

# Maximum Margin Planning

Nathan Ratliff, Drew Bagnell and Martin Zinkevich

Presenters:

Ashesh Jain, Michael Hu

CS6784 Class Presentation

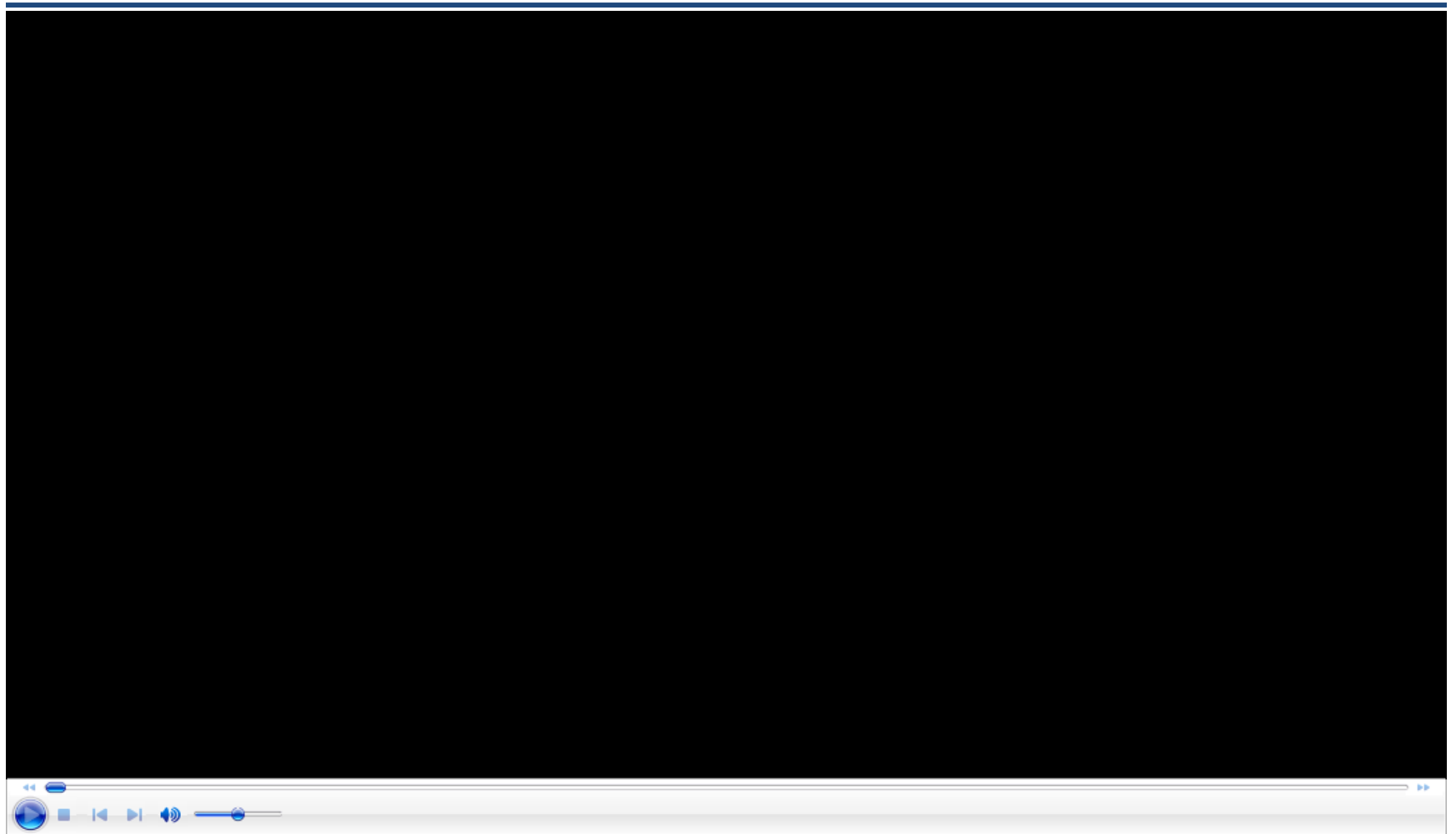
# Theme

---

1. Supervised learning
2. Unsupervised learning
3. Reinforcement learning
  - Reward based learning

## Inverse Reinforcement Learning

# Motivating Example (Abbeel et. al)



# Outline

---

## Basics of Reinforcement Learning

- State and action space
- Activity

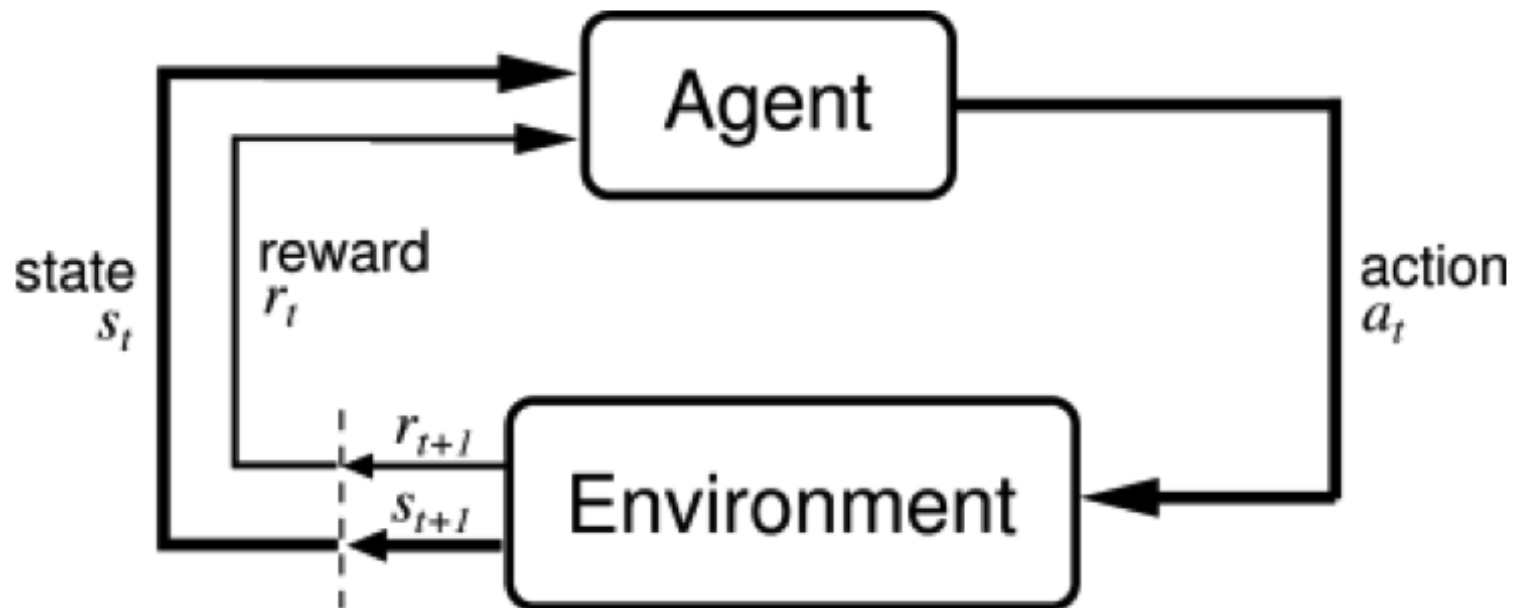
## Introducing Inverse RL

- Max-margin formulation
- Optimization algorithms
- Activity

## Conclusion

# Reward-based Learning

---



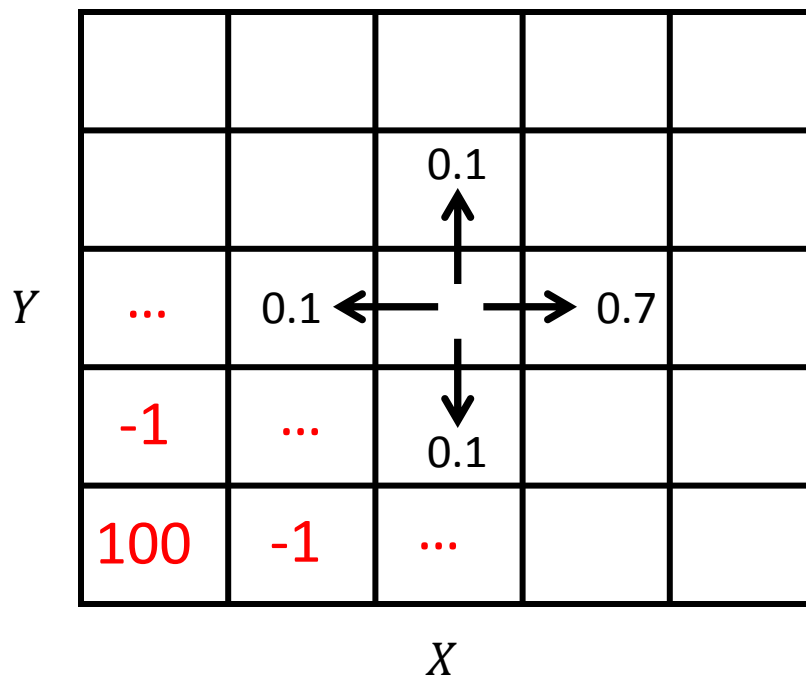
Credit: Sutton and Barto, Reinforcement Learning: An Introduction, 1998

# Markov Decision Process (MDP)

---

MDP is defined by  $(\mathcal{X}, \mathcal{A}, p, R)$

- $\mathcal{X}$  State space
- $\mathcal{A}$  Action space
- $p$  Dynamics model  
 $p(x'|x, a)$
- $R$  Reward function  
 $R(x)$  or  $R(x, a)$



# Markov Decision Process (MDP)

MDP is defined by  $(\mathcal{X}, \mathcal{A}, p, R)$

$$\pi: \mathcal{X} \rightarrow \mathcal{A}$$

$\mathcal{A}$  Action space

$p$  Dynamics model

$$p(x'|x, a)$$

$R$  Reward function

$$R(x) \text{ or } R(x, a)$$

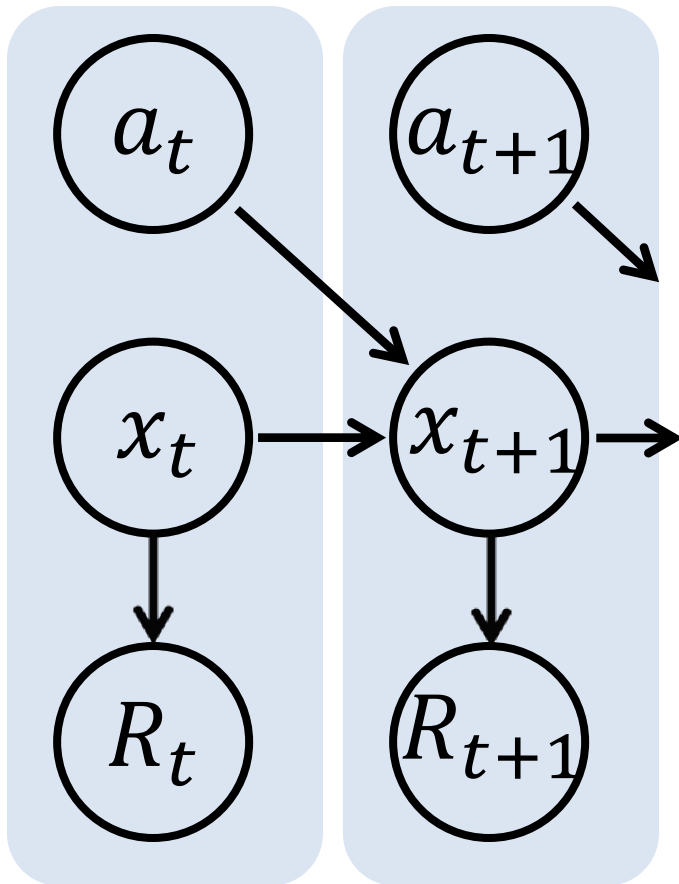
Policy

→	→	→	→	→
↑	↑	↓	→	↑
↑	→	↓	↑	↑
↑	↑	↓	→	↑
↑	↑	→	↑	↑

Example policy

# Graphical representation of MDP

---



Goal: Learn an optimal policy

$$\pi: \mathcal{X} \rightarrow \mathcal{A}$$

$$\pi^* = \arg \max_{\pi} E \left[ \sum_{t=0}^{\infty} \gamma^t R_t(x_t, \pi(x_t)) \right]$$

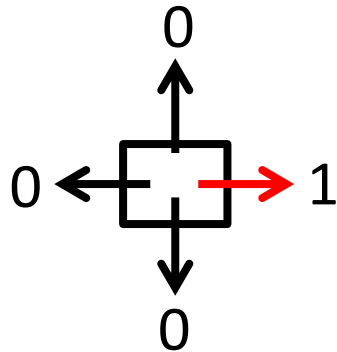
$$(x_t, \pi(x_t)) \rightarrow (x_{t+1}, \pi(x_{t+1})) \rightarrow \dots$$



# Activity – Very Easy!!!

---

Deterministic Dynamics



	-1	-1	-1	-1	 100
	-1	-1	-1	-1	-1
	-1	-1	 -50	-1	-1
	-1	-1	-1	-1	 -80
	-1	-1	-1	-1	-1




## Draw the optimal Policy!

# Activity – Very Easy!!!

---

-1	-1	-1	-1	 100
-1	-1	-1	-1	-1
-1	-1	 -50	-1	-1
-1	-1	-1	-1	 -80
-1	-1	-1	-1	-1

# Activity – Solution

→ -1	→ -1	→ -1	→ -1	 100
↑ -1	↑ -1	↑ -1	↑ -1	↑ -1
↑ -1	↑ -1	 -50	↑ -1	↑ -1
↑ -1	→ -1	→ -1	↑ -1	 -80
↑ -1	→ -1	→ -1	↑ -1	← -1

# Outline

---

## Basics of Reinforcement Learning

- State and action space
- Activity

## Introducing Inverse RL

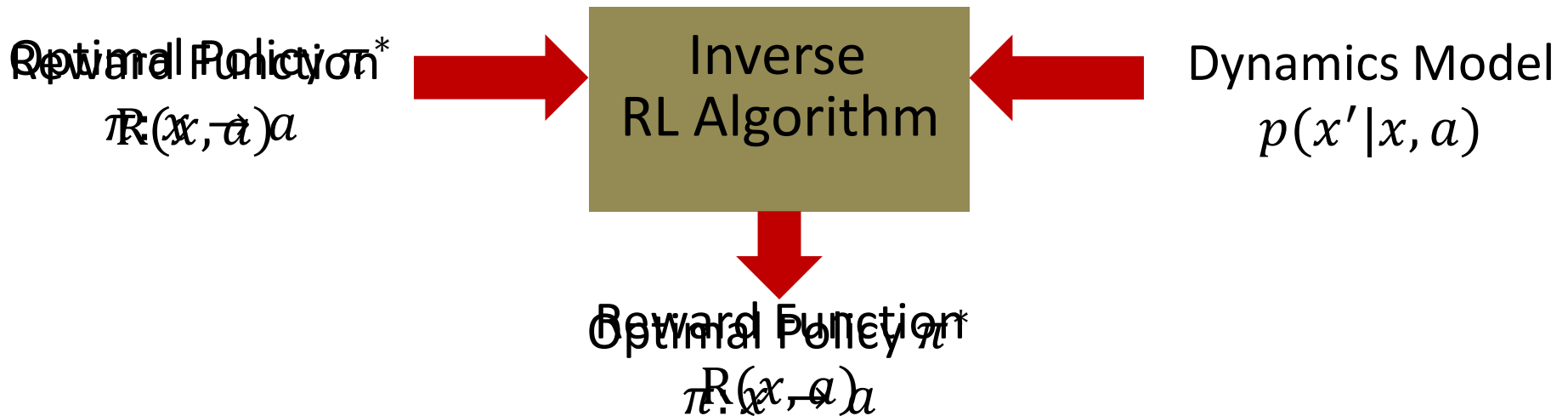
- Max-margin formulation
- Optimization algorithms
- Activities

## Conclusion

# RL to Inverse-RL (IRL)

---

- State space:  $\mathcal{X}$
- Action space:  $\mathcal{A}$



**Inverse problem:** Given optimal policy  $\pi^*$  or samples drawn from it, recover the reward  $R(x, a)$

# Why Inverse-RL?

---

- State space:  $\mathcal{X}$
- Action space:  $\mathcal{A}$

Reward Function  
 $R(x, a)$

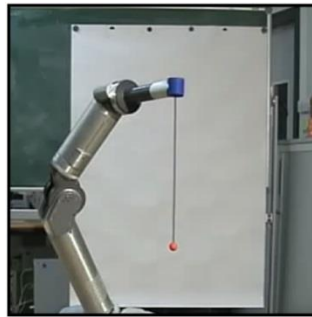


RL Algorithm

Specifying  $R(x, a)$  is hard, but **samples** from  $\pi^*$  are easily available



Ratliff et. al.



Kober et. al.



Abbeel et. al.

# IRL framework

$\pi^*: x \rightarrow a$

Expert

Interacts



Demonstration

$$y^* = (x_1, a_1) \rightarrow (x_2, a_2) \rightarrow (x_3, a_3) \rightarrow \dots \rightarrow (x_n, a_n)$$

$$w^T f(y^*) = w^T \begin{bmatrix} \text{blue} \\ \text{red} \\ \text{dark blue} \\ \text{purple} \\ \text{brown} \end{bmatrix} + w^T \underbrace{\begin{bmatrix} \text{dark blue} \\ \text{light purple} \\ \text{dark blue} \\ \text{dark blue} \\ \text{olive} \end{bmatrix}}_{\text{Reward}} + w^T \begin{bmatrix} \text{blue} \\ \text{tan} \\ \text{dark blue} \\ \text{red} \\ \text{grey} \end{bmatrix} + \dots + w^T \begin{bmatrix} \text{black} \\ \text{light grey} \\ \text{dark blue} \\ \text{grey} \\ \text{black} \end{bmatrix}$$

# Paper contributions

---

1. Formulated IRL as structural prediction
  - Maximum margin approach
2. Two optimization algorithms
  - Problem specific solver
  - Sub-gradient method
3. Robotic application: 2D navigation etc.



# Outline

---

## Basics of Reinforcement Learning

- State and action space
- Activity

## Introducing Inverse RL

- Max-margin formulation
- Optimization algorithms
- Activities

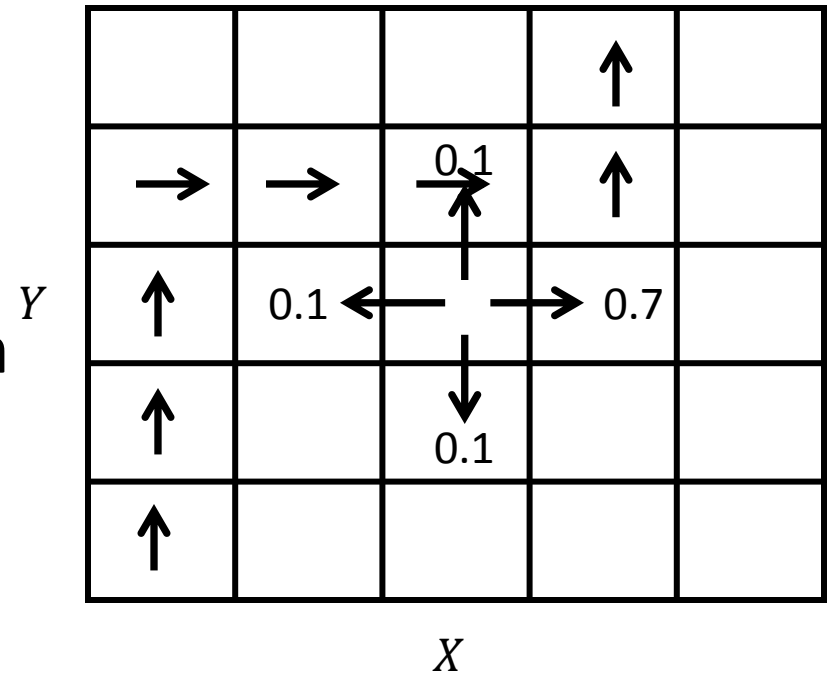
## Conclusion

# Formulation

---

Input  $\{(\mathcal{X}_i, \mathcal{A}_i, p_i, f_i, y_i, \mathcal{L}_i)\}_{i=1}^n$

- $\mathcal{X}_i$  State space
- $\mathcal{A}_i$  Action space
- $p_i$  Dynamics model
- $y_i$  Expert demonstration
- $f_i(\cdot)$  Feature function
- $\mathcal{L}_i(\cdot, y_i)$  Loss function

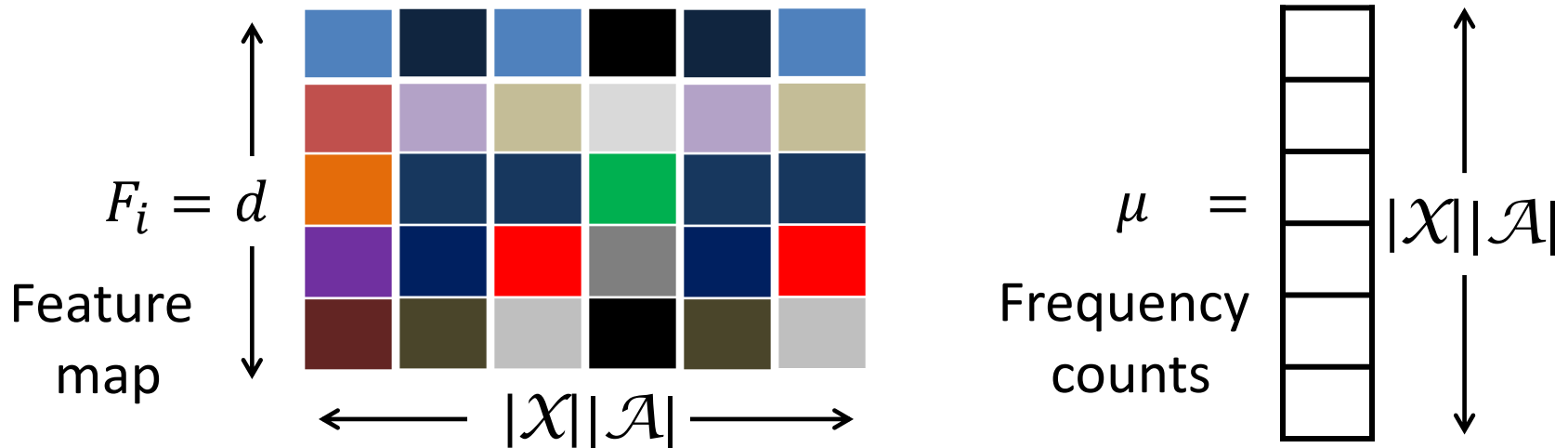


# Max-margin formulation

$$\min_{w, \xi_i} \frac{\lambda}{2} \|w\|^2 + \frac{1}{n} \sum_i \beta_i \xi_i$$

$$s. t. \forall i \quad w^T f_i(y_i) \geq \max_{y \in \mathcal{Y}_i} w^T f_i(y) + \mathcal{L}(y, y_i) - \xi_i$$

Features linearly decompose over path  $f_i(y) = F_i \mu$



# Max-margin formulation

---

$$\min_{w, \xi_i} \frac{\lambda}{2} \|w\|^2 + \frac{1}{n} \sum_i \beta_i \xi_i$$

$$s. t. \quad \forall i \quad w^T F_i \mu_i \geq \max_{\mu \in \mathcal{G}_i} w^T F_i \mu + l_i^T \mu - \xi_i$$

$\mu$  satisfies Bellman-flow constraints

$$\forall x' \quad \underbrace{\sum_{x,a} \mu^{x,a} p_i(x'|x,a)}_{\text{Inflow}} + s_i^{x'} = \underbrace{\sum_a \mu^{x',a}}_{\text{Outflow}}$$

# Outline

---

## Basics of Reinforcement Learning

- State and action space
- Activity

## Introducing Inverse RL

- Max-margin formulation
- Optimization algorithms
  - Problem specific QP-formulation
  - Sub-gradient method
- Activities

## Conclusion

# Problem specific QP-formulation

$$\min_{w, \xi_i} \frac{\lambda}{2} \|w\|^2 + \frac{1}{n} \sum_i \beta_i \xi_i$$

$$\text{s.t. } \forall i \quad w^T F_i \mu_i \geq \max_{\mu \in \mathcal{G}_i} w^T F_i \mu + l_i^T \mu - \xi_i$$

$$\forall i, x, a \quad \underbrace{\sum_{x, a} v_i^{x, a} p_i(x, a)}_{\text{Inflow}} \geq p_i(x, a) \Rightarrow \underbrace{\sum_{x', a} p_i(x', a) v_i^{x'}}_{\text{Outflow}}$$

Can be optimized using off-the-shelf QP solvers

# Outline

---

## Basics of Reinforcement Learning

- State and action space
- Activity

## Introducing Inverse RL

- Max-margin formulation
- Optimization algorithms
  - Problem specific QP-formulation
  - Sub-gradient method
- Activities

## Conclusion

# Optimizing via Sub-gradient

---

$$\min_{w, \xi_i} \frac{\lambda}{2} \|w\|^2 + \frac{1}{n} \sum_i \beta_i \xi_i$$

$$s. t. \quad \forall i \quad w^T F_i \mu_i \geq \max_{\mu \in \mathcal{G}_i} w^T F_i \mu + l_i^T \mu - \xi_i$$

Re-writing  
constraint

$$\xi_i \geq \max_{\mu \in \mathcal{G}_i} (w^T F_i \mu + l_i^T \mu) - w^T F_i \mu_i \geq 0$$

Its tight at  
optimality

$$\xi_i = \max_{\mu \in \mathcal{G}_i} (w^T F_i \mu + l_i^T \mu) - w^T F_i \mu_i$$

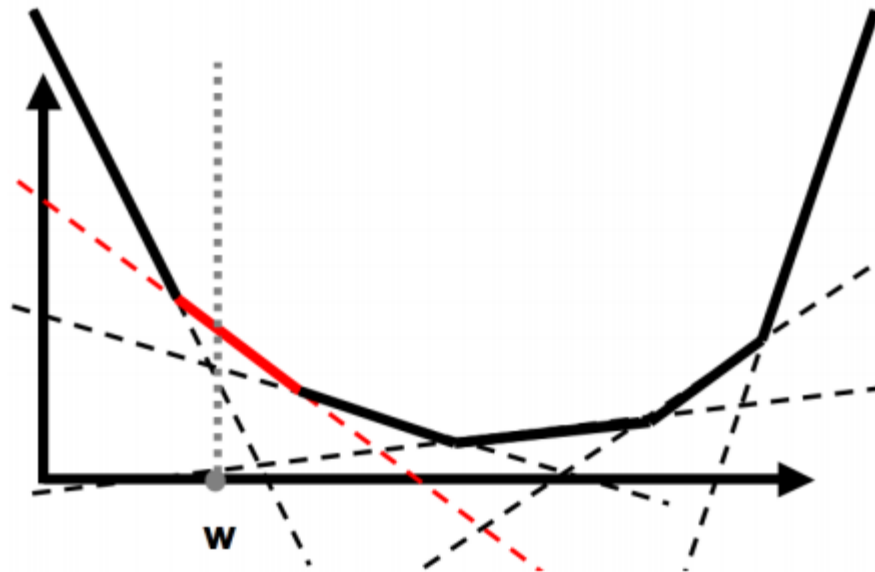
$\mu$  satisfies Bellman-flow constraints



# Optimizing via Sub-gradient

---

$$c(w) = \min_{w, \xi_i} \frac{\lambda}{2} \|w\|^2 + \frac{1}{n} \sum_i \beta_i \left[ \max_{\mu \in \mathcal{G}_i} (w^T F_i \mu + l_i^T \mu) - w^T F_i \mu_i \right]$$



Ratliff, PhD Thesis

# Weight update via Sub-gradient

---

Standard gradient descent update

$$w_{t+1} \leftarrow w_t - \alpha \underbrace{g_{w_t}}_{\nabla_w c(w)}$$

$$\min_{w, \xi_i} \frac{\lambda}{2} \|w\|^2 + \frac{1}{n} \sum_i \beta_i \left[ \max_{\mu \in \mathcal{G}_i} (w^T F_i \mu + l_i^T \mu) - w^T F_i \mu_i \right]$$

$$\mu^* = \operatorname{argmax}_{\mu \in \mathcal{G}_i} \underbrace{(w_t^T F_i + l_i^T)}_{\text{Loss-augmented reward}} \mu$$

$$g_{w_t} = \lambda w_t + \frac{1}{n} \sum_i \beta_i F_i (\mu^* - \mu_i)$$

# Convergence summary

---

$$c(w) = \frac{\lambda}{2} \|w\|^2 + \frac{1}{n} \sum_{i=1}^n r_i(w)$$

$$w \leftarrow w - \alpha g_w$$

How to choose  $\alpha$ ?

If  $\forall w \ \|g_w\| \leq G$  then a constant step size  $0 < \alpha \leq \frac{1}{\lambda}$  gives linear convergence to a region around the minimum

$$\|w_{t+1} - w^*\|^2 \leq \kappa^{t+1} \|w_0 - w^*\|^2 + \frac{\alpha G^2}{\lambda}$$

Optimization error

$$0 < \kappa < 1$$

# Outline

---

## Basics of Reinforcement Learning

- State and action space
- Activity

## Introducing Inverse RL

- Max-margin formulation
- Optimization algorithms
- Activity

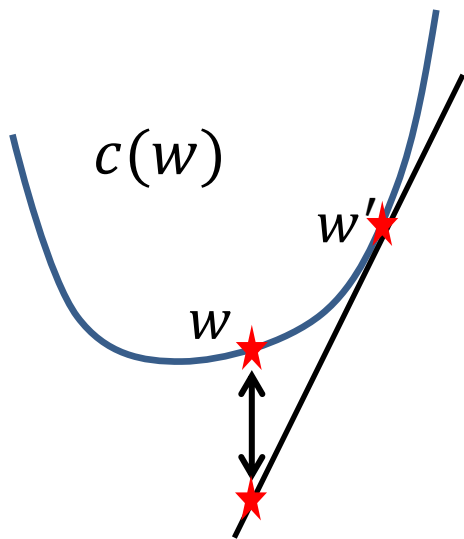
## Conclusion

# Convergence proof

We know  $\forall w \quad \|g_w\| \leq G$

$$\|w_{t+1} - w^*\|^2 = \|w_t - \alpha g_t - w^*\|^2$$

$$= \|w_t - w^*\|^2 + \alpha^2 \|g_t\|^2 - \boxed{\text{?}}$$



$$c(w) \geq c(w') + g_{w'}^T (w - w') + \frac{\lambda}{2} \|w - w'\|^2$$

Substitute  $w = w^*$  and  $w' = w_t$

$$\frac{\lambda}{2} \|w_t - w^*\|^2 \leq g_t^T (w_t - w^*)$$

Use  $c(w^*) \leq c(w_t)$

$$\leq (1 - \alpha\lambda) \|w_t - w^*\|^2 + \alpha^2 G^2$$

# Outline

---

## Basics of Reinforcement Learning

- State and action space
- Activity

## Introducing Inverse RL

- Max-margin formulation
- Optimization algorithms
- Activity

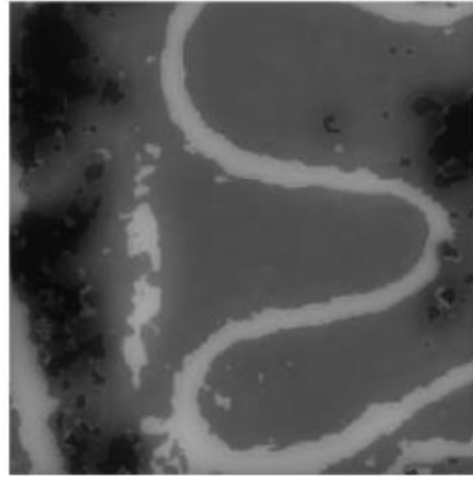
## Conclusion

# 2D path planning

---



Expert's  
Demonstration

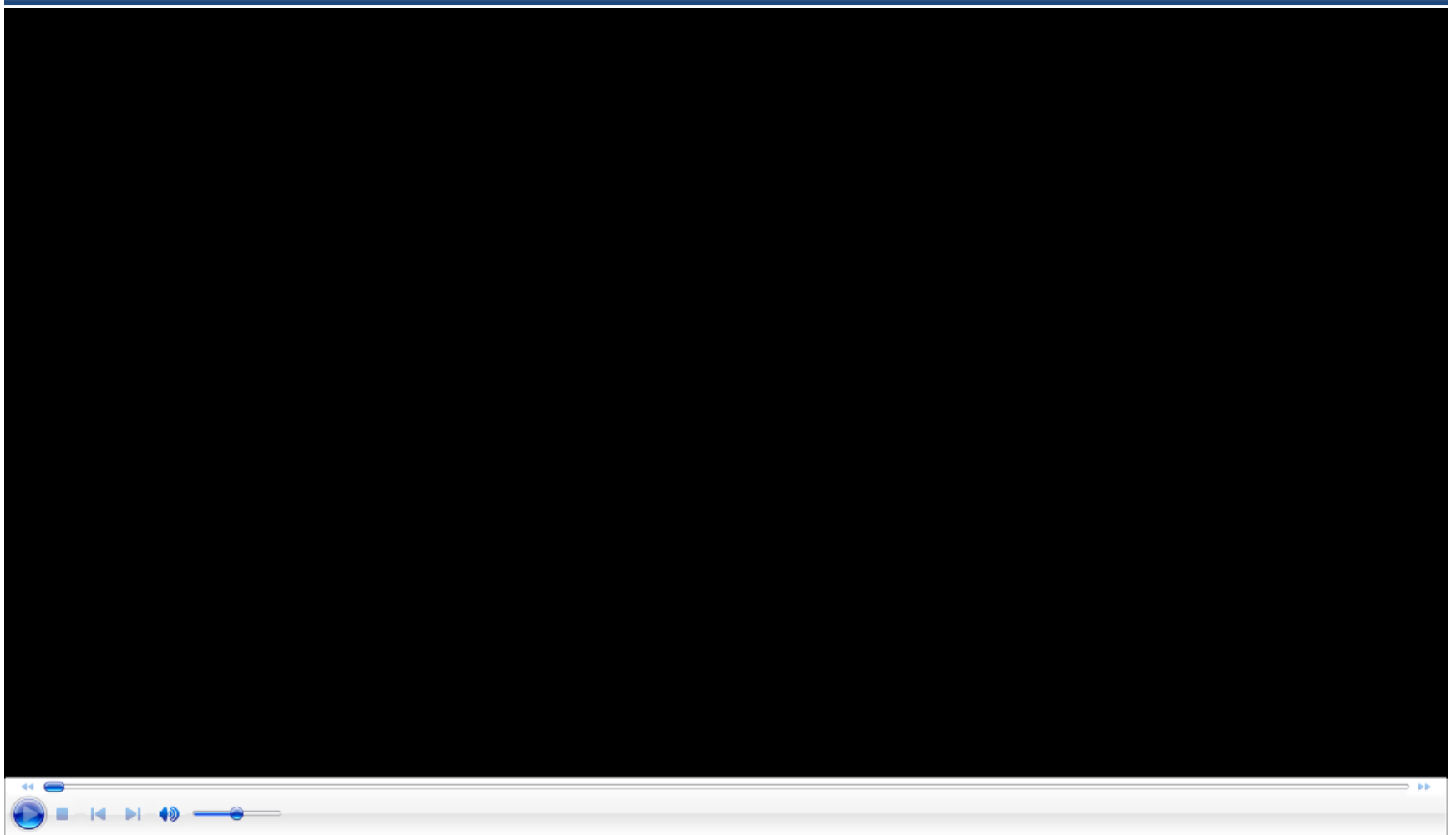


Learned  
Reward



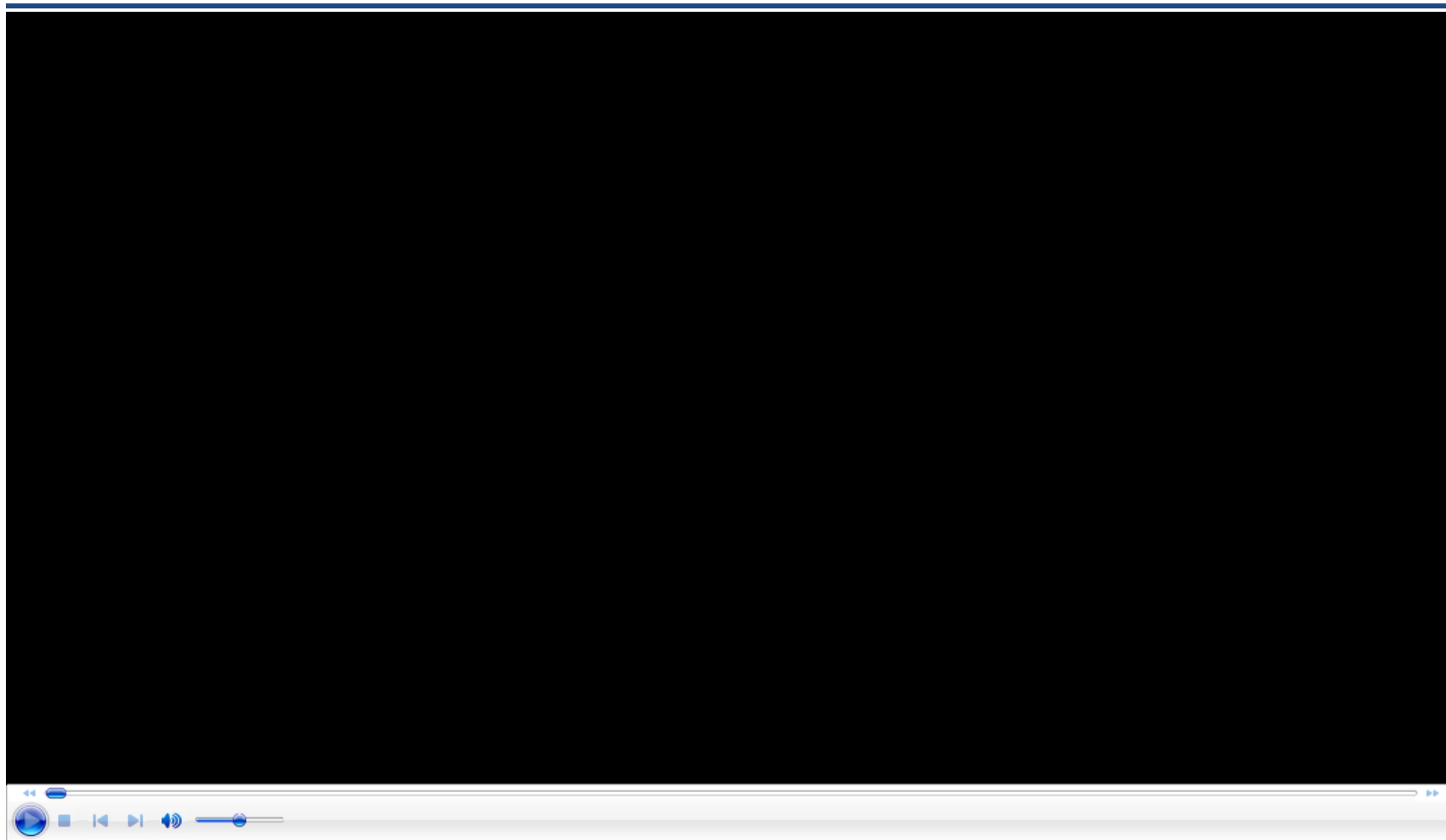
Planning in a new  
Environment

# Car driving (Abbeel et. al)





# Bad driver (Abbeel et. al)



# Recent works

---

1. N. Ratliff, D. Bradley, D. Bagnell, J. Chestnutt. *Boosting Structures Prediction for Imitation Learning*. In NIPS 2007
2. B. Ziebart, A. Mass, D. Bagnell, A. K. Dey. *Maximum Entropy Inverse Reinforcement Learning*. In AAAI 2010
3. P. Abbeel, A. Coats, A. Ng. *Autonomous helicopter aerobatics through apprenticeship learning*. In IJRR 2010
4. S. Ross, G. Gordon, D. Bagnell. *A Reduction of Imitation Learning and Structured Prediction to No-Regret Online Learning*. In AISTATS 2011
5. B. Akgun, M. Cakmak, K. Jiang, A. Thomaz. *Keyframe-based learning from demonstrations*. In IJSR 2012
6. A. Jain, B. Wojcik, T. Joachims, A. Saxena. *Learning Trajectory Preferences for Manipulators via Iterative Improvement*. In NIPS 2013

# Questions

---

- Summary of key points:
  - Reward-based learning can be modeled by the MDP framework.
  - Inverse Reinforcement Learning takes perception features as input and outputs a reward function.
  - The max-margin planning problem can be formulated as a tractable QP.
  - Sub-gradient descent solves an equivalent problem with provable convergence bound.