## Inverse Reinforcement Learning via Deep Gaussian Process

#### Presenter: Ming Jin UAI 2017

#### Ming Jin, Andreas Damianou, Pieter Abbeel, and Costas J. Spanos

#### Understanding people's preference is key to bringing technology to our daily life









### How AI can act like a social/ friendly/ intelligent/ real person?

- People's decision making often involves...
  - Long-term planning vs. short-term gain
  - Risk seeking vs. risk aversion
  - Individual preferences over outcomes





### How AI can act like a social/ friendly/intelligent/real person?

Three learning paradigms...



How to learn an agent's intention by observing a limited number of demonstrations?

### Outline

- Inverse reinforcement learning (IRL) problem formulation
- Gaussian process reward modeling
- Incorporating the representation learning
- Variational inference
- Experiments
- Conclusion



Key challenges:

- Providing a formal specification of the control task.
- Building a good dynamics model.
- Finding closed-loop controllers.

(Adapted from Abbeel's slides on "Inverse reinforcement learning")

## IRL in a nutshell: given demonstrations, infer reward

- Input:
  - Dynamics model  $P_{sa}(s_{t+1}|s_t, a_t)$
  - No reward function  $r^*(s)$
  - Demonstration  $\mathcal{M}$ :  $s_0, a_0, s_1, a_1, s_2, a_2, \ldots$

(= trace of the teacher's policy  $\pi^*$ )

- Desired output:
  - Reward function r(s)
  - Policy  $\hat{\pi}: S \to A$ , which (ideally) has performance guarantees, e.g., expected reward difference (EVD)

$$\mathbb{E}\left[\sum_{t=0}^{\infty} \gamma^{t} r^{*}(x_{t}) | \pi^{*}\right] - \mathbb{E}\left[\sum_{t=0}^{\infty} \gamma^{t} r^{*}(x_{t}) | \hat{\pi}\right]$$

## Parametric vs. nonparametric representation of reward function

- Linear representation:  $r(s) = \mathbf{w}^T \boldsymbol{\phi}(s)^{\text{Limited}}$ 
  - Maximum margin planning (MMP): (Ratliff et al, 2006)
  - Maximum Entropy IRL (MaxEnt): (Ziebart et al., 2008)
- Nonlinear representation:
   Deep NN: (Wulfmeier et al., 2015)
- Needs significant demonstrations to avoid overfitting
- Nonparametric representation:
  - Gaussian Process IRL: (Levine et al., 2011)

Deep Gaussian Process IRL: (Jin et al., 2017)

### Two slides on Gaussian Process

- A Gaussian distribution depends on a mean and a covariance matrix.
- A Gaussian process depends on a mean and a covariance function.
- Let's start with a multivariate Gaussian Distr.:

$$p(f_1, f_2, \dots, f_s, f_{s+1}, f_{s+2}, \dots, f_N) \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{K})$$

$$f_A \qquad f_B$$

$$\boldsymbol{\mu} = \begin{bmatrix} \boldsymbol{\mu}_A \\ \boldsymbol{\mu}_B \end{bmatrix} \text{ and } \boldsymbol{K} = \begin{bmatrix} \boldsymbol{K}_{AA} & \boldsymbol{K}_{AB} \\ \boldsymbol{K}_{BA} & \boldsymbol{K}_{BB} \end{bmatrix}$$

(Adapted from Damianou's slides on "System identification and control with (deep) Gaussian process")

### Two slides on Gaussian Process

• In the GP context, we deal with an infinite dimensional Gaussian distribution:

$$\boldsymbol{\mu} = \begin{bmatrix} \boldsymbol{\mu}_X \\ \dots \end{bmatrix} \text{ and } \boldsymbol{K} = \begin{bmatrix} \boldsymbol{K}_{XX} & \dots \\ \dots & \dots \end{bmatrix}$$

#### Conditional distribution:

**Gaussian distribution** 

#### Gaussian process

Given  $p(f_A, f_B) \sim \mathcal{N}(\mu, K)$ Then the posterior distribution:  $p(f_A | f_B) \sim \mathcal{N}(\mu_{A|B}, K_{A|B})$ with  $\mu_{A|B} = \mu_A + K_{AB} K_{BB}^{-1}(f_B - \mu_B)$  $K_{A|B} = K_{AA} - K_{AB} K_{BB}^{-1} K_{BA}$  Given  $p(f_X, f_*) \sim \mathcal{N}(\mu, K)$ Then the posterior distribution:  $p(f_*|f_X) \sim \mathcal{N}(\mu_{*|X}, K_{*|X})$ with  $\mu_{A|B} = \mu_A + K_{AB}K_{BB}^{-1}(f_B - \mu_B)$  $K_{A|B} = K_{AA} - K_{AB}K_{BB}^{-1}K_{BA}$ 

(Adapted from Damianou's slides on "System identification and control with (deep) Gaussian process")

## Reward function modeled by a Gaussian process

- Reward is a function of states:  $r(x) \in \mathcal{R}$
- We discretize the world into *n* states, each described by a feature vector  $x_i \in \mathbb{R}^m$ :

$$\boldsymbol{X} = \begin{bmatrix} \boldsymbol{x}_1^T \\ \vdots \\ \boldsymbol{x}_n^T \end{bmatrix} \in \mathcal{R}^{n \times m}, \, \boldsymbol{r} = \begin{bmatrix} r(\boldsymbol{x}_1) \\ \vdots \\ r(\boldsymbol{x}_n) \end{bmatrix} \in \mathcal{R}^n$$

• It can be modeled with a zero-mean GP prior:  $r|X, \theta \sim \mathcal{N}(0, K_{XX})$ 

Parameter of the covarince function

## How to find out the reward given demonstrations?

• GPIRL trains the parameters by maximum likelihood estimation (Levine et al., 2011):

$$p(\mathcal{M}|\mathbf{X}, \boldsymbol{\theta}) = \int p(\mathcal{M}|\mathbf{r})p(\mathbf{r}|\mathbf{X}, \boldsymbol{\theta})d\mathbf{r}$$
  
RL term (Ziebart et al., 2008) GP prior

• Prediction of the reward at any test input is found through the conditional:

$$r^* | r, X, x^* \sim \mathcal{N}(K_{x^*} K_{XX}^{-1} r, k_{x^* x^*} - K_{x^* X} K_{XX}^{-1} K_{Xx^*})$$

## Can we improve the *complexity* of reward function w/o *overfitting*?

• Step function example:

 $\overbrace{f^{(1)}(x_1)}^{X_1}$ 

Predictive draws of step function



Overfitting does not appear to be a problem for a deeper architecture...  $f^{(2)}(f^{(1)}(x_1))$ 

...the main challenge is to train such a system!





Figures adapted from http://keyonvafa.com/deep-gaussian-processes/

## Can we improve the *complexity* of reward function w/o *overfitting*?

• IRL with reward modeled as GP: (Levine et al., 2011)

$$(\mathbf{X} \xrightarrow{\mathsf{GP}} \mathbf{r} \xrightarrow{\mathsf{RL} \text{ policy}}_{\text{learning}} \xrightarrow{\mathsf{M}} \mathbf{M}$$
$$p(\mathcal{M} | \mathbf{X}, \boldsymbol{\theta}) = \int p(\mathcal{M} | \mathbf{r}) p(\mathbf{r} | \mathbf{X}, \boldsymbol{\theta}) d\mathbf{r}$$

• IRL with reward modeled as deep GP:

$$\begin{array}{c} & \overbrace{X} \xrightarrow{\mathsf{GP}} \overbrace{X_2} \xrightarrow{\mathsf{GP}} \cdots \xrightarrow{\mathsf{K}_L} \xrightarrow{\mathsf{GP}} \overbrace{r} \xrightarrow{\mathsf{RL policy}} \underset{\text{learning}}{\overset{tarring}}{\overset{$$

## We need a tractable form of the likelihood function for training

• Let's illustrate with a 2-stack GP:

The integral is still intractable...  $p(\mathcal{M}|\mathbf{X})$  . Fil policy  $= \int p(\mathcal{M}, \boldsymbol{r}, \boldsymbol{D}, \boldsymbol{B} | \boldsymbol{X}) d(\boldsymbol{D}, \boldsymbol{B}, \boldsymbol{r})$ GP prior:  $\mathcal{N}(0, \mathbf{K}_{\mathbf{X}\mathbf{X}})$  $p(\mathcal{M}|\mathbf{r}) * p(\tilde{\mathbf{r}}|\mathbf{D}) * p(\mathbf{D}|\mathbf{B}) * p(\mathbf{B}|\mathbf{X})$ RL probability given by Gaussian noise: MaxEnt (Ziebart et al., 2008)  $D_{ij} \sim \mathcal{N}(B_{ij}, \lambda^{-1})$ 

### Introduce inducing points

• Add to each latent layer..



# Use *variational lower bound* to make training tractable

• Introduce variational distribution:

Model distributions  $\mathcal{P}$ 

: cond. Gaussian  $p(\boldsymbol{D}|\boldsymbol{B})$ : Gaussian noise  $p(\boldsymbol{V}|\boldsymbol{W})$ :  $GP(\boldsymbol{0}, \boldsymbol{K}_{\boldsymbol{W}\boldsymbol{W}})$  Variational distributions Q

 $q(\boldsymbol{B}|\boldsymbol{V},\boldsymbol{X},\boldsymbol{W}) = p(\boldsymbol{B}|\boldsymbol{V},\boldsymbol{X},\boldsymbol{W})q(\boldsymbol{f})$ 

$$q(\boldsymbol{D}): \prod \delta(\boldsymbol{d}^m - \boldsymbol{K}_{\boldsymbol{X}\boldsymbol{W}}\boldsymbol{K}_{\boldsymbol{W}\boldsymbol{W}}^{-1}\widetilde{\boldsymbol{v}}^m) q(\boldsymbol{V}): \prod \mathcal{N}(\widetilde{\boldsymbol{v}}^m, \boldsymbol{G}^m)$$

• Jensen's inequality to derive lower bound:

$$\log \int \mathcal{P} \geq \int \mathcal{Q} \, \log \frac{\mathcal{P}}{\mathcal{Q}}$$

### Use variational lower bound to make training tractable

Derive a tractable lower bound:  $\log p(\mathcal{M}|\mathbf{X})$  $= \log \int p(\mathcal{M}|\mathbf{r})p(\mathbf{r}|\mathbf{f},\mathbf{D})p(\mathbf{f})p(\mathbf{D}|\mathbf{B})p(\mathbf{B}|\mathbf{V},\mathbf{X})p(\mathbf{V})$ 

•

 $\geq \int q(\boldsymbol{f})q(\boldsymbol{D})p(\boldsymbol{B}|\boldsymbol{V},\boldsymbol{X})q(\boldsymbol{V})\log\frac{p(\mathcal{M}|\boldsymbol{r})p(\boldsymbol{r}|\boldsymbol{f},\boldsymbol{D})p(\boldsymbol{f})p(\boldsymbol{D}|\boldsymbol{B})p(\boldsymbol{V})}{q(\boldsymbol{f})q(\boldsymbol{D})q(\boldsymbol{V})}$  $=\mathcal{L}_{\mathcal{M}}+\mathcal{L}_{\mathcal{G}}-\mathcal{L}_{KL}+\mathcal{L}_{\mathcal{B}}-\frac{nm}{2}\log(2\pi\lambda^{-1})$ The value can be computed Optimizing the objective We can also **compute the derivative** turns the variational with respect to the parameters: distribution Q into the true Kernel function parameters model posterior Inducing points

## Inducing points provide a succinct summary of data

• Given a new state *x*<sup>\*</sup>, the predicted reward is a function of the *latent representation*:



This can be used for knowledge transfer in a new situation

### The original likelihood for the deep GP IRL is intractable..

Our method introduces inducing points combined with variational inequality to derive a lower bound for inference.

### Outline

- Inverse reinforcement learning (IRL) problem formulation
- Gaussian process reward modeling
- Incorporating the representation learning
- Variational inference

#### • Experiments

Conclusion

## We use the *expected value difference (EVD)* as the metric

• EVD measures the expected reward earned under the optimal policy  $\pi^*$  w/ true reward and the policy derived by IRL  $\hat{\pi}$ 

$$\mathbb{E}\left[\sum_{t=0}^{\infty} \gamma^{t} r(x_{t}) | \pi^{*}\right] - \mathbb{E}\left[\sum_{t=0}^{\infty} \gamma^{t} r(x_{t}) | \hat{\pi}\right]$$

 We also visually compare the true reward with that learned by IRL

## ObjectWorld has a nonlinear reward structure

- Reward: +1: 1step within 
   A 3 step within 
   1 if only 3 steps withi 
   0 otherwise
- · Features: min. dist. to an object of each type
- Nonlinear, but still preserves local distance



### Both DGP-IRL and GPIRL captures the correct reward

Plot of EVD for the world where demos are given



Plot of EVD for a new world where demos are not available



### Both DGP-IRL and GPIRL captures the correct reward

#### **Ground truth**



#### FIRL

$\triangleright$	►	►	►	$\triangleright$	$\bigtriangledown$	►	▲	۲	$\bigtriangledown$	٩	4
$\triangleright$	►	►	►	$\bigtriangledown$	$\bigtriangledown$	►	►	$\triangleright$	۲	٩	∢
-	۲	►	►	$\triangleright$	$\bigtriangledown$	►	►	D			∢
$\bigtriangledown$	$\overline{}$	►	►	$\triangleright$	$\triangleright$	$\triangleright$	$\triangleright$				
-	-	-	►	$\triangleright$	$\triangleright$	$\triangleright$	$\triangleright$				∢
$\bigtriangledown$	$\bigtriangledown$	$\bigtriangledown$	⊲			<b>^</b>	<b>^</b>	<b>^</b>			
$\diamond$	$\bigtriangledown$	٩	⊲	•				►			
$\bigtriangledown$	$\bigtriangledown$	٩	⊲	∢	•	Ŷ	►	▶			
$\bigtriangledown$	$\bigtriangledown$	٩	⊲	∢	•	$\overline{}$	-	•	$\triangleright$	0	∢
٠		٩	⊲	∢	-	$\overline{}$	►	⊳	٩		
$\triangleright$			٩	٩	٩	$\triangleright$	$\triangleright$		¢		
$\triangleright$	∢	∢		⊲	⊲	$\triangleright$	$\triangleright$				

#### **DGP-IRL**



**MaxEnt** 

1

1

 $\bigtriangleup$ 

D

 $\bigtriangleup$ 

#### **GPIRL**

	►	►	$\triangleright$	$\triangleright$	$\bigtriangledown$		$\triangleright$	٠	0	٩	٩
►	•	►	►	⊳	$\bigtriangledown$		$\triangleright$	D	۲	⊲	∢
-	¢	►	►	$\triangleright$	$\bigtriangledown$	$\triangleright$	►	$\triangleright$	$\bigtriangleup$		∢
$\overline{}$	►	►	►	$\triangleright$	$\triangleright$	$\triangleright$	$\triangleright$				
-	►	-		$\triangleright$	$\triangleright$	$\triangleright$	$\bigtriangleup$				
$\bigtriangledown$	$\triangleright$	$\bigtriangledown$	$\bigtriangledown$								
$\diamond$	$\bigtriangledown$	$\bigtriangledown$	٩	∢				►	$\bigtriangleup$		
$\bigtriangledown$	$\bigtriangledown$	٩	٩	∢	∢	¢	٨	►			
$\bigtriangledown$	⊲	٩	٩	∢	$\overline{}$	▼		-		Ô	
٠	$\triangleleft$	⊲	⊲	٩	-	$\overline{}$	٨	0	∢		
0	∢	٩	٩	$\bigtriangledown$	$\bigtriangledown$	$\triangleright$	$\triangleright$		Ô		
	⊲	∢	٩	∢	$\triangleright$	D	$\triangleright$				

#### MMP



BinaryWorld has even more nonlinear reward structure

- Reward: neighborhood
   with 4 blues (+1) and 5 blues (-1)
  - Nonlinear, combinatorial
- Features: ordered list of colors (top left to bottom right)



# DGP-IRL outperforms GPIRL in this more complex case

Plot of EVD for the world where demos are given

Plot of EVD for a new world where demos are not available



## DGP-IRL outperforms GPIRL in this more complex case

#### **Ground truth**



FIRL

<b>&gt;</b>	٠	۲	٠	<b>&gt;</b>	٠	۲	٠	۰	۰	٠	<
	٠	•	<b></b>		٠	٢	٠	٠	<b></b>	٠	4
٠		٠	$\diamond$		٥		٠	۲	٠		<b></b>
$\diamond$	<b>&gt;</b>	<u>ا</u>	<b></b>	٠	$\diamond$	٠	٠	<b>♦</b>	•	۰	•
>	٠	<b>♦</b>	٠	<b>♦</b>	<b>♦</b>	٢	٠	<b>♦</b>	٠	•	<b></b>
٠	۲	0	•	<b></b>	۰	۲			•	٠	1
$\diamond$		٠	٢		<b>&gt;</b>	<b>&gt;</b>	<b>&gt;</b>	٠	٠		<b></b>
	٠	۰	<b></b>	٠	<b></b>	٠		٠	•	•	
٠		٠	٠	٠	<b>&gt;</b>	<b>&gt;</b>	<u>ب</u>	٠	•	•	<b></b>
	٠	٢	<b></b>	٠	<b>&gt;</b>	٠	<b>♦</b>	<b></b>	•		4
<b>•</b>	۲	٠	٥	>	٠	۰	۲	$\diamond$		<b></b>	\$
•					٠		<b>♦</b>			\$	<b></b>

#### DGP-IRL



MaxEnt



#### **GPIRL**



MMP



## DGP-IRL learns the most succinct feature for the reward

The reward is not separable in the original feature space...



...but separable in the latent space



# DGP-IRL also outperforms in learning the driving behavior

The goal is to navigate the robot car as fast as possible, but avoid speeding when the police car is nearby.



### Outline

- Inverse reinforcement learning (IRL) problem formulation
- Gaussian process reward modeling
- Incorporating the representation learning
- Variational inference
- Experiments
- Conclusion

### How to learn an agent's intention by observing *a limited number* of demonstrations?

- Model the reward function as a deep GP to handle complex characteristics.
- This enables simultaneous feature state representation learning and IRL reward estimation.
- Train the model through variational inference to guard against overfitting, thus work efficiently with limited demonstrations.

### Future works

- DGP-IRL enable easy incorporation of side knowledge (through priors on the latent space) to IRL.
- Combine deep GPs with other complicated inference engines, e.g., selective attention models (Gregor et al., 2015).
- Application domains: mobile sensing for health, building and grid controls, multiobjective control, and human-in-the-loop gamification.

#### **Incomplete List of Collaborators**

#### Amazon.com, UK

Andreas Damianou

#### **UC Berkeley**

- Pieter Abbeel
- Costas Spanos











## Additional slides

# DGP-IRL also outperforms in learning the driving behavior

- Three-lane highway, vehicles of specific class (civilian or police) and category (car or motorcycle) are positioned at random, driving at the same constant speed.
- The robot car can switch lanes and navigate at up to three times the traffic speed.
- The state is described by a continuous feature which consists of the closest distances to vehicles of each class and category in the same lane, together with the left, right, and any lane, both in the front and back of the robot car, in addition to the current speed and position.