

Probabilistic Inference and Learning to Control: **pilco**

Carl Edward Rasmussen and Rowan McAllister

July 19, 2015

Abstract

The **pilco** toolbox is an Octave and Matlab implementation of time series inference and controller optimisation for nonlinear dynamical systems in continuous-state, discrete-time settings. The toolbox is both flexible and extensible. Flexibility is achieved by allowing user specification of which dynamics models, inference procedures and policy functional forms to use. Users may choose to provide their own data of a dynamical system, or generate synthetic data using any of the accompanying ‘scenarios’ such as the cart-pole system. Extensibility is a result of the modular code structure, separating class types such as controllers, dynamics models, scenarios, each of which are easily added to.

Contents

1	Overview	2
1.1	Toolbox Components	2
1.2	Program Flow	2
1.3	States	3
1.3.1	State Representation	3
1.3.2	State Variable Ordering	3
1.3.3	State Variable Classes	4
2	Base functions	4
2.1	Loss	5
2.2	Simulate	5
2.3	Propagate	8
2.4	Rollout	9
3	Controllers	11
3.1	Ctrl	11
3.2	CtrlNF	15
3.3	CtrlBF	18

1 Overview

There are several supporting structures, classes and functions the user should be aware of. Each is discussed in further details in later sections. Many functions will have derivative counterparts, named similarly except with an additional suffix ‘d’.

1.1 Toolbox Components

Scenario: Scenarios are computer simulators of dynamical systems. Each scenario represents a particular dynamical system - such as the cart-pole system - and generates synthetic data using the known laws of motion. To begin a simulation, execute the `doit.m` file within a scenario’s directory.

Dynamics Model: A dynamics model object is a probabilistic model which is used to make inferences from observed system evolutions and predictions on future dynamics. Each model is either a Gaussian process (GP), or GP state space model (GP-SSM).

Policy: A policy is a deterministic function, parameterised by θ , that maps state distribution s to action u ; $\pi : s \times \theta \rightarrow u$. Various functional forms of the policy are available, including linear, trigonometric, Gaussian radial basis functions (RBFs) and combinations thereof.

Controller: A controller is a wrapper object of the policy function. In addition, a controller computes various covariances, derivatives and can filter the sequence of system-state observations.

Cost: The cost function takes a state distribution s and returns a scalar value on the unit interval; $c : s \rightarrow [0, 1]$ The output represents the instantaneous cost (or measure of *undesirability*) per unit of time of the system being in state s . The cost specifies what it is the user wants to achieve, a 0 cost being the most desirable kind of state, and a 1 being the least desired.

1.2 Program Flow

Select policy function form, parameterised by (initially random) θ :

$$\pi : s \times \theta \rightarrow u. \quad (1)$$

FOR each trial n in $[1, 2, \dots, N]$:

- **Infer** dynamics given all observed data so far; $p(s^{t+1}|s^t, u^t, \mathcal{D}^{1:n-1})$.
- **Simulate** system from initial state distribution s_0 until final state distribution s_{horizon} , using the dynamics model to make successive predictions from one time step to the next.
- **Evaluate** controller:

$$f(s_0; \theta) = \sum_{t=0}^{\text{horizon}} \gamma^t \text{cost}(s_t), \quad (2)$$

where f is the cumulative discounted cost, θ is the policy parameters. Derivative information $\frac{df}{d\theta}$ is also computed.

- **Optimise** policy (via gradient decent):

$$\theta^* \leftarrow \arg \min_{\theta} f(s_0; \theta). \quad (3)$$

- **Apply** controller and generate more training data \mathcal{D}^n . This may be done using an available scenario, or the user’s hardware.

1.3 States

1.3.1 State Representation

The *state* of the control system can be expressed in various ways. The simplest representation is the *physical state*, consisting:

1. present position vector \mathbf{x}^t ,
2. present velocity vector $\dot{\mathbf{x}}^t$.

For the best results, we incorporate all relevant information such that satisfies the Markov assumption. One way the above Markov assumption is violated is with laggy sensors. Under sensor lag, observations reflect the state of the system in the past, during which time the previous action has still been in effect. Thus previous actions become relevant when inferring the current state. So a better state representation when mitigating the effects of sensor-lag would be:

1. previous action vector \mathbf{u}^{t-1} ,
2. present position vector \mathbf{x}^t ,
3. present velocity vector $\dot{\mathbf{x}}^t$.

The above representation is especially important when sensor lag amounts to a significant proportion of the time discretisation. Now consider the case where, in addition to sensor delay, velocity information $\dot{\mathbf{x}}$ is unavailable. In this case, the 2-Markov state representation is helpful, consisting:

1. previous-previous action vector \mathbf{u}^{t-2} ,
2. previous position vector \mathbf{x}^{t-1} ,
3. previous action vector \mathbf{u}^{t-1} ,
4. present position vector \mathbf{x}^t .

Intuitively, the above state representation works well because one could estimate any relevant velocity information using finite differences. All aforementioned state representations are possible in the **pilco** toolbox, including general N-Markov representations. Note the ordering of state variables must be chronological.

The *state distribution* is approximated as a multivariate Gaussian over all state variables. We encode the state distribution with a structure:

```
state % state structure
state.m % state mean vector
state.s % state covariance matrix
```

1.3.2 State Variable Ordering

Each scenario directory contains its own executable **doit.m** file, used to begin simulations. The **doit** file specifies the *ordering* of the state variables (amongst other things). For example, consider a hypothetical scenario: a 2-Markov cart-pole with an augmented variable. The state variables would be ordered as follows:

1	1	1	oou	even older value of u
2	2	2	ox	old cart position
3	3	3	otheta	old angle of the pendulum

4	4	4	ou	old value of u
5	5	5	x	cart position
6	6	6	theta	angle of the pendulum
7			v	cart velocity
8			dtheta	angular velocity
9			x - ox	relative change in cart position
10			u	force applied to cart
	7		sin(otheta)	sine of old angle of the pendulum
	8		cos(otheta)	cosine of old angle of the pendulum
	9		sin(theta)	sine of angle of the pendulum
	10		cos(theta)	cosine of angle of the pendulum

where the **green** column of numbers are the indices for each state variable in our state (here our state dimensionality is $D = 6$). The three types of indices always agree the first D numbers. The **red** column keeps extending past index D , and includes some extra variables that the simulator will use. The **blue** column, used by the policy, performs some trigonometric operations on some of the state variables, and thus also has an extends past D .

The **red** and **green** indices must refer to variables in a chronological order. The chronological order also takes the understanding that action variables at any time t occur *after* state variables at time t . E.g. index 10 is the final red index.

Typically a **doit** file will only display the **red** and **green** indices. The **blue** indices are specified implicitly as a function of the **green** indices and the set of angular variables discussed below.

1.3.3 State Variable Classes

Let us now group the above variables that we used in the above subsection. There exists three different ways of indexing each state variable in Sec. 1.3.2 (**red** **green** and **blue**) and the sets of indices below have been colour coded w.r.t. which columns of indices they refer to.

```

D    = 6;           % dimensionality of the state
E    = 2;           % number of outputs from the dynamics model
U    = 1;           % dimensionality of the action
angi = [3 6];       % indices for variables treated as angles (sin/cos)
augi = [9];         % indices for variables augmented to ode vars
dyni = [1 2 3 4 5 6]; % indices for input into the dynamics model
dyno = [5 6];       % indices of dynamics model output, and loss input
odei = [5 6 7 8];   % indices for the ode solver
poli = [1 2 4 5 7 8 9 10]; % indices for the inputs to the policy

```

The **dyni** and **dyno** variables are currently redundant, since it is always the case that **dyni** = 1: D , and **dyno** = D - E +1: D .

Note the [7 8 9 10] indices of **poli** refer to the trigonometric terms of whatever **angi** specified were angle variables.

2 Base functions

The **base** directory contains common files that are always used to train a controller. Evaluation of a controller's parameterisation is with **loss.m** (Sec. 2.1). The controller's parameterisation *loss* is a function of the cumulative cost distribution, provided by **simulate.m** (Sec. 2.2). An important

subroutine of `simulate.m` is the `propagate.m` function (Sec. 2.3). The `propagate.m` function computes predictive state distributions at time step $t + 1$ given the state distribution at time t current time.

2.1 Loss

```
5a <loss.m 5a>≡
1 function [f, df] = loss(p, s, dyn, ctrl, cost, H, exp, cc_prev, n)
2 % If the heuristic function 'explore' is not present (or if it is empty) then
3 % the expected cumulative cost is returned. If 'explore' is present, then the
4 % loss is a function of the cumulative cost mean and variance of the current
5 % parameterisation 'cc', and the parameterisation of the previous rollout
6 % 'cc_prev', and the number of rollout trials remaining 'n'. Derivatives of
7 % these quantities are computed when desired.
8 %
9 % [f, df] = loss(p, s, dyn, ctrl, cost, H, exp, cc_prev, n)
10 %
11 % p          policy parameter structure
12 % s          initial state structure
13 % dyn        dynamics model object
14 % ctrl       controlloer object
15 % cost       cost object
16 % H          time-steps horizon
17 % exp        exploration struct
18 % cc_prev    cumulative (discounted) cost structure of previous rollout
19 % n          number of trials remaining (inc. current point in time, so n>0)
20 % f          1x1 loss
21 % df         1xP loss derivative wrt policy parameters
22 %
23 % Copyright (C) 2015 by Carl Edward Rasmussen and Rowan McAllister 2015-06-01
24
25 ctrl.set_policy_p(p);
26 exploring = exist('explore','var') && ~isempty(explore);
27 sargs = {s, dyn, ctrl, cost, H};
28
29 if nargin == 1 % no derivatives
30     [S, A, ~, cc] = simulate(sargs{:});
31     if exploring
32         f = explore(exp, cc_prev, n, S, A, cc, [], sargs{:});
33     else
34         f = cc.m;
35     end
36 else % derivatives
37     [S, A, ~, cc, dcc] = simulate(sargs{:});
38     if exploring
39         [f, df] = explore(exp, cc_prev, n, S, A, cc, dcc, sargs{:});
40     else
41         f = cc.m;
42         df = dcc.m;
43     end
44 end
```

2.2 Simulate

The *cumulative cost* of a distribution over states is defined to be the long term (possibly discounted) total cost of starting in a given state distribution and following a policy π up to a time horizon

H. The simulate function returns states struct arrays, action struct arrays, cost struct arrays, the cumulative costs and possibly cumulative cost derivatives wrt the policy parameters.

5b $\langle \text{simulate.m 5b} \rangle \equiv$

```

1 function [s, a, c, cc, dcc] = simulate(s, dyn, ctrl, cost, H)
2
3 % [s, a, c, cc, dcc] = simulate(s, dyn, ctrl, cost, H)
4 %
5 % s          .      state structure
6 %   m        F*1    mean vector
7 %   s        F*F    covariance matrix
8 %   ?        possibly other fields representing additional information
9 % dyn        .      dynamics model object
10 % ctrl       .      controller object
11 %   is       .      struct indexing vectorized state distributions variables
12 %   m        F*1    indices of mean parameters
13 %   s        F*F    indices of variance parameters
14 %   np       .      number of parameters in the policy
15 %   ns       .      number of state distribution parameters (means and variances)
16 % policy    .      policy structure
17 %   p        .      policy parameters structure
18 % cost       .      cost function object
19 %   cov      @      function returning cost-covariances between two states
20 %   fcn      @      function returning expected and variance of a states's cost
21 %   gamma    .      discount factor
22 % H          .      length of prediction horizon
23 % s (output) H+1    state struct array containing Gaussian state distributions
24 %   m        F*1    mean vector
25 %   s        F*F    covariance matrix
26 % a          H      action struct array containing Gaussian action distributions
27 %   m        U*1    mean vector
28 %   s        U*U    covariance matrix
29 % c          H+1    cost struct array containing Gaussian cost distributions
30 %   m        .      mean scalar
31 %   s        .      variance scalar
32 % cc         .      cumulative (discounted) cost structure
33 %   m        .      mean scalar
34 %   s        .      variance scalar
35 % dcc        .      derivative structure of cc
36 %   m        1xP    derivative cc-mean wrt policy parameters, same fields as p
37 %   s        1xP    derivative cc-variance wrt policy parameters,same fields as p
38 %
39 % Copyright (C) 2008-2015 Carl Edward Rasmussen, 2015-05-31
40
41 global currT; currT=1;
42 if ~isfield(s,'s'); s.s = zeros(length(s.m)); end
43 c = cost.fcn(s); % init c
44 cc.m = c.m; cc.s = c.s; % init cc
45 q = s.s; % init q
46 gamma = cost.gamma; D = ctrl.D; F = length(s.m);
47 if nargin < 5 % no derivatives
48     for t = 1:H % iterate up to horizon
49         compute state action cost 7a
50         if nargin > 3 % then accumulate discounted cost
51             compute cumulative costs 7b
52         end
53     end
54 else % do derivatives

```

55 *<compute cumulative costs and derivatives 7c>*
56 **end**

The states, actions and cost distributions are computed by iterating up to the horizon, propagating forward the state distributions and recording the action distributions and computing the cost.

7a *<compute state action cost 7a>*≡ (5b)
1 `[s(t+1), C, a(t)] = propagate(s(t), dyn, ctrl);` % propagate state forward
2 `c(t+1) = cost.fcn(s(t+1));` % calc cost distribution

The discounted cumulative cost is computed by iteration up to the horizon, propagating the state distribution forward and accumulating the discounted cost

$$f^t = f^{t-1} + \gamma^t c^t, \quad t = 1, \dots, H, \quad \text{and} \quad f^{t=0} = 0, \quad (4)$$

where c^t is the instantaneous cost at time t and γ is the *discount* factor.

7b *<compute cumulative costs 7b>*≡ (5b)
1 `cc.m = cc.m + gamma^t*c(t+1).m;`
2 `cc.s = cc.s + gamma^(2*t)*c(t+1).s;`
3 *% cross-covariance terms:*
4 `q = C'*[q, s(t+1).s];`
5 **for** `j = 1:t`
6 `J = (j-1)*F+(1:D); % j'th block of q`
7 `cc.s = cc.s + 2*gamma^(j+t-1)*cost.cov(s(j), s(t+1), q(1:D,J)');`
8 **end**

To accumulative costs and derivatives we again iterate up to the horizon, but now we need to keep track of the derivatives using the chain rule, as we move forward in time. When calling **propagated.m** (the derivative counterpart function of *<propagate.m 8>*) it computes the state distribution at time t from the state distribution at time $t-1$ and derivatives wrt both the state and the policy parameters:

$$s^t = g(s^{t-1}, p), \quad \frac{\partial s^t}{\partial s^{t-1}} \quad \text{and} \quad \frac{\partial s^t}{\partial p}, \quad (5)$$

where g is the *transition* function, which takes a distribution over states and a policy and returns the distribution over the next state. The derivatives are propagated forward using the chain rule

$$\frac{ds^t}{dp} = \frac{\partial s^t}{\partial s^{t-1}} \frac{ds^{t-1}}{dp} + \frac{\partial s^t}{\partial p}, \quad \text{where} \quad \frac{ds^{t=0}}{dp} = 0, \quad (6)$$

which is iterated forward in time in line 11 below. Finally the derivative of the cumulative discounted cost is the discounted cumulative derivatives in line 12 below

$$\frac{df^t}{dp} = \frac{df^{t-1}}{dp} + \gamma^t \frac{dc_t}{ds^t} \frac{ds^t}{dp} \quad \text{where} \quad \frac{df^{t=0}}{dp} = 0 \quad (7)$$

where f^t is the cumulative discounted cost up to time t , c^t is the instantaneous cost at time t and γ is the discount factor.

7c *<compute cumulative costs and derivatives 7c>*≡ (5b)
1 `is = ctrl.is;`
2 `ic = unwrap([is.m(1:D), is.s(1:D,1:D)]);` % cost indices, depend on real vars
3 `sdp = cell(H+1); sdp{1} = zeros(ctrl.ns, ctrl.np);` % init derivatives
4 `qdp = sdp{1}(is.s,:);`
5 `dp = zeros(2, ctrl.np);` % first row = mean, second row = variance
6 **for** `t = 1:H` % iterate up to horizon
7 `[s(t+1), C, a(t), dsds, dsdp, dCds, dCdp] = propagated(s(t), dyn, ctrl);`
8 `[c(t+1), dcds] = cost.fcn(s(t+1));` % cost and derivative wrt state
9 `cc.m = cc.m + gamma^t*c(t+1).m;`

```

10 cc.s = cc.s + gamma^(2*t)*c(t+1).s;
11 sdp{t+1} = dsds*sdp{t} + dsdp; % chain rule
12 dp = dp + cost.gamma^t*(dcds*sdp{t+1}(ic,:)); % TODO verify gamma power
13 % cross-covariance terms:
14 dCdp = transposed(dCdp,C) + transposed(dCds,C)*sdp{t};
15 qdp = prodd([],dCdp,[q,s(t+1).s]) + prodd(C',[qdp;sdp{t+1}(is.s,:)]);
16 q = C'*[q, s(t+1).s];
17 for j = 1:t % cross terms
18     J = (j-1)*F+(1:D); % j'th column-block of q (real vars only)
19     [cov, dcovdsj, dcovdst, dcovdq] = cost.cov(s(j), s(t+1), q(1:D,J)');
20     cc.s = cc.s + 2*gamma^(j+t-1)*cov;
21     dcov = dcovdsj*sdp{j}(ic,:) + dcovdst*sdp{t+1}(ic,:) + ...
22         dcovdq*transposed(qdp(sub2ind2(F,1:D,J),:),D);
23     dp(2,:) = dp(2,:) + 2*gamma^(j+t-1)*dcov;
24 end
25 end
26 dcc.m = dp(1,:);
27 dcc.s = dp(2,:);

```

Two final subtleties consists in firstly, the derivatives of structures wrt to structures are represented as matrices (ie the structures have been vectorised), such that the products in line 11 and 12 are valid. This requires the structure **dcds** in line 12 to be vectorised (noting that **c** is a struct of two scalars). And secondly, note that in eq. (6) the state may contain fields in addition to **s.m** and **s.s**, but in eq. (7) only the fields **s.m** and **s.s** are considered (as the cost cannot depend directly on additional fields), and this is achieved by addressing directly those relevant indices using the predefined **ctrl.is** index structure in line 2.

2.3 Propagate

Propagate predicts the distribution over states at time $t+1$ given the distribution over states at time t . Predictions are made using a controller object to compute control actions, and then concatenating state and action information for input into a dynamics model.

8 $\langle propagate.m \ 8 \rangle \equiv$

```

1 function [s, C, a] = propagate(s, dyn, ctrl)
2
3 % Propagate the state distribution one time step forward.
4 %
5 % [s, C, a] = propagate(s, dyn, ctrl)
6 %
7 % s          .          state structure
8 %   m        F x 1      mean vector
9 %   s        F x F      covariance matrix
10 %   ?        possibly other fields representing additional information
11 % dyn        .          dynamics model object
12 %   D        dimension of the physical state
13 %   E        dimension of predictions from dyn model
14 %   pred     @          dynamics model function
15 %   pn       E x 1      log std dev process noise
16 % ctrl       .          controller object
17 %   fc       @          controller function
18 %   U        dimension of control actions
19 %   C        F x F      inverse input covariance times input-output covariance
20 %   a        .          action structure
21 %   m        U x 1      mean vector
22 %   s        U x U      covariance matrix
23 %

```



```

24 % Copyright (C) 2008-2015 Carl Edward Rasmussen and Rowan McAllister 2015-06-05
25
26 F = length(s.m); D = ctrl.D; E = dyn.E; U = ctrl.U;           % short hand names
27 Dz = F-D;                                                     % length of predicted information states
28 i = 1:D;                                                       % indices of physical state input variables
29 j = D+1:F;                                                     % indices of information state
30 k = F + (1:U);                                                 % indices of control actions
31 l = max(k) + (1:Dz);                                           % indices of predicted information state
32 m = max([k,l]) + (1:E);                                       % indices of predicted states
33 ij = [i j]; ik = [i k]; kl = [k l]; ijk1 = [ij kl];         % short hand
34 o = [ik(end-D+E+1:end) m l];                                  % ind. to select next state
35 M = zeros(max(m),1); M(ij) = s.m; S = zeros(max(m)); S(ij,ij) = s.s; % init
36
37 [M(kl), S(kl,kl), A, s] = ctrl.fcn(s, dyn);                   % control signal and inf state
38 q = S(ij,ij)*A; S(ij,kl) = q; S(kl,ij) = q';                 % action covariances
39
40 [M(m), S(m,m), B] = dyn.pred(M(ik), S(ik,ik));               % compute distr of next state
41 S(m,m) = S(m,m) + diag(exp(2*dyn.pn));                       % add process noise
42 q = S(ijk1,ik)*B; S(ijk1,m) = q; S(m,ijk1) = q';           % next state covariances
43
44 C = [eye(F) A [eye(F,D) A(:,1:U)]*B];                       % inv input var times cov
45 % C_exact_ZC = [eye(F,D), A(:,U+1:end)]*C_gph;
46
47 s.m = M(o); s.s = (S(o,o)+S(o,o'))/2; C = C(ij,o);          % select next state
48 a.m = M(k); a.s = (S(k,k)+S(k,k'))/2;                       % control action

```

On line 37 the controller outputs the matrix A . A which is the covariance between the input variable (the state, distributed as a Gaussian of mean $s.m$ and variance $s.s$) and the output variable (a concatenation of the action variable and the state filter prediction variables, distributed as a Gaussian of mean $M(kl)$ and variance $S(kl, kl)$) *pre-multiplied* by the inverse variance of the input variable. I.e.

$$A = \mathbb{V}[s^t]^{-1} C[s^t, \{u^t, s_{F-D+1:F}^{t+1}\}], \quad (8)$$

where s^t is the state at time t , u^t is the action taken at time t , and $s_{F-D+1:F}^{t+1}$ is the predicted filter state variables (whose indexes are $[F-D+1, F]$). On line 40 the dynamics model outputs the matrix B . Similarly, B is an input-output covariance matrix pre-multiplied by the inverse of input variance. We have

$$B = \mathbb{V}[\{s^t, u^t\}]^{-1} C[\{s^t, u^t\}, s^{t+1}], \quad (9)$$

where s^t is the state at time t and u^t is the action taken at time t resulting in state s^{t+1} at time $t+1$. The use of an pre-multiplied variance inverse of the input variable has the nice property that the covariance between one variable and any ancestor is simply computed with the product of the ancestor's variance and all such C terms along the path in between.

2.4 Rollout

We define a sampled trajectory of a system's possible evolution up to horizon H as a *rollout*. The system transitions from (point mass) state s^t to (point mass) state s^{t+1} as the system observes point observations y^t and applies point control signals u^t at each point in time t . If the state struct s does not have a variance field (i.e if $s.s$ does not exist), then subroutines will assume they are called by **rollout** opposed to **propagate**.

```

9 <rollout.m 9>≡
1 function [data, latent, L] = rollout(start, ctrl, H, plant, cost, verb)
2 % Compute a state trajectory using an ode solver (and any additional dynamics)
3 % from a particular starting state with either a particular policy or random
4 % actions.

```

```

5 %
6 % [data, latent, L] = rollout(start, ctrl, H, plant, cost, verb)
7 %
8 % start          nX x 1  vector containing start state (without controls)
9 % ctrl           controller structure
10 % fcn            @      function implementing the controller
11 % init           @      function initialising controller's filtered state
12 % policy         policy structure
13 % fcn            @      policy function
14 % p              parameter structure (if empty use random actions)
15 % maxU           nU x 1  vector of control input saturation values
16 % H              rollout horizon in steps
17 % plant          the dynamical system structure
18 % augi           (opt) indices for states passed to augment function
19 % augment        (opt) augment state using a known mapping
20 % constraint      (opt) stop rollout if violated
21 % dyno           indices for states passed to cost
22 % noise          observation noise
23 % odei           indices for states passed to the ode solver
24 % poli           indices for states passed to the policy
25 % cost           cost object
26 % verb           verbosity level
27 %
28 % data           data struct
29 % state          H+1xnX state matrix
30 % action         H x nU action matrix
31 % L              loss incurred at each timestep (1 by H)
32 % latent         matrix of latent states (H+1 by nX)
33 %
34 % Copyright (C) 2012-2015 Carl Edward Rasmussen and Rowan McAllister 2015-06-05
35
36 clear odestep;                                % clear persistent old action function
37 if isfield(plant,'augment'), augi = plant.augi; % sort out indices
38 else plant.augment = @(x)[]; augi = []; end
39 calc_loss = nargout > 2;
40 if nargin < 6; verb = 0; end
41 D = ctrl.D; E = ctrl.E; F = ctrl.F; U = ctrl.U;
42 N = length(start); odei = plant.odei;
43 latent = zeros(H+1, N); y = NaN(H+1, N); L = zeros(1, H+1); u = zeros(H, U);
44 obs_noise = @( ) (randn(1,E)*chol(plant.noise));
45
46 latent(1,:) = start;                            % initialise
47 y(1,1:D-E) = latent(1,1:D-E);
48 y(1,D-E+1:D) = latent(1,D-E+1:D) + obs_noise(); % add noise to observations
49 if calc_loss; L(1) = cost.fcn(struct('m',latent(1,1:D)')).m; end
50
51 s.m = y(1,1:D)'; s = ctrl.reset_filter(s);        % reset filter if exists
52
53 for i = 1:H % ----- run ROLLOUT
54 % Test constraints and stop rollout if violated -----
55 if isfield(plant,'constraint') && plant.constraint(latent(i,:))
56     H = i-1; if verb; disp('state constraints violated...'); end; break;
57 end
58
59 % Apply policy -----
60 s.m(1:D) = y(i,1:D)'; % receive an observation
61 [uzm,~,~,s] = ctrl.fcn(s);
62 u(i,:) = uzm(1:U); % action 'u' component of uzm

```

```

63 s.m(D+1:F) = uzm(U+1:end); % predicted filter 'zm' component of uzm
64
65 latent(i+1,1:D-E) = [latent(i,U+E+1:D), u(i,:)];
66 latent(i+1,odei) = odestep(latent(i,odei), u(i,:), plant);
67
68 y(i+1,1:D-E) = [y(i,U+E+1:D), u(i,:)];
69 y(i+1,D-E+1:D) = latent(i+1,D-E+1:D) +obs_noise(); % TODO: add process noise?
70
71 % Compute Cost -----
72 if calc_loss; L(i+1) = cost.fcn(struct('m',latent(i+1,1:D'))).m; end
73 end
74 if verb; disp(['Trial lasted ',num2str(floor(H)),' steps']); end
75
76 data.state = y(1:H+1,:); data.action = u(1:H,:);
77 latent = latent(1:H+1,:); L = L(1,1:H+1); % trim any trailing zeros
78
79 function xa = augment(x, plant)
80 xa(plant.odei) = x;
81 xa(plant.augi) = plant.augment(xa);

```

3 Controllers

Controller classes act as wrappers to policy functions. The superclass `Ctrl.m` defines the common functions a controller object supports, discussed Sec. 3.1. Then follows two sections on child classes that inherit `Ctrl.m`; the ‘No Filter’ controller `CtrlNF.m` in Sec. 3.2 and the ‘Bayes Filter’ controller `CtrlBF.m` Sec. 3.3.

3.1 Ctrl

```

11 <Ctrl.m 11>≡
1 classdef Ctrl < handle
2 %CTRL, controller superclass, which ctrlBF and ctrlNF inherit.
3 %
4 % Ctrl Properties:
5 % actuate - @ (optional) call this function with the action
6 % angi - indices for vars treated as angles (sin/cos rep)
7 % D - number of real-world state variables
8 % dyn - . dynamics model object (only used by CTRLBF class)
9 % E - number of physical state variables
10 % F - number of state variables (real + info vars)
11 % is - . state index structure
12 % np - number of policy parameters
13 % ns - state parameters (including filter)
14 % on - D x D non-log observation noise
15 % onp - D x D non-log obs. noise (only last E vars non zeros)
16 % poli - indices for the inputs to the policy
17 % policy - policy struct
18 % U - number of control outputs
19 %
20 % Ctrl Methods:
21 % build_state_index - (private) builds 'is' field to index state-struct
22 % clear_filter - clears all filter variables (if they exist)
23 % Ctrl - constructor
24 % fcn - main function, computes control signal
25 % random_action - outputs a random control, independent of state
26 % reset_filter - resets state fields {zs,zc,v} to a broad prior

```

```

27 % set_dynmodel      - sets the dynamics model object
28 % set_on            - sets the N-Markov observations noise
29 % set_policy_opt     - sets the policy optimisation settings
30 % set_policy_p       - sets the policy parameters
31 %
32 % See also CTRLBF.M, CTRLNF.M.
33 % Copyright (C) 2015 by Carl Edward Rasmussen and Rowan McAllister 2015-06-05
34
35 properties (SetAccess = private)
36     actuate
37     angi
38     D
39     dyn
40     E
41     F
42     is
43     np
44     ns
45     on
46     onp
47     poli
48     policy
49     U
50 end
51
52 methods
53
54     % Constructor
55     function self = Ctrl(D, E, policy, angi, poli, actuate)
56         % CTRL is the super-class controller constructor
57         %
58         % ctrl = Ctrl(D, E, policy, angi, poli, actuate)
59         %
60         % INPUTS:
61         %   D           number of state variables
62         %               (or a ctrl object to be copied)
63         %   E           number of predicted state variables
64         %   policy      .   policy struct
65         %   fcn         @   policy function
66         %   maxU        maximum control output magnitudes
67         %   opt         .   optimisation structure
68         %   fh          figure handle for minimize() to display to
69         %   length      how many optimisations steps
70         %   method
71         %   verbosity
72         %   angi        indices for variables treated as angles
73         %   poli        indices for the inputs to the policy
74         %   actuate     @   function to actuate calculated action
75
76         % Special case input:
77         % Another controller object might be the first (and only) input.
78         % This allows easy translations from one controller type to the next
79         if ~isnumeric(D); assert(nargin < 2); ctrl = D;
80             D = ctrl.D;
81             E = ctrl.E;
82             policy = ctrl.policy;
83             angi = ctrl.angi;
84             poli = ctrl.poli;

```

```

85     actuate = ctrl.actuate;
86     self.on = ctrl.on;
87     self.onp = ctrl.onp;
88     self.dyn = ctrl.dyn;
89 end
90
91 self.D = D;
92 self.E = E;
93 self.policy = policy;
94 if ~isfield(policy,'opt'); self.policy.opt = ...
95     struct('length',-1000,'method','BFGS','MFEPLS',20,'verbosity',3); end
96 if ~isfield(self.policy.opt,'fh'); self.policy.opt.fh = 1; end
97 self.U = length(policy.maxU);
98 self.build_state_index();
99 if ~isfield(policy,'type'); policy.type = ''; end
100 if isfield(self.policy,'p');
101     self.np = length(unwrap(self.policy.p));
102 end
103
104 if exist('angi','var'); self.angi = angi; else self.angi = []; end
105 if exist('poli','var'); self.poli = poli; else self.poli = 1:self.D; end
106 if exist('actuate','var') && ~isempty(actuate); self.actuate=actuate; end
107 end
108
109 function [uM,uS,uC,s] = fcn(self, s)
110 % CTRL.FCN is the main function to output control. Sub-classes will
111 % override this function.
112 %
113 % [uM, uS, uC, s, duMds, duSds, duCds, dsds, duMdp, duSdp, duCdp, ...
114 %   dsdp] = CTRL.FCN(s,propdyn)
115 %
116 % self      .      controller structure
117 % actuate   @      function to actuate calculated action
118 % angi      .      indices of angular variabels
119 % dyn       @      controller's dynamics model (for CtrlBF only)
120 % E         .      number of predictive state variables
121 % on        D x D   observation noise
122 % onp       D x D   observation noise (only last E vars non-zero)
123 % poli      .      indices of policy input
124 % policy    .      policy structure
125 % fcn       @      policy function
126 % U         .      number of control outputs
127 % s         .      state structure
128 % m         F x 1   state mean
129 % s         F x F   state variance
130 % v         (F-D) x (F-D) filter variance
131 % propdyn   @      propagates's dynmodel (for CtrlBF)
132 % M         (U+D) x 1 control signal mean vector
133 % S         (U+D) x (U+D) control signal variance matrix
134 % C         F x (U+D) input-output covariance matrix
135 % dMds      (U+D) x S derivatives of outputs wrt input state struct
136 % dSds      (U+D)^2 x S
137 % dCds      F*(U+D) x S
138 % dsds      S x S   ouput state derivative wrt input state
139 % dMdp      (U+D) x P P is the total number of parameters is the policy
140 % dSdp      (U+D)^2 x P
141 % dCdp      F*(U+D) x P
142 % dsdp      S x P   ouput state derivative wrt policy parameters

```

```

143     %
144     % See also CTRLBF.FCN, CTRLNF.FCN.
145     if strcmp(self.policy.type, 'random')
146         uM = self.random_action();
147         uS = zeros(self.U); uC = zeros(self.D,self.U);
148     end
149 end
150
151 % Filter Clearer
152 function s = clear_filter(self, s)
153     % Clears all filter variables (if they exist).
154     % s = CTRL.CLEAR_FILTER(s)
155     % s: state struct
156     i = 1:self.D;
157     s.m = s.m(i);
158     if isfield(s,'s'); s.s = s.s(i,i); end
159     if isfield(s,'v'); s = rmfield(s,'v'); end
160 end
161
162 function u = random_action(self)
163     % u = CTRL.RANDOM_ACTION(), outputs a random action
164     % u: U x 1 control action
165     u = self.policy.maxU.*(2*rand(1,self.U)-1);
166 end
167
168 function s = reset_filter(~, s)
169     % s = CTRL.RESET_FILTER(s), does nothing unless overwritten
170     % s: state struct
171     % See also CTRLBF.RESET_FILTER
172 end
173
174 function set_dynmodel(self, dyn)
175     % CTRL.SET_DYNMODEL(dyn), for updating controller's dynmodel
176     % dyn: dynamics model object
177     % Note: dyn property only ever used by CTRLBF
178     self.dyn = dyn;
179 end
180
181 function set_on(self, onE)
182     % CTRL.SET_ON(onE), sets observation noise
183     % onE: E x 1, log observation noise (physical state-variables only)
184     assert(self.E == length(onE));
185     onD = [-inf(self.U,1); onE(:)];
186     onD = repmat(onD,ceil(self.D/(self.U+self.E)),1);
187     onD = onD(end-self.D+1:end);
188     self.on = diag(exp(2*onD));
189     p = self.D-self.E+1:self.D;
190     self.onp = 0*self.on; self.onp(p,p) = self.on(p,p);
191 end
192
193 function set_policy_p(self, p)
194     % CTRL.SET_POLICY_P(p), updates the policy parameters
195     % p: policy parameter struct
196     self.policy.p = p;
197     self.np = length(unwrap(p));
198 end
199
200 function set_policy_opt(self, opt)

```

```

201     % CTRL.SET_POLICY_OPT(p), updates the policy optimisation settings
202     %   opt:   policy optimisation settings
203     self.policy.opt = opt;
204 end
205
206 end
207
208 methods (Access = private)
209     (function-build-state-index 15)
210 end
211
212 end

```

PILCO stores derivatives of any object as matrices, where rows represent vectorised dependent variables, and columns represent vectorised independent variables. To vectorise any object we use the `unwrap.m` function. The state distribution is stored as a structure. Thus when handling state derivatives, we need to keep track of which elements in the native structure correspond to which index in the vectorisation of that structure. We do so by using a *state index* structure, which has the exact same form of the state structure, except each element within the index structure is a unique integer, corresponding to that element's index when in vectorised form. The state index is built using *(function-build-state-index 15)*:

```

15 (function-build-state-index 15)≡ (11)
1 function build_state_index(self)
2 % CTRL.BUILD_STATE_INDEX(), computes internal field 'is' - a state struct
3 % whose members are a states members' indexes. Requires subclass to
4 % implement function CTRL.RESET_FILTER
5 s.m = nan(self.D,1);
6 s = self.reset_filter(s); % generates possible filter variables
7 s.s = nan(length(s.m));
8 if isfield(s,'reset'); s = rmfield(s,'reset'); end
9 self.ns = length(unwrap(s));
10 self.is = rewrap(s,1:self.ns);
11 self.F = length(s.m);
12 end

```

3.2 CtrlNF

The most straightforward controller class is one which inputs (noisy) signals from sensors directly into the policy function. The 'No Filter' controller class `control/CtrlNF` implements exactly this.

The process of controlling a dynamical system using `CtrlNF.m` is depicted in Fig. 1. It is assumed the latent system state X_t is observed by imperfect sensors as Y_t , where the random variable Y_t is equal to X_t plus some Gaussian noise:

$$Y_t \sim X_t + \epsilon_t, \quad (10)$$

$$\epsilon_t \stackrel{iid}{\sim} \mathcal{N}(0, \Sigma_\epsilon). \quad (11)$$

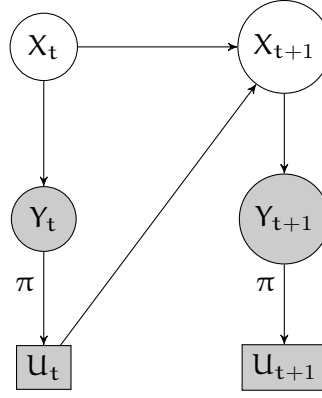


Figure 1: Directed graphical model of a system controller with No Filter from timestep t to timestep $t + 1$. A dynamical system begins in state X_t , which sensors observe noisily as Y_t . The controller uses Y_t to decide control signal U_t via policy function π . Finally, the control signal U_t is applied to the system, resulting in new state X_{t+1} .

16a $\langle CtrlNF.m$ 16a) \equiv

```

1 classdef CtrlNF < Ctrl & handle
2
3 % Controller with No Filter. The state is given either as a point s.m or as a
4 % distribution N(s.m,s.s). First augment with trignometric functions if
5 % necessary, then call the policy, and finally (optionally) call the actuate
6 % function. There is no filter, so no updates to any state structure filter
7 % fields are required.
8 %
9 % See also CTRL.M, CTRLNFT.M.
10 % Copyright (C) 2015 by Carl Edward Rasmussen and Rowan McAllister 2015-06-02
11
12 methods
13
14 % Constructor
15 function self = CtrlNF(varargin)
16     % CTRLNF.CTRLNF is the sub-class constructor of CTRL.CTRL
17     % For help, see also Ctrl.Ctrl
18     self@Ctrl(varargin{:}); % Super Ctrl constructor
19 end
20
21 % Main function
22 function [uM, uS, uC, s, duMds, duSds, duCds, dsds, ...
23     duMdp, duSdp, duCdp, dsdp] = fcn(self, s, ~)
24     % For help, see also CTRL.FCN
25
26     <initialise 16b>
27     <augment state with trig variables 17>
28     <compute control signal 18>
29 end
30
31 end
32
33 end

```

16b $\langle initialise$ 16b) \equiv (16a)

```

1 % CtrlNF only operates on a noisy version of the state:
2 D = self.D; DD = D*D;

```



```

3 s = self.clear_filter(s);           % clear any filter variables in state
4 assert(length(s.m) == D);
5 if isfield(s,'s')
6     sy = s.s + self.on;             % propagate mode (distribution of states)
7 else
8     sy = zeros(D);                  % rollout mode (point-mass state sample)
9 end
10
11 angi = self.angi; poli = self.poli; A = length(angi);
12 derivativesRequested = nargout > 4;
13 ns = self.ns; is = self.is;
14 D1 = D + 2*A;
15 i=1:D;
16 M = zeros(D1,1); M(i) = s.m; S = zeros(D1); S(i,i) = sy;
17 if derivativesRequested
18     idx = @(i,j,I) bsxfun(@plus, I*(i'-1), j);
19     Mds = zeros(D1,ns); Mds(i,is.m) = eye(D);
20     Sds = zeros(D1*D1,ns); Sds(:,is.s) = kron(Mds(:,is.m),Mds(:,is.m));
21     dsds = zeros(self.ns); dsdp = zeros(ns,self.np);
22     dsds(is.m,is.m) = eye(D); dsds(is.s(:),is.s(:)) = eye(D*D);
23     dsds(is.s,is.s) = symmetrised(dsds(is.s,is.s),[1,2]);
24 end

```

The state distribution \mathbf{s} , is then extended the sine and cosine each angle variable (indexed by \mathbf{angi}) using `gTrig.m`. The `gTrig.m` function also outputs covariance information \mathbf{C}_{gtrig} , which relates to the input-output covariance except with pre-multiplicative input variance:

$$\mathbb{C}[\mathbf{Y}_t, \mathbf{Y}_t^\circ] = \mathbb{V}[\mathbf{Y}_t] \mathbf{C}_{gtrig}[\mathbf{Y}_t, \mathbf{Y}_t^\circ] \quad (12)$$

17 *\langle augment state with trig variables 17 $\rangle \equiv$* (16a)

```

1 % augment state with trig functions
2 i = 1:D; k = D+1:D1;
3 if ~derivativesRequested
4     [M(k), S(k,k), cg] = gTrig(M(i), S(i,i), angi);
5 else
6     kk = idx(k,k,D1); ik = idx(i,k,D1); ki = idx(k,i,D1);
7     [M(k), S(k,k), cg, Mds(k,is.m), Sds(kk,is.m), cgdm, ...
8         Mds(k,is.s), Sds(kk,is.s), cgds] = gTrig(M(i), S(i,i), angi);
9     qdm = prodd(S(i,i),cgdm);
10    Sds(ik,is.m) = qdm; Sds(ki',is.m) = qdm;
11    qds = prodd(S(i,i),cgds) + prodd([], 'eye', cg);
12    Sds(ik,is.s) = qds; Sds(ki',is.s) = qds;
13 end
14 q = S(i,i)*cg; S(i,k) = q; S(k,i) = q';

```

The policy, which may be a function of both the regular state variables, and augmented state variables is then computed:

$$\mathbf{U}_t \leftarrow \pi(\{\mathbf{Y}_t, \mathbf{Y}_t^\circ\}) \quad (13)$$

The non-trivial part is computing output \mathbf{uC} which is covariance between state input and control output $\mathbb{C}[\mathbf{X}_t, \mathbf{U}_t]$. We begin noting:

$$\mathbb{C}[\{\mathbf{Y}_t, \mathbf{Y}_t^\circ\}, \mathbf{U}_t] = \mathbb{V}[\{\mathbf{Y}_t, \mathbf{Y}_t^\circ\}] \mathbf{C}_{\text{poli}}[\{\mathbf{Y}_t, \mathbf{Y}_t^\circ\}, \mathbf{U}_t] \quad (14)$$

$$\begin{bmatrix} \mathbb{C}[\mathbf{Y}_t, \mathbf{U}_t] \\ \mathbb{C}[\mathbf{Y}_t^\circ, \mathbf{U}_t] \end{bmatrix} = \begin{bmatrix} \mathbb{V}[\mathbf{Y}_t] & \mathbb{C}[\mathbf{Y}_t, \mathbf{Y}_t^\circ] \\ \mathbb{C}[\mathbf{Y}_t^\circ, \mathbf{Y}_t] & \mathbb{V}[\mathbf{Y}_t^\circ] \end{bmatrix} \mathbf{C}_{\text{poli}}[\{\mathbf{Y}_t, \mathbf{Y}_t^\circ\}, \mathbf{U}_t] \quad (15)$$

$$\therefore \mathbb{C}[\mathbf{Y}_t, \mathbf{U}_t] = \left[\mathbb{V}[\mathbf{Y}_t], \mathbb{C}[\mathbf{Y}_t, \mathbf{Y}_t^\circ] \right] \mathbf{C}_{\text{poli}}[\{\mathbf{Y}_t, \mathbf{Y}_t^\circ\}, \mathbf{U}_t] \quad (16)$$

$$= \mathbb{V}[\mathbf{Y}_t][\mathbf{I}, \mathbf{C}_{\text{gtrig}}[\mathbf{Y}_t, \mathbf{Y}_t^\circ]] \mathbf{C}_{\text{poli}}[\{\mathbf{Y}_t, \mathbf{Y}_t^\circ\}, \mathbf{U}_t] \quad (17)$$

Thus,

$$\mathbf{uC} = \mathbb{C}[\mathbf{X}_t, \mathbf{U}_t] = \mathbb{C}[\mathbf{X}_t, \mathbf{Y}_t] \mathbb{V}[\mathbf{Y}_t]^{-1} \mathbb{C}[\mathbf{Y}_t, \mathbf{U}_t] \quad (18)$$

$$= \mathbb{V}[\mathbf{X}_t][\mathbf{I}, \mathbf{C}_{\text{gtrig}}[\mathbf{Y}_t, \mathbf{Y}_t^\circ]] \mathbf{C}_{\text{poli}}[\{\mathbf{Y}_t, \mathbf{Y}_t^\circ\}, \mathbf{U}_t] \quad (19)$$

```

18  <compute control signal 18>≡ (16a)
1  % compute control signal
2  if ~derivativesRequested
3      [uM, uS, uC] = self.policy.fcn(self.policy, M(poli), S(poli,poli));
4  else
5      [uM, uS, uC, mdm, sdm, cdm, mds, sds, cds, duMdp, duSdp, duCdp] = ...
6          self.policy.fcn(self.policy, M(poli), S(poli,poli));
7  end
8  if isfield(self, 'actuate'), self.actuate(uM); end % actuate controller
9
10 ec = [eye(D) cg]; ecp = ec(:,poli);
11 if derivativesRequested
12     poli2 = idx(poli,poli,D1); ii = sub2ind2(D1,i,i);
13     duMds = mdm*Mds(poli,:) + mds*Sds(poli2,:);
14     duSds = sdm*Mds(poli,:) + sds*Sds(poli2,:);
15
16     duC = cdm*Mds(poli,:) + cds*Sds(poli2,:);
17     decp = [zeros(DD,ns); cgdm*Mds(i,:) + cgds*Sds(ii,:)];
18     decp = decp(sub2ind2(D,1:D,poli),:);
19     duCds = prodd(ecp,duC) + prodd([],decp,uC);
20
21     duCdp = prodd(ecp,duCdp);
22
23     duCds(:,is.s) = symmetrised(duCds(:,is.s),2);
24     duMds(:,is.s) = symmetrised(duMds(:,is.s),2);
25     duSds(:,is.s) = symmetrised(duSds(:,is.s),2);
26 end
27 uC = ecp*uC;

```

3.3 CtrlBF

CtrlBF.m implements control using Bayesian filtering. The *state* of the control system using filtering contains two parts:

1. the *latent state* \mathbf{x} ,

2. the *filter state* distribution $\mathbf{b} \sim \mathcal{N}(\mathbf{z}, \mathbf{V})$.

Thus, the *state distribution* is in principle a distribution over the random variables, \mathbf{x} , \mathbf{z} and \mathbf{V} . However, as an approximation we are going to assume that the distribution on the variance \mathbf{V} is just a delta function (ie, that the variance is some fixed value). Assuming further that the state distribution is Gaussian

$$\begin{bmatrix} \mathbf{x} \\ \mathbf{z} \end{bmatrix} \sim \mathcal{N}\left(\begin{bmatrix} \mathbf{m}_x \\ \mathbf{m}_z \end{bmatrix}, \begin{bmatrix} \Sigma_x & \Sigma_{xz} \\ \Sigma_{zx} & \Sigma_z \end{bmatrix}\right). \quad (20)$$

A Bayes filter (BF) maintains a belief-posterior on \mathbf{x} , denoted $\mathbf{B}_{t:t}$, conditioned on all information available thus far: the entire history of the system observations $\mathbf{y}_{1:t}$ and applied control signals $\mathbf{u}_{1:t-1}$. Conditioning on more information than the current observation \mathbf{y}_t yields a more informed (and smooth) estimate of \mathbf{x}_t . Being a function of all observations, $\mathbf{B}_{t:t}$ is less susceptible to the noise injected into the most recent observation \mathbf{y}_t , and consequently the controller's input is much smoother. To maintain \mathbf{B}_{t+1} the BF makes two recursive updates per timestep:

1. Update step: Compute $\mathbf{B}_{t:t}$ using prior belief $\mathbf{B}_t = \mathbf{p}(\mathbf{x}_t)$ and observation likelihood $\mathcal{L}(\mathbf{x}_t|\mathbf{y}_t) = \mathbf{p}(\mathbf{y}_t|\mathbf{x}_t)$,
2. Predict step: Compute \mathbf{B}_{t+1} by mapping updated belief $\mathbf{B}_{t:t}$ through transition model $\mathbf{p}(\mathbf{x}_{t+1}|\mathbf{b}_{t:t}, \mathbf{u}_t)$.

A directed graphical model of a Bayes filter is shown Fig. 3.3:

