**Lecture 14**

# Reinforcement Learning

Instructor: Ilias Tagkopoulos

iliast@ucdavis.edu

**Step 1. Get enough data!**

Dataset

**Step 2. Do all of the data samples have labels?**

$$\begin{bmatrix} x_{11} & \cdots & x_{1m} \\ \vdots & \ddots & \vdots \\ x_{n1} & \cdots & x_{nm} \end{bmatrix} \begin{matrix} y_1 \\ \cdots \\ y_n \end{matrix}$$

Yes

No

$$\begin{bmatrix} x_{11} & \cdots & x_{1m} \\ \vdots & \ddots & \vdots \\ x_{n1} & \cdots & x_{nm} \end{bmatrix}$$

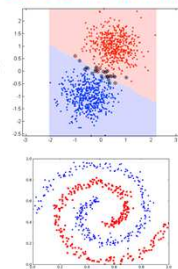**SUPERVISED LEARNING**
**Step 3: The task is to predict a continuous variable, assign a new sample to a class, or perform an optimal action?**

**UNSUPERVISED LEARNING**
**Step 3: The task is to cluster data together, find latent factors or complete missing data?**

Assign to a class

Predict continuous variable

Perform optimal actions

**Clustering**

- K-means
- Hierarchical clustering
- SOM
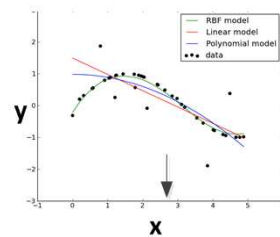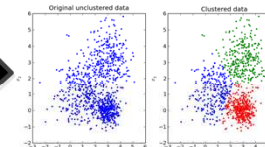
- Bayesian Classification (Naïve Bayes)
- Linear Discriminant Analysis
- Artificial Neural Networks
- Decision Trees
- Support Vector Machines

**CLASSIFICATION**

**REGRESSION**

**REINFORCEMENT LEARNING (*)**

Agent

State

Action    Reward

Environment

y

x

Linear, polynomial, logistic, …

Markov Decision Process (MDP), POMDP, Q-learning, …

**Dimensionality Reduction**

- PCA
- ICA

**Missing Data**

users

movies

- Collaborative filtering
- Market Basket analysis
- …

# What is Reinforcement Learning



- **Define states (s), actions (a), observations (o), rewards (r)**
- **States can be observed (e.g. MDP) or Latent/Hidden (e.g. POMDP)**
- **The task is to maximize the cumulative reward**

# Examples of reinforcement learning (RI)



‣ Games (e.g., G. Tesauro, D. Silver)

‣ Operations research and scheduling (e.g., W. Powell, P. Tadepalli)

‣ Recently: robotics (e.g., P. Abbeel, J. Peters, P. Stone, M. Riedmiller)

# What is the objective in Reinforcement Learning?



▸ Make optimal decisions $a^*$ by maximizing an expected utility

$$a^* \in \arg\max_a \mathbb{E}[r(a)] = \arg\max_a \sum_{j=1}^{m} r(s_j, a) p(s_j)$$

$a$ : decision

$s$ : information about environment/state

▸ Bayesian sequential decision theory (statistics)

▸ Optimal control theory (engineering)

▸ Reinforcement learning (computer science, psychology)

# One example of optimal policy

## Example: Winning the Lottery

| Actions | Outcomes |
|---------|----------|
| $a_1$: play | $s_1$: Win the lottery |
| $a_2$: don't play | $s_2$: Don't win the lottery |

**Optimal action**

$$a^* = \arg\max_{a_i} \sum_{j=1}^{2} r_{ij} p(s_j|a_i)$$

$$p(s_1|a_1) = 10^{-7} \qquad r_{11} = 500,000 \text{ USD}$$
$$p(s_2|a_1) = 1 - 10^{-7} \qquad r_{12} = -1 \text{ USD}$$
$$p(s_1|a_2) = 0 \qquad r_{21} = 0 \text{ USD}$$
$$p(s_2|a_2) = 1 \qquad r_{22} = 0 \text{ USD}$$

## Markov Decision Process: Definition

- $\mathcal{S}$: State space (finite)

- $\mathcal{A}$: Action space (finite)

- $\mathcal{P}$: Transition probability $p(s_{k+1}|s_k, a_k)$

- $r$: Reward function

- $\gamma \in [0, 1)$: Discount factor

- $\pi$: Policy

  - Deterministic: $a = \pi(s)$

  - Stochastic: $a \sim p_\pi(a|s)$ \qquad alternative notation: $p_\pi(a|s) = \pi(a|s)$



### Objective

Find a policy $\pi^*$ that maximizes the expected long-term reward

$$V^\pi(s) = \mathbb{E}\Big[\sum_{k=0}^{\infty} \gamma^k r_{k+1}\big|s_0 = s, \pi\Big], \qquad r_{k+1} = r_{k+1}(s_k, a_k, s_{k+1})$$

## Value Functions

- State Value Function: How good is it to be in a particular state $s$? Well, this depends on the current policy:

$$V^\pi(s) = \mathbb{E}[R|s_0 = s] = \mathbb{E}\Big[\sum_{k=0}^{\infty} \gamma^k r_{k+1}|s_0 = s, \pi\Big]$$

$$= \mathbb{E}[r_1 + \gamma V^\pi(s_1)|s_0 = s, \pi] \qquad \text{Self-consistency}$$

- **What you want is to find and the policy that maximize the cumulative reward, calculated by the value of each state in the trajectory**

$$\pi^*(s) = \arg\max_a \mathbb{E}[r(s, a, s') + \gamma V^*(s')]$$

# This leads to the Bellman equations

## One for state value

$$V^*(s) = \max_a Q^*(s, a)$$

$$= \max_a \mathbb{E}\Big[ \sum_{k=0}^{\infty} \gamma^k r_{k+1} | s_0 = s, a_0 = a, \pi^* \Big]$$

$$= \max_a \mathbb{E}\Big[ r_1 + \gamma \sum_{k=0}^{\infty} \gamma^k r_{k+2} | s_0 = s, a_0 = a, \pi^* \Big]$$

$$= \max_a \mathbb{E}\Big[ r_1 + \gamma V^*(s_1) | s_0 = s, a_0 = a, \pi^* \Big]$$

## Another for state-action value

$$Q^*(s, a) = \mathbb{E}\big[ r_1 + \gamma \max_{a'} Q^*(s_1, a') | s_0 = s, a_0 = a \big]$$

$$= \mathbb{E}\big[ r_1 + \gamma V^*(s_1) | s_0 = s, a_0 = a \big]$$

## This leads to the Bellman equations

One for state value

$$V^*(s) = \max_a Q^*(s, a)$$

$$= \max_a \mathbb{E}\Big[\sum_{k=0}^{\infty} \gamma^k r_{k+1} \big| s_0 = s, a_0 = a, \pi^*\Big]$$

$$= \max_a \mathbb{E}\Big[r_1 + \gamma \sum_{k=0}^{\infty} \gamma^k r_{k+2} \big| s_0 = s, a_0 = a, \pi^*\Big]$$

$$= \max_a \mathbb{E}\Big[r_1 + \gamma V^*(s_1) \big| s_0 = s, a_0 = a, \pi^*\Big]$$

Another for state-action value

$$Q^*(s, a) = \mathbb{E}\Big[r_1 + \gamma \max_{a'} Q^*(s_1, a') \big| s_0 = s, a_0 = a\Big]$$

$$= \mathbb{E}\Big[r_1 + \gamma V^*(s_1) \big| s_0 = s, a_0 = a\Big]$$

Solving them (Dynamic Programming, Monte Carlo) will find the optimal policy

# Example: Modeling Sepsis with Partially Observable MDP



**A** — Data preprocessing and Model formulation

**Database**

**Electronic Health Records** (1492 patients)

Blood test - culture
Temperature
Respiratory Rate
SBP
WBC
MAP

State are defined by tests/vitals

Antibiotics (48 total)
Vancomycin
Ceftriaxone
Vancomycin & Ceftriaxone

Actions defined by combinations of the top 5 drugs

**Vitals Likelihood**

Distribution of vitals given a specific class across all patients.

Temperature

PDF

Temperature (°C)

SIRS
Sepsis
No SIRS

**Partially Observable Markov Decision Process (POMDP) Model**

**States**
Sepsis, SIRS, Septic Shock, Bacteremia, Bacteremia Probable SIRS(BPS), Probable SIRS(PS), No SIRS, Probable Septic Shock(PSS). Absorbing states(Death,Discharge)

**Actions**
Antibiotic administration: combinations of the 5 most commonly used antibiotics

**Observations**
Vitals (Temperature, Respiratory Rate, SBP WBC, MAP)
Tests (Blood culture)

**Rewards**
Costs/rewards are state-specific and are defined by medical experts

# Example: State-action-observation-reward space

**A — Data preprocessing and Model formulation**

Database

Electronic Health Records (1492 patients)
Blood test - culture, Temperature, Respiratory Rate, SBP, WBC, MAP — State are defined by tests/vitals

Antibiotics (48 total): Vancomycin, Ceftriaxone, Vancomycin & Ceftriaxone — Actions defined by combinations of the top 5 drugs

**Partially Observable Markov Decision Process (POMDP) Model**

States: Sepsis, SIRS, Septic Shock, Bacteremia, Bacteremia Probable SIRS(BPS), Probable SIRS(PS), No SIRS, Probable Septic Shock(PSS). Absorbing states(Death,Discharge).

Actions: Antibiotic administration: combinations of the 5 most commonly used antibiotics

Observations: Vitals (Temperature, Respiratory Rate, SBP, WBC, MAP), Tests (Blood culture)

Rewards: Costs/rewards are state-specific and are defined by medical experts

Vitals Likelihood: Distribution of vitals given a specific class across all patients. Temperature — SIRS, Sepsis, No SIRS

**B — Model training**

State - action space

Legend: Absorbing states, Patient states, Actions, Current belief, Observations

Model Training: Initialization → Sample the belief space → Create a sample set of belief points, i.e. probability distributions over all states → Maximize value in sample belief set → Perform Value Iteration to find the actions that maximize the value for each belief point in the sample set → Assign optimal policy → Calculate optimal policy for each belief point in the sample set → Generalize optimal policy for all belief points → Assign optimal policy to any new belief point through neighbor clustering and interpolation → Evaluation

**C — Model execution and parameter update**

Current Belief (bi) — PDF, States S1–S8

Action affects patient status → Environment (Patient State) → Patient status affects observations → Updated Observations → Updated Belief (bi+1)

S1 = No SIRS
S2 = PrSIRS
S3 = SIRS
S4 = Sepsis
S5 = Becteremia
S6 = Bacteremia PS
S7 = Septic Shock
S8 = Septic Shock PS

**D — Policy evaluation**

n-Fold cross validation: Initialization → Split dataset → Split the database into n sets → Train the POMDP leaving out (a different one each time) set for testing → Evaluation of policy on test set by comparison of transitions using the policy drugs with the ones not using them. (Evaluation of derived policy / POMDP train)

Reduced training set evaluation: Initialization → Split dataset → Database split into 3 sets, 2 sets for training, one for testing → Train the POMDP with 2 of the sets, reducing the data by 0,10,15,20,25 and 50%. → Evaluation of policy on test set by comparison of transitions using the policy drugs with the ones not using them. (Evaluation of derived policy / POMDP train)
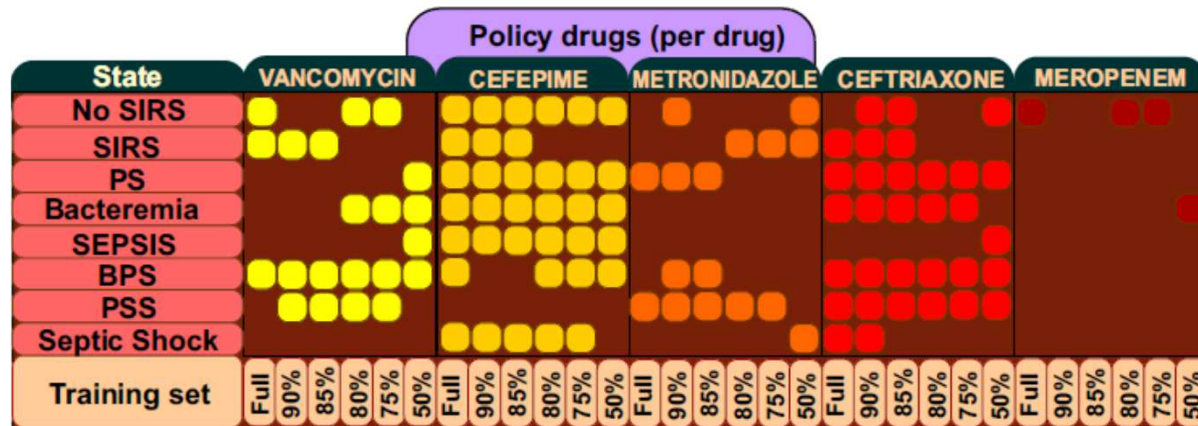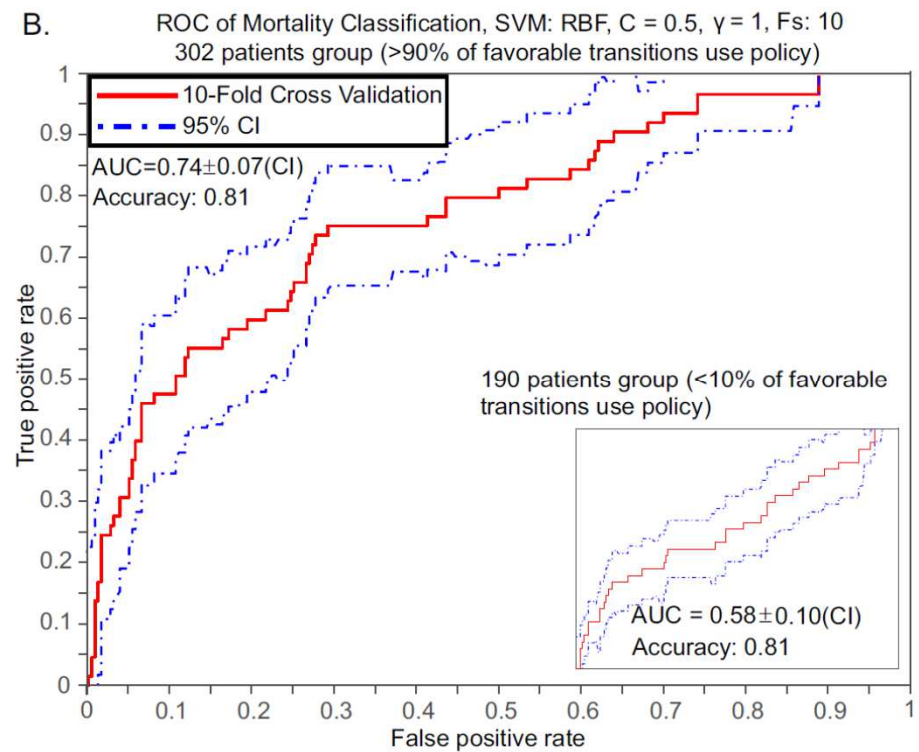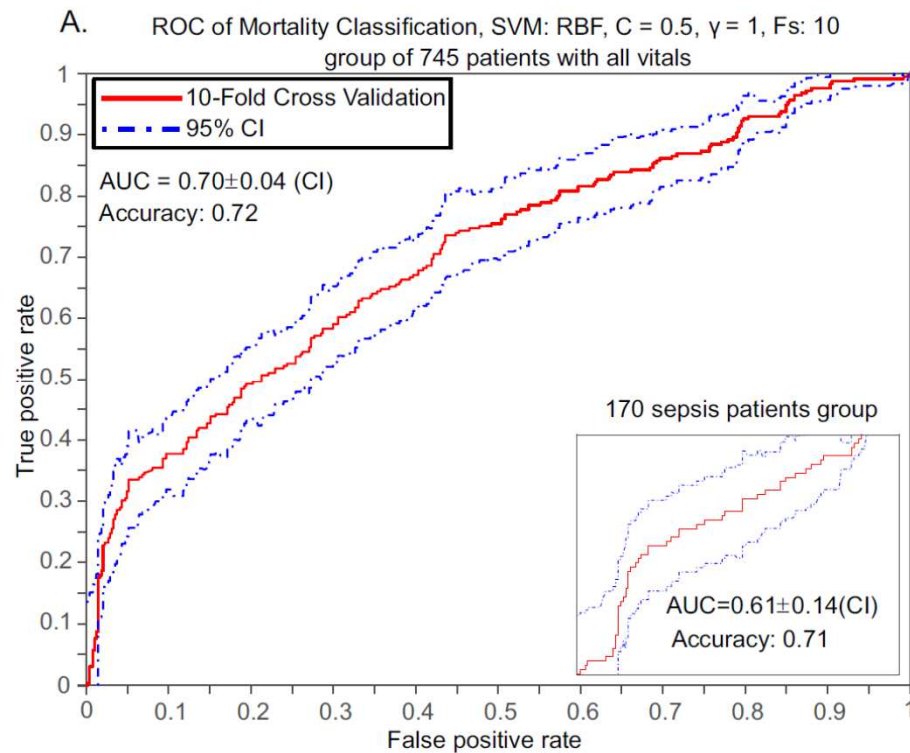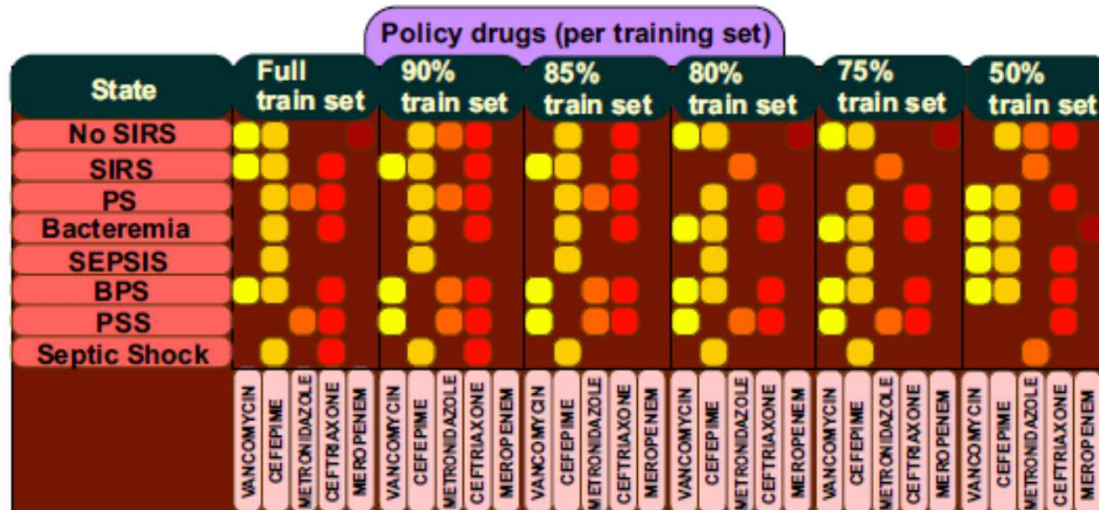
Gultepe E., Green J., Nguyen H., Adams J., Albertson T., **Tagkopoulos I.** (2014) From vital signs to clinical outcomes for patients with sepsis: A machine learning basis for a clinical decision support system. *Journal of American Medical Informatics Association (JAMIA)*, 21(2):315-25, doi:10.1136/amiajnl-2013-1815
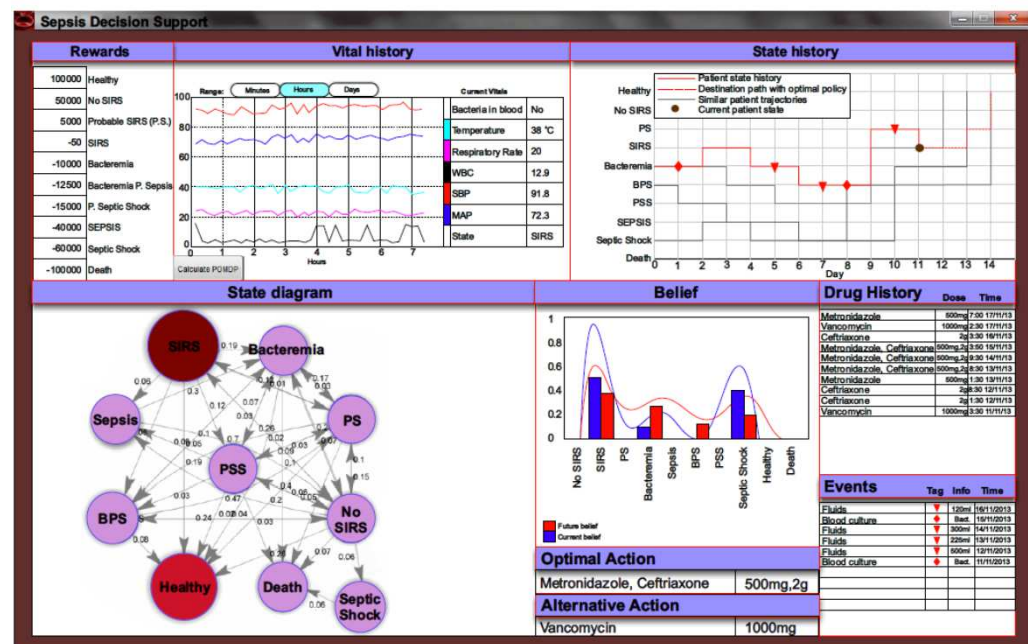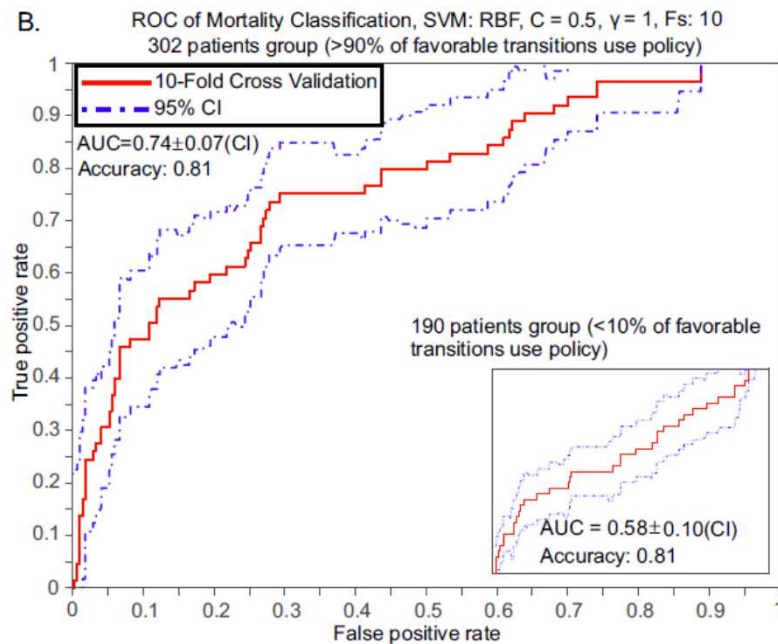
**Data-derived treatment led to better outcomes (47% vs. 36%)**

# 1498 crop fields in CA San Joaquin valley, 1997 to 2008
# Optimal pest management (*Lygus Hesperus*)

**Only ~82% of the farmers use the data-predicted optimal management**

|        | 5-11Jun | 12-18Jun | 19-25Jun | 26Jun-2Jul | 3-9Jul | 10-16Jul | 17-23Jul |
|--------|---------|----------|----------|------------|--------|----------|----------|
| Low    | 0.97    | 0.98     | 0.95     | 0.97       | 0.99   | 0.99     | 1.00     |
| Medium | 0.08    | 0.89     | 0.89     | 0.96       | 0.98   | 0.97     | 0.97     |
| High   | 0.18    | 0.67     | 0.74     | 0.83       | 0.88   | 0.90     | 0.87     |

**Cost from sub-optimal pest management per acre (yield-related cost only)**

cotton.

|        | 5-11Jun | 12-18Jun | 19-25Jun | 26Jun-2Jul | 3-9Jul  | 10-16Jul | 17-23Jul |
|--------|---------|----------|----------|------------|---------|----------|----------|
| Low    | $14.50  | $17.51   | $20.81   | $19.79     | $4.04   | $54.70   | $35.50   |
| Medium | $45.59  | $2.79    | $18.75   | $18.55     | $22.46  | $17.95   | $45.61   |
| High   | $33.81  | $2.64    | $20.06   | $12.96     | $23.89  | $19.99   | $42.99   |

# End of Lecture 14