

Brydon Eastman (b2eastma@uwaterloo.ca)  
Michelle Przedborski  
Mohammad Kohandel  
Template: Felix Breur, 2010

### Perturbed Virtual Patients

We prepared 3 sets of 200 virtual patients at 15%, 20%, and 25% perturbation strength level by perturbing the 5 nominal patient [1] parameter values by latin hypercube sampling a scaling factor up to the strength.



We treated these values as unknowns during the training and testing processes.

Only non-dimensionalizing the objective functional by the theoretical maximum for display purposes.

### Motivation

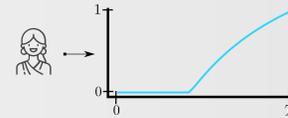
Personalized therapy often requires identifiability of hard to measure parameters.

### Research Question

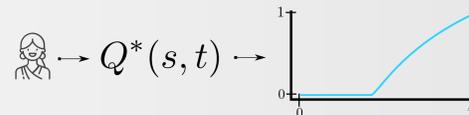
Can we personalize therapy without directly measuring these parameters?

Assuming we know the patient specific parameters, we can find the optimum via optimal control theory [2] or by reinforcement learning (or other methods...)

In optimal control theory, we derive the treatment directly



In RL, we approximate a value function from which we derive the treatment schedule [3]



We can define an optimal treatment by optimizing an objective functional

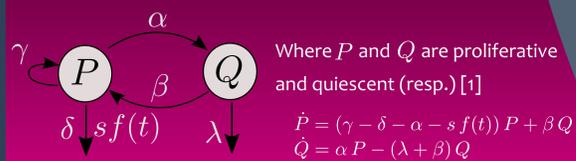
$$J(f) = \int_0^T \underbrace{P_{bm}(t)}_{\text{Preserve Healthy Cells}} + \underbrace{Q_{bm}(t)}_{\text{Deliver Drug}} - \frac{b}{2} (1 - f(t))^2 dt$$

# Reinforcement learning derived chemotherapeutic schedules for robust patient-specific therapy given unknown patient response parameters

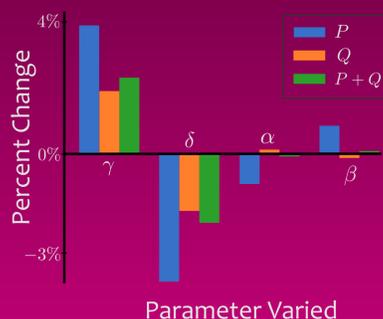
**Brydon Eastman**

University of Waterloo

### Tumor Growth Inhibition Model



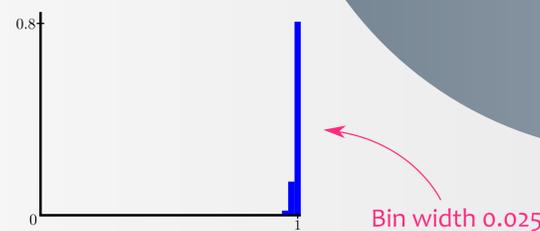
### Local Sensitivity Analysis



Perturbing parameters above the mean by 1% produces large changes in the predicted trajectories.

Sensitivity in the parameters suggests need for accurate identification of patient specific parameters.

The reinforcement learner produced scores clustered around the maximum.



While the optimal controller scores were much more diffuse.



As we perturbed further from the mean, the OC scores became more diffuse still while the RL scores remained clustered between 0.925 and 1.000

Suppose we only know some nominal patient parameters



Which we use to train our learner

$$Q^*(s, t) \equiv Q^*(s, t; \text{patient icon})$$

We want to choose the form of  $S$  such that

$$Q^*(s, t; \text{patient icon})$$

Produces schedules near the optimum for all virtual patients



Our form of  $S$  should

- > Be readily measurable
- > Depend on personal parameters

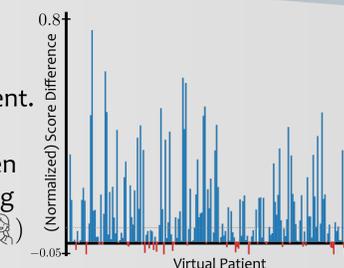
We use a window of daily relative bone marrow densities:

$$\frac{1}{P_{bm}(0) + Q_{bm}(0)} \begin{bmatrix} P_{bm}([t]) + Q_{bm}([t]) \\ P_{bm}([t-1]) + Q_{bm}([t-1]) \\ \vdots \\ P_{bm}([t-9]) + Q_{bm}([t-9]) \end{bmatrix}$$

Each bar represents a virtual patient.

The height is the difference between the RL score, obtained from applying the schedule learned from  $Q^*(s, t; \text{patient icon})$

and the OC score, obtained from applying the nominal control.



The RL derived schedules remain nearly optimal as perturbation strength increases

### Other Personalizations

We also considered OC based personalizations

These were obtained via various interpolation strategies

In all situations we solve the OC problem for 50 training patients and interpolate based on bone-marrow state and time

As in the other case, the RL remains robust as perturbation strength increases.

The difference between the RL and the OC scores can be as high as 0.35 or as low as -0.05. (Again where the swings favour the RL)

### (Some) Works Cited

1. John Carl Panetta. A mathematical model of breast and ovarian cancer treated with paclitaxel. *Mathematical biosciences*, 146(2):89–113, 1997
2. John Carl Panetta and K Renee Fister. Optimal control applied to cell-cycle-specific cancer chemotherapy. *SIAM Journal on Applied Mathematics*, 60(3):1059–1072, 2000.
3. Hado Van Hasselt, Arthur Guez, and David Silver. Deep reinforcement learning with double q-learning. In *Proceedings of the AAAI conference on artificial intelligence*, volume 30, 2016.
4. Brydon Eastman, Michelle Przedborski, and Mohammad Kohandel. Reinforcement learning derived chemotherapeutic schedules for robust patient-specific therapy. *bioRxiv*, 2021.