

Jan Peters Gerhard Neumann







"learning to do an act from seeing it done"

(Thorndike, 1898)

In ALVINN: "reproducing the observed steering action for a given retinal image"









Purpose of this Lecture

Learning From Demonstrations

 How can we teach a robot without programming?

Policy Representations

- Show you important characteristics of commonly used policies
- State space vs. trajectory space view

Introduce the concept of Movement Primitives:

- How can we incorporate modularity?
- Data-driven acquisition of movements





1. Learning Policies from Demonstrations by Supervised Learning

- 2. Policy Representations
 - State-space representations
 - Trajectory-based representations
- 3. Imitation Learning with Movement Primitives
 - Dynamic Movement Primitives
 - Probabilistic Movement Primitives
 - Beyond a single primitive



- Very successful strategy for humans
- Learning controllers from scratch by reinforcement learning is often very time consuming or even too difficult
- the search space may frequently be way too large for the agent to explore it in its lifetime
- an expert takes many years to optimize his policy and a robot could avoid his expensive training by cloning his policy

Already rats can imitate!





- Student rats observe companion actor rats performing different spatial tasks differing according to the experimental requirements.
- After the observational training, surgical ablation to block any further learning in the student rat.
- The observer rats displayed exploration abilities that closely matched the previously observed behaviors.

Legio et al, *Brain Res. Protocols, 2003* Heyes, *Trends in Cog. Sciences,* 2001



... and dolphins...





 Infants as young as 42 minutes old copy several facial actions (e.g., Meltzoff & Moore, 1977).





- **Teleoperation**: Use a joystick to train an RC car, a mouse for training a Quake III player, the steering wheel of the Navlab, data gloves, etc.
- **Kinesthetic Teach-In**: Take the robot by the hand like a tennis teacher teaches a tennis student.
- Vision: Video-based tracking of human beings.
- Marker-based Tracking: With markers and a basic skeleton, very precise human data can be obtained.
- Sensuits: Suits with encoders and accelerometers attached to human beings.



- Behavioral Cloning is the simplest form of learning from demonstration
- An expert is available and supplies data traces:

$$m{s}_1 o m{u}_1 o m{s}_2 o m{u}_2 o m{s}_3 o m{u}_3 o$$

• In our case, often $\, s = \left[egin{array}{c} q \ \dot{q} \end{array}
ight]$

• The student infers a policy from these data traces, i.e.,

$$\boldsymbol{u} = \pi(\boldsymbol{s}) \text{ or } \boldsymbol{u} \sim \pi(\boldsymbol{u}|\boldsymbol{s})$$

• In principle, this can be treated as a supervised learning problem.



Standard ML techniques can simply be applied to the data set

$$\mathcal{D} = \{oldsymbol{s}_i,oldsymbol{u}_i\}$$

to extract a **policy**

$$u = \pi(s) = \phi^T(s)\theta, \text{ or}$$
$$u \sim \pi(u|s) = \mathcal{N}(u|\mu(s), \Sigma(s))$$

... the problem frequently boils down to a regression problem.

The clean-up effect: due to regularization, the noise in the demonstration is no longer exhibited by the reproduction and, hence, the clone often surpasses the quality of the expert.



1. Learning Policies from Demonstrations by Supervised Learning

- 2. Policy Representations
 - State-space representations
 - Trajectory-based representations
- 3. Imitation Learning with Movement Primitives
 - Dynamic Movement Primitives
 - Probabilistic Movement Primitives
 - Beyond a single primitive

Why do we use parametric policies?



A parametric policy is a conditional probability distribution $\pi(u|s;\theta)$ that chooses the actions u depending on the state s of the robot

- Parametric policy naturally incorporates continuous actions
- Estimate from demonstration / imitation learning

Generalize to unseen situations

• Search for improved parameters / reinforcement learning





What are desirable properties?



- Compactness: Low number of parameters
- Learn-ability: Easy to learn from demonstration and by reinforcement learning
- Stochasticity: Can encode exploration and variability
- Optimality: Can encode optimal behavior?
- Scalability: Can be used for a high number of DoFs?
- Modularity:

Adaptability: Reusable for new situations?

Co-activation and Blending of movements

Useable for stroke-based and rhythmic movements
 16





Why use a stochastic policy?

Used for exploration in reinforcement learning (later)

Can also capture variability of movements

Exploration models:

 $\boldsymbol{u} = \pi(\boldsymbol{s}) = f_{\boldsymbol{w}}(\boldsymbol{s}), \quad \boldsymbol{\theta} = \boldsymbol{w}$ No exploration:

Uncorrelated Exploration: $\pi(\boldsymbol{u}|\boldsymbol{s};\boldsymbol{\theta}) = \mathcal{N}(\boldsymbol{u}|f_{\boldsymbol{w}}(\boldsymbol{s}), \sigma^2 \boldsymbol{I}), \quad \boldsymbol{\theta} = \{\boldsymbol{w}, \sigma^2\}$

Correlated Exploration: $\pi(\boldsymbol{u}|\boldsymbol{s};\boldsymbol{\theta}) = \mathcal{N}(\boldsymbol{u}|f_{\boldsymbol{w}}(\boldsymbol{s}),\boldsymbol{\Sigma}), \quad \boldsymbol{\theta} = \{\boldsymbol{w},\boldsymbol{\Sigma}\}$



We also have to learn the variances of the linear models



Exploration might also hurt

State space vs. trajectory space representations



- Policy depends on the state and on the parameters
- Represents a globally valid policy
- Complex non-linear representations are needed

Examples:

- Neural Networks
- RBF Networks
- Gaussian Processes
- Locally Weighted Regression Models





State space representations

Linear controllers: $f_{\boldsymbol{w}}(\boldsymbol{s}) = \phi^T(\boldsymbol{s})\boldsymbol{w}$

Most simple case: linear PD controller

$$\phi(\boldsymbol{s}) = \begin{bmatrix} 1 \\ \boldsymbol{s} \end{bmatrix}, \quad f_{\boldsymbol{w}}(\boldsymbol{s}) = \boldsymbol{K}\boldsymbol{s} + \boldsymbol{k}$$

[-] Good feature representation needs to be known

[+] Very compact representation (low number of parameters)

[+] Easy to learn (linear regression)

Pole Balancing



Widrow and Smith (1964) used supervised learning to acquire a pole balancing policy.



Widrow and Smith (1964) used supervised learning to acquire a pole balancing policy.



Pole Balancing



Sammut's Cessna Pilot







Non-linear state space representations

Radial Basis Function (RBF) networks:

$$f_{\boldsymbol{w}}(\boldsymbol{s}) = \phi^T(\boldsymbol{s})\boldsymbol{\beta}, \quad \phi_i(\boldsymbol{s}) = \exp\left(-0.5\sum_{j=1}^D (s_j - \mu_{ij})^2 / h_{ij}\right)$$

$$m{w} = \{m{eta}, m{\mu}_{1:K}, m{h}_{1:K}\}$$

Normalized RBF:

$$f_{\boldsymbol{w}}(\boldsymbol{s}) = \frac{\sum_{i=1}^{K} \phi_i(\boldsymbol{s})\beta_i}{\sum_{i=1}^{K} \phi_i(\boldsymbol{s})}$$

[-] A high number of parameters

[-] Non-convex optimization

[-] Hard to scale Curse of dimensionality

[+] Automatic feature construction



Alternatives: Gaussian Mixture Models (GMM), Neural Networks

ALVINN & Navlab in 1989-1995!





Pomerlau (1989-1995)



No-Hands-Across-America



ALVINN allowed the Navlab vehicle of CMU's robotics institute to drive 2796km autonomously as part of their 'No-Hands-Across-America' Tour in 1995. 24





State: Camera Image



Action: Steering Wheel,

Brakes, Gas



Intelligence



Function Approximator:

A Two-Layered Neural Network



Intelligence



Function Approximator:

A Two-Layered Neural Network

JUST (nonlinear) REGRESSION!





Video from ALVIN

27



State space representations

for walking



Represent controller in a low-dimensional manifold



(Grimes, Rashid, & Rao, NIPS 2007)

Doubts on Direct Behavioral Cloning from State-Action Pairs



- It becomes brittle for **larger state-spaces** unless you have a taskappropriate representation.
- Frequently leads to **catastrophic failures** if the controller has not been trained in this area of the state-action space (Sammut, 2010) or if there have been small changes in the system (Camacho & Michie, 1995).
- Reproduction of single human teachers always works best (Camacho & Michie, 1995).
- There is no guarantee that the reproduction is meaningful, nor an interpretation of behavior.
- The data is treated as if it was i.i.d.
- We do not know whether we can also reproduce the long-term behavior!
- Only learning of **individual motions is "easy".**

State space vs. trajectory space representations

Time-dependent representation: $\pi(\boldsymbol{u}|\boldsymbol{s},\boldsymbol{t};\boldsymbol{\theta})$

Policy also depends on time, e.g., follow a specific trajectory

For the same time step, the robot is often in similar states

Simple local models are often sufficient!



Time-dependend representations

For example: Time-dependent linear feedback controllers

$$f_{\boldsymbol{w}}(\boldsymbol{s},t) = \sum_{i=1}^{K} \phi_i(t) (\boldsymbol{K}_i \boldsymbol{s} + \boldsymbol{k}_i)$$

- Time dependent basis functions, e.g., normalized RBF functions
- Scales quadratically with # DoF D: KD/2(D+1)
- Equivalent to PD-trajectory tracking with time-varying controller gains
 - Variable stiffness controllers
- Locally optimal representation (why we will see in the next lectures!)



Trajectory Generators:

Directly learn desired trajectory $oldsymbol{q}^*(t;oldsymbol{w})$

Use feedback controller to follow trajectory

$$\pi(\boldsymbol{s},t;\boldsymbol{w}) = \boldsymbol{K}_P(\boldsymbol{q}^*(t;\boldsymbol{w}) - \boldsymbol{q}) + \boldsymbol{K}_D(\dot{\boldsymbol{q}}^*(t;\boldsymbol{w}) - \dot{\boldsymbol{q}})$$

where typically $oldsymbol{K}_P$ and $oldsymbol{K}_D$ are hand tuned diagonal matrices

Possible Trajectory Representations:

- Splines $q^*(t; w) = \sum_{i=0}^5 w_i t^i$
- Linear basis function models (RBFs) $\boldsymbol{q}^*(t; \boldsymbol{w}) = \boldsymbol{\phi}^T(t) \boldsymbol{w}$
- Dynamical Systems



- 1. Learning Policies from Demonstrations by Supervised Learning
- 2. Policy Representations
 - State-space representations
 - Trajectory-based representations
- **3. Imitation Learning with Movement Primitives**
 - Dynamic Movement Primitives
 - Probabilistic Movement Primitives
 - Beyond a single primitive



What are movement primitives?

- Movement primitives are a compact representation of a movement
- Often represented as parametrized trajectory generator

 $\boldsymbol{\tau} = f(\boldsymbol{w}), \quad \boldsymbol{w} \dots$ parameters of the primitive

Imitation Learning with trajectory generators

- By learning the desired trajectory, we also learn the desired long term behavior!
- However, we still have to learn how to follow this trajectory
- ➡ If we do not have good trajectory tracking controllers, it does not work
Dynamical systems as Trajectory Generators

Dynamical systems can be used to represent trajectories

Integrating the dynamical system results in a trajectory

• **Example:** First order linear dynamical system:



 $\dot{y} = \alpha(c-y)$

What movements can a differential equation encode?



Second order linear dynamical system

$$\ddot{y} = \alpha(\beta(c-y) - \dot{y})$$

Linear differential equations:

- well-defined behavior
- But: limited class of movements



How can we make it more representative?



$$\dot{y} = -2\cos x - \cos y$$
$$\dot{x} = -2\cos y - \cos x$$

Repeller / ddles >0 -3 -2 -1 2 3 х

Use non-linear dynamical systems ?

- Can represent more complex behavior
- Can also get unstable! ☺



Different behaviors might emerge...



Attractors

Limit cycles

Chaos

Movements as dynamical systems











Dynamic Movement Primitives (DMPs)

We can encode desirable properties such as:

- stability
- perturbation robustness
- periodic and point-to-point behaviors
- Attractors that have rather complex shape
- Easy to learn
- Coupling of a high number of DoFs
- Timing, temporal scaling
- Generalization (structural equivalence for parameter changes)

Point-to-Point Movements as Dynamic Systems





E.g., for a one degree-of-freedom movement, start with a simple damped spring model

$$- \begin{bmatrix} \mathbf{x} \\ \mathbf{y} \end{bmatrix} = \dot{y} = \alpha(\beta(g - y) - \dot{y})$$

(Ijspeert, Nakanishi & Schaal, NIPS 2003; Schaal, Peters, Nakanishi, Ijspeert, ISRR 2003)



How can we encode a desired behavior?

Add a forcing function to obtain a moving attractor

$$\ddot{y} = \alpha(\beta(g - y) - \dot{y}) + f_{\boldsymbol{w}}(t)$$
$$= \alpha(\beta(g + f_{\boldsymbol{w}}(t)/(\alpha\beta) - y) - \dot{y})$$

- The forcing function encodes the desired additional acceleration profile
 - $f_{\boldsymbol{w}}$... learnable function



How can we encode a temporal scaling?

Add a phase variable z_t to replace time $\ddot{y} = \tau^2 \alpha (\beta(g - y) - \dot{y}/\tau) + \tau^2 f_w z$ $\dot{z} = -\tau \alpha_z z$

Also uses dynamical system to model phase z

 τ ... temporal scaling variable



Adapting the temporal scaling...



higher τ has ner movement speed



How to represent f?

Normalized RBF basis functions
$$\phi_i(z) = \exp(-0.5(z - c_i)^2/h_i)$$

$$f_{\boldsymbol{w}}(z) = \frac{\sum_{i=1}^{K} \phi_i(z) w_i z}{\sum_{j=1}^{K} \phi_j(z)}$$

➡ Matrix Form:

$$f_{\boldsymbol{w}}(z) = \boldsymbol{\psi}^T(z)\boldsymbol{w}, \text{ with } \psi_i(z) = \frac{\phi_i(z)z}{\sum_{j=1}^K \phi_j(z)}$$

$$\Rightarrow \text{ For } t \to \infty, f_{\boldsymbol{w}}(z) \to 0 \text{ as } z \to 0$$

A DMP is stable per construction as the forcing function vanishes it is just a standard PD for $t \to \infty$



Representation of the forcing function



Integrating the dynamical system leads to the trajectory



For multi-DoF robots, we use an individual DMP per DoF

Phase variable z is shared

53

Coupling between joints due to the shared phase

For periodic movements, we can use periodic phase variables





Adapting the meta-parameters...





Given:

- A desired trajectory and its derivatives $m{q}_{1:T}, \dot{m{q}}_{1:T}, \ddot{m{q}}_{1:T}$
- A goal attractor g (e.g. final position of trajectory)
- Parameters: $lpha, eta, lpha_z$ (typically fixed)
- Temporal Scaling au: Adjusted to movement duration

Algorithm:

Compute target values for each time step

$$f_t = \ddot{q}_t / \tau^2 - \alpha (\beta (g - q_t) - \dot{q} / \tau)$$

• Compute shape parameters $\,w$ by linear (ridge) regression

$$\boldsymbol{w} = (\boldsymbol{\Psi}^T \boldsymbol{\Psi} + \sigma^2 \boldsymbol{I})^{-1} \boldsymbol{\Psi}^T \boldsymbol{f}$$

Example: A Tennis Backhand









Rhythmic Motor Primitives



(Kober & Peters, ICRA 2009)



Fast Coupling between System and Gait





Important properties of movement primitive representations

- Data-Driven: easily learnable from demonstrations
- Generalization: easily adaptable to a new situation
- Combination: Co-activate multiple primitives to solve a combination of tasks
- Temporal Scaling: Modulate the execution speed of the movement
- Coupling: Represent the coupling between a high number of joints
- Variability: Reproduce the stochasticity in the demonstrations
- Optimality: Can we represent optimal behavior?
- Can be applied for rhythmic and stroke-based movements

Movement Primitives



What we have so far...

- Data-Driven: Yes
- Generalization: Only adapt final positions
- Combination: No idea how...
- Temporal Scaling: Yes
- Coupling: Yes, but only the mean is coupled, no correlations
- Variability: No
- Optimality: Is following a single trajectory really optimal? No
- Can be applied for rhythmic and stroke-based movements: Yes

Probabilistic Movement Primitives



Stochastic representation of trajectories:

$$\boldsymbol{\tau} \sim p(\boldsymbol{\tau})$$

• Use w to represent a single trajectory

 $\boldsymbol{\tau} = f(\boldsymbol{w}) + \boldsymbol{\epsilon}$

- Learn a distribution $p(\boldsymbol{w})$ over the \boldsymbol{w} vector
- Integrate out \boldsymbol{w} to obtain $p(\boldsymbol{\tau})$

Why is this useful?

- We can also represent uncertainty
- Uncertainty gives us information on importance of time points
- We can apply probabilistic operations





Representation of a single trajectory

$$y_t = \boldsymbol{\psi}_t^T \boldsymbol{w} + \epsilon_y \qquad \epsilon_y \sim \mathcal{N}(0, \sigma^2)$$

Phase-dependent basis:

$$\boldsymbol{\psi}_t = \boldsymbol{\psi}(z_t)$$

For example, normalized Gaussian basis functions

Probabilistic model:

$$p(\boldsymbol{\tau}|\boldsymbol{w}) = \prod_{t} \mathcal{N}(y_t | \boldsymbol{\psi}_t^T \boldsymbol{w}, \sigma^2) = \mathcal{N}(\boldsymbol{\tau}|\boldsymbol{\Psi}\boldsymbol{w}, \sigma^2 \boldsymbol{I}),$$

with $\boldsymbol{\Psi} = [\boldsymbol{\psi}_1, \dots, \boldsymbol{\psi}_T]^T$

Trajectory distribution: distribution over the parameters $p(\boldsymbol{w}|\boldsymbol{\theta})$

$$p(\boldsymbol{\tau}|\boldsymbol{\theta}) = \int_{\boldsymbol{w}} p(\boldsymbol{\tau}|\boldsymbol{w}) p(\boldsymbol{w}|\boldsymbol{\theta}) d\boldsymbol{w}$$

67



You can always rely on old friends...

Lets use a Gaussian: $p(w|\theta) = \mathcal{N}(w|\mu_w, \Sigma_w)$

Computing the trajectory distribution is now easy

$$\begin{split} p(\boldsymbol{\tau}|\boldsymbol{\theta}) &= \int p(\boldsymbol{\tau}|\boldsymbol{w}) p(\boldsymbol{w}|\boldsymbol{\theta}) d\boldsymbol{w} \\ &= \int \mathcal{N}(\boldsymbol{\tau}|\boldsymbol{\Psi}\boldsymbol{w}, \sigma^2 \boldsymbol{I}) \mathcal{N}(\boldsymbol{w}|\boldsymbol{\mu}_{\boldsymbol{w}}, \boldsymbol{\Sigma}_{\boldsymbol{w}}) d\boldsymbol{w} \\ &= \mathcal{N}(\boldsymbol{\tau}|\boldsymbol{\Psi}\boldsymbol{\mu}_{\boldsymbol{w}}, \sigma^2 \boldsymbol{I} + \boldsymbol{\Psi}\boldsymbol{\Sigma}_{\boldsymbol{w}} \boldsymbol{\Psi}^T) \end{split}$$

Hence, we can easily evaluate mean and variance for any time point



Hence, we can easily evaluate mean and variance for any time point





How can we encode a distribution over multiple DoFs?

Use a concatenated weight and trajectory vector and block-diagonal basis matrix

$$oldsymbol{ au} = egin{bmatrix} oldsymbol{ au}_1, \ dots \ oldsymbol{ au}_D \end{bmatrix} \quad oldsymbol{w} = egin{bmatrix} oldsymbol{ au}_1, \ dots \ oldsymbol{ au}_2, \ dots \ oldsymbol{ au}_2, \ dots \ oldsymbol{ au}_2, \ dots \ oldsymbol{ au}_D \end{bmatrix} \quad oldsymbol{\Phi} = egin{bmatrix} oldsymbol{ au} & oldsymbol{ au} & oldsymbol{ au}_2, \ dots \ oldsymbol{ au}_2, \ dots \ d$$

The same linear relation holds: $au = \Phi w$

 \clubsuit We use a distribution $p({m w}|{m \mu}_{{m w}},{m \Sigma}_{{m w}})$ over the parameters of all DoFs

For a single time step: $p(\boldsymbol{y}|\boldsymbol{\theta}) = \mathcal{N}(\boldsymbol{y}_t|\boldsymbol{\phi}_t\boldsymbol{\mu}_{\boldsymbol{w}}, \sigma^2 \boldsymbol{I} + \boldsymbol{\phi}_t\boldsymbol{\Sigma}_{\boldsymbol{w}}\boldsymbol{\phi}_t^T)$

Covariance matrix encodes correlation between the joints



Trajectory distribution tracking

How do we use a trajectory distribution for robot control?

We can obtain a time-varying stochastic feedback-controller in closed form

$$oldsymbol{u}_t = oldsymbol{k}_t + oldsymbol{K}_t oldsymbol{y}_t + oldsymbol{\epsilon}_u, \quad oldsymbol{\epsilon}_u \sim \mathcal{N}(oldsymbol{0}, oldsymbol{\Sigma}_u)$$

from $\boldsymbol{\tau} \sim p(\boldsymbol{\tau}|\boldsymbol{\theta})$ that exactly reproduces the given trajectory distribution (mean and variances)

- Same structure as optimal controllers for linear(ized) systems
- But it needs an accurate model



Generalization via Conditioning



Generalization: Change intermediate or end-point of the movement

We can condition p(w) on reaching position q_t^* at time-step t

- ➡ New trajectory distribution $p(w|q_t = q_t^*)$ is obtained by Bayes theorem
- Closed-form solution for Gaussian trajectory distributions



Combination of Movement Primitives



Modularity: Combine primitives to solve a combination of tasks

Implemented as **product of distributions:**

- "Intersection" of trajectory distributions
- Area, in which all distributions have
 high probability

$$p_{\rm co}(\boldsymbol{q}_t) \propto \prod_{i=1}^N p_i(\boldsymbol{q}_t)^{\alpha_i(t)}$$

 $p_i(\boldsymbol{q}_t) \dots$ i-th movement primitive $\alpha_i(t) \dots$ activation factors

 Computed in closed-form for Gaussian distributions







7-link planar robot arm, controlled by inverse dynamics

- Trained 2 movements for reaching different via-points at different time steps
- Combination of the movements reaches all 2 via-points





7-link KUKA robot arm, playing maracas

- Record rhythmic movements to produce sounds
- Blend between different rhythmic movements



Case Study: Robot Hockey



7-link KUKA robot arm, playing hockey

Train 2 primitives with high variance in shooting angle or in distance



Demonstration 1

- Product of the primitives:
 - **Combination of both tasks**





Demonstration 2

 Conditioning to select the shooting angle





Movement Primitives

What we have so far...

- Data-Driven: Yes
- Generalization: Yes
- Combination: Yes
- Temporal Scaling: Yes
- Coupling: Yes
- Variability: Yes
- Optimality: Yes
- Can be applied for rhythmic and stroke-based movements: Yes

Libraries of Primitives







Imagine the following situation...



What about many primitives?





What you can do with it...






- What to Imitate? The data traces will contain outliers, redundant data, data that is irrelevant to the task. How can the system extract the relevant components? Imitate on which **level of abstraction**?
- How to Imitate? body of the teacher ≠ body of the student
 "Correspondence Problem".
- When to Imitate? Not all behavior in a data stream may be suited for imitation. Untackled questions
- Whom to Imitate? If a scene with several actors is observed, the correct one needs to be extracted. Untackled questions

(Nehaniv & Dautenhan, 2001)

Imitation Learning





Problems of Imitation Learning

- Imperfect demonstrations -> require reinforcement learning

Intent identification -> requires inverse reinforcement learning
 30



What you should know...

- State-space representations versus trajectory-based policy representations
- What is imitation learning and when does it fail?
- What are the main ideas of using movement primitives?
- Why do use dynamical systems? Advantages/Disadvantages?
- Why do use a probabilistic representation?