

The CLP: Constrained Linear Predictors

June 28, 2023

The Goal

Consider a second order random process X , such that at each value of $t \in \mathbb{R}$, we have a random variable X_t . We may randomly sample this vector at n points, gaining a vector $\vec{T} = (t_1, t_2, t_3, \dots)$ of times at which the samples were made, and $\vec{X} = (X_{t_i})$. Strictly speaking these are both random variables in and of themselves, up until the moment that we ‘realise’ them. We can index into these vectors using the integer $0 \leq i < n$, and we assume without loss of generality that the samples are sorted in time, such that $t_i < t_{i+1} \forall i$.

In the case of the BLP, we wish to find a predictor, \hat{X}_t , which will predict the values of X_t on a set of ‘prediction points’, $t \in T$, subject to three further conditions:

- We are willing to present an *a priori* guess at the functional form of the predictor, in the form of a ‘prior function’ $g(t)$.
- The only thing we ‘know’ (or are willing to *ansatz*) about X_t is the second moment kernel (a generalisation of the covariance):

$$\langle (X_t - g(t))(X_s - g(s)) \rangle = k(t, s)$$

- Our predictor should be linear, such that:

$$\hat{X}_t = g(t) + \vec{a}_t \cdot (\vec{X} - \vec{G})$$

Where $G_i = g(t_i)$

We again reiterate that X_t , \vec{X} and \hat{X}_t are - strictly speaking - random variables until we make them into real numbers at the moment we wish to actually make a prediction. \vec{a}_t is a real n -tuple, which takes on different values at each value of t .

These are the ingredients of the standard BLP. The goal of this work is to extend this by adding one further condition:

- The Predictor should obey a number of constraints of the form $h(\{X_t\}) = 0$

We are therefore attempting to formulate the *Constrained Linear Predictor* (CLP)

1 Deriving the CLP

We define the CLP as the linear predictor which minimises the Mean Squared Error, averaged across all realisations of the random variable, computed at the set T of points at which we wish to make predictions, and which obeys our constraints.

Therefore, the CLP minimises the following Lagrangian:

$$\mathcal{L} = \sum_{t \in T} \langle (X_t - \hat{X}_t)^2 \rangle - \sum_j \lambda_j h_j(\{\hat{X}\}) \quad (1)$$

Here $h_j(\{\hat{X}_t\})$ is the j^{th} constraint on the *prediction points*¹, such that $h_j = 0$ when the constraint is met, and is non-zero otherwise, with the sum running over all such constraints. $\lambda_j \in \mathbb{R}$ are the associated Lagrange Multipliers. In the standard BLP we are able to treat the Lagrangian as separable in each element of T - minimising the MSE individually at each $t \in T$ is equivalent to performing a global minimisation: in the CLP this is not true, and we must consider the global case.

The issue at present is that we do not know what the behaviour of X_t is - we might have an initial guess (i.e. our prior, $g(t)$), but the entire purpose of this exercise is that we do not know X_t . However, by expanding out the brackets, we are able to write the Lagrangian in the following form:

$$\begin{aligned} \mathcal{L} &= \left[\sum_{t \in T} \langle X_t'^2 \rangle - 2\vec{a}_t \cdot \langle X_t' \vec{X}' \rangle + \langle (\vec{a}_t \cdot \vec{X}')^2 \rangle \right] - \sum_j \lambda_j h_j(\{\hat{X}\}) \\ &= \left[\sum_{t \in T} \langle X_t'^2 \rangle - 2\vec{a}_t \cdot \vec{k}_t + \vec{a}_t \cdot (K \vec{a}_t) \right] - \sum_j \lambda_j h_j(\{\hat{X}\}) \end{aligned} \quad (2)$$

Where:

$$\begin{aligned} X_t' &= X_t - g(t) \\ \vec{X}' &= \vec{X} - \vec{G} \\ \vec{k}_t &\in \mathbb{R}^n \text{ such that } [\vec{k}_t]_i = k(t, t_i) \\ K &\in \mathbb{R}^{n \times n} \text{ such that } K_{ij} = k(t_i, t_j) \end{aligned} \quad (3)$$

Note that since the kernel is, by definition, symmetric in its arguments, $K^T = K$. Note that we have also taken the explicit step of writing our kernel as a relationship between the *transformed* data - i.e. X' - the imposition of different functions $g(t)$ might therefore warrant different kernels. This is true even if the transform is the (commonly used) constant ‘mean scaling’, $g(t) = \langle X_t \rangle \approx \frac{1}{n} \vec{X} \cdot \mathbb{1}$.

By performing this transform we have placed the incomputable terms - that of $\langle (X_t')^2 \rangle$ into a constant term. Since Lagrangians are invariant under constant scalings, it is possible to find an optimal value of \vec{a}_t using only the remaining computable terms.

However - as we shall see - we are in the uncomfortable position of trying to impose conditions on the predicted values, $P_i = \hat{X}_{t_i} = g(t_i) + \vec{a}_{t_i} \cdot \vec{X}'$ whilst our object of interest is now the vector \vec{a}_{t_i} .

¹For clarity and avoidance of symbol-collision with the other X-s, we will denote the prediction points as $P_i = \hat{X}_{t_i} = g(t_i) + \vec{a}_t \cdot \vec{X}'$

We therefore limit ourselves to the case of *linear constraints*, i.e., those which can be written in the following form:

$$\begin{aligned} h_j(\{P\}) &= c_j - \sum_k d_{jk} P_k \\ &= c_j - \sum_k d_{jk} \left(g(t_k) + \vec{a}_{t_k} \cdot \vec{X}' \right) \end{aligned} \quad (4)$$

We can then take the derivative of the Lagrangian with respect to \vec{a}_{t_i} , and find that:

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial \vec{a}_{t_i}} &= 2K \vec{a}_{t_i} - 2\vec{k}_i - \sum_j \lambda_j \frac{\partial h_j}{\partial \vec{a}_{t_i}} \\ &= 2K \vec{a}_{t_i} - 2\vec{k}_i + \left(\sum_j \lambda_j b_{ji} \right) \vec{X}' \\ &= 2K \vec{a}_{t_i} - 2\vec{k}_i + \eta_i \vec{X}' \end{aligned} \quad (5)$$

Hence, the optimal value of \vec{a}_{t_i} is:

$$\begin{aligned} \vec{a}_{t_i} &= K^{-1} \left(\vec{k}_i - \frac{\eta_i}{2} \vec{X}' \right) \\ &= \vec{v}_i - \frac{\eta_i}{2} \vec{w} \end{aligned} \quad (6)$$

The optimal predicted value is:

$$\begin{aligned} P_i &= g(t_i) + \vec{a}_{t_i} \cdot \vec{X}' \\ &= g(t_i) + \vec{v}_i \cdot \vec{X}' - \frac{\eta_i}{2} \vec{w} \cdot \vec{X}' \\ &= g(t_i) + A_i - \frac{\eta_i}{2} B \end{aligned} \quad (7)$$

1.1 Exact Constraints

In the case where the constraints h_j are exact – i.e. the sets $\{c\}$ and $\{d\}$ are exactly determined, we may therefore analytically solve to find the set of Lagrange multipliers, then $\vec{\eta}$, and hence compute the predictor. We note that $\vec{\eta}$ can be written as:

$$\vec{\eta} = D^T \vec{\lambda} \quad (8)$$

Where $D_{ij} = d_{ij}$ is the constraint matrix, $\vec{\eta}_k = \eta_k$ is a vector on \mathbb{R}^N and $\vec{\lambda}_k = \lambda_k$ is a vector on \mathbb{R}^m , where m is the number of constraints. The requirement that the constraints are met can be written as:

$$D\vec{p} = \vec{c} \quad (9)$$

Where $\vec{p}_i = P_i$ is another vector on \mathbb{R}^n and $\vec{c}_i = c_i \in \mathbb{R}^m$. Writing $g(t_i) + A_i = q_i$, this is then:

$$D \left(\vec{q} - \frac{B}{2} D^T \vec{\lambda} \right) = \vec{c} \iff \vec{\lambda} = \frac{2}{B} (DD^T)^{-1} (D\vec{q} - \vec{c}) \quad (10)$$

Therefore:

$$\vec{p} = (\mathbf{1}_N - D^T(DD^T)^{-1}D) \vec{q} + D^T(DD^T)^{-1}\vec{c} \quad (11)$$

In the case where there is only a single constraint ($m = 1$), this simplifies such that $D \rightarrow \vec{d}^T$:

$$\vec{p} = \vec{q} + \frac{c - \vec{q} \cdot \vec{d}}{\vec{d}^2} \vec{d} \quad (12)$$

1.2 Inexact Constraints

In the case where the constraints are not exact, but serve to enforce bounds – i.e. monotonicity or positivity – there is a problem since the parameters of the constraint are not fixed. We may not care, for example, how much greater X_{i+1} is than X_i , only that it *is* greater.

We could enforce this through slack variables and utilise the KKT conditions, however for our purposes it is better to *parameterise* the constraint.

Various parameterisations are possible, but perhaps the most comprehensible is to consider that the *prediction* points, P_i are a function of some other parameters $\vec{\theta} \in R^m$, such that:

$$\begin{aligned} P_i &= \mathcal{T}_i(\vec{\theta}) \\ h_j(\mathcal{T}_i(\vec{\theta})) &= 0 \quad \forall i, j, \vec{\theta} \end{aligned} \quad (13)$$

For example, in the case of enforcing positivity, we might have that $P_i = e^{z_i}$, which is equivalent to asserting that $d_{ij} = \delta_{ij}$ and $c_i = e^{z_i}$. Rearranging Eq. (7), we are able to write η_i as a function of this Transform, and hence write \vec{a}_{t_i} in the following form:

$$\vec{a}_{t_i} = \vec{v}_i + \frac{P_i(\vec{\theta}) - A_i - g(t_i)}{B} \vec{w} \quad (14)$$

This might seem somewhat tautological – we have written \vec{a}_{t_i} in terms of the prediction values – but the entire purpose of \vec{a}_{t_i} is to make predictions!

The usefulness of this comes evident when we insert Eq. (14) back into the Lagrangian – essentially performing a change of coordinates from $\mathcal{L}(\vec{a}, \vec{\theta})$ to $\mathcal{L}(\vec{\theta})$, since we have now ensured that \vec{a}_t will always be at its optimal value for each value of $\vec{\theta}$.

$$\vec{k}_i \cdot \vec{a}_{t_i} = \vec{v}_i \cdot \vec{k}_i + \frac{P_i(\vec{\theta}) - A_i - g(t_i)}{B} \vec{w} \cdot \vec{k}_i \quad (15)$$

$$\begin{aligned} \vec{a}_{t_i} \cdot (K \vec{a}_{t_i}) &= \left(\vec{v}_i + \frac{P_i(\vec{\theta}) - A_i - g(t_i)}{B} \vec{w} \right) \cdot \left(\vec{k}_i + \frac{P_i(\vec{\theta}) - A_i - g(t_i)}{B} \vec{w} \right) \\ &= \vec{v}_i \cdot \vec{k}_i + \left(\frac{P_i(\vec{\theta}) - A_i - g(t_i)}{B} \right) (\vec{w} \cdot \vec{k}_i + A_i) + \frac{(P_i(\vec{\theta}) - A_i - g(t_i))^2}{B} \end{aligned} \quad (16)$$

Since $\vec{w} \cdot \vec{k}_i = (K^{-1} \vec{X}') \vec{k}_i = (K^{-1} \vec{k}_i) \vec{X}' = \vec{v}_i \cdot \vec{X}' = A_i$ due to the symmetry of K , and the constraints are all automatically satisfied thanks to our parameterisation, we find that the Lagrangian simplifies to:

$$\begin{aligned} \mathcal{L}(\vec{\theta}) &= \sum_i \left(\langle (X'_i)^2 \rangle - \vec{k}_i \cdot \vec{v}_i \right) + \frac{1}{B} (P_i(\theta) - A_i - g(t_i))^2 \\ &= \text{const in } \vec{\theta} + \frac{1}{B} \sum_i (P_i(\theta) - A_i - g(t_i))^2 \\ \mathcal{L}' &= \sum_i P_i (P_i(\theta) - 2(A_i + g(t_i))) \end{aligned} \tag{17}$$

Where in the final line we took the opportunity to perform a rescaling (recalling that $B > 0$ is enforced by the positive definiteness of K) which leaves the optimum invariant. In some cases it is trivial to identify the optimal values of P_i - for example, in the case where $P_i = e^{\theta_i}$, the maximum is evidently:

$$P_i = \begin{cases} A_i + g(t_i) & \text{if this is } > 0 \\ 0 & \text{else} \end{cases} \tag{18}$$

In short, the CLP is equal to the BLP except when the condition is violated, at which point a hard cut is placed on it.

More complex conditions however, can lead to more complex behaviour - the monotonicity constraint, for example, exhibits the obvious behaviour that it again follows the BLP when it is monotonic, and is flat when the BLP has a negative gradient - but the *location* where the CLP becomes flat is non-trivial, with flatness necessarily occurring *before* the BLP changes direction: a tradeoff in following the BLP locally versus becoming too large too early without the ability to decrease due to the monotonic constraint.

In these cases a more complex search is required - where the behaviour of the constraint is evident *a priori* (such as the monotonic constraint), one can limit the space of the search. In the general case, however, a numerical optimisation is required.

The derivative of the Lagrangian with respect to the constraint parameters is:

$$\frac{\partial \mathcal{L}'}{\partial \theta_m} = 2 \sum_i (P_i - A_i - g(t_i)) \frac{\partial P_i}{\partial \theta_m} \tag{19}$$

This can be used to numerically optimise the values of $\vec{\theta}$

1.3 Inexact Constraints (Redux)

We note that we performed a fairly drastic change in approach between the exact constraints and the inexact constraints - is it possible to maintain the same approach for both?

We consider now that the parameters \vec{c} of the constraints are functions of an (unconstrained) external parameter, $\vec{z} \in \mathbb{R}^m$ - letting $\vec{c} = \text{const}$ recovers the condition of the exact equalities. However, in any other case we must still find the values of \vec{z} which optimise the global Lagrangian - and hence we need to rewrite our Lagrangian in terms of \vec{c} .

From Eq. (11), we can rewrite the predicted value-vector (recalling that $\vec{p}_i = P_i = \hat{X}_{t_i}$) as:

$$\begin{aligned}\vec{p} &= \vec{j} + R\vec{c}(\vec{z}) \\ R &= D^T(DD^T)^{-1} \\ \vec{j} = (\mathbb{1}_N - RD)\vec{q} &\iff j_i = g(t_i) + A_i + \sum_{j,k} R_{ij}D_{jk}(g(t_k) + A_k)\end{aligned}\tag{20}$$

We note that from a conceptual standpoint it is not a problem for the ‘mixing’ constraints D_{ij} to be the functions of \vec{z} , but this assumption allows us to precompute many of the otherwise troublesome entities. We can also rewrite \vec{a}_{t_i} as:

$$\begin{aligned}\vec{a}_{t_i} &= \vec{v}_i - \frac{\eta_i}{2}\vec{w} \\ &= \vec{v}_i + \frac{[R(\vec{c} - D\vec{q})] \cdot \hat{e}_i}{B}\vec{w} \\ &= \vec{j}_i + \frac{(R\vec{c}) \cdot \hat{e}_i}{B}\vec{w}\end{aligned}\tag{21}$$

Where

$$\begin{aligned}R &= D^T(DD^T)^{-1} \\ \vec{j}_i &= \vec{v}_i - \frac{(RD\vec{q}) \cdot \hat{e}_i}{B}\vec{w}\end{aligned}\tag{22}$$

We therefore have:

$$\begin{aligned}\vec{k}_i \cdot \vec{a}_{t_i} &= \vec{v}_i \cdot \vec{k}_i + \frac{A_i}{B}((R\vec{c}) \cdot \hat{e}_i - (RD\vec{q}) \cdot \hat{e}_i) \\ &= \text{const in } \vec{c} + \frac{A_i}{B}(R\vec{c}) \cdot \hat{e}_i \\ \vec{a}_{t_i} \cdot (K\vec{a}_{t_i}) &= \left(\vec{j}_i + \frac{(R\vec{c}) \cdot \hat{e}_i}{B}\vec{w}\right) \cdot \left(K\vec{j}_i + \frac{(R\vec{c}) \cdot \hat{e}_i}{B}\vec{X}\right) \\ &= \text{const in } \vec{c} + 2\frac{(R\vec{c}) \cdot \hat{e}_i}{B}\vec{j}_i \cdot \vec{X} + \frac{1}{B}((R\vec{c}) \cdot \hat{e}_i)^2 \\ &= \text{const in } \vec{c} + 2\frac{(R\vec{c}) \cdot \hat{e}_i}{B}(A_i - (RD\vec{q}) \cdot \hat{e}_i) + \frac{1}{B}((R\vec{c}) \cdot \hat{e}_i)^2\end{aligned}\tag{23}$$

Therefore:

$$\begin{aligned}\mathcal{L}' &= \sum_i \vec{a}_{t_i} \cdot K\vec{a}_{t_i} - 2\vec{k}_i \cdot \vec{a}_{t_i} \\ &= \text{const in } \vec{c} + \sum_i ((R\vec{c}) \cdot \hat{e}_i)^2 - 2(R\vec{c}) \cdot \hat{e}_i(RD\vec{q}) \cdot \hat{e}_i \\ &= \text{const in } \vec{c} + (R\vec{c})^2 - 2(R\vec{c}) \cdot (RD\vec{q}) \\ &= \text{const in } \vec{c} + (R\vec{c}(\vec{z}) - RD\vec{q})^2\end{aligned}\tag{24}$$

The derivative with respect to the (unconstrained) vectors \vec{z} is:

$$\begin{aligned}\frac{\partial \mathcal{L}'}{\partial z_m} &= (R\vec{c}(\vec{z}) - RD\vec{q}) \cdot R \frac{\partial \vec{c}}{\partial z_m} \\ &= (\vec{p}(\vec{z}) - \vec{q}) \cdot R \frac{\partial \vec{c}}{\partial z_m}\end{aligned}\tag{25}$$

Since \vec{q} is the BLP prediction we can once again see that the derivative is zero if the BLP obeys the constraints ($\vec{c} - D\vec{q} = 0$), so the CLP will always revert to the BLP if this meets our constraints.

2 The CLUP

In the prior work, we assumed the the function $g(t)$ was ‘handed down’ to us to act as a prior function. However, this may induce biases in our predictor, meaning that:

$$\langle X_t - \hat{X}_t \rangle \neq 0\tag{26}$$

We should ideally search for an *unbiased* predictor. The derivation of the Constrained Linear Unbiased Predictor (CLUP) should follow along similar lines to the standard BLUP, but we reproduce it in full for the sake of rigour.

We suppose that our random variable X_t can be written as:

$$X_t = m(t) + Y_t\tag{27}$$

Where $m : R \rightarrow R$ is the ‘mean function’ and Y_t is a zero-mean random variable. Therefore:

$$\langle X_t \rangle = m(t)\tag{28}$$

The main difference between $m(t)$ and $g(t)$ is that we assumed $g(t)$ was just a prior to ‘help us along’ without any intrinsic relation to X_t – here, however, we are asserting that $m(t)$ is a meaningful function – albeit an incomputable one, since we remain unwilling to assert any properties on $\langle X_t \rangle$. Because of this restriction, we cannot subtract away $m(t)$ from our data to formulate $\vec{X}' = \vec{Y}$ – we must keep everything in terms of our original, untransformed data.

We *can*, however, assert that $m(t)$ can be decomposed into a sum of basis functions, $\vec{\varphi}(t)$, where $\varphi_i(t) : \mathbb{R} \rightarrow \mathbb{R}$ is the i^{th} basis function. We therefore have:

$$\begin{aligned}m(t) &= \sum_{i=0}^{\omega} \phi_i(t) \beta_i \\ &= \vec{\beta} \cdot \vec{\varphi}(t)\end{aligned}\tag{29}$$

We again note that $\vec{\beta}$ is not a known value, however, we continue in the expectation that it will cancel out in future. We also note that without further information we must assume that $\omega \rightarrow \infty$, in practice we can limit the dimensionality by assuming that $m_{\omega}(t) \approx m(t)$ for finite ω . We can also formulate the matrix Φ :

$$\Phi \in \mathbb{R}^{\omega \times n} \text{ such that } \Phi_{ij} = \varphi_i(t_j)\tag{30}$$

Therefore:

$$\vec{X}_i = \sum_j^\omega \Phi_{ji} \beta_j \iff \vec{X} = \Phi^T \vec{\beta} + \vec{Y} \quad (31)$$

Our linear predictor takes the form:

$$\begin{aligned} \hat{X}_t &= \vec{a}_t \cdot \vec{X} \\ &= \vec{a}_t \cdot (\Phi^T \vec{\beta}) + \vec{a}_t \cdot \vec{Y} \\ &= \vec{\beta} \cdot (\Phi \vec{a}_t) + \vec{a}_t \cdot \vec{Y} \end{aligned} \quad (32)$$

We can also note that:

$$X_t - \hat{X}_t = (\Phi \vec{a}_t - \vec{\varphi}) \cdot \vec{\beta} + Y_t - \vec{a}_t \cdot \vec{Y} \quad (33)$$

Since $\langle Y \rangle = \vec{0}$ and $\langle Y_t \rangle = 0$ by definition, we note that the unbiased constraint is equal to:

$$\langle X_t - \hat{X}_t \rangle = 0 \iff (\Phi \vec{a}_t - \vec{\varphi}) \cdot \langle \vec{\beta} \rangle = 0 \quad (34)$$

Therefore if we write our unbiased constraint as $\mathcal{U}_t = \langle X_t - \hat{X}_t \rangle$, such that $\mathcal{U}_t = 0$ when the constraint is met:

$$\begin{aligned} \mu_t \mathcal{U}_t &= \left(\mu_t \langle \vec{\beta} \rangle \right) \cdot (\Phi \vec{a}_t - \vec{\varphi}) \\ &= \tilde{\mu}_t \cdot (\Phi \vec{a}_t - \vec{\varphi}) \end{aligned} \quad (35)$$

A suitable change of coordinates $\mu_t \rightarrow \tilde{\mu}_t$ therefore enables us to bypass the unknown $\langle \vec{\beta} \rangle$. Hence the Lagrangian of the system takes the form:

$$\begin{aligned} \mathcal{L}(\vec{a}, \vec{\lambda}, \{\tilde{\mu}\}) &= \sum_{t \in T} \left(\langle (X_t - \hat{X}_t)^2 \rangle + \tilde{\mu} \cdot (\Phi \vec{a}_t - \vec{\varphi}) \right) + \sum_j \lambda_j h_j(\{\hat{X}\}) \\ &= \sum_{t \in T} \left(\langle X_t^2 \rangle + \vec{a}_t \cdot (K \vec{a}_t) - 2 \vec{k}_t \cdot \vec{a}_t + \tilde{\mu} \cdot (\Phi \vec{a}_t - \vec{\varphi}) \right) + \sum_j \lambda_j h_j(\{\hat{X}\}) \end{aligned} \quad (36)$$

This is identical in form to Eq. (7), with the addition of some additional constraints - those labelled by $\tilde{\mu}$, which act to ensure that $\langle X_t - \hat{X}_t \rangle = 0$ for all $t \in T$, i.e. we are now including *unbiasedness* as a constraint.

We have also introduced \vec{k}_t and K as the second moment matrices on X :

$$\begin{aligned} K &\in \mathbb{R}^{n \times n} & K_{ij} &= \langle X_{t_i} X_{t_j} \rangle \\ \vec{k}_t &\in \mathbb{R}^n & [\vec{k}_t]_i &= \langle X_t X_{t_i} \rangle \end{aligned} \quad (37)$$

We once again impose the linearity condition on our constraints, such that $h_j = c_j - \sum_k d_{jk} \hat{X}_{t_k}$:

$$\vec{h} = \vec{c} - D \vec{p} \quad (38)$$

Where $\vec{p}_i = \hat{X}_{t_i}$. The Lagrangian therefore simplifies to:

$$\mathcal{L}(\vec{a}, \vec{\lambda}, \{\tilde{\mu}\}, \vec{c}) = \sum_{t \in T} \left(\langle X_t^2 \rangle + \vec{a}_t \cdot (K \vec{a}_t) - 2 \vec{k}_t \cdot \vec{a}_t + \tilde{\mu} \cdot (\Phi \vec{a}_t - \vec{\varphi}) \right) + \vec{\lambda} \cdot (\vec{c} - D \vec{p}) \quad (39)$$

The Lagrangian derivatives are:

$$\frac{\partial \mathcal{L}}{\partial \vec{a}_{t_i}} = 2K \vec{a}_{t_i} - 2 \vec{a}_t \cdot \vec{k}_i + \Phi^T \tilde{\mu} + \frac{\partial \vec{p}}{\partial \vec{a}_{t_i}} (D^T \vec{\lambda}) \quad (40)$$

$$= 2K \vec{a}_{t_i} - 2 \vec{k}_i + \Phi^T \tilde{\mu} + Q_i(\vec{X}) D^T \vec{\lambda}$$

$$\frac{\partial \mathcal{L}}{\partial \vec{\lambda}} = \vec{c} - D \vec{p} \quad (41)$$

$$\frac{\partial \mathcal{L}}{\partial \tilde{\mu}} = \Phi_t \vec{a}_t - \vec{\varphi}_t \quad (42)$$

Where $Q_i(\vec{X}) \vec{v} = \vec{v} \cdot \hat{e}_i \vec{X}$, and equivalently $Q_i^T \vec{v} = \vec{v} \cdot \vec{X} \hat{e}_i$. The optimal value of \vec{a}_t is therefore at:

$$\begin{aligned} \vec{a}_{t_i} &= K^{-1} \vec{a}_{t_i} - 2K^{-1} \Phi_i^T \tilde{\mu}_i + K^{-1} Q_i D^T \vec{\lambda} \\ &= \vec{v}_i - 2K^{-1} \Phi_i^T \tilde{\mu}_i + (D^T \vec{\lambda}) \cdot \hat{e}_i \vec{w} \end{aligned} \quad (43)$$

2.0.1 Applying Constraints

We now impose the unbiased constraint - substituting Eq. (43) into Eq. (42), which gives us:

$$\Phi_i \vec{v}_i - 2\Phi_i K^{-1} \Phi_i^T \tilde{\mu}_i + (D^T \vec{\lambda}) \cdot \hat{e}_i \Phi_t \vec{w} = \vec{\varphi} \quad (44)$$

Solving for $\tilde{\mu}_i$:

$$\begin{aligned} \tilde{\mu}_i &= \frac{1}{2} (\Phi_i K^{-1} \Phi_i^T)^{-1} \left(\Phi_i \vec{v}_i + (D^T \vec{\lambda}) \cdot \hat{e}_i \Phi_t \vec{w} - \vec{\varphi} \right) \\ &= \frac{1}{2} M^{i-1} \left(\Phi_i \vec{v}_i + (D^T \vec{\lambda}) \cdot \hat{e}_i \Phi_t \vec{w} - \vec{\varphi} \right) \end{aligned} \quad (45)$$

And therefore

$$\begin{aligned} \vec{a}_{t_i} &= \left(\mathbb{1} - K^{-1} \Phi_i^T M^{i-1} \Phi_i \right) \vec{v}_i + K^{-1} \Phi_i^T M^{i-1} \vec{\varphi} + (D^T \vec{\lambda}) \cdot \hat{e}_i \left(\mathbb{1} - K^{-1} \Phi_i^T M^{i-1} \Phi_i \right) \vec{w} \\ &= B^i \vec{v}_i + C^i \vec{\varphi} + (D^T \vec{\lambda}) \cdot \hat{e}_i B^i \vec{w} \end{aligned} \quad (46)$$

where $C^i = K^{-1} \Phi_i^T M^{i-1}$

and $B^i = (\mathbb{1} - C^i \Phi_i)$

The prediction values are therefore:

$$\begin{aligned} \hat{X}_i &= P_i = \vec{a}_{t_i} \cdot \vec{X} \\ &= (B^i \vec{v}_i + C^i \vec{\varphi}_i) \cdot \vec{X} + (D^T \vec{\lambda}) \cdot \hat{e}_i \vec{X} \cdot B^i K^{-1} \vec{X} \\ &= \alpha_i + (D^T \vec{\lambda}) \cdot \hat{e}_i \beta_i \end{aligned} \quad (47)$$

Where α_i and β_i are generalisations of the A_i and B_i terms used in the CLP. We can then form a vector of $\vec{p}_i = P_i$, such that:

$$\begin{aligned}\vec{p} &= \vec{\alpha} + \vec{\beta} \otimes D^T \vec{\lambda} \\ &= \vec{\alpha} + \mathcal{B} D^T \vec{\lambda}\end{aligned}\tag{48}$$

Here \otimes is the element-wise Hadamard product, and $\mathcal{B} = \text{diag}(\vec{\beta})$. Inserting this into the imposed constraints of Eq. (41), we find:

$$\begin{aligned}D\vec{p} &= \vec{c} \\ D\vec{\alpha} + D\mathcal{B}D^T \vec{\lambda} &= \vec{c} \\ \vec{\lambda} &= (D\mathcal{B}D^T)^{-1} (\vec{c} - D\vec{\alpha})\end{aligned}\tag{49}$$

We can then insert this back into the definition of \vec{a}_{t_i} to find:

$$\vec{a}_{t_i} = B^i \vec{v}_i + C^i \vec{\varphi} + \left(D^T (D\mathcal{B}D^T)^{-1} (\vec{c} - D\vec{\alpha}) \right) \cdot \hat{e}_i B^i \vec{w}\tag{50}$$

Where:

$$\begin{aligned}\vec{v}_i &= K^{-1} \vec{k}_i \\ \vec{w} &= K^{-1} \vec{X} \\ M^i &= \Phi_i K^1 \Phi_i^T \\ C^i &= K^{-1} \Phi_i^T M^{i,-1} \\ B^i &= \mathbb{1} - C^i \Phi_i \\ \alpha_i &= (B^i \vec{v}_i + C^i \vec{\varphi}_i) \cdot \vec{X} \\ \beta_i &= (B^i \vec{w}) \cdot \vec{X}\end{aligned}\tag{51}$$

3 Optimising the Kernel