



BME I5100: Biomedical Signal Processing

Linear Discrimination



Lucas C. Parra
Biomedical Engineering Department
City College of New York





Schedule

Week 1: Introduction

Linear, stationary, normal - the stuff biology is **not** made of.

Week 1-4: Linear systems

Impulse response

Moving Average and Auto Regressive filters

Convolution

Discrete Fourier transform and z-transform

Sampling

Week 5-8: Random variables and stochastic processes

Random variables

Moments and Cumulants

Multivariate distributions

Stochastic processes

Week 9-14: Examples of biomedical signal processing

Probabilistic estimation

Harmonic analysis - **estimation** circadian rhythm and speech

Linear discrimination - **detection** of evoked responses in EEG/MEG

Independent components analysis - **analysis** of MEG signals

Dynamical Models - Kalman filter and Hidden Markov Models

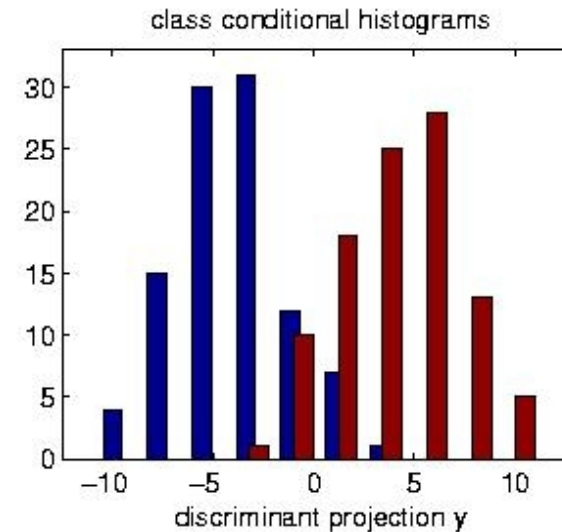
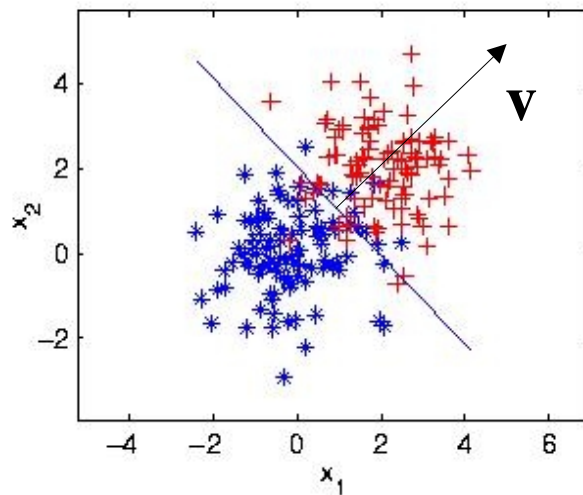
Matched and Wiener filter - **filtering** in ultrasound



Linear Discrimination

Given samples \mathbf{x} from two classes c_1 and c_2 find vector \mathbf{v} so that the projection y separates the two classes:

$$y = \mathbf{v}^T \mathbf{x} + v_0$$



Redefine $\mathbf{x} = [1, x_1, x_2, \dots, x_d]^T$ and $\mathbf{v} = [v_0, v_1, v_2, \dots, v_d]^T$ to write in short:

$$y = \mathbf{v}^T \mathbf{x}$$



Linear Discrimination - Logistic Regression

The goal is to find a mapping from \mathbf{x} to class label c or at least an expected value:

$$\hat{c} = f(y) = f(\mathbf{v}^T \mathbf{x})$$

However, since there is overlap, rather than making a fixed determination on the class label we will build a model that tells us what is the likelihood of a class c given input \mathbf{x} .

$$p(c|y) = p(c|\mathbf{v}^T \mathbf{x})$$

We will now derive an expression of the Likelihood of class labels c given projection y , which are given by input \mathbf{x} and parameters \mathbf{v} . The optimal \mathbf{v} results then from ML.



Linear Discrimination - Logistic Regression

Consider the posterior probability $p(c_1|y)$

$$p(c_1|y) = \frac{p(y|c_1) p(c_1)}{p(y|c_2) p(c_2) + p(y|c_1) p(c_1)} = \frac{1}{1 + \exp(-l)}$$

We introduced here the **Likelihood ratio** l :

$$l = \ln \frac{p(y|c_1) p(c_1)}{p(y|c_2) p(c_2)}$$

For a large class of distributions, $p(y|c)$, called exponential family, the Likelihood ratio simplifies under certain assumptions to

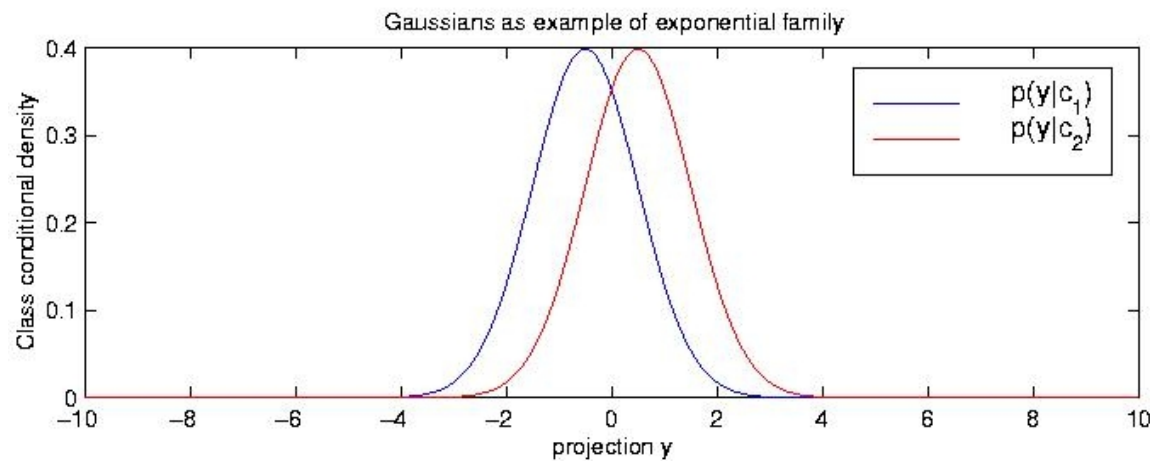
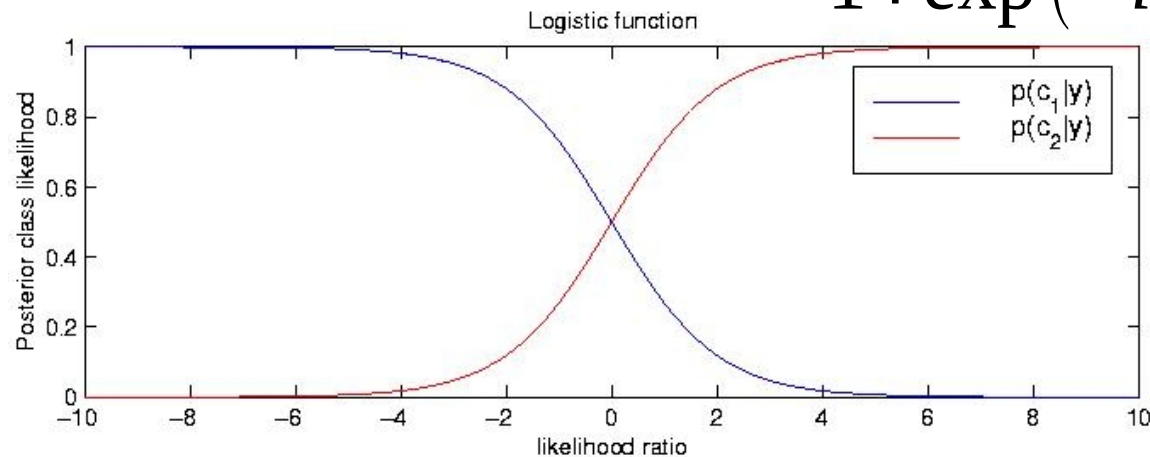
$$l = y + \ln \frac{p(c_1)}{p(c_2)}$$



Linear Discrimination - Logistic Regression

For example Gaussian projections, $p(y | c) = N(y; \mu_c, \sigma)$ with different mean for each class but the same standard deviation.

Logistic function:
$$p(c_1 | y) = \frac{1}{1 + \exp(-l)}$$





Linear Discrimination - Gaussian Data

Example: If the two classes are Gaussian distributed with different mean for each class but the *same covariance matrix* their optimal separation is linear and solution is particularly simple:

$$p(\mathbf{x}|c) \propto \exp\left[-(\mathbf{x}-\boldsymbol{\mu}_c)^T \boldsymbol{\Sigma}^{-1}(\mathbf{x}-\boldsymbol{\mu}_c)/2\right]$$

The likelihood ration is then linear

$$l = \ln \frac{p(\mathbf{x}|c_1)p(c_1)}{p(\mathbf{x}|c_2)p(c_2)} = \mathbf{v}^T \mathbf{x} + v_0$$

where the separation vector and bias are given by the means and covariance:

$$\mathbf{v} = \boldsymbol{\Sigma}^{-1}(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)$$

$$v_0 = -\frac{1}{2}\boldsymbol{\mu}_1^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_1 + \frac{1}{2}\boldsymbol{\mu}_2^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_2 + \ln \frac{p(c_1)}{p(c_2)}$$



Linear Discrimination - Logistic Regression

More generally, assume that y is distributed according to some distribution of the exponential family, which includes Gaussian, Bernoulli, Poisson, and others.

The goal is then to find the optimal projection vector \mathbf{v} such that the projections, $y = \mathbf{v}^T \mathbf{x}$, results in a logistic likelihood for the class labels

$$p(c_1|y) = \frac{1}{1 + \exp(-\mathbf{v}^T \mathbf{x})} = f(\mathbf{v}^T \mathbf{x})$$

Here we have absorbed $\ln p(c_1)/p(c_2)$ into v_0 .



Linear Discrimination - Logistic Regression

To derive the ML solution define, $c=0$, and, $c=1$, to identify the two different classes and denote in short $p(c|y)=f$. We can write then

$$p(c|y) = f^c (1-f)^{1-c}$$

This is called the Bernoulli density of c , and f is the mean:

$$\hat{c} = E[c] = f(\mathbf{v}^T \mathbf{x})$$

Note that the expected value we wanted to compute is therefore given now by the logistic function.



Linear Discrimination - Logistic Regression

Given i.i.d. samples $\mathbf{x}[k]$, $c[k]$ the log-likelihood of the data is

$$\begin{aligned} L(\mathbf{v}) &= \ln \prod_k p(c[k], \mathbf{x}[k] | \mathbf{v}) \\ &= \sum_k \ln p(c[k] | \mathbf{v}^T \mathbf{x}[k]) p(\mathbf{x}[k]) \\ &= \sum_k c[k] \ln f(\mathbf{v}^T \mathbf{x}[k]) + (1 - c[k]) \ln(1 - f(\mathbf{v}^T \mathbf{x}[k])) \\ &\quad + \text{const.} \end{aligned}$$

And the optimum solution according to ML is

$$\underset{\mathbf{v}}{\operatorname{argmin}} L(\mathbf{v})$$

There is not close form solution for this minimum.



Linear Discrimination - Logistic Regression

However, the minimum can be computed using a fast algorithm based on Iteratively Reweighted Least Squares (IRLS). It is a type of Newton-Raphson gradient descent algorithm called Fisher Scoring method (McCullagh, Nelder 1983):

$$\mathbf{v}_{t+1} = \mathbf{v}_t - E \left[\frac{\partial L(\mathbf{v})}{\partial \mathbf{v} \partial \mathbf{v}^T} \right]^{-1} \frac{\partial L(\mathbf{v})}{\partial \mathbf{v}}$$

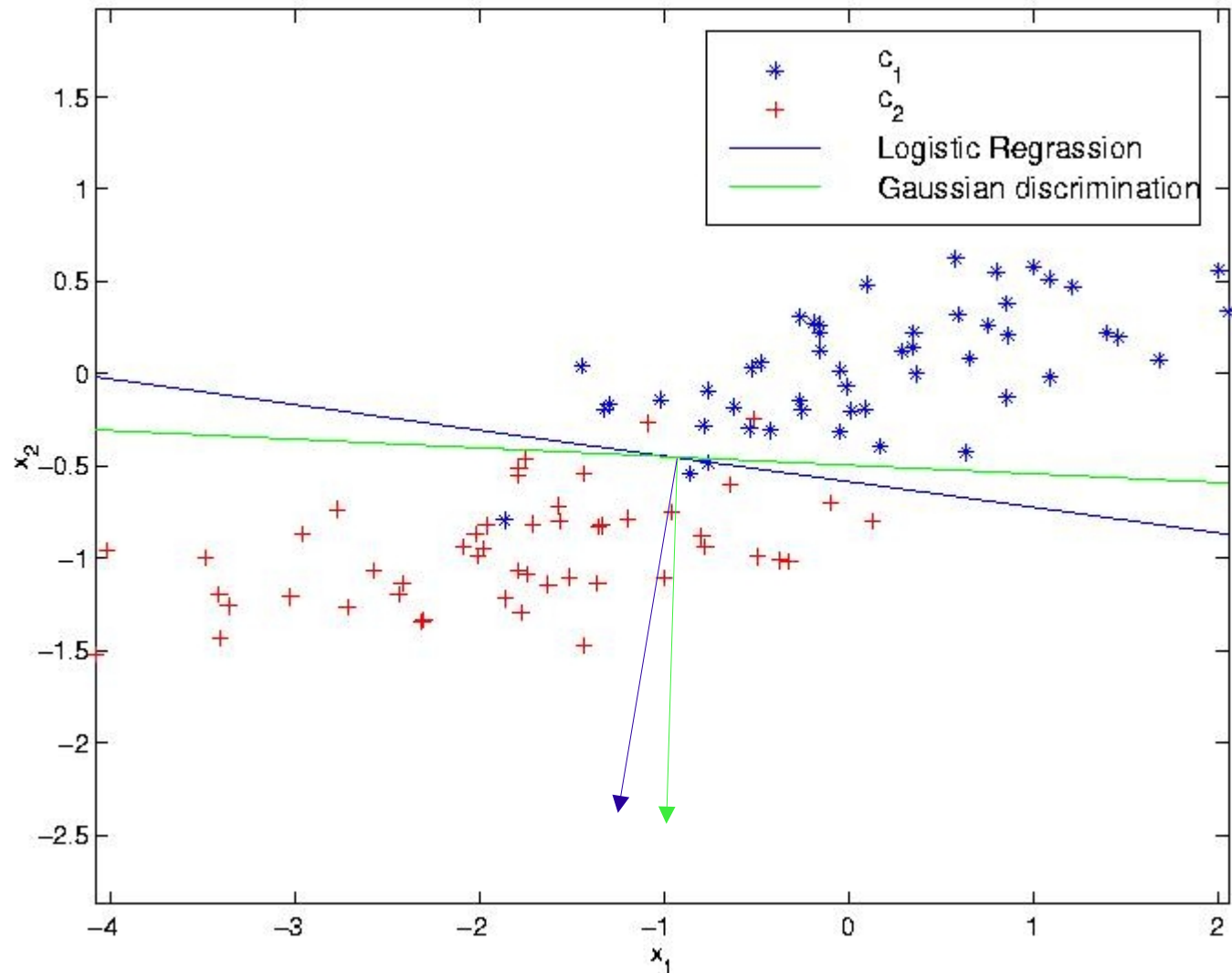
The expected Hessian can be computed fast and converges typically within a few iterations.

```
>> v = logist(x,c); % not part of matlab
```



Linear Discrimination - Comparison

Gaussian solution and minimum $L(\mathbf{v})$ give different results when data not Gaussian or for small sample size.

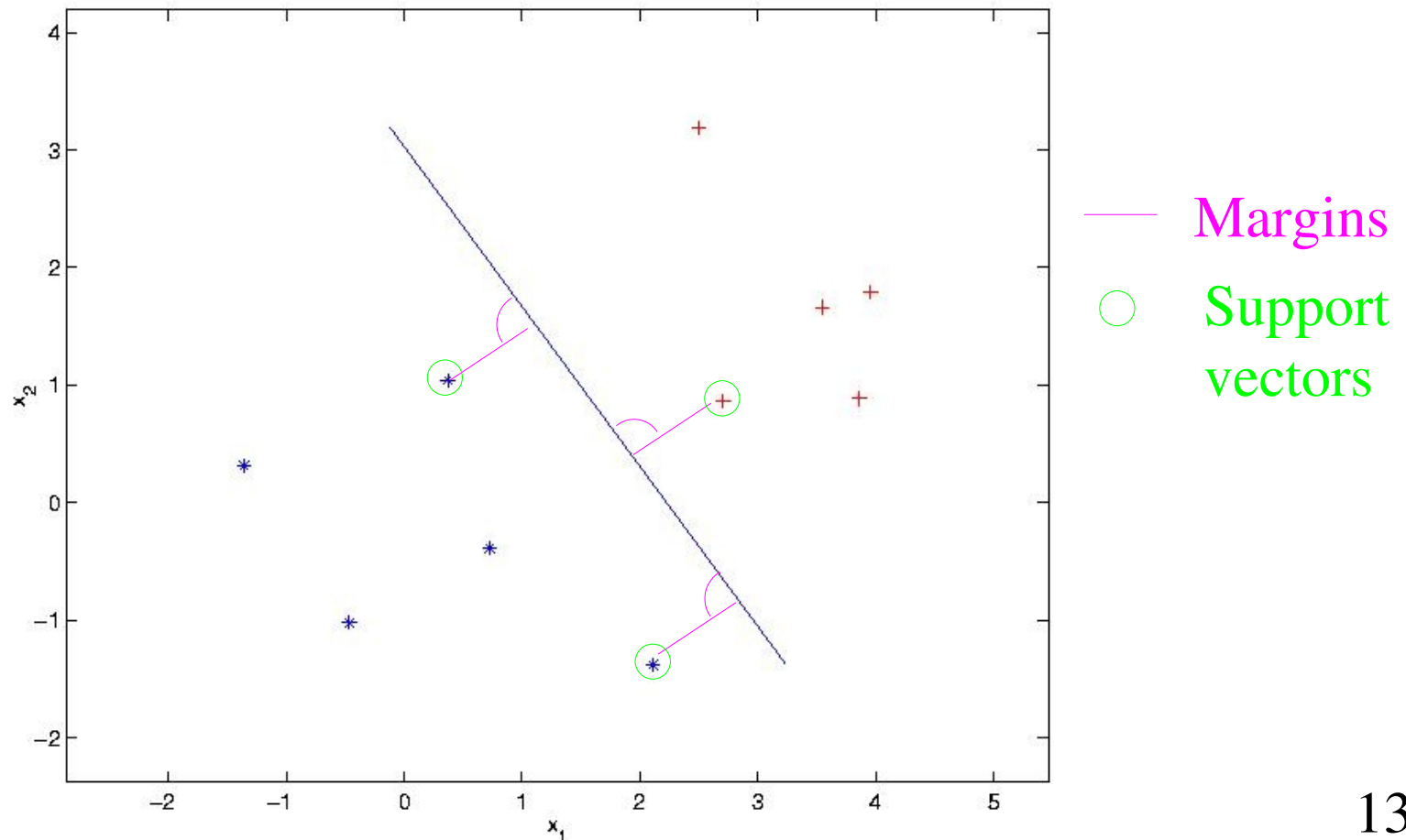




Linear Discrimination - Support Vectors

In case of perfect separability LR is not well defined.

Occurs often when we have very few samples and/or high dimensions. In that case is it better to chose the separation that **maximizes the margin to support vectors**.

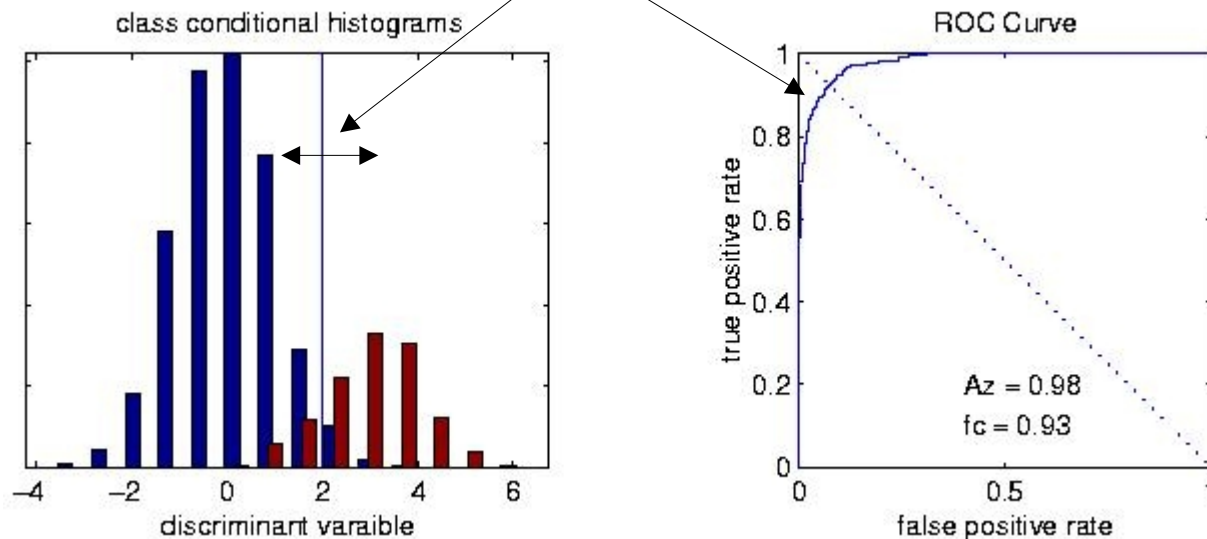




Linear Discrimination - Performance

The performance of a binary classification problem is typically evaluated with a **Receiver Operator Characteristic (ROC)** curve:

Moving threshold sweeps a curve of true positive rate versus false positive rate



A_z : Area under the ROC curve measures performance independent of threshold. It is 0.5 for chance performance.

f_c : Fraction correct (1-error rate) is conventionally measured where $tp=1-fp$. It is 0.5 for chance performance.



Linear Discrimination - Leave-One-Out

Note that the performance on the training data is biased and always better than the performance on unseen data.

Therefore ROC and A_z has to always be computed on unseen test data!

If there is not sufficient data available to separate into training and test set one should use the **leave-one-out procedure**:

1. For each sample k
Find the optimal \mathbf{v} on all data except $c[k], \mathbf{x}[k]$.
Use this optimal \mathbf{v} to compute the leave-one-out class likelihood $p(c[k] | \mathbf{v}^T \mathbf{x}[k])$.
2. Compute A_z using leave-one-out class likelihoods.

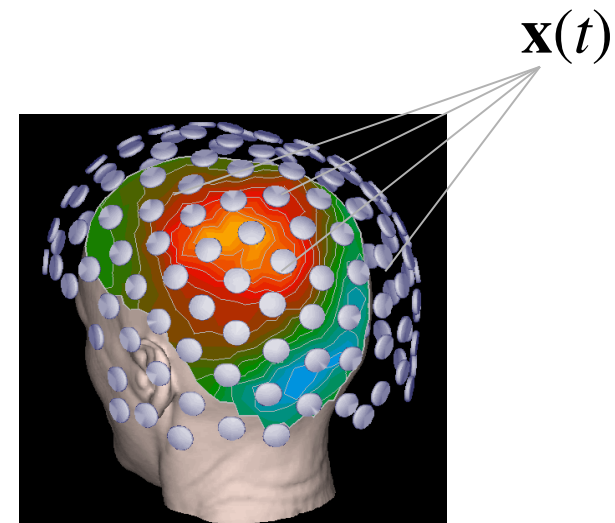


LD - Application to EEG and MEG

Conventional Event Related Potentials (ERP) averages over trials to increase signal to noise ratio.

The goal is to detect single trials without averaging over trials or over time. We substitute trial averaging by spatial integration.

$$y(t) = \mathbf{v}^T \mathbf{x}(t)$$



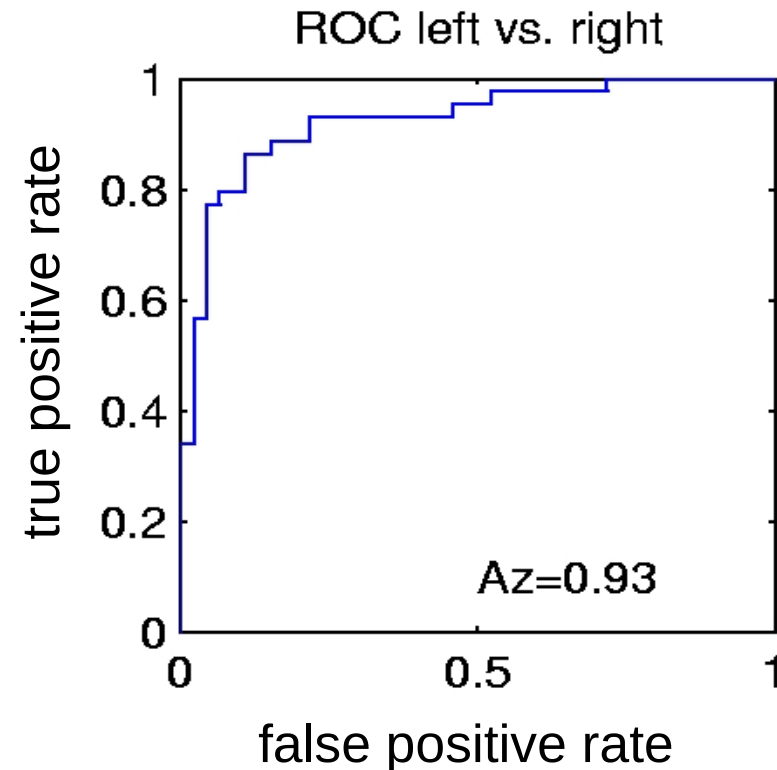
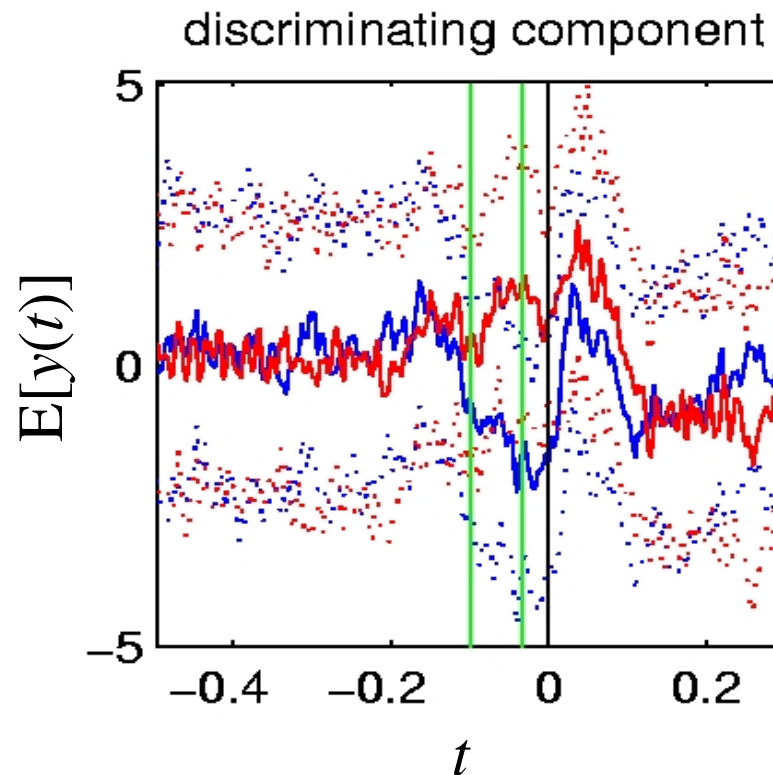
With Linear discrimination we can now compute spatial weights \mathbf{v} which maximally discriminate sensor array signals $\mathbf{x}(t)$ for two different conditions observed at times t_1 and t_2 .



LD - Application to EEG and MEG

Example: Find motor planing activity in MEG

Predict button press from 122 MEG sensors with linear discriminator w such that $y(t)$ differs the most during 100 ms to 30 ms *prior* to left (t_1) and right (t_2) button push.





LD - Application to EEG and MEG

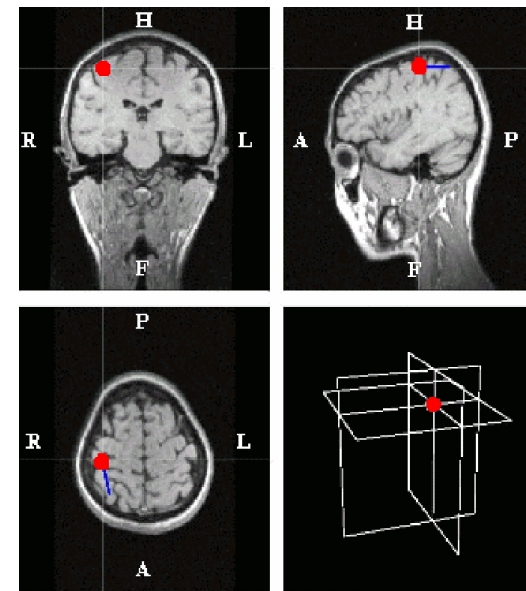
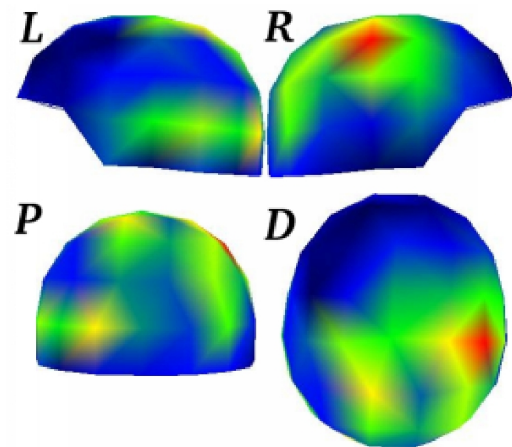
Localization of discriminating component: What is the electrical coupling \mathbf{a} of the hypothetical source \mathbf{y} that explains most of the activity \mathbf{X} ?

Least squares solution:

$$\mathbf{a} = \frac{\mathbf{X} \mathbf{y}}{\mathbf{y}^T \mathbf{y}}$$

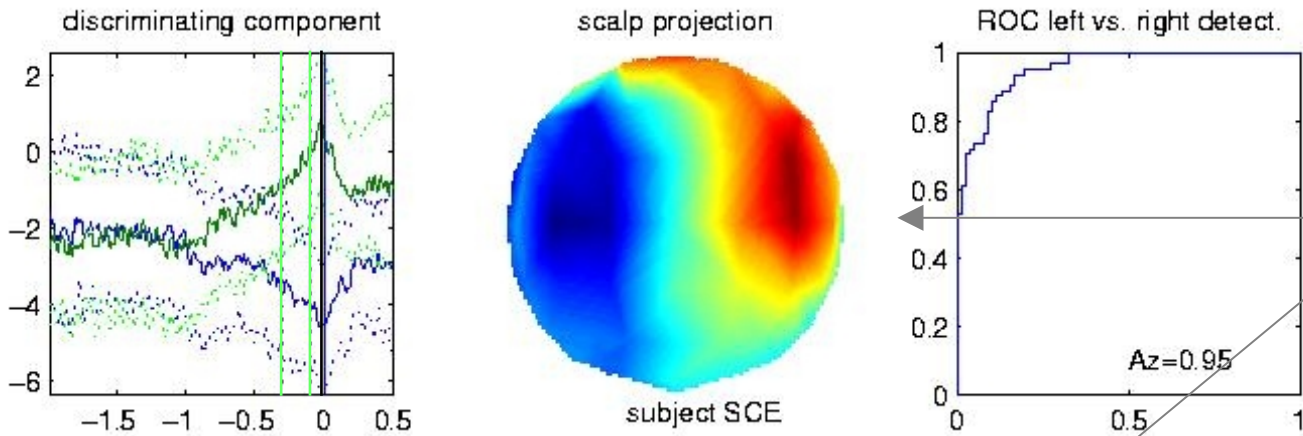
where \mathbf{X} has one column per sample, and \mathbf{y} is a vector with all samples. \mathbf{X} has to be zero mean across samples.

Strong coupling indicates low attenuation. Intensity on these “sensor projections” \mathbf{a} indicates closeness of the source to the sensors.



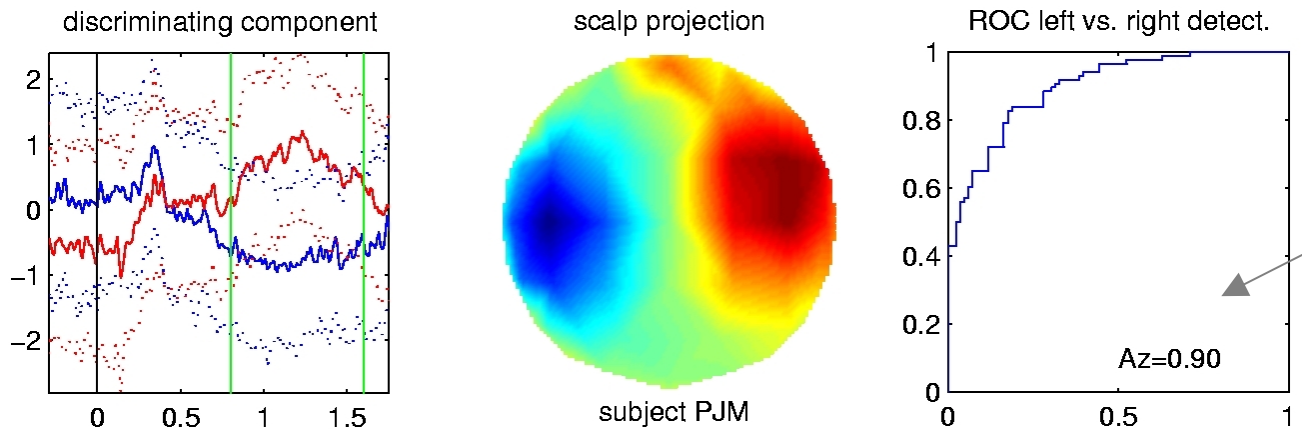
LD - Detection of Motor Planing and Imagery

Prediction of **explicit** finger tap (59 EEG sensors, 250-100ms prior)



Activity of explicit and imagined motor action have similar spatial distribution.

Detection of **imagined** finger tap (59 EEG sensors, 800 ms)



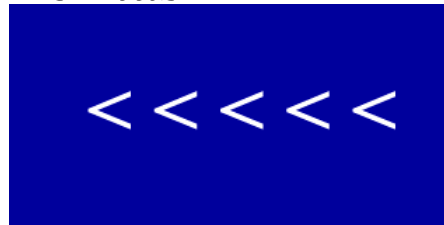
Transmission of binary information with this covert mental imagery results in communication at **12 bits/minute.**



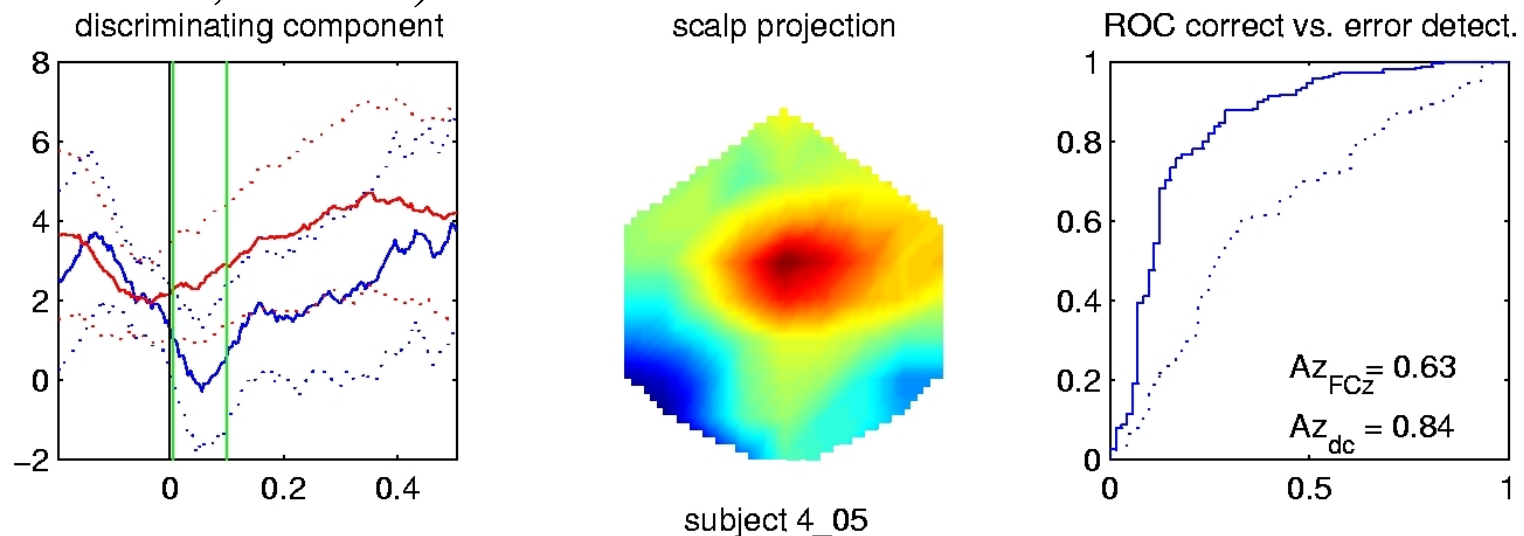
LD - Detection of Error Related Negativity

Error Related Negativity (ERN) occurs following perception of errors. It is hypothesized to originate in Anterior Cingulate and to represent response conflict or subjective loss.

Example: Erikson Flanker task



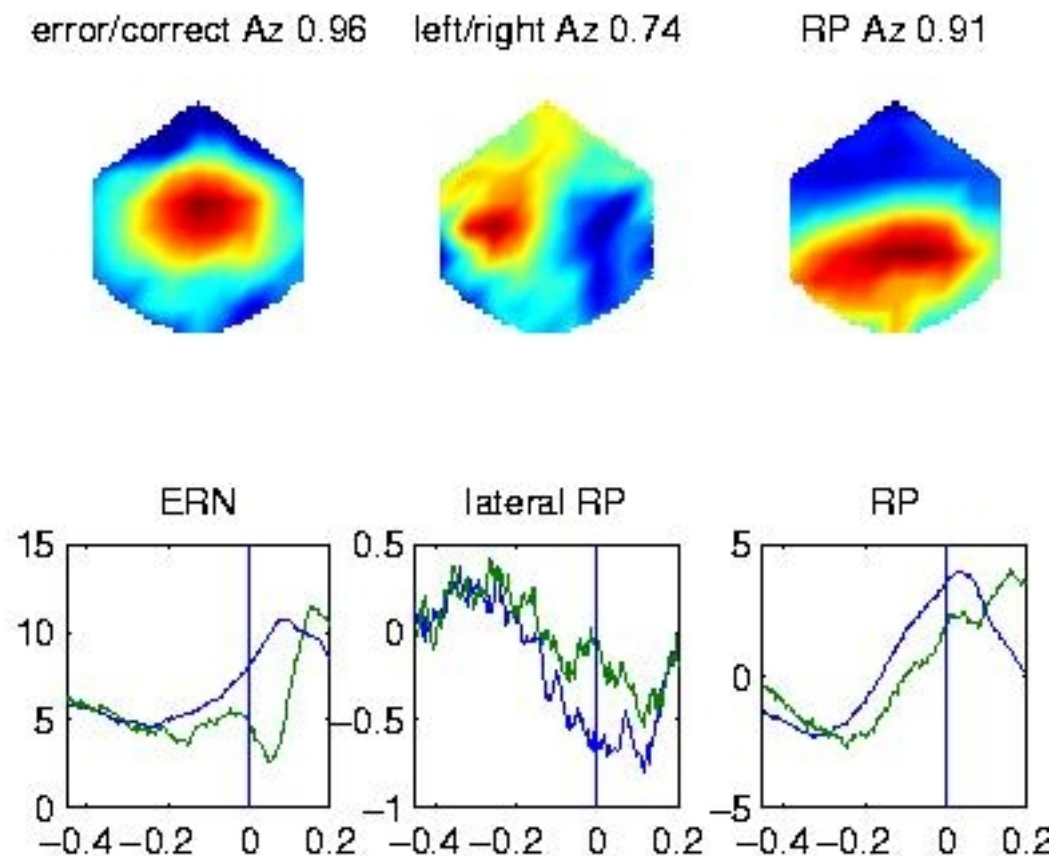
Discrimination of error versus correct response (64 EEG sensors, 100ms)





LD - Detection of Readiness Potential

In preparation of motion there is a potential buildup in the 200 ms prior to motion over motor areas. It is called readiness potential. When differentiating left/right one observed a lateralized readiness potential.





LD - Detection of Readiness Potential

Assignment 11:

Load `eeg-ern.mat` or generate overlapping random variables x_1 and x_2 .

If you use `eeg-ern.mat` note that for every trail (78 for x_1 and 300 for x_2) there are each 25 samples. Consider them all as i.i.d. samples.

Find a linear discriminator that discriminates between x_1 and x_2 assuming Gaussian distributions.

Plot resulting v and show ROC and A_z value on training data.

Optional:

Display scalp projection a using `scalp(coord, a)`

Show ROC curve and A_z with the average y over 25 samples of each trial.