

# Unsupervised Learning in Neural Networks

## LECTURE NOTES

Contributors: Fabian Brix, Renato Kempter

December 28, 2013

## 1 Neuroscience

### 1.1 Visual Cortex

#### 1.1.1 Primary Visual Cortex

The primary visual cortex is the simplest, earliest cortical visual area. It is highly specialized for processing information about static and moving objects and is excellent in pattern recognition.

### 1.2 Receptive Fields

The term receptive field originally described an area of the body surface where a stimulus could elicit a reflex (Sherrington, 1906). The definition was later extended to sensory neurons defining the receptive field as a restricted region of visual space where a luminous stimulus could drive electrical responses in a retinal ganglion cell: *"Responses can be obtained in a given optic nerve fiber only upon illumination of a certain restricted region of the retina, termed the receptive field of the fiber"* (Hartline, 1938).  
[http://www.scholarpedia.org/article/Receptive\\_field](http://www.scholarpedia.org/article/Receptive_field)

**General definition** spanning different types of neurons across sensory modalities: the receptive field is a portion of sensory space that can elicit neuronal responses when stimulated. The sensory space can be defined in a single dimension (e.g. carbon chain length of an odorant), two dimensions (e.g. skin surface) or multiple dimensions (e.g. space, time and tuning properties of a visual receptive field). The neuronal response can be defined as FIRING RATE (i.e. **number of action potentials generated by**

a **neuron**) or include also subthreshold activity (i.e. depolarizations and hyperpolarizations in membrane potential that do not generate action potentials).

### 1.2.1 Visual Receptive Fields

ICA 2, discourse receptive fields.

### Receptive Field Development

## 2 Supervised vs. Unsupervised learning

**Supervised learning** a system is trained for regression or classification via a set of labeled training data. The goal is to optimize the parameters of an objective function using the labeled data (How to correctly predict the labelling "l" or "f" of input patterns according to the training data?). Supervised learning is dealt with in Prof. Seegers class "Pattern Classification and Machine Learning".

**Unsupervised learning** In unsupervised learning one tries to extract patterns from the data without having any information about their classification/outputs. One therefore tries to detect the intrinsic structure of the data. (How does an hand-written "l" or "f" look like?).

**Reinforcement learning** In this context an "agent" is employed and given a certain choice of actions and rewards upon choosing the correct action. The actions exist according to training input patterns (animal conditioning: a mouse sees an "l" or an "f" and chooses action  $a_1/a_2$  accordingly and is given a reward upon success). Reinforcement learning can be either conducted in a supervised or an unsupervised manner.

## 3 Hebbian Learning Rule

change in a given synaptic weight is proportional to both the pre-synaptic input and the output activity of the post-synaptic neuron.

- time-dependent
- local
- strongly interactive

### 3.1 Synaptic plasticity

change in connection strength

Hebb, 1949: "When an *axon* of cell *j* repeatedly or persistently takes part in firing (through *synapses* to *dendrites* of) cell *i*, then *j*'s efficiency as one of the cells firing is increased".

The Hebbian learning rule is a local UNSUPERVISED LEARNING rule, because it depends solely on the states of the respective pre-synaptic *i* and post-synaptic *j* neurons and the present efficacy  $w_{ij}$ , but not on the state of other neurons *k*. The respective pre-synaptic neurons and post-synaptic neuron can however be simultaneously active, so that correlations in their activity lead to firing of the neuronal cell *i*. In Hebbian learning these correlations are used to learn the synaptic transmission efficacies  $w_{ij}$ , so that in essence it is correlation-based learning.

### 3.2 Homeostatic plasticity

Homeostatic plasticity refers to the capacity of neurons to regulate their own excitability relative to network activity, a compensatory adjustment that occurs over the timescale of days. The term homeostatic plasticity derives from two opposing concepts: "homeostatic" (a product of the Greek words for "same" and "state" or "condition") and plasticity (or "change"), thus homeostatic plasticity means "staying the same through change."

### 3.3 Rate-based Hebbian Learning

**Goal** is to find a mathematically formulated learning rule based on Hebb's postulate.

We begin by focusing on a single synapse with efficacy  $w_{ij}$  transmitting signals from a pre-synaptic neuron *i* to a post-synaptic neuron *j*. The activity ((mean) firing rate - average number of spikes per unit time) of *i* is denoted by  $v_i$  and that of *j* by  $v_j$ . Based on the (1) **locality** aspect of Hebbian Learning we can write a general description of the change of the synaptic efficacy.

$$\frac{d}{dt}w_{ij} = \mathbf{F}(w_{ij}; v_j^{pre}, v_i^{post})$$

The membrane potential is further uniquely defined by the postsynaptic firing rate  $v_i = g(u_i)$  with a monotone gain function *g* and therefore doesn't have to be included in *F*.

The (2) **cooperativity** aspect of Hebb's postulates implies that pre- and

postsynaptic neuron have to be active simultaneously for a synaptic weight change to occur and this property can be used to learn something about the function F. Taylor series expansion around  $v_i = v_j = 0$  leads to:

$$\frac{d}{dt}w_{ij} = c_0(w_{ij}) + c_1^{post}(w_{ij})v_i + c_1^{pre}(w_{ij})v_j + c_2^{pre}(w_{ij})v_j^2 + c_2^{post}(w_{ij})v_i^2 + c_2^{corr}v_iv_j + \mathcal{O}(v^3)$$

The term with  $c_2^{corr}$  is bilinear in pre- and postsynaptic activity showing that the change of the synaptic efficacy  $w_{ij}$  is indeed subject to a combination of changes in both activities.

Simplest choice of F is to fix  $c_2^{corr} = c > 0$  and set all other terms of the Taylor expansion to zero resulting in the prototype of Hebbian learning.

$$\frac{d}{dt}w_{ij} = c_2^{corr}v_iv_j$$

This learning rule includes only first-order terms and therefore models non-Hebbian plasticity. More complicated learning rules include higher-order terms of the Taylor expansion. If  $c_2^{corr} = c < 0$ , the learning rule is anti-Hebbian.

F is dependent on the efficacy  $w_{ij}$ , so that it will not grow without limit. A saturation of synaptic weights can be achieved if parameter  $c_2^{corr}$  goes to zero as  $w_{ij}$  approaches its maximum value:

$$c_2^{corr}(w_{ij}) = \gamma_2(1 - w_{ij})$$

Originally Hebb did not consider a rule for decreasing the synaptic weights. In order to have a feasible model however, the synaptic efficacy should decrease in the absence of stimulation.

$$c_0(w_{ij}) = -\gamma_0 w_{ij}$$

The combination results in the simplest feasible learning rule:

$$\frac{d}{dt}w_{ij} = \gamma_2(1 - w_{ij})v_iv_j - \gamma_0 w_{ij}$$

### 3.3.1 Gating

**Postsynaptic gating** A weight change occurs only if the postsynaptic neuron is active,  $v_i > 0$ .  $\Rightarrow$  the weight changes are "gated" by the postsynaptic neuron. Only the efficacies of highly active presynaptic neurons  $v_j > v_\phi$  are strengthened.

$$\frac{d}{dt}w_{ij} = \gamma v_i(v_j - v_\phi(w_{ij})), \gamma > 0$$

**Presynaptic gating** role of pre- and postsynaptic firing rate are exchanged. Now, a change in synaptic weights can only occur if the presynaptic neuron is active,  $v_j > 0$  and the activity of the postsynaptic neuron determines the direction of the change.

$$\frac{d}{dt}w_{ij} = \gamma v_j(v_i - v_\phi)$$

### 3.3.2 BCM rule

The "Bienenstock-Cooper-Munroe" rule is a generalization of the presynaptic gating rule, where  $\Phi$  is a non-linear function and the reference rate  $v_\phi$  is the running average of the postsynaptic activity  $v_i$ .

$$\frac{d}{dt}w_{ij} = \eta \Phi(v_i - v_\phi) v_j - \gamma w_{ij}$$

## 3.4 Learning in Rate Models (functional consequences of hebbian learning)

**Goal:** understand how activity-dependent learning rules influence the formation of connections between neurons in the brain.

Plasticity is controlled by the statistical properties of the presynaptic input for the postsynaptic neuron.

### 3.4.1 Evolution of synaptic weights

For the sake of simplicity we model the presynaptic input as a set of static patterns  $\{x^\mu \in \mathcal{R}^N; 0 \leq \mu \leq p\}$ . At each time step one of the patterns is selected at random and the presynaptic rates are fixed to  $v_i = x_i^\mu$ . The synaptic weights are modified according to a Hebbian learning rule dependent on the correlation of pre- and postsynaptic activity:  $\frac{d}{dt}w_{ij} = a_2^{corr} v_i^{post} v_j^{pre}$ .

In a general rate model the firing rate  $v_i^{post}$  is given by a nonlinear function of the total input:  $v^{post} = g(\sum_i w_i v_i^{pre})$ . For simplification reasons we focus on a linear rate model:

$$v_i^{post} = \sum_k w_{ik} v_k^{pre}$$

We can now combine the learning rule and the linear rate model to form:

$$\frac{d}{dt}w_{ij} = a_2^{corr} \sum_k w_{ik} v_k^{pre} v_j^{pre}$$

We are interested in the long-term behaviour of the synaptic weights and therefore consider the expectation value of the weight vector, i.e., the weight vector averaged over the sequence  $(x^{\mu_1}, x^{\mu_2}, \dots, x^{\mu_n})$  that have so far been presented to the network.

$$\left\langle \frac{d}{dt} w_{ij} \right\rangle = a_2^{corr} \sum_k \langle w_{ik} \rangle \underbrace{\langle x_k^\mu x_j^\mu \rangle}_{C_{kj}} \quad C_{kj} = \langle x_k x_j \rangle = \frac{1}{p} \sum_{\mu=1}^p x_k^\mu x_j^\mu$$

$C_{kj}$  are the correlation matrices for respective components  $i$  and  $j$  in the sum of the learning rule. Using PCA we can detect the direction of maximal variance in the set of patterns.  $\Rightarrow$  Express weight vector  $w$  in eigenvectors of  $C$ . WHAT DO PRINCIPAL COMPONENTS MEAN IN THIS RESPECT?

**PCA** think of neuronal output as projection of weight vector, fill in PCA equations

### 3.5 Oja's rule

The Oja learning rule is a mathematical formalization of the Hebbian learning rule, such that over time the neuron actually learns to compute a principal component of its input stream. The forgetting term added to balance the growth of the weights this time is not only proportional to the value of the weight, but also to the square of the output of the neuron  $v_i$ .

$$\frac{d}{dt} w_{ij} = a_2^{corr} v_j v_i - \underbrace{a_2^{post}}_{a_2^{corr} w_{ij}} (v_i^{post})^2 = a_2^{corr} (v_j v_i - w_{ij} (v_i^{post})^2)$$

Now, the forgetting term balances the growth of the weight. The squared output  $(v_i^{post})^2$  guarantees that the larger the output of the neuron becomes, the stronger is this balancing effect.

#### 3.5.1 Oja's rule and PCA

$$v_i^{post} = \mathbf{w}^T \mathbf{x}^\mu$$

$$\frac{d}{dt} \mathbf{w} = a_2^{corr} (\mathbf{x}^\mu \mathbf{x}^{\mu T} \mathbf{w} - \mathbf{w}^T \mathbf{x}^\mu \mathbf{x}^{\mu T} \mathbf{w} \mathbf{w}), \{x^\mu, 0 \leq \mu \leq p\}$$

This is the incremental change for just one input vector  $x^\mu$ . When the algorithm is run for a long time, changing the input vector at every step, one can look at the average behaviour. An especially interesting question

is what is the value of the weights when the average change in the weight is zero. This is the point of convergence of the algorithm. Averaging right hand side over  $x^\mu$  while  $w$  stays constant yields the following equation.

$$\frac{d}{dt} \mathbf{w} = a_2^{corr} \left( \underbrace{\langle \mathbf{x} \mathbf{x}^T \rangle}_C \mathbf{w} - \underbrace{(\mathbf{w}^T \langle \mathbf{x} \mathbf{x}^T \rangle \mathbf{w})}_{\lambda} \mathbf{w} \right) = 0, \bar{x} = 0$$

Considering that the quadratic form  $\mathbf{w}^T C \mathbf{w}$  is a scalar, this equation clearly is the eigenvalue-eigenvector equation for the covariance matrix  $C$ . This analysis shows that if the weights converge in the Oja learning rule  $C \mathbf{w} = \lambda \mathbf{w}$ , then the weight vector becomes one of the eigenvectors of the input covariance matrix  $\mathbf{w} = \mathbf{e}_1$  (ONLY THIS STATE IS STABLE), and the output of the neuron becomes the corresponding principal component. Principal components are defined as the inner products between the eigenvectors and the input vectors. For this reason, the simple neuron learning by the Oja rule becomes a principal component analyzer (PCA). Although not shown here, it has been proven that the neuron will find the **first principal component** and the norm of the weight vector tends to one. PCA IS USED FOR HEBBIAN LEARNING IN **linear** NEURONS

### 3.6 Independent Component Analysis (ICA)

**Goal:** Compute a transformation  $\mathbf{y} = \mathbf{A} \mathbf{x}$  such that the entries  $y_i$  of the transformed vector (components) are mutually independent. Statistical independence is a stronger constraint than uncorrelatedness required by the PCA and actually implies it. (The two conditions are equivalent *only* if the sources  $x_i$  are generated from Gaussian random variables - one random vector).

**Assumptions:** Random input vector  $\mathbf{x}$  is generated by a linear combination of statistically independent  $p(\mathbf{x}) = p(x_1, \dots, x_n) = p_1(x_1) \dots p_n(x_n)$  and stationary components (sources),  $\mathbf{x} = \mathbf{A} \mathbf{y}$ . The task is now is to find a matrix  $\mathbf{W}$  so as to recover the original signal sources, the components of  $\mathbf{y} = [s_1 \ s_2 \ \dots \ s_n]^T$  by exploiting the information hidden in the samples of the mixed signal patterns  $\mathbf{x}$  (blind source separation).  $\mathbf{A}$  is known as the mixing and  $\mathbf{W}$  as the demixing matrix, respectively.

In contrast to PCA which can always be performed, ICA is meaningful only if the involved random variables are non-Gaussian. Each of the resulting independent components is **uniquely** estimated up to a multiplicative constant.

For this reason the components are considered to be of unit variance. (Multidimensional Gaussians are isometric/rotationally symmetric and therefore the orientation/principal axis of the data can't be recovered). Searching for independent rather than uncorrelated features gives us the means of exploiting a lot more information, hidden in the higher order statistics of the data.

**Ambiguities** permutation of speakers (ICA can't tell which is which) sign of the speakers (ICA can't guarantee outputting the right sign of the signal) for many applications (audio) one doesn't care about these facts

### Preprocessing (whitening)

- Zero mean data
- Whitening (unit variance). Perform PCA on mixed signals and divide each component by the square-root of its eigenvalue  $x_k^n = \frac{1}{\sqrt{\lambda_k}} x_k \Rightarrow C^n = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}$ . After PCA entries of projections  $\mathbf{y}$  of random vectors are uncorrelated and independence can be rewritten to that of projections  $p(y_1, \dots, y_n) = p_1(y_1) \dots p_n(y_n)$

How to find natural axis? We need measure for non-gaussianity. If we look at two-dimensional data, an arbitrary projection will be close to a Gaussian distribution which doesn't tell us anything about our data (neuronal inputs).  $\Rightarrow$  search for direction which is maximally non-Gaussian.

#### 3.6.1 Gaussianity vs. non-gaussianity

This approach in performing ICA is a direct generalization of the PCA technique. The Karhunen-Love transform focuses on the second-order statistics and demands the cross-correlations  $\langle y_i y_j \rangle$  to be zero. In ICA demanding that the components of  $\mathbf{y}$  be statistically independent is equivalent to demanding all the higher order cross-cumulants to be zero (going up to the fourth-order cumulants is sufficient for many applications). careful: Cross-cumulant vs auto-cumulant cumulants vs. moments? These are actually



moments?

$$\mathcal{K}_1(y(i)) = \mathbb{E}[y(i)] = 0, \quad \text{mean}$$

$$\mathcal{K}_2(y(i)y(j)) = \mathbb{E}[y(i)y(j)], \quad \text{(Co)variance}$$

$$\mathcal{K}_3(y(i)y(j)y(k)) = \mathbb{E}[y(i)y(j)y(k)], \quad \text{skewness}$$

$$\begin{aligned} \mathcal{K}_4(y(i)y(j)y(k)y(r)) &= \mathbb{E}[y(i)y(j)y(k)y(r)] - \mathbb{E}[y(i)y(j)] \mathbb{E}[y(k)y(r)] \\ &\quad - \mathbb{E}[y(i)y(k)] \mathbb{E}[y(j)y(r)] - \mathbb{E}[y(i)y(r)] \mathbb{E}[y(j)y(k)] \\ &= \mathbb{E}[y^4] - 3 \mathbb{E}[y^2], \quad \text{kurtosis} \end{aligned}$$

We need a measure for non-gaussianity, which we find in the higher order cumulants, because for Gaussian distributions there is no difference in ICA & PCA. The kurtosis (4-th order cumulant) is such a measure of non-gaussianity. The smaller it is the more it describes a leptokurtic curve with large values around the mean descending fast towards the first standard deviation. The problem can now be reduced to finding a matrix,  $W$ , for that the second-order (covariances) and fourth-order cross-cumulants, kurtosis, of the transformed patterns are zero.

DOES ANYONE UNDERSTAND THE MEANING OF THE FOLLOWING EQUATION?

$$J(w) = [\mathbb{E}(F(y)) - \mathbb{E}(F(v))]^2$$

TODO: MINIMIZE NON-GAUSSIANITY

### 3.6.2 ICA by neuronal rule

batch

$$\Delta \mathbf{w} = \eta \mathbb{E} [\mathbf{x}g(\mathbf{w}^T \mathbf{x})]$$

online

$$\Delta \mathbf{w} = \eta \mathbf{x}g(\mathbf{w}^T \mathbf{x})$$

TODO

### 3.6.3 FastICA

1. Initialize  $\mathbf{w}$
2. Newton step  $\mathbf{w}^{new} = \mathbb{E} [\mathbf{x}(g(\mathbf{w}^T \mathbf{x}) - \mathbf{w}g'(\mathbf{w}^T \mathbf{x}))]$
3. Renormalize  $\mathbf{w} = \frac{\mathbf{w}^{new}}{\|\mathbf{w}^{new}\|}$
4. if algorithm as not converged, go to 2)

### 3.6.4 Temporal ICA

**Goal:** Find unmixing matrix  $W$  such that outputs of time-dependent signals are independent.

**Assumptions:** Signals are time-dependent and the different components may be recorded with delays  $\tau_i$

$$\langle y_i(t)y_k(t - \tau_k) \rangle = \delta_{ik}\lambda(\tau)$$

## 4 Competitive Learning

### 4.1 Clustering

#### 4.1.1 Nearest Neighbour

Initialize prototypes  $p_k$  to represent groups of data points. Assign datapoints to nearest neighbouring prototype  $|\mathbf{p}_k - \mathbf{x}| \leq |\mathbf{p}_i - \mathbf{x}|$  resulting in Voronoi tessellation of input space. For data compression the input  $\mathbf{x}$  is classified and the index  $k$  transmitted. This gives rise to the construction error

$$\mathbb{E}(\mathbf{w}_1, \dots, \mathbf{w}_n) = \sum_k \sum_{\mu \in C_k} (\mathbf{w}_k - \mathbf{x}^\mu)^2$$

But how to choose initialization of prototypes? To minimize the reconstruction error one has discretize the input space adaptively.

#### 4.1.2 K-Means

Standard K-means is based on a random initialization of prototypes.

1. For every data pattern  $x^\mu$  the winning prototype  $\mathbf{p}_k$  is determined with the nearest neighbour rule.
2. The winning prototype is moved closer to the classified  $\mathbf{x}_k^\mu$  by updating it according to  $\Delta \mathbf{p}_k = \eta(\mathbf{p}_k - \mathbf{x}^\mu)$ .

**Showing that the error decreases** Gradient descent on the error surface in the following manner

$$\Delta \mathbf{p}_j = -\eta \frac{\partial \mathbb{E}}{\partial \mathbf{w}_j} = -\eta \frac{\partial}{\partial \mathbf{w}_j} \sum_k \sum_{\mu \in C_k} (\mathbf{w}_k - \mathbf{x}^\mu)^2$$

$$\text{batch update : } \Delta \mathbf{p}_j = -2\eta \sum_{\mu \in C_j} (\mathbf{x}^\mu - \mathbf{w}_j)$$

$$\text{online update : } \Delta \mathbf{p}_j = 2\eta(\mathbf{x}^\mu - \mathbf{w}_j)$$

In the batch update rule distances from the data points in the prototype's class are summed up and the prototype is moved towards the mean with step size  $\eta$ .

After convergence:  $\Delta w_k = 0$

$$\begin{aligned} \sum_{\mu \in C_k} \mathbf{w}_k - \sum_{\mu \in C_k} \mathbf{x}^\mu &= N_k \mathbf{w}_k - \sum_{\mu \in C_k} \mathbf{x}^\mu \\ \Rightarrow \mathbf{w}_k &= \frac{1}{N_k} \sum_{\mu \in C_k} \mathbf{x}^\mu \end{aligned}$$

So with the batch update rule each prototype is at the center of its data cloud after convergence.

The online update has a stochastic component, because one has to choose a data point at random and move into its direction by the step size  $\eta$ . This causes the algorithm to jitter around the mean. To alleviate this problem the trick is to reduce  $\eta$  so that the algorithm eventually freezes.

**Problem of dead units** Dead units are prototypes without data clouds. The error surface allows the algorithm to converge into a state permitting dead units through several local minima. To overcome this problem it is useful to initialize the prototypes to data points in k-means.

## 4.2 Neuronal Implementation of Clustering

(now multiple postsynaptic neurons for the same input?) In the neuronal implementation the data patterns take the form of sensory inputs to neurons in the brain, one postsynaptic neuron for every prototype  $\mathbf{p}_j$ . The activity of these neurons is defined by a non-linear neuron-model  $v_1^{post} = y_1 = g(\sum_k \mathbf{w}_{ik} \mathbf{x}_k^\mu) = g(\mathbf{w}^T \mathbf{x}^\mu)$  which performs a non-linear transform of the scalar product of the prototypes used as weights  $\mathbf{w}_j = \mathbf{w}_k$  which are now synaptic weights and the sensory inputs.

For NORMALIZED WEIGHT VECTORS the clustering expression  $|\mathbf{p}_k - \mathbf{x}| \leq |\mathbf{p}_i - \mathbf{x}|$  can be rewritten to  $\mathbf{w}_k^T \mathbf{x} \geq \mathbf{w}_i^T \mathbf{x}$ . This means that the smallest distance to the prototype  $\mathbf{w}_k$  results in the largest scalar product. We reformulate our clustering aim to finding the postsynaptic neuron with the maximal scalar product ("biggest total stimulation").

**Winner take all circuit** How can it be assured that only one neuron wins, that is to say that the respective input pattern is classified to a single prototype? Simply by introducing a **lateral term** incorporating activities of the other postsynaptic neurons for the other prototypes. The winner is the neuron that receives the most stimulation from outside (strongest outside drive) inhibiting the others which is equivalent to saying input pattern is classified to the prototype with which it forms the largest scalar product  $\mathbf{w}_k \leq \mathbf{x}^\mu$ .

$$y_i(t+1) = g(\text{total drive}) = g\left(\sum_k w(t+1)_{ik} \mathbf{x}_k^\mu + \sum_j B_{ij} y(t)_j\right), \quad B_{ij} < 0$$

$$\Rightarrow y_i(t^{final}) = \begin{cases} 1 & \text{for winner} \\ 0 & \text{for others} \end{cases}$$

How to update weights?

**Hebbian Learning Rule (no inhibition)** For every pattern  $\mathbf{x}^\mu$  find winning prototype  $\mathbf{w}_k = \mathbf{p}_k$  and update it according to the following learning rule:

$$\begin{aligned} \frac{d}{dt} w_{ki} &= \eta(y_k x_i^\mu - \gamma w_{ki} y_k) \\ &= \eta(\delta_{ki} x_i - w_{ki} \delta_{ki}) \\ &= \eta(x_i - w_{ki}) \end{aligned}$$

update for winner k (only neuron active) is same eq. as before (k-means), now derived as a hebbian learning rule:

$$\frac{d}{dt} \mathbf{w}_k = \eta(\mathbf{x}^\mu - \mathbf{w}_k)$$

### Kohonen Learning Rule

1. choose input  $\mathbf{x}^\mu$
2. determine winner k according to current weights for input (scalar product)
3. update weights of winner k

$$\frac{d}{dt} w_{ki} = \eta(y_k x_i^\mu - \gamma w_{ki} y_k)$$

and neighbours  $j$

$$\frac{d}{dt} = \eta \Lambda(j, k)(\mathbf{x}^\mu - \mathbf{w}_j), \quad \text{on-line}$$

Every neighbour  $j$  is inhibited according to inhibition coefficients of  $k$  for  $j$ .

4. decrease size of neighbourhood function  $\Lambda(j, k)$  to ensure convergence?
5. reiterate from 1) until stabilization

Thus, the neuron whose weight vector was closest to the input vector is updated to be even closer. The result is that the winning neuron is more likely to win the competition the next time a similar vector is presented, and less likely to win when a very different input vector is presented. As more and more inputs are presented, each neuron in the layer closest to a group of input vectors soon adjusts its weight vector toward those input vectors. Eventually, if there are enough neurons, every cluster of similar input vectors will have a neuron that outputs 1 when a vector in the cluster is presented, while outputting a 0 at all other times. Thus, the competitive network learns to categorize the input vectors it sees.

## 5 Reinforcement learning

### 5.1 Introduction

**RL vs. supervised learning** Supervised learning, used in statistical pattern recognition and artificial neural networks is learning from examples provided by a knowledgeable external supervisor whereas in RL an agent is instructed to learn from its own experience. It is a "behavioral" learning problem: Through interaction between the learning system and its environment, whereas the system seeks to achieve a specific goal.

**trade-off exploration vs. exploitation** An agent has to exploit what it ALREADY KNOWS in order to obtain a reward, but it also has to EXPLORE in order to make BETTER ACTION SELECTIONS in the future. The agent must try a variety of actions and progressively favour those that appear to be best, because neither schemes can be pursued exclusively without failing at the task. Estimate reward probabilities vs. take action to maximize reward.

**RLs main characteristic** NOT ONE-SHOT DECISION MAKING PROBLEM, harder than supervised learning RL explicitly considers the WHOLE problem of a goal-directed agent INTERACTING with an uncertain environment. All reinforcement learning agents have explicit goals, can sense aspects of their environments, and can choose actions to influence their environments. Moreover, it is usually assumed from the beginning that the agent has to operate despite significant UNCERTAINTY ABOUT THE ENVIRONMENT it faces. When reinforcement learning involves PLANNING, it has to address the *interplay between planning and real-time action selection*, as well as the question of how environmental models are acquired and improved.

An agent's actions are permitted to affect the future state of the environment (e.g., the next chess position, the level of reservoirs of the refinery, the next location of the robot), thereby affecting the options and opportunities available to the agent at later times. Correct choice requires taking into account indirect, delayed consequences of actions, and thus may require foresight or planning. The knowledge the agent brings to the task at the start—either from previous experience with related tasks or built into it by design or evolution—influences what is useful or easy to learn, but interaction with the environment is essential for adjusting behavior to exploit specific features of the task.

#### 5.1.1 Elements of reinforcement learning

there are four main subelements of a reinforcement learning system.

**Policy** commonly denoted by  $\pi$ , it defines the learning agents way of behaving at a given time. Mapping from the perceived states of the environment to the actions that present themselves to the agent in these states. The policy alone is sufficient to determine the behaviour of an agent. In other words, a policy is a rule used by the learning system to decide what to do, given the current state of the environment.

**Reward function** defines the goal in a reinforcement learning problem. Maps each perceived state of the environment to a single number, a reward, which describes the intrinsic desirability of the state. The agent's sole objective is to maximize its total gain from the rewards it receives. The reward function may serve as a basis for altering the agent's policy.

**Value function** While the reward function indicates which actions are immediately beneficial to the agent, a **value function** specifies what is

good for the agent in the long run. The value of a state is the EXPECTED total amount of rewards the agent will accumulate from this state onwards. For example, a state might always yield a low immediate reward but still have a high value because it is regularly followed by other states that yield high rewards.

**Model of the environment** A model mimics the behaviour of the environment. *For example, given a state and action, the model might predict the resultant next state and next reward.* Models are used for PLANNING, by which we mean any way of deciding on a course of action by considering possible future situations before they are actually experienced. The incorporation of models and planning into reinforcement learning systems is a relatively NEW DEVELOPMENT. Early reinforcement learning systems were explicitly trial-and-error learners; what they did was viewed as almost the opposite of planning. Nevertheless, it gradually became clear that reinforcement learning methods are closely related to dynamic programming methods, which do use models, and that they in turn are closely related to state-space planning methods.

## 5.2 Markov Decision Problems (MDP)

Reinforcement learning algorithms model the world as a Markov Decision Problem (MDP)

**MDP**  $(S, A, P_{s \rightarrow s'}^a, \gamma, \mathcal{R})$

$S$ , state space

$A$ , contains possible actions associated with states

$\sum_{s'} P_{s \rightarrow s'}^a = 1$ , outgoing action prob. sum to 1

$\gamma$ , discount factor for future rewards

$\mathcal{R} : S \rightarrow \mathbb{R}$ , reward function maps states to reward values

policy/strategy profile  $\pi(s, a)$ , describes which action  $a$  to take in state  $s$

## 5.3 Reward-based action learning: Q-values

In the case of a **1-step horizon** the agent starts of in an initial state  $s$  from where he can get to several states  $s'_i$  by choosing an associated action  $a_i$ . A reward  $P_{s \rightarrow s'}^a$  is associated with every action.

**Goal:** The agent should find the optimal strategy  $a^*$  resulting in the maximization  $Q(s, a^*) > Q(s, a_j)$  of the expected reward  $Q(s, a) = \sum_{s'} P_{s \rightarrow s'}^a R_{s \rightarrow s'}^a$ .

**Update rule:** The Q values are however not known at the beginning, because the reward probabilities  $P_{s \rightarrow s'}^a$  are not known. Therefore, the agent needs to learn them via iterative updates:

$$\Delta Q(s, a) = \eta [R_{s \rightarrow s'}^a - Q(s, a)] = \eta [r - Q(s, a)]$$

These consist of the difference between received reward and expected reward multiplied by a learning rate  $\eta$ . Over multiple episodes the Q-values converge to the true value of the reward  $r$ , expressed by convergence of the mean update  $\Delta Q(s, a) \rightarrow 0$ .

## 5.4 SARSA and Bellman equation

For **Multi-step horizon** reward-based action learning we reformulate the update rule to take into account the Q-values from the next step multiplied by a discount factor  $\gamma$ :

$$\begin{aligned} \Delta Q(s, a) &= \eta [r - (Q(s, a) - \gamma Q(s', a'))] \\ &= \eta \underbrace{[r + \gamma Q(s', a') - Q(s, a)]}_{\delta_t} = \eta \delta_t \end{aligned}$$

This update rule is called **SARSA** for **State Action Reward State Action** because it uses the state and action of time  $t$  along with the reward, state, and action of time  $t + 1$  to form the so-called **temporal difference (TD) error**  $\delta_t$ . This update rule is called **associative** because the values that it learns are associated with particular states in the environment  $s \in S$ . We can write the same update rule without any  $s$  terms. This update rule would learn the value of actions without associating them with any environment state, thus it is **nonassociative**. *For instance, if a learning agent were playing a slot machine with two handles, it would be learning only action values, i.e. the value of pulling lever one versus pulling lever two.* The associative case of learning state-action values is commonly thought of as the full reinforcement learning problem, although the related state value and action value cases can be learned by the same TD method.

**Bellman equation** The Bellman equation expresses the expected reward  $Q(s, a)$  recursively over an infinite/multi-step horizon:

$$Q(s, a) = \sum_{s'} P_{s \rightarrow s'}^a [R_{s \rightarrow s'}^a + \gamma \sum_{a'} \pi(s', a') Q(s', a')]$$



The  $Q(s', a')$  are weighted differently according to the chosen policy  $\pi(s, a)$  derived from  $Q(s, a)$ .

**SARSA algorithm** The Bellman equation can be used for the SARSA algorithm implementing the SARSA update rule:

1. Being in state  $s$ , chose action  $a$  according to **policy**
2. Observe reward  $r$ , next state  $s'$
3. Choose action  $a'$  in state  $s'$  according to **policy**
4. Update  $\Delta Q(s, a)$
5.  $s' \rightarrow s, a' \rightarrow a$
6. Goto 1)

## 5.5 Strategies for exploration vs. exploitation

There are different strategies for exploitation/**ways to find optimal policy** while correct  $Q$ -values are not known:

### 5.5.1 Choosing a policy

The action choices for all states  $s \in S$  are specified by a policy  $\pi$ . A policy consists of a rule for choosing the next state based on the value or predictions for possible next states (and in our case actions). More specifically, a policy maps state values to actions. Policy choice is very important due to the need to balance exploration and exploitation. This often means taking suboptimal actions in order to learn more about the value of an action or state.

- **Greedy policy:** Consider an agent who takes the highest-valued action  $a^*$ ,  $Q(s, a^*) > Q(s, a)$  at every step. This policy is called a greedy policy, because the agent is only exploiting its current knowledge of state or action values to maximize the expected reward and is not exploring to improve its estimates of other states. *A greedy policy is likely undesirable for learning because there may be a state which is actually good but whose value cannot be discovered because the agent currently has a low value estimate for it, precluding/preventing its selection by a greedy policy.*

- **$\epsilon$ -greedy policy:** take action which looks best with . A common way of balancing the need to explore states with the need to exploit current knowledge in maximizing rewards is to follow an  $\epsilon$ -greedy policy which simply follows a greedy policy with probability  $P = 1 - \epsilon$  and takes a random action with probability  $\epsilon$ . *This method is quite effective for a large variety of tasks.*
- **Softmax strategy:** take action  $a'$  which looks best with probability Another common policy is softmax. You may recognize softmax as the stochastic activation function for Boltzmann machines, often called the Gibbs or Boltzmann distribution. With softmax, the probability of choosing action  $a$  at time  $t$  is

$$P(a') = \frac{\exp(\beta Q(a'))}{\sum_a \exp(\beta Q(a))}$$

where the denominator sums over all the exponentials of all possible action values and  $\tau = \frac{1}{\beta}$  is the temperature coefficient. A high temperature causes all actions to be equiprobable, while a low temperature skews the probability toward a greedy policy. *The temperature coefficient can also be annealed over the course of training, resulting in greater exploration at the start of training and greater exploitation near the end of training.*

- **Optimistic greedy:** start off with  $Q$  values that are too big

**Action Values  $Q$  versus State Values  $V$**  EXPECTED PAYOFFS Almost all reinforcement learning algorithms are based on estimating value functions. These are functions of states  $V$  (or of state-action pairs  $Q$ ) that estimate how good it is for the agent to be in a given state  $V_s$  (or how good it is to perform a given action in a given state  $Q_a$ ). The notion of “how good” here is defined in terms of future rewards that can be expected, or, to be precise, in terms of expected return. Of course the rewards the agent can expect to receive in the future depend on what actions it will take. Accordingly, value functions are defined with respect to particular policies.

Recall that a policy  $\pi$  is a mapping from states  $s \in S$  and actions  $a \in A$  to the probability  $\pi(s, a)$  of taking action  $a$  when in state  $s$ . Informally, the value of a state  $s$  under a policy  $\pi$ , denoted by  $V^\pi(s)$  is the expected return when starting in  $s$  and FOLLOWING  $\pi$  THEREAFTER. The Bellman equation

for state values  $V$  becomes:

$$\begin{aligned}
 V^\pi(s) &= \sum_a \pi(s, a) \overbrace{\sum_{s'} P_{s \rightarrow s'}^a [R_{s \rightarrow s'}^a + \gamma \sum_{a'} \pi(s', a') Q(s', a')]}^{Q(s, a)} \\
 &= \sum_a \pi(s, a) \sum_{s'} P_{s \rightarrow s'}^a [R_{s \rightarrow s'}^a + \gamma V^\pi(s')], \quad \text{in form of itself}
 \end{aligned}$$

The difference between  $Q$  and  $V$  is the following: The value function is the expected return starting from state  $s_t$  and following the policy  $\pi$ . The state-action value function  $Q(s, a)$  is the expected return starting from state  $s_t$ , taking action  $a$  and **only thereafter** following policy  $\pi$ . WE HOWEVER USE THE STATE-ACTION VALUES  $Q$  INSTEAD OF THE MERE STATE VALUES  $V$

## 5.6 On-policy vs off-policy learning

TODO: PART DOESN'T FULLY MATCH SLIDES?

On-policy means SARSA; off-policy means Q-learning. The difference lies in the method used to compute the update for  $Q(s, a)$ . While for SARSA, we use the policy  $\pi$  to select the action  $a'$  for the discounted term ( $\gamma \sum_{a'} \pi(s', a') Q(s', a')$ ), the Q-learning algorithm simply uses the greedy approach to select  $a'$ . Off-policy Q-learning by default uses a greedy update rule and the policy derived from  $Q$  is only employed for selection of the first action  $a$  from the current state:

$$\Delta Q(s, a) = \alpha \left[ r(a) + \gamma \max_{a'} Q(s', a') - Q(s, a) \right]$$

The **Problem** with these algorithms is that learning is slow, that is to say the propagation of Q-values is slow.

## 5.7 Eligibility traces

An eligibility trace is a temporary record of the **occurrence of an event**, such as the visiting of a state or the **taking of an action** (relevant in our case). The trace marks the memory parameters associated with the event as eligible for undergoing learning changes by weighting the learning rule accordingly.

At the moment of the reward update, we want to also update previous action values along trajectories, because otherwise the information diffusion

across several states takes too much time. Every time an action  $a$  is chosen according to a policy, the eligibility traces are evaluated thus:

- If  $a =$  action taken in state  $j$ :  $e(j, a)(t) = \gamma\lambda e(j, a)(t - \Delta t) + 1$
- else:  $e(s, a)(t) = \gamma\lambda e(s, a)(t - \Delta t)$

The  $\gamma$  here is the memory reduction parameter.

Then, for all states  $s$  and actions  $a$  use the respective eligibility trace as an additional coefficient for the SARSA update rule with TD error  $\delta_t = r_t + \gamma Q(s', a') - Q(s, a)$ .

$$\Delta Q(s, a) = \eta \delta_t e(j, a)$$

## 5.8 Continuous states

Q-values for continuous states approximation by a weighted sum of basis functions

$$Q(s, a) = \sum_j w_{aj} * r_j(s) = \sum_j w_{aj} * \phi(s - s_j)$$

$s_j$  center of basis function,  $\Phi$  shape of basis functions - determines how much of the weights are used.

If the time interval between subsequent weight updates approaches zero, the difference between two consecutive states  $s' - s = v(a) * \Delta t$ , where  $v(a)$  is the velocity resulting from the choice of the action  $a$ .  $s' - s$  approximates zero, therefore  $s' \approx s$ . Similarly, the finer the time resolution becomes, the likelier it is that  $a \approx a'$ . Thus, for small time steps, an online gradient descent on the error term  $E_t$  becomes closer to a  $\delta_t$  modulated Hebbian rule. TODO: EXERCISE  $\frac{dQ(s', a')}{dw_{aj}}$

## 5.9 Temporal Difference TD( $\lambda$ ) algorithms

TD algorithms are often used in reinforcement learning to predict a measure of the total amount  $r_{tot.}$  of reward expected over the future, but they can be used to predict other quantities such as the expected reward  $Q$  as well. Eligibility traces are usually implemented in this context by an exponentially-decaying memory trace, with decay parameter  $\lambda$ . This generates a family of TD algorithms **TD**( $\lambda$ ),  $0 \leq \lambda \leq 1$  with **TD**(0) corresponding to updating only the immediately preceding prediction (as described above - schol- arpedia), and **TD**(1) corresponding to equally updating all the preceding

predictions (which is actually Q-learning).

$$\begin{aligned} e(j, a)(t) &= \gamma \lambda e(j, a)(t - \Delta t) + \begin{cases} 1 * r_j & \text{if } a = a' \\ 0 & \text{else} \end{cases} \\ &= \gamma \lambda e(j, a)(t - \Delta t) + r_j \delta_{a, a'} \end{aligned}$$

*Eligibility traces do not have to be exponentially-decaying traces, but these are usually used since they are relatively easy to implement and to understand theoretically.*

### 5.9.1 Neuronal interpretation

**Recap:** Hebbian learning is a theoretical concept exploiting statistical correlations in passive changes of the activations of presynaptic and postsynaptic neurons. It is useful for development, i.e. the wiring of neurons for receptive fields in the brain). Reinforcement learning in contrast relies on conditional changes (reward is received for policy choosing certain actions) to maximize a reward and is therefore useful for learning new behaviours.

**Hebbian Learning of Actions:** example: current states as input and cues as inputs for action neurons. How to model reward on top of action neurons after they have become selective to a certain combination of states and inputs?. We only want to learn the actions that are rewarded! Hebbian Learning fails. Most models of Hebbian Learning/STDP do not take into account neuromodulators, so cannot describe success/shock

**Neuromodulation** Neuromodulation is the physiological process by which a given neuron uses one or more neurotransmitters to regulate diverse populations of neurons. This is in contrast to classical synaptic transmission, in which one presynaptic neuron directly influences a single postsynaptic partner. Neuromodulators secreted by a small group of neurons diffuse through large areas of the nervous system, affecting multiple neurons. Examples of neuromodulators include dopamine, serotonin, acetylcholine, histamine and others.

extend theory should be good for memory formation as well as **action learning**. Therefore global factors have to be included for feedback (neuromodulator)

$$\Delta w_{ij} \propto F(\underbrace{v_i, v_j}_{\text{local}}, \underbrace{\text{SUCCESS}}_{\text{global}})$$

**Timing issues in Reinforcement learning** success signal is delayed (spike time scale: 1-10ms, reward delay: 100-5000ms) → need memory trace (eligibility trace) implement eligibility traces as memory at synapse, same definition as in **TD**( $\lambda$ )

what is winner takes all interaction here? correlation of reward in current state and reward obtained from action

$$\begin{aligned} e(j, a)(t) &= \gamma \lambda e(j, a)(t - \Delta t) + \begin{cases} r_j r_a & \text{if } a = a' \\ 0 & \text{else} \end{cases} \\ &= \gamma \lambda e(j, a)(t - \Delta t) + r_j \delta_{a, a'} \end{aligned}$$

if a is successful the eligibility trace contains its reward

$$e(j, a)(t) = \gamma \lambda e(j, a)(t - \Delta t) + r_j(s_t) r_a(t)$$

**Plasticity Model for Reinforcement learning** spike-based model

$$\Delta w_{ij} \propto \text{success} \cdot (v_i v_j + \dots)$$

in contrast to old spike-timing dependent plasticity

$$\Delta w_{ij} \propto v_i v_j + \dots$$

synaptic update of current action a

$$\begin{aligned} \Delta w_{aj} &= \eta \delta_t r_j \\ \Delta w_{aj} &= \eta \delta_t r_j r_a \quad \text{if } a \text{ taken} \end{aligned}$$

success = reward - exp. reward  $\delta_t$  (dopamine neurons), success is maximized when  $\delta_t$  is minimized → correct prediction

**Molecular mechanism** changes in synaptic connections synaptic terminal (left), dendrite (right): calcium ions pour into the nerve terminal and these release neurotransmitters from small vesicles. Some of these neurotransmitters activate receptors on the target neuron which in turn opens channels letting sodium ions in TODO

TD incorporates eligibility and continuous state space. We do a function approximation using weights and instead of updating  $Q(s, a)$ . TD method with a function approximator like a neural network, we update  $V_t^\pi$  by finding the output gradient with respect to the weights. This step is not captured

in the previous equation  $\delta_t = r_t + \gamma Q(s', a') - Q(s, a)$ . We will discuss these implementation specific details later.

$$Q_w(s, a) = \sum_{i=1}^n w_i^a * r_i(s)$$

$$\Delta w_{aj} = \eta \delta_t e_{aj}$$

### Eligibility and continous TD

1. choose action  $a$  according to policy (greedy)
2. observe reward  $r$  in state  $s'$ , from there choose next action  $a'$
3. calculate TD error  $\delta_t$  in SARSA
4. update eligibility trace of action  $a$  in current state  $e(j, a)$
5. perform weight update  $\Delta w_{ja} = \eta \delta_t e(j, a)(t)$
6. update Q-value approximation for state  $j$   $Q_w(j, a) = \sum_{s' \in S'} w_{s'}^a r_{s'}(j)$   
(rewards we can get through actions  $a$  leading to states  $s'$  in  $s$ ) for next episode.
7.  $a \leftarrow a', s \leftarrow s'$ , return to 2)

eligibility traces implement (short-term) memory The weights grow/connection is reinforced if action  $a$  at state  $s$  is successful in obtaining a reward

**Algorithmic application: backgammon** TODO DOES SOMEBODY HAVE THE BLACKBOARD: REPRESENTATION OF STATE-VALUE  
Lecture 7 (Backgammon)?

### 5.9.2 Learning algorithm for state values $V$

This error will be multiplied by a learning rate  $\alpha$  and then used to update our weights. The general update rule is

$$V_t = V_t + \alpha [r_{t+1} + \gamma V_{t+1} - V_t]$$

## 5.10 Biological Principles of Learning

### 5.10.1 Spatial Learning

Map brain: place neuronal cells - sensitive to spatial location, Environment box: place fields which are then mapped to place cells

Hippocampal place cells, depends on visual cues, works also in the dark

**navigation to a hidden goal (Morris Water maze)** different starting positions, task learning depends on the hippocampus code for position of the animal, learn action (North, East, South, West, etc.) towards goal after many trials the network should have learnt the right actions to take from positions in the water maze (the weights for the right actions should become specific to the appropriate positions.)

### 5.10.2 Unsupervised learning of place cells

from reinforcement learning of the appropriate actions for states to the unsupervised learning of place cells in the hippocampus?

**Goal:** in order to have the code for the position of the animal in the maze we have to learn the place cells. Derive a model that stores views of visited places. Robot in an environment: visual input at each time step:

**Local view:** activation of set of 9000 gabor wavelets  $L(\mathbf{p}, \Phi) = \{F_k\}$  store views of visited places single view cell (VC) stores a local view all local views are stored in an incrementally growing view cells population.

**Extracting direction** views are aligned by current gaze directions  $\Phi_i$  and their differences  $\Delta\Phi = \Phi_i - \Phi_k$

**Extracting position** small difference between local views - spatially close positions  $F_i$  from set of filter responses

Difference:  $\Delta L = |F_i - F(t)|$

Similarity Measure:  $r_i^{VC} = \exp\left(-\frac{(\Delta L)^2}{2\sigma_{VC}^2}\right)$



**Learning place cells** Population of place cells created incrementally during exploration. Activity of a place: place cells are connected to grid cells of view cells?  $\mathbf{w}_{ij}^p$  - weight vector for place cell i

$$r^{pc} = \sum w_{ij}^{pc} r_j^{pre}$$

competitive hebbian learning:

$$\delta w_{ij}^{pc} = \eta |r_i^{pc} - \phi| (r_j^{pre} - w_{ij}^{pc})$$

**Place cell based vs. landmark-based navigation** Place cells: all place cells connected to all possible action cells(nucleus accumbens) (important weights  $w_{ij}$  are learnt). Learnt actions result in correct movement towards goal.

Landmark-based navigation: Visual filters connected to action cells (dorsal striatum), so that agent turns towards landmark

- Unsupervised learning leads to interaction between vision and path integration
- Model of landmark-based navigation
- Model of map-based navigation
- Learning of actions: **reinforcement learning**
- Learning of views and places: **unsupervised learning**

**Biological mechanisms: grid cells** firing grids extend over the whole space. Subpopulation with different grid sizes exist. Grid cells project to the CA1 area. Grid-cell firing fields rotate following the rotation of the cue. Biological locus of the path integration (electrode in dorsomedial entorhinal cortex)

## 5.11 Policy Gradient Methods

Policy gradient methods are a type of reinforcement learning techniques that rely upon optimizing **parametrized policies**  $\pi_\theta$  with respect to the **expected return** (long-term cumulative reward) by gradient descent.

### 5.11.1 Limitations of TD methods (Q-learning & SARSA)

Reinforcement learning is probably the most general framework in which reward-related learning problems of animals, humans or machine can be phrased. However, most of the methods proposed in the reinforcement learning community are not yet applicable to many problems such as robotics, motor control, etc. due to the following reasons. The motivation for policy gradient methods is that they do not suffer from many of the problems that have been marring traditional reinforcement learning approaches:

- **Lack of guarantees of a value function?:** In continuous observation domains (environment) function approximations which are hard to select/tune are necessary for these methods
- **Intractability:** In only partially observable, **non-markovian** settings, TD algorithms are shown to be fundamentally flawed. Because of the uncertain state information the learning systems need to be modeled as partially observable Markov decision problems which often results in excessive computational cost.
- **Convergence:** Even in fully observable, **markovian** settings, TD algorithms, e.g. Q-learning, most traditional reinforcement learning methods using function approximators have no convergence guarantees and there exist even divergence examples.
- **Complexity:** Continuous states and actions in high dimensional spaces cannot be treated by most reinforcement learning approaches using TD methods, because they are difficult to represent using TD methods.

Policy gradient methods differ significantly as they do not suffer from these problems in the same way. For example, uncertainty in the state might degrade the performance of the policy (if no additional state estimator is being used) but the optimization techniques for the policy do not need to be changed. Continuous states and actions can be dealt with in exactly the same way as discrete ones while, in addition, the learning performance is often increased. **Convergence** at least to a local optimum is **guaranteed**.

### 5.11.2 Mathematical Formulation

UNDERSTAND: by associating certain actions with certain stimuli in a stochastic fashion???

**Goal:** In policy optimization the goal is to optimize the expected reward/reward directly

$$J(\theta) = E \left[ \sum_{k=0}^H a_k r_k \right]$$

according to the parameters  $\theta \in \mathbb{R}^K$  of the policy  $\pi_\theta$ , where  $a_k$  denote time-step dependent weighting factors, often set to  $a_k = \gamma y_k$  for discounted reinforcement learning (where  $\gamma$  is in  $[0,1]$ ) or  $a_k = \frac{1}{H}$  for the average reward case.

**Basic approach:** We model the learning system for a policy gradient method in a discrete-time manner and we will denote the current time step by  $k$ . The input to the system is described using a probability distribution  $s_{k+1} \sim (s_{k+1}|s_k, a_k)$  as model where  $s_k \in \mathbb{R}^M$  denotes the current action, and  $s_k, s_{k+1} \in \mathbb{R}^N$  denote the current and next state, respectively. We furthermore assume that actions are generated by a policy  $a_k \sim \pi_\theta(a_k|s_k)$  which is modeled as a probability distribution in order to incorporate exploratory actions. The policy is assumed to be parameterized by  $K$  policy parameters  $\theta \in \mathbb{R}^K$ . The sequence of states and actions forms a **trajectory** denoted by  $\tau = [s_0 : H, a_0 : H]$  where  $H$  denotes the **horizon** which can be infinite. The words trajectory, history, trial, or roll-out are used interchangeably. At each instant of time, the learning system receives a reward denoted by  $r_k = r(s_k, a_k) \in \mathbb{R}$ .

**Optimization:** For real-world applications, any change to the policy parameterization has to be smooth because drastic changes can be hazardous for the actor. Furthermore useful initializations of the policy based on domain knowledge would otherwise vanish after a single update step. Therefore the policy method of choice is steepest descent on the expected return. The policy parametrization is thus updated according to the following gradient update rule

$$\theta_{h+1} = \theta_h + \alpha_h \nabla_{\theta} J|_{\theta=\theta_h}$$

where  $\alpha \in \mathbb{R}^+$  denotes a learning rate and  $h \in \{0, 1, 2, \dots\}$  the current update number. The time step  $k$  and update number  $h$  are two different variables. In actor-critic-based policy gradient methods, the frequency of updates of  $h$  can be nearly as high as of  $k$ . However, in most episodic methods, the policy update  $h$  will be significantly less frequent.

The **main problem** in policy gradient methods is to obtain a good estimator of the policy gradient  $\nabla_{\theta} J|_{\theta=\theta_h}$

### 5.11.3 Likelihood Ratio Method

Assume that trajectories  $\tau$  are generated from the policy distribution of a system  $\tau \sim p_\theta(\tau) = p(\tau|\theta)$  with return  $r(\tau) = \sum_{k=0}^H a_k r_k$  which leads to  $J(\theta) = \mathbb{E}[r(\tau)] = \int_{\mathbf{T}} p_\theta(\tau) r(\tau) d\tau$ . In this case, the policy gradient can be estimated using the likelihood ratio better known as REINFORCE (Williams, 1992) trick, i.e., by using

$$\nabla_\theta p_\theta(\tau) = p_\theta(\tau) \nabla_\theta \log p_\theta(\tau)$$

from standard differential calculus:  $\nabla_\theta \log p_\theta(\tau) = \nabla_\theta p_\theta(\tau) / p_\theta(\tau)$  we obtain

$$\begin{aligned} \nabla_\theta J(\theta) &= \int_{\mathbf{T}} \nabla_\theta p_\theta(\tau) r(\tau) d\tau \\ &= \int_{\mathbf{T}} p_\theta(\tau) \nabla_\theta \log p_\theta(\tau) r(\tau) d\tau \\ &= E[\nabla_\theta \log p_\theta(\tau) r(\tau)]. \end{aligned}$$

As the expectation  $\mathbb{E}[\cdot]$  can be replaced by sample averages, denoted by  $\langle \cdot \rangle$ , only the derivative  $\nabla_\theta \log p_\theta(\tau)$  is needed for the estimator. Importantly, this derivative can be computed without knowledge of the generating distribution  $p_\theta(\tau)$  as

$$p_\theta(\tau) = p(\mathbf{s}_0) \prod_{k=0}^H p(\mathbf{s}_{k+1} | \mathbf{s}_k, \mathbf{a}_k) \pi_\theta(\mathbf{a}_k | \mathbf{s}_k)$$

implies that

$$\nabla_\theta \log p_\theta(\tau) = \sum_{k=0}^H \nabla_\theta \log \pi_\theta(\mathbf{a}_k | \mathbf{s}_k)$$

as only the policy depends on  $\theta$ . By taking the sample average as Monte Carlo (MC) approximation of this expectation by taking  $N$  trial histories we get

$$\nabla_\theta J = \frac{1}{N} \sum_{n=1}^N \nabla_\theta \log p_\theta(\tau_n) r(\tau_n).$$

which is a fast approximation of the policy gradient for the current policy.  
N TRIALS OR N TRIAL HISTORIES OF ONE TRIAL

### 5.11.4 Drawbacks of Policy Gradient

Of course, policy gradients are not the salvation to all problems having significant problems themselves. They are by definition on-policy (note that tricks like importance sampling can slightly alleviate this problem) and need

to forget data very fast in order to avoid the introduction of a bias to the gradient estimator. **Hence, the use of sampled data is not very efficient.** In tabular representations ???, value function methods are guaranteed to converge to a global maximum while policy gradients only converge to a local maximum and there may be many maxima in discrete problems.

## 5.12 Neuronal implementation of Policy Gradient Methods

TODO: BIOLOGY, NEURONS, HEBB, PERCEPTRON → MATHING LAW

# 6 Energy-Based Models (EBM)

A scalar energy is associated to each configuration of the variables of interest in the model (cost function). Learning then corresponds to modifying the energy function so that its shape has desirable properties, for example we would like desirable configurations to have low energy. The probabilistic distribution of for the model is define by an energy function:

$$p(x) = \frac{\exp(-E(x))}{Z}$$

This distribution is the Boltzmann or Gibbs distribution which originates from classical statistical mechanics.

The normalization factor  $Z$  is called the partition function  $Z = \sum_x \exp(-E(x))$ . An energy-based model can be learnt by performing (stochastic) gradient descent on the empirical negative log-likelihood of the training data.

## 6.1 EBMs with hidden units

Often one does not observe the examples  $x$  fully or one wants to introduce some non-observed variables to increase the expressive power of the model. So we consider a hidden part  $y$  additionally to the observed part  $x$ , giving:

$$P(x) = \sum_y P(x, y) = \sum_y \frac{\exp(-E(x, y))}{Z}$$

By introducing the notation of FREE ENERGY

$$\mathcal{F} = -\log \sum_y \exp(-\mathcal{F})$$

we can rewrite the probability distribution in the condensed of the gibbs distribution:

$$P(x) = \frac{\exp(-\mathcal{F}(x))}{Z}, \text{ with } Z = \sum_x \exp(-\mathcal{F}(x))$$

## 6.2 Restricted Boltzmann Machines (RBM)

Boltzmann Machines (BMs) are a particular form of log-linear (Gibbs distribution) Markov Random Fields (MRF), where the energy function is linear in its free parameters. In order to represent complicated distributions (limited parametric setting  $\rightarrow$  non-parametric setting) we consider some of the variables to be hidden (not observed). One can increase the modeling capacity of the Boltzmann Machine (BM) by adding for hidden variables/units. Furthermore a RBM is restricted to a bipartite Graph, because there no lateral connections visible-visible and hidden-hidden are allowed. We introduce the notation of  $v = x$  for visible units and  $h = y$  for hidden units.

The energy function therefore is defined as:

$$E(v, h) = -b^t v - c^t h - h^t W_{u,v} u$$

where  $W_{u,v}$  represents a connection matrix of weights for hidden and visible units and  $b, c$  are the offsets of the visible and hidden layers respectively.

By plugging in the energy function we get the following representation of the free energy formula:

$$F(v) = -b^t v - \sum_i \log \sum_{h_i} \exp(h_i(c_i + W_i v))$$

Because of the specific structure of RBMs, visible and hidden units are conditionally independent of one-another

A particular set of "synaptic" weights is said to constitute a perfect model of the environmental structure if it leads to exactly the same probability distribution of the states of the visible units as when these units are clamped by the environmental input vectors.

### RBM with binary units

Commonly RBMs are studied with binary units  $v_i, h_i \in \{0, 1\} \dots$ . The free energy simplifies to:

$$\mathcal{F}(v) = -b^t v - \sum_i \log(1 + \exp(c_i + W_i v))$$

[deeplearning.net/tutorial/rbm.html](http://deeplearning.net/tutorial/rbm.html)

sigmoid function (sigm):

$$\sigma(t) = \text{sig}(t) = \frac{1}{1 + \exp(-t)} = \frac{1}{2} \cdot (1 + \tanh(\frac{t}{2}))$$