

EE613 - Machine Learning for Engineers

LINEAR REGRESSION

Sylvain Calinon

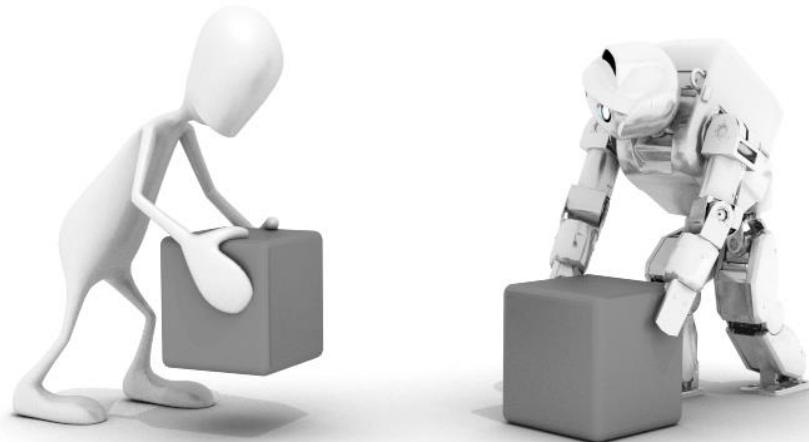
Robot Learning and Interaction Group

Idiap Research Institute

Sep 30, 2021

Robot Learning and Interaction group

Idiap Research Institute



Learning from demonstration

EE613 labs (INM11):
Dr João Silvério



Outline

Linear regression (Sep 30)

- Least squares
- Singular value decomposition (SVD)
- Nullspace projection (kernels in least squares)
- Ridge regression (Tikhonov regularization)
- Weighted least squares
- Iteratively reweighted least squares (IRLS)
- Recursive least squares
- Logistic regression

Tensor regression (Oct 7)

- Tensor factorization
- Tensor-variate regression

Nonlinear regression (Oct 14)

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

Hidden Markov model (HMM) (Nov 11)

Regression

Nonlinear regression:

Approximation function
(predictor function)
(forecasting function)

$$y = f(x)$$

Output
(response variable)
(outcome variable)

Input
(explanatory variable)
(feature variable)

$$\in \mathbb{R}$$

Linear regression:

$$y = Ax$$

Tensor regression:

**Similar to $y = Ax$,
but with x a tensor
instead of a vector**

Nonlinear regression on x reformulated as linear regression on f :

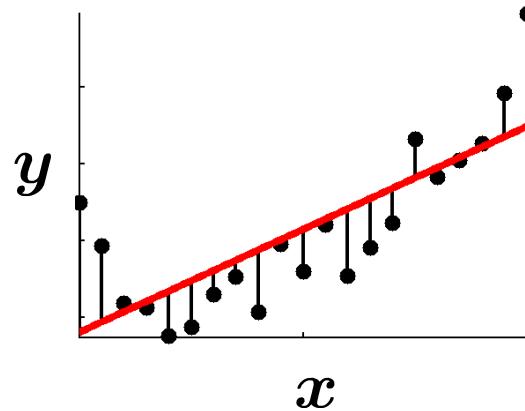
$$y = A \begin{bmatrix} f_1(x) \\ f_2(x) \\ \vdots \end{bmatrix}$$

LEAST SQUARES

circa 1795

Least squares: a ubiquitous tool

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



Fast and robust
implementation?

Recursive computation?

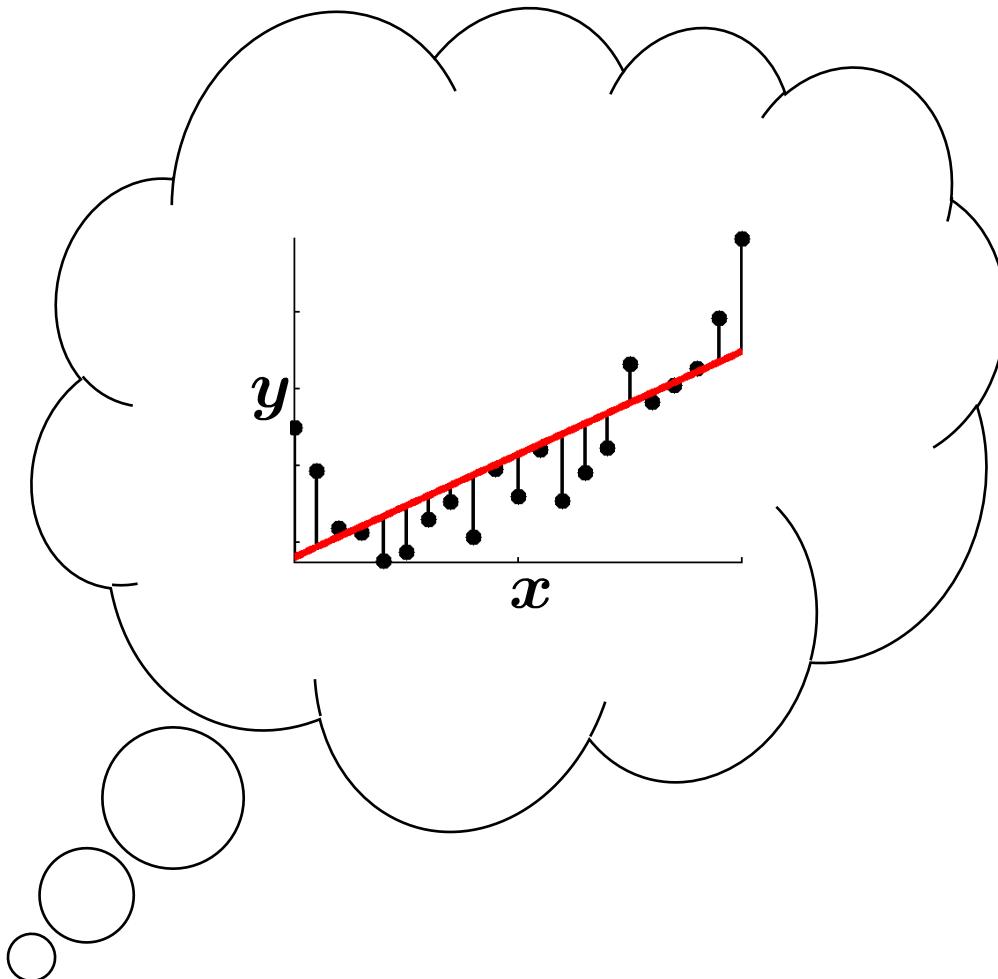
Solution with a secondary objective?

Weighted least squares?

Regularized least squares?

L1-norm instead of L2-norm?

Youtube search for “deep learning regression”



YouTube CH

deep learning regression

Home Trending History

BEST OF YOUTUBE

- Music
- Sports
- Gaming
- Movies
- News
- Live
- 360° Video
- Browse channels

FILTER

Applications of Deep Neural Networks Regression

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 ...

HOW TO DO LINEAR REGRESSION

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 ...

LINEAR REGRESSION WITH GRADIENT DESCENT

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 ...

BEGINNER INTRO TO NEURAL NETWORKS 8

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

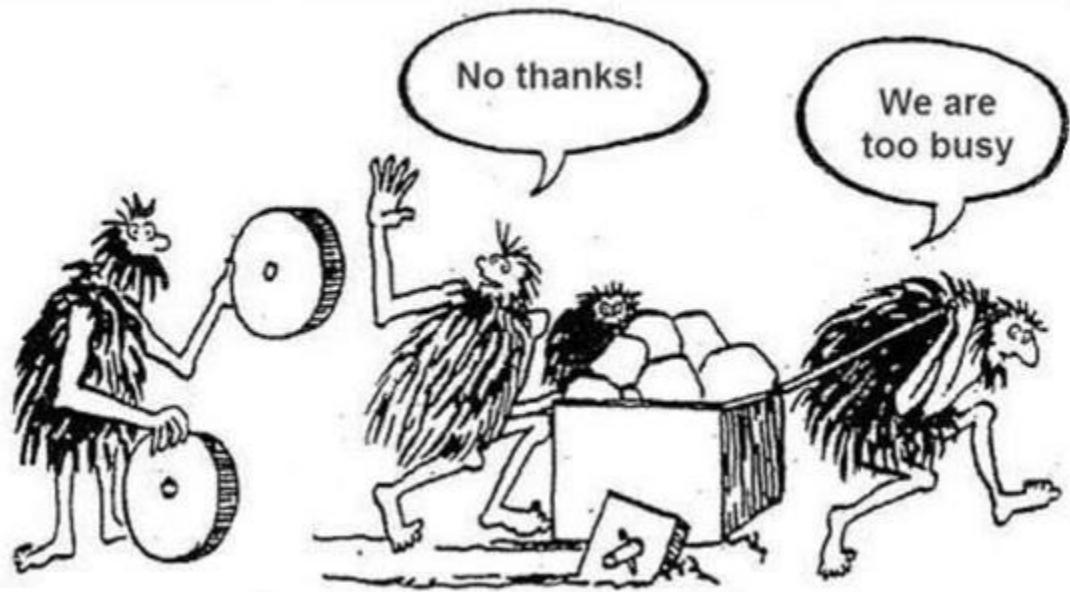
Technology for Noobs • 3.5K views • 1 year ago

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 ...

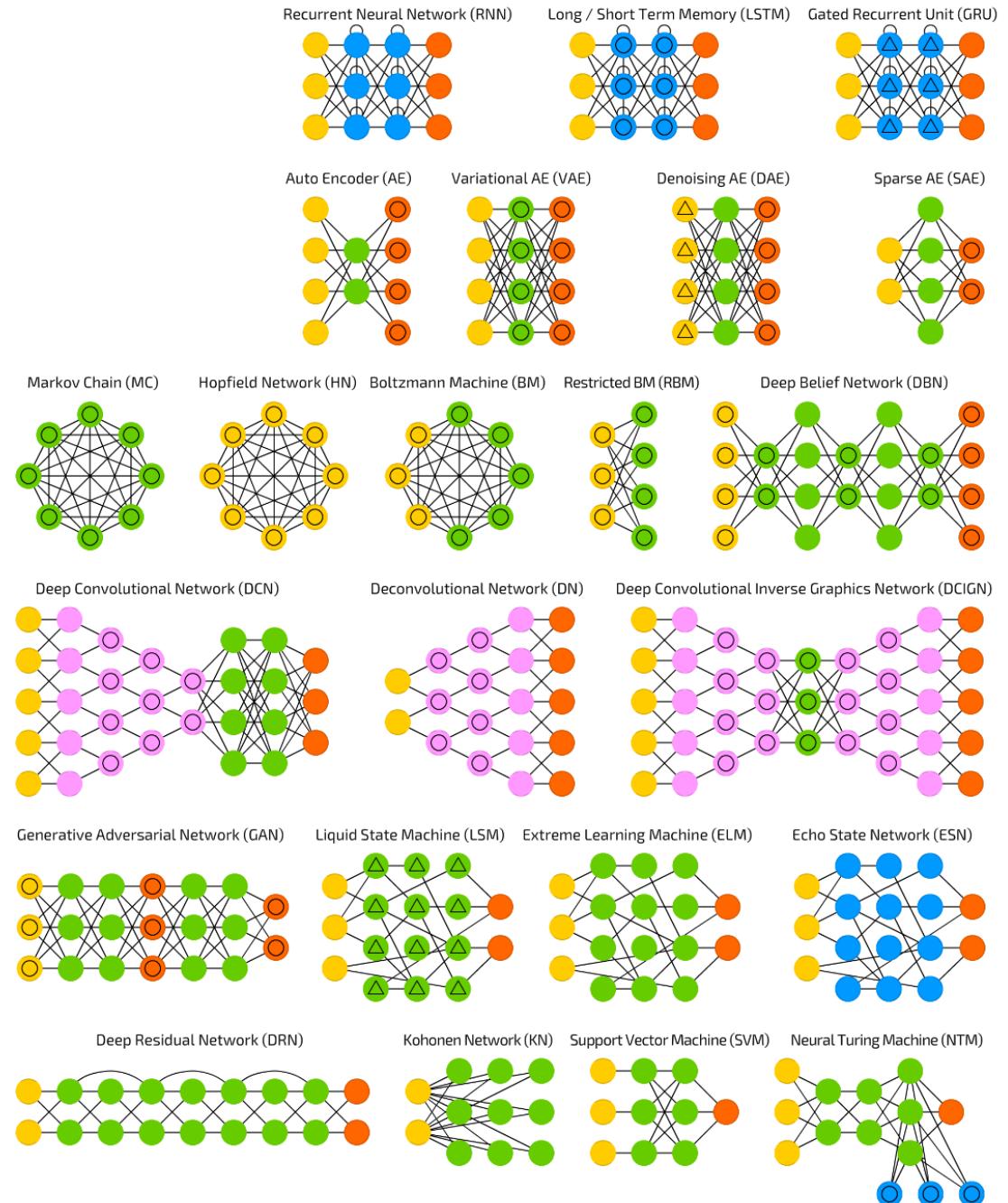


The **pinv-net!** 😊
single layer,
single node,
linear activation!



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

- 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



Multivariate linear regression (ordinary least squares)

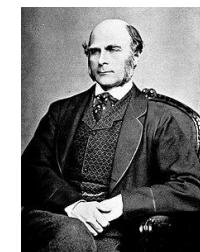
- **Least squares is everywhere:** from simple problems to large scale problems.
- It was the earliest form of regression, which was published by **Legendre** in 1805 and by **Gauss** in 1809. They both applied the method to the problem of determining the orbits of bodies around the sun from astronomical 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.
- **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 $\mathbf{y} = \mathbf{X}\mathbf{a}$.

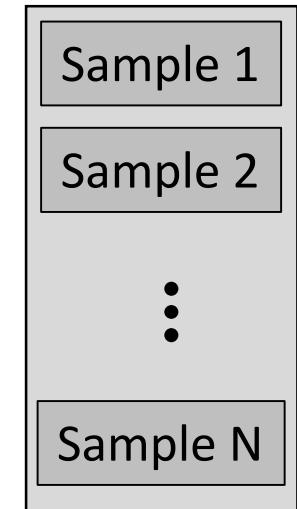
A solution can be found by minimizing the ℓ_2 norm

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

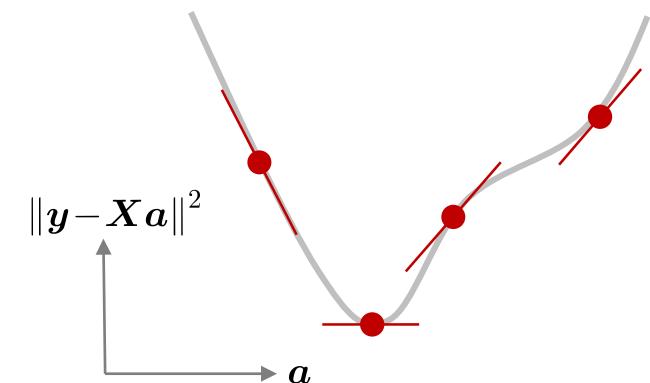
By differentiating with respect to \mathbf{a} and equating to zero

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

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



\mathbf{X}



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 $\mathbf{Y} = \mathbf{X}\mathbf{A}$.

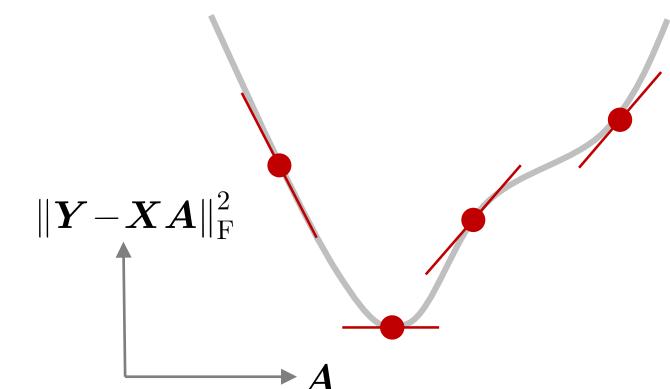
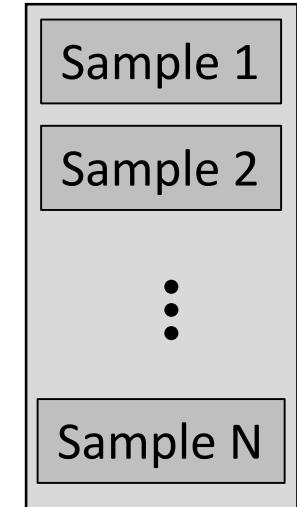
A solution can be found by minimizing the Frobenius norm

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

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

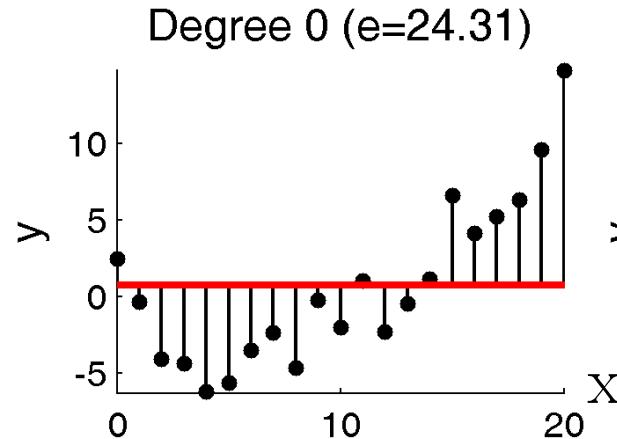
$$-2\mathbf{X}^\top\mathbf{Y} + 2\mathbf{X}^\top\mathbf{X}\mathbf{A} = \mathbf{0} \iff \hat{\mathbf{A}} = (\mathbf{X}^\top\mathbf{X})^{-1}\mathbf{X}^\top\mathbf{Y}$$

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

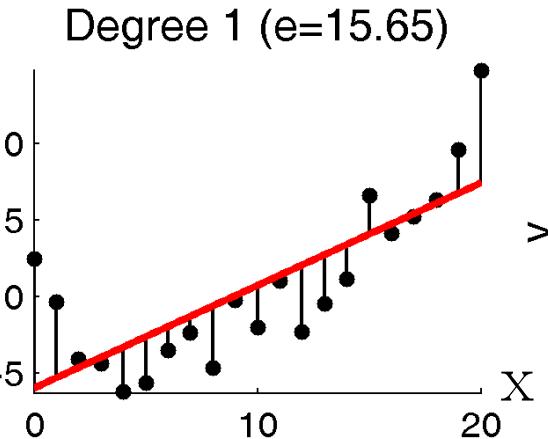


Polynomial fitting with least squares

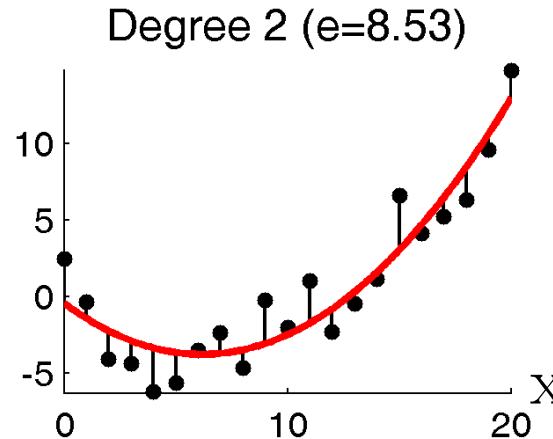
$$\hat{A} = X^\dagger Y$$



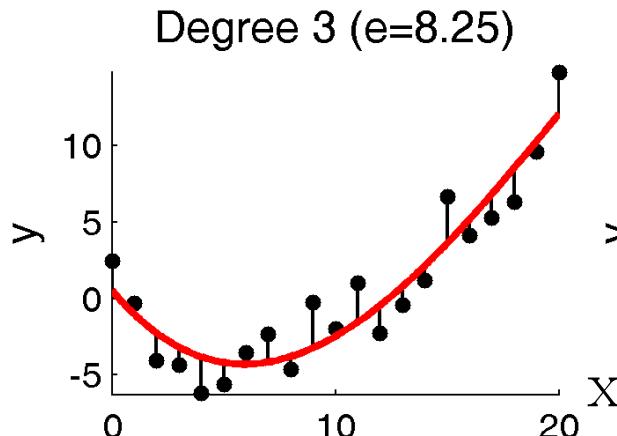
$$x = 1$$



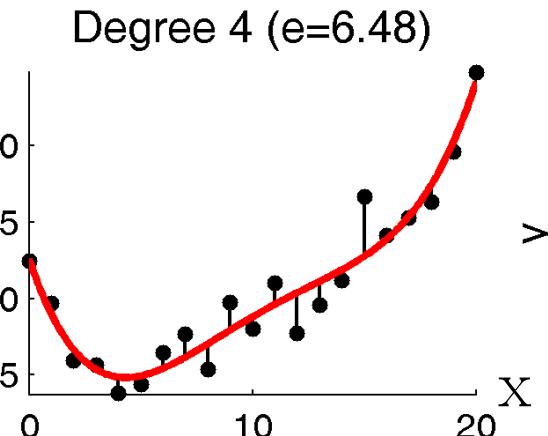
$$x = [1, x]$$



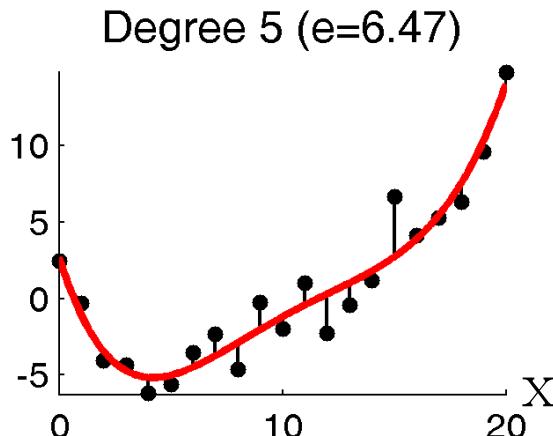
$$x = [1, x, x^2]$$



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



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



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

Least squares computed with SVD

Singular value decomposition (SVD)

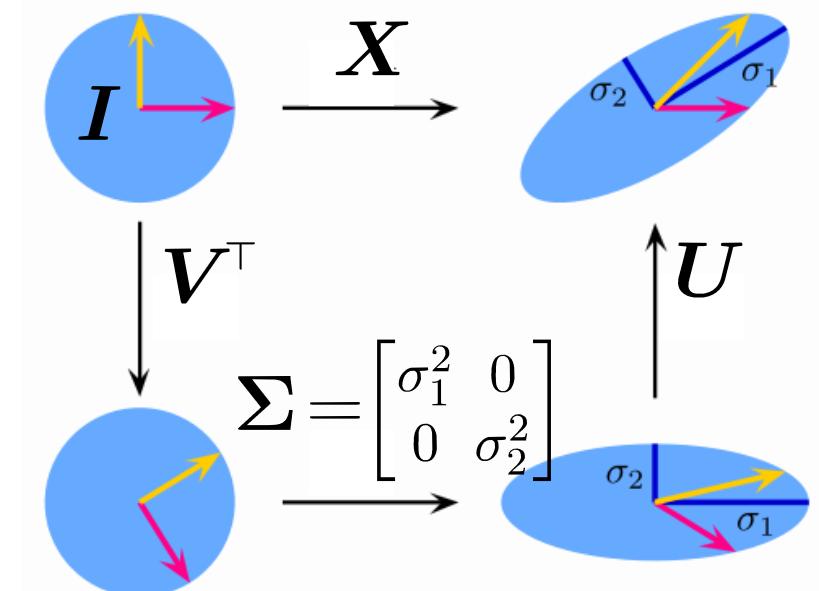
$$X \in \mathbb{R}^{N \times D^I} = U \in \mathbb{R}^{N \times N} \Sigma \in \mathbb{R}^{N \times D^I} V^\top \in \mathbb{R}^{D^I \times D^I}$$

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

Unitary matrix
(orthogonal)

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

Unitary matrix
(orthogonal)



Least squares computed with SVD

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

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

\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 $D^x \times D^x$ orthogonal matrices, and Σ is an $N \times D^x$ matrix with all its elements outside of the main diagonal equal to 0.

$$\Sigma = \begin{bmatrix} 4 & 0 & 0 & 0 & 0 \\ 0 & 2 & 0 & 0 & 0 \\ 0 & 0 & 0.1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \end{bmatrix}$$



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.

$$\Sigma^\dagger = \begin{bmatrix} 0.25 & 0 & 0 & 0 \\ 0 & 0.5 & 0 & 0 \\ 0 & 0 & 10 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix}$$

Nullspace projection (kernels in least squares)

Python notebook:
`demo_LS_polFit.ipynb`

Matlab code:
`demo_LS_polFit_nullspace01.m`

Nullspace projection

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

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

\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.

Nullspace projection computed with SVD

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

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

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

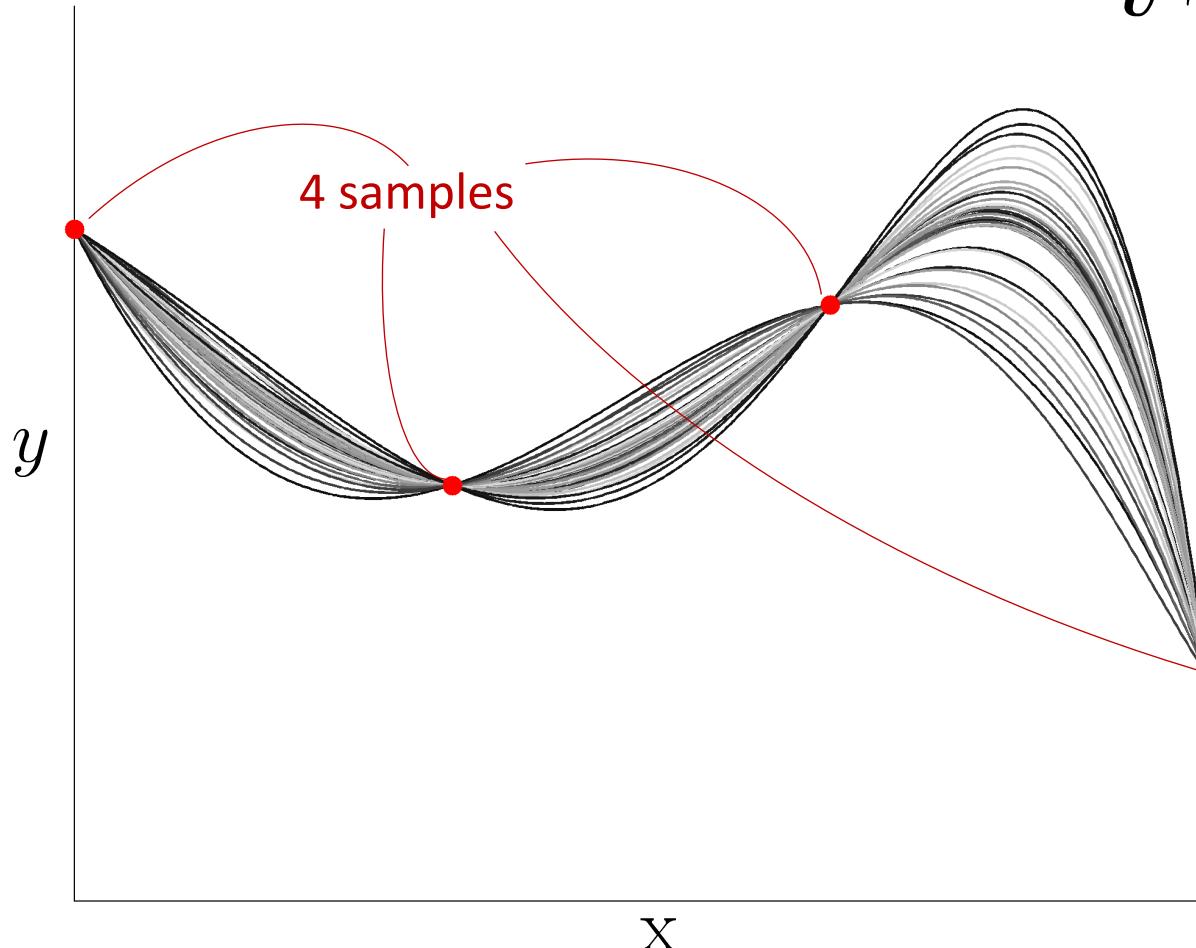
where $\tilde{\mathbf{U}}$ is a matrix formed by the columns of \mathbf{U} that span for the corresponding zero rows in Σ .

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

$\tilde{\mathbf{U}}$

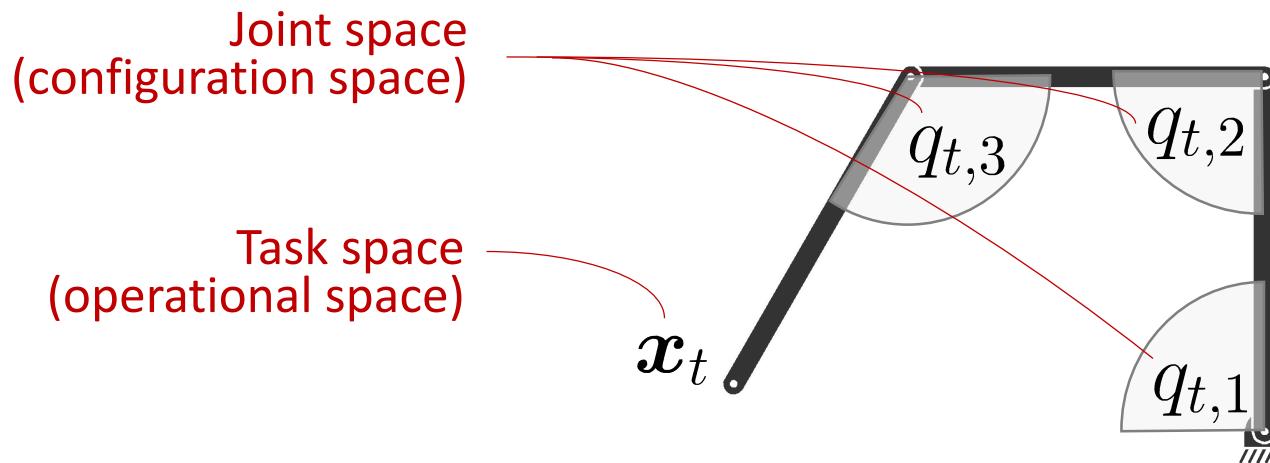
Example: Polynomial fitting

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



$$\begin{aligned}\hat{\mathbf{a}} &\in \mathbb{R}^7 \\ \mathbf{X} &\in \mathbb{R}^{4 \times 7} \\ \mathbf{y} &\in \mathbb{R}^4 \\ \mathbf{N} &\in \mathbb{R}^{7 \times 7} \\ \mathbf{v} &\in \mathbb{R}^7\end{aligned}$$

Example: Robot inverse kinematics



Forward kinematics is computed with

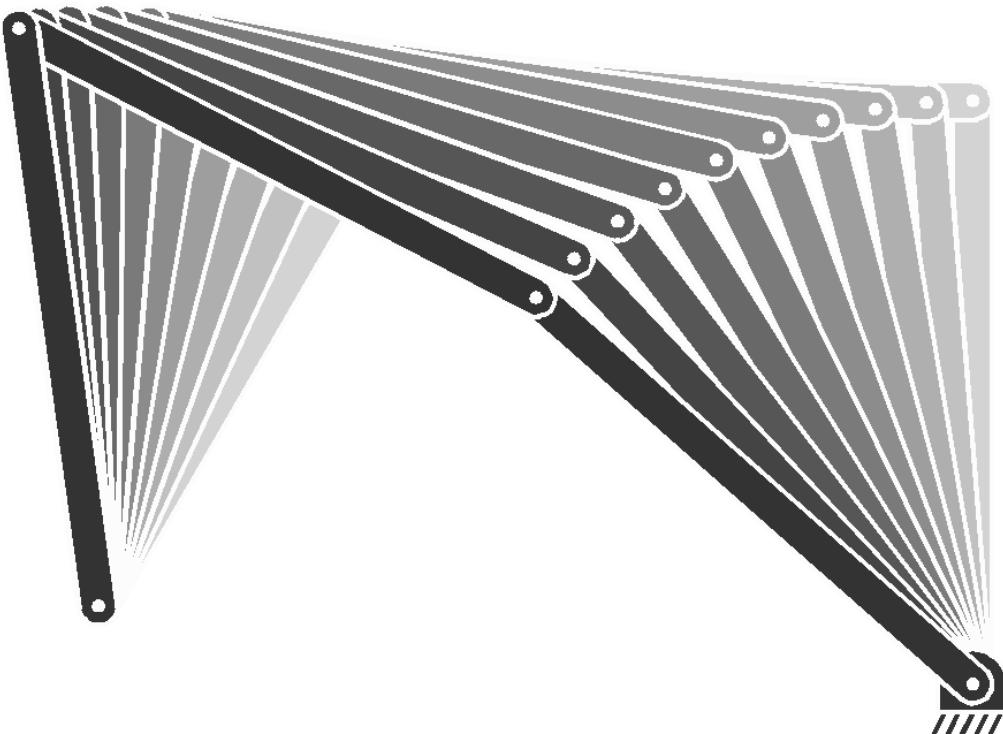
$$\boldsymbol{x}_t = \boldsymbol{f}(\boldsymbol{q}_t) \iff \dot{\boldsymbol{x}}_t = \frac{\partial \boldsymbol{x}_t}{\partial t} = \frac{\partial \boldsymbol{f}(\boldsymbol{q}_t)}{\partial \boldsymbol{q}_t} \frac{\partial \boldsymbol{q}_t}{\partial t} = \boldsymbol{J}(\boldsymbol{q}_t) \dot{\boldsymbol{q}}_t,$$

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

An inverse kinematics solution can be computed with

$$\hat{\dot{\boldsymbol{q}}}_t = \boldsymbol{J}^\dagger(\boldsymbol{q}_t) \dot{\boldsymbol{x}}_t + \boldsymbol{N}(\boldsymbol{q}_t) \boldsymbol{g}(\boldsymbol{q}_t), \quad \text{where } \boldsymbol{N}(\boldsymbol{q}_t) = \boldsymbol{I} - \boldsymbol{J}^\dagger(\boldsymbol{q}_t) \boldsymbol{J}(\boldsymbol{q}_t).$$

Example: Robot inverse kinematics



Primary constraint:
keeping the tip of
the robot still

$$\begin{aligned}\hat{\dot{\mathbf{q}}}_t &= \mathbf{J}^\dagger(\mathbf{q}_t) \dot{\mathbf{x}}_t + \mathbf{N}(\mathbf{q}_t) \mathbf{g}(\mathbf{q}_t) \\ &= \mathbf{J}^\dagger(\mathbf{q}_t) \begin{bmatrix} 0 \\ 0 \end{bmatrix} + \mathbf{N}(\mathbf{q}_t) \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}\end{aligned}$$

**Secondary
constraint:**
moving the
first joint

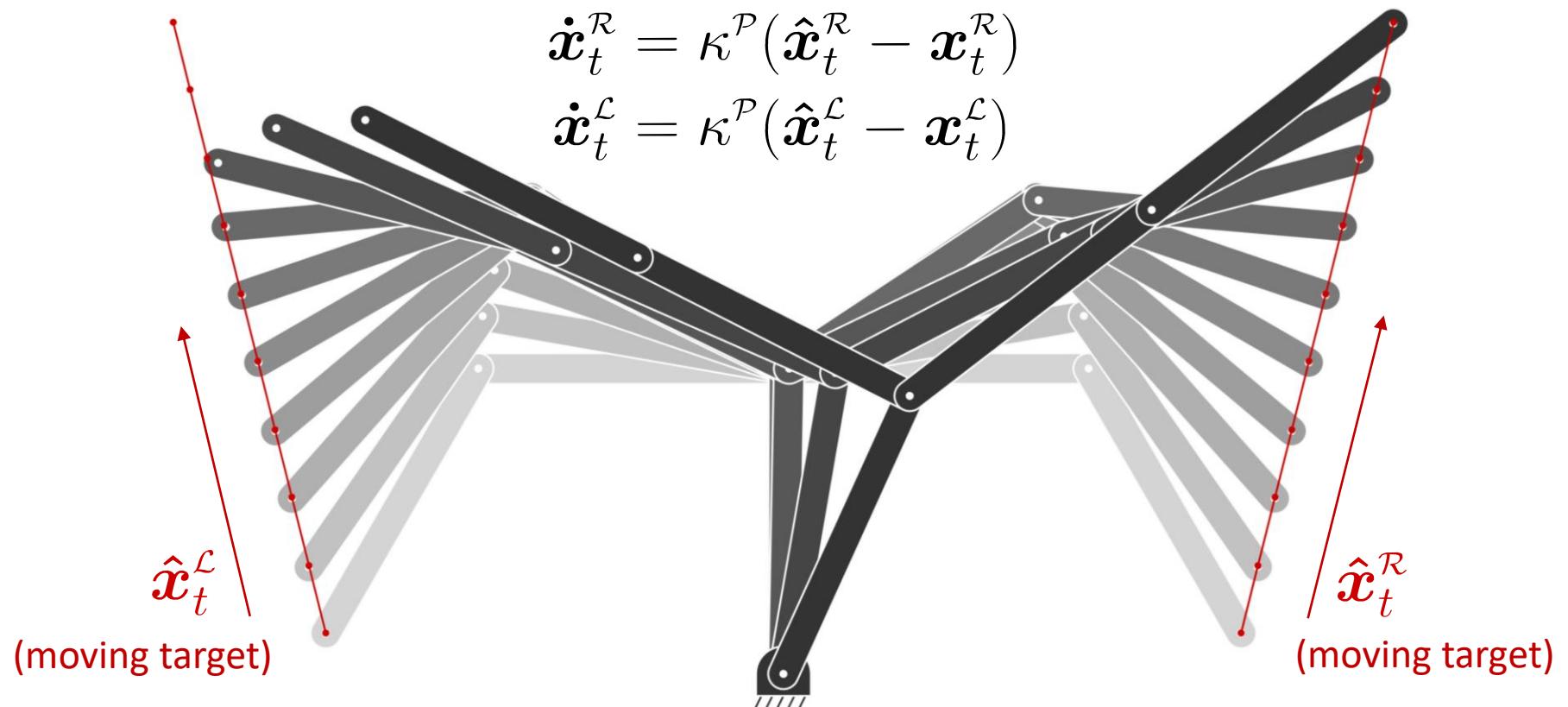
Example: Robot inverse kinematics

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

Tracking target with left hand, if possible

Tracking target with right hand

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



Ridge regression (robust regression, Tikhonov regularization, penalized least squares)

Python notebook:
`demo_LS_polFit.ipynb`

Matlab example:
`demo_LS_polFit02.m`

Ridge regression

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

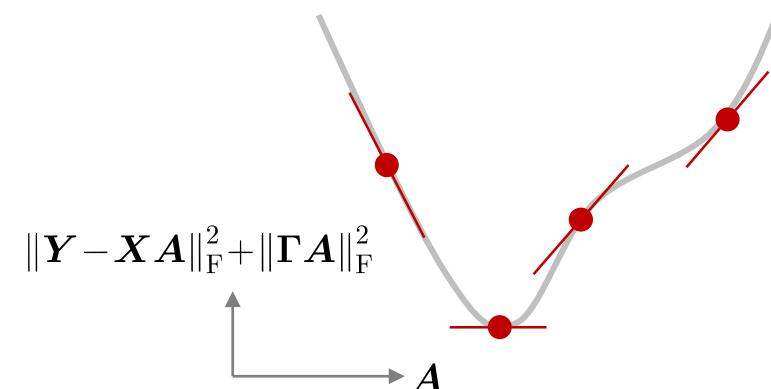
$$\begin{aligned}\hat{\mathbf{A}} &= \arg \min_{\mathbf{A}} \|\mathbf{Y} - \mathbf{X}\mathbf{A}\|_F^2 + \|\boldsymbol{\Gamma}\mathbf{A}\|_F^2 \\ &= \arg \min_{\mathbf{A}} \text{tr}\left((\mathbf{Y} - \mathbf{X}\mathbf{A})^\top(\mathbf{Y} - \mathbf{X}\mathbf{A})\right) + \text{tr}\left((\boldsymbol{\Gamma}\mathbf{A})^\top\boldsymbol{\Gamma}\mathbf{A}\right)\end{a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

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 $0 < \lambda \ll 1$ (i.e., giving preference to solutions with smaller norms), the process is known as **ℓ_2 regularization**.

Ridge regression

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 computed 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

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

and has zeros elsewhere.

$$\Sigma = \begin{bmatrix} 4 & 0 & 0 & 0 & 0 \\ 0 & 2 & 0 & 0 & 0 \\ 0 & 0 & 0.1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \end{bmatrix}$$

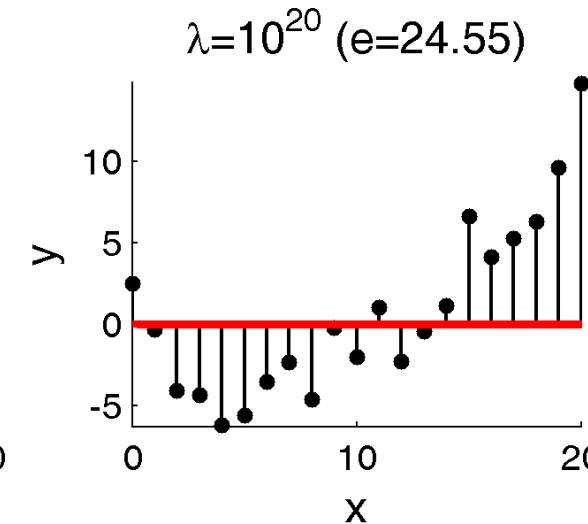
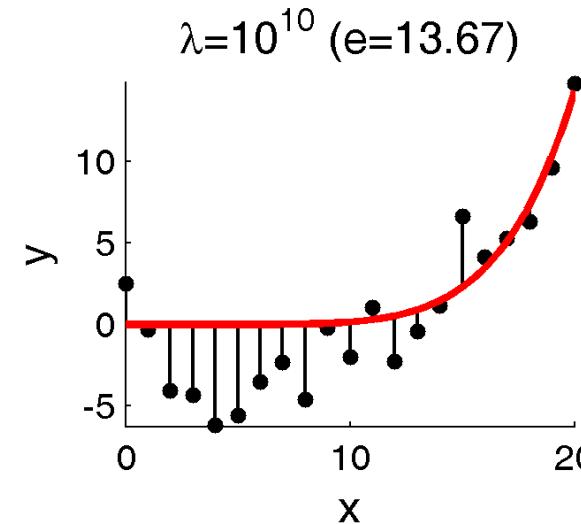
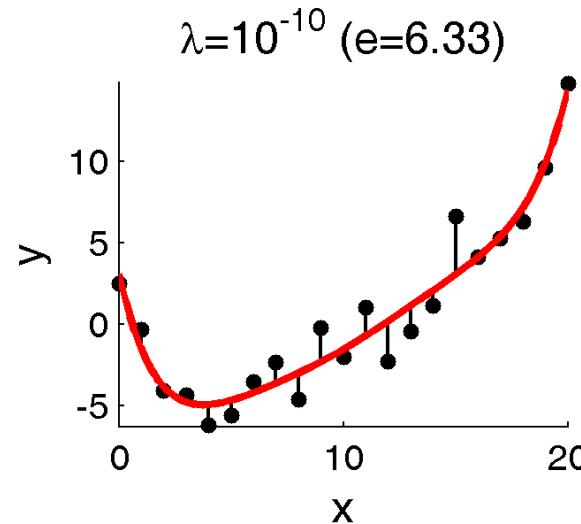


$$\tilde{\Sigma} = \begin{bmatrix} 0.2498 & 0 & 0 & 0 \\ 0 & 0.4988 & 0 & 0 \\ 0 & 0 & 5 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix}$$

(for $\lambda = 0.01$)

Example: Polynomial fitting

$D^x = 7$ (polynomial of degree 6)



Weighted least squares (Generalized least squares)

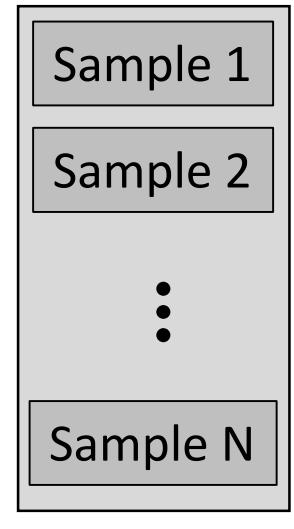
Python notebook:
`demo_LS_weighted.ipynb`

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

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

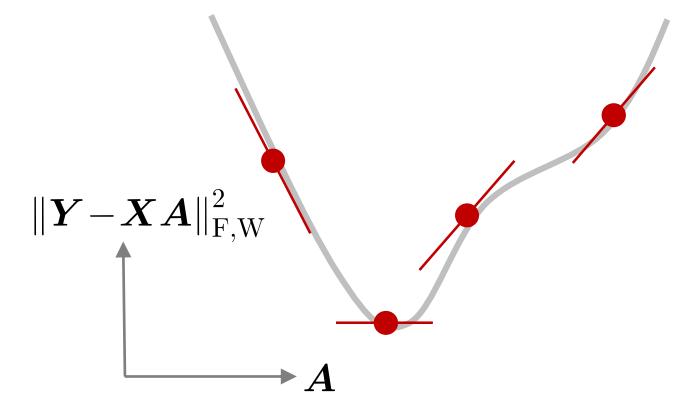


\mathbf{X}

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

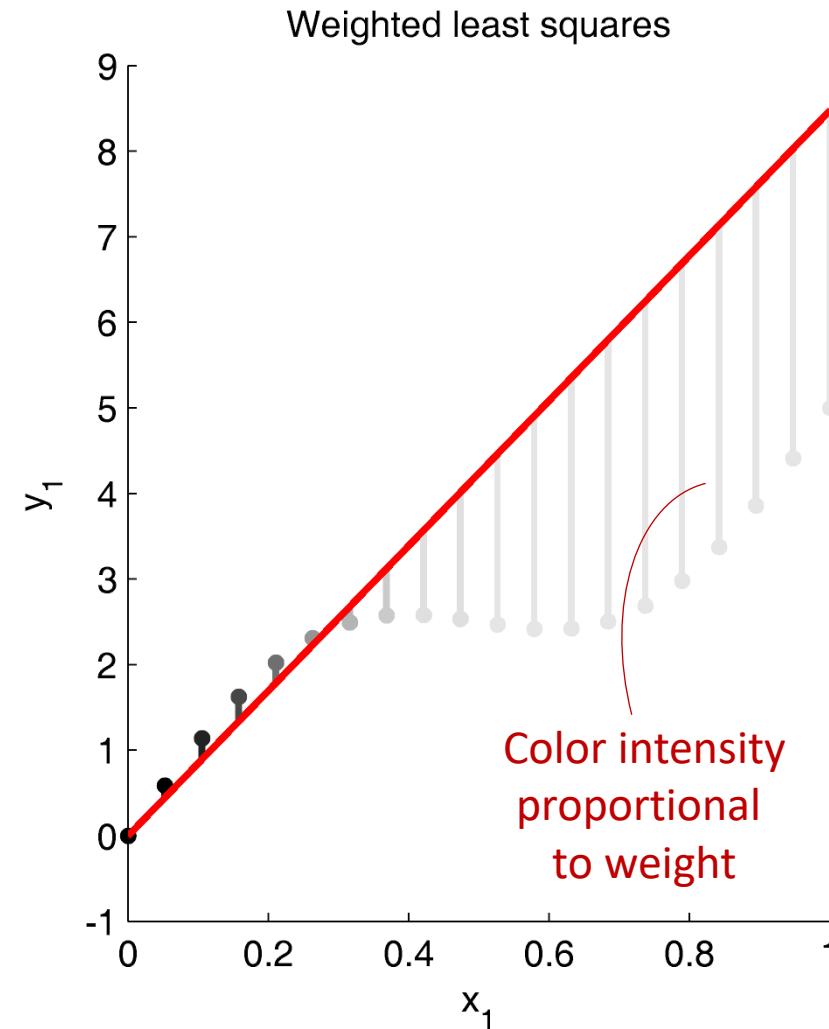
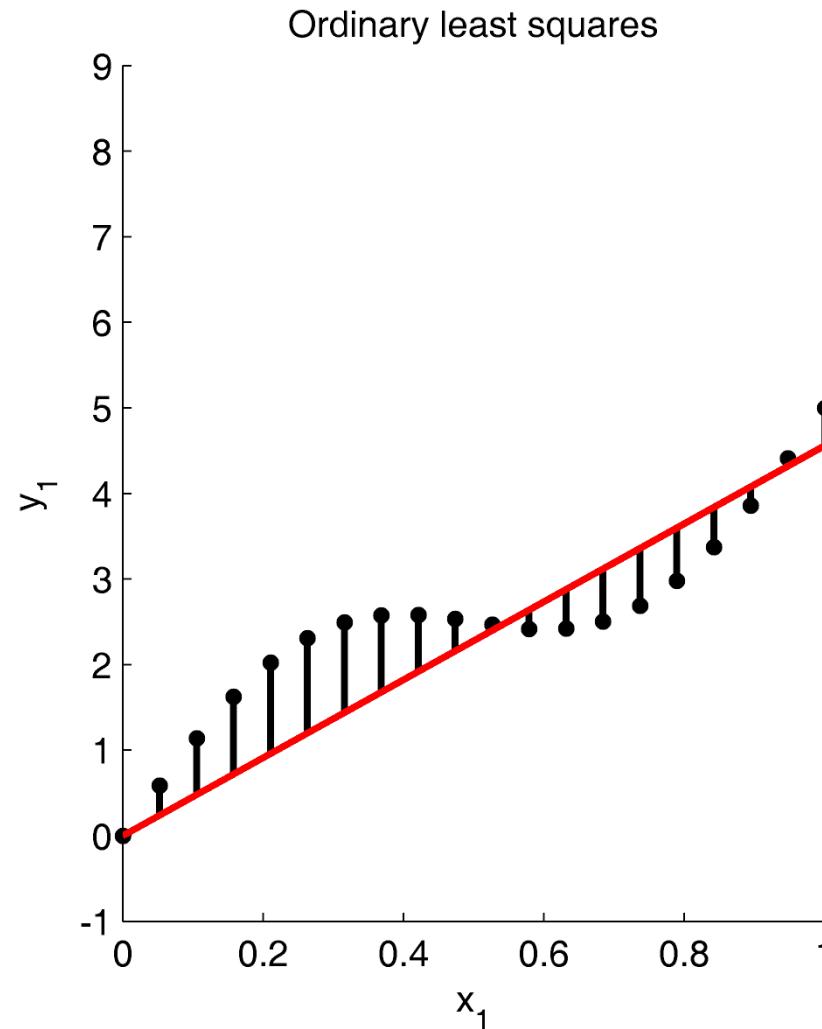
$$-2\mathbf{X}^\top \mathbf{W} \mathbf{Y} + 2\mathbf{X}^\top \mathbf{W} \mathbf{X} \mathbf{A} = \mathbf{0} \iff \hat{\mathbf{A}} = (\mathbf{X}^\top \mathbf{W} \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{W} \mathbf{Y}$$

\mathbf{X}_W^\dagger



Weighted least squares

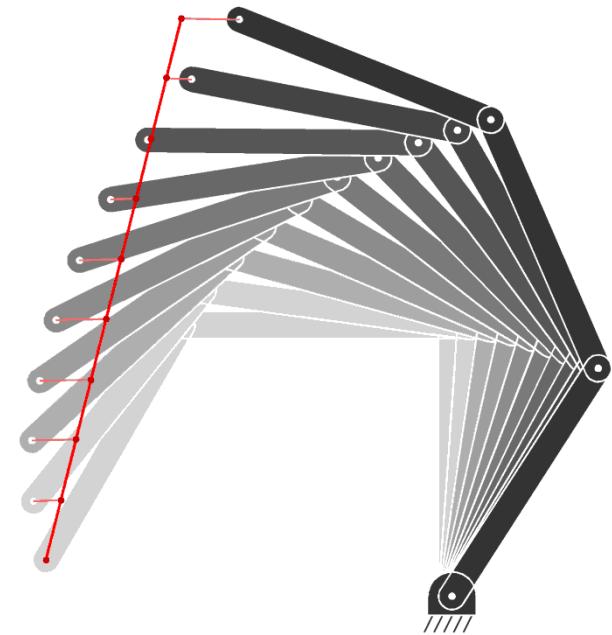
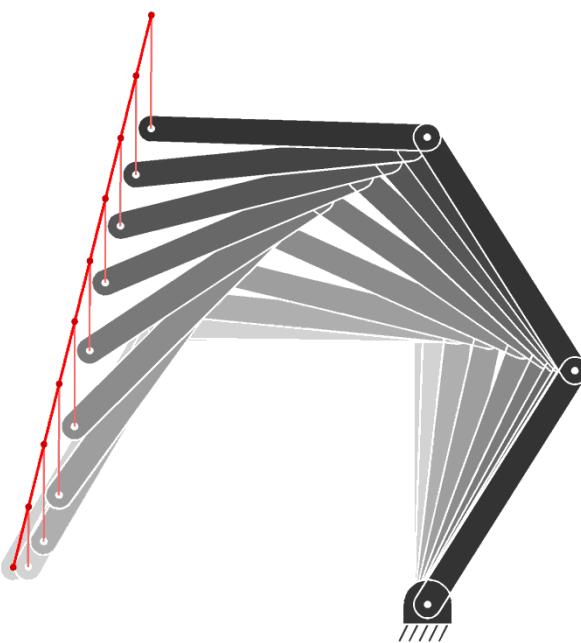
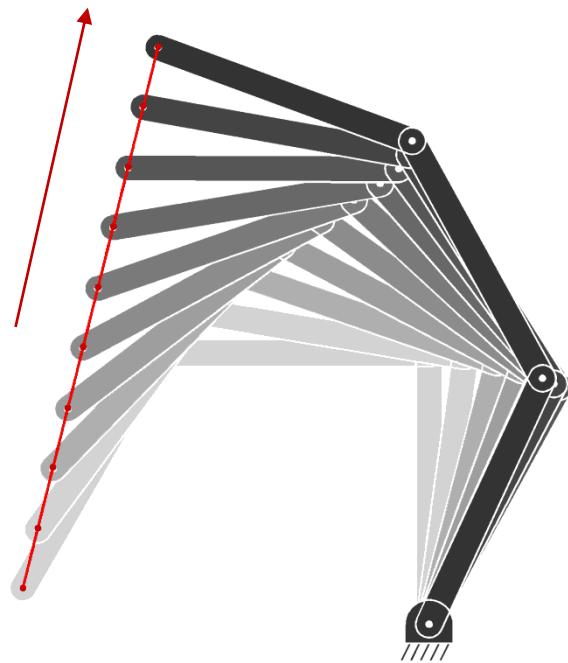
$$\hat{A} = (X^\top W X)^{-1} X^\top W Y$$



Example: Robot inverse kinematics (weights in task space)

$$\hat{\dot{q}}_t = (\mathbf{J}^\top \mathbf{W}^x \mathbf{J} + \lambda \mathbf{I}_q)^{-1} \mathbf{J}^\top \mathbf{W}^x \dot{x}_t$$

(moving target)



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

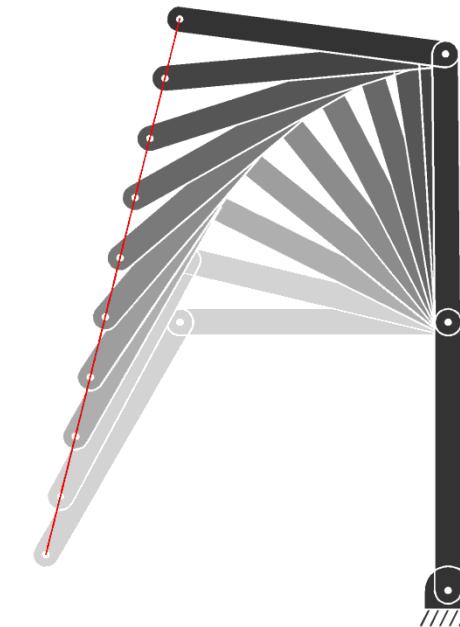
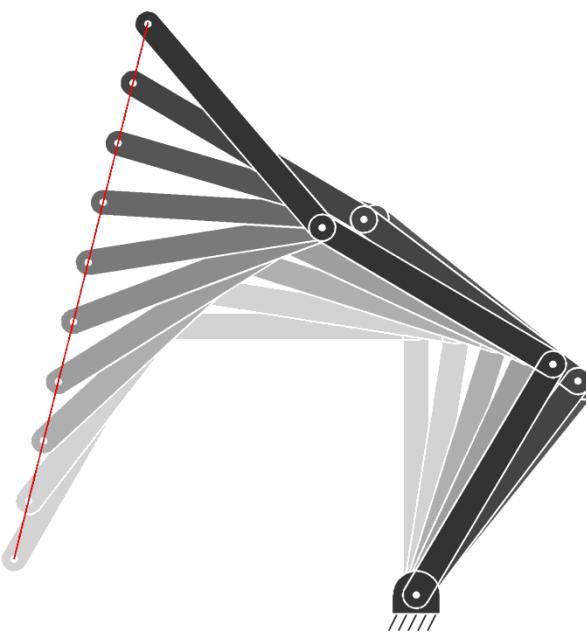
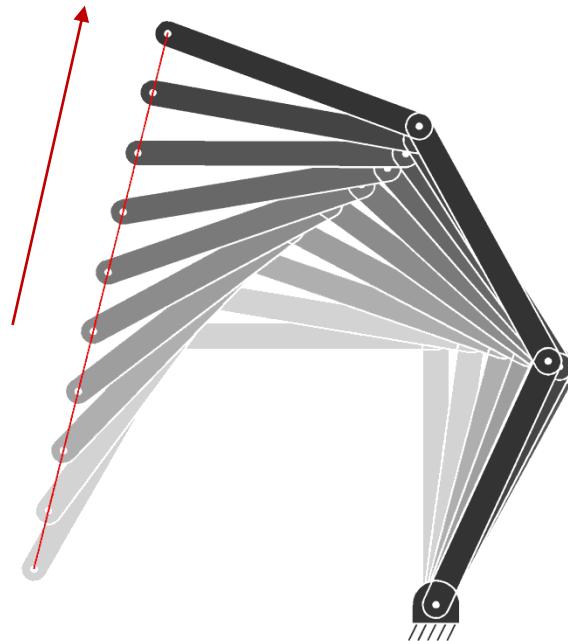
$$\mathbf{W}^x = \begin{bmatrix} 1 & 0 \\ 0 & .01 \end{bmatrix}$$

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

Example: Robot inverse kinematics (weights in joint space)

$$\hat{\dot{q}}_t = \mathbf{W}^Q \mathbf{J}^\top (\mathbf{J} \mathbf{W}^Q \mathbf{J}^\top + \lambda \mathbf{I}_x)^{-1} \dot{x}_t$$

(moving target)



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

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

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

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 (IRLS)** generalizes least squares by raising the error term to a power that is less than 2:
 - the optimization problem can no longer be called “least squares”
- The strategy of IRLS is that an error $|\mathbf{e}|^p$ can be rewritten as $|\mathbf{e}|^p = |\mathbf{e}|^{p-2} \mathbf{e}^2$.
- $|\mathbf{e}|^{p-2}$ can be interpreted as a weight, which is used to minimize \mathbf{e}^2 with **weighted least squares**.
 - we solve a least squares problem at each iteration of the algorithm
- $p=1$ corresponds to **least absolute deviation regression**.

Iteratively reweighted least squares (IRLS)

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

transformed as
weight \mathbf{W}

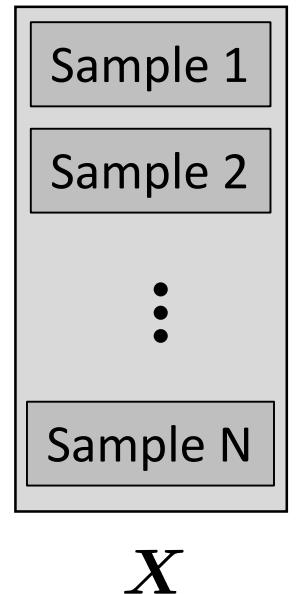
For an ℓ_p norm cost function defined by

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

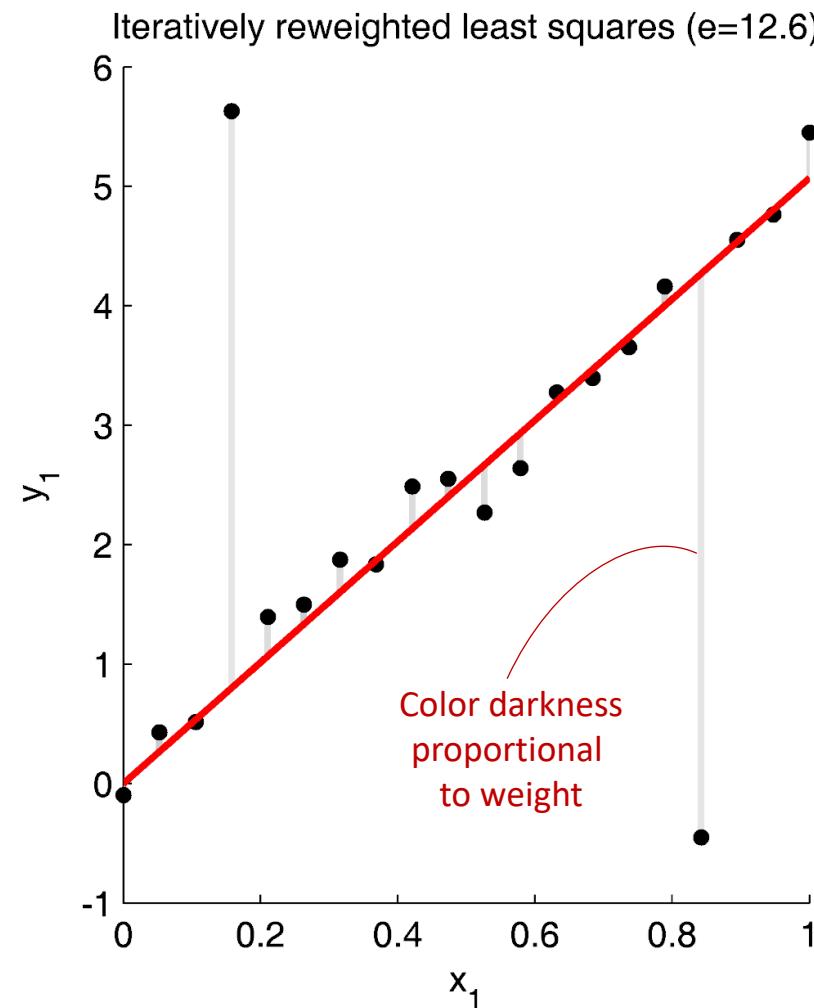
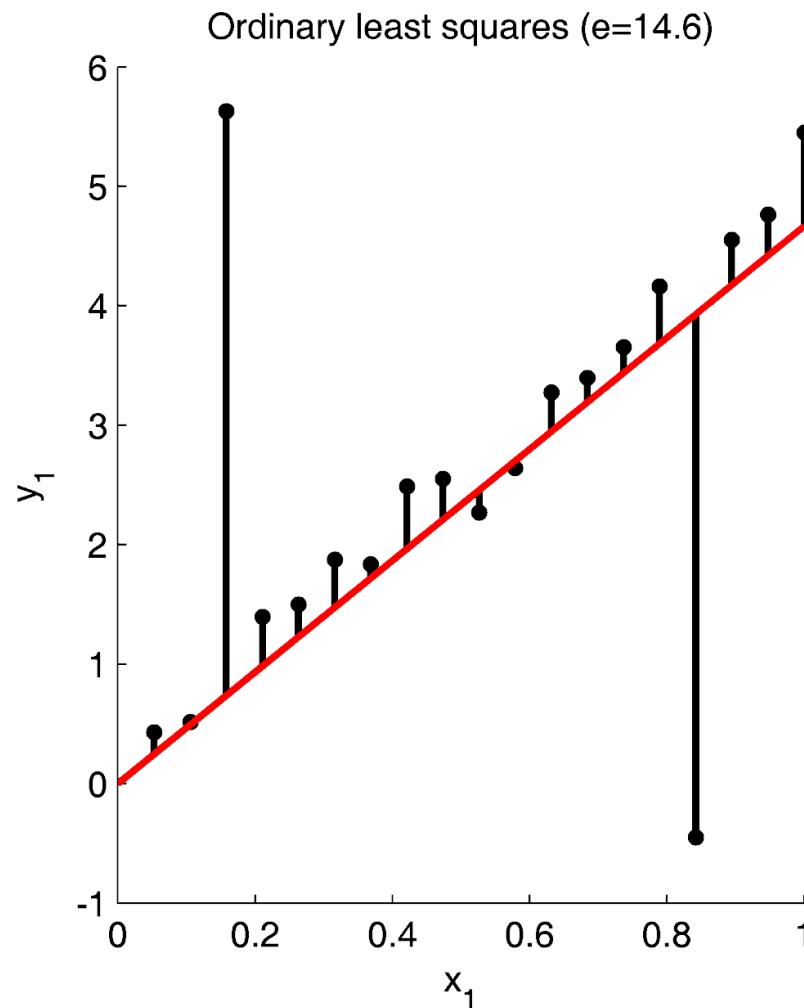
$\hat{\mathbf{A}}$ is estimated by starting from $\mathbf{W} = \mathbf{I}$ and iteratively computing

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

$$\mathbf{W}_{n,n} \leftarrow |\mathbf{Y}_n - \mathbf{X}_n \mathbf{A}|^{p-2} \quad \forall n \in \{1, \dots, N\}$$



Iteratively reweighted least squares (IRLS)



→ regression that can be more 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.

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

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

↓

$$\mathbf{X}_{\text{new}} \in \mathbb{R}^{N_{\text{new}} \times D^T}$$

Recursive least squares

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

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

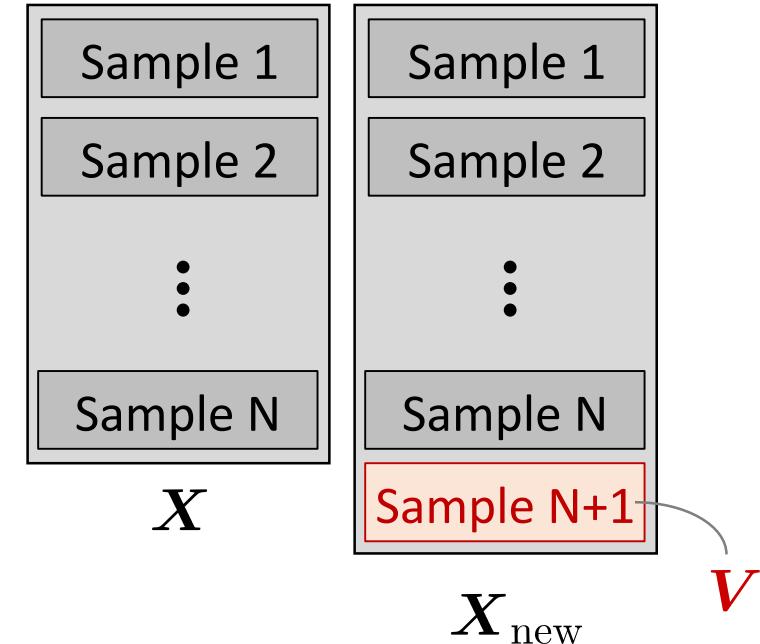
whose inverse can be computed with

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

which is exploited to efficiently compute the least squares update as

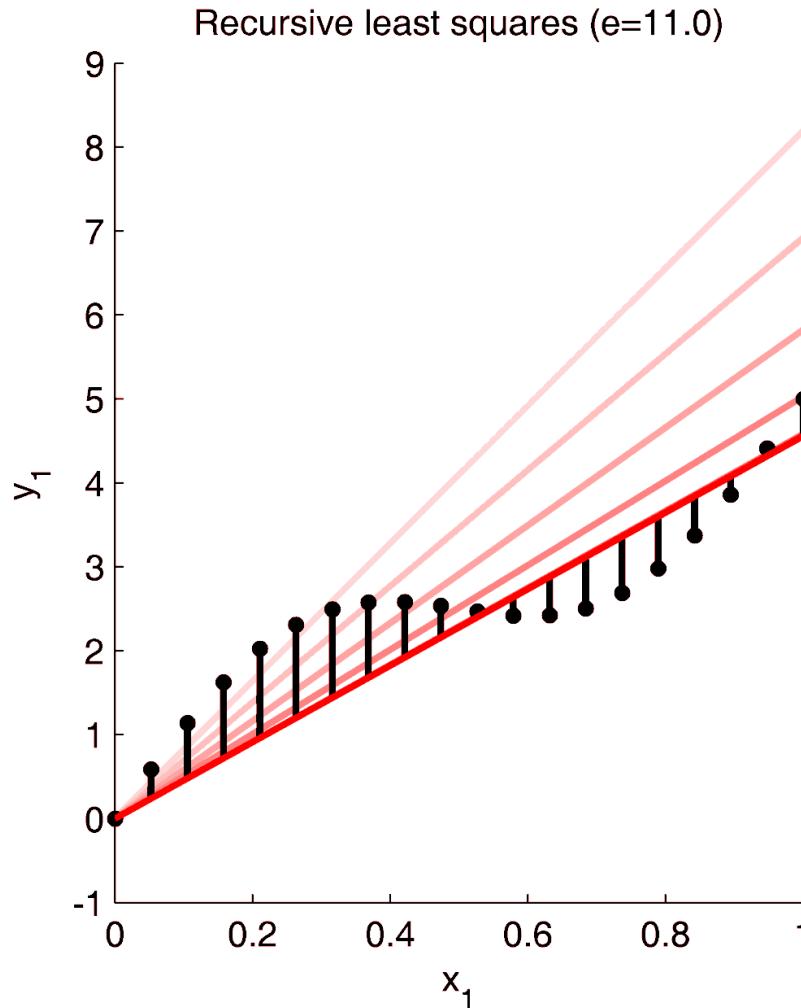
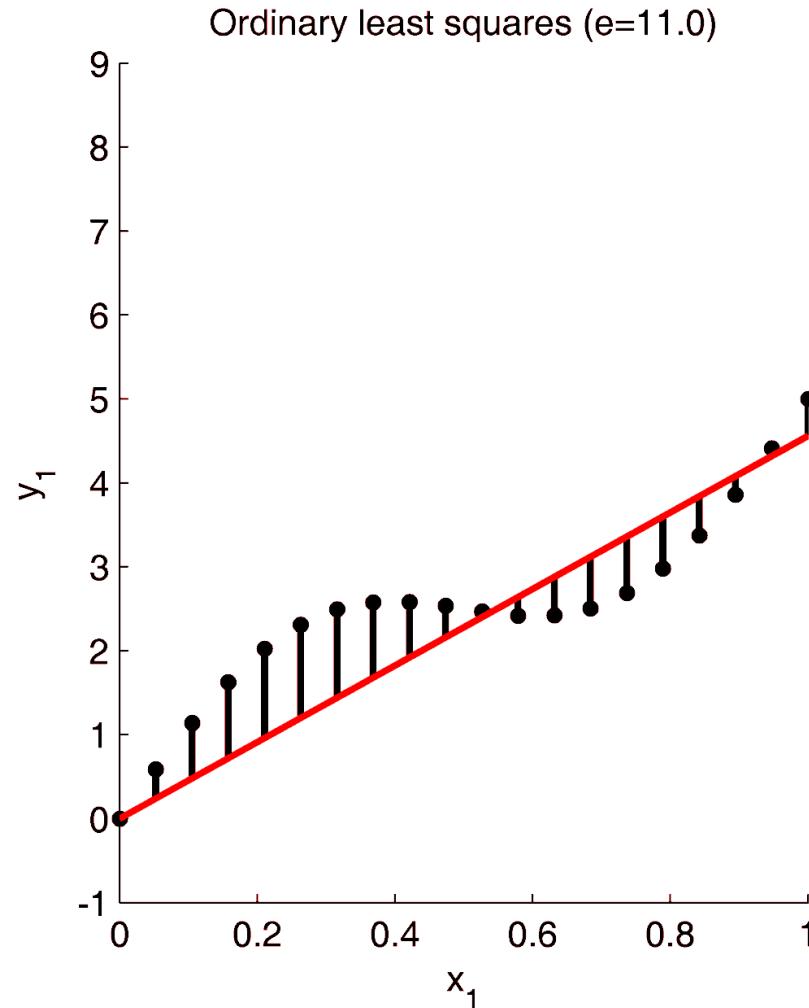
$$\hat{\mathbf{A}}_{\text{new}} = \hat{\mathbf{A}} + \mathbf{K} (\mathbf{C} - \mathbf{V} \hat{\mathbf{A}})$$

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



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

Recursive least squares



→ the least squares estimate is the same in the two cases

Linear regression:

Example of application

Linear quadratic tracking (LQT)

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

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)

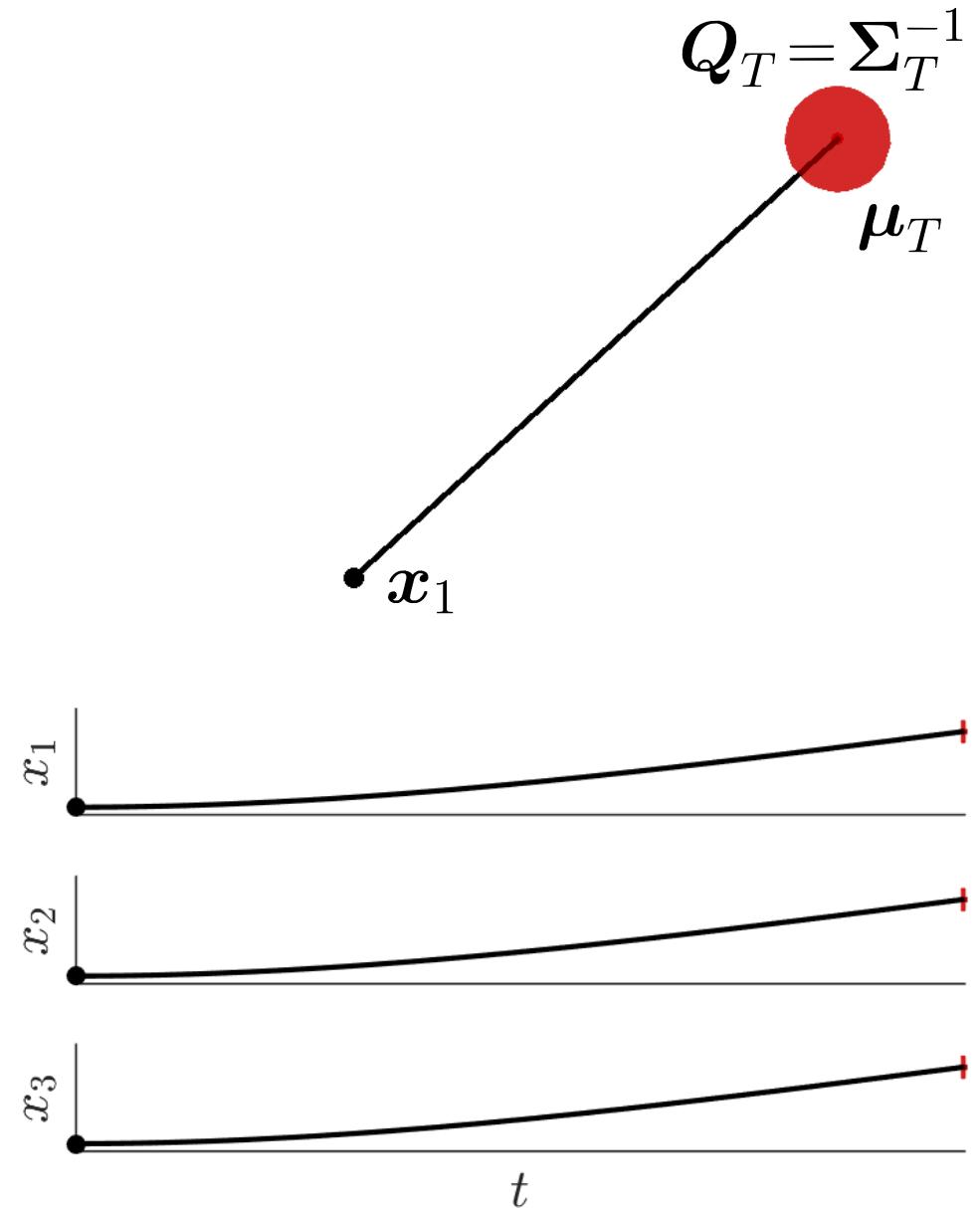
$\boldsymbol{\mu}_t$ desired state

\boldsymbol{u}_t control command (acceleration)

Q_t precision matrix

R_t control weight matrix

$$\boldsymbol{u} = \begin{bmatrix} \boldsymbol{u}_1 \\ \boldsymbol{u}_2 \\ \vdots \\ \boldsymbol{u}_T \end{bmatrix}$$



Linear quadratic tracking (LQT)

Track path! Use low control commands!

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

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)



Linear quadratic tracking (LQT)

$$c = \sum_{t=1}^T \left((\mu_t - x_t)^\top Q_t (\mu_t - x_t) + u_t^\top R_t u_t \right)$$

$$= (\mu - x)^\top Q (\mu - x) + u^\top R u$$

$$Q = \begin{bmatrix} Q_1 & 0 & \cdots & 0 \\ 0 & Q_2 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & Q_T \end{bmatrix}$$

$$\mu = \begin{bmatrix} \mu_1 \\ \mu_2 \\ \vdots \\ \mu_T \end{bmatrix}$$

$$x = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_T \end{bmatrix}$$

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

$$R = \begin{bmatrix} R_1 & 0 & \cdots & 0 \\ 0 & R_2 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & R_T \end{bmatrix}$$



Linear quadratic tracking (LQT)

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

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

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

⋮

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

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

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



Linear quadratic tracking (LQT)

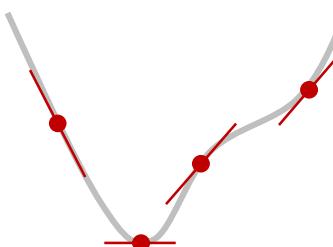
The constraint can then be inserted in the cost function:

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

Solving for u is similar to a **weighted ridge regression** problem, and results in the analytic solution:

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

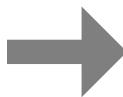
$$\hat{u} = \begin{bmatrix} \hat{u}_1 \\ \hat{u}_2 \\ \vdots \\ \hat{u}_T \end{bmatrix}$$



Linear quadratic tracking (LQT)

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

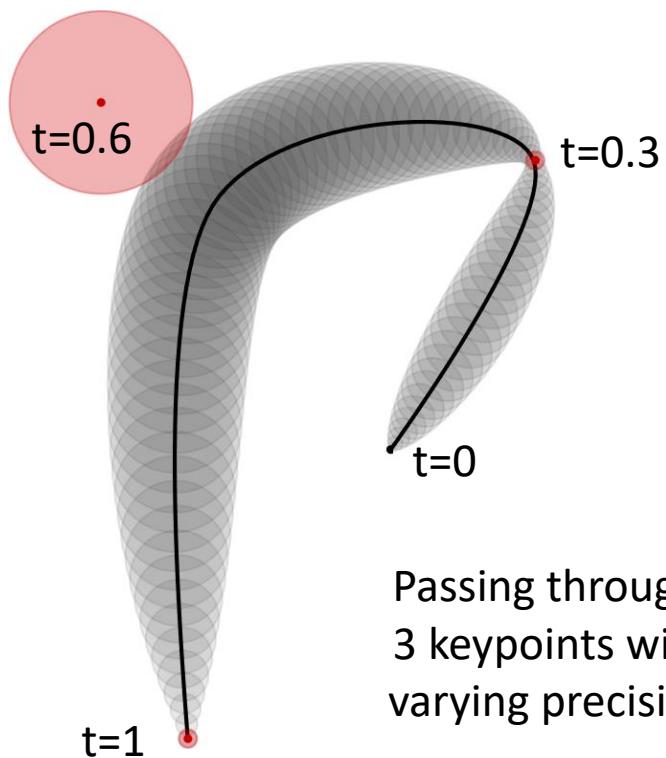
$$\hat{\Sigma}^u = (S^{u\top} Q S^u + R)^{-1}$$



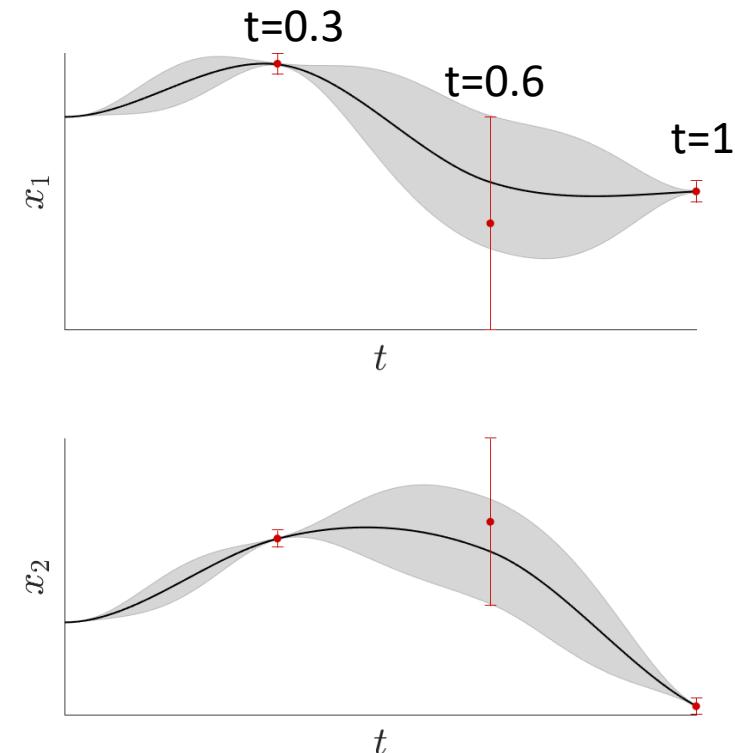
$$\hat{x} = S^x x_1 + S^u \hat{u}$$

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

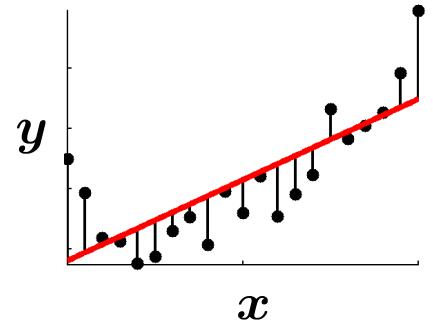
The error on the regression estimate can be used to form a distribution



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



Linear quadratic tracking (LQT) - Example



For $t \leq \frac{T}{4}$ (left hand holding the ball), we have

$$\begin{array}{c}
 \text{Hand} \\
 \text{Hand} \\
 \text{Ball}
 \end{array}
 \begin{bmatrix} \dot{\mathbf{x}}_1,t \\ \ddot{\mathbf{x}}_1,t \\ \dot{\mathbf{f}}_1,t \\ \dot{\mathbf{x}}_2,t \\ \ddot{\mathbf{x}}_2,t \\ \dot{\mathbf{f}}_2,t \\ \dot{\mathbf{x}}_3,t \\ \ddot{\mathbf{x}}_3,t \\ \dot{\mathbf{f}}_3,t \end{bmatrix}_{\dot{\mathbf{x}}_t} =
 \underbrace{\begin{bmatrix} 0 & \mathbf{I} & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & \mathbf{I} & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & \mathbf{I} & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & \mathbf{I} & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & \mathbf{I} & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & \mathbf{I} \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \end{bmatrix}}_{\mathbf{A}_t^c}
 \underbrace{\begin{bmatrix} \mathbf{x}_1,t \\ \dot{\mathbf{x}}_1,t \\ \mathbf{f}_1,t \\ \mathbf{x}_2,t \\ \dot{\mathbf{x}}_2,t \\ \mathbf{f}_2,t \\ \mathbf{x}_3,t \\ \dot{\mathbf{x}}_3,t \\ \mathbf{f}_3,t \end{bmatrix}}_{\mathbf{x}_t} +
 \underbrace{\begin{bmatrix} 0 & 0 \\ \mathbf{I} & 0 \\ 0 & 0 \\ 0 & 0 \\ 0 & \mathbf{I} \\ 0 & 0 \\ 0 & 0 \\ 0 & 0 \\ 0 & 0 \end{bmatrix}}_{\mathbf{B}_t^c}
 \underbrace{\begin{bmatrix} \mathbf{u}_{1,t} \\ \mathbf{u}_{2,t} \end{bmatrix}}_{\mathbf{u}_t}$$

$\mathbf{f}_{i,t} = m_i \mathbf{g}$ with $\mathbf{g} = \begin{bmatrix} 0 \\ -9.81 \end{bmatrix}$

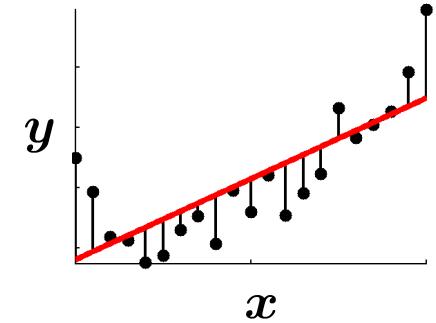
At $t = \frac{T}{2}$ (right hand hitting the ball), we have

$$\begin{array}{c}
 \text{Hand} \\
 \text{Hand} \\
 \text{Ball}
 \end{array}
 \begin{bmatrix} \dot{\mathbf{x}}_1,t \\ \ddot{\mathbf{x}}_1,t \\ \dot{\mathbf{f}}_1,t \\ \dot{\mathbf{x}}_2,t \\ \ddot{\mathbf{x}}_2,t \\ \dot{\mathbf{f}}_2,t \\ \dot{\mathbf{x}}_3,t \\ \ddot{\mathbf{x}}_3,t \\ \dot{\mathbf{f}}_3,t \end{bmatrix}_{\dot{\mathbf{x}}_t} =
 \underbrace{\begin{bmatrix} 0 & \mathbf{I} & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & \mathbf{I} & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & \mathbf{I} & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & \mathbf{I} & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \end{bmatrix}}_{\mathbf{A}_t^c}
 \underbrace{\begin{bmatrix} \mathbf{x}_1,t \\ \dot{\mathbf{x}}_1,t \\ \mathbf{f}_1,t \\ \mathbf{x}_2,t \\ \dot{\mathbf{x}}_2,t \\ \mathbf{f}_2,t \\ \mathbf{x}_3,t \\ \dot{\mathbf{x}}_3,t \\ \mathbf{f}_3,t \end{bmatrix}}_{\mathbf{x}_t} +
 \underbrace{\begin{bmatrix} 0 & 0 \\ \mathbf{I} & 0 \\ 0 & 0 \\ 0 & 0 \\ 0 & \mathbf{I} \\ 0 & 0 \\ 0 & 0 \\ 0 & 0 \\ 0 & 0 \end{bmatrix}}_{\mathbf{B}_t^c}
 \underbrace{\begin{bmatrix} \mathbf{u}_{1,t} \\ \mathbf{u}_{2,t} \end{bmatrix}}_{\mathbf{u}_t}$$

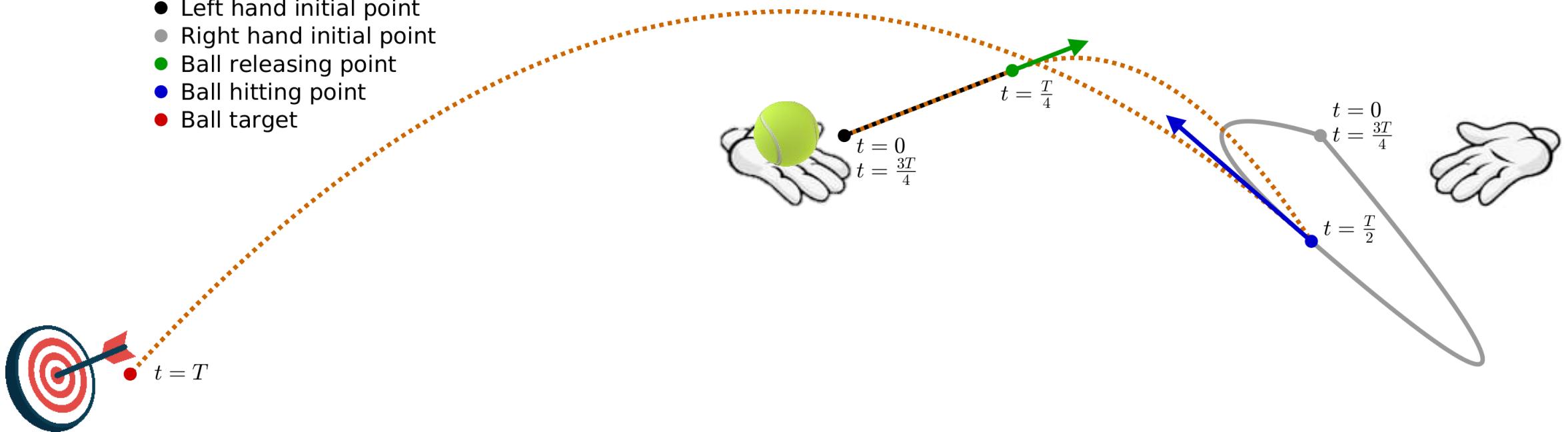
For $\frac{T}{4} < t < \frac{T}{2}$ and $t > \frac{T}{2}$ (free motion of the ball), we have

$$\begin{array}{c}
 \text{Hand} \\
 \text{Hand} \\
 \text{Ball}
 \end{array}
 \begin{bmatrix} \dot{\mathbf{x}}_1,t \\ \ddot{\mathbf{x}}_1,t \\ \dot{\mathbf{f}}_1,t \\ \dot{\mathbf{x}}_2,t \\ \ddot{\mathbf{x}}_2,t \\ \dot{\mathbf{f}}_2,t \\ \dot{\mathbf{x}}_3,t \\ \ddot{\mathbf{x}}_3,t \\ \dot{\mathbf{f}}_3,t \end{bmatrix}_{\dot{\mathbf{x}}_t} =
 \underbrace{\begin{bmatrix} 0 & \mathbf{I} & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & \mathbf{I} & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & \mathbf{I} & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & \mathbf{I} & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \end{bmatrix}}_{\mathbf{A}_t^c}
 \underbrace{\begin{bmatrix} \mathbf{x}_1,t \\ \dot{\mathbf{x}}_1,t \\ \mathbf{f}_1,t \\ \mathbf{x}_2,t \\ \dot{\mathbf{x}}_2,t \\ \mathbf{f}_2,t \\ \mathbf{x}_3,t \\ \dot{\mathbf{x}}_3,t \\ \mathbf{f}_3,t \end{bmatrix}}_{\mathbf{x}_t} +
 \underbrace{\begin{bmatrix} 0 & 0 \\ \mathbf{I} & 0 \\ 0 & 0 \\ 0 & 0 \\ 0 & \mathbf{I} \\ 0 & 0 \\ 0 & 0 \\ 0 & 0 \\ 0 & 0 \end{bmatrix}}_{\mathbf{B}_t^c}
 \underbrace{\begin{bmatrix} \mathbf{u}_{1,t} \\ \mathbf{u}_{2,t} \end{bmatrix}}_{\mathbf{u}_t}$$

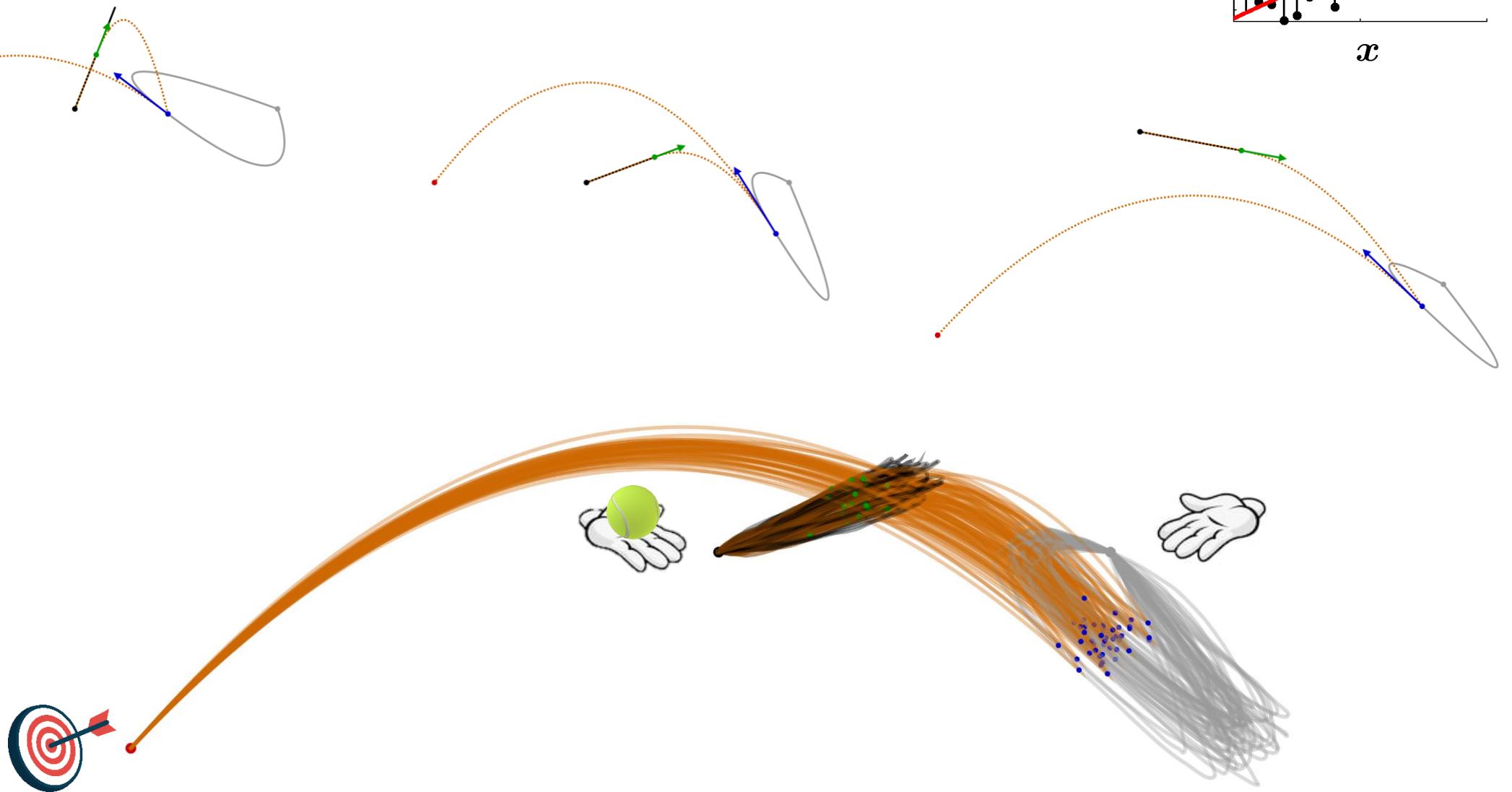
Linear quadratic tracking (LQT) - Example



- Left hand motion (Agent 1)
- Right hand motion (Agent 2)
- Ball motion (Agent 3)
- Left hand initial point
- Right hand initial point
- Ball releasing point
- Ball hitting point
- Ball target



Linear quadratic tracking (LQT) - Example



Logistic regression

Python notebook:
`demo_LS_IRLS_logRegr.ipynb`

Matlab code:
`demo_LS_IRLS_logRegr01.m`

Costs functions and associated regression solutions

Univariate output \mathbf{y} :

$$\hat{\mathbf{a}} = \arg \min_{\mathbf{a}} \|\mathbf{y} - \mathbf{X}\mathbf{a}\|^2 = \underbrace{(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top}_{\mathbf{X}^\dagger} \mathbf{y}$$

$$\hat{\mathbf{a}} = \arg \min_{\mathbf{a}} \|\mathbf{y} - \mathbf{X}\mathbf{a}\|_{F, \mathbf{W}}^2 = (\mathbf{X}^\top \mathbf{W} \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{W} \mathbf{y}$$

$$\hat{\mathbf{a}} = \arg \min_{\mathbf{a}} \|\mathbf{y} - \mathbf{X}\mathbf{a}\|_F^2 + \|\mathbf{\Gamma}\mathbf{a}\|_F^2 = (\mathbf{X}^\top \mathbf{X} + \mathbf{\Gamma}^\top \mathbf{\Gamma})^{-1} \mathbf{X}^\top \mathbf{y}$$

$$\hat{\mathbf{a}} = \arg \min_{\mathbf{a}} f_{\mathbf{a}}(\mathbf{X}, \mathbf{y})$$

Multivariate output \mathbf{y} :

$$\hat{\mathbf{A}} = \arg \min_{\mathbf{A}} \|\mathbf{Y} - \mathbf{X}\mathbf{A}\|^2 = \underbrace{(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top}_{\mathbf{X}^\dagger} \mathbf{Y}$$

$$\hat{\mathbf{A}} = \arg \min_{\mathbf{A}} \|\mathbf{Y} - \mathbf{X}\mathbf{A}\|_{F, \mathbf{W}}^2 = (\mathbf{X}^\top \mathbf{W} \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{W} \mathbf{Y}$$

$$\hat{\mathbf{A}} = \arg \min_{\mathbf{A}} \|\mathbf{Y} - \mathbf{X}\mathbf{A}\|_F^2 + \|\mathbf{\Gamma}\mathbf{A}\|_F^2 = (\mathbf{X}^\top \mathbf{X} + \mathbf{\Gamma}^\top \mathbf{\Gamma})^{-1} \mathbf{X}^\top \mathbf{Y}$$

Logistic regression

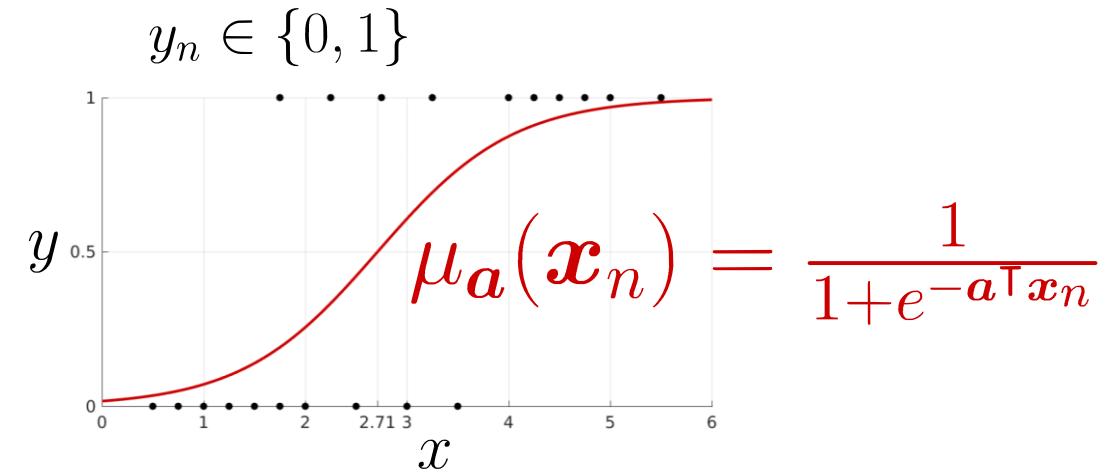
Bernoulli distribution (for binary variables):

$$\mathcal{L}_n = \begin{cases} p & \text{if } y_n = 1, \\ 1 - p & \text{if } y_n = 0. \end{cases}$$

$\mathcal{P}(y_n = 1)$
 $\mathcal{P}(y_n = 0)$

$$= p^{y_n} (1 - p)^{(1-y_n)}$$

Logistic function:



Likelihood of n^{th} datapoint:

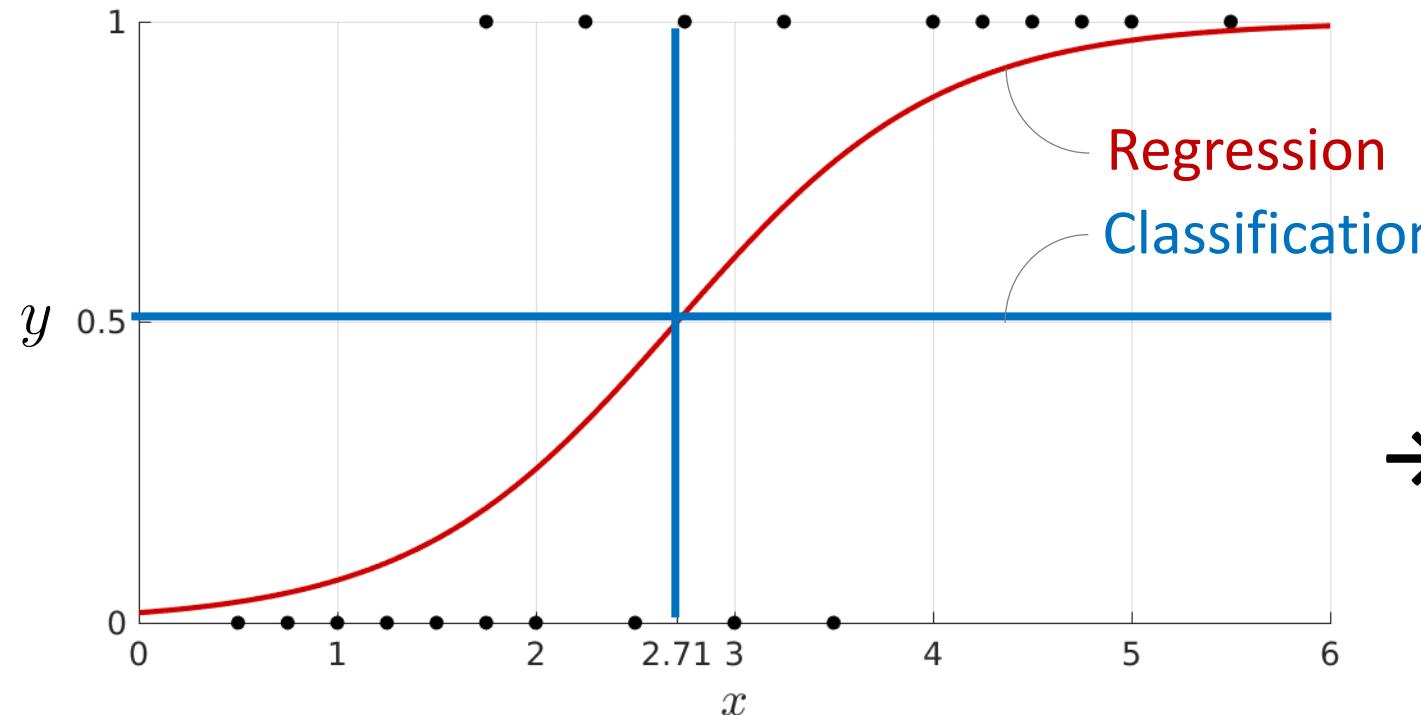
$$\mathcal{L}_n = \mu_a(\mathbf{x}_n)^{y_n} (1 - \mu_a(\mathbf{x}_n))^{(1-y_n)}$$

Likelihood of N datapoints (*independence assumption*):

$$\mathcal{L} = \prod_n \mu_a(\mathbf{x}_n)^{y_n} (1 - \mu_a(\mathbf{x}_n))^{(1-y_n)}$$

Logistic regression

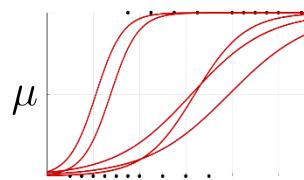
Pass/fail in function of the time spent to study at an exam:



→ **Regression exploited for classification problem**

Logistic function:

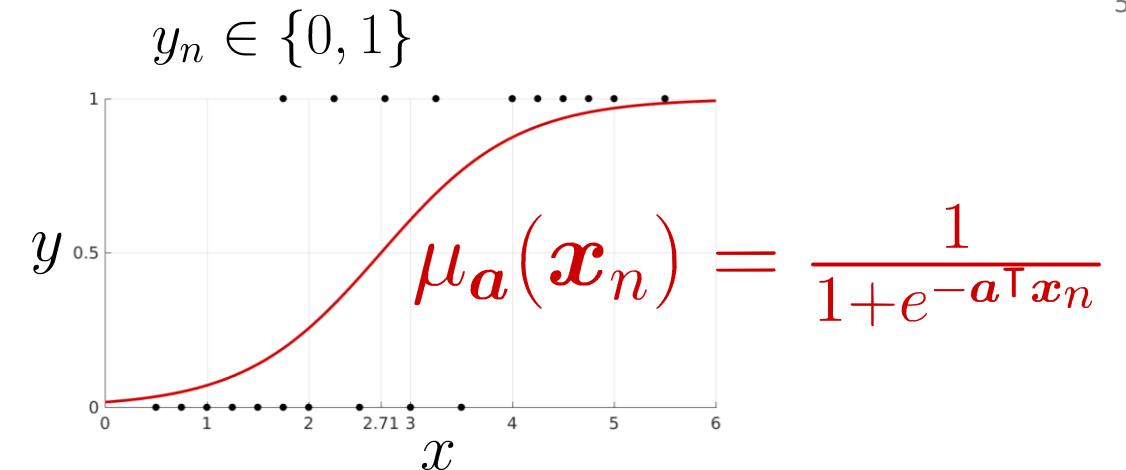
$$\mu_a(x) = \frac{1}{1+e^{-a^\top x}} \quad \mu_a(x) = \frac{1}{1+e^{-(a_1+a_2x)}}$$



Logistic regression

Likelihood of N datapoints:

$$\mathcal{L} = \prod_n \mu_{\mathbf{a}}(\mathbf{x}_n)^{y_n} (1 - \mu_{\mathbf{a}}(\mathbf{x}_n))^{(1-y_n)}$$



Cost function as negative log-likelihood:

$$c = - \sum_n y_n \log (\mu_{\mathbf{a}}(\mathbf{x}_n)) + (1 - y_n) \log (1 - \mu_{\mathbf{a}}(\mathbf{x}_n))$$

$$\begin{aligned} \frac{\partial c}{\partial \mathbf{a}} &= - \sum_n y_n \mu_{\mathbf{a}}^{-1} \mu_{\mathbf{a}} (1 - \mu_{\mathbf{a}}) \mathbf{x}_n - (1 - y_n)(1 - \mu_{\mathbf{a}})^{-1} \mu_{\mathbf{a}} (1 - \mu_{\mathbf{a}}) \mathbf{x}_n \\ &= - \sum_n y_n (1 - \mu_{\mathbf{a}}) \mathbf{x}_n - (1 - y_n) \mu_{\mathbf{a}} \mathbf{x}_n \\ &= \sum_n (\mu_{\mathbf{a}} - y_n) \mathbf{x}_n \end{aligned}$$

$$\mu(x) = \frac{1}{1+e^{-x}}$$

$$\frac{\partial \mu}{\partial x} = \mu(1 - \mu)$$

Logistic regression

It can for example be solved with Newton's method, by iterating

$$\frac{\partial c}{\partial \mathbf{a}} = \sum_n (\mu_{\mathbf{a}} - y_n) \mathbf{x}_n$$

$$\mu_{\mathbf{a}}(\mathbf{x}_n) = \frac{1}{1+e^{-\mathbf{a}^\top \mathbf{x}_n}}$$

$$\mathbf{a} \leftarrow \mathbf{a} - \mathbf{H}^{-1} \mathbf{g},$$

with gradient $\mathbf{g} = \sum_n (\mu_{\mathbf{a}}(\mathbf{x}_n) - y_n) \mathbf{x}_n = \mathbf{X}^\top (\mu_{\mathbf{a}} - \mathbf{y})$ and Hessian $\mathbf{H} = \mathbf{X}^\top \mathbf{W} \mathbf{X}$, with diagonal matrix $\mathbf{W} = \text{diag}(\mu_{\mathbf{a}} * (1 - \mu_{\mathbf{a}}))$.

We then obtain

$$\begin{aligned} \mathbf{a} &\leftarrow \mathbf{a} - \mathbf{H}^{-1} \mathbf{g} \\ &\leftarrow \mathbf{a} - (\mathbf{X}^\top \mathbf{W} \mathbf{X})^{-1} \mathbf{X}^\top (\mu_{\mathbf{a}} - \mathbf{y}) \\ &\leftarrow (\mathbf{X}^\top \mathbf{W} \mathbf{X})^{-1} (\mathbf{X}^\top \mathbf{W} \mathbf{X} \mathbf{a} - \mathbf{X}^\top (\mu_{\mathbf{a}} - \mathbf{y})) \\ &\leftarrow (\mathbf{X}^\top \mathbf{W} \mathbf{X})^{-1} \mathbf{X}^\top (\mathbf{W} \mathbf{X} \mathbf{a} + \mathbf{y} - \mu_{\mathbf{a}}) \\ &\leftarrow (\mathbf{X}^\top \mathbf{W} \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{W} \mathbf{z}, \end{aligned}$$

with *working response* $\mathbf{z} = \mathbf{X} \mathbf{a} + \mathbf{W}^{-1} (\mathbf{y} - \mu_{\mathbf{a}})$.

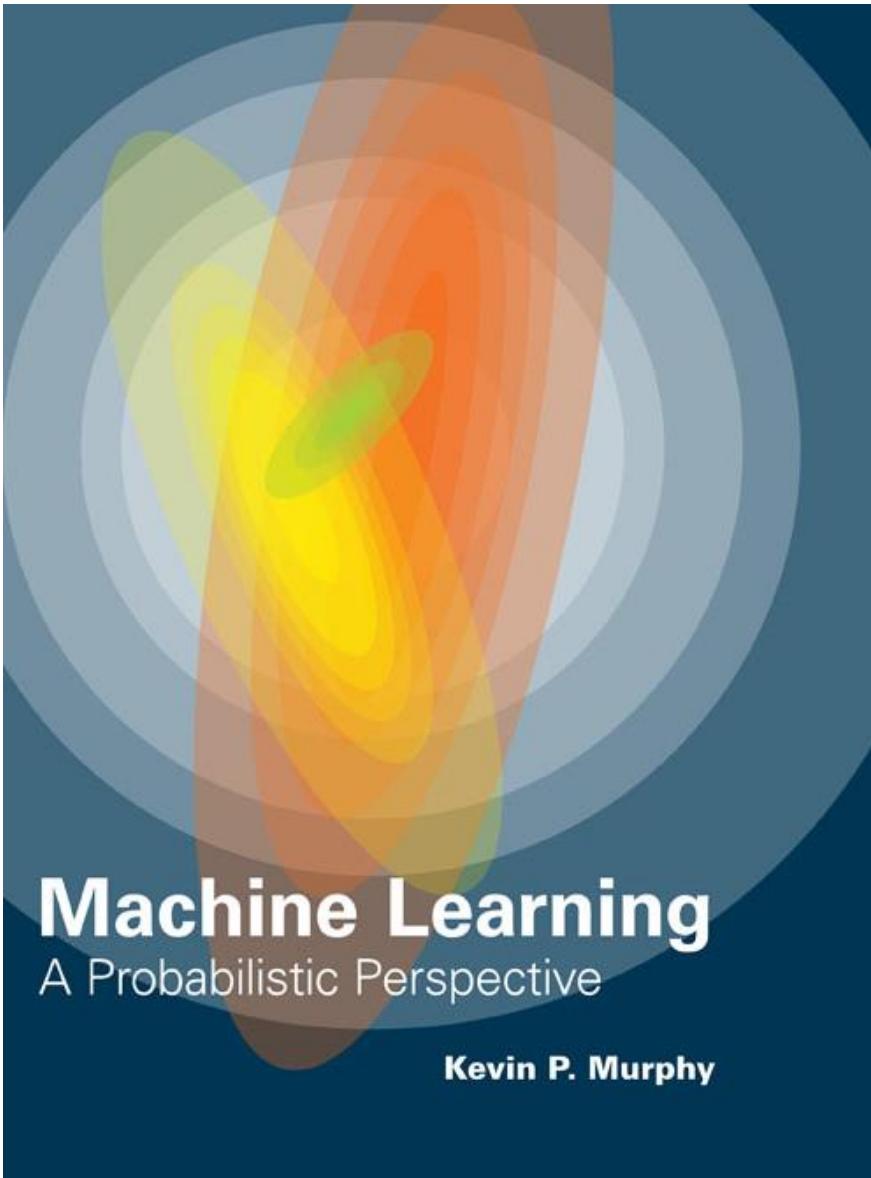
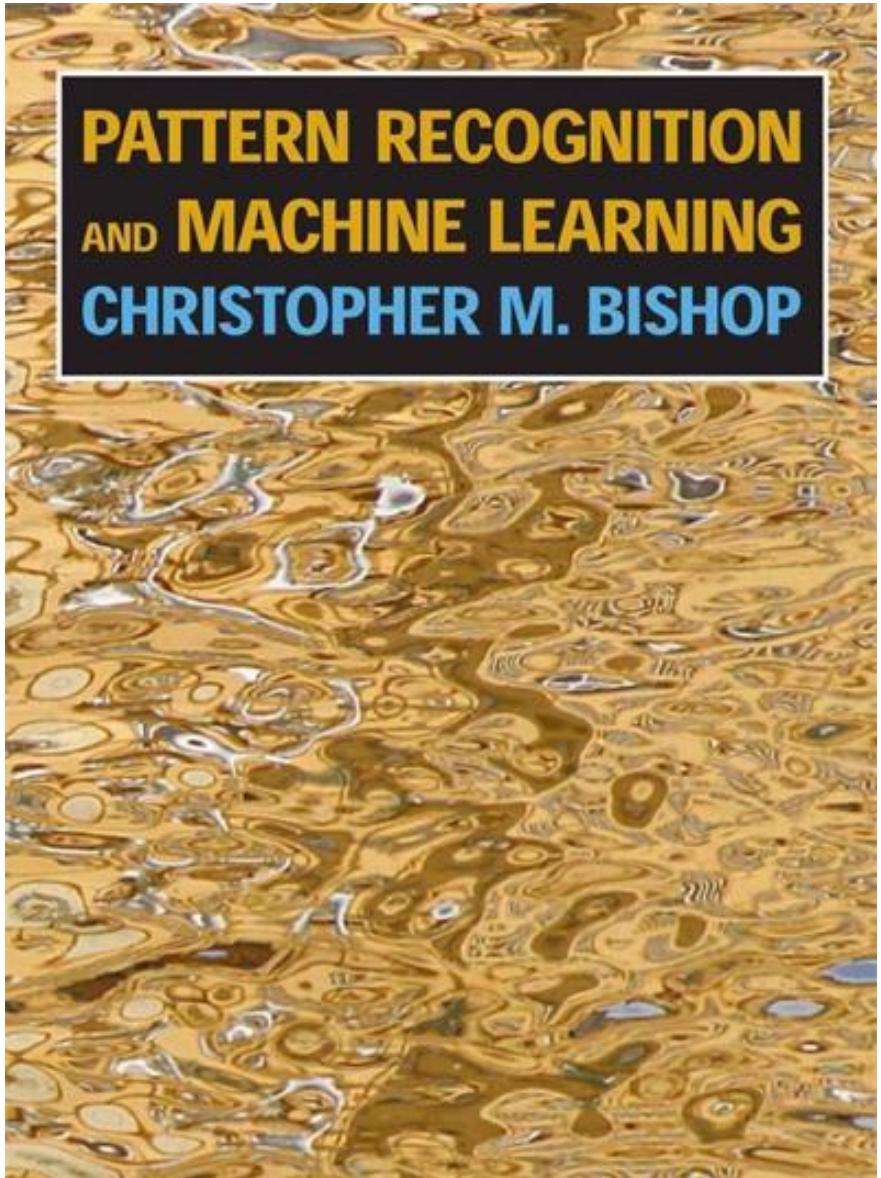
Hadamard (elementwise) product

$$\mu(x) = \frac{1}{1+e^{-x}}$$

$$\frac{\partial \mu}{\partial x} = \mu(1 - \mu)$$

→ IRLS algorithm

Recommended reading material



**The Matrix
Cookbook**

Kaare Brandt Petersen
Michael Syskind Pedersen

Labs this afternoon

<https://moodle.epfl.ch/course/view.php?id=16819>

INM11, 14:15-16:00



Dr João Silvério

