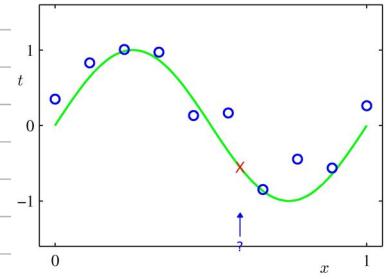


Regression

No Free Lunch Theorem

have to make some assumptions about structure of the data, otherwise generalization is impossible



Model

$$y(x, \vec{w}) = \vec{w}^T \vec{\phi}(x), \quad \vec{\phi}(x) = (\phi_0(x), \dots, \phi_m(x))^T$$

Fitting other approaches are also possible

$$E(\vec{w}) = \frac{1}{N} \sum_n e_n(\vec{w}) = \frac{1}{N} \sum_{n=1}^N \frac{1}{2} (y(x_n, \vec{w}) - t_n)^2 \quad (\text{MSE})$$

$$\frac{dE}{d\vec{w}} = \left[\frac{\partial E}{\partial w_0}, \dots, \frac{\partial E}{\partial w_m} \right]$$

$$= \frac{1}{N} \sum_{n=1}^N (\vec{w}^T \vec{\phi}(x_n) \vec{\phi}(x_n)^T - t_n \vec{\phi}(x_n)^T)$$

$$= \frac{1}{N} \left(\vec{w}^T \sum_{n=1}^N \vec{\phi}(x_n) \vec{\phi}(x_n)^T - \sum_{n=1}^N t_n \vec{\phi}(x_n)^T \right)$$

$$= \frac{1}{N} (\vec{w}^T \underline{\Phi}^T \underline{\Phi} - \vec{t}^T \underline{\Phi})$$

$$= \vec{O}_m^T$$

$$\Rightarrow \vec{w}^T \underline{\Phi}^T \underline{\Phi} = \vec{t}^T \underline{\Phi}$$

nonlinear in the inputs, but linear in parameters $\underline{\Phi}$

$$\vec{w}^T = \vec{t}^T \underline{\Phi} (\underline{\Phi}^T \underline{\Phi})^{-1}$$

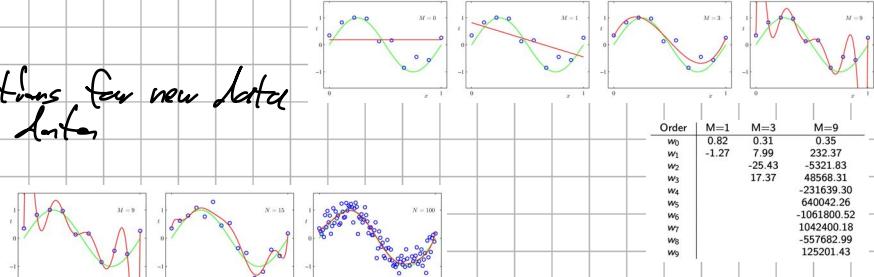
$$\Rightarrow \vec{w}_{\min} = (\vec{t}^T \underline{\Phi} (\underline{\Phi}^T \underline{\Phi})^{-1})^T = ((\underline{\Phi}^T \underline{\Phi})^{-1})^T \underline{\Phi}^T (\vec{t}^T)^T = ((\underline{\Phi}^T \underline{\Phi})^T)^{-1} \underline{\Phi}^T \vec{t} = (\underline{\Phi}^T \underline{\Phi})^{-1} \underline{\Phi}^T \vec{t}$$

Overfitting

overfitted polynomial yields bad predictions for new data tuned to random noise present in data has very large coefficients

ways to reduce:

- more training data (easiest way)
- artificially create e.g. injecting random noise to samples, domain-specific transformations
- limiting number of parameters
- constrain parameters by penalizing their absolute value
 - in case of linear regression e.g. $\hat{E}(\vec{w}) = E(\vec{w}) + \frac{\lambda}{2} \|\vec{w}\|_2^2$ (L2 regularization)
 - closed form solution (Ridge regression): $\vec{w}_{\text{ridge}} = (\lambda I + \underline{\Phi}^T \underline{\Phi})^{-1} \underline{\Phi}^T \vec{t}$



Underfitting

model does not reflect structure of training data

Regularization

changing training criterion of a system to improve generalization

Features

choice of features $\phi_j(x)$ often depends on knowledge of the task
preprocessing data in a way to make relevant information accessible and remove irrelevant content
reduce dimensionality of data

Curse of Dimensionality

number of model parameters rises exponentially with input dimension

example: polynomial fitting we have to take into account all possible products between variables, up to desired order:

$$y(\vec{x}, \vec{w}) = w_0 + \sum_i w_i^{(i)} x^{(i)} + \sum_{\substack{i_1, i_2 \\ i_1 < i_2}} w_i^{(i_1, i_2)} x^{(i_1)} x^{(i_2)} + \dots + \sum_{\substack{i_1, \dots, i_n \\ i_1 < i_2 < \dots < i_n}} w_i^{(i_1, \dots, i_n)} x^{(i_1)} \dots x^{(i_n)}$$

we need to estimate many coefficients, which may be inefficient and can cause severe overfitting

→ number of model parameters rises exponentially with input dimension

ways to tackle the problem:

- real data often lies in a subspace of lower dimensionality
- use dimensionality reduction techniques on data
- define suitable set of features: $\Phi_L(x_1, \dots, x_L)$ where number of features can be much smaller than L^M
- neural networks can intrinsically deal with high-dimensional input and often do not require sophisticated features

Intuitively: #Coefficients in Polynomial (= #Monomials of L variables and degree $\leq M$)

stars and bars argument:

first L bins correspond to the powers of x_1, \dots, x_L and ignore last bin

#stars in first L bins is $\leq M$

→ this produces all possible monomials of L variables and degree $\leq M$

$$\binom{L+M}{M}$$

Classification

discriminant function: $C = f(\vec{\phi})$
 probabilistic modeling: $p(C_k | \vec{\phi}) = f(\vec{\phi}, C_k)$

easy

kNN

non parametric (works directly on data, as opposed to learning parameters of a model)
 majority vote determines class assignment
 a possible extension: weigh influence of kNN (e.g. by distance)

Problems:

- requires to store all training data and compute distances between test sample and all training samples
- does not always "generalize"
- sensitive to outliers
- poor learning for high dimensional feature spaces (Curse of Dimensionality: samples do not cover space well)
- difficult to determine suitable number of neighbors

Linear Classifier

parametric (structure of data is summarized by a set of parameters, independent of number of training samples)

discriminant function:

$$y(\vec{\phi}) = \vec{w}^T \vec{\phi} + w_0 \quad \text{with} \quad C(\vec{\phi}) = \begin{cases} C_0 & \text{if } y(\vec{\phi}) \geq 0 \\ C_1 & \text{if } y(\vec{\phi}) < 0 \end{cases}$$

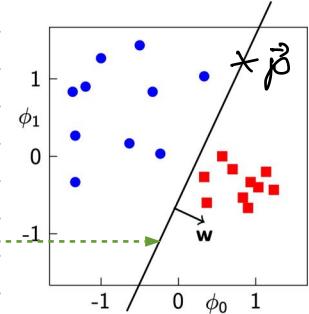
\vec{w} is called the weight vector, w_0 the bias

decision boundary: hyperplane with equation $\vec{w}^T \vec{\phi} + w_0 = 0$

\vec{w} is orthogonal to the decision boundary

distance at any point $\vec{\phi}$ to decision boundary:

$$(\vec{\phi} - \vec{p}) \cdot \frac{\vec{w}}{\|\vec{w}\|} = \frac{\vec{w}^T \vec{\phi} - \vec{w}^T \vec{p}}{\|\vec{w}\|} = \frac{\vec{w}^T \vec{\phi} - \vec{w}^T \vec{p} - w_0 + w_0}{\|\vec{w}\|} = \frac{\vec{w}^T \vec{\phi} + w_0}{\|\vec{w}\|} = \frac{y(\vec{\phi})}{\|\vec{w}\|}$$



just like for regression, we can instead define suitable nonlinear functions

prone to drastic underfitting

multiple classes:

use K linear functions $y_k(\vec{\phi}) = \vec{w}_k^T \vec{\phi} + w_{k0}$ and assign $\vec{\phi}(x)$ to class C_k if $y_k(\vec{\phi}(x)) > y_j(\vec{\phi}(x)) \forall j \neq k$

Logistic Regression

$$\tilde{y}(\vec{\phi}, \vec{w}) = \vec{w}^T \vec{\phi} + w_0$$

$$y(\vec{\phi}, \vec{w}) = \sigma(\tilde{y}(\vec{\phi}, \vec{w})) = \frac{1}{1 + e^{-\vec{w}^T \vec{\phi} + w_0}}$$

$$p(C_1 | \vec{\phi}) = y(\vec{\phi}), \quad p(C_0 | \vec{\phi}) = 1 - y(\vec{\phi})$$

$$E(\vec{w}) = -\log L(\vec{w}) = -\log \left(\prod_{n=1}^N y_n^{t_n} (1-y_n)^{1-t_n} \right) = -\sum_{n=1}^N (t_n \log y_n + (1-t_n) \log (1-y_n)) \rightarrow \text{minimize}$$

no closed form solution

if C_0 and C_1 are linearly separable \rightarrow overfitting

Foundations

assumption: data is described by a prob. distk. $P(\vec{x}, \vec{t})$, $\vec{x} \in \mathbb{R}^m$, $\vec{t} \in \mathbb{R}^n$

→ unknown

let $\vec{y}(\vec{x})$ be a predictor for \vec{t} , $L(\vec{t}, \vec{y}(\vec{x}))$ any loss

True Risk

$$R(\vec{y}) := \mathbb{E}_{P(\vec{x}, \vec{t})} L(\vec{t}, \vec{y}(\vec{x})) = \int L(\vec{t}, \vec{y}(\vec{x})) dP(\vec{x}, \vec{t})$$

goal of ML is to find predictor that minimizes the risk!

$$\vec{y}^* = \underset{\vec{y}: \mathbb{R}^m \rightarrow \mathbb{R}^n}{\operatorname{arg\,min}} R(\vec{y})$$

Irreducible Error

stems from fact that input does not completely determine the target

Empirical Risk

$$R(\vec{y}) \approx \frac{1}{N} \sum_{n=1}^N L(t_n, \vec{y}(\vec{x}_n)) =: R_{\text{emp}}(\vec{y}) \quad \text{with dataset } D = \{(\vec{x}_n, t_n)\}_{n=1\dots N}$$

for pathological functions \vec{y} , the empirical risk does not reflect true risk

↳ restrict predictors we take into account to a subset of functions from \mathbb{R}^m to \mathbb{R}^n :

Hypothesis Class H

we hope this helps us to obtain a predictor whose R_{emp} (on D_{train}) reflects R_{true} (on D_{test}):

$$\vec{y}_H^* = \underset{\vec{y} \in H}{\operatorname{arg\,min}} R(\vec{y}), \quad \vec{y}_{\text{fit}} = \underset{\vec{y} \in H}{\operatorname{arg\,min}} R_{\text{emp}}(\vec{y})$$

in almost all cases, exactly computing the optimal predictor in H is intractable
→ needs to be approximated

all ML algorithms aim at optimally defining the hypothesis class H , such that efficient optimization over a wide range of functions is possible, while retaining generalization ability

Difference between R and R_{emp}

many bounds have been proposed, typically of the form

$$\forall \vec{y} \in H \quad R(\vec{y}) - R_{\text{emp}}(\vec{y}) < \mathcal{O}\left(\sqrt{\frac{\text{cap}(H)}{N}}\right)$$

Kernel Methods

Linear methods in a feature space
 → yet they offer a different view of both linear models, and the concept of features

Consider Linear Regression with L_2 regularization:

$$E(\vec{w}) = \sum_{n=1}^N \frac{1}{2} (\vec{w}^T \vec{\phi}_n - t_n)^2 + \frac{\lambda}{2} \|\vec{w}\|_2^2 \text{ where } \vec{\phi}_n = \vec{\Phi}(\vec{x}_n)$$

Gradient w.r.t. \vec{w} is $\frac{\partial}{\partial \vec{w}} E = \sum_{n=1}^N (\vec{w}^T \vec{\phi}_n - t_n) \vec{\phi}_n + \lambda \vec{w}$ and setting it to 0 yields

$$\begin{aligned} \vec{w} &= \sum_{n=1}^N -\frac{1}{\lambda} (\vec{w}^T \vec{\phi}_n - t_n) \vec{\phi}_n \\ &= \vec{\Phi}^T \vec{\alpha} \end{aligned}$$

$$\vec{\Phi} = \begin{bmatrix} \vec{\phi}_1^T \\ \vdots \\ \vec{\phi}_N^T \end{bmatrix} = \begin{bmatrix} \Phi_0(\vec{x}_1) & \dots & \Phi_N(\vec{x}_1) \\ \vdots & \ddots & \vdots \\ \Phi_0(\vec{x}_N) & \dots & \Phi_N(\vec{x}_N) \end{bmatrix}$$

We see that not only the model is linear, but also that \vec{w} is a linear combination of the (features of the) training data!

Dual Representation

The observation that \vec{w} and features $\vec{\Phi}$ have such a simple relationship leads us to the idea of expressing the whole regression problem in terms of the feature vectors
 → such a reformulation, called **dual representation**, is possible for a variety of linear models

substituting $\vec{w} = \vec{\Phi}^T \vec{\alpha}$ into $E(\vec{w})$ yields:

$$\begin{aligned} E(\vec{\alpha}) &= \frac{1}{2} \sum_{n=1}^N (\vec{\phi}_n^T \vec{w} - t_n)^2 + \frac{\lambda}{2} \|\vec{w}\|_2^2 \\ &= \frac{1}{2} \begin{bmatrix} \vec{\phi}_1^T \vec{w} - t_1 \\ \vdots \\ \vec{\phi}_N^T \vec{w} - t_N \end{bmatrix}^T \begin{bmatrix} \vec{\phi}_1^T \vec{w} - t_1 \\ \vdots \\ \vec{\phi}_N^T \vec{w} - t_N \end{bmatrix} + \frac{\lambda}{2} \vec{w}^T \vec{w} \\ &= \frac{1}{2} (\vec{\Phi} \vec{w} - \vec{t})^T (\vec{\Phi} \vec{w} - \vec{t}) + \frac{\lambda}{2} \vec{w}^T \vec{w} \\ &= \frac{1}{2} (\vec{w}^T \vec{\Phi}^T - \vec{t}^T) (\vec{\Phi} \vec{w} - \vec{t}) + \frac{\lambda}{2} \vec{w}^T \vec{w} \\ &= \frac{1}{2} (\vec{w}^T \vec{\Phi}^T \vec{\Phi} \vec{w} - \vec{w}^T \vec{\Phi}^T \vec{t} - \vec{t}^T \vec{\Phi} \vec{w} + \vec{t}^T \vec{t}) + \frac{\lambda}{2} \vec{w}^T \vec{w} \\ &= \frac{1}{2} (\vec{\alpha}^T \vec{\Phi}^T \vec{\Phi} \vec{\alpha} - \vec{\alpha}^T \vec{\Phi}^T \vec{t} - \vec{t}^T \vec{\Phi} \vec{\alpha} + \vec{t}^T \vec{t}) + \frac{\lambda}{2} \vec{\alpha}^T \vec{\Phi} \vec{\Phi}^T \vec{\alpha} \\ &= \frac{1}{2} \vec{\alpha}^T \vec{\Phi} \vec{\Phi}^T \vec{\Phi} \vec{\alpha} - \vec{\alpha}^T \vec{\Phi} \vec{\Phi}^T \vec{t} + \frac{1}{2} \vec{t}^T \vec{t} + \frac{\lambda}{2} \vec{\alpha}^T \vec{\Phi} \vec{\Phi}^T \vec{\alpha} \end{aligned}$$

Kernel function

$$k(\vec{x}_n, \vec{x}_m) := \vec{\Phi}(\vec{x}_n)^T \vec{\Phi}(\vec{x}_m) = \langle \vec{\Phi}(\vec{x}_n), \vec{\Phi}(\vec{x}_m) \rangle$$

Gram matrix

$$K := \vec{\Phi} \vec{\Phi}^T = (k(\vec{x}_n, \vec{x}_m))_{n,m=1,\dots,N} = \begin{bmatrix} -\vec{\phi}_1^T \\ \vdots \\ -\vec{\phi}_N^T \end{bmatrix} \begin{bmatrix} 1 & \dots & 1 \\ \vec{\phi}_1^T & \cdots & \vec{\phi}_N^T \\ 1 & \dots & 1 \end{bmatrix} \quad \text{so}$$

→ allows us to express the error in the dual representation:

$$E(\vec{\alpha}) = \frac{1}{2} \vec{\alpha}^T K K \vec{\alpha} - \vec{\alpha}^T \vec{t} + \frac{1}{2} \vec{t}^T \vec{t} + \frac{\lambda}{2} \vec{\alpha}^T K K \vec{\alpha}$$

Solving in Dual Representation

setting the gradient of $E(\vec{\alpha})$ w.r.t. $\vec{\alpha}$ to zero, one obtains the following minimum at the error function:

$$\vec{\alpha} = (\mathbf{K} + \lambda \mathbf{I}_N)^{-1} \vec{t} \quad \rightarrow \text{inverting an } N \times N \text{ matrix}$$

plugging this solution into the original model, the prediction function is given by

$$y(\vec{x}) = \vec{k}(\vec{x})^T (\mathbf{K} + \lambda \mathbf{I}_N)^{-1} \vec{t},$$

with $\vec{k}(\vec{x}) = (k(\vec{x}_n, \vec{x}))_{n=1,\dots,N}$,

The Role of the Kernel

the solution is found by inverting a $N \times N$ matrix, whereas in the original formulation, we have inverted an $M \times M$ matrix

advantage of dual formulation: we never use the feature vector ϕ directly

input data appears only inside the kernel function $k(\vec{x}, \vec{y})$

→ allows us to **implicitly** use very high-dimensional (possibly infinite-dimensional) feature spaces!

Example: Quadratic Features

$$\phi(\vec{x})^T = [1, \sqrt{2}x_1, \dots, \sqrt{2}x_k, x_1^2, \dots, x_k^2, \sqrt{2}x_1x_2, \dots, \sqrt{2}x_1x_k, \sqrt{2}x_2x_3, \dots, \sqrt{2}x_2x_k, \dots, \sqrt{2}x_{k-1}x_k]$$

$$k(\vec{x}, \vec{y}) = \vec{\phi}(\vec{x})^T \vec{\phi}(\vec{y})$$

$$= 1 + 2x_1y_1 + \dots + 2x_ky_k + x_1^2y_1^2 + \dots + x_k^2y_k^2 + 2x_1x_2y_1y_2 + \dots + 2x_{k-1}x_ky_{k-1}y_k$$

$$= 1 + 2 \sum_{i=1}^k x_i y_i + \sum_{i=1}^k x_i^2 y_i^2 + 2 \sum_{1 \leq i < j \leq k} x_i x_j y_i y_j$$

$$= 1 + 2 \vec{x}^T \vec{y} + \sum_{i=1}^k x_i y_i + 2 \sum_{1 \leq i < j \leq k} x_i y_i x_j y_j$$

$$= 1 + 2 \vec{x}^T \vec{y} + \left(\sum_{i=1}^k x_i y_i \right)^2$$

$$= 1 + 2 \vec{x}^T \vec{y} + \vec{x}^T \vec{y}^2$$

$$= (1 + \vec{x}^T \vec{y})^2$$

Reinforcement Learning

agent takes actions in some environment and collects scalar rewards
goal of agent: maximize return (cumulative reward)

differs from SL:

- learn actions instead of predictions
- no fixed dataset, data is generated by interacting with environment
- no explicit error correction by terms
- rewards are asynchronous, not clear which action led to which reward
- inherently sequential, focus on online capabilities

often only partial information available (full state of environment is unknown)

Definitions:

S_t state of environment at time t

A_t action taken at time t , after seeing S_t

R_t reward after taking A_t

$$G_t = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots = \sum_{k=0}^H \gamma^k R_{t+k+1}$$

cumulative future reward (return)

γ discount factor, $0 < \gamma \leq 1$

$\gamma = 1$: agent is **farsighted** (all rewards matter equally)
 $\gamma < 1$: agent is **shortsighted**

H horizon, may be finite or infinite we assume $H=\infty$, finite horizon is approximated by setting $\gamma < 1$

π policy tells agent for any state what action to take

stochastic: $\pi(s, a) \rightarrow \pi(s|a)$
deterministic: $a = \pi(s)$

we can reformulate objective of RL: Learn policy which maximizes (expected*) return!

* environment often stochastic (i.e. same action may lead to different successor states and rewards)
policy may also be stochastic
→ return may also be stochastic and we are interested in expected value

expected value
→ in practice we run agent in some way, collect samples and then average over samples

choosing an action may lead to (pos or neg) reward only much later in the future

→ credit assignment problem: how much credit should each action in the sequence of actions get for the outcome?

→ no explicit error correction: reward validates actions, but does not give instructions

exkurs: how RL can be used to simulate SL / how RL is a superset of SL

sequence could contain just 1 step (episodes of 1 timestep)

state: image, action: pick class

whether this is a good idea is a different question

model

- set of states
- ways in which actions lead to state transitions
- rewards which are obtained in each state, or state-action pairs

stochastic
deterministic

→ ex robot: state only partially observable,
does not know what is located

model-based reinforcement learning

- only possible if above information is perfectly known and set of states and actions is finite
- create full map of optimal actions in each state (i.e. set of states (and actions) finite and observable)
- plan actions which lead to high return
- table of state and actions and then plan how to act best

s	a
1	1
2	2

in chess this would be finite but the number of states is already very large (would take too long to print in unison)

however, often:

- number of states very large (e.g. chess) or unlimited
- state cannot be completely observed (e.g. robot)
- even if we know states, but we do not know about expected rewards, i.e. we don't understand the reward function

using full maps from states to actions → model-free RL

model-free reinforcement learning

- focus directly on learning good action from partial observation of the state
- generalize to unseen states / transitions of states

Model-based

finite number of states $\{s_1, \dots, s_N\}$ and we can exactly observe the current state s_t at time step t : $S_t = s_n$

finite set of actions $\{a_1, \dots, a_m\}$

→ each action causes a change of state, which may be probabilistic or deterministic

we know the distribution of rewards in a particular state, or for a particular state-action pair
→ also rewards may be probabilistic

in particular, the whole model has a finite description

formalized as a **Markov Decision Process (MDP)**

→ state transitions have the **Markov property**: $P(S_{t+1} = s' | S_t, S_{t-1}, \dots, S_1) = P(S_{t+1} = s' | S_t)$

→ also reward only depends on current state and current action

formally, an MDP is a tuple $(S, A, P_{ss'}, R_{ss'})$, where

$$S = \{s_1, \dots, s_N\}$$

$$A = \{a_1, \dots, a_m\}$$

$P_{ss'}^a = P(S_{t+1} = s' | S_t = s, A_t = a)$ are the state transition probabilities

$R_{ss'}^a = E[R_{t+1} | S_{t+1} = s', S_t = s, A_t = a]$ is the reward function

we also include the **discount factor** γ in this definition

value function

$$V^\pi(s) = E[G_t | S_t = s] \quad \text{where } G_t = \sum_{k=0}^H \gamma^k R_{t+k+1}$$

$$\text{for } H=\infty: G_t = R_{t+1} + \gamma G_{t+1}$$

reward the agent can expect in future if it continues with its policy from state s

→ indicates how "good" it is to be in a given state, given a policy π

→ we assume a stationary environment here

depends on state and policy

→ in chess a state depends on what you are going to do

action-value function

$$Q^\pi(s, a) = E[G_t | S_t = s, A_t = a]$$

expected return if agent chooses action a in state s , and then continues to follow policy π

for a given policy, we can just average over all possible actions that might be taken to calculate $V^\pi(s)$ from $Q^\pi(s, a)$

$$V^\pi(s) = \sum_{a \in A} \pi(a|s) Q^\pi(s, a)$$

$$V^\pi(s) = \sum_{a=1}^m \pi(a|s) Q^\pi(s, a) ?$$

Otherwise, by summing over all possible successor states s' that can be reached by taking action a , we obtain

$$\begin{aligned} Q^\pi(s, a) &= E(G_t | S_t = s, A_t = a) \\ &= \sum_{s' \in S} P_{ss'}^a E[G_t | S_{t+1} = s', S_t = s, A_t = a] \\ &= \sum_{s' \in S} P_{ss'}^a E[R_{t+1} + \sum_{k=1}^{H-1} \gamma^k R_{t+k+1} | S_{t+1} = s', S_t = s, A_t = a] \\ &= \sum_{s' \in S} P_{ss'}^a (R_{ss'}^a + E[\sum_{k=0}^{H-1} \gamma^k R_{t+k+2} | S_{t+1} = s']) \\ &= \sum_{s' \in S} P_{ss'}^a (R_{ss'}^a + \gamma E[\sum_{k=0}^{H-1} \gamma^k R_{(t+k)+k+1} | S_{t+1} = s']) \\ &\stackrel{H=\infty}{=} \sum_{s' \in S} P_{ss'}^a (R_{ss'}^a + \gamma E[G_{t+1} | S_{t+1} = s']) \\ &= \sum_{s' \in S} P_{ss'}^a (R_{ss'}^a + \gamma V^\pi(s')) \end{aligned}$$

?

$$Q^\pi(s, a) = \sum_{s' \in S} P_{ss'}^a V^\pi(s') \quad \text{assuming reward has constant value...}$$

if environment is deterministic,
there would only be one possible successor state s' if we take action a in state s

Note that this is only exact for infinite horizon ($H=\infty$)

→ $H < \infty$ is rarely used (as rigid)

normally $H=\infty$ with discount factor, which means don't have to look 1000 steps into the future, because with e.g. $\gamma=0.9$, the possible rewards after 20-30 steps are so small we can ignore them.

Bellman equation \rightarrow depends on policy!

\rightarrow foundation for a lot of algorithms in RL

plugging equations into each other, yields recursive relation

\rightarrow memorize!!!!

$$V^\pi(s) = \sum_{a \in A} \pi(a|s) \sum_{s' \in S} P_{ss'}^a (R_{ss'}^a + \gamma V^\pi(s'))$$

$$Q^\pi(s, a) = \sum_{s' \in S} P_{ss'}^a (R_{ss'}^a + \gamma \sum_{a' \in A} \pi(a'|s') Q^\pi(s', a'))$$

Bellman equation for value function

\rightarrow value of state s is fundamental of sum of the value of all the successor states plus the reward for taking the step $s \rightarrow s'$ taking into account the probabilities of the actions and successor states

if π, γ, R, P are fixed, it can be shown, that the true value function of the states is the only solution for the Bellman equation

\hookrightarrow across all states, i.e. if we have N states, we have N equations to solve

agent can now quantify properties of environment, and the quality of its policy

\rightarrow importantly, agent can also use V and/or Q to improve the policy!

which is the goal of RL

greedy policy

assume Q (or V) are given (and plz. i.e. we have a table), then we can derive a policy as follows:

$$\pi(s) = \arg \max_a Q(s, a) \stackrel{?}{=} \arg \max_a \sum_{s'} P_{ss'}^a V(s')$$

bc, because $R_{ss'}^a$ is missing

more generally

\cdot ϵ -greedy policy: select greedy action $1-\epsilon$ of time and some other, random action ϵ of the time
 \cdot stochastic policy: select actions probabilistically (e.g. follows the relative values of Q)
 \rightarrow interesting for learning (exploration)

\hookrightarrow exploitation, on the other hand, would mean always taking best action, as far as I've learned them, so I can exploit my knowledge

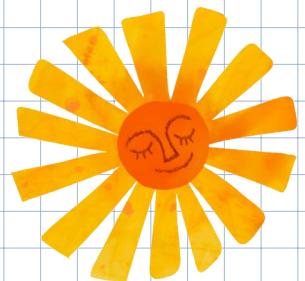
this does not solve the RL objective, because V and Q depend on the policy!!!

optimal (action-)value function

maximize for each state

$$V^*(s) = \max_\pi V^\pi(s)$$

$$Q^*(s, a) = \max_\pi Q^\pi(s, a)$$



V^* is the value function of a policy, so it follows of course the Bellman equation

\hookrightarrow since it is the value function of the best policy (i.e. I always choose the best action), I can also write it like this:

$$V^*(s) = \max_{a \in A} \sum_{s' \in S} P_{ss'}^a (R_{ss'}^a + \gamma V^*(s'))$$

what is beautiful here: no policy mentioned, only depends on model
 \rightarrow if we find way to compute this optimal value function, we can derive a policy from it and as it turns out, that will be the optimal policy A+

if we can compute V^* , it is good to be worth it \rightarrow of course, we still need to approximate V^*

in nutshell: optimal value function is Value function of optimal policy, but can be written in a form that depends only on model

Dynamic Programming

→ can be used if we have a model with finite set of states, finite set of actions, we know the probabilities of state transitions, we know the reward function

Policy Iteration

start with random policy, then iteratively repeat 2 steps → Policy Evaluation: compute value function of current policy
→ converges to optimal policy → Policy Improvement: improve policy wrt current value function

Policy Evaluation

compute V iteratively for given π using Bellman equation as update rule: fixed point iteration

converges to true value function of the policy, because it is the only fixed point to the set of equations

$$V_\pi(s) \leftarrow \sum_{a \in A} \pi(a|s) \sum_{s' \in S} P_{ss'}^a (R_{ss'}^a + \gamma V_\pi(s'))$$

Policy Improvement Theorem

let π and π' be deterministic policies such that $Q^\pi(s, \pi'(s)) \geq V^\pi(s) \forall s \in S \Rightarrow$ then π' is better than π and $V^{\pi'}(s) \geq V^\pi(s) \forall s \in S$
furthermore if inequality on left is strict, the second one is strict as well

Policy Improvement

derive new policy using information from old policy and model:
at each state, select action which appears best

$$\pi'(s) = \operatorname{argmax}_{a \in A} \sum_{s' \in S} P_{ss'}^a (R_{ss'}^a + \gamma V^\pi(s'))$$

due to PIT, new policy is strictly better than old one, unless optimal policy has already been found

in finite MDP only finitely many deterministic policies → algorithm guaranteed to converge (but can be slow for large state spaces)

Value Iteration

start with random value function, then iteratively improve value function (greedily)

→ converges to optimal value function

→ compute optimal policy from optimal value function

initialize $V_0(s) = 0 \forall s$

directly optimizes V using Bellman optimality equation as an update rule:

$$V_k(s) \leftarrow \max_{a \in A} \sum_{s' \in S} P_{ss'}^a (R_{ss'}^a + \gamma V_{k-1}(s'))$$

can be seen as a special case of policy iteration, where policy evaluation step is truncated

example:

robot should reach diamond (reward +1) and avoid bomb (reward -1), states with reward are final
12 states (really 11 and the wall), 4 actions
assume discounting, $\gamma = 0.9$



robot should reach diamond (reward +1) and avoid bomb (reward -1), states with reward are final
12 states (really 11 and the wall), 4 actions, discount $\gamma = 0.9$

initialization

0.0	0.0	0.0	0.0
0.0	XXX	0.0	0.0
0.0	0.0	0.0	0.0

first iteration

0.0	0.0	0.0	1.0
0.0	XXX	0.0	-1.0
0.0	0.0	0.0	0.0

second iteration

0.0	0.0	0.9	1.0
0.0	XXX	0.0	-1.0
0.0	0.0	0.0	0.0

third iteration

0.0	0.81	0.9	1.0
0.0	XXX	0.81	-1.0
0.0	0.0	0.0	0.0

fifth iteration

0.73	0.81	0.9	1.0
0.66	XXX	0.81	-1.0
0.0	0.66	0.73	0.66

finished reached fixed point

0.73	0.81	0.9	1.0
0.66	XXX	0.81	-1.0
0.59	0.66	0.73	0.66

From the optimal value function, we can compute the optimal policy:

optimal value function associated policy

0.73	0.81	0.9	1.0
0.66	XXX	0.81	-1.0
0.59	0.66	0.73	0.66

optimal value function optimal policy

in more complicated scenarios, state value estimates still get updated after the initial estimate →
iterate until little change happens

Power of Dynamic Programming:

planning-based → don't need to run the environment / or simulation

as soon as we have more complicated situations, we don't have access to all the information necessary for DP → need to run env. in simulator or whatever and generate actual experience
→ BUT: we still do a similar thing (evaluate policy etc.)

in practical cases model may be untrustable / inexistent / only be partially observable / continuous

... i.e. noisy
With probability 0.2, robot goes to an arbitrary adjacent state instead of performing the action

first iteration second iteration third iteration

0.0	0.0	0.0	1.0
0.0	XXX	0.0	-1.0
0.0	0.0	0.0	0.0

0.0	0.0	0.05	1.0
0.0	XXX	0.05	-1.0
0.0	0.0	0.0	0.0

0.0	0.59	0.76	1.0
0.0	XXX	0.54	-1.0
0.0	0.0	0.0	0.0

Tabular Model-free

state space is known, finite, and observable

→ we can make a table of state-action values and try to estimate/optimize them

S	A	Q
1		
2		
3		

Sampling-based Approximations

Bellman equation:

$$V^\pi(s) = \sum_{a \in A} \pi(a|s) \sum_{s' \in S} P_{ss'}^a (R_{ss'}^a + \gamma V^\pi(s')) = \sum_{a \in A} \sum_{s' \in S} \pi(a|s) P_{ss'}^a (R_{ss'}^a + \gamma V^\pi(s'))$$

Policy evaluation

$$V_k(s) \leftarrow \sum_{a \in A} \pi(a|s) \sum_{s' \in S} P_{ss'}^a (R_{ss'}^a + \gamma V_{k-1}(s'))$$

in model-free RL we don't have $P_{ss'}^a$ and $R_{ss'}^a$ → replace sum with expectation

$$V_k(s) \leftarrow \mathbb{E}_{a,s'} [R_{ss'}^a + \gamma V_{k-1}(s')]$$

estimate expectation by sampling:

$$\hat{V}(s) = \frac{1}{N} \sum_{n=1}^N G_{\text{smp}}^{(n)}(s) \quad \text{note } V(s) := \mathbb{E}[G_t \mid S_t = s] \text{ where } G_t = \sum_{k=0}^H \gamma^k R_{t+k+1}$$

or using a running average:

$$\hat{V}(s) \leftarrow (1-\alpha) \hat{V}_{\text{old}}(s) + \alpha G_{\text{smp}}(s) \quad \text{what policy is used?}$$

Value Iteration

$$V_k(s) \leftarrow \max_a \sum_{s'} P_{ss'}^a (R_{ss'}^a + \gamma V_{k-1}(s'))$$

issue: doesn't really work for sampling-based framework, because cannot sample across max operator not really works
→ that's why policy iteration, even though it's more complicated and needs more iterations, is so important

Q function

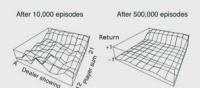
if we do not have access to state transitions (we don't know which actions lead to which states), V is not enough

$$\hat{Q}(s, a) = \frac{1}{N} \sum_{n=1}^N G_{\text{smp}}^{(n)}(s, a)$$

issue: depending on policy, certain state-action pairs never occur
→ stronger requirement than reaching all states

Blackjack example

- Assume a policy, e.g. hit until sum >= 20
- rewards: +1 for win, 0 for draw, -1 for loss
- simulate many games, starting from all possible states, and average rewards



Value function for the example policy, estimated via MC. After 500000 episodes, convergence has been achieved. Do you see why the strategy is not very good?

Monte Carlo Methods

idea: collect data from entire episodes

- start anywhere and let agent collect rewards

- assume agent always reaches terminal state (episodic task), so that we get an estimate for the return

- estimate state values or action-values by averaging over observed returns!

first-visit: only consider first visit to my state

every-visit: consider all visits

conceptually simple but don't converge

Monte Carlo on Actions

agent: If we do not have full knowledge of state transitions, collecting state values is not good enough
→ need to explicitly estimate Q function in order to suggest a policy