

On the sparse estimation

ATR Computational Neuroscience
Laboratories

Masa-aki Sato

Contents

1. Why sparse estimation is necessary?

Generalization ability for ill posed problem

2. Why sparse estimation can be achieved?

Role of Bayesian estimation

3. Example of sparse estimation for real problem

Ill-posed problem

Large number of parameters are estimated
from small number of data

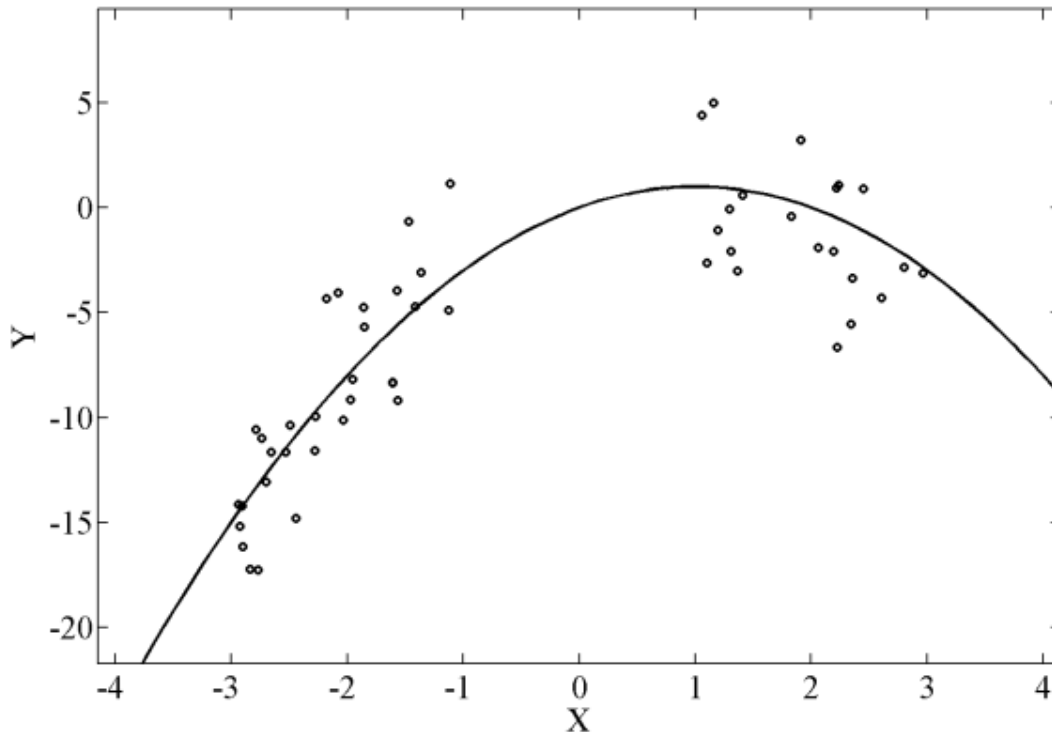
- **Maximum Likelihood (Minimum Squared Error)**
 - Estimate optimal parameter which maximize likelihood
 - Overfitting:
Complex models with many adjustable parameters tend to fit noise in training data and degrade generalization ability
- **Sparse estimation**
 - Extract relevant features and discard irrelevant features to attain good generalization ability

Function approximation by polynomial

$$y = w_0 + w_1x + w_2x^2 + \dots + w_Nx^N$$
$$= W \cdot X \quad \text{Linear parameter model}$$

$$X = [1, x, x^2, \dots, x^N] \quad \text{Input (feature) variable}$$

$$W = [w_0, w_1, w_2, \dots, w_N] \quad \text{Weight parameter vector}$$



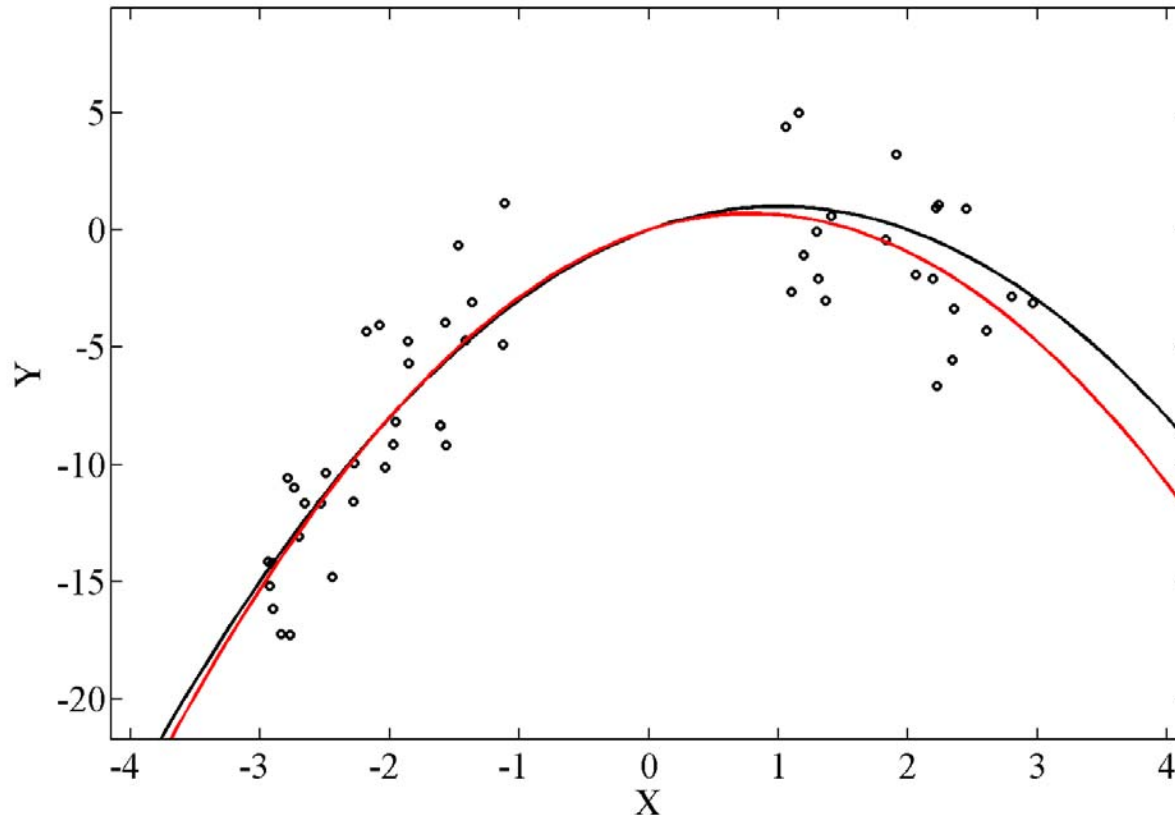
50 data points are
randomly generated
from quadratic function

$$y = 2x - x^2 + \textit{noise}$$

Maximum likelihood (Quadratic function)

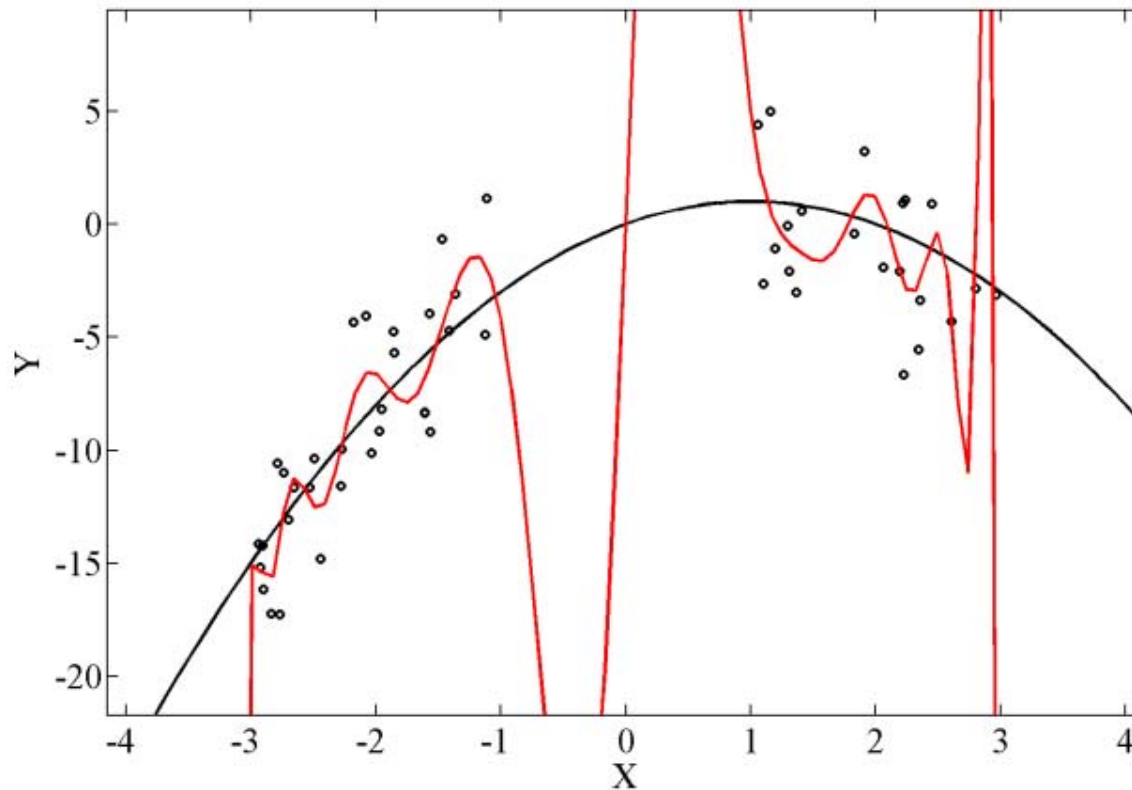
↔ Minimize Squared Error

$$error = \sum_{t=1}^T (y(t) - W \cdot X(t))^2 \quad \text{Find optimal } W$$



Maximum likelihood (20 degree polynomial)

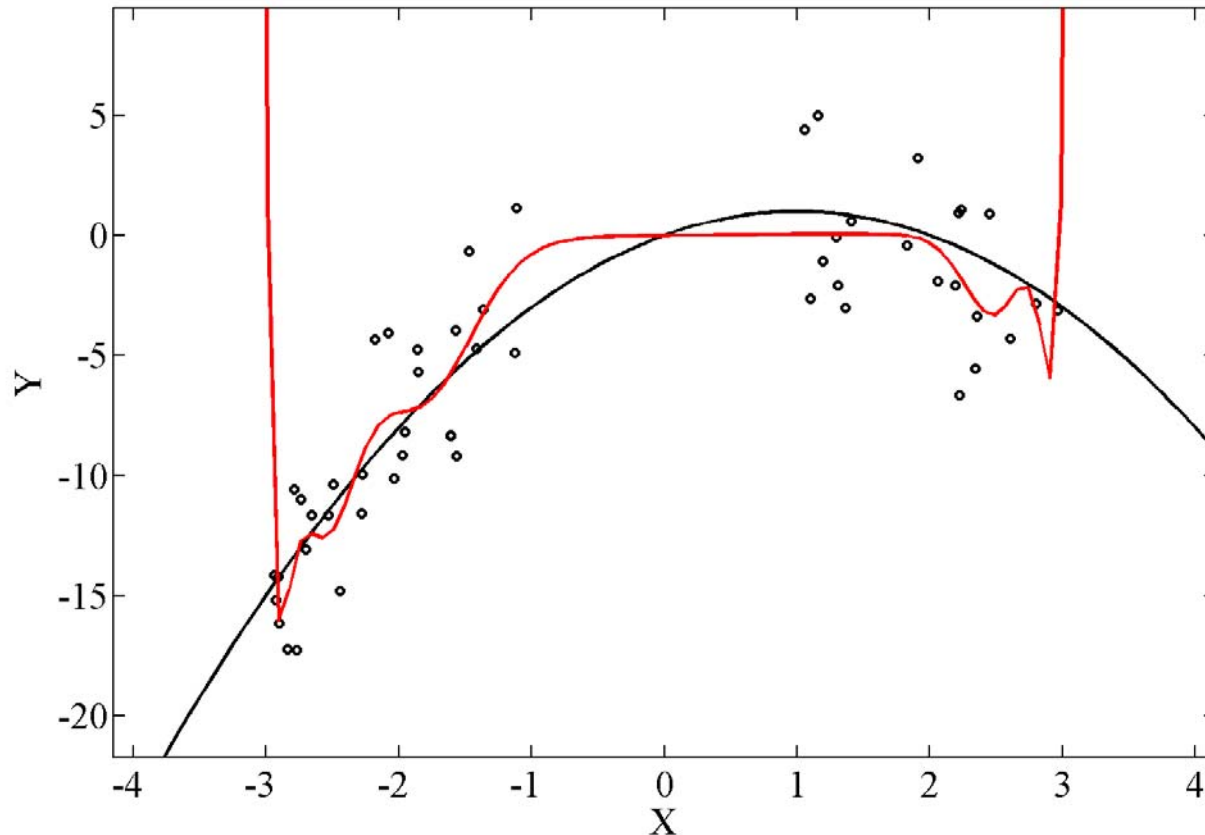
Overfitting : Optimal W fits noise in training data
Generalization is not good



Regularization method (20 degree polynomial)

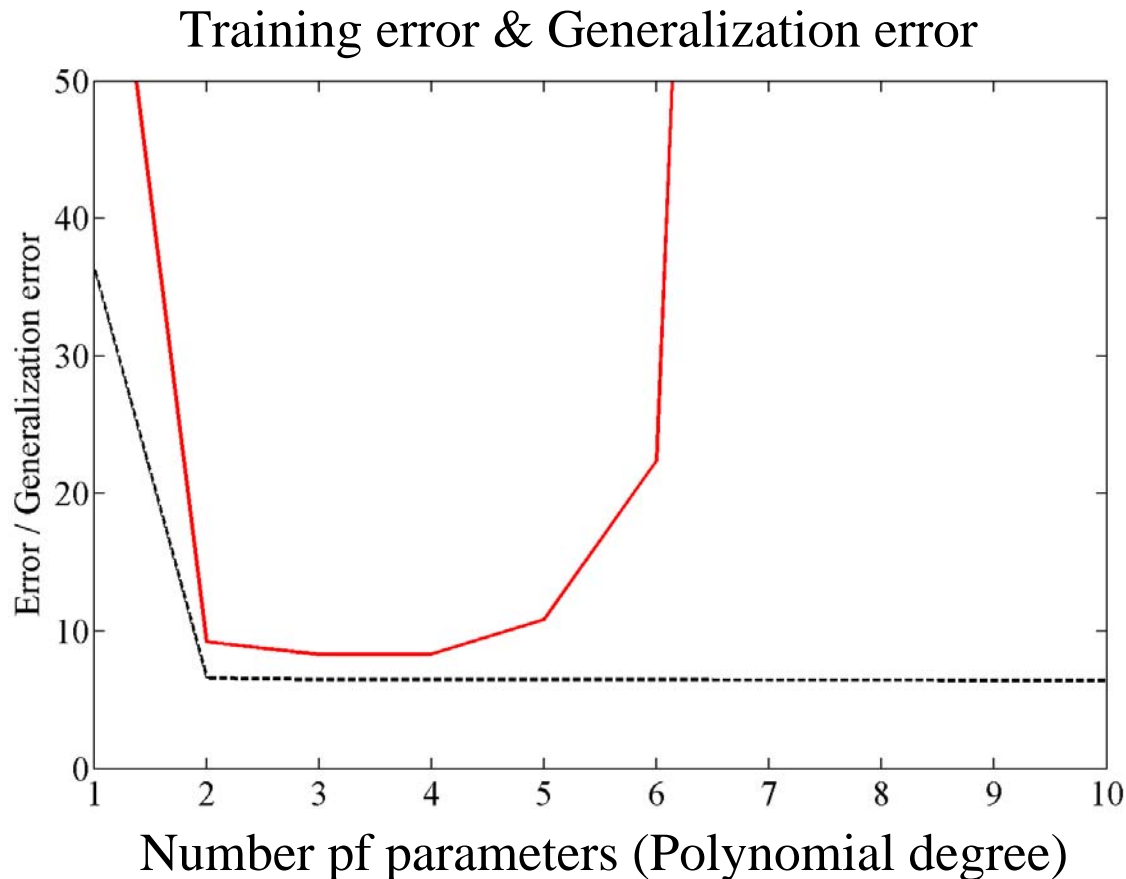
$$\text{error} = \sum_{t=1}^T (y(t) - W \cdot X(t))^2 + \alpha \cdot W^2 \quad \text{Minimization}$$

↔ Prior : $P(W) \propto \exp\left(-\frac{1}{2}\alpha \cdot W^2\right)$



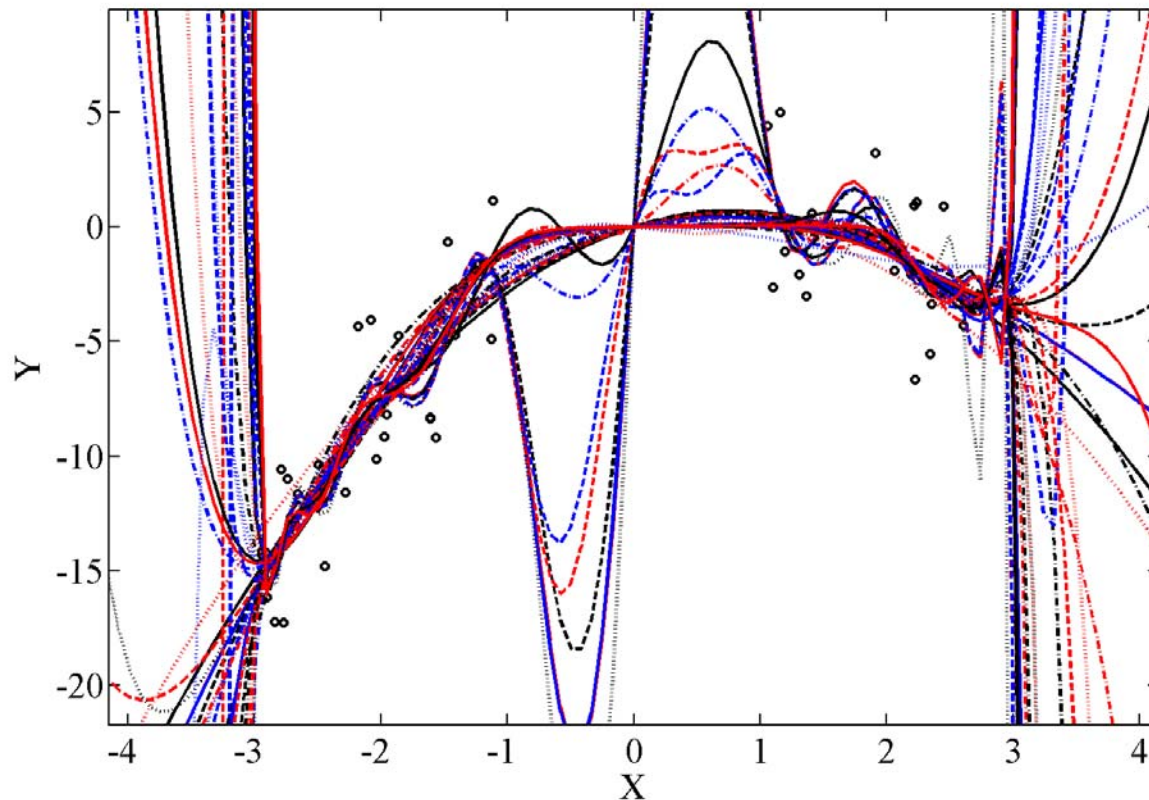
Model selection

- Search best model which gives best generalization error by changing number of parameters (polynomial degree)
- Combinatorial search is almost impossible for large degree of freedom



Sparse estimation by Bayesian method

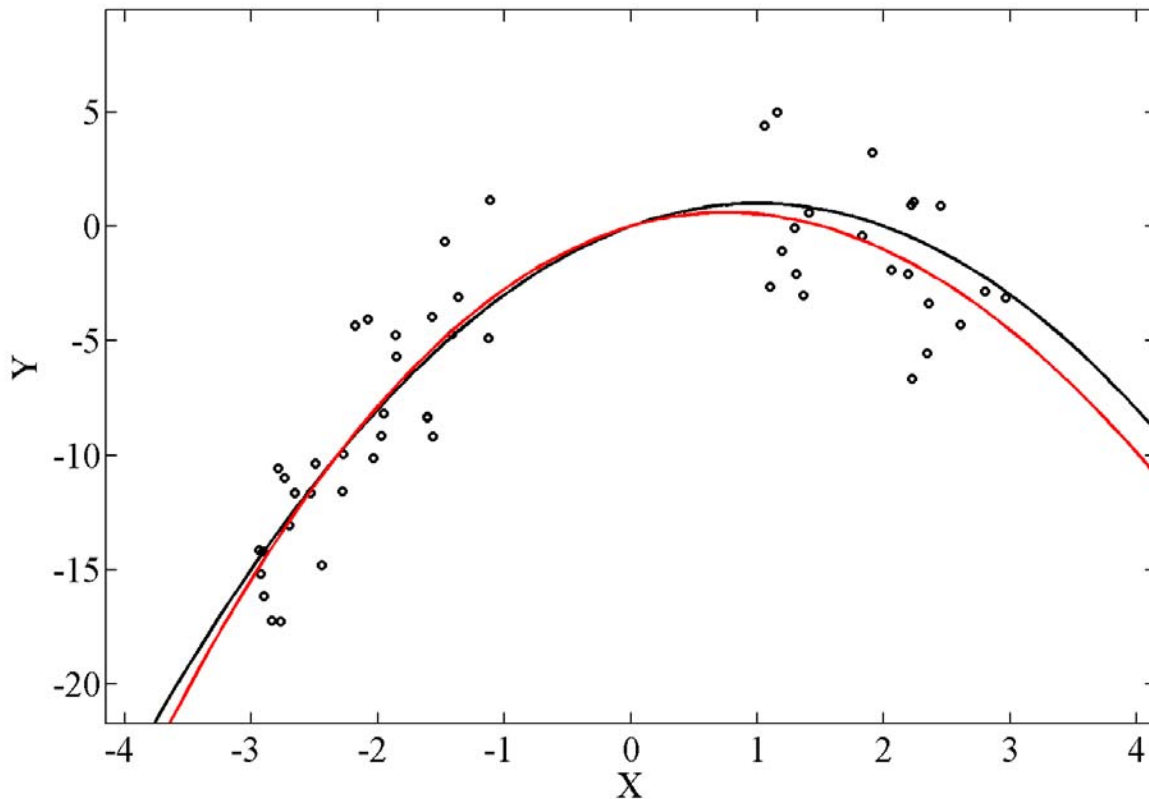
- Parameters are considered as random variable
- Posterior probability is calculated for possible parameter value
- Estimation is done by integrated over possible value according to posterior probability distribution



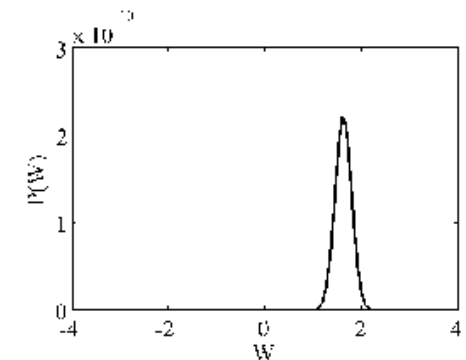
Sparse estimation (20 degree polynomial)

- Precision parameter α_n is introduced for each weight component w_n and estimated from observed data

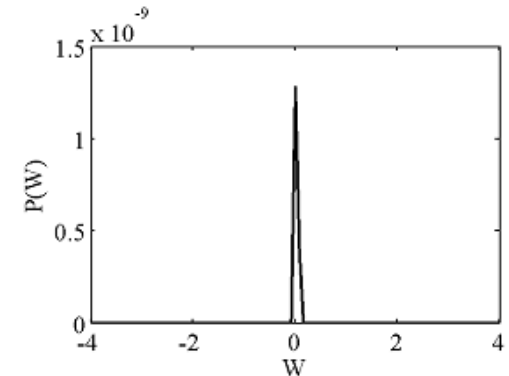
$$\text{Prior } P(w_n) \propto \exp\left(-\frac{1}{2}\alpha_n \cdot w_n^2\right), \quad \alpha_n = \text{不定}$$



Posterior for 2nd order weight



Posterior for 3rd order weight



Sparse Estimation

Prunes irrelevant features in the model
and increase generalization ability

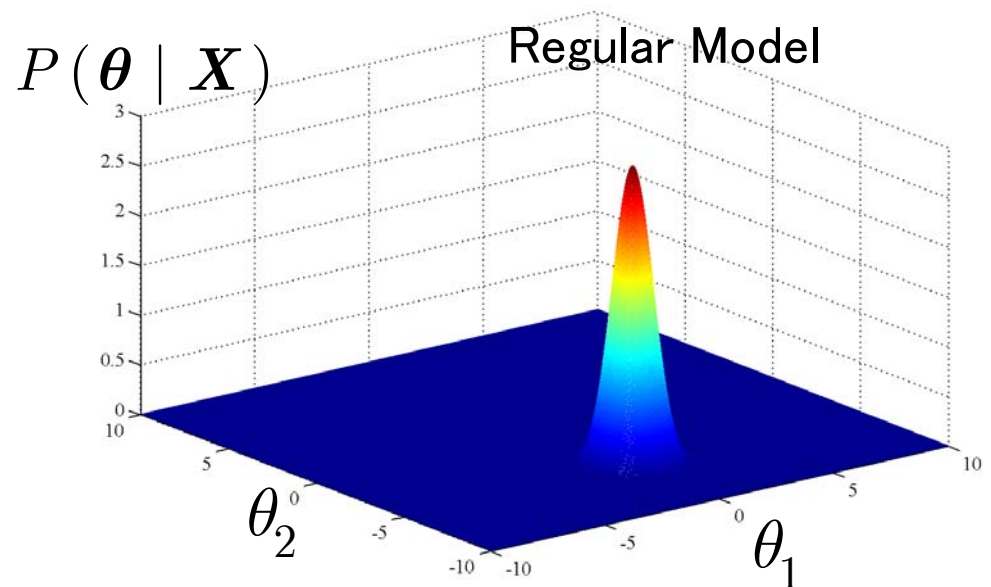
MAP / ML Estimation

(Maximum a Posteriori / Maximum Likelihood)

- Posterior $P(\boldsymbol{\theta} | \mathbf{X}) = \frac{\overset{\text{Likelihood}}{P(\mathbf{X} | \boldsymbol{\theta})} \overset{\text{prior}}{P_0(\boldsymbol{\theta})}}{\underset{\text{Marginal likelihood}}{P(\mathbf{X})}}$

- Estimate optimal parameter

$$\boldsymbol{\theta}_{\text{MAP}} = \arg \max \log(P(\mathbf{X} | \boldsymbol{\theta}) P_0(\boldsymbol{\theta}))$$



Full Bayesian Estimation

• Posterior $P(\boldsymbol{\theta} | \mathbf{X}) = \frac{\overset{\text{Likelihood}}{P(\mathbf{X} | \boldsymbol{\theta})} \overset{\text{prior}}{P_0(\boldsymbol{\theta})}}{\underset{\text{Marginal likelihood}}{P(\mathbf{X})}}$

- Estimate posterior parameter distribution and integrate over parameters according to the posterior.

Marginal likelihood $P(\mathbf{X}) = \int d\boldsymbol{\theta} P(\mathbf{X} | \boldsymbol{\theta}) P_0(\boldsymbol{\theta})$

Estimated parameter $\bar{\boldsymbol{\theta}} = \int d\boldsymbol{\theta} P(\boldsymbol{\theta} | \mathbf{X}) \boldsymbol{\theta}$

Model reduction by Bayesian method

Mixture of Gaussian example

Redundant Model (ill-posed problem)

Estimation model is a Mixture of two Gaussian units

$$P(\mathbf{x} | \boldsymbol{\theta}) = g_0 \cdot N(\mathbf{x} | \theta_1) + (1 - g_0) \cdot N(\mathbf{x} | \theta_2)$$

Assume data is generated by a single Gaussian

Single unit model correspond to three cases :

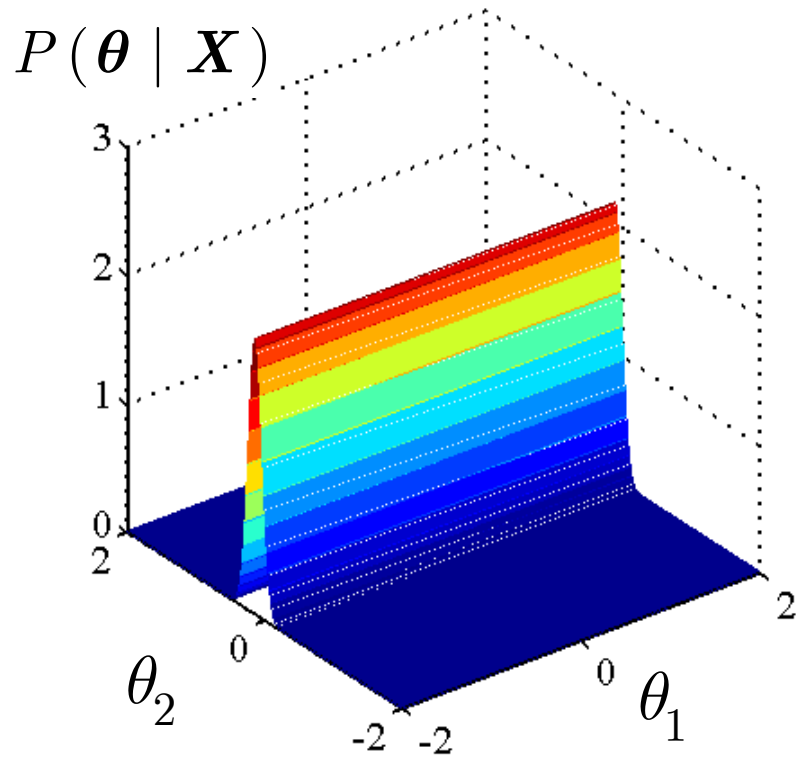
$$P(\mathbf{x} | \boldsymbol{\theta}) = N(\mathbf{x} | \theta_1), \quad g_0 = 1, \quad \theta_2 = \text{arbitrary}$$

$$P(\mathbf{x} | \boldsymbol{\theta}) = N(\mathbf{x} | \theta_2), \quad g_0 = 0, \quad \theta_1 = \text{arbitrary}$$

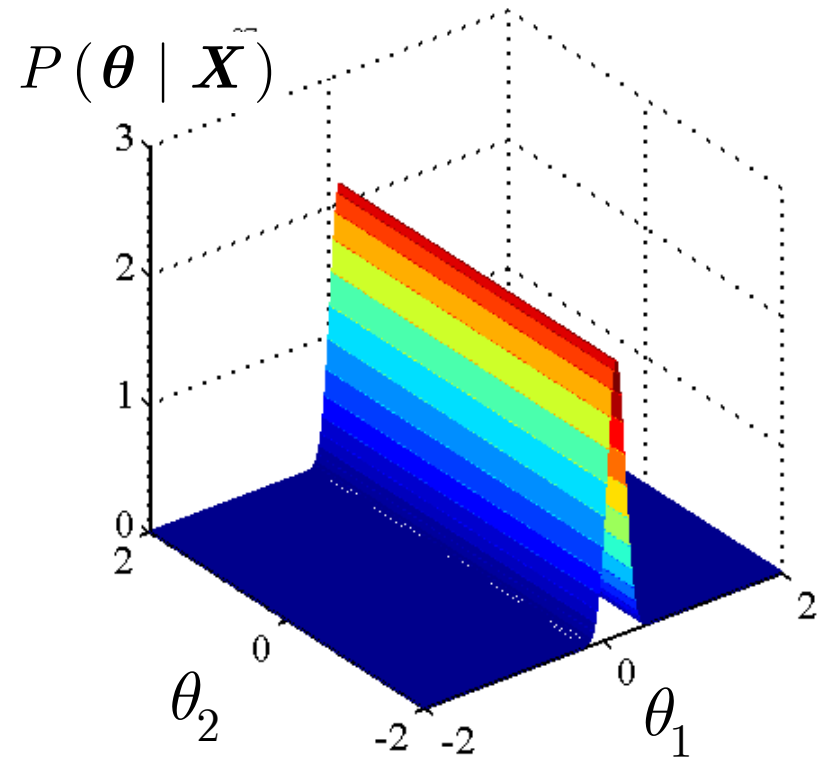
$$P(\mathbf{x} | \boldsymbol{\theta}) = N(\mathbf{x} | \theta_1), \quad \theta_1 = \theta_2, \quad g_0 = \text{arbitrary}$$

Posterior distribution for redundant model

$$g_0 = 1$$



$$g_0 = 0$$



Fisher Information matrix becomes singular

Pruning of redundant parameters

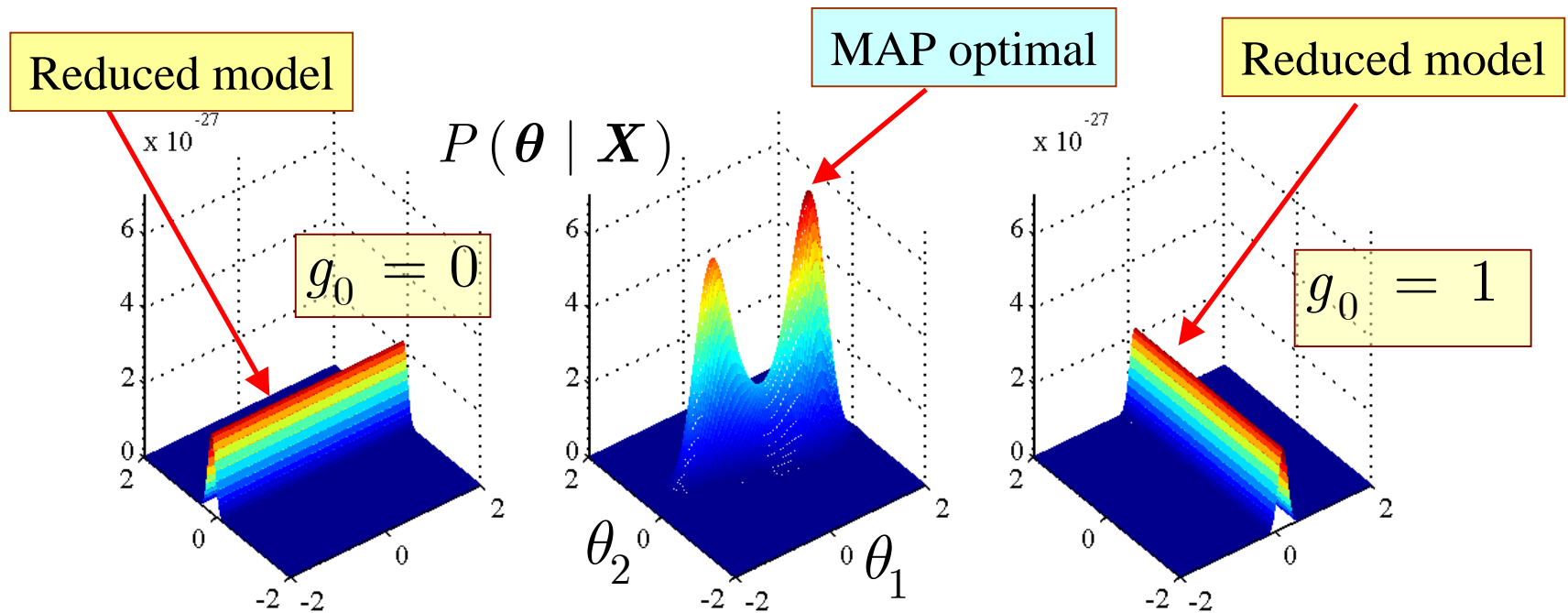
MAP

- Complex models explain a given data better than simpler models and give higher posterior value
- Then all parameters are used for prediction

Full Bayesian

- Reduced simpler model dominates by integration over parameters

Posterior distribution for 100 sample data generated by a single Gaussian model



Parameter pruning in Sparse Estimation

Prior in Sparse Estimation Model

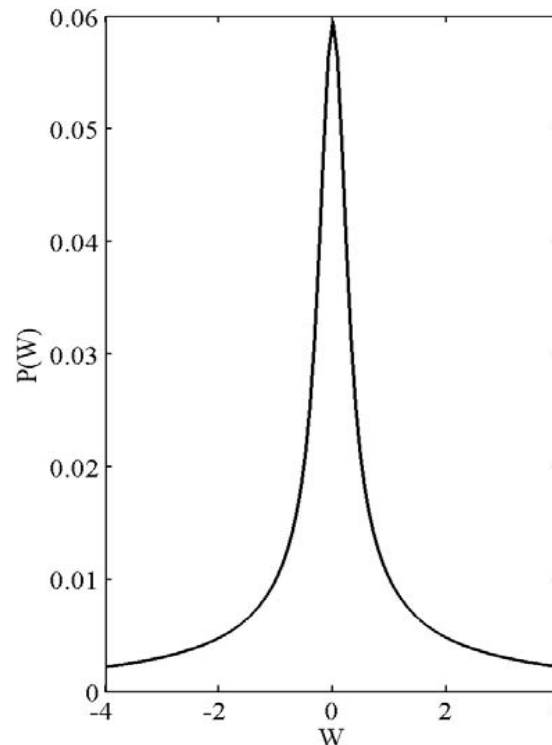
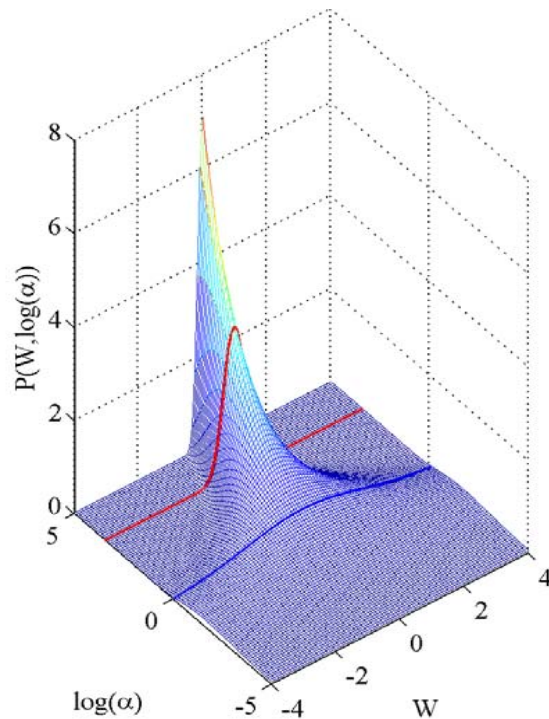
α_n controls a precision (width) of
weight parameter distribution

Prior

$$P(w_n | \alpha_n) \propto \sqrt{\alpha_n} \cdot \exp\left(-\frac{1}{2} \alpha_n \cdot w_n^2\right)$$

$$P(\log(\alpha_n)) = \text{const.} \quad (\text{Non-informative prior})$$

$$\langle w_n \rangle = 0, \quad \alpha_n = \text{Arbitrary}$$

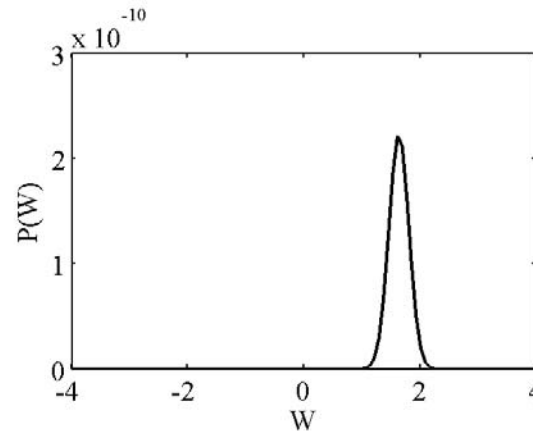
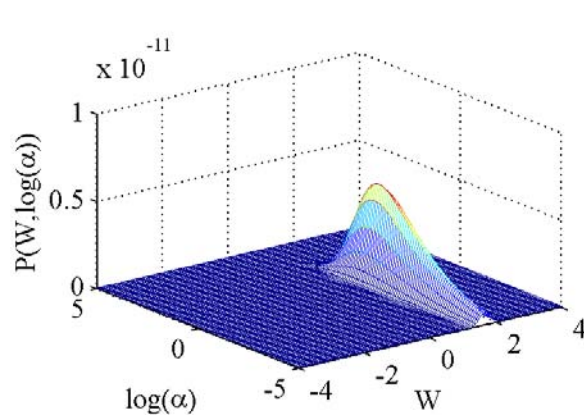


Posterior parameter distribution for (w_n, α_n)

Other parameters are integrated out

And their effects are taken into account

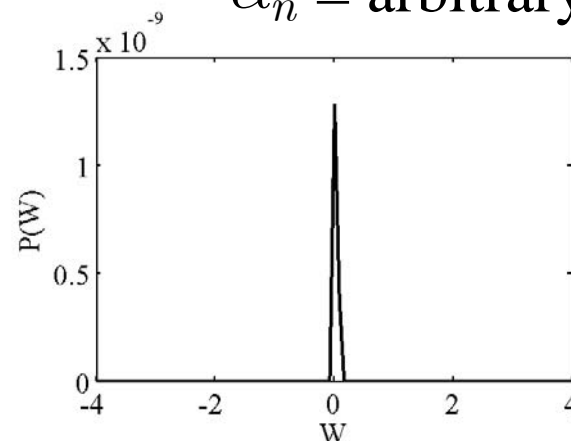
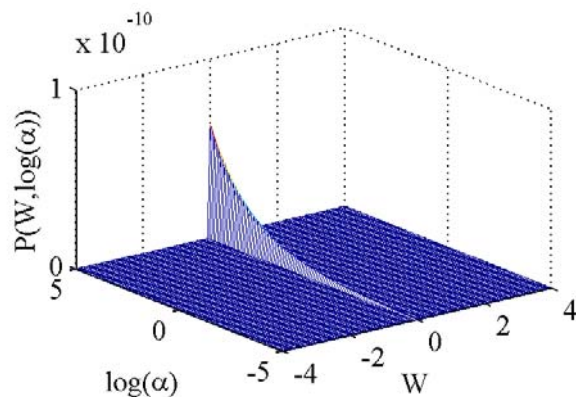
Posterior for relevant parameter



Posterior for irrelevant parameter

$$w_n = 0$$

$$\alpha_n = \text{arbitrary}$$



Calculation of posterior distribution

- Posterior** $P(\mathbf{J}, \alpha | \mathbf{B}) = \frac{\overset{\text{Likelihood}}{P(\mathbf{B} | \mathbf{J})} \overset{\text{(Hierarchical) prior}}{P_0(\mathbf{J} | \alpha) P_0(\alpha)}}{\underset{\text{Marginal likelihood}}{P(\mathbf{B})}}$
- $\langle \mathbf{J} \rangle = \int \mathbf{J} P(\mathbf{J}, \alpha | \mathbf{B}) d\mathbf{J} d\alpha$

Free energy maximization \leftrightarrow Posterior calculation

Distance between trial posterior $Q(\mathbf{J}, \alpha)$ and true posterior $P(\mathbf{J}, \alpha | \mathbf{B})$

$$F[Q(\mathbf{J}, \alpha)] = \ln P(\mathbf{B}) - KL[Q(\mathbf{J}, \alpha) || P(\mathbf{J}, \alpha | \mathbf{B})]$$

Log marginal likelihood (Evidence)

Variational Bayesian (VB) method

Factorization assumption : $Q(\mathbf{J}, \boldsymbol{\alpha}) = Q_{\mathbf{J}}(\mathbf{J})Q_{\boldsymbol{\alpha}}(\boldsymbol{\alpha})$

Maximization of $F(Q)$

Maximization of $F(Q)$ w.r.t. $Q_{\boldsymbol{\alpha}}(\boldsymbol{\alpha})$

Maximization of $F(Q)$ w.r.t. $Q_{\mathbf{J}}(\mathbf{J})$

Repeated until convergence

Posterior distribution $P(\mathbf{J} | \mathbf{B}) \approx Q_{\mathbf{J}}(\mathbf{J})$ at the maximum

Log marginal likelihood $\log(P(\mathbf{B})) \approx \text{maximized } F(Q)$

Free energy

Free energy after integration of current distribution

$$F = -\frac{1}{2} \left[\text{Tr} \left(\boldsymbol{\Sigma}_{VB}^{-1} \cdot \langle \mathbf{B} \cdot \mathbf{B}' \rangle \right) + \log |\boldsymbol{\Sigma}_{VB}| \right]$$

MEG covariance matrix $\langle \mathbf{B} \cdot \mathbf{B}' \rangle = \mathbf{G} \cdot \langle \mathbf{J}_0 \cdot \mathbf{J}_0' \rangle \cdot \mathbf{G}' + \sigma_0 \mathbf{I}$

Estimated MEG covariance $\boldsymbol{\Sigma}_{VB} = \mathbf{G} \cdot \mathbf{W} \cdot \boldsymbol{\alpha} \cdot \mathbf{W}' \cdot \mathbf{G}' + \sigma \mathbf{I}$

Optimal condition (Free energy maximum)

$$\alpha(n) = \frac{\langle \bar{\mathbf{J}} \cdot \bar{\mathbf{J}}' \rangle}{\mathcal{G}_{VB}(n)}$$

Estimation gain

Variational Bayesian method

Posterior calculation is converted to
free energy maximization

Free energy = (Likelihood) + (Model complexity)

Likelihood

$$\begin{aligned} L &= -\frac{1}{2} \left[(Y - W \cdot X)^2 + W' \cdot \alpha \cdot W \right] \\ &= -\frac{1}{2} \left[(Y)^2 - (Y \cdot X)^2 \cdot (X \cdot X' + \alpha)^{-1} \right] \end{aligned}$$

Error

Decreasing
function of α

Model complexity

$$\begin{aligned} H &= -\frac{1}{2} \log |X \cdot X' \cdot \alpha^{-1} + 1| \\ &\rightarrow -\frac{1}{2} N \cdot \log(T) \quad \text{for finite } \alpha, T \gg 1 \end{aligned}$$

Increasing
function of α

BIC

One dimensional case

$$H = -\frac{1}{2} \log |X \cdot X' \cdot \alpha^{-1} + 1|$$

$\rightarrow 0 \quad \text{as} \quad \alpha \rightarrow \infty$

General dimension

$$H = -\frac{1}{2} \log |X \cdot X' \cdot \alpha^{-1} + 1|$$

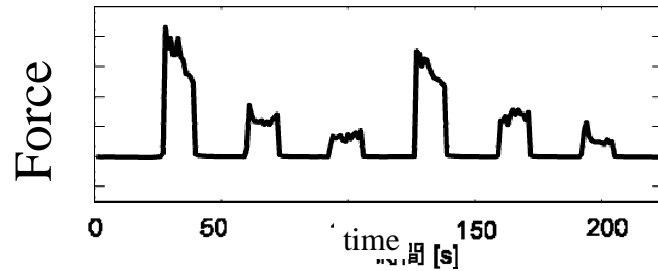
$\rightarrow -\frac{1}{2} N_{eff} \cdot \log(T)$

N_{eff} Number of finite α , $T \gg 1$

Example of sparse estimation for real problem

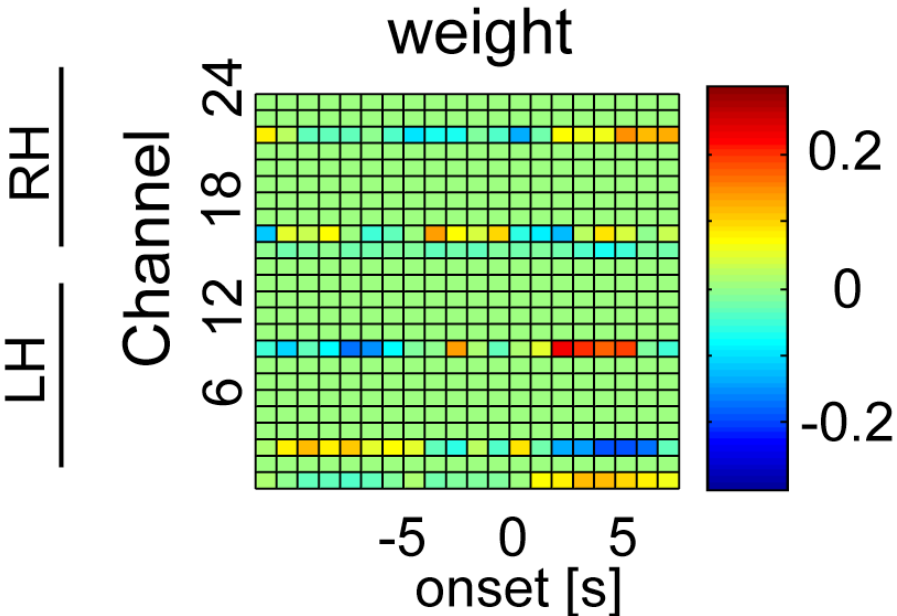
I Nambu, R Osu, M Sato, S Ando, M Kawato, E Naito
Single-trial reconstruction of finger-pinch forces from
human motor-cortical activation measured by near-
infrared spectroscopy (NIRS)
NeuroImage 47 (2009) 628.637

Sparse estimation

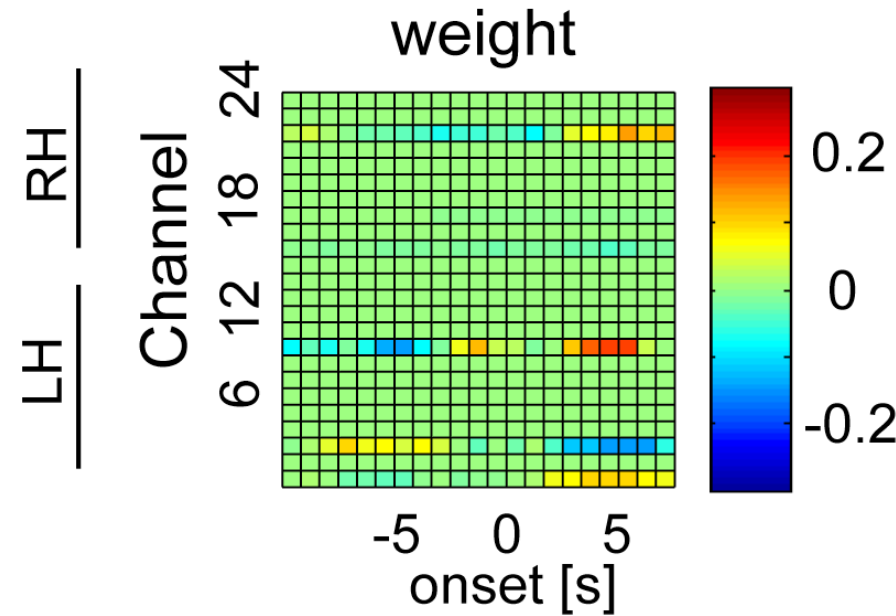


Estimate pinching force
from 24ch x 21 (sec) NIRS data
(Nambu et al)

Estimated weight
by brute force model search



Estimated weight
by sparse estimation



Estimated force from NIRS data

