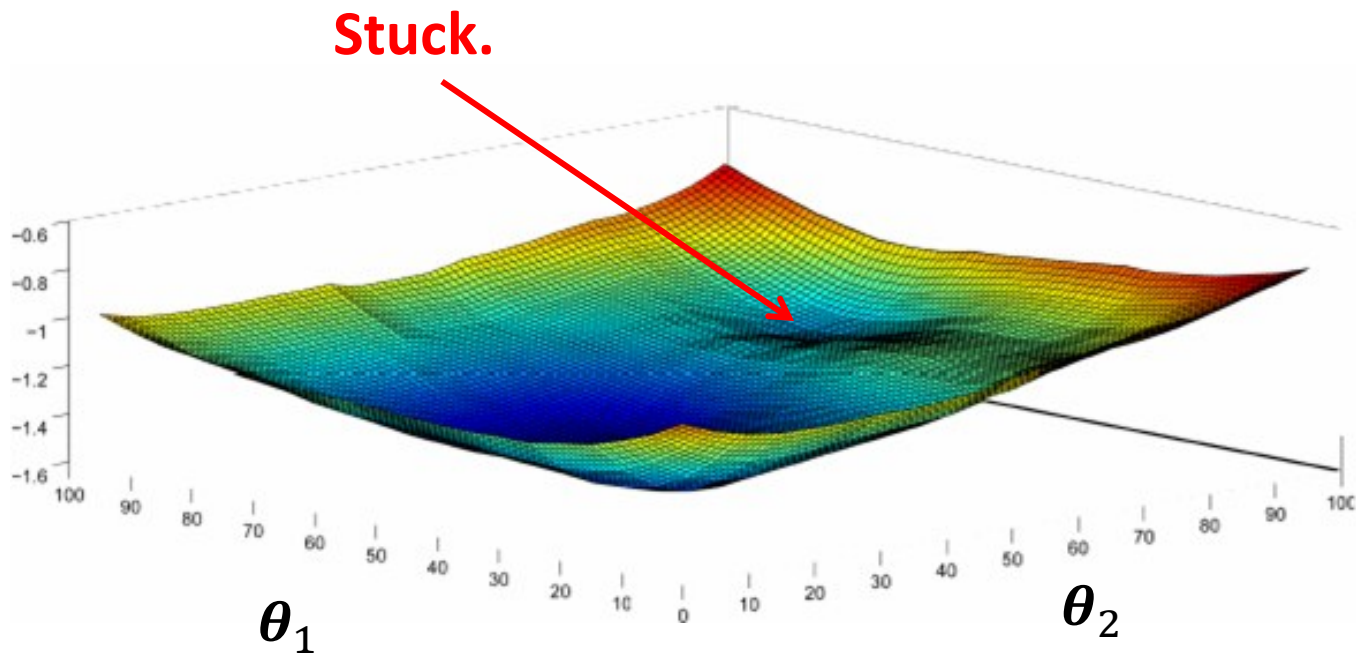
Local Optimization

- **Method 3, Quasi-Newton:** approximates B^{-1} incrementally from gradients gathered along the way (no second derivatives necessary)
- **Practical variant:** L-BFGS (Limited-Memory Broyden-Fletcher-Goldfarb-Shanno) – one of the best known “off-the-shelf” second-order optimizers, *very easy* to use

```
>> xopt = liblbfgs (@myfunc, x0)
```

Problems

- Local Optimization can run into *local minima* or *starve on plateaus*, except if the function has some suitable properties...



Convexity

- Convex functions have exactly one local optimum, also for every \mathbf{x} , \mathbf{y} and $t \in [0,1]$ the following holds:

$$f(t\mathbf{x} + (1 - t)\mathbf{y}) \leq tf(\mathbf{x}) + (1 - t)f(\mathbf{y})$$

- This makes them *exactly solvable* by local optimization methods
- Surprisingly many problem types can be formulated as *convex optimization problems*

Smoothness

- Many relevant BCI problems amount to optimization problems with non-differentiable features (e.g., sparse problems), so smooth optimization is not immediately applicable
- **Fixes:** Can use smooth surrogate functions (e.g., *proximal optimization*), or solve an equivalent problem that is smooth (see *convex duality*) or split the non-smooth terms off (see *operator splitting*), or use non-smooth methods (e.g., *subgradient descent*)





8.2 Going Beyond CSP

(and a bit of equation juggling!)

Transforming CSP

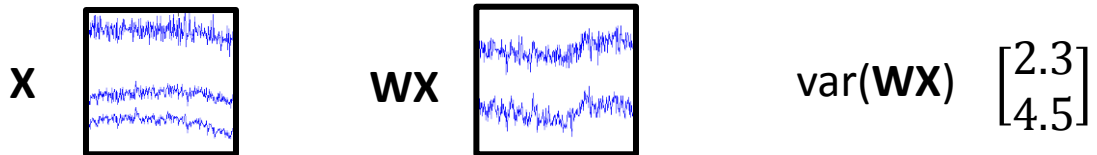
- Consideration: Given a zero-mean trial $\mathbf{X} \in \mathbb{R}^{C \times T}$ with covariance $\mathbf{\Sigma} \in \mathbb{R}^{C \times C}$, spatial filters $\mathbf{W} \in \mathbb{R}^{S \times C}$, linear weights $\boldsymbol{\theta} \in \mathbb{R}^S$ and bias b

Transforming CSP

- Consideration: Given a zero-mean trial $\mathbf{X} \in \mathbb{R}^{C \times T}$ with covariance $\mathbf{\Sigma} \in \mathbb{R}^{C \times C}$, spatial filters $\mathbf{W} \in \mathbb{R}^{S \times C}$, linear weights $\boldsymbol{\theta} \in \mathbb{R}^S$ and bias b

- Omitting the log from CSP, we have:

$$y = b + \boldsymbol{\theta} \text{var}(\mathbf{WX})$$

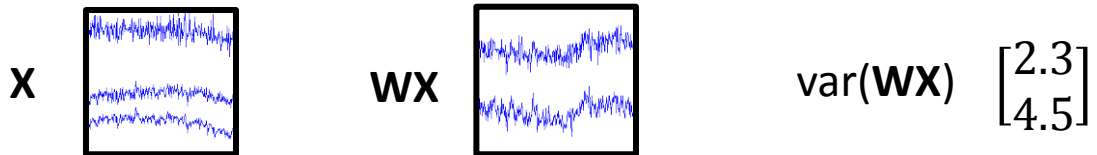


Transforming CSP

- Consideration: Given a zero-mean trial $\mathbf{X} \in \mathbb{R}^{C \times T}$ with covariance $\mathbf{\Sigma} \in \mathbb{R}^{C \times C}$, spatial filters $\mathbf{W} \in \mathbb{R}^{S \times C}$, linear weights $\boldsymbol{\theta} \in \mathbb{R}^S$ and bias b

- Omitting the log from CSP, we have:

$$y = b + \boldsymbol{\theta} \text{var}(\mathbf{W}\mathbf{X})$$



- Rewriting in terms of individual spatial filters \mathbf{W}_k :

$$y = b + \sum_{k=1}^S \boldsymbol{\theta}_k \text{var}(\mathbf{W}_k \mathbf{X})$$

Transforming CSP

- The variance term can be expressed using the covariance matrix Σ of segment \mathbf{X} :

$$y = b + \sum_{k=1}^S \theta_k \text{var}(\mathbf{W}_k \mathbf{X}) = b + \sum_{k=1}^S \theta_k (\mathbf{W}_k \Sigma \mathbf{W}_k^T)$$

Transforming CSP

- The variance term can be expressed using the covariance matrix Σ of segment X :

$$y = b + \sum_{k=1}^S \theta_k \text{var}(\mathbf{W}_k X) = b + \sum_{k=1}^S \theta_k (\mathbf{W}_k \Sigma \mathbf{W}_k^T)$$

- And $\mathbf{W}_k \Sigma \mathbf{W}_k^T$ can be replaced by the inner product between two matrices $\langle \mathbf{W}_k \mathbf{W}_k^T, \Sigma \rangle$

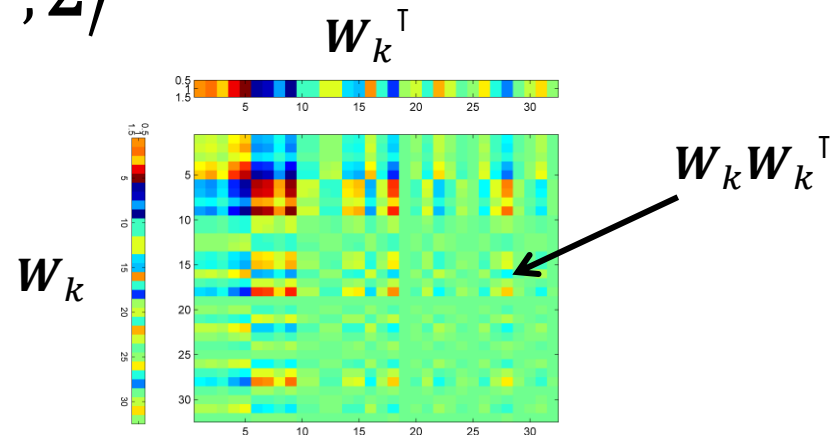
Transforming CSP

- The variance term can be expressed using the covariance matrix Σ of segment X :

$$y = b + \sum_{k=1}^S \theta_k \text{var}(\mathbf{W}_k X) = b + \sum_{k=1}^S \theta_k (\mathbf{W}_k \Sigma \mathbf{W}_k^\top)$$

- And $\mathbf{W}_k \Sigma \mathbf{W}_k^\top$ can be replaced by the inner product between two matrices $\langle \mathbf{W}_k \mathbf{W}_k^\top, \Sigma \rangle$

$$b + \sum_{k=1}^S \theta_k \langle \mathbf{W}_k \mathbf{W}_k^\top, \Sigma \rangle$$



Transforming CSP

- The variance term can be expressed using the covariance matrix Σ of segment X :

$$y = b + \sum_{k=1}^S \theta_k \text{var}(\mathbf{W}_k X) = b + \sum_{k=1}^S \theta_k (\mathbf{W}_k \Sigma \mathbf{W}_k^\top)$$

- And $\mathbf{W}_k \Sigma \mathbf{W}_k^\top$ can be replaced by the inner product between two matrices $\langle \mathbf{W}_k \mathbf{W}_k^\top, \Sigma \rangle$, and regrouped:

$$b + \sum_{k=1}^S \theta_k \langle \mathbf{W}_k \mathbf{W}_k^\top, \Sigma \rangle = b + \left\langle \sum_{k=1}^S \theta_k \mathbf{W}_k \mathbf{W}_k^\top, \Sigma \right\rangle$$

$$= b + \langle \boldsymbol{\Theta}, \Sigma \rangle$$

Transforming CSP

- Thus this form is *linear in the covariance matrix* of X :

$$y = b + \langle \boldsymbol{\theta}, \boldsymbol{\Sigma} \rangle = \mathbf{b} + \tilde{\boldsymbol{\theta}} \text{vec}(\boldsymbol{\Sigma})$$

- Could again learn $\tilde{\boldsymbol{\theta}}$ using a simple linear method (e.g., LDA), but *very* high-dimensional (#parameters= C^2)
- Need a method suitable for large-scale problems

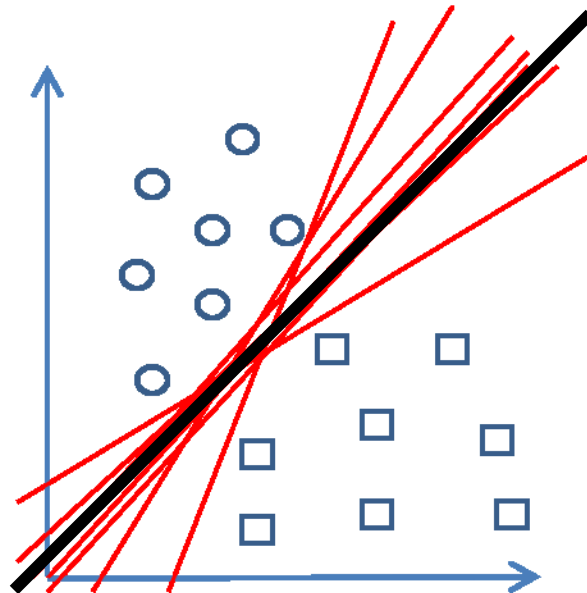




8.3 Large-Scale Machine Learning

Large-Scale Machine Learning

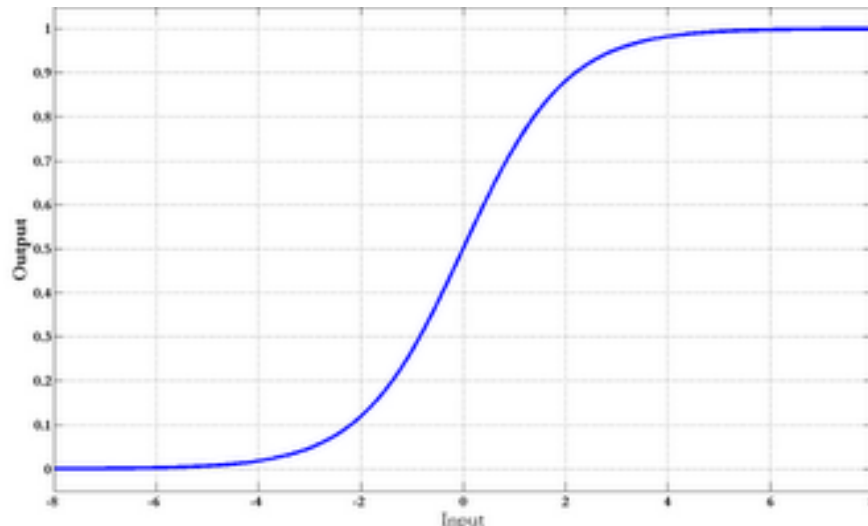
- *Discriminative* learning approaches like Support Vector Machines (SVMs) and Generalized Linear Models (GLMs) are well-adapted to high-dimensional / large-scale problems
- These directly optimize the parameters θ given the data



Large-Scale Machine Learning

- Logistic Regression is a GLM that maps X onto binary outputs via a logistic “link function”

$$\frac{1}{1 + e^{-yf_{\theta}(X)}}, (y \in \{-1, +1\})$$



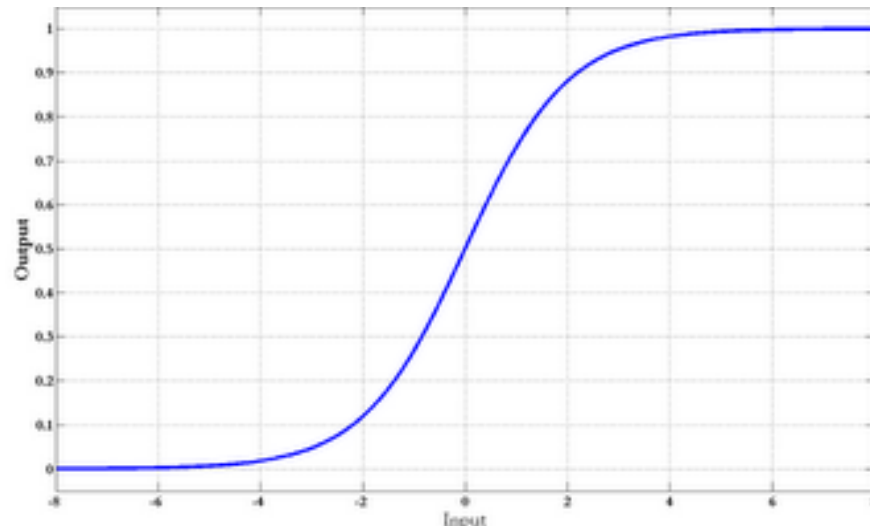
$- f_0(\mathbf{X}) \rightarrow$

Large-Scale Machine Learning

- Logistic Regression is a GLM that maps X onto binary outputs via a logistic “link function”

$$q_{\theta}(Y = y|X) = \frac{1}{1 + e^{-yf_{\theta}(X)}}, (y \in \{-1, +1\})$$

Interpreted as the
probability that $Y=1$
(or $Y=-1$)



$- f_0(X) \rightarrow$

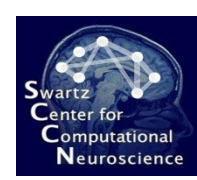


Large-Scale Machine Learning

- Logistic Regression is a GLM that maps X onto binary outputs via a logistic “link function”

$$q_{\theta}(Y = y|X) = \frac{1}{1 + e^{-yf_{\theta}(X)}}, (y \in \{-1, +1\})$$

- ... and linear function $f_{\theta}(X)$
$$f_{\theta}(X) = \theta X + b$$



Large-Scale Machine Learning

- θ can be obtained via off-the-shelf convex optimization tools (such as CVX) by solving the problem

$$\min_{\theta} \log(1 + e^{-y f_{\theta}(X)})$$

- The $\log(\dots)$ term is called the *logistic loss* and quantifies the misfit between predicted labels and true labels, for a particular choice of θ

Large-Scale Machine Learning

- For large problems, solution is still prone to *over-fitting* to random noise in the data – need to plug in some *additional assumptions*

$$\min_{\boldsymbol{\theta}} \log(1 + e^{-y f_{\boldsymbol{\theta}}(X)}) + \lambda \Omega(\boldsymbol{\theta})$$

- Many choices for regularization term Ω
 - $\Omega(\boldsymbol{\theta}) = \|\boldsymbol{\theta}\|^2$ encourages small weights
 - $\Omega(\boldsymbol{\theta}) = \|\boldsymbol{\theta}\|_1 = |\boldsymbol{\theta}_1| + |\boldsymbol{\theta}_2| + \dots$ encourages *sparsity*
 - can also get sparsity on groups of weights
 - combinations of these, ...





8.4 Application to the Spectral Model



Applying to the Spectral Model

- In the previous supervised oscillatory model $y = b + \langle \Theta, \Sigma \rangle$, the matrix-shaped Θ allows for a special matrix norm regularization $\Omega(\Theta)$:

$$\min_{\Theta} \log(1 + e^{-yf_{\Theta}(x)}) + \dots$$

Applying to the Spectral Model

- In the previous supervised oscillatory model $y = b + \langle \Theta, \Sigma \rangle$, the matrix-shaped Θ allows for a special matrix norm regularization $\Omega(\Theta)$:

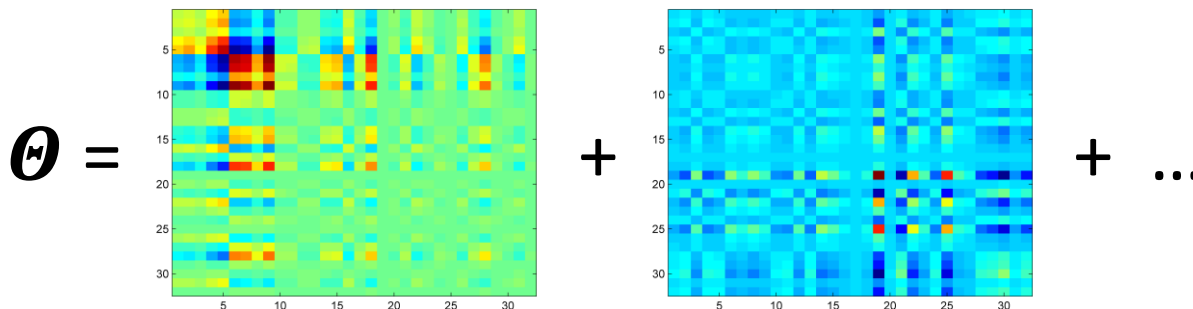
$$\min_{\Theta} \log(1 + e^{-yf_{\Theta}(x)}) + \lambda \sum_{k=1}^{\text{rank}(\Theta)} \sigma_k(\Theta)$$

Applying to the Spectral Model

- In the previous supervised oscillatory model $y = b + \langle \Theta, \Sigma \rangle$, the matrix-shaped Θ allows for a special matrix norm regularization $\Omega(\Theta)$:

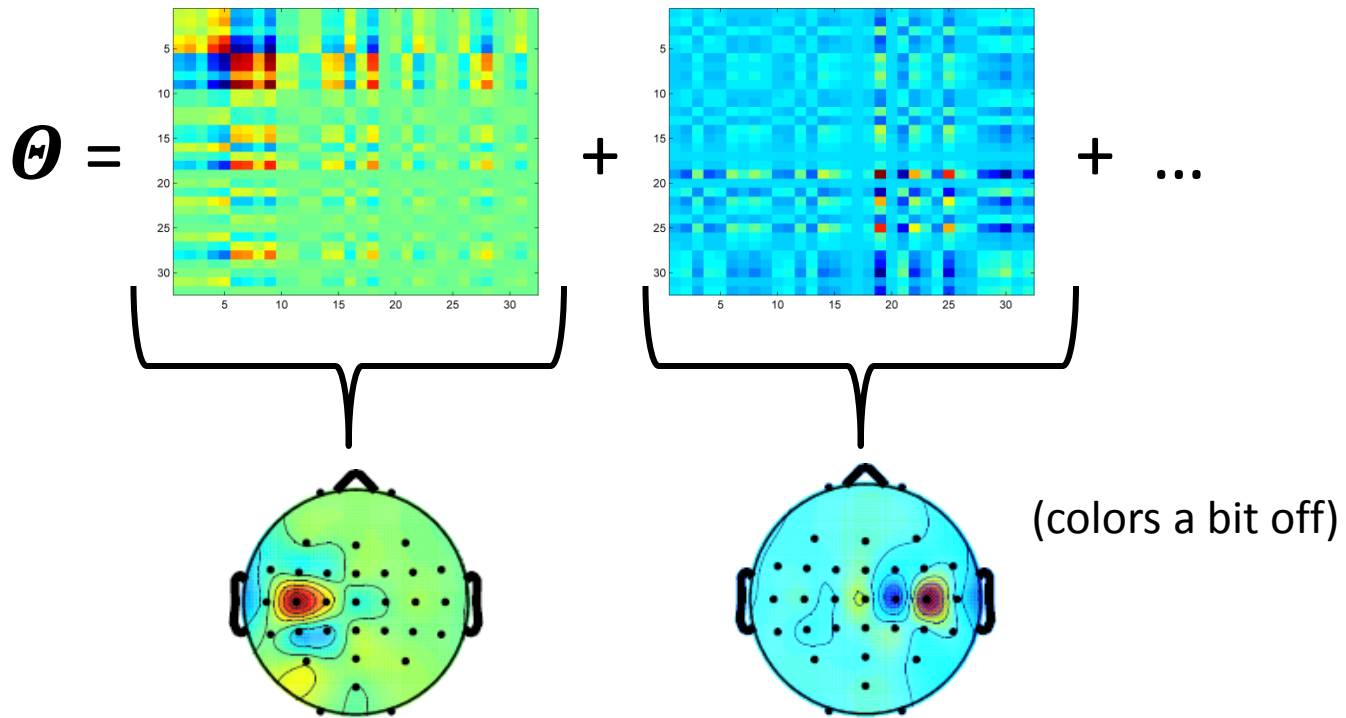
$$\min_{\Theta} \log(1 + e^{-y f_{\Theta}(x)}) + \lambda \sum_{k=1}^{\text{rank}(\Theta)} \sigma_k(\Theta)$$

- This encourages a low-rank structure in Θ , i.e.



Applying to the Spectral Model

- Thus, the weight matrix is equivalent to the weighted sum of a small set of spatial filters applied to the covariance matrix of the signal!







8.5 Application to ERPs

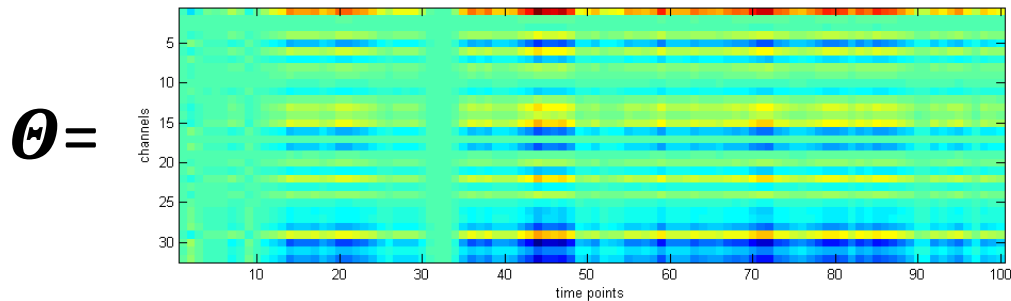
Application to ERPs

- Same approach can be applied to the raw data epoch \mathbf{X} instead of its covariance matrix Σ
- So we optimize for a GLM $y = b + \langle \boldsymbol{\theta}, \mathbf{X} \rangle$

$$\min_{\boldsymbol{\theta}} \log(1 + e^{-y f_{\boldsymbol{\theta}}(\mathbf{x})}) + \lambda \sum_{k=1}^{\text{rank}(\boldsymbol{\theta})} \sigma_k(\boldsymbol{\theta})$$

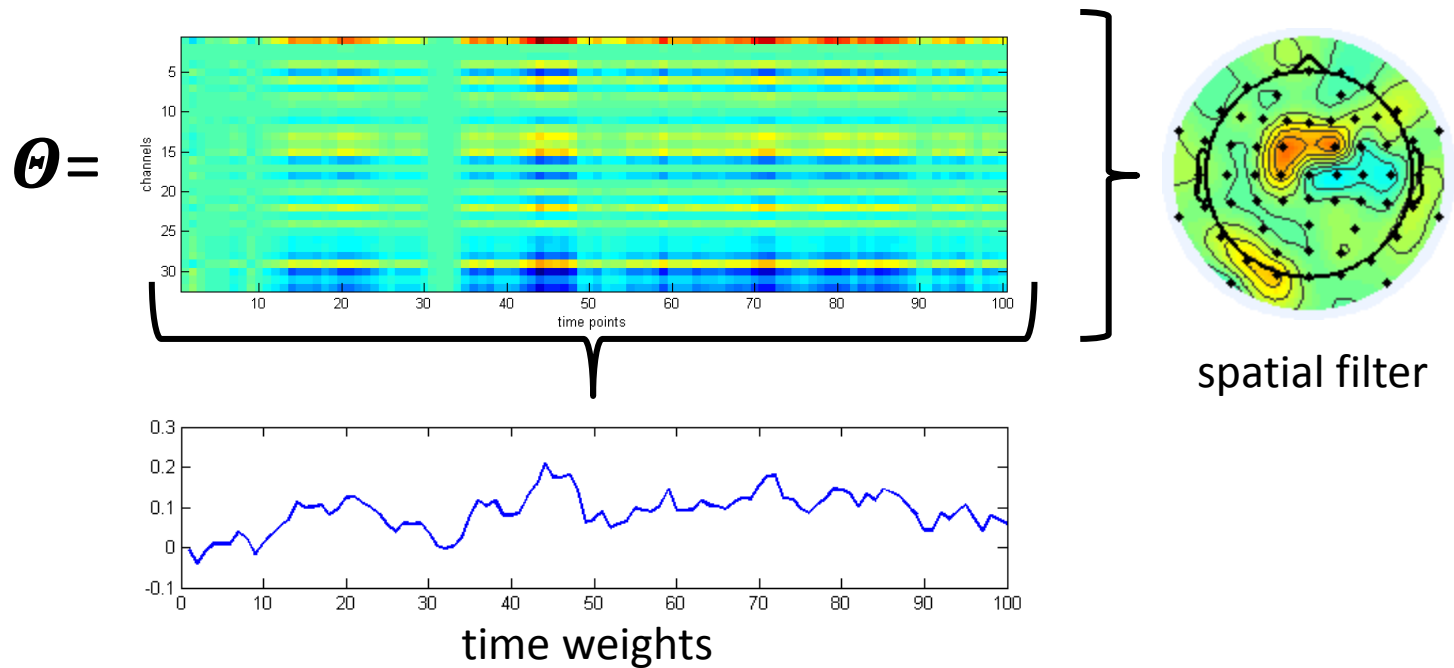
Application to ERPs

- Same approach can be applied to the raw data epoch \mathbf{X} instead of its covariance matrix Σ
- This learns a linear ERP weight matrix with one weight for each channel and time point



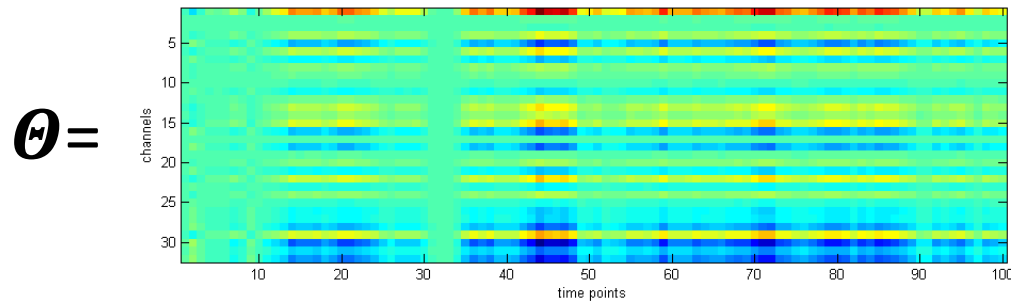
Application to ERPs

- Same approach can be applied to the raw data epoch \mathbf{X} instead of its covariance matrix Σ
- This learns a linear ERP weight matrix with one weight for each channel and time point



Application to ERPs

- Same approach can be applied to the raw data epoch \mathbf{X} instead of its covariance matrix Σ
- This learns a linear ERP weight matrix with one weight for each channel and time point



- Θ is low-rank, so corresponds to a sum of a few spatial filters and their weights over time

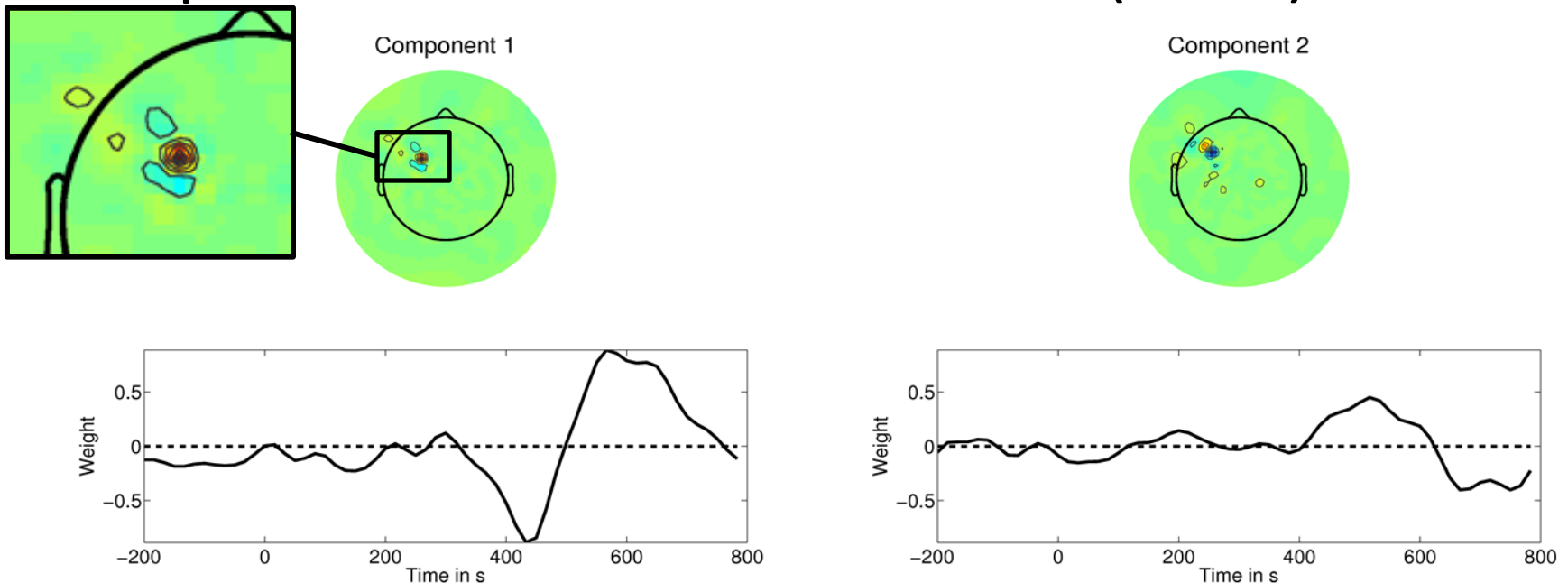


Application to ERPs

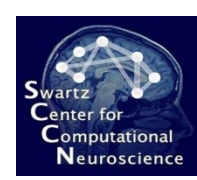
- Thus, no hand-selected time windows needed
- Also, results are regularized to find / pick up few sources and their relevant time courses

Application to ERPs

- Thus, no hand-selected time windows needed
- Also, results are regularized to find / pick up few sources and their relevant time courses
- Rapid Serial Visual Presentation (RSVP) task:



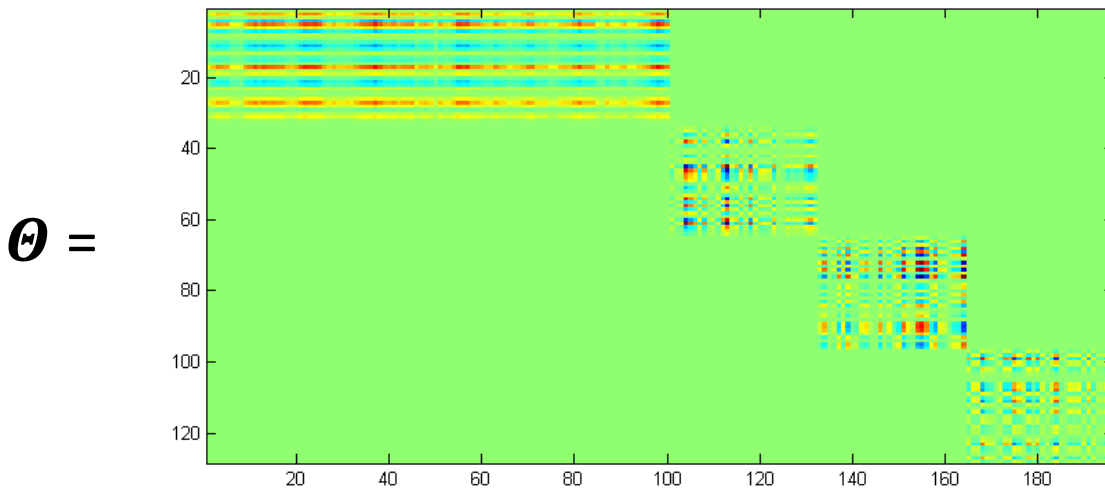




8.6 Learning ERP and Oscillatory Weights Simultaneously

Feature Combination

- If for each trial instead of the covariance matrix or the raw ERP, a *block-diagonal concatenation* of both is used, the method learns a low-rank weight matrix combining both features simultaneously



Feature Combination

- If for each trial instead of the covariance matrix or the raw ERP, a *block-diagonal concatenation* of both is used, the method learns a low-rank weight matrix combining both features simultaneously
- For a block-diagonal weight matrix Θ , it holds that

$$\min_{\Theta} \dots + \lambda \sum_{k=1}^{\text{rank}(\Theta)} \sigma_k(\Theta)$$

is equivalent to solving for its blocks Θ_b :

$$\min_{\Theta} \dots + \lambda \sum_{b=1}^B \sum_{k=1}^{\text{rank}(\Theta_b)} \sigma_k(\Theta_b)$$

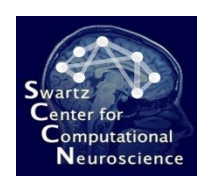


Sparsity and Feature Selection

- Also, covariance matrices for multiple frequency bands and time windows can be concatenated
- For regularizers that are a sum of terms (of a certain type), most terms in the sum will be driven to zero and only a *sparse subset* of terms remains non-zero, i.e., the relevant features are selected automatically

$$\min_{\theta} \dots + \lambda \sum_{b=1} \dots$$

- Recovery of the relevant support is statistically extremely efficient – it holds that the number of irrelevant dimensions under which the support can be accurately recovered is *exponential* in the number of observations (i.e., trials) [Ng 1998]



Sparsity and Feature Selection

- Thus, only the relevant subset of frequencies or time windows (for covariance) or ERP sources is typically learned
- Can be taken even further, e.g., could encourage weights for different time/frequency bins to share a small set of spatial filters (again using rank constraints on concatenated matrices) – this is called multi-task learning

Final Prediction Functions

- **Basic oscillatory case** (assuming \mathbf{X} is band-passed):

$$y = \frac{1}{1 + e^{- (b + \langle \boldsymbol{\theta}, \mathbf{X}\mathbf{X}^\top \rangle)}}$$

- **ERP case** (\mathbf{X} can be band-passed):

$$y = \frac{1}{1 + e^{- (b + \langle \boldsymbol{\theta}, \mathbf{X} \rangle)}}$$

- **Combined cases** (here for temporal filters \mathbf{F}_1 and \mathbf{F}_2):

$$y = \frac{1}{1 + e^{- \left\langle \boldsymbol{\theta}, \begin{bmatrix} \mathbf{X} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{X}\mathbf{F}_1\mathbf{X}^\top & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{X}\mathbf{F}_2\mathbf{X}^\top \end{bmatrix} \right\rangle - b}}$$





8.7 Practical Remarks

Solving It

- Problems of this size are impossible to solve using CVX; need a custom solver
- For a wide range of sparse estimation problems the DAL (Dual-Augmented Lagrangian) solver is applicable and *very fast*
- For an even wider range of problems the ADMM (Alternating Direction Method of Multipliers) framework is applicable and also very fast

ADMM

- Framework for distributed very large scale optimization – leads to parallel algorithms
- Can be done with *very* simple MATLAB code

These are often fairly simple problems

$$x^{k+1} := \operatorname{argmin}_x L_\rho(x, z^k, y^k)$$

$$z^{k+1} := \operatorname{argmin}_z L_\rho(x^{k+1}, z, y^k)$$

$$y^{k+1} := y^k + \rho(Ax^{k+1} + Bz^{k+1} - c)$$



Other Methods in this Framework

- **Support Vector Machines:** use a different loss function (“hinge loss” instead of logistic loss)
- **Multiple Kernel Learning:** using group sparsity on kernel matrices (selecting few kernels)
- **Hierarchical Kernel Learning:** very advanced non-linear feature selection approach using tree-structured sparsity
- **Linear Regression:** Usable for a continuous output space instead of discrete





L8 Questions?