

# Contents

<b>Notation</b>	<b>3</b>
<b>Introduction</b>	<b>4</b>
<b>1 Basics of Recursive Bayesian Estimation</b>	<b>5</b>
1.1 Problem Statement . . . . .	5
1.2 Theoretical solution . . . . .	6
1.3 Kalman Filter . . . . .	7
1.4 Particle Filter . . . . .	9
1.5 Marginalized Particle Filter . . . . .	11
<b>2 Software analysis</b>	<b>12</b>
2.1 Requirements . . . . .	12
2.2 Programming paradigms . . . . .	13
2.2.1 Procedural paradigm . . . . .	14
2.2.2 Object-oriented paradigm . . . . .	14
2.2.3 Functional paradigm . . . . .	15
2.2.4 Other programming language considerations . . . . .	16
2.3 C++ . . . . .	19
2.4 Matlab . . . . .	20
2.5 Python . . . . .	20
2.6 Cython . . . . .	20
2.6.1 Gradual Optimisation . . . . .	20
2.6.2 Parallelisation . . . . .	20
2.6.3 Pure Python mode . . . . .	20
2.6.4 Limitations . . . . .	21
2.6.5 Performance comparison with C and Python . . . . .	21
2.7 Survey of Existing Libraries for Bayesian estimation/decision making	21

2.8	Choice . . . . .	21
<b>3</b>	<b>The PyBayes Library</b>	<b>22</b>
3.1	Interpreted and Compiled . . . . .	22
3.2	Library Layout . . . . .	22
3.2.1	Random Variable Meta-representation . . . . .	22
3.2.2	Probability Density Functions . . . . .	22
3.2.3	Bayesian Filters . . . . .	22
3.3	Documentation, Testing and Profiling . . . . .	23
3.4	Comparison with BDM . . . . .	23
	<b>Conclusion</b>	<b>24</b>
	<b>Bibliography</b>	<b>25</b>

# Notation

Throughout this text, following notation is used

$\mathbb{N}$	set of natural numbers excluding zero
$\mathbb{N}_0$	set of natural numbers including zero
$\mathbb{R}$	set of real numbers
$t$	discrete time moment; $t \in \mathbb{N}_0$
$a_t$	value of quantity $a$ at time $t$ ; $a_t \in \mathbb{R}^n, n \in \mathbb{N}$
$a_{t t'}$	quantity with two indices: $t$ and $t'$ there is no implicit link between $a_t$ , $a_{t t'}$ and $a_{t'}$
$a_{t:t'}$	sequence of quantities $(a_t, a_{t+1} \dots a_{t'-1}, a_{t'})$
$p(a_t)$	probability density function <sup>1</sup> of quantity $a$ at time $t$ (unless noted otherwise)
$p(a_t b_{t'})$	conditional probability density function of quantity $a$ at time $t$ given value of quantity $b$ at time $t'$
$\delta(a)$	Dirac delta function; used exclusively in context of probability density functions to denote discrete distribution within framework of continuous distributions <sup>2</sup>
$\mathcal{N}(\mu, \Sigma)$	multivariate normal (Gaussian) probability density function with mean vector $\mu$ and covariance matrix $\Sigma$

---

<sup>1</sup>for the purpose of this text, probability density function  $p$  is multivariate non-negative function  $\mathbb{R}^n \rightarrow \mathbb{R}$ ;  $\int_{\text{supp } p} p(x_1, x_2 \dots x_n) dx_1 dx_2 \dots dx_n = 1$

<sup>2</sup>so that  $\int_{-\infty}^{\infty} f(x) \delta(x - \mu) dx = f(\mu)$  and more complex expressions can be derived using integral linearity and Fubini's theorem.

# Introduction

TODO motivatin for bayes filtration + a need for a convenient library (rapid prototyping vs. speed)

applications: robotics, navigation, + tracking of toxic plume after radiation accident.

Decision-making being a logical and natural “next step” - beyond the scope of this text.

[proposed citations:[[14](#), [5](#), [7](#), [6](#), [9](#)]]

# Chapter 1

## Basics of Recursive Bayesian Estimation

In following sections the problem of recursive Bayesian estimation is stated and its analytical solution is derived. Later on, due to practical intractability of the solution in its general form, a few methods that either simplify the problem or approximate the solution are shown.

### 1.1 Problem Statement

Assume a dynamic system described by a hidden real-valued *state vector*  $x$  which evolves at discrete time steps according to a known function  $f_t$  (in this text called *process model*) as described by (1.1).

$$x_t = f_t(x_{t-1}, v_{t-1}) \quad (1.1)$$

Variable  $v_t$  in (1.1) denotes random *process noise*, which may come from various sources and is often inevitable. Sequence of  $v_t$  is assumed to be identically independently distributed random variable sequence.

The state of the system is hidden and can only be observed through a real-valued *observation vector*  $y$  that relates to the state  $x$  as in (1.2), but adds further *observation noise*  $w$ .

$$y_t = h_t(x_t, w_t) \quad (1.2)$$

In (1.2)  $h_t$  is known function called *observation model* in this text and  $w_t$  is identically independently distributed random variable sequence that denotes observation noise.

The goal of recursive<sup>1</sup> Bayesian estimation is to give an estimate of the state  $x_t$

---

<sup>1</sup>by the word recursive we mean that it is not needed to keep track of the whole batch of previous observations in practical methods, only appropriate quantities from time moments  $t - 1$  and  $t$  are needed to estimate  $x_t$ . However, this does not apply to the derivation of the solution, where the notation of whole batch of observations  $y_{1:t}$  is used.

given the observations  $y_{1:t}$  provided the knowledge of the functions  $f_t$  and  $h_t$ . More formally, the goal is to find the probability density function  $p(x_t|y_{1:t})$ . Theoretical solution to this problem is known and is presented in next section.

## 1.2 Theoretical solution

At first, we observe that probability density function  $p(x_t|x_{t-1})$  can be derived from the process model (1.1) (given the distribution of  $v_k$ ) and that  $p(y_t|x_t)$  can be derived from the observation model (1.2) respectively. (given the distribution of  $w_k$ )

Because recursive solution is requested, suppose that  $p(x_{t-1}|y_{1:t-1})$  and  $p(x_0)$  are known<sup>2</sup> in order to be able to make the transition  $t-1 \rightarrow t$ .

In the first stage that can be called *prediction*, *a priori* probability density function  $p(x_t|y_{1:t-1})$  is calculated without knowledge of  $y_t$ . We begin the derivation by performing the reverse of the marginalization over  $x_{k-1}$ .

$$p(x_t|y_{1:t-1}) = \int_{-\infty}^{\infty} p(x_t, x_{t-1}|y_{1:t-1}) dx_{t-1}$$

Using chain rule for probability density functions, the element of integration can be split.

$$p(x_t|y_{1:t-1}) = \int_{-\infty}^{\infty} p(x_t|x_{t-1}, y_{1:t-1})p(x_{t-1}|y_{1:t-1}) dx_{t-1}$$

With an assumption that the modelled dynamic system (1.1) possesses *Markov Property*<sup>3</sup>,  $p(x_t|x_{t-1}, y_{1:t-1})$  equals  $p(x_t|x_{t-1})$ . [1] This leaves us with the result (1.3).

$$p(x_t|y_{1:t-1}) = \int_{-\infty}^{\infty} p(x_t|x_{t-1})p(x_{t-1}|y_{1:t-1}) dx_{t-1} \quad (1.3)$$

As we can see, *a priori* probability density function only depends on previously known functions and therefore can be calculated.

We continue with the second stage that could be named *update*, where new observation  $y_t$  is taken into account and a *posteriori* probability density function  $p(x_t|y_{1:t})$  is calculated. Bayes' theorem can be used to derive a posteriori probability density function (1.4).

$$p(x_t|y_{1:t}) = \frac{p(y_t|x_t, y_{1:t-1})p(x_t|y_{1:t-1})}{p(y_t|y_{1:t-1})} \quad (1.4)$$

According to the observation model (1.2) and assuming Markov property,  $y_t$  only depends on  $x_t$ . That is  $p(y_t|x_t, y_{1:t-1}) = p(y_t|x_t)$ . Therefore a posteriori probability

---

<sup>2</sup> $p(x_0)$  can be called initial probability density function of the state vector.

<sup>3</sup>an assumption of independence that states that system state in time  $t$  only depends on system state in  $t-1$  (and is not directly affected by previous states).

density function can be further simplified into (1.5).

$$p(x_t|y_{1:t}) = \frac{p(y_t|x_t)p(x_t|y_{1:t-1})}{p(y_t|y_{1:t-1})} \quad (1.5)$$

While both probability density functions in the numerator of (1.5) are already known,  $p(y_t|y_{1:t-1})$  found in the denominator can be calculated using the formula (1.6), where marginalization over  $x_t$  is preformed. Quantity (1.6) can also be interpreted as *marginal likelihood* (sometimes called *evidence*) of observation. [12]

$$p(y_t|y_{1:t-1}) = \int_{-\infty}^{\infty} p(y_t|x_t)p(x_t|y_{1:t-1}) dx_t \quad (1.6)$$

Computing (1.6) isn't however strictly needed as it does not depend on  $x_t$  and serves as a normalising constant in (1.5). Depending on use-case the normalising constant may not be needed at all or may be computed alternatively using the fact that  $p(x_t|y_{1:t})$  integrates to 1.

We have shown that so called *optimal Bayesian solution*[1] can be easily analytically inferred using only *chain rule for probability density functions*, *marginalization* and *Bayes' theorem*. (equations (1.3), (1.5) and (1.6) forming the main steps of the solution) On the other hand, using this method directly in practice proves difficult because at least one parametric multidimensional integration has to be performed (in (1.3)), which is (in its general form) hardly tractable for greater than small state vector dimensions.

This is a motivation for various simplifications and approximations among which we have chosen Kalman filter described in the next section and particle filter family described later.

## 1.3 Kalman Filter

Kalman filter<sup>4</sup> poses additional set of strong assumptions on modelled dynamic system, but greatly simplifies optimal Bayesian solution (1.3), (1.5) into a sequence of algebraic operations with matrices. On the other hand, when these requirements can be fulfilled, there is no better estimator in Bayesian point of view because Kalman filter computes  $p(x_t|y_{1:t})$  *exactly*.<sup>5</sup>

Assumptions additionally posed on system by Kalman filter are:

1.  $f_t$  in the process model (1.1) is a linear function of  $x_t$  and  $v_t$ .
2.  $v_t \sim \mathcal{N}(0, Q_t)$  meaning that process noise  $v_t$  is normally distributed with zero mean<sup>6</sup> and with known covariance matrix  $Q_t$ .
3.  $h_t$  in the observation model (1.2) is a linear function of  $x_t$  and  $w_t$ .

---

<sup>4</sup>first presented by Rudolf Emil Kalman in 1960.

<sup>5</sup>not accounting for numeric errors that arise in practical implementations.

<sup>6</sup>zero mean assumption is not strictly needed, it is however common in many implementations.

4.  $w_t \sim \mathcal{N}(0, R_t)$  meaning that observation noise  $w_t$  is normally distributed with zero mean and with known covariance matrix  $R_t$ .
5. initial state probability density function is Gaussian.

It can be proved that if above assumptions hold,  $p(x_t|y_{1:t})$  is Gaussian for all  $t > 0$ . [10] Furthermore, given assumptions 1. and 2. the process model (1.1) can be reformulated as (1.7), where  $A_t$  is real-valued matrix that represents  $f_t$ . Using the same idea and assumptions 3. and 4. the observation model (1.2) can be expressed as (1.8),  $C_t$  being real-valued matrix representing  $h_t$ . Another common requirement used below in algorithm description is that  $v_t$  and  $w_t$  are stochastically independent.

$$x_t = A_t x_{t-1} + \hat{v}_{t-1} \quad A_t \in \mathbb{R}^{n,n} \quad n \in \mathbb{N} \quad (1.7)$$

$$y_t = C_t x_t + \hat{w}_t \quad C_t \in \mathbb{R}^{j,n} \quad j \in \mathbb{N} \quad j \leq n \quad (1.8)$$

Note that we have marked noises  $v_t$  and  $w_t$  as  $\hat{v}_t$  and  $\hat{w}_t$  when they are transformed through  $A_t$ , respectively  $C_t$  matrix. Let also  $\hat{Q}_t$  denote the covariance matrix of  $\hat{v}_t$  and  $\hat{R}_t$  denote the covariance matrix of  $\hat{w}_t$  in further text.

At this point we can describe the algorithm of Kalman filter. As stated above, a posteriori probability density function is Gaussian and thus can be parametrised by mean vector  $\mu$  and covariance matrix  $P$ . Let us denote a posteriori mean from previous iteration by  $\mu_{t-1|t-1}$  and associated covariance by  $P_{t-1|t-1}$  as in (1.9).

$$p(x_{t-1}|y_{1:t-1}) = \mathcal{N}(\mu_{t-1|t-1}, P_{t-1|t-1}) \quad (1.9)$$

A priori probability density function (1.10) can then be calculated as follows: [1]

$$p(x_t|y_{1:t-1}) = \mathcal{N}(\mu_{t|t-1}, P_{t|t-1}) \quad (1.10)$$

$$\mu_{t|t-1} = A_t \mu_{t-1|t-1}$$

$$P_{t|t-1} = A_t P_{t-1|t-1} A_t^T + \hat{Q}_{t-1}$$

Before introducing a posteriori probability density function it is useful to establish another Gaussian probability density function (1.11) that is not necessarily needed, but is useful because it represents marginal likelihood (1.6).

$$p(y_t|y_{1:t-1}) = \mathcal{N}(\nu_{t|t-1}, S_{t|t-1}) \quad (1.11)$$

$$\nu_{t|t-1} = C_t \mu_{t|t-1}$$

$$S_{t|t-1} = C_t P_{t|t-1} C_t^T + \hat{R}_t$$

The update phase of Kalman filter can be performed by computing so-called *Kalman gain* matrix (1.12), a posteriori probability density function (1.13) is then derived from a priori one using Kalman gain  $K_t$  and observation  $y_t$ . [1]

$$K_t = P_{t|t-1} C_t^T S_{t|t-1}^{-1} \quad (1.12)$$

$$p(x_t|y_{1:t}) = \mathcal{N}(\mu_{t|t}, P_{t|t}) \quad (1.13)$$

$$\mu_{t|t} = \mu_{t|t-1} + K_t (y_t - \nu_{t|t-1})$$

$$P_{t|t} = P_{t|t-1} - K_t C_t P_{t|t-1}$$



In all formulas above  $A^T$  denotes a transpose of matrix  $A$  and  $A^{-1}$  denotes inverse matrix to  $A$ . As can be seen, formulas (1.3) and (1.5) have been reduced to tractable algebraic operations, computing inverse matrix<sup>7</sup> being the most costly one.

It should be further noted that Kalman filter and described algorithm can be easily enhanced to additionally cope with *intervention* (or control) vector applied to the system, making it suitable for the theory of decision-making. Numerous generalisations of Kalman filter exist, for example *extended Kalman filter* that relaxes the requirement of linear system by locally approximating non-linear system with Taylor series. These are out of scope of this text, but provide areas for subsequent consideration.

On the other hand, the assumption of Gaussian a posteriori probability density function cannot be easily overcome and for systems that show out non-Gaussian distributions of state vector another approach have to be taken. [1] One such approach can be Monte Carlo-based *particle filter* presented in the next section.

## 1.4 Particle Filter

Particle filters represent approximate solution of the problem of recursive Bayesian estimation, thus can be considered *suboptimal* methods. Underlying algorithm described below is most commonly named *sequential importance sampling (SIS)*. The biggest advantage of particle filtering is that requirements posed on modelled system are much weaker than those assumed by optimal methods such as Kalman filter. Simple form of particle filter presented in this section (that assumes that modelled system has Markov property) requires only knowledge of probability density function  $p(x_t|x_{t-1})$  representing the process model and knowledge of  $p(y_t|x_t)$  representing the observation model.<sup>8</sup>

Sequential importance sampling approximates posterior density by weighted empirical probability density function (1.14).

$$p(x_t|y_{1:t}) \approx \sum_{i=1}^N \omega_t^{(i)} \delta(x_t - x_t^{(i)}) \quad (1.14)$$

$$\forall i \in \mathbb{N} \quad i \leq N : \omega_i \geq 0 \quad \sum_{i=1}^N \omega_i = 1$$

In (1.14)  $x_t^{(i)}$  denotes value of  $i$ -th *particle*: possible state of the system at time  $t$ ;  $\omega_t^{(i)}$  signifies weight of  $i$ -th particle at time  $t$ : scalar value proportional to expected

<sup>7</sup>it can be shown that  $S_{t|t-1}$  is positive definite given that  $C_t$  is full-ranked, therefore the inverse in (1.12) exists.

<sup>8</sup>both probability density functions are generally time-varying and their knowledge for all  $t$  is needed, but their representation (parametrised by conditioning variable) is frequently constant in time in practical applications.

probability of the system being in state in small neighbourhood of  $x_t^{(i)}$ ;  $N$  denotes total number of particles<sup>9</sup>, a significant tunable parameter of the filter.

Described particle filter is bootstrapped by sampling  $N$  random particles from initial probability density function  $p(x_0)$ . Let  $i \in \mathbb{N}$   $i \leq N$ , transition  $t-1 \rightarrow t$  can be performed as follows:

1. for each  $i$  compute  $x_t^{(i)}$  by random sampling from conditional probability density function  $p(x_t|x_{t-1})$  where  $x_{t-1}^{(i)}$  substitutes  $x_{t-1}$  in condition. This step can be interpreted as a simulation of possible system state developments.
2. for each  $i$  compute weight  $\omega_t^{(i)}$  using (1.15) by taking observation  $y_t$  into account.  $x_t$  is substituted by  $x_t^{(i)}$  in condition in (1.15). Simulated system states are confronted with reality through observation.

$$\omega_t^{(i)} = p(y_t|x_t)\omega_{t-1}^{(i)} \quad (1.15)$$

3. normalise weights according to (1.16) so that approximation of a posteriori probability density function integrates to one.

$$\omega_t^{(i)} = \frac{\omega_t^{(i)}}{\sum_{j=1}^N \omega_t^{(j)}} \quad (1.16)$$

Relative computational ease of described algorithm comes with cost: first, particle filter is in principle non-deterministic because of the random sampling in step 1, in other words, particle filter is essentially a Monte Carlo method; second, appropriate number of particles  $N$  has to be chosen — too small  $N$  can lead to significant approximation error while inadequately large  $N$  can make particle filter infeasibly time-consuming. It can be proved that particle filter converges to true a posteriori density as  $N$  approaches infinity and certain other assumptions hold [4], therefore the number of particles should be chosen as a balance of accuracy and speed.

Only two operations with probability density functions were needed: sampling from  $p(x_t|x_{t-1})$  and evaluating  $p(y_t|x_t)$  in known point. Sometimes sampling from  $p(x_t|x_{t-1})$  is not feasible<sup>10</sup> and/or better results are expected by taking observation  $y_t$  into account during sampling (step 1). This can be achieved by introducing so-called *proposal density* (sometimes *importance density*)  $q(x_t|x_{t-1}, y_t)$ . Sampling in step 1 then uses  $q(x_t|x_{t-1}, y_t)$  instead, where  $x_{t-1}$  in condition is substituted by  $x_{t-1}^{(i)}$ . Weight computation in step 2 have to be replaced with (1.17) that compensates different sampling distribution (every occurrence of  $x_t$ ,  $x_{t-1}$  in mentioned formula has to be substituted by  $x_t^{(i)}$  and  $x_{t-1}^{(i)}$  respectively). See [1] for derivation of these formulas and for discussion about choosing adequate proposal density.

$$\omega_t^{(i)} = \frac{p(y_t|x_t)p(x_t|x_{t-1})}{q(x_t|x_{t-1}, y_t)}\omega_{t-1}^{(i)} \quad (1.17)$$

---

<sup>9</sup> $N$  is assumed to be arbitrary but fixed positive integer for our uses. Variants of particle filter exist that use adaptive number of particles, these are not discussed here.

<sup>10</sup>but can be replaced by evaluation in known point.

Particle filters also suffer from a phenomenon known as *sample impoverishment* or *degeneracy problem*: after a few iterations all but one particles' weight falls close to zero.<sup>11</sup> A measurement of degeneracy can be obtained by computing an approximate of *effective sample size*  $N_{\text{eff}}$  at given time  $t$  using (1.18). [1]

$$N_{\text{eff}} \approx \left( \sum_{i=1}^N \left( \omega_t^{(i)} \right)^2 \right)^{-1} \quad (1.18)$$

Very small  $N_{\text{eff}}$  compared to  $N$  signifies substantial loss of “active” particles, which is certainly undesirable as it hurts accuracy while leaving computational demands unchanged. One technique to diminish this is based on careful choice of proposal density (as explained in [1]) second one is to add additional *resample* step to above algorithm.<sup>12</sup>

4. for each  $i$  resample  $x_t^{(i)}$  from approximate a posteriori probability density function  $\sum_{i=1}^N \omega_t^{(i)} \delta(x_t - x_t^{(i)})$ . Given that sampling is truly random and independent this means that each particle is in average copied  $n_i$  times, where  $n_i$  is roughly proportional to particle weight:  $n_i \approx \omega_t^{(i)} N$ .

Step 4 therefore facilitates avoidance of particles with negligible weight by replacing them with more weighted ones. Such enhanced algorithm is known as *sequential importance resampling (SIR)*.

Recursive Bayesian estimation using SIR methods can be applied to a wide range of dynamic systems (even to those where more specialised methods fail) and can be tuned with number of particles  $N$  and proposal density  $q$ . On the other hand a method specially designed for a given system easily outperforms general particle filter in terms of speed and accuracy.

## 1.5 Marginalized Particle Filter

TODO: MPF. If in lack of time, write short mention and merge into PF section.

---

<sup>11</sup>it has been shown that variance of particle weights continually raises as algorithm progresses. [1]

<sup>12</sup>this additional step may be performed in every iteration or just when  $N_{\text{eff}}$  falls below certain threshold.

# Chapter 2

## Software analysis

In this chapter general software development approaches and practices will be confronted with requirements posed on desired software library for Bayesian decision making. After stating these requirements, feasibility of various programming paradigms applied to our real-world problem is discussed. Continues brief survey of already available software whose conformance to requirements is studied. The chapter is ended by comparison of suitable features of 3 chosen programming languages: C++, Matlab language and Python. Emphasis is put on Python/Cython combination that was chosen for implementation.

In whole chapter, the term *user* refers to someone who uses the library in order to implement higher-level functionality (such as simulation of dynamic systems); user is essentially a programmer.

### 2.1 Requirements

Our intended audience is broad scientific community interested in the field of recursive Bayesian estimation and decision-making. Keeping this in mind and in order to formalise expectations for desired library for Bayesian estimation, following set of requirements was developed.

Functionality:

- Framework for working with potentially conditional probability density functions should be implemented including support for basic operations such as product and chain rule. Chain rule implementation should be flexible in a way that for example  $p(a_t, b_t | a_{t-1}, b_{t-1}) = p(a_t | a_{t-1}, b_t) p(b_t | b_{t-1})$  product can be represented.
- Basic Bayesian filtering methods such as Kalman and particle filter have to be present, plus at least one of more specialised algorithms — marginalized particle filter or non-linear Kalman filter variants.

General:

- Up-to-date, complete and readable API<sup>1</sup> documentation is required. Such documentation should be well understandable by someone that already understands mathematical background of particular algorithm.
- High level of interoperability is needed; data input/output should be straightforward as well as using existing solutions for accompanying tasks such as visualising the results.
- The library should be platform-neutral and have to run on major server and workstation platforms, at least on Microsoft Windows and GNU/Linux.
- The library should be Free/Open-source software as it is believed by the authors that such licensing/development model results in software of greatest quality in long term. Framework used by the library should make it easy to adapt and extend the library for various needs.

Usability:

- Initial barriers for installing and setting up the library should be lowest possible. For example a necessity to install third-party libraries from sources is considered infeasible.
- Implementation environment used for the library should allow for high programmer productivity; prototyping new solutions should be quick and cheap (in terms of effort) operation. This requirement effectively biases towards higher-level programming languages.

Performance:

- Computational overhead<sup>2</sup> should be kept reasonably low.
- Applications built atop of the library should be able to scale well on multi-processor systems. This can be achieved for example by thread-safety of critical library objects or by explicit parallelisation provided by the library.

It is evident that some of the requirements are antagonistic, most prominent example being demand for *low computational overhead* while still offering *high programmer productivity* and rapid prototyping. The task of finding tradeoffs between contradictory tendencies or developing smart solutions that work around traditional limitations is left upon the implementations.

## 2.2 Programming paradigms

Many programming paradigms exist and each programming language usually suggests particular paradigm, though many languages let programmers choose from or combine multiple paradigms. This section discusses how well could be three most prominent paradigms (procedural, object-oriented and functional) applied to software library for Bayesian estimation. Later on additional features of implementation environments such as interpreted vs. compiled approach or argument passing convention are eval-

---

<sup>1</sup>Application Programming Interface, a set of rules that define how a particular library is used.

<sup>2</sup>excess computational costs not directly involved in solving particular problem; for example interpreter overhead.

uated.

### 2.2.1 Procedural paradigm

Procedural paradigm is the traditional approach that appeared along first high-level programming languages. Procedural programming can be viewed as a structured variant of imperative programming, where programmer specifies steps (in form of orders) needed to reach desired program state. Structured approach that empathises dividing the code into logical and self-contained blocks (procedures, modules) is used to make code more reusable, extensible and modular. Today's most notable procedural languages include C and Fortran.

Most procedural languages are associated with very low overhead (performance of programs compiled using optimising compiler tend to be very close to ideal programs written in assembly code); mentioned languages are also spread and well-known in scientific computing.

On the other hand, while possible, it is considered elaborate task by the author to write a modular and extensible library in these languages. Another disadvantage is that usually only very basic building blocks are provided by the language — structures like lists and strings have to be supplied by the programmer or a third-party library. This only adds to the fact that procedural paradigm-oriented languages are commonly not easy to learn and that programmer productivity associated with these languages may be much lower compared to more high-level languages.

### 2.2.2 Object-oriented paradigm

Object-oriented paradigm extends procedural approach with the idea of *objects* — structures with procedures (called *methods*) and variables (called *attributes*) bound to them. Other feature frequently offered is *polymorphism* (an extension to language's type system that adds the notion of *subtypes* and a rule that subtype of a given type can be used everywhere where given type can be used) most often facilitated through a concept of *classes*, common models for sets of objects with same behaviour but different payload; objects are then said to be *instances* of classes. Subclass *inherits* methods and attributes from its superclass and can *override* them or add its own. *Encapsulation*, language mechanism to restrict access to certain object attributes and methods, may be employed by the language to increase robustness by hiding implementation details. In order to be considered object-oriented, statically typed languages (p. 16) should provide *dynamic dispatch*<sup>3</sup>, an essential complement to polymorphism, for certain or all object methods.

Notable examples of languages that support (although not exclusively) object-oriented paradigm are statically typed C++, Java and dynamically typed (p. 16) MATLAB language, Python, Smalltalk.

---

<sup>3</sup>a way of calling methods where the exact method to call is resolved at runtime based on actual (dynamic) object type (in contrast to static object type).

Object-oriented features typically have very small overhead compared to procedural code with equal functionality, so additional complexity introduced is the only downside, in author’s opinion. We believe that these disadvantages are greatly outweighed by powerful features that object-oriented languages provide (when utilised properly).

It was also determined that a library for Bayesian estimation could benefit from many object-oriented techniques: probability density function and its conditional variant could be easily modelled as classes with abstract methods that would represent common operation such as evaluation in a given point or drawing random samples. Classes representing particular probability density functions would then subclass abstract base classes and implement appropriate methods while adding relevant attributes such as border points for uniform distribution. This would allow for example to create generic form of particle filter (p. 9) that would accept any conditional probability density function as a parameter. Bayesian filter itself can be abstracted into a class that would provide a method to compute a posteriori probability density function from a priori one taking observation as a parameter.

### 2.2.3 Functional paradigm

Fundamental idea of functional programming is that functions have no side effects — their result does not change or depend on program state, only on supplied parameters. A language where each function has mentioned attribute is called *purely functional* whereas the same adjective is applied to such functions in other languages. This is often accompanied by a principle that all data are immutable (apart from basic list-like container type) and that functions are so-called “first-class citizens” — they can be passed to a function and returned. Placing a restriction of no side-effect on functions allows compiler/interpreter to do various transformations: parallelisation of function calls whose parameters don’t depend on each other’s results, skipping function calls where the result is unused, caching return values for particular parameters.

Among languages specially designed for functional programming are: Haskell, Lisp dialects Scheme and Clojure, Erlang. Python supports many functional programming techniques<sup>4</sup>.

While functional programming is popular subject of academic research, its use is much less widespread compared to procedural and object-oriented paradigms. Additionally, in the author’s opinion, transition to functional programming requires significant change of programmer’s mindset. Combined with the fact that syntax of mentioned functionally-oriented languages differs significantly from many popular procedural or object-oriented languages, we believe that it would be unsuitable decision for a library that aims for wide adoption.

---

<sup>4</sup>e.g. functions as first-class citizens, closures, list comprehensions

## 2.2.4 Other programming language considerations

Apart from recently discussed general approaches to programming, we should note a few other attributes of languages or their implementations that significantly affect software written using them. The first distinction is based on type system of a language — we may divide them into 2 major groups:

### **statically typed languages**

bind object types to *variables*; vast majority of type-checking is done at compile-time. This means that each variable can be assigned only values of given type (subject to polymorphism); most such languages require that variable (function parameter, object attribute) types are properly declared.

### **dynamically typed languages**

bind object types to *values*; vast majority of type-checking is done at runtime. Programmer can assign and reassign objects of arbitrary types to given variable. Variables (and object attributes) are usually declared by assignment.

We consider dynamically typed languages more convenient for programmers — we’re convinced that the possibility of sensible variable reuse and lack of need to declare variable types lets the programmer focus more on the actual task, especially during prototyping stage. This convenience however comes with a cost: dynamic typing imposes inevitable computing overhead as method calls and attribute accesses must be resolved at runtime. Additionally, compiling a program written in statically typed language can reveal many simple programming errors such as calling mistyped methods, even in unreachable code-paths; this is not the case for dynamically-typed languages and we suggest compensating this with more thorough test-suite.

Another related property is interpreted vs. compiled nature; we should emphasize that this property refers to language *implementation*, not directly to the language itself, e.g. C language is commonly regarded as compiled one, several C interpreters however exist. We use the term “language is compiled/interpreted” to denote that principal implementation of that language is compiled, respectively interpreted.

### **compiled implementations**

translate source code directly into machine code suitable for given target processor. Their advantages are zero overhead, no requirement to install additional software (interpreter) on target machine.

### **interpreted implementations**

either directly execute commands in source code or, more frequently, translate source code into platform-independent *intermediate representation* which is afterwards executed in a *virtual machine*. We may allow the translate and execute steps to be separated so that Java and similar languages can be included. Advantages include usually shorter write-run-debug cycle that speeds up development and portable distribution options. Interpreted languages have been historically associated with considerable processing overhead, but *just-in-time*



*compilation*<sup>5</sup> along with *adaptive optimisation*<sup>6</sup> present in modern interpreters can minimise or even reverse interpreter overhead: Paul Buchheit have shown<sup>7</sup> that second and onward iterations of fractal-generating Java program were actually 5% faster than equivalent C program. We have reproduced the test with following results: Java program was 10% slower (for second and subsequent iterations) than C program and 1600% slower when just-in-time compilation was disabled. Complete test environment along with instructions how to reproduce it be found in [examples/benchmark\\_c\\_java](#) directory in the PyBayes source code repository.

There exists a historic link between statically typed and compiled languages, respectively dynamically typed and interpreted languages. Java which is itself statically typed and it's major implementation is interpreted and Erlang's (which is dynamically typed) compiled HiPE<sup>8</sup> implementation are some examples of languages that break the rule. We believe that this historic link is the source of a common misconception that interpreted languages are inherently slow. Our findings (see also Python/Cython/C benchmark on p. 21) indicate that the source of heavy overhead is likely to be the dynamic type system rather than overhead of modern just-in-time interpreters. In accordance with these findings, we may conclude that choice of language implementation type should rather be based on development and distribution convenience than on expected performance.

Each programming language may support one or more following function call conventions that determine how function parameters are passed:

#### **call-by-value convention**

ensures that called function does not change variables passed as parameters from calling function by copying them at function call time. This provides clear semantics but incurs computational and memory overhead, especially when large data structures are used as parameters. As a form of optimisation, some language implementations may employ copy-on-write technique so that variables are copied only when they are mutated from within called function, thus saving space and time when some parameters are only read from.

#### **call-by-reference convention**

hands fully-privileged references to parameters to called function. These references can be used to modify or assign to parameters within called function and these changes are visible to calling function. This approach minimises function call overhead but may appear confusing to a programmer when local variable is changed "behind her back" unexpectedly. On the other hand, call-by-reference allows for programming techniques impossible with call-by-value alone (e.g. a function that swaps two values).

#### **call-by-object (call-by-sharing) convention**

---

<sup>5</sup>interpreter feature that translates portions of bytecode into machine code at runtime.

<sup>6</sup>a technique to use profiling data from recent past (collected perhaps when relevant portion of code was run in interpreted mode) to optimise just-in-time compiled code.

<sup>7</sup><http://paulbuchheit.blogspot.com/2007/06/java-is-faster-than-c.html>

<sup>8</sup>The High-Performance Erlang Project: <http://www.it.uu.se/research/group/hipe/>

can be viewed as a compromise between call-by-value and call-by-reference: parameters are passed as references that can be used to modify referred objects (unless marked immutable), but cannot be used to assign to referred objects (or this assignment is invisible to calling function). When an object is marked as immutable, passing this object behaves like call-by-value call without copying overhead (in the calling function point of view). Java and Python use call-by-object as their sole function calling method<sup>9</sup> and both mark certain elementary types (most prominently numbers and strings) as immutable. C's pointer-to-const and C++'s reference-to-const parameters can be viewed as call-by-object methods where referred objects are marked as immutable in called function scope.

We suggest that a language that supports at least one of call-by-reference or call-by-object conventions is used for Bayesian estimation library; while call-by-value-only languages can be simpler to implement, we are convinced that they impose unnecessary restrictions on library design and cause overhead in places where it could be avoided.

Last discussed aspect of programming languages relates to memory management:

#### **garbage-collected languages**

provide memory management in the language itself. This fact considerably simplifies programming as programmer doesn't need to reclaim unused memory resources herself. Another advantage is that automatic memory management prevents most occurrences of several programming errors: memory leaks,<sup>10</sup> dangling pointers<sup>11</sup> and double-frees.<sup>12</sup> Two major approaches to garbage collection exist and both incur runtime computational or memory overhead. *Tracing garbage collector* repeatedly scans program heap<sup>13</sup> memory for objects with no references to them, then reclaims memory used by these objects. Program performance may be substantially impacted while tracings garbage collector performs its scan; furthermore the moment when garbage collector fires may be unpredictable. *Reference counting* memory management works by embedding an attribute, *reference count*, to each object that could be allocated on heap and then using this attribute to track number of references to given object. When reference count falls to zero, the object can be destroyed. Reference counting adds small memory overhead per each object allocated and potentially significant computational overhead as reference counts have to be kept up-to-date. However, techniques exist that minimise this overhead, for example those mentioned in [8].

#### **non garbage-collected languages**

put the burden of memory management on shoulders of the programmer: she

---

<sup>9</sup>python case: <http://effbot.org/zone/call-by-object.htm>

<sup>10</sup>an error condition when a region of memory is no longer used, but not reclaimed.

<sup>11</sup>a pointer to an object that has been already destroyed; such pointers are highly error-prone.

<sup>12</sup>an error condition where a single region of memory is reclaimed twice; memory corruption frequently occurs in this case.

<sup>13</sup>an area of memory used for dynamic memory allocation.

is responsible for correctly reclaiming resources when they are no longer in use. The advantages are clear: no overhead due to memory management, probably also smaller complexity of language implementation. However, as mentioned earlier, languages without automatic memory management make certain classes of programmer errors more likely to occur.

In our view, convenience of garbage-collected languages outweighs overhead they bring for a project like a library for Bayesian estimation targeting wide adoption. We also believe that automatic memory management can simplify library design and its usage as there is no need to specify who is responsible for destroying involved objects on the library side and no need to think about it at the user side.

## 2.3 C++

C++ is regarded as one of the most popular programming languages today, along with Java and C;<sup>14</sup> it combines properties of both low-level and high-level languages, sometimes being described as intermediate-level language. C++ extensively supports both procedural and class-based object-oriented paradigm, forming a multi-paradigm language; generic programming is implemented by means of *templates*, which allow classes and functions to operate on arbitrary data types while still being type-safe. C++ is statically-typed, all major implementations are compiled, supports call-by-value (the default), call-by-reference and a variant of call-by-object function call conventions. C++ lacks implicit garbage collection for heap-allocated data — the programmer must reclaim memory used by those objects manually; use of *smart pointers*<sup>15</sup> may although help with this task. C++ is almost 100% compatible with the C language in a way that most C programs compile and run fine then compiled as C++ programs. C++ also makes it easy to use C libraries without a need to recompile them. [13]

When used as an implementation language for desired library for recursive Bayesian estimation, we have identified potential advantages of the C++ language:

### **widespread**

C/C++ code forms large part of the software ecosystem. Thanks to that

### **low overhead**

C++ was designed to incur minimal overhead possible. In all benchmarks we've seen (e.g. The Computer Language Benchmarks Game<sup>16</sup>), C++ is hard to outperform by a significant margin (

set of drawbacks

BDM

---

<sup>14</sup>TIOBE Programming Community Index for July 2011: <http://www.tiobe.com/index.php/content/paperinfo/tpci/index.html>

<sup>15</sup>a template class that behaves like a pointer through use of operator overloading but adds additional memory management features such as reference counting

<sup>16</sup><http://shootout.alioth.debian.org/>

- advantages of C (speed, C prevalence (many optimised libraries, BLAS, LAPACK..., OpenMP))
- disadvantages of C in our “situation” (steep learning curve, complexity because of low-levelness high initial barriers (need to have compiler, libraries...), inconveniently long edit/build/test process)

## 2.4 Matlab

BDM (partially?)

- advantages (popularity, existing toolboxes, rapid development (high-level))
- disadv: strict copy-on-write, problematic object model (not in original design), difficulties interfacing existing C (F) code

## 2.5 Python

NumPy... parallelisation (approaches, improvements in Py 3.2) - GIL.. Py3k

## 2.6 Cython

general info etc... extension types, building, ease of interfacing C (and F) code, .pxd files, NumPy support

[citations:[3](#), [11](#), [2](#)]

### 2.6.1 Gradual Optimisation

how can optimisation be approached (gradually) and why this approach is superior

[see *cypy*...]

### 2.6.2 Parallelisation

integrate `_python_cython` patched with OpenMP (13x speedup in 16-core system)

`prange` CEP – implemented!

[see *cypy*...]

### 2.6.3 Pure Python mode

About it and why it should be used in a hypothetical bayesian python library

### 2.6.4 Limitations

2 types:

- not-supported code (few cases, but bad, ongoing work)
- not-optimised code (much more work needed, but not hard to fix in most cases)
  - exception handling (functions returning void etc)
  - limitations of pure python mode in regards to traditional .pyx files

### 2.6.5 Performance comparison with C and Python

benchmark<sub>ccypy</sub>

## 2.7 Survey of Existing Libraries for Bayesian estimation/decision making

Brief survey of existing libraries ... not fullfilling all requirements..

The need to implement a new one :-)

## 2.8 Choice

python/cython was choosen ...

# Chapter 3

## The PyBayes Library

Introduction, general directions, future considerations  
+ open development on github, open-source

### 3.1 Interpreted and Compiled

### 3.2 Library Layout

[proposed citation: [\[15\]](#)]

#### 3.2.1 Random Variable Meta-representation

Why it is needed (ref to ProdCPdf)

#### 3.2.2 Probability Density Functions

Nice UML diagrams! (better more smaller UMLs than one big) One for general pdf layout, one for AbstractGaussPdf family, one for AbstractEmpPdf family

#### 3.2.3 Bayesian Filters

UML

Nice graph of a run of a particle filter (Mirda has the plotting code)  
similar of marginalized particle filter? (gausses would be plotted vertically)  
[mention this:[\[12\]](#)]

### 3.3 Documentation, Testing and Profiling

TODO: move above Library Layout?

Documenting PyBayes using Sphinx, approach to documentation (mathematician-oriented), math in documentation

Testing - the separation of

- tests: test one class in isolation, quick, determinism (would be good, not achievable)

- stresses: test a great portion of code at once, run longer, non-determinism..

Profiling python/cython - how, existing support in PyBayes

- how to correct profiling-induced overhead

### 3.4 Comparison with BDM

[skip if in time press]

# Conclusion



# Bibliography

- [1] M. Sanjeev Arulampalam, Simon Maskell, and Neil Gordon. A tutorial on particle filters for online nonlinear/non-gaussian bayesian tracking. *IEEE Transactions on Signal Processing*, 50:174–188, 2002. 6, 7, 8, 9, 10, 11
- [2] S. Behnel, R. Bradshaw, C. Citro, L. Dalcin, D.S. Seljebotn, and K. Smith. Cython: The best of both worlds. *Computing in Science Engineering*, 13(2):31–39, march-april 2011. 20
- [3] Stefan Behnel, Robert W. Bradshaw, and Dag Sverre Seljebotn. Cython tutorial. In Gaël Varoquaux, Stéfan van der Walt, and Jarrod Millman, editors, *Proceedings of the 8th Python in Science Conference*, pages 4–14, Pasadena, CA USA, 2009. 20
- [4] D. Crisan and A. Doucet. A survey of convergence results on particle filtering methods for practitioners. *Signal Processing, IEEE Transactions on*, 50(3):736–746, 2002. 10
- [5] F. Gustafsson, F. Gunnarsson, N. Bergman, U. Forssell, J. Jansson, R. Karlsson, and P.J. Nordlund. Particle filters for positioning, navigation, and tracking. *Signal Processing, IEEE Transactions on*, 50(2):425–437, 2002. 4
- [6] R. Hofman, V. Šmídl, and P. Pecha. Data assimilation in early phase of radiation accident using particle filter. In *The Fifth WMO International Symposium on Data Assimilation*, Melbourne, Australia, 2009. 4
- [7] R. Hofman and Šmídl V. Assimilation of spatio-temporal distribution of radionuclides in early phase of radiation accident. *Bezpečnost jaderné energie*, 18:226–228, 2010. 4
- [8] Yossi Levanoni and Erez Petrank. An on-the-fly reference counting garbage collector for Java. *ACM Transactions on Programming Languages and Systems*, 28(1), January 2006. 18
- [9] P. Pecha, Hofman R., and Šmídl V. Bayesian tracking of the toxic plume spreading in the early stage of radiation accident. In *Proceedings of the 2009 European Simulation and Modelling Conference*, Leicester, GB, 2009. 4

- [10] V. Peterka. Bayesian approach to system identification. In P. Eykhoff, editor, *Trends and Progress in System identification*, pages 239–304. Pergamon Press, Oxford, 1981. 8
- [11] Dag Sverre Seljebotn. Fast numerical computations with cython. In Gaël Varoquaux, Stéfan van der Walt, and Jarrod Millman, editors, *Proceedings of the 8th Python in Science Conference*, pages 15–22, Pasadena, CA USA, 2009. 20
- [12] V. Šmídl. Software analysis unifying particle filtering and marginalized particle filtering. In *Proceedings of the 13th International Conference on Information Fusion*. IET, 2010. 7, 22
- [13] B. Stroustrup. *The C++ programming language, Third Edition*. Addison-Wesley, 2000. 19
- [14] S. Thrun, W. Burgard, and D. Fox. *Probabilistic Robotics (Intelligent Robotics and Autonomous Agents)*. 2005. 4
- [15] V. Šmídl. *Software analysis of Bayesian distributed dynamic decision making*. PhD thesis, University of West Bohemia, Faculty of Applied Sciences, Pilsen, Czech Republic, Plzeň, 2005. 22