

启发式动态规划在污水处理过程控制中的应用

薄迎春^{1,2†}, 乔俊飞²

(1. 中国石油大学(华东) 信息与控制工程学院, 山东 青岛 266580; 2. 北京工业大学 电子信息与控制工程学院, 北京 100124)

摘要: 针对溶解氧及硝态氮浓度的跟踪控制问题, 提出了一种基于回声状态网络的启发式动态规划控制方法, 该方法首先对当前策略进行评价, 然后根据评价结果对当前策略进行调整, 这个过程交替进行, 直至发现最优的控制策略. 评价函数及控制策略的逼近均采用回声状态网络实现. 为保证控制器的可用性, 对控制器学习过程的参数选择范围进行了分析. 污水处理过程的控制实验表明, 该方法能够显著提高系统控制的平稳性及控制精度.

关键词: 启发式动态规划; 污水处理; 溶解氧; 回声状态网络; 收敛性

中图分类号: TP273 文献标识码: A

Application of heuristic dynamic programming to wastewater treatment process control

BO Ying-chun^{1,2†}, QIAO Jun-fei¹

(1. College of Information and Control Engineering, China University of Petroleum (East China), Qingdao Shandong 266580, China;
2. College of Electronic and Control Engineering, Beijing University of Technology, Beijing 100124, China)

Abstract: To control the dissolved oxygen concentration and nitrate concentration of wastewater treatment process (WWTP), we propose a heuristic dynamic programming (HDP) method based on the echo state network (ESN). This scheme evaluates the current policy, and then the HDP controller adjusts the current policy according to the evaluation results. This procedure is continued until the optimal policy is found. Echo state networks are adopted to approximate the evaluation function and the optimal policy. To ensure the availability of the HDP controller, we investigate the parameter ranges of the HDP controller. Experimental results indicate that the ESN-based HDP controller has advantages over the PID controller in control stability and control precision.

Key words: heuristic dynamical programming (HDP); wastewater treatment; dissolved oxygen; echo state network (ESN); convergence

1 引言(Introduction)

污水处理是关系环境可持续发展的关键问题之一. 污水处理过程非线性较强, 各控制目标间耦合严重, 入水流量及污染物浓度频繁变化, 这些因素增加了污水处理过程控制的复杂性, 也导致常规控制效果较差^[1]. 近年来, 模型预测控制(model predictive control, MPC)在该领域得到了一定的研究^[2-3]. 如文献[2]基于MPC方法对污水处理过程的溶解氧浓度进行控制. 为提高系统的除氮性能, 其在MPC的基础上增加了入水氨氮的前馈补偿, 该方法取得了较好的除氮效果. 但是, 由于污水处理过程模型的时变特性和反应过程的不确定性, 建立适合控制需求的污水处理过程模型较为困难^[3], 这也增加了MPC在污水处理领域应用的难度^[4].

启发式动态规划(heuristic dynamic programming,

HDP)^[5-6]是解决模型难以建立的控制问题的有效方法之一. HDP控制器的设计过程对系统的动力学模型依赖较小, 甚至可以在无模型的方式进行^[7]. 典型的HDP控制器主要由评价模块和控制模块两大部分组成^[5]. 首先给定一个控制策略, 系统在该策略下发生状态转移, 评价模块根据状态的反馈信息对当前的控制策略的效果进行评价(策略评价), 然后控制模块根据评价结果对当前的策略进行调整(策略提升), 这个过程不断重复, 直至寻找到最优的控制策略. 由于HDP能够在与环境的交互中逐渐学习到最优的策略, 所以HDP本质上属于一种自适应的学习控制方法. 目前, HDP已经在发电^[8]、复杂非线性系统控制^[9]、交通^[10]等多个领域得到了研究和应用.

针对污水处理过程模型难以建立的问题, 提出了一种基于回声状态网络的HDP控制方法, 并将其应用

于污水处理过程的溶解氧和硝态氮浓度控制。

2 污水处理过程简介(Review of wastewater treatment process)

工业和生活污水经过一系列前置处理措施进入污水处理装置的生化反应区, 经过一系列硝化-反硝化过程, 最终处理后的污水经二次沉淀池沉淀后排入接纳水体(河流、湖泊等)。生化反应装置是污水处理过程的核心环节, 其主要由厌氧池和好氧池两大部分组成(如图1所示)。

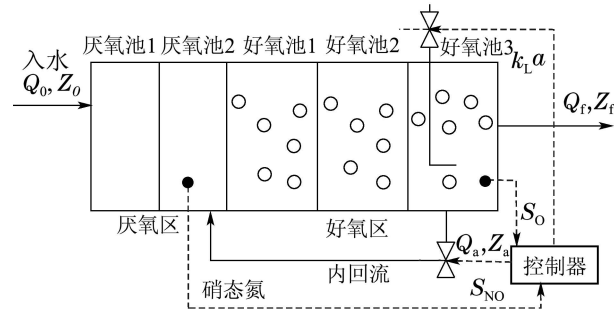


图 1 污水处理装置生化反应过程示意图

Fig. 1 Schematic diagram of biochemical process of wastewater treatment plant

有机氮的去除是生化反应过程的主要目标之一。好氧区溶解氧浓度(S_O)及厌氧区硝态氮浓度(S_{NO})对有机氮的去除效果有重要影响。 S_O 的水平过低, 会使污泥活性降低, 抑制生物对有机物的降解, 易产生污泥膨胀; 而 S_O 浓度过高会加速消耗污水中的有机物, 增加系统能耗^[3]。 S_{NO} 浓度过低, 会影响反硝化过程

的速度, 而其浓度过高则会增加出水硝态氮的含量, 进而增加后续处理成本。所以, 实现 S_O 及 S_{NO} 的平稳控制对于提高出水质量至关重要。

3 启发式动态规划控制器设计(Design of the heuristic dynamic programming controller)

污水处理过程可表示为

$$\mathbf{x}(k+1) = F(\mathbf{x}(k), \mathbf{u}(k)), \quad (1)$$

其中: $\mathbf{x}(k)$ 是系统状态, $\mathbf{u}(k)$ 为控制量, 其维数分别为 n, m , $F(\cdot)$ 是关于 $\mathbf{x}(k)$ 和 $\mathbf{u}(k)$ 的未知非线性函数。这里 $\mathbf{u}(k)$ 是系统状态 $\mathbf{x}(k)$ 的函数, 即

$$\mathbf{u}(k) = h(\mathbf{x}(k)). \quad (2)$$

系统(1)的最优控制问题可归结为选择合适的控制策略 \mathbf{u} 实现式(3)的性能指标最小化^[5]。

$$J(\mathbf{x}(k)) = \sum_{j=k}^{\infty} \gamma^{j-k} U(\mathbf{x}(j)), \quad (3)$$

其中: $J(\mathbf{x}(k))$ 为评价函数, $U(\mathbf{x}(k))$ 为立即评价, 即一个周期内引发的成本值, $0 < \gamma \leq 1$ 是折扣因子。为方便起见, $J(\mathbf{x}(k))$ 简记为 $J(k)$, $U(\mathbf{x}(k))$ 则简记为 $U(k)$ 。式(3)也可写为

$$J(k) = U(k) + \gamma J(k+1). \quad (4)$$

根据Bellman优化原理, 最优的控制策略为^[7]

$$\mathbf{u}^*(k) = \arg \min_{\mathbf{u}} [U(k) + \gamma J(k+1)]. \quad (5)$$

HDP采用逐步逼近的方式对式(5)进行求解。其原理如图2所示。

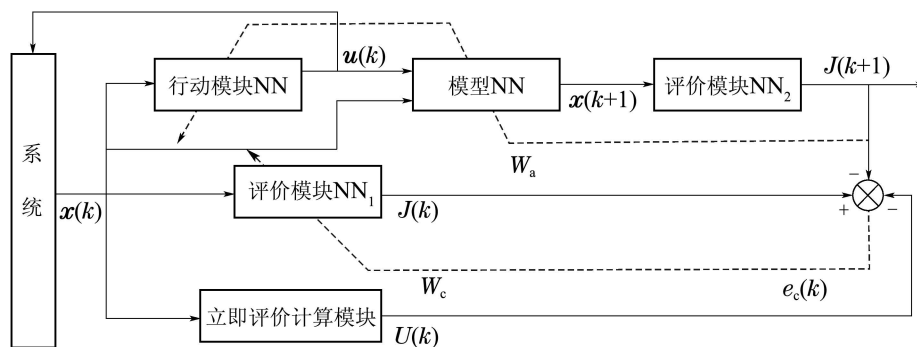


图 2 HDP控制器结构(虚线表示评价网络及行动网络的权值调整路径)

Fig. 2 Architecture of HDP controller (the dotted lines correspond to the learning path of the actor NN and the critic NN)

HDP控制器中, 策略评价和策略提升分别由评价模块(critic)和行动模块(actor)完成。这两个模块采用神经网络实现。其中: actor NN为控制网络; critic NN₁与critic NN₂为评价网络, 其结构及参数均相同, 但是critic NN₁的输入和输出分别为 $\mathbf{x}(k)$ 及 $J(k)$, 而critic NN₂的输入和输出分别为 $\mathbf{x}(k+1)$ 及 $J(k+1)$; 模型NN为辨识网络, 其功能是对系统下一步的状态进行预测; utility为立即评价计算模块。

这里所有的神经网络均采用回声状态网络(echo state network, ESN)。

3.1 回声状态网络(Echo state network)

ESN是Jaeger于2001年提出的一种新型递归神经网络^[11]。其核心是一个动态神经元池(dynamic neurons reservoir, DNR), 神经元池内包含大量神经元, 这些神经元之间采用随机、稀疏的方式连接。其结构如图3所示。

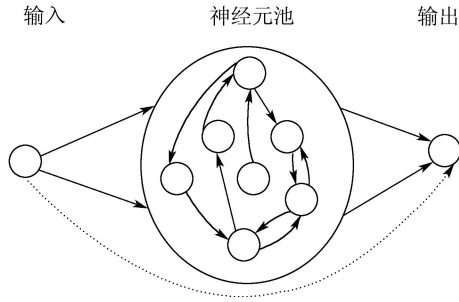


图3 ESN结构

Fig. 3 Schematic diagram of ESN structure

ESN的隐层神经元输出可表示为

$$\mathbf{s}(k+1) = f(W_{IN}\mathbf{u}(k+1) + W_R\mathbf{s}(k)), \quad (6)$$

其中: $\mathbf{s}(k) = [\mathbf{s}_1(k) \cdots \mathbf{s}_N(k)]^T$ 为内部状态, 即ESN隐层神经元输出, $\mathbf{u}(k) = [\mathbf{u}_1(k) \cdots \mathbf{u}_K(k)]^T$ 为ESN输入. W_{IN} , W_R 分别为ESN输入到隐层及隐层神经元之间的连接权值矩阵, K , N 分别为输入和隐层的神经元个数.

ESN的输出可表示为

$$\mathbf{y}(k+1) = W^T\mathbf{s}(k+1), \quad (7)$$

这里: $\mathbf{y}(k) = [\mathbf{y}_1(k) \cdots \mathbf{y}_L(k)]^T$ 为ESN的输出, L 为输出神经元个数, W 为隐层到输出的连接权值矩阵. ESN中, W_{IN} , W_R 随机确定, 并在网络学习过程中保持不变^[11], 也就是说, ESN只需对 W 进行学习^[12], 所以, ESN各输出对应的权值调整是独立的.

3.2 HDP控制器的学习(Learning of HDP controller)

HDP控制器的设计过程实质上是评价和控制网络通过权值参数的调整逐渐逼近最优策略的过程.

3.2.1 评价网络的学习(Learning of critic network)

设 $W_c(k)$ 为评价网络的可调权值, 则 W_c 的调整可通过最小化 E_c 进行,

$$E_c(k) = \frac{1}{2}e_c^2(k), \quad (8)$$

其中

$$e_c(k) = J(k) - U(k) - \gamma J(k+1). \quad (9)$$

当 $E_c(k) = 0$ 时, 根据式(9)可知, $J(k) = U(k) - \gamma J(k+1)$ 成立. 这意味着当 $E_c(k) = 0$ 时, 式(9)与式(3)等价, 即评价网络的输出即为当前策略下的成本函数值. 这说明当评价网络的训练达到稳态时, 评价网络可实现状态到评价值的映射.

根据梯度下降算法, 评价网络的权值增量为

$$\begin{aligned} \Delta W_c(k) &= -\eta_c(k) \left(\frac{\partial E_c(k)}{\partial W_c(k)} \right)^T = \\ &= -\eta_c(k) e_c(k) \left(\frac{\partial e_c(k)}{\partial W_c(k)} \right)^T = \end{aligned}$$

$$-\eta_c(k) e_c(k) \left(\frac{\partial J(k)}{\partial W_c(k)} \right)^T, \quad (10)$$

这里 $\eta_c(k)$ 为评价网络的学习率. 因为 $J(k)$ 是评价网络的输出, 根据式(7), 可知

$$\frac{\partial J(k)}{\partial W_c(k)} = \mathbf{s}_c^T(k), \quad (11)$$

$\mathbf{s}_c(k)$ 为评价网络隐层神经元的输出. 将式(11)代入式(10)可得

$$\Delta W_c(k) = -\eta_c(k) e_c(k) \mathbf{s}_c(k). \quad (12)$$

3.2.2 控制网络的学习(Learning of actor network)

控制网络权值调整的目标是最小化成本函数 $J(k)$, 所以其性能指标可设定为^[6]

$$\frac{\partial J(k)}{\partial W_a(k)} = \frac{\partial U(k)}{\partial W_a(k)} + \gamma \frac{\partial J(k+1)}{\partial W_a(k)} = 0, \quad (13)$$

其中 $W_a(k)$ 为控制网络的可调权值. 由于各控制量对应的权值调整是独立的, 所以, 仅以控制量 \mathbf{u}_i 对应的权值向量 $W_{a,i}$ 为例说明控制网络的权值调整过程. 根据梯度下降算法及链式求导法则, 控制网络的权值增量为

$$\Delta W_{a,i}(k) = -\eta_{a,i}(k) \Xi_i(k) \left(\frac{\partial \mathbf{u}_i(k)}{\partial W_{a,i}(k)} \right)^T, \quad (14)$$

这里

$$\Xi_i(k) = \frac{\partial U(k)}{\partial \mathbf{u}_i(k)} + \gamma \frac{\partial J(k+1)}{\partial \mathbf{u}_i(k)}. \quad (15)$$

因为 $\mathbf{u}_i(k)$ 是评价网络的输出, 根据式(7), 有

$$\frac{\partial \mathbf{u}_i(k)}{\partial W_{a,i}(k)} = \mathbf{s}_a^T(k), \quad (16)$$

$\mathbf{s}_a(k)$ 为控制网络隐层神经元的输出. 将式(16)代入式(14)可得

$$\Delta W_{a,i}(k) = -\eta_{a,i}(k) \Xi_i(k) \mathbf{s}_a(k). \quad (17)$$

3.3 控制器的收敛性分析(Convergence analysis of HDP controller)

定义离散的Lyapunov函数为

$$L(k) = \frac{1}{2}e^2(k), \quad (18)$$

$e(k)$ 为与权值调整过程相应的误差性能指标, 则

$$\begin{aligned} \Delta L(k) &= \frac{1}{2}e^2(k+1) - \frac{1}{2}e^2(k) = \\ &= \Delta e(k) \left[e(k) + \frac{1}{2} \Delta e(k) \right], \end{aligned} \quad (19)$$

这里 $\Delta e(k) = e(k+1) - e(k)$. 设 W 为神经网络权值, 根据全微分定理, 有

$$\Delta e(k) = \frac{\partial e(k)}{\partial W(k)} \Delta W(k). \quad (20)$$

假设评价及控制网络的权值有界. 即存在正数 w_{cM} 及 w_{aM} , 使 $\|W_c(k)\| < w_{cM}$, $\|W_a(k)\| < w_{aM}$.

定理 1 若 $\eta_c(k)$ 满足

$$\eta_c(k) < \frac{2}{\|\mathbf{s}_c(k)\|^2}, \quad (21)$$

则评价网络的权值学习过程是收敛的.

证 定义 $e(k) = e_c(k)$, $W(k) = W_c(k)$. 根据式(9)(11), 有

$$\frac{\partial e(k)}{\partial W(k)} = \frac{\partial J(k)}{\partial W(k)} = \mathbf{s}_c^T(k). \quad (22)$$

将式(12)及式(22)代入式(20), 可得

$$\Delta e(k) = -\eta_c(k)e(k)\|\mathbf{s}_c(k)\|^2, \quad (23)$$

将式(23)代入式(19), 可得

$$\Delta L(k) = -\eta_c(k)e^2(k)\|\mathbf{s}_c(k)\|^2 \cdot \left[1 - \frac{1}{2}\eta_c(k)\|\mathbf{s}_c(k)\|^2\right], \quad (24)$$

所以只要

$$1 - \frac{1}{2}\eta_c(k)\|\mathbf{s}_c(k)\|^2 > 0, \quad (25)$$

则 $\Delta L(k) \leq 0$, 并且仅当 $e(k)=0$ 时, $\Delta L(k)=0$. 解不等式(25)即可得到式(21). 根据离散系统的Lyapunov理论, 当式(21)成立时, 评价网络值学习过程是收敛的. **证毕.**

定理 2 若 $\eta_{a,i}(k)$ 满足

$$\eta_{a,i}(k) < \frac{2(U(k) + \gamma J(k+1))}{\Xi_i^2(k)\|\mathbf{s}_a(k)\|^2}, \quad (26)$$

则控制网络的权值调整过程是收敛的.

证 令

$$e(k) = U(k) + \gamma J(k+1), \quad W(k) = W_{a,i}(k),$$

则

$$\frac{\partial e(k)}{\partial W(k)} = \frac{\partial U(k)}{\partial W_{a,i}(k)} + \gamma \frac{\partial J(k+1)}{\partial W_{a,i}(k)}. \quad (27)$$

根据链式偏导法则, 并将式(16)代入式(27), 可得

$$\frac{\partial e(k)}{\partial W(k)} = \Xi_i(k)\mathbf{s}_a^T(k), \quad (28)$$

将式(17)(28)代入式(20)可得

$$\Delta e(k) = -\eta_{a,i}(k)\Xi_i^2(k)\|\mathbf{s}_a(k)\|^2, \quad (29)$$

将式(29)代入式(19)可得

$$\Delta L(k) = -\eta_{a,i}(k)\Xi_i^2(k)\|\mathbf{s}_a(k)\|^2 \cdot \left[e(k) - \frac{1}{2}\eta_{a,i}(k)\Xi_i^2(k)\|\mathbf{s}_a(k)\|^2\right]. \quad (30)$$

当满足

$$e(k) > \frac{1}{2}\eta_{a,i}(k)\Xi_i^2(k)\|\mathbf{s}_a(k)\|^2, \quad (31)$$

$\Delta L(k) \leq 0$. 解不等式(31), 即可得到式(26). **证毕.**

值得说明的是, 与评价网络不同, 在控制网络的

学习过程中, $e(k) = U(k) + \gamma J(k+1)$ 并不一定趋向于0(这取决于立即评价的定义方式), 而是趋向于一个极小值. 从式(30)可以看出, 在满足式(26)条件下, 仅当 $\Xi_i(k) = 0$ 时, $\Delta L(k) = 0$. 而 $\Xi_i(k) = 0$ 正是控制网络学习的目标(见式(13)).

4 实验研究(Experimental studies)

由于不同污水处理过程的入水情况不同, 这使得控制策略比较十分困难. 为此, 国际水协提出了一个用于污水处理过程控制策略测试及比较的标准平台(benchmark simulation no. 1, BSM1). BSM1对污水处理过程装置的标准模型、入水条件、控制策略的测试方法及控制策略的评价标准都进行了明确的定义, 该平台能够实现不同控制策略在相同条件下的客观比较. 所以, 本文所有实验均基于BSM1进行.

4.1 立即评价设计(Design of utility)

HDP控制器设计中, 首先要确定立即评价指标 $U(k)$ 的形式. 对于污水处理过程的跟踪控制而言, 溶解氧及硝态氮浓度的控制精度对于出水质量有较大影响. 其控制目标是使溶解氧和硝态氮浓度尽可能接近相应的设定值, 所以, 将 $U(k)$ 的形式定义为

$$U(k) = \frac{1}{2}\mathbf{E}^T(k)\mathbf{E}(k), \quad (32)$$

这里:

$$\mathbf{E}(k) = [E_1(k) \ E_2(k)],$$

$$E_1(k) = R_{SO}(k) - S_O(k),$$

$$E_2(k) = R_{SNO}(k) - S_{NO}(k),$$

$$R_{SO} = 2 \text{ mg/L}, \quad R_{SNO} = 1 \text{ mg/L}$$

分别为溶解氧及硝态氮浓度的设定值^[13].

4.2 参数初始化(Parameter initialization)

评价网络的输入为 $\mathbf{x}(k) = [S_O(k) \ S_{NO}(k)]^T$, 输出为 $J(k)$; model ESN的输入为 $[K_{La}(k) \ Q_a(k) \ S_O(k) \ S_{NO}(k)]^T$, 控制网络的输入为

$$[K_{La}(k) \ Q_a(k) \ S_O(k) \ S_{NO}(k)$$

$$X_{BH}(k) \ X_{BA}(k) \ Q_{in}(k)]^T,$$

输出为 $[\Delta K_{La}(k) \ \Delta Q_a(k)]^T$. $K_{La}(k)$ 为好氧区的空气转换系数, $Q_a(k)$ 为内回流流量, 而 $X_{BH}(k)$, $X_{BA}(k)$ 和 $Q_{in}(k)$ 分别为入水异养菌浓度、入水自养菌浓度以及入水流量. 系统的采样周期 $T = 1.25 \times 10^{-2} \text{ h} = 45 \text{ s}$; 评价网络, 控制网络及辨识网络的隐层神经元个数分别为40, 40和50; 其谱半径^[11]约为0.68, 0.68和0.76; 折扣因子 $\gamma = 0.95$. 实验主要将HDP与PID控制器进行对比, PID控制器的参数

为: 对于溶解氧控制回路, $K_P = 200$, $T_I = 15$, $T_D = 2$; 对于硝态氮控制回路,

$$K_P = 20000, T_I = 5000, T_D = 400.$$

4.3 实验结果及分析(Result and analysis)

4.3.1 解耦能力测试(Decoupling ability test)

首先对HDP的解耦能力进行测试, 将入水流量固定在 $Q_{in} = 21471 \text{ m}^3/\text{d}$, 溶解氧及硝态氮浓度的设定值按照如下的规则调整:

$$\begin{cases} R_{S_{O_2}} = 2, R_{S_{NO_3}} = 1, & t < 0.5, \\ R_{S_{O_2}} = 2, R_{S_{NO_3}} = 0.8, & 0.5 \leq t < 0.75, \\ R_{S_{O_2}} = 1.9, R_{S_{NO_3}} = 0.8, & t \geq 0.75, \end{cases} \quad (33)$$

式中 t 表示运行时间(d). 在被控量设定值按照上述规律改变时, HDP及PID的控制效果如图4所示.

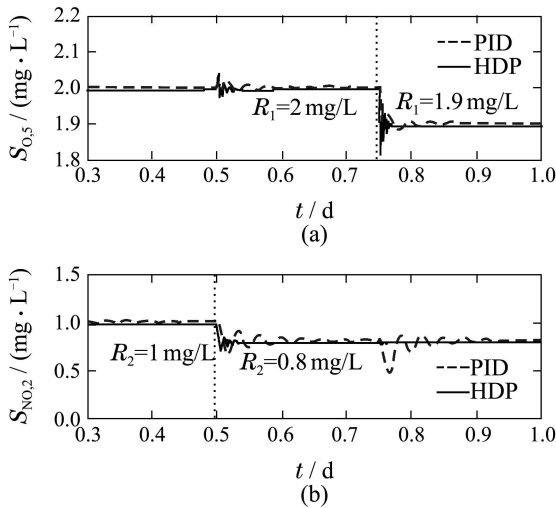


图4 HDP的解耦性能

Fig. 4 Decoupling ability of HDP controller

由图4可看出, 当 $R_{S_{NO_3}}$ 从 1 mg/L 变为 0.8 mg/L 时(图4(b)), S_{O_2} 产生了较小的波动(图4(a)), 这说明硝态氮浓度变化对好氧区溶解氧浓度影响较小. 在PID控制作用下, 当 $R_{S_{O_2}}$ 从 2 mg/L 降低到 1.9 mg/L 时(图4(a)), S_{NO_3} 波动较大(图4(b)). 这说明好氧区溶解氧浓度变化对厌氧区硝态氮浓度有较大影响, 二者耦合作用明显. 在HDP控制器作用下, $R_{S_{O_2}}$ 改变时, S_{NO_3} 只出现了轻微的波动, 这表明HDP控制器解耦能力较强.

4.3.2 控制性能测试(Control performance test)

接下来对HDP控制器的整体控制性能进行测试, BSM1提供了14天的入水数据, 这些数据均来自实际的污水处理厂. 入水流量及部分污染物浓度变化如图5所示.

图6为HDP和PID控制器的控制效果. 可以看出, HDP控制器作用下, S_{O_2} 及 S_{NO_3} 的波动范围明显减

小, 控制精度明显提高. 当较大的干扰出现时(第8-12天, 入水流量及入水污染物浓度发生了较大的波动), HDP的溶解氧浓度控制性能较正常时略有下降(图5(a)第8-12天), 但经过一段时间的学习, HDP的控制性能得到恢复, 这说明HDP自身具有较强的自学习能力.

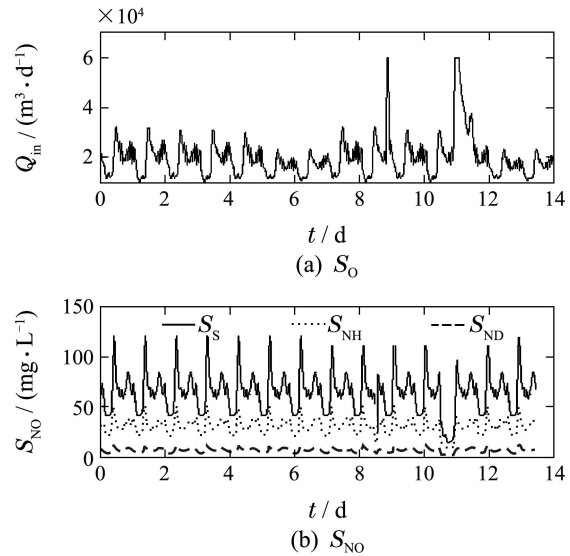


图5 入水流量及部分入水污染物浓度

Fig. 5 Influent flow rate and partial influent pollution concentration

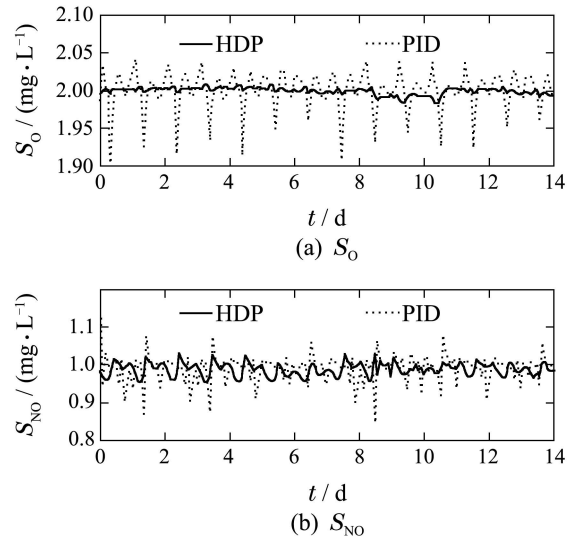


图6 HDP与PID的控制性能比较

Fig. 6 Comparison of HDP controller and PID controller

表1列出了HDP及PID控制器的与控制精度相关的评价指标. 对于溶解氧浓度控制, HDP与PID比较, 绝对平均误差降低了73.4%, 最大偏差降低了79.1%, 误差方差降低了1个数量级. 对于硝态氮浓度控制, HDP与PID相比较, 绝对平均误差降低了56.3%, 最大偏差降低了74.4%, 误差方差降低了1个数量级. 可以看出, HDP在控制精度及控制的平稳性上明显优于PID控制器.

表1 2种控制器的控制精度评价指标

Table 1 Evaluation indices of control precision under two controllers

控制变量	PID		HDP	
	S_{O_2}	S_{NO_3}	S_{O_2}	S_{NO_3}
绝对平均误差/($\text{mg} \cdot \text{L}^{-1}$)	0.0139	0.0268	0.0037	0.0117
绝对最大误差/($\text{mg} \cdot \text{L}^{-1}$)	0.0980	0.1756	0.0204	0.0450
误差方差/($\text{mg} \cdot \text{L}^{-1}$) ²	2.54×10^{-4}	0.0016	4.16×10^{-5}	2.20×10^{-4}

在HDP的在线学习过程中, 学习率的选择为

$$\begin{cases} \eta_c(k) = \frac{1}{\|\mathbf{s}_c(k)\|^2}, \\ \eta_{a,i}(k) = \frac{U(k) + \gamma J(k+1)}{\Xi_i^2(k) \|\mathbf{s}_a(k)\|^2 + \epsilon}, \end{cases} \quad (34)$$

这里 ϵ 为一小的正数.

5 结论(Conclusions)

溶解氧和硝态氮浓度是污水处理过程的两个关键指标, 由于实际污水处理过程的生化反应模型难以建立, 并且两个变量之间存在严重的耦合, 所以其控制难度较大. 启发式动态规划采用策略评价和策略提升的方法逐步逼近最优的控制策略, 在策略搜索的过程中HDP能够综合考虑各指标之间的耦合关系, 所以HDP控制器具有较强的解耦能力. 此外, HDP不需要污水处理的生化反应过程的机理模型, 其通过对输入输出的测量数据的学习即可实现策略评价和策略提升. 即HDP能够以数据驱动的模式运行. 对于污水处理过程这类机理模型难以建立的复杂系统, HDP实现简单, 有较好的应用前景.

参考文献(References):

- [1] 史雄伟, 乔俊飞, 苑明哲. 基于改进粒子群优化算法的污水处理过程优化控制 [J]. 信息与控制, 2011, 40(5): 698 - 703. (SHI Xiongwei, QIAO Junfei, YUAN Mingzhe. Optimal control for wastewater treatment process based on improved particle swarm optimization algorithm [J]. *Information and Control*, 2011, 40(5): 698 - 703.)
- [2] SHEN W H, CHEN X Q, PONS M N. Model predictive control for wastewater treatment process with feedforward compensation [J]. *Chemical Engineering Journal*, 2009, 155(1/2): 161 - 174.
- [3] 丛秋梅, 柴天佑, 余文. 污水处理过程的递阶神经网络建模 [J]. 控制理论与应用, 2009, 26(1): 8 - 14. (CONG Qiumei, CHAI Tianyou, YU Wen. Modeling wastewater treatment plant via hierarchical neural networks [J]. *Control Theory & Applications*, 2009, 26(1): 8 - 14.)
- [4] DELLANA S A, WEST D. Predictive modeling for wastewater applications: Linear and nonlinear approaches [J]. *Environmental Modelling & Software*, 2009, 24(1): 96 - 106.
- [5] WANG F Y, ZHANG H G, LIU D R. Adaptive dynamic programming: an introduction [J]. *IEEE Computational Intelligence Magazine*, 2009, 4(2): 39 - 47.
- [6] PROKHOROV D V, WUNSCH D C. Adaptive critic designs [J]. *IEEE Transactions on Neural Networks*, 1997, 8(5): 997 - 1007.
- [7] LEWIS F L, VRABIE D. Reinforcement learning and adaptive dynamic programming for feedback control [J]. *IEEE Circuits and Systems Magazine*, 2009, 9(3): 32 - 50.
- [8] VENAYAGAMOORTHY G K, HARLEY R G, WUNSCH D C. Comparison of a heuristic dynamic programming and a dual heuristic programming based adaptive critics neurocontroller for a turbo-generator [C] // *Proceedings of the International Joint Conference on Neural Networks*. Como, Italy: IEEE, 2000: 233 - 238.
- [9] 曹宁, 张化光, 罗艳红, 等. 一类非线性奇异摄动系统的近似最优控制 [J]. 控制理论与应用, 2011, 28(5): 688 - 692. (CAO Ning, ZHANG Huaguang, LUO Yanhong, et al. Suboptimal control of a class of nonlinear singularly perturbed systems [J]. *Control Theory & Applications*, 2011, 28(5): 688 - 692.)
- [10] ZHAO D B, BAI X R, WANG F Y, et al. DHP method for ramp metering of freeway traffic [J]. *IEEE Transactions on Intelligent Transportation System*, 2011, 12(4): 990 - 999.
- [11] JAEGER H. Harnessing nonlinearity: predicting chaotic systems and saving energy in wireless communication [J]. *Science*, 2004, 304(5667): 78 - 80.
- [12] 韩敏, 穆大芸. 基于Adaboost算法的回声状态网络预报器 [J]. 控制理论与应用, 2011, 28(4): 601 - 604. (HAN Min, MU Dayun. Improvement of echo state network accuracy with adaboost [J]. *Control Theory & Applications*, 2011, 28(4): 601 - 604.)
- [13] ALEX J, MAGDEBURG V. *Benchmark Simulation Model No. 1 (BSM1)* [S]. IWA Taskgroup on Benchmarking of Control Strategies for WWTPs, 2008.

作者简介:

薄迎春 (1977—), 男, 博士研究生, 目前研究方向为神经计算及智能优化控制, E-mail: boyingchun@sina.com.cn;

乔俊飞 (1968—), 男, 教授, 博士生导师, 目前研究方向为复杂系统建模、智能优化控制、神经网络等, E-mail: isbox@sina.com.cn.