

EE613
Machine Learning for Engineers

LINEAR REGRESSION

Sylvain Calinon
Robot Learning & Interaction Group
Idiap Research Institute
Nov. 4, 2015

Outline

- Multivariate ordinary least squares
- Singular value decomposition (SVD)
- Kernels in least squares (nullspace projection)
- Ridge regression (Tikhonov regularization)
- Weighted least squares (WLS)
- Recursive least squares (RLS)

Multivariate ordinary least squares

demo_LS01.m

demo_LS_polFit01.m

Multivariate ordinary least squares

- Least squares is everywhere: from simple problems to large scale problems.
- It was the earliest form of regression, which was published by **Legendre** in 1805 and by **Gauss** in 1809.
- They both applied the method to the problem of determining the orbits of bodies around the Sun from astronomical observations.
- The term regression was only coined later by **Galton** to describe the biological phenomenon that the heights of descendants of tall ancestors tend to regress down towards a normal average (a phenomenon also known as regression toward the mean).
- **Pearson** later provided the statistical context showing that the phenomenon is more general than a biological context.

Multivariate ordinary least squares

By describing the input data as $\mathbf{X} \in \mathbb{R}^{N \times D^I}$ and the output data as $\mathbf{Y} \in \mathbb{R}^{N \times D^O}$, we want to find $\mathbf{A} \in \mathbb{R}^{D^I \times D^O}$ such that $\mathbf{Y} = \mathbf{X}\mathbf{A}$.

A solution can be found by minimizing the L_2 -norm of the residuals

$$\begin{aligned}\hat{\mathbf{A}} &= \arg \min_{\mathbf{A}} \|\mathbf{Y} - \mathbf{X}\mathbf{A}\|_2 \\ &= \arg \min_{\mathbf{A}} (\mathbf{Y} - \mathbf{X}\mathbf{A})^\top (\mathbf{Y} - \mathbf{X}\mathbf{A}) \\ &= \arg \min_{\mathbf{A}} (\mathbf{Y}^\top \mathbf{Y} - 2\mathbf{A}^\top \mathbf{X}^\top \mathbf{Y} + \mathbf{A}^\top \mathbf{X}^\top \mathbf{X}\mathbf{A})\end{aligned}\quad \text{Moore-Penrose pseudoinverse}$$

By differentiating with respect to \mathbf{A} and equating to zero

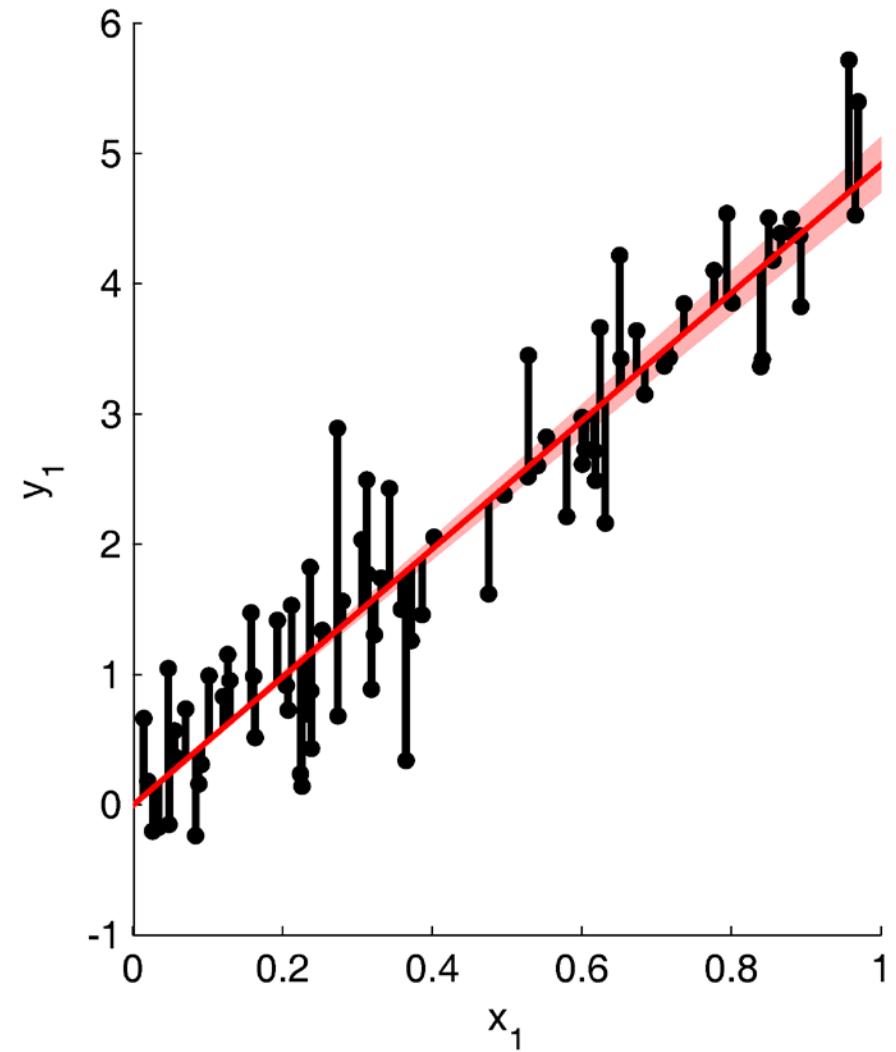
$$-2\mathbf{X}^\top \mathbf{Y} + 2\mathbf{X}^\top \mathbf{X}\mathbf{A} = \mathbf{0} \iff \hat{\mathbf{A}} = \boxed{(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Y}}$$

or $\hat{\mathbf{A}} = \boxed{\mathbf{X}^\top (\mathbf{X}\mathbf{X}^\top)^{-1}} \mathbf{Y}$

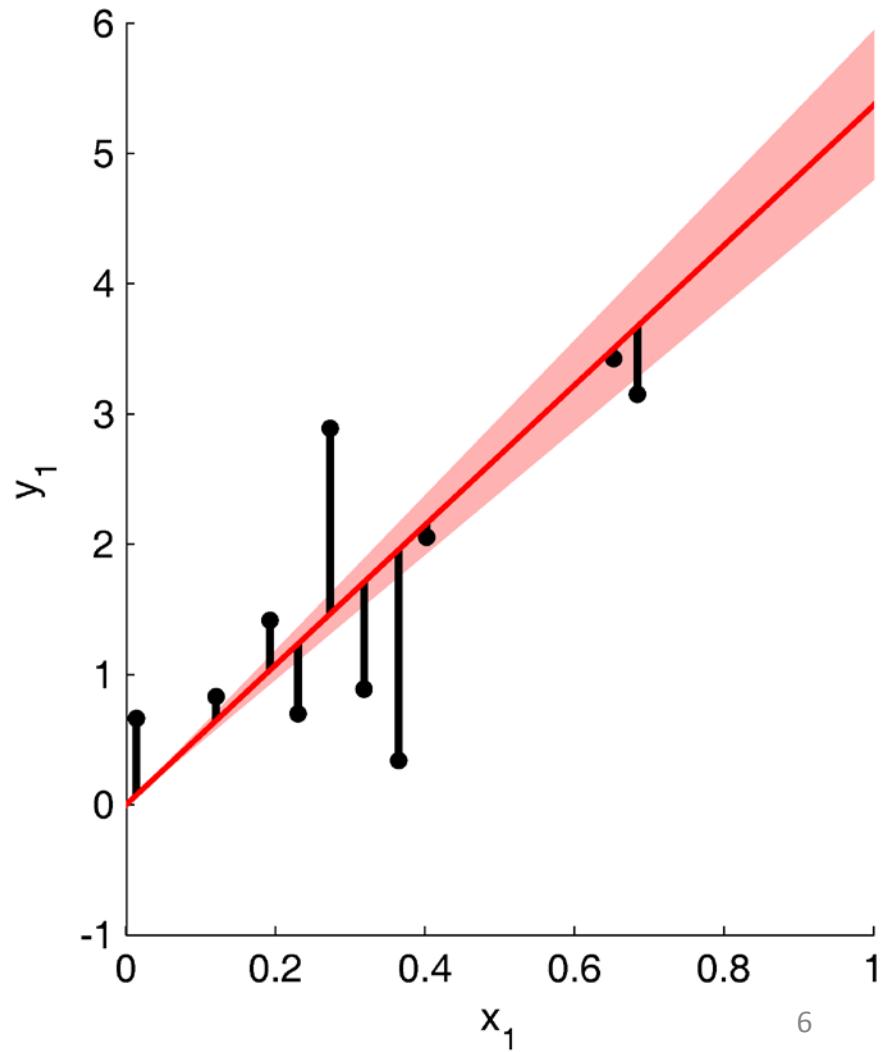
with residuals (parameters errors) given by $\hat{\Sigma}^{\mathbf{A}} = (\mathbf{X}^\top \mathbf{X})^{-1}$

Multivariate ordinary least squares

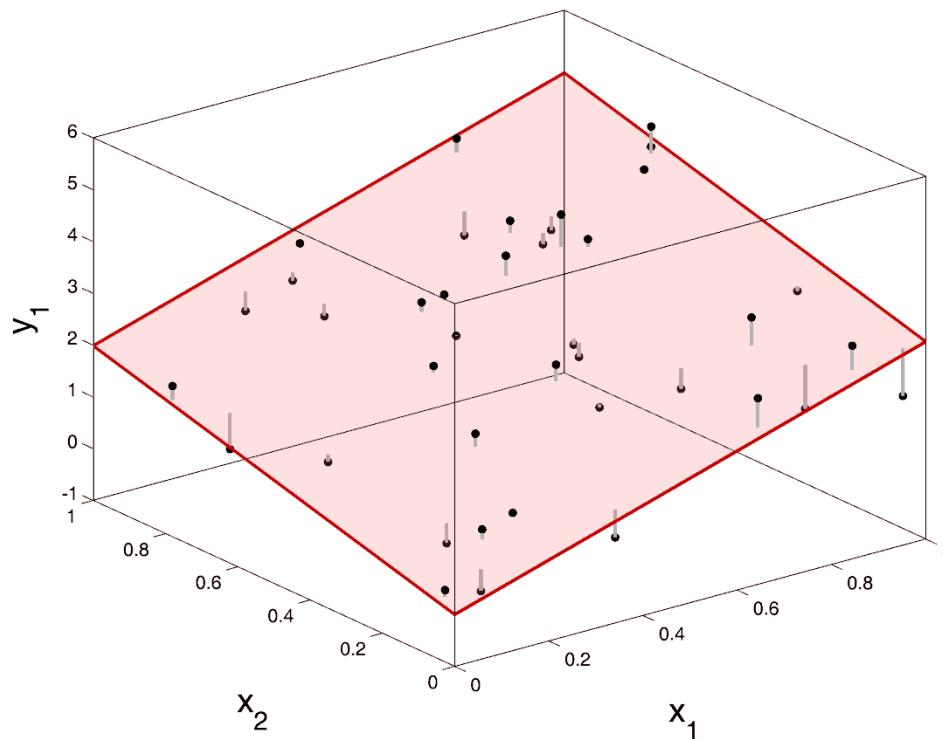
$D^{\mathcal{I}}=1, D^{\mathcal{O}}=1, N=80$



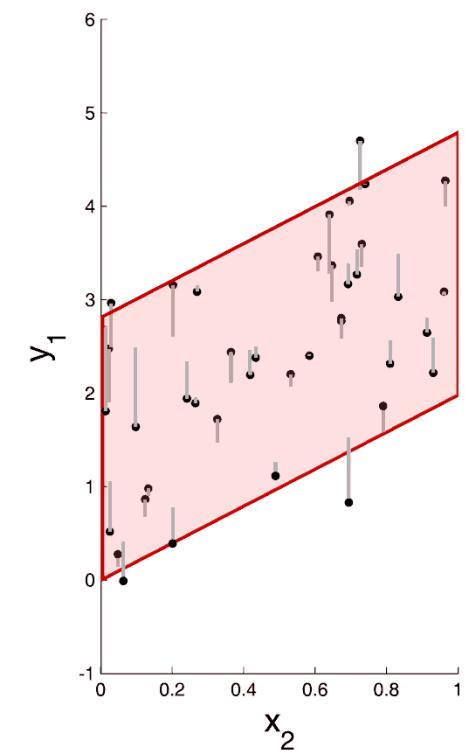
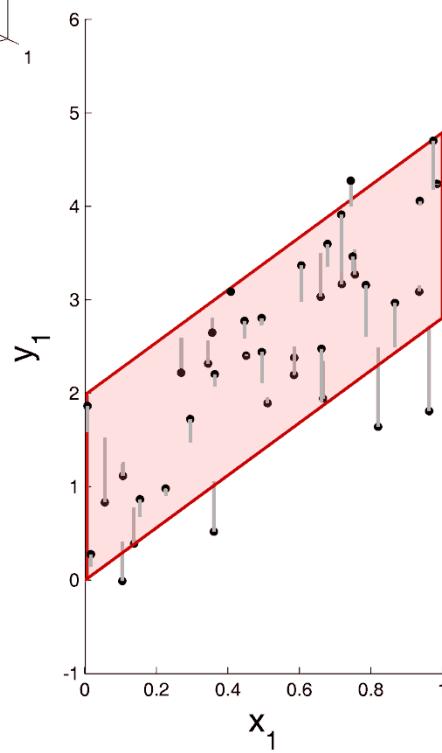
$D^{\mathcal{I}}=1, D^{\mathcal{O}}=1, N=10$



Multivariate ordinary least squares



$$D^{\mathcal{I}} = 2, D^{\mathcal{O}} = 1, N = 40$$



Multivariate ordinary least squares

The result minimizes the sum of squared residuals

$$(\mathbf{Y} - \mathbf{X}\mathbf{A})^\top (\mathbf{Y} - \mathbf{X}\mathbf{A}) = \sum_{t=1}^N (\mathbf{Y}_t - \mathbf{X}_t\mathbf{A})^\top (\mathbf{Y}_t - \mathbf{X}_t\mathbf{A})$$

It is a **frequentist approach**: it assumes that there are enough measurements to say something meaningful about \mathbf{A} .

In the Bayesian approach, the parameters are supplemented with additional information in the form of a prior probability distribution. The prior belief about the parameters is combined with the data likelihood function to yield a posterior belief about \mathbf{A} .

In the computation of the least squares solution, the matrix $\mathbf{X}^\top \mathbf{X}$ in the pseudoinverse is known as the **Gramian matrix** of \mathbf{X} .

It is usually inefficient to invert this Gramian matrix directly, but it can be decomposed in ways that either speed up computation or provide numerically more stable estimates.

Least squares with Cholesky decomposition

If the matrix $\mathbf{S} = \mathbf{X}^\top \mathbf{X}$ is positive definite and well-conditioned (\rightarrow full rank), the problem can be solved directly by using the **Cholesky decomposition** $\mathbf{S} = \mathbf{R}^\top \mathbf{R}$, where \mathbf{R} is an upper triangular matrix

$$\mathbf{R}^\top \mathbf{R} \mathbf{A} = \mathbf{X}^\top \mathbf{Y}$$

The solution is obtained in two stages, a forward substitution step solving for \mathbf{Z}

$$\mathbf{R}^\top \mathbf{Z} = \mathbf{X}^\top \mathbf{Y}$$

followed by a backward substitution solving for \mathbf{A}

$$\mathbf{R} \mathbf{A} = \mathbf{Z}$$

where both substitutions can exploit the triangular nature of \mathbf{R} .

When using *Matlab*, $\hat{\mathbf{A}}$ and $\hat{\Sigma}^{\mathbf{A}}$ can for example be computed with the `lscov` function.

Singular value decomposition (SVD)

$$X \in \mathbb{R}^{N \times D^I} = U \in \mathbb{R}^{N \times N} \cdot \Sigma \in \mathbb{R}^{N \times D^I} \cdot V^\top \in \mathbb{R}^{D^I \times D^I}$$

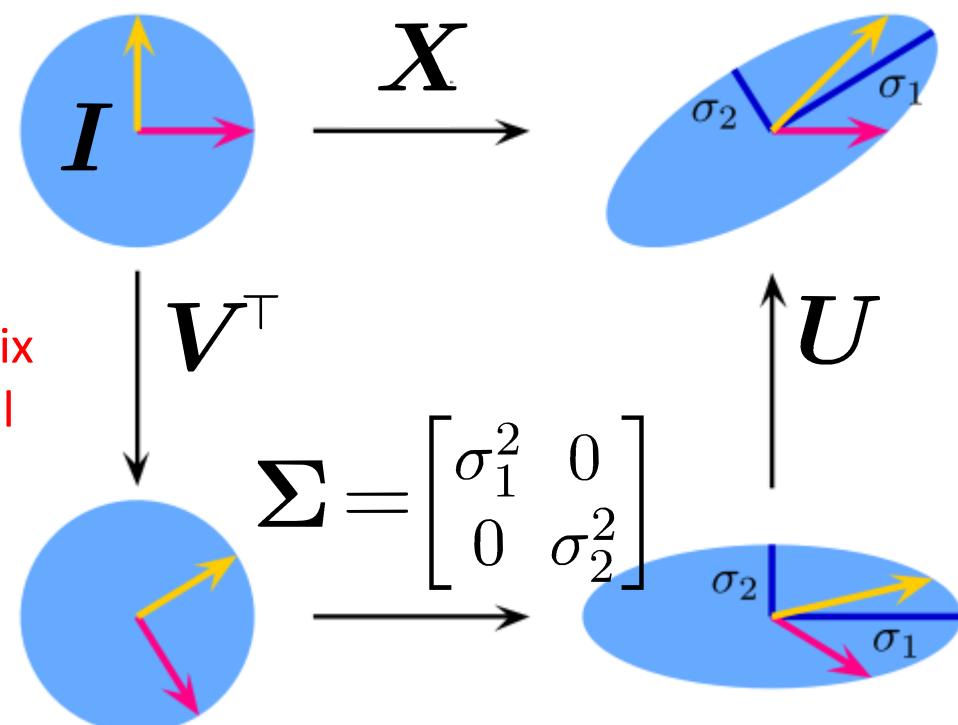
$$\begin{bmatrix} 1 & 0 & 0 & 0 & 2 \\ 0 & 0 & 3 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \\ 0 & 4 & 0 & 0 & 0 \end{bmatrix} = \begin{bmatrix} 0 & 0 & 1 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & -1 \\ 1 & 0 & 0 & 0 \end{bmatrix} \cdot \begin{bmatrix} 4 & 0 & 0 & 0 & 0 \\ 0 & 3 & 0 & 0 & 0 \\ 0 & 0 & \sqrt{5} & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \end{bmatrix} \cdot \begin{bmatrix} 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 \\ \sqrt{0.2} & 0 & 0 & 0 & \sqrt{0.8} \\ 0 & 0 & 0 & 1 & 0 \\ -\sqrt{0.8} & 0 & 0 & 0 & \sqrt{0.2} \end{bmatrix}$$

Matrix with non-negative
diagonal entries
(singular values of X)

Unitary matrix
(orthonormal
bases)

$$X = U \Sigma V^\top$$

Unitary matrix
(orthonormal
bases)



Singular value decomposition (SVD)

SVD provides an explicit representation of the **range** and **null space** of a matrix \mathbf{X} .

The right-singular vectors corresponding to vanishing singular values of \mathbf{X} span the **null space** of \mathbf{X} .

The left-singular vectors corresponding to the non-zero singular values of \mathbf{X} span the **range** of \mathbf{X} .

As a consequence, the **rank** of \mathbf{X} equals the number of non-zero singular values which is the same as the number of non-zero diagonal elements in Σ .

In numerical linear algebra, the singular values can be used to determine the effective **rank of a matrix**, as rounding error may lead to small but non-zero singular values in a rank deficient matrix.

Singular value decomposition (SVD)

SVD is very general in the sense that it can be applied to any $m \times n$ matrix whereas **eigenvalue decomposition** can only be applied to certain classes of square matrices. Nevertheless, the two decompositions are related. Given an SVD of \mathbf{X} , we have

$$\begin{aligned}\mathbf{X}^\top \mathbf{X} &= \mathbf{V} \boldsymbol{\Sigma}^\top \mathbf{U}^\top \mathbf{U} \boldsymbol{\Sigma} \mathbf{V}^\top = \mathbf{V} (\boldsymbol{\Sigma}^\top \boldsymbol{\Sigma}) \mathbf{V}^\top \\ \mathbf{X} \mathbf{X}^\top &= \mathbf{U} \boldsymbol{\Sigma} \mathbf{V}^\top \mathbf{V} \boldsymbol{\Sigma}^\top \mathbf{U}^\top = \mathbf{U} (\boldsymbol{\Sigma} \boldsymbol{\Sigma}^\top) \mathbf{U}^\top\end{aligned}$$

The right-hand sides of these relations describe the eigenvalue decompositions of the left-hand sides. Consequently:

- Columns of \mathbf{V} (right-singular vec.) are eigenvectors of $\mathbf{X}^\top \mathbf{X}$.
- Columns of \mathbf{U} (left-singular vec.) are eigenvectors of $\mathbf{X} \mathbf{X}^\top$.
- The non-zero elements of $\boldsymbol{\Sigma}$ (non-zero singular values) are the square roots of the non-zero eigenvalues of $\mathbf{X}^\top \mathbf{X}$ or $\mathbf{X} \mathbf{X}^\top$.

Singular value decomposition (SVD)

- Applications employing SVD include pseudoinverse computation, least squares fitting, multivariable control, matrix approximation, as well as the determination of the **rank, range and null space of a matrix**.
- The SVD can also be thought of as decomposing a matrix into a weighted, ordered sum of separable matrices (e.g., decomposition of an image processing filter into separable horizontal and vertical filters).
- It is possible to use the SVD of a square matrix \mathbf{A} to determine the orthogonal matrix \mathbf{O} closest to \mathbf{A} . The closeness of fit is measured by the Frobenius norm of $\mathbf{O}-\mathbf{A}$. The solution is the product $\mathbf{U}\mathbf{V}^T$.
- A similar problem, with interesting applications in shape analysis, is the orthogonal **Procrustes problem**, which consists of finding an orthogonal matrix \mathbf{O} which most closely maps \mathbf{A} to \mathbf{B} .
- Extensions to higher order arrays exist, generalizing SVD to a multi-way analysis of the data (multilinear algebra).

Condition number of a matrix with SVD

Using the pseudoinverse and a matrix norm, one can define the **condition number** of a matrix \mathbf{X} as

$$\text{cond}(\mathbf{X}) = \|\mathbf{X}\| \|\mathbf{X}^\dagger\|$$

A large condition number implies that the problem of finding least-squares solutions to the corresponding system of linear equations is **ill-conditioned** in the sense that small errors in the entries of \mathbf{X} can lead to huge errors in the entries of the solution.

The condition number can also be computed from the singular value decomposition $\mathbf{X} = \mathbf{U}\Sigma\mathbf{V}^\top$, where σ_i are the singular values in the diagonal of Σ , with

$$\text{cond}(\mathbf{X}) = \frac{\max_i \sigma_i}{\min_i \sigma_i}$$

Least squares with SVD

\mathbf{X} can be decomposed with the **singular value decomposition**

$$\mathbf{X} = \mathbf{U}\Sigma\mathbf{V}^\top$$

where \mathbf{U} and \mathbf{V} are $m \times m$ and $n \times n$ orthogonal matrices, and Σ is an $m \times n$ matrix with all its elements outside of the main diagonal equal to 0. With this decomposition, a solution to the least squares problem is given by

$$\hat{\mathbf{A}} = \mathbf{V}\Sigma^\dagger\mathbf{U}^\top \mathbf{Y}$$

where the pseudoinverse of Σ can be easily obtained by inverting the non-zero diagonal elements and transposing the resulting matrix. This can for example be computed in Matlab/GNU Octave with

```
[U,S,V] = svd(X)
S(S>0) = S(S>0).^-1
A = V*S'*U' * Y
```

Least squares with SVD

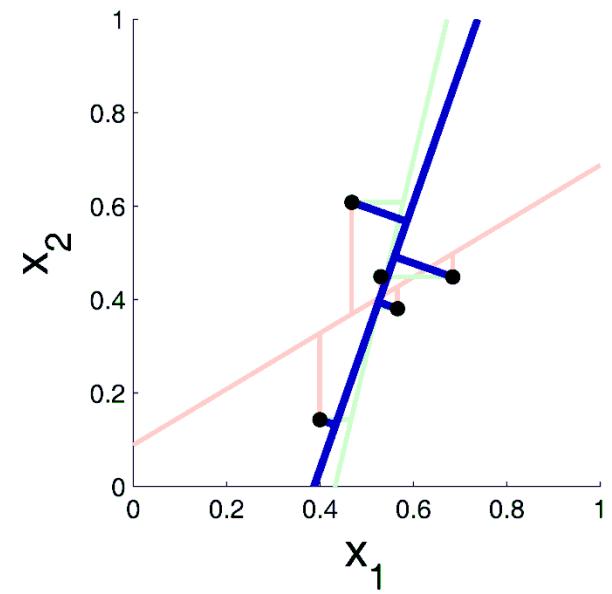
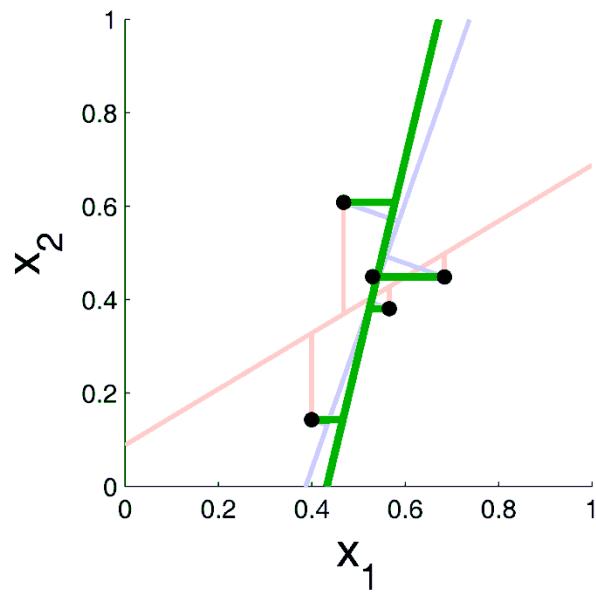
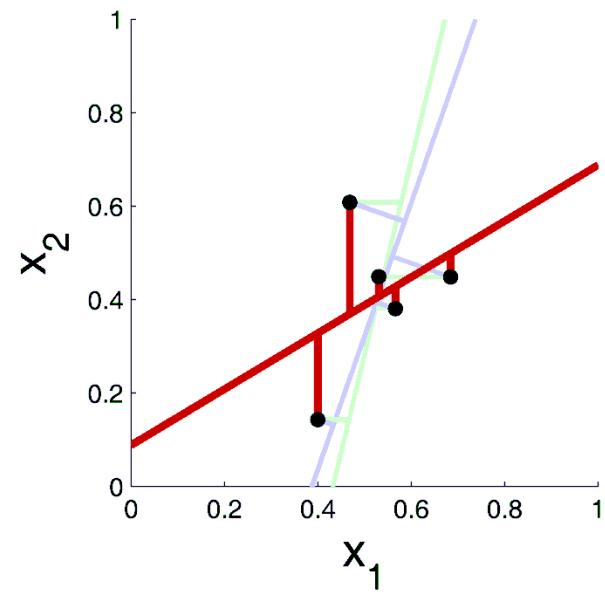
This method is computationally intensive (dominated by the cost of computing the SVD).

However, it is particularly useful if $\mathbf{X}^\top \mathbf{X}$ is **ill-conditioned** (i.e. if its condition number multiplied by the computer relative round-off error is large).

In that case, including the smallest singular values in the inversion merely adds numerical noise to the solution.

This can be cured with the truncated SVD approach, giving a more stable and exact answer, by explicitly setting to zero all singular values below a certain threshold and so ignoring them (a process closely related to factor analysis).

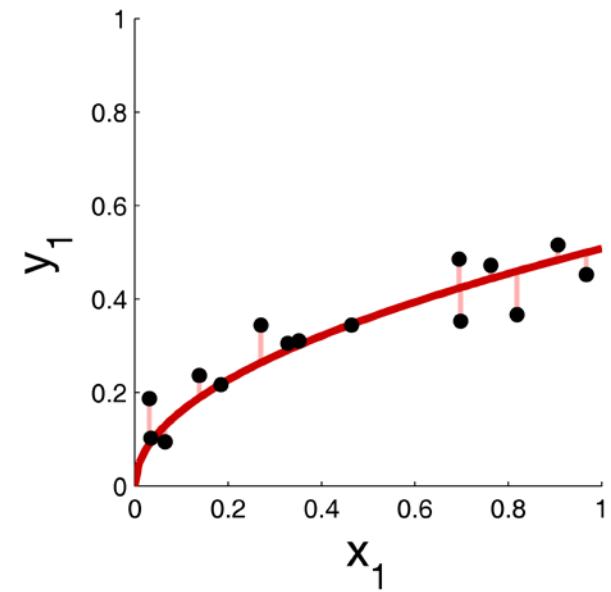
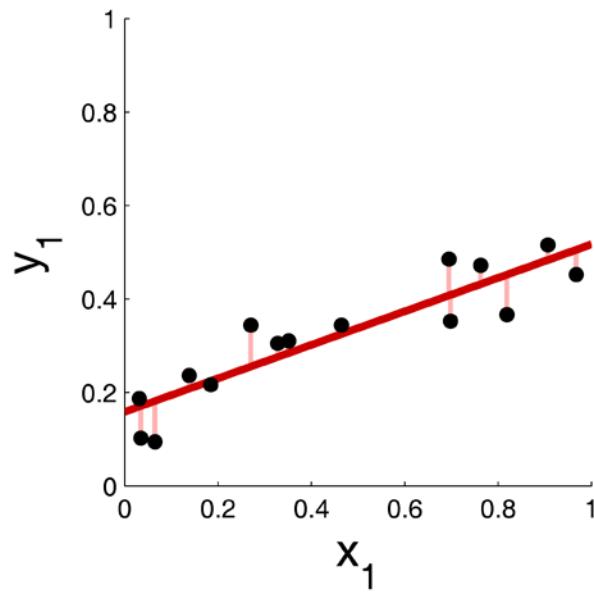
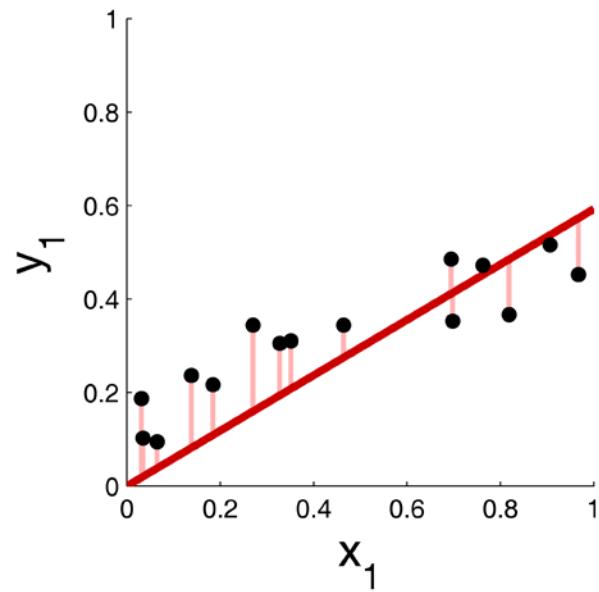
Different line fitting alternatives



Data fitting with linear least squares

- In linear least squares, we are not restricted to using a line as the model as in the previous slide!
- For instance, we could have chosen a quadratic model $Y=X^2A$: This model is still linear in the A parameter.
 - **The function does not need to be linear in the argument:** only in the parameters that are determined to give the best fit.
- Also, not all of X contains information about the datapoints: the first/last column can for example be populated with ones, so that an offset is learned
- This can be used for **polynomial fitting** by treating x, x^2, \dots as being distinct independent variables in a multiple regression model.

Data fitting with linear least squares



$$\mathbf{X} = \begin{bmatrix} x_{1,1} & \cdots & x_{1,D^{\mathcal{I}}} \\ \vdots & \ddots & \vdots \\ x_{N,1} & \cdots & x_{N,D^{\mathcal{I}}} \end{bmatrix}$$

$$= \mathbf{x}$$

$$\mathbf{X} = \begin{bmatrix} x_{1,1} & \cdots & x_{1,D^{\mathcal{I}}} & 1 \\ \vdots & \ddots & \vdots & \vdots \\ x_{N,1} & \cdots & x_{N,D^{\mathcal{I}}} & 1 \end{bmatrix}$$

$$= [\mathbf{x}, 1]$$

$$\mathbf{X} = \begin{bmatrix} \sqrt{x_{1,1}} & \cdots & \sqrt{x_{1,D^{\mathcal{I}}}} \\ \vdots & \ddots & \vdots \\ \sqrt{x_{N,1}} & \cdots & \sqrt{x_{N,D^{\mathcal{I}}}} \end{bmatrix}$$

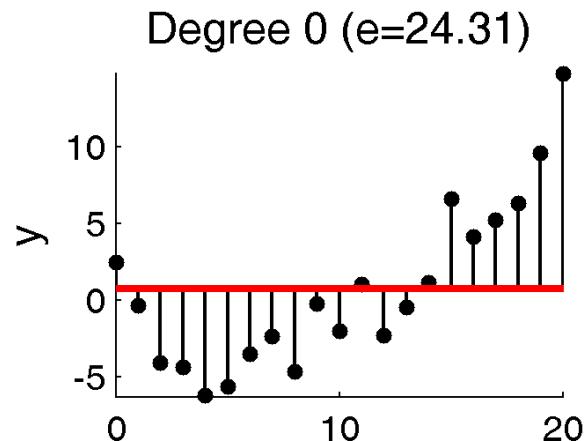
$$= \sqrt{\mathbf{x}}$$

Data fitting with linear least squares

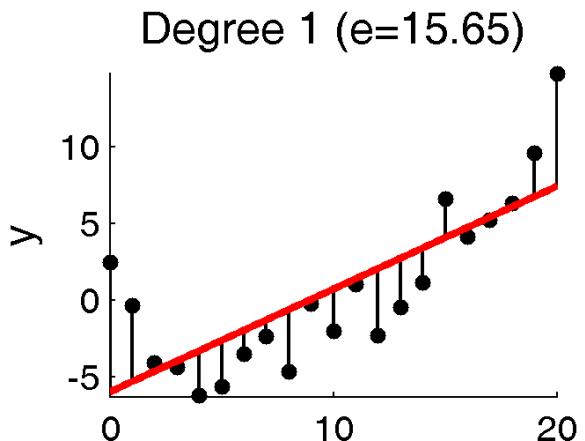
- Polynomial regression is an example of regression analysis using **basis functions** to model a functional relationship between two quantities. Specifically, it replaces x in linear regression with polynomial basis $[1, x, x^2, \dots, x^d]$.
- A drawback of polynomial bases is that the basis functions are "non-local".
- It is for this reason that the polynomial basis functions are often used along with other forms of basis functions, such as splines, radial basis functions, and wavelets.

(We will learn more about this at next lecture...)

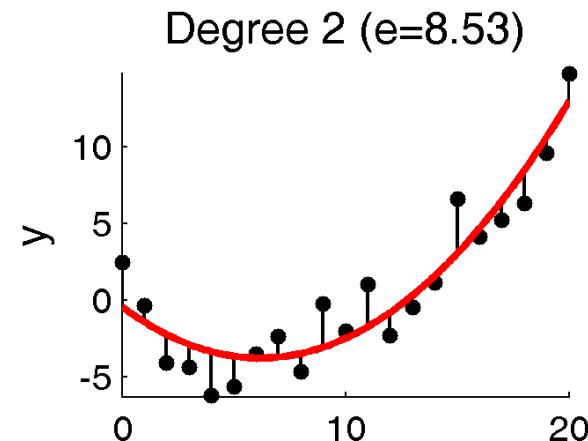
Polynomial fitting with least squares



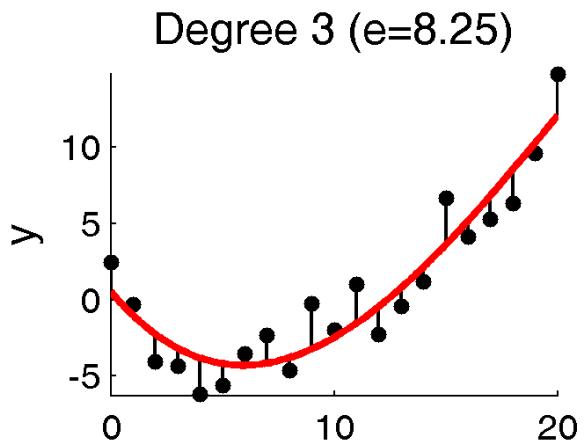
$$X = 1$$



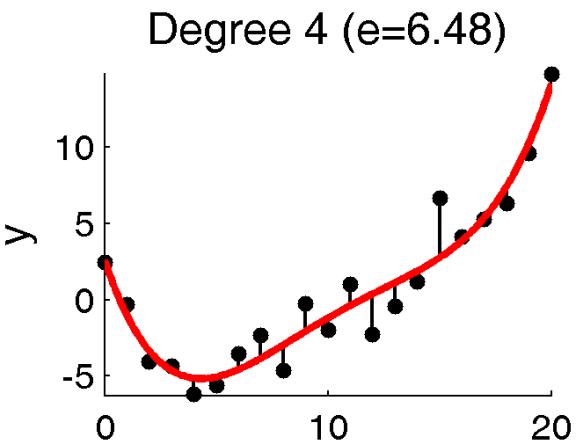
$$X = [x, 1]$$



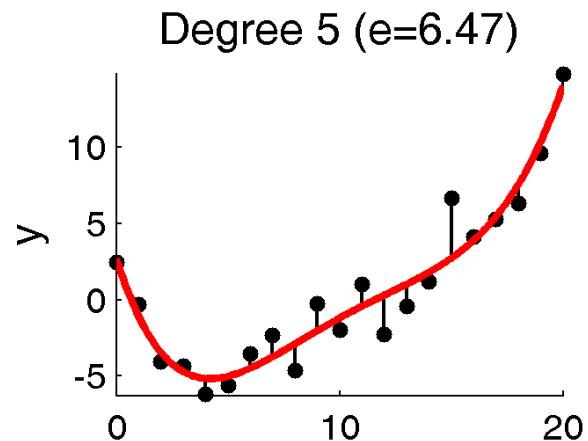
$$X = [x^2, x, 1]$$



$$X = [x^3, x^2, x, 1]$$



$$X = [x^4, x^3, \dots, 1]$$



$$X = [x^5, x^4, \dots, 1]$$

Kernels in least squares (nullspace projection)

demo_LS_nullspace01.m

Kernels in least squares (nullspace)

While the pseudoinverse provides a single least norm solution, we can still obtain all solutions of the linear system by employing a **nullspace projection operator**

$$\hat{\mathbf{A}} = \mathbf{X}^\dagger \mathbf{Y} + \overbrace{(\mathbf{I} - \mathbf{X}^\dagger \mathbf{X})}^{\mathbf{N}} \mathbf{V}$$

In the above, the solution is unique if and only if \mathbf{A} has full column rank, in which case \mathbf{N} is a zero matrix.

\mathbf{V} can be any vector/matrix (e.g., resulting from the minimization of a secondary objective function).

The nullspace projection guarantees that $\|\mathbf{Y} - \mathbf{X}\hat{\mathbf{A}}\|$ is still minimized.

Kernels in least squares (nullspace)

Indeed, we can see that for any vector/matrix \mathbf{V} , the nullspace projection ensures that it does not affect the minimization

$$\begin{aligned} \mathbf{E} &= \mathbf{Y} - \mathbf{X} \hat{\mathbf{A}} \\ &= \mathbf{Y} - \mathbf{X} \left(\mathbf{X}^\top (\mathbf{X} \mathbf{X}^\top)^{-1} \mathbf{Y} + \overbrace{(\mathbf{I} - \mathbf{X}^\top (\mathbf{X} \mathbf{X}^\top)^{-1} \mathbf{X})}^N \mathbf{V} \right) \\ &= \mathbf{Y} - \mathbf{X} \mathbf{X}^\top (\mathbf{X} \mathbf{X}^\top)^{-1} \mathbf{Y} + \mathbf{X} (\mathbf{I} - \mathbf{X}^\top (\mathbf{X} \mathbf{X}^\top)^{-1} \mathbf{X}) \mathbf{V} \\ &= \mathbf{Y} - \mathbf{Y} + (\mathbf{X} - \mathbf{X} \mathbf{X}^\top (\mathbf{X} \mathbf{X}^\top)^{-1} \mathbf{X}) \mathbf{V} \\ &= (\mathbf{X} - \mathbf{X}) \mathbf{V} \\ &= \mathbf{0} \mathbf{V} \end{aligned}$$

Kernels in least squares (nullspace)

An alternative way of computing the nullspace projection matrix is to exploit the singular value decomposition

$$\mathbf{X}^\dagger = \mathbf{U} \boldsymbol{\Sigma} \mathbf{V}^\top$$

to compute

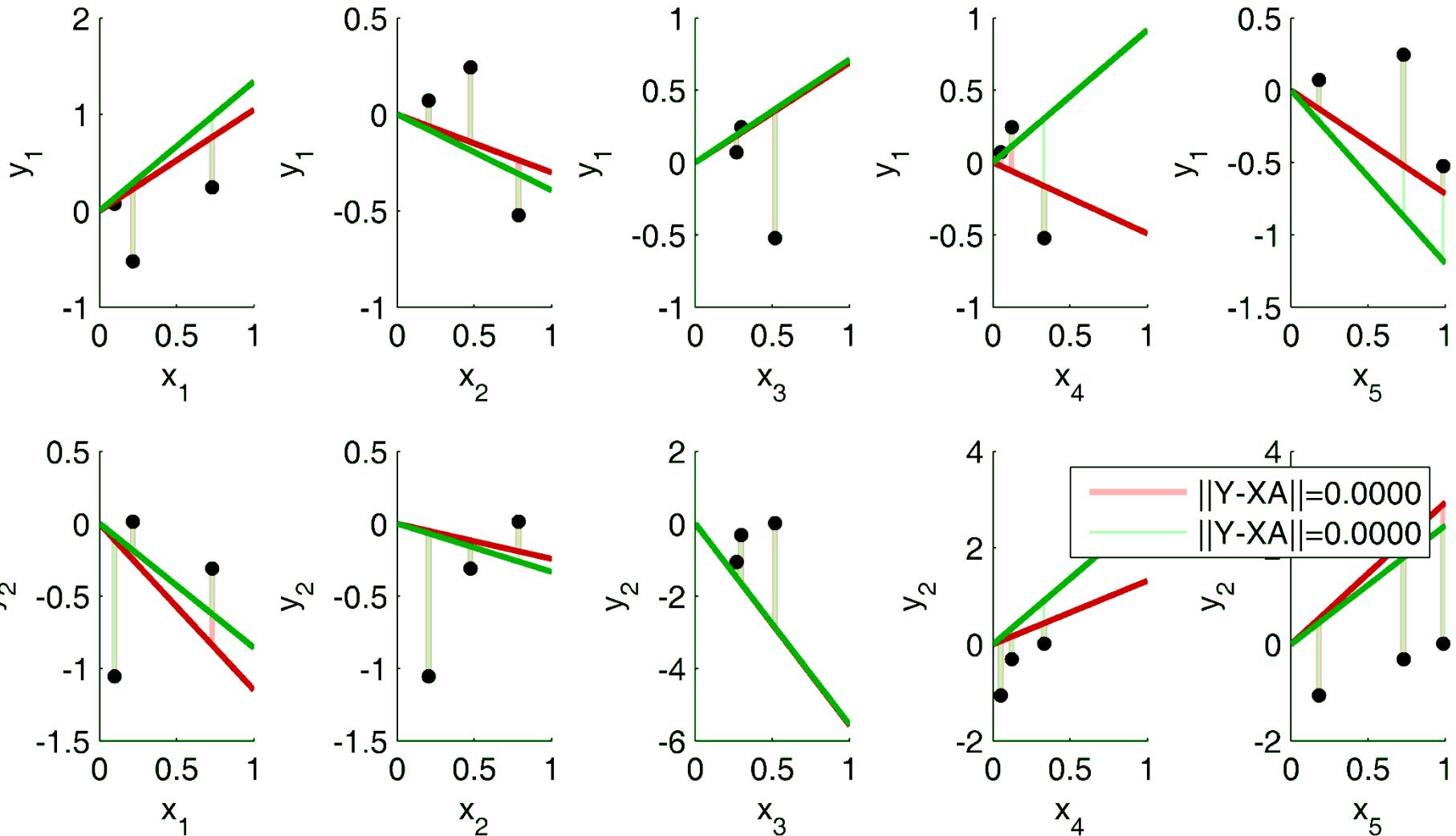
$$\mathbf{N} = \tilde{\mathbf{U}} \tilde{\mathbf{U}}^\top$$

where $\tilde{\mathbf{U}}$ is a matrix formed by the columns of \mathbf{U} that span for the corresponding zero rows in $\boldsymbol{\Sigma}$.

This can for example be implemented in Matlab/Octave with

```
[U,S,V] = svd(pinv(X))
sp = sum(S,2) < 1E-1
N = U(:,sp) * U(:,sp)'
```

Kernels in least squares (nullspace)

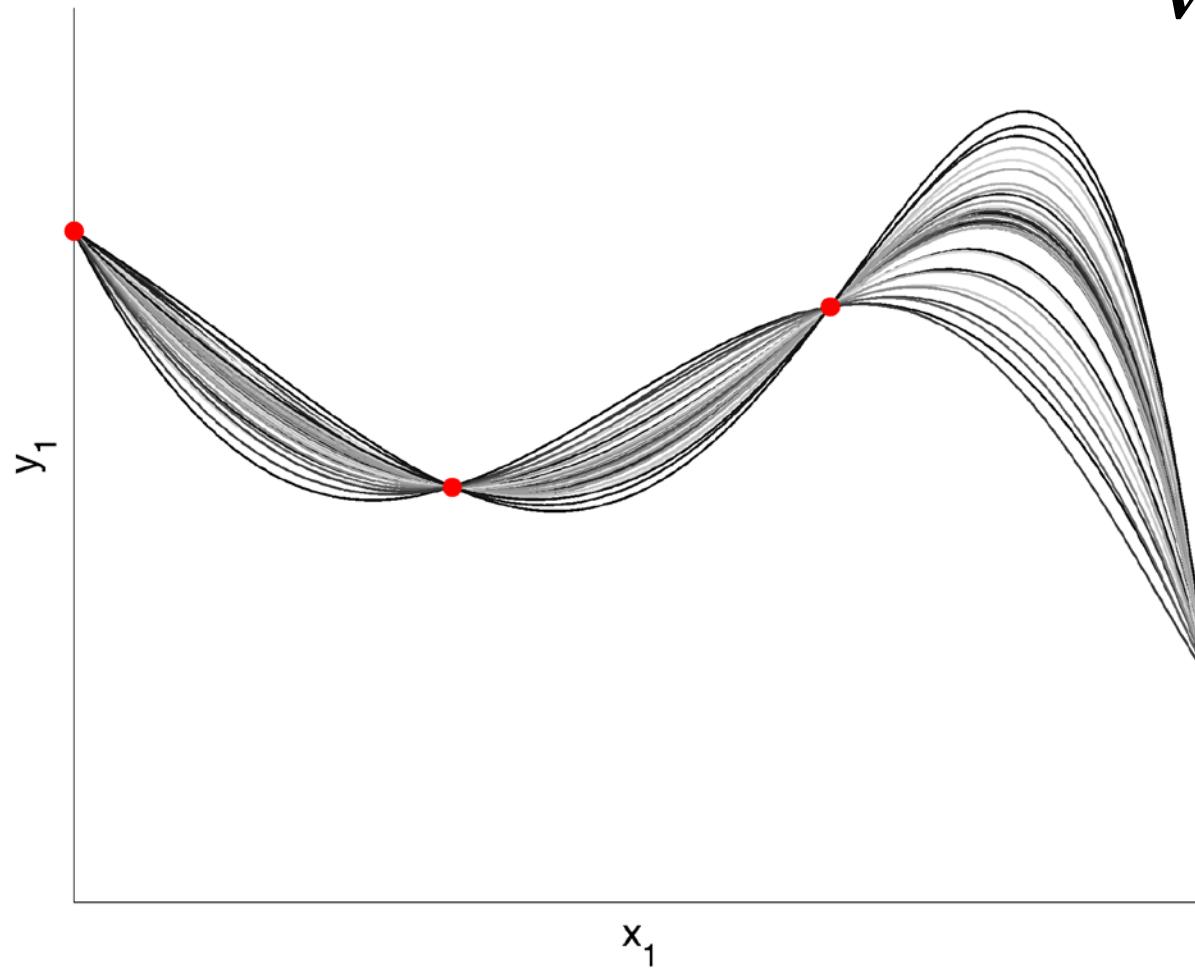


— without nullspace

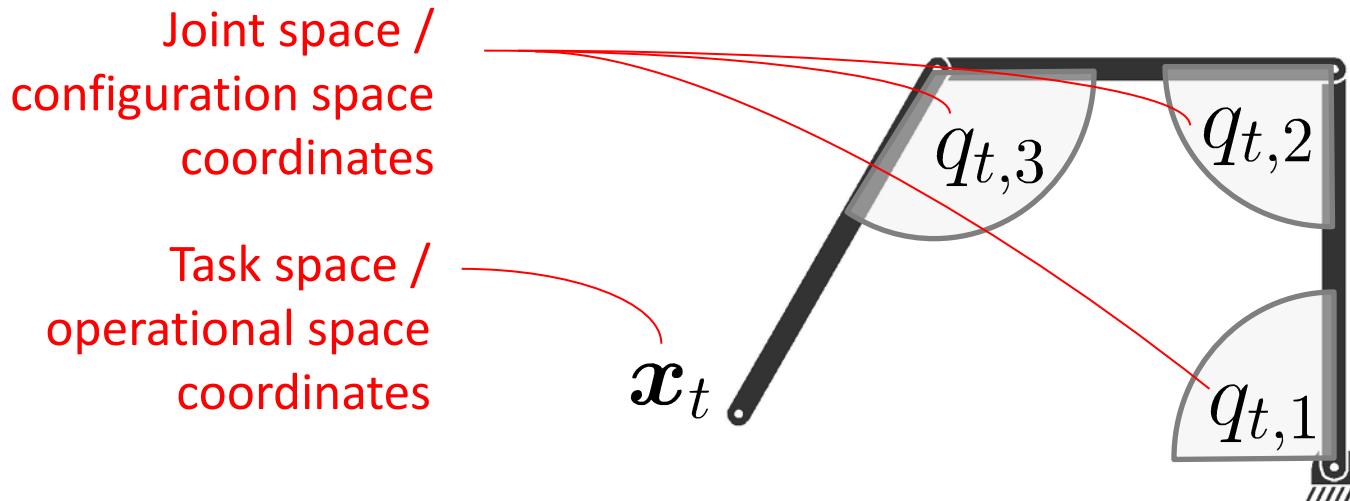
— with nullspace

Example with polynomial fitting

$$\hat{A} = X^\dagger Y + NV \quad \text{with} \quad X = [x^6, x^5, \dots, 1] \\ V \sim \mathcal{N}(0, I)$$



Example with robot inverse kinematics



Forward kinematics is computed with

$$\mathbf{x}_t = f(\mathbf{q}_t) \iff \dot{\mathbf{x}}_t = \frac{\partial \mathbf{x}_t}{\partial t} = \frac{\partial f(\mathbf{q}_t)}{\partial \mathbf{q}_t} \frac{\partial \mathbf{q}_t}{\partial t} = \mathbf{J}(\mathbf{q}_t) \dot{\mathbf{q}}_t$$

where $\mathbf{J}(\mathbf{q}_t) = \frac{\partial f(\mathbf{q}_t)}{\partial \mathbf{q}_t}$ is a Jacobian matrix.

A least norm inverse kinematics solution can be computed with

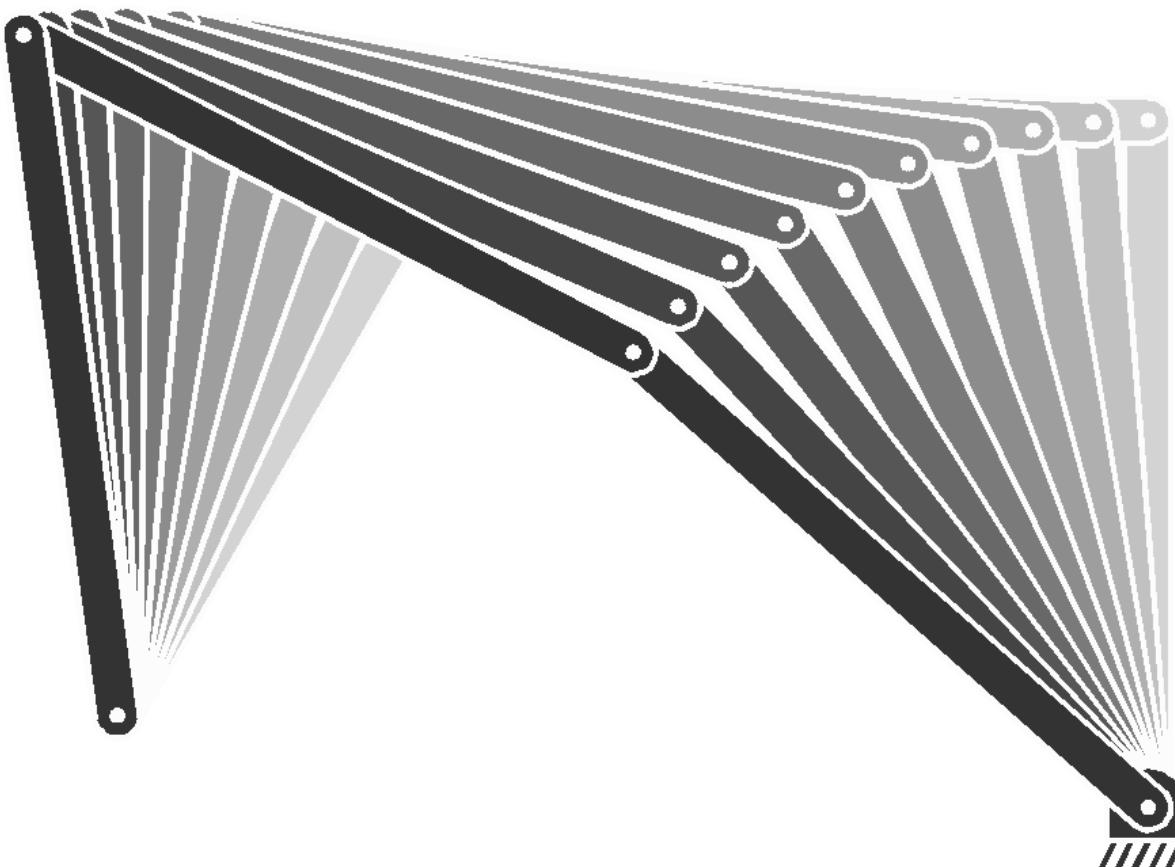
$$\hat{\dot{\mathbf{q}}}_t = \mathbf{J}^\dagger(\mathbf{q}_t) \dot{\mathbf{x}}_t + \mathbf{N}(\mathbf{q}_t) g(\mathbf{q}_t)$$

Example with robot inverse kinematics

$$\hat{\dot{q}}_t = \mathbf{J}^\dagger(\mathbf{q}_t) \dot{x}_t + \mathbf{N}(\mathbf{q}_t) g(\mathbf{q}_t)$$

→ Keeping the tip still
as primary constraint

$$= \mathbf{J}^\dagger(\mathbf{q}_t) \begin{bmatrix} 0 \\ 0 \end{bmatrix} + \mathbf{N}(\mathbf{q}_t) \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}$$



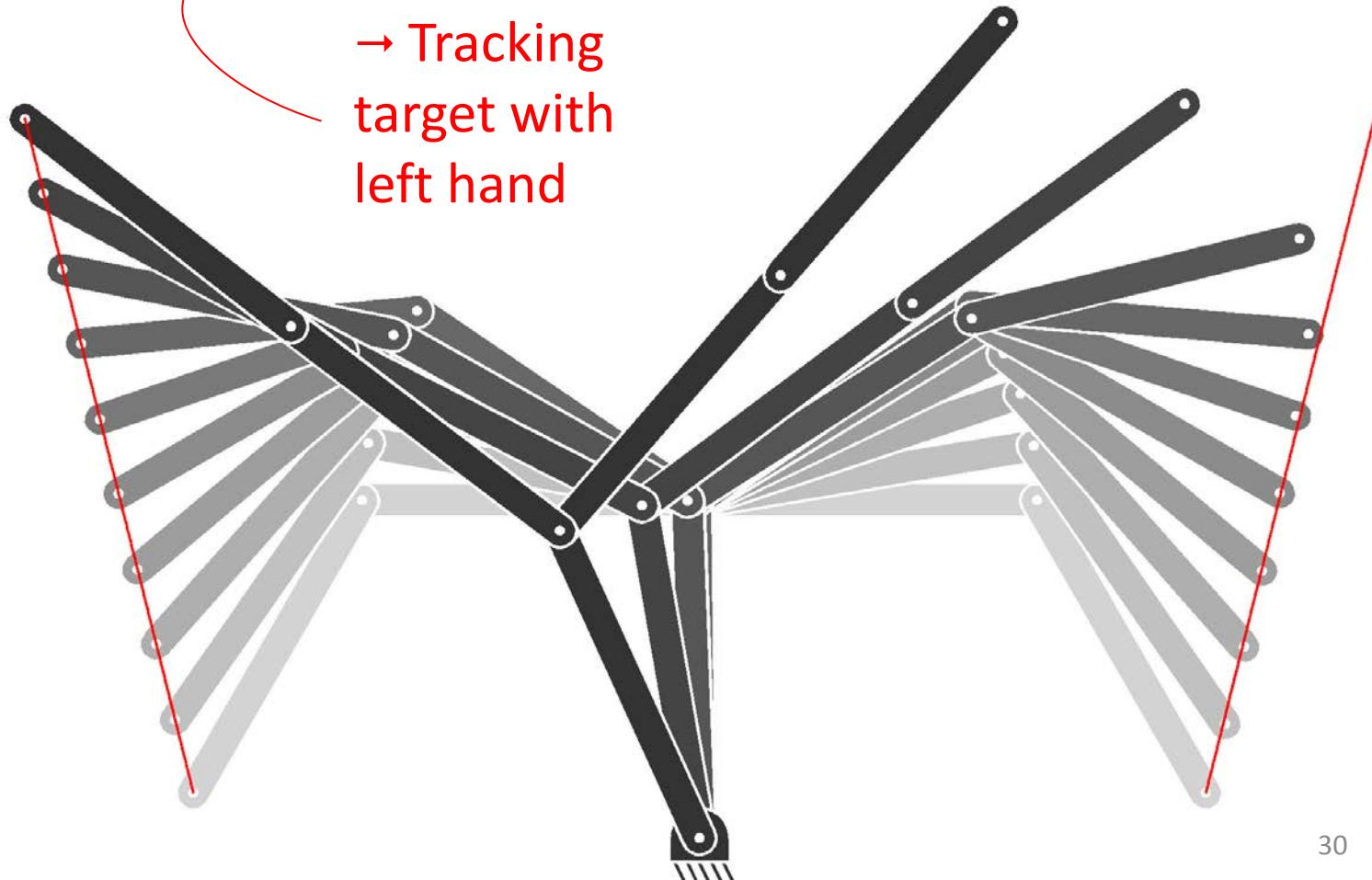
→ Trying to move
the first joint as
secondary constraint

Example with robot inverse kinematics

$$\begin{aligned}\hat{\dot{q}}_t &= \mathbf{J}^{\mathcal{L}\dagger} \dot{x}_t^{\mathcal{L}} + \mathbf{N}^{\mathcal{L}} \mathbf{J}^{\mathcal{R}\dagger} \dot{x}_t^{\mathcal{R}} \\ &= \mathbf{J}^{\mathcal{L}\dagger} (\hat{x}_t^{\mathcal{L}} - x_t^{\mathcal{L}}) + \mathbf{N}^{\mathcal{L}} \mathbf{J}^{\mathcal{R}\dagger} (\hat{x}_t^{\mathcal{R}} - x_t^{\mathcal{R}})\end{aligned}$$

→ Tracking target
with right hand if
possible

→ Tracking
target with
left hand



Ridge regression (Tikhonov regularization, penalized least squares)

demo_LS_polFit02.m

Ridge regression (Tikhonov regularization)

The classical frequentist linear least squares solution is to estimate the matrix of regression coefficients using the Moore-Penrose pseudoinverse.

In the Bayesian approach, additional information are given to the parameters in the form of a **prior probability distribution**.

The prior can take different functional forms depending on the domain and the information that is available a priori.

After specifying the prior, the posterior distribution can be expressed by using the prior mean μ , with the strength of the prior indicated by the prior precision matrix Γ .

The special case with $\mu = \mathbf{0}$ and $\Gamma = \lambda I$ is called **ridge regression**.

Ridge regression (Tikhonov regularization)

The least squares objective can be modified to give preference to a particular solution with

$$\begin{aligned}\hat{\mathbf{A}} &= \arg \min_{\mathbf{A}} \|\mathbf{Y} - \mathbf{X}\mathbf{A}\|_2 + \|\boldsymbol{\Gamma}\mathbf{A}\|_2 \\ &= \arg \min_{\mathbf{A}} (\mathbf{Y} - \mathbf{X}\mathbf{A})^\top(\mathbf{Y} - \mathbf{X}\mathbf{A}) + (\boldsymbol{\Gamma}\mathbf{A})^\top\boldsymbol{\Gamma}\mathbf{A} \\ &= \arg \min_{\mathbf{A}} (\mathbf{Y}^\top\mathbf{Y} - 2\mathbf{A}^\top\mathbf{X}^\top\mathbf{Y} + \mathbf{A}^\top\mathbf{X}^\top\mathbf{X}\mathbf{A} + \mathbf{A}^\top\boldsymbol{\Gamma}^\top\boldsymbol{\Gamma}\mathbf{A})\end{aligned}$$

By differentiating with respect to \mathbf{A} and equating to zero, we can see that

$$-2\mathbf{X}^\top\mathbf{Y} + 2\mathbf{X}^\top\mathbf{X}\mathbf{A} + 2\boldsymbol{\Gamma}^\top\boldsymbol{\Gamma}\mathbf{A} = \mathbf{0}$$

yielding

$$\hat{\mathbf{A}} = (\mathbf{X}^\top\mathbf{X} + \boldsymbol{\Gamma}^\top\boldsymbol{\Gamma})^{-1}\mathbf{X}^\top\mathbf{Y}$$

If $\boldsymbol{\Gamma} = \lambda\mathbf{I}$ with $\lambda \ll 1$ (i.e., giving preference to solutions with smaller norms), the process is known as **L_2 regularization**.

Ridge regression (Tikhonov regularization)

Ridge regression can alternatively be computed with augmented matrices

$$\tilde{\mathbf{X}} = \begin{bmatrix} \mathbf{X} \\ \boldsymbol{\Gamma} \end{bmatrix} \quad \tilde{\mathbf{Y}} = \begin{bmatrix} \mathbf{Y} \\ \mathbf{0} \end{bmatrix}$$

with $\mathbf{0} \in \mathbb{R}^{D^{\mathcal{I}} \times D^{\mathcal{O}}}$ and $\boldsymbol{\Gamma} \in \mathbb{R}^{D^{\mathcal{I}} \times D^{\mathcal{I}}}$, yielding

$$\begin{aligned}\hat{\mathbf{A}} &= (\tilde{\mathbf{X}}^\top \tilde{\mathbf{X}})^{-1} \tilde{\mathbf{X}}^\top \tilde{\mathbf{Y}} \\ &= \left(\begin{bmatrix} \mathbf{X} \\ \boldsymbol{\Gamma} \end{bmatrix}^\top \begin{bmatrix} \mathbf{X} \\ \boldsymbol{\Gamma} \end{bmatrix} \right)^{-1} \begin{bmatrix} \mathbf{X} \\ \boldsymbol{\Gamma} \end{bmatrix}^\top \begin{bmatrix} \mathbf{Y} \\ \mathbf{0} \end{bmatrix} \\ &= (\mathbf{X}^\top \mathbf{X} + \boldsymbol{\Gamma}^\top \boldsymbol{\Gamma})^{-1} \mathbf{X}^\top \mathbf{Y}\end{aligned}$$

$$\begin{aligned}\mathbf{X} &\in \mathbb{R}^{N \times D^{\mathcal{I}}} \\ \mathbf{Y} &\in \mathbb{R}^{N \times D^{\mathcal{O}}} \\ \mathbf{A} &\in \mathbb{R}^{D^{\mathcal{I}} \times D^{\mathcal{O}}}\end{aligned}$$

Ridge regression (Tikhonov regularization)

Ridge regression also has links with SVD. For the singular value decomposition

$$\mathbf{X} = \mathbf{U}\Sigma\mathbf{V}^\top$$

with σ_i the singular values in the diagonal of Σ , a solution to the ridge regression problem is given by

$$\hat{\mathbf{A}} = \mathbf{V}\tilde{\Sigma}\mathbf{U}^\top \mathbf{Y}$$

where $\tilde{\Sigma}$ has diagonal values

$$\tilde{\sigma}_i = \frac{\sigma_i}{\sigma_i^2 + \lambda^2}$$

and has zeros elsewhere.

Ridge regression (Tikhonov regularization)

This regularization improves the conditioning of the problem, thus enabling a direct numerical solution.

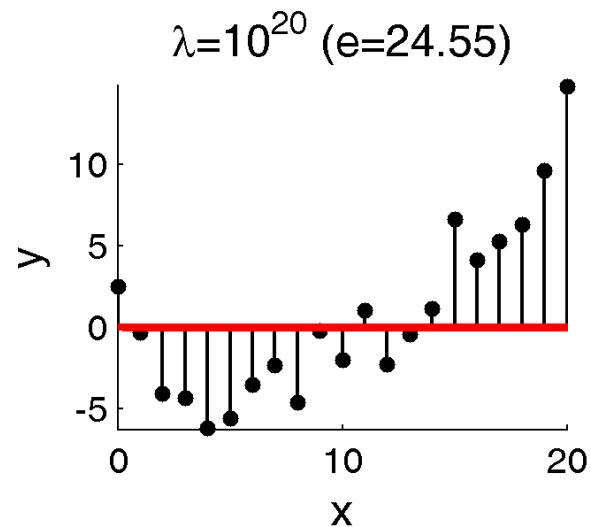
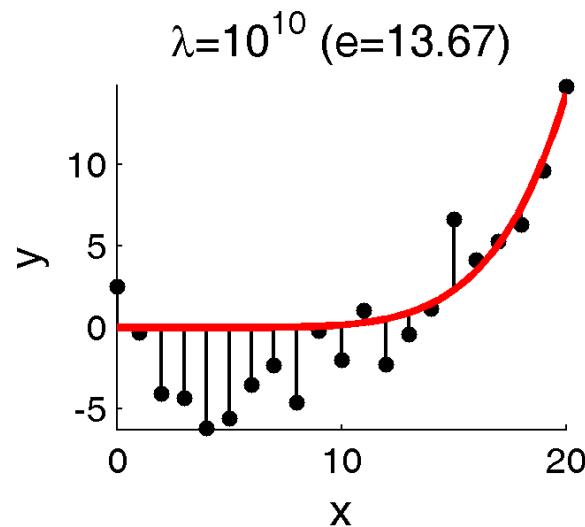
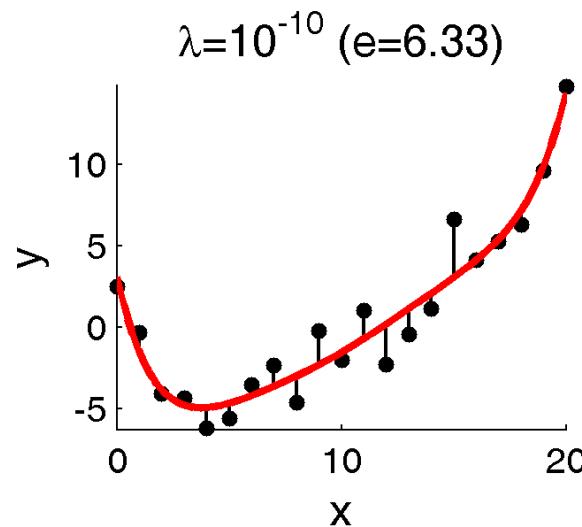
Although at first the choice of the solution to this regularized problem may look artificial and the matrix Γ rather arbitrary, the process can be justified from a **Bayesian point of view**.

Statistically, the prior probability distribution of \mathbf{A} can be taken as a multivariate normal distribution. The assumptions that the means are zero and their components are independent can for example be made, with the components having the same covariance $\Sigma^{\mathbf{A}}$. The data are also subject to errors, with the errors in \mathbf{Y} also assumed to be independent with zero mean and covariance $\Sigma^{\mathbf{Y}}$.

Under these assumptions, the Tikhonov-regularized solution is the most probable solution given the data and the *a priori* distribution of \mathbf{A} , by following Bayes theorem.

Ridge regression (Tikhonov regularization)

$D^x = 7$ (polynomial of degree 7)



LASSO (L_1 regularization)

- An alternative regularized version of least squares is **Lasso (least absolute shrinkage and selection operator)** using the constraint that the L_1 -norm $|A|_1$ is no greater than a given value.
- The **L_1 -regularized formulation** is useful due to its tendency to prefer solutions with fewer nonzero parameter values, effectively reducing the number of variables upon which the given solution is dependent.
- The increase of the penalty term in ridge regression will reduce all parameters while still remaining non-zero, while in Lasso, it will cause more and more of the **parameters to be driven to zero**. This is an advantage of Lasso over ridge regression, as driving parameters to zero **deselects the features from the regression**.

LASSO (L_1 regularization)

- Thus, Lasso automatically selects more relevant features and discards the others, whereas ridge regression never fully discards any features.
- Lasso is equivalent to an unconstrained minimization of the least-squares penalty with $|A|_1$ added. In a Bayesian context, this is equivalent to placing a zero-mean Laplace prior distribution on the parameter vector.
- The optimization problem may be solved using quadratic programming or more general convex optimization methods, as well as by specific algorithms such as the **least angle regression algorithm**.

Weighted least squares (Generalized least squares)

demo_LS_weighted01.m

Weighted least squares

By describing the input data as $\mathbf{X} \in \mathbb{R}^{N \times D^I}$ and the output data as $\mathbf{Y} \in \mathbb{R}^{N \times D^O}$, we want to minimize

$$\begin{aligned}\hat{\mathbf{A}} &= \arg \min_{\mathbf{A}} (\mathbf{Y} - \mathbf{X}\mathbf{A})^\top \mathbf{W} (\mathbf{Y} - \mathbf{X}\mathbf{A}) \\ &= \arg \min_{\mathbf{A}} (\mathbf{Y}^\top \mathbf{W} \mathbf{Y} - 2\mathbf{A}^\top \mathbf{X}^\top \mathbf{W} \mathbf{Y} + \mathbf{A}^\top \mathbf{X}^\top \mathbf{W} \mathbf{X} \mathbf{A})\end{aligned}$$

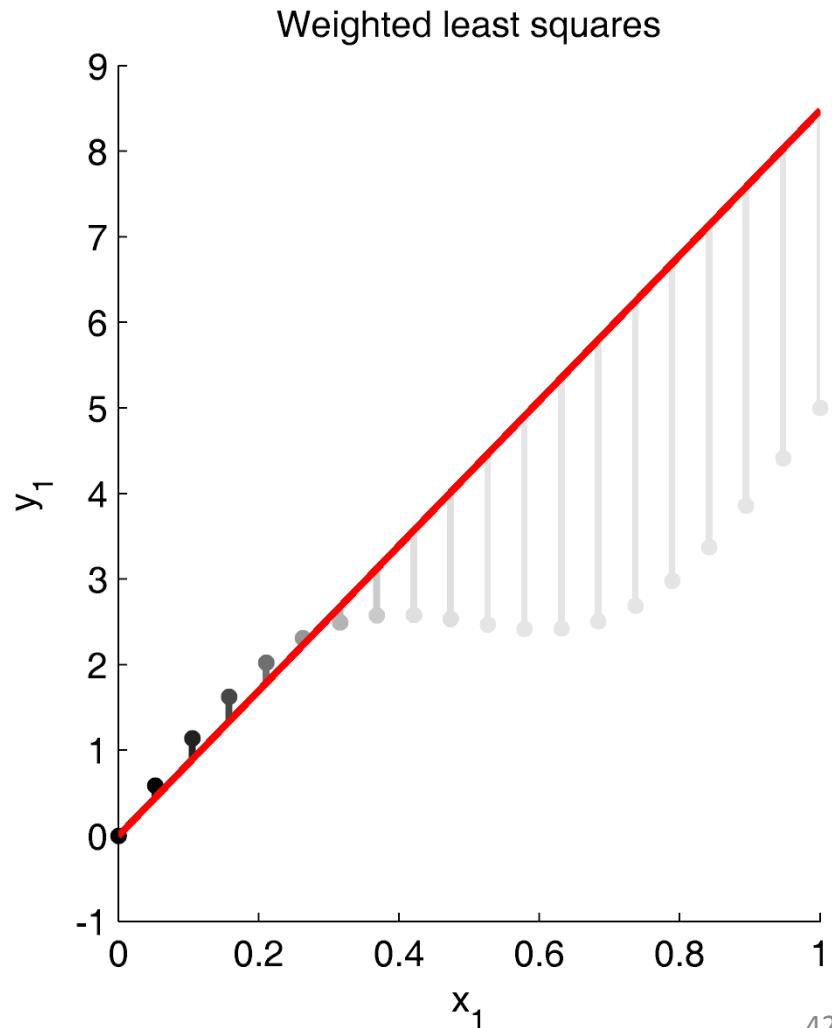
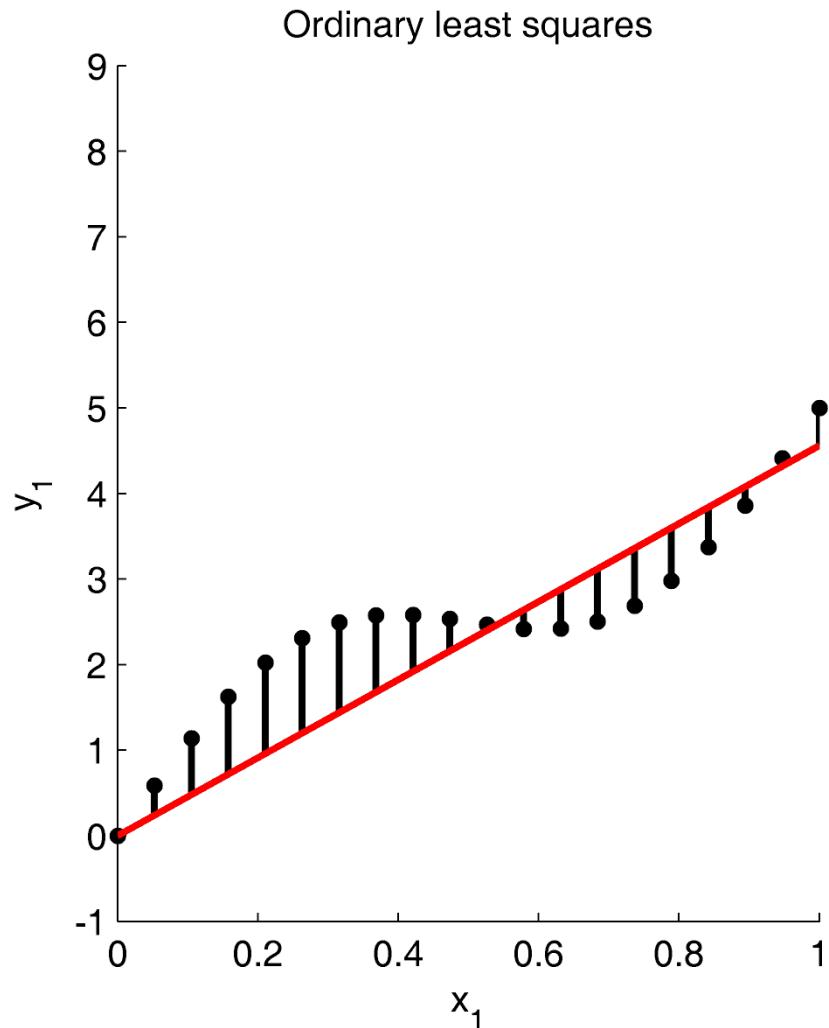
By differentiating with respect to \mathbf{A} and equating to zero

$$\begin{aligned}-2\mathbf{X}^\top \mathbf{W} \mathbf{Y} + 2\mathbf{X}^\top \mathbf{W} \mathbf{X} \mathbf{A} = \mathbf{0} &\iff \hat{\mathbf{A}} = (\mathbf{X}^\top \mathbf{W} \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{W} \mathbf{Y} \\ &\text{or } \hat{\mathbf{A}} = \mathbf{W} \mathbf{X}^\top (\mathbf{X} \mathbf{W} \mathbf{X}^\top)^{-1} \mathbf{Y}\end{aligned}$$

with residuals (parameters errors) given by $\hat{\Sigma}^{\mathbf{A}} = (\mathbf{X}^\top \mathbf{W} \mathbf{X})^{-1}$

Weighted least squares

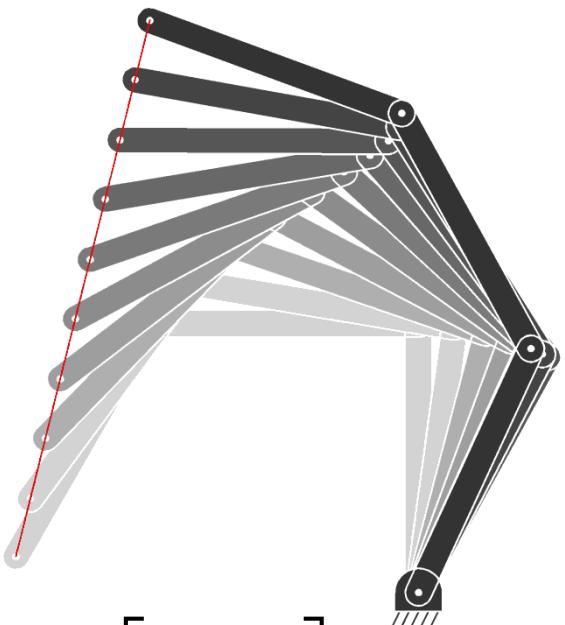
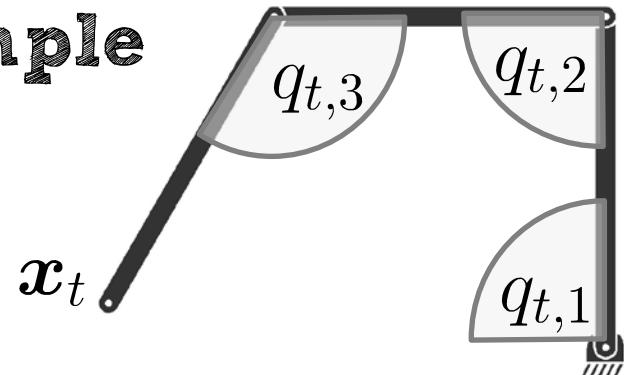
$$\hat{\mathbf{A}} = (\mathbf{X}^\top \mathbf{W} \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{W} \mathbf{Y}$$



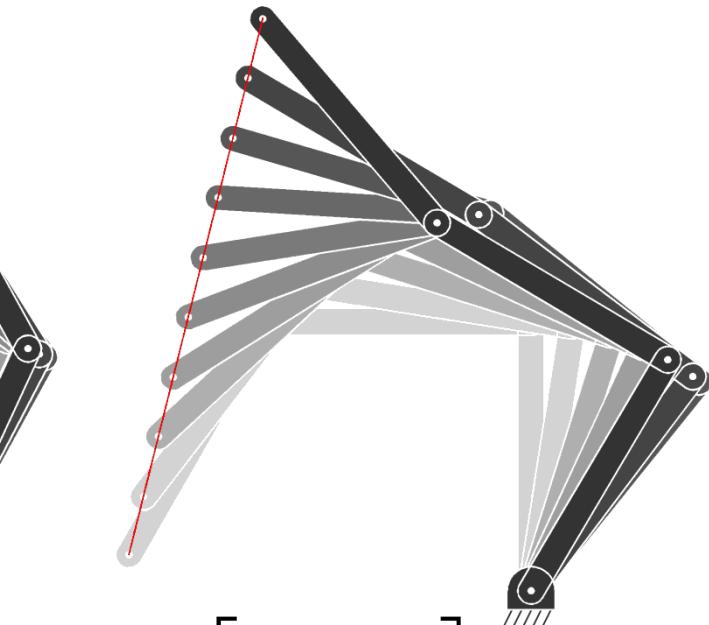
Weighted least squares - Example

$$\hat{A} = W X^\top (X W X^\top)^{-1} Y$$

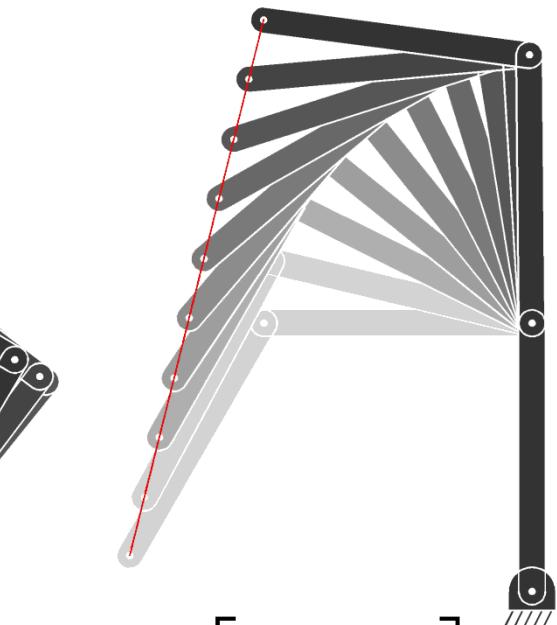
$$\hat{q}_t = W^Q J^\top (J W^Q J^\top)^{-1} \dot{x}_t$$



$$W^Q = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}$$



$$W^Q = \begin{bmatrix} 1 & 0 & 0 \\ 0 & .01 & 0 \\ 0 & 0 & 1 \end{bmatrix}$$

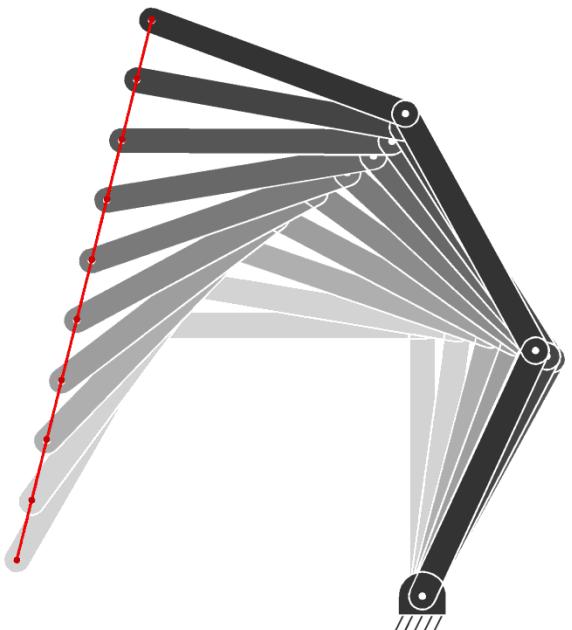
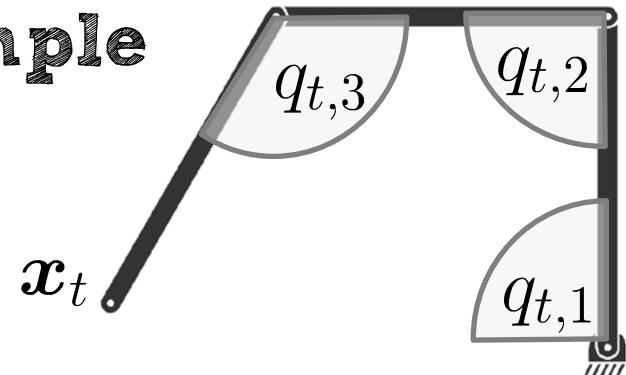


$$W^Q = \begin{bmatrix} .01 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}$$

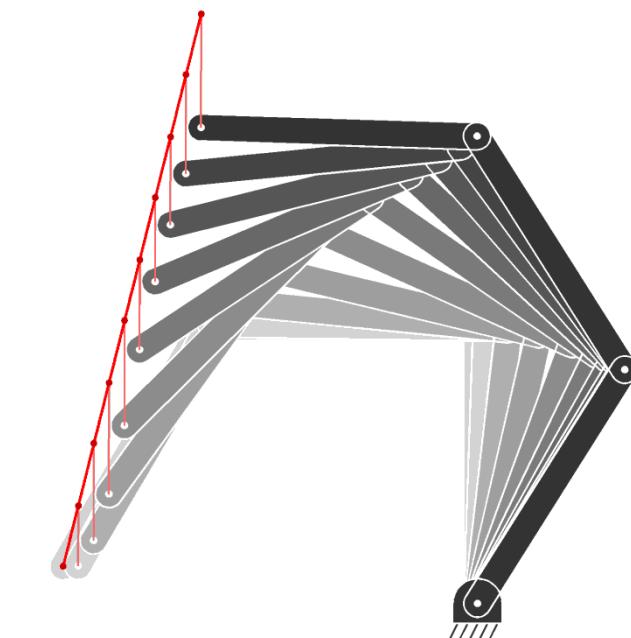
Weighted least squares - Example

$$\hat{A} = (X^\top W X)^{-1} X^\top W Y$$

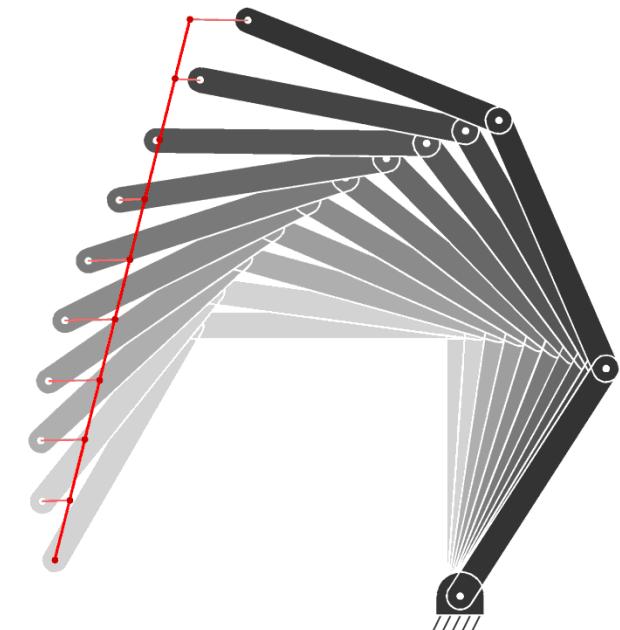
$$\hat{\dot{q}}_t = (J^\top W^\chi J)^{-1} J^\top W^\chi \dot{x}_t$$



$$W^\chi = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$



$$W^\chi = \begin{bmatrix} 1 & 0 \\ 0 & .01 \end{bmatrix}$$



$$W^\chi = \begin{bmatrix} .01 & 0 \\ 0 & 1 \end{bmatrix}$$

Iteratively reweighted least squares (IRLS)

demo_LS_IRLS01.m

Robust regression

- Regression analysis seeks to find the relationship between one or more independent variables and a dependent variable. Methods such as ordinary least squares have favorable properties if their underlying assumptions are true, but can give misleading results if those assumptions are not true
→ **Least squares estimates are highly sensitive to outliers.**
- Robust regression methods are designed to be not overly affected by violations of assumptions through the underlying data-generating process. It down-weights the influence of outliers, which makes their residuals larger and easier to identify.
- A simple approach to robust least squares fitting is to first do an ordinary least squares fit, then identify the k data points with the largest residuals, omit these, perform the fit on the remaining data.

Robust regression

- Another simple method to estimate parameters in a regression model that are less sensitive to outliers than the least squares estimates is to use **least absolute deviations**.
- An alternative parametric approach is to assume that the residuals follow a **mixture of normal distributions**:
 - A contaminated normal distribution in which the majority of observations are from a specified normal distribution, but a small proportion are from a normal distribution with much higher variance.
- Another approach to robust estimation of regression models is to replace the normal distribution with a **heavy-tailed distribution** (e.g., t-distributions).
 - Bayesian robust regression often relies on such distributions.

Iteratively reweighted least squares (IRLS)

- Iteratively Reweighted Least Squares generalizes trimmed least squares by raising the error to a power less than 2.
→ No longer least squares.
- IRLS can be used to find the maximum likelihood estimates of a generalized linear model, and find an M-estimator as a way of mitigating the influence of outliers in an otherwise normally-distributed data set. For example, by minimizing the least absolute error rather than the least square error.
- The strategy is that an error $|e|^p$ can be rewritten as $|e|^p = |e|^{p-2} e^2$. Then, $|e|^{p-2}$ can be interpreted as a weight, which is used to minimize e^2 with **weighted least squares**.
- $p=1$ corresponds to **least absolute deviation regression**.

Iteratively reweighted least squares (IRLS)

The resulting algorithm is simple to implement.

For an L_p norm objective defined by

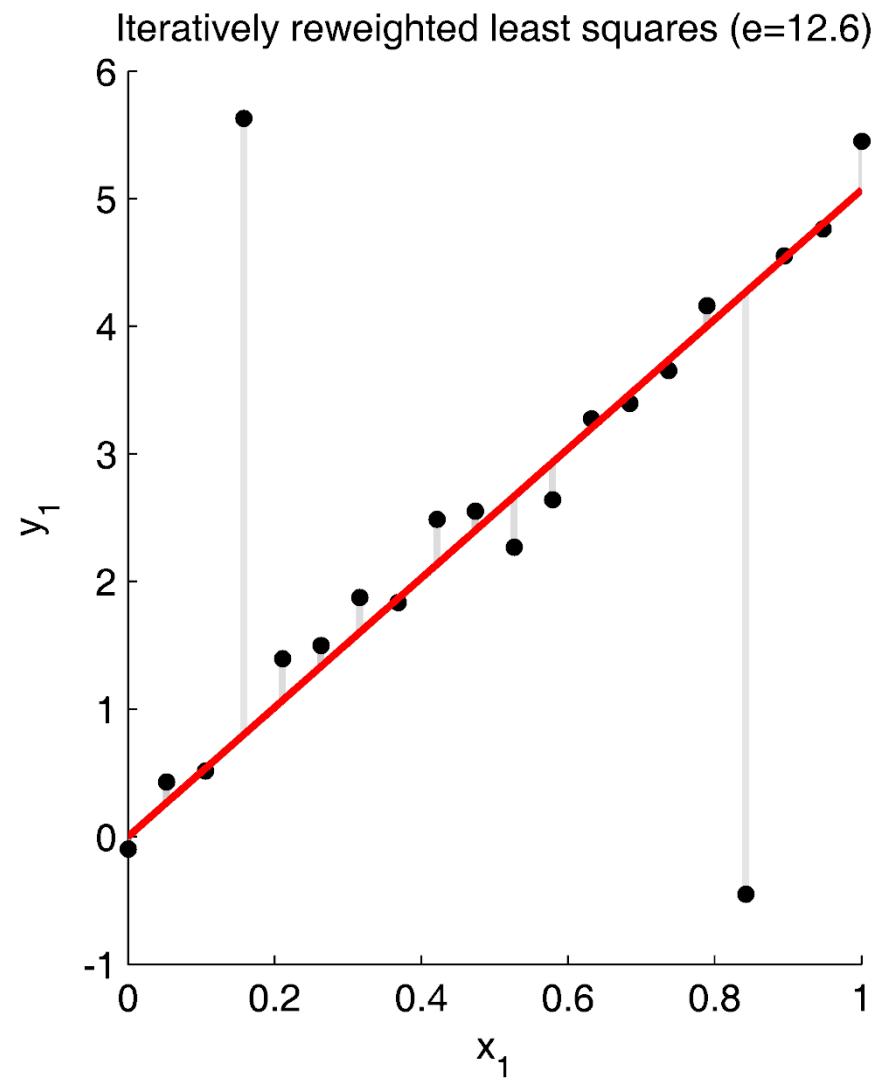
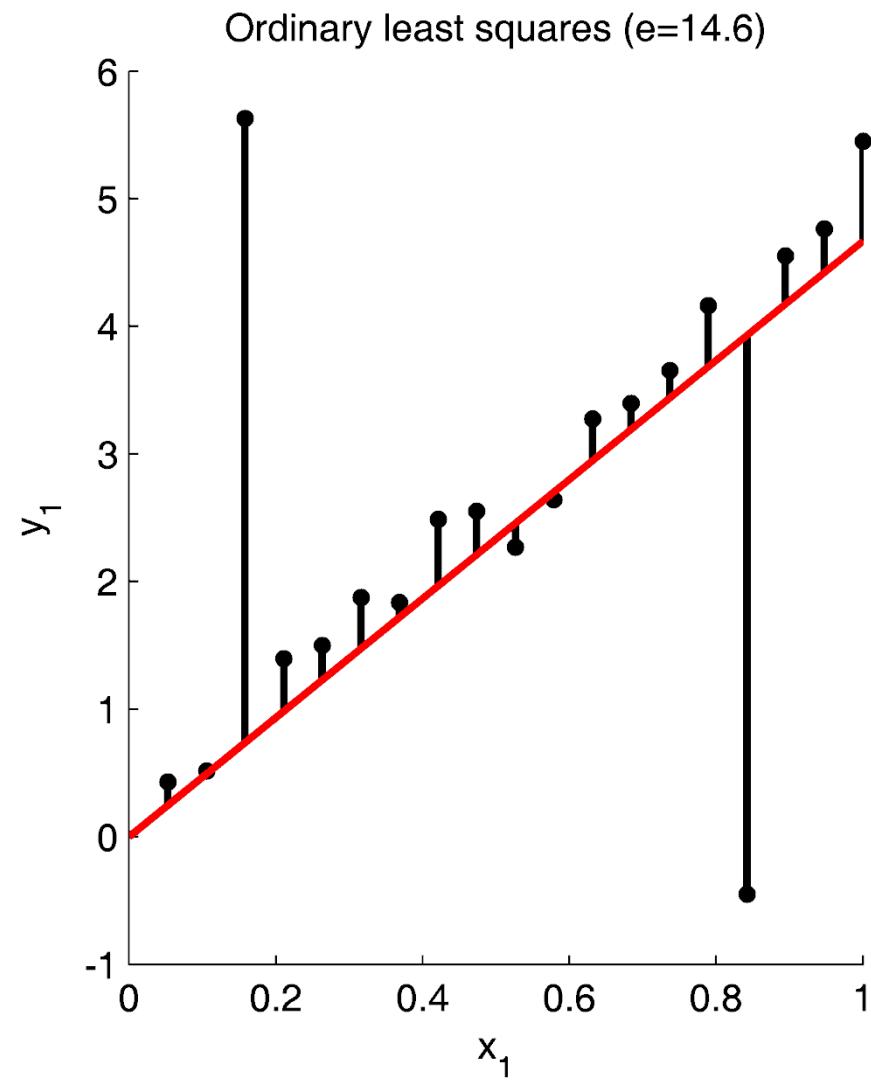
$$\hat{\mathbf{A}} = \arg \min_{\mathbf{A}} \|\mathbf{Y} - \mathbf{X}\mathbf{A}\|_p$$

$\hat{\mathbf{A}}$ is estimated by starting from $\mathbf{W} = \mathbf{I}$ and iteratively computing until convergence

$$\hat{\mathbf{A}} \leftarrow (\mathbf{X}^\top \mathbf{W} \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{W} \mathbf{Y}$$

$$\mathbf{W}_{t,t} \leftarrow |\mathbf{Y}_t - \mathbf{X}_t \mathbf{A}|^{p-2} \quad \forall t \in \{1, \dots, T\}$$

Iteratively reweighted least squares (IRLS)



Recursive least squares

demo_LS_recursive01.m

Recursive least squares

If new datapoints become available after computation of the pseudoinverse, the solution can be updated efficiently by exploiting the **Sherman-Morrison-Woodbury formulas** that relate the inverse of a matrix after a small-rank perturbation to the inverse of the original matrix, namely

$$(\mathbf{B} + \mathbf{U}\mathbf{V})^{-1} = \mathbf{B}^{-1} - \overbrace{\mathbf{B}^{-1}\mathbf{U} (\mathbf{I} + \mathbf{V}\mathbf{B}^{-1}\mathbf{U})^{-1} \mathbf{V}\mathbf{B}^{-1}}^{\mathbf{E}}$$

with $\mathbf{U} \in \mathbb{R}^{n \times m}$ and $\mathbf{V} \in \mathbb{R}^{m \times n}$.

When $m \ll n$, the correction term \mathbf{E} can be computed more efficiently than inverting $\mathbf{B} + \mathbf{U}\mathbf{V}$.

Recursive least squares

This can be proven by showing that multiplying both sides of the equation by $\mathbf{B} + \mathbf{U}\mathbf{V}$ gives the identity.

For the right side, we can see that

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

By defining $\mathbf{B} = \mathbf{X}^\top \mathbf{X}$, the above relation can be exploited to update a least squares solution when new datapoints become available.

Recursive least squares

Indeed, if $\mathbf{X}_{\text{new}} = [\mathbf{X}^\top, \mathbf{V}^\top]^\top$ and $\mathbf{Y}_{\text{new}} = [\mathbf{Y}^\top, \mathbf{C}^\top]^\top$, we then have

$$\begin{aligned}\mathbf{B}_{\text{new}} &= \mathbf{X}_{\text{new}}^\top \mathbf{X}_{\text{new}} \\ &= \mathbf{X}^\top \mathbf{X} + \mathbf{V}^\top \mathbf{V} \\ &= \mathbf{B} + \mathbf{V}^\top \mathbf{V}\end{aligned}$$

whose inverse can be computed with

$$\mathbf{B}_{\text{new}}^{-1} = \mathbf{B}^{-1} - \mathbf{B}^{-1} \mathbf{V}^\top (\mathbf{I} + \mathbf{V} \mathbf{B}^{-1} \mathbf{V}^\top)^{-1} \mathbf{V} \mathbf{B}^{-1}$$

Recursive least squares

This is exploited to estimate the update as

$$\begin{aligned}
 \hat{\mathbf{A}}_{\text{new}} &= \mathbf{B}_{\text{new}}^{-1} \mathbf{X}_{\text{new}}^\top \mathbf{Y}_{\text{new}} \\
 &= \mathbf{B}_{\text{new}}^{-1} (\mathbf{X}^\top \mathbf{Y} + \mathbf{V}^\top \mathbf{C}) \\
 \hat{\mathbf{A}} &= \mathbf{B}^{-1} \mathbf{X}^\top \mathbf{Y} \\
 &= \mathbf{B}_{\text{new}}^{-1} (\mathbf{B} \hat{\mathbf{A}} + \mathbf{V}^\top \mathbf{C}) \\
 \mathbf{B}_{\text{new}} &= \mathbf{B} + \mathbf{V}^\top \mathbf{V} \\
 &= \mathbf{B}_{\text{new}}^{-1} ((\mathbf{B}_{\text{new}} - \mathbf{V}^\top \mathbf{V}) \hat{\mathbf{A}} + \mathbf{V}^\top \mathbf{C}) \\
 &= \mathbf{B}_{\text{new}}^{-1} (\mathbf{B}_{\text{new}} \hat{\mathbf{A}} - \mathbf{V}^\top \mathbf{V} \hat{\mathbf{A}} + \mathbf{V}^\top \mathbf{C}) \\
 &= \hat{\mathbf{A}} + \mathbf{B}_{\text{new}}^{-1} \mathbf{V}^\top (C - \mathbf{V} \hat{\mathbf{A}}) \\
 &= \hat{\mathbf{A}} + \mathbf{K} (C - \mathbf{V} \hat{\mathbf{A}})
 \end{aligned}$$

$\mathbf{X}_{\text{new}} = \begin{bmatrix} \mathbf{X} \\ \mathbf{V} \end{bmatrix}$, $\mathbf{Y}_{\text{new}} = \begin{bmatrix} \mathbf{Y} \\ \mathbf{C} \end{bmatrix}$

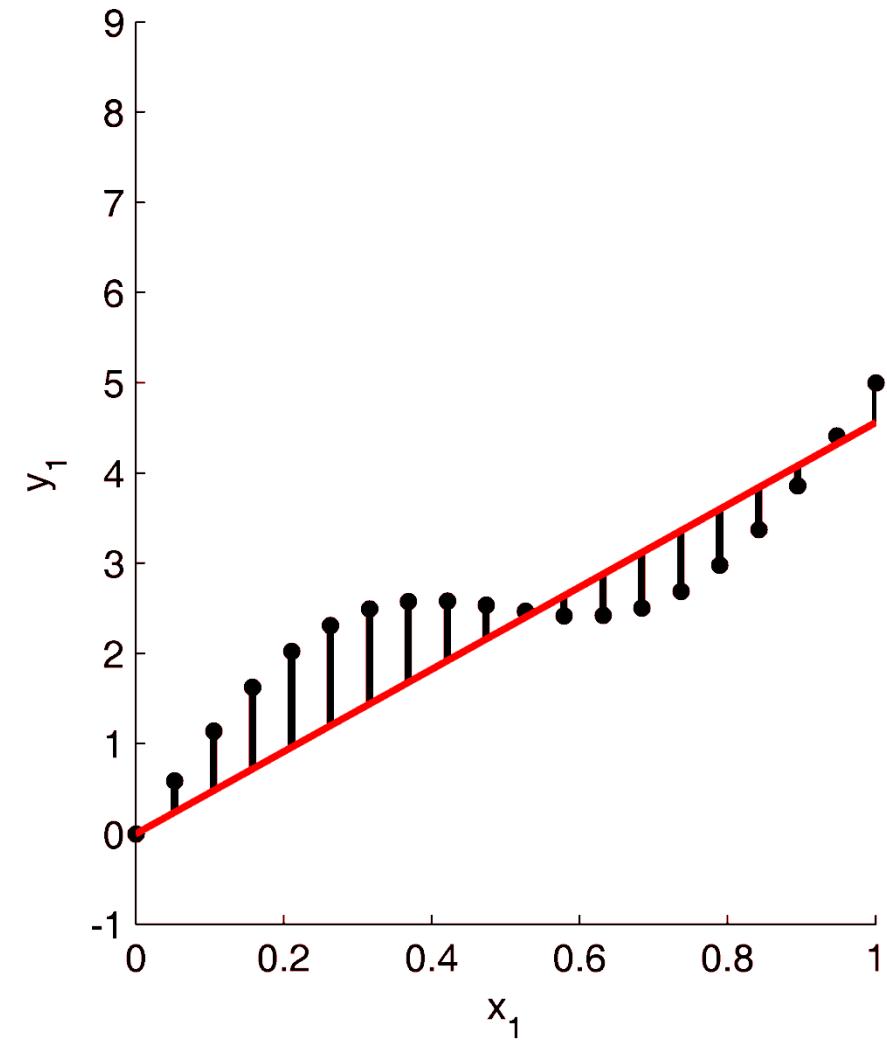
Recursive least squares

with Kalman gain efficiently computed as

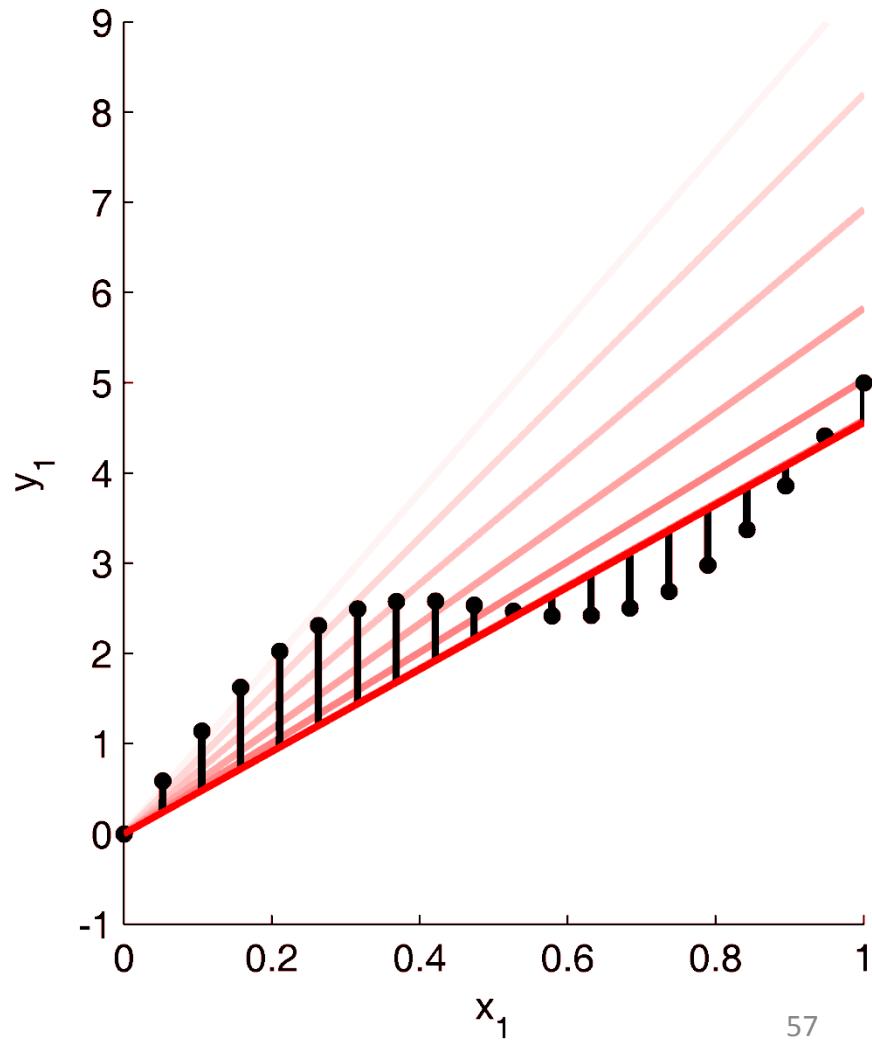
$$\begin{aligned} \mathbf{K} &= \mathbf{B}_{\text{new}}^{-1} \mathbf{V}^\top \\ &= \mathbf{B}^{-1} \mathbf{V}^\top - \mathbf{B}^{-1} \mathbf{V}^\top (\mathbf{I} + \mathbf{V} \mathbf{B}^{-1} \mathbf{V}^\top)^{-1} \mathbf{V} \mathbf{B}^{-1} \mathbf{V}^\top \\ &= \mathbf{B}^{-1} \mathbf{V}^\top (\mathbf{I} - (\mathbf{I} + \mathbf{V} \mathbf{B}^{-1} \mathbf{V}^\top)^{-1} \mathbf{V} \mathbf{B}^{-1} \mathbf{V}^\top) \\ &= \mathbf{B}^{-1} \mathbf{V}^\top (\mathbf{I} + \mathbf{V} \mathbf{B}^{-1} \mathbf{V}^\top)^{-1} ((\mathbf{I} + \mathbf{V} \mathbf{B}^{-1} \mathbf{V}^\top) - \mathbf{V} \mathbf{B}^{-1} \mathbf{V}^\top) \\ &= \mathbf{B}^{-1} \mathbf{V}^\top (\mathbf{I} + \mathbf{V} \mathbf{B}^{-1} \mathbf{V}^\top)^{-1} \end{aligned}$$

Recursive least squares

Ordinary least squares ($e=11.0$)



Recursive least squares ($e=11.0$)



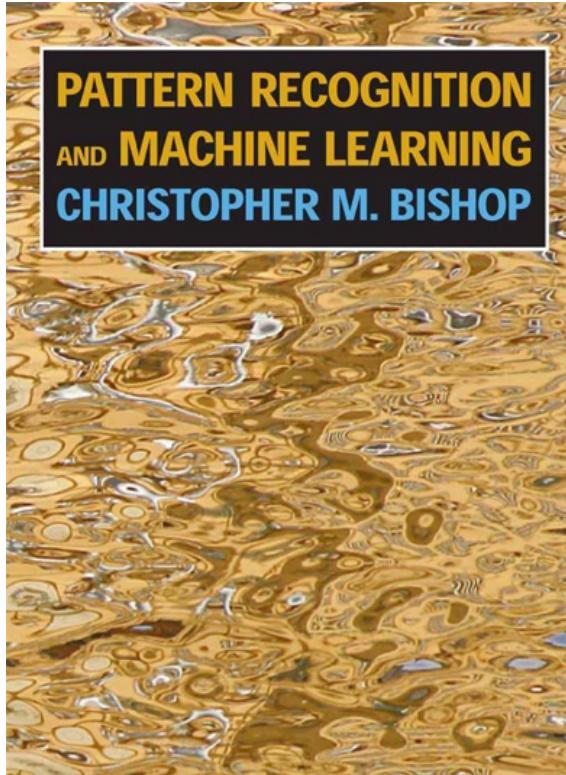
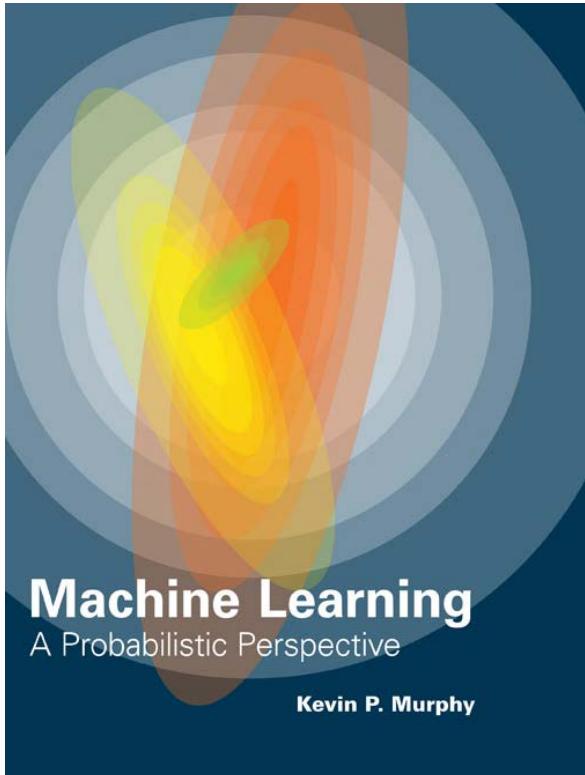
Main references

Regression

F. Stulp and O. Sigaud. Many regression algorithms, one unified model – a review.
Neural Networks, 69:60–79, September 2015

W. W. Hager. Updating the inverse of a matrix. SIAM Review, 31(2):221–239, 1989

G. Strang. Introduction to Applied Mathematics. Wellesley-Cambridge Press, 1986.



**The Matrix
Cookbook**

Kaare Brandt Petersen
Michael Syskind Pedersen