# MDP & POMDP

**Agent**

observe

apply action

**Observation**

**action**

**Environment**

# Approximate Q-learning
nice and simple

Q(s,a0),  Q(s,a1),  Q(s,a2)

model
W = params

image

**Q-values:**

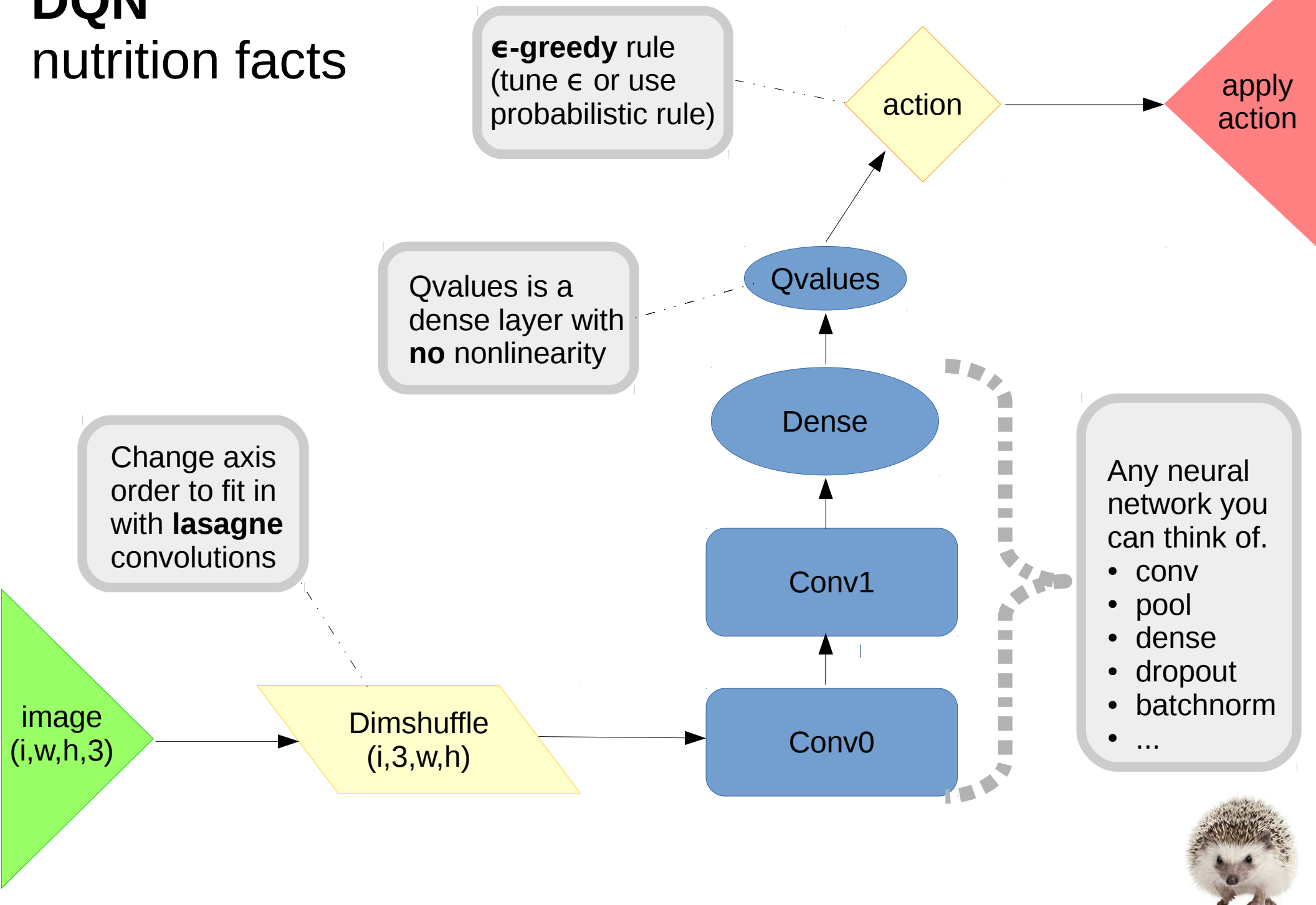$$\hat{Q}(s_t, a_t) = r + \gamma \cdot argmax_a' \hat{Q}(s_{t+1}, a')$$

**Objective:**

$$L = (Q(s_t, a_t) - r + \gamma \cdot argmax_a' Q(s_{t+1}, a'))^2$$

**Gradient step:**

$$w_{t+1} = w_t - \alpha \cdot \frac{\delta L}{\delta w}$$

# DQN
## nutrition facts
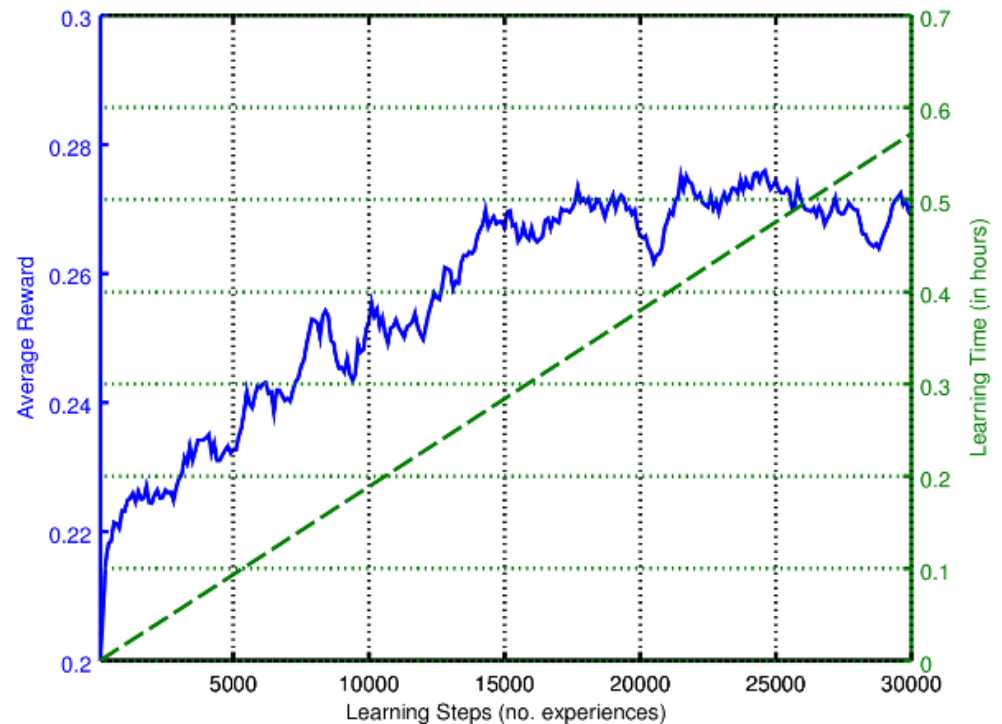
ε-greedy rule
(tune ε or use
probabilistic rule)

action

apply
action

Qvalues is a
dense layer with
**no** nonlinearity

Qvalues

Dense

Change axis
order to fit in
with **lasagne**
convolutions

Conv1

Any neural
network you
can think of.
- conv
- pool
- dense
- dropout
- batchnorm
- ...

image
(i,w,h,3)

Dimshuffle
(i,3,w,h)

Conv0

# Approximate Q-learning
## problems

- Training samples are **not "i.i.d"**,

- Model forgets parts of environment it haven't visited for some time,

- Fallbacks on the learning curve

- **Any ideas?**

# Deep Q-learning
## Multiple agent trick

**Idea:** Throw in several agents with shared **W**.

- Chances are, they will be exploring different parts of the environment,

- More stable training,

- Requires a lot of interaction,

- Alternative to experience replay.