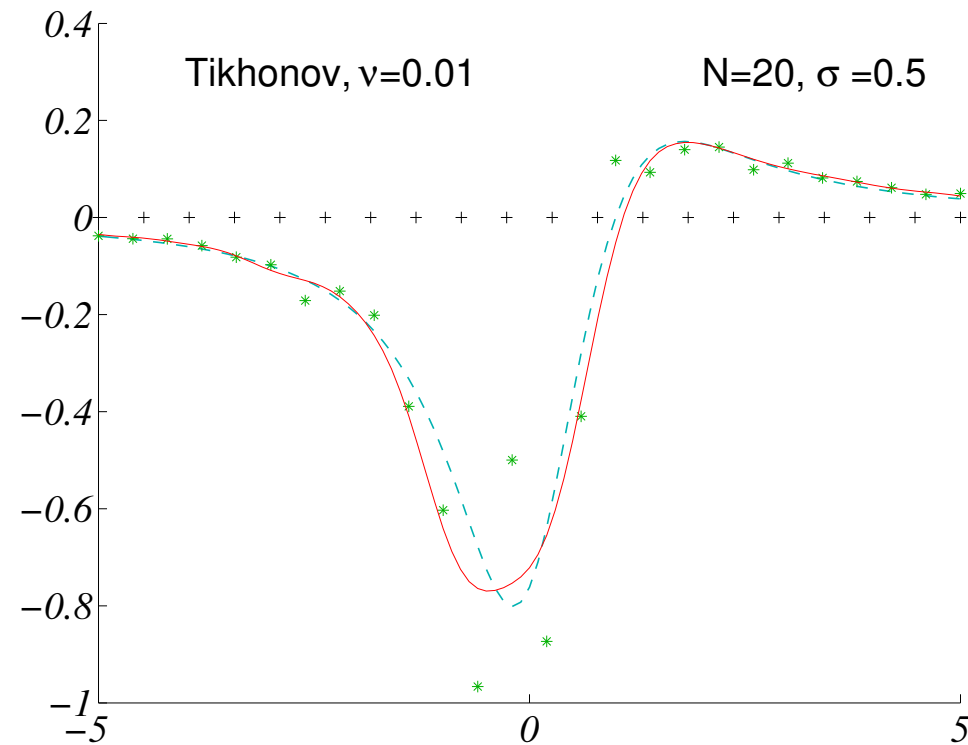


Advanced Machine Learning

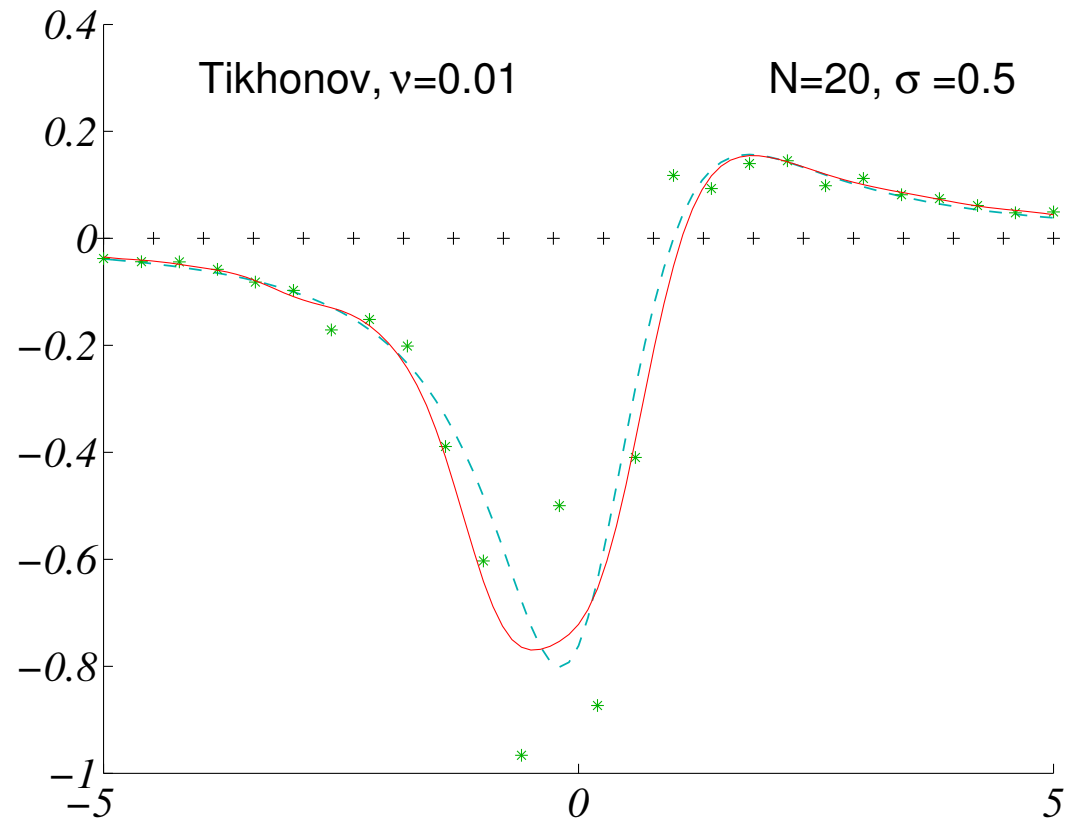
Regularisation



Regularisation, weight decay, Tikhonov

Outline

1. **Training Output Layer**
2. Regularisation
 - Weight Decay
 - Tikhonov



Back to Radial Basis Functions

- Recall in the last lecture that an RBF is a machine of the form

$$f(\mathbf{x}|\mathbf{w}, \{\boldsymbol{\mu}_i, \sigma_i\}_{i=1}^K) = \sum_{i=1}^K w_i \phi_i(\mathbf{x}) + w_0$$

- where $\phi_i(\mathbf{x}) = \psi\left(\left\|\frac{\mathbf{x}-\mathbf{c}_i}{\sigma_i}\right\|\right)$ are the radial basis functions
- Given data $\mathcal{D} = \{(\mathbf{x}_i, y_i)\}$, the weights in the output layer are learnt by linear least squares

$$E(\mathbf{w}|\mathcal{D}) = \frac{1}{P} \sum_{k=1}^P (\boldsymbol{\phi}(\mathbf{x}_k)^\top \mathbf{w} - y_k)^2 = \frac{1}{P} \|\boldsymbol{\Phi}^\top \mathbf{w} - \mathbf{y}\|^2$$

Back to Radial Basis Functions

- Recall in the last lecture that an RBF is a machine of the form

$$f(\mathbf{x}|\mathbf{w}, \{\boldsymbol{\mu}_i, \sigma_i\}_{i=1}^K) = \sum_{i=1}^K w_i \phi_i(\mathbf{x}) + w_0$$

- where $\phi_i(\mathbf{x}) = \psi\left(\left\|\frac{\mathbf{x}-\mathbf{c}_i}{\sigma_i}\right\|\right)$ are the radial basis functions
- Given data $\mathcal{D} = \{(\mathbf{x}_i, y_i)\}$, the weights in the output layer are learnt by linear least squares

$$E(\mathbf{w}|\mathcal{D}) = \frac{1}{P} \sum_{k=1}^P (\boldsymbol{\phi}(\mathbf{x}_k)^\top \mathbf{w} - y_k)^2 = \frac{1}{P} \|\boldsymbol{\Phi}^\top \mathbf{w} - \mathbf{y}\|^2$$

Back to Radial Basis Functions

- Recall in the last lecture that an RBF is a machine of the form

$$f(\mathbf{x}|\mathbf{w}, \{\boldsymbol{\mu}_i, \sigma_i\}_{i=1}^K) = \sum_{i=1}^K w_i \phi_i(\mathbf{x}) + w_0$$

- where $\phi_i(\mathbf{x}) = \psi\left(\left\|\frac{\mathbf{x}-\mathbf{c}_i}{\sigma_i}\right\|\right)$ are the radial basis functions
- Given data $\mathcal{D} = \{(\mathbf{x}_i, y_i)\}$, the weights in the output layer are learnt by linear least squares

$$E(\mathbf{w}|\mathcal{D}) = \frac{1}{P} \sum_{k=1}^P (\boldsymbol{\phi}(\mathbf{x}_k)^\top \mathbf{w} - y_k)^2 = \frac{1}{P} \|\boldsymbol{\Phi}^\top \mathbf{w} - \mathbf{y}\|^2$$

Linear Least Squares

- Same as linear perceptron
- Output $\phi(\mathbf{x})^\top \mathbf{w}$ where $\phi_k(\mathbf{x}) = \phi\left(\frac{\|\mathbf{x} - \boldsymbol{\mu}_k\|}{\sigma_k}\right)$
- Mean squared error

$$E = \frac{1}{P} \sum_{k=1}^P (\phi(\mathbf{x}_k)^\top \mathbf{w} - y_k)^2 = \frac{1}{P} \|\Phi^\top \mathbf{w} - \mathbf{y}\|^2$$

where

$$\Phi = \begin{pmatrix} \phi_1(\mathbf{x}_1) & \phi_1(\mathbf{x}_2) & \cdots & \phi_1(\mathbf{x}_P) \\ \phi_2(\mathbf{x}_1) & \phi_2(\mathbf{x}_2) & \cdots & \phi_2(\mathbf{x}_P) \\ \vdots & \vdots & \ddots & \vdots \\ \phi_N(\mathbf{x}_1) & \phi_N(\mathbf{x}_2) & \cdots & \phi_N(\mathbf{x}_P) \end{pmatrix} \quad \mathbf{y} = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_P \end{pmatrix}$$

Linear Least Squares

- Same as linear perceptron
- Output $\phi(\mathbf{x})^\top \mathbf{w}$ where $\phi_k(\mathbf{x}) = \phi\left(\frac{\|\mathbf{x} - \boldsymbol{\mu}_k\|}{\sigma_k}\right)$
- Mean squared error

$$E = \frac{1}{P} \sum_{k=1}^P (\phi(\mathbf{x}_k)^\top \mathbf{w} - y_k)^2 = \frac{1}{P} \|\Phi^\top \mathbf{w} - \mathbf{y}\|^2$$

where

$$\Phi = \begin{pmatrix} \phi_1(\mathbf{x}_1) & \phi_1(\mathbf{x}_2) & \cdots & \phi_1(\mathbf{x}_P) \\ \phi_2(\mathbf{x}_1) & \phi_2(\mathbf{x}_2) & \cdots & \phi_2(\mathbf{x}_P) \\ \vdots & \vdots & \ddots & \vdots \\ \phi_N(\mathbf{x}_1) & \phi_N(\mathbf{x}_2) & \cdots & \phi_N(\mathbf{x}_P) \end{pmatrix} \quad \mathbf{y} = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_P \end{pmatrix}$$

Linear Least Squares

- Same as linear perceptron
- Output $\phi(\mathbf{x})^\top \mathbf{w}$ where $\phi_k(\mathbf{x}) = \phi\left(\frac{\|\mathbf{x} - \boldsymbol{\mu}_k\|}{\sigma_k}\right)$
- Mean squared error

$$E = \frac{1}{P} \sum_{k=1}^P (\phi(\mathbf{x}_k)^\top \mathbf{w} - y_k)^2 = \frac{1}{P} \|\Phi^\top \mathbf{w} - \mathbf{y}\|^2$$

where

$$\Phi = \begin{pmatrix} \phi_1(\mathbf{x}_1) & \phi_1(\mathbf{x}_2) & \cdots & \phi_1(\mathbf{x}_P) \\ \phi_2(\mathbf{x}_1) & \phi_2(\mathbf{x}_2) & \cdots & \phi_2(\mathbf{x}_P) \\ \vdots & \vdots & \ddots & \vdots \\ \phi_N(\mathbf{x}_1) & \phi_N(\mathbf{x}_2) & \cdots & \phi_N(\mathbf{x}_P) \end{pmatrix} \quad \mathbf{y} = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_P \end{pmatrix}$$

Linear Least Squares

- Same as linear perceptron
- Output $\phi(\mathbf{x})^\top \mathbf{w}$ where $\phi_k(\mathbf{x}) = \phi\left(\frac{\|\mathbf{x} - \boldsymbol{\mu}_k\|}{\sigma_k}\right)$
- Mean squared error

$$E = \frac{1}{P} \sum_{k=1}^P (\phi(\mathbf{x}_k)^\top \mathbf{w} - y_k)^2 = \frac{1}{P} \|\Phi^\top \mathbf{w} - \mathbf{y}\|^2$$

where

$$\Phi = \begin{pmatrix} \phi_1(\mathbf{x}_1) & \phi_1(\mathbf{x}_2) & \cdots & \phi_1(\mathbf{x}_P) \\ \phi_2(\mathbf{x}_1) & \phi_2(\mathbf{x}_2) & \cdots & \phi_2(\mathbf{x}_P) \\ \vdots & \vdots & \ddots & \vdots \\ \phi_N(\mathbf{x}_1) & \phi_N(\mathbf{x}_2) & \cdots & \phi_N(\mathbf{x}_P) \end{pmatrix} \quad \mathbf{y} = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_P \end{pmatrix}$$

Linear Least Squares

- Same as linear perceptron
- Output $\phi(\mathbf{x})^\top \mathbf{w}$ where $\phi_k(\mathbf{x}) = \phi\left(\frac{\|\mathbf{x} - \boldsymbol{\mu}_k\|}{\sigma_k}\right)$
- Mean squared error

$$E = \frac{1}{P} \sum_{k=1}^P (\phi(\mathbf{x}_k)^\top \mathbf{w} - y_k)^2 = \frac{1}{P} \|\Phi^\top \mathbf{w} - \mathbf{y}\|^2$$

where

$$\Phi = \begin{pmatrix} \phi_1(\mathbf{x}_1) & \phi_1(\mathbf{x}_2) & \cdots & \phi_1(\mathbf{x}_P) \\ \phi_2(\mathbf{x}_1) & \phi_2(\mathbf{x}_2) & \cdots & \phi_2(\mathbf{x}_P) \\ \vdots & \vdots & \ddots & \vdots \\ \phi_N(\mathbf{x}_1) & \phi_N(\mathbf{x}_2) & \cdots & \phi_N(\mathbf{x}_P) \end{pmatrix} \quad \mathbf{y} = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_P \end{pmatrix}$$

Least Squares Solution

- In matrix form

$$E = \frac{1}{P} (w^T \Phi \Phi^T w - 2w^T \Phi y + y^T y)$$

- Minimum given by $\nabla E = 0$

$$\nabla E = \frac{2}{P} (\Phi \Phi^T w - \Phi y) = 0$$

- Optimal weight vector

$$w^* = (\Phi \Phi^T)^{-1} \Phi y$$

Least Squares Solution

- In matrix form

$$E = \frac{1}{P} (w^\top \Phi \Phi^\top w - 2w^\top \Phi y + y^\top y)$$

- Minimum given by $\nabla E = 0$

$$\nabla E = \frac{2}{P} (\Phi \Phi^\top w - \Phi y) = 0$$

- Optimal weight vector

$$w^* = (\Phi \Phi^\top)^{-1} \Phi y$$

Least Squares Solution

- In matrix form

$$E = \frac{1}{P} (w^\top \Phi \Phi^\top w - 2w^\top \Phi y + y^\top y)$$

- Minimum given by $\nabla E = 0$

$$\nabla E = \frac{2}{P} (\Phi \Phi^\top w - \Phi y) = 0$$

- Optimal weight vector

$$w^* = (\Phi \Phi^\top)^{-1} \Phi y$$

Singular Value Decomposition

- Remember the singular valued decomposition formula

$$\Phi = \mathbf{U} \mathbf{S} \mathbf{V}^T$$

where

$$\mathbf{S} = \begin{pmatrix} s_1 & \cdots & 0 & \cdots & 0 \\ \vdots & \ddots & \vdots & \ddots & \vdots \\ 0 & \cdots & s_N & \cdots & 0 \end{pmatrix}$$

- \mathbf{U} and \mathbf{V} are orthogonal matrices, i.e. $\mathbf{U}^T = \mathbf{U}^{-1}$ and $\mathbf{V}^T = \mathbf{V}^{-1}$

Singular Value Decomposition

- Remember the singular valued decomposition formula

$$\Phi = \mathbf{U} \mathbf{S} \mathbf{V}^T$$

where

$$\mathbf{S} = \begin{pmatrix} s_1 & \cdots & 0 & \cdots & 0 \\ \vdots & \ddots & \vdots & \ddots & \vdots \\ 0 & \cdots & s_N & \cdots & 0 \end{pmatrix}$$

- \mathbf{U} and \mathbf{V} are orthogonal matrices, i.e. $\mathbf{U}^T = \mathbf{U}^{-1}$ and $\mathbf{V}^T = \mathbf{V}^{-1}$

Singular Value Decomposition

- Remember the singular valued decomposition formula

$$\Phi = \mathbf{U} \mathbf{S} \mathbf{V}^T$$

where

$$\mathbf{S} = \begin{pmatrix} s_1 & \cdots & 0 & \cdots & 0 \\ \vdots & \ddots & \vdots & \ddots & \vdots \\ 0 & \cdots & s_N & \cdots & 0 \end{pmatrix}$$

- \mathbf{U} and \mathbf{V} are orthogonal matrices, i.e. $\mathbf{U}^T = \mathbf{U}^{-1}$ and $\mathbf{V}^T = \mathbf{V}^{-1}$

$\Phi \Phi^T$ Matrix

- Now the 'covariance-like matrix' is equal to

$$\begin{aligned}\Phi \Phi^T &= (\mathbf{U} \mathbf{S} \mathbf{V}^T) (\mathbf{U} \mathbf{S} \mathbf{V}^T)^T \\ &= (\mathbf{U} \mathbf{S} \mathbf{V}^T) (\mathbf{V} \mathbf{S}^T \mathbf{U}^T) \\ &= \mathbf{U} (\mathbf{S} \mathbf{S}^T) \mathbf{U}^T = \mathbf{U} \mathbf{\Lambda} \mathbf{U}^T\end{aligned}$$

- $\mathbf{\Lambda} = \mathbf{S} \mathbf{S}^T = \text{diag}(s_1^2, s_2^2, \dots, s_N^2) = \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_N)$
- $(\Phi \Phi^T)$ is a positive definite matrix with eigenvalues $\lambda_i = s_i^2$
- The inverse covariance matrix is

$$(\Phi \Phi^T)^{-1} = \mathbf{U} \mathbf{\Lambda}^{-1} \mathbf{U}^T$$

- Since $(\Phi \Phi^T)^{-1}(\Phi \Phi^T) = \mathbf{U} \mathbf{\Lambda}^{-1} \mathbf{U}^T \mathbf{U} \mathbf{\Lambda} \mathbf{U}^T = \mathbf{I}$

$\Phi \Phi^T$ Matrix

- Now the 'covariance-like matrix' is equal to

$$\begin{aligned}\Phi \Phi^T &= (\mathbf{U} \mathbf{S} \mathbf{V}^T) (\mathbf{U} \mathbf{S} \mathbf{V}^T)^T \\ &= (\mathbf{U} \mathbf{S} \mathbf{V}^T) (\mathbf{V} \mathbf{S}^T \mathbf{U}^T) \\ &= \mathbf{U} (\mathbf{S} \mathbf{S}^T) \mathbf{U}^T = \mathbf{U} \mathbf{\Lambda} \mathbf{U}^T\end{aligned}$$

- $\mathbf{\Lambda} = \mathbf{S} \mathbf{S}^T = \text{diag}(s_1^2, s_2^2, \dots, s_N^2) = \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_N)$
- $(\Phi \Phi^T)$ is a positive definite matrix with eigenvalues $\lambda_i = s_i^2$
- The inverse covariance matrix is

$$(\Phi \Phi^T)^{-1} = \mathbf{U} \mathbf{\Lambda}^{-1} \mathbf{U}^T$$

- Since $(\Phi \Phi^T)^{-1}(\Phi \Phi^T) = \mathbf{U} \mathbf{\Lambda}^{-1} \mathbf{U}^T \mathbf{U} \mathbf{\Lambda} \mathbf{U}^T = \mathbf{I}$

$\Phi \Phi^T$ Matrix

- Now the 'covariance-like matrix' is equal to

$$\begin{aligned}\Phi \Phi^T &= (\mathbf{U} \mathbf{S} \mathbf{V}^T) (\mathbf{U} \mathbf{S} \mathbf{V}^T)^T \\ &= (\mathbf{U} \mathbf{S} \mathbf{V}^T) (\mathbf{V} \mathbf{S}^T \mathbf{U}^T) \\ &= \mathbf{U} (\mathbf{S} \mathbf{S}^T) \mathbf{U}^T = \mathbf{U} \mathbf{\Lambda} \mathbf{U}^T\end{aligned}$$

- $\mathbf{\Lambda} = \mathbf{S} \mathbf{S}^T = \text{diag}(s_1^2, s_2^2, \dots, s_N^2) = \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_N)$
- $(\Phi \Phi^T)$ is a positive definite matrix with eigenvalues $\lambda_i = s_i^2$
- The inverse covariance matrix is

$$(\Phi \Phi^T)^{-1} = \mathbf{U} \mathbf{\Lambda}^{-1} \mathbf{U}^T$$

- Since $(\Phi \Phi^T)^{-1}(\Phi \Phi^T) = \mathbf{U} \mathbf{\Lambda}^{-1} \mathbf{U}^T \mathbf{U} \mathbf{\Lambda} \mathbf{U}^T = \mathbf{I}$

$\Phi \Phi^T$ Matrix

- Now the 'covariance-like matrix' is equal to

$$\begin{aligned}\Phi \Phi^T &= (\mathbf{U} \mathbf{S} \mathbf{V}^T) (\mathbf{U} \mathbf{S} \mathbf{V}^T)^T \\ &= (\mathbf{U} \mathbf{S} \mathbf{V}^T) (\mathbf{V} \mathbf{S}^T \mathbf{U}^T) \\ &= \mathbf{U} (\mathbf{S} \mathbf{S}^T) \mathbf{U}^T = \mathbf{U} \mathbf{\Lambda} \mathbf{U}^T\end{aligned}$$

- $\mathbf{\Lambda} = \mathbf{S} \mathbf{S}^T = \text{diag}(s_1^2, s_2^2, \dots, s_N^2) = \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_N)$
- $(\Phi \Phi^T)$ is a positive definite matrix with eigenvalues $\lambda_i = s_i^2$
- The inverse covariance matrix is

$$(\Phi \Phi^T)^{-1} = \mathbf{U} \mathbf{\Lambda}^{-1} \mathbf{U}^T$$

- Since $(\Phi \Phi^T)^{-1}(\Phi \Phi^T) = \mathbf{U} \mathbf{\Lambda}^{-1} \mathbf{U}^T \mathbf{U} \mathbf{\Lambda} \mathbf{U}^T = \mathbf{I}$

$\Phi \Phi^T$ Matrix

- Now the 'covariance-like matrix' is equal to

$$\begin{aligned}\Phi \Phi^T &= (\mathbf{U} \mathbf{S} \mathbf{V}^T) (\mathbf{U} \mathbf{S} \mathbf{V}^T)^T \\ &= (\mathbf{U} \mathbf{S} \mathbf{V}^T) (\mathbf{V} \mathbf{S}^T \mathbf{U}^T) \\ &= \mathbf{U} (\mathbf{S} \mathbf{S}^T) \mathbf{U}^T = \mathbf{U} \mathbf{\Lambda} \mathbf{U}^T\end{aligned}$$

- $\mathbf{\Lambda} = \mathbf{S} \mathbf{S}^T = \text{diag}(s_1^2, s_2^2, \dots, s_N^2) = \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_N)$
- $(\Phi \Phi^T)$ is a positive definite matrix with eigenvalues $\lambda_i = s_i^2$
- The inverse covariance matrix is

$$(\Phi \Phi^T)^{-1} = \mathbf{U} \mathbf{\Lambda}^{-1} \mathbf{U}^T$$

- Since $(\Phi \Phi^T)^{-1}(\Phi \Phi^T) = \mathbf{U} \mathbf{\Lambda}^{-1} \mathbf{U}^T \mathbf{U} \mathbf{\Lambda} \mathbf{U}^T = \mathbf{I}$

Least-Squares Solution

- The inverse covariance matrix is

$$(\Phi \Phi^T)^{-1} = \mathbf{U} (\mathbf{S} \mathbf{S}^T)^{-1} \mathbf{U}^T$$

- where $(\mathbf{S} \mathbf{S}^T)^{-1} = \text{diag}(s_1^{-2}, s_2^{-2}, \dots, s_N^{-2})$
- Thus

$$\begin{aligned} (\Phi \Phi^T)^{-1} \Phi &= \left(\mathbf{U} (\mathbf{S} \mathbf{S}^T)^{-1} \mathbf{U}^T \right) (\mathbf{U} \mathbf{S} \mathbf{V}^T) \\ &= \mathbf{U} (\mathbf{S} \mathbf{S}^T)^{-1} \mathbf{S} \mathbf{V}^T \\ &= \mathbf{U} \mathbf{S}^+ \mathbf{V}^T \end{aligned}$$

$$\text{where } \mathbf{S}^+ = (\mathbf{S} \mathbf{S}^T)^{-1} \mathbf{S}$$

Least-Squares Solution

- The inverse covariance matrix is

$$(\Phi \Phi^T)^{-1} = \mathbf{U} (\mathbf{S} \mathbf{S}^T)^{-1} \mathbf{U}^T$$

- where $(\mathbf{S} \mathbf{S}^T)^{-1} = \text{diag}(s_1^{-2}, s_2^{-2}, \dots, s_N^{-2})$

- Thus

$$\begin{aligned} (\Phi \Phi^T)^{-1} \Phi &= \left(\mathbf{U} (\mathbf{S} \mathbf{S}^T)^{-1} \mathbf{U}^T \right) (\mathbf{U} \mathbf{S} \mathbf{V}^T) \\ &= \mathbf{U} (\mathbf{S} \mathbf{S}^T)^{-1} \mathbf{S} \mathbf{V}^T \\ &= \mathbf{U} \mathbf{S}^+ \mathbf{V}^T \end{aligned}$$

$$\text{where } \mathbf{S}^+ = (\mathbf{S} \mathbf{S}^T)^{-1} \mathbf{S}$$

Pseudo-Inverse

- The **pseudo-inverse** is defined as $\Phi^+ = \mathbf{U} \mathbf{S}^+ \mathbf{V}^T$ where

$$\mathbf{S}^+ = (\mathbf{S} \mathbf{S}^T)^{-1} \mathbf{S} = \begin{pmatrix} s_1^{-1} & \dots & 0 & \dots & 0 \\ \vdots & \ddots & \vdots & \ddots & \vdots \\ 0 & \dots & s_N^{-1} & \dots & 0 \end{pmatrix}$$

- Thus the solution to the least squares problem is

$$\mathbf{w}^* = (\Phi \Phi^T)^{-1} \Phi \mathbf{y} = \Phi^+ \mathbf{y} = \mathbf{U} \mathbf{S}^+ \mathbf{V}^T \mathbf{y}$$

- For non-square matrices Matlab uses the pseudo-inverse so in Matlab we can write

```
w = Psi\t
```

Pseudo-Inverse

- The **pseudo-inverse** is defined as $\Phi^+ = \mathbf{U} \mathbf{S}^+ \mathbf{V}^T$ where

$$\mathbf{S}^+ = (\mathbf{S} \mathbf{S}^T)^{-1} \mathbf{S} = \begin{pmatrix} s_1^{-1} & \cdots & 0 & \cdots & 0 \\ \vdots & \ddots & \vdots & \ddots & \vdots \\ 0 & \cdots & s_N^{-1} & \cdots & 0 \end{pmatrix}$$

- Thus the solution to the least squares problem is

$$\mathbf{w}^* = (\Phi \Phi^T)^{-1} \Phi \mathbf{y} = \Phi^+ \mathbf{y} = \mathbf{U} \mathbf{S}^+ \mathbf{V}^T \mathbf{y}$$

- For non-square matrices Matlab uses the pseudo-inverse so in Matlab we can write

```
w = Psi\t
```

Pseudo-Inverse

- The **pseudo-inverse** is defined as $\Phi^+ = \mathbf{U} \mathbf{S}^+ \mathbf{V}^T$ where

$$\mathbf{S}^+ = (\mathbf{S} \mathbf{S}^T)^{-1} \mathbf{S} = \begin{pmatrix} s_1^{-1} & \dots & 0 & \dots & 0 \\ \vdots & \ddots & \vdots & \ddots & \vdots \\ 0 & \dots & s_N^{-1} & \dots & 0 \end{pmatrix}$$

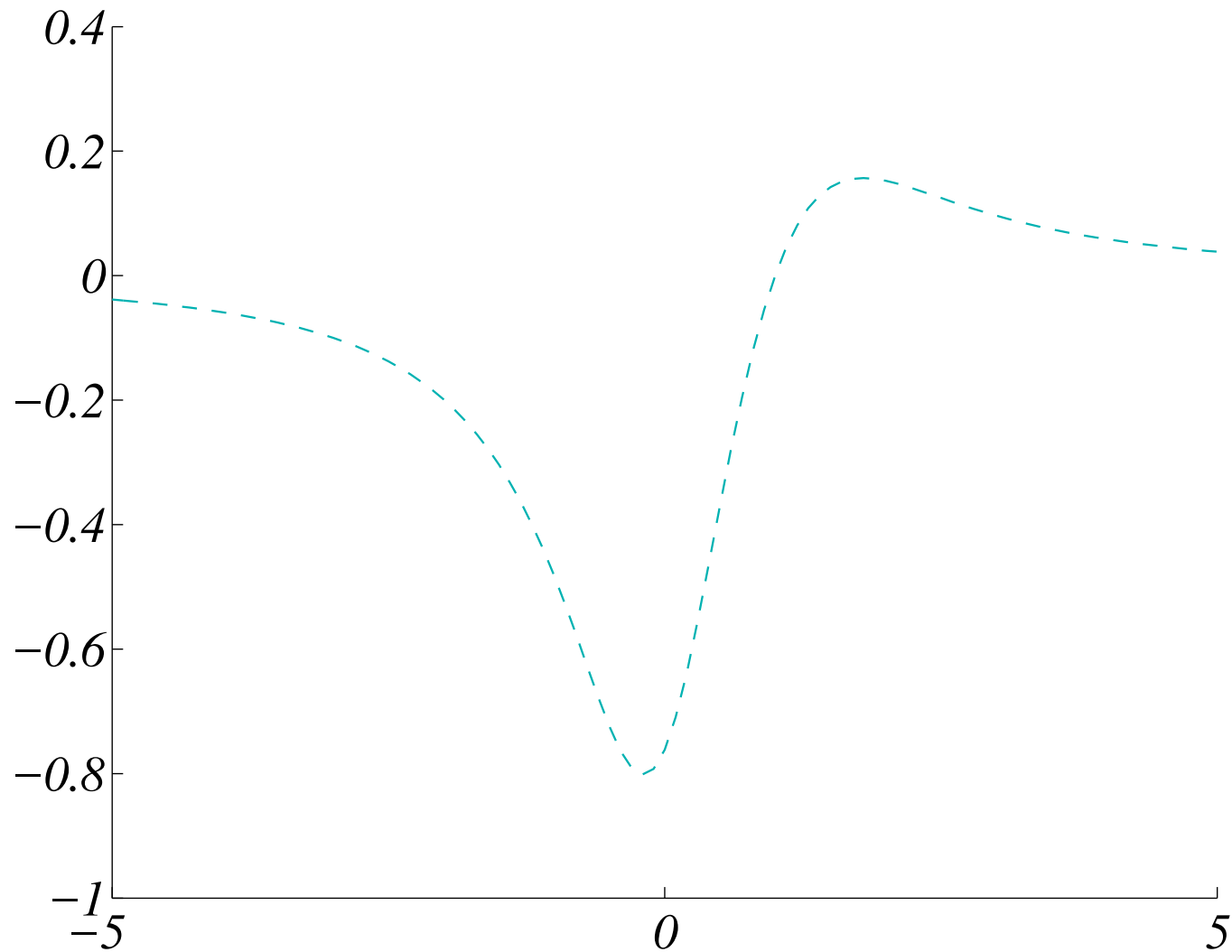
- Thus the solution to the least squares problem is

$$\mathbf{w}^* = (\Phi \Phi^T)^{-1} \Phi \mathbf{y} = \Phi^+ \mathbf{y} = \mathbf{U} \mathbf{S}^+ \mathbf{V}^T \mathbf{y}$$

- For non-square matrices Matlab uses the pseudo-inverse so in Matlab we can write

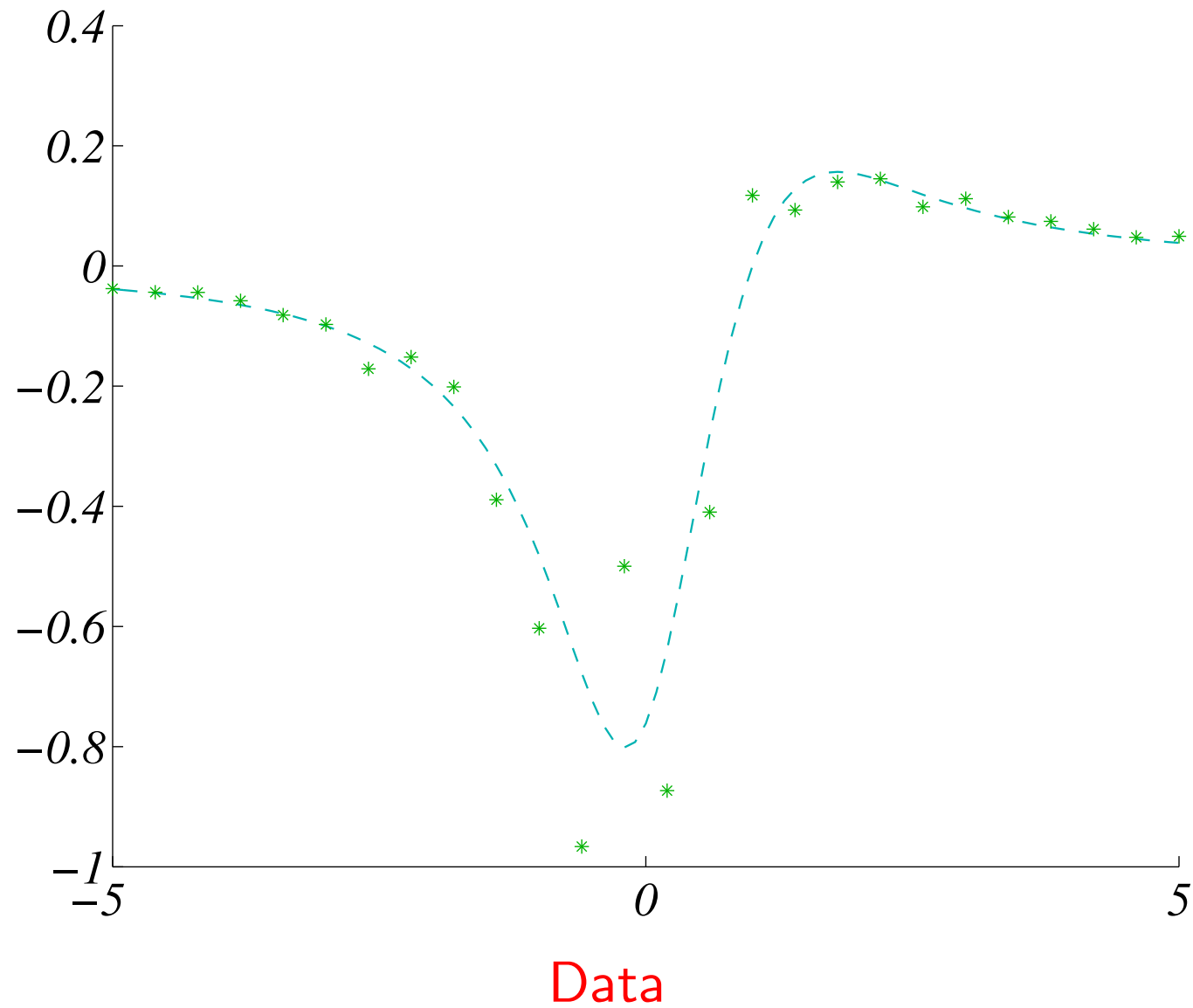
```
w = Psi\t
```

Example

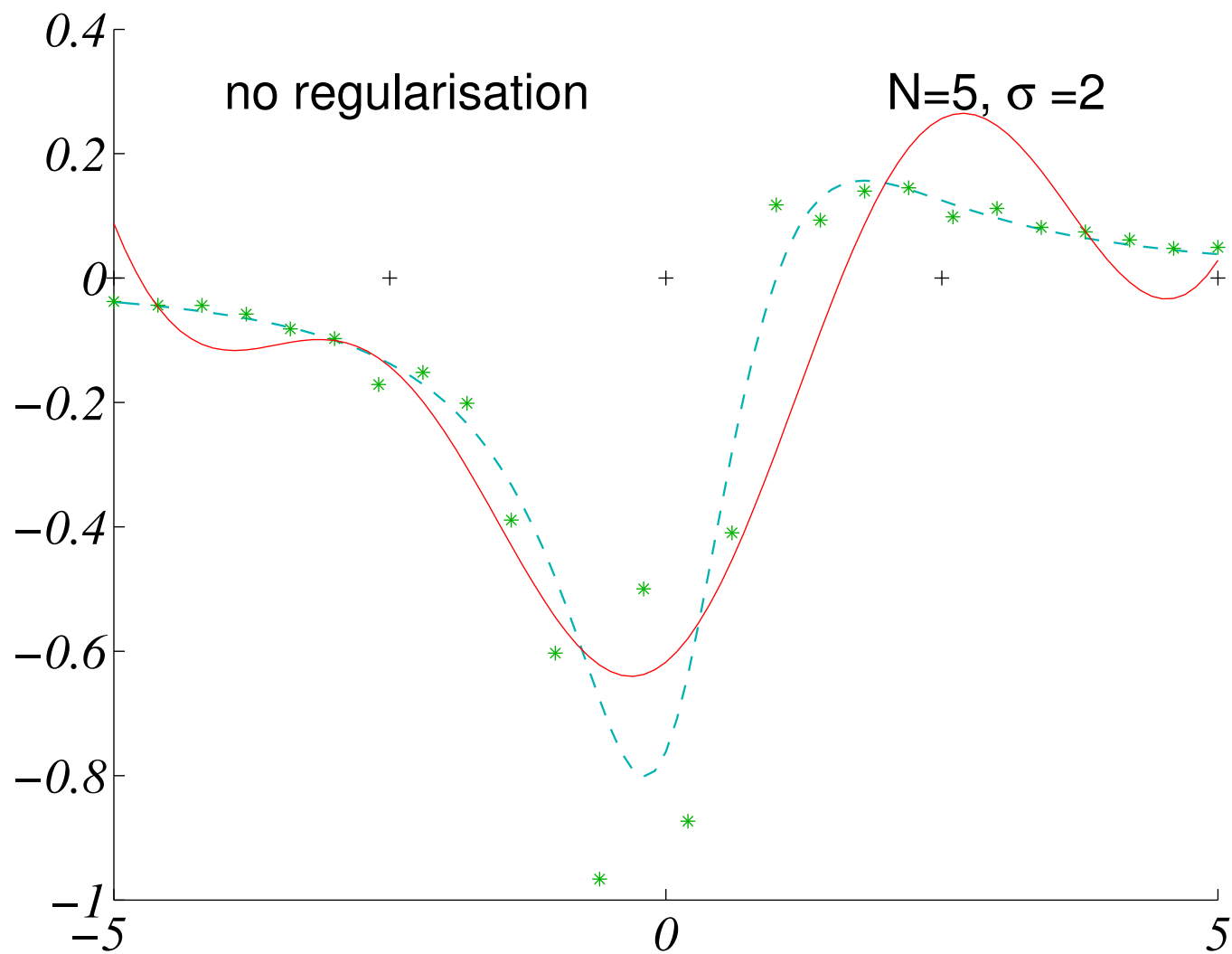


Function

Example

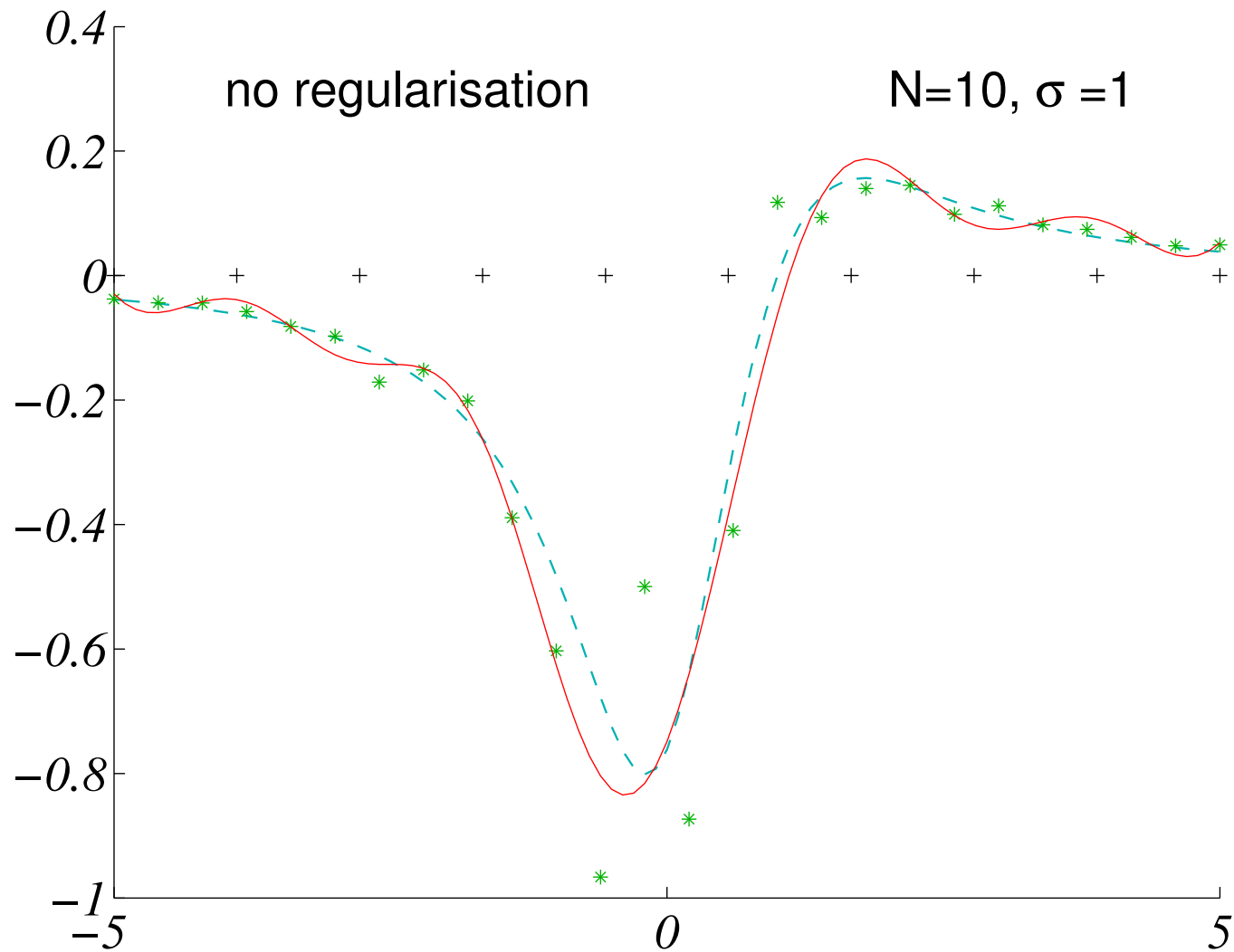


Example



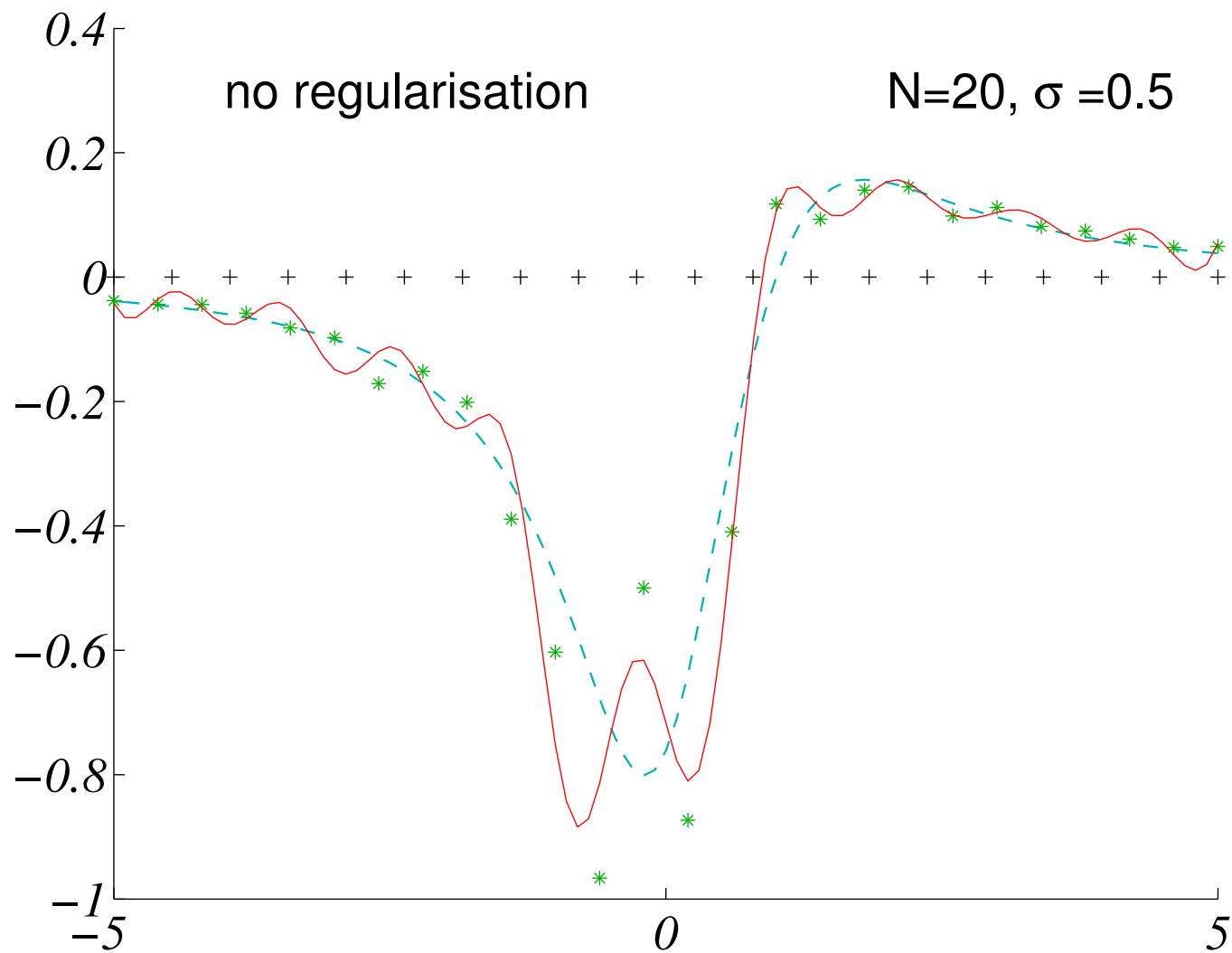
Centres too far apart

Example



Centres about right

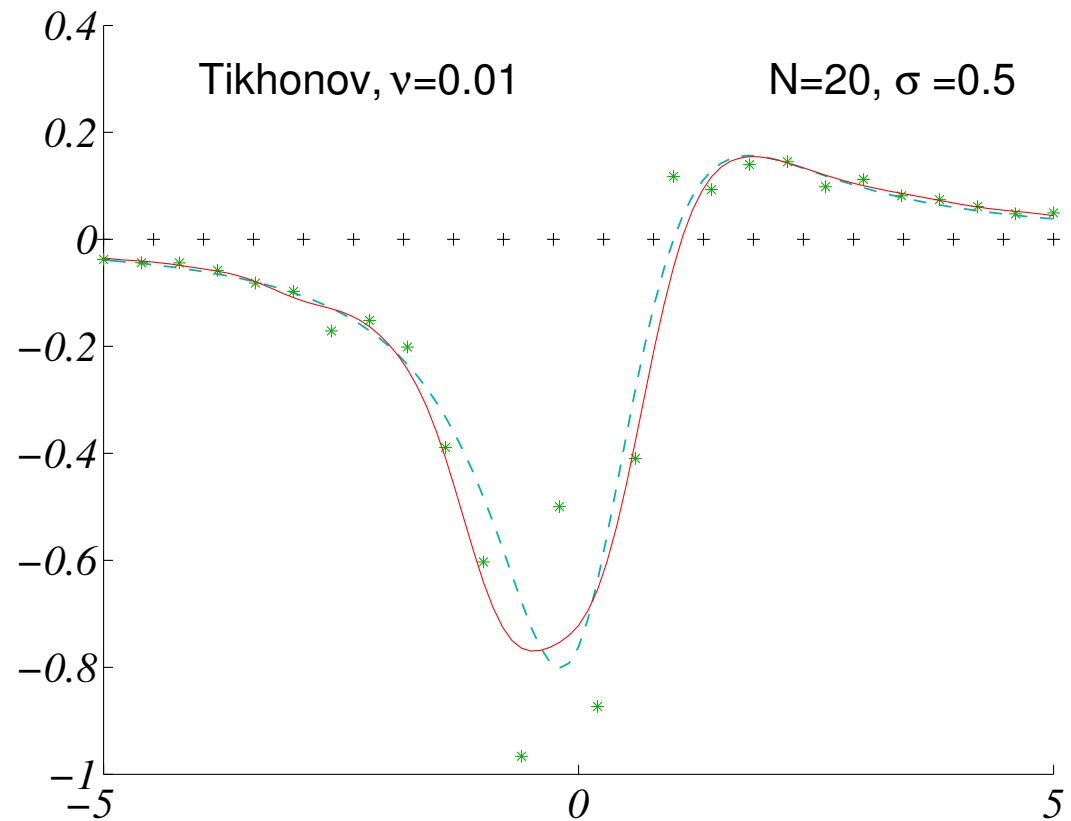
Example



Too many centres

Outline

1. Training Output Layer
2. **Regularisation**
 - Weight Decay
 - Tikhonov



Regularisation

- One way to improve generalisation performance is to bias a learning machine to learn simpler (smoother) functions
- This should reduce the sensitivity of the learning machine on the learning data
- To achieve this we can add **regularisation terms** that punishes complex functions

Regularisation

- One way to improve generalisation performance is to bias a learning machine to learn simpler (smoother) functions
- This should reduce the sensitivity of the learning machine on the learning data
- To achieve this we can add **regularisation terms** that punishes complex functions

Regularisation

- One way to improve generalisation performance is to bias a learning machine to learn simpler (smoother) functions
- This should reduce the sensitivity of the learning machine on the learning data
- To achieve this we can add **regularisation terms** that punishes complex functions

Regularisation Term

- We can add a regularisation term to force smoothness

$$E = \|\Phi^T \mathbf{w} - \mathbf{y}\|^2 + \nu R(\mathbf{w})$$

- Common regularisation terms are

- ★ Weight decay $R(\mathbf{x}) = \|\mathbf{w}\|^2$

- * Easy to implement

- * *Ad hoc*

- ★ Tikhonov regularisation $R(\mathbf{x}) = \int \sum_i \left(\frac{\partial^2 f(\mathbf{x}; \mathbf{w})}{\partial x_i^2} \right)^2 p(\mathbf{x}) d\mathbf{x}$

- * Slightly more complicated to implement

- * Punishes rapid changes in gradient

Regularisation Term

- We can add a regularisation term to force smoothness

$$E = \|\Phi^\top \mathbf{w} - \mathbf{y}\|^2 + \nu R(\mathbf{w})$$

- Common regularisation terms are

- ★ Weight decay $R(\mathbf{x}) = \|\mathbf{w}\|^2$

- * Easy to implement

- * *Ad hoc*

- ★ Tikhonov regularisation $R(\mathbf{x}) = \int \sum_i \left(\frac{\partial^2 f(\mathbf{x}; \mathbf{w})}{\partial x_i^2} \right)^2 p(\mathbf{x}) d\mathbf{x}$

- * Slightly more complicated to implement

- * Punishes rapid changes in gradient

Regularisation Term

- We can add a regularisation term to force smoothness

$$E = \|\Phi^T \mathbf{w} - \mathbf{y}\|^2 + \nu R(\mathbf{w})$$

- Common regularisation terms are

- ★ Weight decay $R(\mathbf{x}) = \|\mathbf{w}\|^2$

- * Easy to implement

- * *Ad hoc*

- ★ Tikhonov regularisation $R(\mathbf{x}) = \int \sum_i \left(\frac{\partial^2 f(\mathbf{x}; \mathbf{w})}{\partial x_i^2} \right)^2 p(\mathbf{x}) d\mathbf{x}$

- * Slightly more complicated to implement

- * Punishes rapid changes in gradient

Weight Decay

- Define the 'error'

$$E = \|\Phi^T \mathbf{w} - \mathbf{y}\|^2 + \nu \|\mathbf{w}\|^2$$

- First term is proportional to means squared error
- Second term is **regularisation** term
- ν (Greek letter nu) controls balance between minimising the training error and preferring smooth solutions

Weight Decay

- Define the 'error'

$$E = \|\Phi^T \mathbf{w} - \mathbf{y}\|^2 + \nu \|\mathbf{w}\|^2$$

- First term is proportional to means squared error
- Second term is **regularisation** term
- ν (Greek letter nu) controls balance between minimising the training error and preferring smooth solutions

Weight Decay

- Define the 'error'

$$E = \|\Phi^T \mathbf{w} - \mathbf{y}\|^2 + \nu \|\mathbf{w}\|^2$$

- First term is proportional to means squared error
- Second term is **regularisation** term
- ν (Greek letter nu) controls balance between minimising the training error and preferring smooth solutions

Weight Decay

- Define the ‘error’

$$E = \|\Phi^T \mathbf{w} - \mathbf{y}\|^2 + \nu \|\mathbf{w}\|^2$$

- First term is proportional to means squared error
- Second term is **regularisation** term
- ν (Greek letter nu) controls balance between minimising the training error and preferring smooth solutions

Least Mean Squared Solution

- Regularised learning error

$$\begin{aligned} E &= \|\Phi^T \mathbf{w} - \mathbf{y}\|^2 + \nu \|\mathbf{w}\|^2 \\ &= \mathbf{w}^T \Phi \Phi^T \mathbf{w} - 2\mathbf{w}^T \Phi \mathbf{y} + \mathbf{y}^T \mathbf{y} + \nu \mathbf{w}^T \mathbf{w} \\ &= \mathbf{w}^T (\Phi \Phi^T + \nu \mathbf{I}) \mathbf{w} - 2\mathbf{w}^T \Phi \mathbf{y} + \mathbf{y}^T \mathbf{y} \end{aligned}$$

- Gradient

$$\nabla E = 2 (\Phi \Phi^T + \nu \mathbf{I}) \mathbf{w} - 2 \Phi \mathbf{y}$$

- Minimum $\nabla E = 0$

$$\mathbf{w}^* = (\Phi \Phi^T + \nu \mathbf{I})^{-1} \Phi \mathbf{y}$$

Least Mean Squared Solution

- Regularised learning error

$$\begin{aligned} E &= \|\Phi^\top \mathbf{w} - \mathbf{y}\|^2 + \nu \|\mathbf{w}\|^2 \\ &= \mathbf{w}^\top \Phi \Phi^\top \mathbf{w} - 2\mathbf{w}^\top \Phi \mathbf{y} + \mathbf{y}^\top \mathbf{y} + \nu \mathbf{w}^\top \mathbf{w} \\ &= \mathbf{w}^\top (\Phi \Phi^\top + \nu \mathbf{I}) \mathbf{w} - 2\mathbf{w}^\top \Phi \mathbf{y} + \mathbf{y}^\top \mathbf{y} \end{aligned}$$

- Gradient

$$\nabla E = 2 (\Phi \Phi^\top + \nu \mathbf{I}) \mathbf{w} - 2 \Phi \mathbf{y}$$

- Minimum $\nabla E = 0$

$$\mathbf{w}^* = (\Phi \Phi^\top + \nu \mathbf{I})^{-1} \Phi \mathbf{y}$$

Least Mean Squared Solution

- Regularised learning error

$$\begin{aligned} E &= \|\Phi^\top \mathbf{w} - \mathbf{y}\|^2 + \nu \|\mathbf{w}\|^2 \\ &= \mathbf{w}^\top \Phi \Phi^\top \mathbf{w} - 2\mathbf{w}^\top \Phi \mathbf{y} + \mathbf{y}^\top \mathbf{y} + \nu \mathbf{w}^\top \mathbf{w} \\ &= \mathbf{w}^\top (\Phi \Phi^\top + \nu \mathbf{I}) \mathbf{w} - 2\mathbf{w}^\top \Phi \mathbf{y} + \mathbf{y}^\top \mathbf{y} \end{aligned}$$

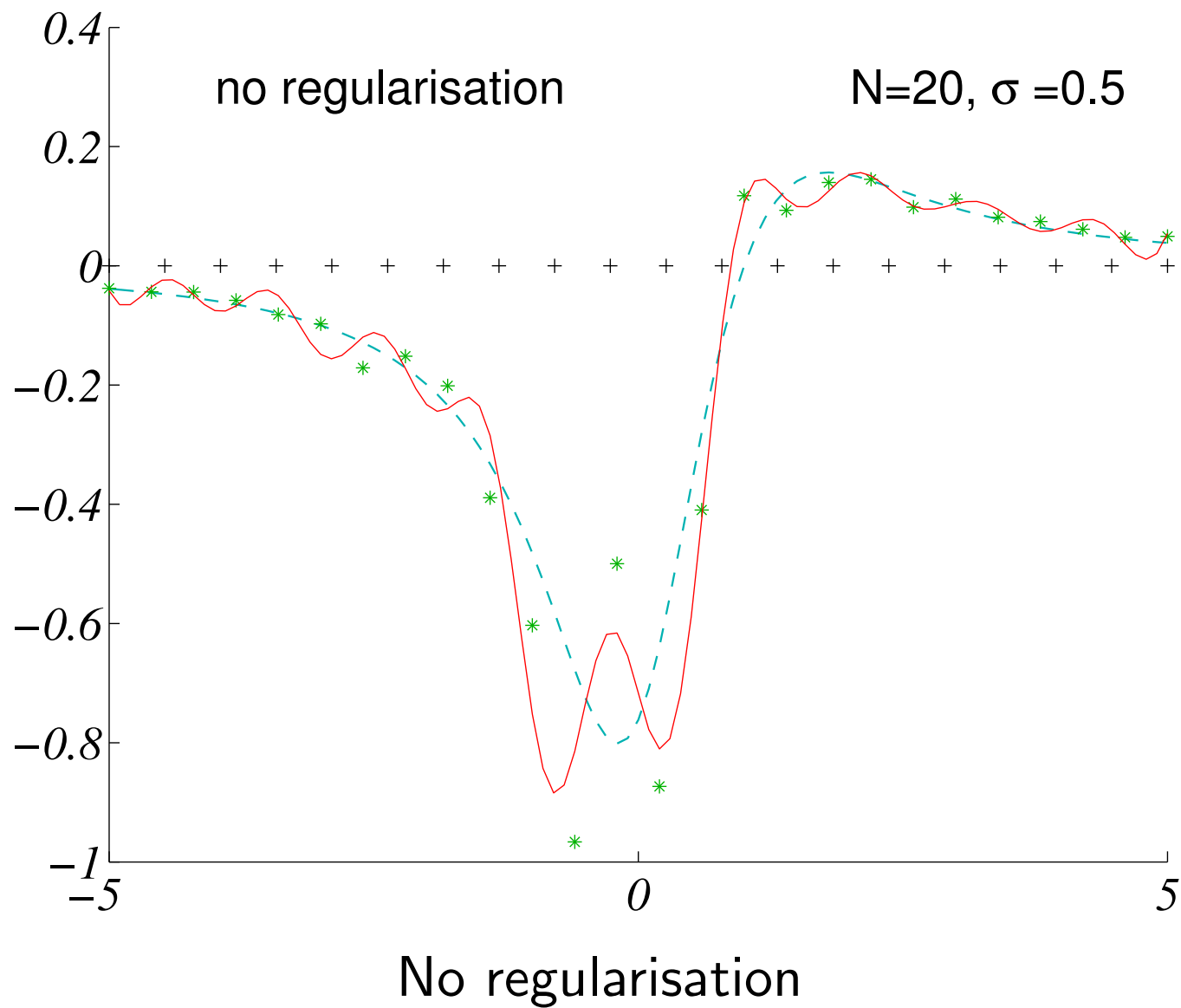
- Gradient

$$\nabla E = 2 (\Phi \Phi^\top + \nu \mathbf{I}) \mathbf{w} - 2 \Phi \mathbf{y}$$

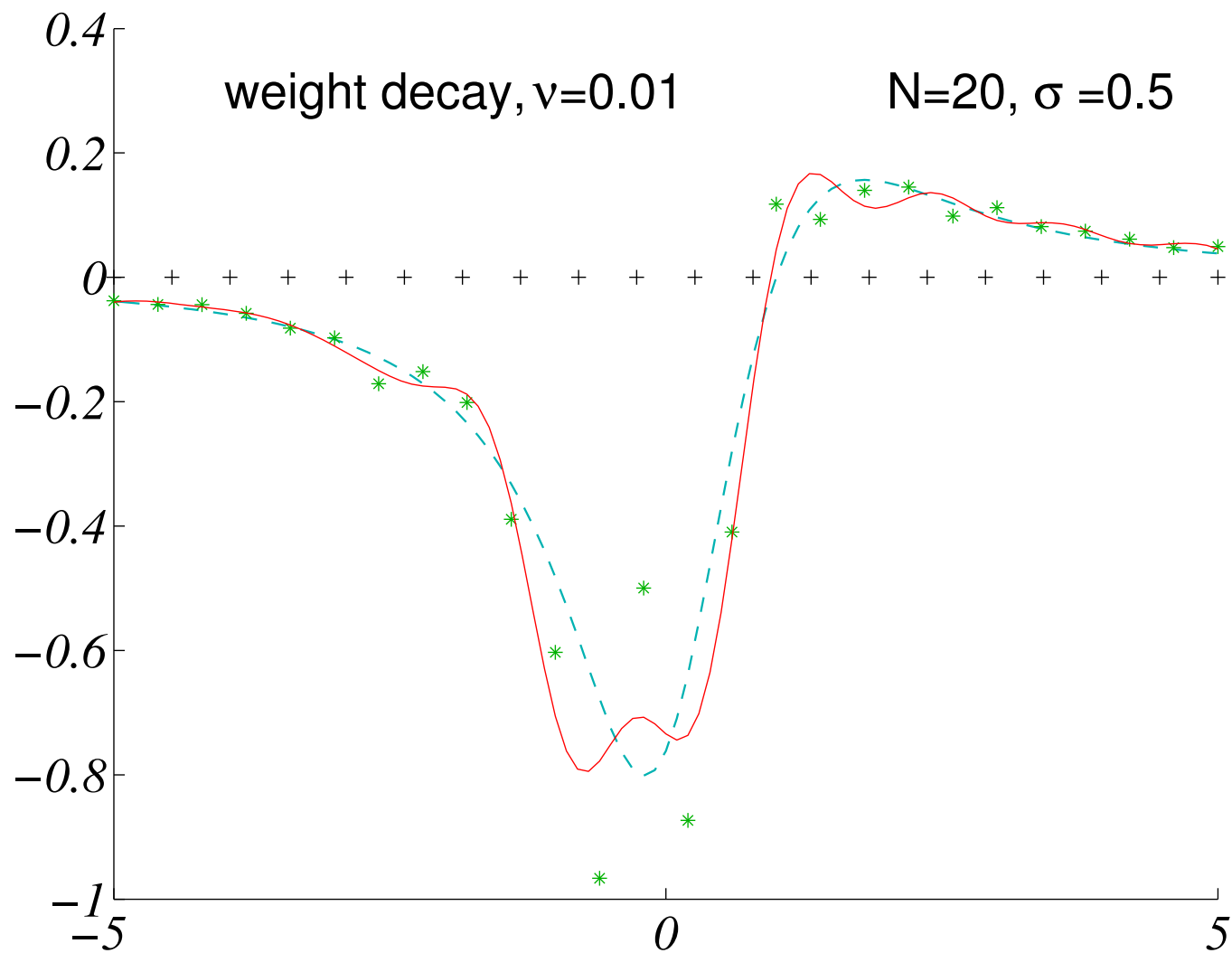
- Minimum $\nabla E = 0$

$$\mathbf{w}^* = (\Phi \Phi^\top + \nu \mathbf{I})^{-1} \Phi \mathbf{y}$$

Example



Example



Weight decay regularisation

Eigen Analysis

- Regularised least mean squared solution

$$\mathbf{w}^* = (\Phi \Phi^\top + \nu \mathbf{I})^{-1} \Phi \mathbf{y}$$

- Now $\Phi \Phi^\top$ is a symmetric matrix which can be written

$$\Phi \Phi^\top = \mathbf{U} \Lambda \mathbf{U}^\top$$

- $\Lambda = \mathbf{S} \mathbf{S}^\top = \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_N)$, $\mathbf{U} = (\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_N)$
- $\Phi \Phi^\top$ is positive semi-definite so $\lambda_i \geq 0$

Eigen Analysis

- Regularised least mean squared solution

$$\mathbf{w}^* = (\Phi \Phi^\top + \nu \mathbf{I})^{-1} \Phi \mathbf{y}$$

- Now $\Phi \Phi^\top$ is a symmetric matrix which can be written

$$\Phi \Phi^\top = \mathbf{U} \Lambda \mathbf{U}^\top$$

- $\Lambda = \mathbf{S} \mathbf{S}^\top = \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_N)$, $\mathbf{U} = (\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_N)$
- $\Phi \Phi^\top$ is positive semi-definite so $\lambda_i \geq 0$

Eigen Analysis

- Regularised least mean squared solution

$$\mathbf{w}^* = (\Phi \Phi^\top + \nu \mathbf{I})^{-1} \Phi \mathbf{y}$$

- Now $\Phi \Phi^\top$ is a symmetric matrix which can be written

$$\Phi \Phi^\top = \mathbf{U} \Lambda \mathbf{U}^\top$$

- $\Lambda = \mathbf{S} \mathbf{S}^\top = \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_N)$, $\mathbf{U} = (\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_N)$
- $\Phi \Phi^\top$ is positive semi-definite so $\lambda_i \geq 0$

Eigen Analysis

- Regularised least mean squared solution

$$\mathbf{w}^* = (\Phi \Phi^\top + \nu \mathbf{I})^{-1} \Phi \mathbf{y}$$

- Now $\Phi \Phi^\top$ is a symmetric matrix which can be written

$$\Phi \Phi^\top = \mathbf{U} \Lambda \mathbf{U}^\top$$

- $\Lambda = \mathbf{S} \mathbf{S}^\top = \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_N)$, $\mathbf{U} = (\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_N)$

- $\Phi \Phi^\top$ is positive semi-definite so $\lambda_i \geq 0$

Properties of Regularised Matrix

- Regularised matrix

$$\begin{aligned}\Phi \Phi^T + \nu \mathbf{I} &= \mathbf{U} (\Lambda + \nu \mathbf{I}) \mathbf{U}^T \\ (\Phi \Phi^T + \nu \mathbf{I})^{-1} &= \mathbf{U} (\Lambda + \nu \mathbf{I})^{-1} \mathbf{U}^T\end{aligned}$$

where

$$(\Lambda + \nu \mathbf{I})^{-1} = \text{diag} \left(\frac{1}{\lambda_1 + \nu}, \dots, \frac{1}{\lambda_N + \nu} \right)$$

- Regularisation makes inverse well posed if it wasn't
- Regularisation improves the conditioning (i.e. sensitive to changes in data)

$$\frac{|\lambda_{max}|}{|\lambda_{min}|} \longrightarrow \frac{|\lambda_{max} + \nu|}{|\lambda_{min} + \nu|} < \frac{|\lambda_{max}|}{|\lambda_{min}|}$$

Properties of Regularised Matrix

- Regularised matrix

$$\begin{aligned}\Phi \Phi^T + \nu \mathbf{I} &= \mathbf{U} (\Lambda + \nu \mathbf{I}) \mathbf{U}^T \\ (\Phi \Phi^T + \nu \mathbf{I})^{-1} &= \mathbf{U} (\Lambda + \nu \mathbf{I})^{-1} \mathbf{U}^T\end{aligned}$$

where

$$(\Lambda + \nu \mathbf{I})^{-1} = \text{diag} \left(\frac{1}{\lambda_1 + \nu}, \dots, \frac{1}{\lambda_N + \nu} \right)$$

- Regularisation makes inverse well posed if it wasn't
- Regularisation improves the conditioning (i.e. sensitive to changes in data)

$$\frac{|\lambda_{max}|}{|\lambda_{min}|} \longrightarrow \frac{|\lambda_{max} + \nu|}{|\lambda_{min} + \nu|} < \frac{|\lambda_{max}|}{|\lambda_{min}|}$$

Properties of Regularised Matrix

- Regularised matrix

$$\begin{aligned}\Phi \Phi^T + \nu \mathbf{I} &= \mathbf{U} (\Lambda + \nu \mathbf{I}) \mathbf{U}^T \\ (\Phi \Phi^T + \nu \mathbf{I})^{-1} &= \mathbf{U} (\Lambda + \nu \mathbf{I})^{-1} \mathbf{U}^T\end{aligned}$$

where

$$(\Lambda + \nu \mathbf{I})^{-1} = \text{diag} \left(\frac{1}{\lambda_1 + \nu}, \dots, \frac{1}{\lambda_N + \nu} \right)$$

- Regularisation makes inverse well posed if it wasn't
- Regularisation improves the conditioning (i.e. sensitive to changes in data)

$$\frac{|\lambda_{max}|}{|\lambda_{min}|} \longrightarrow \frac{|\lambda_{max} + \nu|}{|\lambda_{min} + \nu|} < \frac{|\lambda_{max}|}{|\lambda_{min}|}$$

Properties of Regularised Matrix

- Regularised matrix

$$\begin{aligned}\Phi \Phi^T + \nu \mathbf{I} &= \mathbf{U} (\Lambda + \nu \mathbf{I}) \mathbf{U}^T \\ (\Phi \Phi^T + \nu \mathbf{I})^{-1} &= \mathbf{U} (\Lambda + \nu \mathbf{I})^{-1} \mathbf{U}^T\end{aligned}$$

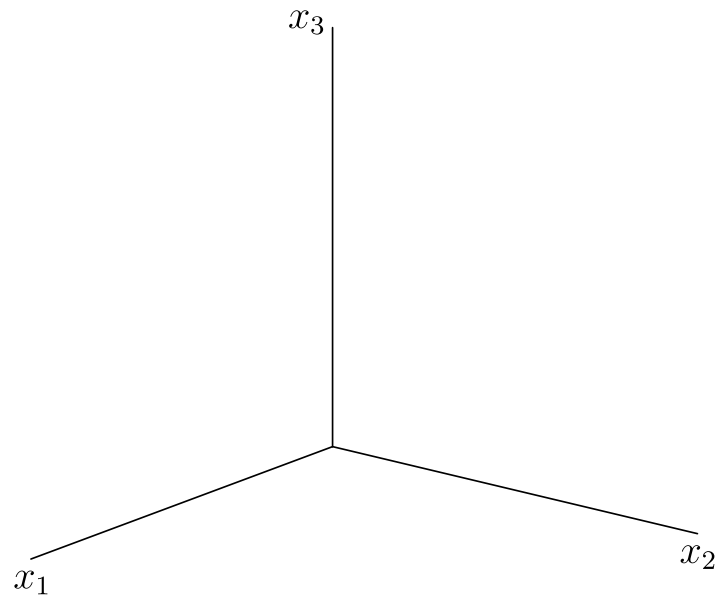
where

$$(\Lambda + \nu \mathbf{I})^{-1} = \text{diag} \left(\frac{1}{\lambda_1 + \nu}, \dots, \frac{1}{\lambda_N + \nu} \right)$$

- Regularisation makes inverse well posed if it wasn't
- Regularisation improves the conditioning (i.e. sensitive to changes in data)

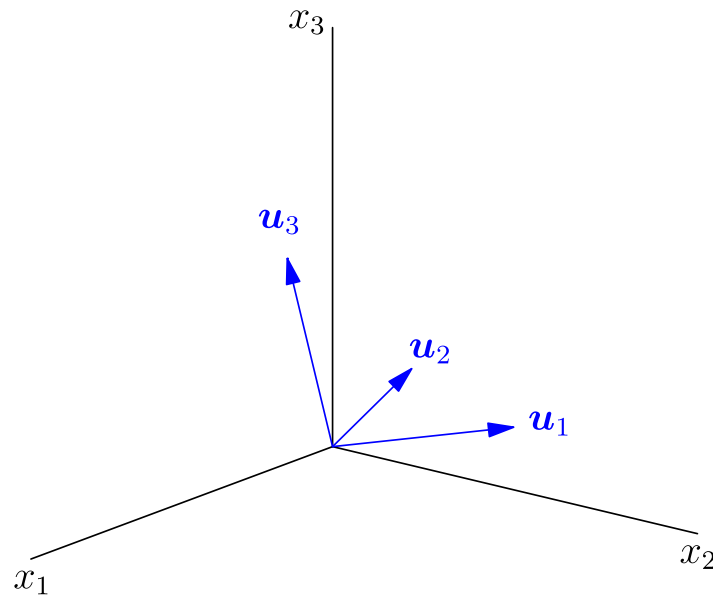
$$\frac{|\lambda_{max}|}{|\lambda_{min}|} \longrightarrow \frac{|\lambda_{max} + \nu|}{|\lambda_{min} + \nu|} < \frac{|\lambda_{max}|}{|\lambda_{min}|}$$

III-Conditioning



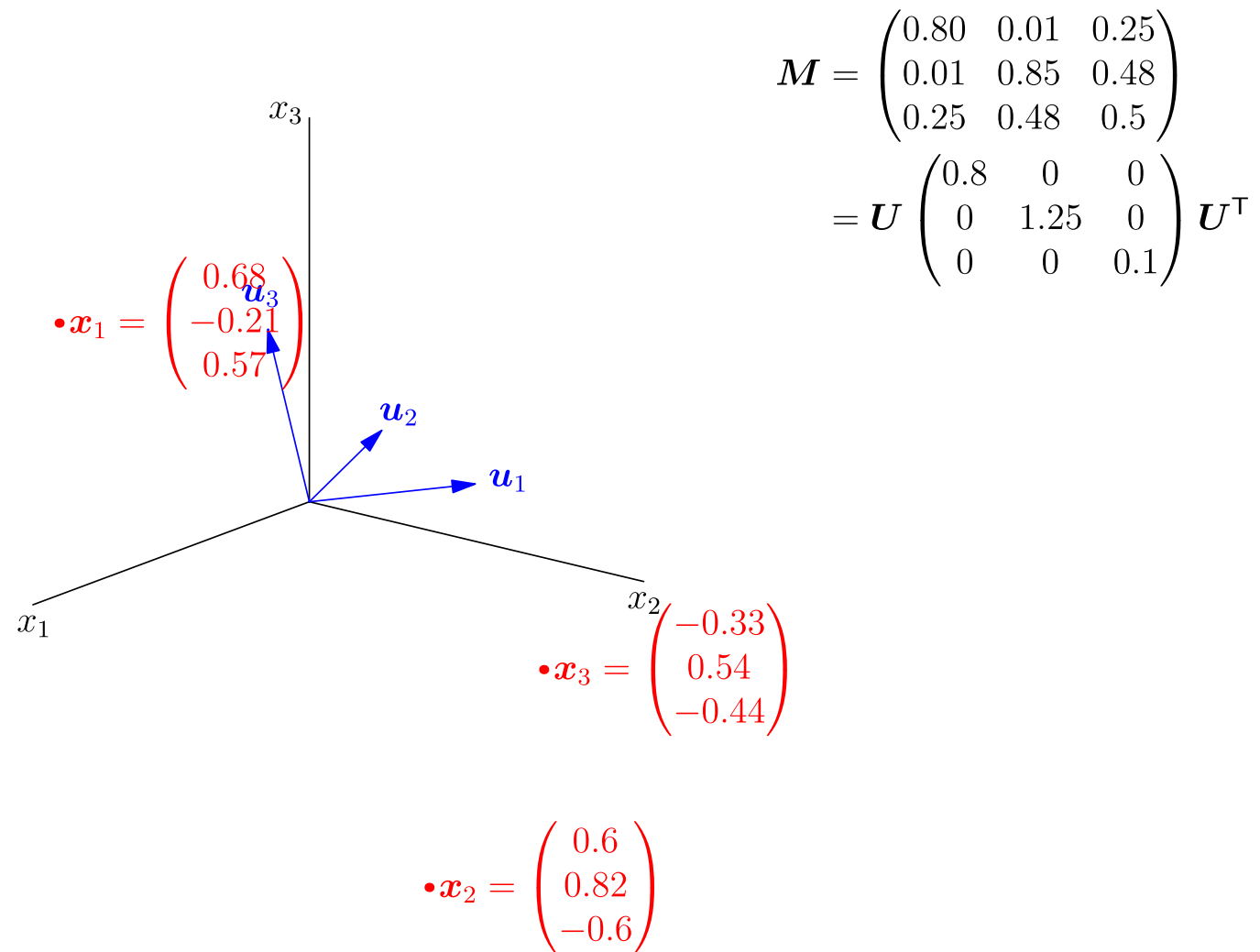
$$\begin{aligned} \mathbf{M} &= \begin{pmatrix} 0.80 & 0.01 & 0.25 \\ 0.01 & 0.85 & 0.48 \\ 0.25 & 0.48 & 0.5 \end{pmatrix} \\ &= \mathbf{U} \begin{pmatrix} 0.8 & 0 & 0 \\ 0 & 1.25 & 0 \\ 0 & 0 & 0.1 \end{pmatrix} \mathbf{U}^T \end{aligned}$$

III-Conditioning

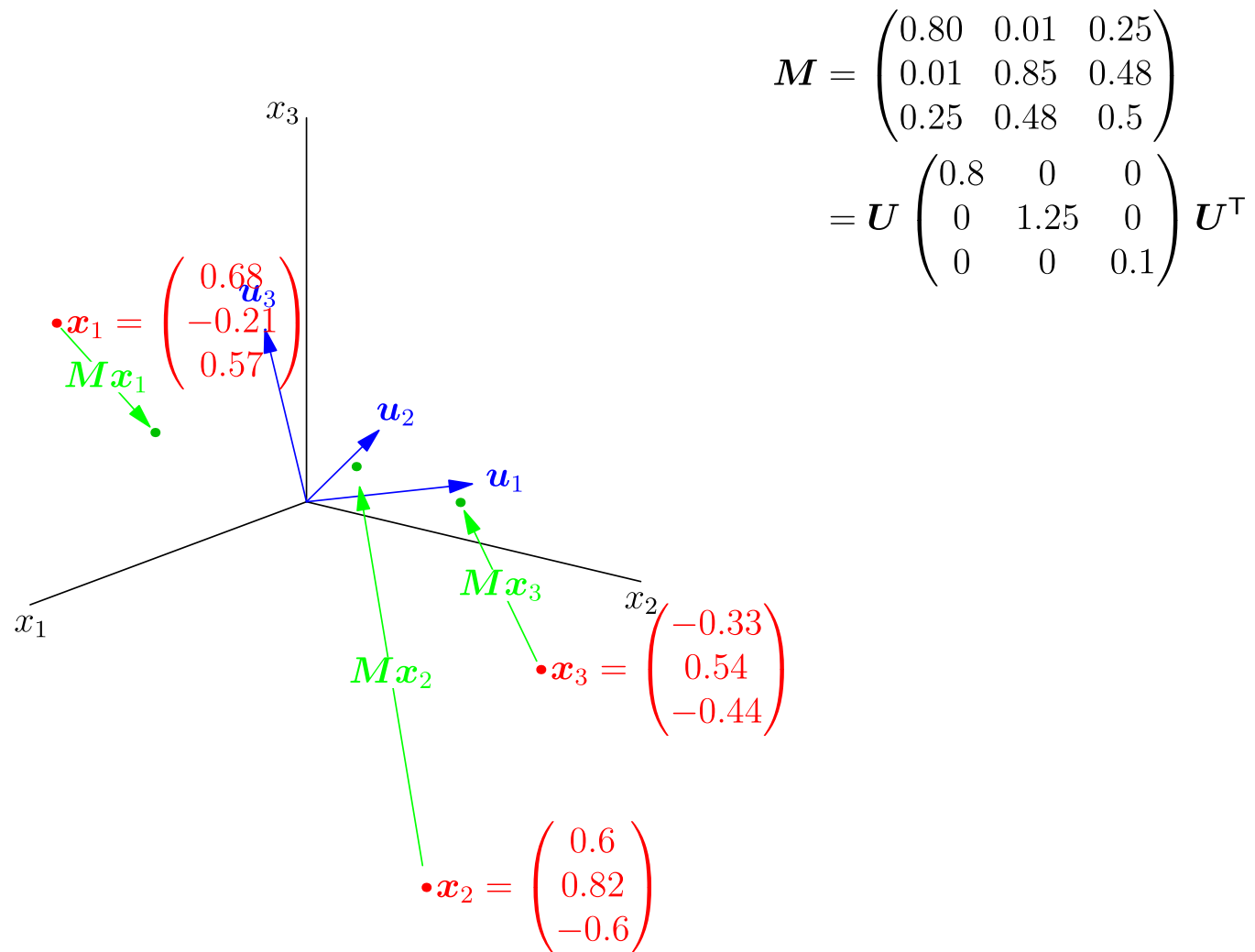


$$\begin{aligned} M &= \begin{pmatrix} 0.80 & 0.01 & 0.25 \\ 0.01 & 0.85 & 0.48 \\ 0.25 & 0.48 & 0.5 \end{pmatrix} \\ &= U \begin{pmatrix} 0.8 & 0 & 0 \\ 0 & 1.25 & 0 \\ 0 & 0 & 0.1 \end{pmatrix} U^T \end{aligned}$$

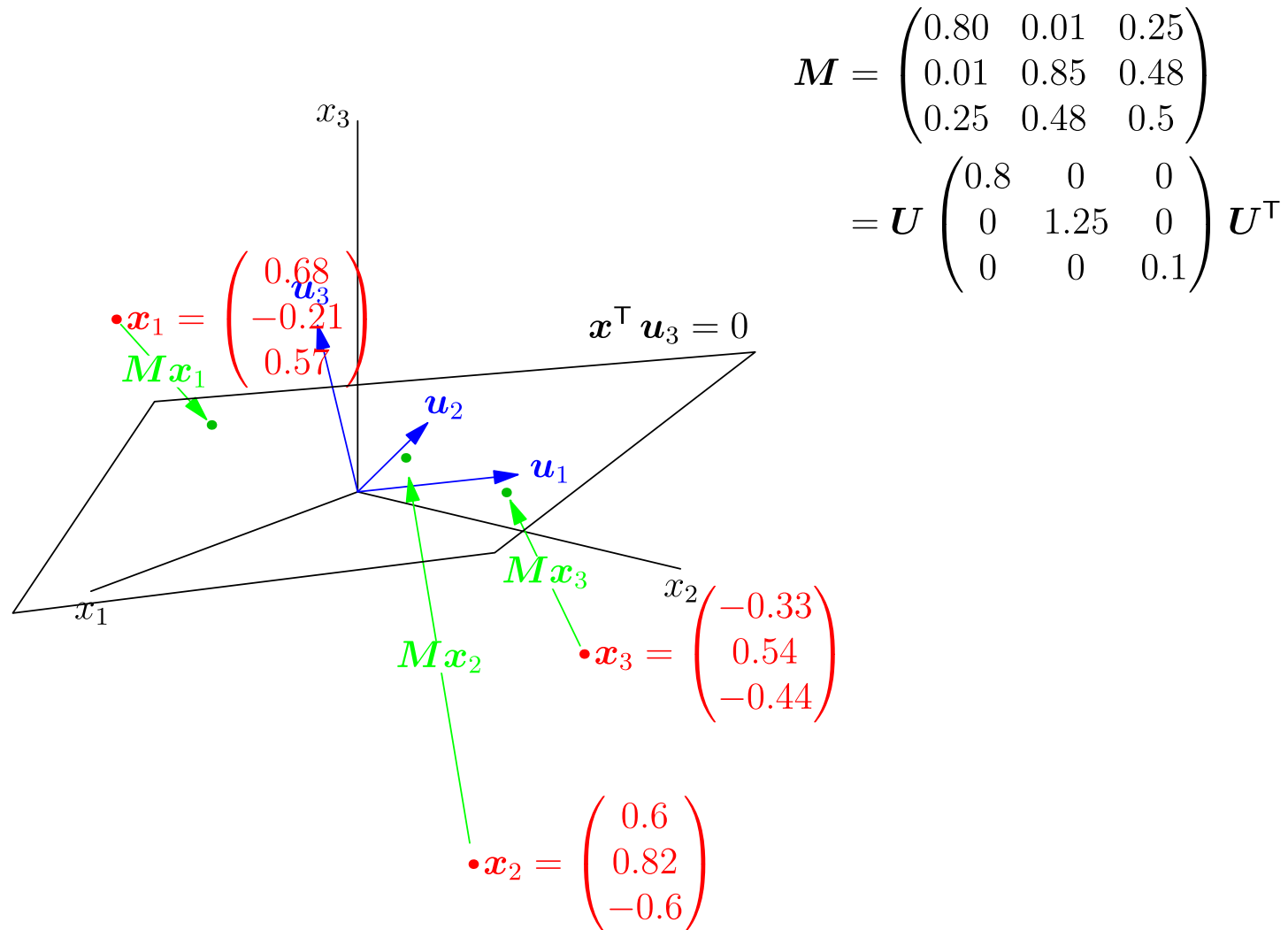
III-Conditioning



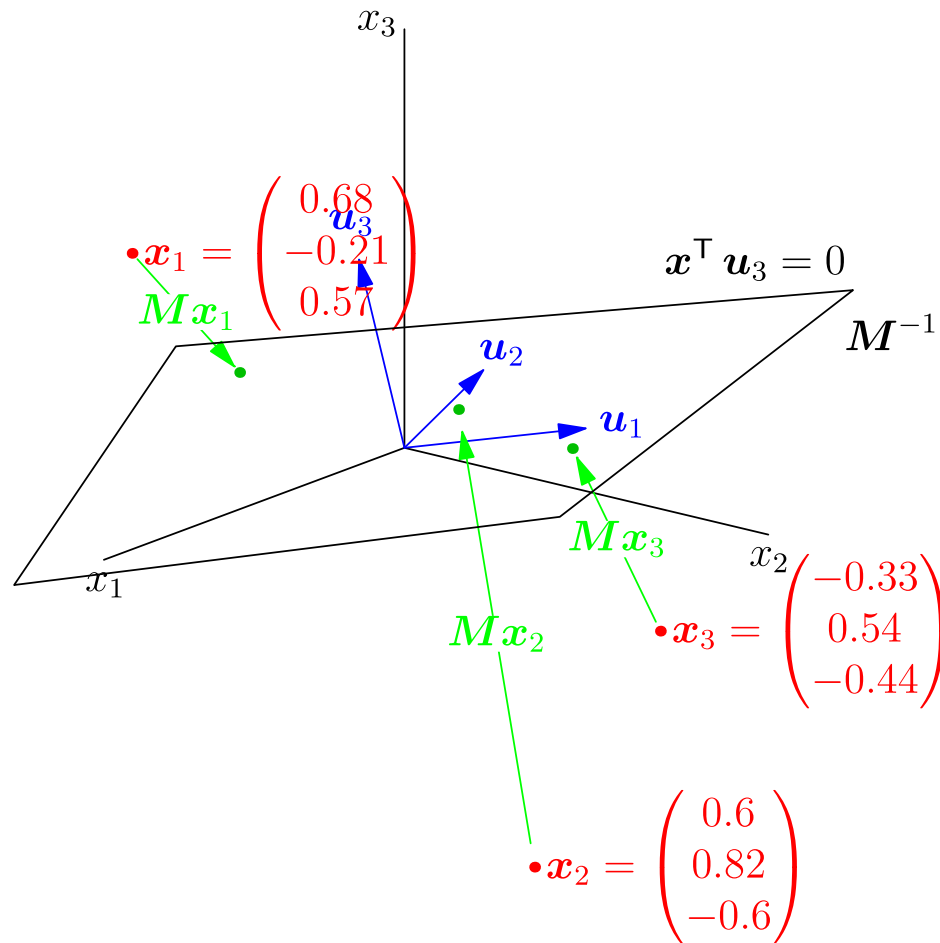
III-Conditioning



III-Conditioning



III-Conditioning



$$M = \begin{pmatrix} 0.80 & 0.01 & 0.25 \\ 0.01 & 0.85 & 0.48 \\ 0.25 & 0.48 & 0.5 \end{pmatrix}$$

$$= U \begin{pmatrix} 0.8 & 0 & 0 \\ 0 & 1.25 & 0 \\ 0 & 0 & 0.1 \end{pmatrix} U^T$$

$$M^{-1} = \begin{pmatrix} 1.9 & 1.17 & -2.1 \\ 1.17 & 3.34 & -3.8 \\ -2.1 & -3.8 & 6.8 \end{pmatrix}$$

$$= U \begin{pmatrix} 1.25 & 0 & 0 \\ 0 & 0.8 & 0 \\ 0 & 0 & 10 \end{pmatrix} U^T$$

What Does It Mean?

- Without regularisation the optimal weights are

$$\mathbf{w}^* = (\Phi \Phi^\top)^{-1} \Phi \mathbf{y}$$

- The matrix $\Phi \Phi^\top$ depends on the data
- If $\Phi \Phi^\top$ is ill-conditioned the inverse is unreliable
- With a weight decay regulariser the optimal weights are

$$\mathbf{w}^* = (\Phi \Phi^\top + \nu \mathbf{I})^{-1} \Phi \mathbf{y}$$

- These don't fit the data quite so well, but they are less sensitive to the data

What Does It Mean?

- Without regularisation the optimal weights are

$$\mathbf{w}^* = (\Phi \Phi^\top)^{-1} \Phi \mathbf{y}$$

- The matrix $\Phi \Phi^\top$ depends on the data
- If $\Phi \Phi^\top$ is ill-conditioned the inverse is unreliable
- With a weight decay regulariser the optimal weights are

$$\mathbf{w}^* = (\Phi \Phi^\top + \nu \mathbf{I})^{-1} \Phi \mathbf{y}$$

- These don't fit the data quite so well, but they are less sensitive to the data

What Does It Mean?

- Without regularisation the optimal weights are

$$\mathbf{w}^* = (\Phi \Phi^\top)^{-1} \Phi \mathbf{y}$$

- The matrix $\Phi \Phi^\top$ depends on the data
- If $\Phi \Phi^\top$ is ill-conditioned the inverse is unreliable
- With a weight decay regulariser the optimal weights are

$$\mathbf{w}^* = (\Phi \Phi^\top + \nu \mathbf{I})^{-1} \Phi \mathbf{y}$$

- These don't fit the data quite so well, but they are less sensitive to the data

What Does It Mean?

- Without regularisation the optimal weights are

$$\mathbf{w}^* = (\Phi \Phi^\top)^{-1} \Phi \mathbf{y}$$

- The matrix $\Phi \Phi^\top$ depends on the data
- If $\Phi \Phi^\top$ is ill-conditioned the inverse is unreliable
- With a weight decay regulariser the optimal weights are

$$\mathbf{w}^* = (\Phi \Phi^\top + \nu \mathbf{I})^{-1} \Phi \mathbf{y}$$

- These don't fit the data quite so well, but they are less sensitive to the data

What Does It Mean?

- Without regularisation the optimal weights are

$$\mathbf{w}^* = (\Phi \Phi^\top)^{-1} \Phi \mathbf{y}$$

- The matrix $\Phi \Phi^\top$ depends on the data
- If $\Phi \Phi^\top$ is ill-conditioned the inverse is unreliable
- With a weight decay regulariser the optimal weights are

$$\mathbf{w}^* = (\Phi \Phi^\top + \nu \mathbf{I})^{-1} \Phi \mathbf{y}$$

- These don't fit the data quite so well, but they are less sensitive to the data

Tikhonov Regularisation

- We can use a more direct penalty for non-smooth functions

$$\int \sum_{i=1}^N \left(\frac{\partial^2 f(\mathbf{x}; \mathbf{w})}{\partial x_i^2} \right)^2 p(\mathbf{x}) d\mathbf{x}$$

- In practice we cannot compute this
- We can estimate this quantity by

$$\frac{1}{P} \sum_{p=1}^P \sum_{i=1}^N \left(\frac{\partial^2 f(\mathbf{x}_p; \mathbf{w})}{\partial x_i^2} \right)^2$$

- Average over all input patterns

Tikhonov Regularisation

- We can use a more direct penalty for non-smooth functions

$$\int \sum_{i=1}^N \left(\frac{\partial^2 f(\mathbf{x}; \mathbf{w})}{\partial x_i^2} \right)^2 p(\mathbf{x}) d\mathbf{x}$$

- In practice we cannot compute this
- We can estimate this quantity by

$$\frac{1}{P} \sum_{p=1}^P \sum_{i=1}^N \left(\frac{\partial^2 f(\mathbf{x}_p; \mathbf{w})}{\partial x_i^2} \right)^2$$

- Average over all input patterns

Tikhonov Regularisation

- We can use a more direct penalty for non-smooth functions

$$\int \sum_{i=1}^N \left(\frac{\partial^2 f(\mathbf{x}; \mathbf{w})}{\partial x_i^2} \right)^2 p(\mathbf{x}) d\mathbf{x}$$

- In practice we cannot compute this
- We can estimate this quantity by

$$\frac{1}{P} \sum_{p=1}^P \sum_{i=1}^N \left(\frac{\partial^2 f(\mathbf{x}_p; \mathbf{w})}{\partial x_i^2} \right)^2$$

- Average over all input patterns

Tikhonov Regularisation

- We can use a more direct penalty for non-smooth functions

$$\int \sum_{i=1}^N \left(\frac{\partial^2 f(\mathbf{x}; \mathbf{w})}{\partial x_i^2} \right)^2 p(\mathbf{x}) d\mathbf{x}$$

- In practice we cannot compute this
- We can estimate this quantity by

$$\frac{1}{P} \sum_{p=1}^P \sum_{i=1}^N \left(\frac{\partial^2 f(\mathbf{x}_p; \mathbf{w})}{\partial x_i^2} \right)^2$$

- Average over all input patterns

Tikhonov Regulariser

- Now, $f(\mathbf{x}; \mathbf{w}) = \sum_k w_k \phi_k(\mathbf{x})$ so

$$\frac{\partial^2 f(\mathbf{x}; \mathbf{w})}{\partial x_i^2} = \sum_k w_k \frac{\partial^2 \phi_k(\mathbf{x})}{\partial x_i^2}$$

- And

$$\frac{1}{P} \sum_{p=1}^P \sum_{i=1}^N \left(\frac{\partial^2 f(\mathbf{x}_p; \mathbf{w})}{\partial x_i^2} \right)^2 = \mathbf{w}^\top \mathbf{T} \mathbf{w}$$

- where

$$T_{jk} = \frac{1}{P} \sum_{p=1}^P \sum_{i=1}^N \frac{\partial^2 \phi_j(\mathbf{x}_p)}{\partial x_i^2} \frac{\partial^2 \phi_k(\mathbf{x}_p)}{\partial x_i^2}$$

Tikhonov Regulariser

- Now, $f(\mathbf{x}; \mathbf{w}) = \sum_k w_k \phi_k(\mathbf{x})$ so

$$\frac{\partial^2 f(\mathbf{x}; \mathbf{w})}{\partial x_i^2} = \sum_k w_k \frac{\partial^2 \phi_k(\mathbf{x})}{\partial x_i^2}$$

- And

$$\frac{1}{P} \sum_{p=1}^P \sum_{i=1}^N \left(\frac{\partial^2 f(\mathbf{x}_p; \mathbf{w})}{\partial x_i^2} \right)^2 = \mathbf{w}^\top \mathbf{T} \mathbf{w}$$

- where

$$T_{jk} = \frac{1}{P} \sum_{p=1}^P \sum_{i=1}^N \frac{\partial^2 \phi_j(\mathbf{x}_p)}{\partial x_i^2} \frac{\partial^2 \phi_k(\mathbf{x}_p)}{\partial x_i^2}$$

Tikhonov Regulariser

- For Gaussian RBFs, $\phi_k(\mathbf{x}) = e^{-\|\mathbf{x} - \boldsymbol{\mu}\|^2 / 2\sigma^2}$

$$\frac{\partial^2 \phi_k(\mathbf{x})}{\partial x_i^2} = \left(\frac{(x_i - \mu_i)^2 - \sigma^2}{\sigma^4} \right) \phi_k(\mathbf{x})$$

- Minimise

$$E = \|\boldsymbol{\Phi}^\top \mathbf{w} - \mathbf{y}\|^2 + \nu \mathbf{w}^\top \mathbf{T} \mathbf{w}$$

- $\nabla E = 0$ gives

$$\mathbf{w}^* = (\boldsymbol{\Phi} \boldsymbol{\Phi}^\top + \nu \mathbf{T})^{-1} \boldsymbol{\Phi} \mathbf{y}$$

Tikhonov Regulariser

- For Gaussian RBFs, $\phi_k(\mathbf{x}) = e^{-\|\mathbf{x} - \boldsymbol{\mu}\|^2 / 2\sigma^2}$

$$\frac{\partial^2 \phi_k(\mathbf{x})}{\partial x_i^2} = \left(\frac{(x_i - \mu_i)^2 - \sigma^2}{\sigma^4} \right) \phi_k(\mathbf{x})$$

- Minimise

$$E = \|\boldsymbol{\Phi}^\top \mathbf{w} - \mathbf{y}\|^2 + \nu \mathbf{w}^\top \mathbf{T} \mathbf{w}$$

- $\nabla E = 0$ gives

$$\mathbf{w}^* = (\boldsymbol{\Phi} \boldsymbol{\Phi}^\top + \nu \mathbf{T})^{-1} \boldsymbol{\Phi} \mathbf{y}$$

Tikhonov Regulariser

- For Gaussian RBFs, $\phi_k(\mathbf{x}) = e^{-\|\mathbf{x} - \boldsymbol{\mu}\|^2 / 2\sigma^2}$

$$\frac{\partial^2 \phi_k(\mathbf{x})}{\partial x_i^2} = \left(\frac{(x_i - \mu_i)^2 - \sigma^2}{\sigma^4} \right) \phi_k(\mathbf{x})$$

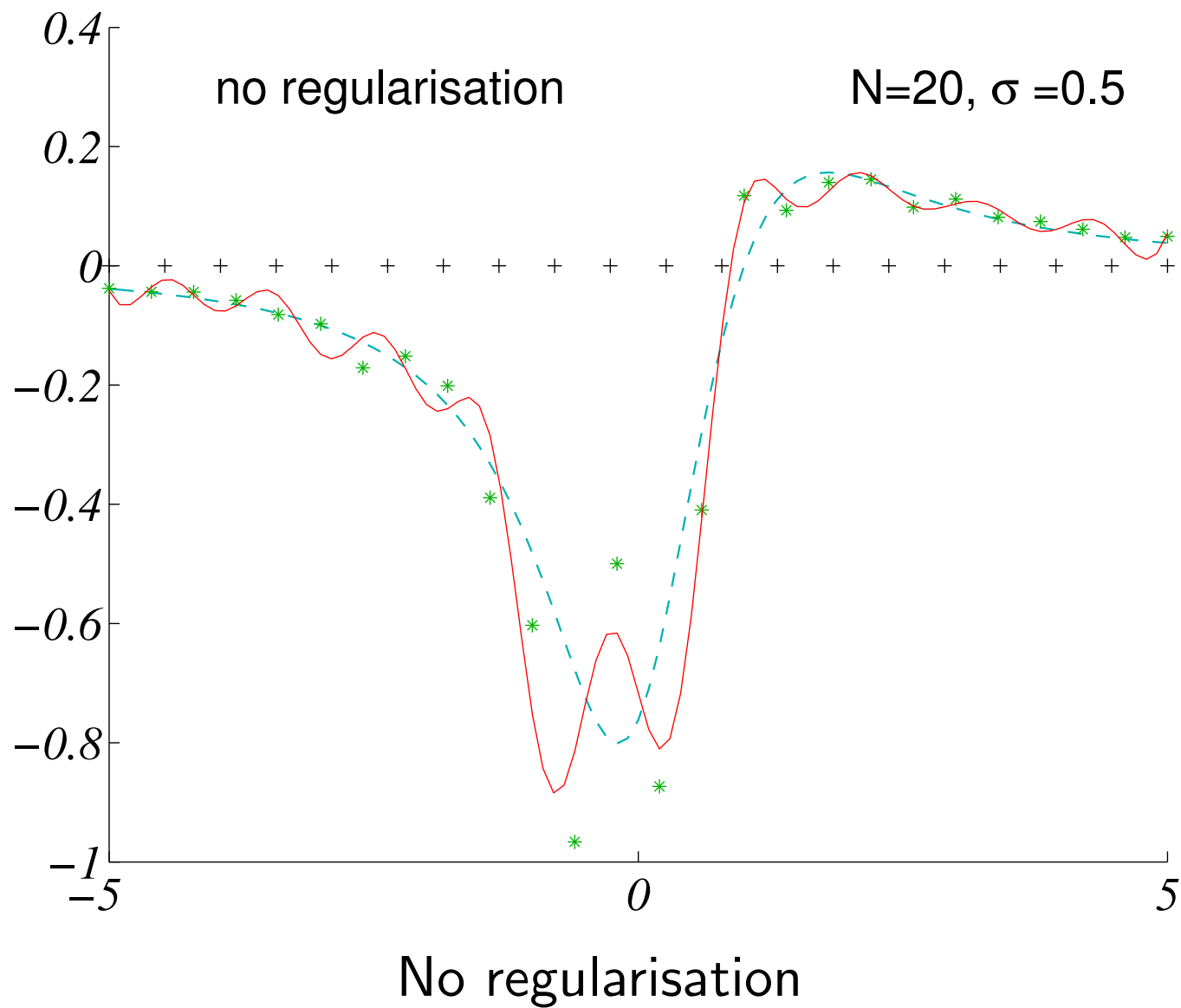
- Minimise

$$E = \|\boldsymbol{\Phi}^\top \mathbf{w} - \mathbf{y}\|^2 + \nu \mathbf{w}^\top \mathbf{T} \mathbf{w}$$

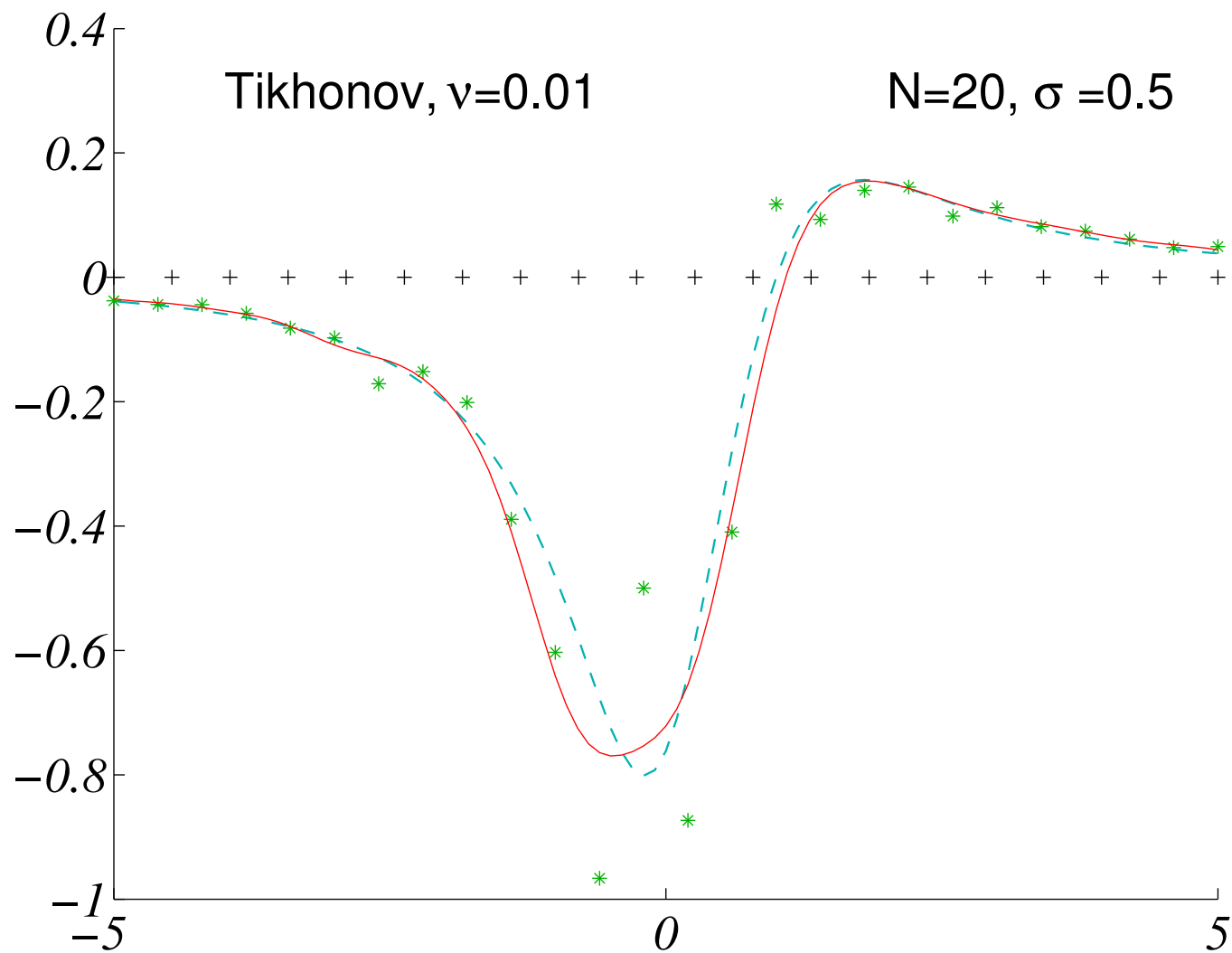
- $\nabla E = 0$ gives

$$\mathbf{w}^* = (\boldsymbol{\Phi} \boldsymbol{\Phi}^\top + \nu \mathbf{T})^{-1} \boldsymbol{\Phi} \mathbf{y}$$

Example



Example



Tikhonov Regularisation

Using Regularisers

- Regularisation terms are widely used—they really improve performance
- They can also be used in MLPs
- They prevent a learning machine with lots of parameters from **overfitting** the data
- They raise a new question “How do we choose the regularisation parameters, ν , etc.”
- Usually, we use another set of unseen data to test for the best regularisation parameters

Using Regularisers

- Regularisation terms are widely used—they really improve performance
- They can also be used in MLPs
- They prevent a learning machine with lots of parameters from **overfitting** the data
- They raise a new question “How do we choose the regularisation parameters, ν , etc.”
- Usually, we use another set of unseen data to test for the best regularisation parameters

Using Regularisers

- Regularisation terms are widely used—they really improve performance
- They can also be used in MLPs
- They prevent a learning machine with lots of parameters from **overfitting** the data
- They raise a new question “How do we choose the regularisation parameters, ν , etc.”
- Usually, we use another set of unseen data to test for the best regularisation parameters

Using Regularisers

- Regularisation terms are widely used—they really improve performance
- They can also be used in MLPs
- They prevent a learning machine with lots of parameters from **overfitting** the data
- They raise a new question “How do we choose the regularisation parameters, ν , etc.”
- Usually, we use another set of unseen data to test for the best regularisation parameters

Using Regularisers

- Regularisation terms are widely used—they really improve performance
- They can also be used in MLPs
- They prevent a learning machine with lots of parameters from **overfitting** the data
- They raise a new question “How do we choose the regularisation parameters, ν , etc.”
- Usually, we use another set of unseen data to test for the best regularisation parameters

Summary

- In machine learning we care about performance on unseen data
- Over complicated machines will tend to learn the data set well, but will have poor generalisation performance
- We often wish to tailor the complexity of the machine to the data set
- One way to do this automatically is by introducing regularisation terms

Summary

- In machine learning we care about performance on unseen data
- Over complicated machines will tend to learn the data set well, but will have poor generalisation performance
- We often wish to tailor the complexity of the machine to the data set
- One way to do this automatically is by introducing regularisation terms

Summary

- In machine learning we care about performance on unseen data
- Over complicated machines will tend to learn the data set well, but will have poor generalisation performance
- We often wish to tailor the complexity of the machine to the data set
- One way to do this automatically is by introducing regularisation terms

Summary

- In machine learning we care about performance on unseen data
- Over complicated machines will tend to learn the data set well, but will have poor generalisation performance
- We often wish to tailor the complexity of the machine to the data set
- One way to do this automatically is by introducing regularisation terms