

**EE613**  
**Machine Learning for Engineers**

**LINEAR REGRESSION I**

**Sylvain Calinon**  
**Robot Learning & Interaction Group**  
**Idiap Research Institute**  
**Oct. 31, 2019**

# **EE613 - List of courses**

19.09.2019 (JMO) Introduction

26.09.2019 (JMO) Generative I

03.10.2019 (JMO) Generative II

10.10.2019 (JMO) Generative III

17.10.2019 (JMO) Generative IV

24.10.2019 (JMO) Decision-trees

**31.10.2019 (SC) Linear regression I**

07.11.2019 (JMO) Kernel SVM

**14.11.2019 (SC) Linear regression II**

21.11.2019 (FF) MLP

28.11.2019 (FF) Feature-selection and boosting

**05.12.2019 (SC) HMM and subspace clustering**

**12.12.2019 (SC) Nonlinear regression I**

**19.12.2019 (SC) Nonlinear regression II**

# Outline

## Linear Regression I (Oct 31)

- Least squares
- Singular value decomposition (SVD)
- Kernels in least squares (nullspace)
- Ridge regression (Tikhonov regularization)
- Weighted least squares
- Iteratively reweighted least squares (IRLS)
- Recursive least squares

## Linear Regression II (Nov 14)

- Logistic regression
- Tensor-variate regression

## Hidden Markov model (HMM) & subspace clustering (Dec 5)

## Nonlinear Regression I (Dec 12)

- Locally weighted regression (LWR)
- Gaussian mixture regression (GMR)

## Nonlinear Regression II (Dec 19)

- Gaussian process regression (GPR)

# Labs

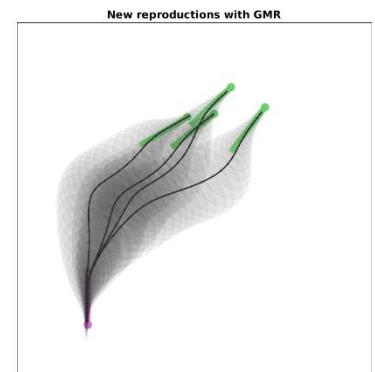
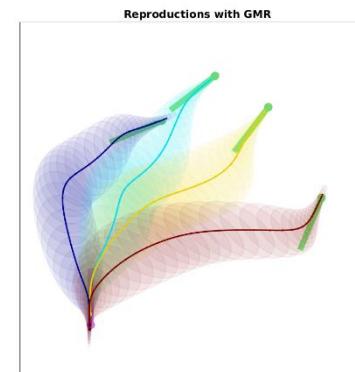
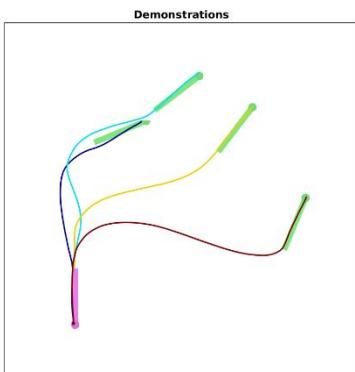
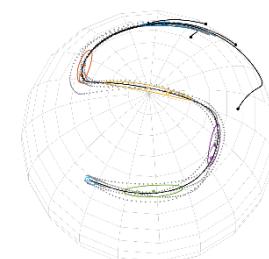
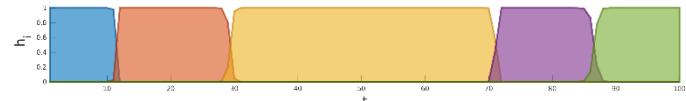
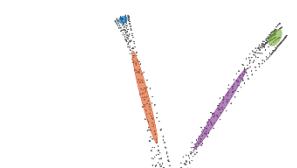
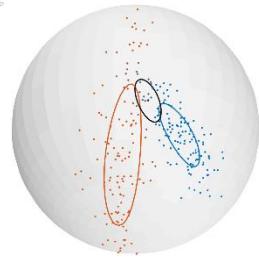
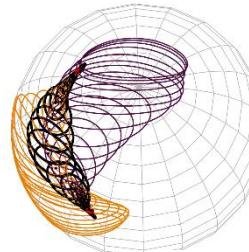
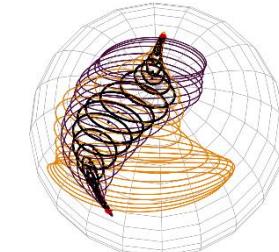
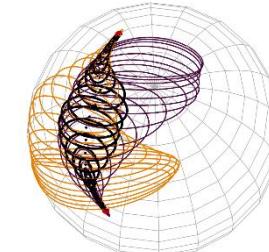
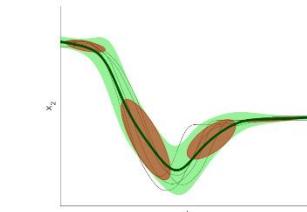
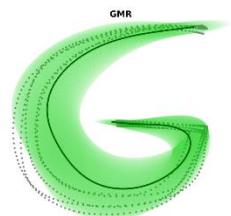
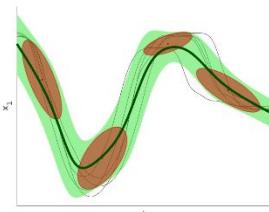
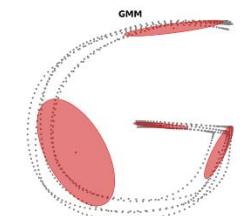
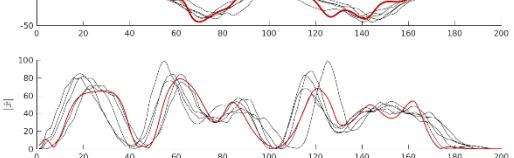
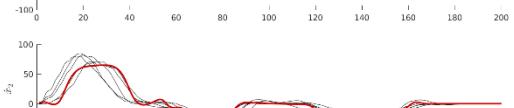
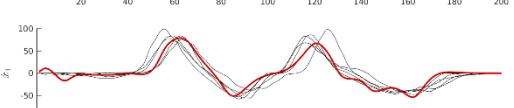
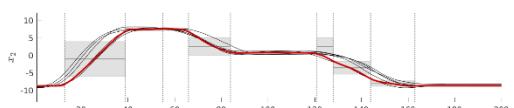
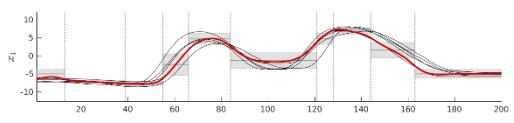
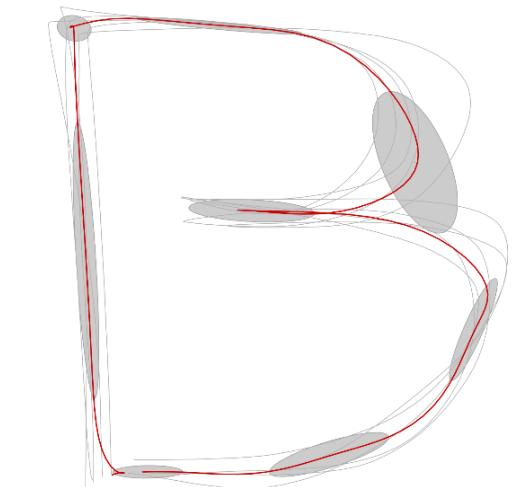


## Teguh Lembono

**Python notebooks and labs exercises:**  
<https://github.com/teguhSL/ee613-python>

Branch: master	ee613-python / python_notebooks / linear_regression_1 /	Create new file	Find file	History
 teguhSL minor edits				Latest commit c1da0e0 9 hours ago
..				
 Ex1.ipynb	minor edits			9 hours ago
 Ex2.ipynb	minor edits			9 hours ago
 Ex3.ipynb	minor edits			9 hours ago
 demo_LS.ipynb	minor edits			9 hours ago
 demo_LS_polFit.ipynb	minor edits			9 hours ago
 demo_LS_recursive.ipynb	minor edits			9 hours ago
 demo_LS_weighted.ipynb	minor edits			9 hours ago

# PbDlib

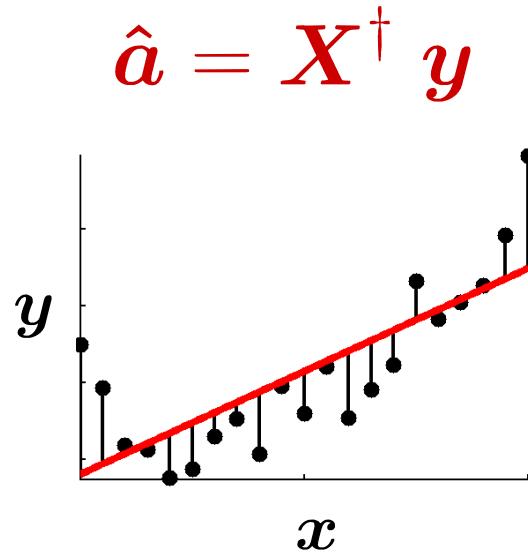


<http://www.idiap.ch/software/pbdlib/>

# LEAST SQUARES

circa 1795

# Least squares: a ubiquitous tool



Weighted least squares?

Regularized least squares?

L1-norm instead of L2-norm?

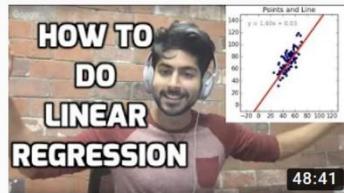
Nullspace structure?

Recursive computation?

[Home](#)[Trending](#)[History](#)**BEST OF YOUTUBE**[Music](#)[Sports](#)[Gaming](#)[Movies](#)[News](#)[Live](#)[360° Video](#)[Browse channels](#)[FILTER](#)**Applications of Deep Neural Networks  
Regression****5.3: Regression Neural Networks for Keras and TensorFlow (Module 5, Part 3)**

Jeff Heaton • 6K views • 1 year ago

Performing regression with keras neural networks. Producing a lift chart. This video is part of a course that is taught in a hybrid ...

**Linear Regression Machine Learning (tutorial)**

Siraj Raval • 87K views • Streamed 2 years ago

I'll perform linear regression from scratch in Python using a method called 'Gradient Descent' to determine the relationship ...

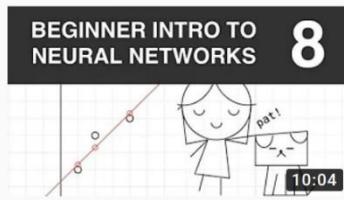
CC

**3.4: Linear Regression with Gradient Descent - Intelligence and Learning**

The Coding Train • 64K views • 1 year ago

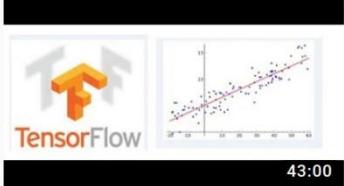
In this video I continue my Machine Learning series and attempt to explain Linear Regression with Gradient Descent. My Video ...

CC

**Beginner Intro to Neural Networks 8: Linear Regression**

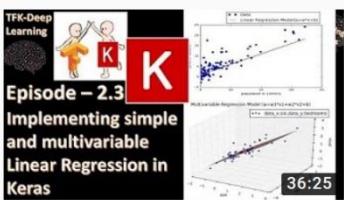
giant\_neural\_network • 53K views • 1 year ago

Hey everyone! In this video we're going to look at something called linear regression. We're really just adding an input to our ...

**Learning Tensorflow with linear regression**

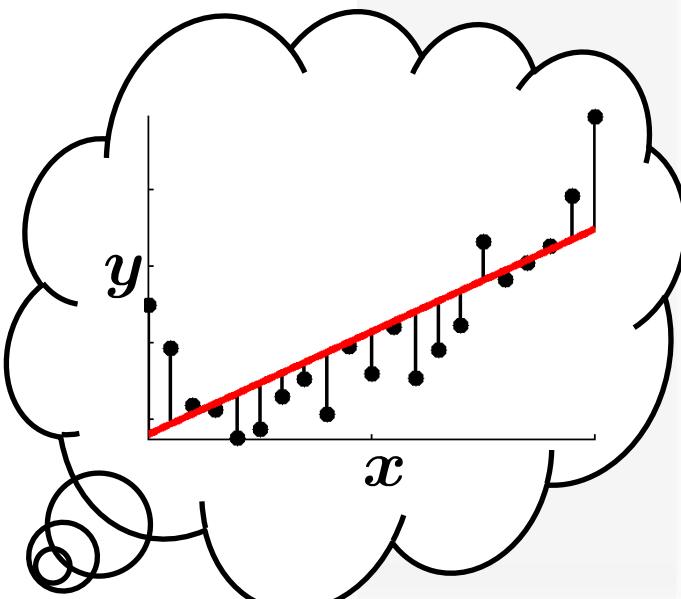
Technology for Noobs • 3.5K views • 1 year ago

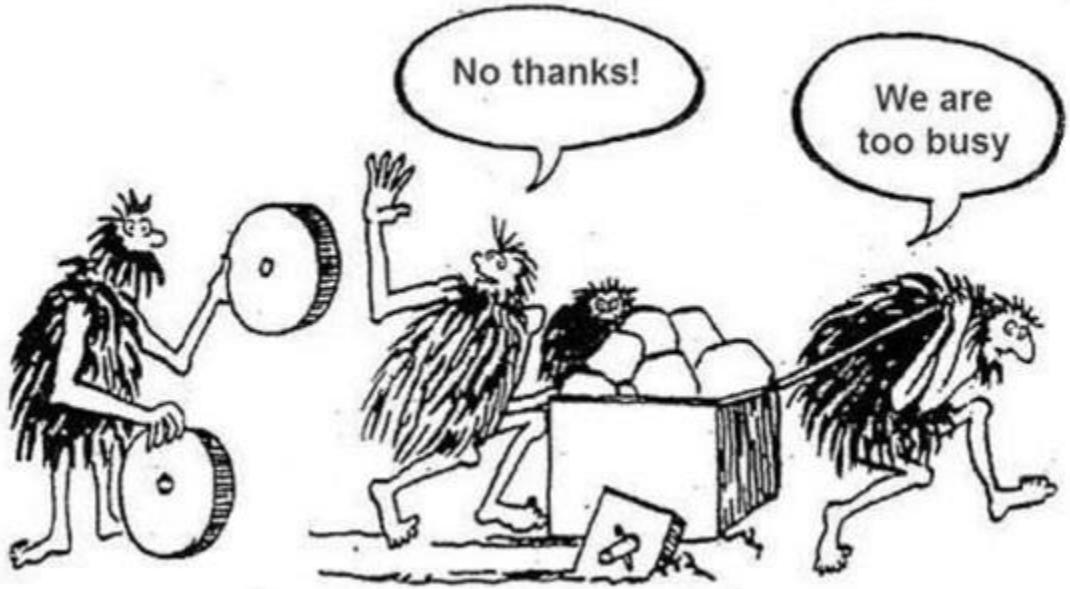
In this video, I will cover basics of tensorflow. Below are the topics that will be covered: 1. Basic of linear regression 2. Basics of ...

**Ep-2.3: Linear Regression in Keras || TFK-Deep Learning || Exploring Neurons**

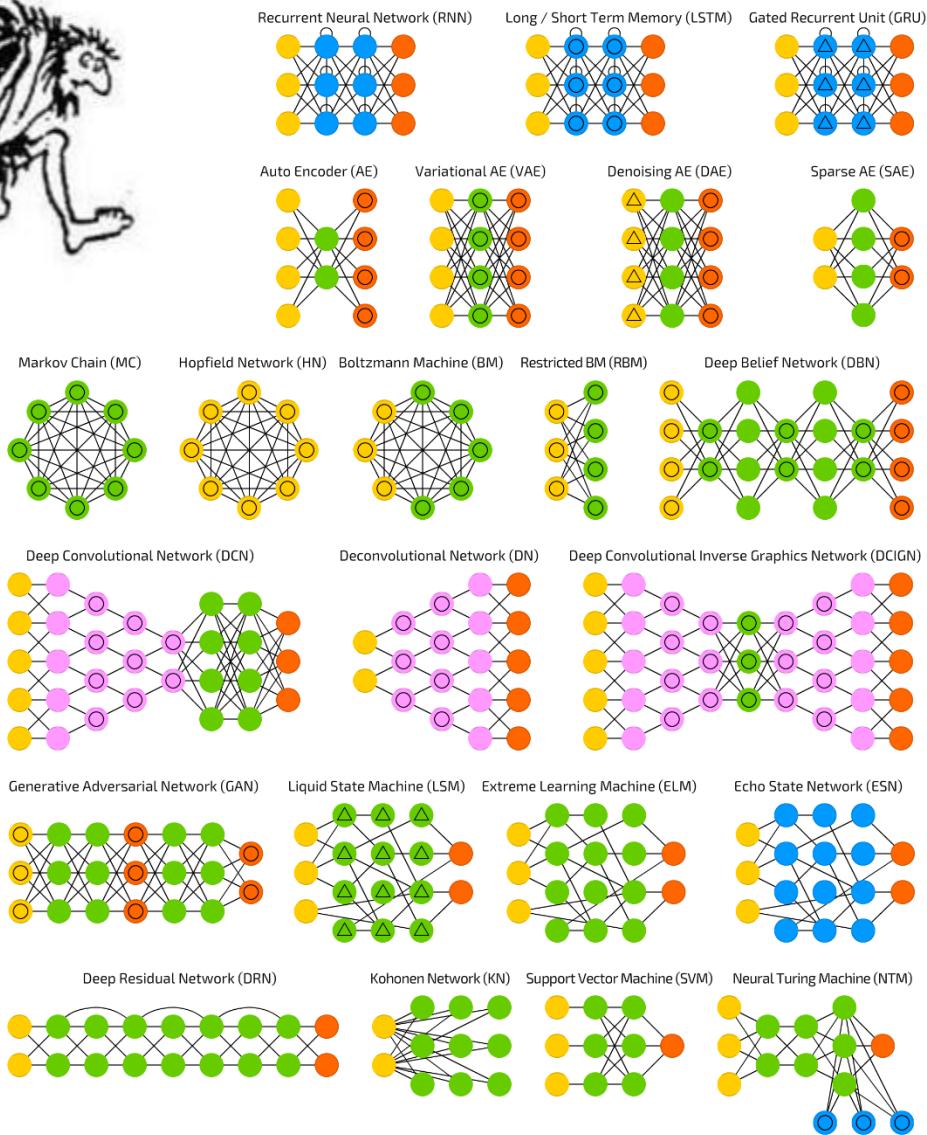
Anuj shah • 1.2K views • 1 year ago

This video explains the implementation of simple and multiple linear regression in keras. The theoretical discussion of linear ...





- Backfed Input Cell
- Input Cell
- △ Noisy Input Cell
- Hidden Cell
- Probabilistic Hidden Cell
- △ Spiking Hidden Cell
- Output Cell
- Match Input Output Cell
- Recurrent Cell
- Memory Cell
- △ Different Memory Cell
- Kernel
- Convolution or Pool



# **Linear regression**

**Python notebooks:**

**demo\_LS.ipynb, demo\_LS\_polFit.ipynb**

**Matlab codes:**

**demo\_LS01.m, demo\_LS\_polFit01.m**

# Linear regression

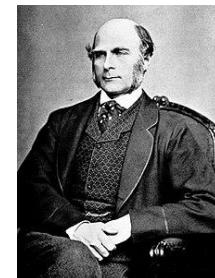
- **Least squares is everywhere:** from simple problems to large scale problems.
- It was the earliest form of regression, which was published by **Legendre** in 1805 and by **Gauss** in 1809. They both applied the method to the problem of determining the orbits of bodies around the Sun from astronomical 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.
- **Pearson** later provided the statistical context showing that the phenomenon is more general than a biological context.



Adrien-Marie Legendre



Carl Friedrich Gauss



Francis Galton



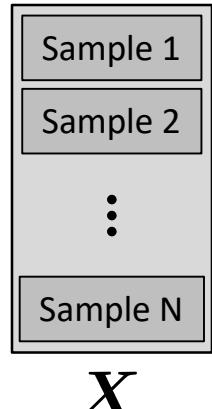
Karl Pearson

# Multivariate linear regression

By describing the input data as  $\mathbf{X} \in \mathbb{R}^{N \times D^I}$  and the output data as  $\mathbf{y} \in \mathbb{R}^N$ , we want to find  $\mathbf{a} \in \mathbb{R}^{D^I}$  to have  $\mathbf{y} = \mathbf{X}\mathbf{a}$ .

A solution can be found by minimizing the  $\ell_2$  norm

$$\begin{aligned}\hat{\mathbf{a}} &= \arg \min_{\mathbf{a}} \|\mathbf{y} - \mathbf{X}\mathbf{a}\|^2 \\ &= \arg \min_{\mathbf{a}} (\mathbf{y} - \mathbf{X}\mathbf{a})^\top (\mathbf{y} - \mathbf{X}\mathbf{a}) \\ &= \arg \min_{\mathbf{a}} \mathbf{y}^\top \mathbf{y} - 2\mathbf{a}^\top \mathbf{X}^\top \mathbf{y} + \mathbf{a}^\top \mathbf{X}^\top \mathbf{X}\mathbf{a}\end{aligned}$$



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

$$-2\mathbf{X}^\top \mathbf{y} + 2\mathbf{X}^\top \mathbf{X}\mathbf{a} = \mathbf{0} \iff \hat{\mathbf{a}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}$$

Moore-Penrose  
pseudoinverse  $\mathbf{X}^\dagger$

# Multiple multivariate linear regression

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  $\mathbf{A} \in \mathbb{R}^{D^I \times D^O}$  to have  $\mathbf{Y} = \mathbf{X}\mathbf{A}$ .

A solution can be found by minimizing the Frobenius norm

$$\begin{aligned}\hat{\mathbf{A}} &= \arg \min_{\mathbf{A}} \|\mathbf{Y} - \mathbf{X}\mathbf{A}\|_F^2 \\ &= \arg \min_{\mathbf{A}} \text{tr}\left((\mathbf{Y} - \mathbf{X}\mathbf{A})^\top(\mathbf{Y} - \mathbf{X}\mathbf{A})\right) \\ &= \arg \min_{\mathbf{A}} \text{tr}(\mathbf{Y}^\top\mathbf{Y} - 2\mathbf{A}^\top\mathbf{X}^\top\mathbf{Y} + \mathbf{A}^\top\mathbf{X}^\top\mathbf{X}\mathbf{A})\end{aligned}$$

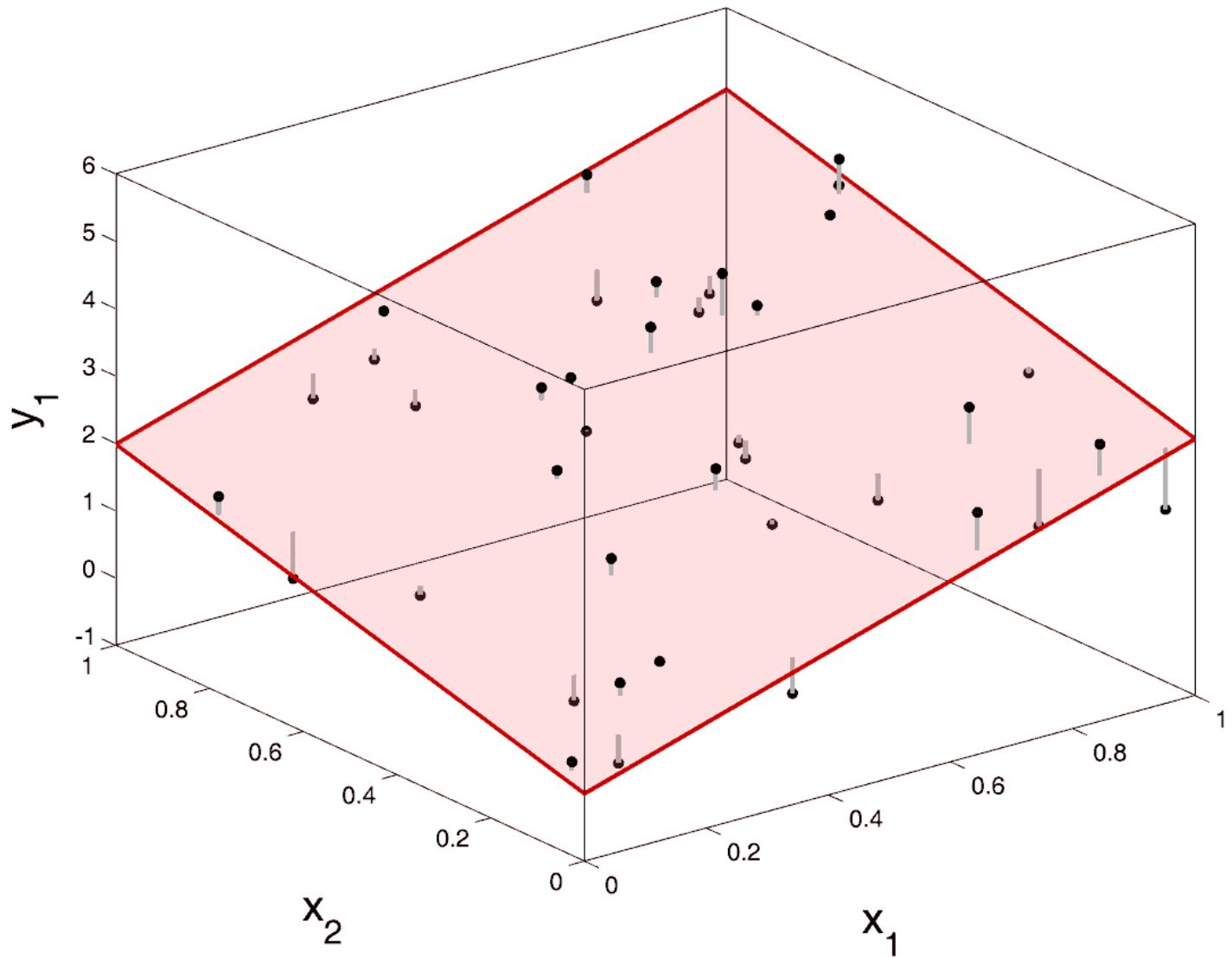
By differentiating with respect to  $\mathbf{A}$  and equating to zero

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

Moore-Penrose  
pseudoinverse  $\mathbf{X}^\dagger$



# Example of multivariate linear regression



$$\boldsymbol{x} = [x_1, x_2]$$

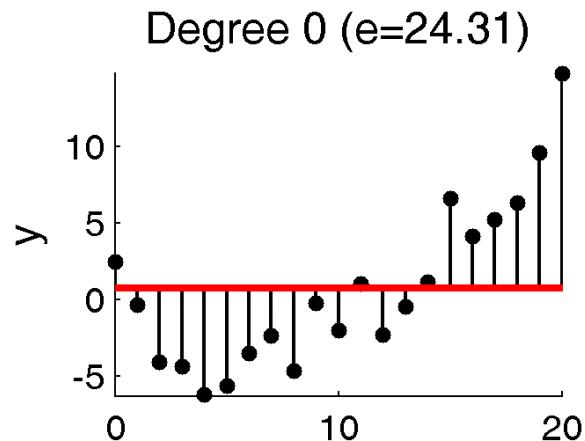
$$N = 40$$

$$D^{\mathcal{I}} = 2$$

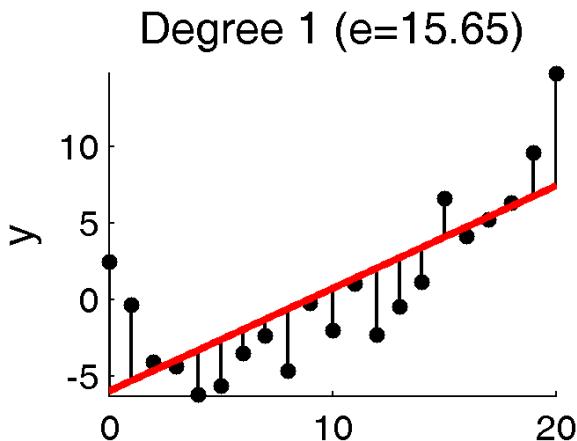
$$D^{\mathcal{O}} = 1$$

# Polynomial fitting with least squares

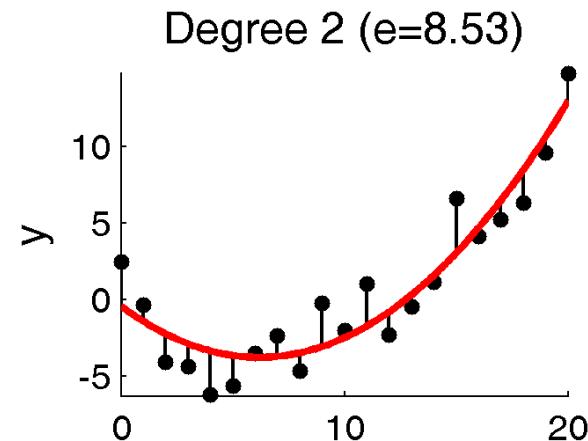
$$\hat{A} = X^\dagger Y$$



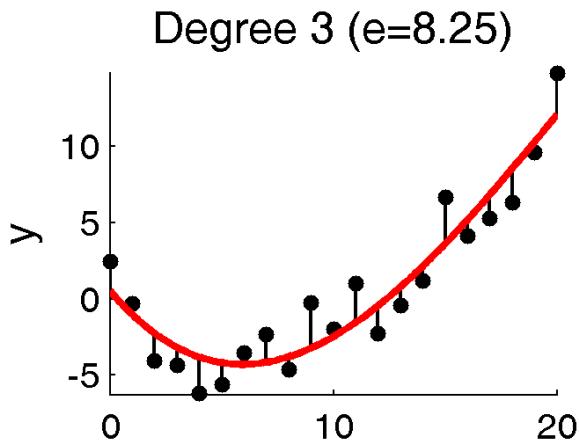
$$x = 1$$



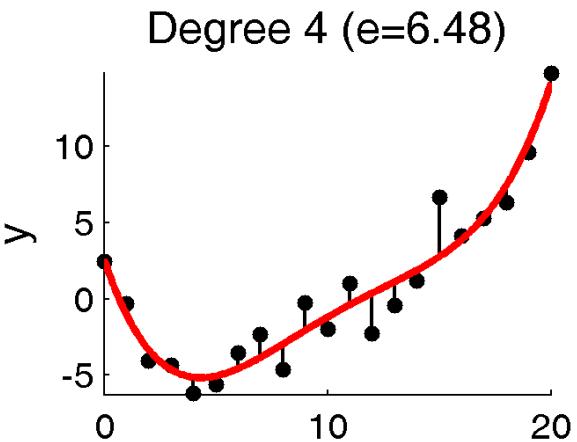
$$x = [1, x]$$



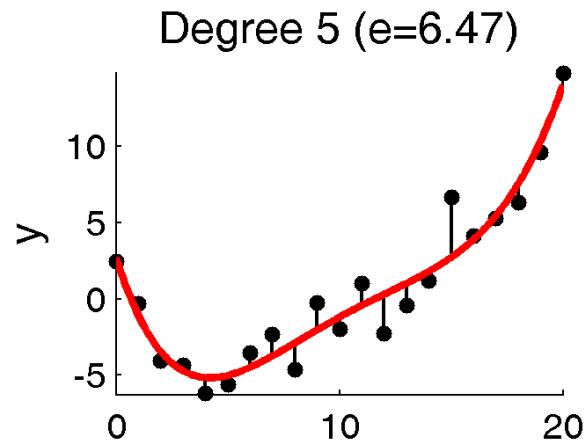
$$x = [1, x, x^2]$$



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



$$x = [1, x, x^2, x^3, x^4]$$



$$x = [1, x, x^2, x^3, x^4, x^5]$$

# Singular value decomposition (SVD)

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

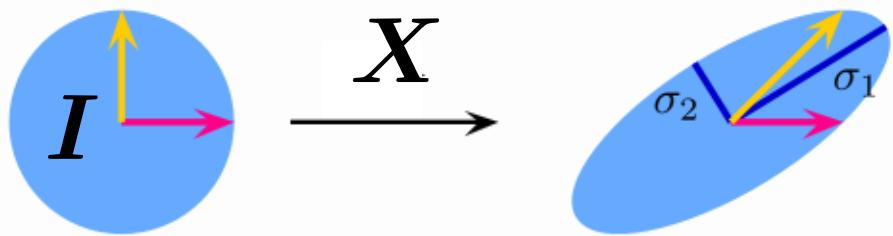
$$\begin{bmatrix} 1 & 0 & 0 & 0 & 2 \\ 0 & 0 & 3 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \\ 0 & 4 & 0 & 0 & 0 \end{bmatrix} = \begin{bmatrix} 0 & 0 & 1 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & -1 \\ 1 & 0 & 0 & 0 \end{bmatrix} \cdot \begin{bmatrix} 4 & 0 & 0 & 0 & 0 \\ 0 & 3 & 0 & 0 & 0 \\ 0 & 0 & \sqrt{5} & 0 & 0 \\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

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

Unitary matrix  
(orthogonal)

Unitary matrix  
(orthogonal)

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



# Least squares with SVD

$$\mathbf{X} \in \mathbb{R}^{N \times D^I}$$

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

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

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

where  $\mathbf{U}$  and  $\mathbf{V}$  are  $N \times N$  and  $D^I \times D^I$  orthogonal matrices, and  $\Sigma$  is an  $N \times D^I$  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  $\Sigma$  can be easily obtained by inverting the non-zero diagonal elements and transposing the resulting matrix.

# **Kernels in least squares (nullspace projection)**

**Python notebook:**  
**demo\_LS\_polFit.ipynb**

**Matlab code:**  
**demo\_LS\_polFit\_nullspace01.m**

# Kernels in least squares (**nullspace**)

The pseudoinverse provides a single least norm solution, but we can sometimes obtain other solutions by employing a **nullspace projection operator  $N$**

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

$\mathbf{V}$  can be any vector/matrix (typically, a gradient minimizing a secondary objective function).

The nullspace projection guarantees that  $\|\mathbf{Y} - \mathbf{X}\hat{\mathbf{A}}\|_F^2$  is still minimized.

# Kernels in least squares (nullspace)

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

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

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

to compute

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

$$\begin{matrix} \mathbf{X}^\dagger \in \mathbb{R}^{D^T \times N} \\ \begin{bmatrix} 1 & 0 & 0 & 0 & 2 \\ 0 & 0 & 3 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \\ 0 & 4 & 0 & 0 & 0 \end{bmatrix} \end{matrix} = \begin{matrix} \mathbf{U} \in \mathbb{R}^{D^T \times D^T} \\ \begin{bmatrix} 0 & 0 & 1 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & -1 \\ 1 & 0 & 0 & 0 \end{bmatrix} \end{matrix} = \begin{matrix} \Sigma \in \mathbb{R}^{D^T \times N} \\ \begin{bmatrix} 4 & 0 & 0 & 0 & 0 \\ 0 & 3 & 0 & 0 & 0 \\ 0 & 0 & \sqrt{5} & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \end{bmatrix} \end{matrix} = \begin{matrix} \mathbf{V}^\top \in \mathbb{R}^{N \times N} \\ \begin{bmatrix} 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 \\ \sqrt{0.2} & 0 & 0 & 0 & \sqrt{0.8} \\ 0 & 0 & 0 & 1 & 0 \\ -\sqrt{0.8} & 0 & 0 & 0 & \sqrt{0.2} \end{bmatrix} \end{matrix}$$

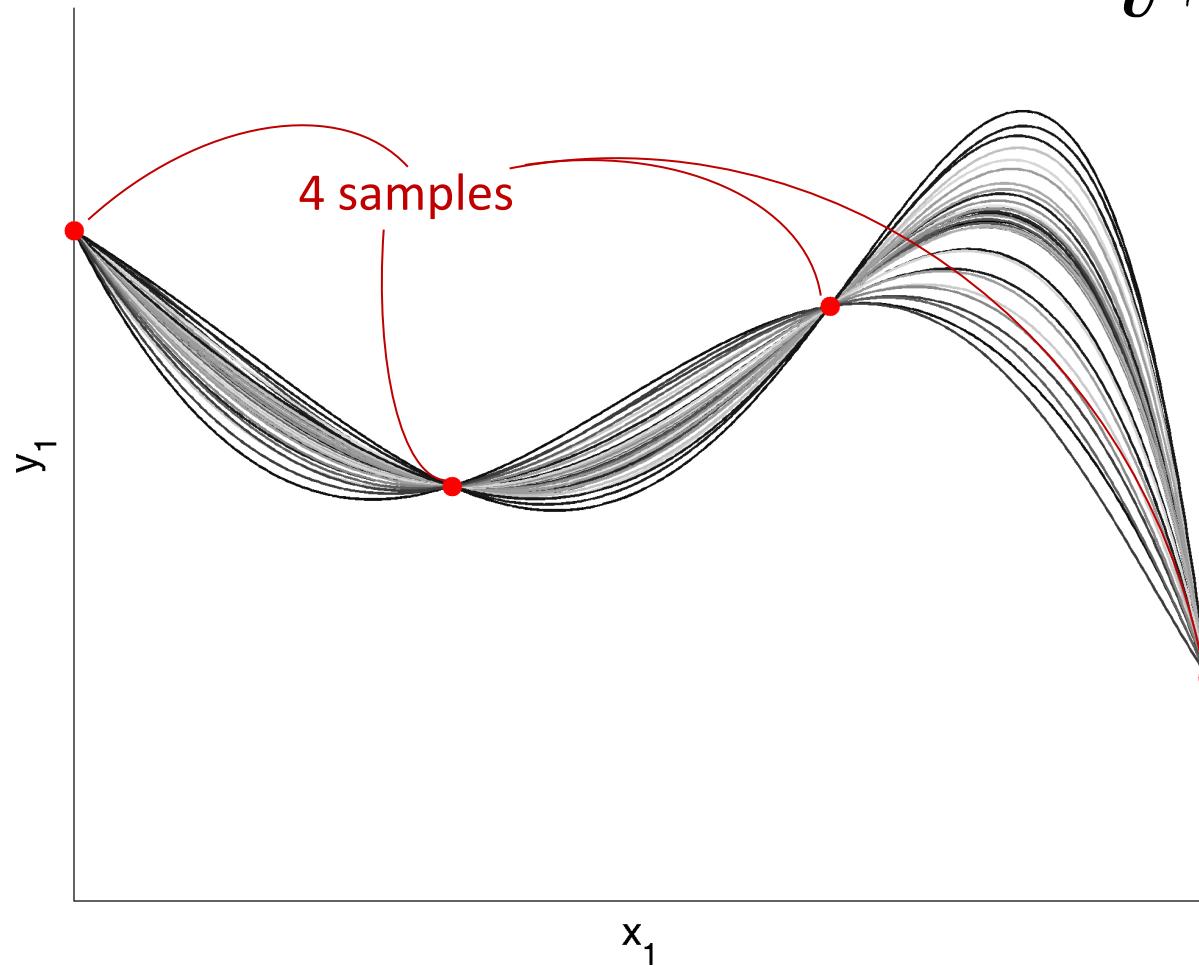
where  $\tilde{\mathbf{U}}$  is a matrix formed by the columns of  $\mathbf{U}$  that span for the corresponding zero rows in  $\Sigma$ .

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)'
```

# Example with polynomial fitting

$$\hat{\mathbf{a}} = \mathbf{X}^\dagger \mathbf{y} + N\mathbf{v} \quad \text{with} \quad \mathbf{x} = [1, x, x^2, \dots, x^6]$$
$$\mathbf{v} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$$

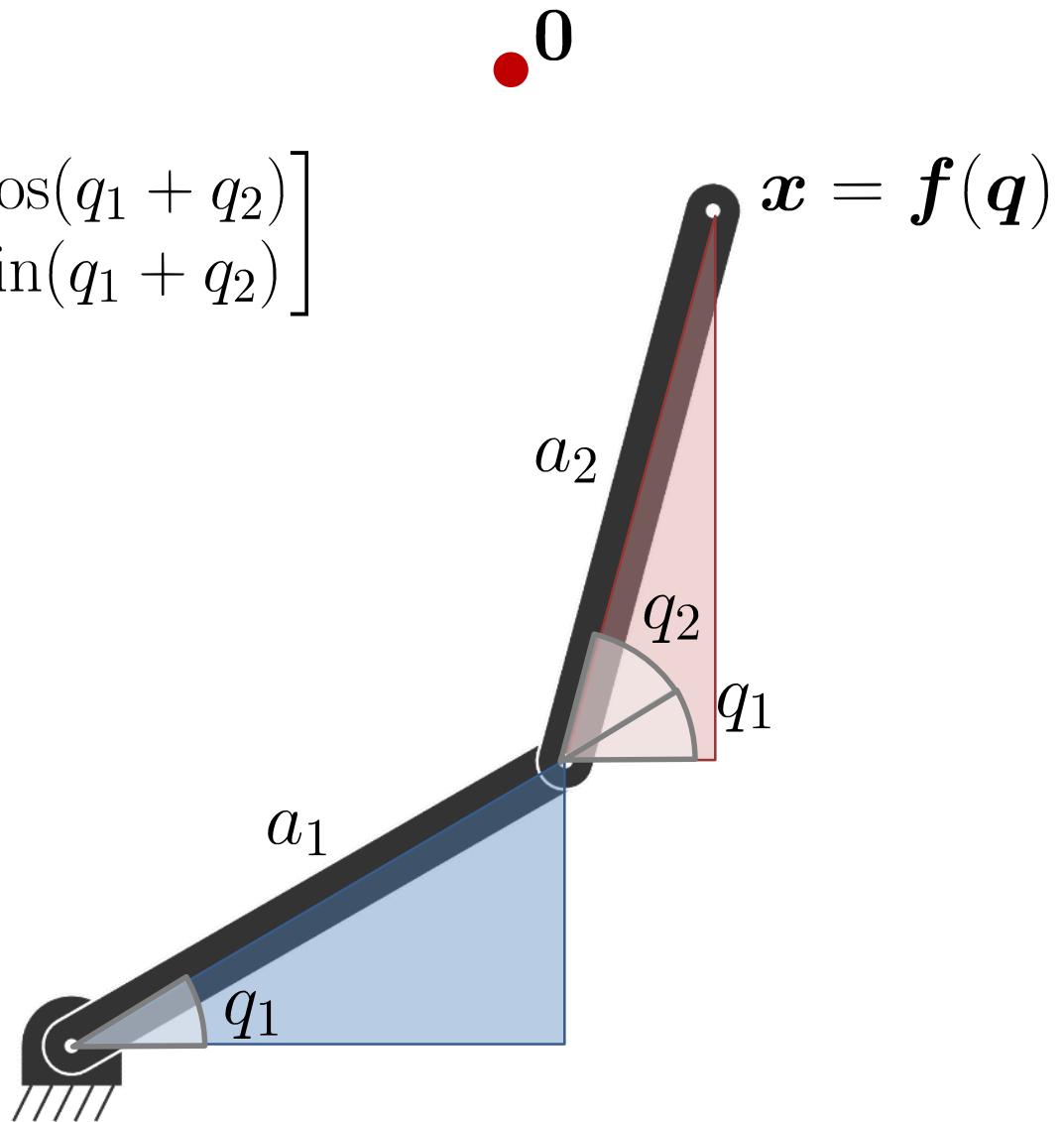


$$\mathbf{X} \in \mathbb{R}^{4 \times 7}$$
$$\mathbf{y} \in \mathbb{R}^4$$
$$\hat{\mathbf{a}} \in \mathbb{R}^7$$

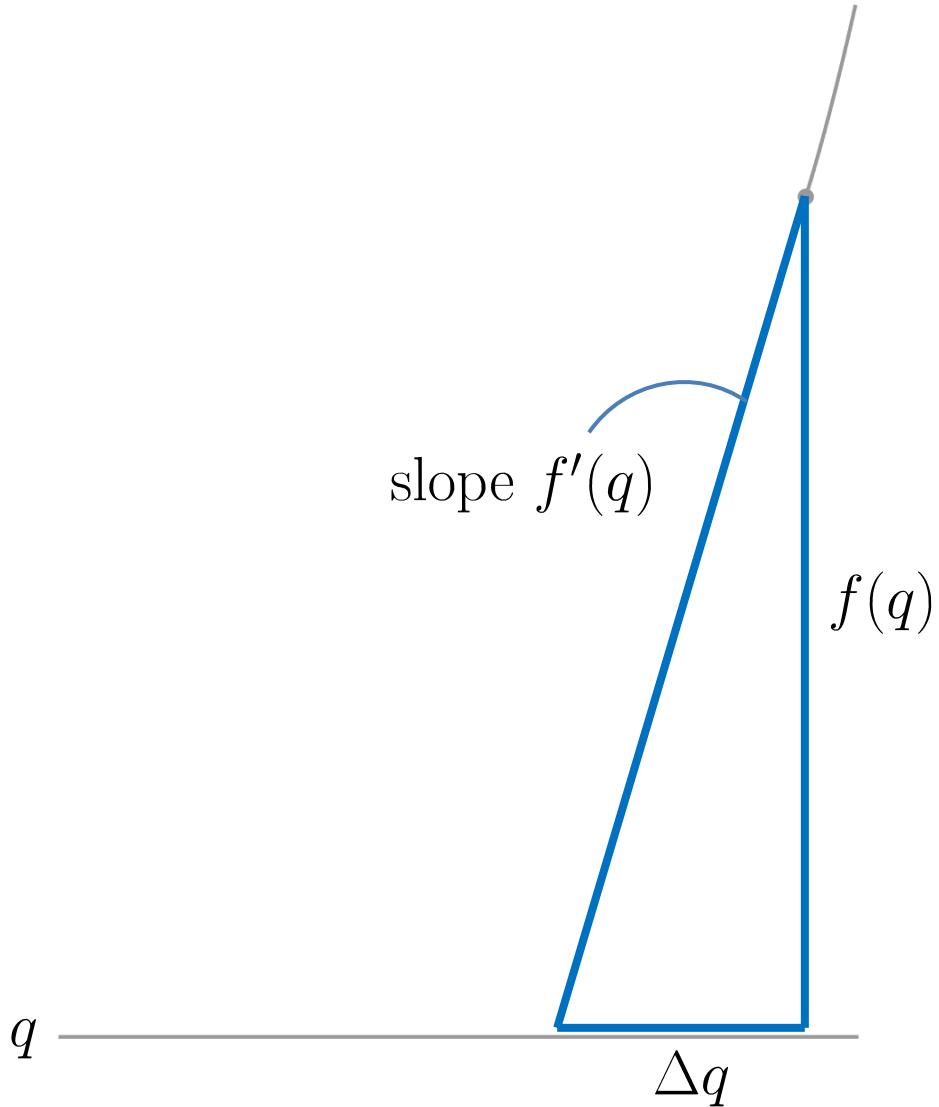
# Example with robot inverse kinematics

Forward kinematics

$$\begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} a_1 \cos(q_1) + a_2 \cos(q_1 + q_2) \\ a_1 \sin(q_1) + a_2 \sin(q_1 + q_2) \end{bmatrix}$$

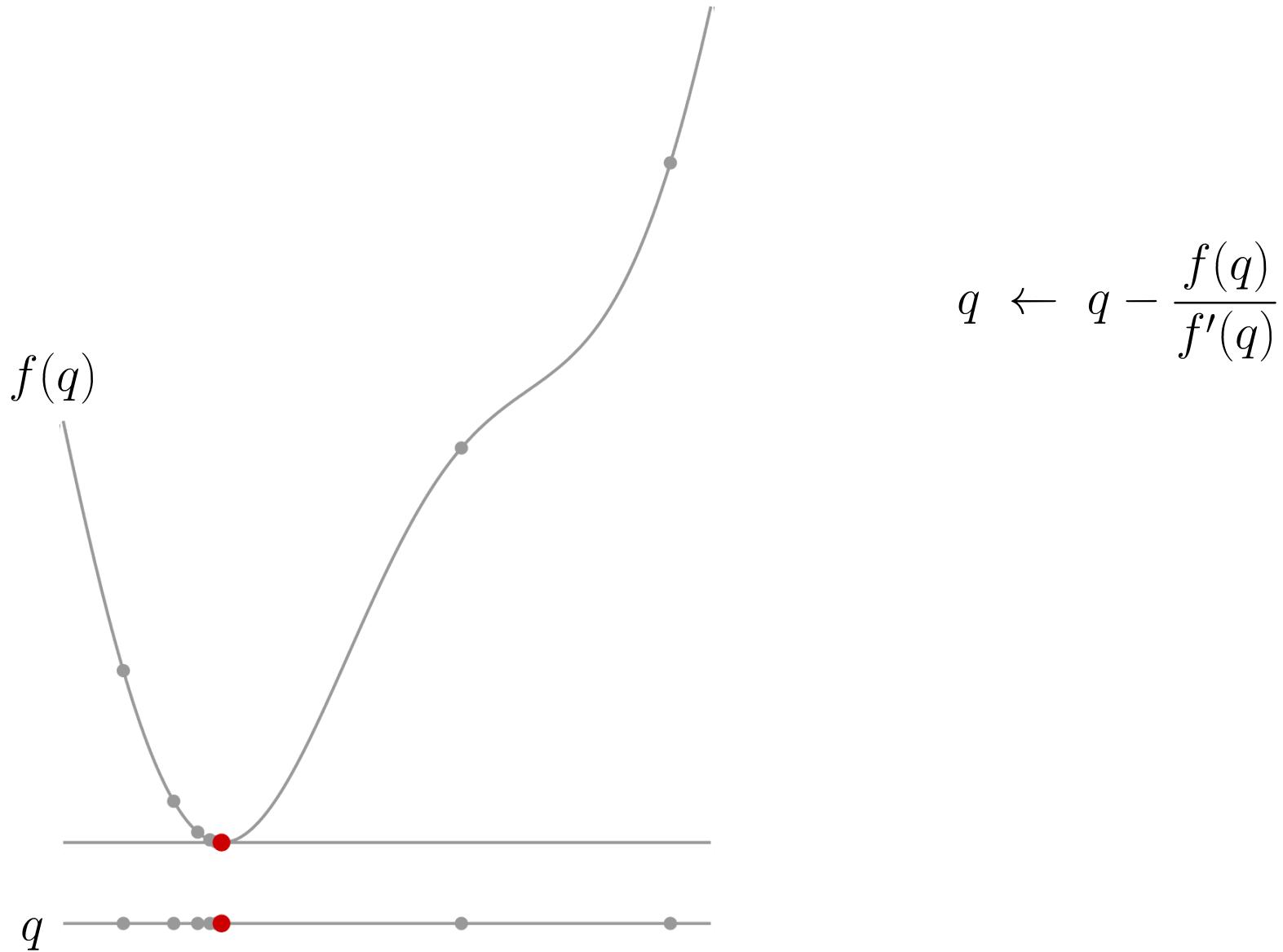


**Find  $q$  to have  $f(q)=0$**



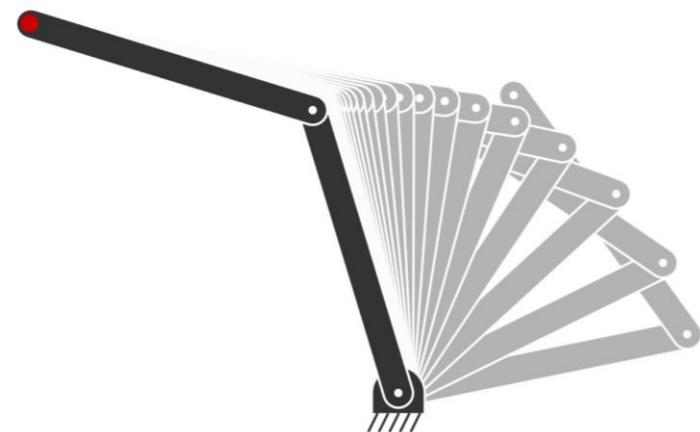
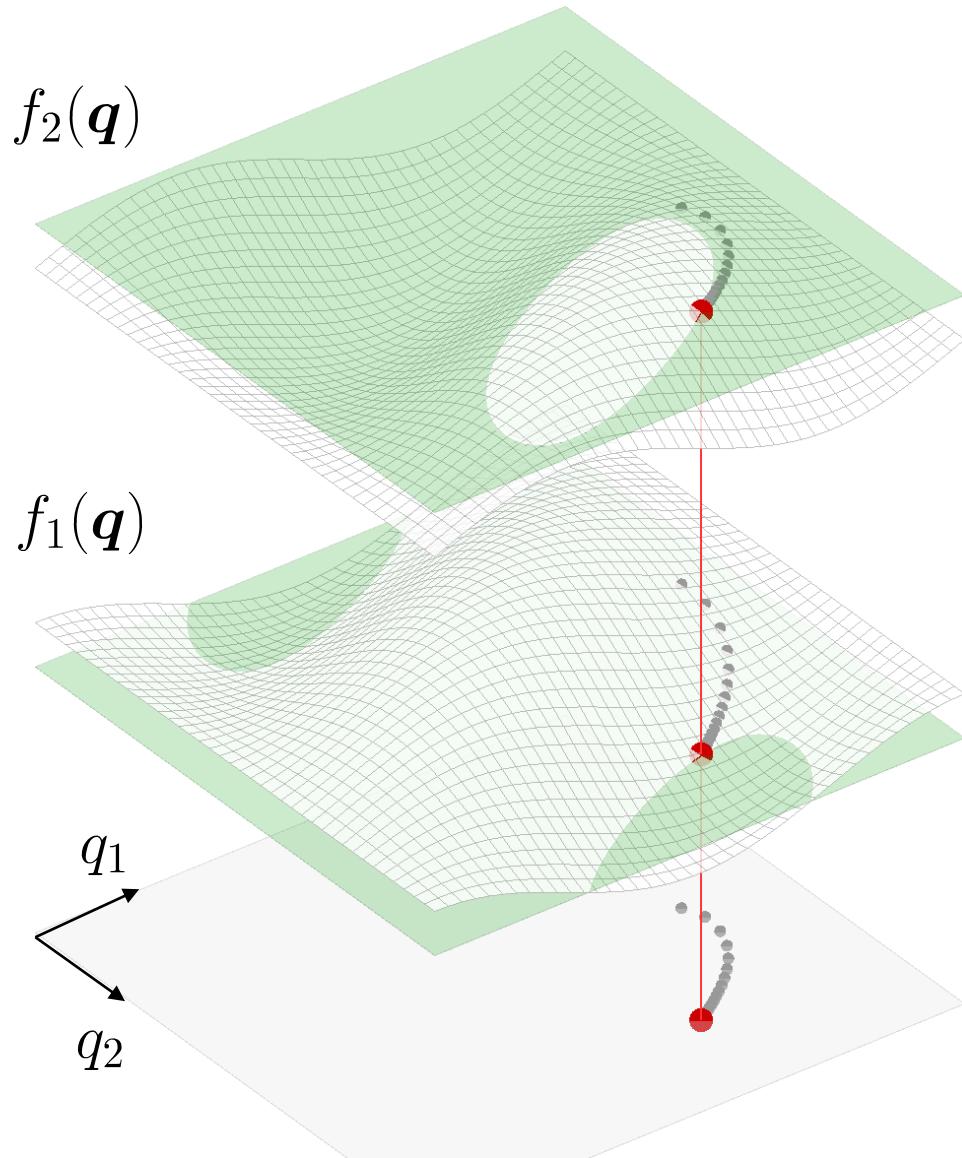
$$f'(q) = \frac{f(q)}{\Delta q}$$
$$\iff \Delta q = \frac{f(q)}{f'(q)}$$

# Gauss-Newton algorithm



# Example with robot inverse kinematics

Gauss-Newton algorithm

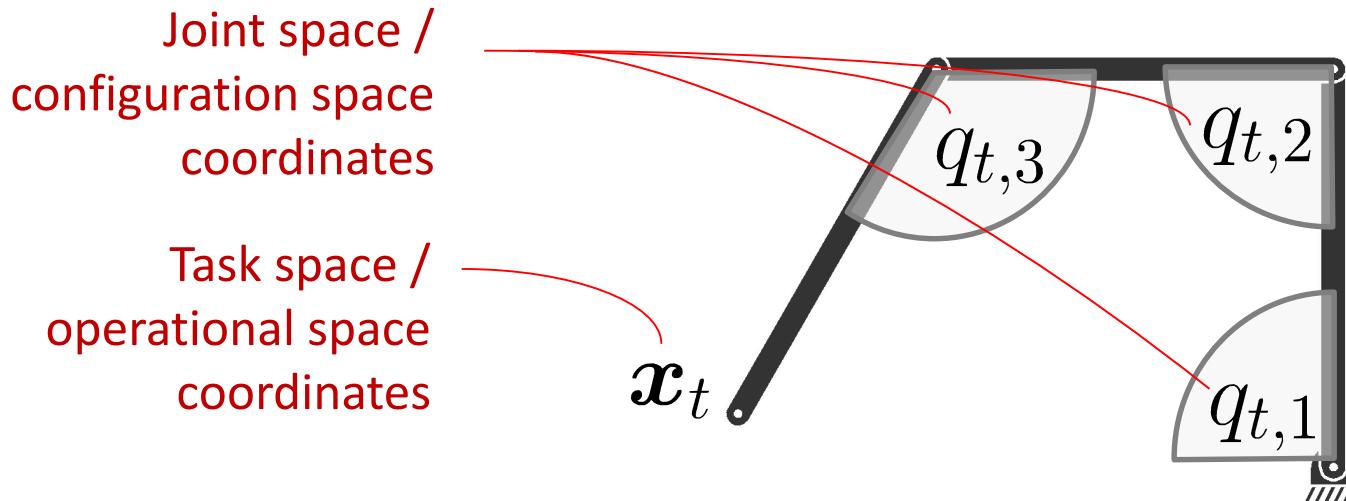


$$\mathbf{q} \leftarrow \mathbf{q} - \alpha \mathbf{J}^\dagger(\mathbf{q}) \mathbf{f}(\mathbf{q})$$

$$\mathbf{J}(\mathbf{q}) = \begin{bmatrix} \frac{\partial f_1(\mathbf{q})}{\partial q_1} & \frac{\partial f_1(\mathbf{q})}{\partial q_2} \\ \frac{\partial f_2(\mathbf{q})}{\partial q_1} & \frac{\partial f_2(\mathbf{q})}{\partial q_2} \end{bmatrix}$$

$$\in \mathbb{R}^{2 \times 2}$$

# Example with robot inverse kinematics



Forward kinematics is computed with

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

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

An inverse kinematics solution can be computed with

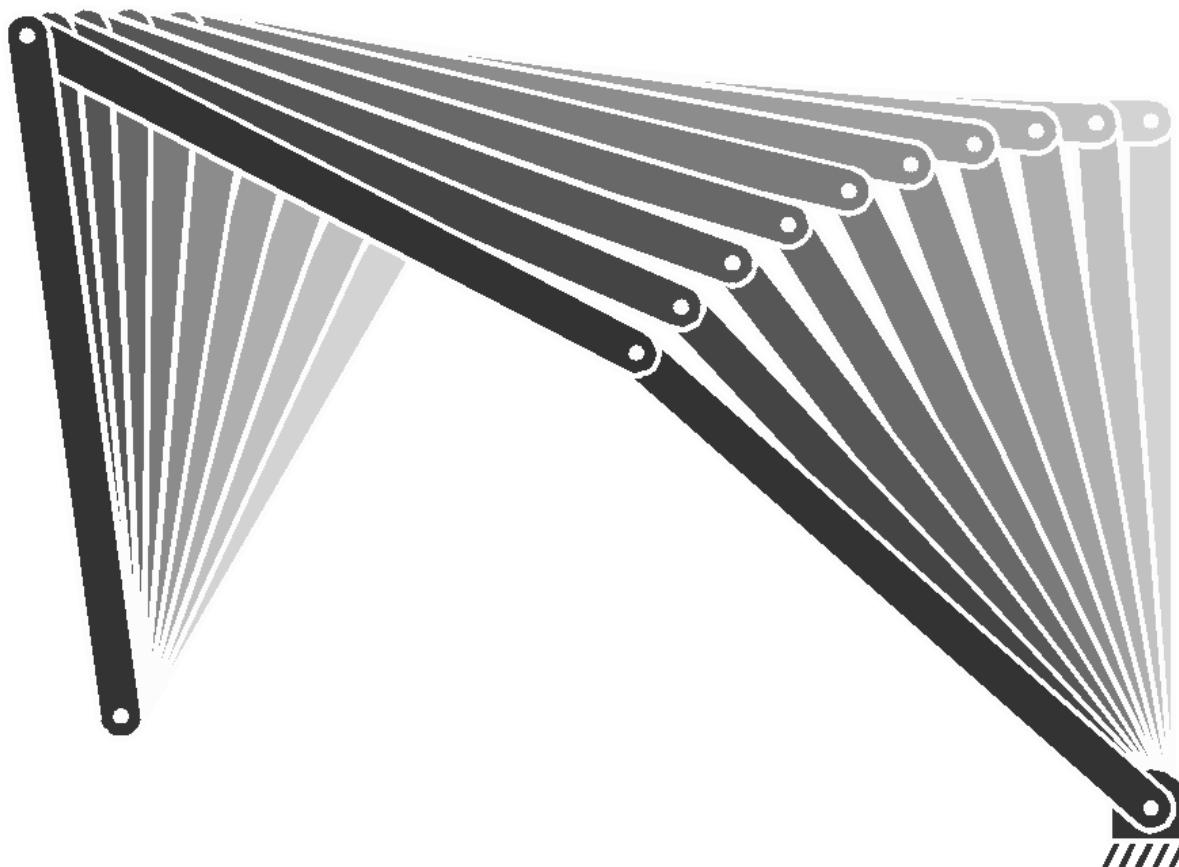
$$\hat{\dot{\boldsymbol{q}}}_t = \boldsymbol{J}^\dagger(\boldsymbol{q}_t) \dot{\boldsymbol{x}}_t + \boldsymbol{N}(\boldsymbol{q}_t) g(\boldsymbol{q}_t)$$

# Example with robot inverse kinematics

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

→ Primary constraint:  
keeping the tip  
of the robot still

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



→ Secondary  
constraint:  
trying to move  
the first joint

# Example with robot inverse kinematics

$$\hat{\dot{q}}_t = \mathbf{J}_t^{\mathcal{R}\dagger} \dot{x}_t^{\mathcal{R}} + \mathbf{N}_t^{\mathcal{R}} \dot{q}_t^{\mathcal{L}}$$

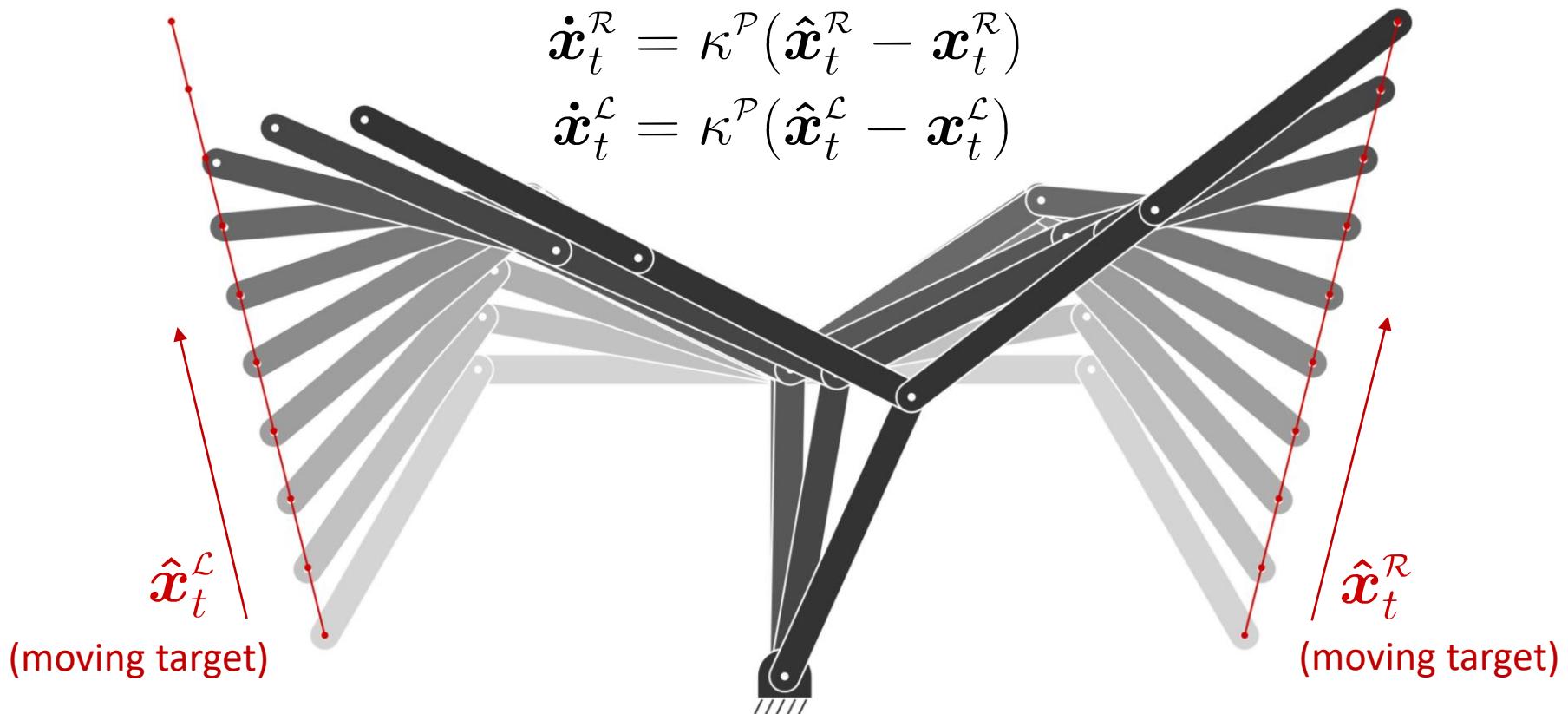
Tracking target  
with left hand,  
if possible

Tracking  
target with  
right hand

$$\dot{q}_t^{\mathcal{L}} = (\mathbf{J}_t^{\mathcal{L}} \mathbf{N}_t^{\mathcal{R}})^{\dagger} (\dot{x}_t^{\mathcal{L}} - \mathbf{J}_t^{\mathcal{L}} \mathbf{J}_t^{\mathcal{R}\dagger} \dot{x}_t^{\mathcal{R}})$$

$$\dot{x}_t^{\mathcal{R}} = \kappa^{\mathcal{P}} (\hat{x}_t^{\mathcal{R}} - x_t^{\mathcal{R}})$$

$$\dot{x}_t^{\mathcal{L}} = \kappa^{\mathcal{P}} (\hat{x}_t^{\mathcal{L}} - x_t^{\mathcal{L}})$$



# Ridge regression (Tikhonov regularization, penalized least squares)

Python notebook:  
**demo\_LS\_polFit.ipynb**

Matlab example:  
**demo\_LS\_polFit02.m**

# Ridge regression (Tikhonov regularization)

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

$$\begin{aligned}\hat{\mathbf{A}} &= \arg \min_{\mathbf{A}} \|\mathbf{Y} - \mathbf{X}\mathbf{A}\|_{\text{F}}^2 + \|\boldsymbol{\Gamma}\mathbf{A}\|_{\text{F}}^2 \\ &= \arg \min_{\mathbf{A}} \text{tr}\left((\mathbf{Y} - \mathbf{X}\mathbf{A})^\top(\mathbf{Y} - \mathbf{X}\mathbf{A})\right) + \text{tr}\left((\boldsymbol{\Gamma}\mathbf{A})^\top\boldsymbol{\Gamma}\mathbf{A}\right)\end{aligned}$$

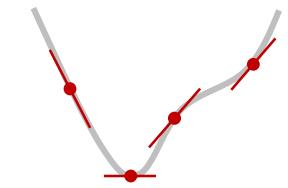
By differentiating with respect to  $\mathbf{A}$  and equating to zero, we can see that

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

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  **$\ell_2$  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  $\mathbf{0} \in \mathbb{R}^{D^I \times D^O}$  and  $\boldsymbol{\Gamma} \in \mathbb{R}^{D^I \times D^I}$ , yielding

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

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

# 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  $\sigma_i$  the singular values in the diagonal of  $\Sigma$ , a solution to the ridge regression problem is given by

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

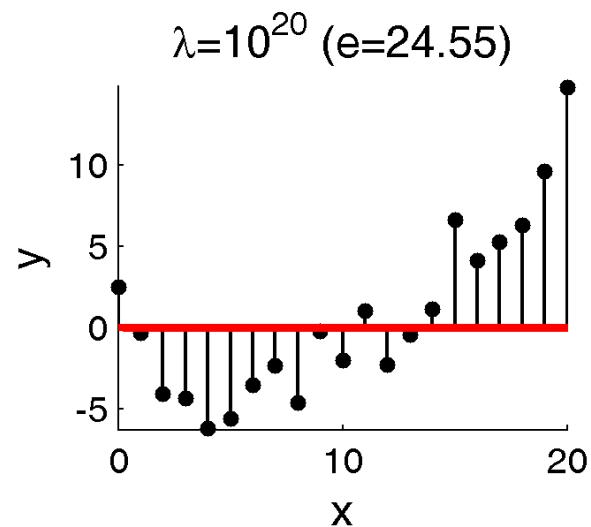
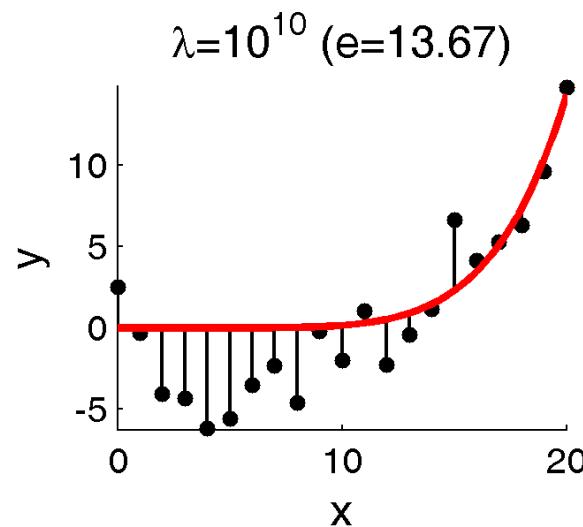
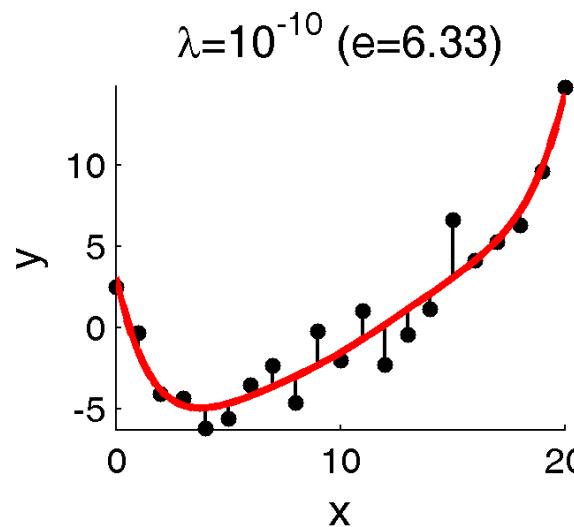
where  $\tilde{\Sigma}$  has diagonal values

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

and has zeros elsewhere.

# Ridge regression (Tikhonov regularization)

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



# **Weighted least squares (Generalized least squares)**

**Python notebook:  
demo\_LS\_weighted.ipynb**

**Matlab example:  
demo\_LS\_weighted01.m**

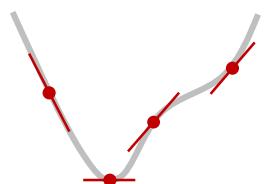
# Weighted least squares

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

$$\begin{aligned}\hat{\mathbf{A}} &= \arg \min_{\mathbf{A}} \|\mathbf{Y} - \mathbf{X}\mathbf{A}\|_{F,\mathbf{W}}^2 \\ &= \arg \min_{\mathbf{A}} \text{tr}\left((\mathbf{Y} - \mathbf{X}\mathbf{A})^\top \mathbf{W}(\mathbf{Y} - \mathbf{X}\mathbf{A})\right) \\ &= \arg \min_{\mathbf{A}} \text{tr}(\mathbf{Y}^\top \mathbf{W} \mathbf{Y} - 2\mathbf{A}^\top \mathbf{X}^\top \mathbf{W} \mathbf{Y} + \mathbf{A}^\top \mathbf{X}^\top \mathbf{W} \mathbf{X} \mathbf{A})\end{aligned}$$

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

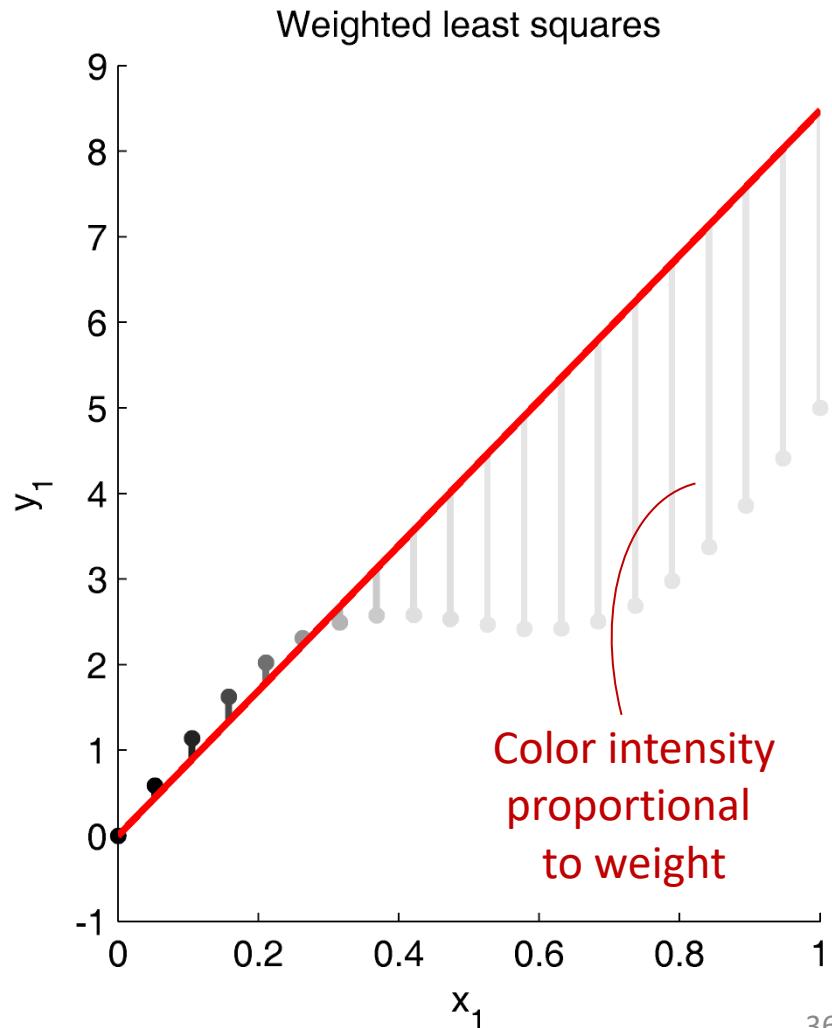
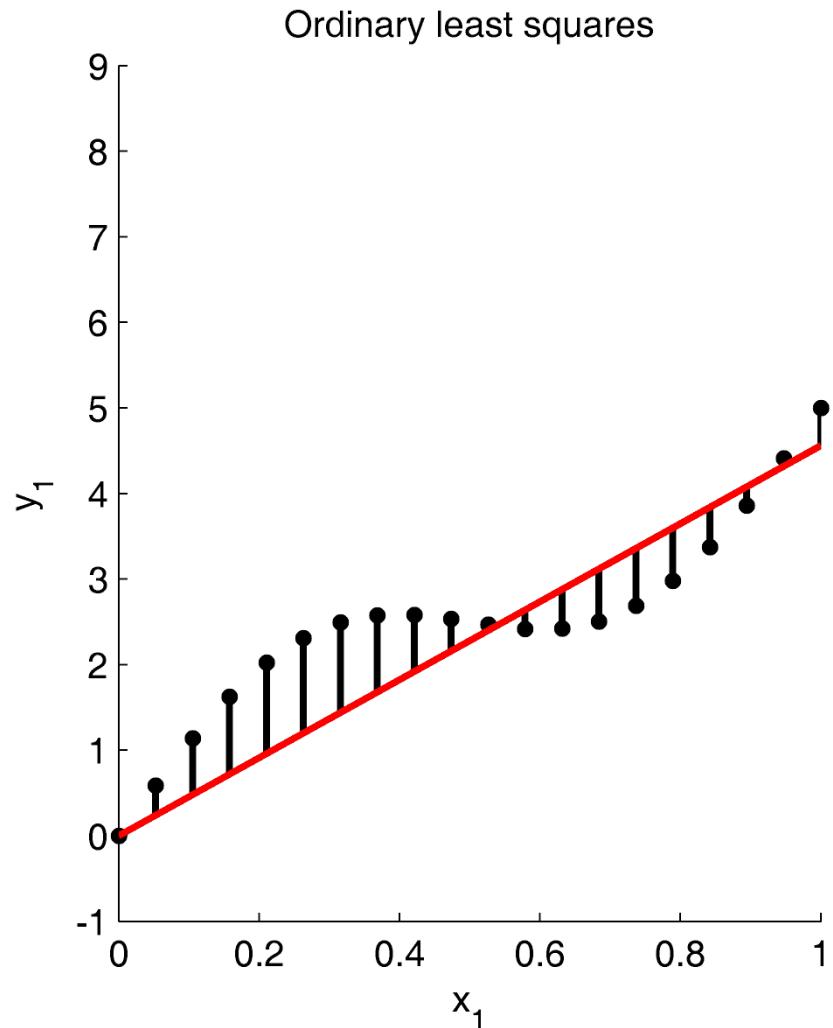
$$-2\mathbf{X}^\top \mathbf{W} \mathbf{Y} + 2\mathbf{X}^\top \mathbf{W} \mathbf{X} \mathbf{A} = \mathbf{0} \iff \hat{\mathbf{A}} = (\mathbf{X}^\top \mathbf{W} \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{W} \mathbf{Y}$$



$$\mathbf{X}_\mathbf{W}^\dagger$$

# Weighted least squares

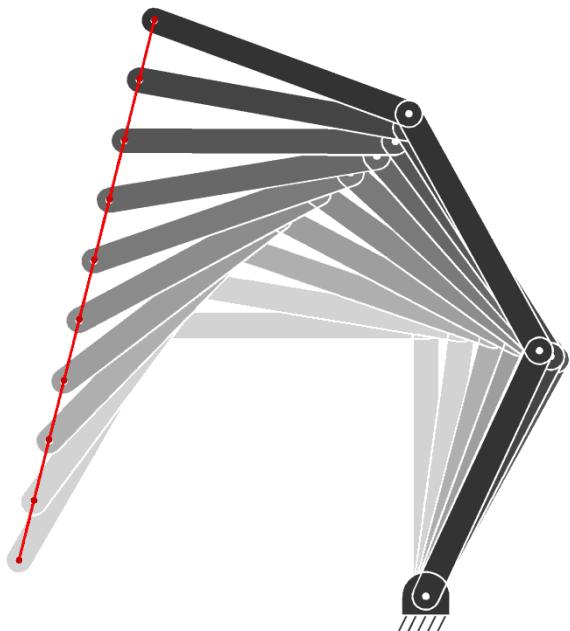
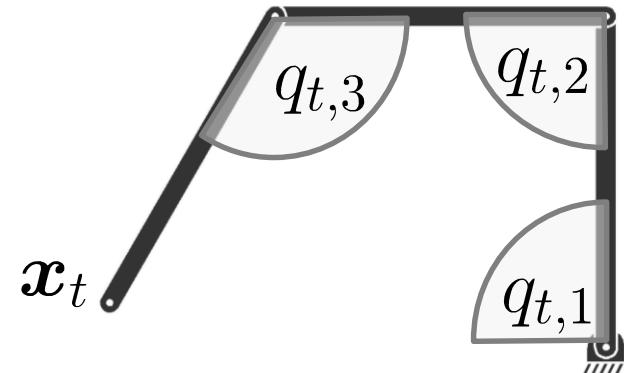
$$\hat{\mathbf{A}} = (\mathbf{X}^\top \mathbf{W} \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{W} \mathbf{Y}$$



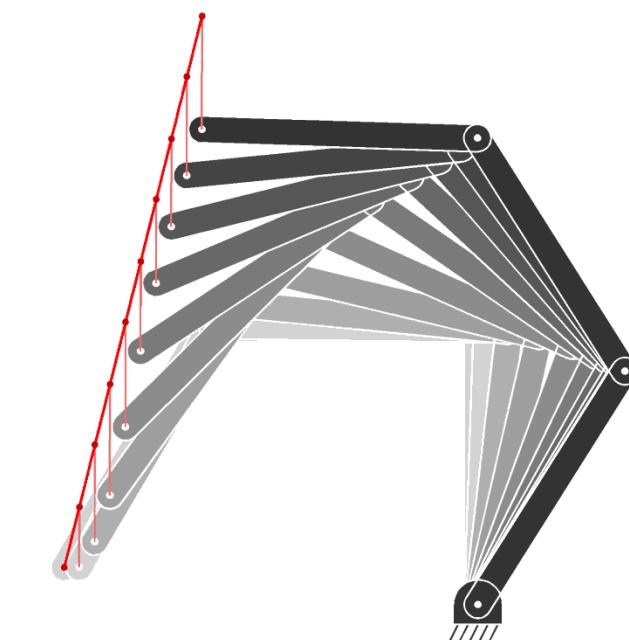
# Weighted least squares - Example I

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

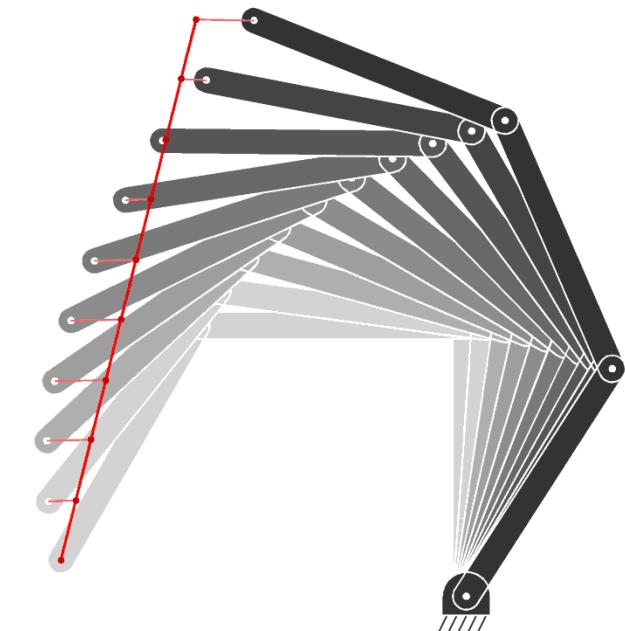
$$\hat{\dot{\mathbf{q}}}_t = (\mathbf{J}^\top \mathbf{W}^\chi \mathbf{J})^{-1} \mathbf{J}^\top \mathbf{W}^\chi \dot{\mathbf{x}}_t$$



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



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

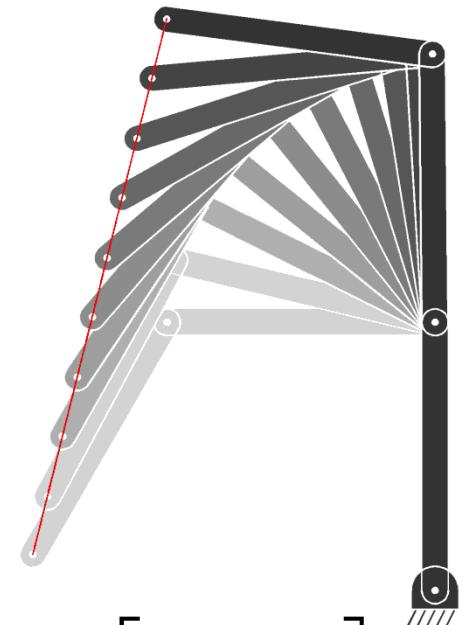
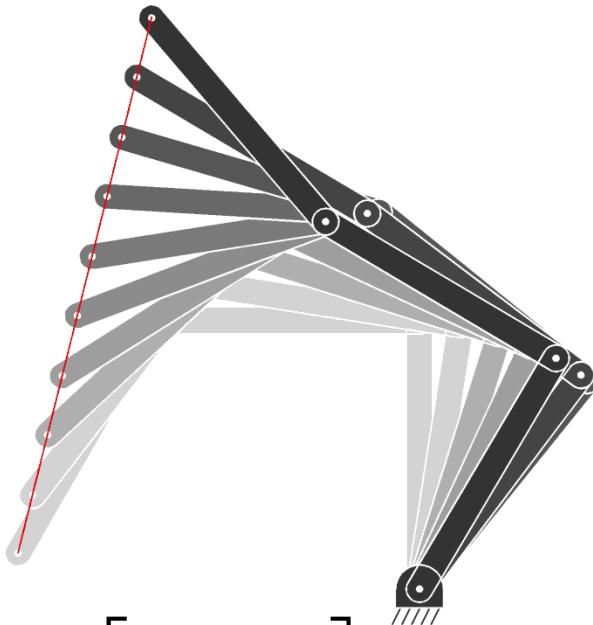
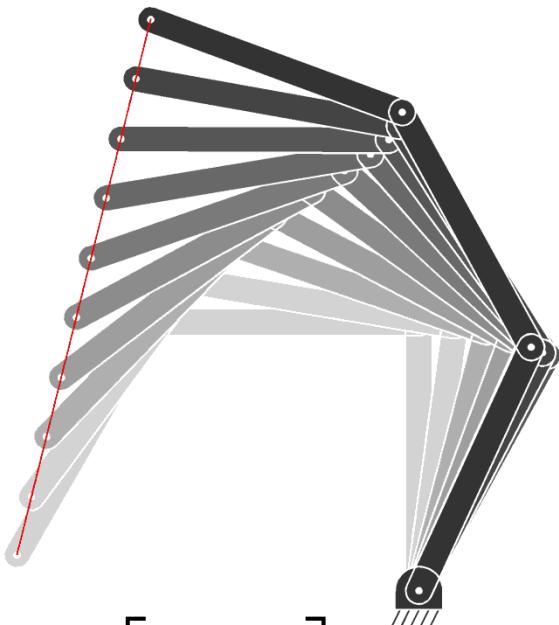
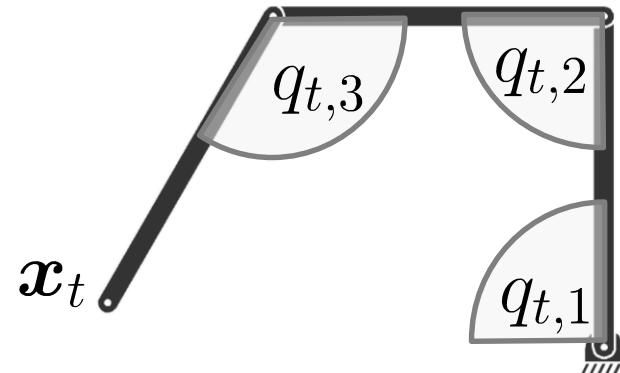


$$\mathbf{W}^\chi = \begin{bmatrix} .01 & 0 \\ 0 & 1 \end{bmatrix}$$

# Weighted least squares - Example II

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

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



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

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

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

# **Iteratively reweighted least squares (IRLS)**

**Python notebook:**  
**demo\_LS\_weighted.ipynb**

**Matlab code:**  
**demo\_LS\_IRLS01.m**

# Iteratively reweighted least squares (IRLS)

- **Iteratively Reweighted Least Squares** generalizes least squares by raising the error to a power that is less than 2:  
→ can no longer be called “least squares”
- The strategy is that an error  $|e|^p$  can be rewritten as  $|e|^p = |e|^{p-2} e^2$ .
- $|e|^{p-2}$  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)

$$|\mathbf{e}|^p = |\mathbf{e}|^{p-2} \mathbf{e}^2$$

For an  $\ell_p$  norm objective defined by

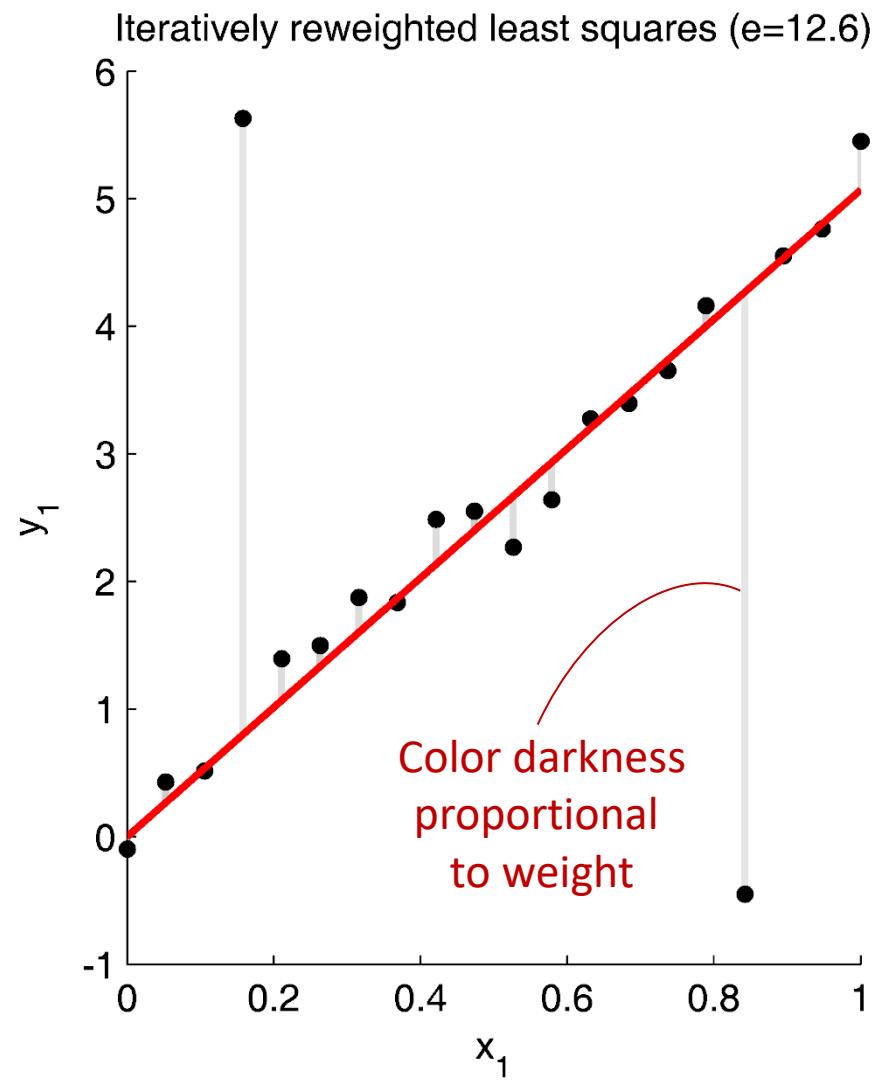
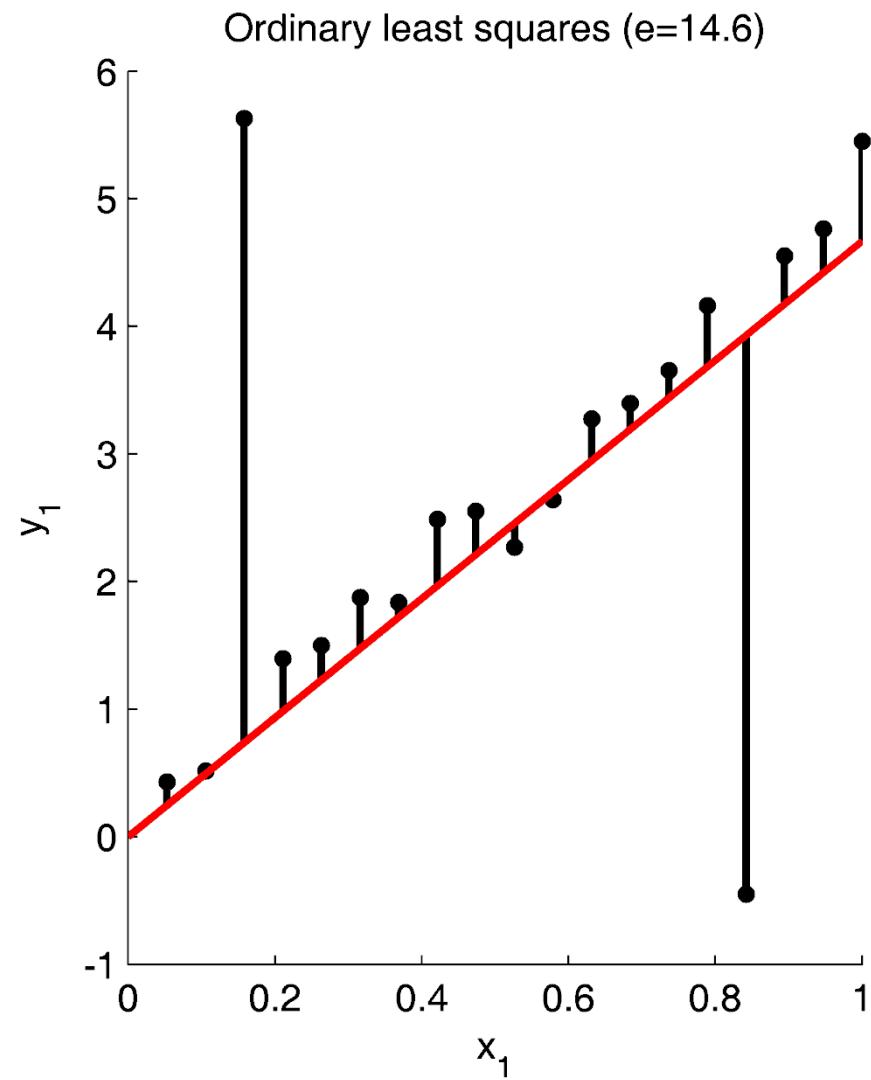
$$\hat{\mathbf{A}} = \arg \min_{\mathbf{A}} \|\mathbf{Y} - \mathbf{X}\mathbf{A}\|_{\text{F},p}^2$$

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

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

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

# IRLS as regression robust to outliers



# **Recursive least squares**

**Python notebook:**  
**demo\_LS\_recursive.ipynb**

**Matlab code:**  
**demo\_LS\_recursive01.m**

# Recursive least squares

Sherman-Morrison-Woodbury relation:

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

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  $\mathbf{B} + \mathbf{U}\mathbf{V}$ .

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

# Recursive least squares

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

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

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

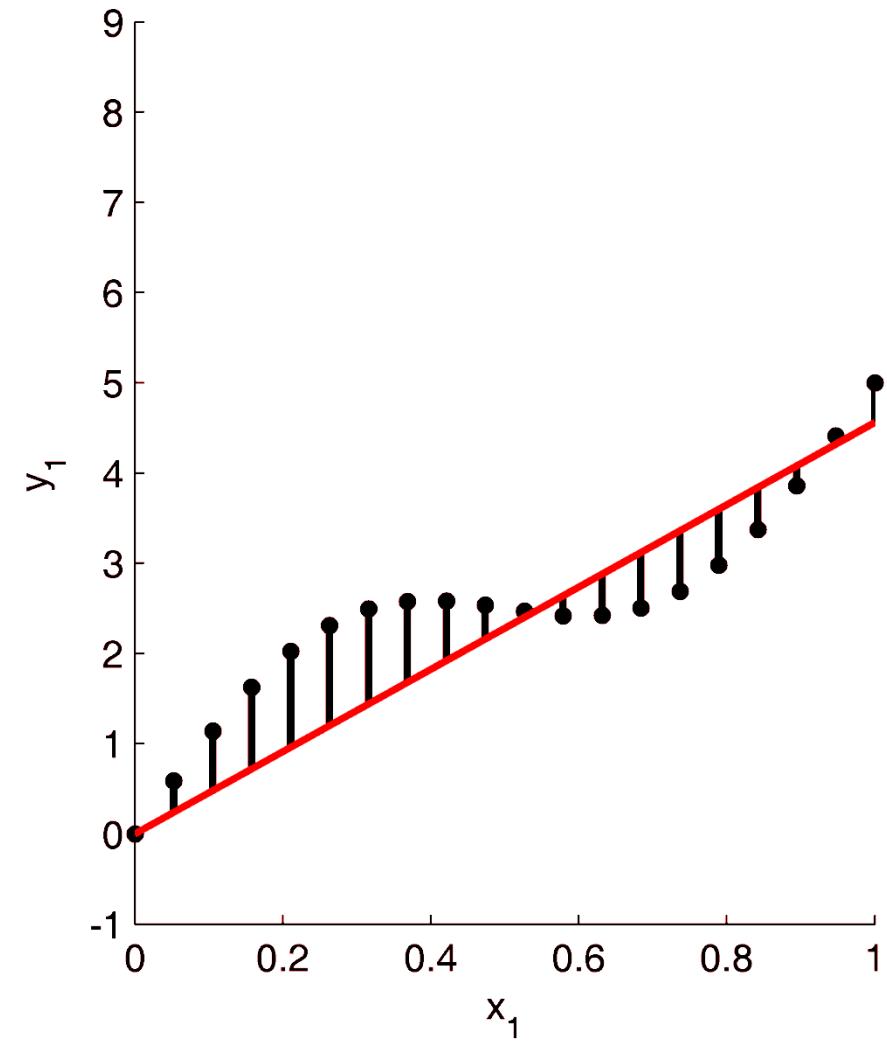
which is exploited to efficiently compute the update as

$$\hat{\mathbf{A}}_{\text{new}} = \hat{\mathbf{A}} + \mathbf{K} (\mathbf{C} - \mathbf{V} \hat{\mathbf{A}})$$

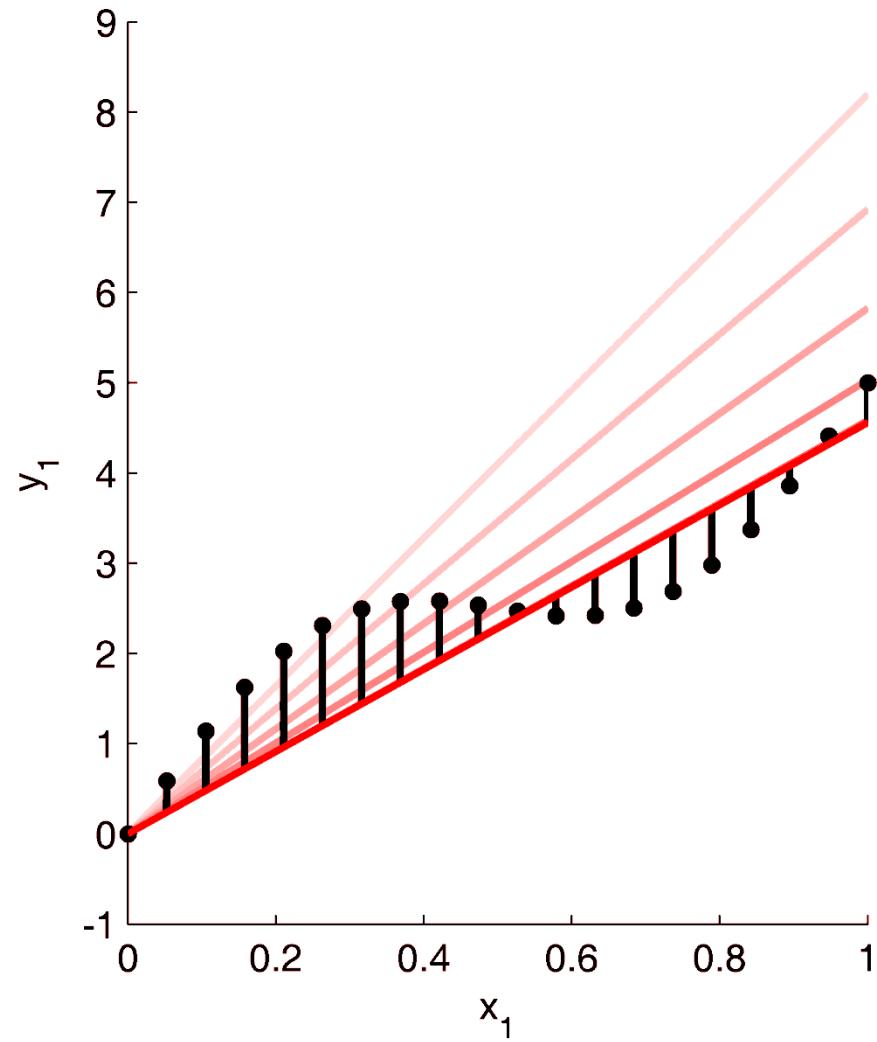
with Kalman gain  $\mathbf{K} = \mathbf{B}^{-1} \mathbf{V}^\top (\mathbf{I} + \mathbf{V} \mathbf{B}^{-1} \mathbf{V}^\top)^{-1}$

# Recursive least squares

Ordinary least squares ( $e=11.0$ )



Recursive least squares ( $e=11.0$ )



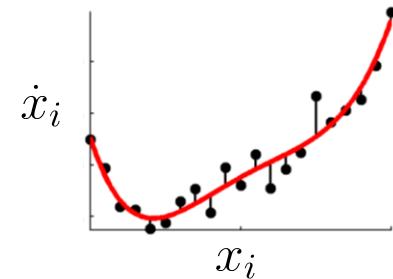
# **Linear regression: Examples of applications**

# Koopman operators in control

Nonlinear

$$\dot{x} = f(x)$$

$$\begin{bmatrix} \dot{x}_1 \\ \dot{x}_2 \end{bmatrix} = \begin{bmatrix} \lambda_1 x_1 \\ \lambda_2(x_2 - x_1^2) \end{bmatrix}$$



$$\dot{y} = A y$$

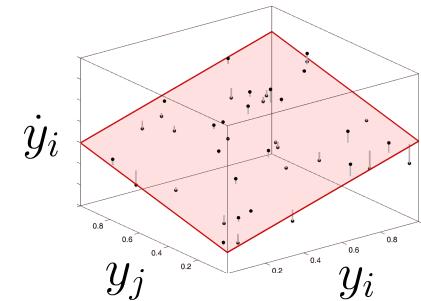
$$\begin{bmatrix} \dot{y}_1 \\ \dot{y}_2 \\ \dot{y}_3 \end{bmatrix} = \begin{bmatrix} \lambda_1 & 0 & 0 \\ 0 & \lambda_2 & -\lambda_2 \\ 0 & 0 & 2\lambda_1 \end{bmatrix} \begin{bmatrix} y_1 \\ y_2 \\ y_3 \end{bmatrix}$$

with

$$\begin{bmatrix} y_1 \\ y_2 \\ y_3 \end{bmatrix} = \begin{bmatrix} x_1 \\ x_2 \\ x_1^2 \end{bmatrix}$$

$$\begin{aligned} \dot{y}_3 &= \frac{\partial y_3}{\partial x_1} \dot{x}_1 \\ &= 2x_1 \lambda_1 x_1 \\ &= 2\lambda_1 y_3 \end{aligned}$$

Linear in state space  
of higher dimension



Main challenge in Koopman analysis:  
How to find these basis functions?

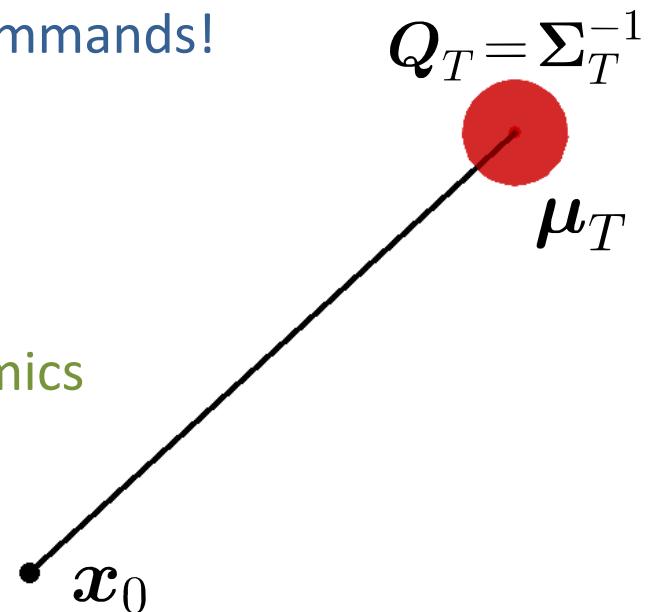
# Linear quadratic tracking (LQT)

$$\min_{\boldsymbol{u}} \sum_{t=1}^T \left\| \boldsymbol{\mu}_t - \boldsymbol{x}_t \right\|_{Q_t}^2 + \left\| \boldsymbol{u}_t \right\|_{R_t}^2$$

Track path!    Use low control commands!

$$\text{s.t. } \boldsymbol{x}_{t+1} = \boldsymbol{A}\boldsymbol{x}_t + \boldsymbol{B}\boldsymbol{u}_t$$

System dynamics



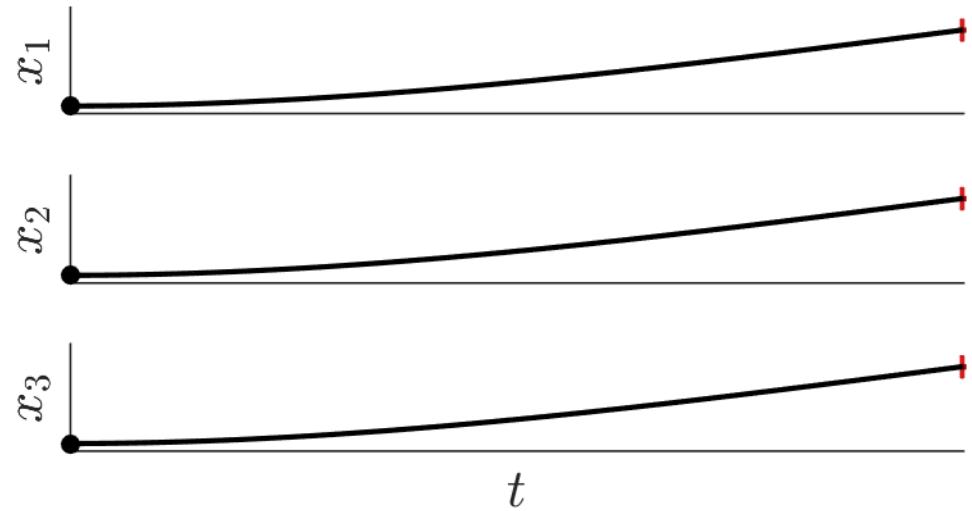
$\boldsymbol{x}_t$  state variable (position+velocity)

$\boldsymbol{\mu}_t$  desired state

$\boldsymbol{u}_t$  control command (acceleration)

$\boldsymbol{Q}_t$  precision matrix

$\boldsymbol{R}_t$  control weight matrix



# How to solve this objective function?

$$\min_{\boldsymbol{u}} \sum_{t=1}^T \left\| \boldsymbol{\mu}_t - \boldsymbol{x}_t \right\|_{Q_t}^2 + \left\| \boldsymbol{u}_t \right\|_{R_t}^2$$

Track path!

Use low control commands!

s.t.  $\boldsymbol{x}_{t+1} = \boldsymbol{A}\boldsymbol{x}_t + \boldsymbol{B}\boldsymbol{u}_t$

System dynamics

**Pontryagin's max. principle,  
Riccati equation,  
Hamilton-Jacobi-Bellman**  
*(the Physicist perspective)*



**Dynamic programming**  
*(the Computer Scientist perspective)*



**Linear algebra**  
*(the Algebraist perspective)*



# Let's first re-organize the objective function...

$$\begin{aligned}
 c &= \sum_{t=1}^T \left( (\mu_t - x_t)^\top Q_t (\mu_t - x_t) + u_t^\top R_t u_t \right) \\
 &= (\mu - x)^\top Q (\mu - x) + u^\top R u
 \end{aligned}$$



$$Q = \begin{bmatrix} Q_1 & 0 & \cdots & 0 \\ 0 & Q_2 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & Q_T \end{bmatrix} \quad R = \begin{bmatrix} R_1 & 0 & \cdots & 0 \\ 0 & R_2 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & R_T \end{bmatrix}$$

$$\mu = \begin{bmatrix} \mu_1 \\ \mu_2 \\ \vdots \\ \mu_T \end{bmatrix} \quad x = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_T \end{bmatrix}$$

$$u = \begin{bmatrix} u_1 \\ u_2 \\ \vdots \\ u_T \end{bmatrix}$$

Let's then re-organize the constraint...

$$\mathbf{x}_{t+1} = \mathbf{A} \mathbf{x}_t + \mathbf{B} \mathbf{u}_t$$



$$\mathbf{x}_2 = \mathbf{A}\mathbf{x}_1 + \mathbf{B}\mathbf{u}_1$$

$$\mathbf{x}_3 = \mathbf{A}\mathbf{x}_2 + \mathbf{B}\mathbf{u}_2 = \mathbf{A}(\mathbf{A}\mathbf{x}_1 + \mathbf{B}\mathbf{u}_1) + \mathbf{B}\mathbf{u}_2$$

$$\vdots$$

$$\mathbf{x}_T = \mathbf{A}^{T-1}\mathbf{x}_1 + \mathbf{A}^{T-2}\mathbf{B}\mathbf{u}_1 + \mathbf{A}^{T-3}\mathbf{B}\mathbf{u}_2 + \cdots + \mathbf{B}_{T-1}\mathbf{u}_{T-1}$$

$$\begin{bmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \\ \mathbf{x}_3 \\ \vdots \\ \mathbf{x}_T \end{bmatrix} = \underbrace{\begin{bmatrix} \mathbf{I} \\ \mathbf{A} \\ \mathbf{A}^2 \\ \vdots \\ \mathbf{A}^{T-1} \end{bmatrix}}_{\mathbf{S}^x} \mathbf{x}_1 + \underbrace{\begin{bmatrix} \mathbf{0} & \mathbf{0} & \cdots & \mathbf{0} & \mathbf{0} \\ \mathbf{B} & \mathbf{0} & \cdots & \mathbf{0} & \mathbf{0} \\ \mathbf{AB} & \mathbf{B} & \cdots & \mathbf{0} & \mathbf{0} \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ \mathbf{A}^{T-2}\mathbf{B} & \mathbf{A}^{T-3}\mathbf{B} & \cdots & \mathbf{B} & \mathbf{0} \end{bmatrix}}_{\mathbf{S}^u} \begin{bmatrix} \mathbf{u}_1 \\ \mathbf{u}_2 \\ \vdots \\ \mathbf{u}_T \end{bmatrix}$$

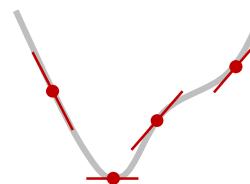
$$\mathbf{x} = \mathbf{S}^x \mathbf{x}_1 + \mathbf{S}^u \mathbf{u}$$

# Linear quadratic tracking (LQT)

The constraint can then be put into the objective function:

$$\begin{aligned}x &= S^x x_1 + S^u u \\c &= (\mu - x)^\top Q (\mu - x) + u^\top R u \\&= (\mu - S^x x_1 - S^u u)^\top Q (\mu - S^x x_1 - S^u u) + u^\top R u\end{aligned}$$

Solving for  $u$  results in the analytic solution:



$$\hat{u} = (S^{u^\top} Q S^u + R)^{-1} S^{u^\top} Q (\mu - S^x x_1)$$

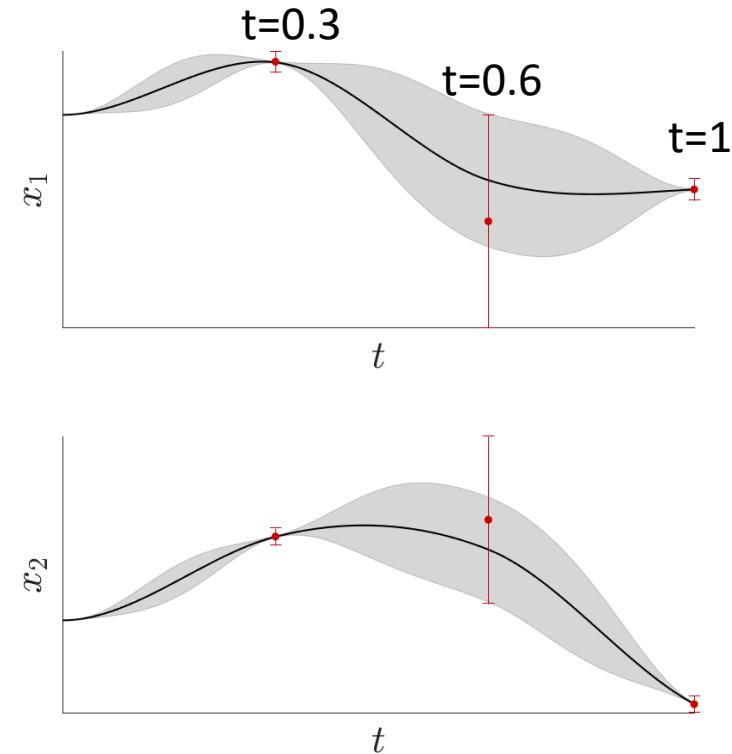
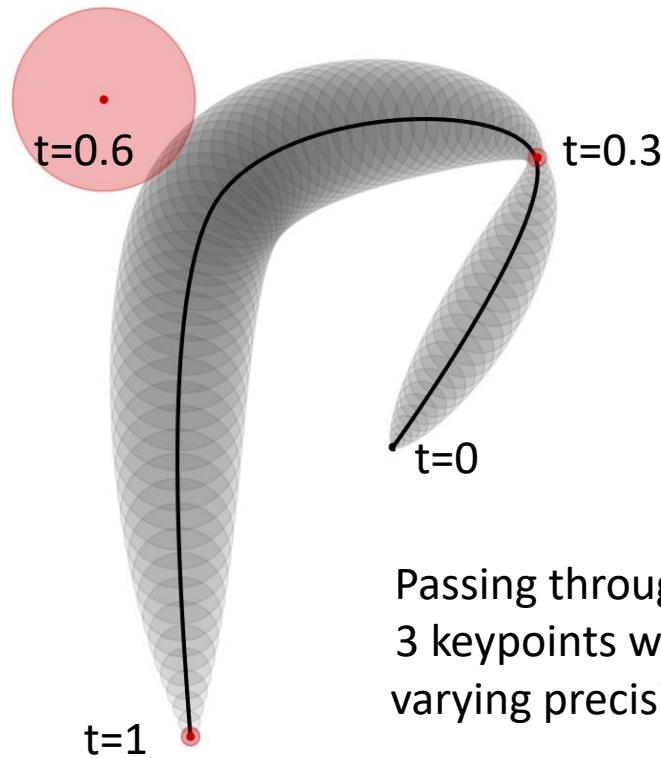
# Linear quadratic tracking (LQT)

$$\hat{u} = (S^{u^\top} Q S^u + R)^{-1} S^{u^\top} Q (\mu - S^x x_1)$$
$$\hat{\Sigma}^u = (S^{u^\top} Q S^u + R)^{-1}$$



$$\hat{x} = S^x x_1 + S^u \hat{u}$$
$$\hat{\Sigma}^x = S^u (S^{u^\top} Q S^u + R)^{-1} S^{u^\top}$$

The distribution in control space can  
be projected back to the state space



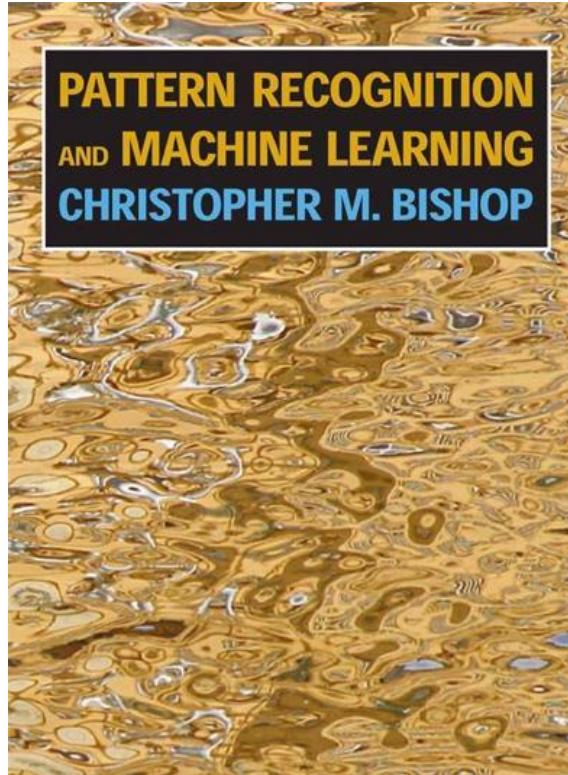
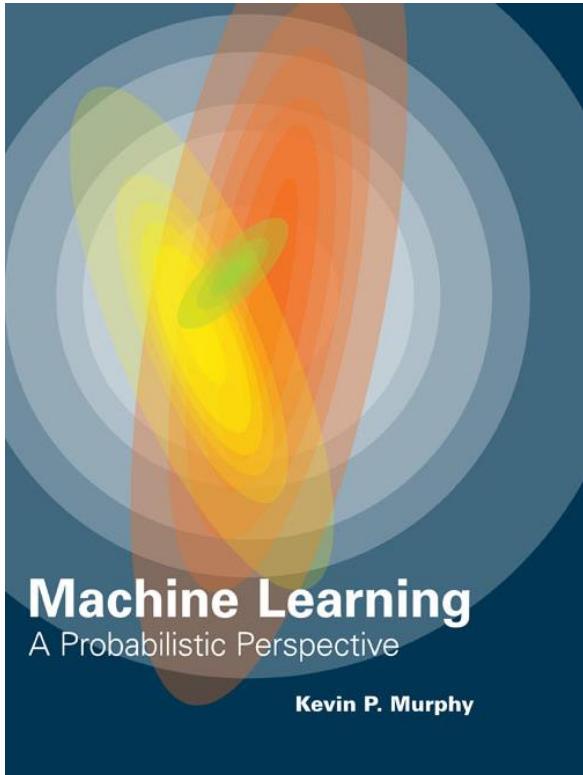
# References

## 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