

Machine Learning

Reinforcement Learning and Decisions

Carl Henrik Ek - carlhenrik.ek@bristol.ac.uk

November 21, 2017

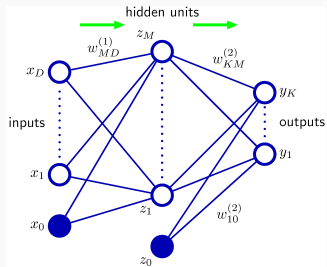
<http://www.carlhenrik.com>

Introduction

Today

- Last lecture of material
- Second hour: doing an ml project

Neural Network



$$y_k(\mathbf{x}, \mathbf{w}) = \sigma \left(\sum_j^M w_{kj}^{(2)} h \left(\sum_j^D w_{ji}^{(1)} x_i + w_{j0}^{(1)} \right) + w_{k0}^{(2)} \right)$$

Why are composite functions attractive?

$$y = g(x) = f_K(f_{K-1}(f_{K-2}(\dots f_1(x) \dots)))$$

- Kernel of a function

$$\text{Kern}(f_k) = \{(\mathbf{x}, \mathbf{x}') | f_k(\mathbf{x}) = f_k(\mathbf{x}')\}$$

- Kernel of a function

$$\text{Kern}(f_k) = \{(\mathbf{x}, \mathbf{x}') | f_k(\mathbf{x}) = f_k(\mathbf{x}')\}$$

- Image of a function

$$\text{Im}(f_k(\mathbf{x})) = \{\mathbf{y} \in Y | \mathbf{y} = f_k(\mathbf{x}), \mathbf{x} \in X\}$$

- Kernel of function

$$\text{Kern}(f_1) \subseteq \text{Kern}(f_{k-1} \circ \dots \circ f_2 \circ f_1) \subseteq \text{Kern}(f_k \circ f_{k-1} \circ \dots \circ f_2 \circ f_1)$$

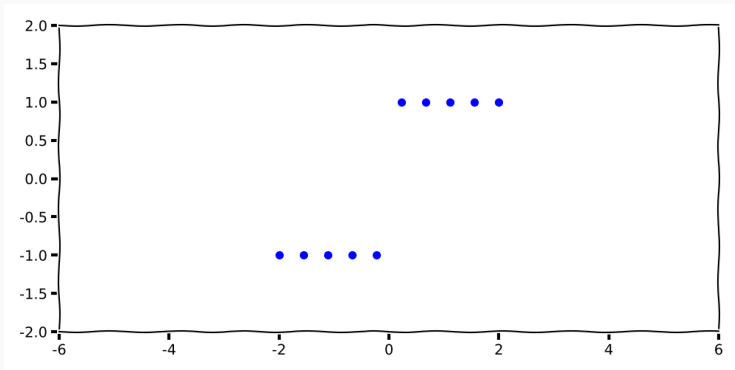
- Kernel of function

$$\text{Kern}(f_1) \subseteq \text{Kern}(f_{k-1} \circ \dots \circ f_2 \circ f_1) \subseteq \text{Kern}(f_k \circ f_{k-1} \circ \dots \circ f_2 \circ f_1)$$

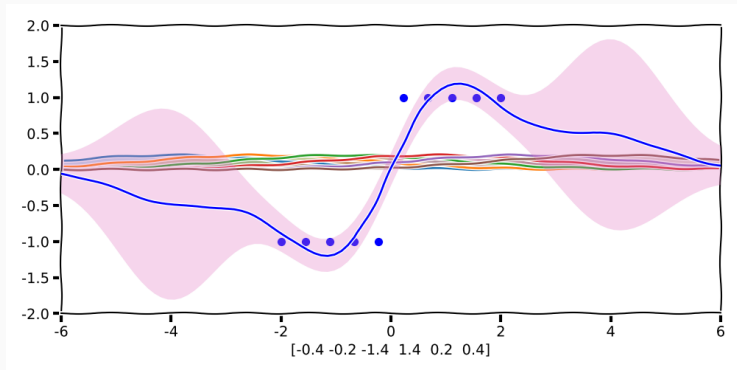
- Image of a function

$$\text{Im}(f_k \circ f_{k-1} \circ \dots \circ f_2 \circ f_1) \subseteq \text{Im}(f_k \circ f_{k-1} \circ \dots \circ f_2) \subseteq \dots \subseteq \text{Im}(f_k)$$

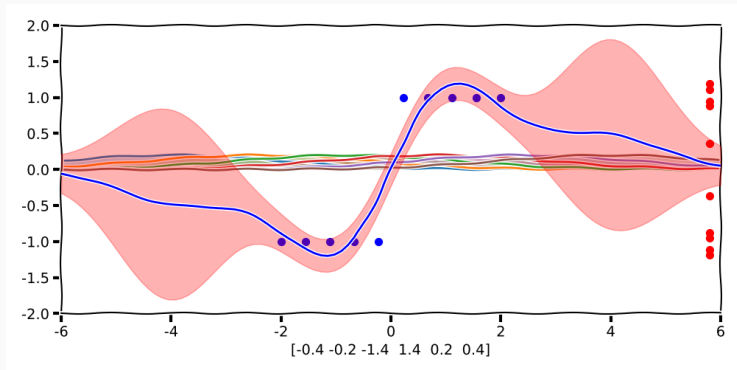
Composite Functions



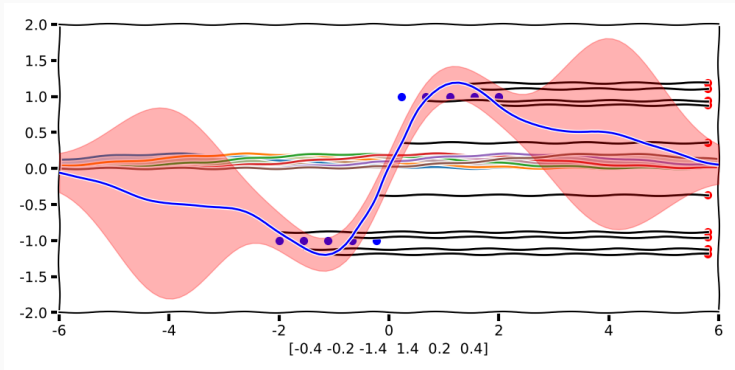
Composite Functions



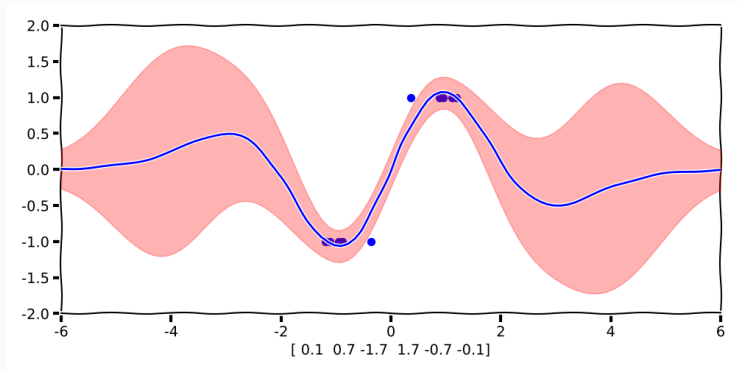
Composite Functions



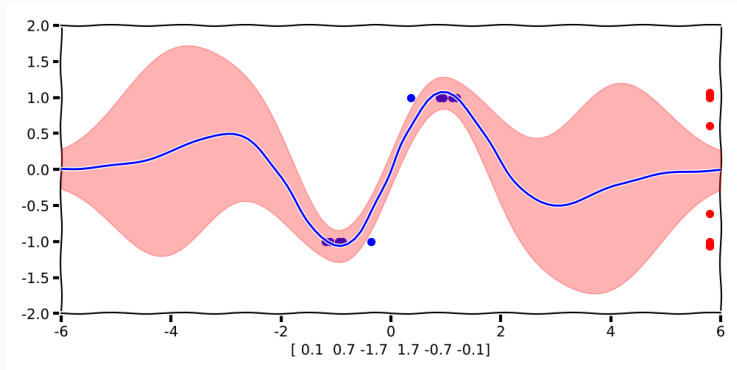
Composite Functions



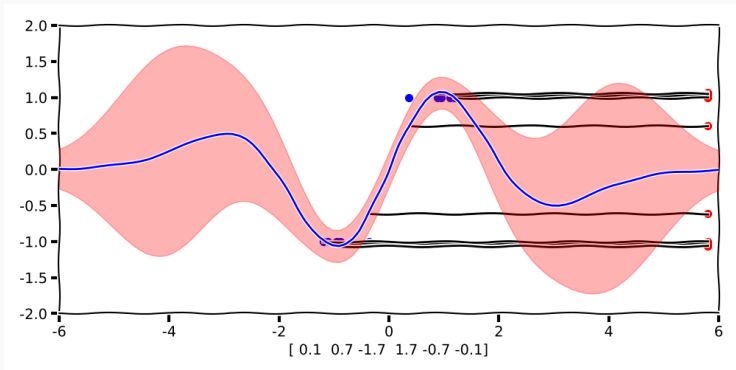
Composite Functions



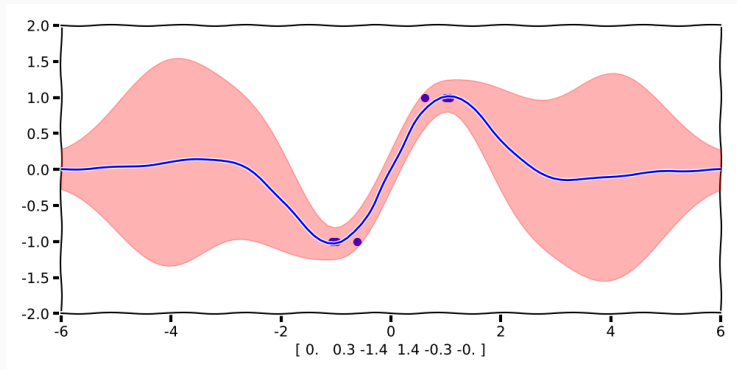
Composite Functions



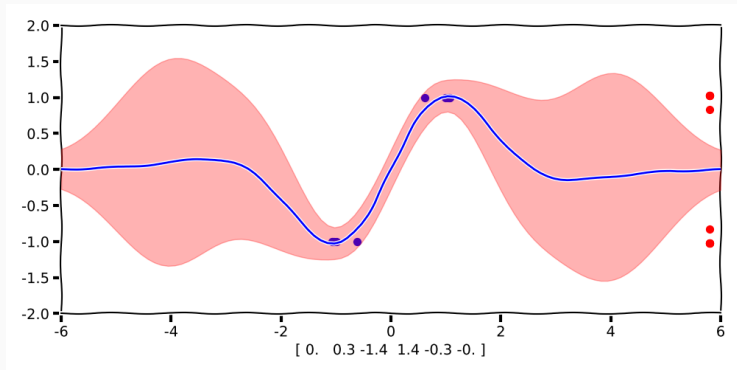
Composite Functions



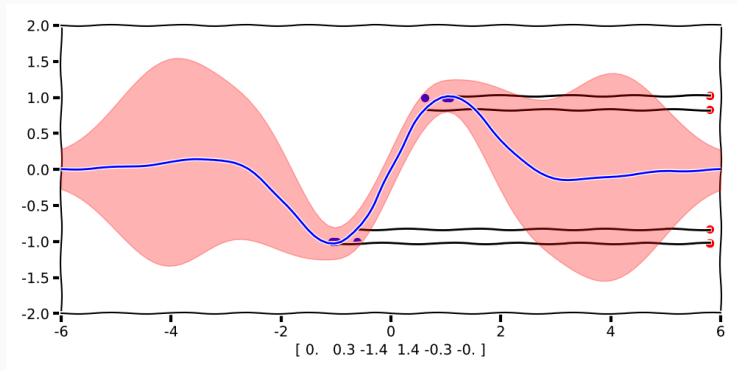
Composite Functions



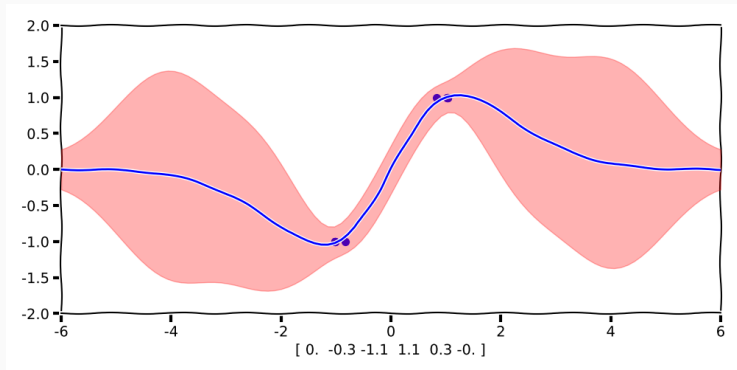
Composite Functions



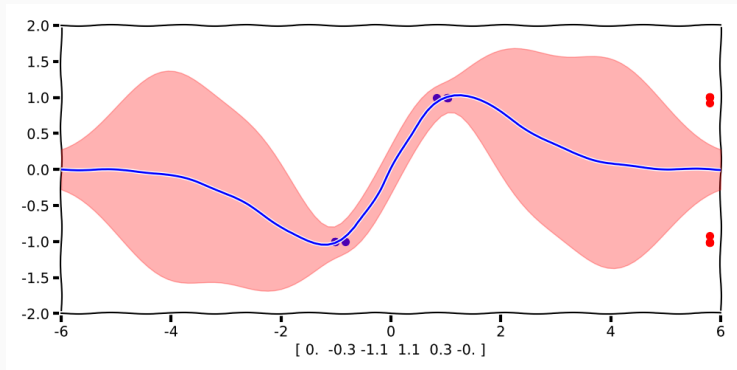
Composite Functions



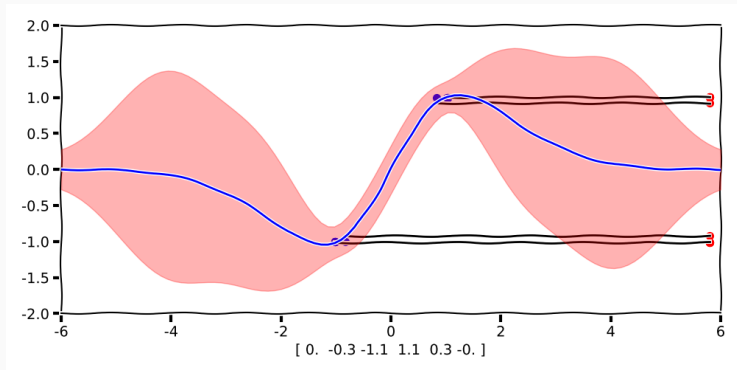
Composite Functions



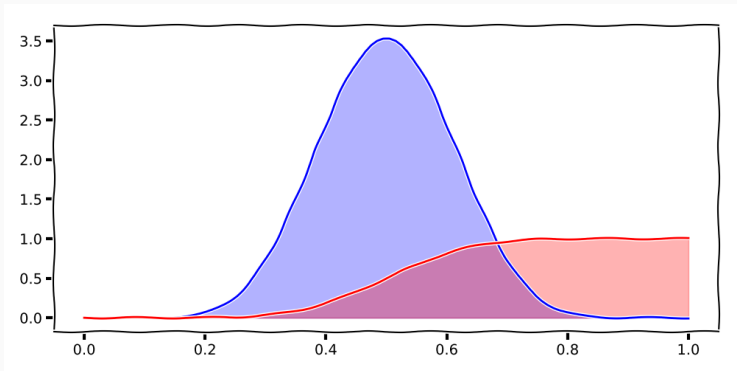
Composite Functions



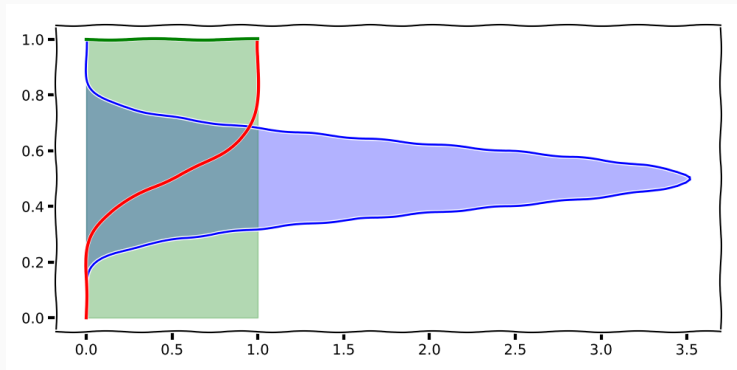
Composite Functions



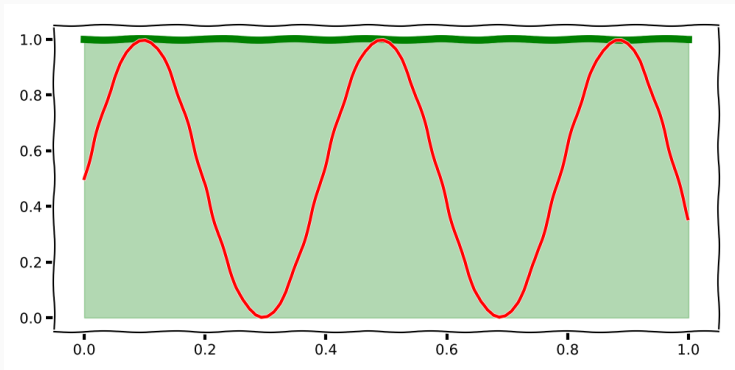
Sampling



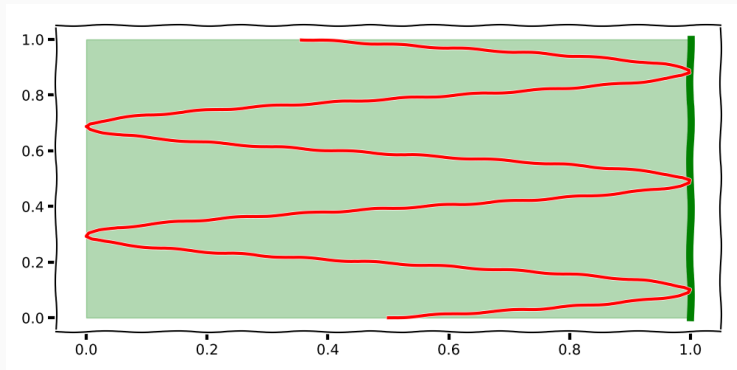
Sampling



Change of Variables



Change of Variables



Reinforcement Learning

Supervised Learning predict output from input

$$\mathcal{D} = \{y_i, x_i\}_{i=1}^N$$

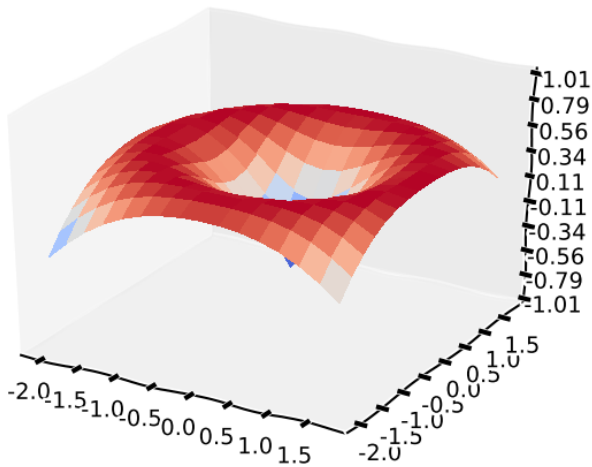
$$p(\mathbf{y}|\mathbf{x}, \theta)$$

Unsupervised Learning model the data

$$\mathcal{D} = \{y_i\}_{i=1}^N$$

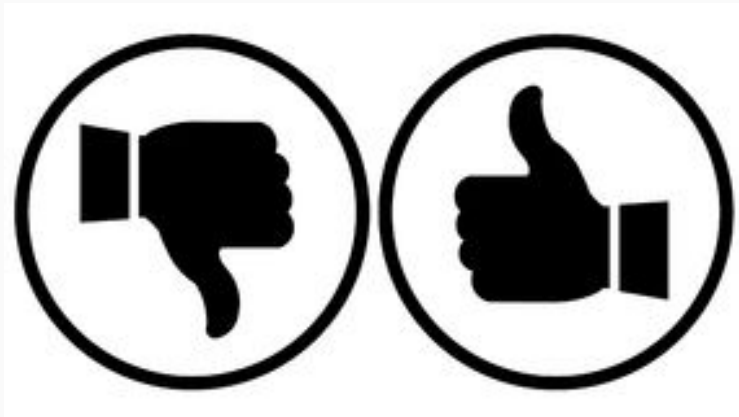
$$p(\mathbf{y}|\theta)$$

Structure



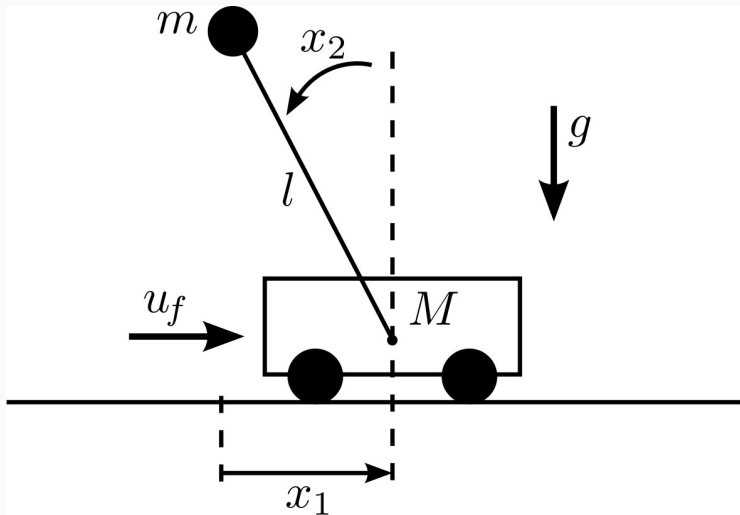
Bicycles





Can we learn without specifying how the task should be achieved by providing, rewards (positive) and punishment (negative)?

Inverted Pendulum¹

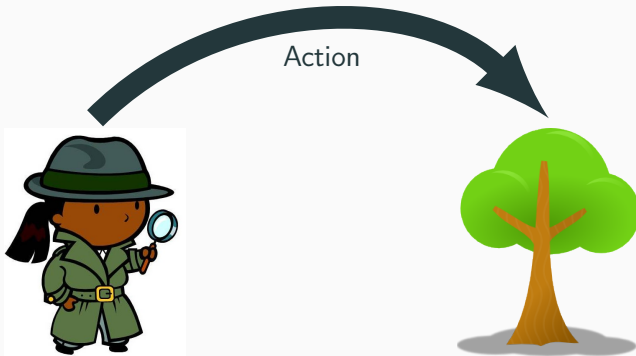


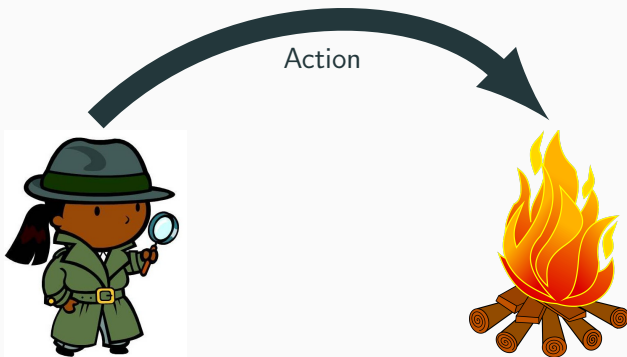
¹<https://www.youtube.com/watch?v=XiigTGKZfks>

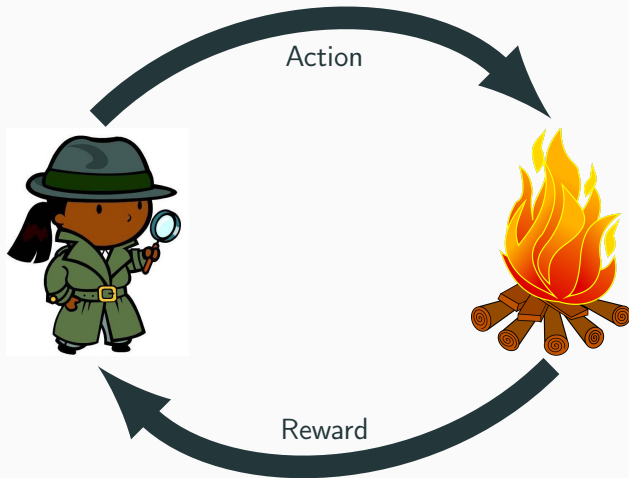
Pong











Formalism [1]

\mathcal{S} a discrete set of environment states

\mathcal{A} a discrete set of agent actions

r a scalar set of reinforcement signals, real line, or $\{0, 1\}$

I input function, how the agent views the state of the environment

π policy, mapping from state to action that maximises a long-run measurement of reinforcement

- The world is non-deterministic
 - we can be in the same state and do the same action and different things happens
- The world is stationary
 - the probabilities of the state transitions do not change
- Input function
 - if the agent can see the state of the world we call this fully observable
 - if the agent can only see part of the state, its partially observable

Example

Environment you are in state 65 you have 4 possible actions

Example

Environment you are in state 65 you have 4 possible actions

Agent I take action 2

Example

Environment you are in state 65 you have 4 possible actions

Agent I take action 2

Environment you recieved reinforcement of 7 units, you are now
in state 15 you have 2 possible actions

Example

Environment you are in state 65 you have 4 possible actions

Agent I take action 2

Environment you recieved reinforcement of 7 units, you are now
in state 15 you have 2 possible actions

Agent I take action 1

Example

Environment you are in state 65 you have 4 possible actions

Agent I take action 2

Environment you recieved reinforcement of 7 units, you are now
in state 15 you have 2 possible actions

Agent I take action 1

Environment you recieved reinforcement of -4 units, you are now
in state 65 you have 4 possible actions

Example

Environment you are in state 65 you have 4 possible actions

Agent I take action 2

Environment you recieved reinforcement of 7 units, you are now
in state 15 you have 2 possible actions

Agent I take action 1

Environment you recieved reinforcement of -4 units, you are now
in state 65 you have 4 possible actions

Agent I take action 2

Example

Environment you are in state 65 you have 4 possible actions

Agent I take action 2

Environment you recieved reinforcement of 7 units, you are now
in state 15 you have 2 possible actions

Agent I take action 1

Environment you recieved reinforcement of -4 units, you are now
in state 65 you have 4 possible actions

Agent I take action 2

Environment you recieved reinforcement of 5 units, you are now
in state 44 you have 5 possible actions

Optimal Behaviour



- Finite time horizon

$$E \left[\sum_{t=0}^h r_t \right]$$

Optimal Behaviour

- Finite time horizon

$$E \left[\sum_{t=0}^h r_t \right]$$

- Average reward

$$\lim_{h \rightarrow \infty} E \left[\frac{1}{h} \sum_{t=0}^h r_t \right]$$

Optimal Behaviour

- Finite time horizon

$$E \left[\sum_{t=0}^h r_t \right]$$

- Average reward

$$\lim_{h \rightarrow \infty} E \left[\frac{1}{h} \sum_{t=0}^h r_t \right]$$

- Infinite horizon (discounted reward)

$$E \left[\sum_{t=0}^{\infty} \gamma^t r_t \right]$$

$$0 \leq \gamma \leq 1$$

- State Transition

$$p(s(t)|s(t-1), a(t))$$

- Input function

$$p(y(t)|x(t), a(t))$$

- Reward function

$$r(t) = E(s(t), a(t), s(t+1))$$

- Policy

$$a(t) = \pi(s(t))$$

Markov Decision Process

- Fully observable system

$$x(t) = y(t) = s(t)$$

- Transition matrix

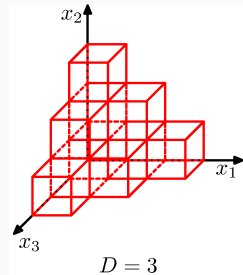
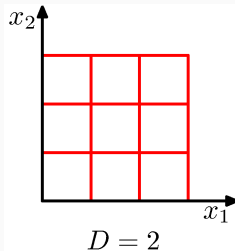
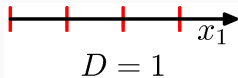
$$T(s, a, s') = p(s(t+1) = s' | s(t) = s, a(t) = a)$$

- Reward matrix

$$R(s, a, s') = E(s(t) = s, a(t) = a, s(t+1) = s')$$

- For this set-up the optimal policy can be computed using Dynamic Programming

Markov Decision Process



We cannot enumerate all the states and actions

Reinforcement Learning

- As we cannot evaluate everything we want to learn from data
- Reinforcement function
 - mapping between state and action pair to reward
- Environment function
 - mapping between observables to state
- Value function
 - the total reward of a specific policy

Reinforcement Learning

- We can use any available machine learning method to learn these functions
- Gaussian processes
- Linear regression
- Composite functions (neural networks)
- The tricky thing is how to get the data?

Bandits



Bandit problems



- You are in a room with k slot machines
- Each have a different (unknown) probability of pay-off
- You are permitted h different pulls
- Whats the optimal strategy?

Exploration vs. Exploitation

- When thinking of a strategy to come up with a policy we need to balance
 - exploring the environment
 - exploiting what we know

Exploration vs. Exploitation

- When thinking of a strategy to come up with a policy we need to balance
 - exploring the environment
 - exploiting what we know
- Re-inforcement learning can be thought of in terms of bandit problems
 - sequential bandits
 - delayed reward
 - betting in poker (I think)

Exploration vs. Exploitation

- When thinking of a strategy to come up with a policy we need to balance
 - exploring the environment
 - exploiting what we know
- Re-inforcement learning can be thought of in terms of bandit problems
 - sequential bandits
 - delayed reward
 - betting in poker (I think)
- Designing this strategy is the main challenge

<https://www.youtube.com/watch?v=faDKMMwOS2Q>

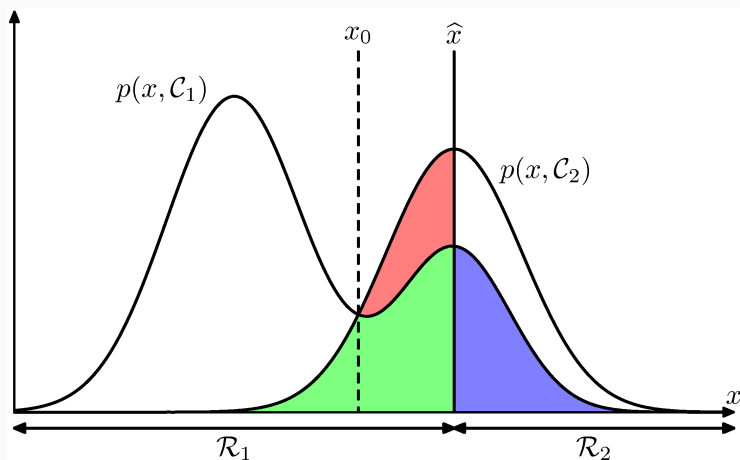
`https://youtu.be/QHcAlAprFxA`

Decisions

- So far we have just described how to model data
 1. make assumptions
 2. derive posterior
- Everything has been stochastic as we propagate uncertainty
- *you can't marginalise over menu items in a restaurant*

²<http://inverseprobability.com/2017/11/15/decision-making>

Decisions



Loss Functions

	Cancer	\neg Cancer
Cancer	0	100
\neg Cancer	1	0

$$\mathbb{E}[L] = \sum_k \sum_j \int_{\mathcal{R}_j} L_{kj} p(\mathbf{x}, C_k) d\mathbf{x}$$

Utilitarian theory

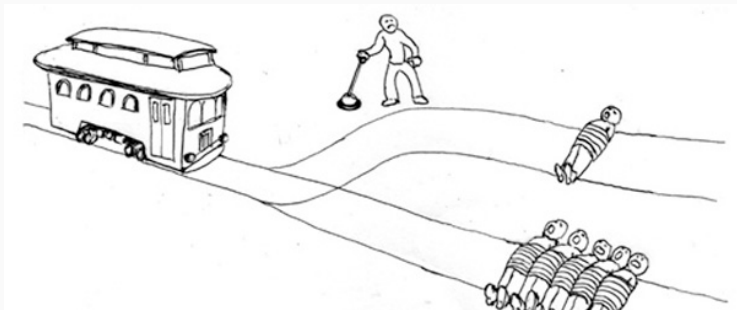
Utilitarianism is an ethical theory which states that the best action is the one that maximizes utility. "Utility" is defined in various ways, usually in terms of the well-being of sentient entities. Jeremy Bentham, the founder of utilitarianism, described utility as the sum of all pleasure that results from an action, minus the suffering of anyone involved in the action.

– Wikipedia

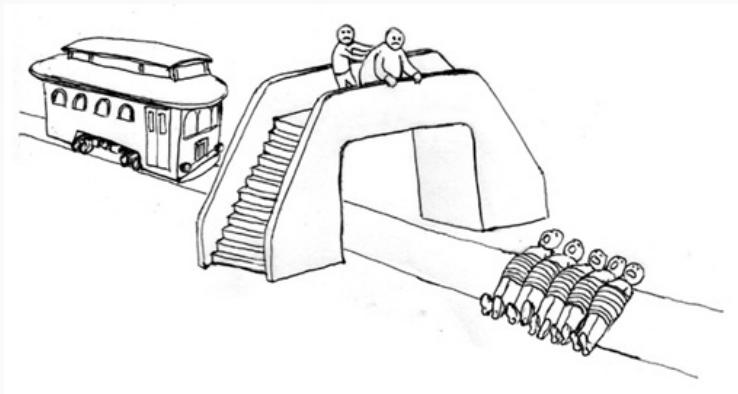
Utilitarian theory

- We can make formulate decisions as a mathematical principle
- Very attractive as it allows decisions to be made explicitly it allows for accountability
- Does it work?

Trolley



Trolley



	Cancer (1%)	\neg Cancer (99%)
Positive	80%	9.6%
Negative	20%	90.4%

- *What is the probability that you have cancer given a positive test?*

$$p(\text{cancer}=\text{true}|\text{test}=\text{pos}) = \frac{p(\text{test}=\text{pos}|\text{cancer}=\text{true})p(\text{cancer}=\text{true})}{p(\text{test}=\text{pos})}$$

$$p(\text{cancer}=\text{true}|\text{test}=\text{pos}) = \frac{p(\text{test}=\text{pos}|\text{cancer}=\text{true})p(\text{cancer}=\text{true})}{p(\text{test}=\text{pos})}$$

$$\begin{aligned} p(\text{test}=\text{pos}) &= \int p(\text{test}=\text{pos}|\text{cancer})p(\text{cancer}) = \\ & p(\text{test}=\text{pos}|\text{cancer}=\text{false})p(\text{cancer}=\text{false}) + \\ & p(\text{test}=\text{pos}|\text{cancer}=\text{true})p(\text{cancer}=\text{true}) \end{aligned}$$

$$p(\text{cancer}=\text{true}|\text{test}=\text{pos}) = \frac{p(\text{test}=\text{pos}|\text{cancer}=\text{true})p(\text{cancer}=\text{true})}{p(\text{test}=\text{pos})}$$

$$p(\text{test}=\text{pos}) = \int p(\text{test}=\text{pos}|\text{cancer})p(\text{cancer}) = \\ p(\text{test}=\text{pos}|\text{cancer}=\text{false})p(\text{cancer}=\text{false}) + \\ p(\text{test}=\text{pos}|\text{cancer}=\text{true})p(\text{cancer}=\text{true})$$

$$p(\text{cancer}=\text{true}|\text{test}=\text{pos}) = \frac{0.8 \cdot 0.01}{0.096 \cdot 0.99 + 0.8 \cdot 0.01} = 0.078$$

$$p(\text{cancer}=\text{true}|\text{test}=\text{pos}) = \frac{p(\text{test}=\text{pos}|\text{cancer}=\text{true})p(\text{cancer}=\text{true})}{p(\text{test}=\text{pos})}$$

$$p(\text{test}=\text{pos}) = \int p(\text{test}=\text{pos}|\text{cancer})p(\text{cancer}) = \\ p(\text{test}=\text{pos}|\text{cancer}=\text{false})p(\text{cancer}=\text{false}) + \\ p(\text{test}=\text{pos}|\text{cancer}=\text{true})p(\text{cancer}=\text{true})$$

$$p(\text{cancer}=\text{true}|\text{test}=\text{pos}) = \frac{0.8 \cdot 0.01}{0.096 \cdot 0.99 + 0.8 \cdot 0.01} = 0.078$$

- Only 15% of (medical) doctor answers this correctly

Self driving cars



- We might argue that there is such a thing as a single consistent utility function
 - Most likely we will then over simplify things
- Different people have different utility functions
 - what makes me happy does not make Donald Trump happy
- We need to treat utility functions with uncertainty

$$p(\text{decision}) = \int p(\text{decision} | \text{utility function}) p(\text{utility function})$$

- You know how to do this, you know the theory of this, you can move this up to any level

Summary

- Reinforcement learning
 - how can we learn when we do not know how to do something
 - exploration exploitation
 - if we knew all states, all rewards, all actions it can be done exact
- Decisions
 - one single utility functions is a too simplistic assumption
 - add uncertainty and marginalise

eof

References



Leslie Pack Kaelbling, Michael L. Littman, and Andrew W. Moore.

Reinforcement learning: A survey.

J. Artif. Intell. Res., 4:237–285, 1996.