

Robust Timeseries Analysis in the Context of Multiple Asynchronous Input Channels

Ryan Michael
kerinin@gmail.com

March 10, 2009

Abstract

Spectral Analysis is a common form of analysing timeseries data. We propose a method of spectral analysis based on sets probability distributions generated from subsets of the observed data. By decomposing the observations into a set of observed behaviors, the spectral analysis can be used more accurately to predict the timeseries' behavior.

1 Introduction

1.1 Existing Work

http://en.wikipedia.org/wiki/Linear_discriminant_analysis

1.2 General Overview

The goal is to create a method of statistical inference capable of processing timeseries data which is both multi-variate and exhibits different behaviors at different times. This type of data is common, and developing a robust method of analysis has applications in many domains. The general approach is to create a series of estimates using subsets of the observed data, and to then combine these estimates in an intelligent manner which captures the relevance of each estimate to the current prediction task. By using a set of 'typical' estimates, we are able to reduce the computational demands of the system, as each estimate is a condensed representation of the data from which it was derived. This approach also allows us to reduce data redundancy by only using distinct estimates.

The most basic operation used in this system is the estimation of probability densities. Based on a set of observations drawn from some random process, we generate an estimate of the underlying probability distribution. This estimate tells us the probability of each point in the input space being observed. Areas of the input space in which a dense set of observations are observed are given high probability, while areas of the input space with few observations are given low probability. This

basic operation, in combination with some basic laws of statistics allow us to build the full inference system.

We assume that observations are pulled from multiple independent sources, all of which respond to some underlying phenomena. For instance one set of observations could be from a microphone and another from a light detector. We do not know how the inputs are related or to what extent their behavior changes over time. For example one input could be a steady sine wave and another could be based on the decay of a radioactive isotope. In the former case previous observations of the source are useful, in the latter they're not. Alternately two inputs could be light detectors in the same room or they could be on different continents; in the former their input would be highly similar, in the latter not as much.

Our general strategy has three phases; generating estimates, correlating estimates, and applying estimates to a given prediction task.

1.3 Generating Estimates

Estimates are generated by picking a time window at random and only dealing with observations which take place in that window. Using these observations we create an estimate of the probability distribution underlying the observations (this distribution would include time as a variable). This process is repeated until we feel we have a reasonable sample of the system as a whole, at which time we can start to correlate the estimates.

1.4 Correlating Estimates

The goal of correlating estimates is to determine the relationships between estimates at different times and from different sources. Estimates are correlated by treating them as variables whose value for a given time window is determined by the extent to which the observed data corresponds to the estimate. For each time window there exist a set of estimates which have some value describing their accuracy at predicting the observed value. We can treat the accuracy measurement of each estimate as a multi-dimensional point, each dimension determined by an estimate's accuracy. Using a set of these points taken at different times, we create a probability distribution estimate. The domain of this probability distribution has the same dimensionality as the number of estimates we have generated. This probability density allows us to predict the probability of an estimate in of one source based on the probability of an estimate in of another source because frequent combinations of accuracy values will have higher probability than other combinations. This 'correlation density' also tells us which estimates are most commonly observed - information we can use in conjunction with estimate similarity to determine which estimates to use and which to discard.

1.5 Making Predictions

Once the estimates have been correlated, we are able to generate predictions. For simplicity, we'll assume that the first two phases occur on a set of 'training' data which is representative of the underlying data, while prediction takes place continuously using new sets of observations which occur in some time window. We make predictions for a given source by combining the existing estimates we have for that source based on their accuracy at predicting the given observations. Because we have determined the correlation between estimate accuracy of different sources, we can use other sources to refine our confidence in each estimate of the given source; the influence of estimates of the given source which do not correspond to a likely 'point' in the correlation density are suppressed while the influence of estimates which correspond to likely points in the correlation density are enhanced.

2 Formal Description

2.1 Problem Setting

We begin with a hidden random variable

$$X = (\Omega, \mathcal{F}, \mathcal{P})$$

Our knowledge of X comes from a set of independent sources which we treat as random processes generated by X :

$$\begin{aligned} X^n : \Omega &\mapsto \Omega^n \in \mathbb{R}^d \times \mathbb{R}_+ \\ X^n &= (\Omega^n, \mathcal{F}^n, \mathcal{P}^n) \\ \mathbf{X} &= [X^0, \dots, X^c] \end{aligned}$$

We refer to each of these sources as a channel, and refer to each channel as the i^{th} element of the set \mathbf{X} :

$$C^n = [\Omega^n, \mathcal{F}^n, \mathcal{P}^n]$$

For each channel we are given a set of ℓ observations of dimension d :

$$\begin{aligned} \mathbf{x}_i^n &= (x_i^n, t_i^n) \\ X^n &= [\mathbf{x}_0^n, \dots, \mathbf{x}_\ell^n] \in \mathbb{R}^d \times \mathbb{R}_+ \end{aligned}$$

Given a set of time durations $\theta \in \Theta$, we define a set of time windows:

$$\mathcal{T}^n = \begin{bmatrix} (t_0^n, \theta_0) & \dots & (t_0^n, \theta_z) \\ \vdots & \ddots & \vdots \\ (t_\ell^n, \theta_0) & \dots & (t_\ell^n, \theta_z) \end{bmatrix}$$

We assume that \mathcal{P} can be approximated in a given time window using a scaled and shifted set of weighted distributions defined over some set of time intervals.

$$\mathcal{P}(t, \theta) = \sum_i \delta_i \cdot f_i(t - t_i) \frac{\theta}{\theta_i} \quad (1)$$

where δ_n is the weight corresponding to f_n . Finally, we assume that a similar mixture of distributions can be determined for each channel:

$$\mathcal{P}^n(t, \theta) = \sum_i \delta_i^n \cdot f_i^n(t - t_i) \frac{\theta}{\theta_i} \quad (2)$$

2.2 Single Channel Setting

We begin by considering the case where only one channel exists, so for now we will omit the superscript and refer to $X^n = (\Omega^n, \mathcal{F}^n, \mathcal{P}^n)$ as $X = (\Omega, \mathcal{F}, \mathcal{P})$. Given a set of test data \hat{X} , our goal is to estimate the probability distribution of X over some time window $\mathcal{T}_{\hat{x}}$:

$$\begin{aligned} \hat{\mathbf{x}}_i &= (\hat{x}_i, t_i) \\ \hat{X} &= [\hat{\mathbf{x}}_0, \dots, \hat{\mathbf{x}}_k] \\ \mathcal{T}_{\hat{x}} &= (t_{\hat{x}}, \theta_{\hat{x}}), \quad t_{\hat{x}} < \min_t \hat{X} < \max_t \hat{X} < t_{\hat{x}} + \theta_{\hat{x}} \end{aligned}$$

We begin by defining the subset of the training observations which fall into each time window, scaled and shifted to the interval $t \in [0, 1]$:

$$\begin{aligned} S_{t, \theta} &= \left[\left(x_i, \frac{t_i - t}{\theta} \right) \mid x_i \in X, t \leq t_i < t + \theta \right] \\ \mathcal{S} &= \begin{bmatrix} S_{t_0, \theta_0} & \dots & S_{t_0, \theta_z} \\ \vdots & \ddots & \vdots \\ S_{t_\ell, \theta_0} & \dots & S_{t_\ell, \theta_z} \end{bmatrix} \end{aligned} \quad (3)$$

We treat \mathcal{S} as a random process:

$$\mathcal{S} = [\Omega^{\mathcal{S}}, \mathcal{F}^{\mathcal{S}}, \mathcal{P}^{\mathcal{S}}]$$

and observe that \hat{X} can be treated as an observation of \mathcal{S} :

$$S_{\hat{x}} = \left[\left(x_i, \frac{t_i - t_{\hat{x}}}{\theta_{\hat{x}}} \right) \mid (x_i, t_i) \in \hat{X} \right] \quad (4)$$

We can now frame the task of estimating the probability distribution of \hat{X} as a task of estimating a probability density function φ for the random process \mathcal{S} :

$$\Pr(X = \mathbf{x} \mid \hat{X}) \mapsto \Pr(\mathcal{S} = \{S_{\hat{x}} \cup \mathbf{x}\}) \simeq \varphi(\mathbf{x}, S_{\hat{x}}) \quad (5)$$

2.3 Multiple Channel Setting

In order to extend this result to settings in which multiple channels exist, we return to 3 and extend the scope to multiple channels:

$$S_{t,\theta} = \left[\left(x_i^n, \frac{t_i - t}{\theta} \right) \mid x_i^n \in \mathbf{X}, t \leq t_i < t + \theta \right] \quad (6)$$

In this case, we treat each channel as an orthonormal basis of the abstract space Ω^S . Equation 4 is likewise extended in the same manner:

$$S_{\hat{x}} = \left[\left(x_i^n, \frac{t_i - t_{\hat{x}}}{\theta_{\hat{x}}} \right) \mid (x_i^n, t_i) \in \hat{\mathbf{X}} \right] \quad (7)$$

3 Parzen Window Estimation

One method of evaluating 5 is by using the Parzen Window method. We choose the Parzen Window method because it allows us to estimate probabilities of unordered sets, provided they have an addition operation and a kernel function exists to provide a distance metric.

3.1 Single Channel Parzen Window

We will again begin by considering the single-channel case, then extend the resulting equations as necessary. The Parzen Window method requires the definition of a metric over the abstract space Ω^S . Such a metric can be defined using the symmetric Kullback Liebler divergence with probability measures $\phi_n(\mathbf{x}) \simeq Pr(X = \mathbf{x} \mid S_n)$:

$$\begin{aligned} D_{KL}(S_n \| S_m) &= \sum_{\mathbf{x} \in \{S_n \cup S_m\}} \phi_n(\mathbf{x}) \log \frac{\phi_n(\mathbf{x})}{\phi_m(\mathbf{x})} \\ &= - \sum_{\mathbf{x} \in \{S_n \cup S_m\}} \phi_n(\mathbf{x}) \log \phi_m(\mathbf{x}) + \sum_{\mathbf{x} \in \{S_n \cup S_m\}} \phi_n(\mathbf{x}) \log \phi_n(\mathbf{x}) \end{aligned} \quad (8)$$

$$\|S_n - S_m\|_{KL} = D_{KL}(S_n \| S_m) + D_{KL}(S_m \| S_n) \quad (9)$$

The Parzen Window estimation of \mathcal{P}^S is defined as:

$$\Pr(\mathcal{S} = S) \simeq \sum_{i \in \mathcal{S}} \frac{1}{|\mathcal{S}|} K_\gamma(S, S_i) \quad (10)$$

where $|\cdot|$ denotes the cardinality of (\cdot) and K is some kernel function, for instance the Radial Basis Function:

$$K_\gamma(S_i, S_j) = \frac{1}{\gamma\sqrt{2\pi}} e^{-\frac{1}{\gamma}\|S_i - S_j\|_{KL}^2} \quad (11)$$

The same method can be used to estimate $\phi_n(\mathbf{x})$ for a given subset S_n :

$$\phi_n(\mathbf{x}) = \sum_{\mathbf{x}_i \in S_n} \frac{1}{|S_n|} K_\gamma(\mathbf{x}, \mathbf{x}_i) \quad (12)$$

Substituting 10 into 5 our probability distribution estimate becomes:

$$\begin{aligned} \varphi_P(\mathbf{x}, S) &= \sum_{i \in S} \frac{1}{|S|} K_\gamma(S_{\mathbf{x}}, S_i) \\ S_{\mathbf{x}} &= \{S \cup \mathbf{x}\} \end{aligned} \quad (13)$$

3.2 Multiple Channel Parzen Window

Extending the Parzen Window approach requires the realization that 8 requires that S_n and S_m both be defined over the same abstract space¹. As mentioned earlier, in the Multiple Channel context each channel is treated as an orthonormal basis of Ω^S . An obvious approach to defining a measure over Ω^S for multiple channels is to use the Euclidean norm of the Kullback Liebler divergence of each channel considered independently:

$$S_n^c = [\mathbf{x} \mid \mathbf{x} \in \{S_n \cap X^c\}]$$

This requires the following minor extension of 11:

$$K_\gamma(S_i, S_j) = \prod_c \frac{1}{\gamma\sqrt{2\pi}} e^{-\frac{1}{\gamma}\|S_i^c - S_j^c\|_{KL}^2} \quad (14)$$

4 Support Vector Estimation

The Parzen Window method is neither sparse nor computationally efficient, and as the number of observations grows, these deficiencies quickly become prohibitive. We now investigate the use of Support Vector Machines to generate $\varphi(\mathbf{x}, S)$.

¹ For instance if $X_n \in \mathbb{R}^2$ and $X_m \in \mathbb{R}^3$, it is impossible to calculate $\phi_n((x_m, t))$ because the quantity $\|x_m - x_n\|^2$ from using 12 is ambiguous.

4.1 Random Process Estimation

Support Vector Machines are usually used to estimate probability distributions by solving the related problem of estimating the cumulative distribution function of the random variable in question. This reduces the problem to one of estimating a non-linear mapping from observations to cumulative distribution values, which can be formulated as an optimization problem over a linear operator equation. Unfortunately, these methods depend on the ability to calculate an empirical distribution for each observation:

$$F_\ell(x) = \frac{1}{\ell} \sum_i \theta(x - x_i) \quad (15)$$

where $\theta(x)$ is the indicator function. To evaluate this function, the abstract space Ω^S must be ordered. While we have described a distance metric over Ω^S , it is not clear what a meaningful ordering relation would be.

Rather than calculating the cumulative probability distribution of Ω^S , we begin with the assumption that the Parzen Window estimate of the probability distribution is accurate and attempt to minimize the square loss between the Support Vector estimate and the Parzen Window estimate. Because we are hoping to generate a sparse representation of the probability distribution, we add a regularizing term Ω which penalizes similar S . The Support Vector approach seeks a solution in the following form:

$$\begin{aligned} \varphi_{SV}(\mathbf{x}, S : \beta) &= \sum_i \beta_i K_\gamma(S_{\mathbf{x}}, S_i) \\ S_{\mathbf{x}} &= \{S \cup \mathbf{x}\} \end{aligned} \quad (16)$$

So we can express the Support Vector optimization problem as:

$$\begin{aligned} W(\beta) &= \sum_{\mathbf{x} \in X} \left(\varphi_P(\mathbf{x}, S) - \varphi_{SV}(\mathbf{x}, S : \beta) \right)^2 + \lambda \Omega(\beta, S) \\ &= \sum_{\mathbf{x} \in X} \left(\varphi_{SV}(\mathbf{x}, S : \beta)^2 - 2\varphi_P(\mathbf{x}, S)\varphi_{SV}(\mathbf{x}, S : \beta) \right) + \lambda \Omega(\beta, S) \\ &= \sum_{i,j} \beta_i \beta_j \sum_{\mathbf{x} \in X} K_\gamma(S_{\mathbf{x}}, S_i) K_\gamma(S_{\mathbf{x}}, S_j) - \sum_i \beta_i \sum_{\mathbf{x} \in X} \sum_j \frac{2}{|S|} K_\gamma(S_{\mathbf{x}}, S_i) K_\gamma(S_{\mathbf{x}}, S_j) + \lambda \sum_i \beta_i \sum_{\mathbf{x} \in X} K_\gamma(S_{\mathbf{x}}, S_i)^{-1} \\ &= \sum_{i,j} \beta_i \beta_j \sum_{\mathbf{x} \in X} K_\gamma(S_{\mathbf{x}}, S_i) K_\gamma(S_{\mathbf{x}}, S_j) + \sum_i \beta_i \sum_{\mathbf{x} \in X} \left(\lambda K_\gamma(S_{\mathbf{x}}, S_i)^{-1} - \frac{2}{|S|} \sum_j K_\gamma(S_{\mathbf{x}}, S_i) K_\gamma(S_{\mathbf{x}}, S_j) \right) \\ \text{subject to } & \sum_i \beta_i = 1, \quad \beta_i \geq 0, \quad i = 1, \dots, |S| \end{aligned} \quad (17)$$

Notice the regularizer selected is defined as:

$$\Omega(\beta, S) = \sum_i \beta_i \sum_{\mathbf{x} \in X} K_\gamma(S_{\mathbf{x}}, S_i)^{-1} \quad (18)$$

Equation 17 is a quadratic optimization problem defined as:

$$\begin{aligned} W(\beta) &= \frac{1}{2} \beta^T P \beta + q^T \beta \\ P_{i,j} &= \sum_{\mathbf{x} \in X} K_\gamma(S_{\mathbf{x}}, S_i) K_\gamma(S_{\mathbf{x}}, S_j) \\ q_i &= \sum_{\mathbf{x} \in X} \left(\lambda K_\gamma(S_{\mathbf{x}}, S_i)^{-1} - \frac{2}{|S|} \sum_j K_\gamma(S_{\mathbf{x}}, S_i) K_\gamma(S_{\mathbf{x}}, S_j) \right) \end{aligned}$$

$$P = \langle \mathbf{K}^T \cdot \mathbf{K} \rangle \quad (19)$$

$$q = \lambda \langle \mathbf{K}^T \cdot \mathbf{1}_{(|K|,1)} \rangle^{-1} - \frac{2}{|S|} \langle \mathbf{K}^T \cdot \mathbf{K} \cdot \mathbf{1}_{(|K|,1)} \rangle \quad (20)$$

4.2 Random Process Kernel Definition

When developing the Parzen Window algorithm, we used kernel function 8 which is based on the Kullback Liebler divergence of a probability estimate 12 defined at the sets being evaluated. Our motivation in developing a Support Vector approach is to generate sparse representations of \mathcal{P}^S , in part to reduce the computational demands of evaluating $\varphi(\mathbf{x}, S)$ for test data sets. Unfortunately, kernel function 8 requires probability estimates of both sets being compared - it would be helpful to develop a kernel function capable of evaluating a distance between a test set $S_{\hat{X}}$ and a training set S_i without first calculating an estimate of the probability distribution of $S_{\hat{X}}$.

Note that 17 is formulated in such a way that the first kernel argument is always $S_{\mathbf{x}}$, which allows us to define a kernel which takes a set of points as its first argument and a probability distribution as its second. The definition of 19 also allows us to define non-symmetric kernels without sacrificing convexity or monotonicity in the optimization objective function. We can therefore use any measure of the divergence between a set of points and a probability distribution. We consider the Renyi entropy with $\alpha = 1$ and select the Shannon Entropy as our kernel metric, as it is equivalent to the Kullback Liebler divergence.

$$\|X - \phi\|_H = \sum_{x \in X} \phi(x) \log \phi(x) \quad (21)$$

Recalling that $S_{\mathbf{x}} = \{S \cup \mathbf{x}\}$, we observe that:

$$\|S_{\mathbf{x}} - S_i\|_H = \|S - S_i\|_H + \|\mathbf{x} - S_i\|_H \quad (22)$$

We restate the kernel function as:

$$K_\gamma(X, S_i) = \frac{1}{\gamma\sqrt{2\pi}} e^{-\frac{1}{\gamma}\|X - S_i\|_H^2} \quad (23)$$

In situations where probabilities are being calculated for a uniform set of points in Ω^S , this can considerably reduce the computational time needed to evaluate \mathbf{K} , as $\|\mathbf{x} - S_i\|_H$ can be computed once and then added to each $\|S - S_i\|_H, S \in \mathcal{S}$.

4.3 Random Vector Estimation

Evaluating 17 over a training set X produces a set $\mathcal{S}_{SV} \in \mathcal{S}$ referred to as Support Vectors. While 17 produces a sparse representation of \mathcal{P}^S by eliminating some S , further sparseness can be achieved by refining the definition of $\phi(\mathbf{x})$ to allow the elimination of some \mathbf{x} in each $S \in \mathcal{S}_{SV}$. In the context of computing the kernel matrices used in optimizing 17 the Parzen Window definition of $\phi(\mathbf{x})$ is computationally acceptable, as it results in a good approximation with a minimal amount of computation². In the context of evaluating 16 over a testing set \hat{X} , the Parzen Window algorithm for $\phi(\mathbf{x})$ is sub-optimal due to the fact that it requires the full set of observations for each $S \in \mathcal{S}_{SV}$. In this context, a sparse algorithm for $\phi(\mathbf{x})$ would reduce both the data required to store \mathcal{S}_{SV} and the computational resources required to evaluate ??.

We return now to the empirical cumulative distribution function method of Support Vector density estimation and search for a solution in the following form, which is able to use multiple non-symmetric kernel functions simultaneously:

$$\phi(x) = \sum_{i=1}^{\ell} (\beta_i^1 \mathcal{K}_1(x_i, x) + \dots + \beta_i^\kappa \mathcal{K}_\kappa(x_i, x)) \quad (24)$$

We select values of β which minimize the square loss of the empirical cumulative probability distribution using some regularizer:

$$\min \left(\sum_{i=1}^{\ell} \left(y_i - \sum_{j=1}^{\ell} \sum_{n=1}^{\kappa} \beta_j^n k_n(x_i, x_j) \right)^2 + \lambda \sum_{i=1}^{\ell} \sum_{n=1}^{\kappa} \frac{1}{\gamma_n} \beta_i^n \right) \quad (25)$$

$$\text{subject to } \sum_{i=1}^{\ell} \sum_{n=1}^{\kappa} \beta_i^n = 1, \quad \beta_i \geq 0 \quad (26)$$

$$(27)$$

² The Parzen Window estimate at a point requires distance computations for each observation and a summation. By contrast, the SV estimate requires the same number of distance computations as well as the solving of a quadratic optimization problem whose complexity increases exponentially with the number of observations.

given a kernel function $k(x, x')$ from the sigmoid family to approximate the cumulative probability distribution and it's derivative which we refer to as the cross-kernel $\mathcal{K}(x, x')$ which we use to construct an estimate of the probability distribution:

$$k(x, x') = \frac{1}{1 + e^{-\gamma(x-x')}} \quad (28)$$

$$\mathcal{K}(x, x') = -\frac{\gamma}{2 + e^{\gamma(x-x')} + e^{-\gamma(x-x')}} \quad (29)$$

4.4 Data Pre-Processing

5 Results

The architecture has been tested against several data sets. In all cases the system parameters are left unchanged to eliminate the possibility of optimizing the system performance to best match the known outcomes.

5.1 Eunite Competition Data

5.2 Santa Fe Data

5.3 CATS Benchmark Data

5.4 Results Summary

6 Further Research

6.1 Data Retention

6.2 Ensemble System

6.3 Weighted Influence

6.4 Sliding Prediction Offset

7 Conclusion

Eat it, bitches