

# 1 Summary of the unwritten research.

## 1. Still on Baehrens

Doing more reading on Gaussian Processes and overall a review of Markov Chain and other stochastic processes to get a better Idea. This is mostly through Tsitsilkis's Probability and other references or lectures found online.

Some more Baehrens rabbit hole:

1. obrezanova2009. States that GPs tend to produce more predictive models than their counterparts.

- GPs seem to be inherently resistant against overfitting(?)
- It offers the ability to show uncertainty in predictions.
- Tends to be robust for nonlinear regression.
- can be applied to both regression and classification.
- **Methods and Data:**
  - Goal for the GP Gaussian Process Model: Is to model the probability distribution of the class  $Y$  for a molecule, given its descriptor vector  $\mathbf{x}$ ,  $p(y|\mathbf{x})$ .

2. claywood2017 where they talk about using Gaussian Processes for more interpretability of the tasks at hand.

## 2. Bayesian Inference from good old Tsitsilkis & Bertsekas

- Bayesian Statistical Inference: This is the process of extracting information about an unknown variable or model by using available data.
- What makes statistics a bit different from probability: Probability relies on axioms, some assumptions and the consequences of their combinations; probability, on the other hand, can yield different answers to the same question (all of them could be reasonable).
- Within the field of statistics there are two main schools of thought:
  - **Bayesian:** Treats unknown variables as random variables with known (prior) distributions. This is done by assigning a random variable  $\Theta$  that characterizes the model and by postulating a **prior probability** distribution  $p_{\Theta}(\theta)$ . We then would use *Baye's rule* to derive a **posterior probability** distribution  $p_{\Theta|X}(\theta|x)$
  - **Classical/Frequentist:** Unknown variables are treated as *quantities* that happen to be unknown
  - They differ in the way they view of the nature of unknown models or variables

- Two types of inferences:
  - **Model inference:** we study a real phenomenon or process. Model can then **make predictions of the future**. Model of some planetary trajectory
  - **Variable Inference:** Estimate the value of one or more unknown variables, by using some info. e.g. values sent via noisy channel
  - Some blurred lines in this distinction.
- Classification of **Statistical Inference Problems**
  1. **Estimation:** Model is known but we only want to estimate a (possibly multidimensional) parameter  $\theta$  (which could be viewed as a random variable).
  2. **Binary Hypothesis Testing** problem: Two hypothesis exists and data is used to decide which of the two is true. More generally the **m-ary hypothesis testing** for  $m$  hypothesis to be tested.
- There are problems where the uncertain object **cannot** be described by a **fixed** number of **parameters**. These are called **nonparametric** problems/models.
- **Inference Methods:**
  1. **Maximum a posteriori probability** (MAP): Select the parameter/hypothesis with maximum conditional/posterior probability given the data.
  2. **Least Mean Squares** (LMS) estimation: Select estimator/model that *minimizes* the mean squared error between parameter and estimate.
  3. **Linear least mean squares:** Select an estimator which is a linear function of the data and minimizes the mean squared error between the parameter and its estimate.
- **Bayesian Inference and the Posterior Distribution**
  - Assume we know:
    1. The joint distribution of  $\Theta$  and  $X$
    2. Prior distribution  $p_{\Theta}$  (discrete) or  $f_{\Theta}$  (continuous)
    3. Conditional Distribution  $p_{X|\Theta}$
  - Once we get a sample  $x$  from  $X$  we can use posterior distribution  $p_{\Theta|x}$  of  $\Theta$  to solve Bayesian inference problem.

## 2 Gaussian Process Regression From First Principles

article

Found gold baby. The paper linked is a very quick and concise intro to gaussian processes for regression.

One way to view Gaussian Processes is as a distribution *over functions*. Given the parameters that define the GP (Mean and covariance matrices) we can sample a function at the point  $\mathbf{x} \in R^d$  (R for reals)

The sampling goes as follows:

$$f(\mathbf{x}) \sim GP(m(\mathbf{x}), k(\mathbf{x}, \mathbf{x}'))$$

Where  $m(\cdot)$  is for mean and  $k(\cdot, \cdot)$  for covariance

So it seems to be that if we have our set

$$\mathbf{X}_1, \dots, \mathbf{X}_j \sim N(\mu, \Sigma)$$

Then those random variables are used to sample  $f(\mathbf{X})$  on every step. Remember that we read that functions can sort of be thought as mapping one vector to another. Maybe we are mapping independent random variables to their samples and that forms a function. When we sample from the entire set once again then the new vector of outputs that we get will yield a different function.

No, no. I think that its rather that each random variable can yield a function.

Given an a dataset, GPR predicts a posterior Gaussian Distribution for targets over test points. BY computing the parameters of this Gaussian distribution given observed training data.

This *non-parametric* property of GPR seems to be all the rage

### 3 Some explanation of the Kernel Trick

Source

Goal would be to linearly separate previously unseparable elements in higher dimensions. Best example in article is  $\theta(x) = x \bmod 2$

Basically its the art of selecting the right transformation to make some boundary that is a  $n-1$  hyperplane of the extended  $n$  dimensions.

The support vectors are the vectors that describe the points that allow us to specify a hyperplane that defines our decision boundary.

#### 3.1 The actual explanation now

SO whats the problem with this ? It seems that this higher dimensional computations would naturally become expensive and prohibitive. The **Kernel Trick** solves this by using methods that represent our data using a set of pairwise similarity comparisons between the original data observations instad of explicitly applying the transformations and representing the data by these transformed coordinates in higher dimensional feature space.

## 4 Limit Theorems

These theorems are related to the asymptotic behavior of sequence of random variables

Im not sure if im understanding this kernel trick .