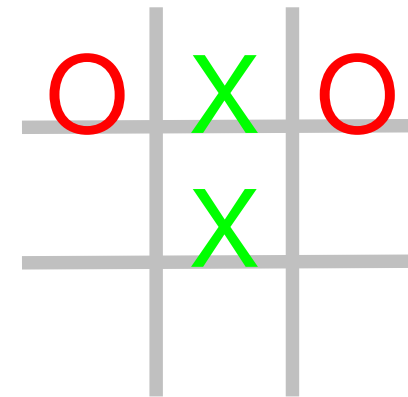# Reinforcement Learning

- Introduction
  - MENACE (Michie 1963)
- Formalization
  - Policies
  - Value Function
  - Q-Function
- Model-based Reinforcement Learning
  - Policy Iteration
  - Value Iteration
- Model-free Reinforcement Learning
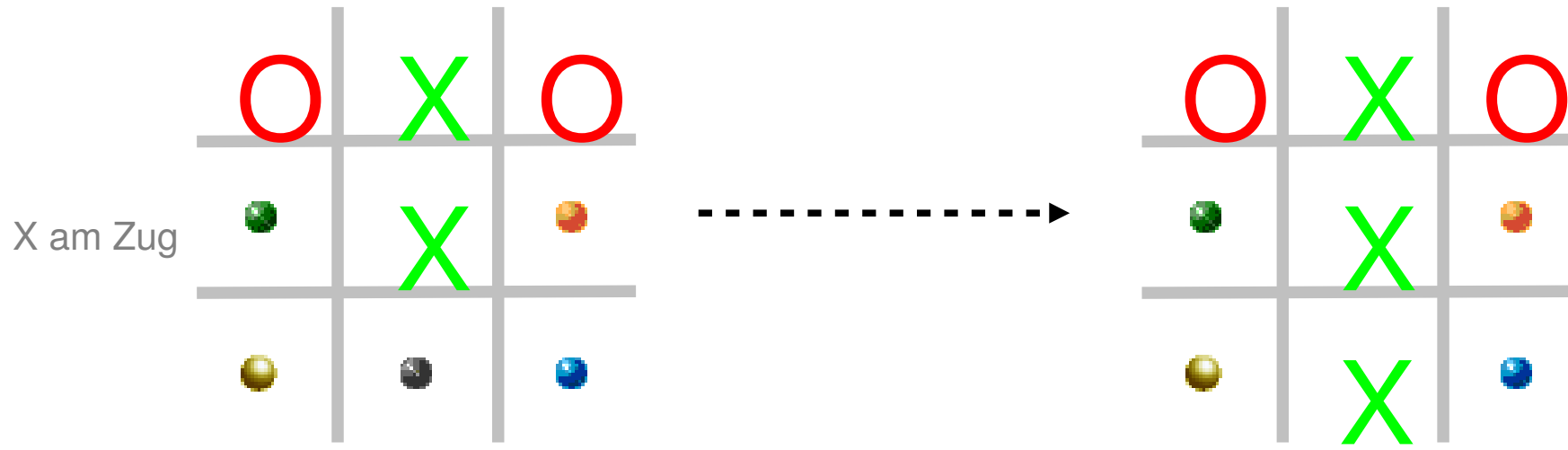  - Q-Learning
  - extensions
- Application Examples

# Reinforcement Learning

- ## Ziel:
  - Lernen von (guten) Entscheidungen durch Feedback (Reinforcement) der Umwelt (z.B. Spiel gewonnen/verloren).

- ## Anwendungen:
  - **Spiele:**
    - Tic-Tac-Toe: MENACE (Michie 1963)
    - Backgammon: TD-Gammon (Tesauro 1995)
    - Schach: KnightCap (Baxter et al. 2000)
  - **Andere:**
    - Elevator Dispatching
    - Robot Control
    - Job-Shop Scheduling
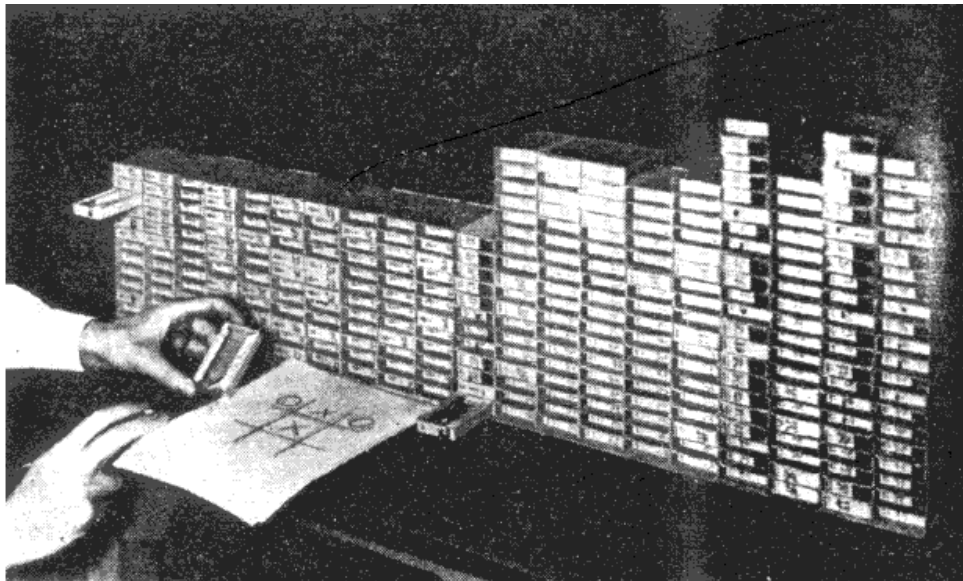
# MENACE (Michie, 1963)

- **Lernt Tic-Tac-Toe zu spielen**

- **Hardware:**
  - 287 Zündholzschachteln
    (1 für jede Stellung)
  - Perlen in 9 verschiedenen Farbe
    (1 Farbe für jedes Feld)

- **Spiel-Algorithmus:**
  - Wähle Zündholzschachtel, die der Stellung entspricht
  - Ziehe zufällig eine der Perlen
  - Ziehe auf das Feld, das der Farbe der Perle entspricht

- **Implementation:** http://www.codeproject.com/KB/cpp/ccross.aspx

X am Zug

Zur Stellung passende
Schachtel auswählen

Den der Farbe der
gezogenen Kugel
entsprechenden
Zug ausführen

Eine Kugel aus
der Schachtel ziehen

# Reinforcement Learning in MENACE

- **Initialisierung**
  - alle Züge sind gleich wahrscheinlich, i.e., jede Schachtel enthält gleich viele Perlen für alle möglichen Züge
- **Lern-Algorithmus:**
  - Spiel verloren → gezogene Perlen werden einbehalten (*negative reinforcement*)
  - Spiel gewonnen → eine Perle der gezogenen Farbe wird in verwendeten Schachteln hinzugefügt (*positive reinforcement*)
  - Spiel remis → Perlen werden zurückgelegt (keine Änderung)
- **führt zu**
  - Erhöhung der Wahrscheinlichkeit, daß ein erfolgreicher Zug wiederholt wird
  - Senkung der Wahrscheinlichkeit, daß ein nicht erfolgreicher Zug wiederholt wird

# Credit Assignment Problem

- **Delayed Reward**
  - Der Lerner merkt erst am Ende eines Spiels, daß er verloren (oder gewonnen) hat
  - Der Lerner weiß aber nicht, welcher Zug den Verlust (oder Gewinn verursacht hat)
    - oft war der Fehler schon am Anfang des Spiels, und die letzten Züge waren gar nicht schlecht

- **Lösung in Reinforcement Learning:**
  - Alle Züge der Partie werden belohnt bzw. bestraft (Hinzufügen bzw. Entfernen von Perlen)
  - Durch zahlreiche Spiele konvergiert dieses Verfahren
    - schlechte Züge werden seltener positiv verstärkt werden
    - gute Züge werden öfter positiv verstärkt werden

# Reinforcement Learning - Formalization

- **Learning Scenario**
  - $s \in S$ : state space
  - $a \in A$  : action space
  - $s_0 \in S_0$ : initial states
  - a state transition function  $\delta : S \times A \to S$
  - a reward function  $r : S \times A \to \mathbb{R}$

- **Markov property**
  - rewards and state transitions only depend on last state
  - not on how you got into this state

# Reinforcement Learning - Formalization

- State and action space can be
  - Discrete: $S$ and/or $A$ is a set
  - Continuous: $S$ and/or $A$ are infinite (not part of this lecture!)

- State transition function can be
  - Stochastic: Next state is drawn according to $\delta(s'|s,a)$
  - Deterministic: Next state is fixed $\delta(s,a) = s'$

# Reinforcement Learning - Formalization

- Enviroment:

  - the agent repeatedly chooses an action according to some *policy* $\pi(a|s)$ *or* $\pi(s) = a$

  - this will put the agent in state s into a new state s' according to

    stochastic: $\mathrm{Pr}^{\pi}(s'|s) = \delta(s'|s,a)\pi(a|s)$
    deterministic: $s' = \delta(s, \pi(s))$

  - in some states the agent receives feedback from the environment (reinforcement)

# MENACE - Formalization

- Framework
  - states = matchboxes, discrete
  - actions = moves/beads, discrete
  - policy = prefer actions with higher number of beads, stochastic
  - reward = game won/ game lost
    - *delayed* reward: we don't know right away whether a move was good or bad+
  - transition function: choose next matchbox according to rules, deterministic

- Task
  - Find a policy that maximizes the sum of future rewards

# Learning Task

- ## delayed reward
    - reward for actions may not come immediately (e.g., game playing)
    - modeled as: every state $s_i$ gives a reward $r_i$, but most $r_i=0$
- ## goal: maximize cumulative reward (return) for trajectories a policy is generating
    - reward from "now" until the end of time

$$R(\pi) = R(\tau^\pi) = \sum_{t=0}^{\infty} \gamma^t r(s_t, a_t)$$

    - immediate rewards are weighted higher, rewards further in the future are discounted (discount factor $\gamma$)
    - sum to infinty could be infinite without discount

# Learning Task

- How can we compute $R(\tau^\pi)$ ?

$$R(\tau^\pi) = \sum_{t=0}^{\infty} \gamma^t r(s_t, a_t)$$

$$= r(s_0, a_0) + \gamma r(s_1, a_1) + \gamma^2 r(s_2, a_2) \cdots$$

$$= r(s_0, \pi(s_0)) + \sum_{t=1}^{\infty} \gamma^t r(\delta(s_{t-1}, \pi(s_{t-1})), \pi(s_t))$$

$$= V^\pi(s_0)$$

- A deterministic policy and transition function creates a single trajectory.
- Sum the observed rewards (with decay)
- Also called the value for the first state
- Value function = return when starting in state s and following policy π afterwards

# Optimal Policies and Value Functions

- ## Optimal policy

  - the policy with the highest expected value for all states

  $$\pi^*(s) = \arg\max_{\pi} V^{\pi}(s)$$

  $$= \arg\max_{a \in A} r(s,a) + \gamma V^{\pi^*}(\delta(s,a))$$

  - Always select the action that maximizes the value function for the next step, when following the optimal policy afterwards

  - But we don't know the optimal policy...

# Policy Iteration

- **Policy Improvement Theorem**
  - if it is true that selecting the first action in each state according to a policy $\pi'$ and continuing with policy $\pi$ is better than always following $\pi$ then $\pi'$ is a better policy than $\pi$

- **Policy Improvement**
  - always select the action that maximizes the value function of the current policy
  $$\pi'(s) = \arg\max_{a \in A} r(s,a) + \gamma V^{\pi}(\delta(s,a))$$

- **Policy Evaluation**
  - Compute the value function for the new policy

- **Policy Iteration**
  - Interleave steps of policy evaluation with policy improvement
  $$\pi^0(s) \to V^{\pi^0}(s) \to \pi^1(s) \to \cdots \to \pi^*(s)$$
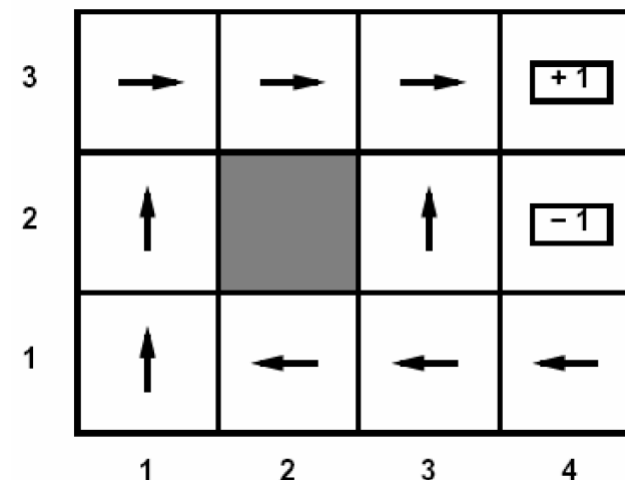
# Policy Evaluation

- We need the value of all states, but can only initiate in $S_0$

  - Update all states along the trajectory

- We assumed the transition function to be deterministic, that is not realistic in many settings

  - Monte Carlo approximation
  - Create k samples and average

$$V^{\pi}(s_0) = \mathbb{E}_{s_t} \sum_{t=0}^{\infty} \gamma^t r(s_t, a_t)$$

$$= r(s_0, \pi(s_0)) \sum_{t=1}^{\infty} \gamma^t \mathbb{E}_{s_t} \delta(s_t | s_{t-1}, \pi(s_{t-1})) r(s_t, \pi(s_t))$$

$$= r(s_0, \pi(s_0)) + \frac{1}{k} \sum_{i=0}^{k} \sum_{t=1}^{\infty} \gamma^t r(s_t^i, \pi(s_t^i))$$

# Policy Evaluation - Example

- ## Simplified task
  - ### we don't know $\delta$
  - ### we don't know $r$
  - ### but we are given a policy $\pi$
    - #### i.e., we have a function that gives us an action in each state

- ## Goal:
  - ### learn the value of each states

- ## Note:
  - ### here we have no choice about the actions to take
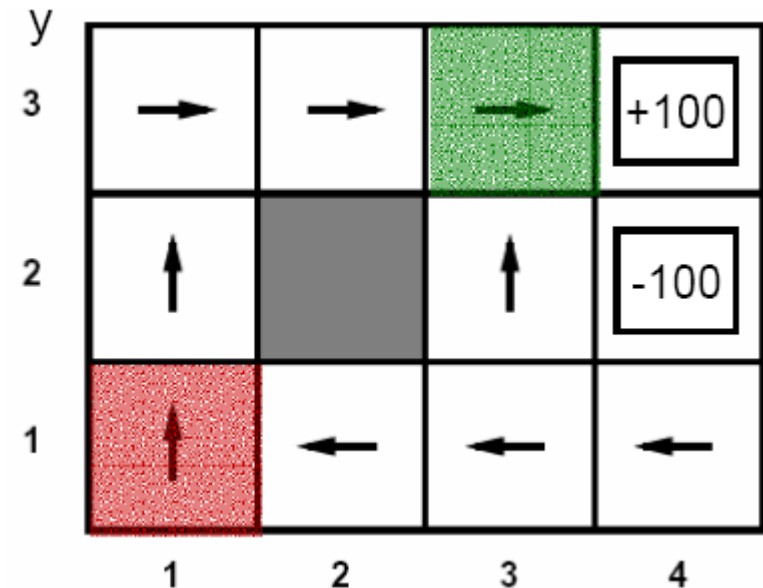  - ### we just execute the policy and observe what happens

# Policy Evaluation – Example

- Episodes:

| (1,1) up -1 | (1,1) up -1 |
|---|---|
| (1,2) up -1 | (1,2) up -1 |
| (1,2) up -1 | (1,3) right -1 |
| (1,3) right -1 | (2,3) right -1 |
| (2,3) right -1 | (3,3) right -1 |
| (3,3) right -1 | (3,2) up -1 |
| (3,2) up -1 | (4,2) exit -100 |
| (3,3) right -1 | (done) |
| (4,3) exit +100 | |
| (done) | |

Transitions are indeterministic!

$\gamma = 1,$

$$V^{\pi}(1,1) \leftarrow (92 + -106)/2 = -7$$

$$V^{\pi}(3,3) \leftarrow (99 + 97 + -102)/3 = 31.3$$

# Policy Improvement

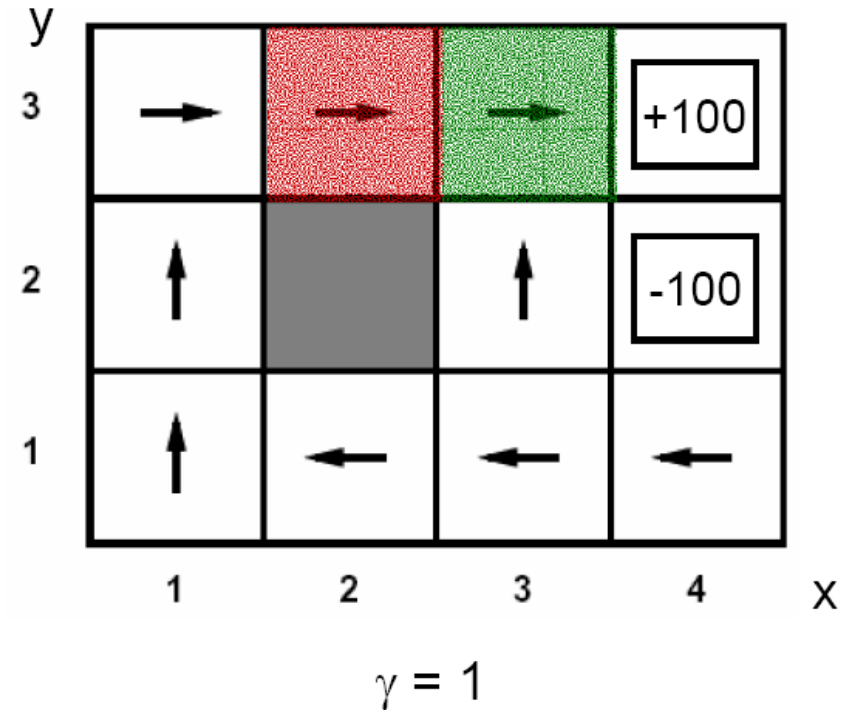- Compute the value for every state
- Update the policy according to

$$\pi'(s) = \arg\max_{a \in A} r(s, a) + \gamma \mathbb{E}_{s'} \delta(s'|s, a) V^\pi(s')$$

- But here we need the transition function we don't know ?

# Simple Approach:
# Learn the Model from Data

- ## Episodes:

| | |
|---|---|
| (1,1) up -1 | (1,1) up -1 |
| (1,2) up -1 | (1,2) up -1 |
| (1,2) up -1 | (1,3) right -1 |
| (1,3) right -1 | (2,3) right -1 |
| (2,3) right -1 | (3,3) right -1 |
| (3,3) right -1 | (3,2) up -1 |
| (3,2) up -1 | (4,2) exit -100 |
| (3,3) right -1 | (done) |
| (4,3) exit +100 | |
| (done) | |



$\gamma = 1$

$\mathbf{P}((4,3)|(3,3),\mathrm{right})=1/3$

$\mathbf{P}((3,3)|(2,3),\mathrm{right})=2/2$

But do we really need to learn the transition model?

# Q-function

- the Q-function does not evaluate states, but evaluates state-action pairs
- The Q-function for a given policy $\pi$
  - is the cumulative reward for starting in $s$, applying action $a$, and, in the resulting state $s'$, play according to $\pi$

$$Q^{\pi}(s_0, a_0) = r(s_0, a_0) + \sum_{t=1}^{\infty} \gamma^t \mathbb{E}_{s_t} \delta(s_t | s_{t-1}, a_{t-1}) r(s_t, \pi(s_t))$$

$$= r(s_0, a_0) + \frac{1}{k} \sum_{i=0}^{k} \sum_{t=1}^{\infty} \gamma^t r(s_t^i, a_t^i) \mid s_t \sim \delta(s_t | s_{t-1}, \pi(s_{t-1}))$$

- Now we update the policy without the transition function

$$\pi'(s) = \arg\max_a Q^{\pi}(s, a)$$

# Exploration vs. Exploitation

- The current approach requires us to evaluate every action
    - We need to sample each state (that is reachable from $S_0$)
    - We need to compute argmax a over all available actions

- Exhaustive sampling is unrealistic
    - The state/action space may be very large, even infinite (continuous)
    - We approximate an expectation, hence multiple samples for every state/action are required

- We need to decide where to sample the transition function
    - Interesting = visited by the optimal policy
    - But we don't know the optimal policy till the end

# Exploration vs. Exploitation

- Exploit
  - Use the action we assume to be the best
  - Approximate the optimal policy

- Explore
  - Optimal action may be wrong due to approximation errors
  - Try a suboptimal action

- Define probabilities for exploration and exploitation
  - Policy evaluation with stochastic policy

$$Q^\pi(s_0, a_0) = r(s_0, a_0) + \frac{1}{k} \sum_{i=0}^{k} \sum_{t=1}^{\infty} \gamma^t r(s_t^i, a_t^i) \mid s_t^i \sim \mathrm{Pr}^\pi(s_t^i | s_{t-1}^i)$$

  - Well defined tradeoff can reduce samplecounts substantially
  - Most relevant problem for reinforcement learning

# Exploration vs. Exploitation

- ϵ-greedy
  - Fixed probability for selecting a suboptimal action

$$\pi'(a|s) = \begin{cases} 1 - \epsilon + \frac{\epsilon}{|A|} & \textbf{if } a = \arg\max_{a \in A} Q^\pi(s,a) \\ \frac{\epsilon}{|A|} & \textbf{otherwise} \end{cases}$$

- Soft-Max
  - Action probability related to expected value

$$\pi'(a|s) = \frac{e^{Q^\pi(s,a)/t}}{\int e^{Q^\pi(s,a)/t}}$$

- High exploration in the beginning
- Pure exploitation at the end
- Tradeoff must change over time

# Drawbacks

- Policy Iteration with Monte Carlo evaluation works well in practise with small state spaces
    - Don't learn a policy for each state, but learn the policy as a function
    - Especially well suited for continuous state spaces
    - Amount of function parameters usually much smaller than the amount of states
    - Requires well defined function space
    - Direct Policy Search (not part of this lecture)

- Alternative: Bootstrapping
    - Evaluate policy based on estimates
    - May induce errors
    - But requires much lower amount of samples

# Optimal Q-function

- the optimal Q-function is the cumulative reward for starting in $s$, applying action $a$, and, in the resulting state $s'$, play optimally (derivation: deterministic policy)

$$Q^*(s_0, a_0) = r(s_0, a_0) + \sum_{t=1}^{\infty} \gamma^t \mathbb{E}_{s_t} \delta(s_t|s_{t-1}, \pi^*(s_{t-1})) r(s_t, \pi^*(s_t))$$

$$= r(s_0, a_0) + \gamma \mathbb{E}_{s_1} \delta(s_1|s_0, a_0) r(s_1, \pi^*(s_1)) + \gamma^2 \mathbb{E}_{s_2} \delta(s_2|s_1, \pi^*(s_1)) r(s_2, \pi^*(s_2)) + \cdots$$

$$= r(s_0, a_0) + \gamma \left( \mathbb{E}_{s_1} \delta(s_1|s_0, a_0) r(s_1, \pi^*(s_1)) + \gamma \mathbb{E}_{s_2} \delta(s_2|s_1, \pi^*(s_1)) r(s_2, \pi^*(s_2)) + \cdots \right)$$

$$= r(s_0, a_0) + \gamma \mathbb{E}_{s_1} \delta(s_1|s_0, a_0) Q^*(s_1, \pi^*(s_1))$$

$$= r(s_0, a_0) + \gamma \mathbb{E}_{s_1} \delta(s_1|s_0, a_0) \max_{a_1 \in A} Q^*(s_1, a_1)$$

- *Bellman equation*: $Q(s, a) = r(s, a) + \gamma \mathbb{E}_{s'} \delta(s'|s, a) \max_{a' \in A} Q(s', a')$

  - the value of the Q-function for the current state $s$ and an action $a$ is the same as the sum of
    - the reward in the current state $s$ for the chosen action $a$
    - the (discounted) value of the Q-function for the best action that I can play in the successor state $s'$

# Better Approach:
# Directly Learning the Q-function

- Basic strategy:
  - start with some function $\hat{Q}$, and update it after each step
  - in MENACE: $\hat{Q}$ returns for each box $s$ and each action $a$ the number of beads in the box
- update rule:
  - the Bellman equation will in general not hold for Q i.e., the left side and the right side will be different

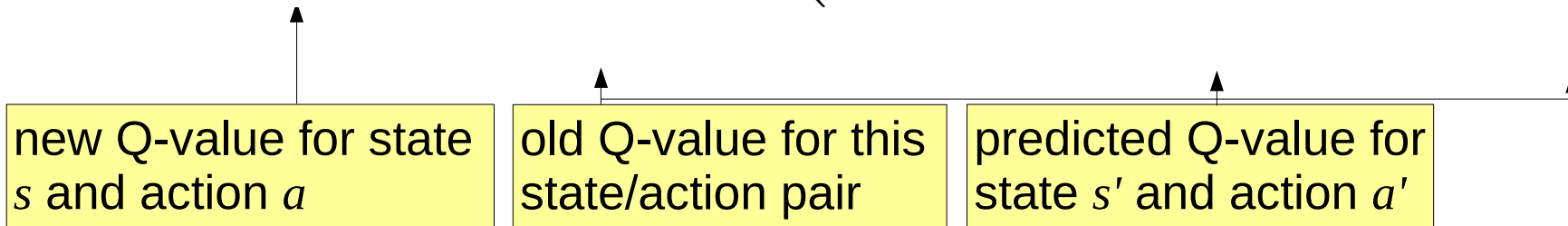$$Q(s,a) = r(s,a) + \gamma \mathbb{E}_{s'} \delta(s'|s,a) \max_{a' \in A} Q(s', a')$$

  - We can not easily compute the expectation
  - But we have multiple samples that contribute to the expectation

# Better Approach:
# Directly Learning the Q-function

- Update Q-Function whenever we observe a transition s,a,r,s'

- Weighted update by a learning rate α

$$\hat{Q}(s,a) \leftarrow (1-\alpha)\hat{Q}(s,a) + \alpha(r(s,a) + \gamma \max_{a' \in A} \hat{Q}(s',a'))$$

$$\leftarrow \hat{Q}(s,a) + \alpha \left( r(s,a) + \gamma \max_{a' \in A} \hat{Q}(s',a') - \hat{Q}(s,a) \right)$$

| new Q-value for state $s$ and action $a$ | old Q-value for this state/action pair | predicted Q-value for state $s'$ and action $a'$ |
|---|---|---|

# Q-learning (Watkins, 1989)

1. initialize all $\hat{Q}(s,a)$ with 0

2. observe current state $s$

3. loop

   1. select an action $a$ and execute it

   2. receive the immediate reward and observe the new state $s'$

   3. update the table entry

   $$\hat{Q}(s,a) \leftarrow \hat{Q}(s,a) + \alpha[(r(s,a) + \gamma \, max_{a'} \hat{Q}(s',a')) - \hat{Q}(s,a)]$$

   4. $s = s'$

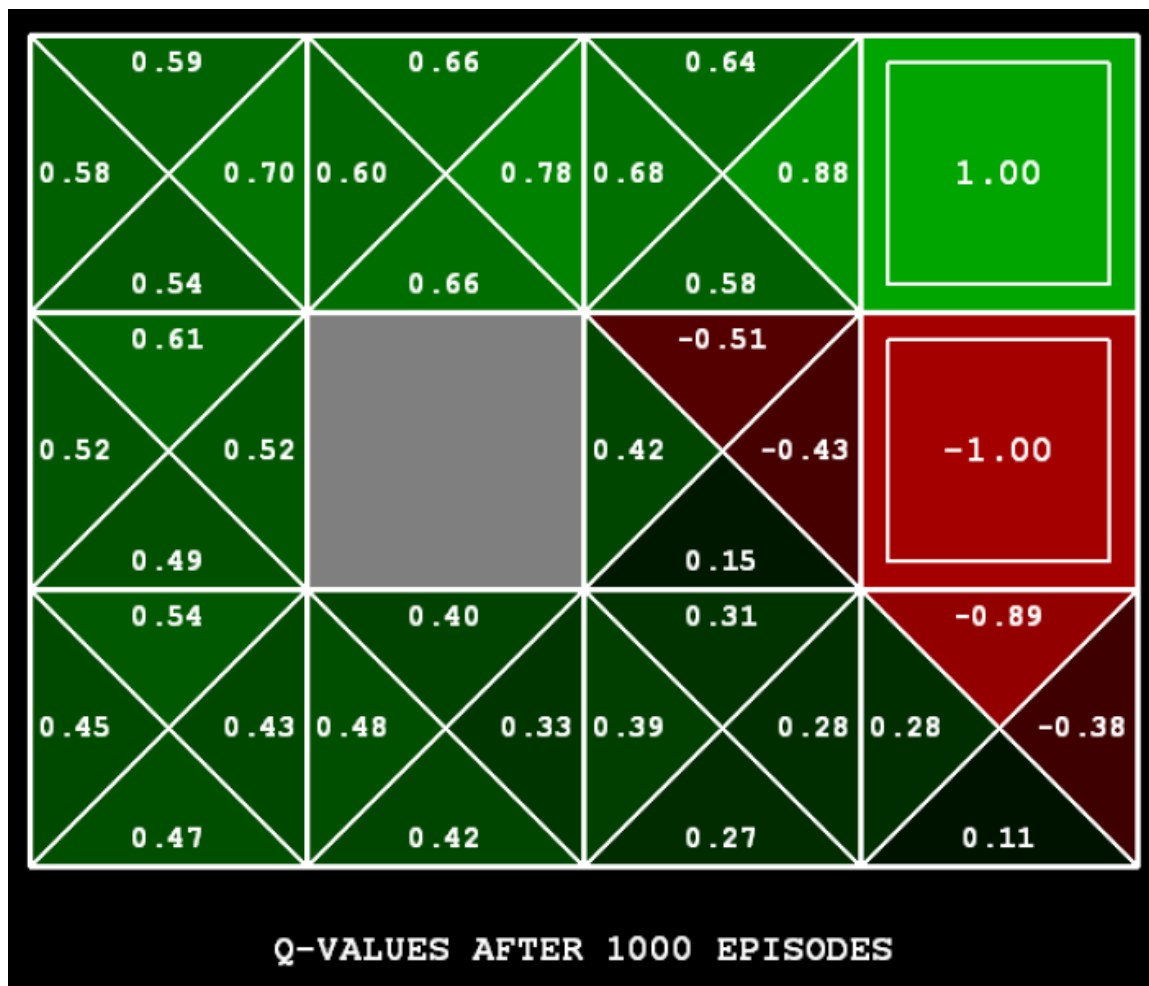**Temporal Difference:**
Difference between the estimate of the value of a state/action pair before and after performing the action.
→ Temporal Difference Learning

# Example: Maze

- Q-Learning will produce the following values

# Miscellaneous

- **Weight Decay:**
    - $\alpha$ decreases over time, e.g. $\alpha = \dfrac{1}{1 + visits(s,a)}$

- **Convergence:**
  it can be shown that Q-learning converges
    - if every state/action pair is visited infinitely often
        - not very realistic for large state/action spaces
        - but it typically converges in practice under less restricting conditions

- **Representation**
    - in the simplest case, $\hat{Q}(s,a)$ is realized with a look-up table with one entry for each state/action pair
    - a better idea would be to have trainable function, so that experience in some part of the space can be generalized
    - special training algorithms for, e.g., neural networks exist

# Drawbacks of Q-Learning

- We still need to compute argmax a, requiring estimates for all actions
    - argmax a is the optimal policy
    - Our policy converges to the optimal policy
    - Don't use argmax a, but the action from the current policy

- perform *on-policy updates*
    - update rule assumes action $a'$ is chosen according to current policy
    - Update whenever observing a sample s,a,r,s',a'

    $$\hat{Q}(s,a) \leftarrow \hat{Q}(s,a) + \alpha \left( r(s,a) + \gamma \hat{Q}(s',a') - \hat{Q}(s,a) \right)$$

    - convergence if the policy gradually moves towards a policy that is greedy with respect to the current Q-function
    - SARSA

# Batch TD Learning

- ■ We try to minimize the <span style="color:blue">Bellman Error</span>

$$Q^\pi(s,a) = \arg\min_Q ||r(s,a) + \gamma \mathbb{E}_{s'} \delta(s'|s,a) Q^\pi(s',a') - Q^\pi(s,a)||, a' \sim \pi(a'|s')$$

- ■ We don't need a weighted update, but can minimize the error globally
  - ▪ Uses multiple samples at once to compute the expectation
  - ▪ Store samples s,a,r,s',a' from current iteration
  - ▪ Minimize error over all obtained samples s,a,r,s',a'

$$Q^\pi(s,a) = \arg\min_Q \sum_{s,a,r,s'} ||r(s,a) + \gamma Q^\pi(s',a') - Q^\pi(s,a)||, a' \sim \pi(a'|s')$$

# Properties of RL Algorithms

- **Transition Function**
  - Model-based: Assumed to be know or approximated
  - Model-free

- **Sampling**
  - On-Policy: Samples must be from the policy we want to evaluate
  - Off-Policy: Samples obtained from any policy

- **Policy Evaluation**
  - Value-based: Computes a state/action value function (this lecture)
  - Direct: Compute expected return for a policy

- **Exploration**
  - Directed: Method guides to a specific trajectory/state/action
  - Undirected: Method allows random sampling close to the expected maximum

# Discussion

- Q-Learning: Model-based/free? On-/Off-Policy?

# Discussion

- Q-Learning: Model-based/free? On-/Off-Policy?
    - Model-free
    - Off-Policy – argmax a, not current policy
    - Q(s',a')-Q(s,a) only depends on the (static) transition function

# Discussion

- Q-Learning: Model-based/free? On-/Off-Policy?
    - Model-free
    - Off-Policy – argmax a, not current policy
    - Q(s',a')-Q(s,a) only depends on the (static) transition function

- SARSA: Model-based/free? On-/Off-Policy?

# Discussion

- Q-Learning: Model-based/free? On-/Off-Policy?
    - Model-free
    - Off-Policy – argmax a, not current policy
    - Q(s',a')-Q(s,a) only depends on the (static) transition function

- SARSA: Model-based/free? On-/Off-Policy?
    - Model-free
    - On-Policy

# Discussion

- Q-Learning: Model-based/free? On-/Off-Policy?
    - Model-free
    - Off-Policy – argmax a, not current policy
    - Q(s',a')-Q(s,a) only depends on the (static) transition function

- SARSA: Model-based/free? On-/Off-Policy?
    - Model-free
    - On-Policy

- Batch TD: Model-based/free? On-/Off-Policy?

# Discussion

- Q-Learning: Model-based/free? On-/Off-Policy?
    - Model-free
    - Off-Policy – argmax a, not current policy
    - Q(s',a')-Q(s,a) only depends on the (static) transition function

- SARSA: Model-based/free? On-/Off-Policy?
    - Model-free
    - On-Policy

- Batch TD: Model-based/free? On-/Off-Policy?
    - Model-free
    - On-Policy

# Discussion

- Q-Learning: Model-based/free? On-/Off-Policy?
    - Model-free
    - Off-Policy – argmax a, not current policy
    - Q(s',a')-Q(s,a) only depends on the (static) transition function

- SARSA: Model-based/free? On-/Off-Policy?
    - Model-free
    - On-Policy

- Batch TD: Model-based/free? On-/Off-Policy?
    - Model-free
    - On-Policy
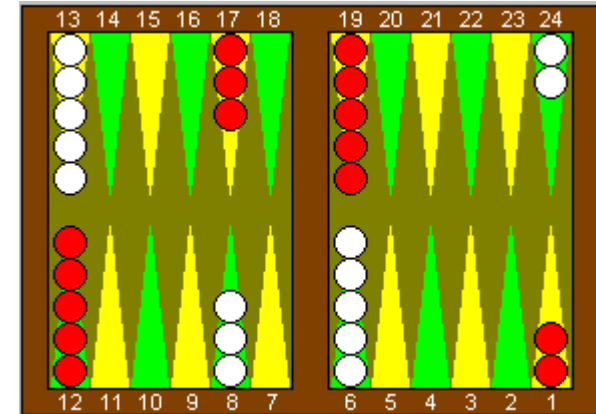
- Batch TD: Off-Policy possible?

# Discussion

- Q-Learning: Model-based/free? On-/Off-Policy?
    - Model-free
    - Off-Policy – argmax a, not current policy
    - Q(s',a')-Q(s,a) only depends on the (static) transition function

- SARSA: Model-based/free? On-/Off-Policy?
    - Model-free
    - On-Policy

- Batch TD: Model-based/free? On-/Off-Policy?
    - Model-free
    - On-Policy

- Batch TD: Off-Policy possible?
    - Store s,a,r,s' from **any** iteration
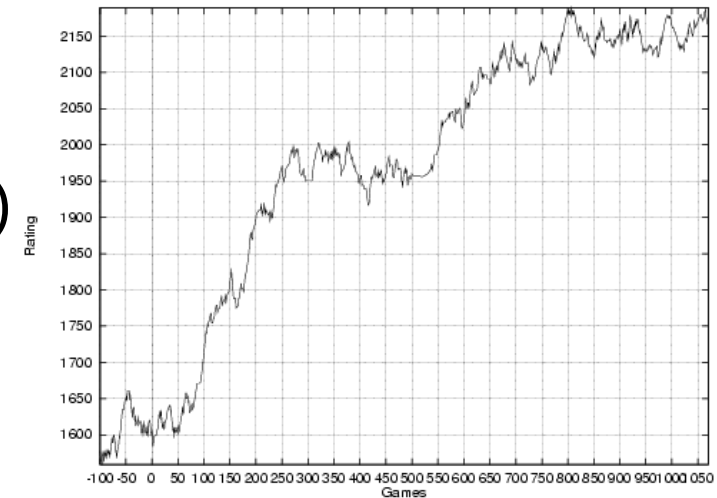    - Compute $a' \sim \pi(a'|s')$ for evaluating the according policy

# TD-Gammon (Tesauro, 1995)

- **weltmeisterliches Backgammon-Programm**

  - Entwicklung von Anfänger zu einem weltmeisterlichen Spieler nach 1,500,000 Trainings-Spiele gegen sich selbst (!)

  - Verlor 1998 WM-Kampf über 100 Spiele knapp mit 8 Punkten

  - Führte zu Veränderungen in der Backgammon-Theorie und ist ein beliebter Trainings- und Analyse-Partner der Spitzenspieler

- **Verbesserungen gegenüber MENACE:**

  - Schnellere Konvergenz durch Temporal-Difference Learning

  - Neurales Netz statt Schachteln und Perlen erlaubt Generalisierung
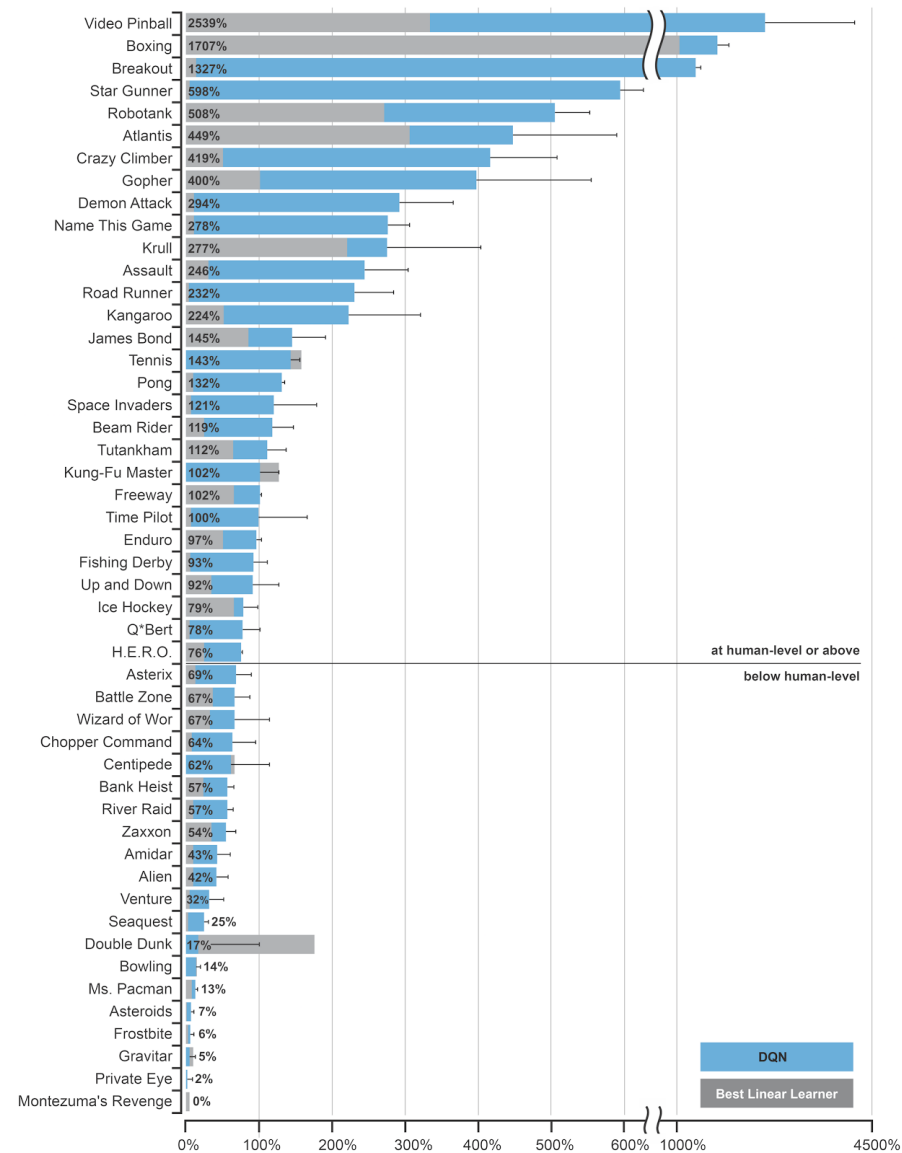
  - Verwendung von Stellungsmerkmalen als Features

# KnightCap (Baxter et al. 2000)

- ## Lernt meisterlich Schach zu spielen
    - ### Verbesserung von 1650 Elo (Anfänger) auf 2150 Elo (guter Club-Spieler) in nur ca. 1000 Spielen am Internet



- ## Verbesserungen gegenüber TD-Gammon:
    - ### Integration von TD-learning mit den tiefen Suchen, die für Schach erforderlich sind
    - ### Training durch Spielen gegen sich selbst → Training durch Spielen am Internet

# Super Human ATARI playing
## (Minh et al. 2013)

- **Reinforcement Learning with Deep Learning**

- **State-of-the-Art**

- **Better than humans in 29/49 ATARI games**

- **Extremely high computation times**

# Reinforcement Learning Resources

- Book
  - On-line Textbook on Reinforcement learning
    - http://www.cs.ualberta.ca/~sutton/book/the-book.html
- More Demos
  - Grid world
    - http://thierry.masson.free.fr/IA/en/qlearning_applet.htm
  - Robot learns to crawl
    - http://www.applied-mathematics.net/qlearning/qlearning.html
- Reinforcement Learning Repository
  - tutorial articles, applications, more demos, etc.
    - http://www-anw.cs.umass.edu/rlr/
- RL-Glue (Open Source RL Programming framework)
    - http://glue.rl-community.org/