Learning to Imitate a Fighter Pilot

Shie Mannor

Department of Electrical Engineering Technion

ICML July 2011

S. Mannor (Technion)

Learning to Imitate a Fighter Pilot

ICML July 2011

イロト イポト イヨト イヨト

1/18

э



S. Mannor (Technion)

Learning to Imitate a Fighter Pilo

ICML July 2011 2 / 18

We want to imitate



S. Mannor (Technion)

Learning to Imitate a Fighter Pilc

э

イロト イヨト イヨト イヨト

The system we have



э

イロト イヨト イヨト イヨト

Versatility.

- Accuracy.
- Randomness.
- Multi-agency.
- Concurrency.

э

5/18

(I) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1))

- Versatility.
- Accuracy.
- Randomness.
- Multi-agency.
- Concurrency.

э

(I) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1))

5/18

- Versatility.
- Accuracy.
- Randomness.
- Multi-agency.
- Concurrency.

э

イロト イヨト イヨト イヨト

- Versatility.
- Accuracy.
- Randomness.
- Multi-agency.
- Concurrency.

э

A (10) > A (10) > A (10)

- Versatility.
- Accuracy.
- Randomness.
- Multi-agency.
- Concurrency.

э

A (10) > A (10) > A (10)

RL view

• Single reward: ±1

• State space: Coordinates X, Y, Z, Velocity V_X , V_Y , V_Z Differentials ΔX , ΔY , ΔZ , Velocity ΔV_X , ΔV_Y , ΔV_Z Attack angle, speed, fuel, sun location Missiles/canons Defensive props

• Episodic task: 30-60 seconds per fight

RL view

- Single reward: ±1
- State space:

Coordinates X, Y, Z, Velocity V_X , V_Y , V_Z Differentials ΔX , ΔY , ΔZ , Velocity ΔV_X , ΔV_Y , ΔV_Z Attack angle, speed, fuel, sun location Missiles/canons Defensive props

• Episodic task: 30-60 seconds per fight

< 回 ト < 三 ト < 三

RL view

- Single reward: ±1
- State space:

Coordinates X, Y, Z, Velocity V_X , V_Y , V_Z Differentials ΔX , ΔY , ΔZ , Velocity ΔV_X , ΔV_Y , ΔV_Z Attack angle, speed, fuel, sun location Missiles/canons Defensive props

• Episodic task: 30-60 seconds per fight

< 回 ト < 三 ト < 三

Pilots are good!

- Very fast responses
- Established flight methodology
- Principle: always better to be behind the enemy
- A bit like chess

- Pilots are good!
- Very fast responses
- Established flight methodology
- Principle: always better to be behind the enemy
- A bit like chess

- Pilots are good!
- Very fast responses
- Established flight methodology
- Principle: always better to be behind the enemy
- A bit like chess

(I) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1))

- Pilots are good!
- Very fast responses
- Established flight methodology
- Principle: always better to be behind the enemy
- A bit like chess

不得下 イヨト イヨ

- Pilots are good!
- Very fast responses
- Established flight methodology
- Principle: always better to be behind the enemy
- A bit like chess

不得 とくき とくき

RL: A Naive Approach

- Train agent against itself
- Train agent against a human player
- Temporal credit assignment problem
- Never win against even a moderately good player

A (10) A (10) A (10)

RL: A Naive Approach

- Train agent against itself
- Train agent against a human player
- Temporal credit assignment problem
- Never win against even a moderately good player

A (10) A (10) A (10)

RL: A Naive Approach

- Train agent against itself
- Train agent against a human player
- Temporal credit assignment problem
- Never win against even a moderately good player

Just hack the damn thing

- Very hard to get good at the game
- Pilots training is 3-5 years
- Many maneuvers to consider

(I) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1)) < ((1))

- Just hack the damn thing
- Very hard to get good at the game
- Pilots training is 3-5 years
- Many maneuvers to consider

9/18

A (10) A (10) A (10)

- Just hack the damn thing
- Very hard to get good at the game
- Pilots training is 3-5 years
- Many maneuvers to consider

A (10) A (10) A (10)

- Just hack the damn thing
- Very hard to get good at the game
- Pilots training is 3-5 years
- Many maneuvers to consider

< 3

< A

A sample maneuver



S. Mannor (Technion)

Learning to Imitate a Fighter Pilo

<ロト < 合ト < 言ト < 言ト ICML July 2011

æ

Our approach: Plural Reinforcement Learning

We can learn from many sources:

- Simulation
- Viewing "good" traces of human VS human Annotated or not
- Human VS agent
- Ask pilot what to do?
- Ask pilot what went wrong

11/18

Our approach: Plural Reinforcement Learning

We can learn from many sources:

- Simulation
- Viewing "good" traces of human VS human Annotated or not
- Human VS agent
- Ask pilot what to do?
- Ask pilot what went wrong

11/18

Our approach: Plural Reinforcement Learning

We can learn from many sources:

- Simulation
- Viewing "good" traces of human VS human Annotated or not
- Human VS agent
- Ask pilot what to do?
- Ask pilot what went wrong

Learn sequence of actions = options = macro-actions from traces

Call them maneuvers

- Tried unsupervised learning on traces ⇒ failed. Flight is too versatile.
- Got verbal description from expert: Expert inaccurate Lots of free parameters Articulating maneuvers is hard
- Parameterized maneuvers.

4 E 5

Learn sequence of actions = options = macro-actions from traces

- Call them maneuvers
- Tried unsupervised learning on traces ⇒ failed. Flight is too versatile.
- Got verbal description from expert: Expert inaccurate Lots of free parameters Articulating maneuvers is hard
- Parameterized maneuvers.

Learn sequence of actions = options = macro-actions from traces

- Call them maneuvers
- Tried unsupervised learning on traces ⇒ failed. Flight is too versatile.
- Got verbal description from expert: Expert inaccurate Lots of free parameters Articulating maneuvers is hard
- Parameterized maneuvers.

< E

Learn sequence of actions = options = macro-actions from traces

- Call them maneuvers
- Tried unsupervised learning on traces ⇒ failed. Flight is too versatile.
- Got verbal description from expert: Expert inaccurate Lots of free parameters Articulating maneuvers is hard
- Parameterized maneuvers.

Value function approach

• Illicit value from expert.

- Full advantage
- Full disadvantage
- Equality
-

< 3

э

Value function approach

- Illicit value from expert.
- Full advantage
- Full disadvantage
- Equality
-

13/18

Value function elicitation

- Greedy over expert value function \Rightarrow trembling \Rightarrow smoothing
- Features for linear approximation—unnatural
- Obtaining more constraints from traces
- Value function approach does not seem to cut it

4 Th

Value function elicitation

- Greedy over expert value function ⇒ trembling ⇒ smoothing
- Features for linear approximation—unnatural
- Obtaining more constraints from traces
- Value function approach does not seem to cut it

Value function elicitation

- Greedy over expert value function ⇒ trembling ⇒ smoothing
- Features for linear approximation—unnatural
- Obtaining more constraints from traces
- Value function approach does not seem to cut it

Policy gradients

- Policy is complex
- Parametrization is tricky
- Had some success with isolated maneuvers
- Pilots do not use a pre-specified policy

< 3

▲ 同 ▶ → 三 ▶

Policy gradients

- Policy is complex
- Parametrization is tricky
- Had some success with isolated maneuvers
- Pilots do not use a pre-specified policy

4 E

Reward shaping

- Reward extremely rare
- Use expert for reward signals
- Hard to find a consistent reward function
- Randomization
- Risk aversion

Reward shaping

- Reward extremely rare
- Use expert for reward signals
- Hard to find a consistent reward function
- Randomization
- Risk aversion

Reward shaping

- Reward extremely rare
- Use expert for reward signals
- Hard to find a consistent reward function
- Randomization
- Risk aversion

16/18

Using traces for better features

Main challenge: when to start which maneuver and with which params.

How to combine all forms of input from pilots?

• How to use rare feedback from expert?

- Stop and ask
- Initiating and defining maneuvers
- Better use of traces
- Reward shaping/learning maneuvers

18/18

< 6 b

- How to use rare feedback from expert?
- Stop and ask
- Initiating and defining maneuvers
- Better use of traces
- Reward shaping/learning maneuvers

< 6 b

18/18

- How to use rare feedback from expert?
- Stop and ask
- Initiating and defining maneuvers
- Better use of traces
- Reward shaping/learning maneuvers

- How to use rare feedback from expert?
- Stop and ask
- Initiating and defining maneuvers
- Better use of traces
- Reward shaping/learning maneuvers

18/18

- How to use rare feedback from expert?
- Stop and ask
- Initiating and defining maneuvers
- Better use of traces
- Reward shaping/learning maneuvers