

EE613 - Machine Learning for Engineers

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

LINEAR REGRESSION

Sylvain Calinon

Robot Learning and Interaction Group

Idiap Research Institute

Nov 2, 2023

EE613 schedule

Thu. 21.09.2023	(C) 1. ML introduction
Thu. 28.09.2023	(C) 2. Bayesian 1 (C) 3. Bayesian 2
Thu. 12.10.2023	(C) 4. Hidden Markov Models
Thu. 19.10.2023	(C) 5. Dimensionality reduction
Thu. 26.10.2023	(C) 6. Decision trees
Thu. 02.11.2023	(C) 7. Linear regression
Thu. 09.11.2023	(C) 8. Nonlinear regression
Thu. 16.11.2023	(C) 9. Kernel Methods - SVM
Thu. 23.11.2023	(C) 10. Tensor factorization
Thu. 30.11.2023	(C) 11. Deep learning 1
Thu. 07.12.2023	(C) 12. Deep learning 2
Thu. 14.12.2023	(C) 13. Deep learning 3
Thu. 21.12.2023	(C) 14. Deep learning 4

Regression

Nonlinear regression:

Approximation function
(predictor function)
(forecasting function)

$$Y = f(X)$$

Output
(response variable)
(outcome variable)

Input
(explanatory variable)
(feature variable)

Linear regression:

$$Y = XA$$

Tensor regression:

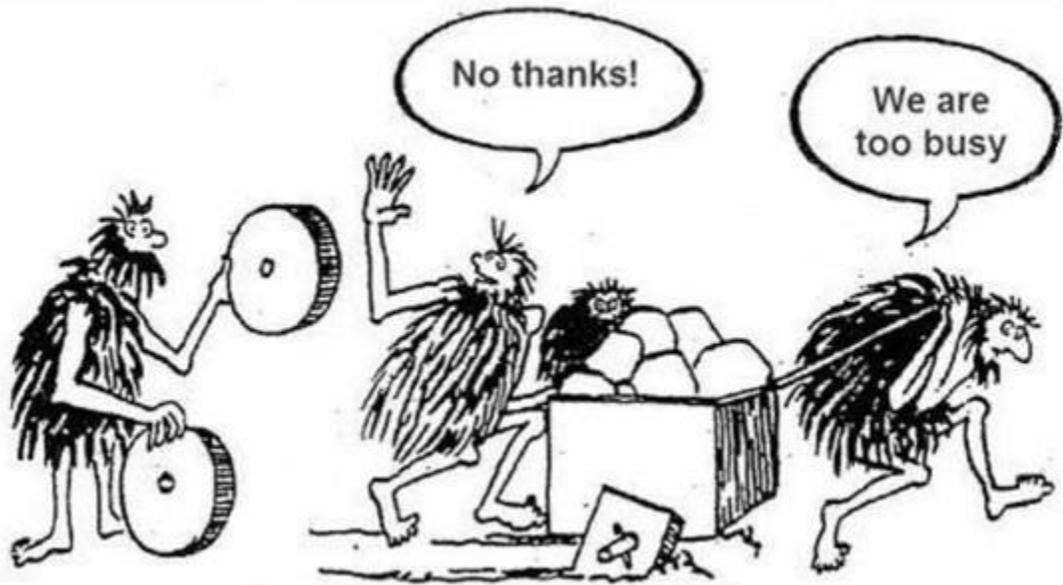
Similar to $Y = XA$,
but with X a tensor
instead of a matrix

Nonlinear regression on x reformulated as linear regression on f :

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

LEAST SQUARES

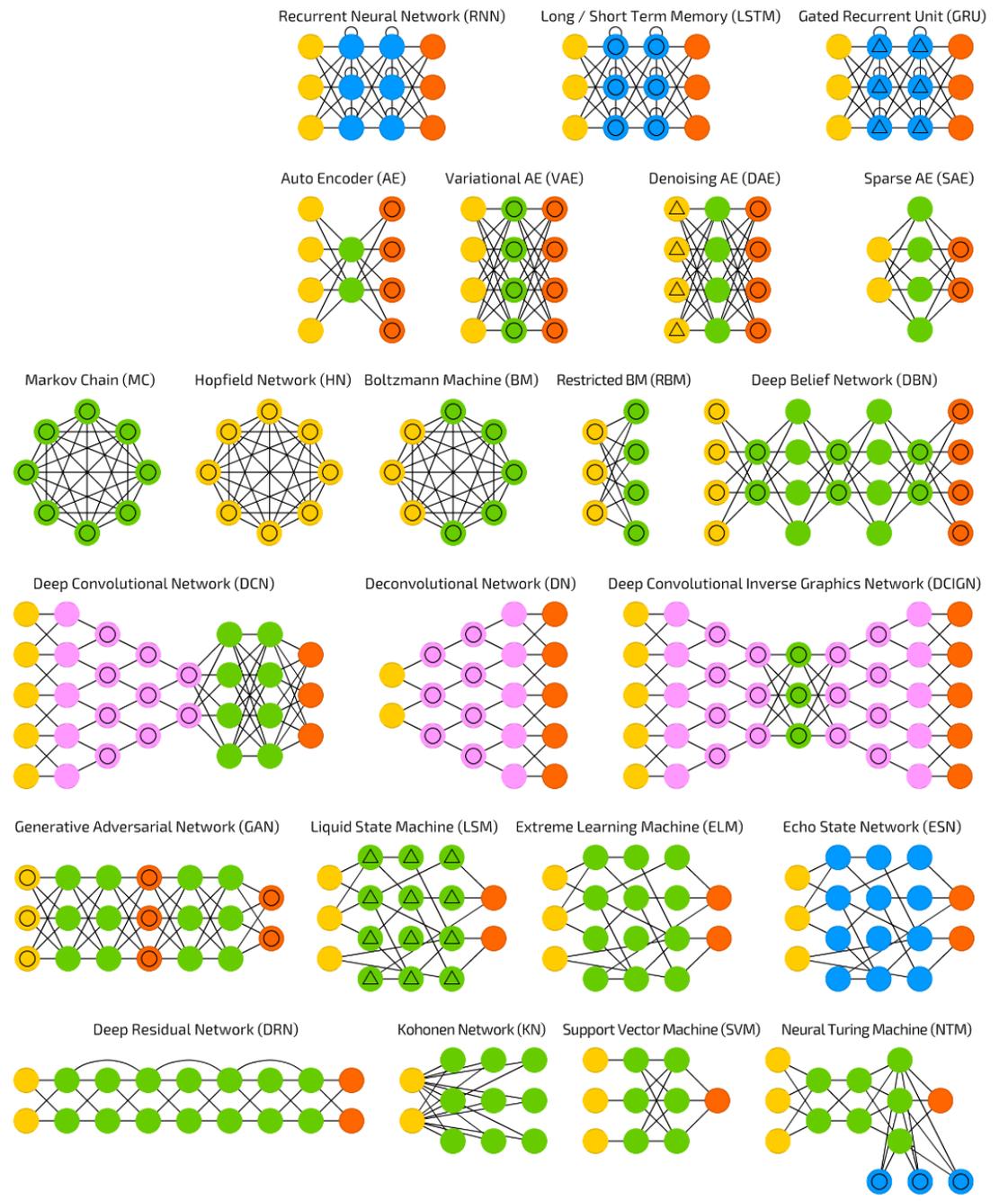
circa 1795



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

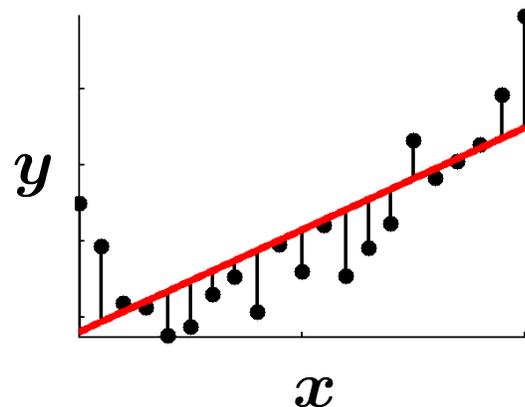
$$\hat{A} = X^\dagger Y$$

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



Least squares: a ubiquitous tool

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



Weighted least squares?

Regularized least squares?

L1-norm instead of L2-norm?

Fast and robust
implementation?

Recursive computation?

Solution with a secondary objective?

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 $\mathbf{y} = \mathbf{X}\mathbf{a}$.

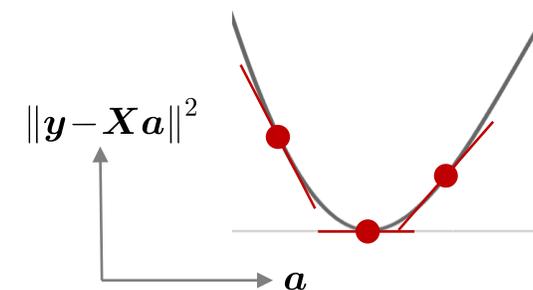
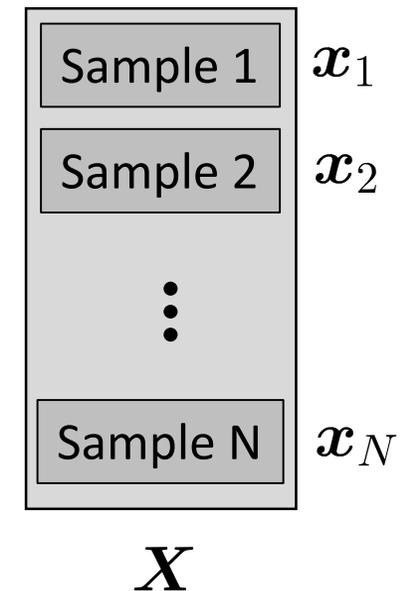
A solution can be found by minimizing the ℓ_2 norm

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

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

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

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

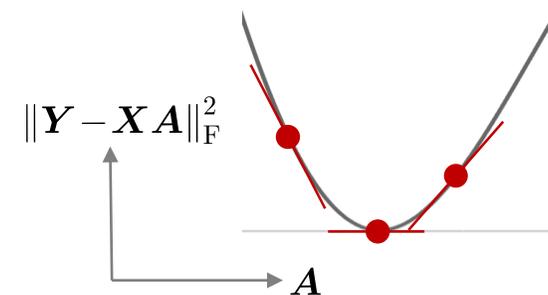
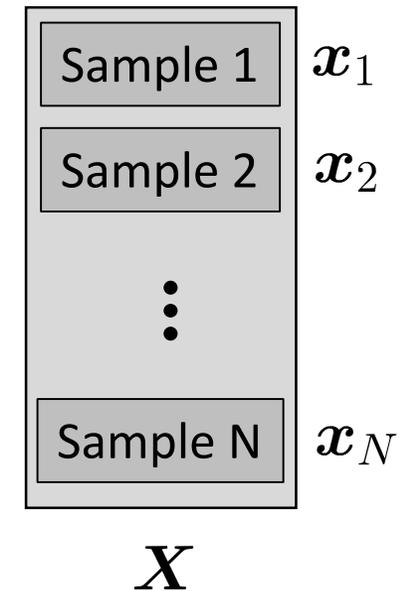
A solution can be found by minimizing the Frobenius norm

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

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

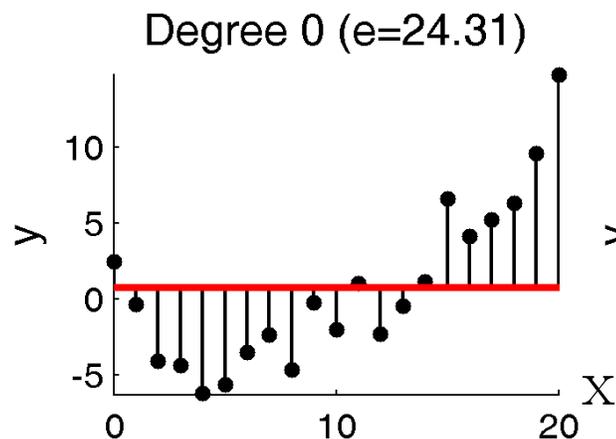
$$-2 \mathbf{X}^{\top} \mathbf{Y} + 2 \mathbf{X}^{\top} \mathbf{X} \mathbf{A} = \mathbf{0} \quad \iff \quad \hat{\mathbf{A}} = (\mathbf{X}^{\top} \mathbf{X})^{-1} \mathbf{X}^{\top} \mathbf{Y}$$

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

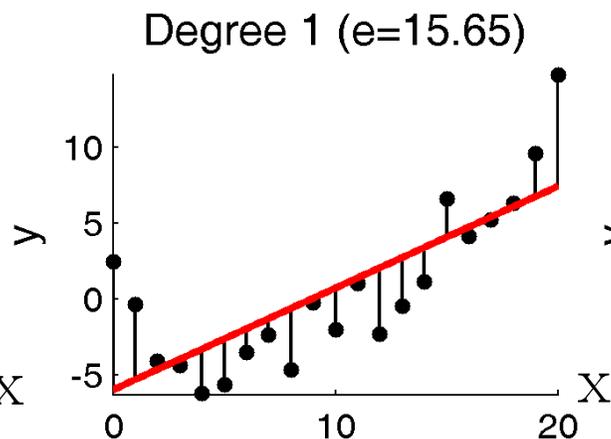


Polynomial fitting with least squares

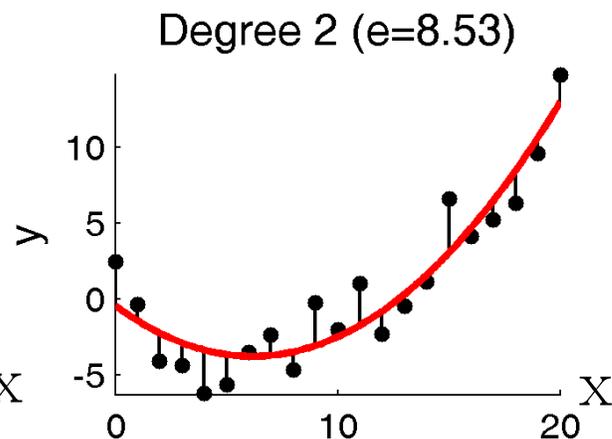
$$\hat{A} = X^\dagger Y$$



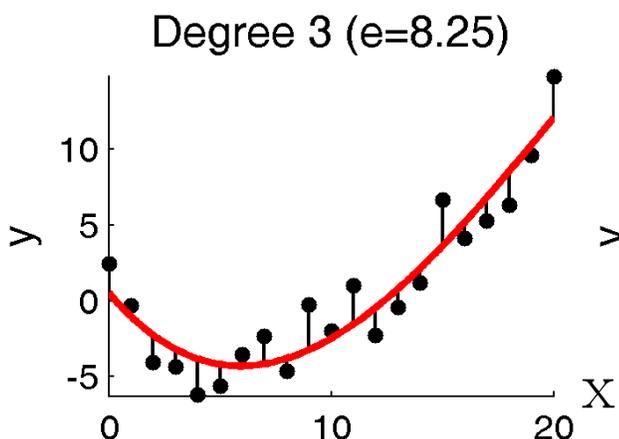
$$\mathbf{x} = 1$$



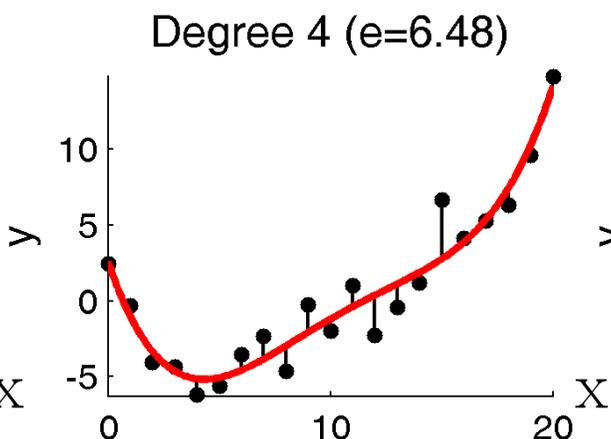
$$\mathbf{x} = [1, x]$$



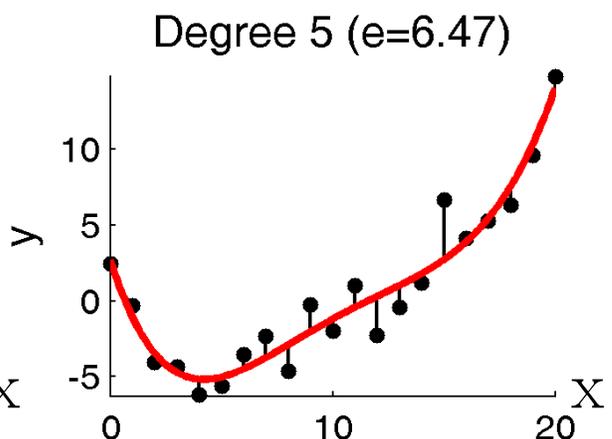
$$\mathbf{x} = [1, x, x^2]$$



$$\mathbf{x} = [1, x, x^2, x^3]$$



$$\mathbf{x} = [1, x, x^2, x^3, x^4]$$



$$\mathbf{x} = [1, x, x^2, x^3, x^4, x^5]$$

Least squares computed with SVD

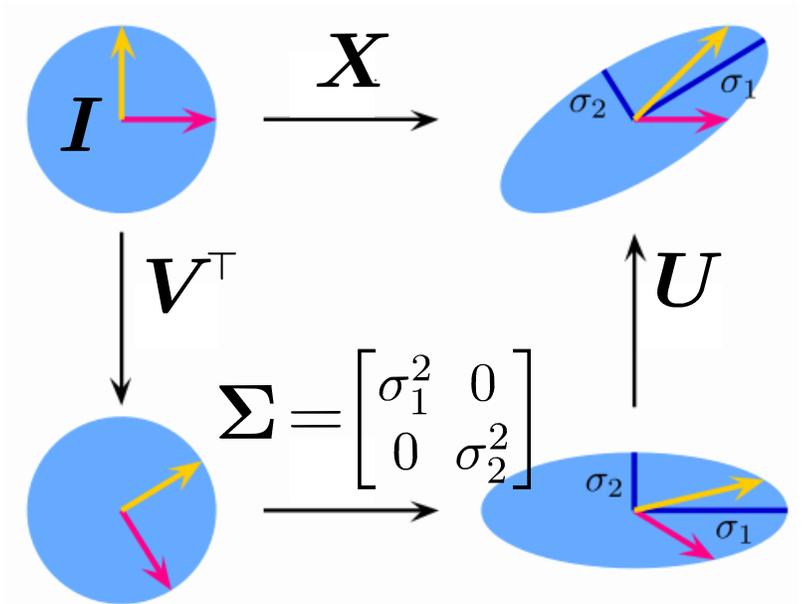
Singular value decomposition (SVD)

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

Unitary matrix
(orthogonal)

Unitary matrix
(orthogonal)

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



Least squares computed with SVD

$$\hat{\mathbf{A}} = \overbrace{(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top}^{\mathbf{X}^\dagger} \mathbf{Y} \quad \mathbf{X} \in \mathbb{R}^{N \times D^{\mathcal{I}}}$$

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

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

where \mathbf{U} and \mathbf{V} are $N \times N$ and $D^{\mathcal{I}} \times D^{\mathcal{I}}$ orthogonal matrices, and $\mathbf{\Sigma}$ is an $N \times D^{\mathcal{I}}$ matrix with all its elements outside of the main diagonal equal to 0.

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



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

With this decomposition, a solution to the least squares problem is given by

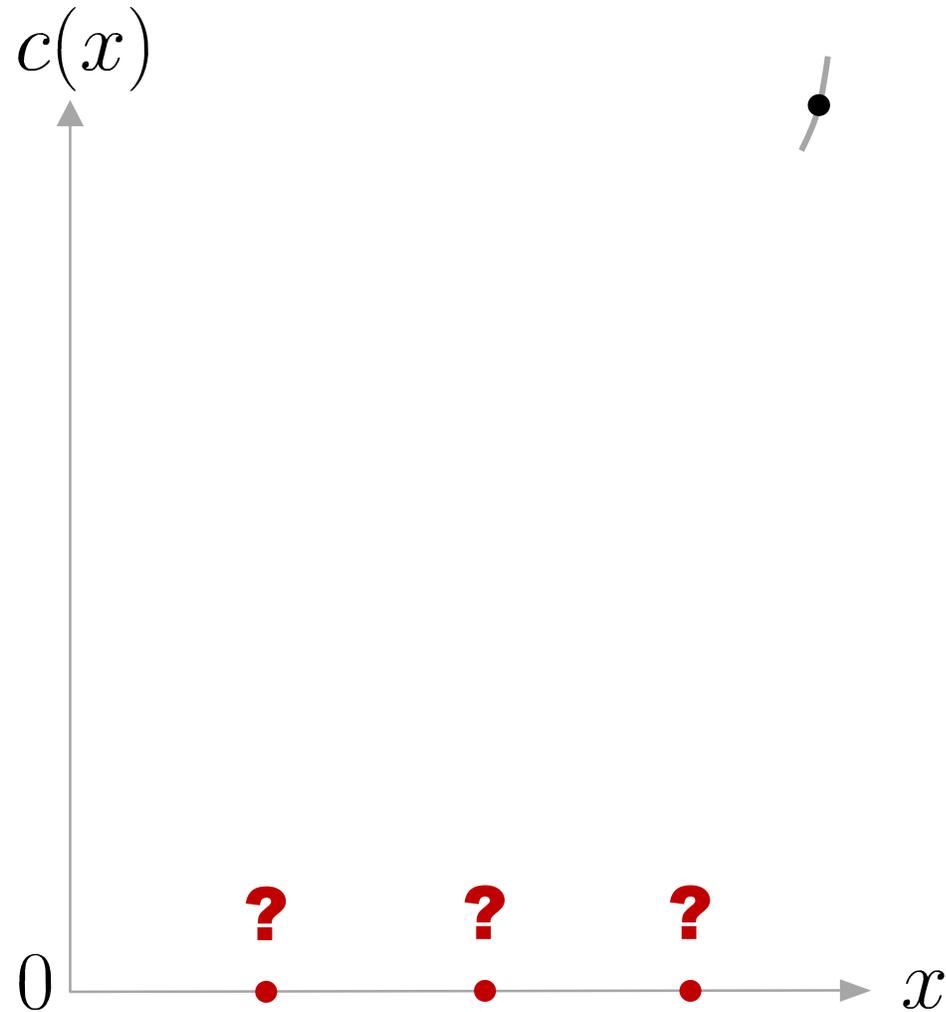
$$\hat{\mathbf{A}} = \mathbf{V} \mathbf{\Sigma}^\dagger \mathbf{U}^\top \mathbf{Y}$$

where the pseudoinverse of $\mathbf{\Sigma}$ can be easily obtained by inverting the non-zero diagonal elements and transposing the resulting matrix.

Newton's method

(least squares problem solved at each iteration of the optimization process)

Newton's method for **minimization**



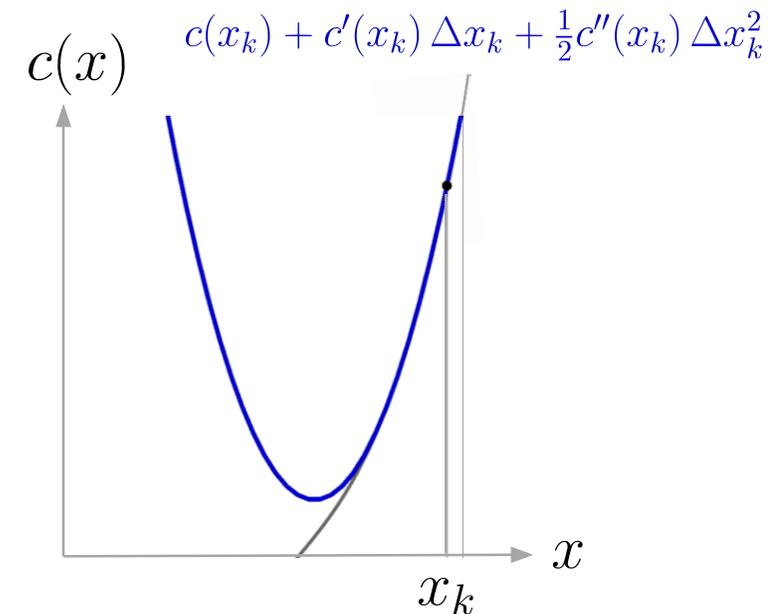
Newton's method for minimization

Newton's method attempts to solve $\min_x c(x)$ or $\max_x c(x)$ from an initial guess x_1 by using a sequence of second-order Taylor approximations of $c(x)$ around the iterates.

The second-order Taylor expansion of $c(x)$ around x_k is

$$c(x_k + \Delta x_k) \approx c(x_k) + c'(x_k) \Delta x_k + \frac{1}{2} c''(x_k) \Delta x_k^2.$$

The next iterate $x_{k+1} = x_k + \Delta x_k$ is defined so as to minimize this quadratic approximation in Δx_k .



Newton's method for minimization

If the second derivative is positive, the quadratic approximation is a convex function of Δx_k , and its minimum can be found by setting the derivative to zero.

Since

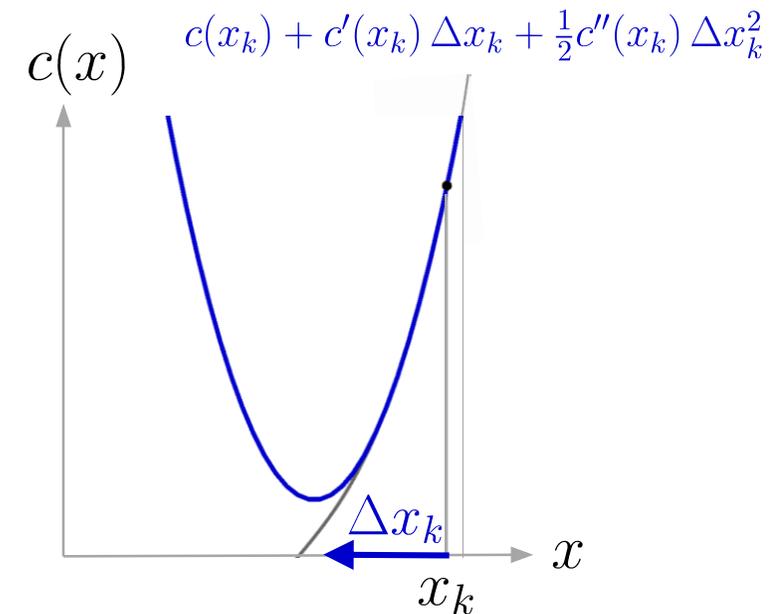
$$\frac{d}{d\Delta x_k} \left(c(x_k) + c'(x_k) \Delta x_k + \frac{1}{2} c''(x_k) \Delta x_k^2 \right) = c'(x_k) + c''(x_k) \Delta x_k = 0,$$

the minimum is achieved for

$$\Delta x_k = -\frac{c'(x_k)}{c''(x_k)}.$$

Newton's method thus performs the iteration

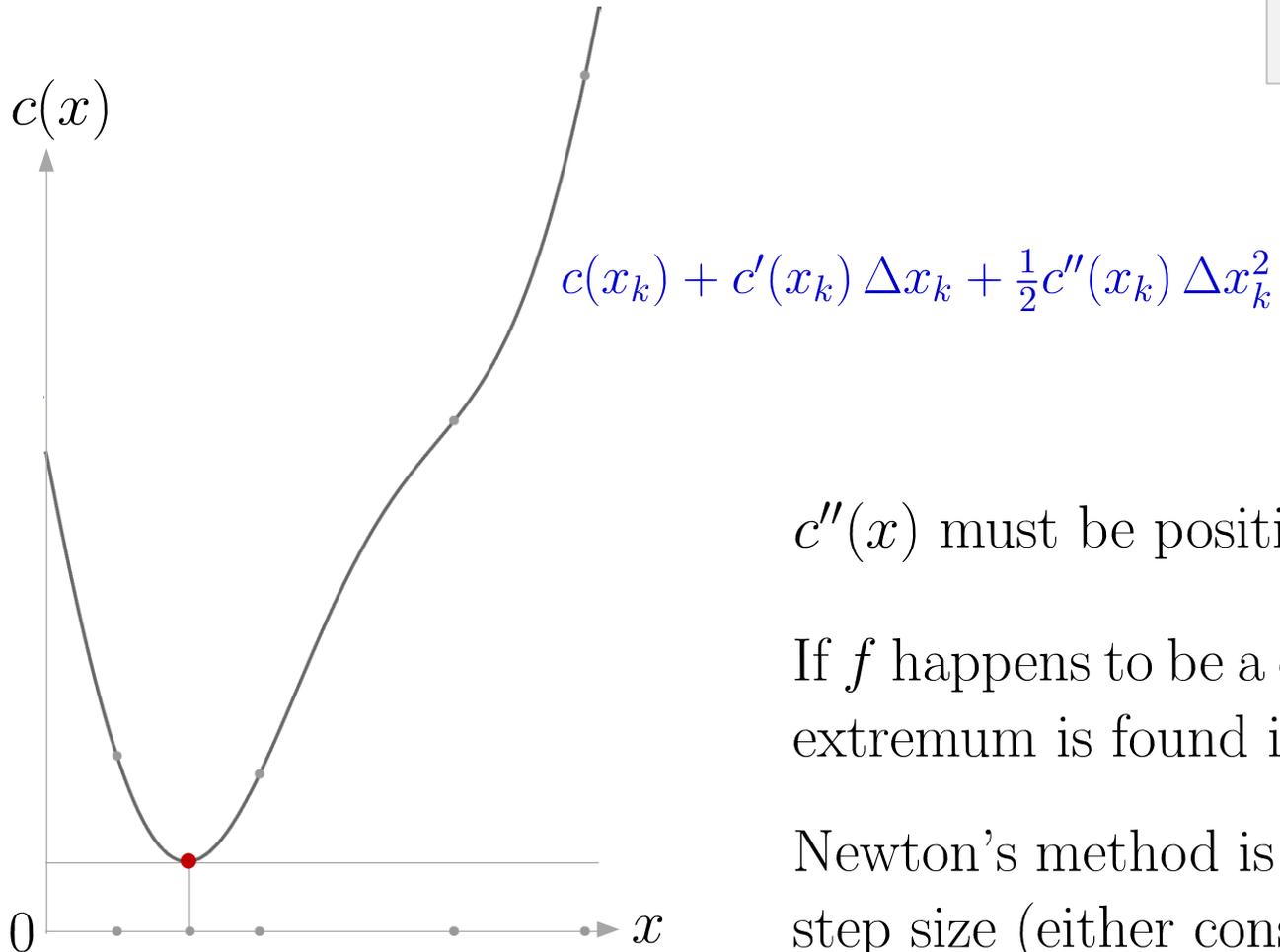
$$x_{k+1} = x_k - \frac{c'(x_k)}{c''(x_k)}.$$



Newton's method for **minimization**

$$x_{k+1} = x_k - \frac{c'(x_k)}{c''(x_k)}$$

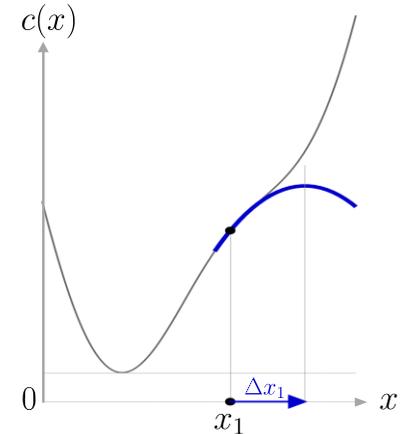
The geometric interpretation of Newton's method is that at each iteration, the goal is to **fit a paraboloid** to the surface of $c(x)$ at x_k , having the same slope and curvature as the surface at that point, and then move to the maximum or minimum of that paraboloid.



$c''(x)$ must be positive to find a minimum.

If f happens to be a quadratic function, then the exact extremum is found in one step.

Newton's method is often modified to include a small step size (either constant or estimated).



Newton's method for **minimization**

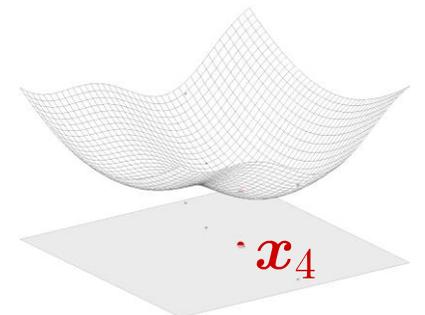
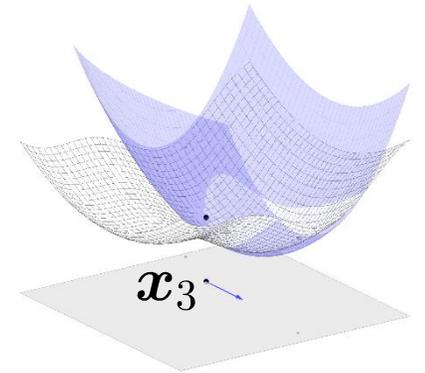
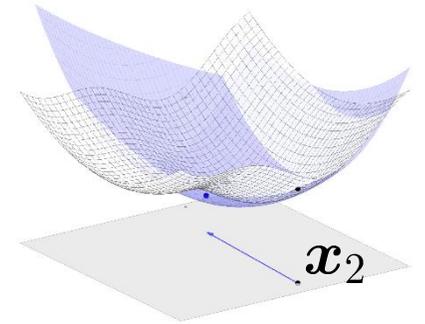
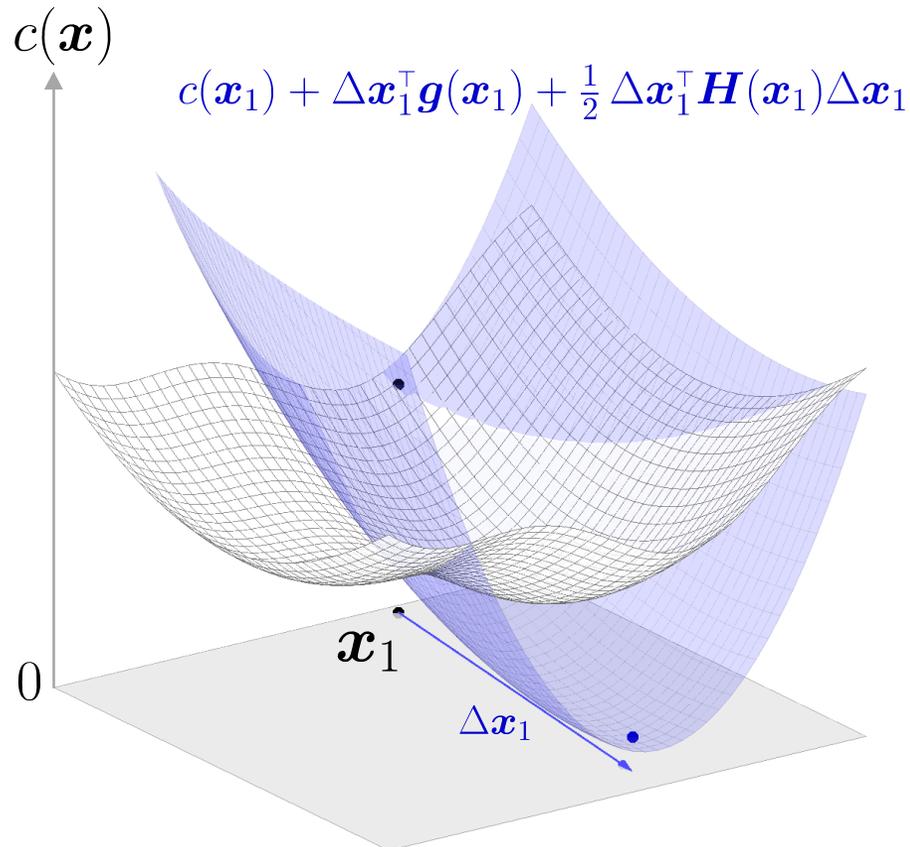
The multidimensional case similarly provides

$$\mathbf{x}_{k+1} = \mathbf{x}_k - \mathbf{H}(\mathbf{x}_k)^{-1} \mathbf{g}(\mathbf{x}_k),$$

with \mathbf{g} and \mathbf{H} the **gradient** and **Hessian** matrix of f .

(1D case)

$$x_{k+1} = x_k - \frac{c'(x_k)}{c''(x_k)}$$



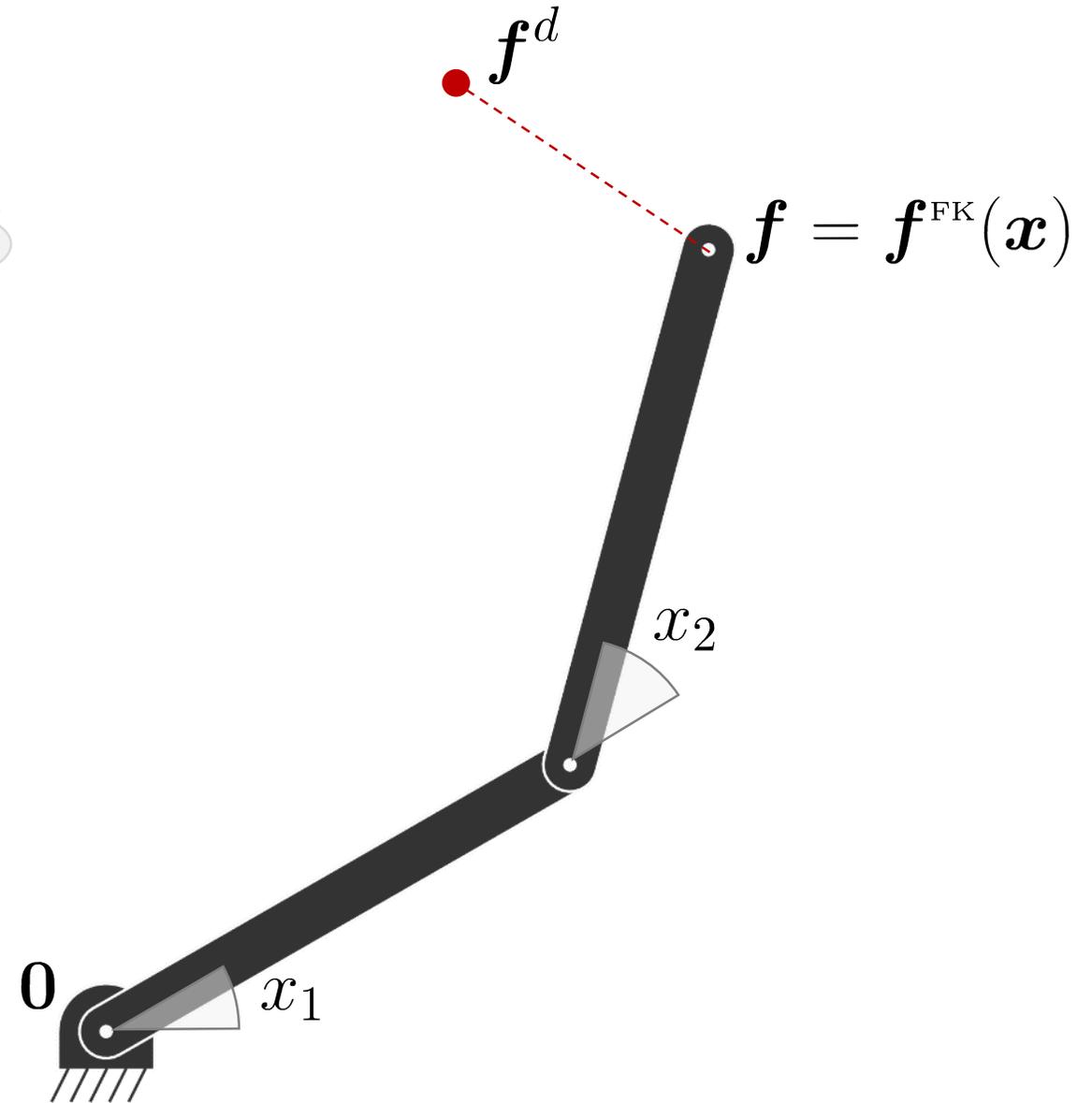
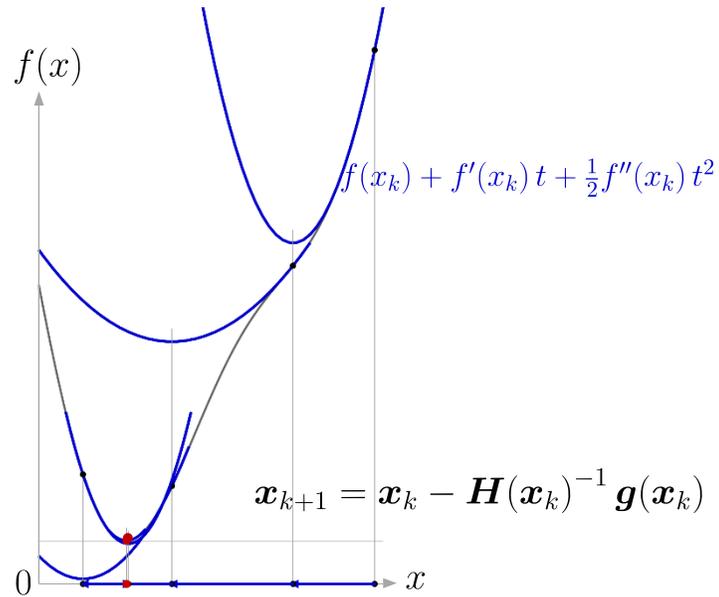
Example: Robot inverse kinematics (IK)

Forward kinematics (FK): $f = f^{\text{FK}}(x)$

Inverse kinematics (IK): $x = f^{\text{FK}-1}(f)$?

IK in practice: Minimizing $\|f^d - f^{\text{FK}}(x)\|^2$

→ Newton's method



Newton's method applied to IK: Gauss-Newton algorithm

The **Gauss-Newton algorithm** is a special case of Newton's method in which the cost is quadratic (sum of squared function values), with $c(\mathbf{x}) = \sum_{i=1}^{D_f} r_i^2(\mathbf{x}) = \mathbf{r}^\top(\mathbf{x}) \mathbf{r}(\mathbf{x})$, and by ignoring the second-order derivative terms of the Hessian.

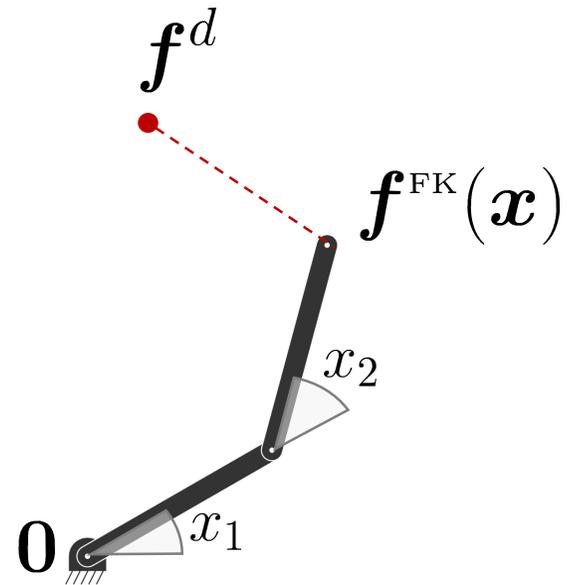
The gradient and Hessian can in this case be computed with

$$\mathbf{g}(\mathbf{x}) = 2 \mathbf{J}^\top(\mathbf{x}) \mathbf{r}(\mathbf{x}), \quad \mathbf{H}(\mathbf{x}) \approx 2 \mathbf{J}^\top(\mathbf{x}) \mathbf{J}(\mathbf{x}),$$

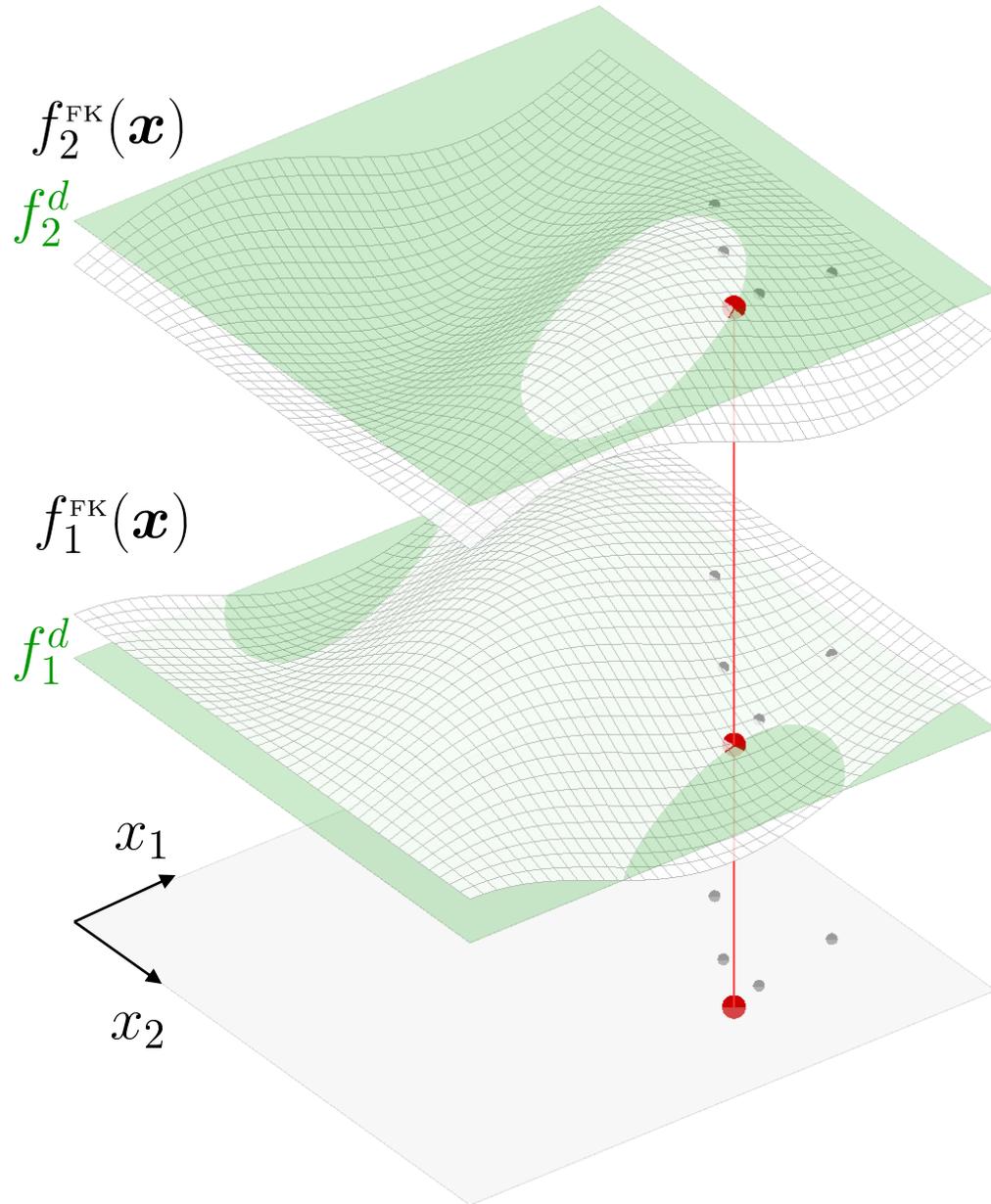
with $\mathbf{J}(\mathbf{x}) = \frac{\partial \mathbf{r}(\mathbf{x})}{\partial \mathbf{x}} \in \mathbb{R}^{D_f \times D_x}$ the Jacobian matrix of $\mathbf{r}(\mathbf{x}) \in \mathbb{R}^{D_f}$.

By considering $\mathbf{r}(\mathbf{x}) = \mathbf{f}^{\text{FK}}(\mathbf{x}) - \mathbf{f}^d$ for our IK problem, the update at each iteration k then becomes

$$\begin{aligned} \mathbf{x}_{k+1} &= \mathbf{x}_k - \mathbf{H}(\mathbf{x}_k)^{-1} \mathbf{g}(\mathbf{x}_k) \\ &= \mathbf{x}_k - \underbrace{\left(\mathbf{J}^\top(\mathbf{x}_k) \mathbf{J}(\mathbf{x}_k) \right)^{-1} \mathbf{J}^\top(\mathbf{x}_k)}_{\mathbf{J}^\dagger(\mathbf{x}_k)} \mathbf{r}(\mathbf{x}_k) \\ &= \mathbf{x}_k + \mathbf{J}^\dagger(\mathbf{x}_k) (\mathbf{f}^d - \mathbf{f}^{\text{FK}}(\mathbf{x})). \end{aligned}$$

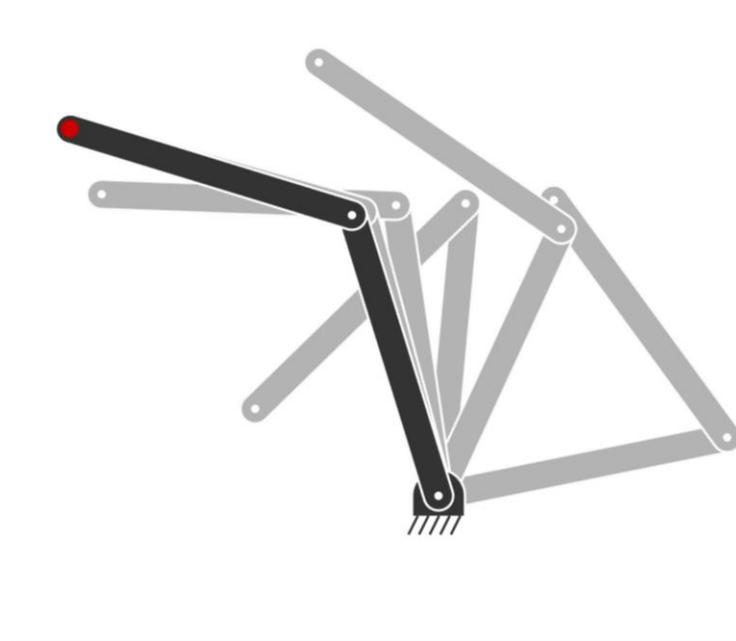


Example: Robot inverse kinematics (IK)



$$\mathbf{x}_{k+1} = \mathbf{x}_k + \mathbf{J}^\dagger(\mathbf{x}_k) \overbrace{\left(\mathbf{f}^d - \mathbf{f}^{\text{FK}}(\mathbf{x}) \right)}^{r(\mathbf{x})}$$

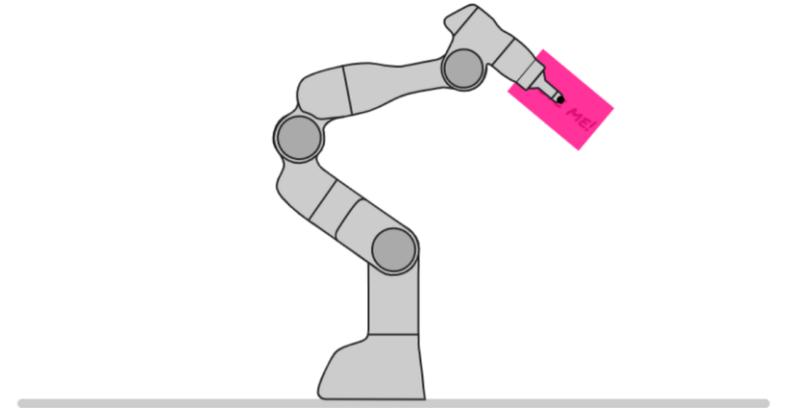
$$\mathbf{J}(\mathbf{x}_k) = \begin{bmatrix} \frac{\partial r_1(\mathbf{x}_k)}{\partial f_{1,k}} & \frac{\partial r_1(\mathbf{x}_k)}{\partial f_{2,k}} \\ \frac{\partial r_2(\mathbf{x}_k)}{\partial f_{1,k}} & \frac{\partial r_2(\mathbf{x}_k)}{\partial f_{2,k}} \end{bmatrix} \in \mathbb{R}^{2 \times 2}$$



Example: Robot inverse kinematics (IK)

```
IK  Damped IK  Weighted IK  Prioritized IK
1  x = [-np.pi/4, np.pi/2, np.pi/4] # Initial robot state
2  def controlCommand(x, param):
3      f = fkin(x, param)
4      J = Jkin(x, param)
5      u = np.linalg.pinv(J) @ logmap(param.Mu, f) # Position & orientation tracking
6      #u = J.T @ logmap(param.Mu, f) * 1E-4 # Gradient-based tracking
7      #u = np.linalg.pinv(J[:2,:]) @ (param.Mu[:2] - f[:2]) # Position tracking
8      #u = np.linalg.pinv(J[2:,:]) @ (param.Mu[2:] - f[2:]) # Orientation tracking
9      #u = np.zeros(param.nbVarX) # Zero control commands
10 return u
```

Object orientation



<https://robotics-codes-from-scratch.github.io/>

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 $\|\mathbf{Y} - \mathbf{X}\hat{\mathbf{A}}\|_{\text{F}}^2$ is still minimized.

Nullspace projection *computed with SVD*

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

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

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

to compute

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

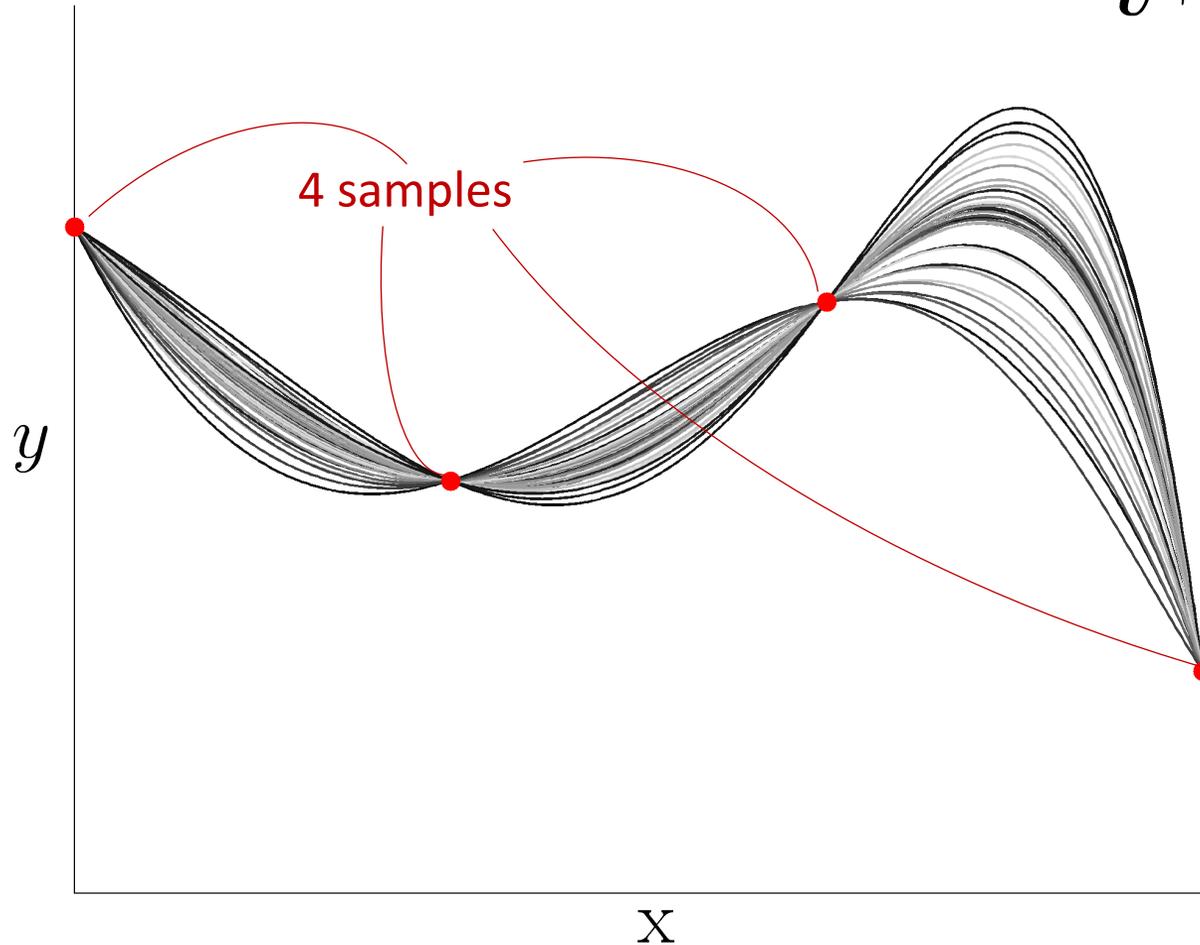
where \tilde{U} is a matrix formed by the columns of U that span for the corresponding zero rows in Σ .

$$\begin{array}{c}
 \overbrace{X^\dagger \in \mathbb{R}^{D^I \times N}} \\
 \left[\begin{array}{ccccc} 1 & 0 & 0 & 0 & 2 \\ 0 & 0 & 3 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \\ 0 & 4 & 0 & 0 & 0 \end{array} \right] = \overbrace{U \in \mathbb{R}^{D^I \times D^I}} \\
 \left[\begin{array}{ccccc} 0 & 0 & 1 & 0 & \\ 0 & 1 & 0 & 0 & \\ 0 & 0 & 0 & -1 & \\ 1 & 0 & 0 & 0 & \end{array} \right] \underbrace{\tilde{U}} \\
 \overbrace{\Sigma \in \mathbb{R}^{D^I \times N}} \\
 \left[\begin{array}{ccccc} 4 & 0 & 0 & 0 & 0 \\ 0 & 3 & 0 & 0 & 0 \\ 0 & 0 & \sqrt{5} & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \end{array} \right] \overbrace{V^\top \in \mathbb{R}^{N \times N}} \\
 \left[\begin{array}{ccccc} 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 \\ \sqrt{0.2} & 0 & 0 & 0 & \sqrt{0.8} \\ 0 & 0 & 0 & 1 & 0 \\ -\sqrt{0.8} & 0 & 0 & 0 & \sqrt{0.2} \end{array} \right]
 \end{array}$$

Example: Polynomial fitting

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

$$\mathbf{v} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$$



$$\hat{\mathbf{a}} \in \mathbb{R}^7$$

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

$$\mathbf{y} \in \mathbb{R}^4$$

$$\mathbf{N} \in \mathbb{R}^{7 \times 7}$$

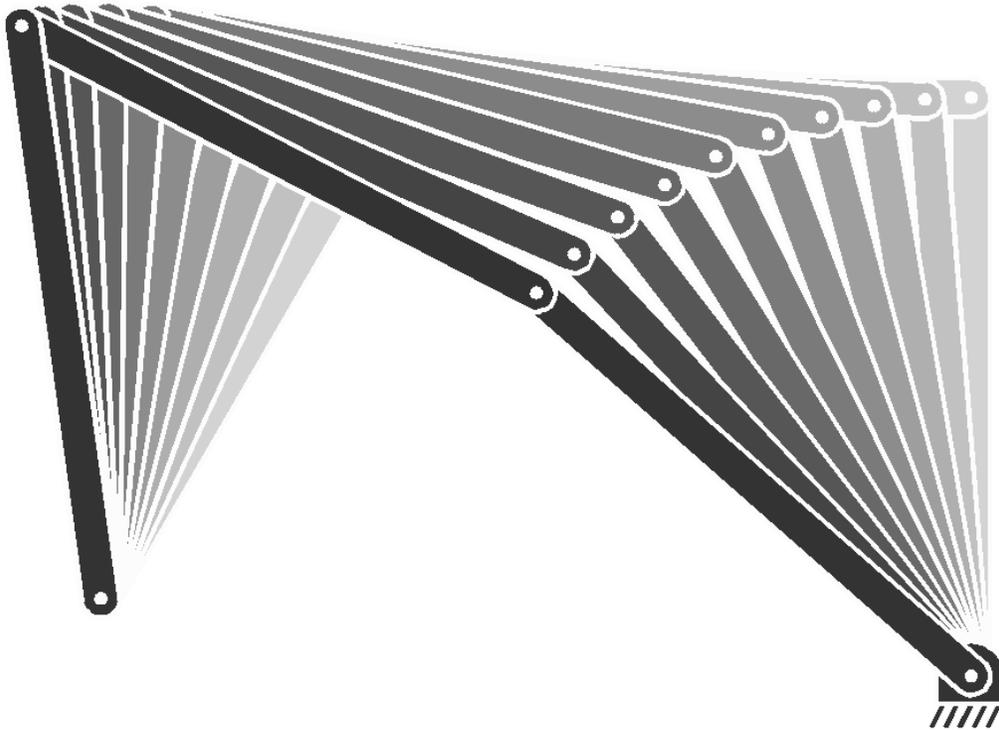
$$\mathbf{v} \in \mathbb{R}^7$$

Example: Robot inverse kinematics

Primary objective:
keeping the tip of
the robot still

$$\Delta \mathbf{x} = \mathbf{J}^\dagger(\mathbf{x}) \Delta \mathbf{f} + \mathbf{N}(\mathbf{x}) \Delta \mathbf{x}^{\text{secondary}}$$

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



Secondary objective:
moving the
first joint

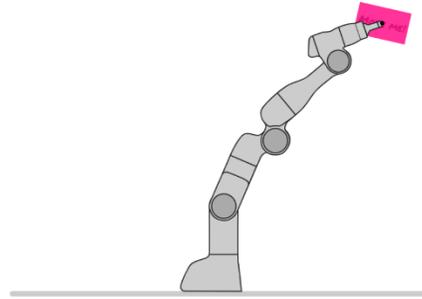
Example: Robot inverse kinematics (single and dual arms)

```

IK Damped IK Weighted IK Prioritized IK
1 x = [-np.pi/4, np.pi/2, np.pi/4] # Initial robot state
2 def controlCommand(x, param):
3     f = fkin(x, param)
4     J = Jkin(x, param)
5
6     #N = np.eye(param.nbVarX) - np.linalg.pinv(J[:2,:]) @ J[:2,:] # Nullspace project
7     #u = u + N @ [1, 0, 0] # Prioritized tracking
8
9     # Prioritized control (position tracking prioritized over orientation tracking)
10    dfp = (param.Mu[:2] - f[:2]) * 10 # Position correction
11    dfo = (param.Mu[2:] - f[2:]) * 10 # Orientation correction
12    Jp = J[:2,:] # Jacobian for position
13    Jo = J[2:,:] # Jacobian for orientation
14    pinvJp = np.linalg.inv(Jp.T @ Jp + np.eye(param.nbVarX) * 1e-2) @ Jp.T # Damped ps
15    Np = np.eye(param.nbVarX) - pinvJp @ Jp # Nullspace projection operator
16    up = pinvJp @ dfp # Command for position tracking
17    JoNp = Jo @ Np
18    pinvJoNp = JoNp.T @ np.linalg.inv(JoNp @ JoNp.T + np.eye(1) * 1e1) #
19    uo = pinvJoNp @ (dfo - Jo @ up) # Command for orientation tracking (
20    u = up + Np @ uo # Control commands
21
22    return u

```

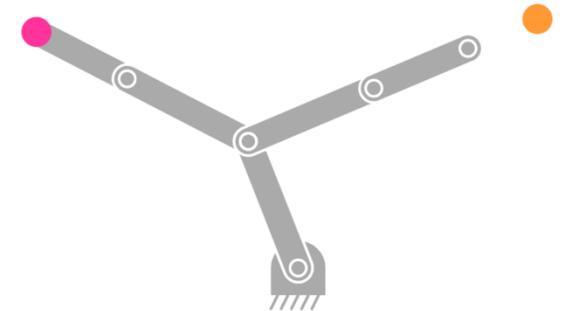
Object orientation



```

IK Prioritized IK
1 def controlCommand(x, param):
2     f = fkin(x, param)
3     J = Jkin(x, param)
4     # Prioritized control (left tracking as main objective)
5     df1 = (param.Mu[:2] - f[:2,0]) * 10 # Left hand correction
6     dfr = (param.Mu[2:] - f[2:,0]) * 10 # Right hand correction
7     J1 = J[:2,:] # Jacobian for left hand
8     Jr = J[2:,:] # Jacobian for right hand
9     pinvJ1 = np.linalg.inv(J1.T @ J1 + np.eye(param.nbVarX) * 1e1) @ J1.T # Damped pse
10    N1 = np.eye(param.nbVarX) - pinvJ1 @ J1 # Nullspace projection operator
11    u1 = pinvJ1 @ df1 # Command for position tracking
12    JrN1 = Jr @ N1
13    pinvJrN1 = JrN1.T @ np.linalg.inv(JrN1 @ JrN1.T + np.eye(2) * 1e4) # Damped pseudo
14    ur = pinvJrN1 @ (dfr - Jr @ u1) # Command for right hand tracking (with left hand
15    u = u1 + N1 @ ur # Control commands
16    return u

```



<https://robotics-codes-from-scratch.github.io/>

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

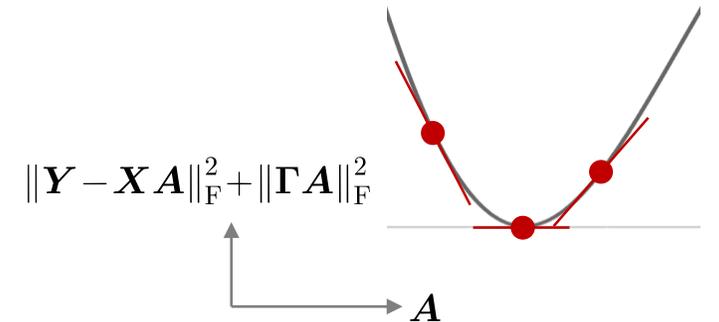
$$\begin{aligned}\hat{\mathbf{A}} &= \arg \min_{\mathbf{A}} \|\mathbf{Y} - \mathbf{X}\mathbf{A}\|_{\mathbb{F}}^2 + \|\mathbf{\Gamma}\mathbf{A}\|_{\mathbb{F}}^2 \\ &= \arg \min_{\mathbf{A}} \text{tr}\left((\mathbf{Y} - \mathbf{X}\mathbf{A})^\top(\mathbf{Y} - \mathbf{X}\mathbf{A})\right) + \text{tr}\left((\mathbf{\Gamma}\mathbf{A})^\top\mathbf{\Gamma}\mathbf{A}\ri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

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

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

yielding

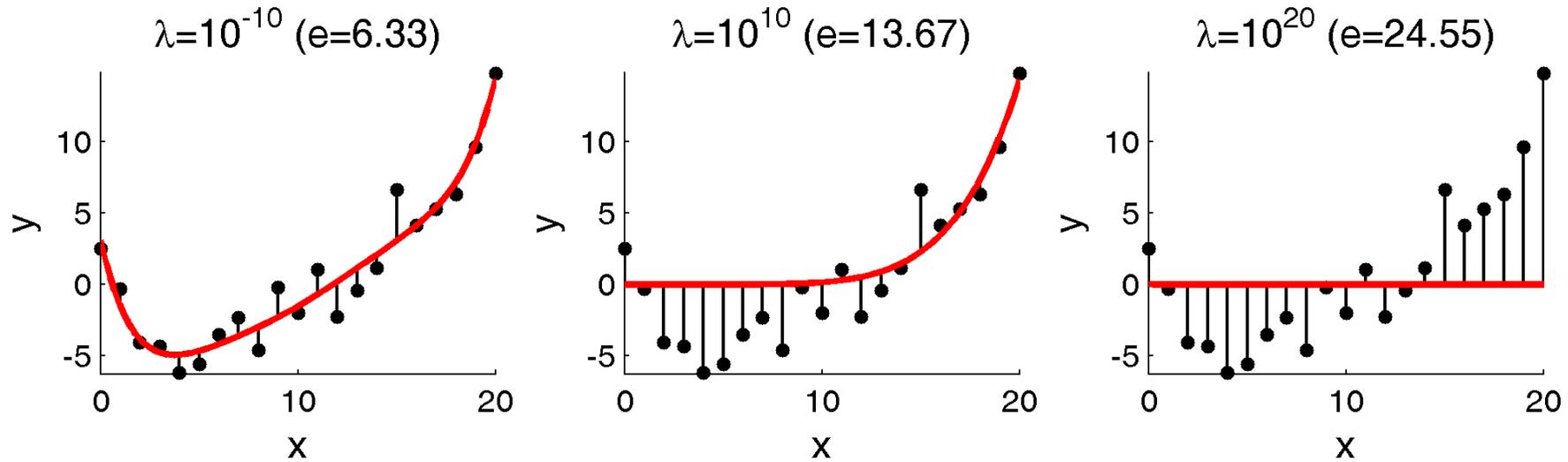
$$\hat{\mathbf{A}} = (\mathbf{X}^\top\mathbf{X} + \mathbf{\Gamma}^\top\mathbf{\Gamma})^{-1}\mathbf{X}^\top\mathbf{Y}$$



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

Example: Polynomial fitting

$D^{\mathcal{I}} = 7$ (polynomial of degree 6)



Ridge regression *computed with SVD*

For the singular value decomposition

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

with σ_i the singular values in the diagonal of $\mathbf{\Sigma}$, a solution to the ridge regression problem is given by

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

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

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

and has zeros elsewhere.

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



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

(for $\lambda = 0.01$)

Weighted least squares (Generalized least squares)

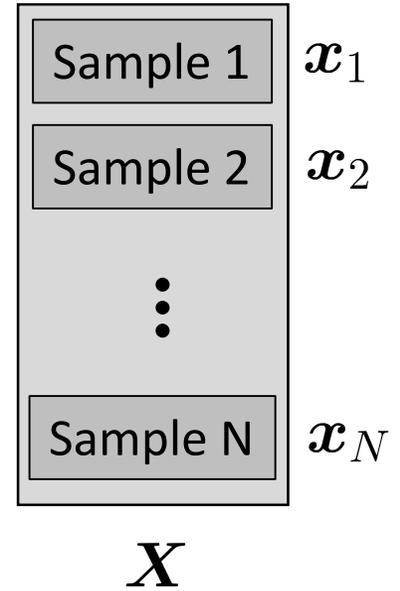
Python notebook:
demo_LS_weighted.ipynb

Matlab example:
demo_LS_weighted01.m

Weighted least squares

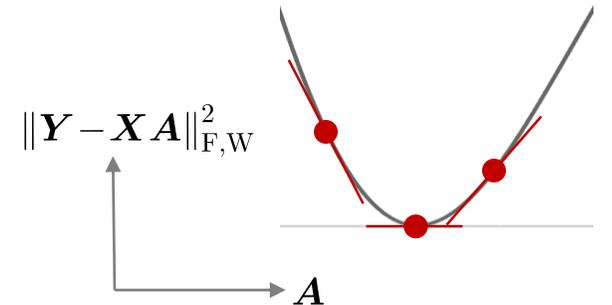
By describing the input data as $\mathbf{X} \in \mathbb{R}^{N \times D^I}$ and the output data as $\mathbf{Y} \in \mathbb{R}^{N \times D^O}$, with a weight matrix $\mathbf{W} \in \mathbb{R}^{N \times N}$, we want to minimize

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



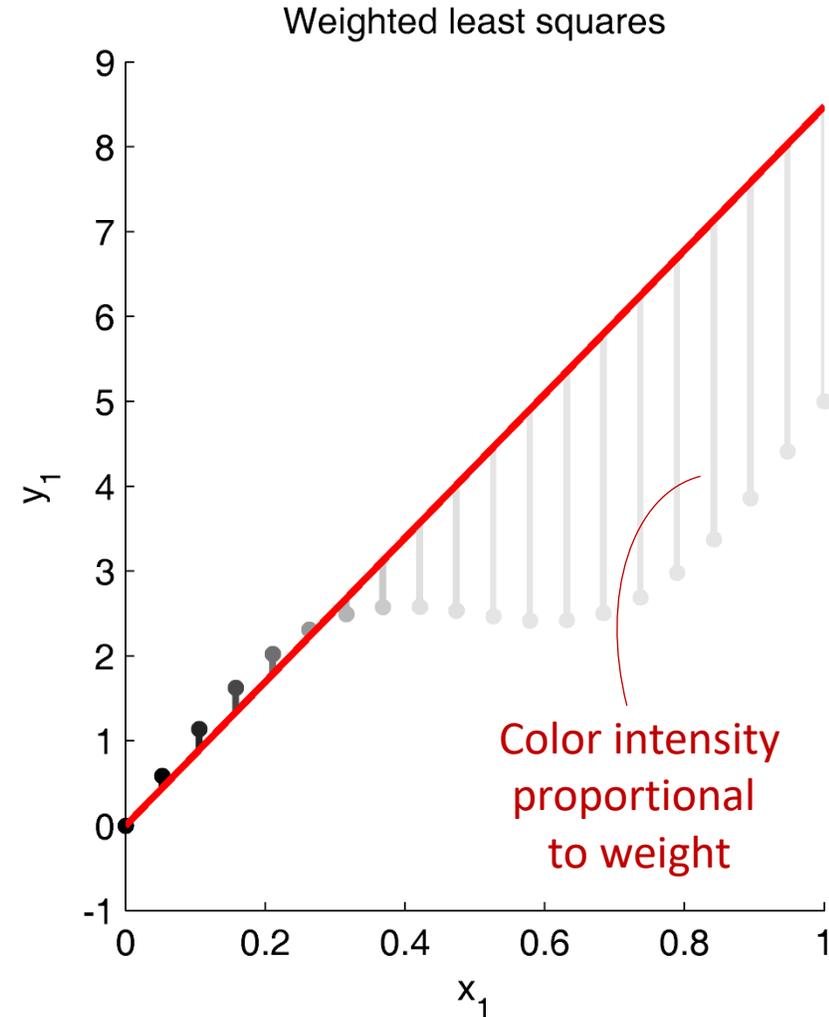
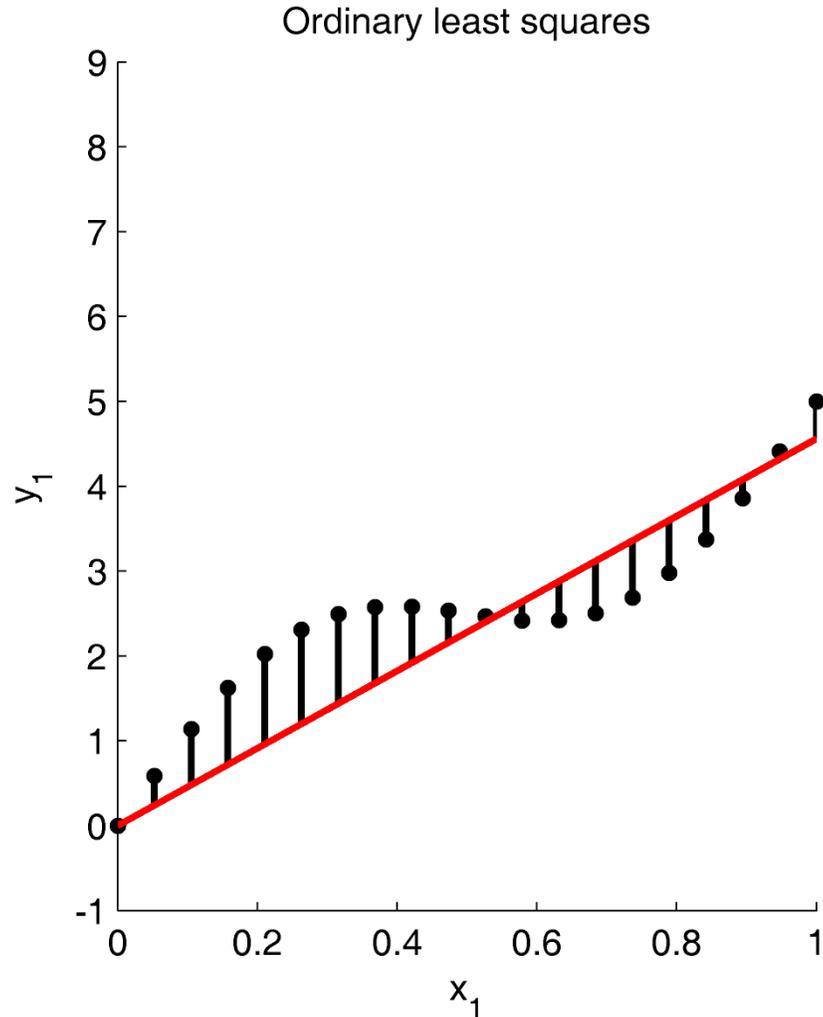
By differentiating with respect to \mathbf{A} and equating to zero

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



Weighted least squares

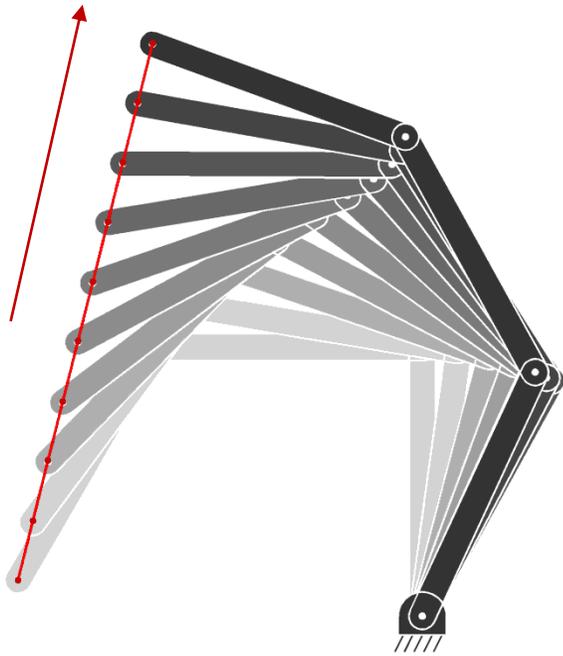
$$\hat{\mathbf{A}} = (\mathbf{X}^T \mathbf{W} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{W} \mathbf{Y}$$



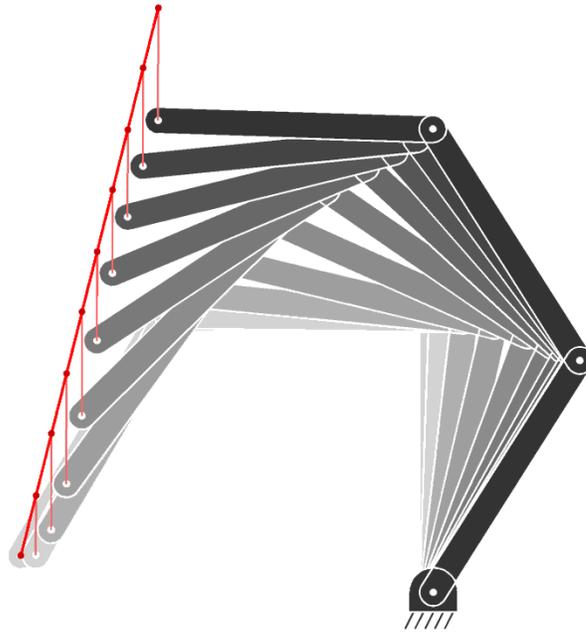
Example: Robot inverse kinematics (weights in task space)

$$\Delta x = (J^T W^F J + \lambda I_x)^{-1} J^T W^F \Delta f$$

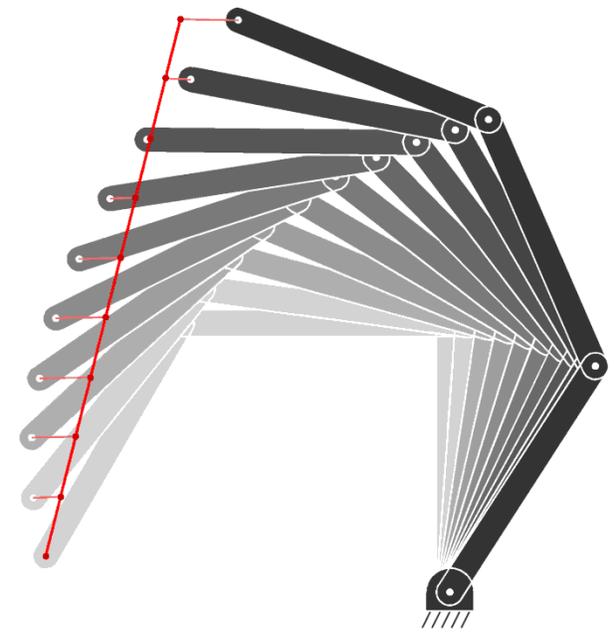
(moving target)



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



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



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

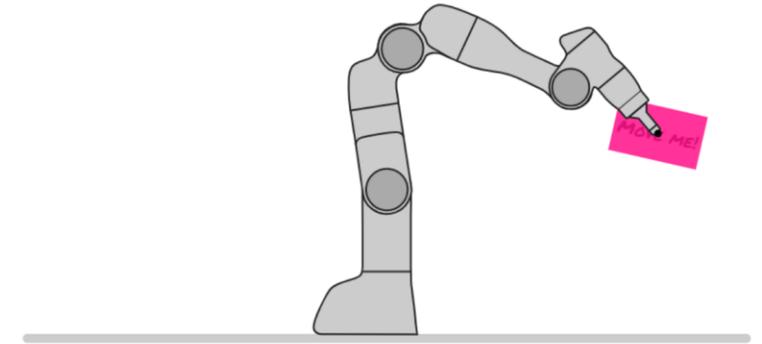
Example: Robot inverse kinematics (weights in task space)

IK Damped IK **Weighted IK** Prioritized IK

```

1 x = [-np.pi/4, np.pi/2, np.pi/4] # Initial robot state
2 def controlCommand(x, param):
3     f = fkin(x, param)
4     J = Jkin(x, param)
5
6     # Weights in task space
7     Wf = np.diag([1, 1, 0])
8     pinvWJ = np.linalg.inv(J.T @ Wf @ J + np.eye(param.nbVarX) * 1E-2) @ J.T @ Wf # We
9     u = pinvWJ @ logmap(param.Mu, f) # Position & orientation tracking
10
11 # # Weights in configuration space
12 # Wx = np.diag([0.01, 1, 1])
13 # pinvWJ = Wx @ J[:2,:].T @ np.linalg.inv(J[:2,:] @ Wx @ J[:2,:].T + np.eye(2) * 1E-2)
14 # u = pinvWJ @ (param.Mu[:2] - f[:2]) # Position tracking
15
16 return u

```



Object orientation

<https://robotics-codes-from-scratch.github.io/>

Recursive least squares

Python notebook:

`demo_LS_recursive.ipynb`

Matlab code:

`demo_LS_recursive01.m`

Recursive least squares

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

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



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

Sherman-Morrison-Woodbury relation:

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

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 $\mathbf{B} + \mathbf{UV}$.

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

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

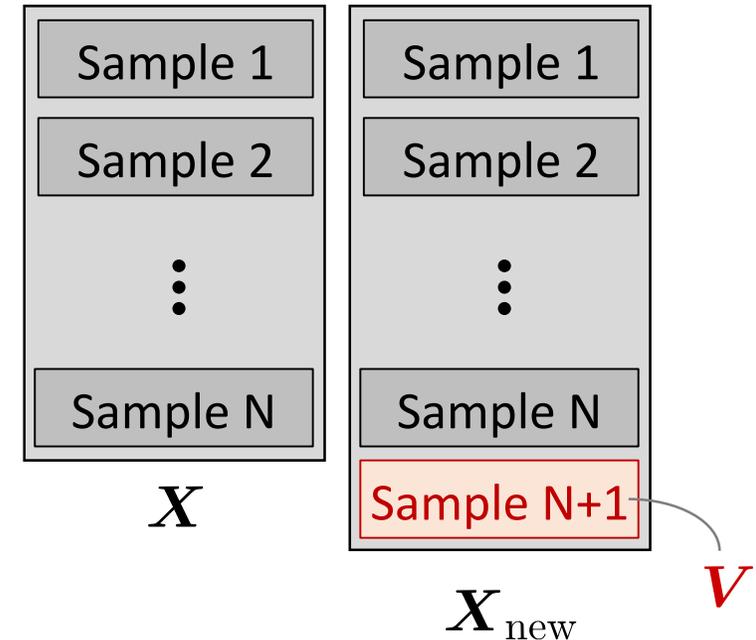
$$\mathbf{B}_{\text{new}}^{-1} = \mathbf{B}^{-1} - \mathbf{B}^{-1} \mathbf{V}^\top (\mathbf{I} + \mathbf{V} \mathbf{B}^{-1} \mathbf{V}^\top)^{-1} \mathbf{V} \mathbf{B}^{-1}$$

which is exploited to efficiently compute the least squares update as

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

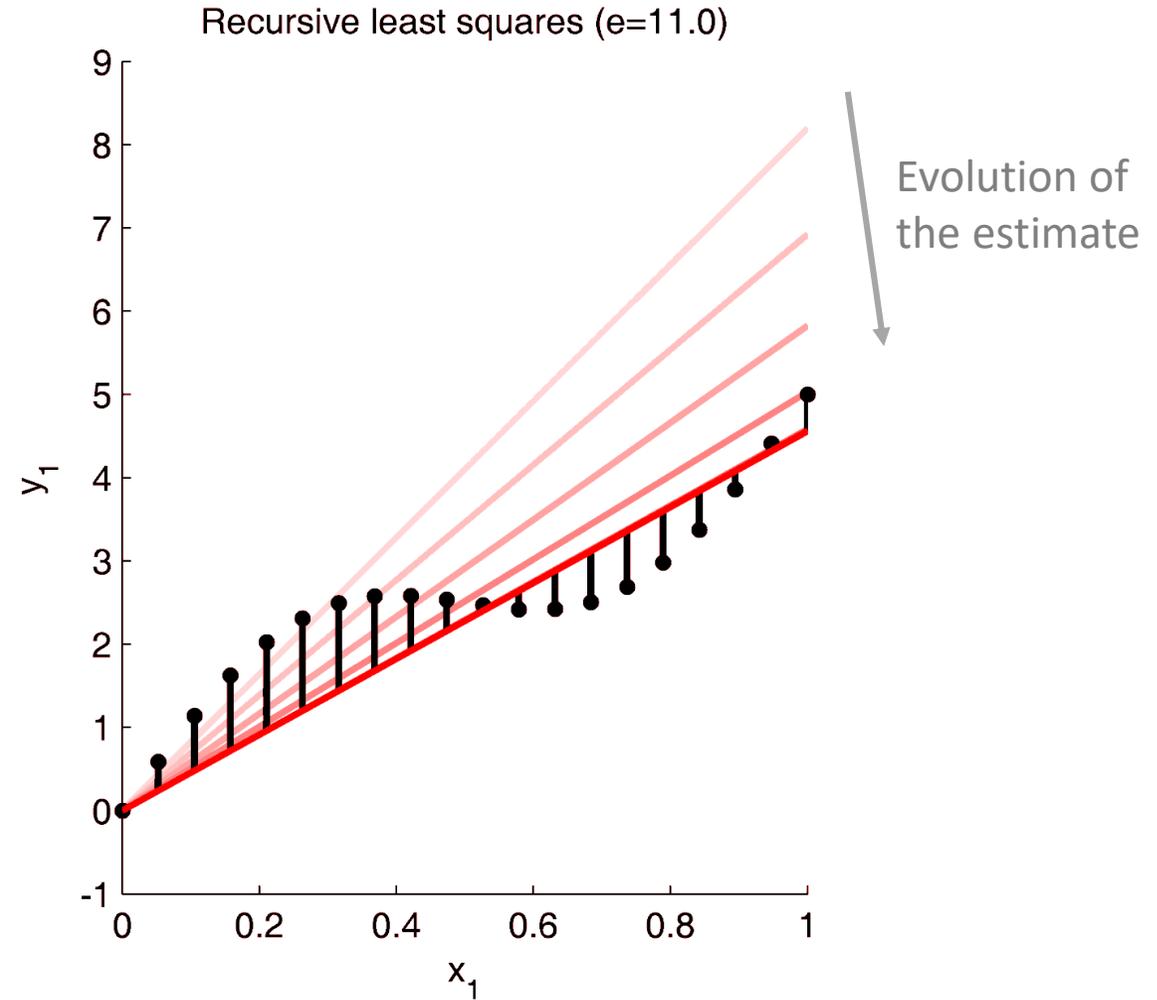
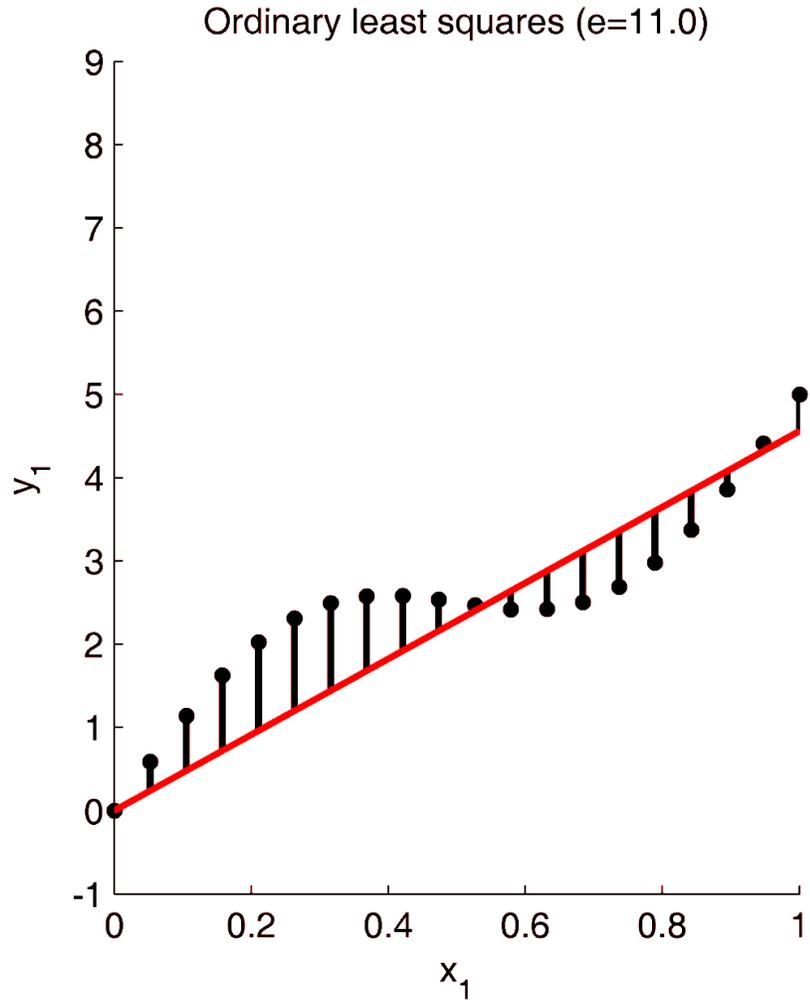
with Kalman gain $\mathbf{K} = \mathbf{B}^{-1} \mathbf{V}^\top (\mathbf{I} + \mathbf{V} \mathbf{B}^{-1} \mathbf{V}^\top)^{-1}$

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



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

Recursive least squares



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

Linear regression:

A more elaborated example

(but still only least squares!)

Linear quadratic regulator (LQR)

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

Track path! Use low control commands!

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

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

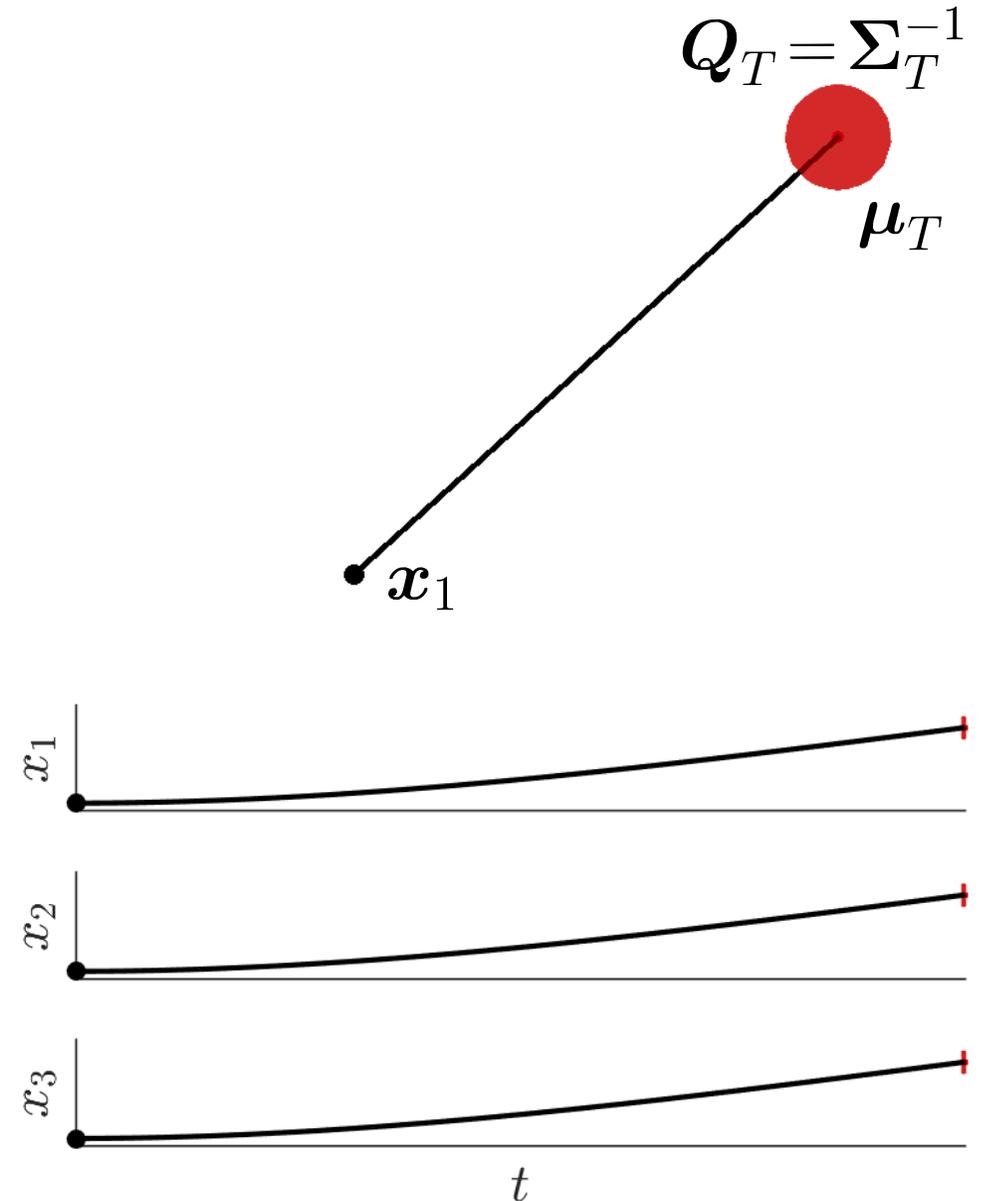
$\boldsymbol{\mu}_t$ desired state

\boldsymbol{u}_t control command (acceleration)

\boldsymbol{Q}_t precision matrix

\boldsymbol{R}_t control weight matrix

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



Linear quadratic regulator (LQR)

$$\min_{\mathbf{u}} \sum_{t=1}^T \left(\|\boldsymbol{\mu}_t - \mathbf{x}_t\|_{\mathbf{Q}_t}^2 + \|\mathbf{u}_t\|_{\mathbf{R}_t}^2 \right)$$

Track path! Use low control commands!

$$\text{s.t. } \mathbf{x}_{t+1} = \mathbf{A}\mathbf{x}_t + \mathbf{B}\mathbf{u}_t \quad \text{System dynamics}$$

**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 regulator (LQR)

$$c = \sum_{t=1}^T \left((\boldsymbol{\mu}_t - \boldsymbol{x}_t)^\top \boldsymbol{Q}_t (\boldsymbol{\mu}_t - \boldsymbol{x}_t) + \boldsymbol{u}_t^\top \boldsymbol{R}_t \boldsymbol{u}_t \right)$$

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

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

$$\boldsymbol{Q} = \begin{bmatrix} \boldsymbol{Q}_1 & \mathbf{0} & \cdots & \mathbf{0} \\ \mathbf{0} & \boldsymbol{Q}_2 & \cdots & \mathbf{0} \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{0} & \mathbf{0} & \cdots & \boldsymbol{Q}_T \end{bmatrix}$$

$$\boldsymbol{R} = \begin{bmatrix} \boldsymbol{R}_1 & \mathbf{0} & \cdots & \mathbf{0} \\ \mathbf{0} & \boldsymbol{R}_2 & \cdots & \mathbf{0} \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{0} & \mathbf{0} & \cdots & \boldsymbol{R}_T \end{bmatrix}$$

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

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



Linear quadratic regulator (LQR)

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

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

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

$$\vdots$$

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

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

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



Linear quadratic regulator (LQR)

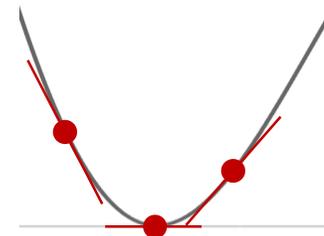
The constraint can then be inserted in the cost function:

$$\begin{aligned}
 c &= (\boldsymbol{\mu} - \boldsymbol{x})^\top \boldsymbol{Q} (\boldsymbol{\mu} - \boldsymbol{x}) + \boldsymbol{u}^\top \boldsymbol{R} \boldsymbol{u} \\
 &= (\boldsymbol{\mu} - \boldsymbol{S}^x \boldsymbol{x}_1 - \boldsymbol{S}^u \boldsymbol{u})^\top \boldsymbol{Q} (\boldsymbol{\mu} - \boldsymbol{S}^x \boldsymbol{x}_1 - \boldsymbol{S}^u \boldsymbol{u}) + \boldsymbol{u}^\top \boldsymbol{R} \boldsymbol{u}
 \end{aligned}$$

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

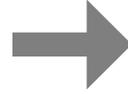
$$\hat{\boldsymbol{u}} = (\boldsymbol{S}^{u^\top} \boldsymbol{Q} \boldsymbol{S}^u + \boldsymbol{R})^{-1} \boldsymbol{S}^{u^\top} \boldsymbol{Q} (\boldsymbol{\mu} - \boldsymbol{S}^x \boldsymbol{x}_1)$$

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



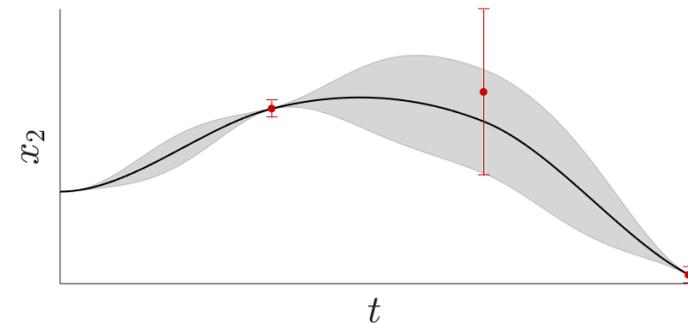
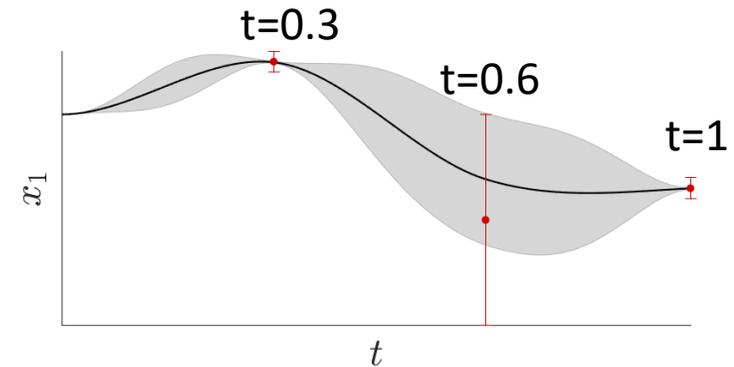
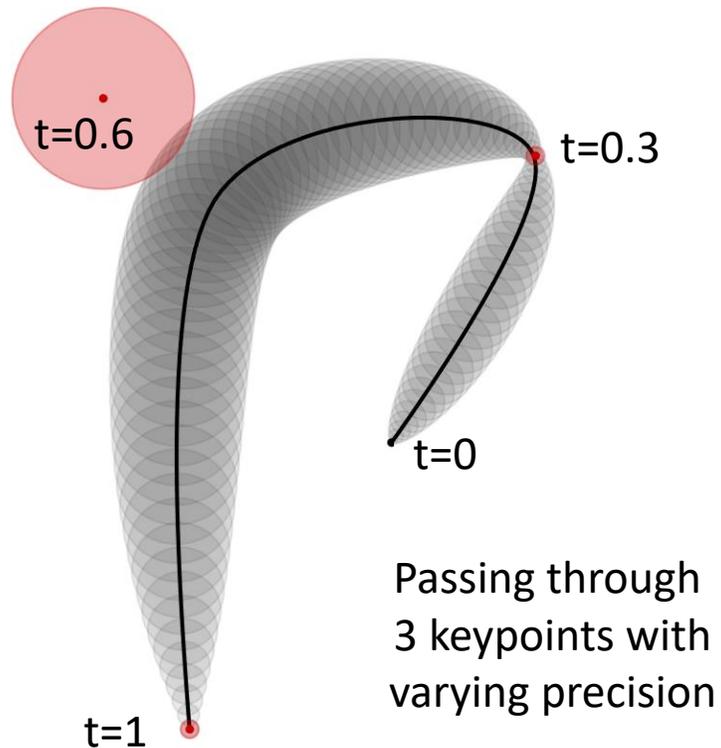
Linear quadratic regulator (LQR)

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

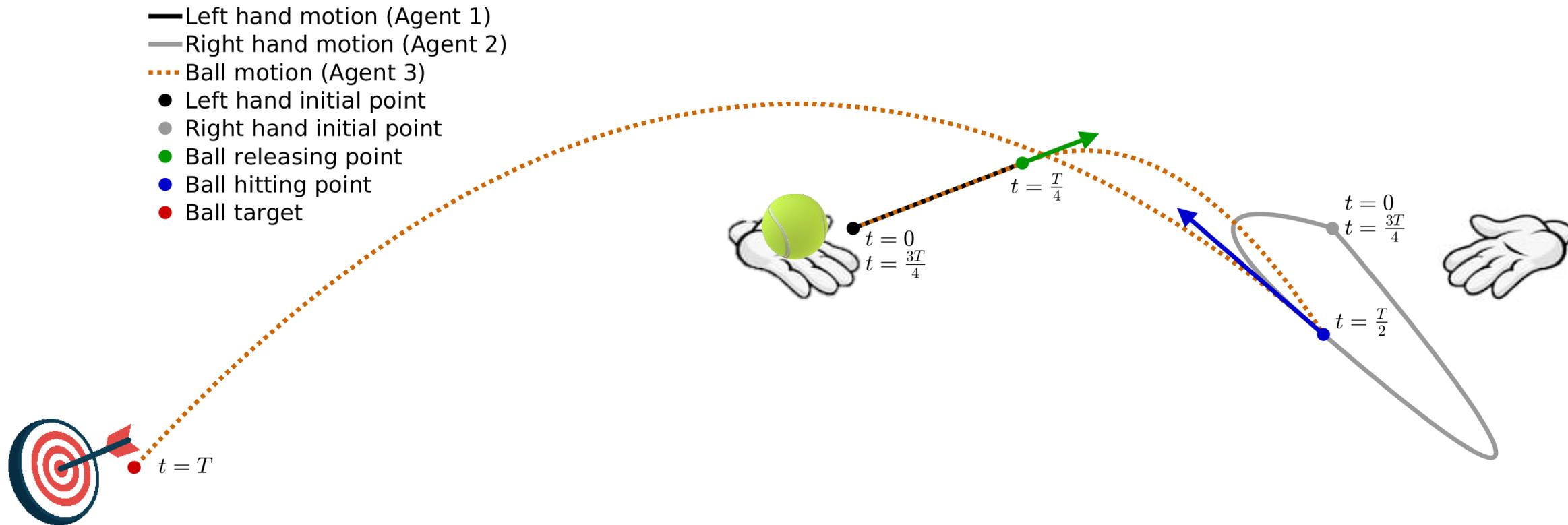
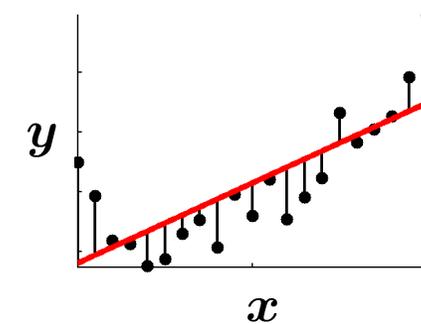


$$\hat{x} = S^x x_1 + S^u \hat{u}$$

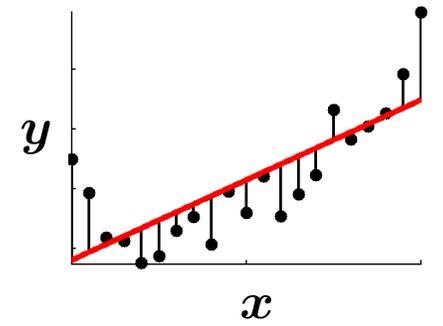
The control trajectories can then be converted to state trajectories



Linear quadratic regulator (LQR) - Example



Linear quadratic regulator (LQR) - Example



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

$$\begin{array}{c} \text{Hand} \\ \text{Hand} \\ \text{Ball} \end{array} \begin{bmatrix} \dot{x}_{1,t} \\ \ddot{x}_{1,t} \\ \dot{f}_{1,t} \\ \dot{x}_{2,t} \\ \ddot{x}_{2,t} \\ \dot{f}_{2,t} \\ \dot{x}_{3,t} \\ \ddot{x}_{3,t} \\ \dot{f}_{3,t} \end{bmatrix} = \underbrace{\begin{bmatrix} 0 & I & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & I & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ \hline 0 & 0 & 0 & 0 & I & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & I & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ \hline 0 & I & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & I \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \end{bmatrix}}_{A_t^c} \underbrace{\begin{bmatrix} x_{1,t} \\ \dot{x}_{1,t} \\ f_{1,t} \\ x_{2,t} \\ \dot{x}_{2,t} \\ f_{2,t} \\ x_{3,t} \\ \dot{x}_{3,t} \\ f_{3,t} \end{bmatrix}}_{x_t} + \underbrace{\begin{bmatrix} 0 & 0 \\ I & 0 \\ 0 & 0 \\ 0 & 0 \\ 0 & I \\ 0 & 0 \\ 0 & 0 \\ 0 & 0 \\ 0 & 0 \end{bmatrix}}_{B_t^c} \underbrace{\begin{bmatrix} u_{1,t} \\ u_{2,t} \end{bmatrix}}_{u_t}$$

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

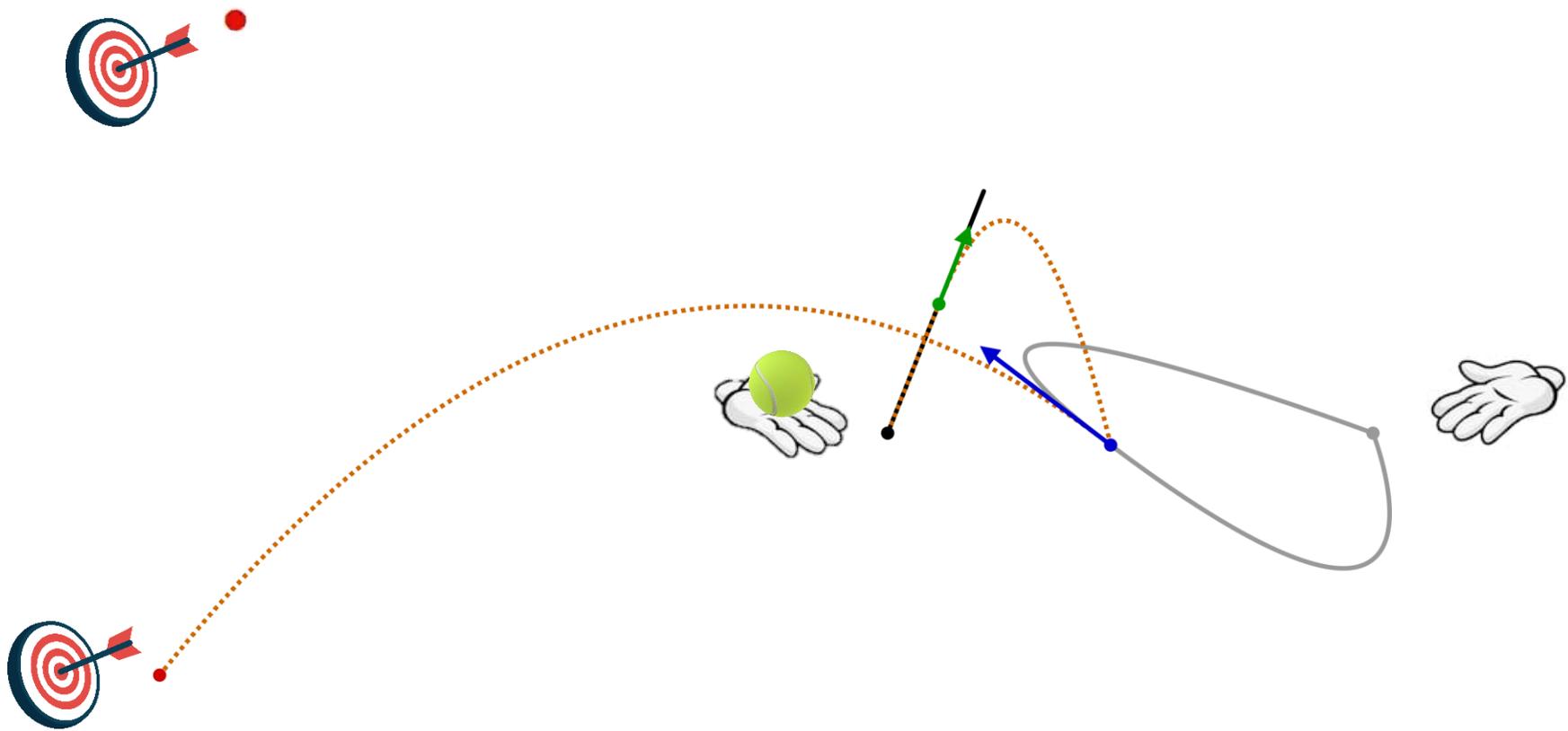
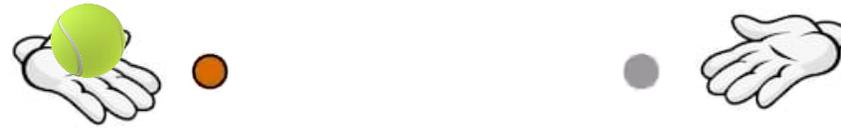
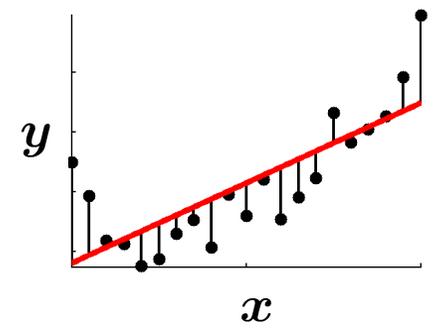
At $t = \frac{T}{2}$ (right hand hitting the ball), we have

$$\begin{array}{c} \text{Hand} \\ \text{Hand} \\ \text{Ball} \end{array} \begin{bmatrix} \dot{x}_{1,t} \\ \ddot{x}_{1,t} \\ \dot{f}_{1,t} \\ \dot{x}_{2,t} \\ \ddot{x}_{2,t} \\ \dot{f}_{2,t} \\ \dot{x}_{3,t} \\ \ddot{x}_{3,t} \\ \dot{f}_{3,t} \end{bmatrix} = \underbrace{\begin{bmatrix} 0 & I & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & I & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ \hline 0 & 0 & 0 & 0 & I & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & I & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ \hline 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & I \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \end{bmatrix}}_{A_t^c} \underbrace{\begin{bmatrix} x_{1,t} \\ \dot{x}_{1,t} \\ f_{1,t} \\ x_{2,t} \\ \dot{x}_{2,t} \\ f_{2,t} \\ x_{3,t} \\ \dot{x}_{3,t} \\ f_{3,t} \end{bmatrix}}_{x_t} + \underbrace{\begin{bmatrix} 0 & 0 \\ I & 0 \\ 0 & 0 \\ 0 & 0 \\ 0 & I \\ 0 & 0 \\ 0 & 0 \\ 0 & 0 \\ 0 & 0 \end{bmatrix}}_{B_t^c} \underbrace{\begin{bmatrix} u_{1,t} \\ u_{2,t} \end{bmatrix}}_{u_t}$$

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

$$\begin{array}{c} \text{Hand} \\ \text{Hand} \\ \text{Ball} \end{array} \begin{bmatrix} \dot{x}_{1,t} \\ \ddot{x}_{1,t} \\ \dot{f}_{1,t} \\ \dot{x}_{2,t} \\ \ddot{x}_{2,t} \\ \dot{f}_{2,t} \\ \dot{x}_{3,t} \\ \ddot{x}_{3,t} \\ \dot{f}_{3,t} \end{bmatrix} = \underbrace{\begin{bmatrix} 0 & I & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & I & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ \hline 0 & 0 & 0 & 0 & I & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & I & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ \hline 0 & 0 & 0 & 0 & 0 & 0 & 0 & I & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & I \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \end{bmatrix}}_{A_t^c} \underbrace{\begin{bmatrix} x_{1,t} \\ \dot{x}_{1,t} \\ f_{1,t} \\ x_{2,t} \\ \dot{x}_{2,t} \\ f_{2,t} \\ x_{3,t} \\ \dot{x}_{3,t} \\ f_{3,t} \end{bmatrix}}_{x_t} + \underbrace{\begin{bmatrix} 0 & 0 \\ I & 0 \\ 0 & 0 \\ 0 & 0 \\ 0 & I \\ 0 & 0 \\ 0 & 0 \\ 0 & 0 \\ 0 & 0 \end{bmatrix}}_{B_t^c} \underbrace{\begin{bmatrix} u_{1,t} \\ u_{2,t} \end{bmatrix}}_{u_t}$$

Linear quadratic regulator (LQR) - Example



Recap: Costs functions and associated solutions

Univariate output \mathbf{y} :

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

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

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

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

Multivariate output \mathbf{Y} :

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

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

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

Iteratively reweighted least squares (IRLS)

Python notebook:
demo_LS_weighted.ipynb

Matlab code:
demo_LS_IRLS01.m

Iteratively reweighted least squares (IRLS)

$$\|e\|_p = \left(\sum_{n=1}^N |e_n|^p \right)^{1/p}$$

- IRLS is useful to minimize ℓ_p norms with $\arg \min \|e\|_p = \arg \min \sum_{n=1}^N |e_n|^p$
- The strategy of IRLS is that $|e_n|^p$ can be rewritten as $|e_n|^{p-2} e_n^2$
- $|e_n|^{p-2}$ can be interpreted as a weight, which is used to minimize e_n^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)

$$|e_n|^p = \underbrace{|e_n|^{p-2}}_{\text{transformed as weight } \mathbf{W}} e_n^2$$

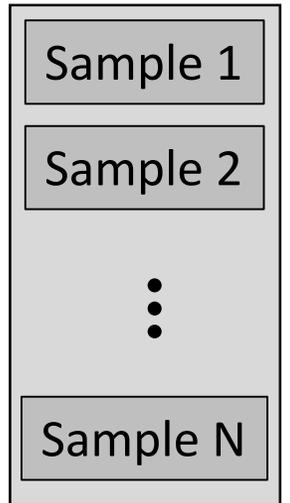
For an ℓ_p norm cost function defined by

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

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

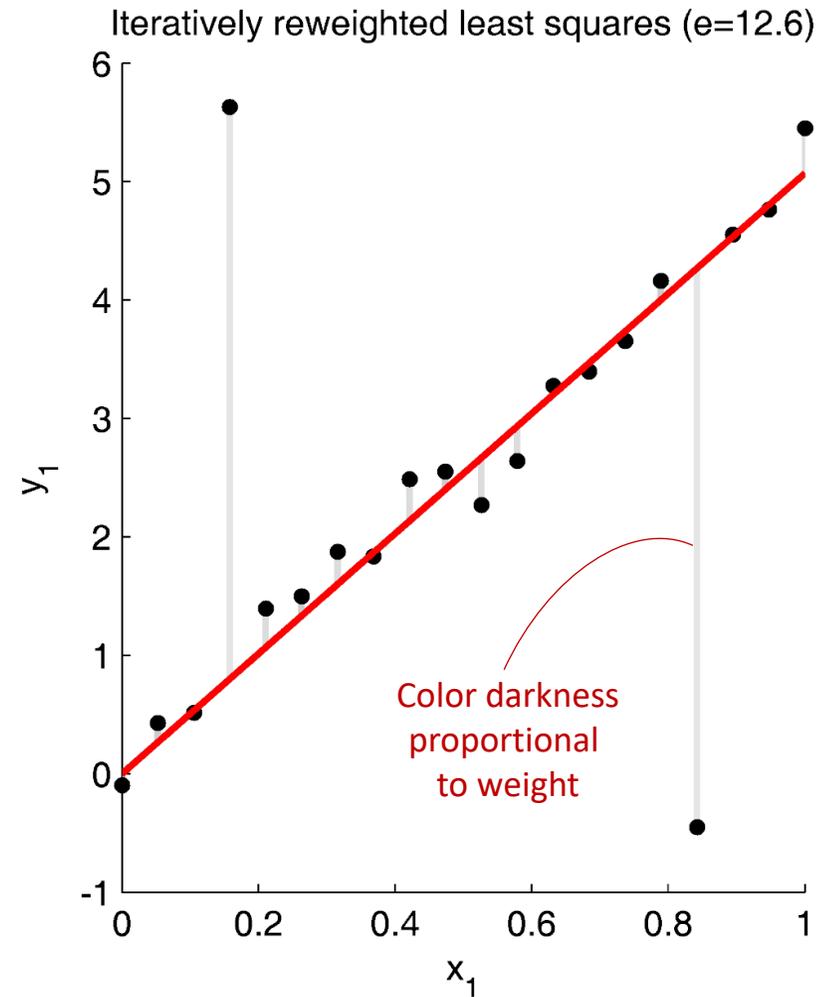
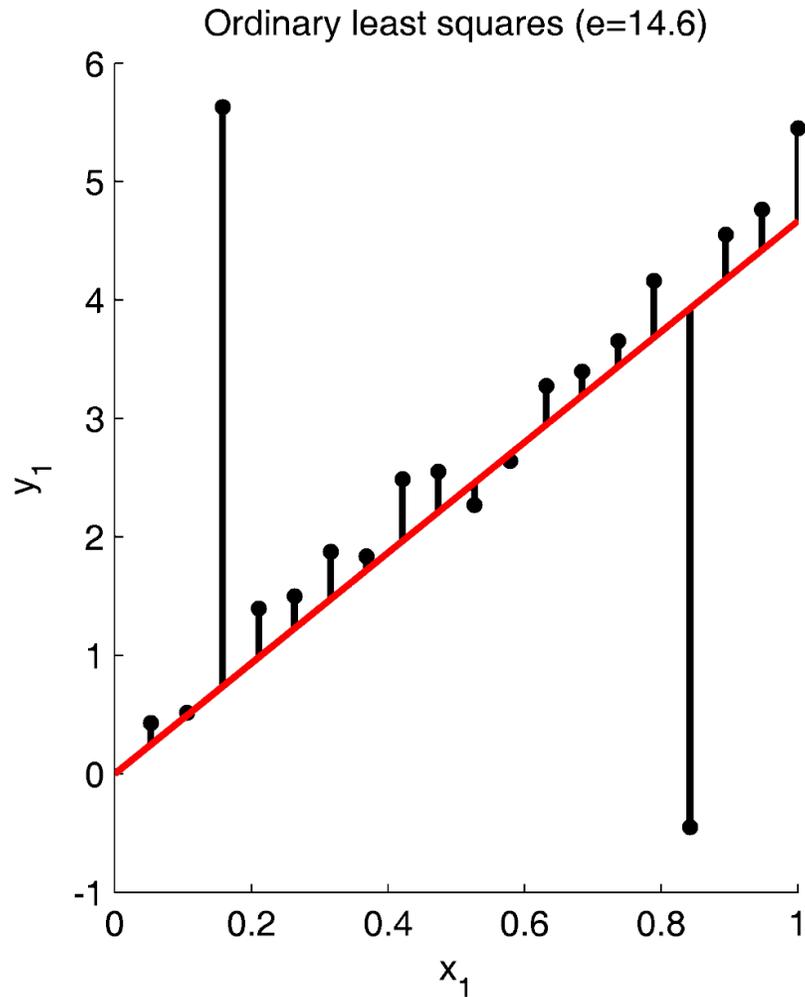
$$\hat{\mathbf{A}} \leftarrow (\mathbf{X}^\top \mathbf{W} \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{W} \mathbf{Y} \quad \hat{\mathbf{A}} = \arg \min_{\mathbf{A}} \|\mathbf{Y} - \mathbf{X}\mathbf{A}\|_{\mathbf{F}, \mathbf{W}}^2$$

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



\mathbf{X}

Iteratively reweighted least squares (IRLS)



→ regression that can sometimes be more robust to outliers

Logistic regression

Python notebook:

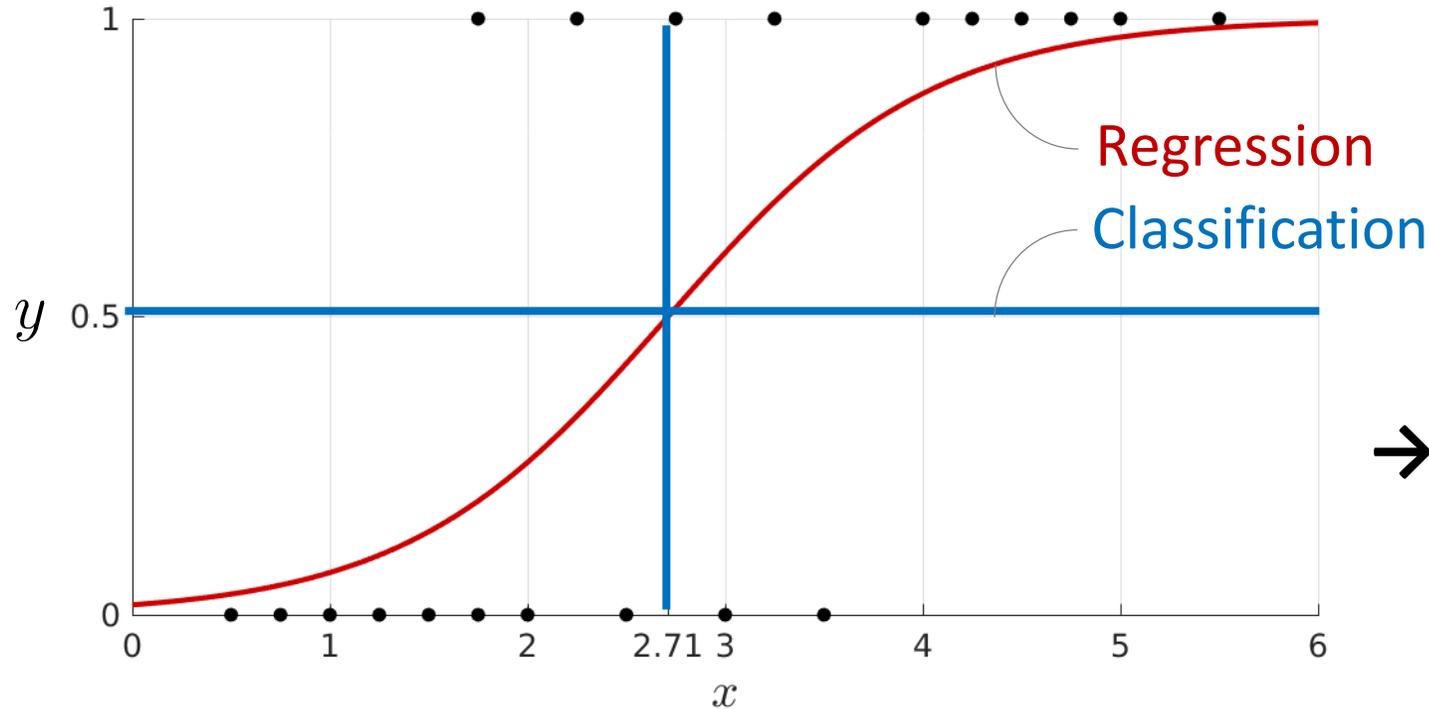
`demo_LS_IRLS_logRegr.ipynb`

Matlab code:

`demo_LS_IRLS_logRegr01.m`

Logistic regression

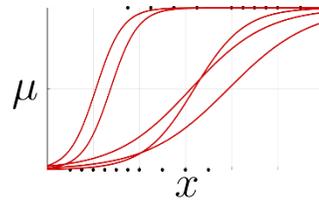
Example: Pass/fail in function of the time spent to study at an exam:



→ **Regression** exploited for **classification** problem

Logistic function:

$$\mu_{\mathbf{a}}(\mathbf{x}) = \frac{1}{1 + e^{-\mathbf{a}^T \mathbf{x}}} \quad \mu_{\mathbf{a}}(x) = \frac{1}{1 + e^{-(a_1 + a_2 x)}}$$



Logistic regression

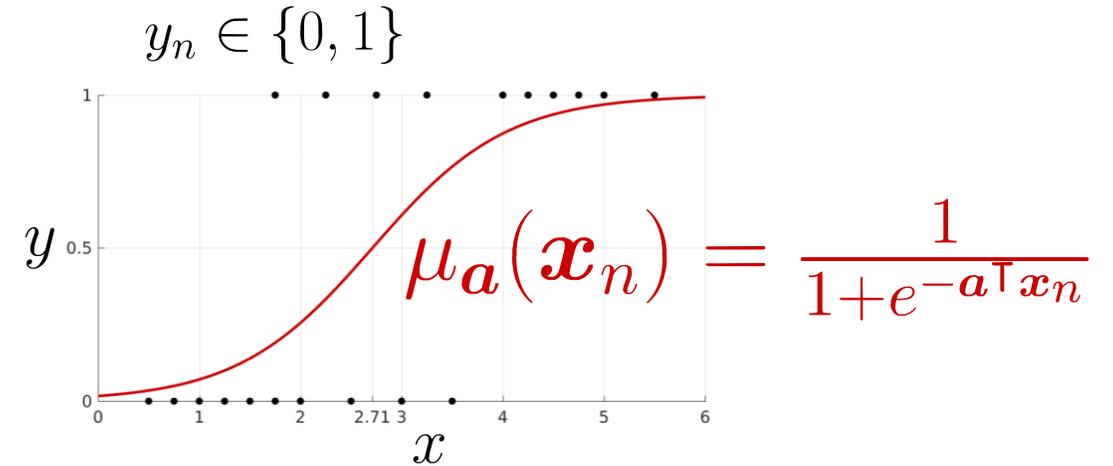
Bernoulli distribution (for binary variables):

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

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

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

Logistic function:



Likelihood of n^{th} datapoint:

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

Likelihood of N datapoints (independence assumption):

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

Logistic regression

Likelihood of N datapoints:

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

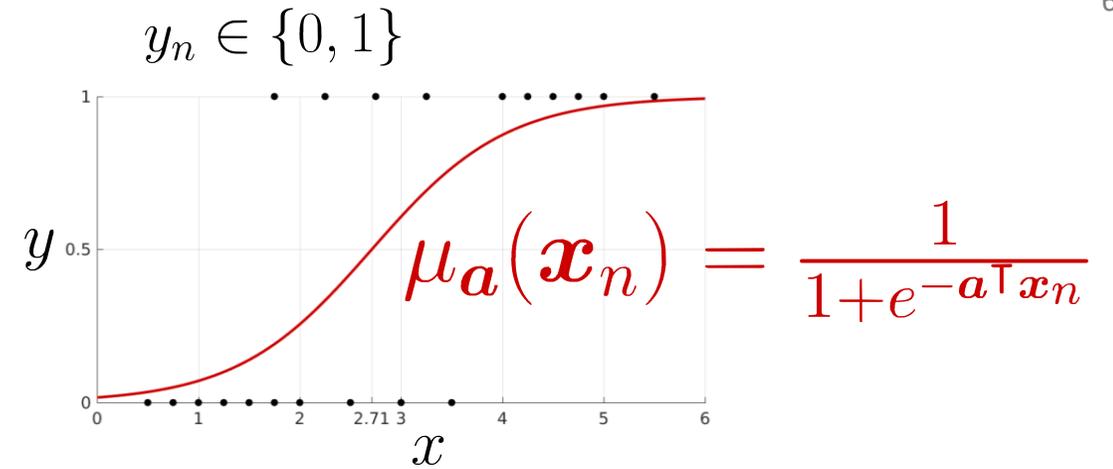
Cost function as negative log-likelihood:

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

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

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

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



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

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

Logistic regression

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

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

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

Hadamard (elementwise) product

We then obtain

$$\begin{aligned} \mathbf{a} &\leftarrow \mathbf{a} - \mathbf{H}^{-1} \mathbf{g} \\ &\leftarrow \mathbf{a} - (\mathbf{X}^\top \mathbf{W} \mathbf{X})^{-1} \mathbf{X}^\top (\boldsymbol{\mu}_{\mathbf{a}} - \mathbf{y}) \end{aligned}$$

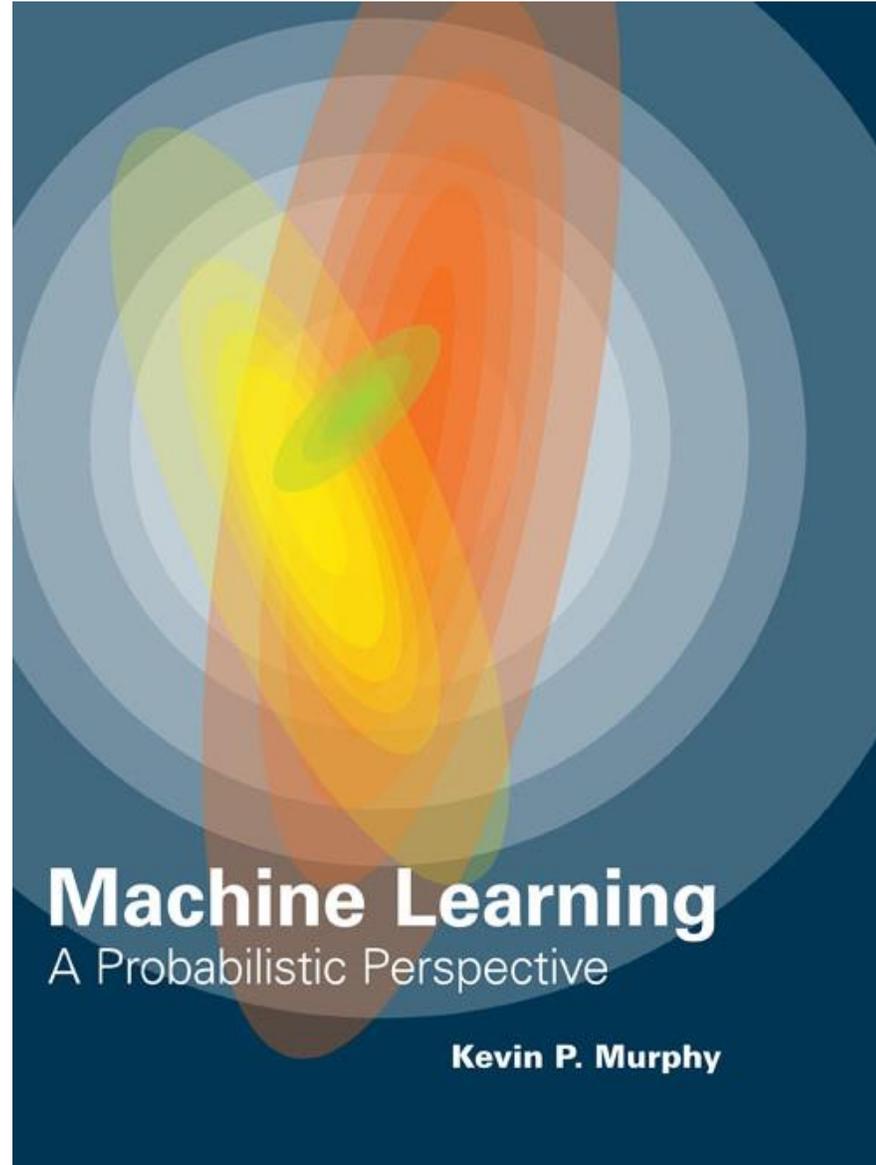
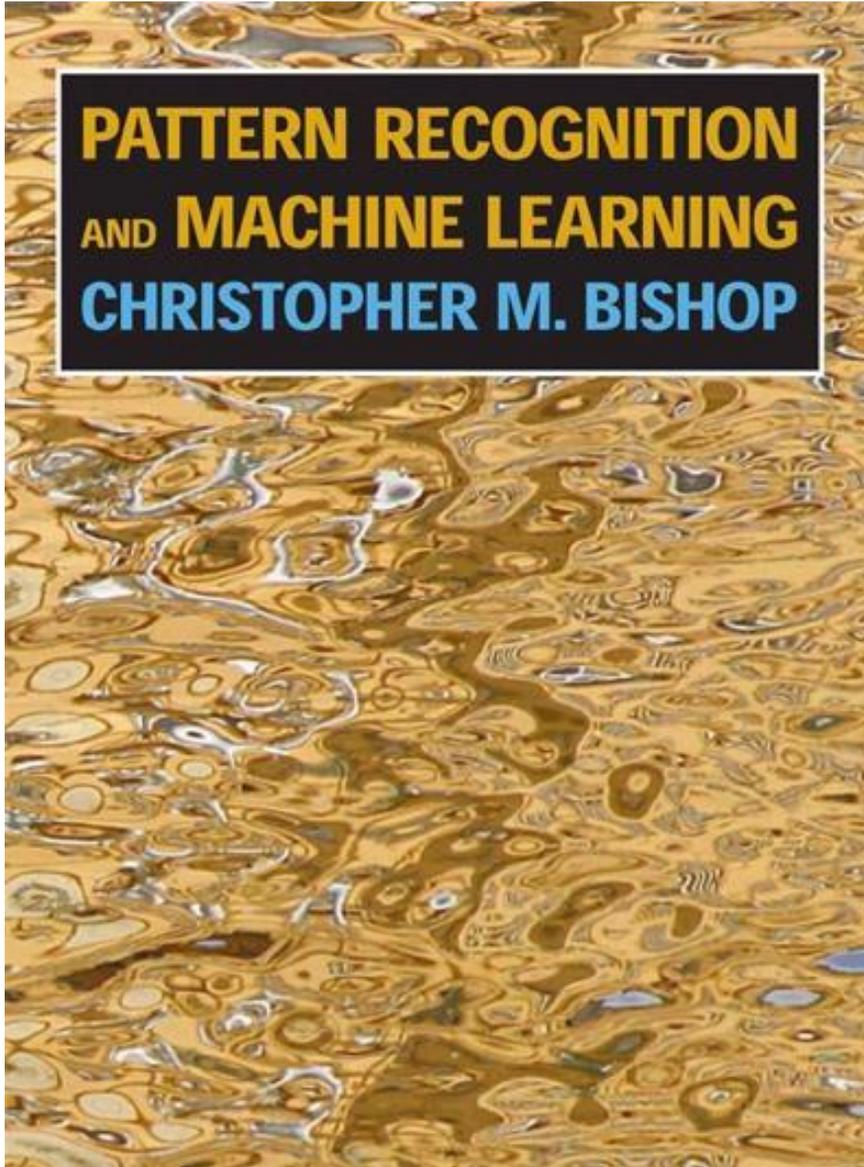
$$\frac{\partial c}{\partial \mathbf{a}} = \sum_n (\mu_{\mathbf{a}} - y_n) \mathbf{x}_n$$

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

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

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

General references



The Matrix Cookbook

Kaare Brandt Petersen
Michael Syskind Pedersen