

# Optimal and Learning Control for Autonomous Robots Lecture 10



**A D R L**

Jonas Buchli  
Agile & Dexterous Robotics Lab



# Class logistics

Have you signed up for the Interview for Ex 2?

<http://doodle.com/w2ahzwpdrwa5p5a9>  
(using your group ID!)



# Lecture 10 Goals

- ★ Function approximation, basis functions
- ★ Path integral stochastic optimal control
- ★ Path integral RL



# (Back to) Continuous state action spaces

Function approximation



# Mountain Car Problem

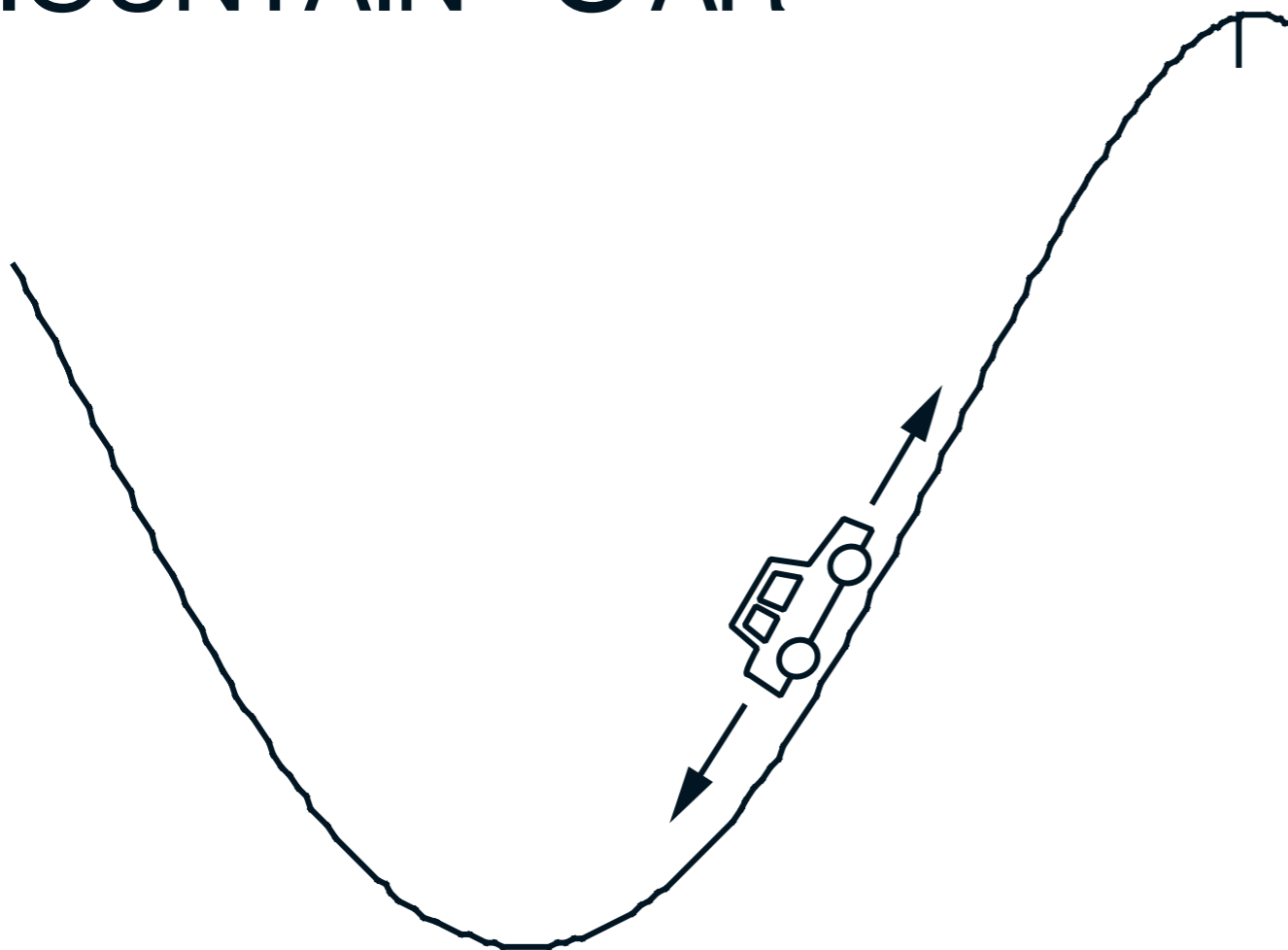
A continuous-state problem

MOUNTAIN CAR

Goal

Reward

- Goal: +10
- Step: -1



# Solving Mountain Car

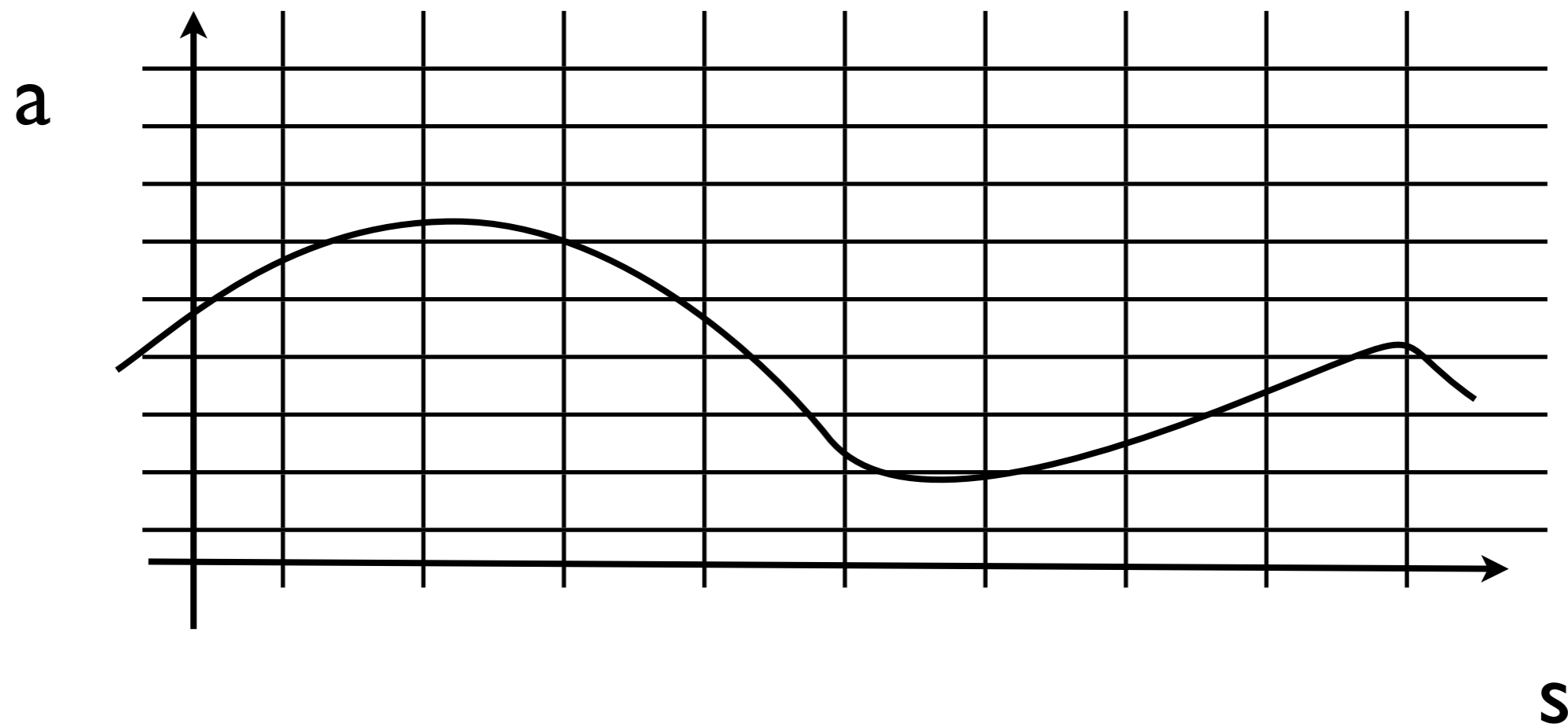
Practically continuous problems are (almost?)  
always solved by some sort of discretization!!!

Let's look at the problem of discretization in a  
bit more detail!



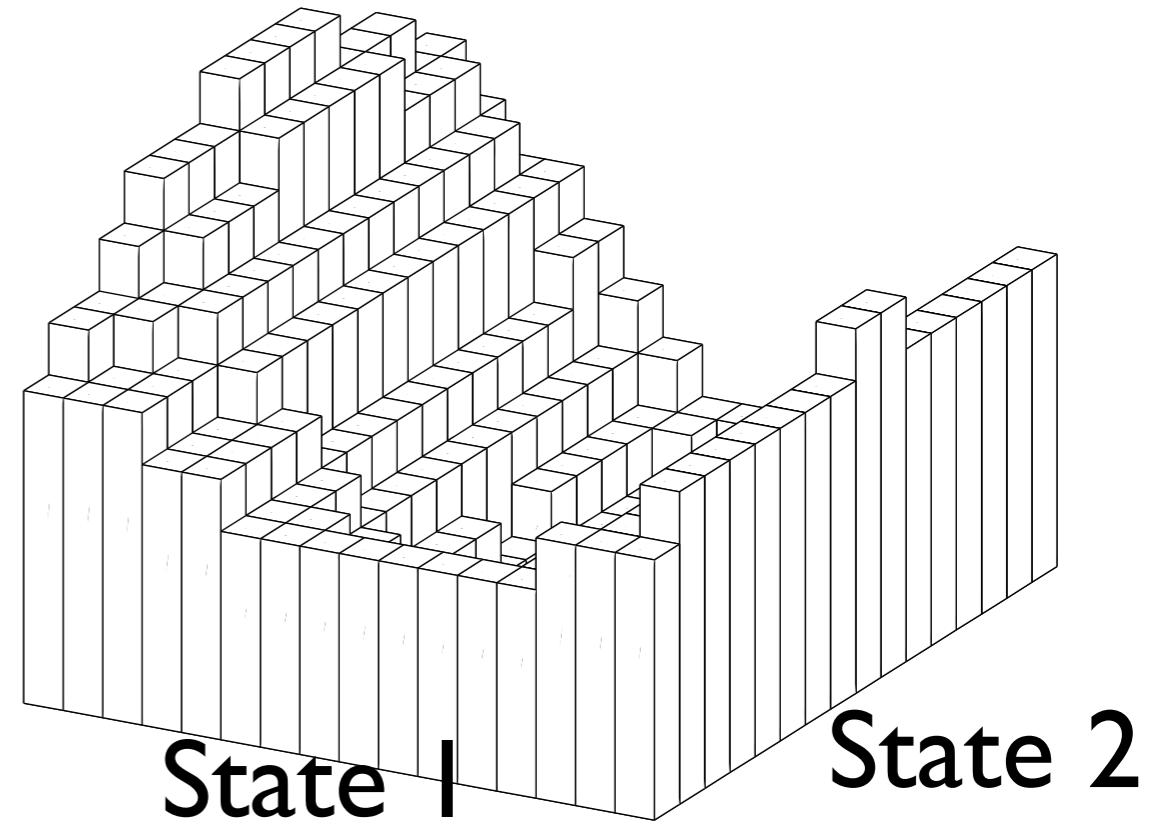
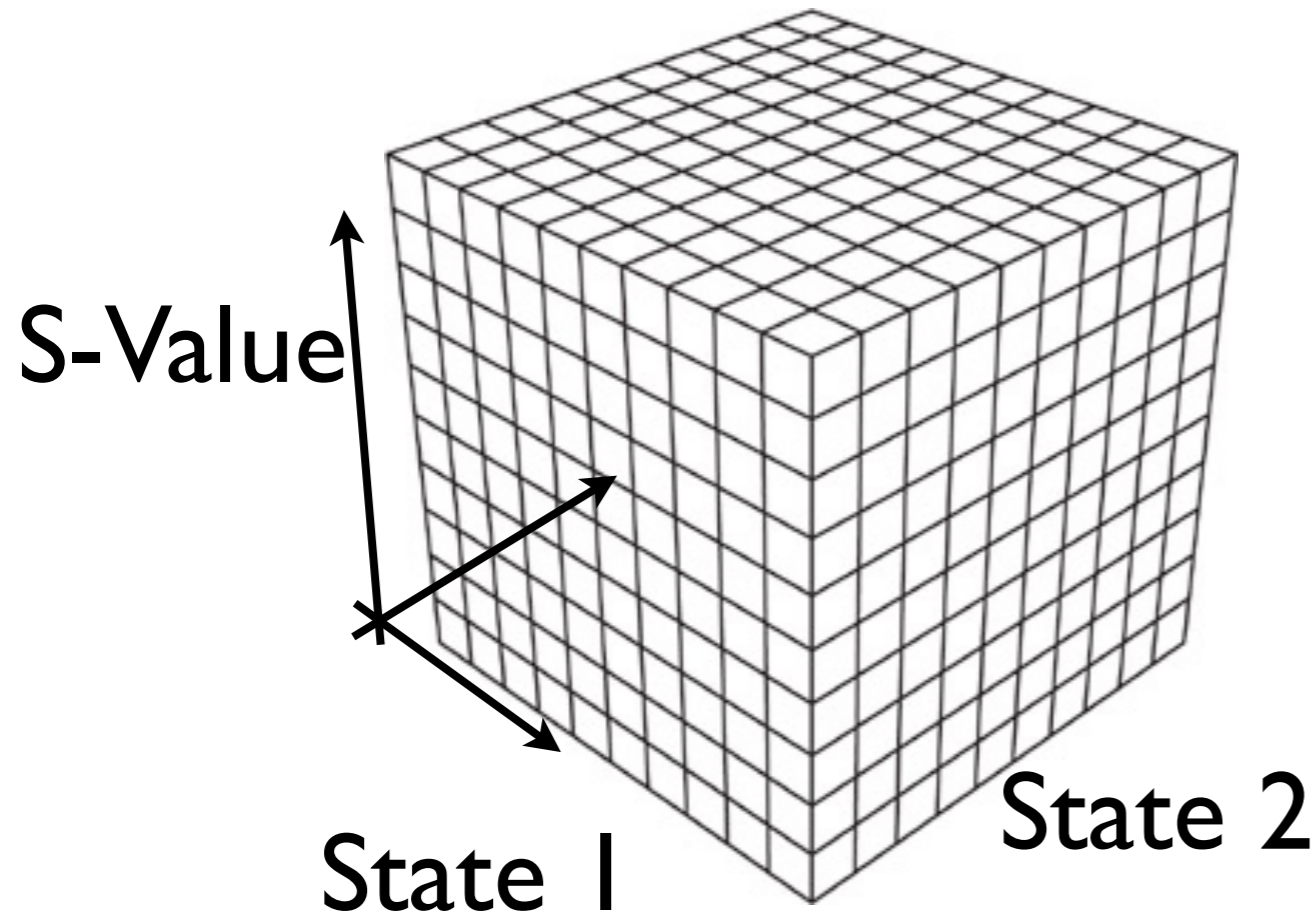
# Control Policies

state - action mapping

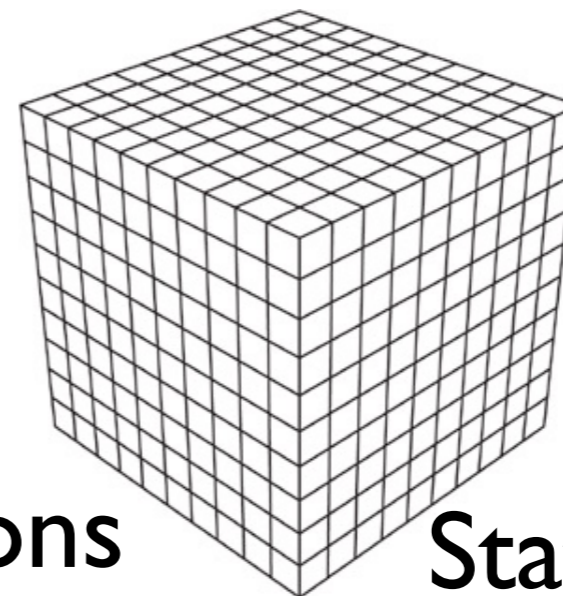


Problem: dimensions!

# Value discretization

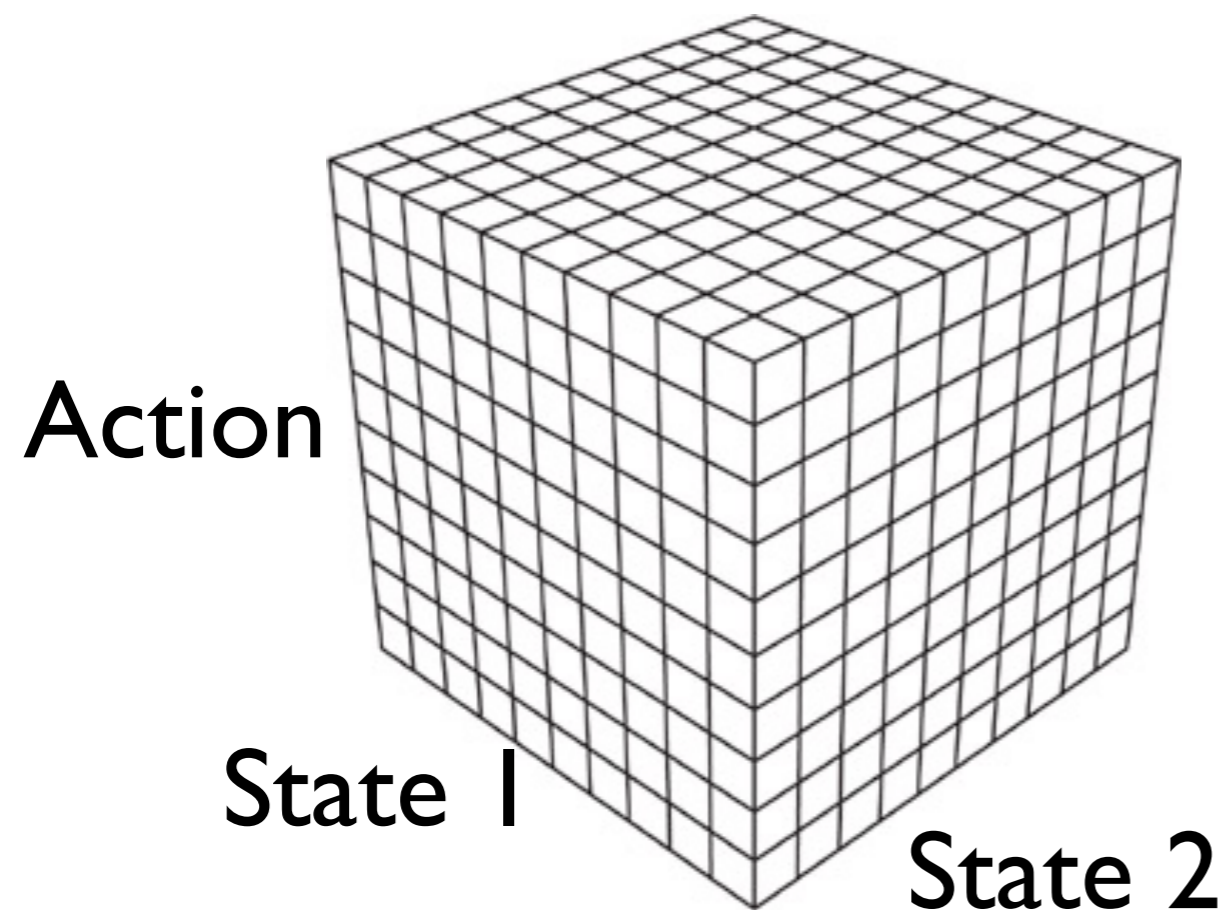


SA-Value

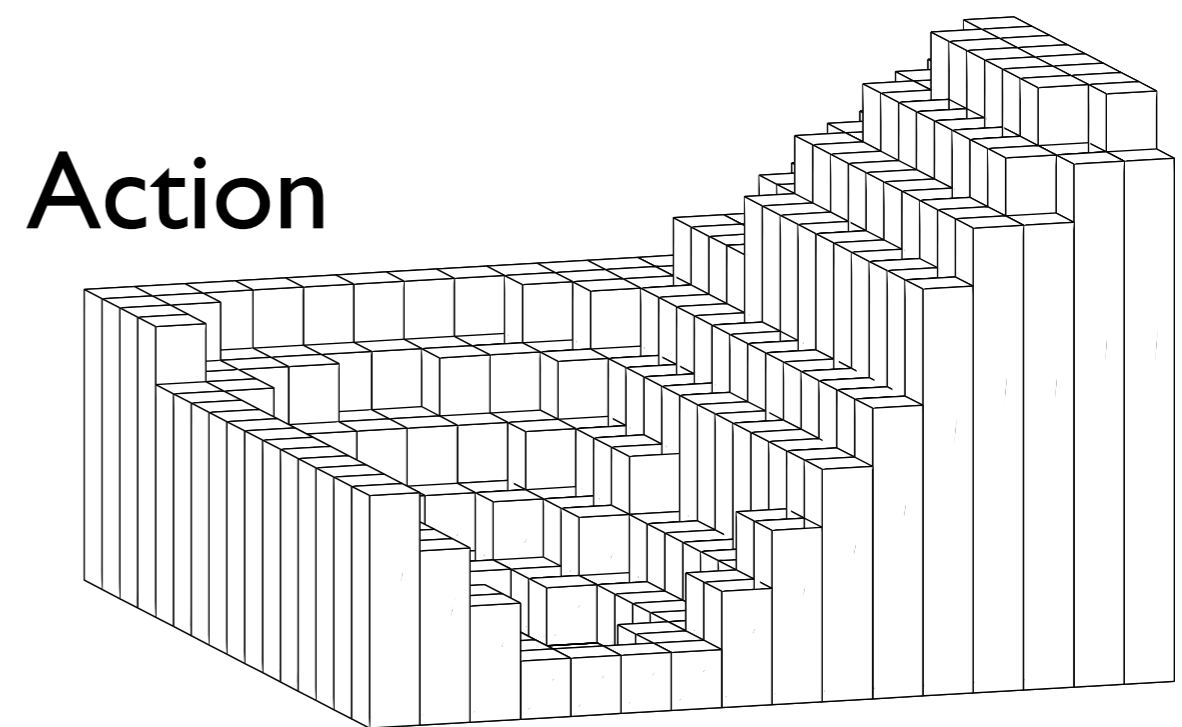




# Action discretization



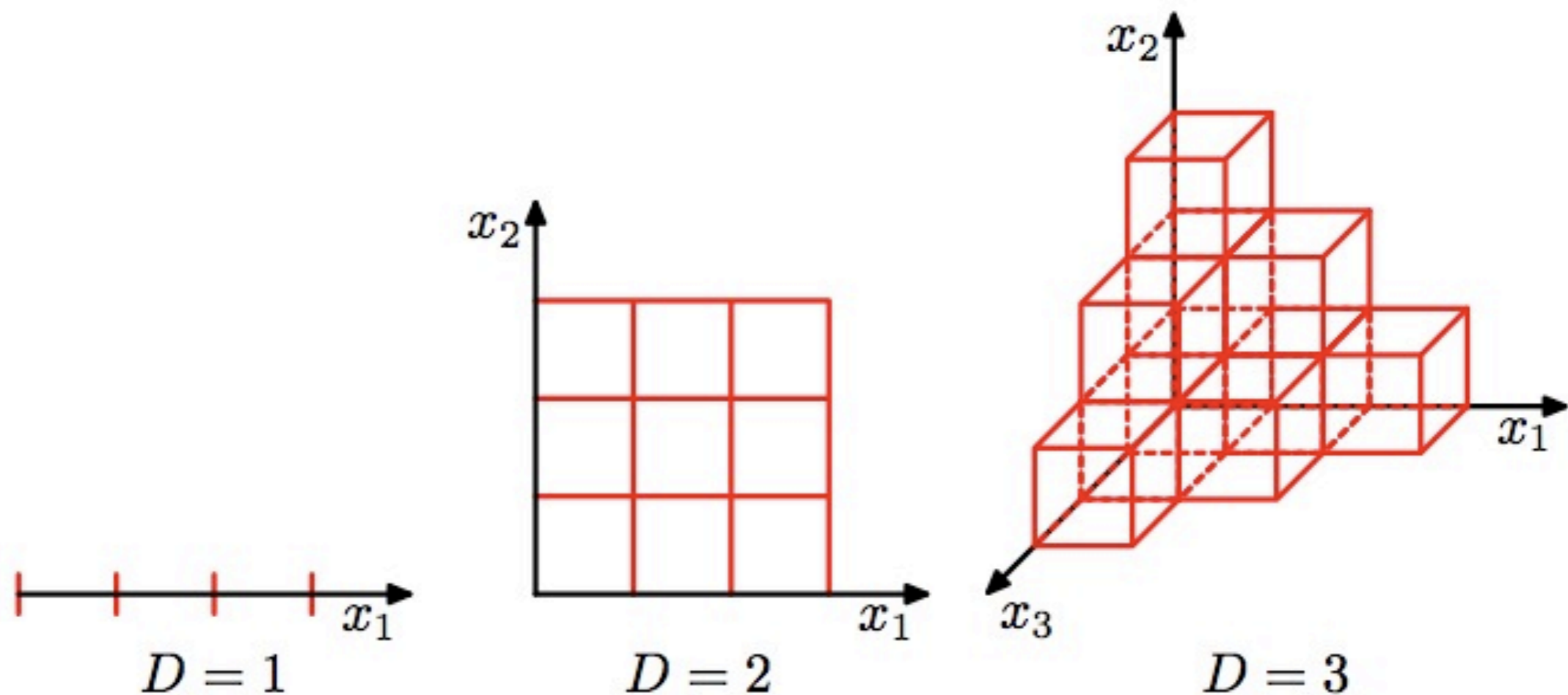
Grid



Policy



# Discretization of large/ high-dim. state spaces



[Bishop]

cf. Discussion on exploration



# Examples

High dimensional continuous state actions spaces with stochastic dynamics

Optimal(?) control in fluids

Approach: Computational Fluid dynamics & Evolutionary Algorithm



# Stefan Kern and Petros Koumoutsakos

ETH Zurich

# Kristina Eschler

hkg Zurich

The logo for ETH Zurich, consisting of the letters 'ETH' in a bold, italicized, sans-serif font.

Eidgenössische Technische Hochschule Zürich  
Swiss Federal Institute of Technology Zurich

<http://www.cse-lab.ethz.ch>



$C = \text{'don't eat me!'}$

$$\frac{\partial C}{\partial \theta} = \frac{\partial \text{'don't eat me!'}}{\partial \text{'how to flap???'}}$$

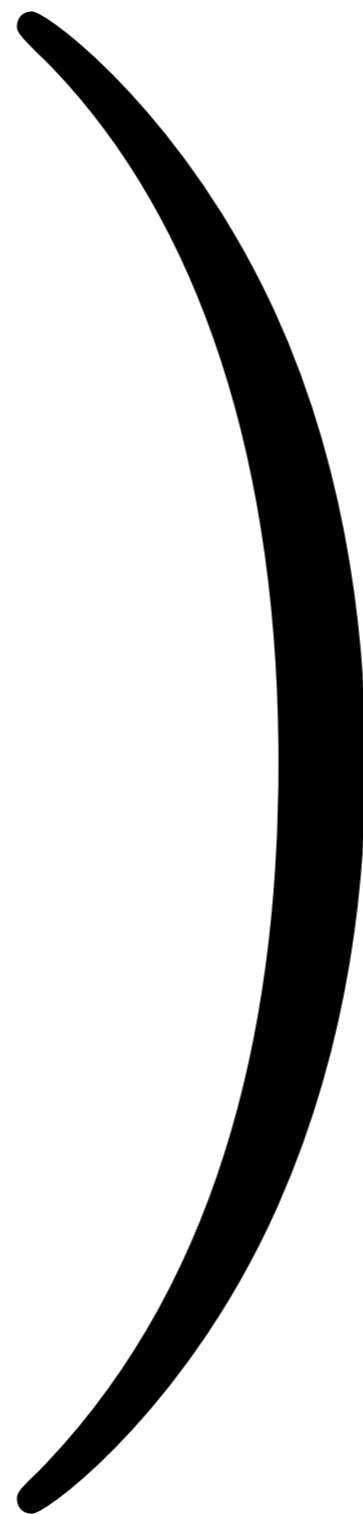


$$\rho \left( \frac{\partial \mathbf{v}}{\partial t} + \mathbf{v} \cdot \nabla \mathbf{v} \right) = -\nabla p + \nabla \cdot \mathbf{T} + \mathbf{f}$$

# Reinforcement Learning: real-world- sampling based optimal control

Why do high dimensional systems appear  
'repeatable/low-dimensional to us?'





Buchli - OLCAR - 2015





# Discretization issues

- Inflexible (need to decide division ahead of time)
- Inefficient (e.g. if slow varying function but division was decided to be fine)
- or not precise enough... if tiling is too coarse

Is there a way to avoid the issues of tiling  
and get a handle on the complexity?

Ideally: parameter(s) controlling  
complexity (in this class)

even more ideal: complexity adjusted  
automatically (not addressed in this  
class)



# Function approximation

Goal: approximate a given  
(arbitrary) function

Need:  
'Basis functions'  
Parameters



# Function approximation



# Function approximation

Function approximation:  $f(x, \theta) \approx y(x)$

function approximator/model

target function (e.g. observed)

$$\min (f(x) - y(x)) \{ \forall x \}$$

does not work, no finite  
minimum

$$\min \{ \|f(x) - y(x)\| \}$$

$\theta$  ?

need 'approximators' that can  
'easily' be tuned and can  
express arbitrary functions



# Norms

inner product

$$\|\mathbf{x}\| := \sqrt{\mathbf{x} \cdot \mathbf{x}}$$

Euclidian Norm

$$\|\mathbf{x}\| := \sqrt{x_1^2 + \dots + x_n^2}$$

$\mathbb{R}^n$

only defined on Euclidean spaces

p-Norm

$$\|\mathbf{x}\|_p := \left( \sum_{i=1}^n |x_i|^p \right)^{1/p}$$

l-Norm

$$\|\mathbf{x}\|_1 := \sum_{i=1}^n |x_i|$$

max-Norm

$$\|\mathbf{x}\|_\infty := \max(|x_1|, \dots, |x_n|)$$

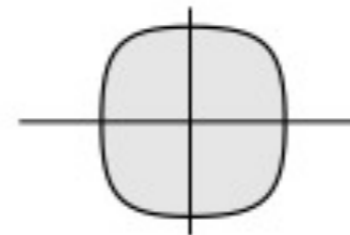
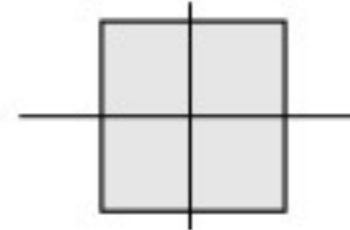
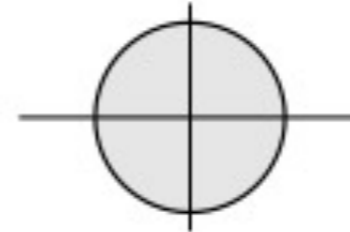
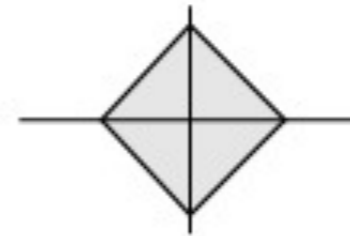


$$\|x\|_1 = \sum_{i=1}^m |x_i|,$$

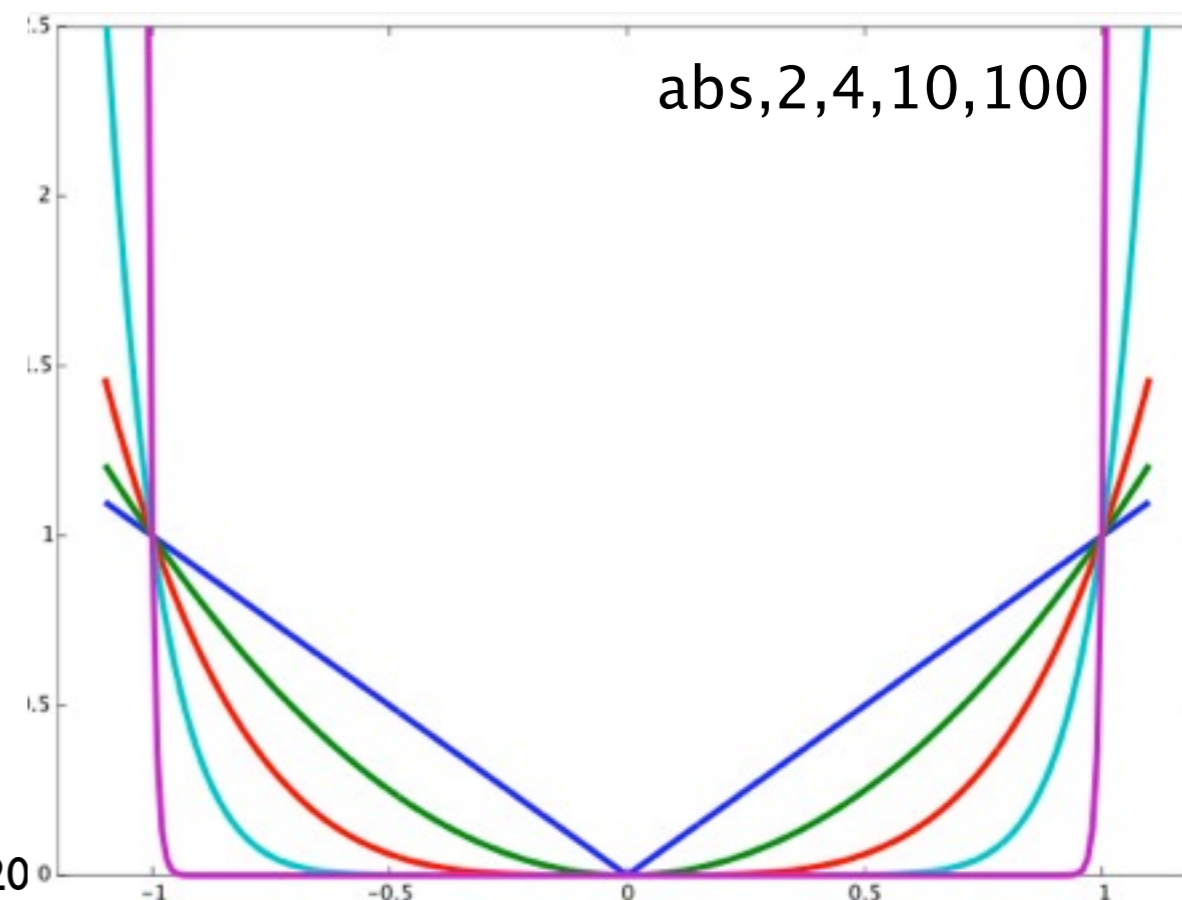
$$\|x\|_2 = \left( \sum_{i=1}^m |x_i|^2 \right)^{1/2} = \sqrt{x^* x},$$

$$\|x\|_\infty = \max_{1 \leq i \leq m} |x_i|,$$

$$\|x\|_p = \left( \sum_{i=1}^m |x_i|^p \right)^{1/p} \quad (1 \leq p < \infty).$$



$$\|\mathbf{x}\|_p = \left( x_1^p + \dots + x_n^p \right)^{\frac{1}{p}}$$



# Norms

inner product

$$\|\mathbf{x}\| := \sqrt{\mathbf{x} \cdot \mathbf{x}}$$

Euclidian Norm

$$\|\mathbf{x}\| := \sqrt{x_1^2 + \dots + x_n^2}$$

$\mathbb{R}^n$

only defined on Euclidean spaces

p-Norm

$$\|\mathbf{x}\|_p := \left( \sum_{i=1}^n |x_i|^p \right)^{1/p}$$

l-Norm

$$\|\mathbf{x}\|_1 := \sum_{i=1}^n |x_i|$$

max-Norm

$$\|\mathbf{x}\|_\infty := \max(|x_1|, \dots, |x_n|)$$

Example:

$$x = \begin{bmatrix} p \\ \alpha \end{bmatrix}$$

$$[\mathbf{x}] = \begin{bmatrix} m \\ rad \end{bmatrix}$$

~~$$[m^2] + [rad^2]$$~~

no 'natural' definition of distance





# Weighted p-norms

Weighted Euclidean

$$\|x\|_W = \|Wx\|$$

W: diagonal weighting matrix

$$\|x\|_W = \left( \sum_{i=1}^m |w_i x_i|^2 \right)^{1/2}$$

# Function approximation as optimization

e.g. linear regression:

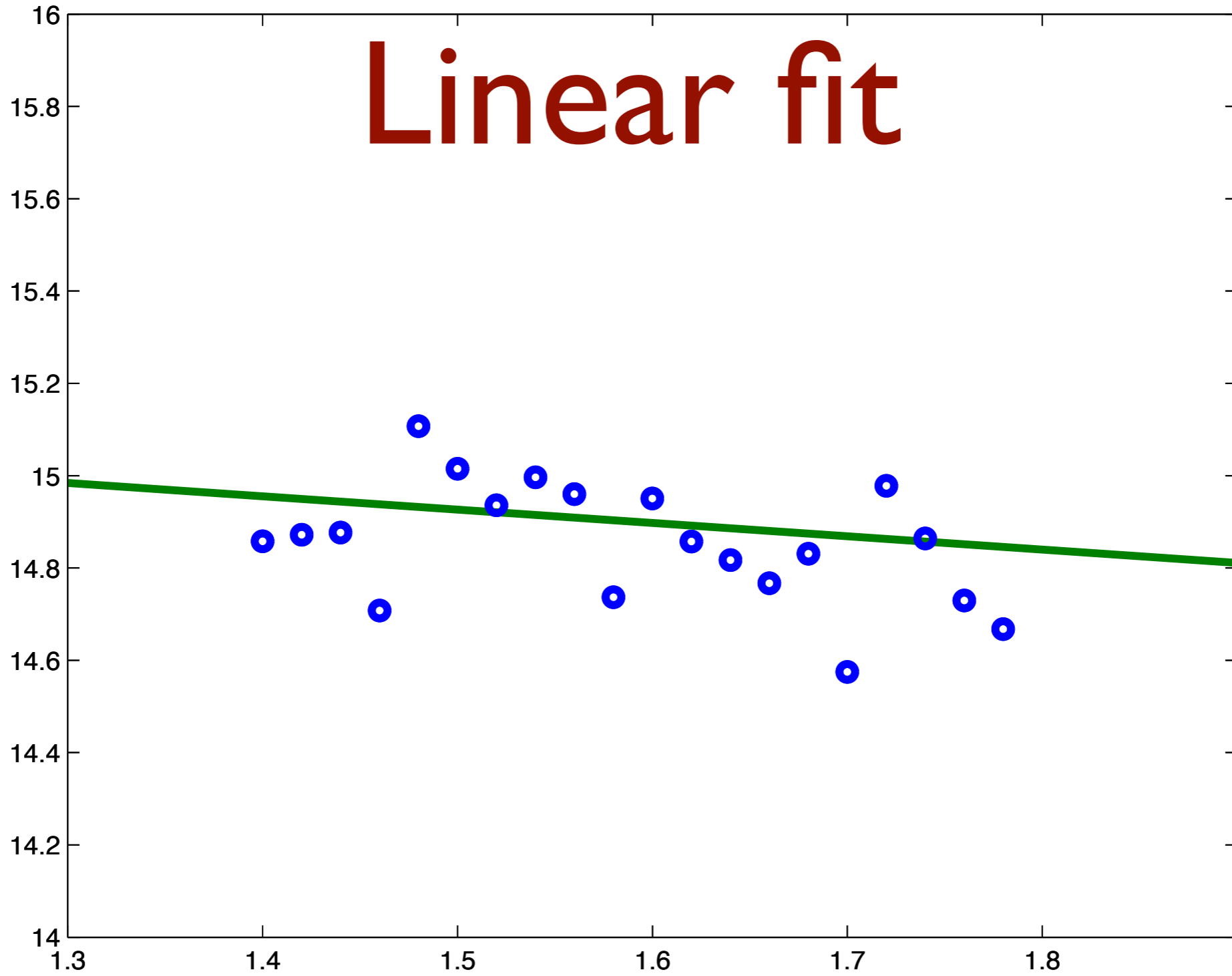
opt. in least squares sense...

Cost function = error function

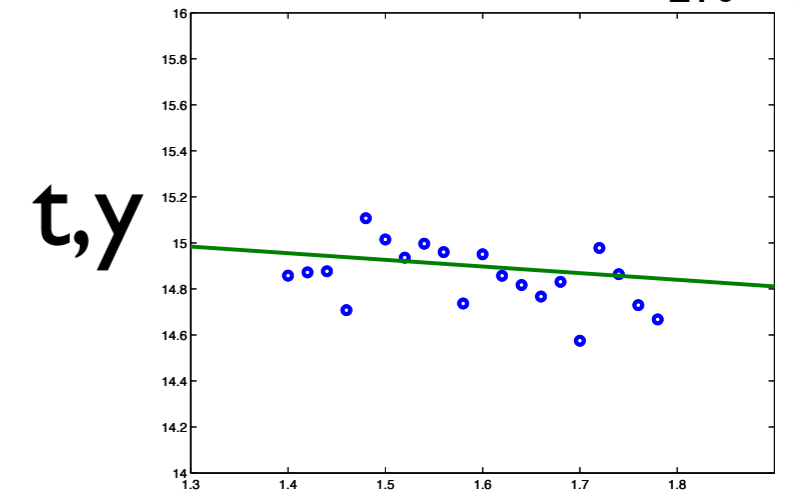
Look at the gradient, set it to 0

This is key for lots of learning methods





# Least squares linear regression



Input variables:  $\mathbf{x} = [x_1, \dots, x_n]$

Observations:  $t$  with gaussian noise:  $t = y + \epsilon$

Parameters:  $\mathbf{b} = [b_1, \dots, b_n]$

Linear model  $y = \mathbf{b}\mathbf{x}^T$

Find  $\mathbf{b}$  such that

$p$  observations:  $\mathbf{t} = \begin{bmatrix} t_1 \\ \dots \\ t_p \end{bmatrix}$

$$\min \|\mathbf{X}\mathbf{b}^T - \mathbf{t}\|$$

at input:

$$\mathbf{X} = \begin{bmatrix} [x_1, \dots, x_n]_1 \\ \dots \\ [x_1, \dots, x_n]_p \end{bmatrix}$$



# Solving for LS fit

$$\min \|\mathbf{X}\mathbf{b}^T - \mathbf{t}\| \stackrel{\text{same minimum}}{\Leftrightarrow} \min \left[ (\mathbf{X}\mathbf{b}^T - \mathbf{t})^T (\mathbf{X}\mathbf{b}^T - \mathbf{t}) \right] \begin{array}{l} \text{'quadratic form'} \\ \text{'cost function!'} \\ \text{'distance'} \end{array}$$

$$E = (\mathbf{X}\mathbf{b}^T - \mathbf{t})^T (\mathbf{X}\mathbf{b}^T - \mathbf{t})$$

$$\min E \Leftrightarrow \nabla E = 0$$

$$\nabla E = 2\mathbf{X}^T (\mathbf{X}\mathbf{b}^T - \mathbf{t})$$

$$(\mathbf{X}^T \mathbf{X}) \quad n \times n$$

$$\mathbf{b}^T = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{t}$$

Pseudoinverse!

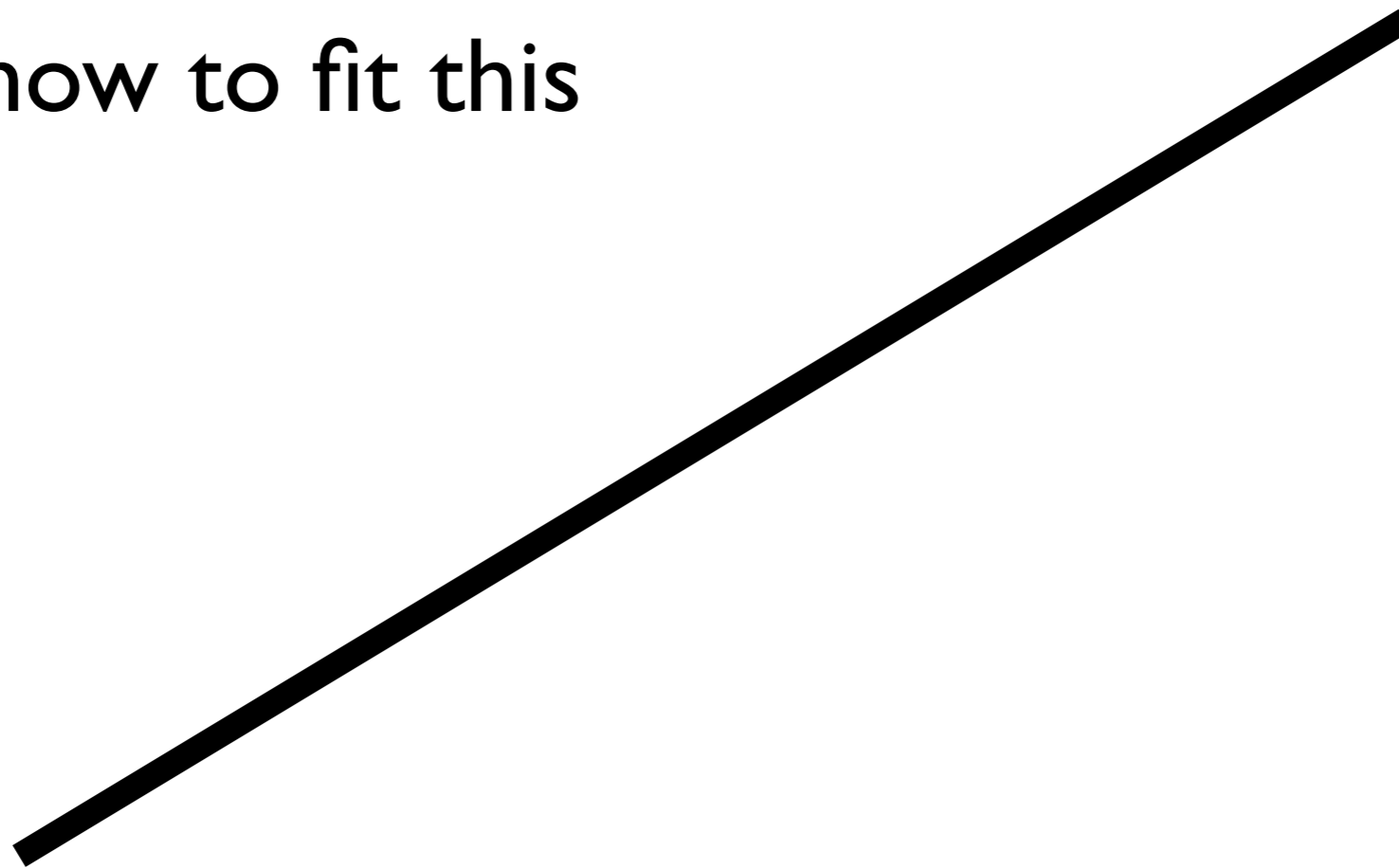


Matlab: pinv

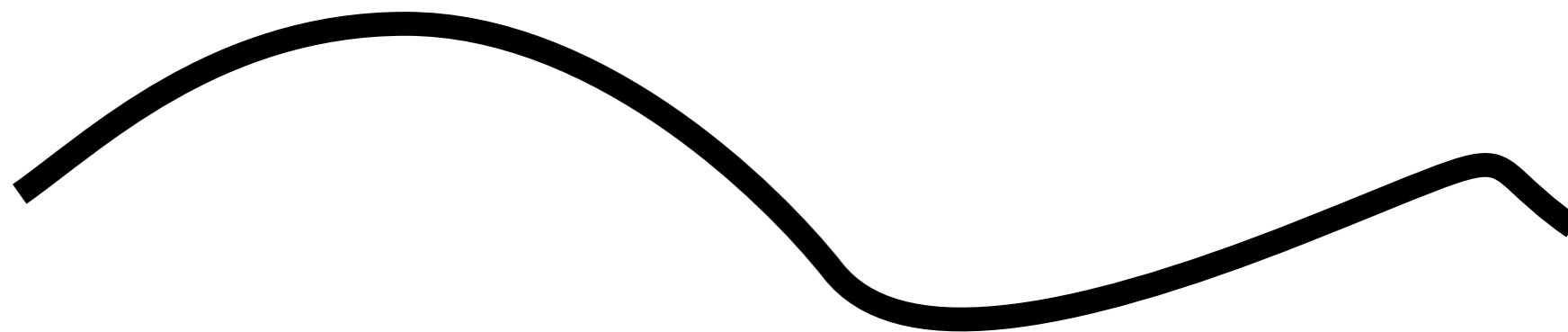
Buchli - OLCAR - 2015

**ETH** zürich

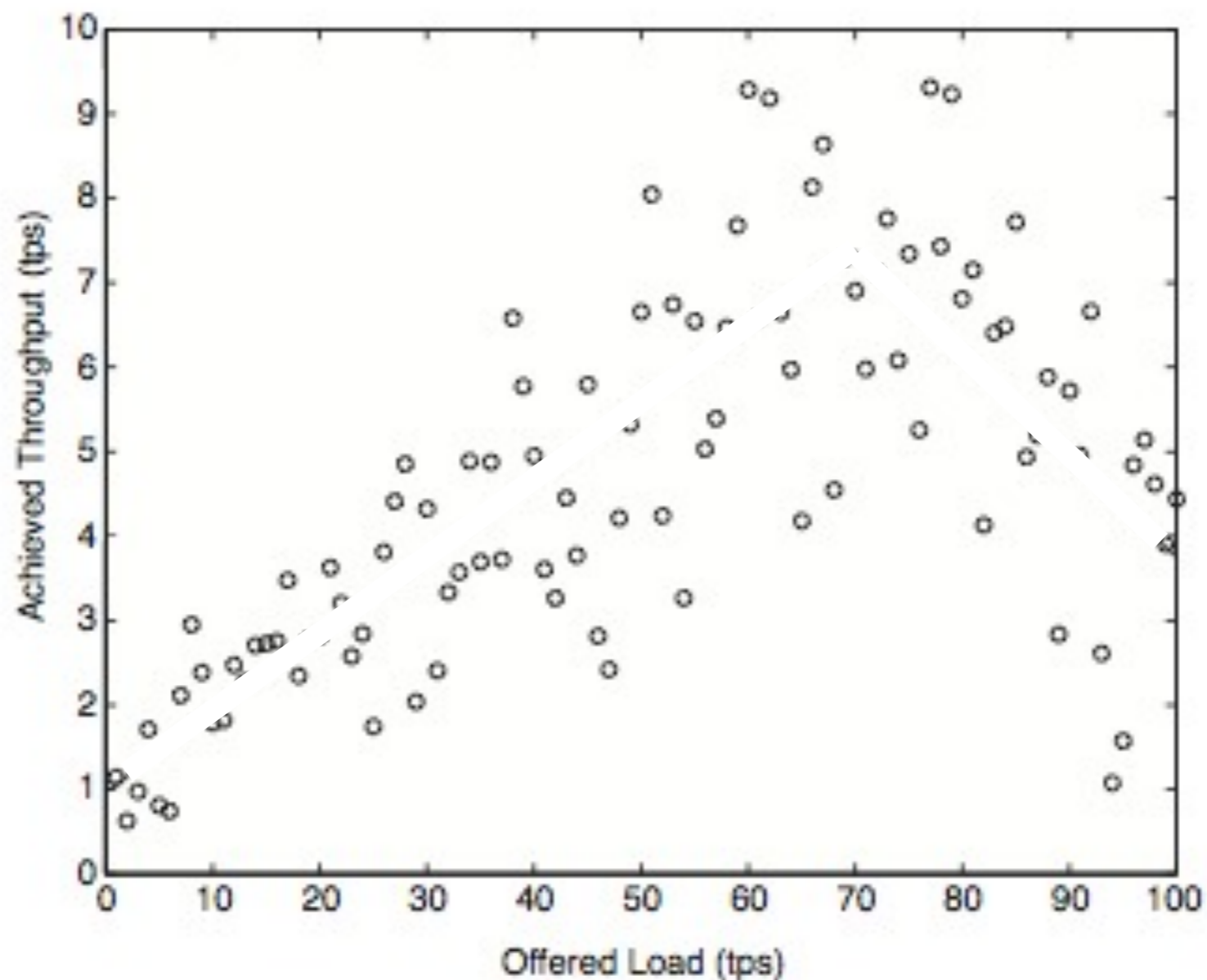
know how to fit this



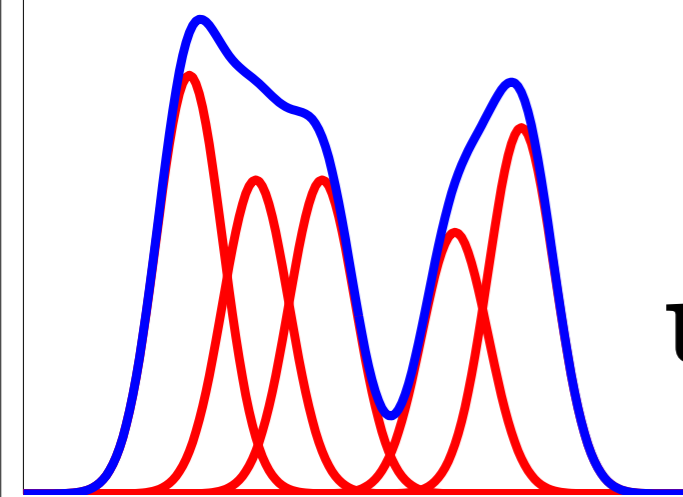
but what about this?



# Fit with linear model?



# Function approximation using basis functions



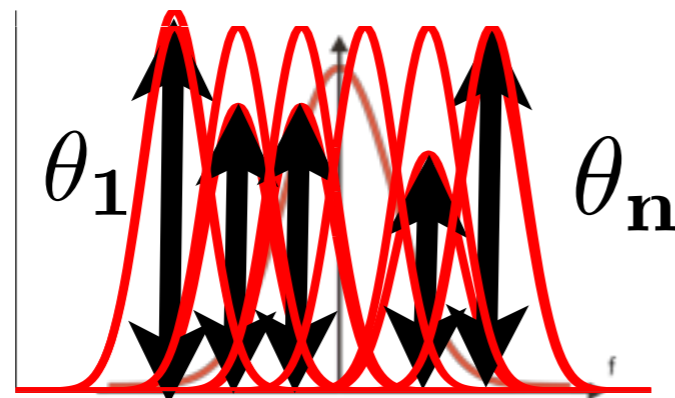
$$\mathbf{u}(\mathbf{t}) = \boldsymbol{\phi}(\mathbf{t})^T \boldsymbol{\theta}$$

Basis  
functions

Torques  
Current  
ref. position  
'Heading'  
...

Parameters

Learned





'inner product'

$$u(t) = \phi(\mathbf{t})^T \boldsymbol{\theta}$$

$$u(t) = \boldsymbol{\theta}^T \phi(\mathbf{t})$$

Let's look at this expression in more detail

$$\phi(\mathbf{t}) = \begin{bmatrix} \phi_0(t) \\ \cdot \\ \cdot \\ \cdot \\ \phi_n(t) \end{bmatrix} \quad \boldsymbol{\theta} = \begin{bmatrix} \theta_0 \\ \cdot \\ \cdot \\ \cdot \\ \theta_n \end{bmatrix} \quad \boldsymbol{\theta} = [\theta_0, \dots, \theta_n]^T$$

$$u(t) = \theta_0 \phi_0(t) + \dots + \theta_n \phi_n(t)$$

# Example: $ax+b$

$$y = ax + b$$

$$\theta = [a \quad b]^T$$

$$\phi(\mathbf{x}) = \begin{bmatrix} x \\ 1 \end{bmatrix}$$

# Example: $ax+b$

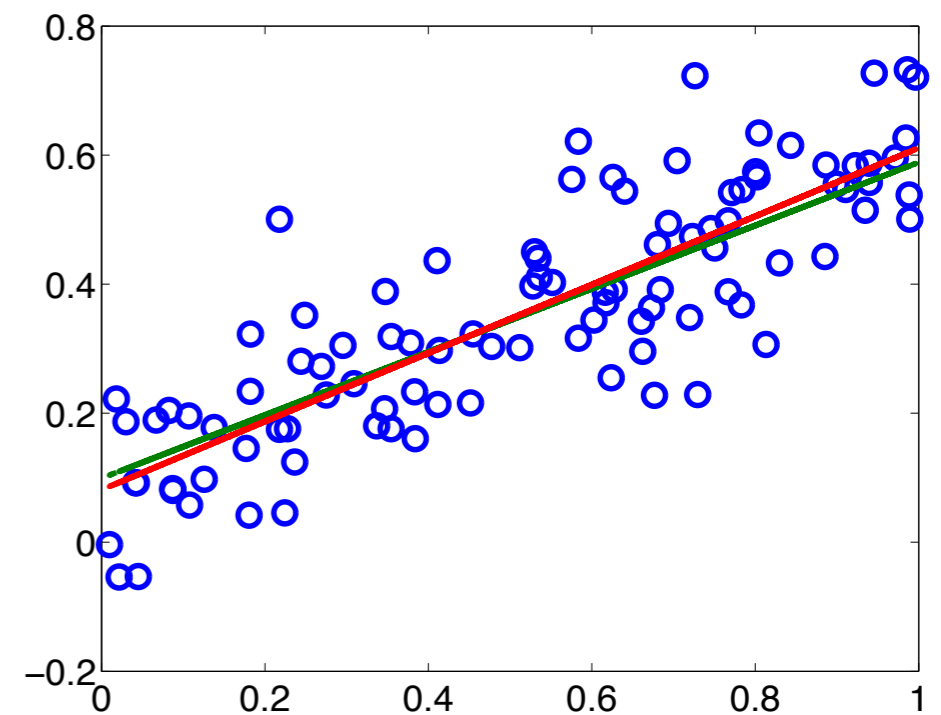
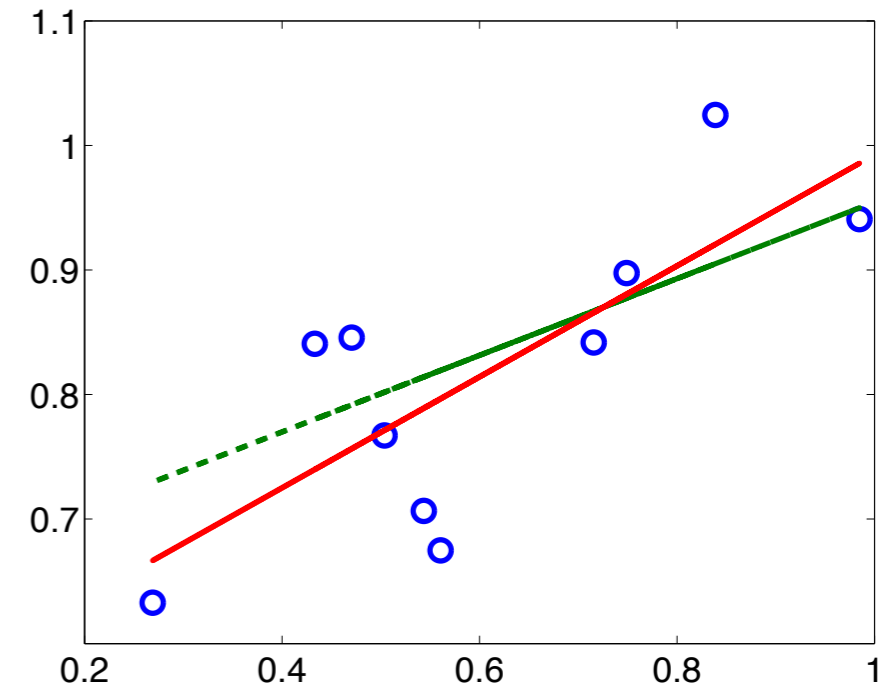
```
% how many observations
p = 10;

x = rand(p,1);
b1 = rand
b2 = rand
eps = 0.1 * randn(p,1);
t = b2 * x + b1 + eps;
x1 = [x, ones(p,1)];

b_est = pinv(x1) * t

h = plot(x,t,'o',x,t-eps,x,b_est' * x1');
```

Try  $p = 10$   
 Try  $p = 100$



# Basis functions

Goal: function approximation of arbitrary functions

Wanted: 'Good' function approximator

Good: easy to find parameters, expressive, ...

$$f(x) = \sum_i w_i \Psi_i(x)$$

Idea: push nonlinearity into the basis functions, independent of parameters

Linear in parameters!!!!

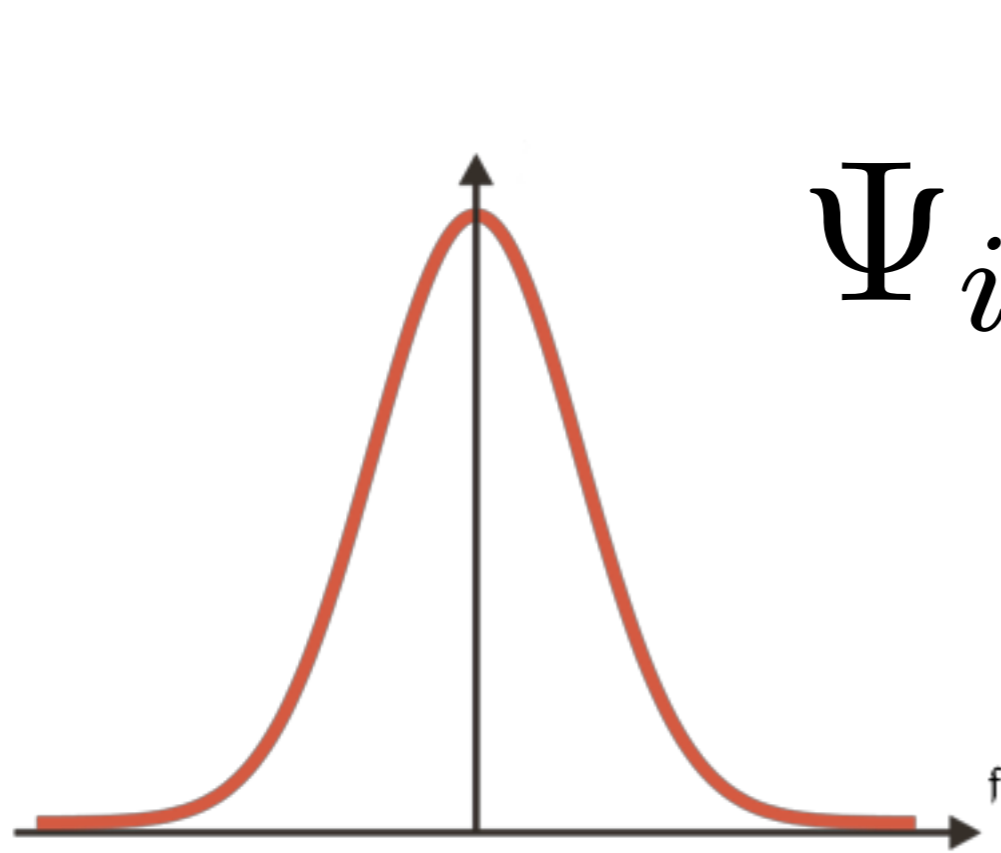


# Fourier basis

$$\frac{a_0}{2} + \sum_{n=1}^{\infty} [a_n \cos(nx) + b_n \sin(nx)]$$

Problem: infinite validity, need infinitely many basis functions to approximate a non-periodic function

# Gaussian basis



$$\Psi_i(x) = e^{-\frac{(x - c_i)^2}{b_i}}$$

Idea: push nonlinearity into the basis functions

basis: more localized than sines

$$f(x) = \sum_i w_i \Psi_i(x)$$

Note: can take derivatives of any order of the gaussian basis...

# Polynomial basis

$$f(x, \theta) = \theta_0 + \theta_1 x + \theta_2 x^2 + \dots$$

$$f(x, \theta) = \sum_{n=0}^{\infty} \theta_n x^n$$

# Other bases

- Index functions

$$f(x, \theta) = \sum_{n=0}^{\infty} \theta_n b_n(x)$$

$$b(x) = 1 \quad \forall x \in [a, b], \quad 0 \text{ otherwise}$$

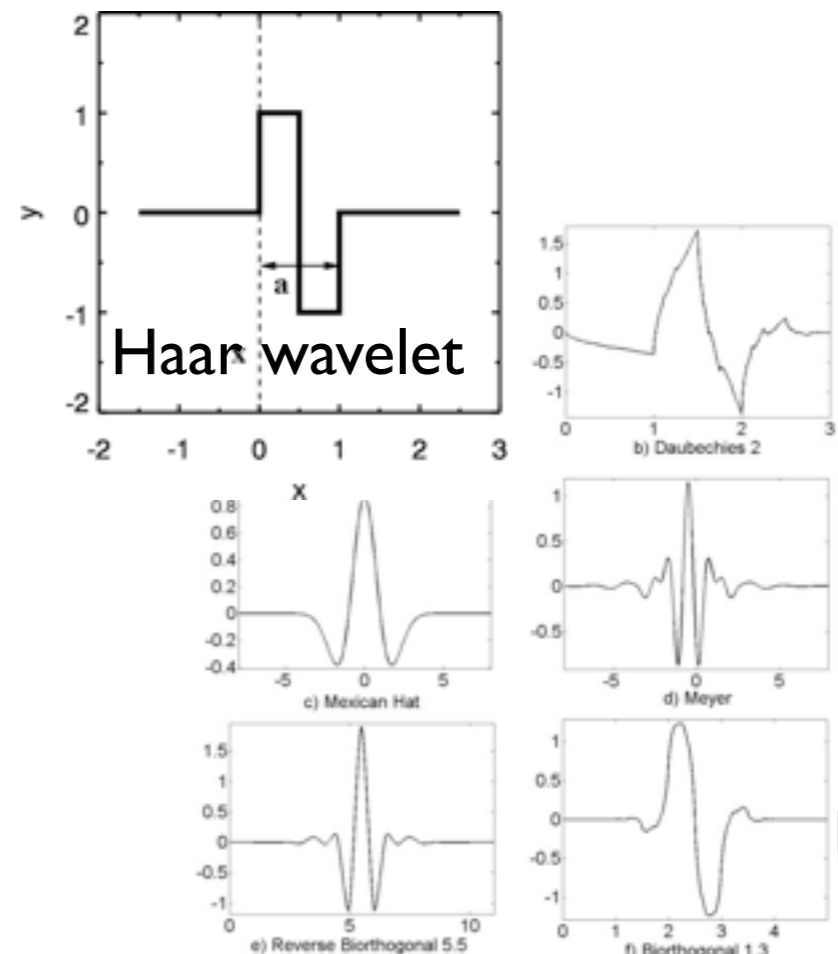
- Index functions multiplied with other functions  
(e.g. local linear models)

$$b(x) = g(x) \quad \forall x \in [a, b], \quad 0 \text{ otherwise}$$

$$b(x) = x \quad \forall x \in [a, b], \quad 0 \text{ otherwise}$$

- Wavelets: lots of different bases

Basis functions can be defined on  $\mathbb{C}$





# The function basis zoo

All bases are equal,  
but some bases are more equal than others!



# Fit nonlinear functions

Linear model: 
$$y(\mathbf{x}, \mathbf{w}) = \sum_{j=0}^{M-1} w_j \phi_j(\mathbf{x}) = \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x})$$

Gaussian basis function: 
$$\phi_j(x) = \exp \left\{ -\frac{(x - \mu_j)^2}{2s^2} \right\}$$

Observation:  $t = y(\mathbf{x}, \mathbf{w}) + \epsilon$

$$p(t|\mathbf{x}, \mathbf{w}, \beta) = \mathcal{N}(t|y(\mathbf{x}, \mathbf{w}), \beta^{-1})$$

precision (inv. variance)  $\beta$

Max. likelihood treatment leads to:

$$\begin{aligned} \ln p(\mathbf{t}|\mathbf{w}, \beta) &= \sum_{n=1}^N \ln \mathcal{N}(t_n | \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}_n), \beta^{-1}) \\ &= \frac{N}{2} \ln \beta - \frac{N}{2} \ln(2\pi) - \beta E_D(\mathbf{w}) \end{aligned}$$

$$E_D(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N \{t_n - \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}_n)\}^2.$$



Gradient: 
$$\nabla \ln p(\mathbf{t}|\mathbf{w}, \beta) = \sum_{n=1}^N \{t_n - \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}_n)\} \boldsymbol{\phi}(\mathbf{x}_n)^T.$$

Gradient:

$$\nabla \ln p(\mathbf{t}|\mathbf{w}, \beta) = \sum_{n=1}^N \{t_n - \mathbf{w}^T \phi(\mathbf{x}_n)\} \phi(\mathbf{x}_n)^T.$$

Gradient = 0

$$0 = \sum_{n=1}^N t_n \phi(\mathbf{x}_n)^T - \mathbf{w}^T \left( \sum_{n=1}^N \phi(\mathbf{x}_n) \phi(\mathbf{x}_n)^T \right).$$

$$\mathbf{w}_{\text{ML}} = (\Phi^T \Phi)^{-1} \Phi^T \mathbf{t}$$

$$\Phi^\dagger \equiv (\Phi^T \Phi)^{-1} \Phi^T$$

Pseudoinverse!

$$\Phi = \begin{pmatrix} \phi_0(\mathbf{x}_1) & \phi_1(\mathbf{x}_1) & \cdots & \phi_{M-1}(\mathbf{x}_1) \\ \phi_0(\mathbf{x}_2) & \phi_1(\mathbf{x}_2) & \cdots & \phi_{M-1}(\mathbf{x}_2) \\ \vdots & \vdots & \ddots & \vdots \\ \phi_0(\mathbf{x}_N) & \phi_1(\mathbf{x}_N) & \cdots & \phi_{M-1}(\mathbf{x}_N) \end{pmatrix} \quad N \times M$$

# Why not using pseudo-inverse?

What happens if you have millions of datapoints/observations?  
... or lots of dimensions in the problem?

$$\mathbf{w} = (\Phi^T \Phi)^{-1} \Phi^T \mathbf{t}$$



$$\Phi^\# = (\Phi^T \Phi)^{-1} \Phi^T$$

**dimensions?**

square matrix  $M \times M$   
(each entry  $N$  multiplications)

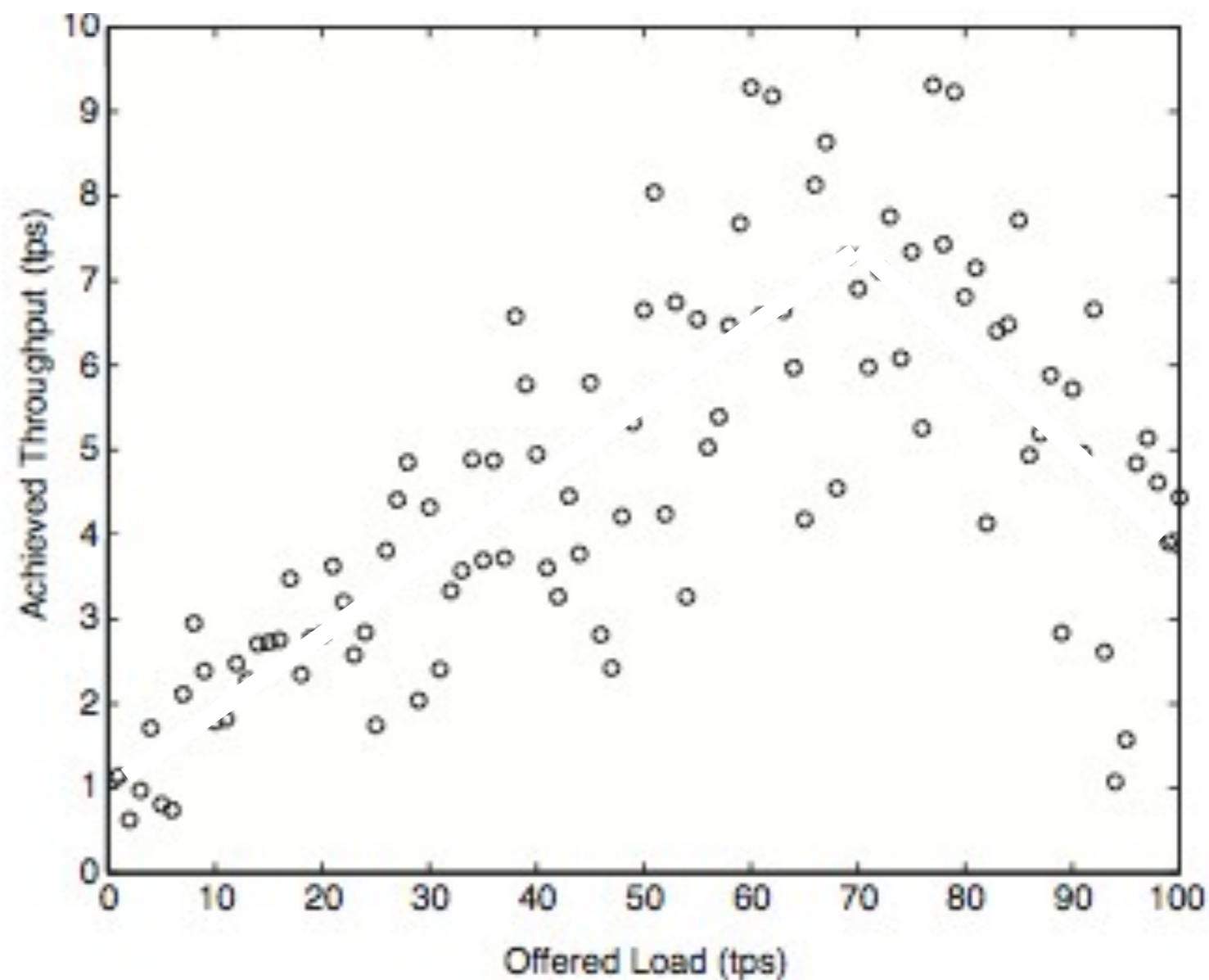
- use SVD (found iteratively)
- Iterative least squares

Matlab: `svd`



# What basis function is used in Le Boudec's example?

$$Y_i = (a + bx_i)1_{x_i \leq \xi} + (c + dx_i)1_{\{x_i > \xi\}} + \epsilon_i$$



# Nonlinear fitting

It's easy to fit a linear function

... or anything where the parameters show up  
linear!!!



# Example: Sampling

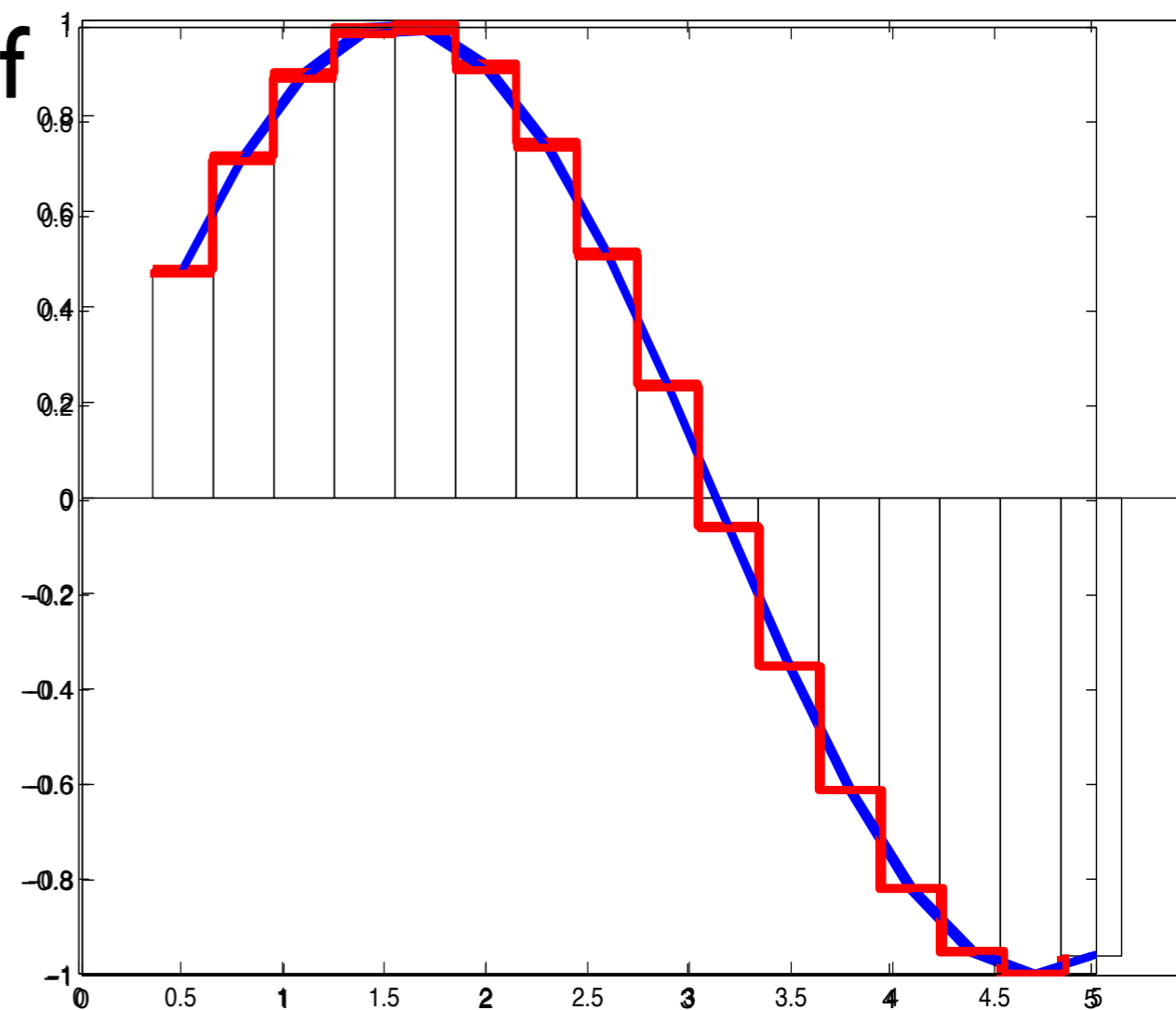


# Function approximation view of sampling

discretization is special case of  
function approximation

can make the same argument  
with dirac pulse sampling

remember: no overlap / no  
generalization





Function approximation as 'bridge' between  
continuous and discrete world

Can lower dimensionality of the problem  
(dimensionality becomes an open parameter!!!)

Tradeoff: high  $N$ , good approximation, optimal policy,  
curse of dimensionality



# Approximate what?

Supervised learning: Target function known → difference is cost

Reinforcement learning: Cost is sampled from examples  
But think of it this way: There is a 'target function', i.e. the optimal control/or value function and the approximator has to minimize the distance between the 'guess' and this function

- Direct policy learning (e.g. policy gradients)
- Value function approximation

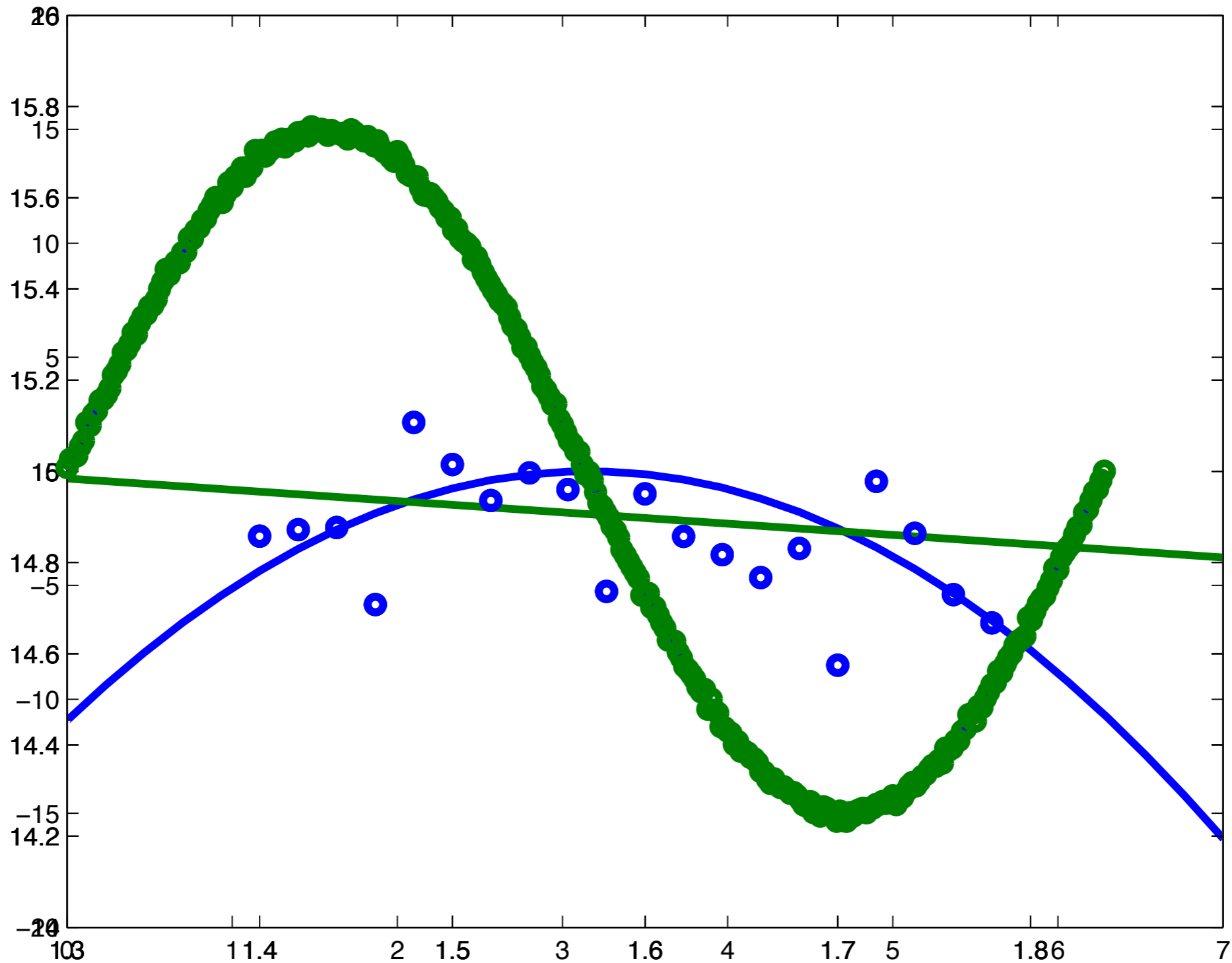


# Some issues...



# Extrapolation vs. Interpolation



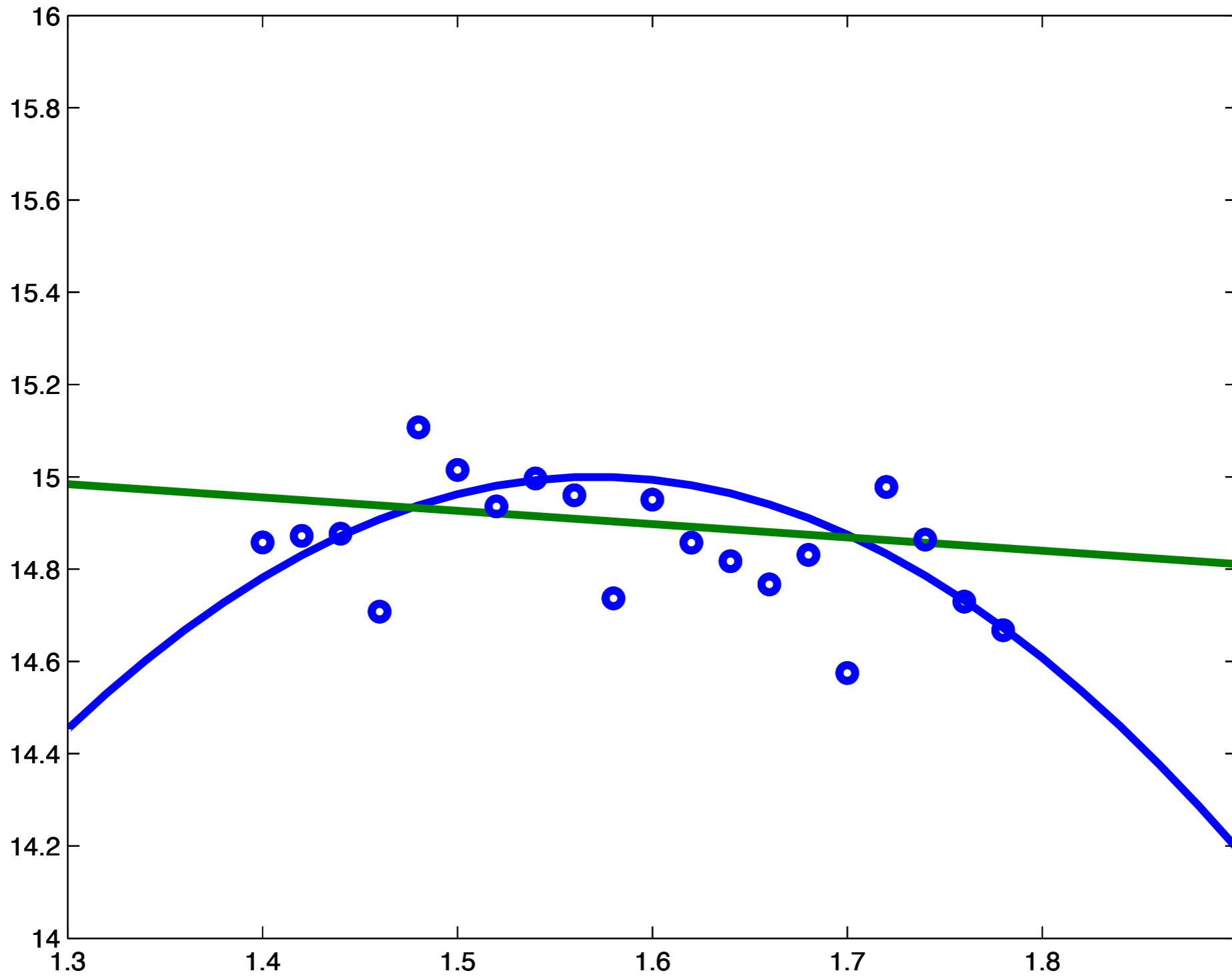


# Extrapolate?

## Need a model.. let's fit one..

Buchli - OLCAR - 2015





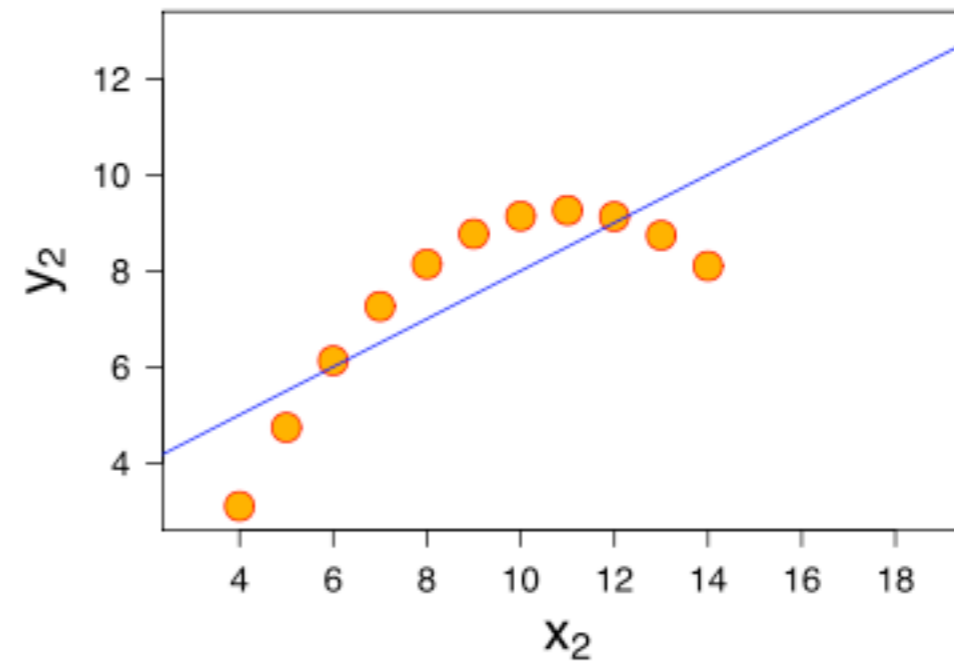
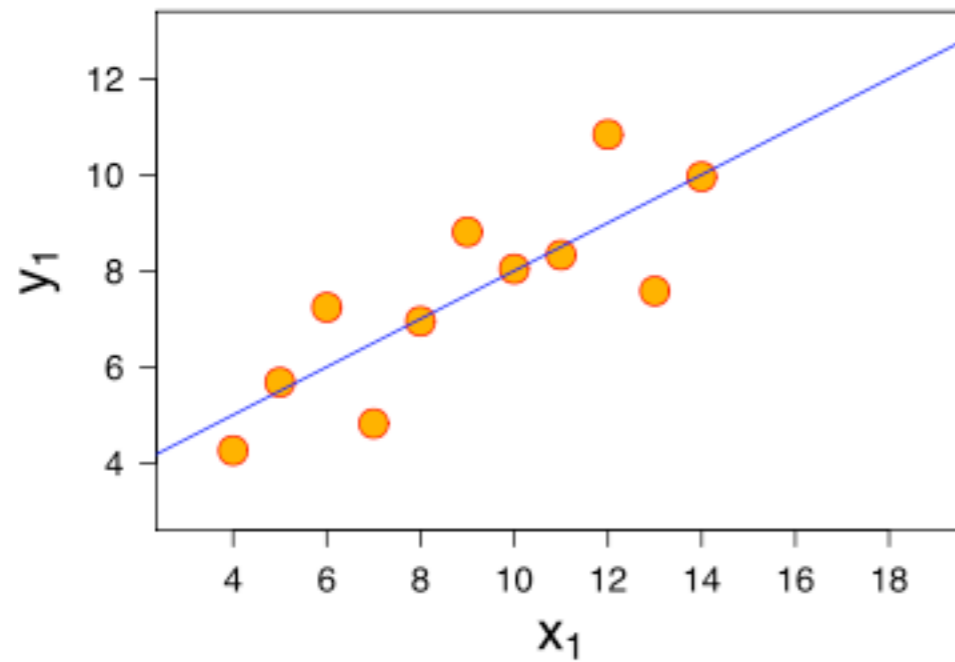
# Interpolation?

Buchli - OLCAR - 2015

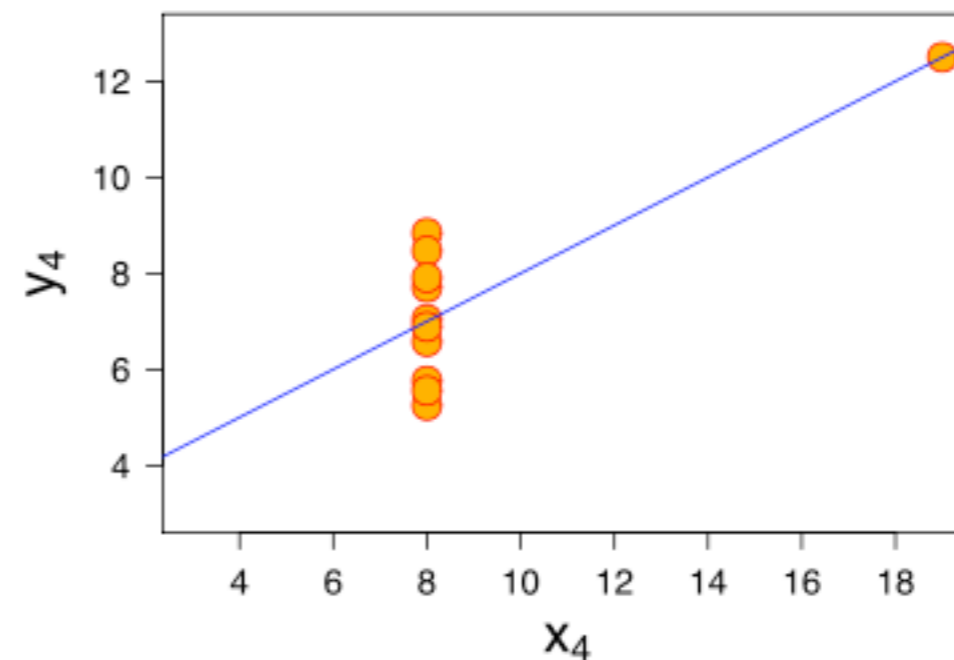
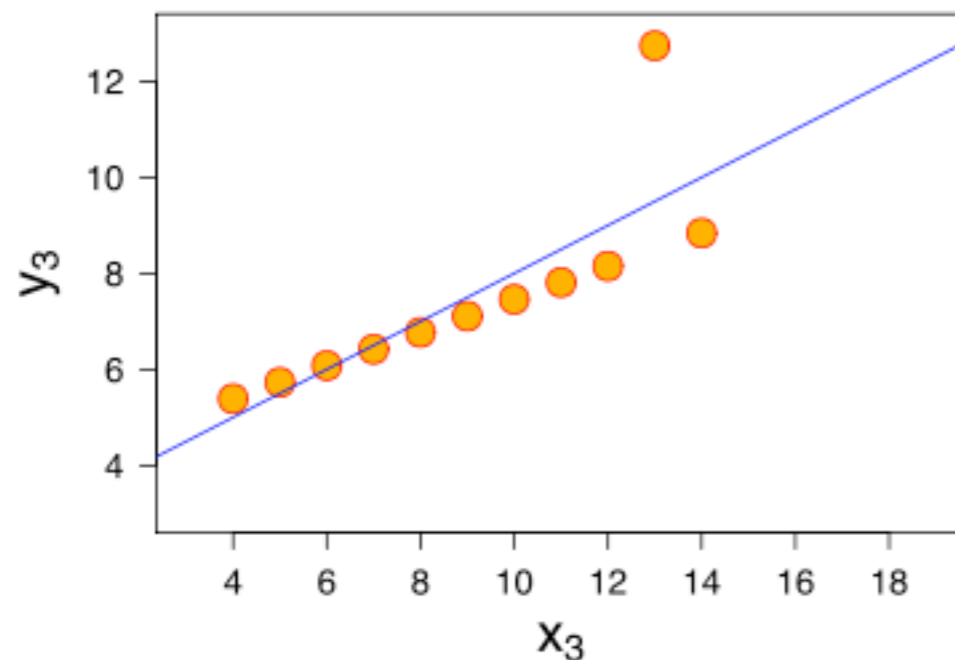


# Model choice, model bias and explanatory power of models

# Anscombe's quartet



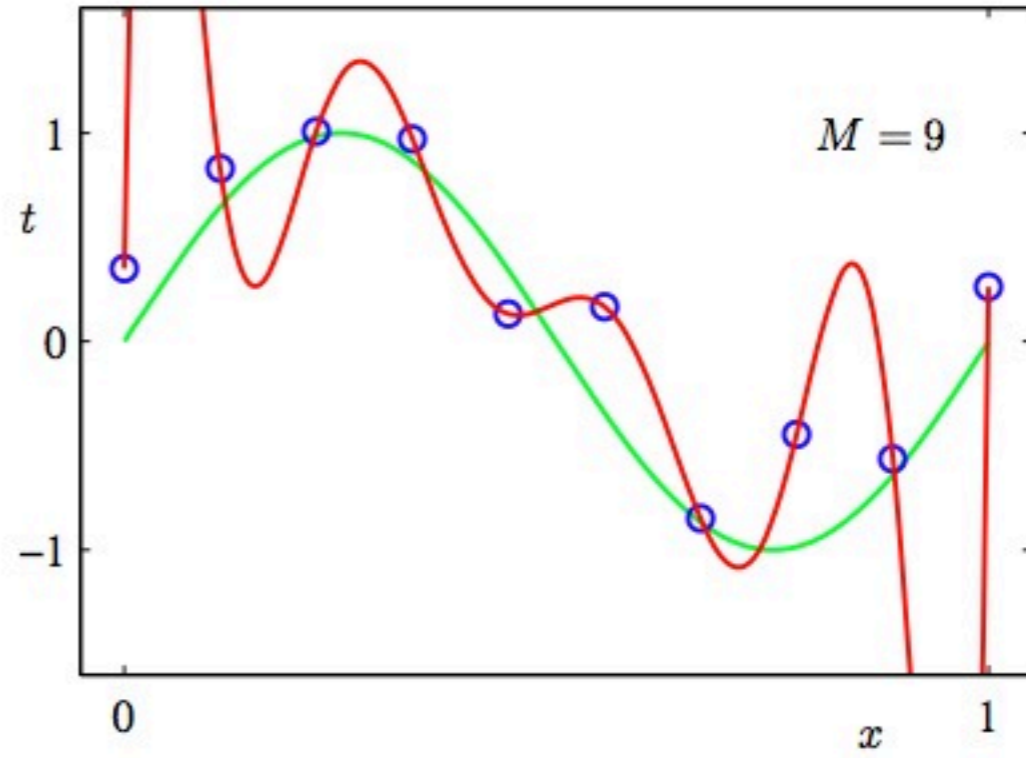
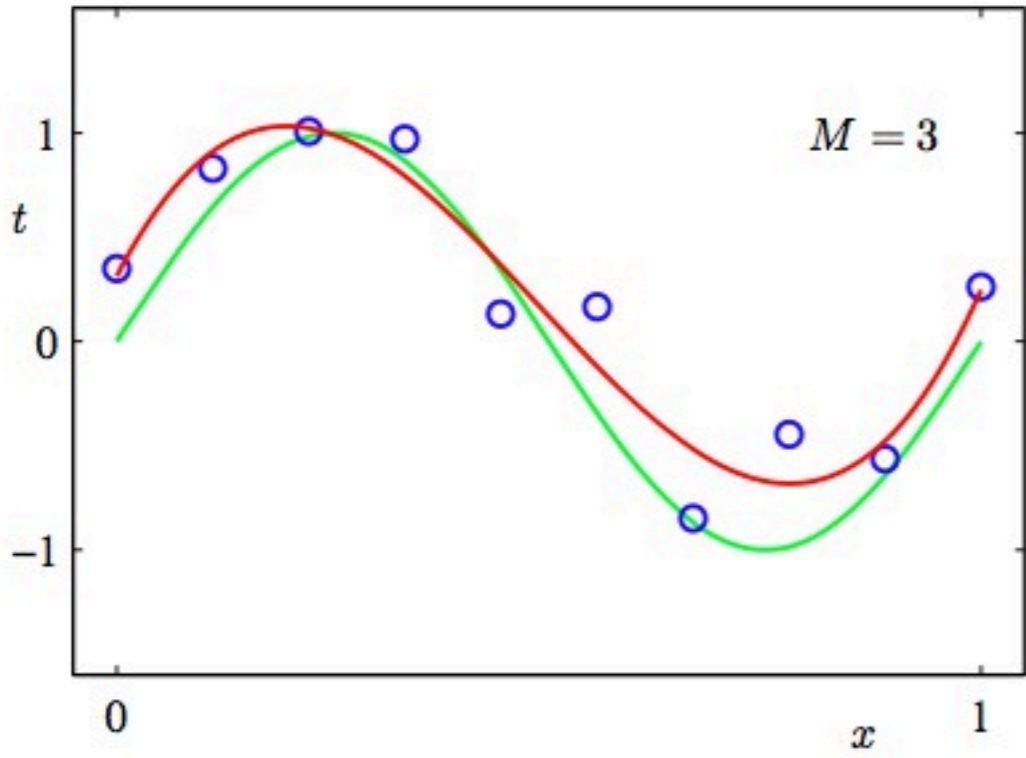
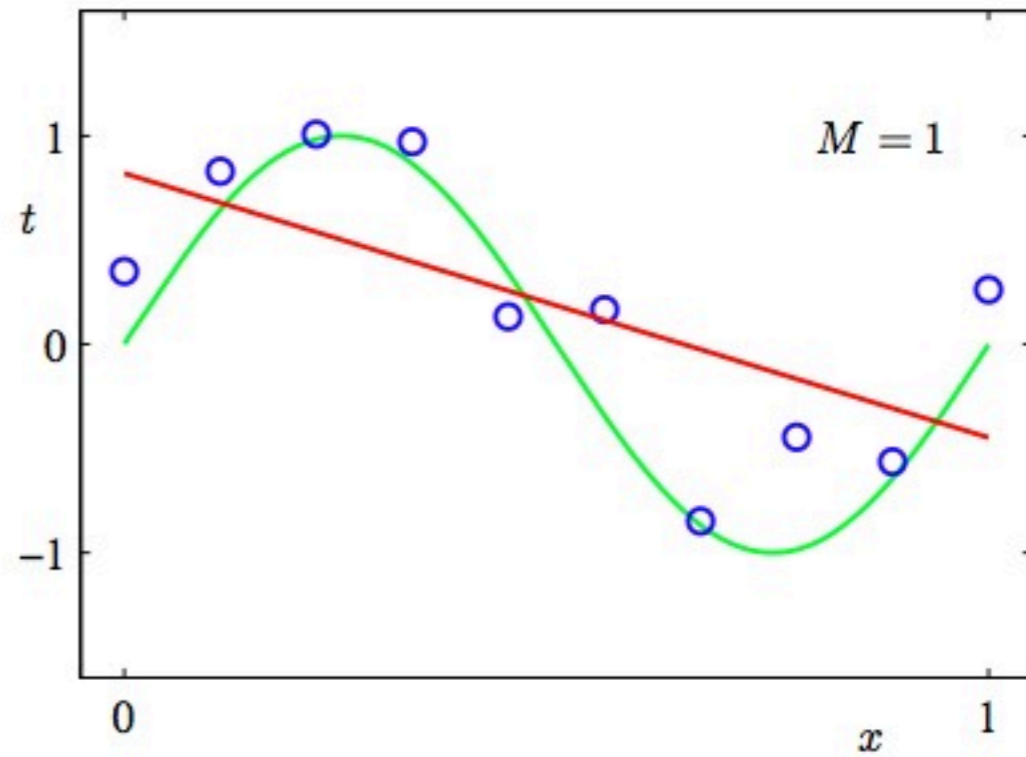
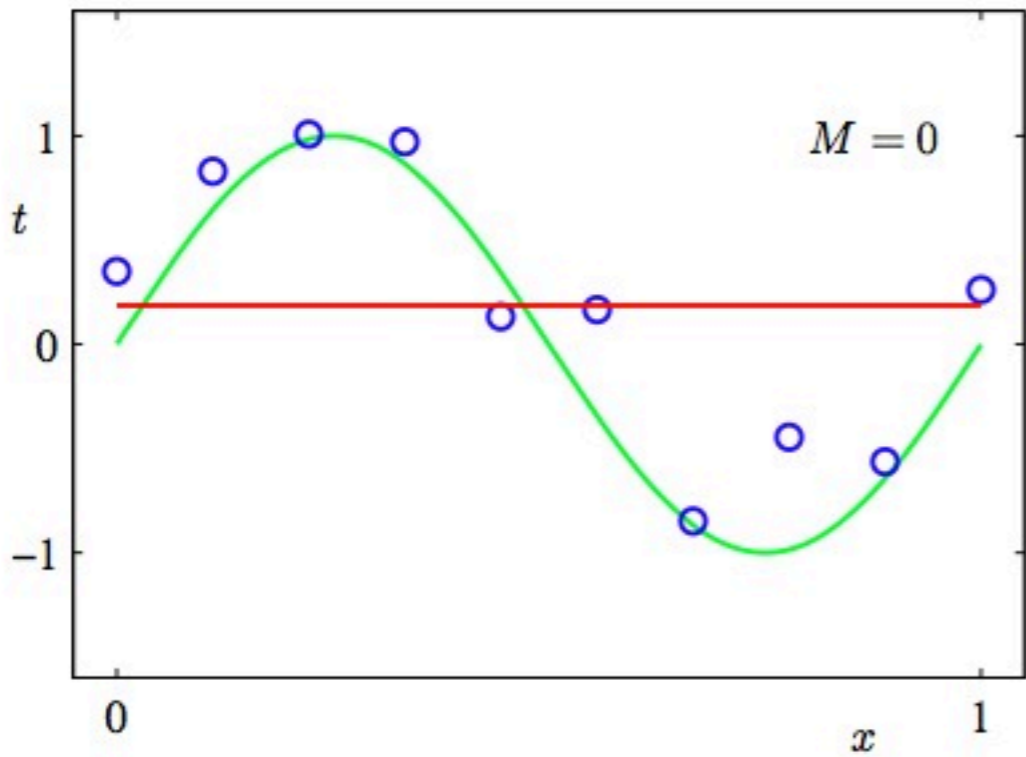
least squares very prone to outliers





# Overfitting

The approximator has to have enough expressive power  
to capture the essentials,  
... but it should not try to over-interpret the data!



# State vs time based policies

Policy:

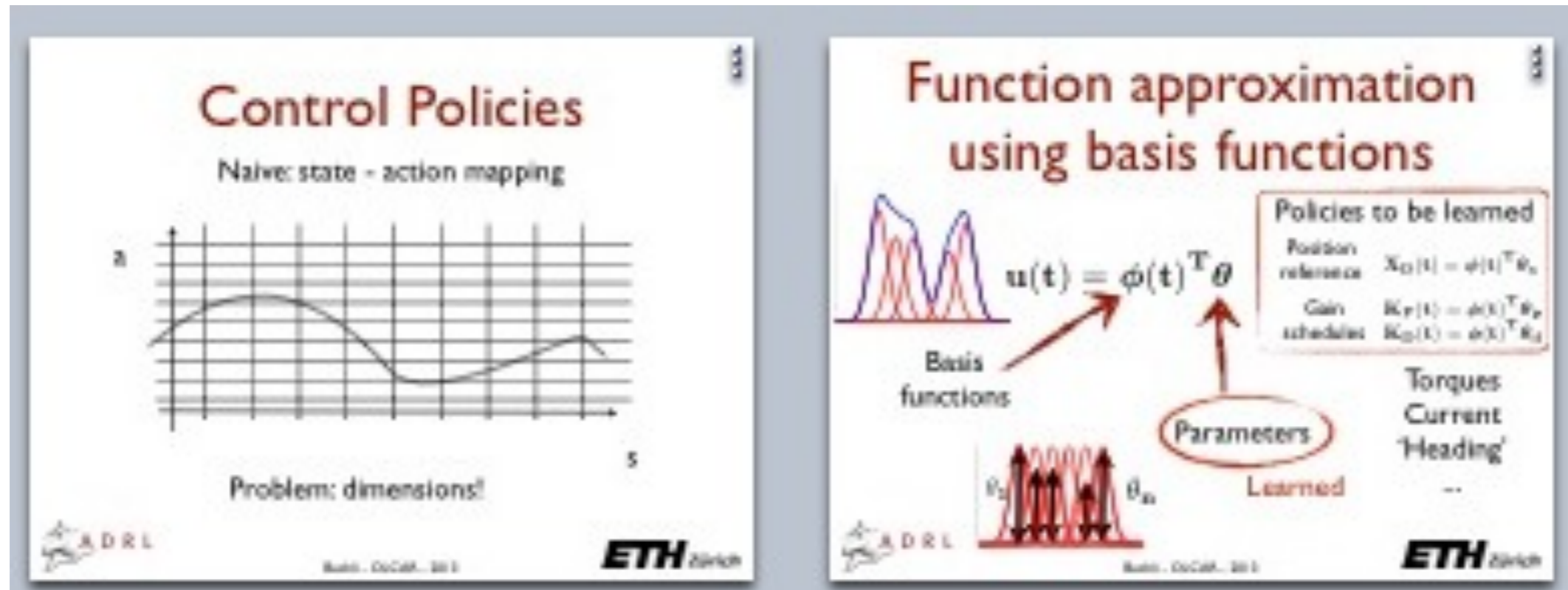
function of states

function of time

basis functions

$$\Psi(x)$$

$$\Psi(t)$$



## 'classic RL' Sutton & Barto

p 9 on efficiency of 'evolutionary' vs. direct value based methods

p17 on how opt. control IS learning and emphasis on 'incremental'!

## Parametrized Policy Learning (applied RL in Robotics in last ~10y)

see Examples/Papers

⇔ Planning



# Path integral RL

Goal: Solve continuous time stochastic optimal control problem by sampling

⇒ First look at Path Integral Optimal Control



(Reminder from L3)

# Reminder

# Continuous time optimal control



(Reminder from L3)

# Continuous time system

System dynamics

$$\dot{\mathbf{x}}(t) = \mathbf{f}_t(\mathbf{x}(t), \mathbf{u}(t))$$

Cost

$$J = e^{-\beta(t_f - t_0)} \Phi(\mathbf{x}(t_f)) + \int_{t_0}^{t_f} e^{-\beta(t - t_0)} L(\mathbf{x}(t), \mathbf{u}(t)) dt$$

$0 \leq \beta$  discount / decay rate

'exponential decay'



# Continuous time optimal control problem

(Reminder from L3)

Find control  $u^*(t) = \mu^*(t, x(t))$  minimizing

control (input)                      policy

$$J = e^{-\beta(t_f - t_0)} \Phi(\mathbf{x}(t_f)) + \int_{t_0}^{t_f} e^{-\beta(t - t_0)} L(\mathbf{x}(t), \mathbf{u}(t)) dt$$

Given constraints

$$\dot{\mathbf{x}}(t) = \mathbf{f}_t(\mathbf{x}(t), \mathbf{u}(t))$$

Goal: Optimal policy

$$\mu^* = \arg \min_u J$$





# Hamilton Jacob Bellman Equation



Carl Gustav Jacob Jacobi  
(1804-1851)



William Rowan Hamilton  
(1805-1865)



Richard Bellman  
1920-84

$$\frac{\partial V^*}{\partial t} = \beta V^* - \min_{\mathbf{u} \in \mathbf{U}} \left\{ L(\mathbf{x}, \mathbf{u}) + \left( \frac{\partial V^*}{\partial \mathbf{x}} \right)^T \mathbf{f}(\mathbf{x}, \mathbf{u}) \right\}$$

$$\beta V^* - \frac{\partial V^*}{\partial t} = \min_{\mathbf{u} \in \mathbf{U}} \left\{ L(\mathbf{x}, \mathbf{u}) + \left( \frac{\partial V^*}{\partial \mathbf{x}} \right)^T \mathbf{f}(\mathbf{x}, \mathbf{u}) \right\}$$

In general: Nonlinear, Partial Differential Equation  
Has no analytical solution... : (

Backwards in time!  $V^*(t_f, \mathbf{x}) = \Phi(\mathbf{x})$



(Reminder from L3)

# Reminder

# Stochastic optimal control



(Reminder from L3)

# Stochastic system

$$\dot{\mathbf{x}}(t) = \mathbf{f}_t(\mathbf{x}(t), \mathbf{u}(t)) + \mathbf{B}(t)\mathbf{w}(t), \quad \mathbf{x}(0) = \mathbf{x}_0$$

Mean:  $E[\mathbf{w}(t)] = \bar{\mathbf{w}} = 0$

mean-free

Co-variance:  $E[\mathbf{w}(t)\mathbf{w}(\tau)^T] = \mathbf{W}(t)\delta(t - \tau)$

uncorrelated over time

$$E[\mathbf{w}(t)\mathbf{w}(\tau)^T] = 0 \\ t \neq \tau$$

Expected cost:

$$J = E \left\{ e^{-\beta(t_f - t_0)} \Phi(\mathbf{x}(t_f)) + \int_{t_0}^{t_f} e^{-\beta(t' - t_0)} L(\mathbf{x}(t'), \mathbf{u}(t')) dt' \right\}$$



# Stoch. HJB

(Reminder from L3)

$$\beta V^*(t, \mathbf{x}) - V_t^*(t, \mathbf{x}) = \min_{\mathbf{u}(t)} \left\{ L(\mathbf{x}, \mathbf{u}(t)) + V_{\mathbf{x}}^{*T} \mathbf{f}_t(\mathbf{x}, \mathbf{u}(t)) + \frac{1}{2} \text{Tr}[V_{\mathbf{xx}}^* \mathbf{B}(t) \mathbf{W}(t) \mathbf{B}^T(t)] \right\}$$

Hamilton Jacobi Bellman Equation

$$V^*(t_f, \mathbf{x}) = \Phi(\mathbf{x})$$



# Towards Path-integral SOC

Goal: Solve continuous time stochastic optimal control problem by sampling

- ★ Can solve SOC with certain assumptions, eg. LQ.
- ★ Here a less stringent form of the dynamics (ctl. affine) and (almost) NO assumptions on cost

Why important?

Cost is key design 'handle' to tell the system what to do...



# Control affine opt control problem

Note similarity to stoch. system (L3)

$$\dot{\mathbf{x}}(t) = \mathbf{f}_t(\mathbf{x}(t), \mathbf{u}(t)) + \mathbf{B}(t)\mathbf{w}(t)$$

$$\dot{\mathbf{x}}_t = \mathbf{f}(\mathbf{x}_t, t) + \mathbf{G}(\mathbf{x}_t)\mathbf{u}_t$$

$$V(\mathbf{x}_0) = \min_{\mathbf{u}_{t_0:t_N}} \mathbb{E}[R(\tau)]$$

$$R(\tau) = \phi(t_N) + \int_{t_0}^{t_N} r_t dt \quad r_t = q(\mathbf{x}_t) + \frac{1}{2} \mathbf{u}_t^T R \mathbf{u}_t$$



$$V(x_t) = \min_{u_t} E_\tau [R(\tau)]$$

$$\int p(\tau) R(\tau) d\tau \quad \int p(\tau) \left( \frac{1}{\lambda} \phi + \frac{1}{\lambda} \int r dt \right) d\tau$$

Discretize and use EM-like idea: PoWER -  
 Problem: pseudo-probability - restriction on  
 cost function

Other idea: Treat probability as a diffusion  
 process - Connection with statistical physics  
 Forward dynamics! Sampling (Monte Carlo)!



# Continuous Random Processes

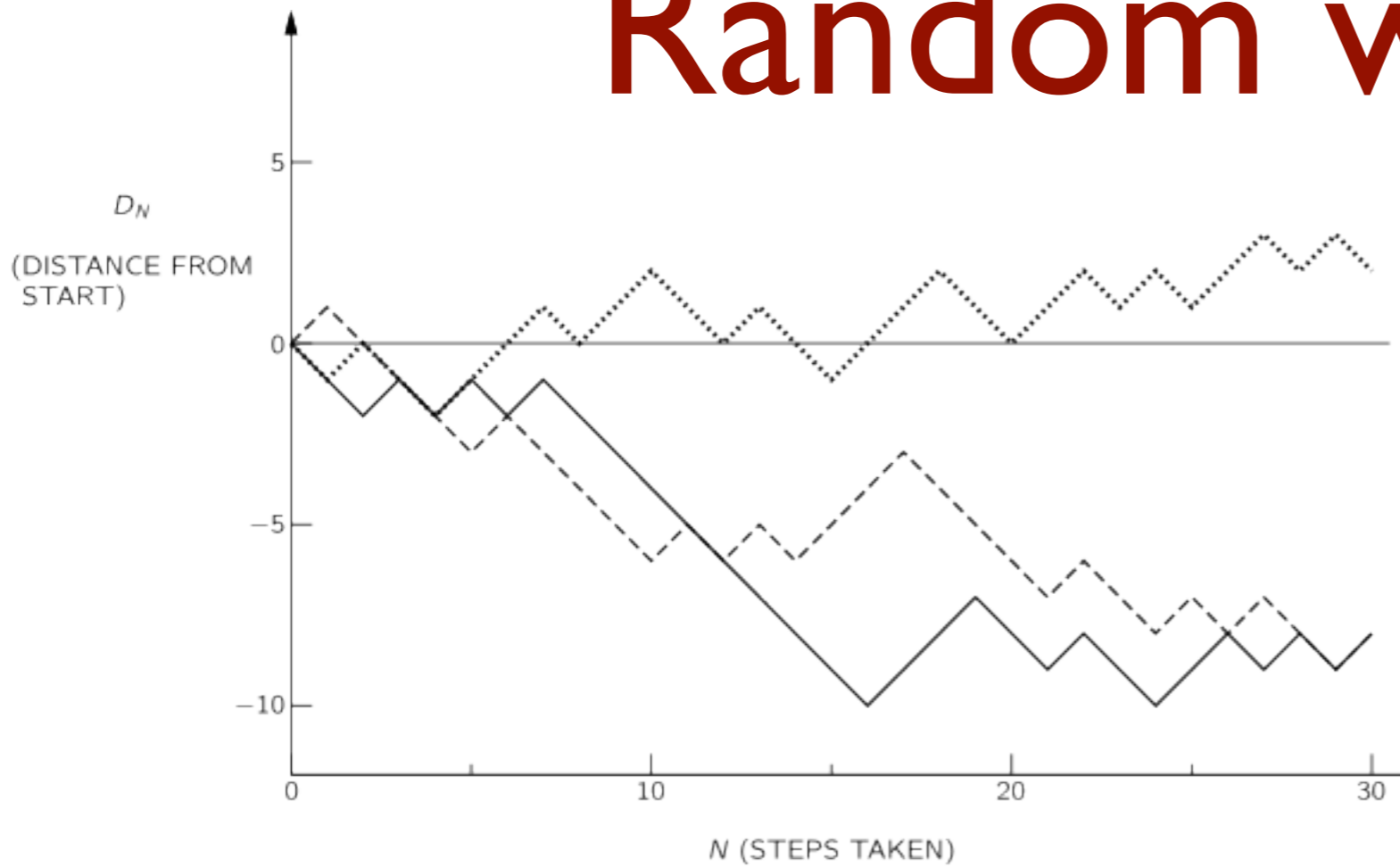
$$d\mathbf{x} = \mathbf{f}(\mathbf{x}, \mathbf{u}) dt + F(\mathbf{x}, \mathbf{u}) d\omega$$



# Major difficulty: Definition of stochastic processes in continuous time!

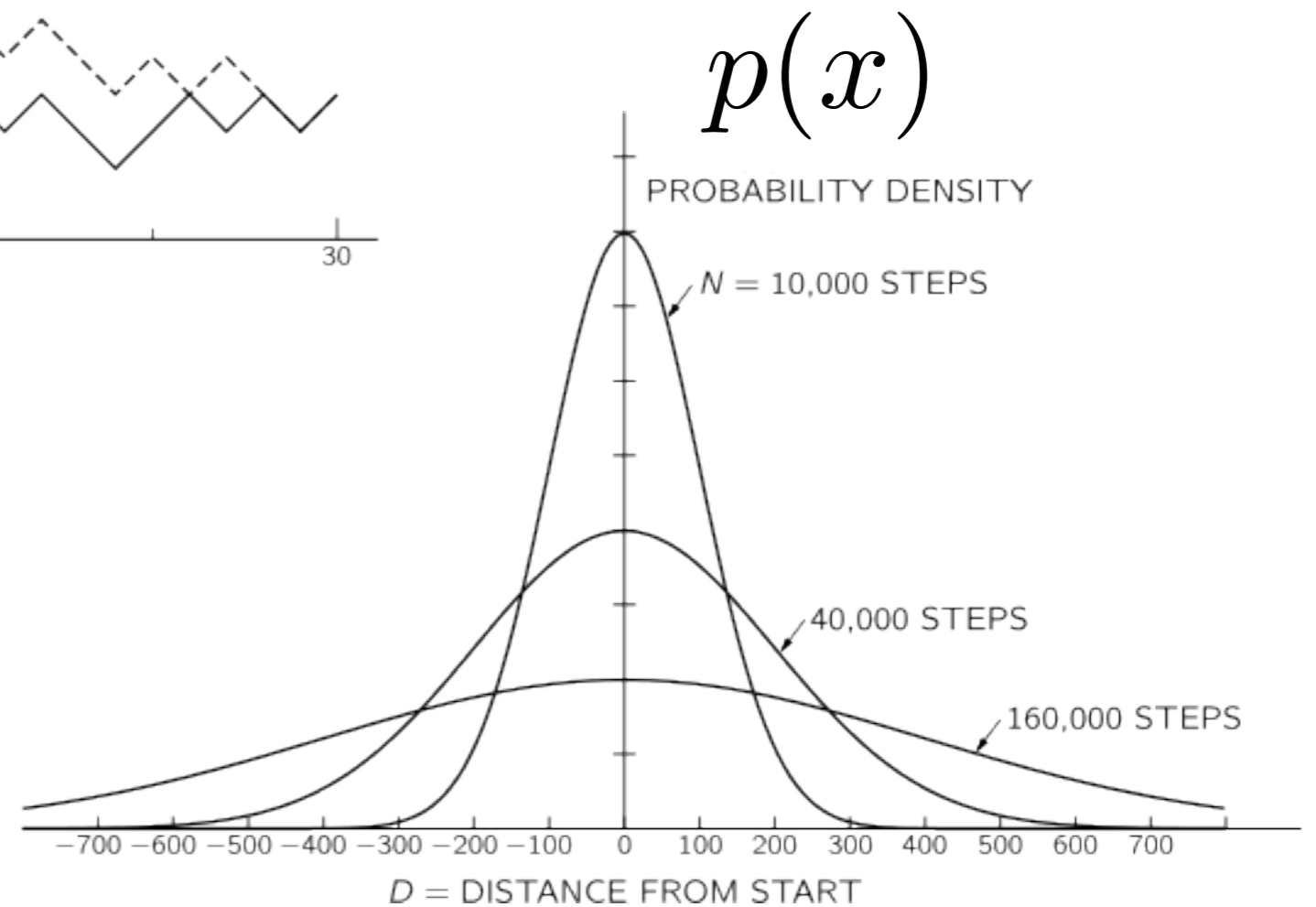
		states	
		discrete	continuous
time	discrete	$\begin{array}{c} dx \rightarrow 0 \\ \leftarrow \\ \int \rightarrow \Sigma \end{array}$	$\int \rightarrow \Sigma \begin{array}{c} \uparrow \\ \downarrow \end{array} dt \rightarrow 0$
	continuous		Lots of gnarly math!

# Random walks



‘paths diffuse over time’

‘density’



# Expectations over paths

$$\Psi_{t_i} = E_{\tau_i} \left( \Psi_{t_N} e^{-\int_{t_i}^{t_N} \frac{1}{\lambda} q_t dt} \right) = E_{\tau_i} \left[ \exp \left( -\frac{1}{\lambda} \phi_{t_N} - \frac{1}{\lambda} \int_{t_i}^{t_N} q_t dt \right) \right]$$

forward!

... but stochastic

$$\int p(\tau) \exp \left( -\frac{1}{\lambda} \phi - \frac{1}{\lambda} \int q dt \right) d\tau$$

$$\tau = x(t \dots t_N) \sim p(x, u)$$

an instance of a random path segment (a random 'number', but in spaces of functions)

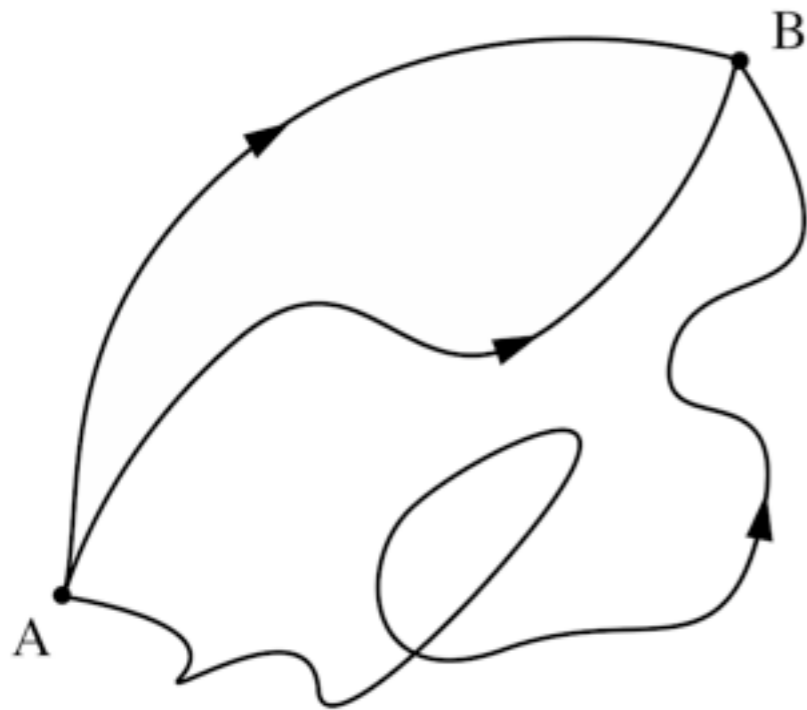
$$E[X] = \int x p(x) dx$$

Continuous time,  $x$  is function of time  $x = f(t)$

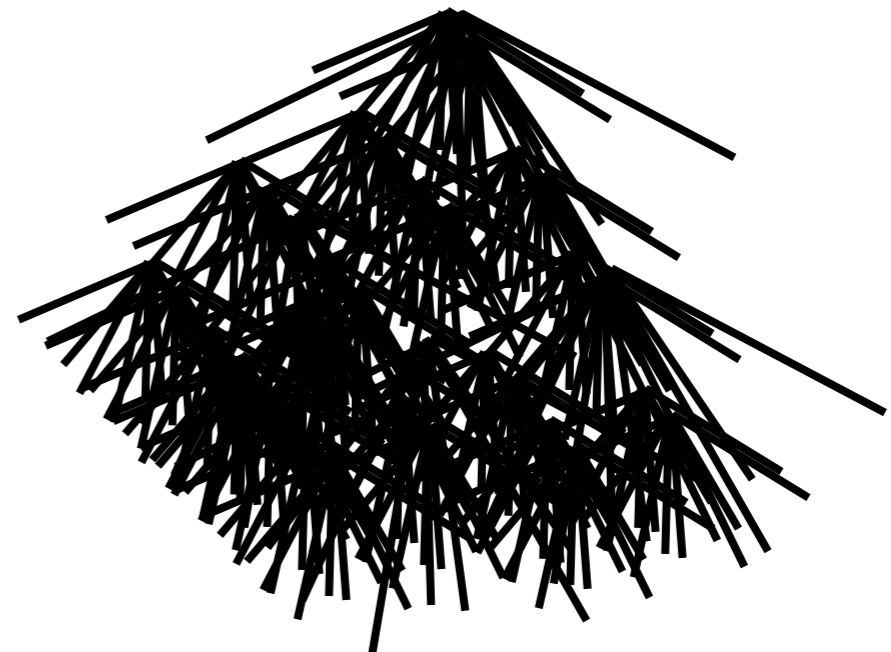


# Comparison to graphs

can think of all possibilities of a random walk as graph



There are several ways to end up in a certain state, each path has an associated probability



When does 'branching' occur?

Idea: do discrete time and take limit



$$dt \rightarrow 0$$

# Continuous decision processes

Take random walk and take limits

$dx \rightarrow 0$       probability densities

$dt \rightarrow 0$       probability flow

for all times  $\int p(x) dx = 1$       conservation law!

Conserved flow?

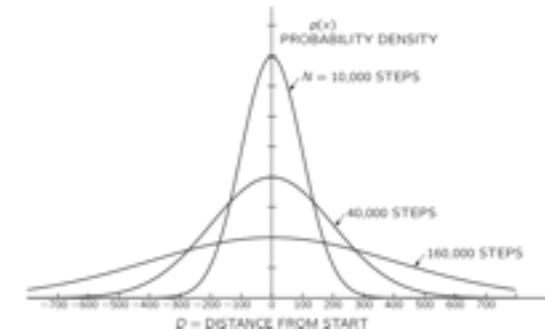
You know how to do that!



# Pre-requisite I: Brownian motion

Assume process with probability distribution

$$\mathbb{P}_{\mathbf{w}}(t, w) = \frac{1}{\sqrt{2\pi\sigma^2 t}} \exp\left(-\frac{(w - \mu t)^2}{2\sigma^2 t}\right)$$



at any time:

$$\mathbb{E}\{w(t)\} = \mu t$$

$$\text{Var}\{w(t)\} = \sigma^2 t$$

Defined via increment process:

$$dw(t) = \lim_{\Delta t \rightarrow 0} w(t + \Delta t) - w(t)$$

1. The increment process,  $dw(t)$ , has a Gaussian distribution with the mean and the variance,  $\mu\Delta t$  and  $\sigma^2\Delta t$  respectively.
2. The increment process,  $dw(t)$ , is statistically independent of  $w(s)$  for any  $s \leq t$ .



# Simulate Brownian Motion

Integrate discretized increment process:

$$w(t + \Delta t) = w(t) + \mu\Delta t + \sqrt{\Delta t}\sigma\varepsilon, \quad w(0) = 0$$

$$\varepsilon \sim N(0, 1)$$

... this is a discretized SDE

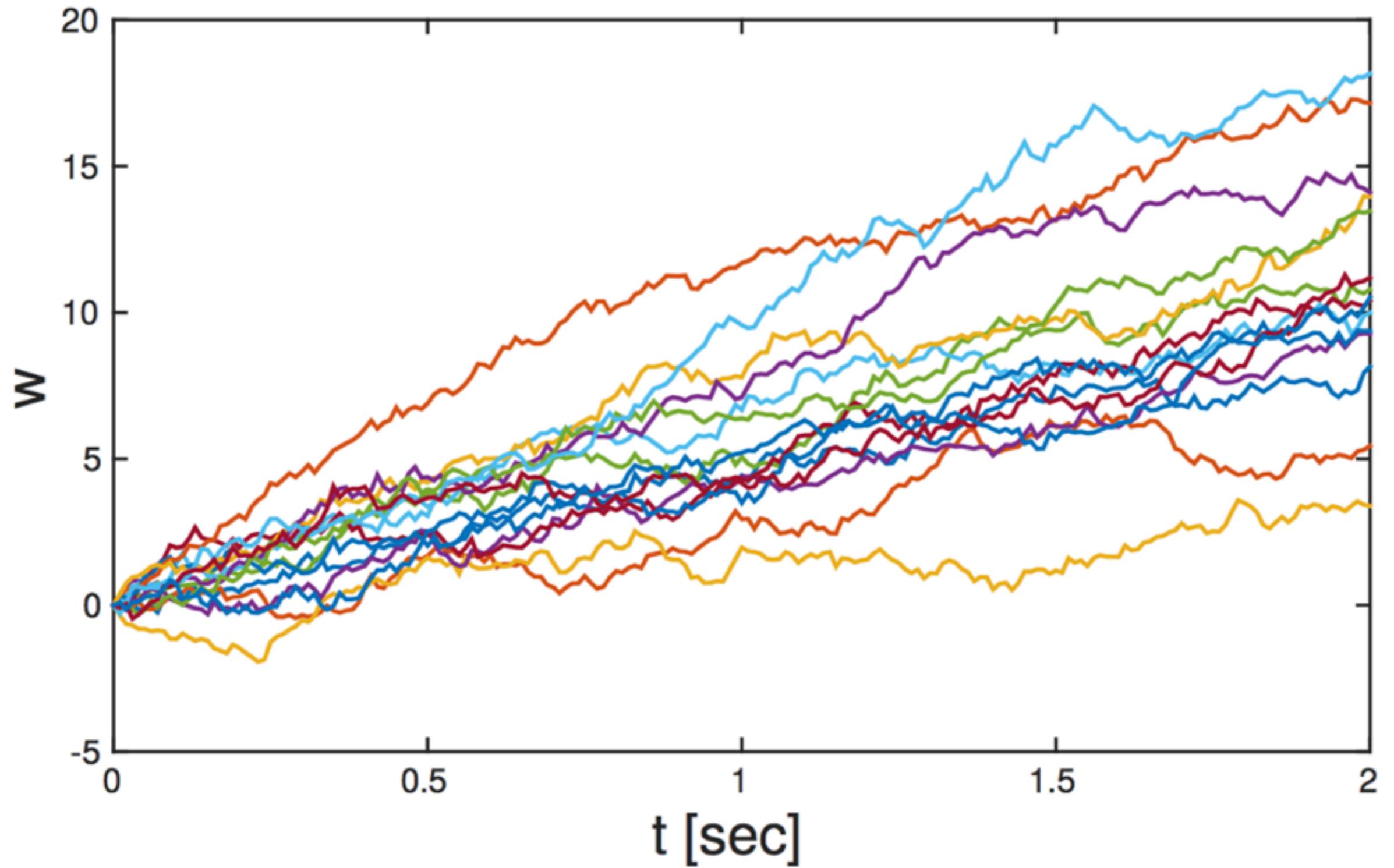


Figure 3.1: Brownian Motion with  $\mu = 5$  and  $\sigma^2 = 4$



# Pre-requisite 2: Stochastic Differential Equations (SDE)

SDE: 
$$d\mathbf{x} = \mathbf{f}(t, \mathbf{x})dt + \mathbf{G}(t, \mathbf{x})d\mathbf{w}$$

$\mathbf{w}(t)$

Brownian motion (Wiener Process)  
(zero mean, covariance =  $I$ )

# SDE Integration

$$d\mathbf{x} = \mathbf{f}(t, \mathbf{x})dt + \mathbf{G}(t, \mathbf{x})d\mathbf{w}$$

Numerical integration:

use:      small time step  $\Delta t$        $d\mathbf{w} = \sqrt{\Delta t}\boldsymbol{\varepsilon}$

assumes constant increment:  $w(t + \Delta t) = w(t) + \mu\Delta t + \sqrt{\Delta t}\sigma^2\varepsilon, \quad w(0) = 0$

$$\mathbf{x}(t_{n+1}) = \mathbf{x}(t_n) + \mathbf{f}(t_n, \mathbf{x}(t_n))\Delta t + \mathbf{G}(t_n, \mathbf{x}(t_n))\sqrt{\Delta t}\boldsymbol{\varepsilon}, \quad \boldsymbol{\varepsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$$

Process is nonlinear through  $\mathbf{f}(t, \mathbf{x})$ : Non-gaussian

But, for  $\Delta t \rightarrow 0$

$$\mathbb{P}_{\mathbf{x}}(t + \Delta t, \mathbf{x} \mid t, \mathbf{y}) = \mathcal{N}\left(\mathbf{y} + \mathbf{f}(t, \mathbf{y})\Delta t, \mathbf{G}(t, \mathbf{y})\mathbf{G}^T(t, \mathbf{y})\Delta t\right)$$



# Probabilistic Dynamics

Discrete time:  
Markov chains

☞ Master Equation

Continuous time:

Jumps: Continuous-time Markov chain

Smooth: Markov Process

☞ Fokker-Planck



# Pre-requisite 3: The Fokker-Planck Equation

Time evolution of probability distribution?  
(Not a Gaussian process)

Probability distribution is solution of a  
nonlinear PDE: The Fokker-Planck Equation



# Fokker-Planck Equation

## formal definitions

Assuming stochastic process:  $d\mathbf{x} = \mathbf{f}(t, \mathbf{x})dt + \mathbf{G}(t, \mathbf{x})d\mathbf{w}$

$\mathbb{P}_{\mathbf{x}(t)}(t = s, \mathbf{x} \mid \tau, \mathbf{y})$  conditional probability distribution of  $\mathbf{x}(t)$

given that process at initial time  $s$  has value  $\mathbf{y} = \mathbf{x}(s)$

Time evolution of prob. distribution:

$$\partial_t \mathbb{P} = -\nabla_x^T (\mathbf{f} \mathbb{P}) + \frac{1}{2} \text{Tr} [\nabla_{xx} (\mathbf{G} \mathbb{P})]$$



$$\nabla_x() = \begin{bmatrix} \frac{\partial}{\partial x_1} \\ \frac{\partial}{\partial x_2} \\ \vdots \\ \frac{\partial}{\partial x_n} \end{bmatrix}$$

$$\nabla_{xx}() = \begin{bmatrix} \frac{\partial}{\partial x_{11}} & \cdots & \frac{\partial}{\partial x_{1n}} \\ \vdots & \ddots & \vdots \\ \frac{\partial}{\partial x_{n1}} & \cdots & \frac{\partial}{\partial x_{nn}} \end{bmatrix}$$

# Fokker-Planck Equation

PDE for time evolution of probability distribution

$$\frac{\partial}{\partial t} p(x, t) = - \underbrace{\frac{\partial}{\partial x} [\mu(x, t) p(x, t)]}_{\text{Drift}} + \frac{\partial^2}{\partial x^2} \underbrace{[D(x, t) p(x, t)]}_{\text{Diffusion}}$$

$$dx = \mathbf{f}(x, t) dt + \mathbf{G}(x) d\omega$$

brownian motion, no drift  $\mathbf{f}(\mathbf{x}_t, t) dt = 0$   $\mathbf{G}(\mathbf{x}_t) = 1$

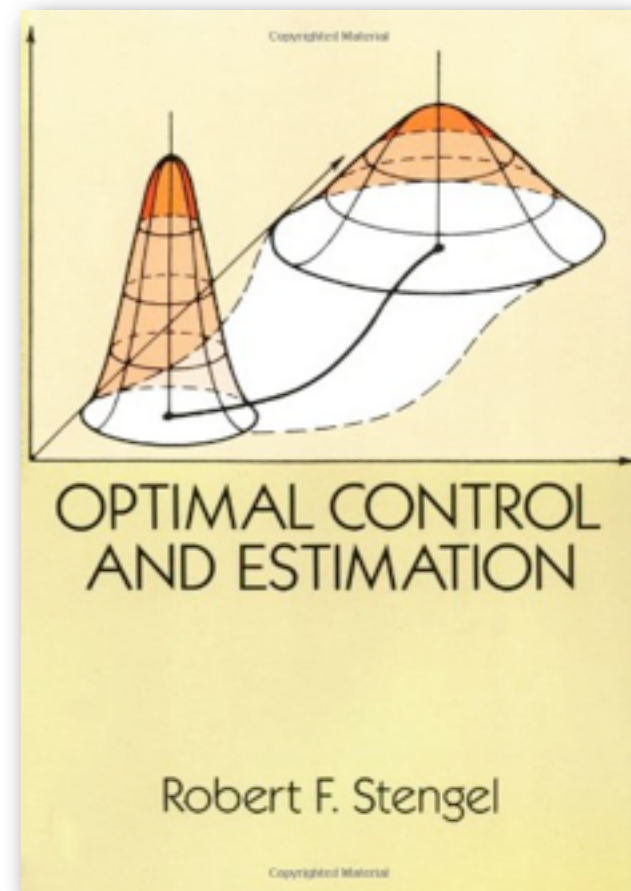
$$d\mathbf{x} = d\omega$$

$$\Rightarrow \frac{\partial p(x, t)}{\partial t} = \frac{1}{2} \frac{\partial^2 p(x, t)}{\partial x^2}$$

$$\Rightarrow p(x, t) = \frac{1}{\sqrt{2\pi t}} e^{-\frac{x^2}{2t}}$$

conservation law!

$$\int p(x) dx = 1$$



[BB EXAMPLE]



# Fokker-Planck Equation

PDE for time evolution of probability distribution

$$\frac{\partial}{\partial t} p(x, t) = \underbrace{-\frac{\partial}{\partial x} [\mu(x, t)p(x, t)]}_{\text{Drift}} + \underbrace{\frac{\partial^2}{\partial x^2} [D(x, t)p(x, t)]}_{\text{Diffusion}}$$

cf. Fluid Dynamics  
Heat and Charge diffusion  
cf. Particle filters



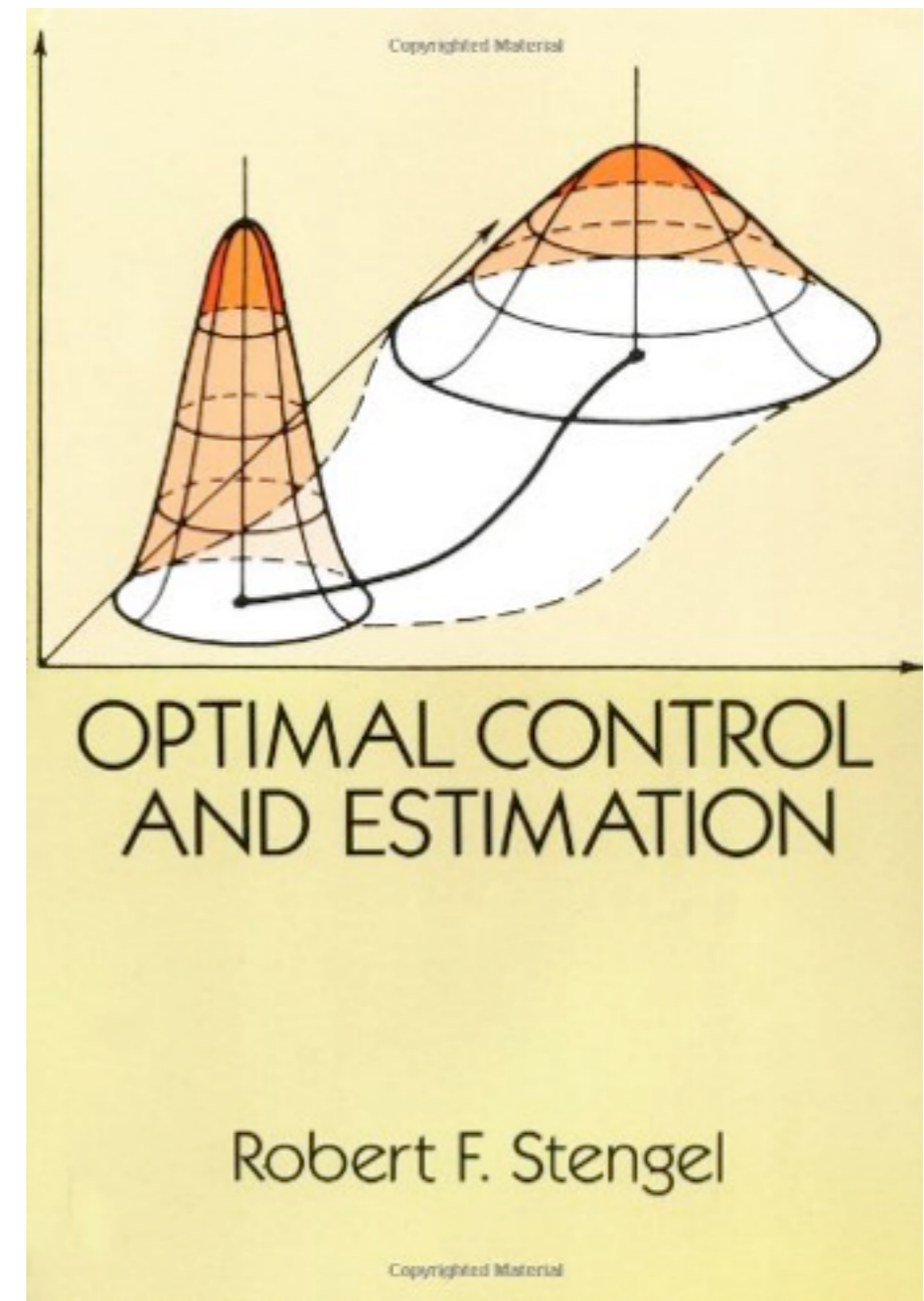
Finance, Biology, Chemistry, Physics, Sociology,  
Anthropology, Control & Machine Learning





# Stochastic Control

## 'Controlled Diffusion' Controlled Brownian Motion



Fokker-Planck modeling is conceptually very useful

... e.g. to develop algorithms

But, naively applying the concept is often not practical



# Linearly-Solvable Markov Decision Process

Class of stochastic optimal control problems for  
which HJB is linear (in Value)



# From SOC to LMDP

Assume SDE with control/noise affine form

$$d\mathbf{x} = \mathbf{f}(t, \mathbf{x})dt + \mathbf{g}(t, \mathbf{x}) (\mathbf{u}dt + d\mathbf{w}), \quad d\mathbf{w} \sim \mathcal{N}(\mathbf{0}, \Sigma dt)$$

divide by  $dt$       substitute  $\frac{d\mathbf{w}}{dt}$  by  $\varepsilon$

$$\dot{\mathbf{x}} = \mathbf{f}(t, \mathbf{x}) + \mathbf{g}(t, \mathbf{x}) (\mathbf{u} + \varepsilon), \quad \varepsilon \sim \mathcal{N}(\mathbf{0}, \Sigma)$$

**Cost**       $J = E \left\{ \Phi(\mathbf{x}(t_f)) + \int_{t_0}^{t_f} q(t, \mathbf{x}) + \frac{1}{2} \mathbf{u}^T \mathbf{R} \mathbf{u} dt \right\}$

**Not LMDP (nonlinear HJB)**



# Towards linear HJB

of control affine SOC problem

Use HJB for general stoch. cont. time control problem

$$\beta V^*(t, \mathbf{x}) - V_t^*(t, \mathbf{x}) = \min_{\mathbf{u}(t)} \left\{ L(\mathbf{x}, \mathbf{u}(t)) + V_{\mathbf{x}}^{*T} \mathbf{f}_t(\mathbf{x}, \mathbf{u}(t)) + \frac{1}{2} \text{Tr}[V_{\mathbf{xx}}^* \mathbf{B}(t) \mathbf{W}(t) \mathbf{B}^T(t)] \right\}$$

Substitute (assumptions/definitions):

$$\beta \leftarrow 0$$

$$L(\mathbf{x}, \mathbf{u}(t)) \leftarrow q(t, \mathbf{x}) + \frac{1}{2} \mathbf{u}^T \mathbf{R} \mathbf{u}$$

$$\mathbf{f}_t(\mathbf{x}, \mathbf{u}(t)) \leftarrow \mathbf{f}(t, \mathbf{x}) + \mathbf{g}(t, \mathbf{x}) \mathbf{u}$$

$$\mathbf{B}(t) \leftarrow \mathbf{g}(t, \mathbf{x})$$

$$\mathbf{W}(t) \leftarrow \Sigma$$

$$-\partial_t V^*(t, \mathbf{x}) = \min_{\mathbf{u}} \left\{ q(t, \mathbf{x}) + \frac{1}{2} \mathbf{u}^T \mathbf{R} \mathbf{u} + \nabla_x^T V^*(t, \mathbf{x}) (\mathbf{f}(t, \mathbf{x}) + \mathbf{g}(t, \mathbf{x}) \mathbf{u}) + \frac{1}{2} \text{Tr}[\nabla_{xx} V^*(t, \mathbf{x}) \mathbf{g}(t, \mathbf{x}) \Sigma \mathbf{g}^T(t, \mathbf{x})] \right\}$$



# HJB Equation

of control affine opt. ctrl problem

$$-\partial_t V^*(t, \mathbf{x}) = \min_{\mathbf{u}} \left\{ q(t, \mathbf{x}) + \frac{1}{2} \mathbf{u}^T \mathbf{R} \mathbf{u} + \nabla_x^T V^*(t, \mathbf{x}) (\mathbf{f}(t, \mathbf{x}) + \mathbf{g}(t, \mathbf{x}) \mathbf{u}) + \frac{1}{2} \text{Tr}[\nabla_{xx} V^*(t, \mathbf{x}) \mathbf{g}(t, \mathbf{x}) \Sigma \mathbf{g}^T(t, \mathbf{x})] \right\}$$

gradient of RHS = 0 yields

$$\mathbf{u}^*(t, \mathbf{x}) = -\mathbf{R}^{-1} \mathbf{g}^T(t, \mathbf{x}) \nabla_x V(t, \mathbf{x})$$

# Optimal control

‘improve value’  
negative gradient

$$\mathbf{u}^*(t, \mathbf{x}) = -\mathbf{R}^{-1} \mathbf{g}^T(t, \mathbf{x}) \nabla_x V(t, \mathbf{x})$$

‘coord. trafo.’  
‘states to control’

‘Cost of controls to get improvement’

# Optimal HJB

substitute opt. control back into HJB  $\Rightarrow$

$$\mathbf{u}^*(t, \mathbf{x}) = -\mathbf{R}^{-1} \mathbf{g}^T(t, \mathbf{x}) \nabla_{\mathbf{x}} V(t, \mathbf{x})$$

$$-\partial_t V^* = q - \frac{1}{2} \nabla_{\mathbf{x}}^T V^* \mathbf{g} \mathbf{R}^{-1} \mathbf{g}^T \nabla_{\mathbf{x}} V^* + \nabla_{\mathbf{x}}^T V^* \mathbf{f} + \frac{1}{2} \text{Tr}[\nabla_{\mathbf{x}\mathbf{x}} V^* \mathbf{g} \Sigma \mathbf{g}^T]$$

Nonlinear PDE!



# Optimal HJB

substitute opt. control back into HJB  $\Rightarrow$

$$\mathbf{u}^*(t, \mathbf{x}) = -\mathbf{R}^{-1} \mathbf{g}^T(t, \mathbf{x}) \nabla_{\mathbf{x}} V(t, \mathbf{x})$$

$$-\partial_t V^* = q - \frac{1}{2} \nabla_{\mathbf{x}}^T V^* \mathbf{g} \mathbf{R}^{-1} \mathbf{g}^T \nabla_{\mathbf{x}} V^* + \nabla_{\mathbf{x}}^T V^* \mathbf{f} + \frac{1}{2} \text{Tr}[\nabla_{\mathbf{x}\mathbf{x}} V^* \mathbf{g} \Sigma \mathbf{g}^T]$$

Nonlinear PDE!

short notation, replace:  $\mathbf{g} \mathbf{R}^{-1} \mathbf{g}^T$  by  $\Xi$

add'l assumption: control cost linked to noise

$$\mathbf{R} \Sigma = \lambda \mathbf{I}$$

$$\Rightarrow \mathbf{g} \Sigma \mathbf{g}^T = \lambda \Xi$$



$$\lambda > 0$$

# log transform

$$-\partial_t V^* = q - \frac{1}{2} \nabla_x^T V^* \Xi \nabla_x V^* + \nabla_x^T V^* \mathbf{f} + \frac{\lambda}{2} \text{Tr}[\nabla_{xx} V^* \Xi]$$

Nonlinear PDE!

Desirability  $\Psi$

Log transform

$$V^*(t, \mathbf{x}) = -\lambda \log \Psi(t, \mathbf{x})$$

$\Rightarrow$

$$\begin{aligned} \partial_t V^*(t, \mathbf{x}) &= -\lambda \frac{\partial_t \Psi}{\Psi} \\ \nabla_x V^*(t, \mathbf{x}) &= -\lambda \frac{\nabla_x \Psi}{\Psi} \\ \nabla_{xx} V^*(t, \mathbf{x}) &= \frac{1}{\lambda} \nabla_x V^* \nabla_x^T V^* - \lambda \frac{\nabla_{xx} \Psi}{\Psi} \end{aligned}$$

rewrite scalar expression:

$$-\frac{1}{2} \nabla_x^T V^* \Xi \nabla_x V^* = -\frac{1}{2} \text{Tr}[\nabla_x^T V^* \Xi \nabla_x V^*] = -\frac{1}{2} \text{Tr}[\nabla_x V^* \nabla_x^T V^* \Xi]$$

substitute into HJB:

$$\lambda \frac{\partial_t \Psi}{\Psi} = q - \cancel{\frac{1}{2} \text{Tr}[\nabla_x V^* \nabla_x^T V^* \Xi]} - \lambda \mathbf{f}^T \frac{\nabla_x \Psi}{\Psi} + \cancel{\frac{1}{2} \text{Tr}[\nabla_x V^* \nabla_x^T V^* \Xi]} - \frac{\lambda^2}{2} \text{Tr} \left[ \frac{\nabla_{xx} \Psi}{\Psi} \Xi \right]$$



multiply both sides by  $-\Psi/\lambda$

# Linear HJB for desirability

$$-\partial_t \Psi = -\frac{1}{\lambda} q \Psi + \mathbf{f}^T \nabla_x \Psi + \frac{\lambda}{2} \text{Tr}[\mathbf{\Xi} \nabla_{xx} \Psi]$$

Linear PDE in desirability  $\Psi$

equivalent form:

$$\begin{aligned} -\partial_t \Psi &= \mathbf{H}[\Psi] & \mathbf{H} &= -\frac{1}{\lambda} q + \mathbf{f}^T \nabla_x + \frac{\lambda}{2} \text{Tr}[\mathbf{\Xi} \nabla_{xx}] \\ & & &= -\frac{1}{\lambda} q + \sum_i \mathbf{f}_i \frac{\partial}{\partial x_i} + \frac{\lambda}{2} \sum_{i,j} \mathbf{\Xi}_{ij} \frac{\partial^2}{\partial x_i \partial x_j} \end{aligned}$$

linear, but still no analytic solution for arbitrary  $q(\mathbf{x}, t)$



could solve **backward**

terminal condition  $\Psi(t_f, \mathbf{x}) = \exp\left(-\frac{1}{\lambda} \Phi(\mathbf{x})\right)$

# Forward solution through diffusion process

$$-\partial_t \Psi = -\frac{1}{\lambda} q \Psi + \mathbf{f}^T \nabla_x \Psi + \frac{\lambda}{2} \text{Tr}[\mathbf{\Xi} \nabla_{xx} \Psi]$$

Can solve this equation using 'forward diffusion process'

$$\Psi_{t_i} = E_{\tau_i} \left( \Psi_{t_N} e^{-\int_{t_i}^{t_N} \frac{1}{\lambda} q_t dt} \right) = E_{\tau_i} \left[ \exp \left( -\frac{1}{\lambda} \phi_{t_N} - \frac{1}{\lambda} \int_{t_i}^{t_N} q_t dt \right) \right]$$

'expectation in respect to path  $\tau$ '

path drawn from random process

$$\tau = x(t \dots t_N) \sim p(x, u)$$

$$d\mathbf{x} = \mathbf{f}(t, \mathbf{x}) dt + \mathbf{g}(t, \mathbf{x}) d\mathbf{w} \quad d\mathbf{w} \sim N(0, \Sigma)$$

forward!

... but stochastic

Remember the forward search in the discrete state, discrete time problem (Lect. I)



# Credits

material from:

Sutton & Barto's book: [http://  
webdocs.cs.ualberta.ca/~sutton/book/the-  
book.html](http://webdocs.cs.ualberta.ca/~sutton/book/the-book.html)

Bishop: Pattern Recognition and Machine Learning

Feynman Lectures

